

We are Accepting NULL HYPOTHESIS for 'IL confirmed' because absolute value of 2 1 sample testing(Z) is 1.342031
566493386 and it's greater than 2 value(at 0.025) = 1.96

We are rejecting NULL HYPOTHESIS for 'IN confirmed' because absolute value of 2 1 sample testing(Z) is 2.755986
312150163 and it's greater than 2 value(at 0.025) = 1.96

We are rejecting NULL HYPOTHESIS for 'IL deaths' because absolute value of 2 1 sample testing(Z) is 3.243840728
1961175 and it's greater than 2 value(at 0.025) = 1.96

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of 2 1 sample testing(Z) is 7.730917844
623442 and it's greater than 2 value(at 0.025) = 1.96

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=30) or dataset is normal : Satisfied because n >= 30

Hence the test is applicable

T 1 sample testing for mean of Mar '21

```
In [87]: def T_test_single(feature, sample_mean_feb, sample_variance_mar):
    lambda_mar = mar_month[feature].mean()
    lambda_mar = sample_mean_feb/np.sqrt(sample_variance_mar/len(mar_month))

    T = abs(T)

    if (T > 2.042):
        print("We are rejecting NULL HYPOTHESIS for " + feature + " because absolute value of T 1 sample testing(Z) is 1.2191958440379214 and it's greater than Z value(at 0.025) = 1.96")
    else:
        print("We are Accepting NULL HYPOTHESIS for " + feature + " because absolute value of T 1 sample testing(Z) is 0.6754221197921421 and it's greater than Z value(at 0.025) = 1.96")

    # print(T)
    print('\n\n')

for feature in features:
    sample_mean_feb = feb_month[feature].mean()
    sample_variance_mar = variance(mar_month[feature].values)
    # print(np.sqrt(sample_variance_mar))
    T_test_single(feature, sample_mean_feb, sample_variance_mar)
```

We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of T 1 sample testing(Z) is 2.174513
441351876 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 11.25498
3019981796 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 8.100718477
990918 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 14.65317550
39723 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 37.5838474514719 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 = 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

```
In [88]: def walds_test_double(feature):
    lambda_feb = feb_month[feature].mean()
    lambda_mar = mar_month[feature].mean()

    delta = lambda_mar - lambda_feb
    W = (delta)/(np.sqrt((lambda_feb/len(feb_month) + lambda_mar/(len(mar_month))))

    W = abs(W)

    if (W > 1.96):
        print("We are rejecting NULL HYPOTHESIS for " + feature + " because absolute value of Walds 2 sample testing(W) is 187.4665744482369 and it's greater than Z value(at 0.025) = 1.96")
    else:
        print("We are Accepting NULL HYPOTHESIS for " + feature + " because absolute value of Walds 2 sample testing(W) is 125.21234378566539 and it's greater than Z value(at 0.025) = 1.96")

    # print(W)
    print('\n\n')
```

We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of Walds 2 sample testing(W) is 27.25888884959588 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 47.077654405046 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 13.0349
3216155413 and it's greater than T(57,0.025) = 2.0025

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of Walds 2 sample testing(W) is 16.0118
696108 and it's greater than T(57,0.025) = 2.0025

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 125.21234378566539 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 31.386159361564246 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=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

```
In [89]: def T_test_double(feature, sample_variance_feb, sample_variance_mar):
    lambda_feb = feb_month[feature].mean()
    lambda_mar = mar_month[feature].mean()

    delta = lambda_mar - lambda_feb
    T = (delta)/(np.sqrt((sample_variance_feb/len(feb_month) + sample_variance_mar/(len(mar_month))))

    T = abs(T)

    if (T > 2.0025):
        print("We are rejecting NULL HYPOTHESIS for " + feature + " because absolute value of T 2 sample testing(Z) is 4.111130199927878 and it's greater than T(57,0.025) = 2.0025")
    else:
        print("We are Accepting NULL HYPOTHESIS for " + feature + " because absolute value of T 2 sample testing(Z) is 12.023210757714896 and it's greater than T(57,0.025) = 2.0025")

    # print(T)
    print('\n\n')
```

We are rejecting NULL HYPOTHESIS for 'IL confirmed' because absolute value of T 2 sample testing(Z) is 4.111130199927878 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 4.171270232
31417 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 5.888125932
31417 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 5.888125932
31417 and it's greater than T(57,0.025) = 2.0025

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 Z value(at 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 Z value(at 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 Z value(at 0.025) = 2.0025

We are rejecting NULL HYPOTHESIS for 'IN deaths' because absolute value of T 2 sample testing(Z) is 6.49063693322033 and it's greater than Z value(at 0.025) = 2.0025

Is the test applicable?

