# Reddit: Analysis of Suspect Accounts

STAT 418 Project
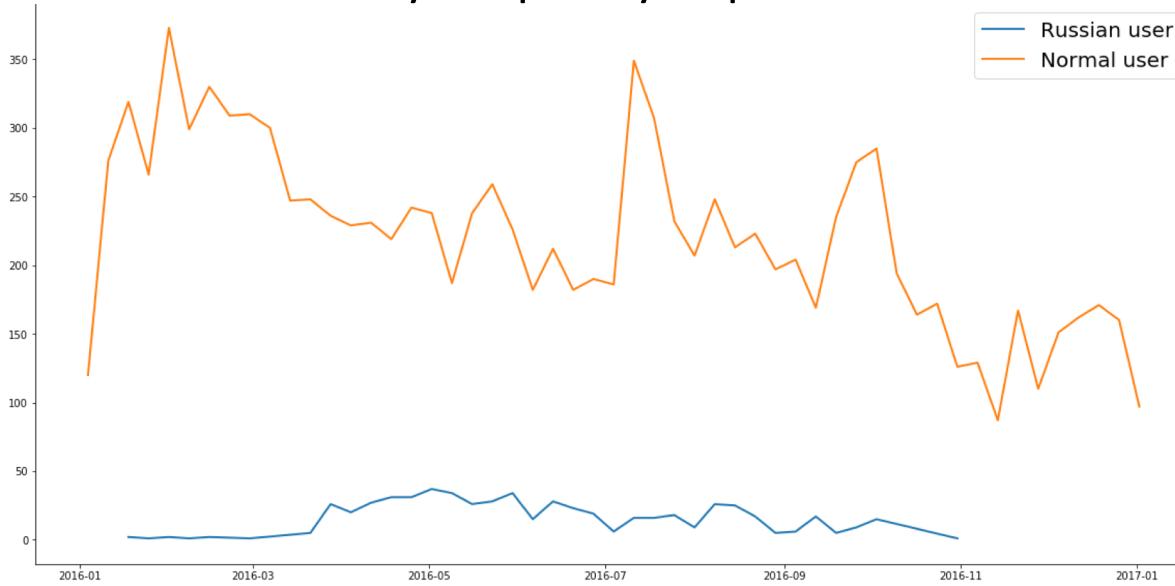
Janella Shu

# Background and Dataset

- Reddit's 2017 transparency report: 944 accounts suspected to have originated from IRA
  - IRA employed fake accounts to influence the 2016 United States presidential elections
- Subreddit: Bad_Cop_No_Donut
  - Currently has 203k members
  - "Law enforcement abuse stories regarding: abuse of power, corruption, and other misfortunes in developing police states. "
- Date range: January 1, 2016 – December 31, 2016
  - 12,272 posts: 11,688 from normal accounts; 584 from Russian accounts
- Submissions
  - id: ID of the submission
  - author
  - created_utc : time the submission was created, Unix Time
  - is_self: whether or not the submission is a selfpost
  - selftext: the submission's selftext
  - title: the title of the submission
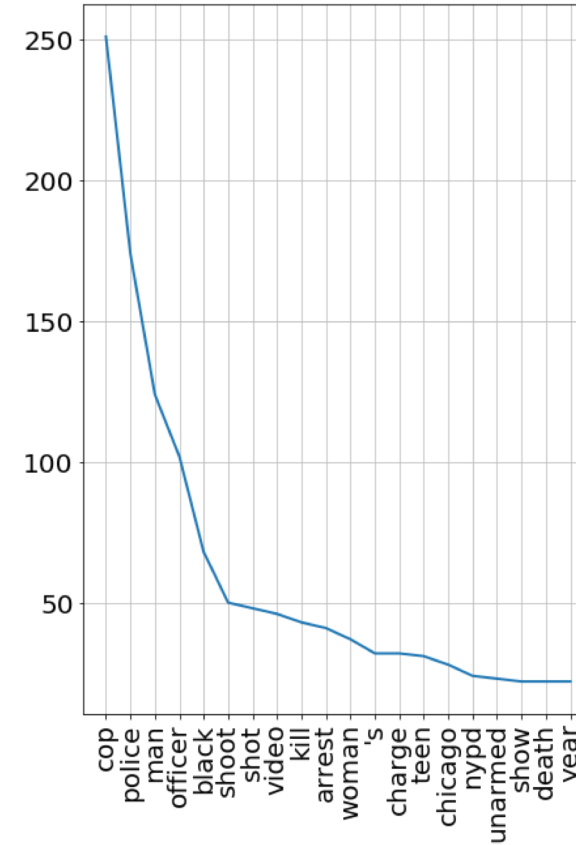  - url: the URL the submission links to, or the permalink if a self post

