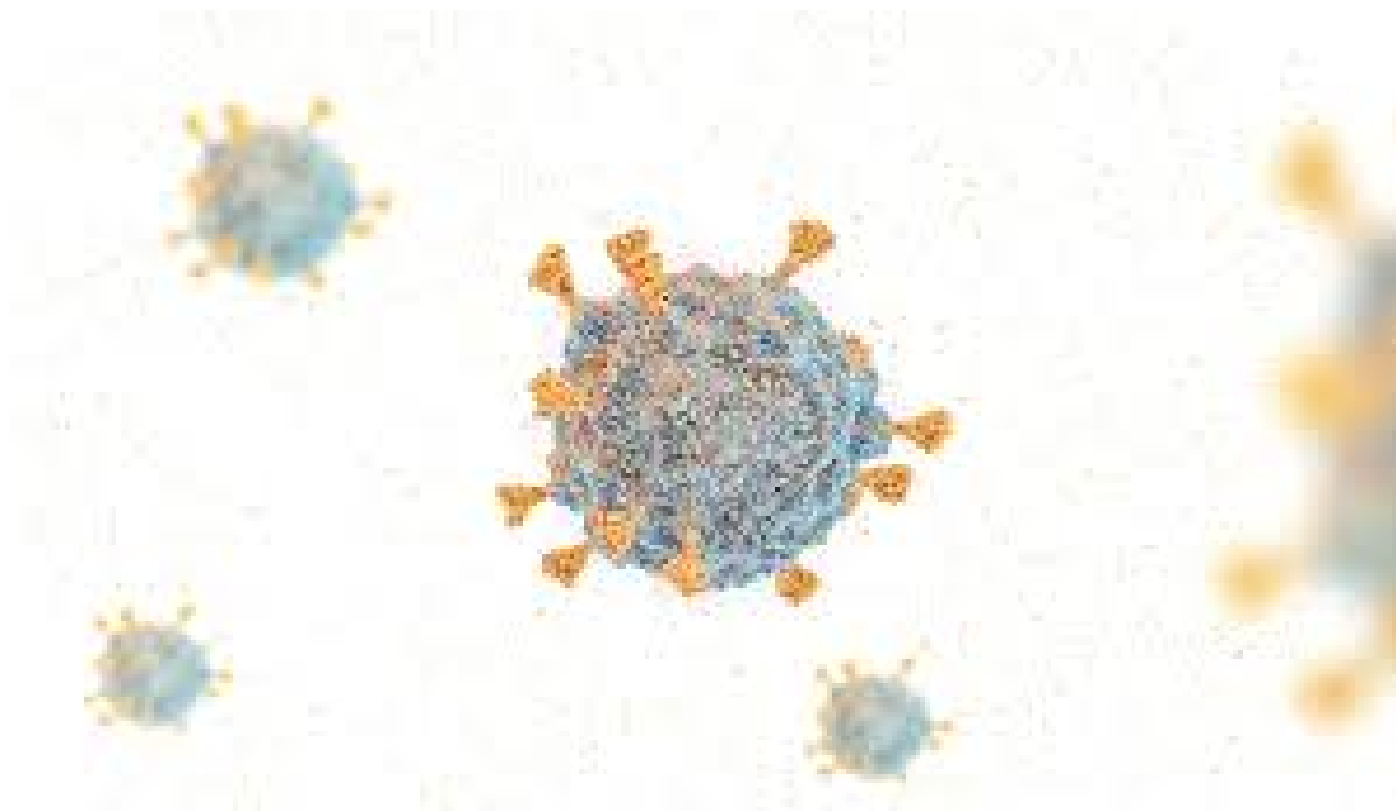**the Dataset was gotten from Kaggle on Omicron daily cases by country (COVID-19 variant) the last update and download for this execerise was on 03/02/2022**



In [1]:
```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
```

In [2]:
```python
cvd=pd.read_csv('covid-variants.csv')
```

In [49]:
```python
cvd.head()
```

Out[49]:

| | location | variant | num_sequences | perc_sequences | num_sequences_total | month | year | day |
|---|---|---|---|---|---|---|---|---|
| **0** | Angola | Alpha | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **1** | Angola | B.1.1.277 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **2** | Angola | B.1.1.302 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **3** | Angola | B.1.1.519 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **4** | Angola | B.1.160 | 0 | 0.0 | 3 | 7 | 2020 | 6 |

In [22]:
```python
cvd.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100416 entries, 0 to 100415
Data columns (total 6 columns):
 #   Column              Non-Null Count   Dtype
---  ------              --------------   -----
 0   location            100416 non-null  object
 1   date                100416 non-null  datetime64[ns]
 2   variant             100416 non-null  object
 3   num_sequences       100416 non-null  int64
 4   perc_sequences      100416 non-null  float64
 5   num_sequences_total 100416 non-null  int64
dtypes: datetime64[ns](1), float64(1), int64(2), object(2)
memory usage: 4.6+ MB
```

In [5]:
```python
cvd.variant.value_counts()
```

Out[5]:
```
Alpha            4184
B.1.1.277        4184
others           4184
S:677P.Pelican   4184
S:677H.Robin1    4184
Omicron          4184
Mu               4184
Lambda           4184
Kappa            4184
Iota             4184
Gamma            4184
Eta              4184
```

```
Epsilon            4184
Delta              4184
Beta               4184
B.1.620            4184
B.1.367            4184
B.1.258            4184
B.1.221            4184
B.1.177            4184
B.1.160            4184
B.1.1.519          4184
B.1.1.302          4184
non_who            4184
Name: variant, dtype: int64
```

In [6]:
```python
cvd.num_sequences.value_counts()
```

Out[6]:
```
0        84173
1         2753
2         1405
3          905
4          631
         ...
1690         1
1719         1
2156         1
1184         1
862          1
Name: num_sequences, Length: 1563, dtype: int64
```

In [7]:
```python
cvd.isnull().any()
```

Out[7]:
```
location              False
date                  False
variant               False
num_sequences         False
perc_sequences        False
num_sequences_total   False
dtype: bool
```

In [8]:
```python
cvd.describe()
```

Out[8]:

| | num_sequences | perc_sequences | num_sequences_total |
|---|---|---|---|
| count | 100416.000000 | 100416.000000 | 100416.000000 |
| mean | 72.171676 | 6.154355 | 1509.582457 |
| std | 1669.262169 | 21.898989 | 8445.291772 |
| min | 0.000000 | -0.010000 | 1.000000 |
| 25% | 0.000000 | 0.000000 | 12.000000 |
| 50% | 0.000000 | 0.000000 | 59.000000 |
| 75% | 0.000000 | 0.000000 | 394.000000 |
| max | 142280.000000 | 100.000000 | 146170.000000 |

# The Exploratort Data Analysis [EDA]

In [9]:
```python
plt.figure(figsize=(14,10))
plt.xticks(rotation=90)
sns.countplot(data=cvd,x='date',hue='variant')
plt.title('Dates')
```

Out[9]:  Text(0.5, 1.0, 'Dates')

