

Table of Contents

- 1 Phase 4 Project
 - 1.1 Business Objective
 - 1.2 Methodology
- 2 Dataset
- 3 Analysis
 - 3.1 Pre-processing
- 4 Trending topic on Twitter
 - 4.1 Wordcloud
 - 4.2 FreqDist
 - 4.3 Removing *stopwords*
- 5 What is the most popular product?
- 6 What do customers like/dislike in a product?
 - 6.1 Ipad
 - 6.2 Apple
 - 6.3 iPhone
 - 6.4 Ipad and iPhone apps
 - 6.5 Google
 - 6.6 Android
- 7 Model to predict company from tweets
 - 7.1 Vectorizers
 - 7.1.1 CountVectorizer
 - 7.2 Tf-IDF Vectorizer
- 8 Tuning LogisticRegression with Countvectorizer
 - 8.1 min_df and max_df values
 - 8.2 n-gram
 - 8.3 Stemming using PorterStemmer
- 9 Word2Vec
 - 9.1 Experimentation
 - 9.2 Mean Embeddings
 - 9.3 Modelling
- 10 Next Steps

Phase 4 Project

Business Objective

To help Acme Online, an online electronics store, analyze customer tweets from their Twitter page about Apple and Google products. The result of this analysis will be used to find out which

company's product has more favourable reviews and the reasons behind it - this will help Acme Online adjust their inventory accordingly.

Methodology

1. Analyze tweets to check what customers are talking about.
2. Analyze tweets to identify the most popular product - pts 1&2 can be used to tweak Acme Online's inventory accordingly.
3. For each product, we will look to see what customers like/dislike to identify opportunities for improvement, if applicable.
4. Since human intervention was used identify products based on tweets, we will attempt to build a model using NLP to automate this. We will use the f1-score for model evaluations since minimizing False Positive and False Negatives is desirable .

Dataset

Dataset sourced from CrowdFlower via data.world: <https://data.world/crowdflower/brands-and-product-emotions>

Analysis

```
In [1]: #import relevant libraries
import pandas as pd
import matplotlib.pyplot as plt
from wordcloud import WordCloud
from nltk.stem import PorterStemmer
import nltk
from nltk.corpus import stopwords
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import MultinomialNB
from sklearn.ensemble import RandomForestClassifier
from sklearn import svm
from sklearn.metrics import classification_report, plot_confusion_matrix
from sklearn.metrics import f1_score, accuracy_score
from sklearn.pipeline import Pipeline
import numpy as np
from nltk import word_tokenize
from gensim.models import Word2Vec
from nltk.tokenize import RegexpTokenizer
from nltk import FreqDist
import warnings
warnings.filterwarnings('ignore')
```

Importing the dataset:

```
In [2]: df= pd.read_csv('tweets.csv', encoding='unicode escape',)
df.head()
```

```
Out[2]:      tweet_text  emotion_in_tweet_is_directed_at  is_there_an_emotion_directed_at_a_brand_or_product
```

	tweet_text	emotion_in_tweet_is_directed_at	is_there_an_emotion_directed_at_a_brand_or_product
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	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	Negative emotion
4	@sxtxstate great stuff on Fri #SXSW: Marissa M...	Google	Positive emotion

Pre-processing

To save ourselves from lot's of keystrokes, let's rename the columns:

```
In [3]: #renaming the columns to make it less cumbersome
df.rename(columns={'emotion_in_tweet_is_directed_at':'product_service',
                  'is_there_an_emotion_directed_at_a_brand_or_product':'emotion'},inpl
df.head()
```

```
Out[3]:
```

	tweet_text	product_service	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	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	Negative emotion
4	@sxtxstate great stuff on Fri #SXSW: Marissa M...	Google	Positive emotion

```
In [4]: #getting some info
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9093 entries, 0 to 9092
Data columns (total 3 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   tweet_text      9092 non-null   object
1   product_service  3291 non-null   object
2   emotion         9093 non-null   object
```

```
dtypes: object(3)
memory usage: 213.2+ KB
```

There are no numeric values in our df which is what we'd expect given that we're analyzing tweets.

```
In [5]: #checking for null values
df.isna().sum()
```

```
Out[5]: tweet_text      1
product_service  5802
emotion         0
dtype: int64
```

From the above, we can see that the `product_service` column has a large number of missing values; more than 50%. Let's leave it for now and remove the one empty row in `tweet_text`

```
In [6]: #removing the null value in the tweet_text column
df = df[df['tweet_text'].notnull()]
```

Let's take a look at the `product_service` column to see the different kinds of products that are involved:

```
In [7]: #examining the product_service column
df['product_service'].value_counts()
```

```
Out[7]: 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: product_service, dtype: int64
```

Let's group some of the categories to facilitate easier analysis:

```
In [8]: #Let's group product/services that resemble each other for both brands. This will make :

df['product_service'].replace('Other Google product or service', 'Google', inplace=True)
df['product_service'].replace('Other Apple product or service', 'Apple', inplace=True)
df['product_service'].replace('Android App', 'Android', inplace=True)
df['product_service'].fillna('Not Applicable', inplace=True)

#checking
df['product_service'].value_counts()
```

```
Out[8]: Not Applicable  5801
iPad          946
Google        723
Apple         696
iPad or iPhone App  470
iPhone        297
Android       159
Name: product_service, dtype: int64
```

Let's apply some of the common pre-processing steps when it comes to working with text data:

1. Remove capitalization
2. Remove punctuations and special characters

3. Tokenizing

```
In [9]: # Removing capitalization
df['tweet_text'] = df['tweet_text'].str.lower()

#removing punctuations using the default pattern in sklearn and tokenizing
basic_token_pattern = r"(?u)\b\w\w+\b"
tokenizer = RegexpTokenizer(basic_token_pattern)

#applying the tokenizer to the df and creating a new column
df['text_token'] = df['tweet_text'].apply(tokenizer.tokenize)
df.head(10)
```

```
Out[9]:
```

	tweet_text	product_service	emotion	text_token
0	.@wesley83 i have a 3g iphone. after 3 hrs twe...	iPhone	Negative emotion	[wesley83, have, 3g, iphone, after, hrs, tweet...
1	@jessedee know about @fludapp ? awesome ipad/i...	iPad or iPhone App	Positive emotion	[jessedee, know, about, fludapp, awesome, ipad...
2	@swonderlin can not wait for #ipad 2 also. the...	iPad	Positive emotion	[swonderlin, can, not, wait, for, ipad, also, ...
3	@sxsw i hope this year's festival isn't as cra...	iPad or iPhone App	Negative emotion	[sxsw, hope, this, year, festival, isn, as, cr...
4	@sxtxstate great stuff on fri #sxsw: marissa m...	Google	Positive emotion	[sxtxstate, great, stuff, on, fri, sxsw, maris...
5	@teachntech00 new ipad apps for #speechtherapy...	Not Applicable	No emotion toward brand or product	[teachntech00, new, ipad, apps, for, speechthe...
7	#sxsw is just starting, #ctia is around the co...	Android	Positive emotion	[sxsw, is, just, starting, ctia, is, around, t...
8	beautifully smart and simple idea rt @madebyma...	iPad or iPhone App	Positive emotion	[beautifully, smart, and, simple, idea, rt, ma...
9	counting down the days to #sxsw plus strong ca...	Apple	Positive emotion	[counting, down, the, days, to, sxsw, plus, st...
10	excited to meet the @samsungmobileus at #sxsw ...	Android	Positive emotion	[excited, to, meet, the, samsungmobileus, at, ...

Trending topic on Twitter

By answering this question, we can understand what customers are tweeting about. We can visualize this using a feature called **WordCloud**.

