

```

1 #data-set 8.csv
2 # https://github.com/michaelofsbu/CSE-544-Datasets
3
4
5 #imports
6 from google.colab import drive
7 import numpy as numpy
8 import numpy as np
9 import pandas as pd
10 import seaborn
11 import matplotlib.pyplot as plt
12 from scipy.stats import poisson
13 from scipy.stats import geom
14 from scipy.stats import binom
15 from numpy.linalg import inv
16 from numpy import dtype
17 from scipy.stats import gamma
18 import itertools
19 import math
20 plt.style.use('ggplot')

1 drive.mount('/content/gdrive')

```

Drive already mounted at /content/gdrive; to attempt to forcibly remount, call d

## ▼ Basic Info about Datafile 8.csv

```

1 df = pd.read_csv('gdrive/MyDrive/Colab Notebooks/8.csv')
2
3 df_cumm = df
4
5
6
7
8
9 print(df.head())
10 print()
11 print("Data start date: "+ str(df['Date'].min()))
12 print("Data start date: "+ str(df['Date'].max()))
13 print("Total days in the data: "+ str(len(df)))
14
15 totalrows = df.shape[0]
16

```

	Date	IL confirmed	IN confirmed	IL deaths	IN deaths
0	2020-01-22	0	0	0	0
1	2020-01-23	0	0	0	0

2	2020-01-24	1	0	0	0
3	2020-01-25	1	0	0	0
4	2020-01-26	1	0	0	0

Data start date: 2020-01-22

Data start date: 2021-04-03

Total days in the data: 438

## ▼ PART1 : Data Pre-Processing

### Let's look at various statistics of the data

Step1 : Converting cumulative datapoints to daily datapo

Step2 : Searching for missing values as a part of cleaning data

Step3 : Removing outliers if any

## ▼ Converting cummulative data points to daily data points

```

1 df1 = df.iloc[:,0]
2 #print(df1)
3 df.drop(df.columns[[0]], axis = 1, inplace = True)
4 #print(df)
5 for key, value in df.iteritems():
6     po = value
7     #print(type(po))
8     ope = []
9     for j in range(len(po)):
10         if(j>0):
11             ope.append(po[j]-po[j-1])
12         else:
13             ope.append(po[0])
14     df[key] = ope
15 df["Date"] = df1.values
16
17
18 features = []
19
20 for i in range(len(df.columns)-1):
21     features.append(df.columns[i])
22
23 print("Dataset features : " + str(features))

```

Dataset features : ['IL confirmed', 'IN confirmed', 'IL deaths', 'IN deaths']

## ▼ Removing missing values

```

1 data = pd.DataFrame([])
2 for feature in features:
3     tempDict = {}
4     tempDict['Feature'] = feature
5     tempDict['Mean'] = df[feature].mean()
6     tempDict['Std Dev'] = df[feature].std()
7     tempDict['Missing values count'] = totalrows - len(df[df[feature].notna()])
8     tempDict['Missing values %'] = tempDict['Missing values count']
9     data = data.append(tempDict, ignore_index=True)
10
11 print(data)
12
13
14

```

	Feature	Mean	...	Missing values count	Std Dev
0	IL confirmed	2863.442922	...	0.0	3183.408897
1	IN confirmed	1583.433790	...	0.0	1980.552477
2	IL deaths	48.812785	...	0.0	48.251834
3	IN deaths	29.858447	...	0.0	78.234074

[4 rows x 5 columns]

There are no missing values in any of the columns. So we need not care much about cleaning the data. Let's check for outliers in the next step using Tukey's rule.

## ▼ Tukey's rule for detecting outliers

```

1 #outliers
2 outlier_index = set()
3 for feature in features:
4     Q1 = np.percentile(df[feature], 25)
5     Q3 = np.percentile(df[feature], 75)
6     IQR = Q3 - Q1
7
8     outlier_ = df[(df[feature] < Q1 - IQR*1.5) | (df[feature] > Q3 + 1.5*IQR)].index
9     outlier_index = outlier_index.union(outlier_)
10
11 print(outlier_index)
12
13 print(str(len(outlier_index))+" outliers are found\n")
14 df = df.drop(outlier_index, axis=0)
15
16

```

```
17 print("After dropping the outliers, we have "+str(df.shape[0])+" datapoints")

{282, 284, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300,
76 outliers are found

After dropping the outliers, we have 362 datapoints
```

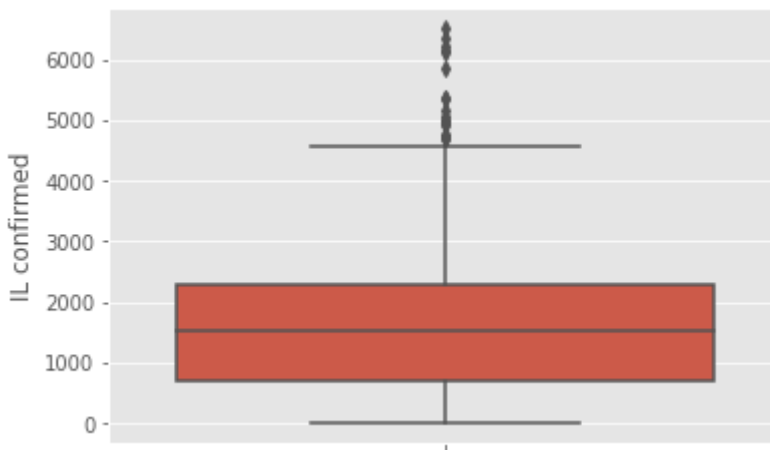
I have used Tukey's rule on each column individually. Even if one attribute classify it as outlier, the whole data point is considered as outlier. This I've attained using a set() datastructure and update it whenever I detect a outlier.

Finally as you can see, the set() is notempty with 76 index values. So we can conclude there are 76 outliers. I'm removing all the outliers.

## ▼ Data Analysis

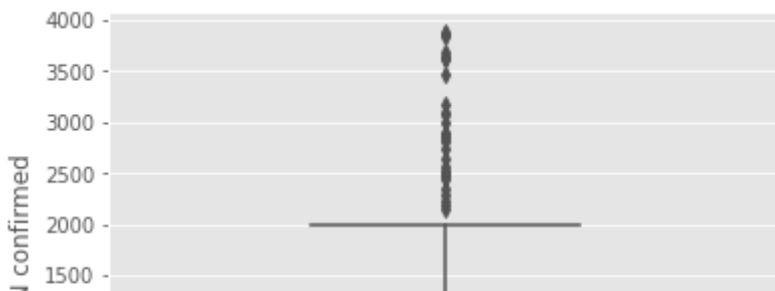
### ▼ Box plot of 'IL confirmed'

```
1 fig1 = seaborn.boxplot(y=df[features[0]])
```



### ▼ Box plot of 'IN confirmed'

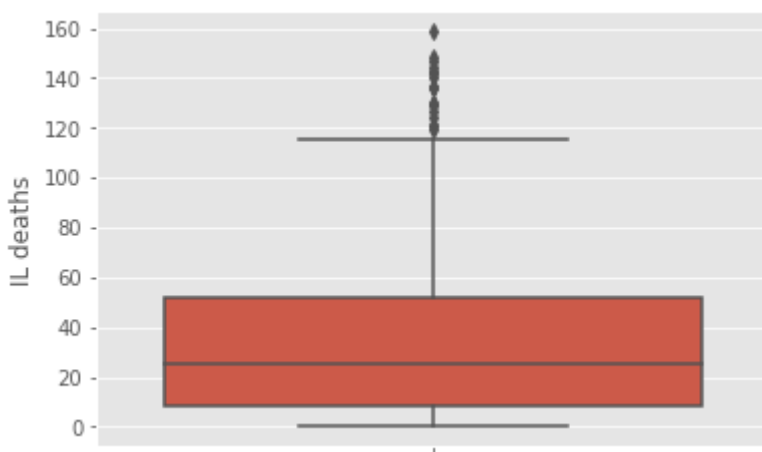
```
1 fig2 = seaborn.boxplot(y=df[features[1]])
```



### ▼ Box plot of 'IL deaths'

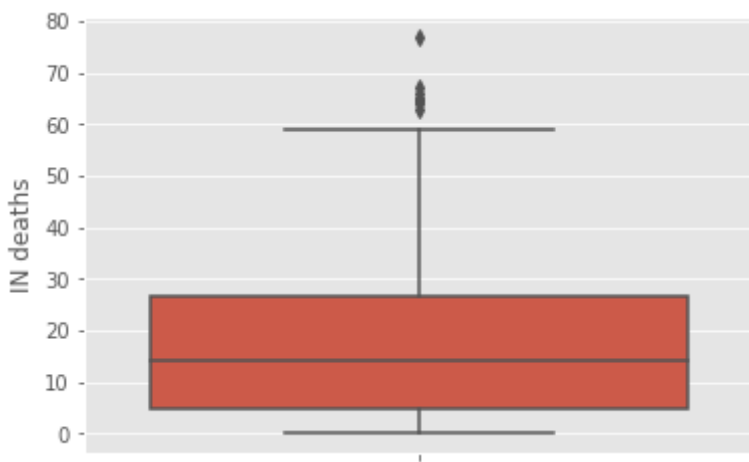


```
1 fig3 = seaborn.boxplot(y=df[features[2]])
```



### ▼ Box plot of 'IN deaths'

```
1 fig4 = seaborn.boxplot(y=df[features[3]])
```

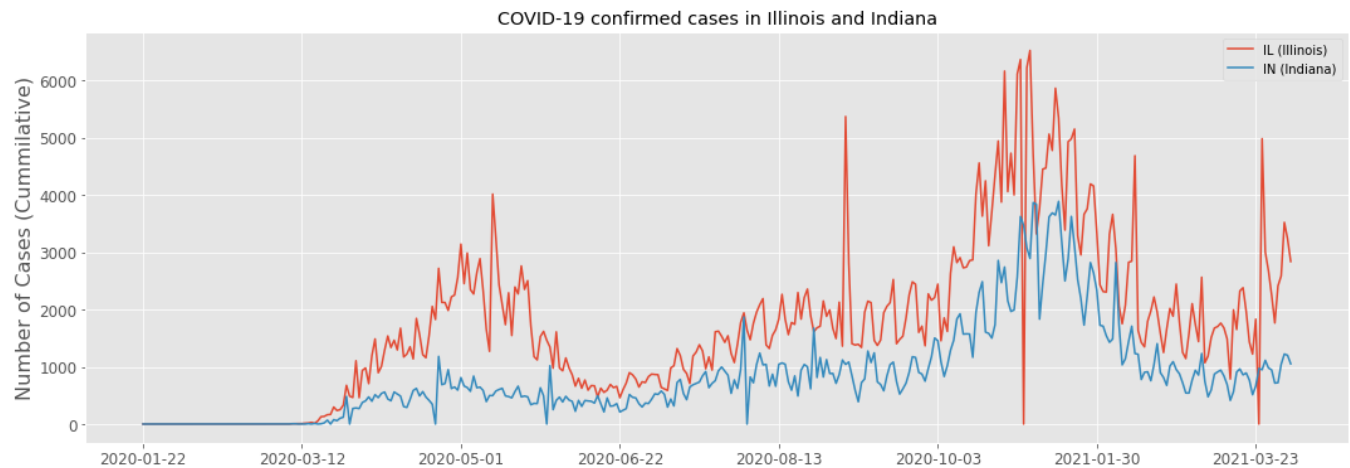


### ▼ Comparison of cases in IL and IN states

```

1 fig, ax = plt.subplots(figsize=(18,6))
2 plt.plot(df['Date'], df['IL confirmed'], label="IL (Illinois)")
3 plt.plot(df['Date'], df['IN confirmed'], label="IN (Indiana)")
4 plt.legend(loc='best')
5 plt.title('COVID-19 confirmed cases in Illinois and Indiana')
6 plt.ylabel('Number of Cases (Cumulative)', fontsize = 16)
7 plt.xticks(np.arange(0, len(df), 50))
8 plt.xticks(fontsize = 12)
9 plt.yticks(fontsize = 12)
10 plt.show()
11
12

