**This task involves performing exploratory data analysis on a dataset.Create**
⌄ **visualizations to understand the distribution of variables, identify outliers, and check for correlations between variables.**

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns


data = pd.read_csv("/content/USvideos.csv")


data.shape
```

```
(40949, 16)
```

```python
data = data.drop_duplicates()


data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 40901 entries, 0 to 40948
Data columns (total 16 columns):
 #   Column                 Non-Null Count  Dtype
---  ------                 --------------  -----
 0   video_id               40901 non-null  object
 1   trending_date          40901 non-null  object
 2   title                  40901 non-null  object
 3   channel_title          40901 non-null  object
 4   category_id            40901 non-null  int64
 5   publish_time           40901 non-null  object
 6   tags                   40901 non-null  object
 7   views                  40901 non-null  int64
 8   likes                  40901 non-null  int64
 9   dislikes               40901 non-null  int64
 10  comment_count          40901 non-null  int64
 11  thumbnail_link         40901 non-null  object
 12  comments_disabled      40901 non-null  bool
 13  ratings_disabled       40901 non-null  bool
 14  video_error_or_removed 40901 non-null  bool
 15  description            40332 non-null  object
dtypes: bool(3), int64(5), object(8)
memory usage: 4.5+ MB
```

```python
columns_to_remove = ['thumbnail_link','description']
data = data.drop(columns = columns_to_remove)
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 40901 entries, 0 to 40948
Data columns (total 14 columns):
 #   Column                 Non-Null Count  Dtype
---  ------                 --------------  -----
 0   video_id               40901 non-null  object
 1   trending_date          40901 non-null  object
 2   title                  40901 non-null  object
 3   channel_title          40901 non-null  object
 4   category_id            40901 non-null  int64
 5   publish_time           40901 non-null  object
 6   tags                   40901 non-null  object
 7   views                  40901 non-null  int64
 8   likes                  40901 non-null  int64
```

```
 9   dislikes              40901 non-null  int64
10   comment_count         40901 non-null  int64
11   comments_disabled     40901 non-null  bool
12   ratings_disabled      40901 non-null  bool
13   video_error_or_removed 40901 non-null  bool
dtypes: bool(3), int64(5), object(6)
memory usage: 3.9+ MB
```

```
from datetime import datetime
import datetime
```

```
data['trending_date'] = data['trending_date'].apply(lambda x : datetime.datetime.strptime(x,'%y.%d
data.head(3)
```

| | video_id | trending_date | title | channel_title | category_id |
|---|---|---|---|---|---|
| 0 | 2kyS6SvSYSE | 2017-11-14 | WE WANT TO TALK ABOUT OUR MARRIAGE | CaseyNeistat | 22 |
| 1 | 1ZAPwfrtAFY | 2017-11-14 | The Trump Presidency: Last Week Tonight with J... | LastWeekTonight | 24 |
| 2 | 5qpjK5DgCt4 | 2017-11-14 | Racist Superman | Rudy Mancuso, King Bach & Le... | Rudy Mancuso | 23 |

```
data['publish_time'] = pd.to_datetime(data['publish_time'])
```

```
data['publish_month'] = data['publish_time'].dt.month
data['publish_day'] = data['publish_time'].dt.day
data['publish_hour'] = data['publish_time'].dt.hour
data.head(2)
```

| | video_id | trending_date | title | channel_title | category_id |
|---|---|---|---|---|---|
| 0 | 2kyS6SvSYSE | 2017-11-14 | WE WANT TO TALK ABOUT OUR MARRIAGE | CaseyNeistat | 22 |
| 1 | 1ZAPwfrtAFY | 2017-11-14 | The Trump Presidency: Last Week Tonight with J... | LastWeekTonight | 24 |

```
print(sorted(data['category_id'].unique()))
[1, 2, 10, 15, 17, 19, 20, 22, 23, 24, 25, 26, 27, 28, 29, 30, 43]
```

```
[1, 2, 10, 15, 17, 19, 20, 22, 23, 24, 25, 26, 27, 28, 29, 43]
[1, 2, 10, 15, 17, 19, 20, 22, 23, 24, 25, 26, 27, 28, 29, 30, 43]
```

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

data.head()
```

