DSC 681-002 - Fall 2024: Applied Machine Learning

Group Project Hackathon

Professor Yue Han

Group 1:

Alma Monreal

Minh Khue Nguyen

Riya Shah

Shakif Farhan Shah

Table of Contents

Project goal and objective	3
Data Description	4
Analyses and Results	8
Discussion	9
Code	12

Project goal and objective

The goal of this project is to develop a model to predict the popularity of a tweet by cleaning, comparing and analyzing the data provided. The model will be developed using the general machine learning process.

The objective of this project and of creating the model is to be able to process data provided to better understand the most relevant factors that can make a tweet go viral.

The data provided contains items related to tweet information and user profile information.

Data Description

B1. What is the dataset for your project? If there is a link to your dataset, please provide the link. If not, please provide the first 5 rows including the headers of your dataset:

There are three datasets provided. 1) tweet_info 2) tweet_more_info and 3) user_profile.

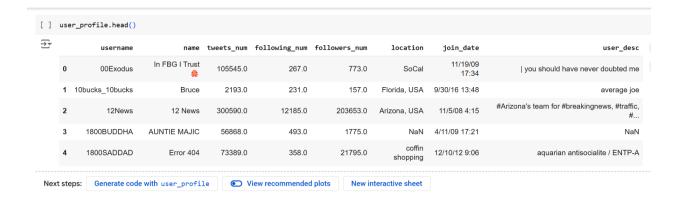
1) tweet_info



2) tweet_more_info

tweet_more_info.head() \rightarrow tweet_id character hashtag mention term 2 993672051763351552

3) user_profile



B2. What is the dependent/output variable? What are the independent/input variables? Why do you choose these variables?

Dependent/Output Variable: The number of retweets (rt_num). This variable measures the popularity of a tweet.

Independent/Input Variables:

- hashtag: Number of hashtags in the tweet.
- img_num: Number of images included in the tweet.
- has_vid: Whether the tweet included a video.
- mention: Number of mentions in the tweet.
- followers_num: Number of followers of the tweet's creator.
- tweets_num: Total number of tweets posted by the user.

We chose these variables because the number of retweets is the dependent variable that we want to know. Which depends on other factors which are the independent variables mentioned above. Those independent variables are important things to consider towards the popularity of a tweet and therefore it helps to analyze what makes a tweet go viral in order to predict future outcomes of the popularity of a tweet or to recreate those same characteristics in future tweets to make them go viral.

B3. Please provide the statistics table of your cleaned data. Below is an example, yours can be a little different.

Variable	Number of Observatio	Min	Max	Mean	Standard deviation	
Retweets (Dependent variable)	10000	0	21536	853.83 7	7046.603497	
Hashtags (independent variable)	10000	0	7	0.0746	0.358117	

Followers number (independent variable)	1.000000e+0 4	0	8.8266 66e+0 7	1.7012 20e+06	8.553986e+06	
Fav Number (independent variable)	10000	0	55994 8	2207.5 9	14517.38	
Tweets Number (independent variable)	10000	0	51224	29302. 4	48569.86	

Analyses and Results

C1. What models are you using for this project and why?

We use a support vector machine with regression to help us predict the number of retweets the tweet will get. SVR is effective for predicting continuous numerical values, such as the number of retweets. It is also a great tool for handling data and creating predictions even with a limited sample size.

C2. What are the results? Which model works better?

Due to the numerical nature of rt_num, it is more suitable for a regression model like SVR, designed for predicting continuous outputs.

If the dependent variable were categorical, classification models would be more appropriate. For this project, predicting tweet popularity relies on numerical engagement metrics, justifying the use of regression.

Discussion

D1. What are the implications of your analyses and results?

The implications of our analyses include:

Key Factors of Popularity: The analysis identifies that features like the number of followers, favorites, hashtags, and total tweets by the user have significant influence on the number of retweets.

Predictive Power: The SVR model has the ability to predict the popularity of tweets, providing insights for marketers and influencers to optimize their content for the best engagement.

Broader Applications: These results can guide businesses in creating social media strategies and content by focusing on measurable attributes like hashtags, mentions, and visuals to increase their engagement and growth.

D2. What are some potential ways to improve your models in the future?

Add More Features: Include new data points like sentiment analysis or keywords in tweets.

Use Better Models: Try other models like Random Forest or Neural Networks and compare performance.

Fine-Tune Parameters: Optimize settings like SVR's kernel or regularization.

Expand the Dataset: Use more tweets or create data for training.

Engineer Features: Create new variables like retweets per follower or normalized hashtag counts.

