# Machine Learning and Cybercrime

Rebekah Overdorf





#### What is Machine Learning?

Machine learning uses statistical techniques to give computers the ability to "learn" with data, without being explicitly programmed.



## Traditional Programming



- Know something about the task
- Explicitly tell the machine how to complete it
- Test the model on something you know
- Deploy your program!

## Traditional Programming Example: Hate Speech Detection

- Know something about the task
- Explicitly tell the machine how to complete it
- Test the model on something you know
- Deploy your program!

- Hate speech usually contains the words in this list I made: list.txt
- Write a program that looks for these words
- Test the model on some hate speech that I found. "Wow it works on a bunch of these!"
- Deploy!

#### Machine Learning?



- Data Mining: Get data that you know something about
- Train the Model: "Teach" the machine about that data
- Test the model on something you know
- Deploy the program!

## Machine Learning Example: Hate Speech Detection

- Data Mining: Get data that you know something about
- Train the Model: "Teach" the machine about that data
- Test the model on something you know
- Deploy the program!

- I have 2 data sets, one of hate speech and one of not hate speech
- I throw all of this data at a "classifier" that uses statistics to figure out what is hate speech and what isn't.
- Test the model on some more hate speech that I found. "Wow it works on a bunch of these!"
- Deploy!

## This "machine learning" is super useful

- For understanding cybercrime to
  - Predict their communications
    - Overdorf et. al. 2018 (arxiv
  - Link accounts using writing style
    - Afroz et. al. 2012
  - Detect Phishing Sites and Emails
    - Xiang et. al. 2011, Basnet et. al. 2008
  - Identify items for sale/transactions
    - Portnoff et. al. 2017

- Identify Scams
  - McCoy et. al. 2016
- Identify Hate Speech
  - Davidson et. al. 2017
- Identify Child Porn
  - Peersman et. al. 2012, 2016
- Study their structure
  - Garg et. al 2015



Page 1 of 2 1 2 next	🔒 This topic is locked

#### PRIVATE INVESTIGATION SERVICE Dox Anyone Accurate Results Advanced Techniques BTC Accepted

previously Rigzorra MEMBER Posts: 25 Joined: Jun 25, 2017 Reputation: 7 Likes: 6	CYBER I In-depth de-	INVESTIGATIONS anonymization & investigation	
Credits: 0 Leecher level: 43 HALF YEAR REGISTERED			
		*	
	Everyone on the web has a unique fingerprint. We tak anyone online. With years of experience in the private right ser	que fingerprint. We take efficient and drastic measures into uncovering the true identity behind perience in the private investigation field, you can feel rest assured that you're ordering from the right service for your doxing needs.	
		<b>₽</b>	₩
			SWIFT DELIVERY
	O	RDER A PACKAGE	
	STANDARD DOX	PREMIUM DOX	
	\$14.99+ per dox	\$19.99+ per dax	

ol 2 1 2 next		
PRIVATE INVEST	IGATION SERVICE Dox Anyone Accurate Results Advanced Techniques BTC Accepted	PREMIUM DOX
StackFramed	Posted 23 August 2017-09:47 PM	¢10 00+
$\frown$		per dox
StackFramed		
$\bigvee$		Name
previously Rigzorra		Age
osts: 25 bined: Jun 25, 2017	CYBER INVESTIGA	Email Address
eputation: 7 kes: 6 wellin: 0	GET STARTED	Social Media Profiles
echer level: 43 HALF YEAR REGISTERED		IP Address
		Location
	REVEAL THEIR TRUE IDEN	Home Address
	Everyone on the web has a unique fingerprint. We take efficient and drastic r anyone online. With years of experience in the private investigation field, you a cloth secure for your device on	Household Information
		Phone Number
		Family Members
	ACCURATE DOXES EFFICIENT METHODS HIGHLY EXPERIENCED	Database Dumps
		And More
	ORDER A PACKAGE	ARRANGE A <b>PREMIUM DOX</b>

#### Public

Reply to this post

#### Private

Send message

#### Public

Reply to this post

#### Private

Send message

We always have this

We only have this when someone leaks data



## Potential Issues with deploying AI Systems

- ML often finds correlation, not causation
- Bias
- Base Rate Fallacy
- Privacy (e.g. membership attacks)
- Adversarial ML threats
- Result in antisocial and negative environmental outcomes
- Have adverse side effects
- Are built to only benefit a subset of users

- Externalize risks
- Produce errors due to ds distributional shift
- Result in systems that exploit states that fulfill the objective function, but do not complete the intended task
- Distribute errors unfairly
- Lack of transparency

•

- Example: Hate Speech Detection
- Hate Speech: 1 in every 1000 posts



- Example: Hate Speech Detection
- Hate Speech: 1 in every 1000 posts
- False Positive rate of 5%
- True Positive rate of 100%



- Example: Hate Speech Detection
- Hate Speech: 1 in every 1000 posts
- False Positive rate of 5%
- True Positive rate of 100%
- A new tweet is posted, and our ML method says it's hate speech.
- What is the probability that it is?