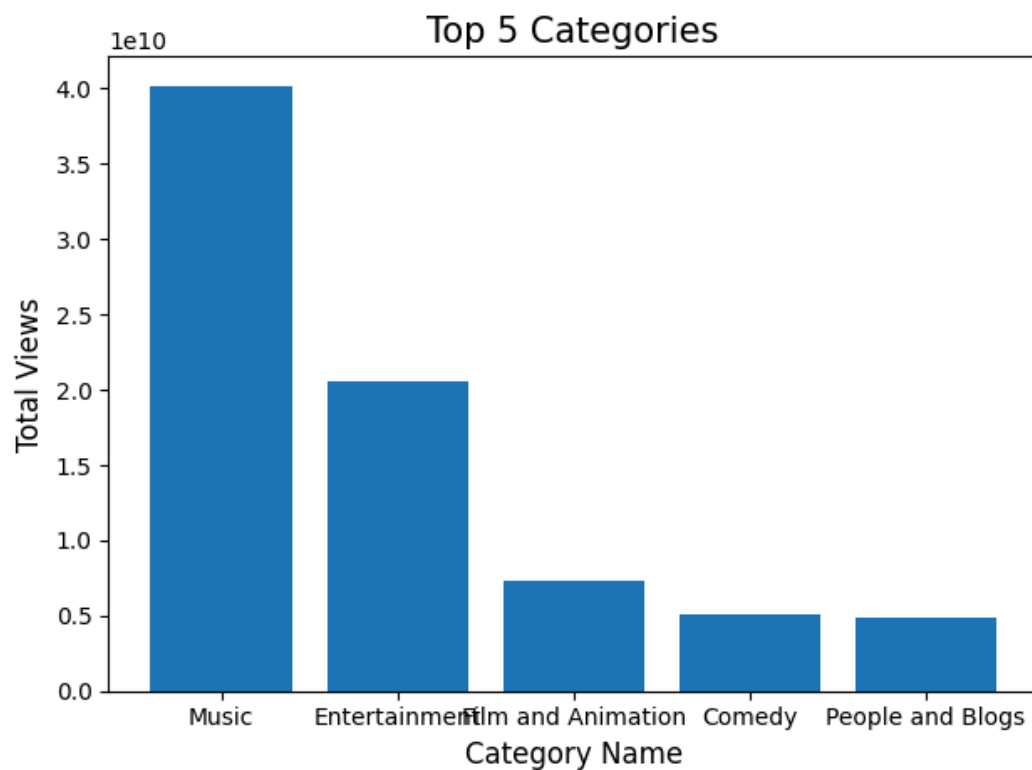| | video_id | trending_date | title | channel_title | category_id |
|---|---|---|---|---|---|
| 0 | 2kyS6SvSYSE | 2017-11-14 | WE WANT TO TALK ABOUT OUR MARRIAGE | CaseyNeistat | 22 |
| 1 | 1ZAPwfrtAFY | 2017-11-14 | The Trump Presidency: Last Week Tonight with J... | LastWeekTonight | 24 |
| 2 | 5qpjK5DgCt4 | 2017-11-14 | Racist Superman \| Rudy Mancuso, King Bach & Le... | Rudy Mancuso | 23 |
| 3 | puqaWrEC7tY | 2017-11-14 | Nickelback Lyrics: Real or Fake? | Good Mythical Morning | 24 |
| 4 | d380meD0W0M | 2017-11-14 | I Dare You: GOING BALD!? | nigahiga | 24 |

```python
data['year'] = data['publish_time'].dt.year
yearly_counts = data.groupby('year')['video_id'].count()
yearly_counts.plot(kind = 'bar', xlabel = 'Year', ylabel = 'Total Publish Video Per Year')
plt.show()
```
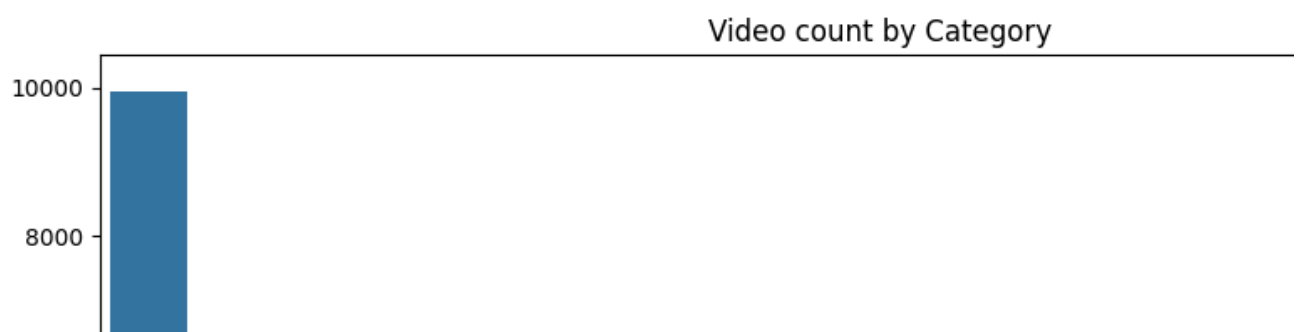
```
category_views = data.groupby('category_name')['views'].sum().reset_index()
