



Predicting Tweet Popularity through a Feature-based Approach

Group 12:
Mojan Majid and Huayin Luo

What type of tweets are receiving the most engagement?

1. Engagement with content
2. **"Fake news", misinformation**

1) Motivation

It is an interesting and important question to determine what type of tweets are receiving the most engagement and popularity, from multiple viewpoints

Motivation - Existing Work

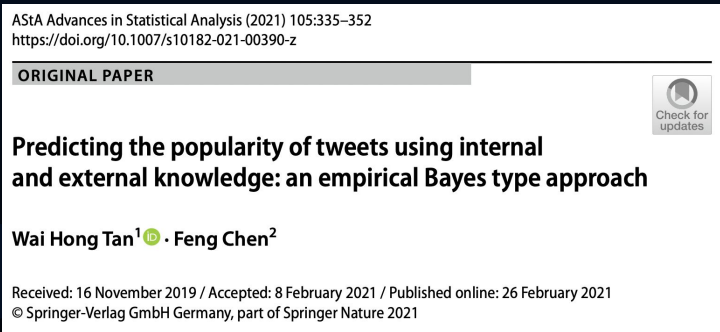
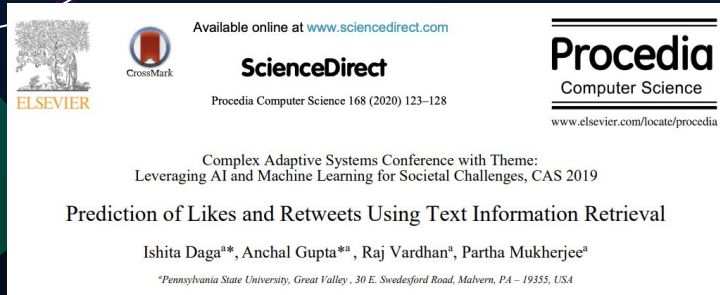
Predicting popularity of user content is a perplexing problem that captivates many.

- **Limited text length**

280 char (*as of Feb 2023, 4000 char for Twitter Blue users)

- **Popularity is complex & ambiguous**

Most practical approaches take into account time-based data



Reward Variables? ✨

1. favourite count
2. retweet count
3. creation time
4. tweet text
5. user screen name
6. user status count
7. user followers count

CSV



✨
Features

#CDNPOLI, no retweets (avoid duplicate tweets)

Limitations: truncated text

2) Data Description

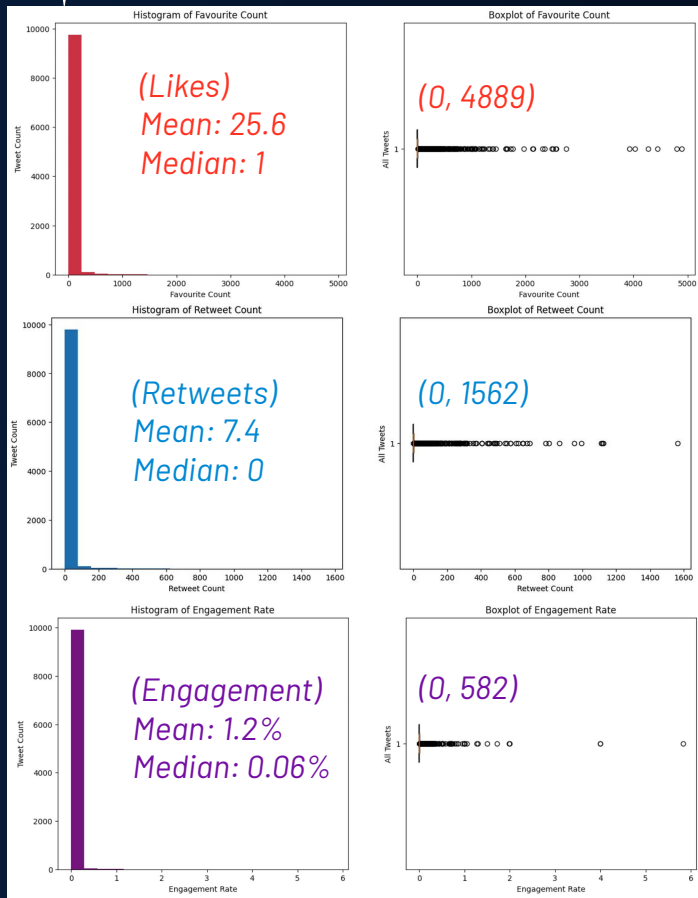
We used Cait's dataset that had 10,000 unique tweets. This data was extracted using the Twitter Developer API. ✨

Data Description - New Variables

We implemented some additional variables based on the existing variables:

Variable Name	Data Type	
Engagement Rate	Continuous (from 0 to 1)	$\frac{100(\# \text{ likes} + \# \text{ retweets})}{\# \text{ followers}}$
Engagement Score	Binary (1 if engagement rate >0.01, 0 otherwise)	
Low Engagement Score	Binary (1 if engagement rate >0.001, 0 otherwise)	
Time of Day Created	Categorical (Morning, Afternoon, Night, Evening) <i>Limitation: Created Time (UTC), we divide based on EST</i>	
Sentiment Score	Continuous (from -1 to 1)	