# Exploratory Data Analysis
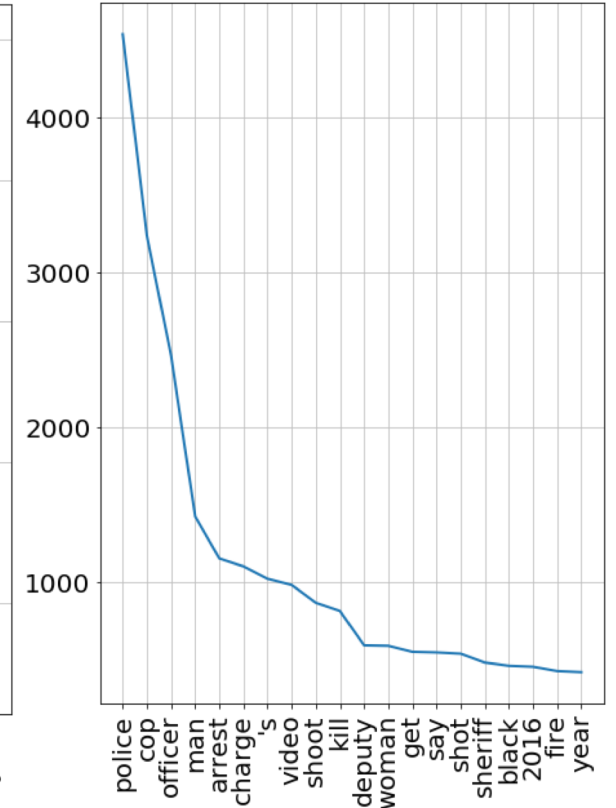
## Weekly frequency of posts



## # of posts by hour of the day



## Top 20 most frequently used words in title

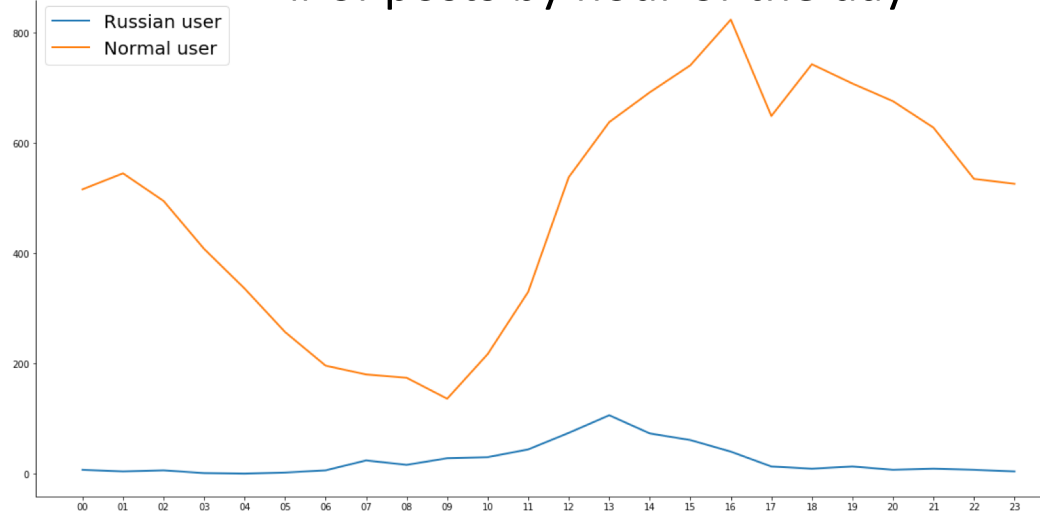### Russian users



### Normal users
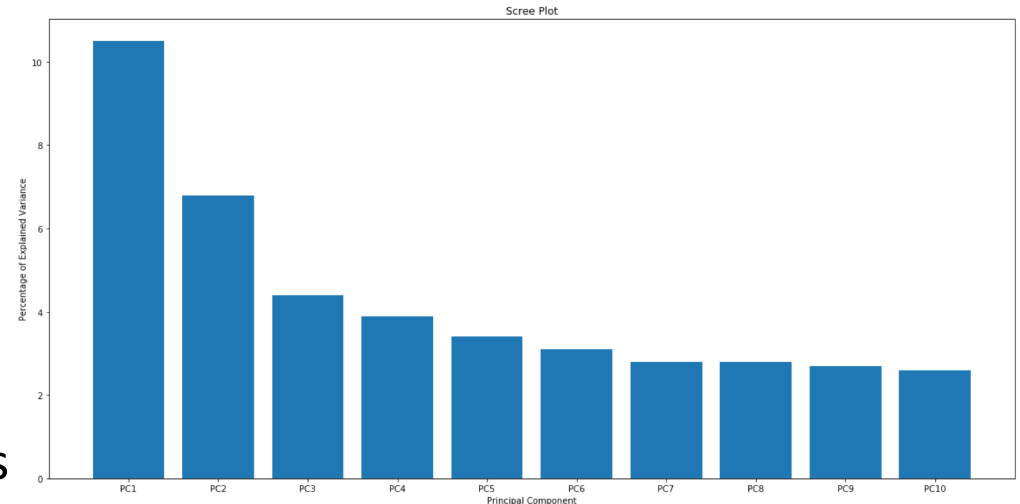


- Submission title was tokenize, removed punctuation and stop words, and lemmatized each word

# Data Analysis

- Converted text to vector using Term Frequency – Inverse Document (TF-IDF)
  - TF: how often the word appears within a document
  - IDF: measure of how significant that word is in the whole corpus
- Split data into approx. 70/30 for training and testing set



- Logistic Regression Model
  - Model 1: -3.181 + 1.066 *PC1 + 0.489*PC2
    - PCA to reduce dimensionality of td-idf scores
  - Model 2: -3.573 + 1.279 *cop + 1.643*black + 1.306*man + 1.596*teen
  - Model 3: -3.88+ 6.41 *russian_prop + 0.038*black_man
