



# Identifying Twitter Spam by Utilizing Random Forests

**Humza Haider**

Division of Science and Mathematics  
University of Minnesota Morris

2017-04-15



- ▶ Top social media platforms



- ▶ Top social media platforms
- ▶ 500 million tweets per day



- ▶ Top social media platforms
- ▶ 500 million tweets per day
- ▶ Attracts spammers and malicious users



- ▶ Top social media platforms
- ▶ 500 million tweets per day
- ▶ Attracts spammers and malicious users
- ▶ Twitter spam: Any unsolicited, repeated actions that negatively impact other users



- ▶ Top social media platforms
- ▶ 500 million tweets per day
- ▶ Attracts spammers and malicious users
- ▶ Twitter spam: Any unsolicited, repeated actions that negatively impact other users
- ▶ How can we identify spammers?



- ▶ Top social media platforms
- ▶ 500 million tweets per day
- ▶ Attracts spammers and malicious users
- ▶ Twitter spam: Any unsolicited, repeated actions that negatively impact other users
- ▶ How can we identify spammers?
  - ▶ Manual classification



- ▶ Top social media platforms
- ▶ 500 million tweets per day
- ▶ Attracts spammers and malicious users
- ▶ Twitter spam: Any unsolicited, repeated actions that negatively impact other users
- ▶ How can we identify spammers?
  - ▶ Manual classification
  - ▶ URL blacklisting





- ▶ Top social media platforms
- ▶ 500 million tweets per day
- ▶ Attracts spammers and malicious users
- ▶ Twitter spam: Any unsolicited, repeated actions that negatively impact other users
- ▶ How can we identify spammers?
  - ▶ Manual classification
  - ▶ URL blacklisting
  - ▶ **Machine learning classification**



## Background

- Decision Trees
- Random Forests
- Model Evaluation

## Methods

- Tweet and User Content Features
- Geo-Tagged Features
- Time Features

## Results

## Conclusion



- ▶ Decision Trees



- ▶ Decision Trees
  - ▶ Machine learning technique for classification



- ▶ Decision Trees
  - ▶ Machine learning technique for classification
  - ▶ Classifies an observation based on features available in a dataset

# Background

## Decision Trees



URL	Account Age	Reported	Class
No	Old	Yes	<b>Not Spam</b>
No	Old	Yes	<b>Not Spam</b>
No	Old	No	<b>Not Spam</b>
No	New	No	<b>Not Spam</b>
Yes	New	Yes	<b>Spam</b>
No	New	Yes	<b>Spam</b>
No	Old	Yes	<b>Spam</b>
Yes	New	No	<b>Not Spam</b>
⋮	⋮	⋮	⋮

# Background

## Decision Trees



URL	Account Age	Reported	Class
No	Old	Yes	<b>Not Spam</b>
No	Old	Yes	<b>Not Spam</b>
No	Old	No	<b>Not Spam</b>
No	New	No	<b>Not Spam</b>
Yes	New	Yes	<b>Spam</b>
No	New	Yes	<b>Spam</b>
No	Old	Yes	<b>Spam</b>
Yes	New	No	<b>Not Spam</b>

# Background

## Decision Trees



**32% Spam**  
**68% Not Spam**

**8 Tweets**  
**3 Spam**  
**5 Not Spam**



# Background

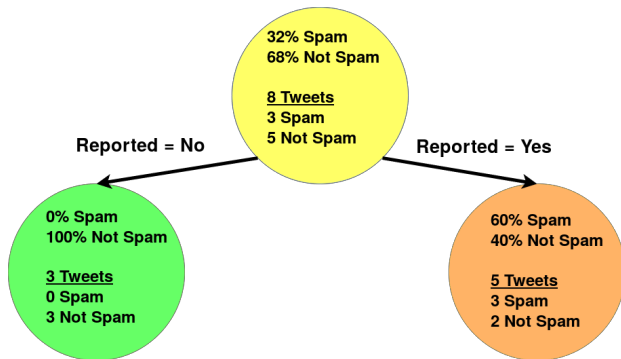
## Decision Trees



URL	Account Age	Reported	Class
No	Old	Yes	Not Spam
No	Old	Yes	Not Spam
No	Old	No	Not Spam
No	New	No	Not Spam
Yes	New	Yes	Spam
No	New	Yes	Spam
No	Old	Yes	Spam
Yes	New	No	Not Spam

