



Sentiment Analysis Project.

Business Overview

Introduction

In today's digital age, social media platforms such as Twitter have become powerful sources of real-time customer feedback and opinions. Understanding the sentiment expressed by customers toward specific brands and products is essential for businesses to make informed decisions, enhance customer satisfaction, and maintain a positive brand reputation.

The goal of this project aims to develop a sentimental analysis model specifically tailored to analyze Twitter data related to Google, Apple, and other products.

Business Problem Statement

As a consulting firm, Twitter has assigned us the task of building a model that can rate the sentiment of a Tweet based on its content that can correctly categorize Twitter sentiment about Apple and Google products into positive, negative, or neutral categories and gain valuable insights into public perception, that will be used for informed decision-making in business strategies and customer satisfaction enterprise.

Data Understanding

For this project, we utilized a dataset from CrowdFlower via Data.world, consisting of approximately 9,000 tweet sentiments about Apple and Google products. The dataset includes columns such as `tweet_text`, `emotion_in_tweet_is_directed_at`, and

is_there_an_emotion_directed_at_a_brand_or_product. Our objective is to develop a sentiment analysis model that should accurately classify tweets into positive, negative, or neutral sentiment categories.

Objectives for the project:

1. To develop a binary classifier that can classify tweets into positive or negative sentiment categories. This will serve as a proof of concept and provide a foundation for further analysis. The classifier will be a Logistic regression model and the benchmark accuracy will be 85%.
2. To expand to a multiclass classifier, thereby incorporating the neutral tweets to create a multiclass classifier that can accurately classify tweets as positive, negative, or neutral. This will provide a more comprehensive sentiment analysis of the tweets. The classifier will be a XGBoost model and MultinomialNB the benchmark accuracy will be 70%.
3. To compare sentiment between Apple and Google products by analyzing the sentiment distribution of tweets mentioning Apple, Google, and other products.

Data Preparation

To prepare the data, we performed various preprocessing steps.

These included; removing duplicates, handling missing values by filling them with "none" for the emotion_in_tweet_is_directed_at column, and dropping the row with a missing tweet_text. We also employed tokenization, lowercase conversion, stopword removal, and lemmatization to refine the text data.

Modeling

The packages/libraries utilized for data preparation and analysis included NLTK, sklearn, and pandas.

- NLTK was used for tokenization, stopword removal, and lemmatization.
- Sklearn's CountVectorizer was employed for vectorization.
- Pandas was used for data manipulation and handling missing values.

For modeling, we employed logistic regression as the binary classifier to classify tweets into positive or negative sentiment categories. As a benchmark, we aimed for an accuracy of 70%.

In addition, we expanded the model to a multiclass classifier to include the neutral sentiment category.

Testing Model

Evaluation

To evaluate the performance of the final model, we employed accuracy. This metric provided insight into the model's ability to correctly classify instances and predict positive and negative sentiment.

Accuracy is a commonly used metric for evaluating sentiment analysis models because it provides a straightforward measure of the overall performance. In sentiment analysis, the primary goal is to classify a given text or document into positive, negative, or neutral sentiment categories accurately. Accuracy, in this context, represents the proportion of correctly classified instances compared to the total number of instances. It gives a clear indication of how well a sentiment analysis model is able to predict sentiment labels correctly across the entire dataset.

In conclusion, the sentiment analysis model developed using logistic regression demonstrated satisfactory performance in classifying tweet sentiments. The pre-processing steps and modeling choices contributed to accurate predictions and meaningful insights. However, limitations such as missing values and data quality issues should be considered when interpreting the results. Recommendations include further refining the model and expanding the dataset for improved sentiment analysis accuracy.

Deployment

We deployed the model using the framework STREAMLIT. The application takes a sentence and returns the sentiment. Link to web app <https://moringaphase4project.streamlit.app/> (<https://moringaphase4project.streamlit.app/>)

Data understanding

Load Libraries

We will start by importing all libraries used throughout the notebook.

In [1]:

```
1 #Loading libraries
2
3 import re
4 import pandas as pd
5 import numpy as np
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8 %matplotlib inline
9 import string
10 import warnings
11 warnings.filterwarnings(action='ignore', category=FutureWarning)
12
13 import nltk
14 nltk.download('punkt')
15 nltk.download('stopwords')
16 nltk.download('wordnet')
17 from nltk.stem import WordNetLemmatizer
18 from sklearn.preprocessing import LabelEncoder
19 from nltk.tokenize import word_tokenize, sent_tokenize
20 from nltk.corpus import stopwords
21 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
22 from nltk import FreqDist
23
24 # modelling libraries
25 from sklearn.model_selection import train_test_split
26 from sklearn.linear_model import LogisticRegression
27 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
28 from sklearn import metrics
29 from sklearn.naive_bayes import MultinomialNB
30 import xgboost as xgb
31
32 #imbalance libraries
33 from imblearn.over_sampling import RandomOverSampler
34 from imblearn.under_sampling import RandomUnderSampler
35 import xgboost as xgb
36 from sklearn.metrics import accuracy_score
```

```
[nltk_data] Downloading package punkt to
[nltk_data]     C:\Users\user\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]     C:\Users\user\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]     C:\Users\user\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```

Loading the data

In [2]:

```
1 # Loading and previewing data
2 data = pd.read_csv('data.csv', encoding='ISO-8859-1')
3 data.head()
```

Out[2]:

	tweet_text	emotion_in_tweet_is_directed_at	is_there_an_emotion_directed_at_a_brand_or_pro	
0	@wesley83 I have a 3G iPhone. After 3 hrs twe...	iPhone		Negative em
1	@jessedee Know about @fludapp ? Awesome iPad/i...	iPad or iPhone App		Positive em
2	@swonderlin Can not wait for #iPad 2 also. The...	iPad		Positive em
3	@sxsw I hope this year's festival isn't as cra...	iPad or iPhone App		Negative em
4	@sxtxstate great stuff on Fri #SXSW: Marissa M...	Google		Positive em



Statistic Describing the dataset

```
In [3]:  
1 # class to describe dataset  
2  
3 class Describer:  
4  
5     # initialize object  
6     def __init__(self, df):  
7         self.df = df  
8  
9     # method to check shape of data  
10    def shape(self):  
11        out = print(f"The DataFrame has:\n\t* {self.df.shape[0]} rows\n\t")  
12        return out  
13  
14    # method to check info on dataset  
15    def data_info(self):  
16        out = print(self.df.info(), '\n')  
17        return out  
18  
19    # method to describe numerical columns  
20    def data_describe(self):  
21        out = self.df.describe()  
22        return out
```

```
In [4]:  
1 # creating an instance of the class describer  
2 describe_df = Describer(data)  
3  
4 # lets view the shape of the data  
5 describe_df.shape()
```

The DataFrame has:

* 9093 rows
* 3 columns

In [5]:

```

1 # Lets print summary infomation on the dataset
2 print('Summary infomation on dataset')
3 print('-----')
4 describe_df.data_info()

```

Summary infomation on dataset

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9093 entries, 0 to 9092
Data columns (total 3 columns):
 #   Column           Non-Null Count  Dtype  
pe
---  -- 
-- 
0   tweet_text        9092 non-null    obj    
1   emotion_in_tweet_is_directed_at 3291 non-null    obj    
2   is_there_an_emotion_directed_at_a_brand_or_product 9093 non-null    obj    
dtypes: object(3)
memory usage: 213.2+ KB
None

```

Summary of Data Understanding

- The dataset has 3 columns namely:tweet_text, emotion_in_tweet_is_directed_at ,is_there_an_emotion_directed_at_a_brand_or_product
- The dataset has 9093 rows
- All the dataset's columns datatypes are objects

3. Data Preparation

Introduction

We will be preparing our data for analysis by checking for attributes such as;

Missing values

Duplicates

Other inconsistencies as computed below

Type *Markdown* and *LaTeX*: α^2

Missing values

```
In [6]: # function to identify missing values
2
3 # identify missing
4 def identify_missing_values(data):
5     """Identify if the data has missing values"""
6     # identify if data has missing values(data.isnull().any())
7     # empty dict to store missing values
8     missing = []
9     for i in data.isnull().any():
10         # add the bool values to empty list
11         missing.append(i)
12     # convert list to set (if data has missing value, the list should have
13     missing_set = set(missing)
14     if (len(missing_set) == 1):
15         out = print("The Data has no missing values")
16     else:
17         out = print("The Data has missing values.")
18
19     return out
20
21 identify_missing_values(data)
```

