

Parallel Computation Framework for Discrete Copula Modeling

Dhyey Mavani



Amherst College

Submitted to the Department of Mathematics and Statistics
of Amherst College in partial fulfillment of the requirements
for the degree of Bachelor of Arts with honors.

Advisor(s):
Professor Shu-Min Liao

February 3, 2025

Table of contents

Abstract	iv
Acknowledgements	v
1 Introduction	1
2 Unraveling Notion of Dependence through Copulas	2
2.1 Tables	4
2.2 Chapter 2 Code	5
3 Checkerboard Copula and Regression Association Measure	15
3.1 Chapter 2 Code	15
4 Applications of Parallel Computing	18
4.1 Chapter 2 Code	18
5 Software (Package) Implementation and Testing	20
5.1 Chapter 2 Code	20
6 Conclusion	23
6.1 Chapter 2 Code	23
References	25
Appendices	27
A Code availability	27
B Corrections	30

Abstract

Understanding regression dependencies among discrete variables in categorical data—especially with ordinal responses—is a significant challenge in fields like finance, where the natural order of variables can unlock deeper insights into underlying distributions and data generating processes (DGPs). While numerous model-based methods have been developed to examine these structures, there is a notable lack of flexible, model-free approaches. To address this gap, a novel model-free measure based on the checkerboard copula, was introduced by Wei and Kim (2021) to identify and quantify regression dependence in multivariate categorical data involving both ordinal and nominal variables. Building upon this foundation, my thesis focuses on developing scalable and modularized implementations of discrete checkerboard copula modeling in R and Python, utilizing parallel computing to enhance efficiency and accessibility for large-scale data analysis. Initial experimentation and deployment confirm the effectiveness of these tools, providing researchers with a powerful resource for exploratory modeling and a deeper investigation into regression dependence structures within complex categorical datasets.

Acknowledgements

I want to thank everyone who made my experience unforgettable during my thesis journey. Firstly, I want to convey my gratitude to Professor Shu-Min Liao for advising me throughout my time at Amherst College and believing in me to take on the challenge of developing a statistical software component encompassing her most recent research work. From my first research experience on campus building R-Blocks to introducing me to her research collaborator (Professor Daeyoung Kim), Prof. Liao played a pivotal role in my development. Additionally, I am indebted to Dr. Kim for his continuous encouragement and feedback while developing the software this past year.

I am also incredibly grateful to my college advisor and statistics major advisor, Professor Nicholas Horton, for always advocating for me and supporting me throughout the Amherst College experience. I also thank Professor Jun Ishii for teaching me Advanced Econometrics and Professor Amy Wagaman for teaching me Advanced Data Analysis, which helped me gain a clear and solid understanding of the foundational tools I could build on in this work.

Finally, I would like to express gratitude towards my family for their constant belief in my abilities. Special thanks to my mom, dad, sister, grandfather, and grandmother for making me capable of the opportunity to study abroad. Last but not least, I would like to thank my

friends, peers, and colleagues on campus, who took courses, worked, and played sports with me. This academic, personal, and professional growth journey would not be possible without their support.

Chapter 1

Introduction

xxx

Chapter 2

Unraveling Notion of Dependence through Copulas

In this chapter, we will formalize some aspects of phenomena such as dependence and association. In order to aid us in our understanding, we will use two bivariate random vectors. Along the way, we will visualize and analyze various aspects of their dependence and association through Python code.

Consider (X_1, X_2) and (Y_1, Y_2) be bivariate random vectors, each consisting of 10000 independent data-points, which are distributed with the joint distributions F_X and F_Y respectively. Given these bivariate vectors, one might ask: How can I compare the relationship between (X_1, X_2) to the relationship between (Y_1, Y_2) ? One of the measures that can help us compare and contrast these relationships is Pearson correlation coefficient (commonly denoted as $\rho_{pearson}$). After preliminary calculations on a Python3 kernel, we can see that $\rho_{pearson}(X_1, X_2) \approx 0.802$, but on the other hand, the correlation between $\rho_{pearson}(Y_1, Y_2) \approx 0.755$. From these measure-values, it seems that the dependence between

(X_1, X_2) is stronger than the dependence between (Y_1, Y_2) . Although this agrees with our scatter plots in Figure 2.1, it is vital to note that $\rho_{pearson}$ only captures the linear dependence between the underlying random variables at hand.