T test is applicable when n value is very small, but here n = 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)

```
In [90]: def calculate_MME(y1, y2):
    return numpy.sum((y1 - y2) * (y1 - y2)) / len(y1)

def plot_ecdf(input_list, label, color):
    input_list.sort()
    # sort input array
    n = len(input_list)

    # initialize x and y to plot CDF
    x = input_list[0]
    y = 0

    for point in input_list:
        value = y/(len(y) - 1) + 1 / n
        # update x and y values
        x = x + (point - point)
        y = y + (y/(len(y) - 1), value)

    # ecdf step function plot
    plt.plot(x, y, label=label, color=color)
    return x[1:], y[1:]

data = df.iloc[0:,0:4].to_numpy()
data = data[253:345, 0:1].astype(numpy.float64)

# NUMBER OF CASES
cases_ill = data[:, 0]
cases_states_sample_mean = numpy.mean(cases_ill)
cases_states_sample_variance = numpy.var(cases_ill)
cases_states_mme_poisson = cases_states_sample_mean
cases_states_mme_geometric = 1/cases_states_sample_mean
cases_states_mme_binomial = 1/cases_states_sample_variance / cases_states_sample_mean
cases_states_mme_n_binomial = cases_states_sample_mean**2 / (cases_states_sample_mean - cases_states_sample_variance)

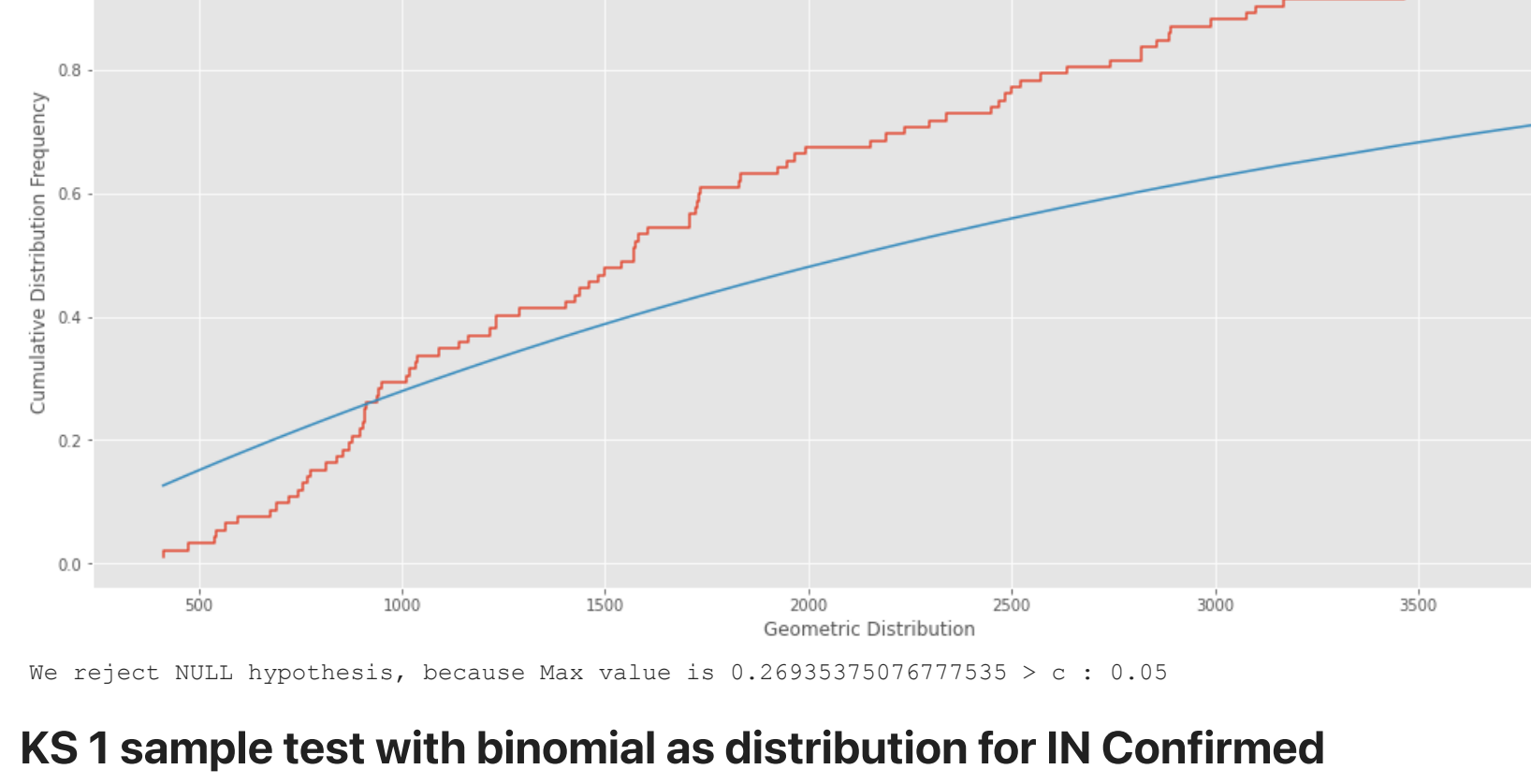
# NUMBER OF DEATHS
deaths_ill = data[:, 2]
deaths_sample_mean = numpy.mean(deaths_ill)
deaths_sample_variance = numpy.var(deaths_ill)
deaths_mme_poisson = deaths_sample_mean
deaths_mme_geometric = 1 / deaths_sample_mean
deaths_mme_p_binomial = 1 - deaths_sample_variance / deaths_sample_mean
deaths_mme_n_binomial = deaths_sample_mean**2 / (deaths_sample_mean - deaths_sample_variance)
```

KS 1 sample test with poisson as distribution for IN Confirmed

```
In [91]: # print(cases_states_mme_poisson)
cases_in = data[:, 1]
cases_in = numpy.sort(cases_in)
cdf_y = numpy.array([])
cdf = 0
n = len(cases_in)
max_diff = 0
poisson_cdf = numpy.array([])
poisson_cdf = numpy.append(poisson_cdf, poisson_point)
if max_diff < numpy.abs(cdf - poisson_point):
    max_diff = numpy.abs(cdf - poisson_point)
poisson_cdf = numpy.append(poisson_cdf, poisson_point)
cdf = 1/n
cdf_y = numpy.append(cdf_y, cdf)

plt.figure('Cases Case', figsize=(18,8))
plt.xlabel('Poisson Distribution')
plt.ylabel('Cumulative Distribution Frequency')
plt.step(cases_in, cdf_y, label = "IN")
plt.plot(cases_in, poisson_cdf, label = "IL")
plt.show()

if max_diff > 0.05:
    print("We reject NULL hypothesis, because Max value is " + str(max_diff) + " > c : 0.05")
else:
    print("We accept NULL hypothesis, because Max value is " + str(max_diff) + " <= c : 0.05")
```



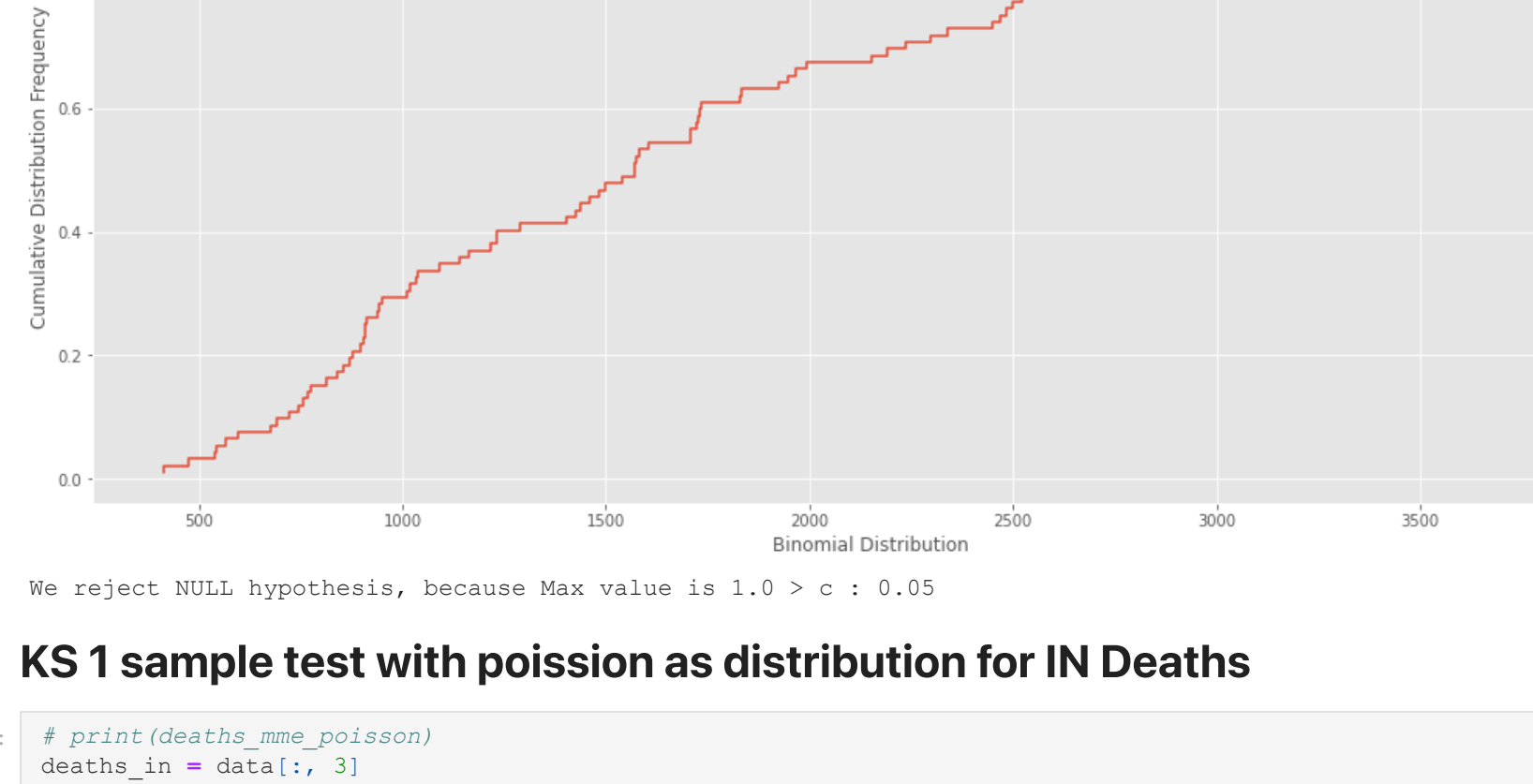
We reject NULL hypothesis, because Max value is 0.8465489533002235 > c : 0.05

