**NLP Model For Twitter Sentiment Analysis on Google and Apple products**

The goal of this project is to develop a Natural Language Processing (NLP) model that analyzes and rate sentiment in tweets depending on their content, related to Google and Apple products.The dataset sourced from Kaggle, contains tweets labeled as positive, negative and Neutral

### Why Sentiment Analysis Matters     ¶

In today's digital landscape,sentiment analysis plays a vital role in understanding public perception. Companies use this method to gain insightful information form massive volumes of textual data such as social media comments,news stories and customer evaluations. We can improve decision-making and business intelligence by implementing NLP algorithms like Support Vector Machines(SVM) and Naive bayes.

The main objective of this project is to :

1. Preprocess Twitter data(tokenization, cleaning and vectorization)
2. Apply machine learning models for sentiment classification
3. Asses model performance using important metrics

Our goal is to create an effective sentiment analysis model by the project's conclusion that offers insightful information about consumer perception of Apple and Google products

## Summary

### Dataset Overview

This dataset contains information about tweets, including their date, classification, and other relevant features. For this task, I extracted only two columns:

- **Feature Variable (Text/Tweets)**: The actual tweet content.
- **Sentiment**: The rating of each tweet, categorized as negative, positive, or neutral.

---

### Data Preparation

The data preparation phase included **data cleaning** and **data preprocessing** using the pandas and NLTK libraries.

**1. Data Cleaning (Using Pandas)**

- The dataset had **no missing values**, so no imputation was required.
- **Removed duplicate tweets**—for instance, one tweet appeared 304 times.
- Dropped irrelevant sentiment labels**, such as 'not_relevant', to retain only the core sentiments **(negative, positive, and neutral)**.
- **Converted sentiments to integer values** to ensure compatibility with machine learning models.

**2. Data Preprocessing (Using NLTK)**

- **Removed hyperlinks, usernames, single-character words, and hashtags** (including their values) using **regular expressions**, as they do not meaningfully contribute to sentiment analysis.
- **Eliminated stopwords and punctuation ** using the **NLTK corpus library** for stopwords and the **string library** for punctuation, as these do not add significant meaning to sentences.
- **Applied lemmatization ** using the **WordNet Lemmatizer**, converting words to their root form (e.g., "running" → "run").

---

### Data Visualization

To explore and understand the dataset, I used:

- **Seaborn's** countplot to visualize the distribution of target sentiment classes.
- **WordCloud** to generate a visual representation of the most common words in the dataset.

---

### Modeling

- Used **Scikit-learn's model_selection** library to split the dataset into training and testing sets.
- Implemented **pipelines** to streamline **vectorization, SMOTE (Synthetic Minority Over-sampling Technique), and classification models**.
- Applied **TF-IDF Vectorizer** to convert text data into numerical representations.

- Used **SMOTE** to address **class imbalance**, as one sentiment category comprised **more than 50%** of the dataset.
- Evaluated various **machine learning algorithms** for classification to determine the best-performing model.
- Created **custom functions** to automate repetitive processes such as model fitting and prediction.

---

### Evaluation Metrics

To assess model performance, I used the following evaluation metrics:

- **Accuracy Score**: The proportion of correctly classified instances out of the total instances.
- **Precision Score**: The ratio of correctly predicted positive instances to total predicted positive instances.
- **Recall Score**: The ratio of correctly predicted positive instances to actual positive instances in the dataset.
- **ROC Curve**: A graphical representation of the true positive rate versus the false positive rate.

---

### Model Evaluation & Predictions

- The **test set** obtained from `train_test_split (X_test, y_test)` was used for model evaluation and making predictions.
- Additionally, I developed a **custom function** that allows classification of sentiment based on user input.

## Step 1: Importing libraries

In [68]:
```python
1  #Basic Python Libraries
2  import pandas as pd
3  import numpy as np
4  import seaborn as sns
5  import matplotlib.pyplot as plt
6  #Natural Language processing libraries
7  from nltk.corpus import stopwords
8  import string
9  from nltk.stem.wordnet import WordNetLemmatizer
10 from nltk.tokenize import RegexpTokenizer,word_tokenize
11 import re
12 from nltk import FreqDist
13 from wordcloud import WordCloud
14 #scikit-learn
15 from sklearn.model_selection import train_test_split,RandomizedSearchCV,GridSearchCV
16 from sklearn.feature_extraction.text import TfidfVectorizer
17 from sklearn.linear_model import LogisticRegression
18 from sklearn.pipeline import Pipeline
19 from sklearn.ensemble import RandomForestClassifier,GradientBoostingClassifier
20 from xgboost import XGBClassifier
21 from sklearn.tree import DecisionTreeClassifier
22 from sklearn.naive_bayes import MultinomialNB
23 from sklearn import svm
24 from sklearn.preprocessing import label_binarize
25 #Evaluation metrics
26 from sklearn.metrics import accuracy_score,recall_score,precision_score,f1_score,roc_curve,auc,r2_score,f1_sc
27 #imbalanced-learn
28 from imblearn.pipeline import Pipeline  # Use imbalanced-learn's Pipeline
29 from imblearn.over_sampling import SMOTE
```

## Step 2: Understanding the dataset

```
In [2]:    1  #Read the data from the csv file as a dataframe and dispaly the first five rows
           2  data = pd.read_csv('data/Apple-Twitter-Sentiment-DFE.csv',encoding='latin-1')
           3  data.head()
```

Out[2]:

| | _unit_id | _golden | _unit_state | _trusted_judgments | _last_judgment_at | sentiment | sentiment:confidence | date | id | query | s |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 623495513 | True | golden | 10 | NaN | 3 | 0.6264 | Mon Dec 01 19:30:03 +0000 2014 | 5.400000e+17 | #AAPL OR @Apple | |
| 1 | 623495514 | True | golden | 12 | NaN | 3 | 0.8129 | Mon Dec 01 19:43:51 +0000 2014 | 5.400000e+17 | #AAPL OR @Apple | |
| 2 | 623495515 | True | golden | 10 | NaN | 3 | 1.0000 | Mon Dec 01 19:50:28 +0000 2014 | 5.400000e+17 | #AAPL OR @Apple | |
| 3 | 623495516 | True | golden | 17 | NaN | 3 | 0.5848 | Mon Dec 01 20:26:34 +0000 2014 | 5.400000e+17 | #AAPL OR @Apple | |
| 4 | 623495517 | False | finalized | 3 | 12/12/14 12:14 | 3 | 0.6474 | Mon Dec 01 20:29:33 +0000 2014 | 5.400000e+17 | #AAPL OR @Apple | |

For this project we will use the text column as the feature variable and sentiment column as the target variable

```
In [3]:    1  # extracting the text and sentiment column and previewing first five rows
           2  data = data[['text','sentiment']]
           3  data.head()
```