The Data has missing values.

```
In [7]: 1 # function to display missing values
2
3 def missing_values(data):
4     """A simple function to identify data has missing values"""
5     # identify the total missing values per column
6     # sort in order
7     miss = data.isnull().sum().sort_values(ascending = False)
8
9     # calculate percentage of the missing values
10    percentage_miss = (data.isnull().sum() / len(data)).sort_values(ascen
11
12    # store in a dataframe
13    missing = pd.DataFrame({"Missing Values": miss, "Percentage)": perc
14
15    # remove values that are missing
16    missing.drop(missing[missing["Percentage(%)" == 0].index, inplace =
17
18    return missing
19
20
21 missing_values(data)
22
```

Out[7]:

	Missing Values	Percentage(%)
emotion_in_tweet_is_directed_at	5802	0.638073
tweet_text	1	0.000110

In [8]:

```
1 # dropping missing row in tweet_text
2
3 data.dropna(axis=0,subset=['tweet_text'],inplace=True)
```

In [9]:

```
1 # confirmation
2 data.isna().any()
```

Out[9]:

tweet_text	False
emotion_in_tweet_is_directed_at	True
is_there_an_emotion_directed_at_a_brand_or_product	False
dtype: bool	

In [10]:

```
1 # replacing null values in recipient with none
2 data['emotion_in_tweet_is_directed_at'].fillna('none', inplace =True)
```

```
In [11]: 1 #confirming we replaced the 'nan' with 'none'  
2 data['emotion_in_tweet_is_directed_at'].value_counts()
```

```
Out[11]: none          5801  
iPad           946  
Apple          661  
iPad or iPhone App 470  
Google          430  
iPhone          297  
Other Google product or service 293  
Android App      81  
Android          78  
Other Apple product or service 35  
Name: emotion_in_tweet_is_directed_at, dtype: int64
```

```
In [12]: 1 data.isna().any()
```

```
Out[12]: tweet_text          False  
emotion_in_tweet_is_directed_at  False  
is_there_an_emotion_directed_at_a_brand_or_product  False  
dtype: bool
```

Justification

- We discovered that there is 0.01% missing values in **tweet_text** column and 63.81% missing values in **emotion_in_tweet_is_directed_at** column.
- We replaced the missing values in the **emotion tweet is directed at** column with **none** for computational purposes and dropped the missing values in the **tweet text** column

Duplicates

```
In [13]: 1 # checking for duplicates
2
3 # Duplicated entries
4 def identify_duplicates(data):
5     """Simple function to identify any duplicates"""
6     # identify the duplicates (dataframe.name.duplicated() , can add .sum()
7     # empty list to store Bool results from duplicated
8     duplicates = []
9     for i in data.duplicated():
10         duplicates.append(i)
11     # identify if there is any duplicates. (If there is any we expect a T
12     duplicates_set = set(duplicates)
13     if (len(duplicates_set) == 1):
14         print("The Data has no duplicates")
15     else:
16         no_true = 0
17         for val in duplicates:
18             if (val == True):
19                 no_true += 1
20         # percentage of the data represented by duplicates
21         duplicates_percentage = np.round(((no_true / len(data)) * 100), 3
22         print(f"The Data has {no_true} duplicated rows.\nThis constitutes
23
24
25
26 identify_duplicates(data)
```

The Data has 22 duplicated rows.
 This constitutes 0.242% of the data set.

```
In [14]: 1 data.duplicated().sum()
```

Out[14]: 22

```
In [15]: 1 # handling the duplicates
2 def remove_duplicated_rows(data):
3     """Simple Function to remove duplicated rows"""
4     data.drop_duplicates(subset=None, keep="first", inplace=True)
5     # confirm if the duplicated rows have been removed
6     confirm = identify_duplicates(data)
7
8     return confirm
9
10
11 remove_duplicated_rows(data)
```

The Data has no duplicates

Summary

We found out that the **tweet text** column had 22 duplicates (0.242%) and went ahead and dropped them.

Renaming long column names

In [16]:

```
1 # renaming the columns
2 data = data.rename(columns={'emotion_in_tweet_is_directed_at':'recipient'
3                           'is_there_an_emotion_directed_at_a_brand_or_product':
4 data.head()
```

Out[16]:

	tweet_text	recipient	emotion
0	@wesley83 I have a 3G iPhone. After 3 hrs twe...	iPhone	Negative emotion
1	@jessedee Know about @fludapp ? Awesome iPad/i... iPad or iPhone App	iPad or iPhone App	Positive emotion
2	@swonderlin Can not wait for #iPad 2 also. The...	iPad	Positive emotion
3	@sxsw I hope this year's festival isn't as cra... iPad or iPhone App	iPad or iPhone App	Negative emotion
4	@sxtystate great stuff on Fri #SXSW: Marissa M...	Google	Positive emotion

Justification

We renamed the column `emotion_in_tweet_is_directed_at` to `recipient` and the column `is_there_an_emotion_directed_at_a_brand_or_product` to `emotion` to enhance clarity and facilitate their usage in our subsequent computation

4. Exploratory Data Analysis (EDA)

Introduction

We will conduct Univariate and Bivariate analysis of the sentiments and create visualizations to see how they relate with each other and individually.

Univariate analysis

We will plot a count plot to show the distribution of the column 'emotion'

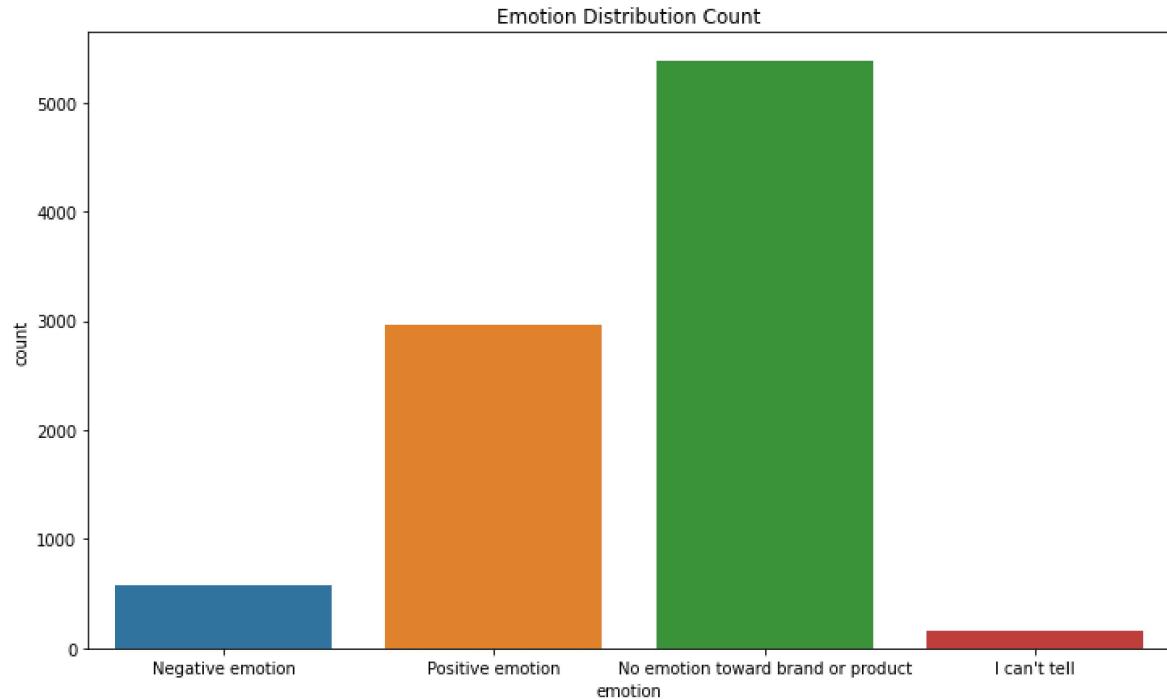
In [17]:

```

1 # checking the distribution of emotions
2 percentage = data['emotion'].value_counts(normalize=True)
3 print(percentage * 100)
4
5 # count plot
6 plt.figure(figsize=(12, 7))
7 sns.countplot(data['emotion'])
8 plt.title('Emotion Distribution Count')
9 plt.show();

```

No emotion toward brand or product	59.261301
Positive emotion	32.745314
Negative emotion	6.273429
I can't tell	1.719956
Name: emotion, dtype: float64	



Observation

We see that majority of people had No emotion toward brand or product at 59.26% followed by people with Positive emotions at 32.75%

I can't tell makes up less than 2% of our dataset, and doesn't offer much more information in the way of word significance than the tweets labeled No emotion toward brand or product.

