

Name : Shreyans Tatiya

Roll No. : 16010123325

Batch : E-2

Branch : Computer Engineering

## Tutorial 4

### Correlation, Regression and Probability Distribution

#### Q1

For the following data set {(25,70), (28,80), (32,85), (36,75), (38,59), (40,65), (39,78), (42,50), (41,54), (45,66)} (i) Draw the scatter diagram  
(ii) Find the correlation coefficients (iii) Find both the regression lines (iv) Plot both regression lines together (v) Find the error for both regression lines

```
In [1]: # i)
import numpy as np
import matplotlib.pyplot as plt

# Data
data = [(25,70), (28,80), (32,85), (36,75), (38,59), (40,65), (39,78), (42,50), (41,54), (45,66)]
x = np.array([point[0] for point in data]) # Independent variable
y = np.array([point[1] for point in data]) # Dependent variable

plt.figure(figsize=(8, 6))
plt.scatter(x, y, color='blue', label='Data Points')
```

Out[1]: <matplotlib.collections.PathCollection at 0x118d3cf40>



```
In [4]: # ii)
mean_x = np.mean(x)
mean_y = np.mean(y)

# Calc correlation coefficient
corr_cf = np.corrcoef(x, y)[0,1]
print("Correlation Coefficient:", corr_cf)
```

Correlation Coefficient: -0.5764311756246668

```
In [6]: # iii)
# Regression coefficients for Y on X
b_yx = np.sum((x - mean_x) * (y - mean_y)) / np.sum((x - mean_x)**2)
c_yx = mean_y - b_yx * mean_x

# Regression coefficients for X on Y
b_xy = np.sum((x - mean_x) * (y - mean_y)) / np.sum((y - mean_y)**2)
c_xy = mean_x - b_xy * mean_y

print(f"Regression line Y on X: y = {b_yx:.2f}x + {c_yx:.2f}")
print(f"Regression line X on Y: x = {b_xy:.2f}y + {c_xy:.2f}")

Regression line Y on X: y = -1.04x + 106.27
Regression line X on Y: x = -0.32y + 58.39
```

```
In [8]: # iv)
# Generate points for regression lines
x_vals = np.linspace(min(x), max(x), 100)

y_vals_yx = b_yx * x_vals + c_yx # Y on X
y_vals_xy = (x_vals - c_xy) / b_xy # X on Y

# Plotting
def plot_regression():
    plt.figure(figsize=(8, 6))

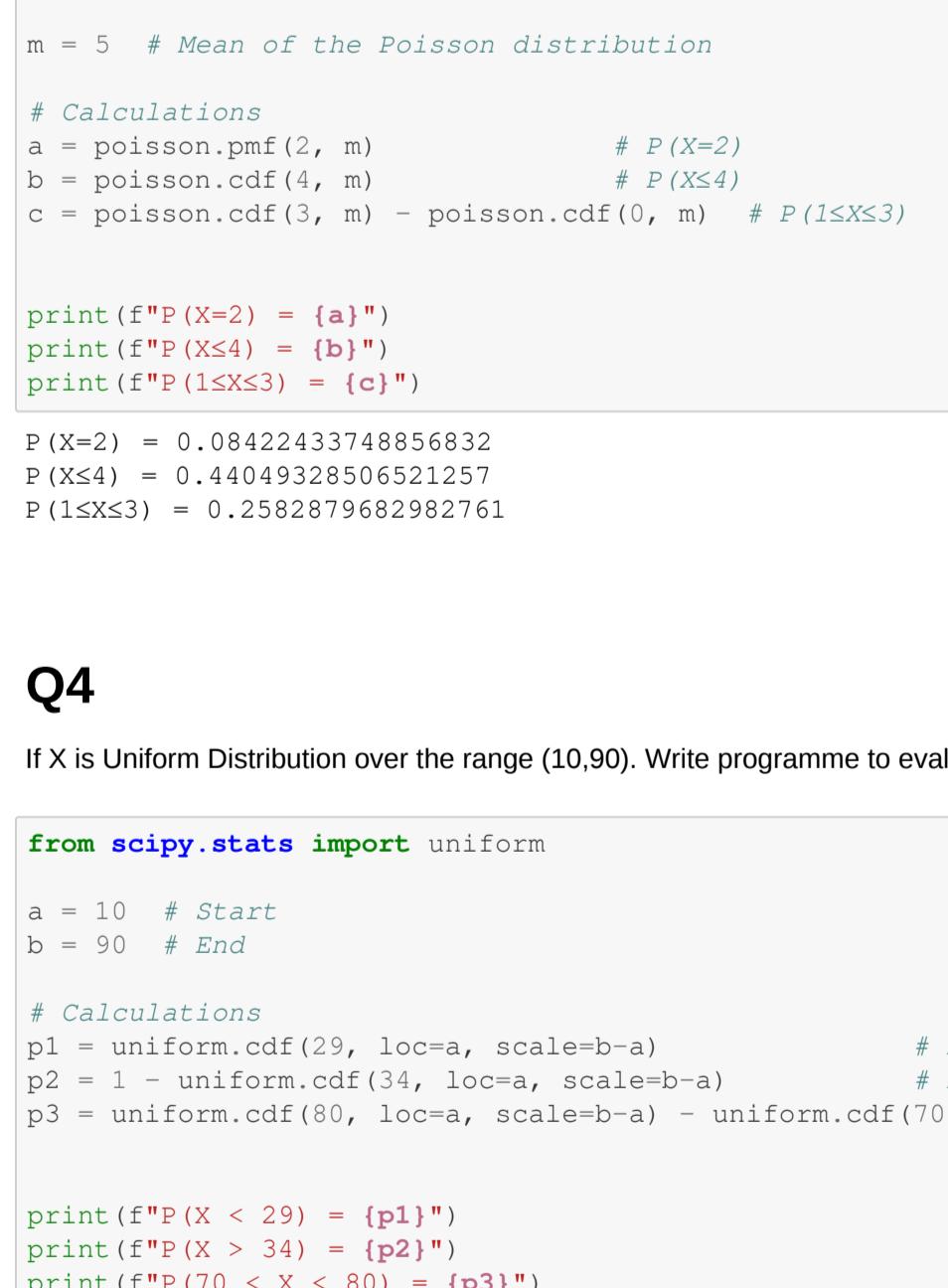
    # Scatter plot of data points
    plt.scatter(x, y, color='blue', label='Data Points')

    # Y on X
    plt.plot(x_vals, y_vals_yx, color='red', label='Y on X: y = {:.2f}x + {:.2f}'.format(b_yx, c_yx))

    # X on Y
    plt.plot(x_vals, y_vals_xy, color='green', label='X on Y: x = {:.2f}y + {:.2f}'.format(b_xy, c_xy))

    plt.xlabel('X')
    plt.ylabel('Y')
    plt.title('Scatter Plot with Regression Lines')
    plt.legend()
    plt.grid()
    plt.show()

plot_regression()
```



```
In [9]: # v) MSS for each value
import pandas as pd
n = len(data)

# For Y on X
pred_y = b_yx * x + c_yx
err_yx = (y - pred_y)**2
mss_yx_values = err_yx / n

# For X on Y
pred_x = b_xy * y + c_xy
err_xy = (x - pred_x)**2
mss_xy_values = err_xy / n

# Create a table showing MSS for each data point
data_table = {
    "X": x,
    "Y": y,
    "Predicted Y (Y on X)": pred_y,
    "Error (Y on X)^2": err_yx,
    "MSS (Y on X)": mss_yx_values,
    "Predicted X (X on Y)": pred_x,
    "Error (X on Y)^2": err_xy,
    "MSS (X on Y)": mss_xy_values,
}
```