Out[3]:

| | text | sentiment |
|---|---|---|
| 0 | #AAPL:The 10 best Steve Jobs emails ever...htt... | 3 |
| 1 | RT @JPDesloges: Why AAPL Stock Had a Mini-Flas... | 3 |
| 2 | My cat only chews @apple cords. Such an #Apple... | 3 |
| 3 | I agree with @jimcramer that the #IndividualIn... | 3 |
| 4 | Nobody expects the Spanish Inquisition #AAPL | 3 |

```
In [4]:    1  #Check the number of records and features using the shape method
           2  data.shape
           3  print(f'This dataset contains {data.shape[0]} rows')
           4  print(f'This dataset contains {data.shape[1]} columns')
```

```
This dataset contains 3886 rows
This dataset contains 2 columns
```

```
In [5]:    1  #Access information about our dataset
           2  data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3886 entries, 0 to 3885
Data columns (total 2 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   text       3886 non-null   object
 1   sentiment  3886 non-null   object
dtypes: object(2)
memory usage: 60.8+ KB
```

Both our columns have non-null values, there are no missing data points ensuring our dataset is ready for futher processing and model

In [6]:
```python
1  # conforming that there are no null values
2  data.isna().sum()
```

Out[6]:
```
text         0
sentiment    0
dtype: int64
```

In [7]:
```python
1  #get summary statistics of our data
2  data.describe()
```

Out[7]:

|        | text | sentiment |
|--------|------|-----------|
| count  | 3886 | 3886 |
| unique | 3219 | 4 |
| top    | RT @OneRepublic: Studio at 45,000 ft. One out... | 3 |
| freq   | 304  | 2162 |

```
1  we see a tweet that appears 304 times we may need to drop duplicates
2  sentiment has 4 categories with 3 (Neutral) taking up 55.64% this created an imbalance in the dataset which we
   will need to handle.There is 4 sentiment values instead of the expected 3 we may need to check out the exra
   one
```

## step 3: Data cleaning

In [8]:
```python
1  # check for and remove duplicates
2  data[data.duplicated()]
```

Out[8]:

|      | text | sentiment |
|------|------|-----------|
| 32   | RT @thehill: Justice Department cites 18th cen... | 3 |
| 34   | RT @thehill: Justice Department cites 18th cen... | 3 |
| 38   | RT @thehill: Justice Department cites 18th cen... | 3 |
| 42   | RT @thehill: Justice Department cites 18th cen... | 3 |
| 45   | RT @thehill: Justice Department cites 18th cen... | 3 |
| ...  | ... | ... |
| 3846 | RT @TeamCavuto: Protesters stage #DieIn protes... | 3 |
| 3852 | RT @TeamCavuto: Protesters stage #DieIn protes... | 3 |
| 3855 | RT @Ecofantasy: Thinking of upgrading to #Yose... | 1 |
| 3878 | RT @shannonmmiller: Love the @Apple is support... | 5 |
| 3885 | RT @SwiftKey: We're so excited to be named to ... | 5 |

643 rows × 2 columns

In [9]:
```python
1  #Drop the above duplicates
2  data.drop_duplicates(inplace=True)
```

In [10]:
```python
1  #check sentiment labels
2  data['sentiment'].value_counts()
```

Out[10]:
```
3               1681
1               1102
5                379
not_relevant      81
Name: sentiment, dtype: int64
```

In [11]:
```python
1  #drop the not_relevant label
2  data = data.query('~(sentiment=="not_relevant")') #data = data.query('sentiment !="not_relevant"')
3  #confirm sentiment 'not_relevant' has been removed
4  data['sentiment'].value_counts()
```

Out[11]:
```
3    1681
1    1102
5     379
Name: sentiment, dtype: int64
```

```
In [12]:    1  #convert sentiment column to int
            2  data['sentiment'] = data['sentiment'].astype(int)
            3  #confirm removal using info method
            4  data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 3162 entries, 0 to 3884
Data columns (total 2 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   text       3162 non-null   object
 1   sentiment  3162 non-null   int32
dtypes: int32(1), object(1)
memory usage: 61.8+ KB
```

### Step 3:Preprocess text

```
In [13]:    1  def clean_text(text):
            2      stopword_list = stopwords.words('english')
            3      stopword_list += string.punctuation
            4      lemmatizer = WordNetLemmatizer()
            5
            6      #remove hyperlinks,usermames,words with 1 character,hashtags and their values
            7      text =re.sub(r"https?:[^\s]+|@[\S]+|\b\w\b|\#\w+|\.\.+", "", text)
            8
            9      # Tokenize the text
           10      text = word_tokenize(text)
           11
           12      # Lowercase the text and remove stopwords
           13      text = [word.lower() for word in text if word.lower() not in stopword_list]
           14
           15       #lemmatize
           16      text = [lemmatizer.lemmatize(word)for word in text]
           17
           18      # join processed text as a single string
           19      text = ' '.join(text)
           20      return text
```

Here we are performing some pre-processing on the data befoee converting it into vectors and passing it to the machine learning model.

we create a function that iterates through each record,and does the following

1. using regular expression
   - we get rid of hyperlinks in the data that was showing the link where the tweet came from,
     - we also remove '@username' since the username of the person who tweeted doesn't help in rating a sentiment
     - Remove twitter hashtags together with the value since they do not contribute to a sentiment meaningfully
     - Remove single character words
2. Remove stop words such as 'i','me' which do not add much value and convert them to lowercase so that eg word bad and BAd are treated as the same thing
3. Lemmatize words to reduce words to their base so that word like run and running are treated as the same thing

```
In [14]:    1  #Apply the function to the text column
            2  data['cleaned_text'] = data['text'].apply(clean_text)
```

```
In [15]:    1  ### confirm data procesed
            2  data[['text','cleaned_text']].head()
```
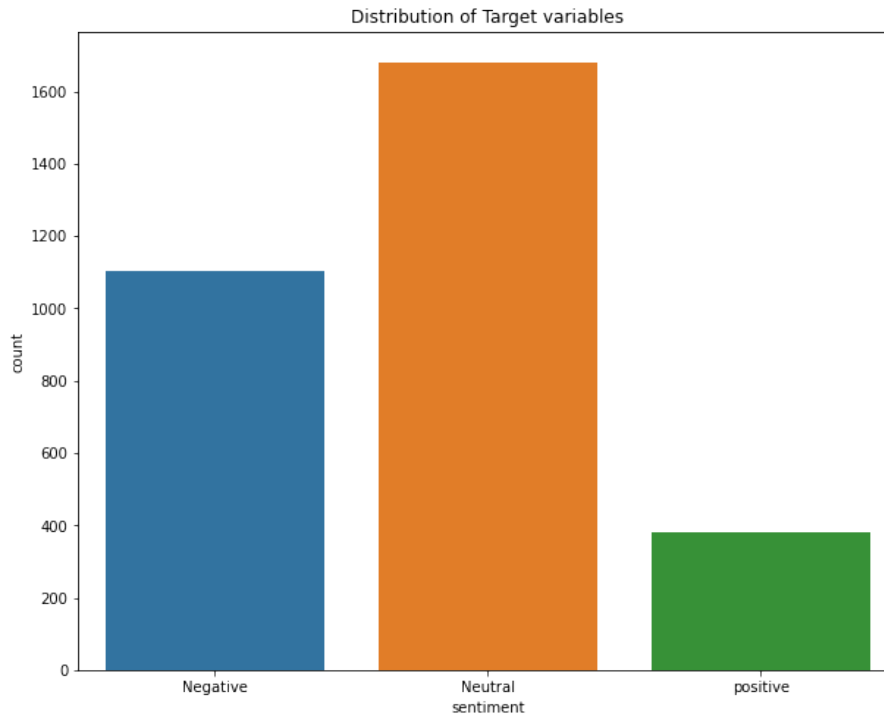
Out[15]:

|   | text | cleaned_text |
|---|---|---|
| 0 | #AAPL:The 10 best Steve Jobs emails ever...htt... | 10 best steve job email ever |
| 1 | RT @JPDesloges: Why AAPL Stock Had a Mini-Flas... | rt aapl stock mini-flash crash today aapl |
| 2 | My cat only chews @apple cords. Such an #Apple... | cat chew cord |
| 3 | I agree with @jimcramer that the #IndividualIn... | agree trade extended today pullback good see |
| 4 | Nobody expects the Spanish Inquisition #AAPL | nobody expects spanish inquisition |

