# AML-2203 Advanced Python AI and ML tools

# **Final Project**



Submitted To:

Vahid Hadavi

# Team:

Mohammed Muzammil Lakdawala (C0872315)

Gurinder Singh (C0871566)

Maryam Shokrollahi (C0864597)

Mohammed Fuwad Anjum (C0867421)

Gursimran Kaur (C0867039)

#### Introduction

In recent years, social media has become a prominent platform for people to express their thoughts and opinions, and it has also become a vital source of information for businesses and individuals alike. Twitter, in particular, has gained immense popularity as a platform for news, politics, and entertainment. With over 330 million monthly active users, Twitter has become a valuable source of data for researchers and data scientists.

In this project, we aimed to predict the number of likes that a tweet from Elon Musk's Twitter account would receive. Elon Musk is a well-known entrepreneur, investor, and visionary who is known for his innovative ideas and outspoken opinions. His tweets often go viral, and we wanted to explore the factors that contribute to their popularity.

To achieve this goal, we collected data on Elon Musk's tweets and their corresponding number of likes using SNScrape. We then preprocessed the data and extracted relevant features such as the number of retweets, the length of the tweet, and the time of day it was posted. We used machine learning techniques such as regression analysis to build a predictive model to estimate the number of likes a tweet would receive.

The insights gained from this project could be useful for businesses and individuals who want to improve their social media strategy and increase their engagement on Twitter. By understanding the factors that contribute to tweet popularity, they can create more effective content and ultimately reach a wider audience.

In the following sections, we will discuss our data collection process, data preprocessing techniques, feature engineering, and model building process in detail. We will also present our results and discuss the implications of our findings.

#### Methodology

After scraping the data, it looks something like this

df.head() ✓ 0.0s									
	Unnamed: 0	User	Date Created	Number of Likes	Source of Tweet	Tweet	ReTweet Count	Reply Count	View Count
0		elonmusk	2023-04-06 20:31:24+00:00	1418	Twitter for iPhone	@DirtyTesLa We are gradually reducing it, prop	126	331	206454.0
1		elonmusk	2023-04-06 20:26:19+00:00	1174	Twitter for iPhone	@GRDecter Any highly successful investor with	84	142	59256.0
2		elonmusk	2023-04-06 20:17:48+00:00	1947	Twitter for iPhone	@GRDecter Probably Buffett. He could do it usi		161	72716.0
3		elonmusk	2023-04-06 20:15:32+00:00	1604	Twitter for iPhone	@cb_doge Actually, will add this app to my nam	114	283	78871.0
4	4	elonmusk	2023-04-06 20:13:36+00:00	9968	Twitter for iPhone	@KanekoaTheGreat Is this real?	1260	1524	1290796.0

We first prepare the data by cleaning and preprocessing the text and then made a function that creates a Bag of Words (BOW) representation of the data.

```
def toBOW(sentence, words):
    bag = []
    for word in words:
        bag.append(1) if word in sentence else bag.append(0)
    return bag

    0.0s
```

Now we add a new column to the DataFrame called "Engagement" that is the sum of the "Number of Likes" and "ReTweet Count" columns.

```
df['Engagement'] = df['Number of Likes'] + df['ReTweet Count']

$\square$ 0.0s
```

We then extract the distinct words from the tweets using the prepareSentence function, which applies stemming, removes stop words, and removes non-alphabetic characters.

At this point we have too many words (1,67,144) to work with so we need to reduce it. So, we then set lower and upper threshold values for word counts and only include words that fall within this range.

```
distinct words = set(words)
   lower_threshold = 10
   upper threshold = 350
   counts = []
   final_words = []
   for word in distinct words:
       counts.append(words.count(word))
       if words.count(word) > lower threshold and words.count(word) < upper threshold:
           final_words.append(word)
   print(len(words))
   print(len(distinct words))
   print(len(final words))
 √ 1m 8.3s
167144
19303
2124
```

The Bag of Words (BOW) representation is created by converting each tweet into a binary vector representing the presence or absence of the selected words.

We then split the data into inputs and outputs for training the neural network. Inputs are the BOW vectors and outputs are scores for each tweet that indicate whether it received any engagement (1) or not (0).

```
inputs = []
outputs = []

for index, tweet in df.iterrows():
    sentence = prepareSentence(tweet['Tweet'])
    bag = toBOW(sentence, final_words)
    inputs.append(bag)
    #Calculate a score, 1 if any engagement, 0 if none
    score = min(tweet["Engagement"],1)

    outputs.append(score)

    14.7s
```

