

Adipurush Sentiment Analysis: Harnessing Machine Learning to Understand Twitter Buzz



```
In [3]: ┏ 1 import pandas as pd
```

```
In [4]: ┏ 1 import matplotlib.pyplot as plt  
2 import seaborn as sns
```

```
In [5]: ┏ 1 import warnings  
2 warnings.filterwarnings('ignore')
```

```
In [6]: ┏ 1 df = pd.read_csv('adipurush_tweets.csv')
```

In [7]: 1 df

Out[7]:

	Date Created	Number of Likes	Source of Tweet	Tweets
0	2023-06-30 09:21:00+00:00	0	NaN	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipuru...
1	2023-06-30 09:20:57+00:00	0	NaN	Now Playing!! Book Your Ticket Now!! 🎟️ 👉凛@go...
2	2023-06-30 09:20:22+00:00	0	NaN	@ponilemova #Adipurush
3	2023-06-30 09:20:00+00:00	3	NaN	Adipurush VS 72 Hoorain VS The Kerala Story Co...
4	2023-06-30 09:15:22+00:00	3	NaN	ST: #Adipurush https://t.co/lsgKcgQuKL
...
9996	2023-06-23 10:08:50+00:00	0	NaN	Adipurush 1st Week WW Box Office Collections: ...
9997	2023-06-23 10:08:49+00:00	0	NaN	#GodMorningFriday\nवास्तव में #Adipurush यानि ...
9998	2023-06-23 10:08:17+00:00	3101	NaN	Let the empowering lyrics of #Shivoham elevate...
9999	2023-06-23 10:08:01+00:00	0	NaN	When it comes to choosing a service or product...
10000	2023-06-23 10:07:45+00:00	0	NaN	A film about #Ramayana, our greatest epic coul...

10001 rows × 4 columns

In [8]: 1 df.shape

Out[8]: (10001, 4)

In [9]: 1 df.columns

Out[9]: Index(['Date Created', 'Number of Likes', 'Source of Tweet', 'Tweets'], dtype='object')

In [10]: 1 df.duplicated().sum()

Out[10]: 1

In [11]: 1 df = df.drop_duplicates()

In [12]: 1 df.isnull().sum()

```
Out[12]: Date Created      0
          Number of Likes    0
          Source of Tweet    10000
          Tweets              0
          dtype: int64
```

In [13]: 1 df = df.drop('Source of Tweet', axis = 1)

In [14]: 1 df

Out[14]:

	Date Created	Number of Likes	Tweets
0	2023-06-30 09:21:00+00:00	0	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipurush...
1	2023-06-30 09:20:57+00:00	0	Now Playing!! Book Your Ticket Now!! 🎟️🏠🎟️ \n@go...
2	2023-06-30 09:20:22+00:00	0	@ponilemova #Adipurush
3	2023-06-30 09:20:00+00:00	3	Adipurush VS 72 Hoorain VS The Kerala Story Co...
4	2023-06-30 09:15:22+00:00	3	ST: #Adipurush https://t.co/lGKcgQuKL
...
9996	2023-06-23 10:08:50+00:00	0	Adipurush 1st Week WW Box Office Collections: ...
9997	2023-06-23 10:08:49+00:00	0	#GodMorningFriday\nवास्तव में #Adipurush यानि ...
9998	2023-06-23 10:08:17+00:00	3101	Let the empowering lyrics of #Shivoham elevate...
9999	2023-06-23 10:08:01+00:00	0	When it comes to choosing a service or product...
10000	2023-06-23 10:07:45+00:00	0	A film about #Ramayana, our greatest epic coul...

10000 rows × 3 columns

In [15]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 10000 entries, 0 to 10000
Data columns (total 3 columns):
 #   Column            Non-Null Count  Dtype  
--- 
 0   Date Created      10000 non-null   object 
 1   Number of Likes   10000 non-null   int64  
 2   Tweets             10000 non-null   object 
dtypes: int64(1), object(2)
memory usage: 312.5+ KB
```

In [16]: 1 df.describe()

Out[16]:

Number of Likes	
count	10000.000000
mean	73.544500
std	369.705682
min	0.000000
25%	0.000000
50%	1.000000
75%	9.000000
max	14778.000000

In [17]: 1 df.nunique()

Out[17]: Date Created 9831
Number of Likes 718
Tweets 9874
dtype: int64

In [18]: 1 df_sorted = df.sort_values(by='Number of Likes', ascending=False)

In [19]: 1 df_sorted.head(10)

Out[19]:

	Date Created	Number of Likes	Tweets
5036	2023-06-26 02:51:52+00:00	14778	Pan India Star #Prabhas clearly said NO for #A...
2975	2023-06-27 12:35:31+00:00	8266	#Breaking: Comments by Allahabad high court to...
8180	2023-06-24 09:10:09+00:00	8112	#Adipurush #Prabhas #BhushanKumar https://t.co...
3593	2023-06-27 02:23:37+00:00	7010	आदिपुरुष निर्माताओं को लगा एक और झटका, इलाहाबा...
6069	2023-06-25 07:20:59+00:00	5580	👉#AdiPurush Telugu Version Hits 100CR SHARE 
3601	2023-06-27 01:59:03+00:00	5149	#Adipurush WW BO\n\nZOOMS past ₹40000 cr.\n...
4744	2023-06-26 06:20:58+00:00	4912	#Adipurush goes from strength to strength at t...
4716	2023-06-26 06:30:01+00:00	4788	We are incredibly touched by the overwhelming ...
5636	2023-06-25 13:08:17+00:00	4741	Witness the epic saga unfold! 🎟️\nBook your tic...
1559	2023-06-28 14:40:47+00:00	4561	कुरान पर गलत तथ्यों के साथ एक छोटी सी डॉक्यूमें...

In [20]: 1 df['Date Created'] = pd.to_datetime(df['Date Created'])

In [21]: 1 df

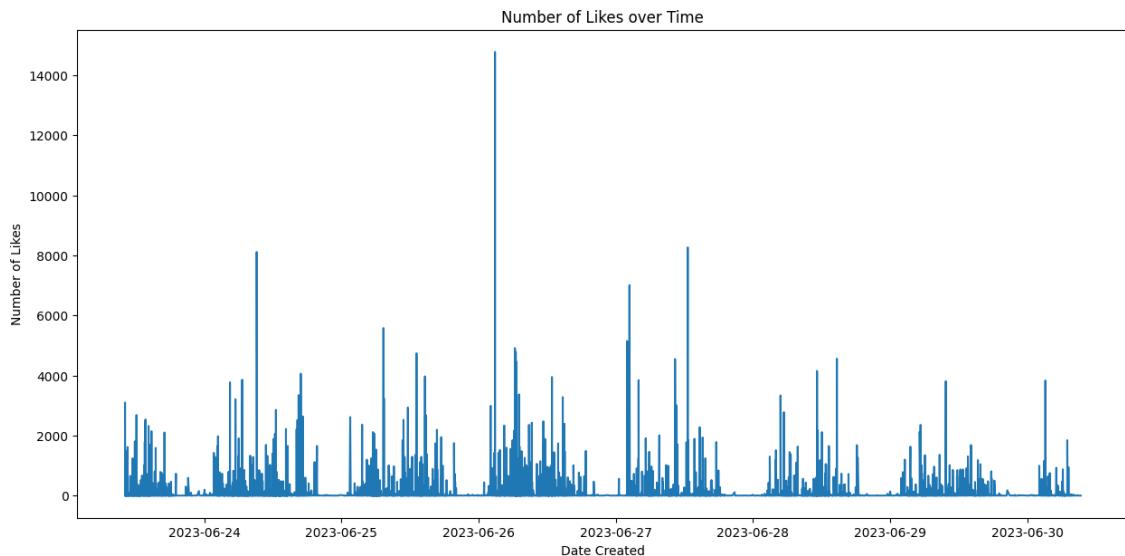
Out[21]:

	Date Created	Number of Likes	Tweets
0	2023-06-30 09:21:00+00:00	0	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipuru...
1	2023-06-30 09:20:57+00:00	0	Now Playing!! Book Your Ticket Now! 🎟️🍿副院长...
2	2023-06-30 09:20:22+00:00	0	@ponilemova #Adipurush
3	2023-06-30 09:20:00+00:00	3	Adipurush VS 72 Hoorain VS The Kerala Story Co...
4	2023-06-30 09:15:22+00:00	3	ST: #Adipurush https://t.co/lsgKcgQuKL
...
9996	2023-06-23 10:08:50+00:00	0	Adipurush 1st Week WW Box Office Collections: ...
9997	2023-06-23 10:08:49+00:00	0	#GodMorningFriday\nवारस्तव में #Adipurush यानि ...
9998	2023-06-23 10:08:17+00:00	3101	Let the empowering lyrics of #Shivoham elevate...
9999	2023-06-23 10:08:01+00:00	0	When it comes to choosing a service or product...
10000	2023-06-23 10:07:45+00:00	0	A film about #Ramayana, our greatest epic coul...

10000 rows × 3 columns

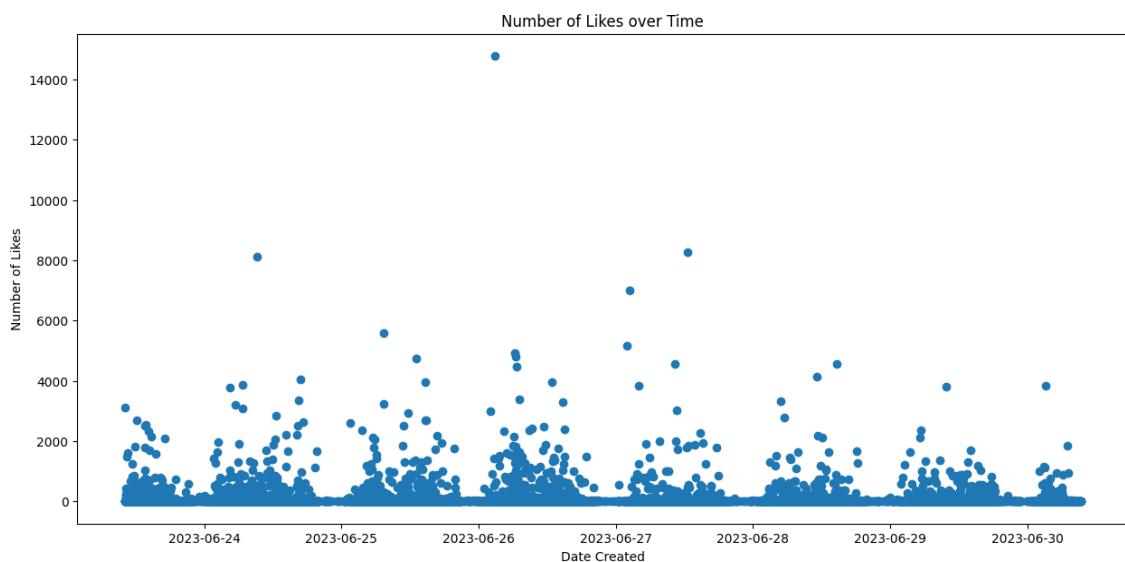
In [22]:

```
1 df_sorted_date = df.sort_values('Date Created')
2 plt.figure(figsize=[15,7],)
3 plt.plot(df_sorted_date['Date Created'], df_sorted_date['Number of Likes'])
4 plt.xlabel('Date Created')
5 plt.ylabel('Number of Likes')
6 plt.title('Number of Likes over Time')
7 plt.show()
```



In [23]:

```
1 plt.figure(figsize=[15,7],)
2 plt.scatter(df_sorted_date['Date Created'], df_sorted_date['Number of Likes'])
3 plt.xlabel('Date Created')
4 plt.ylabel('Number of Likes')
5 plt.title('Number of Likes over Time')
6 plt.show()
```



```
In [24]: ┌ 1 import re
  2 import string
  3 from tqdm.notebook import tqdm
  4 from datetime import datetime
  5 import dateutil.parser
```

```
In [25]: ┌ 1 import nltk
  2 from spellchecker import spellchecker
  3 from nltk.sentiment.vader import SentimentIntensityAnalyzer as SIA
```

```
In [26]: ┌ 1 from wordcloud import wordcloud, ImageColorGenerator
  2 from nltk.corpus import stopwords
  3 import random
```

```
In [27]: ┌ 1 nltk.download('vader_lexicon')
  2 nltk.download('stopwords')
```

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

Out[27]: True

```
In [28]: ┌ 1 languages = stopwords.fileids()
  2
  3 # Print the number of supported Languages
  4 print("Number of supported languages:", len(languages))
  5
  6 # Print the list of supported Languages
  7 print("Supported languages:", languages)
```

```
Number of supported languages: 29
Supported languages: ['arabic', 'azerbaijani', 'basque', 'bengali', 'catalan', 'chinese', 'danish', 'dutch', 'english', 'finnish', 'french', 'german', 'greek', 'hebrew', 'hinglish', 'hungarian', 'indonesian', 'italian', 'kazakh', 'nepali', 'norwegian', 'portuguese', 'romanian', 'russian', 'slovene', 'spanish', 'swedish', 'tajik', 'turkish']
```

```
In [29]: ┌ 1 from nltk.tokenize import TweetTokenizer
```

```
In [30]: ┌ 1 english_stopwords = stopwords.words('english')
  2 hinglish_stopwords = stopwords.words('hinglish')
```

In [31]:

```
1 def clean_tweet(tweet):
2     # Remove URLs, hashtags, mentions, and special characters
3     tweet = re.sub(r"http\S+|www\S+|@\w+|\#\w+", "", tweet)
4     tweet = re.sub(r"[\^\\w\\s]", "", tweet)
5
6     # Tokenize the tweet
7     tokenizer = TweetTokenizer(preserve_case=False, reduce_len=True, s
8     tokens = tokenizer.tokenize(tweet)
9
10    # Remove stopwords for English and Hinglish
11    tokens = [token for token in tokens if token not in english_stopwo
12
13    # Remove punctuation and convert to Lowercase
14    tokens = [token.translate(str.maketrans(' ', ' ', string.punctuation
15    tokens = [token.lower() for token in tokens]
16
17    # Join tokens back into a string
18    cleaned_tweet = ' '.join(tokens)
19
20    return cleaned_tweet
```

In [32]:

```
1 df['Cleaned_Tweets'] = df['Tweets'].apply(clean_tweet)
```

In [33]: 1 df

Out[33]:

	Date Created	Number of Likes	Tweets	Cleaned_Tweets
0	2023-06-30 09:21:00+00:00	0	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipuru...	womens ashes 2023 live streaming broadcast tv ...
1	2023-06-30 09:20:57+00:00	0	Now Playing!! Book Your Ticket Now!! 🎟️🍿🎥\n@go...	playing book ticket
2	2023-06-30 09:20:22+00:00	0	@ponilemova #Adipurush	
3	2023-06-30 09:20:00+00:00	3	Adipurush VS 72 Hoorain VS The Kerala Story Co...	adipurush vs 72 hoorain vs kerala story contro...
4	2023-06-30 09:15:22+00:00	3	ST: #Adipurush https://t.co/lGKcgQuKL	st
...
9996	2023-06-23 10:08:50+00:00	0	Adipurush 1st Week WW Box Office Collections: ...	adipurush 1st week ww box office collections ए...
9997	2023-06-23 10:08:49+00:00	0	#GodMorningFriday\nवास्तव में #Adipurush यानि ...	वस्तव म यन सबस पहल भगवन जसन सरव सषट क रचन क ह व...
9998	2023-06-23 10:08:17+00:00	3101	Let the empowering lyrics of #Shivoham elevate...	empowering lyrics elevate spirit envelop world...
9999	2023-06-23 10:08:01+00:00	0	When it comes to choosing a service or product...	choosing service product beneficial opt authen...
10000	2023-06-23 10:07:45+00:00	0	A film about #Ramayana, our greatest epic cou...	film greatest epic earn boc worth budget shame...

10000 rows × 4 columns

In [34]: 1 def clean_text(text):
2 text = text.lower()
3 return text.strip()

In [35]: 1 df.Cleaned_Tweets = df.Cleaned_Tweets.apply(lambda x: clean_text(x))

In [36]: 1 def tokenization(text):
2 tokens = re.split('W+',text)
3 return tokens

In [37]: 1 df.Cleaned_Tweets = df.Cleaned_Tweets.apply(lambda x: tokenization(x))

In [38]: 1 from nltk.stem import WordNetLemmatizer
2 wordnet_lemmatizer = WordNetLemmatizer()

In [39]: █ 1 nltk.download('wordnet')

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

Out[39]: True

In [40]: █ 1 nltk.download('omw-1.4')