**Step5: EDA**

**Data Visualization of Target variables**

In [16]:
```python
fig, ax = plt.subplots(figsize=(10,8))
ax= sns.countplot(data = data, x='sentiment')
#using xtickslabels to label values 1 as negative,3 as neutral and 5 as positive
ax.set(title='Distribution of Target variables',xticklabels=['Negative','Neutral','positive']);
```



```
1  The label data looks imbalanced where neutral takes more than 50 % of the values and positive values very
   less(11%)
```

## Check most common words in the dataset using word cloud

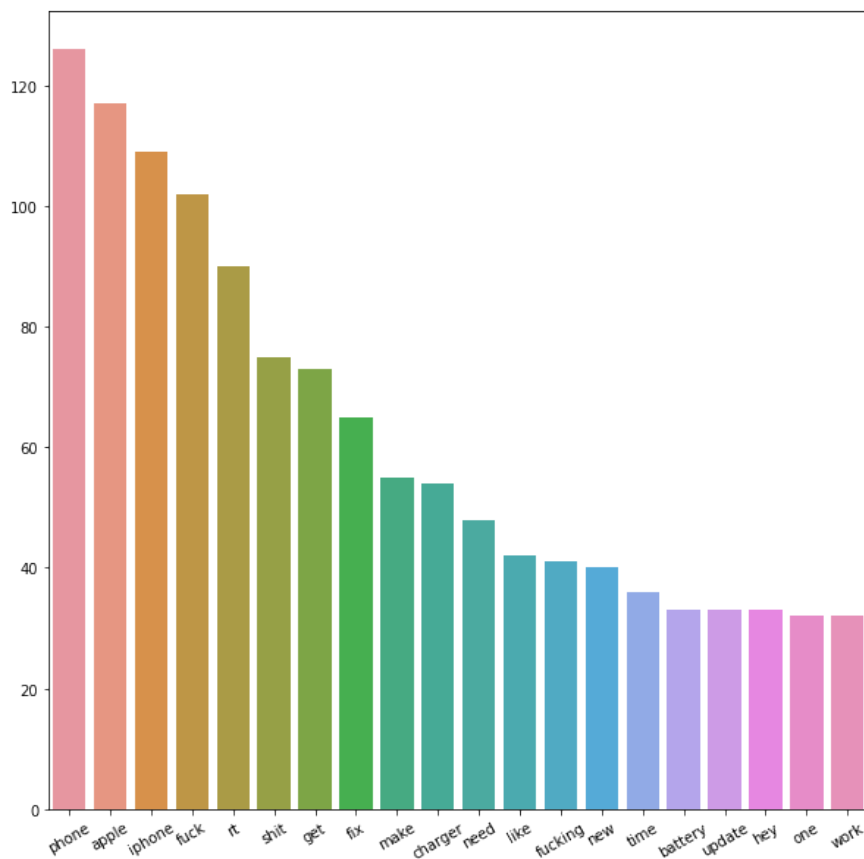visualization where frequent words appear enlarged as compared to less frequent words

**common words for the negative sentiments**

In [17]:
```python
#Negative words
negative_words = ' '.join(data.query("sentiment==1")['cleaned_text'])
wordcloud = WordCloud(width=800, height=500, random_state=21, max_font_size=110).generate(negative_words)

plt.figure(figsize=(10, 7))
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis('off');
```

```
In [18]:    1  ## Creating FreqDist for whole BoW , keeping the 20 most common tokens
            2  negative_tokens = word_tokenize(negative_words) #tokenize data
            3  all_fdist = FreqDist(negative_tokens).most_common(20)
            4
            5  ## Conversion to Pandas series via Python Dictionary for easier plotting
            6  all_fdist = pd.Series(dict(all_fdist))
            7
            8  ## Setting figure, ax into variables
            9  fig, ax = plt.subplots(figsize=(10,10))
           10
           11  ## Seaborn plotting using Pandas attributes + xtick rotation for ease of viewing
           12  all_plot = sns.barplot(x=all_fdist.index, y=all_fdist.values, ax=ax)
           13  plt.xticks(rotation=30);
```
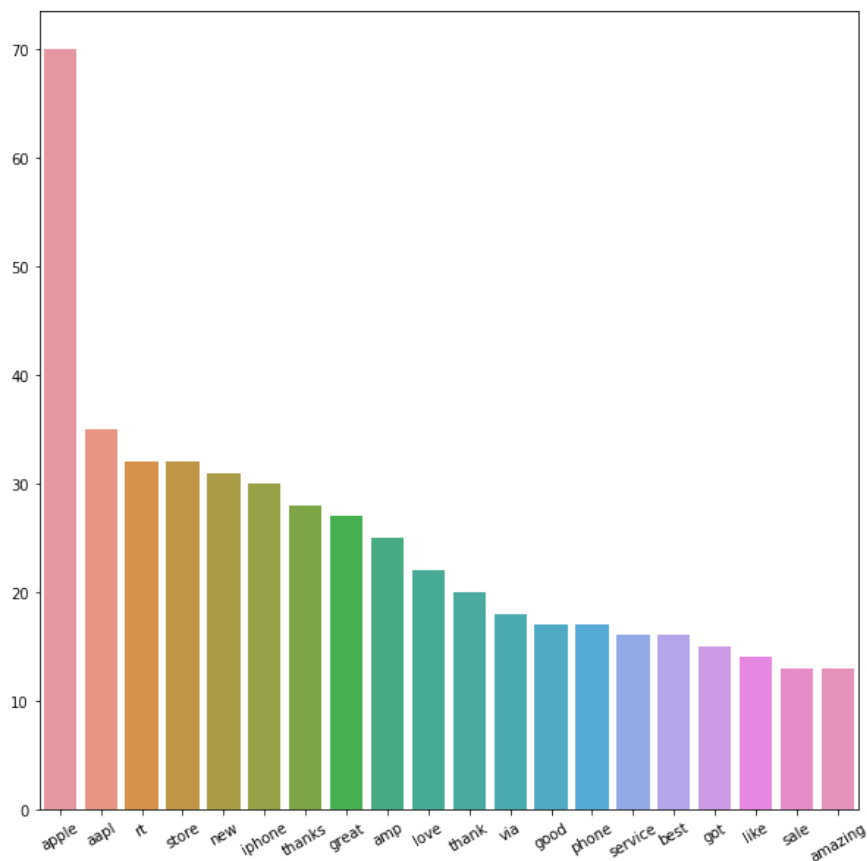


