# Unit-1

## Introduction

## IntrotoCourse

Welcome to Privacy and Security in Online Social Media course on NPTEL. I am PK. I am faculty at IIIT, Delhi. I received my PhD from Carnegie Mellon University and my primary area of interest is Privacy, Security and Computational Social Science, Data Science, Social Computing and topics surrounded. I am a part of Cyber Security Education and Research Center at IIIT,Delhi. I am also a part of Research Group called Precog, which primarily works on privacy and security in online social media, computational social science, data science, social computing and usable technologies which are around these topics.
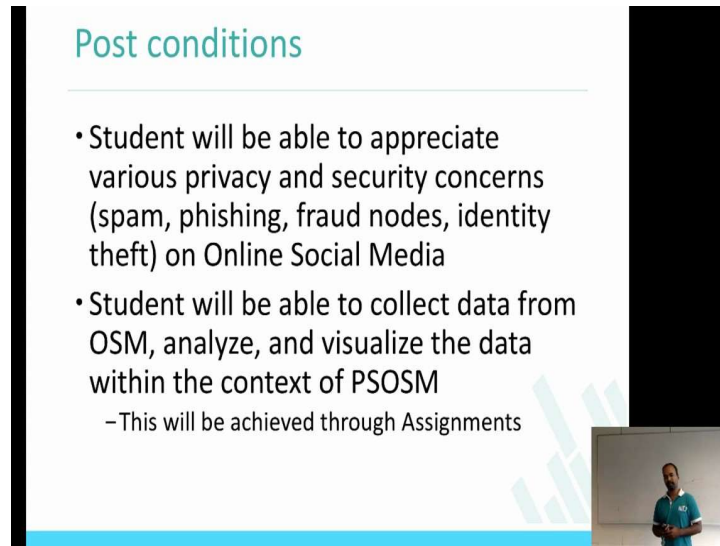
(ReferSlideTime:00:52)



This is the Facebook post, I did about a year and half back which is regarding the course feedback that students give us at IIIT, Delhi. I have been teaching this course PSOSM Privacy and Security in Online Social Media at IIIT,Delhi for a couple of times. So, this isthefeedbackandweareretherewiththecourse.We actuallyhaveapostsessionwhere

students actually present their work, what they have done over the semester in the form of a poster, in the form of a demo. So, this is the picture with all the students from the class, wearing the same T shirt and with actually the external evaluators who to came to evaluate these projects.
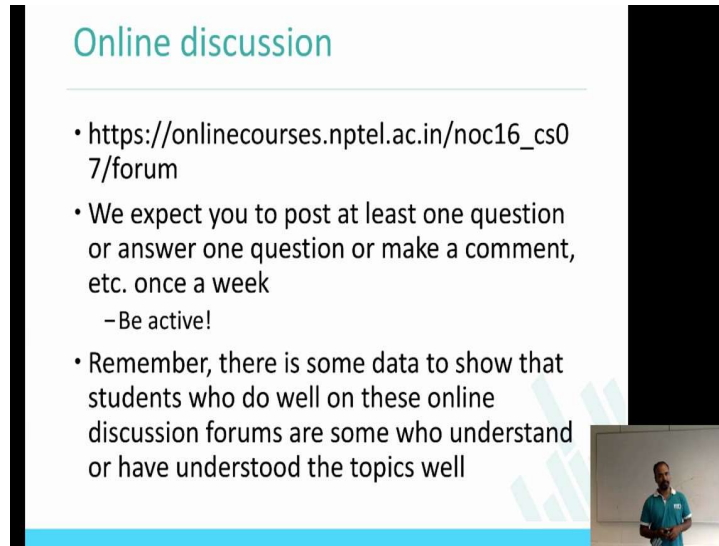
(ReferSlideTime:01:33)



So, what do you get out of this course? Some of the post conditions for this course is going to be at the end of the course, you will be able to appreciate various privacy and securityconcerns, spam, phishing,fraud, identitytheftand related issues on online social networks. Then the primary focus of this course is going to be different aspects of security and privacy on online social media.

Throughout the course, you will also be exposed to actually collecting data from online socialnetworkslikeFacebook,Twitter analyzingthese content andvisualizing thisdata in terms of the question thatyou are trying to ask, forexample, 1percent could be Iwant to understandwhether thefollowersthatIhaveonTwitterareactuallylegitimateorfake. Wecan do actually have achieve this goal, this post condition in terms of able to collect data, analyze data and visualize data through the assignments that you could be getting across these course.
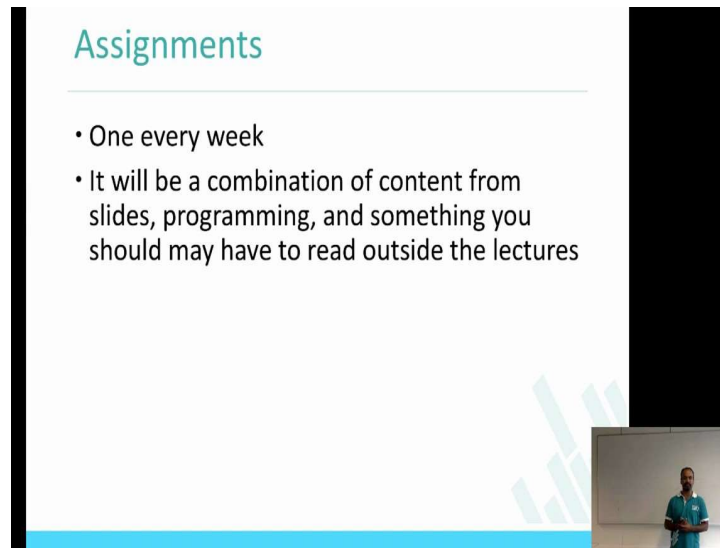
So, one of things that we primarily want to actually focus on is also about discussion around the topics that we will be discussing in this course. Here is a link to the online forum; we hope that you would actually participate. We expect you to actually post at least one question or answer or have a question, make comment, etcetera once per week. There is already research literature to show that people were active outside the class to perform well in the topics that are actually discussed inside the class.
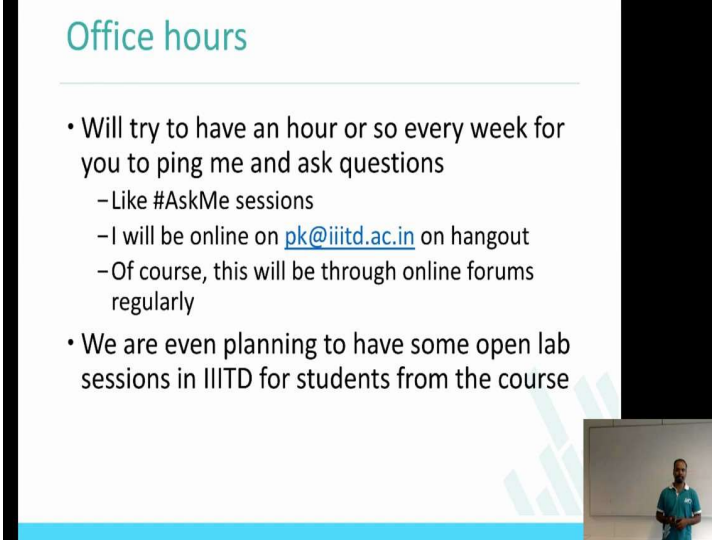
(ReferSlideTime:03:11)



So, what we plan to do in some of assignments is, we plan to have one homework assignment per weekcapturingthe topics thatwehave coveredin theclass,in thelecture and we will actually have homework questions around that every week and we hope to actually get you a sense of these kind of different topics same as collecting data from onlinesocialnetwork,whatkindof analysiscanbedonewiththisdata?Howtovisualize a data and things like that? Mostly these questions will be from slides, some programming; something that should actually be able to answer if you actually read content outside the lectures in the pointers that we discuss in the forum also.
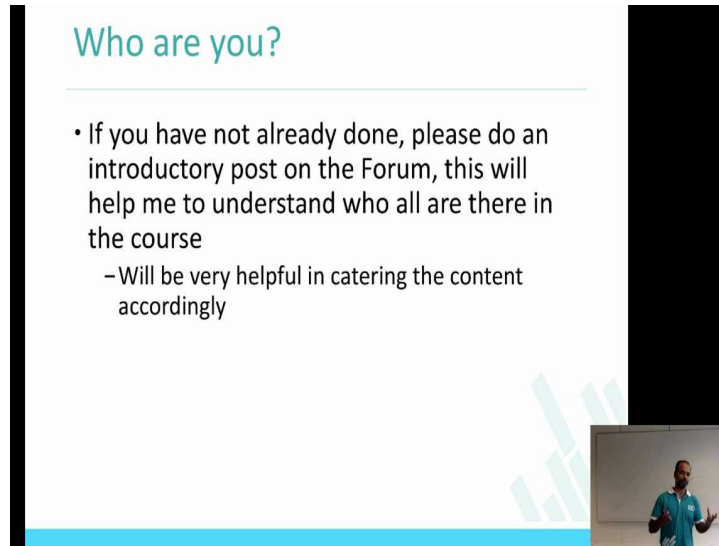
(ReferSlideTime:04:00)



Ialsoplantohavesomethingtomakeitmoreexiting,tomakeitmore interestingforyou to actually participate in the topic and I also plan to have some office hours where I would be online on Hangout, where you can actually ask me questions or of course, the forum also can be used for asking these questions, it's basically like hashtag AskMe sessions that you may have heard about in the past or read it on twitter and other social networks. This will also allow you to interact with me directly, probably on hangout it could even try video sessions.

Iamalsoplanningtohave someopenlab sessions at IIIT,Delhi, where you canactually, this is only for students who are going to be mostly in and around Delhi. If you are in Delhi, you could actually show up on campus sometime at decided time that we will let you know. Join the open lab sessions where we could actually have the TAand others answer some questions for you, help you do the course better,help you understand some concepts even more deeply.

(ReferSlideTime:05:14)



So, I already saw some of you introducing itself on the mailing list. It is actually great to see that whether more than 1500 students who have registered for this course. It will be nice to actually have most of you introducing yourself on the forum and the main reason forme tounderstand whoyouareis actuallyhelpingme tocater thecontentaccordingly. It will help me to create content, it will help me actuallygive you appropriate pointers, if I know the proportion of the distribution of the students who are taking this class.

(ReferSlideTime:05:52)



Why should I teach this class? As I said earlier, I have been teaching this class at IIIT, Delhi couple of times, but before that; starting to teach this class I did things which actually makes very interesting way into teaching this course. I did a work shop at UFMG which is Universidade Federal de Minas Gerais in Brazil. I have done some workshops around the topic of privacy and security in online social media whose work converted into a conference called aconference on online social network. I also have taught this class in Brazil, which is a full credit course over the summers in 2012 and 2013.

Teaching Assistants, it is not going to be just me, you are not going to just listen to me over the entire course. Teaching assistants will help us in creating the questions, doing some lab sessions helping us in giving you more content wherever necessary,helping us even creating the home works and things like that and for now we have full fabulous TA's, who are all my PhD students at IIIT, Delhi that is Anupama Agarwal, whose primary interest is actually understanding social reputation on social networks.

Thatis Srishti Gupta,whoseprimaryinterestisstudyingtheonlinesocialmediaandwith the phone numbers and OTT kind of technologies that are available. That is Prateek, whoseinterestis onstudying malicious contentonFacebook.Thatis NiharikaSachdeva, whose interest is primarily on visible security and studying how technology and social networks have been used by police organizations around the world and particularly in India.So,these4TA'swillhelpusindoing lab sessions,settingupquestionsanswersfor the homework and helping us in generalmaking the coursemore exciting and interesting for us.

Now, let us look at topics that we will cover over this course. Initially, I will describe differentaspectsofonlinesocialnetworks,whatarethedifferentsocialnetworksthatare available which are popular? What kind of terms you need to know before we actually start delving into these social networks more deeply? Then we will also get hands-on experience with setting up python; understand how to collect data from Twitter API? How do we store the data in Mongo DB, MySQL? This will be more like a lab session where we walk you through on how to set this up.

Later in the course we will start looking at trust and credibility which is how much can you actually delete the content that are posted on Twitter or Facebook? What kind of problems exists? What kind of techniques are available to actually identify whether the post is credible or not.

Then we will also look at privacy issues on online social network. Privacy is becoming such a big topic because of the proliferation of online social networks, what information is leaked? What information can be actually collated? What information can be stitched together to create a profile user or which can be actually misused against. Then we also look at social network analysis, text analytics that can be done using the content for socialnetwork.NLTKisoneoftheplatformswhichwewillalsoexposeyoutoanalyze

thecontentfromsocialnetworks.Thisagainwouldbeahands-onsessionwhereyouwill get experience on using these tools, techniques from the content in the course. E-crime which is also much related and very important topic in the context of online social networks is something we had also covered.

In this part of the course you will actually look at phishing. Wewill actuallylook at fake content, fake accounts and related topics.Wewillalso giveyou somehands-on details in terms of actually drawing graph with plotly,analyzing the data with highcharts, creating graph with high charts and also geo-location analysis because some of the content thatwe will be analyzing during the course and looking at during the course will be actually, will have information of geo-location, which is latitude longitude from a particular location where the post has been done. So, tools like these which is probably high chart, NLTK and social network analysis tools like ora will be actually very, very hands-on experience for you where you get a whole lot of ways to actually analyze the data anything that is relevant to online social networks.

Next, we look at policing which is in India, particularly if you see online social network has become such a big platform for police organizations to use in terms of interacting with citizens. So, we will actually study how police organizations are using this online social media for increasing their effectiveness of keeping this society safely and we will also look at how citizens have beem using social networks to interact with police organizations. We will not just look at only Indian context, we will also look at broader context in the world, how organizations are actually using it. There is also this whole topic of identity resolution which we will cover, which is in my case my Facebook handle is ponnurangam.kumaraguru, my Twitter handle is ponguru and my YouTube account is PK.

If you were to understand whether these three accounts are actually same is actually avery hard problem. So, we will actually look at some of the identity resolutiontechniques that people have created, how we can actually stitch these accounts together? It can be actually very useful for multiple reasons, one it could be useful for advertising agencies to actually present the ads appropriately. You could also be useful for making decisions on whether it is the same person talking about in multiple social networks.
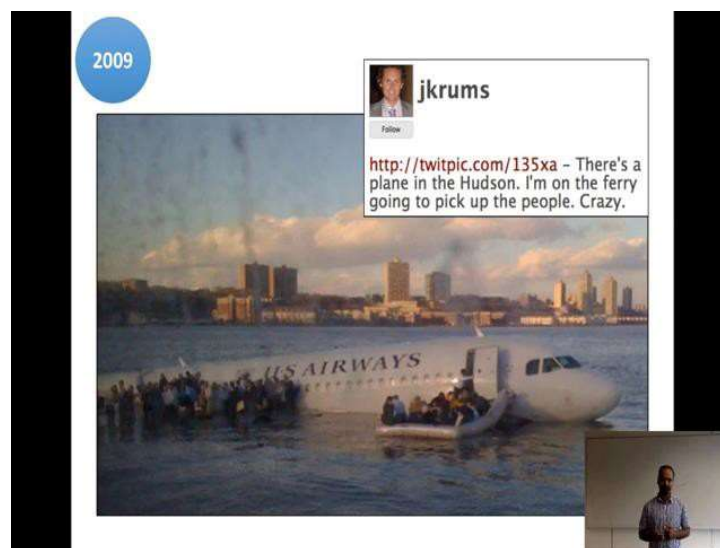
At the end of the course, we will actually review with some very broad questions and verybroadtopicswhichareconnectedtosocialnetworkswhichwewillnotabletocover in this course like deep learning, machine learning, national language processing, image analysis. These are the topics are becoming very, very popular in terms of using social network data because of the proliferation of the social network and understanding what data is available and how we can use this data is becoming a very important topic.

Mostly if you all see, image analysis is also becoming an important one because mostly these days theposts arecomingwith images. Itis notjust text only; it is actuallytext and images or sometimes only images. So, this will be a very broad tour of these new topics that are popping up around the online social network topic.

# Incidents

Welcome back, until now in the course, we have seen some logistics about the course what online social media is? What is the impact and numbers and some ways in whichwe can actually use this social media services and some examples of these social media services? What I will cover today is actually looking at some of the incidences, both positive and negative where social media has actually a played role.

(ReferSlideTime:00:42)



Here is a first one, as I have put in the slide, it happened in 2009. This is the first time ever social media service like Twitter was actually used for crisis management. Here is a tweet which actually jkrums, j k rums actually posted, which reads as, 'There is a planein the Hudson. I am on the ferry going to pick up the people. Crazy'. Until then Twitter was basicallyused for conversations, basicallyused for saying what I am doing in life in the morning, that we are posting, 'Morning Monday', 'I am having coffee', 'I am traveling here' and things like that.

Whereas first time in 2009, jkrums actually posted this tweet, where when the US Airways flight landed in the Hudson river this post came out and before the first responders could reach, there was actually public, citizens who were actually helping in the situation. Therefore, Twitter and social media services have started being used in many different ways and I am going to talk to you about some examples, someincidenceswere socialmedias' playedabout positive andthe negativerole. Inthiscase it is verypositive because it helps in actually solving crisis and help in first responder.
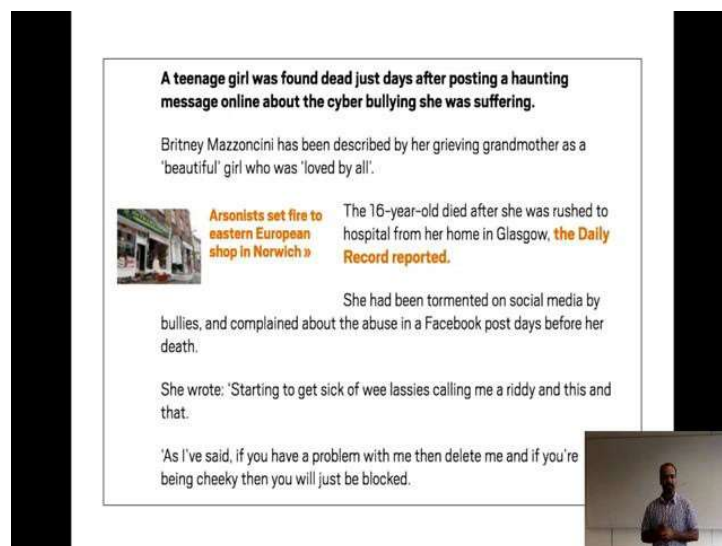
(ReferSlideTime:02:06)



Here is an example, where in India there was a kid who is actually lost in a railway station and somebody took a picture of this kid with railway the police officer and in 20 minutes she was actually able to connect with her parents. The primary way by which it was done was actually the picture was posted on social media particularly tagging, mentioning the concern that I said earlier in my lecture mentioning the Indian railways minister and therefore this tweet got viral, this picture got viral and the kid was able to connect with her parents within 20 minutes.

(ReferSlideTime:02:50)



Keeping some of these examples, particularly with missing child there has been also organization which has been started only to help finding out missing child and parents through social media. So, here is one example which is 'find your missing child' using only social media services, and this website actually provides lot oftips on how one can use these social media services to connect with children to find out the missing children.

(ReferSlideTime:03:21)



Now, here is another example where a teenager was found dead just days after posting a messageonsocialmedia.So,thisismorelikeanegativething,wherethegirlposted

information on her Facebook account, where she was saying about, she is being bullied. Cyber bullying is one of the big problems on social media also. We will talk about itlittle later in the course, but because of the cyber bullying people have actually killed themselves and many a times they actually leave a message on their own socialnetworks. In this case, she is being posting information about being cyberbullied done on Facebook, at the end of the day, she actually gives her life.

(ReferSlideTime:04:08)



Also, inthe world if you lookat it whensocialmedia was started using, were being used for crisis, this is the first time when social media was actually used to create or to propagate an incident. In the 2009 Hudson River, it was actually used for solving a problem, whereas in this case UK riots made worse by rolling news on Twitter and Facebook. So, the messages were sent on Twitter or Facebook saying that lets go to this site in this street.

(ReferSlideTime:04:47)



There have been many incidences like this which you may have come across, where it is not that content on social media is talking about an incident that happens on in the real world, but these day social media itself as being used for organizing an event.
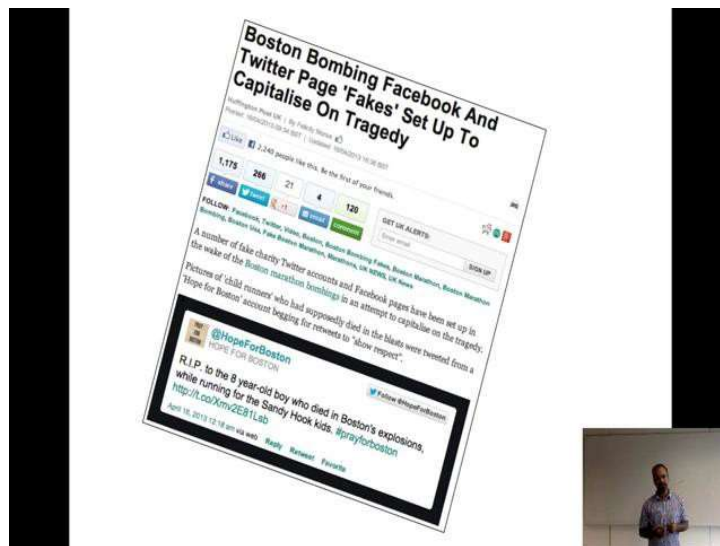
So in this example, if you see news articles which came of talking about 'Nepal earthquake: Government using social media to connect and provide relief'. Particularlyin India, if you see there is lot more usage of social media for interacting with citizens and this is being done for the last at least a year or year and half, where it has become more, the social networks have proliferated the way that the government is organizing themselves and interacting with citizens.

(ReferSlideTime:05:39)



Another incident that happened in the world which also is something that, if you read about the effect of social media, this incident would actually be one of the most phenomenal event that happened in terms of using social networks. So this was main social media that was used in this case was actually Twitter, where they were actually using twitter to connect with citizens giving them to one place and revolution happened becauseofusingthesesocialnetworkwhere it went veryviralthroughtheseservices like Twitter and Facebook.

(ReferSlideTime:06:19)

It is not only that social media services are being used or misused in situations like the Egypt or situations like finding kids. There is also a lot of misinformation that are actually floating around on social media services, for example, in this case Boston Bombing that happened, there was a tweet which said that, 'RIP to the 8 year old boy who died in Boston's explosions, while running for the Sandy Hook kids'. There was no 8 year old kid, who was actually participating in the Boston marathon. There was also another tweet during this Boston marathon said that, please RT this tweet, which is retweet this tweet, we will actually pay 1 dollars to Boston marathon. There was no money that was actually transferred to Boston marathon and this tweet get retweeted many many times, the Boston the RT tweet got retweeted more than few thousand times in couple of hours.

So, studying this misinformation is one of the main focus on this course. Also we talked about cyber bullying, cyber bullying will come back again the course. We also talked about incidences like these incidences, like Egypt revolution. We will look at these kind of incidences later in the course where we can collect data for these incidences and do some analysis around it. I also mentioned about how Indian government is using social media for their interactions with citizens. A specific module in the course we will also study about how police organization are actually using social networks to interact with citizens.
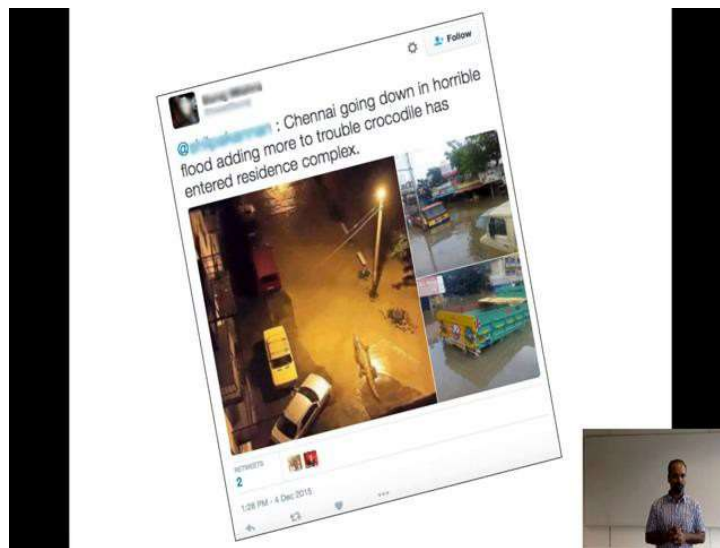
(ReferSlideTime:08:01)

Here is another problem. Now, it is not about only this fake content, it is not only about the cyber bullying content that are being posted on social networks. Here is The Associated Press, which isa verified account; verified account means it is legitimate and it is actuallyassociatedtrust,twitter actuallyverifiesthe account if youare acelebrity, if you have lot more followers and if you are actuallya marketing companywhichcanpay for it, you can actuallyget your accounts verified. Inthis case, if you seeThe Associated Press is actually posting this tweet called, 'Breaking: Two explosions in the white house and Barack Obama is injured'.

I am sure you already understand the implications of this specific tweet. This tweet is verified and therefore, people thought that it is a legitimate post, but unfortunately this account wascompromised, for a little bit ofthetimeand that iswhenthispost wasdone. Therefore, itis not only that the contents that are posted are not credible, butyou can also have these problems like compromised accounts. So, I am just enumeratingdifferent, using these examples and what I am trying to do is to enumerate the different problems that happen on social media.

(ReferSlideTime:09:18)



Also here is another one that I was also following very closely. This is Chennai floods November-December 2015 and when the floods happen there was a lot of uproar about one of these pictures, which is crocodile on the streets and there is no crocodile, there is no alligator on streets. So, this is another way by which actually suchmis-information is

beingspreadandthere have beenincidences inthe past wherehurricanesandyintheUS, where they had a picture with the ==shark in the water== and this is crocodile in the water in the street in Chennai.

So, these kind of problems which is not only the text problem, there is also with the imageproblemwhere inthe fake imagesare being postedonsocialmedia whichbecome viral which also has impact in the society.

(ReferSlideTime:10:12)



So, here is another interesting problem that happened about a year or two years back. Robin Williams, when he passed away, there was a good bye message for him, which lookedasheactuallypostedthis videoonFacebook, but thiswasactuallya fakevideo, it was not a video that was taken before his death. So, situations like these, which are floods, incidences like these to death to Robin Williams, Boston marathon blast, people use these situations to create malicious content, content that is not credible on social media. So, this is another problem which we will actually study in detail.

(ReferSlideTime:10:52)



So,this is againexample, that Italked fewsecondsbackwhich isthethree images inthis slide. The one on the left which is McDonalds the tweet, which is 'McDonalds in Virginia beach flooded, which is actually a picture of McDonalds, but it was not taken from Virginia beach, it is actually taken somewhere else in the other part of the world few years before and the image in the middle is the image that I will talk about, which looks like a shark in the water, while the hurricane sandy was going on and the third image on the right hand top is actually the image from a movie which was used to say that is how it looks now, when this hurricane is going on. So, all these are fake images, but they were all used in an incident like hurricane sandy to propagate malicious intent,to propagate these kinds of information which is not credible.

(ReferSlideTime:11:50)



So, until now we saw some events that where having talking about malicious content, how people connect with government organizations, how these information can be actually used to reconnect with parents, and things like that, but here is an examplewhich also happened couples of years, some years back which is MI6, which is military intelligence chief could not take his job because his wife actually posted some pictureson Facebook, which pictures also became public and the picture showed that this MI6 chief who was going to become a chief his being with some people whom he should not have been with.

Therefore, there is also a privacy issue which we will talk about which is how much information about people around you is getting revealed through these social networks. So, there is this whole idea of privacy the whole idea of using policing and online social networks, there is also credible information, misinformation on social network. So, this are the different topics as I said in one of my early lectures about the topics that we will cover.

(ReferSlideTime:13:07)



Now, let me talk to you about whyFacebook id knowing what throughyour Facebookis important because I can actually use the information, who you are and where you are from,thingslikethat fromyourFacebookaccountandactuallysendrelevant information to you and that is why Facebook, understanding Facebook account, understandingFacebook handle becomes such an important topic also which is personalization, whichis understanding user behavior and things like that targetted advertisement and topics around it.

(ReferSlideTime:13:43)

A few more implications about social media also, if you look at this one where there are many people around the world who lost the job because they have been using Facebook and Twitter too much.

So, therefore, this implication which is about usage of social networks itself and organizations keeping track about what you are posting and there have been incidences also where people are posted employee ofcompaniesposted about projectstheyworking on. They should not be talking about on social media and when these contents goes in hands of people whom they should not be looking and this information can actually be used against the company, the company actually takes very a strong thing around thing against the employee itself.

(ReferSlideTime:14:31)



Now,just wrappinguptheweek1content,what allwehavedone, wehavedoneactually growth of online social media, which is how large it is? What is it? How much of contents is getting generated on 60 seconds something like that and we also looked at velocity,volume,variety,valueand veracity.So,thoseare5V'sthat wetalkedaboutand in different social networks like Facebook, Twitter, Linkedin, Google plus, Whisper, Periscope, Tinder, these arethe different social networks also we saw about. And now in this lecture I talked about some incidences that where social media is used for both positive and negative implications. So, this is the topics that we covered.

(ReferSlideTime:15:23)



We have also uploaded content on setting up your machine to use Linux and python, these are two tutorials that we have uploaded for this week. So, as I was sayingin the first part of the introduction of the course we will actually have some hands-on tutorials overthesemester. So,thiswillactuallyhelp youto gethands-onand inparticularlysome of things that we are going talking early in the course will help you actually tounderstand and do things by yourself.

Aswe move forwardinthecourse, forexample, youneedto write Python codeto collect the data fromTwitter and Facebook, you will have to get yourself familiarized withthat, it is not that the difficult the tutorials about 25 to 30 minutes together Linux and python. If you can set things up if you can understand little bit about, how to set up Linux and python and set it up as earlier as possible in a machine that you would have a access to which you can use it for the entire course time.

(ReferSlideTime:16:32)



And last two slides I had was, there is Facebook page that we have, the group which is Precog where we actually talk about lot of things that we going to be discussing in the course also, which is in this page manyofthings that we talk about or share is the things activities that are around the topics that I have actually have been talking about in this course and things that the students do and then ifthere is interesting, that happens onthe topic,wewouldalso share it here. So,Iwould highlyrecommend youto lookatthispage Precog dot IIIT-D on Facebook.

(ReferSlideTime:17:10)

It is very similar to the page, this is the website that we have actually maintain precogdot iiitd dot edu, this is a website where we actually collate all the information that the group, the work ofthe group that we do together. So, I highly recommend you to look atthese two, Facebook page and this web page for any updates on this topics around if not only for this course, beyond this course also because I think it will always be good to have updated content or the latest contents on this topics.

(ReferSlideTime:17:43)



With that I will wrap this for week 1. So, what we will do is the week is, next week we willtalk little bit aboutoneofthetopicsthat I coveredverybrieflywithincidencestoday and we will also have some hands-on session tutorials also uploaded for next week. So, there will be also homework upload for week 1.

Please tryto attempt it and if there are any questions, feel free to post it on the forumon NPTEL online course website and I'll be happy to actually answer some of them. If it is something that we could actually answer, you cannot be asking what is the answer to these questions. Therefore, I will sign-off from here and see you in the next lecture.

**OSMAPIsandtoolsfordata collection**

WelcomebacktothecourseonPrivacyandSecurityinOnlineSocialMedia,week2. (Refer Slide

Time: 00:16)



I hope you are participating in the online forum that we have in the course. I already seea
lotofpeopleaskingquestions, andtryingto answer.Mysincererequest will be,please,
pleasereadthepostsbeforeyouactuallyaskthequestion;thatis,readthepoststhathave been
already asked, the questions that have already been asked and the answers that has been
already given, before asking the question. And, please participate also in the online
forum, not just only asking questions; if you know the answers for the questions that
others are asking, please try and answer them also.

(ReferSlideTime:00:52)



I hope most of you got to see the assignment 1 that we had posted. So, I think, the weekendoftheweek2isthedeadlinefortheassignment1.Pleasetrytoworkitout.The assignment 1 is actually pretty simple. We have just captured some questions from the slides that we did, and some from the tutorials.

(ReferSlideTime:01:15)



So,letmejustgiveyouaquicksummaryofwhatwehaveseenuntilnowandthen,Iwill goaheadwiththetopicsthatIwantedtocovertoday.So,first,wesawwhatsocialmedia is;differenttypesofsocialnetworks,differenttypes ofcontentthatgetsgeneratedonour

social network; classical online social media services, and then some, which are more like ephemeral social networks and anonymous social networks.

(ReferSlideTime:01:41)



We also saw what online social media means in 60 seconds; so, how much of data is getting generated on social media in 60 seconds. We saw 400 hours of videos uploaded on YouTube, and 3.3 million posts are done on Facebook and things like that. This basically shows us that, large amount of content that are getting generated on online social media services.

(ReferSlideTime:02:09)

We also saw what 4 or 5 V's of online social media are, they are volume, velocity, veracity, variety and value - those are the 5 V's of online social media.

(ReferSlideTime:02:25)



And then, I looked at, I showed you some events, where online social media has played an important role in the real world and in the society also.

(ReferSlideTime:02:37)



Telling you about different issues on online social media, for example, in this case, it is compromised account; an account was compromised, where the post said 'Two explosionsinWhiteHouseandBarrackObamainjured',and,therewas,therewasafter

<mark>effects</mark> ofthis tweet. So,welookedat different issues thatare happening ononline social media; compromised account, fake content in this case; and image of a crocodile on the streets of Chennai, while Chennai floods was going on in December 2015, caused panic among citizens.

(ReferSlideTime:02:59)



(ReferSlideTime:03:11)



And, there are also people who lose jobs and others issues, because of the usage ofonline social media.

(ReferSlideTime:03:19)



And, in the week 1, we also covered a little bit about Linux and python; hopefully, you are all set, in terms of using the platforms, because, I think, there were some questions about, 'can we use windows?'You should be able to use windows, and do programs on python, but it just said our support will be mostly on Linux. And, of course learning Linux will also be good for you.

(ReferSlideTime:03:44)



So, what I want to cover today is a couple of things; one, I wanted to actually look at differentframeworksorplatforms,thatyouwouldgettoknowwhiledoingthiscourse,

or <mark>in another</mark> terms, you should know while doing this course, and collecting data from online social media, analyzing and making inferences.

We will look at what an API is; different kinds of APIs that are available for Facebook and Twitter.Then, we will also look at programming language. There has been a tutorial on python. So, I will just quickly go over.In any case, my work for this week 2.1, about these topics, are onlygenerally,to introduceand then, wewill haveatutorials, which are specifically focused on some of them.

Then, we look at programming languages; and then, we will also look at a little bit of database, how this data is stored, what kind of format that the data is coming out; and a little bit about visualization tool.

(ReferSlideTime:04:45)



First, API, which is Application Programming Interface; this basically enables you to interact with the online social media, programmatically, and collect data from there. What does this mean? This basically means that, you can actually have a tunnel that is from your program to the social media services, to collect data. It just creates a tunnel between your program and the online social media services, where you are going to ask some data and then, the social media service is going to respond with saying, here is the data that you asked for, right.

Particularly, in our case, we will actually look at APIs for Facebook and Twitter, which will help you to collect data from Facebook and Twitter. There is other APIs also; all other social media services or majority of the social media services provide you with an API. We can't cover everything in this course. So, we are going to start looking at only the most popular ones, or the ones that we can actually use for this course, which will help you to understand howAPIs work, what data can be collected. So, you can actually do it for other social media services, yourself.

So, one of the important thing that you want to also keep in mind is that, about the rate limit, which is that in social media services when we want to collect data you cannot collect the data everything that is available on social, on these services. Because, I am sure, the companies do not want to give you all the data also. They have set it up, you know, by saying that, they have a rate limits for every social media service, and every piece of data that we want to collect from them. So, we will look at something in the tutorials about rate limits, particularly about each of the social media API , but I wanted to just give an idea about, there is going to be a rate limit, in terms of the data you can collect from these services.

(ReferSlideTime:06:33)



Also then in python, since you have already done a tutorial on python, I will keep it reallyshort.Itisbasicallyaprogramminglanguage, thatisusedtocollectdata andisone

of the popular languages currently used in terms of writing API requests to the social media services.

And, it also has a lot of libraries for reading URLs, parsing data, interact with API, and understanding the JSON objects, and things like that.

(ReferSlideTime:07:00)



Data format, so particularlytheAPI, when you send the request to Facebook saying that, 'please give me all the data about friends that PK has', or, aboutthe date of birth of PK, orabout myfriends'network. So,what it is going togive you back is actually,it is going to give you in some format. One of the formats that it gives you is actually a JSON format, which we will see in brief what this format means and how we actually interpret the data that is coming back from Facebook, or Twitter. XML, which is also a format withsomesocialmediaservices give,ortheJSON,isalsoalittlebitlikeanXML,which is Extended Mark-up Language.

(ReferSlideTime:07:49)



So, here is what a JSON means. JSON means, JavaScript Object Notation, which is a data that you get back from the social media services. So, here is an example that I have in this slide, which just shows you about the JSON object that is returned, when you are askingforidandnameofaparticularuser.So,thisistheGraphAPIExplorer,whichyou will see in the tutorial in more detail but, it is essentially a through by browser you can actually look at the data, look at the JSON objects of the Facebook data of yourself, or whatever the FacebookAPI allows, which we will be able to see through this graphAPI.

So, again, that we emphasize JSON is the JavaScript Object Notation, which is the way that the data is stored in Facebook, data is stored in twitter when you request through the API, for saying, 'give me this data about PK', it is returning the data in JSON format. It is basically the format that most social media services use today.

So, when you take the data from JSON, and when you want to interpret the data that is available in this JSON, data that is coming back from Facebook or Twitter, you can actually use JSON dot viewer dot stack dot hu. This is only for you to see visually,what data is coming back; you can take the data that is coming out of Facebook, copy paste it into this JSON viewer,and you will beable to see, whatthe fields are. When you look at the data that is coming back from Facebook, it is generally a block of data; it is just a lot of data that comes back. So, you can actually take it, and put it into the JSON viewer,to see what are the fields that it is actually giving you. Wego through this slowly,when we do the tutorials.

And, of course when you collect the data, so first is API which is a way by which you wanttocollectthedata,andthedataiscomingbackinJSON.Whenyoucollectthedata, you have to store it in some format, right. So, the format that majority of the times, the data is stored, is in MySQL. Basically, it is a relational database to store the data, and data is stored in rows and columns, and simple queries, you could use to get the data.

For example, in this case, I am just selecting user id, user screen name from the data that is being collected through Facebook, right. So, that helps meaning, again I am emphasizing that, this is not a course on MySQLitself; we will onlylook at some simple queries on how to look at the data that you have actually stored through the programsthat you have written.

(ReferSlideTime:10:28)



MongoDB is one of the popular ones, more recently we have started looking at and people are actually using this. So, MongoDB is another way by which the data is stored and the data that is collected from Facebook is actually stored.

(ReferSlideTime:10:47)



So, again, let me emphasize which is API; then, there is programming language; then, there is MySQL database or MongoDB, which is data is coming through an API, collected and dumpedinto this MySQL or MongoDB. So, now, we also need a way by whichtolookatthedatathatisbeingstored.So,oneofthewaysyoucouldusethis

actually phpMyAdmin, which actually allows you to look at the data that you have in your own database.

(ReferSlideTime:11:20)



So MySQLphpMyAdmin can look at the data from MySQL, and RoboMongo will help you to look at the data from a MongoDB. So, essentially, these are the ways by which you can collect the data, store the data and look at the data that is available with you.

(ReferSlideTime:11:38)



So, this is another view of RoboMongo, which shows you what are the different fields that are available; what data is stored in those fields.

All content on Facebook is actually stored in a graph format; that is, user - the friends that I would have, the pictures that I upload, the videos that I upload, and the status updates that I do, everything is actually a node in the graph. And, every interaction, which is basically like the comments, likes and things like that, becomes edges in this graph. Facebook actually stores all interactions, of all data that they have within the graph format; that is why the API that they have is also called as a graphAPI.

Here is the another view ofthe same message, which is,all objects are stored as nodes in the graph; connections like friends, friendships, likes are edges and all nodes have a uniquenumericID,whichisusers,pagesandposts.And,wewillbetalkingmostlyabout    users; we shall later talk about also, <mark>pages</mark>, which is one of the ways by which contentcan be generated on Facebook.

(ReferSlideTime:12:52)



In tutorials this week, you will actually look at in detail about what a Facebook API is, how do you actually create the secret key,what kind of authentication that you will have to provide Facebook, in terms of collecting data, what data can be collected and things like that.

**TrustandCredibilityonOSM**

So now, Let us look at week 2.2 of Privacy and Security in Online Social Media course on NPTEL.

(ReferSlideTime:00:12)



In this part, what I am planning to cover is actually getting deeper into the topic called Trust and Credibility.This is the slide that I will probably tweet over the course multiple times just to tell you where we are and where we are going. We finished about overview of online social media and we in the lab sessions, we have done Linux and Python and you will actually get to see a little bit about Facebook, Twitter API's and now we will actually look at the topic trust and credibility in detail, and later we cover some topics like privacy in social network analysis, e-crime, policing, using online social media and also identityresolution.

(ReferSlideTime:00:53)



Let us take a look at this graph. In this graph just to read the graph at the x-axis, this is the number of hours after the Boston blast, the data basically from the Boston blast that happened in the US. The x-axis is number of hours after the blast and the y-axis is log of tweetsbasicallywhatdoesitshow ,itshowsthatatanygivenpointintimewhichisafter the Boston blast, how much of tweets is being uploaded on Twitter.

So, there are 3 different colors in this graph which is blue, green and red. So, the one which is in the red is actually legitimate information, which you can call as the true information which is posted on Twitter. Green which is the rumor which is information that is not legitimate or untrustworthy, the non-credible content that was being postedlike the example, like the crocodile example that I mentioned in earlier lecture and the blue one which is the sum of the rumor on the true information.

It clearly shows the messages, some implications from the step one that is the true information is actually coming later; it is taking much more time than the rumors that started. In this example, there was in this event Boston blast, there was actually multiple post which were related, which to this event was not actually legitimate, for example,one postwhich said that8year old kid was actuallypartof this Boston blastwhich when therewasnokidinvolvedintheBostonblast.Therewasalsoanothertweetwhichsaid

that please RT this tweet and we will actually pay 1 dollar to Boston marathon league which also was not true.

There are many examples like this and these tweets got retweeted for more than thousands of times when Boston blast happened. This actually shows that there is multiple things one can actually look at one; how do you actually reduce the time in which the true information is coming, which is from currently it is about 9 hours or so, how can you actually get this true information come on to the social network as early as possible.

Theother solutionthat you could alsothink of this,how can you quicklyreduce thefalse information that is going on social media from, to reduce, for example in this case the green one is actually peaking in couple of hours and then its actually higher than the true information, how can you actually quickly reduce the effect of or the flow or the information propagation of this particular rumor on social networks.

So, those are the two things that you could actually do, atleast do to reduce the effect of rumorsyouneedtoactuallyunderstand,whattherumoris?Howcanweactuallyidentify    these rumors on a Twitter that is what we basically look at in the section of this course. Which is to identify ways by which I will look at the tweets and identify whether theyare legitimate or not.

(ReferSlideTime:04:04)



So now, we will look at misinformation on social media, which are some examples of the misinformation that was on Twitter. Here is one example, which actually took a lot of effect in social media. When Ebola was going on there were a lot of messages saying Ebola hoax which causes deaths and there was also discussion on the post about how salt water could be used to actually reduce Ebola and things like that.

(ReferSlideTime:04:36)

Boston marathon, hereis a specifictweet thatIjustnow mentioned, which is <mark>R.I.Pto the</mark> 8 year boy, who died in Boston explosions while running for the Sandy Hook kids and that was not true at all.

(ReferSlideTime:04:50)



<mark>There's been</mark> many, many examples I am just going to give you some examples as motivationforthissection ofthiscourse,tweetsoffalseshootoutscausepanicinMexico city, this is one of the incident.

And some tweets, some images that I actually <mark>talked</mark> about even in my first lecture,which is McDonalds in Virginia Beach flooded, the image of the left where the image was actually the real image, but it was not taken during Virginia Beach flood, but it was actually taken many years before and they associated first and the past also. Here is a rumor in the right hand bottom, which is London riots, here it reports the London zoo was broken into and large amount of animals have escaped that is again a rumor. There have been many rumors like this in many events that have happened in the past.

What we are going to do is we are going to actually take one specific example. We will actually do multiple examples over the course, over the entire course. Take up this event and look at the actuallythe topic ofmisinformation in this case and other topics in future to study how we can actually analyze this data and make some inferences out of it in the context of trust and credibility. We will do the similar way in the future also, for any topicwe'lltakeanevent,wewill  takesomedatathathasbeencollecteddosome  analysis on the content and make some inferences of the topic that we are interested in.

Inthiscasewearegoingtotake Hurricane  Sandy.Hurricane  Sandyhappened  in October  22- 31,2012andtheideaforusinganeventisthatyouwillabletorelatetoitandmostof    the    times analysis is done looking at the particular event, for example, now many people are interested in studying elections in the US and I know there are people also interested in studying elections in India when it happens. So, the damages that were totally worth for Hurricane Sandy was about 75 billion and the Hurricane Sandy basically in the north eastern part of the US.

(ReferSlideTime:06:53)



There has been many, many fake images that was floating during the Hurricane Sandy. The one McDonalds as I said before and one middle has shark in the water and people were actually,there was panic among a topic and the right hand topic which is also from Hurricane Sandyand which there was an image froma picture and it was actuallyposted on Twitter saying that is how it is looking in the US now.

(ReferSlideTime:07:24)

So, particularly for Hurricane Sandy, if you see that there is, we also know effect of these fake information that was going on Twitter. 'Hurricane Sandy brings storm of fake news and photos to New York', 'Man faces fallout for spreading a fake Sandy reports on Twitter'. These are some incidence which is happening around the event Hurricane Sandy.

(ReferSlideTime:07:46)



So, what we are going to look at is again some methodology things that I will be talking about in this course is generic. I will try to actually emphasize on this methodology. So, that we would takes this methodology and apply it in any scenario that we were interested in, sometimes even in the homework and assignments that we will do as part of this course which is first we start collecting data from Twitter about the Hurricane Sandy and then some kind of data characterization which is understanding, how much of data is come? What data is come and things like? That future generation obtaining the ground truth and then evaluating the results.

This is a very high level probably 30-40000ft high level view of what the majority of the analysis on social media data would be going on. We ourselves in the course will look at different levels of view of this slide, which is later in that course we will also look at something more detail in terms of actually this the whole process. The simple process is

collect some data and do some characterization, understand some features, use those features to create a model, use that modelto actually study the larger amount of data and evaluate the results that is the general. If you have taken any machine learning or aninformation retrieval course that is the kind of a simple process that people follow.

(ReferSlideTime:09:06)



So, the data we are talking about and this is one the most exciting thing that I feel about thestudying and researching inthearea ofonlinesocialmedia, isthesizeofthedatathat we are talking about. In this case, in the Hurricane Sandy thing we are talking about 1,782,526 tweets that were collected, while Hurricane Sandy happened. Total unique users were about 1 million users and tweets with URLs was about 622,000. So, thatgives you a sense of how much of data was collected in terms of the hurricane sandy.

So, again please keep this as a template when you are doing some kind of analysis of events. These kinds of attempt, these kinds of analysis, and these kinds of data description will help actually to look at the data to understand what the data is and in other sense it will also help for somebody who's going to do it again. These kinds of the sameanalysisorsomething similar theywould beable toactuallytakeawaysome points from the data that you would describe.

Also, in this case the map in the bottom also shows that where the tweet is come from of course, these tweets have geotagged the information therefore we are able to actually mark it on the whole map on where the tweet has come from.

(ReferSlideTime:10:31)



Of course the big question is that how will you get the ground truth because now, if we were to look at this tweets and say that which are fake, which are legitimate you need to know what the fake tweets are. So, the multiple ways that people have tried, multiple techniques that people try which is, we look at it some in the course in some probably I will just mention it as we look at the slides.