# Convert to DataFrame for display

result\_df = pd.DataFrame(data\_table)

# Print results table

print(result\_df)

Total\_MSS\_y\_on\_x = np.sum(mss\_yx\_values)

Total\_MSS\_x\_on\_y = np.sum(mss\_xy\_values)

print("Total MSS (Y on X) :", Total\_MSS\_y\_on\_x)

print("Total MSS (X on Y) :", Total\_MSS\_x\_on\_y)

X	Y	Predicted Y (Y on X)	Error (Y on X)^2	MSS (Y on X)
0	25	70	80.260615	10.539107
1	28	80	77.145494	8.148204
2	32	85	72.984799	144.365052
3	36	75	68.824104	38.141689
4	38	59	66.427357	59.965769
5	40	65	64.663409	0.113293
6	39	78	65.201870	15.120187
7	42	50	62.583062	158.333447
8	41	54	63.623236	92.606664
9	45	66	59.462541	42.738374

Total MSS (Y on X) : 80.10054288816502

Total MSS (X on Y) : 24.599066355451814

#### Q2

If X is Binomial Distribution B(n,p) where n=15 p=0.45 Write program to evaluate and print (i) P(X=10) (ii) P(X≤12) (iii) P(X≥9)

```
In [13]: from scipy.stats import binom

n = 15
p = 0.45

a = binom.pmf(10, n, p) # P(X=10)
b = binom.cdf(12, n, p) # P(X≤12)
c = 1 - binom.cdf(8, n, p) # P(X≥9)

# Print results
print(f"P(X=10) = {a}")
print(f"P(X≤12) = {b}")
print(f"P(X≥9) = {c}")
```