Wordcloud

```
In [10]: #importing the stopwords list to pass onto the WC generator
stopwords_list = stopwords.words('english')
stopwords_list.append('mention')

#dropping null values
df.dropna(inplace=True)
```


give very little information about the contents of the text itself. Since they are most likely to occur a lot more than nouns and adjectives, they also distort our analysis and are hence best removed from the corpus for analysis.

```
In [11]: #defining a function to remove stopwords
def remove_stopwords(token_list):
    stopwords_removed = [token for token in token_list if token not in stopwords_list]
    return stopwords_removed

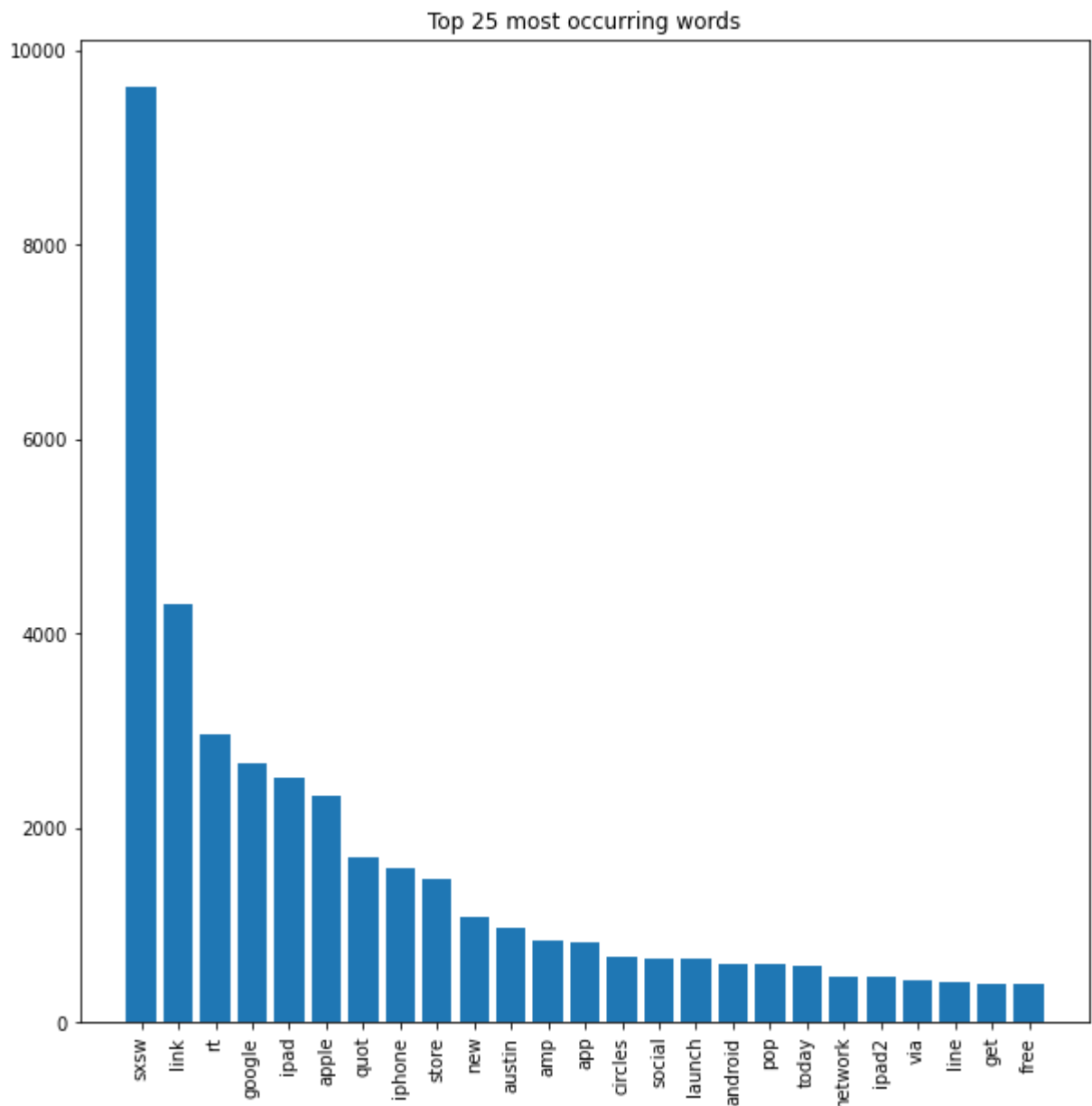
#applying the function to the text_token column
df['text_token'] = df['text_token'].apply(remove_stopwords)
```

```
In [12]: #defining a function to plot the top_25 most occurring words

def plot_freq_dist(words):
    freq_dist = FreqDist(df[words].explode())
    # listing out the top 25 most occurring words and ther respective counts
    top_25 = list(zip(*freq_dist.most_common(25)))

    #creating a plot of the top_25 words
    fig,ax=plt.subplots(figsize=(10,10))
    ax.bar(top_25[0],top_25[1])
    ax.set_title('Top 25 most occurring words')
    ax.tick_params(axis='x', rotation=90)
```

```
In [13]: #plotting the freq_dist
plot_freq_dist('text_token')
```



We can see from the above that the words `SXSW`, `Google`, `iPad` are some of the most tweeted words. A google search of `SXSW` reveals it to be arts and music festival held in Austin,TX. Hence, we can reasonably conclude that tweets collected for the analysis was from the city of Austin,TX and also coincided when the festival was running. It is also quite possible that people were streaming it on their iPads with great success!

What is the most popular product?

This is to answer the first question : What are the emotional responses for each product_service category listed in the data? For eg: for the category 'Apple' how many positive,negative and neutral responses are there?

By comparing the responses for each category, we can gauge customer sentiment

Pivot Tables can help better organize the data for the analysis

In [14]: *#creating a pivot table to organize the data*


```
df_pivot = df.pivot_table(index='product_service',aggfunc='count',columns='emotion')
df_pivot
```

Out[14]:

	text_token							tweet_text	
	emotion	I can't tell	Negative emotion	No emotion toward brand or product	Positive emotion	I can't tell	Negative emotion	No emotion toward brand or product	Positive emotion
product_service									
Android	NaN	16.0	2.0	141.0	NaN	16.0	2.0	141.0	
Apple	2.0	97.0	22.0	575.0	2.0	97.0	22.0	575.0	
Google	2.0	115.0	24.0	582.0	2.0	115.0	24.0	582.0	
Not Applicable	147.0	51.0	5297.0	306.0	147.0	51.0	5297.0	306.0	
iPad	4.0	125.0	24.0	793.0	4.0	125.0	24.0	793.0	
iPad or iPhone App	NaN	63.0	10.0	397.0	NaN	63.0	10.0	397.0	
iPhone	1.0	103.0	9.0	184.0	1.0	103.0	9.0	184.0	

In [15]:

```
df_pivot.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 7 entries, Android to iPhone
Data columns (total 8 columns):
#   Column                                                                 Non-Null Count  Dtype
---  -
0   (text_token, I can't tell)                                           5 non-null     float64
1   (text_token, Negative emotion)                                       7 non-null     float64
2   (text_token, No emotion toward brand or product)                   7 non-null     float64
3   (text_token, Positive emotion)                                       7 non-null     float64
4   (tweet_text, I can't tell)                                           5 non-null     float64
5   (tweet_text, Negative emotion)                                       7 non-null     float64
6   (tweet_text, No emotion toward brand or product)                   7 non-null     float64
7   (tweet_text, Positive emotion)                                       7 non-null     float64
dtypes: float64(8)
memory usage: 504.0+ bytes
```

Since the text_column is not relevant right now, let's remove it:

In [16]:

```
#dropping the columns
df_pivot.drop(df_pivot.columns[[0,1,2,3]],axis=1,inplace=True)
```

Renaming the columns for better understanding:

In [17]:

```
#renaming the columns
df_pivot.columns = ["I can't tell",'Negative emotion','No emotion toward brand or produ
df_pivot
```

Out[17]:

	I can't tell	Negative emotion	No emotion toward brand or product	Positive emotion
product_service				
Android	NaN	16.0	2.0	141.0

	I can't tell	Negative emotion	No emotion toward brand or product	Positive emotion
product_service				
Apple	2.0	97.0	22.0	575.0
Google	2.0	115.0	24.0	582.0
Not Applicable	147.0	51.0	5297.0	306.0
iPad	4.0	125.0	24.0	793.0
iPad or iPhone App	NaN	63.0	10.0	397.0
iPhone	1.0	103.0	9.0	184.0

```
In [18]: #dropping 'Not Applicable' since it is not relevant here
df_pivot.drop('Not Applicable',axis=0,inplace=True)

#rearranging the columns for better visualization
df_pivot=df_pivot[['Positive emotion','Negative emotion', 'No emotion toward brand or p

#sorting the values for better visuzalization
df_pivot.sort_values('Positive emotion',ascending=False,inplace=True)

df_pivot
```

