# Python Portfolio Project - TikTok

February 27, 2024

## 0.1 *TIK-TOK PROJECT*

### 0.1.1 Analyzing the view-statistics

```python
[1]: # Import packages for data manipulation
     import pandas as pd
     import numpy as np

     # Import packages for data visualization
     import matplotlib.pyplot as plt
     import seaborn as sns
```

```python
[2]: # Load dataset into dataframe
     data = pd.read_csv("/Users/samantarana/Downloads/tiktok_dataset.csv")
```

```python
[3]: # Display and examine the first few rows of the dataframe
     data.head()
```

```
[3]:    # claim_status     video_id  video_duration_sec  \
     0  1         claim  7017666017                  59
     1  2         claim  4014381136                  32
     2  3         claim  9859838091                  31
     3  4         claim  1866847991                  25
     4  5         claim  7105231098                  19

                              video_transcription_text verified_status  \
     0  someone shared with me that drone deliveries a…    not verified
     1  someone shared with me that there are more mic…    not verified
     2  someone shared with me that american industria…    not verified
     3  someone shared with me that the metro of st. p…    not verified
     4  someone shared with me that the number of busi…    not verified

       author_ban_status  video_view_count  video_like_count  video_share_count  \
     0      under review          343296.0           19425.0              241.0
     1            active          140877.0           77355.0            19034.0
     2            active          902185.0           97690.0             2858.0
     3            active          437506.0          239954.0            34812.0
     4            active           56167.0           34987.0             4110.0
```

```
     video_download_count  video_comment_count
0                     1.0                   0.0
1                  1161.0                 684.0
2                   833.0                 329.0
3                  1234.0                 584.0
4                   547.0                 152.0
```

[4]: `# Get the size of the data`
`data.size`

[4]: 232584

[5]: `# Get the shape of the data`
`data.shape`

[5]: (19382, 12)

[6]: `# Get basic information about the data`
`data.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 19382 entries, 0 to 19381
Data columns (total 12 columns):
 #   Column                   Non-Null Count  Dtype
---  ------                   --------------  -----
 0   #                        19382 non-null  int64
 1   claim_status             19084 non-null  object
 2   video_id                 19382 non-null  int64
 3   video_duration_sec       19382 non-null  int64
 4   video_transcription_text 19084 non-null  object
 5   verified_status          19382 non-null  object
 6   author_ban_status        19382 non-null  object
 7   video_view_count         19084 non-null  float64
 8   video_like_count         19084 non-null  float64
 9   video_share_count        19084 non-null  float64
 10  video_download_count     19084 non-null  float64
 11  video_comment_count      19084 non-null  float64
dtypes: float64(5), int64(3), object(4)
memory usage: 1.8+ MB
```

[7]: `# Generate a table of descriptive statistics`
`data.describe()`

[7]:
```
                   #       video_id  video_duration_sec  video_view_count  \
count   19382.000000  1.938200e+04        19382.000000      19084.000000
mean     9691.500000  5.627454e+09           32.421732     254708.558688
std      5595.245794  2.536440e+09           16.229967     322893.280814
```
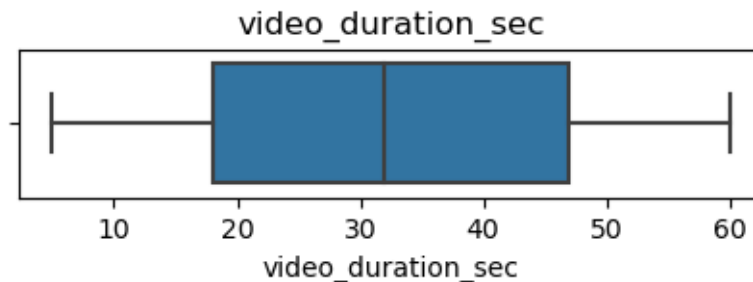
2