KS 1 sample test with geometric as distribution for IN Confirmed

```
In [92]: # print(cases_states_mme_geometric)
cdf_y = numpy.array([])
cdf = 0
max_diff = 0
geom_cdf = numpy.array([])
geom_cdf = numpy.append(geom_cdf, geom_point)
if max_diff < numpy.abs(cdf - geom_point):
    max_diff = numpy.abs(cdf - geom_point)
geom_cdf = numpy.append(geom_cdf, geom_point)
cdf = 1 / n
cdf_y = numpy.append(cdf_y, cdf)

plt.figure('Cases Case', figsize=(18,8))
plt.xlabel('Geometric Distribution')
plt.ylabel('Cumulative Distribution Frequency')
plt.step(cases_in, cdf_y, label = "IN")
plt.plot(cases_in, geom_cdf, label = "IL")
plt.show()

if max_diff > 0.05:
    print("We reject NULL hypothesis, because Max value is " + str(max_diff) + " > c : 0.05")
else:
    print("We accept NULL hypothesis, because Max value is " + str(max_diff) + " <= c : 0.05")
```



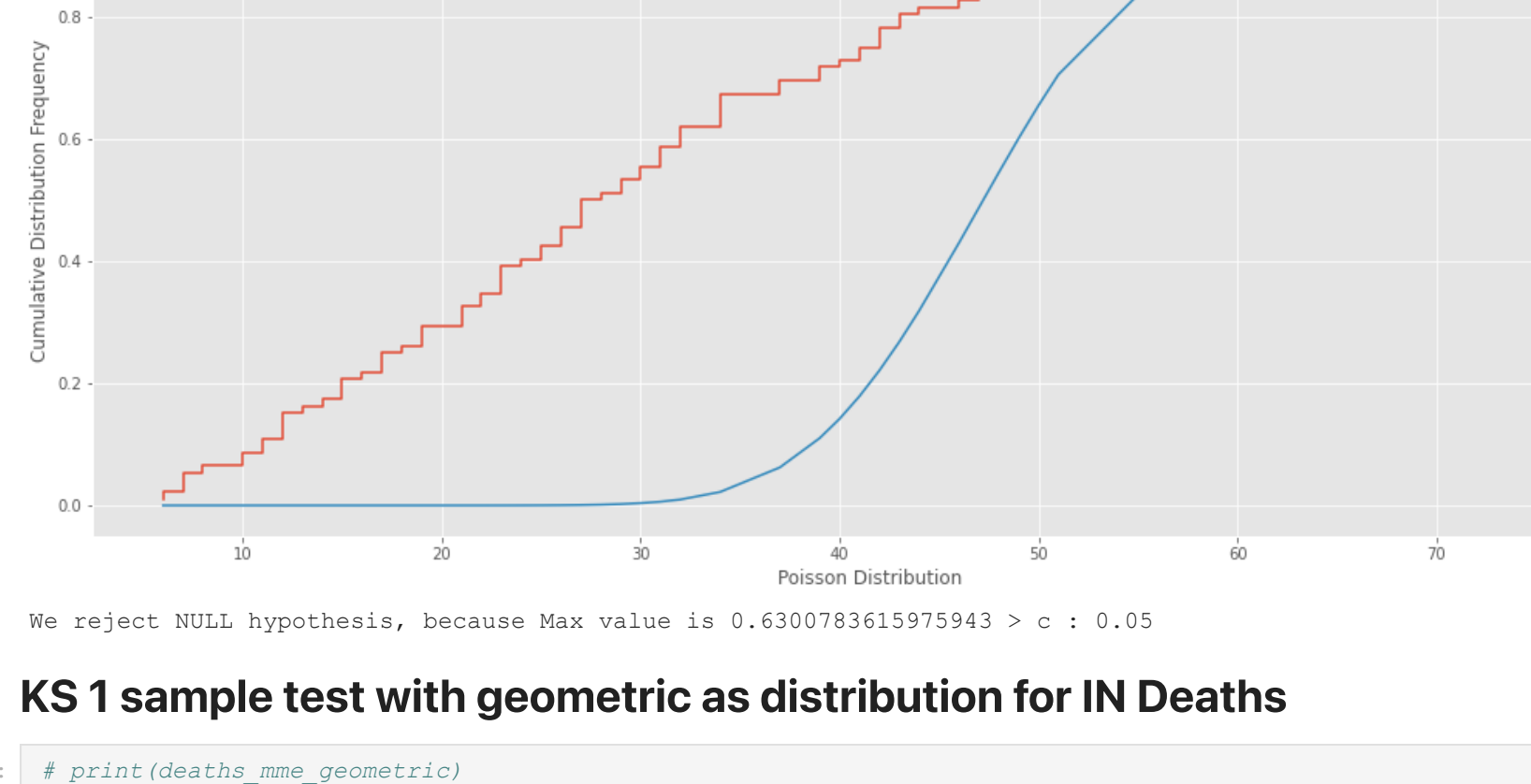
We reject NULL hypothesis, because Max value is 0.269357391507677535 > c : 0.05

KS 1 sample test with binomial as distribution for IN Confirmed

```
In [93]: cdf_y = numpy.array([])
cdf = 0
max_diff = 0
binom_cdf = numpy.array([])
for i in cases_in:
    binom_point = binom.cdf(i, cases_states_mme_n_binomial, cases_states_mme_p_binomial)
    binom_cdf = numpy.append(binom_cdf, binom_point)
    if max_diff < numpy.abs(cdf - binom_point):
        max_diff = numpy.abs(cdf - binom_point)
    binom_cdf = numpy.append(binom_cdf, binom_point)
cdf = 1 / n
cdf_y = numpy.append(cdf_y, cdf)

plt.figure('Cases Case', figsize=(18,8))
plt.xlabel('Binomial Distribution')
plt.ylabel('Cumulative Distribution Frequency')
plt.step(cases_in, cdf_y, label = "IN")
plt.plot(cases_in, binom_cdf, label = "IL")
plt.show()

if max_diff > 0.05:
    print("We reject NULL hypothesis, because Max value is " + str(max_diff) + " > c : 0.05")
else:
    print("We accept NULL hypothesis, because Max value is " + str(max_diff) + " <= c : 0.05")
```