Avoid Overfitting: Use techniques like cross-validation and regularization to make the model generalize better.

D3. What are your main takeaways of this project (or any of the projects for this course)?

Main takeaways from this course are the process of building and evaluating a machine learning model, from data cleaning to model interpretation and presenting them in an insightful manner. Learning the strengths and limits of Support Vector Regression, especially in predicting numerical outcomes like retweets.

Also, to have the ability for critical thinking. Developing a mindset for questioning, going the extra mile and improving models to achieve better accuracy.

Lastly, gaining insights into how data science can be applied to real-world problems like optimizing content of various business entities for their social media platforms.

D4. What are your suggestions/feedback on how to improve the class experience of this course?

It is a well designed course with a lot of useful information, but I feel that 2 months is too short of a time frame to learn the basics applied machine learning in depth. Would be better if it was a 4 month course.

Code

HACKATHON.ipynb

```
# Step 1: Setting working directory and importing libraries
from google.colab import drive
drive.mount('/content/drive')
import os
os.chdir('/content/drive/MyDrive/Colab Notebooks/Applied Machine
Learning/DSC 681 Hackthon Project Data')
import pandas as pd
from sklearn.model selection import train test split
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVR
from sklearn.metrics import mean absolute error, mean squared error,
r2 score
import seaborn as sb
# Step 2: Load datasets
tweet_info = pd.read_csv('tweet_info.csv')
tweet more info = pd.read csv('tweet more info.txt', delimiter='\t') #
Assuming tab-separated file
```

```
user profile = pd.read csv('user profile.csv')
tweet info.head()
tweet More info.head()
user profile.head()
# Step 3: Check for missing values
print(f"tweet info missing values: \n{tweet info.isnull().sum()}")
print(f"\ntweet more info missing values:
\n{tweet more info.isnull().sum()}")
print(f"\nuser_profile missing values: \n{user_profile.isnull().sum()}")
# Step 4: Remove missing values and duplicates
tweet info = tweet info.drop(columns=['text', 'quote']) # Drop rows where
'text' or 'quote' columns have NaN values
user profile = user profile.drop(columns=['name', 'location', 'join date',
'user desc']) # Drop rows where specified columns have NaN values
# Step 5: Show the first few rows of each dataset
print(f"tweet info: \n{tweet info.head()}")
print(f"\ntweet more info: \n{tweet more info.head()}")
print(f"\nuser profile: \n{user profile.head()}")
# Step 6: Check for duplicates in the DataFrames
print(tweet info.duplicated().sum())
print(tweet more info.duplicated().sum())
print(user profile.duplicated().sum())
# Step 7: Consolidate the data
merged data = tweet info.merge(tweet more info, on='tweet id', how='left')
```

```
merged data = merged data.merge(user profile, on='username', how='left')
merged data.head()
# Step 8: Fill missing values and check merged data
print(merged data.isnull().sum())
# Step 9: Sampling the data
data = merged data[['followers num', 'rt num', 'fav num', 'tweets num',
'hashtag']]
sampledata = data.sample(n=10000, replace=False)
print(f"sampledata:\n{sampledata.info()}")
print(f"\nsampledatadescribe: \n{sampledata.describe(include='all')}")
# Step 10: Split the data into dependent (y) and independent (X) variables
y = sampledata['rt num']
# Safely select only the available columns in merged data
X = sampledata[['followers_num', 'fav_num', 'tweets_num', 'hashtag']]
# Step 11: Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random state=42)
print(f" X_train: \n{X_train}")
print(f"\ny train: \n{y train}")
# Step 12: Correlation analysis
print(sampledata['rt num'].corr(sampledata['followers num']))
samplecor = sampledata.corr(method='pearson')
```

```
print(samplecor)
#Step 13: Producing correlation through heatmap.
sb.heatmap(samplecor,
           xticklabels=samplecor.columns,
           yticklabels=samplecor.columns,
           cmap='RdBu r',
           annot=True,
           linewidth=0.5)
# Step 14: Standardize the features
scaler = StandardScaler()
X train = scaler.fit transform(X train)
X test = scaler.transform(X test)
# Step 15: Train the SVM model
model = SVR(kernel='linear')
model.fit(X_train, y_train)
# Step 16: Make predictions
y_pred = model.predict(X_test)
y pred
# Step 17: Evaluate the model
mae = mean absolute error(y test, y pred)
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print("\nModel Evaluation:")
```

```
print(f"Mean Absolute Error: {mae}")
print(f"Mean Squared Error: {mse}")
print(f"R2 Score: {r2}")
```