```
min          1.000000  1.234959e+09            5.000000           20.000000
25%       4846.250000  3.430417e+09           18.000000         4942.500000
50%       9691.500000  5.618664e+09           32.000000         9954.500000
75%      14536.750000  7.843960e+09           47.000000       504327.000000
max      19382.000000  9.999873e+09           60.000000       999817.000000

       video_like_count  video_share_count  video_download_count  \
count      19084.000000       19084.000000          19084.000000
mean       84304.636030       16735.248323           1049.429627
std       133420.546814       32036.174350           2004.299894
min            0.000000           0.000000              0.000000
25%          810.750000         115.000000              7.000000
50%         3403.500000         717.000000             46.000000
75%       125020.000000       18222.000000           1156.250000
max       657830.000000      256130.000000          14994.000000

       video_comment_count
count         19084.000000
mean            349.312146
std             799.638865
min               0.000000
25%               1.000000
50%               9.000000
75%             292.000000
max            9599.000000
```
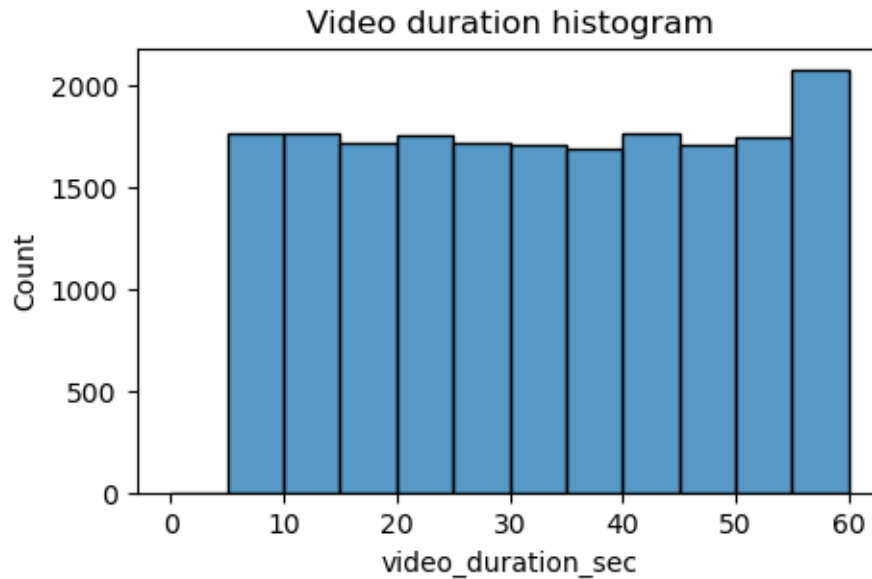
```python
[8]: # Create a boxplot to visualize distribution of `video_duration_sec`
     plt.figure(figsize=(5,1))
     plt.title('video_duration_sec')
     sns.boxplot(x=data['video_duration_sec']);
```
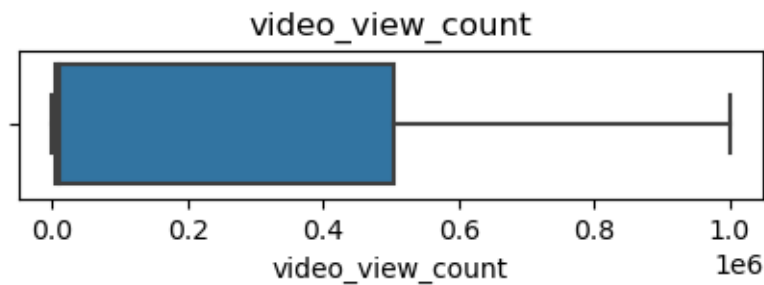


```python
[9]: plt.figure(figsize=(5,3))
     sns.histplot(data['video_duration_sec'], bins=range(0,61,5))
     plt.title('Video duration histogram');
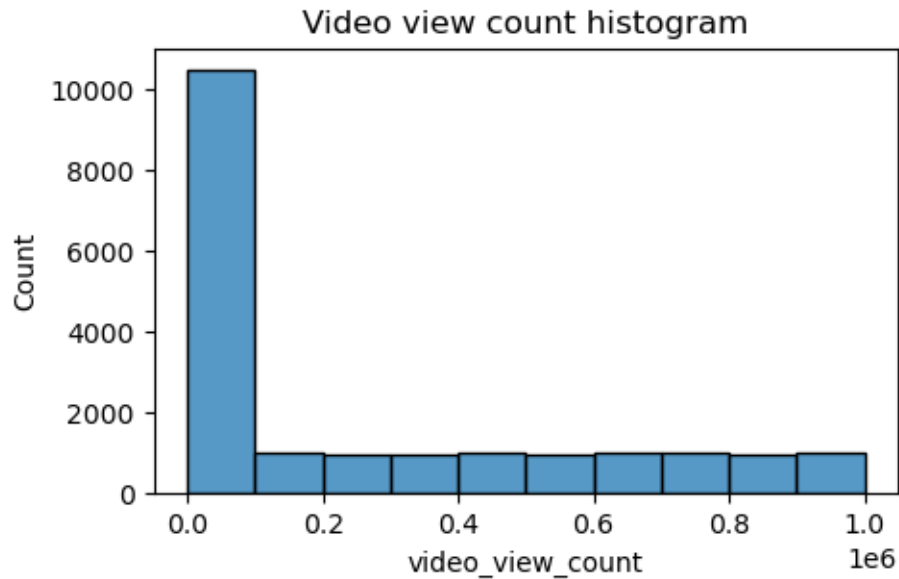```