In [10]:    `# Let's check the variant wise with top 10 countries with maximum virus`

```python
for virus in cvd.variant.unique():
    most_cases = cvd.loc[cvd['variant'] == virus].groupby('location')[
        'num_sequences'].agg('sum').sort_values(ascending=False)[:10]

    most_cases = pd.DataFrame({'Location':most_cases.index, 'Number of Case':most_cases.values})

    plt.figure(figsize=(20,8))
    sns.barplot(y='Location',x="Number of Case",data=most_cases,palette="plasma_r")
    plt.title('COUNTRIES HAVE MORE {} CASES THAN OTHERS'.format(virus).upper(),loc='center',fontweight="bold")
```
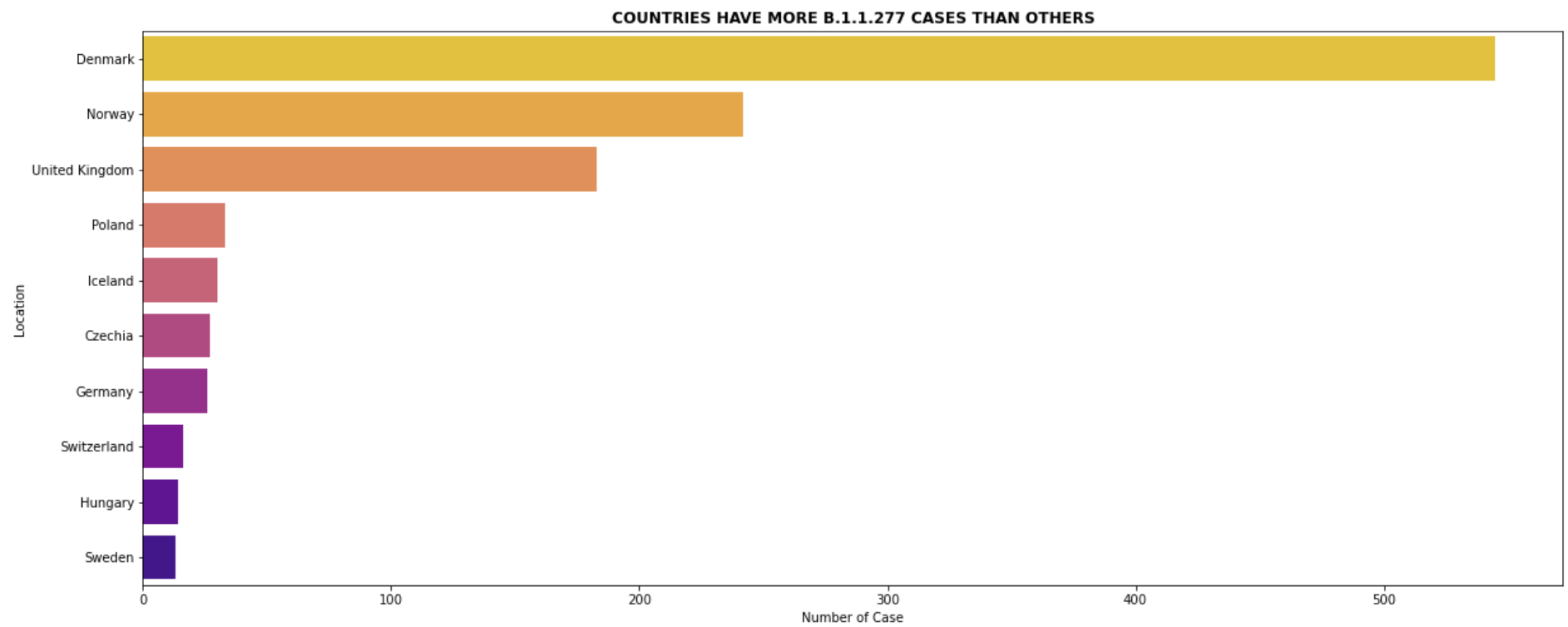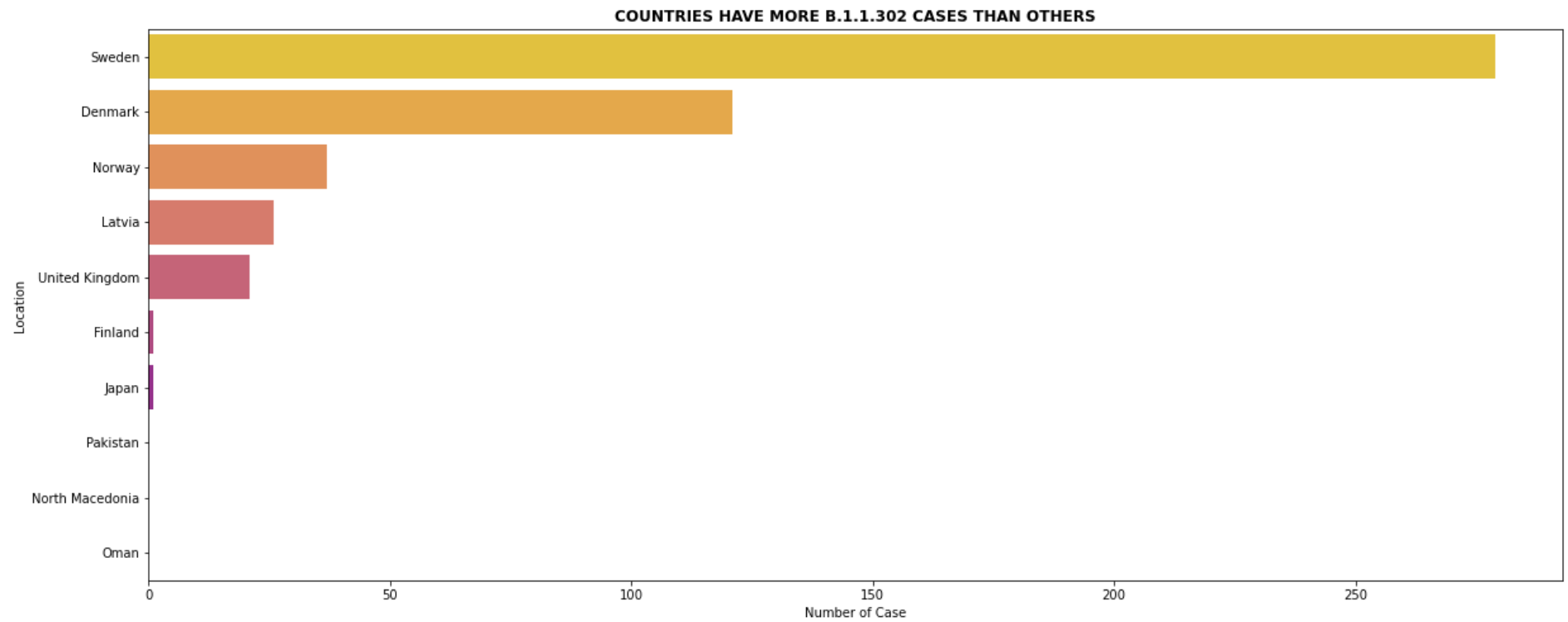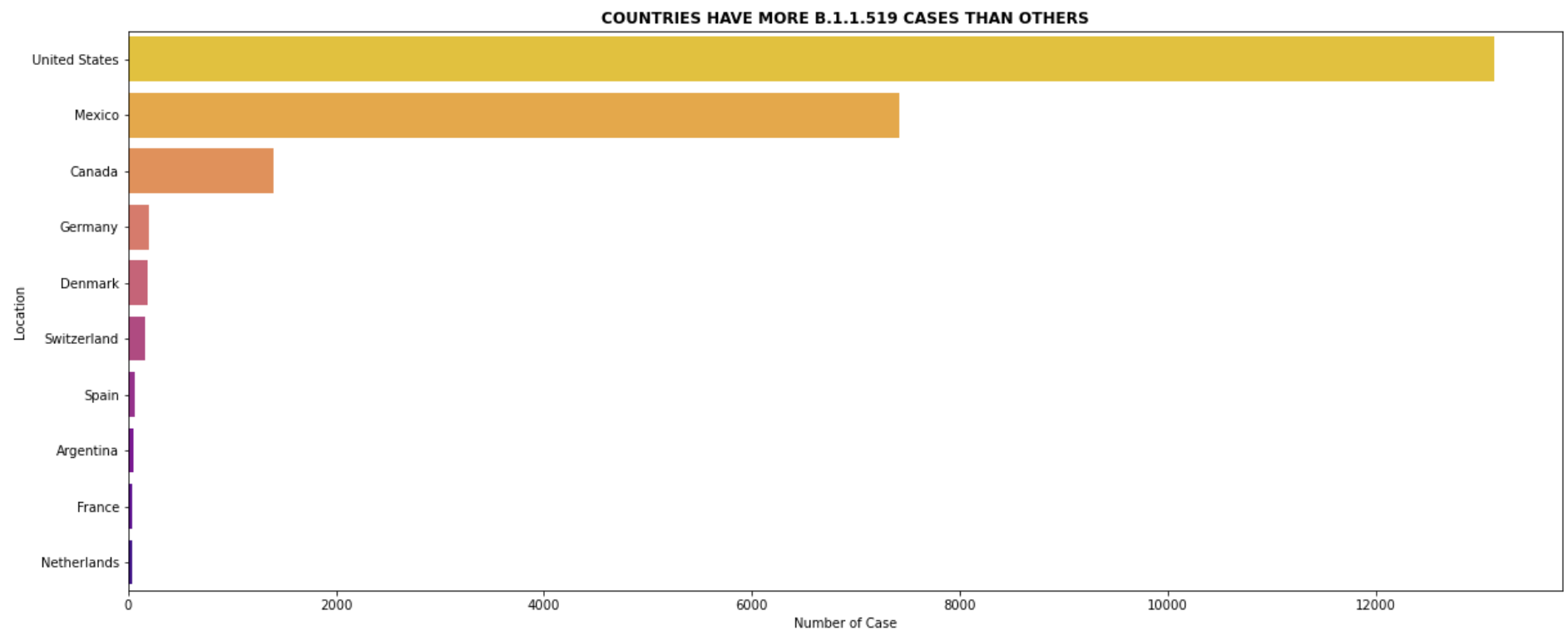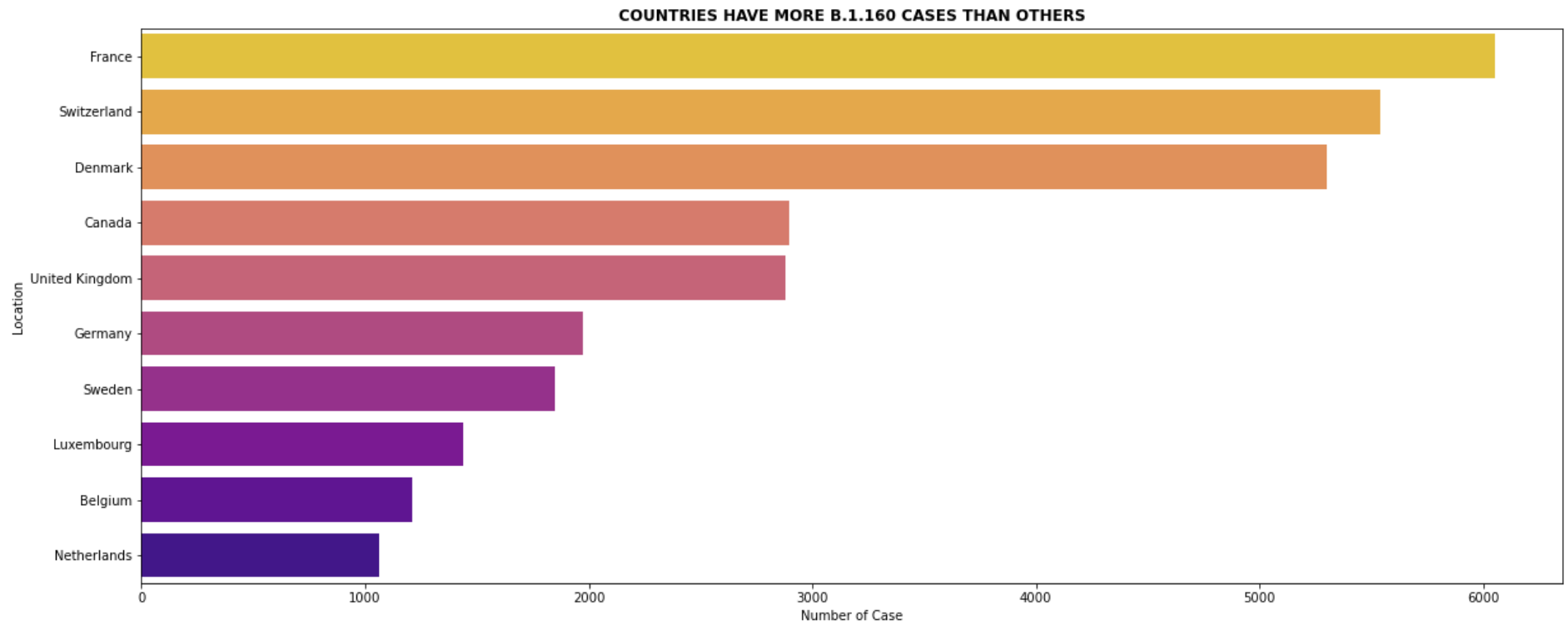
C:\Users\BUSINE~2\AppData\Local\Temp/ipykernel_18476/186801648.py:9: RuntimeWarning: More than 20 figures have been opened. Figure
