

CSE648: Privacy and Security in Online Social Media

Homework - 2 (150 Marks)

Deadline: 23:59 Hrs April 15, 2021

Instructions:

- Use Jupyter Notebooks / Google Colab for doing this assignment and showing analysis, graphs and code. Name the file as **name_rollnumber_section.ipynb**.
 - Write all inferences and include the plots in a separate file named **name_rollnumber.pdf**.
 - Please cite any sources that you might use in the process.
 - Please write your own code. All code will be tested for plagiarism and if found, institute policy will be followed.
 - *Attempt the homework individually.*
 - Submit your code and report as individual files. **Do not** zip your submissions.
-

Section - I (Online Policing)

[95 Marks]

In this section, you are required to select any one of the twitter handles of police in India. Preferably, pick handles that are verified or have a significant number of followers and tweets for a better quality of analysis. Feel free to pick the handles that you are interested in analyzing, e.g. if you are from North Delhi, please analyze [@DCPNorthDelhi](#). This will give you an understanding how your location specific handles are doing. Collect as many tweets from their account as possible within the rate limits of Twitter API. You are free to use an API wrapper for the same. We recommend using [Tweepy](#).

Deliverables: {Code, Data Collected}

[10 + 5 = 15 Marks]

Q1. In this question, consider the tweets from their account which are replies to tweets by other users mentioning the handle. Retrieve these **tweets by other users** using the Twitter API.

- a. Use regular expressions to detect PII from the text data in tweets. (for example - emails, phone numbers, aadhaar number etc.) Report the different PII detected in a table (columns: PII type, # of PII detected etc.) with photos of a few examples.
- b. How many of these tweets contain images, videos or other media? Report the PII detected in tweets which contain media (for example - phone numbers, vehicle number plates etc.) and attach screenshots of examples.

Deliverables: {PII Table, Examples, # of Tweets containing media and PII with screenshots}

[10 + 5 + 10 = 25 Marks]

Q2. Using the data collected from the handle and tweets by other users mentioning the handle, analyze the following:

- a. Look at the response time of the handle in replying to tweets by other users. Report minimum, maximum, mean and standard deviation of response time. Can you use this information to infer whether the handle is being controlled by a human or by a bot? Give reasons for your inference.
- b. Make a time series plot for their response times to tweets made by other users. What do you observe? Comment on whether the response times are quicker during certain periods of the day, week or month?

Deliverables: {Response Time Statistics, Inferences, Justifications, Time Series Plot, Comments}

[4 + 3 + 3 + 10 + 5 = 25 Marks]

Q3. Select a subsample ($n = 30$) of the tweets you collected from the users and thematically code them as “report” (reporting a crime), “requests” (assistance etc.), “appreciation” (tweets appreciating the police), “accountability” (holding the police accountable) and “other”. Note that the coding process can be highly subjective and you are free to come up with more codes to categorize tweets mentioning the police account.

- a. Using examples from the subsample, demonstrate how you can model different categories of tweets to provide actionable information that is useful for the police.
- b. Comment on how this can bridge the communication gap between citizens and police and help improve their accountability.

Deliverables: {Table of 30 Tweets with Codes, Modelling, Comments}

[10 + 10 + 10 = 30 Marks]

Section - II (Trending Topics on Twitter)

[55 Marks]

In this section, you are required to select a trending topic from Twitter at the time of data collection. Use the most prominent hashtag of that trending topic and collect tweets containing that hashtag. Now, from these set of tweets, extract the top 10 hashtags other than the most prominent hashtag you used to obtain these tweets. Collect data on these top 10 hashtags similar to how you collected for the first hashtag.

Deliverables: {Code, Data Collected}

[15 + 5 = 20 Marks]

Q1. Plot the top 10 hashtags in your data in terms of (i) most number of occurrences (count multiple occurrences of a hashtag in the same tweet as 1), (ii) most number of likes and (iii) most number of retweets. Read about traffic manipulation on twitter [here](#). Calculate the coefficient of traffic manipulation for each of the three sets of top 10 hashtags and represent them in a table. Can you comment if the data was generated inorganically for some of these hashtags? Do you find some hashtags unrelated to the topic in your data? How do you infer that?

Deliverables: {3 Plots of Top 10 Hashtags, CTM Table, Inferences and Comments}

[5*3 + 5*3 + 5 = 35 Marks]