Video duration histogram

All videos are 5-60 seconds in length, and the distribution is uniform.

```
[10]:  # Create a boxplot to visualize distribution of `video_view_count`
       plt.figure(figsize=(5, 1))
       plt.title('video_view_count')
       sns.boxplot(x=data['video_view_count']);
```
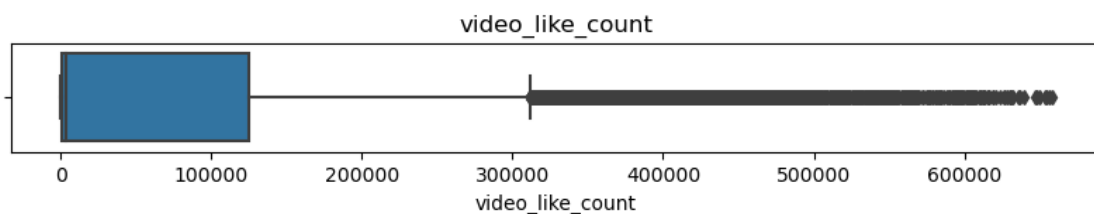


video_view_count

```
[11]:  plt.figure(figsize=(5,3))
       sns.histplot(data['video_view_count'], bins=range(0,(10**6+1),10**5))
       plt.title('Video view count histogram');
```

Video view count histogram

This variable has a very uneven distribution, with more than half the videos receiving fewer than 100,000 views. Distribution of view counts > 100,000 views is uniform.

```
[12]:  # Create a boxplot to visualize distribution of `video_like_count`
       plt.figure(figsize=(10,1))
       plt.title('video_like_count')
       sns.boxplot(x=data['video_like_count']);
```
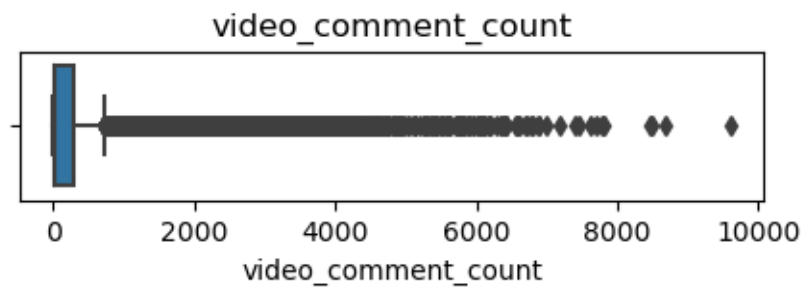