```

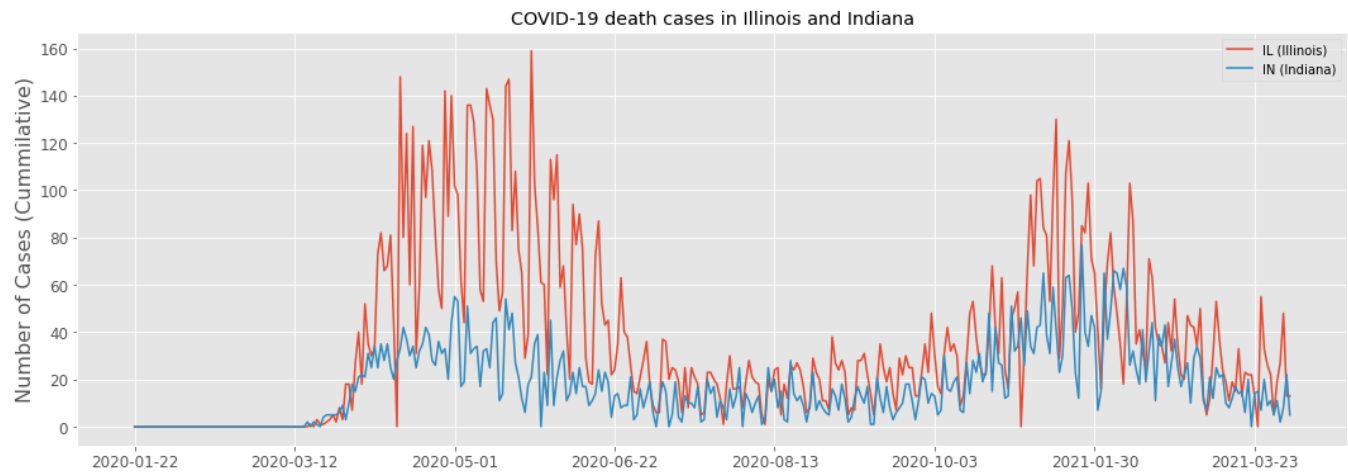


## ▼ Comparison of deaths in IL and IN states

```

1 fig, ax = plt.subplots(figsize=(18,6))
2 plt.plot(df['Date'], df['IL deaths'], label="IL (Illinois)")
3 plt.plot(df['Date'], df['IN deaths'], label="IN (Indiana)")
4 plt.legend(loc='best')
5 plt.title('COVID-19 death cases in Illinois and Indiana')
6 plt.ylabel('Number of Cases (Cumulative)', fontsize = 16)
7 plt.xticks(np.arange(0, len(df), 50))
8 plt.xticks(fontsize = 12)
9 plt.yticks(fontsize = 12)
10 plt.show()

```



## ▼ PART 2a : Prediction in August '20 month

### Error Functions for reporting the accuracy of my predictions

```

1 #Error Functions MAPE and MSE
2 def MSE(Y_act, Y_pred):
3     SSE = 0
4     Y_sse = Y_act-Y_pred
5     for i in range(0,len(Y_pred)):
6         SSE = SSE + Y_sse[i]*Y_sse[i]
7     return SSE/len(Y_act)
8
9 def MAPE(Y_act,Y_pred):
10    err = (abs(Y_act-Y_pred))/Y_act
11    MAPE = 0
12    for i in range(0, len(Y_act)):
13        MAPE = MAPE + err[i]
14    return 100*MAPE/len(Y_act)

1
2 #Beta function and Auto Regression
3 def Beta(X,Y):
4     return np.matmul(np.matmul(inv(np.matmul(np.transpose(X),X)),np.transpose(X)),Y)
5
6 def AR(aug_month, p, col):
7     aug_3, aug_1 = aug_month[0:-7], aug_month[-7:0]
8     aug_3 = aug_3.to_numpy()
9     aug_2, aug_1 = aug_3[:p], aug_3[p:]

```

```

9     aug_3_cases = aug_3[:,col]
10    X_train = []
11    Y_train = []
12    for i in range(len(aug_3_cases)-p):
13        X_train = X_train + [aug_3_cases[i:i+p]]
14        Y_train = Y_train + [aug_3_cases[i+p]]
15    X_train, Y_train = np.array(X_train,dtype=float),np.array(Y_train,dtype=float)
16    one_app = np.ones((len(X_train),), dtype=float)
17    one_app = np.reshape(one_app,(len(one_app),1))
18    X_train = np.append(one_app , X_train, axis = 1)
19    beta = Beta(X_train,Y_train)
20    aug_np = aug_month.to_numpy()
21    aug_cases = aug_np[:,col]
22    x = []
23    Y_pred = []
24    Y_act = []
25    for i in range(len(aug_cases)-7,len(aug_cases)):
26        x = [aug_cases[i-p:i]]
27        x = np.append([1],x)
28        y = np.matmul(x,beta)
29        Y_pred = Y_pred + [y]
30        Y_act = Y_act + [aug_cases[i]]
31    Y_pred, Y_act = np.array(Y_pred),np.array(Y_act)
32    Y_pred = np.reshape(Y_pred,(len(Y_act),))
33    return Y_pred,Y_act
34
35 def EWMA(aug_month, alpha, col):
36     aug_np = aug_month.to_numpy()
37     aug_cases = aug_np[:,col]
38     Y_pred = [aug_cases[0]]
39     Y_pred = np.array(Y_pred,dtype=float)
40     Y_act = []
41     for i in range(1,len(aug_cases)):
42         Y_pred = np.append(Y_pred , [(alpha*aug_cases[i-1] +(1-alpha)*Y_pred[i-1])])
43     Y_act = aug_cases
44     Y_pred,Y_act = np.array(Y_pred[-7:]),np.array(Y_act[-7:])
45     Y_pred = np.reshape(Y_pred,(len(Y_act),))
46     return Y_pred,Y_act

```

```

1 aug_month = df[(df['Date'] >= '2020-08-01') & (df['Date'] <= '2020-08-28')]
2 print("----- TIME SERIES ANALYSIS FOR THE ILLIONIS STATE----- ")
3 col = [0,2,1,3]
4 for c in col:
5     print( "Time Series Analysis using : AR ")
6     ar = [3, 5]
7     for a in ar:
8         Y_pred, Y_act = AR(aug_month,a,c)
9         if( c == 1 or c ==2):
10             print( "    MSE using AR(",a,") for # of COVID cases : ",MSE(Y_act, Y_p
11             print( "    MAPE using AR(",a,") for # of COVID cases : ",MAPE(Y_act, Y_
12         else:
13             print( "    MSE using AR(",a,") for # of COVID fatalities : ",MSE(Y_act,

```



```

13         print( "      MSE using AR(",a," ) for # of COVID fatalities : ",MSE(Y_act,
14         print( "      MAPE using AR(",a," ) for # of COVID fatalities : ",MAPE(Y_act,
15     print( " Using EWMA time series analysis ")
16     ar = [0.5, 0.8]
17     for a in ar:
18         Y_pred, Y_act = EWMA(aug_month,a,c)
19         if( c == 1 or c ==2):
20             print( "      MSE using EWMA(",a," ) for # of COVID cases : ",MSE(Y_act, Y_
21             print( "      MAPE using EWMA(",a," ) for # of COVID cases : ",MAPE(Y_act,
22         else:
23             print( "      MSE using EWMA(",a," ) for # of COVID fatalities : ",MSE(Y_act,
24             print( "      MAPE using EWMA(",a," ) for # of COVID fatalities : ",MAPE(Y_act,
25     if(c == 3):
26         print("----- TIME SERIES ANALYSIS FOR THE INDIANA STATE-----")
27
28 #col = [0,2,1,3]
29

```

----- TIME SERIES ANALYSIS FOR THE ILLIONIS STATE-----

Time Series Analysis using : AR

```

MSE using AR( 3 ) for # of COVID fatalities : 122872.72205712622
MAPE using AR( 3 ) for # of COVID fatalities : 12.720051390973742
MSE using AR( 5 ) for # of COVID fatalities : 164291.9795117699
MAPE using AR( 5 ) for # of COVID fatalities : 16.740613726507025

```

Using EWMA time series analysis

```

MSE using EWMA( 0.5 ) for # of COVID fatalities : 93981.0120588206
MAPE using EWMA( 0.5 ) for # of COVID fatalities : 14.35864988607015
MSE using EWMA( 0.8 ) for # of COVID fatalities : 92653.51727834626
MAPE using EWMA( 0.8 ) for # of COVID fatalities : 12.893428483645723

```

Time Series Analysis using : AR

```

MSE using AR( 3 ) for # of COVID cases : 50.61458137982926
MAPE using AR( 3 ) for # of COVID cases : 41.931729930360724
MSE using AR( 5 ) for # of COVID cases : 66.20707908762164
MAPE using AR( 5 ) for # of COVID cases : 33.69360565652969

```

Using EWMA time series analysis

```

MSE using EWMA( 0.5 ) for # of COVID cases : 91.89441084356743
MAPE using EWMA( 0.5 ) for # of COVID cases : 62.68739372822626
MSE using EWMA( 0.8 ) for # of COVID cases : 95.6192281461625
MAPE using EWMA( 0.8 ) for # of COVID cases : 52.16492605323872

```

Time Series Analysis using : AR

```

MSE using AR( 3 ) for # of COVID cases : 147223.00351167255
MAPE using AR( 3 ) for # of COVID cases : 25.505557020173747
MSE using AR( 5 ) for # of COVID cases : 172647.89981848048
MAPE using AR( 5 ) for # of COVID cases : 30.957691736776912

```

Using EWMA time series analysis

```

MSE using EWMA( 0.5 ) for # of COVID cases : 177565.48419546496
MAPE using EWMA( 0.5 ) for # of COVID cases : 33.97358098284393
MSE using EWMA( 0.8 ) for # of COVID cases : 246036.86266481175
MAPE using EWMA( 0.8 ) for # of COVID cases : 39.380694490480856

```

Time Series Analysis using : AR

```

MSE using AR( 3 ) for # of COVID fatalities : 31.21891969688699
MAPE using AR( 3 ) for # of COVID fatalities : 97.93763823096482
MSE using AR( 5 ) for # of COVID fatalities : 19.945044550772376
MAPE using AR( 5 ) for # of COVID fatalities : 63.32794190309236

```

Using EWMA time series analysis

```

MSE using EWMA( 0.5 ) for # of COVID fatalities : 58.14365420638482

```

MAPE using EWMA( 0.5 ) for # of COVID fatalities : 101.94576585691568

MSE using EWMA( 0.8 ) for # of COVID fatalities : 71.93430252540321

MAPE using EWMA( 0.8 ) for # of COVID fatalities : 108.91794860666978

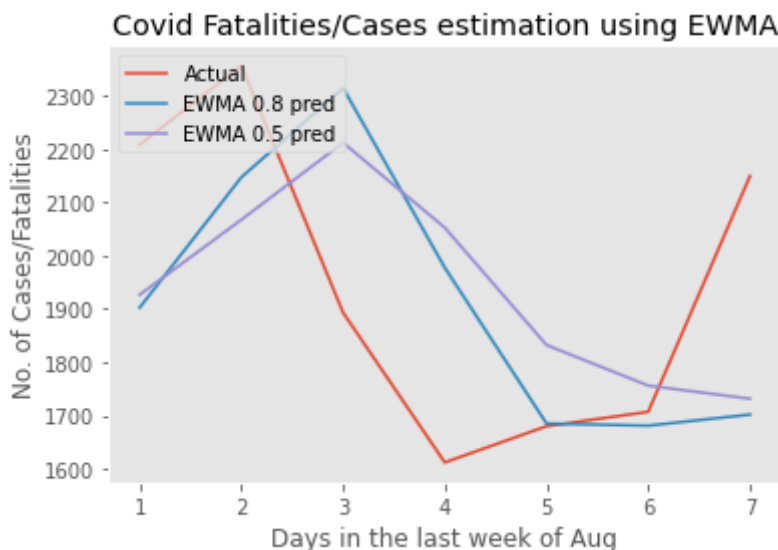
----- TIME SERIES ANALYSIS FOR THE INDIANA STATE-----

## ▼ EWMA with aplha=0.5,0.8 for IL Confirmed

```

1
2 columns = [1,3,0,2]
3 col = 0
4
5 Y_pred_8,Y_act = EWMA(aug_month,0.8,col)
6
7
8 Y_pred_5,Y_act = EWMA(aug_month,0.5,col)
9
10 X = [i for i in range(1,8)]
11 plt.figure('EWMA')
12 plt.plot(X, Y_act ,label='Actual')
13 plt.plot(X, Y_pred_8 ,label='EWMA 0.8 pred')
14 plt.plot(X, Y_pred_5 ,label='EWMA 0.5 pred')
15 plt.xlabel('Days in the last week of Aug')
16 plt.ylabel('No. of Cases/Fatalities')
17 plt.title('Covid Fatalities/Cases estimation using EWMA')
18 plt.legend(loc="upper left")
19 plt.grid()
20 plt.show()
21
22
23

