

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from datetime import datetime

In [2]: df = pd.read_csv('D:\Vinayak\INTER\USVideos.csv')

In [3]: df.head()

Out[3]:
```

	video_id	trending_date	title	channel_title	category_id	publish_time	tags	views	likes	dislikes	comment_count	thumbnail_link	comments_disabled	ratings_disabled	video_error
0	2y5S6svSYSE	17-11-11	WE WANT TO TALK ABOUT OUR MARRIAGE	CaseyNeistat	22	2017-11-13T17:13:01.000Z	SHANeNeil marlin	748374	57527	2966	15954	https://i.ytimg.com/vi/2y5S6svSYSE/default.jpg	False	False	
1	1ZAPwHAFY	17-11-11	The Trump Presidency: Last Week Tonight with J...	LastWeekTonight	24	2017-11-13T07:30:00.000Z	last week tonight trump presidency"last week ...	2418783	97185	6146	12703	https://i.ytimg.com/vi/1ZAPwHAFY/default.jpg	False	False	
2	5qgK5DqJC4	17-11-11	Racist Superman Rudy Mancuso, King Bach & Le...	Rudy Mancuso	23	2017-11-19 05:24:00Z	superman"rudy"mancuso"king"bach"...	3191434	146033	5339	8181	https://i.ytimg.com/vi/5qgK5DqJC4/default.jpg	False	False	
3	puqWIECTY	17-11-11	Nickelback Lyrics: Real or Fake?	Good Mythical Morning	24	2017-11-13T11:00:04.000Z	rhet and link"gmmt""good mythical morning"Ti...	343168	10172	666	2146	https://i.ytimg.com/vi/puqWIECTY/default.jpg	False	False	
4	4B8OmEODVM	17-11-11	I Dare You: GONG BALDI?	rigahiga	24	2017-11-12T18:01:41.000Z	ryan"thga"ThgaW"ThgaWga"Ti dare you"Ti...	2095731	132235	1989	17518	https://i.ytimg.com/vi/4B8OmEODVM/default.jpg	False	False	

```
In [4]: df.shape
Out[4]: (48949, 16)

In [5]: dfn = df.drop_duplicates()
dfn.shape
Out[5]: (48981, 16)

In [6]: dfn.describe()

Out[6]:
```

	category_id	views	likes	dislikes	comment_count
count	40901.000000	4.090100e+04	4.090100e+04	4.090100e+04	4.090100e+04
mean	19.870588	2.306078e+06	7.427173e+04	3.171722e+03	8.484956e+03
std	7.560362	7.397719e+06	2.289999e+05	2.904624e+04	3.745139e+04
min	3.000000	5.490000e+02	0.000000e+00	0.000000e+00	0.000000e+00
25%	17.000000	2.418702e+05	5.480000e+03	2.020000e+02	6.130000e+02
50%	24.000000	6.516040e+05	1.699000e+04	6.300000e+02	1.855000e+03
75%	25.000000	1.821820e+06	5.532800e+04	1.890000e+03	5.752000e+03
max	43.000000	2.752119e+08	5.613827e+06	1.674420e+06	1.361900e+06

```
In [7]: dfn.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 48981 entries, 0 to 48948
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  --
0   video_id              48981 non-null    object
1   trending_date         48981 non-null    object
2   title                 48981 non-null    object
3   channel_title         48981 non-null    object
4   category_id           48981 non-null    int64
5   publish_time          48981 non-null    object
6   tags                  48981 non-null    object
7   views                 48981 non-null    int64
8   likes                 48981 non-null    int64
9   dislikes              48981 non-null    int64
10  comment_count         48981 non-null    int64
11  thumbnail_link        48981 non-null    object
12  comments_disabled     48981 non-null    bool
13  ratings_disabled      48981 non-null    bool
14  video_error_or_remed  48981 non-null    bool
15  description           48332 non-null    object
dtypes: bool(3), int64(5), object(8)
memory usage: 4.5+ MB

In [8]: columns_to_remove = ['thumbnail_link', 'description']
dfn = dfn.drop(columns=columns_to_remove)
dfn.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 48981 entries, 0 to 48948
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  --
0   video_id              48981 non-null    object
1   trending_date         48981 non-null    object
2   title                 48981 non-null    object
3   channel_title         48981 non-null    object
4   category_id           48981 non-null    int64
5   publish_time          48981 non-null    object
6   tags                  48981 non-null    object
7   views                 48981 non-null    int64
8   likes                 48981 non-null    int64
9   dislikes              48981 non-null    int64
10  comment_count         48981 non-null    int64
11  comments_disabled     48981 non-null    bool
12  ratings_disabled      48981 non-null    bool
13  video_error_or_remed  48981 non-null    bool
dtypes: bool(3), int64(5), object(6)
memory usage: 3.9+ MB