So, in this particular Hurricane Sandy analysis that was done in the way it was done was the Guardian, which is actually a media house they collected actually fake information, thenmanuallyannotatedthattheseimages arefakebecausetheyarearepositoryoflotof content that gets generated on social media, they were able to actually annotate and produce the dataset which is, which say that in Hurricane Sandy these are the fake post right. So, the reputableonlineresource tofilterfake andreal images, 'Guardiancollected and publicly distributed a list of fake and true images' what did they distribute, tweets with fake images 10,350 tweets; users with fake images 10,215; tweets with real images and users with real images.

So, using the reputable guardian data, we actually looked at the data that was collected in this and then form how many tweets have these posts, how many images that was posted that were actually fake and how many unique users and things like that, that is what that was done.

(ReferSlideTime:12:06)



So, when you do these kinds of analysis when generally look at online social media analysis, it is always best to look at analysis like this; who, when, where, what, why and how. These kinds of analysis is will actually help you to answer some interesting questions; who posted it? When did they post? Where did they post from? What did they post about? Why did they post and how did they post? So, why and how was slightly trickier here, it is hard to get, why did the person post a rumor? It is hard to tell unless the user, unless the person who posted and itself actually confesses, how do you, how did the course is probably is I mean probably possible to get which is to look at, what they why did they use how did they post, into the social media.

Now, let us look at more specifically the analysis on who, when, where and what and how and what. So,in this graph what I am showing you is a network analysis which is to show you, who is the person who posted the tweet and how the information is getting diffused. The one of the left is the user who posted content on this particular event and when the user actually post this content obviously, in Twitter there is going to be this retweets, favorites and mentions of the user. So, in this case the blue the red dot is theone of a particular user who posted this content and the blue dots are the ones where the users are actually retweeted this content. So, if we just look at this the users' content is actually spreading among the other users in the network.

So,theoneintheleftis givingyouthenthhour,theonetherightisgivingyouthenplus 1th hour what is the difference here. So, in the first one there is only one user where as content is getting spread, whereas in the one on the right if you see the post is actually diffused so heavily in the network, within one hour. But there is also other observations that you can actually have, if you look at the number in the left of the user id. So, this basicallytheuser id is the onethrough which wecan actuallycollectdatafromTwitter it is a unique for a every particular user. So, the number in the left and then one that numbers that are prominent in the right are actually going different.

What does this show? This shows that the content, that some body starts, let us take ifPK starts content and his content gets diffused in the network, he may or may not be the one who is actually more popular at after given point in time. This basically shows that the information is, we can actually draw multiple inferences from this analysis which is, who is posting the content? How the information is getting diffused, for example, if you say 3 plus see, on the left last 3 digits are 443 which is the prominent user, whereas ifyou look at in the right it is 199 the user which is in the center of the network.

So,thatisoneimportantanalysisthatyoucandoinference,thatyoucandrawfromthese kind of analysis, this is called network analysis. We later in the course will actually see some tools where you can actually draw these graphs with the data that you collect from Twitter or a facebook. So, in this graph you can also see in one of the user on the right hand top corner which also has more number retweets. There are some users which is bottom ofthe graph, onthe right which is, whoare also more the tweets are getting more retweeted.

(ReferSlideTime:15:43)



So, now let us look at a different kind of analysis from the data that we collect from Twitter,oneoftheproblemsthatwecanactuallysolvethis istoactuallyclassifywhether thepostispostthatisgiventous,butthatcomesTwitterisactuallyfakeorreal.So,in

that context, we will actually apply a technique called classification given that this is not a machinelearning or an information retrival class. Wewillnot go into detail about what the classification is? What are different techniques? I will just only look at techniques that we applied with the data that we collect from Twitter.

So, if you look at this the different kinds of features that we actually get from the post that we get from Twitter or user features and tweet features. There are actually three kinds of content that you can actually look at from Twitter, one is the user profile which is, who am I in my case the faculty at IIIT, Delhi got my PhD from ==Carnegie Mellon== university and things like that. Second, the people who I am connected with, that is my network mynetwork would befacultythatare ==around== theworld, students thatarearound the world who are doing ==cyber security==, people who are, people who studied on social media and things like that is my network.

Third is actually the content that I post itself what I am talking about I am talking about my students I am talking about ==PSOSM== course, I am talking about some random things on social media right. So, those of the three broad categories of content that you can actuallydrawfromthesocialmedia datawhichisuserprofile,thecontentthatsomebody post and the third one being the network that somebody is connected to.

So,inthiscaseweareactuallylookingattheuserfeatures whichismorelikefeel,which ismoreliketheprofiletweetfeatures.==Tweet==featuresarefromthetweetitself;letmejust go through few of them that I have listed in this slide, which is in terms of user features number of friends. Somebody has number of followers, follower-friend ratio is also one of the important things that we can actually use while making the decision on whether this user isactuallylegitimateorfake, for example, if somebodyis verypopular user,the number of followings that they would have is actually much lower that is the people that PK will follow is actually lower than the number of people who would actually follow PK.

So, that ratio can be actually used to make a judgment on whether the network, whether the user is actually legitimate or not, number of times listed list is another feature in Twitterwhereletustake,ifIwanttocreatealistofallthestudentswhoaretaking

PSOSM course I will create a list as PSOSM course on NPTEL and I will actually addall the Twitter users on to this list.

So,thatiscalledalistandparticularlythislisthasbeenusedindifferentinterestingways ifIweretofindtheexpertsaroundtheworldonparticulartopiclistcouldbeactuallyone of the good ways to find out which is if I were to create a list or somebody else creates a list on cyber security, and if they add PK onto it there is high probability that the person believes the PK is actually an expert and therefore, he is adding him or her into that list.

So, number of times listed is one feature that you can use, please go through the Twitter network, do play around with the list and other features that I am taking about, user has the url which is in my profile I will actually say that I am faculty at IIIT, Delhi and this URL called precog dot iiitd dot edu dot in, that URL is there how could the you can actually use that feature to predict user is a verified himself, verified user is another important feature that you could use because of their total number of users on Twitter there is only few a hundreds and thousands of users, who are actually verified, verified takes some process and you have to be you are to have a larger followers and things like that.

So, verified user can be a good feature to decide whether the user is legitimate or a malicioususer.Ageoftheuseraccount, andthishasbeenafeaturethatpeoplehaveused in traditional internet security methodologies, where they have actually used age of aweb site, age of the domain registration to actually find out whether the domain is legitimate or not. It is same feature which is PK created an account 5 years back it is more legitimate and there is a PK account which means or there is Amitabh Bachchan account or Rajinikanth account which is created recently which may not be actually legitimateaccount.So,thatistheintuitionbehindusingageofauseraccount.So,nextis tweet features let us just look at little bit about tweet features itself, in tweet features lengthofthetweetisagoodinformationthatyoucanactuallyusetofindoutwhetherthe post is legitimate or not.

So, these are the features that you can use in general from many different analysis I am onlyusingitfortheproblemoftrustandcredibilitythatwearetalkingabout.Also

length of the tweet number of words in the tweet, contains question mark, contains exclamation marks, number of question marks, number of exclamation marks, contains happy emoticon, contains sad emoticon, and things like that. So, this these are the different features that you can actually draw from tweets and the features that I told earlier which were actually user features. So, five fold cross validation is a technique which is used to make sure that the confidence on the classification accuracy that we are building is higher that is the reason why we use actually five fold cross validation and there are many other techniques, I am not going to into details of different other techniques that are available.

(ReferSlideTime:21:22)



Now, let us look at the results from the classification that we did. So, as I saidF1 is a user feature F2 is a tweet features, and in classification techniques we can actually useF1 F2 separately and also create a features of F1 plus F2 the two techniques that were applied one is Naive bayes, which I actually usesbayes theorem to find out whether particular post is legitimate or fake and which features actually influence a lot in making the decision. That is also another technique which is a graph based technique which is decision tree, where all the outcomes and the probabilities are actually layed down on as in the form of a graph and it is a very popular machine learning technique which is appliedtomakedecisions.Andthisparticularcasedecisiontreeactuallyseemtoawork

better the people while using the tweet features, while the efficiency was about 97.65 percent, in predicting whether the post is fake or real. The tweet features I have choosed seems to a you have played well in both a naive bayes and decision tree where as use the features did not that play that much well in making the decision.

(ReferSlideTime:22:32)



Now, let us look at an event Boston blast again, we will use this technique could be looking at events through the events I will actually inject a lot of techniques that we will actually study in this class, and terminologies also that we will see. It is a twin blast that happened in 2013 in April and 3 people were killed and 264 were injured, first imagethat come on Twitter was within 4 minutes. It is basically a Boston marathon that was going on and the blast of the finish line was the event that happened.

(ReferSlideTime:23:01)



And there were actually multiple fake tweets here I am just showing you two popular fake tweets that were actually floating around, the first one I have showed this tweet in thepastalso 'R.I.P.to the 8year old boywho died in Boston'sexplosions,whilerunning for the Sandy Hook kids'. There was no kid who has participated in the marathon and then other post the at the bottom you see, for every retweet we will donate one dollar to the Boston marathon victims, and it is posted by an account called underscore Boston marathon something that you want to keep in mind which was not a legitimate account and this post was retweeted for about 50,000 times and these are the two popular tweets that were floating around during the event which were fake

(ReferSlideTime:23:47)



Data Description

| | |
|---|---|
| Total tweets | 7,888,374 |
| Total users | 3,677,531 |
| Tweets with URLs | 3,420,228 |
| Tweets with Geo-tag | 62,629 |
| Retweets | 4,464,201 |
| Replies | 260,627 |
| Time of the blast | Mon Apr 15 18:50 2013 |
| Time of first tweet | Mon Apr 15 18:53 2013 |
| Time of first image | Mon Apr 15 18:54 2013 |
| Time of last tweet | Thu Apr 25 01:23 2013 |

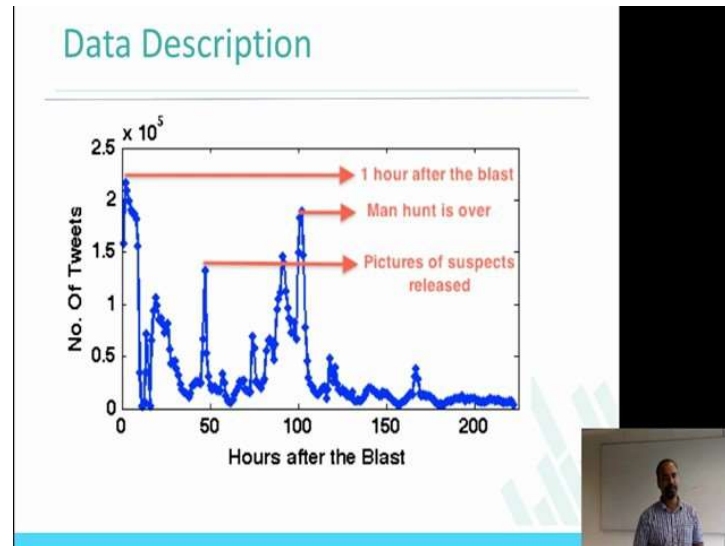Datathatwascollectedduringthiseventswasactuallyabout 7.8milliontweetsthatwere collected, 3.6 million users posted this tweet and if you look at the advantage of actually working in this space of online social media is actually this large numbers that we look at, tweets with URLs is about 3.4 million, 62,000 people are posts had geo tag and about 1 percent is what Twitterclaims that the tweets that are posted on Twitter are geo tagged tweets, about 4 point 4 million replies 260,000 in the timeline of the blast. First tweet, first image, and last tweet, all of that is capture in this slide I will show it you because when you actually present and an analyze events. analyze a particular topic and I think studying this analyzing unit is only one way, but we are actually, you can actually adapt this to studying any topic.

For example in this case how do we collect the data we take hashtag Boston marathon which is actually trending and start collecting tweets, which has hashtag Boston marathon, look at other words that are in the post that has hashtag Boston marathon and use those key words to start collecting other tweets like query expansion concept and thereby we collect that post from Twitter. And this methodology can be used for collecting any data, data could be hashtag macbook pro, hashtag apple hashtag india and things like that. Not necessarily it has to be hashtag also, it could be any other words.

HereIamactuallyshowingyou a graphwhichis onthe x-axis is the hours afterthe blast in the y axis is number of tweets. So, the the crux of this slide is to show you that the spikes that happen on social media is actually very very connected or correlated to the events that happened in real world, for example, here are the first spike that happened is one hour after the blast then there was a spike in Twitter tweets which is pictures ofsuspects released and man hunt is over.

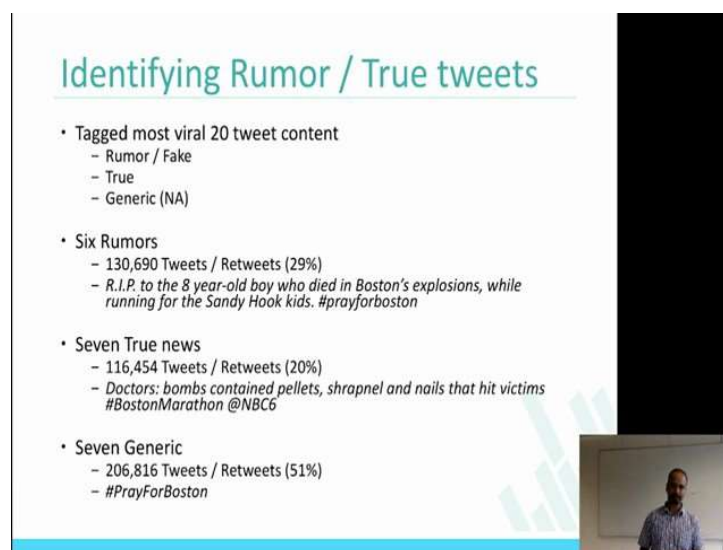So, if you reallylook at it, that is the waythat the content gets generated on social media isactuallybehavingand we've tried lookingatthesekindofblastformany,manyevents and itlooksverysimilar,interms ofactuallyspiking correlated totherealworldincident or an event that happens. And here is another slide which actually shows you the geo tagged coordinates of the tweets that were posted on Twitter.

(ReferSlideTime:26:16)



Particularly for this Boston marathon blast, and every dot in this slide show you the tweets that have come from that particular location, understandably the <mark>posts have</mark> come mostly from the US.

(ReferSlideTime:26:39)



So,nowinthiseventonBostonmarathon<mark>thetechniquethat</mark>wasappliedforfindingout

whether a post is actually rumor, true or fake, first tagged most viral 20 tweet content which is whether it is rumor fake true or generic rumors, 6 rumors were actuallycollected from the posts that were talking about Boston marathon and seven true news was collected, which is doctored bomb contained pellets, shrapnel and nails that hit victims Boston marathon hashtag NBC6. So, those kinds of tweets were collected which is about true news that that was getting posted during the Boston marathon and six rumors werecollected and seven generic posts that had prayfor that the Boston. Prayfor Boston also was actually trending during that time.

So, essentially what we were trying to do is we were trying to study, look at the rumors thatwereposted,truenewsthatwerecoming andsome genericinformation,genericpost that happened during in the Boston marathon event. In this kind of you see a generic sense of what are the different post that happened in an event like this and rumors about 29 percent were actually retweeted, true news about 20 percent were retweeted and 51 percent of the generic content was actually fit for retweeted.

(ReferSlideTime:28:06)



Here is another view of the data which is to show you the fake content user profiles. So, every time such event happens incidentally what happens is many of the fake user profilesgetgeneratedduringtheevent,fakeaccountsgetsgeneratedtoactuallyusethe
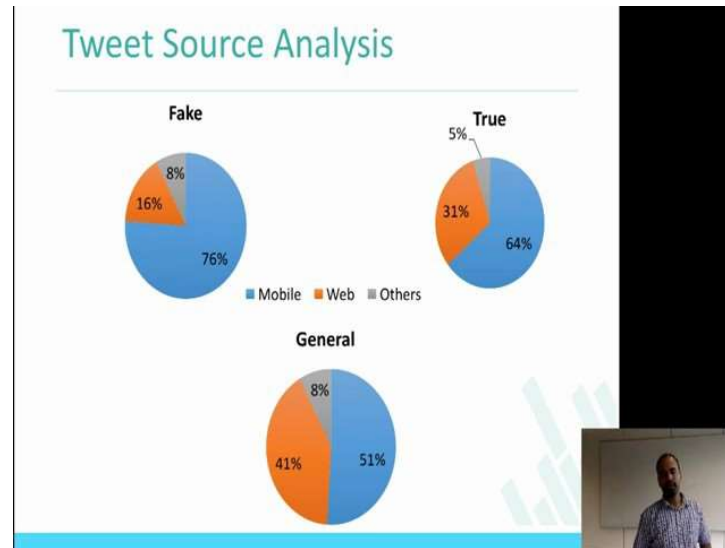
eventtopropagatemaliciouscontent.

So, in this example I am showing you account 1, account 2, account 3, account 4 different accounts were actually created, you will actually see that the accounts has number of followers which is account 4 being very high, account 1 being very low. So, this is and if you see the account 1 which is actually created on March 24, 2013 which waspostingfakecontent,account2whichwascreatedonOctober15,2013verycloseto the Boston marathon and account 3 February 2013 and account 4 2008. Total number of status this is the updates that they created and number of fake tweets that they posted is2, 2, 1 and 1 respectively.

And if you seesome of theaccounts where getactuallysuspended and these suspensions happen because people report about this handle to Twitter in a multiple ways to actually keep down a particular account, while one is large number of people actually report a particular handle to Twitter and it gets suspended and there are through government processes you can actually apply for suspending an account.

Some if you see the last column, it is interesting that some of the user handles which are posting fake content on events like Boston marathon actually are active, even when we were actually collecting the data. This shows you that fake content propagated by fake user handles and these user handles created just after the event or just before the event, just after the event happens.

(ReferSlideTime:30:04)



Now, let us look at different view of the analysis which is tweet source analysis thisgives you an insight about what devices do people use while posting the content, 76 percentofthe postthat was identifiedas fake was posted throughmobile,whereas the 64 in true and 51 in general content that were posted. This insight about what device isbeing used, while posting this content can be very useful in making decisions, for example, if you wanted totargeted advertisement, what kind of devices are being used can be very useful in making the decision. So, the device this information is available in theJSON thatyou collectfromTwitter foreverytweet. So, you can usethatto makethis judgement.
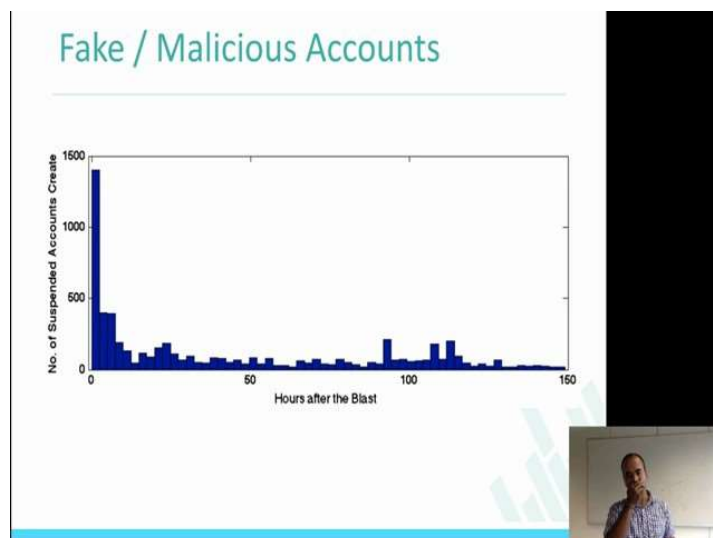
(ReferSlideTime:30:47)



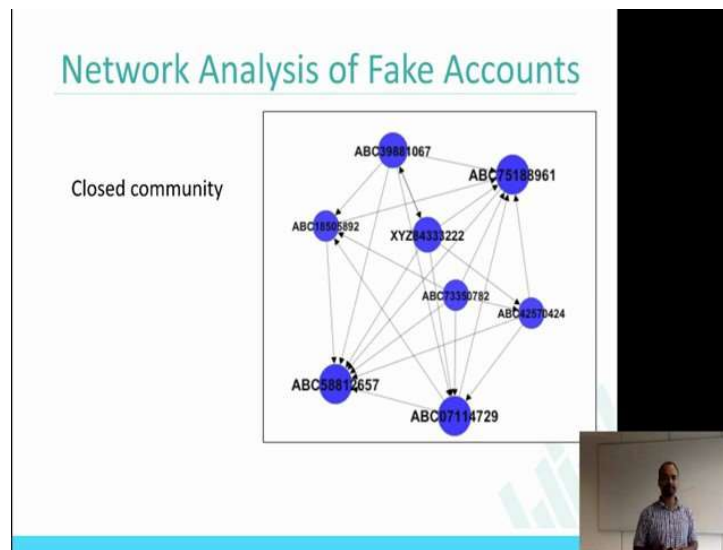So, if you look at the number of accounts that were created during this event, it wasabout 32,000 new Twitter accounts were created during this event, which were actually talking about this particular event. Out of this 19 percent were deleted or suspended by Twitter which again could have happen for multiple reasons and 19 percent of accounts that were created were actually suspended.

(ReferSlideTime:31:12)

So, this is graph to show you hours after the blast; x-axis being the hours after the blast and y-axis being number of suspended accounts that were created, which is to just to show you that the number of accounts that gets created immediately after the account is also high, in addition to that number of accounts getting deleted also high the fake or maliciousaccounts thatweresuspended,thatwerecreated andsuspended wereveryhigh immediately after the event and after the event it kind of reduces little bit.
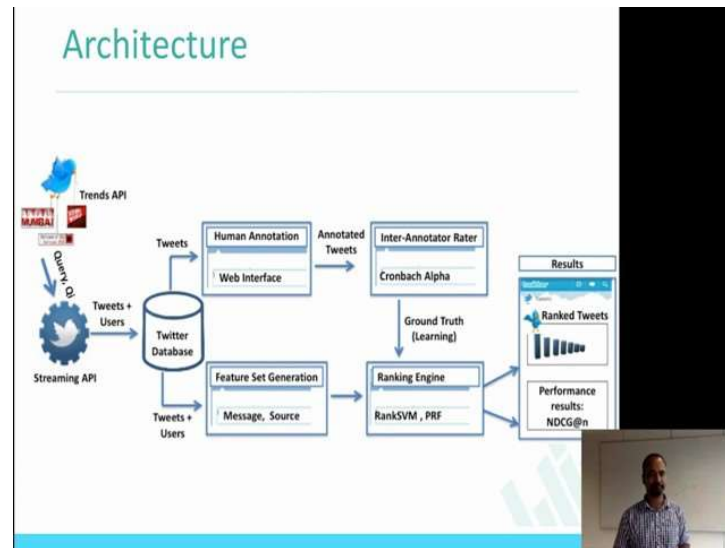
(ReferSlideTime:31:45)



Let us look at the connection between the fake accounts itself. So, again just keep in mind the kinds of analysis that were doing is who, when, where and what, whyand how. So, again that one insight into the analysis is this how were the people who were posting this content which is fake are connected. So, one insight is that they are actually pretty closed and this is not only in this domain you can actually seen this kind of analysis in many other domains also, for example, in classical security problem like phishing.

The number of groups, number of accounts, number of sets of people who do this actually by it is small and they are all very well connected, similar kind of inferences is derived from this particular analysis also where for the Boston marathon if you look at people,thenodehereis theuserandtheedgebetweenthenodesaretheactionofretweet followingandfollowers.So,therefore,thereisathereisaclosedcommunitythatis

actuallyoperatingin termsofpostingthisfake content.

(ReferSlideTime:32:51)



So, while doing this you could actually think about. So, earlier I think in week one I actually showed you some very high level slide about how machine learning and about how these kind of approaches of identifying fake and legitimate can be done, this is just to slightly zoomed in view of the same slide, which is to show you that the data is comingfromTwitter throughstreamingAPI's.Assumeyouknow,everybodyknowsnow what is streaming API is, which is to collect data from the social networks, tweets are dumped to the database, human annotations which we saw earlier also in terms of annotating the post even now we saw about fake, generic and true our post those are all sometimes human annotated, sometimes you could actually use some simple techniques to do the annotation, one of the important thing that you also want to do in this annotations are done are inter annotator agreement which is if I say something that is legitimate and if you say something that is legitimate then probably more people saying post is legitimate, then the post must be legitimate, that is the kind of intuition that the Cronbach's Alpha, which is value that you may get for finding out inter annotator agreement and at the Cronbach's Alpha is generally about 0.7.

Itis actuallyunderstoodthat the data has youcan have more confidence intheinferences

that you are drawing from this particular data. Cronbach's Alpha is the value that you will calculate while finding out inter annotator agreement and so as we discussed earlier also feature extraction, feature extraction is a technique by which we took F1, F2 and those things you will use that to find the model here.

I will just describe a little bit in the later slides about what model can be generated and you use that to find out whether this particular post is legitimate or not and then you can actually show that to the user also, that is the architecture that is presented here again is a very simple machine learning approach which is take the posts, use the post, do some feature extraction, use the feature extraction to create a model, use the model to actually predict whether the post is legitimate or not.

(ReferSlideTime:35:14)



## Data Statistics

| Events | Tweets | Trending Topics |
|---|---|---|
| UK Riots | 542,685 | #ukriots, #londonri- ots, #prayforlondon |
| Libya Crisis | 389,506 | libya, tripoli |
| Earthquake in Virginia | 277,604 | #earthquake, Earth- quake in SF |
| JanLokPal Bill Agitation | 182,692 | Anna Hazare, #jan- lokpal, #anna |
| Apple CEO Steve Jobs resigns | 158,816 | Steve Jobs, Tim Cook, Apple CEO |
| US Downgrading | 148,047 | S&P, AAA to AA |
| Hurricane Irene | 90,237 | Hurricane Irene, Tropical Storm Irene |
| Google acquires Motorola Mobility | 68,527 | Google, Motorola Mobility |
| News of the World Scandal | 67,602 | Rupert Murdoch, #murdoch |
| Abercrombie & Fitch stocks drop | 54,763 | Abercrombie & Fitch, A&F |
| Muppets Bert and Ernie were gay | 52,401 | Bert and Ernie |
| Indiana State Fair Tragedy | 49,924 | Indiana State Fair |
| Mumbai Blast, 2011 | 32,156 | #mumbaiblast, Dadar, #needhelp |
| New Facebook Messenger | 28,206 | Facebook Messenger |

So, in using this architecture just taking, instead of just doing one or two events, multiple events data were collected and used to find out whether a particular technique, technology can be identified where this post is legitimate or not and here are the events UK riots, Libya crisis, earthquake in Virginia and US downgrading there are many, many events data were collected and in as I said before the column in the third column here which is trending topics.

These were the topics that were trending using which the data was collected; the column 2 shows you the number of tweets, again alarge number of datawas used in while doing this analysis.

(ReferSlideTime:35:58)
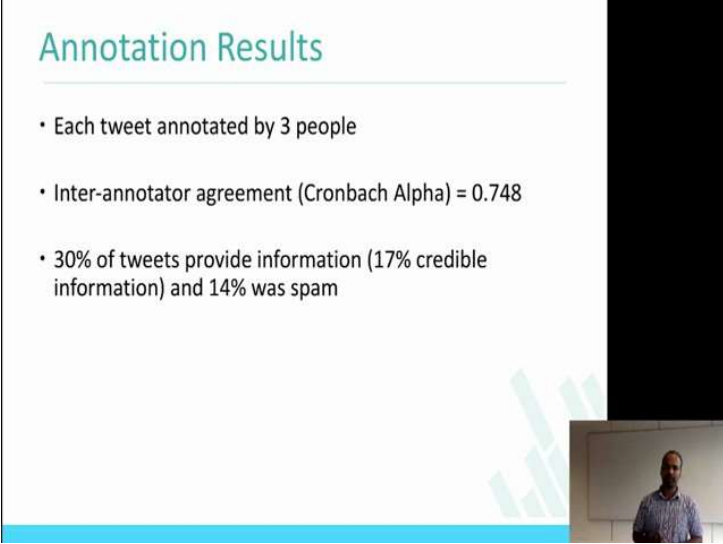


As discussed before, one of the methods used to find out whether particular post is legitimateornot,annotationwasdonetherearemultipleways todothisannotation.Also you could get 3 of your friends to sit down together and I will tell you whether this every post is legitimate or fake, you can do through mechanical turk, mechanical turk is a crowd source mechanism by which you could actually show these posts in theplatform of mechanical turk. The turkers also to look at it, turkers are basically people all around the world who are doing this task for a small money, step one in this case contains information about the tweet, postis shown to the user and in the user actually decides on one of these four characteristics which is contains information,is related to the event not only related to the events, skip.

Ifin thesteponeandsays thatcontainstheuserdecides thatthereisainformationinthis post,thentheuserisaskedaboutdefinitelycredible,seemscredible,definitelyincredible and skip tweet, again here, I am only going through the methodology which is post is taken.Itisannotatedyoucouldannotateitforanyparticulartopicthatyouwouldwhere

you want to study, in this case it is credibility, but you could also think about it whether this post has phising URL or not, if this post is talking about a particularly event or not, this post is sensitive or not you could do many, many things in terms of annotations and in the topics that you are interested in studying from the post that is being collected.

So, from step one you take the data and then you ask the users to classify, it as definitely credible, seems credible, definitely credible and if there is nothing the user cannot makea decision, skip the tweet.

(ReferSlideTime:37:41)



AndthatiswhyIsaidabout Cronbach'sAlphawhichissomethingIwillemphasizehere, each tweet should be annotated by at least three people because that will give you more confidence in the data and when you do this it is called inter annotator agreement or Cronbach's Alpha. If you calculate that and if it is more than about 0.7 it is generally accepted that the data has more value or confidence in it, 30 percent of the tweetsprovide information which is in the step one users agree that 30 percent of the tweets are shown to them add information, only 17 percent have credible information and 14 percent was spam.

(ReferSlideTime:38:21)



So, feature sets, we now discussed just now some time back about different features F1 and F2 in that slide here message based features and source base features which is again, if you lookatitItold you aboutfeaturesfromthepostswhich is messagebased features. So, features from the profile which is a source based features.

(ReferSlideTime:38:46)

Using these features we used a metric called the NDCG, which is normalized discounted cumulative gain it is nothing, but way by which you can actually mention the efficiency of the search, NDCG is being commonly used in finding, how good a search engine is performing in this case, we are using it to find out how what is the quality of the classification that we make whether it is legitimate or fake in this metric called NDCG.

(ReferSlideTime:39:20)



Also here is a graph for looking at the content from the tweets that we collected looking at the post it that we collected, recency, tweet, user, Twitter plus user, these are the features that we used. If you remember to find out whether a post is legitimate or not we are here we are drawinggraph of n,which is on the x-axis and NDCG value, which is on the y-axis, you can clearly see that the tweet plus user which is at the top of the graph doing well in terms of the NDCG values. This basically helps to understand that what you post and who you are are a good features to make judgment on whether the contentis legitimate or not, that is the kind of inferences that you should be chasing while you are analyzing the content from social media which helps in some actionable information also.

For example, here what you post and who you are helps to find out whether the post is credible or not which helps in making lot of decisions.

(ReferSlideTime:40:34)



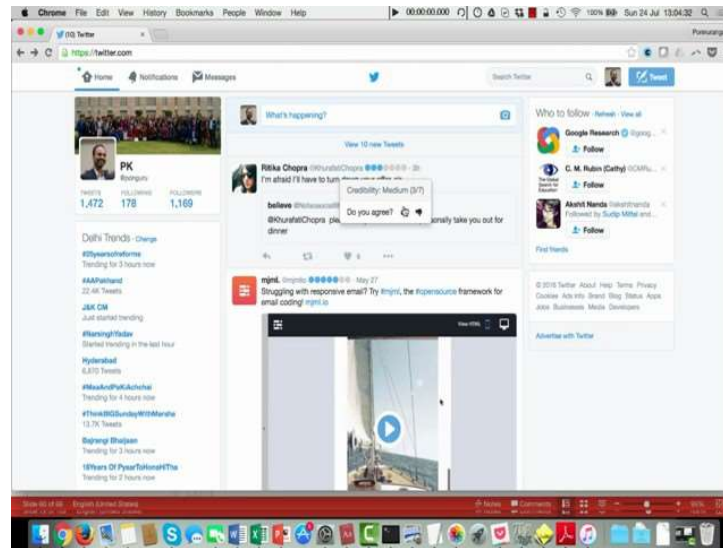Using this understanding of what feature works and what features do not work and what feature actually influences in finding out whether the post is credible or not, the TweetCred, a chrome browser extension was built and this extension it helps you to identify whether this particular post is credible or not. I'll just show you a light demo of the TweetCred extension and then I will walk you through, what this available in the chrome extension.
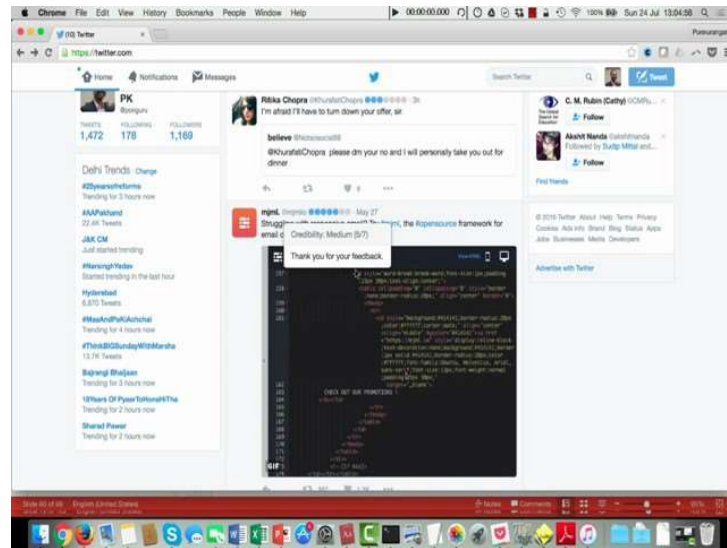
What I am showing you here is the Tweetcred browser plugin, chrome browser plugin which actually helps you in making a decision whether a particular tweet is credible or not. It just gives you; it basically uses all the features that we discussed until now,where it is being bundled with this chrome extension.

So, look at this tweet, if you go to Twitter dot com in your timeline, this information about whether a post is how credible is it will not be there this is coming from TweetCred. If you look at this it actually gives you a value of 3 on 7, it is calculating the value of credibility on a scale of 1 to 7 and in this case it is showing that this is my timeline in this case it is showing that this post is 3 on 7, this post is 5 on 7 and values like that. So, it is going to work for all the post that are in a timeline, it is going to work for what in your search its going to work for post in dm and things like that.
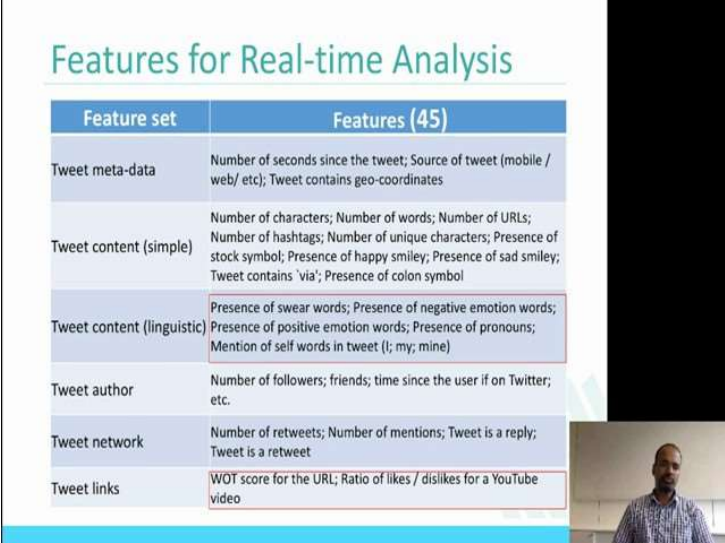
So, let us look at the values that it is presenting also. So, if you see here it is actually showing you a value of 3 on 7 and then when I hover it. It actuallygives me information called credibility medium 3 on 7, do youagree? So, this is the way if you remember,the machinelearningmodelthatpeopleusingthefeatureyoutakethatmoduleandwhenever we get feedback like this, we can actually go and update that model to make it more efficient.

(ReferSlideTime:42:55)



So, inthiscase you could actuallysaythat no,I actuallyagree with the value of3 on7, it gives you message saying thank you for the feedback. In this case, let us take it if I were to say that the value I do not agree with value then it actually asking what you are agreeing with. So, I say no this is actually more credible it should be actually 7, when it says thank you for feedback. So, what it basically does is, it is capturing these details from you and it is going to make use of it when we end up updating the model that was built at the back end for the TweetCred. This information can be used in making the judgment. So, that is the chrome browser extension of TweetCred, which basically takes the features in real time and makes the judgments and presents it to the user with the values of 1 to 7.
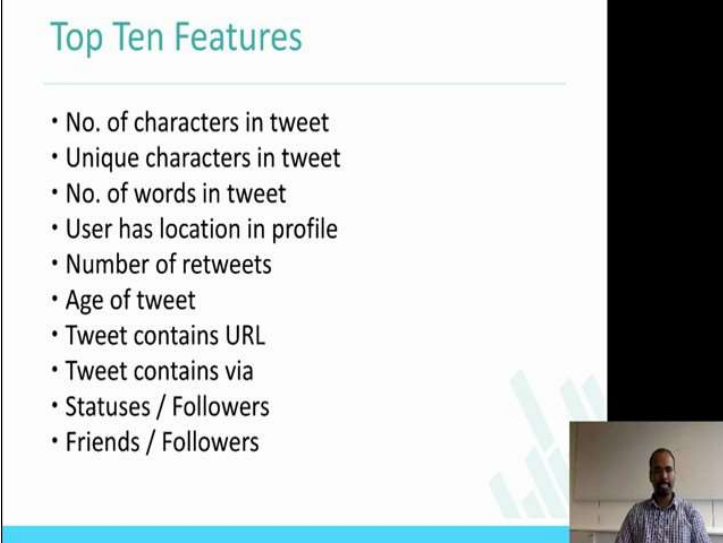
(ReferSlideTime:43:27)



## Features for Real-time Analysis

| Feature set | Features (45) |
|---|---|
| Tweet meta-data | Number of seconds since the tweet; Source of tweet (mobile / web/ etc); Tweet contains geo-coordinates |
| Tweet content (simple) | Number of characters; Number of words; Number of URLs; Number of hashtags; Number of unique characters; Presence of stock symbol; Presence of happy smiley; Presence of sad smiley; Tweet contains 'via'; Presence of colon symbol |
| Tweet content (linguistic) | Presence of swear words; Presence of negative emotion words; Presence of positive emotion words; Presence of pronouns; Mention of self words in tweet (I; my; mine) |
| Tweet author | Number of followers; friends; time since the user if on Twitter; etc. |
| Tweet network | Number of retweets; Number of mentions; Tweet is a reply; Tweet is a retweet |
| Tweet links | WOT score for the URL; Ratio of likes / dislikes for a YouTube video |

So, you may remember the features that we discussed in this lecture, but unfortunatelyall features cannot be actually used while doing it in real time, for example, finding out all the followers that you have and using some scores on the followers is actually hard. So, here are the 45 features that were actually used to while doing the real time analysis itself, specially I wanted to highlight on the presence of swear words, presence of negative emotion words, presence of positive emotion words, web of trust score, whichis WOT scorefor the URL and ratio of likes and dislikes fromthe YouTubevideo which has links to the YouTube. So, these are the features that we did not discuss before. So, I kind of thought we'll highlight them when I'm presenting the slide. So, these features were used in building tweetcred demo that I showed you.

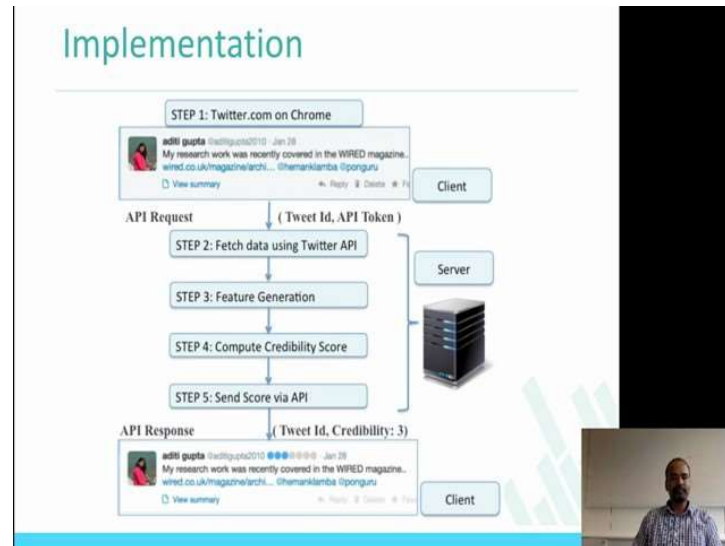Of course, the common question is what are the top 10 features that actually makes the decision or which influences in identifying whether post is a credible or not more efficiently.Itisnumberofcharactersinthetweet,uniquecharactersinthetweet,number of words in the tweet, user has location in the profile, number of retweets, age of tweet, tweet contains URL, tweet contains via which is through, how the post was done, status and followers, friends and followers, thoseare thetop 10 features of fromTwitter, which can be actually used to make a judgment on whether a post is legitimate or not. Please keepinmindthisisonlyforTwitter,thefeaturesthatyoulookforFacebook,thefeatures that you may look for Instagram, in other social network may be very different.

(ReferSlideTime:45:16)



Here is just a slide to show you how the implementation for the tweetcred was done which is chrome browser extension. There is a post that is on your timeline, it takes the post fetches data using Twitter API, which is the architecture that I showed you earlier, where feature extractions were done. Then, the model which is built, tweet is taken through API feature extracted the credibility score is computed with the techniques that we discussed until now and the values assigned back to the API, and then tweet ID and the credibility value comes back, it is presented in your timeline saying this value is actually 3 on 7, the demo I showed you. It is simple chrome extension that was about to show these values.

(ReferSlideTime:46:05)



So, users can also give feedback to the system and that is showed in your demo TweetCred actually ask user to say agree or disagree with the values that is presented. If you agree that is okay, if you don't agree please provide the information, please providea value that it should be what you think it should be, that is what is presented in the top left which is when you agree , bottom right is actually saying if you disagree.

(ReferSlideTime:46:35)

Different types of users that can be <mark>foreseen</mark> using TweetCred type of tools. You can at leastremember TweetCredisonlyoneexamplethatIampresentinghere,therearemany other tools that one to think of while analyzing social media content and information presented to the user. In this case emergency responders, fire fighters, journalists and news media and general users also have started using tools like this.

(ReferSlideTime:47:01)



Let us do a quick summary of week 2, when we started we actually looked at API's, programming interfaces and we you also have a tutorial for Facebook in this week. Then we looked at very, very briefly what Python programming language, MySQL, Mongo DB, PhpMyAdmin and when you collect the data you are going to actually get the rate limits. Please remember that there is always going to be rate limits when you are collecting the data from these social media services and we talked a little bit about the format in which the social media service is <mark>store</mark> the data, which is JSON, when you collect the data, you are going to get JSONs which you have to analyze through your scripts and Facebook stores all the data in terms of a graph. We looked at that brieflythen westarted <mark>digging</mark> deeper intotrustincredibilityasafocusarea.Welookedatthese concepts of trust and credibility through events being Boston marathon was one of the events, Hurricane Sandy is another event that we looked at we looked at these events. ThroughtheseeventsthatIwastryingtotellyouhowdataisbeinganalyzed,whatkind
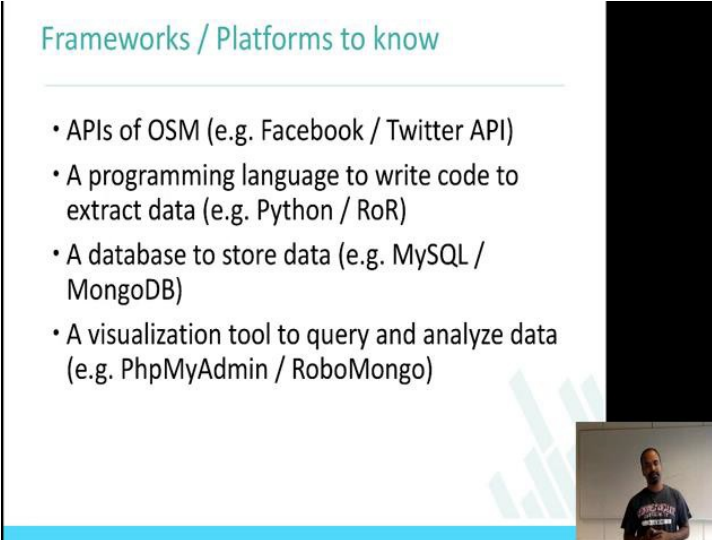
oftechniquesare beingapplied onthisdata.

We looked at classification is one of the major technique that is used while designing whether a post is credible or not and during this analysis, I also told you about who, when, where and what, why and how are the basic questions that you can actually analyze using in the social media content and specifically we have also looked at some social network analysis techniques inputs. So, that is the week 2, hopefully you will go through the content and if you have any questions please go and ask in the forum, we'll be happy to actually answer there.

# Unit-2

## MisinformationonSocialMedia

Welcomeback to thecoursePrivacyandSecurityin OnlineSocialMedia, this isweek 3. Let me put you go over what we covered in week 2.
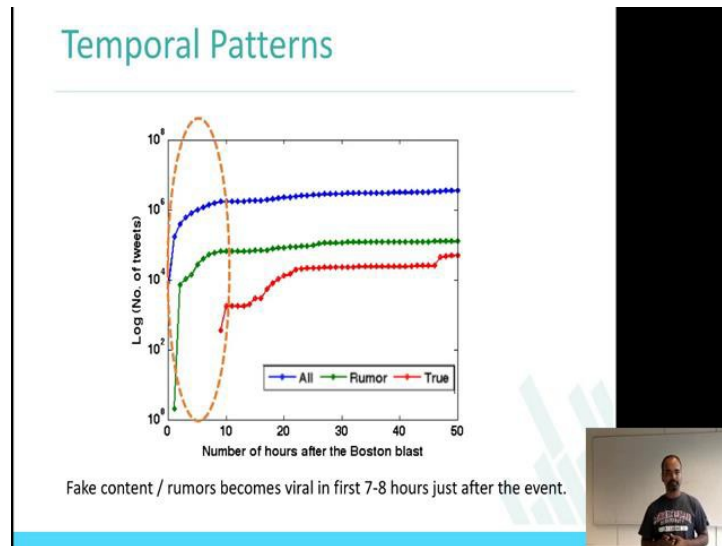
(ReferSlideTime:00:20)



We started looking at what an API is, and then we looked at what Python Programming Languages, MongoDB, how the data is stored, how we can actually visualize the data using PhpMyAdmin or RoboMongo. I hope all of you have already done hands on exercises and practices with Facebook API and twitter API.

(ReferSlideTime:00:45)



After that we looked at the topic Trust and Credibility, in that we looked at multiple events; through the events we looked at some concepts. Here is a slide that I used inweek 2, where we showed that the truthful information is coming into the social media slower than the rumours, and there are multiple techniques by which you can actually attack this problem.
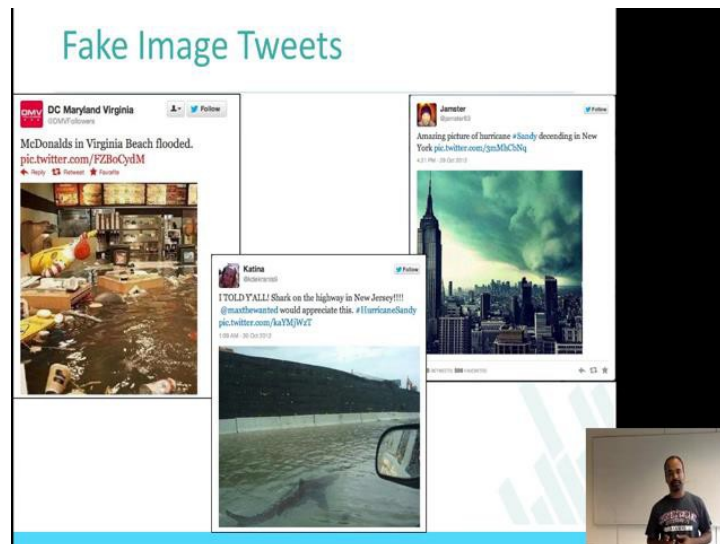
(ReferSlideTime:01:10)

And I also showed you some examples about Misinformation; tweets that were being posted on social media and there have been multiple effects of it; fake content getting viral, some values on the stock market getting affected and of course rumours also.

