# Homework 7

1. (Spider on a wall) There's a spider living on a wall of your living room that has a painting behind which the spider likes to hide. Figure 1 shows a diagram of the wall; it is 10 feet high and 10 feet wide.
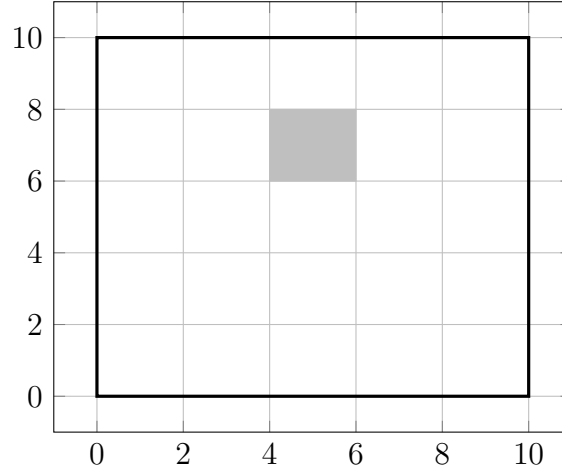


Figure 1: Wall and painting (in gray).

After observing the spider for a while you determine that (1) it spends twice the time behind the painting than on the rest of the wall, (2) it never crawls on the painting or leaves the wall, (3) if it is not behind the painting then it is equally likely to be anywhere on the wall. Since you cannot see it behind the painting, you assume that when it is there it is also equally likely to be at any spot.

(a) Model the position of the spider as a bivariate random variable and give its pdf.

We use $\tilde{x}, \tilde{y}$ to model the position of the spider. Then according to the problem statement we have:

$$\int_4^6 \int_6^8 f_{\tilde{x},\tilde{y}}(x,y)dydx = 2\left(1 - \int_4^6 \int_6^8 f_{\tilde{x},\tilde{y}}(x,y)dydx\right) \tag{1}$$

$$\iint_{(x,y)\notin \text{ Painting}} f_{\tilde{x},\tilde{y}}(x,y)dydx + \iint_{(x,y)\in \text{ Parting}} f_{\tilde{x},\tilde{y}}(x,y)dydx = 1 \tag{2}$$

By assumption, we have $f_{\tilde{x},\tilde{y}}(x,y) = \begin{cases} a, & (x,y) \in (4,6] \times [6,8] \\ b, & (x,y) \in [0,10] \times (0,10] \cap (x,y) \notin [4,6] \times [6,8]. \\ 0, & \text{otherwise.} \end{cases}$

1

Plugging into the formula (1) we get:

$$\int_4^6 \int_6^8 a\,dy\,dx = 2\left(1 - \int_4^6 \int_6^8 a\,dy\,dx\right)$$
$$4a = 2(1 - 4a)$$
$$a = \frac{1}{6}$$

Plugging into formula (2) we get:

$$(10 \times 10 - 2 \times 2)b + 4a = 1$$
$$b = \frac{1}{288}$$

Thus

$$f_{\tilde{x},\tilde{y}}(x,y) = \begin{cases} \frac{1}{6}, & (x,y) \in [4,6] \times [6 \times 8] \\ \frac{1}{288}, & (x,y) \in [0,10] \times [0,10] \cap (x,y) \notin [4,6] \times [6,8]. \\ 0, & \text{otherwise.} \end{cases}$$

(b) Compute the pdf of the height at which the spider is located and sketch it.

The pelf of the height should be the marginal distribution of $\tilde{y}$.