Upon observing the Figure 2.1 closely, we note that the marginal distributions of X_1 and X_2 are close to normal, unlike the marginals of Y_1 and Y_2 . Moreover, we can see that the relationship between Y_1 and Y_2 is non-linear. This vast difference in marginals takes away our trust from the appropriateness of the use of $\rho_{pearson}$ as a measure to compare dependence between the data vectors at hand.

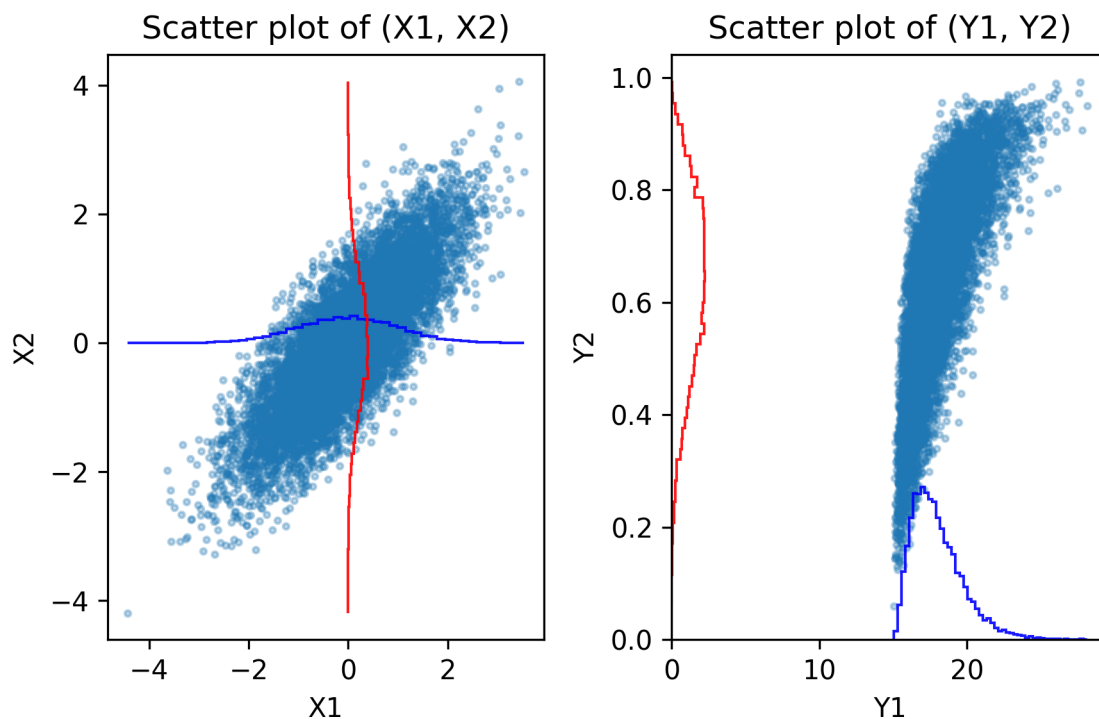


Figure 2.1: Scatter plots of 10000 independent observations of (X_1, X_2) and (Y_1, Y_2) with overlaid curves depicting respective marginal distributions.

Let's introduce a lemma that will help us transform the marginals so that the resulting marginals are more similar, and try to only capture or extract the “dependence” components, which will allow us to make fairer comparisons.

Lemma 2.1 (Probability Integral Transformation). *(Hofert et al. 2018) Let F be a continuous distribution function and let $X \sim F$, then $F(X)$ is a standard uniform random variable, that is, $F(X) \sim U(0, 1)$.*

Lemma 2.1 allows us to transform a continuous random variable to a random variable which has standard uniform distribution. So, by using this transformation, we can now convert our marginals X_1, X_2, Y_1, Y_2 individually to be distributed $\text{Uniform}(0, 1)$. And, since now the resulting marginals will all be of the same type, it will allow us to compare the dependence between random variables on fairer grounds.

For instance, if we know that $X_1 \sim N(0, 1) = F_1$, $X_2 \sim N(0, 1) = F_2$, $Y_1 \sim \text{Gamma}(3, 15) = G_1$, and $Y_2 \sim \text{Beta}(5, 3) = G_2$, where F_1, F_2, G_1, G_2 denote the distribution functions of the respective random variables. By Lemma 2.1, we can say that $F_1(X_1), F_2(X_2), G_1(Y_1)$, and $G_2(Y_2)$ are each distributed $\text{Uniform}(0, 1)$.

