

# **CCRVAM : a Python Package for Model-Free Exploratory Analysis of Multivariate 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

April 9, 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 Copulas and Association Measures</b>	<b>7</b>
2.1 Unraveling the Notion of Dependence . . . . .	7
2.2 Copulas as a Unified Framework for Dependence . . . . .	12
2.3 Sklar’s Theorem and Invariance Principle . . . . .	14
2.4 Copulas for Continuous and Discrete Data . . . . .	16
2.5 Checkerboard Copula . . . . .	18
2.6 Copula-based Measures of Association and Estimation . . . . .	21
<b>3 Checkerboard Copula Regression, its Visualization and Association measure for Model-Free Regression Dependence Analysis of Multivariate Discrete Data</b>	<b>32</b>
3.1 Data . . . . .	33
3.2 Checkerboard Copula and its Density . . . . .	34
3.3 Checkerboard Copula Score . . . . .	34
3.4 Checkerboard Copula Regression, Prediction and Visualization . . . . .	37
3.5 Checkerboard Copula Regression Association Measure . . . . .	43
3.6 Visualization Methods . . . . .	48
<b>4 Software (Package) Implementation and Testing</b>	<b>50</b>
4.1 Set-up and Example Data . . . . .	50
4.2 Types of Input Data Supported . . . . .	51
4.3 Checkerboard copula score (especially an ordinal response variable) . . . . .	54
4.4 Checkerboard copula Regression (CCR) . . . . .	56
4.5 CCR Predicted Category Visualization . . . . .	58
4.6 CCR Prediction Uncertainty Evaluation Using Nonparametric Bootstrap Resampling . . . . .	61
4.7 (S)CCRAM Estimation . . . . .	63
4.8 (S)CCRAM Uncertainty Evaluation Using Bootstrap Resampling . . . . .	64
4.9 Statistical Significance Testing for (S)CCRAM Using Permutation Test . . . . .	68
4.10 Software Architecture and Design Principles . . . . .	71
4.11 Testing, Validation, and Performance Evaluation . . . . .	74
4.12 User Documentation and Example Workflows . . . . .	76

<b>5</b>	<b>Real Data Analysis</b>	<b>78</b>
5.1	Dataset Overview . . . . .	78
5.2	Data Preparation and Loading . . . . .	79
5.3	Exploratory Data Analysis . . . . .	87
5.4	Calculating Checkerboard Copula Scores (CCS) . . . . .	87
5.5	Checkerboard Copula Regression (CCR) Analysis . . . . .	88
5.6	Quantifying Association with (S)CCRAM . . . . .	89
5.7	Uncertainty Quantification Using Bootstrap . . . . .	91
5.8	Statistical Significant Testing Using Permutation Tests . . . . .	92
5.9	Bootstrap Analysis for CCR Predictions . . . . .	94
5.10	Discussion and Clinical Interpretation . . . . .	96
<b>6</b>	<b>Conclusion and Future Work</b>	<b>100</b>
	<b>References</b>	<b>104</b>
	<b>Appendices</b>	<b>109</b>
<b>A</b>	<b>Code availability</b>	<b>109</b>
A.1	Chapter 2 Code . . . . .	109
A.2	Chapter 3 Code . . . . .	118
<b>B</b>	<b>Corrections</b>	<b>123</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, we explore Wei and Kim (2021) and Liao et al. (2024)’s model-free dependence measure based on the checkerboard copula to robustly identify and quantify regression relationships in multidimensional contingency tables containing both ordinal and nominal variables. The work then introduces the development of novel, scalable, modularized implementations—primarily in Python, using complementary libraries 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, *ccrvam*, 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 Prof. Kim for his continuous encouragement and feedback throughout this past year.

I am also incredibly grateful to my college and statistics major advisor, Professor Nicholas Horton, for always advocating for me and supporting me throughout the Amherst College experience. I also thank Prof. Jun Ishii, Prof. Amy Wagaman, Prof. Katharine Correia, and Prof. Pamela Matheson for teaching me Advanced Econometrics, Advanced Data Analysis, Missing Data Analysis, and Intermediate Statistics, which helped me gain a clear and solid understanding of the foundational tools.

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, and peers 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

Exploratory Data Analysis (EDA), a concept first articulated by John Tukey in 1977 (Tukey 1977), forms an indispensable foundation for statistical analyses and modern data science workflows. Tukey underscored the significance of letting data guide hypothesis generation, rather than relying exclusively on pre-defined models or confirmatory statistical tests. According to Tukey, exploring data should precede formal modeling efforts, as it reveals patterns, structures, outliers, and underlying assumptions that would otherwise remain hidden or overlooked (Tukey (1977); Buja et al. (2009); Gelman and Vehtari (2021)). The classic example of Anscombe's quartet (Anscombe 1973) vividly illustrates how datasets with identical summary statistics—means, variances, and correlation coefficients—can exhibit dramatically different relationships once explored visually, reinforcing the necessity of visual exploration in the preliminary stages of data analysis.

In practical terms, EDA employs a variety of visualization techniques and statistical summaries designed to enhance intuitive understanding and detect anomalies, thereby ensuring the validity and robustness of subsequent analyses. (Tufte (1983); Donoho (2017)) Modern ana-

lytical frameworks, such as the Cross-Industry Standard Process for Data Mining (CRISP-DM) (Shearer 2000), explicitly incorporate EDA to inform data preprocessing, feature selection, and model construction. Particularly in high-dimensional or complex datasets, effective EDA is crucial for identifying relevant variables, understanding multivariate interactions, and preventing misleading conclusions.

As the dimensionality and complexity of datasets increase, traditional EDA methods face significant challenges, particularly in effectively capturing complex dependence structures within multivariate categorical data, which is typically presented as multi-dimensional contingency table. Not only, there have been efforts to extend continuous data visualization methods to the discrete case, but also many graphical methods were developed specifically for categorical data. (Liao et al. 2024)

While visualization remains a powerful tool, statistical methodologies capable of systematically modeling and quantifying dependence become increasingly necessary as complements. Copula theory provides a promising approach for addressing these complexities. A copula is a mathematical function that combines marginal distributions of individual variables into a comprehensive joint distribution, enabling distinct and flexible modeling of dependence structures independently from marginal behaviors (Sklar (1959); Nelsen (2006); Joe (2014)). Copulas have found widespread application across diverse fields, notably finance, engineering, and environmental studies, particularly due to their ability to capture intricate and non-linear dependence patterns beyond traditional correlation measures.

Despite the advantages copulas offer in modeling continuous data, their application to discrete or categorical variables presents unique challenges. For discrete data, the copula representation is generally not unique, complicating standard analytical procedures. Addressing these

issues, the checkerboard copula has been proposed as a pragmatic and effective method for modeling discrete dependencies. This technique assigns uniform probabilities within discrete segments or ‘blocks,’ preserving the dependency structure intrinsic to categorical data while minimizing additional assumptions (Wei and Kim 2021). The checkerboard copula thus facilitates capturing complex dependencies that are otherwise challenging to quantify using conventional copula approaches.

One key aspect of statistical modeling that has consistently intrigued researchers across various disciplines—including social sciences, healthcare, economics, and behavioral research—is regression dependence. (Wei and Kim 2021) Specifically, regression dependence aims to quantify and interpret relationships where one variable, designated as the response, depends on a set of explanatory variables. Ordinal response variables present unique analytical challenges, necessitating methods specifically tailored to their inherent ordered structure.

Traditionally, assessing regression dependence in ordinal categorical data has relied primarily on parametric, model-based methods. Popular approaches include cumulative logit (proportional odds) models, cumulative link models with alternative links (probit, log-log, complementary log-log), adjacent-categories logit models, continuation ratio logit models, stereotype models, latent variable and association models, canonical correlation analyses, and correspondence analysis models (Wei and Kim 2021). While these approaches offer robust frameworks for explanatory or predictive modeling, they impose specific structural assumptions about relationships within the data, potentially limiting their exploratory power.

Alternatively, non-model-based approaches have been developed to provide flexible assessments of association without predefined functional forms. Commonly utilized methods include generalized Cochran-Mantel-Haenszel procedures, ordinal odds ratios, Kendall’s tau

and its variants, Goodman and Kruskal's gamma, Spearman's rank correlation, Somers' D, and Kendall's partial tau (Wei and Kim 2021). Despite their widespread use, these traditional methods suffer from significant limitations. They primarily address pairwise associations, often treat variables symmetrically without distinguishing clearly between explanatory and response roles, and inadequately leverage the ordinal information inherent within categorical data. Specifically, it is not straightforward to use non-model-based methods when our objective is to understand regression dependence in multivariate ordinal data with an ordinal response variable. (Wei and Kim 2021)

To address these limitations, recent methodological innovations by Wei and Kim (2021) have introduced the Checkerboard Copula Regression Association Measure (CCRAM) and its scaled counterpart, Scaled CCRAM (SCCRAM). The foundational tool of this method is the checkerboard (multilinear extension) copula (Genest et al. 2014; Schweizer and Sklar 1974), constructed through multilinear interpolation of discrete marginal distributions. Checkerboard copulas uniquely capture and represent the dependence structure among discrete ordinal variables, offering a robust basis for exploring multivariate contingency tables (Denuit and Lambert 2005; Nešlehová 2007).

CCRAM is a non-parametric, model-free measure explicitly designed to quantify regression-type dependencies involving an ordinal response and multiple categorical explanatory variables. By leveraging checkerboard copula scores—which incorporate the informative ordinal rankings of variables—CCRAM provides a comprehensive, robust measure analogous to the R-squared statistic commonly used in linear regression but tailored specifically for categorical contexts. It flexibly accommodates high-dimensional data, providing interpretable summaries of dependencies within complex contingency tables. (Wei and Kim 2021)

The introduction of SCCRAM further enhances interpretability by scaling CCRAM values onto a standardized 0–1 interval. This scaling facilitates straightforward, intuitive comparisons of association strengths across diverse research contexts and datasets. Additionally, CCRAM and SCCRAM enable detailed decomposition of multivariate dependencies, allowing researchers to isolate and analyze contributions from individual variables or subsets of predictors—an advantage notably lacking in traditional non-model-based measures (Wei et al. 2023).

This thesis is based upon these recent methodological innovations highlighted in Wei and Kim (2021) and Liao et al. (2024), aims to better understand these novel methods, and for the first time ever, introduce a Python package encapsulating the end-to-end EDA workflow that allow researchers to run their experiments correctly, and efficiently without worrying about implementational details of the methods all the time. In this work, we limit our focus to development of EDA workflow for multivariate-categorical data with an ordinal response variable and a set of independent categorical (ordinal or nominal) variables in a multi-way contingency table.

The remainder of the thesis is structured as follows. Chapter 2 provides an extensive review of copula theory and key association measures, beginning with an exploration of the fundamental concepts of dependence and copula functions, emphasizing continuous data contexts. It includes a discussion of Sklar’s theorem and the invariance principle, elaborates on copula modeling applicable to both continuous and discrete data, and covers essential copula-based association measures such as Kendall’s tau and Spearman’s rho, laying the theoretical foundation necessary for subsequent developments. Chapter 3 introduces Checkerboard Copula Regression (CCR), addressing the unique challenges posed by multivariate categorical data. It presents the theoretical formulation of checkerboard copulas, defines the checkerboard copula score, and illustrates the CCR methodology, its predictive capabilities, visualization

techniques, and the novel Checkerboard Copula Regression Association Measure (CCRAM). Chapter 4 focuses on the practical implementation and testing of CCR within CCRVAM, the novel software package that we developed, detailing the data requirements, software architecture, computational strategies, visualization tools, and rigorous validation processes. Chapter 5 provides a comprehensive empirical application of the proposed methods to real-world data, demonstrating the effectiveness of CCR in uncovering hidden dependence structures and improving analytical insights. Finally, Chapter 6 summarizes the key contributions, outlines the significance of integrating exploratory data analysis, copula theory, and novel association metrics, acknowledges current limitations, and suggests promising directions for future research, including methodological extensions and advanced visualization approaches.

## Chapter 2

# Copulas and Association Measures

### 2.1 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.

#### 2.1.1 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 python 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.

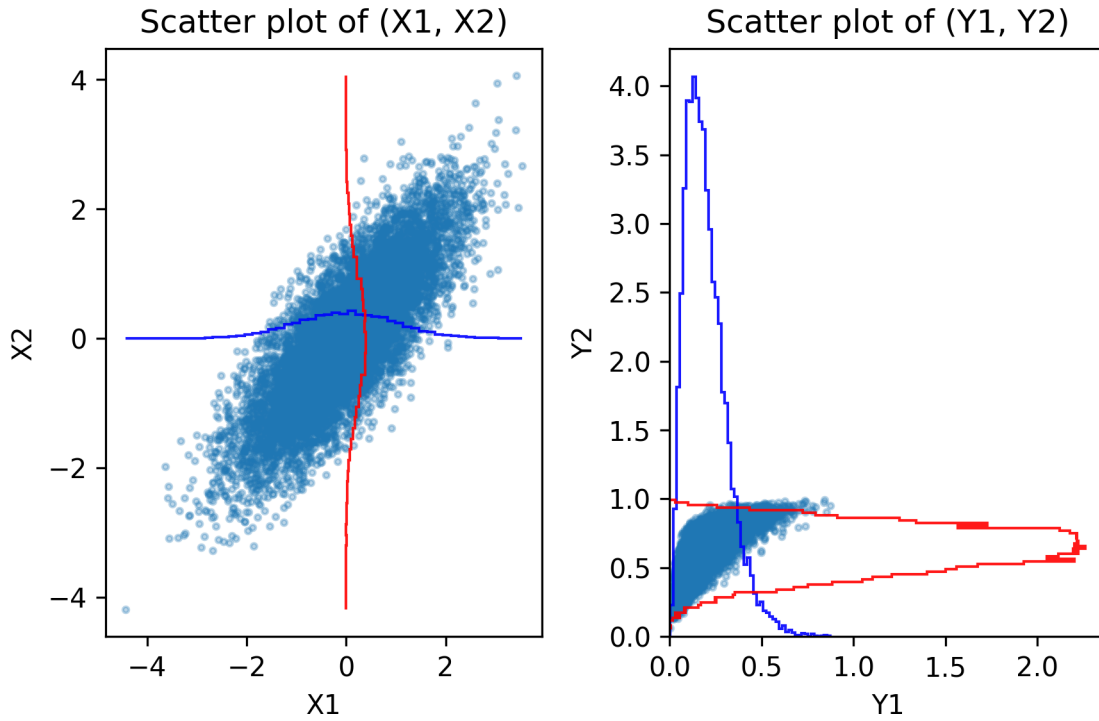


Figure 2.1: Scatter plots of 10000 independent observations of  $(X_1, X_2)$  and  $(Y_1, Y_2)$  with overlaid curves depicting respective marginal distributions.

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  $\rho_{\text{pearson}}$  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 2.1** (Quantile Function). *(Hofert et al. 2018)  $F^{\leftarrow}$  (Quantile Function) is defined*

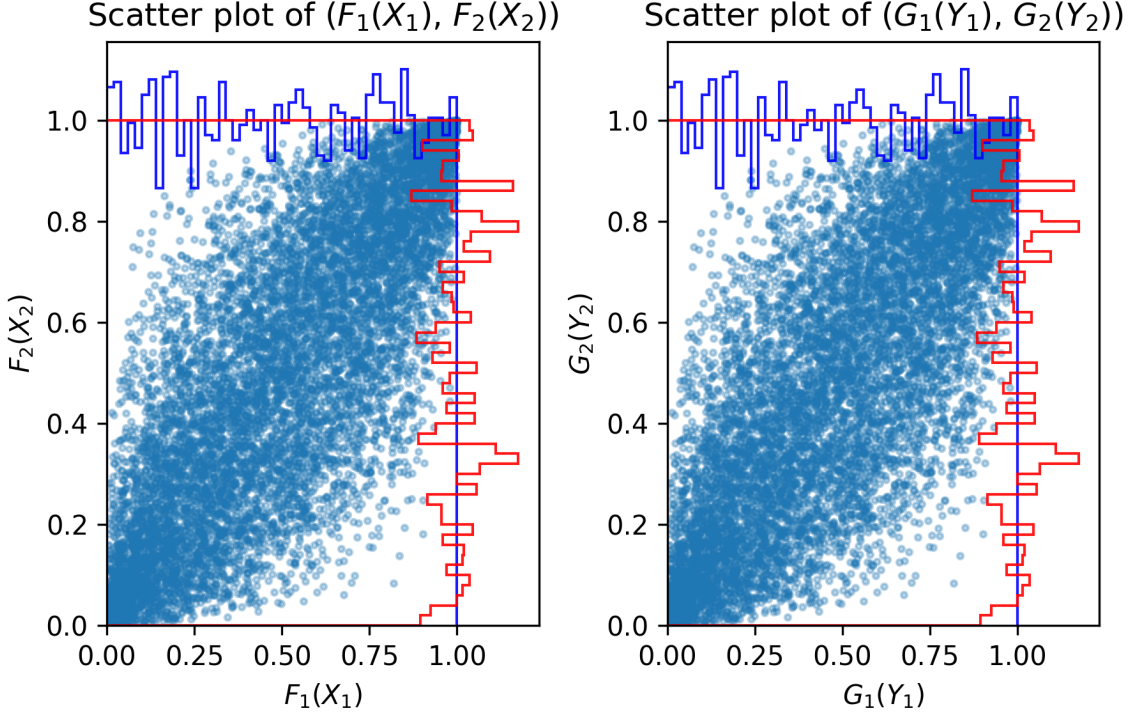


Figure 2.2: Scatter plots of 10000 independent observations of  $(F_1(X_1), F_2(X_2))$  and  $(G_1(Y_1), G_2(Y_2))$  with overlaid curves depicting respective marginal distributions.