```

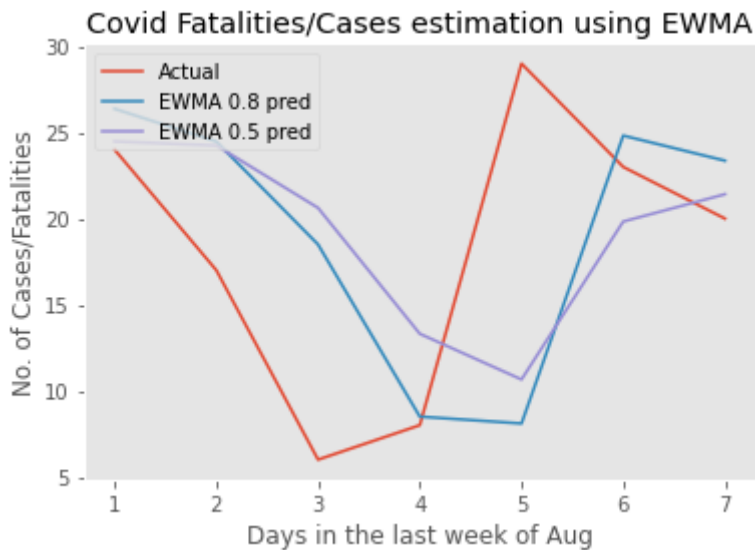


## ▼ EWMA with alpha=0.5,0.8 for IL Deaths

```

1 col = 2
2 Y_pred_8,Y_act = EWMA(aug_month,0.8,col)
3 Y_pred_5,Y_act = EWMA(aug_month,0.5,col)
4
5 X = [i for i in range(1,8)]
6 plt.figure('EWMA')
7 plt.plot(X, Y_act ,label='Actual')
8 plt.plot(X, Y_pred_8 ,label='EWMA 0.8 pred')
9 plt.plot(X, Y_pred_5 ,label='EWMA 0.5 pred')
10 plt.xlabel('Days in the last week of Aug')
11 plt.ylabel('No. of Cases/Fatalities')
12 plt.title('Covid Fatalities/Cases estimation using EWMA')
13 plt.legend(loc="upper left")
14 plt.grid()
15 plt.show()

```



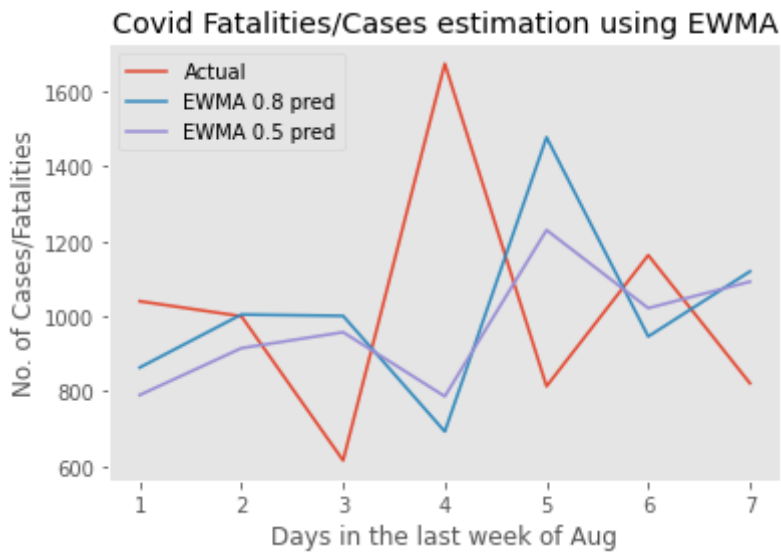
## ▼ EWMA with alpha=0.5,0.8 for IN Confirmed

```

1 col = 1
2 Y_pred_8,Y_act = EWMA(aug_month,0.8,col)
3 Y_pred_5,Y_act = EWMA(aug_month,0.5,col)
4
5 X = [i for i in range(1,8)]
6 plt.figure('EWMA')
7 plt.plot(X, Y_act ,label='Actual')
8 plt.plot(X, Y_pred_8 ,label='EWMA 0.8 pred')
9 plt.plot(X, Y_pred_5 ,label='EWMA 0.5 pred')
10 plt.xlabel('Days in the last week of Aug')
11 plt.ylabel('No. of Cases/Fatalities')
12 plt.title('Covid Fatalities/Cases estimation using EWMA')
13 plt.legend(loc="upper left")

```

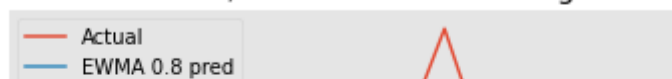
```
14 plt.grid()
15 plt.show()
```



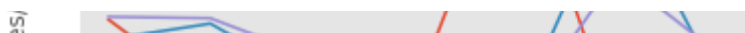
## ▼ EWMA with alpha=0.5,0.8 for IN Deaths

```
1 col = 3
2 Y_pred_8,Y_act = EWMA(aug_month,0.8,col)
3 Y_pred_5,Y_act = EWMA(aug_month,0.5,col)
4
5 X = [i for i in range(1,8)]
6 plt.figure('EWMA')
7 plt.plot(X, Y_act ,label='Actual')
8 plt.plot(X, Y_pred_8 ,label='EWMA 0.8 pred')
9 plt.plot(X, Y_pred_5 ,label='EWMA 0.5 pred')
10 plt.xlabel('Days in the last week of Aug')
11 plt.ylabel('No. of Cases/Fatalities')
12 plt.title('Covid Fatalities/Cases estimation using EWMA')
13 plt.legend(loc="upper left")
14 plt.grid()
15 plt.show()
```

## Covid Fatalities/Cases estimation using EWMA



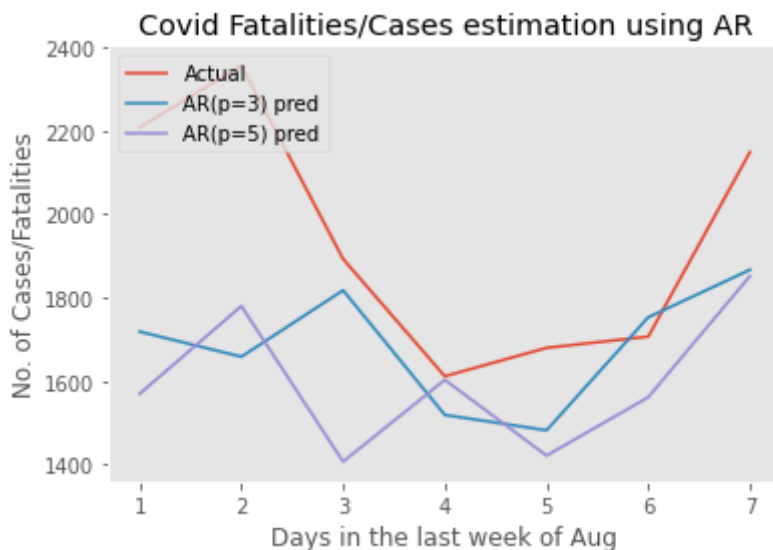
## Auto Regression with p=3,5 for IL Confirmed



```

1 col = 0
2 Y_pred_ar5,Y_act = AR(aug_month,5,col)
3 Y_pred_ar3,Y_act = AR(aug_month,3,col)
4 plt.figure('AR')
5 plt.plot(X, Y_act ,label='Actual')
6 plt.plot(X, Y_pred_ar3 ,label='AR(p=3) pred')
7 plt.plot(X, Y_pred_ar5 ,label='AR(p=5) pred')
8 plt.xlabel('Days in the last week of Aug')
9 plt.ylabel('No. of Cases/Fatalities')
10 plt.title('Covid Fatalities/Cases estimation using AR')
11 plt.legend(loc="upper left")
12 plt.grid()
13 plt.show()

```



## Auto Regression with p=3,5 for IL Deaths

```

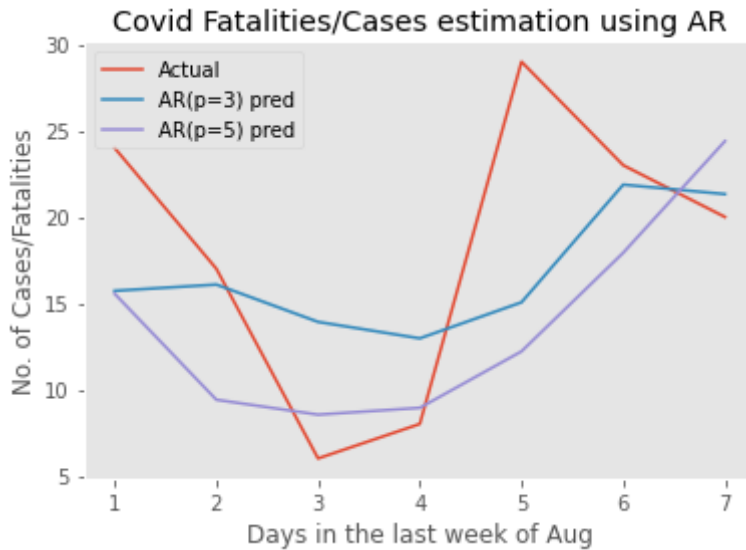
1 col = 2
2 Y_pred_ar5,Y_act = AR(aug_month,5,col)
3 Y_pred_ar3,Y_act = AR(aug_month,3,col)
4 plt.figure('AR')
5 plt.plot(X, Y_act ,label='Actual')
6 plt.plot(X, Y_pred_ar3 ,label='AR(p=3) pred')
7 plt.plot(X, Y_pred_ar5 ,label='AR(p=5) pred')
8 plt.xlabel('Days in the last week of Aug')
9 plt.ylabel('No. of Cases/Fatalities')

```

```

10 plt.title('Covid Fatalities/Cases estimation using AR')
11 plt.legend(loc="upper left")
12 plt.grid()
13 plt.show()

```



## ▼ Auto Regression with p=3,5 for IN Confirmed

```

1 col = 1
2 Y_pred_ar5,Y_act = AR(aug_month,5,col)
3 Y_pred_ar3,Y_act = AR(aug_month,3,col)
4 plt.figure('AR')
5 plt.plot(X, Y_act ,label='Actual')
6 plt.plot(X, Y_pred_ar3 ,label='AR(p=3) pred')
7 plt.plot(X, Y_pred_ar5 ,label='AR(p=5) pred')
8 plt.xlabel('Days in the last week of Aug')
9 plt.ylabel('No. of Cases/Fatalities')
10 plt.title('Covid Fatalities/Cases estimation using AR')
11 plt.legend(loc="upper left")
12 plt.grid()
13 plt.show()

```



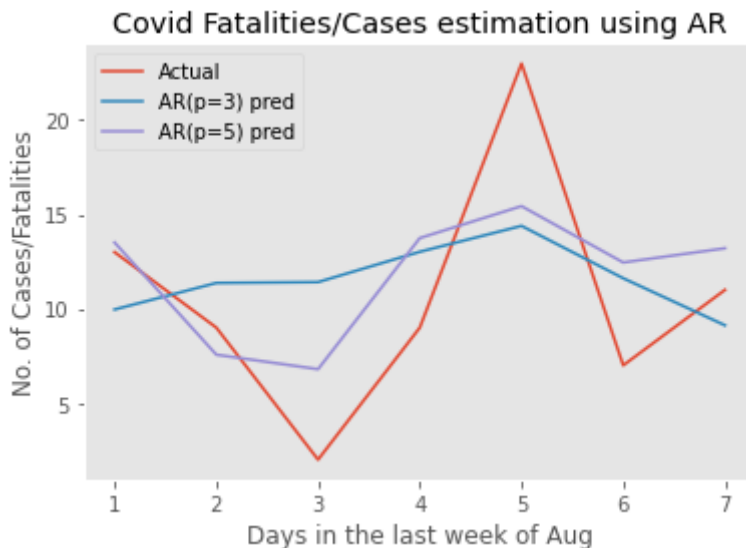
## ▼ Auto Regression with p=3,5 for IN Deaths



```

1 col = 3
2 Y_pred_ar5,Y_act = AR(aug_month,5,col)
3 Y_pred_ar3,Y_act = AR(aug_month,3,col)
4 plt.figure('AR')
5 plt.plot(X, Y_act ,label='Actual')
6 plt.plot(X, Y_pred_ar3 ,label='AR(p=3) pred')
7 plt.plot(X, Y_pred_ar5 ,label='AR(p=5) pred')
8 plt.xlabel('Days in the last week of Aug')
9 plt.ylabel('No. of Cases/Fatalities')
10 plt.title('Covid Fatalities/Cases estimation using AR')
11 plt.legend(loc="upper left")
12 plt.grid()
13 plt.show()