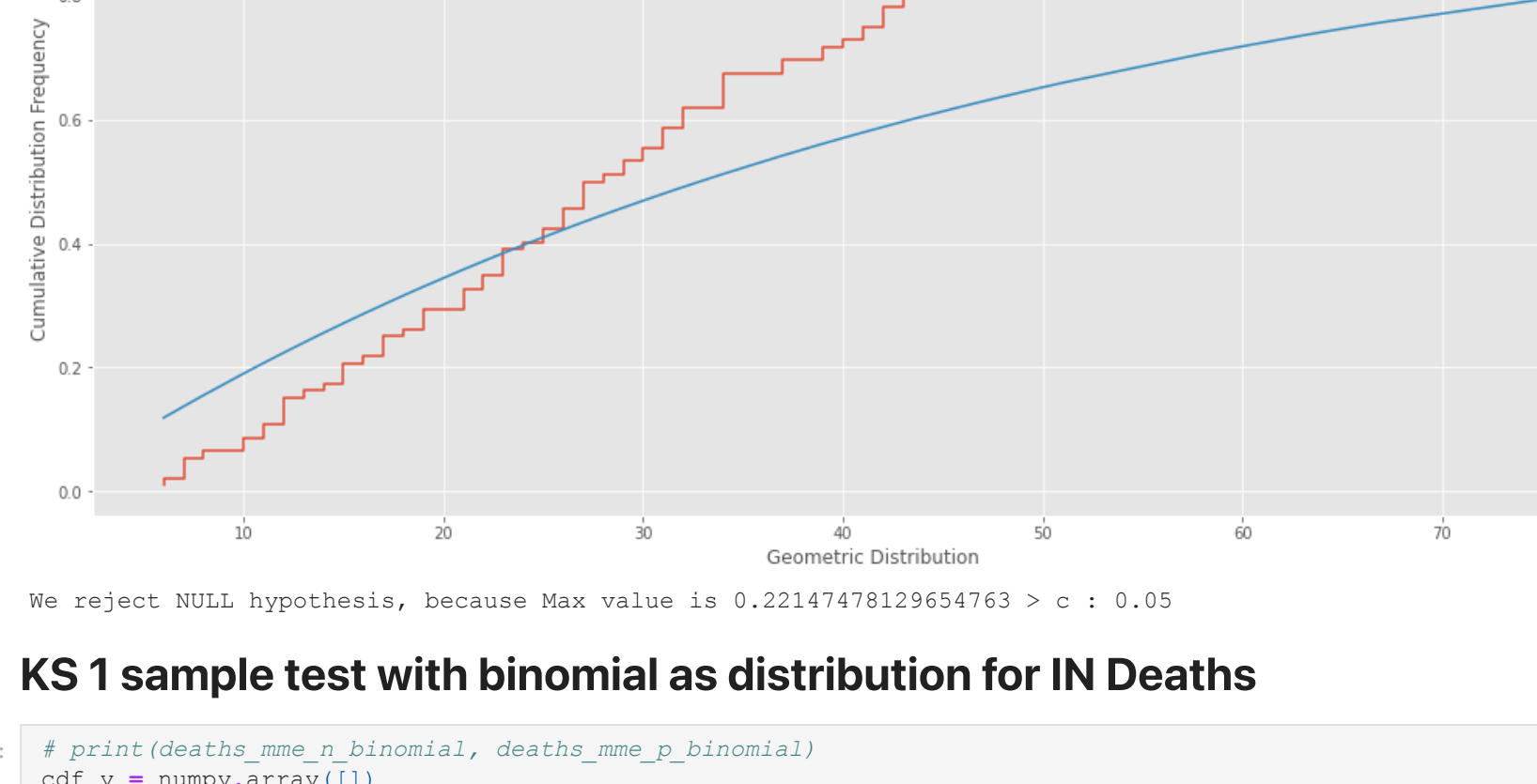
We reject NULL hypothesis, because Max value is 1.0 > c : 0.05

KS 1 sample test with poisson as distribution for IN Deaths

```
In [94]: # print(deaths_mme_poisson)
deaths_in = data[:, 3]
deaths_in = numpy.sort(deaths_in)
cdf_y = numpy.array([])
cdf = 0
n = len(deaths_in)
max_diff = 0
poisson_cdf = numpy.array([])
poisson_cdf = numpy.append(poisson_cdf, poisson_point)
if max_diff < numpy.abs(cdf - poisson_point):
    max_diff = numpy.abs(cdf - poisson_point)
poisson_cdf = numpy.append(poisson_cdf, poisson_point)
cdf = 1/n
cdf_y = numpy.append(cdf_y, cdf)

plt.figure('Deaths Case', figsize=(18,8))
plt.xlabel('Poisson Distribution')
plt.ylabel('Cumulative Distribution Frequency')
plt.step(deaths_in, cdf_y, label = "IN")
plt.plot(deaths_in, poisson_cdf, label = "IL")
plt.show()

if max_diff > 0.05:
    print("We reject NULL hypothesis, because Max value is " + str(max_diff) + " > c : 0.05")
else:
    print("We accept NULL hypothesis, because Max value is " + str(max_diff) + " <= c : 0.05")
```



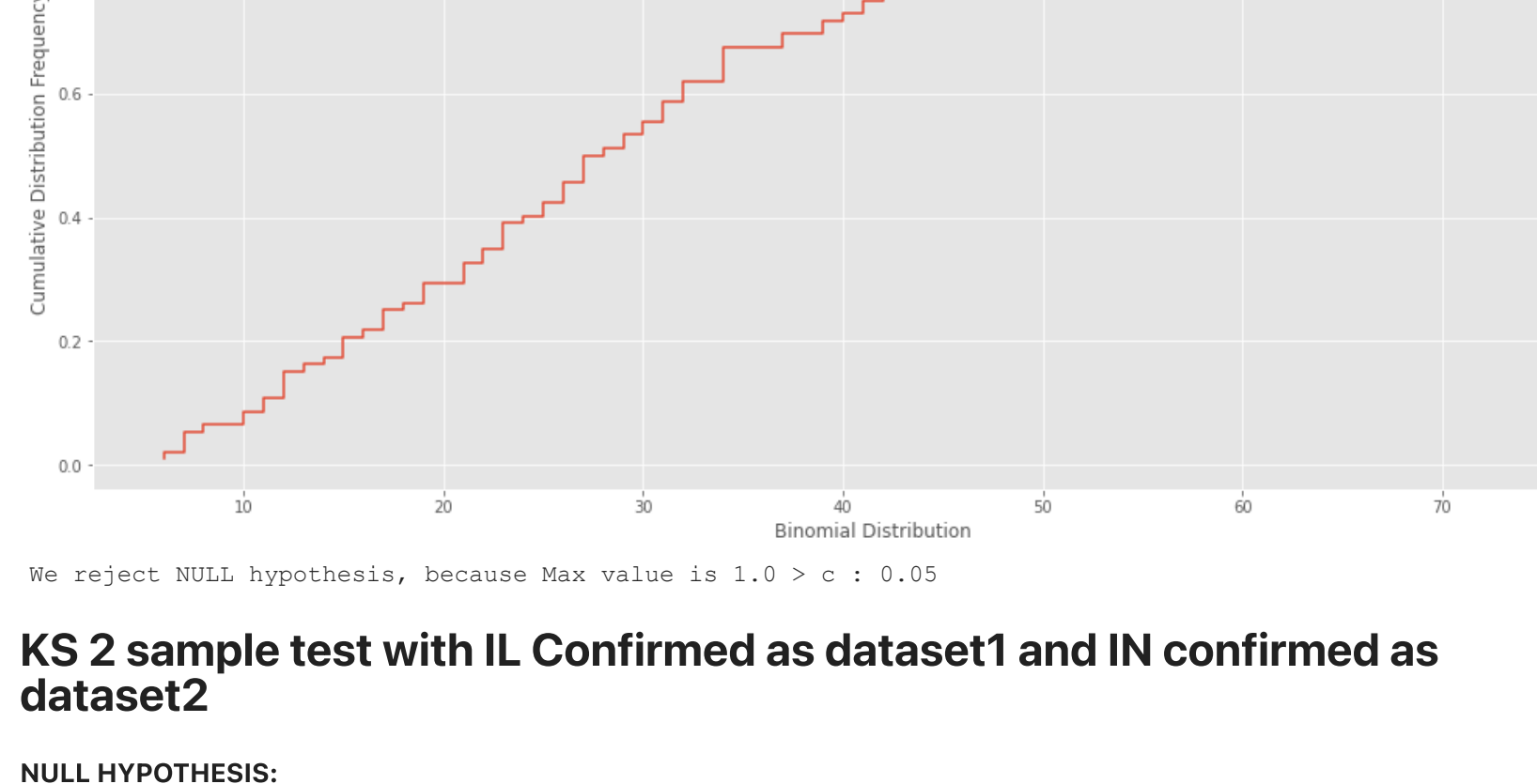
We reject NULL hypothesis, because Max value is 0.6300783619759343 > c : 0.05