s created through the pyplot interface (`matplotlib.pyplot.figure`) are retained until explicitly closed and may consume too much
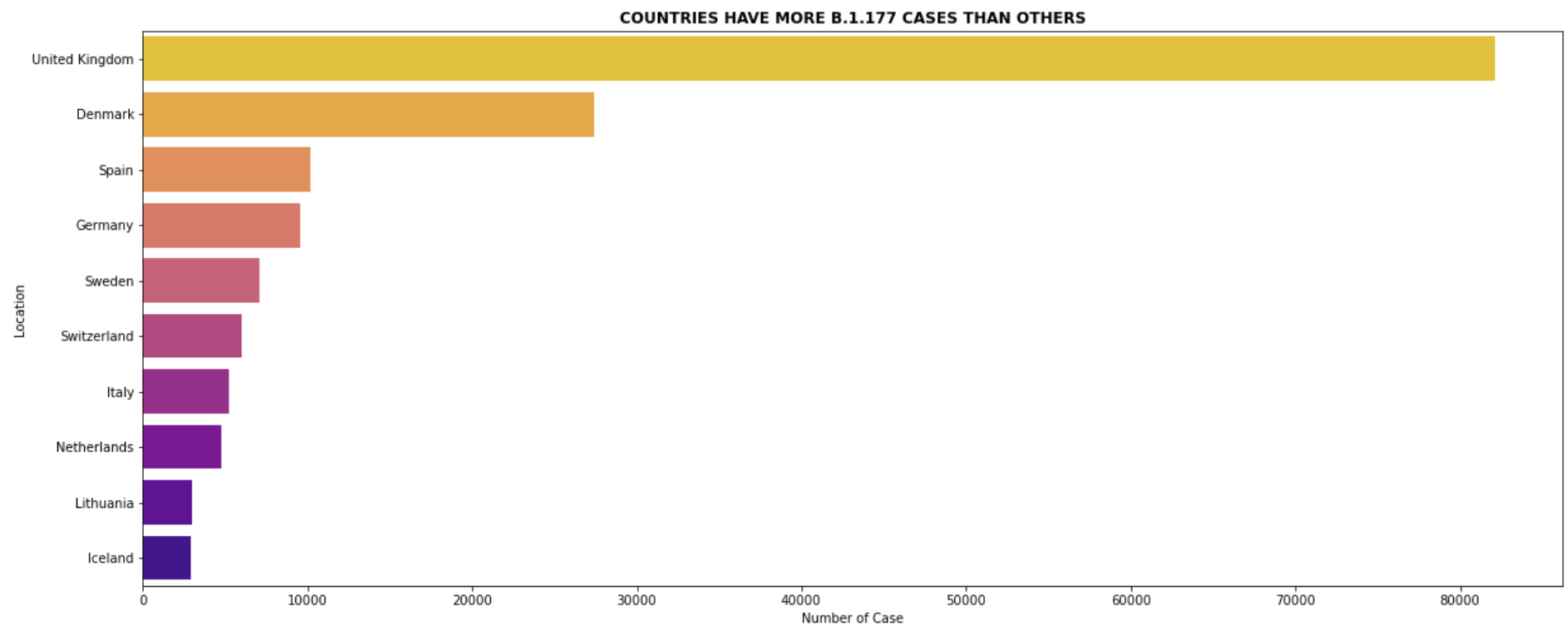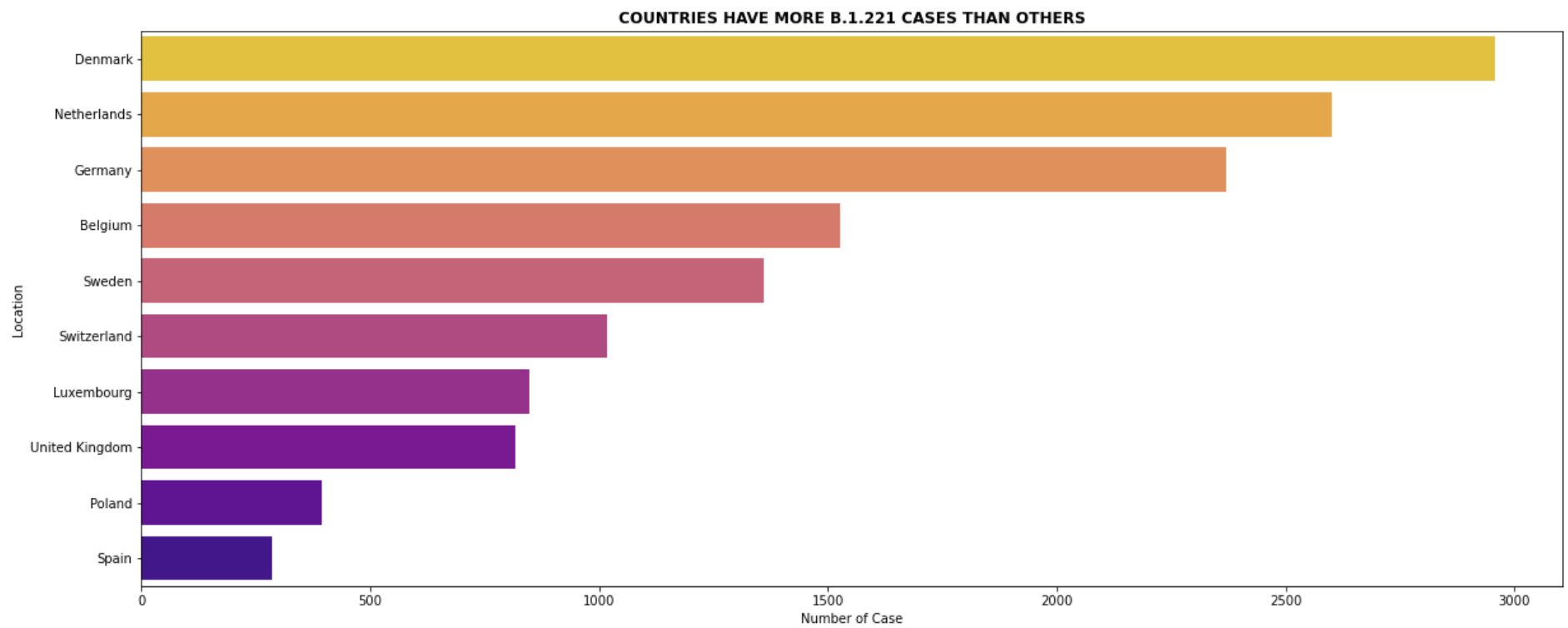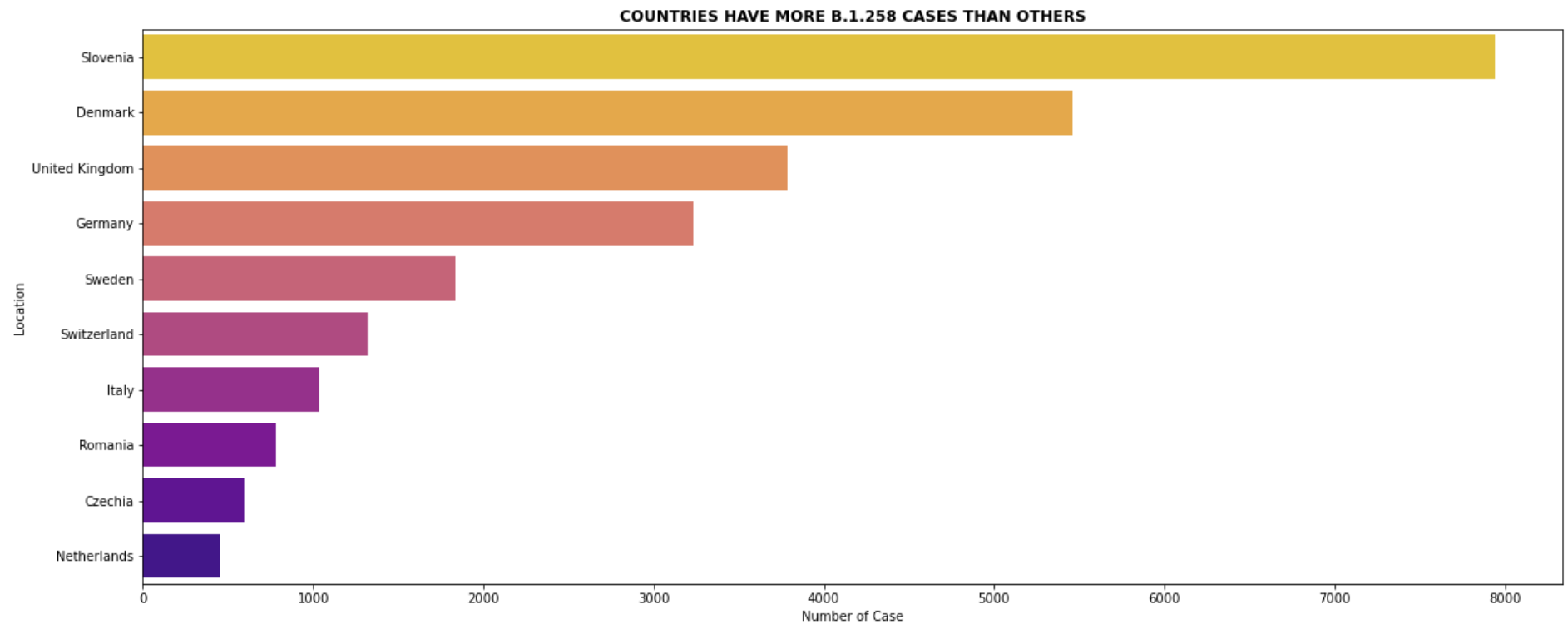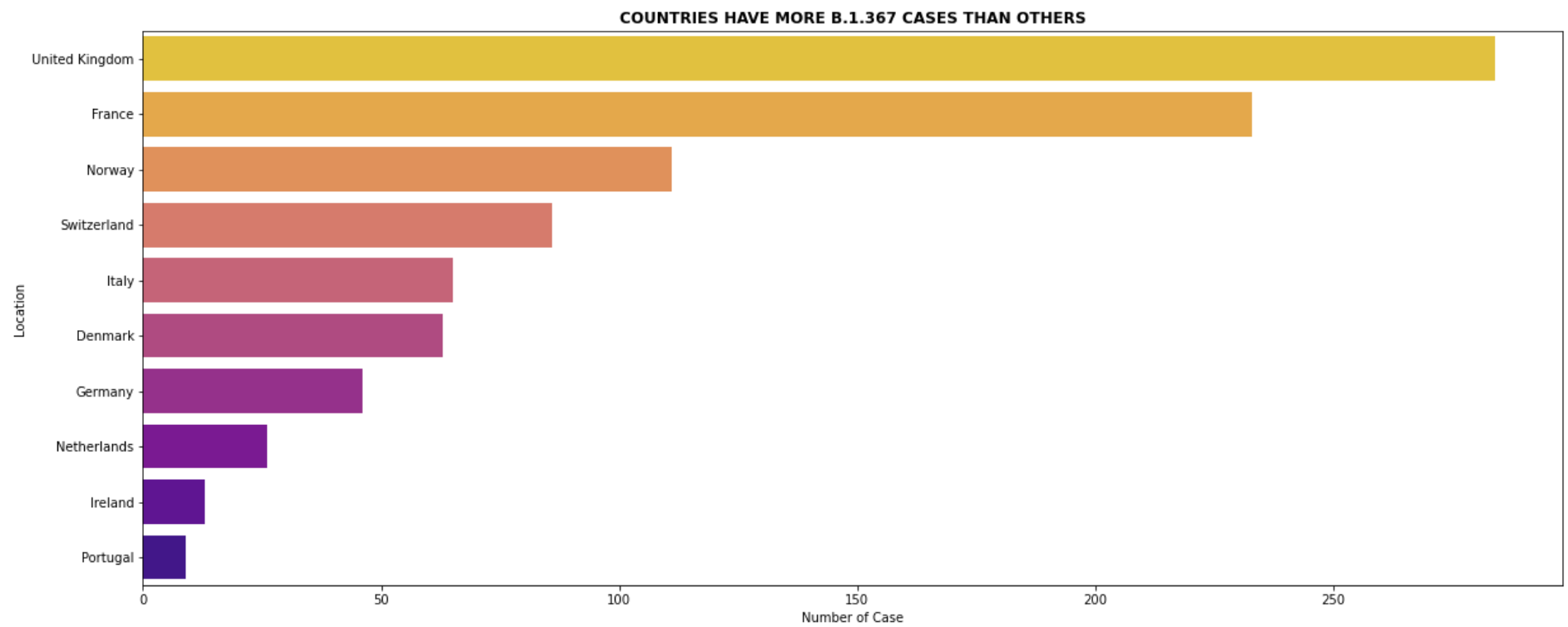memory. (To control this warning, see the rcParam `figure.max_open_warning`).
  plt.figure(figsize=(20,8))