```



```

1 df_out = pd.DataFrame()
2 df_out['Date'] = aug_month['Date'][-7:]
3 df_out['Actual_Cases'] = Y_act
4 df_out['Predicted_cases_AR(3)'] = Y_pred_ar3
5 df_out['Predicted_cases_AR(5)'] = Y_pred_ar5
6 df_ar = df_out.reset_index()
7 del df_ar['index']
8 print('PREDICTIONS USING AR(3) and AR(5)')
9 print(df_ar)

```

PREDICTIONS USING AR(3) and AR(5)

	Date	Actual_Cases	Predicted_cases_AR(3)	Predicted_cases_AR(5)
0	2020-08-21	13	9.962786	13.510117

1	2020-08-22	9	11.369561	7.564766
2	2020-08-23	2	11.414789	6.796103
3	2020-08-24	9	13.021625	13.748345
4	2020-08-25	23	14.396010	15.442694
5	2020-08-27	7	11.618570	12.455057
6	2020-08-28	11	9.123464	13.208006

## ▼ PART 2b : Stats between Feb and Mar '21

```

1 feb_month = df[(df['Date'] >= '2021-02-01') & (df['Date'] <= '2021-02-28')]
2 mar_month = df[(df['Date'] >= '2021-03-01') & (df['Date'] <= '2021-03-31')]
3
4
5 def variance(list_feature):
6     sum_of_squared = 0
7     mean = list_feature.mean()
8     n = len(list_feature)
9     for i in list_feature:
10         sum_of_squared = sum_of_squared + (i - mean)*(i - mean)
11     return sum_of_squared/(n-1)
12

```

## ▼ Walds 1 sample testing for mean of Mar '21

```

1 def walds_test_single(feature, sample_mean_feb):
2     lambda_mar = mar_month[feature].mean()
3     W = (lambda_mar - sample_mean_feb)/(np.sqrt(lambda_mar/len(mar_month)))
4
5     W = abs(W)
6
7     if(W > 1.96):
8         print("We are rejecting NULL HYPOTHESIS for '" + feature + "' because absol
9     else:
10         print("We are Accepting NULL HYPOTHESIS for '" + feature + "' because absol
11
12     # print(W)
13     print('\n\n')
14
15
16
17 for feature in features:
18     sample_mean_feb = feb_month[feature].mean()
19
20     # print(mar_month[feature].mean())
21     walds_test_single(feature, sample_mean_feb)
22

```



We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of Wald

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of Wald

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of Walds

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of Walds

## Result of Walds 1 sample testing for mean of cases and death

### Null hypothesis (H0):

Mean of cases/deaths in Mar '21 = Mean of cases/deaths in Feb '21.

### Alternate hypothesis(H1):

Mean of cases/deaths in Mar '21 not equals Mean of cases/deaths in Feb '21.

### Result:

We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of Walds 1 sample testing(W) is 268.9736251180282 and it's greater than Z value(at 0.025) = 1.96

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of Walds 1 sample testing(W) is 179.87439845233195 and it's greater than Z value(at 0.025) = 1.96

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of Walds 1 sample testing(W) is 38.92535543737809 and it's greater than Z value(at 0.025) = 1.96

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of Walds 1 sample testing(W) is 44.35531245905998 and it's greater than Z value(at 0.025) = 1.96

### Is the test applicable?

Since we used MLE as estimator, it is Asymptotically Normal as n tends to infinity (n>=30 is satisfied). This is the only condition, hence the test is applicable.

## ▼ Z 1 sample testing for mean of Mar '21

1

```
2 def z_test_single(feature, sample_mean_feb, true_variance):
```

```

3     lambda_mar = mar_month[feature].mean()
4     Z = (lambda_mar - sample_mean_feb)/np.sqrt(true_variance/len(mar_month))
5
6     Z = abs(Z)
7
8     if(Z > 1.96):
9         print("We are rejecting NULL HYPOTHESIS for '" + feature + "' because absol
10    else:
11        print("We are Accepting NULL HYPOTHESIS for '" + feature + "' because absol
12
13    # print(Z)
14    print('\n\n')
15
16
17
18 for feature in features:
19     sample_mean_feb = feb_month[feature].mean()
20     true_variance = variance(df[feature].values)
21     # print(np.sqrt(true_variance))
22     Z_test_single(feature, sample_mean_feb, true_variance)
23
24
25
26
27
28

```

We are Accepting NULL HYPOTHESIS for 'IL confirmed' because absolute value of Z is

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of Z is

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of Z is

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of Z is

## Result of Z 1 sample testing for mean of cases and death

### Null hypothesis (H0):

Mean of cases/deaths in Mar '21 = Mean of cases/deaths in Feb '21.

### Alternate hypothesis(H1):

Mean of cases/deaths in Mar '21 not equals Mean of cases/deaths in Feb '21.

**Result:**

We are Accepting NULL HYPOTHESIS for 'IL confirmed' because absolute value of Z 1 sample testing(Z) is 0.6754221197921421 and it's greater than Z value(at 0.025) = 1.96

We are Accepting NULL HYPOTHESIS for 'IN confirmed' because absolute value of Z 1 sample testing(Z) is 0.5955376536202984 and it's greater than Z value(at 0.025) = 1.96

We are Accepting NULL HYPOTHESIS for 'IL deaths' because absolute value of Z 1 sample testing(Z) is 0.8191724665851586 and it's greater than Z value(at 0.025) = 1.96

We are Accepting NULL HYPOTHESIS for 'IN deaths' because absolute value of Z 1 sample testing(Z) is 1.2191958440379214 and it's greater than Z value(at 0.025) = 1.96

**Is the test applicable?**

There are 2 requirements for the Z test to be applicable:

1. std dev of the distribution is known : Satisfied
2. Either n tends to infinity ( $n \geq 30$ ) or dataset is normal : Satisfied because  $n \geq 30$

Hence the test is applicable

**▼ T 1 sample testing for mean of Mar '21**

```

1 def T_test_single(feature, sample_mean_feb, sample_variance_mar):
2     lambda_mar = mar_month[feature].mean()
3     T = (lambda_mar - sample_mean_feb)/np.sqrt(sample_variance_mar/len(mar_month))
4
5     T = abs(T)
6
7     if(T > 2.042):
8         print("We are rejecting NULL HYPOTHESIS for '" + feature + "' because absol
9     else:
10        print("We are Accepting NULL HYPOTHESIS for '" + feature + "' because absol
11
12    # print(T)
13    print('\n\n')
14
15
16
17 for feature in features:
18     sample_mean_feb = feb_month[feature].mean()
19     sample_variance_mar = variance(mar_month[feature].values)
20     # print(np.sqrt(sample_variance_mar))
21     T_test_single(feature, sample_mean_feb, sample_variance_mar)

```

We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of T

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of T 1 sample testing(Z) is 20.298201597921185 and it's greater than  $T(30,0.025) = 2.042$

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of T 1 sample testing(Z) is 25.70299698367066 and it's greater than  $T(30,0.025) = 2.042$

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of T 1 sample testing(Z) is 37.58384741574119 and it's greater than  $T(30,0.025) = 2.042$

## Result of T 1 sample testing for mean of cases and death

### Null hypothesis (H0):

Mean of cases/deaths in Mar '21 = Mean of cases/deaths in Feb '21.

### Alternate hypothesis(H1):

Mean of cases/deaths in Mar '21 not equals Mean of cases/deaths in Feb '21.

### Result:

We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of T 1 sample testing(Z) is 17.790801279678224 and it's greater than  $T(30,0.025) = 2.042$

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of T 1 sample testing(Z) is 20.298201597921185 and it's greater than  $T(30,0.025) = 2.042$

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of T 1 sample testing(Z) is 25.70299698367066 and it's greater than  $T(30,0.025) = 2.042$

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of T 1 sample testing(Z) is 37.58384741574119 and it's greater than  $T(30,0.025) = 2.042$

### Is the test applicable?

T test is applicable when n value is very small, but here  $n \geq 30$  which means n tends to infinity.

Hence the test is not applicable

## ▼ Walds 2 sample testing for mean of Feb '21 and Mar '21

```
1 def walds_test_double(feature):
2     lambda_feb = feb_month[feature].mean()
3     lambda_mar = mar_month[feature].mean()
```

```

4
5
6     delta = lambda_mar - lambda_feb
7     W = (delta)/(np.sqrt(lambda_feb/len(feb_month) + lambda_mar/(len(mar_month))))
8
9     W = abs(W)
10
11     if(W > 1.96):
12         print("We are rejecting NULL HYPOTHESIS for '" + feature + "' because absol
13     else:
14         print("We are Accepting NULL HYPOTHESIS for '" + feature + "' because absol
15
16     # print(W)
17     print('\n\n')
18
19
20
21
22
23 for feature in features:
24     walds_test_double(feature)

```

We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of Wa

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of Wa

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of Walds

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of Walds

## Result of Walds 2 sample testing for mean of cases and death

### Null hypothesis (H0):

Mean of cases/deaths in Mar '21 = Mean of cases/deaths in Feb '21.

### Alternate hypothesis(H1):

Mean of cases/deaths in Mar '21 not equals Mean of cases/deaths in Feb '21.

### Result:

We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of Walds 2 sample testing(W) is 187.4665744482369 and it's greater than Z value(at 0.025) = 1.96

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of Walds 2 sample testing(W) is 125.21233437856539 and it's greater than Z value(at 0.025) = 1.96

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of Walds 2 sample testing(W) is 27.16225354349577 and it's greater than Z value(at 0.025) = 1.96

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of Walds 2 sample testing(W) is 31.136159361564246 and it's greater than Z value(at 0.025) = 1.96

### Is the test applicable?

Since we used MLE in case of estimator, it is Asymptotically Normal as n tends to infinity ( $n \geq 30$  is satisfied). This is the only condition, hence the test is applicable.

## ▼ T 2 sample unpaired testing for mean of Feb '21 and Mar '21

```

1 def T_test_double(feature, sample_variance_feb, sample_variance_mar):
2     lambda_feb = feb_month[feature].mean()
3     lambda_mar = mar_month[feature].mean()
4
5
6     delta = lambda_mar - lambda_feb
7     T = (delta)/(np.sqrt(sample_variance_feb/len(feb_month) + sample_variance_mar/len(mar_month)))
8     T = abs(T)
9
10    if(T > 2.0025):
11        print("We are rejecting NULL HYPOTHESIS for '" + feature + "' because absolute value of T is greater than 2.0025")
12    else:
13        print("We are Accepting NULL HYPOTHESIS for '" + feature + "' because absolute value of T is less than 2.0025")
14
15    # print(T)
16    print('\n\n')
17
18
19 for feature in features:
20     sample_variance_feb = variance(feb_month[feature].values)
21     sample_variance_mar = variance(mar_month[feature].values)
22     # print(np.sqrt(sample_variance_mar))
23     T_test_double(feature, sample_variance_feb, sample_variance_mar)
24

```

We are Accepting NULL HYPOTHESIS for 'IL confirmed' because absolute value of T :

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of T :

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of T 2 s:

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of T 2 s:

## Result of Walds 2 sample testing for mean of cases and death

### Null hypothesis (H0):

Mean of cases/deaths in Mar '21 = Mean of cases/deaths in Feb '21.

### Alternate hypothesis(H1):

Mean of cases/deaths in Mar '21 not equals Mean of cases/deaths in Feb '21.

### Result:

We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of T 2 sample testing(Z) is 12.023210757714896 and it's greater than  $T(57, 0.025) = 2.0025$

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of T 2 sample testing(Z) is 11.557631363866554 and it's greater than  $T(57, 0.025) = 2.0025$

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of T 2 sample testing(Z) is 12.279530786054828 and it's greater than  $T(57, 0.025) = 2.0025$

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of T 2 sample testing(Z) is 6.490636933322033 and it's greater than  $T(57, 0.025) = 2.0025$

### Is the test applicable?

T test is applicable when n value is very small, but here  $n \geq 30$  which means n tends to infinity.