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 2.2** (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 2.2 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 2.2, since we have uniform densities there. Applying Lemma 2.2 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  $F_2^{-1}(G_2(Y_2)) \sim F_2$ .

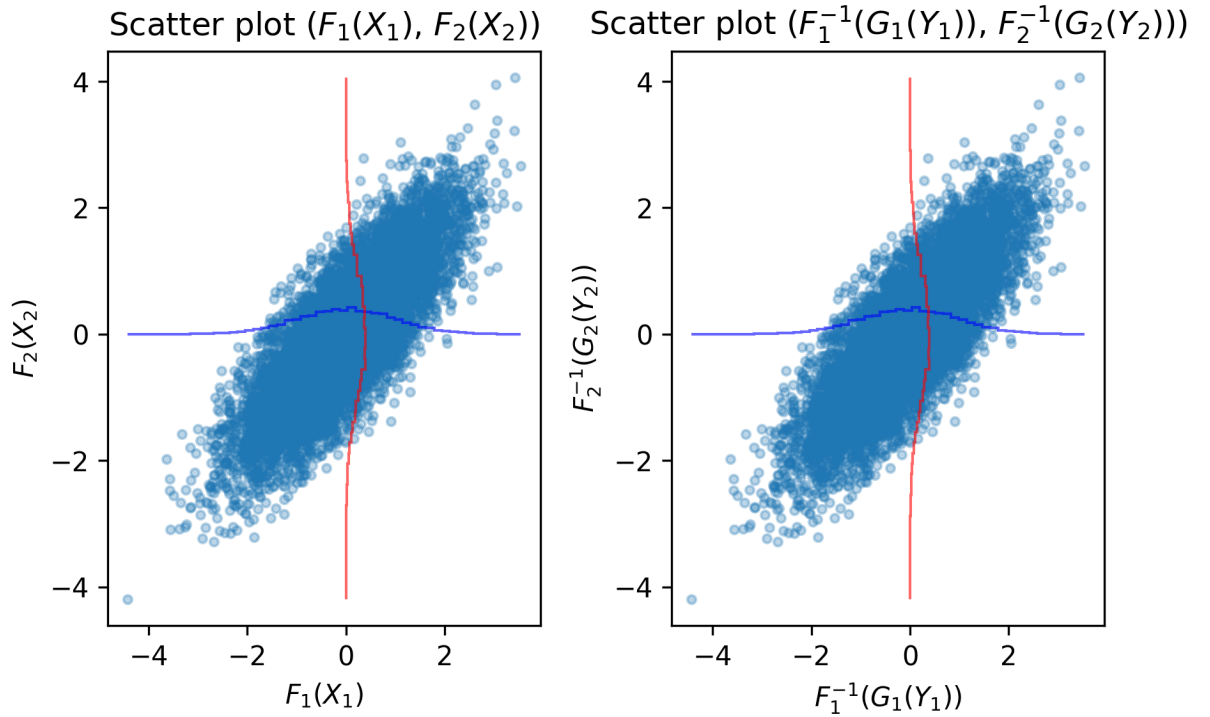


Figure 2.3: Scatter plots of 10000 independent observations of  $(X_1, X_2)$  and  $(F_1^{-1}(G_1(Y_1)), F_2^{-1}(G_2(Y_2)))$  with overlaid curves depicting respective marginal distributions.

Notice in Figure 2.3 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 2.1 and Lemma 2.2 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)

## 2.2 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 2.1 and Lemma 2.2 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 2.2** (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 2.3** (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).

### 2.2.1 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 2.1** (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  $u, v \in [0, 1]$ .

## 2.3 Sklar's Theorem and Invariance Principle

Theorem 2.2 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 2.2** (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 2.2 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 2.2 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 2.2 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 2.2 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.

### 2.3.1 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 2.3** (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  $C$ .*

### **i** Note

Theorem 2.3 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.

## 2.4 Copulas for Continuous and Discrete Data

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 the continuous case, 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.*



### 2.4.1 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 2.3. 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). 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 the multilinear extension of the subcopula, which leads to a copula that is commonly known as *Checkerboard Copula*.

## 2.5 Checkerboard Copula

In this section, we will define Checkerboard Copula (Genest et al. 2014), which is just a multilinearly interpolated copula on  $[0, 1]^d$  constructed from any subcopula  $C^S$ . Before we dive into this definition, let's settle notations for some of the tools that we will be using.

Let  $\mathbf{X} = (X_1, \dots, X_d)^T$  be a  $d$ -dimensional random vector with joint cumulative distribution function (c.d.f.)  $H$ , and arbitrary marginal c.d.f.s  $F_1, \dots, F_d$ .

We know by Theorem 2.2, that there exists at least one copula  $C : [0, 1]^d \rightarrow [0, 1]$  such that  $H(x_1, \dots, x_d) = C(F_1(x_1), \dots, F_d(x_d))$ , where  $x_1, \dots, x_d \in \mathbb{R}$ . As mentioned before, if the marginal c.d.f.s are continuous, then  $C$  is unique, and that  $C = \mathbf{F}(\mathbf{X}) = (F_1(X_1), \dots, F_d(X_d))$ . But in the case of discrete data, it is only uniquely determined on the domain  $\Pi_{j=1}^d \text{Ran}(F_j)$ .

Now, let's suppose that each  $X_j$  in  $\mathbf{X}$  is an ordinal variable with finite number of categories  $I_j$ , which we can denote with numbers through a fixed ordering  $\{x_1^j < \dots < x_{i_j}^j < \dots < x_{I_j}^j\}$ . We need this ordering to represent our data in a  $d$ -way contingency table format that classifies our observations with respect to the  $d$ -variables each representing an axis, and a fixed ordering making sure that we can unambiguously locate the observations with a certain combination of categories within the contingency table.

**Definition 2.4** (Joint p.m.f. in a  $d$ -way contingency table). (Wei and Kim 2021) A **joint probability mass function (p.m.f.)** of  $\mathbf{X}$  in the  $d$ -way contingency table is an array of size  $\Pi_{j=1}^d I_j$  defined as:

$$P = \{p_{\mathbf{i}} = p_{i_1, \dots, i_d} = \Pr(X_1 = x_{i_1}^1, \dots, X_d = x_{i_d}^d) | i_j \in \{1, \dots, I_j\}, j \in \{1, \dots, d\}, \mathbf{i} \in \Pi_{j=1}^d \mathbb{I}_j\}$$

where  $\mathbb{I}_j = \{1, \dots, I_j\}$  is the index set for  $X_j$ ,  $\mathbf{1} = (1, \dots, 1)^T$  and  $\mathbf{I} = (I_1, \dots, I_d)^T$  denote

index vectors of length  $d$ .

**i** Note

As a consequence of the law of total probability, based on our notation, we also have

$$\text{that: } \sum P = \sum_{\mathbf{i}=1}^{\mathbf{I}} p_{\mathbf{i}} = \sum_{i_1=1}^{I_1} \cdots \sum_{i_d=1}^{I_d} p_{i_1, \dots, i_d} = 1$$

**Definition 2.5** (Marginal p.m.f. in a  $d$ -way contingency table). (Wei and Kim 2021) A **marginal probability mass function (p.m.f.)** of  $i_j$ -th entry in  $X_j$  is:

$$p_{+i_j+} = \sum_{i_1=1}^{I_1} \cdots \sum_{i_{j-1}=1}^{I_{j-1}} \sum_{i_{j+1}=1}^{I_{j+1}} \cdots \sum_{i_d=1}^{I_d} p_{i_1, \dots, i_d} = \sum_{\mathbf{i}_{-j}=\mathbf{1}_{-j}}^{\mathbf{I}_{-j}} p_{\mathbf{i}}$$

where  $\mathbf{i}_{-j}$ ,  $\mathbf{1}_{-j}$ , and  $\mathbf{I}_{-j}$  denote index vectors of  $\mathbf{i}$ ,  $\mathbf{1}$ , and  $\mathbf{I}$  with  $j$ -th entry omitted, respectively.

Thus, **marginal p.m.f.** for the  $(d - 1)$ -dimensional random vector without  $X_j$  ( $\mathbf{X}_{-j} = (X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_d)^T$ ) is denoted by

$$p_{i_1, \dots, +j, \dots, i_d} = \sum_{i_j=1}^{I_j} p_{i_1, \dots, i_d}$$

**Definition 2.6** (Conditional p.m.f. in a  $d$ -way contingency table). (Wei and Kim 2021) A **conditional probability mass function (p.m.f.)** of  $X_j$  given  $\mathbf{X}_{-j}$  is:

$$p_{i_j|\mathbf{i}_{-j}} = \frac{p_{i_1, \dots, i_d}}{p_{i_1, \dots, +j, \dots, i_d}}$$

Also, let's denote (finite and discrete) range of the marginal distribution of  $X_j$  to be  $D_j = \{u_0^j, \dots, u_{i_j}^j, \dots, u_{I_j}^j\}$ . Then,  $u_0^j = 0$ ,  $u_{I_j}^j = 1$ , and  $u_{i_j}^j = \sum_{k_j=1}^{i_j} p_{+k_j+}$ .

**i** Note

If the superscript in  $u_{ij}^j$  can be trivially deduced from the subscript, then for notational ease, the superscript can be omitted.

Then, by Theorem 2.2, the unique subcopula  $C^S$  associated with  $d$ -dimensional random vector  $\mathbf{X}$  over  $\prod_{j=1}^d D_j$  is given by  $H(x_{i_1}^1, \dots, x_{i_d}^d) = \sum_{k_1 \leq i_1} \dots \sum_{k_d \leq i_d} p_{k_1, \dots, k_d} = C^S(u_{i_1}^1, \dots, u_{i_d}^d)$ .

As mentioned before, any subcopula  $C^S$  on  $\prod_{j=1}^d D_j$  can be extended to a copula  $C$  on  $[0, 1]^d$  by multilinear interpolation, Wei and Kim (2021) we define checkerboard copula as follows:

**Definition 2.7** (Checkerboard Copula). (Wei and Kim 2021)

Let  $C^S$  be a subcopula on  $\prod_{j=1}^d D_j$  satisfying as defined above. For any  $\mathbf{u} = (u_1, \dots, u_d) \in [0, 1]^d$ , let  $u_j^\ell$  and  $u_j^u$  be, respectively, the least and greatest elements of  $\bar{D}_j$  (the closure of set  $D_j$ ) satisfying  $u_j^\ell \leq u_j \leq u_j^u$ . Note that if  $u_j$  is in  $D_j$ , then  $u_j^\ell = u_j^u$ . Furthermore, for any  $S \subseteq \{1, \dots, d\}$ , let

$$\lambda_j(u_j) = \begin{cases} \frac{u_j - u_j^\ell}{u_j^u - u_j^\ell}, & \text{if } u_j^\ell < u_j^u \\ 1, & \text{if } u_j^\ell = u_j^u \end{cases} \quad \text{and} \quad \lambda_S(u_1, \dots, u_d) = \prod_{i \in S} \lambda_i(u_i) \prod_{i \notin S} (1 - \lambda_i(u_i)).$$

Then, the **checkerboard copula**  $C^+$  of the ordinal random vector  $\mathbf{X}$  is defined as

$$C^+(\mathbf{u}) = C^+(u_1, \dots, u_d) = \sum_{S \subseteq \{1, \dots, d\}} \lambda_S(u_1, \dots, u_d) C^S(u_{s_1}, \dots, u_{s_d}),$$

where  $u_{s_j} = u_j^u$  if  $j \in S$  and  $u_{s_j} = u_j^\ell$  otherwise.

Additionally, by taking the derivatives of  $C^+$  with respect to  $u_1, \dots, u_d$ , the **checkerboard**

**copula density function** is defined to be

$$c^+(\mathbf{u}) = c^+(u_1, \dots, u_d) = \frac{p_{i_1, \dots, i_d}}{\prod_{j=1}^d p_{+i_j+}}, \quad \text{where } u_{i_j-1}^j < u_j \leq u_{i_j}^j.$$

## 2.6 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 ( $\rho_{\text{pearson}}$ ). 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:

### 2.6.1 Pearson's Correlation Coefficient ( $\rho_{\text{pearson}}$ ) & its Properties

**Definition 2.8** (Pearson correlation coefficient). Given a random vector  $(X, Y)$  with  $\text{Var}(X) < \infty$  and  $\text{Var}(Y) < \infty$ , then:

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

, 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  $Var(X) = \mathbb{E}((X - \mathbb{E}(X))^2)$ .

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

1.  $\rho_{pearson} \in [-1, 1]$
2.  $|\rho_{pearson}(X, Y)| = 1$  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  $\rho_{pearson}(X, Y) = -1$ , and  $a > 0$  if and only if  $\rho_{pearson}(X, Y) = 1$ .  
In both cases,  $X, Y$  are called *perfectly linearly dependent*
3. If  $X$  and  $Y$  are independent, then  $\rho_{pearson}(X, Y) = 0$ .
4.  $\rho_{pearson}$  is invariant under *strictly increasing linear* transformations.

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

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

#### 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  $\text{Ran}Y$ .

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

Let  $X \sim U(0, 1)$  and define  $Y = \log(X)$ . Pearson's correlation is:  $\rho_{pearson}(X, Y) = \frac{\text{Cov}(X, \log X)}{\sigma_X \sigma_Y}$ .

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.

**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:  $\text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] = \mathbb{E}[X^3] - (0)(\frac{1}{3})$ . Since  $\mathbb{E}[X^3] = 0$ , we get  $\text{Cov}(X, Y) = 0$ . 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.

### 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:  $\rho_{\text{pearson}}(X_1, X_2) = 0.5$ .

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

5. **Unattainability of Certain Correlations:** Given margins  $F_1, F_2$ , some  $\rho_{\text{pearson}} \in [-1, 1]$  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 ( $\rho_{\text{spearman}}$ ) and Kendall's Tau ( $\tau_{\text{kendall}}$ ) 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  $(x_1 - x_2)(y_1 - y_2) < 0$ .

### 2.6.3 Kendall's Tau

**Definition 2.9** (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 2.9, 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  $\tau_{kendall}(C) = 4 \int_{[0,1]^2} C(u, v) d(C(u, v)) - 1$ .

### 2.6.4 Spearman's Rho

**Definition 2.10** (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 2.1.

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

#### 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 2.1

So far, we have seen that in the continuous case, we have association measures that are already model-free, and margin-free, but this is not the case in discrete situations especially because the representations of Spearman's rho and Kendall's tau in terms of copula  $C$  are not well-defined due to non-uniqueness of copula.

There has been significant research towards the development of model-free and margin-free measures for the discrete case. In the next sub-section, we will discuss some such measures for 2-dimensional discrete data from Denuit and Lambert (2005), Genest et al. (2014), Genest et al. (2017), and Nešlehová (2007).

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

Before understanding the 2-D modified versions of Kendall's tau and Spearman's rho that overcome the limitations induced by having “gaps” in the discrete case, let's define our a simpler notation for checkerboard copula in 2-D.

For a bivariate random vector  $(X_1, X_2)$  with discrete marginals  $F_{X_1}$  and  $F_{X_2}$ , we can construct the checkerboard copula  $C^+$  through the following process:

1. Start with the empirical copula  $C$  defined on the range of  $(F_{X_1}(X_1), F_{X_2}(X_2))$
2. Extend this copula to the entire unit square by creating a piecewise constant function over rectangles defined by the discontinuity points of the marginals
3. The resulting in a (visualisable) “checkerboard” pattern, which gives the approach its name

Mathematically, the resulting 2-D checkerboard copula can be represented as:

$$C^+(u, v) = \sum_i \sum_j C(a_i, b_j) \cdot \mathbf{1}_{(a_i, a_{i+1}] \times (b_j, b_{j+1}]}(u, v)$$

where  $a_i$  and  $b_j$  are the jump points of the marginal distributions, and  $\mathbf{1}$  is the indicator function.

### Checkerboard Version of Kendall's Tau

Denuit and Lambert (2005) and Nešlehová (2007) have shown that a modified version of Kendall's tau based on the checkerboard copula can be defined as:

$$\tau^+(X_1, X_2) = 4 \int_{[0,1]^2} C^+(u, v) dC^+(u, v) - 1$$

This measure maintains many desirable properties of the continuous case tau while being properly normalized for discrete data.

### Checkerboard Version of Spearman's Rho

Similarly, a checkerboard version of Spearman's rho can be defined as:

$$\rho^+(X_1, X_2) = 12 \int_{[0,1]^2} C^+(u, v) dudv - 3$$

### Strengths and Limitations of 2-D Checkerboard Approaches

The checkerboard-based association measures for discrete data represent a significant advancement in the field of dependence modeling. By addressing the “gap problem” inherent in discrete distributions, these measures provide more accurate and reliable quantification of dependence structures compared to their classical counterparts.

