Privacy in the Modern Era: The Cases of Online Social Network and Machine Learning Model

Yang Zhang







Michael Backes



Pascal Berrang



Cheng-Te Li



Jun Pang





Mario Fritz



Mathias Humbert



Tahleen Rahman



Ahmed Salem



About Me

- Postdoc at CISPA Helmholtz Center i.G. working with Michael Backes
- From January 2019, research group leader at CISPA Helmholtz Center for Information Security
- Data privacy
 - Biomedical data, machine learning models, social network
- Ph.D. positions and summer interns available









CISPA-Stanford Program

- <u>https://www.cispa-stanford.org/</u>
- Elite scholar program for doctors
 - 1 or 2 years at CISPA
 - 2 years at Stanford University as a visiting professor
 - 3 years at CISPA as a research group leader
 - Drop me an email if you are interested





The Advancement of ICT













Privacy!!!











Outline

- Social network privacy
- Machine learning privacy





Outline

- Social network privacy
- Machine learning privacy





Social Network Privacy

- In basic setting, users
 - Articulate their personal attributes and their social relations -> many attacks exist
- De-facto way for communication
 - Texts, images... -> some attacks exist
 - Cooler information
 - Location check-in, hashtags... -> privacy?







Social Network Privacy

- Location check-in to infer social relation
 - walk2friends: Inferring Social Links from Mobility Profiles (CCS 2017)
- Hashtag to infer location
 - Tagvisor: A Privacy Advisor for Sharing Hashtags (WWW 2018)







Social Network Privacy

- Location check-in to infer social relation
 - walk2friends: Inferring Social Links from Mobility Profiles (CCS 2017)
- Hashtag to infer location
 - Tagvisor: A Privacy Advisor for Sharing Hashtags (WWW 2018)







Location Check-in











Location Check-in

yelp&

hot+new nearby



Novela	0.08 mi
* 🗙 🗙 📩 🔛 41 i	review \$\$
662 Mission St, Fi Cocktail Bars	nancial Distric



0.28 mi Sushirrito 🖈 🖈 🖈 🔝 57 review 226 Kearny St, Financial District Japanese, Sushi Bars,...



MKT Restaurant... 0.14 mi 🔀 🔂 🚼 🔝 7 reviews 🛛 \$\$\$\$ Four Seasons Hotel 757 Market American (New)



Via Moto 0.18 mi 😭 🔂 🚼 🚼 🔀 7 reviews Metreon 135 4th St, Financial D Pizza, Italian, Sandwiches



7:25

а

Si

th









Location Privacy

- 4 spatial-temporal points can identify 95% of the individuals
- Mobility traces can be effectively de-anonymized
- You are where you go
 - Demographics
 - Social relations





Location Privacy

- 4 spatial-temporal points can identify 95% of the individuals
- Mobility traces can be effectively de-anonymized
- You are where you go
 - Demographics
 - Social relations





- Social relations can be sensitive, e.g., office romance
- 17.2% -> 56.2% (Facebook users in New York)
- NSA's co-traveler program



Research Question

Can two users' check-ins be used to predict their social relations?







Existing Approach

- Solution 1: common locations two users have visited
- Almost all data mining approaches take this way
- Location entropy
- Can't work when two users share no common locations









Existing Approach

- Solution 2: mobility profiles/features
- Summarize each user's mobility profiles
 - Friends share similar mobility profiles than strangers
- Feature engineering
 - Tedious efforts and domain expert knowledge
 - Time consuming





Every Single Time!!!



Representation Learning

- Learning features (representation/deep learning)
 - Follow a general object (unsupervised)
- Graph representation learning (graph embedding)
 - Preserve each user's neighbors in a social network
- Mobility feature learning





walk2friends

- Assumption: A user's mobility neighbors can reflect her mobility profile/features
- Define each user's mobility neighbors 1.
- 2. Learn mobility features/profiles
- Infer two users' social relation 3.





