

In [1]: 1 `#pip install --upgrade pip`

In [2]: 1 `#pip install vaderSentiment`

In [3]: 1 `#pip install wordcloud`

```
In [4]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from colorama import Fore, init
6 import plotly.express as px
7
8 import string
9 import re
10 import nltk
11 from nltk.corpus import stopwords
12 from nltk.tokenize import word_tokenize
13 from nltk.stem import PorterStemmer
14 from nltk.sentiment import SentimentIntensityAnalyzer
15 from nltk import tokenize
16 from nltk.tokenize import sent_tokenize
17 from nltk.tokenize import word_tokenize
18 from tqdm.notebook import tqdm
19 from collections import Counter
20 from wordcloud import WordCloud
21
22 nltk.download('vader_lexicon')
23 nltk.download('punkt')
24 nltk.download('stopwords')
25
26 import warnings
27 warnings.filterwarnings('ignore')
```

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

```
In [5]: 1 df = pd.read_csv(r"C:\Users\Kishore\OneDrive\Desktop\Total Machine Learning\CSV Files\sentimentdataset.csv")  
        2 df
```

Out[5]:

	Unnamed: 0.1	Unnamed: 0	Text	Sentiment	Timestamp		User	Platform	Hashtags	Retweets
0	0	0	Enjoying a beautiful day at the park! ...	Positive	2023-01-15 12:30:00		User123	Twitter	#Nature #Park	15.0
1	1	1	Traffic was terrible this morning. ...	Negative	2023-01-15 08:45:00		CommuterX	Twitter	#Traffic #Morning	5.0
2	2	2	Just finished an amazing workout! 🏋️ ...	Positive	2023-01-15 15:45:00		FitnessFan	Instagram	#Fitness #Workout	20.0
3	3	3	Excited about the upcoming weekend getaway! ...	Positive	2023-01-15 18:20:00		AdventureX	Facebook	#Travel #Adventure	8.0
4	4	4	Trying out a new recipe for dinner tonight. ...	Neutral	2023-01-15 19:55:00		ChefCook	Instagram	#Cooking #Food	12.0
...
727	728	732	Collaborating on a science project that receiv...	Happy	2017-08-18 18:20:00	ScienceProjectSuccessHighSchool		Facebook	#ScienceFairWinner #HighSchoolScience	20.0
728	729	733	Attending a surprise birthday party organized ...	Happy	2018-06-22 14:15:00	BirthdayPartyJoyHighSchool		Instagram	#SurpriseCelebration #HighSchoolFriendship	25.0
729	730	734	Successfully fundraising for a school charity ...	Happy	2019-04-05 17:30:00	CharityFundraisingTriumphHighSchool		Twitter	#CommunityGiving #HighSchoolPhilanthropy	22.0
730	731	735	Participating in a multicultural festival, cel...	Happy	2020-02-29 20:45:00	MulticulturalFestivalJoyHighSchool		Facebook	#CulturalCelebration #HighSchoolUnity	21.0

	Unnamed: 0.1	Unnamed: 0	Text	Sentiment	Timestamp		User	Platform	Hashtags	Retweets
731	732	736	Organizing a virtual talent show during challe...	Happy	2020-11-15 15:15:00		VirtualTalentShowSuccessHighSchool	Instagram	#VirtualEntertainment #HighSchoolPositivity	24.0

732 rows × 15 columns

```
In [6]: 1 def null_count():
2         return pd.DataFrame({'features': df.columns,
3                               'dtypes': df.dtypes.values})
4         null_count()
```

Out[6]:

	features	dtypes
0	Unnamed: 0.1	int64
1	Unnamed: 0	int64
2	Text	object
3	Sentiment	object
4	Timestamp	object
5	User	object
6	Platform	object
7	Hashtags	object
8	Retweets	float64
9	Likes	float64
10	Country	object
11	Year	int64
12	Month	int64
13	Day	int64
14	Hour	int64

```
In [7]: 1 # drop column Unnamed: 0.1
        2 df.drop(columns='Unnamed: 0.1',inplace=True)
```

```
In [8]: 1 # Rename column Unnamed: 0 to id
        2 df.rename(columns={'Unnamed: 0':'Id'},inplace=True)
```

```
In [9]: 1 df.isnull().sum()
```

```
Out[9]: Id          0
        Text        0
        Sentiment   0
        Timestamp   0
        User        0
        Platform    0
        Hashtags    0
        Retweets    0
        Likes       0
        Country     0
        Year        0
        Month       0
        Day         0
        Hour        0
        dtype: int64
```

```
In [10]: 1 df.dtypes
```

```
Out[10]: Id                int64
Text                object
Sentiment           object
Timestamp           object
User                object
Platform            object
Hashtags            object
Retweets            float64
Likes               float64
Country             object
Year                int64
Month               int64
Day                 int64
Hour                int64
dtype: object
```

```
In [11]: 1 df.duplicated().sum()
```

```
Out[11]: 0
```

```
In [12]: 1 df.columns
```

```
Out[12]: Index(['Id', 'Text', 'Sentiment', 'Timestamp', 'User', 'Platform', 'Hashtags',
               'Retweets', 'Likes', 'Country', 'Year', 'Month', 'Day', 'Hour'],
              dtype='object')
```

```
In [13]: 1 for column in df.columns:
          2     num_distinct_values = len(df[column].unique())
          3     print(f"{column}: {num_distinct_values} distinct values")
```

```
Id: 732 distinct values
Text: 707 distinct values
Sentiment: 279 distinct values
Timestamp: 683 distinct values
User: 685 distinct values
Platform: 4 distinct values
Hashtags: 697 distinct values
Retweets: 26 distinct values
Likes: 38 distinct values
Country: 115 distinct values
Year: 14 distinct values
Month: 12 distinct values
Day: 31 distinct values
Hour: 22 distinct values
```

```
In [14]: 1 ## drop Columns
```

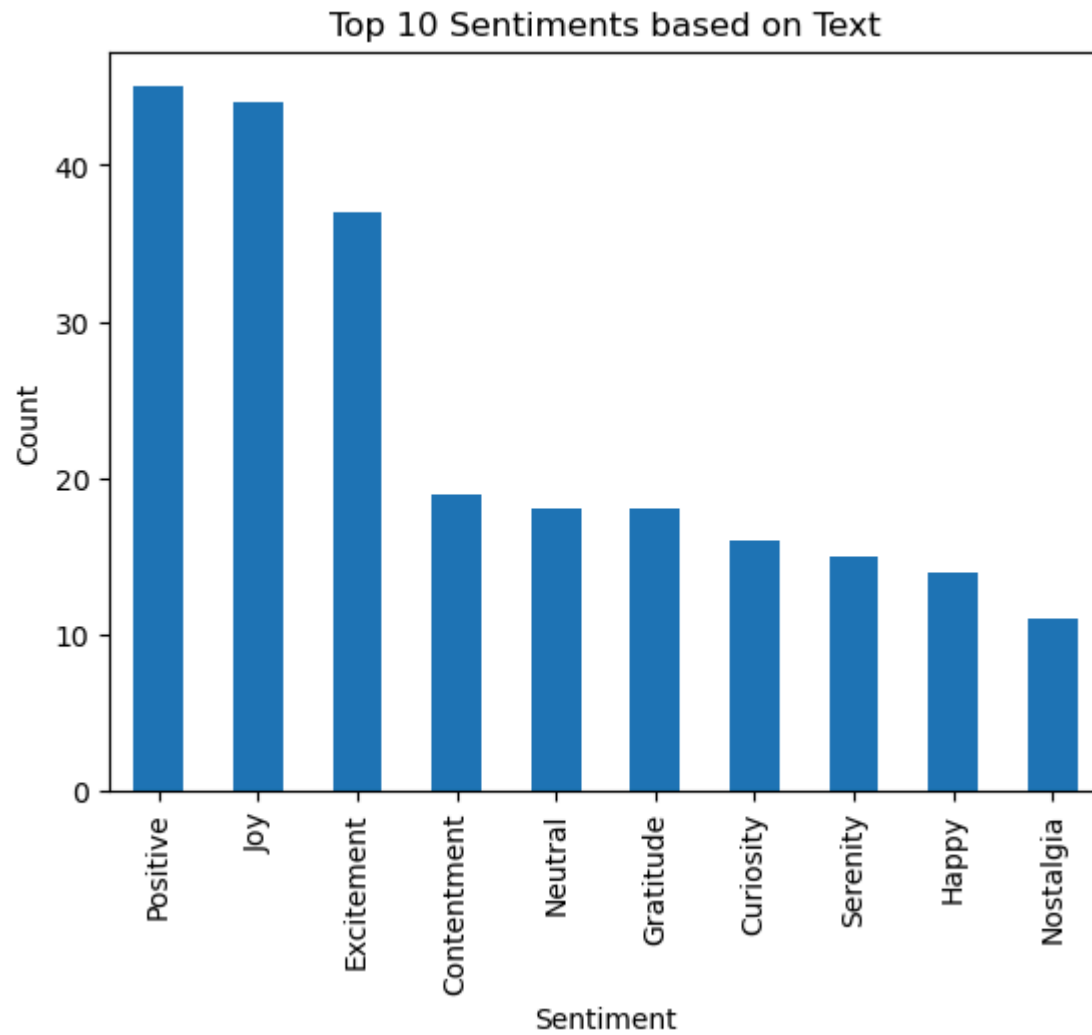
```
In [15]: 1 ##df = df.drop(columns=['Id', 'Hashtags', 'Day', 'Hour', 'Sentiment'])
```

```
In [16]: 1 df['Timestamp'] = pd.to_datetime(df['Timestamp'])
          2 df['Day'] = df['Timestamp'].dt.day
          3 df['Month'] = df['Timestamp'].dt.month
          4 df['Year'] = df['Timestamp'].dt.year
          5 df['Day_of_Week'] = df['Timestamp'].dt.day_name()
```

```
In [17]: 1 df['Text'] = df['Text'].str.strip()
          2 df['Sentiment'] = df['Sentiment'].str.strip()
          3 df['User'] = df['User'].str.strip()
          4 df['Platform'] = df['Platform'].str.strip()
          5 df['Hashtags'] = df['Hashtags'].str.strip()
          6 df['Country'] = df['Country'].str.strip()
```



```
In [18]: 1 df['Sentiment'].value_counts().nlargest(10).plot(kind='bar')
2 plt.title('Top 10 Sentiments based on Text')
3 plt.xlabel('Sentiment')
4 plt.ylabel('Count')
5 plt.show()
```