(ReferSlideTime:01:30)



More examples that I showed you in week 2, showing that there is a shock in the hurricane sandywhere this picture got viral and it had effects on the public also.

(ReferSlideTime:01:45)

Most specifically I was trying to give you an intuition about how, what analysis can be done. Particularly; who, when, where, what, why, and how. These are the kind of analysis that you should be interested in doing while looking at the social media content.

(ReferSlideTime:02:04)



And then, we later looked at different features that are available in tweets particularly user features and that tweet features and we tell detail about what these features mean, how these features can be put together to create a classifier which can look at tweet and then say that whether it is legitimate or fake tweet.

(ReferSlideTime:02:24)

Some more examples, particularlythis is fromthe BostonMarathonwhere I showed you that a tweet which said, RIP to the 8 year-old who died in Boston explosion was retweeted more than 30000 times and malicious user used this spread or occurence ofan events to actually spread the content and get victims to go to malicious URL's which share malicious information.

(ReferSlideTime:02:51)



This is one slide when I talked about, how the data can be represented, what data has been collected for doing these kinds of analysis.

(ReferSlideTime:03:02)

I also mentioned about the spikes in the social media data content that is generated on social media is very, very populated with the actually event that happens in the real world. This is one example where man hunt is over, a lot of people are talking about it and therefore there is spike in the tweets that are showing up.

(ReferSlideTime:03:23)



This is Geo-Located tweets where each dotis a tweet which has a geo tag attached withit and such kind ofgraphs can be helpful in saying where these tweets are coming from.

(ReferSlideTime:03:36)

Wealso talkedabout howthecommunityofusers who arepostingthis fakecontent, who created fake accounts, how they are connected. Interestingly, they are all connected very closely and it is a closed community.

(ReferSlideTime:03:53)



And I also walked you through a multiple architecture diagrams mentioning about how data is collected from social media, what kind of annotations and how do you actually verify with the data that you annotated is actually of high quality and, what kind of feature generations can be done, what is the model that we developed and what is the model that one can develop, and now what are the evaluation matrix to actually find out whether the technique that we have applied and the model that we have created isactually good.

So, this is an architecture that I tell in detail talking about each block and explaining all this block helps in creating some interesting solutions for the problems in the trust and credibility space.

(ReferSlideTime:04:39)



And then I showed you about plugin which is called TweetCred, I hope some of youhave played around with the tweetcred plugin to find out how the tweets are evaluated and the value of x on 7 is presented to the users.

(ReferSlideTime:04:58)



So, now what I want to actually cover is little bit about how one can actually take this understanding of twitter and then apply into other social networks, because this is a privacy and security in online social media course so I thought that it will be interesting tofind outhowthesekindoftechniquesthatwelearntfromtwitter canbeacquiringinto

other social networks particularly I will talk about Facebook. If you think about it initially I talked about how Facebook and twitter are different in my week 1 lecture, where we said that Facebook is a bi-directional network and twitter is a unidirectional network and the structure itself is very different.

And, the features that are available in these two social networks to study are also different. In twitter it is followers and following and Facebook, it is actually friends and the informationthat these networks provide through API are also verydifferent. And the structureofthe networks have different, particularlyI wanted to highlight this friendship thing in Facebook;theconnections are morepersonal and ifthere is apost that showsup by your friend, there is some tendency that it is more likely to be truthful and then weyou believe that your friend's post is actually more truthful than a random person's post.

So, that is the one of the differences between the Facebook and twitter network, particularly keeping this trust and credibility as the space of discussion. I wanted to highlight this difference, and now given this difference we should also look at how we can actually use the modelthat we have understood in twitter to apply it into Facebook.
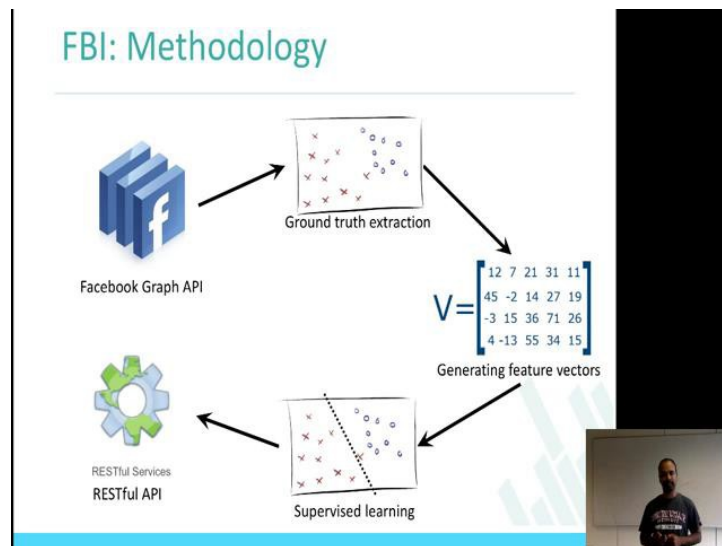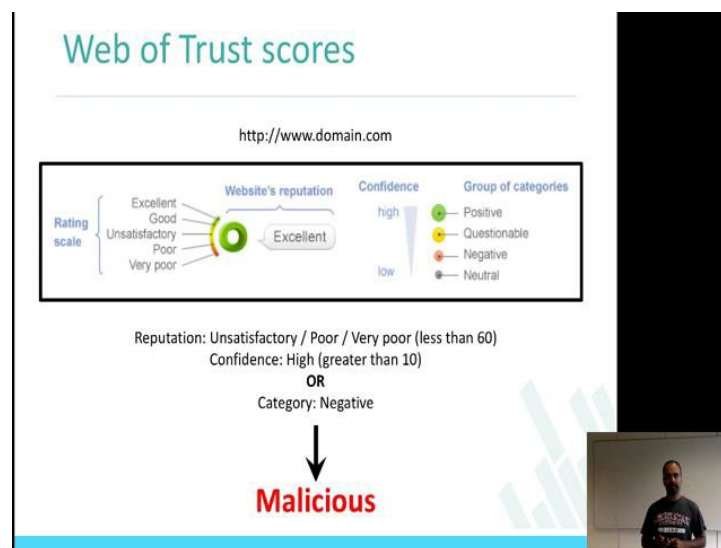
(ReferSlideTime:06:37)



The architecture if you see it is almost the same, in this case it is just presented slightly differently. FBI: stands for Facebook Inspector, a similar tool that is like tweetcredwhich takesthe Facebook post fromFacebook graph API and then looksat the postsand

make some judgement on how whether these post are malicious or not, credible or not, trust worthy are not.

In this case itis the same architecture which takes the post, do some feature extraction, do some ground truth understanding ofthe post, then creates some feature vectors out of it, create a model out of it, in this case supervised learning model because we actually have data from the posts that we are collecting and then create a RESTful API through which you can actually find out whether this post is malicious or not. Same architecture, very similar to tweetcred so I do not think we should actually spend a lot of time in understanding more details of this. If there is any question please feel free to ask in theforum for sure.

(ReferSlideTime:07:43)



So, one thing that was also mentioned in the tweetcred or the twitter trust and credibility slides is a Web of Trust that is called WOT. Then I thought I would just mention it briefly what does it mean. It basicallytakes adomain and producesanoutput which says that a value similar to tweetcred, similar to other services that you may have seen where input is a domain and the output is score, which you can use to say that whether it is a malicious domain or not. Then in the past also I mentioned about how long the domain has been registered, who registered the domain and things like that. These features canbe used to make the judgement.

So web of trust basically gives you value of excellent, good, satisfactory, poor and very poor. If you give domain saying iiit dot edu dot in, it will actually come back with the rating scale and a confidence scale. We use this in Facebook inspector because in Facebook inspector it is also going to look at URL as the feature or particularly the domain as a feature from the post that we are analyzing.

(ReferSlideTime:08:57)



So, here is the pointer to the plugin. It will be interesting if you canactuallydownload it and play around with it. These are links to the Chrome extension and to the Firefox. Let me just walk you through how this plugin works, what does it do, how is it different or how is it very similar to tweetcred.

(ReferSlideTime:09:21)



This is Facebook inspector on the Chrome store, and basically you can add this to your browser, I meanI already have it on my Chrome otherwise it should say add to Chrome. When you add it, when you go to your Facebook timeline, you should be able to see some difference in the post that you are seeing.

(ReferSlideTime:09:47)



For example here is my Facebook timeline for now and newsfeed, if you see there is no many,ifyoulook atthepost thereisnoannotationsdoneinthe poststhatyou canseeon

my timeline. Whereas let me show you some examples where the Facebook inspector is actually showing you some information.

(ReferSlideTime:10:09)



Here is an example of a post which is done by VK Choudhary and the post we just click here for Bollywood updates. Hema Malini congratulates Deepika. Is Deepika Padukone engaged? In this post if you see, there is some annotation done by the Facebook inspector, it says confidence is low, the decision that was made with the model that was generated is low and it is using features, click on the image for more details. If you are interested, you can actually click on the image, see for more details. Here is another example that I will show you.

(ReferSlideTime:10:42)



In the first case it is probably a rumour, that is why it is actually finding out and saying Facebook inspector is producing this result with red mark. Here is an example which could be a spam, whichisasofnowwereallydo notknowiPhone7designsand features that are available, but this post says about iPhone 7 is awesome and amazing whichcould have one. And this post is being annotated by Facebook inspector saying it is a malicious post. That is how Facebook inspector works.

(ReferSlideTime:11:20)

And here is also a plugin that you can use if you are a Firefox user. Therefore, if a Facebook inspector is available as a Chrome browser plugin and as a Firefox plugin add on which you can use on your browser.

So, that is the way you could think about taking way and understandings from twitter, where we studied about howto build techniques using the features from twitter to create an understanding of whether the post is a credible or not, here I showed you about Facebook.

**PrivacyandPictureonOnlineSocialMedia**

Welcome back to the course. I hope you are enjoying the course in terms of studying some new concepts, new ideas, and new solutions. This is the week 4 of the course Privacyand Securityin Online Social Media, what I will do now is ==continue== the topic on privacy that we were talking last time.

(ReferSlideTime:00:32)



Now,just ==let to let you== know we are in the topic of privacy for now,we just covered the trust and credibility, and I assume by now you are all very well versed with little bit of Linux little bit of a Python, how to collect data from twitter, how to store the data, what kind of MySQL queries you should write and collecting data and all that.

(ReferSlideTime:00:54)



**Westin's 3 categories**

- Fundamentalists, 25%
- Pragmatists, 60%
- Unconcerned, 15%

In the last week we saw about how Westin categorized all the US citizens into 3 categories; Fundamentalist, Pragmatists and Unconcerned. Fundamentalist is being 25 percent, pragmatists is being 60 percent and unconcerned being 15 percent. Fundamentalist are the people who actually do not give away any personal information. Pragmatists makedecision aboutprivacykeeping the situation in mind. Unconcerned are the set of people who gave away personal information and be part of revealing personal information is about 15 percent in the US.

(ReferSlideTime:01:27)



**Internet & Social Media**

What do you feel about privacy of your personal information on your OSN?

|  | Q42, N = 6,855 |
|---|---|
| It is not a concern at all | 19.30 |
| Since I have specified my privacy settings, my data is secure from a privacy breach | 42.13 |
| Even though, I have specified my privacy settings, I am concerned about privacy of my data | 23.84 |
| It is a concern, but I still share personal information | 8.02 |
| It is a concern; hence I do not share personal data on OSN | 6.71 |

I kind of asked you couple of questions last time about some data that was collected among large set of population in India. So this is one of the questions that I asked which is what you feel about privacy of your personal information on your online social network, which is about Facebook. About 42 percent, the highest was about 42 percent who said that specified my privacy settings my data is secured from a privacy breach.

(ReferSlideTime:02:00)



Another question that I asked you also is about if you receive a friendship request on yourmostfrequentlyusedonlinesocialnetwork,whichisFacebookinthiscasewhichof the following people will you add as friends. And the highest was actually person of opposite gender.I am pretty sure in the last couple of weeks going through the class that youaretakingonthesocialnetworknow,even yourownbehaviormaybe changing,youmay be looking at some of these requests more closely, you may be devising your mechanism by which any friend request that you get, how you are going to accept it or how you are going to deny it.

(ReferSlideTime:02:35)



http://precog.iiitd.edu.in/research/privacyindia/

Now,thedataispubliclyavailablepleasefeelfreetoactuallyplayaroundwiththedata. (Refer

Slide Time: 02:40)



## Hard to define

"Privacy is a value so complex, so entangled in competing and contradictory dimensions, so engorged with various and distinct meanings, that I sometimes despair whether it can be usefully addressed at all."

Robert C. Post, *Three Concepts of Privacy*, 89 Geo. L.J. 2087 (2001).

Last time I left you withthe question saying; what are the kind of privacyissues that you haveonFacebook,Twitter?Howyoudefineprivacy?Ithinkitisnicetoseesomeofyou     posting information about various Facebook privacy issues or your own questions about Facebook privacy issues on the forum. We should actually make the forum more active because I think there are some very repeated questions that comes up, we're tying to answerassuchaspossiblebutwhentheyareveryrepeatedwecanavoidactually

answering also. I strongly recommend you to ask, check the forum before posting the questions.

So, let us look at what privacy is a little bit and then give a little detail about some research that was goes down in terms of analyzing the privacy status on Facebook. One of the definitions that was given earlier about privacy was that "Privacy is a value so complex, so entangle in competing and contradictory dimensions, so engorged with various and distinct meanings, that I sometimes despair whether it can be usefully addressed at all." So that was Robert talking about privacy in his book 'Three Concepts of Privacy.'

But I think the privacy by definitionsis actually thought. I mean, if you were to look at what privacy is for you, why are you sitting and listening to this lecture, versus privacyin your school, privacy at home, privacy at work is very different. It is very hard todefine what privacy is for a particular individual across various situations, that is what this definition is actually trying to capture. Contradictory dimensions, so entangled and competing and contradictory dimensions.
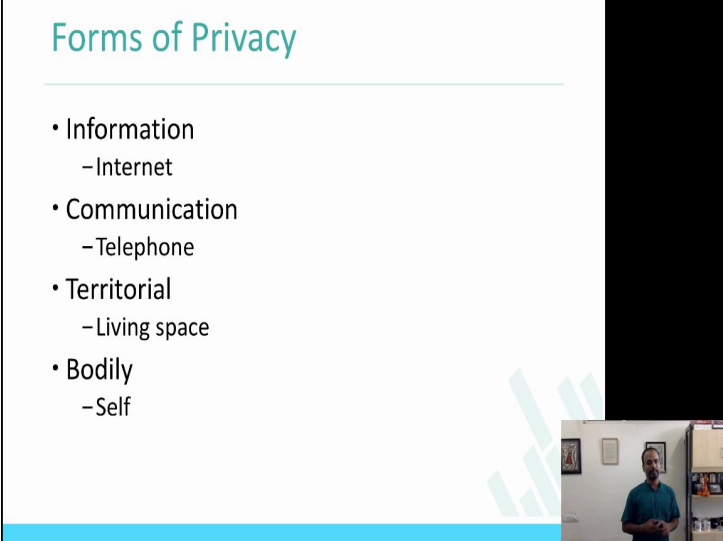
(ReferSlideTime:04:31)



Fundamentally privacy is been always talked about control over information, here are two definitions of Alan Westin actually tried defining in his book in a 'Privacy and Freedom'in1967."Privacyistheclaimofindividuals,groupsorinstitutionsto

determine themselves when, how and what extent information about them is communicated to others."

So it is basically about to determine for themselves, how much of my information I can actually share with others. "Each individual is continually engaged in personal adjustment process in which he balances the desire for privacy with the desire for disclosure and communication." How much do I want to reveal about myself, how much doIwanttoactuallyanonymizeinformationaboutmyself,howmuchdoIwanttoreveal about myself, is the way that the word privacy is defined and is the way by which youare controlling the information that you are actually spreading.

So, I am sure you kind of get the definition privacy which is very hard to define and also it is very difficult to actually come up with the list of privacy expectations for any individual in all given contexts. They strictly convey privacy is about control over information. It sometimes could be actually a group information also, given that idea is more or collective society we generally talk about a privacy of a group, instead of individual privacy, that the society is where its individualistic society where the privacy information of the individuals are more protected than the privacy information of the group.
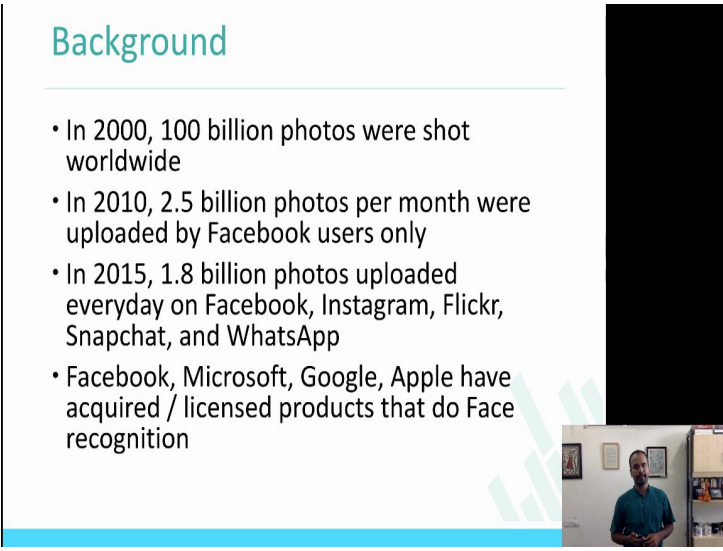
(ReferSlideTime:06:11)



Some forms of privacy that people have come up with; information privacy, communicationprivacy,territorialprivacyandbodilyprivacy.Majorityofthetimes

when we talk about privacy particularly in courses like these it is always referred to as information privacy and particularly the internet privacy.

There is also communication privacy which is telephones and other forms of communication. Territorial privacy is about my living space, my home, my city, my country and, the topics around that. Bodily privacy is about self. So, information about my own physical presence is actually also discussed in the concept of privacy. For example, a CCTVcamera is one example where bodily privacy can be actually attacked.

(ReferSlideTime:07:04)



Now let us look at some specific studies that are being done in terms of analyzing the privacy in online social networks. Here is the study that I will walk you through the referencetothe studyis attheendofthedayoftheslides,butwewalkyou throughwhat theydid, what they find, how revealing the information are, how good the studywas and how the privacy is being actually studied in the context of Facebook and social networks and publicly available information.

Somebackgroundaboutpicturesthatwereuploadedonsocialnetworksitself.Intheyear 2000, 100 billion photos were shot worldwide. In 2010, 2.5 billion photos per month were uploaded by Facebook users only. Whereas, if you remember the first lecture 1 where I actually showed you a infographic about what among the information is uploadedonsocialnetworksin1minute,weactuallysawthat1.8billionphotoswere

uploaded everyday on Facebook, Instagram, Flickr Snapchat, and Whatsapp together.So there is a lot of information, lot of pictures that are actuallyuploaded on social networks.

Companies like Facebook, Microsoft, Google,Apple have actually acquired a lot of face recognition companies in the last few years, to study, to understand, to use these technologies to identify faces on pictures that are being uploaded on the all social networks or online services. It has become very, very important to apply these kind techniques like, machine learning, deep learning and concepts around that into these images to study what is happening on online social networks, I actually recently wrote also a blog about the importance of images on online social networks. I'll actuallyshared it on the forum just after this lecture.

(ReferSlideTime:09:07).



If you really look at what is going on currently in terms of these pictures that were uploaded and the privacy about individuals, increasing public self disclosures through onlinesocialnetworkshappen,whichisItakeapictures,Itakeaselfiestandingnearone of the very important spots let us take in Delhi I upload this picture you know that I amin Delhi,orletustakeapicturenexttoTaj Mahalandupload itonmyFacebookaccount you know that I am actually traveling to Taj Mahal now.

There used to be actually a site called please rob me dot com I do not think so thewebsite is active now. This website what did they did was its called please rob me dot com, what we interestingly did was let us take it if I have a twitter account and I created

it from Delhi and posting about weather in Chennai or Hyderabad or California they would actually pick this tweet and post it on please rob me dot com saying that this account was originally created from Delhi and whereas now this post is actually talking about weather in California, so probably you are not at home and therefore your homes should be locked.
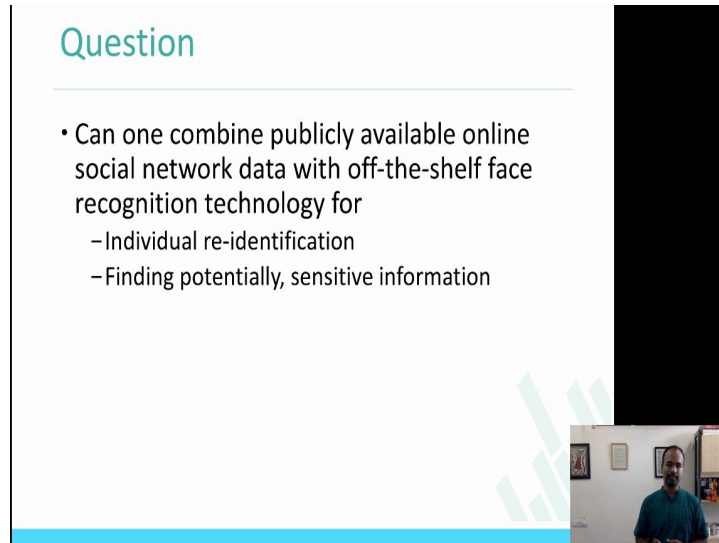
It got flacked a lot, but I think it is an interesting idea that they actually picked up to makeuseoftheinformationthattheusersofsocialnetworkaredisclosingbythemselves about their location. As a self-disclosure through online social networks and there are many manyissues that are going all around because of self-disclosure of information on Twitter, Facebook, Instagram and other networks.

Parallely in one side this increase in public information is going on. In parallel there is also increase in face recognition accuracy. In earlier the accuracy which lower now the techniques, technologies that are actually improved. In particular if you look at networks like Facebook it is actually pretty high it is because they search space that they have to search for a particular face in the picture that you are uploading is actually only your friends, majorityof the times you're going to be taking pictures with thefriends to whom you are already are connected with or probably they are in a one, and one and half hour or two hours away from here.

So, that is happening on one side. And also this is whole idea of cloud, storing information on the cloud, easily able to compute, computing cost is becoming lower and lower for doing any of these analysis. On the fourth dimension, problem is that identification of this users, who they are, what kind of information they are valuing is also getting better. Meaning, the concepts like k-anonymity came in 15 or 20 years before, but certain many further and advance techniques that have been developed to identify users, to identify faces, to identify information about users, to re-identify people on social network, people on other networks.

Those are four different things that are eluding; one, increasing self-disclosure, improving the accuracy of face recognition techniques, the whole idea of cloud and ubiquitous computing, and the techniques for re-identification of users is actually getting better and better.

(ReferSlideTime:12:33)



Theoneimportantquestionandoneinterestingquestionthatpeoplecouldaskis,canone combine publicly available online social network data with the off the shelf face recognition technology which is something that is already available, and be able to re-identifying individuals and finding potentially sensitive information. So that is the question thatweweretalkingaboutinthe nextdeckof slides which is,canwetakesome publiclyavailableinformationwhichisthatthethingsthatIhaduploadonFacebook,the things that I had upload on Twitter.

Can you use that and connect it with the off the shelf face recognition technology which is some tools like tensorflow that I will also mention later in the slides. Use these techniques to identify just basis and be able to actually re-identify the person and or also find out sensitive information about the users themselves. That is the question that we will be talking about right now.

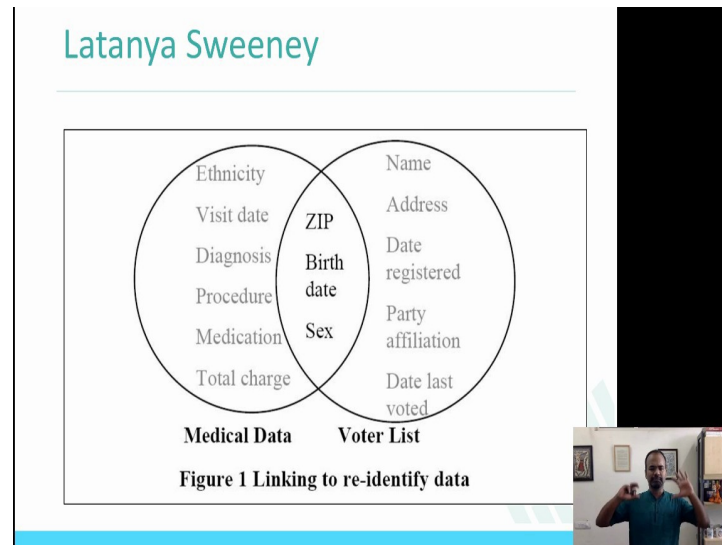Here is a goal. Goal is to use un-identified sources which is any websites that you canthinkof,matchdotcom,shaadidotcom,photosfromFlickr,CCTVfeedsandthingslike that, which is impossible to identify or its very hard, the user themselves are not disclosing who they are in these websites. It could be either they have psuedonyms and names that you cannot identify or re-identify to that particular person. Can we actually take these sources, shaadi dot com and pictures from Flickr and Facebook, connected to identify sources which are on Facebook, I would actually reveal that I am so and so on.

On Linkedin I will put this as I am so and so, on government website and other services that are available. Which is un-identified sources like, shaadi dot com, identified sources which is where I am disclosing that I am so and so, and I upload a picture my account is actually ponnurangam.kumaraguru, can we actually put these two together to get some sensitive information of the individual. For example, gender orientation like example SocialSecurityNumber,likeexampleAdhaar card number and theinformation like that. It can be pretty nasty if you can actually put this together and the get some personal information. So that is what we will be studying in our next slots.

Justto giveyou someverybroad old view ofsome phenomenonalwork thatwas donein this topic Latanya Sweeney, who did this word called k-anonymity, where she actually picked up the medical data and connected to the voter list which is publicly available. If you look at the medical data she has ethnicity, visit date, diagnosis, procedure, medication and the total charges that was paid by the patient. Name, address, date registered, party of affiliation, date last voted. Taking this information which is from voters list and from the medical data putting it together she had found actually zip code, birth date and gender was actually common among both of them.

She was able to identify if you give the system that she built birth day and gender she was able to re-identify a lot of US citizens uniquely. So that is the idea that built on to create something called as k anonymity, but the problem she highlighted was that bringing these two different sets of data which is independent medical data and voter data, you could actually re-identify users uniquely.

(ReferSlideTime:16:18)



In experiment one, they actually connected the online data to the online data. Theyinterestingly mined publicly available images from Facebook and they going to re-identify profiles just on one of the most popular dating sites in the US. They used this tool called pittpatt dot com, which was face recognizing tool. Well, after the study was done the tool was actually acquired by Google it is doing face detection and face recognition.Youcould actuallyuseTensorflownow.Tensorflowisaopen source library for machine learning techniques. Please consider exploring tensorflow little bit and how it works and what are the libraries that are available inside tensor flow.
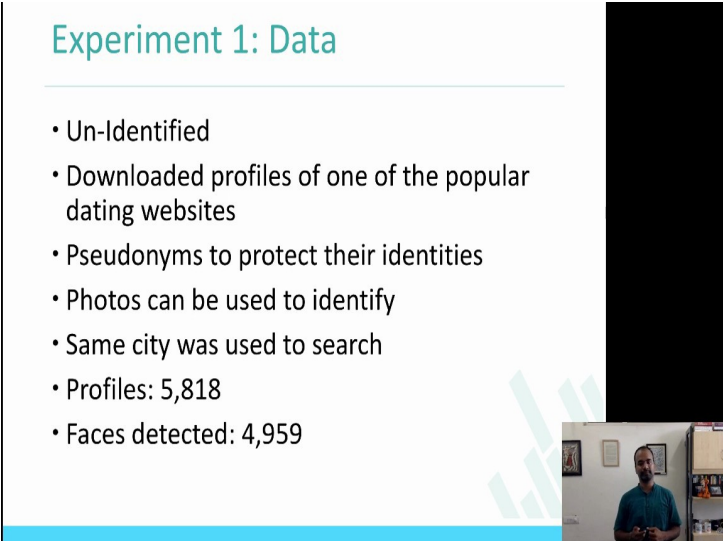
(ReferSlideTime:17:06)

The data that theyused was first as I said; they took the identified data, theydownloaded the Facebook profiles from one city in the US which is possible in the way that youknow about Facebook data collection now you could actually collect data from a particular city. Profiles that they collected were about 270,000, images that were collected around 274,000. The faces that are detected were about 110,000 faces. This is the data that they had for the identified data set, which is where you could actually say these are the names; these are profiles that are connected to these pictures.

(ReferSlideTime:17:50)



Un-identified data, they downloaded the pictures of one of the popular dating websites. So first identified, take a back; the first is the identified data, now we are talking about un-identified data, which is like the CCTV camera, publicly available information or from match dot com, shaadi dot com. Theydownloaded the profiles and the pseudonyms of their,to protect their identities, of course the names were notgoing to be revealed, the accounts may actually have pseudonyms also.

Thephotosthatweredownloadedfromthesewebsiteswhereactuallyusedtoidentifythe profile. To make the connection appropriate they actually use the same city for the search, they download data from Facebook and the city from this un-identified data set. The profiles that were collected here were about close to 6000 and the faces that were detected were about closed to 5000. So that is identified and that is un-identified data.

(ReferSlideTime:18:51)



The approach that was taken was un-identified data, dating website, identified data, Facebook profiles and the re-identification was to be done. More than 500 million pairs were actually compared, because if each picture and each of the profile, each of the data set were compared with each of the pictures in the other data set, from the un-identified to the identified and the reverse also. What they did was, they did only used the best matching pairforeachofdatingsite picture, andpittpatt andIamsure intensorflow also itgivesyou inspecificvalues,itactuallyproducesvaluesinsomerangetheyusethebest value that they could get in terms of comparing two pictures.

And to confirm, to get ground truth when this pictures are just the same data sometimesif the techniques that are machine learning techniques are not going to be fool proof and they are not going to make 100 percent right prediction. Therefore, they are actually showed these pictures to Mturkers, the users who are part of mechanical turk which is a crowdsourced mechanismwhere you can actuallyputa small task of likethis identifying where these two pictures are same people and you could actually pay them small money for doing the task.

And there were asked to rate the pictures on the likert scale of 1 to 7, at least 5 Turkers for each pair.Again please tryandlook at what are MechanicalTurkers,mechanical turk is a crowdsourced mechanism. For example, if I were do a task in identifying whether a givenemailisphishingornotIwouldactuallyitshowtotheMturkers,Iwouldcreate

the taskonmechanicalturk andget users toactuallylook atthe image andsaywhether it is phishing or not. Look at the profile and Twitter to say whether it is fake or not, they would actually go to the profile, they would click on the link in go to the profile in Twitter look at the profile and then make a judgment whether it is legitimate or not.

So it is the very popular and there are many many services like this, crowd flower which is mechanism in which many of these services come together, it is also very popular crowd flower is one - c r o w d f l o w e r, is one of the popular services like this - Mechanical turk which is from Amazon is also very popular. They took these two pictures showed to users, mechanical turkers asking to actually compare the images and make the decision. So, at least 5 Turkers for each pair because then we'll see more confidence, more and more people say that, more and more people take a image and say that this is the chair and there is high confidence that is going to be a chair.

(ReferSlideTime:21:37)



What they were able to find out was highly likely, which on the likert scale, is highly likely matches where about 6.3 percent that images that they took from this un-identified and identified and randomly they compared using the pittpatt tool and showed in the mechanical turkers. The comparison highly matches were about 6.3 percent and highly likely and likely matches were about 10.5 percent. Which basically says that 1 on 10 fromthedatingsitecanbeidentified,becausethedatingsite isanun-indentifieddataset, whereas Facebook is my identified.

So every time I see one of the pictures in the 10 pictures that I see, I will be able to actually clearly exactly identify who this person is, because I have the Facebook data, this is done of the same city and therefore it should be probably correct and mechanical turkers actuallyconfirmed that. So, you can seethat10 percentof thetimes theuserscan be actually identified.

(ReferSlideTime:22:40)



One question to you and I hope this question since there will be some discussion inforum also is that; what can you do better if you were the attacker? And if you weremake use of this information and do something to increase the rate of the efficiency or use this information to do something against the user what kind of things would you do Because as an attacker you making one this percentage to be more right, because it is 10 percent you're getting a hit rate of only 10 percent, or 1 and 10 pictures. Whereas, if you were to have a better attack or threat mechanism you could actually do things by which you can increase this percentage to more, so more and more pictures are actually re- identified and therefore it can be actually used maliciously.

Experiment 2 as I said there are 3 things. So the second one what they did was they connected the offline and the online. First one, they compare online versus online which is the dating website and Facebook, now what they did was they did the offline and the online.

Pictures from Facebook, one of the Facebook college network data was collected to identify students who are in campus and it was actually compared to the offline pictures also. What was stated when the students were actually participating in the study. So this is the experiment number 2; all connecting to the same questions which is can we actually take images, pictures from these social networks like Facebook and re-identify peoplewhoconnectedtonetworks,to datawhereuserscannot getinfrom,CCTVsource in.

(ReferSlideTime:24:21)



So, what they did was they actually put a booth in the university, took 3 pictures of the participant, they basically were standing and collecting data of the college students in this university took 3 pictures for participant, collected data over 3 days. They collected about 25 percent profiles, images were about 26,262 and the face is detected were about 114000, so Facebook data for that university. So, the data that were collected from Facebook whichisonlineisabout25000,profileswereabout25000,pictureswereabout 26000, faces were about 114000 thousand.

(ReferSlideTime:24:59)

Just to summarize or just to look at the whole experimental set up itself is that, pictures taken of individuals walking in campus, asked them to fill the survey. Next slide I also have a image to actually show you what was the process of the study. But now pictures were taken of the individuals walking on the campus, they were asked to fill an online survey.Pictures matched from cloud while they are filling the survey,because what they did was they ask that you want to participate in the study, ok I will take you 3 pictures, when they took the pictures then they asked into fill on online survey.

While they were actually filling the online survey, technique the system that they are acted would go compare this pictures what they are took to the Facebook pictures that they are already collected from the university itself and bring back the comparison and showedtothem.Lastpageofthesurveywithoptionsofthatpictures,sobythetimethey actually fill the survey they were actually shown the pictures, saying what this is the picture that we got from Facebook, do you actually agree to it. Asked to select the pics which matched closely, produce by the recognizer.

So, that is the process of the study,please understand how the study was done, collected pictures were taken individually walking in campus, they were asked to fill the survey, while filling the survey the data the system was comparing the pictures on Facebook, pictures were brought back to the survey showed to the user and saying tell us if these pictures are right about you.

(ReferSlideTime:26:28)

Same thing is captured here in the process format in the background, which is upload pictures of the users, pictures are taken which is 1 and then respon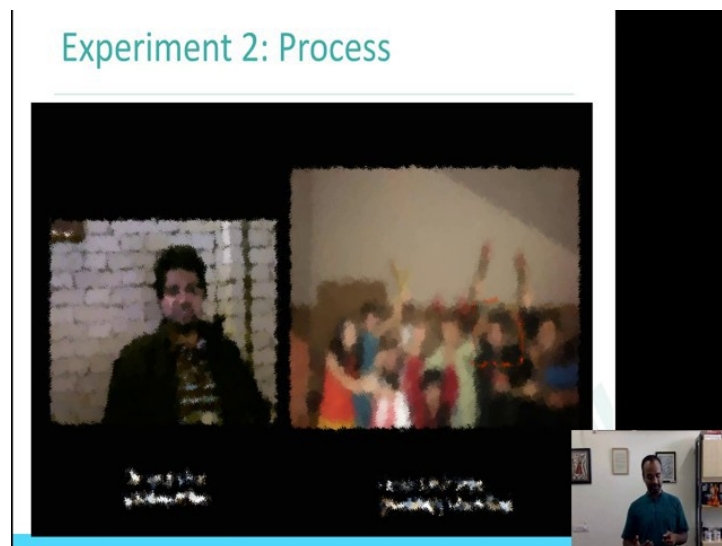ses coming from the server, start survey which is 3 and then 4 is generated survey token, so that through this survey token you will actually be able to say that comparing the images and bringing it back, which is 5is looking at custom surveytokens sendto the user who can actuallyfill the survey. And then by the time of 6 is happening which is face recognition results are being produced and then survey results both the images that are actually used which is given to 7.

Sothatistheprocessofthestudy,notaverydifficult,notaverycomplicatedstudybutit is actually collecting some very interesting data.

(ReferSlideTime:27:22)



This is the result what they did from the data collection. The left picture is the picture autonomous to the picture for the purpose of just re-identification of the user itself. The pictureontheleftisthepicturethattheytook whiletheuser wasactuallyparticipatingin the study. So when the user logged in they took the picture that on the left.

Using that the picture they are able to actually identify the picture on the right which is the picture from Facebook where this user was actually identified. So, that is the output so to say, the input is the picture with the survey and output is the image from Facebook which is re-identified this person in particular pictures. This can be actually pretty revealing the pictures compared on Facebook.

**Experiment 2: Results**

- 98 participants
  - All students and had FB accounts
- 38.18% of participants were matched with correct FB profile
  - Including a participant who mentioned that he did not have a picture on FB
  - Average computation less than 3 seconds

In about 98 participants all students in the study, there were about 98 participants, all students were the ones who participated they were collecting it from the university setup and they all had Facebook accounts also. The results were 38 percent of participantswere matched with correct Facebook profiles, which is the pictures that were taken, 38 percent of the people who took the pictures in the study were exactly matched with the Facebookprofileandtheiraccount,theirinformationisactuallybroughtbacktocompare to confirm it with the user.

Interestingly there was also a participant who mentioned that he did not have the picture onFacebook, actuallyinformationofthatparticular person,ofthatparticularparticipants was also brought back. Of course, it was actually taking very less time to do this comparison. I hope the study is making sense which is 38 percent of the times the users that were taken pictures from the university campus were identified from the Facebook profile.

Experiment 3 is interesting because they actually tried using the experiment understandings from experiment 1 and 2 to take this personally identifiable with information likes Social Security Number. In this experiment 3 they wanted to predict Social Security Number from public data. So, they used the faces and the Facebook data that werecollected from the experiment 1 and2 with the public data to predict the Social Security Number. 27 percent of subjects' first 5 Social Security Number digits were identified with four attempts.

So essentially what is this means, this means that every time I took up a face from the database, I was able to identify the first 5 digits of the Social Security Number, 27 percent of the times. That revealing, that is not a very good sign, were 27 percent of the subjects were able to find out five SSN digits of them. So that is the third experiment. And I am keeping the third experiment little light because this is in total the interesting things were pictures, un-identified data sets, identified data set and at the end they were able actually do connected to social security member also.

(ReferSlideTime:30:37)



Interestingly I am sure you could also think about how these kinds of techniques can be applied in terms of identifying Adhaar number in India also and other personal details. The study was done in the US and therefore if you were to repeat this study and find out Adhaar number or others details of Indian Citizens it will be actually interesting to look at that. If there is anyideas, if there is any questions that you have in terms of how study could be performed in India, it will be interesting to talk about it in the forum.

(ReferSlideTime:31:09)



Hereare the pointerstostudythatIjust nowdiscussedabout.

(ReferSlideTime:31:17)



And with this I will actuallywrap-up the 4.1 week. I hope you understood what we were talking about, we just talking about the Privacy Issues in Online Social Networks particularly focused on collecting images and identifying users using the face, pictures, using the images that are uploaded on social networks.

# Unit-3

## Policing and Online Social Media Part-I

Welcome back tothe course Privacyand SecurityinOnline SocialMedia, this isweek 5. Hope you have enjoyed the course by doing the home works, doing some of theexercises that we are talking in the class and generally also getting excited about the course.

(ReferSlideTime:00:27)



Thisisthe generaloverview ofthe course,wearelookingat thetopicofprivacy now.

Last time intheweekwesawthecontent aboutprivacyand howprivacyimplicationsare going on Online Social Media, in particular we saw some experiments where they are done, an analysis of data from online and online, offline and online and sometimes looking at predicting the Social Security Number also. The online and online study that they did was to take up data from Shaadi dot com type of sites like Match dot com compare it to Facebookpictures, thensee whether theywere able to identifypeople from this matches. Offline and online, they took pictures of students walking around on campus and compared to their Facebook profile and then, said that this is your profileand then got some confirmation about the picturesthat theydownloaded fromFacebook.

Theyuse this publiclyavailable dataonline, offline, allthe datato predict actuallySocial Security Number and they were successful in predicting 1 in 10 percent of the data that theysaw with some confidence ofthe five digits ofSocial SecurityNumber. So,that is a kind of thing that we have seen until now.

Let me tell you few things more about privacy on Online Social Media and then wemove on to another topic this week. It is not only that pictures, and you can use these information from social media, but you can actually use something more specific information like a location also to find out where you are and where you have been moving around and things like that.

Here is one study that researchers have done to show that privacy information from Foursquare which is one of thelocation based social networks, one of the very popular location based social network. Data from Foursquare can be actually used to find out where you live.

These two research which wass done back to back is to inferring home location fromthe check insthat you do infoursquare. Iamnot going toget into detailsofstudythat iswhy I put thepointersto thepapers, but I willtalk you ingeneral howthiswasdone, howthis could be done, and how you can actually look at some of the data yourself also.

Foursquare is one of the popular online social networks, just for locations it is called location based social network. The different topics and different concepts in foursquare are check ins. Check ins is, you check in to the hotel, you check in to the airports and similarly you check in to a location in foursquare and you can actually alsoleave a tip,let us take if I go to Sarvana Bhavan, Connaught Place in Delhi. I have food there and I can leave a tip saying that food was pretty good. And you can also become a mayor in foursquare which is, if I visit this place, if I visit this location in foursquare the most number of times in the last 60 days, I become the mayor of this location.

The mayor informationcanbeactuallyprettyuseful.Todayorganizationsare monetizing this check ins and mayorship in foursquare.Also someone actually is providing you free parking spots if you are a mayor of that location for the week. This information that is you check in can actually be used to find out your location, your home location also. Peoplehavestudiedotherthings,peoplehavestudiedactuallyfromthepicturesthatyou

upload can I actually find out your home location. This work is specifically focused on finding outthe home location fromsocialnetworks like foursquare, and there was a high confidence in finding out the home location with the foursquare check ins that people have done.

Mobility of people is actually not that much. Another conclusion that they also foundwas people do not move a lot from their current location. With this information like check in, mayorship, they were able to actually find out with the high confidence the home location of the person within few kilometers of distance of error.

Therefore, social network data privacy, initially we saw some survey where people actually said about their information of the social networks, then we looked at some studies where pictures uploaded on social media and pictures uploaded onthese publicly available websites which they called as unidentified sources can be actually used to finda person specifically or uniquely identify an individual. So here I am saying that your location can be also inferred from the social networks like foursquare.

(ReferSlideTime:05:26)



IfyouareinterestedmoreinthistopicIactuallyhighrecommendyoutogetsomeof thesepointers,whichisthesearealltheconferenceswherethetopicsofprivacyand

securityinOnlineSocialMedia ingeneralalso getspublished, but also specificallyabout these topics like, image analysis connected to the data from the web, pictures which is what wetalk about and locationand things like that.WWW - which is one ofthe toptier conferences in internet space;then,there isconference onweband socialmedia which is ICWSM. Then there is conference on online social networks and there is also collaborative work which is CSCW.

These are not the comprehensive lists, these are just to tell you that if any of you is interested in looking at these topics more, you should probably be looking into these pointers.And of course, if there is anyquestions feel free to drop it onthe forum.

(ReferSlideTime:06:30)



Now, I thought I will actually move on to another topic given that we are in the week 5 we should actually start looking at other topics also in the course.

So, until now you have done general overview of Online Social Media some Linux, Python, technicaltopics liketrust andcredibilitywhichjust lookedat privacy. So,what I thought I will do now is about spending this week's content mostly on topic called policing.

(ReferSlideTime:07:00)



I would definitelylike to hear fromstudents inthe class about how you use online social media particularly the kind of questions that I have in the slide. How many friends and followers do you have on Facebook and Twitter? If you can post all this in forum it will be nice to see howthe participants ofthis course are actuallyusing Online SocialMedia.

How many of you are friends with the police on your social network? Police, I mean Police Organizations; how often do you use social networkto post comments or interact with police? And of course, the question that I'm going to be trying to address in this week content is actually what has police been using Online Social Media for, what havethey been doing, what can they do, how we can actually help, how people using social media like you and me can actuallyparticipate with the police in online social networks.

I do not know how many of you have seen this picture before, but this is the first time ever the socialmedia was actuallyused to attack crisis, to actuallyfind out what is going oninrealworld before people actually, before the first responders, before police actually got to it. This is US airways airplane that landed on Hudson river and you could see people are actually trying to get out or people trying to help. But the first picture and the first tweet that came out, but this picture was actually a person walking on the riverside posting a picture saying that 'I am going to actually go help them.'

First untilthen this was in 2009 untilthen, the social networkTwitter or anyother social network was actually mostly used for talking to each other, saying what they are doing, hashtag Mondaymorning, things likethat, weretheoncethat were used to talk onsocial networks like Twitter and Facebook. Whereas, first time it was used to solve the problems, solve the crisis is actually this one.

And of course, there the police is being part of the social media, socialnetworks for some time now. Here is my example where it became very interesting. So, hash tag myNYPD, is the hashtagthat NYPDpolice posted saying withthis tweet 'do you have a photo with a member of NYPD, tweet us and tag it, with myNYPD. It may be featuredon our Facebook page. This came out on Twitter and they said that we can actually post the pictures on our Facebook account. Interestingly, with this hash tag they startedgetting pictures like this.

(ReferSlideTime:10:08)



ItisnotonlylikethattheygotpictureslikethiswhichisfrommyNYPD.

(ReferSlideTime:10:12)



But interestingly there were also pictures, this hash tag went from myNYPD tomyLAPD.
The point here is that, these kinds of strategies that police use is also
happeningonOnlineSocialMediawhichistotakeupthehashtag,getpublicto

participate and uploading the picturesand ofcourse sometimes you get different kindsof pictures <mark>than</mark> what they meant or expected.

(ReferSlideTime:10:43)



Let uscome backtotheIndiancontext.Inthe lastfew yearsprobablylast oneandhalfor 2 years there is a lot of adoption of Online Social Networks, Online Social Media services in India also. Here are some examples, I am going to go through a lot of accounts, handles which <mark>has been</mark> using Online Social Media in the last few years. And then, over the process of actually talking about these handles I am also going to inject some technicaltopics interms ofwhat arethe problems, how actuallypeople <mark>talking</mark> this course and people understanding the social networks technically can also help in identifying these problems. The one on the top leftis actually Delhi traffic police, theone on the bottom is actually Hyderabad city police. I will actually show you moredetails about examples also pretty <mark>soon</mark>.

(ReferSlideTime:11:37)



Here is the general way by which Police Organizations are actually using the social media services. This is a Facebook page of Bangalore City Police. In India now, Bangalore City Police, Bangalore Traffic Police, Delhi Traffic Page, these are the very popular handles inthe countrynow.And that iswhy you willget,teaching thisparticular course is actually pretty exciting for me, it is because we are actually looking at topics which are very rather relevant, I mean just open your Facebook now and look at Bangalore City Police you will actually look at some of the things I am talking now. Bangalore City Police is the verified page you can see a blue tick next to the top left of Bangalore City Police.

And there is about thousands and thousands of likes this picture that was taken sometime back I am sure the likes have changed now.And in the bottomthe picture showing you an example of a typical post that comes from these kinds of police pages. The post says that 'we are taking up traffic signals synchronization on 10 corridors in the city for smooth traffic flow' and it is coming from handle AddICPTraffic. So the idea is that citizens can inform about what is going on in the decisions that theyare making.

(ReferSlideTime:13:07)



Andyoucan of course see a numberof likes,peopleinteractingwith these posts,with the Police Organizations.

(ReferSlideTime:13:20)



SoCPBLR,atCPBLR isahandlethatgotverypopularaboutoneandhalfyearsbefore.

