

# **CCRVAM: A Model-Free Framework for Exploratory Analysis of Multi-Dimensional Discrete Data with Ordinal Response Variable**

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 (Amherst College), Professor Daeyoung Kim (University of  
Massachusetts Amherst)

February 22, 2025



# Table of contents

<b>Abstract</b>	<b>v</b>
<b>Acknowledgements</b>	<b>vi</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Motivation and Exploratory Data Analysis (EDA)</b>	<b>2</b>
<b>3 Copulas and Association Measures</b>	<b>3</b>
3.1 Unraveling the Notion of Dependence . . . . .	3
3.2 Copulas as a Unified Framework for Dependence . . . . .	8
3.3 Sklar’s Theorem and Invariance Principle . . . . .	10
3.4 Copulas for Continuous and Discrete Data . . . . .	12
3.5 Copula-based Measures of Association and Estimation . . . . .	14
<b>4 Checkerboard Copula Regression, its Visualization and Association measure for Model-Free Regression Dependence Analysis of Multivariate Discrete Data</b>	<b>21</b>
4.1 Set-up for Multivariate Categorical Data . . . . .	22
4.2 Checkerboard Copula and its Density . . . . .	22
4.3 Checkerboard Copula Score . . . . .	22
4.4 Checkerboard Copula Regression, Prediction and Visualization . . . . .	22
4.5 Checkerboard Copula Regression Association Measure . . . . .	23
<b>5 Software (Package) Implementation and Testing</b>	<b>24</b>
5.1 Set-up and Example Data . . . . .	24
5.2 Types of Input Data Supported . . . . .	24
5.3 Checkerboard copula score (especially an ordinal response variable) . . . . .	25
5.4 Checkerboard copula Regression (CCR) . . . . .	25
5.5 CCR Prediction and Visualization . . . . .	25
5.6 (Scaled) Checkerboard copula Regression Association Measure . . . . .	25
5.7 Visualization of Dependence Structures . . . . .	26
5.8 Software Architecture and Design Principles . . . . .	26
5.9 Testing, Validation, and Performance Evaluation . . . . .	26
5.10 User Documentation and Example Workflows . . . . .	26
<b>6 Real Data Analysis</b>	<b>27</b>
6.1 Dataset . . . . .	27

<b>7 Conclusion and Future Work</b>	<b>28</b>
<b>References</b>	<b>29</b>
<b>Appendices</b>	<b>31</b>
<b>A Code availability</b>	<b>31</b>
A.1 Chapter 2 Code . . . . .	31
A.2 Chapter 3 Code . . . . .	35
A.3 Chapter 4 Code . . . . .	39
A.4 Chapter 5 Code . . . . .	43
A.5 Chapter 6 Code . . . . .	47
<b>B Corrections</b>	<b>51</b>

# Abstract

Understanding regression dependencies among discrete variables—particularly when dealing with ordinal responses—remains a challenging yet vital task for uncovering the underlying structure in complex datasets. Traditional exploratory data analysis (EDA) methods and continuous copula models offer valuable insights for continuous data, but they often fall short when applied directly to categorical data, leading to issues with interpretability and generalizability. This thesis first revisits these traditional approaches, critically examining their limitations in the context of discrete data analysis. Building on this foundation, a novel model-free dependence measure based on the checkerboard copula is introduced to robustly identify and quantify regression relationships in multidimensional contingency tables containing both ordinal and nominal variables. The work further details the development of scalable, modularized implementations—primarily in Python, with complementary tools like NumPy, Pandas, SciPy, and Matplotlib to enhance efficiency for large-scale analyses. Through extensive experimentation and real-world case studies, the proposed framework and accompanying software package, Discopula, are shown to provide researchers with a powerful and flexible resource for exploratory modeling, paving the way for deeper insights into regression dependence structures in 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

# Motivation and Exploratory Data Analysis (EDA)

Start with motivation and connect with how papers and book mention it.

For this chapter 2, you need to read some references concerning two topics :

- 1) Importance of EDA in statistical analysis and data science
- 2) non-model-based (or model-free) association measures for categorical data.

To this end, you may want to read references in cited in Wei and Kim (2021) and a book chapter that Professor Liao sent you.

2.3 limitations (ref book and JMA2021), and mention how our thing CCRAM and SCCRAM are valid for n-dimensions (extensibility)

## Chapter 3

# Copulas and Association Measures

### 3.1 Unraveling the Notion of Dependence

In this section, we aim to formalize concepts of dependence and association. To facilitate our understanding, we will use two bivariate random vectors and visualize their relationships through Python code.

#### 3.1.1 Motivating Example

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 3.1, it is vital to note that  $\rho_{pearson}$  only captures the linear dependence between the underlying random variables at hand.

Upon observing the Figure 3.1 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.

