



Correlation Between Major Crimes Indicators and Twitter Trends

BDA 106: Capstone Final Report

Bangar P., Brzozowski M., &
Khudoyarova O.



Background & Research Project

- Toronto Police Services (TPS) is dedicated to delivering police services, in partnership with their communities, to keep Toronto the best and safest place to be
- With over 500 million public tweets per day, twitter can be a great resource in understanding people's opinions on the events happening in the world.
- This report explores the possibility of correlations between Major Crime Indicators (MCI) and public opinions expressed on Twitter.




Data mining & description

Twitter Trending

- Twitter Trends are based on an algorithm that determines which topics are popular at a particular time and place.
- Trendogate pulls trends from Twitter
 - List of top 50 Trending Topics per day in Toronto
 - Unordered

Trends in Toronto / Canada

 Share

Trends for 2019-07-20

#killjoys

#ElTrafico

#ForTheW

#TheWitcher

#FridayThoughts

#senalg

#TheGiftAlbum

#SDCC

#SignsYourHouseIsHaunted

Data mining & description

Twitter Tweets

- Downloaded from Archive.org
- Semi-structured Data
- 1% of all tweets
- Nested JSON format
- 19.5GB for 1 day

```
e_S0R4","location":"18 177","url":"https://youtu.be/Vn  
http://twitter.com/download/android\" rel=\"nofollow  
,\"location\":null,\"url\":null,\"description\":null,\"translat  
51841,\"id_str\":\"2403051841\",\"name\":\"Arianna\", \"screen_name  
0e22\\u0e42\\u0e14\\u0e21 \\u0e1b\\u0e23\\u0e30\\u0e40\\u0e17\\u  
\":null,\"in_reply_to_user_id\":null,\"in_reply_to_user_id_s  
c! \\u2729\", \"screen_name\":\"alpheccalight\", \"location\":null  
screen_name\":\"Moooooogu0\", \"user\":{\"id\":4706331134,\"id_str  
status_id_str\":\"1150964096247816192\", \"in_reply_to_user_i  
een_name\":\"MrShikharMisra\", \"location\":\"LUCKNOW, INDIA\",  
rita\", \"screen_name\":\"0laaf\", \"location\":\"Rio de Janeiro,  
ser\":{\"id\":317662657,\"id_str\":\"317662657\", \"name\":\"SEFA \\  
22|i|\\u2022 \\u00d8 i+! ].[\", \"url\":null,\"description\":\"i  
str\":\"1151701797238558720\", \"name\":\"\\uce00\\ud78c\\ub9c8\\uc  
\"\\u1d1b\\u029c\\u1d07 s\\u029f\\u1d1c\\u1d1b\\u1d1b\\u028f \\u1d  
593466731163648\", \"name\":\"Mago Negro \\ud83c\\udf41\\ud83d\\u  
\", \"screen_name\":\"aimeeholland_\", \"location\":null,\"url\":\"h  
y_to_user_id_str\":\"22262052\", \"in_reply_to_screen_name\":\"  
tan\", \"user\":{\"id\":1089241232314896384,\"id_str\":\"10892412  
n_name\":\"FatArellano\", \"location\":\"quezon city\", \"url\":nul  
84,\"id_str\":\"897368868884291584\", \"name\":\"Karim Kajo\", \"sc  
tr\":\"1308885715\", \"name\":\"ste\", \"screen_name\":\"hoseokwine'  
40\\u049c\\u0394RI\", \"screen_name\":\"kari01105\", \"location\":\"  
a\\u30b3\\u30cdR https://t.co/BbPUDkfLm5\", \"source\":\"\\u0  
ray_m16\", \"location\":\"Cachoeiras de Macacu, Brasil\", \"url'  
y_to_screen_name\":\"jtcope4\", \"user\":{\"id\":2988286017,\"id_  
\"\\u0618\", \"screen_name\":\"purplemondler\", \"location\":\"Je r\\  
l\", \"in_reply_to_screen_name\":null, \"user\":{\"id\":1041422282  
ser id\":null, \"in reply to user id str\":null, \"in reply to
```