video_like_count

```
[13]:  # plt.figure(figsize=(5,3))
       ax = sns.histplot(data['video_like_count'], bins=range(0,(7*10**5+1),10**5))
       labels = [0] + [str(i) + 'k' for i in range(100, 701, 100)]
       ax.set_xticks(range(0,7*10**5+1,10**5), labels=labels)
       plt.title('Video like count histogram');
```
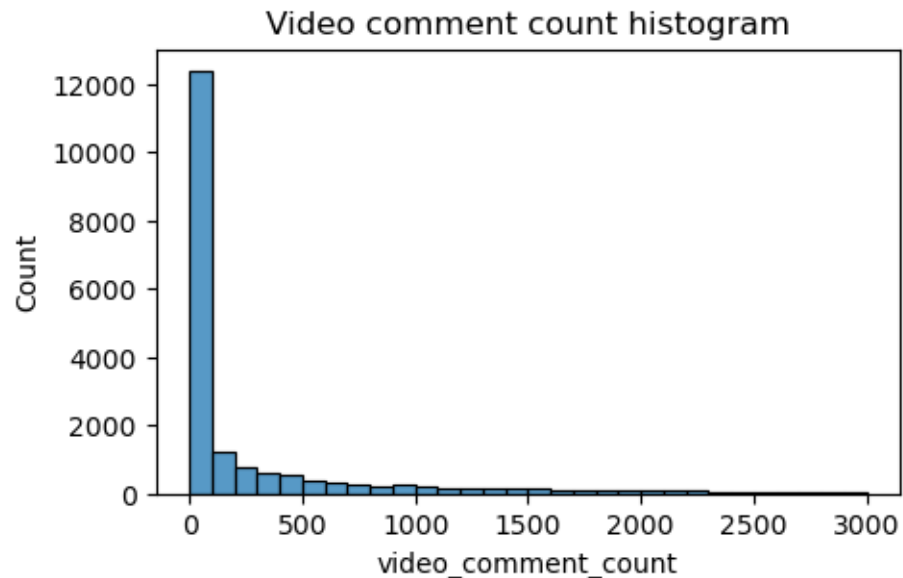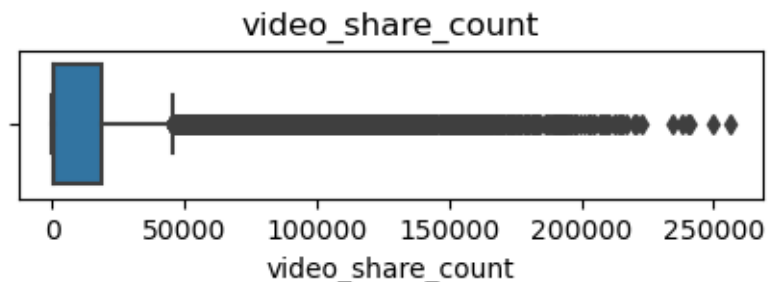
5

## Video like count histogram



Similar to view count, there are far more videos with $< 100{,}000$ likes than there are videos with more. However, in this case, there is more of a taper, as the data skews right, with many videos at the upper extremity of like count.

```
[14]: # Create a boxplot to visualize distribution of `video_comment_count`
      plt.figure(figsize=(5,1))
      plt.title('video_comment_count')
      sns.boxplot(x=data['video_comment_count']);
```

## video_comment_count

```
[15]: plt.figure(figsize=(5,3))
      sns.histplot(data['video_comment_count'], bins=range(0,(3001),100))
      plt.title('Video comment count histogram');
```
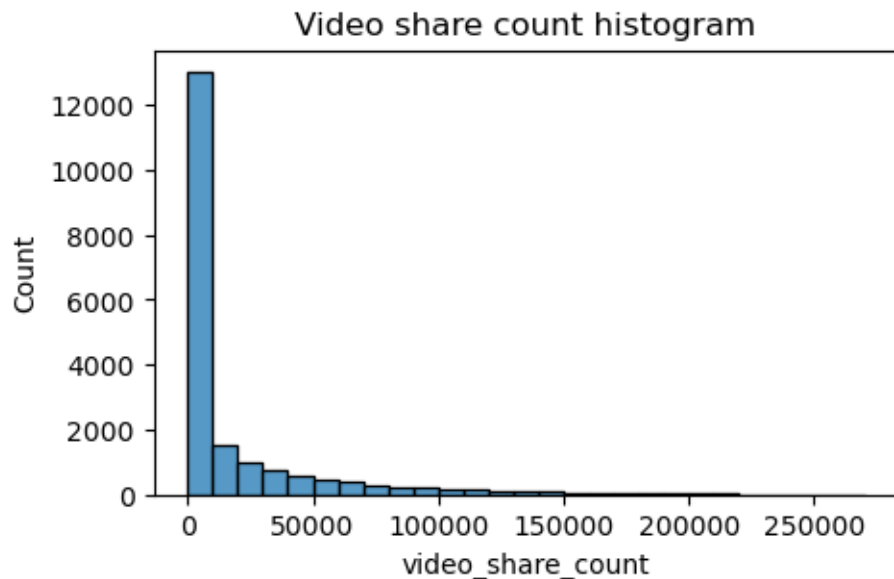

