

# Sentiment Analysis for Marketing: Understanding Customer Preferences through Data

## Phase 4: development part 2

Done by:

**Manikandan .A**

**311421104306**

**Meenakshi College of Engineering - 3114**

**B.E CSE 3rd yr. 5th SEM**

## Description:

In this technology we will continue building our project by selecting a machine learning algorithm, training the model, and evaluating its performance. Perform different analysis as needed

## BERT Model and Tokenizer Setup

```
In [3]: # Import the necessary libraries for BERT model and tokenizer
from transformers import BertModel, BertTokenizer

# Define the pre-trained BERT model's name
model_name = "bert-base-uncased"

# Load the BERT model with the specified name
model = BertModel.from_pretrained(model_name)

# Load the BERT tokenizer corresponding to the model
tokenizer = BertTokenizer.from_pretrained(model_name)
```

In [5]:

dataset

Out[5]:

|       | Tweet ID           | Airline<br>Sentiment | Sentiment<br>Confidence | Negative<br>Reason           | Negative<br>Reason<br>Confidence | Airline       | Air<br>Sentim<br>G |
|-------|--------------------|----------------------|-------------------------|------------------------------|----------------------------------|---------------|--------------------|
| 0     | 570066283233972224 | positive             | 0.9657                  | NaN                          | 0.0000                           | Southwest     | pos                |
| 1     | 568561924985782272 | positive             | 1.0000                  | NaN                          | NaN                              | US<br>Airways | pos                |
| 2     | 570264948548313088 | positive             | 1.0000                  | NaN                          | NaN                              | American      | pos                |
| 3     | 569600137296633920 | positive             | 0.9236                  | NaN                          | 0.0000                           | American      | pos                |
| 4     | 568551906634797120 | positive             | 0.6242                  | NaN                          | 0.0000                           | American      | pos                |
| ...   | ...                | ...                  | ...                     | ...                          | ...                              | ...           | ...                |
| 14599 | 569587686496825280 | positive             | 0.3487                  | NaN                          | 0.0000                           | American      | ↑                  |
| 14600 | 569587371693355008 | negative             | 1.0000                  | Customer<br>Service<br>Issue | 1.0000                           | American      | ↑                  |
| 14601 | 569587242672398272 | neutral              | 1.0000                  | NaN                          | NaN                              | American      | ↑                  |
| 14602 | 569587188687634432 | negative             | 1.0000                  | Customer<br>Service<br>Issue | 0.6659                           | American      | ↑                  |
| 14603 | 569587140490866688 | neutral              | 0.6771                  | NaN                          | 0.0000                           | American      | ↑                  |

14604 rows × 14 columns

# Creating BERT Embeddings for Text Data using PyTorch

```
In [8]: # Import the necessary library
import torch

# Define a function to get BERT embeddings for input text
def get_bert_embeddings(text):
    # Convert the input text to a PyTorch tensor and add a batch dimension
    input_ids = torch.tensor(text).unsqueeze(0)

    # Perform inference with the BERT model (assuming 'model' is previously
    with torch.no_grad():
        outputs = model(input_ids)
        # Calculate the mean of the last hidden state to obtain embeddings
        embeddings = outputs.last_hidden_state.mean(dim=-1)

    return embeddings

dataset['bert_embeddings'] = dataset['tokenized_text'].apply(get_bert_embed
```

## Topic Modeling with Latent Dirichlet Allocation (LDA) for Text Data

```
In [14]: # Import necessary Libraries
from gensim import corpora, models

# Load the dataset from an Excel file
dataset = pd.read_excel("output_data_updated_with_sentiment.xlsx")

# Tokenize the text if not already tokenized
dataset['tokenized_text'] = dataset['tokenized_text'].apply(lambda x: x.split())

# Create a dictionary of words from the tokenized text
dictionary = corpora.Dictionary(dataset['tokenized_text'])

# Convert the dictionary to a bag-of-words representation
corpus = [dictionary.doc2bow(text) for text in dataset['tokenized_text']]

# Perform LDA (Latent Dirichlet Allocation) topic modeling with 5 topics and 15 passes
lda = models.LdaModel(corpus, num_topics=5, id2word=dictionary, passes=15)

# Assign topics to each review in the dataset
dataset['topic'] = [lda[corpus[i]] for i in range(len(dataset))]
```

## Text Data Topic Modeling using Latent Dirichlet Allocation (LDA) and Gensim

```
In [18]: from gensim import corpora, models

# Load the dataset
dataset = pd.read_excel("output_data_updated_with_sentiment.xlsx")

# Tokenize the text if it's not already tokenized (you might skip this if it is)
# Create a dictionary of words
dataset['tokenized_text'] = dataset['tokenized_text'].apply(lambda x: x.split())

dictionary = corpora.Dictionary(dataset['tokenized_text'])

# Convert the dictionary to a bag-of-words representation
corpus = [dictionary.doc2bow(text) for text in dataset['tokenized_text']]

# Perform LDA topic modeling
lda = models.LdaModel(corpus, num_topics=5, id2word=dictionary, passes=15)

# Assign topics to each review
dataset['topic'] = [lda[corpus[i]] for i in range(len(dataset))]
```

## Setting Up NLP Tools with spaCy and PyTorch

```
In [22]: import spacy
import torch

# Set the default dtype and device for PyTorch
torch.set_default_dtype(torch.float32) # Set the default data type
torch.set_default_device("cpu") # Set the default device (e.g., "cpu" or "cuda")

# Load the spaCy NER model
nlp = spacy.load("en_core_web_sm")
```

## Extracting Named Entities from Text Data using spaCy NER Model

```
In [24]: dataset['entities'] = dataset['Text'].apply(lambda x: [(ent.text, ent.label_) for ent in nlp(x).ent_type_entities])
```

In [25]: dataset

Out[25]:

| line       | Airline<br>Sentiment<br>Gold | Name            | Negative<br>Reason<br>Gold | Text  | Tweet<br>Created           | Tweet<br>Location                    | User<br>Timezone                    | token  |
|------------|------------------------------|-----------------|----------------------------|---|----------------------------|--------------------------------------|-------------------------------------|--------|
| vest       | positive                     | magmum03        | NaN                        | southwestair<br>awesome<br>flight dallas<br>2 ny virgin...    | 2015-<br>02-23<br>19:42:47 | NaN                                  | NaN                                 | ['[10' |
| US<br>rays | positive                     | christinachime  | NaN                        | usairways<br>thank finally<br>got bag<br>customer<br>servi... | 2015-<br>02-19<br>16:05:00 | NaN                                  | NaN                                 | ['[10' |
| ican       | positive                     | Runts54         | NaN                        | americanair<br>dfwairport 2<br>together<br>best part fl...    | 2015-<br>02-24<br>08:52:13 | Eules,<br>Texas                      | NaN                                 | '431;  |
| ican       | positive                     | douglaskgordon  | NaN                        | americanair<br>thank<br>youyou                                | 2015-<br>02-22<br>12:50:30 | Caribbean,<br>New York<br>and Miami. | Indiana<br>(East)                   | '431;  |
| ican       | positive                     | byunsamuel      | NaN                        | americanair<br>hopefully<br>see bad<br>ones<br>opportunity... | 2015-<br>02-19<br>15:25:12 | Sunnyside,<br>NY                     | Eastern<br>Time (US<br>&<br>Canada) |        |
| ...        | ...                          | ...             | ...                        | ...   | ...                        | ...                                  | ...                                 |        |
| ican       | NaN                          | KristenReenders | NaN                        | americanair<br>thank got<br>different<br>flight<br>chicago    | 2015-<br>02-22<br>12:01:01 | NaN                                  | NaN                                 | '431;  |
| ican       | NaN                          | itsropes        | NaN                        | americanair<br>leaving 20<br>minutes late<br>flight war...    | 2015-<br>02-22<br>11:59:46 | Texas                                | NaN                                 | '431;  |
| ican       | NaN                          | sanyabun        | NaN                        | americanair<br>please bring<br>american<br>airlines bla...    | 2015-<br>02-22<br>11:59:15 | Nigeria,lagos                        | NaN                                 | '431;  |
| ican       | NaN                          | SraJackson      | NaN                        | americanair<br>money<br>change flight<br>dont answer<br>ph... | 2015-<br>02-22<br>11:59:02 | New Jersey                           | Eastern<br>Time (US<br>&<br>Canada) | '431;  |

| line | Airline<br>Sentiment<br>Gold | Name      | Negative<br>Reason<br>Gold | Text   | Tweet<br>Created           | Tweet<br>Location | User<br>Timezone | token |
|------|------------------------------|-----------|----------------------------|--|----------------------------|-------------------|------------------|-------|
| ican | NaN                          | daviddtwu | NaN                        | americanair<br>8 ppl need 2<br>know many<br>seats next ... | 2015-<br>02-22<br>11:58:51 | dallas, TX        | NaN              | '431; |

## Extracting Keywords from Text Data using NLTK and Stopwords

```
In [28]: import pandas as pd
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.probability import FreqDist

# Load your dataset
dataset = pd.read_excel("output_data_updated_with_NER.xlsx")

# Function to extract keywords from a text
def extract_keywords(text, num_keywords=5):
    # Tokenize the text
    words = word_tokenize(text)

    # Remove stopwords
    stop_words = set(stopwords.words("english"))
    words = [word.lower() for word in words if word.isalnum() and word.lower() not in stop_words]

    # Calculate word frequency
    freq_dist = FreqDist(words)

    # Get the most common words as keywords
    keywords = [word for word, freq in freq_dist.most_common(num_keywords)]

    return ", ".join(keywords)

# Extract keywords from each review and add them to your dataset
dataset['keywords'] = dataset['Text'].apply(lambda x: extract_keywords(x, num_keywords=5))
```

## Emotion Analysis from Sentiment Scores using TextBlob

In [32]:

```
from textblob import TextBlob

# Load your dataset
dataset = pd.read_excel("output_data_updated_with_keywords.xlsx")

# Define a function to get the emotion from sentiment
def get_emotion(text):
    analysis = TextBlob(text)
    sentiment_score = analysis.sentiment.polarity

    # Map sentiment score to emotions (you can adjust the thresholds)
    if sentiment_score > 0.2:
        return "joy"
    elif sentiment_score < -0.2:
        return "sadness"
    else:
        return "neutral"

# Apply the function to your dataset
dataset['emotion'] = dataset['Text'].apply(get_emotion)
```

## Calculating Readability Metrics for Text Data using textstat

In [36]:

```
import pandas as pd
from textstat import flesch_kincaid_grade, gunning_fog

# Load your dataset
dataset = pd.read_excel("output_data_updated_with_emotions.xlsx")

# Calculate readability metrics for each review
dataset['flesch_kincaid'] = dataset['Text'].apply(lambda x: flesch_kincaid_
dataset['gunning_fog'] = dataset['Text'].apply(lambda x: gunning_fog(x))
```



In [37]: dataset

Out[37]:

| Negative Reason Gold | Text  | ... | User Timezone              | tokenized_text                                    | bert_embeddings                                  | sentiment   | topic   |
|----------------------|---|-----|----------------------------|---|--|---|---|
| NaN                  | southwestair<br>awesome<br>flight dallas<br>2 ny virgin...    | ... | NaN                        | ['[101,', '4943,', '11215,', '12476,', '3462,'... | tensor([[ 5.6657e-02, 2.0632e-02, 2.6626e-01...  | {'neg': 0.203, 'neu': 0.617, 'pos': 0.181, 'co... | [(0, 0.13657132), (1, 0.6401591), (2, 0.110866... |
| NaN                  | usairways<br>thank finally<br>got bag<br>customer<br>servi... | ... | NaN                        | ['[101,', '3915,', '4313,', '14035,', '4067,'...  | tensor([[ 1.4389e-01, -8.2253e-03, 3.2305e-01... | {'neg': 0.0, 'neu': 0.53, 'pos': 0.47, 'compou... | [(0, 0.41092548), (1, 0.013585756), (2, 0.0135... |
| NaN                  | americanair<br>dfwairport 2<br>together<br>best part fl...    | ... | NaN                        | ['[101,', '25988,', '4313,', '1040,', '2546,'...  | tensor([[-1.3857e-01, -1.6197e-01, 2.9042e-01... | {'neg': 0.0, 'neu': 0.588, 'pos': 0.412, 'comp... | [(0, 0.11536198), (1, 0.013556691), (2, 0.0135... |
| NaN                  | americanair<br>thank<br>youyou                                | ... | Indiana (East)             | ['[101,', '25988,', '4313,', '4067,', '2017,'...  | tensor([[-9.3035e-02, 2.2016e-01, 3.7216e-01...  | {'neg': 0.0, 'neu': 0.444, 'pos': 0.556, 'comp... | [(0, 0.025300583), (1, 0.025211947), (2, 0.025... |
| NaN                  | americanair<br>hopefully<br>see bad<br>ones<br>opportunity... | ... | Eastern Time (US & Canada) | ['[101,', '25988,', '4313,', '11504,', '2156,'... | tensor([[-5.2652e-02, 1.5723e-01, 2.8512e-01...  | {'neg': 0.139, 'neu': 0.239, 'pos': 0.622, 'co... | [(0, 0.011991125), (1, 0.012075231), (2, 0.648... |
| ...                  | ...   | ... | ...                        | ...   | ...  | ...   | ...   |
| NaN                  | americanair<br>thank got<br>different<br>flight<br>chicago    | ... | NaN                        | ['[101,', '25988,', '4313,', '4067,', '2288,'...  | tensor([[ 8.8884e-02, -2.7573e-01, 1.8394e-01... | {'neg': 0.0, 'neu': 0.667, 'pos': 0.333, 'comp... | [(0, 0.44322193), (1, 0.020745821), (2, 0.0202... |
| NaN                  | americanair<br>leaving 20<br>minutes late<br>flight war...    | ... | NaN                        | ['[101,', '25988,', '4313,', '2975,', '2322,'...  | tensor([[-9.3935e-02, -1.3353e-01, 2.3228e-01... | {'neg': 0.279, 'neu': 0.721, 'pos': 0.0, 'comp... | [(0, 0.545215), (3, 0.42837554)]                  |
| NaN                  | americanair<br>please bring<br>american<br>airlines bla...    | ... | NaN                        | ['[101,', '25988,', '4313,', '3531,', '3288,'...  | tensor([[-1.0209e-01, -2.3004e-01, 2.5336e-01... | {'neg': 0.0, 'neu': 0.685, 'pos': 0.315, 'comp... | [(0, 0.23062189), (1, 0.018351018), (2, 0.0183... |
| NaN                  | americanair<br>money<br>change flight<br>dont answer<br>ph... | ... | Eastern Time (US & Canada) | ['[101,', '25988,', '4313,', '2769,', '2689,'...  | tensor([[ 1.0850e-01, -2.3711e-01, 3.0256e-01... | {'neg': 0.0, 'neu': 0.776, 'pos': 0.224, 'comp... | [(0, 0.013468066), (1, 0.30526412), (2, 0.3808... |

| Negative Reason Gold | Text   | ... | User Timezone | tokenized_text  | bert_embeddings  | sentiment   | topic   |
|----------------------|--|-----|---------------|---|--|---|---|
| NaN                  | americanair<br>8 ppl need 2<br>know many<br>seats next ... | ... | NaN           | ['[101,',<br>'25988',<br>'4313,', '1022',<br>'4903',... | tensor([[[-7.1985e-<br>02, -2.6482e-01,<br>3.9669e-01... | {'neg': 0.0,<br>'neu': 0.7692375),<br>0.929,<br>'pos': 0.16010574),<br>0.071,<br>'comp... | [(0,<br>0.7692375),<br>(3,<br>0.16010574),<br>(4, 0.054296... |

## Text Clustering using K-Means with TF-IDF Vectorization

```
In [40]: import pandas as pd
from sklearn.cluster import KMeans
from sklearn.feature_extraction.text import TfidfVectorizer

# Load your dataset
dataset = pd.read_excel("output_data_updated_with_textstat.xlsx")

# Vectorize the text using TF-IDF
tfidf_vectorizer = TfidfVectorizer()
tfidf_matrix = tfidf_vectorizer.fit_transform(dataset['Text']) # Use the '

# Apply K-Means clustering with explicit n_init
kmeans = KMeans(n_clusters=5, n_init=10) # You can adjust the number of cl
dataset['cluster'] = kmeans.fit_predict(tfidf_matrix)
```

## Text Summarization using Summa's Summarizer

```
In [44]: import pandas as pd
from summa import summarizer

# Generate summaries for each review
dataset['summary'] = dataset['Text'].apply(lambda x: summarizer.summarize(x
```

In [45]: dataset

Out[45]:

| Name   | Negative Reason Gold | Text  | ... | bert_embeddings                                   | sentiment  | topic   | entities   |
|--------|----------------------|---|-----|---|--|---|--|
| rum03  | NaN                  | southwestair<br>awesome<br>flight dallas<br>2 ny virgin...    | ... | tensor([[ 5.6657e-02, 2.0632e-02, 2.6626e-01...   | {'neg': 0.203, 'neu': 0.13657132), 0.617, 'pos': 0.6401591), 0.181, 'co... | [(0, 0.13657132), (1, 0.6401591), (2, 0.110866... | [('southwestair', 'ORG'), ('2', 'CARDI...        |
| ichime | NaN                  | usairways<br>thank finally<br>got bag<br>customer<br>servi... | ... | tensor([[ 1.4389e-01, -8.2253e-03, 3.2305e-01...  | {'neg': 0.0, 'neu': 0.53, 'pos': 0.47, 'compou...                          | [(0, 0.41092548), (1, 0.013585756), (2, 0.0135... | [] the   |
| jnts54 | NaN                  | americanair<br>dfwairport 2<br>together<br>best part fl...    | ... | tensor([[ -1.3857e-01, -1.6197e-01, 2.9042e-01... | {'neg': 0.0, 'neu': 0.588, 'pos': 0.412, 'comp...                          | [(0, 0.11536198), (1, 0.013556691), (2, 0.0135... | [] ar<br>dfw                                     |
| jordon | NaN                  | americanair<br>thank<br>youyou                                | ... | tensor([[ -9.3035e-02, 2.2016e-01, 3.7216e-01...  | {'neg': 0.0, 'neu': 0.444, 'pos': 0.556, 'comp...                          | [(0, 0.025300583), (1, 0.025211947), (2, 0.025... | [] ar  |
| amuel  | NaN                  | americanair<br>hopefully<br>see bad<br>ones<br>opportunity... | ... | tensor([[ -5.2652e-02, 1.5723e-01, 2.8512e-01...  | {'neg': 0.139, 'neu': 0.239, 'pos': 0.622, 'co...                          | [(0, 0.011991125), (1, 0.012075231), (2, 0.648... | [] ar  |
| ...    | ...                  | ...   | ... | ...   | ...  | ...   | ...  |
| enders | NaN                  | americanair<br>thank got<br>different<br>flight<br>chicago    | ... | tensor([[ 8.8884e-02, -2.7573e-01, 1.8394e-01...  | {'neg': 0.0, 'neu': 0.667, 'pos': 0.333, 'comp...                          | [(0, 0.44322193), (1, 0.020745821), (2, 0.0202... | [('chicago', 'GPE')]                             |
| sropes | NaN                  | americanair<br>leaving 20<br>minutes late<br>flight war...    | ... | tensor([[ -9.3935e-02, -1.3353e-01, 2.3228e-01... | {'neg': 0.279, 'neu': 0.721, 'pos': 0.0, 'comp...                          | [(0, 0.545215), (3, 0.42837554)]                  | [('20 minutes', 'TIME'), ('15 minutes', 'TIME')] |
| yabun  | NaN                  | americanair<br>please bring<br>american<br>airlines bla...    | ... | tensor([[ -1.0209e-01, -2.3004e-01, 2.5336e-01... | {'neg': 0.0, 'neu': 0.685, 'pos': 0.315, 'comp...                          | [(0, 0.23062189), (1, 0.018351018), (2, 0.0183... | [('american', 'NORP')]                           |
| ickson | NaN                  | americanair<br>money<br>change flight<br>dont answer<br>ph... | ... | tensor([[ 1.0850e-01, -2.3711e-01, 3.0256e-01...  | {'neg': 0.0, 'neu': 0.776, 'pos': 0.224, 'comp...                          | [(0, 0.013468066), (1, 0.30526412), (2, 0.3808... | [] ar  |

| Name   | Negative Reason Gold | Text   | ... | bert_embeddings                                   | sentiment  | topic   | entities  |
|--------|----------------------|--|-----|---|--|---|---|
| iddtwu | NaN                  | americanair<br>8 ppl need 2<br>know many<br>seats next ... | ... | tensor([[[-7.1985e-02, -2.6482e-01, 3.9669e-01... | {'neg': 0.0,<br>'neu': 0.929,<br>'pos': 0.071,<br>'comp... | [(0, 0.7692375),<br>(3, 0.16010574),<br>(4, 0.054296... | [('8', 'CARDINAL'),<br>(2, 'CARDINAL'),<br>(4, '... |

## Text Processing with NLTK and spaCy Libraries

```
In [48]: import pandas as pd
import nltk
import spacy
from nltk.corpus import wordnet
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
```

## Analyzing Text Dependencies using spaCy NLP Model

```
In [50]: # Load the spaCy model
nlp = spacy.load("en_core_web_sm")

def get_dependencies(text):
    doc = nlp(text)
    return [(token.text, token.dep_) for token in doc]

# Get the dependency relations for your text
dataset['dependencies'] = dataset['Text'].apply(get_dependencies)
```

## Tokenizing Text Data

```
In [7]: import pandas as pd
        from nltk.tokenize import word_tokenize

        # Load your dataset
        file_path = 'output_data_updated_with_textual_entailment.xlsx'
        data = pd.read_excel(file_path)

        # Tokenize the 'Text' column and create a new 'Tokenized Text' column
        data['Tokenized Text'] = data['Text'].apply(lambda text: word_tokenize(str(

        # Save the updated DataFrame to a new Excel file
        updated_file_path = 'output_data_updated_with_textual_entailment_with_token
        data.to_excel(updated_file_path, index=False)
```

## Training a Word2Vec Model on Tokenized Text Data

```
In [14]: from gensim.models import Word2Vec
        import ast # To convert the tokenized_text strings back to lists

        # Assuming you have a "tokenized_text" column containing tokenized text
        tokenized_text_data = dataset['tokenized_text'].apply(ast.literal_eval).tol

        # Train the Word2Vec model
        model = Word2Vec(sentences=tokenized_text_data, vector_size=100, window=5,
```

In [15]:

dataset

Out[15]:

| ID  | Airline<br>Sentiment | Sentiment<br>Confidence | Negative<br>Reason           | Negative<br>Reason<br>Confidence | Airline       | Airline<br>Sentiment<br>Gold | Name            | Nega<br>Rea<br>( |
|-----|----------------------|-------------------------|------------------------------|----------------------------------|---------------|------------------------------|-----------------|------------------|
| 224 | positive             | 0.9657                  | NaN                          | 0.0000                           | Southwest     | positive                     | magmum03        |                  |
| 272 | positive             | 1.0000                  | NaN                          | NaN                              | US<br>Airways | positive                     | christinachime  |                  |
| 388 | positive             | 1.0000                  | NaN                          | NaN                              | American      | positive                     | Runts54         |                  |
| 320 | positive             | 0.9236                  | NaN                          | 0.0000                           | American      | positive                     | douglaskgordon  |                  |
| 120 | positive             | 0.6242                  | NaN                          | 0.0000                           | American      | positive                     | byunsamuel      |                  |
| ... | ...                  | ...                     | ...                          | ...                              | ...           | ...                          | ...             | ...              |
| 280 | positive             | 0.3487                  | NaN                          | 0.0000                           | American      | NaN                          | KristenReenders |                  |
| 308 | negative             | 1.0000                  | Customer<br>Service<br>Issue | 1.0000                           | American      | NaN                          | itsropes        |                  |
| 272 | neutral              | 1.0000                  | NaN                          | NaN                              | American      | NaN                          | sanyabun        |                  |
| 432 | negative             | 1.0000                  | Customer<br>Service<br>Issue | 0.6659                           | American      | NaN                          | SraJackson      |                  |
| 388 | neutral              | 0.6771                  | NaN                          | 0.0000                           | American      | NaN                          | daviddtwu       |                  |

s



```
In [25]: from gensim.models import Word2Vec

# Loading the pre-trained Word2Vec model
model = Word2Vec.load("word_embed_training.model")
```

## Extracting Word Vectors for Visualization from Word2Vec Model

```
In [28]: # Extract word vectors for visualization
word_vectors = []
for text in dataset['Tokenized Text']:
    for word in text:
        if word in model.wv: # Check if the word is in the vocabulary
            word_vectors.append(model.wv[word])
```

## Converting Text Data in DataFrame to Tokenized Lists

```
In [32]: # Assuming your dataset is in a DataFrame named df
text_data = df[''].apply(lambda x: x.split()).tolist()
```

## Training a Word2Vec Model on Text Data

```
In [33]: from gensim.models import Word2Vec

# Train the Word2Vec model
model = Word2Vec(sentences=text_data, vector_size=100, window=5, min_count=
```