Empirical studies by Genest et al. (2017) demonstrate that the checkerboard versions consistently outperform traditional measures when applied to discrete data with varying degrees of ties. Their simulation results showed that  $\tau^+$  and  $\rho^+$  maintain proper type I error rates in hypothesis testing scenarios, whereas the traditional measures can be overly conservative or,

in certain cases, anti-conservative.

Furthermore, Nešlehová (2007) established the asymptotic properties of these estimators, proving their consistency and deriving their limiting distributions under the null hypothesis of independence. This theoretical foundation enhances the statistical validity of inference procedures based on these measures.

The checkerboard-based measures mentioned above offer several important advantages:

1. **Proper Normalization:** Unlike traditional dependence measures, checkerboard-based measures are inherently bounded between -1 and 1. This boundedness simplifies comparison and interpretation, clearly indicating the direction and strength of dependence. The extreme values (-1 or 1) are specifically attained under conditions of perfect monotone dependence, making them particularly useful for identifying maximal dependence scenarios accurately.
2. **Invariance:** A critical strength of checkerboard-based measures is their invariance under strictly increasing transformations of the marginal distributions. This means that these measures remain unaffected by transformations such as scaling or monotone reshaping of the data. Consequently, analyses based on these measures are robust, facilitating comparisons across datasets and ensuring reliability regardless of the chosen scale.
3. **Interpretability:** Checkerboard-based dependence measures consistently maintain interpretability irrespective of data discreteness. Traditional measures often struggle to offer clear interpretations when data exhibits discrete characteristics, leading to ambiguity. In contrast, checkerboard-based approaches offer intuitive interpretations, ensuring clarity whether the data are discrete, continuous, or mixed, thus providing

valuable insights in diverse practical applications.

4. **Compatibility:** An essential advantage of checkerboard-based measures is their compatibility with continuous marginal distributions. When applied to purely continuous data, these checkerboard-based measures naturally simplify to the standard, widely-used measures of dependence. This property ensures seamless integration of these methods within existing analytical frameworks and allows for straightforward comparisons between traditional and checkerboard-based results.

Additionally, when implementing these checkerboard-based association measures it is important to note that computational complexity increases with the number of distinct values in each marginal.

However, it is important to note that despite these advancements, the measures discussed thus far are primarily designed for bivariate (2-dimensional) data. In many real-world applications across fields such as bioinformatics, econometrics, and multivariate time series analysis, we encounter high-dimensional discrete data where pairwise association measures may not capture the full complexity of the dependence structure.

### 2.6.6 Looking Forward: Beyond Bivariate Measures

Although these 2-dimensional checkerboard approaches represent significant progress in the field, modern data analysis increasingly demands tools that can handle multivariate discrete data effectively. In Chapter 3, we will review Wei and Kim (2021)'s groundbreaking work on association measures that are valid for  $n$ -dimensional discrete data. Their framework extends the checkerboard concept to higher dimensions while preserving the desirable properties of margin-freedom and model-independence.

Chapter Chapter 4 will then present our contribution to the field: the first-ever Python package implementation of these n-dimensional measures. By making these advanced statistical tools accessible to practitioners through an open-source, well-documented, and computationally efficient Python package, we aim to bridge the gap between theoretical advances and practical applications.

This implementation serves to enable researchers across disciplines to analyze complex discrete multivariate data without the limitations imposed by traditional correlation measures or the need for parametric assumptions. The package includes comprehensive visualization tools, statistical testing procedures, and integration capabilities with popular data science frameworks.

By developing this implementation, we hope to facilitate wider adoption of these powerful association measures in fields ranging from genomics to social network analysis, where discrete multivariate data is abundant but appropriate analysis tools have been limited.

## **Chapter 3**

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

In this chapter, we will review and combine concepts from Wei and Kim (2021) and Liao et al. (2024) in order to ultimately define regression based on checkerboard copula along with some visualization and association measures for model-free regression dependence analysis of multivariate discrete data.



### 3.1 Data

Before we dive into details of the method and some pre-requisites, let's understand what is the type of data of interest. The methods described in this chapter are applicable on multivariate categorical data in the form of multi-dimensional contingency table with an ordinal response variable and a set of categorical (nominal/ordinal) predictors.

#### 3.1.1 2-D Example

(Adapted from Wei and Kim 2021)

Consider a dataset that contains two ordinal variables: the dose of a treatment drug for acute migraine ( $X_1$ ) with  $I_1 = 5$  categories ( $(x_1^1, x_2^1, x_3^1, x_4^1, x_5^1) = (\text{very low, low, medium, high, very high})$ ), and the severity of migraine pain recorded after treatment ( $X_2$ ) with  $I_2 = 3$  categories ( $(x_1^2, x_2^2, x_3^2) = (\text{mild, moderate, severe})$ ).

Table 3.1: Joint p.m.f of  $X_1$  and  $X_2$ ,  $P = \{p_{i_1 i_2}\}$ .

$X_1 \backslash X_2$	$x_1^2$	$x_2^2$	$x_3^2$
$x_1^1$	0	0	2/8
$x_2^1$	0	1/8	0
$x_3^1$	2/8	0	0
$x_4^1$	0	1/8	0
$x_5^1$	0	0	2/8

This example has been carefully constructed so that  $X_2$  has a quadratic relationship with  $X_1$  (as level of  $X_2$  decreases, the level of  $X_1$  increases). We can also observe that  $X_2$  is a function of  $X_1$  with probability 1, but not the other way around (for a given category of  $X_1$ , there is

one and only one category of  $X_2$  whose corresponding joint probability is non-zero).

## 3.2 Checkerboard Copula and its Density

Recall that checkerboard copula and the corresponding copula density is defined in Chapter 2's Definition 2.7. In order to understand it better, let's continue building upon the example data above by defining the copula and its corresponding density.

### 3.2.1 2-D Example (continued...)

(Adapted from Wei and Kim 2021)

Upon application of Definition 2.5, we can see that the marginal p.m.f.s of  $X_1$  and  $X_2$  are  $p_{i_1+} \in \{2/8, 1/8, 2/8, 1/8, 2/8\}$  and  $p_{+i_2} \in \{2/8, 2/8, 4/8\}$ , respectively. Furthermore, we can see that the ranges of the marginal c.d.f.s of  $X_1$  and  $X_2$  are  $D_1 = \{u_0^1, u_1^1, u_2^1, u_3^1, u_4^1, u_5^1\} = \{0, 2/8, 3/8, 5/8, 6/8, 1\}$  and  $D_2 = \{u_0^2, u_1^2, u_2^2, u_3^2\} = \{0, 2/8, 4/8, 1\}$ , respectively. As per Definition 2.7, we can visualize checkerboard copula density of  $X_1$  and  $X_2$  in Figure 3.1.

## 3.3 Checkerboard Copula Score

Ordinal variables contain categories with natural ordering but unknown distances between them. When analyzing associations in ordinal contingency tables, we can leverage this inherent ordering information. To achieve this, Wei and Kim (2021) introduce checkerboard copula scores derived from the checkerboard copula.

As established in Definition 2.7, the checkerboard copula represents a smoothed version of the subcopula associated with ordinal random vector  $X$ . It distributed probability mass uniformly across d-dimensional hyperrectangles in  $[0, 1]^d$ , specifically within intervals  $[u_{i_{j-1}}^j, u_{i_j}^j]$ , where

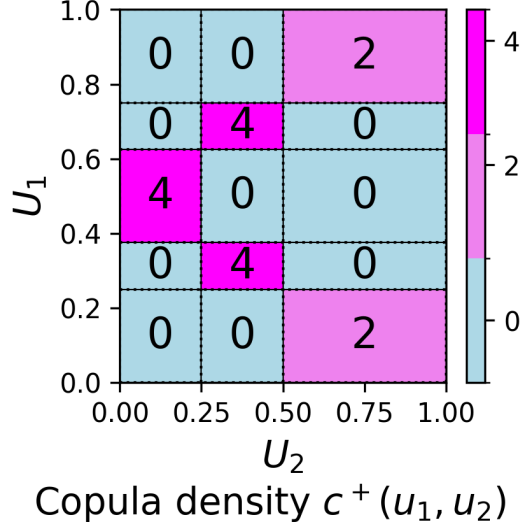


Figure 3.1: Checkerboard Copula Density Visualization for 2-D Example (This figure's styling is designed to reproduce the work presented by Wei and Kim, 2021 in supplementary materials)

$u_{i_j}^j$  is defined by the marginal cumulative distribution functions. Furthermore, we can define a transformation of  $X_j$  via  $U_j$  as  $S_j = \mathbb{E}[U_j|X_j]$ , where  $j \in \{1 \dots d\}$ . Note that here,  $S_j$  is an ordinal random variable with numerical support values  $\{s_1^j, \dots, s_{i_j}^j, \dots, s_{I_j}^j\}$ , where  $s_{i_j}^j = (u_{i_{j-1}}^j + u_{i_j}^j)/2$ .

**Definition 3.1** (Checkerboard Copula Scores (CCS)). (Wei and Kim 2021) The **checkerboard copula scores (CCS)** of ordinal variable  $X_j$  are  $\{s_1^j, \dots, s_{i_j}^j, \dots, s_{I_j}^j\}$ , where  $s_{i_j}^j = (u_{i_{j-1}}^j + u_{i_j}^j)/2$  for  $i_j \in \{1, \dots, I_j\}$  and  $u_{i_j}^j$  as defined in Section 2.5. In other words, CCS is a set of the average of the marginal distributions evaluated at every two consecutive categories of  $X_j$ .

There are several interesting properties of these scores as proven by Wei and Kim (2021). One that is of our interest is the formula for mean and variance of the support vector  $S_j$ .

**Lemma 3.1** (Mean and Variance of Support Vector ( $S_j$ )). (Wei and Kim 2021) The probability-weighted mean (or expected value) of  $S_j$  is  $\mu_{S_j} = 0.5$ . The variance of  $S_j$  is  $\sigma_{S_j}^2 = \frac{1}{4} \sum_{i_j=1}^{I_j} u_{i_{j-1}}^j u_{i_j}^j p_{+i_j+}$ .

### 3.3.1 2-D Example (continued...)

(Adapted from Wei and Kim 2021)

Upon application of Definition 3.1, we obtain the checkerboard copula scores of  $X_1$  and  $X_2$  as  $(2/16, 5/16, 8/16, 11/16, 14/16)$  and  $(2/16, 6/16, 12/16)$  respectively. Furthermore, by Lemma 3.1, we can say that the  $(\mu_{S_j}, \sigma_{S_j}^2)$  of  $S_1$  and  $S_2$  are  $(0.5, 81/1024)$  and  $(0.5, 9/128)$  respectively.

### 3.3.2 Empirical Estimation of CCS

Since we are dealing with discrete count data, at times, we generally do not have access to the joint probability matrix, instead we have to base our analysis on the counts of observations with various categorical combinations of variables of interest. Thus, below, we establish the missing link between count data and pre-defined distributions as introduced by Wei and Kim (2021).

Let  $\{n_{i_1, \dots, i_d}\}$ ,  $i_j \in \{1, \dots, I_j\}$ ,  $j \in \{1, \dots, d\}$ , denote counts in a  $d$ -way contingency table obtained by classifying  $n = \sum_{i_1=1}^{I_1} \dots \sum_{i_d=1}^{I_d} n_{i_1, \dots, i_d}$  observations (or **cases**) into categories of  $d$  variables,  $X_1, \dots, X_d$ .

Let's define marginal sums of  $i_j$ -th category in  $X_j$  as  $n_{+i_j+} = \sum_{i_1=1}^{I_1} \dots \sum_{i_{j-1}=1}^{I_{j-1}} \sum_{i_{j+1}=1}^{I_{j+1}} \dots \sum_{i_d=1}^{I_d} n_{i_1, \dots, i_d}$ , and  $(d-1)$ -variate marginal frequencies of  $\mathbf{X}_{-j}$  as  $n_{i_1, \dots, +j, \dots, i_d} = \sum_{i_j=1}^{I_j} n_{i_1, \dots, i_d}$ .

In terms of these, we can define estimators for the probabilities as follows:

$$\hat{p}_{i_1, \dots, i_d} = \frac{n_{i_1, \dots, i_d}}{n}, \hat{p}_{+i_j+} = \frac{n_{+i_j+}}{n}, \hat{p}_{i_1, \dots, +j, \dots, i_d} = \frac{n_{i_1, \dots, +j, \dots, i_d}}{n}, \hat{p}_{i_j|i_{-j}} = \frac{\hat{p}_{i_1, \dots, i_d}}{\hat{p}_{i_1, \dots, +j, \dots, i_d}}$$

Moreover, the range of marginal c.d.f. of  $X_j$  is estimated by  $[\hat{u}_0^j, \dots, \hat{u}_{i_j}^j, \dots, \hat{u}_{I_j}^j]$  with  $\hat{u}_0^j = 0$

and  $u_{i_j}^j = \sum_{k_j=1}^{i_j} \hat{p}_{+k_j+}$ .

Using above-established pre-requisites, Definition 3.1, and Lemma 3.1 we can estimate the checkerboard copula scores  $\hat{s}_1^j, \dots, \hat{s}_{I_j}^j$  with  $\hat{s}_{i_j}^j = (\hat{u}_{i_j-1}^j + \hat{u}_{i_j}^j)/2$  and  $\hat{\sigma}_{\hat{s}_j}^2 = \sum_{i_j=1}^{I_j} \hat{u}_{i_j-1}^j \hat{u}_{i_j}^j \hat{p}_{+i_j+}/4$ .

### 3.4 Checkerboard Copula Regression, Prediction and Visualization

Let  $\mathbf{U}$  be a uniform random vector on  $[0, 1]^d$  associated with the checkerboard copula  $C^+$  for a  $d$ -way ordinal contingency table.

**Definition 3.2** ( $(d-1)$ -Marginal Density). (Wei and Kim 2021) The  $(d-1)$ -**marginal density** for  $\mathbf{U}_{-j} = (U_1, \dots, U_{j-1}, U_{j+1}, \dots, U_d)^T$  is defined as

$$c^+(\mathbf{u}_{-j}) = \frac{p_{i_1, \dots, +j, \dots, i_d}}{\prod_{k=1, k \neq j}^d p_{+i_k+}}$$

where  $\mathbf{u}_{-j} = (u_1, \dots, u_{j-1}, u_{j+1}, \dots, u_d)^T$  in  $[0, 1]^{d-1}$  and  $u_{i_k-1}^k < u_k < u_{i_k}^k$ . Here,  $k \in \{1, \dots, j-1, j+1, \dots, d\}$ ,  $j \in \{1, \dots, d\}$ ,  $p_{i_1, \dots, +j, \dots, i_d}$ , and  $p_{+i_k+}$  as in Definition 2.5.

**Definition 3.3** (Conditional Density of  $U_j$  given  $\mathbf{U}_{-j}$ ). (Wei and Kim 2021) The **conditional density of  $U_j$  given  $\mathbf{U}_{-j}$** , where  $\mathbf{U}_{-j} = (U_1, \dots, U_{j-1}, U_{j+1}, \dots, U_d)^T$  is defined as

$$c^+(u_j | \mathbf{u}_{-j}) = \frac{c^+(\mathbf{u})}{c^+(\mathbf{u}_{-j})} = \frac{p_{i_j | i_{-j}}}{p_{+i_j+}}$$

where  $\mathbf{u}_{-j} = (u_1, \dots, u_{j-1}, u_{j+1}, \dots, u_d)^T$  in  $[0, 1]^{d-1}$  and  $u_{i_k-1}^k < u_k < u_{i_k}^k$ . Here,  $j \in \{1, \dots, d\}$ ,  $p_{i_j | i_{-j}}$ , and  $p_{+i_j+}$  as in Definition 2.6 and Definition 2.5, respectively.

Now, as mentioned in Wei and Kim (2021), we can define checkerboard copula regression

function as follows

**Definition 3.4** (Checkerboard Copula Regression (CCR)). (Wei and Kim 2021) The **checkerboard copula regression function** of  $U_j$  on  $\mathbf{U}_{-j}$  is defined as

$$r_{U_j|\mathbf{U}_{-j}}(\mathbf{u}_{-j}) \equiv E_{c^+}(U_j|\mathbf{U}_{-j} = \mathbf{u}_{-j}) = \int_0^1 u_j c^+(u_j|\mathbf{u}_{-j}) du_j = \sum_{i_j=1}^{I_j} p_{i_j|\mathbf{i}_{-j}} s_{i_j}^j$$

In other words, CCR function is the mean checkerboard score of  $X_j$  with respect to the conditional distribution at the category  $\mathbf{i}_{-j}$  of  $(d - 1)$  explanatory variables  $\mathbf{X}_{-j}$ .

### 3.4.1 2-D Example (continued...)

(Adapted from Wei and Kim 2021)

Upon application of the above definitions, we obtain the following tabular representations of conditional p.m.f.s and checkerboard copula regressions.

