# POLYTECHNIC UNIVERSITY OF THE PHILIPPINES

**SENTIMENT ANALYSIS IN PYTHON**

In Partial Fulfillment of the Requirements for the
Bachelor of Science in Computer Engineering

By:

**ESCOBAL, Anne Margaret O.**

To:

**Engr. Simon Salvador E. Tidon**

2022

**POLYTECHNIC UNIVERSITY OF THE PHILIPPINES**

**TABLE OF CONTENTS**

# POLYTECHNIC UNIVERSITY OF THE PHILIPPINES

**INTRODUCTION**

Sentiment analysis is the process of classifying whether a block of text is positive, negative, or, neutral. Sentiment analysis is contextual mining of words which indicates the social sentiment of a brand and also helps the business to determine whether the product which they are manufacturing is going to make a demand in the market or not. The goal which Sentiment analysis tries to gain is to analyze people's opinion in a way that it can help the businesses expand.

Python sentiment analysis is a methodology for analyzing a piece of text to discover the sentiment hidden within it. It accomplishes this by combining machine learning and natural language processing (NLP). Sentiment analysis allows you to examine the feelings expressed in a piece of text.

**LIBRARIES USED**

- **Plotly**

The Plotly Python library is an interactive open-source library. This can be a very helpful tool for data visualization and understanding the data simply and easily. plotly graph objects are a high-level interface to plotly which are easy to use. It can plot various types of graphs and charts like scatter plots, line charts, bar charts, box plots, histograms, pie charts, etc.

- **Pandas**

Python Pandas is defined as an open-source library that provides high-performance data manipulation in Python.

- **Matplotlib**

Matplotlib is an amazing visualization library in Python for 2D plots of arrays. Matplotlib is a multi-platform data visualization library built on NumPy arrays and designed to work with the broader SciPy stack.

- **Seaborn**

Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.

- **NLTK (National Language Toolkit)**

NLTK is a leading platform for building Python programs to work with human language data.

- **WORDCLOUD**

Word Cloud is a data visualization technique used for representing text data in which the size of each word indicates its frequency or importance. Significant textual data points can be highlighted using a word cloud. Word clouds are widely used for analyzing data from social network websites.
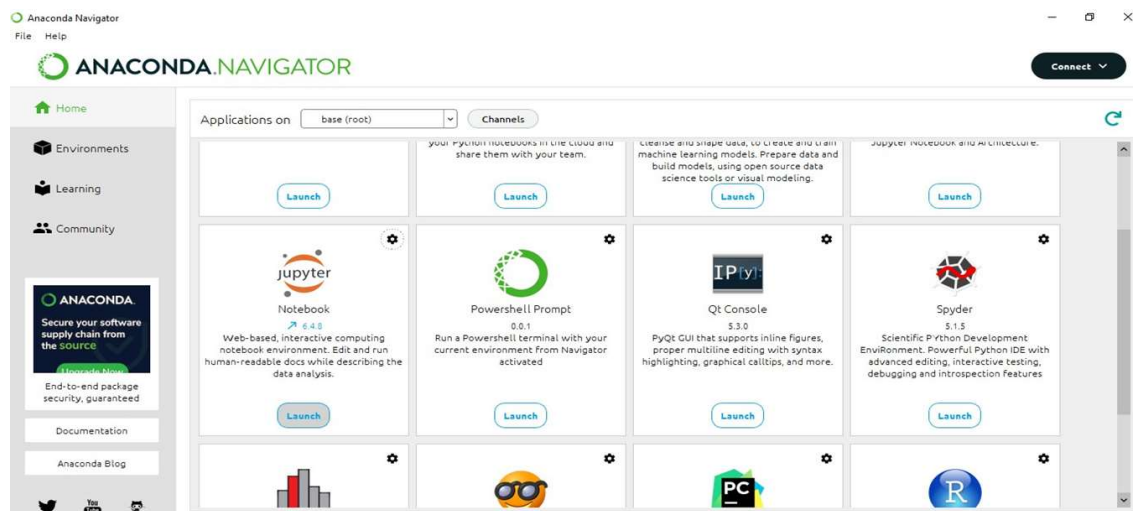
- **SKLEARN**