Hence the test is not applicable

## ▼ PART 2c : Stats of last 3 months in 2020

### ▼ Calculating MME parameters required

**\*\*** Calculating all the MME parameters for IL(ILLINOIS) state and using them for all the 1 sample KS tests below as guess values on IN

## state(INDIANA) \*\*

```

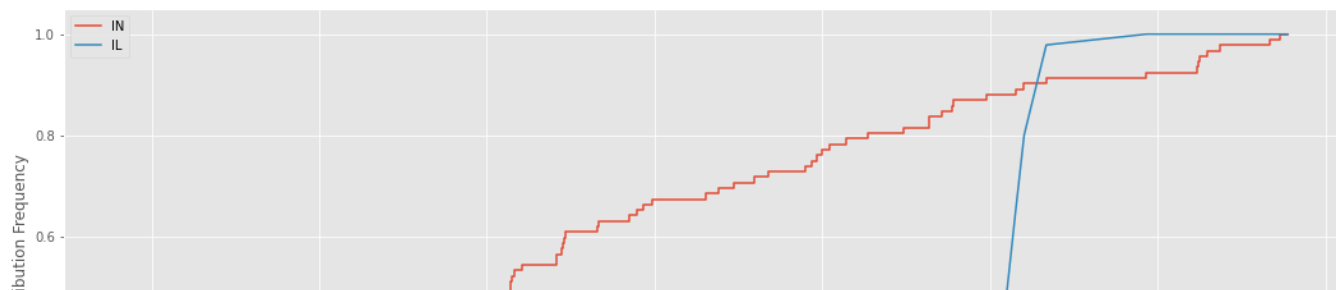
1 def calculate_MSE(y1, y2):
2     return numpy.sum((y1 - y2) * (y1 - y2)) / len(y1)
3
4
5 def plot_eCDF(input_list, label, color):
6     input_list.sort()
7     # sort input array
8     n = len(input_list)
9
10    # initialize x and y to plt CDF
11    x = [input_list[0]]
12    y = [0]
13    for point in input_list:
14        value = y[len(y) - 1] + 1 / n
15        # update x and y values
16        x = x + [point, point]
17        y = y + [y[len(y) - 1], value]
18
19    # eCDF step function plot
20    plt.plot(x, y, label=label, color=color)
21    return x[1:], y[1:]
22
23
24 data = df.iloc[0:,0:4].to_numpy()
25 data = data[253:345, 0:].astype(numpy.float64)
26
27 # NUMBER OF CASES
28 cases_il = data[:, 0]
29 cases_states_sample_mean = numpy.mean(cases_il)
30 cases_states_sample_variance = numpy.var(cases_il)
31 cases_states_mme_poisson = cases_states_sample_mean
32 cases_state_mme_geometric = 1/cases_states_sample_mean
33 cases_states_mme_p_binomial = 1 - cases_states_sample_variance / cases_states_samp
34 cases_state_mme_n_binomial = cases_states_sample_mean**2/(cases_states_sample_mean-
35
36 # NUMBER OF DEATHS
37 deaths_il = data[:, 2]
38 deaths_sample_mean = numpy.mean(deaths_il)
39 deaths_sample_variance = numpy.var(deaths_il)
40 deaths_mme_poisson = deaths_sample_mean
41 deaths_mme_geometric = 1 / deaths_sample_mean
42 deaths_mme_p_binomial = 1 - deaths_sample_variance / deaths_sample_mean
43 deaths_mme_n_binomial = deaths_sample_mean ** 2 / (deaths_sample_mean - deaths_sam
44
45

```

## ▼ KS 1 sample test with poisson as distribution for IN Confirmed



```
1 # print(cases_states_mme_poisson)
2 cases_in = data[:, 1]
3 cases_in = numpy.sort(cases_in)
4 cdf_y = numpy.array([])
5 cdf = 0
6 n = len(cases_in)
7 max_diff = 0
8 poisson_xpoint = 0
9 poisson_cdf = numpy.array([])
10 for i in cases_in:
11     poisson_point = poisson.cdf(i, cases_states_mme_poisson)
12     poisson_cdf = numpy.append(poisson_cdf, poisson_point)
13     if max_diff < numpy.abs(cdf-poisson_point):
14         max_diff = numpy.abs(cdf-poisson_point)
15         poisson_xpoint = poisson_point
16     cdf += 1/n
17     cdf_y = numpy.append(cdf_y, cdf)
18
19 plt.figure('Cases Case', figsize=(18,8))
20 plt.xlabel('Poisson Distribution')
21 plt.ylabel('Cumulative Distribution Frequency')
22 plt.step(cases_in, cdf_y, label = "IN")
23 plt.plot(cases_in, poisson_cdf, label = "IL")
24 plt.legend(loc="upper left")
25 plt.show()
26
27 if max_diff > 0.05:
28     print("We reject NULL hypothesis, because Max value is "+str(max_diff)+" > c :")
29 else:
30     print("We accept NULL hypothesis, because Max value is "+str(max_diff)+" <= c :
```



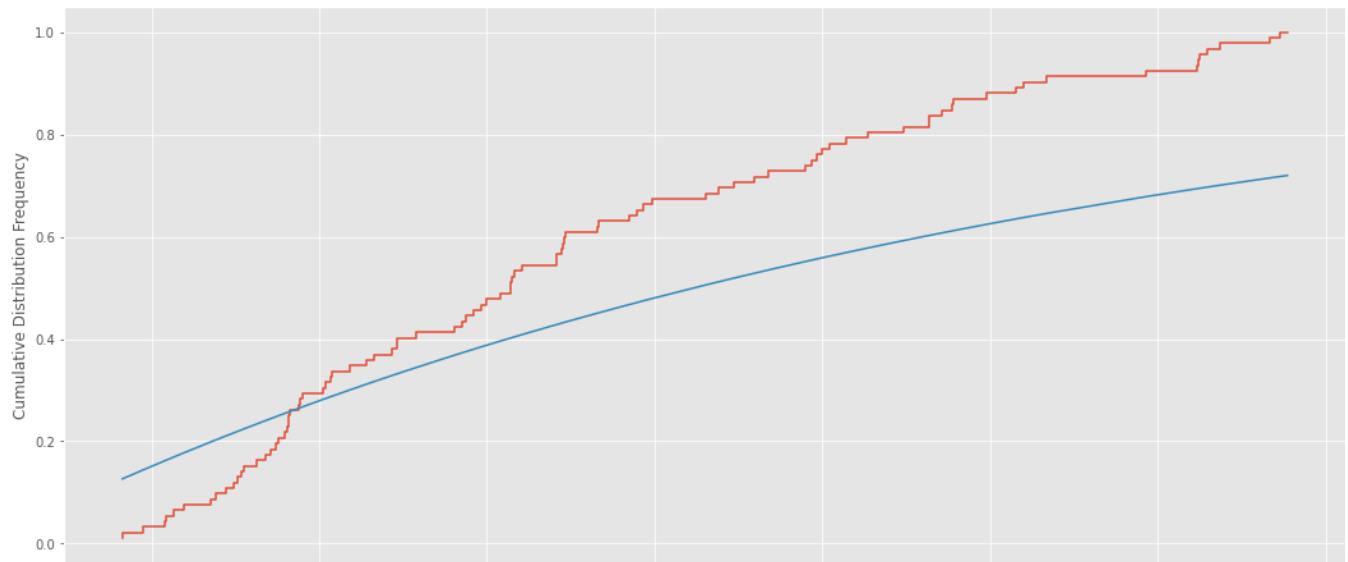
## ▼ KS 1 sample test with geometric as distribution for IN Confirmed



```

1 # print(cases_state_mme_geometric)
2 cdf_y = numpy.array([])
3 cdf = 0
4 max_diff = 0
5 geom_cdf = numpy.array([])
6 geometric_dipoint = 0
7 for i in cases_in:
8     geom_point = geom.cdf(i, cases_state_mme_geometric)
9     geom_cdf = numpy.append(geom_cdf, geom_point)
10    if max_diff < numpy.abs(cdf - geom_point):
11        max_diff = numpy.abs(cdf - geom_point)
12        geometric_dipoint = geom_point
13    cdf += 1 / n
14    cdf_y = numpy.append(cdf_y, cdf)
15
16 plt.figure('Cases Case', figsize=(18,8))
17 plt.xlabel('Geometric Distribution')
18 plt.ylabel('Cumulative Distribution Frequency')
19 plt.step(cases_in, cdf_y, label = "IN")
20 plt.plot(cases_in, geom_cdf, label = "IL")
21 plt.show()
22
23 if max_diff > 0.05:
24     print("We reject NULL hypothesis, because Max value is "+str(max_diff)+" > c :")
25 else:
26     print("We accept NULL hypothesis, because Max value is "+str(max_diff)+" <= c :")

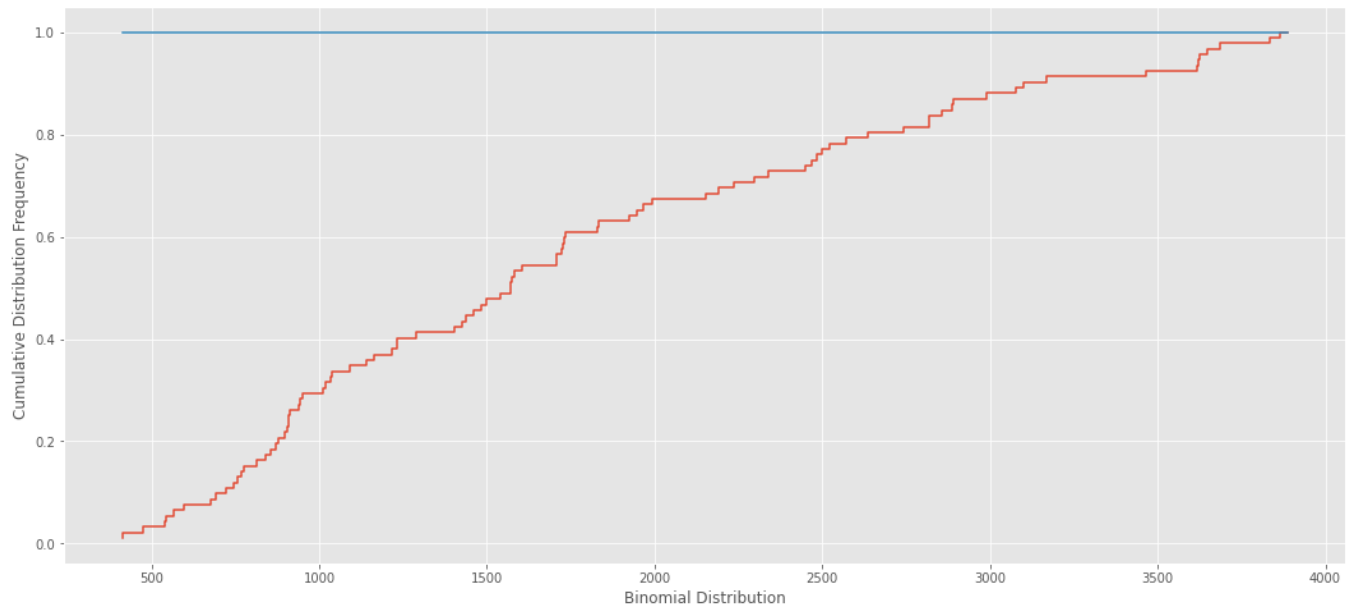
```



## ▼ KS 1 sample test with binomial as distribution for IN Confirmed

```

1 cdf_y = numpy.array([])
2 cdf = 0
3 max_diff = 0
4 binom_cdf = numpy.array([])
5 for i in cases_in:
6     binom_point = binom.cdf(i, cases_state_mme_n_binomial, cases_states_mme_p_binom
7     binom_cdf = numpy.append(binom_cdf, binom_point)
8     if max_diff < numpy.abs(cdf - binom_point):
9         max_diff = numpy.abs(cdf - binom_point)
10    cdf += 1 / n
11    cdf_y = numpy.append(cdf_y, cdf)
12
13 plt.figure('Cases Case', figsize=(18,8))
14 plt.xlabel('Binomial Distribution')
15 plt.ylabel('Cumulative Distribution Frequency')
16 plt.step(cases_in, cdf_y, label = "IN")
17 plt.plot(cases_in, binom_cdf, label = "IL")
18 plt.show()
19
20 if max_diff > 0.05:
21     print("We reject NULL hypothesis, because Max value is "+str(max_diff)+" > c :
22 else:
23     print("We accept NULL hypothesis, because Max value is "+str(max_diff)+" <= c :
```