```
[nltk_data] Downloading package omw-1.4 to
[nltk_data]     C:\Users\admin\AppData\Roaming\nltk_data...
[nltk_data]     Package omw-1.4 is already up-to-date!
```

Out[40]: True

In [41]: █ 1 def lemmatizer (text):
2 lemm_text = ''.join([wordnet_lemmatizer.lemmatize(word) for word in
3 return lemm_text

In [42]: █ 1 df.Cleaned_Tweets=df.Cleaned_Tweets.apply(lambda x: lemmatizer(x))

In [43]: █ 1 def remove_digits(text):
2 clean_text = re.sub(r"\b[0-9]+\b\s*", "", text)
3 return(text)

In [44]: █ 1 df.Cleaned_Tweets=df.Cleaned_Tweets.apply(lambda x: remove_digits(x))

In [45]: █ 1 def remove_digits1(sample_text):
2 clean_text = " ".join([w for w in sample_text.split() if not w.isd
3 return(clean_text)

In [46]: █ 1 df.Cleaned_Tweets=df.Cleaned_Tweets.apply(lambda x: remove_digits1(x))

In [47]: █ 1 from langdetect import detect
2
3 def detect_language(text):
4 try:
5 lang = detect(text)
6 return lang
7 except:
8 return None
9
10 df['Language'] = df['Cleaned_Tweets'].apply(detect_language)

In [48]: 1 df

Out[48]:

	Date Created	Number of Likes	Tweets	Cleaned_Tweets	Language
0	2023-06-30 09:21:00+00:00	0	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipuru...	womens ashes live streaming broadcast tv chann...	en
1	2023-06-30 09:20:57+00:00	0	Now Playing!! Book Your Ticket Now!! 🎟️🍿🎟...	playing book ticket	en
2	2023-06-30 09:20:22+00:00	0	@ponilemova #Adipurush		None
3	2023-06-30 09:20:00+00:00	3	Adipurush VS 72 Hoorain VS The Kerala Story Co...	adipurush vs hoorain vs kerala story controver...	en
4	2023-06-30 09:15:22+00:00	3	ST: #Adipurush https://t.co/lzGKcgQuKL		st no
...
9996	2023-06-23 10:08:50+00:00	0	Adipurush 1st Week WW Box Office Collections: ...	adipurush 1st week ww box office collections 😊...	en
9997	2023-06-23 10:08:49+00:00	0	#GodMorningFriday\nवास्तव में #Adipurush यानि ...	वस्तव म यन सबस पहल भगवन जसन सरव सष्ट क रचन क ह व...	hi
9998	2023-06-23 10:08:17+00:00	3101	Let the empowering lyrics of #Shivoham elevate...	empowering lyrics elevate spirit envelop world...	en
9999	2023-06-23 10:08:01+00:00	0	When it comes to choosing a service or product...	choosing service product beneficial opt authen...	en
10000	2023-06-23 10:07:45+00:00	0	A film about #Ramayana, our greatest epic coul...	film greatest epic earn boc worth budget shame...	en

10000 rows × 5 columns

In [49]: 1 df1 = df.copy()

In [50]: 1 df1['english_tweets'] = df[df['Language'] == 'en']['Cleaned_Tweets']

In [51]: 1 df1

Out[51]:

	Date Created	Number of Likes	Tweets	Cleaned_Tweets	Language	englisht
0	2023-06-30 09:21:00+00:00	0	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipuru...	womens ashes live streaming broadcast tv chann...	en	wome live s bro
1	2023-06-30 09:20:57+00:00	0	Now Playing!! Book Your Ticket Now!! 🎟️🍿🎥\n@go...	playing book ticket	en	play
2	2023-06-30 09:20:22+00:00	0	@ponilemova #Adipurush		None	
3	2023-06-30 09:20:00+00:00	3	Adipurush VS 72 Hoorain VS The Kerala Story Co...	adipurush vs hoorain vs kerala story controver...	en	adipurush kecc
4	2023-06-30 09:15:22+00:00	3	ST: #Adipurush https://t.co/lsgKcqQuKL	st	no	
...	
9996	2023-06-23 10:08:50+00:00	0	Adipurush 1st Week WW Box Office Collections: ...	adipurush 1st week ww box office collections	en	adipurush wee collect
9997	2023-06-23 10:08:49+00:00	0	#GodMorningFriday\nवास्तव में #Adipurush यानि ...	वास्तव म यन सबस पहल भगवन जसन सरव सषट क रचन क ह व...	hi	
9998	2023-06-23 10:08:17+00:00	3101	Let the empowering lyrics of #Shivoham elevate...	empowering lyrics elevate spirit envelop world...	en	em lyric spiri
9999	2023-06-23 10:08:01+00:00	0	When it comes to choosing a service or product...	choosing service product beneficial opt authen...	en	service benef
10000	2023-06-23 10:07:45+00:00	0	A film about #Ramayana, our greatest epic coul...	film greatest epic earn poc worth budget shame...	en	film epic worl

10000 rows × 6 columns



In [52]: 1 df1 = df1.dropna()

In [53]: 1 df1

Out[53]:

	Date Created	Number of Likes	Tweets	Cleaned_Tweets	Language	english_tweets
0	2023-06-30 09:21:00+00:00	0	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipuru...	womens ashes live streaming broadcast tv chann...	en	womens ashes live streaming broadcast tv chann...
1	2023-06-30 09:20:57+00:00	0	Now Playing!! Book Your Ticket Now!! 🎟️🍿🎟️\n@go...	playing book ticket	en	playing book ticket
3	2023-06-30 09:20:00+00:00	3	Adipurush VS 72 Hoorain VS The Kerala Story Co...	adipurush vs hoorain vs kerala story controver...	en	adipurush vs hoorain vs kerala story controver...
5	2023-06-30 09:08:27+00:00	1	This is how the story should be told. @omraut ...	story told learn hotstar india graphic india g...	en	story told learn hotstar india graphic india g...
8	2023-06-30 09:04:09+00:00	0	@VikasAgarwall Milord says: If my compatriots...	milord compatriots backstab ie end exposing fa...	en	milord compatriots backstab ie end exposing fa...
...
9995	2023-06-23 10:09:41+00:00	1	S Rangarajan garu, main poojari of chilkur bal...	rangarajan garu poojari chilkur balaji appreci...	en	rangarajan garu poojari chilkur balaji appreci...
9996	2023-06-23 10:08:50+00:00	0	Adipurush 1st Week WW Box Office Collections: ...	adipurush 1st week ww box office collections ☺...	en	adipurush 1st week ww box office collections ☺...
9998	2023-06-23 10:08:17+00:00	3101	Let the empowering lyrics of #Shivoham elevate...	empowering lyrics elevate spirit envelop world...	en	empowering lyrics elevate spirit envelop world...
9999	2023-06-23 10:08:01+00:00	0	When it comes to choosing a service or product...	choosing service product beneficial opt authen...	en	choosing service product beneficial opt authen...
10000	2023-06-23 10:07:45+00:00	0	A film about #Ramayana, our greatest epic coul...	film greatest epic earn boc worth budget shame...	en	film greatest epic earn boc worth budget shame...

5088 rows × 6 columns

