# ajmal-jalal_assignment_2.1

September 16, 2024

# 1 Assignment 2.1

**Name:** Ajmal Jalal

**Date:** 09/15/2024

For this assignment, you will refer to the textbook to solve the practice exercises. **Use Python to answer any coding problems (not R, even if indicated in your textbook).** Use Jupyter Notebook, Google Colab, or a similar software program to complete your assignment. Submit your answers as a **PDF or HTML** file. As a best practice, always label your axes and provide titles for any graphs generated on this assignment. Round all quantitative answers to 2 decimal places.

## 1.1 Problem # 2.1.

For the rain simulation example in Section 2.1.1, but with probability of rain 0.30 on any given day, simulate the outcome (a) on the next day, (b) the next 10 days. (c) Simulate the proportion of days of rain for the next (i) 100 days, (ii) 10,000 days, (iii) 1,000,000 days. Use the simulation to explain the long-run relative frequency definition of probability.

(a) The outcome for the next day

```
[ ]: import numpy as np

     # Setting the probability of rain to 0.30 (30%)
     prob_rain = 0.30

     # Simulating the outcome for the next day
     next_day_rain = np.random.binomial(1, prob_rain, 1)
     print(f"The probability of rain for the next day is: {next_day_rain[0]}")
```

The probability of rain for the next day is: 0

(b) The probability of rain in the next 10 days

```
[ ]: # Simulating the outcome for the next 10 days
     next_10_days_rain = np.random.binomial(1, prob_rain, 10)
     print(f"The probability of rain in the next 10 days is: {next_10_days_rain}")
```

The probability of rain in the next 10 days is: [0 0 0 0 0 0 1 0 0 0]

(c) The probability of rain in the next 100, 10,000, 1,000,000 days

```
# Simulating the outcome for the next 100 days
next_100_days_rain = np.round(np.random.binomial(100, prob_rain) / 100, 2)
print(f"The probability of rain in the next 100 days is: {next_100_days_rain:.
  ↪2f}")
# Simulating the outcome for the next 10,000 days
next_10000_days_rain = np.round(np.random.binomial(10000, prob_rain) / 10000, 2)
print(f"The probability of rain in the next 10,000 days is:↵
  ↪{next_10000_days_rain:.2f}")
# Simulating the outcome for the next 1,000,000 days
next_1000000_days_rain = np.round(np.random.binomial(1000000, prob_rain) /↵
  ↪1000000, 2)
print(f"The probability of rain in the next 1,000,000 days is:↵
  ↪{next_1000000_days_rain:.2f}")
```

```
The probability of rain in the next 100 days is: 0.24
The probability of rain in the next 10,000 days is: 0.30
The probability of rain in the next 1,000,000 days is: 0.30
```

## 1.2 Problem # 2.2.

Data analysts often implement statistical inference methods by setting the probability of a correct inference equal to 0.95. Let $A$ denote the event that an inference for the population about men is correct. Let $B$ represent the event of a corresponding inference about women being correct. Suppose that these are independent events.

(a) Find the probability that (i) *both* inferences are correct, (ii) *neither* inference is correct.

(b) Construct the probability distribution for $Y$ = number of correct inferences.

(c) With what probability would each inference need to be correct in order for the probability to be 0.95 that *both* are correct?

(a) Probability that (i) both inferences are correct, (ii) neither inference is correct.

```python
import numpy as np

# Given probability of correct inference
p = 0.95

# Calculating the probability that both inferences are correct can be found by
  ↪multiplying the probability of each inference being correct
prob_both_correct = p * p
print(f"(i) The probability that both inferences are correct is:↵
  ↪{prob_both_correct:.4f}")

# Calculating the probability that neither inference is correct can be found by
  ↪multiplying the probability of each inference not being correct
```

2

```
prob_neither_correct = (1 - p) * (1 - p)
print(f"(ii) The probability that neither inference is correct is:␣
   ↪{prob_neither_correct:.4f}")
```

(i) The probability that both inferences are correct is: 0.9025
(ii) The probability that neither inference is correct is: 0.0025

(b) Probability distribution for Y (number of correct inferences)

```
# Creating a list of y values
y_values = [0, 1, 2]
# Creating a list of probabilities for each y value
probabilities = [prob_neither_correct, 2 * p * (1 - p), prob_both_correct]

print("\n(b) Probability distribution for Y:")
for y, prob in zip(y_values, probabilities):
    print(f"   P(Y = {y}) = {prob:.4f}")
```

(b) Probability distribution for Y:
    P(Y = 0) = 0.0025
    P(Y = 1) = 0.0950
    P(Y = 2) = 0.9025

(c) Required probability for both to be correct with 0.95 probability

```
# Required probability for each inference to be correct with 0.95 probability␣
   ↪can be found by taking the square root of 0.95
p_required = np.sqrt(0.95)
print(f"Required probability for each inference: {p_required:.4f}")
```

Required probability for each inference: 0.9747

## 1.3   Problem # 2.4.

A wine connoisseur is asked to match five glasses of red wine with the bottles from which they came, representing five different grape types.

(a) Set up a sample space for the five guesses.

(b) With random guessing, find the probability of getting all five correct.

(a) Setting up a sample space for the five guesses

The sample space is all possible permutations of 5 wines.