Table 3.2: Conditional p.m.f of  $X_2$  given  $X_1$

$X_1 \backslash X_2$	$x_1^2$	$x_2^2$	$x_3^2$
$x_1^1$	0	0	1
$x_2^1$	0	1	0
$x_3^1$	1	0	0
$x_4^1$	0	1	0
$x_5^1$	0	0	1

Table 3.3: Conditional p.m.f of  $X_1$  given  $X_2$

$X_1 \backslash X_2$	$x_1^2$	$x_2^2$	$x_3^2$
$x_1^1$	0	0	1/2
$x_2^1$	0	1/2	0
$x_3^1$	1	0	0
$x_4^1$	0	1/2	0
$x_5^1$	0	0	1/2

Table 3.4: Checkerboard copula regression of  $U_2$  on  $U_1$

$u_1$	$r_{U_2 U_1}(u_1)$
$[0, 2/8]$	12/16
$(2/8, 3/8]$	6/16
$(3/8, 5/8]$	2/16
$(5/8, 6/8]$	6/16
$(6/8, 1]$	12/16

Table 3.5: Checkerboard copula regression of  $U_1$  on  $U_2$

$u_2$	$r_{U_1 U_2}(u_2)$
$[0, 2/8]$	1/2
$(2/8, 4/8]$	1/2
$(4/8, 1]$	1/2

### 3.4.2 Point Prediction Using CCR

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. Thus, we can use Definition 3.4 for predicting the category of response variable for a given combination of categories of explanatory variables, while describing the dependence structure between them.

Suppose that  $X_j$  is the response variable, and all the remaining variables in the table (denoted by  $\mathbf{X}_{-j}$ ) are to be used as predictors. Recall Definition 2.6, where we denote (finite and discrete) range of the marginal distribution of  $X_j$  to be  $D_j = \{u_0^j, \dots, u_{i_j}^j, \dots, u_{I_j}^j\}$ . Then,  $u_0^j = 0$ ,  $u_{I_j}^j = 1$ , and  $u_{i_j}^j = \sum_{k=1}^{i_j} p_{+k_j+}$ . As mentioned in Wei and Kim (2021), we can use this to find  $\mathbf{u}_{-j}^*$  from  $\text{Ran}(\mathbf{X}_{-j}) = \prod_{k=1, k \neq j}^d D_k$ . Using this along with Definition 3.4 gives us the estimated value of the checkerboard copula regression,  $u_j^* = r_{U_j|\mathbf{U}_{-j}}(\mathbf{u}_{-j}^*)$ . Now, using this we can obtain  $i_j^*$  and  $u_{i_j^*}^j$  such that  $u_{i_j^*-1}^j < u_j^* < u_{i_j^*}^j$ . This finally leads us to the predicted category  $x_{i_j^*}^j$  of the response variable  $X_j$ .

In order to better understand this, let's walk through the example at hand.

### 3.4.3 2-D Example (continued...)

(Adapted from Wei and Kim 2021)

Upon application of the method detailed above, we can predict the category of  $X_2$  for each category of  $X_1$ . For instance, given that  $X_1 = x_3^1 = x_{i_1=3}^{1*}$ , the corresponding  $u_3^{1*} = 5/8$ , and thus the predicted value of the CCR is  $u_2^* = r_{U_2|U_1}(5/8) = 1/8 \in [0, 2/8]$ . This implies that  $i_2^* = 1$  and  $u_{i_2^*=1}^2 = 2/8$  because  $u_0^2 = 0 < u_2^* = 1/8 \leq u_1^2 = 2/8$ . Hence, the predicted category of  $X_2$  given  $X_1 = x_3^1$  is  $f_{X_2|X_1}(x_3^1) = x_1^2$ .



After applying this method to all combinations of predictors and response, we obtain the following tabular representations of point predictions through CCR.

Table 3.6: Point prediction through CCR of  $X_2$  on  $X_1$

$X_1$	$u_2^*$	$f_{X_2 X_1}$
$x_1^1$	6/8	$x_3^2$
$x_2^1$	3/8	$x_2^2$
$x_3^1$	1/8	$x_1^2$
$x_4^1$	3/8	$x_2^2$
$x_5^1$	6/8	$x_3^2$

Table 3.7: Point prediction through CCR of  $X_1$  on  $X_2$

$X_2$	$u_1^*$	$f_{X_1 X_2}$
$x_1^2$	1/2	$x_3^1$
$x_2^2$	1/2	$x_3^1$
$x_3^2$	1/2	$x_3^1$

We can clearly see from the above tables how the prediction results reflect the quadratic relationship shown when we first established this 2-D example at the start of this chapter.

#### 3.4.4 Empirical Estimation of CCR and Point Prediction

Now, continuing from the notation established in Section 3.3.2, we can estimate CCR for

$$k \in 1, \dots, j-1, j+1, \dots, d,$$

$$\hat{r}_{U_j|U_{-j}}(\mathbf{u}_{-j}) = \sum_{i_j=1}^{I_j} \hat{p}_{i_j|I_{-j}} \hat{s}_{i_j}^j \text{ for } \hat{u}_{i_k-1} < u_k \leq \hat{u}_{i_k}^k$$

Now, using the above alongside the steps we mentioned in previous subsection, we can obtain the predicted category of a response variable for each combination categories of predictors. That is, for a given combination categories of the  $(d - 1)$ -predictors  $\mathbf{X}_{-j}$ , we find the corresponding  $\hat{\mathbf{u}}_{-j}^*$  from the estimated ranges of  $\mathbf{X}_{-j}$  and then obtain the estimated value of the CCR as mentioned above,  $\hat{u}_j^* = \hat{r}_{U_j|U_{-j}}(\hat{\mathbf{u}}_{-j}^*)$ . From the estimated range of a response variable  $X_j$ , we get  $i_j^*$  and  $\hat{u}_{i_j^*}^j$  such that  $\hat{u}_{i_j^*-1}^j < \hat{u}_j^* \leq \hat{u}_{i_j^*}^j$ . This implies that the predicted category of  $X_j$  is  $\hat{x}_{i_j^*}^j$ .

More details on these estimators including asymptotic analysis are outside the scope of this work, but are discussed at length in Wei and Kim (2021).

### 3.4.5 Uncertainty Evaluation of the CCR prediction using nonparametric bootstrap

To quantify the uncertainty of the predicted category obtained from CCR, we can use nonparametric bootstrap. This involves generating multiple bootstrap samples from the original contingency table, computing the checkerboard copula regression for each resampled dataset, and then predicting the category of the response variable for each combination of categories of the explanatory variables.

For each bootstrap sample, we follow these steps:

1. Resample with replacement from the original data to create a bootstrap sample with the same size as the original data.

2. Estimate the checkerboard copula regression based on this bootstrap sample
3. Predict the category of the response variable using the estimated regression
4. Repeat steps 1-3 multiple times (e.g., 1000 times)

The distribution of predicted categories across bootstrap samples provides a measure of the prediction uncertainty. We can calculate the proportion of times each category is predicted for a given combination of explanatory variables, with higher proportions indicating greater confidence in the prediction.

For example, in our 2-D contingency table, using 1000 bootstrap resamples, we quantified the uncertainty of the predicted category of  $X_2$  for each category of  $X_1$ . The results showed that the proportion of bootstrap samples where the predicted category matched our original prediction was 100%, indicating high confidence in our predictions.

### 3.5 Checkerboard Copula Regression Association Measure

Using CCR discussed above, Wei and Kim (2021) proposed the checkerboard copula regression based association measure (CCRAM) for a multi-way contingency table with an ordinal response variable and categorical (ordinal or nominal) explanatory variables.

**Definition 3.5** (Checkerboard Copula Regression-Based Association Measure (CCRAM)).

(Wei and Kim 2021) The **checkerboard copula regression-based association measure (CCRAM)** of  $X_j$  on  $\mathbf{X}_{-j} = (X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_d)^\top$  is

$$\rho_{(\mathbf{X}_{-j} \rightarrow X_j)}^2 \equiv \frac{\text{Var}[r_{U_j|\mathbf{U}_{-j}}(\mathbf{U}_{-j})]}{\text{Var}(U_j)} = \frac{\text{E} \left[ \left( r_{U_j|\mathbf{U}_{-j}}(\mathbf{U}_{-j}) - 1/2 \right)^2 \right]}{1/12} = 12 \sum_{\mathbf{i}_{-j}=1}^{I_{-j}} \left( \sum_{i_j=1}^{I_j} p_{i_j|\mathbf{i}_{-j}} s_{i_j}^j - 1/2 \right)^2 p_{i_1, \dots, i_j, \dots, i_d}$$

where  $j \in \{1, \dots, d\}$ , and  $U_j$  and  $\mathbf{U}_{-j} = (U_1, \dots, U_{j-1}, U_{j+1}, \dots, U_d)^\top$  are the random variables on  $[0, 1]^d$  associated with the checkerboard copula density  $c^*(\mathbf{u})$  in Definition 2.7.

Extending what is proven for this new measure, Wei and Kim (2021) provides a proposition with proof containing several properties of CCRAM concluding that

- CCRAM can identify linear and non-linear relationship between a response variable and several explanatory variables. CCRAM can also be applied when any predictors are nominal and/or a response variable is binary.
- CCRAM is lower bounded by 0 and upper bounded by  $12\sigma_{S_j}^2$ , where 0 means no contribution of predictors to the construction of the checkerboard copula regression function.

In order to provide a normalized measure that is independent of marginal distribution of  $X_j$ , Wei and Kim (2021) proposes SCCRAM, which is a scaled version of CCRAM, and it is mathematically defined as follows:

**Definition 3.6** (Scaled Checkerboard Copula Regression-Based Association Measure (SC-CRAM)). (Wei and Kim 2021) The **scaled checkerboard copula regression-based association measure (SCCRAM)** of  $X_j$  on  $\mathbf{X}_{-j} = (X_1, \dots, X_{j-1}, X_{j+1}, \dots, X_d)^\top$  is defined as

$$\rho_{(\mathbf{X}_{-j} \rightarrow X_j)}^{2*} = \frac{\rho_{(\mathbf{X}_{-j} \rightarrow X_j)}^2}{12\sigma_{S_j}^2}$$

where  $\rho_{(\mathbf{X}_{-j} \rightarrow X_j)}^2$  and  $\sigma_{S_j}^2$  are defined as in Definition 3.5 and Lemma 3.1 respectively.

(S)CCRAM is designed to quantify the regression association identified by Checkerboard Copula Regression and its prediction.

### 3.5.1 2-D Example (continued...)

(Adapted from Wei and Kim 2021)

Using the above definitions, we obtain:  $(\rho_{(X_1 \rightarrow X_2)}^2, 12\sigma_{S_2}^2 \rho_{(X_1 \rightarrow X_2)}^{2*}) = (27/32, 27/32, 1)$  and  $(\rho_{(X_2 \rightarrow X_1)}^2, 12\sigma_{S_1}^2 \rho_{(X_2 \rightarrow X_1)}^{2*}) = (0, 243/256, 0)$ .

$\rho_{(X_1 \rightarrow X_2)}^{2*} = 1$  implies that  $X_1$  perfectly explains the variation in  $X_2$  induced by its checkerboard copula score and its marginal distribution and this result agrees with the observation that  $r_{U_2|U_1}(u_1)$  equals one and only one of the checkerboard score of  $X_2$ . Note that this result also supports the fact that  $X_2$  is functionally dependent on  $X_1$  with probability 1. On the other hand,  $\rho_{(X_2 \rightarrow X_1)}^{2*} = 0$  means that  $r_{U_1|U_2}(u_2) = E(U_1) = 1/2, \forall u_2$ . Thus,  $X_2$  has no contribution arising from its score and its marginal distribution in explaining the variation in  $X_1$ .

### 3.5.2 Empirical Estimation of (S)CCRAM

Now, continuing from the notation established in Section 3.3.2, the estimators for the CCRAM and SCCRAM as defined in Definition 3.5 and Definition 3.6 respectively are given below:

$$\hat{\rho}_{(\mathbf{X}_{-j} \rightarrow X_j)}^2 = 12 \sum_{i_{-j}=1}^{I_{-j}} \left( \sum_{i_j=1}^{I_j} \hat{p}_{i_j|i_{-j}} \hat{s}_{i_j}^d - \frac{1}{2} \right)^2 \hat{p}_{i_1, \dots, +j, \dots, i_d}, \quad \hat{\rho}_{(\mathbf{X}_{-j} \rightarrow X_j)}^{2*} = \frac{\hat{\rho}_{(\mathbf{X}_{-j} \rightarrow X_j)}^2}{12\hat{\sigma}_{S_j}^2}$$

### 3.5.3 Uncertainty Evaluation of the estimated (S)CCRAM using nonparametric bootstrap distribution and its confidence interval

To assess the uncertainty of the estimated (S)CCRAM, we can employ nonparametric bootstrap to generate an empirical sampling distribution. This approach involves:

1. Generate B bootstrap samples of same size as original data (typically  $B = 1000$ ) by resampling with replacement from the original data.

2. For each bootstrap sample  $b = 1, \dots, B$ , compute the estimated (S)CCRAM, denoted as

$$\hat{\rho}_{(\mathbf{X}_{-j} \rightarrow X_j), b}^2 \text{ and } \hat{\rho}_{(\mathbf{X}_{-j} \rightarrow X_j), b}^{2*}.$$

3. Construct the empirical bootstrap distribution of the estimates  $\{\hat{\rho}_b^2\}_{b=1}^B$ .

From this bootstrap distribution, we can:

Calculate the bootstrap standard error as the standard deviation of the bootstrap estimates, and construct confidence intervals using various methods:

- Percentile Method (Davison and Hinkley 1997): The  $(1 - \alpha) \times 100\%$  confidence interval is given by the  $\alpha/2$  and  $1 - \alpha/2$  quantiles of the bootstrap distribution:  $\left[ \hat{\rho}_{\alpha/2}^2, \hat{\rho}_{1-\alpha/2}^2 \right]$
- Basic (Reverse Percentile) Method (Hesterberg 2014): This method reflects the bootstrap quantiles around the original estimate to produce the interval:  $\left[ 2\hat{\rho}^2 - \hat{\rho}_{1-\alpha/2}^2, 2\hat{\rho}^2 - \hat{\rho}_{\alpha/2}^2 \right]$ , where  $\hat{\rho}^2$  is the (S)CCRAM estimate from the original dataset.
- BCa (Bias-Corrected and Accelerated) Method (Efron 1987): This method adjusts the interval to correct for both bias and skewness in the bootstrap distribution, using acceleration and bias-correction terms derived from the data. It tends to offer better coverage properties, especially in small samples or skewed distributions.

These confidence intervals provide a principled way to quantify the precision of the (S)CCRAM estimates, offering insight into the variability due to sampling and helping assess the statistical significance of regression dependence.

### 3.5.4 Statistical significance of the estimated (S)CCRAM using Permutation distribution and its hypothesis testing

In case of CCRAM, Wei and Kim (2021) proposes permutation testing for null hypothesis  $H_0 : \rho^2_{(\mathbf{x}_{-j} \rightarrow X_j)} = 0$ , which indicates no association between the response variable and the explanatory variables.

The permutation testing procedure involves:

1. Calculate the observed (S)CCRAM value for the original data
2. Generate  $M$  permutation samples (typically  $M = 10^6$ ) by randomly permuting the response variable values while keeping the explanatory variables fixed, thus breaking any potential association
3. For each permutation sample  $m = 1, \dots, M$ , compute the (S)CCRAM, denoted as  $\hat{\rho}^2_{(\mathbf{x}_{-j} \rightarrow X_j),m}$  and  $\hat{\rho}^{2*}_{(\mathbf{x}_{-j} \rightarrow X_j),m}$
4. Construct the empirical permutation distribution of the estimates under the null hypothesis
5. Calculate the p-value as the proportion of permutation statistics that are as extreme as or more extreme than the observed statistic:

$$p\text{-value} = \frac{1}{M} \sum_{m=1}^M I(\hat{\rho}^2_{(\mathbf{x}_{-j} \rightarrow X_j),m} \geq \hat{\rho}^2_{(\mathbf{x}_{-j} \rightarrow X_j),obs})$$

If the p-value is less than a predetermined significance level (e.g., 0.05), we reject the null hypothesis and conclude that there is a significant association between the response variable and the explanatory variables.

This permutation approach provides a distribution-free method for hypothesis testing, analogous to testing  $R^2 = 0$  in linear regression, but appropriate for the categorical data context of the checkerboard copula regression.

## **3.6 Visualization Methods**

As mentioned before, effective visualization is crucial for understanding and interpreting the regression associations identified by checkerboard copula regression. Liao et al. (2024) details several visualization approaches that are particularly suited for displaying the dependence structures in multi-dimensional contingency tables with an ordinal response variable.

### **3.6.1 Cross-tabulation**

For simpler two-way contingency tables (as seen in ice cream study example within Liao et al. (2024)), cross-tabulation provides a straightforward approach to visualize the predicted categories of the response variable for each category of the explanatory variable. One can enhance these tables by color-coding the predicted categories and including bootstrap proportions to represent prediction uncertainty. This allows for clear comparison between the observed regression pattern and the pattern expected under independence.