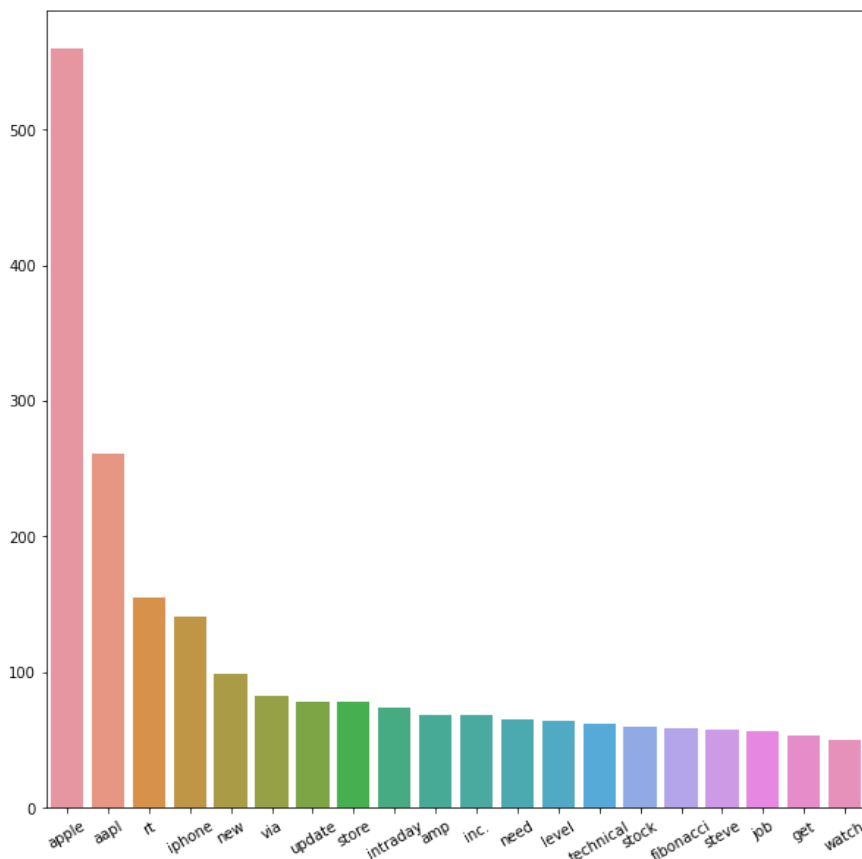**common words for the positive sentiments**

```python
In [19]:    1  #positive words
            2  positive_words = ' '.join(data.query('sentiment==5')['cleaned_text'])
            3  wordcloud = WordCloud(width=800, height=500, random_state=21 ,max_font_size=110).generate(positive_words)
            4
            5  plt.figure(figsize=(10,7))
            6  plt.imshow(wordcloud,interpolation='bilinear')
            7  plt.axis('off');
```

In [20]:

```python
## Creating FreqDist for whole BoW , keeping the 20 most common tokens
positive_tokens = word_tokenize(positive_words)
all_fdist = FreqDist(positive_tokens).most_common(20)

## Conversion to Pandas series via Python Dictionary for easier plotting
all_fdist = pd.Series(dict(all_fdist))

## Setting figure, ax into variables
fig, ax = plt.subplots(figsize=(10,10))

## Seaborn plotting using Pandas attributes + xtick rotation for ease of viewing
all_plot = sns.barplot(x=all_fdist.index, y=all_fdist.values, ax=ax)
plt.xticks(rotation=30);
```

```python
In [21]:    1  #neutral
            2  neutral_words = ' '.join(data.query("sentiment==3")['cleaned_text'])
            3  from wordcloud import WordCloud
            4  wordcloud = WordCloud(width=800, height=500, random_state=21, max_font_size=110).generate(neutral_words)
            5
            6  plt.figure(figsize=(10, 7))
            7  plt.imshow(wordcloud, interpolation="bilinear")
            8  plt.axis('off')
            9  plt.show()
```

```python
In [22]: ▾    1  ## Creating FreqDist for whole BoW , keeping the 20 most common tokens
             2  neutral_tokens = word_tokenize(neutral_words)
             3  all_fdist = FreqDist(neutral_tokens).most_common(20)
             4
             5  ## Conversion to Pandas series via Python Dictionary for easier plotting
             6  all_fdist = pd.Series(dict(all_fdist))
             7
             8  ## Setting figure, ax into variables
             9  fig, ax = plt.subplots(figsize=(10,10))
            10
            11  ## Seaborn plotting using Pandas attributes + xtick rotation for ease of viewing
            12  all_plot = sns.barplot(x=all_fdist.index, y=all_fdist.values, ax=ax)
            13  plt.xticks(rotation=30);
```



## step 6: Modelling

In this section i will conduct a train_test split for modelling then use pipelines to streamline vectorization, smote and modelling process. I will then evaluate the models using accuray_score,precision,recall,f1 score and roc_auc The idea behind choosing these models is that we want to run all the classifiers on the dataset ranging from simple ones to complex models, and then try to find the one which gives the best performance among them and tune it

```python
In [23]: ▾    1  #split data into x and y
             2  X = data['cleaned_text']
             3  y = data['sentiment']
             4
             5  #split data into a training and test set
             6  X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,stratify=y,shuffle=True,random_state=42)
```

We select our feature variable as X and target variable as y. Then split our data into train and test data for modelling training and evaluation purposes.Shuffle and stratify ensures a well-balanced dataset split,which improves model performance and evaluation reliability

```
In [24]:    1  #create pipeline to streamline vectorization and for modelling, start with base model
            2  #(Logistic Regression)
            3  pipeline = Pipeline([
            4      ('tdif',TfidfVectorizer(ngram_range=(1,2))),
            5      ('smote',SMOTE()),
            6      ('model',LogisticRegression())
            7  ])
```