KS 1 sample test with geometric as distribution for IN Deaths

```
In [95]: # print(deaths_mme_geometric)
cdf_y = numpy.array([])
cdf = 0
max_diff = 0
geom_cdf = numpy.array([])
geom_cdf = numpy.append(geom_cdf, geom_point)
if max_diff < numpy.abs(cdf - geom_point):
    max_diff = numpy.abs(cdf - geom_point)
geom_cdf = numpy.append(geom_cdf, geom_point)
cdf = 1 / n
cdf_y = numpy.append(cdf_y, cdf)

plt.figure('Deaths Case', figsize=(18,8))
plt.xlabel('Geometric Distribution')
plt.ylabel('Cumulative Distribution Frequency')
plt.step(deaths_in, cdf_y, label = "IN")
plt.plot(deaths_in, geom_cdf, label = "IL")
plt.show()

if max_diff > 0.05:
    print("We reject NULL hypothesis, because Max value is " + str(max_diff) + " > c : 0.05")
else:
    print("We accept NULL hypothesis, because Max value is " + str(max_diff) + " <= c : 0.05")
```



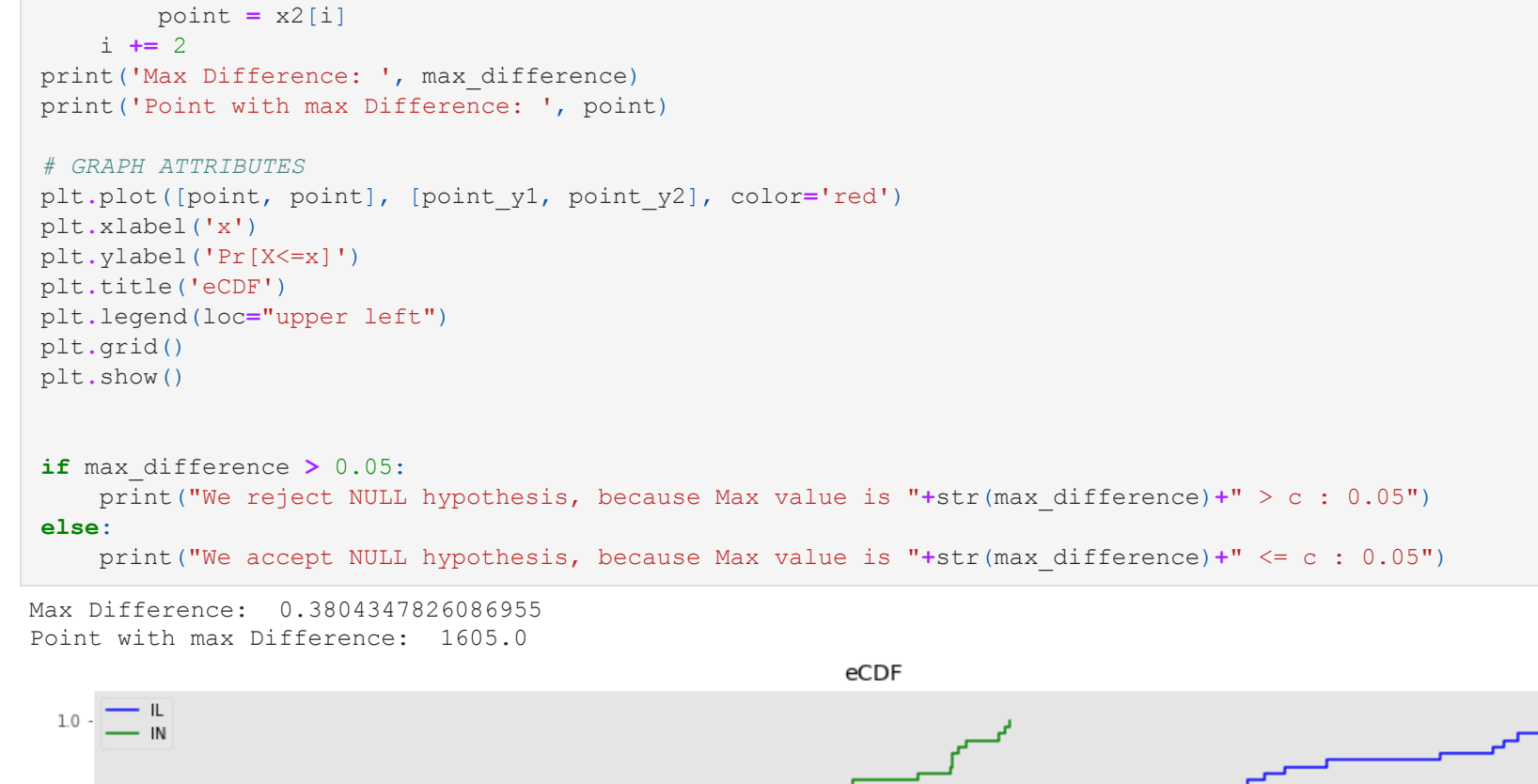
We reject NULL hypothesis, because Max value is 0.22147478129654763 > c : 0.05

KS 1 sample test with binomial as distribution for IN Deaths

```
In [96]: # print(deaths_mme_n_binomial, deaths_mme_p_binomial)
cdf_y = numpy.array([])
cdf = 0
max_diff = 0
binom_cdf = numpy.array([])
for i in deaths_in:
    binom_point = binom.cdf(i, deaths_mme_n_binomial, deaths_mme_p_binomial)
    binom_cdf = numpy.append(binom_cdf, binom_point)
    if max_diff < numpy.abs(cdf - binom_point):
        max_diff = numpy.abs(cdf - binom_point)
    binom_cdf = numpy.append(binom_cdf, binom_point)
cdf = 1 / n
cdf_y = numpy.append(cdf_y, cdf)

plt.figure('Deaths Case', figsize=(18,8))
plt.xlabel('Binomial Distribution')
plt.ylabel('Cumulative Distribution Frequency')
plt.step(deaths_in, binom_cdf, label = "IN")
plt.plot(deaths_in, binom_cdf, label = "IL")
plt.show()

if max_diff > 0.05:
    print("We reject NULL hypothesis, because Max value is " + str(max_diff) + " > c : 0.05")
else:
    print("We accept NULL hypothesis, because Max value is " + str(max_diff) + " <= c : 0.05")
```