- Example: Hate Speech Detection
- Hate Speech: 1 in every 1000 posts
- False Positive rate of 5%
- True Positive rate of 100%
- A new tweet is posted, and our ML method says it's hate speech.
- What is the probability that it is?
  - ~2%
  - $p(hate|H) = \frac{p(H|hate)p(hate)}{p(H)}$ •  $p(hate|H) = \frac{1 * 0.001}{0.05095}$



- ML can assist in decisions related to cybercrime, but commonly the base rate is skewed. Most things aren't crime.
- Adding a person to inspect the decisions is crucial.



 Occurs when a all of the mistakes of the classifier are distributed to a subpopulation/group



- Example: Email Phishing Detection
- Yellow = Yes phishing email
- Blue = Benign



- Example: Email Phishing Detection
- Yellow = Yes phishing email
- Blue = Benign



- Example: Email Phishing Detection
- Yellow = Yes phishing email
- Blue = Benign



- Example: Email Phishing Detection
- Yellow = Yes phishing email
- Blue = Benign
- Squares = Emails in English
- Stars = Emails in German

 Occurs when a classifier is trained in one area and deployed in another.



- Example: Bot or Not?
- Squares = Bots in the US
- Stars = Bots in Central Asia

• Occurs when a classifier is trained in one area and deployed in another.

 $\bigstar$ 



- Example: Bot or Not?
- Squares = Bots in the US
- Stars = Bots in Central Asia

• Occurs when a classifier is trained in one area and deployed in another.

 $\bigstar$ 

![](_page_26_Figure_2.jpeg)

- Example: Bot or Not?
- Squares = Bots in the US
- Stars = Bots in Central Asia

• Occurs when a classifier is trained in one area and deployed in another.

![](_page_27_Figure_2.jpeg)

#### Conclusion

• Machine Learning can be a great tool for studying and preventing cybercrime, but is prone to adverse side effects that are often invisible to those deploying them.

#### Contact Info

- Rebekah@esat.kueluven.be
- @bekah\_Overdorf
- Computer-Supported Cooperative Crime
  - Garg, Afroz, Overdorf, Greenstadt.
- Under the Underground
  - Overdorf, Troncoso, McCoy, Greenstadt
- POTs: Protective Optimization Technologies
  - Overdorf, Balsa, Troncoso, Gurses
- Blogs, Twitter Feeds, and Reddit Comments: Cross-domain Authorship Attribution
  - Overdorf, Greenstadt
- How Unique is Your. onion? An Analysis of the Fingerprintability of Tor Onion Services
  - Overdorf, Juarez, Acar, Greenstadt, Diaz.

## Backups

## Automated Labeling

#### Labeling

![](_page_32_Figure_1.jpeg)

#### Post Example

![](_page_33_Figure_1.jpeg)

#### Post Example

![](_page_34_Figure_1.jpeg)

#### Post Example

![](_page_35_Figure_1.jpeg)
# Post Example



# Post Example



# Post Example

























 $1 \text{ if } f(\underline{f}) + f(\underline{f}) + f(\underline{f}) > \theta$ 0 otherwise













# Structure

# PM Data Analysis

- Who's important?
- Who has the most influence?
- What do communities look like?
- Is forum-enforced banning effective?
- Can we decide which user to remove to be the most disruptive?
- How does money move?
- How much is a product actually sold for?
- Which items actually sell?
- How does trust flow/scale?
- Can we link accounts?

## The Data



# The Data

- Carders
  - German
  - Carding
- L33tCrew
  - German
  - Carding
- BlackhatWorld
  - English
  - Young at leak
- Nulled
  - English
  - Huge
  - Varried



## Motivation

- How does trust scale?
  - How are the forums organized?
- Who is important?
  - Leaders, central members
- How can we disrupt them?
  - Or how can't we?

## Motivation

- How does trust scale?
  - How are the forums organized?
- Who is important?
  - Leaders, central members
- How can we disrupt them?
  - Or how can't we?

# **Community Detection**

- How are the users organized into communities?
- What do these communities look like?
- What does each specialize in, if at all?

# What do these communities look like?

- Dunbar Number (150)
- Structure
  - Mob-like vs Gang-like
- Topics
  - Topics are varied, not uniform, meaning communities specialize.

- Louvain Method for community detection
- LDA for topic modeling









# Motivation

- How does trust scale?
  - How are the forums organized?
- Who is important?
  - Leaders, central members
- How can we disrupt them?
  - Or how can't we?