```
In [54]: 1 df1['Year'] = df1['Date Created'].dt.year
2 df1['Month'] = df1['Date Created'].dt.month
3 df1['Day'] = df1['Date Created'].dt.day
```

```
In [55]: 1 df1
```

Out[55]:

	Date Created	Number of Likes	Tweets	Cleaned_Tweets	Language	english_tweets
0	2023-06-30 09:21:00+00:00	0	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipuru...	womens ashes live streaming broadcast tv chann...	en	womens ashes live streaming broadcast tv chann...
1	2023-06-30 09:20:57+00:00	0	Now Playing!! Book Your Ticket Now!! 🎟️🎟️\n@go...	playing book ticket	en	playing book ticket
3	2023-06-30 09:20:00+00:00	3	Adipurush VS 72 Hoorain VS The Kerala Story Co...	adipurush vs hoorain vs kerala story controver...	en	adipurush vs hoorain vs kerala story controver...
5	2023-06-30 09:08:27+00:00	1	This is how the story should be told. @omraut ...	story told learn hotstar india graphic india g...	en	story told learn hotstar india graphic india g...
8	2023-06-30 09:04:09+00:00	0	@VikasAgarwalli Milord says: If my compatriots...	milord compatriots backstab ie end exposing fa...	en	milord compatriots backstab ie end exposing fa...
...
9995	2023-06-23 10:09:41+00:00	1	S Rangarajan garu, main poojari of chilkur bal...	rangarajan garu poojari chilkur balaji appreci...	en	rangarajan garu poojari chilkur balaji appreci...
9996	2023-06-23 10:08:50+00:00	0	Adipurush 1st Week WW Box Office Collections: ...	adipurush 1st week ww box office collections ☺...	en	adipurush 1st week ww box office collections ☺...
9998	2023-06-23 10:08:17+00:00	3101	Let the empowering lyrics of #Shivoham elevate...	empowering lyrics elevate spirit envelop world...	en	empowering lyrics elevate spirit envelop world...
9999	2023-06-23 10:08:01+00:00	0	When it comes to choosing a service or product...	choosing service product beneficial opt authen...	en	choosing service product beneficial opt authen...
10000	2023-06-23 10:07:45+00:00	0	A film about #Ramayana, our greatest epic coul...	film greatest epic earn boc worth budget shame...	en	film greatest epic earn boc worth budget shame...

5088 rows × 9 columns

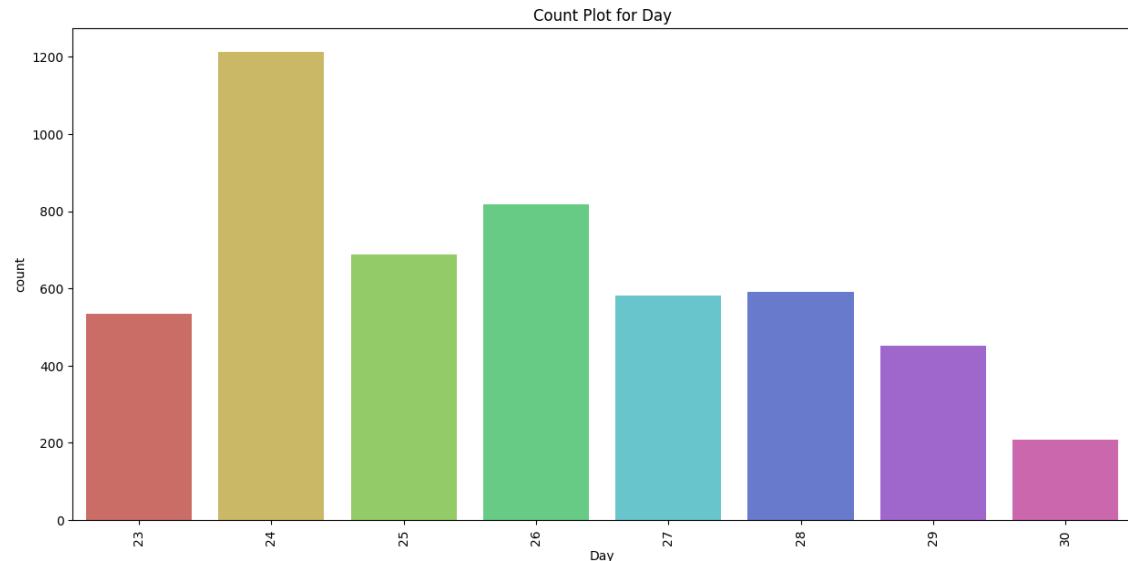
```
In [56]: 1 df1.nunique()
```

```
Out[56]: Date Created      5051
Number of Likes       555
Tweets                 5041
Cleaned_Tweets        4456
Language                1
english_tweets         4456
Year                     1
Month                    1
Day                      8
dtype: int64
```

```
In [57]: 1 df1['Time'] = df1['Date Created'].dt.time
```

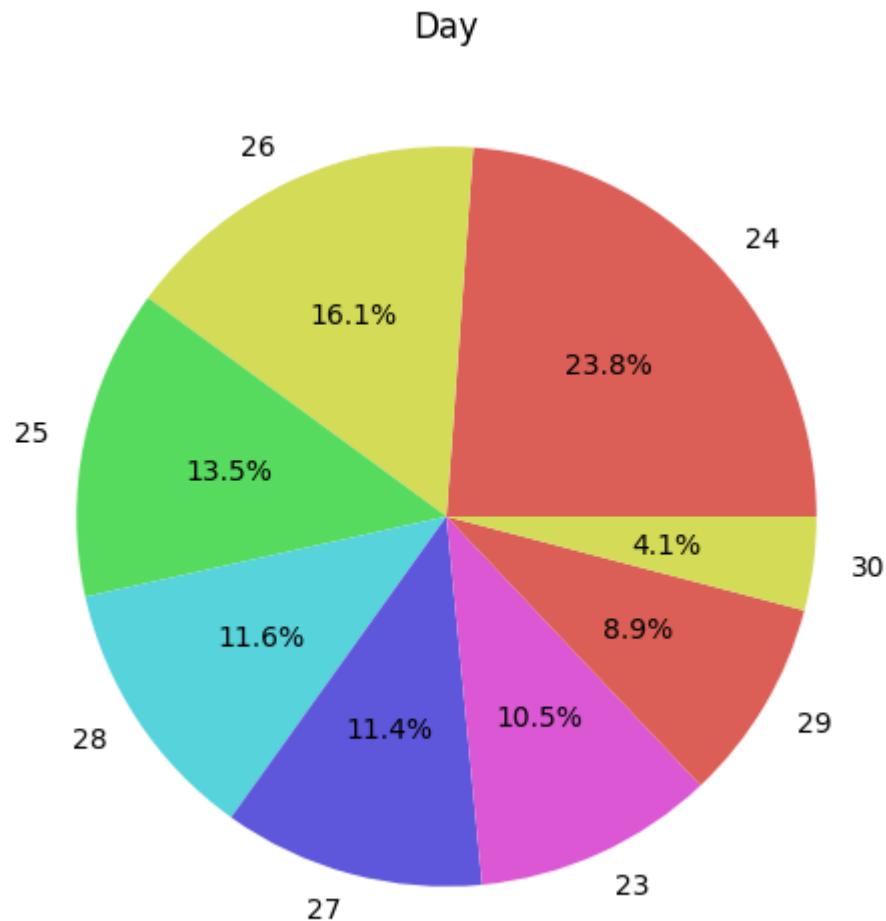
```
In [58]: 1 df1['Tweet_Length'] = df1['english_tweets'].str.len()
```