(ReferSlideTime:13:32)



This is the current account of CPBLR, it has a part I just took the screenshot a few minutes before actually preparing this deck of slides. It is about 10000 tweets, the account is following 60 people and about 688,000 people are actually following this account. Then this shows that this account is actually very popular.

For example, probably I took the screenshot a year or sometime beforeI took on 28thlast year probably and it says about 335,000. And now it is actually August 2016 it is about 688,000.The activity of this account is actually very high and the number of followers; the interaction that this account has with citizens is also pretty high.Therefore, they gain a lot of followers.

So, what we are going to be looking at this is, how these accounts are doing, in specific we will also look at some data that was collected to analyze Bangalore data itself.

(ReferSlideTime:14:40)



## Popular departments

| Police Departments | Likes | Followers | Post | Joined |
|---|---|---|---|---|
| USA | | | | |
| New York | 383,372 | 147,000 | No | 2012 |
| Boston | 137,403 | 312,000 | No | 2010 |
| Baltimore | 36,530 | 70,400 | Yes | 2012 |
| Metropolitan, Columbia | 16,071 | 56,900 | Yes | 2008 |
| Seattle | 12,912 | 103,000 | No | 2010 |
| UK | | | | |
| Greater Manchester* | 98,193 | 205,000 | Yes | 2011 |
| West Midlands* | 86,904 | 115,000 | No | 2008 |
| Essex* | 66,461 | 85,300 | No | 2011 |
| London* | 46,889 | 267,000 | No | 2011 |
| Northern Ireland* | 26,173 | 71,300 | No | 2009 |

Just to give you a sense of how the Police Organizations around the world are. Police organizations like New York, Boston, Baltimore and many, many other cities have actually adopted Facebook's and Twitter's. Here we are just showing you the Facebook which is likes and followers. The common post and joined is that joined is the year that they started. Post is whether they are allowed to actually post the content on theFacebook page, because on Facebook you can actually control whether the people who have liked the post or people who are accessing this page can actually post on to thepage.

In UK also there has been a lot of adoption of Facebook. Of course, if one is interested you could actually look at this topic more closely and do a lots of analysis on analyzing the post, analyzing the content that is going on online social network, both in the policein context. If any of you are interested in doing it,if any of you are interested in looking at this problem more closely I will be happyto actually chat with you.

(ReferSlideTime:15:55)



Popular departments - India

| Police Departments | Likes | Followers | Post | Joined |
|---|---|---|---|---|
| Bangalore Traffic | 2,49,968 | 8,045 | Yes | 2012 |
| Delhi Traffic | 2,02,858 | 2,59,000 | Yes | 2011 |
| Hyderabad Traffic | 1,88,480 | 1,361 | Yes | 2012 |
| Bangalore City | 1,05,463 | 12,100 | Yes | 2011 |
| Kolkata Traffic | 63,789 | - | Yes | 2010 |
| Chennai | 50,979 | 1,108 | Yes | 2013 |
| Gurgaon | 43,901 | 718 | Yes | 2013 |
| Gurgaon Traffic | 24,475 | - | Yes | 2010 |
| Hyderabad | 13,602 | 537 | Yes | 2014 |
| UP Police PR | 8,486 | 4,585 | Yes | 2013 |
| Guwahati Police | 3,255 | 295 | Yes | 2011 |

So, here is the situation what is in India now? This is slightly outdated data but I think this is Bangalore, Delhi Traffic Police which are actually popular, and of course many cities have actuallyadopted getting onto Online SocialMedia and particularlyFacebook and Twitter.

There is a little bit of YouTube but I think YouTube is pretty small. In our case the accounts have probably started after 2010 or 2011 and all the pages that are created in India for using for policing allows citizens to actuallypost onthese pages.

(ReferSlideTime:16:38)



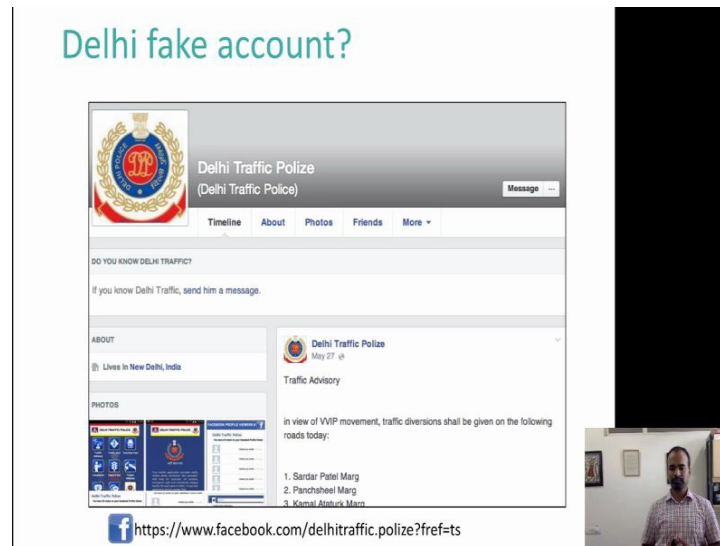Let me walk you through this some city police handles and then look at what kind of things that they do on these handles. This is Bangalore City Police, Bangalore Traffic Police, Facebook and Twitter page, this is a typical police I told you about traffic earlier the other post that they could do is something like this which is cash reward of Rupees 10 lakh for helping them.

(ReferSlideTime:17:04)

This is Delhitraffic.Again Delhitraffic, Twitter and Facebook pages.This is the kind of post thattheDelhipolicealso doesintermsoftraffic, intermsofactuallyappreciationof their police force, their police officers.

(ReferSlideTime:17:22)



Interestingly, there is not only that these handles have been managed by these Police Organizations, there are also fake handles that are being generated because of using of social media byPolice Organizations also. I think it is typicalproblemthat technologists like students who are taking this class should probably look at.

What are the ways to actually identify this kind of fake content? First of all even to identify whether there are fake handles, 500 accounts which are doing the similar things that the Delhi Traffic Police page is doing or any other Police Organizations. Wherever you are from meaning I am sure the cities like you are sitting in and looking at these lectures you could probably look at whether the Police Organizations in your state, in your cityhas a Facebook page or aTwitter page and see what kind ofactivities that they are doing.

It will be interesting if some of you actually post about your own city or state police organizationactivitiesontheforum.ItwillbeinterestingtolookatwhatPolice

Organizations are actually doing. Maybe, by looking at these topics for sometime you willget a sense but Ithink ifyou are a localite you willprobablyalso understand what is going on. Inthis case, the slide that I have here which is DelhiTraffic Police- po l iz e is a fake account and there are many, many accounts like these for a different Police Organizations.

(ReferSlideTime:18:46)



Here is UPPolice. UPPolice also has a Facebook page and aTwitter page and theyalso seem to be talking more about the activities in terms of traffic in terms of helping them and in terms of generally interacting with the citizens.

(ReferSlideTime:19:03)



That isHyderabadPoliceofcourse.

(ReferSlideTime:19:08)



This is the Hyderabad Police recent Twitter handle; I just took the screenshot few minutes before again. This is also a verified account; I will come back to this verified accountinfewminutes.ButIthinkletmejusthighlightwhataverifiedaccounthereis,

the blue tick mark next to Hyderabad Police here is the verified account; says that this account is verified that is, Twitter has verified this is reallythe Hyderabad Police.

Itisnotthat easyto get averifiedaccount, not allaccountsonTwitterareverified;onlya little fraction of users on Twitter are actually verified, I think it is about a million users now.These verified accountsare veryhelpfulbecause Ithink for Police Organizationsto say that who they are is actually very, very necessary.

(ReferSlideTime:20:09)



Just to show you this is Hyd City Police which is the real verified account, but as thereare actually accounts like Hyderabad police which does not even have profile picture changed; this is called an egg profile in Twitter terms, and you can see this handle hasbeen from 2012. Sometimes, this handle may not exist now but I think when we wereanalyzing, we keep track of some of these handles also.

(ReferSlideTime:20:30)



Sometimes, some states have taken the decision of all the states and all the city police will have just one page, and then everybodywill be posting content and interacting only throughthat. Some organizations have takenthe decision ofhaving multiple pages for at the district level for different activities.

(ReferSlideTime:20:52)

GuwahatiPolice.

(ReferSlideTime:20:56)



Kolkata Police. I think Kolkata Police also there is a interesting handle which is on the second from the top on the left which says Kolkata Police, fake Kolkata Police twitter,and then a smiley. Therefore, there are these kind of problems where Police Organizationshandles have been mimicked, mastibated to create accountsand see ifyou cansee it actuallyhas at Kolkata Police;the profile picturethat it isusing isprobablythe same as a legitimate Kolkata Police account also.

(ReferSlideTime:21:27)



PunePolice.AgainPuneTrafficPolice, PunePoliceand it isactuallyveryhardto verify, even for us, even for anybody, even when you go to get your own cities or state police pages as I said now, you yourself will find it hard to get the actual real account which unless it is verified is going to be a veryhard to actually justifyor find out whether they are legitimate.

(ReferSlideTime:21:54)

ThisisTrafficPoliceJaipur.

(ReferSlideTime:21:59)



I am just giving you vulnerant tour of some of the handles that we have found. The thought here is that whether you like it or not, whether as citizens want to see their accounts Police Organization accounts in these networks, the accounts are created like I have taken police there are many other polices; it's not clear whether these are actually the Police Organization handles.

(ReferSlideTime:22:24)



Here is a quick one in terms of Mumbai Police, at Mumbai cops, it iays protecting and serving Mumbai; keep updated with latest news on how we are combating crime, fearanddisorderinthecommunity. It isnot clear whether it isa legitimateaccount.Andthen there is Chetan Gavali at Mumbai Police, Delhi Police too there is a handle.

(ReferSlideTime:22:47)

There are multiple handles like this. So, one thing when we were looking at thesehandles we realized is that, hard to find out which is <mark>legitimate and</mark> which is not legitimate. That is <mark>when we</mark> created this page in capturing all the State Police Organizations and City Organizations, their Facebook pages, Twitter handles and the source from where we are actually getting this information.

(ReferSlideTime:23:17)



So, now let me walk you through the website that I have just now mentioned which is looking at different handles. If you see here, we have been capturing different State Police Organizations, their Facebook pages and their Twitter handle. Source column is basically to show that whether we were able to confirm whether this is the legitimate handle.

(ReferSlideTime:23:48)



| | | | | |
|---|---|---|---|---|
| Sikkim | Sikkim Police | SikkimPolice | sikkimpolice | Taken from website |
| Sikkim | Traffic Police Gangtok Sikkim | Traffic-Police-Gangtok-Sikkim-127599730651346 | - | Yet to be confirmed |
| Tamil Nadu | Chennai City Police | Chennai.Police | chennaipolice_ | Yet to be confirmed |
| Tamil Nadu | Chennai Traffic Police | chennaitrafficpolice | cctpolice | Yet to be confirmed |
| Telangana | Hyderabad Traffic Police | HYDTP | HYDTraffic | Taken from website |
| Telangana | Hyderabad City Police | hyderabadpolice | hydcitypolice ✓ | Taken from website |
| Telangana | PS Narayanguda Hyderabad City Telangana State | psnarayanguda | - | Taken from website |
| Telangana | Cyber Crime Police Hyderabad | cybercrimepolice.gov.in | - | Taken from website |
| Telangana | Cyberabad Police Commissionerate | Cyberabad-Police-Commissionerate-1564780007078190 | - | Yet to be confirmed |
| Telangana | Cyberabad Traffic Police | cyberabadtp | - | Yet to be confirmed |
| Tripura | Traffic Police Agartala | Traffic-Police-Agartala-154008624746928 | - | Yet to be confirmed |
| Uttar Pradesh | UP Police PR | UpPolicePr | uppolicepr | Taken from website |
| Uttar Pradesh | UP Traffic Police | - | dirtraffic | Yet to be confirmed |
| Uttar Pradesh | Igzone Lucknow | Igzone-lucknow-160549427441397 | - | Yet to be confirmed |
| Uttar Pradesh | Noida Traffic Police | - | noidatrafficpol | Yet to be confirmed |
| Uttar Pradesh | Agra Traffic Poilice | - | Agrtraffic | Yet to be confirmed |
| West Bengal | Kolkata Traffic Police | KolkataTrafficPolice | KolkataPolice ✓ | Taken from website |
| West Bengal | West Bengal Police | West-Bengal-Police-136547649780044 | - | Yet to be confirmed |

Soyoucanclearlyseethereabout 80pages in80accounts, 80rowsinthispage.Youcan clearlyseethatthereareonlyfewofthemwhichareverified, because verified istheonly way to actually confirm that this page is actually legitimately the Police Organization page. And the other way that Police Organizations are doing is to connect the actual website which is Delhi Police dot gov dot in and probably link it to the actual Facebook or Twitter handle from that page. This just helps us to confirm that the legitimate Facebook page or Twitter handle.

You cansee that there areonlyfew ofthemare verified, Kolkata Police, Hyderabad City Police, Bangalore City Police, Bangalore City Traffic, both of them in Facebook and Twitter, and then Delhi Police and CPDelhi. We trykeep this page updated ifand when we get to know that if there is any page which we are missing.

Againif there's any page that you think we are missing please let us now in the forum,we will be actually happy. The idea here is for you to get a sense of how social media is being used in these context. So, keeping the topic of the course in mind in terms of the data collection analysis, let me show you how one can actually collect the data fromthese social networks and actually do some analysis which could be very useful for making some decisions.

(ReferSlideTime:25:12)



(ReferSlideTime:25:20)



Here isa pagethat wehave set upfor sometime which isbasicallya list ofpagesthat are taken fromthe earlier page handles and then it is not only the link to those pages. So let me show you.

(ReferSlideTime:25:27)



(ReferSlideTime:25:29)



Here what we have done is we just have linked to these pages, it just opens up the CP Delhipage for nowfromthis link.Whereas, herewhat we did waswethought we should collect the data and do analysis on the data itself. There is going to capture some things that you have done in the tutorials which is taking data from these networks.

So now let me go to the Karnataka and let me go to Bangalore CityPolice to show you. Essentially the India map on top is capturing all the handles that we have in the database and fromthere if you click ona state it is going to take youto the state and show you all the pages that are there in that state, from there you can go down, go to a specific account. What I am doing here is, specifically looking at a Twitter handle of Bangalore CityPolice, so here isthe citypoliceTwitterhandle and wehave also capturedthe actual website which is dcp dot gov dot in which is useful for us.

(ReferSlideTime:26:28)



Nowwhat we arelooking at isthe Facebookand youcangotoTwitteralso.InFacebook if you remember inthe tutorialwe talked about data collectionand so now collecting the data from these accounts, putting them into the database from the json and doing the graphright. So wearegoing to slicker some other graphs that wecanactuallyprettyuse. Here if you see, you can actually zoom in to the data to see that how these handles have been posting.

(ReferSlideTime:26:58)



For example here, it says let us look at it here, Bangalore CityPolice, the likes that they got were1 and thenthepost that theydid was12 and the commentsthat theygot was32. This kind of gives the sense of how active these handles are. You could do interesting analysis with these kinds of data from these social networks.

You have enough skills now in the course to collect such data and look at the data also. This week even the course questions that we are going to be looking at for the course is going to be based on something, some data that you can collect and some analysis that you can do. So that gives you a good sense.

(ReferSlideTime:27:44)



Then let us look at only the post analysis, which is the post that they have done but if theydo, what time did theydo, what post you can keep zooming into the datato look at, what kind ofpost, when did they do, what number ofpoststheydo. For example, here itsays in December 29-2014 they did 59 posts.

(ReferSlideTime:28:10)

Similarly, you could look at likes and comments also. This kind of gives you a good sense of how these Police Organizations are using the data, how this data can be collected, what analysis can be done, and I am sure you can do more analysis with this data also. I am going to be also specifically talking about one specific set of questions that we were trying to answer with the data that we collected from this social media services for the Police Organizations.

# Policing and Online Social Media Part-II

(ReferSlideTime:00:30)



Welcome back. So, let me now talk a little bit about the specifics of how the data from police organizations can be collected, and what kind of analysis can be done to find out some interesting things.

Here is one research question; here is one question that you can think about - objectives of the study. And then I am going to be taking about whether online social media can support police to get actionable information about crime and residents' opinion about policing activities in urban cities yeah, so that is the goal. So, let us try and, see if youcanactuallyteachthisobjectiveto studysomedata fromFacebook andTwitterand make some useful inferences.

So, let me just breakthisobjective into pieces, which is, canweuse Facebookto support policetogetactionableinformation?Whatisanactionableinformation,actionable

information is something like do this, can you actually get this done, I mean I am having a problem inthe streetthat is traffic issues inthe roadthere is a pot hole which is broken on this street, a car broken down. So, these are actionable information that police organizations cantake fromthe post and that is actuallyuseful for decision making.And residents'opinion, ofcourse, what people think about police, what arethey talking about police is also useful information for police organizations.

(ReferSlideTime:01:46)



So, before I get into further I think there was a question in the forum asking aboutwhatisre-identification. IthoughtIwillactuallymentionit here ratherthanactuallypacking it in the forum itself. So, re-identification is nothing but, take some information, you want to,actually, yougot ainformationabout PK,thereissomepubliclyavailable information which has no reference to PK. For example, if you remember the LatanyaSweeny(Refer Time: 02:14) slide where we talked about voters record and medical record. Justin medical records they are actually identifiable. Just in voter records also they are identifiable. If we put together the identification actually becomes much stronger, youare able to uniquely identify more people with more data put together.

For example, again, let me go to myown example, you cantake some publiclyavailable information about me on some websites. Say oh faculty at IIIT and things like that. And

you go back to Facebook, and then you use the Facebook pictures that are publicly available about me, take those pictures connectit with these posts you can say it oh this is actually PK, this is how I willalso I mean insides your faculty and IIIT, NPTEL.

So, re-identification of information of a particular individual, of a particular thing is actually the concept that we discussed last week. I hope that makes it clear which is unidentified datasetswhich inthe class that wetalked about max.comand identified data sets we stored which where there is if your time you know one can find out the this is here. So, taking some unidentified data and using some identified data putting them together and identifying the users actually is an (Refer Time: 03:26).

(ReferSlideTime:03:48)



So, wedidthisworkonsocialnetworksforpolice andresidents inIndiaexploringonline communications. So, this is the paper that with I am going to be talking about, but this actually more than a paper that the data that I will be talking about right now.

So, in general, actually in the last two years or so social media has been used for crime prevention. It can be used effectively for finding out what people are saying, you can actually collect the information lot of things about what is going on in the society, because it is going to be a very hard to have police organization, police personnel ateverygivenpoint intime, at anygivenpoint inthe societyalso. So you canactuallyget a lot of information frompublic throughthe social networks, which can be used to prevent crime. So, essentially you can build societies which are safer if were to actually analyze use social media services.

(ReferSlideTime:04:26)



(ReferSlideTime:04:30)



So,intermsofactuallythethemeitselfthedatathatwelookedatisactuallyfrom Bangalore,
Karnataka.

So, in terms of methodology, what kind of data did we collect? So, keeping the goal for studying whether we can actually collect actionable information from social media we started looking at this data, we collected the data from the Facebook page of Bangalore City Police in 2014. Looking at what are the posts that was done. And we filtered the posts and comments, because we wanted to studywhat public said to the police in terms of what post that they did, what comments that they say on the Facebook page. Andabout 1600 comments and 255 posts were actually collected.

(ReferSlideTime:05:19)



So, interms of methodology, there are actually multiple ways the people actuallylook at this data type. We are looking at the post and we are looking at the comments we can analyze indifferent ways. So,oneapproachthat wetookwasfindingout what peopleare talkingaboutwhichis misinformation,query,trafficdetailsthat isaboutthecontent.And then we looked at for the style of I think which is formal or informal.

And interms oftypes ofpostthat wereshowing up which is acknowledge to, like, or say thanks, reply to, such suggest a solution and the follow up by asking further details, ignored byno reply, because these are the ones that are coming fromthe police side. So, citizens post and what do police do about it.

So, if you look at in the right hand side, it is says twenty four categories for the public post and two categories for the style, and four categories of the police responses. So, again given that we are talking about content and injecting some analysis that you could do with the data yourself also in terms of lexical analysis in terms of actually (Refer Time: 06:27) the content itself that is from the post.

(ReferSlideTime:06:33)



So, if you look at the results, some of the results are very interesting in terms of what kind of post were done by citizens for these on this page. Majority of the things were actually from the neighborhood concerns right. Then it is appreciation which is talking about thanks to police and appreciating the things that police does. And it kind of goesdown. suggestions, auto driver related, fraud, till traffic issues.

And, if you look at the comments and the likes, the comments for actions like appreciation are actually higher than the comments for satisfaction; whereas, if you look at the likes, the likes for satisfaction, appreciation and success stories are actually veryhighe. It is probably very intuitive that (Refer Time: 07:20) how the police post gets reactions fromthe society, the likes are actually pretty high for satisfaction, appreciation post and for success stories compared to some of the other ones.

So,this isgives you asenseoftheanalysisthat you cando withanykind ofdatathat you collect.You remember we talked about Bostonblasts and Hurricane Sandy (Refer Time: 07:44) and those kind of events in the context of credibility and trust (Refer Time:07:50). So here, we aredoing these similar kinds of analysis, similar kinds of questions that we are asking, but we are actually using different sets of data and different kinds of graphsthatweareproducing.So,thiswillhelpyoutogetasenseofwhatarekindof

posts are actually showing up on these pages, and what kind of reactions are being seen on for these posts also.

(ReferSlideTime:08:13)



Similar to the analysis that we did in Boston and Hurricane Sandy, we can do the geospatial analysis also with this data here (Refer Time: 08:22). The one on the top is showing you the poststhat arecoming fromthe different partsofBangalore forthe posts that we saw in the page.And ofcourse, one could do some heat map, one could find out where are the places fromwhere majorityofthe posts are coming and you could use that for making some decisions.

So, given the goal was actionable information, we were actually focused on finding out fromthe content what kind ofinformationcan be drawn. So, here is one post whichtalks about temporaldata at least which can be drawn, (Refer Time: 09:00)time between 5.30 and 6 pm. Location, blah, blah, blah, not a single police postedhere,I waswaiting for an auto at the circle blah, blah, blah right. So, this gives them, this gives the police organizations a good sense of what time is it, what location is it, what should be done, what is the problem, it is easy to actually collect this information. If it was not given in this form, if this information was not there, the police has to actually ask saying what location is it, what time is it, and things like that.

 So, if you look at the communication style, the style is also interesting (Refer Time: 09:37) that lot ofdiscussions that happen on from the police side is actually very formal; formal versus in informal. Dear Sir, Request to take action on Railway Station this is from the citizens, parking contractors they are not issuing parking slips right. Kudos to the BanasawadiTraffic Police Team. My salute and this is post for appreciation.

So, from the police if it comes, it is (Refer Time: 10:07) almost going to be always formal.And stayvisible ofcourse,this isthe point Ihave said earlier, which isFacebook and any social network for that matter can become the way by which police can actually connect with society most strongly.

And I am sure as you are going through the course you will also start looking at, I hope you will also start looking at the police pages of a local city from your location are actually start saying, what kind of post they are doing, what kind of things that they are looking for, what kind of interactions are theyhaving. So, the whole bodyofknowledge, body of research, body of work is to actually look at increase the community policing right. So, you can actually increase the interactions with the society to get more information from them.

(ReferSlideTime:10:53)



Of course, these are some details, I will go through them slightly quickly. Average time response 30 hours, maximum time was 211 hours, minimum time was about 4 minutes. Showing that there is large variance in terms of actually responses that they get 4minutesto 211hours, sothatisa lotofdifference intermsofthetimesresponsethat they get.

(ReferSlideTime:11:18)

So, herearethedifferent typesofpost that come frompoliceand thekind ofengagement that they have. So, acknowledged 21 percent of the post police actually acknowledged. Dear XXX, we will take all possible legal measures in this regard, thank you. And as a reply Dear XXX, Please lodge a complaint – 22 percent. And Dear XXX, This post has been forwarded to appropriate police station. And about follow up, Dear XXX, Please provide the police station details. Thank you.

So, this kind tells you what kind of interactions of police organizations having. Andabout anything postednot have getting response. So, the goalis to find out theoneofthe one ofthe interesting questions that you could also think about is how to actually have a post whichwill have the response frompolice, that would be also an interesting question to look at.

(ReferSlideTime:12:20)



So, if you look at the concept of finding out what citizens are worried about, what citizens are talking about, I will go through some tweets what we so to say in terms of actually looking at what poststhe citizens are doing, how we can actually take out some useful information from these posts. So, in this case, we aretalking about worried as the starting pointwhich isfrom the postyou can actually look at worried,if somebody will

misuse my bike, worried at the person who is duplicated my registration number will commit, blah, blah, blah, worried about they coming back to attack me.

Forinstance, thisconceptofidentifyingthecontent,textualcontent andseeingwhat kind of posts that citizens are doing can be extremely useful. If only if one can generate these trees in real time it can be very useful for police to make some decisions.And if this can be done in real time to showing up, oh, currently there is a post on Facebook which hasactionable information and the actionable information is the time, the details and this citizen's post is actually having about worried about few things can be veryuseful.

(ReferSlideTime:13:37)



Of course, the direct versus indirect information drawing from this post - Direct information, it going on with me, I am actually going through the problem or I am actually part of the situation that I am talking about. Sometimes it could be indirectwhich is 'Dear BCP, though I stayat JPNagar, but being part ofKSFC layout near blah, blah, blah. Iamnot fromthere, but Isee aproblemthere, so I'mletting you know (Refer Time:14:05)It could bethat myfriendsaysthis,myfriend livesthereontheregardingly I post, post on Facebook, I do a post on Facebook about the friend that who lives in a different location not about myself.

So, there are recent posts also which is directly about myself and indirectly about somebody else. And of course, being it is probably very intuitive or to realize that you could take the content from the social networks, and actually took for accountability of both sides for example, accountabilityofpolice and accountabilityofcitizens also. How we candothat we could lookatthe post and see howfast theyare responding, what kind of responses they are it is coming and how citizens are also responding to these queries that the police is making.

So, accountability can be good question to ask from the post that is collected from the social network. And of course, the little that we have seen little that is being looked at there is also mutual accountability that is going on citizens think that police should be doing and (Refer Time: 15:21) police think that citizens should be doing it, there is a accountability, in terms of, because this platforms publicly available.

(ReferSlideTime:15:24)



And ofcourse, police organizations respond to this post, and request for informationand follow up onthings also making themselves accountable for the activitythat is going on.

(ReferSlideTime:15:43)

A citizen accepts that they are also accountable to make the city safe. Citizens also believethattheyshouldbeparticipating intheseactivities intermsofposting, interacting with police, giving information and making sure that the city is safe.

(ReferSlideTime:16:01)



So, if you look at the tree again earlier the example was worried now if you look at the other concept ofwhy. So, whythese illegal practices are not being stopped?Why do not you stop tobacco?Why this, whythat, right. So,this could also be a good wayto look at the content and cull out the actionable information from these posts right. So, these are things that you can do this is the types of analysis that you can do in terms of what citizens are talking about, what police organizations are actually posting.

Of course, police can also understand needs and wants right police can encounter fearand anxiety if they know resident expectations like needs and wants I want something living in this place I want some specific safety, I know that this is happening, I'm complaining, please take care of it.If only all this can be done using the content, using the information that is coming on social media, it could be very helpful. Of course, it is notthat onlythis is the onlysource for making allthese judgements (ReferTime: 17:10).

(ReferSlideTime:17:11)



Looking at few more examples in terms of need to be punished blah, blah, blah, need to be so and so, need to hang this guy, need to do more research on why that is going on, need such information for doing this. So, this kind of tree(Refer Time: 17:28)information can be actually very helpful, I think I have emphasised enough about thistree, I'llgo through.

(ReferSlideTime:17:31)

So, this is about needs which is what this needed by the citizens, and what did theywant also. Want to hear more of these, want to see the punishment of xyz, want to and ==deletetherest== (ReferTime:17:46) want tosaythankstoBCPSIRright.Thiskindofanalysis in terms of wants and needs which is also connected to the actionable information that we talked about is very helpful.

(ReferSlideTime:17:59)



So, now just keeping these things in mind the data that we have collected from the Facebook Bangalore Citypolice what arethe things that we canthink about. Just aquick summary of ==what we looked at also, right,== (Refer Time: 18:10) in terms of the data one could actually look at collecting allthese information, and helping understand actionable information.

Actionable information, in the sense that I showed you it was just a ==tree==, but how you actually take this ==and give it== to police organizations to look at. It could be that the same ==tree could be shown, but I think== highlighting some post saying that here is the post that you should lookat more carefully.And probablywhenproposing what kind ofpost to be givenand for a specific post here is a template for the replythat you should produce and ==thingslikethat==.Increasingtheproductivityofthepoliceorganizationslookingatthis

post canbe very, veryuseful. Ofcourse, wesawthat bothcitizens andpoliceareactually accountable because they are actually interacting on this public forum.

And of course, that is also understanding of fear; understanding of wants, understanding of needs from the citizens for police also. With that I will stop with this part of thelecture which is so to in this week, we looked at how initially we just started off with privacy, closing up the topic on privacy, then we look at different police organizations Facebook handles why they should do, what kind of post show up on these Facebook pagesandTwitterhandles,what kindofhandlesexist.Andthenwe lookedat specifically analyzing the post for identifying actionable information. With this, I will stop this lecture.

# Policing and Online Social Media Part-III

Welcome back. Letuscontinueonthetopicofhow manySocialMediaand thePolicing. Now,I am going to look at different set of questions within the context of <mark>using</mark> content from social media from the Police Organizations itself.

(ReferSlideTime:00:25)



So, the question here is can we explore the feasibility of social media in quantifying attributes of communication, which is how the communication happens between Citizen to Citizen, Police toCitizen, Police to Police and Citizens to Police, in the frame workof Facebook, in the frame work of a Facebook page of Police Organization. So that is the question which is, can we actually find out what kind of attributes, how the communication happens <mark>within</mark> these four sets of interactions.

Then, identifying behavioral attributes like affective expression, engagement social and cognitiveresponseprocess;insimpletermitisbasicallylookingatfindingouthow

positive or negative the interactions are, how frequentlythe engagement happens, how is the response between the parties involved in the interaction.

(ReferSlideTime:01:20)



Since, specific questions that were look at are; Topical Characteristics, Engagement Characteristics, Emotional Exchanges, and Cognitive and Social Orientation. Let me walk you through little bit about what these are, and then show you some data and try to finds some results with the data that was collected, so nature of content and topics that characterize social media discussion threads. That is the topical characteristics which is what kind of content and topics people talking about.

And how the citizens and police engaged in social media discussions, how the engagement happens, what level of the engagement is it, how much engagement is there in the interaction. The third one the emotional exchanges, is a nature of emotions and affective expression that manifest in social media, which is more like how much emotions positive, negative, some concepts that are actually talk about later is violence, arousal and things like that. Cognitive and social orientation, what are the linguistic attributes that characterize cognitive and social response processes in the context of the policing Facebook pages.

These arethe fourthedifferentaspects that wewilllook atin terms ofanalyzing thedata and seeing what we can draw from the content that is getting generated. Again please keep in mind the way that I am driving this whole course, whole content is to actually take a problem try to ask some questions, look at some data, make some, do some analysis, make some inferences, and come back to the question that we have asked and say what did we learn from those questions.

If we see from the starting point that is how we have been doing credibility and trust, Boston Marathon, Hurricane Sandy kind of topics and then on privacy to doing some analysis of the face recognition and even in policing earlier part of this week we lookedat the content how the police and citizen, only Bangalore City Police interactions were going on.

(ReferSlideTime:03:22)



If you remember in the second part of this week I actually showed you a page where we were collecting, where they were pages of different Police Organization that was listed. It's the same list that wetook which is 85 Police Organizations list, ofthere are 85public and official police departments of course some of them are verified some of them are not verified, so we took this data and we took the data for the pages that were at least about averageageof3yearsbetween2010andApril2015,wherethedatacollectedandthere

were47,474wallpostsandstatusupdates.

In the wall post itself, wall post or the status updates I actually do status update from the Facebook page in my account. Wall posts are the once that are showing upon my wall from others who post on the page. And status updates from me, wall posts from the people interacting with me on the page.

(ReferSlideTime:04:26)



## Data Categorization

|  | Total Posts | w/ ≥ 1 Comment | P&C | C |
|---|---|---|---|---|
|  | 85,408 | 46,845 | 5,519 $P_{P\&C}$ | 41,326 $P_C$ |
|  | 47,474 | 24,984 | 17,196 $C_{P\&C}$ | 7,788 $C_C$ |

So,totalpoststhatweredonebyapolice.Theicons thatwillconsistentlygothroughthis police is for one that on the second row in this table and citizens is for the third row on this table. So, the total posts done by police were about 85,408, by the citizenswere about 47,474.This number would be same as in the last slide which is citizen walls posts people are posting on the page and the status updates which is police themselves the posting on the page. And if you looked at the post which had one or more comment is about 46,000. The posts which had one more comment from the citizens among the 47,000 is about 24,000.

So, if we see the interactions where police and citizens are interacting, the comments which is Pand C is when the police post, the way to read this table is row two which is thepolicepostandthecolumnnumber3whichisPandCmeansthatpoliceisposting

thecontentandthepoliceandcitizensbothareinteractingwhichisP,PandC.AndtheC is only citizens are interacting, which shows that the police <mark>posts,</mark> citizens are actually interacting with police more.

And when you look at total posts in terms of citizens which is 47,000 which has posts, comment, number of post which are more than one comment, one more comment is about 24,000, and the wayto read this is citizens are posting and then police and citizens arecommenting onitwhichisabout17,000andonlycitizenscommentingisabout 7000.

This just gives you sense of what the data set is, which is how much of interactions are happening, when police post with citizens and police, and how much of interactions are happening in citizen post, but police and citizens and citizens are interacting.

(ReferSlideTime:06:32)



Also the things that, the measures were done in terms looking at what analysis could be done <mark>were</mark> the topics engagement, emotional, social cognitive evaluation. These are the four things that I showed you guys research questions can be started this lecture. Topics, N Gram Analysis was done which is how frequently the words are actually appear in 1 gram would be the single word analysis, <mark>bigram</mark> would be 2 words appearing together.

K-means Clusters, clusters are basically way in when we see users together interactingon the same post more number of people, there is a cluster there are. We can actually look in at how much of clusters happen. Clusters for example, the students who are taking this class from anywhere in the country and let us take if you look at only the students who are taking it from one particular state that would be a cluster. Students taking this course from one state could be treated as a cluster, within the network of all the students taking this course from all around the country,that is the way I have lookedat clusters.

Number of police in citizen who comment in posts: distinct citizens who comment in posts, average number of likes and comments. These are the things we look at, next of couple of slides, we'll just go through in only weeks. So, valence is another way ofsaying about the positive and negative sentiment. Arousal is about intensity of the posts itself, how strong the post is. And then answers could be done in many ways.

In thecoursealsowewillgethereintroduceinthetutorial, wewillgetyou introduced toNLTK, which is a tool kit for doing some of these text, word analysis also. Here LIWC which is again, a tool, which could be used, takes input as text paragraph let us assumeor a post from Facebook, given this input to LIWC, LIWC will produce different categories of concepts that are appearing in the content that was given to the LIWC. Which means if you get a paragraph which as a lot of positive words, lot of happywords it will produce a category as happy which will be what we have.

So, it has being very frequently used in analyzing text content, but specifically within these categories that LIWC, I think that it has about 35 or 37 categories that has if youare interested in looking at the categories in the text which is posted, LIWC could actually give you that. And again similar there are ANEW dictionary also. Interpersonal focus, social orientation and cognition, again this is analyzed through LIWC content we shall walk you through now.

(ReferSlideTime:09:27)



First let us look at the topic characteristics. Four things we said, the first one is the topic that were being discussed. So focus on, so when police post again if you look at icon on the top it is police, when the police post are focusing on advisories and the status of different cases being investigated and information that the citizens will be interested in, that is what police is posting about.

If you look at the content analysis which is in police post interactions in police and citizens there is always discussions about rules, safety, violations, these are looking at only unigrams now. So you look at the concept of rules people, police posting aboutthese are the rules do not cut the signal, safety rules, wear helmet, violations and topics like that.

In terms of only citizens talking about it is actually following the content, notice, prosecuted, the interactions that citizens have, is mostly focused on those topics. Thisjustgivesyouasenseandthehighlightedisbasically totellyou howthedifferencesare, where in terms of the topics. Following, notice and prosecution versus rules, safety and violations; this gives you a sense of what kind of topics are discussed when police post and the interactions happen between police and citizens and between only citizens.

(ReferSlideTime:10:50)



Nowifyouflipthequestionandsee,whencitizensdothepostwhathappens.Mostposts request the police to a take action, which is why probably citizens are using these pages, whichiscitizenspostonthepageforaskingsomeaction.Ifyoulookatonlytheunigram

analysis,weseethatpleasetakeaction, veryevidently,pleasetaketheactionseemstobe the most frequent thing which will come when citizens post and the interactions between citizens and police.

In the citizens again if you take at least one, please, one, take, action, people consist unigrams that are done within the citizens. It is basically showing you that the interactions happens between police and citizens, when citizens posts for action, police and citizens are talking about, take actions on these and citizens among themselves are also talking about take action and please help and things like that. These two topic analysis just helps us to understand how the interactions or what type of topics are being discussed when citizens and police post on content, post on these Facebook pages.

(ReferSlideTime:12:11)



Now, when you look at the Clusters of Topic, here when police initiated discussions are more focused than citizen initiated, let us see how; and why we are actually making this conclusion which is, when police starts the discussion it is more specific about topics, awareness drive, safety campaigns which is specifically talking about. If you look at the categories that are going to be shown up for the police created discussions, it's actually very small where as if you look at citizens it is actually quite disperse.

Here in terms of police, awareness drive and safety campaigns, road senses, the offering of courtesy, and the parent of safety, things the posts that police do. Prosecuted and action report, action taken blah blah blah regarding your tweet petition and the details about it. Advisories on situations, good morning to all commuters of Shillong City, there is heavy movement over NH-40 to 44 and blah blah blah. So, the topics are pretty small.

(ReferSlideTime:13:24)



If you look at the citizens one, the topics are quiet diverse. Understandably again, it is very intuitive than the topics that are coming from citizens are much wider than the police itself. But the police initiated discussions are more focused than citizen initiated.

In the conclusion here is an appreciation, newspaper articles, citizen tips and complaints, neighborhood problems, missing people and this list goes up. We are appreciating the police, sharing some news articles that they would like to share it on these pages, citizen tips, complaints, driving in wrong side at blah blah blah, neighborhood problems,missing people. These are the buckets of topics that we saw in the post that were coming from citizens. So, you can clearly see the topics have been very diverse and the police discussions are much more focused.

Now let us move on to the next question. So the first problem, first saying to look forward is the topics that are being discussed on these Facebook pages between citizens and between police. This is very different from the 5.1 and 5.2 topics that our week lecture that we saw, which were we were interested in more like within the Bangalore City Police can we actually get some actionable information.

So with the questions, the intention for looking at these kinds of data is very different than these two, one is looking for content, actionable information here we're lookingwhat the interaction says, can be actually learn something specific. Engagement characteristics;howdothecitizensandpoliceengageinsocialmediadiscussionthreads?

(ReferSlideTime:15:15).



Now let us look at the engagement or comments characteristics which is, when police do the post how the comments are discussed, when citizens do a post how comments are being spread on Facebook. Ifyou lookat the firstset of columns which is the second and the third column, this is basically showing that when police does the post what happensto the comments if police and citizens are involved in the comments. So, it shows it's about 55,000 within the police and citizens group.

As you look at citizens, if citizens' post and then the comments are from police and citizens it is seems to be much much lower. It is about 26 percent lower in terms of the commentsontheFacebookpages.Ifthepostarecomingfromcitizensandthecomments are coming from police and citizens.

This probably is because when the police does the post there is much more interactions because the information that they are giving is probably more useful, probably more engaging in terms of actually the interactions on this platform. When police does and only citizens, citizens only interact it's higher and when citizens does only citizens interactions of the comments are also lower than the police side.

This slide is just a inside view of the data of the last slide. Last slide was this. This istotal number where as this is average number. This is average number of comments for the posts that we saw in the last slide, which is if the police and citizens are interacting there the average 3.3 comments are there for post, where as likes is about 9 likes, only citizens is about3.69 and about13.38 is on the average. So, the total gives you apicture, whereas when you look at an average it is actually the picture is slightly different which is C PC is actually 9.49 percent lower than the Cc which is the comments by police and the citizens versus comments by only citizens.

In the same way in the likes also the Cc which is comments by citizens only if the post comes from police is actually higher compared to the comments by police and citizens. This is per post average number of likes and average number of comments. In example which we kept there, a Citizen post: my family and I are getting the unwanted calls from blah blah blah blah. Police replied: dear please visit at your nearest police stationand give the details. Police suggest an appropriate action and the discussion tends to close early when the police are starting to interact. That is the reason why C P C is actually lower is because when policestarts looking at the post, police interacts they can actually easily or quickly close the loop of the discussion on the post.

Hope that make sense which is, the second question that we saw is engagement characteristics how police and citizens are engaging. So the kind of conclusions you can see is average number of comments and average number of likes when police and citizensarethereisactuallylower thanonlyacitizensbeingthere.Asimplereasoncould be police can actually quickly close in the loop of the discussion.

(ReferSlideTime:19:07)



Thenextoneweloookatisemotionalexchange,whichisnatureofemotionsandaffective expression that are being discussed on social media.

(ReferSlideTime:19:18)



As a one conclusion that we can look at is when citizens initiated threads there is always highernegativesentiment.HereisatablewhichactuallyconnectstothetopicsthatIsaid earlier which is negative effect, anxiety, anger and arousal. In the same one you can see thatnegativeaffectishigherwhenpoliceisactuallyinvolvedintheinteraction,whichis 0.021 versus 0.018. The reason for this would be that citizens are expressing their views about crimes, expressing the views about things that are going on and around them, even in the some of the graph seen before like neighborhood problems and they want to actually express the views so therefore the negative affect is higher.

(ReferSlideTime:20:04)



If you look at the next one where the anxiety is actually lower when police starts interactinginthese posts whichareinitiatedbythecitizens,youcanseethevaluetobe 0.001 versus Cc to be 0.003.And the reason could be is that police starting to interact in these post, police could actually complete the loop, as I said few slides before also completing the discussion and making the discussion closure which lets the anxiety ofthe citizens to go low.

And in the third row, we actually show you can anger is also higher in terms of C P C versus C C which is 0.006 to 0.005, where the anger is higher is again could be that citizens are actually expressing the views and they really want to get things done andthey are approaching the police to actually get things done as quickly as possible. And arousal is of course higher in C PC. Again when the citizens initiated threads, the C PC is higher, the arousal because the intensive discussion is higher because you want actually get things done from police again.

Now let us look at the last part of this <mark>deck of</mark> slides in terms of actually the questions thatwestartedwith.Westartedthefirstquestiontobetopicalcharacteristics,secondone to betheengagementcharacteristics, third oneto beemotional exchanges, and thefourth one to be cognitive and social orientation. Here the question that we are trying to post was; what are the linguistic attributes or characterize cognitive in social response process? I think this is a very broad and question that people could study in multiple ways.

Theoneapproachthatwecoulditwaslookingatthepropernounsi'sandhe'sandthey's,       words used in the post that were initiated by the citizens. And we see that most of the post that the citizens initiate, are actually self driven or mentions more of themselves. You can clearly see 'i' being the highest in this table. An example for the post is, 'I have lived in the UK and all the time I have never heard anyone honking blah blah blah'. So, these kinds of post are mostly presented in the Facebook data that we collect.

Now let us look at why all this matters. For the entire week for this course, this week I kind of took you around the data that are created on Facebook, that are created on social media for Police Organizations doing some analysis, seeing what kind of interactions that are going on between police and citizens, why does it all matter. One it definitely tells police improve policing and community sensing, which is to collect data from these kind of social networks can be used to record and sense behavioral attributes, such as engagement, emotions, social support. Which is what we are trying to capture in 5.3. And it will enable police incidence community to enhance emotional support to residents experiencing safety issues also.

These are the thing that we talked about again in the analysis. Discussion thread with police and citizens where gives the level of anxiety when police talks interacting right, a couples of slides back showing you about C P C, in C P C the value for anxiety is actually lower is just because that they, when police staffs interacting the citizens tend to be more, anxious level being it is lower or it tends to be lower because police starts actually answering some of the questions that people are asking.

For example, I say that there is some problem here by in my neighborhood and police comes and says, Oh good, thank you for letting us know about it and we will actually do

the appropriate thing forwarding this to the right police officer all that. So, there is some level of satisfaction to the citizens while the CP responses.

(ReferSlideTime:24:02)



Interestingly we can actually develop some technologies also, helping communities to make consensus based decisions regarding support and actions they seek from police. Theycan actuallybuild technologies which will help these kinds of interactions and help decisions also to be made which can help have more safer cities, help our more interactions between police and citizens. It can help actually understand the change in emotions of people also.

For example, I manage the police page today, I see what kind of views that people have about a particular problem today and what kind of content they are generating today which is from my jurisdiction and month later what happens, a year later what happens, understanding this will actually help in terms of making some predictive analytics solutions also. Predictive analytics around this data and actually helping Police Organizations make their job better, make cities more safer.

Itcanalsohelpsenseandtherecordthereactionsofcitizensandsharetheserecordswith decisionmakers.Itcanactuallyget,meaningsomeorganizationsalsocandothis.

Publicly saying that these are the number of posts that we got this month this week and this is what we get. It can actually act as an early warning system; it can act as a predictive usage of this data, it can use for predictive analysis also. Essentially all of this can help both societal impact in terms of police and citizens interacting better. Technologies can be built to help police make their decision better; it can help citizens have safer lives.

With that I will stop the content for week 5. So, essentially week 5 is totally about policing. And we looked at initially how these organizations create, what kind of pages do they have, and what content are getting generated on this pages, from there we wenton to ask some research questions are on that. And these kind of research is not onlydone in terms of just Indianwebpages, Indian Police Organizations, this is done allacross the world.

Thankyou.

# eCrime on Online Social Media Part-I

WelcomebacktothecoursePrivacyandSecurityinOnlineSocialMedia,thisisweek6.

(ReferSlideTime:00:17)



So, what we have seen until now is generally, overview of online socialmedia. We have had a lot of hands on tutorials about Linux, Python, Twitter API, Mongo DB, MySQL and then I went into topics like Trust and Credibility, then we saw Privacy, last week we saw what is Policing how online social media is being used by police organizations specifically in India and what research problems, what questions that you can actually studyfromthe data that you collect fromthese socialmedia services.

(ReferSlideTime:00:56)



Let me just quickly tell you what we saw. Multiple police organizations have actually adopted using Facebook, Twitter, for sharing for interacting with the citizens and that is the topic that we saw inthe context ofpolicing.Aspecific questionthat we saw washow we canactuallyuse this data fromsocialmedia to collect actionable information.

(ReferSlideTime:01:19)

Isitpossible <mark>wecancollect</mark> someactionableinformation?Isitpossibletousethis information for making any interesting judgments?

In this context we also saw that how we can use the text content that is posted on these social media services to take some actionable information. For example, <mark>this tree shows</mark> how you can understand the needs of citizens who are posting on these networks. Likefor example it says, need to be punished, need to hang this guy. So, these are the needs from the citizens who are posting this content on Facebook or Twitter or other social media services.

(ReferSlideTime:02:10)



Then you can also look at understanding wants, which is what is that citizens are interested in wanting from police, this we like want to hear more of these, here want to see the punishment of such people, want to saythanks to BCPSir. This is something we saw earlier and we are just going to quicklybrush it only.

(ReferSlideTime:02:30)

Thefour questionsthat wesaw specificallywhere,topicalcharacteristics,what topicsare being discussed, how the engagement between police and citizens are happening, what emotional exchanges are happening between citizens and police, specifically we also looked at arousal, violence and topics around that. Finally, we looked at cognitive and socialorientation, linguistic attributes, unigram, and bigram, and topics around that.

(ReferSlideTime:03:04)



In this data specifically what we saw was collecting data from 85 publicly and official departments between this period of a 2010 and 2015, the analysis was done on 47,474 wall posts and 85,000 status updates.

