

Assignment -3

Farhan Alam

February 2025

Question 1

Let X and Y have a bivariate normal distribution with parameters:

$$\mu_X = 3, \quad \mu_Y = 1, \quad \sigma_X^2 = 16, \quad \sigma_Y^2 = 25, \quad \rho_{XY} = \frac{3}{5}$$

where μ_X and μ_Y represent the means of X and Y , σ_X^2 and σ_Y^2 represent the variances of X and Y , and ρ_{XY} represents the correlation coefficient between X and Y .

Determine the following probabilities:

1. $P(3 < Y < 8)$.
2. $P(3 < Y < 8 \mid X = 7)$.
3. $P(-3 < X < 3)$.
4. $P(-3 < X < 3 \mid Y = -4)$.

1 Introduction

Bivariate normal distributions are widely used in statistical modeling to describe the relationship between two continuous random variables. In this report, we analyze a given bivariate normal distribution and compute specific probability values associated with it.

2 Data

We consider two normally distributed random variables, X and Y , with the following parameters:

$$\begin{aligned} \mu_X &= 3, & \mu_Y &= 1, \\ \sigma_X^2 &= 16, & \sigma_Y^2 &= 25, \\ \rho_{XY} &= \frac{3}{5} \end{aligned}$$

where μ_X and μ_Y represent the means, σ_X^2 and σ_Y^2 represent the variances, and ρ_{XY} is the correlation coefficient.

3 Methodology

The given problem involves computing probabilities from a bivariate normal distribution. The probability computations are based on standard normal transformations:

1. Convert the given limits to standard normal form using $Z = \frac{X-\mu}{\sigma}$.
2. Use the cumulative distribution function (CDF) of the normal distribution to evaluate probabilities.
3. For conditional probabilities, apply properties of conditional normal distributions:

$$Y|X = x \sim \mathcal{N}\left(\mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(x - \mu_X), \sigma_Y^2(1 - \rho^2)\right).$$

4 Results

We determine the following probabilities:

1. $P(3 < Y < 8)$
2. $P(3 < Y < 8|X = 7)$
3. $P(-3 < X < 3)$
4. $P(-3 < X < 3|Y = -4)$

These probabilities can be computed using statistical tables or numerical integration methods. The code used as follows:-

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 import numpy as np
4 from scipy.stats import norm
5 import numpy as np
6 import math
7 ux=3
8 uy=1
9 varx=4
10 vary=4
11 rho=3/5
12 cdfy2 = norm.cdf(8, loc=uy, scale=vary)
13 cdfy1 = norm.cdf(3, loc=uy, scale=vary)
14 print(f"The probability from p(3<y<8) is {cdfy2-cdfy1}")

```