The neural network (MLPRegressor) is then defined and trained using the training data with some of the following hyperparameters.

```
# Define and training the network
nnet = MLPRegressor(activation='relu', alpha=0.0001, hidden_layer_sizes=(int(len(final_words)*0.5),int(len(final_words)*0.25)),solver='adam', max_iter=400)
nnet.fit(inputs, outputs)

/ 4m 11.7s

Python

MLPRegressor
MLPRegressor
MLPRegressor(hidden_layer_sizes=(1062, 531), max_iter=400)
```

Finally, the trained neural network is used to predict the engagement of potential tweets in a list called feed. The potential\_posts dictionary is used to store the predicted engagement score for each post, and the posts are sorted by their predicted engagement score to identify the top and bottom five posts.

```
Top 5 posts
("Cybertruck is going to be a game-changer in the world of electric trucks. Can't wait to start production!", 1.0163000704977478)
("Who needs a boring old house when you can live in a sustainable, high-tech Mars colony! #ColonizeMarsThe Falcon Heavy launch was absolutely epic. Can't wait
("The Hyperloop is going to change the way we travel forever. Who's ready for a 700 mph ride?", 1.0081185127648657)
("Boring Company flamethrowers are the ultimate tool for clearing brush and entertaining guests. #FlameOn', 1.0042685905810447)
("Just landed on Mars. Time to start building a new world!', 1.0040125732489746)
Bottom 5 posts
("Sustainable energy is the only way forward. We need to invest in clean technology now more than ever.', 0.980505997344562)
("Spacex's Starlink internet service is going to connect the entire world. No more dead zones!", 0.9806338733758235)
("Life is too short to waste time on things that don't matter. Follow your dreams and make a difference.", 0.9876500442131579)
("Tf you're not pushing the limits of what's possible, you're not really living.", 0.9904028318852714)
("Neuralink is going to revolutionize the way we interact with technology. The possibilities are endless.', 0.991656579321063)
```

This was a success but to truly call this a success we need evidence and to provide it we calculated the errors for it. We calculated the Mean Absolute Error (MAE), Root Mean Square Error (RMSE) and Mean Square Error (MSE) and all of them are very low.

```
print(f'MME train: {metrics.mean_absolute_error(y_train, y_train_pred)}, test: {metrics.mean_absolute_error(y_test, y_test_pred)}')
print(f'RMSE train: {metrics.mean_squared_error(y_train, y_train_pred)}}, test: {metrics.mean_squared_error(y_test, y_test_pred)}')
print(f'RMSE train: {(mean_squared_error(y_train, y_train_pred))}, test: {(mean_squared_error(y_test, y_test_pred))}')

> 0.0s

Python

MAE train: 0.01191221498055484, test: 0.025335159499178834

MSE train: 0.0030998510092899604, test: 0.0005539037326517877
```

But this type of prediction is only useful to us when we have a list of tweets and we want to know which one will perform better than the others and it won't have the exact number of likes that tweet is going to receive. So, we did an alternative approach as well.

We first performed some preprocessing on the tweet column like removing the numbers, punctuation and small words as follows

We also performed sentiment analysis on it using SentimentIntensityAnalyzer which gives us 4 scores which are compound, negative, neutral and positive. We will be using this for training our model as well to get some more features.

Then, we used TFIDF Vectorizer on the tweet column, we set the max\_features as 20,000 as we have a lot of tweets and going lower might cause us to lose a lot of our features.

We then used Truncated SVD to reduce the number of dimensions but we have to be careful as going too low might cause us to lose variance and a lot of information along with it.

```
# Using Truncated SVD to reduce the dimensions of the data
number_of_dimensions = 10000
svd = TruncatedSVD(n_components=number_of_dimensions)
reduced_data = svd.fit_transform(vectorized_dataframe)
print(reduced_data.shape)

# List of explained variances
tsvd_var_ratios = svd.explained_variance_ratio_

var_explained = svd.explained_variance_ratio_.sum()
print(var_explained)

✓ 25m 5.3s

(21601, 10000)
0.9661877116226325
```

We went with 10,000 features as this gave us 0.966 and when we tried with 5,000 rows, we were getting a variance of around 0.6 which was very low.

We then divide the DataFrame into target and features as follows

```
df_copy = df.copy(deep=True)
y = df_copy['Number of Likes']
x = df_copy[['neg', 'pos', 'compound', 'neu']]

$\square 0.1s

x = pd.concat([x, vectorized_dataframe], axis=1, join="inner")
$\square 7.8s$
```

Here, x are the features and y is the target.

We wanted to perform Polynomial Feature Extraction on it but due to hardware restrictions we were unable to do so as only a degree of 2 was requiring memory of 22 Terabytes.

So, we moved on to training our Regression model.

We trained our Random Forest Regressor on the train set.