**COUNTRIES HAVE MORE ALPHA CASES THAN OTHERS**

**COUNTRIES HAVE MORE B.1.1.277 CASES THAN OTHERS**

COUNTRIES HAVE MORE B.1.1.302 CASES THAN OTHERS

**COUNTRIES HAVE MORE B.1.1.519 CASES THAN OTHERS**

**COUNTRIES HAVE MORE B.1.160 CASES THAN OTHERS**

COUNTRIES HAVE MORE B.1.177 CASES THAN OTHERS

**COUNTRIES HAVE MORE B.1.221 CASES THAN OTHERS**

**COUNTRIES HAVE MORE B.1.258 CASES THAN OTHERS**

**COUNTRIES HAVE MORE B.1.367 CASES THAN OTHERS**

**COUNTRIES HAVE MORE B.1.620 CASES THAN OTHERS**

## COUNTRIES HAVE MORE BETA CASES THAN OTHERS

COUNTRIES HAVE MORE DELTA CASES THAN OTHERS

**COUNTRIES HAVE MORE EPSILON CASES THAN OTHERS**

**COUNTRIES HAVE MORE ETA CASES THAN OTHERS**

**COUNTRIES HAVE MORE GAMMA CASES THAN OTHERS**

**COUNTRIES HAVE MORE IOTA CASES THAN OTHERS**

**COUNTRIES HAVE MORE KAPPA CASES THAN OTHERS**

**COUNTRIES HAVE MORE LAMBDA CASES THAN OTHERS**

**COUNTRIES HAVE MORE MU CASES THAN OTHERS**

COUNTRIES HAVE MORE OMICRON CASES THAN OTHERS

**COUNTRIES HAVE MORE S:677H.ROBIN1 CASES THAN OTHERS**

**COUNTRIES HAVE MORE S:677P.PELICAN CASES THAN OTHERS**

## COUNTRIES HAVE MORE OTHERS CASES THAN OTHERS

**COUNTRIES HAVE MORE NON_WHO CASES THAN OTHERS**



In [11]:
```python
plt.figure(figsize=(14,10))
plt.xticks(rotation=90)
sns.countplot(data=cvd,x='variant', color='Red')
plt.title('variant')
```

Out[11]:  Text(0.5, 1.0, 'variant')

variant

```
In [12]:  plt.figure(figsize=(20,12))
```

```python
sns.histplot(data=cvd,x='perc_sequences',hue='variant')
plt.title('perc_sequences')
```

Out[12]:  Text(0.5, 1.0, 'perc_sequences')



```python
plt.figure(figsize=(13,9))
sns.histplot(data=cvd,x='num_sequences',hue='variant')
```

```
plt.title('num_sequences')
```

Out[13]:  Text(0.5, 1.0, 'num_sequences')



In [14]:
```
plt.figure(figsize=(30,15))
plt.xticks(rotation=89, horizontalalignment='right', fontsize=12)
plt.tight_layout()
#sns.histplot(data=cvd,x='location',bins=50)
```

```
#sns.set_context('notebook')
sns.countplot(data=cvd,x='location',color='Red')

plt.title('location')
```

Out[14]:  Text(0.5, 1.0, 'location')



In [37]:

```
#cvd['start date'] =pd.to_datatime(cvd['start date'])
```

```
#cvd['date'] = pd.to_datetime(cvd['date'], errors='ignore')
cvd
```

Out[37]:

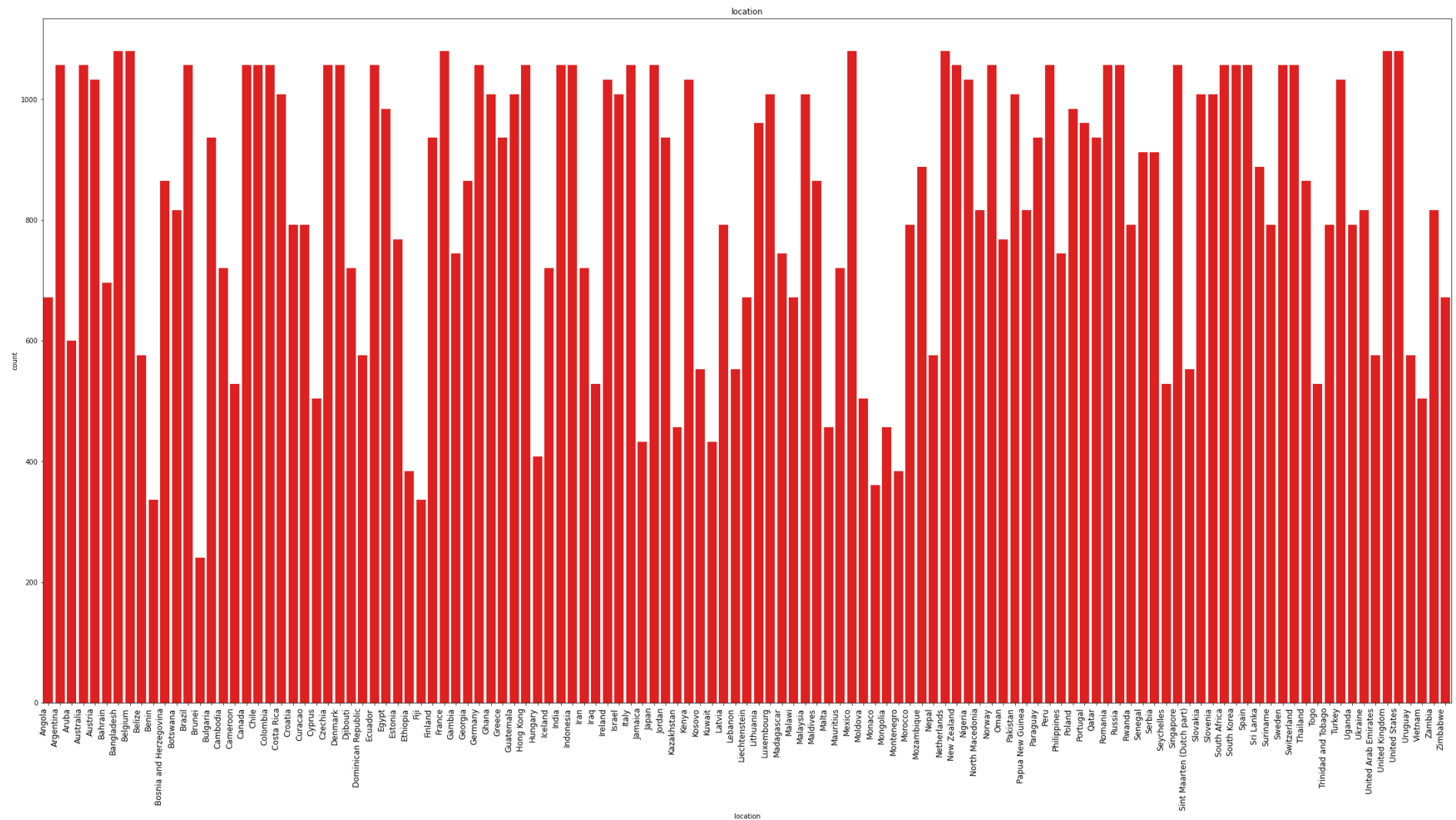| | location | variant | num_sequences | perc_sequences | num_sequences_total | month | year | day |
|---|---|---|---|---|---|---|---|---|
| **0** | Angola | Alpha | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **1** | Angola | B.1.1.277 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **2** | Angola | B.1.1.302 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **3** | Angola | B.1.1.519 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **4** | Angola | B.1.160 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| **100411** | Zimbabwe | Omicron | 0 | 0.0 | 6 | 11 | 2021 | 1 |
| **100412** | Zimbabwe | S:677H.Robin1 | 0 | 0.0 | 6 | 11 | 2021 | 1 |
| **100413** | Zimbabwe | S:677P.Pelican | 0 | 0.0 | 6 | 11 | 2021 | 1 |
| **100414** | Zimbabwe | others | 0 | 0.0 | 6 | 11 | 2021 | 1 |
| **100415** | Zimbabwe | non_who | 0 | 0.0 | 6 | 11 | 2021 | 1 |

