

1. Overview

This Jupyter notebook focuses on analyzing the sentiment of tweets related to technological products, specifically Apple and Google.

The data is first explored to understand its structure, including the distribution of categorical variables and the handling of missing values. Various preprocessing steps are applied, such as cleaning the tweet text, tokenization, and removing stopwords. A function is developed to identify the brand of the technological product mentioned in each tweet so it will be an additional feature in our dataset.

2. Data Understanding

2.1 Data Description

2.2 Import necessary libraries

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 %matplotlib inline
        5 import seaborn as sns
        6 import re # Import regular expressions library
        7
        8 import nltk
        9 nltk.download('punkt')
       10 from nltk.corpus import stopwords
       11 from nltk.tokenize import word_tokenize
```

[nltk_data] Downloading package punkt to C:\Users\Usuario\nltk_data...

[nltk_data] Package punkt is already up-to-date!

2.3 Define global variables

```
In [2]: 1 path = '../data/judge-1377884607_tweet_product_company.csv'
```

3.4 Functions

```
In [3]: 1 def plot_categorical_proportions(df):
2         """
3         Plots bar charts for each categorical variable in a DataFrame, showing
4         ordered by proportion in descending order. Each bar is labeled with its
5         proportion.
6
7         Inputs:
8         df (pd.DataFrame): The DataFrame to analyze.
9
10        Outputs:
11        None
12
13        Description:
14        This function identifies categorical variables, calculates the proportion
15        and plots a bar chart for each categorical variable. Labels on the bars
16        excluding the 'tweet_text' column.
17        """
18        # Excluding 'tweet_text' column
19        df = df.drop(columns=['tweet_text'])
20
21        for col in df.columns:
22            # Calculating proportions
23            value_counts = df[col].value_counts(normalize=True).sort_values(ascending=False)
24            percentages = value_counts * 100 # Convert proportions to percentages
25
26            # Plotting
27            plt.figure(figsize=(10, 6))
28            ax = percentages.plot(kind='bar')
29            ax.set_title(f'Proportion of Categories in {col}')
30            ax.set_ylabel('Percentage')
31
32            # Adding percentage labels on the bars
33            for p in ax.patches:
34                ax.annotate(f'{p.get_height():.2f}%', (p.get_x() + p.get_width() / 2,
35                                                            p.get_height() * 1.05),
36                            ha='center', va='bottom', xytext=(0, 10), textcoords='offsetpoints')
```

In [4]:

```
1 def plot_grouped_charts(df):
2     """
3     Creates combined plots for each column in the DataFrame based on their
4     For numeric columns, histograms for all statuses are combined in one p
5     For categorical columns, grouped bar charts are created.
6
7     Args:
8         df (pd.DataFrame): The DataFrame containing the data.
9     """
10    status_col = 'is_there_an_emotion_directed_at_a_brand_or_product'
11    unique_statuses = df[status_col].unique()
12    colors = plt.get_cmap('tab10') # Fetches a colormap with distinct col
13
14    for col in df.columns:
15        if col not in [status_col, 'tweet_text']:
16            if df[col].dtype in ['int64', 'float64']: # Numeric Columns
17                plt.figure(figsize=(12, 6))
18
19                # Histogram for all statuses
20                for i, status in enumerate(unique_statuses):
21                    sns.histplot(df[df[status_col] == status][col], kde=True,
22                                stat='density', label=str(status), color=
23
24                plt.title(f'Combined Histogram of {col} by {status_col}')
25                plt.legend(title=status_col)
26                plt.show()
27
28                # Boxplot for all statuses
29                plt.figure(figsize=(12, 6))
30                sns.boxplot(x=status_col, y=col, data=df, palette='tab10')
31                plt.title(f'Combined Boxplot of {col} by {status_col}')
32                plt.show()
33
34            elif df[col].dtype == 'object': # Categorical Columns
35                plt.figure(figsize=(10, 6))
36                sns.countplot(data=df, x=status_col, hue=col)
37                plt.title(f'Grouped Bar Chart of {status_col} by {col}')
38                plt.ylabel('Count')
39                plt.xlabel(status_col)
40                plt.legend(title=col, loc='upper right')
41                plt.xticks(rotation=45)
42                plt.show()
43
```

In [5]:

```
1 def txt_clean(txt):
2     """
3     Clean and preprocess text data for further analysis.
4
5     Parameters:
6
7         txt (str): The text string that needs to be cleaned and tokenized.
8
9     Returns:
10
11         list: A list of cleaned and tokenized words, where punctuation and
12              is converted to lowercase, stopwords and Twitter mentions are removed,
13              strings are filtered out.
14
15     """
16
17     # List of additional strange characters to remove
18     strange_chars = '!"$%&\'()*+,-./:;<=>?[\\]^_`{|}~“”!#•0a'
19
20     sw = stopwords.words('english')
21     sw.extend(['link', 'rt', 'get'])
22     no_accents_re = re.compile('^[a-z]+$')
23     twitter_re = re.compile('@[a-zA-Z]*')
24
25     # Replace punctuation and strange characters with spaces
26     txt = txt.translate(str.maketrans(strange_chars, ' ' * len(strange_chars)))
27
28     # Tokenize the text
29     tokens = word_tokenize(txt)
30
31     # Convert to lowercase
32     tokens = [w.lower() for w in tokens]
33     # Remove @ mentions
34     tokens = [w for w in tokens if not twitter_re.match(w)]
35     # Remove words with accents
36     tokens = [w for w in tokens if no_accents_re.match(w)]
37     # Remove stopwords
38     tokens = [w for w in tokens if w not in sw]
39     # Remove empty strings
40     tokens = [w for w in tokens if w]
41
42     return tokens
```

In [6]:

```
1  # Define the function to identify if the tweet is about a Google or Apple
2  def identify_product(tweet_text):
3      """
4      Identify if the tweet is about a Google or Apple product.
5
6      Parameters:
7      tweet_text (str): The text of the tweet.
8
9      Returns:
10     str: 'Google' if the tweet mentions a Google product, 'Apple' if the t
11         'Both' if the tweet mentions both, 'Unknown' if it mentions neith
12     """
13     google_keywords = ['google', 'pixel', 'pixels', 'nexus', 'nexuses', 'a
14                        'chromebook', 'chromebooks', 'nest', 'nests', 'stad
15     apple_keywords = ['apple', 'apples', 'iphone', 'iphones', 'ipad', 'ipa
16                      'macbooks', 'imac', 'imacs', 'watch', 'watches', 'ai
17                      'appstore', 'ios', 'itunes']
18
19     # Ensure tweet_text is a string
20     if not isinstance(tweet_text, str):
21         return 'Unknown'
22
23     # Replace "app store" with "appstore" before tokenization
24     tweet_text = tweet_text.replace("app store", "appstore")
25
26     # Remove all numbers from the tweet text
27     tweet_text = re.sub(r'\d+', '', tweet_text)
28
29     # Clean the text and obtain tokens
30     tokens = txt_clean(tweet_text)
31
32     # Check if any keyword exists as a substring within the tokens
33     google_mentioned = any(keyword in token for keyword in google_keyw
34     apple_mentioned = any(keyword in token for keyword in apple_keywor
35
36     if google_mentioned and apple_mentioned:
37         return 'Both'
38     elif google_mentioned:
39         return 'Google'
40     elif apple_mentioned:
41         return 'Apple'
42     else:
43         return 'Unknown'
```

```
In [7]: 1 def plot_emotion_distribution(df, product_column='product_mention', emotion_column='emotion_type')
2         """
3         Plots a bar chart showing the distribution of emotion types by product
4         with annotations displaying the counts and percentages.
5
6         Parameters:
7         df (pandas.DataFrame): The DataFrame containing the data.
8         product_column (str): The name of the column that contains product mention
9         emotion_column (str): The name of the column that contains emotion type
10        """
11        # Calculate the counts and normalize to get percentages
12        counts = df.groupby([product_column, emotion_column]).size().reset_index()
13        total_counts = df[product_column].value_counts().reset_index()
14        total_counts.columns = [product_column, 'total']
15
16        # Merge counts with totals to calculate percentages
17        counts = counts.merge(total_counts, on=product_column)
18        counts['percentage'] = (counts['counts'] / counts['total']) * 100
19
20        # Plotting
21        plt.figure(figsize=(10, 6))
22        ax = sns.barplot(data=counts, x=product_column, y='percentage', hue=emotion_column)
23
24        # Annotate each bar with the corresponding count and percentage
25        for p in ax.patches:
26            height = p.get_height()
27            ax.annotate(f'{int(height)} ({height:.1f}%)',
28                      xy=(p.get_x() + p.get_width() / 2, height),
29                      xytext=(0, 5), # Offset label position above the bar
30                      textcoords='offset points',
31                      ha='center', va='center')
32
33        plt.title('Distribution of Emotion Types by Product Mention')
34        plt.ylabel('Percentage')
35        plt.xlabel('Product Mention')
36        plt.xticks(rotation=45)
37        plt.legend(title=emotion_column)
38        plt.show()
39
```

```
In [8]: 1 def get_contingency_table_with_percentage_sign(df, product_column='product'
2         """
3         Generates a contingency table for the given DataFrame based on product
4         showing the proportions as percentages with a percentage sign, rounded
5
6         Parameters:
7         df (pandas.DataFrame): The DataFrame containing the data.
8         product_column (str): The name of the column that contains product men
9         emotion_column (str): The name of the column that contains emotion typ
10
11        Returns:
12        pd.DataFrame: A contingency table with proportions (as percentages) of
13        """
14        # Create a contingency table with counts
15        contingency_table = pd.crosstab(df[product_column], df[emotion_column]
16
17        # Convert counts to proportions by dividing each cell by the row sum a
18        contingency_table_percentage = contingency_table.div(contingency_table
19
20        # Round the percentages to 2 decimal places and add the percentage sig
21        contingency_table_percentage = contingency_table_percentage.round(2).a
22
23        return contingency_table_percentage
```

3.5 Code

3.5.1 Exploratory Analysis

3.5.1.1 Looking at the dataset