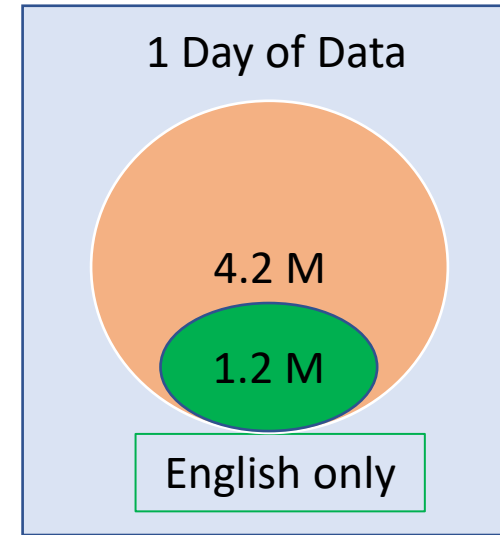
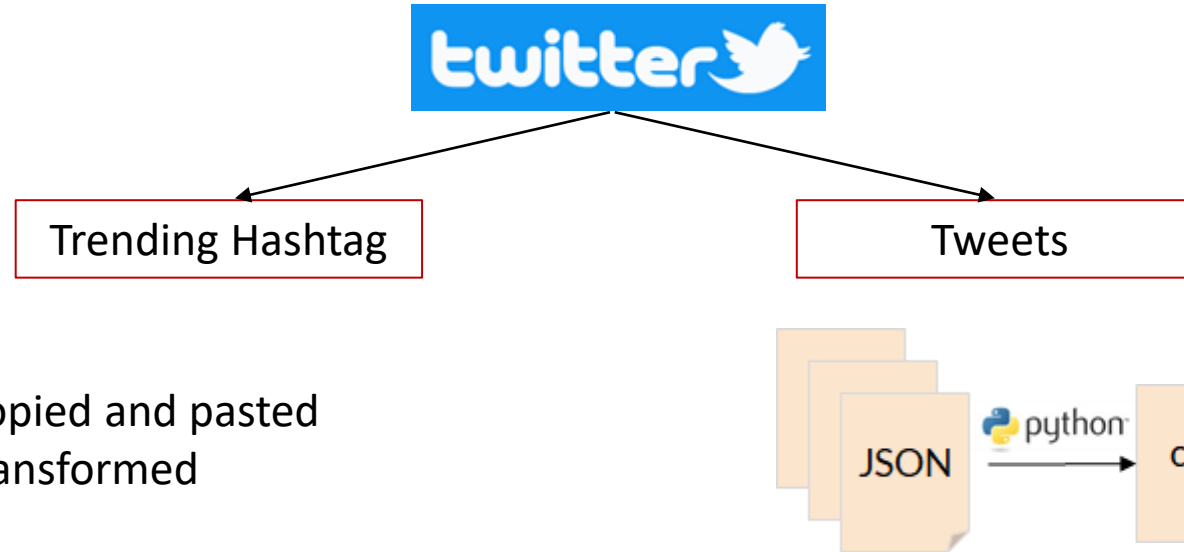
Data mining & description

MCI

- Downloaded from the Toronto Police Services Website
- Structured Data
- 29 columns
- Major Crime details
 - Location
 - Offense
 - Date
- 206,435 rows
- 2014-2019

X	float64
Y	float64
Index_	int64
event_unique_id	object
occurrencedate	object
reporteddate	object
premisetype	object
ucr_code	int64
ucr_ext	int64
offence	object
reportedyear	int64
reportedmonth	object
reportedday	int64
reporteddayofyear	int64
reporteddayofweek	object
reportedhour	int64
occurrenceyear	float64
occurrencemonth	object
occurrenceday	float64
occurrencedayofyear	float64
occurrencedayofweek	object
occurrencehour	int64
MCI	object
Division	object
Hood_ID	int64
Neighbourhood	object
Long	float64
Lat	float64
ObjectId	int64
dtype:	object