100416 rows × 8 columns

In [48]:

```
# Seprate date with apply function

#cvd['month'] = cvd['date'].apply(lambda date: date.month)
#cvd['year'] = cvd['date'].apply(lambda date: date.year)
#cvd['day'] = cvd['date'].apply(lambda date: date.day)



# We will drop date column as we don't need keep it in our dataframe

#cvd.drop('date',axis=1, inplace=True)
cvd.head()
```

Out[48]:

|   | location | variant | num_sequences | perc_sequences | num_sequences_total | month | year | day |
|---|----------|---------|---------------|----------------|---------------------|-------|------|-----|
| **0** | Angola | Alpha | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **1** | Angola | B.1.1.277 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **2** | Angola | B.1.1.302 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **3** | Angola | B.1.1.519 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **4** | Angola | B.1.160 | 0 | 0.0 | 3 | 7 | 2020 | 6 |

In [47]:
```python
cvd.head()
```

Out[47]:

|   | location | variant | num_sequences | perc_sequences | num_sequences_total | month | year | day |
|---|----------|---------|---------------|----------------|---------------------|-------|------|-----|
| **0** | Angola | Alpha | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **1** | Angola | B.1.1.277 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **2** | Angola | B.1.1.302 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **3** | Angola | B.1.1.519 | 0 | 0.0 | 3 | 7 | 2020 | 6 |
| **4** | Angola | B.1.160 | 0 | 0.0 | 3 | 7 | 2020 | 6 |

In [50]:
```python
# Let's check all summed variant with montly ratio

cvd_val1 = cvd.loc[cvd["variant"]== virus].groupby('month')['num_sequences'].agg('sum').sort_values(ascending=False)
cvd_val1 = pd.DataFrame({'Month':cvd_val1.index, 'Number of Cases':cvd_val1.values})

plt.figure(figsize=(14,8))
sns.barplot(x='Month', y='Number of Cases',data=cvd_val1);
plt.title('Monthly Cases Ratio Of All Summed Variant',fontweight="bold");
```
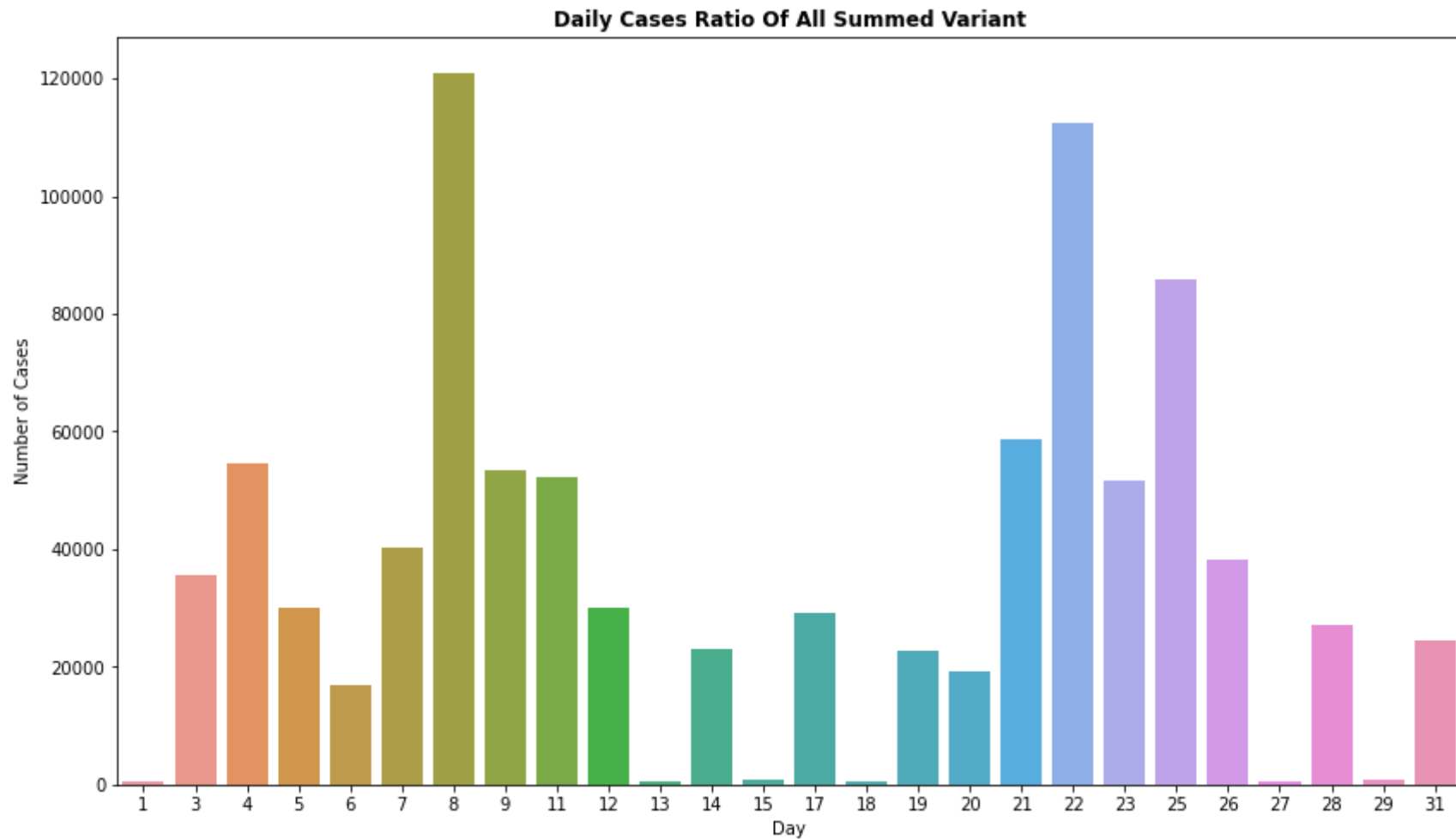
**Monthly Cases Ratio Of All Summed Variant**