top_categories = category_views.sort_values(by='views', ascending = False).head(5)
plt.bar(top_categories['category_name'], top_categories['views'])
plt.xlabel('Category Name', fontsize = 12)
plt.ylabel('Total Views', fontsize = 12)
plt.title('Top 5 Categories', fontsize = 15)
plt.tight_layout()
plt.show()
```



```
plt.figure(figsize = (12,6))
sns.countplot(x = 'category_name', data=data, order=data['category_name'].value_counts().index)
plt.xticks(rotation=90)
plt.title('Video count by Category')
plt.show()
```
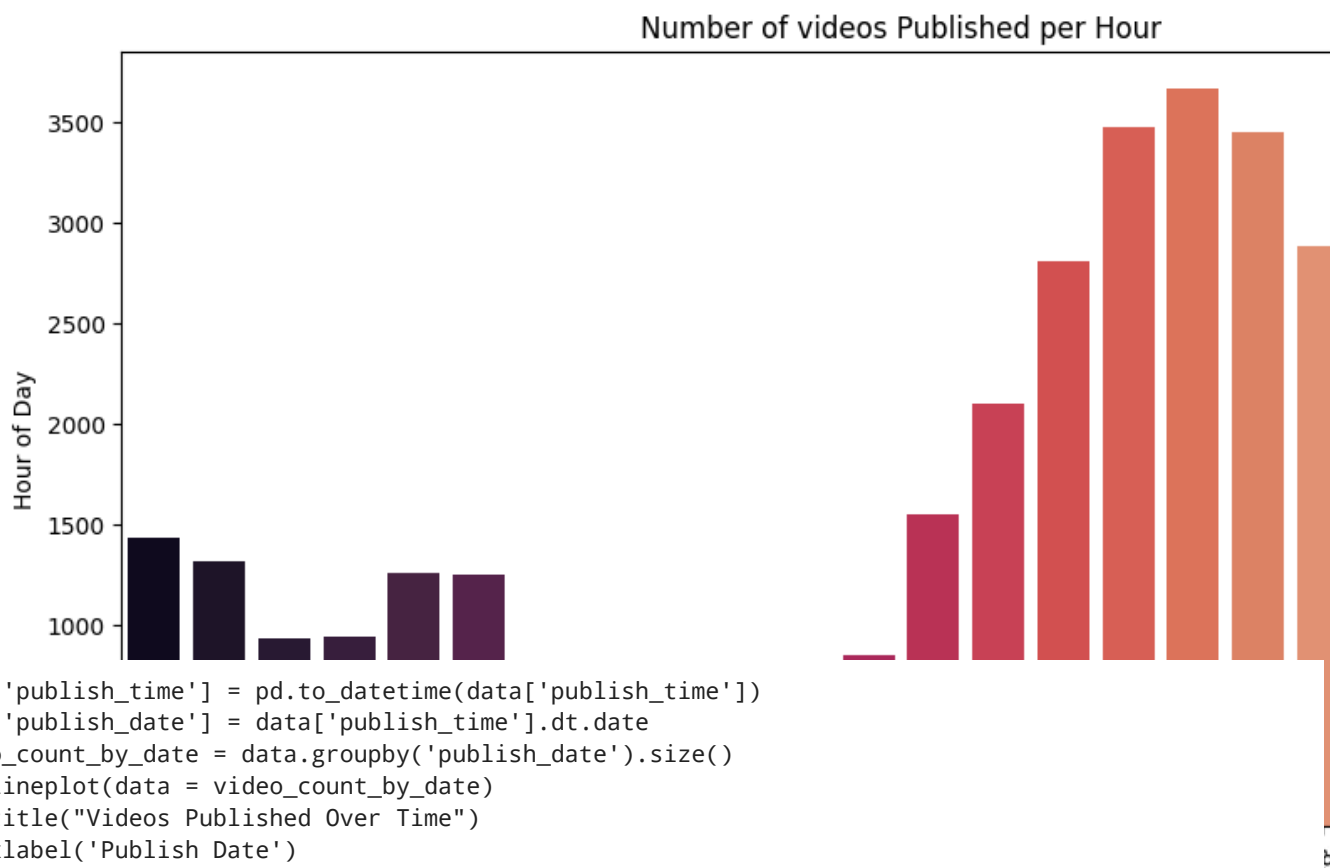
```python
videos_per_hour = data['publish_hour'].value_counts().sort_index()

plt.figure(figsize=(12,6))
sns.barplot(x= videos_per_hour.index, y = videos_per_hour.values, palette = 'rocket')
plt.title('Number of videos Published per Hour')
plt.xlabel('Number of Videos Published Per Hour')
plt.ylabel('Hour of Day')
plt.xticks(rotation = 45)
plt.show()
```
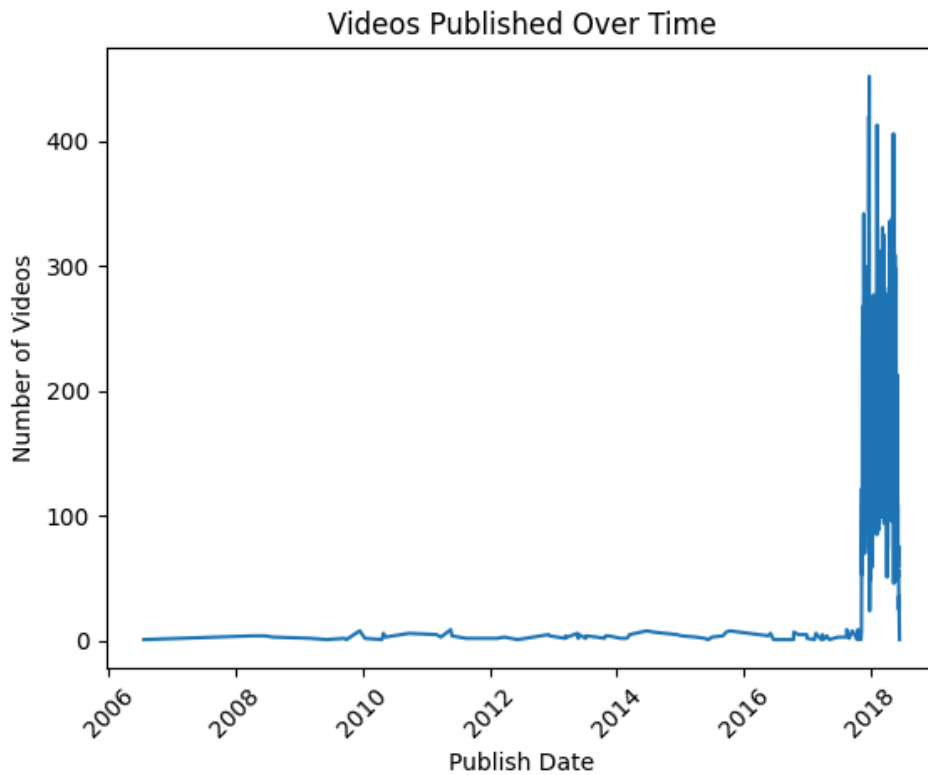
<ipython-input-34-242e26f9b13c>:4: FutureWarning:

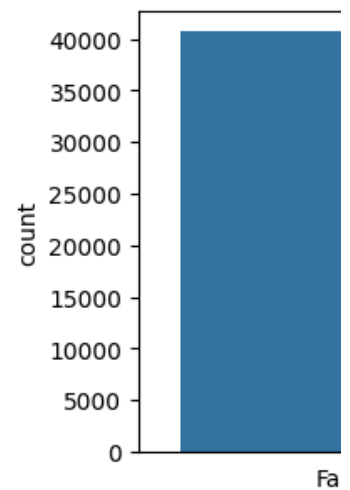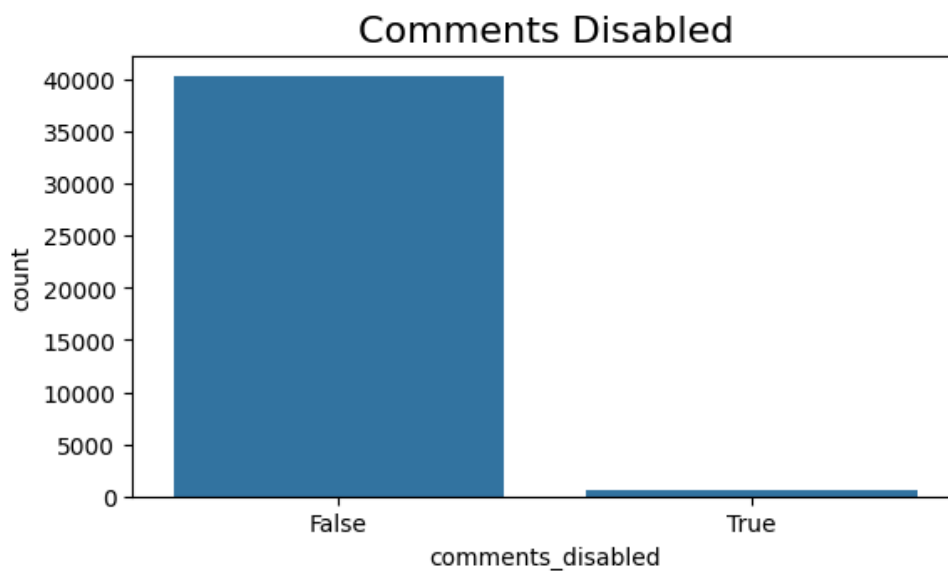Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign

  sns.barplot(x= videos_per_hour.index, y = videos_per_hour.values, palette = 'rocket')



```python
data['publish_time'] = pd.to_datetime(data['publish_time'])
data['publish_date'] = data['publish_time'].dt.date
video_count_by_date = data.groupby('publish_date').size()
sns.lineplot(data = video_count_by_date)
plt.title("Videos Published Over Time")
plt.xlabel('Publish Date')
plt.ylabel('Number of Videos')
plt.xticks(rotation = 45)
plt.show()
```

Videos Published Over Time

```
plt.figure(figsize = (14,8))
plt.subplots_adjust(wspace = 0.2,hspace = 0.4, top = 0.9)
plt.subplot(2,2,1)
g = sns.countplot(x ='comments_disabled', data = data)
g.set_title("Comments Disabled",fontsize= 16)
plt.subplot(2,2,2)
g1 = sns.countplot(x = 'ratings_disabled', data = data)
g1.set_title("Rating Disabled",fontsize = 16)
plt.subplot(2,2,3)
g2 = sns.countplot(x = 'video_error_or_removed',data = data)
g2.set_title("Video Error or Removed",fontsize = 16)
plt.show()
```



Comments Disabled

Video Error or Removed

```
corr_matrix = data['views'].corr(data['likes'])
corr_matrix
```

0.8491785476230508