Features



Limitations:

- (1) Range of values
- (2) Distribution of values

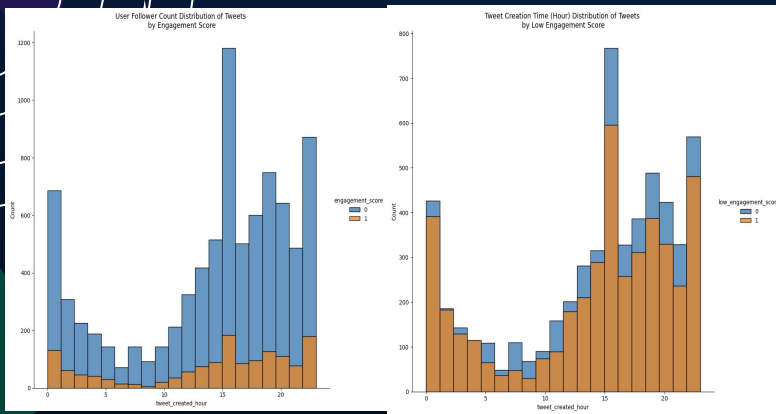
3) EDA

Some exploratory graphs & tables

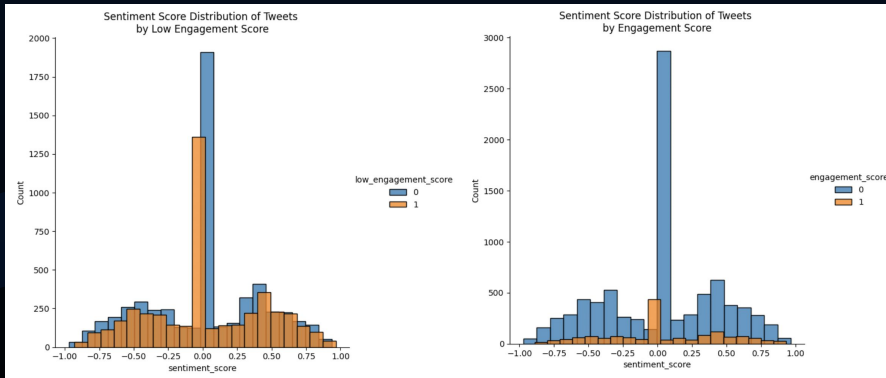
$$\frac{100(\# \text{ likes} + \# \text{ retweets})}{\# \text{ followers}}$$

[Figure 1: Distribution of Favourite Count, Retweet Count, & Engagement Rate]

EDA - More Visualizations



Not Significant: Tweet Creation Time, by Engagement Scores



Significant: Sentiment Score, by Engagement Scores

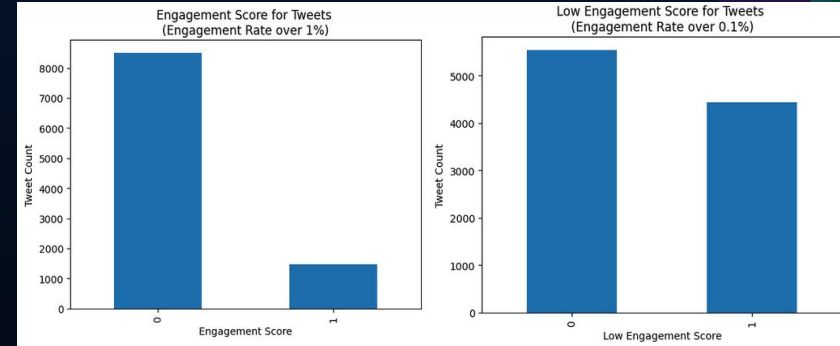


Figure 2: Engagement Score & Low Engagement Score Distribution

We looked at the difference in the distribution of variables, by engagement score, to identify potential features.

1002 Features

1. Sentiment Score for tweet $[-1, 1]$
2. Tweet Text Vectors (1000 words)

+

1. Username length

=

Outcome

Engagement Score
Low Engagement Score

NLTK VADER
Pretrained Sentiment Model

Preprocessing + TF-IDF

4) Our Model

We used a two step approach, with
Text Features for Non-Text Prediction

Logistic Regression

Binary: 0 or 1

Engagement Score

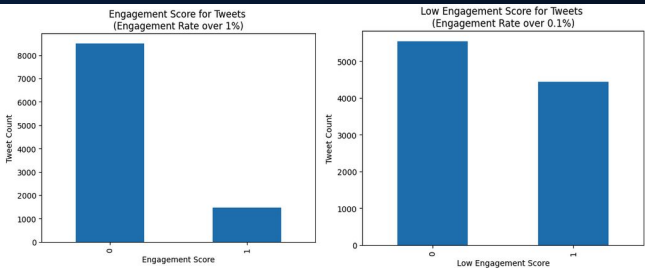
Train Accuracy: 84.7%

Test Accuracy: 85.7%

Low Engagement Score

Train Accuracy: 67.4%

Test Accuracy: 67.0%



(Recall)