In [52]:

```python
# Let's check all summed variant with daily ratio

cvd_val1 = cvd.loc[cvd["variant"]== virus].groupby('day')['num_sequences'].agg('sum').sort_values(ascending=False)
cvd_val1 = pd.DataFrame({'Day':cvd_val1.index, 'Number of Cases':cvd_val1.values})

plt.figure(figsize=(14,8))
sns.barplot(x='Day', y='Number of Cases',data=cvd_val1);
plt.title('Daily Cases Ratio Of All Summed Variant',fontweight="bold");
```
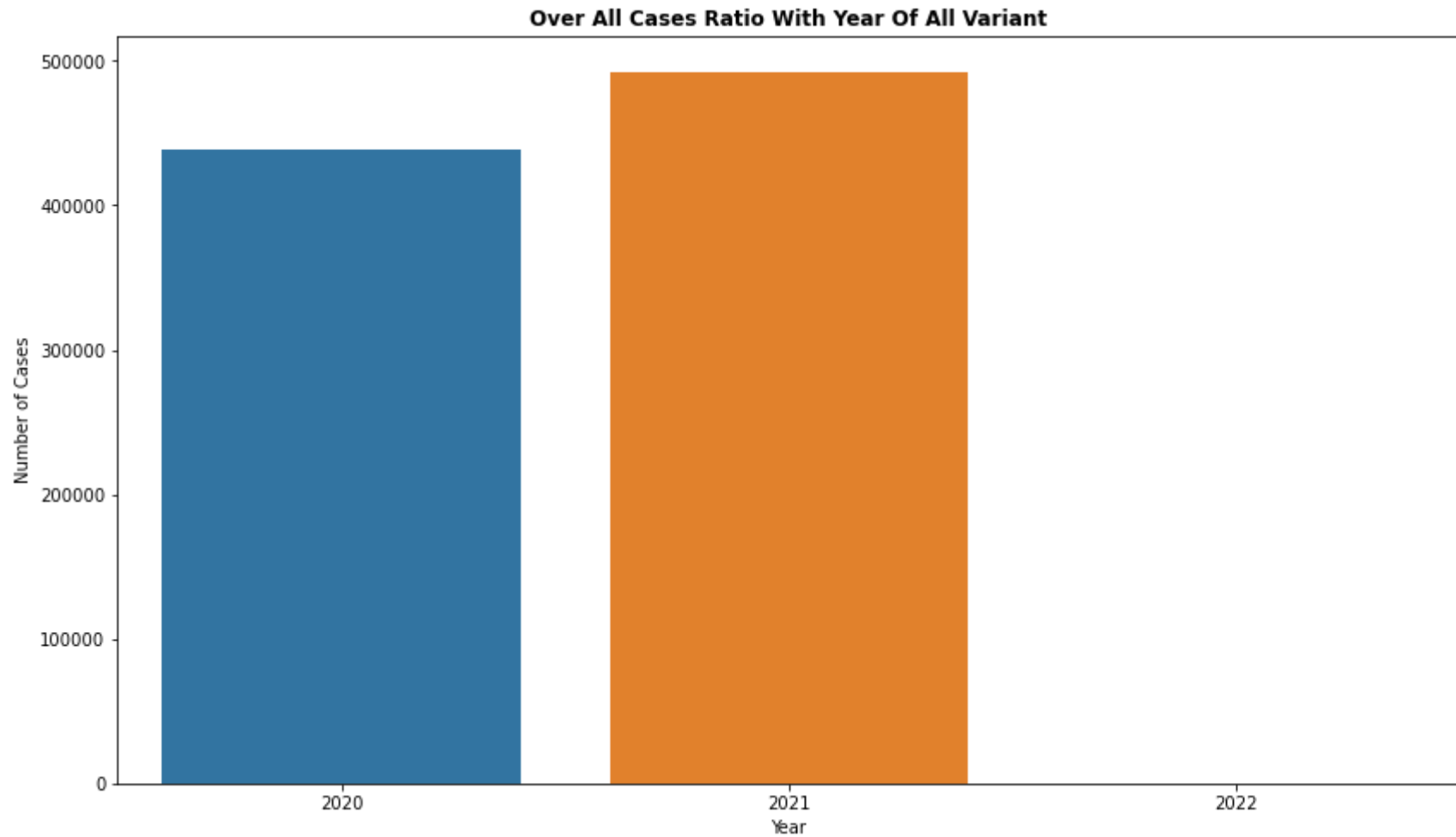
**Daily Cases Ratio Of All Summed Variant**



```
In [54]:    # Let's check all summed variant with yearly ratio

            cvd_val1 = cvd.loc[cvd["variant"]== virus].groupby('year')['num_sequences'].agg('sum').sort_values(ascending=False)
            cvd_val1 = pd.DataFrame({'Year':cvd_val1.index, 'Number of Cases':cvd_val1.values})

            plt.figure(figsize=(14,8))
            sns.barplot(x='Year', y='Number of Cases',data=cvd_val1);
            plt.title('Over All Cases Ratio With Year Of All Variant',fontweight="bold");
```
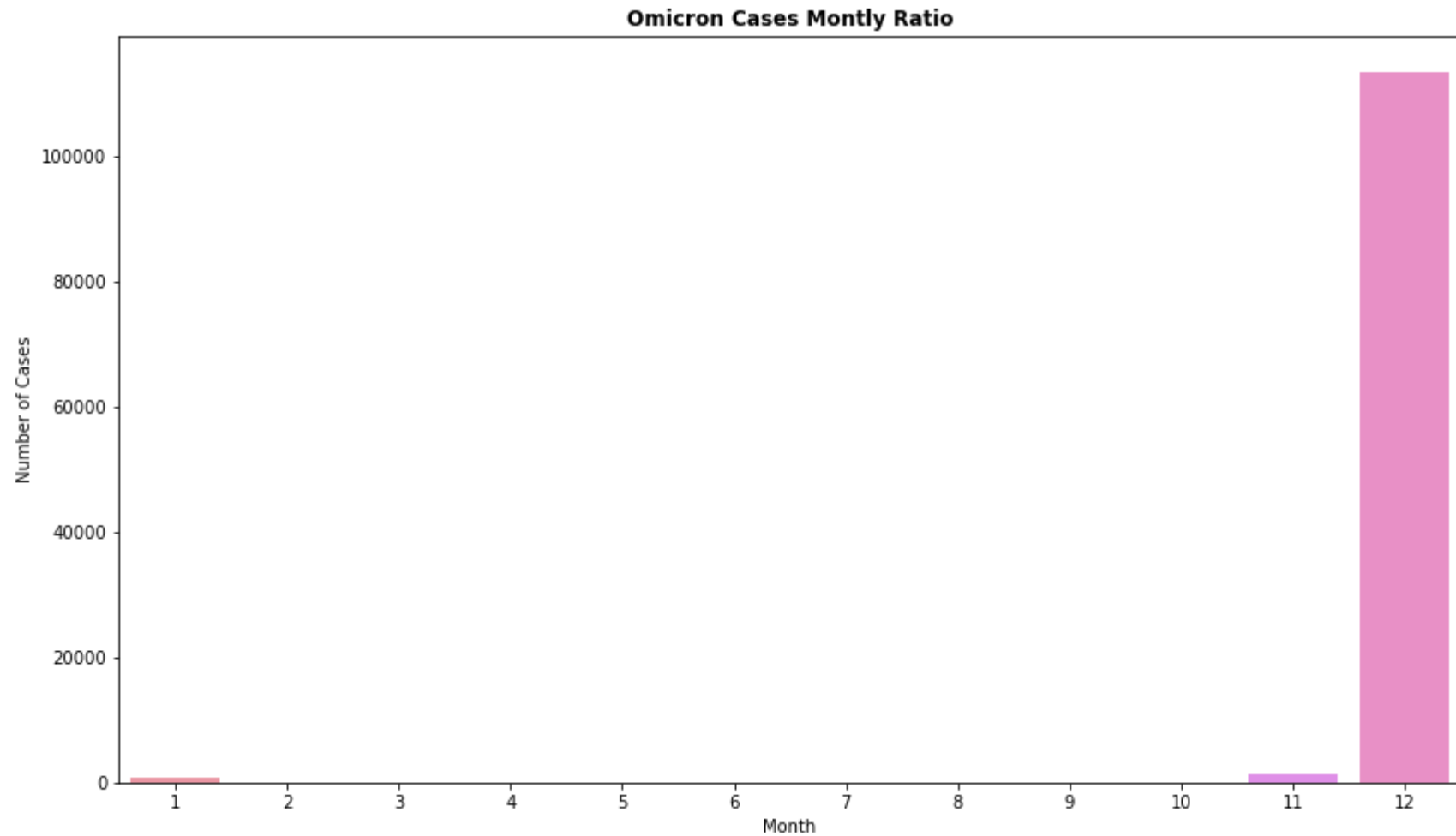
**Over All Cases Ratio With Year Of All Variant**



Now Let's check the Omicron variant

In [57]:

```python
cvd_val1 = cvd.loc[cvd["variant"]== 'Omicron'].groupby('month')['num_sequences'].agg('sum').sort_values(ascending=False)
cvd_val1 = pd.DataFrame({'Month':cvd_val1.index, 'Number of Cases':cvd_val1.values})

plt.figure(figsize=(14,8))
sns.barplot(x='Month', y='Number of Cases',data=cvd_val1);
plt.title('Omicron Cases Montly Ratio',fontweight="bold");
```