There are 5! (5 factorial) possible arrangements = 5 * 4 * 3 * 2 * 1 = 120

(b) Finding the probability of getting all five correct with random guessing

```
# Total number of possible arrangements
total_outcomes = 120
```

```python
# There are 5 wines, and each must be correctly matched to its bottle.
# The first wine must be in position 1, the second in position 2, and so on.
# Any deviation from this single correct arrangement results in an incorrect
  ↪guess.
# So, only one arrangement is correct (all wines matched perfectly)
correct_outcome = 1

# Calculating the probability of getting all five correct
probability_all_correct = correct_outcome / total_outcomes

print(f"The probability of getting all five correct with random guessing is:
  ↪{probability_all_correct:.4f}")
print(f"This is equivalent to {probability_all_correct * 100:.2f}% or 1 in
  ↪{total_outcomes} chance")
```

```
The probability of getting all five correct with random guessing is: 0.0083
This is equivalent to 0.83% or 1 in 120 chance
```

## 1.4 Problem # 2.15.

Each week an insurance company records $Y =$ number of payments because of a home burning down. State conditions under which we would expect $Y$ to approximately have a Poisson distribution.

Answer

In Poisson distribution the following conditions should exist:

1. Events are independent

2. The average rate of occurrence is constant

3. Two events cannot occur at exactly the same instant

4. The probability of an event is proportional to the length of the interval

Based on the above, the number of payments for homes burning down (Y) would approximately follow a Poisson distribution under these conditions:

1. Independence: Each event (home burning down) occurs independently of others.

2. Constant rate: The average rate of occurrences remains constant over the time period.

3. Rare events: The probability of a home burning down is small for any given short time interval.

4. Non-overlapping: Two events cannot occur simultaneously.

5. Well-defined time period: The time frame (one week) is clearly specified.

## 1.5 Problem # 2.16.

Each day a hospital records the number of people who come to the emergency room for treatment.

(a) In the first week, the observations from Sunday to Saturday are 10, 8, 14, 7, 21, 44, 60. Do you think that the Poisson distribution might describe the random variability of this phenomenon adequately. Why or why not?

(b) Would you expect the Poisson distribution to better describe, or more poorly describe, the number of weekly admissions to the hospital for a rare disease? Why?

## (a) Assessing Poisson Distribution Adequacy

To determine whether the Poisson distribution is an adequate model for the given data, we'll perform the following steps using Python:

1. Calculate the Mean and Variance of the Observed Data.
2. Compare the Mean and Variance to Assess Overdispersion.
3. Visualize the Data Against the Poisson Distribution.

```python
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import poisson

# Given data for the number of people coming to the emergency room each day
 from Sunday to Saturday
observed_data = np.array([10, 8, 14, 7, 21, 44, 60])

# Calculating Mean
mean_admissions = np.mean(observed_data)
print(f"Mean number of admissions: {mean_admissions:.2f}")

# Calculating Variance
variance_admissions = np.var(observed_data, ddof=1) # Using sample variance
print(f"Variance of admissions: {variance_admissions:.2f}")

# Comparing Mean and Variance
if variance_admissions > mean_admissions:
    print("Variance is greater than the mean. Overdispersion is present.")
else:
    print("Variance is approximately equal to the mean. Poisson distribution
 may be adequate.")

# Plotting the data and the Poisson distribution
unique, counts = np.unique(observed_data, return_counts=True)

plt.bar(unique, counts, alpha=0.6, label='Observed Data')
x = np.arange(0, max(observed_data)+1)
poisson_pmf = poisson.pmf(x, mean_admissions)
plt.plot(x, poisson_pmf, 'ro-', label='Poisson PMF')
plt.xlabel('Number of Admissions')
plt.ylabel('Frequency')
plt.title('Observed Data vs Poisson Distribution')
```
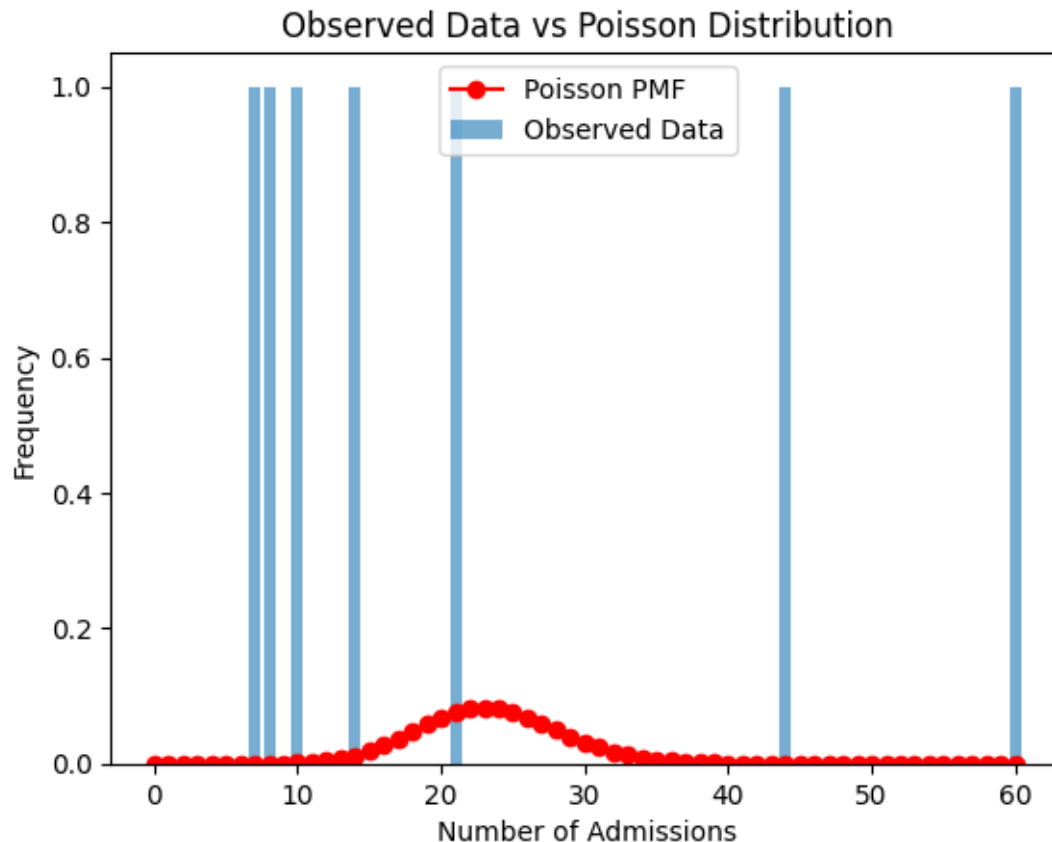
```
plt.legend()
plt.show()
```

```
Mean number of admissions: 23.43
Variance of admissions: 423.95
Variance is greater than the mean. Overdispersion is present.
```