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

```
In [97]: # TWO POPULATION TEST
X = data[:, 2]
Y = data[:, 3]

# PLOT THE GRAPHS
plt.figure('ECDF', figsize=(18,8))
x1, y1 = plot_ecdf(X, 'IN', color='blue')
x2, y2 = plot_ecdf(Y, 'IN', color='green')

max_difference = 0
point = 0
point_y2 = 0
i = 0
j = 0
while i < len(x2):
    y2_left, y2_right = y2[i], y2[i + 1]
    while j + 2 < len(x1) and x1[j + 2] < x2[i]:
        j += 2
    if x2[i] == x1[j]:
        y1_left, y1_right = y1[j], y1[j + 1]
    else:
        y1_left, y1_right = y1[j + 1], y1[j + 1]
    max_difference = numpy.max([max_difference, numpy.absolute(y1_left - y2_left), numpy.absolute(y1_right - y2_right)])
    point_y2 = y2_left
    point = x2[i]
    i += 1
    j += 1
print('Max Difference: ', max_difference)
print('Point with max Difference: ', point)
```



We reject NULL hypothesis, because Max value is 0.380437826086955 > c : 0.05

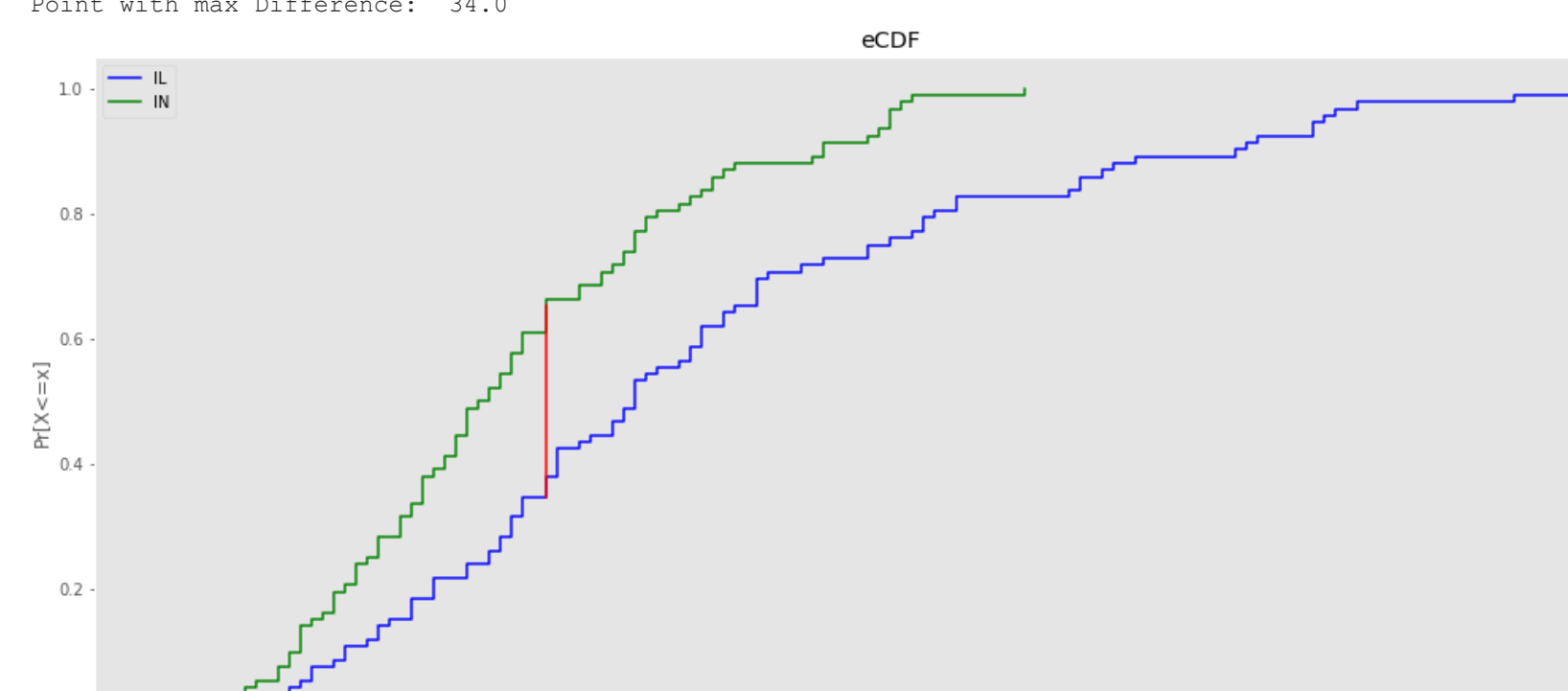
KS 2 sample test with IL Deaths as dataset1 and IN Deaths as dataset2

NULL HYPOTHESIS:
Distribution of dataset1 = Distribution of dataset2

```
In [98]: # TWO POPULATION TEST
X = data[:, 2]
Y = data[:, 3]

# PLOT THE GRAPHS
plt.figure('ECDF', figsize=(18,8))
x1, y1 = plot_ecdf(X, 'IN', color='blue')
x2, y2 = plot_ecdf(Y, 'IN', color='green')

max_difference = 0
point = 0
point_y2 = 0
i = 0
j = 0
while i < len(x2):
    y2_left, y2_right = y2[i], y2[i + 1]
    while j + 2 < len(x1) and x1[j + 2] < x2[i]:
        j += 2
    if x2[i] == x1[j]:
        y1_left, y1_right = y1[j], y1[j + 1]
    else:
        y1_left, y1_right = y1[j + 1], y1[j + 1]
    max_difference = numpy.max([max_difference, numpy.absolute(y1_left - y2_left), numpy.absolute(y1_right - y2_right)])
    point_y2 = y2_left
    point = x2[i]
    i += 1
    j += 1
print('Max Difference: ', max_difference)
print('Point with max Difference: ', point)
```



We reject NULL hypothesis, because Max value is 0.1512373913043484 > c : 0.05

Permutation test

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

NULL HYPOTHESIS:
Distribution of dataset1 = Distribution of dataset2


```
data = df.iloc[0:4,0:4].to_numpy()
data = data[233:345, 0:1].astype(numpy.float64)
```

```
# CASES
x1 = data[:, 0]
y1 = data[:, 1]

X_avg = numpy.mean(x1)
f_avg = numpy.mean(y1)
t_obs = numpy.absolute(X_avg - Y_avg)

number = 0
combined = numpy.append(x1, y1, axis=0)
for i in range(1000):
    permutation = numpy.random.permutation(combined)
    X_permutation = permutation[: len(x1)]
    Y_permutation = permutation[len(x2):]
    t_predict = numpy.absolute(numpy.mean(X_permutation) - numpy.mean(Y_permutation))
    if t_predict > t_obs:
        number += 1

p_value = number / 1000
print(p_value)
if p_value > 0.05:
    print('Reject Null Hypothesis as p value is '+str(p_value)+' which is greater than 0.5')
else:
    print('Accept Null Hypothesis as p value is '+str(p_value)+' which is greater than 0.5')
```