```

15 | xo=7
16 |
17 | cdfy2 = norm.cdf(8, loc=uy+rho*(uy/ux)*(xo-ux), scale=math.
    |         sqrt((1-rho*rho))*vary)
18 | cdfy1 = norm.cdf(3, loc=uy+rho*(uy/ux)*(xo-ux), scale=math.
    |         sqrt((1-rho*rho))*vary)
19 | print("The probability from P(3<y<8|ux=7) is ",{cdfy2-cdfy1
    |     })
20 | cdfy2 = norm.cdf(3, loc=ux, scale=varx)
21 | cdfy1 = norm.cdf(-3, loc=ux, scale=varx)
22 | print(f"The probability from P(-3<x<3) is {cdfy2-cdfy1}")
23 | yo=-4
24 | cdfy2 = norm.cdf(3, loc=ux+rho*(ux/uy)*(yo-uy), scale=math.
    |         sqrt((1-rho*rho))*varx)
25 | cdfy1 = norm.cdf(-3, loc=ux+rho*(ux/uy)*(yo-uy), scale=math.
    |         sqrt((1-rho*rho))*varx)
26 | print("The probability from P(-3<x<3|uy=-4) is ",{cdfy2-
    |     cdfy1})

```

The outputs were as follows:-

1. The probability $P(3 < Y < 8)$ is 0.26847838186216977.
2. The probability $P(3 < Y < 8|X = 7)$ is 0.32748810663813477.
3. The probability $P(-3 < X < 3)$ is 0.4331927987311419.
4. The probability $P(-3 < X < 3|Y = -4)$ is 0.1717928107053457.

Question 2

- (a) Write a program to generate P samples from a multinomial random variable $X \in \mathbb{R}^n$, having a multivariate normal distribution $\mathcal{N}_n(\mu, \Sigma)$, where $\mu \in \mathbb{R}^n$ denotes the mean vector and $\Sigma \in \mathbb{R}^{n \times n}$ is the covariance matrix of X .
- (b) Using P generated samples in part (a), obtain new samples using the following equation:

$$Y = (X - \mu)^T \Sigma^{-1} (X - \mu)$$

Observe the distribution of Y for different values of n and P .

- (c) Compute the probability that

$$\mathbb{P}[(x - \mu)^T \Sigma^{-1} (x - \mu) \leq c^2]$$

for a given c .

Introduction

In this work, we aim to generate samples from a multivariate normal distribution and analyze the quadratic form transformation:

$$Y = (X - \mu)^T \Sigma^{-1} (X - \mu).$$

The objective is to observe the distribution of Y for varying dimensions n and sample sizes P . Additionally, we compute the probability that:

$$\mathbb{P}[(X - \mu)^T \Sigma^{-1} (X - \mu) \leq c^2]$$

for a given threshold c . This probability is closely related to the chi-squared distribution with n degrees of freedom.

Methodology

Generating Multivariate Normal Samples

We generate P samples from an n -dimensional multivariate normal distribution $N_n(\mu, \Sigma)$. The mean vector $\mu \in \mathbb{R}^n$ is randomly selected, and the covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix.

Computing Quadratic Form Transformation

For each sample X_i , we compute the quadratic form:

$$Y_i = (X_i - \mu)^T \Sigma^{-1} (X_i - \mu).$$

This transformation follows a chi-squared distribution with n degrees of freedom.

Computing the Probability

Given a threshold c , we estimate:

$$\mathbb{P}[(X - \mu)^T \Sigma^{-1} (X - \mu) \leq c^2]$$

by computing the fraction of samples satisfying the condition and comparing it with the cumulative distribution function (CDF) of the chi-squared distribution. The code used is as follows:-

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import random
4 import seaborn as sns
5 import scipy.stats as st
6
7 p = [2, 3, 6, 8]
```

```

8 s = [10, 100, 1000]
9
10 for i in p:
11     for k in s:
12
13         mu = np.random.random(i) * 10
14         cov = np.zeros((i, i), dtype=float)
15
16
17         for f in range(i):
18             for j in range(f, i):
19                 if f == j:
20                     cov[f][j] = random.uniform(0.1, 50)
21                 else:
22                     cov[f][j] = random.uniform(-50, 50)
23                     cov[j][f] = cov[f][j]
24
25         cov = cov.T @ cov
26
27
28         x = np.random.multivariate_normal(mu, cov, k)
29
30         Y = []
31         for sample in x:
32             x_minus_mu = sample - mu # Shape (i,)
33
34             y = x_minus_mu @ np.linalg.inv(cov) @ x_minus_mu
35                 .T
36             Y.append(y)
37
38         plt.figure()
39         plt.hist(Y, bins=20, density=True, alpha=0.7, label=
40             f'Samples:_{k}')
41
42         x_vals = np.linspace(min(Y), max(Y), 100)
43         plt.plot(x_vals, st.chi2.pdf(x_vals, df=i), 'r-',
44             label=f'Chi-squared_{df={i}}')
45
46         plt.title(f'Distribution_{of}_{Y}_{for}_{k}_{samples}_{(
47             Dimensions:_{i})}')
48         plt.xlabel('Y')
49         plt.ylabel('Density')
50         plt.legend()
51         plt.show()

```

The program follows these steps:

1. Generate P samples from an n -dimensional multivariate normal distribu-

tion $N_n(\mu, \Sigma)$, where:

- $\mu \in \mathbb{R}^n$ is a randomly generated mean vector.
- $\Sigma \in \mathbb{R}^{n \times n}$ is a symmetric positive semi-definite covariance matrix.

2. Compute the quadratic form transformation:

$$Y = (X - \mu)^T \Sigma^{-1} (X - \mu).$$

3. Plot a histogram of computed Y values and overlay the probability density function (PDF) of the chi-squared distribution with n degrees of freedom.