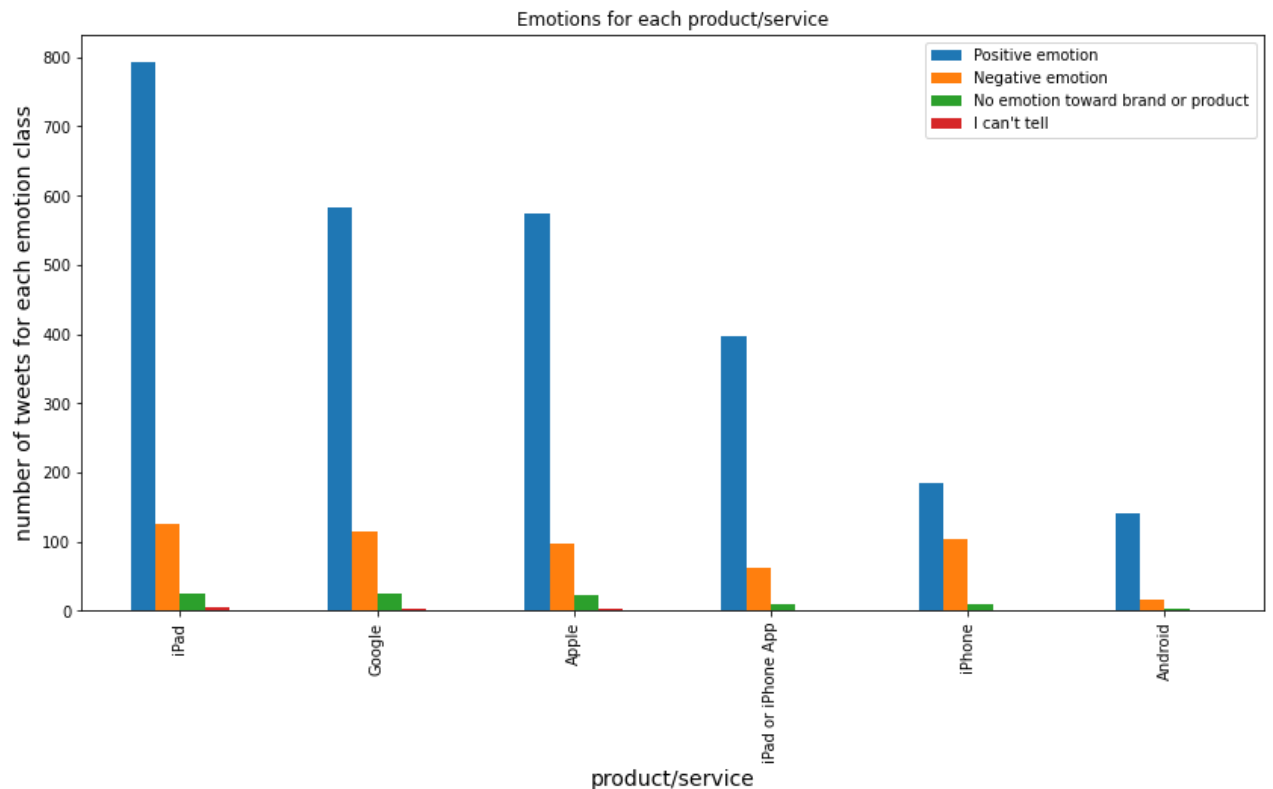
```
Out[18]:
```

	Positive emotion	Negative emotion	No emotion toward brand or product	I can't tell
product_service				
iPad	793.0	125.0	24.0	4.0
Google	582.0	115.0	24.0	2.0
Apple	575.0	97.0	22.0	2.0
iPad or iPhone App	397.0	63.0	10.0	NaN
iPhone	184.0	103.0	9.0	1.0
Android	141.0	16.0	2.0	NaN

Now that we've formatted the table to our liking, let's plot a bar chart see the different emotions for each product. This will give us an idea about how customers feel about each product

```
In [19]: # bar chart listing emotion class for each product_service

df_pivot.plot(kind='bar',figsize=(14,7));
plt.title('Emotions for each product/service');
plt.ylabel('number of tweets for each emotion class',fontsize=14);
plt.xlabel('product/service',fontsize=14);
```



We have a clear winner in **iPad!** i.e. the iPad is the most popular product among customers and Android the least. We can also see that the number of negative tweets seem to be somewhat level across all products except for Android.

What do customers like/dislike in a product?

Here, we are looking to answer the second question. We can do this by breaking down for each product, the different emotions to see if there are any key words that stand out. For eg: we can list out tweets by positive and negative emotions for iPad and analyze separately to gauge sentiment.

```
In [20]: #updating stopwords list to include SXSX since it appears nearly 10,000 times
stopwords_list.append('SXSX')
```

Let's define some functions to make things easier:

Function to get only **positive tweets**

```
In [21]: def get_positive(df, category, emotion):
positive_df = df.loc[(df['product_service'] == category) & (df['emotion'] == 'Positive')]
return positive_df
```

Function to get only **negative tweets**

```
In [22]: def get_negative(df, category, emotion):
negative_df = df.loc[(df['product_service'] == category) & (df['emotion'] == 'Negative')]
return negative_df
```

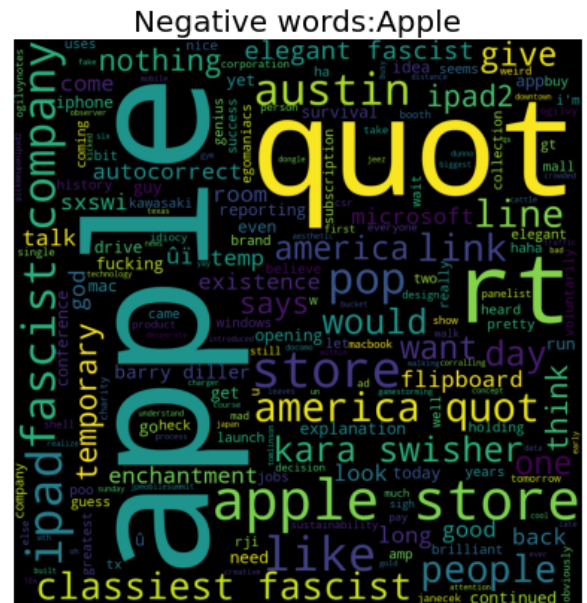
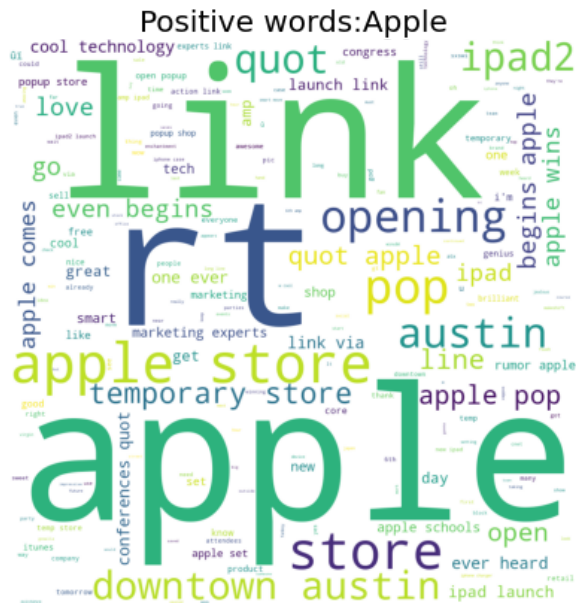
Function to generate **Wordclouds for positive and negative tweets**

```
In [23]: def wordcloud_gen(df, category):
```


button are some of the words that feature prominently thus illustrating displeasure of the users regarding some of the features of the iPad . iPad2 is also mentioned quite a lot.