```
In [59]: 1 plt.figure(figsize=[15,7],)
2 plt.title('Count Plot for Day')
3 sns.countplot(x = 'Day', data = df1, palette = 'hls')
4 plt.xticks(rotation = 90)
5 plt.show()
```



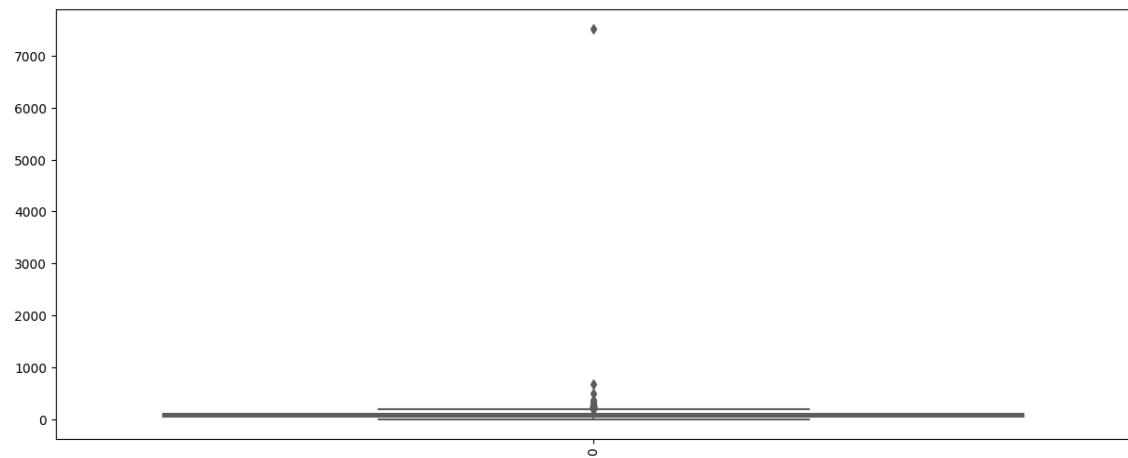
In [60]:

```
1 plt.figure(figsize=(15, 6))
2 counts = df1['Day'].value_counts()
3 plt.pie(counts, labels=counts.index, autopct='%1.1f%%', colors=sns.col
4 plt.title('Day')
5 plt.show()
```



In [69]:

```
1 plt.figure(figsize=(15,6))
2 sns.boxplot(df1['Tweet_Length'], palette='hls')
3 plt.xticks(rotation=90)
4 plt.show()
```



```
In [70]: ┌─ 1 def label_sentiment(x:float):
      2     if x < -0.05 : return 'negative'
      3     if x > 0.35 : return 'positive'
      4     return 'neutral'
```

```
In [71]: ┌─ 1 sia = SIA()
```

```
In [72]: ┌─ 1 df1['sentiment'] = [sia.polarity_scores(x)['compound'] for x in tqdm(d
      2 df1['overall_sentiment'] = df1['sentiment'].apply(label_sentiment);
```

100%

5088/5088 [00:01<00:00, 3896.58it/s]

In [74]: 1 df1

Out[74]:

	Date Created	Number of Likes	Tweets	Cleaned_Tweets	Language	english_tweets
0	2023-06-30 09:21:00+00:00	0	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipuru...	womens ashes live streaming broadcast tv chann...	en	womens ashes live streaming broadcast tv chann...
1	2023-06-30 09:20:57+00:00	0	Now Playing!! Book Your Ticket Now!! 🎟️🎟️🎟️\n@go...	playing book ticket	en	playing book ticket
3	2023-06-30 09:20:00+00:00	3	Adipurush VS 72 Hoorain VS The Kerala Story Co...	adipurush vs hoorain vs kerala story controver...	en	adipurush vs hoorain vs kerala story controver...
5	2023-06-30 09:08:27+00:00	1	This is how the story should be told. @omraut ...	story told learn hotstar india graphic india g...	en	story told learn hotstar india graphic india g...
8	2023-06-30 09:04:09+00:00	0	@VikasAgarwall Milord says: If my compatriots...	milord compatriots backstab ie end exposing fa...	en	milord compatriots backstab ie end exposing fa...
...
9995	2023-06-23 10:09:41+00:00	1	S Rangarajan garu, main poojari of chilkur bal...	rangarajan garu poojari chilkur balaji appreci...	en	rangarajan garu poojari chilkur balaji appreci...
9996	2023-06-23 10:08:50+00:00	0	Adipurush 1st Week WW Box Office Collections: ...	adipurush 1st week ww box office collections ☺...	en	adipurush 1st week ww box office collections ☺...
9998	2023-06-23 10:08:17+00:00	3101	Let the empowering lyrics of #Shivoham elevate...	empowering lyrics elevate spirit envelop world...	en	empowering lyrics elevate spirit envelop world...
9999	2023-06-23 10:08:01+00:00	0	When it comes to choosing a service or product...	choosing service product beneficial opt authen...	en	choosing service product beneficial opt authen...
10000	2023-06-23 10:07:45+00:00	0	A film about #Ramayana, our greatest epic coul...	film greatest epic earn boc worth budget shame...	en	film greatest epic earn boc worth budget shame...

5088 rows × 13 columns

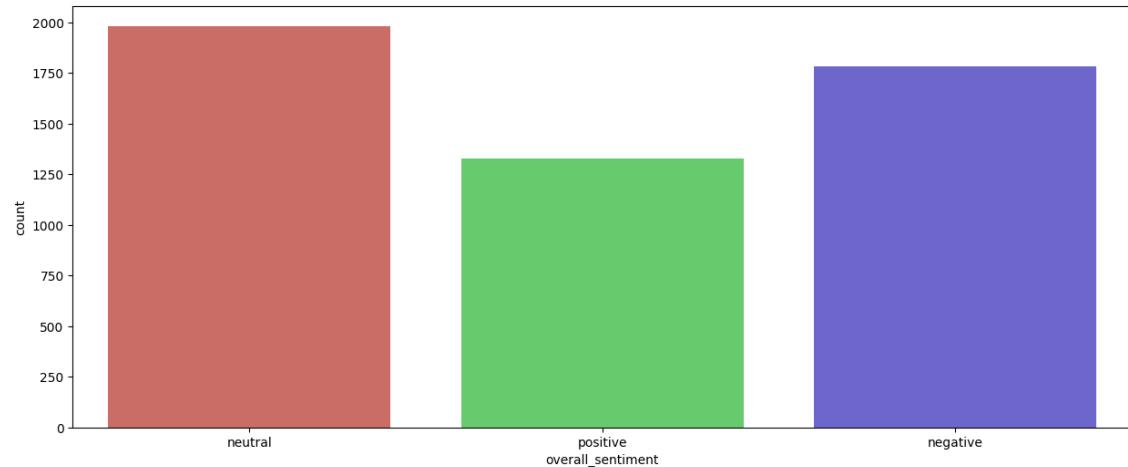
In [75]: 1 df1['overall_sentiment'].unique()

Out[75]: array(['neutral', 'positive', 'negative'], dtype=object)