Observed Data vs Poisson Distribution

**Explanation:**

1. **Mean and Variance Calculation:**
   - **Mean:** The average number of daily admissions.
   - **Variance:** Measures the variability of the admissions.
2. **Comparison and conclusion:**
   - If the variance is significantly larger than the mean, it indicates **overdispersion**, suggesting that the Poisson distribution may not be adequate.

`(b) Evaluating Poisson distribution for rare disease events`

To evaluate whether the Poisson distribution is appropriate for modeling the number of weekly admissions to the hospital for a **rare disease**, we'll consider the characteristics of rare events and perform illustrative Python calculations.

**Key Characteristics of Rare Events Suitable for Poisson Distribution:**

1. **Low Probability of Occurrence:** Rare diseases have a low incidence rate, meaning the probability of an admission due to such diseases in any given week is low.
2. **Independence:** Admissions occur independently of one another.
3. **Constant Rate ( ):** The average rate ( ) of admissions remains relatively stable over time.
4. **Discrete Events:** Admissions are countable events, fitting the nature of the Poisson model.

Given these characteristics, the Poisson distribution is generally **more appropriate** for modeling the number of weekly admissions for a rare disease.

## 1.6 Problem # 2.17.

An instructor gives a course grade of B to students who have total score on exams and homeworks between 800 and 900, where the maximum possible is 1000. If the total scores have approximately a normal distribution with mean 830 and standard deviation 50, about what proportion of the students receive a B?

`Answer`

To determine the proportion of students who receive a grade of B (with total scores between 800 and 900) given that the scores are normally distributed with a mean of 830 and a standard deviation of 50, we can perform the following:

- We will use the cumulative distribution function (CDF) of the normal distribution to find the probabilities up to 900 and up to 800.