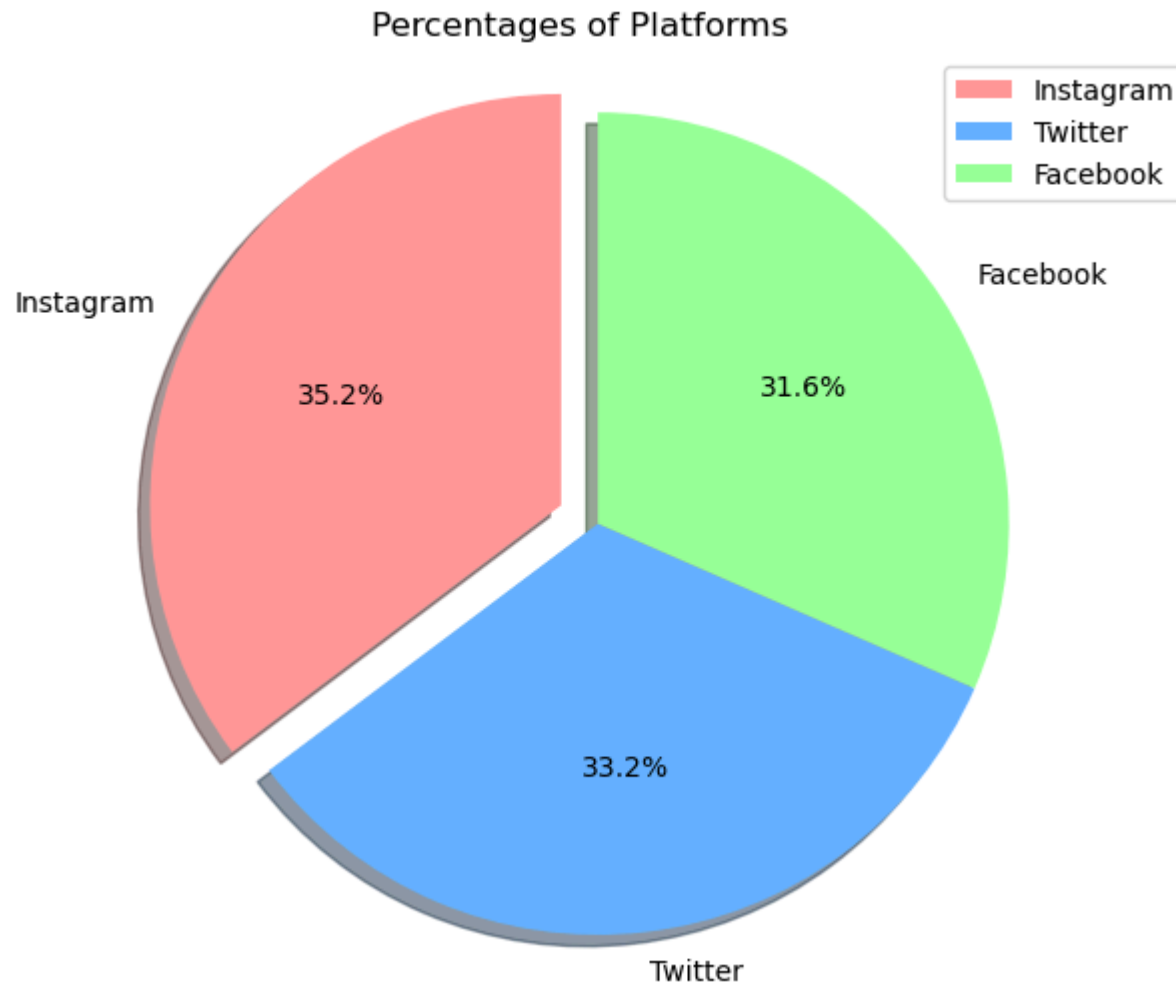


```
In [19]: 1 df['Platform'].value_counts()
```

```
Out[19]: Instagram    258  
Twitter      243  
Facebook     231  
Name: Platform, dtype: int64
```

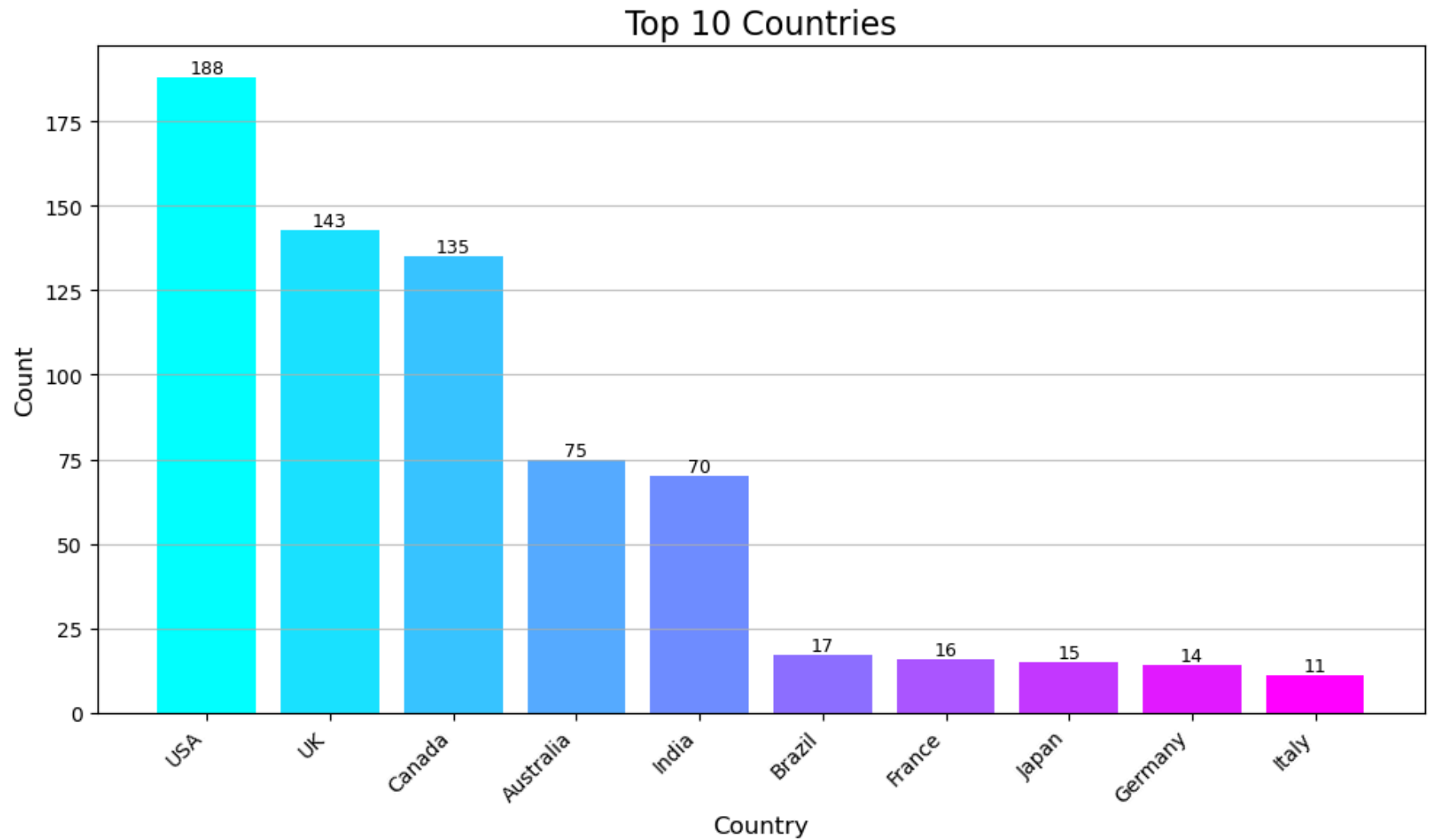
```
In [20]: 1 #df['Platform'].value_counts().plot(kind='pie', autopct='%1.1f%%')  
2 #plt.title('Percentages of Platforms')  
3 #plt.legend()  
4 #plt.show()
```

```
In [21]: 1 # Get the value counts of 'Platform' column
2 platform_counts = df['Platform'].value_counts()
3
4 # Define colors for the pie chart
5 colors = ['#ff9999', '#66b3ff', '#99ff99', '#ffcc99']
6
7 # Define explode to highlight the first platform
8 explode = [0.1 if i == 0 else 0 for i in range(len(platform_counts))]
9
10 # Create the pie chart
11 plt.figure(figsize=(8, 6))
12 plt.pie(platform_counts, labels=platform_counts.index, autopct='%1.1f%%', colors=colors, startangle=90, shadow=True)
13
14 # Equal aspect ratio ensures that pie is drawn as a circle
15 plt.axis('equal')
16 plt.title('Percentages of Platforms')
17 plt.legend(platform_counts.index, loc="best")
18 plt.show()
19
```



```
In [22]: 1 #df['Country'].value_counts().nlargest(10).plot(kind='bar')
          2 #plt.title('Top 10 Country')
          3 #plt.legend()
          4 #plt.show()
```

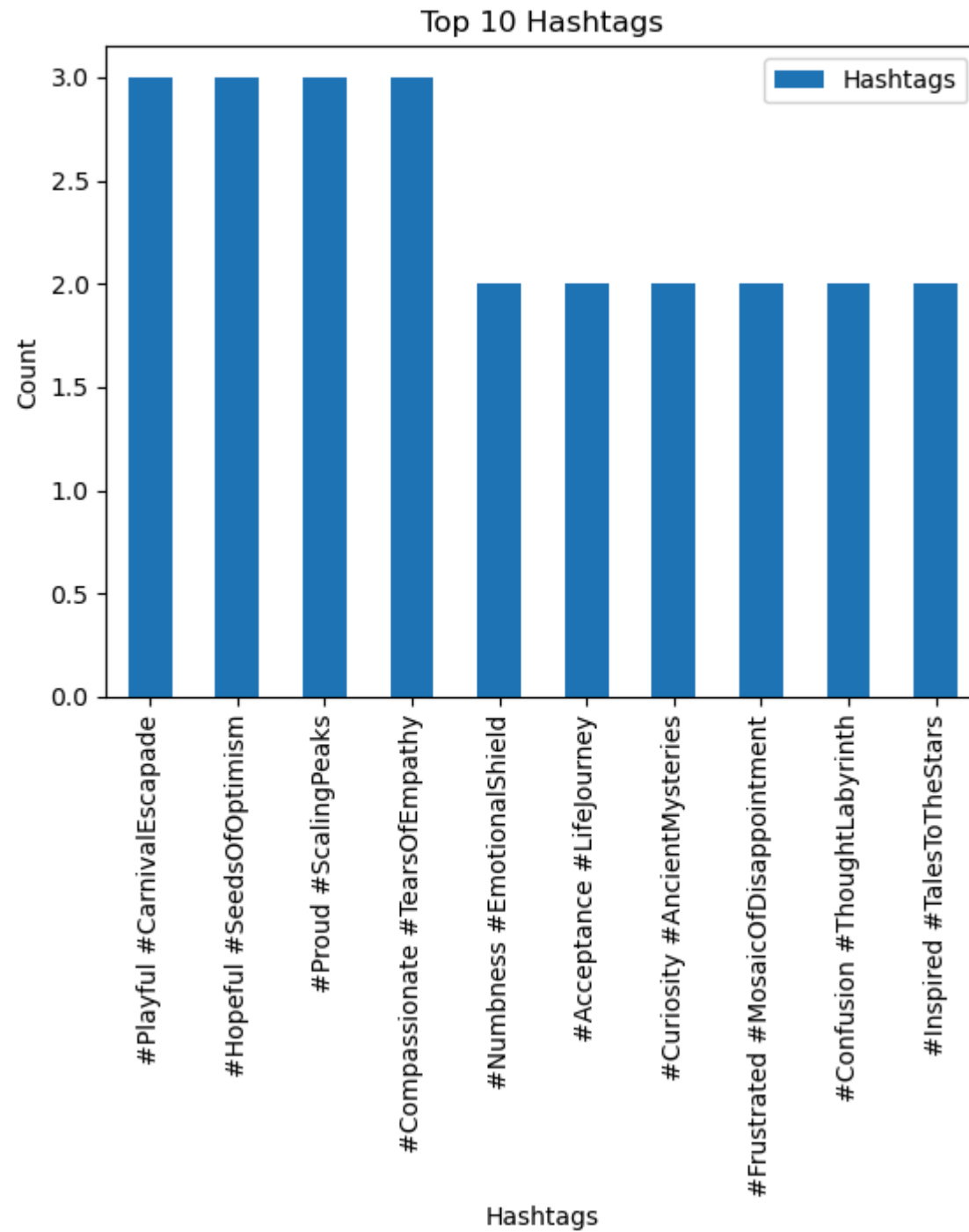
```
In [23]: 1 # Get the top 10 countries by value counts
2 top_countries = df['Country'].value_counts().nlargest(10)
3
4 # Define colors for the bar chart
5 colors = ['#1f77b4', '#ff7f0e', '#2ca02c', '#d62728', '#9467bd', '#8c564b', '#e377c2', '#7f7f7f', '#bcbd22', '#17
6
7 # Create a color gradient
8 gradient = np.linspace(0, 1, len(top_countries))
9 colors = plt.cm.cool(gradient)
10
11 # Create the bar chart with enhanced aesthetics
12 plt.figure(figsize=(10, 6))
13 bars = plt.bar(top_countries.index, top_countries.values, color=colors)
14
15 # Add title and Labels
16 plt.title('Top 10 Countries', fontsize=16)
17 plt.xlabel('Country', fontsize=12)
18 plt.ylabel('Count', fontsize=12)
19
20 # Customize ticks and grid lines
21 plt.xticks(rotation=45, ha='right', fontsize=10)
22 plt.yticks(fontsize=10)
23 plt.grid(axis='y', linestyle='--', alpha=0.7)
24
25 # Add data Labels
26 for bar in bars:
27     plt.text(bar.get_x() + bar.get_width() / 2, bar.get_height(),
28             round(bar.get_height(), 2), ha='center', va='bottom', fontsize=9)
29
30 # Show the plot
31 plt.tight_layout()
32 plt.show()
33
```



In []:

1

```
In [24]: 1 df['Hashtags'].value_counts().nlargest(10).plot(kind='bar')
          2 plt.title('Top 10 Hashtags')
          3 plt.xlabel('Hashtags')
          4 plt.ylabel('Count')
          5 plt.legend()
          6 plt.show()
```




```
In [25]: 1 # Describe data numerical
         2 df.describe()
```

Out[25]:

	Id	Retweets	Likes	Year	Month	Day	Hour
count	732.000000	732.000000	732.000000	732.000000	732.000000	732.000000	732.000000
mean	369.740437	21.508197	42.901639	2020.471311	6.122951	15.497268	15.521858
std	212.428936	7.061286	14.089848	2.802285	3.411763	8.474553	4.113414
min	0.000000	5.000000	10.000000	2010.000000	1.000000	1.000000	0.000000
25%	185.750000	17.750000	34.750000	2019.000000	3.000000	9.000000	13.000000
50%	370.500000	22.000000	43.000000	2021.000000	6.000000	15.000000	16.000000
75%	553.250000	25.000000	50.000000	2023.000000	9.000000	22.000000	19.000000
max	736.000000	40.000000	80.000000	2023.000000	12.000000	31.000000	23.000000

```
In [26]: 1 numerical_columns = df[['Day', 'Month', 'Year', 'Likes', 'Retweets']]
         2
         3 for col in numerical_columns.columns:
         4     print(f"Minimum {col}: {df[col].min()} | Maximum {col}: {df[col].max()}")
```

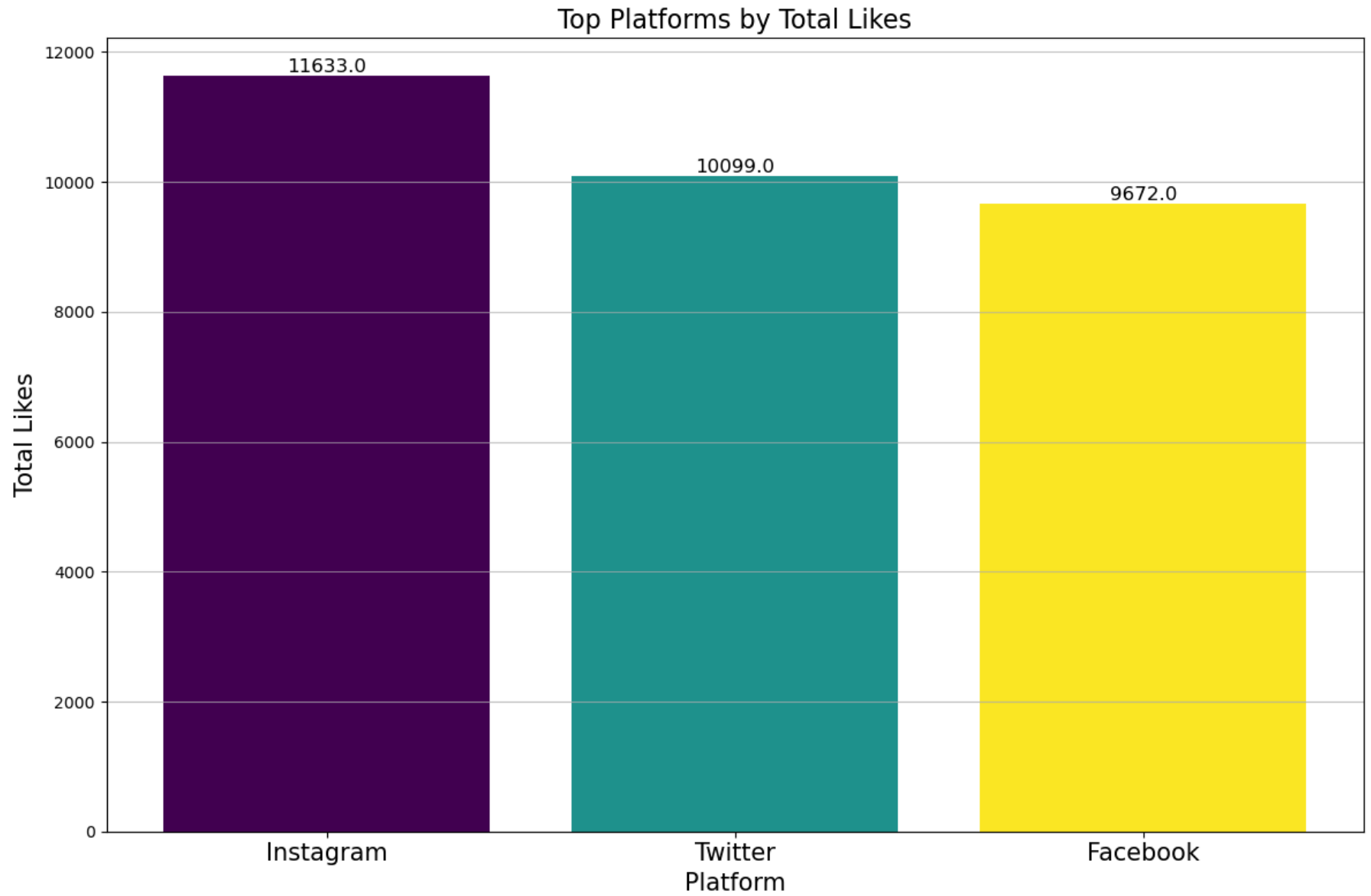
Minimum Day: 1 | Maximum Day: 31
 Minimum Month: 1 | Maximum Month: 12
 Minimum Year: 2010 | Maximum Year: 2023
 Minimum Likes: 10.0 | Maximum Likes: 80.0
 Minimum Retweets: 5.0 | Maximum Retweets: 40.0

Platform Top liked by Users ?

```
In [27]: 1 #top_likes_platform = df.groupby('Platform')['Likes'].sum().nlargest(10)
          2 #top_likes_platform.plot(kind='bar')
          3 #plt.title('Top Platforms by Total Likes')
          4 #plt.xlabel('Platform')
          5 #plt.ylabel('Total Likes')
          6 #plt.show()
```

In [28]:

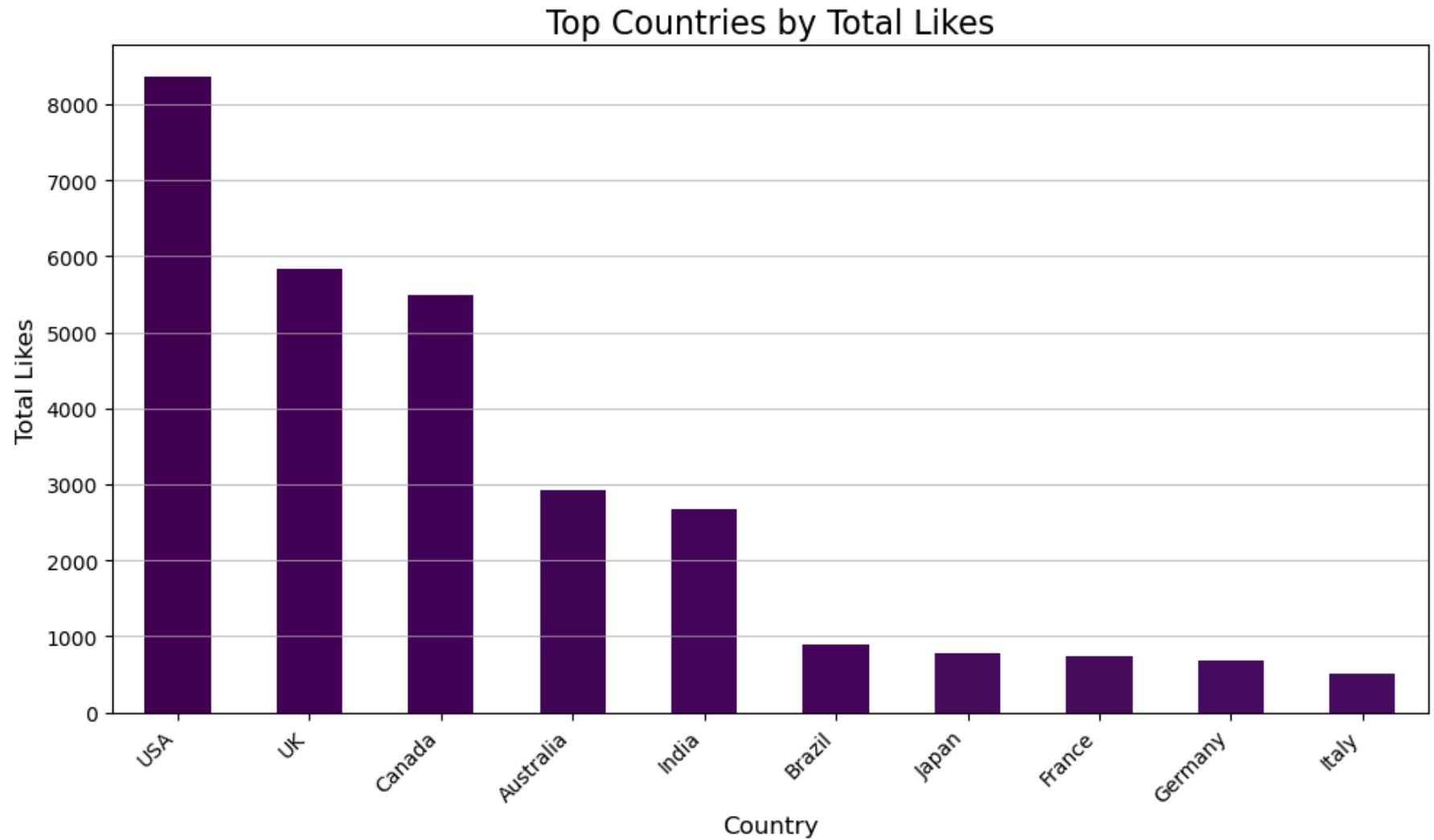
```
1
2 # Group by 'Platform' and calculate the sum of 'Likes', then select the top 10 platforms
3 top_likes_platform = df.groupby('Platform')['Likes'].sum().nlargest(10)
4
5 # Create a color gradient
6 colors = plt.cm.viridis(np.linspace(0, 1, len(top_likes_platform)))
7
8 # Create the bar chart with enhanced aesthetics
9 plt.figure(figsize=(12, 8))
10 bars = plt.bar(top_likes_platform.index, top_likes_platform.values, color=colors)
11
12 # Add title and labels
13 plt.title('Top Platforms by Total Likes', fontsize=16)
14 plt.xlabel('Platform', fontsize=15)
15 plt.ylabel('Total Likes', fontsize=15)
16
17 # Customize ticks and grid lines
18 plt.xticks(rotation=0, fontsize=15)
19 plt.yticks(fontsize=10)
20 plt.grid(axis='y', linestyle='--', alpha=0.7)
21
22 # Add data labels
23 for bar in bars:
24     plt.text(bar.get_x() + bar.get_width() / 2, bar.get_height(),
25             round(bar.get_height(), 2), ha='center', va='bottom', fontsize=12)
26
27 # Show the plot
28 plt.tight_layout()
29 plt.show()
30
```



Most Country Liked posts

```
In [29]: 1 #top_country_likes=df.groupby('Country')['Likes'].sum().nlargest(10)
          2 #top_country_likes.plot(kind='bar')
          3 #plt.title('Top country likes')
          4 #plt.xlabel('Country')
          5 #plt.ylabel('count')
          6 #plt.show()
```

```
In [30]: 1 import matplotlib.pyplot as plt
2
3 # Group by 'Country' and calculate the sum of 'Likes', then select the top 10 countries
4 top_country_likes = df.groupby('Country')['Likes'].sum().nlargest(10)
5
6 # Selecting colors from the 'viridis' colormap
7 colors = plt.cm.viridis(range(len(top_country_likes)))
8
9 # Create the bar chart with enhanced aesthetics
10 plt.figure(figsize=(10, 6))
11 top_country_likes.plot(kind='bar', color=colors)
12
13 # Add title and labels
14 plt.title('Top Countries by Total Likes', fontsize=16)
15 plt.xlabel('Country', fontsize=12)
16 plt.ylabel('Total Likes', fontsize=12)
17
18 # Customize ticks and grid lines
19 plt.xticks(rotation=45, ha='right', fontsize=10)
20 plt.yticks(fontsize=10)
21 plt.grid(axis='y', linestyle='--', alpha=0.7)
22
23 # Show the plot
24 plt.tight_layout()
25 plt.show()
26
```



