# Task 2 - Let's read the CSV file with some more analytical and statistical analysis like head(), describe() etc

## Importing Necesary dependencies

```
In [1]:
1  import matplotlib.pyplot as plt
2  import seaborn as sns
3  from wordcloud import WordCloud, STOPWORDS
4  import numpy as np
5  import pandas as pd
6  import nltk
7  nltk.download('stopwords')
8  from scipy import stats
9  import warnings,os
10 import re
11
12
13
14 from sklearn.model_selection import train_test_split
15 from sklearn.pipeline import Pipeline
16 from sklearn.feature_extraction.text import TfidfVectorizer
17 from sklearn.linear_model import LogisticRegression
18
19 from subprocess import check_output
20 import warnings
21 warnings.filterwarnings('ignore')
```

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

```
In [2]:
1  df = pd.read_csv('combined_data.csv')
2  df.head()
```

Out[2]:

| | author | date | Comment | like_count | reply_count | comment_length |
|---|---|---|---|---|---|---|
| 0 | @maryamadirie8722 | 2024-04-19T09:43:45Z | Cadaab ka rabi bilaash uma san mujrimiin Rabiy... | 0 | 0 | 61 |
| 1 | @hagiohaji7937 | 2024-04-19T03:35:24Z | Beenaale gacan ku dhiigle tuug, hala dilo. | 0 | 0 | 42 |
| 2 | @fadumawarsamemumin4232 | 2024-04-18T12:17:27Z | All Shaba a le lagu Haye Polis dada cun ah Qam... | 0 | 0 | 96 |
| 3 | @user-we3ty2ct8f | 2024-04-17T22:02:09Z | Hadu I yaqaano ha ii abtiriyo kulahaa kk ma ho... | 0 | 0 | 125 |
| 4 | @abdiali4991 | 2024-04-17T08:07:17Z | Sarkaal kan bistoolda loo dibtay waa ina asaga... | 0 | 0 | 117 |

```
In [3]:  1  df.describe()
```

Out[3]:

|       | like_count  | reply_count | comment_length |
|-------|-------------|-------------|----------------|
| count | 2527.000000 | 2527.000000 | 2527.000000    |
| mean  | 2.000000    | 0.237831    | 96.521567      |
| std   | 9.466614    | 1.211593    | 112.602332     |
| min   | 0.000000    | 0.000000    | 1.000000       |
| 25%   | 0.000000    | 0.000000    | 36.000000      |
| 50%   | 0.000000    | 0.000000    | 68.000000      |
| 75%   | 1.000000    | 0.000000    | 118.500000     |
| max   | 171.000000  | 27.000000   | 1775.000000    |

```
In [4]:  1  len(df)
```

Out[4]: 2527

# Task 3 - Let's clean the data by removing the extra data and outliers, dropping or filling the missing values, etc. Create the final dataframe for further analysis

## Cheking For duplicate values

```
In [5]:  1  df.duplicated().sum()
```

Out[5]: 1

## Droping Duplicate Values

```
In [6]:  1  df.drop_duplicates(inplace=True)
```

```
In [7]:  1  df.duplicated().sum()
```

Out[7]: 0

## Check for missing values

```
In [8]:  1  df.isna().sum()
```

Out[8]:
```
author           0
date             0
Comment          0
like_count       0
reply_count      0
comment_length   0
dtype: int64
```

```
In [9]:   1  # Remove rows with missing values
          2  df_cleaned = df.dropna()
```

```
In [10]:  1  df=df_cleaned
```

# Task 4: Lets analyze and visualize the distribution of post length and word counts?

**Calculate word count**

```
In [11]:    1  word_count = df['word_count'] = df['Comment'].apply(lambda x: len(x.split()))
```

```
In [12]:    1  word_count
```

```
Out[12]:  0        10
          1         7
          2        18
          3        22
          4        19
                   ..
          2522      3
          2523      3
          2524     17
          2525      4
          2526     12
          Name: Comment, Length: 2526, dtype: int64
```

```
In [13]:    1  df.head()
```

Out[13]:

| | author | date | Comment | like_count | reply_count | comment_length | word_count |
|---|---|---|---|---|---|---|---|
| **0** | @maryamadirie8722 | 2024-04-19T09:43:45Z | Cadaab ka rabi bilaash uma san mujrimiin Rabiy... | 0 | 0 | 61 | 10 |
| **1** | @hagiohaji7937 | 2024-04-19T03:35:24Z | Beenaale gacan ku dhiigle tuug, hala dilo. | 0 | 0 | 42 | 7 |
| **2** | @fadumawarsamemumin4232 | 2024-04-18T12:17:27Z | All Shaba a le lagu Haye Polis dada cun ah Qam... | 0 | 0 | 96 | 18 |
| **3** | @user-we3ty2ct8f | 2024-04-17T22:02:09Z | Hadu I yaqaano ha ii abtiriyo kulahaa kk ma ho... | 0 | 0 | 125 | 22 |
| **4** | @abdiali4991 | 2024-04-17T08:07:17Z | Sarkaal kan bistoolda loo dibtay waa ina asaga... | 0 | 0 | 117 | 19 |

## Summary Statistics

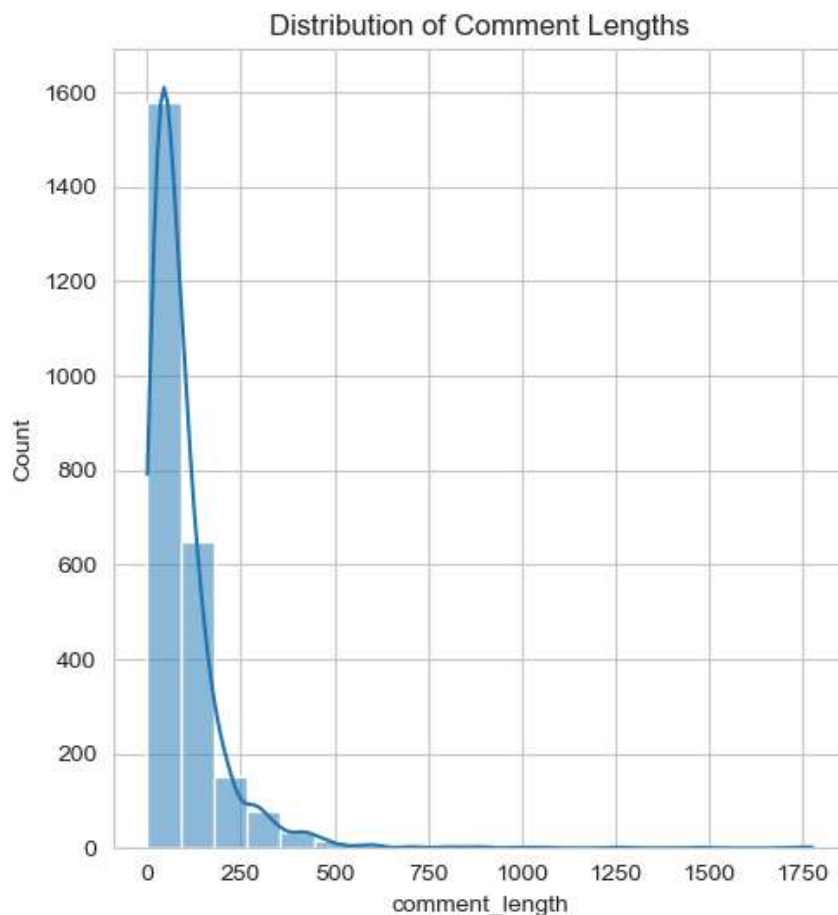**Before visualizing, it's often useful to look at summary statistics to get a sense of the distribution**