Video comment count histogram

Again, the vast majority of videos are grouped at the bottom of the range of values for video comment count. Most videos have fewer than 100 comments. The distribution is very right-skewed.

```
[16]: # Create a boxplot to visualize distribution of `video_share_count`
      plt.figure(figsize=(5,1))
      plt.title('video_share_count')
      sns.boxplot(x=data['video_share_count']);
```
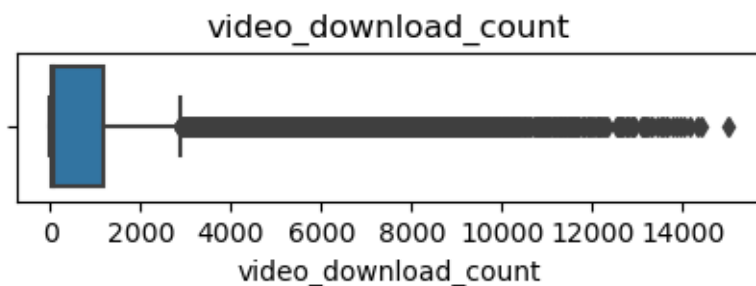

video_share_count

```
[17]: plt.figure(figsize=(5,3))
      sns.histplot(data['video_share_count'], bins=range(0,(270001),10000))
      plt.title('Video share count histogram');
```
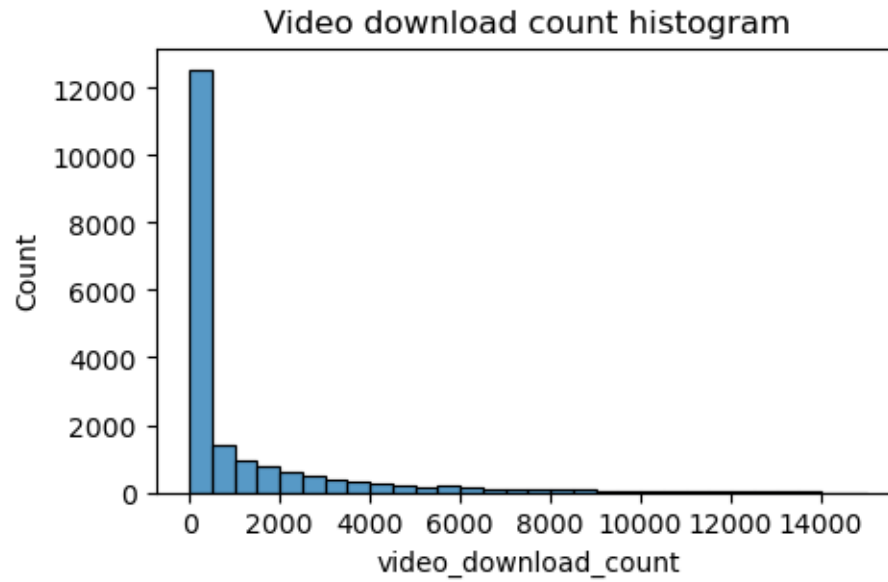
Video share count histogram

The overwhelming majority of videos had fewer than 10,000 shares. The distribution is very skewed to the right.

```
[18]: # Create a boxplot to visualize distribution of `video_download_count`
      plt.figure(figsize=(5,1))
      plt.title('video_download_count')
      sns.boxplot(x=data['video_download_count']);
```