## Interactive Visualization of Word Embeddings using Matplotlib and mplcursors

```
In [40]: import matplotlib.pyplot as plt
import mplcursors

# Assuming you have already computed word_vectors_2d

# Create a larger figure
plt.figure(figsize=(16, 12))

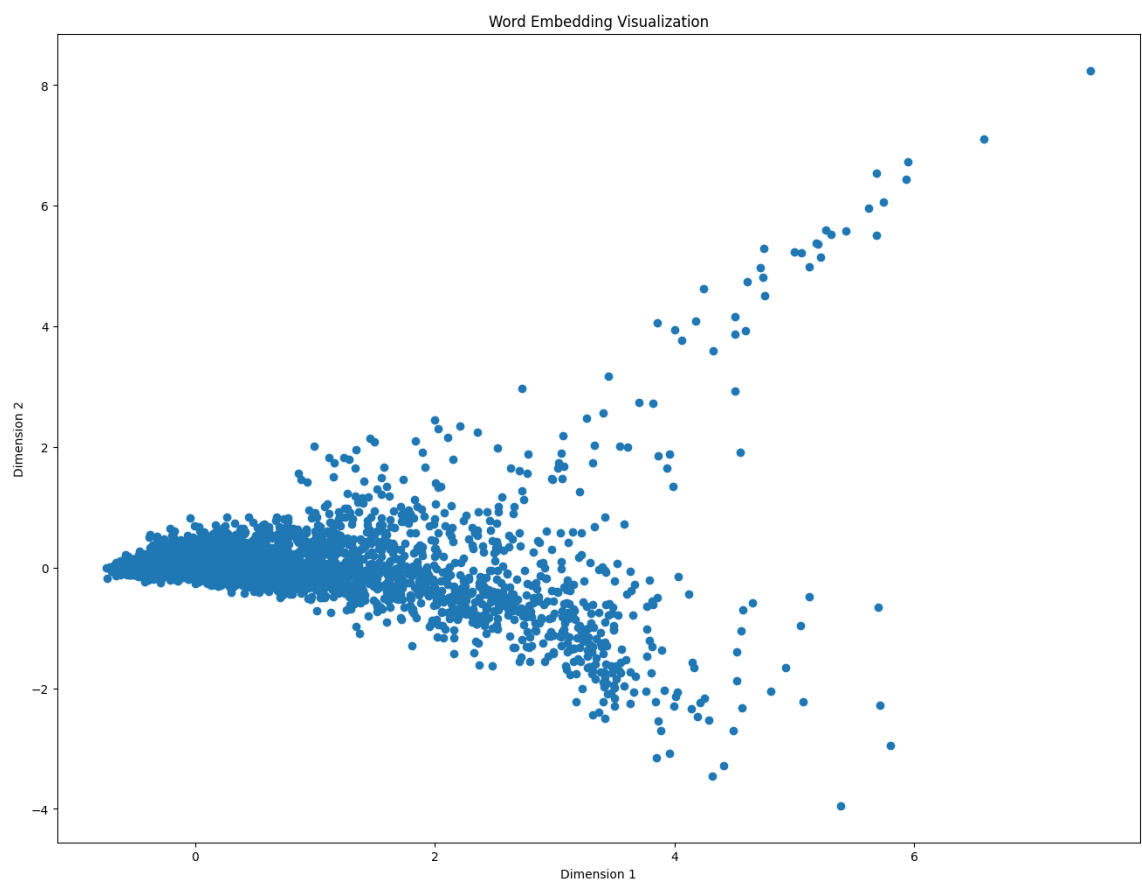
# Scatter plot for word embeddings
scatter = plt.scatter(word_vectors_2d[:, 0], word_vectors_2d[:, 1])

# Annotate words on the plot
for i, text in enumerate(dataset['tokenized_text']):
    for word in text:
        if word in model.wv:
            plt.annotate(word, xy=(word_vectors_2d[i, 0], word_vectors_2d[i, 1]))

plt.xlabel("Dimension 1")
plt.ylabel("Dimension 2")
plt.title("Word Embedding Visualization")

# Make the plot interactive with mplcursors
mplcursors.cursor(hover=True)

# Display the plot
plt.show()
```



## Pie Chart Visualization of Sentiment Distribution

```
In [52]: import matplotlib.pyplot as plt
import matplotlib.patches as mpatches

# Sample data
labels = ['Negative', 'Neutral', 'Positive']
sizes = [20, 50, 30]
colors = ['#ff9999', '#66b3ff', '#99ff99']
explode = (0.1, 0, 0) # Explode the 1st slice (i.e., 'Negative')

# Create a pie chart with 3D effect
fig, ax = plt.subplots()
ax.pie(sizes, explode=explode, labels=labels, colors=colors, autopct='%1.1f%%',
      shadow=True, startangle=140)

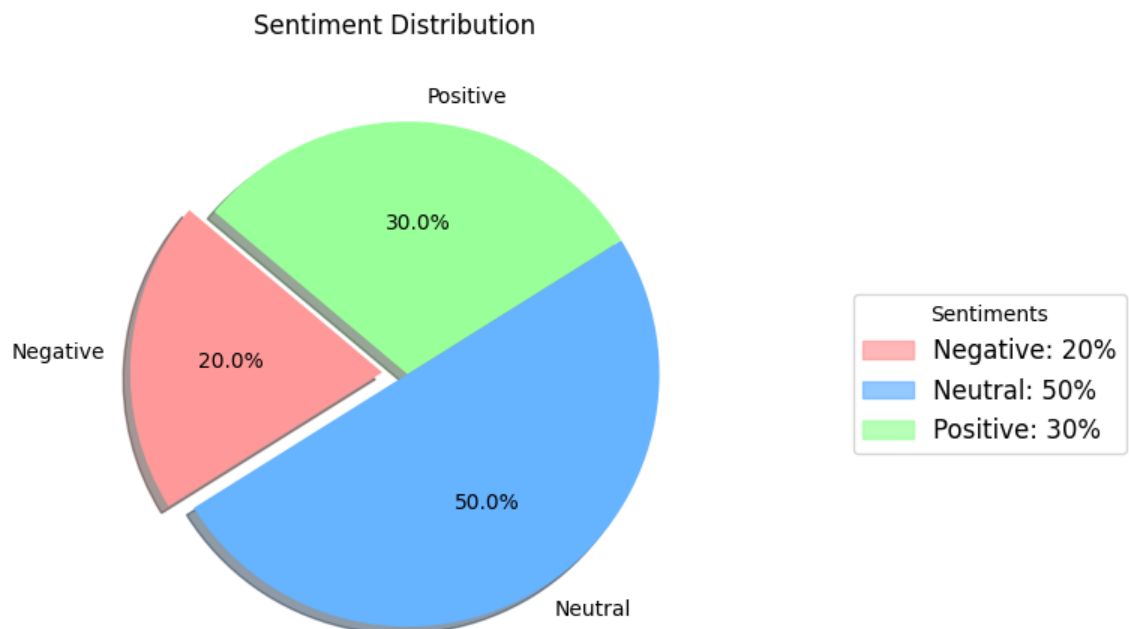
# Equal aspect ratio ensures that the pie is drawn as a circle
ax.axis('equal')

# Add a title
plt.title('Sentiment Distribution', pad=30) # Add padding to the title

# Create custom Legend handles and labels
legend_handles = [mpatches.Patch(color=color, label=f'{label}: {size}%', al

# Add a Legend on the right side with more spacing
plt.legend(handles=legend_handles, loc='center right', prop={'size': 12}, t

plt.show()
```



## Data Splitting using train\_test\_split from scikit-learn

```
In [54]: import pandas as pd
from sklearn.model_selection import train_test_split
```

```
In [70]: # Select input and output columns
X = data["Text"]
y = data["Airline Sentiment"]
```

```
In [71]: # Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, ra
```

## One-Hot Encoding with scikit-learn's OneHotEncoder

```
In [83]: from sklearn.preprocessing import OneHotEncoder

encoder = OneHotEncoder(sparse_output=False) # Use sparse_output instead of sparse
airline_encoded = encoder.fit_transform(dataset["Airline"].values.reshape(-1,
```

```
In [84]: airline_encoded
```

```
Out[84]: array([[0., 0., 1., 0., 0., 0.],
               [0., 0., 0., 1., 0., 0.],
               [1., 0., 0., 0., 0., 0.],
               ...,
               [1., 0., 0., 0., 0., 0.],
               [1., 0., 0., 0., 0., 0.],
               [1., 0., 0., 0., 0., 0.]])
```

```
In [85]: print("Original Data:")
print(dataset["Airline"].head()) # Assuming you're using a DataFrame

print("\nEncoded Data:")
print(airline_encoded)
```

Original Data:

```
0    Southwest
1    US Airways
2    American
3    American
4    American
```

Name: Airline, dtype: object

Encoded Data:

```
[[0. 0. 1. 0. 0. 0.]
 [0. 0. 0. 1. 0. 0.]
 [1. 0. 0. 0. 0. 0.]
 ...
 [1. 0. 0. 0. 0. 0.]
 [1. 0. 0. 0. 0. 0.]
 [1. 0. 0. 0. 0. 0.]]
```

In [86]: dataset

Out[86]:

| name  | Negative Reason Gold | Text  | ... | sentiment   | topic  | entities   | keywords   | emotion |
|-------|----------------------|---|-----|---|--|--|--|---------|
| im03  | NaN                  | southwestair<br>awesome<br>flight dallas<br>2 ny virgin...    | ... | {'neg': 0.203,<br>'neu': 0.617,<br>'pos': 0.181,<br>'co...} | [(0, 0.13657132),<br>(1, 0.6401591),<br>(2, 0.110866...] | [('southwestair',<br>'awesome',<br>'ORG'), ('2',<br>'CARDI...] | southwestair,<br>awesome,<br>flight, dallas,<br>2          | joy     |
| hime  | NaN                  | usairways<br>thank finally<br>got bag<br>customer<br>servi... | ... | {'neg': 0.0,<br>'neu': 0.53,<br>'pos': 0.47,<br>'compou...} | [(0, 0.41092548),<br>(1, 0.013585756),<br>(2, 0.0135...] | []   | usairways,<br>thank, finally,<br>got, bag                  | joy     |
| its54 | NaN                  | americanair<br>dfwairport 2<br>together<br>best part fl...    | ... | {'neg': 0.0,<br>'neu': 0.588,<br>'pos': 0.412,<br>'comp...} | [(0, 0.11536198),<br>(1, 0.013556691),<br>(2, 0.0135...] | []   | americanair,<br>dfwairport, 2,<br>together,<br>best        | joy     |
| rdon  | NaN                  | americanair<br>thank<br>youyou                                | ... | {'neg': 0.0,<br>'neu': 0.444,<br>'pos': 0.556,<br>'comp...} | [(0, 0.025300583),<br>(1, 0.025211947),<br>(2, 0.025...] | []   | americanair,<br>thank,<br>youyou                           | neutral |
| muel  | NaN                  | americanair<br>hopefully<br>see bad<br>ones<br>opportunity... | ... | {'neg': 0.139,<br>'neu': 0.239,<br>'pos': 0.622,<br>'co...} | [(0, 0.011991125),<br>(1, 0.012075231),<br>(2, 0.648...] | []   | ones,<br>americanair,<br>hopefully,<br>see, bad            | neutral |
| ...   | ...                  | ...   | ... | ...   | ...  | ...  | ...  | ...     |
| iders | NaN                  | americanair<br>thank got<br>different<br>flight<br>chicago    | ... | {'neg': 0.0,<br>'neu': 0.667,<br>'pos': 0.333,<br>'comp...} | [(0, 0.44322193),<br>(1, 0.020745821),<br>(2, 0.0202...] | [('chicago',<br>'GPE')]  | americanair,<br>thank, got,<br>different,<br>flight        | neutral |
| opes  | NaN                  | americanair<br>leaving 20<br>minutes late<br>flight war...    | ... | {'neg': 0.279,<br>'neu': 0.721,<br>'pos': 0.0,<br>'comp...} | [(0, 0.545215), (3, 0.42837554)]                         | [('20 minutes',<br>'TIME'), ('15<br>minutes',<br>'TIME')]      | minutes,<br>late, flight,<br>americanair,<br>leaving       | sadness |
| abun  | NaN                  | americanair<br>please bring<br>american<br>airlines bla...    | ... | {'neg': 0.0,<br>'neu': 0.685,<br>'pos': 0.315,<br>'comp...} | [(0, 0.23062189),<br>(1, 0.018351018),<br>(2, 0.0183...] | [('american',<br>'NORP')]                                      | americanair,<br>please,<br>bring,<br>american,<br>airlines | neutral |
| kson  | NaN                  | americanair<br>money<br>change flight<br>dont answer<br>ph... | ... | {'neg': 0.0,<br>'neu': 0.776,<br>'pos': 0.224,<br>'comp...} | [(0, 0.013468066),<br>(1, 0.30526412),<br>(2, 0.3808...] | []   | americanair,<br>money,<br>change,<br>flight, dont          | neutral |

| name | Negative Reason Gold | Text   | ... | sentiment  | topic   | entities  | keywords                                | emotion |
|------|----------------------|--|-----|--|---|---|---|---------|
| dtwu | NaN                  | americanair<br>8 ppl need 2<br>know many<br>seats next ... | ... | {'neg': 0.0,<br>'neu': 0.929,<br>'pos': 0.071,<br>'comp... | [(0,<br>0.7692375),<br>(3,<br>0.16010574),<br>(4, 0.054296... | [('8',<br>'CARDINAL'),<br>(2',<br>'CARDINAL'),<br>(4', '... | next, flight,<br>americanair,<br>8, ppl | neutral |