### **3.6.2 Bubble Plots**

When dealing with higher-dimensional contingency tables (such as the back pain data with three explanatory variables in Wei and Kim (2021)), bubble plots offer an effective visualization approach. In these plots, the x-axis represents different combinations of categories of explanatory variables, while the y-axis shows the categories of the response variable. The predicted category for each combination is indicated by dark dots, and bubbles (circles) with



varying sizes represent the proportion of times each category is predicted across bootstrap samples. This visualization clearly reveals complex association patterns, such as potential interaction effects among explanatory variables.

### **3.6.3 Doubledecker Plots**

For multi-dimensional contingency tables with temporal or hierarchical structure (as in Three Mile Island data within Liao et al. (2024)), doubledecker plots provide a particularly insightful visualization. These plots display vertical splits for explanatory variables and horizontal splits for the response variable. The width of each bar is proportional to the observed frequency of the corresponding combination of explanatory variables, while the heights of color-coded blocks within each bar represent the proportions of predicted categories across bootstrap samples. This approach effectively visualizes both the magnitude and uncertainty of predictions while accounting for the natural ordering or hierarchy among variables.

All these visualization methods can be paired with corresponding “null reference” plots generated through permutation methods, allowing researchers to visually assess whether the detected regression patterns significantly differ from those expected under independence. This combination of exploratory visualization and resampling-based calibration provides a comprehensive framework for understanding regression dependence in categorical data without relying on parametric assumptions.

## Chapter 4

# Software (Package) Implementation and Testing

In this chapter, we introduce `ccrvam`, a Python package that implements the Checkerboard Copula Regression-based Visualization and Association Measure (CCRVAM) techniques discussed in previous chapters. Despite the growing importance of multivariate categorical data analysis with ordinal response variables in disciplines like medicine, social sciences, and economics, there has been a notable lack of user-friendly, well-tested software implementations that scale efficiently to higher-dimensional problems. The `ccrvam` package addresses this gap by providing a comprehensive suite of tools for analyzing multivariate discrete data using the checkerboard copula approach.

### 4.1 Set-up and Example Data

The `ccrvam` package is built for analyzing multi-dimensional contingency tables with an ordinal response variable and a set of categorical (nominal/ordinal) explanatory

variables/predictors. This aligns perfectly with the theoretical frameworks established in Chapter 3. The package is designed with ease of installation and use in mind, particularly within Jupyter Notebook environments common in data analysis workflows.

#### 4.1.1 Installation

For quick use in Jupyter Notebooks, the package can be installed directly from PyPI with:

```
1 pip install ccrvam==0.9.6
```

For more containerized production-heavy work, it's recommended to use a custom virtual environment. Instructions for setting that up can be found in [ccrvam/.github/README.md](https://ccrvam.github.io/README.md).

## 4.2 Types of Input Data Supported

The ccrvam package is designed to be flexible with respect to input data formats. The package supports these as main formats:

#### 4.2.1 Loading from Contingency Table Format (In-Place)

The ccrvam package implements the theoretical framework outlined in Chapter 3 by operating on multi-dimensional contingency tables. These tables represent the joint distribution of categorical variables where: one variable is explicitly designated as an ordinal response (dependent variable); one or more variables function as categorical predictors (independent variables); and table entries contain frequency counts of observations. The package accepts these contingency tables directly as NumPy arrays, with each dimension corresponding to the categories of a particular variable. This structure allows for efficient computation of the checkerboard copula scores and subsequent association measures while preserving the

natural ordering of the response variable categories and accommodating multiple categorical predictors simultaneously.

For example, consider the 2-D example from the previous chapter:

```
1 import numpy as np
2 from ccrvam import GenericCCRVAM
3
4 # Migraine treatment example (dose vs. pain severity)
5 contingency_table_2d = np.array([
6     [0, 0, 20],
7     [0, 10, 0],
8     [20, 0, 0],
9     [0, 10, 0],
10    [0, 0, 20]
11 ])
12
13 # Create a CCRVAM object
14 ccrvam_obj = GenericCCRVAM.from_contingency_table(contingency_table_2d)
15 # Dimension of the inferred joint probability matrix P:
16 print(ccrvam_obj.P.shape)
```

(5, 3)

```
1 # Joint Probability matrix P:
2 print(ccrvam_obj.P)
```

```
[[0.    0.    0.25 ]
 [0.    0.125 0.    ]
 [0.25  0.    0.    ]
 [0.    0.125 0.    ]
 [0.    0.    0.25 ]]
```

For higher dimensions, the package supports multi-dimensional NumPy arrays representing contingency tables across multiple predictors. More examples in this regard can be found in [Chapter 5](#).

#### 4.2.2 Loading from External Data Files

Through `DataProcessor` class in our package (`ccrvam`), we support flexible loading of categorical data in multiple formats, accommodating diverse data structures commonly encountered in statistical analysis. The class provides a unified interface for importing data regardless of its original format, making it accessible for CCRVAM analysis without requiring extensive preprocessing.

The `DataProcessor` class supports three primary data formats:

1. **Case-form data:** Where each row represents an individual case with categorical variables organized in separate columns, allowing for straightforward representation of raw survey or experimental results.
2. **Frequency-form data:** Where each row to contain a unique combination of categorical variables along with their corresponding frequency count, offering a more compact representation when many observations share identical category combinations.
3. **Table-form data:** Where direct contingency tables are represented as multi-dimensional

arrays, which provides the most computationally efficient input format when data has already been aggregated into contingency tables by other statistical software.

The implementation also includes robust handling of several key data management features. The package supports custom variable naming schemes, allowing users to define meaningful labels for their variables rather than relying on default numeric identifiers. It provides automatic mapping of non-integer category values to integer indices, enabling seamless processing of categorical data with text or mixed-type labels. The software accommodates custom delimiters for text-based input files, offering flexibility when importing data from various sources with different formatting conventions. Additionally, the implementation supports dimensional specification for proper array structuring, ensuring that contingency tables are correctly shaped according to the number of categories in each variable, which is essential for accurate computation of the checkerboard copula scores and associated.

This flexible approach to data loading ensures compatibility with various data collection methodologies and storage formats, allowing researchers to focus on analysis rather than data conversion. The `DataProcessor` integrates seamlessly with `GenericCCRVAM` class to initialize model objects directly from imported data, creating a streamlined workflow from raw data to statistical inference. Detailed examples demonstrating each data format, along with corresponding code snippets and implementation considerations, are provided in Chapter 5.

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

Following the theoretical foundation in Definition 3.1, the package implements checkerboard copula scores for ordinal variables. These scores represent a transformation that leverages

the inherent ordering information in ordinal variables.

The implementation maintains fidelity to the mathematical definitions while providing a computationally efficient vectorized implementation:

```
1 # Calculate and display CCS for X1
2 scores_X1 = ccrvam_obj.calculate_ccs(1)
3 print(scores_X1)
```

```
[0.125, 0.3125, 0.5, 0.6875, 0.875]
```

```
1 # Calculate and display Variance of CCS for X1
2 variance_ccs_X1 = ccrvam_obj.calculate_variance_ccs(1)
3 print(variance_ccs_X1)
```

```
0.0791015625
```

Recall that for each ordinal variable  $X_j$  with categories  $i_j \in \{1, \dots, I_j\}$ , the scores  $s_{i_j}^j = (u_{i_j-1}^j + u_{i_j}^j)/2$  are calculated where  $u_{i_j}^j$  is defined by the marginal cumulative distribution. This implementation follows directly from the empirical estimation procedure detailed in Section 3.3.2. We can see how the output from the above code-chunk matches with the result in our running example from Chapter 3, which further verifies the reproducible functionality of our package.

ccrvam employs vectorized operations through NumPy (Developers 2025) to ensure computational efficiency, which becomes particularly important for higher-dimensional tables. The variance calculation implements Lemma 3.1, providing a measure of dispersion that is essential for the scaled association measures discussed later.

## 4.4 Checkerboard copula Regression (CCR)

The Checkerboard Copula Regression functionality follows the definition provided in Definition 3.4, computing the conditional expectation of the copula score for the response variable given values of the predictor variables.

The prediction functionality follows the empirical estimation procedure described in #sec-empirical-ccr, where the predicted category  $\hat{x}_{i_j}^j$  is determined by finding the interval containing the estimated regression value  $\hat{u}_j^* = \hat{r}_{U_j|U_{-j}}(\hat{\mathbf{u}}_{-j}^*)$ .

The `get_category_predictions_ccr()` method performs the essential function of predicting the categories of the response variable (specified through the response input argument) based on given predictor values (enumerated in the predictors input argument). This method implements the core predictive capability of the checkerboard copula regression approach, translating theoretical associations into practical category predictions. The method returns these predictions in an easy-to-read Pandas (McKinney 2011) DataFrame format, making it straightforward for researchers to examine and interpret the results in a familiar tabular structure. Additionally, the method supports custom variable names for enhanced interpretation, allowing users to replace default numeric identifiers with meaningful labels that reflect the actual variables being analyzed in their specific domain context.

The implementation also allows for multiple conditioning axes, supporting complex multivariate analyses, which we can be seen in the examples mentioned in Section 4.12.2.

```
1 # Predictions from X1 to X2:
2 predictions_X1_to_X2 = ccrvam_obj.get_predictions_ccr(
3     predictors=[1],
```



```

4     response=2
5 )
6 print(predictions_X1_to_X2)

```

	X1 Category	Predicted Response Category
0	1	3
1	2	2
2	3	1
3	4	2
4	5	3

```

1 # Example: Showcasing the use of custom variable names for the output
2 # Predictions from Education Level to Income Bracket:
3 variable_to_name_dict = {
4     1: "Income",
5     2: "Education"
6 }
7 predictions_Education_to_Income = ccrvam_obj.get_predictions_ccr(
8     predictors=[2],
9     response=1,
10    variable_names=variable_to_name_dict
11 )
12 print(predictions_Education_to_Income)

```

	Education Category	Predicted Income Category
0	1	3

1	2	3
2	3	3

The package also provides reference prediction under joint independence, which is important for interpreting the substantive meaning of predictions by comparing against what would be expected if no association existed.

Hence, we can also obtain the response category prediction under the assumption of joint independence between X1 and X2 as follows:

```
1 # Response category prediction under the joint independence between X1 and X2
2 print(ccrvam_obj.get_prediction_under_indep(2))
```

2

## 4.5 CCR Predicted Category Visualization

The `ccrvam` package includes a comprehensive set of visualization tools for exploring dependence structures in multivariate ordinal data. The package provides a built-in visualization method for CCR predictions:

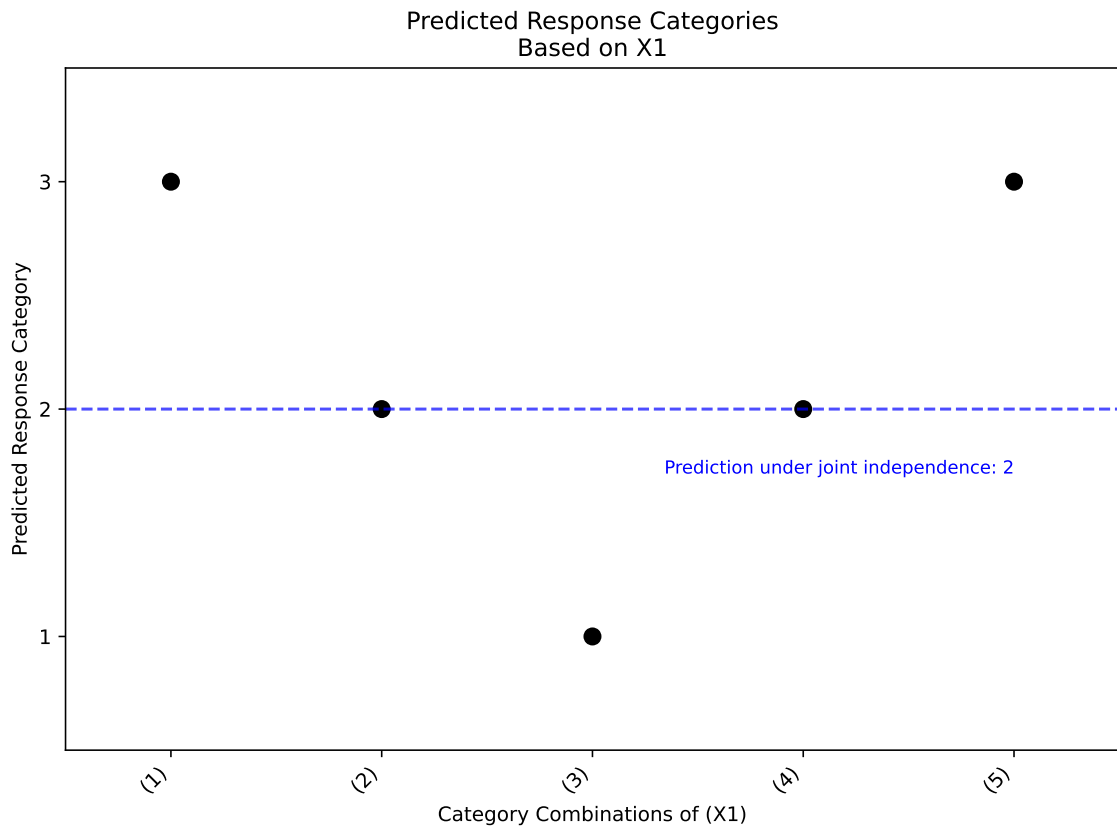
This visualization approach creates a heatmap-style plot showing the predicted categories of the response variable for different combinations of predictor variable categories. It includes markers for predicted categories and optional reference lines for predictions under joint independence.

The visualization methods support various color schemes for different visual preferences, customizable figure sizes and resolutions, text annotations showing prediction values, different legend styles for handling many predictor combinations, and exportable high-resolution

graphics for publications.

These visualizations help researchers understand and communicate the complex dependence structures detected by the CCR approach, making the results more accessible and interpretable.

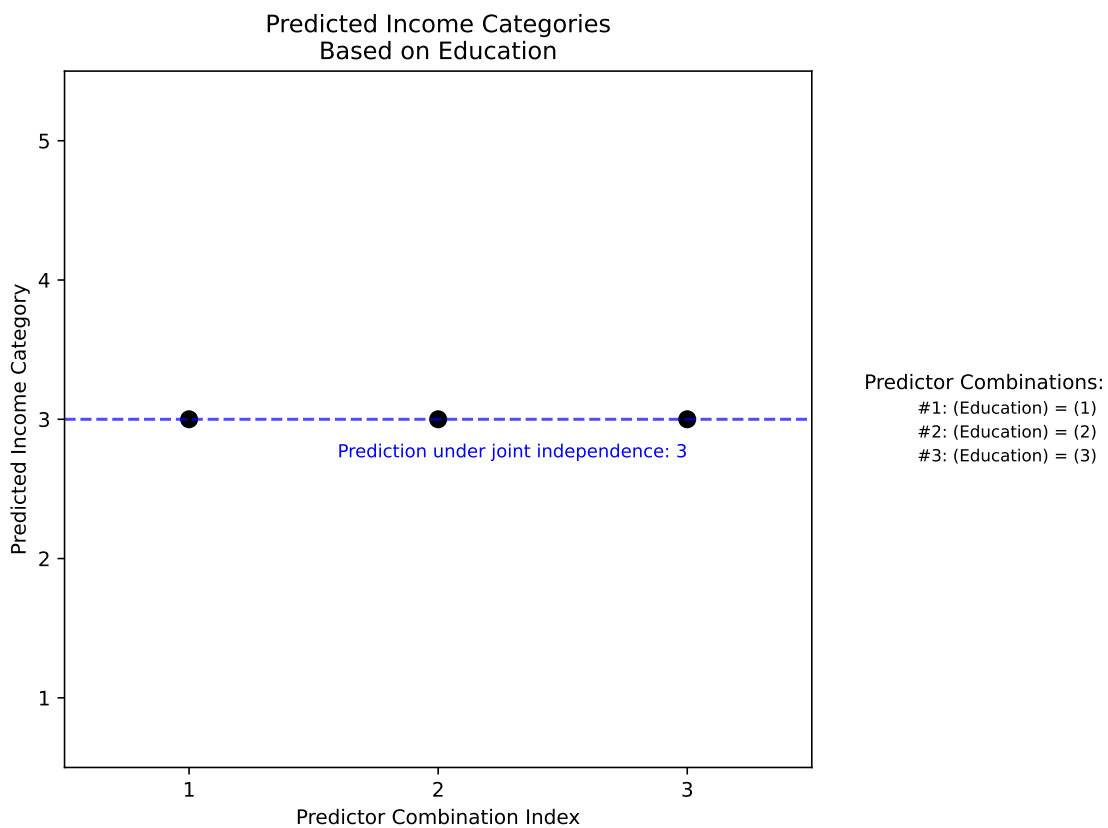
```
1 # Plotting with default naming scheme with tuple labels on x-axis
2 ccrvam_obj.plot_ccr_predictions(
3     predictors=[1],
4     response=2,
5     legend_style="xaxis"
6 )
```



```

1 # Plotting with custom naming scheme with legend of category combinations
2 var_names={1:"Income", 2:"Education"}
3
4 ccrvam_obj.plot_ccr_predictions(
5     predictors=[2],
6     response=1,
7     legend_style="side",
8     variable_names=var_names
9 )

