

```
In [1]: 1 #Importing Modules
2 import numpy as np
3 import pandas as pd
4 import scipy.stats as stats
5 import matplotlib.pyplot as plt
6 from math import sqrt
7 import datetime as dt
8 from scipy.stats import kde
```

## 1. Use the state data generated in Stage II to fit a distribution to the number of COVID-19 new cases using any of MoM, MLE, and KDE methods.

- Graphically plot the distribution
- Describe the type of distribution and its statistics(Moments of distribution - center, variance, skewness, kurtosis) in the notebook
- Compare the distribution and its statistics to 3 other states of your choosing. Describe if the distributions look different and what does that imply.

```
In [2]: 1 #Reading Transformed Data from Stage II with new cases column
2 transformed_data = pd.read_csv('../Stage2/transformed_data_with_new_values.csv', parse_dates=['Date'])
3 transformed_data.head()
```

```
Out[2]:
```

	Date	Week	countyFIPS	County_Name	State	StateFIPS	population	Cases	New_Cases	Deaths	New_Deaths
0	2022-06-01	22	1001	Autauga County	AL	1	55869	15969	6	216	0
1	2022-06-01	22	1003	Baldwin County	AL	1	223234	56580	68	683	0
2	2022-06-01	22	1005	Barbour County	AL	1	24686	5710	3	99	0
3	2022-06-01	22	1007	Bibb County	AL	1	22394	6508	8	105	0
4	2022-06-01	22	1009	Blount County	AL	1	57826	15077	4	244	0

```
In [3]: 1 transformed_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 672388 entries, 0 to 672387
Data columns (total 11 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Date            672388 non-null  datetime64[ns]
1   Week            672388 non-null  int64
2   countyFIPS      672388 non-null  int64
3   County_Name     672388 non-null  object
4   State           672388 non-null  object
5   StateFIPS       672388 non-null  int64
6   population      672388 non-null  int64
7   Cases           672388 non-null  int64
8   New_Cases       672388 non-null  int64
9   Deaths         672388 non-null  int64
10  New_Deaths      672388 non-null  int64
dtypes: datetime64[ns](1), int64(8), object(2)
memory usage: 56.4+ MB
```

Let us take the data for NC State first and plot the histogram for New\_Cases value to identify the distribution

```
In [4]: 1 #Filtering NC State data
2 NC_transformed_data = transformed_data.query("State=='NC']").reset_index().drop(columns='index').copy()
3 NC_transformed_data.head()
```

```
Out[4]:
```

	Date	Week	countyFIPS	County_Name	State	StateFIPS	population	Cases	New_Cases	Deaths	New_Deaths
0	2022-06-01	22	37001	Alamance County	NC	37	169509	49188	462	488	0
1	2022-06-01	22	37003	Alexander County	NC	37	37497	10600	41	139	0
2	2022-06-01	22	37005	Alleghany County	NC	37	11137	3041	27	16	0
3	2022-06-01	22	37007	Anson County	NC	37	24446	6672	43	101	0
4	2022-06-01	22	37009	Ashe County	NC	37	27203	6575	43	79	0

```
In [5]: 1 NC_transformed_data.New_Cases.min()
```

```
Out[5]: 0
```

```
In [6]: 1 NC_transformed_data.New_Cases.max()
```

```
Out[6]: 5266
```

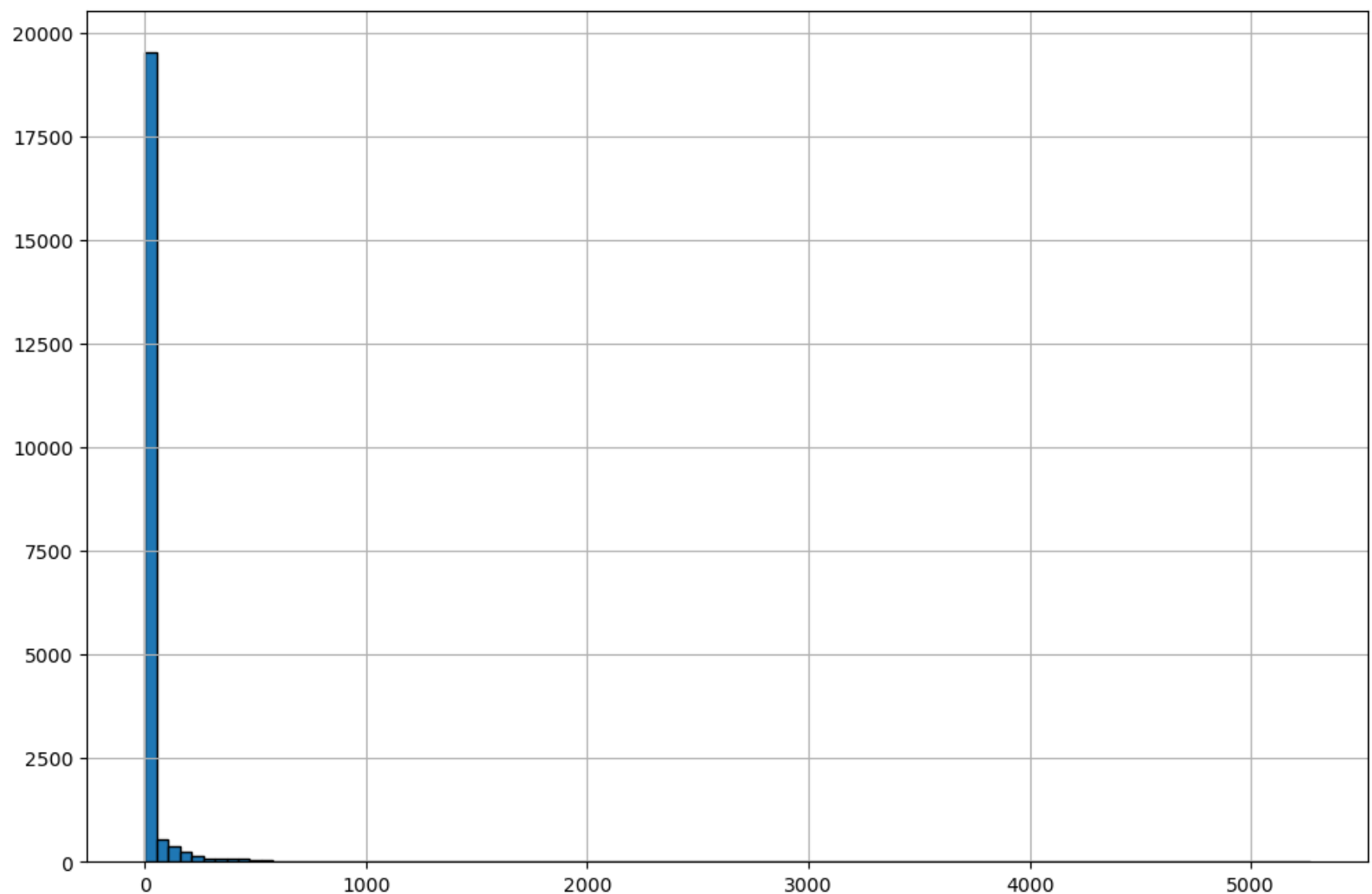
```
In [7]: 1 NC_transformed_data.shape
```

```
Out[7]: (21400, 11)
```

So we have 21400 records of new cases for NC state with values ranging from 0 to 5266

```
In [8]: 1 #Plotting histogram
2 NC_transformed_data.New_Cases.hist(bins=100, ec='black', figsize=(12,8))
```

```
Out[8]: <AxesSubplot:>
```



Here we can see that

1. the data is highly positive skewed
2. With tail to the right
3. Data contains non negative discrete values (Since covid cases can only be an integer  $n$  such that  $n \geq 0$ )

With all these information, Poisson Distribution suits well for this case. We can term it as the number of cases occurred during a given day period.

### Fitting Poisson Distribution using MLE

To use Maximum Likelihood Estimator, we need to solve the maximization value for poisson distribution

$$\frac{\lambda^x e^{-\lambda}}{x!}$$

From the link you provided, [MLE for Different Distributions](https://nbviewer.org/github/rasbt/pattern_classification/blob/master/parameter_estimation_techniques/max_likelihood_est_distributions.ipynb?create=1)