Apple

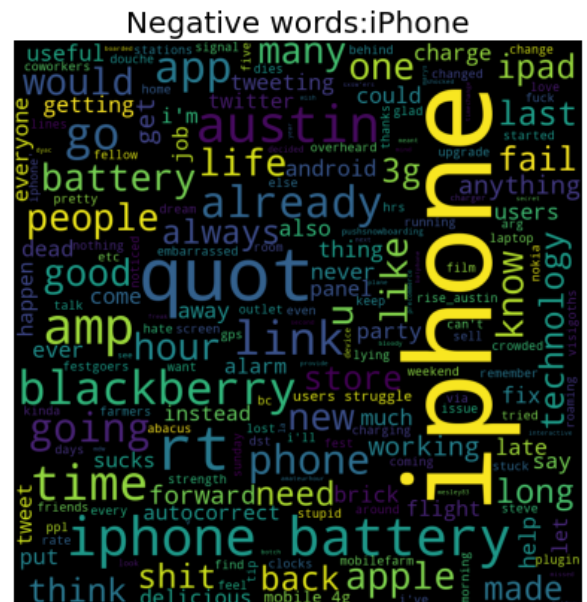
```
In [25]: wordcloud_gen(df, 'Apple')
```



Sentiment against Apple seems to be quite severe given the high number of tweets featuring the word `fascist`!

iPhone

```
In [26]: wordcloud_gen(df, 'iPhone')
```

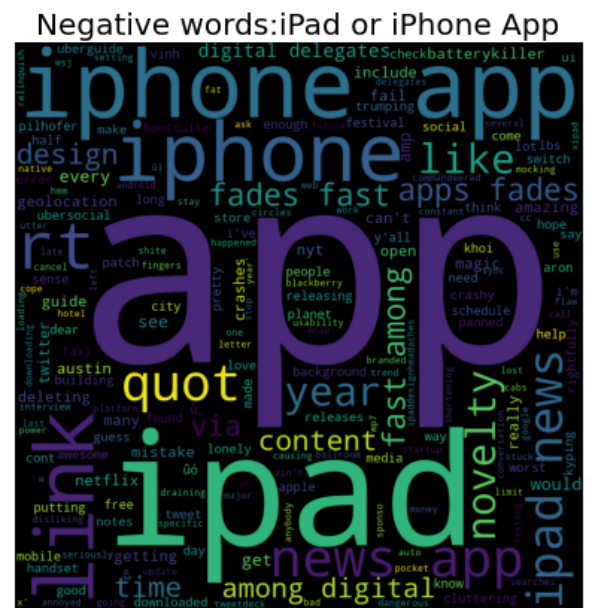
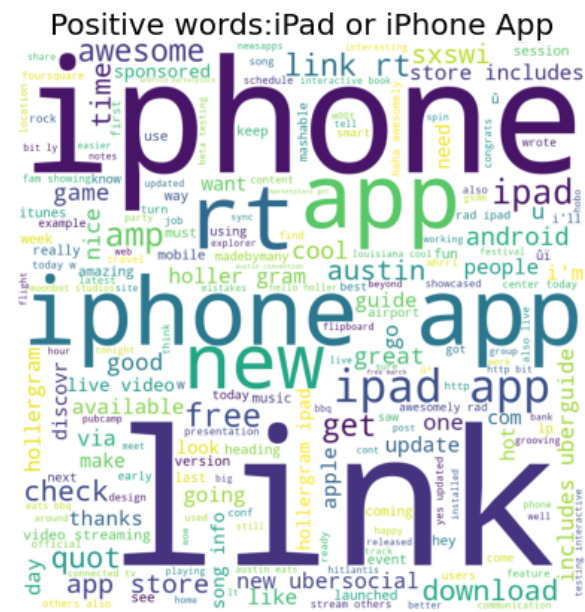


From positive emotions, verizon stands out suggesting their superiority from the other carriers. iphone battery, battery from the negative emotions illustrate unequivocally where the

problem lies with the iPhone.

Ipad and iPhone apps

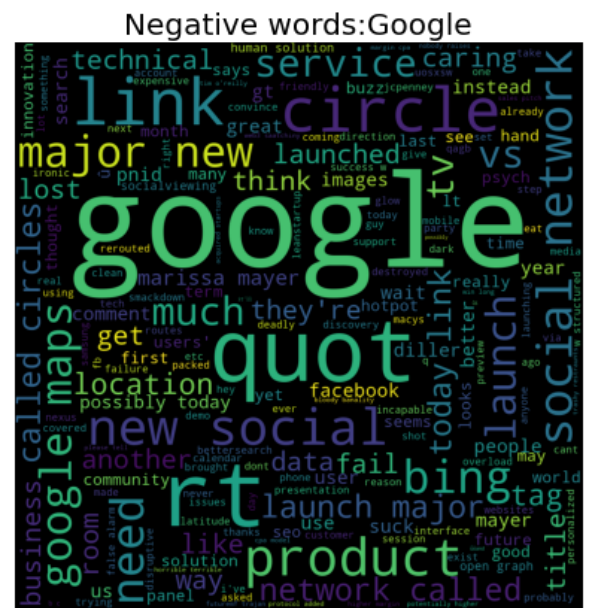
```
In [27]: wordcloud_gen(df, 'iPad or iPhone App')
```



Nothing really stands out here

Google

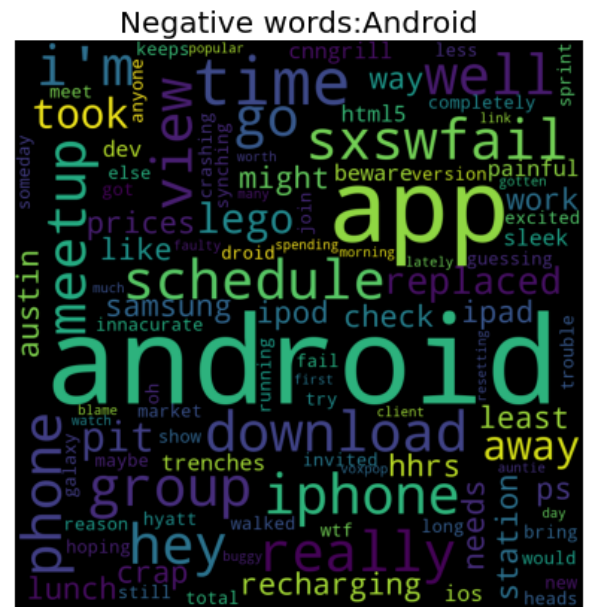
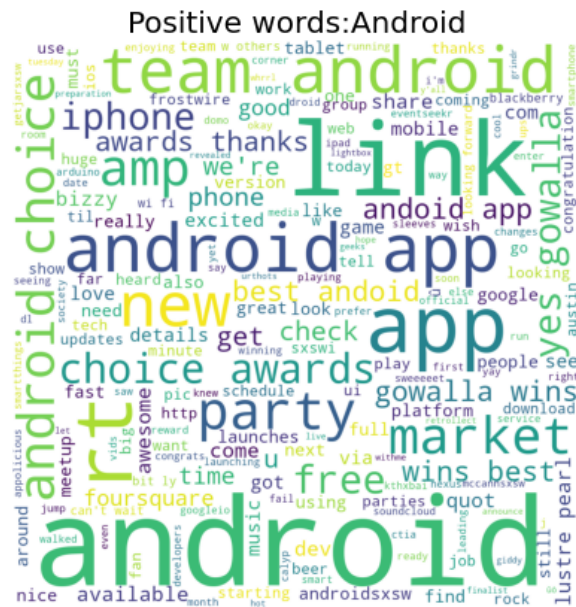
```
In [28]: wordcloud_gen(df, 'Google')
```



google maps seems to be equally represented in both positive and negative tweets. Some people like it and some don't. Same thing with social network . How that ties in with google needs some more exploration.

Android

```
In [29]: wordcloud_gen(df, 'Android')
```



The words `awards`, `wins`, `best` from the positive tweets maybe point towards someone from the festival winning or an app on the Android platform winning some award! `Samsung` pops up in negative tweets suggesting issues with apps running on Samsung phones.

Model to predict company from tweets

Since we're only concerned about which company's is favoured by customers let's further group the products as Apple and Google.

```
In [30]: #Let's group product/services that resemble each other for both brands. This will make

df['product_service'].replace('Android', 'Google', inplace=True)
df['product_service'].replace(['iPhone', 'iPad', 'iPad or iPhone App'], 'Apple', inplace=True)

#checking
df['product_service'].value counts()
```

```
Out[30]: Not Applicable    5801
         Apple             2409
         Google             882
         Name: product service, dtype: int64
```

```
In [31]: # creating target values for product_service
new_map = {'Not Applicable':0,
           'Apple':1,
           'Google':2}

df['target'] = df['product_service'].map(new_map)
df.head()
```

```
Out[31]:
```

	tweet_text	product_service	emotion	text_token	target
0	.@wesley83 i have a 3g iphone. after 3 hrs twe...	Apple	Negative emotion	[wesley83, 3g, iphone, hrs, tweeting, rise_aus...	1
1	@jessedee know about @fludapp ? awesome ipad/i...	Apple	Positive emotion	[jessedee, know, fludapp, awesome, ipad, iphon...	1
2	@swonderlin can not wait for #ipad 2 also. the...	Apple	Positive emotion	[swonderlin, wait, ipad, also, sale, sxsw]	1
3	@sxsw i hope this year's festival isn't as cra...	Apple	Negative emotion	[sxsw, hope, year, festival, crashy, year, iph...	1
4	@sxtxstate great stuff on fri #sxsw: marissa m...	Google	Positive emotion	[sxtxstate, great, stuff, fri, sxsw, marissa, ...	2

Vectorizers

To be able to apply ML models to text data, we must first convert them into a numeric form.