Mobility Neighbor

- A user's mobility neighbors include
 - Locations a user has visited
 - Others who have visited similar locations and their locations
- Breadth first search
 - Not considering the visiting frequencies
- Random walk sampling





Mobility Neighbor











Mobility Feature Learning



 $\underset{\theta}{\operatorname{arg\,max}} p(\widehat{\boldsymbol{\boldsymbol{\mu}}} | \boldsymbol{\boldsymbol{\lambda}}; \theta) \cdot p(\boldsymbol{\boldsymbol{\boldsymbol{\omega}}} | \boldsymbol{\boldsymbol{\lambda}}; \theta)$ $\boldsymbol{\mathbb{Z}};\boldsymbol{\theta})\cdot p(\boldsymbol{\mathbb{M}}|\boldsymbol{\mathbb{Z}};\boldsymbol{\theta})$ $\mathbb{Z}: \theta) \cdot p($ $p(\mathbf{m}|\mathbf{Q};\theta) \cdot p(\mathbf{Q};\theta)$ $p(\boldsymbol{\rho} | \boldsymbol{\rho}; \theta) \cdot p(\boldsymbol{\rho} | \boldsymbol{\rho}; \theta)$ $p(\mathbf{m}|\mathbf{w};\theta) \cdot p(\mathbf{w}|\mathbf{w};\theta)$



- Learn a function:
- Each node to predict it's neighbors
- $p(\cdot | \cdot; \theta)$ Softmax

$$) \cdot p(\mathbf{a} | \mathbf{a}; \theta) \cdot p(\mathbf{a}; \theta) \cdot p(\mathbf{a$$





Social Relation Inference

s(2, 2) = 0.9s(2, 2) = 0.8s(2,2) = 0.6s(2, 2) = 0.4s(2,2) = 0.3s(2,2) = 0.2



- Cosine similarity • Unsupervised
 - Predict any social relation





Dataset

- Instagram users' check-ins
 - New York, Los Angeles and London
 - Foursquare (location semantics)
- Social relations (two users follow each other)
- Dataset available!







	New York	Los Angeles	Lon
No. check-ins	1,843,187	1,301,991	500,
No. locations	25,868	22,260	10,0
No. users	44,371	30,679	13,
No. social links	193,995	129,004	25,4















































Defense

- Hiding
 - Delete certain proportion of check-ins
- Replacement
 - Random walk to replace locations







Defense

- Generalization
 - Geo-coordinate and location semantics
 - MoMA -> art (40.76N, -73.97W)
- Recover location first
 - art (40.76N, -73.97W) -> MoMA or Tom Otterness Frog?





Utility Metric

- Each user's check-in distribution
 - Both original and obfuscated
- Jensen-Shannon divergence
- Average over all users







Defense











Defense

	AUC		Utility		Recovery rate	
	ls	hs	ls	hs	ls	hs
lg	0.77	0.75	0.57	0.30	52%	23%
hg	0.73	0.67	0.20	0.06	14%	2%






Defense









- A novel social relation inference attack with mobility profiles
 - Unsupervised and predict any social relations
 - Outperforms baseline models
- Three general defense mechanisms
 - Replacement and hiding outperform generalization





Social Network Privacy

- Location check-in to infer social relation
 - walk2friends: Inferring Social Links from Mobility Profiles (CCS 2017)
- Hashtag to infer location
 - Tagvisor: A Privacy Advisor for Sharing Hashtags (WWW 2018)





















Chris Messina[™] @chrismessina

#barcamp [msg]?

12:25 PM - 23 Aug 2007

146 RETWEETS 288 FAVORITES





how do you feel about using # (pound) for groups. As in





#ShareaCoke









#ImWithHer



al Can





#like4like #foodporn #tbt







#privacy

#locationprivacy







Research Question

Can hashtags a user shares be used to infer her location?









Tagvisor

- Attack: location inference with hashtags
- Defense: Tagvisor, a privacy advisor to mitigate the privacy threat by hashtags





#dataset

- Collected through Instagram's APIs
- New York, Los Angeles, and London
- Hashtags + locations (check-ins)

	New York	Los Angeles	London
No. of posts	144,263	61,767	34,018
No. of hashtags	8,552	4,600	2,395
No. of users	3,911	1,625	992
No. of locations	498	268	141







17 likes

#sunday #sun #reading #rva #tan #light #relax #girl #me #outside #spring #warm #instagood #photooftheday #iphonesia #instamood #igers #instagramhub #picoftheday #instadaily #bestoftheday #igdaily #instagramers #webstagram #all_shots #statigram #popular #photography #art #iphoneography