We reject NULL hypothesis, because Max value is 1.0 >  $\alpha$  : 0.05

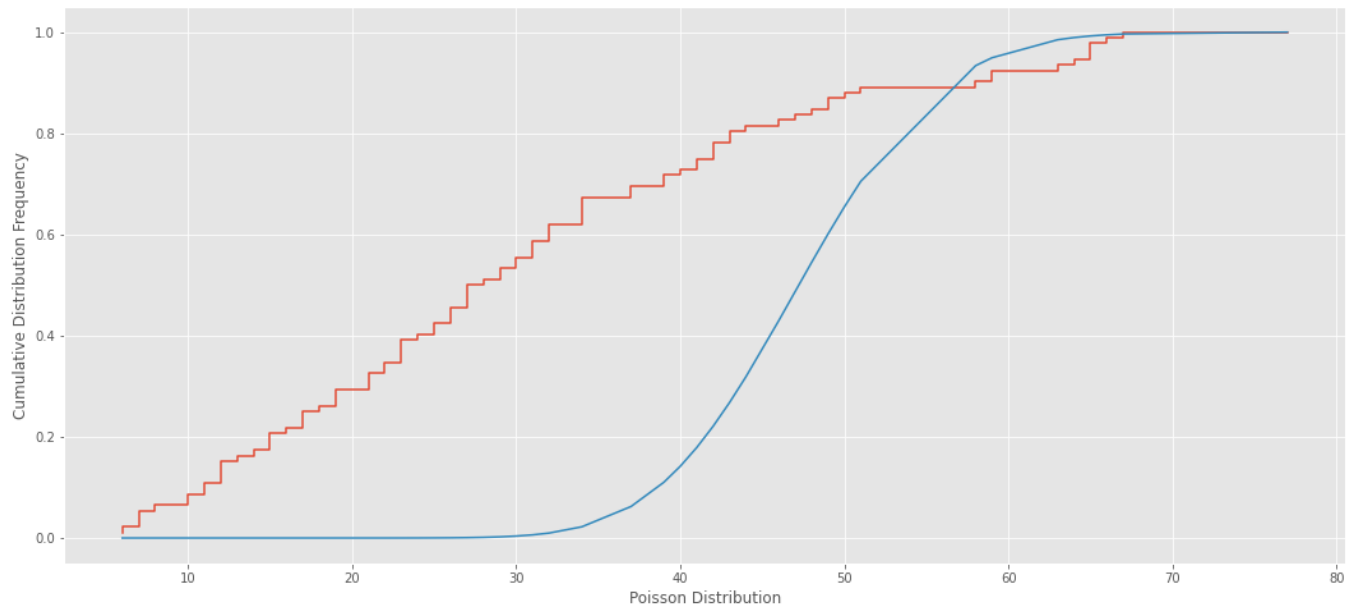
## ▼ KS 1 sample test with poisson as distribution for IN Deaths

```

1
2 # print(deaths_mme_poisson)
3 deaths_in = data[:, 3]
4 deaths_in = numpy.sort(deaths_in)
5 cdf_y = numpy.array([])
6 cdf = 0
7 n = len(deaths_in)
8 max_diff = 0
9 poisson_cdf = numpy.array([])
10 for i in deaths_in:
11     poisson_point = poisson.cdf(i, deaths_mme_poisson)
12     poisson_cdf = numpy.append(poisson_cdf, poisson_point)
13     if max_diff < numpy.abs(cdf - poisson_point):
14         max_diff = numpy.abs(cdf - poisson_point)
15     cdf += 1 / n
16     cdf_y = numpy.append(cdf_y, cdf)
17
18 plt.figure('Deaths Case', figsize=(18,8))
19 plt.xlabel('Poisson Distribution')
20 plt.ylabel('Cumulative Distribution Frequency')
21 plt.step(deaths_in, cdf_y, label = "IN")
22 plt.plot(deaths_in, poisson_cdf, label = "IL")
23 plt.show()
24
25 if max_diff > 0.05:
26     print("We reject NULL hypothesis, because Max value is "+str(max_diff)+" >  $\alpha$  :
27 else:
28     print("We accept NULL hypothesis, because Max value is "+str(max_diff)+" <=  $\alpha$  :

```

```
20 print( we accept NULL hypothesis, because max value is 0.6300783615975943 > c : 0.05 )
```



We reject NULL hypothesis, because Max value is 0.6300783615975943 > c : 0.05

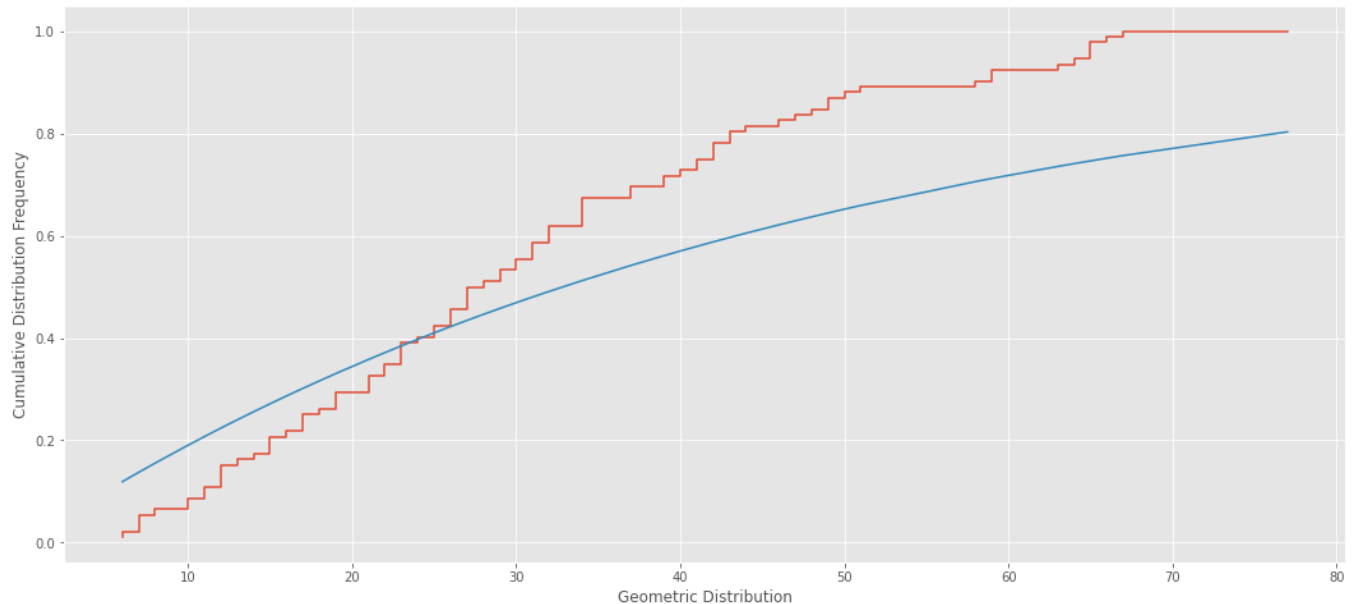
## ▼ KS 1 sample test with geometric as distribution for IN Deaths

```
1
2 # print(deaths_mme_geometric)
3 cdf_y = numpy.array([])
4 cdf = 0
5 max_diff = 0
6 geom_cdf = numpy.array([])
7 for i in deaths_in:
8     geom_point = geom.cdf(i, deaths_mme_geometric)
9     geom_cdf = numpy.append(geom_cdf, geom_point)
10    if max_diff < numpy.abs(cdf - geom_point):
11        max_diff = numpy.abs(cdf - geom_point)
12    cdf += 1 / n
13    cdf_y = numpy.append(cdf_y, cdf)
14
15 plt.figure('Deaths Case', figsize=(18,8))
16 plt.xlabel('Geometric Distribution')
17 plt.ylabel('Cumulative Distribution Frequency')
18 plt.step(deaths_in, cdf_y, label = "IN")
```

```

19 plt.plot(deaths_in, geom_cdf, label = "IL")
20 plt.show()
21
22
23
24 if max_diff > 0.05:
25     print("We reject NULL hypothesis, because Max value is "+str(max_diff)+" > c :
26 else:
27     print("We accept NULL hypothesis, because Max value is "+str(max_diff)+" <= c :

```



We reject NULL hypothesis, because Max value is 0.22147478129654763 > c : 0.05

## ▼ KS 1 sample test with binomial as distribution for IN Deaths

```

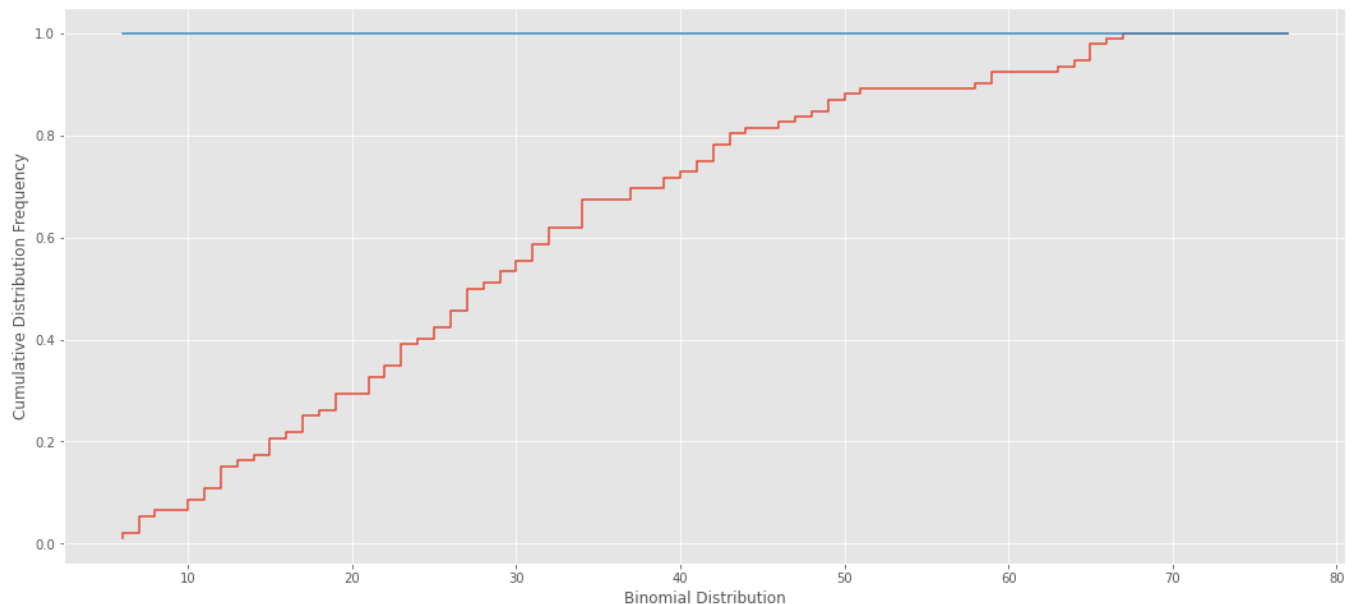
1
2 # print(deaths_mme_n_binomial, deaths_mme_p_binomial)
3 cdf_y = numpy.array([])
4 cdf = 0
5 max_diff = 0
6 binom_cdf = numpy.array([])
7 for i in cases_in:
8     binom_point = binom.cdf(i, deaths_mme_n_binomial, deaths_mme_p_binomial)
9     binom_cdf = numpy.append(binom_cdf, binom_point)
10     if max_diff < numpy.abs(cdf - binom_point):

```

```

11     max_diff = numpy.abs(cdf - binom_point)
12     cdf += 1 / n
13     cdf_y = numpy.append(cdf_y, cdf)
14
15 plt.figure('Deaths Case', figsize=(18,8))
16 plt.xlabel('Binomial Distribution')
17 plt.ylabel('Cumulative Distribution Frequency')
18 plt.step(deaths_in, cdf_y, label = "IN")
19 plt.plot(deaths_in, binom_cdf, label = "IL")
20 plt.show()
21
22
23 if max_diff > 0.05:
24     print("We reject NULL hypothesis, because Max value is "+str(max_diff)+" > c :
25 else:
26     print("We accept NULL hypothesis, because Max value is "+str(max_diff)+" <= c :