$$f_{\tilde{y}}(y) \quad = \int_0^{10} f_{\tilde{x},\tilde{y}}(x,y)dx$$
$$= \begin{cases} \int_0^{10} \frac{1}{288}dx = \frac{5}{144} & , y \in [0,6] \cup [8,10] \\ \int_0^4 \frac{1}{288}dx + \int_6^{10} \frac{1}{288}dx + \int_4^6 \cdot \frac{1}{6}dx = \frac{13}{36}, & y \in [6,8] \\ 0 & , \text{ otherwise} \end{cases}$$
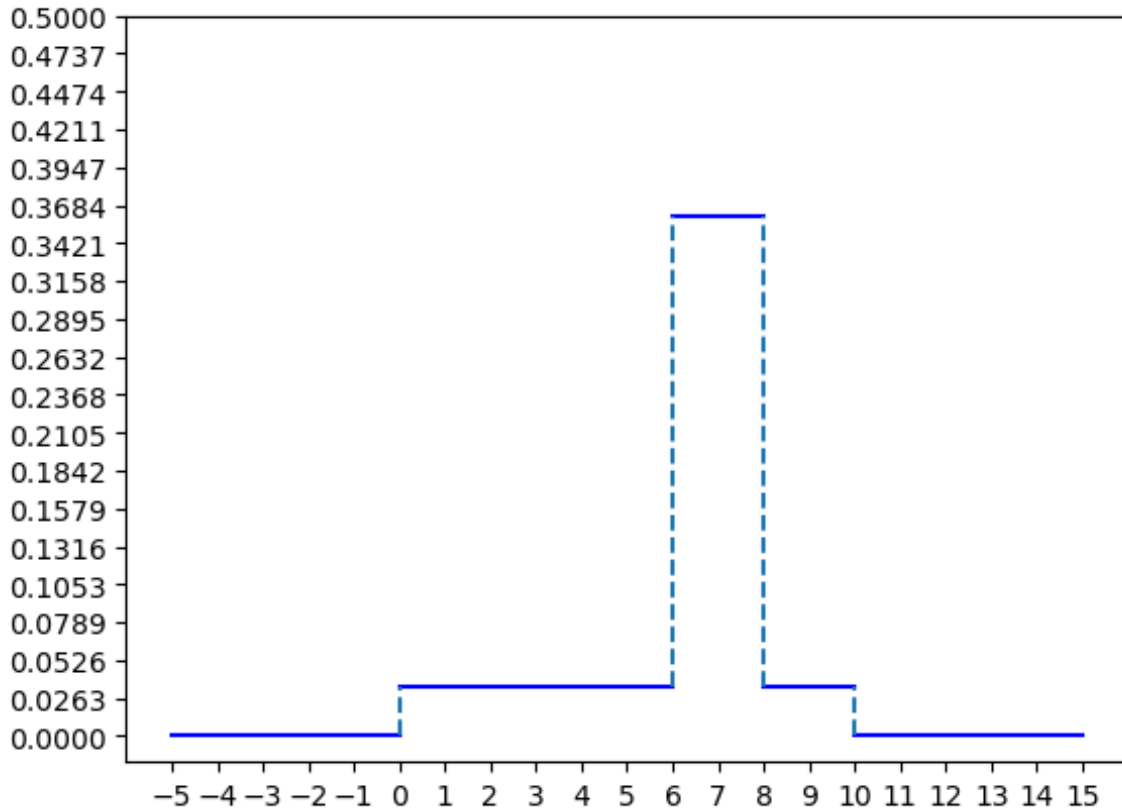


Figure 2: Marginal PMF Matrix

Attched is the link to the code(click it): Jupyter notebook link

```
1   # Problem 1(b)
2   x_1 = np.linspace(-5,0,100)
3   x_2 = np.linspace(0, 6,100)
4   x_3 = np.linspace(6, 8,100)
5   x_4 = np.linspace(8, 10, 100)
6   x_5 = np.linspace(10,15, 100)
7
8   plt.plot(x_1, 0 * np.ones(len(x_1)), c = "b", linestyle="-")
9   plt.plot(x_2, 5 / 144 * np.ones(len(x_2)), c = "b", linestyle="-")
10  plt.plot(x_3, 13 / 36 * np.ones(len(x_3)), c = "b", linestyle="-")
11  plt.plot(x_4, 5 / 144 * np.ones(len(x_4)), c = "b", linestyle="-")
12  plt.plot(x_5, 0 * np.ones(len(x_5)), c = "b", linestyle="-")
13  plt.vlines(0, ymin = 0, ymax = 5 /144, linestyle="—")
14  plt.vlines(6, ymin = 5 / 144 , ymax = 13 / 36, linestyle="—")
15  plt.vlines(8, ymin = 5 / 144, ymax = 13 / 36, linestyle="—")
16  plt.vlines(10, ymin = 0, ymax = 5 /144, linestyle="—")
17  plt.yticks(np.linspace(0,0.5,20))
18  plt.xticks(np.linspace(-5,15,21))
19  plt.show()
```

3

(c) Compute the conditional cdf of the height at which the spider is located, given that you can see it (i.e. it's not under the painting) and sketch it.