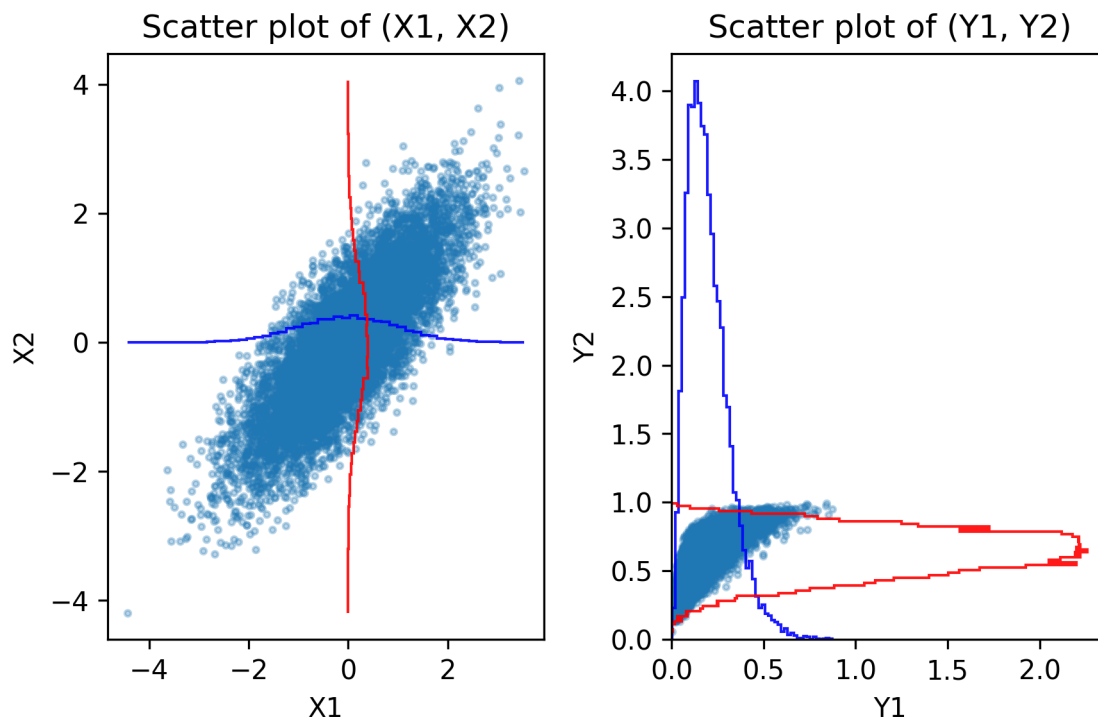


Figure 3.1: Scatter plots of 10000 independent observations of  $(X_1, X_2)$  and  $(Y_1, Y_2)$  with overlaid curves depicting respective marginal distributions.

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 3.1** (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 3.1 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 3.1, 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 3.2, we can see that the transformed data vectors appear to be significantly similar. We can computationally verify this by quickly calculating the  $\rho_{\text{pearson}}$  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, meaning that both have same dependence structures.

An alternative way to approach the problem (of comparing dependence of distinct pairs of marginals), is by transforming the marginals of  $(Y_1, Y_2)$  to be normal (same as marginals of  $(X_1, X_2)$ ). As one can predict, in order to accomplish this transformation, we would need to “undo the current distributional mappings on  $(Y_1, Y_2)$ ”, which we can formally define as generalized inverse as follows:

**Definition 3.1** (Quantile Function). *(Hofert et al. 2018)  $F^{\leftarrow}$  (Quantile Function) is defined*

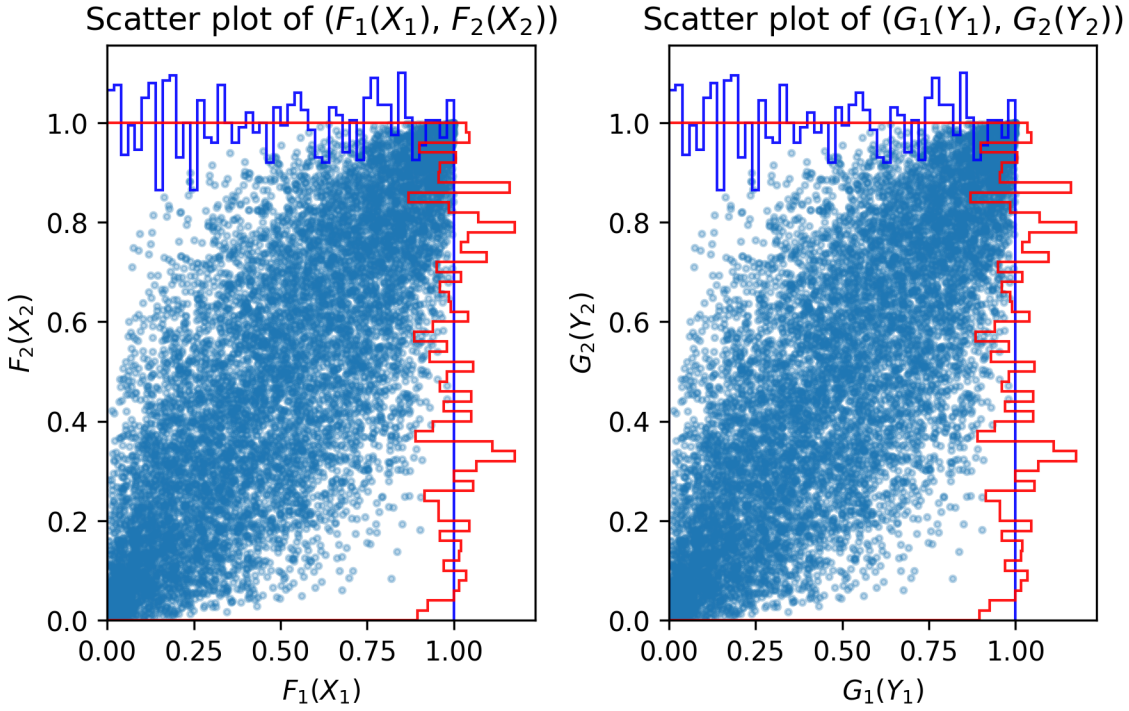


Figure 3.2: Scatter plots of 10000 independent observations of  $(F_1(X_1), F_2(X_2))$  and  $(G_1(Y_1), G_2(Y_2))$  with overlaid curves depicting respective marginal distributions.

as  $F^{\leftarrow}(y) = \inf\{x \in \mathbb{R} | F(x) \geq y\}$ , where  $y \in [0, 1]$ , and  $\inf$  is the infimum of a set.

**Warning**

The quantile function  $F^{\leftarrow} = F^{-1}$  only when  $F$  is continuous and strictly increasing. Thus it is important to note that, in other cases, the ordinary inverse  $F^{-1}$  need not exist. (Hofert et al. 2018)

With the above definition of  $F^{\leftarrow}$ , let's introduce a lemma from (Hofert et al. 2018) that will help us perform the transformation to normal.

**Lemma 3.2** (Quantile Transformation). *(Hofert et al. 2018) Let  $U \sim \text{Unif}(0, 1)$  and let  $F$  be any distribution function be a distribution function. Then  $F^{\leftarrow}(U) \sim F$ , that is,  $F^{\leftarrow}(X)$  is distributed with density  $F$ .*

**Note**

Lemma 3.2 is valid for non-continuous densities  $F$  as well. (Hofert et al. 2018)

Let's start with the transformations where we left off in Figure 3.2, since we have uniform densities there. Applying Lemma 3.2 on  $G_1(Y_1)$  and  $G_2(Y_2)$  using quantile functions  $F_1^{\leftarrow} = F_1^{-1}$  and  $F_2^{\leftarrow} = F_2^{-1}$  respectively gives us that  $F_1^{-1}(G_1(Y_1)) \sim F_1$  and  $F_2^{-1}(G_2(Y_2)) \sim F_2$ .