And of course, the technicalimplications ofdoing allthis is helping communities to help thepoliceorganizations, buildtechnologieswhichcanbeusedbycitizenstointeractwith policebetter,buildtechnologiesthatpolicecanuseforinteractingwithcitizens<mark>better</mark>and making the society <mark>a safer</mark> place to live.

So, that is the broader goal of studying these concepts on social media. As I said in the last week also I would really like to see people talk about their citypolice organizations andinteractionsifany,ontheforum,I <mark>havenot</mark> seenanythingmuchuntilnow,butIthink fornowmanyofyoumayjust understandingthecontent<mark>thatisjustthe</mark>contentitself.But it is actually great to see some going to see some interesting questions, students are asking in the forum.

So, what we willdo now is we willmove onto another topic fromhere. The topic now I want to look at is e-crime; e-crime, cyber crime anything that isaround electroniccrime, but focus it only on the social media context. Crimes happen allaround the places usingthe internet, using the web, but people will focus on these kind of crimes only that is happening on social media.

And asthe patternonthecourse,wewilldo somebasicsnowinthefirst part ofthisweek then I willget into some research questions or questions that one could answer using the data that is been collected on crimes from these social media services itself. We will do these hands on tutorials also, which is looking at social network analysis tools and then NLTK. And there are other hands on tutorials also that are we have planned over the course of next few weeks.

(ReferSlideTime:05:21)



So here is a list of not a comprehensive list, here is a list of crimes that I thought I will cover before getting into details of any one particular topic. We are going to look at one or two topics in detail, but before that let me just walk you through some crimes that happenonononlinesocialmedia.Iamsuresomeofthiswearealreadyawarebutletmejust brush it to get your sense of what the crimes that are going on in these social media services.

The first one which is phishing, and again these are not arranged in any particular order andtheyarenot comprehensiveatall.Thephishingproblemonsocialmediaservices,the act oftrickingsomeoneto into handling orlogging detailswhichisbasicallythereisa,in traditionalwaysinemails,inemaildomainsyougetemailswhichsayspleaseclickonthislink or please click onthe links to update a password or your account is expired click on this link to activate your account.

Whenyouclickonthislinkyouare taken to afakewebsitewhichsometimeslookslikea legitimatewebsite, but sometimesit doesnot needto belooking likealegitimatewebsite also. And when you go there it is asking for username, password and when you give the username, password, you are basicallysharing the credentials to someone else.

And these kind of emailshave been playing around for a long time and there are many sophisticated attacks that has happened using these emails; phishing itself, just phishing, sphere phishing. Sphere phishing is a way by which you target a set of people. For example, in this course I could just target only the people who are taking this course saying as though it is email coming from PK at IIIT, saying please click on this link to know further links that I have actually posted on the web about the course.And then, of course,someof youmay beinterestedin whatIam speakingaboutthecourseandyou

willclickonthelink,butitisnotactuallyalegitimateemailor alegitimatelink.Sothatis about sphere phishing, but then there are other types of phishing also, which is whaling where the specific CEO's of a company are targeted while sending out these phishing emails. There are many different types of phishing attacks that have been going on. So, that istraditional.

But now when you move on to the socialnetwork, these attacks have also calculated the socialmediaservicesalso. For example, alinkontotheTwitter timelinewilltellyouthat please click onthislinkto get somemoneyandthenwhenyouclickonthislinkorplease clickonthislinktochangetheFacebookpasswordthatyou havecreated,therewassome problemin your access with Facebook; click on this link to update the password.

As in the traditionalwayifyou click onthis link you willend up actuallygoing to a fake websiteand giving awaythecredentials,that isphishingandImeanyoucanthinkofit as a phishing as in the traditional ways in emails itself, but spreading on the social media services. There could be a linkonFacebook, there couldbealinkonTwitter;therecould be anemailto say, please click onthis image to get some more informationabout atopic and it could actuallytake you to a fake website. So that is phishing.

(ReferSlideTime:08:44)

So, specifically the examples in phishing that are going on now or have been around for sometime is Facebook technicalsupport sent you a notificationsaying that, there issome probleminyouraccountpleasego verify.Facebooknewloginsystemthatisemailsgoing aroundwhichsaysthatFacebookhasinventedanewloginsystemandclickonthislinkto create your account on this new login system or merge this account to the Facebook account and things like that, these emails have been going around.

Andifyoureallylookat it,Facebookcredentialsarebecomingmoreandmoreimportant, because if I know your Facebook credentials I actually get to know your friends, I actuallyget to know your pattern ofusage, interest and topicsthat you maybe interested in spending time. These things can be used against you. So that is the reason why Facebook credentials are also becoming more and more popular, compared to the email address,compareto thefinancialaccount detailsthat one wasalso chasingbefore.Thatis phishing.

(ReferSlideTime:09:50)



Let me walk you through some fake things that are going on online social media also. Here isone whichisfake customer service accountswhichis, IhaveaproblemIactually post a tweet saying I have problem with this bank. For example in this case everytime I havebeenonmybank'swebsitelately,ithasnotbeenworking,frustrating.Theyare

actually tagging the right bank; they are tagging the right organization. For example, we could actually think of the same thing tagging HDFC Bank, ICICI Bank. These kinds of organizations have legitimate accounts and people using Twitter can tag them.

So, what happens now? This is a real tweet and real customer asking for real problem. What the fraudsters do is they look at this tweet, they have mechanisms to figure out these kind of tweets are going on. They actually reply to these tweets as though it is the bankwhichisreplyingto thistweet.And theywillcreateaccountswhichareverysimilar to the real account and reply to the post as though the real account, real organization is actually talking to you. In this example the Usual Studio Dear Charlee, We sincerely apologize for this - loginto your account via secure signonchannelblah blahblah.

This is the customer service account; fake things that is going on, whichis realcustomer tagging or connecting to a real bank organization or an organization. The fraudsters createaccountswhichareverycloseto therealaccount andactuallystart interactingwith the customer.Theseisfakecustomerserviceaccount problem.Thisparticularexampleis on Twitter, but one could think of such problems being on all social networks also. Because all of these legitimate organizations are actually using social media to interact with their customers.

(ReferSlideTime:11:48)

Here is the second one, fake comments on popular post. I think some of the posts thatthat become very popular. Let us take the prime minister was talking about it, if it was Obamawho istalkingabout somethingthesepostsbecomeverypopular.Andwhenthese posts become verypopular there arealso lot ofcomments. For example, nowI amsureif you look at the Olympics Facebook page or the twitter handle or the hashtag, people are actually talking a lot about things that are going on in the Olympics in the context of Facebook page and Twitter accounts also.

Therefore, what scammers do is that they actually pretend to be Facebook users so they can comment on this. For example, I could create an account, I could create Facebook account which looks very legitimate, I can start posting on these Olympics relevant post which are very popular and I will kick you from that to a fake website, and get your information.And if you too click on this link I will give you also down somemalware inyour code and things like that.

So that is a second type of a fake thing that is going on, which is fake comments on popular posts because the reason why it is popular post is that it gets in more and more fashion, and more people actuallyget to see it, it isconnected to the topic that people are more interested on.

(ReferSlideTime:13:12)

The third one is fake live streaming videos, which is particularly in the context of Olympics and cricket matches, world cups and things like that, there is tendency of actually looking for these matches in live. Here is an example where this post is actually saying live video for thismatch, right?Ifyou are interested inwatching it inyour laptop, in your phone you tend to actually look at these pages, look at these links which talks about this game and tend to actuallytaking into a fake website.

Fake live streaming videos, which is there is no video, there is no real video which is connected, but the scammers actually tend to take the users to fake things. And they do this in the context ofsome games that are going on, some eventsthat are going on, some shows that are going on. For example, currently in terms of Rio somebody says that currently India has won medal and here is the video of the match. So, that is the kind of scamthat is going on in the context oflive streaming videos.

(ReferSlideTime:14:27)



The next one is fake online discounts which is, scammers take the real account, real organization in this case - Netflix, it could be anything Facebook, it could be Flipkart, it couldbeanyrealorganization.Theycreatefakeaccountsthat lookslike realbusinessand they are actually carry out business using these fake list, but giving you discounts. Like forexampleNetflixcouldsaythat,thispagewhichisafakepage,itcouldsaythatthereis

a 10 percent discount in Netflex account that you open now. 40 percent discount for the next 6 months, if you open the account right now. These kind of posts can actually lure people into using these fake accounts, fake pages, fake services. So, that is the next fake crime that I thought we will talk about.

(ReferSlideTime:15:25)



Next type is, Fake Online Surveys and Contests. These kind of scams have been around foralong,longtime,wherethecriminalsofthesescammersget youto get survey,fillthe surveyto get somemoney, to getsomeinformation.Forexample,howdoyouknowyour personality? Personality test and find out other people who are bonding your date, who hasthesamepersonalityand thingslikethat, whilethesekindofthingshavebeenaround for a long time.And there were also contests, win1000 Rupees for filling onthis survey. So these have being there in traditional ways now these have moved on to the social media services. Here is an example where, what is your opinion, we would liketo know, participate in our research surveys and enter to win prizes, here is the link.

Again this is a fake claim, this could actually be malicious, and this could actually be collecting personal information. But the source of starting this is getting you to click on the link a survey or contest. So, that is the last cyber fake version that I thought I would actuallymentionittoyou.Quicklyafewfakecrimesthatyoucanthinkof-fake

customer service account, fake comments on popular post, fake live streaming videos, fake online discounts, and fake online surveys and contests. So, these are the different types ofcrimes that can be go scams that can be happen on socialmedia services.

(ReferSlideTime:17:01)



Hereisfewmore,ImeanIthinkifyoulookatmyfirstslidewhereIshowedyoudifferent types of fake crime things I was going to talk about, fake was that part. Now here is another one which is a Fake Tip: Foursquare is the most popular location based social network. Inthisfoursquarefor exampleI couldactuallywalkinto IIITandthensaythat I have checked into IIITDelhi. So that isthecheckin, andyoucanalso leaveatip,I go to Saravana Bhavan. I eat food at Saravana Bhavan and I say that the food is pretty good. So, in that tip, people actually, the scammers and the criminals actuallypost information that can take you to a fake website.

For example, here by the original XanGo and mangosteen juice at best price, this link. This is the tip posted on particular location, so it is taking you to link which could be actually phishing so, it is also studying, giving you information, advertising about such certainproductd. So that is the fake tip, the informationthat are posted ina tip that isnot relevant to that particular venue, andtaking youto afakecontent istheproblemhere.So, that is the fake account foursquare.

(ReferSlideTime:18:23)



Social reputation has become such a big deal now, everybody talks about I have 2.5 million users and then 2.5 million followers, then the number of likes that you have on your page is becoming the way that people measure your influence in the society. Even among friends, it does not have to be the celebrities, politicians, evenamong friends you are more, more curious about how many friends other have. The social status is now being measured by the presence in social media; by the number of likes that you get on posts, number offriends that you have, it is becoming more and more popular.

For example, Facebook likes andAmazon reviews,YouTube likes, the endorsement that happens on LinkedIn where you are endorsed for a particular topic, how many people have endorsed you, what topics have you been endorsed. These are becoming a measure ofinfluenceinthesociety,numberoftweets,numberoffollowers,numberoffollowings, all of them become a measure by which people think of your social reputation. But the problemisallofthemalsohaveproblems,becauseofthesewaysbywhichcreatingsocial reputation has happened, you can actuallymanipulate these socialreputation also.

In this case, some examples that I had put in here is Flipkart, social reputation can be manipulated by actually writing good reviews about product. So, reviews become a big way by which you can actually manipulate the social reputation of the product, of the company, of the seller, all of them can actually be manipulated. It is actually a very big problemintermsofstudyingAmazon's reviewsor Flipkart'sreviewsalso for products.

Here is a case;Amazon sues 1000 people over fake reviews. People have been studying withreviewsproblemforalong,longtime.Itisnot onlyreviews,itisalsoaboutstudying the fake followers, studying the fake endorsements, all of them are actually relevant problems. If anybody is interested in taking up some of these, these are actually very interestingproblems,verychallengingproblemsalso andveryrealworldproblemswhich is, you canactuallylook at the solutions that you build, becoming /influencingpeople's thinking.

Here is another problem in terms of crimes on socialmedia. Clickbaiting, where you are actual director your keeping the website, so you go read a particular page of news or something,theretheypresentyouwithinformationwhichissometimesrelevantsometime    not relevant and they take you to a fake website. So, here in this case also, the link here, this information was actually presented in one of the social media services where it was taking it to a fake website. Clickbaiting - getting you to click on links which are not legitimate.

Hashtaghijacking;hashtaghijackingisalsobecomingabigissuethesedays,Iassumeall of you know what a hashtag is. Hashtag is the waybywhich a particular set oftweets, if you want to talk about now Olympics you use hashtag Rio 2016. So that is the way of using hashtag Rio 2016, you are saying that the content that I amposting isconnected to this topic, so Twitter canactuallybring inallthesepostswhichhashashtagRio 2016and show it to people who are interested init. So, that is the logic behind using a hashtag.

InthisexamplewhereCocoColahasactuallyposted tweet whichsays, 'TimeforaRoyal Celebration hashtag Royalbaby'. Here what coke is doing is, coke is actually using a hashtag which is very popular otherwise for actually selling their product. That is hijacking right, royal baby is nothing relevant to coke. They are kind of using it to promote their products. So that is one wayof hijacking the hashtag.

Here is another example also. WhyI stayed was a hashtag that was trending, was getting popular so this pizza.DiGiorno Pizza thought of using this popular hash tag actually to sell to mention about pizza. They used this hashtag why I stayed because you had the pizza, but unfortunately this also back fire, here is the post that they had to actually apologize for doing this post. A million apologies, did not read what the hashtag was about. The hashtag was actually used in some of the context where people were actually using thishashtag talk about a particular situation.Therefore, taking thehashtagwhichis not relevant to thistopic, using it for selling a product isactuallyhijacking.

Nowjust furthertalkingabout forexample,youwouldsaysomethingabout whatyouare doing now with the hashtag Rio 2016 which will actually show up on people who are looking at timeline for the posts which has Rio 2016. So that is the problem in with hashtag hijacking.

(ReferSlideTime:23:49)



Compromised account, I have actually shown this particular tweet even in my trust and credibility section, but I brought this back just to tell you different problem, I think I explained the problem then but I will explain it in the context of e-crimes also. Compromisedaccount, whereTheAssociatedPressisaverifiedaccount andthisaccount was compromised for sometime which is, somebody else had access to this account and the tweet was, Breaking: Two Explosions in the White House and Barack Obama is injured'I amsure you can allagree that the effect this tweet must have had.

This is account compromised, somebody else getting access to your account because of leak of username password and getting that to misuse, getting the account to be misused also. That is compromised account.

Impersonation: Impersonationisalso another problemwhichisIcantakeanaccount like for example, any of you in the class I can take some details of you that I know pictures, and your city, and the information that I could collate from online sources, use that to actually create an account which as though looks like it is you. Here is a complaint that Kiran Rao has actually filed saying that fake account has been created, and there are many, many fake accounts like this. If you know remember the policing section I also showed you about the fake account of police organizations also.And it is not just about individuals, even organization's accounts are actuallycreated fake.

Here is another interesting problem which is, Work from home scam.Againthese things have been in traditional ways for example, if you are driving down somewhere in the signal, you will see a post which says, 'want to won 1000 Rupees a day sitting at home pleasecallthisnumber'thesekindofscamsarebeingthere.Hereisanexampleofascam that went popular in Pinterest where this image was actually floating around, 'want to makeanextrasalarysimplybyfillingoutsurveyformajorcompanies,hereisawebsiteto go to.Youget paid5to 40dollarsper survey.Thisisworkfromhomescam.Againthere isalot ofscamswhicharesimilarto theseworkfromhomescams,different versionsthat are verypopular on socialnetworks. So, this is an important scamalso.

With that I will actually wrap up my first part of the week 6, where I thought I will just introduce you to different scams, different crimes, because we will talk about crimes in this week, looking at different crimes some data was collected, what kind of analysis could be done, what kind of solutions that we could build in reducing these problems of crimes on social networks.

# eCrime on Online Social Media Part-II.

Once again welcome back. This is Privacy and Security in Online Social Media, week 6 the second part.

(ReferSlideTime:00:18)



Also inthe first part, wegenerallysawabout whate-crimesare, different typesofcrimes on social networks, specific examples about different crimes; how that affects yor social reputation? How malicious users are actually making use of these social network interactions and topics are answered.

So, what we will see now is some specific problem in socialnetworks, crimes and issues on social networks and we will takesome one data set and answer some questions with that data set. Also search engines rank websites basically the Pagerank idea where every page is linked to everypage and Pagerank ofthe rank ofeverypage increases depending on the links that it has with the pages and. So, essentially if you have more of high in-degree helps in increasing the Pagerank.

So, Googleworksonthistechnology, where youactuallyhave;youcreateawebsite, you link it to, let us take to ==IIITD's== website and ==IIITD== links it back to you, then I think your Pagerank increases heavily, so that is simple idea for Pagerank, but link forming in onthe web is basically an idea where websites exchange reciprocal lengths with other==sites== to improve the rank. So, the idea of making the links between websites which is not otherwise there; creating links or increasing the links ofthe websites to other websites is actually link farming.

Same ideas are connected to Pagerank. ==Pagerank is benign or legitimate links that youcreate. link farming is the idea in which these links are created which are not benignones.==

(ReferSlideTime:02:12)



So,hereisasimplediagramto showthat what,how link farmingorwhat link farming is. A link farm is a form of spamming the index of the search engine which is essentially increasing the links between different websites like, for example, website and website A and B. Allthe, if you start creating links betweenthese websites, iftheydo not exist and that is called actually link farming. Sometimes, it is also called spamdexing and spamexing. So, that is the idea for link farming. Link farming is a way which non legitimate links are created betweenthe websites. The idea for doing this is when you do this and when you increase the in-degree which is the links that are coming into the website increases then the Pagerank of the website automatically increases.

(ReferSlideTime:03:04)



So, whylink farming inTwitter?So, whatwearegoingto studyis, wearegoingtostudy the ideaoflink farmingspecificallyonlyinthecontext ofTwitter.So,whylink farming? What does it help? Who benefits because of actually link farming on Twitter? So, essentially Twitter, I mean given its nature, given amount of data that is gettinggenerated on Twitter, it is basically a web within the web.

If you want to go to, go and look at the out breaking news, Twitter is the place to start with. It has large amount of data on real time news. There is multiple research done to show that Twitter is where news breaks. If you want to look at the latest in now, year 2016, Twitter is probably the best place to look at and people start searchingfor a topicin Twitter actuallymeaning Twitter definitelyis a micro blogging website, where people push content, but if you look at the pattern in which Twitter is being used, it also being used for search for a topic, live search into people.

So, when you search for atopic inTwitter,thewaythat theresultsarepresented depends on many factors. So, search engines in Twitter can rank, actually follow ranks. It can actually present the results depending on the connections that you have with the person who is talking about that topic like, for example, ifI search for a topic like hashtag year 2016. If any of my friends are talking, it could show up on top. I mean itcould show up onthenumber offollowersthat people have. Ifit isa verified account, it should showup on top.

So, the search results can be; search results can be used, these kinds of techniques. PageRank would be one, which is how many people are actually connected to this particular tweet and who are the users, who are connected to this particular tweet. All of this information can be used to actually decide on presenting the search results. Of course, the way that the search results are presented is actually going to bias the users to goto, iftwitter is showing youtheresultsontop5there is more likelythat youaregoing to actually go look at those particular tweets.

And of course, high in-degree which I think when mentioned the part Pagerank, I said, high in-degree which has number of followers seen as a matter of influence on. number of followers that you have is a measure of social reputation. I also mentioned this in the last part of the week 6. There is a score called Klout. So, I mean I would recommend again you people look at what Klout score is, Klout score is essentially a way by which Klout collates all your online presence, particularly in the social media, and gives numbers toit. For example, my score would be 24, which basically says that on scale of1 to 100 what kind of influence are you having on the social media services.

Klout is an interesting mechanism that also, a paper which talks about how Kloutactually finds out these values and researchers have used Klout score as a way tomeasure the influence of the users also. Of course, the topic of influencer by itself is actually hard because you are defining who an influencer is; it is becoming more and more difficult. So, while link farming in twitter is basically a large amount of data is getting generated, real time information is spread there and when users search for topic, the information is actuallypresented depending onthe links, depending onthe Pagerank, depending on the links that follow a rank depending on the links that users are.

And particularly link farming in Twitter is basically, spammers follow other users and attempt to get themto follow back also. Essentially, how do they increase the in-degree, the in-degree is increased if I am a spammer, I start following thousands and thousandsof people and there is a probability that you will actually; the people that I am trying to follow, now will actually follow me back. And again, there is multiple researchers, people who have shown, how the reciprocity can be, there is a high probability that if I follow you, you will follow me back and, giving, with that effect, the link farming actually increases on increases and therefore, twitter can be used to increase the linkfarm.

So, here is a slide to show the differences and similarities of link farming in web and Twitter. In the web increasing my Pagerank, increasing my in-degree, increases my probability of showing up in the search results. In the Twitter space, increasing the in-degree actually increases the gain; similarly, to show on mytweets onthe search results.

In the webs, spammers actually use link farming. In Twitter, spammers do actually link farming, but it is also done bylegitimate and popular users, I think that is the whole idea withwhere youactuallyincreasethe in-degree bymaking your numberoffollowers high and therefore, you can actually your content can actually be presented to a large number largeset ofusers. And ofcourse, inthecontext ofTwitter, inthecontext ofweb, it is not necessary that if I link your website you are probably going to link back to my website;hyperlinks are not created in that way.

Whereas in the context of Twitter there is a high probability that, let us take if I am, I actually follow one of the students who are taking this class, there is a high probability that the studentis going to follow me back again and the same way if I follow aprofessor and the professor probably there is a high probability there the professor will follow me back.

I thought I will walkyou through some literature in the context of spam in Twitter, which isin this case I do not know I think I am going to talk about 2 or 3 research results. Just to tell you, the context of where link farming is going to be kept which is spam in a broader sense. So spam campaigns, here is a paper which is titles 'SuspendedAccounts inRetrospect:AnanalysisofTwitter Spam'. So, 5spamcampaignscontrolling 145 thousand accounts combined are able to persist for months at a time. So, here is a zoomed in version ofthe abstract which reads as, we identify about 1.1 million accounts suspended by Twitter for disruptive activities over the course of 7 months. In theprocess, we collect a dataset of 1.8 billion tweets and 80 million of which belongs to spam accounts.

The problem with that comes with the spam is that it is not only high in these social networks, but thereareactuallytheyalso stayfor longer. So, here it is7 spamcampaigns controlling145000accountsandtheypersists formultiple monthsandtheanotherpartof the abstract reads as, our results show that 77 percent of spam accounts identified by Twitter aresuspended within adayoftheir first tweet.So,thesearesomenumbers, some idea to get a sense of what is happening in Twitter in the context of spam.

Here is another one and this paper is also a popular paper in the context of spam on Twitter. So,this is '@spam: TheUndergroundon140 Charactersor Less'. So, whatthey found is, they found 8 percent of the 25 million URLs posted on the site pointing to phishing, malware and scams. So, that is a lot of URLs which are actually malicious, 8 percent of 25 million URLs, where they are pointing to malware, phishing, scams and malicious ones and they have also found that the click rate is actually higher.

So, they say that, we find that Twitter is a highly successful platform for coercing usersto visit spam pages with a click through rate of 0.13 percent compared to, much lower ratespreviouslyreported for email spam. So, the idea isthat if you get anemailofwhich is which says that please click on this line for buying a product with 10 percent discount that is low probability of you clicking on this link at the email, whereas if the same postis coming fromsomebodywhomyou are following onTwitterthere is a highprobability that you are going to actually check this up. So, that is the probably why this rate is actually high in the context of social networks because it is these posts are coming from your friends.

You remember the Associated Press example that we talked about earlier, where they, where I showed that the post had mentions about White House blast and then they cost actually, why because it is actually coming from Associated Press and there are many peopleactuallyfollowitanditisalsoverifiedaccount.So, thatis why theclickrateis

highand astheclickratesare higher, it willactuallybecome moreand moreasuccessful spam campaign.



Here is the third one; third research which shows the detecting, which is titled as, 'Detecting and Analyzing Automated Activity on Twitter'. So, what it shows is that 16 percent ofactive accountsexhibit a high degree of automation. There are again, there are multiple people working on the space in terms of actually identifying automated post on social networks, particularly on Twitter. One simple technique that people try and researchers have tried and it is also being used in some other products is to actually look at the frequency of the post that somebody does, as a human being, you and I probably will not be posting 50 tweets or 45 tweets and, whereas a bot, an automated service would actually do that.

So, people do the graph of hour of the day and minute of the hour; x axis can be the minute ofthe hour, yaxis can be the hour ofthe dayand if you draw the plot, iftheyare very, very close to each other then there is a high probability that it is actually aautomated service that is did this post.

16 percent ofthe active accounts exhibit a highdegree ofco-ordination. Theyalso found that the11 percent ofaccountsthat appearto publishexclusivelythroughthebrowser are in fact, they are automated accounts that spoof the source of the updates, that is also interesting right. Now, it is not only just the content which is spammed, it is also the spoofing of the source, how the post was done.

(ReferSlideTime:14:16)



So, what we are going to look at inthe specific questions that we are going to answer is, we are going to actually look at research that was done which is using the entire data set of Twitter, which was collected in 2009. It has 54 million users, I think it will be extremely hard to collect its data today because of the number of users, the connections and probably also the infrastructurethat you may need to collect this data.

So, this 2009, 54 million users, 1.9 billion links between the users and it is probably one of the largest set on data largest data set on Twitter.

(ReferSlideTime:14:54)

So, how the definition of the follower and followings is used in this context is, A and B, if there is a link between A and B and going from B to A, they are marked as towards that, then B is, A is B's following and B is A's follower. The error marks towards it which isontheside. So,this isB, which isthe followerofAand Aisthe followingofB, I think wetalked about what Twitter is, basic terminologies, very early in the course. So, thatisfollowerand following for you. Andoverthegraphatthebottomtalksabout spam targets which is the targets where spam is going to be sent, targeted follower and spam follower.

So, B and C is basically showing you the spam followers which are which are the ones that are going to be following S, and A and B are the spam targets they are going to be getting the spam from S. So, the terms are going to be follower, following, spam targets wheretime is going to bespent, sent, spamfollowers, thosearetheoneswhich aregoing to be following. The base which is B and C or its equivalent to B and first part of the graph spam followers and of course, targeted follower and that makes sense.

(ReferSlideTime:16:35)



So, what they formed was they formed 379,340 accounts that has been suspended in the interval of this period August 2009 to February 2011, spam activity of course, these accounts were suspended because there was a high spam activity and login activity becauseif you donotlogin toyour accountfor sometime, Twitter can actually suspend

your account. 41,352 suspended accounts posted at least one blacklisted URL shortened by bitly or tinyurl. So, there is a set of URL shorteners called bitly, tinyurl.

Some of you may have used it, if not please go look at them. The idea for the URL shorteners is, if you want to share a long URL, particularly because of the social networks' influence, social networks' growth, these kinds of URL shorteners have actuallybecome very, verypopular becausewhenIwantto do apost inTwitter, which is only 140 characters, I do not want to really spend a lot of a space in just posting the URL, insteadIwouldactuallysend it totheURLshortener, whichwillreducethe link, if it was like 100 characters it will just give me into bitly, bitly dot com slash some 6 or 8 unique characters which would redirect me to the actual website. So, 41 thousand suspended accounts posted at least one URL, which was actually shortened.

(ReferSlideTime:18:02)



Let us look at just the spam. So, if you remember the terminology spam targets, spam followers, entire followers. So, number of spam targets followers that the graph at the bottom, actually the Venn diagram at the bottom shows you spam targets were about 13 million, targeted followers, were about 1 million and the spam followers were about 248 thousand, 82 percent of the spam followers overlap with the spam targets alright, whichis the followers, spam followers, who are going to be following some accounts are actually part ofthe spam targets itself, 82 percent ofthe spam followers overlap with the spam targets.

So, what isthegoodway?So,here isan interesting wayofactually looking at it thedata. So, this is cumulative, CDF, which actually shows you node rank at the x axis and cumulative number ofspammers inthe yaxis. There is some interesting conclusions thatyou can actually draw from this, which is to say that the number of spammers who rank within the top k according to the Pagerank. This is ranked according to the rank of the user. So,ifyou seewithinthe first 10,000usersthereareactually7spammers, what does this mean? This means that, if you would actually list down rank of all the users of Twitter, look at the followers, look at the node rank which is the Pagerank, in-degrees of those followers,

you can actually see 7 spammers within the top 10,000 users and 304 within the 100,000 and 2131 within the first 1 million users, which is if you take 10 lakh users, the top 10 lakh users with PageRank, with the in-degree as high, list them, you will see that two thousand users of this 10 lakh users are actually spammers alright. So, that gives you sense of you know the spammers are actually very popular alright. It essentially shows that these users have high in-degree which isthe context ofthe problemthat wetrying to study which is link farming.

(ReferSlideTime:20:30)



With that, I will actually stop the second part ofthe week 6. I will continue with the rest
of the results and analysis soon.

.

**Unit-4**

**Link Farming in Online Social Media**

Welcome back to the course Privacy and Security in Online Social Media. So, this is week7, Ihopeyou gotachance tolookatthe content inweek6where welookedat link farming spam in Twitter and some kind of work which was able to find out what link farmers are what is the characteristic of link farmers. So, today we will continue, this week we will continue a little bit about the same topic, we will finish it and move onto something else.

(ReferSlideTime:00:44)



If you remember last week, I showed you about what is link farming and somedata about link farming. So, here is a graph which has on x axis spam follower node rank whichiswhatistheprobabilitythat spamfollowersaretheaccountswhichactually follow spam.IfyourememberthegraphthatIshowedyoua,b,c,d,e,fwheretherewas something as spam follower and then something that was spam followings.

So, x axis is spam follower node rank which is the rank of the account, which is being a spamfollower and then onthe yaxis is fraction reciprocated in links, which is I think we had briefly mentioned this last week also, the probability of you following me when I actually follow you, reciprocity so to say.If I follow you what is the probability that you will follow me back that is what is put on the y axis. So, fraction of reciprocated in links from spammers versus spam follower node ranks.

(ReferSlideTime:01:56)



That is if you look at the top 100,000 spam follower accounts for 60 percent of all links acquired by the spammers. So, what does this mean? This means top 100,000 spam follower accounts for 60 percent of all links acquired by the spammers. So, if there were 100 links that were created for the spammers, 60 percent of them are coming from the spam followers, which is an account which actually follows back the spammers. Top spam followers tend to reciprocate all links established to them by spammers.

So, if you look at this graph the way to read this graph is on the x axis is the log scale, which is 1, 10, 100, 1000, 10000, 100000 and 1 million. So, that is how it is written on the x axis. y axis is the probability of you following me back, if I follow you. So, if you look at the first let us take a look at the first value 1 or 1 to 10. So, here the probabilities isalmostclosetoonewhichisthetoprankedspamfollowerswhichistheaccounts

which are having the chances of following you back is very high. If you ==arrange== the spam followers in the node rank which is the number of followers that they may have is actually very high.

So, what ==does== this mean? This basically means that the probability of a spam follower following a spammer is very high, which is what spammers actually make use of, which is if there is a probability, if there is a chance that you will follow me back, the spammers will keep following people like you and there is a high chance that you will follow me back and therefore, spammers increase their in-links, which means on the topic that we discussed about link farming which means ==that== my node rank is increasing, which is, spammers' node rank is increasing which is what they want because of that their content will show ==up== on search, their content becomes more popular and therefore, they will probably benefit from it, ==I hope that== connects the dots.

So, I am going to use the same ==data==, to actually ==emphasize== on what is the behavior of these spammer? And what the link framers do? What do the spam followers do?

(ReferSlideTime:04:38)



Let us look at another graph. This graph is showing the probability of response which is in terms of just responding to a request with the indegree, indegree is number of links

thattheyhave.So,this is showingyou probabilityofresponsewas in-degreeforallusers targeted by spammers. If there are maybe, 10 to the power 7 users were targeted by the spammers, what is the probability that they are going to actually respond?

Users with low integrity do not reciprocate to links from spammers. If you look at the graph let us take less than 100, less than 1000 which is if I have followers which are less than 1000,, there is very less probability that I will actually follow the spammer back.Let us look at in the graph again, 10 to the power 3. Let us look at the value 10 to the power of 3, the probability of that users even then 10, 100 to 1000 if you see, the probability of somebody following back - the spammer is actually about 40 percent, 50 percent around 60 percent.

Ifyoulookatthelaterpartofthegraphonthexaxisandintheyaxis whichisabout0.7, responsivenessincreaseswithnumber offollowers,in-degreeisthenumberoffollowers. Asthefollowersincrease,thechancesofsomebodyfollowing you backwhenyou senda requestishigh.Ihopethatissinkingin,again letmereiteratethepointwhichis,onthex axis in this graph we are seeing the in-degree which is the number of followers, y axis is the probability of response when a request is sent for following or when somebody follows you, chances of you following me back. Users of less in-degree do not reciprocate to spammers; whereas, users with larger in-degree which is larger number of followers tend to follow back more.

So, what they did was after looking at these two things, which is the probability of somebody following spammer part is high, they actually looked at the top five link farmers and looked at the bios, what are the accounts and here is a sense of what these accounts are Larry Wentz, Internet affiliate, Marketing; Judy Rey Wasserman, artist, founder.

So, these are all the accounts which had the links to spammers; top five link farmers according to the links to spammers according to the Pagerank; the word Pagerank is nothing, but the links that you have to be out which is the in degree and the out degree that is what is PageRank is. Chris Latko, interested in Tech, will follow back and Paul Merriwether,helpingothers,letustalksoon;Aaronlee,socialmediamanager.So,itjust basically shows you what kind of users of the top five link farmers which is creating these links to others and then getting others to follow you back. Internet, social media manager will follow back; these are the kinds of accounts.

Ifyoulookatthepagerankwhichis,higherthelinksofhowyouarelinkedtoothersalso        this Barack Obama, Obama 2012 campaign staff; Britney Spears; NPR politics; UK prime minister; JetBlue Airways. So, this is also showing that it is not just the real spammers,butmaliciousintention,theyactuallydoinglinkfarming;evenlegitimate

accounts even more so popular accounts are actually part of the link farming ecosystem andtheyincreasetheirfollowers.Thatisakindofrevolutionthattheywantedtoactually       get whichis tolook atthe bio,the accounts oflinkfarmers andhave some understanding of what kind of users these are.

(ReferSlideTime:08:57)



Interestingly top link farmers are not the spammers. They looked at top 100,000 link farmers at the point of analysis of which of course, 18826 were suspended. Twitter did something, figured that all these are actually spammers, these are actually malicious accounts. Therefore, they suspended it. 4768 accounts were not found which is that they probably deactivated the account the account does not exist, whereas, 76 percent of the account of 100,000 link farmers which is, how do they get this hundred thousand link farmers; they do the graph of node rank, they do the graph of what, who were the most popular link farmers and they got this hundred thousand link farmers and of this 76 percent were still alive.

Interestingly of that, 235 were verified accounts, if you remember verified accounts are the accounts which has a blue tick next to the account and these are the accounts which arelegitimatethatistheyknow,theyshowthattheyaretherealpeoplewhichisAmitabh Bachchanhas the real account,verifiedaccount; Obama has a verifiedaccount whichare

the real peoplewho they saytheyare.

They manually checked 100 random users of 235, but volunteers of course, they got some user volunteers to verify whether to look at these 100 random users and said somethingabouttheusers. They foundthat86wererealaccounts,theyactuallygotmore than one people to look at it, therefore if more than one person says that it is a real account, there's a high probability that its a real account. They actually found that of 86 real users, people were like had the account as business, internet marketing, entrepreneurship, money and social media. These are the topics that the 86 real accounts are talking about. It just gives you the sense of, it also connects very well to the Twitter account bios that we saw in the slide which is top five link farmers.

So, this shows that the top link farmers are not really the ones who are standing in real world, but they could be actually, they are actually real accounts, they are actually verified accounts

(ReferSlideTime:11:27)



Let us delve little bit more into the node degree distribution which is how the in-degree followers are for top 100,000 link farmers, for spammers and for random sample. Well whatisthegoalhere?Thegoalhereistotryandcomparethehundredthousandlink

farmers they ==found== with spammers, which are real spammers; their accounts that are recorded and random sample of users, if they compare ==these three== types of users, the observation can be very helpful to understand what is the property of ==these== 100,000 link farmers..

If you look at this graph, this graph basically shows that top link farmers have very high ==in-degree== compared to spammers and random sample. So, let us go through again the graph in ==detail==; x axis is in-degree, which is again log scale ==which is 10 to the power of1== to 10 to the power of 7, that is the number of followers, cumulative distributionfrequency CDF is on the y axis. The way you look at this is that the more the value onthe red graph, red line is which is if you look at the 1000 users which has 1000followers which is in degree which is very high. If you look at that, the CDF is about 0.1.

So, top link farmers which is, if you arranged ==the in-degree== in particular order, the top ==link== farmers have very high in-degree compared to spammers which is the ==graph== for the spammers and the random samples is very different and also the CDF is actually verylow within the in degree, which is 10 to the power of 2, 10 to the power of 3. Top link farmers ==have== very high in degree compared to spammers and random sample.

(ReferSlideTime:13:22)

Interestingly, they also found that for the out-degree also, the graph looks very similar which is, the number of followings that I have and the top link farmers have very high out degree compared to the spammers and random sample and slightly very different from, slightly different from the in degree graph, but still again the 100,000 link farmers graph is much higher in terms of out degree compared to the spammers and the random sample.

(ReferSlideTime:13:57)



They also did an interesting analysis on finding the ratio between in degree and the out degree, this ratio is actually very useful because if you look at really large followers account like Amitabh Bachchan or Obama, the number of accounts that follows themwill be very high versus the number of followings that they have is actually very low. That is one parameter, one way to look at the legitimate accounts. If you look at again legitimate accounts like mine, probably the number of followers and the number of followings\ are actually very close to each other.

So, this is what they want to find out which is, if you find the ratio of in verses out, how is this;whatisthepatternoftheusers;mostofthetoplinkfarmershavearationearto1, if you look at the graph 10 to the 0 is 1 and interestingly, the value for the top link farmers,whyistheratiobetweenfollowersandthefollowingsisequalto1,whichisthe

pattern that I was telling you for real users like mine. So, therefore, link farmers also have this similar behavior. So, the other point to take away from it is given all this, it is going to be hard to find out who is the link farmer that is the kind of intuition that isbuild behind all this analysis.

Again, if you see this graph x axis is so, 10 to the power of minus 1, minus 2, minus 3 is where the in-degree verses out-degree ,when the out is much larger than the n would be the one that are less than 1. So, you can clearly see that the top link farmers which is the red color has the ratio of one, whereas, if you look at spammers and the random sample they are not really one, there is some difference with the examples that I took like Amitabh Bachchan, Obama and myself. So, I hope that makes sense in terms of what distribution is in-degree, out-degree distribution? What is the ratio of in degreeversus out degree? It give you a sense, go through the slides, go through the materials and if there is any confusion or any more clarifications needed, feel free to post it on forum, I will be actually happy to help in understanding these content also.

(ReferSlideTime:16:16)



Now,ifyou              lookatthebiooftop100,000linkfarmersandrandomsample,justtogetan comparison ofwhat are the people, what are the accounts which are actuallylink farmers andwhataretheaccountswhichare randomsampletalkingabout?Theleftoneis

actually the link farmers, the right one is actually random sample. You can clearly see here that the left one is talking more about market, online, internet, social, love which probably is in random sample also - life, music, live, love, right.

So, theone some conclusion thattheydo fromthisanalysis is thaton theleftyou seeLF, you can see that promoting their own business or content or trends in a domain, links to legitimate external sources. Of course, they are basically talking about some business, talking about some links that are outside twitter, outside the network that they are promoting this content. Thatis isright, donot tweetto external sources which is,thereis not a lot of links to other sources.

Again this will be a pattern that this research formed, but the pattern if you like to study this today, it may be very different also. I am kind of looking at some of these classical worlds which looks at some of the question that we should be asking in privacy and security in online social media topic.

(ReferSlideTime:17:43)



So, the final conclusion for this part of the work is, this part of the course / lecture is characteristics of link farmers we found the in degree verses out degree ratio is actually prettyhigh.Thein-degreeisveryhighcomparedtospammersandrandomsample.The

out-degree is also very high; the probability of spam followers is also very high in termsof requests sent to them or if Ifollow you, there is a high probabilitythat you will follow me if you are a spam follower. Surprisingly, legitimate popular and highly active users such as bloggers and experts, most likely engage in link farming, these are accounts like Britney Spears, Obama all of these accounts actually have link farming behavior.

So, the problem is that this increases the social capital and the influence because if the link farming engage, if the concept, I will go back to the first slide again on this topic. If there is high probability of, if there are chances that you let your social reputation, the links between the users are higher, their social recognition are I mean today social reputation, influence is all measured by number of followers you have and the kind propagation of of the content that you have and therefore, link farming can be pretty effective in terms of increasing your social capital and influence.

(ReferSlideTime:19:10)



WiththatIwillactuallystopthisparticularpartofthelecture,whichisweek7.1.

**Nudges**

Welcome back to the course Privacy and Security in Online Social Media, week number 7 - second section of the week number seven. So I hope you got a chance to look at the content that we made available for the week 7.1.

(ReferSideTime:00:27)



I do not know,how many of you actually read privacy policies. Researchers have shown that people do not read privacy policies. Just let us do some exercise now. So, some of you may have actually done online transactions in the last one week or almost all of you would have done some kind of online transactions, logging to your Facebook, logging to your Twitter, transactions of buying things online or between banks. How many of you haveeverreadanyprivacypolicyofthewebsitesthatyouhaveinteractedwithinthelast week or a month, must you really low. And it is not just what I am saying, people have actually studied this exhaustively.

So, here is one piece of work which says, if we were to make people read privacy policies,whatitwoulditcost.Thestudy was done in the US(ReferTime:01:19) anditis actually part of PhD thesis work, where the question was, what would happen ifeveryone read the privacy policy for each website they visited once each month. Time taken would be about 244 hours per year, which is basically, that questions more on the lines of economics of what would it cost. It is 244 hours,convert it into money,

, and in total number, if we just look at the US population, national opportunity cost for reading privacypolicywould be 781 billion US dollars, which is, if I get everycitizen of the US to read the privacy policy for least one month for the websites they have seen, it would cost in some 781 billion dollars. That is a lot of money, and that is a lot of opportunity that is being lost because they are spending time on reading this privacy policy. While they could make the pair, which they would (Refer Time: 02:14) use to make the decision.

(ReferSideTime:02:18)



So, keeping that context in mind, which is reading the privacy policies, researchers started asking some questions, and broadly also, there is this whole area whereresearches are working on, and technologies is being built to help users make informed decisions,whichishowcanIhelpuserstomakeinformeddecisionswiththe

information that is presented to them, with the information that they can actually use from these services.

Thespecificgoalistohelpindividuals avoidregrettableonlinedisclosures,whichis,can we actually build technology, can we actually build something for the users to use, so thattheycanactuallybehappyaboutthecontentthatthattheyarepostingoravoidbeing regrettable for the information that they are disclosing online.

(ReferSideTime:03:10)



Facemail from MIT

Sohere is onetechnologythat was builtat (ReferTime:03:14) MIT,meaning evennow when we are using this discussion forum for this class, when you type in the discussion forum, add the course name, you will actually find out in this plot, when if you use this technology,you do not really know who are the people who are getting this email, right, (Refer Time: 03:36), because the groups are set up differently, people have been signed up into the groups, you really don't get to know who you are interacting with (Refer Time: 03:44). For example, in a mailing list in your company and mailing list in your college or mailing list that you maintain for yourself with your friends, so all of this you do not really get to see who is getting these emails (Refer Time: 03:57).

So,whattheseresearchersatMITdidwas,theyareactuallysaid,okay,wheneverthe

email is going to go, because this would, this could, be a problem also, right, because you do not know who is getting the email. So, it could be a problem the information that you are sharing could actually go to people whom you do not want them to see this content.

So what they said was, okay, if there is a email address the state called (Refer Time: 04:19) psosmnptel2016 at abc dot com, they would actually show you the profilepictures of the people who are going to get this email, so that is a way by which to show that how many people are actually going to get the email, that is the information on the slide also, their profile pictures.And it would also show you who are the people who are gettingthe,becauseoftheprofileyou gettoknowmoreaboutthepeoplewhoaregetting to see this email.

So, this just helps users to make a decision on sending this email, because let us take if you were to send an email, and if you wrongly typed the mailing list address, it could actually end up going to a wrong mailing list. You wanted to send it only to 10 people, whereas the email list that you send is actually going to 100 people, so this can actually help avoid (Refer Time: 05:10).

In another version of it that they built also where they were actually showing you thebottom one, when they were actually showing you the profile picture of the person who you have been interacted more slightly bigger than others ones. This is also helping you to make a judgement on who is getting the email, who were the people who are getting these e mails and who you have interacted with more frequently than others. So this is justananotherexampleofhelpinguserstomakeinformeddecision,whetheryouwantto send the email or not.

(ReferSideTime:05:48)

So, keeping that in mind, which is to help users make this decision - informed decision - here is a piece of work that was done to look at the context of just Facebook. So in this experiments that they did, what they did was they built Chrome browser extension,which would actually show you, you go to facebook.com, when you are going to do a post it is actually going to nudge you with some information. They have different set of sub nudges, I will walk through what they did.

And we will also see how effective it was, which was effective, which was not effective. Because today, you could actually, meaning, I am sure some of you have experienced that you did some post and which you did not want somebody to, somebody in your friends networktosee,andtheygottoseethepost,andtherehavebeenincidences inthe past also. If you remembered even in the week 1 or week 2, I showed you someexamples about MI6 chiefs (Refer Time: 06:44) while posting some content that wentpublic, where it was not intended to go public.

So, these kind of incidents have actually made people to think about building technologies that will help them, help users make that decisions. So here is the first idea called picture nudge. What it does is if you wanted to do a post, when you are doing a post it is actually going to stop you and saythat these people, which is the profile picture connected to the face mail, it is the same thing as in the email mailing list. These people yourfriendsandfriendsofyourfriendscanseeyourpost,itisgoingtostopyouandtell

you this information. Again in the bottom screenshot which showing you these people and anyone on the internet can see your post. So, this actually helps you to make a decision on whether you want to actually do this post or not.

(ReferSideTime:07:38)



Here is a second experimental set up which is timer nudge. Here it is not showing youthe profile pictures. But it showing you that you have ten seconds to cancel your post. This information just lets you to say that essentially, you are doing this post, you really want to do this post, wait for ten seconds, if you want to change your mind, do it now, and then do the update. That is the timer update.

(ReferSideTime:08:10)

## Experimental setup

- Sentiment nudge

Here is next one which is sentiment nudge. Again there have been incidences all around the world where people have actually posted the things on Facebook and it is actually backfired on them. Backfired in terms of actually very negative effects also for the content that they have posted.

So, to avoid such things here is a setup where it shows you and in this case I am angry,so it shows you that it is negative; other people can perceive your post just negative. So, negativesentimentisattachedtothepost,sopleasebecareful,doyoureallywantpostit, and things like that, so that is the sentiment nudge. Let us look at all the three - picture nudge with pictures, timer nudge with actually ten seconds time, and sentiment nudge which gives you the sentiment of the post that you are making.

(ReferSideTime:09:04)

So methodology of the study that they did was Chrome browser as they said, they did exitsurveys, which iswhenpeoplecompleted thestudy,theyasked themsomequestions and also asked them questions in terms of both quantitative and also interviews with the participants. So, IRB is institutional review board, which is basically to say that if youareinteracting orcollecting datafromusers,itwillbehuman subjects,you reallywantto make sure that things do not go wrong when they are actually doing the study, and they should not feel offended and things like that. And that is why IRB approval is necessary when you are interacting with human subjects.

And of course, users were recruited in multiple ways. (Refer Time: 09:52) Of course,puttingflyersallaroundtheplacesendingoutemails,Craigslist,allof that.Theygot21 participants who completed the field study, because it is the Chrome browser plug-in, they could actually use it at home or wherever, and 13 people participated in the interviews, which is the exit survey.