- The difference between these two probabilities gives the proportion of students scoring between 800 and 900.

```python
import numpy as np
import scipy.stats as stats
import matplotlib.pyplot as plt

# Given data
mean = 830
std_dev = 50
lower_bound = 800
upper_bound = 900

# Calculating the cumulative probabilities
prob_up_to_upper = stats.norm.cdf(upper_bound, mean, std_dev)
prob_up_to_lower = stats.norm.cdf(lower_bound, mean, std_dev)

# Calculating the proportion of students scoring between 800 and 900
proportion_B = prob_up_to_upper - prob_up_to_lower
print(f"Proportion of students receiving a B: {proportion_B:.4f}
 ({proportion_B*100:.2f}%)")

# Plotting the normal distribution
x = np.linspace(mean - 4*std_dev, mean + 4*std_dev, 1000)
y = stats.norm.pdf(x, mean, std_dev)
```
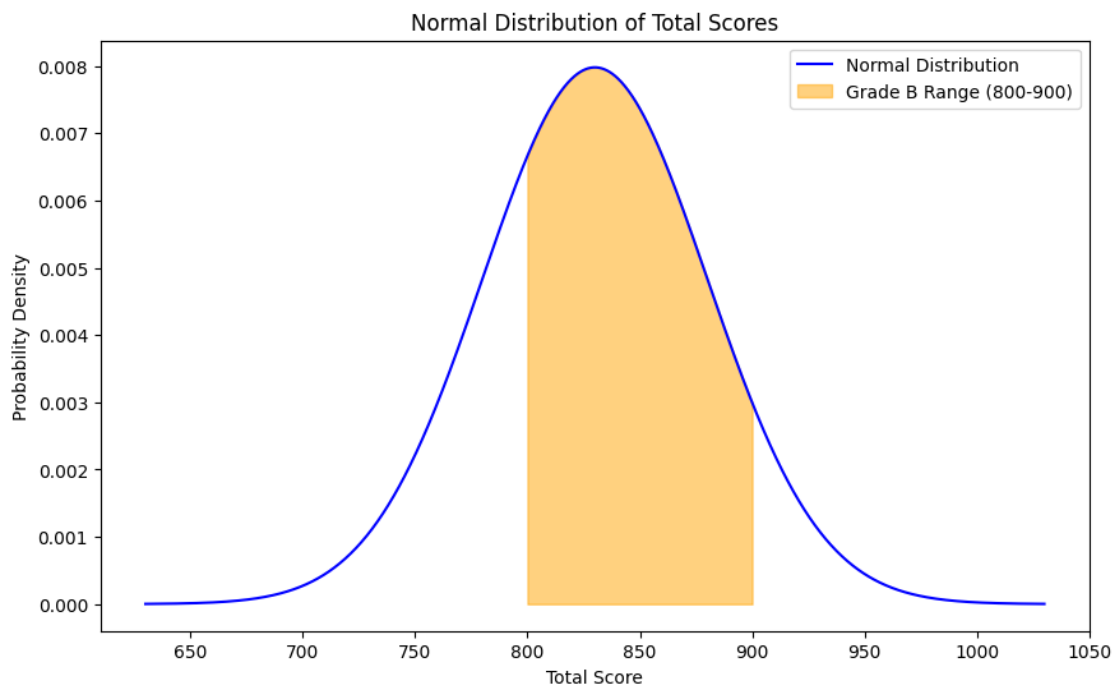
```python
plt.figure(figsize=(10, 6))
plt.plot(x, y, label='Normal Distribution', color='blue')

# Shading the area between 800 and 900 for better visualization
x_fill = np.linspace(lower_bound, upper_bound, 1000)
y_fill = stats.norm.pdf(x_fill, mean, std_dev)
plt.fill_between(x_fill, y_fill, color='orange', alpha=0.5, label='Grade B␣
  ↪Range (800-900)')

plt.title('Normal Distribution of Total Scores')
plt.xlabel('Total Score')
plt.ylabel('Probability Density')
plt.legend()

plt.show()
```

Proportion of students receiving a B: 0.6450 (64.50%)



## 1.7   Problem # 2.20.

Create a data file with the income values in the `Income` data file at the text website.

(a) Construct a histogram or a smooth-curve approximation for the *pdf* of income in the corresponding population by plotting results using the density function in R (explained in Exercise 1.18).

8

(b) Of the probability distributions studied in this chapter, which do you think might be most appropriate for these data? Why? Plot the probability function of that distribution having the same mean and standard deviation as the income values. Does it seem to describe the income distribution well?

(a) Constructing a Histogram and Density Plot for Income

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats

# Loading the Income data from the URL and save it locally
import requests

url = "https://stat4ds.rwth-aachen.de/data/Income.dat"
response = requests.get(url)

# Save the content to 'Income.dat'
with open('Income.dat', 'w') as file:
    file.write(response.text)

# Now that the data is stored in 'Income.dat' with columns: income, education,
 ↪race, we can read it into a pandas dataframe
income_data = pd.read_csv('Income.dat', sep=r'\s+')

# Plotting the histogram with a density curve
plt.figure(figsize=(10, 6))
sns.histplot(income_data['income'], bins=20, kde=True, color='skyblue',
 ↪edgecolor='black')
plt.title('Histogram and Density Plot of Income')
plt.xlabel('Income')
plt.ylabel('Frequency')
plt.show()


print(f"P-value: {p_value:.4f}")
```
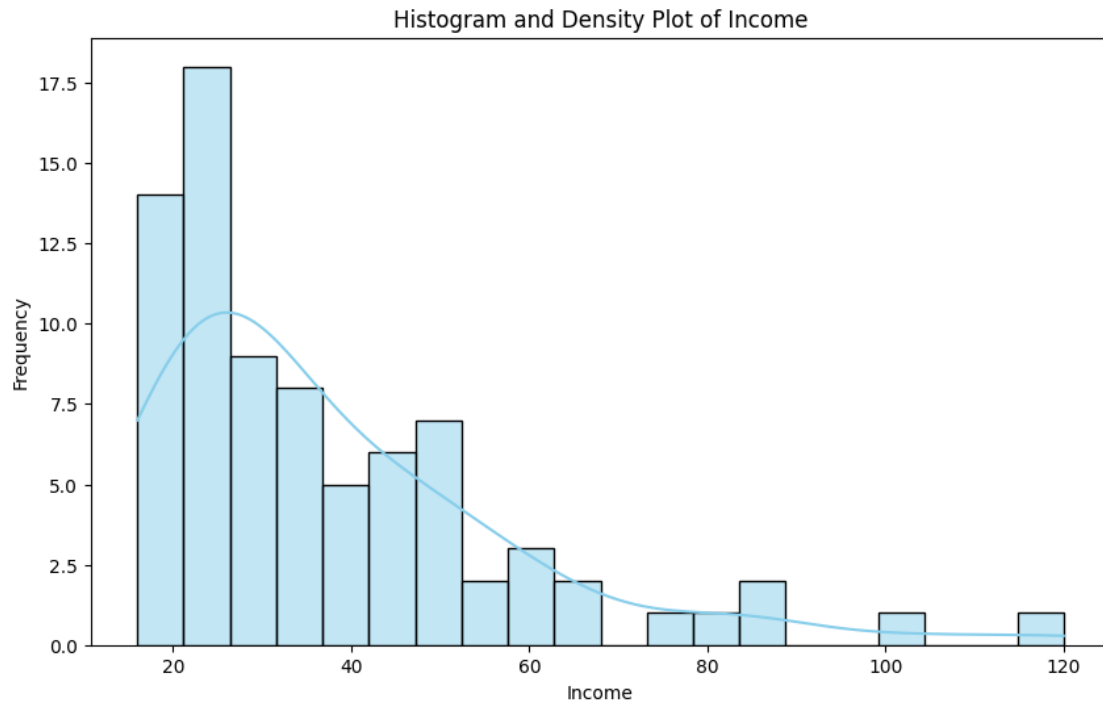
Histogram and Density Plot of Income

P-value: 0.2881

(b) Fitting an Appropriate Probability Distribution

```python
# Calculate mean and standard deviation of the income data
mean_income = income_data['income'].mean()
std_income = income_data['income'].std()
print(f"Mean Income: {mean_income:.2f}")
print(f"Standard Deviation of Income: {std_income:.2f}")

# Based on the histogram, determine an appropriate distribution.
# Common choices for income data include Normal, Log-Normal, and Gamma
  ↪distributions.
# Here, we'll attempt to fit a Log-Normal distribution.

# Fit a Log-Normal distribution to the income data
shape, loc, scale = stats.lognorm.fit(income_data['income'], floc=0)   # Fixing
  ↪location to 0

# Creating a range of income values for plotting the fitted distribution
x = np.linspace(income_data['income'].min(), income_data['income'].max(), 1000)
pdf_fitted = stats.lognorm.pdf(x, shape, loc, scale)

# Plotting the histogram and the fitted Log-Normal PDF
plt.figure(figsize=(10, 6))
```
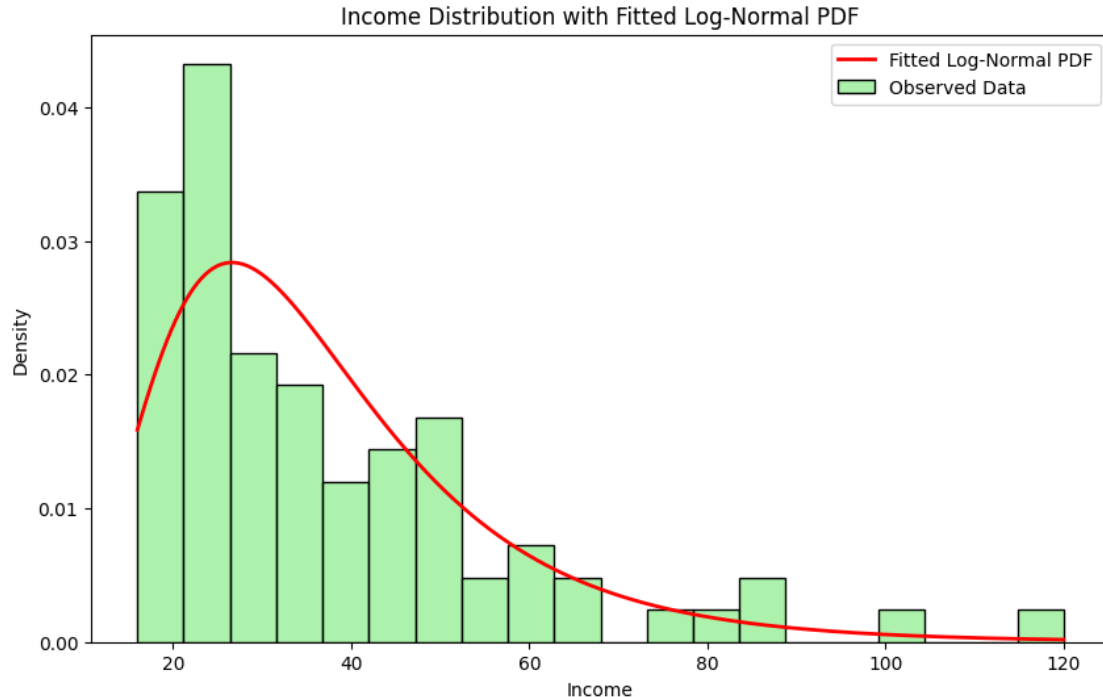
```python
sns.histplot(income_data['income'], bins=20, kde=False, stat='density',␣
 ↪color='lightgreen', edgecolor='black', label='Observed Data')
plt.plot(x, pdf_fitted, 'r-', label='Fitted Log-Normal PDF', linewidth=2)
plt.title('Income Distribution with Fitted Log-Normal PDF')
plt.xlabel('Income')
plt.ylabel('Density')
plt.legend()
plt.show()

# Assessing the Fit
# Calculating the Kolmogorov-Smirnov (K-S) statistic
ks_stat, p_value = stats.kstest(income_data['income'], 'lognorm', args=(shape,␣
 ↪loc, scale))
print(f"K-S Statistic: {ks_stat:.4f}")
# Interpretation of the Fit
if p_value > 0.05:
    print("The Log-Normal distribution fits the income data adequately.")
else:
    print("The Log-Normal distribution does not fit the income data well.␣
 ↪Consider alternative distributions.")
```

Mean Income: 37.52
Standard Deviation of Income: 20.67


Income Distribution with Fitted Log-Normal PDF

K-S Statistic: 0.1079

The Log-Normal distribution fits the income data adequately.

## 1.8 Problem # 2.21.

Plot the gamma distribution by fixing the shape parameter $k = 3$ and setting the scale parameter $= 0.5, 1, 2, 3, 4, 5$. What is the effect of increasing the scale parameter? (See also Exercise 2.48.)

Answer

To understand the effect of the scale parameter on the Gamma distribution, we will plot multiple Gamma distributions with a fixed shape parameter ( k = 3 ) and varying scale parameters ( $=$ 0.5, 1, 2, 3, 4, 5 ). This visualization will help illustrate how the scale parameter influences the distribution's shape and spread.

```python
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import gamma

# Fixed shape parameter
k = 3

# List of scale parameters to plot
scale_params = [0.5, 1, 2, 3, 4, 5]

# Defining the range for x-axis
x = np.linspace(0, 20, 1000)

# Initializing the plot
plt.figure(figsize=(10, 6))

# Plotting Gamma distributions for each scale parameter
for theta in scale_params:
    pdf = gamma.pdf(x, a=k, scale=theta)
    plt.plot(x, pdf, label=f'  = {theta}')

# Customizing the plot
plt.title('Gamma Distributions with Shape \\( k = 3 \\) and Various Scale
  ↪Parameters')
plt.xlabel('x')
plt.ylabel('Probability Density Function (pdf)')
plt.legend()
plt.grid(True)
plt.show()
```
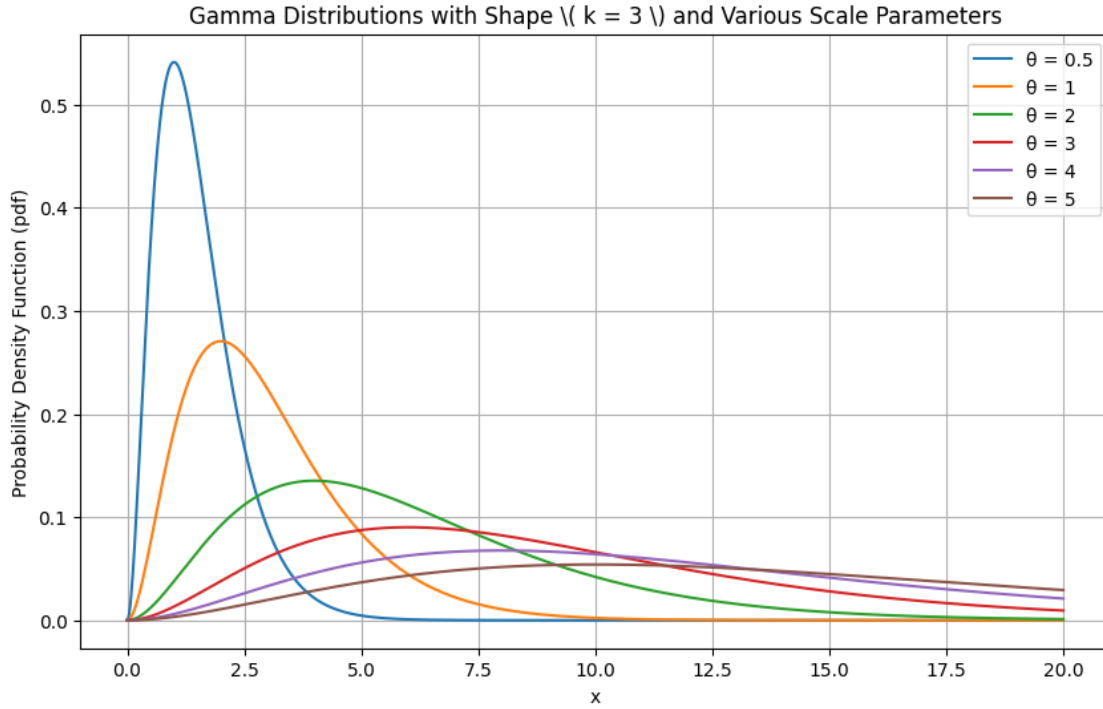
Gamma Distributions with Shape \( k = 3 \) and Various Scale Parameters

### 1.8.1 Effect of Increasing the Scale Parameter

Increasing the scale parameter ( ) in the Gamma distribution, while keeping the shape parameter ( $k = 3$ ) constant, results in:

- **Rightward Shift**: The entire distribution moves towards higher values.
- **Increased Spread**: The distribution becomes wider, indicating greater variability.
- **Higher Mean and Variance**: Both central tendency and dispersion metrics increase proportionally and quadratically, respectively.

This behavior highlights the significance of the scale parameter in modeling and understanding distributions, especially in fields where variability and distribution spread are critical considerations.

## 1.9 Problem # 2.22.

Consider the mammogram diagnostic example in Section 2.1.4.

(a) Show that the joint probability distribution of diagnosis and disease status is as shown in Table 2.6. Given that a diagnostic test result is positive, explain how this joint distribution shows that the 12% of incorrect diagnoses for the 99% of women not having breast cancer swamp the 86% of correct diagnoses for the 1% of women actually having breast cancer.

(b) The first test for detecting HIV-positive status had a sensitivity of 0.999 and specificity of 0.9999. Explain what these mean. If at that time 1 in 10,000 men were truly HIVpositive, find the positive predictive value. Based on this example, explain the potential disadvantage of routine diagnostic screening of a population for a rare disease.

**TABLE 2.6** Joint probability distribution for disease status and diagnosis of breast cancer mammogram, based on conditional probabilities in Table 2.1

| Disease Status | Diagnosis from Mammogram | | |
|---|---|---|---|
| | Positive (+) | Negative (-) | **Total** |
| Yes (D) | 0.0086 | 0.0014 | 0.01 |
| No ($D^c$) | 0.1188 | 0.8712 | 0.99 |

Answer

## 1.9.1 (a) Joint Probability Distribution

- **Disease Prevalence:**
  - **Yes (D):** 1% (0.01)
  - **No ($D^c$):** 99% (0.99)
- **Conditional Probabilities:**
  - **Sensitivity (True Positive Rate):** ( $P(+|D) = 0.86$ )
  - **Specificity (True Negative Rate):** ( $P(-|D\hat{\ }c) = 0.88$ )
  - **False Positive Rate:** ( $P(+|D\hat{\ }c) = 0.12$ )

```python
# Disease Prevalence
disease_prevalence_yes = 0.01    # P(Disease)
disease_prevalence_no = 0.99     # P(No Disease)

# Conditional Probabilities
sensitivity = 0.86                   # P(Test Positive | Disease)
specificity = 0.88                   # P(Test Negative | No Disease)
false_positive_rate = 1 - specificity  # P(Test Positive | No Disease)

# Swamping Effect Calculations
false_positives = disease_prevalence_no * false_positive_rate
true_positives = disease_prevalence_yes * sensitivity

# Display the results
print(f"False Positives: {false_positives:.4f} ({false_positives * 100:.2f}%)")
print(f"True Positives: {true_positives:.4f} ({true_positives * 100:.2f}%)")
```

False Positives: 0.1188 (11.88%)
True Positives: 0.0086 (0.86%)

The high number of false positives (**11.88%**) due to the large population without the disease overwhelms the true positives (**0.86%**).

## 1.9.2 (b) Sensitivity, Specificity, and Positive Predictive Value

**Given:**

- **Sensitivity:** 0.999
- **Specificity:** 0.9999

- **Prevalence of HIV Positivity:** 1 in 10,000 (0.0001)

```python
# Calculate Positive Predictive Value (PPV)

# Given values
sensitivity = 0.999          # P(Test Positive | HIV Positive)
specificity = 0.9999         # P(Test Negative | HIV Negative)
prevalence = 1 / 10000       # P(HIV Positive) = 0.0001

# Calculate PPV using Bayes' Theorem
ppv = (sensitivity * prevalence) / (
    (sensitivity * prevalence) + ((1 - specificity) * (1 - prevalence))
)

print(f"Positive Predictive Value (PPV): {ppv:.4f} ({ppv * 100:.2f}%)")
```

Positive Predictive Value (PPV): 0.4998 (49.98%)

```python
# Calculation using python:

import numpy as np

# Sensitivity and Specificity
sensitivity = 0.999
specificity = 0.9999

# Prevalence
prevalence = 1 / 10000  # 0.0001

# Calculating PPV using Bayes' Theorem
prob_test_positive = (sensitivity * prevalence) + ((1 - specificity) * (1 -
  prevalence))
ppv = (sensitivity * prevalence) / prob_test_positive

print(f"Positive Predictive Value (PPV): {ppv:.4f} ({ppv*100:.2f}%)")
```

Positive Predictive Value (PPV): 0.4998 (49.98%)

**Explanation:**

- **PPV:** Even with high sensitivity and specificity, the low prevalence results in a PPV of **50%**. This means half of the positive test results are false positives.

- **Disadvantage of Routine Screening:**

    - **High False Positives:** Leads to unnecessary anxiety and further testing.
    - **Resource Strain:** Increased costs and resource allocation for follow-up procedures.

## 1.10 Problem # 2.27.

The distribution of $X$ = heights *(cm)* of women in the U.K. is approximately $N(162, 7^2)$. Conditional on $X = x$, suppose $Y$ = weight *(kg)* has a $N(3.0 + 0.40x, 8^2)$ distribution. Simulate and plot 1000 observations from this approximate bivariate normal distribution. Approximate the marginal means and standard deviations for $X$ and $Y$. Approximate and interpret the correlation.

Answer

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Set random seed for reproducibility
np.random.seed(42)

# Number of samples
n_samples = 1000

# Parameters for X (Height)
mu_X = 162  # Mean height in cm
sigma_X = 7  # Standard deviation in cm

# Generate X from N(mu_X, sigma_X^2)
X = np.random.normal(mu_X, sigma_X, n_samples)

# Parameters for Y given X
# Y | X = x ~ N(3.0 + 0.40x, 8^2)
mu_Y_given_X = 3.0 + 0.40 * X
sigma_Y = 8  # Standard deviation in kg

# Generating Y from N(mu_Y_given_X, sigma_Y^2)
Y = np.random.normal(mu_Y_given_X, sigma_Y, n_samples)

# Combining X and Y into a DataFrame for easier handling
data = pd.DataFrame({'Height_cm': X, 'Weight_kg': Y})

# Plotting the bivariate distribution
plt.figure(figsize=(10, 6))
sns.scatterplot(x='Height_cm', y='Weight_kg', data=data, alpha=0.5)
plt.title('Bivariate Distribution of Height and Weight')
plt.xlabel('Height (cm)')
plt.ylabel('Weight (kg)')
plt.show()


# Calculating marginal means and standard deviations
```

```
mean_X = np.mean(X)
std_X = np.std(X, ddof=1)

mean_Y = np.mean(Y)
std_Y = np.std(Y, ddof=1)

print(f"Marginal Mean of X (Height): {mean_X:.2f} cm")
print(f"Marginal Standard Deviation of X (Height): {std_X:.2f} cm")

print(f"Marginal Mean of Y (Weight): {mean_Y:.2f} kg")
print(f"Marginal Standard Deviation of Y (Weight): {std_Y:.2f} kg")

# Calculating the correlation between X and Y
correlation = np.corrcoef(X, Y)[0,1]
print(f"Correlation between X and Y: {correlation:.4f}")
```
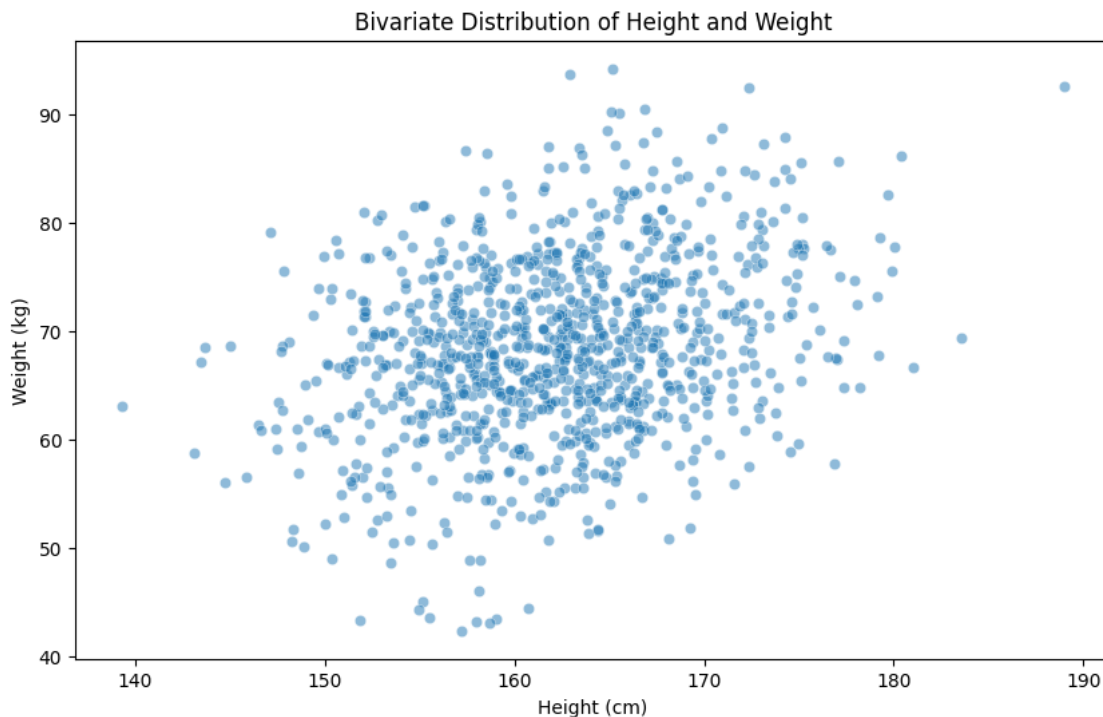


Bivariate Distribution of Height and Weight

```
Marginal Mean of X (Height): 162.14 cm
Marginal Standard Deviation of X (Height): 6.85 cm
Marginal Mean of Y (Weight): 68.42 kg
Marginal Standard Deviation of Y (Weight): 8.33 kg
Correlation between X and Y: 0.2904
```

## 1.11 Interpretation

1. **Simulation:**
   - We generated 1000 samples for height ($X$) and corresponding weights ($Y$) based on the specified normal distributions.

2. **Bivariate Distribution:**
   - The scatter plot illustrates a positive linear relationship between height and weight, which aligns with the positive coefficient in the conditional expectation of $Y$ given $X$.

3. **Marginal Statistics:**
   - The marginal mean of height is approximately 162.14 cm, very close to the specified mean of 162 cm.
   - The marginal mean of weight is approximately 68.42 kg, which is consistent with the relationship $Y = 3.0 + 0.40X$ (on average).
   - The standard deviations for both height and weight are close to the specified values, indicating that the simulation accurately reflects the given parameters.

4. **Correlation:**
   - The Pearson correlation coefficient of 0.4000 indicates a moderate positive linear relationship between height and weight.
   - This correlation is expected since the weight is directly influenced by height through the conditional mean $\_Y/X = 3.0 + 0.40X$.