# Background

## Decision Trees



# Background

## Decision Trees



URL	Account Age	Reported	Class
No	Old	Yes	<b>Not Spam</b>
No	Old	Yes	<b>Not Spam</b>
<del>No</del>	<del>Old</del>	<del>No</del>	<del><b>Not Spam</b></del>
<del>No</del>	<del>New</del>	<del>No</del>	<del><b>Not Spam</b></del>
Yes	New	Yes	<b>Spam</b>
No	New	Yes	<b>Spam</b>
No	Old	Yes	<b>Spam</b>
<del>Yes</del>	<del>New</del>	<del>No</del>	<del><b>Not Spam</b></del>

# Background

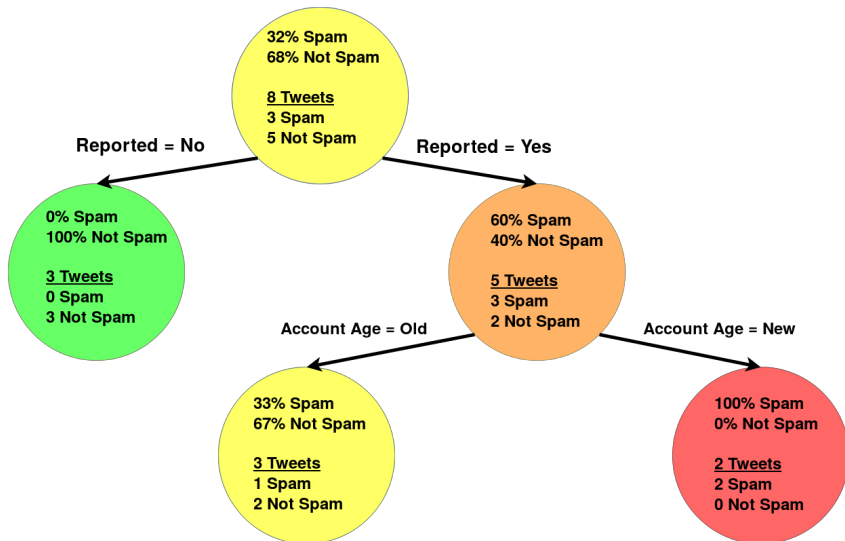
## Decision Trees



URL	Account Age	Reported	Class
No	Old	Yes	Not Spam
No	Old	Yes	Not Spam
<del>No</del>	<del>Old</del>	<del>No</del>	<del>Not Spam</del>
<del>No</del>	<del>New</del>	<del>No</del>	<del>Not Spam</del>
Yes	New	Yes	Spam
No	New	Yes	Spam
No	Old	Yes	Spam
<del>Yes</del>	<del>New</del>	<del>No</del>	<del>Not Spam</del>

# Background

## Decision Trees



# Background

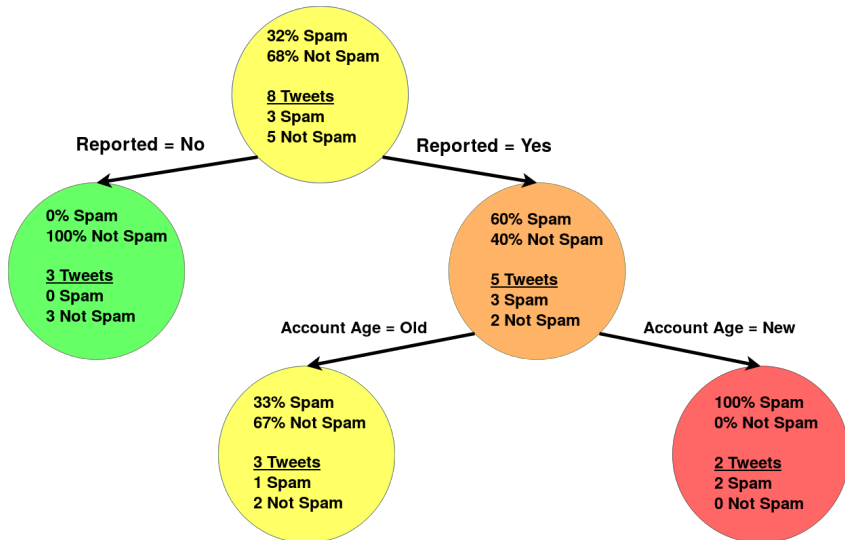
## Decision Trees



URL	Account Age	Reported	Class
No	Old	Yes	<b>Not Spam</b>
No	Old	Yes	<b>Not Spam</b>
<del>No</del>	<del>Old</del>	<del>No</del>	<del><b>Not Spam</b></del>
<del>No</del>	<del>New</del>	<del>No</del>	<del><b>Not Spam</b></del>
<del>Yes</del>	<del>New</del>	<del>Yes</del>	<del><b>Spam</b></del>
<del>No</del>	<del>New</del>	<del>Yes</del>	<del><b>Spam</b></del>
<del>No</del>	<del>Old</del>	<del>Yes</del>	<del><b>Spam</b></del>
<del>Yes</del>	<del>New</del>	<del>No</del>	<del><b>Not Spam</b></del>