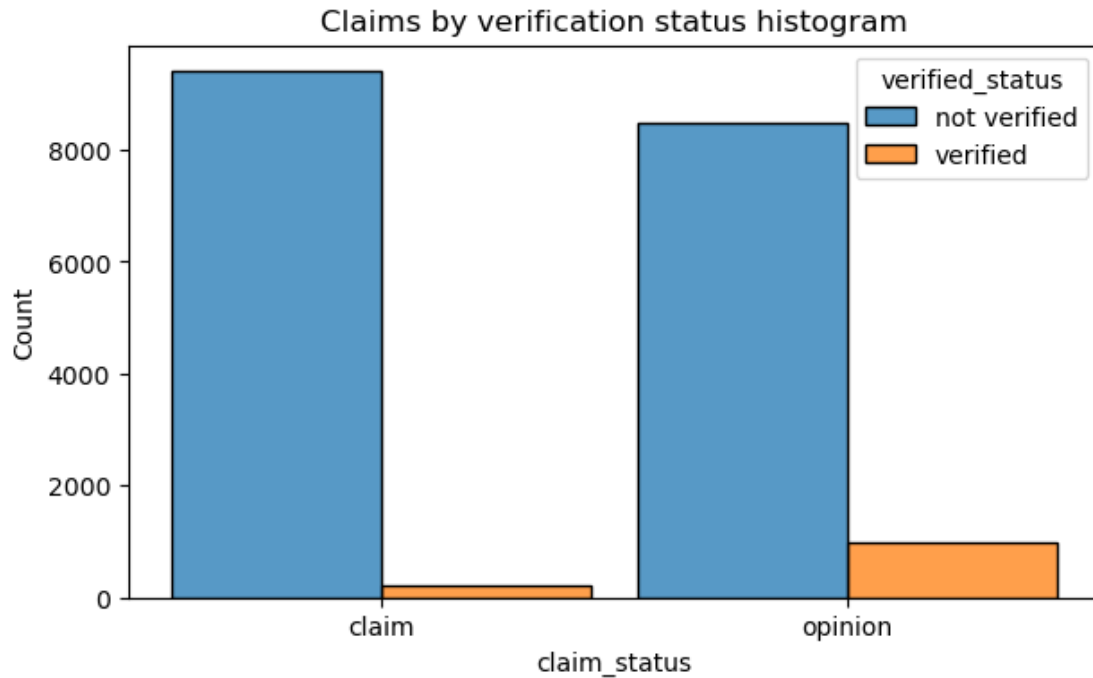

video_download_count

```
[19]: plt.figure(figsize=(5,3))
      sns.histplot(data['video_download_count'], bins=range(0,(15001),500))
      plt.title('Video download count histogram');
```
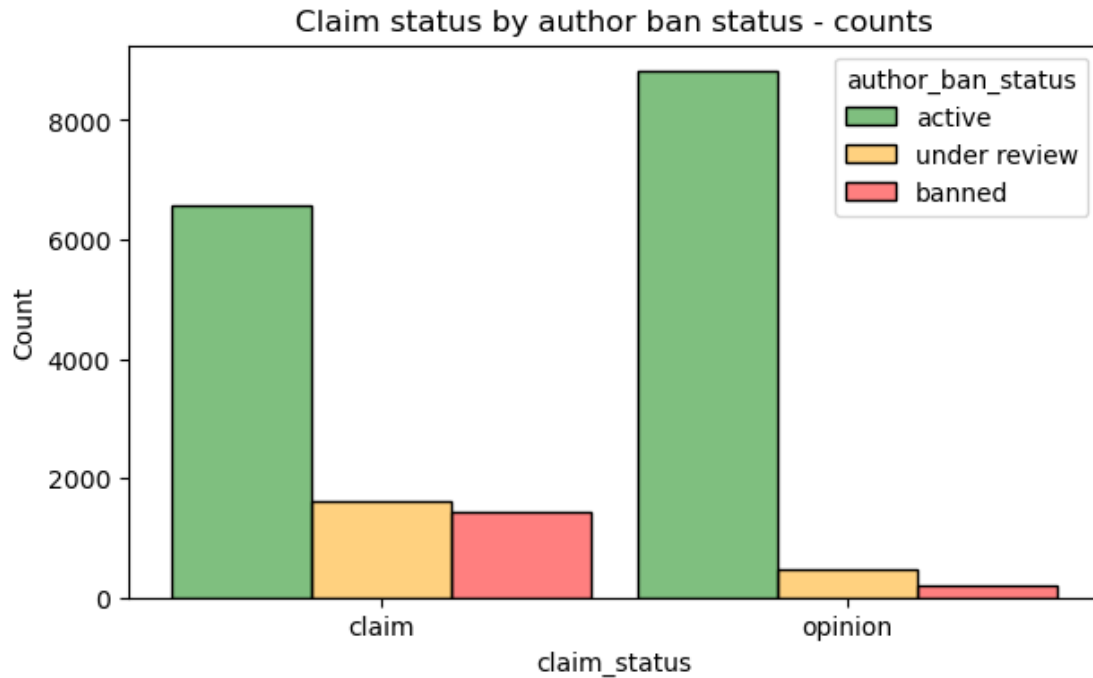
Video download count histogram

The majority of videos were downloaded fewer than 500 times, but some were downloaded over 12,000 times. Again, the data is very skewed to the right.

```
[20]: plt.figure(figsize=(7,4))
      sns.histplot(data=data,
                  x='claim_status',
                  hue='verified_status',
                  multiple='dodge',
                  shrink=0.9)
      plt.title('Claims by verification status histogram');
```