```
In [9]: 1 df_tweets = pd.read_csv(path, encoding='ISO-8859-1')
        2 df_tweets.head()
```

Out[9]:

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

```
In [10]: 1 df_tweets.shape
```

Out[10]: (9093, 3)

3.5.1.2 Looking at the data types

```
In [11]: 1 # Let's start by having a look at the type of each column
        2 df_tweets.dtypes
```

```
Out[11]: tweet_text                object
emotion_in_tweet_is_directed_at    object
is_there_an_emotion_directed_at_a_brand_or_product  object
dtype: object
```

3.5.1.3 Null values

```
In [12]: 1 # Let's see how the proportion of null values
        2 (df_tweets.isna().sum()/len(df_tweets))*100
```

```
Out[12]: tweet_text                0.010997
emotion_in_tweet_is_directed_at    63.807324
is_there_an_emotion_directed_at_a_brand_or_product  0.000000
dtype: float64
```

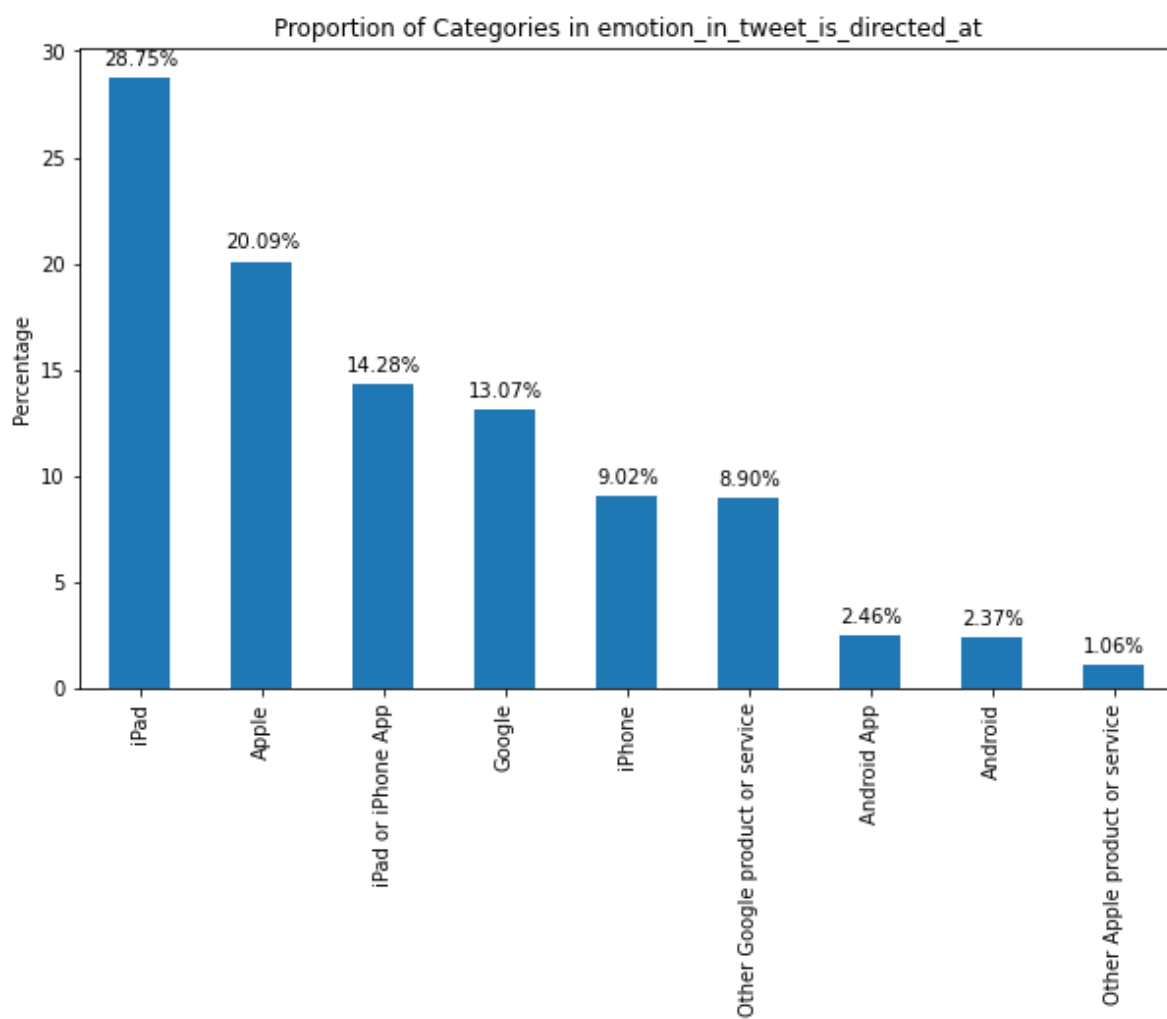

