

Project Description: Twitter US Airline Sentiment

Data Description:

A sentiment analysis job about the problems of each major U.S. airline. Twitter data was scraped from February of 2015 and contributors were asked to first classify positive, negative, and neutral tweets, followed by categorizing negative reasons (such as "late flight" or "rude service").

Dataset:

The project is from a dataset from Kaggle.

Link to the Kaggle project site: <https://www.kaggle.com/crowdflower/twitter-airline-sentiment>

The dataset has to be downloaded from the above Kaggle website.

The dataset has the following columns:

- tweet_id
- airline_sentiment
- airline_sentiment_confidence
- negativereson
- negativereson_confidence
- airline
- airline_sentiment_gold
- name
- negativereson_gold
- retweet_count
- text
- tweet_coord
- tweet_created
- tweet_location
- user_timezone

Objective:

To implement the techniques learnt as a part of the course.

Learning Outcomes:

- Basic understanding of text pre-processing.
- What to do after text pre-processing:
 - Bag of words
 - Tf-idf
- Build the classification model.
- Evaluate the Model.

Steps and tasks:

1. Import the libraries, load dataset, print shape of data, data description. (5 Marks)
2. Understand of data-columns: (5 Marks)
 - a. Drop all other columns except "text" and "airline_sentiment".
 - b. Check the shape of data.
 - c. Print first 5 rows of data.
3. Text pre-processing: Data preparation. (20 Marks)
 - a. Html tag removal.
 - b. Tokenization.
 - c. Remove the numbers.
 - d. Removal of Special Characters and Punctuations.
 - e. Conversion to lowercase.
 - f. Lemmatize or stemming.
 - g. Join the words in the list to convert back to text string in the dataframe. (So that each row contains the data in text format.)
 - h. Print first 5 rows of data after pre-processing.
4. Vectorization: (10 Marks)
 - a. Use CountVectorizer.
 - b. Use TfidfVectorizer.
5. Fit and evaluate model using both type of vectorization. (6+6 Marks)
6. Summarize your understanding of the application of Various Pre-processing and Vectorization and performance of your model on this dataset. (8 Marks)

Happy Learning!