In [14]:
```python
# Display summary statistics for post length and word count
summary_stats = df[['comment_length', 'word_count']].describe()
print(summary_stats)
```

```
       comment_length    word_count
count     2526.000000   2526.000000
mean        96.545131     15.173793
std        112.618396     17.951119
min          1.000000      1.000000
25%         36.000000      6.000000
50%         68.000000     11.000000
75%        118.750000     19.000000
max       1775.000000    298.000000
```
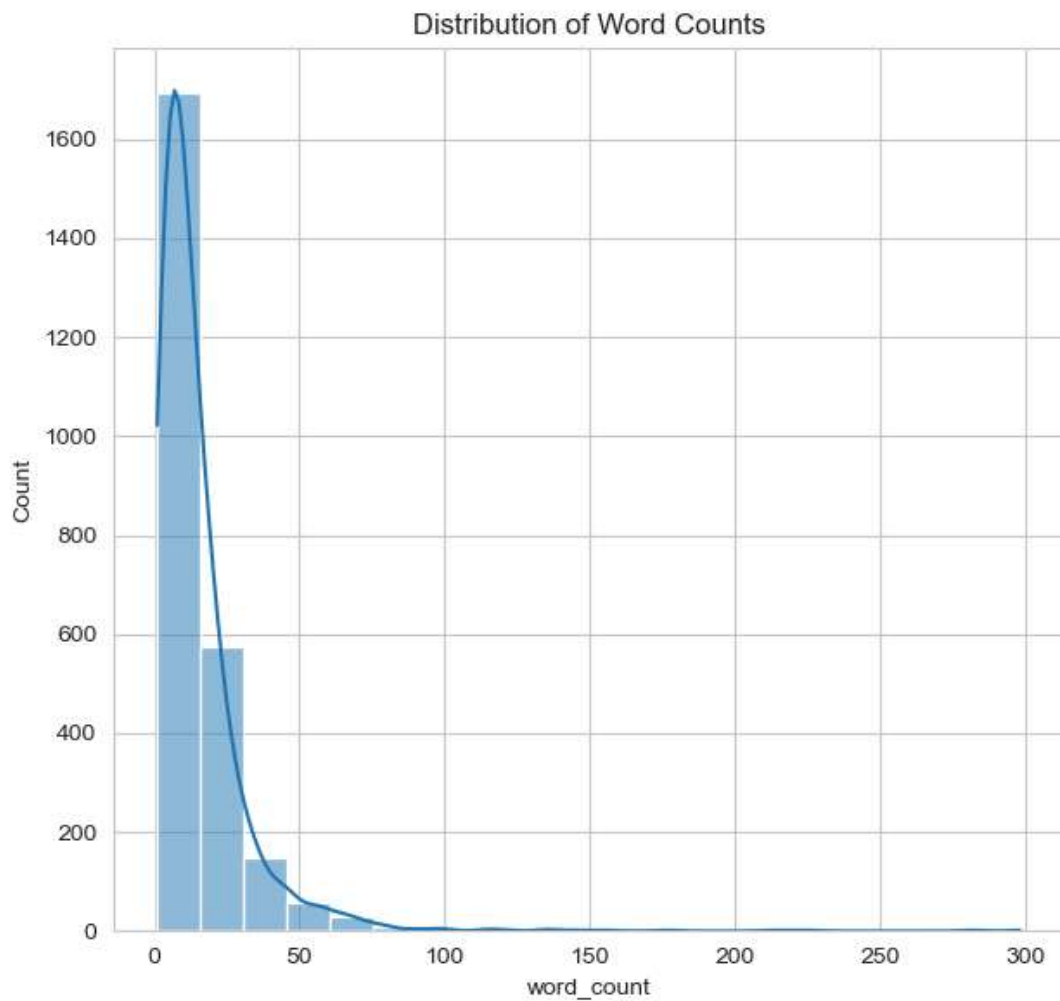
## Visualizing the Distribution

In [15]:
```python
# Set the aesthetic style of the plots
sns.set_style("whitegrid")

# Plot distribution of post lengths
plt.figure(figsize=(12, 6))
plt.subplot(1, 2, 1)   # 1 row, 2 columns, 1st subplot
sns.histplot(df['comment_length'], kde=True, bins=20)
plt.title('Distribution of Comment Lengths')
```

Out[15]: Text(0.5, 1.0, 'Distribution of Comment Lengths')

```
In [16]:    1  # Plot distribution of word counts
            2  plt.figure(figsize=(12, 6))
            3  plt.subplot(1, 2, 2)   # 1 row, 2 columns, 2nd subplot
            4  sns.histplot(df['word_count'], kde=True, bins=20)
            5  plt.title('Distribution of Word Counts')
            6
            7  plt.tight_layout()
            8  plt.show()
```



Distribution of Word Counts

# Task 5 - Let's analyze and visualize the most trending comments

## the most trending comments Sort by Likes and Get Top 5

In [17]:
```python
df.sort_values(by='like_count', ascending=False)[0:5]
```

Out[17]:

| | author | date | Comment | like_count | reply_count | comment_length | word_count |
|---|---|---|---|---|---|---|---|
| **999** | @jimaleahmed1444 | 2024-04-09T22:18:00Z | Taliye Shuute iyo Xeer ilaalintuba waaa dad Mu... | 171 | 15 | 164 | 24 |
| **1192** | @Cadar49 | 2024-04-09T18:10:22Z | Libaan waa in lasoo qabtaa asgaa ugu danbi wee... | 165 | 27 | 132 | 24 |
| **1249** | @sadiyoabdullahi5076 | 2024-04-09T17:44:23Z | Wiilka kalena weey ku dari lahaayen ee tiisa g... | 162 | 10 | 79 | 14 |
| **1371** | @Aniga-qn9iv | 2024-04-09T16:08:41Z | Dad Xamar jooga way adag tahay ama waqti badan... | 155 | 27 | 96 | 17 |
| **1284** | @poom5180 | 2024-04-09T17:26:29Z | Shabakooyin kale ayaa jira ee sifiican halooga... | 126 | 4 | 65 | 9 |

In [18]:
```python
# Extracting the information for visualization
top_liked_comments_users = df[['author', 'Comment', 'like_count']].sort_values(by='like_
                                                                                
# Visualization
plt.figure(figsize=(10, 8))
sns.barplot(x='like_count', y='Comment', hue='author', dodge=False, data=top_liked_commen
plt.title('Top 10 Liked Comments and Their Authors')
plt.xlabel('Likes')
plt.ylabel('Comments')
plt.legend(title='Author', loc='lower right')
plt.show()
```