(ReferSideTime:10:12)

## Analysis metrics

- Number of changes in inline privacy settings
- Number of cancelled or edited posts
- Posts frequency
- Topic sensitivity

So, in the metrics that they are actually used to analyze what is going on within this context of providing nudges of these things. Number of changes in inline privacy settings, number of canceled or edited posts, post frequency and topic sensitivity, essentially they were trying to understand by giving these nudges are people changingthe behavior. And if they change the behavior what are they changing, so that is the conceptthattheywereactuallytryingto study,that is themetrics that theyweretrying to study.

(ReferSideTime:10:43)



## Profile picture nudge

- One participant changed from "Friends" to "Friends except acquaintances" when she posted "Survived one of the craziest, most exhausting days ever!"
- Another participant ended up cancelling "a couple of posts" because of the profile picture nudge

Profile picture nudge. So the first one, I think I have one slide for a nudge to tell you what happened in this study. One participant changed from friends to which is probably all of them, friends except acquaintances, when she posted survived one of the craziest, most exhausting days ever. So essentially people are changing the group in which they are sharing the content depending on theinformation that thenudge is actuallyproviding them. Another participant ended up canceling a couple of posts, because of the profile picture nudge basically saying that oh it is going to actually lot more people than what I thought, what I think, so let me not do (Refer Time: 11:22) the post.

(ReferSideTime:11:24)



Timer nudge, one participant actually said at times annoying, and at times handy. Because I am pretty sure, right, because if you are doing, let's take, 10 posts a day or, like, 5 post a day also, it is going to stop you for every post ten seconds and then only togo. Waiting for a timer to expire or hit the 'post now'. So it is essentially a feature that could be provided. Make it more public, when it was venting type, right, make it more, so, it can be more suggestive. Another participant said, made be think about the posts, which is the same, which is the behavior that is actually important for these nudges to actually make within the users(Refer Time: 12:05). Changing the user behavior, it is abig need while using this technology.

(ReferSideTime:12:15)



Sentiment nudge, nudge was missing the context, of course, right, because I think the whole context of, in what context I am doing the post is actually very important to find out sentiment. And even human beings, it is actually hard to get the sentiments, so, the nudge was, the tool was actually making, browser plug-in was actually making errors while the calculating or finding out what these sentiments are.

Many participants canceled the posts, because, I think it is because if the posts are negative, and many people are going to be actually offended by the post within your network or in public, it is actually going to be bad for you. So, people are actually canceled the post. And the post frequency also reduces; they were actually doing 13 posts, it went down to 7.

(ReferSideTime:13:07)

**Conclusion**

- Interventions help users make better decision
- More work is needed in order to understand which type of nudge works in which context

So, essentially all of this is helpful in making the, helping the users to make better decisions. And particularly when it comes to posting something that others are going to be offended, posting something that people are going to, people whomyou do not intend to actually share the information, all of this is actually helpful, these technologies are helpful, for the users to make a better decision.

And of course, more work is needed to understand which type of nudges work in which context, because the contexts could be very different - I am just doing you a quick update, saying, I am, for example, I was actually doing a lot of updates, yesterday,about the convocation at IIIT-Delhi with the hash tag IIIT Convo 5. And if I wait, if it was (Refer Time: 13:55) going to stop me, and probably I did like thirty, forty posts on Twitter. But every time if it is going to stop me for ten seconds, that is not going to be good, so, and also, in terms of the, in terms of the sentiment that it is showing,everything has to be done slightly better.

So, I think more work is needed, I am sure there are people in the class who areinterested in taking some of this, it may be interesting projects to work on. With that I will stop the 7.2 part of the week. And I will continue on something which is in the context of phishing; and phishing in the context of social media in the next part of this week also.

**Semanticattacks:Spearphishing**

Welcomebacktoweek7andthisisthethirdpartoftheweek7.Inthisclass,inthis section what will see is, we will see about PhishingAttacks in online social networks.

(ReferSlideTime:00:18)



So, this is a slide from an MIT PhD thesis which actually looked at what a semanticattack is? Semantic attacks are attacks that happens where humans are targeted. So, for example, Bruce Schneierwho is supposed to be a expert in security classified the different types of attacks that could happen as physical, syntactic and semantic, but physical attacks are the were happening like 15-20 years before where the attackers would actually get physical access to the machine.

Whereas in syntactic attacks are the attacks that were happening around the programs, around the systems that are built, which is more like the denial of service attacks, buffer overflowattacksandattackslikethat,butthesemanticattacksareattackswhichtarget

the way we as humans assign meaning to the content which is that what do we because the specific attack that we will be talking about is phishing.

For example, ifyou get anemail frompk at iiitd dot ac dot in now,talking about NPTEL course and whichhas a link saying please give your user name and a passwordto seethe content heremost likelythat you'regonnaactuallyclickthe linkandgivethe information which may be a phishing link also. So, that is started the way you actually think you are seeing an e mail that is coming from legitimately pk at iiitd, and the system thinks that you are actually going to this free website forum psosm on NPTELdot come slash logindot html, but actually it is a phishing website. it is targeted phishing website.

So, system and mental model, what this is in this PhD thesis they actually nicelyput itthat semantic barrier, which is the difference between what system thinks we are doing and what you think that system is doing will actually be called semantic barrier in the larger the barrier is it is because actually difficult to not fall for such types of products. So, ifyou lookat the mentalmodelit says, who is the other partywhat is the meaning of the message. So, the example that I said also what is the meaning of the email we got what is the meaning of the who is sending the message and information is all mental model ofthe user, but a system modelwhois the remote machine where booked website I am going to access and information like that.

So, user model or the mental model which what users think that is happening, system modelwhich is what, the systemthinks that the user are doing, the barrier between them the difference between them is actually called semantic barrier. and the larger the barrier is, its actually hard to actually fix the problem. So, this is what we will use this is what we will actually talk about mostly in the section called phishing.

(ReferSlideTime:03:23)



Here is the broad category of semantic attacks: Security attacks - physical, semantic and syntactic which is what Bruce Schneier did and if you look at Semantic attacks, you can actually go through multipe categories Phishing, Mules, Nigerian, 4 1 scams and attacks like that, and in phishing also there are multiple categories – update your information, banks and in your ICICI banks sending you a message saying that, please update your information within next 24 hours or your account would be closed. Verification, saying that, we want to verify whether it is really you, please click and verify. Security alert, Microsoft is updating the latest version of MacOS, thereis an update.

Here is the link, please go and update. Mortgage information, meaning your mortgage, the due is coming closer, please click this link and do something. All of these kinds of categories ofattacksare called phishing attacks and almost allcompanies todayprobably are undergoing, are part of, or being victims of this attack of phishing. Even academicinstitutes probably are victims of phishing attacks.

Here is a simple example, which is an email that the 3 parts of the email whichisactually makes the legitimate email and the difference between the legitimate email and the phishing email, which is the subject line, subject line and urgency in the message on there are there is the line. These are the three things that happen that is a part of the phishing email which at least one wants to keep attention on. Subject eBay urgent notification from billing department. We regret to inform you that your eBay account could be suspended if you do not update your account information and then there is alinkthereandthenwhenweclickonthelinkittakes youtoawebsitecalledkusidotorg.

(ReferSlideTime:05:15)



Which supposedly should be takingyou toeBay sign in page, so thatis the sharing that is a very classical phishing attack when there are multiple ways ofa changing these kind of phishing attacks.

(ReferSlideTime:05:34)

Here is some cost again, economics about phishing, some costs that is relevant to the topic ofphishing. Costsofhundred thousand employees organization, which is the, ifthe phishing attack happens what would be the cost to contain the malware, the cost to contain a malware, the cost of malware not contained. So, if you look at the cost its actually pretty high in terms of actually even the phishing attack.Total extrapolated cost is 3 million 76; 3 million plus dollars, all right. So, it needs a lot of money that is spent every year, FTC and many other organizations in US actually try to course aboutphishing and there is an organization called anti-phishing working group which actually specifically works on the problem of phishing and how to actually reduce it.

(ReferSlideTime:06:28)



## Types of Phishing Attacks

- Phishing
- Context-aware phishing / spear phishing
- Whaling
- Vishing
- Smsishing
- Social Phishing?

So, here are some kinds of Phishing Attacks I think we probably briefly mention this in thepast also.So,I'llgoover quickly,phishingwhichisaclassicalonethat Ishowedyou, Context-aware phishing the email that I talked about sending in to the students takingthis course, Whaling is an attack which is sent to the chief executive officers of the company, Vishing is over the phone, Smsishing is over the SMS. So, what is social phishing?That is what the topic that we have been discussing for the rest ofthis week.

SocialPhishing;doesanybodyknowwhatsocialphishingis?

(ReferSlideTime:07:10)



Social phishing is nothing but looking at the information from the social context then using that to actually phish, it is not about finding whether you are taking a course that could be many other information that I defined on from a Facebook page, from the facebook account, things thatyou've done and things like that. So, the topics thatwehave seen until now are using older data we saw Latanya Sweeney's work using medical health data, again Latanya Sweeney's work, using pictures from FB voter data.

We also saw the work that was down in collecting pictures from the university campus collecting this information and making some judgments about the user. Finding the people who they see in the campus whether they will get the right profile from the Facebook. So, those are the topics that we saw, but we never saw about what social phishing is.

So, here is a goalthe goal is to see how phishing attacks can be performed by collecting personal information from social networks right. So, it is not about sending into the CEO's, it is not about sending to the students of this class. It is about actually can I collect some informationabout you fromyour social network behavior and use it against you, how easily or effectively can phisher use this information. Again there is a very classicalworkthat was done some years back. So,it'llbe nice to actuallyknow howthey did it some years back. I am pretty sure these studies can be done again to seehow itgoes and this study was done in the US.

(ReferSlideTime:08:41)



Here aresome examples. So,Ilove you virus, someofyou mayknow this. Kindlycheck the attached love letter, see attached files, this is an email that comes, coming from me right. It was one of the first virus that was actually spread.

(ReferSlideTime:09:02)

So, what they did was they collected publicly available personal information usingsimple tools you could actually collect now information from Facebook. So, we will actually have some tutorialsalso about NLTK, how to use it howto analyze thistext that are coming fromthese post, all right. So, you could collect this information and find out what is the date of birth mentioned there. This was done in Indiana university. I wasreferring to Indiana university. Coerrelated this data with Indiana University's addressbook.

Which is theycollected allthe posts done bystudentsofIndiana universityand then they launched the study in April 2005, they launched for the age group between 18 and 24 which is the student population in campus most of the times.

(ReferSlideTime:09:53)



So, here is the slide whichactuallyI think thenext slide we have also, go throughthis in a text form, but here is the slide that actually walks you through in terms of what the procedure that they follow right. First they actually look at public data, blogging, social networking sites, they collected the data which is stored into the social database, social network database, they use this data to create an email which is - FromAlice at indiana dot edu,To subject Bobat indiana dot edu.This is cool, heycheck thisout withthe URLthere, right.

So, from there the information is sent as an email, bob to friends at Indian University,and when the user clicks on the link. It goes to the Indiana University's website and it actually checks authentication web and the authentication logs, it tracks and takes into a user name and a password page, then when they give user name and password it checks the user name and the password whether that is appropriate which is checked.

Controls authenticator and comes back and server overloaded try again message is sent, authentication failed. So, those are two outcomes of the whole process which is success, server overloaded try again and an authorization failed right. That is a process that theyfollowed in terms of finding out information from social networks, sending out this emails and getting some uses to get to that page. Many people that have done this study after that, even if you go look out my own work in 2007, 8 and 9, I have done similar kind of sending out phishing emails and seeing how people behave.

(ReferSlideTime:11:35)



So, two things that theydid, ofcourse it was experimentalreseatch. So, theywere trying to compare how the email from Indian university email id, but from an unknown person that is a control goal, if I get an email from Indiana email id which in my case I get an email from somebody who is from IIIT, Delhi with the iiit email.

But I do not knowthe person because that's something I can actually get from the social context, experimental group is from a friend in Indiana university itself, which is if Ihave already connected to the friends in some way you saw the Facebook posts and following this person in twitter, I mention them in the post on twitter, so all that it actually helps to find out that is the experimental setup.

(ReferSlideTime:12:24)



So, the same methodology, the chart that was there here is the verbose of it, blogging social network and other public data is harvested, data is correlated and stored in a relational database, heuristics are used to craft the spoofed email messages, message is sent to Bob, Bob follows the link contained within the email and is sent an unshared redirect, bob is sent to an attacker whuffo dot com.

Bob has prompted for his university credentials, bobs credentials are verified with the university's authenticator and bob is successfully phished, bob is not phished in this session, he could try again all right. That is the verbose of the architecture that wasshown or the experimental methodology that was shown in the slide before.

Continuing on the analysis of data that researchers collected in terms of social phishing here are the results. So, in this table what we are seeing is control condition and social condition experimental condition, as in the rows columns being successful targeted percentage and confidence intervals. Successful meaning how many people got those emails who actually fellfor the it, targetedthe number ofpeoples who were actuallysent this email to percentage of course, is the ratio of targeted versus successful. So, let us look at some results.

(ReferSlideTime:13:49)



So clearly control group is, high which is that 16 percent of the participants falling for these kind ofemails itsvery, veryhigh ingeneralit isnot abad levelwhen inIthink it is hard to believe that16 percent of the participants actually fell for this kind of theseemail, but the advantage here or the context to keep in mind for the data is that sender email was fromIndiana university itself, I think that is the reason why this percentage is very high.

For example, if you get an email frompk at iiitd dot ac dot in versus if you get an email from pk at lets takes some abc dot com, the higher chance of you clicking and going through what a emailsaying, asking you to do would be pk at iiitd dot comor pk at iiitd dot acdot in, is veryhighand ofcourse, 72percentofparticipants inthe socialcondition clicking the email and doing whatever is asked in the email is actually pretty consistent with other studies that have been studies where they have shown that this percentage is even higher than 72 percent.

Moreresults which is to seethe success rateofhow people, authenticatortothis website. So, here are some interesting results again. Seventy percent of authentications. So, what doesthis graph show, this graph has xaxis being the time, date, the dates which is 6 pm, 6am,12noon,6pmasinthexaxis,yaxistobethepercentageofpeoplewhoactually

clicked for authentication, so green is showing you cumulative authentications, red line showing you authentications per hour and then blue line is visits per hour. So, you essentially what does it mean, blue line is showing you that number ofpeople who went to this website, red line is showing the number of people who actually authenticated which means it'll always be below.

This clearlyshows that 70 percent ofauthentications inthe first 12 hours. So, ifyou look at the first part of graph, 70 percent of authentications which is the red line which is actually in the first 12 hours itself. The problem is that people fall for these kind of attacks immediately when they get these emails, right because there is a sense ofurgency, there is a sense of completion, completing it immediately and that level of urgency is put in this email. If you remember the example that I showed you from eBay website, an email that they sent has a subject line also has urgent notification, urgent verification, but this actually puts a challenge on solving the problem.

Phishing, which is takedown has to be successful, which is if the websites are takendown as early as possible as soon as possible, then there is a high chance of this several users who were going to this website can be actually stopped. If the websites are not taken down unfortunately these users just actually end up actually going to the fake website and giving away other personal information right. Again success rate, how people react to these emails what level of authentication, whatis the percentage ofpeople who are actually giving their account details, is what theyshow in this graph.

(ReferSlideTime:17:17)



So, here is another interesting analysis that this research actually shows which is that subject, subjects actually tried, participants tried multiple times to actually authenticate which is the blue line is actually showing you repeated authentication and the red line is showing the refreshes of authenticated users, which is they trying to refresh and see whether they are able to log in to system. If you remember the architecture that final output is two. So, let me go back, finaloutcome of the study is at two levels, where this authentication failed, server overloaded, so they try again. So, when user sees this they feel something is wrong with the system. Let me just refresh it and let me just try itagain, that is what is happening in here.

So,triedagain becauseoverload messagewasshownthat lotofpeoplewho actuallytried because the overload message was shown. So, this is basically showing you the lower bound of users to fall and continued to be deceived and if you we look the blue line which is people are actually authenticating to this website repeatedly, even it is actually showing you that authentication error, some people actually seem to tried to 80 times.

So, the x axis is here showing you the count which is the log scale, y axis showing you the number of subjects you can clearly see that about 80 percent of the 80, some people eventrieditfor80timestogetintothewebsite,andthisisnottheonlystudy,whichis

showing this, this is probably the classical study, one of the first studies which showed this, but later there have been many studies who showed that such kind of repeated authentications happen with the users.

(ReferSlideTime:19:02)

## Gender

|  | To Male | To Female | To Any |
|---|---|---|---|
| From Male | 53% | 78% | 68% |
| From Female | 68% | 76% | 73% |
| From Any | 65% | 77% | 72% |

- 18,294 Ms and 19,527 Fs
- Overall F more victims
- More successful if it came from opposite gender

Here is the ratio, here is the analysis of the gender, because they had the Indianauniversity's, university student details they could actually find out male versus female, the gender details of the participants. So, this table actually shows you on the rows, it shows you frommale, fromfemale, fromanyone, to male, to femaleandto anyone.This basically says that if the email is coming from, again if you remember the study was set up, they collected the data, they have crafted the email, and while they are crafted the email they were actually doing all these experiments to see, if I send it from male what happens to when the email goes to a male versus male to female, all right.

So, this shows that overall female were more victims which is you can see on the third column, which isto female being much higher thanto male, it doesnot matterwhere the email is coming from, female seems to be more vulnerable to these kind of attacks. 18,294 males and 19,527 females were actually being part of the study, more successful if it came fromthe opposite gender. You can clearly see that from male to female which istherow2andthenthecolumn3,78percentandfromfemaletomalewhichis68

percent. So, this number which is 68 is the highest in the column of the to male, 78 which is the highest in the column of to female, which basically shows that if a male gets an email from female and the female gets the email from a male, it is actually high.

The percentage of chance of actually authenticating and giving away the information is much higher, if the email came from the opposite gender, that is a very interesting conclusion that to show that phishing attacks, or vulnerable phishing attacks are successful, but is also more successful if the emails come from the opposite gender, and sure these results can be repeated even in non email context.

(ReferSlideTime:21:18)



• Younger targets more vulnerable

Again given that they had a lot of demographics data here are some things about age group, things about the departments that they were part of. You can clearly see here that the youngertargetsaremorevulnerable, whichisthe youngerandtheparticipantsarethe more vulnerable that they are to, for authenticating to the study. You can see that freshman, the difference here is that the orange and the blue, the orange is showing you the social phishing which is the experimental setup; the blue is showing you the control condition,while youcanseethedifference betweenthe freshman,difference betweenthe social and the control is the highest in terms of freshman.

And as you go up it keeps reducing, so fromor junior and senior in sophomore, ==it lookslike little high,== but if you put the sophomore and freshman together it basically says that the younger the people the more vulnerable they are to these kind of attacks.

(ReferSlideTime:22:18)



Similarly,theyalso ==had==whichdepartment==the==participantscame from.So,==theywereableto== actually create a graph which ==looks== like this, again the same color of color scheme, which is orange being the social condition and blue being the control condition, which basically shows that all majors significant difference between control and experimental, which is any department of the campus it does not matter, the difference between social and controls ==is very high== which is social people fall more compared to the control condition.

It also showed that the science department had the maximum difference, if we look at science 80 percent is for the social and 0 percent is for the control condition. So, which ==shows== that the science department had the maximum difference between the social and the control.And it was also ==evident== that the technologyhas the smallest, which is people ==who== study technology that have probably are less vulnerable to these kind of attacks which is about 36 percent here, difference between the social and the control condition right. So, this is way by which research is actually found, which kind of department and

the studentsgoing to whichkind ofdepartmentsare actuallyvulnerable tothese phishing attacks.

(ReferSlideTime:23:38)



In general, this studygot a lot of negative reactions fromthe participants which is like it was unethical, inappropriate, illegal and it was also fraudulent, researchers fired researchers were fired, psychological participants claimed that there were psychological cost, because I think they were under pressure, they did not know about the study happening and things along that, and interestingly there were people who wrote blogs, people who wrote reactions about the study and they said that they were not part of the study and they did not fall for these attacks with somebody elsefell for, which also shows that admitting that I am vulnerable is actually is also hard; I think that is a misunderstanding over spoofing emails, underestimation of publicly available information.

So, participants did not, meaning generally also you and I will not perceive how bad the publicly available information about you can be used against you, since this was one of the first studies these reactions were actually interesting, but there are people who have done the studies after this which were again you studied how people fall for phishing emails.

Essentially, what does the results show is that extensive education campaigns is necessary, browser solutions of course,take downhas to be much faster, digitallysigned emails have to become more prevalent and of course, online social media provides lot more information for making these attacks more successful, all right. So, people should stop sharing a lot more personal information on social networks, digitally signed emails should become more prevalent, browser solutions should be built.

To say that this website is actually, this email is actually phishing and this website is actually malicious and this website is a fraudulent website, and of course all of this can ccome to education campaigns.

(ReferSlideTime:25:39)





Some reference is for this research that I discussed. With that I stop actuallythe week 7; we will actually look at some more exciting topics in week 8.

**ProfileLinkingonOnlineSocial Media**

(ReferSlideTime:00:19)



Welcome to the course of Privacy and Security in Online Social Media. This is week 8, the first part of week 8. So, just look at the profiles on the screen, it has first one Facebook handle called ponnurangam dot kumaraguru (ponnurangam.kumaraguru), the second one which is Twitter profile called ponguru, and the third one which is LinkedIn, which is ponguru again.

So, the question is, can you actually match all these 3 URLs or all these 3 profiles and saythat it is the same profile. That is the question that, we are going to try and answer in this, this part of, this week of the course. Which is I have handles of ponnurangam.kumaraguru from Facebook, ponguru from Twitter, and ponguru from LinkedIn. Can I actuallyuse this? What do I need to do to make sure that these 3 profiles are same or to understand that whether these 3 profiles are same. There are multiple actually test cases, scenarios for it; I will actually discuss a little bit later in the lecture.

(ReferSlideTime:01:36)



Here is, the top one is my Facebook profile, the one at the bottom is my Twitter profile the one.

(ReferSlideTime:01:44)



Thenextslide,thisismyLinkedInprofile,publicLinkedInprofile.So,ifyoulookat thesethreeimages,youcanactuallyseeoryoucanactuallythinkaboutsomefeatures

that you can use for deciding whether these 3 profiles are mine. For example, you can look at my profile picture in both, they seem to be the same thing, you can look at probablysome friends that I have on Facebook and people who are following me or then the accounts that I am following on Twitter, you can look at some of these features to make the decision. Unfortunately, in the public profile that I have on LinkedIn, there is no profile picture, but there are details like associate professor at IIIT Delhi, Data Security Council of India, Carnegie Mellon University and connections like that.

For example, my personal website, the personal website from here may be actuallylinked to my website at IIIT Delhi. So, you can actually make all these connections to find out whether this is actually the same details is both in Facebook and the Twitter. I am sure many of you are listening to this lecture also have multiple accounts. So, the question that you can ask yourself is, how do I put, how do you put your own accounts together to find out whether they are same or not. So, that is the problem that we look at.

(ReferSlideTime:03:21)

## This lecture

• Tracking social footprint / identities across different social network

So, tracking social footprint identities across different social networks, which is finding out whether they are the same.

(ReferSlideTime:03:33)



Other Times, Other Values:
Leveraging Attribute History to Link User Profiles across Online Social Networks

Paridhi Jain
Indraprastha Institute of
Information Technology (IIIT-D),
India
paridhij@iiitd.ac.in

Ponnurangam Kumaraguru
Indraprastha Institute of
Information Technology (IIIT-D),
India
pk@iiitd.ac.in

Anupam Joshi
University of Maryland, Baltimore
County (UMBC),
USA
joshi@cs.umbc.edu

**ABSTRACT**

Profile linking is the ability to connect profiles of a user on different social networks. Linked profiles can help companies like Disney to build psychographics of potential customers and segment them for targeted marketing in a cost-effective way. Existing methods link profiles by observing high similarity between most recent (current) values of the attributes like name and username. However, for a section of users observed to evolve their attributes over time and choose dissimilar values across their profiles, these current values have low similarity. Existing methods then falsely conclude mation, lists her friends and later creates content to share with her friends. The quality, quantity and veracity of the information created and shared by her vary with the OSN, thereby resulting in dissimilar profiles of the same user, scattered on the world wide web, with no explicit links directing to one another. These disparate profiles liberate her from any privacy concerns that could emerge if the profiles were implicitly collated. However, linking these disparate unlinked profiles can benefit various stakeholders.

Companies like Disney and PepsiCo carry out psychographic segmentation based upon customers' activities, interests, opinions

Jain, P., Kumaraguru, P., and Joshi, A. Other Times, Other Values:
Leveraging Attribute History to Link User Profiles across Online Social Networks

And as always in the past also, in the lectures I have said many of these topics that I am discussing in the class it is all connected to some research done. So, here is a paper that I am going to be talking in detail today, which is 'Other times, Other values: Leveraging Attribute History to Link User Profiles across Online Social Networks'.

(ReferSlideTime:04:02)



Knowing this can be useful!

So,bigadvantageofactually<mark>knowing</mark>theseconnections,whethertheyaresameis actually very, very useful.

(ReferSlideTime:04:10)



Let us look at this slide, please look at this slide, which says about duplicating audience. Where if I were, so, in this case there are 437,000 likes on a Facebook page and about 153,000 followers that the account has and about 800,000 followers that the handle has onLinkedIn.So,the questionis,ifIwere tosendanadvertisement,ifIwhere toactually send some information to these users, will it be same, <mark>will it</mark> be a sum of all of them or will it be something smaller, because that could be some of these 470,000 profiles, the sameusersareactually153,000inTwitterandthesameuserswereactuallyonLinkedIn.

For example, I am sure some of you in listening to this lecture will have accounts on Facebook,TwitterandLinkedIn.Ifyouhaveaccountonallthe3andifPKwantstosend you about information on PSOSM on NPTEL, it is actually useless to send the information to the same handle, which is, let's take, Sonu Gupta in Facebook, Sonu GuptainTwitterandsonudotgupta24inLinkedIn,becauseitisthesameperson.We're actually wasting our resources in sending this information about (Refer Time: 05:47) PSOSM on NPTEL to the same person in all 3 accounts.

So, that is the problem to actually look at. So, the question is, people have multiple accounts on social media and sending information to all of them, you want to send information to the people only once. So, that is the goal then, but there are many test cases for this problem. At the end of this lecture I actually talk about some other test cases in law enforcement (Refer Time: 06:17) agencies and in other situations.

(ReferSlideTime:06:22)



A technical challenge for actually putting them together is also harder because if you look at some networks, you get actually details, which are something more personal, in some networks, they are not actually that personal.

Forexample,in thein thetop part,Iamshowingyou here aboutYouTube,beinga video sharing service, you can get actually opinions, you can get what they like, what kind of videos they actually saw, in Tinder, which is the dating side, little bit of personal information is available, connecting in to LinkedIn, which is professional and Facebook, which is also personal details, right. So, the question is, what information can you actually collect from these different social networks, which have different types of information, how do you put them together and create answer the question that westarted off with, finding out whether multiple handles are same or different, right. I hope that is clear.

(ReferSlideTime:07:30)



So, the question about profile linking, what are the approaches that we can take? The approaches that we can take is list out common attributes, which is Facebook has my gender, my age, my university that I work at, places that I got my degrees from. Twitter has my followers, my profile again, the website that I am connected to, the place that I work, all that information. We can actually list on all common attributes, compare the attributes, which I think in the example that I showed you, I showed you profile picture being same, profile picture being same on Twitter and on Facebook, we can actually compare that.

Compare attribute values using syntactic, semantic or graph based, which is what I am typing in, on, the social networks, what content are am I posting and what will, what is the details in my profiles and the graph is basically my networks - my friends in Facebook, my followers, followings in Twitter..

Andthen highsimilarity,if there is, inmycase in the example that I showed you it is the exact the same picture profile picture on both the places. If things are like that, it mostly likely the same person. And then the question is also, you can, so, one thing that I will talk about few slides later is not just that you want to look at these details only that is now,

butyou canactuallylookatdetailsthatarepastalso,whichis,youdonothavetolookat onlythe post that I did now or the profile picture that I had now or the handle that I have now.You can actually go back in time and look at the post that I have done and you can derive some information even from that. For example, one thing I will show you also is people actually change their user handles sometimes. So, can you actually use that information to actually derive whether it is the same profile?

(ReferSlideTime:09:47)



Now, if you look at this graph, this graph is actually showing you the changes that has happened in terms of just username, the point which I just now said, which is, some details of the profiles can actually change over time. It is not that you have to look at the details that are now,but you can actually look at the past - that is the problem, that is the question that weare trying to answer there. So,here 376 million users where tracked and the graph is showing you x axis to be the time and y axis to be the percentage of users with multiple screen names, which is names that they have changed.

For example, in my case, currently I have ponguru, whereas in the past, let's take if Iwouldhavehadponnuguru123orprofessor@iiitdelhi,allthosethingsareactually

getting captured here; which means some 7 percent of the people had different, change their usernames sometimes, in the data that we collected in around January 2011. That is the way to infer.And then there is another peak around January 2000 or February 2010, the two peaks in this graph are basically showing you, that, what the first peak isshowing you, about between 5 and 6 percent of users that we were tracking, the handles where changed, and about 7 percent of the account user handles where changed around January 2011 or December 2010.

So, this basically shows you that people change accounts, people change their handles. (Refer Time: 11:39) Again for people listening to this lecture, think about yourself,how many people have actually changed the handles that you use. In Facebook I think youcan change it with only once, but in Twitter you can change it as many number of times you want, which means it is actually possible to keep changing your account every now and then.

(ReferSlideTime:12:04)



So, continuing on the same thread, which is about the details changing for the users. So, in this case if you see, the x axis is the detail of the user, which is username, name, description, location, language, zone and profile picture, just basically showing you, and theyaxisisshowingyouthepercentageofusers,whichisinthiscase8millionusers,

were actually seen for period of 2 months. What does it mean to for the username, when about6to7percentoftheuserschangetheirusernameatleastonce,whichis,therewere          two values for these users, that is the way to read the graph.

Let us go to the good one, or the or the one that is higher, in terms of profile picture, if you see, of the 8 million people that were tracked about 40 percent, 35 percent, of the people change the profile picture 3 times at least, right. So, that is the one that is there in the blue. Just on top of it, which is yellow,which is about 40 to 20 percent of the people changed profile pictures for 4 times, and about 10 percent of the people changed it 5 times.

Which means in the period of two months, 40 percent, 10 percent of the people changed their profile pictures at least 5 times. I am sure you can relate it to the behavior that you have, which is just how many times that you change, in my case, probably I change my profilepictureonceayearoronceinayearandhalforso.Butprofilepicturechanging,I have seen many people change their profile picture pretty often. So, that is what is reflecting on this, the right most thing. And left most thing, where username, similarlyfor name people have changed the names. And if you look at, about 35 percent of the people are changing their description, which is say, professor at IIIT Delhi, at least two times in the data that was collected.

Nobody changes language, nobody is changing, very few people are changing the zone, time zone, that they are in, right that. So, basically this graph and this graph, the slide 11 and slide 12, is basically showing you the change in information in the account.

(ReferSlideTime:14:54)



So his is just an example to show you how people change their handles. For example, in my case, I have Twitter now, that is how registered, which is t1, whereas later I could changemyuserhandlesaspongurutobecomeexplorerunderscorepk(explorer_pk)and at time three, I could change my account as logical Tamilian, for that matter.In that case first one, it was actually identifiable, ponguru, we can probably derive it from my name, second one when I had explorer PK explorer, professor, something like that, slightly getting anonymized (Refer Time: 15:43) and same thing as logical Tamilian also, it isgetting anonymized. (Refer Time: 15:44). And it is also unmatching, the point that is expressed in the slide is also to show that the handles ponguru and ponguru at t1 versus ponguru and logical Tamilian at t3 is actually not possible to put together and find the answer. So, there is difficulty in putting this handles together.

(ReferSlideTime:16:17)



## Problem Statement

*Given two user profiles and the respective **username** sets,*
*each composed of past and current **usernames,***
*find if profiles refer to a single individual?*

Given two users - this is more scientific way of asking the question - given two user profiles and the respective usernames sets, each composed of past and currentusernames, find if profiles refer to the same individual. That is the question that we are trying to ask, which is, I give you ponguru, and I give you ponguru's current user handle and the past user handles, can you put them together and say that whether it is the same ponguru, which is a (Refer Time: 16:48) professor at IIIT Delhi and ponurangam dot kumaraguru (ponurangam.kumaraguru) in Facebook, ponguru in Twitter,and ponguru in LinkedIn.

So, in this slide the point that is described is that why only usernames, why should we look at only the usernames as the change, as the history, use the information from the history to actually study this profile linking. Because it is unique attribute for a user universally and publicly available attribute, because it is not, you cannot make your user handle private. And sometimes the lines of the handles are also restricted. So, it is not infinite space that I have to actually look for. And of course, in terms data collection, in terms of details that we can actually collect from social media, it is easy for collecting user handles.

So, that is the reason why studying usernames is the way that we looked at.

(ReferSlideTime:18:00)



## Methodology

So, this is slightly a dense slide. Let us see how we can actually get this slide across. So, what is given?

(ReferSlideTime:18:15)



## Problem Statement

*Given two user profiles and the respective **username** sets, each composed of past and current **usernames**, find if profiles refer to a single individual?*

So, we would let us go back to the problem statement, the problem statement is, given two user profiles and the respective username sets, which is, I am giving you the handle of ponguru and can l actually find out.

(ReferSlideTime:18:25)



So, that is what is actually explained in this slide which is SN A. SN A, which is in our case let us take it as Twitter, SN B which is something like, let's take we keep it as Facebook. We are going to look at handles in these two networks and find out whether the handles that we are looking at are same, right. So, we look at features, for example, thatwesaidearlier,profilepicture,locationoftheaccount,usethisdetailsasthefeatures and find out whether they are actually the same user.

So, here is a one very, very interesting way and very easy way actually to find out whether it is the same image and there is no probability that this feature may not be useful. But majority of the time this feature is actually very very helpful, what is this? This is Twitter handle, which is, which says, in this case, l u z y, and the user is actually connecting her own Tumblr account in this page, right. So, this basically allows you to say that if I were to find out this handle's, luzy's, Tumblr account, I should just look at the profile. Same way, in my case, if you go, I think, my LinkedIn or my Facebook, has my Twitter account also there. Which means I have explicitly specified, that self identification,whichis Iamidentifyingmyself,thatIamthisinTwitter,Iamalsothis in Tumblr, which I'm sure some of you may have done, in my case, I definitely have my precog.iiitd.edu.in URL in my Twitter account.

There is also another way this self identification happens, which again I do it very often,I post pictures on Facebook, I take the link of the album, and then I go post the link tothe album in my Twitter account. Which now, if you see, you can actually connect that ponguru,@ponguruinTwitter account, isthesameaccount whichisactuallypostingthe pictures on Facebook, which is this album. And therefore, they should be actually the same people. Evenwithout theprofilepicture, evenifmyprofile picture is different, you canusethistomakethedecisionthatitisactually thesameuser. Ihopethatismaking

sense.As I said, there is a small probability,that this may notbe true, but majorityof the times this is actually predictable.

Student: Sir (Refer Time: 21:23) say if a person is not very much active on Twitter, but active on Facebook, then how can we link?

Ponnurangam Kumaraguru: Not active, I do not think so, activity frequency actually matters here, right, because let us take -

Student: No, activity, I mean to say if a person is updating something on Twitter (ReferTime: 21:46) he need not update the same thing -

Ponnurangam Kumaraguru: Oh, sure, sure, if the person does not update the same thing on itthen itis okay,thereisthereisgoing to bealways aproblem. Butif thepersondoesit, it does not have to be the same like what I am saying, linking of the pictures, it could be the same post at the same time and for both pages, both accounts. Even that is useful right? If the person is not posting then I cannot help you, but if the person has the same profile picture, same description and things like that, I can put them together.

But if the person doing the same, like for example, if you see my post, right, I'll do 10-30onFacebook,10-30onTwitterand10-30onLinkedIn,allatthesametimeanditwill allbethesame content. So,now,itiseasyto find out, right, even though my,let'stakeif I change my Twitter account to instead of ponguru I change it as professor at IIIT, still you can actually make it because it's the same content at the same time.

Student: (Refer Time: 22:55) If I have a Facebook account and I am using it very often, but I never update anything on Twitter, then there is any possibility that we can link the accounts? If there is nothing on the account, the Twitter account, to link it?

Ponnurangam Kumaraguru: Sure. So, if there is no content on Twitter then I think it is slightly hard, but if the account has some details about the account, let us take description,likeuserprofilepicture,thenit'seasy.Butintheletmestretchthe(Refer

Time: 23:27) even slightly further away, which is that if the person does not even have a Twitter account, ==what do you do?== It's also a problem, right. So, you just went to the extentofsayingofthepersondoesnothaveactivity,butifthepersondoesnotevenhave a handle, then it's even harder, right.

Student: It may happen that have my actual Twitter account or some else (Refer Time: 23:54).

Ponnurangam Kumaraguru: Yeah yeah. So, here right Ponguru's my account as in my in in Twitter and I have my Facebook account of Sonu Gupta, what do you do? You just cannot put them together, right. At least from the handles and this, you cannot put them together, but that is why we need to use all these features to put them together. In terms of profilepicture, posts that you do, butif you areconscious enough to keep this account two independent, then I think it is impossible to do it, and, but companies like Facebook can do which is beyond what we are taking about in this lectures, because they can actuallylookat it fromthe IPaddress,theycan actuallylook atitfromthe time ofaccess and still ==make it== (Refer Time: 24:39).

(ReferSlideTime:24:42)

So, if you look at, here, this is showing you ways by which you can actually collect this data,soherethe oneonthe topis showing you that, keeptrack ofahandle.Forexample, you keep track of ponguru at now, what my handle is, what post am I doing, every day youcomeandlookatthis handle,andthenyousay,oh,suddenlyhechangestoprofessor IIIT, professor at IIIT, you can say that the change has happened, right. Because what have, how is the data that is getting stored in Twitter? The basic idea on Twitter is that they give you a unique id - that does not change, right - that id is associated with the handle, you can change that handle.

So, now,you keep track of this id and you know that this id 24 - just making it simple – 24, is actually ponguru. Now you can actually keep track of this 24 always and then ponguru changes to professor at IIIT.Then you add, update, it to your into your database saying that, oh, this handle actually changed. And another way of looking at this is the URL change, which is, the person actually changes the URL in terms of connecting tothe other accounts. Like the Tumblr one that I said, somebody is actually going it is in, I am keeping track of Tumblr accounts also and I am keeping track of other accounts. There, the profile is actually changing their description to say that my Tumblr account changed, or in Tumblr they are saying my Twitter account changed. So, this track again <mark>you can actually set</mark> (Refer Time: 26:38).

(ReferSlideTime:26:39)

## Sample

- **User ID**: 595929421
- **Past usernames on Twitter**:
  - ["bigeasye_", "reezy11_", "epiceric_", "soulanola", "swampson_", "hebetheeeric", "swampkidd_"]
- **Past Usernames on Instagram**:
  - ["bigeasye_", "epiceric17", "swampson", "hebetheeeric"]}

So, here is an example of users whose accounts have changed, whose handles have changed. User ids, as I said before, 24, in this case is 595929421, that is the handle. That isthe user id. Wekeep like a track of that, and if you look at the names that this handle has changed it, when it from bigeasye underscore, to reezy11, to epiceric underscore, to something else, to swampson x y, swamkidd, right. It changed 1 2 3 4 5 6 7 times, it has actually changed the handles.

And if you look at the same user account in Instagram, this seems to also changed four times, but there is a connection between the users handles that the person had in Twitter and then Instagram also. Right, you can actually look at this to also make sense,oh, that thereis this swampsoninInstagram,then swampsonunderscore inTwitter,is itthesame person? Like, for example, I am sure many of you were, who have common names, for example, Ponnuragan Kumaraguru is not so common, so, if you want to create anaccount probably you are the only one, you can get the handle, but let's take Shristi Gupta it is so common.

That if you want to create an account now any of the social networks, you are not goingto get (Refer Time: 28:21) Shristi Gupta. So, you are probably going to get Shristi Gupta 123,SonuGupta246,19o7,things likethat.This was,therefore;youcanactuallyuse this information also, that some parts of the handle is very similar, so, are they the same people?You could usethis Jacqard'sdistance, and there aremany other - Editdistance– there are many other measures by which you can find out whether the, how far is the handle from each other can be also used to say whether it is the same person.

(ReferSlideTime:28:59)



So, here is one version of the same slide that I showed you earlier, which is, usernames are collected, which is what we discussed now. Now we look at some features that you can actually use to put them together and then we'll find out what the predictions are. This is the same slide that I had about 5 or 8 slides before on the whole process of actually identifying whether the handles are same.

(ReferSlideTime:29:33)

So, here these sets of features that probably actually used in our work. Usernamecreationbehavior,she,shejust(ReferTime:29:46)categorizedthefeaturesintodifferent buckets. Similar length in terms of username, similar choice of characters, similar arrangement of characters ponguru in Twitter and t o n g o r e a 24 in LinkedIn and temporal behavioral feature also, evolution of length, I started with 6, now at 7, now it is 8, what kind of characters are changing, evolution of choice of characters.

(ReferSlideTime:30:22)

## Sample

- **User ID**: 595929421
- **Past usernames on Twitter**:
  - ["**bigeasye_**", "reezy11_", "epiceric_", "soulanola", "**swampson_**", "**hebetheeeric**", "swampkidd_"]
- **Past Usernames on Instagram**:
  - ["bigeasye_", "epiceric17", "swampson", "hebetheeeric"]}

If you see here, all these features, some of these features can be discussed here. If you look at the account details then probably this user started, both Twitter and Instagram is just the same, and then after some point in time the person had epiceric (Refer Time: 30:39) the third in Twitter is the same as the second one in Instagram. Fifth in Twitter is very similar to the third in Instagram. So, you can find this evolution and make something out of this also. Occasional reuse patterns.

So, common username, the same username being used and features, temporal ordering again,characters,howtheyareplaced;youcanuseallofthesefeatures,whichis, <mark>Paridhiis</mark> calling it more as the behavioral patterns across usernames.You can use these features to say whether the handles are same and the number of features sets that we had was about 56. And if you remember the account, if you remember the trust and credibility section, which is I think week 1 or week 2 that we saw,then we actually saw 45 features in TweetCred and in trust content in Twitter, we saw about 45 features thatAditi used in terms of actually finding out whether this particular content that is posted on Twitter is credible <mark>or not</mark> (Refer Time: 31:59).

So, now, we have the details from the users, details from the handles, what are the changes in the features, and we can all, we can put them all together to create a set of users, candidates sets, so to call, and then actually make a judgment, give the, give the output as, here is a probability of Sonu Gupta and Sonu Gupta 1 2 3 being the same is 0.9, versus Sonu Gupta and Sonu Gupta 0917 being the probabilityis about 0.4. So, you can actually make that output, that is the last part.

## Datasets

- Linking profiles
  - Twitter – Instagram
  - Twitter – Tumblr
  - Twitter - Facebook
- Past usernames available for both profiles:
  - 21,446 positive pairs, 21,449 negative pairs
- Past usernames available only on Twitter but current username available on other profile:
  - 112,451 positive pairs, 112,451 negative pai

Now, details. So until now it is more theoretical about how this could be done. Now let just look at specific things what Parishi did in terms of actually getting the data from multiple social networks and finding out how much we can actually do well. So, let us look at some specific examples. In this case we, are looking at data collected between Twitter and Instagram, Twitter and Tumblr,Twitter and Facebook. Past usernames were collected. 21,000 positive pairs, which is, details that collected from these social networks, and about past usernames available only on Twitter, but current usernames available on other profiles is about 140,000. So, essentially the idea is that the data was collected between multiples social networks of the current and the past user handles and how this was put together, and what kind of mechanisms was used to find out, whether these handles are same.

## Supervised Classification

**1. Independent Supervised Framework**

$U_S$: {'eenjolrass', 'isabelnevills', 'giuliettacapuleti', 'tobsregbo'}

$U_C$: {'enjoolras', 'isabelnevilles'}

$u_C$: {'isabelnevilles'}

$U_S - U_C$
$(U_S - u_C)$

**Username Set Features**
[Naive Bayes, SVM, Decision Tree, Random Forest]

Positive? → Same User

Negative? → Different Users

**2. Cascaded Supervised Framework**

$U_S$: {'eenjolrass', 'isabelnevills', 'giuliettacapuleti', 'tobsregbo'}

$U_C$: {'enjoolras', 'isabelnevilles'}

$u_C$: {'isabelnevilles'}

{'tobsregbo' 'isabelnevilles'}

$U_S - U_C$
(or $U_S - u_C$)

**Classifier I**
Current Username Features
[Exact Match, Substring Match]

Positive? → Same User

↓ Negative?

**Classifier II**
Username Set Features
[Naive Bayes, SVM, Decision Tree, Random Forest]

Negative? → Differe

And it is the same diagram that I showed you before, take the handles, understand some features, put the features together, and create the score. And that is what is done here - two methods are done, oneis you just do only the features, and the other method that we did was, do a classifier and then apply it and to find out whether it is the same user.

## Prediction

| Framework Config. | Accuracy | FNR | FPR |
|---|---|---|---|
| Exact Match (b1) | 55.38 | 89.34 | 0.00 |
| Substring Match (b2) | 60.99 | 78.46 | 0.00 |
| Independent [Naive Bayes] | 72.19 | 55.86 | 0.13 |
| Cascaded [b1→Naive Bayes] | 72.48 | 55.27 | 0.14 |
| Cascaded [b1 → SVM [Linear] | 76.74 | 45.16 | 1.65 |
| Cascaded [b2 → Naive Bayes] | 72.51 | 54.97 | 0.1 |
| Cascaded [b2 → SVM [Linear] | 76.84 | 45.16 | 1.2 |

Herearethedifferentmethodsthatwasused,whichis,exactmatch,right,exactmatchof thehandles,substringmatch,andthenclassifiers,differentclassifierswereappliedandif youlookatthefirstoneandthefifthone,aretheoneswhichhadthemaximumaccuracy, whichbasicallysays thatif youlookatthehandles,thewaythatthehandleslook similar,

,and iftheyare exactmatch,is veryhigh probabilitythattheyarethe sameusers,that is, veryless probability,that theywould actually be different users.And if we use the SVM classifier and then apply it on to find out whether there is same users using all the, using all the 26 features that we talked about, there is high probability that will be able to, about 76 percent is the accuracy to find out whether they are actually the same handles.

(ReferSlideTime:35:29)



## Prediction

A comparison of cascaded framework accuracy with and without Twitter-Tumblr instances

| Framework Config. [History on Both or One] | Accuracy | FNR | FPR |
|---|---|---|---|
| Exact Match (b1) | 55.38 | 89.34 | 0.00 |
| Cascaded [all network] | 76.74 | 45.16 | 1.65 |
| Exact Match without Tumblr (b1) | 66.17 | 67.51 | 0.00 |
| Cascaded [without Tumblr] | 91.20 | 16.60 | |

(ReferSlideTime:35:31)



## Measuring Volume of Sentiments

So, let me show you some, why, so, if you remember the motivation that I started off with is this sending this advertisement to people, I do not want to waste my money. Just let us go back and connect to the motivation also. There are, if PK or NPTEL wants to send out advertisement to all the students who are on Facebook and Twitter, we have to send this information to onlyone user onlyonce, you do not want to duplicate and waste money, that is the motivation I started, but the motivation can be many other things. There are some examples also.

So, in this case if you see, the sentiments of the user for a, let us take any topic that you take, the one on the left is on Facebook, the one on the right is on Twitter, you really want to understand whether sentiment, for example, you just look at this, you want to understand whether the sentiments of these people are expressed on Facebook and Twitter, are they same, and if they are same or if they are different, are they are same user.