([https://nbviewer.org/github/rasbt/pattern\\_classification/blob/master/parameter\\_estimation\\_techniques/max\\_likelihood\\_est\\_distributions.ipynb?create=1](https://nbviewer.org/github/rasbt/pattern_classification/blob/master/parameter_estimation_techniques/max_likelihood_est_distributions.ipynb?create=1)) We can see the Maximum Likelihood Estimator for Poisson distribution is as follows:

$$\theta = \frac{\sum_{k=1}^n x_k}{n}$$

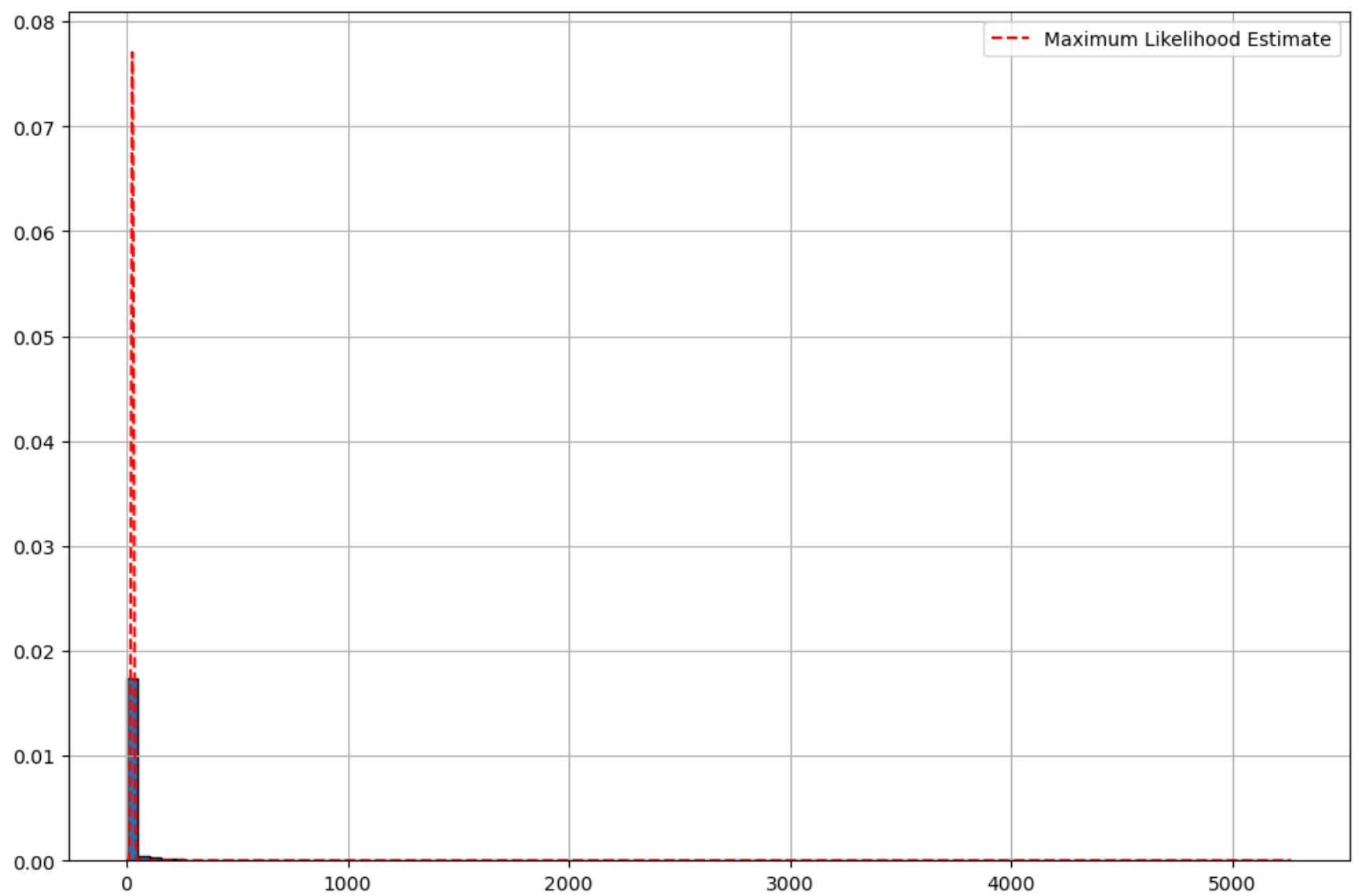
This is simply mean of the values

```
In [9]: 1 # Finding Mean of all the new cases
2 theta = NC_transformed_data.New_Cases.mean()
3 theta
```

```
Out[9]: 26.796728971962615
```

```
In [10]: 1 x= np.arange(0,NC_transformed_data.New_Cases.max())
2 NC_transformed_data.New_Cases.hist(density=True, bins=100, ec='black', figsize=(12,8))
3 l1=plt.plot(x, stats.poisson.pmf(x, theta), 'r--', label="Maximum Likelihood Estimate")
4 plt.legend(handles=[l1])
```

Out[10]: <matplotlib.legend.Legend at 0x1bb81020d90>

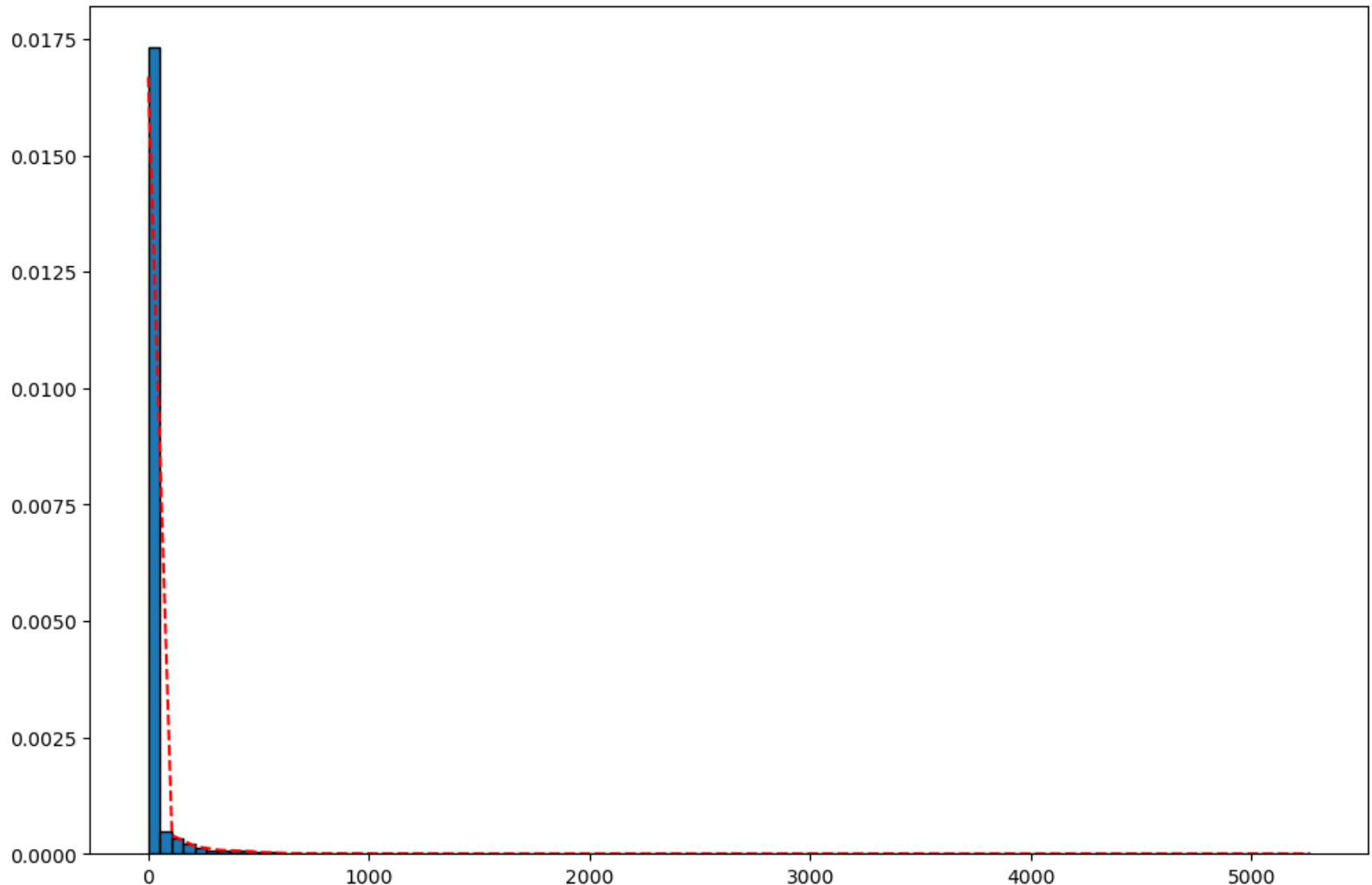


Here we can see that the distribution plot looks almost similar to the Histogram data. But the density plot didn't fit well. So let us try to fit using KDE

```
In [11]: 1 # Fitting density plot using KDE
2 NC_transformed_data.New_Cases.hist(density=True, bins=100, grid=False, ec='black', figsize=(12,8))
3 x = np.linspace(0, NC_transformed_data.New_Cases.max())
4 density = kde.gaussian_kde(NC_transformed_data.New_Cases, bw_method=None)
5 xgrid = np.linspace(x.min(), x.max(), 100)
6 plt.plot(x, density(x), 'r--')
```

C:\Users\Dell\AppData\Local\Temp\ipykernel\_4452\3686681080.py:4: DeprecationWarning: Please use `gaussian\_kde` from the `scipy.stats` namespace, the `scipy.stats.kde` namespace is deprecated.  
density = kde.gaussian\_kde(NC\_transformed\_data.New\_Cases, bw\_method=None)

Out[11]: [<matplotlib.lines.Line2D at 0x1bb81193b80>]