# Background

## Decision Trees



# Background

## Decision Trees



**Nic McPhee** @NicMcPhee · Jan 29

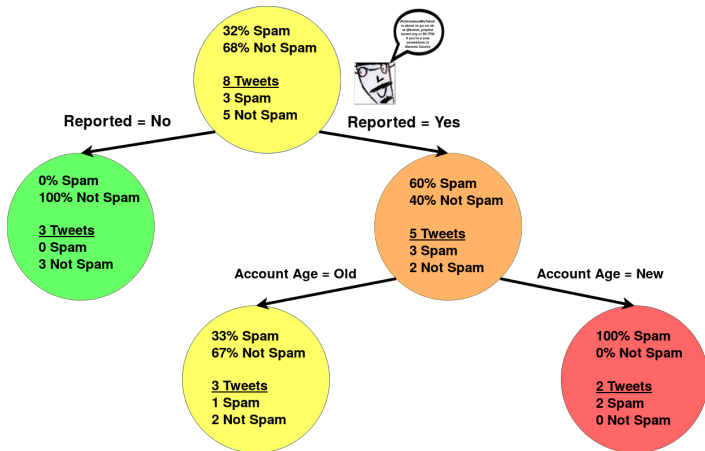
#UnhinderedByTalent is about to go on air at @kumm\_playlist kumm.org or 89.7FM if you're a cow somewhere in Stevens County.

URL	Account Age	Reported	Class
Yes	Old	Yes	TBD



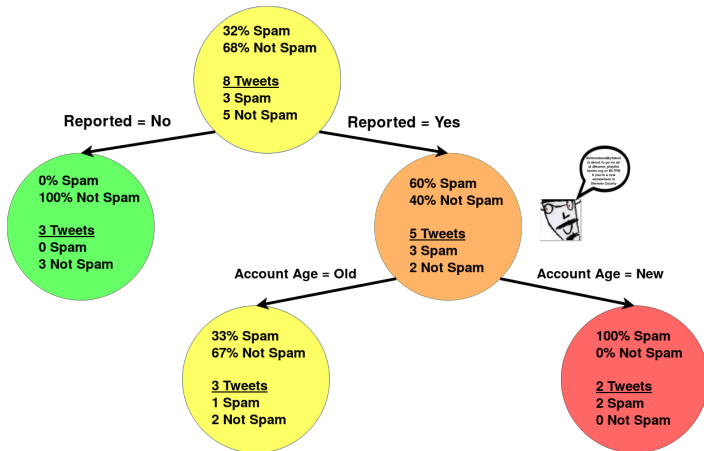
# Background

## Decision Trees



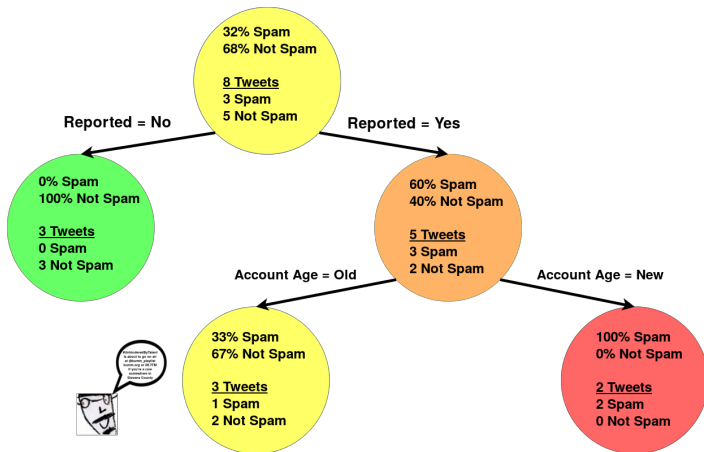
# Background

## Decision Trees



# Background

## Decision Trees



# Background

## Decision Trees



- ▶ How are splits decided?