Looking at Figure 2.2, we can see that the transformed data vectors appear to be significantly similar. We can computationally verify this by quickly calculating the ρ_{pearson} for $(F_1(X_1), F_2(X_2))$ and $(G_1(Y_1), G_2(Y_2))$, which turns out to be 0.788 for both data vector pairs, implying that both have same dependence.

2.1 Tables

Your tables should be publication quality. Consider using [gt](#) (Iannone et al. 2024) or [kable-Extra](#) (Zhu 2024) to customize your tables. The [gtsummary](#) package (Sjoberg et al. 2021) may also come in handy.

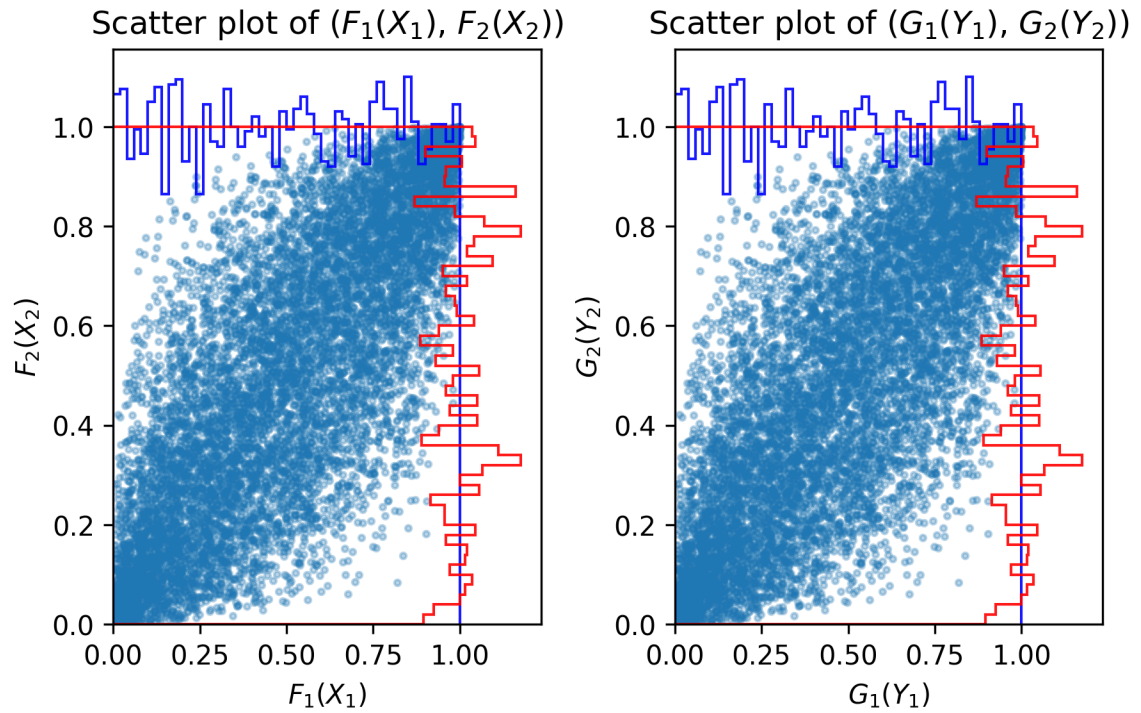


Figure 2.2: Scatter plots of 10000 independent observations of $(F_1(X_1), F_2(X_2))$ and $(G_1(Y_1), G_2(Y_2))$ with overlaid curves depicting respective marginal distributions.

2.2 Chapter 2 Code

The following code was used to create Chapter 2.