```



## 4.6 CCR Prediction Uncertainty Evaluation Using Nonparametric Bootstrap Resampling

To quantify prediction uncertainty, the package implements nonparametric bootstrap methods:

```
1  from ccrrvam import bootstrap_predict_ccr_summary
2
3  prediction_matrix = bootstrap_predict_ccr_summary(
4      contingency_table_2d,
5      predictors=[1],
6      predictors_names=["X"],
7      response=2,
8      response_name="Y",
9      n_resamples=9999
10 )
11
12 # Predictions Summary Matrix
13 print(prediction_matrix)
```

	Y=1	Y=2	Y=3
X=1	0.0	0.0	100.0
X=2	0.0	100.0	0.0
X=3	100.0	0.0	0.0
X=4	0.0	100.0	0.0
X=5	0.0	0.0	100.0

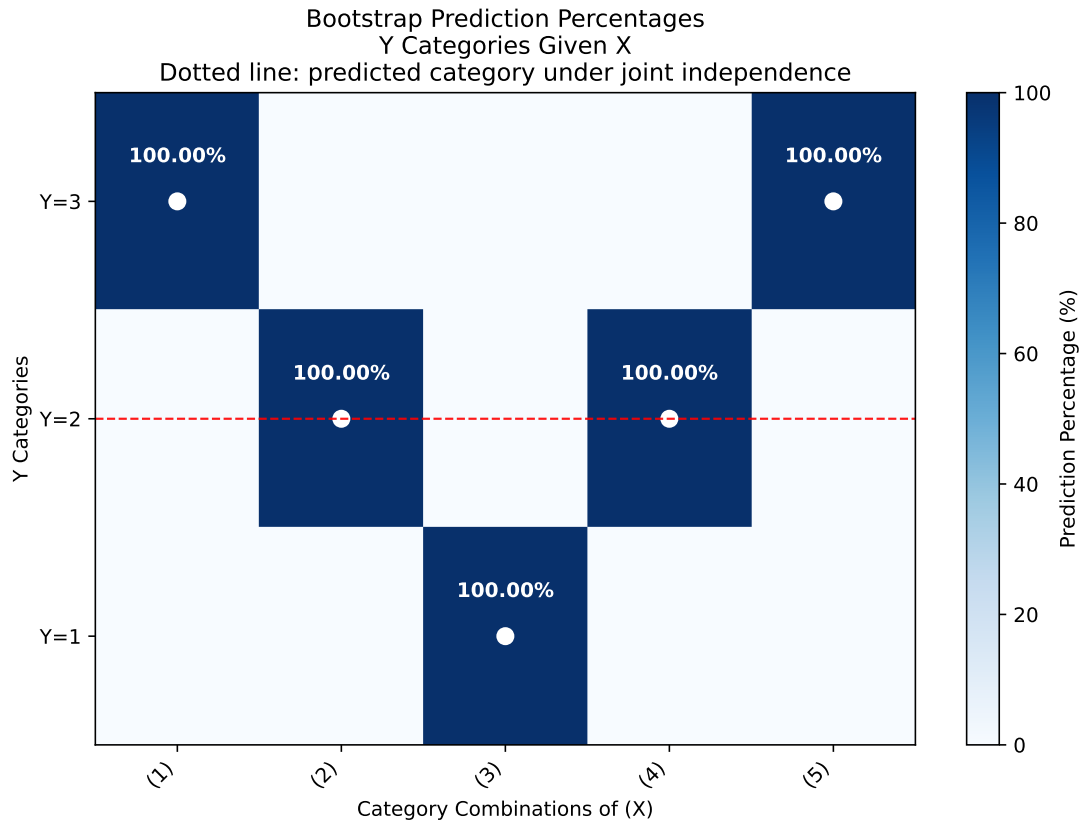
```
1 print(prediction_matrix.predictions)
```

	Predicted
X=1	3
X=2	2
X=3	1
X=4	2
X=5	3

This implementation follows the nonparametric bootstrap procedure outlined in the previous chapter, where multiple bootstrap samples are generated from the original contingency table, and predictions are made for each resampled dataset. The distribution of predicted categories provides a measure of prediction uncertainty, represented as percentages in the resulting heatmap visualization.

The visualization component employs a color gradient to represent the confidence in predictions, with darker colors indicating higher prediction percentages (greater confidence). Dotted lines indicate predictions under joint independence, providing a reference point for interpretation. More input arguments and options for customization can be explored further in <https://ccrvam.readthedocs.io/>, which hosts detailed documentation for our ccrvam package.

```
1 # You can also visualize the results with the attached plotting method
2 prediction_matrix.plot_prediction_heatmap()
```



## 4.7 (S)CCRAM Estimation

The package implements both the unscaled (CCRAM) and scaled (SCCRAM) versions of the checkerboard copula regression association measure, as defined in Definition 3.5 and Definition 3.6:

```

1 ccram_X1_to_X2 = ccrvam_obj.calculate_CCRAM(
2     predictors=[1],
3     response=2
4 )
5 print(f"CCRAM X1 to X2: {ccram_X1_to_X2:.4f}")

```

CCRAM X1 to X2: 0.8438

```
1 sccram_X1_to_X2 = ccrvam_obj.calculate_CCRAM(  
2     predictors=[1],  
3     response=2,  
4     scaled=True  
5 )  
6 print(f"SCCRAM X1 to X2: {sccram_X1_to_X2:.4f}")
```

SCCRAM X1 to X2: 1.0000

The implementation follows the empirical estimation procedures outlined in the previous chapter, with CCRAM measuring the proportion of variance in the response variable's checkerboard copula score that can be explained by the predictor variables. SCCRAM normalizes this measure to be bounded between 0 and 1, making it easier to interpret and compare across different datasets.

Both measures quantify the strength of even the nonlinear regression relationship between the ordinal response variable and categorical predictors, going beyond traditional correlation measures that primarily detect linear relationships.

## 4.8 (S)CCRAM Uncertainty Evaluation Using Bootstrap Resampling

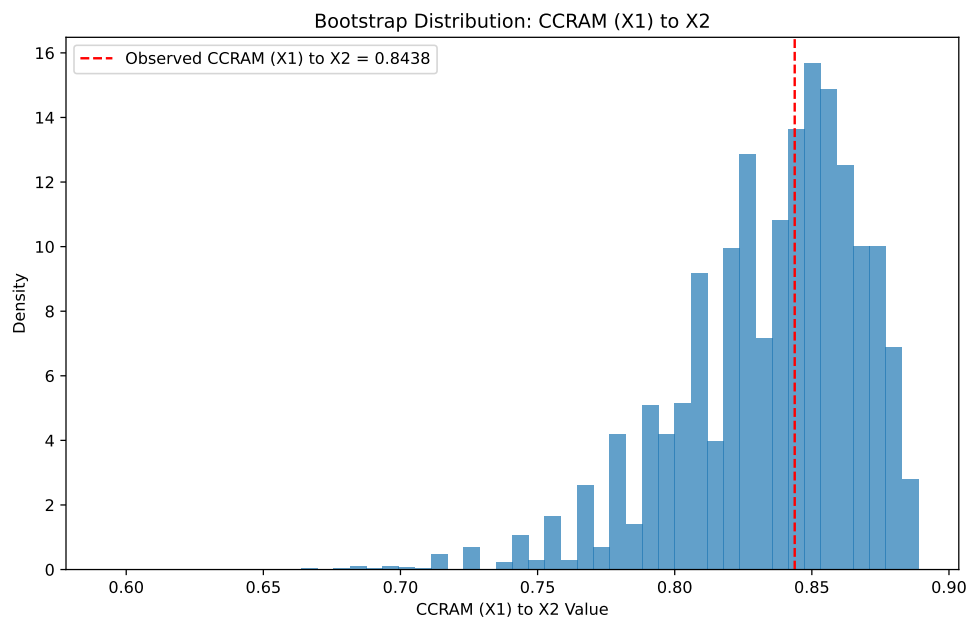
As mentioned in Section 3.5.3, in order to assess the precision of CCRAM and SCCRAM estimates, the package implements nonparametric bootstrap procedures:



```

1 from ccrvam import bootstrap_ccram
2
3 ccram_result = bootstrap_ccram(
4     contingency_table_2d,
5     predictors=[1],
6     response=2,
7     n_resamples=9999,
8     scaled=False,
9     confidence_level=0.95,
10    method="percentile",
11    random_state=None
12 )

```



```
1 # Metric Name
2 print(ccram_result.metric_name)
```

CCRAM (X1) to X2

```
1 # Observed Value
2 print(f"{ccram_result.estimated_value:.4f}")
```

0.8438

```
1 # 95% Confidence Interval
2 lower_CI_bound = ccram_result.confidence_interval[0]
3 upper_CI_bound = ccram_result.confidence_interval[1]
4 print(f"({lower_CI_bound:.4f}, {upper_CI_bound:.4f})")
```

(0.7553, 0.8815)

```
1 # Standard Error
2 print(f"{ccram_result.standard_error:.4f}")
```

0.0327

```
1 # Bootstrap Estimates
2 bootstrap_estimates = ccram_result.bootstrap_distribution
3 print(f"{type(bootstrap_estimates)}")
```

<class 'numpy.ndarray'>

```

1 # Calculate bootstrap bias
2 bootstrap_mean = np.mean(bootstrap_estimates)
3 bootstrap_bias = bootstrap_mean - ccram_result.observed_value
4
5 # Calculate bootstrap standard error
6 bootstrap_std_error = np.std(bootstrap_estimates, ddof=1)
7
8 # Calculate ratio of bias to standard error
9 bias_to_se_ratio = bootstrap_bias / bootstrap_std_error
10
11 # Additional Bootstrap Statistics
12 print(f"Bootstrap Mean: {bootstrap_mean:.4f}")

```

Bootstrap Mean: 0.8353

```

1 print(f"Bootstrap Bias: {bootstrap_bias:.4f}")

```

Bootstrap Bias: -0.0085

```

1 print(f"Bootstrap Standard Error: {bootstrap_std_error:.4f}")

```

Bootstrap Standard Error: 0.0327

```

1 print(f"Bias to Standard Error Ratio: {bias_to_se_ratio:.4f}")

```

Bias to Standard Error Ratio: -0.2587

The bootstrap procedure generates multiple resamples from the original contingency table, calculates the (S)CCRAM for each resample, and constructs confidence intervals based on

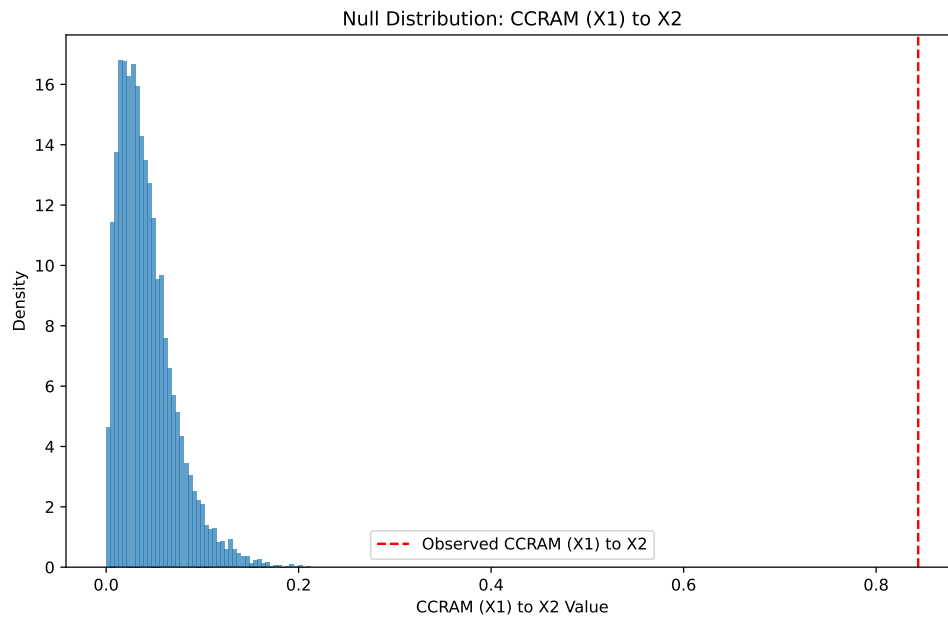
the resulting distribution. This provides a measure of the sampling variability and precision of the (S)CCRAM estimate.

The visualization component plots the bootstrap distribution with the observed (S)CCRAM value highlighted, providing a graphical representation of the uncertainty in the estimate. This analysis can be repeated for SCCRAM by setting `scaled = True` as an input argument of the `bootstrap_ccram()` function.

## 4.9 Statistical Significance Testing for (S)CCRAM Using Permutation Test

As mentioned in Section 3.5.4, in order to assess the statistical significance of CCRAM or SCCRAM, the package implements permutation testing as well. We demonstrate the usage of the same in the case of our 2-D example below:

```
1  from ccram import permutation_test_ccram
2
3  perm_result = permutation_test_ccram(
4      contingency_table_2d,
5      predictors=[1],
6      response=2,
7      scaled=False,
8      alternative='greater',
9      n_resamples=9999
10 )
```



```
1 print(f"Metric Name: {perm_result.metric_name}")
```

Metric Name: CCRAM (X1) to X2

```
1 print(f"Observed Value: {perm_result.observed_value:.4f}")
```

Observed Value: 0.8438

```
1 print(f"P-Value: {perm_result.p_value:.4f}")
```

P-Value: 0.0001

```
1 # Permutation Distribution
2 permutation_distribution = perm_result.null_distribution
3 print(f"Permutation Distribution (Type): {type(permutation_distribution)}")
```

Permutation Distribution (Type): <class 'numpy.ndarray'>

```

1 # Calculate quantiles
2 q01 = np.quantile(permutation_distribution, 0.01)
3 # 0.5-th quantile (median)
4 median = np.median(permutation_distribution)
5 # 0.99-th quantile
6 q99 = np.quantile(permutation_distribution, 0.99)
7
8 # Calculate interquartile range (IQR)
9 q25 = np.quantile(permutation_distribution, 0.25)
10 q75 = np.quantile(permutation_distribution, 0.75)
11 iqr = q75 - q25
12
13 # Permutation Distribution Summary Statistics:
14 print(f"0.01-th Quantile: {q01:.4f}")

```

0.01-th Quantile: 0.0032

```

1 print(f"0.5-th Quantile (Median): {median:.4f}")

```

0.5-th Quantile (Median): 0.0360

```

1 print(f"0.99-th Quantile: {q99:.4f}")

```

0.99-th Quantile: 0.1363

```

1 print(f"Interquartile Range (IQR): {iqr:.4f}")

```

Interquartile Range (IQR): 0.0366

This implementation follows the permutation testing procedure outlined in the Chapter 3, where the response variable values are randomly permuted to break any association with the predictor variables, thus generating a null distribution under the hypothesis of no association. The p-value is calculated as the proportion of permutation statistics that are as extreme as or more extreme than the observed statistic.

The visualization component plots the null distribution with the observed CCRAM value highlighted, providing a graphical representation of the statistical significance. This analysis can be repeated for SCCRAM by setting `scaled = True` as an input argument of the `permutation_test_ccram()` function.

## 4.10 Software Architecture and Design Principles

The `ccrvam` Python package was developed following modern software engineering principles to ensure reliability, maintainability, and extensibility. We used `PyPi-Template` (VG 2024) to initialize the skeleton of our software package.

### 4.10.1 Component Structure and Object-Oriented Design

The package follows an object-oriented design, encapsulating related functionality within classes. This design allows users to work with a unified interface while hiding the pesky implementation details from the user, therefore making the package intuitive to use while maintaining flexibility.

The package is organized into three main components:

1. **Core CCRVAM Implementation** (`GenericCCRVAM` class within `gencopula` module):

This is the central object that implements the fundamental calculations for checkerboard

copula regression, while handling internal data representation and transformation. Through this, we provide our users several methods for prediction and association measures.

2. **Statistical Simulation Framework** (`genstatsim` module): This module implements bootstrap and permutation testing procedures, while providing uncertainty quantification for predictions and measures. Through this, we also provide users the flexibility through visualization and exporting methods for statistical results
3. **Data Processing Utilities** (`utils` module): This module handles user-facing data loading and formatting methods, providing conversion between different data representations such as table form, case form, and frequency form as outlined in Chapter 5. Through this, we also provide users basic functionality for data validation and pre-processing.

We can visualize the package structure and user-experience-workflow through the images below (powered by Mermaid (Sveidqvist and Mermaid 2014)):

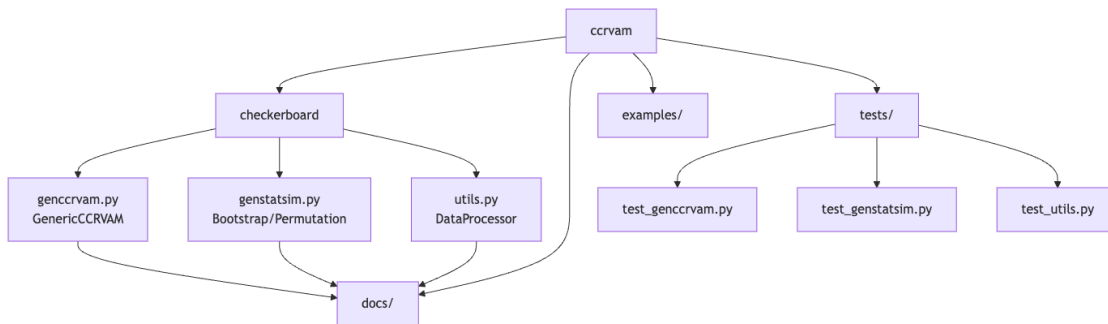


Figure 4.1: High-Level Package Structure of CCRVAM

Figure 4.1 illustrates the modular organization of the `ccrvam` codebase. The core functionality is encapsulated within the `checkerboard` subpackage, which houses the main analytical



engine (`genccrvam.py`), statistical simulation tools (`genstatsim.py`), and data preprocessing utilities (`utils.py`). The top-level package initialization file (`__init__.py`) exposes these modules for external use, while supplementary materials such as documentation, examples, and tests are organized into their respective directories.