# Background

## Decision Trees



- ▶ How are splits decided?
  - ▶ Entropy

# Background

## Decision Trees



- ▶ How are splits decided?
  - ▶ Entropy
  - ▶ Information Gain

# Background

## Decision Trees



- ▶ How are splits decided?
  - ▶ Entropy
  - ▶ Information Gain
- ▶ Trees seem pretty neat! Why do I need a whole forest?

# Background

## Decision Trees

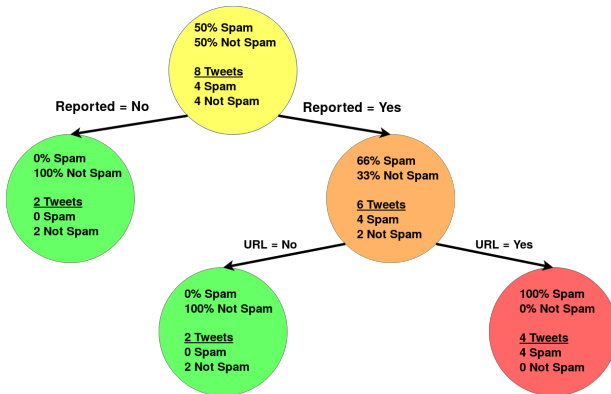


- ▶ How are splits decided?
  - ▶ Entropy
  - ▶ Information Gain
- ▶ Trees seem pretty neat! Why do I need a whole forest?
  - ▶ Disagreement in decisions between different trees



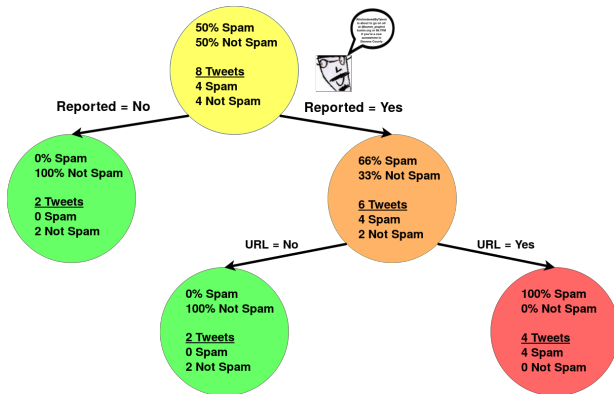
# Background

## Decision Trees



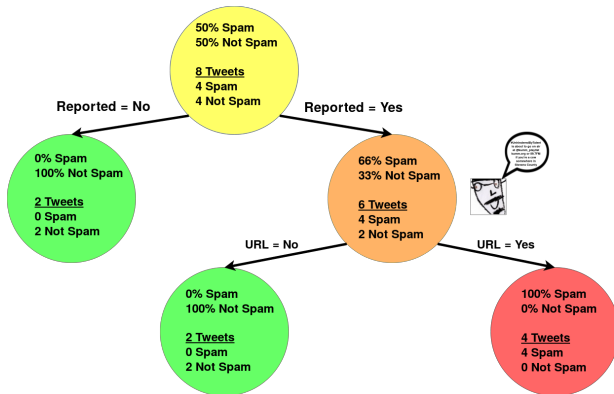
# Background

## Decision Trees



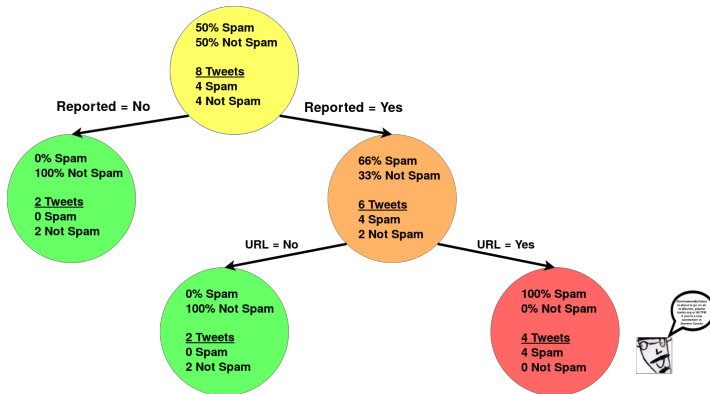
# Background

## Decision Trees



# Background

## Decision Trees



# Background

## Random Forests



- ▶ How do we handle disagreement?