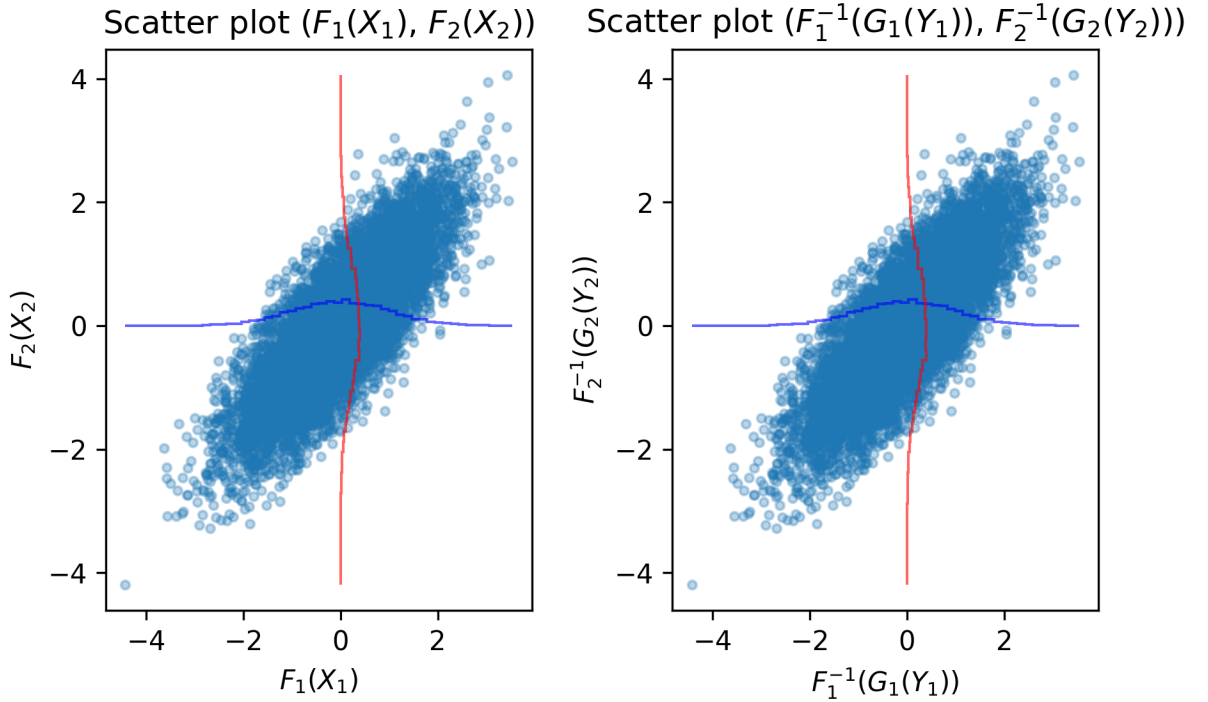


Figure 3.3: Scatter plots of 10000 independent observations of  $(X_1, X_2)$  and  $(F_1^{-1}(G_1(Y_1)), F_2^{-1}(G_2(Y_2)))$  with overlaid curves depicting respective marginal distributions.

Notice in Figure 3.3 that the resulting transformed distribution through this alternative method resembles that of  $(X_1, X_2)$ . Hence, we can conclude that they have the same dependence. Furthermore, through a quick calculation, we can see that  $\rho_{\text{pearson}}(F_1^{-1}(G_1(Y_1)), F_2^{-1}(G_2(Y_2))) = 0.802$ , which is the same as the Pearson correlation coefficient between  $X_1$  and  $X_2$ . This is the level of flexibility that a combination of transformations presented in Lemma 3.1 and Lemma 3.2 can lend us.

**i** Note

“(X<sub>1</sub>, X<sub>2</sub>) and (Y<sub>1</sub>, Y<sub>2</sub>) have the same dependence”  $\iff$  “(X<sub>1</sub>, X<sub>2</sub>) and (Y<sub>1</sub>, Y<sub>2</sub>) have the same copula” (Hofert et al. 2018)

## 3.2 Copulas as a Unified Framework for Dependence

Copulas are a class of multivariate distribution functions with  $Unif(0, 1)$  marginals. The motivating example in the previous section explains the usage of copulas as the structures capturing margin-independent dependence between random variables.

**i** Note

The choice of  $Unif(0, 1)$  as a post-transformation margin for the data at hand is somewhat arbitrary although it does simplify further results. One can use modifications of Lemma 3.1 and Lemma 3.2 to define copulas with respect to any margin of choice without affecting the final conclusions about the dependence between the data at hand. (Hofert et al. 2018)

In order to understand copulas better, for now, let’s restrict ourselves to the 2-D (2-dimensional) case. Firstly, let’s introduce the definition of a broader class of functions called subcopulas as a preliminary, which will help us mathematically define copulas as a special case. (Nelsen 2006)

**Definition 3.2** (2-Dimensional Subcopula). (Erdely 2017) A **two-dimensional subcopula** (2-subcopula) is a function  $C^S : D_1 \times D_2 \rightarrow [0, 1]$ , where  $\{0, 1\} \subseteq D_i \subseteq [0, 1]$  for  $i \in \{1, 2\}$  with the following conditions satisfied:

- *Grounded*:  $C^S(u, 0) = 0 = C^S(0, v)$ ,  $\forall u \in D_1, \forall v \in D_2$ .



- *Marginal Consistency*:  $\forall u \in D_1$  and  $\forall v \in D_2$ ,  $C^S(u, 1) = u$  and  $C^S(1, v) = v$ .
- *2-increasing*:  $\forall u_1, u_2 \in D_1$  and  $\forall v_1, v_2 \in D_2$  such that  $u_1 \leq u_2$  and  $v_1 \leq v_2$ ,  $C^S(u_1, v_1) - C^S(u_2, v_1) + C^S(u_2, v_2) - C^S(u_1, v_2) \geq 0$ .

**Definition 3.3** (2-Dimensional Copula). (Erdely 2017) A **two-dimensional copula** (2-copula) is a function  $C : [0, 1] \times [0, 1] \rightarrow [0, 1]$ , with the following conditions satisfied:

- *Grounded*:  $C(u, 0) = 0 = C(0, v)$ ,  $\forall u \in [0, 1]$ ,  $\forall v \in [0, 1]$ .
- *Marginal Consistency*:  $\forall u \in [0, 1]$  and  $\forall v \in [0, 1]$ ,  $C(u, 1) = u$  and  $C(1, v) = v$ .
- *2-increasing*:  $\forall u_1, u_2 \in [0, 1]$  and  $\forall v_1, v_2 \in [0, 1]$  such that  $u_1 \leq u_2$  and  $v_1 \leq v_2$ ,  $C(u_1, v_1) - C(u_2, v_1) + C(u_2, v_2) - C(u_1, v_2) \geq 0$ .