This is accomplished by using *Vectorizers*. *Vectorizers* convert each word in the corpus into a feature and create vectors for each. There are different vectorizers and here, we will use the following three:

1. CountVectorizer
2. Tf-IDF Vectorizer
3. Word2vec Vectorizer

CountVectorizer

CountVectorizer builds on the *Bag Of Words* concept. All the words in the corpus are taken and their frequencies are calculated. The output of the CountVectorizer is a sparse matrix where each feature is a word and the column is the vector of it's frequencies in each document.

```
In [32]: #setting up X,y train and test sets
X= df['text_token']
y = df['target']

X_train,X_test,y_train,y_test = train_test_split(X,y,random_state=123)
```

Since we have already pre-processed our text data, we have to circumvent *CountVectorizer's* preprocessing and tokenizing parameters. We do this by creating a dummy function:

```
In [33]: def dummy(doc):
return doc
```

Let's build baseline models using LogisticRegression, Naive-Bayes, Random Forest and SVM . We can build pipelines for each model and calculate the f1-score for each by creating a loop.

```
In [34]: # building a pipeline of LogisticRegression, Naive-Bayes, SVM and RandomForest models

pipe_lr = Pipeline([('vectorizer',CountVectorizer(stop_words=stopwords_list,preprocess
('model',LogisticRegression(random_state=123,solver='liblinear'))
]))

pipe_nb = Pipeline([('vectorizer',CountVectorizer(stop_words=stopwords_list,preprocess
```



```

        ('model', MultinomialNB())
    ])

pipe_rf = Pipeline([('vectorizer', CountVectorizer(stop_words=stopwords_list, preprocess
        ('model', RandomForestClassifier(random_state=123))
    ])

pipe_svm = Pipeline([('vectorizer', CountVectorizer(stop_words=stopwords_list, preprocess
        ('model', svm.SVC(random_state=123))
    ])

#setting up names for the classification report
names_dict = dict(df['product_service'].value_counts())
names = [name for name in names_dict]

#build a list of tuples to build a df
models = ['LogReg', 'MultiNB', 'RForest', 'SVM']
f1 = []

#fitting the models on the train sets
pipelines = [pipe_lr, pipe_nb, pipe_rf, pipe_svm]

for pipe in pipelines:
    pipe.fit(X_train, y_train)
    predictions = pipe.predict(X_test)
    # print(pipe)
    # print(classification_report(y_test, predictions, target_names=names))
    f1.append(f1_score(y_test, predictions, average='macro'))

#building a df of the f1_scores
scores = list(zip(models, f1))
scores_df = pd.DataFrame(data=scores, columns=['model', 'f1_score_cv'])

```

Tf-IDF Vectorizer

Tfidf Vectorizer takes into account the relative importance of the word to the corpus. It combines *term frequency* and *inverse document frequency*. It calculates how often a word occurs in a document (term frequency) and also how many documents contain the word (inverse document frequency). The output is again a sparse matrix like with CountVectorizer.

```

In [35]: # repeating the same processes as above
pipe_lr = Pipeline([('vectorizer', TfidfVectorizer(stop_words=stopwords_list, preprocess
        ('model', LogisticRegression(random_state=123, solver='liblinear'))
    ])

pipe_nb = Pipeline([('vectorizer', TfidfVectorizer(stop_words=stopwords_list, preprocess
        ('model', MultinomialNB())
    ])

pipe_rf = Pipeline([('vectorizer', TfidfVectorizer(stop_words=stopwords_list, preprocess
        ('model', RandomForestClassifier(random_state=123))
    ])

pipe_svm = Pipeline([('vectorizer', TfidfVectorizer(stop_words=stopwords_list, preprocess
        ('model', svm.SVC(random_state=123))
    ])

#setting up names for the classification report

```

```

names_dict = dict(df['product_service'].value_counts())
names = [name for name in names_dict]

#build a list of tuples to build a df
models=['LogReg', 'MultiNB', 'RForest', 'SVM']
f1 = []

#fitting the models on the train sets
pipelines = [pipe_lr,pipe_nb,pipe_rf,pipe_svm]

for pipe in pipelines:
    pipe.fit(X_train,y_train)
    predictions = pipe.predict((X_test))
    # print(pipe)
    # print(classification_report(y_test,predictions,target_names=names))
    f1.append(f1_score(y_test,predictions,average='macro'))

#adding the tf f1-scores to the scores_df
scores_df['f1_score_tf'] = f1

```

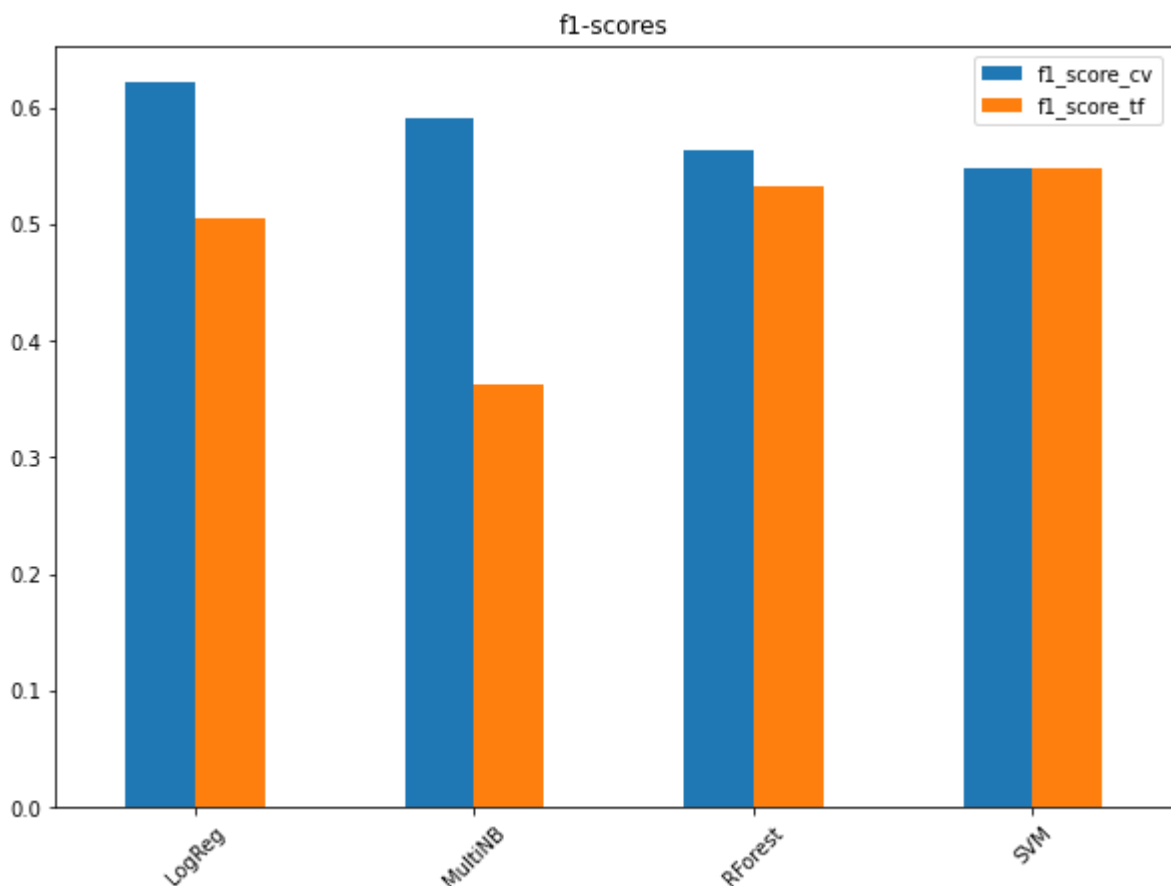
Now that we've run models with 2 CV and Tfidf vectorizers, let's visualize the f1-scores of each