3.5.2 Descriptive Analysis

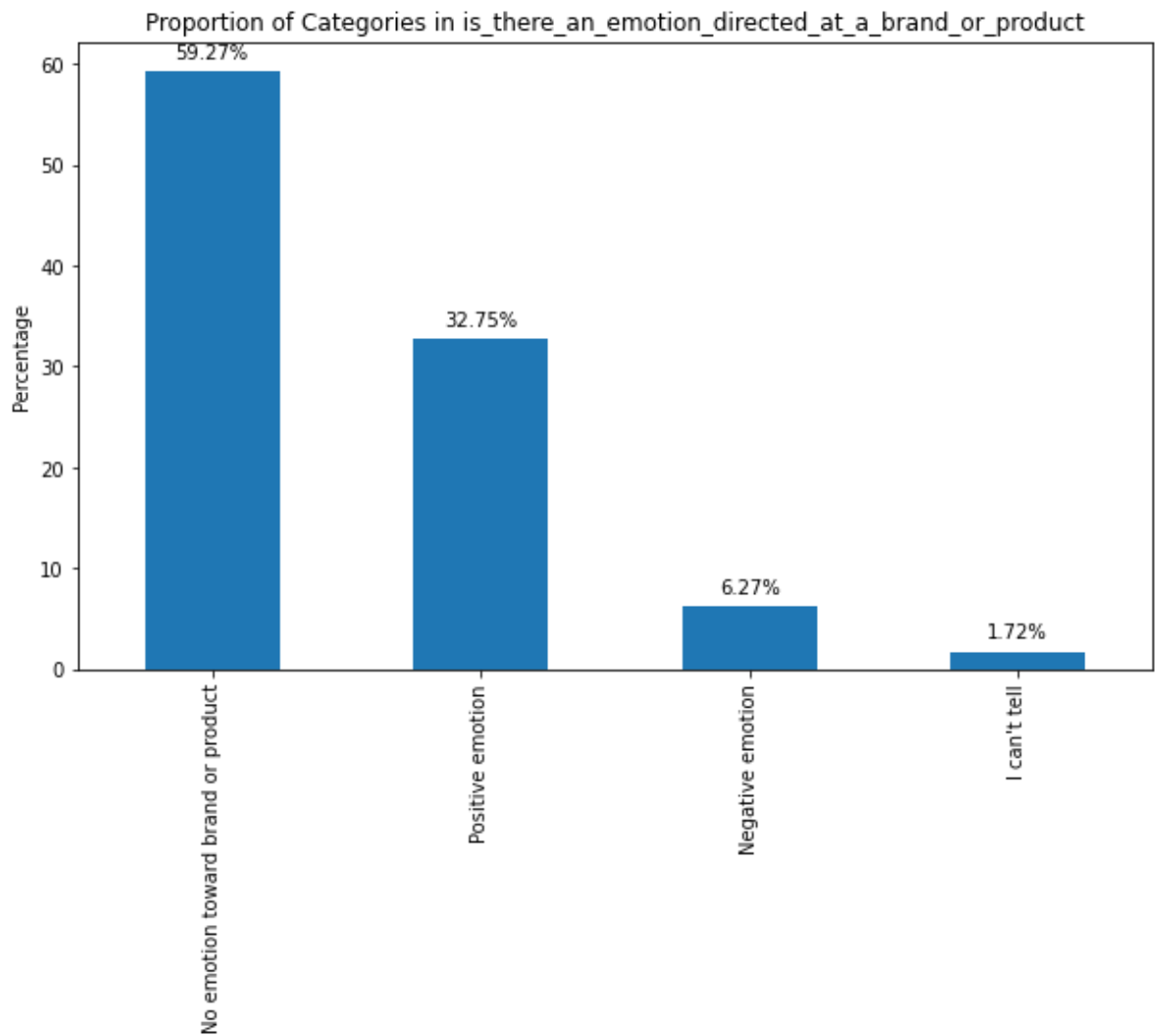
All the columns from our dataset are categorical

3.5.2.1 Unvaried Analysis

Let's see the Unvaried results of the initial dataset without any changes

```
In [13]: 1 plot_categorical_proportions(df_tweets)
```





As we can see in the column `emotion_in_tweet_is_directed_at`, most of the tweets are about iPad or Apple devices. However, there are labels marked as different but that are actually talking about the same category. For example, Android and Android App.

Regarding the `is_there_an_emotion_directed_at_a_brand_or_product` variable, we can see that only 39.02% of the tweets have actually either a positive or negative emotion.

Considering the miss classification in the column `emotion_in_tweet_is_directed_at`, let's start by looking at the tweets related to the null values of that column.