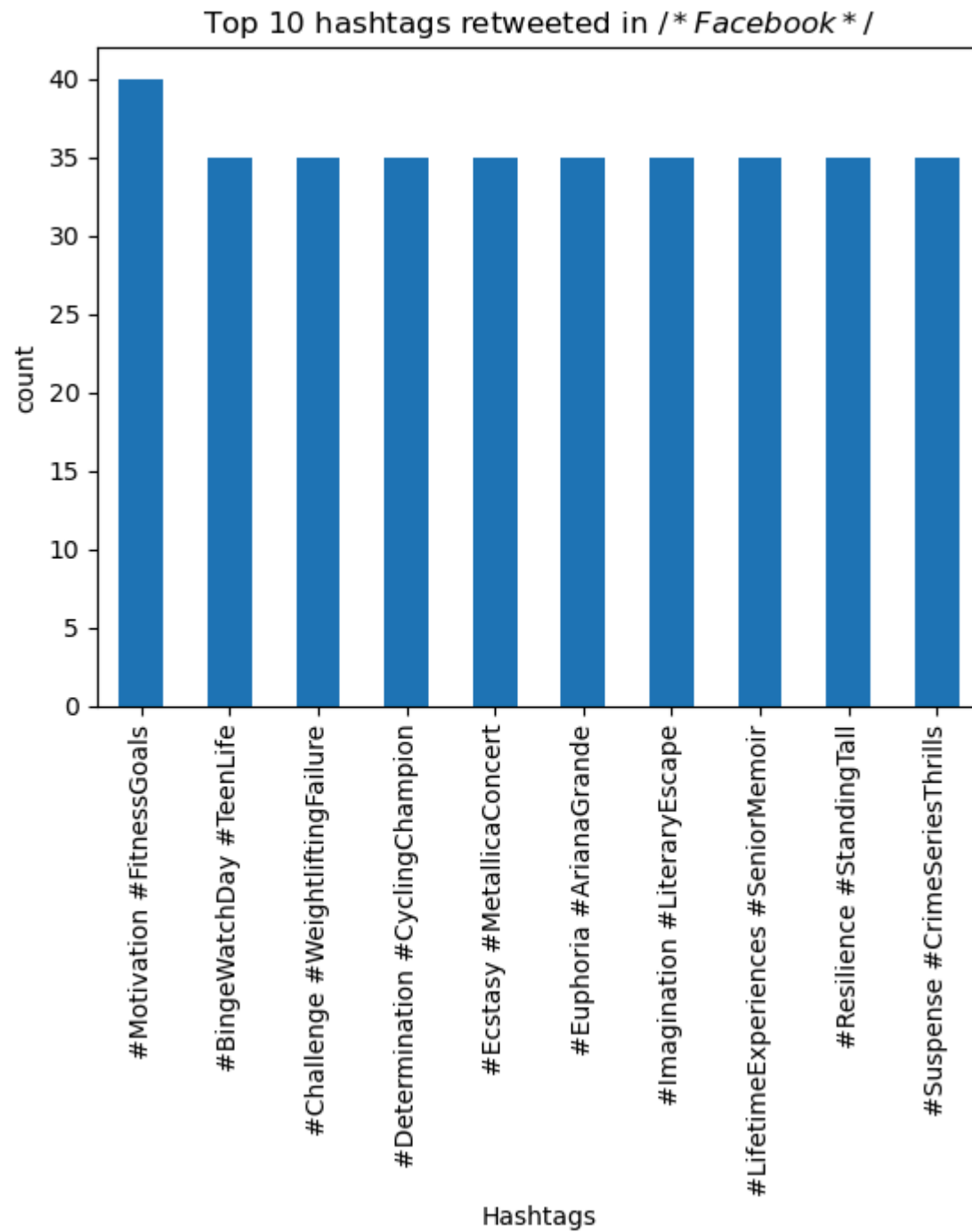
Segmenting users for each platform

```
In [31]: 1 Facebook=df[df['Platform']=='Facebook']  
        2 Twitter=df[df['Platform']=='Twitter']  
        3 Instagram=df[df['Platform']=='Instagram']
```


Facebook

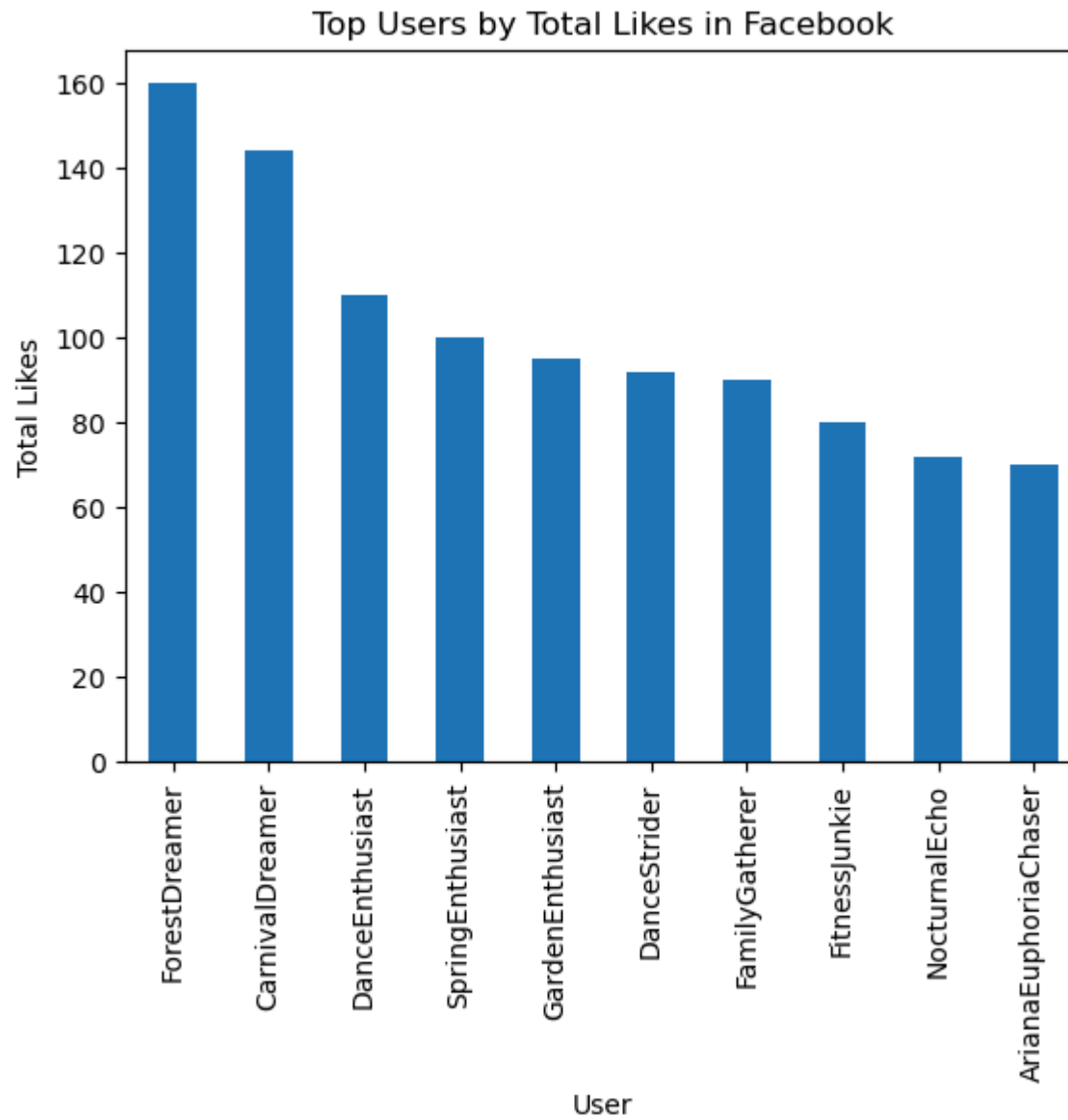
Top 10 hashtags retweeted

```
In [32]: 1 H_R_f=Facebook.groupby('Hashtags')['Retweets'].max().nlargest(10).sort_values(ascending=False)
          2 H_R_f.plot(kind='bar')
          3 plt.title('Top 10 hashtags retweeted in $/*Facebook*$/')
          4 plt.xlabel('Hashtags')
          5 plt.ylabel('count')
          6 plt.show()
```



Which User liked mostly

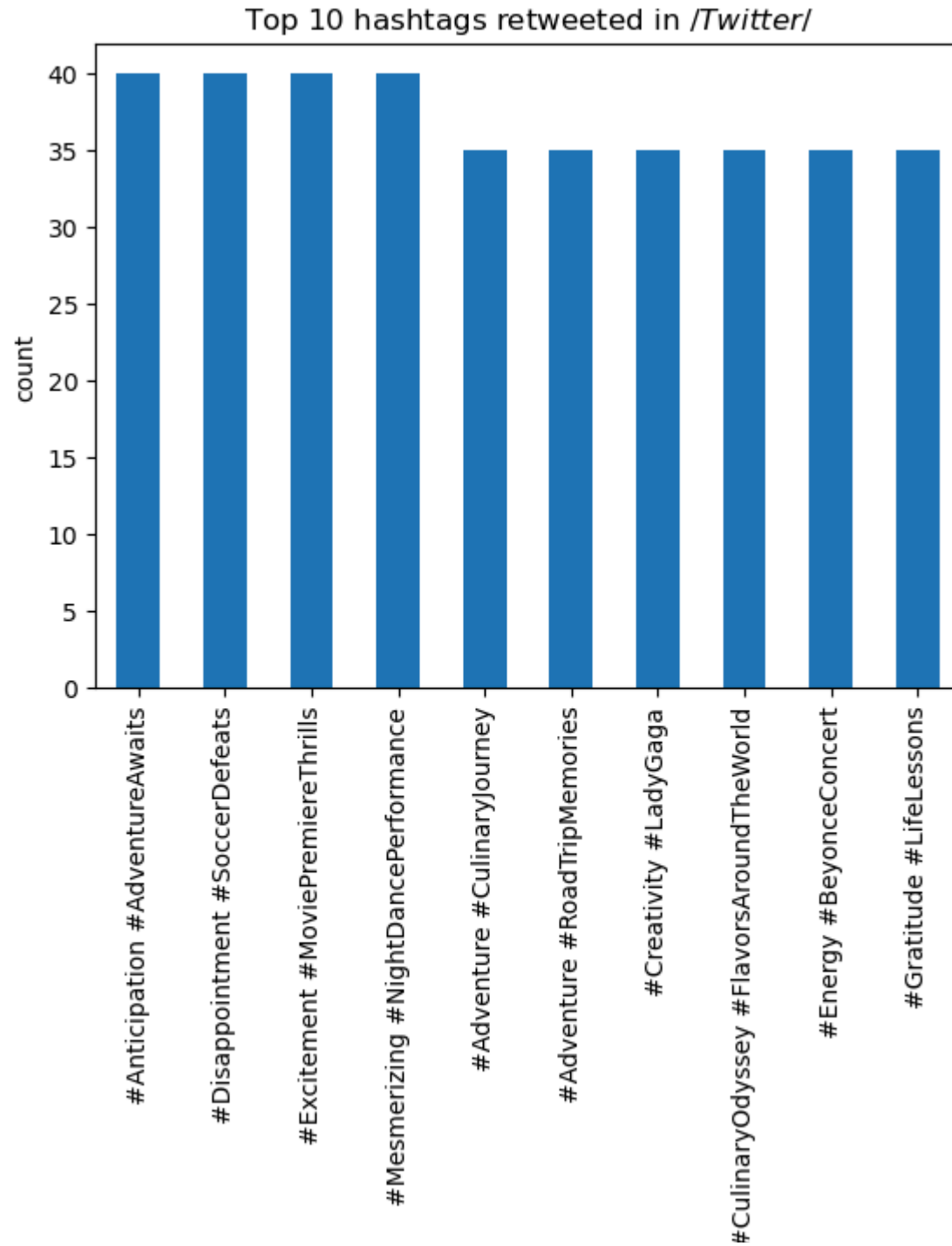
```
In [33]: 1 top_likes_platform_F = Facebook.groupby('User')['Likes'].sum().nlargest(10)
          2 top_likes_platform_F.plot(kind='bar')
          3 plt.title('Top Users by Total Likes in Facebook')
          4 plt.xlabel('User')
          5 plt.ylabel('Total Likes')
          6 plt.show()
```