```
In [76]: 1 df1['overall_sentiment'].value_counts()
```

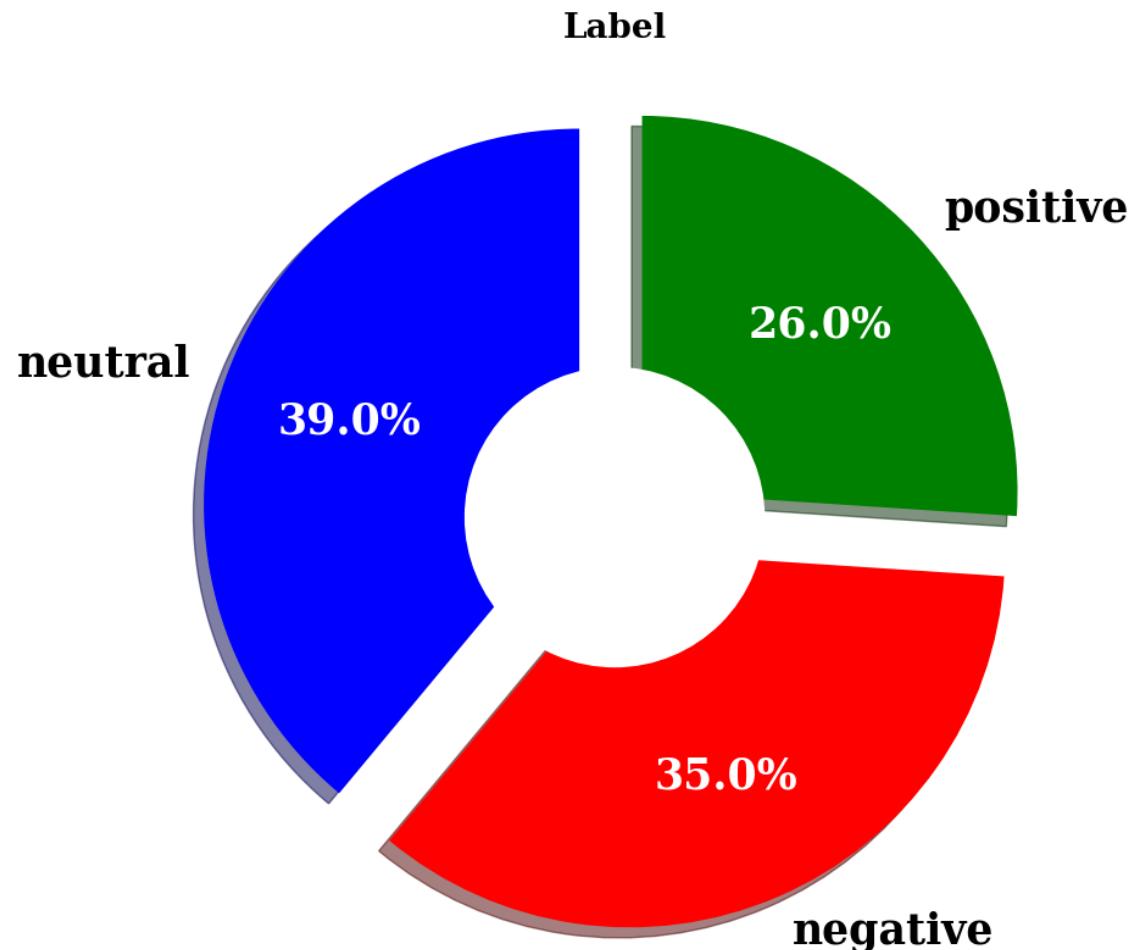
```
Out[76]: neutral    1982  
negative   1781  
positive    1325  
Name: overall_sentiment, dtype: int64
```

```
In [77]: 1 plt.figure(figsize=(15,6))  
2 sns.countplot(x='overall_sentiment', data = df1, palette = 'hls')  
3 plt.xticks(rotation = 0)  
4 plt.show()
```



In [78]:

```
1 label_data = df1['overall_sentiment'].value_counts()
2
3 explode = (0.1, 0.1, 0.1)
4 plt.figure(figsize=(14, 10))
5 patches, texts, pcts = plt.pie(label_data,
6                                 labels = label_data.index,
7                                 colors = ['blue', 'red', 'green'],
8                                 pctdistance = 0.65,
9                                 shadow = True,
10                                startangle = 90,
11                                explode = explode,
12                                autopct = '%1.1f%%',
13                                textprops={ 'fontsize': 25,
14                                            'color': 'black',
15                                            'weight': 'bold',
16                                            'family': 'serif' })
17 plt.setp(pcts, color='white')
18
19 hfont = {'fontname':'serif', 'weight': 'bold'}
20 plt.title('Label', size=20, **hfont)
21
22 centre_circle = plt.Circle((0,0),0.40,fc='white')
23 fig = plt.gcf()
24 fig.gca().add_artist(centre_circle)
25 plt.show()
```



In [79]: 1 df1

Out[79]:

	Date Created	Number of Likes	Tweets	Cleaned_Tweets	Language	english_tweets
0	2023-06-30 09:21:00+00:00	0	#ENGvAUS #ENGvsAUS #AUSvENG #AUSvsENG #Adipuru...	womens ashes live streaming broadcast tv chann...	en	womens ashes live streaming broadcast tv chann...
1	2023-06-30 09:20:57+00:00	0	Now Playing!! Book Your Ticket Now!! 🎟️🍿🎟️\n@go...	playing book ticket	en	playing book ticket
3	2023-06-30 09:20:00+00:00	3	Adipurush VS 72 Hoorain VS The Kerala Story Co...	adipurush vs hoorain vs kerala story controver...	en	adipurush vs hoorain vs kerala story controver...
5	2023-06-30 09:08:27+00:00	1	This is how the story should be told. @omraut ...	story told learn hotstar india graphic india g...	en	story told learn hotstar india graphic india g...
8	2023-06-30 09:04:09+00:00	0	@VikasAgarwall Milord says: If my compatriots...	milord compatriots backstab ie end exposing fa...	en	milord compatriots backstab ie end exposing fa...
...
9995	2023-06-23 10:09:41+00:00	1	S Rangarajan garu, main poojari of chilkur bal...	rangarajan garu poojari chilkur balaji appreci...	en	rangarajan garu poojari chilkur balaji appreci...
9996	2023-06-23 10:08:50+00:00	0	Adipurush 1st Week WW Box Office Collections: ...	adipurush 1st week ww box office collections ☺...	en	adipurush 1st week ww box office collections ☺...
9998	2023-06-23 10:08:17+00:00	3101	Let the empowering lyrics of #Shivoham elevate...	empowering lyrics elevate spirit envelop world...	en	empowering lyrics elevate spirit envelop world...
9999	2023-06-23 10:08:01+00:00	0	When it comes to choosing a service or product...	choosing service product beneficial opt authen...	en	choosing service product beneficial opt authen...
10000	2023-06-23 10:07:45+00:00	0	A film about #Ramayana, our greatest epic coul...	film greatest epic earn boc worth budget shame...	en	film greatest epic earn boc worth budget shame...

5088 rows × 13 columns

In [80]: 1 df2 = df1[['english_tweets', 'overall_sentiment']]

In [81]: 1 df2

Out[81]:

	english_tweets	overall_sentiment
0	womens ashes live streaming broadcast tv chann...	neutral
1	playing book ticket	neutral
3	adipurush vs hoorain vs kerala story controver...	neutral
5	story told learn hotstar india graphic india g...	neutral
8	milord compatriots backstab ie end exposing fa...	positive
...
9995	rangarajan garu poojari chilkur balaji appreci...	positive
9996	adipurush 1st week ww box office collections ☺...	neutral
9998	empowering lyrics elevate spirit envelop world...	positive
9999	choosing service product beneficial opt authen...	positive
10000	film greatest epic earn boc worth budget shame...	positive

5088 rows × 2 columns

In [82]: 1

