

10 Academy

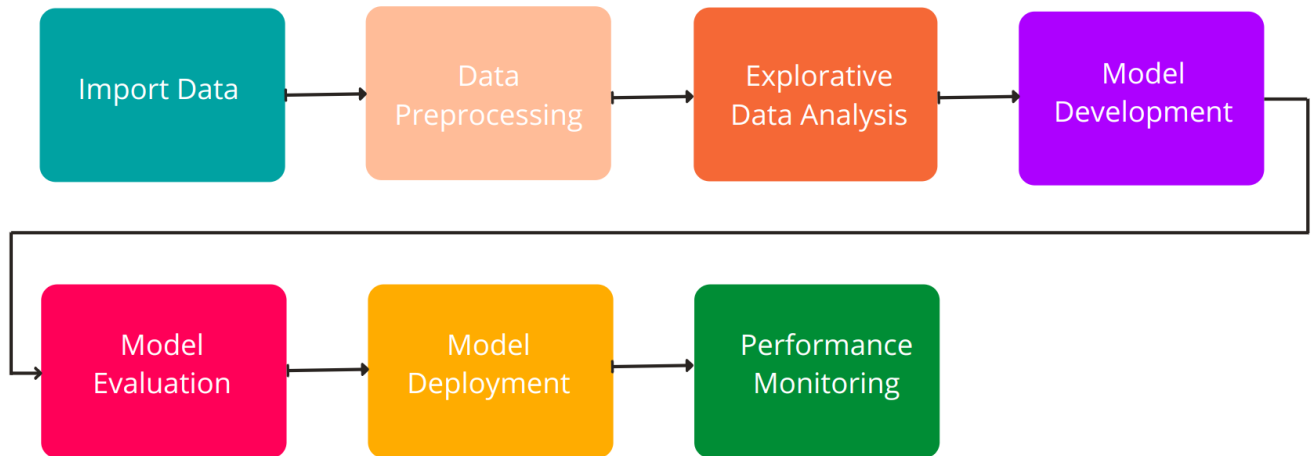
Week 0

Report 3 - Interim Report

Name: Yishak Tadele

ML Workflow Chart and Description

Diagram



In the following section, I will explain the work I have done and the new insights I obtained.

1. Data Importing

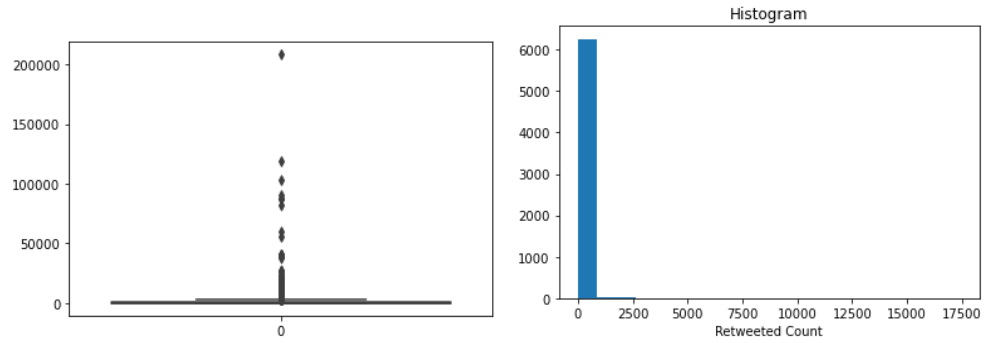
- After loading the JSON dump into the pandas' data frame I have saved it to the pickle file to make it available for later access.
- The pickle file is directly loaded and further preprocessed in the next stages.

2. Data Preprocessing

- In this stage, I have cleaned the dataset further, especially the Twitter content extracted further. For example
 - Stop words, usernames and links and non-English characters are removed.
 - Lemmatization and vectorized representation is applied to the text feature.

3. Data Exploration /EDA

Plots: most of the features with numerical values have the following distribution in a boxplot. These features are followers_count, friends_count, retweet_count and favorite_count. The distribution reflects that a small number of users get a lot of likes and retweets and most of the users account likability is below average. The same argument can be built for the histogram plots which can be shown.



Analytical Results:

```
1 df.describe()
```

| | polarity | subjectivity | favorite_count | retweet_count | followers_count | friends_count |
|-------|-------------|--------------|----------------|---------------|-----------------|---------------|
| count | 6326.000000 | 6326.000000 | 6326.000000 | 6326.000000 | 6.326000e+03 | 6326.000000 |
| mean | 0.062901 | 0.313011 | 0.399621 | 45.417325 | 1.016427e+04 | 1543.008536 |
| std | 0.229075 | 0.279022 | 1.534209 | 353.598765 | 2.162880e+05 | 4968.226632 |
| min | -1.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000e+00 | 0.000000 |
| 25% | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 5.500000e+01 | 95.250000 |
| 50% | 0.000000 | 0.300000 | 0.000000 | 0.000000 | 2.960000e+02 | 390.500000 |
| 75% | 0.150000 | 0.500000 | 0.000000 | 4.000000 | 1.273000e+03 | 1427.500000 |
| max | 1.000000 | 1.000000 | 59.000000 | 17409.000000 | 1.024910e+07 | 208360.000000 |