KDE plot seems to fit more appropriately for this data. Now let us look at the statistics of Poisson distribution for the given data

```
In [12]: 1 # Calculating statistics of the distribution using mathematical formulae
2 # mean = theta
3 # variance = theta
4 # skewness = 1/sqrt(theta)
5 # kurtosis = 1/theta
6 print(f"mean of the poisson distribution : {theta}")
7 print(f"Variance of the poisson distribution : {theta}")
8 print(f"Skewness of the poisson distribution : {1/sqrt(theta)}")
9 print(f"Kurtosis of the poisson distribution : {1/theta}")
```

mean of the poisson distribution : 26.796728971962615  
Variance of the poisson distribution : 26.796728971962615  
Skewness of the poisson distribution : 0.19317864172521376  
Kurtosis of the poisson distribution : 0.0373179876187985

```
In [13]: 1 # calculating statistics of the distribution using scipy library.
2 NC_mean, NC_var, NC_skew, NC_kurt = stats.poisson.stats(theta, moments='mvsk')
3 print(f"mean of the poisson distribution : {NC_mean}")
4 print(f"Variance of the poisson distribution : {NC_var}")
5 print(f"Skewness of the poisson distribution : {NC_skew}")
6 print(f"Kurtosis of the poisson distribution : {NC_kurt}")
```

mean of the poisson distribution : 26.796728971962615  
Variance of the poisson distribution : 26.796728971962615  
Skewness of the poisson distribution : 0.19317864172521376  
Kurtosis of the poisson distribution : 0.0373179876187985

Now let us do the similar analysis for three other states and compare the values

```
In [14]: 1 #Filtering CA, WA, NY State data
2 Three_state_transformed_data = transformed_data.query("State in ['CA','WA','NY']").reset_index().drop(columns
3 Three_state_transformed_data.head()
```

Out[14]:

	Date	Week	countyFIPS	County_Name	State	StateFIPS	population	Cases	New_Cases	Deaths	New_Deaths
0	2022-06-01	22	6001	Alameda County	CA	6	1671329	285709	658	1870	0
1	2022-06-01	22	6003	Alpine County	CA	6	1129	128	0	0	0
2	2022-06-01	22	6005	Amador County	CA	6	39752	8820	3	87	0
3	2022-06-01	22	6007	Butte County	CA	6	219186	34122	17	427	0
4	2022-06-01	22	6009	Calaveras County	CA	6	45905	7522	8	121	0

```
In [15]: 1 Three_state_transformed_data.New_Cases.min()
```

Out[15]: 0

```
In [16]: 1 Three_state_transformed_data.New_Cases.max()
```

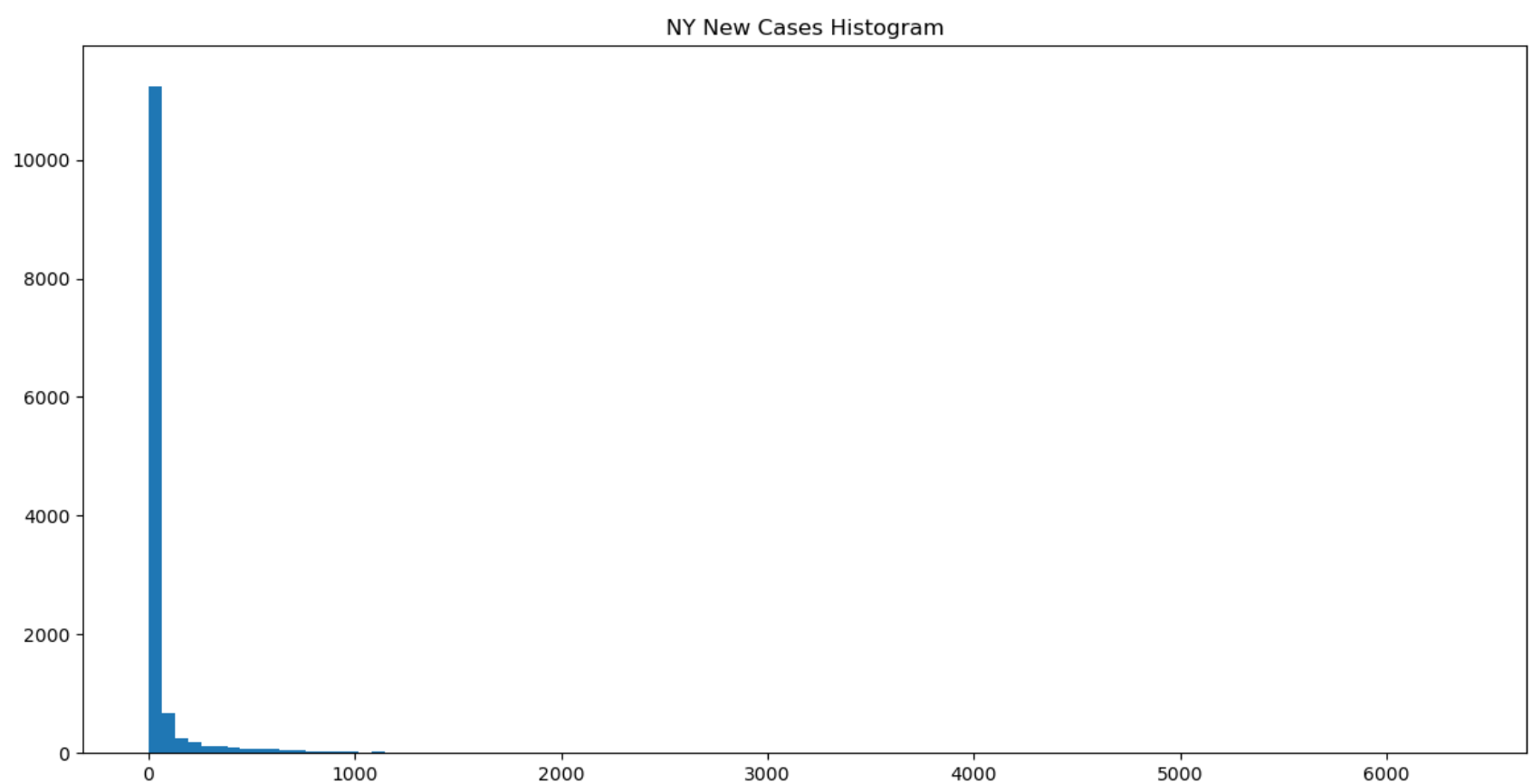
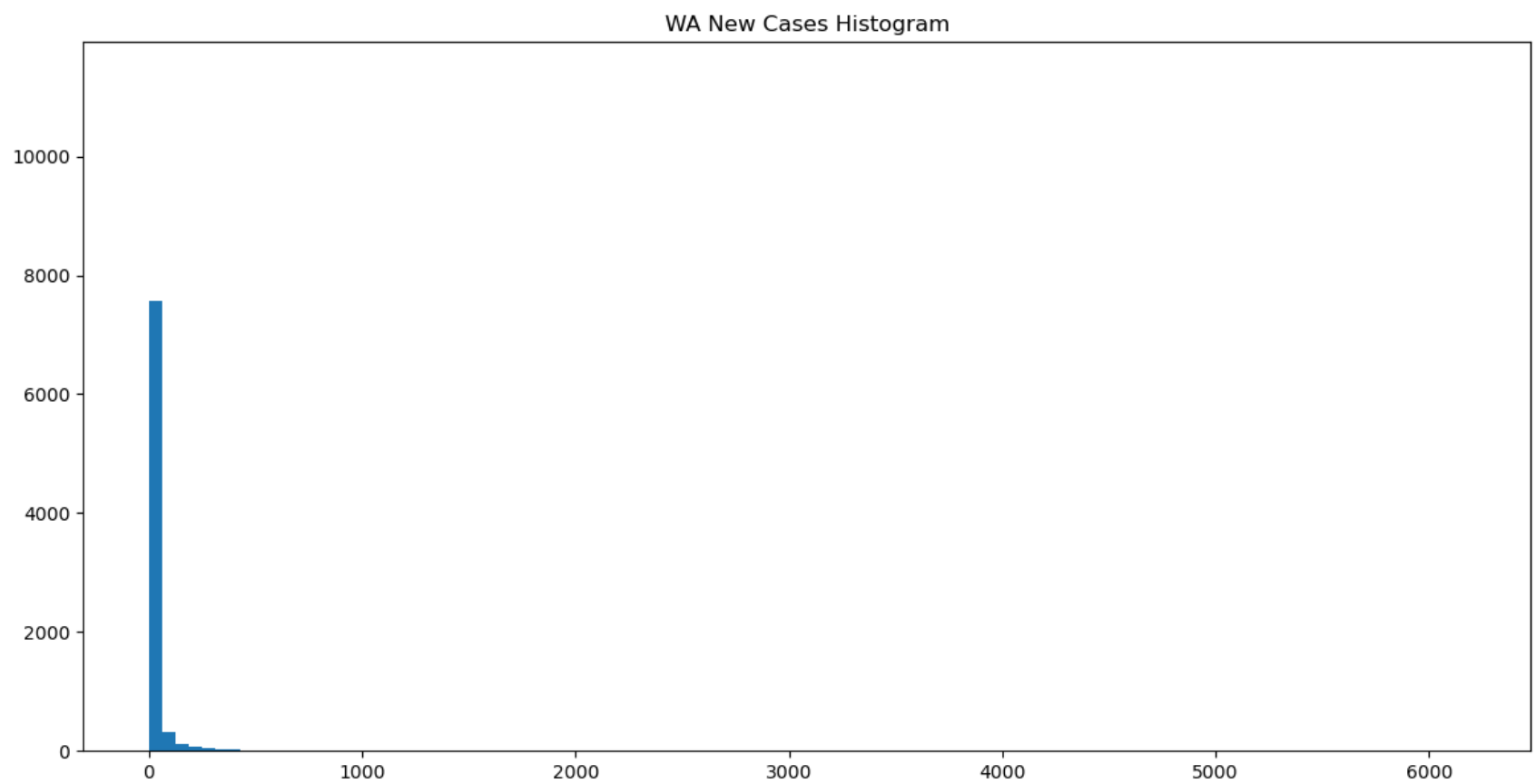
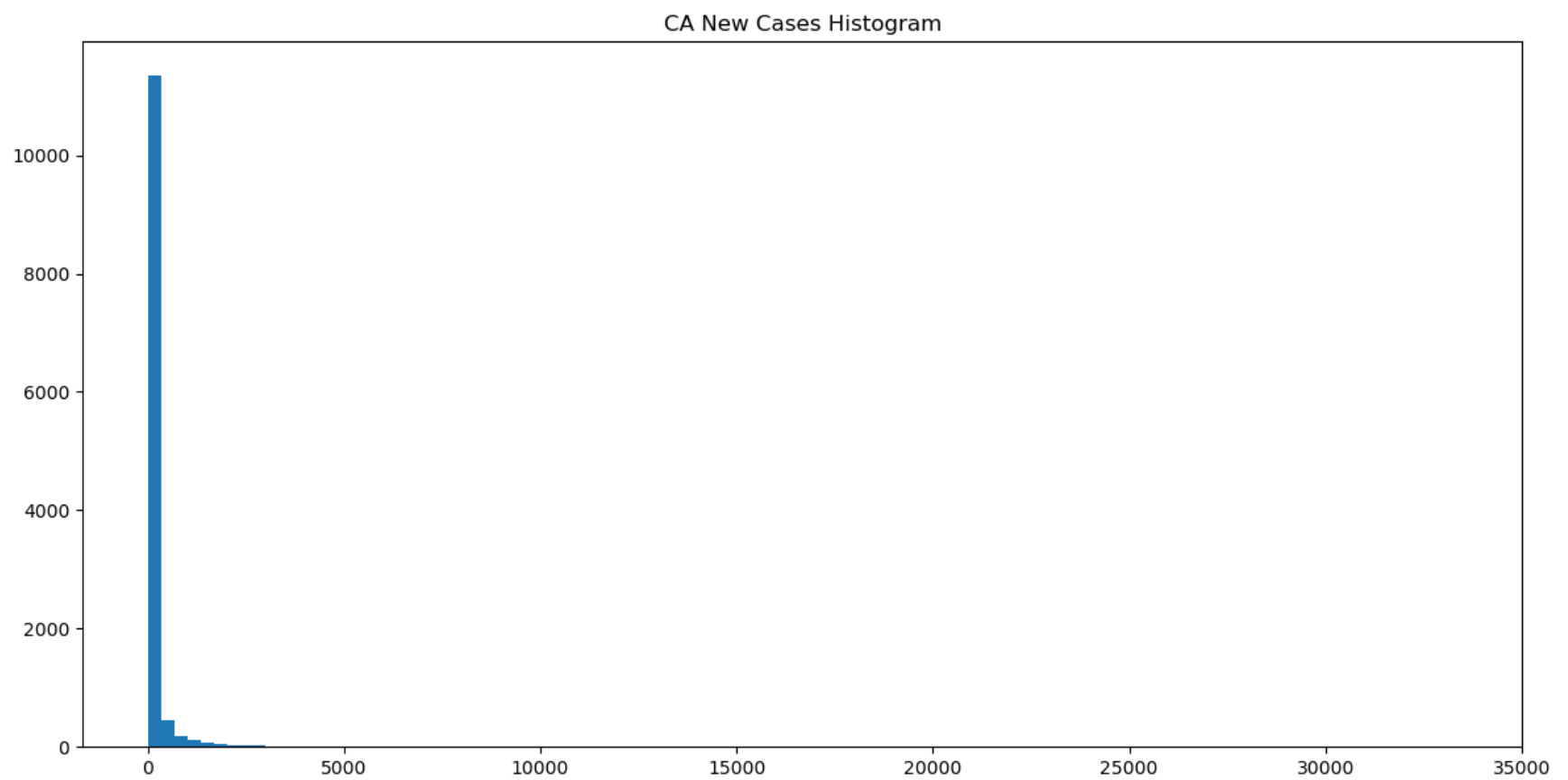
Out[16]: 33349

```
In [17]: 1 Three_state_transformed_data.shape
```