Data Cleansing

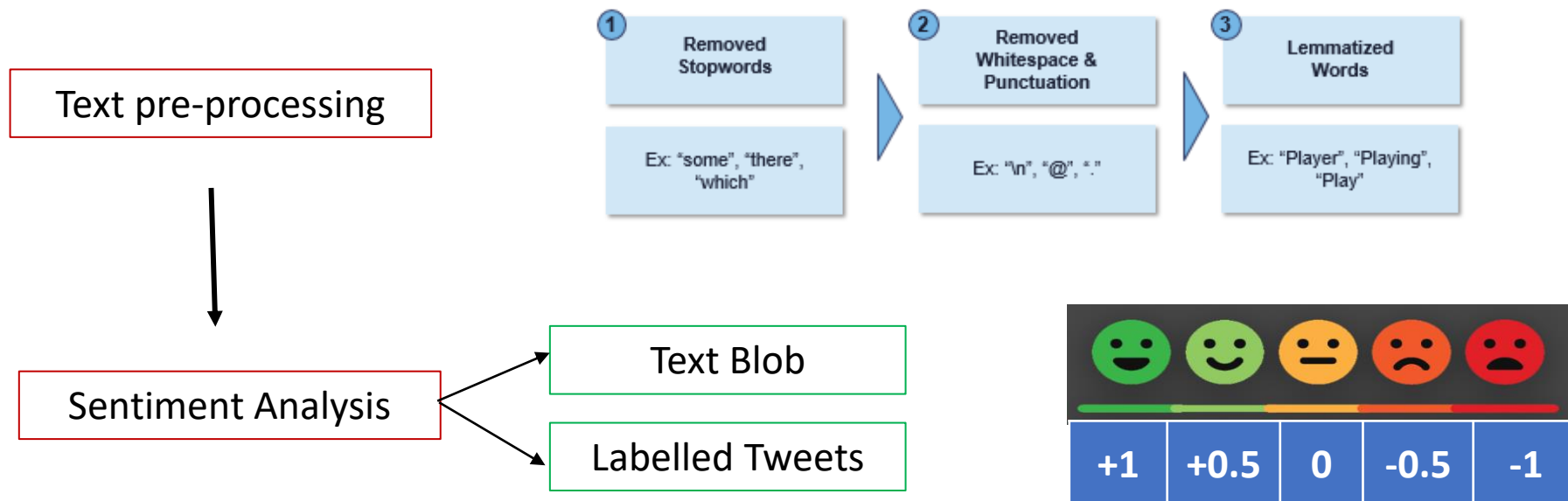


- Manually copied and pasted
- Hashtags transformed

Top 10

A	B	C
Date	Trending Posts (Toronto Only)	Cleaned Trending Post
2019-07-20	#killjoys	kill joys
	#ElTrafico	el trafico
	#ForTheW	for the w
	#TheWitcher	witcher
	#FridayThoughts	friday thoughts
	#senalg	senalg
	#TheGiftAlbum	gift album
	#SDCC	sdcc
	#SignsYourHouseIsHaunted	signs house haunted
	#therealityhouse	reality house
	#QueerEye4	queer eye
	#VeronicaMars	veronica mars
	#FridayFeeling	friday feeling
	#thingsstupidpeoplesay	things stupid people say
	#ParksDay	park day
	#UFCSanAntonio	ufc san antonio
	#SMSociety	sm society