So that, you can actuallymeasure that the negative sentiment of anytopic is nota sum of all negative sentiments in all social networks, but only the unique people that you wantto take a note of, (Refer Time: 36:52) right. Because if I say something positive in Facebook,andIamsamepersonwhoissayingpositiveinTwitter,itisnot,ifyou,you

cannot measure the positivity as twice, but it is only once because it is only one persons sentiment, right. So, that is another motivation. The other motivation also that, the other reason why this ==identity== resolution is an interesting problem is because you can actually look at, even law enforcement can actually use this. Which is, somebody uploaded a malicious video on YouTube and in there is another handle, which uploaded the same video on Twitter.

Now I want to find out whether ==it is the same person who is posting it.== Somebody is actually speaking against some people or some organization on Twitter. And there is another handle which is speaking against some persons or some organization in a different network. I want to know whether it is the same person, because they do not want to be wasting the time in assuming that it'stwo people and wasting time in finding two people, but it is only one person that they have to chase and catch, right. So, that is the motivation, that is also another motivation to actually find out whether these two handles are same or not, right. So, there is very interesting motivation for doing this work, and there is a lot of interesting things one could actually try out.

(ReferSlideTime:38:29)

## Conclusion

- Profile linking may be necessary for many organizations / needs
- Better profile linking is possible with past history of user handles

So, here the last take away from this part, which is profile linking may be necessary for manyorganizationsasthequestionsthatwesaid,Idonotwanttowastemymoney,I

want to actually understand, whether,how many people have posted, what is the volume of actually positivity or negativity, or I want to find out who is actually speaking online and to <mark>link</mark> users.

And the conclusions from this work are that essentially you do not have to only bank on the current handles that people have, current information that people have; even usingtheonesthatfromthepastcanactuallyprovetheefficiency,accuracy,ofthe–andthatis Paridhi, who is in the picture, who just graduated with a PhD from the work that she did on this topic.

(ReferSlideTime:39:29)



I am going leave you for this week just to try this out. We <mark>can probably connect this</mark> to the quiz that we have or homework that we have for the course also, but here is what I want to try,you to try.Take two of your accounts on your two different social networks, which is Facebook and Twitter. <mark>Let's just stick</mark> to only Facebook and Twitter. Just take these twohandles,list downall thethings thatyou can actuallydo,findoutvarious ways in which you can actually link these two accounts. List the features, features just we talked about, right? These ones or there could be many others also. List down these, and list down things that you will change in the profile to make them look two different account networks also. You could do both ways.

You could do list down things that you will change to make it two different account, or list down things that you will do to make it the same account. Right? Share it in the forum, let's see what you people actually come up with. ==I hope the activity is== clear. It is that take your Facebook account, take your, take your Twitter account, list down the features that are available, that you think you can actually connect with the accounts.And list - that is the first output - second output is list down all the things that you willdo to make it look same - that is the second account, second output. The third output is, list down all the things that you will do which will make that it is two different accounts.

(ReferSlideTime:41:16)

## References

- Paridhi Jain's Ph.D. thesis work

(ReferSlideTime:41:19)



So, that is all I had for this week that this 8.1. I will actually continue in a different topic when I start off with 8.2.

# AnonymousNetworks

Welcome back to the Privacy and Security in Online Social Media course on NPTEL. This is week 8, and this is the second part of the week.

(ReferSlideTime:00:25)



So, until now, the social networks that we are seeing is generally popular networks like Facebook, Twitter and these are called online social networks.And particularly we have also looked at Foursquare, which is a location based social network. Then I think briefly we have also talked about ephemeral social networks, which are networks where the contents that is getting generated can be actually removed after some period of time, where the contents are ephemeral, which it is like a snapshot network; where you post some content and after sometimes that content get's deleted right.

What we are going to look at this part of the lecture is something called anonymous network. Anonymous networks are networks, where it is not clearly visible or it is not possible to find out who is actually posting the content. So, we will go in detail about whatanonymousnetworksarewithsomeexamplesandIwillalsoshowyousome

research done, some work done, on finding out how anonymous network behaves,compared to normal networks like facebook or twitter. Some examples of anonymous networks are 4chan, Whisper, Secret, Yik Yak, Wickr, these are the different types of anonymous social network, there are many, there are many such networks that are available there here is only a small list.

(ReferSlideTime:02:10)



WhydoyouneedanonymoussocialnetworkSo,wealreadyhaveFacebook,wealready have twitter, why you might need a network or network of the category of anonymous network or network that gives the preference or gives the facility for having anonymity. Increasing awareness of privacy, so people are getting to know more and more about privacy,people are getting to or people want to have more privacy on online social networks. So, therefore people are looking for networks, that will give more anonymity. And there werealso incidences like Snowden; projects like PRISM were the information that ispubliclyavailable ortheinformation that isavailable totheseorganizations canbe used for other reasons also.

And of course, there is an incident in India, where the some post was done and that post called actually viral and there were consequences of the post also. So, therefore many many incidences around the world, which are happening, which is expecting, which is making users who use social networks expect more privacy, expect anonymity in the networks.Becausefor example,IfI doapostonfacebook,ifI doa poston twitteritis

actually very clear that it is pk ponnurangam dot kumaraguru dot or ponguru in twitter is actually doing the post. In fact, if I wanted to say something on social networks, but I do not want to be attributed to the post then I would actually use these anonymous social network.

(ReferSlideTime:04:13)



So, here is thatslide screen shotof the website whisper dot sh in thecontent is organized asinthetoporderinthisimage    popular,latest,Lol,confessions,relationship,Ohmygod    and they create these categories so that the content that is uploaded on whisper gets into1 of these categories and the URLis whisper dot sh. I'll let you to actually play around a little on of the website, create an account and see how the accounts work.

(ReferSlideTime:04:55)



Here is a URL, here is a video, which describes some features of YouTube. We take a lookatvideonowandthenIwilldescribesomedetailswhichisfromthevideoandother features of whisper.

(ReferSlideTime:05:22)



Secret,s lies and plenty of spam on this super popular mobile app - whisper.Whisper is a confessional app that encourages you to post your secrets behind a screen of anonymity. You know kinda like that other website - Post Secret, where in users can submit those deepdarksecretstheywouldnoteventelltheirbestfriend,likeIsecretlytooknude

pictures of my best friend. But whisper, which been around for a couple years and it is being steadily gaining popularity is just as much way to people to connect around us on flattering,embarrassing,tabooorsometimesdisturbingconfessionorsometimesnoneof those things pretty often.

Whisper does not harvest your email or contacts and screen names are less prominent within the app you can also change it whenever you want along with the pin that takes place of the password. So, there's definitely increased premium on anonymity than most social networks. As per your deep dark secrets, those you can post by hitting the plus. Type your whisper and it will auto generate a stock photo to go along with it. So, yes is not just a confessional, but it is meme generator which you can share via email, SMS or social.

(ReferSlideTime:06:18)



And you can private message, there is premium messaging for certain users, trolls and spammers according to the F.A.Q, but for everyone else it is free as long as you play by the rules.Forthat reasonyou have all seen posts like these, alot ofthem peopletrying to hookup or these - message your favorite singer and if you are dumb enough to fall forthat,you areprobablyundertheageof10andgodhelpyou,butwhereverthereissecrets there are also bad apples and sometimes lies.

In September for example, someone posted a supposed murder confession on PostSecret, promptinga frenziedreddit searchfortheselfprofessedcriminal.Andtilldateno

crimehas ever been found associated with that post. And this week in Arizona, a cop was arrested for having sex with a minor he met on whisper app after she posted that she wanted to get pregnant. But considering the allure of posting secrets and knowing other people's, it's unlike to slow down whisper for now.As always you can let me know what you think I am onTwitter,Facebook, Google orVK on anniegaus and you can get a free netflix trial with a signup on Netflix dot com slash wtbd, thanks for watching.

(ReferSlideTime:07:16)



Now, that we are seeing the video, the video actually talks about how whisper is being used, what kind of users get on whisper and what kind posts they do and how whisper actually works in creating some content, it actually gets merged ontoimages andyou get posts and creating memes in other terms. So, the way that people react to the post on whisper is by hearts and also you can chats on the post that you make. Again please remember all of this is going to be anonymous.

(ReferSlideTime:08:02)



So, terminologies that we'll see to understand the rest of the lecture, we need to understand some terminologies, whispers or the posts, replies or I do a post and you are actually replying like a comment in Facebook or a reply in Twitter. And the posts are anonymous you really do not get to see it is ponguru . I may have an account, which is called professor, teaching computer science or anything that I wanted to keep that is the username and interestingly whisper also allows you to back with probably have seen video also, whisper also allows you to change the usernames as anonymous as you want and more number of times also. So, that makes it much more difficult to go back andlook at the person who posted the content.

And whisper does not associated any personal information of the user id, it is not collecting any information and does not archive any user history, which at least that'swhat they claim, it does not support persistent social links between users. The person who hearts at that, the person who replies it, the links of the users are not kept, where as if you remember the homework and thequestionsthat you haveseen in thepastwherein the context of facebook or twitter.

The content for all the relationship between the users are stored as a graph and you can analyze those graph, also retrieving the graph from twitter or facebook and use these graph to make some inferences. Heart a message anonymously may also use just in (ReferTime:10:02).AheartisbasicallytheonethatIshowedyouintheslide,likethe

==like== in facebook. If in the private messages against or this in the video that I had a few minutes before, which showed private messages also you can actually post private messages between the users.

(ReferSlideTime:10:28)



Thatisthescreenshotfromwhisper.

(Refer Slide Time: 10:36)



So, what we are going to look at is we are going to try an answer these four questions. How do whisper users interact in an anonymous environment, how is the interactions on whisper?Dousersformcommunitiessimilartothoseintraditionalnetworks,likefor

examplepeopleinteractiononfacebook,howisthisdifferentfrompeopleinteractionson whisper or twitter. Does whisper's lack of identities eliminate strong ties between users, which is if I do not have strong ties which is if you do not know that PK is talking toyou, does it eliminate the strong relationship that you and I would have. Let's take bothof are on whisper,you doyour postand I come react toit, Ido areplyto it, if you donot know it is PK, who is the faculty at IIIT Delhi or some profile that I have, if you do not knowthatitis me,willyou continuetalkingtome? Isthereisastrongerrelationship that happens.

For example, you could also see in twitter or in facebook that some people are very stronglyconnected. Forexample, ifIdoanypostthataresome setsofpeoplewhowould always like it, who would always accurately make a comment or reply or retweet. So, thosethatarebasicallycalled strongtiesandthatdoesitexistson whisperisthequestion we have to look at.

Now also whisper, because of being anonymous does it eliminates stickiness critical to long term engagement. Stickiness is basically is a factor bywhich you are actually gluedon to the network, more and more people get connected to it or single person is actually spending more time on the network. I know that is clear those are the four goals that we have for the rest of this lecture where we will take one click one particular network inthis case whisper, we will actually try an answer for these four questions using somedata, using some inferences that we draw.

(ReferSlideTime:13:14)



So, data that was collected for doing this analysis is from 2014. And of course, that whisper does not have an apiso, data was scraped and what all they include. They included whisper id, which is like a post id, time stamp when the post was done, plain text of the whisper - thetext that was on thepicture, author's nickname which is the your handle, names so to say in the traditional sense, alocation tag if it was available, number of replies for the whisper and of course, the likes is the hearts that we talked about. So, that is clear simple to collection I think you all of you have seen this kind of data collection the past in all the networks that we have seen. Take it a network collect some basic data, do some analysis and answer interesting questions that actually makes sense.

(ReferSlideTime:14:26)



Data collection again, so 9 million whispers, 15 million replies 1 million GUIDs, which is global user universal identifier, which is the id for every user like you've seen in the twitter also.So,theusersgetoneuniqueidwhichiswhatwascollected. So,interestingly the team that worked on this work also interacted with the whisper team. Where they actually talked to them about the data that that they were collecting and about this universal identifier, which they were able to convince the whisper team thatusing this id you could actually go back and find out which user did what. So, the user GUID concept was removed on June 2014.

(ReferSlideTime:15:27)

(ReferSlideTime:15:34)



## Data collection

- 9,343,590 whispers
- 15,268,964 replies
- 1,038,364 GUIDs
  - Global Universal Identifier
  - Makes it possible to track user, but was removed in June 2014
- Interacted with Whisper team

So, they collected the data, the researchers looked at what is going on whisper and then went and had discussion with whisper team to remove this. That for they actually wroteabout.

(ReferSlideTime:15:40)



## Data

- 55% of whispers receives no replies
- 25% have a chain of at least 2 replies

So,now lookat analysis again, we havedone this inthe past soI'mgoingto goslowlyin terms of what analysis, first time using this data, what kind of inferences to be true.And allconnectedtothesefourquestionsthatwe have.Inthexaxissuchasthetime,itisthe

time, inthis casebetweenFebruaryandMay,thatis wheretheycollected thedata.Andy axis is the number of posts per day.

And they actually look at 3 different types of posts, which is one as whisper, so to say what is content that is getting generated, one is the replies, which ishow many repliesare being posted for the particular whisper. And there is also third category of whispers being deleted. We get to this deletion later, which is also interesting problem, which is that, when in twitter also, we have more recent studies in 2016, people have seen that lot of content that are posted on the social network gets actually deleted for whateverreasons that the users are deciding to.

So, in this case, in whisper case 55 percent of whispers receives no replies, people just post content and nobody even replies to these posts. 25 percent have a chain of at least 2 replies.Only25percentsotosayactually weshouldreaditthatway,the25percenthave achain of atleast2 replies. 55 percentof thereplies, 55 percentof thewhispersdon'tgeta reply.

(ReferSlideTime:17:27)



Time between original replies this is also interesting thing how quickly are theresponses towhisperthatisposted.Howquicklydothepeopleactuallylookatthepostthatisdone and how they reply and what they reply. In this case, we are only looking at the time we are not looking at the content. So, if your x axis is time again less than one minute, one minutetoonehour,onehourtooneday,onedaytooneweekandgreaterthan1week.

That is the x axis, y axis is fraction of replies, fraction of replies it shows you whatis the proportion of replies that the whisper gets. 54 percent of replies are within hour of those original whispers.

You can add the first two bars which is less than one minute and one minute to an hour, thiswillshowyou 54percent,54percentoftherepliesarrivewithinonehour,94percent within oneday.Basically, shows that if they do notget aresponsein one day they do not get it. More than half of them get a response within one hour.One point three percent of replies arrive within a week or more that is the last bar on the graph. So, essentially the conclusion is that if awhisper does notget attention shortlyafter posting, it is unlikelyto get attention later. Understandably, that because it is an anonymous people kind of post content, it gets little bit of attention and then dies off. This is similar to other networks also that we have seen.

(ReferSlideTime:07:16)



So,apostper user, which isjustthexaxis iswhispersandreplies peruser.Whispersand replies per user, which is how many times user is doing it and there is two lines in the graph, one the dark line, which is the whisper and the dotted line, which is the reply, they axis is the CDF you have seen many of the CDFs before cumulative frequency of user. Here it is basically showing that 80 percent of the users post less than 10 total whispers, which is actually pretty bad if you just look at the networks, 80 percent of the users post less than 10 total whispers and replies.

Which is if you flip it and see it is probably looking at 20 percent of the users areactually the people who are actually very active or in ==another== sense less percentage of people are the ones who are actually doing maximum number of activities in thenetwork, which wehavealreadyseenin othernetworksalso. 15percentoftheusersonly post replies, but no original whispers, which basically again shows that less fraction of people posts replies that they do not create original content, only ==look at what users are== doing and then they are replying to it.Thirty percent of the users only post whispers. But no replies, they're just the people who are creating the original content, but they actually do not reply to any of the content, reply or react to it.

(ReferSlideTime:21:32)



And that is clear, that is basically has two analysis that we saw, one is how much you attention is the content posted ==on whisper is== getting.

(ReferSlideTime:21:44)



And how much time this it get take to get their attention and then what are the level of activity do users have on these networks.

(ReferSlideTime:21:51)



So, now,we are going to look at the topic that we have seen more in the past also which is network analysis. We have also have tutorials on this topic looking at how you can actually use metrics, which are developed in network analysis to make some inferences. Here they took whisper, they also took random users from facebook, compared also to twitter.So,thefirstcolumnisgraphwhisper,facebook,twitter,secondcolumnisnumber

of nodes - 690,000 nodes in whisper, 707,000 in facebook, 4,317,000 nodes in twitter, number of edges, average degree, clustering coefficient, average path length and assortativity coefficient. I will go through average degree clustering coefficient and the rest and tell you what does this is mean.

So, first if you look at the column four, which shows you average degree, the average degree is actually very high for whisper compared to facebook and twitter, what does it mean? This means that I am connecting to lot more. So, which is also connected to the clustering coefficient, but this says that users interact to the large sample of further users which means any user in whisper is not restricted only to a set of people.

But they otherwise interact, they interact with the large set of people. When you compare it to the facebook or twitter that we talk about the interactions are much closely connected, it's mostly with the followers that you have or probably people who mention you or probably the hashtag that you interested in. Facebook is mostly of friends. So, if you look at facebook, where it is only 1.78, for twitter is it is 3.93. Now let us look at whisper, it is 9.47 the degree in which they interact with the users in whisper is prettylarge.

Whisper users are likely to interact with complete strangers, look at the clustering coefficient. If you remember, what clustering coefficient tells, clustering coefficient just lets you to say how the graph looks like, whispers or whisper users are likely to interact with complete strangers who are highly unlikely to interact with each other also. So, if you look at the values it is pretty low,0.033 compared to 0.059 and 0.048 in twitter. So, theyhave also looked at 100 random nodes,average path length calculated, shortest path was the shortest average path among the 3 is actually for the whisper; if you look at the column average path length 4.28, 10.13 for facebook and 5.52 for twitter.

So, this just says that average length in the graph, if you take the whisper graph is actuallythelowestandthereisaveragepathlength.Basicallywhatdoesitmean itmeans average degree being highest, clustering coefficient being lowest, average path length being lowest is inferred that is it is the random graph. People interactions are completely random, there is no specific small world phenomenon that happens in a network like whisper. That is a good difference from the traditional networks that we have seen like facebook and twitter.

(ReferSlideTime:26:17)



## Network analysis

| Graph | # of Nodes | # of Edges | Avg. Degree | Clustering Coef. | Avg. Path Length | Assortativity Coef. |
|---|---|---|---|---|---|---|
| Whisper | 690K | 6,531K | 9.47 | 0.033 | 4.28 | -0.011 |
| Facebook | 707K | 1,260K | 1.78 | 0.059 | 10.13 | 0.116 |
| Twitter | 4,317K | 16,972K | 3.93 | 0.048 | 5.52 | -0.025 |

- Assortativity measures the probability for nodes in a graph to link to other nodes of similar degrees.
- Close to zero → random graph

Now, let us look at assortativity. Assortativity measures the probability of nodes in a graph to link to other nodes of similar degrees. So, the more the value that is closer to 0, or the less the value is, it is actually assumed that the graph is a random or you can infer the graph is a random graph. If you look at it, it is the lowest value minus 0.011 and this is the assortativity coefficient of the whisper, for all the 3 graphs. It basically says that it is a random graph. I know that is clear that. So, essentially the conclusion from this network analysis that you can draw is that whisper network is a random graph, whisper network people actually interact with the random people and the graph is actually pretty sparse.

(ReferSlideTime:27:27)



So, now another interesting thing that they did is to study what content was getting deleted on whisper. So, for this they collected the 1 point 7 million whispers, that have been deleted in 3 months, 18 percent of the content deleted. 18 percent of the total generated content was deleted from whisper where as compared to twitter, which is only 4 percent.

Andthisbegsthequestionwhich isthatfor whisperwhyisthispercentage high,because anonymous content you posted today you feel like there is some problem you feel like you created the some contents which others do not like to see or you do not want any attribution to you, even though it is an anonymous network, still you want to get itdeleted. So, there is higher proportion of content generated on whisper which is getting deleted.

Content moderation, so what moderation in the context of whisper is, the analysis that they did was, they extracted keywords from all whispers, which put all the text that was created that was drawn from or collected from whisper, removed the common stop words, removed words that appear in less than 0.05 percent of whisper, that remove all the words that actually people care about or people have used it, compute deletion ratio they take calculated a value, which actually says that number of deleted.

Whispers with these words, by all whispers with this words. Which is essentially to say that what is the chance if the word appears in the post; and what is ratio for this word getting deleted; what is the ratio that whisper that has this word getting deleted. And it ranks the words with deletion ratio, they rank basically all the words which are with the deletion ratio and they looked at top ten and bottom, top keywords and the bottom keywords, here is the table which actually shows you the top keywords and the bottom keywords.

We will see it in the next slide, they ran this methods for 9 million original whispers. They saw the 1.7 million are deleted, 2324 keywords ranked by deletion ratio, manually they put them in categories to see which categories are largest amongst of deletion the lowest number they rank them in the tables here.

(ReferSlideTime:30:30)



The categories that they had a sexting, selfie, chat, topic, emotion, the top 50 keywords most related to the deleted whispers. And the top 50 keywords least related to deleted whispers. The top points 50 keywords that are in the top, the bottom of the table gives you the bottom that was there on the deleted whispers. Essentially showing that you sexting, selfie and chat are the categories, which were most frequently deleted, and emotion, religion, entertain, life story, work, politics and others were the least deleted categories from the whispers.

(ReferSlideTime:31:34)

This graphis actually showingyou, Imean let us lookat thewayinwhich thedeletion is happening, how much time and relationship for actually deletion. 70 percent of the deleted whispers are deleted within one week after posting. So, that is the first graphfrom the left, which is 70 percent of the post, x axis is week y axis probability of getting deleted, proportion of whispers getting deleted, 70 percent of the whispers are deleted within one week after posting. The right side shows you delay before whisper is getting deleted, that this one, 2 percent of the whispers stay for more than a month, if you see it had a four weeks that is the graph from the left.

(ReferSlideTime:32:44)



Now, lets look at the content analysis on the right. fine grained analysis, recrawled for 200 thousand latest whispers, they were actually interested in trying to understand how many hours, this was a week the first graph, what is the analysis in the hours that is what they are interested. They actually found that 32153 was deleted, peak deletion was between 3 and 9 hours, which is any post on, if it is was both get deleted is between 3 to 9 hours. Majority deletes within 24 hours. So, it is even if youzoom in to the data forthis one week, majority of them are actually getting deleted within the first 24 hours.

User interactions is another interesting analysis that they did, which is how frequently how users actually interact, which is 2 handles in this case are actually interacting between them. This graph is showing you on the x axis, geo distance between the pairedusers, of course these are the locations that people actually disclose, percentage of user pairs, what is a number, what is the percentage of user pairs which are actually interacting.

The colors are blues is two interactions, yellows are two to five interactions, red is 6 to10interactions, anything thatwas above10interactions, which wasactuallygiven green. Youcould alreadyguess,that thenumber ofinteractions whichis higher,is actuallyvery lowinthewhisper,whichagainusingthenetworkanalysis, usingthethingsthatwehave alreadyseen,youcouldactuallymaketheinference.Thatiswhygreenisverylowonthe graph. 90 percent of that the two users are co located in the same state, 75 percent have their distance which is less than 40 miles. So, this basically shows that the users are also co located very closely, within 75 percent have their distance less than 40 miles.

Smaller user population in same nearby area, higher chance of encounter, so if you look at a graph less than 10 to 100, 100 to 100,000, greater than 1000 on the left and then combined post of paired user which is on the right. The left side is showing you user population in nearby region. The right is showing you combined number of posts of paired users, More whispers two users post, more likely they encounter with each other. If the users are likely to post more; they are likely to interact more also. So, that'sactually looking at a the right graph which is combined number of posts of paired users, more whispers two users post which says if you and I want to interact, if you and I actually generating more whispers that's more likely that you and I interact, that is the inference that you can draw.

(ReferSlideTime:36:37)



Now,look at how users engage in thesein this network on whisper.So, here x axis is the time of week, y axis is the accumulated number of users and the two data points that are drawninthisgraphis theexistingandthenewcontent that isgetting generated. Roughly 80,000 user per week are interacting daily,new posts in the entire network remain stable thatwewillactuallyseeinthenextgraphalso.Howthisrenames same isactually,ifyou see the new content that is generated per week on the graph, they actually seeing to be very same across.

Even though there are more users that are getting added to the network, it does not look likethenewcontentthatisgettinggeneratedisactuallyincreasing.Thatistheinteresting conclusion that you can actually see in this graph, daily new posts in the entire network remain stable despite new users joining. That you can see actually accumulated number of users is increasing.

So,thisbasicallyshows thateventhoughthe users areincreasing,whichmeans theposts should be increasing, the engagement should be increasing. But it is not, this basically shows that that are lot of people who are getting into the network, generating some content and then lot of people, who are known to be already in the network are not generating the content. That is why the proportion of the net the content is getting generated is always remaining the same. Even though there are more users are added toit.

(ReferSlideTime:38:37)



User engagement, so, here is the graph that I will actually show you. Number ofwhispers and replies this is weeks,this is also in weeks.So, it is kind of the same kind of graph and you will see number of whispers and replies by both new and old users if you see, the top, the new users make a twenty percent of the contribution in the content. Content by new users do not grow, right. So, this is time and number of whispers and replies for that particular week.

(ReferSlideTime:39:17)

The earlier graph was the cumulative number of users right, accumulated number of users. So, the existing user which is the light without the check that the graph the bar is actually, which is risen which is showing you that the number of users are actually risen.

(ReferSlideTime:39:38)



The next graph is essentially showing you that the post and the replies that, that is there in the network is actually pretty constant even though the number of users are increased, correct. So, that is the conclusion that they had in the user engagement, that is basically kind of addresses the question that we started off with, which is do user's engaged differently in a network like whisper.

So, the conclusion is clearly different from traditional social networks. We saw that the average path length is different, clustering coefficient is different, we saw that the deletion is actually pretty high and inferences like that. Without strong user identities or persistent social links users interact with strangers which is also derived from the conclusion that user is interacting with any random people on the network right. There was not a persistent relationship between the users, moderation is of course necessary because content is getting deleted very highly.

So, that is all I had for this week, in this week, we saw what 8.2 we saw anonymous networks,whicharenetworkswhereyoucanpostcontentwhereyoucan maintain ahigh anonymity. In 8.1, we saw how to actually do identical resolution with multiple accounts given to us. That is we get. I will see you in next week.

# Unit-5

## Privacy in Location Based Social Networks Part-I

WelcomebacktoPrivacyandSecurityinOnlineSocialMediacourseonNPTEL.Thisis week 9. So, what we will do today and in this week is that we will look at some of the research which was done in terms of just looking at the papers itself and going through the paper in terms of different techniques that are applied. The goal here is that next couple of hours, what we will do is we will get you to actually look at research papers written on the topic and we will go through the same analysis that you have done across the course for you to get a sense of how the analysis that you have done fit into actually making some interesting inferences.

Where do people start, how do theywriteapaper whatallthings fits into thepaper,what allanalysis that theyhave doneessentially,itis looking atthecontentthatyou have seen in the past, but in terms of the structure of the paper itself we will go through some of them.

(ReferSlideTime:01:22)

So, location, the first topic that we will go through is location based privacy problems. Location based services on online social media, there are many actually and there are some which are very popular which are like Foursquare, Yelp, Gowalla, Facebook, Twitter these are the different social media services that are actually pretty popular in terms of giving the location based services, for example, in Foursquare you could actually see where is the next, let us take petrol pump, in the directions that you're travelling. Secondly, in the Yelp you can look at where the restaurant or places that you are interested in, what kind of reviews do they have, Facebook you can actually looked you can actually do check in into a location in Facebook.

Similarly,you can do the geo location information shared on Twitter.So, essentially you must have seen check-ins from people on Facebook like XYZ is on T3 Indira Gandhi AirportinDelhitravelingtoXYZplace.So,thatislocationbasedservice,youopenyour phone, you check into this location or post this status updatewith thelocation <mark>updated init.</mark> Foursquare, Yelp, these are very, very popular services which actually do location based services, these are called location based social networks.

(ReferSlideTime:03:18)



So,ofcourse,allofthemhavesomesortofprivacyconcernstoo.Iam,Iamsurebynow youalreadyrealizedtheprivacyconcernsinlocationbasedserviceswhichare,wheremy

location is shared, if my location is shared there are going to be concerns accordingly, that is, somebody else will get to know where I am, if the information is given public, then many more people actually, more than just your friends get access to your location at that given point in time, right.

So, that is the privacy concern and of course, every each of these social media services will have its own privacy concerns also.

(ReferSlideTime:04:08)



So, just to give you a sense of what are the perceptions in terms of actually these location based services. If you remember, there was a study that I referred earlier in the lectures also, where 10427 people were asked some questions about different aspects of privacy. One section in that was online social media, here is a question that was asked in the study, the question reads, what privacy settings do you have for the following information on Facebook?

Please provide your responses to the best of your knowledge. Location - not shared, friends, friends of friends, network, everyone, and 'I have customized it', those are the different options that was given for the question. I do not know where each of you will fit in, but if you look at it, here is a distribution of responses that was got for this particular

question - 33 percent to friends, friends of friends is 12.5 percent, network is 5.7, everyone is 39 percent and customized is 3.1.

It just says that location, is been. what privacy settings has been set up is for everyone, right, location information is shared maximum to everyone on Facebook, and if you put everyone in friends that is a lot of percentage of responses which where the information is been shared. So, that is a point I wanted to get across, which is that information that is shared through location can be actually used for things beyond what you think for now also.

(ReferSlideTime:06:26)



## Perceptions

- While travelling (i.e. in roaming), the mobile service providers use regional languages to present information e.g. user busy, phone switched off. For example, if your phone connection is from Delhi and if you are traveling in Mumbai, the messages are presented in Marathi. Would you consider this feature as privacy invasive?
- Strongly Agree: 10%
- Agree: 44%
- Neutral: 22.9%
- Disagree: 19.2%
- Strongly disagree: 3.8%

Perceptions,again.Anotherquestioninthesamesurvey,samedatacollection,whichwas asked to get some responses from participants - While traveling the mobile service providers use regional languages to present information, that is, user busy, phone switched off.

For example if your phone connection is from Delhi and if you are traveling in Mumbai the messages are presented in Marathi. Would you consider this feature as privacy invasive, I am sure you got the question, the question is simply that I have a number whichIboughtitinDelhiandIamtravelinginMumbaiandmyphoneisswitchedoff.

When somebody calls me, the messages are actually want to be saying that the phone is switched off, but it is going to be giving that message in Marathi.

The same thing changes when I go to Kerala, this number is going to be actually this messageis going tobein Malayalam, which isbasicallyrevealing theinformation where you are at when somebody is trying to call you. This can be privacy intrusive because it just says that the exact location where you are at least in the regional language the state where you are. So, if you see, would you consider this feature as privacy invasive? 44 percentof themsay,thatit is privacyinvasive, 42 percent neutral, disagree is 19 percent, strongly disagree is 3 percent, and strongly agree is 10 percent.

So, essentially if you look at just the people who are agreeing, it is 54 percent, people who are disagreeing is about 22- 23 percent and rest are in neutral that just to share that and allofthisdatawascollectedonlyinIndia.Theperceptionsherewearetalkingabout person living in India who is thinking about these problems.

(ReferSlideTime:08:34)



So, that is basically a sense of what location based privacy is. Location based privacy services are what kind of issues you can have, but what we will do is as I said we will actuallylookatsomespecificresearchthathavedoneintermsofinformation.Interms

of analysis, data collection, what the topic is, we will take the paper now and will go through the paper only in detail about what was done what kind of analysis was done, what kind of inferences was drawn.

This will actually it help you to get a better sense of how to use the social media data techniques that we have learnt in the in this course until now in terms of actuallymaking some inferences. So, now, here is a paper that we will look at and spend some time onthe paper.The title of the paper is 'WeKnow WhereYou Live: PrivacyCharacterizationof Foursquare behavior'.We Know Where You Live: Privacy Characterization ofFoursquare behavior

(ReferSlideTime:09:52)



This is what we will do. We will go through the paper content and look at what is mentioned in the paper and go through them. I will actually give some insights about what is going on. So, as I said, we are going to look at location based social network. Foursquare is one of the popular ones. The way Foursquare works is, it gives incentives to users to use the check-ins specific places.

So, a visit is a check-in.You go to a place, for example, IIITDelhi is a venue, you come toIIITDelhi,whichisavenue,andthenyoudoa check-in,whichisyouarevisitingthe

place and mayorship is for frequent visits. Mayorship is nothing, but it is the person who is being to that place for most number of times in the last 60 days, which is. he is checked-in that location for the most number of times in the last 60 days. that is the mayorship and people do mayorship for many reasons, there is actually incentives that are done.

For example, if you go to any malls, if you go to places they are actually giving you -and if you are a mayor of that location - you can actually get more discounts or you can actually get some parking spots free for a week or so. So, being a mayor is actually giving you some incentives. So, check-ins is the action to be in that place, venues is the location or the places, mayorships.

They can also leave tips, tips at specific menus are the is the information that people put in to the Foursquares saying I will checked in into this location, for examples, ==SaravanaBhavan,== I had food, food was pretty good. That is a tip and if you agree on the tips, like the like in Facebook or like the re tweet or the like in Twitter again, there is something called as done, d o n e in Foursquare which is actually again saying that I like this tip.So, if you see check-ins are actually available only for your friends, but the list of mayorships, tips and done of users are publicly available to everyone.

So, this basically allows us to collect this information ==and do some== interesting analysis which is what we will actually look at in this part of the lecture - collecting data from Foursquare,analyzingthemandmakingsomeinterestinginferences,particularlylooking at privacy issues in Foursquare and given the title of the paper, particularly looking at, can we actually find out where people live, from just the check-in, from just their Foursquare behavior.

This paper basically explores these publicly available features – mayorships, tips, dones, and their usage for informational leakage, but interesting part of this particular paper is that it actually uses the data of entire Foursquare. At that point in time, which is 13 million users, and the paper kind of concludes that - there are many interesting conclusions in this paper we will look at all of them in detail - but in the abstract, they talkabout,our==results==indicatethat,oneeasily,onecaneasilyinferthehomecityof

around 75 percent of the analyzed users within 58 hours. So, that is actually is privacy intrusive if I can actually tell you where you live just by looking at your check-in location and the Foursquare information then it is actually privacy intrusive.

In this case, only publicly available information is used. There used to be a websitecalled please rob me dot com (pleaserobme.com), which they took it down after some time. This websitedidsomethinginteresting-whichispleaserobdotrobbedandme dot com please rob me dot com - the creators of this websites basically looked at the tweets and if a user talks about location x or if the user created the account user location that is the information that you will actually get in Twitter, with that information they were actually saying that a person created an account in Delhi and he is talking about a weather in Chennai that is a probability that he is not he is not at home.

Or if you have checked in, if you have done a post in tweet with your geolocation on from Delhi and in another hour or so, you are talking about actually weather and couple ofhoursyouaretalkingaboutweatherinChennaiagain,thatisaprobabilitythatyou are notathome,you movedfromDelhitoChennai.So,theywereusingthisinformationand the tweets that were of this category which is user from location x and weather or other information, for example, even traffic right, you are from Chennai and you actually post you're posting traffic about Delhi and the tweet is also saying that I am going to this place and that is a heavy traffic, there is a higher probability that you are actually in Delhi.

Using all ==this== information they find out that this particular user who is posting this tweet is not in his or her home in location and therefore, they would actually take the post and show it in this website called please rob me dot com forburglars to goand robthe home. And again another leak of information, another impact of information about your location, can be actually seen from the example of please rob me dot com.

So, then I think the paper talks about in general about ==what== Foursquare is and for the benefit of the students in the class I have highlighted the parts that have actually something that I am going to be looking at, but feel free to look at the entire paper, but for thebenefitof timeconstraintandthelecturesalsoI am goingtolookatonly theones

thatIhave highlightedhere.

(ReferSlideTime:17:54)



So for people who are curious about writing papers like this, the structure it is – abstract, then there is introduction, introduction you generally talk about what the problem of attack in the paper is, give some background about the domain - in this case it is online social media - then talking about location based services, then giving some information about, so it has to be both, giving information about the topic - location based services, give more of quantitative numbers also saying how many people are using it, what level of impact is it making and things like that.

(ReferSlideTime:18:20)

Andthenyou quicklymention aboutthe methodologythatyou did.

(ReferSlideTime:18:45)



So, essentially introduction would be just a shorter abridged version of the rest of the paper which is covering methodology, inferences, contributions, and conclusions. And the question that the paper specifically addresses is, despite being a private data that the user may choose not to reveal, can we still infer the home city of a user in Foursquare from our mayorships, tips, and done, which are publicly available information.

So, that is the question and then in the paper there is something called as related work, which is focused on only the question that we are asking in introduction, you kind of generallymotivatefrom30,000feetheightabouttheproblemsaying,thisistheproblem, this is what we are planning to work. In related work, you focus on only the work that you are going to show in the paper and talk about past literature which are connected to that.

(ReferSlideTime:19:42)



studies have focused on geographically referenced information addressing aspects such as understanding why users share their location [16], human mobility patterns [2, 3, 14], user profile identification [10, 17], event detection [15] and analysis of a city urban development through check ins [5].

The information sharing in LBSNs and online social networks in general also raises concerns about exposure of user private data, touching privacy related issues. For instance, some studies have shown that it is possible to infer user implicit data through explicit information shared in such sys-

So,inthiscase,authorsaretalkingabouthowpeoplehaveusedhumanmobilitypatterns, even direction analysis of a city urban development, through check-ins. There is a lot of workactuallyinterms ofusingFoursquare datatofindoutthenatureofthecityandhow people have actually used the social network in a location based social network like Foursquare. People have actually used, researchers have used the information of Foursquare to design a city.

So, you can actually look at -if you are interested - you can actually look at this project called livelihoods.org. This is a project where they are actually using Foursquare information to see how people actually move in a given city and can we actually re-design a city keeping this information in mind.

Now,welookatsomethingmorecloselyinterms ofjustFoursquareDatasetitself.

(ReferSlideTime:21:13)



To understand the analysis better you need to actually be clearer about some of the terminologies. So, here I am going to actually explain them, something we have already done. So, check-in is, the, where members can share the location with friends and followers through check-ins, that is the check-in. Check-ins are performed via devices with GPS when a user is close to a specific location, which is a venue, which I have said before. So,venuescanbeairport,restaurantand monuments, you comeand checkin that particular location.

(ReferSlideTime:21:55)



So, the interesting aspect of Foursquare is that they have made it more - gamified it - which is basically turning the platform where users actually get more incentive, more addicted through this gamification nature, which is, user get badges, mayorships here. Badges are first time person who came gets a newbie badge, so to say, and somebody who checks in to the system late in the night, early in the morning, there were different badges that were that could be actually given in system slide force for mayorship.

Secondly, works already explain right. So, the text says badges or like medals given to userswhocheckinataspecificvenuesorachievesomepredefinednumberofcheck-ins.

(ReferSlideTime:23:07)



Tips as I said, users can post tips at specific venues commenting on their previous experiences when visiting the corresponding physical places. Tips can also serve as feedback recommendation or review to help others user other users choose places tovisit. So, the idea is you got to restaurant you have food you want to actually give feedback to others saying the food was good or the other way the food was bad. So,when others want to get to this restaurant they can actually use feedback to make it thatis a tip.

So, when visiting a venues page after reading a previously posted tip, the user may mark it as done or to do in sign for agreement. When the tips content or intention to visit the location in the future. So, it is essentiallysaying that I saw this tip the tip is actuallyvery useful for me. So, I take a button which is done or do you say that I want to keep it is to do which is I want to go to this place in future.

Now,letuslookatthedatasetthatwascollected.

(Refer Slide Time: 00:12)



So, the dataset is about total of about 13 million users collected ran from August to October 2011, the total of 13 million users. It is almost close to the entire Foursquare in terms of the number users that are using it, and the total dataset contains 10 million tips and what we are interested in the data particularly the questions that we are asking as I said, we are interested in mostly the tips, dones and mayorships, because that is the information that is publicly available.

What can you use this information for in terms of the actually getting the location of the particular user. It is the total number of tips that are available in the dataset is about 10 million, total number of dones is about 9million, and mayor ship is about 15 million and different venues that are available are about 15 million again.

(ReferSlideTime:01:56)



Table 1. Dictionary. GI = geographic information. UHC = User Home City. VL = Venue Location.

| Statistics | UHC | VL |
|---|---|---|
| # in dataset | 13,570,060 | 15,898,484 |
| # valid GI | 13,299,112 | 11,683,813 |
| # valid but ambiguous GI | 359,543 | 2,868,636 |
| # non-GI | 244,233 | 4,214,671 |
| # empty entries | 26,715 | 0 |

Table 2. Quality of Geographic Information.

| Quality | # Users | # Venues |
|---|---|---|
| Continent | 107 | 61 |
| Country | 602,932 | 294,596 |
| State | 390,224 | 93,513 |
| County | 251,383 | 276,097 |

So, let us just look at just characterization of this data, which is what kind ofin generally about the dataset that is available. So, we will look at everytable and every figure in this paper.Also here if you look at the number in the dataset is 13 million UHC. UHC stands for user home city; VL stands for venue location, the number of venues that areavailable; and GI stands for geographic information right.

Of course, there is going to be some information, some locations which are not going to be valid right for example, something in the middle of sea, you are not going to get any location, and there are locations that may be generated which is somebody's heart right,h e a r t. So, these kind of locations has to be removed that is what happened between number in the dataset and valid GI; valid, but ambiguous which is it is valid, but we are not able to figure out the exact location, reverse look up, and find out the location that falls into the third row.

Non- geographic information and empty entries, so essentially the dataset was pruned to get data which the researchers can actuallyuseto do the analysis right.This is what even youwoulddoforthehomeworksthatyoudidyoucollected somedata,butyouprobably did not do the way to actually prune the data to get more accurate, more specific data.

(ReferSlideTime:04:00)

| Quality | # Users | # Venues |
|---|---|---|
| # non-GI | 244,233 | 4,214,671 |
| # empty entries | 26,715 | 0 |

Table 2. Quality of Geographic Information.

| Quality | # Users | # Venues |
|---|---|---|
| Continent | 107 | 61 |
| Country | 602,932 | 294,596 |
| State | 390,224 | 93,513 |
| County | 251,383 | 276,097 |
| City | 10,354,058 | 6,937,523 |
| Neighborhood | 981,139 | 1,060,124 |
| Area of Interest/Airport | 27,307 | 47,896 |
| Street | 326,751 | 95,543 |
| Point of Interest | 5,607 | 9,792 |
| Coordinate | 61 | 32 |

Thus, in order to standardize the home city and venue lo-

Also the quality of geographic information is continent, countries, state and this information that is available in this dataset is number of users, number of venues. Country, state, county, city, neighborhood, area of interest or airport, street, point of interest and coordinates. So, all this pieces of information you will get in your json when you collect data from foursquare.And this was basically pruned to get more quality data which can be used for analysis, is that making sense. So, these are called exploratorydata analysis, here we just explaining the data itself describing the data in terms of what is available in the data that was collected.

(ReferSlideTime:05:02)

identify locations at the finer granularity of streets. More-over, note that the use of standardized city name allows us to

⁵http://developer.yahoo.com/geo/placefinder/

Multiple tools were used. So, let me just show you one of them which are developer dot yahoo dot com slash geo slash placefinder. These kind of tools lets you actually reverse look up a place and find out where they are in the map; if you give them location, it can actually give you the latitude, longitude even the other way round, you give the latitude longitude it will give you the location.

(ReferSlideTime:03:35)



Figure 1. Cumulative Distribution of the Number of Mayorships, Tips and Dones per User.

Many of such tools were used in this particular research to find out the location of the checkins,locationoftips,andotherinformationthatwascollected.So,ifyoulookatthe rest of the analysis, so here is the two figures that we will also talk about; first let us talk about the figure 1, which is shown on the left, but let us look at the content.

(ReferSlideTime:06:03)



For the figure 1, as shown in the figure 1 and consistent with previous analysis of Foursquarethedistributionofnumberofmayorships,tipsanddones.So,thisiswhatyou do first when you collect such data. So, one other things that you would have generally seen in social media analysis is to show whether the data is a power law or show overthat the data that you have collected from social media follows the pareto principle. which is to say that 20 percent of the users only actually contribute to the 80 percent of the content that is generated on the social media.

One of the similar type of graph was drawn here which is to see the distribution of mayorships, tips and dones per users and the inferences that they are skewed, with a heavy tail, implying that few users have many mayorships - tips or dones, while vast majority of them have only one mayorship, tip or done. So, that is essentially what a Pareto principle is that is essentially, what power law is also.

(ReferSlideTime:07:30)



Figure 1. Cumulative Distribution of the Number of Mayorships, Tips and Dones per User.

So, if you look at the graph which is on the left for the figure 1, you will see the same thing which is large amount of, so large amount of users actually have small amount of users have so that is what this here. So, let us take figure 90 percent of the data is getting generated by small set of users; majority of the users over here do not contribute to anyof the mayorships, tips or dones, so that is the graph that you would read.

(ReferSlideTime:08:07)



Similarly,let us go to figure 2.Figure 2shows the distribution ofnumber of mayorships, tipsandddonespercity,consideringonlycitieswithatleastoneinstanceoftheattribute.

So, it is actually wanted to see whether from a city, whether you are able to get a lot of mayorship, tips and dones, this can actually help some, help find out how the data is.As shown the distributions are also very skewed, with a few cities having as many as 100 mayorship tips and dones.

(ReferSlideTime:08:47)



So, again if you look at figure 2, the distribution is very similar in terms of the small number of cities having large number of mayorship, tips and dones; and large number of cities, do not have these. So, these are all large, so if you look at it somewhere around 80 or90yearssomethingthatisonlyasmallsetofcitieshere.Largeamountofcitiesdonot have mayorship tips and dones or very little as many little mayorship, tips and dones because the condition was that they were considering only cities with at least one instance of the attribute, which is they should have had one tip, mayorship or done.

(ReferSlideTime:09:39)



So,thenlookingatcorrelationbetweenthenumberofmayorship,tipsanddonespercity, they found that a high correlation between the number of mayorship and the number tips across cities, with the coefficient of 0.78. Similarly, the correlation is also high between number of mayorships and the number of dones, which is if there are more mayorships there in a city that is high chance that there will more of tips and dones also. This is helping us to understand that where if I find cities which I have a high mayorship, I should be able to find, there should be more of tips and dones also there.

(ReferSlideTime:10:34)

So, if you look at the tips, tips are concentrated in different location around the Earth; which is if you remember tips of the content that people post for a particular venue. The top 3 cities in the number tips are New York, Jakarta and Sao Paulo with the total of 600,000 tips. Dones, on the other hand, tend to be concentrated in venues in the US, in cities like first New York, Chicago and San Francisco, and total they have about 1 million dones, so which is to show that some cities, some popular cities have a lot of these tips and dones NewYork being common in both tips and dones. Once again to say that cities generate a lot of these tips mayorships and dones and of course, this wouldalso probably lead in to check ins also.

(ReferSlideTime:11:55)



So, here is another graph which we will see in figure 3 and figure 3 for that we will see. So,figure3showsthesedistributionsinmapsoftheglobewitheachpointerrepresenting a city with venues, with at least one mayorship tip and done. So, essentially until now we only saw per city what is happening? Now when you look at it in a map, the figure 3 actually shows the results. As the maps show, Foursquare venues are spread all over the world including remote places such as Svalbard, an archipelago in the Arctic Ocean.

(ReferSlideTime:12:41)



(a) Mayorships.      (b) Tips.

Figure 3. Global Distribution of Mayo

So,forexample,letusgolookatthefigure3a,bandc.So,thisisthemayorship.Thisis basically showing you every dot in this graph, every blue dot in this graph, figure 3(a) shows you the distribution of the mayorships that are available around the world, thatwere done around the world.

(ReferSlideTime:13:00)



(b) Tips.      (c) D

Figure 3. Global Distribution of Mayorships, Tips and Dones.

Figure 3(b) shows you the tips that were done. So, if you see earlier, we saw that the there is high correlation between mayorship, tips and dones. So, therefore, when mayorshipsarehighthereisgoingtobetipsanddonesalsowhicharehigh.So,youcan

clearly see heavy concentration on many places in the world. And the last one is dones, the figure 3(c) shows you green dot, every green dot is a done from that particular locations. So, this basically shows you mayorships of the blue dot, tips - the red dot, and the green dot being dones. So that is figure 3 is denser with the total number of unique cities.