    - russian_prop: indicator of if submission linked to Russian propaganda website(i.e. DoNotShoot.us and BlackMattersUS.com), https://www.thewrap.com/russian-propaganda-us-election-donald-trump-hillary-clinton-wikileaks-drudge-report-info-wars/
    - black_man: tf-idf score for the bigram "black man"

# Model Evaluation and Dash App

**Model 1**

predicted

| | | 0 | 1 |
|---|---|---|---|
| | | 3246 | 201 |
| actual | 0 | | |
| | 1 | 209 | 12 |

**Model 2**

predicted

| | | 0 | 1 |
|---|---|---|---|
| | 0 | 3370 | 77 |
| actual | | | |
| | 1 | 195 | 26 |

**Model 3**

predicted

| | | 0 | 1 |
|---|---|---|---|
| | 0 | 3433 | 14 |
| actual | | | |
| | 1 | 106 | 115 |

| | | | |
|---|---|---|---|
| Sensitivity (TP/actual Y) | 12/221 = 0.054 | 27/221 = 0.122 | 115/221 = 0.52 |
| Precision (TP/predicted Y) | 12/213= 0.056 | 26/113 = 0.23 | 115/119 = 0.96 |

- Dash app: http://18.236.183.235:8050

# Next Steps

- Collect more submissions so that I can try other methods like doc2vec

- Try other subreddits or date range

- Analyze user comments

- Add more figures to my Dash app