**the most trending comments Sort by replies and Get Top 5**

In [19]:
```
1  df.sort_values(by='reply_count', ascending=False)[0:5]
```

Out[19]:

| | author | date | Comment | like_count | reply_count | comment_length | word_c... |
|---|---|---|---|---|---|---|---|
| 1192 | @Cadar49 | 2024-04-09T18:10:22Z | Libaan waa in lasoo qabtaa asgaa ugu danbi wee... | 165 | 27 | 132 | |
| 1371 | @Aniga-qn9iv | 2024-04-09T16:08:41Z | Dad Xamar jooga way adag tahay ama waqti badan... | 155 | 27 | 96 | |
| 1144 | @CoofavBulle-og5ld | 2024-04-09T18:48:58Z | Qoslaaye\nLiibaan\nQoone \nseddexdaas nin hada... | 16 | 18 | 112 | |
| 999 | @jimaleahmed1444 | 2024-04-09T22:18:00Z | Taliye Shuute iyo Xeer ilaalintuba waaa dad Mu... | 171 | 15 | 164 | |
| 1354 | @heeganbuuxa-tx1te | 2024-04-09T16:33:14Z | kuligood waa cabtooy mid midkale dhaama malaha... | 75 | 11 | 274 | |

In [20]:
```python
1   # Extracting the information for visualization
2   top_replied_comments_users = df[['author', 'Comment', 'reply_count']].sort_values(by='rep
3
4   # Visualization
5   plt.figure(figsize=(10, 8))
6   sns.barplot(x='reply_count', y='Comment', hue='author', dodge=False, data=top_replied_col
7   plt.title('Top 10 Liked Comments and Their Authors')
8   plt.xlabel('Likes')
9   plt.ylabel('Comments')
10  plt.legend(title='Author', loc='lower right')
11  plt.show()
12
```



# Task 6 - Let's analyze and visualize the top 5 users by number of comments

## Unique Authors

In [21]:
```python
1  unique_authors = df['author'].nunique()
2  unique_authors
```
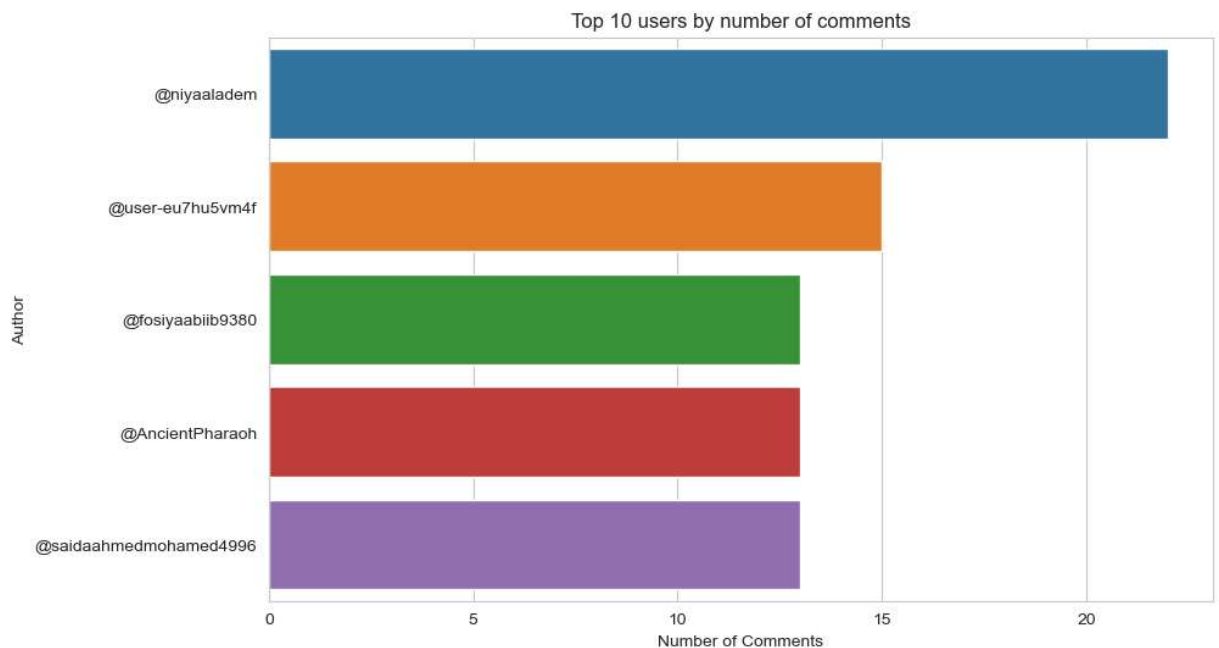
Out[21]: 1705

In [22]:
```python
1  # Author analysis
2  top_contributors = df['author'].value_counts().head(5)
3  top_contributors
```

Out[22]:
```
@niyaaladem              22
@user-eu7hu5vm4f         15
@fosiyaabiib9380         13
@AncientPharaoh          13
@saidaahmedmohamed4996   13
Name: author, dtype: int64
```

```
In [23]:  1  # Top Contributors
          2  plt.figure(figsize=(10, 6))
          3  sns.barplot(x=top_contributors.values, y=top_contributors.index)
          4  plt.title('Top 10 users by number of comments')
          5  plt.xlabel('Number of Comments')
          6  plt.ylabel('Author')
          7  plt.show()
```



## Task 7 - Let's visualize the prevalent words in the comments using WordCloud

```
In [24]:   1
           2  # defining my own list completely
           3  my_stopwords = {'ku', 'iyo', 'uu', 'ha', 'ma', 'laga', 'ugu', 'waa', 'in', 'ee', 'aa', '
           4               'inay', 'la', 'ah', 'ka', 'ayaa', 'iska', 'wax', 'oo', 'soo', 'ayuu', 'ba
           5               'wa', 'i', 'buu', 'inu', 'loo', 'waxaa', 'waxaan', 'ah', 'lama', 'maxay'
           6               'a', 'lagu', 'maxaa', 'inuu', 'wada', 'wuxuu', 'hala', 'e', 'waxan', 'sid
           7               'ah', 'wuu', 'ama', 'sii', 'hadii', 'ay' 'aya', 'siduu', 'yaa', 'ayu', '
           8               'ba', 'aha', 'iga', 'baa', 'ay', 'muxuu', 'maa', 'is', 'ayay', 'so', 'lah
           9               'haa', 'he', 'ilaa', 'hada', 'mida', 'may', 'waxa', 'waan', 'hadaad', 'wa
          10               'inaad', 'og', 'sow', 'inta', 'haku', 'lahaa', 'inay', 'alx', 'kugu', 'ya
          11               'miya', 'ayey', 'maxa', 'haduu', 'leh', 'ayan', 'hadi', 'iney', 'isku', '
          12               'muu', 'kama', 'aamiin', 'Asc', 'asaga', 'no', 'aniga', 'ahayn', 'lahayn
          13               'kulaha', 'xamar', 'kan', 'ilahow', 'ey', 'asaga' 'wll', 'to' }
          14
          15
```

```python
# Word cloud for comments
text = ' '.join(comment for comment in df.Comment)
stopwords = set(my_stopwords)
wordcloud = WordCloud(stopwords=stopwords, background_color="white", width=1000, height=
            min_font_size = 10).generate(text)

plt.figure(figsize=(10, 8), facecolor = None)
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.show()
```

## Task 8 - Let's visualize in what STATE/COUNTRY were the top posts posted that mention your TOPIC

unfortunately YouTube's Data API restricts access to directly access to the location data of its users due to their privacy and policy reasons.

