# Building a Tweet Classifier (Binary)

Pang Hong Xiang
18th Mar 2023
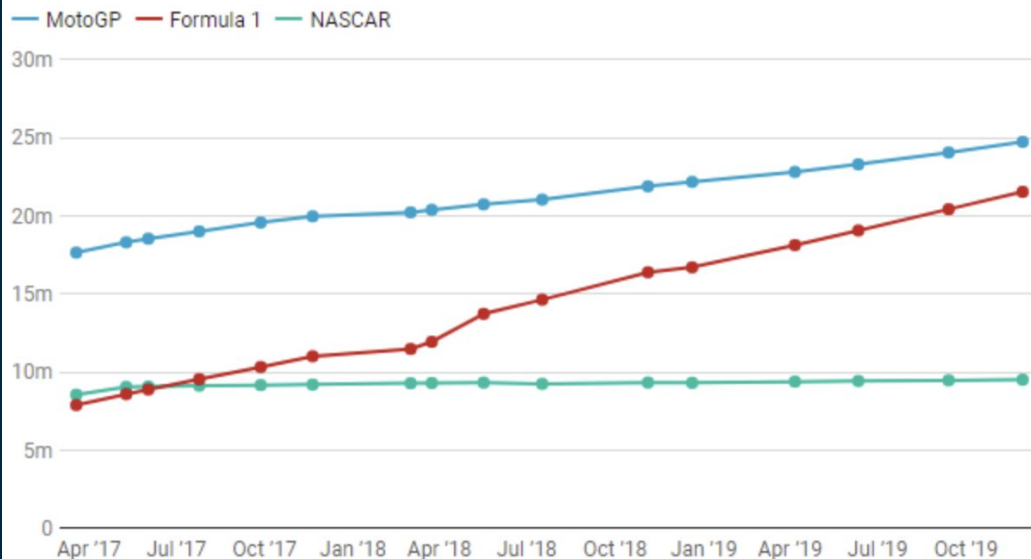
Formula 1's Libe... stock (NASDAQ: FWONA) has also
in... in price by a...argin, returning +62% since *Drive To*
S...remiered...21

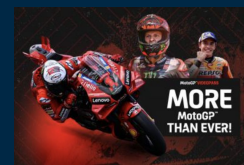US Grand Prix attendance number...
attendance numbers i...

The cherry on top? Social media growth: 49 million followers with 40% growth annually and 1.5 billion aggregate social media engagements – in other words, ... all F1 content is practically viral content

Races now average approximate...
2021 season (Abu Dhabi Grand Prix) beating Supe...
million simultaneous viewers vs 101 million simultaneous viewers

2

MotoGP vs F1 vs NASCAR on social media since 2017

— MotoGP — Formula 1 — NASCAR

Source: Motorsport Broadcasting • Get the data • Created with Datawrapper

# Understanding the Problem
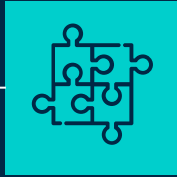
↑ twitter users following topic of Formula 1

↑ auto-racing related accounts and tweets

**Fans are flocking to Twitter to discuss the twists and turns of motorsport's most prestigious competition.**

↑ complicated to differentiate tweet topics

To maximize relevant tweets on news feeds and minimize spam,

*Build a classifier that is able to differentiate Formula 1 related tweets from MotoGP related tweets.*

# Problem–Solving Process

**01**

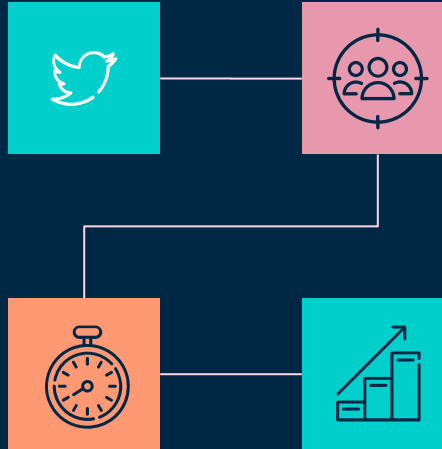Data Collection, Exploration & Processing