The results obtained were as follows:-

For c part we used the following code-

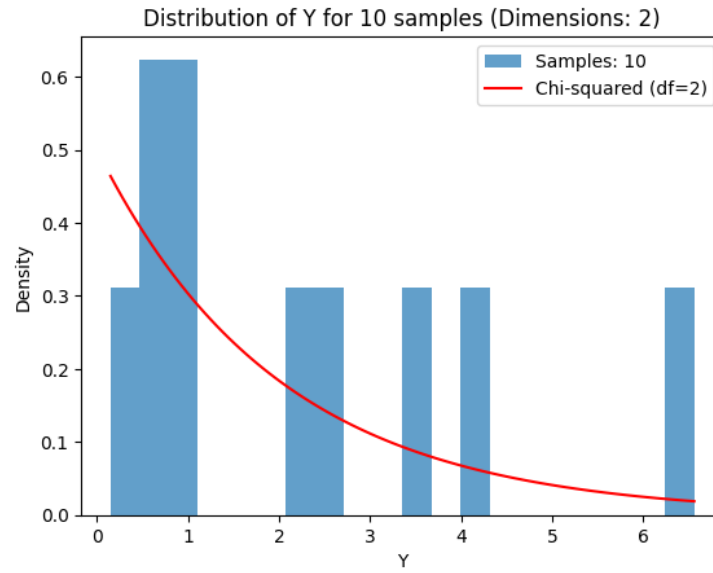


Figure 1:

```

1 import numpy as np
2 import scipy.stats as st
3 import matplotlib.pyplot as plt
4 c=int(input("Enter value of c:"))
5 for df in p:
6
7
8
9     pdf = st.chi2.pdf(c, df)
10    print(f"the value of (y-mu)^t*(sigma)^-1*(y-u) <={c*c} is {pdf}")

```

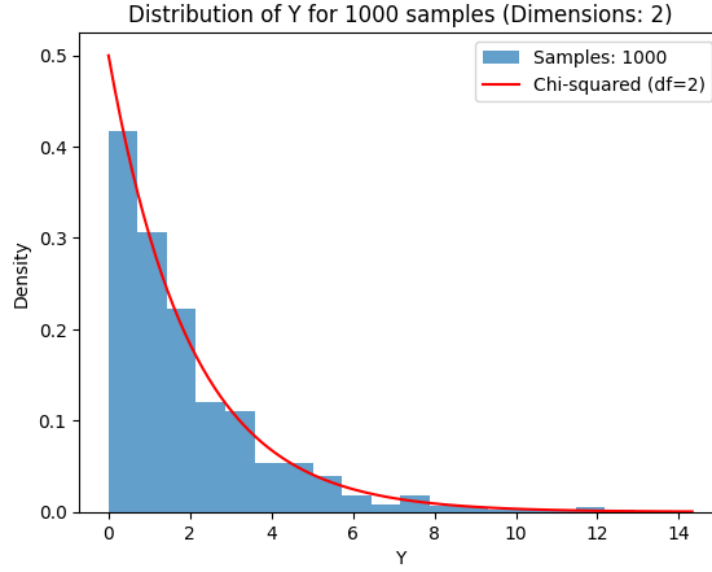


Figure 2:

Methodology

The program follows these steps:

1. The user provides a value for c .
2. The probability density function (PDF) of the chi-squared distribution is computed at c for different degrees of freedom n .
3. The probability is printed for each n .

The results obtained for $c=4$ were as follows:-

$$\Pr [(X - \mu)^T \Sigma^{-1} (X - \mu) \leq 16] = 0.06766764161830634$$

$$\Pr [(X - \mu)^T \Sigma^{-1} (X - \mu) \leq 16] = 0.10798193302637613$$

$$\Pr [(X - \mu)^T \Sigma^{-1} (X - \mu) \leq 16] = 0.1353352832366127$$

$$\Pr [(X - \mu)^T \Sigma^{-1} (X - \mu) \leq 16] = 0.09022352215774178$$

Question 3

3. The probability distributions of two different classes are known to follow a Normal distribution with the following parameters:

$$\mu_1 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \quad \mu_2 = \begin{bmatrix} -2 \\ -3 \end{bmatrix}$$