#attack



- Bag-of-words for feature representation
- Random forest classifier
- Multiple-class classification, e.g., 498 classes (locations) in New York
- All posts are trained together







#attack

	New York		Los Angeles		London		All cities	
	attack	baseline	attack	baseline	attack	baseline	attack	baseline
Correctness	0.613	0.015	0.685	0.015	0.686	0.020	0.624	0.010
Distance (km)	0.917	4.198	1.870	11.275	0.857	4.518	211.471	3563.082
Accuracy	0.697	0.053	0.758	0.048	0.761	0.051	0.712	0.045







#attack







#tagvisor

- A privacy advisor for sharing hashtags
- Fool the attacker's location inferencer (ML classifier)
- Three defense mechanisms
 - Hiding
 - Replacement
 - Generalization (location category)
- Utility: preserving the semantical meaning of hashtags







#hiding

successful attack



delete one hashtag (can be more)

hide #a



hide #b

hide #c









#utility



Hashtag vectors d2 d1 #a: [3.1, 1.3] #b: [2.5, 1.9] #c: [4.0, 5.1]



- Semantical meaning
- Skip-gram, aka word2vec
- Skip-gram over all posts' hashtags







#replacement

successful attack



- Replace each hashtag with all the possible hashtag
 - Search space is very large
- Bound to the most closest hashtags (with word2vec)
 - Reduce the search space
 - Semantical meaning can be preserved









#generalization

- Location category from foursquare
 - #centralpark -> #park
- Do not apply to all hashtags
 - e.g., #tbt #love







#tagvisor

- Check whether the post's location is inferred correctly
 - If no, then publish
 - Else, consider the three defense mechanisms
 - Pick the hashtag set with the highest utility









#tagvisor

Obfuscating bounded number of hashtags



Obfuscating 2 hashtags is enough!









Summary on Tagvisor

- First location inference attack with hashtags
 - Sharing hashtags is not safe!!!
- A privacy advisor to mitigate this risk
 - Minimal risk and maximal utility
 - Fit for the real-world setting







Outline

- Social network privacy
- Machine learning privacy





Machine Learning

















ML Security and Privacy

- Many ML models are used in critical infrastructures
- ML models are trained on sensitive data
 - Biomedical data, emojis (Apple's differential privacy)
- Largely overlooked







Research Question

Does an ML model trained on privacy-sensitive data leak information of the data?

ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models (to appear in NDSS 2019)









Membership Inference

- Determine whether a data point is inside something
 - Biomedical data, case and reference group
 - Location data, NDSS 18'
- Machine learning models
 - Oakland 17'
- One of the most "popular" attacks in the community









ML Routine

Get some data







Train the model





Membership Inference

Get some data















Membership Inference

- Why membership matters?
- A cliché example: a ML model for medical diagnosis, if a person is in the training set, then she has the corresponding disease
- Security implications, IP implications







Threat Model













Attack by Shokri et al.









Our Attack 1

- One shadow model
- One attack model
- Same data distribution







Attack 1







Our Attack 1








Our Attack 2

Can we do better?

- No assumption on the dataset
- Data transferring attack
- Train shadow model on a different dataset, and attack on the target model