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

```
# DEATHS
x2 = data[:, 2]
y2 = data[:, 3]

X_avg = numpy.mean(x2)
Y_avg = numpy.mean(y2)
t_obs = numpy.absolute(X_avg - Y_avg)

number = 0
combined = numpy.append(x2, y2, axis=0)
for i in range(1000):
    permutation = numpy.random.permutation(combined)
    X_permutation = permutation[: len(x2)]
    Y_permutation = permutation[len(x2):]
    t_predict = numpy.absolute(numpy.mean(X_permutation) - numpy.mean(Y_permutation))
    if t_predict > t_obs:
        number += 1

p_value = number / 1000
print(p_value)
if p_value > 0.05:
    print('Reject Null Hypothesis as p value is '+str(p_value)+' which is greater than 0.5')
else:
    print('Accept Null Hypothesis as p value is '+str(p_value)+' which is greater than 0.5')
```

0.0
Accept Null Hypothesis as p value is 0.0 which is greater than 0.5

PART 2d : Prediction using Bayesian Inference

```
INdeaths = df[df['Date'] > '2020-05-31']['IN deaths']
ILdeaths = df[df['Date'] > '2020-05-31']['IL deaths']

# print(sum(ILdeaths[:35])+ILdeaths[:35]))

INandILdeaths=list(INdeaths+ILdeaths)

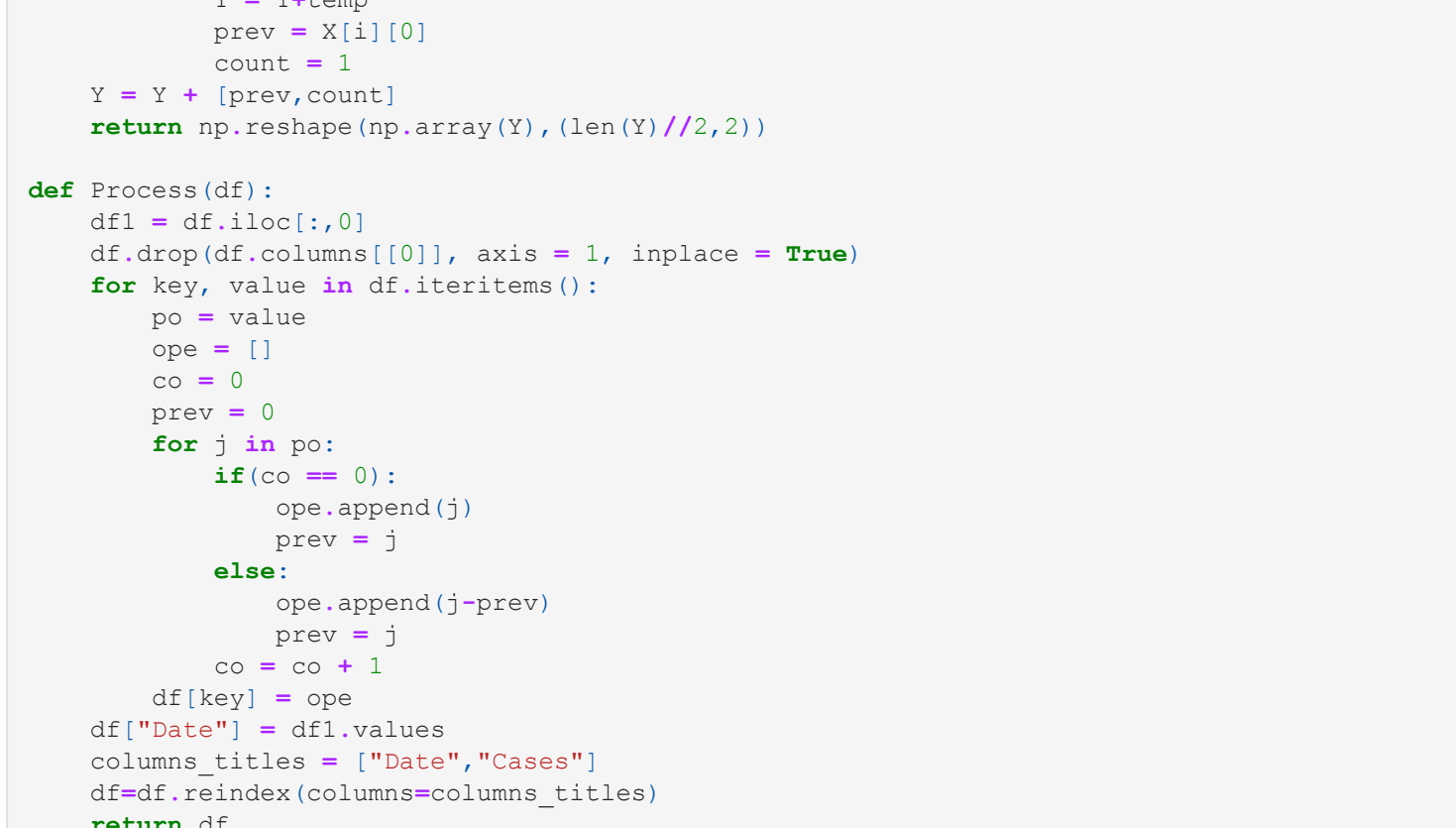
# print(INandILdeaths)

def plot_gamma_distribution(alpha=1, beta=1, label=""):
    plt.plot(gamma_distribution(alpha, beta=1, label=label))
    plt.title("Posterior gamma distributions")
    label = label + " MAP: " + str(round((alpha/beta),3))
    plt.xlabel("Deaths")
    plt.ylabel("Gamma distribution")
    plt.plot(x, gamma.pdf(x, alpha, scale=1/beta), label=label)
    plt.legend()

beta=np.mean(INandILdeaths[:28])
# print(beta)
# print(sum(INandILdeaths[:35]))
# print(sum(INandILdeaths[:42]))
plt.figure(figsize=(12,8))

for i in range(4,8):
    alpha = sum(INandILdeaths[:i+1]*7)+1
    # print(beta)
    scale = (i+1)*7 + 1/beta
    plot_gamma_distribution(alpha, scale, 'Week-'+str(i+1))

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 SanFrancisco crime and California Covid data

```
def CorrelationCoeff(X,Y):
    X = X[57:385,1]
    Y = Y[57:385,1]
    X = X.astype(np.int)
    Y = Y.astype(np.int)
    return np.corrcoef(X, Y)
```

```
def Preprocessdata( X ) :
    l = len(X)
    prev = X[0][0]
    count = 1
    Y = []
    for i in range(0,l):
        if(X[i][0]!=prev):
            count=count+1
            temp = [prev,int(count)]
            Y = Y+temp
            prev = X[i][0]
            count = 1
            Y = Y + [prev,count]
    return np.reshape(np.array(Y), (len(Y)//2,2))

def Process(df):
    df1 = df.iloc[:,0]
    df.drop(df.columns[[0]], axis = 1, inplace = True)
    for key, value in df.iteritems():
        po = value
        ope = []
        co = 0
        prev = 0
        for i in po:
            if(co == 0):
                ope.append(i)
                prev = j
            else:
                ope.append(j-prev)
                prev = j
                co = co + 1
        df[key] = ope
    df['Date'] = df1.values
    columns_titles = ["Date","Cases"]
    df=df.reindex(columns=columns_titles)
    return df
```