**i Note**

A 2-D copula is essentially a 2-subcopula with a full unit square as domain ( $D_1 = D_2 = [0, 1]$ ). Furthermore, copula and subcopula are the same within a domain with continuous variables. Later in this chapter, we will discuss why this doesn't hold when one of the variables is discrete.

In this work, we will mainly deal with 2-D copulas and subcopulas, but the definitions above can be generalized to n-D case with some notable exceptions detailed (with proofs) in section 2.10 of Nelsen (2006). Moreover, there are many different families of copulas bearing peculiar properties and corresponding margins, we are not covering them in detail since that is not the focus of this work, and a comprehensive summary of many of these families can be found in chapter 3 of Hofert et al. (2018).

### 3.2.1 Fréchet-Hoeffding Bounds

For any distribution function, boundedness is always a desired property. In the case of copulas, we have a famous theorem that provides us the upper and lower pointwise bounds.

**Theorem 3.1** (Fréchet-Hoeffding Bounds). *(Hofert et al. 2018)* Given a 2-D copula  $C$ ,  $W(u, v) = \max\{0, u + v - 1\} \leq C(u, v) \leq \min\{u, v\} = M(u, v)$ , where  $u, v \in [0, 1]$ .

## 3.3 Sklar's Theorem and Invariance Principle

Theorem 3.2 by (Sklar 1959) is one of the seminal results in copula theory, which extended the applications of copulas, and explained why copulas captures the dependence by relating the joint distributions to univariate margins.

**Theorem 3.2** (Fréchet-Hoeffding Bounds). *(Hofert et al. 2018)*

1. Let  $H$  be a joint distribution function with univariate margins  $F$  and  $G$ . Then there exists a copula  $C$  such that  $\forall x, y \in \mathbb{R}, H(x, y) = C(F(x), G(y))$ . Furthermore,  $C$  is **unique** in the case when  $F, G$  are continuous; otherwise, in the general case,  $C$  is uniquely determined on  $\text{Ran}F \times \text{Ran}G$ , where  $\text{Ran}F, \text{Ran}G$  denote the ranges of  $F, G$  respectively. That copula  $C$  is given by:  $C(u, v) = H(F^{\leftarrow}(u), G^{\leftarrow}(v))$  such that  $(u, v) \in \text{Ran}F \times \text{Ran}G$ .
2. Conversely,  $H$  is defined as a 2-D distribution function with marginals  $F, G$ , if we are given copula  $C$  along with the univariate marginals  $F, G$ .

In this work, we will mainly deal with two dimensions, but Theorem 3.2 above can be generalized to n-D case as detailed in section 2.10 of Nelsen (2006). Below, we include a few insights drawn from (Hofert et al. 2018) that will be important to our ongoing discussion:

**i** Note

Theorem 3.2 gives us an insight into the name copula as in how it “couples” a joint distribution function to its marginal distributions. This coupling effect and two parts of Theorem 3.2 show us how we can separate (or combine) multivariate dependence structure and univariate margins.

**!** Spoiler Alert

In the case of continuous random variables, there is only one **unique** copula that characterizes the multivariate dependence structure, which is very convenient for reasons we will discuss later in this chapter. This is not the case with discrete variables, which make the direct use of continuous copulas intractable.

**i** Note

Theorem 3.2 can be used to verify the existence of a continuous distribution function  $H$  in case of a multivariate dataset if and only if we are sure of the existence of corresponding continuous univariate marginals for each variable in the dataset.

### 3.3.1 The Invariance Principle

As we saw in the motivating example, the underlying dependence structure did not change over a certain type of transformations. This was very convenient for us, and thus is a favorable property for a copula to have. This property is often formally referred to as “invariance”, which we will formalize in the following theorem from (Hofert et al. 2018)

**Theorem 3.3** (Invariance Principle). *Let  $(X, Y) \sim H$  with continuous margins  $F, G$  and copula  $C$ . If  $T_X, T_Y$  are **strictly increasing** transformations on  $\text{Ran}X, \text{Ran}Y$ , respectively, then  $(T_X(X), T_Y(Y))$  also has copula  $C$ .*

#### Note

Theorem 3.3 was implicitly in action during our analysis for the motivating example because the transformations that we used were of two kinds, namely, probability integral transformation and quantile transformation, and in both of the cases, we were dealing with continuous and **strictly increasing** mappings on the respective ranges of random variables.

### 3.4 Copulas for Continuous and Discrete Data

- you can discuss Challenges and Pitfalls When Applying Continuous Copulas to Discrete Data
- in discrete case, copula is not as flexible as in continuous case

Up to this point, our discussion has centered on continuous random variables. Many of the results and definitions we have used rely on continuity, which ensures that the probability integral transform (PIT) maps each variable to a uniform distribution on  $[0, 1]$ . This property, in turn, guarantees the uniqueness of the copula associated with a joint distribution via Sklar's theorem. In our earlier work, we have taken this uniqueness for granted.

However, real-world data are often **discrete**. When dealing with discrete random variables, the marginal distribution functions are not continuous, and the PIT no longer produces uniform random variables on the full interval  $[0, 1]$ . Instead, we obtain what is known as a **subcopula**—a function defined only on a proper subset of  $[0, 1]^2$ , namely on the ranges of the marginal distributions.

#### **i** Example: Bivariate Bernoulli Distribution

*Imagine a bivariate distribution where each variable follows a Bernoulli law. In this setting, the only possible values for each variable are 0 and 1. The resulting subcopula is then defined on the set of points. Because this set is a proper subset of  $[0, 1]^2$ , the corresponding copula is not uniquely determined by the joint distribution of the variables.*

### 3.4.1 Unidentifiability Issue

Now, let us examine the unidentifiability problem in more detail. To illustrate the issue, consider the following adapted example in the two-dimensional case, inspired by Geenens (2020). Suppose we have a subcopula  $C^S$  defined on a discrete domain, where  $D_1 = \text{Ran}(F)$  and  $D_2 = \text{Ran}(G)$  with the marginal distribution functions  $F$  and  $G$ , respectively. In the continuous case, a two-dimensional (sub)copula is defined on the entire unit square  $[0, 1]^2$ . By contrast, for discrete random variables, the subcopula  $C^S$  is only uniquely specified on the domain  $D_1 \times D_2 = \text{Ran}(F) \times \text{Ran}(G)$ .