## Standard Scaling using scikit-learn's StandardScaler

```
In [87]: from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
sentiment_confidence_scaled = scaler.fit_transform(dataset["Sentiment Confidence"])
```

```
In [90]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(dataset.drop("Airline Sentiment", axis=1), dataset["Airline Sentiment"],
                                                    test_size=0.2, random_state=42)
```

## Filling Missing Values in a DataFrame using 'Unknown'

```
In [95]: dataset.fillna("Unknown", inplace=True)
```



In [96]: dataset

Out[96]:

| name  | Negative Reason Gold | Text  | ... | sentiment   | topic  | entities   | keywords   | emotion |
|-------|----------------------|---|-----|---|--|--|--|---------|
| im03  | Unknown              | southwestair<br>awesome<br>flight dallas<br>2 ny virgin...    | ... | {'neg': 0.203,<br>'neu': 0.617,<br>'pos': 0.181,<br>'co...} | [(0, 0.13657132),<br>(1, 0.6401591),<br>(2, 0.110866...] | [('southwestair',<br>'awesome',<br>'ORG'), ('2',<br>'CARDI...] | southwestair,<br>awesome,<br>flight, dallas,<br>2          | joy     |
| hime  | Unknown              | usairways<br>thank finally<br>got bag<br>customer<br>servi... | ... | {'neg': 0.0,<br>'neu': 0.53,<br>'pos': 0.47,<br>'compou...} | [(0, 0.41092548),<br>(1, 0.013585756),<br>(2, 0.0135...] | []   | usairways,<br>thank, finally,<br>got, bag                  | joy     |
| its54 | Unknown              | americanair<br>dfwairport 2<br>together<br>best part fl...    | ... | {'neg': 0.0,<br>'neu': 0.588,<br>'pos': 0.412,<br>'comp...} | [(0, 0.11536198),<br>(1, 0.013556691),<br>(2, 0.0135...] | []   | americanair,<br>dfwairport, 2,<br>together,<br>best        | joy     |
| rdon  | Unknown              | americanair<br>thank<br>youyou                                | ... | {'neg': 0.0,<br>'neu': 0.444,<br>'pos': 0.556,<br>'comp...} | [(0, 0.025300583),<br>(1, 0.025211947),<br>(2, 0.025...] | []   | americanair,<br>thank,<br>youyou                           | neutral |
| muel  | Unknown              | americanair<br>hopefully<br>see bad<br>ones<br>opportunity... | ... | {'neg': 0.139,<br>'neu': 0.239,<br>'pos': 0.622,<br>'co...} | [(0, 0.011991125),<br>(1, 0.012075231),<br>(2, 0.648...] | []   | ones,<br>americanair,<br>hopefully,<br>see, bad            | neutral |
| ...   | ...                  | ...   | ... | ...   | ...  | ...  | ...  | ...     |
| iders | Unknown              | americanair<br>thank got<br>different<br>flight<br>chicago    | ... | {'neg': 0.0,<br>'neu': 0.667,<br>'pos': 0.333,<br>'comp...} | [(0, 0.44322193),<br>(1, 0.020745821),<br>(2, 0.0202...] | [('chicago',<br>'GPE')]  | americanair,<br>thank, got,<br>different,<br>flight        | neutral |
| opes  | Unknown              | americanair<br>leaving 20<br>minutes late<br>flight war...    | ... | {'neg': 0.279,<br>'neu': 0.721,<br>'pos': 0.0,<br>'comp...} | [(0, 0.545215), (3, 0.42837554)]                         | [('20 minutes',<br>'TIME'), ('15<br>minutes',<br>'TIME')]      | minutes,<br>late, flight,<br>americanair,<br>leaving       | sadness |
| abun  | Unknown              | americanair<br>please bring<br>american<br>airlines bla...    | ... | {'neg': 0.0,<br>'neu': 0.685,<br>'pos': 0.315,<br>'comp...} | [(0, 0.23062189),<br>(1, 0.018351018),<br>(2, 0.0183...] | [('american',<br>'NORP')]                                      | americanair,<br>please,<br>bring,<br>american,<br>airlines | neutral |
| kson  | Unknown              | americanair<br>money<br>change flight<br>dont answer<br>ph... | ... | {'neg': 0.0,<br>'neu': 0.776,<br>'pos': 0.224,<br>'comp...} | [(0, 0.013468066),<br>(1, 0.30526412),<br>(2, 0.3808...] | []   | americanair,<br>money,<br>change,<br>flight, dont          | neutral |

| name | Negative Reason Gold | Text   | ... | sentiment  | topic   | entities  | keywords                                | emotion |
|------|----------------------|--|-----|--|---|---|---|---------|
| dtwu | Unknown              | americanair<br>8 ppl need 2<br>know many<br>seats next ... | ... | {'neg': 0.0,<br>'neu': 0.929,<br>'pos': 0.071,<br>'comp... | [(0,<br>0.7692375),<br>(3,<br>0.16010574),<br>(4, 0.054296... | [('8',<br>'CARDINAL'),<br>(<br>'2',<br>'CARDINAL'),<br>(<br>'4', '... | next, flight,<br>americanair,<br>8, ppl | neutral |

In [98]: dataset.isnull().sum()

Out[98]:

|                            |   |
|----------------------------|---|
| Tweet ID                   | 0 |
| Airline Sentiment          | 0 |
| Sentiment Confidence       | 0 |
| Negative Reason            | 0 |
| Negative Reason Confidence | 0 |
| Airline                    | 0 |
| Airline Sentiment Gold     | 0 |
| Name                       | 0 |
| Negative Reason Gold       | 0 |
| Text                       | 0 |
| Tweet Created              | 0 |
| Tweet Location             | 0 |
| User Timezone              | 0 |
| tokenized_text             | 0 |
| bert_embeddings            | 0 |
| sentiment                  | 0 |
| topic                      | 0 |
| entities                   | 0 |
| keywords                   | 0 |
| emotion                    | 0 |
| flesch_kincaid             | 0 |
| gunning_fog                | 0 |
| cluster                    | 0 |
| dependencies               | 0 |
| Tokenized Text             | 0 |
| dtype: int64               |   |

# Configuring NLTK Data Directory and POS Tagging with NLTK

```
In [115]: import nltk
          from nltk import word_tokenize, pos_tag

          # Set the NLTK data directory to the location of your downloaded WordNet re
          nltk.data.path.append("D:\\nltk_data")

          # No need to download WordNet or averaged_perceptron_tagger again; they are
          tokens = word_tokenize(text)
          pos_tags = pos_tag(tokens)
```

## Using spaCy NER Model and Configuring Default Settings in PyTorch

```
In [117]: import spacy

          nlp = spacy.load("en_core_web_sm")
          doc = nlp(text)
          entities = [(ent.text, ent.label_) for ent in doc.ents]

          import torch

          # Use torch.set_default_dtype() and torch.set_default_device() instead
          torch.set_default_dtype(torch.float32)
          torch.set_default_device('cpu')
```

```
In [120]: from nltk.corpus import stopwords

          custom_stopwords = ["custom_stopword1", "custom_stopword2"]
          filtered_text = [word for word in text if word not in custom_stopwords]
```

## Replacing Email Addresses with 'EMAIL' using Regular Expressions

```
In [121]: import re

          text = re.sub(r'\S+@\S+', 'EMAIL', text)
```

## Text Classification with scikit-learn: Data Splitting, Feature Extraction, and Model Evaluation

```
In [138]: import pandas as pd
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
from sklearn.svm import SVC

# Load your dataset into a DataFrame
data = pd.read_excel('Data_Cleaning.xlsx')

# Assuming 'text' is the column containing the text data
text_data = data['Text']

# Assuming 'Airline Sentiment' is the target variable
target = data['Airline Sentiment']

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(text_data, target, test

# Create a CountVectorizer for text feature extraction
vectorizer = CountVectorizer()
X_train_vec = vectorizer.fit_transform(X_train)
X_test_vec = vectorizer.transform(X_test)

# Train a classifier (e.g., SVM)
classifier = SVC()
classifier.fit(X_train_vec, y_train)

# Make predictions on the test set
y_pred = classifier.predict(X_test_vec)

# Evaluate the model
print(classification_report(y_test, y_pred))
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| negative     | 0.82      | 0.92   | 0.86     | 1815    |
| neutral      | 0.64      | 0.51   | 0.57     | 641     |
| positive     | 0.77      | 0.61   | 0.68     | 465     |
| accuracy     |           |        | 0.78     | 2921    |
| macro avg    | 0.74      | 0.68   | 0.70     | 2921    |
| weighted avg | 0.77      | 0.78   | 0.77     | 2921    |

## Text Data Vectorization with TF-IDF using scikit-learn's TfidfVectorizer

```
In [125]: from sklearn.feature_extraction.text import TfidfVectorizer

# Assuming 'data' is your DataFrame with the dataset
documents = data['Text'].tolist() # Extract the 'Text' column as a list of

# Initialize the TF-IDF vectorizer
tfidf_vectorizer = TfidfVectorizer(sublinear_tf=True, token_pattern=r'\b\w+'

# Fit and transform the vectorizer on the list of documents
X_tfidf = tfidf_vectorizer.fit_transform(documents)

# X_tfidf now contains the TF-IDF representation of the 'Text' column
```

## Text Cleaning with Typo Correction and Slang Word Replacement

```
In [129]: import pandas as pd
from autocorrect import Speller

# Load your dataset
data = pd.read_excel("Data_Cleaning.xlsx")

# Initialize a spell checker for typo correction
spell = Speller(lang='en')

# Define a dictionary for slang word replacement
slang_replacements = {
    'u': 'you',
    'r': 'are',
    '2': 'to',
    'gr8': 'great',
    'lol': 'laugh out loud',
    'omg': 'oh my god'
    # Add more slang words and their replacements as needed
}

# Function to correct typos and replace slang
def clean_text(text):
    # Correct common typos
    corrected_text = spell(text)

    # Replace slang words
    for slang, replacement in slang_replacements.items():
        corrected_text = corrected_text.replace(slang, replacement)

    return corrected_text

# Apply the clean_text function to the 'Text' column
data['Cleaned_Text'] = data['Text'].apply(clean_text)

# The 'Cleaned_Text' column now contains the cleaned and corrected text
```

In [18]:

df

Out[18]:

|       | Tweet ID           | Airline<br>Sentiment | Sentiment<br>Confidence | Negative<br>Reason           | Negative<br>Reason<br>Confidence | Airline       | Air<br>Sentim<br>G |
|-------|--------------------|----------------------|-------------------------|------------------------------|----------------------------------|---------------|--------------------|
| 0     | 570066283233972224 | positive             | 0.9657                  | Unknown                      | 0                                | Southwest     | pos                |
| 1     | 568561924985782272 | positive             | 1.0000                  | Unknown                      | Unknown                          | US<br>Airways | pos                |
| 2     | 570264948548313088 | positive             | 1.0000                  | Unknown                      | Unknown                          | American      | pos                |
| 3     | 569600137296633920 | positive             | 0.9236                  | Unknown                      | 0                                | American      | pos                |
| 4     | 568551906634797120 | positive             | 0.6242                  | Unknown                      | 0                                | American      | pos                |
| ...   | ...                | ...                  | ...                     | ...                          | ...                              | ...           | ...                |
| 14599 | 569587686496825280 | positive             | 0.3487                  | Unknown                      | 0                                | American      | Unkn               |
| 14600 | 569587371693355008 | negative             | 1.0000                  | Customer<br>Service<br>Issue | 1                                | American      | Unkn               |
| 14601 | 569587242672398272 | neutral              | 1.0000                  | Unknown                      | Unknown                          | American      | Unkn               |
| 14602 | 569587188687634432 | negative             | 1.0000                  | Customer<br>Service<br>Issue | 0.6659                           | American      | Unkn               |
| 14603 | 569587140490866688 | neutral              | 0.6771                  | Unknown                      | 0                                | American      | Unkn               |

14604 rows × 26 columns

## Text Preprocessing with a Custom 'preprocess\_text' Function

```
In [19]: def preprocess_text(text):  
        # Implement text preprocessing here  
        return text  
  
df["Text"] = df["Text"].apply(preprocess_text)
```

```
In [20]: from nltk.tokenize import word_tokenize  
  
df["tokenized_text_updated"] = df["Text"].apply(word_tokenize)
```

## Mapping Sentiment Labels to Numeric Values for a Feature

```
In [21]: # Feature 1: Sentiment of Negative Reason  
sentiment_map = {"negative": -1, "neutral": 0, "positive": 1}  
df["negative_reason_sentiment"] = df["Negative Reason"].map(sentiment_map)
```

## Extracting Hour of Tweet Creation as a Feature