```
In [14]: 1 # Filter rows where the column 'emotion_in_tweet_is_directed_at' is null
          2 null_emotion_tweets = df_tweets[df_tweets['emotion_in_tweet_is_directed_at']
          3
          4 # Display the first 5 records of the 'tweet_text' column where 'emotion_in
          5 null_emotion_tweets['tweet_text']
```

```
Out[14]: 5      @teachntech00 New iPad Apps For #SpeechTherapy...
          6                                             NaN
          16      Holler Gram for iPad on the iTunes App Store -...
          32      Attn: All #SXSW frineds, @mention Register fo...
          33      Anyone at #sxsw want to sell their old iPad?

          ...
          9087      @mention Yup, but I don't have a third app yet...
          9089      Wave, buzz... RT @mention We interrupt your re...
          9090      Google's Zeiger, a physician never reported po...
          9091      Some Verizon iPhone customers complained their...
          9092      Ĩ;ÿäü_ÊÏÖÉÁâ_£â_ÛâRT @...
          Name: tweet_text, Length: 5802, dtype: object
```

Looking at the data, we identify that some labels regarding the device were miss classified. So based on our business problem, we decide to create a function capable of identifying, based on the tweet, which device it's talking about (ie, Google or Apple). And we create a new column with that classification

```
In [15]: 1 # Apply the function to the DataFrame and create a new column with the res
          2 df_tweets['product_mention'] = df_tweets['tweet_text'].map(identify_produc
```

As defined in our business problem, we have decided not to filter any of the labels in the product_mention field because we believe that the values Both and Unknown can provide us with relevant information.

```
In [16]: 1 # We will now drop the column emotion_in_tweet_is_directed_at
          2 df_tweets = df_tweets.drop('emotion_in_tweet_is_directed_at', axis=1)
```

We are going to get a sample of the data just to verify that the sentiments labels are correct.

In [17]:

```
1 # Filter the records for Apple and Google with emotions
2 apple_tweets = df_tweets[(df_tweets['product_mention'] == 'Apple') &
3                             (df_tweets['is_there_an_emotion_directed_at_a_brand_or_product'] == 'Positive emotion')]
4
5 google_tweets = df_tweets[(df_tweets['product_mention'] == 'Google') &
6                             (df_tweets['is_there_an_emotion_directed_at_a_brand_or_product'] == 'Positive emotion')]
7
8 # Get a random sample of 100 records for each product with emotions
9 apple_sample = apple_tweets.sample(n=100, random_state=42)
10 google_sample = google_tweets.sample(n=100, random_state=42)
11
12 # Combine the two samples
13 combined_sample = pd.concat([apple_sample, google_sample])
14
15 # Display the combined sample
16 combined_sample
```

Out[17]:

	tweet_text	is_there_an_emotion_directed_at_a_brand_or_product	product_mention
	RT @mention Apple Pop Up Store for #SXSW, Why ...	Positive emotion	Apple
5455			
	No, I didn't get an iPad 2 :(No, I'm not at #...	Positive emotion	Apple
4473			
	ZOMG its iPad 2 :p RT @mention Look everyone! ...	Positive emotion	Apple
3808			
	@mention Nice move on Apple #SXSW will be cent...	Positive emotion	Apple
4646			
	@mention talk about #sxsw and iPad is acceptab...	Positive emotion	Apple
8341			
...
	#smtravel hasn't heard - Google plans to launc...	Positive emotion	Google
609			
	Kinda giddy about #TheIndustryParty with #Goog...	Positive emotion	Google
2579			
	Google Art Project - like street view, except ...	Positive emotion	Google
8273			
	Just met Jared at the Android meetup #sxsw. Bi...	Positive emotion	Google
3947			
	Learned from sponsored #sxsw software: 1) Aust...	Positive emotion	Google
7717			

200 rows × 3 columns

After a rigorous review of each tweet, we conclude that we can trust the `is_there_an_emotion_directed_at_a_brand_or_product` labels in the dataset.

Let's rename the column `is_there_an_emotion_directed_at_a_brand_or_product` to `emotion_type`

```
In [18]: 1 df_tweets.rename(columns={'is_there_an_emotion_directed_at_a_brand_or_pro
```

```
In [19]: 1 df_tweets.head()
```

Out[19]:

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