**Omicron Cases Montly Ratio**
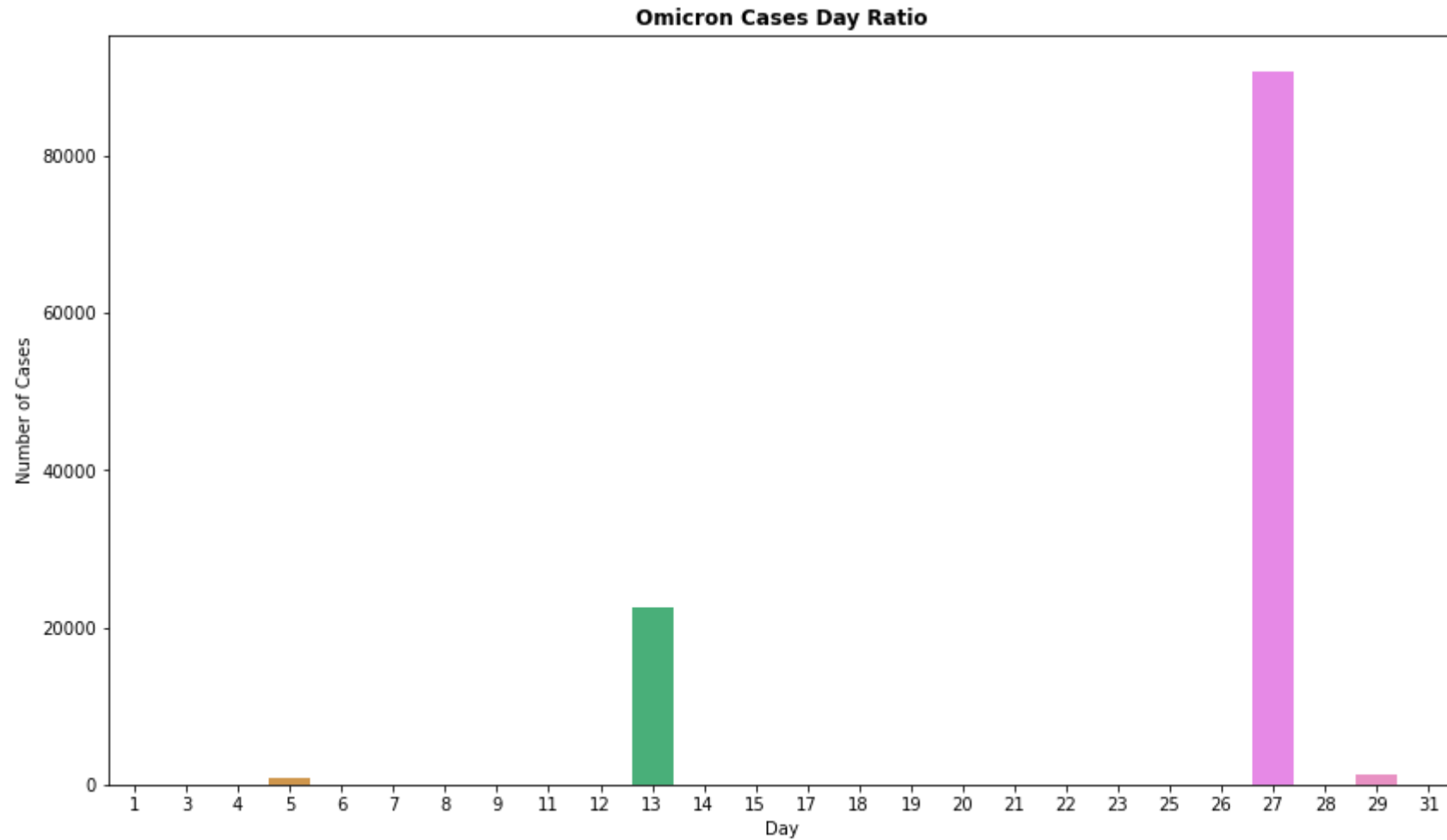


## Days

```
In [58]:   cvd_val1 = cvd.loc[cvd["variant"]== 'Omicron'].groupby('day')['num_sequences'].agg('sum').sort_values(ascending=False)
           cvd_val1 = pd.DataFrame({'Day':cvd_val1.index, 'Number of Cases':cvd_val1.values})

           plt.figure(figsize=(14,8))
           sns.barplot(x='Day', y='Number of Cases',data=cvd_val1);
           plt.title('Omicron Cases Day Ratio',fontweight="bold")
```

```
Out[58]:   Text(0.5, 1.0, 'Omicron Cases Day Ratio')
```
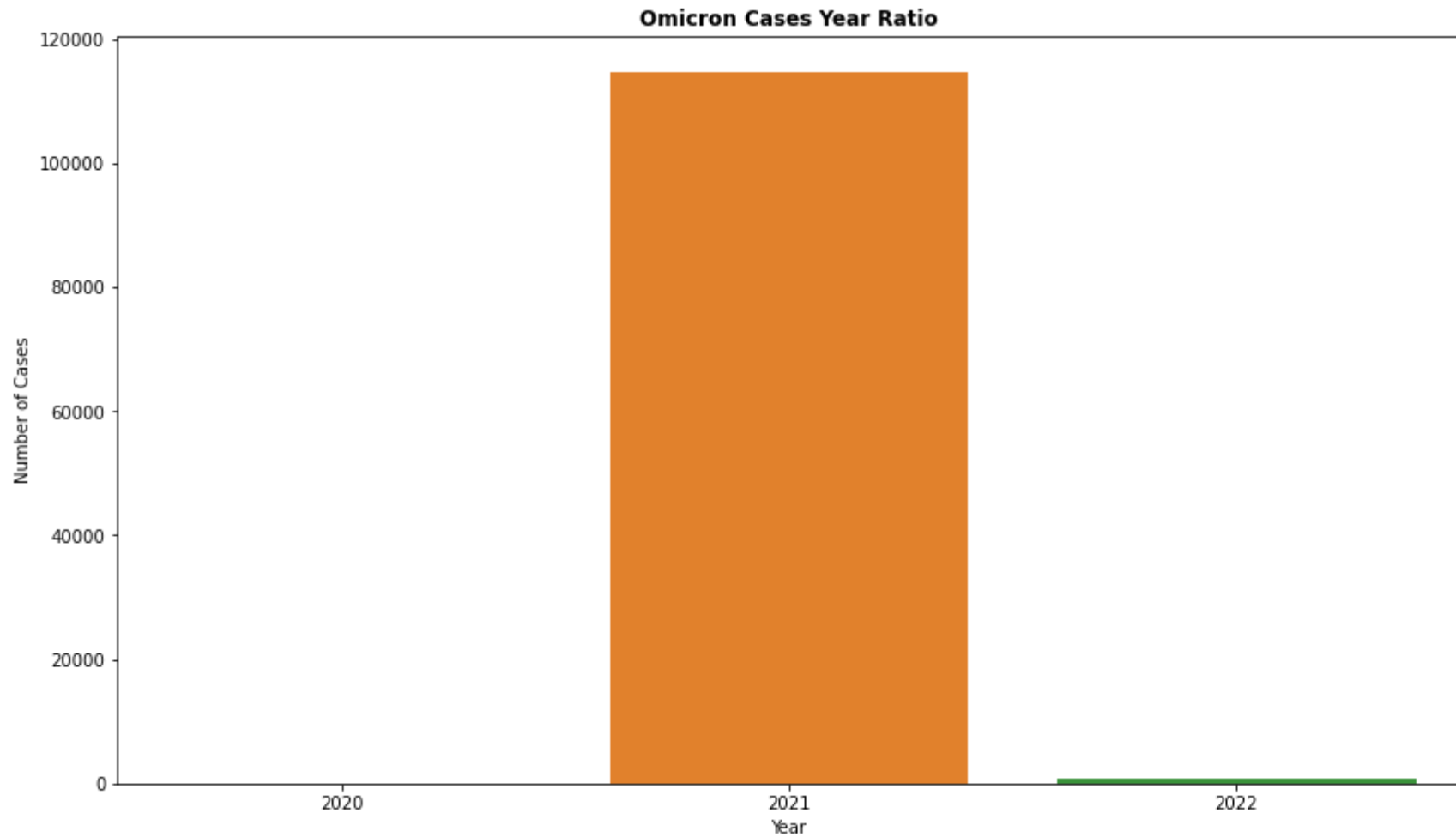
**Omicron Cases Day Ratio**



## year

```
In [60]:  cvd_val1 = cvd.loc[cvd["variant"]== 'Omicron'].groupby('year')['num_sequences'].agg('sum').sort_values(ascending=False)
          cvd_val1 = pd.DataFrame({'Year':cvd_val1.index, 'Number of Cases':cvd_val1.values})

          plt.figure(figsize=(14,8))
          sns.barplot(x='Year', y='Number of Cases',data=cvd_val1);
          plt.title('Omicron Cases Year Ratio',fontweight="bold")
```

```
Out[60]:  Text(0.5, 1.0, 'Omicron Cases Year Ratio')
```

**Omicron Cases Year Ratio**



In [66]:
```python
cvd_val1 = cvd.loc[cvd["variant"]== 'Omicron'].groupby('location')['num_sequences'].agg('sum').sort_values(ascending=False)[:12]
cvd_val1 = pd.DataFrame({'Location':cvd_val1.index, 'Number of Cases':cvd_val1.values})

plt.figure(figsize=(16,8))
sns.barplot(x='Location', y='Number of Cases',data=cvd_val1);
plt.title('Highest Omicron Cases Location',fontweight="bold");
plt.xticks(rotation=90);
```

**Highest Omicron Cases Location**