## Task 9 - Let's visualize on what day(s) of the month was your TOPIC talked about the most on the selected social media platform

checking my data types to see if date column is datetime

```
In [26]:   1  df.dtypes
```

```
Out[26]: author          object
         date            object
         Comment         object
         like_count       int64
         reply_count      int64
         comment_length   int64
         word_count       int64
         dtype: object
```

```
In [27]:   1  # Convert the 'date' column to datetime
           2  df['date'] = pd.to_datetime(df['date'])
           3
```

```
In [28]:   1  # Verify the conversion by checking the data type again
           2  print(df['date'].dtype)
```

```
datetime64[ns, UTC]
```

```
In [29]:   1  df.dtypes
```

```
Out[29]: author                    object
         date            datetime64[ns, UTC]
         Comment                   object
         like_count                 int64
         reply_count                int64
         comment_length             int64
         word_count                 int64
         dtype: object
```

### Extracting Day of the Month

```
In [30]:   1  df['day_of_month'] = df['date'].dt.day
```

### Aggregate Data by Day of the Month

Now, count how many comments were made on each day of the month. This involves grouping the 'data by the day_of_month' and counting the number of comments.
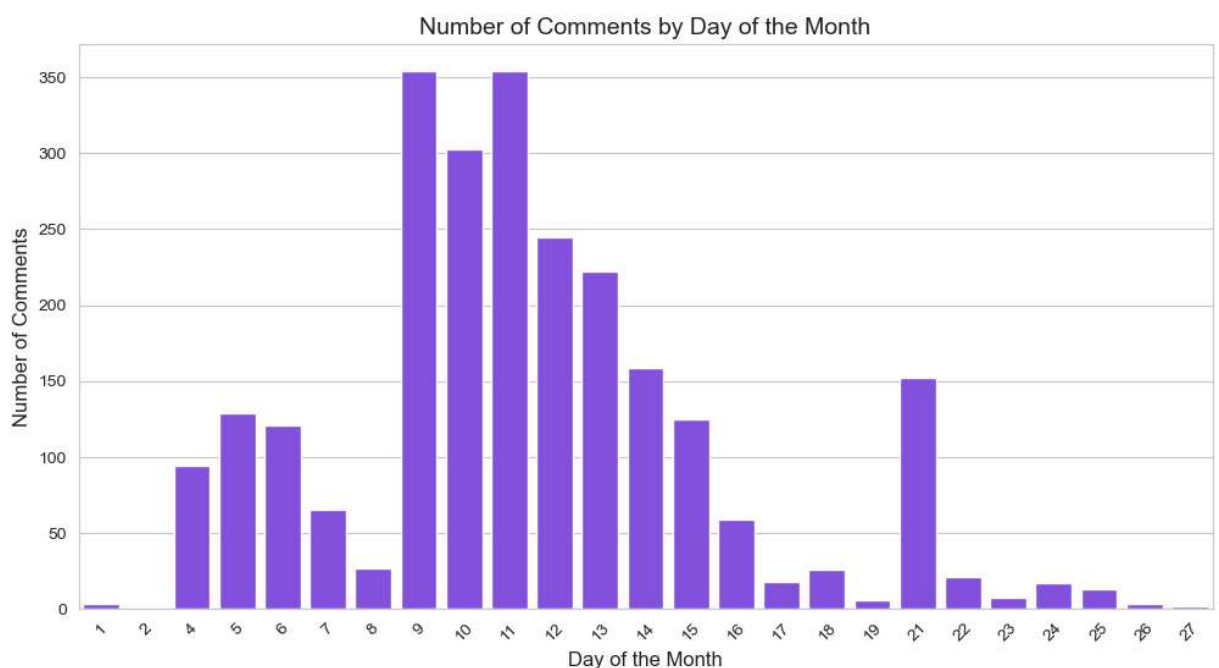
```
1 comments_by_day = df.groupby('day_of_month')['Comment'].count()
2 comments_by_day
```

Out[31]: day_of_month
```
1        3
2        1
4       94
5      129
6      121
7       65
8       27
9      354
10     303
11     354
12     245
13     222
14     159
15     125
16      59
17      18
18      26
19       6
21     152
22      21
23       7
24      17
25      13
26       3
27       2
Name: Comment, dtype: int64
```

## Visualize the Data

**Finally, use a bar chart to visualize the number of comments per day of the month.**

In [32]:

```
1 plt.figure(figsize=(12, 6))
2 sns.barplot(x=comments_by_day.index, y=comments_by_day.values, color='#7D3CF8')  # Electi
3 plt.title('Number of Comments by Day of the Month', fontsize=14)
4 plt.xlabel('Day of the Month', fontsize=12)
5 plt.ylabel('Number of Comments', fontsize=12)
6 plt.xticks(rotation=45) # Helps with readability if there are many days
7 plt.show()
8
9
```

## Task 10: Let's collect 3,000 statements across Somali public pages on Social Media Outlet and annotate them into three sentiment labels:- Positive-wanaag, Negative-xumaan or Neutral-dhexdhexaad

```python
In [33]:   1  Sentiments = pd.read_csv('Sentiment_data.csv')
           2  Sentiments.head()
```

Out[33]:

|   | Statements | Label |
|---|---|---|
| 0 | naxariis | 1 |
| 1 | jaceyl | 1 |
| 2 | amaan | 1 |
| 3 | jiidasho | 1 |
| 4 | raali galin | 1 |

```python
In [34]:   1  Sentiments.isna().sum()
```

Out[34]:  Statements    66
          Label          0
          dtype: int64

```python
In [35]:   1  Sentiments = Sentiments.dropna()
```

```python
In [36]:   1  Sentiments.isna().sum()
```

Out[36]:  Statements    0
          Label         0
          dtype: int64

```python
In [37]:   1  Sentiments.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 3313 entries, 0 to 3378
Data columns (total 2 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   Statements  3313 non-null   object
 1   Label       3313 non-null   int64
dtypes: int64(1), object(1)
memory usage: 77.6+ KB
```

```python
In [38]:   1  def remove_punctuation(text):
           2      return "".join([t for t in text if t not in string.punctuation])
           3
```

```python
In [39]:   1  def words_with_more_than_three_chars(text):
           2      return " ".join([t for t in text.split() if len(t)>3])
           3
```

```python
In [40]:   1  import string
           2
           3  def remove_punctuation(x):
           4      # Check if x is a string instance
           5      if isinstance(x, str):
           6          # If x is a string, remove punctuation
           7          return x.translate(str.maketrans('', '', string.punctuation))
           8      else:
           9          # If x is not a string (e.g., NaN or a number), return it unchanged
          10          return x
```

```
In [41]:  1  # Assuming Sentiments is your DataFrame and 'Statement' is a column in it
          2  Sentiments['Statements'] = Sentiments['Statements'].apply(lambda x: remove_punctuation(x
          3
          4  # Display the first 10 rows of the DataFrame
          5  Sentiments.head(10)
          6
```