```



We reject NULL hypothesis, because Max value is 1.0 > c : 0.05

## KS 2 sample test with IL Confirmed as dataset1 and IN confirmed as dataset2

### NULL HYPOTHESIS:

## Distribution of dataset1 = Distribution of dataset2

```

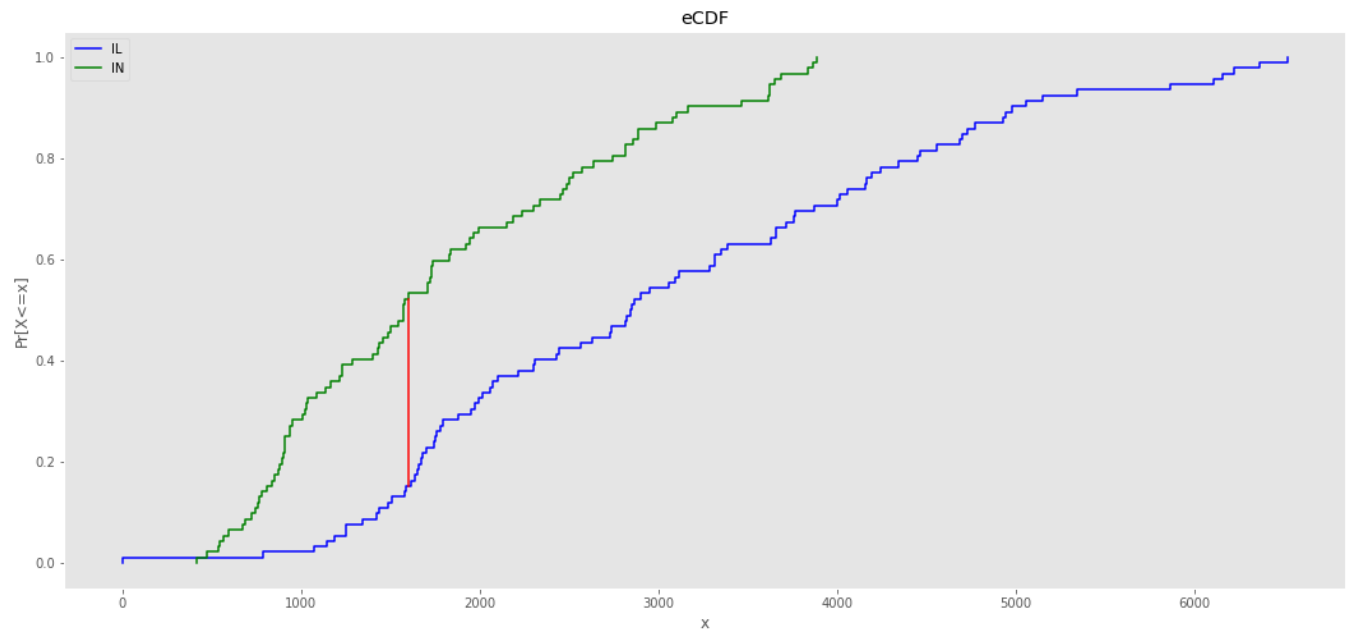
1
2 # TWO POPULATION TEST
3 X = data[:, 0]
4 Y = data[:, 1]
5
6 # PLOT THE GRAPHS
7 plt.figure('eCDF', figsize=(18,8))
8 x1, y1 = plot_eCDF(X, 'IL', color='blue')
9 x2, y2 = plot_eCDF(Y, 'IN', color='green')
10
11 max_difference = 0
12 point = 0
13 point_y1 = 0
14 point_y2 = 0
15 i = 0
16 j = 0
17 while i < len(x2):
18     y2_left, y2_right = y2[i], y2[i + 1]
19     while j + 2 < len(x1) and x1[j + 2] < x2[i]:
20         j += 2
21     if x2[i] == x1[j]:
22         y1_left, y1_right = y1[j], y1[j + 1]
23     else:
24         y1_left, y1_right = y1[j + 1], y1[j + 1]
25     if max_difference < numpy.max([max_difference, numpy.absolute(y1_left - y2_left)
26         max_difference = numpy.max([max_difference, numpy.absolute(y1_left - y2_left)
27         point_y1 = y1_left
28         point_y2 = y2_left
29         point = x2[i]
30     i += 2
31 print('Max Difference: ', max_difference)
32 print('Point with max Difference: ', point)
33
34 # GRAPH ATTRIBUTES
35 plt.plot([point, point], [point_y1, point_y2], color='red')
36 plt.xlabel('x')
37 plt.ylabel('Pr[X<=x]')
38 plt.title('eCDF')
39 plt.legend(loc="upper left")
40 plt.grid()
41 plt.show()
42
43
44 if max_difference > 0.05:
45     print("We reject NULL hypothesis, because Max value is "+str(max_difference)+"
46 else:
47     print("We accept NULL hypothesis, because Max value is "+str(max_difference)+"

```



Max Difference: 0.3804347826086955

Point with max Difference: 1605.0



We reject NULL hypothesis, because Max value is 0.3804347826086955 >  $\alpha$  : 0.05

## KS 2 sample test with IL Deaths as dataset1 and IN Deaths as dataset2

### NULL HYPOTHESIS:

Distribution of dataset1 = Distribution of dataset2

```

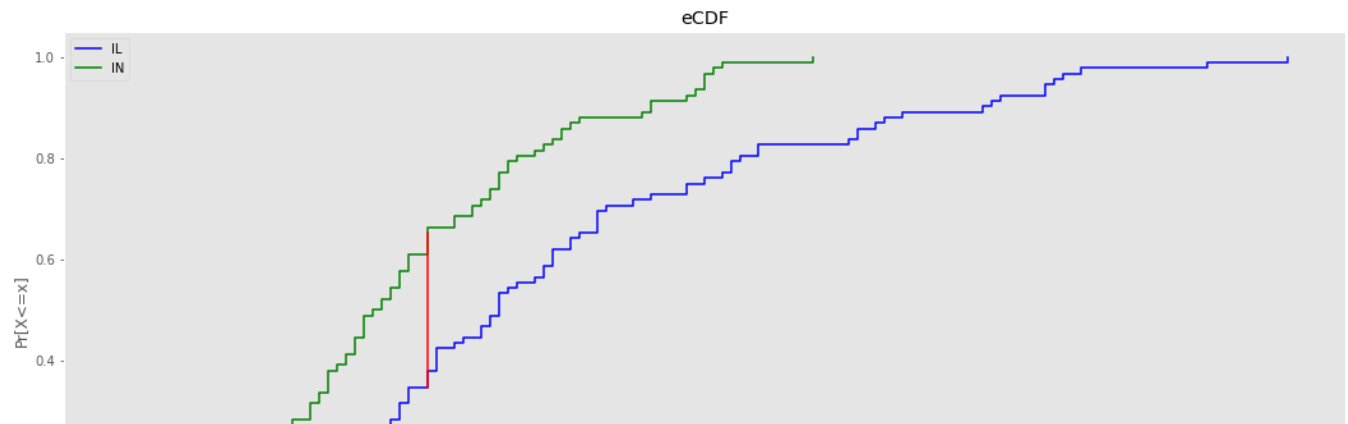
1
2 # TWO POPULATION TEST
3 X = data[:, 2]
4 Y = data[:, 3]
5
6 # PLOT THE GRAPHS
7 plt.figure('eCDF', figsize=(18,8))
8 x1, y1 = plot_eCDF(X, 'IL', color='blue')
9 x2, y2 = plot_eCDF(Y, 'IN', color='green')
10
11 max_difference = 0
12 point = 0
13 point_y1 = 0

```

```
14 point_y2 = 0
15 i = 0
16 j = 0
17 while i < len(x2):
18     y2_left, y2_right = y2[i], y2[i + 1]
19     while j + 2 < len(x1) and x1[j + 2] < x2[i]:
20         j += 2
21     if x2[i] == x1[j]:
22         y1_left, y1_right = y1[j], y1[j + 1]
23     else:
24         y1_left, y1_right = y1[j + 1], y1[j + 1]
25     if max_difference < numpy.max([max_difference, numpy.absolute(y1_left - y2_left)
26     max_difference = numpy.max([max_difference, numpy.absolute(y1_left - y2_left)
27     point_y1 = y1_left
28     point_y2 = y2_left
29     point = x2[i]
30     i += 2
31 print('Max Difference: ', max_difference)
32 print('Point with max Difference: ', point)
33
34 # GRAPH ATTRIBUTES
35 plt.plot([point, point], [point_y1, point_y2], color='red')
36 plt.xlabel('x')
37 plt.ylabel('Pr[X<=x]')
38 plt.title('eCDF')
39 plt.legend(loc="upper left")
40 plt.grid()
41 plt.show()
42
43
44 if max_difference > 0.05:
45     print("We reject NULL hypothesis, because Max value is "+str(max_difference)+"
46 else:
47     print("We accept NULL hypothesis, because Max value is "+str(max_difference)+"
```

Max Difference: 0.3152173913043484

Point with max Difference: 34.0



## Permutation test

0.0 -

## Permutation test with IL Confirmed as dataset1 and IN Confirmed as dataset2

### NULL HYPOTHESIS:

Distribution of dataset1 = Distribution of dataset2

```

1
2 data = df.iloc[0:,0:4].to_numpy()
3 data = data[253:345, 0:].astype(numpy.float64)
4
5 # CASES
6 x1 = data[:, 0]
7 y1 = data[:, 1]
8
9 X_avg = numpy.mean(x1)
10 Y_avg = numpy.mean(y1)
11 t_obs = numpy.absolute(X_avg - Y_avg)
12
13 number = 0
14 combined = numpy.append(x1, y1, axis=0)
15 for i in range(1000):
16     permutation = numpy.random.permutation(combined)
17     X_permutation = permutation[: len(x1)]
18     Y_permutation = permutation[len(x1):]
19     t_predict = numpy.absolute(numpy.mean(X_permutation) - numpy.mean(Y_permutation))
20     if t_predict > t_obs:
21         number += 1
22
23 p_value = number / 1000
24 print(p_value)
25 if p_value > 0.05:

```

```

25 if p_value > 0.05:
26     print('Reject Null Hypothesis as p value is '+str(p_value)+' which is greater t
27 else:
28     print('Accept Null Hypothesis as p value is '+str(p_value)+' which is greater t
29
30
0.0
Accept Null Hypothesis as p value is 0.0 which is greater than 0.5

```

## Permutation test with IL Deaths as dataset1 and IN Deaths as dataset2

### NULL HYPOTHESIS:

Distribution of dataset1 = Distribution of dataset2

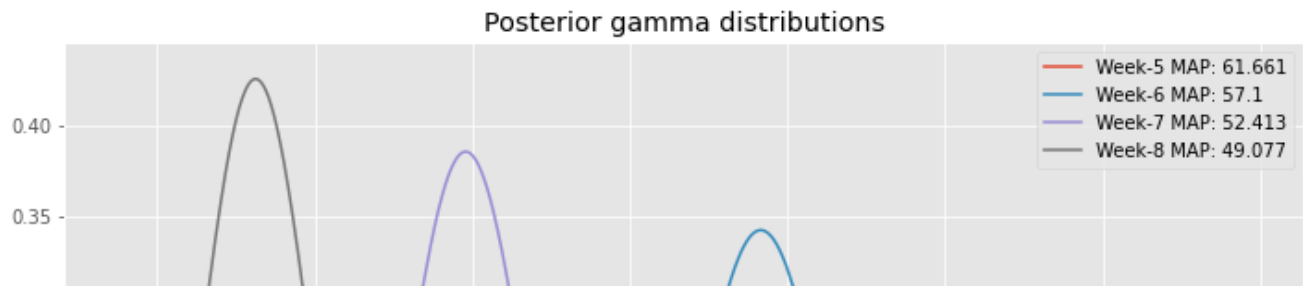
```

1 # DEATHS
2 x2 = data[:, 2]
3 y2 = data[:, 3]
4
5 X_avg = numpy.mean(x2)
6 Y_avg = numpy.mean(y2)
7 t_obs = numpy.absolute(X_avg - Y_avg)
8
9 number = 0
10 combined = numpy.append(x2, y2, axis=0)
11 for i in range(1000):
12     permutation = numpy.random.permutation(combined)
13     X_permutation = permutation[: len(x2)]
14     Y_permutation = permutation[len(x2):]
15     t_predict = numpy.absolute(numpy.mean(X_permutation) - numpy.mean(Y_permutation))
16     if t_predict > t_obs:
17         number += 1
18
19 p_value = number / 1000
20 print(p_value)
21 if p_value > 0.05:
22     print('Reject Null Hypothesis as p value is '+str(p_value)+' which is greater t
23 else:
24     print('Accept Null Hypothesis as p value is '+str(p_value)+' which is greater t
25
0.0
Accept Null Hypothesis as p value is 0.0 which is greater than 0.5

```

## PART 2d : Prediction using Bayesian Inference

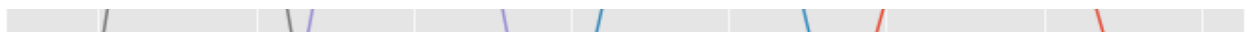
```
1 INdeaths = df[df['Date'] > '2020-05-31']['IN deaths']
2 ILdeaths = df[df['Date'] > '2020-05-31']['IL deaths']
3
4 # print(sum(INdeaths[:35]+ILdeaths[:35]))
5
6 INandILdeaths=list(INdeaths+ILdeaths)
7
8 # print(INandILdeaths)
9
10 def plot_gamma_distribution(alpha=1, beta=1, label="0"):
11     x = np.linspace(gamma.ppf(0.01, alpha, scale=1/beta), gamma.ppf(0.99, alpha, scale=1/beta), 100)
12     plt.title("Posterior gamma distributions")
13     label= label + " MAP: " + str(round((alpha/beta),3))
14     plt.xlabel("Deaths")
15     plt.ylabel("Gamma distribution")
16     plt.plot(x, gamma.pdf(x, alpha, scale=1/beta), label=label)
17     plt.legend()
18
19 beta=np.mean(INandILdeaths[:28])
20 # print(beta)
21 # print(sum(INandILdeaths[:35]))
22 # print(sum(INandILdeaths[:42]))
23 plt.figure(figsize=(12,8))
24
25 for i in range(4,8):
26     alpha = sum(INandILdeaths[(i+1)*7])+1
27     # print(beta)
28     scale = (i+1)*7 + 1/beta
29     plot_gamma_distribution(alpha, scale, 'Week-'+str(i+1))
30
31 plt.show()
```



These diagrams portrays that as the weeks advances, the quantity of passings are expanding and hence the MAP for the Lambda boundary is expanding.