Out[17]: (34026, 11)

We have 34026 records of new cases for three states CA, NY and WA combinedly with values ranging from 0 to 33349

```
In [18]: 1 #plotting histogram subplots state wise
2 states = ['CA', 'WA', 'NY']
3 fig, axs = plt.subplots(len(states), 1, sharey=True, tight_layout=True, figsize=(12,18))
4 for i, state in enumerate(states):
5     axs[i].hist(Three_state_transformed_data.query(f"State=='{state}'").New_Cases, bins=100)
6     axs[i].set_title(f"{state} New Cases Histogram")
7 fig.tight_layout()
```

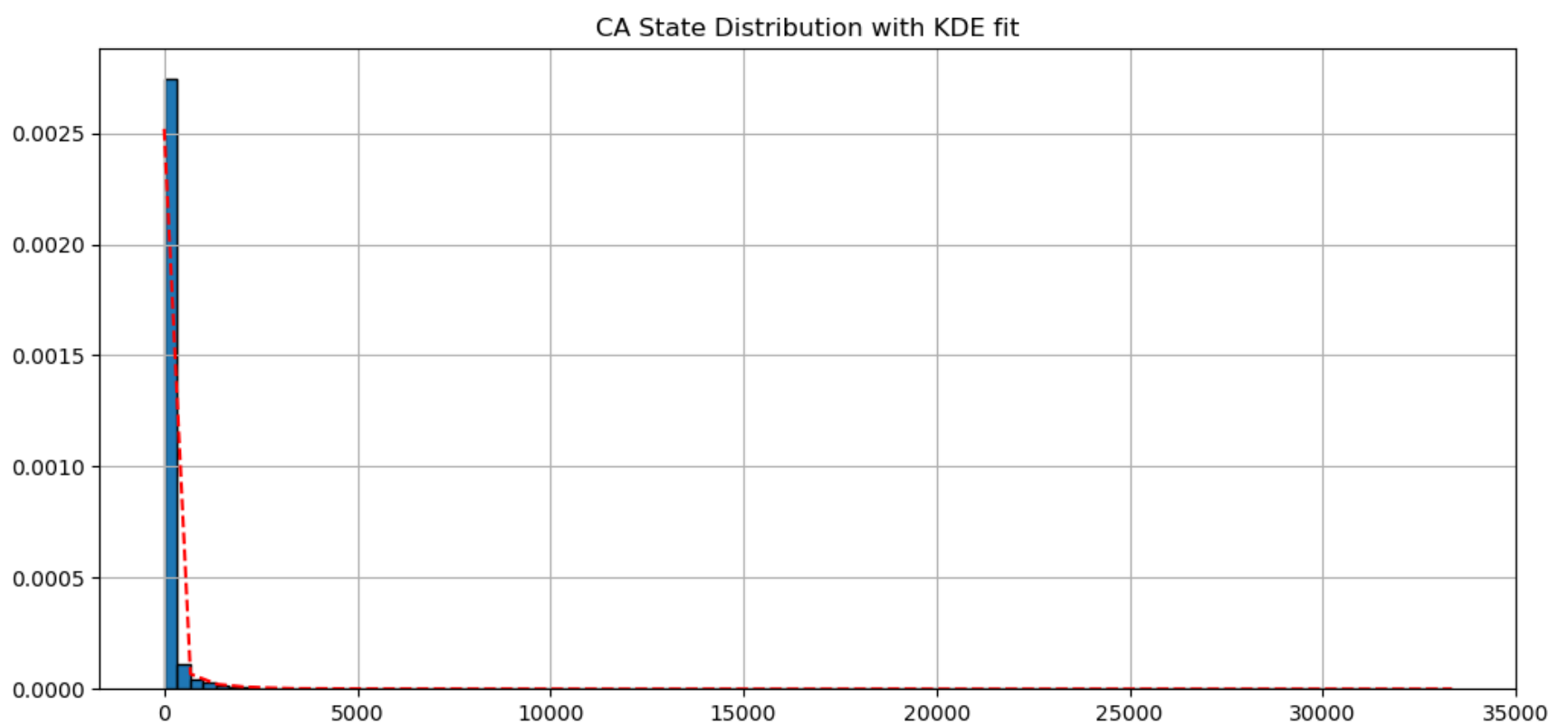
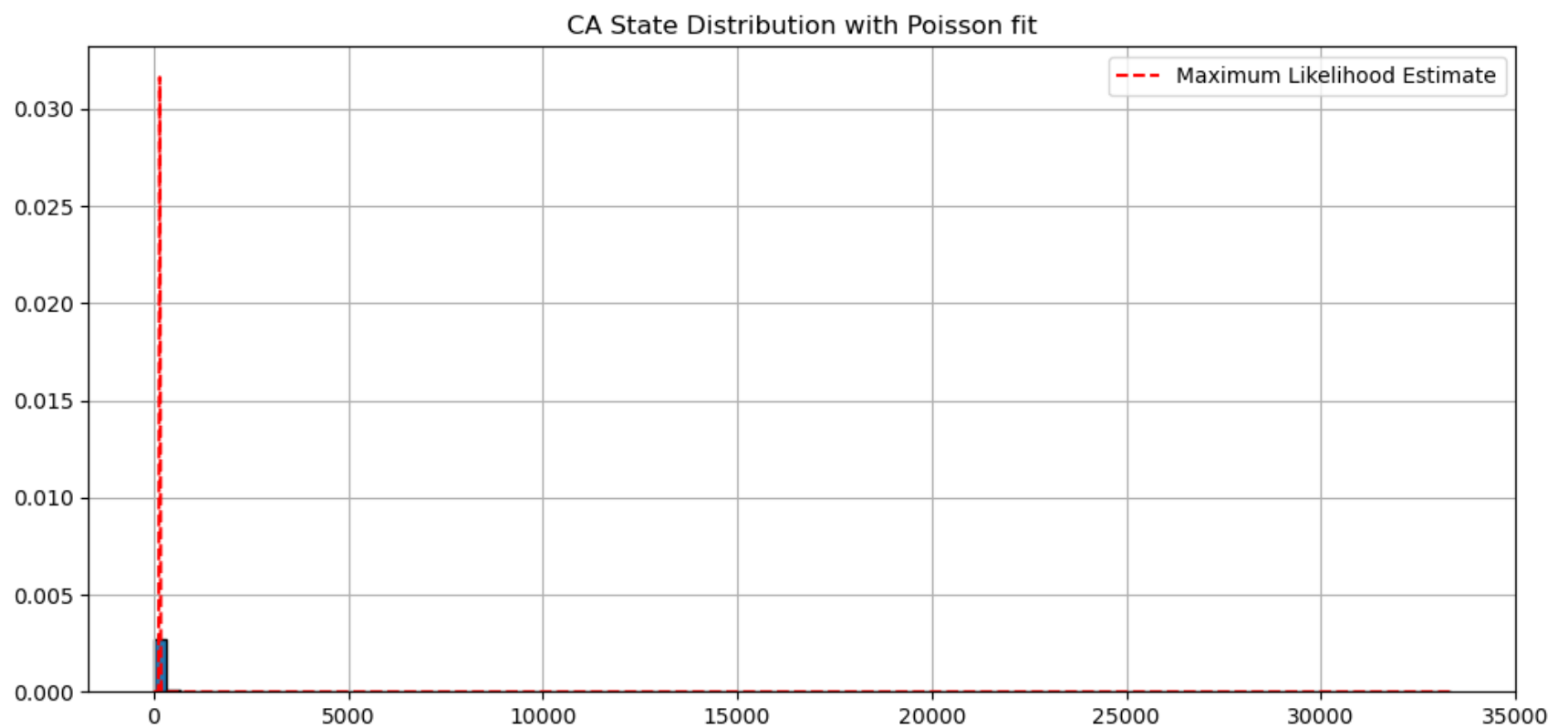