Denote the event $A$ : The spider isn't under the painting; then the conditional CDF of interest could be written as:

$$F_{\tilde{y}|A}(y \mid A) = \frac{P(\tilde{y} \le y \cap A)}{P(A)} \text{ where } P(A) = 1 - \int_4^6 \int_6^8 f_{\tilde{x},\tilde{y}}(x,y)dydx = \frac{1}{3}$$

$$= \begin{cases} 0 & , y \in (-\infty, 0) \\ \frac{P(\tilde{y} \le y)}{P(A)} = \frac{\int_0^y y(\tilde{y}(y)}{1/3} = \frac{5}{48}y & , y \in (0,6) \\ \frac{P(\tilde{y} \le y) - \int_4^6 \int_6^y \frac{1}{288}dydx}{1/3} = \frac{1}{12}y + \frac{1}{8}, & y \in [6,8) \\ \frac{P(\tilde{y} \le y) - \int_4^6 \int_6^8 \frac{1}{288}dydx}{1/3} = \frac{5}{48}y - \frac{1}{24}, & y \in [8,10) \\ 1 & y \in [10, \infty) \end{cases}$$
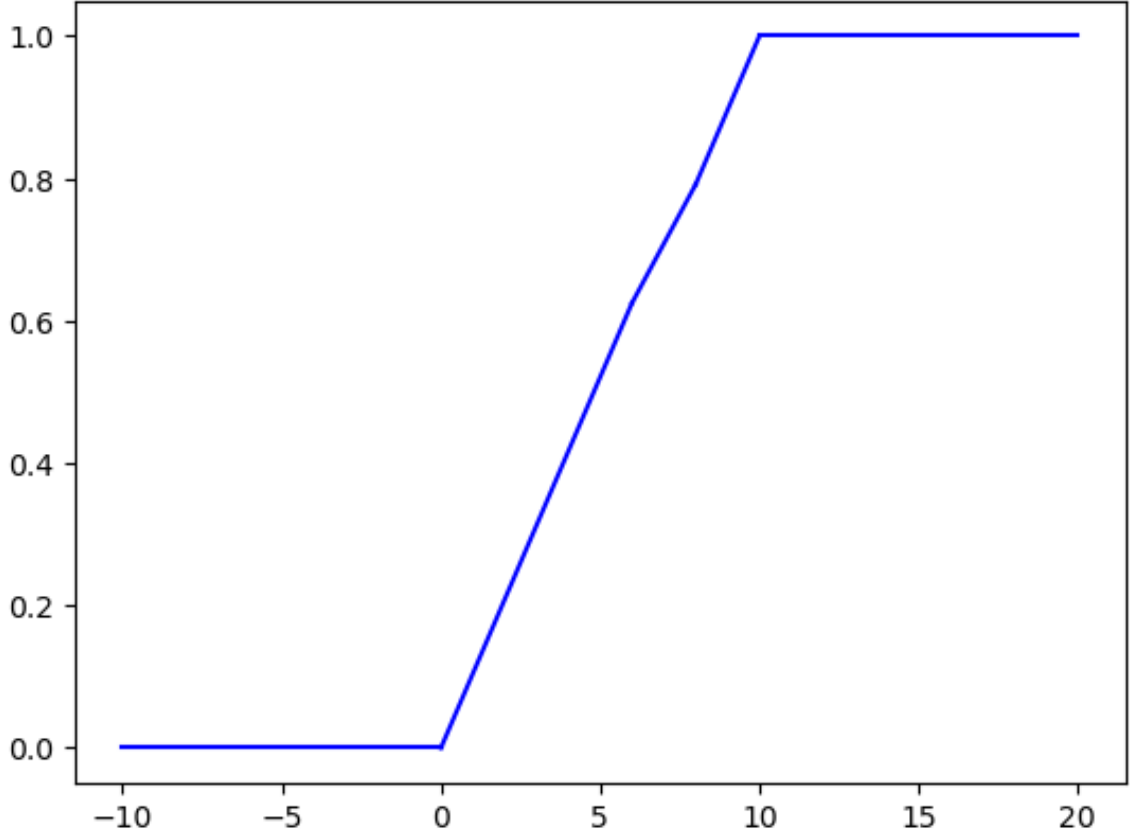


Figure 3: Marginal PMF Matrix

Attched is the link to the code(click it): Jupyter notebook link

```
1   # Problem 1(c)
2   def generate_piece_wise_function(v):
3
```