(ReferSlideTime:13:48)



So, the distribution of mayorships shown in figure 3(a) is denser with a total number of unique cities being 79,000. So, 79,000 cities have higher check ins, have mayorships in the figure 3(a) in the total data. And if you look at the figure 3(b), somewhat sparser tip map for figure 3(b) indicates that there are manycities, particularlyin Canada,Australia, and Central Asia, and Africa, where despite their insistence of venues and mayors' users do not post tips.

So, therefore, there is a chance that there are mayorships in that location, but not tips. Figure 3(c), reveals an even sparser map, with most activity concentrated in touristic or developed areas, such as USA, Western Europe and Southeast Asia. So, essentially even though there is a correlation between mayorships, tips and dones, there is actually some places which are sparser for a tips and dones.

(ReferSlideTime:15:33)



So, now let us look at figure 4; figure 4 here, so this is figure 4. Figure 4 is showing you the cumulative distribution of time interval between consecutive tips and dones posted per user. Why is this interesting? This is interesting to find out about the activity or the frequency of activity of the users.

(ReferSlideTime:16:07)



This is figure 4, figure 4 shows the cumulative distributions of these four measures. We note that the distribution of minimum inter-activity times is very skewed towards short periodsoftimes,withthealmost50percentoftheuserspostingconsecutivetipsand

dones 1 hour apart, got it. So, that is it shows the there is a lot of content that are generated, lot of tips are generated byusers,tips and dones are generated byusers within apart 1 hour. However, an average, median and maximum users do tend to experience very long periods of times between consecutive tips and dones.

So, essentially what this shows is again it is going back to the same power law concept, there are some set of users where there is consecutive tips and dones are done very frequently. There is set of population where this distribution is actually pretty skewed, which is long set of long time taken between two consecutive tips and dones. For instance, around 50 percent of the users have an average interactivitytime of at least 450 hours that is close to about 20 days, whereas around 80 percent of the users have the maximum interactivity of 167 hours - roughly a week.

(ReferSlideTime:17:51)



Figure 4. Cumulative Distribution of Time Interval Between Consecutive Tips/Dones Posted per User.

So, here is the graph for cumulative distribution of time interval between consecutivetips and dones posted per user, it's the same thing as what we saw in the text.

(ReferSlideTime:18:22)



Now, let us look at the next figure, figure 5. So, this part we already saw. Figure 5 is basically looking at the same question which is now we are analyzing the displacement between two venues. The last figure that we saw was looking at two different timing in whichthepostwasdone.Nowwearelookingattwodifferentvenues,whicharedoneby the same user consecutively.

(ReferSlideTime:18:58)



Figure 5. Cumulative Distribution of Displacements Between Consecutive Tips/Dones Posted per User.

So, let us look at the figure 5, first and I will tell you what the figure 5 all means. So, figure5isthecumulativedistributionofdisplacementbetweenconsecutivetipsand

dones posted per user.So, on thex-axis itis showing you thedistance; on they axis, itis showing you the number of the distribution.

(ReferSlideTime:19:27)



Figure 5 shows the distributions of these measures for all analyzed users. Around 36 percent of the users have average and maximum displacements of about 0 kilometers, right, indicating very short distances - within a few meters.

(ReferSlideTime:19:49)



70 percentof theusershavean averagedisplacementofatmost150 kilometers, which is basicallysomebodymovingbetweencities.So,IprobablyfromDelhiIgotoAgra,and

then I do it check in all Mathura, I do check in Mathura, I do a tip or a done, that is what is actually capturing within 150 kilometers.

And about ten percent of the users have a maximum displacement of about 6000 kilometers, this is probably international travel between two consecutive tips or dones, so that shows what is the distribution of the users who we have in the dataset, the consecutive tips and dones that they do on Foursquare. Let us go to the figure again. So, this is basically showing you that 70 percent of the users are about 150 kilometers and the 10 percent is about the 7000 kilometers that is what you will see in this figure.

(ReferSlideTime:21:09)



Figure 5. Cumulative Distribution of Displacements Between Consecutive Tips/Dones Posted per User.

So, if you seen here 7000 kilometers, so it is about the last 10 percent of the users, who are about 7000 kilometers and then very short is about 70 percent. Average 70 percent is that the green one is the average, the green square is the average. The red plus symbol is a median, triangle is the minimum, and the circle is the maximum right, so that gives you a sense of.

(ReferSlideTime:21:51)



So, now,what all we have seen, we have seen the time, consecutive, tips or dones that is done with respect to time, consecutive tips and dones with respect to distance.

(ReferSlideTime:22:18)



Now we will see the next figure, next figure is actually very interesting. Next analysis is actually a very interesting analysis, where they saw how frequently that the check ins,the tips or the dones are coming back for that particular location. So, this is distribution of the returning time. So, if I do tip or a done inIIIT, how frequently do I actually do a tip or a done in that location.

(ReferSlideTime:22:52)



And the figure 6 shows the distribution focusing on returning times under 360 hours, which account for 69.7 percent of all the measured observations. The curve showsclearly daily patterns with returning times often being multiples of 24 hours which isvery similar to the distribution of returning times computed based on the check ins.

So, if you really look at what does it mean why is it 24 hours? If somebody checks into office in the morning today that they have got into the office, shop, institute and everything they do the same check in the next day, so that is what this means right. Checkins,thereferencesisgiventoanotherresearchwherethecheckinswhereseenbut in this case we are also looking at the tips or the dones. That is a very interesting conclusion to know or basically it is complimenting the real world behavior that you could expect from the users.

Here is a graph, which you see here this is about 360 hours.And you can clearly see that it's coming back. So, this is for everyday24, 48, 72, 96 120, so it is kind of coming back every 24 hours and sometimes the frequency is also increasing for something happens. So, this is 168 should be the week. So, therefore, there is a slighting increase from the day that which is which is the 7th day of tips and done. I hope that is making sense essentially the conclusion there is that people come back to that same location with the examples like me doing it in IIIT Delhi, that is figure 6 that is an interesting analysis.

So, now what we will do is, now we will attack the question that we started off with which is to find out how much can be actually inferred about the users' home using this data. So, in this case, we are going to actually use data for most popular location among mayorships, tips and dones of a user or home location using a majority voting scheme. I am going to explain to you at this voting scheme is there are two tables that we will look at and,

(ReferSlideTime:26:01)



Table 3. Home Location Inference.

| | Home City | | | Home State | | | Home | |
|---|---|---|---|---|---|---|---|---|
| | Class 0 | Class 1 | Class 2 | Class 0 | Class 1 | Class 2 | Class 0 | C |
| | 727,179 | 847,876 | 239,129 | 707,953 | 913,166 | 110,110 | 727,179 | 1,0 |
| | 725,073 | 671,576 | 192,781 | 702,583 | 727,219 | 99,672 | 725,073 | 83 |
| | 546,815 | 541,795 | 106,297 | 524,137 | 561,165 | 55,115 | 546,815 | 63 |

We will also see what mechanisms the authors followed in terms of generating this information about what is the possible location that this person's home would be.

Here is the scheme that they followed. They actually put the data into 3 buckets, class 0, class 1 and class 2. The class 0 is of the users who have single activity either the mayorship, the tip or done. They only have one activity in the dataset, whereas class 1 consistsofuserswhohavemultipleactivitieswithpredominantlocationacrossthem.So, for example, I have multiple tips and dones, mayorships in my account, but there is one which is very, very high which is IIIT Delhi for that matter. Class 2 is consists of users with multiple activities in which there is no single location that stands out.

So, again just to understand, if you understand this I think the inferences become much simple, the logic the authors followed is that take all the users who are been doing this tip,donesandmayorshipsinthedata.Class0orthepeoplewhohaveonlysingleactivity either a tip, or a mayorship or a done, class 1 is set of people who were where this one majority location that shows up for them. Class 2 is a set of people where multiple activities are done, butno single location is actually predominant in their activity,so that is the kind of classification that they made with the users.

**Table 3. Home Location Inference.**

**Classes Distribution**

| Features | Home City | | | Home State | | |
|---|---|---|---|---|---|---|
| | Class 0 | Class 1 | Class 2 | Class 0 | Class 1 | Class |
| Mayorship | 727,179 | 847,876 | 239,129 | 707,953 | 913,166 | 110,1 |
| Tip | 725,073 | 671,576 | 192,781 | 702,583 | 727,219 | 99,67 |
| Done | 546,815 | 541,795 | 106,297 | 524,137 | 561,165 | 55,11 |
| Mayorship+Tip | 898,293 | 1,322,214 | 300,831 | 878,578 | 1,398,351 | 146,5 |
| Mayorship+Done | 825,009 | 1,213,917 | 270,974 | 805,029 | 1,278,784 | 130,4 |
| Tip+Done | 831,759 | 1,038,268 | 223,093 | 807,091 | 1,089,638 | 116,5 |
| All | 939,888 | 1,573,471 | 310,045 | 919,938 | 1,643,825 | 153,9 |

**Accuracy**

| Features | Home City | | | Home State | | |
|---|---|---|---|---|---|---|
| | Class 0 | Class 1 | Total | Class 0 | Class 1 | Tota |
| Mayorship | 51.61% | 67.41% | 60.12% | 71.27% | 80.92% | 76.70 |

Now I will show you two tables, one is the number of people, number of the data points from this dataset. If you just consider the class 0, class 1 and class 2, how many people are actually where you can infer home city,home state, and home country. So, how do I read this graph this how do you read this table. This table is this column it showing you features mayorships, tip and done, mayorship plus tip, mayorship plus done, tip plusdone and all of them, right. So, which is if you take onlythe mayorship, what is the class 0,which is onlyifIconsider mayorshipand thereis only1activitybythis user,there are about 727,000 data points in class 0; 847,000 where that are users where one location is actually predominant; and 239,000 there is no location that is predominant, that is how you read this table, correct.

So, if you look at mayorships 127,000; mayorship plus tip 898,000; obviously,mayorship plus tip will be higher, all will be higher for all of them, bigger than all of them and that is for the city.So, for the state 700,000 is for class 0; 900,000 are for class 100,000 is for class 2. Similarly, for the country, so that is giving you a sense of in the data points or the pieces of information that is available for each of these features.

(ReferSlideTime:30:09)



| Features | Home City | | | Home State | | |
|---|---|---|---|---|---|---|
| Done | 546,815 | 541,795 | 106,297 | 524,137 | 561,165 | 55,115 |
| Mayorship+Tip | 898,293 | 1,322,214 | 300,831 | 878,578 | 1,398,351 | 146,526 |
| Mayorship+Done | 825,009 | 1,213,917 | 270,974 | 805,029 | 1,278,784 | 130,439 |
| Tip+Done | 831,759 | 1,038,268 | 223,093 | 807,091 | 1,089,638 | 116,549 |
| All | 939,888 | 1,573,471 | 310,045 | 919,938 | 1,643,825 | 153,955 |

| | Accuracy | | | | | |
|---|---|---|---|---|---|---|
| | Home City | | | Home State | | |
| Features | Class 0 | Class 1 | Total | Class 0 | Class 1 | Total |
| Mayorship | **51.61%** | **67.41%** | 60.12% | **71.27%** | **80.92%** | **76.70%** |
| Tip | 51.52% | 67.29% | 59.11% | 70.29% | 80.59% | 75.53% |
| Done | 50.09% | 61.74% | 55.89% | 70.16% | 78.38% | 74.41% |
| Mayorship+Tip | 51.57% | 66.24% | **60.31%** | 70.21% | 80.27% | 76.39% |
| Mayorship+Done | 51.05% | 65.27% | 59.51% | 70.01% | 79.89% | 76.07% |
| Tip+Done | 51.18% | 64.16% | 58.38% | 69.76% | 79.28% | 75.23% |
| All | 51.46% | 64.86% | 59.85% | 69.74% | 79.53% | 76.02% |

we consider only users whose home city attributes

Now, let us look at accuracy, which is if I were to use this information and find out that mayorship, only using the mayorship, I am able to find the home city 51.6 percentage of the times, wherever the percentage is higher it is been actually kept in bold. So, if you look at mayorship it is class 0 on class 1, we cannot do class 2 because it is actually the places where a particular location cannot be, one single location cannot be inferred, which is because our goal is to infer actually the home location.

So, if I have multiple locations, I am not going to use that column, that is why class 2 does not exist in high accuracy. So, which 67 percentage of the times, home city can be inferred, if I look at class 1 category of people. Which is I do a lot of tips and dones, but my predominant place where I do a tip or a done, tip or a done is actually my home location. Home state becomes higher, seventy percent and home countries even higher,of course, the percentage for the country is going to be, so the percentage of city will always be lesser than state, will always be lesser than country. Because here to get the country that I am from India more difficult to get that I'm from Tamilnadu as a state it is even more difficult to get that IamfromChennai, that is about the inference of the home location.

There's another interesting graph that authors have which is figure 7, which I show you the graph and then I will try to explain it.

Figure 7. Cumulative Distribution of Distances Between Inferred and Declared User Home City.

So, here figure 7 is cumulative distribution of distances between inferred and the declared home city, which is that because people actually declared that what the home location in these services also. If you go to myaccount, you will find that is, if you go to my facebook account I probably say that I am from the current location is from New DelhiandmyhomelocationisfromChennai,sothatinformationyoucanusetomake

the difference which is what did we predict from the table that I showed you know, which is prediction of my home location with class 0 or class 1.

And then use it for finding the difference between what did I say and what I actually have, that is the graph here. So, x-axis is the distance of inferred and declared user home city, and y-axis is the probability. So, this is how you will read the graph.

(ReferSlideTime:33:28)



Which is, distribution of these distances are shown in figure 7. Wefound that 46 percent of the distances are under 50 kilometers that is what the authors did.

(ReferSlideTime:33:40)

They actually zoomed in. So, this is meaning, this is 5000 kilometers, but as if you look at this graph the inside graph is only for 100 kilometers. So, you can clearly see that 46 percent of the distances are 150 kilometers. So, here is 50 kilometers and if you see 46 percent should be somewhere here correct. So, 46 percent of the users are actuallyhaving the error between finding the home location and the actual location is about 50 kilometers, that is pretty small, if I were able to actually use this information with only50 kilometers of error which is I just getting it from the general behavior tips and dones, that's quite effective.

(ReferSlideTime:34:37)



So, if you just take this model, and then if you look at theauthor's claim that 78 percent of the users within 50 kilometers of distance, which is whattheyaresaying is combining these results with the correct inferences produced by our model, we find that we can correctly infer the city of around 78 percentage of the users within 50 kilometers. So, whatever your city is we will be able to make an inference of that city of about 78percent within the 50 kilometers of distance, correct.

(ReferSlideTime:35:27)



So, that gives you a sense of how also if you go if you remember even in the abstract I showed you this 78 percent of accuracy that the authors claimed that you can find the home location. And of course, in the paper structure, you finish off with the conclusions and future work and probably have some limitations if there are any data limitations in data methodology, any limitations in the paper, right.

(ReferSlideTime:35:49)

So, that gives you a sense of how a paper meaning the things that you have seen in the class until now which is to look at take some data do some analysis, make some inferences, how these inferences are put into paper is what we saw in this particular lecture. And the focus was actually taking foursquare and finding the home location. With that, I will stop here for this paper, and I will see you soon.

# Beware of What You Share Inferring Home Location in Social Networks- On the dynamics of username change behavior on Twitter

(ReferSlideTime:00:12)



Welcome back to the course Privacy and Security in Online Social Media. Continuingthe trend that I did last lecture, I am going to continue actually looking at some papers which are basically addressing the problem of privacy leakages from location based services.Ifyou rememberlastlecture,wehadpaper whichlooked atFoursquare,andthe paper analyzed, how they can actually identify, where a person lives. That was onlyusing the Foursquare mayorship, tips and dones. And what we are going to see now is almost the same topic, but we are going to actually compare it with different social networks.

(ReferSlideTime:01:11)



So, ifyou see here inthis paper the authors performa large scale inference studyinthree of the currently most popular social networks like Foursquare, Google plus and Twitter. So, the goal in this paper is verysimilar to the paper that wesaw last time, but it is going to be looking at different social networks not just only Foursquare. So, in this authors lookedatFoursquare,GoogleplusandTwitter.Youknowaspartofthiscourse,youhave already seen all three social networks in terms of their content, in terms of the data collection that is done and information that you can actually collect from the social networks.

The authors actually find that it is possible to infer the user home city with the high accuracyaround 67percent, 72percent and 82percent in the case ofFoursquare, Google plus and Twitter, which is 67 percent for Foursquare, 72 percent for Google plus and 82 percent for Twitter. I am sure as we move along; you will actually understand why Twitter is actually high in terms of finding out the home location.

(ReferSlideTime:02:55)



So, now let us look at the paper in terms of the same structure as we saw last time, introduction, talking about what a location based social networks are.

(ReferSlideTime:03:04)



Talking about collecting data from a different social media services like Google plus Foursquare and Twitter and different research that are done in the context of Foursquare Google plus and Twitter.And I am talking about what information was collected, and a little bit of conclusions of the paper itself, and then talking about how the paper is out maxed.

(ReferSlideTime:03:32)



So, this is generally the structure that we saw even the last time, meaning almost all papers appears see the structure would be the same a paragraph about the 30,000 feet high view of the problem. Then the paragraph about the current problem and what is missing, then the paragraph about what is, what was done in this paper and then some kind ofacontributionfromthis paper.Relatedworkagain Iamnot goingtodetail inthis particular related work.

(ReferSlideTime:04:02)

But related work generally talks about these kind of privacy leakages from locationbased services and work done on collecting data from these three social networks and inferences that were done.

(ReferSlideTime:04:19)



Thenmanyatimesresearchersactuallytendtowritedetailsaboutthesocialnetworkthat isbeing discussed interms ofjustintroducingtheterminologieswhich wesawin thelast paper also. Here it is talking about Foursquare then there would be about Google plus and then Twitter.

(ReferSlideTime:04:42)

So, here if you look at it, the dataset that was used in the study is the same as the last paper.Datasetcrawled betweenOctober 2011,through thesystemanditcomprisesof13 millionusersandaboutcloseto16milliondifferentvenues.Andtheuserhomecityisan optional open text field limited to about 100 characters. For venue, is the location must be defined filling the open text fields, namely city and address limited to 30 and 127 characters respectively.That is the kind of data and that is the kind of information that is available when you collect these data for venue, tips and dones.

(ReferSlideTime:05:42)



So, the entire dataset is about 15 million mayorships, close to 11 million tips, and closeto ten million likes. All right? So, likes is basically dones in terms of Foursquare terminology.

(ReferSlideTime:06:05)



So, now in terms of Google plus data, Google plus is basically a network that is very similar to, I mean, if you have a Gmail account, you essentially have a Google plus account. In total 27 million profile pages through HTTP request were crawled, and 7 million defined at least one place where they lived, and 5000 provided address information and about 7 million filled their education and about close to 6 million filled their employment.

So, these are details, meaning, if you remember, if you just recollect the social network that you use more often, which is like Facebook, you have all these details at the right places that you live, education, colleges that you study, places that you worked, places that you have lived, all of these information are taken from the users and that is what is mentioned here. Which is 7 million people have explicitly stated their education and about 6 million people have explicitly stated their employment details which is I work at IIIT Delhi.

(ReferSlideTime:07:34)



Now,letuslookatthedatafromTwitter.So,thedatathedatafromTwitterwascollected using Streaming API, which all of you are aware of. And the crawl was done for 120 milliontweetspostedbyaboutcloseto20millionuniqueusersfromApriltoJune2012. 0.5percentofthepostsposted byuniqueusersweregeographically tagged.

So, what does this mean, this means that there are only 0.5 percent of the total posts that were collected where there is geo tagged information for the post which is geo tagged information for the users also. There were about 700,000 tweets and about 300,000 unique users. That is the exact location, which we have discussed in the past, which is latitude, longitude of the post from where the post is coming.

(ReferSlideTime:08:47)



So, that is the background about the dataset. Essentially all of them are talking about in millions in Twitter it's about 0.5 million geographically tagged tweets geo tag tweets.

(ReferSlideTime:09:00)



In Google plus that are about 27 million profiles that were crawled and about 7 millioneducation and 6 million employment.

(ReferSlideTime:09:07)



And Foursquare has details of about 16 million mayorship; 11 million tips and close to 10 million likes. That is the dataset we are going to play around with to do the analysis, to find inferences about the home location of the person.

(ReferSlideTime:09:38)



So, in any data set, when you analyze, first you know you want to actually provide exploratory data analysis, and you want to provide what the data set looks like, because this will help reproducibility of that research. This will help others to actually collect data, if they were to, the point here is that if others want to collect the data which is very

similar to what you collected; and if others want to do the same analysis that you did the results should be the same. That is the idea for reproducibility of the research. So, explaining how you did collected the data, explaining what the data looks like is extremely important in terms of actually writing these research papers.

(ReferSlideTime:10:45)



TABLE I
AVAILABILITY OF GEOGRAPHIC INFORMATION (GI) IN VARIOUS ATTRIBUTES IN OUR DATASETS. ALL VALUES PRESENTED ARE IN PERCENTAGE (%).

| | Foursquare | | Google+ | | | | Twitter | |
|---|---|---|---|---|---|---|---|---|
| Statistics | User Home City | Venue City | Places Lived | Address | Education | Employment | User Location | Geo-tagged Tweet |
| % valid UGI | 95.35 | 55.45 | 61.85 | 0.01 | 52.95 | 34.52 | 73.28 | 100.00 |
| % valid AGI | 2.65 | 18.04 | 6.66 | 0.002 | 11.01 | 14.67 | 9.70 | 0.00 |
| % non-GI | 1.80 | 26.51 | 31.48 | 0.01 | 36.04 | 50.81 | 11.90 | 0.00 |
| % empty | 0.20 | 0.00 | 0.00 | 99.98 | 0.00 | 0.00 | 5.12 | 0.00 |

indistinctly, being *Yahoo!* unable to decide which is correct – e.g., "Springfield" is the name of ten different cities, only in The United States. For Foursquare, due to space constraints in the paper, we group tips, likes and mayorships as venue attributes, while users attributes correspond only to the home city field. Note that the vast majority of Foursquare users (98% of 13,570,060) provided valid home city locations, with only a tiny fraction leaving it blank (0.2%) or filling it with non-geographic information (1.8%). Moreover, 11.6 million venues have valid locations associated, although a substantial fraction of all venues have non-valid locations (around 26%) or valid but ambiguous location (18%). This large fraction of non-valid or ambiguous venue locations comes as a surprise, particularly considering that, unlike the user home city field, the venue location information is a mandatory attribute.

In comparison with Foursquare, the fraction of valid locations in our Google+ dataset is much lower for all considered education, employment and address attributes are more often provided at finer granularities, i.e., street level for employment and address, and Point Of Interest (POI) for education. Finally, the "quality" of the location provided in users' tweets is either at the street (18.05%) or at the geographic coordinate (81.95%) levels. The availability of public finer-grained location information opens an opportunity for more specific inferences regarding user home location, such as user residence, as discussed in Section V-C.

*B. Attribute Characterization*

In the previous section, we analyzed the availability of valid and unambiguous geographic information as well as the "quality" of this information across all analyzed attributes. Now, we focus on the usage of these attributes and analyze their distributions across users in each dataset. We aim at assessing the potential of exploiting these attributes for inference

So, table one provides the distribution of geographic information of all considered attributes in each dataset in each dataset. Wepresent for each attribute the percentage of it that corresponds to the valid geographic. Let us look at the tables. So, this table is the one that is referred. This table talks about availability of geographic information in various attributes in the datasets. So, if you look at the second column, which is Foursquare. So, the columns are referred three different networks Foursquare, Google plus and Twitter.And if you look at the statistics which are in the rows, it is valid UGI which is user geographic information, valid AGI, valid geographic information and that is empty.

So, this basically would help you to find out, what is the amount of data that is available which is valid geographic information, valid and ambiguous geographic information, valid ambiguous geographic information and valid non geographic information and empty. What is this all mean I will I will try to explain this. Valid and ambiguous it is actually latitude or longitudinal, it is actually New Delhi; there is no ambiguity in it. Validambiguityitisnotclear,soitsaysnearTajMahalornearGovindpurimetro
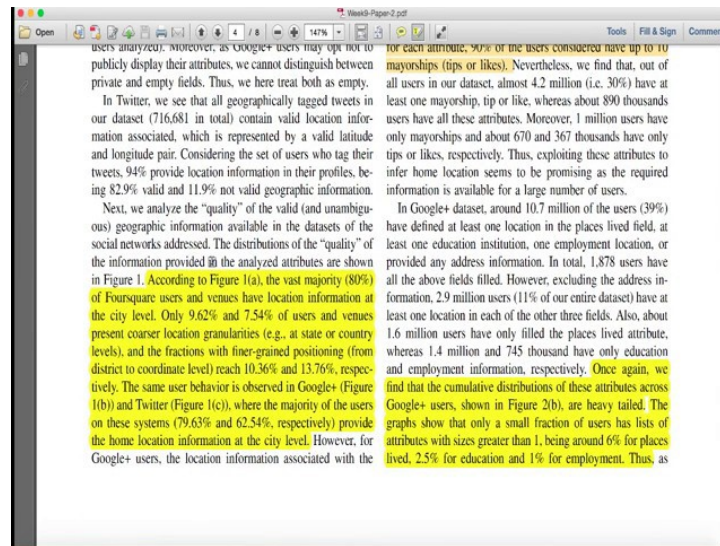
station, so these things are ambiguous.And non-GI – non-geographic information which could be I think as I said before, it could be somebody's heart, h e a r t.

And information like that is actually it is not geographic information at all, and sometimes it could be actually empty. So, essentially that is what is been given in the values. They are all percentages which says user home city is about 95 percent, ambiguous is about 2.6 percent, and non-GI is 1.8 percent, and empty is about 0.2 percent.

So, in Foursquare, it is user home city and venue city. In Google plus, it is places lived address and education and employment. In Twitter, it is the user location geo tagged tweet right. So, this basically tells you different types of information are collected from different networks; I mean that is a whole body of research in terms of actually using these different sets of information from different social networks.
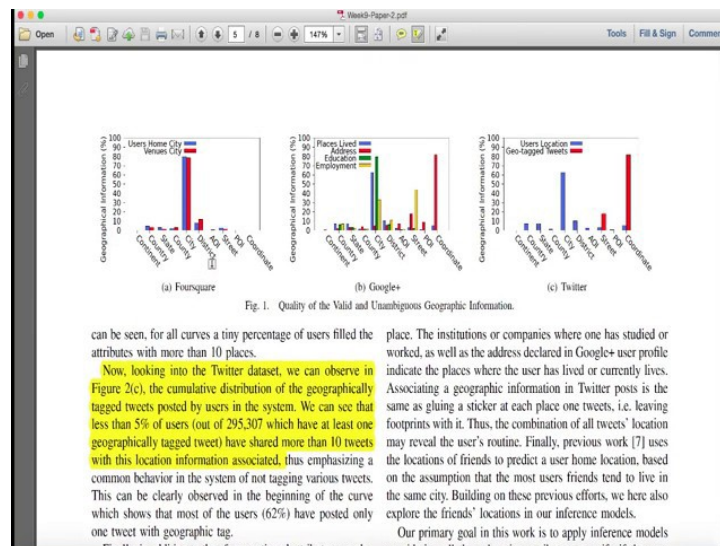
Then Foursquare it is user home city venue and venue. In Google plus, it is places lived address education and employment; in Twitters, it is user location and geo tagged tweet. So, if you look at the unambiguous geographic information for geo tagged tweet from Twitter it is about 100 percent. It is because all the tweets that where collected where actually at the 0.5 percent had geo tagged information in it, so that basically gives you a sense of what kind of data is collected. In terms of Google plus, 53 percent has an unambiguous education, so I studied that Carnegie Mellon University, so that is very very precise, there is unambiguity, there is no problem and actually recoding it or decoding it to a specific university.

(ReferSlideTime:14:57)



So, let us look at different figures, different analysis that is been done using this data. Figure 1 the vast majority 80 percent of Foursquare users and venues have location information at the city level. 9.6 percent and 7.4 percent of users and venues present coarser location granularities at a state or the country levels.
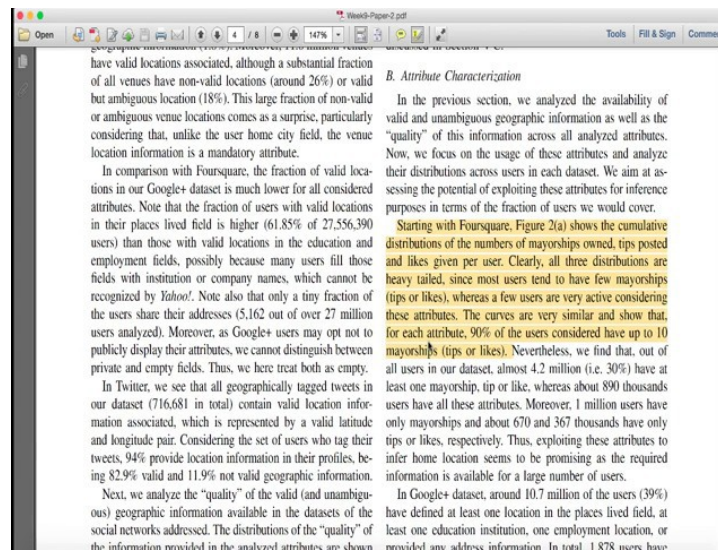
(ReferSlideTime:15:27)



So, we look at the figures. So, this is figure 1(a), 1(b) and 1(c). So, if you look at here quality of valid and unambiguous geographic information. Foursquare, if you look at city,citygives youtheusers'homecityandvenuecity;itisaboutclosetoeightypercent.

So, that is what is written here this says about vast majority 80 percent of Foursquare users and venues have location information at the city level. Some have at the country level, some have at the state level, some have at the street level, so that is the different level of details that the geographic information is available for the location from Foursquare.

So, if you look at Google plus, the information is maximum available for example, it is education that is available at a city level about 70 percent or 70 plus percentage, so thatis what is its written here.The same user behavior is observed in Google plus figure 1(b) and figure 1(c) where the majority of the users of the system 79.63, 62.54 respectively provide thehome location atthe citylevel. Cityis highestinterms ofplaces lived, cityis highest in terms of user location also, so that basically says that we should be able to actually get the city level of information without any problem, because large amount of data for the information about the users is available at the city level.
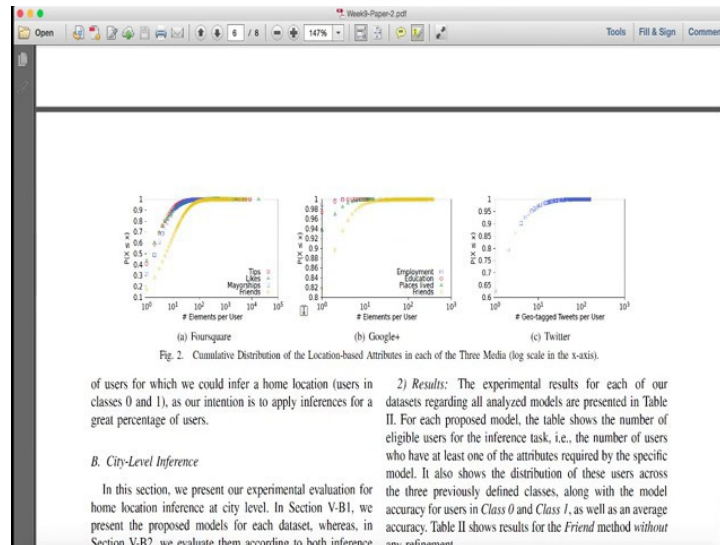
(ReferSlideTime:17:36)



Now, let us look at more analysis with this data. So, now, figure 2(a) shows the cumulative distributions of the numbers of mayorships owned, tips, posted and likes. If you remember even in the last paper, we saw this kind of graphs, which is to show the cumulative distribution of the number of the mayors, tips and dones right. So, if you rememberthegraphthereweresmallsetofpeoplewhohadalotofmayorships,anda

large set of people who had less number of mayorships, so that is the kind of general social media behavior also.

(ReferSlideTime:18:40)

Fig. 2. Cumulative Distribution of the Location-based Attributes in each of the Three Media (log scale in the x-axis).
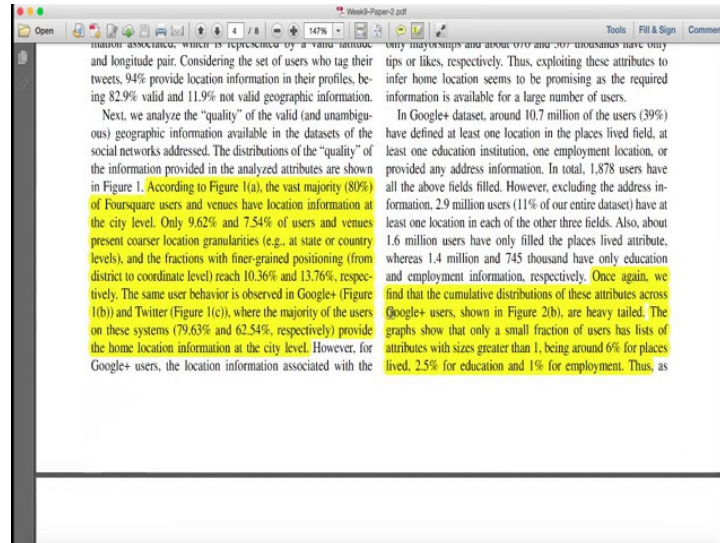
So, let us look at figure 2(a), which is giving you the distribution of the number of mayorships, tips and dones. So, if you see here, the figure 2(a) is giving you the cumulative distribution of location based attributes in each of the three media; a is for Foursquare, b is for Googleplus and c is for Twitter. Given that Twitter has all of them as geo tagged you can clearly seen there is only one line, whereas in the other one there is, friends, mayorships, likes and tips that is in the Foursquare; in Google plus – employment, education, places lived and friends right. So, this is the graph and you can clearly see the graph is very similar to what we have seen in the past in terms of social media data.

So, clearly all three distributions are heavy-tailed (Refer Time: 19:31) which is what I just now said which is social media looking data, since most users tend to have few mayorships whereas a few users have very active considering these attributes. The curves are very similar and shows that each attributes 90 percent of the users considered have up to 10 mayorships right. So, it is the same principle Pareto principle that we talked about a power law that we talked about in the course all of that is playing into this data also. And this is very, very important to show because the reviewers and the readers can actually
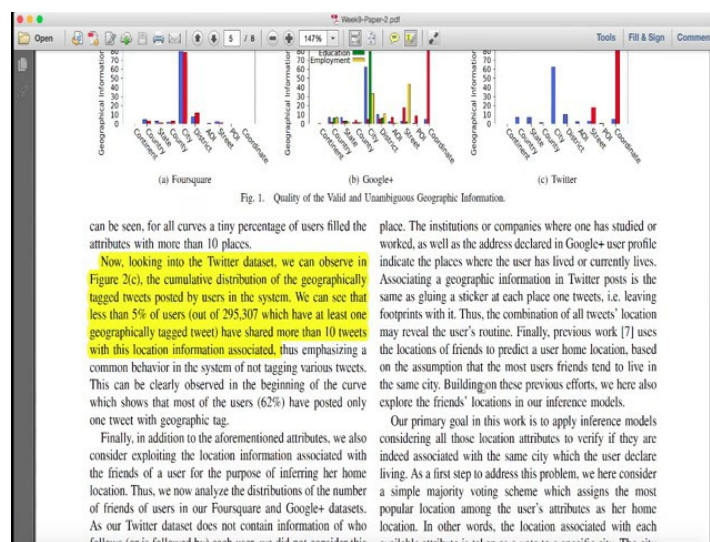
believe that this data is actually representative of other social media research that has been done and analysis that has been done.

(ReferSlideTime:20:25)



Also if you see 2(b), again it's describing all the different social networks; figure 2(b) is showing you for Google plus the graph shows that only a small fraction of users has list of attributes with sizes greater than one being around 6 percent of places lived 2.5 percentage of education and 1 percent of employment.
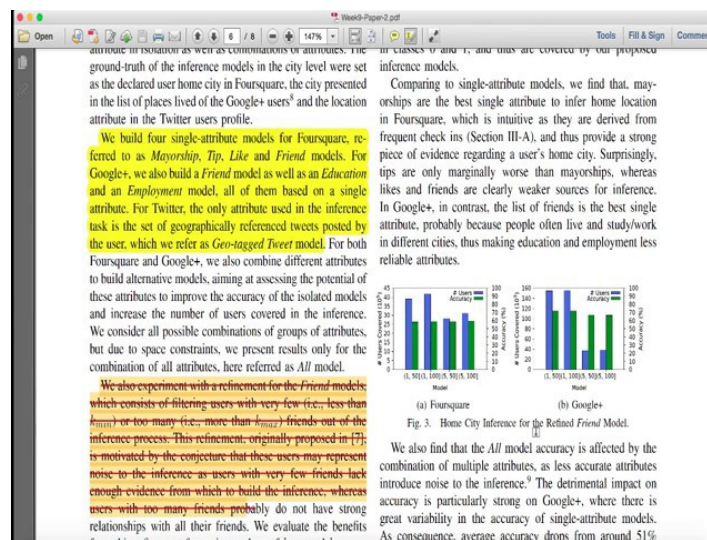
(ReferSlideTime:20:57)

If you look at the Twitter data, we can observe that figure 2(c), the cumulative distribution of geographically tagged tweets posted by users in the system. We can see that less than 5 percent of the users have shared more than 10 tweets with this location information associated, which is again small percentage of people doing location information sharing, more than 10 tweets with the location associated with them.

Also now let us get into inferring location. The methodology that this paper uses is very same to the last methodology, which is written in this paragraph. We group users into three classes, class 0 consists of users who have only one vote that is only one location information that is predominant, and that is only one. Thus, allowing only a unique option tobeassigned fortheuser'shomecity.Class1containstheuserwhohasmultiple votes with the predominant location across them.
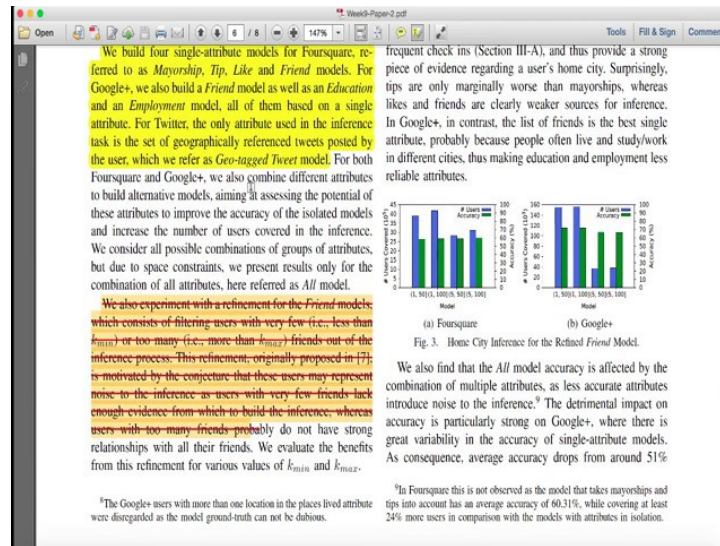
And the class 2 as we have seen in the last paper also consists of users with multiple votes in which there is no single location that stands out. So, three categories of classes three classes that they are made 0, 1 and 2. Wewill see the table with 0, 1 and 2 that lets this locatethehowthedatawas collected, howtheanalysis was done,how thebucketing was done. The results of our experimental evaluation are assessed using two metrics which measure the effectiveness of the proposed model. Accuracy is the fraction of correct inferences of users of class 0, or class 1, right yet again there are (Refer Time: 22:58) the current thing with class 2 will not work.

 (ReferSlideTime:23:04)

So, the model that was built was four single-attribute models for Foursquare, referred to as mayorship, tip, like and friends. For Google plus, friend model and education and employment model all of them based on single-attribute. For Twitter, the only attribute used in the inferences task is the geo tagged location right. So, basically this explains what details were used in collect and making the inference about the location.

(ReferSlideTime:23:44)



Also this is also information we also experiment with the refinement of the friends' model which consists of filtering users. So, another wayjust think about the another way of looking at the friends model is the use the friends from that location to make the decision, so that is filtering with a very few, less than k kilometers, or too many that is more than k max friends out of the inference process.

The refinement originally proposed is motivated by the conjecture that these users may representnoisetotheinferenceasuserswithfewfriendslackenoughevidenceforwhich to build the inference, whereas users with too many friends probably do not have strong relationships with all their friends. It is basically saying that we build a model where we take the users with very few, less than few kilometers; because they are not going to be connected.Alot of friends who may be connected from that location also will not lead a lot of information.

(ReferSlideTime:25:03)



So, this table is the most important analysis orinference from this paper which is to see the summary of results obtained for inference models for a home and home city inference. Remember, we did for three networks - Foursquare, Google plus and Twitter, the inference models that was used for mayorships, tip, like, friend all, education, employment, friend all geo tagged tweets.

Classes distributed 0, 1 and 2; classes 0 and 1 are the only two things that can be done withthisdata(ReferTime:25:45),soitisdone51.So,thewaytoreaditisthatjustusing mayorship, in the class 0, 51.61 percent you can identifythe home cityfor that particular user in the category who has Foursquare account and mayorship data. Class 1, 67 percent; class 1 is basically there are multiple locations, one being predominant. Google plus, the highest seems to be with friend, no refinement; and in tweets, it is since the geo located the accuracy is also being more than anything else.

(ReferSlideTime:26:43)



So, you can that is what I have said the accuracy for the Twitter seems to be higher than the restofit. Let us look atfigure 3,andthen wewillgoback tothe description ofit. So, if you look at figure 3, figure 3 shows the total accuracy which considers the inferences for users in class 0, and 1.And the number of users covered bythe refined friend model, for the various values of k min and k max specified in the x-axis of the graph. Often comparing the results with those in table 2, we see that the refinement improves model accuracy particularly for Google plus where the gain is about 21 percent.

(ReferSlideTime:27:30)

So, essentially this is x-axis is the model that was used within the kilometers and y-axisis the percentage of users that were covered. So, if you see Google plus there is if you will infer the home city for a refined friend model, which is what we said where the min and the max were removed Google plus seems to be doing much better than the added advantage of removing these friends is higher for Google plus compare to Foursquare. That is the inference there.
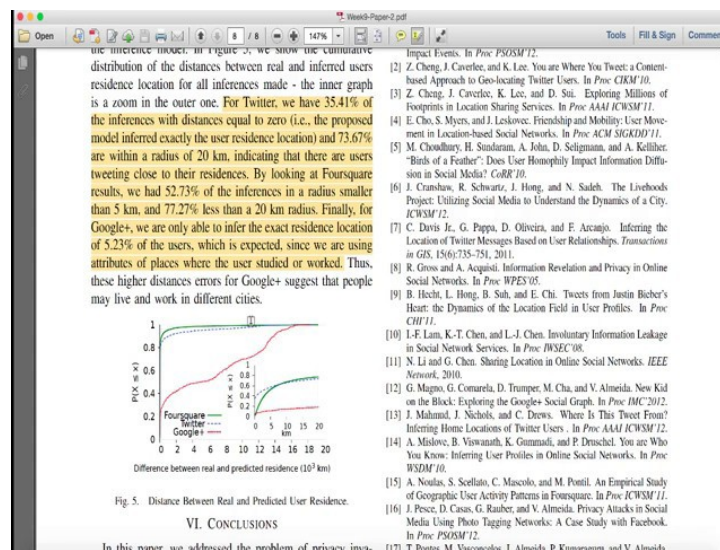
(ReferSlideTime:28:19)



See there is another interesting inference, similar graph we saw on the last paper also. TherewesawonlyforFoursquare;herewe areseeingitforFoursquare,Googleplusand Twitter. Figure 4, corresponds only to the incorrect inferences and the inner graph is basically zoomed into the outer graph. It shows that 46 percent of the distances in Foursquare, and also 27 percent in Google plus, and Twitter are under 50 kilometers. So, that is what is this here 50 kilometers is this part of the inside graph. 50 kilometers is reasonable distance between neighboring cities.

Thus combining these results with the correct inference produced we can make correct inferences in a radius of 50 kilometers with the accuracies that achieves 78.5 percent in Foursquare, 64 in Google plus and 87 in Twitter,which are the things that we saw in the table earlier. This is the representation in the graph which is x-axis is the distance of inferred and declared home city which is something that I declared. And something we wereabletomyaccountPKponguruaccounthasalocationandIinferredthroughthe

process the location what is the distance between these two, the lower the difference the better.

The inside graph is just showing you assumed immersion which shows that about 50 kilometers we were able to get about 46 percent, where if you see here, 50 percent and if this is about 46 percent. So, they just shows you that we are able to actually identify 46 percentofthedistanceinFoursquareisactuallylessthanerrorof50kilometers,which isjustneighboringcities,neighboringplaces, orsometimesitcouldbejustinthesame city.
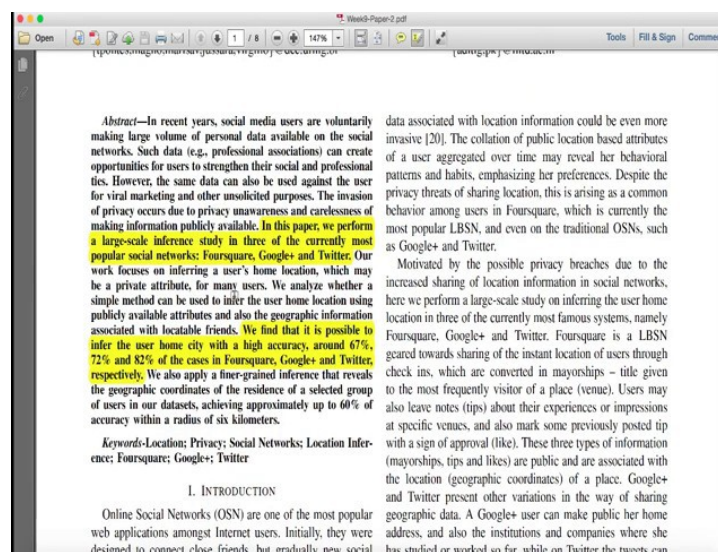
(ReferSlideTime:30:46)



So, now the same thing you can actually do it for the residence right. Now what we didin this graph is basically showing you only the home city, whereas this graph is actually showing you the home residence. So, here is the graph for residence; red is Google plus; blueisTwitter;andgreenisFoursquare.Youcanseetheinsidethegraphalsoherewhich is from 0 to 20 kilometers, whereas this is 0 to 20, but they are all ten to the power of 1000 kilometers is the distance here.

So, you can see that for Twitter we have 35 percent, where is Twitter, so Twitter is blue line,blue line is here,35percentofthe inferences withdistanceequalto0.So,thatis the starting point here if you see, that is the proposed model inferred exactly the user residence. And the reason why this is so high and this is so accurate is because, we are, tweetswerecollectedwhichwereactuallygeotaggedright.And73.67percentarewithin the 20 kilometers radius.

So, if we see here this is 20 kilometers and if we see the blue line it is here that is about 76 percent, 73 percent, which is within the 20 kilometer difference, we were able to find out where the home is which is pretty good. Indicating that that there are users tweeting close to their residences, because I could be living in Okhla, I could be tweeting from somewhere near Jawaharlal Nehru Stadium in Delhi which is less than 20 kilometers, I went to watch a match and I actuallyposted tweet which is also geo tagged,sothat is the kind of 20 kilometers that we can get.

By looking at Foursquare results, we find that Foursquare is green. We find that 52 percent of the inferences in the radius smaller than 5 kilometers. So, if you see heregreen one, if you go at that point it is about 52 percent, 52.73 percent less than 5 kilometers; 77 percent less than 20 kilometers, that is here, 77 percent. Finally, for Google plus, we are only able to infer the exact residence of 5.23 percent of the users, which is expected since we are using attributes of places where the user studied or worked, because here we are only using their employment and education details right.

(ReferSlideTime:33:55)



So, that is how this paper ends, which is to show that let us go to the abstract again, which is to show that they used they used data from Foursquare, Google plus and Twitter. They used this data to infer the home location. This is an extension or the next step for the last paper that we saw which was done only on Foursquare. And they were abletoactuallyshowthatabout67,72and82percentwiththataccuracytheywereable

tofindoutthehomecity,andhomelocation,for,withahighaccuracyintermsofTwitter and then Foursquare, but with less data in a Google plus.

Withthat,Iwillstopthisparticularpaper.Iwillseeyou soon.