Out[41]:

| | Statements | Label |
|---|---|---|
| 0 | naxariis | 1 |
| 1 | jaceyl | 1 |
| 2 | amaan | 1 |
| 3 | jiidasho | 1 |
| 4 | raali galin | 1 |
| 5 | hambalyo | 1 |
| 6 | faraxad | 1 |
| 7 | qurux | 1 |
| 8 | quruxsan | 1 |
| 9 | jecel yahay | 1 |

```
In [42]:  1  # Convert non-string values to strings
          2  Sentiments['Statements'] = Sentiments['Statements'].astype(str)
```

```
In [43]:  1  Sentiments['Statements']=Sentiments['Statements'].apply(lambda x:remove_punctuation(x))
          2  Sentiments['Statements']=Sentiments['Statements'].apply(lambda x:words_with_more_than_thi
          3  Sentiments['Statements']=Sentiments['Statements'].apply(lambda x: ' '.join([word for word
```
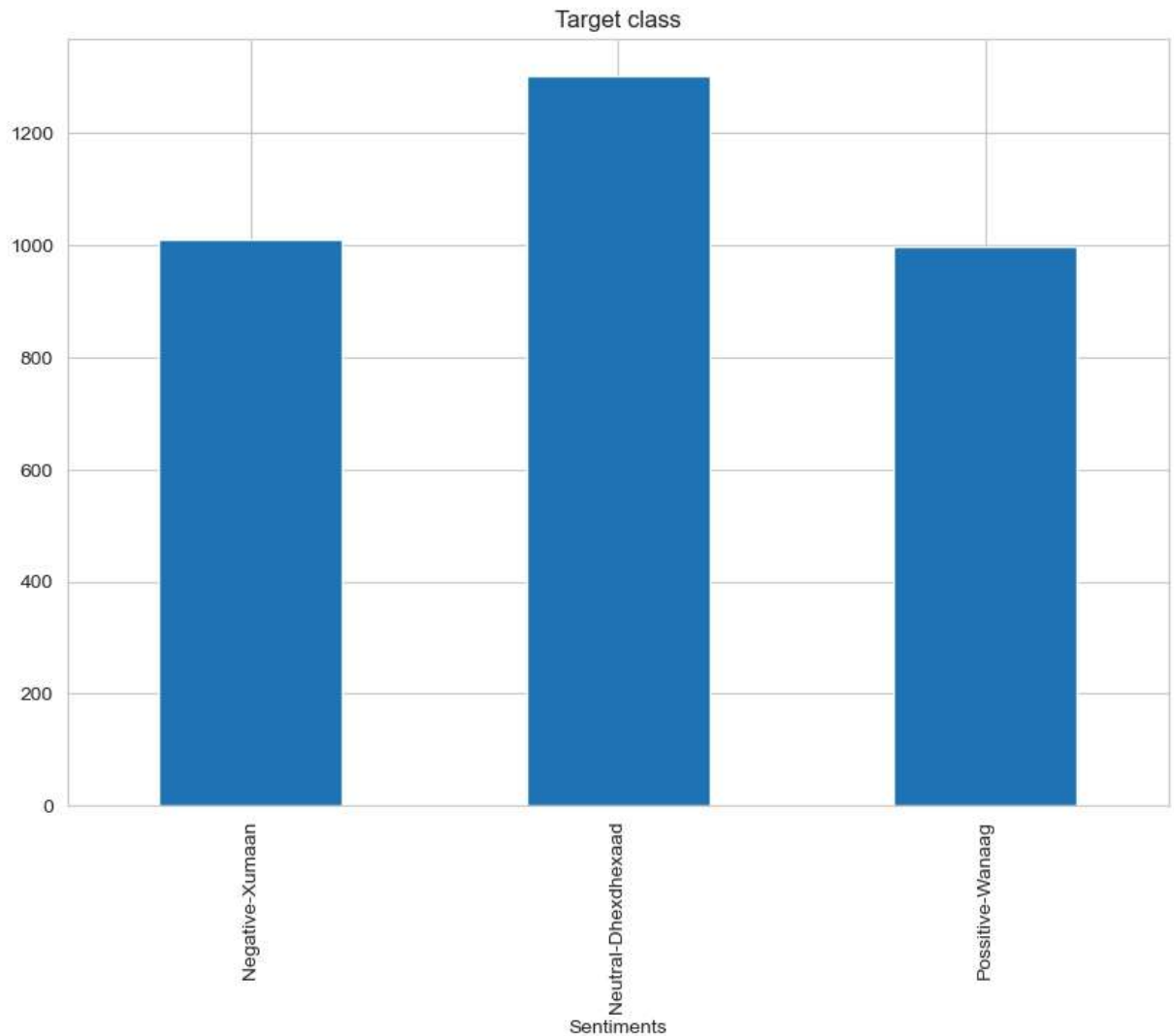
```
In [44]:  1  Sentiments['Sentiments'] = Sentiments['Label'].map({1:'Possitive-Wanaag',2:'Neutral-Dhex
          2  Sentiments.tail()
```

Out[44]:

| | Statements | Label | Sentiments |
|---|---|---|---|
| 3374 | dhiladii dhashay | 3 | Negative-Xumaan |
| 3375 | intaas lagugu badnasaa loodaawado dagaal kiina... | 3 | Negative-Xumaan |
| 3376 | Dhocil masersan xnuun | 3 | Negative-Xumaan |
| 3377 | Fartuuun tuugo kasoo sameey part2 wlhi qoslaay... | 3 | Negative-Xumaan |
| 3378 | Fartuuntuu sheegaayo maba fahmine dhib badnaa ... | 3 | Negative-Xumaan |

```
In [45]:  1  Sentiments.groupby('Sentiments')['Sentiments'].count().plot(kind='bar',title='Target cla
          2
```

Out[45]: `<Axes: title={'center': 'Target class'}, xlabel='Sentiments'>`



```
In [46]:  1  # Convert stop_words_to_lower to a set for faster membership testing
          2  stop_words_set = set(my_stopwords)
          3
          4  # Define a function to clean and process each comment
          5  def process_comment(comment):
          6      review = re.sub('[^a-zA-Z]', ' ', str(comment))
          7      review = review.lower().split()
          8      review = [word for word in review if word not in stop_words_set]
          9      return ' '.join(review)
```

```
In [47]:  1  # Apply the function to each comment in the DataFrame
          2  Sentiments['ProcessedStatement'] = Sentiments['Statements'].apply(process_comment)
          3
          4  Sentiments.head()
```