4

```
 4          if v == 0:
 5              return lambda x:0
 6          elif v == 1:
 7              return lambda x: (5 / 48) * x
 8          elif v == 2:
 9              return lambda x: (1 / 12) * x + 1 / 8
10          elif v == 3:
11              return lambda x: (5 / 48) * x − 1 / 24
12          else:
13              return lambda x: 1
14
15
16  def plot_piece_wise_function(pieces):
17
18      for index, piece in enumerate(pieces):
19          func = generate_piece_wise_function(index)
20          plt.plot(piece, [func(y) for y in piece],c = "blue")
21
22      plt.show()
23
24
25  pieces = [np.linspace(−10,0)
26          ,np.linspace(0,6)
27          ,np.linspace(6,8)
28          ,np.linspace(8,10)
29          ,np.linspace(10,20)]
30
31  plot_piece_wise_function(pieces)
```

2. (Sonar) A scientist is trying to determine the depth of the sea at a certain location. She knows that it must be deeper than 5 km but that is all she knows. To capture this uncertainty she models the depth as uniformly distributed between 5 and 10 km. In order to measure the depth she uses sonar, taking 2 measurements, which we model as two random variables $_1$ and $_2$. If the depth is equal to $x$ then each sonar measurement is uniformly distributed between $x - 0.25$ and $x + 0.25$. The two measurements are conditionally independent given the depth.

(a) Compute and sketch the pdf of the first sonar measurement $_1$.

Denote $\tilde{x}$ with $f_{\tilde{x}} = \begin{cases} \frac{1}{5}, & x \in [5, 10] \\ 0, & \text{otherwise.} \end{cases}$ Denote $\tilde{m}_1$ as measurement $1, \tilde{m}_2$ as measurement 2 . Then by problem statement we have:

$$f_{\tilde{m}_1|\tilde{x}}(m_1 \mid x) = \begin{cases} 2, & m_1 \in [x - 0.25, x + 0.25] \\ 0, & \text{otherwise.} \end{cases}$$

$$f_{\tilde{m}_2|\tilde{x}}(m_2 \mid x) = \begin{cases} 2, & m_2 \in [x - 0.25, x + 0.25] \\ 0, & \text{otherwise.} \end{cases}$$

$$f_{\tilde{m}_1,\tilde{m}_2|\tilde{x}}(m_1, m_2 \mid x) = f_{\tilde{m}_1|\tilde{x}}(m_1 \mid x) f_{\tilde{m}_2|\tilde{x}}(m_2 \mid x).$$

In order to compute $f_{\tilde{m}_1}(m_1)$ ( Similarly for $\tilde{m}_2$), we have:

$$f_{\tilde{m}_1}(m_1) = \int_5^{10} f_{\tilde{m}_1|\tilde{x}}(m_1 \mid x) f_{\tilde{x}}(x) \cdot dx$$

$$= \begin{cases} \int_{m_1-0.25}^{m_1+0.25} 2 \times \frac{1}{5} dx = \frac{1}{5}, & 5.25 \leqslant m_1 \leqslant 9.75 \\ \int_5^{m_1+0.25} 2 \times \frac{1}{5} dx = \frac{2}{5}(m_1 - 4.75), & 4.75 \leqslant m_1 \leqslant 5.25 \\ \int_{m_1-0.25}^{10} 2 \times \frac{1}{5} dx = \frac{2}{5}(10.25 - m_1), & 9.75 \leqslant m \leqslant 10.25 \\ 0 & \textit{otherwise.} \end{cases}$$

(b) Compute the conditional pdf of the depth conditioned on the measurements being equal to 7 km and 7.1 km.

The conditional probability of interest is:

$$
\begin{aligned}
f_{\tilde{x}|\tilde{m}_1,\tilde{m}_2}(x \mid 7, 7.1) &= \frac{f_{\tilde{m}_1,\tilde{m}_2|\tilde{x}}(7, 7.1 \mid x) f_{\tilde{x}}(x)}{f_{\tilde{m}_1,\tilde{m}_2}(7, 7.1)} \\
&= \frac{f_{\tilde{m}_1|\tilde{x}}(7 \mid x) \cdot f_{\tilde{m}_2|\tilde{x}}(7.1 \mid x) f_{\tilde{x}}(x)}{\int_5^{10} f_{\tilde{m}_1|\tilde{x}}(7 \mid x) f_{\tilde{m}_2|\tilde{x}}(7.1 \mid x) f_{\tilde{x}}(x) dx} \\
&= \begin{cases} 0 & , x \in [5, 6.85) \\ \frac{2\times 2 \times \frac{1}{5}}{\int_{6.85}^{7.25} 2\times 2 \times \frac{1}{5} dx} = \frac{4/5}{4/5 \times 0.4} = 2.5, & x \in [6.85, 7.25) \\ 0, & x \in [7.25, 10] \end{cases}
\end{aligned}
$$