2.2.1 Code within chapter

```

1 # Load knitr package
2 library(knitr)
3
4 # Python Engine Setup
5 knit_engines$set(python3 = knit_engines$get("python"))
6
7 # Load packages

```

```

8  library(tidyverse)

9  library(gt)

10

11 # Set default ggplot theme for document

12 theme_set(theme_classic())

13 # If using kableExtra tables, print blank cells instead of `NA`

14 options(knitr.kable.NA = "")

15

16 # Load NBA Data

17 load("data/temp_wnba.RData")

18 import numpy as np

19 import matplotlib.pyplot as plt

20 from scipy.stats import beta, expon, norm, gamma, binom

21 import os

22

23 # Create directory if not exists

24 fig_dir = "fig"

25 os.makedirs(fig_dir, exist_ok=True)

26

27 # Generate Data

28 np.random.seed(8990)

29 n = 10000

30 mean = [0, 0]

```

```

31 cov = [[1, 0.8], [0.8, 1]]
32 X = np.random.multivariate_normal(mean, cov, size=n)
33 X1, X2 = X[:, 0], X[:, 1]
34
35 # Transform X1 and X2 to uniform [0, 1] using the CDF of the normal distribution
36 U1 = norm.cdf(X1)
37 U2 = norm.cdf(X2)
38
39 # Transform U1 and U2 into Gamma and Beta distributions
40 Y1 = gamma.ppf(U1, 3, 15)
41 Y2 = beta.ppf(U2, 5, 3)
42
43 # Calculate Pearson Correlation Coefficients
44 rho_X = np.corrcoef(X1, X2)[0, 1]
45 rho_Y = np.corrcoef(Y1, Y2)[0, 1]
46 print("Pearson correlation for (X1, X2):", rho_X)
47 print("Pearson correlation for (Y1, Y2):", rho_Y)
48
49 # Create Layout design and Set Size-Ratio
50 fig, axes = plt.subplots(1, 2, figsize=(6, 4))
51
52 # Scatter plot for (X1, X2)
53 axes[0].scatter(X1, X2, alpha=0.3, s=5)

```

```

54 axes[0].set_title("Scatter plot of (X1, X2)")
55 axes[0].set_xlabel("X1")
56 axes[0].set_ylabel("X2")
57
58 # Add marginal histograms
59 axes[0].hist(X1, bins=50, density=True, alpha=0.9, color='blue', orientation='vertical', histtype='step')
60 axes[0].hist(X2, bins=50, density=True, alpha=0.9, color='red', histtype='step', orientation='vertical')
61
62 # Scatter plot for (Y1, Y2)
63 axes[1].scatter(Y1, Y2, alpha=0.3, s=5)
64 axes[1].set_title("Scatter plot of (Y1, Y2)")
65 axes[1].set_xlabel("Y1")
66 axes[1].set_ylabel("Y2")
67
68 # Add marginal histograms
69 axes[1].hist(Y1, bins=50, density=True, alpha=0.9, color='blue', orientation='vertical', histtype='step')
70 axes[1].hist(Y2, bins=50, density=True, alpha=0.9, color='red', histtype='step', orientation='vertical')
71
72 # Organize into a tight layout as per matplotlib
73 plt.tight_layout()
74
75 # Save figure instead of showing it
76 fig_path = os.path.join(fig_dir, "motivating_example.png")

```

```

77 plt.savefig(fig_path, dpi=300, bbox_inches='tight')
78
79 # Close the figure to prevent rendering output
80 plt.close(fig)
81 knitr::include_graphics("fig/motivating_example.png")
82
83 # Set random seed for reproducibility
84 np.random.seed(8990)
85
86 # Number of samples
87 n = 10000
88
89 # Generate bivariate normal data (X1, X2)
90 mean = [0, 0]
91 cov = [[1, 0.8], [0.8, 1]] # Correlation ~ 0.8
92 X = np.random.multivariate_normal(mean, cov, size=n)
93 X1, X2 = X[:, 0], X[:, 1]
94
95 # Transform X1 and X2 to uniform [0, 1] using normal CDF
96 U_X1 = norm.cdf(X1)
97 U_X2 = norm.cdf(X2)
98
99 # Transform uniform marginals to desired distributions

```

```

100 Y1 = gamma.ppf(U_X1, a=3, scale=1/15)
101 Y2 = beta.ppf(U_X2, a=5, b=3)
102
103 # Apply probability integral transformation to all variables to make them uniform
104 U_Y1 = gamma.cdf(Y1, a=3, scale=1/15)
105 U_Y2 = beta.cdf(Y2, a=5, b=3)
106
107 # Calculate Pearson Correlation Coefficients
108 rho_U_X = np.corrcoef(U_X1, U_X2)[0, 1]
109 rho_U_Y = np.corrcoef(U_Y1, U_Y2)[0, 1]
110 print("Pearson correlation for ($F_1(X_1)$, $F_2(X_2)$):", rho_U_X)
111 print("Pearson correlation for ($G_1(Y_1)$, $G_2(Y_2)$):", rho_U_Y)
112
113 # Combine transformed data
114 uniform_data = np.vstack([U_X1, U_X2, U_Y1, U_Y2]).T
115
116 # Verify the uniformity of transformed data (Should be 0.5 in value)
117 print("U_X1 mean:", U_X1.mean(), "U_X2 mean:", U_X2.mean())
118 print("U_Y1 mean:", U_Y1.mean(), "U_Y2 mean:", U_Y2.mean())
119
120 # Create Layout design and Set Size-Ratio
121 fig, axes = plt.subplots(1, 2, figsize=(6, 4))
122