Pipeline that does text vectorization i.e converting raw text into numerical features using tfidfVectorizer,smote to balance the target class 'sentiment' and using Logistic regression to classify text
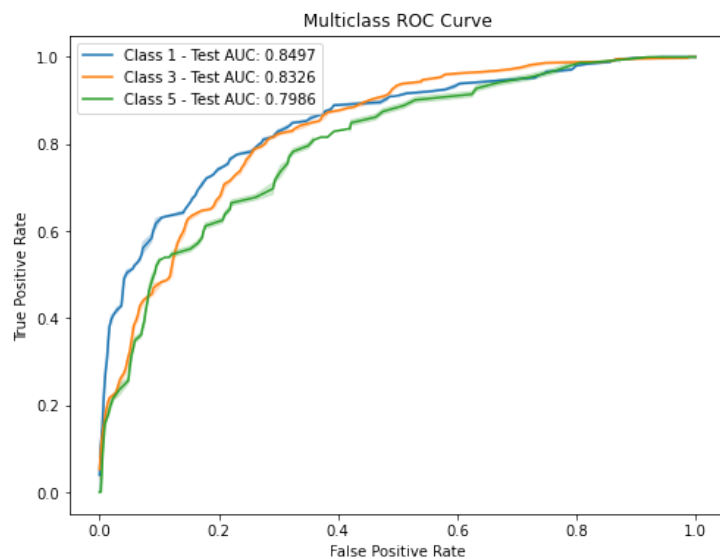
```
In [24]:    1  #create pipeline to streamline vectorization and for modelling, start with base model
            2  #(Logistic Regression)
            3  pipeline = Pipeline([
```

```python
In [25]:
1   #Function to train the model ,make predictions and calculate evaluation metrics
2   def modelling(pipe):
3       pipe.fit(X_train, y_train)
4
5       # Predict train and test data
6       y_hat_train = pipe.predict(X_train)
7       y_hat_test = pipe.predict(X_test)
8
9       # Get accuracy, precision, recall, and F1-score
10      base_train_accuracy = accuracy_score(y_train, y_hat_train)
11      base_test_accuracy = accuracy_score(y_test, y_hat_test)
12      base_train_precision = precision_score(y_train, y_hat_train, average='weighted')
13      base_test_precision = precision_score(y_test, y_hat_test, average='weighted')
14      base_train_recall = recall_score(y_train, y_hat_train, average='weighted')
15      base_test_recall = recall_score(y_test, y_hat_test, average='weighted')
16      base_train_f1 = f1_score(y_train, y_hat_train, average='weighted')
17      base_test_f1 = f1_score(y_test, y_hat_test, average='weighted')
18
19      # Binarize labels for multiclass ROC curve
20      classes = sorted(set(y_train))  # Get unique classes
21      y_train_bin = label_binarize(y_train, classes=classes)
22      y_test_bin = label_binarize(y_test, classes=classes)
23
24      # Get prediction scores
25      if hasattr(pipe, "decision_function"):
26          y_score_train = pipe.decision_function(X_train)
27          y_score_test = pipe.decision_function(X_test)
28      else:
29          y_score_train = pipe.predict_proba(X_train)
30          y_score_test = pipe.predict_proba(X_test)
31
32      # Compute ROC curve and AUC for each class
33      train_auc_list, test_auc_list = [], []
34      plt.figure(figsize=(8, 6))
35
36      for i in range(len(classes)):
37          train_fpr, train_tpr, _ = roc_curve(y_train_bin[:, i], y_score_train[:, i])
38          test_fpr, test_tpr, _ = roc_curve(y_test_bin[:, i], y_score_test[:, i])
39
40          train_auc = auc(train_fpr, train_tpr)
41          test_auc = auc(test_fpr, test_tpr)
42
43          train_auc_list.append(train_auc)
44          test_auc_list.append(test_auc)
45
46          sns.lineplot(x=test_fpr, y=test_tpr, label=f'Class {classes[i]} - Test AUC: {test_auc:.4f}')
47
48      # Average AUC
49      avg_train_auc = sum(train_auc_list) / len(train_auc_list)
50      avg_test_auc = sum(test_auc_list) / len(test_auc_list)
51
52      plt.xlabel("False Positive Rate")
53      plt.ylabel("True Positive Rate")
54      plt.title("Multiclass ROC Curve")
55      plt.legend()
56      plt.show()
57
58      return {
59          'Training Accuracy': base_train_accuracy,
60          'Test Accuracy': base_test_accuracy,
61  #       'Training precision': base_train_precision,
62          'Test precision': base_test_precision,
63  #       'Training recall': base_train_recall,
64          'Test recall': base_test_recall,
65  #       'Training f1_score': base_train_f1,
66          'Test f1_score': base_test_f1,
67          'Average Train AUC': avg_train_auc,
68          'Average Test AUC': avg_test_auc
69      }
70
```
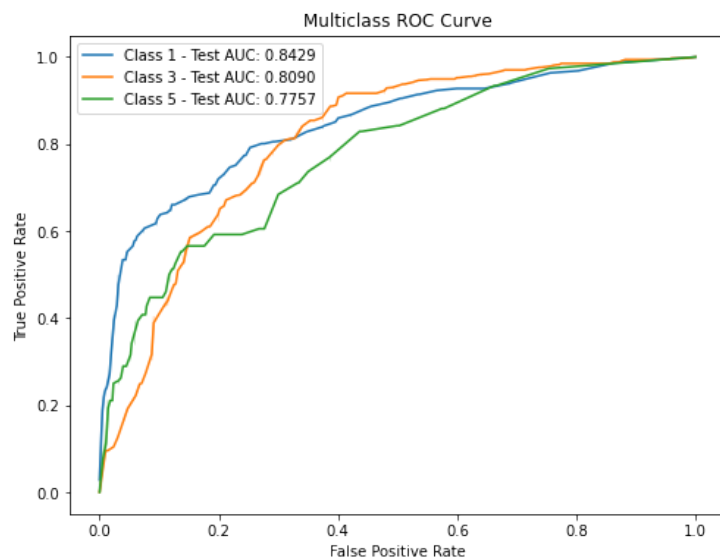
To avoid repetition for various models that we will be creating, we've defined a function that fits, predicts and calculates the accuracy_score

```
In [26]:    1  #logistic Regression evaluation
            2  logreg = modelling(pipeline)
            3  logreg
```
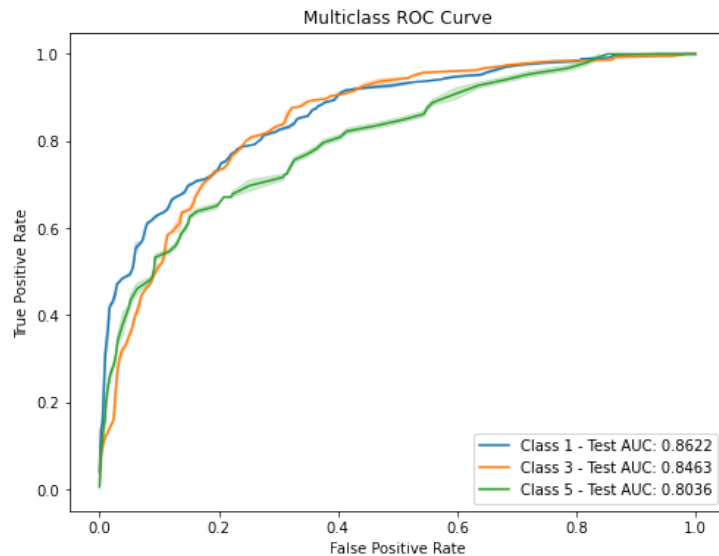
Multiclass ROC Curve



Out[26]: {'Training Accuracy': 0.9636219849742982,
          'Test Accuracy': 0.707740916271722,
          'Test precision': 0.7037161065532846,
          'Test recall': 0.707740916271722,
          'Test f1_score': 0.7054753826807029,
          'Average Train AUC': 0.9932681573749033,
          'Average Test AUC': 0.8269461973387789}

```
In [27]:    1  #randomForest
            2  pipeline.set_params(model=RandomForestClassifier(random_state=42))
            3  rdf = modelling(pipeline)
            4  rdf
```
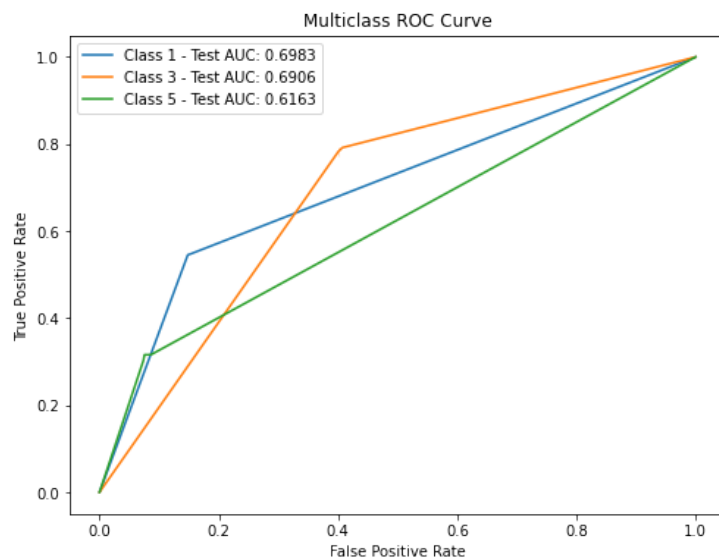
Multiclass ROC Curve