A	B	C
Rank	Hashtag	Body
0	SDCC	RT @ComixVillain: He has it almost right. They have been trying to tell new stories....but much of it is unreadable. I mean, there are pan...
0	SDCC	RT @danielwarrenart: Batman sketch cover at #SDCC https://t.co/GW3leloKLr
0	SDCC	RT @MCU_Direct: Here's the officially-released slate of upcoming Phase Four movies and series set to release in the next two years! #S
0	SDCC	RT @MarvelStudios: Just announced in Hall H at #SDCC, Marvel Studios' THE FALCON AND THE WINTER SOLDIER, an original series
0	SDCC	RT @starwars: 5 things we learned from the Star Wars fashion collaborations panel at #SDCC: https://t.co/gjLBJf18sy https://t.co/9tqD2M0
0	SDCC	RT @SFXmagazine: Black Panther 2
0	SDCC	RT @Fandomopolis: Lena: I don't want to kill Supergirl. I just want to inflict on her the same pain she inflicted on me. #SupergirlSDCC
0	SDCC	RT @MarvelStudios: Just announced in Hall H at #SDCC, Marvel Studios' SHANG-CHI AND THE LEGEND OF THE TEN RINGS, with Sir
0	SDCC	RT @MarvelStudios: Just announced in Hall H at #SDCC, Marvel Studios' BLACK WIDOW with Scarlett Johansson, David Harbour, Florer
0	SDCC	RT @_RyanGajewski: Tessa Thompson on Valkyrie: "First of all, as king, she needs to find her queen." #MarvelSDCC https://t.co/XR1I2ws
0	SDCC	RT @scarlettsgreys: THAT'S MY GIRL! #SDCC2019 https://t.co/n2aCpvtUel
0	SDCC	RT @MarvelStudios: Just announced in Hall H at #SDCC, Marvel Studios' WHAT IF...?, the first animated series in the MCU, with Jeffrey \
0	SDCC	RT @MarvelStudios: Just announced in Hall H at #SDCC, Marvel Studios' THOR: LOVE AND THUNDER with Chris Hemsworth, Tessa Th



	Rank	Hashtag	Body	Sentiment	AVG_Sentiment
0	0	SDCC	RT @ComixVillain: He has it almost right. They...	-1	-0.179081
1	0	SDCC	RT @danielwarrenart: Batman sketch cover at #S...	1	-0.179081
2	0	SDCC	RT @MCU_Direct: Here's is the officially-relea...	1	-0.179081
3	0	SDCC	RT @MarvelStudios: Just announced in Hall H at...	-1	-0.179081
4	0	SDCC	RT @starwars: 5 things we learned from the Sta...	1	-0.179081
...
21099	9	TheGiftAlbum	RT @TheGift_Franco: Y'all spent 7 years callin...	-1	-0.437340
21100	9	TheGiftAlbum	RT @USATODAY: Fans are loving Beyoncé's new so...	1	-0.437340
21101	9	TheGiftAlbum	RT @TheGift_Franco: Y'all spent 7 years callin...	-1	-0.437340
21102	9	TheGiftAlbum	RT @TheGift_Franco: Y'all spent 7 years callin...	-1	-0.437340
21103	9	TheGiftAlbum	RT @Sakpo007: Yessssss!!!!\n\nWe have a um IT...	1	-0.437340

21104 rows × 5 columns



Major Crime Indicators



Location	Date	MCIs	Premise	...



Feature Engineering

Feature Engineering resulted in nine new columns needed for analysis.

Final dataset for analysis:

- 2 months of data (June-July 2019)
- MCI + sentiment analyzed Twitter Tweets
- Important factor – Align them based on date



Data Analysis

- Analysis were carried out using both Microsoft Azure Machine Learning Services and Python Methods.
- Main Three Analysis Performed
 - Simple Two Column Analysis
 - Crime Sum Tier Classification
 - Crimes Stratified by Premise Type