Twitter

Top 10 hashtags retweeted

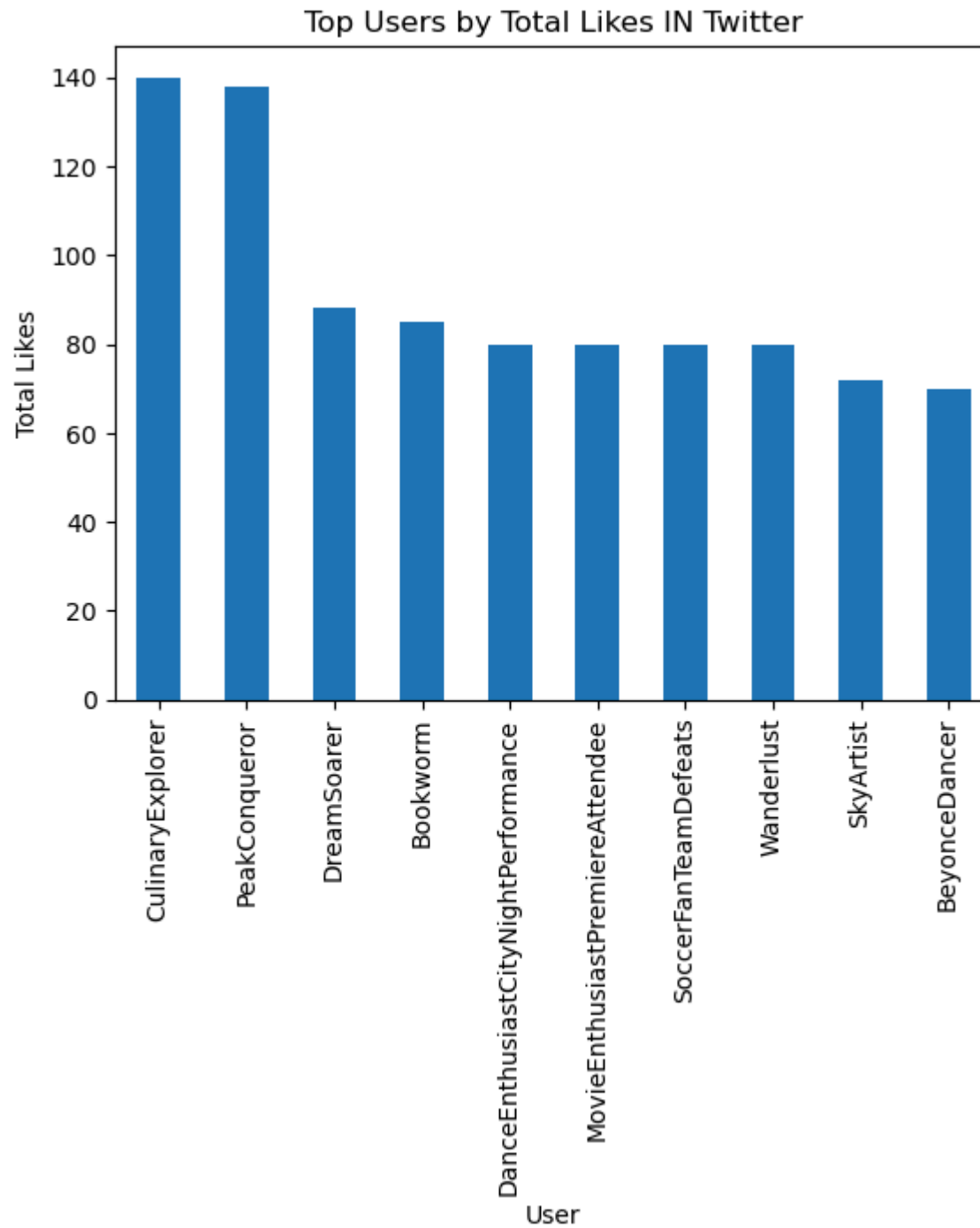
```
In [34]: 1 H_R_t=Twitter.groupby('Hashtags')['Retweets'].max().nlargest(10).sort_values(ascending=False)
          2 H_R_t.plot(kind='bar')
          3 plt.title('Top 10 hashtags retweeted in $/ Twitter $/')
          4 plt.xlabel('Hashtags')
          5 plt.ylabel('count')
          6 plt.show()
```

Hashtags

Who User liked mostly

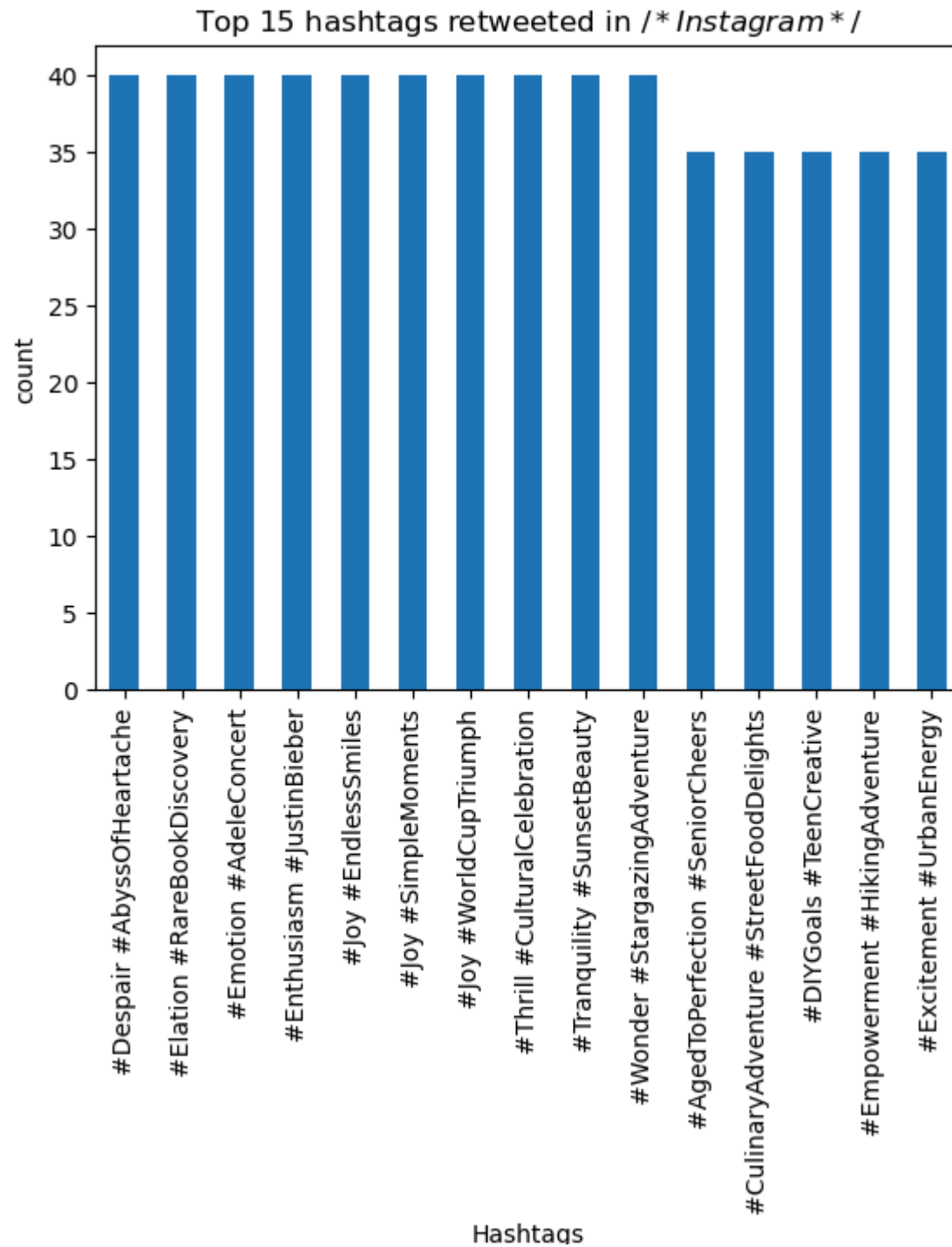
```
In [35]: 1 top_likes_platform_t = Twitter.groupby('User')['Likes'].sum().nlargest(10)
          2 top_likes_platform_t.plot(kind='bar')
          3 plt.title('Top Users by Total Likes IN Twitter')
          4 plt.xlabel('User')
          5 plt.ylabel('Total Likes')
          6 plt.show()
```



Instagram

Top 10 hashtags retweeted

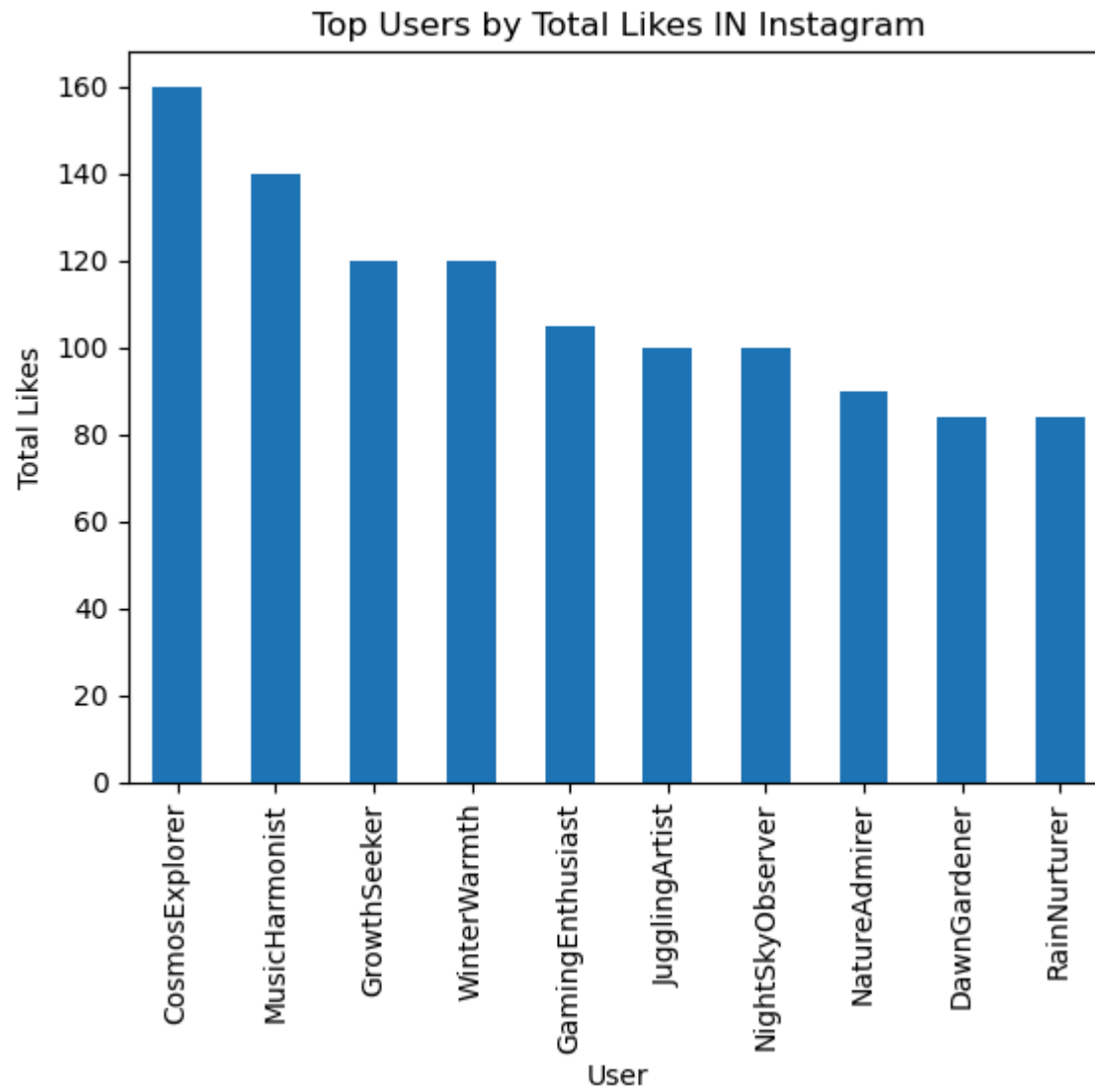
```
In [36]: 1 H_R_i=Instagram.groupby('Hashtags')['Retweets'].max().nlargest(15).sort_values(ascending=False)
          2 H_R_i.plot(kind='bar')
          3 plt.title('Top 15 hashtags retweeted in $/ *Instagram* $/')
          4 plt.xlabel('Hashtags')
          5 plt.ylabel('count')
          6 plt.show()
```