**02**

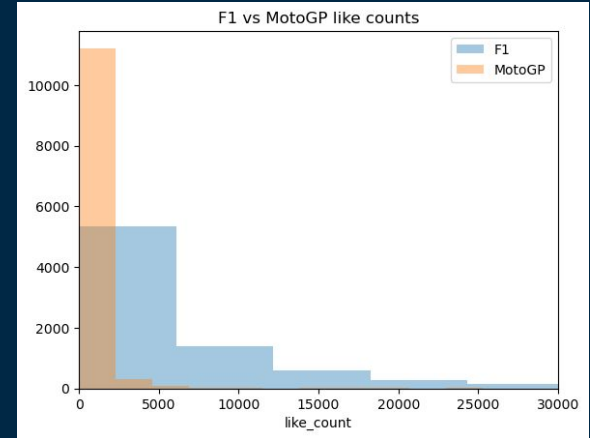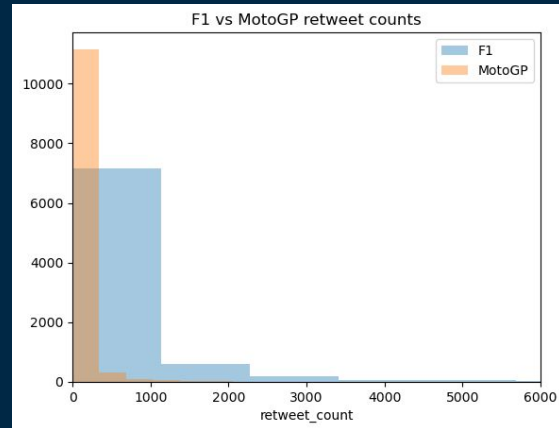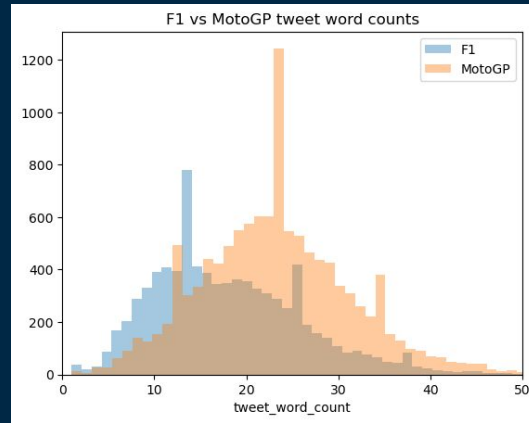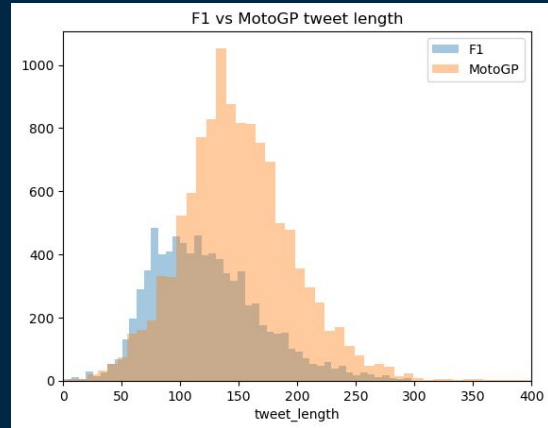Modeling & Evaluation

**03**

Conclusion & Recommendations

# Data Collection

19,850 tweets

@F1 @MotoGP

1st Jan 2022 – 28th Feb 2023

Tweets, likes, retweets

# Data Exploration

# Data Exploration

# Data Exploration

# Data Processing

Create tweet length and tweet word count

Split tweets into text-only, hashtag, mentions

STEP 02

STEP 04

STEP 01

STEP 03

Remove newlines, stopwords, emoticons, quotations, urls

Revisit EDA using processed data to see if anything changed

# Modeling Techniques

## MultiNomial Naive-Bayes

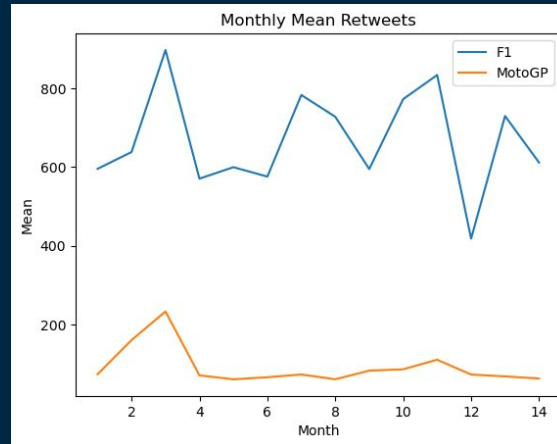Probability of a class, given the occurrence rate of features

## Logistic Regression

Models the relationship between the class and the features

# Modeling Approach

Build a model (M1) using tweet as feature and evaluate the model.

Build a model (M2) using hashtag as feature and evaluate the model.

If hashtag is overly dominant as a predictor, re-evaluate performance of M1 on text-only.

Improve on M1 (M3) by incorporating other features such as tweet length, tweet word count, number of likes, and number of retweets.

Determine best threshold to maximize recall and F1 scores.

# Best Naive-Bayes Model

Using tweet, tweet length, word count, number of likes, number of retweets

**Model (M3, 0.5)**

Tweet: 91.9%
Text-Only: 91.6%

**Accuracy**

Tweet: 87.9%
Text-Only: 87.5%

**Recall**

Tweet: 91.9%
Text-Only: 91.6%

**Precision**

Tweet: 89.9%
Text-Only: 89.5%

**F1 Score**

Tweet: 96.1%
Text-Only: 95.7%

**ROC AUC**

# Best Log-Regression Model

### Model (M3, 0.3)

Using tweet, tweet length, word count, number of likes, number of retweets

### Accuracy

Tweet: 99.4%
Text-Only: 94.0%

### Recall

Tweet: 99.1%
Text-Only: 87.1%

### Precision

Tweet: 99.4%
Text-Only: 98.1%

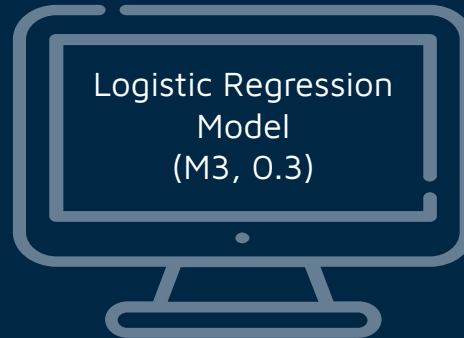### F1 Score

Tweet: 99.2%
Text-Only: 92.3%

### ROC AUC

Tweet: 99.9%
Text-Only: 99.3%

# Recommendation

**To maximize relevant tweets on news feeds and minimize spam,**

Build a classifier that is able to differentiate Formula 1 related tweets from MotoGP related tweets.

Logistic Regression
Model
(M3, 0.3)

# Conclusion

Key Limitation

-   Tweet length, word count, number of likes, number of retweets are all features which could vary greatly among individual users, hence performance may drop further.

Suggested area for improvement

-   Scrape tweets from individual users instead and manually classify for training data.

# THANK YOU