(c) Compute the joint pdf of the two sonar measurements $_1$ and $_2$. Are the two measurements independent? Justify your answer mathematically and explain it intuitively.

$$f_{\tilde{m}_1, \tilde{m}_2}(m_1, m_2) = \int_5^{10} f_{\tilde{m}_1|\tilde{x}}(m_1 \mid x) f_{\tilde{m}_2|\tilde{x}}(m_2 \mid x) f_{\tilde{x}}(x) dx.$$

$$= \begin{cases} 0, & |m_1 - m_2| \geqslant 0.5 \cup \min(10, \min(m_1, m_2) + 0.25) < \max(5, \max(m_1, m_2) - 0.25). \\ \int_{\max(m_1,m_2)-0.25}^{\min(m_1,m_2)+0.25} f_{\tilde{m}_1,|\tilde{x}}(m_1 \mid x) f_{\tilde{m}2}|\tilde{x}(m_2 \mid x) f_{\tilde{x}}(x) dx, & |m_1 - m_2| < 0.5 \end{cases}$$

If we pick $m_1 = 6, m_2 = 8$, then

$$f_{m_1}(m_1) = f_{\tilde{m}_2}(m_2) = \frac{1}{5}$$

However $f_{\tilde{m}_1, \tilde{m}_2}(6, 8) = 0$

$$\therefore f_{\tilde{m}_1}(m_1) f_{\tilde{m}_2}(m_2) \neq f_{\tilde{m}_1}, \tilde{m}_2(m_1, m_2)$$

$\therefore$ Mathematically specking, $\tilde{m}_1, \tilde{m}_2$ aren't independent. Intuitively speaking, if we know the measurement 1, we will know that the depth will likely be within the range of $(m_1 - 0.25, m_1 + 0.25)$. Then this depth will give clues to the range of values within which $m_2$ would likely be.

3. **(Exotic fruit)** A pomologist is studying a newly-discovered plant that produces edible fruit. They have available 10 specimens, which she measures to obtain the following data:

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Fruit weight (kg) | 1.5 | 2.3 | 0.8 | 1.2 | 2.0 | 1.2 | 0.7 | 2.7 | 2.3 | 0.6 |
| Stem height (cm) | 20 | 15 | 18 | 19 | 17 | 22 | 21 | 14 | 17 | 22 |
| Stem radius (mm) | 8 | 12 | 6 | 10 | 17 | 12 | 9 | 14 | 10 | 8 |
| Number of leafs | 110 | 94 | 152 | 123 | 78 | 60 | 111 | 83 | 85 | 90 |
| Median leaf length (cm) | 12 | 21 | 9 | 14 | 19 | 15 | 7 | 29 | 22 | 15 |

The pomologist wants to use these data to model the conditional pdf of the fruit weight given the other four features.

(a) Why is it difficult to estimate the conditional pdf using a simple nonparametric estimator where we apply kernel density estimation to a relevant subset of the data?

Denote the conditional probability that we want to estimate by $f_{\widetilde{\omega}|\tilde{h},\tilde{r},\tilde{l},\tilde{m}}(\omega \mid h, r, l, m)$, where $\widetilde{\omega}$ corresponds to weight, $\tilde{h}$ to height, $\tilde{r}$ to radius, $\widetilde{l}$ to # of leafs, $\tilde{m}$ to median leaf length. Non parametric model falls victims to curse of dimensionality. In this example, in order to estimate that conditional probability, we have to know the joint distribution of $(\tilde{h}, \tilde{r}, \tilde{l}, \tilde{m})$. This is tractable when each random variable only takes finite values. But $\tilde{h}, \tilde{r}, \tilde{l}, \tilde{m}$ are all continuous RV's. thus estimating $f_{\tilde{h},\tilde{r},\tilde{l},\tilde{m}}$ is just impossible given that we only have finite amount of data.

(b) Model all the features as jointly Gaussian and use the model to compute the conditional pdf of the fruit weight for a plant that has a stem height of 15 cm, a stem radius of 20 mm, 120 leafs and median leaf length of 8 cm.