Which User liked mostly

```
In [37]: 1 top_likes_platform_i = Instagram.groupby('User')['Likes'].sum().nlargest(10)
          2 top_likes_platform_i.plot(kind='bar')
          3 plt.title('Top Users by Total Likes IN Instagram')
          4 plt.xlabel('User')
          5 plt.ylabel('Total Likes')
          6 plt.show()
```



Month

```
In [38]: 1 month_mapping = {  
2         1: 'January',  
3         2: 'February',  
4         3: 'March',  
5         4: 'April',  
6         5: 'May',  
7         6: 'June',  
8         7: 'July',  
9         8: 'August',  
10        9: 'September',  
11       10: 'October',  
12       11: 'November',  
13       12: 'December'  
14     }  
15  
16 df['Month'] = df['Month'].map(month_mapping)  
17  
18 df['Month'] = df['Month'].astype('object')
```

Text

```
In [39]: 1 stemmer = PorterStemmer()
2 stop_words = set(stopwords.words('english'))
3
4 def clean(text):
5     text = str(text).lower()
6     text = re.sub('\[.*?\]', '', text)
7     text = re.sub('https?://\S+|www\.\S+', '', text)
8     text = re.sub(r'\s+', ' ', text.strip())
9     text = re.sub('<.*?>+', '', text)
10    text = re.sub('%s' % re.escape(string.punctuation), '', text)
11    text = re.sub('\n', '', text)
12    text = re.sub('\w*\d\w*', '', text)
13    text = re.sub(r'^\x00-\x7F+', '', text)
14    text = " ".join(text.split())
15    tokens = word_tokenize(text)
16
17    cleaned_tokens = [stemmer.stem(token) for token in tokens if token.lower() not in stop_words]
18
19    cleaned_text = ' '.join(cleaned_tokens)
20
21    return cleaned_text
22
23 df["Clean_Text"] = df["Text"].apply(clean)
```

Sentiment Analysis

```
In [40]: 1 df1 = df.copy()
```

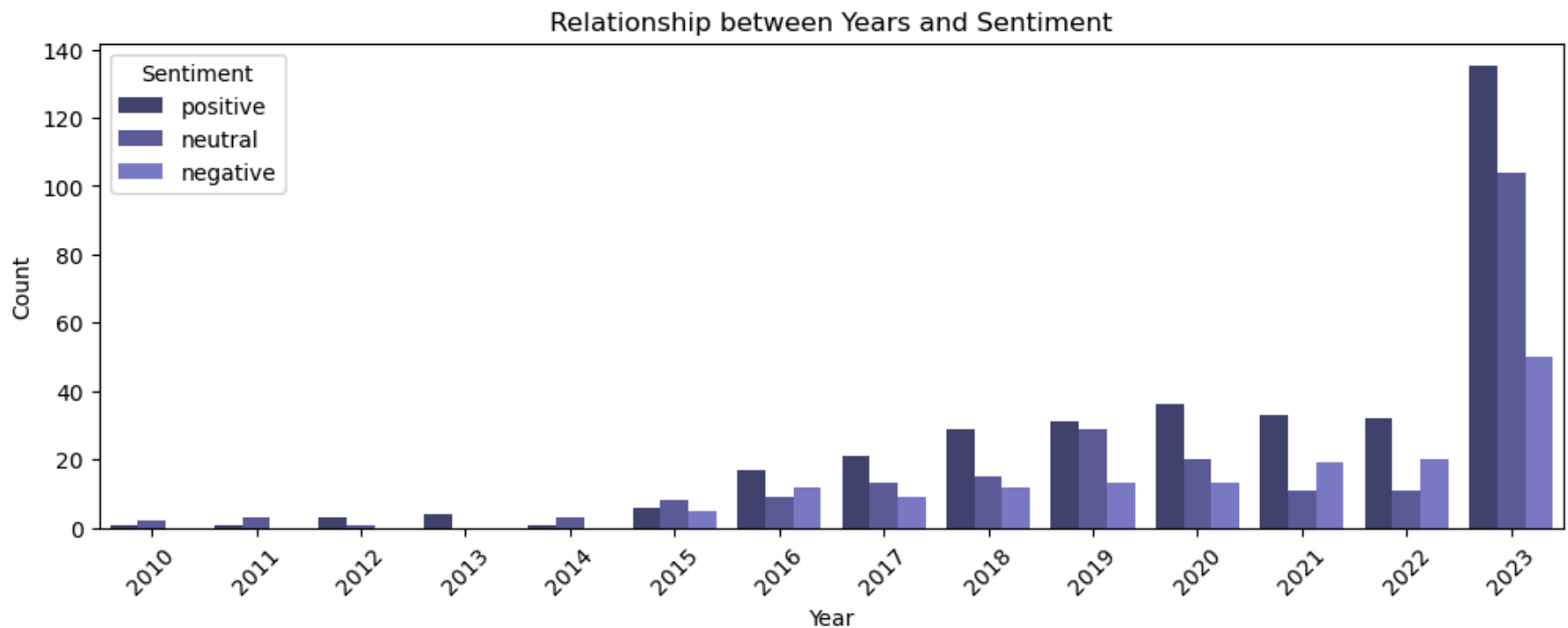
```
In [41]: 1 analyzer = SentimentIntensityAnalyzer()
2
3 df1['Vader_Score'] = df1['Clean_Text'].apply(lambda text: analyzer.polarity_scores(text)['compound'])
4
5 df1['Sentiment'] = df1['Vader_Score'].apply(lambda score: 'positive' if score >= 0.05 else ('negative' if score <
6
7 print(df1[['Clean_Text', 'Vader_Score', 'Sentiment']].head())
```

	Clean_Text	Vader_Score	Sentiment
0	enjoy beauti day park	0.4939	positive
1	traffic terribl morn	0.0000	neutral
2	finish amaz workout	0.0000	neutral
3	excit upcom weekend getaway	0.0000	neutral
4	tri new recip dinner tonight	0.0000	neutral

Relationship B/w Years & Sentiments

```
In [42]: 1 #df['Year'] = df['Timestamp'].dt.year
```

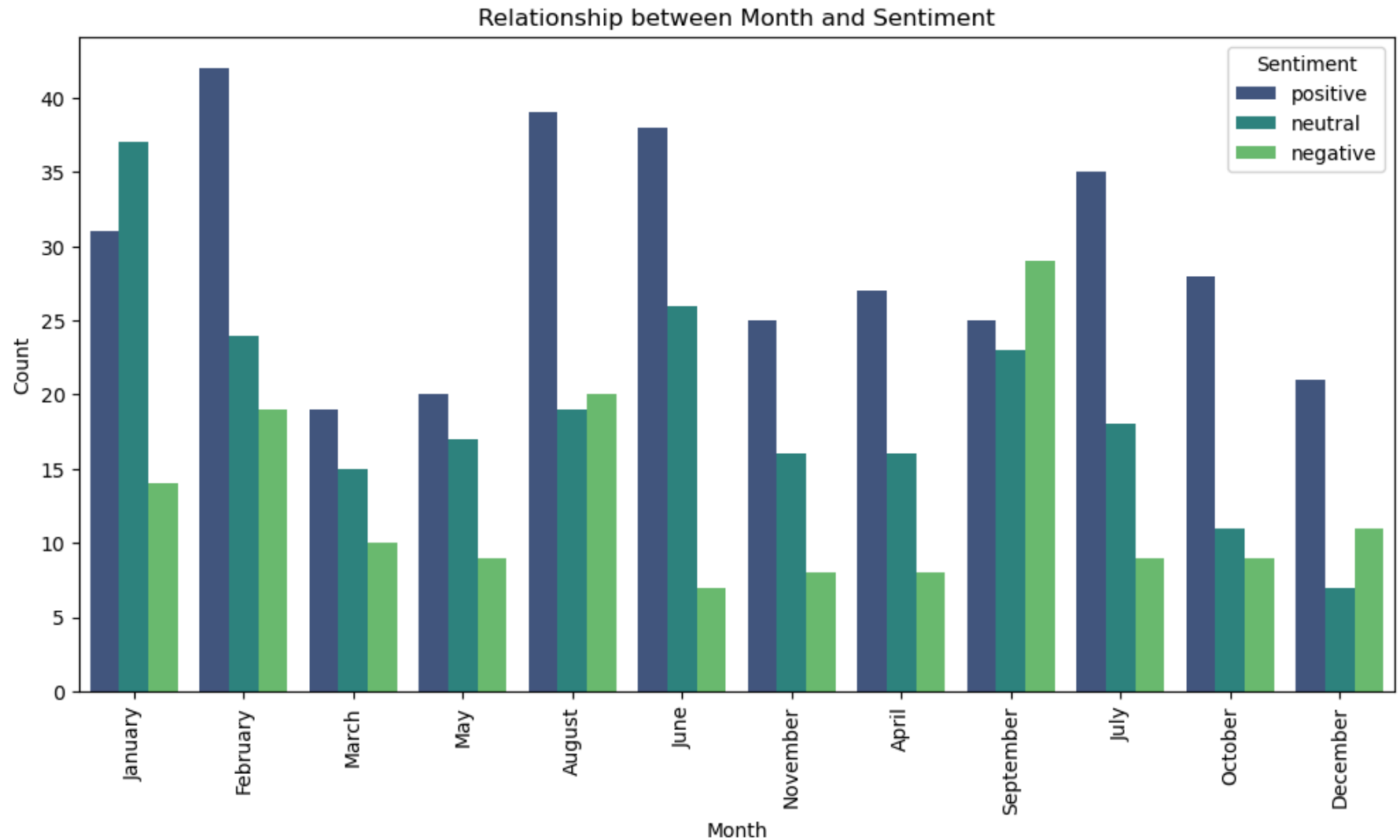
```
In [43]: 1 plt.figure(figsize=(12, 4))
2 sns.countplot(x='Year', hue='Sentiment', data=df1, palette='tab20b')
3 plt.title('Relationship between Years and Sentiment')
4 plt.xlabel('Year')
5 plt.ylabel('Count')
6 plt.xticks(rotation=45)
7 plt.show()
```



Relationship B/w Months & Sentiments

```
In [44]: 1 #df['Month'] = df['Timestamp'].dt.month
```

```
In [45]: 1 plt.figure(figsize=(12, 6))
2 sns.countplot(x='Month', hue='Sentiment', data=df1, palette='viridis')
3 plt.title('Relationship between Month and Sentiment')
4 plt.xlabel('Month')
5 plt.ylabel('Count')
6 plt.xticks(rotation=90)
7 plt.show()
```