Out[27]: {'Training Accuracy': 0.9905100830367735,
          'Test Accuracy': 0.7045813586097947,
          'Test precision': 0.7196297646534613,
          'Test recall': 0.7045813586097947,
          'Test f1_score': 0.681570616869526,
          'Average Train AUC': 0.9995458299107819,
          'Average Test AUC': 0.8091915324598483}

In [28]:
```python
1  #multinomial naive bayes
2  pipeline.set_params(model = MultinomialNB())
3  nb =  modelling(pipeline)
4  nb
```

Multiclass ROC Curve

Class 1 - Test AUC: 0.8622
Class 3 - Test AUC: 0.8463
Class 5 - Test AUC: 0.8036

Out[28]: {'Training Accuracy': 0.9418742586002372,
          'Test Accuracy': 0.6824644549763034,
          'Test precision': 0.7133060787791613,
          'Test recall': 0.6824644549763034,
          'Test f1_score': 0.6913711284026903,
          'Average Train AUC': 0.9889344354190913,
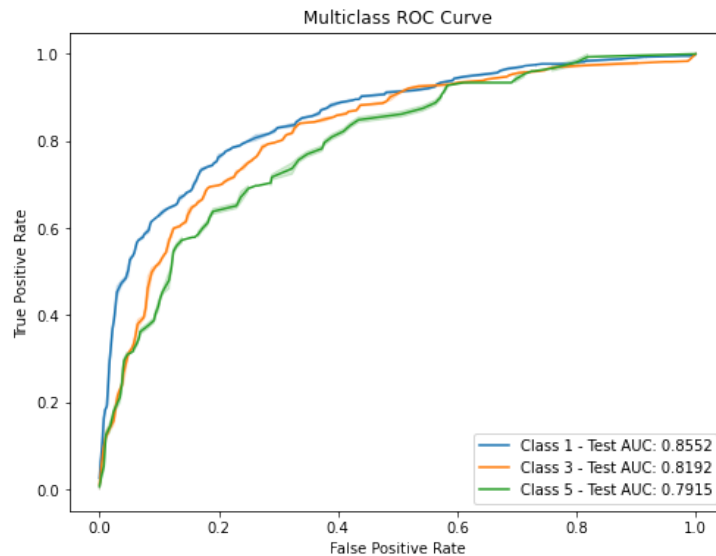          'Average Test AUC': 0.8373466293718987}

In [29]:
```python
1  #Decision tree
2  pipeline.set_params(model = DecisionTreeClassifier(random_state=42))
3  dt = modelling(pipeline)
4  dt
```

Multiclass ROC Curve

Class 1 - Test AUC: 0.6983
Class 3 - Test AUC: 0.6906
Class 5 - Test AUC: 0.6163

Out[29]: {'Training Accuracy': 0.9905100830367735,
          'Test Accuracy': 0.6477093206951027,
          'Test precision': 0.6397823451988728,
          'Test recall': 0.6477093206951027,
          'Test f1_score': 0.6397023128888808,
          'Average Train AUC': 0.999876025147195,
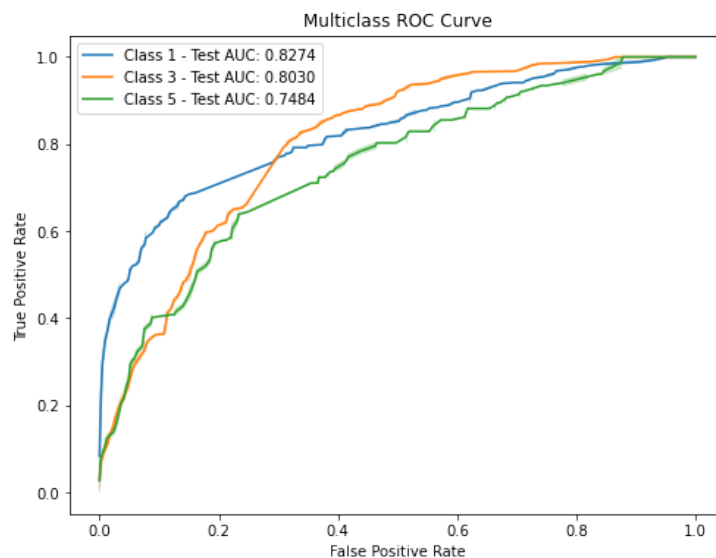          'Average Test AUC': 0.6684032590271273}

In [30]:
```
1  #svc
2  pipeline.set_params(model = svm.SVC())
3  svm = modelling(pipeline)
4  svm
```

Multiclass ROC Curve

Class 1 - Test AUC: 0.8552
Class 3 - Test AUC: 0.8192
Class 5 - Test AUC: 0.7915

Out[30]: {'Training Accuracy': 0.9806247528667458,
 'Test Accuracy': 0.7061611374407583,
 'Test precision': 0.6933178216904188,
 'Test recall': 0.7061611374407583,
 'Test f1_score': 0.6850618166156861,
 'Average Train AUC': 0.9953301054763962,
 'Average Test AUC': 0.8219852591406278}

In [31]:
```
1  pipeline.set_params(model = XGBClassifier())
2  xgb = modelling(pipeline)
3  xgb
```

Multiclass ROC Curve

Class 1 - Test AUC: 0.8274
Class 3 - Test AUC: 0.8030
Class 5 - Test AUC: 0.7484

Out[31]: {'Training Accuracy': 0.8912613681296956,
 'Test Accuracy': 0.6919431279620853,
 'Test precision': 0.7020677275931184,
 'Test recall': 0.6919431279620853,
 'Test f1_score': 0.6728253975054899,
 'Average Train AUC': 0.9733294583555869,
 'Average Test AUC': 0.7929194002687371}

```
In [32]:    1  # Dictionary of model results
            2  model_results = {
            3      "Logistic Regression": logreg,
            4      "Random Forest": rdf,
            5      "Naïve Bayes":nb,
            6      "Decision Tree": dt,
            7      "SVM": svm,
            8      "XGBoost": xgb,
            9  }
           10
           11  # Convert dictionary to DataFrame
           12  df_results = pd.DataFrame.from_dict(model_results, orient='index')
           13
           14  # Display the DataFrame
           15  df_results
           16
```

