# Sentiment Analysis of Tweets about Georgia State University

A Natural Language Processing Use-Case

David King Agbi
Joel Dassoundo
Ifeanyi Moses Uzowuru

J. Mack Robinson College of Business, Georgia State University

# Introduction

- Sentiment analysis is simply identifying and classifying sentiments (which are the point of view or emotions) that are expressed in the message or text source

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Introduction

- Sentiment analysis is simply identifying and classifying sentiments (which are the point of view or emotions) that are expressed in the message or text source
- A Natural Language Processing is simply preparing computers to understand and process the language we speak, either through audio or texts. Siri, Cortana, Bixby and google translator are all examples of this.

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Introduction

- Sentiment analysis is simply identifying and classifying sentiments (which are the point of view or emotions) that are expressed in the message or text source
- A Natural Language Processing is simply preparing computers to understand and process the language we speak, either through audio or texts. Siri, Cortana, Bixby and google translator are all examples of this.
- Build a machine learning model that can analyze several tweets made about GSU with the ability to estimate the sentiments of those tweets.
  Thus, we would classify these tweets as positive, neutral or negative.

GEORGIA STATE

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# WHY?

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

Figure: Process involved in the project

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Task 1 - Exploratory Data Analysis & Visualization

Figure: Output of extracted tweets showing the various columns

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Task 1 - Exploratory Data Analysis & Visualization

■ We do this by first extracting the tweets using Twitter API. The figure below shows the flow diagram of the process



Figure: Extracting Tweets using Twitter API

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Task 1

- Next thing is to import all the necessary libraries we employed in the project

- Next thing is to import all the necessary libraries we employed in the project

- We then analyse or explore the dataset such as getting the information, describing the data, shape of the data, checking for null elements and so on. We visualized the sentiments using the countplot and visualized the null elements by using the heatmap and got these results:

GEORGIA STATE

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

In [416]:
```
##EXPLORING THE DATASET
sns.heatmap(df.isnull(),yticklabels= False, cbar = False, cmap = "crest")

## The isnull here checks if there are any null elements. So the heatmap
## is plain which indicates there are no null elements.
```

Out[416]: <AxesSubplot:>



David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Task 1 - Exploratory Data Analysis & Visualization

As we can see, this looks not balanced

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Task 1 - Exploratory Data Analysis & Visualization

■ Therefore, we applied the augmentation technique to populate the negative tweets to make it more significant. This is the new plot of the distribution of sentiments



David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

- The following are the processing techniques that we employed in the listed order to normalize the text:

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

- The following are the processing techniques that we employed in the listed order to normalize the text:
- Removal of punctuations.

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

- The following are the processing techniques that we employed in the listed order to normalize the text:
- Removal of punctuations.
- Removal of stop words.

GEORGIA STATE

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

- The following are the processing techniques that we employed in the listed order to normalize the text:
- Removal of punctuations.
- Removal of stop words.
- Tokenization

**GEORGIA STATE.**

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

- The following are the processing techniques that we employed in the listed order to normalize the text:
- Removal of punctuations.
- Removal of stop words.
- Tokenization
- Stemming/Lemmatization

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

- The following are the processing techniques that we employed in the listed order to normalize the text:
- Removal of punctuations.
- Removal of stop words.
- Tokenization
- Stemming/Lemmatization
- Converting capital letters to small letters

GEORGIA STATE

# Task 3 - Feature Extraction

- Data in the form of text is not suitable for training a machine learning model. For this reason, we had to convert the tweets into numerical features, ensuring that the inherent learnable pattern is conserved in the best way possible.

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

- Data in the form of text is not suitable for training a machine learning model. For this reason, we had to convert the tweets into numerical features, ensuring that the inherent learnable pattern is conserved in the best way possible.

- To do this, we used the count vectorizer or vectorization (Term Frequency-Inverse Document Frequency (TF-IDF)) to perform textual transformations into vectors

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Task 4 - Machine Learning & Model Evaluation

- Now that the feature extraction is done, the data is ready to be fed into a model

**GEORGIA STATE**

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

- Now that the feature extraction is done, the data is ready to be fed into a model
- Finally, the data is split into training and test using stratified sampling so that the split follows the population distribution and then we applied the machine learning algorithms. Thus, we compared the performance of the:
  1. Multinomial Naive Baye's Model
  2. Multinomial Logistic regression Model
  3. The Decision Tree (Random Forest) Model

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Multinomial Naive Baye's Model

- We had the following confusion matrix and classification report after using the Multinomial Naive Baye's Model

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Multinomial Naive Baye's model

```
#Evaluation scores for multinomial naive bayes
print(classification_report(y_test, naive_bayes))
```

|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| -1.0       | 0.92      | 0.79   | 0.85     | 466     |
| 0.0        | 0.85      | 0.91   | 0.88     | 1063    |
| 1.0        | 0.80      | 0.78   | 0.79     | 572     |
| accuracy   |           |        | 0.85     | 2101    |
| macro avg  | 0.86      | 0.83   | 0.84     | 2101    |
| weighted avg | 0.85    | 0.85   | 0.85     | 2101    |

GEORGIA
STATE

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Multinomial Logistic Regression model

■ We had the following confusion matrix and classification report after using the Multinomial Logistic Regression model



David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Multinomial Logistic Regression model

```python
#Evaluation scores for Multinomial Logistics Regression
from sklearn.metrics import classification_report
print(classification_report(y_test,logistics_pred))
```

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| -1.0 | 0.96 | 0.88 | 0.92 | 466 |
| 0.0 | 0.88 | 0.95 | 0.92 | 1063 |
| 1.0 | 0.90 | 0.83 | 0.86 | 572 |
| accuracy |  |  | 0.90 | 2101 |
| macro avg | 0.91 | 0.89 | 0.90 | 2101 |
| weighted avg | 0.91 | 0.90 | 0.90 | 2101 |

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Random Forest Model

■ We had the following confusion matrix and classification report after using the Random Forest Model

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Random Forest Model

```
#Evaluation scores for random forest
print(classification_report(y_test, rand_f))
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| -1.0         | 0.87      | 0.85   | 0.86     | 466     |
| 0.0          | 0.83      | 0.93   | 0.88     | 1063    |
| 1.0          | 0.87      | 0.70   | 0.77     | 572     |
| accuracy     |           |        | 0.85     | 2101    |
| macro avg    | 0.86      | 0.83   | 0.84     | 2101    |
| weighted avg | 0.85      | 0.85   | 0.85     | 2101    |

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Conclusion

- In terms of negative prediction, Multinomial logistics Regression performed the best in classifying correctly negative tweets and in its precision power, which is why it has the best F1 score and we thereby confirm the algorithm as our model of choice.

David King Agbi  Joel Dassoundo  Ifeanyi Moses Uzowuru

# Thank You!