## Claims by verification status histogram



There are far fewer verified users than unverified users, but if a user is verified, they are much more likely to post opinions.
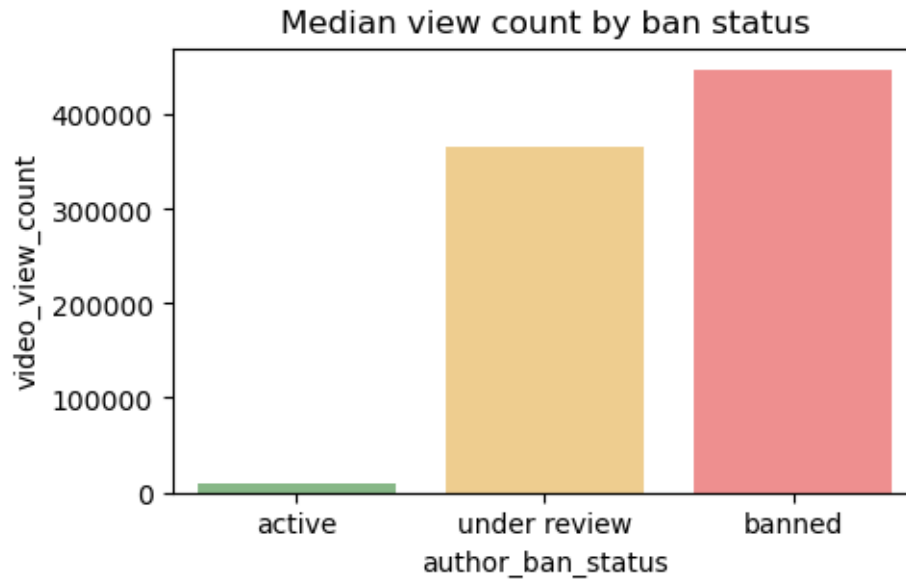
```
[21]: fig = plt.figure(figsize=(7,4))
sns.histplot(data, x='claim_status', hue='author_ban_status',
             multiple='dodge',
             hue_order=['active', 'under review', 'banned'],
             shrink=0.9,
             palette={'active':'green', 'under review':'orange', 'banned':
 ↪'red'},
             alpha=0.5)
plt.title('Claim status by author ban status - counts');
```

Claim status by author ban status - counts

For both claims and opinions, there are many more active authors than banned authors or authors under review; however, the proportion of active authors is far greater for opinion videos than for claim videos. Again, it seems that authors who post claim videos are more likely to come under review and/or get banned.

```
[22]: ban_status_counts = data.groupby(['author_ban_status']).median(
          numeric_only=True).reset_index()

      fig = plt.figure(figsize=(5,3))
      sns.barplot(data=ban_status_counts,
                  x='author_ban_status',
                  y='video_view_count',
                  order=['active', 'under review', 'banned'],
                  palette={'active':'green', 'under review':'orange', 'banned':'red'},
                  alpha=0.5)
      plt.title('Median view count by ban status');
```
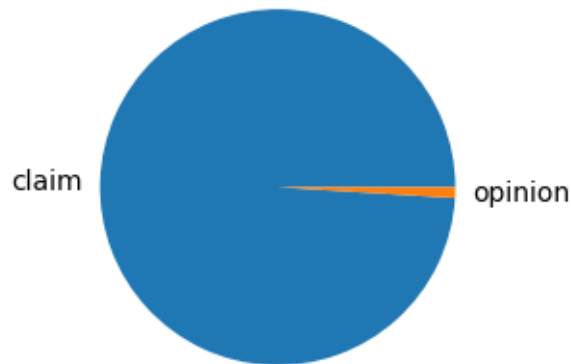
Median view count by ban status

The median view counts for non-active authors are many times greater than the median view count for active authors. Since you know that non-active authors are more likely to post claims, and that videos by non-active authors get far more views on aggregate than videos by active authors, then video_view_count might be a good indicator of claim status.

```
[23]: data.groupby('claim_status')['video_view_count'].median()
```

```
[23]: claim_status
      claim      501555.0
      opinion      4953.0
      Name: video_view_count, dtype: float64
```

```
[24]: fig = plt.figure(figsize=(3,3))
      plt.pie(data.groupby('claim_status')['video_view_count'].sum(),␣
        ↪labels=['claim', 'opinion'])
      plt.title('Total views by video claim status');
```

## Total views by video claim status



The overall view count is dominated by claim videos even though there are roughly the same number of each video in the dataset.