```
In [20]: 1 df_tweets['emotion_type'].value_counts()
```

```
Out[20]: No emotion toward brand or product    5389
Positive emotion                             2978
Negative emotion                             570
I can't tell                                156
Name: emotion_type, dtype: int64
```

Let's rename the value 'No emotion toward brand or product' for 'unknown'

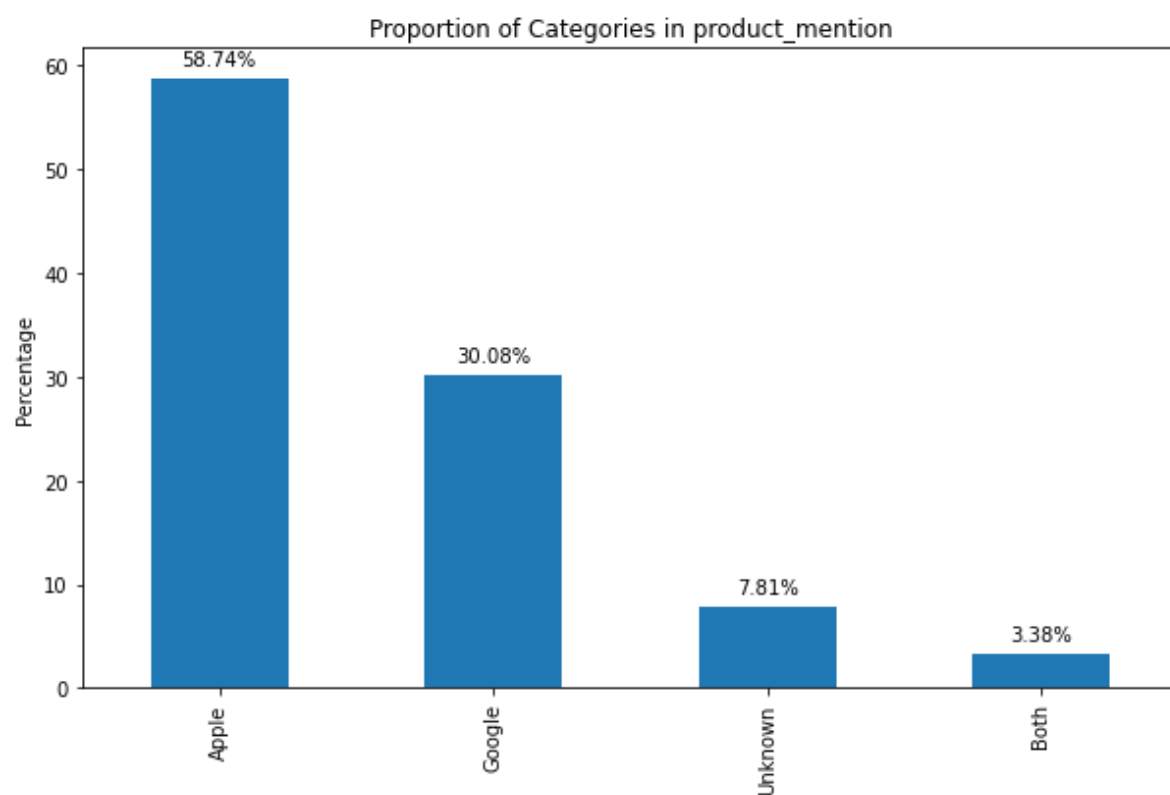
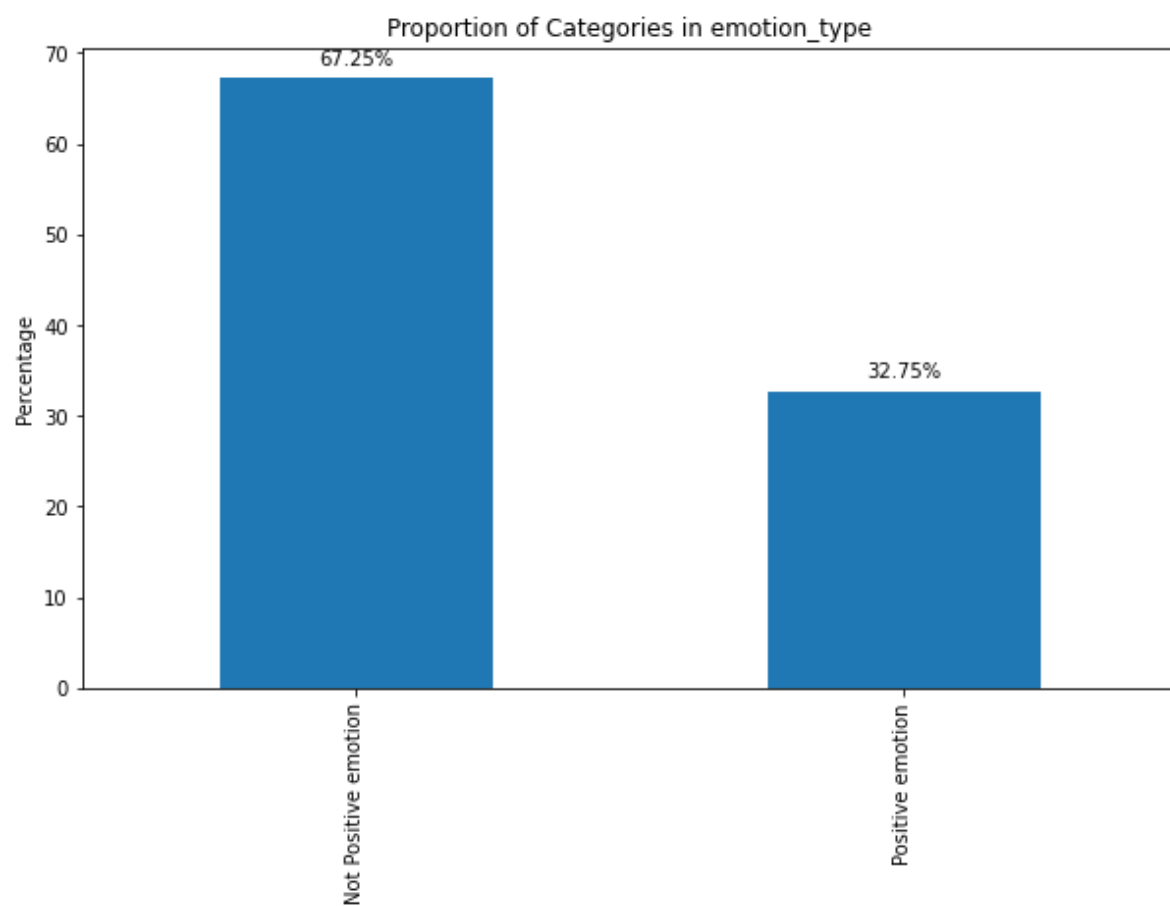
```
In [21]: 1 # Replace values in column 'emotion_type' where the value is "No emotion t
2 df_tweets['emotion_type'] = df_tweets['emotion_type'].replace("No emotion
```

We are going to do a replace of every emotion that is not positive to 'no positive' and leave the positive emotions as they are