Trust in MAP esteem is expanding as the weeks are expanding.

We can likewise see that the pace of increment of MAP of Lambda is diminishing. Along these lines we can induce that increment in passings each week is going towards immersion.



## ▼ PART 3 : EXPLORATORY TASK



Our data set is the San Francisco crime data set. California State COVID data is used for the inferences.

## ▼ Inference 1 : Correlation between SanFransisco crime and California Covid data

```

1 def CorrelationCoeff(X,Y):
2     X = X[57:88,1]
3     Y = Y[57:88,1]
4     X = X.astype(np.int)
5     Y = Y.astype(np.int)
6     return np.corrcoef(X, Y)
7
8
9 def Preprocessdata( X ):
10    l = len(X)
11    prev = X[0][0]
12    count = 1;
13    Y = []
14    for i in range(0,l):
15        if(X[i][0]==prev):
16            count=count+1
17        else:
18            temp = [prev,int(count)]
19            Y = Y+temp
20            prev = X[i][0]

```

```

21         count = 1
22     Y = Y + [prev,count]
23     return np.reshape(np.array(Y),(len(Y)//2,2))
24
25 def Process(df):
26     df1 = df.iloc[:,0]
27     df.drop(df.columns[[0]], axis = 1, inplace = True)
28     for key, value in df.iteritems():
29         po = value
30         ope = []
31         co = 0
32         prev = 0
33         for j in po:
34             if(co == 0):
35                 ope.append(j)
36                 prev = j
37             else:
38                 ope.append(j-prev)
39                 prev = j
40             co = co + 1
41         df[key] = ope
42     df["Date"] = df1.values
43     columns_titles = ["Date","Cases"]
44     df=df.reindex(columns=columns_titles)
45     return df
46
47 def CorrelationCoeffWholeYear(X,Y):
48     X = X[:,1]
49     Y = Y[:,1]
50     X = X.astype(np.int)
51     Y = Y.astype(np.int)
52     return np.corrcoef(X, Y)

1 cols_list = ['Incident Date']
2 df = pd.read_csv('gdrive/MyDrive/Colab Notebooks/SF_CrimeDataset.csv', usecols=cols
3
4 Covid_data = pd.read_csv('gdrive/MyDrive/Colab Notebooks/US_confirmed.csv')
5 Covid_data.head()
6
7 CA_Covid_data = Covid_data.iloc[4].to_frame().reset_index().iloc[1:]
8 CA_Covid_data.columns = ['Date', 'Cases']
9 CA_Covid_data['Date'] = pd.to_datetime(CA_Covid_data['Date'])
10 CA_Covid_data = CA_Covid_data[(CA_Covid_data['Date'] >= '2020/01/01') & (CA_Covid_c
11 CA_Covid_data = Process(CA_Covid_data)
12
13
14 start2019 = "2019/01/01"
15 end2019 = "2019/12/31"
16 bef_covid_data = df[(df['Incident Date'] >= start2019) & (df['Incident Date'] <= er
17
18 start20 = "2020/01/01"

```

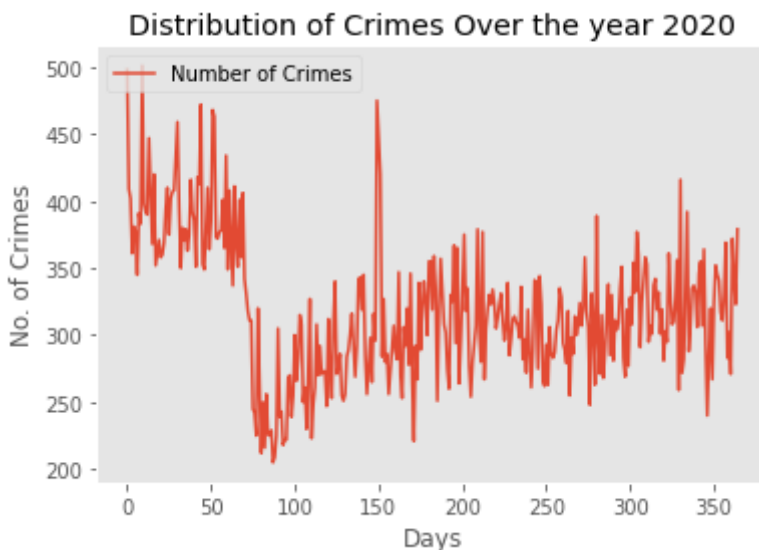
```

18 start20 = "2020/01/01"
19 end20 = "2020/12/31"
20 aft_covid_data = df[(df['Incident Date'] >= start20) & (df['Incident Date'] <= end20)]
21
22 bef_covid_data = np.sort(bef_covid_data.to_numpy(), axis = 0)
23
24 bef_covid_data = Preprocessdata(bef_covid_data)
25
26 aft_covid_data = np.sort(aft_covid_data.to_numpy(), axis = 0)
27 aft_covid_data = Preprocessdata(aft_covid_data)

1 print("Correlation Coefficient considering the whole Year 2020 is : ", CorrelationCoeff)
2
3 X = [i for i in range(365)]
4 plt.figure('Correlation of Covid and Number of Crimes')
5 plt.plot(X, aft_covid_data[:,1].astype(int) ,label='Number of Crimes')
6 plt.xlabel('Days')
7 plt.ylabel('No. of Crimes')
8 plt.title('Distribution of Crimes Over the year 2020')
9 plt.legend(loc="upper left")
10 plt.grid()
11 plt.show()
12 print("Correlation Coefficient considering the month of March 2020 is : ", CorrelationCoeff)
13

```

Correlation Coefficient considering the whole Year 2020 is : -0.0312750345318851



Correlation Coefficient considering the month of March 2020 is : -0.768499925181

We tried to find the effect of COVID on the San Francisco Crime. For this, we computed the correlation between the COVID data set of California state with San Francisco.

### **Correlation in the Year of Covid 2020 is found to be -0.03**

This tells us that there is no correlation between crimes and COVID. We are expecting a strong negative correlation of cases with the COVID due to the factors like lockdown etc. So, we observed



the plot for the crimes in the year 2020. It is found that the cases went down suddenly in the months of MARCH and then slowly rose. So, in total there is not much correlation throughout the year. So, we tried to get the correlation in the month of MARCH.

### **Correlation in the Month of MARCH is found to be -0.77**

This can be assumed as a good negative correlation. Key Point observed is the Emergency lockdown has happened in the month of MARCH resulting in the reduction of number of crimes. So, we can conclude that our estimate is reasonable.

## **Inference 2 : Permutation test to check whether the distribution of**

### **Sanfransico crime before covid and after covid is same or not bold text**

```

1 def PermutationTest(X, Y, rounds):
2     X = X.astype(np.int)
3     Y = Y.astype(np.int)
4     count = 0
5     diff = abs(np.mean(X) - np.mean(Y))
6     tot = X+Y
7     n = len(tot)
8     for i in range(0,rounds):
9         a=np.random.permutation(tot)[:n//2]
10        b=np.random.permutation(tot)[n//2:]
11        if (abs(a.mean() - b.mean())>diff):
12            count+=1
13    return (count/rounds)
14 print("P-value in Permutation Test", PermutationTest(bef_covid_data[:,1],aft_covid_

P-value in Permutation Test 0.0

```

**Null Hypothesis (Ho)** Distribution of Crimes an year before Covid 19 equals with after Covid 19

**Alternate Hypothesis (H1)** Distribution of Crimes an year before Covid 19 not equals with after Covid 19

#### **Procedure:**

We have permuted all the data for 365 days of Crimes before the Covid i.e for 2019 and for all the data for 365 days of Crimes after Covid start i.e for 2020 for 1000 permutations.

#### **Result:**

We have got a permutation test value of 0.0 which is much less than the threshold of 0.05. So we reject the null Hypothesis and accept the Alternate hypothesis.

Hence with the Permutation test, we can infer that the distribution of crimes before and after covid are not equal.

## ▼ Inference 3 : Chi-Square Independence test

```

1 def chi_square(matrix_covid_vehicle):
2     rows = matrix_covid_vehicle.shape[0]
3     cols = matrix_covid_vehicle.shape[1]
4     df = (rows-1)*(cols-1)
5     total_row1,total_row2 = np.sum(matrix_covid_vehicle,axis=0)
6     total_col1,total_col2 = np.sum(matrix_covid_vehicle,axis=1)
7     total = total_row1+total_row2
8     expected_values = np.zeros([2,2])
9     expected_values[0][0] = (float(total_col1)*total_row1)/(total)
10    expected_values[0][1] = (float(total_col2)*total_row1)/(total)
11    expected_values[1][0] = (float(total_col1)*total_row2)/(total)
12    expected_values[1][1] = (float(total_col2)*total_row2)/(total)
13    q_expected = 0.0
14    for i in range(rows):
15        for j in range(cols):
16            q_expected = q_expected + ((expected_values[i][j] - matrix_covid_vehicle[i][j])**2)/expected_values[i][j]
17    return (q_expected,df)
18
19 def initial():
20     df = pd.read_csv('gdrive/MyDrive/Colab Notebooks/SF_CrimeDataset.csv')
21     start2019 = "2019/01/01"
22     end2019 = "2019/12/31"
23     bef_covid_data = df[(df['Incident Date'] >= start2019) & (df['Incident Date'] < end2019)]
24
25
26     start20 = "2020/01/01"
27     end20 = "2020/12/31"
28     aft_covid_data = df[(df['Incident Date'] >= start20) & (df['Incident Date'] <= end20)]
29
30     features = ["Incident Datetime"]
31     X1 = bef_covid_data.loc[:,features].values
32     Y1 = aft_covid_data.loc[:,features].values
33
34     am_count_precovid = 0
35     pm_count_precovid = 0
36     am_count_postcovid = 0
37     pm_count_postcovid = 0
38     for i in range(len(X1)):
39         stry = str(X1[i])
40         if stry.find("PM") != -1:
41             pm_count_precovid = pm_count_precovid + 1
42         else:
43             am_count_precovid = am_count_precovid + 1

```

```

43     am_count_precovid = am_count_precovid + 1
44
45     for i in range(len(Y1)):
46         stry = str(Y1[i])
47         if stry.find("PM") != -1:
48             pm_count_postcovid = pm_count_postcovid + 1
49         else:
50             am_count_postcovid = am_count_postcovid + 1
51
52
53     matrix_covid_crime = np.zeros([2,2])
54     matrix_covid_crime[0][0] = am_count_precovid
55     matrix_covid_crime[0][1] = am_count_postcovid
56     matrix_covid_crime[1][0] = pm_count_precovid
57     matrix_covid_crime[1][1] = pm_count_postcovid
58
59     q_observed,df = chi_square(matrix_covid_crime)
60     print(q_observed,df)
61
62 initial()

```

90645.90707603919 1

In chi-square test we will check whether the Covid-19 affects the crimes taking place in the AM and PM timings.

In this we check whether the 2 sets (**X = Crimes Count in AM/PM, Y = Covid19 Dataset**) are independent or not.

### **Our null hypothesis is that X is dependent on Y**

Ideally, the X has to be dependent on Y because due the Covid 19 there were a lot of restrictions on movement of people and hence it should have reduced the number of crimes that takes place either during AM or PM.

We calculate the chi-square value to find this.

If p-value > alpha, we will reject the null hypothesis, assuming alpha to be 0.05

We will have 2 rows and 2 columns. The columns will be Pre-Covid and Post-Covid, while the rows are Crimes during AM time and Crimes during PM time

**Result:** Given alpha = 0.05. Since Q statistic is really large, from the table we find that the p-value will be really small. i.e p-value <<<< alpha

**So we fail to reject the null hypothesis and accept the null.** \*Hence, X is dependent on Y \* With this we can infer that the crimes during AM time and PM time both have a significant change due to Covid 19.

---

✓ 8s    completed at 16:38

● ×