Bivariate analysis

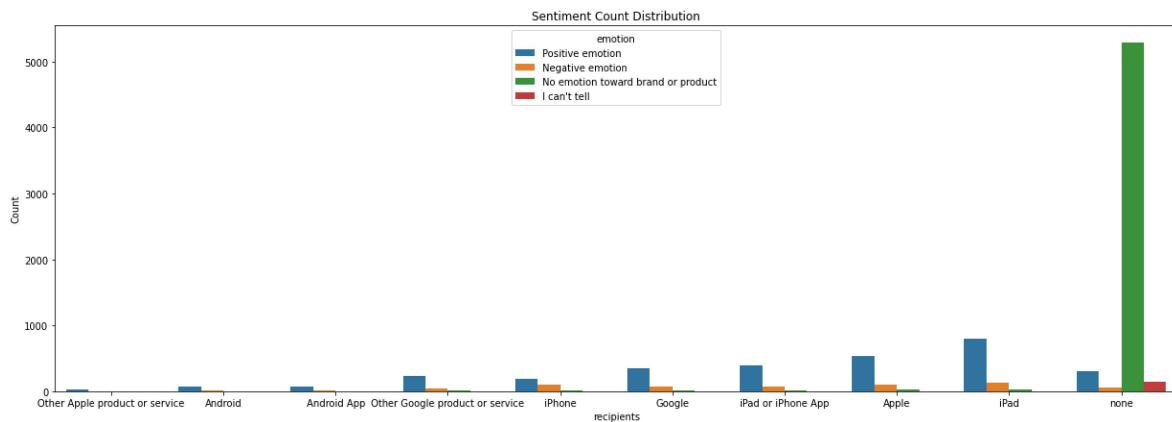
We plot a bar plot to visualize the distribution of sentiments'emotion' relative to the recipient

In [18]:

```

1 # Let's see sentiment distribution product-wise
2 recipient_group = data.groupby('recipient')['emotion'].value_counts().res
3 # Calculate the ascending order of states based on churn count
4 ascending_order = recipient_group.groupby('recipient')['count'].sum().sort
5
6 plt.figure(figsize=(21,7))
7 sns.barplot(x='recipient', y='count', hue='emotion', data=recipient_group)
8 plt.xlabel('recipients')
9 plt.ylabel('Count')
10 plt.title('Sentiment Count Distribution')
11 plt.show()

```



Interpretation

From the graph above Ipad had the highest positive emotions followed by Apple

From the observation above since we similar product brand names we will map products to their brand

Mapping products to Brands

In [19]:

```

1 #What are value counts of recipient column
2 data['recipient'].value_counts()

```

Out[19]:

none	5788
iPad	945
Apple	659
iPad or iPhone App	469
Google	428
iPhone	296
Other Google product or service	293
Android App	80
Android	77
Other Apple product or service	35
Name: recipient, dtype: int64	

```
In [20]: 1 # overall brand distribution
2 # feature engineering
3 data['brand'] = data['recipient'].map({'Other Apple product or service':
4                                     'Android': 'Google',
5                                     'Android App': 'Google',
6                                     'Other Google product or servi
7                                     'iPhone': 'Apple',
8                                     'Google':'Google',
9                                     'iPad or iPhone App': 'Apple',
10                                    'Apple':'Apple',
11                                    'iPad':'Apple',
12                                    'none': 'none'})
```

```
In [21]: 1 # confirm brand column
2 data.head(3)
```

Out[21]:

	tweet_text	recipient	emotion	brand
0	@wesley83 I have a 3G iPhone. After 3 hrs twe...	iPhone	Negative emotion	Apple
1	@jessedee Know about @fludapp ? Awesome iPad/i... iPad or iPhone App	iPad or iPhone App	Positive emotion	Apple
2	@swonderlin Can not wait for #iPad 2 also. The...	iPad	Positive emotion	Apple

Plot brand Distributions

Univariate Analysis

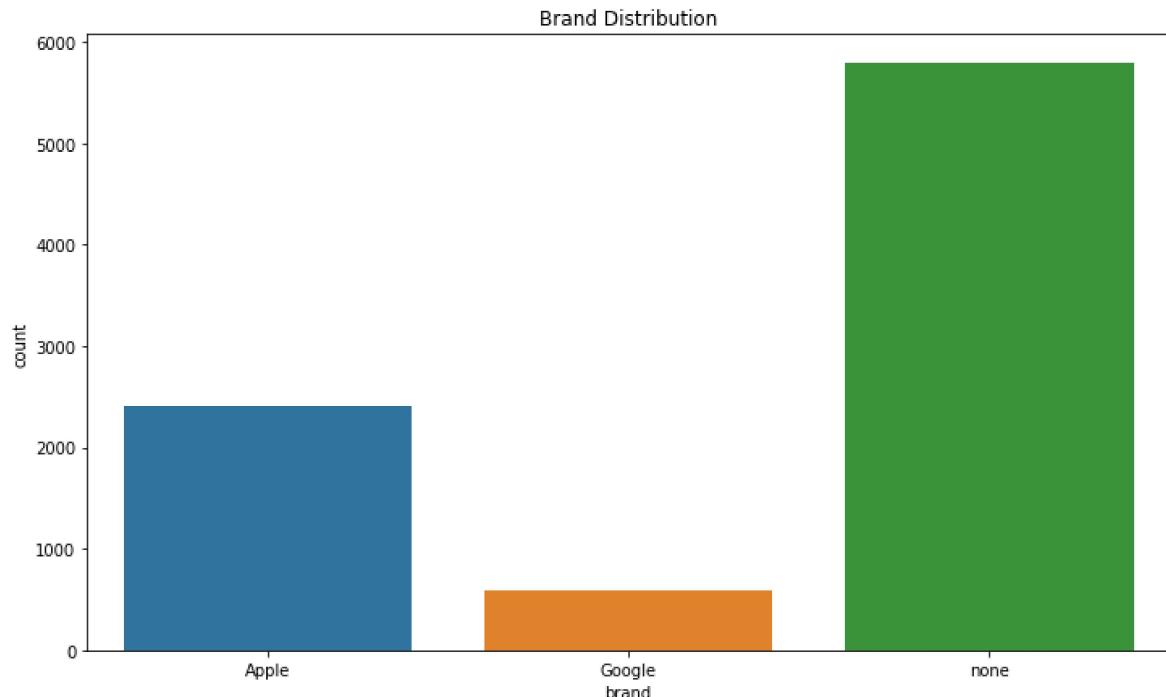
In [22]:

```

1 # checking the distribution of brand sentiment
2 percentage = data['brand'].value_counts(normalize=True)
3 print(percentage * 100)
4
5 # count plot
6 plt.figure(figsize=(12, 7))
7 sns.countplot(data['brand'])
8 plt.title('Brand Distribution')
9 plt.show();

```

none 65.945084
 Apple 27.389769
 Google 6.665148
 Name: brand, dtype: float64



Insight

From the graph above Majority of people had not specified the brand, and for those who did, Apple had the highest brand feedback as compared to Google

Bivariate analysis

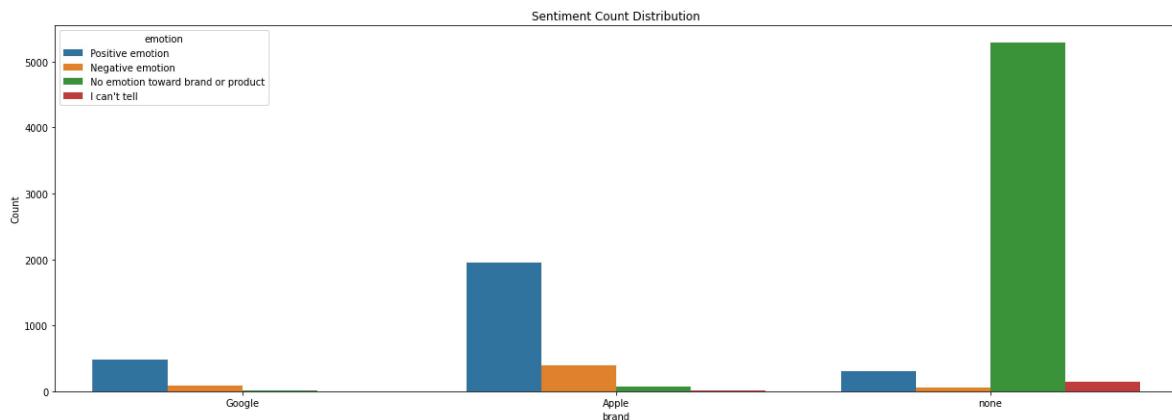
We plot a bar plot to visualize the distribution of sentiments '**emotion**' relative to the **Brand**