Scikit-learn (Sklearn) is the most useful and robust library for machine learning in Python. It provides a selection of efficient tools for machine learning and statistical modeling including classification, regression, clustering and dimensionality reduction via a consistence interface in Python.
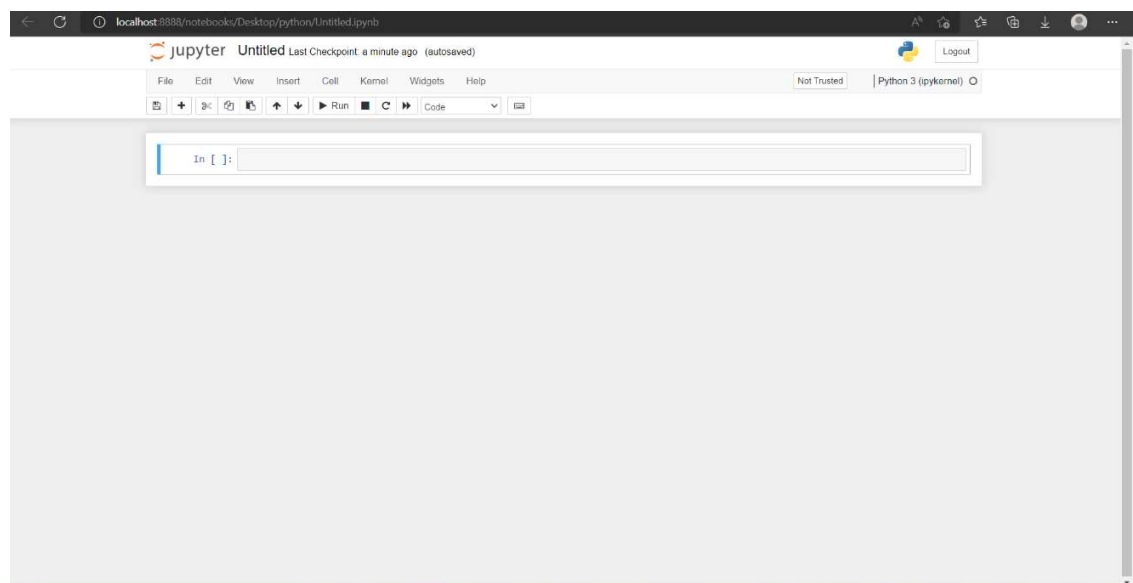
**Steps to implement Sentiments in Python**

1. Launch Jupyter from Anaconda Navigator.



2. Create a new Python Notebook.

3. Import modules/libraries needed.

- pip install plotly==5.8.0
- pip install pandas
- pip install matplotlib
- pip install seaborn
- pip install nltk
- pip install wordcloud
- pip install sklearn

4. Analysis

I used Review.csv file from Kaggle's Amazon Fine Food Reviews dataset to perform the analysis.



We can see that the data frame contains some product, user and review information.

The data that we will be using most for this analysis is "*Summary*", "*Text*", and "*Score.*"

*Text* — This variable contains the complete product review information.

*Summary* — This is a summary of the entire review.

*Score* — The product rating provided by the customer.

5. Data Analysis





Now, we can create some ***wordclouds*** to see the most frequently used words in the

reviews.

6. Classifying Tweets



7. More Data Analysis

Wordcloud Positive Sentiment:



Wordcloud Negative Sentiment



8. Building the Model.

```
In [10]: def remove_punctuation(text):
             final = "".join(u for u in text if u not in ("?", ".", ";", ":", "!",'"'))
             return final
         df['Text'] = df['Text'].apply(remove_punctuation)
         df = df.dropna(subset=['Summary'])
         df['Summary'] = df['Summary'].apply(remove_punctuation)

         dfNew = df[['Summary','sentiment']]
         dfNew.head()
```

Out[10]:

|   | Summary | sentiment |
|---|---|---|
| 0 | Good Quality Dog Food | 1 |
| 1 | Not as Advertised | -1 |
| 2 | Delight says it all | 1 |
| 3 | Cough Medicine | -1 |
| 4 | Great taffy | 1 |

## 9. Testing



```
In [11]: import numpy as np
         index = df.index
         df['random_number'] = np.random.randn(len(index))
         train = df[df['random_number'] <= 0.8]
         test = df[df['random_number'] > 0.8]

         from sklearn.feature_extraction.text import CountVectorizer
         vectorizer = CountVectorizer(token_pattern=r'\b\w+\b')
         train_matrix = vectorizer.fit_transform(train['Summary'])
         test_matrix = vectorizer.transform(test['Summary'])

         from sklearn.linear_model import LogisticRegression
         lr = LogisticRegression()

         X_train = train_matrix
         X_test = test_matrix
         y_train = train['sentiment']
         y_test = test['sentiment']

         lr.fit(X_train,y_train)
         predictions = lr.predict(X_test)

         from sklearn.metrics import confusion_matrix,classification_report
         new = np.asarray(y_test)
         confusion_matrix(predictions,y_test)

         print(classification_report(predictions,y_test))
```

```
               precision    recall  f1-score   support

          -1       0.00      0.00      0.00         1
           1       0.98      0.92      0.95        49

    accuracy                           0.90        50
   macro avg       0.49      0.46      0.47        50
weighted avg       0.96      0.90      0.93        50
```

In [ ]:

8