**Determining the Outleirs**

```
[25]: count_cols = ['video_view_count',
                     'video_like_count',
                     'video_share_count',
                     'video_download_count',
                     'video_comment_count',
                     ]

      for column in count_cols:
          q1 = data[column].quantile(0.25)
          q3 = data[column].quantile(0.75)
          iqr = q3 - q1
          median = data[column].median()
          outlier_threshold = median + 1.5*iqr

          # Count the number of values that exceed the outlier threshold
          outlier_count = (data[column] > outlier_threshold).sum()
          print(f'Number of outliers, {column}:', outlier_count)
```

```
Number of outliers, video_view_count: 2343
Number of outliers, video_like_count: 3468
Number of outliers, video_share_count: 3732
Number of outliers, video_download_count: 3733
Number of outliers, video_comment_count: 3882
```
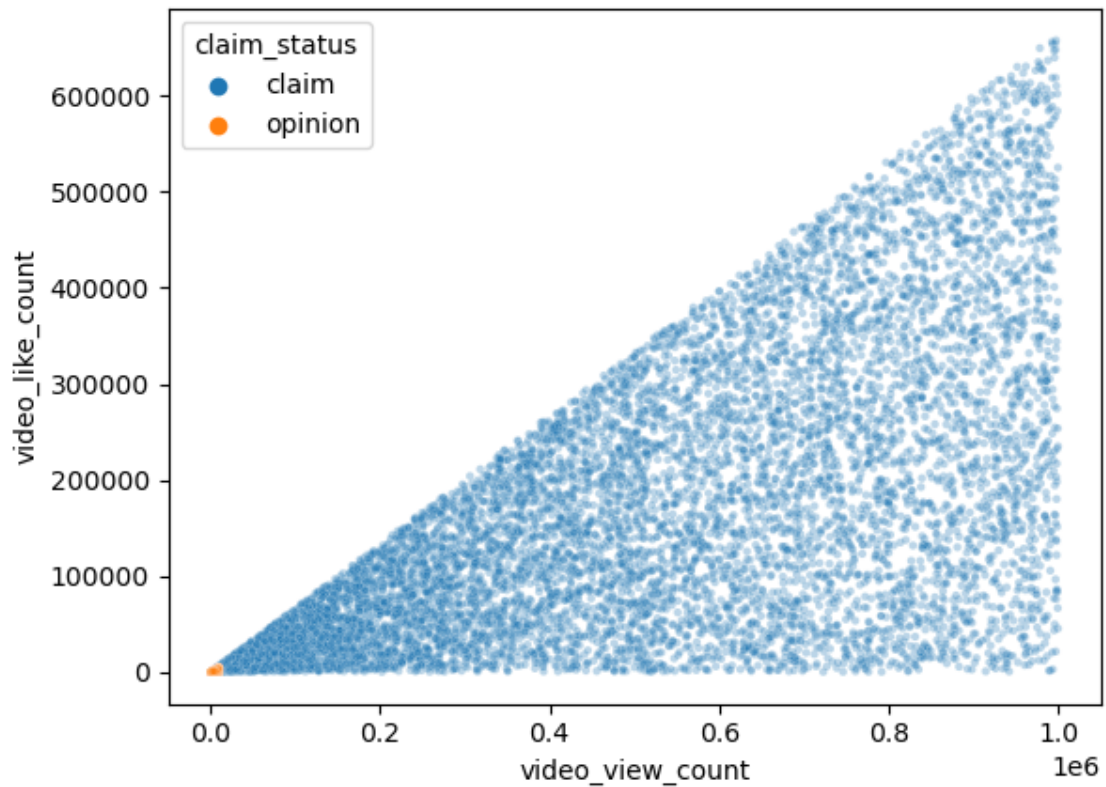
```
[26]: # Create a scatterplot of `video_view_count` versus `video_like_count`␣
      ↪according to 'claim_status'
```

13

```
sns.scatterplot(x=data["video_view_count"], y=data["video_like_count"],
                hue=data["claim_status"], s=10, alpha=.3)
plt.show()
```



[27]:
```
# Create a scatterplot of `video_view_count` versus `video_like_count` for␣
 ↪opinions only
opinion = data[data['claim_status']=='opinion']
sns.scatterplot(x=opinion["video_view_count"], y=opinion["video_like_count"],
                s=10, alpha=.3)
plt.show()
```