In [23]:

```

1 # Let's see sentiment distribution brand-wise
2 brand_group = data.groupby('brand')['emotion'].value_counts().reset_index
3 # Calculate the ascending order of states based on churn count
4 ascending_order = brand_group.groupby('brand')['count'].sum().sort_values
5
6 plt.figure(figsize=(21,7))
7 sns.barplot(x='brand', y='count', hue='emotion', data=brand_group, order=
8 plt.xlabel('brand')
9 plt.ylabel('Count')
10 plt.title('Sentiment Count Distribution')
11 plt.show()

```



Interpretation

From the graph above Apple had the highest positive emotions compared to Google

5. Preprocessing

Introduction

In this section we will;

- Convert the tweet text to lower case
- Remove html tags
- Remove the Url
- Expand the contractions
- Remove the punctuations
- Tokenize
- Remove stopwords
- Lemmatize the tweet

So as to prepare the data for modeling.

```
In [24]: 1 # what are the columns?
          2 data.columns
```

Out[24]: Index(['tweet_text', 'recipient', 'emotion', 'brand'], dtype='object')

Changing text to lowercase

```
In [25]: 1 # normalizing the data to Lowecase
          2 data['tweet_text'] = data['tweet_text'].str.lower()
          3 data.head(2)
```

Out[25]:

	tweet_text	recipient	emotion	brand
0	@wesley83 i have a 3g iphone. after 3 hrs twe...	iPhone	Negative emotion	Apple
1	@jessedee know about @fludapp ? awesome ipad/i... iPad or iPhone App		Positive emotion	Apple

Removing Html tags and URL

```
In [26]: 1 #Identifying Html tag
          2 data['tweet_text'][14]
```

Out[26]: 'great #sxsw ipad app from @madebymany: <http://tinyurl.com/4nqv921>' (<http://tinyurl.com/4nqv921>)

```
In [27]: 1 # removing the html tags
          2 def remove_html(review):
          3     pattern = re.compile('<.*?>')
          4     return pattern.sub(r'', review)
          5
          6 data['tweet_text'] = data['tweet_text'].apply(remove_html)
```

```
In [28]: 1 # removing URL and @ sign
          2 def preprocess_text_removingq_URLand_atsign(text):
          3     # Remove URLs
          4     clean_text = re.sub(r"http\S+|www\S+|https\S+", "", text)
          5     text = re.sub(r'@[^\s]+', 'user', clean_text)
          6     # Other preprocessing steps like removing punctuation, converting to
          7     # ...
          8     return text
          9
          10 data['tweet_text'] = data['tweet_text'].apply(preprocess_text_removingq_U
```

```
In [29]: 1 #lets confirm the url has been replaced with an empty
          2 data['tweet_text'][14]
```

Out[29]: 'great #sxsw ipad app from user '

In [30]: 1 data.head()

Out[30]:

	tweet_text	recipient	emotion	brand
0	.user i have a 3g iphone. after 3 hrs tweeting...	iPhone	Negative emotion	Apple
1	user know about user ? awesome ipad/iphone app...	iPad or iPhone App	Positive emotion	Apple
2	user can not wait for #ipad 2 also. they shoul...	iPad	Positive emotion	Apple
3	user i hope this year's festival isn't as cras...	iPad or iPhone App	Negative emotion	Apple
4	user great stuff on fri #sxsw: marissa mayer (...	Google	Positive emotion	Google

Expanding Words

We convert abbreviated words to their full form

In [31]:

```

1 # expanding the contractions (is-nots)
2 import contractions
3 #from nltk.stem import WordNetLemmatizer
4
5 def expand(text):
6     # Expand contractions
7     expanded_text = contractions.fix(text)
8
9     return expanded_text
10
11 data['tweet_text'] = data['tweet_text'].apply(expand)

```

In [32]: 1 data.tweet_text[3]

Out[32]: "user i hope this year's festival is not as crashy as this year's iphone ap
p. #sxsw"

Removing Punctuations

In [33]:

```

1 # punctuation
2 exclude = string.punctuation
3 print(exclude)

```

!"#\$%&'()*+,-./:;<=>?@[\\]^_`{|}~

In [34]:

```

1 # remove punctuations
2 def remove_punctuations(tweet):
3     return tweet.translate(str.maketrans(' ', ' ', exclude))

```

In [35]:

```
1 data['tweet_text']=data['tweet_text'].apply(remove_punctuations)
2 data.head()
```

Out[35]:

	tweet_text	recipient	emotion	brand
0	user i have a 3g iphone after 3 hrs tweeting a...	iPhone	Negative emotion	Apple
1	user know about user awesome ipadiphone app t...	iPad or iPhone App	Positive emotion	Apple
2	user can not wait for ipad 2 also they should ...	iPad	Positive emotion	Apple
3	user i hope this years festival is not as cras...	iPad or iPhone App	Negative emotion	Apple
4	user great stuff on fri sxsw marissa mayer goo...	Google	Positive emotion	Google

In [36]:

```
1 data.tweet_text[20]
```

Out[36]:

```
'need to buy an ipad2 while i am in austin at sxsw not sure if i will need t
o q up at an austin apple store'
```

Tokenization and Removing Stopwords

In [37]:

```
1 # tokenize the tweets
2 def tokenize_text(tweet):
3     return word_tokenize(tweet)
4
5 data['tweet_text'] = data['tweet_text'].apply(tokenize_text)
6 data.head()
```

Out[37]:

	tweet_text	recipient	emotion	brand
0	[user, i, have, a, 3g, iphone, after, 3, hrs, ...]	iPhone	Negative emotion	Apple
1	[user, know, about, user, awesome, ipadiphone,...]	iPad or iPhone App	Positive emotion	Apple
2	[user, can, not, wait, for, ipad, 2, also, the...]	iPad	Positive emotion	Apple
3	[user, i, hope, this, years, festival, is, not...]	iPad or iPhone App	Negative emotion	Apple
4	[user, great, stuff, on, fri, sxsw, marissa, m...]	Google	Positive emotion	Google

```
In [38]:  
1 # removing stop words  
2 stop_words = set(stopwords.words('english'))  
3  
4 def remove_stopwords(tweet):  
5     # Use List comprehension for efficient List creation  
6     new_tweet = [word for word in tweet if word not in stop_words]  
7     return " ".join(new_tweet)  
8  
9 # Apply the function to the 'review' column  
10 data['tweet_text'] = data['tweet_text'].apply(remove_stopwords)  
11 data.head()
```

Out[38]:

	tweet_text	recipient	emotion	brand
0	user 3g iphone 3 hrs tweeting riseaustin dead ...	iPhone	Negative emotion	Apple
1	user know user awesome ipadiphone app likely a... iPad or iPhone App	iPad or iPhone App	Positive emotion	Apple
2	user wait ipad 2 also sale sxsw	iPad	Positive emotion	Apple
3	user hope years festival crashy years iphone a... iPad or iPhone App	iPad or iPhone App	Negative emotion	Apple
4	user great stuff fri sxsw marissa mayer google...	Google	Positive emotion	Google

In [39]: 1 data.emotion.value_counts()

Out[39]: No emotion toward brand or product 5375
Positive emotion 2970
Negative emotion 569
I can't tell 156
Name: emotion, dtype: int64