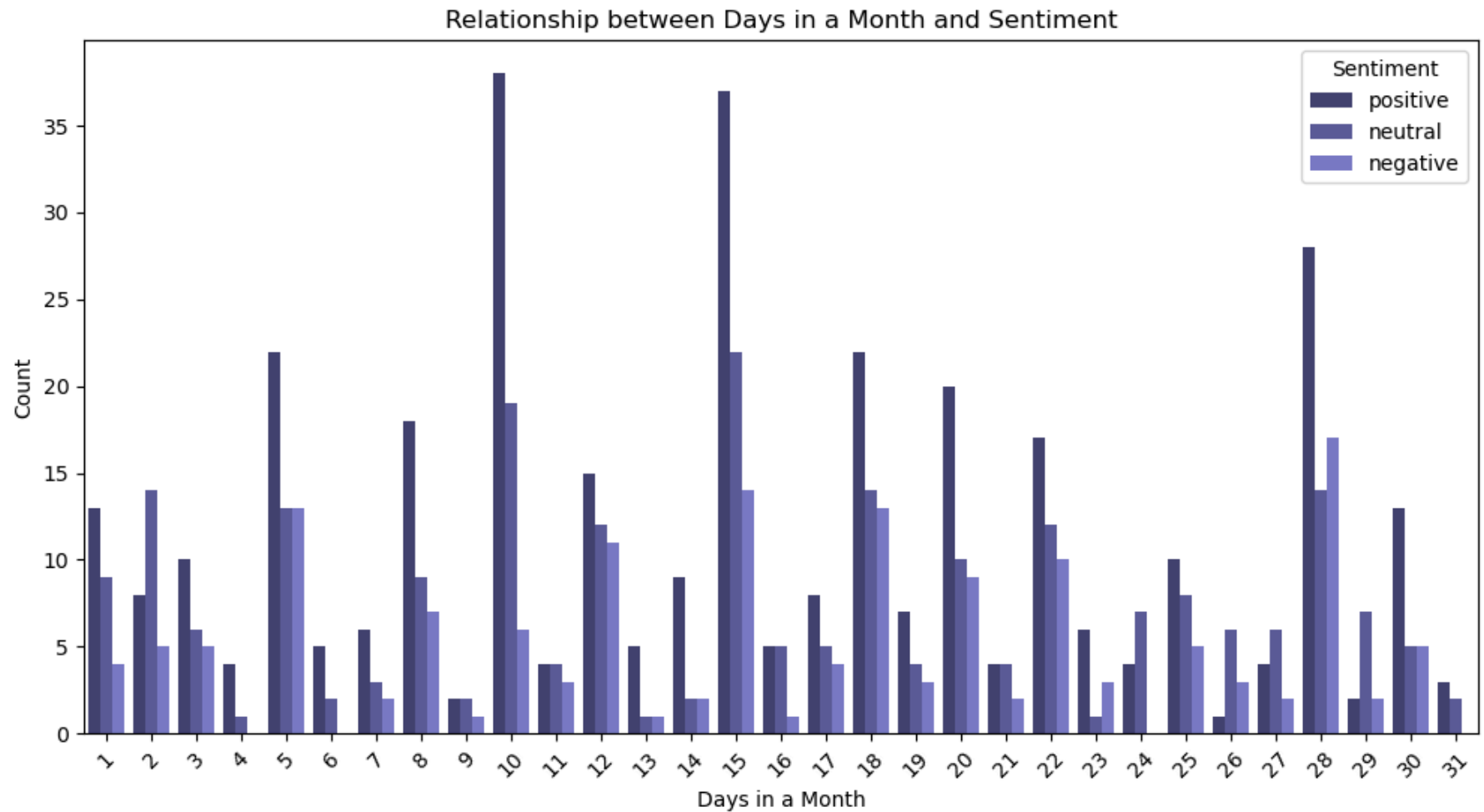


Relationship B/w Days & Sentiments

```
In [46]: 1 #df['Day'] = df['Timestamp'].dt.day
```



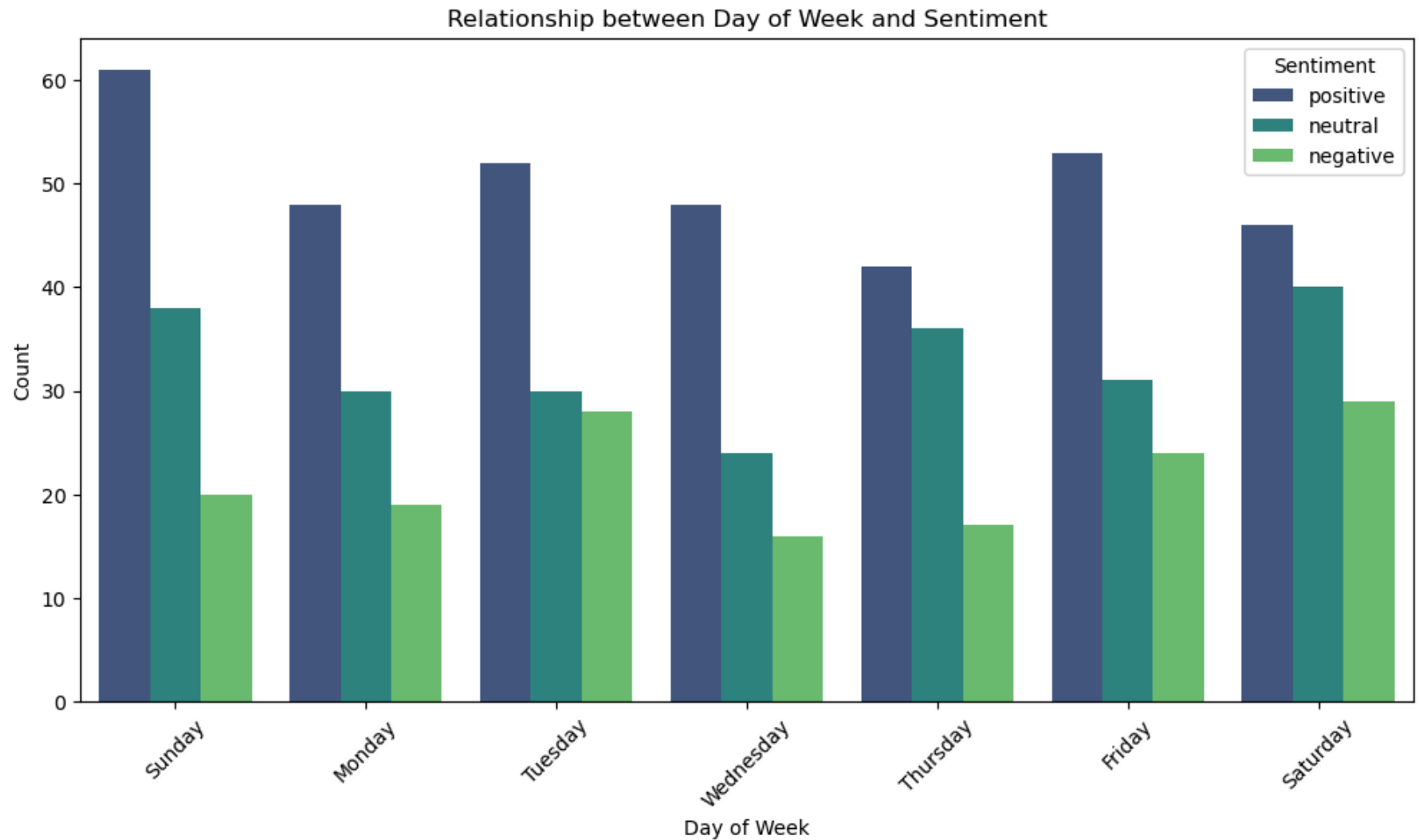
```
In [47]: 1 plt.figure(figsize=(12, 6))
2 sns.countplot(x='Day', hue='Sentiment', data=df1, palette='tab20b')
3 plt.title('Relationship between Days in a Month and Sentiment')
4 plt.xlabel('Days in a Month')
5 plt.ylabel('Count')
6 plt.xticks(rotation=45)
7 plt.show()
```



Relationship B/w Day of Week & Sentiments

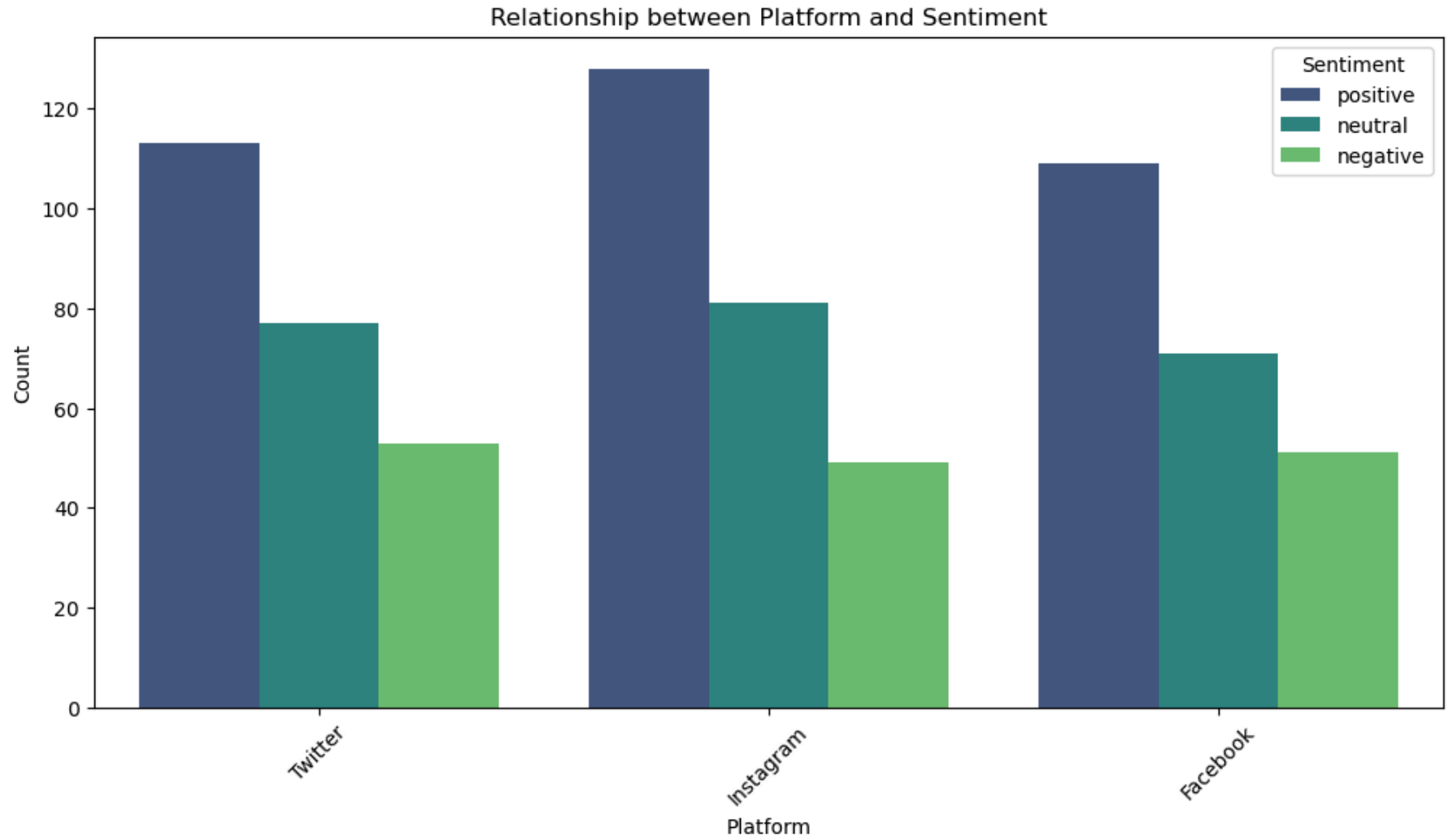
```
In [48]: 1 #df['Day_of_Week'] = df['Timestamp'].dt.day_name()
```

```
In [49]: 1 plt.figure(figsize=(12, 6))
2 sns.countplot(x='Day_of_Week', hue='Sentiment', data=df1, palette='viridis')
3 plt.title('Relationship between Day of Week and Sentiment')
4 plt.xlabel('Day of Week')
5 plt.ylabel('Count')
6 plt.xticks(rotation=45)
7 plt.show()
```



Relationship B/w Platform & Sentiments

```
In [50]: 1 plt.figure(figsize=(12, 6))
          2 sns.countplot(x='Platform', hue='Sentiment', data=df1, palette='viridis')
          3 plt.title('Relationship between Platform and Sentiment')
          4 plt.xlabel('Platform')
          5 plt.ylabel('Count')
          6 plt.xticks(rotation=45)
          7 plt.show()
```



Common Words

In [57]:

```
1 df1['temp_list'] = df1['Clean_Text'].apply(lambda x: str(x).split())
2 top_words = Counter([item for sublist in df1['temp_list'] for item in sublist])
3 top_words_df = pd.DataFrame(top_words.most_common(20), columns=['Common_words', 'count'])
4
5 top_words_df.style.background_gradient(cmap='Blues')
```

Out[57]:

	Common_words	count
0	new	43
1	life	37
2	challeng	34
3	joy	31
4	danc	30
5	day	29
6	feel	29
7	like	28
8	dream	28
9	moment	27
10	emot	27
11	friend	26
12	heart	26
13	explor	25
14	echo	25
15	beauti	24
16	laughter	24
17	embrac	24
18	night	23
19	hope	23