Denote $(\widetilde{w}, \tilde{h}, \tilde{r}, \tau, \widetilde{m}) = \widetilde{\vec{x}} \in R^5$, which follows multivariate gaussian distribution with pal:

$$f_{\vec{x}}(\vec{x}) = \frac{1}{(2\pi)^{\frac{5}{2}} |\Sigma_5|} \exp\left\{-(\vec{x} - \vec{\mu})^\top \Sigma_5^{-1}(\vec{x} - \vec{\mu})\right\}$$

$$\vec{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^{n} \vec{x}_i$$

$$\Sigma_{ML} = \frac{1}{n} \sum_{i=1}^{n} (\vec{x}_i - \vec{\mu}_{ML})(\vec{x}_i - \vec{\mu}_{ML})^\top$$

We then compute the conditional pdf of $\vec{x}[1]$ given $\vec{x}[2:5]$, which is parametrized by;

$$\mu_{\text{cond}} = \mu_{\vec{x}[1]} + \Sigma_{\vec{x}[2:5],\vec{x}[1]} \Sigma_{\vec{x}[2:5]}^{-1} \left(\vec{x}[2:5] - \vec{\mu}_{\vec{x}[2:5]}\right)$$

$$\Sigma_{\text{cond}} = \Sigma_{\vec{x}[1]} - \Sigma_{\vec{x}[2:5],\vec{x}(1)}^\top \Sigma_{\vec{x}[2:5]}^{-1} \Sigma_{\vec{x}[2:5],\vec{x}[1]}$$

By python code we get $\mu_{\text{cond}} = 1.57772945$ and $\Sigma_{\text{cond}} = 0.03747546420415787$ where $\vec{x}[2:5] = [15, 20, 120, 8]^\top$ and get the conditional pdf at this point. Attached is the code and pdf result.
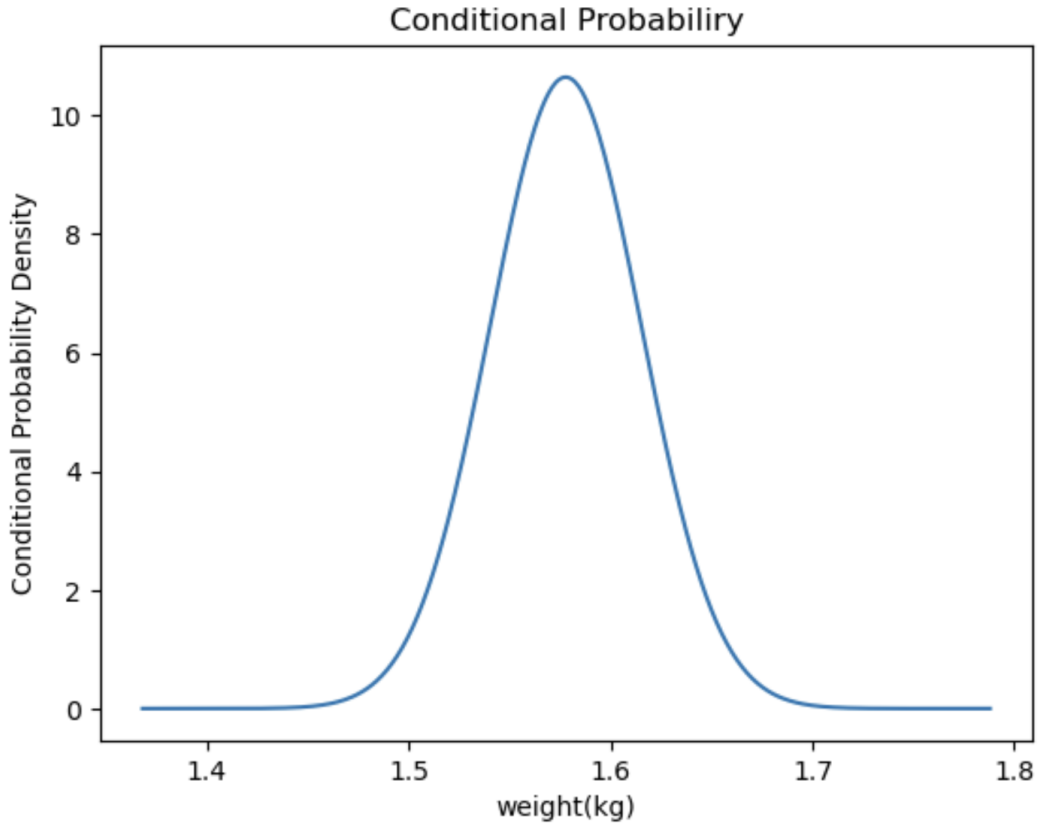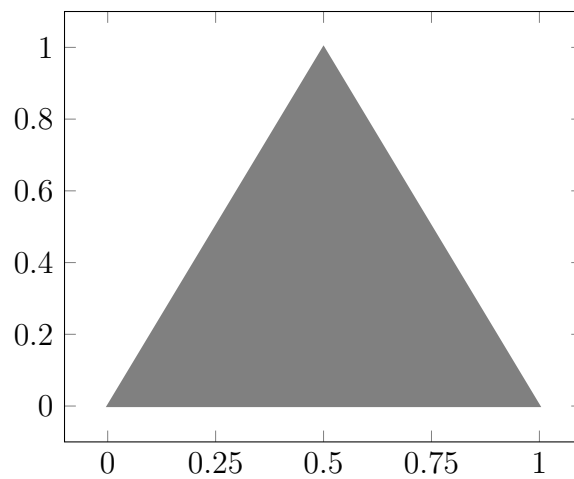