To obtain a full copula  $C$  on  $[0, 1]^2$ , one must “fill in” the gaps—that is, extend the definition of  $C^S$  to those parts of the unit square not covered by  $D_1 \times D_2$ . Unfortunately, there are uncountably many ways to perform this extension while still satisfying the fundamental properties required of a copula in its Definition 3.3. This leads to a **non-uniqueness** (or **unidentifiability**) issue, which complicates both the development and the application of copula-based models for discrete data. This unidentifiability has been examined in depth in the literature such as Geenens (2020), and it calls into question the straightforward (direct) application of copula methods when at least one margin is discrete.

One of the ways to fill in the gaps is by performing a Distributional Transform, which

basically serves to add random “noise” to each of the gaps in parent distribution as described by Rüschendorf (2009) and Faugeras (2017). Formally, considering a random variable  $X \sim F$  and independently, consider  $V \sim U(0, 1)$ , then the distributional transform of  $X$  is  $F(X, V) = P(X < x) + V * P(X = x)$ . After applying this, we can directly proceed to apply results from continuous copula modeling as we have smoothened out the discontinuities. Another method that also accomplishes this goal is described in the next chapter.

- here, you can introduce Checkerboard Copula

### 3.5 Copula-based Measures of Association and Estimation

Now that we have built an object (copula) that allows us to just capture the multivariate dependence structure between variables, we would like to encode certain pieces of this information into a set of robust measures or metrics. We would call these measures, the **measures of association**. There are two types of measures of association: parametric and non-parametric. As discussed briefly for our motivating example, a common (parametric) measure of association is the Pearson correlation coefficient ( $\rho_{pearson}$ ). Although it is really efficient to calculate, it only captures linear dependence between the random data vectors at hand. Let’s discuss this metric in more detail along with its limitations:

#### 3.5.1 Pearson’s Correlation Coefficient ( $\rho_{pearson}$ ) & its Properties

**Definition 3.4** (Pearson correlation coefficient). Given a random vector  $(X, Y)$  with  $Var(X) < \infty$  and  $Var(Y) < \infty$ , then:

$$\rho_{pearson}(X, Y) = \frac{Cov(X, Y)}{\sqrt{Var(X)}\sqrt{Var(Y)}}$$

, where covariance is defined as:

$$\text{Cov}(X, Y) = \mathbb{E}((X - \mathbb{E}(X))(Y - \mathbb{E}(Y)))$$

, and the variance is defined as  $\text{Var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2)$ .

Let's start by going over some commonly-used properties of  $\rho_{\text{pearson}}$  as mentioned in Hofert et al. (2018):

1.  $\rho_{\text{pearson}} \in [-1, 1]$
2.  $|\rho_{\text{pearson}}(X, Y)| = 1$  if and only if  $\exists a, b \in \mathbb{R}$ , with  $a \neq 0$  such that  $Y = aX + b$  almost surely  
with  $a < 0$  if and only if  $\rho_{\text{pearson}}(X, Y) = -1$ , and  $a > 0$  if and only if  $\rho_{\text{pearson}}(X, Y) = 1$ .  
In both cases,  $X, Y$  are called *perfectly linearly dependent*
3. If  $X$  and  $Y$  are independent, then  $\rho_{\text{pearson}}(X, Y) = 0$ .
4.  $\rho_{\text{pearson}}$  is invariant under *strictly increasing linear* transformations.

### 3.5.2 Limitations of Pearson's Correlation Coefficient ( $\rho_{\text{pearson}}$ )

Although Pearson's correlation coefficient  $\rho_{\text{pearson}}$  is useful in many cases, it only captures **linear dependence** and ignores non-linear relationships. Below, we summarize its key limitations along with illustrative examples.

1. **Non-Existence of  $\rho_{\text{pearson}}$ :** Pearson's correlation does not exist for every random vector  $(X, Y)$ , particularly when variances (or other higher order moments) are undefined.

### **i** Example: Heavy-Tailed Distributions

Consider two independent random variables  $X_1, X_2$  drawn from a **Pareto(3)** distribution with  $F(x) = 1 - x^{-3}$ ,  $x \geq 1$ . Define  $X = X_1$ , and  $Y = X_1^2$ . The covariance is given by  $Cov(X, Y) = Cov(X_1, X_1^2) = \mathbb{E}(X_1^3) - \mathbb{E}(X_1)\mathbb{E}(X_1^2)$ . For Pareto(3), it is well-known (and can be easily proven) that  $\mathbb{E}(X_1^3)$  **does not exist** (as the integral diverges). Since Pearson's formula rely on this moment,  $\rho_{pearson}(X, Y)$  **doesn't exist**. On the other hand, we can observe that  $Y = X^2$  shows a **perfect functional dependence**, since  $Y$  can be represented as a deterministic (quadratic) function of  $X$ .

2. **Non-Invariance Under Non-Linear Transformations:**  $\rho_{pearson}$  is not necessarily invariant under all strictly increasing transformations on  $\text{Ran}X$  or  $\text{Ran}Y$ .

### **i** Example: Logarithmic Transformation on $U(0, 1)$

Let  $X \sim U(0, 1)$  and define  $Y = \log(X)$ . Pearson's correlation is:  $\rho_{pearson}(X, Y) = \frac{Cov(X, \log X)}{\sigma_X \sigma_Y}$ . Even though  $Y = \log(X)$  is a **strictly increasing function**,  $\rho_{pearson}$  changes under this transformation. Thus, Pearson's correlation is **not invariant** under (non-linear) monotonic transformations such as log in certain situations.

3. **Uncorrelatedness Does Not Imply Independence:**  $\rho_{pearson} = 0$  does NOT necessarily imply that  $(X, Y)$  are independent.

### **i** Example: Quadratic Transformation on $U(-1, 1)$

Let  $X \sim U(-1, 1)$  and define:  $Y = X^2$ . We can compute:  $\mathbb{E}[X] = 0$ ,  $\mathbb{E}[Y] = \mathbb{E}[X^2] = \frac{1}{3}$ . Now, consider the covariance:  $Cov(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] = \mathbb{E}[X^3] - (0)(\frac{1}{3})$ . Since  $\mathbb{E}[X^3] = 0$ , we get  $Cov(X, Y) = 0$ . Thus,  $\rho_{pearson}(X, Y) = 0$ , but  $X$  **and**  $Y$  **are clearly dependent**, since knowing  $X$  exactly determines  $Y$ . This example demonstrates that a



zero Pearson correlation does **not** imply statistical independence.

4. **Non-Uniqueness of the Joint Distribution Given Marginals and  $\rho_{pearson}$ :** The marginal distributions and the correlation coefficient do not uniquely determine the joint distribution.