Out[47]:

|   | Statements | Label | Sentiments | ProcessedStatement |
|---|---|---|---|---|
| 0 | naxariis | 1 | Possitive-Wanaag | naxariis |
| 1 | jaceyl | 1 | Possitive-Wanaag | jaceyl |
| 2 | amaan | 1 | Possitive-Wanaag | amaan |
| 3 | jiidasho | 1 | Possitive-Wanaag | jiidasho |
| 4 | raali galin | 1 | Possitive-Wanaag | raali galin |

```python
In [48]:   1  # Task 11 - Let's train a sentiment analysis model using machine learning
           2
           3  X = Sentiments['ProcessedStatement']  # The column containing text data
           4  y =  Sentiments['Sentiments']  # The column containing sentiment labels
```

```python
In [49]:   1  # Split the dataset into training and testing sets
           2  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42
```

```python
In [50]:   1
           2  # Define a machine learning pipeline
           3  pipeline = Pipeline([
           4      ('tfidf', TfidfVectorizer(lowercase=True)),
           5      ('classifier', LogisticRegression(random_state=42)),
           6  ])
```

```python
In [51]:   1  # Train the model with the corrected data
           2  pipeline.fit(X_train, y_train)
```

Out[51]:
```
  ▸      Pipeline

   ▸ TfidfVectorizer

  ▸ LogisticRegression
```

```python
In [52]:   1  # Accuracy Testing on the Model
           2
           3  from sklearn.metrics import accuracy_score
           4  y_pred = pipeline.predict(X_test)
           5  accuracy = accuracy_score(y_test, y_pred)
           6  print(f"Accuracy: {accuracy * 100:.2f}%")
           7
```

```
Accuracy: 82.20%
```

```python
In [54]:   1  # Convert stop_words_to_lower to a set for faster membership testing
           2  stop_words_set = set(my_stopwords)
           3
           4  # Define a function to clean and process each comment
           5  def process_comment(comment):
           6      review = re.sub('[^a-zA-Z]', ' ', str(comment))
           7      review = review.lower().split()
           8      review = [word for word in review if word not in stop_words_set]
           9      return ' '.join(review)
          10
          11  # Apply the function to each comment in the DataFrame
          12  df['ProcessedComment'] = df['Comment'].apply(process_comment)
          13
```

```python
In [55]:   1  # Task 12 - Let's now apply the trained model on the selected topic textual contents in d
           2  # Now, you can predict with confidence that the pipeline is fitted
           3  predicted_sentiments = pipeline.predict(df['ProcessedComment'])
           4  df['Predicted Sentiment'] = predicted_sentiments
           5
           6  # Displaying the original comments with their predicted sentiments
           7  print(df[['ProcessedComment', 'Predicted Sentiment']].head())
```

```
                                ProcessedComment Predicted Sentiment
0  cadaab rabi bilaash uma san mujrimiin rabiyaw ...    Negative-Xumaan
1                  beenaale gacan dhiigle tuug dilo    Negative-Xumaan
2  all shaba le haye polis dada cun qamriga qaram...    Negative-Xumaan
3  hadu yaqaano ii abtiriyo kulahaa kk hoyadaaba ...    Negative-Xumaan
4  sarkaal bistoolda dibtay ina asagana talabo la...  Possitive-Wanaag
```

```python
## Applying  the trained model to predict sentiments of the selected topic textual conten
# Applying the trained model to predict sentiments of the selected topic textual contents
predicted_sentiments = pipeline.predict(df['ProcessedComment'])

# Add the predictions to the original dataframe for review
df['Predicted Sentiment'] = predicted_sentiments

# Displaying the original comments with their predicted sentiments
print(df[['ProcessedComment', 'ProcessedComment']].head())
```

```
                              ProcessedComment  \
0  cadaab rabi bilaash uma san mujrimiin rabiyaw ...
1                    beenaale gacan dhiigle tuug dilo
2  all shaba le haye polis dada cun qamriga qaram...
3  hadu yaqaano ii abtiriyo kulahaa kk hoyadaaba ...
4  sarkaal bistoolda dibtay ina asagana talabo la...

                              ProcessedComment
0  cadaab rabi bilaash uma san mujrimiin rabiyaw ...
1                    beenaale gacan dhiigle tuug dilo
2  all shaba le haye polis dada cun qamriga qaram...
3  hadu yaqaano ii abtiriyo kulahaa kk hoyadaaba ...
4  sarkaal bistoolda dibtay ina asagana talabo la...
```

```python
df.head()
```

| | author | date | Comment | like_count | reply_count | comment_length | word_count | da |
|---|---|---|---|---|---|---|---|---|
| 0 | @maryamadirie8722 | 2024-04-19 09:43:45+00:00 | Cadaab ka rabi bilaash uma san mujrimiin Rabiy... | 0 | 0 | 61 | 10 | |
| 1 | @hagiohaji7937 | 2024-04-19 03:35:24+00:00 | Beenaale gacan ku dhiigle tuug, hala dilo. | 0 | 0 | 42 | 7 | |
| 2 | @fadumawarsamemumin4232 | 2024-04-18 12:17:27+00:00 | All Shaba a le lagu Haye Polis dada cun ah Qam... | 0 | 0 | 96 | 18 | |
| 3 | @user-we3ty2ct8f | 2024-04-17 22:02:09+00:00 | Hadu I yaqaano ha ii abtiriyo kulahaa kk ma ho... | 0 | 0 | 125 | 22 | |
| 4 | @abdiali4991 | 2024-04-17 08:07:17+00:00 | Sarkaal kan bistoolda loo dibtay waa ina asaga... | 0 | 0 | 117 | 19 | |

```
1  # Task 13 - Let's visualize the sentiment percentages (positive, negative, neutral) for
2
3  ## before the visualization let's calculating sentiment distribution in the predicted se
4
5  # Calculating sentiment distribution in the predicted sentiments
6  predicted_sentiments_distribution = df['Predicted Sentiment'].value_counts(normalize=Tru
7  predicted_sentiments_distribution
8
```
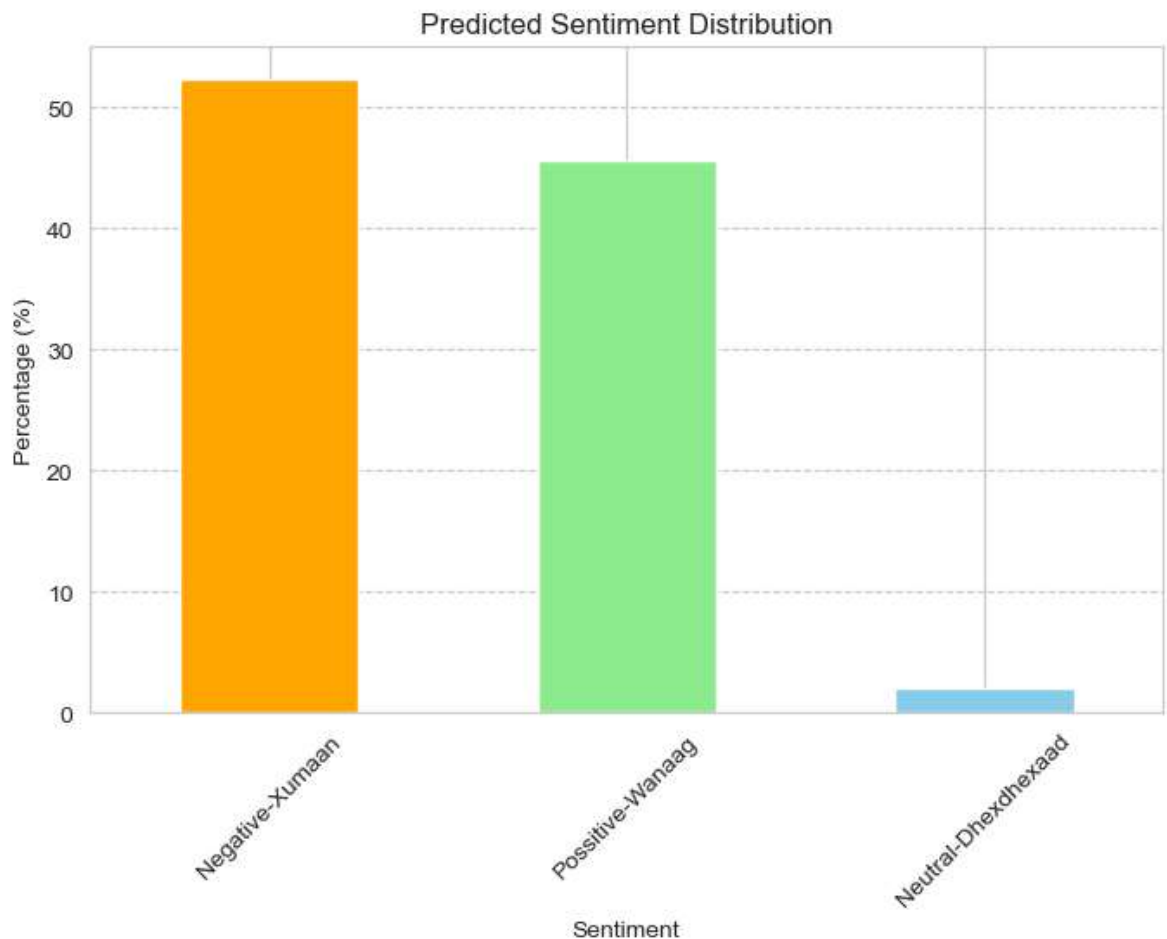
Out[58]: Negative-Xumaan        52.335709
         Possitive-Wanaag       45.645289
         Neutral-Dhexdhexaad     2.019002
         Name: Predicted Sentiment, dtype: float64