All these looks similar to the NC state data with a positive skewness and tail on the right side. Hence fitting poisson distribution to all these states as well

```
In [19]: 1 def plot_state_distributions(state):
2         new_cases = Three_state_transformed_data.query(f"State=='{state}'").New_Cases
3         theta = new_cases.mean()
4         x = np.arange(new_cases.max())
5         plt.figure(figsize=(12,24))
6
7         plt.subplot(2,1,1)
8         new_cases.hist(density=True, bins=100, ec='black', figsize=(12,12))
9         l1=plt.plot(x, stats.poisson.pmf(x, theta), 'r--', label="Maximum Likelihood Estimate")
10        plt.title(f'{state} State Distribution with Poisson fit')
11        plt.legend(handles=[l1])
12
13
14        plt.subplot(2,1,2)
15        new_cases.hist(density=True, bins=100, ec='black', figsize=(12,12))
16        x = np.linspace(0, new_cases.max())
17        density = kde.gaussian_kde(new_cases, bw_method=None)
18        xgrid = np.linspace(x.min(), x.max(), 100)
19        plt.plot(x, density(x), 'r--')
20        plt.title(f'{state} State Distribution with KDE fit')
```

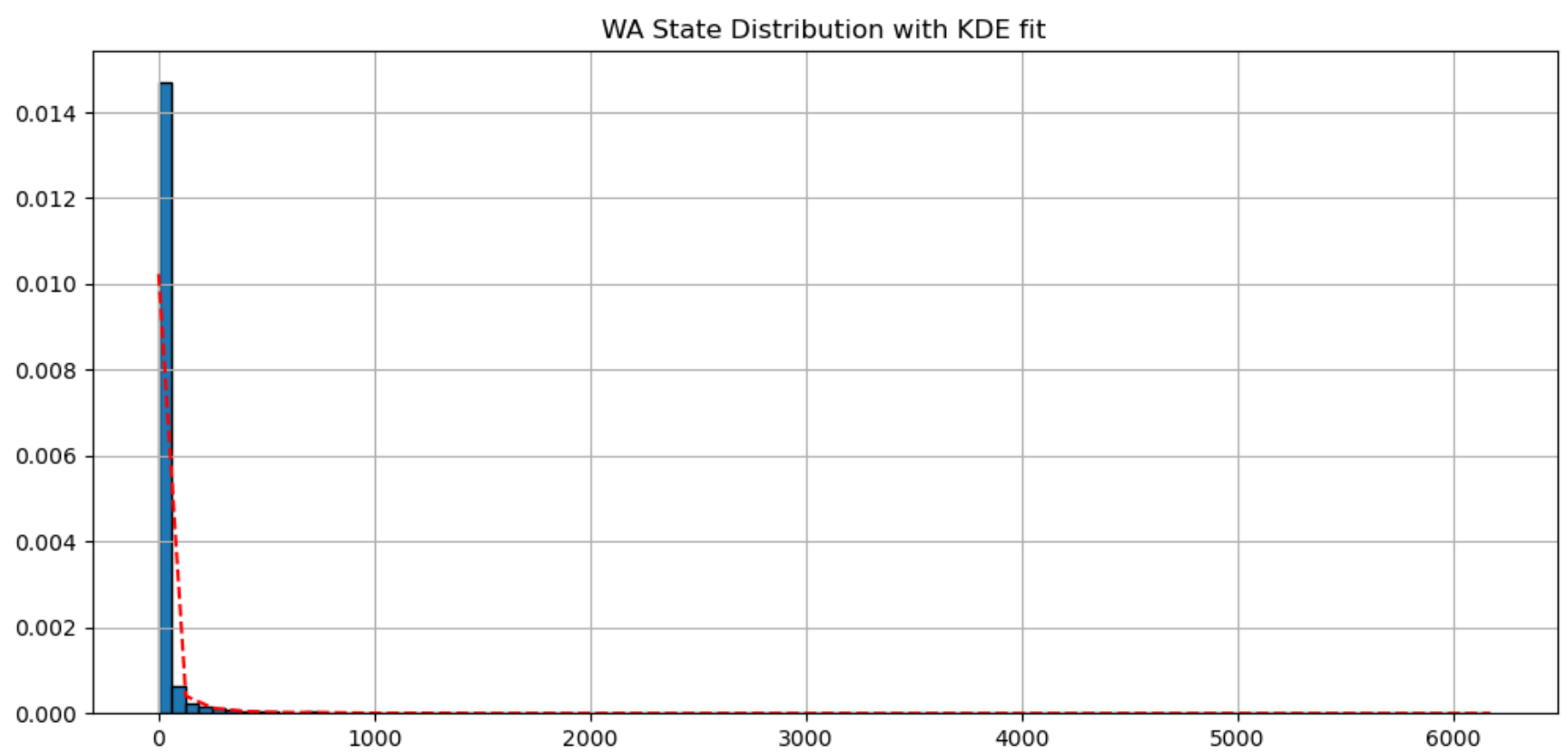
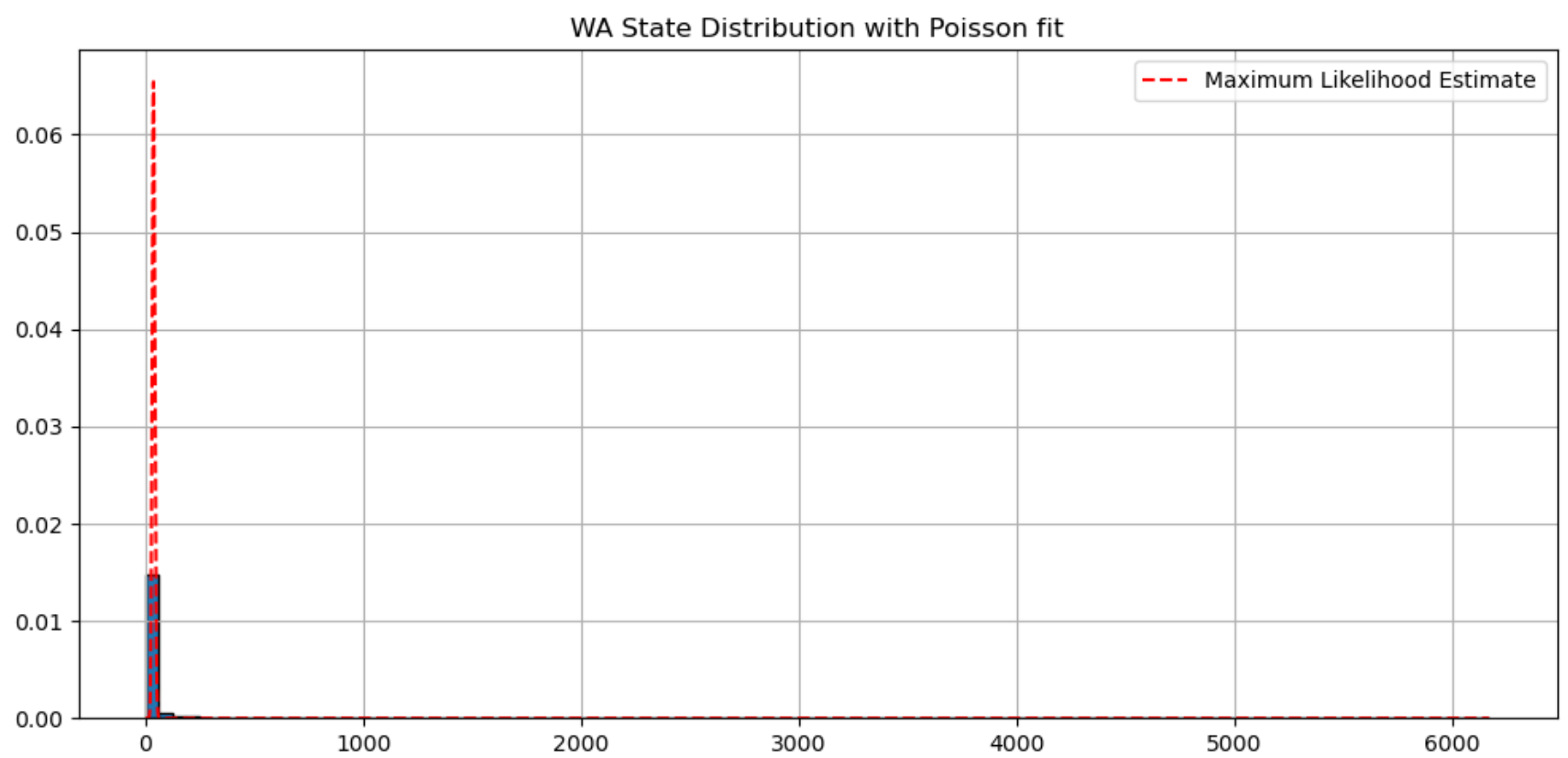
```
In [20]: 1 #Plotting Distribution for CA state
2         plot_state_distributions('CA')
```