**i Example: Bivariate Normal and Mixture Distributions**

Consider two bivariate distributions:

1. **Bivariate Normal Distribution:**

$$(X_1, X_2) \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right).$$

2. **Bivariate Mixture Distribution** (Same Marginals, Different Dependence):

$$X_1 \sim N(0, 1), \quad X_2 = \begin{cases} X_1, & \text{with probability 0.75,} \\ -X_1, & \text{with probability 0.25.} \end{cases}$$

Both cases yield:  $\rho_{pearson}(X_1, X_2) = 0.5$ .

However, their **joint distributions are completely different**, meaning  $\rho_{pearson}$  **does not uniquely determine dependence**.

5. **Unattainability of Certain Correlations:** Given margins  $F_1, F_2$ , some  $\rho_{pearson} \in [-1, 1]$  values cannot be attained by choosing any possible copula for  $(X_1, X_2)$ . An example demonstrating this can be found in Hofert et al. (2018) p.46

In order to circumvent some of the limitations of pearson coefficient, we now consider rank-based correlation measures such as Spearman's Rho ( $\rho_{spearman}$ ) and Kendall's Tau ( $\tau_{kendall}$ )

as they only depend on the underlying copula  $C$  at least in the case of continuous random variables. Again, we will discuss the peculiarities of the discrete case later in this chapter.

These rank-based measures are also known as **measures of concordance**. (Hofert et al. 2018)

In order to better understand this, we would first need to define *concordance*. Consider two points in  $\mathbb{R}^2$ ,  $(x_1, y_1)$  and  $(x_2, y_2)$ . These points are defined as concordant if  $(x_1 - x_2)(y_1 - y_2) > 0$  and discordant if  $(x_1 - x_2)(y_1 - y_2) < 0$ .

### 3.5.3 Kendall's Tau

**Definition 3.5** (Kendall's Tau). Given a bivariate random vector  $(X_1, X_2)$  with continuous marginals  $F_1$  and  $F_2$ , let's define  $(X'_1, X'_2)$  as an independent copy of  $(X_1, X_2)$ . Then the population version of Kendall's tau is defined by:

$$\tau_{kendall}(X_1, X_2) = \mathbb{E}(\text{sign}((X_1 - X'_1)(X_2 - X'_2)))$$

Here,  $\text{sign}(x)$  is the sign-function defined in a piecewise manner as follows:

$$\text{sign}(x) = \begin{cases} -1, & \text{if } x < 0, \\ 0, & \text{if } x = 0, \\ 1, & \text{if } x > 0. \end{cases}$$

Using the above-mentioned notion of concordance, definition of an expected value, and Definition 3.5, we can equivalently define Kendall's Tau as  $\tau_{kendall} = (1)\mathbb{P}((X_1 - X'_1)(X_2 - X'_2) > 0) + (0)\mathbb{P}((X_1 - X'_1)(X_2 - X'_2) = 0) + (-1)\mathbb{P}((X_1 - X'_1)(X_2 - X'_2) < 0) = \mathbb{P}((X_1 - X'_1)(X_2 - X'_2) > 0) - \mathbb{P}((X_1 - X'_1)(X_2 - X'_2) < 0)$ , since in the case of continuous distributions, probability at

any given point is 0, specifically  $\mathbb{P}((X_1 - X'_1)(X_2 - X'_2) = 0) = 0$ .

As mentioned in Hofert et al. (2018) p.53, we can represent  $\tau_{kendall}$  in terms of an underlying copula  $C$  as  $\tau_{kendall}(C) = 4 \int_{[0,1]^2} C(u, v) d(C(u, v)) - 1$ .

- write on how issues with this copula representation in discrete case

### 3.5.4 Spearman's Rho

**Definition 3.6** (Spearman's Rho). Given a bivariate random vector  $(X_1, X_2)$  with continuous marginals  $F_1$  and  $F_2$ , then the population version of Spearman's rho is defined by:

$$\rho_{spearman}(X_1, X_2) = \rho_{pearson}(F_1(X_1), F_2(X_2))$$

We can observe that the Spearman's rho is nothing but Pearson's correlation coefficient of the transformed variables obtained after performing the Probability Integral Transformation defined earlier in Lemma 3.1.

As mentioned in Hofert et al. (2018) p.53, we can represent  $\rho_{spearman}$  in terms of an underlying copula  $C$  as  $\rho_{spearman}(C) = 12 \int_{[0,1]^2} C(u, v) d((u, v)) - 3$ .

#### **i** Note:

$\tau_{kendall}$  and  $\rho_{spearman}$  both overcome the significant limitations of  $\rho_{pearson}$  with the following properties as summarized in Hofert et al. (2018):

- These measures always exist, and are invariance under all (not just linear) strictly increasing transformations
- These measures attain all values in  $[-1, 1]$ , and they specifically attain -1 and 1 when the copula  $C$  attains the Fréchet-Hoeffding bounds  $W$  and  $M$  as defined in

### Theorem 3.1

- review commonly used measures of association and discuss their relation with copulas
- continuous case measures are already model-free, margin free and all, but this is not the case in discrete situations (use this for transition)

### 3.5.5 Checkerboard Copula-based Association Measures for 2-Dimensional Data

- you can find in the references [5,11,12,38] in my 2021 paper.

## **Chapter 4**

# **Checkerboard Copula Regression, its Visualization and Association measure for Model-Free Regression Dependence Analysis of Multivariate Discrete Data**

This Chapter is a review of Wei and Kim (2021) and a book chapter that Professor Liao sent you. Along the way, you can reference chapter 5 for 1-on-1 user correspondence to the package.

## 4.1 Set-up for Multivariate Categorical Data

Multivariate Categorical Data of interest for us is Multi-dimensional Contingency Table with an Ordinal Response Variable and a set of categorical (nominal/ordinal) Predictors

## 4.2 Checkerboard Copula and its Density

XXX

## 4.3 Checkerboard Copula Score

XXX

## 4.4 Checkerboard Copula Regression, Prediction and Visualization

- You can use CCR as an acronym for Checkerboard Copula Regression
- The CCR and its prediction is designed to explore and identify the potential regression association between an ordinal response variable and a set of categorical predictors of interest
- Point Prediction of the category of the ordinal dependent variable
- Uncertainty Evaluation of the CCR prediction using nonparametric bootstrap

## 4.5 Checkerboard Copula Regression Association Measure

- You can use SCCRAM as an acronym for (Scaled) Checkerboard copula Regression Association Measure
- (S)CCRAM is designed to quantify the regression association identified by Checkerboard Copula Regression and its prediction.
- Uncertainty Evaluation of the estimated (S)CCRAM using nonparametric bootstrap distribution and its confidence interval
- Statistical significance of the estimated (S)CCRAM using Permutation distribution and its hypothesis testing