In [9]: from datetime import datetime

In [10]: dfn['trending_date'] = dfn['trending_date'].apply(lambda x: datetime.datetime.strptime(x, '%y.%d.%a'))

In [11]: dfn.head(3)

Out[11]:
```

	video_id	trending_date	title	channel_title	category_id	publish_time	tags	views	likes	dislikes	comment_count	comments_disabled	ratings_disabled	video_error	removed
0	2y5S6svSYSE	2017-11-14	WE WANT TO TALK ABOUT OUR MARRIAGE	CaseyNeistat	22	2017-11-13T17:13:01.000Z	SHANeNeil marlin	748374	57527	2966	15954	False	False	False	False
1	1ZAPwHAFY	2017-11-14	The Trump Presidency: Last Week Tonight with J...	LastWeekTonight	24	2017-11-13T07:30:00.000Z	last week tonight trump presidency"last week ...	2418783	97185	6146	12703	False	False	False	False
2	5qgK5DqJC4	2017-11-14	Racist Superman Rudy Mancuso, King Bach & Le...	Rudy Mancuso	23	2017-11-12T19:05:24.000Z	superman"rudy"mancuso"king"bach"...	3191434	146033	5339	8181	False	False	False	False

```
In [12]: dfn['publish_time'] = pd.to_datetime(dfn['publish_time'])
dfn.head(2)

Out[12]:
```

	video_id	trending_date	title	channel_title	category_id	publish_time	tags	views	likes	dislikes	comment_count	comments_disabled	ratings_disabled	video_error	removed
0	2y5S6svSYSE	2017-11-14	WE WANT TO TALK ABOUT OUR MARRIAGE	CaseyNeistat	22	2017-11-13 17:13:01+00:00	SHANeNeil marlin	748374	57527	2966	15954	False	False	False	False
1	1ZAPwHAFY	2017-11-14	The Trump Presidency: Last Week Tonight with J...	LastWeekTonight	24	2017-11-13 07:30:00+00:00	last week tonight trump presidency"last week ...	2418783	97185	6146	12703	False	False	False	False

```
In [13]: dfn['publish_month'] = dfn['publish_time'].dt.month
dfn['publish_day'] = dfn['publish_time'].dt.day
dfn['publish_hour'] = dfn['publish_time'].dt.hour

In [14]: dfn.head(2)

Out[14]:
```

	video_id	trending_date	title	channel_title	category_id	publish_time	tags	views	likes	dislikes	comment_count	comments_disabled	ratings_disabled	video_error	removed	publish_month	publish_day	publish_hour
0	2y5S6svSYSE	2017-11-14	WE WANT TO TALK ABOUT OUR MARRIAGE	CaseyNeistat	22	2017-11-13 17:13:01+00:00	SHANeNeil marlin	748374	57527	2966	15954	False	False	False	False	11	13	17
1	1ZAPwHAFY	2017-11-14	The Trump Presidency: Last Week Tonight with J...	LastWeekTonight	24	2017-11-13 07:30:00+00:00	last week tonight trump presidency"last week ...	2418783	97185	6146	12703	False	False	False	False	11	13	7

```
In [15]: print(sorted(dfn["category_id"].unique()))
[1, 2, 10, 15, 17, 19, 20, 22, 23, 24, 25, 26, 27, 28, 29, 43]

In [16]: dfn['category_name'] = np.nan

dfn.loc[dfn['category_id'] == 1], 'category_name' = 'Film and Animation'
dfn.loc[dfn['category_id'] == 2], 'category_name' = 'Autos and Vehicles'
dfn.loc[dfn['category_id'] == 10], 'category_name' = 'Music'
dfn.loc[dfn['category_id'] == 15], 'category_name' = 'Pets and Animals'
dfn.loc[dfn['category_id'] == 17], 'category_name' = 'Sports'
dfn.loc[dfn['category_id'] == 19], 'category_name' = 'Travel and Events'
dfn.loc[dfn['category_id'] == 20], 'category_name' = 'Gaming'
dfn.loc[dfn['category_id'] == 22], 'category_name' = 'People and Blogs'
dfn.loc[dfn['category_id'] == 23], 'category_name' = 'Comedy'
dfn.loc[dfn['category_id'] == 24], 'category_name' = 'Entertainment'
dfn.loc[dfn['category_id'] == 25], 'category_name' = 'News and Politics'
dfn.loc[dfn['category_id'] == 26], 'category_name' = 'How to and Style'
dfn.loc[dfn['category_id'] == 27], 'category_name' = 'Education'
dfn.loc[dfn['category_id'] == 28], 'category_name' = 'Science and Technology'
dfn.loc[dfn['category_id'] == 29], 'category_name' = 'Non Profits and Activism'
dfn.loc[dfn['category_id'] == 30], 'category_name' = 'Movies'
dfn.loc[dfn['category_id'] == 43], 'category_name' = 'Shows'