```
def CorrelationCoeffWholeYear(X,Y):
    X = X[:,1]
    Y = Y[:,1]
    X = X.astype(np.int)
    Y = Y.astype(np.int)
    return np.corrcoef(X, Y)

cols_list = ['Incident Date']
df = pd.read_csv('gdrive/MyDrive/Colab Notebooks/SF_CrimeDataset.csv', usecols=cols_list)
Covid_data = pd.read_csv('gdrive/MyDrive/Colab Notebooks/US_confirmed.csv')
Covid_data.head()

CA_Covid_data = Covid_data.iloc[4].to_frame().reset_index().iloc[1:]
CA_Covid_data.columns = ['Date', 'Cases']
CA_Covid_data['Date'] = pd.to_datetime(CA_Covid_data['Date'])
CA_Covid_data = CA_Covid_data[(CA_Covid_data['Date'] >= '2020/01/01') & (CA_Covid_data['Date'] <= '2020/12/31')]
CA_Covid_data = Process(CA_Covid_data)

start2019 = "2019/01/01"
end2019 = "2019/12/31"
bef_covid_data = df[(df['Incident Date'] >= start2019) & (df['Incident Date'] <= end2019)]

start20 = "2020/01/01"
end20 = "2020/12/31"
aft_covid_data = df[(df['Incident Date'] >= start20) & (df['Incident Date'] <= end20) & (df['Incident Date'] != CA_Covid_data)]

aft_covid_data = np.sort(bef_covid_data.to_numpy(), axis = 0)

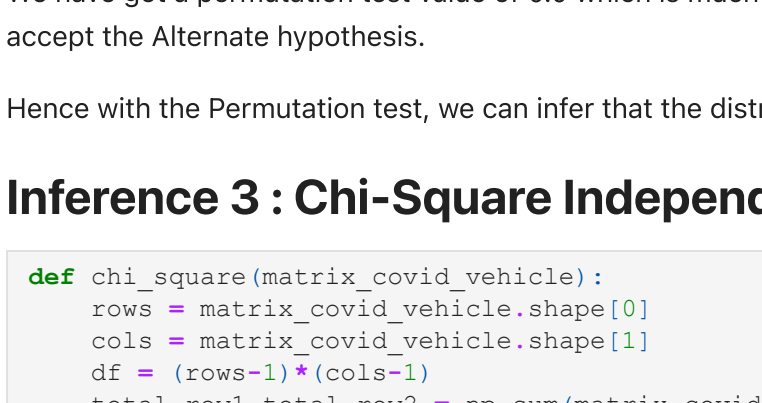
bef_covid_data = Preprocessdata(bef_covid_data)

aft_covid_data = np.sort(aft_covid_data.to_numpy(), axis = 0)
aft_covid_data = Preprocessdata(aft_covid_data)
```

```
print("Correlation Coefficient considering the whole Year 2020 is : ", CorrelationCoeffWholeYear(aft_covid_data, aft_covid_data))

X = []
for i in range(365):
    plt.figure('Correlation of Covid and Number of Crimes')
    plt.plot(X, aft_covid_data[:,1].astype(int), label='Number of Crimes')
    plt.xlabel('Days')
    plt.ylabel('No. of Crimes')
    plt.title('Distribution of Crimes Over the year 2020')
    plt.legend(loc='upper left')
    plt.grid()
```

Correlation Coefficient considering the whole Year 2020 is : -0.031275034531885064



Correlation Coefficient considering the whole Year 2020 is : -0.7684999251821152

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 Sanfransisco crime before covid and after covid is same or not bold text

```
def PermutationTest(X, Y, rounds):
    X = X.astype(np.int)
    Y = Y.astype(np.int)
    count = 0
    diff = abs(np.mean(X) - np.mean(Y))
    tot = X+Y
    n = len(tot)
    for i in range(0,rounds):
        ampp.random.permutation(tot)[:n//2]
        bmpp.random.permutation(tot)[n//2:]
        if (abs(a.mean() - b.mean())>diff):
            count+=1
    return (count/rounds)

p_value = PermutationTest", PermutationTest(bef_covid_data[:,1],aft_covid_data[:,1],1000 ))
```

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

```
def chi_square(matrix_covid_vehicle):
    rows = matrix_covid_vehicle.shape[0]
    cols = matrix_covid_vehicle.shape[1]
    df = (rows-1)*(cols-1)
    total_row, total_row2 = np.sum(matrix_covid_vehicle,axis=0)
    total_col, total_col2 = np.sum(matrix_covid_vehicle,axis=1)
    total = total_row+total_row2
    expected_values = np.zeros((2,2))
    expected_values[0][0] = (float(total_col1)*total_row1)/(total)
    expected_values[0][1] = (float(total_col2)*total_row1)/(total)
    expected_values[1][0] = (float(total_col1)*total_row2)/(total)
    expected_values[1][1] = (float(total_col2)*total_row2)/(total)
    q_expected = 0.0
    for i in range(rows):
        for j in range(cols):
            q_expected = q_expected + ((expected_values[i][j] - matrix_covid_vehicle[i][j])**2)/float(expected_values[i][j])
    return (q_expected,df)
```

```
def Initial():
    df = pd.read_csv('gdrive/MyDrive/Colab Notebooks/SF_CrimeDataset.csv')
    start2019 = "2019/01/01"
    end2019 = "2019/12/31"
    bef_covid_data = df[(df['Incident Date'] >= start2019) & (df['Incident Date'] <= end2019)]

    start20 = "2020/01/01"
    end20 = "2020/12/31"
    aft_covid_data = df[(df['Incident Date'] >= start20) & (df['Incident Date'] <= end20) & (df['Incident Date'] != CA_Covid_data)]

    features = ["Incident Date"]
    X1 = bef_covid_data.loc[:,features].values
    Y1 = aft_covid_data.loc[:,features].values

    am_count_preCovid = 0
    pm_count_preCovid = 0
    am_count_postCovid = 0
    pm_count_postCovid = 0
    for i in range(len(X1)):
        str = str(X1[i])
        if str.find("PM") != -1:
            pm_count_preCovid = pm_count_preCovid + 1
            am_count_preCovid = am_count_preCovid + 1
        else:
            am_count_preCovid = am_count_preCovid + 1

    for i in range(len(Y1)):
        str = str(Y1[i])
        if str.find("PM") != -1:
            pm_count_postCovid = pm_count_postCovid + 1
        else:
            am_count_postCovid = am_count_postCovid + 1

    matrix_covid_crime = np.zeros((2,2))
    matrix_covid_crime[0][0] = am_count_preCovid
    matrix_covid_crime[0][1] = am_count_postCovid
    matrix_covid_crime[1][0] = pm_count_preCovid
    matrix_covid_crime[1][1] = pm_count_postCovid

    q_observed,df = chi_square(matrix_covid_crime)
    print(q_observed,df)

Initial()

90645.90707803919 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.