```

1 def clean_text(text):
2     # Remove non-alphabetic characters and convert to Lowercase
3     cleaned_text = re.sub('[^a-zA-Z]', ' ', text).lower()
4     # Remove extra white spaces
5     cleaned_text = re.sub('\s+', ' ', cleaned_text).strip()
6     # Split the text into words
7     words = cleaned_text.split()
8     # Join the words back into a string
9     cleaned_text = ' '.join(words)
10    return cleaned_text
11
12 # Apply the clean_text function to the 'english_tweets' column
13 df2['Cleaned_English_Tweets'] = df2['english_tweets'].apply(clean_text)

```

In [83]: 1 df2

Out[83]:

	english_tweets	overall_sentiment	Cleaned_English_Tweets
0	womens ashes live streaming broadcast tv chann...	neutral	womens ashes live streaming broadcast tv chann...
1	playing book ticket	neutral	playing book ticket
3	adipurush vs hoorain vs kerala story controver...	neutral	adipurush vs hoorain vs kerala story controver...
5	story told learn hotstar india graphic india g...	neutral	story told learn hotstar india graphic india g...
8	milord compatriots backstab ie end exposing fa...	positive	milord compatriots backstab ie end exposing fa...
...
9995	rangarajan garu poojari chilkur balaji appreci...	positive	rangarajan garu poojari chilkur balaji appreci...
9996	adipurush 1st week ww box office collections ↗...	neutral	adipurush st week ww box office collections st
9998	empowering lyrics elevate spirit envelop world...	positive	empowering lyrics elevate spirit envelop world...
9999	choosing service product beneficial opt authen...	positive	choosing service product beneficial opt authen...
10000	film greatest epic earn boc worth budget shame...	positive	film greatest epic earn boc worth budget shame...

5088 rows × 3 columns

In [84]: 1 df3 = df2[['Cleaned_English_Tweets', 'overall_sentiment']]

In [85]: 1 df3

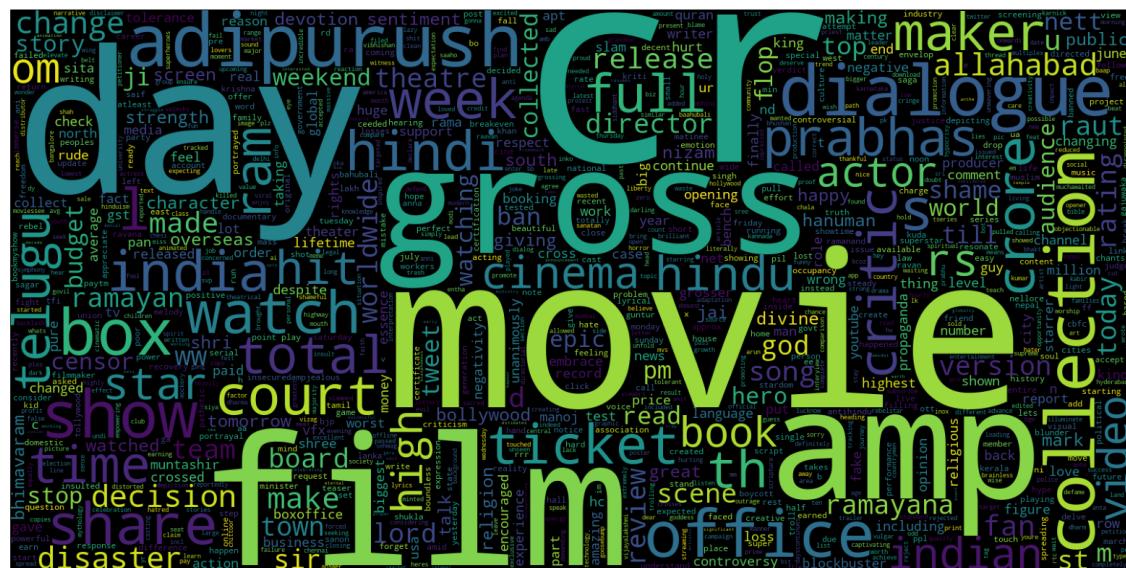
Out[85]:

	Cleaned_English_Tweets	overall_sentiment
0	womens ashes live streaming broadcast tv chann...	neutral
1	playing book ticket	neutral
3	adipurush vs hoorain vs kerala story controver...	neutral
5	story told learn hotstar india graphic india g...	neutral
8	milord compatriots backstab ie end exposing fa...	positive
...
9995	rangarajan garu poojari chilkur balaji appreci...	positive
9996	adipurush st week ww box office collections st	neutral
9998	empowering lyrics elevate spirit envelop world...	positive
9999	choosing service product beneficial opt authen...	positive
10000	film greatest epic earn boc worth budget shame...	positive

5088 rows × 2 columns

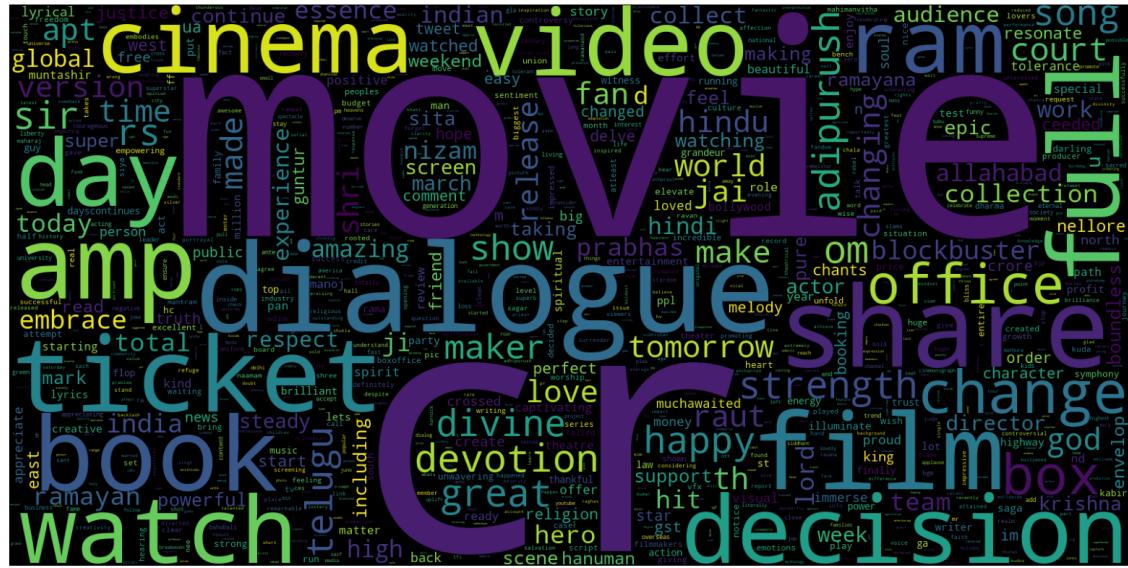
In [86]: ► 1 import wordcloud

```
In [87]: 1 from wordcloud import WordCloud  
2 data = df3['Cleaned_English_Tweets']  
3 plt.figure(figsize = (20,20))  
4 wc = WordCloud(max_words = 1000 , w  
5                 collocations=False).g  
6 plt.imshow(wc)  
7 plt.axis('off')  
8 plt.show()
```



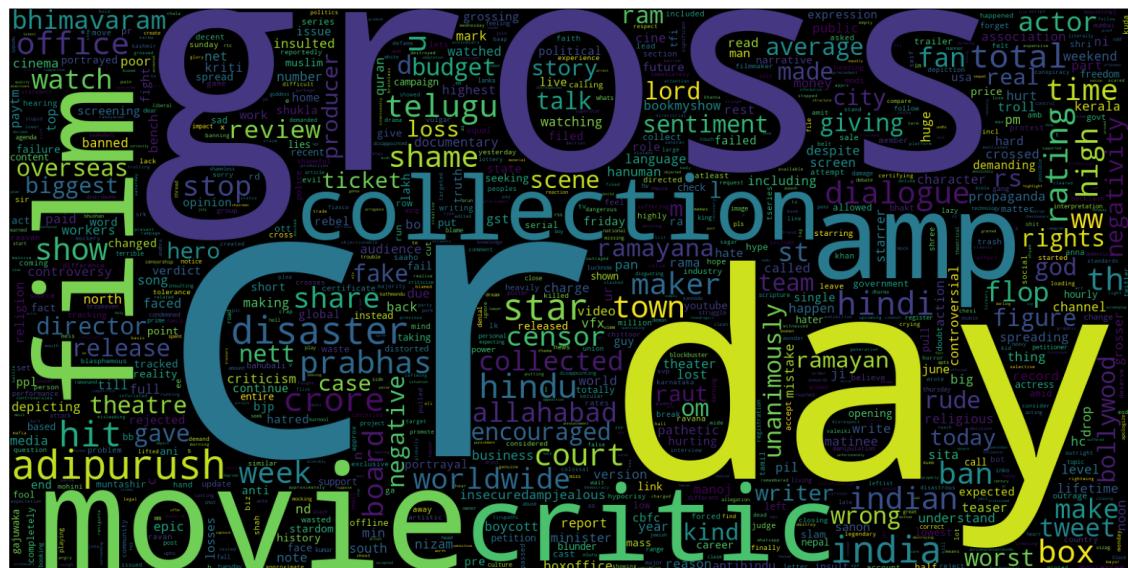
In [88]: ►

