

Assignment 1 Rui Qin 30874157

Question 1 (4 marks)

1. Creating an image of a gorilla riding a motorcycle: **Generative AI systems**
2. Discovering the shopping preferences of customers at Coles: **Recommendation systems**
3. Predicting the amount of Panadol that will be purchased in the next week: **Forecasting**
4. Using genomic and lifestyle factors to determine the chance of an individual contracting diabetes in the next two years: **Risk prediction**

Question 2 (6 marks)

1. Marginal probability of winning ($P(W = 1)$):
$$P(W = 1)$$
$$= P(W = 1, H = 0, P = 0) + P(W = 1, H = 0, P = 1) + P(W = 1, H = 1, P = 0) + P(W = 1, H = 1, P = 1)$$
$$= 0.588$$
2. Probability of winning a game given that the team lost their previous game ($P(W = 1 \mid P = 0)$):
$$P(W = 1 \mid P = 0)$$
$$= (P(W = 1, H = 0, P = 0) + P(W = 1, H = 1, P = 0)) / (P(P = 0, H = 0) + P(P = 0, H = 1))$$
$$= (0.235 + 0.058) / (0.176 + 0.235 + 0.117 + 0.058)$$
$$= 0.5$$
3. Probability of winning a game given that the team won their previous game ($P(W = 1 \mid P = 1)$):
$$P(W = 1 \mid P = 1)$$
$$= (P(W = 1, H = 0, P = 1) + P(W = 1, H = 1, P = 1)) / (P(P = 1, H = 0) + P(P = 1, H = 1))$$
$$= (0.117 + 0.178) / (0.060 + 0.117 + 0.059 + 0.178)$$
$$= 0.713$$
4. Yes. We can compare the probabilities in questions 2 and 3:
$$P(W = 1 \mid P = 1) > P(W = 1 \mid P = 0)$$

It shows that the team is more likely to win after winning their previous game compared to winning after losing their previous game.

5. Probability of winning next home game after previous win:

$$\begin{aligned}
 &P(W = 1 \mid P = 1, H = 1) \\
 &= P(W = 1, P = 1, H = 1) / (P(W = 1, H = 1, P = 1) + P(W = 0, H = 1, P = 1)) \\
 &= 0.178 / (0.178 + 0.059) \\
 &= 0.75
 \end{aligned}$$

Probability of winning next game after previous home game win:

$$\begin{aligned}
 &P(W = 1 \mid P = 1, H = 0) \\
 &= P(W = 1, P = 1, H = 0) / (P(W = 1, P = 1, H = 0) + P(W = 0, P = 1, H = 0)) \\
 &= 0.117 / 0.177 \\
 &= 0.661
 \end{aligned}$$

Probability of winning next two games:

$$\begin{aligned}
 &P(W = 1 \mid P = 1, H = 1) \cdot P(W = 1 \mid P = 1, H = 0) \\
 &= 0.75 \cdot 0.661 \\
 &= 0.496
 \end{aligned}$$

Question 3 (8 marks)

1. Expected value of S:

$$\begin{aligned}
 &E[S] \\
 &= 2 \cdot E[X1] - E[Y1] \\
 &= 2 \cdot ((1 + 2 + 3 + 4 + 5 + 6) / 6) - ((1 + 2 + 3 + 4) / 4) \\
 &= 4.5
 \end{aligned}$$

2. Variance of S:

$$\begin{aligned}
 &\text{Var}[S] \\
 &= E[S^2] - E[S]^2 \\
 &= E[(2 \cdot X1 - Y1)^2] - 20.25 \\
 &= E[4 \cdot X1^2 - 4 \cdot X1Y1 + Y1^2] - 20.25 \\
 &= 4 \cdot E[X1^2] - 4 \cdot E[X1Y1] + E[Y1^2] - 20.25 \\
 &\quad \bullet \quad E[X1^2] = (1^2 + 2^2 + 3^2 + 4^2 + 5^2 + 6^2) / 6 = 91/6 \\
 &\quad \bullet \quad E[Y1^2] = (1^2 + 2^2 + 3^2 + 4^2) / 4 = 7.5 \\
 &\quad \bullet \quad E[X1 \cdot Y1] = 3.5 \cdot 2.5 = 8.75 \\
 &= 12.917
 \end{aligned}$$

3. Probability distribution of S:

For $S = -2$:

- $P(S = -2) = P(2 \cdot X_1 - Y_1 = -2)$ if:
 - $X_1 = 1, Y_1 = 4$
- $P(S = -2) = 1/6 \cdot 1/4 = 1/24$

For $S = -1$:

- $P(S = -1) = P(2 \cdot X_1 - Y_1 = -1)$ if:
 - $X_1 = 1, Y_1 = 3$
- $P(S = -1) = 1/24$

... ..

The final result:

For $P(S = s) = 1/24: s \in \{-2, -1, 10, 11\}$

For $P(S = s) = 1/12: s \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$

4. Expected value of S^3 :

$$E(S^3)$$

$$= \sum [P(S) \cdot S^3]$$

$$= (-2)^3 \cdot P(S=-2) + (-1)^3 \cdot P(S=-1) + 0^3 \cdot P(S=0) + 1^3 \cdot P(S=1) + \dots + 11^3 \cdot P(S=11)$$

$$= (-2)^3 \cdot 1/24 + (-1)^3 \cdot 1/24 + 1^3 \cdot 1/12 + 2^3 \cdot 1/12 + 3^3 \cdot 1/12 + 4^3 \cdot 1/12 + 5^3 \cdot 1/12 + 6^3 \cdot 1/12 + 7^3 \cdot 1/12 + 8^3 \cdot 1/12 + 9^3 \cdot 1/12 + 10^3 \cdot 1/24 + 11^3 \cdot 1/24$$

$$= 265.5$$

5. Approximate value of $E(S^3)$ using the Taylor series procedure:

$$f(x) = x^3, f'(x) = 3x^2, f''(x) = 6x, f'''(x) = 6$$

$$E[S] = 4.5, E[X_1] = 3.5, E[Y_1] = 2.5$$

$$E[S^3]$$

$$= E[f(4.5) + f'(4.5) \cdot (S-4.5) + f''(4.5) / 2! \cdot (S-4.5)^2 + f'''(4.5) / 3! \cdot (S-4.5)^3]$$

$$= 91.125 + 60.75 \cdot E[S - 4.5] + 27 / 2! \cdot E[(S - 4.5)^2] + 6 / 3! \cdot E[(S - 4.5)^3]$$

\therefore

$$\text{Var}[S] = E[(S-E[S])^2] = 12.917$$

$$E[S - 4.5] = E[S] - 4.5 = 0$$

\therefore

$$= 91.125 + 60.75 \cdot 0 + 13.5 \cdot 12.917 + E[(S - 4.5)^3]$$

$$= 91.125 + 60.75 \cdot 0 + 13.5 \cdot 12.917 + E[(S - 4.5)^2] \cdot 0$$

$$= 265.5045$$

6. Expected value of $(2X_1 - Y_1 + 2Y_2)^2$:

$$(2X_1 - Y_1 + 2Y_2)^2$$

$$= 4X_1^2 - 4X_1Y_1 + 8X_1Y_2 + Y_1^2 - 4Y_1Y_2 + 4Y_2^2$$

$$E[(2X_1 - Y_1 + 2Y_2)^2] = 4 \cdot E[X_1^2] - 4 \cdot E[X_1 \cdot Y_1] + 8 \cdot E[X_1 \cdot Y_2] + E[Y_1^2] - 4 \cdot E[Y_1 \cdot Y_2] + 4 \cdot E[Y_2^2]$$

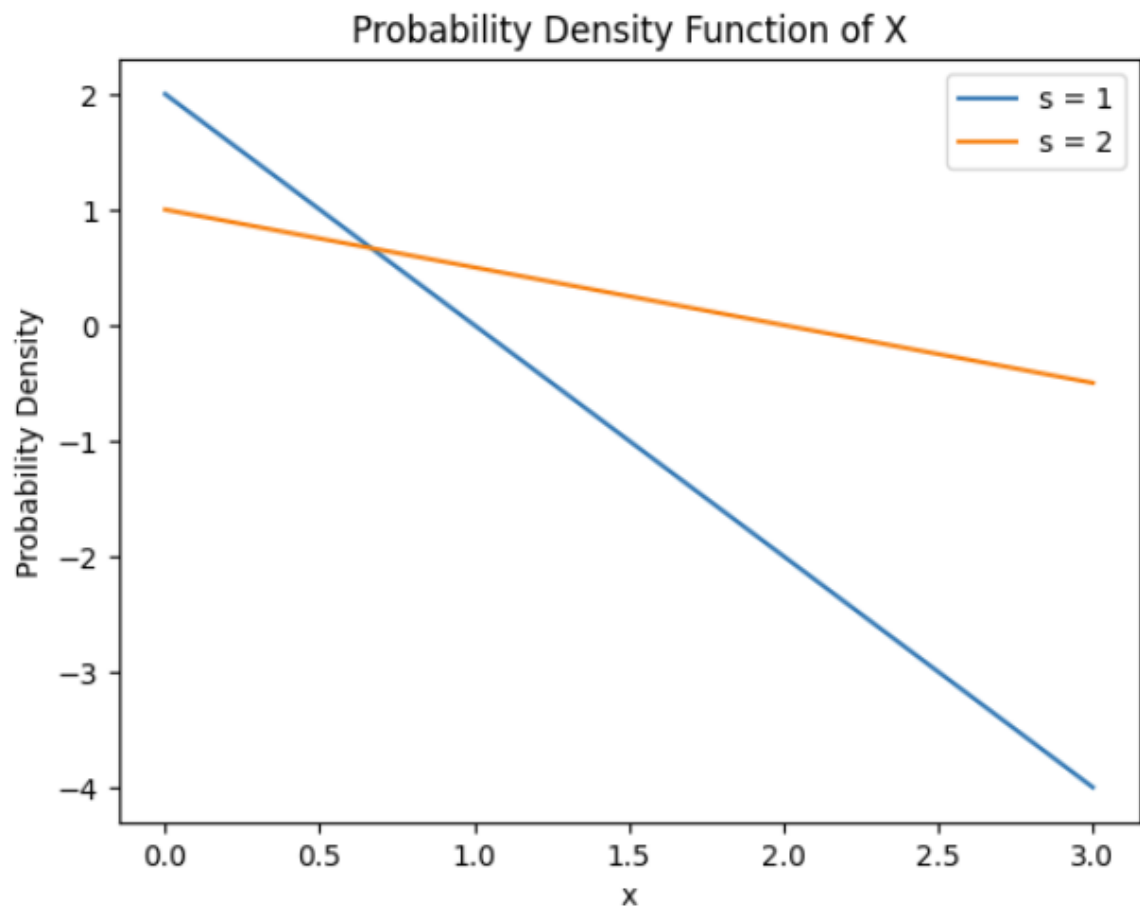
- $E[X_1^2] = 91/6$
- $E[Y_1^2] = E[Y_2^2] = E[Y_1 \cdot Y_2] = 7.5$
- $E[X_1 \cdot Y_1] = E[X_1 \cdot Y_2] = E[X_1] \cdot E[Y_1] = 3.5 \cdot 2.5 = 8.75$

$$= 4 \cdot 91/6 - 4 \cdot 8.75 + 8 \cdot 8.75 + 7.5 - 4 \cdot 7.5 + 4 \cdot 7.5$$

$$\approx 103.167$$

Question 4 (7 marks)

1. Plotting the Probability Density Function (PDF):



2. Expected Value of X ($E[X]$):

$$\begin{aligned} E[X] &= \int_{[0,s]} x \cdot 2(s-x) / s^2 dx \\ &= s/3 \end{aligned}$$

3. Expected Value of \sqrt{x} :

$$\begin{aligned} E[\sqrt{x}] &= \int_{[0,s]} \sqrt{x} \cdot 2(s-x) / s^2 dx \\ &= 4\sqrt{s} / 3 - 1 \end{aligned}$$

4. Variance of X ($V[X]$):

$$\begin{aligned} \text{Var}[x] &= E[(x - E[x])^2] \\ E[X^2] &= \int_{[0,s]} x^2 \cdot 2(s-x) / s^2 dx \\ &= s^2/6 \\ V[x] &= s^2/6 - (s/3)^2 \\ &= s^2/18 \end{aligned}$$

5. Median of X:

To find the median of the X with the given probability density function, we need to use cumulative distribution function = 0.5.

$$\text{CDF}(x | s) = \int_{[-\infty, x]} p(X = t | s) dt = 0.5$$

$$\therefore p(X = t | s) = 2(s-t)/s^2$$

$$\therefore \text{CDF}(x | s)$$

$$= \int_{[0,x]} 2(s-t)/s^2 dt$$

$$= (2sx - x^2) / s^2 = 0.5$$

$$\therefore x^2 - 2sx + 0.5s^2 = 0$$

$$x_1 = \frac{1}{2} (\sqrt{2s} + 2s)$$

$$x_2 = \frac{1}{2} (-\sqrt{2s} + 2s)$$

$$\therefore 0 < x < s, \frac{1}{2} (\sqrt{2s} + 2s) > s$$

$$\therefore \text{Median}(x) = \frac{1}{2} (-\sqrt{2s} + 2s)$$

Question 5 (7 marks)

1. It takes around 15.556 days for a patient to recover from COVID-19 based on Poisson distribution

```
# Question 1
library(MASS)
fit <- fitdistr(df$Days, "Poisson")
estimated_lambda <- fit$estimate["lambda"]
print(estimated_lambda)
```

```
> print(estimated_lambda)
lambda
15.556
```

2.

- a. The possibility can be calculated with R, and final answer is 0.0938464

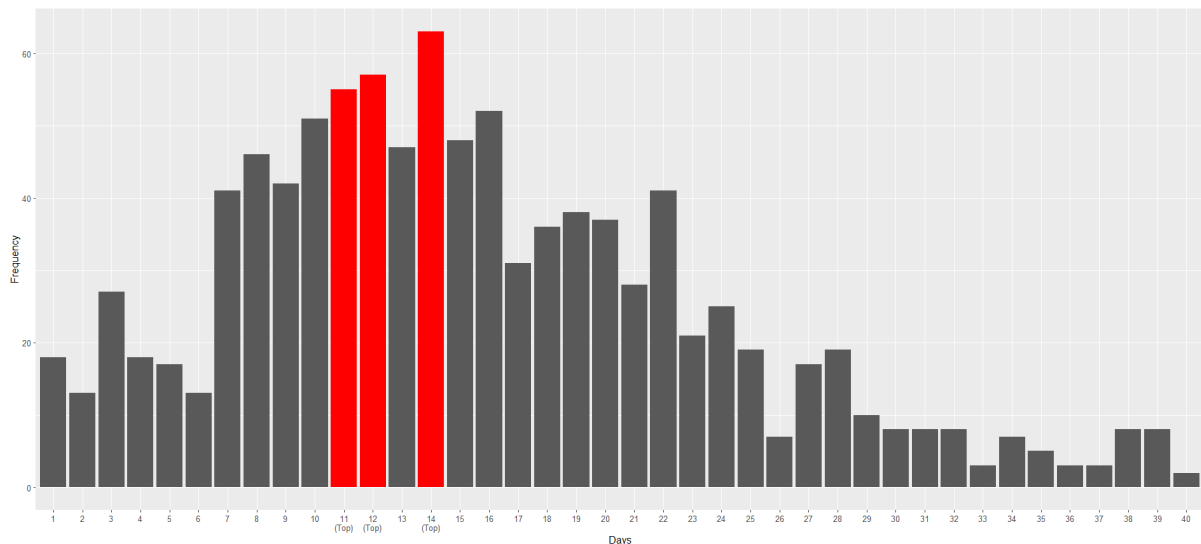
```
# Question 2A
lambda <- 15.556
probability_10_or_less <- ppois(10, lambda)
print(probability_10_or_less)
```

```
> print(probability_10_or_less)
[1] 0.0938464
```

- b. We can know though the graph that, 11, 12 and 14 are the three most likely number of days it takes a patient to recover

```
#Question 2B
freq_df <- data.frame(table(df$Days))
top_values <- head(freq_df[order(-freq_df$Freq), ], 3)$Var1

ggplot(df, aes(x = factor(Days))) +
  geom_bar(stat = "count") +
  geom_bar(data = subset(df, Days %in% top_values),
           aes(x = factor(Days)),
           stat = "count", fill = "red") +
  scale_x_discrete(labels = function(x)
    ifelse(x %in% top_values, paste(x, "\n(Top)", sep = ""),
    as.character(x))) +
  labs(x = "Days", y = "Frequency")
```



- c. We can use function $P(X = k) = (\lambda^k * e^{(-\lambda)}) / k!$ to solve this problem in R, and the result is 0.6115397

```
# Question 2C
# Function to calculate Poisson probability base on PMF
poisson_prob_range <- function(lambda_f, k_min, k_max) {
  probabilities <- numeric(k_max - k_min + 1)
  for (k in k_min:k_max) {
    probabilities[k - k_min + 1] <- (lambda_f^k * exp(-lambda_f)) /
factorial(k)
  }
  return(sum(probabilities))
}
# Total rate parameter for five individuals, Range is 60 to 80 days
probability <- poisson_prob_range(5 * lambda, 60, 80)
print(probability)
```

```
> print(probability)
[1] 0.6115397
```

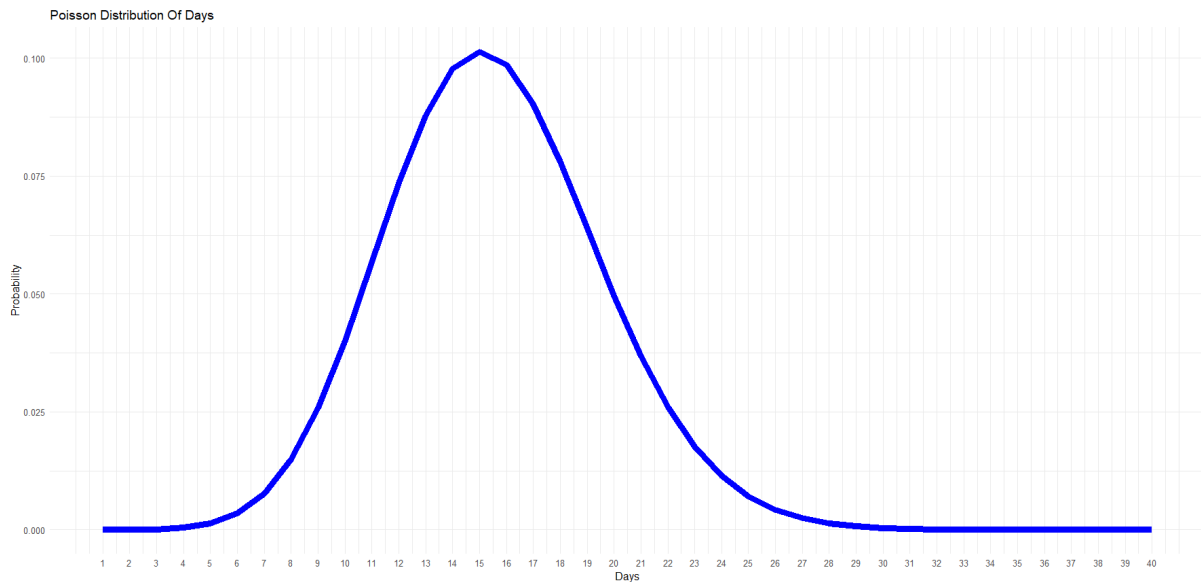
- d. Using binomial distribution since we have multiple patients and based on the formula: $P(X \geq 3) = P(X = 3) + P(X = 4) + P(X = 5)$, the result is 0.8205065

```
# Question 2D
prob_after_14 <- 1 - ppois(13, lambda)
probability_three_or_more <- sum(dbinom(3:5, 5, prob_after_14))
print(probability_three_or_more)
```

```
> print(probability_three_or_more)
[1] 0.8205065
```

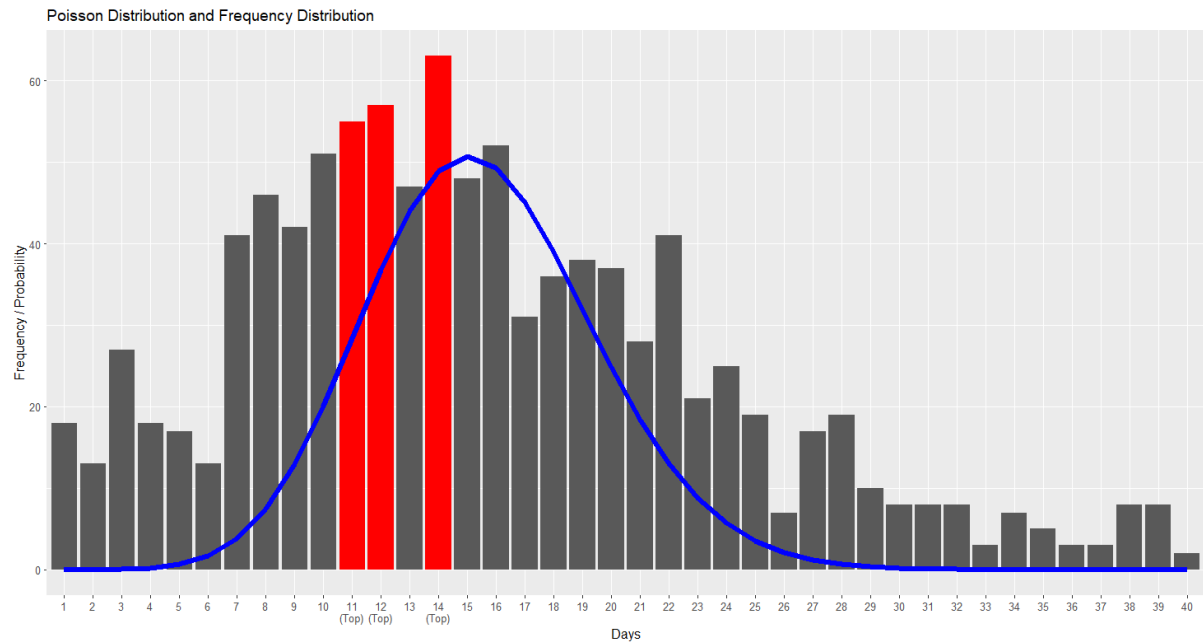
3. Start with drawing Poisson Distribution graph

```
# Question 3
ggplot() +
  geom_line(aes(x = df$Days, y = dpois(df$Days, lambda)), color = "blue", size = 3) +
  labs(title = "Poisson Distribution Of Days",
       x = "Days",
       y = "Probability") +
  scale_x_continuous(breaks = seq(min(df$Days), max(df$Days), by = 1))
```



We can force the Poisson distribution to scale and then merge it with the data distribution map

```
ggplot(df, aes(x = factor(Days))) +
  geom_bar(stat = "count") +
  geom_bar(data = subset(df, Days %in% top_values),
          aes(x = factor(Days)),
          stat = "count", fill = "red") +
  geom_line(aes(x = df$Days,
                y = 500*dpois(df$Days, lambda)),
            color = "blue", size = 2) +
  scale_x_discrete(labels = function(x)
    ifelse(x %in% top_values, paste(x, "\n(Top)", sep = ""), as.character(x))) +
  labs(title = "Poisson Distribution and Frequency Distribution",
       x = "Days",
       y = "Frequency / Probability")
```

As we can see, the Poisson distribution does not provide a perfect fit to the data, most people recover after 11-14 days but the top point of Poisson distribution is 15. However, the two distributions are generally similar, and the Poisson distribution does capture the general trend of the data. Therefore, I believe that the Poisson distribution is an appropriate model for the COVID recovery data.