5) Results

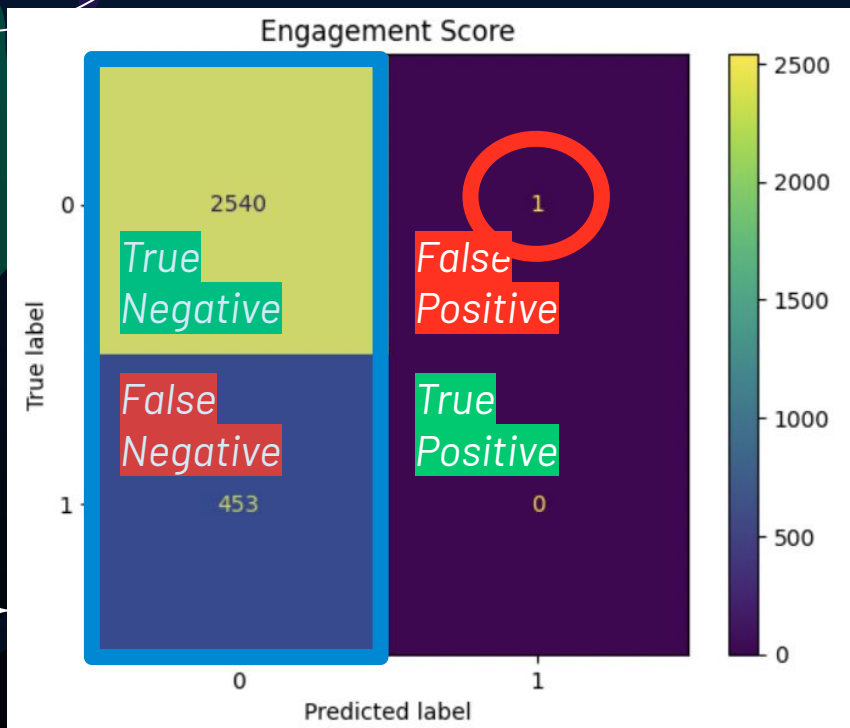
Engagement Score: model performs similarly to a "naive" model.

Low Engagement Score: performs better than both random & naive baseline models

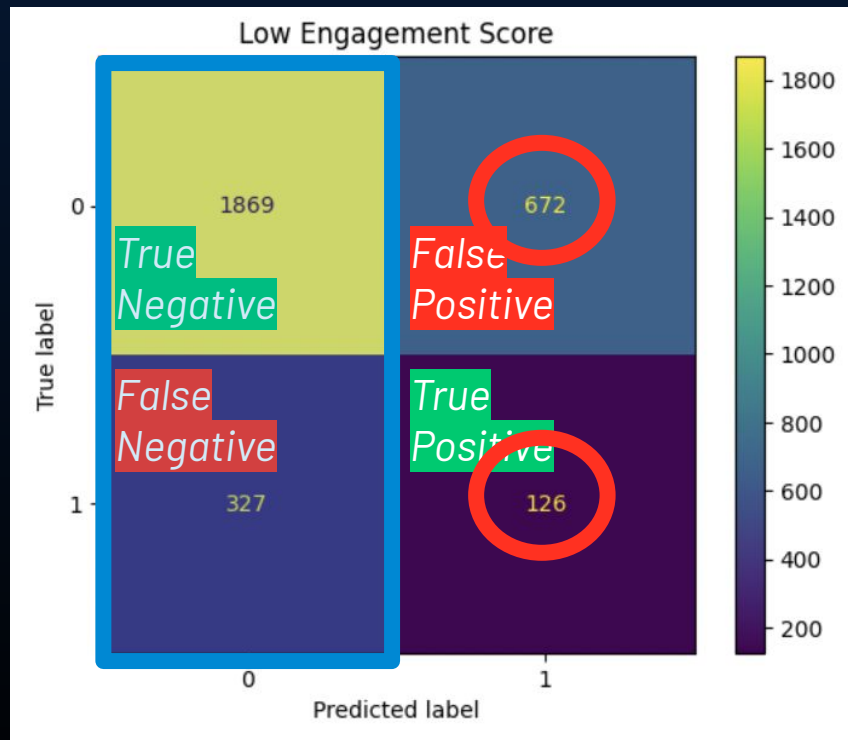
- *Difference = Balanced dataset?*

"Naive" approach?

Results - Conclusions & Next Steps



FNR:0.15, FPR: 1.0



FNR:0.15, FPR: 0.84

Better! Trying.

Results - Conclusions & Next Steps

Our model shows an interesting starting point for combining word features from NLP approaches with other informative variables in predicting popularity and engagement of a tweet. There are many ways in which we can improve:

- (1) Full tweet text (instead of truncated)
 - *More insightful NLP vectors**
- (2) More balanced dataset
 - *balanced amount of popular/not popular tweets*
 - *more very popular/viral tweets**
- (3) Other features...*

Thank you!