In [17]: dfn.head()

Out[17]:
```

	video_id	trending_date	title	channel_title	category_id	publish_time	tags	views	likes	dislikes	comment_count	comments_disabled	ratings_disabled	video_error	removed	publish_month	publish_day	publish_hour
0	2y5S6svSYSE	2017-11-14	WE WANT TO TALK ABOUT OUR MARRIAGE	CaseyNeistat	22	2017-11-13 17:13:01+00:00	SHANeNeil marlin	748374	57527	2966	15954	False	False	False	False	11	13	
1	1ZAPwHAFY	2017-11-14	The Trump Presidency: Last Week Tonight with J...	LastWeekTonight	24	2017-11-13 07:30:00+00:00	last week tonight trump presidency"last week ...	2418783	97185	6146	12703	False	False	False	False	11	13	
2	5qgK5DqJC4	2017-11-14	Racist Superman Rudy Mancuso, King Bach & Le...	Rudy Mancuso	23	2017-11-12 19:05:24+00:00	superman"rudy"mancuso"king"bach"...	3191434	146033	5339	8181	False	False	False	False	11	12	
3	puqWIECTY	2017-11-14	Nickelback Lyrics: Real or Fake?	Good Mythical Morning	24	2017-11-13 11:00:04+00:00	rhet and link"gmmt""good mythical morning"Ti...	343168	10172	666	2146	False	False	False	False	11	13	
4	4B8OmEODVM	2017-11-14	I Dare You: GONG BALDI?	rigahiga	24	2017-11-12 18:01:41+00:00	ryan"thga"ThgaW"ThgaWga"Ti dare you"Ti...	2095731	132235	1989	17518	False	False	False	False	11	12	

```
In [18]: dfn['year'] = dfn['publish_time'].dt.year
yearly_counts = dfn.groupby('year')['video_id'].count()
yearly_counts.plot(kind='bar', xlabel='Year', ylabel='Total Publish Year', title='Total Publish Video Per Year')
plt.show()

Total Publish Video Per Year

In [19]: yearly_views = dfn.groupby('year')['views'].sum()
yearly_views.plot(kind='bar', xlabel='Year', ylabel='Total views', title = 'Total Views per year')
plt.show()

Cell In[19], line 4
plt.
SyntaxError: invalid syntax

In [20]: yearly_views = dfn.groupby('year')['views'].sum()
yearly_views.plot(kind='bar', xlabel='Year', ylabel='Total views', title = 'Total views per year')
plt.xticks(rotation = 0)
plt.tight_layout()
plt.show()

Total views per year

In [21]: category_views = dfn.groupby('category_name')['views'].sum().reset_index()
#sort the categories by views in desc order
top_categories = category_views.sort_values(by='views', ascending=False).head(5)
#creating visualization by using matplotlib
plt.bar(top_categories['category_name'], top_categories['views'])
plt.xlabel('Category Name', fontsize = 12)
plt.xticks(rotation=90)
plt.ylabel('Total Views', fontsize = 12)
plt.title('Top 5 Categories', fontsize = 15)
plt.show()

Top 5 Categories

In [22]: plt.figure(figsize=(12,16))
sns.countplot(x='category_name', data=dfn, order=dfn['category_name'].value_counts().index)
plt.xticks(rotation=90)
plt.title('Video Count by category')
plt.show()

Video Count by category

In [23]: videos_per_hour = dfn['publish_hour'].value_counts().sort_index()
plt.figure(figsize=(25,8))
sns.barplot(x=videos_per_hour.index, y=videos_per_hour.values, palette='rocket')
plt.title('Number of videos published per hour')
plt.xlabel('hour of a day')
plt.ylabel('number of videos')
plt.xticks(rotation=6)
plt.show()

Number of videos published per hour

In [ ]:

In [24]: sns.scatterplot(data=dfn, x='views', y='likes')
plt.title('views vs likes')
plt.xlabel('views')
plt.ylabel('likes')
plt.show()

views vs likes

In [25]: plt.figure(figsize=(14,8))
plt.subplots_adjust(wspace = 0.2, hspace = 0.4, top=0.9)
g1 = sns.countplot(x='comments_disabled', data=dfn)
g1.set_title('Comments Disabled', fontsize=16)
plt.subplot(2,2,3)
g2 = sns.countplot(x='ratings_disabled', data=dfn)
g2.set_title('Ratings Disabled', fontsize=16)
plt.subplot(2,2,3)
g3 = sns.countplot(x='video_error_or_removed', data=dfn)
g3.set_title('Video Error or Removed', fontsize=16)
plt.show()

comments_disabled

video_error_or_removed
```

```
In [27]: corr_matrix[dfn['views'].corr(dfn['likes'])]
corr_matrix

Out[27]: 0.8491785476236589

In [ ]:
```