```



```

123 # Scatter plot for (U_X1, U_X2)

124 axes[0].scatter(U_X1, U_X2, alpha=0.3, s=5)

125 axes[0].set_title("Scatter plot of ( $F_1(X_1)$ ,  $F_2(X_2)$ )")

126 axes[0].set_xlabel(" $F_1(X_1)$ ")

127 axes[0].set_ylabel(" $F_2(X_2)$ ")

128

129 # Add marginal histograms

130 axes[0].hist(U_X1, bins=50, density=True, alpha=0.9, color='blue', orientation='vertical')

131 axes[0].hist(U_X2, bins=50, density=True, alpha=0.9, color='red', histtype='step', orientation='vertical')

132

133 # Scatter plot for (U_Y1, U_Y2)

134 axes[1].scatter(U_Y1, U_Y2, alpha=0.3, s=5)

135 axes[1].set_title("Scatter plot of ( $G_1(Y_1)$ ,  $G_2(Y_2)$ )")

136 axes[1].set_xlabel(" $G_1(Y_1)$ ")

137 axes[1].set_ylabel(" $G_2(Y_2)$ ")

138

139 # Add marginal histograms

140 axes[1].hist(U_Y1, bins=50, density=True, alpha=0.9, color='blue', orientation='vertical')

141 axes[1].hist(U_Y2, bins=50, density=True, alpha=0.9, color='red', histtype='step', orientation='vertical')

142

143 # Organize into a tight layout as per matplotlib

144 plt.tight_layout()

145

```

```

146 # Save figure instead of showing it
147 fig_path = os.path.join(fig_dir, "transformed_motivating_example.png")
148 plt.savefig(fig_path, dpi=300, bbox_inches='tight')
149
150 # Close the figure to prevent rendering output
151 plt.close(fig)
152 knitr::include_graphics("fig/transformed_motivating_example.png")
153
154 # =====
155 # Sample R script for thesis template
156 #
157 # Cleans temp_raw_wnba.csv dataset, which contains data pulled from
158 # https://www.espn.com/wnba/stats/player on 2024/06/19
159 #
160 # Last updated: 2024/06/19
161 # =====
162 library(tidyverse)
163
164 wnba <- read_csv("data/temp_raw_wnba.csv") |>
165   janitor::clean_names() |>
166   # Pull jersey numbers off of names and
167   # turn height text into msmt (6'4" = 6.3333)
168   mutate(jersey = str_extract(name, "[0-9]+$"),

```

```

169     name = str_remove(name, "[0-9]+$"),
170     ht_ft = parse_number(str_extract(ht, "^[0-9]")),
171     ht_in = parse_number(str_extract(ht, "[0-9]+\\\"$")),
172     height = ht_ft * 12 + ht_in,
173     weight = parse_number(wt),
174     position = factor(pos,
175                       levels = c("G", "F", "C"),
176                       labels = c("Guard", "Forward", "Center")) |>
177   select(-c(ht, wt, ht_ft, ht_in, pos))
178
179   save(wnba, file = "data/temp_wnba.RData")

```

2.2.2 Code sourced from external scripts

```

1  # =====
2  # Sample R script for thesis template
3  #
4  # Cleans temp_raw_wnba.csv dataset, which contains data pulled from
5  # https://www.espn.com/wnba/stats/player on 2024/06/19
6  #
7  # Last updated: 2024/06/19
8  # =====
9  library(tidyverse)
10
11 wnba <- read_csv("data/temp_raw_wnba.csv") |>