P(X=10) = 0.05146285992538625

P(X≤12) = 0.998892975853391

P(X≥9) = 0.18176046503835142

#### Q3

If X is Poisson Distribution with mean 5 Write program to evaluate and print (i) P(X=2) (ii) P(X≤4) (iii) P(1≤X≤3)

```
In [17]: from scipy.stats import poisson

m = 5 # Mean of the Poisson distribution

# Calculations
a = poisson.pmf(2, m) # P(X=2)
b = poisson.cdf(4, m) # P(X≤4)
c = poisson.cdf(3, m) - poisson.cdf(0, m) # P(1≤X≤3)

# Print results
print(f"P(X=2) = {a}")
print(f"P(X≤4) = {b}")
print(f"P(1≤X≤3) = {c}")
```

P(X=2) = 0.08422433748856832

P(X≤4) = 0.44049328506521257

P(1≤X≤3) = 0.2582879682982761

#### Q4

If X is Uniform Distribution over the range (10,90). Write programme to evaluate and print (i) P(X<29) (ii) P(X>34) (iii) P(70<X<80)

```
In [18]: from scipy.stats import uniform

a = 10 # Start
b = 90 # End

# Calculations
p1 = uniform.cdf(29, loc=a, scale=b-a) # P(X < 29)
p2 = 1 - uniform.cdf(34, loc=a, scale=b-a) # P(X > 34)
p3 = uniform.cdf(80, loc=a, scale=b-a) - uniform.cdf(70, loc=a, scale=b-a) # P(70 < X < 80)

# Print results
print(f"P(X < 29) = {p1}")
print(f"P(X > 34) = {p2}")
print(f"P(70 < X < 80) = {p3}")
```

P(X < 29) = 0.2375

P(X > 34) = 0.7

P(70 < X < 80) = 0.125

#### Q5

If X is Exponential Distribution with mean 20. Write programme to evaluate and print (i) P(X<10) (ii) P(X>7) (iii) P(11<X<16). Find value of k such that P(X<k) = 0.6.

```
In [19]: from scipy.stats import expon

mean = 20
lambda_p = 1 / mean # Rate parameter

# Calculations
a = expon.cdf(10, scale=1/lambda_p) # P(X < 10)
b = 1 - expon.cdf(7, scale=1/lambda_p) # P(X > 7)
c = expon.cdf(16, scale=1/lambda_p) - expon.cdf(11, scale=1/lambda_p) # P(11 < X < 16)
k = expon.ppf(0.6, scale=1/lambda_p) # k such that P(X < k) = 0.6

# Print results
print(f"P(X < 10) = {a}")
print(f"P(X > 7) = {b}")
print(f"P(11 < X < 16) = {c}")
print(f"The value of k is {k}")
```

P(X < 10) = 0.3934693402873666

P(X > 7) = 0.7046880897187133

P(11 < X < 16) = 0.1276208462632651

The value of k is 18.3258146374831

#### Q6

If X is Normal Distribution with mean 40 and standard deviation 10. Write programme to evaluate and print (i) P(X<38) (ii) P(X>55)

(iii) P(20<X<70). Find value of k1 such that P(X<k1) = 0.3. Also find k2 such that P(X>k2) = 0.8

```
In [2]: from scipy.stats import norm

mean = 40
std_dev = 10

# Calculations
a = norm.cdf(38, loc=mean, scale=std_dev) # P(X < 38)
b = 1 - norm.cdf(55, loc=mean, scale=std_dev) # P(X > 55)
c = norm.cdf(70, loc=mean, scale=std_dev) - norm.cdf(20, loc=mean, scale=std_dev) # P(20 < X < 70)

# Print results
print(f"P(X < 38) = {a}")
print(f"P(X > 55) = {b}")
print(f"P(20 < X < 70) = {c}")
print(f"Value of k1 such that P(X < k1) = 0.3 is {k1}")
print(f"Value of k2 such that P(X > k2) = 0.8 is {k2}")
```

P(X < 38) = 0.42074056089696

P(X > 55) = 0.0668072012688509

P(20 < X < 70) = 0.97589970021907

Value of k1 such that P(X < k1) = 0.3 is 34.559594872919594

Value of k2 such that P(X > k2) = 0.8 is 31.583787664270858