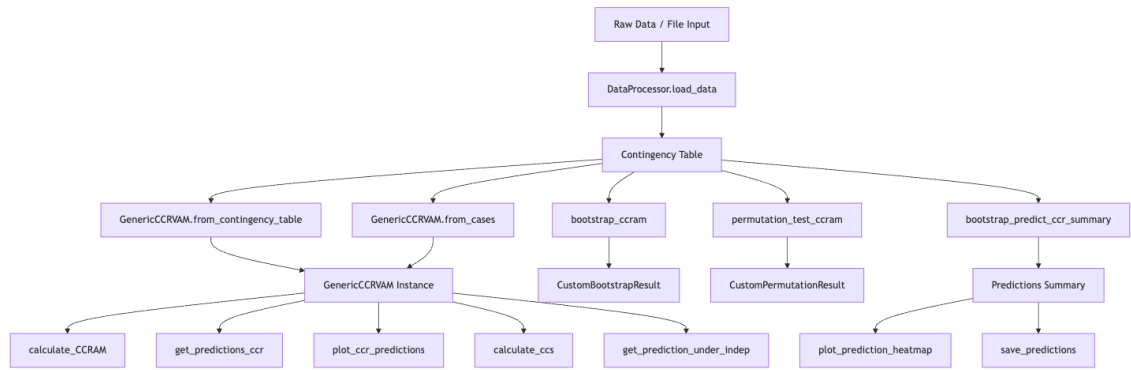


Figure 4.2: Functional Workflow for Checkerboard Copula Analysis

Figure 4.2 outlines the end-to-end computational pipeline within the `ccrvam` package. Raw categorical data is processed into a contingency table using the `DataProcessor`. This table can then be passed into the `GenericCCRVAM` class for association analysis, prediction, and visualization. Alternatively, the same input can be used in bootstrapping and permutation testing workflows to produce confidence intervals, p-values, and prediction heatmaps. The system supports both analytical modeling and robust statistical inference.

#### 4.10.2 Vectorized Implementations, and Error Handling

Performance optimization was a key consideration in the design, particularly for higher-dimensional tables. By leveraging NumPy’s (Developers 2025) vectorized operations, Pandas’s (Mckinney 2011) effective data-handling, SciPy’s (Virtanen et al. 2020) bootstrapping function-calls, and Matplotlib’s (Hunter 2007) efficient graphing APIs (Application Programming Interfaces) the package achieves significantly better performance than naive loop-based

implementations, enabling analysis of larger datasets.

The package includes comprehensive input validation and error handling to provide informative messages when issues arise. This approach helps users identify and fix problems quickly, improving the overall user experience. On the developer-side, this allows for easy debugging, and faster development of new features.

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

### **4.11.1 Comprehensive Test Suite**

The `ccrvam` package includes a comprehensive test suite to ensure correctness and reliability across all implemented functionality. The test suite encompasses unit tests for all core functionality, providing verification of individual components in isolation, while integration tests confirm proper behavior in end-to-end workflows that simulate typical user interactions. Edge case testing rigorously examines boundary conditions where algorithms are most likely to fail, and dimensional-invariant testing validates consistent performance across 2D, 3D, and 4D contingency tables of varying complexity. The suite also incorporates regression tests to prevent the reintroduction of previously fixed bugs as the codebase evolves.

The package achieves over 93% code coverage, ensuring that most code paths and user experiences are well-tested and safe for production use in statistical analysis. In order to achieve better observability into our code and maintain a check ensuring that our tests pass irrespective of machine environments, we leveraged `pytest` (Krekel 2025) and `coverage` (Batchelder and Coverage.py 2025) Python libraries as our testing infrastructure. These tools provide a robust framework for automated test execution and detailed reporting on test coverage. An example test from our suite is shown below:

```

1  @pytest.mark.parametrize(
2      "predictors, response, expected_sccram", [
3          ([1], 2, 0.84375/(12*0.0703125)), # Single axis X1->X2
4          ([2], 1, 0.0), # Single axis X2->X1
5      ]
6  )
7
8  def test_calculate_SCCRAM(
9      generic_ccrvam, predictors,
10     response, expected_sccram
11 ):
12     """Test SCCRAM calculations with multiple conditioning axes."""
13     calculated = generic_ccrvam.calculate_CCRAM(
14         predictors, response, scaled=True
15     )
16     np.testing.assert_almost_equal(calculated, expected_sccram)

```

Additional tests covering various aspects of the package's functionality can be found in the GitHub repository at <https://github.com/DhyeyMavani2003/ccrvam/tree/main/tests>.

#### 4.11.2 Continuous Integration (CI)

The development workflow includes CI testing through GitHub Actions and `.yaml` files. This configuration ensures that tests are run automatically on multiple Python versions (3.8, 3.9, 3.10, 3.11, 3.12, 3.13) for every code change, maintaining compatibility and reliability.

### 4.11.3 Performance Benchmarking

The package includes three main performance optimizations for handling larger contingency tables: vectorized implementations for core calculations, caching of intermediate results such as conditional distributions to avoid redundant computation, efficient data structures for sparse representation where appropriate.

For a 4D contingency table (2x3x2x6) with 112 cases, operations like CCRAM calculation and bootstrap simulations complete in seconds on modern hardware, allowing for interactive analysis.

## 4.12 User Documentation and Example Workflows

The `ccrvam` package includes comprehensive documentation and example workflows to help users get started:

### 4.12.1 API Documentation

The package provides detailed API documentation for all user and developer facing functions and classes through Sphinx (Turner 2025). The documentation is hosted on ReadTheDocs at <https://ccrvam.readthedocs.io/>, and include function signatures, input arguments descriptions, outputs documentation, warnings/errors log, usage examples on 2D and 4D sample datasets, and cross-references to related functions.

### 4.12.2 Example Workflows

The package includes example workflows to demonstrate common analysis patterns. These examples (located at <https://github.com/DhyeyMavani2003/ccrvam/tree/main/examples>

[/jupyter](#)) demonstrate complete analysis workflows from data loading to visualization and statistical testing, helping users understand how to apply the package to their own research questions.

In the next chapter, we will use our `ccrvam` package to perform EDA on some real-world datasets.

## Chapter 5

# Real Data Analysis

In this chapter, we demonstrate the practical application of the Checkerboard Copula Regression-based Visualization and Association Measure (CCRVAM) techniques introduced in previous chapters. We will analyze a real-world dataset using the `ccrvam` package implementation described in Chapter 4. This analysis will showcase how our methods can be used to explore associations between categorical predictors and an ordinal response variable, quantify the strength of these associations, and visualize prediction patterns.

### 5.1 Dataset Overview

The dataset we analyze contains information from a clinical study on back pain treatments, originally presented by Anderson (1984). This dataset is particularly suitable for our methodology as it includes an ordinal response variable (pain relief outcome) and multiple categorical predictor variables.

The dataset consists of 4 categorical variables:

Variable	Description	Categories
$X_1$	Length of Previous Attack	1=Short, 2=Long
$X_2$	Pain Change	1=Better, 2=Same, 3=Worse
$X_3$	Lordosis	1=Absent/Decreasing, 2=Present/Increasing
<i>Pain</i>	Back Pain Outcome	1=worse (W), 2=same (S), 3=slight improvement (SI), 4=moderate improvement (MODI), 5=marked improvement (MARI), 6=complete relief (CR)

This dataset represents a common scenario in medical and social science research, where the goal is to understand how multiple categorical factors influence an ordinal outcome. The pain outcome variable has a natural ordering (from worse to complete relief), making it an ideal candidate for our checkerboard copula approach.

## 5.2 Data Preparation and Loading

The `ccrvam` package provides flexible data loading capabilities through `DataProcessor` class. As described in Chapter 4, this class supports multiple data formats including case-form, frequency-form, and table-form data. Here, we demonstrate how to load the back pain dataset using each of these approaches.

First, we need to import the necessary libraries and define our variable structure:

```

1 import numpy as np
2 from ccrvam import GenericCCRVAM, DataProcessor
3
4 # Define the ordered list of variable names
5 var_list_4d = ["x1", "x2", "x3", "pain"]
6
7 # Define the dimension tuple representing
8 # the number of categories for each variable
9 data_dimension = (2, 3, 2, 6)
10
11 # Create a category mapping for non-integer categories
12 # (required for 'pain' variable)
13 category_map_4d = {
14     "pain": {
15         "worse": 1,
16         "same": 2,
17         "slight.improvement": 3,
18         "moderate.improvement": 4,
19         "marked.improvement": 5,
20         "complete.relief": 6
21     },
22 }

```

The `var_list_4d` defines the order of variables in our analysis. The `data_dimension` tuple



specifies the number of categories for each variable in the same order. The `category_map_4d` provides a mapping from text labels to numeric indices for non-integer categories, which is necessary for the “pain” variable in this dataset.

The CCRVAM package supports three different data loading formats, providing flexibility based on how your data is structured. We’ll demonstrate each method:

### 5.2.1 Case Form Data Loading

Case form represents individual observations, where each row contains the category values for all variables for a single observation:

```
1 # Loading data from case form file
2 contingency_table_4d = DataProcessor.load_data(
3     "./data/caseform.pain.txt",
4     data_form="case_form",
5     dimension=data_dimension,
6     var_list=var_list_4d,
7     category_map=category_map_4d,
8     named=True,
9     delimiter="\t"
10 )
11
12 # Initialize the GenericCCRVAM object
13 rda_ccrvam = GenericCCRVAM.from_contingency_table(
14     contingency_table_4d
```

```
15 )
```

### 5.2.2 Frequency Form Data Loading

Frequency form data contains the category values for all variables along with a count of how many times that combination appears:

```
1  # Loading data from frequency form file
2  contingency_table_4d_from_freq = DataProcessor.load_data(
3      "./data/freqform.pain.txt",
4      data_form="frequency_form",
5      dimension=data_dimension,
6      var_list=var_list_4d,
7      category_map=category_map_4d,
8      named=True,
9      delimiter="\t"
10 )
11
12 # Initialize the GenericCCRVAM object
13 rda_ccrvam_from_freq = GenericCCRVAM.from_contingency_table(
14     contingency_table_4d_from_freq
15 )
```

### 5.2.3 Contingency Table Form Data Loading

Table form represents the data as a multidimensional contingency table with counts directly:

```

1  # Define the 4D contingency table as a NumPy array
2  rda_contingency_table = np.array([
3      # X1=1 (Short)
4      [
5          # X2=1 (Better)
6          [
7              # X3=1 (Absent)
8              [0, 1, 0, 0, 2, 4], # Counts for each Pain outcome
9              # X3=2 (Present)
10             [0, 0, 0, 1, 3, 0]
11         ],
12         # X2=2 (Same)
13         [
14             # X3=1 (Absent)
15             [0, 2, 3, 0, 6, 4],
16             # X3=2 (Present)
17             [0, 1, 0, 2, 0, 1]
18         ],
19         # X2=3 (Worse)
20         [
21             # X3=1 (Absent)
22             [0, 0, 0, 0, 2, 2],
23             # X3=2 (Present)

```

```

24         [0, 0, 1, 1, 3, 0]
25     ]
26 ],
27 # X1=2 (Long)
28 [
29     # X2=1 (Better)
30     [
31         # X3=1 (Absent)
32         [0, 0, 3, 0, 1, 2],
33         # X3=2 (Present)
34         [0, 1, 0, 0, 3, 0]
35     ],
36     # X2=2 (Same)
37     [
38         # X3=1 (Absent)
39         [0, 3, 4, 5, 6, 2],
40         # X3=2 (Present)
41         [1, 4, 4, 3, 0, 1]
42     ],
43     # X2=3 (Worse)
44     [
45         # X3=1 (Absent)
46         [2, 2, 1, 5, 2, 0],

```

```

47         # X3=2 (Present)
48         [2, 0, 2, 3, 0, 0]
49     ]
50 ]
51 ])
52
53 # Load data from the table
54 contingency_table_4d_from_array = DataProcessor.load_data(
55     rda_contingency_table,
56     data_form="table_form",
57     dimension=data_dimension,
58     var_list=var_list_4d,
59     category_map=category_map_4d
60 )
61
62 # Initialize the GenericCCRVAM object
63 rda_ccrvam_from_array = GenericCCRVAM.from_contingency_table(
64     contingency_table_4d_from_array
65 )

```

```

1 # Check if the Resulting Joint Probability Matrices are the same
2 # after loading data using various different methods mentioned above
3 same_1_2 = np.array_equal(rda_ccrvam.P, rda_ccrvam_from_freq.P)
4 same_2_3 = np.array_equal(rda_ccrvam_from_freq.P, rda_ccrvam_from_array.P)

```

```
5
6 # Are P matrices are the same across methods?
7 print(same_1_2 and same_2_3)
```

True

The output of each loading method is a 4-dimensional joint probability matrix with shape (2, 3, 2, 6) corresponding to the number of categories for each variable. This matrix contains the estimated joint probability distribution for all possible combinations of the categorical variables.

The values in the matrix represent the probability of observing each specific combination of categories. For example, the value at position  $[0, 0, 0, 1] = 0.00990099$  represents the probability of observing:  $X_1 = 1$  (Short previous attack),  $X_2 = 1$  (Better pain change),  $X_3 = 1$  (Absent/Decreasing Lordosis), and  $Pain = 2$  (Same pain outcome).

Note that all three loading methods should produce the same joint probability matrix if the data sources are consistent, which we can observe from the identical outputs in the example.

**i** Note:

In the interest of brevity, for further real data analysis covered in this chapter, we will not be walking through each code-chunk. If you are interested, please feel free to checkout the code for the real data analysis in the Jupyter notebooks at <https://github.com/DhyeyMavani2003/ccrvam/tree/main/examples/jupyter>.

### 5.3 Exploratory Data Analysis

Before applying our advanced statistical methods, we first examine the basic probability distributions in the data. These distributions provide insights into the prevalence of each category in our dataset. We observe the following marginal probability density functions (pdfs):

- **Length of Previous Attack ( $X_1$ ):** 38.61% of patients had short previous attacks, while 61.39% had long previous attacks.
- **Pain Change ( $X_2$ ):** 20.79% of patients experienced better pain change, 51.49% had the same level of pain, and 27.72% experienced worse pain change.
- **Lordosis ( $X_3$ ):** 63.37% of patients had absent or decreasing lordosis, while 36.63% had present or increasing lordosis.
- **Back Pain Outcome ( $Pain$ ):** The distribution shows that 4.95% of patients experienced worse pain after treatment (W), 13.86% reported no change (S), 17.82% experienced slight improvement (SI), 19.80% had moderate improvement (MODI), 27.72% reported marked improvement (MARI), and 15.84% experienced complete relief (CR).

These findings provide valuable context for interpreting our subsequent analyses. For the back pain outcome specifically, we observe that treatments were generally effective, with more than 60% of patients experiencing at least moderate improvement (combining the moderate improvement, marked improvement, and complete relief categories).

### 5.4 Calculating Checkerboard Copula Scores (CCS)

Following the methodology described in Chapter 3, we calculate the checkerboard copula scores (CCS) for each variable in our dataset. We compute the CCS for all variables and

determine their respective variances:

- **Length of Previous Attack ( $X_1$ ):** Two distinct scores (0.193, 0.693) with a variance of 0.059.
- **Pain Change ( $X_2$ ):** Three distinct scores (0.104, 0.465, 0.861) with a variance of 0.069.
- **Lordosis ( $X_3$ ):** Two distinct scores (0.317, 0.817) with a variance of 0.058.
- **Back Pain Outcome ( $Pain$ ):** Six distinct scores ranging from 0.025 to 0.921, with a variance of 0.080.

The results show distinct patterns of scores across the variables, with Back Pain Outcome demonstrating the most granular distribution with six distinct scores ranging from approximately 0.025 to 0.921. The variance calculations reveal that the Back Pain Outcome variable has the highest variance at approximately 0.080, while Lordosis ( $X_3$ ) has the lowest at 0.058. These variance values are critical inputs for our subsequent analysis of scaled association measures, as they provide normalization factors that allow for meaningful comparisons across different variable relationships.

## 5.5 Checkerboard Copula Regression (CCR) Analysis

Next, we apply the Checkerboard Copula Regression (CCR) to predict the back pain outcome categories based on the predictor variables: Length of Previous Attack ( $X_1$ ), Pain Change ( $X_2$ ), and Lordosis ( $X_3$ ).