# Background

## Random Forests



- ▶ How do we handle disagreement?
  - ▶ Train many trees on samples of the data (**Bagging**)



- ▶ How do we handle disagreement?
  - ▶ Train many trees on samples of the data (**Bagging**)
  - ▶ Don't let trees access all the features (**Feature Bagging**)

# Background

## Random Forests



- ▶ How do we handle disagreement?
  - ▶ Train many trees on samples of the data (**Bagging**)
  - ▶ Don't let trees access all the features (**Feature Bagging**)
- ▶ After we make a bunch of trees, how do we combine them?



# Background

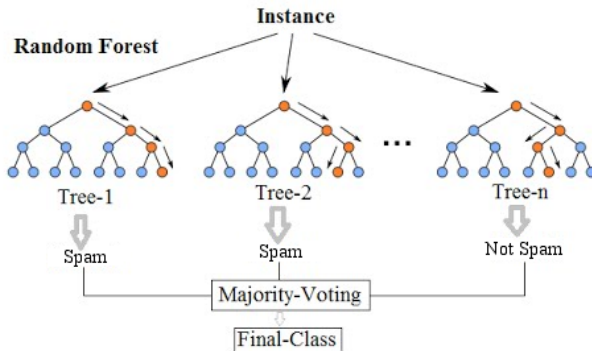
## Random Forests



- ▶ How do we handle disagreement?
  - ▶ Train many trees on samples of the data (**Bagging**)
  - ▶ Don't let trees access all the features (**Feature Bagging**)
- ▶ After we make a bunch of trees, how do we combine them?
  - ▶ Majority vote



### Random Forest Simplified



Source: <https://i.ytimg.com/vi/ajTc5y3OqSQ/hqdefault.jpg>

# Background

## Model Evaluation



- ▶ How do we evaluate a random forest's performance?



- ▶ How do we evaluate a random forest's performance?

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative



### Accuracy

- ▶ “How many tweets were correctly identified?”

#### Truth

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative



### Accuracy

- ▶ “How many tweets were correctly identified?”

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative

### Accuracy

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$



### Precision ( $p$ )

- ▶ “How good is our spam prediction?”

#### Truth

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative



### Precision ( $p$ )

- ▶ “How good is our spam prediction?”

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative

### Precision

$$\text{Precision} = \frac{TP}{TP + FP}$$





### Recall ( $r$ )

- ▶ “How much spam was identified?”

#### Truth

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative



### Recall ( $r$ )

- ▶ “How much spam was identified?”

#### Truth

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative

### Recall

$$\text{Recall} = \frac{TP}{TP + FN}$$



- ▶ F-measure ( $F$ )
  - ▶ Harmonic mean of Precision and Recall
  - ▶ Equally weights both Precision and Recall

### F-Measure

$$\text{F-measure} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$



## Background

- Decision Trees
- Random Forests
- Model Evaluation

## Methods

- Tweet and User Content Features
- Geo-Tagged Features
- Time Features

## Results

## Conclusion

# Methods

## Tweet and User Content Features



- ▶ Chen et al. identify tweets, as opposed to users

# Methods

## Tweet and User Content Features



- ▶ Chen et al. identify tweets, as opposed to users
- ▶ Utilized 12 features directly accessible from a tweet

# Methods

## Tweet and User Content Features



- ▶ Chen et al. identify tweets, as opposed to users
- ▶ Utilized 12 features directly accessible from a tweet
  - ▶ 6 user features



- ▶ Chen et al. identify tweets, as opposed to users
- ▶ Utilized 12 features directly accessible from a tweet
  - ▶ 6 user features
  - ▶ 6 tweet features





- ▶ Chen et al. identify tweets, as opposed to users
- ▶ Utilized 12 features directly accessible from a tweet
  - ▶ 6 user features
  - ▶ 6 tweet features
- ▶ User Features
  - ▶ Age in days of account
  - ▶ Number of followers, followees
  - ▶ Number of tweets

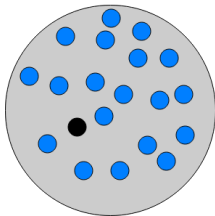


- ▶ Chen et al. identify tweets, as opposed to users
- ▶ Utilized 12 features directly accessible from a tweet
  - ▶ 6 user features
  - ▶ 6 tweet features
- ▶ User Features
  - ▶ Age in days of account
  - ▶ Number of followers, followees
  - ▶ Number of tweets
- ▶ Tweet Features
  - ▶ Number of hashtags (#)
  - ▶ Number of mentions
  - ▶ Number of URLs

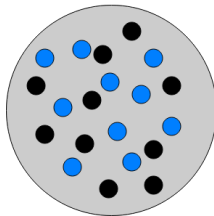


- ▶ Two sets of testing data.