C:\Users\Dell\AppData\Local\Temp\ipykernel\_4452\3660962002.py:17: DeprecationWarning: Please use `gaussian\_kde` from the `scipy.stats` namespace, the `scipy.stats.kde` namespace is deprecated.  
density = kde.gaussian\_kde(new\_cases, bw\_method=None)



```
In [21]: 1 #Plotting Distribution for WA state
        2 plot_state_distributions('WA')
```

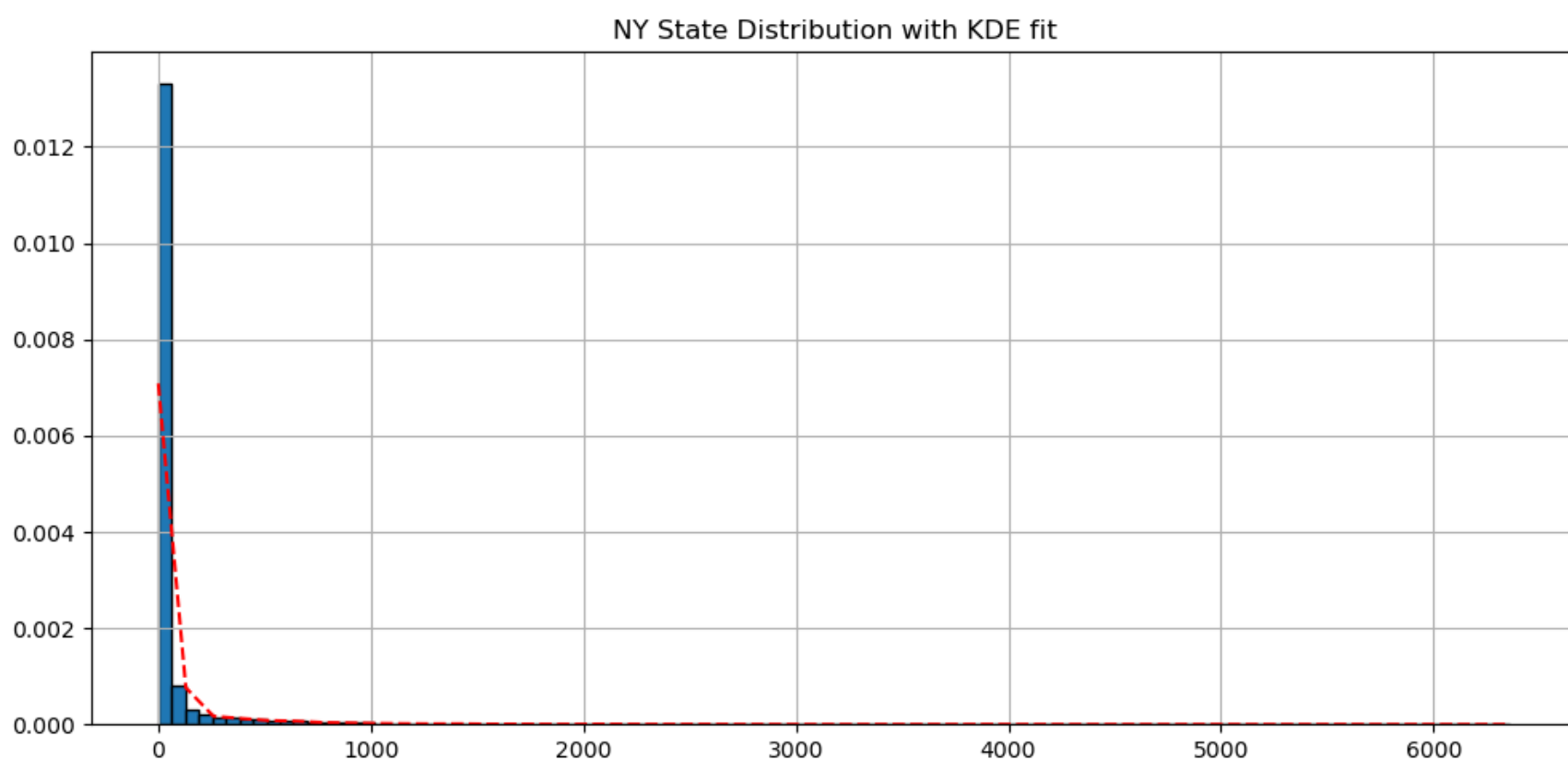
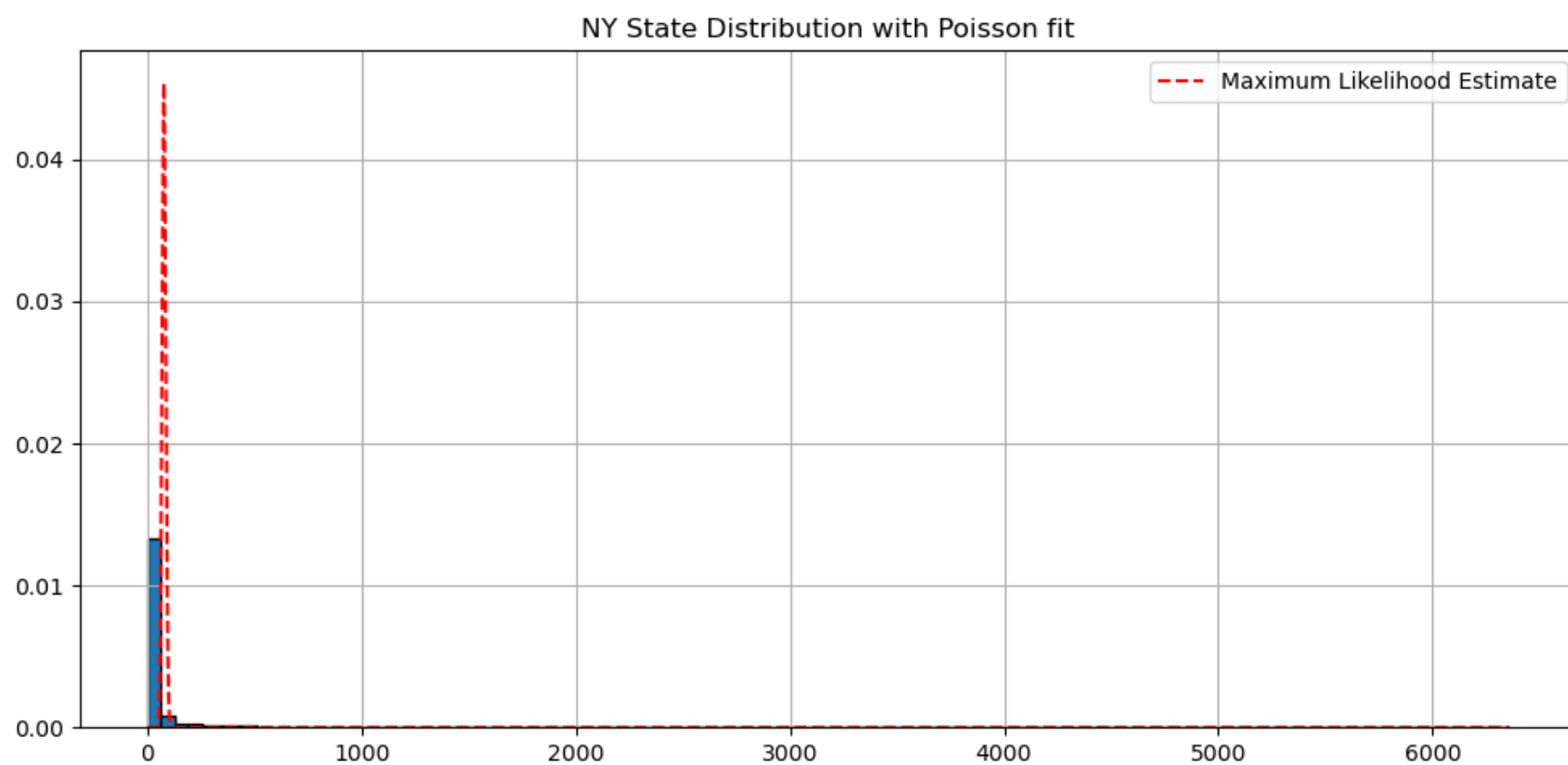
C:\Users\Dell\AppData\Local\Temp\ipykernel\_4452\3660962002.py:17: DeprecationWarning: Please use `gaussian\_kde` from the `scipy.stats` namespace, the `scipy.stats.kde` namespace is deprecated.  
density = kde.gaussian\_kde(new\_cases, bw\_method=None)





```
In [22]: 1 #Plotting Distribution for NY state
2 plot_state_distributions('NY')
```