## Chapter 5

# Software (Package) Implementation and Testing

Motivation + mention how we address the gap in presence of scalable well-tested tool

### 5.1 Set-up and Example Data

Multi-dimensional contingency table with an ordinal response variable and a set of categorical (nominal/ordinal) explanatory variables/predictors

### 5.2 Types of Input Data Supported

XXX

dat, txt, csv - for cases form

np.array directly with adjacent dicts - for contingency table



### **5.3 Checkerboard copula score (especially an ordinal response variable)**

XXX

### **5.4 Checkerboard copula Regression (CCR)**

XXX

### **5.5 CCR Prediction and Visualization**

- Point Prediction of the category of the ordinal dependent variable
- Uncertainty Evaluation of the CCR prediction using nonparametric bootstrap

### **5.6 (Scaled) Checkerboard copula Regression Association Measure**

- (S)CCRAM : quantify the regression association identified by Checkerboard Copula Regression and its prediction.
- Uncertainty Evaluation of the estimated (S)CCRAMs using nonparametric bootstrap distribution and its confidence interval
- Statistical significance of the estimated (S)CCRAMs using Permutation distribution and its hypothesis testing

## **5.7 Visualization of Dependence Structures**

XXX

## **5.8 Software Architecture and Design Principles**

XXX

## **5.9 Testing, Validation, and Performance Evaluation**

XXX

## **5.10 User Documentation and Example Workflows**

XXX

## Chapter 6

# Real Data Analysis

XXX

### 6.1 Dataset

Describe data at hand

- give reminder of different ways of initializing cc

## **Chapter 7**

# **Conclusion and Future Work**

Like introduction, pull everything together and conclude the work!

# References

- Erdely, A. (2017), “[A subcopula based dependence measure](#),” *Kybernetika*, Institute of Information Theory; Automation AS CR, 53, 231–243.
- Faugeras, O. P. (2017), *Dependence Modeling*, 5, 121–132. <https://doi.org/doi:10.1515/demo-2017-0008>.
- Geenens, G. (2020), “Copula modeling for discrete random vectors,” *Dependence Modeling*, 8, 417–440. <https://doi.org/doi:10.1515/demo-2020-0022>.
- Hofert, M., Kojadinovic, I., Maechler, M., and Yan, J. (2018), *Elements of Copula Modeling with r*, Springer Use R! Series.
- Nelsen, R. B. (2006), *An introduction to copulas*, Springer Science & business media.
- Rüschendorf, L. (2009), “On the distributional transform, sklar’s theorem, and the empirical copula process,” *Journal of Statistical Planning and Inference*, 139, 3921–3927. <https://doi.org/https://doi.org/10.1016/j.jspi.2009.05.030>.
- Sklar, M. (1959), “Fonctions de repartition an dimensions et leurs marges,” *Publ. inst. statist. univ. Paris*, 8, 229–231.
- Ushey, K., and Wickham, H. (2024), *Renv: Project environments*.



# 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](https://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/).

### A.1 Chapter 2 Code

The following code was used to create Chapter 2.

### A.1.1 Code within chapter

```
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!")
10
11 # =====
12 # Sample R script for thesis template
13 #
14 # Cleans temp_raw_wnba.csv dataset, which contains data pulled from
15 # https://www.espn.com/wnba/stats/player on 2024/06/19
16 #
17 # Last updated: 2024/06/19
18 # =====
19 library(tidyverse)
20
21 wnba <- read_csv("data/temp_raw_wnba.csv") |>
22   janitor::clean_names() |>
23   # Pull jersey numbers off of names and
```



```

24   # turn height text into msmt (6'4" = 6.3333)
25   mutate(jersey = str_extract(name, "[0-9]+$"),
26          name = str_remove(name, "[0-9]+$"),
27          ht_ft = parse_number(str_extract(ht, "^[0-9]")),
28          ht_in = parse_number(str_extract(ht, "[0-9]+\\\"$")),
29          height = ht_ft * 12 + ht_in,
30          weight = parse_number(wt),
31          position = factor(pos,
32                             levels = c("G", "F", "C"),
33                             labels = c("Guard", "Forward", "Center")))) |>
34   select(-c(ht, wt, ht_ft, ht_in, pos))
35
36   save(wnba, file = "data/temp_wnba.RData")
37
38
39
40   # =====
41   # Sample R script for thesis template
42   #
43   # Doesn't do anything useful
44   #
45   # Last updated: 2024/08/24
46   # =====

```

```
47
48 print("Hello, Amherst!")
49
50 # =====
51 # Sample R script for thesis template
52 #
53 # Doesn't do anything useful
54 #
55 # Last updated: 2024/08/24
56 # =====
57
58 print("Hello, Amherst!")
59
60 # =====
61 # Sample R script for thesis template
62 #
63 # Doesn't do anything useful
64 #
65 # Last updated: 2024/08/24
66 # =====
67
68 print("Hello, Amherst!")
```

### A.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!")
```

## A.2 Chapter 3 Code

The following code was used to create Chapter 3.

### A.2.1 Code within chapter

```
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!")

10

11  # =====

12  # Sample R script for thesis template

13  #

14  # Cleans temp_raw_wnba.csv dataset, which contains data pulled from

15  # https://www.espn.com/wnba/stats/player on 2024/06/19

16  #

17  # Last updated: 2024/06/19

18  # =====

19  library(tidyverse)

20

21  wnba <- read_csv("data/temp_raw_wnba.csv") |>

22    janitor::clean_names() |>

23    # Pull jersey numbers off of names and

24    # turn height text into msmt (6'4" = 6.3333)

25    mutate(jersey = str_extract(name, "[0-9]+$"),

26           name = str_remove(name, "[0-9]+$"),

27           ht_ft = parse_number(str_extract(ht, "^[0-9]")),

28           ht_in = parse_number(str_extract(ht, "[0-9]+\\\"$')),

29           height = ht_ft * 12 + ht_in,

30           weight = parse_number(wt),

31           position = factor(pos,

```

```

32         levels = c("G", "F", "C"),
33         labels = c("Guard", "Forward", "Center")))) |>
34     select(-c(ht, wt, ht_ft, ht_in, pos))
35
36 save(wnba, file = "data/temp_wnba.RData")
37
38
39
40 # =====
41 # Sample R script for thesis template
42 #
43 # Doesn't do anything useful
44 #
45 # Last updated: 2024/08/24
46 # =====
47
48 print("Hello, Amherst!")
49
50 # =====
51 # Sample R script for thesis template
52 #
53 # Doesn't do anything useful
54 #

```

```

55 # Last updated: 2024/08/24
56 # =====
57
58 print("Hello, Amherst!")
59
60 # =====
61 # Sample R script for thesis template
62 #
63 # Doesn't do anything useful
64 #
65 # Last updated: 2024/08/24
66 # =====
67
68 print("Hello, Amherst!")

```

### A.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")

```

## A.3 Chapter 4 Code

The following code was used to create Chapter 4.