```

```

12  janitor::clean_names() |>
13  # Pull jersey numbers off of names and
14  # turn height text into msmt (6'4" = 6.3333)
15  mutate(jersey = str_extract(name, "[0-9]+$"),
16         name = str_remove(name, "[0-9]+$"),
17         ht_ft = parse_number(str_extract(ht, "^[0-9]")),
18         ht_in = parse_number(str_extract(ht, "[0-9]+\\\"$")),
19         height = ht_ft * 12 + ht_in,
20         weight = parse_number(wt),
21         position = factor(pos,
22                           levels = c("G", "F", "C"),
23                           labels = c("Guard", "Forward", "Center")))) |>
24  select(-c(ht, wt, ht_ft, ht_in, pos))
25
26  save(wnba, file = "data/temp_wnba.RData")

```

Chapter 3

Checkerboard Copula and Regression Association Measure

Start with motivation and connect with earlier chapters, then define concepts and weave in examples to solidify readers understanding.

3.1 Chapter 2 Code

The following code was used to create Chapter 2.

3.1.1 Code within chapter

```
1 # Load packages
2 library(tidyverse)
3 library(gt)
4
5 # Set default ggplot theme for document
```

```

6  theme_set(theme_classic())

7  # If using kableExtra tables, print blank cells instead of `NA`

8  options(knitr.kable.NA = "")

9

10 # Load data

11 load("data/temp_wnba.RData")

12 # =====

13 # Sample R script for thesis template

14 #

15 # Doesn't do anything useful

16 #

17 # Last updated: 2024/08/24

18 # =====

19

20 print("Hello, Amherst!")

```

3.1.2 Code sourced from external scripts

```

1  # =====

2  # Sample R script for thesis template

3  #

4  # Doesn't do anything useful

5  #

6  # Last updated: 2024/08/24

7  # =====

```

8

9 `print("Hello, Amherst!")`

Chapter 4

Applications of Parallel Computing

Start with motivation, connect, and dive into details + examples of parallel computing along with use-cases.

4.1 Chapter 2 Code

The following code was used to create Chapter 2.

4.1.1 Code within chapter

```
1 # Load packages
2 library(tidyverse)
3 library(gt)
4
5 # Set default ggplot theme for document
6 theme_set(theme_classic())
7 # If using kableExtra tables, print blank cells instead of `NA`
```



```

8  options(knitr.kable.NA = "")
9
10 # Load data
11 load("data/temp_wnba.RData")
12 # =====
13 # Sample R script for thesis template
14 #
15 # Doesn't do anything useful
16 #
17 # Last updated: 2024/08/24
18 # =====
19
20 print("Hello, Amherst!")

```

4.1.2 Code sourced from external scripts

```

1  # =====
2  # Sample R script for thesis template
3  #
4  # Doesn't do anything useful
5  #
6  # Last updated: 2024/08/24
7  # =====
8
9  print("Hello, Amherst!")

```

Chapter 5

Software (Package) Implementation and Testing

Start basic with motivation, and breakdown the implementation and testing phases properly while making sure that they are accessible.

5.1 Chapter 2 Code

The following code was used to create Chapter 2.

5.1.1 Code within chapter

```
1 # Load packages
2 library(tidyverse)
3 library(gt)
4
5 # Set default ggplot theme for document
```

```

6  theme_set(theme_classic())

7  # If using kableExtra tables, print blank cells instead of `NA`

8  options(knitr.kable.NA = "")

9

10 # Load data

11 load("data/temp_wnba.RData")

12 # =====

13 # Sample R script for thesis template

14 #

15 # Doesn't do anything useful

16 #

17 # Last updated: 2024/08/24

18 # =====

19

20 print("Hello, Amherst!")

```

5.1.2 Code sourced from external scripts

```

1  # =====

2  # Sample R script for thesis template

3  #

4  # Doesn't do anything useful

5  #

6  # Last updated: 2024/08/24

7  # =====

```

8

9 `print("Hello, Amherst!")`

Chapter 6

Conclusion

Like introduction, pull everything together and conclude the work!

6.1 Chapter 2 Code

The following code was used to create Chapter 2.

6.1.1 Code within chapter

```
1 # Load packages
2 library(tidyverse)
3 library(gt)
4
5 # Set default ggplot theme for document
6 theme_set(theme_classic())
7 # If using kableExtra tables, print blank cells instead of `NA`
8 options(knitr.kable.NA = "")
```

```

9
10 # Load data
11 load("data/temp_wnba.RData")
12 # =====
13 # Sample R script for thesis template
14 #
15 # Doesn't do anything useful
16 #
17 # Last updated: 2024/08/24
18 # =====
19
20 print("Hello, Amherst!")