In [59]:

```
1  ## Visualizing the Predicted sentiment
2
3  plt.figure(figsize=(8, 5))
4  predicted_sentiments_distribution.plot(kind='bar', color=['orange', 'lightgreen', 'skybl
5  plt.title('Predicted Sentiment Distribution')
6  plt.xlabel('Sentiment')
7  plt.ylabel('Percentage (%)')
8  plt.xticks(rotation=45)
9  plt.grid(axis='y', linestyle='--')
10 plt.show()
11
```
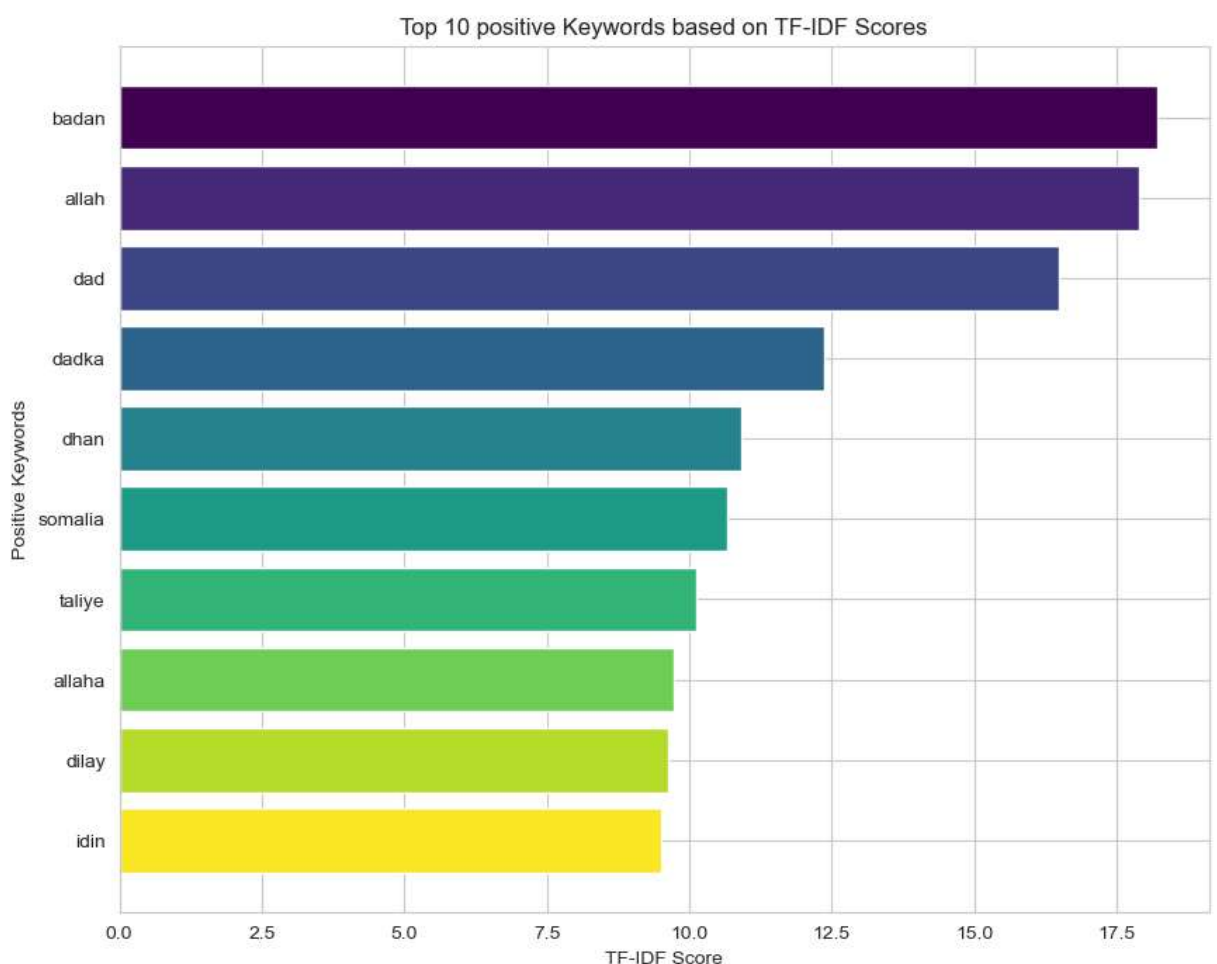
```
In [60]:  1  # Task 14 - Let's visualize the top positive keywords associated with the selected topic
          2
          3  from sklearn.feature_extraction.text import TfidfVectorizer
          4  import numpy as np
          5  import matplotlib.pyplot as plt
          6  import pandas as pd
          7
          8  # Filter the dataset for comments predicted as positive
          9  positive_comments = df[df['Predicted Sentiment'] == 'Possitive-Wanaag']['ProcessedCommen
         10
         11  # Apply TF-IDF vectorization
         12  tfidf_vectorizer = TfidfVectorizer(stop_words='english')  # Adjust or remove stop_words
         13  tfidf_matrix = tfidf_vectorizer.fit_transform(positive_comments)
         14
         15  # Get feature names and tfidf scores
         16  feature_names = tfidf_vectorizer.get_feature_names_out()
         17  scores = np.asarray(tfidf_matrix.sum(axis=0)).ravel()
         18  tfidf_scores = dict(zip(feature_names, scores))
         19
         20  # Sort the scores in descending order
         21  sorted_tfidf_scores = sorted(tfidf_scores.items(), key=lambda x: x[1], reverse=True)
         22
         23  # Convert to DataFrame for easy handling and visualization
         24  df_tfidf_scores = pd.DataFrame(sorted_tfidf_scores, columns=['Term', 'Score']).head(10)
         25
         26  # Visualization using a colormap
         27  plt.figure(figsize=(10, 8))
         28  colors = plt.cm.viridis(np.linspace(0, 1, len(df_tfidf_scores)))
         29  plt.barh(df_tfidf_scores['Term'], df_tfidf_scores['Score'], color=colors)
         30  plt.xlabel('TF-IDF Score')
         31  plt.ylabel('Positive Keywords')
         32  plt.title('Top 10 positive Keywords based on TF-IDF Scores')
         33  plt.gca().invert_yaxis()
         34  plt.show()
         35
         36
```
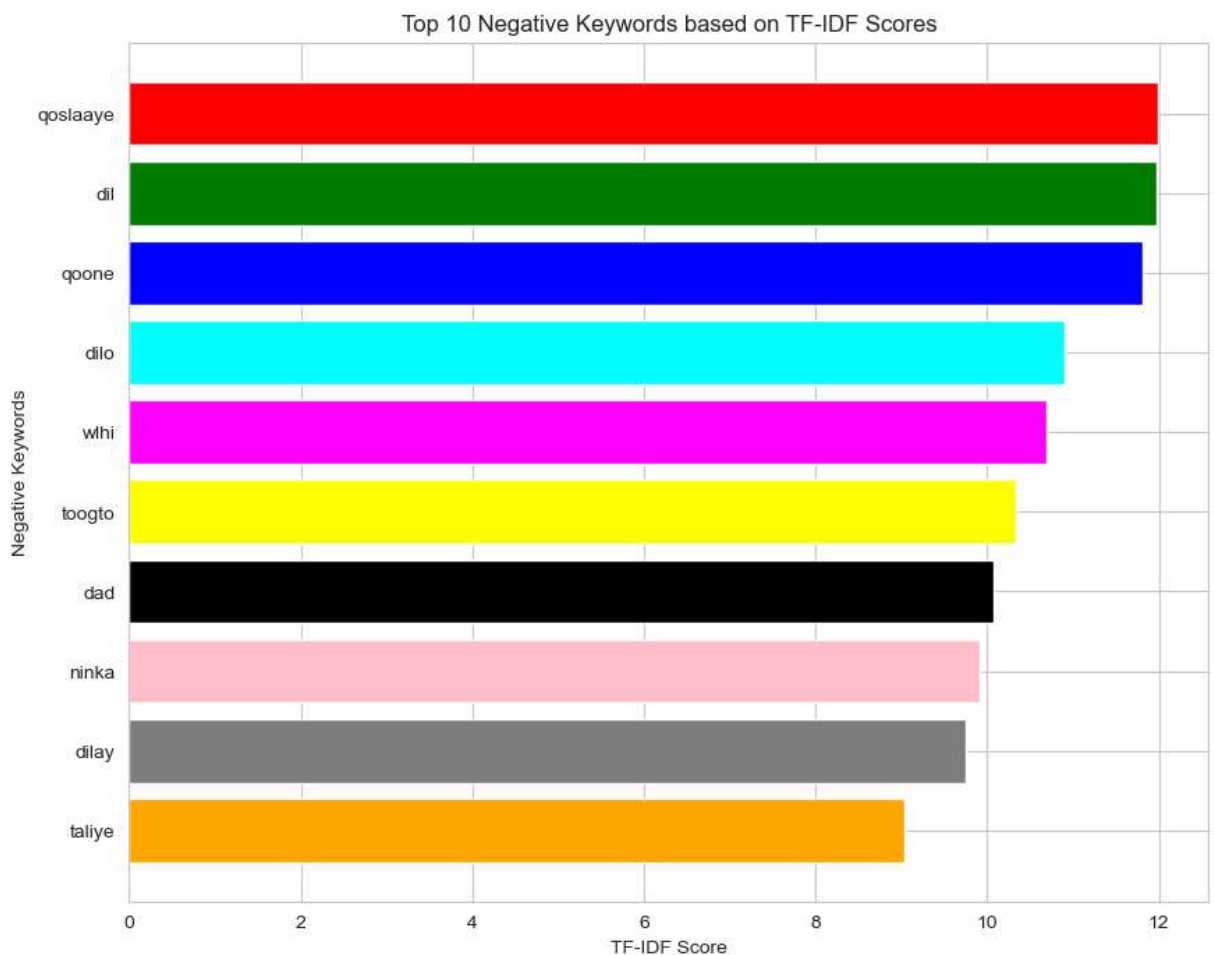