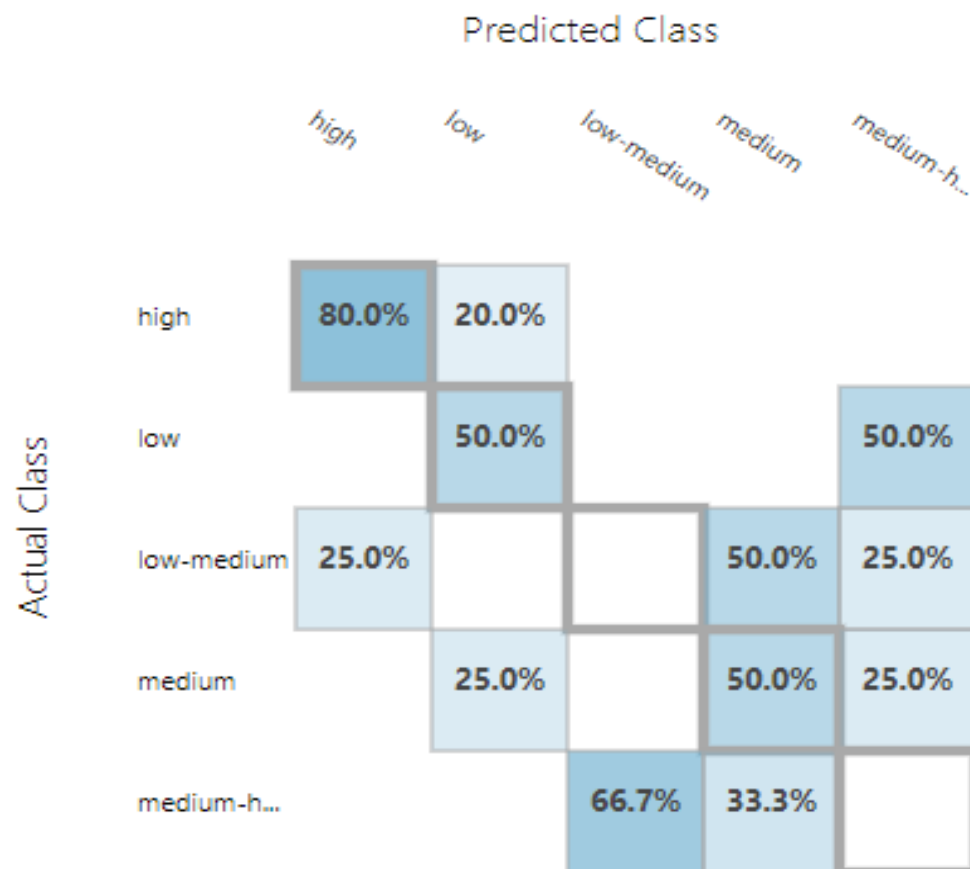
Simple Two Column Analysis

- Two columns considered
 - Total of 61 rows, one row for each day.
- Sentiment dictionary: labelled tweet
- Feature: **AverageHashtagSentiment**
- Target: **Tier_SumofTotalCrimesPerDay**
 - Divided into 5 equal bins depending on the sum of crimes per day
 - High, Medium-High, Medium, Low-Medium, Low
- Multiclass Decision Forest Model
 - Accuracy, Prediction, Recall produced a value of 0.39.

Metrics

Overall accuracy	0.388889
Average accuracy	0.755556
Micro-averaged precision	0.388889
Macro-averaged precision	0.306667
Micro-averaged recall	0.388889
Macro-averaged recall	0.36