C:\Users\Dell\AppData\Local\Temp\ipykernel\_4452\3660962002.py:17: DeprecationWarning: Please use `gaussian\_kde` from the `scipy.stats` namespace, the `scipy.stats.kde` namespace is deprecated.  
density = kde.gaussian\_kde(new\_cases, bw\_method=None)



From all the above plots we can see the KDE fits better for this data.

Let us calculate the statistics for the other states as well

```
In [23]: 1 def print_stats_for_state(state):
2     new_cases = Three_state_transformed_data.query(f"State=='{state}'").New_Cases
3     theta = new_cases.mean()
4     mean, var, skew, kurt = stats.poisson.stats(theta, moments='mvsk')
5     print(f"Stats for the state {state}")
6     print("=====")
7     print(f"mean : {mean}")
8     print(f"Variance : {var}")
9     print(f"Skewness : {skew}")
10    print(f"Kurtosis : {kurt}")
11    print("=====")
12
```

```
In [24]: 1 print_stats_for_state('CA')
2 print_stats_for_state('WA')
3 print_stats_for_state('NY')
4 print("Stats for the state NC")
5 print("=====")
6 print(f"mean : {NC_mean}")
7 print(f"Variance : {NC_var}")
8 print(f"Skewness : {NC_skew}")
9 print(f"Kurtosis : {NC_kurt}")
10 print("=====")
```

```
Stats for the state CA
=====
mean : 159.12649049307123
Variance : 159.12649049307123
Skewness : 0.07927363225942749
Kurtosis : 0.006284308771602944
=====
Stats for the state WA
=====
mean : 37.15660196501318
Variance : 37.15660196501318
Skewness : 0.16405217964334098
Kurtosis : 0.026913117645731018
=====
Stats for the state NY
=====
mean : 77.33283087126922
Variance : 77.33283087126922
Skewness : 0.11371507596100318
Kurtosis : 0.012931118500816722
=====
Stats for the state NC
=====
mean : 26.796728971962615
Variance : 26.796728971962615
Skewness : 0.19317864172521376
Kurtosis : 0.0373179876187985
=====
```

When comparing the statistics of all 4 states,

1. we see that Mean and variance of the distribution is high in California. Indicating more number of cases in CA.
2. Skewness and Kurtosis values are not very high in all of the states, even though the data appears to be skewed, this result shows that the skewness is less. This seems to be valid since the tail we are having is because of a couple of high numbered values. Which are outliers and most part of the data is close to the mean.

**Based on the results in Stage I, perform corelation between Enrichment data variables and COVID-19 cases to observe any patterns.**

You can compare either within your chosen specific state or among different states with the different enrich ment variables. Within the state you can compare the county based covid data to enrichment data for correlat ion. Between states you would need to aggregate to state level data and then perform correlation. Both covid and enrichment data will need to be normalized for population. For number of covid cases you can use a measu re of center value (median or mean) to compare the number of cases.

```
In [25]: 1 president_county = pd.read_csv('../../DATASETS/ENRICHMENT DATASETS/ELECTION_Datasets/president_county.csv')
2 president_county.head()
```

Out[25]:

	state	county	current_votes	total_votes	percent
0	Delaware	Kent County	87025	87025	100
1	Delaware	New Castle County	287633	287633	100
2	Delaware	Sussex County	129352	129352	100
3	District of Columbia	District of Columbia	41681	41681	100
4	District of Columbia	Ward 2	32881	32881	100

```
In [26]: 1 president_county_candidate = pd.read_csv('../../DATASETS/ENRICHMENT DATASETS/ELECTION_Datasets/president_c
2 president_county_candidate.head()
```

Out[26]:

	state	county	candidate	party	total_votes	won
0	Delaware	Kent County	Joe Biden	DEM	44552	True
1	Delaware	Kent County	Donald Trump	REP	41009	False
2	Delaware	Kent County	Jo Jorgensen	LIB	1044	False
3	Delaware	Kent County	Howie Hawkins	GRN	420	False
4	Delaware	New Castle County	Joe Biden	DEM	195034	True

In the above data we can see that in each county Democratic, Republic and other party candidates, Let us transform the data to create three separate columns for Dem\_total\_votes, Rep\_total\_votes and Oth\_total\_votes to find correlation with covid cases

```
In [27]: 1 transformed_president_df = president_county_candidate.copy()
2 transformed_president_df["Dem_total_votes"] = transformed_president_df.apply(lambda x: x.total_votes if x.par
3 transformed_president_df["Rep_total_votes"] = transformed_president_df.apply(lambda x: x.total_votes if x.par
4 transformed_president_df["Oth_total_votes"] = transformed_president_df.apply(lambda x: x.total_votes if x.par
5 transformed_president_df.head()
```

Out[27]:

	state	county	candidate	party	total_votes	won	Dem_total_votes	Rep_total_votes	Oth_total_votes
0	Delaware	Kent County	Joe Biden	DEM	44552	True	44552	0	0
1	Delaware	Kent County	Donald Trump	REP	41009	False	0	41009	0
2	Delaware	Kent County	Jo Jorgensen	LIB	1044	False	0	0	1044
3	Delaware	Kent County	Howie Hawkins	GRN	420	False	0	0	420
4	Delaware	New Castle County	Joe Biden	DEM	195034	True	195034	0	0

```
In [28]: 1 transformed_president_df=transformed_president_df.groupby(by=['state','county']).sum(numeric_only=True).reset
2 transformed_president_df.head()
```

Out[28]:

	state	county	total_votes	Dem_total_votes	Rep_total_votes	Oth_total_votes
0	Alabama	Autauga County	27770	7503	19838	429
1	Alabama	Baldwin County	109679	24578	83544	1557
2	Alabama	Barbour County	10518	4816	5622	80
3	Alabama	Bibb County	9595	1986	7525	84
4	Alabama	Blount County	27588	2640	24711	237

```
In [29]: 1 super_covid_data = pd.read_csv('../../DATASETS/SUPER DATASETS/superCovidDS.csv')
2 super_covid_data = super_covid_data.rename(columns={"County Name":"County_Name"})
3 super_covid_data.head()
```

Out[29]:

	countyFIPS	County_Name	State	StateFIPS	2020-01-22_x	2020-01-23_x	2020-01-24_x	2020-01-25_x	2020-01-26_x	2020-01-27_x	...	2023-01-08_y	2023-01-09_y	2023-01-10_y	2023-01-11_y	2023-01-12_y	2023-01-13_y	2023-01-14_y
0	1001	Autauga County	AL	1	0	0	0	0	0	0	...	230	230	230	230	230	230	230
1	1003	Baldwin County	AL	1	0	0	0	0	0	0	...	719	719	719	719	721	721	721
2	1005	Barbour County	AL	1	0	0	0	0	0	0	...	103	103	103	103	103	103	103
3	1007	Bibb County	AL	1	0	0	0	0	0	0	...	108	108	108	108	108	108	108
4	1009	Blount County	AL	1	0	0	0	0	0	0	...	260	260	260	260	261	261	261

5 rows × 2187 columns

```
In [30]: 1 print(super_covid_data.columns[4])
2 print(super_covid_data.columns[-2])
```