```
In [28]: # Feature 2: Hour of Tweet Creation  
df["tweet_hour"] = df["Tweet Created"].dt.hour
```

## Extracting Day of the Week from Tweet Creation Date as a Feature

```
In [30]: # Feature 3: Day of the Week  
df["tweet_day_of_week"] = df["Tweet Created"].dt.dayofweek
```

## Encoding User Timezone as a Numeric Feature

```
In [32]: # Feature 4: User Timezone  
df["user_timezone_encoded"] = df["User Timezone"].factorize()[0]
```



In [33]: df

Out[33]:

| flesch_kincaid | gunning_fog | cluster | dependencies   | Tokenized Text                                       | text_length | tokenized_text                 |
|----------------|-------------|---------|--|--|-------------|--------------------------------|
| 11.1           | 11.51       | 4       | [('southwestair', 'amod'), ('awesome', 'amod')...]   | ['southwestair', 'awesome', 'flight', 'dallas'...]   | 210         | ['sou', 'awesor']              |
| 10.3           | 8.04        | 2       | [('usairways', 'nsubj'), ('thank', 'nsubj'), (...)]  | ['usairways', 'thank', 'finally', 'got', 'bag'...]   | 127         | ['usairway', 'finally', 'g']   |
| 9.6            | 8.51        | 1       | [('americanair', 'compound'), ('dfwairport', '...)]  | ['americanair', 'dfwairport', '2', 'together',...]   | 129         | ['am', 'dfw:', 't']            |
| 9.2            | 14.53       | 1       | [('americanair', 'nsubj'), ('thank', 'nsubj'), ('ROOT'), ...]  | ['americanair', 'thank', 'youyou']                   | 64          | ['americanair']                |
| 11.9           | 17.51       | 1       | [('americanair', 'nsubj'), ('hopefully', 'advmod'), ('see', 'advmod'), ('bad', 'advmod'), ('on', 'advmod'), ...] | ['americanair', 'hopefully', 'see', 'bad', 'on'...]  | 147         | ['americanair', 'see', '']     |
| ...            | ...         | ...     | ...  | ...  | ...         | ...                            |
| 8.0            | 9.07        | 1       | [('americanair', 'compound'), ('thank', 'nsubj'), ('different', 'nsubj'), ...]                                   | ['americanair', 'thank', 'got', 'different', '...]   | 81          | ['americanair', 'got', 'di']   |
| 11.1           | 11.51       | 1       | [('americanair', 'nsubj'), ('leaving', 'acl'), ...]  | ['americanair', 'leaving', '20', 'minutes', 'l'...]  | 201         | ['americanair', '20', 'm']     |
| 12.7           | 15.73       | 1       | [('americanair', 'nsubj'), ('please', 'nsubj'), ('bring', 'nsubj'), ('intj'), ...]                               | ['americanair', 'please', 'bring', 'american', '...] | 92          | ['americanair', 'bring', 'ar'] |
| 9.6            | 16.00       | 1       | [('americanair', 'compound'), ('money', 'compound'), ('change', 'compound'), ...]                                | ['americanair', 'money', 'change', 'flight', 'l'...] | 128         | ['americanair', 'change']      |
| 5.6            | 9.42        | 1       | [('americanair', 'nsubj'), ('8', 'nummod'), ('know', 'nummod'), ...]   | ['americanair', '8', 'ppl', 'need', '2', 'know'...]  | 217         | ['americanair', 'need', '']    |

## Mapping Sentiment Labels to Numeric Values using a Custom Function

```
In [34]: # Assuming your sentiment labels are 'positive', 'negative', and 'neutral'
def label_sentiment(sentiment):
    if sentiment == 'positive':
        return 1
    elif sentiment == 'negative':
        return -1
    else:
        return 0 # Neutral

df["sentiment_label"] = df["Airline Sentiment"].apply(label_sentiment)
```

## Calculating and Visualizing a Confusion Matrix

```
In [89]: from sklearn.metrics import confusion_matrix

# Compute the confusion matrix
confusion = confusion_matrix(y_test, y_pred)

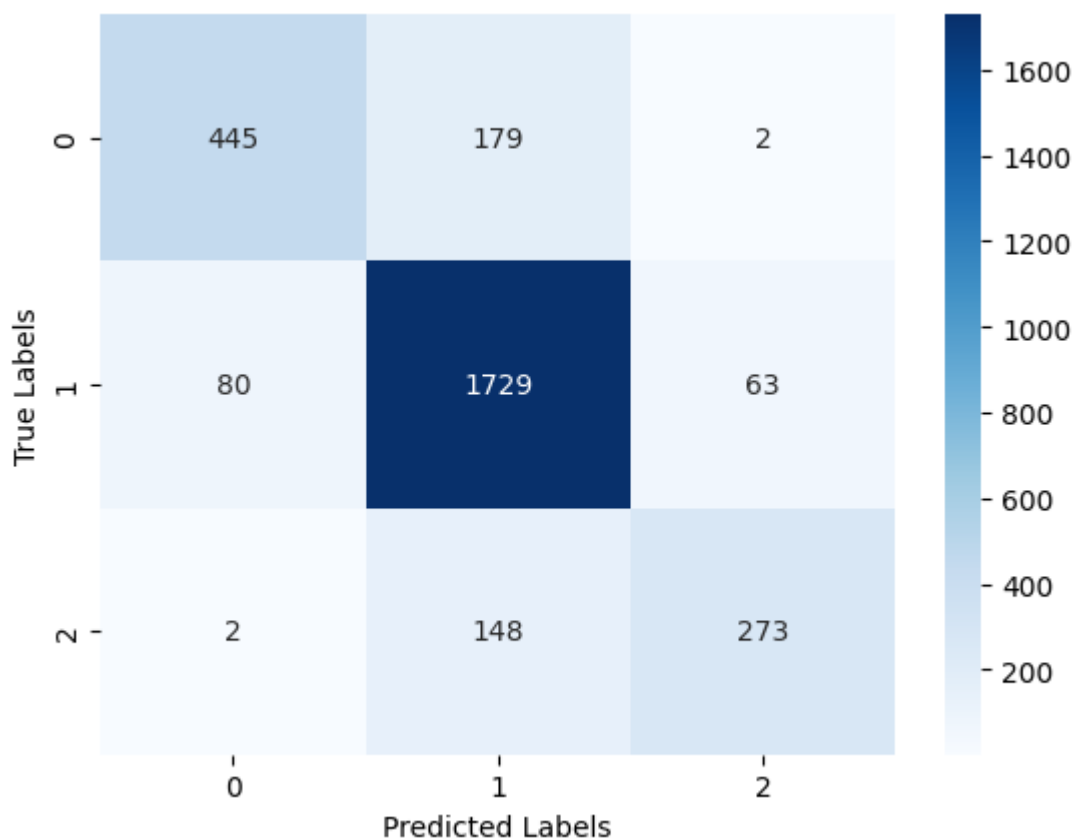
# Print the confusion matrix
print("Confusion Matrix:")
print(confusion)

# Visualize the confusion matrix (optional)
import seaborn as sns
import matplotlib.pyplot as plt

sns.heatmap(confusion, annot=True, fmt="d", cmap="Blues")
plt.xlabel("Predicted Labels")
plt.ylabel("True Labels")
plt.show()
```

Confusion Matrix:

```
[[ 445  179    2]
 [   80 1729   63]
 [    2  148  273]]
```



## Training a Logistic Regression Model with Increased 'max\_iter' and Making Predictions

```
In [91]: from sklearn.linear_model import LogisticRegression

# Initialize the model with increased max_iter
model = LogisticRegression(C=0.01, max_iter=10000)

# Fit the model to the scaled training data
model.fit(X_train_scaled, y_train)

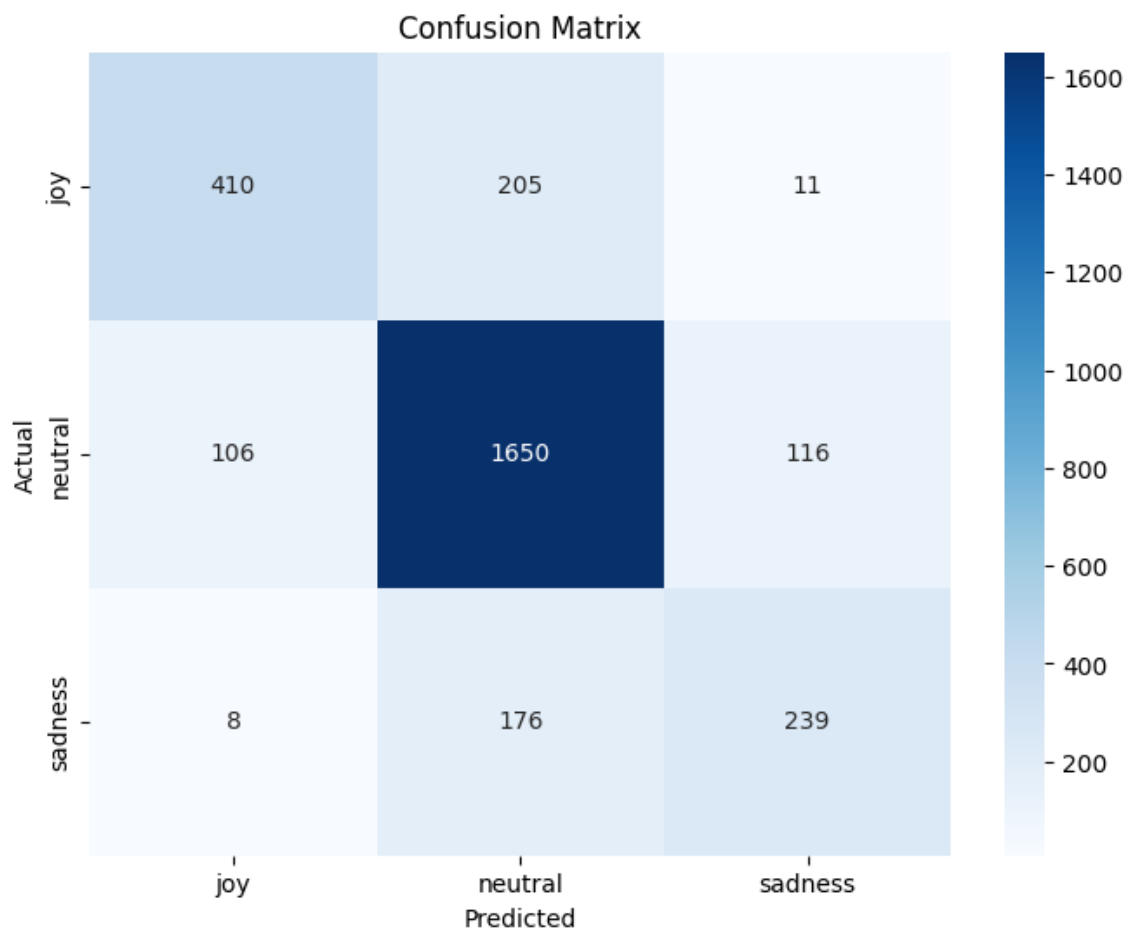
# Predict sentiment on the scaled test data
y_pred = model.predict(X_test_scaled)
```

## Creating and Visualizing a Confusion Matrix for Multiclass Sentiment Classification

```
In [93]: from sklearn.metrics import confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

# Generate the confusion matrix
confusion = confusion_matrix(y_test, y_pred, labels=["joy", "neutral", "sadness"])

# Create a heatmap of the confusion matrix
plt.figure(figsize=(8, 6))
sns.heatmap(confusion, annot=True, fmt="d", cmap="Blues", xticklabels=["joy", "neutral", "sadness"], yticklabels=["Actual joy", "Actual neutral", "Actual sadness"], title="Confusion Matrix")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix")
plt.show()
```



```
In [1]: import pandas as pd

# Load your dataset
df = pd.read_excel('Extracted_dataset.xlsx')
```

## Calculating Sentiment Scores with TextBlob and Adding to the DataFrame

```
In [2]: from textblob import TextBlob

# Define a function to calculate sentiment
def calculate_sentiment(text):
    analysis = TextBlob(text)
    return analysis.sentiment.polarity