Figure 4: Marginal PMF Matrix

Attched is the link to the code(click it): Jupyter notebook link

```
1    # Problem 3(b)
2
3    # A trick here, we reverse the order of the feature
4    #  , only to apply the multi-dimensional gaussian
5    # estimate quickly
6    raw_data = np.array(list(reversed([[1.5,  2.3,  0.8,  1.2,  2.0,  1.2,  0.7,  2.7,  2.3,  0.6]
7                        ,[20,  15,  18,  19 ,  17 ,  22 ,  21 ,  14 ,  17 ,  22]
8                        ,[8,  12,  6,   10,   17,   12,   9,   14,   10,   8 ]
9                        ,[110,  94,  152,  123,  78,  60,  111,  83,  85,  90]
10                       ,[12,  21,  9,  14,  19,  15,  7,  29,  22,  15]
11                       ])))
12   data = raw_data.T
13   n = data.shape[0]
14   condition_x = np.array([15,  20,  120,  8])
15
16
17   # Mean Vector
18   mu_MLE = 1 / n * data.sum(axis = 0)
19   # Covariance matrix
20   Sigma_MLE = 1 / n * (data − mu_MLE).T @ (data − mu_MLE)
21
22
23   sub_mu_MLE_y = mu_MLE[4:]   # (1,)    y
24   sub_mu_MLE_x = mu_MLE[:4]   # (4,)    x
25
26   sub_Sigma_MLE_y = Sigma_MLE[4:,4:]                    # (4,1)   y
27   sub_Sigma_MLE_yx = Sigma_MLE[:4,4:]                   # (x,y)
28   sub_Sigma_MLE_x = np.linalg.inv(Sigma_MLE[:4,:4])     # (4,4)    x
29
30   # Compute conditional model, parametrized by mu_cond and sigma_cond
31   mu_cond = sub_mu_MLE_y + sub_Sigma_MLE_yx.T
32               @ sub_Sigma_MLE_x
33               @ (condition_x − sub_mu_MLE_x)
34   sigma_cond = (sub_Sigma_MLE_y − sub_Sigma_MLE_yx.T
35               @ sub_Sigma_MLE_x
36               @ sub_Sigma_MLE_yx)[0][0]
37
38
39
40   # Generate the conditional probability
41   us = np.linspace(1.578 − 0.03 * 7,1.578 + 0.03 * 7,500)
42   plt.plot(us,  norm.pdf(us,  mu_cond,  sigma_cond))
43   plt.xlabel("weight(kg)")
44   plt.ylabel("Conditional-Probability-Density")
45   plt.title("Conditional-Probabiliry")
46   plt.show()
```