2020-01-22\_x  
2023-01-16\_y

```
In [31]: 1 # Creating transformed Covid data for a given state
2 def get_transformed_covid_data_for_state(state):
3     filtered_super_covid_data = super_covid_data.query(f"State=='{state}'")
4     transformed_super_covid_df = pd.DataFrame(columns=['Date', 'County_Name', 'State', 'population', 'Cases'])
5     start_date = dt.datetime(2020,1,22)
6     end_date = dt.datetime(2023,1,16)
7     date_series = pd.date_range(start_date, end_date, freq='d')
8     date_delta = dt.timedelta(days=1)
9     for date in date_series:
10         data = []
11         for _ , row in filtered_super_covid_data.iterrows():
12             temp = [date, getattr(row, 'County_Name').strip(), getattr(row, 'State'), getattr(row, 'population')]
13             cases_column = date.strftime('%Y-%m-%d_x')
14             temp.append(getattr(row, cases_column))
15             data.append(temp)
16         transformed_super_covid_df = pd.concat([transformed_super_covid_df, pd.DataFrame(data, columns=transformed_super_covid_df.columns)])
17     return transformed_super_covid_df
```

```
In [32]: 1 NC_transformed_covid = get_transformed_covid_data_for_state('NC')
2 NC_transformed_covid['State']='North Carolina'
3 display(NC_transformed_covid.head())
4 NC_transformed_covid.shape
```

	Date	County_Name	State	population	Cases
0	2020-01-22	Alamance County	North Carolina	169509	0
1	2020-01-22	Alexander County	North Carolina	37497	0
2	2020-01-22	Alleghany County	North Carolina	11137	0
3	2020-01-22	Anson County	North Carolina	24446	0
4	2020-01-22	Ashe County	North Carolina	27203	0

Out[32]: (109100, 5)

```
In [33]: 1 transformed_president_df.query("state=='North Carolina'")
```

Out[33]:

	state	county	total_votes	Dem_total_votes	Rep_total_votes	Oth_total_votes
3115	North Carolina	Alamance County	86091	38825	46056	1210
3116	North Carolina	Alexander County	20236	4145	15888	203
3117	North Carolina	Alleghany County	6076	1486	4527	63
3118	North Carolina	Anson County	11194	5789	5321	84
3119	North Carolina	Ashe County	15814	4164	11451	199
...	...	...	...	...	...	...
3210	North Carolina	Wayne County	55537	24215	30709	613
3211	North Carolina	Wilkes County	35466	7511	27592	363
3212	North Carolina	Wilson County	40735	20754	19581	400
3213	North Carolina	Yadkin County	19923	3763	15933	227
3214	North Carolina	Yancey County	11352	3688	7516	148

100 rows × 6 columns

```

In [34]: 1 NC_covid_president_data = pd.merge(NC_transformed_covid, transformed_president_df.query("state=='North Carolina'"),
2       how='left', left_on='County_Name', right_on='county')
3 NC_covid_president_data = NC_covid_president_data.drop(columns=['county', 'state'])
4 # Grouping by data to find the central value for the Cases
5 NC_covid_president_data = NC_covid_president_data.groupby(by=['County_Name', 'State', 'population']).aggregate(
6       'total_votes': max,
7       'Dem_total_votes': max,
8       'Rep_total_votes': max,
9       'Oth_total_votes': max
10      }).reset_index()
11 NC_covid_president_data['Cases'] = NC_covid_president_data['Cases'].apply(lambda x: int(x))
12 display(NC_covid_president_data.head())
13 NC_covid_president_data.shape

```

	County_Name	State	population	Cases	total_votes	Dem_total_votes	Rep_total_votes	Oth_total_votes
0	Alamance County	North Carolina	169509	25361	86091	38825	46056	1210
1	Alexander County	North Carolina	37497	5625	20236	4145	15888	203
2	Alleghany County	North Carolina	11137	1543	6076	1486	4527	63
3	Anson County	North Carolina	24446	3535	11194	5789	5321	84
4	Ashe County	North Carolina	27203	3210	15814	4164	11451	199

Out[34]: (100, 8)

```

In [35]: 1 NC_covid_president_data = pd.merge(NC_covid_president_data, president_county.query("state=='North Carolina'"),
2       how='left', left_on='County_Name', right_on='county')
3 NC_covid_president_data = NC_covid_president_data.rename(columns={'total_votes_y': 'total_votes'})
4 NC_covid_president_data = NC_covid_president_data.drop(columns=["total_votes_x", "county"])
5 NC_covid_president_data["Vote_Percent"] = 100*NC_covid_president_data["current_votes"]/NC_covid_president_data["total_votes"]
6 NC_covid_president_data.head()

```

Out[35]:

	County_Name	State	population	Cases	Dem_total_votes	Rep_total_votes	Oth_total_votes	current_votes	total_votes	Vote_Percent
0	Alamance County	North Carolina	169509	25361	38825	46056	1210	86091	86091	100.00000
1	Alexander County	North Carolina	37497	5625	4145	15888	203	20236	20236	100.00000
2	Alleghany County	North Carolina	11137	1543	1486	4527	63	6076	6076	100.00000
3	Anson County	North Carolina	24446	3535	5789	5321	84	11194	11205	99.90183
4	Ashe County	North Carolina	27203	3210	4164	11451	199	15814	15814	100.00000

```

In [36]: 1 #Normalizing the data with population for 1M people
2 NC_covid_president_data['Cases'] = (1000000 * NC_covid_president_data['Cases'])/NC_covid_president_data['population']
3 NC_covid_president_data['total_votes'] = (1000000 * NC_covid_president_data['total_votes'])/NC_covid_president_data['population']
4 NC_covid_president_data['current_votes'] = (1000000 * NC_covid_president_data['current_votes'])/NC_covid_president_data['population']
5 NC_covid_president_data['Dem_total_votes'] = (1000000 * NC_covid_president_data['Dem_total_votes'])/NC_covid_president_data['population']
6 NC_covid_president_data['Rep_total_votes'] = (1000000 * NC_covid_president_data['Rep_total_votes'])/NC_covid_president_data['population']
7 NC_covid_president_data['Oth_total_votes'] = (1000000 * NC_covid_president_data['Oth_total_votes'])/NC_covid_president_data['population']
8 NC_covid_president_data.head()

```

Out[36]:

	County_Name	State	population	Cases	Dem_total_votes	Rep_total_votes	Oth_total_votes	current_votes	total_votes	Vote_Percent
0	Alamance County	North Carolina	169509	149614	229043	271702	7138	507884	507884	100.00000
1	Alexander County	North Carolina	37497	150012	110542	423713	5413	539669	539669	100.00000
2	Alleghany County	North Carolina	11137	138547	133429	406482	5656	545568	545568	100.00000
3	Anson County	North Carolina	24446	144604	236807	217663	3436	457907	458357	99.90183
4	Ashe County	North Carolina	27203	118001	153071	420946	7315	581332	581332	100.00000

In [39]:

1

NC\_covid\_president\_data.corr()

Out[39]:

	population	Cases	Dem_total_votes	Rep_total_votes	Oth_total_votes	current_votes	total_votes	Vote_Percent
population	1.000000	0.035008	0.381556	-0.346206	0.363087	-0.015851	-0.016276	0.032280
Cases	0.035008	1.000000	-0.333223	-0.039499	-0.491957	-0.486522	-0.487556	-0.025898
Dem_total_votes	0.381556	-0.333223	1.000000	-0.741393	0.147517	0.167597	0.168128	-0.005377
Rep_total_votes	-0.346206	-0.039499	-0.741393	1.000000	0.174355	0.536930	0.536394	0.171888
Oth_total_votes	0.363087	-0.491957	0.147517	0.174355	1.000000	0.471662	0.473106	-0.012776
current_votes	-0.015851	-0.486522	0.167597	0.536930	0.471662	1.000000	0.999936	0.241139
total_votes	-0.016276	-0.487556	0.168128	0.536394	0.473106	0.999936	1.000000	0.230120
Vote_Percent	0.032280	-0.025898	-0.005377	0.171888	-0.012776	0.241139	0.230120	1.000000

From the above correlation matrix, we can identify the below correlation values

- Cases - Vote\_Percent : -0.03 approximately
- Cases - Dem\_total\_votes : -0.33 approximately
- Cases - Rep\_total\_votes : -0.04 approximately
- Cases - Oth\_total\_votes : -0.49 approximately

Hypothesis

1. Here **Vote\_Percent** is the percentage of people came to vote of the total people with vote eligibility. The vote percent is negligibly negatively correlated with covid cases. **So the data shows that there is not a significant difference in the change in vote percent with the number of covid cases**
2. Other parties total votes is moderately negatively correlated with Covid cases. **So the data shows that in the regions with high covid cases, other parites received moderately less votes when compared to the regions with less covid cases**
3. Republic party total votes is negligibly negatively correlated with covid cases. **So the data shows that there is not a significant difference in the change of total votes received by republic party in regions with high covid cases when compared to the regions with fewer covid cases**
4. Democratic party total votes is moderately negatively correlated with covid cases. **So the data shows that in the regions with high covid cases, democratic party received moderately less votes when compared to the regions with fewer covid cases**

In [ ]:

1