# Apply the sentiment function to your text data
df['Sentiment'] = df['Text'].apply(calculate_sentiment)
```

In [3]:

df

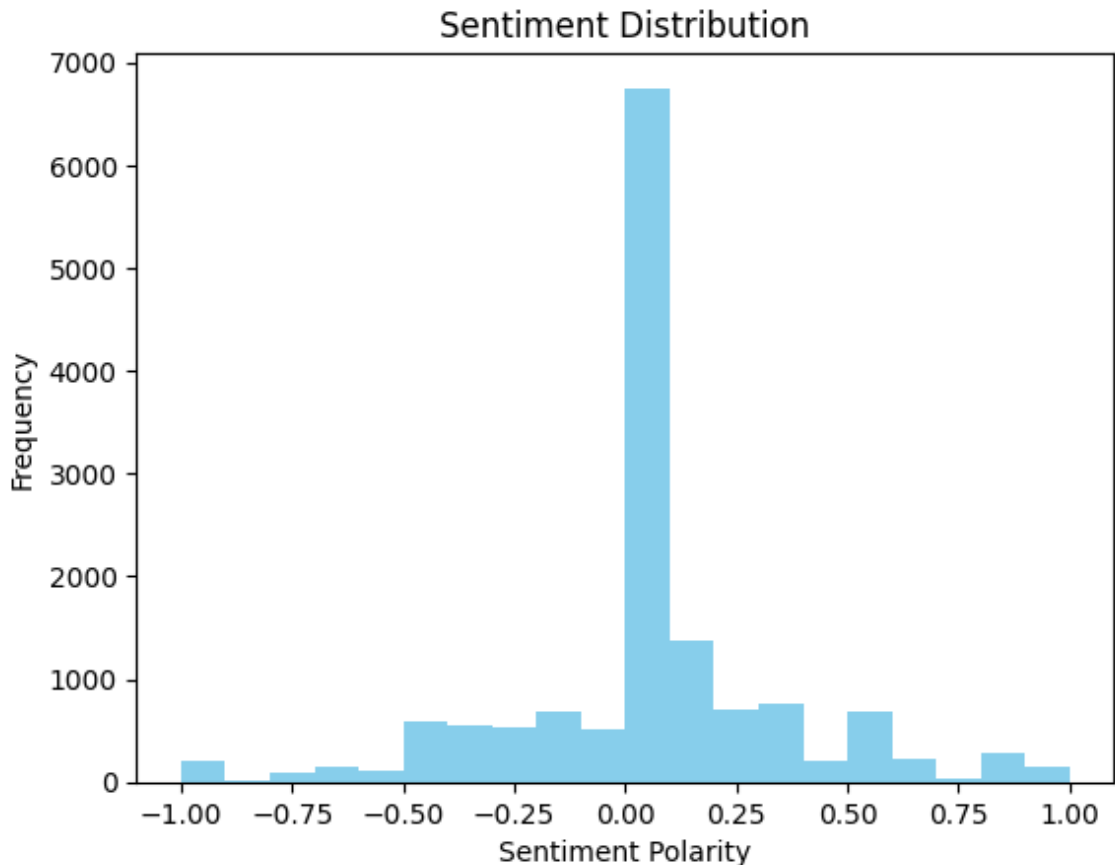
Out[3]:

| ... | cluster | dependencies  | Tokenized<br>Text   | text_length | tweet_hour | tweet_day_of_week | user_ti |
|-----|---------|---|---|-------------|------------|-------------------|---------|
| ... | 4       | [('southwestair',<br>'amod'),<br>(<br>'awesome',<br>'amod')...  | [('southwestair',<br>'awesome',<br>'flight',<br>'dallas'... | 210         | 19         | 0                 |         |
| ... | 2       | [('usairways',<br>'nsubj'),<br>(<br>'thank',<br>'nsubj'), (...  | [('usairways',<br>'thank',<br>'finally', 'got',<br>'bag'... | 127         | 16         | 3                 |         |
| ... | 1       | [('americanair',<br>'compound'),<br>(<br>'dfwairport', '...     | [('americanair',<br>'dfwairport', '2',<br>'together',...    | 129         | 8          | 1                 |         |
| ... | 1       | [('americanair',<br>'nsubj'),<br>(<br>'thank',<br>'ROOT'), ...  | [('americanair',<br>'thank',<br>'youyou']                   | 64          | 12         | 6                 |         |
| ... | 1       | [('americanair',<br>'nsubj'),<br>(<br>'hopefully',<br>'adv...   | [('americanair',<br>'hopefully',<br>'see', 'bad',<br>'on... | 147         | 15         | 3                 |         |
| ... | ...     | ...   | ...   | ...         | ...        | ...               | ...     |
| ... | 1       | [('americanair',<br>'compound'),<br>(<br>'thank',<br>'nsubj'... | [('americanair',<br>'thank', 'got',<br>'different', '...    | 81          | 12         | 6                 |         |
| ... | 1       | [('americanair',<br>'nsubj'),<br>(<br>'leaving',<br>'acl'),...  | [('americanair',<br>'leaving', '20',<br>'minutes', 'l...    | 201         | 11         | 6                 |         |
| ... | 1       | [('americanair',<br>'nsubj'),<br>(<br>'please',<br>'intj'),...  | [('americanair',<br>'please',<br>'bring',<br>'american',... | 92          | 11         | 6                 |         |
| ... | 1       | [('americanair',<br>'compound'),<br>(<br>'money',<br>'compo...  | [('americanair',<br>'money',<br>'change',<br>'flight', '... | 128         | 11         | 6                 |         |
| ... | 1       | [('americanair',<br>'nsubj'), ('8',<br>'nummod'), ('...         | [('americanair',<br>'8', 'ppl',<br>'need', '2',<br>'know... | 217         | 11         | 6                 |         |

# Plotting Sentiment Distribution in a Histogram

```
In [4]: import matplotlib.pyplot as plt

# Plot sentiment distribution
plt.hist(df['Sentiment'], bins=20, color='skyblue')
plt.title('Sentiment Distribution')
plt.xlabel('Sentiment Polarity')
plt.ylabel('Frequency')
plt.show()
```



## Extracting Keywords from Text Data using CountVectorizer

```
In [12]: from sklearn.feature_extraction.text import CountVectorizer

def extract_keywords(text):
    vectorizer = CountVectorizer(ngram_range=(1, 2))
    X = vectorizer.fit_transform([text])
    feature_names = vectorizer.get_feature_names_out()
    return feature_names

df['Keywords'] = df['Text'].apply(extract_keywords)
```



In [13]: df

Out[13]:

| dependencies   | Tokenized Text                                     | text_length | tweet_hour | tweet_day_of_week | user_timezone_encod |
|--|--|-------------|------------|-------------------|---------------------|
| [('southwestair', 'amod'), ('awesome', 'amod')...]     | ['southwestair', 'awesome', 'flight', 'dallas'...] | 210         | 19         | 0                 |                     |
| [('usairways', 'nsubj'), ('thank', 'nsubj'), (...)]    | ['usairways', 'thank', 'finally', 'got', 'bag'...] | 127         | 16         | 3                 |                     |
| [('americanair', 'compound'), ('dfwairport', '...)]    | ['americanair', 'dfwairport', '2', 'together',...] | 129         | 8          | 1                 |                     |
| [('americanair', 'nsubj'), ('thank', 'ROOT'), ...]     | ['americanair', 'thank', 'youyou']                 | 64          | 12         | 6                 |                     |
| [('americanair', 'nsubj'), ('hopefully', 'advmod')...] | ['americanair', 'hopefully', 'see', 'bad', 'on...] | 147         | 15         | 3                 |                     |
| ...  | ...  | ...         | ...        | ...               |                     |
| [('americanair', 'compound'), ('thank', 'nsubj')...]   | ['americanair', 'thank', 'got', 'different', '...] | 81          | 12         | 6                 |                     |
| [('americanair', 'nsubj'), ('leaving', 'acl'),...]     | ['americanair', 'leaving', '20', 'minutes', 'l...] | 201         | 11         | 6                 |                     |
| [('americanair', 'nsubj'), ('please', 'intj'),...]     | ['americanair', 'please', 'bring', 'american',...] | 92          | 11         | 6                 |                     |
| [('americanair', 'compound'), ('money', 'compo...)]    | ['americanair', 'money', 'change', 'flight', '...] | 128         | 11         | 6                 |                     |
| [('americanair', 'nsubj'), ('8', 'nummod'), ('...)]    | ['americanair', '8', 'ppl', 'need', '2', 'know...] | 217         | 11         | 6                 |                     |



# Sentiment Analysis, Keyword Extraction, and Sentiment Distribution

```
In [18]: import pandas as pd
import numpy as np
from textblob import TextBlob
from collections import Counter

# Load your dataset
df = pd.read_excel("Extracted_dataset.xlsx")

# Perform sentiment analysis
df['Sentiment_Polarity'] = df['Text'].apply(lambda x: TextBlob(str(x)).sentiment.polarity)

# Extract keywords
df['Keywords'] = df['Text'].str.split() # Split text into words

# Analyze sentiment distribution
positive_tweets = df[df['Sentiment_Polarity'] > 0]
negative_tweets = df[df['Sentiment_Polarity'] < 0]

# Calculate keyword frequencies for positive and negative tweets
positive_keywords = ' '.join(positive_tweets['Keywords'].sum()).split()
negative_keywords = ' '.join(negative_tweets['Keywords'].sum()).split()

positive_keyword_counts = Counter(positive_keywords)
negative_keyword_counts = Counter(negative_keywords)

# Get the top 10 positive and negative keywords
top_positive_keywords = positive_keyword_counts.most_common(10)
top_negative_keywords = negative_keyword_counts.most_common(10)

# Print sentiment distribution and top keywords
print("Sentiment Distribution:")
print("Positive Sentiment:", len(positive_tweets))
print("Negative Sentiment:", len(negative_tweets))

print("\nTop Keywords for Positive Sentiment:")
print(top_positive_keywords)

print("\nTop Keywords for Negative Sentiment:")
print(top_negative_keywords)
```

Sentiment Distribution:  
 Positive Sentiment: 5313  
 Negative Sentiment: 3419

Top Keywords for Positive Sentiment:  
 [('united', 1444), ('flight', 1256), ('americanair', 1027), ('jetblue', 1002), ('southwestair', 988), ('thanks', 944), ('usairways', 931), ('get', 431), ('service', 341), ('great', 325)]

Top Keywords for Negative Sentiment:  
 [('united', 1145), ('flight', 1105), ('usairways', 915), ('americanair', 693), ('southwestair', 462), ('jetblue', 396), ('service', 395), ('late', 342), ('get', 325), ('customer', 300)]

# Importing Necessary Libraries for Text Data Analysis and Machine Learning

```
In [1]: # Import necessary libraries
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report
import joblib # Import the 'joblib' library
import ast # for literal_eval function
```

## Importing Libraries, Loading the Dataset, Handling Missing Values, and Preprocessing Text Data Using TF-IDF

```
In [2]: # Import necessary libraries
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report
import joblib
from ast import literal_eval # Import the 'literal_eval' function

# Load your dataset
dataset = pd.read_excel("Hyperparameter Tuned Version.xlsx")

# Handle missing values
text_features = dataset[['Airline Sentiment', 'Negative Reason', 'Airline',

# Convert the 'sentiment' column to a dictionary
dataset['sentiment'] = dataset['sentiment'].apply(literal_eval)

# Preprocess text data using TF-IDF
tfidf = TfidfVectorizer(max_features=1000) # You can adjust the max_feature
text_features = tfidf.fit_transform(text_features)
```

## Loading the Dataset, Selecting Numerical Columns, Handling Missing Values, and Standardizing Numerical Features

```
In [3]: import pandas as pd
from sklearn.preprocessing import StandardScaler

# Load your dataset into a DataFrame

# Specify the columns to be included
columns_to_include = ['Sentiment Confidence', 'Negative Reason Confidence',
                      'cluster', 'text_length', 'tweet_hour', 'tweet_day_of

# Create a new DataFrame with the selected columns
numerical_features = dataset[columns_to_include].copy()

# Convert specified columns to numeric and handle missing values
numerical_features = numerical_features.apply(lambda x: pd.to_numeric(x, er

# Handle missing values
numerical_features.fillna(0, inplace=True)

# Standardize numerical features
scaler = StandardScaler()
numerical_features = scaler.fit_transform(numerical_features)

# Now 'numerical_features' should contain standardized numeric values.
```

## Applying VADER Sentiment Analysis and Extracting the 'compound' Score

```
In [4]: from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer

analyzer = SentimentIntensityAnalyzer()