4. (Simulating a random vector) Explain how to simulate a random two-dimensional vector with a joint pdf $f(x)$ that is uniformly distributed in the shaded region. Assume that you have access to independent samples from a uniform distribution in $[0, 1]$. Explain your method, justifying why it works. Then implement it and submit the code along with a scatterplot of 1,000 samples.
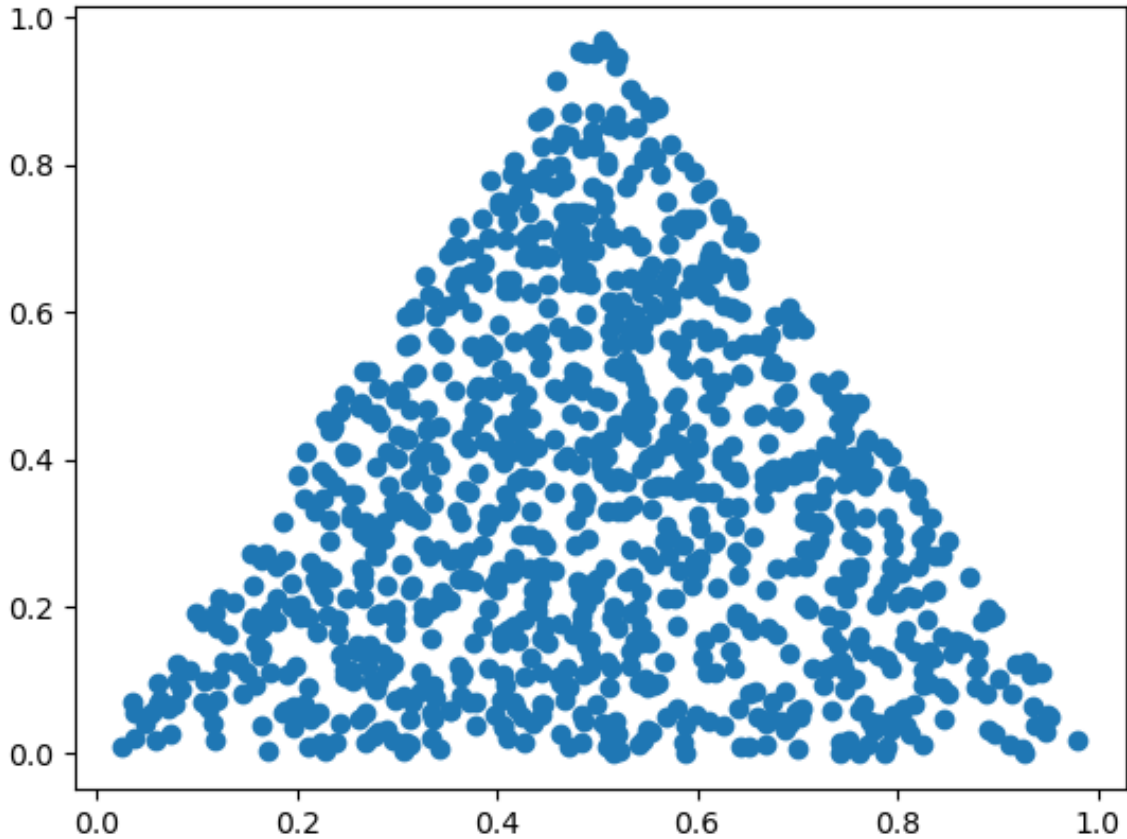
Figure 5: Marginal PMF Matrix

```
1   import numpy as np
2   import scipy.stats
3   import matplotlib.pyplot as plt
4   %matplotlib inline
5
6   # Generate jointly distributed samples from independent uniform distribution
7   def generate_one_sample(inverse_CDF_u1, inverse_CDF_u2):
8       u1 = generate_u1(inverse_CDF_u1)
9       u2 = generate_u2(inverse_CDF_u2, u1)
10      return u1, u2
11
12  # Just for this problem
13  def define_inverse_CDF():
14      func_u1 = lambda x: 1 - np.sqrt(1 - x)
15      func_u2 = lambda x, y: (1 - y) * x + (1 / 2) * y
16
17      return func_u1, func_u2
18
19  # Apply Inverse Sampling
20  def generate_u1(inverse_CDF):
21      u = np.random.uniform(0,1,1)
22      return inverse_CDF(u)
23
24  # Applying Inverse Sampling, Conditional Probability
25  def generate_u2(inverse_CDF, u1):
26      u = np.random.uniform(0,1,1)
27      return inverse_CDF(u, u1)
28
29  def draw_scatter_plot(data):
30      # Since we first generate y, then x.
31      plt.scatter(data[:,1], data[:,0])
32      plt.show()
33
34  def main_procedure(sample_size = 1000):
35      sample_points = []
36      inverse_CDF_u1, inverse_CDF_u2 = define_inverse_CDF()
37      for i in range(sample_size):
38          sample_points.append([*generate_one_sample(inverse_CDF_u1, inverse_CDF_u2)])
```

```
39
40        data = np.array(sample_points)
41        draw_scatter_plot(data)
42
43  main_procedure(1000)
```