### A.3.1 Code within chapter

```
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!")
10
11 # =====
12 # Sample R script for thesis template
13 #
14 # Cleans temp_raw_wnba.csv dataset, which contains data pulled from
15 # https://www.espn.com/wnba/stats/player on 2024/06/19
16 #
17 # Last updated: 2024/06/19
18 # =====
19 library(tidyverse)
20
21 wnba <- read_csv("data/temp_raw_wnba.csv") |>
22   janitor::clean_names() |>
23   # Pull jersey numbers off of names and
```



```

24   # turn height text into msmt (6'4" = 6.3333)
25   mutate(jersey = str_extract(name, "[0-9]+$"),
26          name = str_remove(name, "[0-9]+$"),
27          ht_ft = parse_number(str_extract(ht, "^[0-9]")),
28          ht_in = parse_number(str_extract(ht, "[0-9]+\\\"$")),
29          height = ht_ft * 12 + ht_in,
30          weight = parse_number(wt),
31          position = factor(pos,
32                             levels = c("G", "F", "C"),
33                             labels = c("Guard", "Forward", "Center")))) |>
34   select(-c(ht, wt, ht_ft, ht_in, pos))
35
36   save(wnba, file = "data/temp_wnba.RData")
37
38
39
40   # =====
41   # Sample R script for thesis template
42   #
43   # Doesn't do anything useful
44   #
45   # Last updated: 2024/08/24
46   # =====

```

```

47
48 print("Hello, Amherst!")
49
50 # =====
51 # Sample R script for thesis template
52 #
53 # Doesn't do anything useful
54 #
55 # Last updated: 2024/08/24
56 # =====
57
58 print("Hello, Amherst!")
59
60 # =====
61 # Sample R script for thesis template
62 #
63 # Doesn't do anything useful
64 #
65 # Last updated: 2024/08/24
66 # =====
67
68 print("Hello, Amherst!")

```

### A.3.2 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!")
```

## A.4 Chapter 5 Code

The following code was used to create Chapter 5.

### A.4.1 Code within chapter

```
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!")

10

11  # =====

12  # Sample R script for thesis template

13  #

14  # Cleans temp_raw_wnba.csv dataset, which contains data pulled from

15  # https://www.espn.com/wnba/stats/player on 2024/06/19

16  #

17  # Last updated: 2024/06/19

18  # =====

19  library(tidyverse)

20

21  wnba <- read_csv("data/temp_raw_wnba.csv") |>

22    janitor::clean_names() |>

23    # Pull jersey numbers off of names and

24    # turn height text into msmt (6'4" = 6.3333)

25    mutate(jersey = str_extract(name, "[0-9]+$"),

26           name = str_remove(name, "[0-9]+$"),

27           ht_ft = parse_number(str_extract(ht, "^[0-9]")),

28           ht_in = parse_number(str_extract(ht, "[0-9]+\\\"$")),

29           height = ht_ft * 12 + ht_in,

30           weight = parse_number(wt),

31           position = factor(pos,

```

```

32         levels = c("G", "F", "C"),
33         labels = c("Guard", "Forward", "Center")))) |>
34   select(-c(ht, wt, ht_ft, ht_in, pos))
35
36   save(wnba, file = "data/temp_wnba.RData")
37
38
39
40   # =====
41   # Sample R script for thesis template
42   #
43   # Doesn't do anything useful
44   #
45   # Last updated: 2024/08/24
46   # =====
47
48   print("Hello, Amherst!")
49
50   # =====
51   # Sample R script for thesis template
52   #
53   # Doesn't do anything useful
54   #

```

```

55 # Last updated: 2024/08/24

56 # =====

57

58 print("Hello, Amherst!")

59

60 # =====

61 # Sample R script for thesis template

62 #

63 # Doesn't do anything useful

64 #

65 # Last updated: 2024/08/24

66 # =====

67

68 print("Hello, Amherst!")

```

#### A.4.2 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!")
```

## A.5 Chapter 6 Code

The following code was used to create Chapter 6

### A.5.1 Code within chapter

```
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!")
10
11 # =====
12 # Sample R script for thesis template
13 #
14 # Cleans temp_raw_wnba.csv dataset, which contains data pulled from
15 # https://www.espn.com/wnba/stats/player on 2024/06/19
16 #
17 # Last updated: 2024/06/19
```

```

18 # =====
19 library(tidyverse)
20
21 wnba <- read_csv("data/temp_raw_wnba.csv") |>
22   janitor::clean_names() |>
23   # Pull jersey numbers off of names and
24   # turn height text into msmt (6'4" = 6.3333)
25   mutate(jersey = str_extract(name, "[0-9]+$"),
26          name = str_remove(name, "[0-9]+$"),
27          ht_ft = parse_number(str_extract(ht, "^[0-9]")),
28          ht_in = parse_number(str_extract(ht, '[0-9]+\\\\"$')),
29          height = ht_ft * 12 + ht_in,
30          weight = parse_number(wt),
31          position = factor(pos,
32                             levels = c("G", "F", "C"),
33                             labels = c("Guard", "Forward", "Center"))) |>
34   select(-c(ht, wt, ht_ft, ht_in, pos))
35
36 save(wnba, file = "data/temp_wnba.RData")
37
38
39
40 # =====

```



```

41 # Sample R script for thesis template
42 #
43 # Doesn't do anything useful
44 #
45 # Last updated: 2024/08/24
46 # =====
47
48 print("Hello, Amherst!")
49
50 # =====
51 # Sample R script for thesis template
52 #
53 # Doesn't do anything useful
54 #
55 # Last updated: 2024/08/24
56 # =====
57
58 print("Hello, Amherst!")
59
60 # =====
61 # Sample R script for thesis template
62 #
63 # Doesn't do anything useful

```

```
64 #  
65 # Last updated: 2024/08/24  
66 # =====  
67  
68 print("Hello, Amherst!")
```

### A.5.2 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!")
```

## 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. ...