In [36]: *#visualizing the f1-scores of all the models for the two vectorizers*

```

fig,ax = plt.subplots(figsize=(10,7))
scores_df.plot(kind='bar',ax=ax);
ax.set_xticklabels(models,rotation=45);
ax.set_title('f1-scores');

```



Since the LogisticRegression model with the CountVectorizer has the highest f1-score among all models, let's use that for optimizations

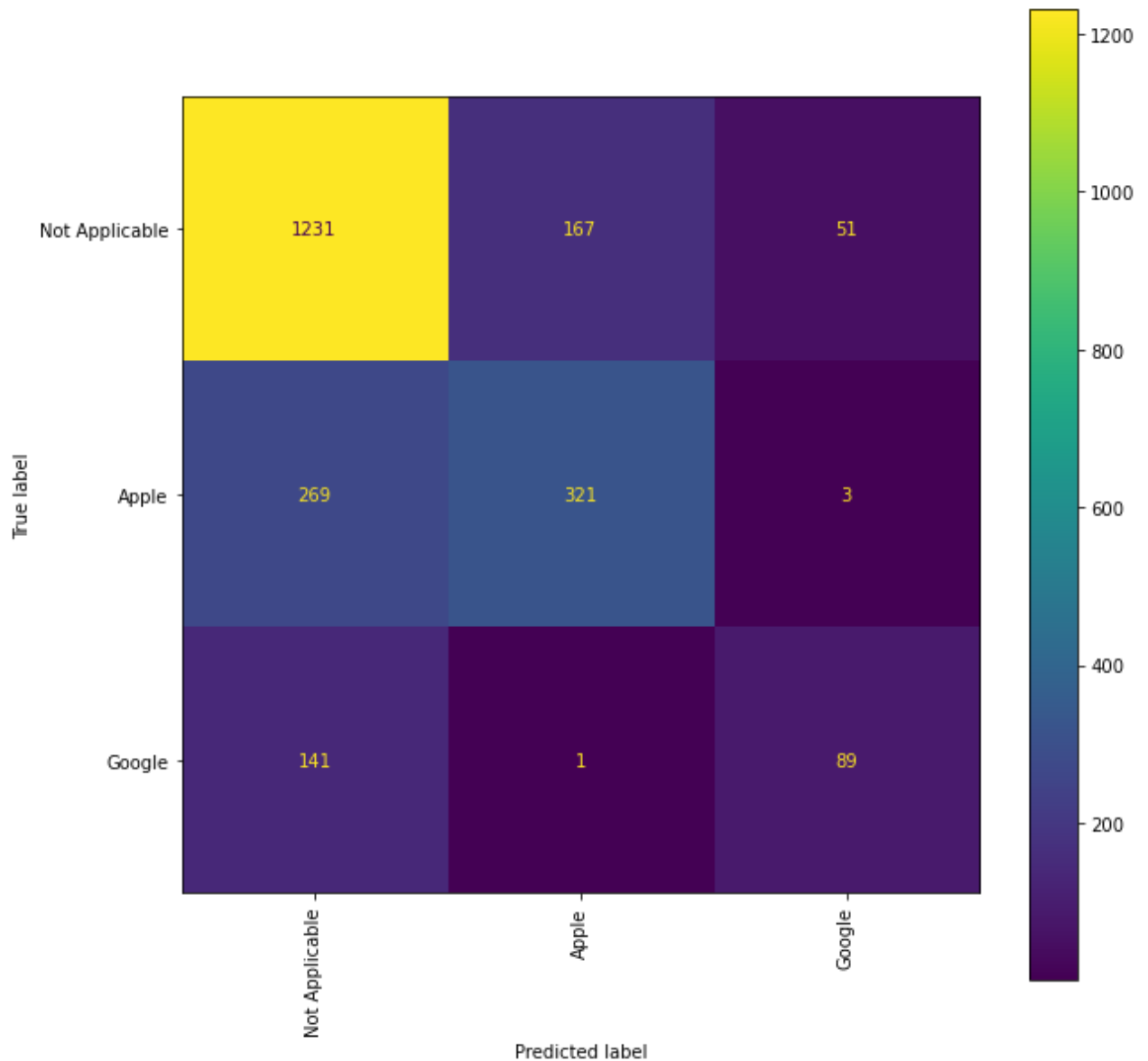
Tuning LogisticRegression with Countvectorizer

Let's get a closer look at the LR with CV model:

```
In [37]: lr_pipe = Pipeline([('vectorizer',CountVectorizer(stop_words=stopwords_list,preprocess
                        ('model',LogisticRegression(random_state=123,solver='liblinear')))
                        ])

lr_pipe.fit(X_train,y_train)
predictions=lr_pipe.predict(X_test)
print(classification_report(y_test,predictions,target_names=names))
fig,ax=plt.subplots(figsize=(10,10))
plot_confusion_matrix(lr_pipe,X_test,y_test,display_labels=names,ax=ax,xticks_rotation=
```

	precision	recall	f1-score	support
Not Applicable	0.75	0.85	0.80	1449
Apple	0.66	0.54	0.59	593
Google	0.62	0.39	0.48	231
accuracy			0.72	2273
macro avg	0.68	0.59	0.62	2273
weighted avg	0.71	0.72	0.71	2273



From the above, we can see that the model does relatively well in the `Not Applicable` class compared to `Apple` and `Google`. This is unsurprising given the imbalance in the data.

min_df and max_df values

We can try tuning some of the hyperparameters of the model to improve performance. `min_df` and `max_df` are two that we can look. `min_df` removes words that appear rarely. Since they are rare, it is possible that they will not provide a lot of information. `max_df` is the opposite of `min_df`. If the words are too frequent, chances are they too do not provide a lot of information. By creating a range for each parameter, we can run a loop to see if model performance improves

```
In [38]: # setting a range for min_df and max_df
min_df_value = np.arange(1,5) # words that appear less than the min_df value in all doc
max_df_value = np.arange(1500,1505) # words that appear more than the max_df value in a

#initiating lists to use for plotting
f_score = []
min_value=[]
max_value=[]
```

```

#setting up the loop for min_df and max_df values
for i in min_df_value:
    for j in max_df_value:
        min_value.append(i)
        max_value.append(j)
        #instantiate pipeline
        new_pipe = Pipeline([('vectorizer',CountVectorizer(stop_words=stopwords_list,m
                                                             ('model',LogisticRegression(random_state=123,solver=
                                                             ]))
        new_pipe.fit(X_train,y_train)
        preds = new_pipe.predict(X_test)

        #getting f1 score
        score = round(f1_score(y_test,preds,average='macro'),3)
        f_score.append(score)
#         print(f'min_df = {i}, max_df = {j}, f1_score={score}')

#visualizing the accuracy score for the different combinations
d = list(zip(min_value,max_value))
fig,ax=plt.subplots(figsize=(15,5))
plt.tick_params(bottom=False)
ax.plot(f_score,marker='o',markerfacecolor='r',ls='--');
ax.set_xticklabels([],[]);
ax.set_title('F1 Score of the Logistic Regression Model with CountVectorizer');

```



Clearly, we've made the model worse. Our initial f1-score was 0.62 but here we're maxed out at 0.574

n-gram

The idea behind n-grams is that sometimes word pairings or short phrases are better. For eg: 'black sheep' is more informative than 'black' and 'sheep' separately

```

In [39]: lr_pipe = Pipeline([('vectorizer',CountVectorizer(stop_words=stopwords_list,preprocess
                                                             ngram_range=(1,2),min_df=1,max_df=200
                                                             ('model',LogisticRegression(random_state=123,solver='liblinear'))
                                                             ]))

lr_pipe.fit(X_train,y_train)
predictions=lr_pipe.predict(X_test)
print(classification_report(y_test,predictions,target_names=names))
# fig,ax=plt.subplots(figsize=(10,10))
# plot_confusion_matrix(lr_pipe,X_test,y_test,display_labels=names,ax=ax,xticks_rotatio

```

	precision	recall	f1-score	support
Not Applicable	0.75	0.86	0.80	1449
Apple	0.67	0.54	0.60	593
Google	0.63	0.35	0.45	231
accuracy			0.73	2273
macro avg	0.68	0.59	0.62	2273
weighted avg	0.72	0.73	0.71	2273