```
In [22]: 1 df_tweets['emotion_type'] = df_tweets['emotion_type'].map(lambda x: 'Posit
```

Let's see the unvaried results after this new updates on `df_tweets`

```
In [23]: 1 plot_categorical_proportions(df_tweets)
```



```
In [24]: 1 df_tweets.columns
```

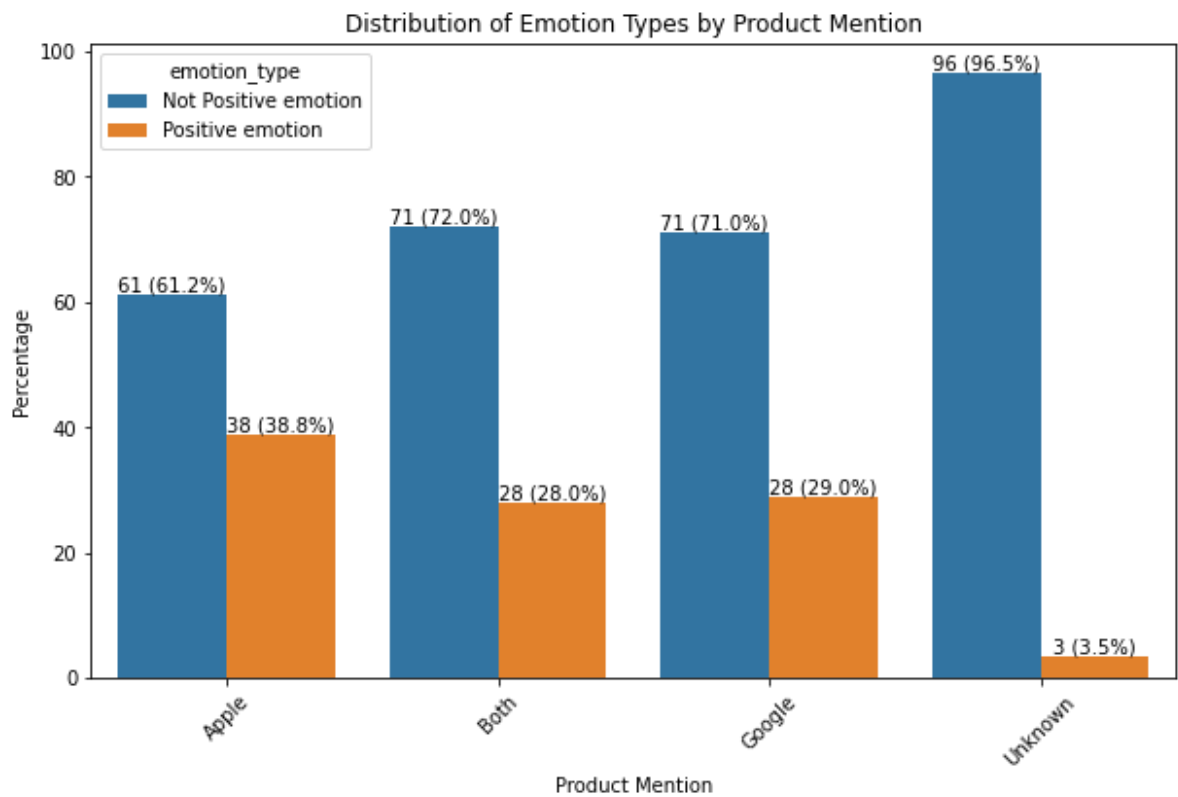
```
Out[24]: Index(['tweet_text', 'emotion_type', 'product_mention'], dtype='object')
```

3.5.2.2. Multivaried Analysis

Let's see the distribution of the sentiments for each company

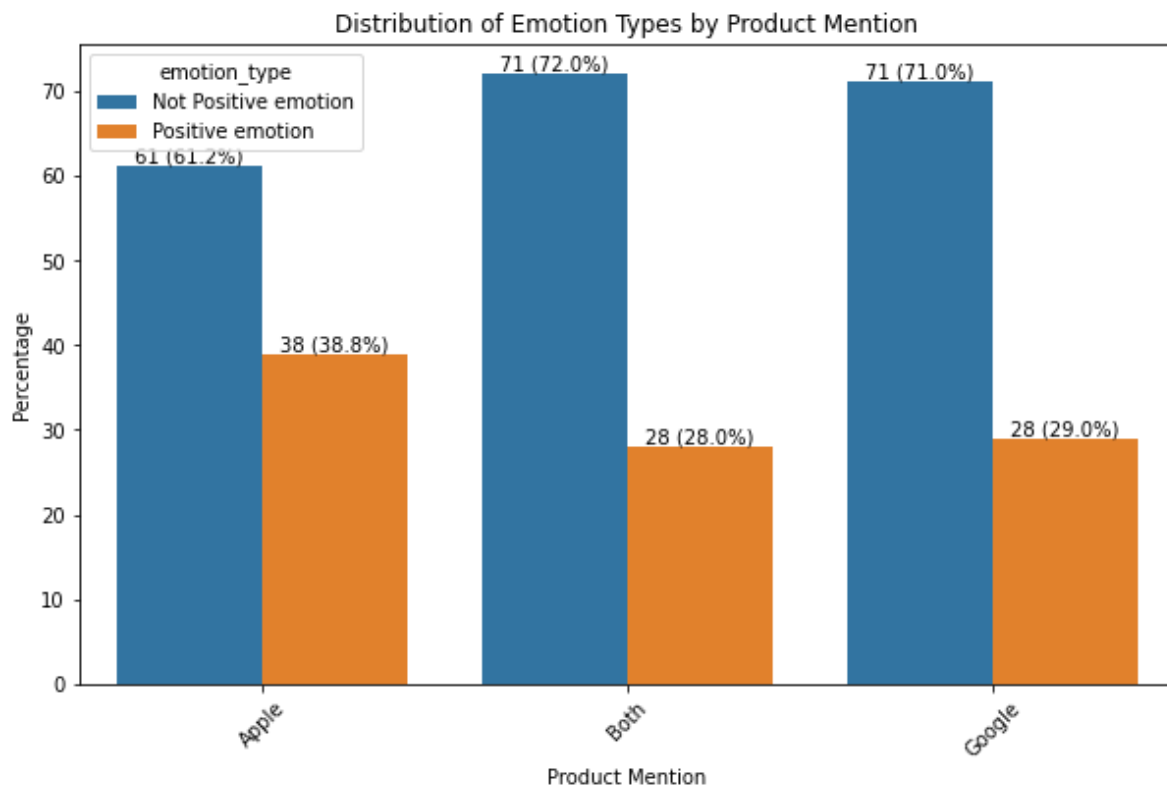
Bar graph

```
In [25]: 1 plot_emotion_distribution(df_tweets)
```



After looking at this graph, we are going to filter the unknowns because it's not giving us much information and it's not in our interests


```
In [26]: 1 # Filter out rows where product_mention is 'Unknown'
2 df_tweets = df_tweets[df_tweets['product_mention'] != 'Unknown']
3
4 plot_emotion_distribution(df_tweets)
```



As is observable in the image above, we don't have a problem of imbalance in our data as we usually use 3% as the threshold.

Contingency Tables

```
In [27]: 1 contingency_table_percentage = get_contingency_table_with_percentage_sign(
2
3 contingency_table_percentage
```

Out[27]:

emotion_type	Not Positive emotion	Positive emotion
product_mention		
Apple	61.17%	38.83%
Both	71.99%	28.01%
Google	71.01%	28.99%

6. Exporting the data

```
In [28]: 1 # Export the DataFrame to a CSV file
          2 # index=False ensures that the DataFrame index is not saved as a separate
          3 df_tweets.to_csv('..\data\df_tweets_clean.csv', index=False)
```