74









Our Attack 2

XUIX	0.50	0.25	0.75	0.87	0.25		0.78	0.25			0.77	0.82		1.0
AC AC	0.50	0.87	0.90	0.85	0.65	0.74	0.92	0.77	0.79	0.80	0.78	0.82		
CIFAT 00	0.50	0.83	0.95	0.87	0.75	0.75	0.89	0.77	0.78	0.79	0.83	0.87		0.8
CIFARCiace	0.50	0.83	0.95	0.88	0.79	0.75	0.88	0.77	0.78	0.79	0.82	0.87		
xion	0.50	0.81	0.92	0.83	0.88	0.75	0.85	0.76	0.77	0.78	0.80	0.83		0.6
LOCAL	0.50	0.86	0.72	0.55	0.68	0.65	0.92	0.54	0.51	0.54	0.84	0.67		
Mariens	0.50	0.84	0.95	0.87	0.77	0.75	0.88	0.77	0.78	0.79	0.83	0.88		
Reit	0.50	0.87	0.88	0.80	0.65	0.71	0.90	0.73	0.77	0.60	0.73	0.73		0.4
Purcha:10	0.50	0.87	0.84	0.77	0.66	0.73	0.93	0.71	0.77	0.75	0.78	0.86		
ourchase 20	0.50	0.87	0.89	0.84	0.66	0.74	0.92	0.76	0.79	0.80	0.82	0.83		0.2
ourchase 50	0.50	0.86	0.93	0.87	0.67	0.75	0.92	0.77	0.79	0.81	0.85	0.86		
ourchase	0.50	0.85	0.95	0.88	0.69	0.75	0.91	0.77	0.79	0.80	0.84	0.89		0.0
aurchase	adult	R. IO	2-100	Face	ation	VIS7	Vews	3Ser	se-In	se-2n	Se-5n	007-2	1	0.0
X		CIFA	CIFAL		^{koc}	14		Purch	urcha	urcha	urcha	Irch _{ase}		
									~ ~			2		

Precision



LUK	0.50	0.50	0.52	0.83	0.50	0.50	0.69	0.50	0.47	0.50	0.57	0.73			1.0
AC AC	0.50	0.82	0.89	0.84	0.54	0.53	0.92	0.59	0.66	0.69	0.76	0.82			
CIFATOO	0.50	0.75	0.95	0.82	0.72	0.52	0.88	0.57	0.62	0.64	0.73	0.83		().8
UHAR se	0.50	0.75	0.95	0.87	0.78	0.52	0.86	0.56	0.61	0.64	0.73	0.82			
xion	0.50	0.68	0.91	0.75	0.86	0.51	0.82	0.54	0.57	0.60	0.66	0.75		(0.6
LOCAL	0.49	0.84	0.55	0.52	0.51	0.53	0.92	0.53	0.51	0.54	0.79	0.62			
Mains	0.50	0.76	0.95	0.83	0.74	0.52	0.86	0.57	0.62	0.65	0.74	0.84			
he y	0.50	0.82	0.86	0.80	0.54	0.53	0.90	0.59	0.66	0.60	0.73	0.71		().4
Purchas 10	0.50	0.84	0.80	0.76	0.55	0.53	0.92	0.59	0.66	0.68	0.76	0.85			
ourchase 20	0.50	0.83	0.88	0.83	0.53	0.53	0.92	0.59	0.66	0.69	0.78	0.83		().2
Y Thase 50	0.50	0.81	0.92	0.85	0.57	0.53	0.91	0.59	0.65	0.69	0.78	0.85			
Yurchase	0.50	0.79	0.95	0.85	0.61	0.53	0.90	0.58	0.64	0.67	0.77	0.86			
Inchaser ,	ldult	P, ZO	100	Face	tion	1157	Vews	کمی	e-10	e.20	e.50)	(J.U
\mathcal{Q}^{\vee}	X	CIFA,	CIFAR		Lo _{ci}	M		^U rch _ā	urchas	urchas	urchas	rchase			
			-						2^{2}	2^{-}	2 4	5			







Sounds Magic, Why?







77

Can we do better?

- Get rid of the shadow model
- Take the maximum, std, or entropy of the posterior as the score
 - The simplest attack
 - Unsupervised
 - Reason: overfitting





















member or non-member





Unsupervised attack

- 1. Statistical measures over posterior
 - Maximum, Std, Entropy
- Decide a threshold for the attack 2.
 - Above ??% maximal posterior is member















Threshold Picking









All Together











Defense

Layer 2

- Dropout
- Model stacking

Layer 1









Defense















Summary of ML-Leaks

- Machine learning models suffer from security and privacy attacks
 - Membership inference
- Three attacks with weak attacker assumption (more practical)
- How to evaluate a machine learning model?
 - Accuracy is enough?
 - Security and privacy matter
 - Just like buying a car
- A very promising direction for our research community!









Summary

- Social Network Privacy
 - walk2friends
 - Tagvisor
- Machine Learning Privacy
 - ML-Leaks



Thank you for your attention! Questions? http://yangzhangalmo.github.io/ @yangzhangalmo yang.zhang@cispa.saarland

 Ph.D. positions and summer interns at CISPA
CISPA-Stanford Program