We can see that there is no discernible change in model performance

Stemming using PorterStemmer

With stemming, we use the use root of the word. For eg: ran,runs,running all stem from the word run. This way we reduce the number of features and can improve accuracy of the model.

```
In [40]: #initializing the stemmer
ps=PorterStemmer()

#creating a function to tokenize and stem the tokens
def stem_and_tokenize(document):
    tokens = tokenizer.tokenize(document)
    return [ps.stem(token) for token in tokens]
```

```
In [41]: df['stemmed_tokens'] = df['tweet_text'].apply(stem_and_tokenize)
df.head()
```

```
Out[41]:
```

	tweet_text	product_service	emotion	text_token	target	stemmed_tokens
0	.@wesley83 i have a 3g iphone. after 3 hrs twe...	Apple	Negative emotion	[wesley83, 3g, iphone, hrs, tweeting, rise_aus...	1	[wesley83, have, 3g, iphon, after, hr, tweet, ...
1	@jessedee know about @fludapp ? awesome ipad/i...	Apple	Positive emotion	[jessedee, know, fludapp, awesome, ipad, iphon...	1	[jessedee, know, about, fludapp, awesom, ipad, ...
2	@swonderlin can not wait for #ipad 2 also. the...	Apple	Positive emotion	[swonderlin, wait, ipad, also, sale, sxsw]	1	[swonderlin, can, not, wait, for, ipad, also, ...
3	@sxsw i hope this year's festival isn't as cra...	Apple	Negative emotion	[sxsw, hope, year, festival, crashy, year, iph...	1	[sxsw, hope, thi, year, festiv, isn, as, crash...
4	@sxtxstate great stuff on fri #sxsw: marissa m...	Google	Positive emotion	[sxtxstate, great, stuff, fri, sxsw, marissa, ...	2	[sxtxstate, great, stuff, on, fri, sxsw, maris...

```
In [42]: #running logistic regression on the stemmed tokens
#re-defining X and y
X = df['stemmed_tokens']
y = df['target']

X_train2,X_test2,y_train2,y_test2 = train_test_split(X,y,random_state=123)
```

```
In [43]: lr_pipe = Pipeline([('vectorizer',CountVectorizer(stop_words=stopwords_list,preprocess
```

```

ngram_range=(1,2),min_df=1,max_df=200
('model',LogisticRegression(random_state=123,solver='liblinear'))
])

lr_pipe.fit(X_train2,y_train2)
predictions=lr_pipe.predict(X_test2)
print(classification_report(y_test2,predictions,target_names=names))

```

	precision	recall	f1-score	support
Not Applicable	0.75	0.85	0.80	1449
Apple	0.66	0.55	0.60	593
Google	0.61	0.37	0.46	231
accuracy			0.72	2273
macro avg	0.67	0.59	0.62	2273
weighted avg	0.71	0.72	0.71	2273

Word2Vec

Word2vec is another vectorization method and falls under the category called **Word Embeddings**. It is essentially a neural network with an i/p layer, hidden layer and an o/p layer. The vectors are created in an **embedding space** and are used to capture the semantic relationships between words.

Here, we will import the Word2vec vector from the open source **gensim** library and use the **skip gram** architecture for modelling. The gensim library has vectors built in that we will use to base our model off of.

```

In [44]: #instantiate the vect
model = Word2Vec(df['text_token'], vector_size=100, window=2, min_count=5, sg=1)

#train the model
model.train(df['text_token'], epochs=15, total_examples=model.corpus_count)

```

Out[44]: (955884, 1574070)

Experimentation

The calculated vectors are stored in the `Word2VecKeyedVectors` instance stored in the `wv` attribute. Let's assign it to a different variable to save ourselves from lot's of keystrokes

```

In [45]: wv=model.wv

```

Checking the vector for the word `battery`. This will display the weights that the model has calculated for the context that the word 'battery' will most likely used in

```

In [46]: wv.most_similar('battery')

```

```

Out[46]: [('double', 0.7639098763465881),
('backup', 0.7617971301078796),
('realized', 0.7609479427337646),
('brightness', 0.7556403279304504),
('charged', 0.7361781597137451),
('woke', 0.7357349991798401),
('charge', 0.7338263988494873),
('fully', 0.7286348938941956),

```

```
('size', 0.7285921573638916),  
('extended', 0.7192755937576294)]
```

```
In [47]: #the vector associated with the word 'battery'  
         wv['battery']
```

```
Out[47]: array([-0.31062588,  0.14848433,  0.39587373, -0.3254254 ,  0.3428913 ,  
                -0.12781687,  0.21243063,  0.20278853,  0.13866153,  0.00916754,  
                -0.42157844, -0.5654535 , -0.16671962, -0.33284813,  0.3532581 ,  
                -0.58793986, -0.17072922, -0.4518122 ,  0.07336178, -0.6427151 ,  
                 0.2295601 ,  0.25446963,  0.44342545,  0.23790124, -0.10359608,  
                 0.37241188, -0.2534862 , -0.1442654 , -0.21276833,  0.06025869,  
                -0.43979532,  0.3085918 , -0.09884585, -0.09907828,  0.16408782,  
                -0.19602658, -0.12746182, -0.31128547, -0.2668944 , -0.3510607 ,  
                -0.15442765, -0.14438877,  0.03245652, -0.07728881,  0.34669074,  
                -0.41281095, -0.08171452,  0.35391945, -0.29985476,  0.65394104,  
                 0.10159953, -0.31657976,  0.24982709,  0.19805153, -0.49896583,  
                 0.01723918, -0.51610565, -0.08741776, -0.04401156, -0.17922996,  
                -0.25911343, -0.02756198,  0.13036403, -0.18349285, -0.662885 ,  
                 0.07961139, -0.21506497,  0.28852466, -0.55803657,  0.46534306,  
                 0.15629855,  0.5971035 ,  0.521865 , -0.40409046, -0.63010347,  
                 0.4598462 , -0.06739046, -0.7279858 , -0.81242514, -0.09288336,  
                -0.36208814,  0.00136437, -0.2624425 ,  0.5682431 , -0.09854691,  
                 0.3997231 ,  0.26411197, -0.05258031,  0.16847298, -0.09261093,  
                -0.1375739 , -0.25194216, -0.10749363, -0.04541309,  0.20824674,  
                 0.3025923 ,  0.02866755, -0.27534315,  0.17304133, -0.11297138],  
          dtype=float32)
```

```
In [48]: #getting the list of words from the model  
         words = list(model.wv.index_to_key)  
         words[0:10] #looking at the first 10 words
```

```
Out[48]: ['sxsx',  
          'link',  
          'rt',  
          'google',  
          'ipad',  
          'apple',  
          'quot',  
          'iphone',  
          'store',  
          'new']
```

```
In [49]: #getting the vectors associated with each of those words  
         vector_list = [model.wv[word] for word in words]  
         vector_list[0] #examining the vector for the first word
```

```
Out[49]: array([ 0.29971173,  0.42509407,  0.20842887,  0.08963112,  0.12931864,  
                -0.2616976 , -0.1723097 ,  0.26966655, -0.79576176, -0.33104157,  
                 0.21922581, -0.13764803,  0.4696644 ,  0.11016657, -0.24328412,  
                -0.19033964, -0.06619477,  0.07306217, -0.22873306, -0.4607964 ,  
                 0.4568384 ,  0.64830726,  0.5703138 , -0.3012564 ,  0.23925447,  
                 0.18719737, -0.04412404,  0.11028476, -0.24879463,  0.05252969,  
                -0.17805837, -0.14992486,  0.4711132 , -0.3932661 ,  0.2882445 ,  
                -0.33420452,  0.32569024, -0.08540811,  0.08160986, -0.31962585,  
                -0.11755685,  0.32137683, -0.18738702,  0.18039149,  0.12187223,  
                -0.07080109, -0.1451052 ,  0.44288433, -0.07770481,  0.37076578,  
                -0.28124377,  0.13016275, -0.22774702,  0.19805995,  0.38397837,  
                 0.4204803 ,  0.14457156, -0.41504878, -0.36513576,  0.07666832,  
                -0.32914212, -0.196766 , -0.02029376, -0.44120765,  0.16723628,  
                 0.10151558, -0.06786009,  0.25020754,  0.07110384,  0.0903364 ,  
                 0.57346267,  0.25396067, -0.25935817, -0.19316117,  0.48967972,  
                -0.1134918 , -0.23788472,  0.04348994, -0.11112805, -0.2847708 ,  
                -0.04455198,  0.37639058,  0.5855522 ,  0.34901792, -0.20244057,
```