Out[32]:

|  | Training Accuracy | Test Accuracy | Test precision | Test recall | Test f1_score | Average Train AUC | Average Test AUC |
|---|---|---|---|---|---|---|---|
| **Logistic Regression** | 0.963622 | 0.707741 | 0.703716 | 0.707741 | 0.705475 | 0.993268 | 0.826946 |
| **Random Forest** | 0.990510 | 0.704581 | 0.719630 | 0.704581 | 0.681571 | 0.999546 | 0.809192 |
| **Naïve Bayes** | 0.941874 | 0.682464 | 0.713306 | 0.682464 | 0.691371 | 0.988934 | 0.837347 |
| **Decision Tree** | 0.990510 | 0.647709 | 0.639782 | 0.647709 | 0.639702 | 0.999876 | 0.668403 |
| **SVM** | 0.980625 | 0.706161 | 0.693318 | 0.706161 | 0.685062 | 0.995330 | 0.821985 |
| **XGBoost** | 0.891261 | 0.691943 | 0.702068 | 0.691943 | 0.672825 | 0.973329 | 0.792919 |

```
In [33]:    1  #checking gaps between train and test accuracy
            2  df_results['Training Accuracy'] - df_results['Test Accuracy']
```

```
Out[33]:  Logistic Regression    0.255881
          Random Forest          0.285929
          Naïve Bayes            0.259410
          Decision Tree          0.342801
          SVM                    0.274464
          XGBoost                0.199318
          dtype: float64
```

```
1  From the models tested we see varying degrees of overfitting and generalization ability across different
   models.Decision trees and random Forest performance might fail in real word due to overfitting
2
3  Higher AUC is seen in the models indicating models performs well in separating the classes
4  XGBoost has the best generalization as it has the smallest gap between the train and test accuracy
5
6  on simple models logistic regression performs well with a good balance between precision and recall
7
8  I will compare XGBoost and Logistic regression by tuning to determine the best model
```

## Step:7 Hyperparameter tuning

We'll tune the following parameters: selecting optimal parameters for the model

- **n_estimators**: Number of boosting rounds
- **max_depth**: maximum tree depth (higher=more complex model)
- **learning_rate**: Step size to prevent overfitting
- **subsmaple**: Fraction of samples used per tree(lower=more regularizationO
- **colsample_bytree**: fraction of features used per tree
- **gamma**: minimum loss reduction reqquired for further splitting

```
In [34]:    1  #Define parameter grid
            2  param_grid = {
            3      'model__n_estimators': [100, 300, 500],
            4      'model__max_depth': [3, 5, 7],
            5      'model__learning_rate': [ 0.1, 0.2,0.4],
            6      'model__subsample': [0.7, 0.8, 1.0],
            7      'model__colsample_bytree': [0.7, 0.8, 1.0],
            8      'model__gamma': [0, 0.1, 0.2]
            9  }
```

In [35]: ▾
```
1  #check pipeline
2  pipeline
```

Out[35]:
```
▸    Pipeline

▸ TfidfVectorizer

      ▸ SMOTE

▸ XGBClassifier
```

In [36]: ▾
```
1  grid_search = GridSearchCV(
2                              estimator = pipeline,
3                              param_grid = param_grid,
4                              scoring = 'accuracy',
5                              verbose=2,
6                              cv = 5,
7                              n_jobs=-1
8  )
9  grid_search.fit(X_train, y_train)
10
11 #Get Best parameters
12 grid_search.best_params_
13
```

```
Fitting 5 folds for each of 729 candidates, totalling 3645 fits
```
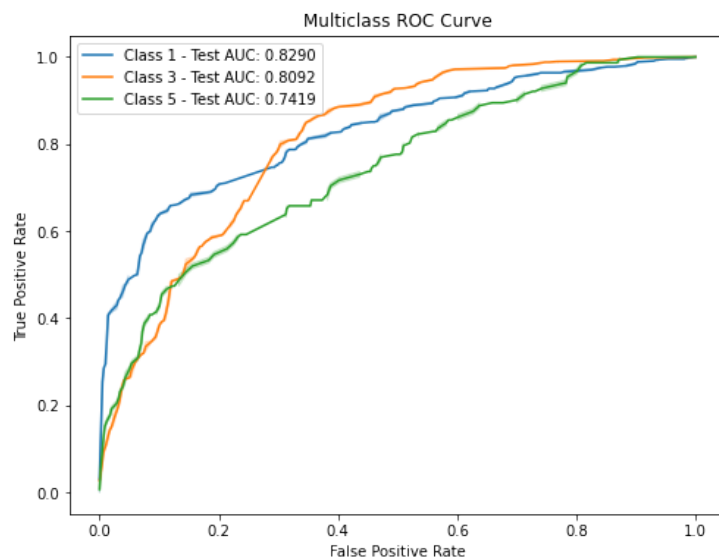
Out[36]:
```
{'model__colsample_bytree': 0.7,
 'model__gamma': 0,
 'model__learning_rate': 0.2,
 'model__max_depth': 3,
 'model__n_estimators': 500,
 'model__subsample': 1.0}
```

## Retraining the model with the best parameters

In [51]: ▾
```
1  pipeline.set_params(model=XGBClassifier(
2                              learning_rate= 0.2,
3                              max_depth= 3,
4                              n_estimators= 500,
5                              colsample_bytree =  0.7,
6                              subsample = 1.0
7  ))
```

Out[51]:
```
▸    Pipeline

▸ TfidfVectorizer

      ▸ SMOTE

▸ XGBClassifier
```

```
In [52]:   1  #Evaluate the model with the best parameters
           2  xgb_tuned = modelling(pipeline)
           3  xgb_tuned
```

Multiclass ROC Curve



```
Out[52]: {'Training Accuracy': 0.9161724001581653,
          'Test Accuracy': 0.7014218009478673,
          'Test precision': 0.7031791809042994,
          'Test recall': 0.7014218009478673,
          'Test f1_score': 0.6826564554491085,
          'Average Train AUC': 0.9858007103440465,
          'Average Test AUC': 0.7933539582783885}
```

```
In [53]:   1  #comparing with the results before the tuning
           2  df_results
```

Out[53]:

|  | Training Accuracy | Test Accuracy | Test precision | Test recall | Test f1_score | Average Train AUC | Average Test AUC |
|---|---|---|---|---|---|---|---|
| **Logistic Regression** | 0.963622 | 0.707741 | 0.703716 | 0.707741 | 0.705475 | 0.993268 | 0.826946 |
| **Random Forest** | 0.990510 | 0.704581 | 0.719630 | 0.704581 | 0.681571 | 0.999546 | 0.809192 |
| **Naïve Bayes** | 0.941874 | 0.682464 | 0.713306 | 0.682464 | 0.691371 | 0.988934 | 0.837347 |
| **Decision Tree** | 0.990510 | 0.647709 | 0.639782 | 0.647709 | 0.639702 | 0.999876 | 0.668403 |
| **SVM** | 0.980625 | 0.706161 | 0.693318 | 0.706161 | 0.685062 | 0.995330 | 0.821985 |
| **XGBoost** | 0.891261 | 0.691943 | 0.702068 | 0.691943 | 0.672825 | 0.973329 | 0.792919 |

```
In [54]:   1  xgb_tuned['Average Test AUC'] - df_results['Average Test AUC']
```

```
Out[54]: Logistic Regression    -0.033592
         Random Forest          -0.015838
         Naïve Bayes            -0.043993
         Decision Tree           0.124951
         SVM                    -0.028631
         XGBoost                 0.000435
         Name: Average Test AUC, dtype: float64
```
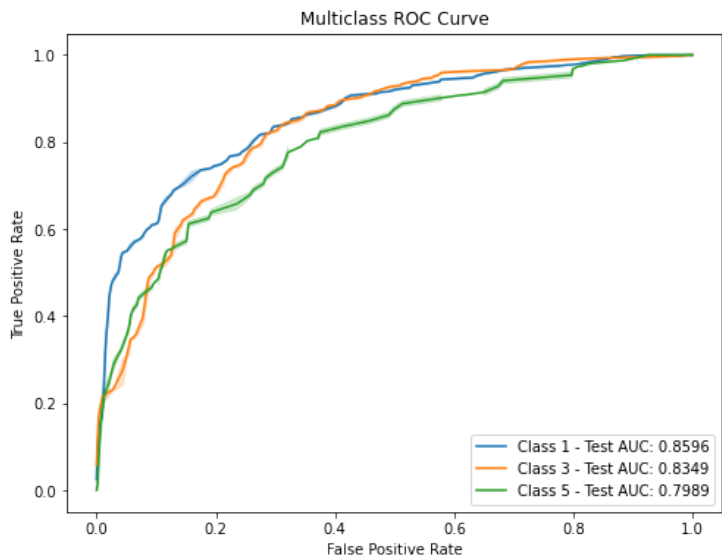
```
           1  No major improvement in accuracy – The test accuracy and F1 score increased by 1% or less
           2  AUC improved slightly – The Test AUC actually increased by just  0.0004
           3  Tuning reduced slight overfitting –  meaning the model is generalizing slightly better.
```