- The following picture shows the range, average and different quartiles for the numerical features.
- The possibly_sensitive feature shows missing values for 3286 entries

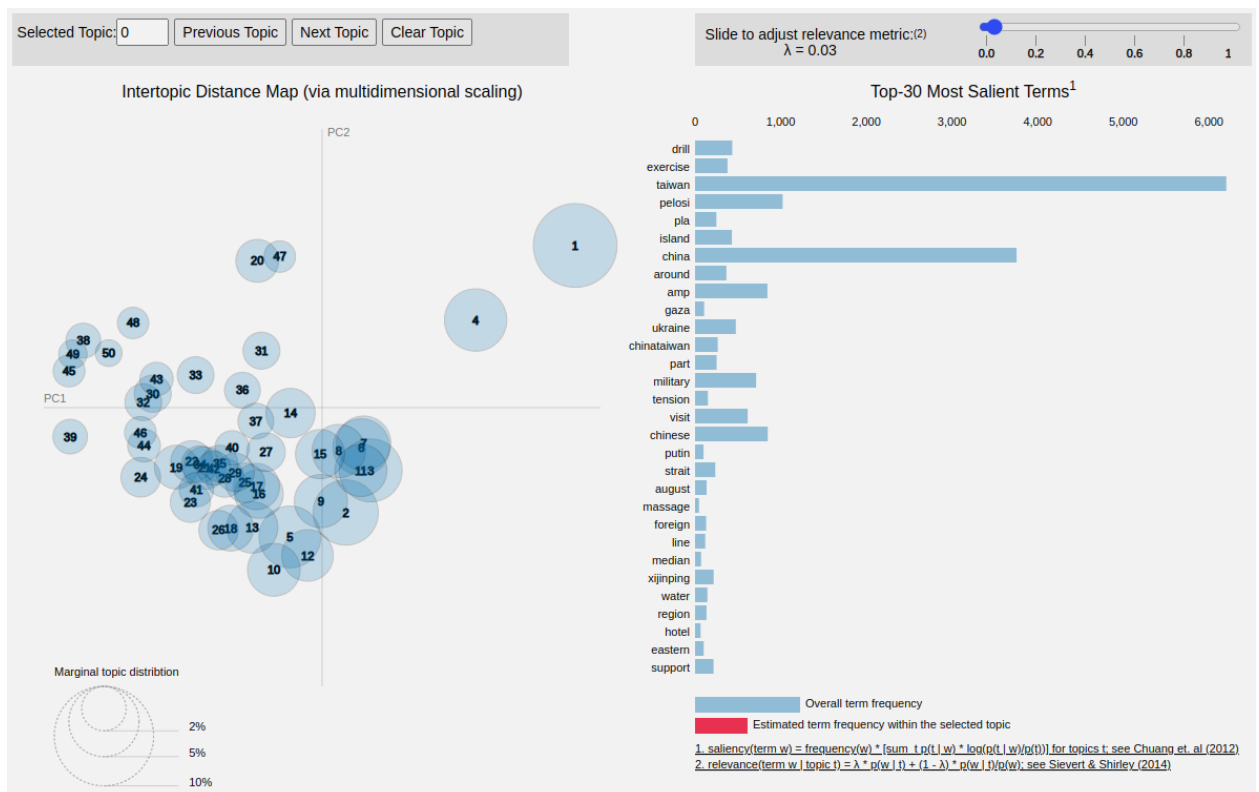
4. Model Development

- a. Topic Modelling: here I have used Latent Dirichlet Allocation/LDA to model the particular topics. It is popularly used to extract topics from a given corpus. LDA categorizes the text into a document and the words per topic, these are modelled based on [Dirchilet distributions and processes](#).

LDA makes two key assumptions the first is the documents are a mixture of topics and the second is topics are a mixture of words.

After developing the model with LDA I have shown the results using pyLDavis package with an interactive display.

- Each circle is a topic and the size represents the abundance of that topic in the corpus.



- After the LDA model training, the generated topics are shown in the notebook I have been working on. But to grasp the content I have plotted word clouds using the four topics as follows.

- The dataset is split into train-test with a 0.6/0.4 ratio and the text is vectorized with a unigram, bigram and trigram vectorizer to evaluate which choice could perform better. These represent the stack of words that can be found in a single vectorized set.
- I have used several classification algorithms to experiment with the results. These are SVM, Naive Bayes, Decision Tree, Random Forest, KNN and Logistic regression classification algorithms.
- Some of the Experiments are shown in the notebook I provided (EDA_and_Model_development).