```

6.1.2 Code sourced from external scripts

```

1 # =====
2 # Sample R script for thesis template
3 #
4 # Doesn't do anything useful
5 #
6 # Last updated: 2024/08/24
7 # =====
8
9 print("Hello, Amherst!")

```

References

- Hofert, M., Kojadinovic, I., Maechler, M., and Yan, J. (2018), *Elements of Copula Modeling with r*, Springer Use R! Series.
- Iannone, R., Cheng, J., Schloerke, B., Hughes, E., Lauer, A., Seo, J., Brevoort, K., and Roy, O. (2024), *Gt: Easily create presentation-ready display tables*.
- Sjoberg, D. D., Whiting, K., Curry, M., Lavery, J. A., and Larmarange, J. (2021), “Reproducible summary tables with the gtsummary package,” *The R Journal*, 13, 570–580. <https://doi.org/10.32614/RJ-2021-053>.
- Ushey, K., and Wickham, H. (2024), *Renv: Project environments*.
- Wei, Z., and Kim, D. (2021), “On exploratory analytic method for multi-way contingency tables with an ordinal response variable and categorical explanatory variables,” *Journal of Multivariate Analysis*, 186, 104793. <https://doi.org/10.1016/j.jmva.2021.104793>.
- Zhu, H. (2024), *kableExtra: Construct complex table with 'kable' and pipe syntax*.

Appendix A

Code availability

This thesis is written using Quarto with **renv** (Ushey and Wickham 2024) to create a reproducible environment. All materials (including the data sets and source files) required to reproduce this document can be found at the Github repository github.com/GITHUB-USERNAME/THESIS-REPO-NAME.

This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-nc-sa/4.0/).

```
1 # =====
2 # Sample R script for thesis template
3 #
4 # Cleans temp_raw_wnba.csv dataset, which contains data pulled from
5 # https://www.espn.com/wnba/stats/player on 2024/06/19
6 #
7 # Last updated: 2024/06/19
```

```

8 # =====
9 library(tidyverse)
10
11 wnba <- read_csv("data/temp_raw_wnba.csv") |>
12   janitor::clean_names() |>
13   # Pull jersey numbers off of names and
14   # turn height text into msmt (6'4" = 6.3333)
15   mutate(jersey = str_extract(name, "[0-9]+$"),
16          name = str_remove(name, "[0-9]+$"),
17          ht_ft = parse_number(str_extract(ht, "^[0-9]")),
18          ht_in = parse_number(str_extract(ht, "[0-9]+\\\"$")),
19          height = ht_ft * 12 + ht_in,
20          weight = parse_number(wt),
21          position = factor(pos,
22                             levels = c("G", "F", "C"),
23                             labels = c("Guard", "Forward", "Center"))) |>
24   select(-c(ht, wt, ht_ft, ht_in, pos))
25
26 save(wnba, file = "data/temp_wnba.RData")

```

```

1 # =====
2 # Sample R script for thesis template
3 #
4 # Doesn't do anything useful

```

```
5 #  
6 # Last updated: 2024/08/24  
7 # =====  
8  
9 print("Hello, Amherst!")
```

Appendix B

Corrections

This section may be excluded if no corrections are made to your thesis after initial submission to the department and before final submission to the college.

Per the [Statistics Honors Thesis Regulations](#):

Corrections to theses may be made after the date on which they are due in the Department's hands. Corrections may be made to the body of the thesis, but every such correction will be acknowledged in a list under the heading "Corrections," along with the statement "When originally submitted, this honors thesis contained some errors which have been corrected in the current version. Here is a list of the errors that were corrected." This list will be given on a sheet or sheets to be appended to the thesis. Corrections to spelling, grammar, or typography may be acknowledged by a general statement such as "30 spellings were corrected in various places in the thesis, and the notation for definite integral was changed in approximately 10 places." However, any correction that affects the meaning of a sentence or paragraph should be described in careful detail, and substantial

additions to the thesis will not be allowed. Questions about what should appear in the “Corrections” should be directed to the Chair. Electronic versions of the thesis, technical appendix, and necessary data and supplemental files must all be updated at the time of correction as well.

When originally submitted, this honors thesis contained some errors which have been corrected in the current version. Here is a list of the errors that were corrected.

1. ...
2. ...