**Tune Logistic Regression**

In [56]:
```
1  pipeline.set_params(model = LogisticRegression())
2  param_grid = {
3      'model__C': [0.01, 0.1, 1, 10, 100],  # Regularization strength
4      'model__penalty': ['l1', 'l2'],  # Type of regularization
5      'model__solver': ['liblinear']  # Required for L1 penalty
6  }
7
8  grid = GridSearchCV(pipeline, param_grid, cv=5, scoring='accuracy',n_jobs=-1)
9  grid.fit(X_train, y_train)
10
11 print("Best Parameters:", grid.best_params_)
12 print("Best accuracy Score:", grid.best_score_)
```

```
Best Parameters: {'model__C': 10, 'model__penalty': 'l2', 'model__solver': 'liblinear'}
Best accuracy Score: 0.7018667084099715
```

In [57]:
```
1  pipeline.set_params(model = LogisticRegression(
2                                          C= 10,
3                                          penalty = 'l2',
4                                          solver = 'liblinear'
5  ))
6  logistic_tuned= modelling(pipeline)
7  logistic_tuned
```


Multiclass ROC Curve

Legend:
- Class 1 - Test AUC: 0.8596
- Class 3 - Test AUC: 0.8349
- Class 5 - Test AUC: 0.7989

Out[57]:
```
{'Training Accuracy': 0.9837880585211546,
 'Test Accuracy': 0.7251184834123223,
 'Test precision': 0.7219159277378586,
 'Test recall': 0.7251184834123223,
 'Test f1_score': 0.7233076101045702,
 'Average Train AUC': 0.9985778434373468,
 'Average Test AUC': 0.8311354483998646}
```

In [58]:
```
1  #compare with previous results
2  df_results
```

Out[58]:

|  | Training Accuracy | Test Accuracy | Test precision | Test recall | Test f1_score | Average Train AUC | Average Test AUC |
|---|---|---|---|---|---|---|---|
| **Logistic Regression** | 0.963622 | 0.707741 | 0.703716 | 0.707741 | 0.705475 | 0.993268 | 0.826946 |
| **Random Forest** | 0.990510 | 0.704581 | 0.719630 | 0.704581 | 0.681571 | 0.999546 | 0.809192 |
| **Naïve Bayes** | 0.941874 | 0.682464 | 0.713306 | 0.682464 | 0.691371 | 0.988934 | 0.837347 |
| **Decision Tree** | 0.990510 | 0.647709 | 0.639782 | 0.647709 | 0.639702 | 0.999876 | 0.668403 |
| **SVM** | 0.980625 | 0.706161 | 0.693318 | 0.706161 | 0.685062 | 0.995330 | 0.821985 |
| **XGBoost** | 0.891261 | 0.691943 | 0.702068 | 0.691943 | 0.672825 | 0.973329 | 0.792919 |

```
1  There is an increase by 2% in accuracy, precision, recall and F1 score.
```

```
2  lets compare XG Boost and logistic regression
```

In [61]:
```python
1  model_results = {
2      "Logistic Regression_tuned": logistic_tuned,
3      "XGBoost_tuned": xgb_tuned,
4  }
5
6  # Convert dictionary to DataFrame
7  df_results_tuned = pd.DataFrame.from_dict(model_results, orient='index')
8
9  # Display the DataFrame
10 df_results_tuned
```

Out[61]:

| | Training Accuracy | Test Accuracy | Test precision | Test recall | Test f1_score | Average Train AUC | Average Test AUC |
|---|---|---|---|---|---|---|---|
| **Logistic Regression_tuned** | 0.983788 | 0.725118 | 0.721916 | 0.725118 | 0.723308 | 0.998578 | 0.831135 |
| **XGBoost_tuned** | 0.916172 | 0.701422 | 0.703179 | 0.701422 | 0.682656 | 0.985801 | 0.793354 |

```
1  Logistic regression performs much better than XG Boost.We will logistic regression model to make our
   predictions
```

## Step8: Make predictions

In [63]:
```python
1  #create function to predict sentiment
2  pipeline.set_params(model = LogisticRegression(
3                                      C= 10,
4                                      penalty = 'l2',
5                                      solver = 'liblinear'))
6  pipeline.fit(X_train,y_train)
7  def sentiment_check(tweet):
8      tweet_processed = clean_text(tweet)
9      print(f'tweet:{tweet_processed}')
10     prediction = pipeline.predict([tweet_processed])
11     return "Negative Statement" if prediction[0] == 1 else "Neutral Statemnet" if prediction[0]==3 else "Posit
12
13 print(sentiment_check('this phones suck,its not what is marketed'))
```

```
tweet:phone suck marketed
Negative Statement
```

In [64]:
```python
1  print(sentiment_check("I bought a new phone and it's so good"))
```

```
tweet:bought new phone good
Positive statement
```

In [65]:
```python
1  print(sentiment_check("great item, good job"))
```

```
tweet:great item good job
Positive statement
```

In [66]:
```python
1  print(sentiment_check("what a stupid cover"))
```

```
tweet:stupid cover
Negative Statement
```

**Our model performs relatively well**

we can look into using deep learning and other vectorization techniques to see if we can improve our models performance

## Step 9: Conclusions and recommendations

## Conclusion

- Logistic Regression outperformed XGBoost, achieving better accuracy and F1 score, making it the preferred model for sentiment classification.
- Hyperparameter tuning for XGBoost led to only marginal improvements (accuracy & F1 score increased by ≤1%, AUC by 0.0004).
- Tuning helped reduce overfitting slightly, meaning XGBoost generalized better than before, but still did not surpass Logistic Regression.

**Recommendations**

- Use Logistic Regression for final predictions since it performs better than XGBoost.
- Consider feature engineering (e.g., word embeddings like Word2Vec or BERT) to improve model performance further.
- Explore deep learning models (e.g., LSTMs or Transformers) if higher accuracy is required.
- Continue tuning XGBoost or test alternative ensemble methods if needed for comparison.

In [ ]:    1