The metrics that we got for it are as follows

```
print(f'MAE train: {metrics.mean_absolute_error(y_train, y_train_pred)}, test: {metrics.mean_absolute_error(y_test, y_test_pred)}')
print(f'RMSE train: {(np.sqrt(metrics.mean_squared_error(y_train, y_train_pred))}, test: {(np.sqrt(metrics.mean_squared_error(y_test, y_test_pred))}')
print(f'RMSE train: {(np.sqrt(metrics.mean_squared_error(y_test, y_test_pred))}')
print(f'MSE train: {(mean_squared_error(y_train, y_train_pred))}, test: {(mean_squared_error(y_test, y_test_pred))}')

v 0.0s

MAE train: 51169.32994514804, test: 50413.39415588119
RMSE train: 110654.17187098303, test: 107664.16521526489
R^2 train: 0.1508096085523517, test: 0.001510765004574921
MSE train: 12244345752.453053, test: 11591572471.499855
```

The R<sup>2</sup> score was only 0.15 for the train and 0.0015 for the test and we thought it was because of the model that we used so we trained some more model.

We started off with the same MLP Regression neural network as before.

```
nnet = MLPRegressor(activation='relu', alpha=0.0001, hidden_layer_sizes=(int(len(x_train)*0.5),int(len(x_train)*0.25)),solver='adam', max_iter=400)

nnet.fit(x_train, y_train)

46m 132s

Python

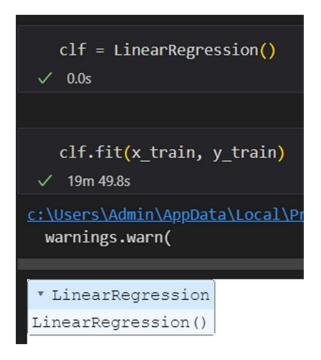
62i\User=\ladam\interpolata\local\Programs\Python\Python310\lib\site_packages\sklearn\utils\validation.py:1858: FutureWarning: Feature names only support names that are all strings. Got feat

warnings.warn(

1.\Users\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\under\un
```

But we got even worse results with it as the R<sup>2</sup> score was even lower at 0.042 for training and 0.044 for the test. Although here the score was higher for the test but it was still extremely low.

At this point we understood that the complexity of the model was way too much for the data we were training it for so, we tried using a simpler model like Linear Regression.



And the R<sup>2</sup> score for the train set came out as 0.796 which was an extreme delight to us as we thought that the model finally learned something but at the same time the score for the test set came out as -6.58 which meant that the model became severely overfit.

```
print(f'MAE train: {metrics.mean_absolute_error(y_train, y_train_pred)}, test: {metrics.mean_absolute_error(y_test, y_test_pred)}')
print(f'RMSE train: {np.sqrt(metrics.mean_squared_error(y_train, y_train_pred))}, test: {np.sqrt(metrics.mean_squared_error(y_test, y_test_pred))}')
print(f'RMSE train: {(r2_score(y_train, y_train_pred))}, test: {(r2_score(y_test, y_test_pred))}')
print(f'MSE train: {(mean_squared_error(y_train, y_train_pred))}, test: {(mean_squared_error(y_test, y_test_pred))}')

✓ 00s

MAE train: 24700.946257073476, test: 6.2021337090668744e+16
RMSE train: 54221.64044023593, test: 2.7640308401127754e+17
R^2 train: 0.7961011640755955, test: -6.580922873859397e+24
MSE train: 2939986292.0302286, test: 7.639866485094535e+34
```

Adding something of only slight relevance, we even tried using XGBoost but it came as a failure as well, as it was running for more than 3 hours and did not finish training.

```
clf = GradientBoostingClassifier(n_estimators=100, learning_rate=1.0, max_depth=1, random_state=0).fit(x_train, y_train)
clf.score(x_test, y_test)
@ 205m 51.7s
```

## Conclusion

In conclusion, we were successful in scoring which tweets are more likely to perform better with respect to other tweets but were unsuccessful when it came to predicting their likes with concrete accuracy.

Some of the reasons for the low R<sup>2</sup> score can be presence of outliers in the data, too less features considered while performing vectorization and lack of Polynomial Features while training the model.

But overall this has been a good learning experience for us as to where to look for when dealing with such underfit models and understanding that sometimes lack of hardware might also be an issue.

## References

Yener, Y. (2020, November 7). *Step by step: Twitter sentiment analysis in Python*. Medium. Retrieved April 20, 2023, from https://towardsdatascience.com/step-by-step-twitter-sentiment-analysis-in-python-d6f650ade58d

Pedregosa, F., Varoquaux, Ga"el, Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... others. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12(Oct), 2825–2830.

Koehrsen, W. (2018, January 28). Overfitting vs. underfitting: A complete example. Medium. Retrieved April 20, 2023, from https://towardsdatascience.com/overfitting-vs-underfitting-a-complete-example-d05dd7e19765