Confusion Matrix



Crime Sum Tier Classification

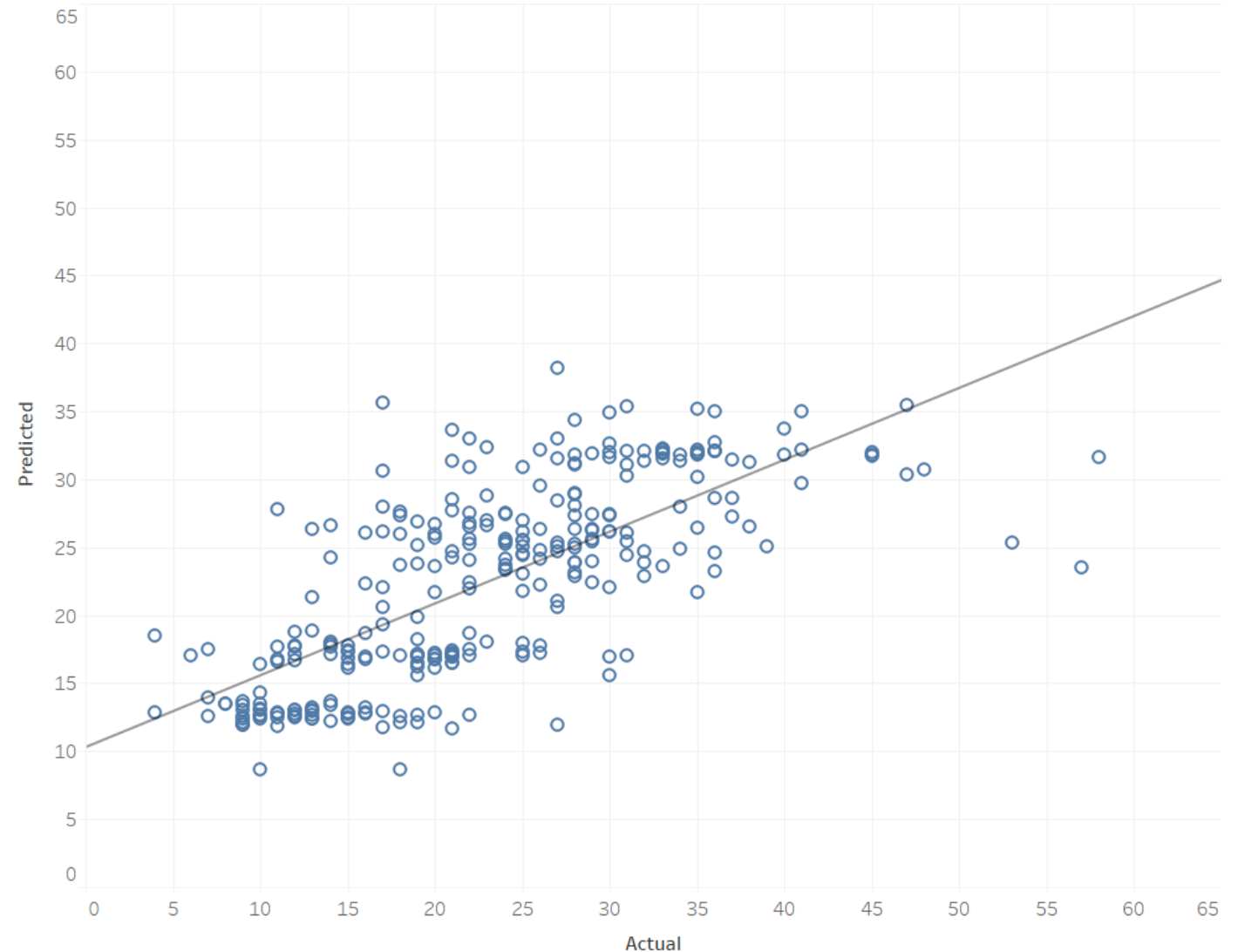
- Twelve columns from considered
 - Total of 7000 rows, one row for each crime.
- Sentiment dictionary: labelled tweet
- Features:
 - **Top 10 Hashtag Sentiment**
 - **MCI labels**
- Target: **Tier_SumofTotalCrimesPerDay**
- Multiclass Neural Network Model
 - Overfit
 - Lack of variance with trending data.

Tier	Average Precision	Average Recall	Average LogLoss
High	1	1	0.0013
Medium-High	1	1	0.0029
Medium	0.8982	0.9876	0.0843
Low-Medium	0.9801	0.9263	0.1032
Low	0.9921	0.9418	0.0829

Regression Analysis of Crimes Stratified by Premise Type

- Features: **Top 10 Hashtag Sentiments, PremiseType**
- Target: **Sum of crimes per premise type per day**
- Neural Network Regression
 - Average coefficients of determination after 10 folds: 0.5
 - Standard deviation of coefficients of determination: 0.14

Predicted vs Actual Number of Crimes per Premise Type



Implications & Recommendations

- We found a minor correlation between Twitter trends and major crime occurrence; however, more data is needed to confirm our findings.
- Expand on the timeframe of data to analyze from two months to a year. Use trending hashtags by the hour if available.
- Assign weight values based on hashtag ranks and sentiment values. Weights can also be assigned to retweets during the sentiment analysis stage to prevent them from overcrowding regular tweets.
- A non-binary sentiment analysis representing emotion on a higher dimension and high-fidelity scale.





Questions?

Special thanks to Guido and Tom