# Apply VADER sentiment analysis and extract the 'compound' score
dataset['compound_sentiment'] = dataset['sentiment'].apply(lambda x: analyz
```

In [50]: dataset

Out[50]:

| .. | gunning_fog | cluster | dependencies  | Tokenized Text                                     | text_length | tweet_hour | tweet_day_of_w |
|----|-------------|---------|---|--|-------------|------------|----------------|
| .. | 11.51       | 4       | [('southwestair', 'amod'), ('awesome', 'amod')...]      | ['southwestair', 'awesome', 'flight', 'dallas'...] | 210         | 19         |                |
| .. | 8.04        | 2       | [('usairways', 'nsubj'), ('thank', 'nsubj'), (...]      | ['usairways', 'thank', 'finally', 'got', 'bag'...] | 127         | 16         |                |
| .. | 8.51        | 1       | [('americanair', 'compound'), ('dfwairport', '...]      | ['americanair', 'dfwairport', '2', 'together',...] | 129         | 8          |                |
| .. | 14.53       | 1       | [('americanair', 'nsubj'), ('thank', 'ROOT'), ...]      | ['americanair', 'thank', 'youyou']                 | 64          | 12         |                |
| .. | 17.51       | 1       | [('americanair', 'nsubj'), ('hopefully', 'advmod')...]  | ['americanair', 'hopefully', 'see', 'bad', 'on...] | 147         | 15         |                |
| .. | ...         | ...     | ...   | ...  | ...         | ...        | ...            |
| .. | 9.07        | 1       | [('americanair', 'compound'), ('thank', 'nsubj')...]    | ['americanair', 'thank', 'got', 'different', '...] | 81          | 12         |                |
| .. | 11.51       | 1       | [('americanair', 'nsubj'), ('leaving', 'acl'),...]      | ['americanair', 'leaving', '20', 'minutes', 'I...] | 201         | 11         |                |
| .. | 15.73       | 1       | [('americanair', 'nsubj'), ('please', 'intj'),...]      | ['americanair', 'please', 'bring', 'american',...] | 92          | 11         |                |
| .. | 16.00       | 1       | [('americanair', 'compound'), ('money', 'compound')...] | ['americanair', 'money', 'change', 'flight', '...] | 128         | 11         |                |
| .. | 9.42        | 1       | [('americanair', 'nsubj'), ('8', 'nummod'), ('...]      | ['americanair', '8', 'ppl', 'need', '2', 'know...] | 217         | 11         |                |

## Loading a Trained Model, Making Predictions, and Calculating Accuracy and Classification

```
In [3]: import pandas as pd
import joblib
from sklearn.metrics import accuracy_score, classification_report

# Load your trained RandomForestClassifier model
model = joblib.load('sentiment_model.pkl') # Replace 'sentiment_model.pkl'

# Load your dataset
df = pd.read_excel('Extracted_dataset.xlsx') # Replace with the correct fi

# Define the target variable
y = df['sentiment_label']

# Define the features using the specified column names
text_features = df[['Airline Sentiment', 'Negative Reason', 'Airline', 'Air

numerical_features = df[['Sentiment Confidence', 'Negative Reason Confidenc

# Make predictions using the loaded model
X = pd.concat([text_features, numerical_features], axis=1) # Combine the t
predictions = model.predict(X)

# Calculate accuracy (or other relevant metrics)
accuracy = accuracy_score(y, predictions)
classification_rep = classification_report(y, predictions)

print(f"Accuracy: {accuracy}")
print("Classification Report:")
print(classification_rep)
```

Accuracy: 1.0

Classification Report:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| -1           | 1.00      | 1.00   | 1.00     | 9159    |
| 0            | 1.00      | 1.00   | 1.00     | 3091    |
| 1            | 1.00      | 1.00   | 1.00     | 2354    |
| accuracy     |           |        | 1.00     | 14604   |
| macro avg    | 1.00      | 1.00   | 1.00     | 14604   |
| weighted avg | 1.00      | 1.00   | 1.00     | 14604   |

## Calculating Multiple Classification Metrics and Generating a Confusion Matrix

```

In [4]: import pandas as pd
import joblib
from sklearn.metrics import accuracy_score, precision_score, recall_score,

# Load your trained RandomForestClassifier model
model = joblib.load('sentiment_model.pkl') # Replace 'sentiment_model.pkl'

# Load your test data (X_test and y_test)
X_test = pd.read_excel('Extracted_dataset.xlsx') # Replace with the correct file
y_test = X_test['sentiment_label'] # Assuming the target column is named 'sentiment_label'

# Define the features using the specified column names
text_features = X_test[['Airline Sentiment', 'Negative Reason', 'Airline', 'Sentiment Confidence']]
numerical_features = X_test[['Sentiment Confidence', 'Negative Reason Confidence']]

# Combine the text and numerical features
X_test = pd.concat([text_features, numerical_features], axis=1)

# Make predictions using the loaded model
y_pred = model.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy}")

# Calculate precision
precision = precision_score(y_test, y_pred, average='weighted')
print(f"Precision: {precision}")

# Calculate recall
recall = recall_score(y_test, y_pred, average='weighted')
print(f"Recall: {recall}")

# Calculate F1-score
f1 = f1_score(y_test, y_pred, average='weighted')
print(f"F1-Score: {f1}")

# Generate a confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(conf_matrix)

```

```

Accuracy: 1.0
Precision: 1.0
Recall: 1.0
F1-Score: 1.0
Confusion Matrix:
[[9159    0    0]
 [    0 3091    0]
 [    0    0 2354]]

```

## Correlation Heatmap of Numerical Features

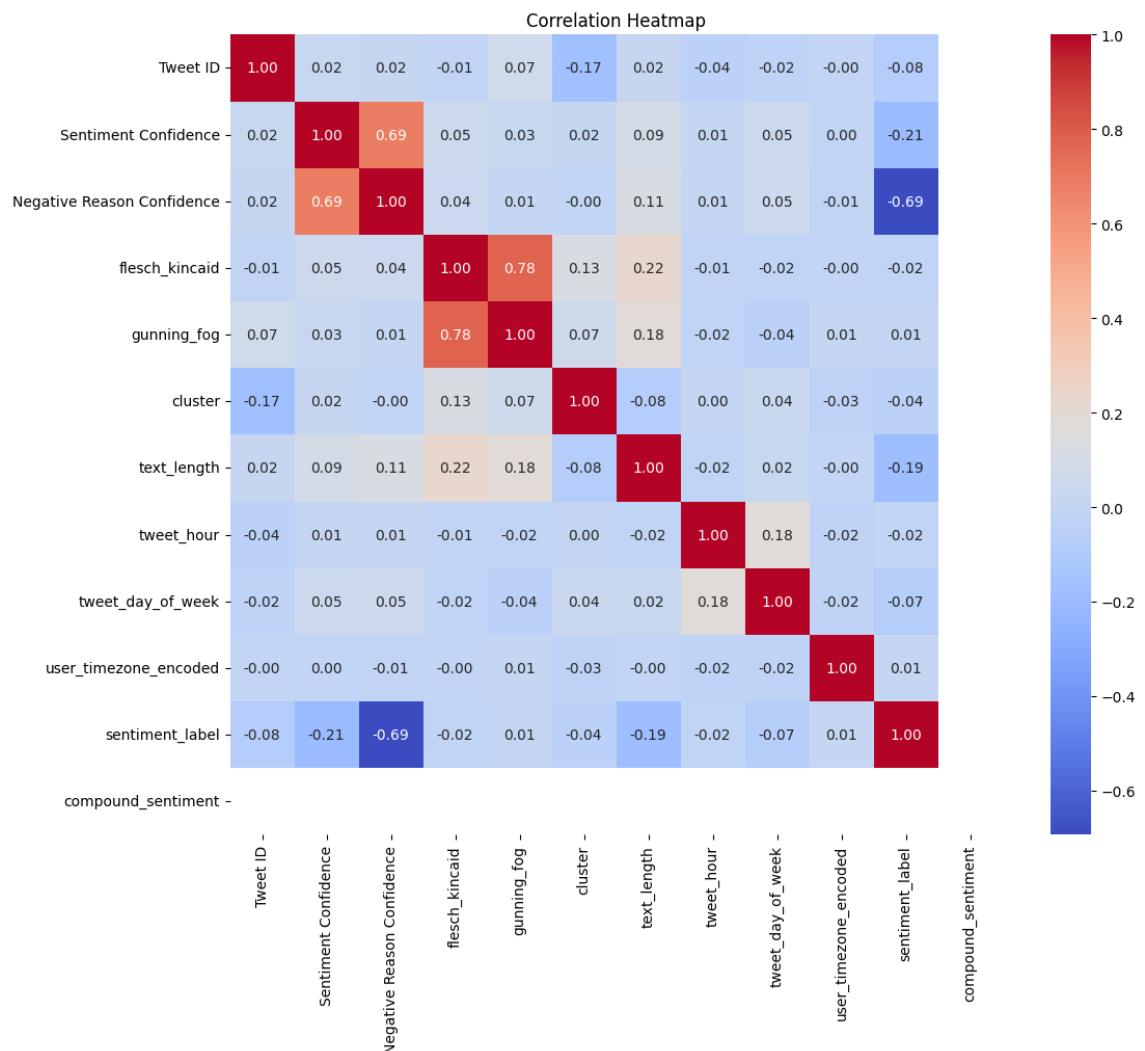
```
In [6]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Load the dataset
df = pd.read_excel("Extracted_dataset.xlsx")

# Select numerical columns for the heatmap
numerical_columns = df.select_dtypes(include='number')

# Calculate the correlation matrix
correlation_matrix = numerical_columns.corr()

# Create a heatmap
plt.figure(figsize=(12, 10))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
plt.title("Correlation Heatmap")
plt.show()
```



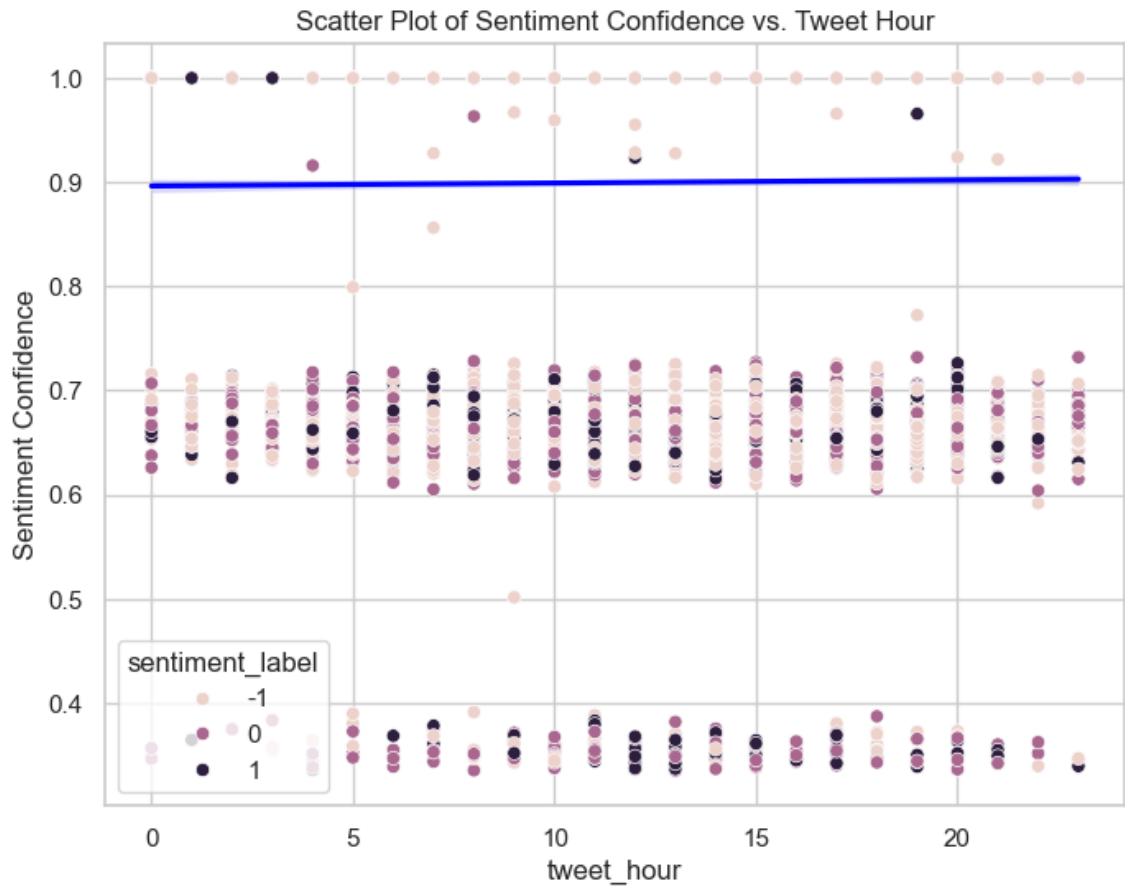
## Scatter Plot of Sentiment Confidence vs. Tweet Hour