Words Distribution in Tweet Text

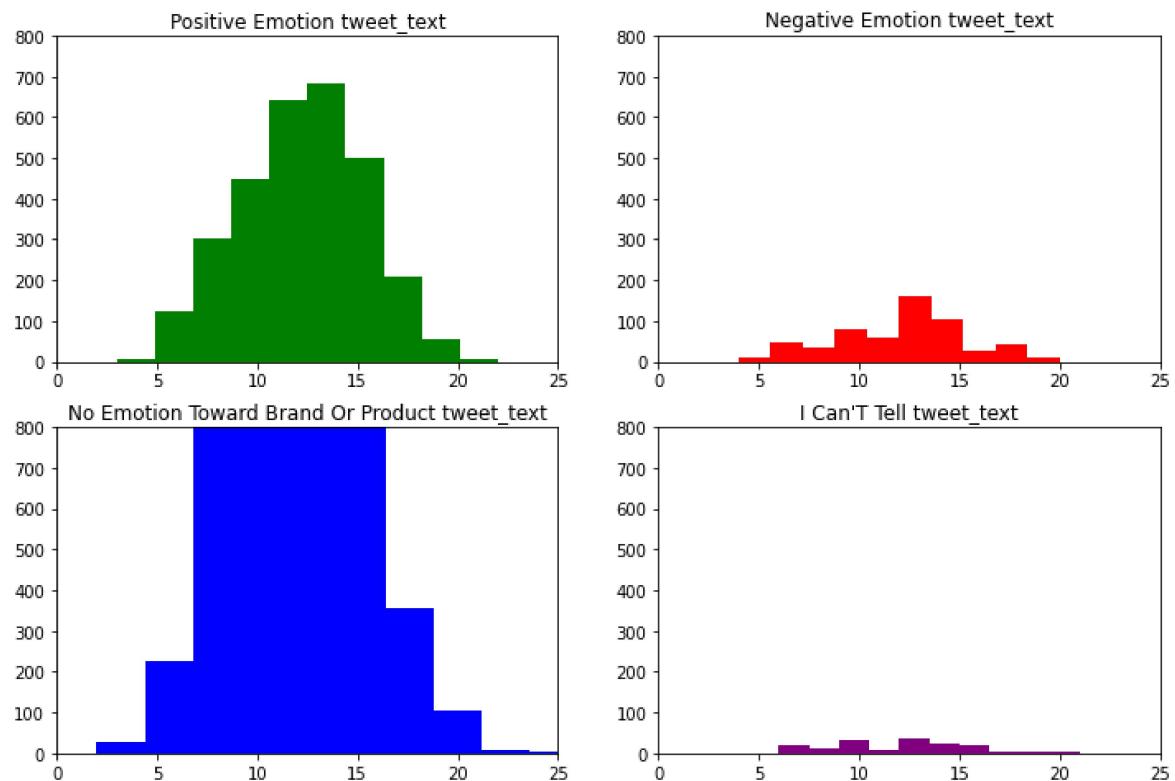
In [40]:

```

1 def generate_review_hist_plot(data, label, color, ax):
2     tweet_len = data[data['emotion']==label].tweet_text.str.split().map(len)
3     ax.set_xlim(0, 25)
4     ax.set_ylim(0, 800)
5     ax.hist(tweet_len, color=color)
6     ax.set_title(f'{label.title()} tweet_text')
7
8 fig, ax = plt.subplots(nrows=2, ncols=2, figsize=(12,8))
9 generate_review_hist_plot(data, 'Positive emotion', 'green',ax[0,0])
10 generate_review_hist_plot(data, 'Negative emotion', 'red',ax[0,1])
11 generate_review_hist_plot(data, 'No emotion toward brand or product', 'blue')
12 generate_review_hist_plot(data,"I can't tell",'purple',ax[1,1])
13
14 fig.suptitle("Words Distribution in Tweet Text", fontsize=18)
15
16 plt.show()

```

Words Distribution in Tweet Text



Interpretation

Majority of the tweets in all the target emotions classes were 10-15 words long.

Re-tokenizing

In [41]:

```
1 # re-tokenizing the data
2 data['tweet_text'] = data['tweet_text'].apply(tokenize_text)
3 data.head()
```

Out[41]:

	tweet_text	recipient	emotion	brand
0	[user, 3g, iphone, 3, hrs, tweeting, riseausti...]	iPhone	Negative emotion	Apple
1	[user, know, user, awesome, ipadiphone, app, l...]	iPad or iPhone App	Positive emotion	Apple
2	[user, wait, ipad, 2, also, sale, sxsw]	iPad	Positive emotion	Apple
3	[user, hope, years, festival, crashy, years, i...]	iPad or iPhone App	Negative emotion	Apple
4	[user, great, stuff, fri, sxsw, marissa, mayer...]	Google	Positive emotion	Google

In [42]:

```
1 data.tweet_text[0]
```

Out[42]:

```
['user',
 '3g',
 'iphone',
 '3',
 'hrs',
 'tweeting',
 'riseaustin',
 'dead',
 'need',
 'upgrade',
 'plugin',
 'stations',
 'sxsw']
```

Justification for Retokenizing

We retokenized since after removing the stopwords the tweets were no longer split as strings.

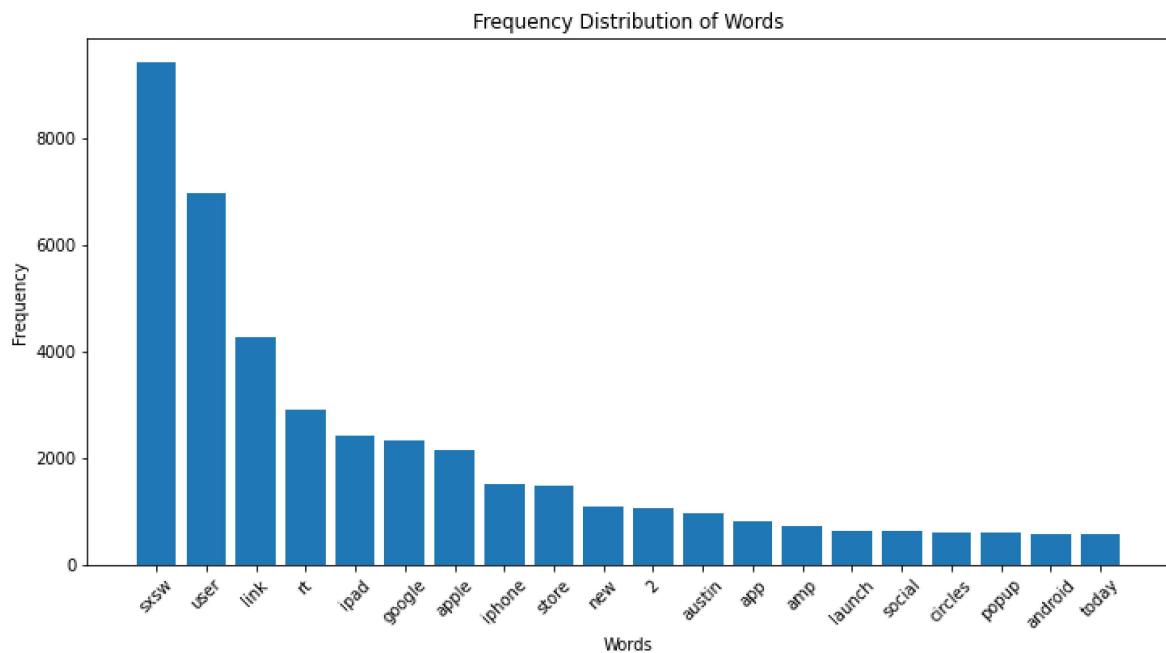
Frequency Distribution after Removing Stopwords

In [43]:

```

1 tweet_texts = data(tweet_text
2 # Flatten the list of tokens into a single list
3 all_tokens = [tweet for sublist in tweet_texts for tweet in sublist]
4
5 # Create the frequency distribution
6 freq_dist = FreqDist(all_tokens)
7
8 # Get the most common words
9 most_common = freq_dist.most_common(20) # Example: Top 20 most common wo
10
11 # Extract the words and frequencies
12 words, frequencies = zip(*most_common)
13
14 # Create a bar plot
15 plt.figure(figsize=(12, 6))
16 plt.bar(words, frequencies)
17 plt.xlabel('Words')
18 plt.ylabel('Frequency')
19 plt.title('Frequency Distribution of Words')
20 plt.xticks(rotation=45)
21 plt.show()

```



Observation

Just by looking at this, we can see there is a high occurrence of

- Venue (**south by south west**) specific words such as `sxsw`
- Twitter specific words(**re-tweet** such as `rt`
- Brand or product specific words such as `iphone`

These words probably occur across our various sentiments and may not necessarily add any

Lemmatization

```
In [44]: 1 # Lemmatization - example cats - cat
2 # Create Lemmatizer
3 word_lem = WordNetLemmatizer()
4
5 # function
6 def lem_words(tweet):
7     return [word_lem.lemmatize(word) for word in tweet]
```

```
In [45]: 1 data['tweet_text'] = data['tweet_text'].apply(lem_words)
2 data['tweet_text'].head()
```

```
Out[45]: 0 [user, 3g, iphone, 3, hr, tweeting, riseaustin...
1 [user, know, user, awesome, ipadiphone, app, l...
2 [user, wait, ipad, 2, also, sale, sxsw]
3 [user, hope, year, festival, crashy, year, iph...
4 [user, great, stuff, fri, sxsw, marissa, mayer...
Name: tweet_text, dtype: object
```

```
In [46]: 1 data.head()
```

```
Out[46]:
```

	tweet_text	recipient	emotion	brand
0	[user, 3g, iphone, 3, hr, tweeting, riseaustin...	iPhone	Negative emotion	Apple
1	[user, know, user, awesome, ipadiphone, app, l...	iPad or iPhone App	Positive emotion	Apple
2	[user, wait, ipad, 2, also, sale, sxsw]	iPad	Positive emotion	Apple
3	[user, hope, year, festival, crashy, year, iph...	iPad or iPhone App	Negative emotion	Apple
4	[user, great, stuff, fri, sxsw, marissa, mayer...	Google	Positive emotion	Google

Dealing with Emojis

We create a function that replaces emojis in tweet text with their corresponding meanings.

```
In [47]: 1 # Defining dictionary containing all emojis with their meanings.
2 emojis = {':)': 'smile', ':-)': 'smile', ';d': 'wink', ':-E': 'vampire',
3           ':-(': 'sad', ':-<': 'sad', ':P': 'raspberry', ':O': 'surprised'
4           ':-@': 'shocked', '@@': 'shocked', '-$': 'confused', ':\\\'': 'an
5           '#': 'mute', ':X': 'mute', '^)': 'smile', ':-&': 'confused',
6           '@@': 'eyeroll', ':-!': 'confused', ':-D': 'smile', ':-0': 'yel
7           '<(-_-)>': 'robot', 'd[-_-]b': 'dj', ":'-)": 'sadsmile', ')':
8           ';-)': 'wink', '0:-)': 'angel', '0*-)': 'angel', '(:-D': 'gossip'
```

In [48]:

```

1 #function to remove emojis and sequence of letters like heeeeey)
2 def process(tweets):
3
4     processed_tweet = []
5
6
7     # Defining regex patterns.
8
9     sequencePattern = r"(.)\1\1+" #matches three or more consecutive oc
10    seqReplacePattern = r"\1\1"
11
12    for tweet in tweets:
13
14        # Replace all emojis.
15        for emoji in emojis.keys():
16            tweet = tweet.replace(emoji, "EMOJI" + emojis[emoji])
17            # Replace 3 or more consecutive letters by 2 letter.
18            tweet = re.sub(sequencePattern, seqReplacePattern, tweet)
19
20        processed_tweet.append(tweet)
21
22    return processed_tweet

```

In [49]:

```

1 data['tweet_text'] = data['tweet_text'].apply(process)
2 data.head()

```

Out[49]:

	tweet_text	recipient	emotion	brand
0	[user, 3g, iphone, 3, hr, tweeting, riseaustin...]	iPhone	Negative emotion	Apple
1	[user, know, user, awesome, ipadiphone, app, l...	iPad or iPhone App	Positive emotion	Apple
2	[user, wait, ipad, 2, also, sale, sxsw]	iPad	Positive emotion	Apple
3	[user, hope, year, festival, crashy, year, iph...	iPad or iPhone App	Negative emotion	Apple
4	[user, great, stuff, fri, sxsw, marissa, mayer...	Google	Positive emotion	Google

In [50]:

```
1 data.tweet_text[0]
```

Out[50]:

```

['user',
 '3g',
 'iphone',
 '3',
 'hr',
 'tweeting',
 'riseaustin',
 'dead',
 'need',
 'upgrade',
 'plugin',
 'station',
 'sxsw']

```

vectorization

CountVectorization

Techniques is used for converting text documents into numerical representations

In [51]:

```

1 # create instance 1 vectorizer
2 bow = CountVectorizer()
3 #create Lemmatized_review column
4 data['lemmatized_review'] = data['tweet_text'].str.join(" ")
5 x = bow.fit_transform(data['lemmatized_review'])

```

In [52]:

```

1 #checking feature names
2 feature_name = bow.get_feature_names_out()
3 feature_name

```

Out[52]: array(['02', '03', '0310', ..., 'úómy', 'úóthe', 'úóuser'], dtype=object)

Feature engineering

Label Encoding the Target

Here we label encode the target feature to transform the values to numerical values.

In [53]:

```

1 #instantiate Labelencoder
2 le = LabelEncoder() # initializing the lib
3 #fit transform
4 data['emotion_code'] = le.fit_transform(data.emotion)
5 le.classes_ # viewing the classes

```

Out[53]: array(['I can't tell', 'Negative emotion',
'No emotion toward brand or product', 'Positive emotion'],
dtype=object)

Previewing to ensure encoding is present

In [54]:

```
1 # preview
2 data.head()
```

Out[54]:

	tweet_text	recipient	emotion	brand	lemmatized_review	emotion_code
0	[user, 3g, iphone, 3, hr, tweeting, riseaustin...]	iPhone	Negative emotion	Apple	user 3g iphone 3 hr tweeting riseaustin dead n...	1
1	[user, know, user, awesome, ipadiphone, app, l...	iPad or iPhone App	Positive emotion	Apple	user know user awesome ipadiphone app likely a...	3
2	[user, wait, ipad, 2, also, sale, sxsw]	iPad	Positive emotion	Apple	user wait ipad 2 also sale sxsw	3
3	[user, hope, year, festival, crashy, year, iph...	iPad or iPhone App	Negative emotion	Apple	user hope year festival crashy year iphone app...	1
4	[user, great, stuff, fri, sxsw, marissa, mayer...	Google	Positive emotion	Google	user great stuff fri sxsw marissa mayer google...	3

In [55]:

```
1 # displaying the encoding scheme
2 data[['emotion', 'emotion_code']].head(10)
```

Out[55]:

	emotion	emotion_code
0	Negative emotion	1
1	Positive emotion	3
2	Positive emotion	3
3	Negative emotion	1
4	Positive emotion	3
5	No emotion toward brand or product	2
7	Positive emotion	3
8	Positive emotion	3
9	Positive emotion	3
10	Positive emotion	3

6. Modeling

- The problem at hand is a classification problem.
- We will explore 2 models: a binary logistic regression model, a multi-class XGBoost model and MultinomialNB.
- Model accuracy will be the metric for evaluation.
- Justification: Accuracy to get a verdict if a tweet is positive or negative.

- Accuracy of 70% will be the threshold to deem the model as successful.

6.1 Binary classification

- In this section we create a base model to identify if a tweet is 'Positive' or 'Negative'.
- LogisticRegression will be used for the classification.
- The normal preprocessing of vectorization and train test split will be implemented.

In [56]:

```
1 # creating a copy of the original data so as to remove unwanted rows i.e.  
2 data_copy = data.copy()
```

In [57]:

```
1 # Define the values to drop  
2 values_to_drop = [0, 2]  
3  
4 # Drop rows that have the values in column 'B'  
5 data_copy = data_copy[~data_copy['emotion_code'].isin(values_to_drop)]
```

In [58]:

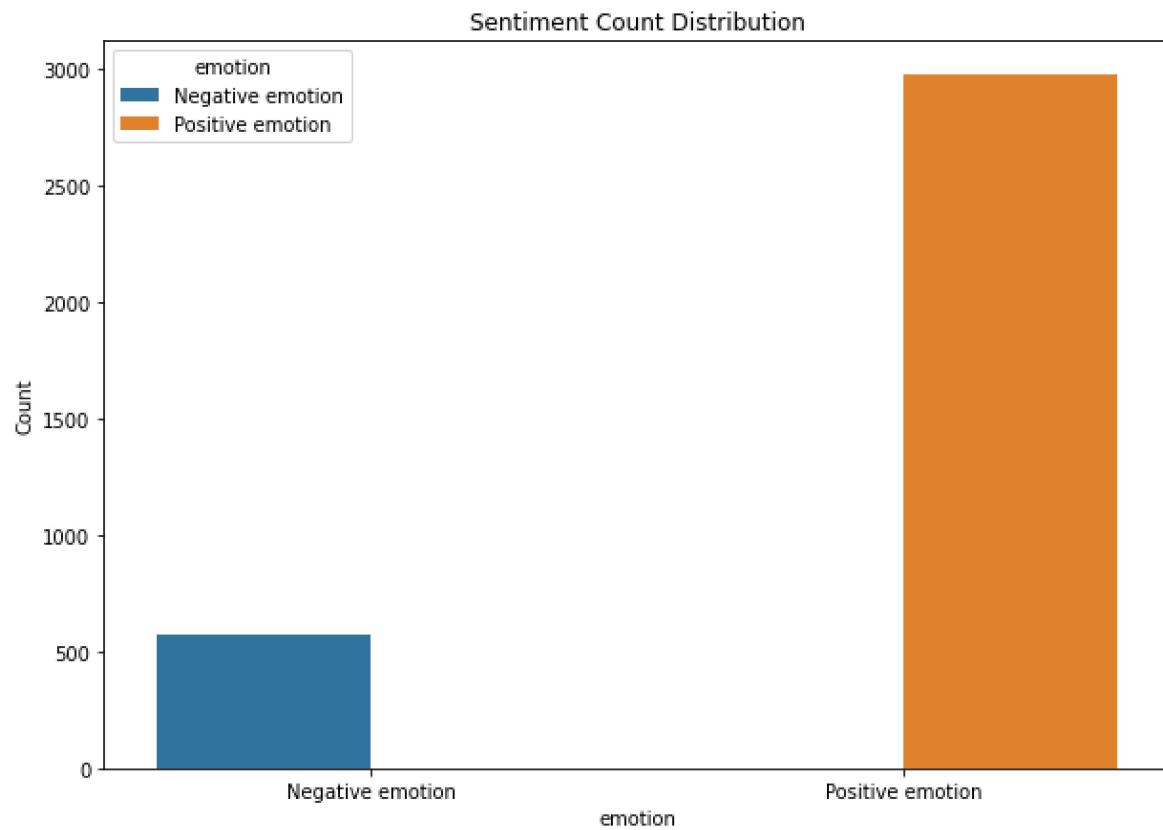
```
1 # value counts  
2 data_copy['emotion_code'].value_counts()
```

Out[58]:

```
3    2970  
1    569  
Name: emotion_code, dtype: int64
```

In [59]:

```
1 # lets see see binary sentiment distribution emotion-wise
2 emotion_group = data_copy.groupby('emotion')['emotion_code'].value_counts
3 # Calculate the ascending order of states based on churn count
4 ascending_order = emotion_group.groupby('emotion')['count'].sum().sort_va
5
6 plt.figure(figsize=(10,7))
7 sns.barplot(x='emotion', y='count', hue='emotion', data=emotion_group, or
8 plt.xlabel('emotion')
9 plt.ylabel('Count')
10 plt.title('Sentiment Count Distribution')
11 plt.show()
```



Intepretation

- From the graph above it is clear that most of the sentiments were positive in nature

Vectorizing Data Copy

CountVectorizer

```
In [60]: 1 # create instance of vectorizer
2 bow2 = CountVectorizer()
3 data_copy['lemmatized_review'] = data_copy['tweet_text'].str.join(" ")
4 # our X
5 x2 = bow2.fit_transform(data_copy['lemmatized_review'])
```

```
In [61]: 1 # preview
2 data_copy.head()
```

Out[61]:

	tweet_text	recipient	emotion	brand	lemmatized_review	emotion_code
0	[user, 3g, iphone, 3, hr, tweeting, riseaustin...]	iPhone	Negative emotion	Apple	user 3g iphone 3 hr tweeting riseaustin dead n...	1
1	[user, know, user, awesome, ipadiphone, app, l...]	iPad or iPhone App	Positive emotion	Apple	user know user awesome ipadiphone app likely a...	3
2	[user, wait, ipad, 2, also, sale, sxsw]	iPad	Positive emotion	Apple	user wait ipad 2 also sale sxsw	3
3	[user, hope, year, festival, crashy, year, iph...]	iPad or iPhone App	Negative emotion	Apple	user hope year festival crashy year iphone app...	1
4	[user, great, stuff, fri, sxsw, marissa, mayer...]	Google	Positive emotion	Google	user great stuff fri sxsw marissa mayer google...	3

Train-test split

```
In [62]: 1 #x = bow.fit_transform(data['Lemmatized_review'])
2 # our y
3 y = data_copy['emotion_code']
```

```
In [63]: 1 # train-test split
2 X_train, X_test, y_train, y_test = train_test_split(x2, y, test_size=0.2,
```

```
In [64]: 1 # confirming shape
2
3 X_train.shape, y_train.shape
```

Out[64]: ((2831, 5871), (2831,))

Fitting logistic regression model on split data

```
In [65]: 1 # instantiating model 1
          2 log_reg = LogisticRegression(random_state=42)
```

```
In [66]: 1 # fitting
          2 log_reg.fit(X_train, y_train)
```

Out[66]:

```
LogisticRegression
LogisticRegression(random_state=42)
```

```
In [67]: 1 # predicting
          2 y_pred = log_reg.predict(X_test)
          3 y_train_pred = log_reg.predict(X_train)
```

Evaluation

```
In [68]: 1 #Accuracy score
          2 print('train',accuracy_score(y_train, y_train_pred))
          3 print('test',accuracy_score(y_test, y_pred))
```

```
train 0.9699752737548569
test 0.903954802259887
```

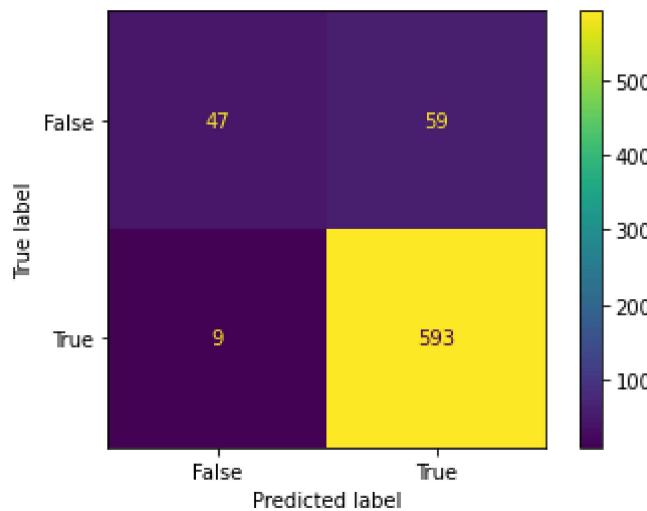
```
In [69]: 1 # Classification metric Report
          2 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
1	0.84	0.44	0.58	106
3	0.91	0.99	0.95	602
accuracy			0.90	708
macro avg	0.87	0.71	0.76	708
weighted avg	0.90	0.90	0.89	708

```
In [70]: 1 #Confusion Matrix
          2 confusion_matrix = confusion_matrix(y_test, y_pred)
```

In [71]:

```
1 #Display Confusion Matrix
2 cm_display = metrics.ConfusionMatrixDisplay(confusion_matrix = confusion_
3 cm_display.plot()
4 plt.show()
```



Machine Learning Communication

Rationale why modeling was implemented.

- While simpler forms of data analysis, such as descriptive statistics or basic data visualization, can provide initial insights, they may not be sufficient for complex problems or large datasets. Machine learning leverages advanced algorithms to uncover hidden patterns.

Results.

- Training accuracy: 96%
- Testing accuracy: 90%

- The accuracy means that the model can predict with an accuracy of 90% whether a tweet is positive or negative.
- The current model is fit for prediction since it is generalizing well to new data even with high accuracy.

Limitations of binary model.

- Not fit for multiclass datasets.

Testing Logistic Regression Model here

- We take a sample sentence, preprocess it then pass it to the model to see the output.

```
In [72]: 1 # tweet
          2 review = 'Awwww!! the google was sooo sweeet!! I loved it'
          3 #review = 'i hate apple'
```

```
In [73]: 1 # Lowercase
          2 review = review.lower()
          3 review
```

Out[73]: 'awwww!! the google was sooo sweeet!! i loved it'

```
In [74]: 1 # removing punctuations
          2 review = remove_punctuations(review)
          3 review
```

Out[74]: 'awwww the google was sooo sweeet i loved it'

```
In [75]: 1 # tokenizing the data
          2 review = tokenize_text(review)
          3 review
```

Out[75]: ['awwww', 'the', 'google', 'was', 'sooo', 'sweeet', 'i', 'loved', 'it']

```
In [76]: 1 # removing stopwords
          2 review = remove_stopwords(review)
          3 review
```

Out[76]: 'awwww google sooo sweeet loved'

```
In [77]: 1 review = tokenize_text(review)
          2 review
```

Out[77]: ['awwww', 'google', 'sooo', 'sweeet', 'loved']

```
In [78]: 1 # Lemmatization
          2 review = lem_words(review)
          3 review
```

Out[78]: ['awwww', 'google', 'sooo', 'sweeet', 'loved']

```
In [79]: 1 # joining to one sentence
          2 review = ' '.join(review)
          3
          4 review
```

Out[79]: 'awwww google sooo sweeet loved'

```
In [80]: 1 # putting to list
          2 review = [review]
          3 review
```

Out[80]: ['awwww google sooo sweeet loved']

```
In [81]: 1 # vectorizing
          2 review_x = bow2.transform(review)
          3 review_x
```

Out[81]: <1x5871 sparse matrix of type '<class 'numpy.int64'>'
with 2 stored elements in Compressed Sparse Row format>

```
In [82]: 1 test_predict = log_reg.predict(review_x)
          2 test_predict
```

Out[82]: array([3])

6.2 Multiclass Classifier

- Here we work with the original dataset.
- We build a multi-class classifier.
- MultinomialNB and XGBoost model will be tested.

MultinomialNB model

```
In [83]: 1 # target variable
          2 y = data['emotion_code']
```

```
In [84]: 1 # splitting the data
          2 X_train_multi, X_test_multi, y_train_multi, y_test_multi = train_test_spl
```

```
In [85]: 1 X_train_multi.shape, y_train_multi.shape
```

Out[85]: ((7256, 9867), (7256,))

```
In [86]: 1 # instantiating the multi class model
2 multi_class = MultinomialNB()
3 # fitting the model
4 multi_class.fit(X_train_multi, y_train_multi)
```

Out[86]:

```
▼ MultinomialNB
  MultinomialNB()
```

```
In [87]: 1 # predicting multi-class
2 y_pred_multi = multi_class.predict(X_test_multi)
3 y_train_pred_multi = multi_class.predict(X_train_multi)
```

```
In [88]: 1 # checking the accuracy
2 print('train',accuracy_score(y_train_multi, y_train_pred_multi))
3 print('test',accuracy_score(y_test_multi, y_pred_multi))
```

train 0.8036108048511577
test 0.6543550165380375

Machine Learning Communication - MultiNomialmodel.

Rationale why MultiNomial modeling was implemented.

- While simpler forms of modeling can do , they may not be able to work with multi-class datasets.

Results.

- Accuracy on the training set: 80%
- Accuracy on the testing set: 66%
- The model is overfitting

Limitations of multinomial model.

- The current model is not fit for prediction since it is not generalizing well to new data even with high accuracy. The model is overfitting.

XGBoost model

```
In [89]: 1 # Create an XGBoost classifier instance
2 xgb_classifier = xgb.XGBClassifier()
3
4 # Fit the model on the training data
5 xgb_classifier.fit(X_train_multi, y_train_multi)
6
```

Out[89]:

```
XGBClassifier
XGBClassifier(base_score=0.5, booster='gbtree', callbacks=None,
              colsample_bylevel=1, colsample_bynode=1, colsample_bytree=
1,
              early_stopping_rounds=None, enable_categorical=False,
              eval_metric=None, gamma=0, gpu_id=-1, grow_policy='depthwi
se',
              importance_type=None, interaction_constraints='',
              learning_rate=0.300000012, max_bin=256, max_cat_to_onehot=
4,
              max_delta_step=0, max_depth=6, max_leaves=0, min_child_wei
```

```
In [90]: 1 # Make predictions on the testing data
2 #predictions = xgb_classifier.predict(test_x)
3 y_pred_xgb = xgb_classifier.predict(X_test_multi)
4 y_train_pred_xgb = xgb_classifier.predict(X_train_multi)
```

```
In [91]: 1 # checking the accuracy
2 print('train',accuracy_score(y_train_multi, y_train_pred_xgb))
3 print('test',accuracy_score(y_test_multi, y_pred_xgb))
```

train 0.7848676957001103
 test 0.685226019845645

Machine Learning Communication XGBoost.

Rationale why ensemble modeling was implemented.

- While simpler forms of data analysis, such as descriptive statistics or basic data visualization, can provide initial insights, they are not sufficient for complex problems or large datasets such as this one. Ensemble models leverages advanced algorithms to uncover hidden patterns, make accurate predictions.

Results.

- Accuracy on the training set: 78%
- Accuracy on the testing set: 68%
- The accuracy means that the model can predict with an accuracy of 78% whether a tweet falls within the specified labels.

Limitations.

- The current model is not fit for prediction since it is overfitting.
- This we see from the difference between train and test accuracy.

Resampling to handle the class imbalance issue.

We try to battle overfitting by handling imbalance issues in the classes, then modelling again.

Trying Resampling with xgboost

In [92]:

```

1 # Creating an instance of the resampler
2 resampler = RandomOverSampler() # or RandomUnderSampler()
3
4 # Perform resampling on the training data
5 resampled_features, resampled_labels = resampler.fit_resample(X_train_multi)
6
7 # Create an XGBoost classifier instance
8 xgb_classifier = xgb.XGBClassifier()
9
10 # Fit the model on the resampled training data
11 xgb_classifier.fit(resampled_features, resampled_labels)
12
13 # Make predictions on the testing data
14 predictions = xgb_classifier.predict(X_test_multi)
15
16 # Evaluate the model
17 accuracy = accuracy_score(y_test_multi, predictions)
18 print("Accuracy:", accuracy)

```

Accuracy: 0.6278941565600882

ML Communication

Result

The resampled model has an accuracy of 61% which is not an improvement over prior models.

Limitations

The resampled model is not fit for modelling since its modelling accuracy is low.

Final note on modeling

- Despite implementing resampling techniques to address class imbalance, the model's accuracy has not improved significantly. This suggests that other factors within the dataset may be limiting predictive performance.

- The binary logistic regression model performs best.
- For Multiclass XGBoost was better

Limitation and Challenges

- Class Imbalance Issue: The dataset suffers from class imbalance, where one sentiment class is dominant while others are underrepresented. This can result in biased models that are more accurate for the majority class but perform poorly on the minority classes. Addressing this issue is important to ensure fair and balanced sentiment analysis.
- Limited Dataset Size: The dataset used for sentiment analysis is relatively small, which can limit the model's ability to capture the full complexity of sentiments expressed in text. A larger and more diverse dataset would provide a broader representation of sentiments and improve the model's performance and generalization.
- Language Ambiguity and Sarcasm Detection: Language can be inherently ambiguous, and detecting sarcasm in text adds an extra layer of complexity. Sarcasm detection is challenging due to the subtleties and nuances involved. Developing robust strategies to handle language ambiguity and detect sarcasm is crucial for accurate sentiment analysis

7. Findings

1. Most tweets were directed to None brand category. This may indicate that customers were not engaging with the brand.
2. Positive sentiments had the highest count compared to Negative sentiments, indicating that most people in general liked respective brands(Google and Apple)
3. Most of the positive tweets were directed to Apple brands
4. In the field of sentiment analysis, one of the significant challenges is dealing with language ambiguity and sarcasm detection. Natural language is complex and often subjective, making it difficult to accurately interpret sentiments from text.
5. On average most of the tweets were 10-15 words long - more words increase ambiguity.
6. NLP is a difficult task to model accurately.

8. Recommendations

We recommend that there be more customer engagement.

Probably check on this areas;

- Churn ratio - rate at which customers discontinue their relationship with a product company within a given time period
- Social media influencers through brand or product endorsement
- Customer feedback - The brands can introduce a rating system to accurately capture the sentiments of their customers

9. Next steps

- In our future work, we plan to explore advanced techniques such as incorporating attention mechanisms, using ensemble methods to further enhance the model's performance by incorporating domain-specific and fine-tuning the model on industry-specific datasets could improve its accuracy and adaptability.
- By considering these evaluation metrics, addressing limitations, and planning for future improvements, we aim to develop a robust NLP sentiment analysis solution that effectively captures sentiment

In [93]:

```
1 import joblib
2
3 # Save the trained model as a pickle file
4 joblib.dump(log_reg, 'logistic_regression_model.pkl')
5 joblib.dump(xgb_classifier, 'xgb.pkl')
```

Out[93]: ['xgb.pkl']