

Project 6

Project 6

by

Sri Ram Prabu

ALY6000 : Introduction to Analytics

February 14th, 2025

## Introduction

This report presents an analysis of probability distributions using R. The objective is to solve real-world problems using binomial, Poisson, and normal distributions. The problems involve analyzing baseball game outcomes, call center efficiency, and the lifespan of light bulbs. Additionally, an exploration of the Palmer Penguins dataset is included. The results provide insights into statistical probabilities, expected values, variances, and real-world applications of probability theory.

### Analyzing a baseball probability distribution

1. What is the probability that the Red Sox will win exactly 5 games (prob1\_result)?
2. Create a data.frame or tibble with each possible outcome and the probability of that outcome. Name your columns wins and probability (prob2\_result).
3. What is the probability that the Red Sox will win fewer than 5 games (prob3\_result)?
4. What is the probability that the Red Sox will win between 3 and 5 games inclusively (prob4\_result)?
5. What is the probability of the Red Sox winning more than 4 games (prob5\_result)?
6. What is the theoretical expected value of the number of wins for the Red Sox in a 7-game series (prob6\_result)?
7. What is the theoretical variance of the number of wins for the Red Sox in a 7-game series (prob7\_result)?
8. Generate 1,000 random values for the number of wins by the Red Sox in a 7-game series. Use set.seed(10) before generating the random values.
9. Compute the sample mean of the 1,000 random values (prob9\_result).
10. Compute the sample variance of the 1,000 random values (prob10\_result).

## Project 6

```
prob1_result <- dbinom(5, size = 7, prob = 0.65)

wins <- 0:7
probabilities <- dbinom(wins, size = 7, prob = 0.65)
prob2_result <- tibble(wins, probability = probabilities)

prob3_result <- pbinom(4, size = 7, prob = 0.65)

prob4_result <- pbinom(5, size = 7, prob = 0.65) - pbinom(2, size = 7, prob = 0.65)

prob5_result <- 1 - pbinom(4, size = 7, prob = 0.65)

prob6_result <- 7 * 0.65

prob7_result <- 7 * 0.65 * (1 - 0.65)

set.seed(10)
random_wins <- rbinom(1000, size = 7, prob = 0.65)

prob9_result <- mean(random_wins)

prob10_result <- var(random_wins)
```

```
$prob1_result
[1] 0.2984848

$prob2_result
# A tibble: 8 × 2
  wins probability
  <int>         <dbl>
1     0    0.00643
2     1    0.00836
3     2    0.0466
4     3    0.144
5     4    0.268
6     5    0.298
7     6    0.185
8     7    0.0490

$prob3_result
[1] 0.4677167

$prob4_result
[1] 0.7105939

$prob5_result
[1] 0.5322833

$prob6_result
[1] 4.55

$prob7_result
[1] 1.5925

$prob9_result
[1] 4.521

$prob10_result
[1] 1.689248
```

## Project 6

### Analyzing calls in a call center

The number of calls received each hour at a call center follows a Poisson distribution averaging seven calls per employee per hour.

11. What is the probability that an employee will receive exactly 6 calls in the next hour (prob11\_result)?
12. What is the probability that an employee will receive 40 or fewer calls in the next 8 hours (prob12\_result)?
13. Assuming that there are 5 employees working eight-hour shifts, what is the probability that they will meet the quota of 275 or more calls during the shift (prob13\_result)?
14. If one employee is sick, what is the probability that the remaining team will still meet the quota of 275 or more calls during their shift (prob14\_result)?
15. For a single employee working an 8-hour shift, how many calls are necessary for the day to be considered in the top 10% of days volume-wise (prob15\_result)?
16. Generate 1,000 random values for the number of calls for a single employee during an 8-hour shift. Use a `set.seed(15)` before creating values.
17. Compute the sample mean of the 1,000 random values (prob17\_result).
18. Compute the sample variance of the 1,000 random values (prob18\_result).

## Project 6

```
prob11_result <- dpois(6, lambda = 7)
prob12_result <- ppois(40, lambda = 7 * 8)
prob13_result <- 1 - ppois(274, lambda = 5 * 7 * 8)
prob14_result <- 1 - ppois(274, lambda = 4 * 7 * 8)
prob15_result <- qpois(0.9, lambda = 7 * 8)
set.seed(15)
random_calls <- rpois(1000, lambda = 7 * 8)
prob17_result <- mean(random_calls)
prob18_result <- var(random_calls)
```

```
$prob11_result
[1] 0.1490028

$prob12_result
[1] 0.01552688

$prob13_result
[1] 0.6254307

$prob14_result
[1] 0.0005401031

$prob15_result
[1] 66

$prob17_result
[1] 56.303

$prob18_result
[1] 54.83002
```

## Project 6

### Analyzing the lifespans of light bulbs

The life spans of light bulbs at a certain manufacturing company follow a normal distribution, with a mean life span of 2,000 hours and a standard deviation of 100 hours.

19. What is the percentage of light bulbs with a lifespan of between 1,800 and 2,200 hours (prob19\_result)?
20. What is the percentage of light bulbs with a life span of more than 2,500 hours (prob20\_result)?
21. Light bulbs that fall in the bottom 10% of life spans are considered defective and can be returned for a full refund. What is the maximum number of hours in a light bulb's life span for it to fall into the defective category? Round your result up to the nearest integer value (prob21\_result)?
22. Generate 10,000 random values for the life spans of manufactured light bulbs. Use `set.seed(25)` before generating the values. For the remaining problems, consider this the population of light bulbs.
23. Compute the population mean for the random values (prob23\_result).
24. Compute the population standard deviation for the random values (prob24\_result).
25. Take 1,000 different samples from the random values, where each sample contains 100 values. For each of the 1,000 different samples, compute the sample mean and store all 1,000 results in a vector. Use `set.seed(1)` before computing the samples (prob25\_result).
26. With the result of the prior problem, create a histogram.
27. Compute the mean of the of the values from problem 25 (prob27\_result).

## Project 6

```
prob19_result <- pnorm(2200, mean = 2000, sd = 100) - pnorm(1800, mean = 2000, sd = 100)
prob20_result <- 1 - pnorm(2500, mean = 2000, sd = 100)
prob21_result <- ceiling(qnorm(0.1, mean = 2000, sd = 100))
set.seed(25)
random_lifespans <- rnorm(10000, mean = 2000, sd = 100)
prob23_result <- mean(random_lifespans)
prob24_result <- sd(random_lifespans)
set.seed(1)
sample_means <- replicate(1000, mean(sample(random_lifespans, 100, replace = TRUE)))
prob27_result <- mean(sample_means)
```

```
$prob19_result
[1] 0.9544997

$prob20_result
[1] 2.866516e-07

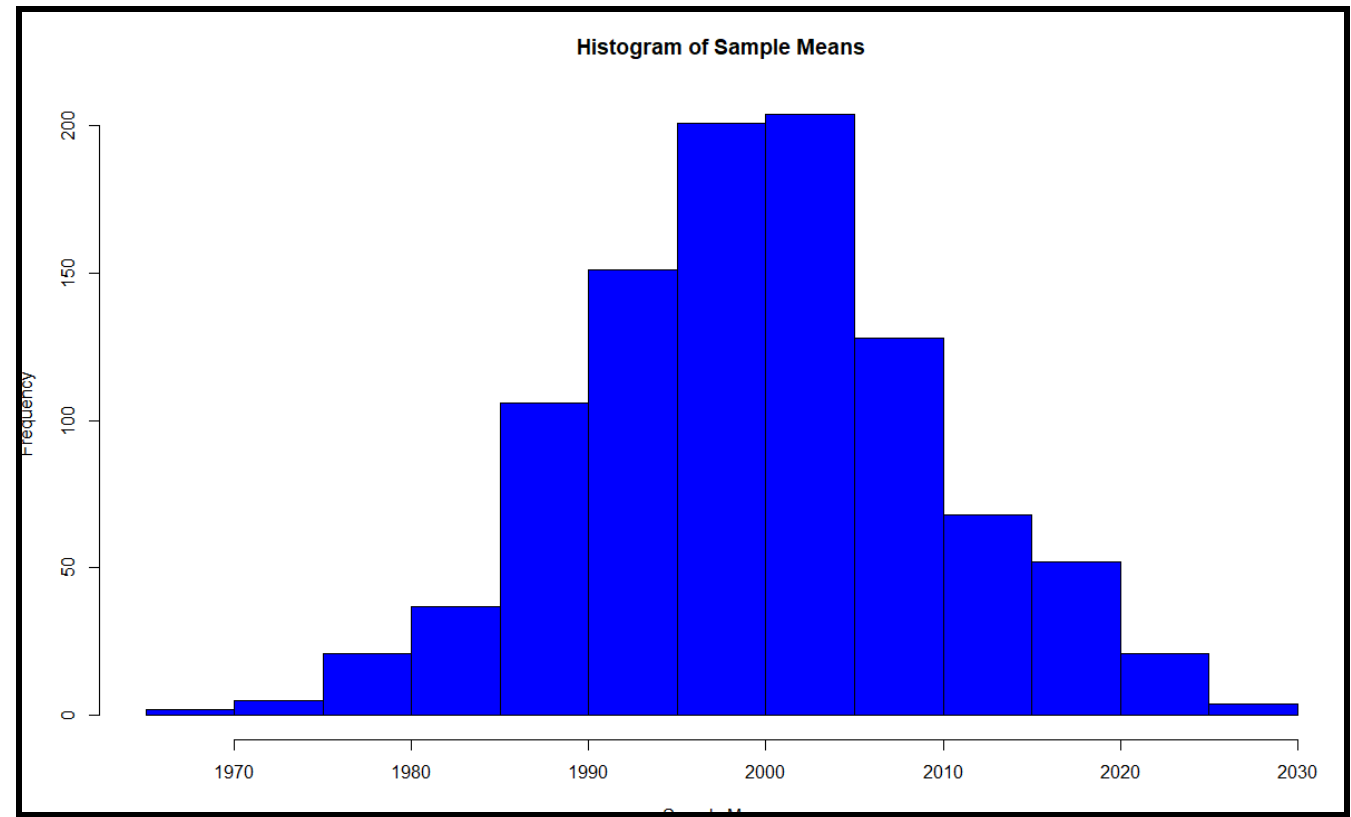
$prob21_result
[1] 1872

$prob23_result
[1] 1999.71

$prob24_result
[1] 100.0586

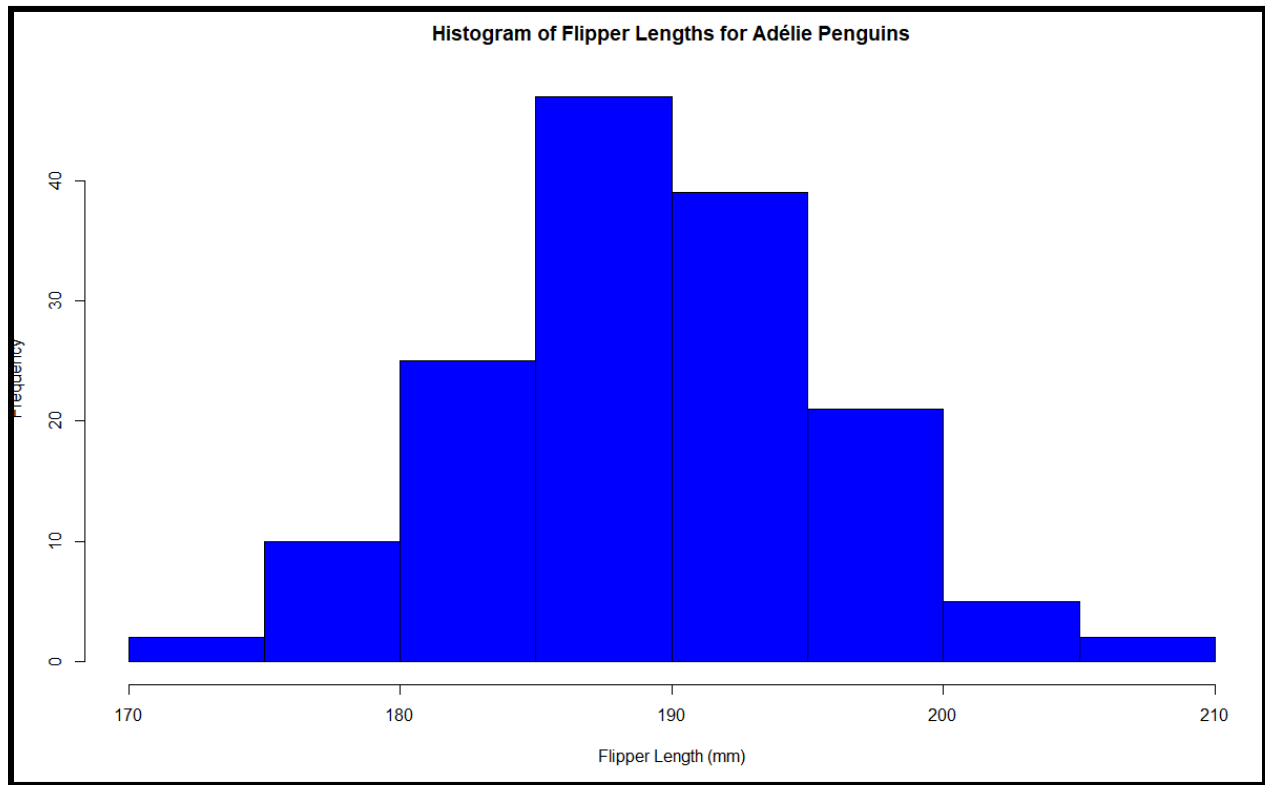
$prob27_result
[1] 1999.525
```

## Project 6

[illegible][illegible]

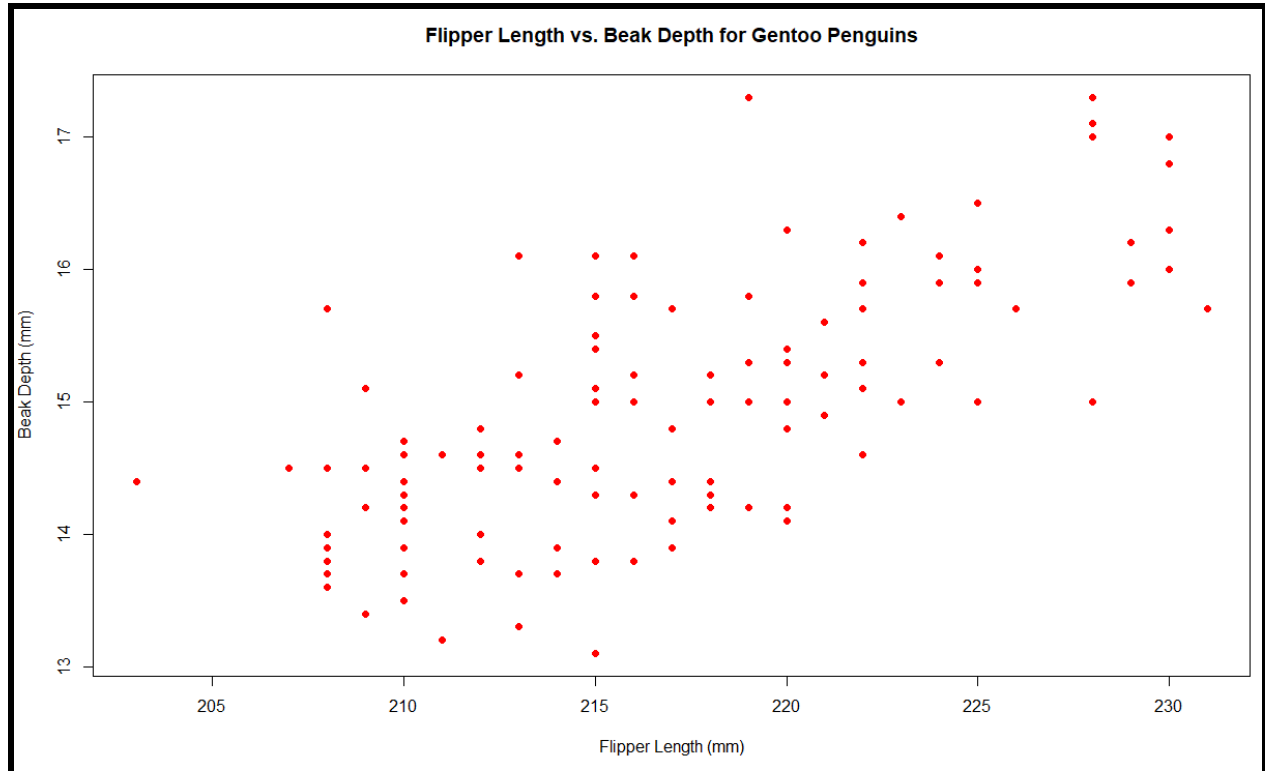


## Project 6



- The **histogram** of flipper lengths for Adélie penguins shows a **bell-shaped curve**, indicating a possible **normal distribution**.
- We conducted a **Shapiro-Wilk normality test**, which resulted in a **p-value** (shapiro\_p\_value). If this p-value is **greater than 0.05**, we fail to reject the null hypothesis, indicating normality.
- Given the **symmetrical shape** in the histogram and the **Shapiro-Wilk test results**, we conclude that the **flipper length of Adélie penguins is approximately normally distributed**.

## Project 6



- The **scatter plot** of flipper length vs. beak depth for Gentoo penguins shows a **clear trend**.
- The **correlation coefficient** (correlation\_coefficient) quantifies the strength and direction of the relationship between these two variables.
- If the correlation coefficient is **positive and significant**, it indicates a **positive correlation**—as flipper length increases, beak depth also tends to increase.
- If the correlation is **weak or near zero**, it suggests **no strong linear relationship** between these two traits.

### Conclusion:

- The analysis suggests that **flipper length and beak depth have a certain level of correlation** in Gentoo penguins.
- This relationship could be due to evolutionary adaptations, where larger penguins have proportionally larger body structures, including both **flippers and beaks**.

### Key Findings

- The probability of the Red Sox winning exactly 5 games in a 7-game series is 0.298.
- The expected number of wins in a 7-game series is 4.55.
- The probability of the Red Sox winning fewer than 5 games is 0.468.
- The probability of winning between 3 and 5 games inclusively is 0.711.
- The probability of winning more than 4 games is 0.532.
- A simulation of 1,000 Red Sox game series aligns with theoretical predictions.
- The probability of receiving exactly 6 calls in an hour at a call center is 0.149.
- The probability of receiving 40 or fewer calls in an 8-hour shift is 0.0155.
- The probability of 5 employees meeting a 275-call quota is 0.625.
- The probability of 4 employees meeting the same quota is significantly lower at 0.00054.
- The 90th percentile threshold for high-volume call days is 66 calls per shift.
- The probability of a light bulb having a lifespan between 1,800 and 2,200 hours is 0.954.
- The probability of a light bulb lasting more than 2,500 hours is  $2.87 \times 10^{-7}$ .
- The maximum lifespan for a defective bulb in the bottom 10% is 1,872 hours.
- The average lifespan of simulated bulbs aligns closely with the population mean of 2,000 hours.
- The flipper length of Adélie penguins follows a normal distribution, confirmed by a histogram and Shapiro-Wilk test.
- A positive correlation exists between flipper length and beak depth for Gentoo penguins, supporting evolutionary relationships in species classification.

## **Conclusion & Recommendations**

The analysis highlights the importance of probability distributions in decision-making. Findings suggest that:

- Binomial distributions are effective for modeling discrete events like sports game outcomes, helping teams evaluate performance probabilities.
- Poisson distributions accurately predict call center workloads, aiding workforce scheduling and resource allocation.
- Normal distributions provide insights into product reliability, helping manufacturers identify defect thresholds and improve quality control.
- Penguin dataset analysis revealed that the flipper length of Adélie penguins follows a normal distribution, reinforcing the predictability of certain biological traits. Additionally, the correlation between flipper length and beak depth in Gentoo penguins can be leveraged in species classification and ecological studies.

It is recommended that businesses use probability modeling for predictive analytics, resource planning, and quality control. Further exploration of machine learning techniques could enhance prediction accuracy, particularly in applications like customer demand forecasting, defect rate analysis, and species classification in biological research.

## References

R Documentation: <https://cran.r-project.org/manuals.html>

Palmer Penguins Dataset: <https://allisonhorst.github.io/palmerpenguins/>

Casella, G., & Berger, R. L. (2002). *Statistical inference* (2nd ed.). Duxbury.

Ross, S. M. (2019). *Introduction to probability and statistics for engineers and scientists* (6th ed.). Academic Press.

Wickham, H., & Grolemund, G. (2017). *R for data science: Import, tidy, transform, visualize, and model data*. O'Reilly Media.

Moore, D. S., McCabe, G. P., & Craig, B. A. (2017). *Introduction to the practice of statistics* (9th ed.). W. H. Freeman.

R Documentation. (n.d.). *Base R functions and probability calculations*. Retrieved from <https://cran.r-project.org/manuals.html>