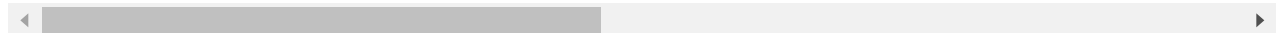
```
-0.2699105 , 0.27420467, 0.19090655, -0.4409254 , 0.23152538,
-0.2975374 , -0.18320633, 0.26953068, 0.14157954, 0.20940313,
-0.4676908 , 0.31250137, 0.14555797, -0.033759 , 0.02853949],
dtype=float32)
```

```
In [50]: #creating a df
word_vect_zip = dict(zip(words,vector_list))
word_vect_df = pd.DataFrame(word_vect_zip)
word_vect_df.head()
```

```
Out[50]:
```

	sxsw	link	rt	google	ipad	apple	quot	iphone	store	r
0	0.299712	-0.114142	-0.198704	0.068581	-0.646271	0.028013	0.173619	0.276399	-0.442767	-0.513
1	0.425094	0.663902	0.293445	0.796443	-0.179417	0.199373	0.532154	0.153646	0.233683	0.539
2	0.208429	0.402341	-0.213680	0.033861	0.376090	0.202648	0.221409	0.778955	0.082375	0.055
3	0.089631	-0.132784	0.138557	-0.199975	0.032881	-0.270201	-0.283521	0.007370	0.521897	-0.002
4	0.129319	0.138335	0.458824	0.466028	0.538263	-0.348458	0.155434	0.286640	0.575911	0.315

5 rows × 2383 columns



The df is used to illustrate the different words created by the model and their corresponding vectors.

Mean Embeddings

Now, we need to get a vector representation for each document to be able to apply a ML model. For this, we will take the mean of the vectors in each document for the words that are in the model vocabulary. Let's define a function to get the mean vector for each document:

```
In [51]: #to get the vector representation of each document, you will define a fuction that will
#vector for each word in the document that is in the model's vocab, beacuse obviously yo
#vectors for words that are not there. Then you will take the average of the vectors an
#vector will the vector for that document

def doc_vector(token):
    vector_size=model.wv.vector_size #getting the size of the vectors created
    resultant_vector = np.zeros(vector_size) # initializing a vector of zeros of the sa
    ctr=1 #counter
    for w in token:
        if w in words:
            ctr += 1
            resultant_vector += wv[w]
    resultant_vector = resultant_vector/ctr
    return resultant_vector
```

```
In [52]: #applying the function to the stemmed_tokens columes
df['vectors'] = df['stemmed_tokens'].apply(doc_vector)
df.head()
```

```
Out[52]:
```

	tweet_text	product_service	emotion	text_token	target	stemmed_tokens	vectors
--	------------	-----------------	---------	------------	--------	----------------	---------

	tweet_text	product_service	emotion	text_token	target	stemmed_tokens	vectors
0	.@wesley83 i have a 3g iphone. after 3 hrs twe...	Apple	Negative emotion	[wesley83, 3g, iphone, hrs, tweeting, rise_aus...	1	[wesley83, have, 3g, iphon, after, hr, tweet, ...	[-0.16743164722408568, 0.194734980485269, 0.18...
1	@jessedee know about @fludapp ? awesome ipad/i...	Apple	Positive emotion	[jessedee, know, fludapp, awesome, ipad, iphon...	1	[jessede, know, about, fludapp, awesom, ipad, ...	[-0.328149801492691, -0.04962060189573094, 0.4...
2	@swonderlin can not wait for #ipad 2 also. the...	Apple	Positive emotion	[swonderlin, wait, ipad, also, sale, sxsw]	1	[swonderlin, can, not, wait, for, ipad, also, ...	[-0.2854916701714198, 0.10605820392568906, 0.2...
3	@sxsw i hope this year's festival isn't as cra...	Apple	Negative emotion	[sxsw, hope, year, festival, crashy, year, iph...	1	[sxsw, hope, thi, year, festiv, isn, as, crash...	[-0.002816371353609221, 0.2870309816374044, 0....
4	@sxtxstate great stuff on fri #sxsw: marissa m...	Google	Positive emotion	[sxtxstate, great, stuff, fri, sxsw, marissa, ...	2	[sxtxstate, great, stuff, on, fri, sxsw, maris...	[-0.21151037141680717, 0.057256707921624184, -...

Modelling

Now that we have a vecotrized representation of each document, we can apply different ML models and check performance

```
In [53]: #redefining X&y
X= df['vectors'].to_list()
y = df['target'].to_list()

X_train,X_test,y_train,y_test = train_test_split(X,y,random_state=123)

In [54]: # repeating the same processes as above
pipe_lr_wv = Pipeline([('model',LogisticRegression(random_state=123,solver='liblinear'))

# pipe_nb_wv = Pipeline([('model',MultinomialNB())])
#cannot use NB on negative values. Values will have to be normalized instead

pipe_rf_wv = Pipeline([('model',RandomForestClassifier(random_state=123))])

pipe_svm_wv = Pipeline([('model',svm.SVC(random_state=123))])

#build a list of tuples to build a df
models=['LogReg_wv', 'RForest_wv', 'SVM_wv']
f1 = []

#fitting the models on the train sets
```

```

pipelines = [pipe_lr_wv, pipe_rf_wv, pipe_svm_wv]

for pipe in pipelines:
    pipe.fit(X_train, y_train)
    predictions = pipe.predict(X_test)
    # print(pipe)
    # print(classification_report(y_test, predictions, target_names=names))
    f1.append(f1_score(y_test, predictions, average='macro'))

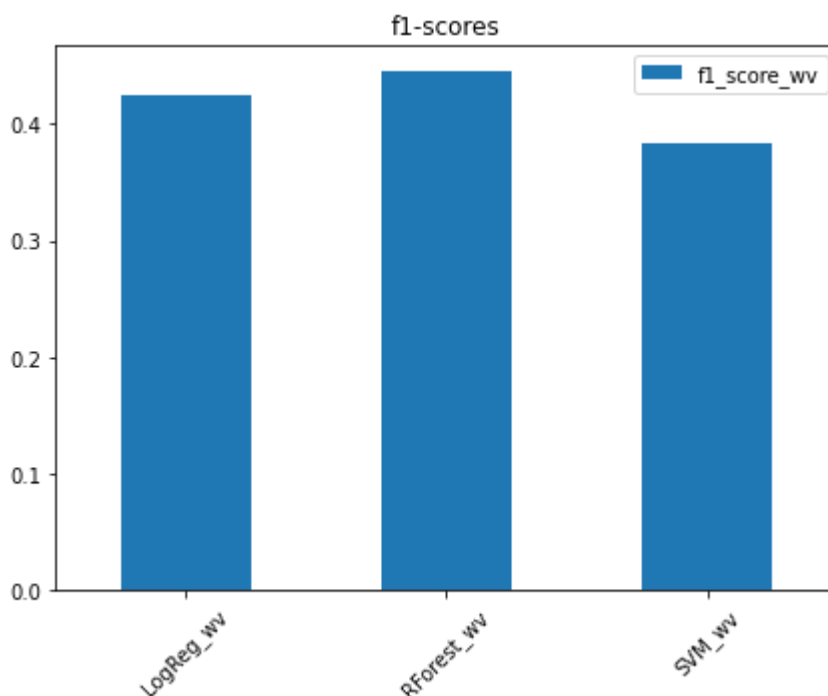
#building a df of the f1_scores
scores_wv = list(zip(models, f1))
scores_df_wv = pd.DataFrame(data=scores_wv, columns=['model', 'f1_score_wv'])

```

```

In [55]: fig, ax = plt.subplots(figsize=(7,5))
scores_df_wv.plot(kind='bar', ax=ax);
ax.set_xticklabels(models, rotation=45);
ax.set_title('f1-scores');

```



As we can see, with Word2vec, the max f1-score that we can achieve is only 0.45, much less than our highest score of 0.62

Next Steps

1. Since the iPad is the most popular product, Acme Online could look for opportunities to boost sales. Acme Online could also maybe expand their portfolio buy offering tablets from other manufacturers to see if they will gain any traction.
2. More data is definitely recommended. Current data is very imbalanced impacting model performance.
3. The hyper parameters of the Word2Vec vectorizer i.e number of epochs, size of the vectors etc. can be tuned to see if results improve.
4. Part-of-Speech tagging can be used to create more features.

5. Ensemble methods like XGBoost and Adaboost can also be trialed for modelling along with other word embedding techniques like fastText and Glove.