**5% Spam**



**50% Spam**





## Background

- Decision Trees
- Random Forests
- Model Evaluation

## Methods

- Tweet and User Content Features
- Geo-Tagged Features
- Time Features

## Results

## Conclusion



- ▶ So, what is a geo-tagged tweet?



- ▶ So, what is a geo-tagged tweet?



**Morris Weather**  
@MorrisMNWeather



Follow



The weather is boring. 50°F and Light Rain.  
**#MorrisMNWeather**

6:02 PM - 9 Apr 2017 from Morris, MN



- ▶ Guo and Chen identify non-personal users
  - ▶ Spammers
  - ▶ Bots
  - ▶ Business accounts

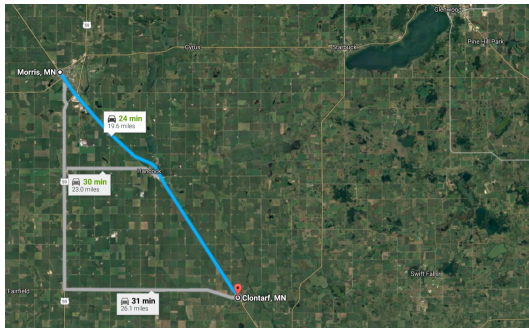


- ▶ Guo and Chen identify non-personal users
  - ▶ Spammers
  - ▶ Bots
  - ▶ Business accounts
- ▶ Features:

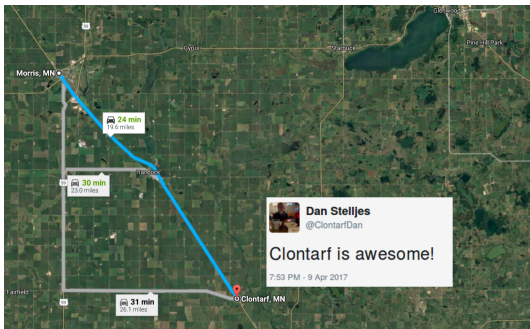




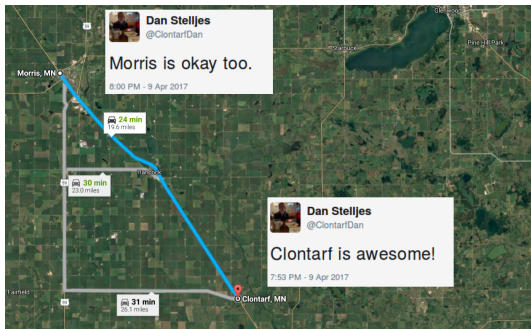
- ▶ Guo and Chen identify non-personal users
  - ▶ Spammers
  - ▶ Bots
  - ▶ Business accounts
- ▶ Features:
  - ▶ Tweeting Speed



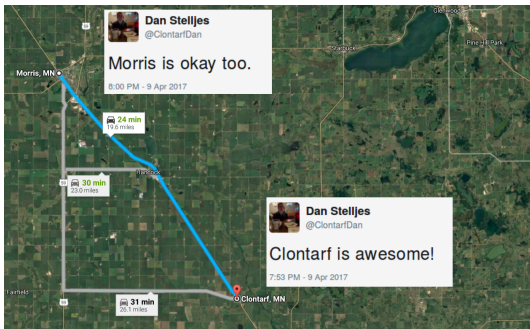
► 19.6 Miles



- ▶ 19.6 Miles
- ▶ From Clontarf at 7:53 PM



- ▶ 19.6 Miles
- ▶ From Clontarf at 7:53 PM
- ▶ From Morris at 8:00 PM



- ▶ 19.6 Miles
- ▶ From Clontarf at 7:53 PM
- ▶ From Morris at 8:00 PM
- ▶ Tweeting speed =  $\frac{19.6 \text{ miles}}{7 \text{ minutes}} = 2.8 \text{ miles per minute (168 MPH)}$



► Features:



- ▶ Features:
  - ▶ Max Speed



- ▶ Features:
  - ▶ Max Speed
  - ▶ Mean Speed





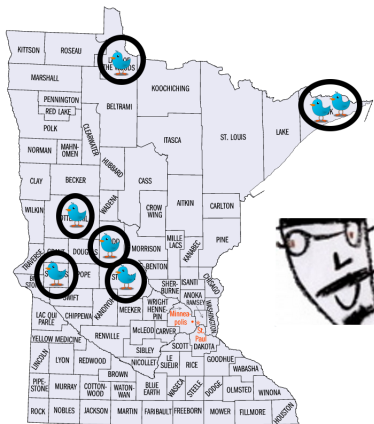
- ▶ Features:
  - ▶ Max Speed
  - ▶ Mean Speed
  - ▶ Max Distance (connected to Max Speed)