```
1 data = df3[df3['overall_sentiment']=="positive"]['Cleaned_English_Tweet']
2 plt.figure(figsize = (20,20))
3 wc = WordCloud(max_words = 1000 , width = 1600 , height = 800,
4                 collocations=False).generate(" ".join(data))
5 plt.imshow(wc)
6 plt.axis('off')
7 plt.show()
```

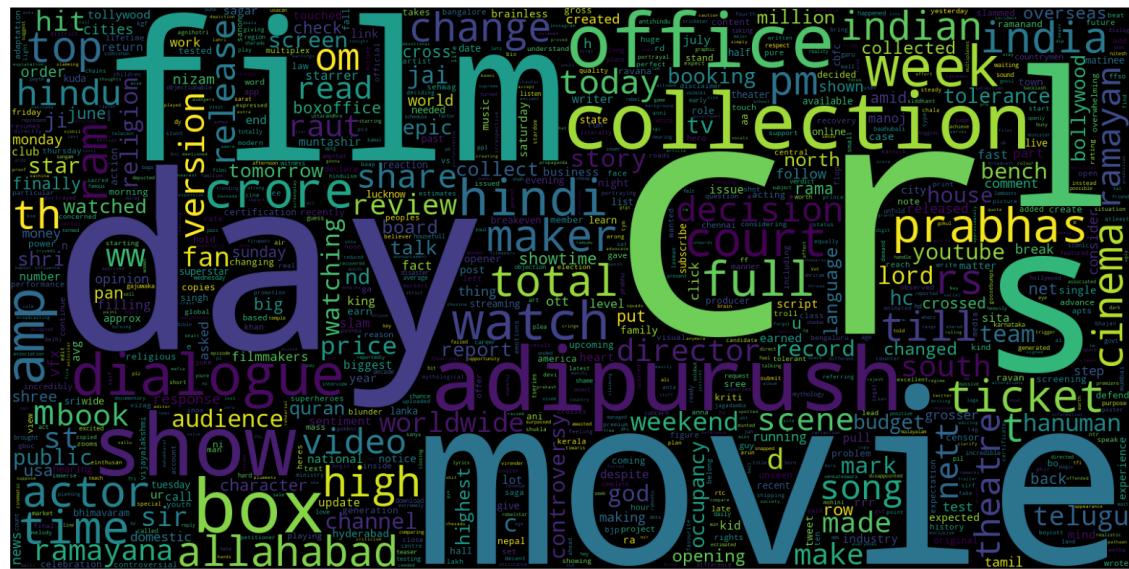


In [89]:

```
1 data = df3[df3['overall_sentiment']=="negative"]['Cleaned_English_Tweet']
2 plt.figure(figsize = (20,20))
3 wc = WordCloud(max_words = 1000 , width = 1600 , height = 800,
4                 collocations=False).generate(" ".join(data))
5 plt.imshow(wc)
6 plt.axis('off')
7 plt.show()
```



```
In [90]: 1 data = df3[df3['overall_sentiment']=="neutral"]['Cleaned_English_Tweet']
2 plt.figure(figsize = (20,20))
3 wc = WordCloud(max_words = 1000 , width = 1600 , height = 800,
4                 collocations=False).generate(" ".join(data))
5 plt.imshow(wc)
6 plt.axis('off')
7 plt.show()
```



```
In [91]: # Load the libraries required for performing classification
          1 from sklearn.naive_bayes import MultinomialNB
          2 from sklearn.metrics import confusion_matrix, accuracy_score, f1_score
          3 from sklearn.model_selection import train_test_split
          4 from nltk.tokenize import TweetTokenizer
          5 from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer
          6 from sklearn.feature_extraction.text import TfidfVectorizer
```

```
In [92]: # Split the data into training and testing data sets
          # Use cleaned english as independent variable and overall sentiment as
          # dependent variable
          X = df3["Cleaned_English_Tweets"].values
          y = df3['overall_sentiment'].values
          X_train, X_test, y_train, y_test = train_test_split(X,y, random_state=
```

```
In [93]: # Extract features using TFIDF Vectorizer
          1
          2
          3 vectorizer = TfidfVectorizer(max_features=1000)
          4 X_train_idf = vectorizer.fit_transform(X_train)
          5 X_test_idf = vectorizer.transform(X_test)
```

In [94]:

```

1 # Print idf values
2 df_idf = pd.DataFrame(vectorizer.idf_, index=vectorizer.get_feature_na
3 # Sort ascending
4 df_idf.sort_values(by=['idf_weights'], ascending = False).head()

```

Out[94]:

	idf_weights
narratives	7.791783
kashmir	7.791783
artistic	7.791783
sree	7.791783
lks	7.568640

In [95]:

```

1 # Perform Multinomial Naive Bayes Classification
2 # Apply MultinomialNB on training data
3 mnb = MultinomialNB()
4 mnb.fit(X_train_idf, y_train)

```

Out[95]:

```

▼ MultinomialNB
MultinomialNB()

```

In [96]:

```

1 # Predict polarity by fitting the model to testing data
2 pred_mnb = mnb.predict(X_test_idf)
3
4 # Calculate accuracy of predicted values
5 acc = accuracy_score(y_test, pred_mnb)
6
7
8 results = pd.DataFrame([['Multinomial Naive Bayes', acc]],
9                         columns = ['Model', 'Accuracy'])
10
11 print(results)

```

	Model	Accuracy
0	Multinomial Naive Bayes	0.667322

In [97]:

```

1 # Perform Random Forest classification on the processed data and compare
2
3 # Random Forest Classifier with 'gini'
4
5 from sklearn.ensemble import RandomForestClassifier
6 clf_rf = RandomForestClassifier()
7 clf_rf.fit(X_train_idf, y_train)
8
9 # Predict using testing data
10 y_pred_rf = clf_rf.predict(X_test_idf)
11
12 # Calculate accuracy
13 acc = accuracy_score(y_test, y_pred_rf)
14
15 model_results = pd.DataFrame([['Random Forest(Gini)', acc]],
16                             columns = ['Model', 'Accuracy'])
17
18 results = results.append(model_results, ignore_index = True)
19 print(results)

```

	Model	Accuracy
0	Multinomial Naive Bayes	0.667322
1	Random Forest(Gini)	0.747217

In [98]:

```

1 # Random Forest Classifier with 'entropy'
2
3 from sklearn.ensemble import RandomForestClassifier
4 clf_rf = RandomForestClassifier(criterion='entropy')
5 clf_rf.fit(X_train_idf, y_train)
6
7 # Predict using testing data
8 y_pred_rf = clf_rf.predict(X_test_idf)
9
10 # Calculate accuracy
11 acc = accuracy_score(y_test, y_pred_rf)
12
13 model_results = pd.DataFrame([['Random Forest(Entropy)', acc]],
14                             columns = ['Model', 'Accuracy'])
15
16 results = results.append(model_results, ignore_index = True)
17 print(results)

```

	Model	Accuracy
0	Multinomial Naive Bayes	0.667322
1	Random Forest(Gini)	0.747217
2	Random Forest(Entropy)	0.747872

In [99]:

```

1 # Display confusion matrix for Random Forest
2
3 confusion_matrix(y_test,y_pred_rf) ### Confusion matrix for Random For

```

Out[99]:

```
array([[386, 124,    7],
       [ 46, 510,   46],
       [ 32, 130, 246]], dtype=int64)
```

In []: █ 1

In []: █ 1