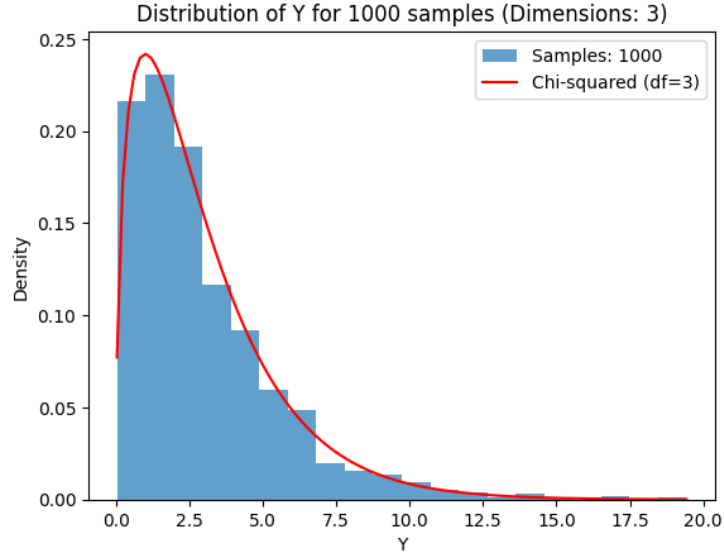


Figure 3:

$$\Sigma_1 = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 2 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 2 & -0.3 \\ -0.3 & 1 \end{bmatrix}$$

Use Bayes' Theorem to perform classification for the datapoints given in the attached file "File Datapoints.txt".

Demonstrate the result using a 2D diagram, illustrating the classes with different colors.

Introduction

In this study, we classify datapoints into two different classes using Bayes' Theorem. The probability distributions of the two classes follow a multivariate normal distribution with given mean vectors and covariance matrices. Given a set of datapoints, our objective is to determine the most likely class for each datapoint based on the posterior probabilities. The results will be demonstrated using a 2D plot, where different colors represent different classes.

Data

The dataset is provided in the file `File.Datapoints.txt`, which contains a collection of datapoints in a two-dimensional space. Each row in the file represents a datapoint with two numerical features. These datapoints will be classified into one of the two classes based on the Bayesian classification approach.

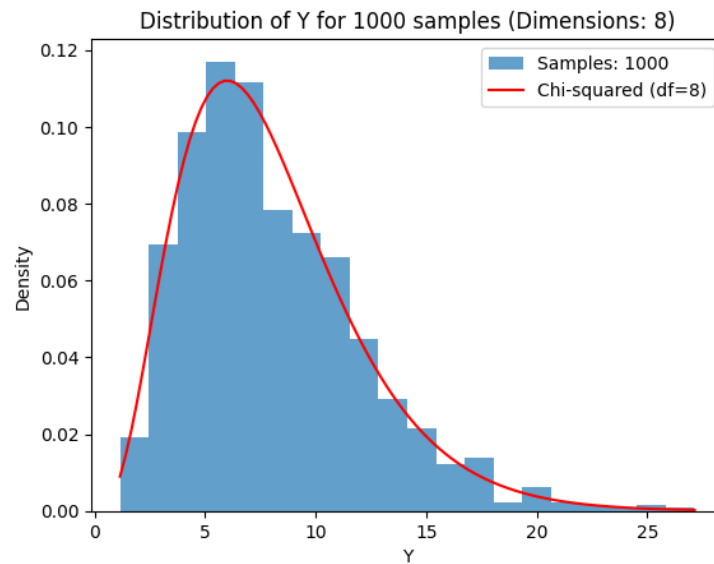


Figure 4:

Methodology

The code used is as follows:-

```

1 import re
2 import pandas as pd
3 from collections import Counter
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 with open('/content/File_Datapoints.txt','r') as file:
7     text=file.read()
8     words= text.split()
9
10
11 l=[]
12 for i in words:
13     l.append(i)
14 p=['s.no','x','y',]
15 for i in p:
16     l.remove(i)
17 l
18 t=0
19 x=[]
20 y=[]
21 for i in l:

```

```

22     if t==0:
23         t=(t+1)%3
24     elif t==1:
25         t=(t+1)%3
26         x.append(i)
27     else:
28         t=(t+1)%3
29         y.append(i)
30 df=pd.DataFrame(x,columns=['x'])
31 df['y']=pd.DataFrame(y)
32 df.index = df.index + 1
33 df = df.astype(float)
34 from scipy.stats import multivariate_normal
35 mu1=[2,3]
36 mu2=[-2,-3]
37 sig1=[[1,0.5],[0.5,2]]
38 sig2=[[2,-0.3],[-0.3,1]]
39
40 data=df.iloc[:,]
41 pl=[]
42 for i,j in data.iterrows():
43     x1=multivariate_normal.pdf(j,mu1,sig1)
44     x2=multivariate_normal.pdf(j,mu2,sig2)
45     if x1>x2:
46         pl.append(1)
47     else:
48         pl.append(2)
49 plt.figure(figsize=(8, 6))
50 plt.scatter(data[data.columns[0]], data[data.columns[1]], c=
    pl, cmap='viridis', s=50)
51 plt.scatter(mu1[0], mu1[1], c='orange', marker='*', s=200,
    label='Class_1_mean')
52 plt.scatter(mu2[0], mu2[1], c='red', marker='*', s=200,
    label='Class_2_mean')
53 plt.title("Data_Points_Classification_using_Bayes'_Theorem")
54 plt.xlabel("x")
55 plt.ylabel("y")
56 plt.legend()
57 plt.grid(True)
58 plt.show()

```

5 Implementation of Bayesian Classification

5.1 Reading and Processing the Data

The dataset is provided in the file `File_Datapoints.txt`. The following steps are performed to read and preprocess the data:

1. The file is opened and read as a text string.

2. The text is split into words, and column headers "s.no", "x", and "y" are removed.
3. The remaining values are separated into two lists corresponding to feature values x and y .
4. A Pandas DataFrame is created using these values, ensuring the data is stored in a numerical format.

5.2 Bayesian Classification

Bayes' theorem is used to classify each datapoint based on two known multivariate normal distributions corresponding to two different classes. The steps for classification are as follows:

1. Define the parameters for the two normal distributions:

- Mean vectors:

$$\mu_1 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \quad \mu_2 = \begin{bmatrix} -2 \\ -3 \end{bmatrix}$$

- Covariance matrices:

$$\Sigma_1 = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 2 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 2 & -0.3 \\ -0.3 & 1 \end{bmatrix}$$

2. For each datapoint x , compute the probability density function (PDF) of the multivariate normal distribution for both classes:

$$p_1 = \mathcal{N}(x|\mu_1, \Sigma_1), \quad p_2 = \mathcal{N}(x|\mu_2, \Sigma_2)$$

3. Assign the class label based on the higher probability density:

$$\text{Class} = \begin{cases} 1, & \text{if } p_1 > p_2 \\ 2, & \text{otherwise} \end{cases}$$

5.3 Visualization of Classification Results

A scatter plot is generated to visualize the classification results:

- Each datapoint is plotted using different colors corresponding to its assigned class.
- The means of both classes are highlighted using star markers ("*"), with orange representing Class 1 and red representing Class 2.
- The axes are labeled accordingly, and a legend is added to distinguish between classes.

The resulting plot illustrates the decision boundary between the two classes based on the Bayesian classification approach.

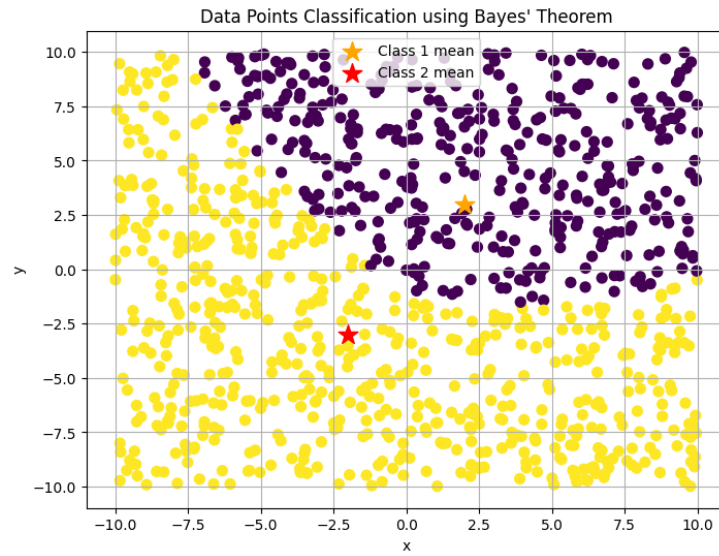


Figure 5: Data Points Classification using Bayes' Theorem

Results and observation

We can observe that a single curve differentiates both the classes the below portion of that curve falls into class 2 and above are in class 1 .