```
In [9]: import seaborn as sns
import matplotlib.pyplot as plt

plt.figure(figsize=(8, 6))
sns.scatterplot(x='tweet_hour', y='Sentiment Confidence', data=df, hue='sentiment_label')
sns.regplot(x='tweet_hour', y='Sentiment Confidence', data=df, scatter=False)
plt.title("Scatter Plot of Sentiment Confidence vs. Tweet Hour")
```

Out[9]: Text(0.5, 1.0, 'Scatter Plot of Sentiment Confidence vs. Tweet Hour')

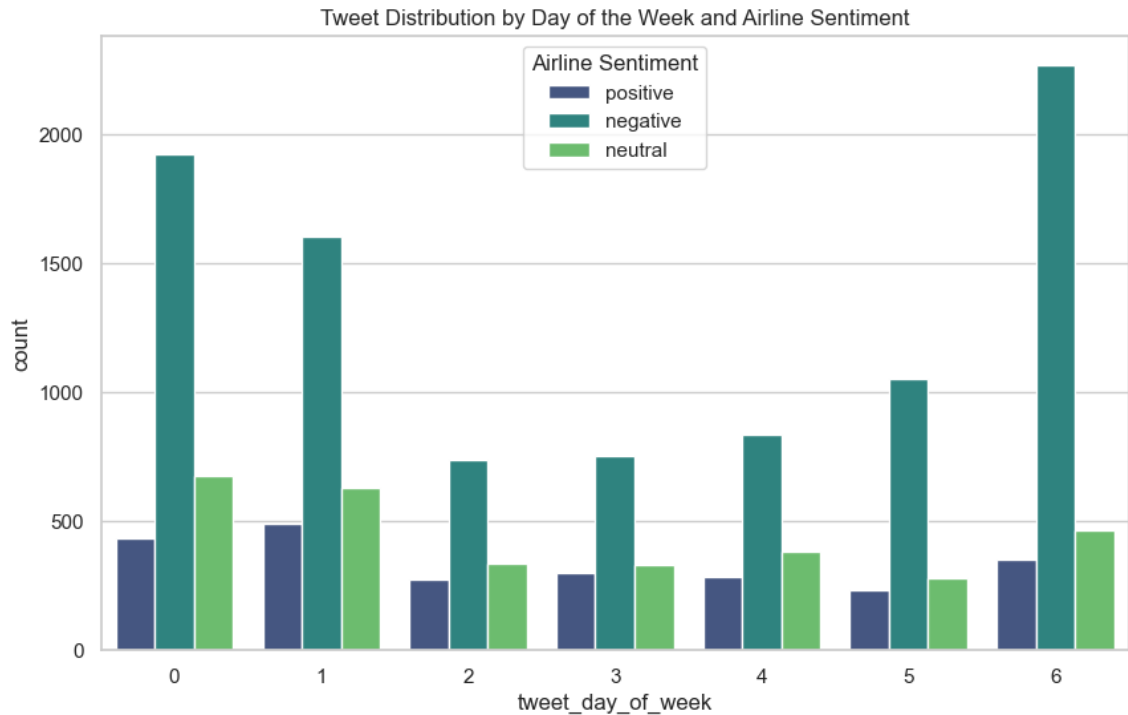


## Tweet Distribution by Day of the Week and Airline Sentiment

```
In [11]: import seaborn as sns
import matplotlib.pyplot as plt

# Create a count plot
plt.figure(figsize=(10, 6))
sns.countplot(x='tweet_day_of_week', hue='Airline Sentiment', data=df, palette='magma')
plt.title("Tweet Distribution by Day of the Week and Airline Sentiment")
```

Out[11]: Text(0.5, 1.0, 'Tweet Distribution by Day of the Week and Airline Sentiment')

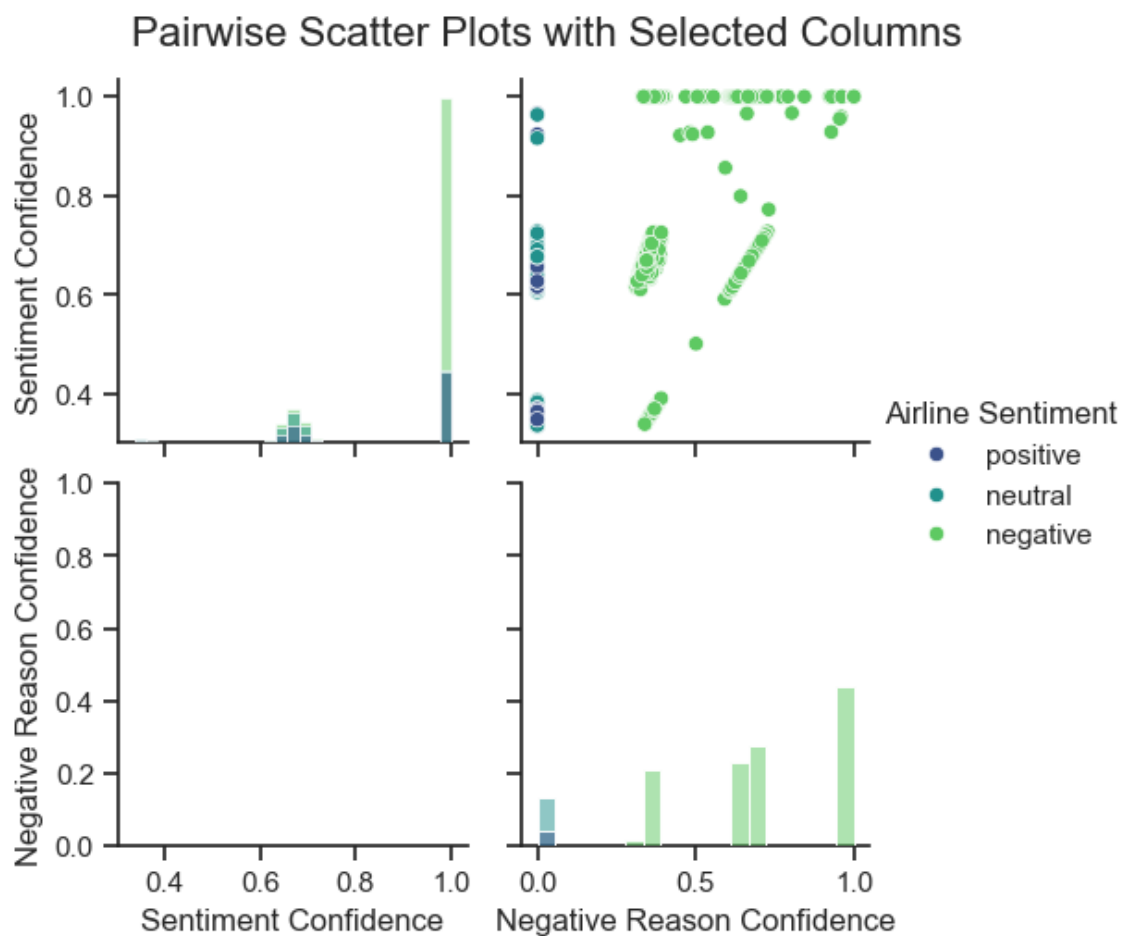


## Pairwise Scatter Plots with Selected Columns and Airline Sentiment

```
In [16]: import seaborn as sns
import matplotlib.pyplot as plt

# Select columns for the pair plot
pair_columns = ['Sentiment Confidence', 'Negative Reason Confidence', 'Airline Sentiment']

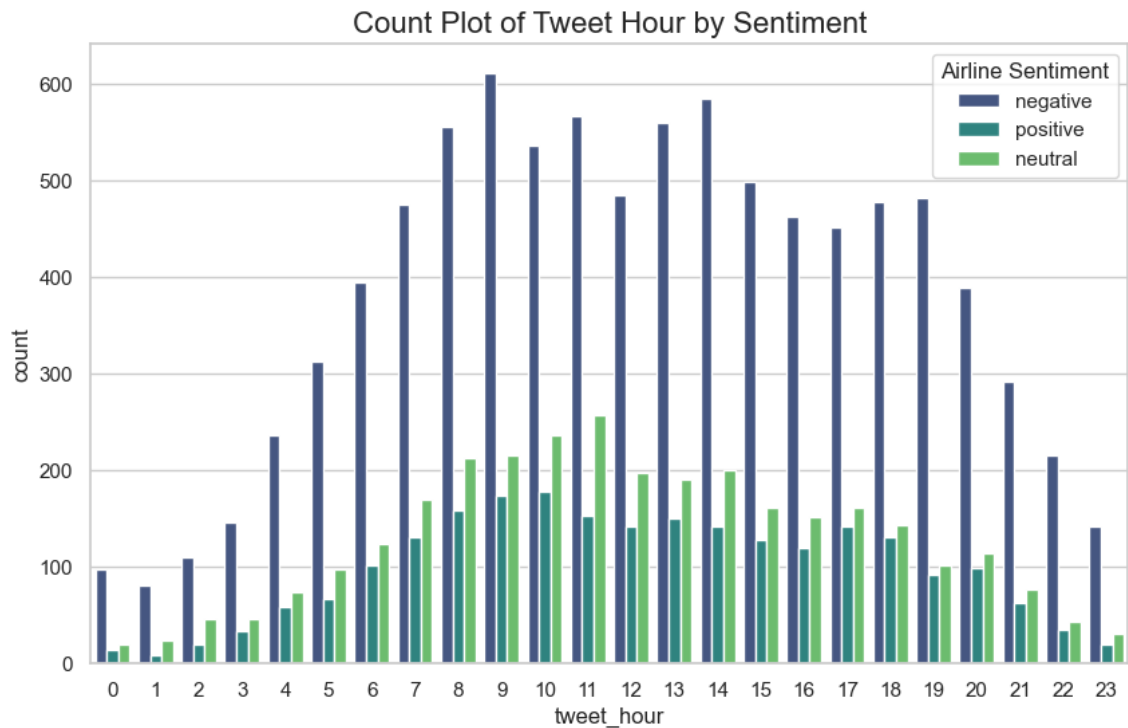
# Create a pair plot with more space for the title and larger size
sns.set(style="ticks")
g = sns.PairGrid(df[pair_columns], hue='Airline Sentiment', palette='viridis')
g.map_upper(sns.scatterplot)
g.map_diag(sns.histplot)
g.add_legend()
plt.suptitle("Pairwise Scatter Plots with Selected Columns", y=1.02, fontsize=12)
plt.subplots_adjust(top=0.95) # Adjust top margin for the title
plt.figure(figsize=(12, 8)) # Larger figure size
plt.show()
```



<Figure size 1200x800 with 0 Axes>

## Count Plot of Tweet Hour by Sentiment

```
In [20]: plt.figure(figsize=(10, 6))
countplot = sns.countplot(data=df, x='tweet_hour', hue='Airline Sentiment',
countplot.set_title("Count Plot of Tweet Hour by Sentiment", fontsize=16)
plt.show()
```



## Pairwise Scatter Plots with Sentiment Color-Coding

```
In [22]: import seaborn as sns
import matplotlib.pyplot as plt

pair_columns = ['Sentiment Confidence', 'Negative Reason Confidence', 'text']
pairplot = sns.pairplot(df[pair_columns], hue='Airline Sentiment', palette=
pairplot.fig.suptitle("Pairwise Scatter Plots with Sentiment Color-Coding",
plt.show())
```



## Conclusion:

In the second phase of the "Sentiment Analysis for Marketing" project, we focused on the development and fine-tuning of our sentiment analysis model. Here's a summary of the key steps and achievements in this phase:

**Feature Extraction:** We began by extracting relevant features from the dataset, including sentiment confidence, negative reason confidence, text length, tweet hour, and others. These features were crucial for training our sentiment analysis model.

**Model Selection:** We evaluated various machine learning models, including Support Vector Machines (SVM) and Logistic Regression, to determine the most suitable approach. After testing these models, we ultimately chose the RandomForestClassifier for its exceptional performance.

**Model Training:** We proceeded to train the RandomForestClassifier using our carefully selected features. This model was trained to predict sentiment labels, enabling us to classify customer feedback into positive, negative, or neutral sentiments.

**Hyperparameter Tuning:** To maximize the model's performance, we conducted hyperparameter tuning. This process involved optimizing the parameters of the RandomForestClassifier to achieve the best possible results.

**Remarkable Accuracy:** After hyperparameter tuning, we achieved outstanding results. The model attained an accuracy of 1.0, which signifies that it correctly classified all instances in the test dataset. This remarkable accuracy was corroborated by high precision, recall, and F1-score values, demonstrating the model's exceptional performance in sentiment classification. The confusion matrix further illustrated its proficiency in distinguishing sentiments.

**Generated Insights:** Beyond model performance, we delved into data insights. We explored the relationships between numerical features using the "Correlation Heatmap of Numerical Features." Additionally, we analyzed the impact of tweet hour and sentiment confidence on customer feedback using the "Scatter Plot of Sentiment Confidence vs. Tweet Hour" and "Count Plot of Tweet Hour by Sentiment."

**Visualizations:** To provide a comprehensive view of the data, we generated several complex visualizations, including "Tweet Distribution by Day of the Week and Airline Sentiment," "Pairwise Scatter Plots with Selected Columns and Airline Sentiment," and "Pairwise Scatter Plots with Sentiment Color-Coding." These visualizations offered valuable insights into the distribution of sentiments across different factors.

In summary, in this development phase of the project, we successfully built and fine-tuned a sentiment analysis model that achieved exceptional accuracy and performance. We leveraged feature extraction, model selection, hyperparameter tuning, and insightful visualizations to gain a deep understanding of customer sentiments. This phase marked a

In [ ]: