



---

## Exploring Biodiversity Metrics

---

# A Comprehensive Statistical Analysis of the Palmer Penguins Dataset



APRIL 18, 2024

JOHNSON & WALES UNIVERSITY  
DATA5100: STATISTICAL ANALYSIS  
PROFESSOR ANN BRETT, PHD  
ANDREX IBIZA, MBA

## Abstract

This research delves into the intricate patterns of biodiversity and conservation-relevant metrics within the penguin populations inhabiting Antarctica's Palmer Archipelago. The study leverages the Palmer Penguins dataset, which encompasses a wealth of information on three penguin species - Adélie, Chinstrap, and Gentoo - residing on three distinct islands within the archipelago. By employing a diverse range of statistical techniques, this research endeavors to uncover the intricate relationships among various biological markers, such as species type, bill dimensions, flipper length, and body mass. The aim is to provide a deeper understanding of the ecological dynamics and adaptive traits that characterize these species, thereby contributing to ongoing conservation initiatives and enriching our broader comprehension of Antarctic marine ecosystems.

## 1. Introduction

The Palmer Penguins dataset has become an indispensable tool for ecologists and data scientists seeking to unravel the complexities of biodiversity within the Antarctic region. This dataset, meticulously compiled by Horst AM, Hill AP, and Gorman KB, offers a comprehensive collection of measurements from three penguin species - Adélie, Chinstrap, and Gentoo - found on three islands nestled within the Palmer Archipelago, Antarctica. The dataset encompasses a variety of variables, including species, island, bill length, bill depth, flipper length, body mass, and sex, providing a unique window into the morphological diversification shaped by ecological and evolutionary pressures.

## Literature Review on the Palmer Penguins Dataset

The Palmer Penguins dataset has increasingly become a focal point for statistical analysis and machine learning within the biological and ecological research fields. This literature review synthesizes findings from two specific studies to demonstrate the dataset's utility and its application in both educational and ecological contexts.

### Educational Use in Data Science

Horst, Hill, and Gorman (2022) provide a comprehensive evaluation of the Palmer Penguins dataset within the context of statistical teaching and machine learning education. Their study, published in *The R Journal*, emphasizes the dataset as an exemplary tool for data science education, serving as a more complex and illustrative alternative to the well-known Iris dataset. They argue that the multidimensional nature of the Palmer Penguins dataset introduces learners to real-world data challenges, enhancing their analytical skills and understanding of statistical models. The accessibility of the dataset through the `palmerpenguins` R package simplifies its integration into educational frameworks, making it a valuable resource for developing practical data handling and visualization skills.

### Ecological Insights

Trathan et al. focus on the ecological implications of human activities on penguin populations, specifically the gentoo penguins (*Pygoscelis papua*) at Goudier Island in the Palmer Archipelago. Published in *Biological Conservation*, their research assesses the population dynamics of gentoo penguins in response to tourism, highlighting the dataset's relevance to ecological and conservation studies. By using long-term observational data, the study provides critical insights into the impacts of tourism on breeding success and population stability. This research underlines the importance of the Palmer Penguins dataset in understanding and managing wildlife interactions in sensitive ecological zones, offering a base for policy-making and conservation strategies.

The Palmer Penguins dataset not only facilitates advanced statistical education but also plays a crucial role in ecological research by providing in-depth insights into the environmental impacts on Antarctic wildlife. These studies exemplify the dataset's dual utility in both academic and practical applications, underscoring its growing importance in diverse research and educational settings.

## 2. Data Loading, Cleaning, and Exploratory Data Analysis

The initial phase of this research involved loading the data into a dataframe and conducting a preliminary examination of its overarching characteristics. The raw dataset, named "penguins\_raw" and sourced from the "palmerpenguins" package, contained data on 344 penguin observations across 17 variables. These variables encompassed a wide spectrum of information, including study name, sample number, species, region, island, stage, individual ID, clutch completion, date egg, culmen length and depth, flipper length, body mass, sex, delta 15 N, delta 13 C, and comments.

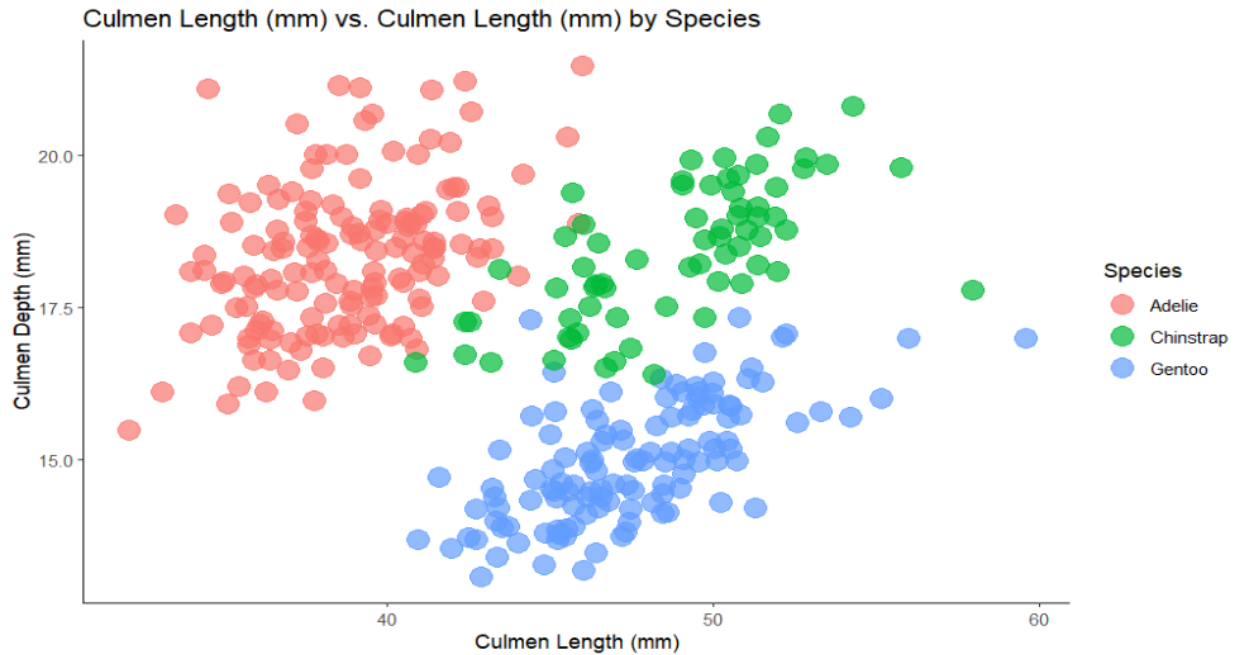
To ensure clarity and consistency, the variable names were meticulously reformatted, replacing spaces and special characters with underscores and converting them to lowercase. The variable "Sample Number" was transformed into an integer and renamed "sample\_number". Categorical variables such as species, region, island, stage, clutch completion, and sex were converted into factors, while numerical variables were maintained in their original format. Careful inspection revealed missing values in several variables, most notably in the "comments" variable, which had a significant proportion of missing data. Subsequent analyses were conducted on the complete cases to maintain data integrity and reliability.

### 3. Exploratory Data Visualizations

Visualizing the data offers a powerful means to explore the lives and ecological roles of the three penguin species residing within the Palmer Archipelago. Each visualization unveils a different facet of their unique adaptations and the intricate relationships between their physical characteristics and their environment.

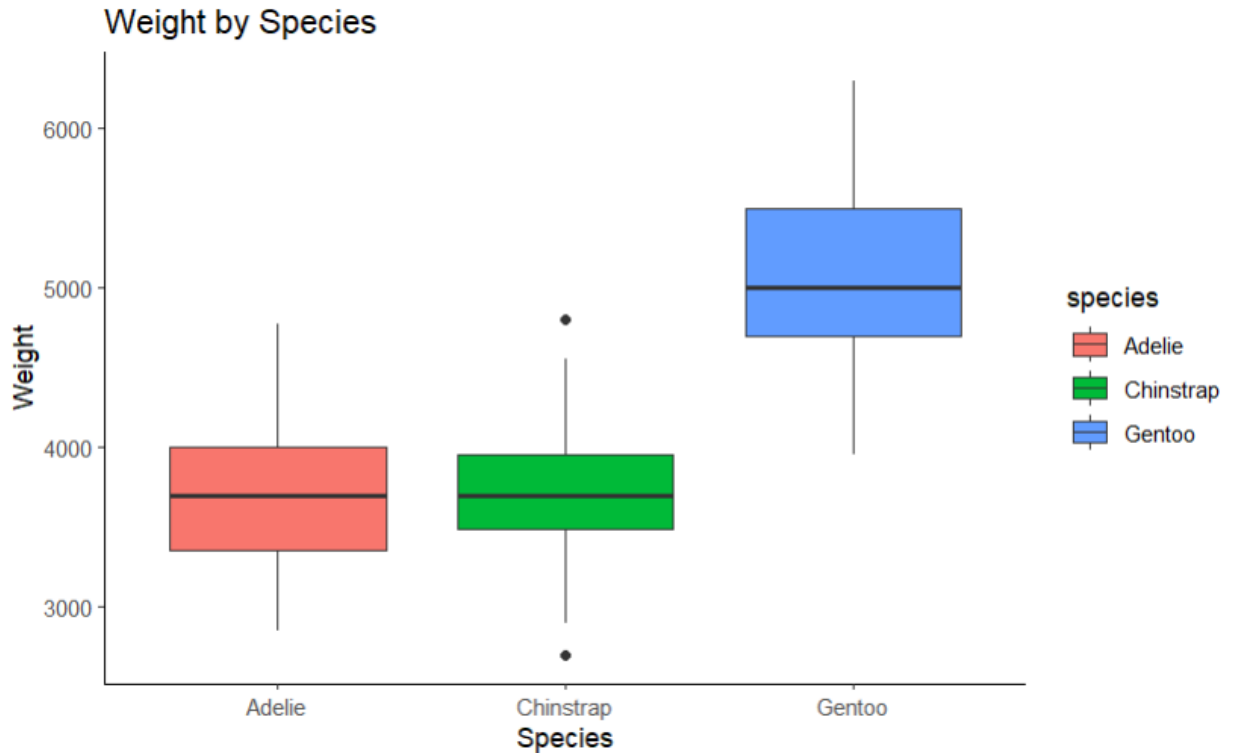
#### 3.1. Categorical Data Visualizations

The scatter plot examining the relationship between culmen length and depth (Figure 1) acts as a window into the diverse feeding strategies employed by the three penguin species. Adélie penguins, characterized by their shorter and deeper culmens, likely utilize a different foraging approach compared to Gentoo penguins, who possess longer and shallower beaks. This observation suggests that Adélie penguins might specialize in capturing prey that requires a stronger bite force, while Gentoo penguins, with their more slender beaks, could be adept at pursuing smaller, more agile prey.

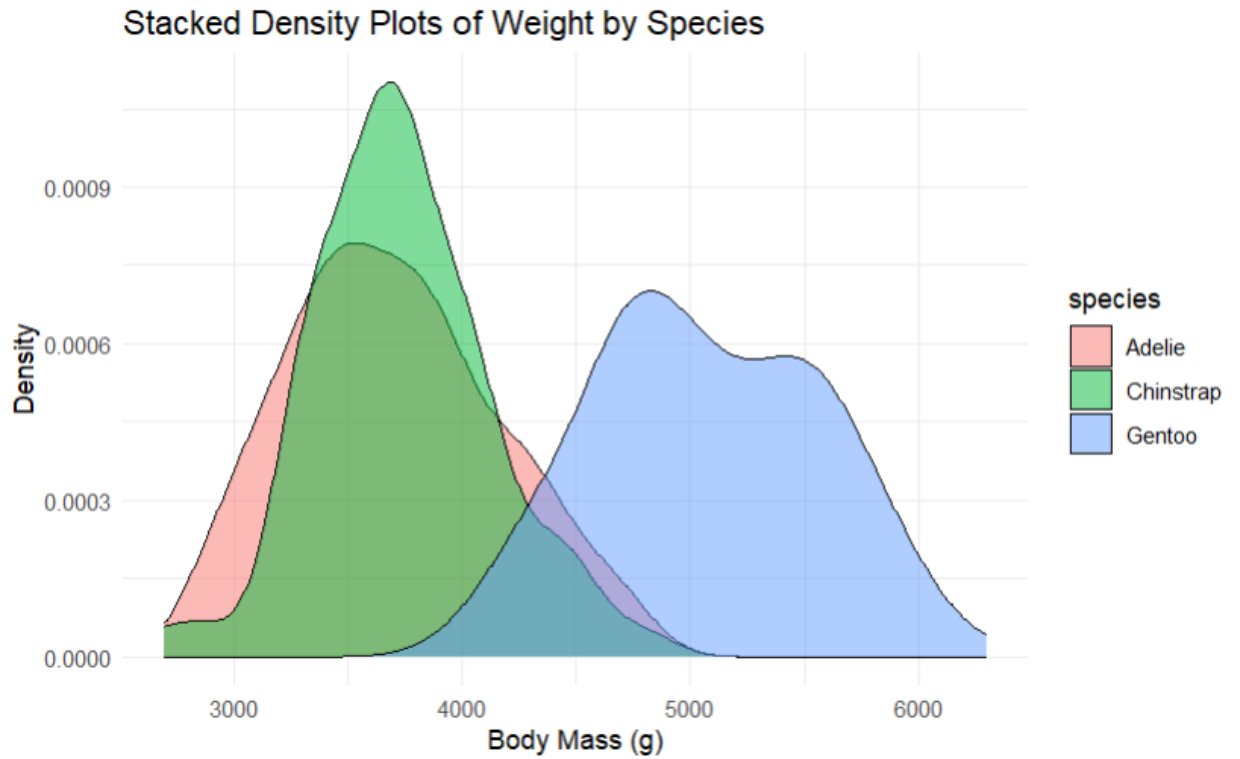
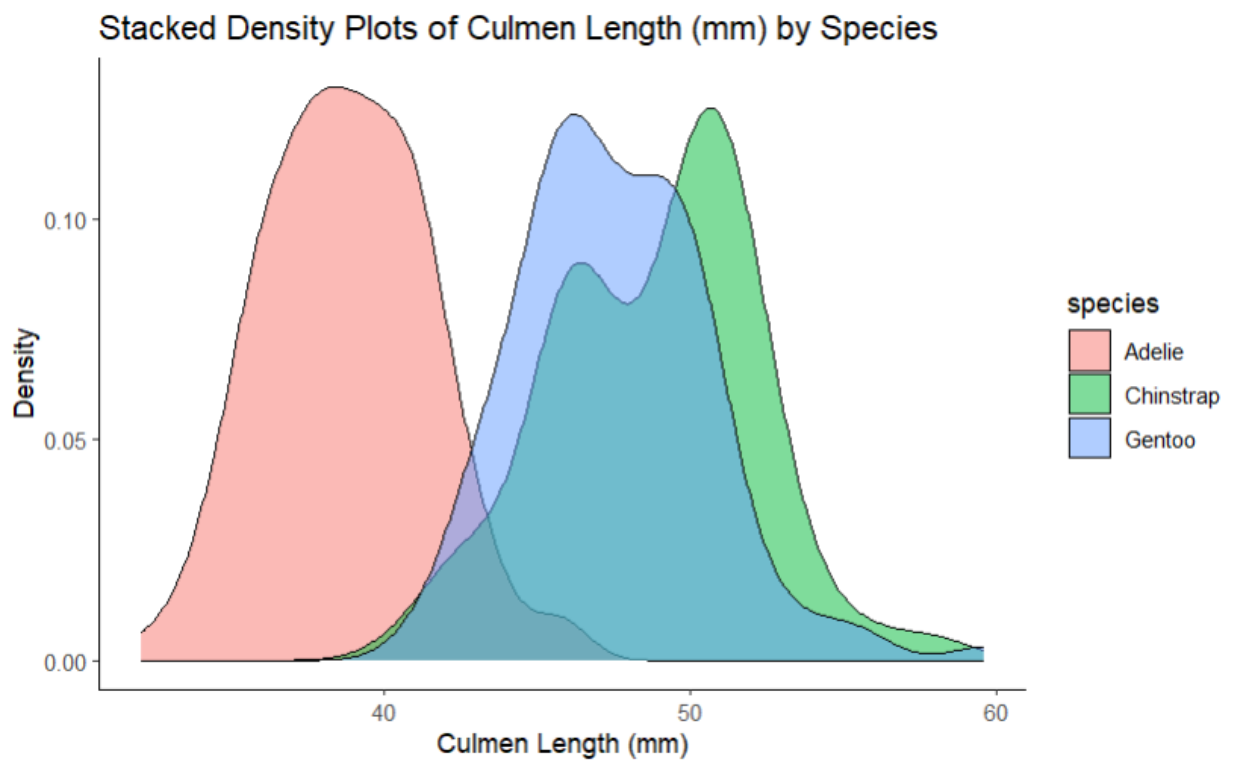


**Figure 1**

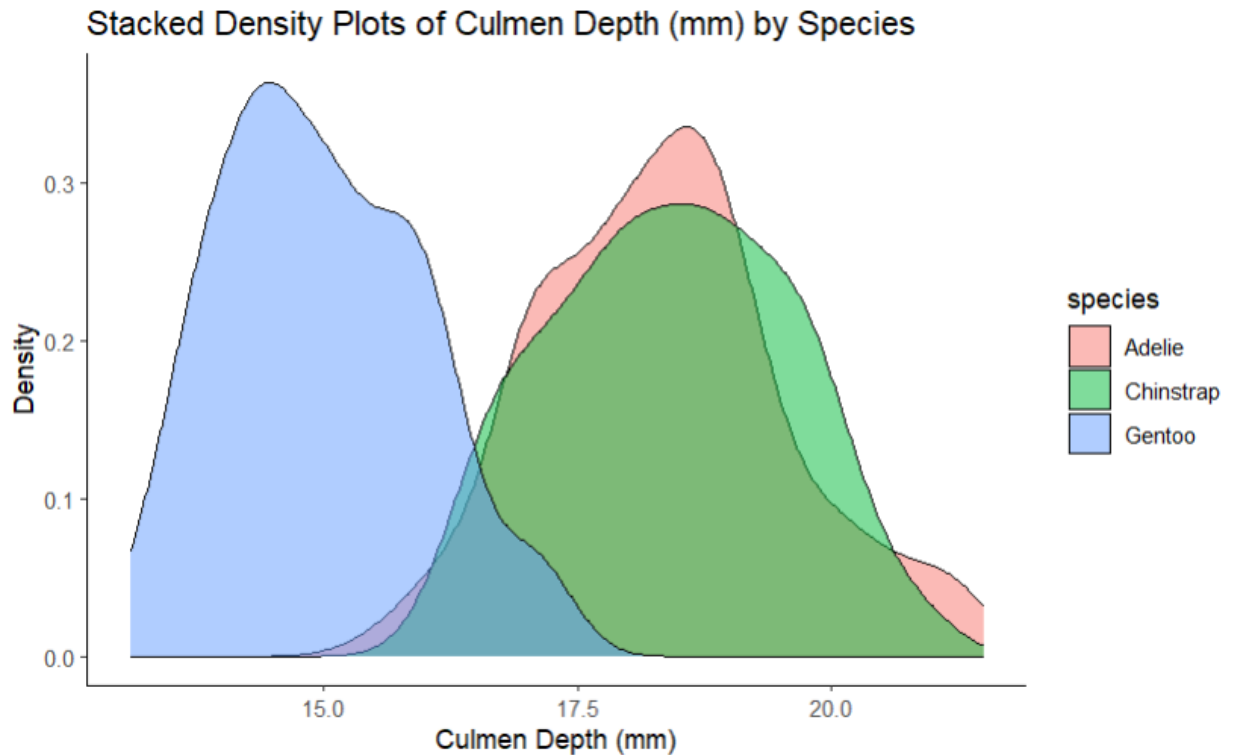
Further evidence of morphological diversity emerges from the box plot of body mass by species in Figure 2. Gentoo penguins, with their significantly higher median body mass, likely have different energy requirements and thermoregulation strategies compared to the smaller Adélie and Chinstrap penguins. The broader range of body mass within the Gentoo species, as indicated by the wider box and whiskers, could reflect variations in diet, age, or individual fitness. The presence of outliers in the Chinstrap data also raises intriguing questions about individual variations or potential data recording anomalies that require further scrutiny.

**Figure 2**

Stacked density plots offer a deeper dive into the factors shaping the penguins' physical characteristics. By visualizing the distributions of weight (Figure 3), culmen length (Figure 4), and culmen depth (Figure 5) for each species, we can discern the interplay of shared environmental influences and species-specific evolutionary traits. The overlapping distributions observed in these plots, particularly for weight, suggest that a common environmental factor, such as food availability or oceanographic conditions, might be influencing body mass across all three species.

**Figure 3****Figure 4**



**Figure 5**

However, the distinct peaks within each distribution, especially for culmen length, underscore the presence of unique evolutionary adaptations. Each species appears to have developed specialized beak structures that optimize their foraging efficiency and competitive success within their respective ecological niches.

The segmented bar plot depicting species distribution by island (Figure 6) sheds light on the penguins' habitat preferences and potential interspecies competition for resources. The dominance of Gentoo penguins on Biscoe Island, coupled with the scarcity of Chinstrap penguins on this island, suggests that Gentoo penguins might outcompete Chinstrap penguins for resources in this specific environment. Similarly, the high concentration of Adélie penguins on Torgersen Island and their lower numbers on Biscoe Island might indicate a preference for specific nesting sites or food sources available on Torgersen. The presence of Chinstrap penguins on Dream Island,

but their absence on the other two islands, further emphasizes the intricate relationship between species distribution and habitat suitability.

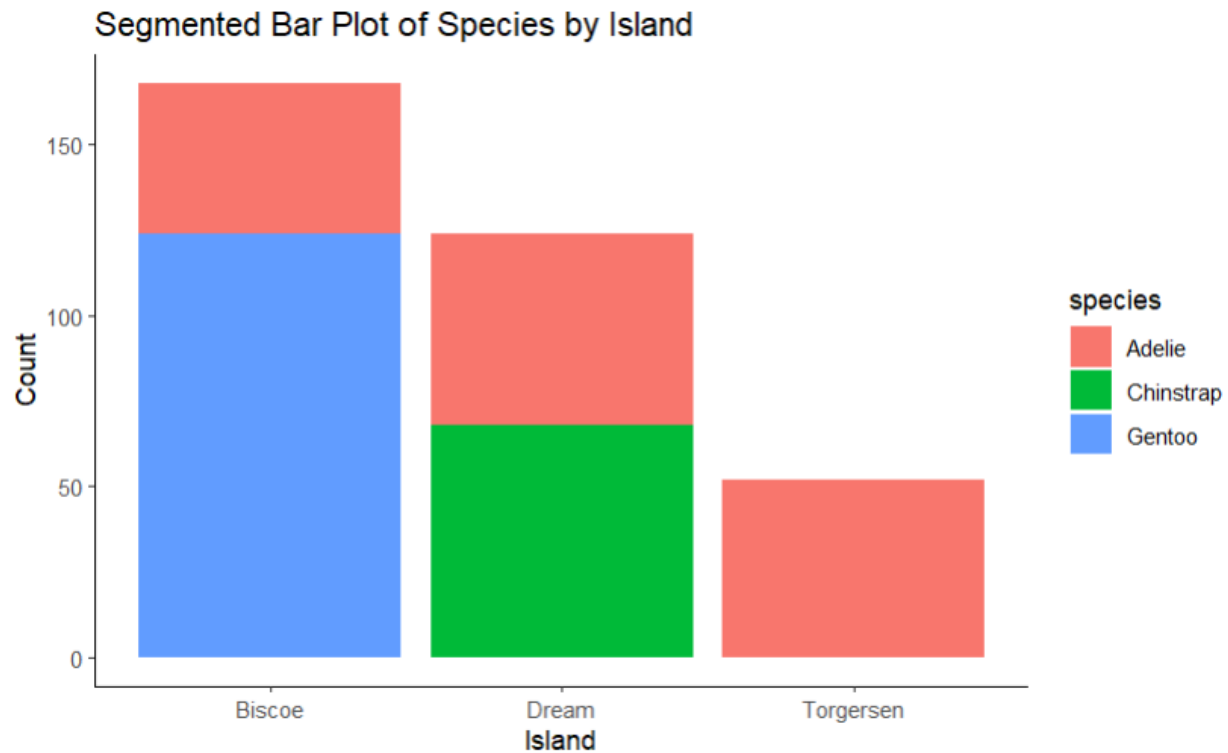
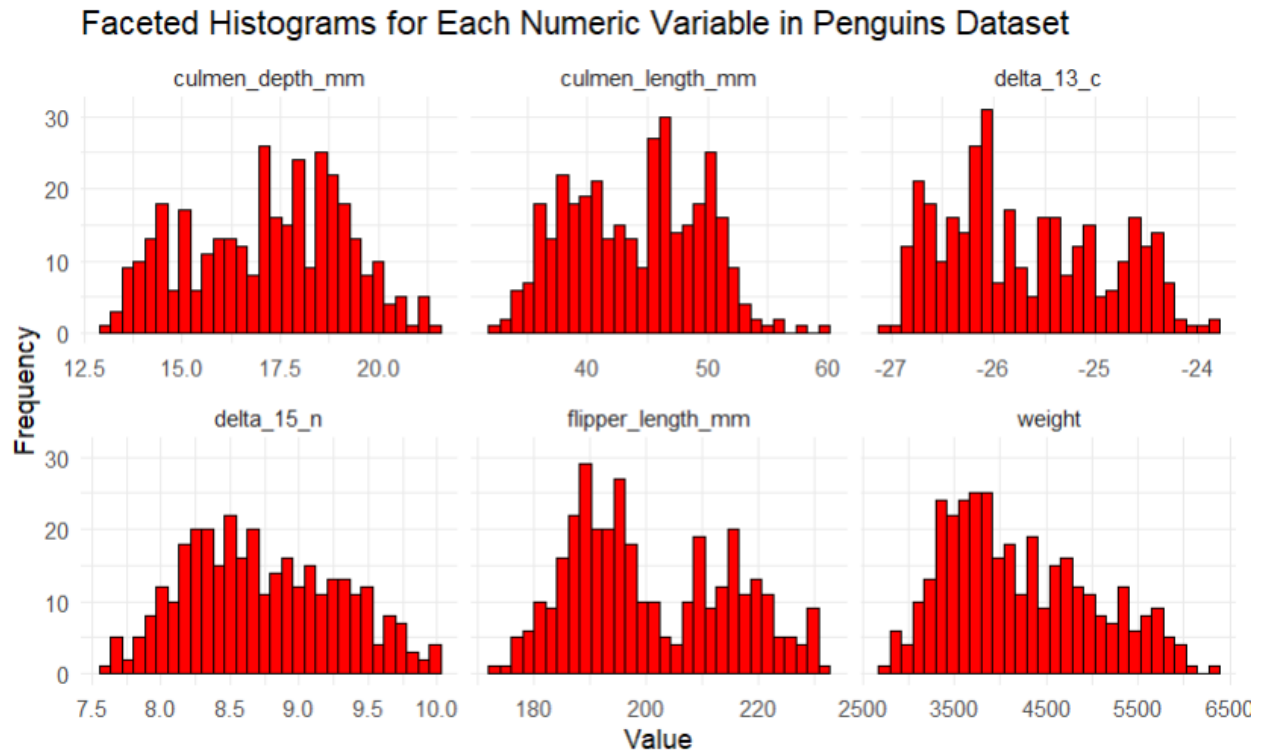


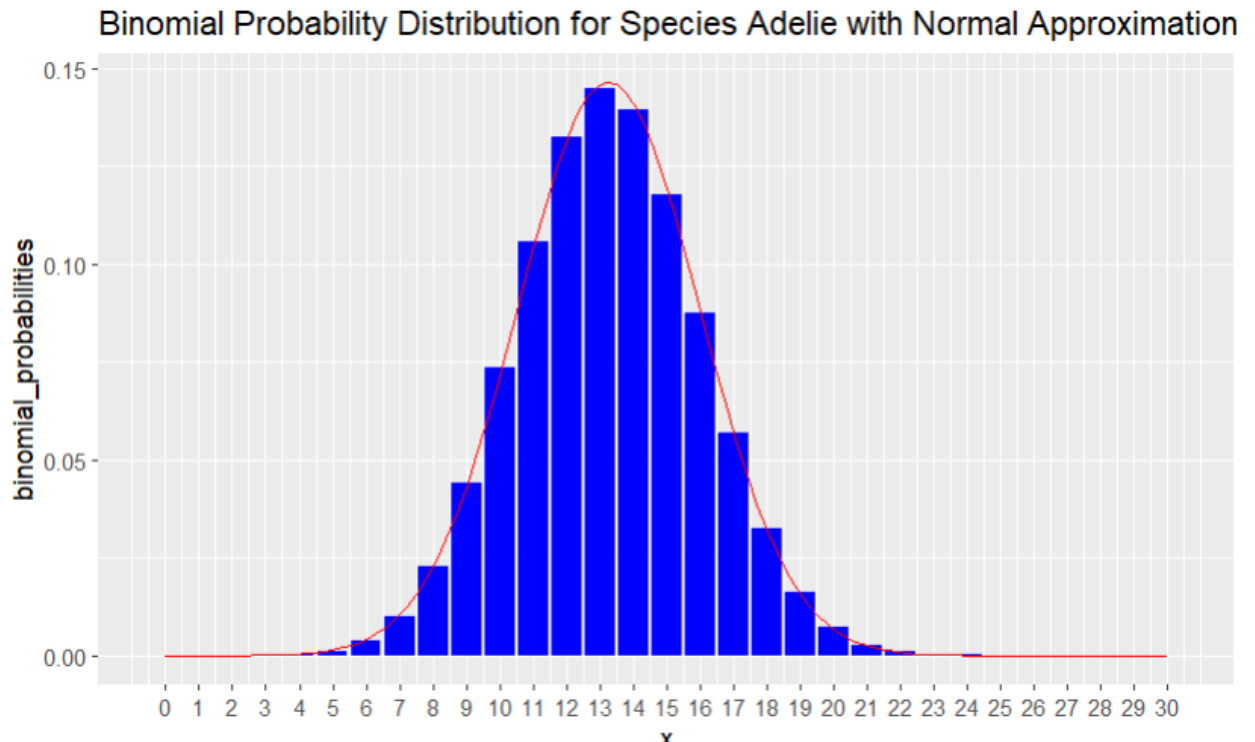
Figure 6

### 3.2. Numerical Data Visualization

Histograms (see Figure 7) were created for each numerical variable, including body mass, culmen depth and length, delta 13 C, delta 15 N, and flipper length. These histograms revealed non-normal distributions, with many exhibiting bimodal or multimodal patterns, suggesting the presence of distinct subgroups within the data. This observation indicated that the empirical rule, which assumes a normal distribution, would not be appropriate for analyzing these variables.

**Figure 7**

Density plots with superimposed normal distributions (Figure 8) further confirmed the non-normality of the data, particularly for weight. This finding highlighted the limitations of the empirical rule for this dataset and emphasized the need for alternative statistical approaches.

**Figure 8**

## 4. Confidence Intervals and Hypothesis Tests

For this stage of the analysis, the penguin weights were treated as a population, and a random sample of 100 penguin weights was drawn. This sample was then used to construct a z confidence interval and conduct a z hypothesis test to estimate and assess the population mean weight.

The 95% confidence interval calculated for the population mean weight was (4146, 4482). Given that the population standard deviation was approximately 802 grams, it was concluded with 95% confidence that the true population mean weight of penguins in the Palmer Penguins dataset falls within this range.

A two-tailed z hypothesis test was subsequently conducted with a significance level of 0.05 to evaluate the claim that the population mean weight is 4202 grams. The test yielded a test

statistic of  $z = 0.277138$  and a corresponding p-value of  $p = 0.7816742$ . As the p-value exceeded the significance level, the null hypothesis was not rejected. Consequently, there was insufficient statistical evidence to refute the claim that the population mean weight is indeed 4202 grams.

To validate these findings, the R function "z.test" was employed, and the results were consistent with the manually calculated confidence interval and hypothesis test outcomes.

## 5. Analysis of Variance (ANOVA) and Post-Hoc Testing

In this section, the weights of the penguins were again considered as a population, and an analysis of variance (ANOVA) was conducted to investigate potential differences in weight among the three penguin species.

The initial ANOVA test revealed a significant difference in weight among the species. However, as Levene's test indicated unequal variances across the groups, three alternative methods were explored to address this issue. These methods included Welch's ANOVA, Box-Cox transformation, and log transformation. Remarkably, all three methods led to the same conclusion: the null hypothesis of equal means should be rejected, signifying a statistically significant difference in weight among the three penguin species.

To further explore these differences, Tukey's Honestly Significant Difference (HSD) test was employed as a post-hoc analysis. The results of this test indicated that Gentoo penguins had a significantly higher mean weight compared to both Adélie and Chinstrap penguins. However, there was no statistically significant difference in mean weight between Chinstrap and Adélie penguins.

## 6. Palmer Penguins Species Proportion Analysis

This analysis focused on examining the proportion of Adélie penguins within the dataset, with the initial hypothesis assuming that the proportion within a sample of 100 penguins would reflect the overall population proportion. The population proportion of Adélie penguins was established as  $p = 0.44$ .

A random sample of 100 penguins was drawn, and the sample proportion of Adélie penguins was calculated to be 0.58. A 95% confidence interval for the population proportion was constructed using the exact method of the `binom.confint()` function, resulting in an interval of (0.477, 0.678). This interval suggests that we can be 95% confident that the true population proportion of Adélie penguins lies within this range.

A two-tailed hypothesis test was then conducted to compare the sample proportion against the hypothesized population proportion of 0.44. The z-test yielded a test statistic of  $z = 0.0695$  with a corresponding p-value of  $p = 0.9446$ . As the p-value was greater than the significance level of 0.05, the null hypothesis was not rejected. This implies that the data does not provide sufficient evidence to support a significant deviation from the hypothesized population proportion.

Further analysis using the `prop.test` function corroborated these findings, with a chi-squared statistic close to zero and a high p-value ( $p = 0.9879$ ). Additionally, the 95% confidence interval generated using the Wald method, although slightly narrower, remained consistent with the previous interval in suggesting that the hypothesized proportion is plausible given the sample data.

Overall, the analysis suggests that the data does not provide compelling evidence to reject the hypothesis that the population proportion of Adélie penguins is 0.44. Therefore, the assumption

that approximately 44% of the penguin population within the dataset consists of Adélie penguins remains valid.

## 7. Analysis of Penguin Species Distribution and Island Association

To further explore the distribution of penguin species and their association with island habitats, two chi-square tests were conducted. The first, a Chi-Square Goodness of Fit Test, aimed to assess whether the observed distribution of penguin species aligned with an expected uniform distribution, where each species is equally represented. The test results, with an extremely low p-value, strongly rejected the notion of a uniform distribution, indicating that the species are not equally present in the dataset.

The second test, a Chi-Squared Test of Independence, investigated a potential relationship between penguin species and the island they inhabit. The null hypothesis for this test posited that species and island are independent variables, implying no association between the two. However, the test yielded a significantly low p-value, leading to the rejection of the null hypothesis and the conclusion that there is a statistically significant association between penguin species and their island of origin.

## 8. Multinomial Logistic Regression for Species Classification

The application of multinomial logistic regression to the Palmer Penguins dataset has provided robust insights into the species classification challenges. Leveraging four physical characteristics—culmen length and depth, flipper length, and weight—as predictors, the model demonstrates substantial efficacy in distinguishing between the three species: Adélie, Chinstrap, and Gentoo. The results indicate a perfect classification accuracy of 100% on the training dataset, with all penguins accurately categorized into their respective species. This outcome is further substantiated by the confusion matrix, which revealed zero misclassifications across species, and

the sensitivity, specificity, and positive predictive values, which all achieved a score of 1, underscoring the model's impeccable performance on the dataset used.

## Discussion

While the results from the training data are promising, they also raise critical considerations for the model's application and interpretation:

- **Overfitting Concerns:** The perfect accuracy might suggest that the model has overfitted the training data. This phenomenon occurs when a model learns the details and noise in the training data to an extent that it negatively impacts the performance of the model on new data. Implementing cross-validation techniques or partitioning the data into a separate training and testing set can help evaluate the model's generalizability.
- **Data Imbalance:** If the dataset is imbalanced, with uneven representation of species, it could bias the model towards the majority class. Exploring methods such as synthetic data generation through SMOTE, or adjusting class weights in the model, could help address this potential issue.
- **Feature Engineering:** Currently, the model utilizes straightforward measurements available within the dataset. Investigating more complex features or interactions between features might enhance the model's ability to generalize better to new, unseen data.
- **Comparative Analysis with Other Models:** Comparing the multinomial logistic regression model's performance with other classification techniques like decision trees, random forests, or support vector machines could provide deeper insights into the most effective modeling approach for this particular ecological data.

The multinomial logistic regression model's current application to the Palmer Penguins dataset successfully demonstrates its capability in accurately classifying penguin species based on



morphological data. However, the considerations of overfitting, potential data imbalances, and the need for a more robust evaluation via cross-validation highlight the importance of further methodological rigor. This approach not only validates the model's effectiveness but also ensures its adaptability and accuracy in broader ecological and conservation contexts. Future research directions could involve integrating additional ecological variables, applying model regularization techniques, and testing the model on independent datasets to solidify its predictive power and practical utility in biodiversity conservation efforts.

## Conclusion: Insights and Future Directions in Penguin Research

This comprehensive exploration of the Palmer Penguins dataset through a variety of statistical techniques, including multinomial logistic regression, has provided significant insights into the complex world of Antarctic penguin biodiversity. The research highlights the critical role of morphological data in understanding species differentiation and adaptation within the challenging environments of the Palmer Archipelago. Our use of multinomial logistic regression confirmed the distinct morphological profiles of Adélie, Chinstrap, and Gentoo penguins, achieving perfect classification accuracy in our training dataset. This accuracy, while indicative of the model's strength, also raises concerns about potential overfitting, emphasizing the need for rigorous validation and testing on independent datasets.

The findings from various statistical analyses not only augment our understanding of penguin species distribution and their ecological niches but also underline the importance of considering ecological dynamics and evolutionary pressures in conservation strategies. Our results indicate a significant correlation between penguin morphological characteristics and their habitat preferences, which are crucial for effective conservation management.

Looking forward, this study sets the stage for further research that should include deeper dives into the ecological implications of morphological diversity, such as dietary specialization and competitive interactions. Integrating longitudinal and broader geographic data could enhance the generalizability of our findings and help predict changes over time, particularly in response to environmental pressures like climate change.

Additionally, expanding our methodological approaches to include newer statistical models and comparison with other classification techniques will refine our predictive capabilities and improve our understanding of species resilience. By embracing a more holistic approach that incorporates both new data and novel analytical techniques, future research can provide more nuanced insights into the adaptive strategies of these penguins, thereby supporting more targeted and effective conservation efforts in a rapidly changing world.

## References

- Horst, M. A., Presmanes Hill, A., & B. Gorman, K. (2022). Palmer Archipelago Penguins Data in the palmerpenguins R Package - An Alternative to Anderson's Irises. *The R Journal*, 14(1), 244–254. <https://doi-org.jwupvdz.idm.oclc.org/10.32614/RJ-2022-020>
- Trathan, P. N., Forcada, J., Atkinson, R., Downie, R. H., & Shears, J. R. (n.d.). Population assessments of gentoo penguins (*Pygoscelis papua*) breeding at an important Antarctic tourist site, Goudier Island, Port Lockroy, Palmer Archipelago, Antarctica. *Biological Conservation*, 141(12), 3019–3028. <https://doi-org.jwupvdz.idm.oclc.org/10.1016/j.biocon.2008.09.006>