Word Count

```
In [58]: 1 Positive_sent = df1[df1['Sentiment'] == 'positive']  
2 Negative_sent = df1[df1['Sentiment'] == 'negative']  
3 Neutral_sent = df1[df1['Sentiment'] == 'neutral']
```

Positive Common Words

```
In [59]: 1 top = Counter([item for sublist in df1[df1['Sentiment'] == 'positive']['temp_list'] for item in sublist])  
2 temp_positive = pd.DataFrame(top.most_common(10), columns=['Common_words', 'count'])  
3 temp_positive.style.background_gradient(cmap='Blues')
```

Out[59]:

	Common_words	count
0	joy	30
1	friend	24
2	laughter	24
3	new	21
4	challeng	20
5	life	20
6	hope	20
7	dream	20
8	embrac	19
9	like	19

Neutral Common Words

```
In [60]: 1 top = Counter([item for sublist in df1[df1['Sentiment'] == 'neutral']['temp_list'] for item in sublist])
          2 temp_positive = pd.DataFrame(top.most_common(10), columns=['Common_words', 'count'])
          3 temp_positive.style.background_gradient(cmap='Greens')
```

Out[60]:

	Common_words	count
0	new	22
1	explor	13
2	excit	12
3	life	12
4	beauti	12
5	night	12
6	danc	12
7	attend	11
8	seren	11
9	feel	10

Negative Common Words

```
In [61]: 1 top = Counter([item for sublist in df1[df1['Sentiment'] == 'negative']['temp_list'] for item in sublist])
          2 temp_positive = pd.DataFrame(top.most_common(10), columns=['Common_words', 'count'])
          3 temp_positive.style.background_gradient(cmap='Reds')
```

Out[61]:

	Common_words	count
0	despair	14
1	lost	14
2	emot	13
3	feel	11
4	bitter	10
5	storm	10
6	day	9
7	like	9
8	grief	8
9	heart	8

Data Preparation

```
In [62]: 1 df2 = df1.copy()
```

```
In [68]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.linear_model import PassiveAggressiveClassifier
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import accuracy_score, classification_report
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.svm import SVC
8 from sklearn.naive_bayes import MultinomialNB
9 from sklearn.model_selection import RandomizedSearchCV
10 from sklearn.metrics import confusion_matrix
```

Splitting the Data

```
In [69]: 1 X = df2['Clean_Text'].values
2 y = df2['Sentiment'].values
```

```
In [70]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Modeling

```
In [71]: 1 vectorizer = TfidfVectorizer(max_features=5000)
2 X_train_tfidf = vectorizer.fit_transform(X_train)
3 X_test_tfidf = vectorizer.transform(X_test)
```

SVM Classifier

```
In [72]: 1 svm_classifier = SVC(random_state=42)
2 svm_classifier.fit(X_train_tfidf, y_train)
```

```
Out[72]: SVC
SVC(random_state=42)
```

```
In [73]: 1 y_pred_svm = svm_classifier.predict(X_test_tfidf)
          2 accuracy_svm = accuracy_score(y_test, y_pred_svm)
          3 classification_rep_svm = classification_report(y_test, y_pred_svm)
```

```
In [74]: 1 print("Support Vector Machine Results:")
          2 print(f"Accuracy: {accuracy_svm}")
          3 print("Classification Report:\n", classification_rep_svm)
```

Support Vector Machine Results:

Accuracy: 0.5986394557823129

Classification Report:

	precision	recall	f1-score	support
negative	1.00	0.38	0.55	32
neutral	0.94	0.31	0.47	55
positive	0.50	0.98	0.67	60
accuracy			0.60	147
macro avg	0.82	0.56	0.56	147
weighted avg	0.78	0.60	0.57	147

Random Forest Classifier

```
In [75]: 1 random_forest_classifier = RandomForestClassifier(random_state=42)
          2 random_forest_classifier.fit(X_train_tfidf, y_train)
```

```
Out[75]: ▼ RandomForestClassifier
          RandomForestClassifier(random_state=42)
```

```
In [76]: 1 y_pred_rf = random_forest_classifier.predict(X_test_tfidf)
          2 accuracy_rf = accuracy_score(y_test, y_pred_rf)
          3 classification_rep_rf = classification_report(y_test, y_pred_rf)
```

```
In [77]: 1 print("\nRandom Forest Results:")
          2 print(f"Accuracy: {accuracy_rf}")
          3 print("Classification Report:\n", classification_rep_rf)
```

Random Forest Results:

Accuracy: 0.6530612244897959

Classification Report:

	precision	recall	f1-score	support
negative	0.86	0.59	0.70	32
neutral	0.77	0.44	0.56	55
positive	0.56	0.88	0.69	60
accuracy			0.65	147
macro avg	0.73	0.64	0.65	147
weighted avg	0.71	0.65	0.64	147

Multinomial Naive Bayes

```
In [78]: 1 nb_classifier = MultinomialNB()
          2 nb_classifier.fit(X_train_tfidf, y_train)
```

```
Out[78]: ▾ MultinomialNB
          MultinomialNB()
```

```
In [79]: 1 y_pred_nb = nb_classifier.predict(X_test_tfidf)
          2 accuracy_nb = accuracy_score(y_test, y_pred_nb)
          3 classification_rep_nb = classification_report(y_test, y_pred_nb)
```

```
In [80]: 1 print("\nMultinomial Naive Bayes Results:")
2 print(f"Accuracy: {accuracy_nb}")
3 print("Classification Report:\n", classification_rep_nb)
```

Multinomial Naive Bayes Results:

Accuracy: 0.6190476190476191

Classification Report:

	precision	recall	f1-score	support
negative	1.00	0.38	0.55	32
neutral	0.90	0.35	0.50	55
positive	0.53	1.00	0.69	60
accuracy			0.62	147
macro avg	0.81	0.57	0.58	147
weighted avg	0.77	0.62	0.59	147

Best Modeling : Random Forest Classifier

```
In [81]: 1 param_dist = {
2         'n_estimators': [100, 200, 300],
3         'max_depth': [None, 10, 20, 30]
4     }
5
```

```
In [82]: 1 RF_classifier = RandomForestClassifier(random_state=42)
          2
          3 randomized_search = RandomizedSearchCV(RF_classifier, param_distributions=param_dist, n_iter=10, cv=5, scoring='a
          4 randomized_search.fit(X_train_tfidf, y_train)
```

```
Out[82]: RandomizedSearchCV
          estimator: RandomForestClassifier
              RandomForestClassifier
```

```
In [83]: 1 best_params_randomized = randomized_search.best_params_
          2 best_params_randomized
```

```
Out[83]: {'n_estimators': 100, 'max_depth': None}
```

```
In [84]: 1 best_RF_classifier_randomized = RandomForestClassifier(random_state=42, **best_params_randomized)
          2 best_RF_classifier_randomized.fit(X_train_tfidf, y_train)
```

```
Out[84]: RandomForestClassifier
          RandomForestClassifier(random_state=42)
```

```
In [85]: 1 y_pred_best_RF_randomized = best_RF_classifier_randomized.predict(X_test_tfidf)
```

```
In [86]: 1 accuracy_best_RF_randomized = accuracy_score(y_test, y_pred_best_RF_randomized)
          2 classification_rep_best_RF_randomized = classification_report(y_test, y_pred_best_RF_randomized)
          3 conf_matrix_test = confusion_matrix(y_test, y_pred_best_RF_randomized)
```

```
In [87]: 1 print("Best RandomForestClassifier Model (RandomizedSearchCV):")
2 print(f"Best Hyperparameters: {best_params_randomized}")
3 print(f"Accuracy: {accuracy_best_RF_randomized}")
4 print("Classification Report:\n", classification_rep_best_RF_randomized)
```

Best RandomForestClassifier Model (RandomizedSearchCV):

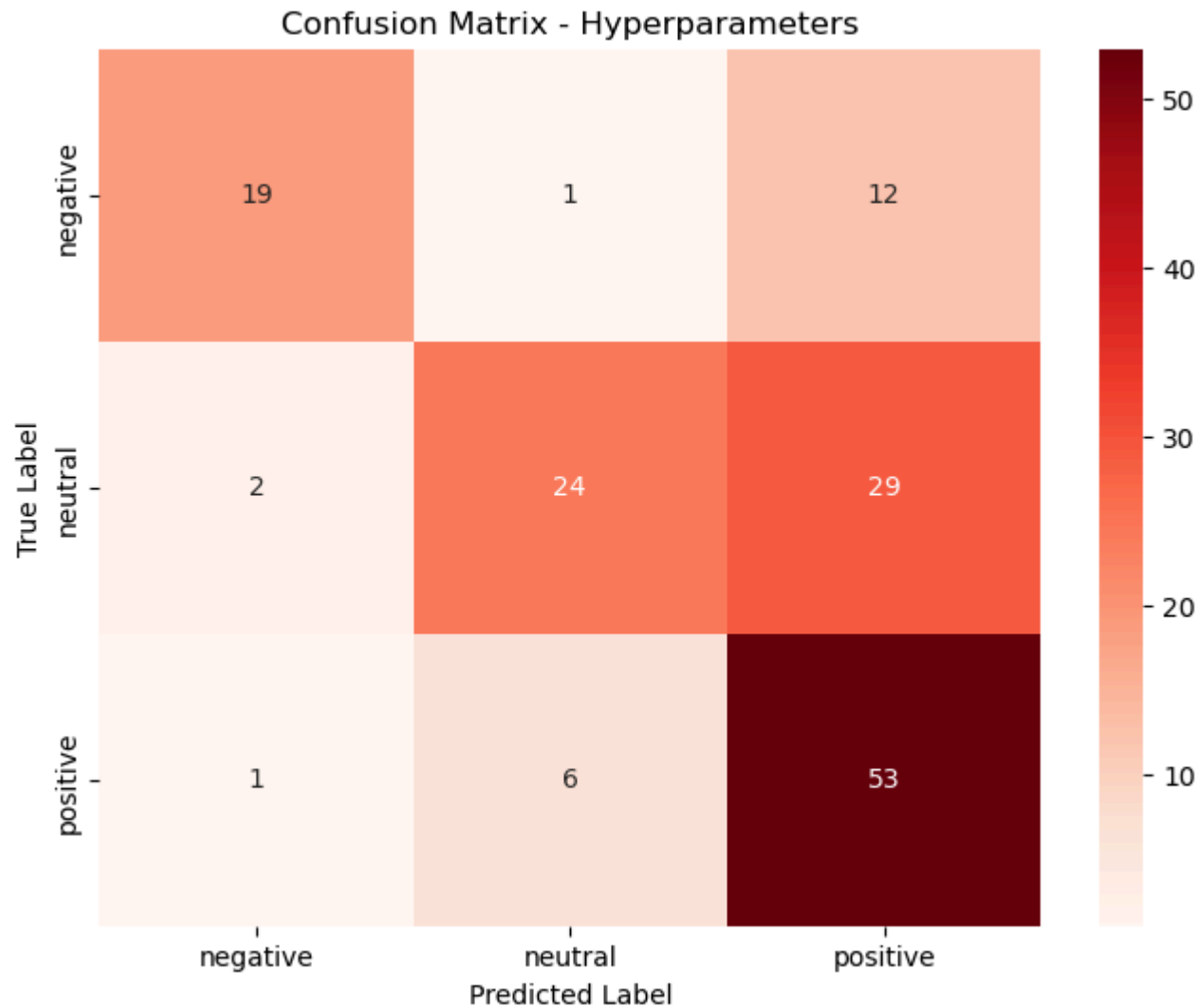
Best Hyperparameters: {'n_estimators': 100, 'max_depth': None}

Accuracy: 0.6530612244897959

Classification Report:

	precision	recall	f1-score	support
negative	0.86	0.59	0.70	32
neutral	0.77	0.44	0.56	55
positive	0.56	0.88	0.69	60
accuracy			0.65	147
macro avg	0.73	0.64	0.65	147
weighted avg	0.71	0.65	0.64	147


```
In [88]: 1 plt.figure(figsize=(8, 6))
2 sns.heatmap(conf_matrix_test, annot=True, fmt='d', cmap='Reds', xticklabels=['negative', 'neutral', 'positive'],
3 plt.title('Confusion Matrix - Hyperparameters')
4 plt.xlabel('Predicted Label')
5 plt.ylabel('True Label')
6 plt.show()
```



In []:

1