Top 10 positive Keywords based on TF-IDF Scores

```
In [61]:   1  # Task 15 - Let's visualize the top negative keywords associated with the selected topic
           2
           3  from sklearn.feature_extraction.text import TfidfVectorizer
           4  import numpy as np
           5  import matplotlib.pyplot as plt
           6  import pandas as pd
           7
           8  # Filter the dataset for comments predicted as positive
           9  positive_comments = df[df['Predicted Sentiment'] == 'Negative-Xumaan']['ProcessedComment
          10
          11  # Apply TF-IDF vectorization
          12  tfidf_vectorizer = TfidfVectorizer(stop_words='english')   # Adjust or remove stop_words
          13  tfidf_matrix = tfidf_vectorizer.fit_transform(positive_comments)
          14
          15  # Get feature names and tfidf scores
          16  feature_names = tfidf_vectorizer.get_feature_names_out()
          17  scores = np.asarray(tfidf_matrix.sum(axis=0)).ravel()
          18  tfidf_scores = dict(zip(feature_names, scores))
          19
          20  # Sort the scores in descending order
          21  sorted_tfidf_scores = sorted(tfidf_scores.items(), key=lambda x: x[1], reverse=True)
          22
          23  # Convert to DataFrame for easy handling and visualization
          24  df_tfidf_scores = pd.DataFrame(sorted_tfidf_scores, columns=['Term', 'Score']).head(10)
          25
          26  # Visualization with custom colors
          27  plt.figure(figsize=(10, 8))
          28  custom_colors = ['red', 'green', 'blue', 'cyan', 'magenta', 'yellow', 'black', 'pink', '
          29  plt.barh(df_tfidf_scores['Term'], df_tfidf_scores['Score'], color=custom_colors)
          30  plt.xlabel('TF-IDF Score')
          31  plt.ylabel('Negative Keywords')
          32  plt.title('Top 10 Negative Keywords based on TF-IDF Scores')
          33  plt.gca().invert_yaxis()
          34  plt.show()
          35
```

```
In [ ]:  1
```