# Who's important?

- Degree Centrality
  - Raw number of connections
  - Associated with higher trust
- Betweenness Centrality
  - Number of shortest paths that pass through the node
  - More information
- Closeness Centrality
  - How far is this node from all other nodes
  - Lowest transaction costs
- Eigenvector Centrality
  - How much influence does this node and it's neighbors have

### Results

		Blac	khatW	orld		Carders					L33tCrew					
Cent.	С	B	ID	OD	D	С	B	ID	OD	D	С	B	ID	OD	D	
Е	0.08	0.66	0.81	0.50	0.71	-0.43	0.79	0.91	0.62	0.77	-0.55	0.85	0.95	0.84	0.91	
С		0.33	0.18	0.51	0.37		-0.19	-0.33	-0.11	-0.21		-0.39	-0.51	-0.35	-0.41	
В			0.81	0.84	0.88			0.90	0.83	0.90			0.91	0.92	0.94	
ID				0.56	0.85				0.71	0.88				0.88	0.96	
OD					0.87					0.94					0.96	

### Results

		BlackhatWorld						Carders						L33tCrew				
Cen	t. C		B	ID	OD	D	С	B	ID	OD	D	С	B	ID	OD	D		
E	0.0	8	0.66	0.81	0.50	0.71	-0.43	0.79	0.91	0.62	0.77	-0.55	0.85	0.95	0.84	0.91		
С			0.33	0.18	0.51	0.37		-0.19	-0.33	-0.11	-0.21		-0.39	-0.51	-0.35	-0.41		
В				0.81	0.84	0.88			0.90	0.83	0.90			0.91	0.92	0.94		
ID					0.56	0.85				0.71	0.88				0.88	0.96		
OD						0.87					0.94					0.96		

- On BlackhatWorld Everything is correlated.
- On L33tCrew and Carders Everything but closeness centrality is correlated.
  - Closeness Centrality How far is this node from all other nodes
#### Motivation

- How does trust scale?
  - How are the forums organized?
- Who is important?
  - Leaders, central members
- How can we disrupt them?
  - Or how can't we?

# Banning

- Duplicate Accounts
- Ripping
- Spamming

# What happens when members are banned?



# What happens when members are banned?



# What happens when members are banned?



	BlackhatWorld		Carders		L33tCrew	7
СМ	⊿ACC	⊿APL	$\triangle ACC$	<b>AAPL</b>	⊿ACC	<b>AAPL</b>
Betweenness (B)	-0.39	0.32	-0.12***	-0.05*	-0.05	0.11
Closeness (C)	0.07	-0.12	-0.07**	-0.05*	-0.19*	0.11
Degree (D)	-0.15	0.22	-0.19***	-0.03	-0.06	0.10
Eigenvector (E)	0.07	-0.12	-0.14***	-0.04	-0.01	0.004
p-value: 0.05> * > 0.01 > ** > 0.001 > ***						

	BlackhatWorld		Carders		L33tCrew	
СМ	⊿ACC	⊿APL		<b>AAPL</b>		⊿APL
Betweenness (B)	-0.39	0.32	-0.12***	-0.05*	-0.05	0.11
Closeness (C)	Feŵ Banne	ed-0.12	-0.07**	-0.05*	-0.19*	
Degree (D)	-(Mémbers	0.22	-0.19***	-0.03	-0.06	
Eigenvector (E)	0.07	-0.12	-0.14***	-0.04	-0.01	0.004
p-value: $0.05 > * > 0.01 > ** > 0.001 > ***$						

	BlackhatWorld		Carders		L33tCrew	
СМ		⊿APL	$\triangle ACC$	⊿APL		⊿APL
Betweenness (B)	-0.39	0.32	-0.12***	-0.05*	-0.05	0.11
Closeness (C)	Few Banne	ed 0.12	-0.07**	-0.05*	-0.19*	
Degree (D)	-(Mémbers	s 0.22	-0.19***	-0.03	-0.06	
Eigenvector (E)	0.07	-0.12	-0.14***	-0.04	-0.01	0.004
p-value: $0.05 > * > 0.01 > ** > 0.001 > ***$						

	BlackhatWorld		Carders		L33tCrew	7
СМ	⊿ACC	⊿APL		⊿APL	⊿ACC	⊿APL
Betweenness (B)	-0.39	0.32	-0.12***	-0.05*	-0.05	0.11
Closeness (C)	Feŵ Bann	ed-0.12	-0.07**	-0.05*	-0.19*	
Degree (D)	-(Mémber	s 0.22	-0.19***	-0.03	-0.06	
Eigenvector (E)	0.07	-0.12	-0.14***	-0.04	-0.01	0.004
p-value: 0.05> * > 0.01 > ** > 0.001 > ***						

• Individuals being banned are not close to other nodes.