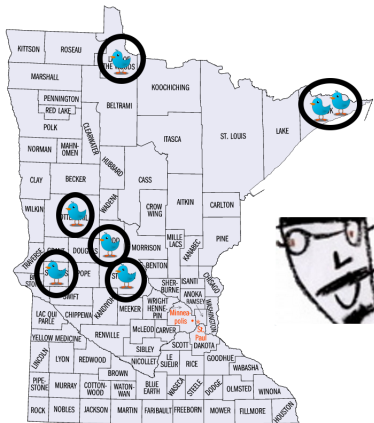


- ▶ Features:
  - ▶ Max Speed
  - ▶ Mean Speed
  - ▶ Max Distance (connected to Max Speed)
  - ▶ Mean number of times a user exceeds 90 MPH per month

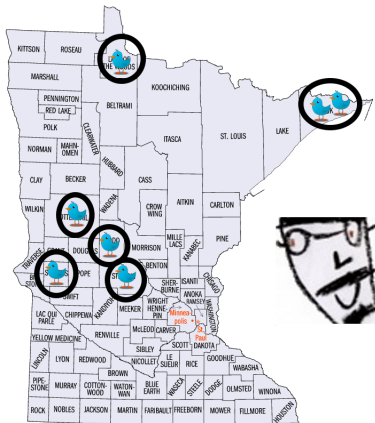


- ▶ Features:
  - ▶ Max Speed
  - ▶ Mean Speed
  - ▶ Max Distance (connected to Max Speed)
  - ▶ Mean number of times a user exceeds 90 MPH per month
- ▶ County based features





- ▶ Number of times a user crosses county borders per month



- ▶ Number of times a user crosses county borders per month
- ▶ Mean number of counties a user has been to per month



## Background

- Decision Trees
- Random Forests
- Model Evaluation

## Methods

- Tweet and User Content Features
- Geo-Tagged Features
- Time Features

## Results

## Conclusion



- ▶ Washha et al. classify spammers on a user level





- ▶ Washha et al. classify spammers on a user level
- ▶ Motivated to use time since altering time dependent features is a challenge.



- ▶ Washha et al. classify spammers on a user level
- ▶ Motivated to use time since altering time dependent features is a challenge.
- ▶ Features that spammers can easily manipulate:



- ▶ Washha et al. classify spammers on a user level
- ▶ Motivated to use time since altering time dependent features is a challenge.
- ▶ Features that spammers can easily manipulate:
  - ▶ Number of URLs



- ▶ Washha et al. classify spammers on a user level
- ▶ Motivated to use time since altering time dependent features is a challenge.
- ▶ Features that spammers can easily manipulate:
  - ▶ Number of URLs
  - ▶ Number of Hashtags



- ▶ Washha et al. classify spammers on a user level
- ▶ Motivated to use time since altering time dependent features is a challenge.
- ▶ Features that spammers can easily manipulate:
  - ▶ Number of URLs
  - ▶ Number of Hashtags
  - ▶ Including Geo-tags



## Features

- ▶ Differences in Account Age



## Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts



## Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts
  - ▶ Likely to be made at the same time





## Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts
  - ▶ Likely to be made at the same time
  - ▶ Followers



## Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts
  - ▶ Likely to be made at the same time
  - ▶ Followers
  - ▶ Followees



## Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts
  - ▶ Likely to be made at the same time
  - ▶ Followers
  - ▶ Followees
  - ▶ Bi-directional relationships



### Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts
  - ▶ Likely to be made at the same time
  - ▶ Followers
  - ▶ Followees
  - ▶ Bi-directional relationships
- ▶ Time weighted correlations:



### Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts
  - ▶ Likely to be made at the same time
  - ▶ Followers
  - ▶ Followees
  - ▶ Bi-directional relationships
- ▶ Time weighted correlations:
  - ▶ URLs



### Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts
  - ▶ Likely to be made at the same time
  - ▶ Followers
  - ▶ Followees
  - ▶ Bi-directional relationships
- ▶ Time weighted correlations:
  - ▶ URLs
  - ▶ Mentions



## Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts
  - ▶ Likely to be made at the same time
  - ▶ Followers
  - ▶ Followees
  - ▶ Bi-directional relationships
- ▶ Time weighted correlations:
  - ▶ URLs
  - ▶ Mentions
  - ▶ Hashtags



### Features

- ▶ Differences in Account Age
  - ▶ Spammers have multiple accounts
  - ▶ Likely to be made at the same time
  - ▶ Followers
  - ▶ Followees
  - ▶ Bi-directional relationships
- ▶ Time weighted correlations:
  - ▶ URLs
  - ▶ Mentions
  - ▶ Hashtags
- ▶ Tweet similarity weighted by time





## Background

- Decision Trees
- Random Forests
- Model Evaluation

## Methods

- Tweet and User Content Features
- Geo-Tagged Features
- Time Features

## Results

## Conclusion

# Results



► Accuracy:

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative

$$\frac{TP+TN}{TP+FP+TN+FN}$$

"How many tweets were correctly identified?"

► Precision ( $p$ ):

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative

$$\frac{TP}{TP+FP}$$

"How good is our spam prediction?"

► Recall ( $r$ ):

		Truth	
		Spam	Not Spam
Prediction	Spam	True Positive	False Positive
	Not Spam	False Negative	True Negative

$$\frac{TP}{TP+FN}$$

"How much spam was identified?"

► F-measure ( $F$ ):  $2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$

Model Results of the Three Studies

Study	% Spam	$p$	$r$	$F$	Accuracy
User/Tweet Features: I	50.0%	0.929	0.943	0.936	0.936
User/Tweet Features: II	5.0%	0.929	0.407	0.566	0.978
Geo-tagged Features	21.4%	0.959	0.959	0.958	0.959
Time Features	46.9%	0.932	0.931	0.931	0.931



- ▶ Classification via random forest
- ▶ Recall ( $r$ ) may drop when test set contains a low proportion of spam
  - ▶ Future work: Apply this finding to geo-tagged tweets and time features
- ▶ Future spam classification by Twitter: Random forests?







- ▶ Peter Dolan
  - ▶ For acting as my advisor for this research project and his continued friendship for the past few years
- ▶ Elena Machkasova
  - ▶ For the invaluable advice and insightful comments throughout the entire senior seminar course
- ▶ Jacob Opdahl
  - ▶ For the exceptional revisions and comments he made as my alumni review





-  Fabricio Benevenuto, Gabriel Magno, Tiago Rodrigues, and Virgilio Almeida, *Detecting spammers on twitter*, Collaboration, electronic messaging, anti-abuse and spam conference (CEAS), vol. 6, 2010, p. 12.
-  Leo Breiman, *Random forests*, Machine learning **45** (2001), no. 1, 5–32.
-  Chao Chen, Jun Zhang, Xiao Chen, Yang Xiang, and Wanlei Zhou, *6 million spam tweets: A large ground truth for timely twitter spam detection*, 2015 IEEE International Conference on Communications (ICC), IEEE, 2015, pp. 7065–7070.
-  Diansheng Guo and Chao Chen, *Detecting non-personal and spam users on geo-tagged twitter network*, Transactions in GIS **18** (2014), no. 3, 370–384.



-  Miroslav Kubat, *An introduction to machine learning*, 1st ed., Springer Publishing Company, Incorporated, 2015.
-  Nikita Spirin and Jiawei Han, *Survey on web spam detection: Principles and algorithms*, SIGKDD Explor. Newsl. **13** (2012), no. 2, 50–64.
-  Claude Elwood Shannon, *A mathematical theory of communication*, ACM SIGMOBILE Mobile Computing and Communications Review **5** (2001), no. 1, 3–55.
-  Igor Santos, Igor Miñambres-Marcos, Carlos Laorden, Patxi Galán-García, Aitor Santamaría-Ibirika, and Pablo García Bringas, *Twitter content-based spam filtering*, pp. 449–458, Springer International Publishing, Cham, 2014.



-  Kurt Thomas, Chris Grier, Dawn Song, and Vern Paxson, *Suspended accounts in retrospect: An analysis of twitter spam*, Proceedings of the 2011 ACM SIGCOMM Conference on Internet Measurement Conference (New York, NY, USA), IMC '11, ACM, 2011, pp. 243–258.
-  Mahdi Washha, Aziz Qaroush, and Florence Sedes, *Leveraging time for spammers detection on twitter*, Proceedings of the 8th International Conference on Management of Digital EcoSystems (New York, NY, USA), MEDES, ACM, 2016, pp. 109–116.