Our analysis generates predictions for each possible combination of predictor variables. The predictions show the expected pain outcome category for each combination of predictor variables. For example, patients with a short previous attack ( $X_1 = 1$ ), better pain change ( $X_2 = 1$ ), and absent lordosis ( $X_3 = 1$ ) are predicted to have marked improvement (category



5) in pain outcomes. The results reveal several important patterns:

1. Patients with a short previous attack ( $X_1 = 1$ ) generally have better outcomes (categories 4-5) than those with long previous attacks ( $X_1 = 2$ )
2. Within the short previous attack group, those with better or worse pain change ( $X_2 = 1$  or  $X_2 = 3$ ) tend to have marked improvement (category 5)
3. For patients with long previous attacks ( $X_1 = 2$ ), those with worse pain change ( $X_2 = 3$ ) generally have the poorest outcomes (category 3: slight improvement)
4. Under the assumption of joint independence between *Pain* and the predictor variables ( $X_1, X_2, X_3$ ), the predicted pain outcome category is 4 (moderate improvement), which serves as a reference point for our analysis.

These findings highlight the complex interrelationships between previous attack duration, pain change, and lordosis in predicting back pain treatment outcomes. The visualizations generated by `ccrvam` in Figure 5.1 further enhance our understanding of these relationships and provide valuable clinical insights.

## 5.6 Quantifying Association with (S)CCRAM

We now quantify the strength of the association between our predictors (Length of Previous Attack, Pain Change, and Lordosis) and the back pain outcome using CCRAM and SCCRAM. The CCRAM (Checkerboard Copula Regression Association Measure) value of 0.2576 indicates that approximately 25.76% of the variation in back pain outcomes can be explained by the three predictor variables. This provides a meaningful measurement of how well our predictor variables collectively explain the pain outcomes observed in patients.

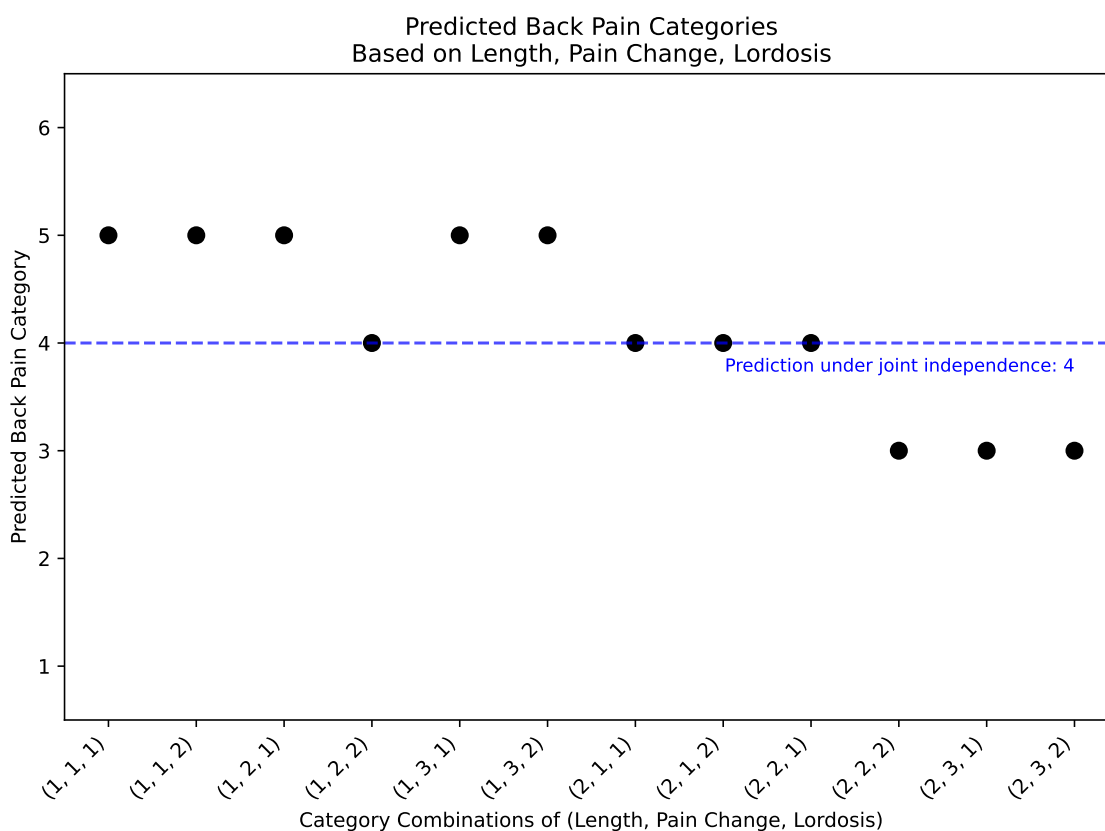


Figure 5.1: This visualization illustrates the predicted back pain outcomes based on combinations of prior attack length, pain change, and lordosis severity using the ccrvam framework. Each tuple on the x-axis represents a unique combination of predictor categories, with the predicted back pain category indicated by the position of the black dot. This plot reveals nuanced interdependencies among clinical factors and emphasizes the role of multivariate interactions in shaping patient outcomes.

For a more standardized interpretation, we calculate the SCCRAM (Scaled Checkerboard Copula Regression Association Measure), which yields a value of 0.2687. This scaled measure accounts for the theoretical maximum association possible in this dataset structure, making it easier to interpret and compare across different studies with varying data characteristics. These association measures provide important quantitative validation of the relationships we observed in our earlier analyses and help establish the overall predictive power of our model.

## 5.7 Uncertainty Quantification Using Bootstrap

To assess the uncertainty in our CCRAM and SCCRAM estimates, we utilize nonparametric bootstrap methods with 9,999 resamples. This approach allows us to estimate confidence intervals and standard errors without making distributional assumptions about our data.

For the CCRAM measure quantifying the association between our predictors (Length of Previous Attack, Pain Change, and Lordosis) and Back Pain outcomes, the bootstrap analysis yields an observed CCRAM value of 0.2576 with a 95% BCa confidence interval of (0.1849, 0.4762) and a standard error of 0.0748. For the scaled measure (SCCRAM), which normalizes the association for better interpretability, we observe a value of 0.2687 with a 95% BCa confidence interval of (0.0691, 0.3509) and a standard error of 0.0775.

These results reveal important insights about our analysis. The confidence intervals indicate that while there is uncertainty in the exact value of association, we can be reasonably confident that the true association is substantial. The positive bias in the bootstrap estimates (0.0666 for CCRAM and 0.0718 for SCCRAM) suggests that our observed values may be conservative estimates of the true association. The relatively high bias-to-standard-error ratios (0.8904 for CCRAM and 0.9267 for SCCRAM) indicate some potential complexity in the underlying

distribution, which further justifies our use of robust bootstrap methods for uncertainty quantification.

The visualizations generated by `ccrvam` in Figure 5.2 and Figure 5.3 further enhance our understanding of the uncertainty in (S)CCRAM by providing an intuitive representation.

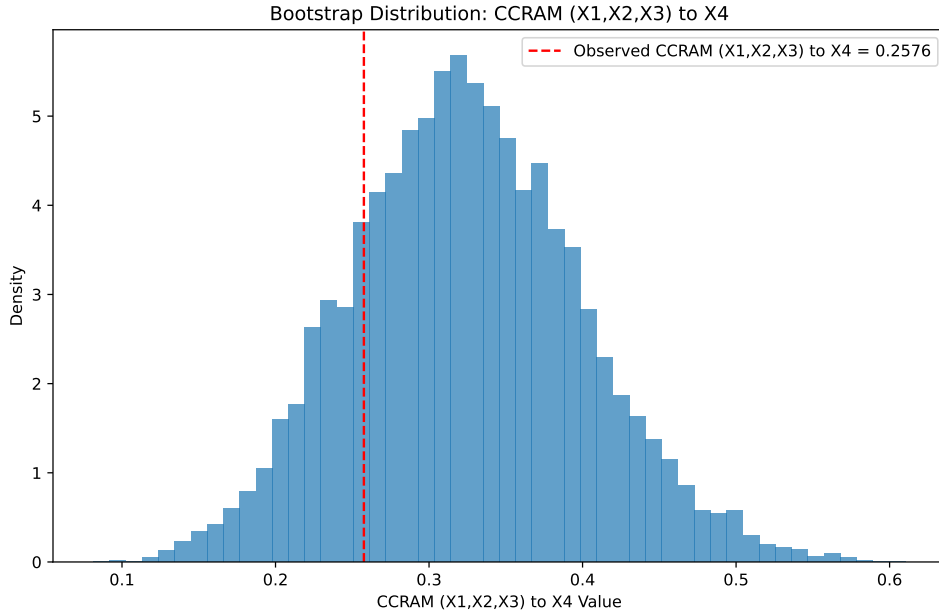


Figure 5.2: Bootstrap distribution of CCRAM ( $X_1, X_2, X_3 \rightarrow X_4$ ). The red dashed line marks the observed CCRAM value of 0.2576. This plot visualizes variability and supports estimation of confidence intervals and bias for the measure of association.

## 5.8 Statistical Significant Testing Using Permutation Tests

To assess whether the observed associations could have occurred by chance, we conduct permutation tests with 9,999 resamples. This approach allows us to construct empirical null distributions for both CCRAM and SCCRAM metrics under the hypothesis of no association between predictors and the pain outcome.

For the CCRAM measure, we observe a value of 0.2576 with a p-value of 0.0016. The per-

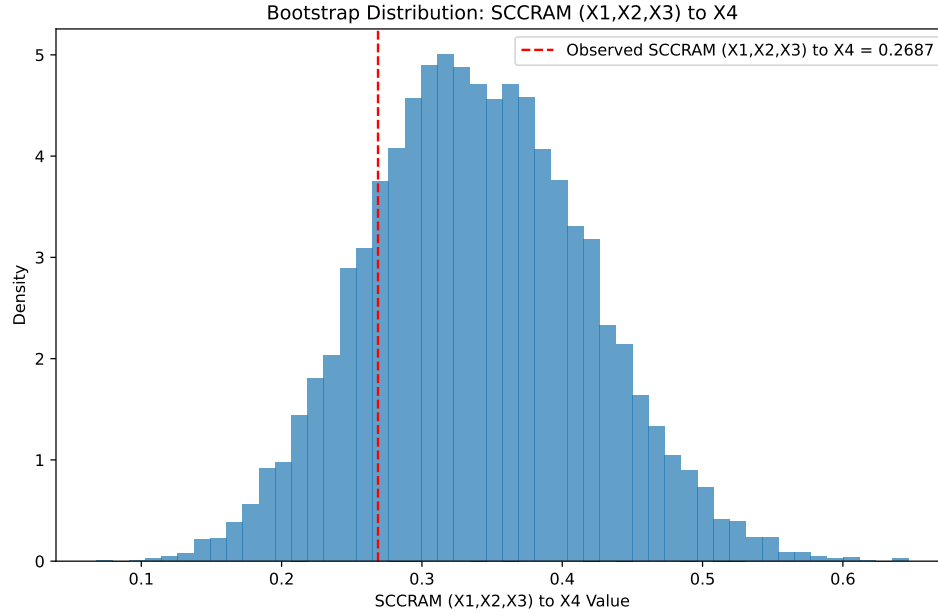


Figure 5.3: Bootstrap distribution of SCCRAM ( $X_1, X_2, X_3 \rightarrow X_4$ ). The red dashed line marks the observed SCCRAM value of 0.2687. This normalized version of CCRAM accounts for maximum possible association and highlights uncertainty in scaled estimates.

mutation distribution exhibits a median of 0.0998, with the 99th percentile at 0.2214. Our observed CCRAM value exceeds even the 99th percentile of the null distribution, providing strong evidence against the null hypothesis of no association. Similarly, for the SCCRAM measure, we observe a value of 0.2687 with an even smaller p-value of 0.0011. The permutation distribution for SCCRAM shows a median of 0.1046 with the 99th percentile at 0.2255. Again, our observed value exceeds the 99th percentile of values that would be expected by chance.

These permutation test results provide strong statistical evidence that the observed associations between our predictor variables (Length of Previous Attack, Pain Change, and Lordosis) and Back Pain outcomes are not due to random variation. The extremely small p-values confirm that these relationships are statistically significant, further validating the clinical relevance of our findings.

The visualizations generated by `ccrvam` in Figure 5.4 and Figure 5.5 further enhance our understanding of the relative association strength in (S)CCRAM by providing an intuitive representation.

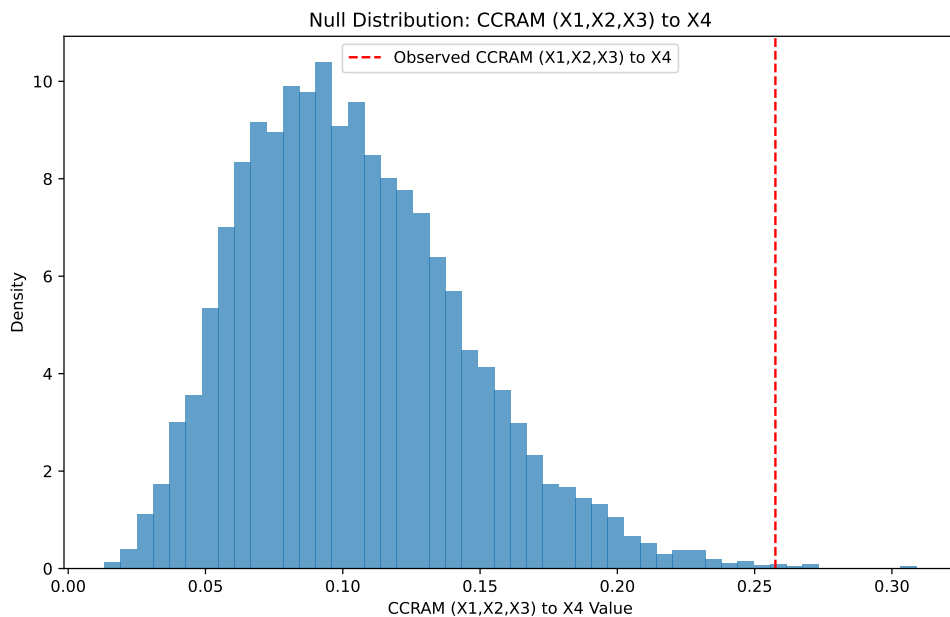


Figure 5.4: Permutation distribution of CCRAM ( $X_1, X_2, X_3 \rightarrow X_4$ ). The red dashed line marks the observed CCRAM value of 0.2576. The empirical null distribution illustrates that such a value is highly unlikely under the assumption of no association ( $p = 0.0016$ ).

## 5.9 Bootstrap Analysis for CCR Predictions

We can also use bootstrap methods to assess the uncertainty in our category predictions. By generating 9,999 bootstrap samples, we obtain a prediction matrix that shows the percentage of bootstrap samples predicting each pain category for each combination of predictor values.

This approach provides a measure of confidence in our predictions. For example, for patients with short previous attack, better pain change, and absent lordosis ( $X_1 = 1, X_2 = 1, X_3 = 1$ ), the prediction of category 5 (marked improvement) occurs in approximately 70.90% of

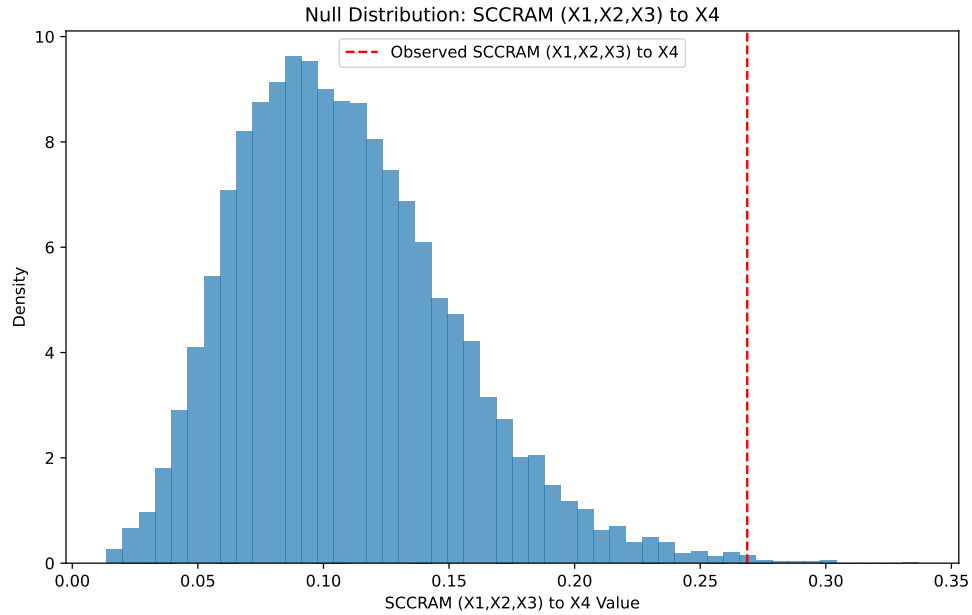


Figure 5.5: Permutation distribution of SCCRAM ( $X_1, X_2, X_3 \rightarrow X_4$ ). The red dashed line marks the observed SCCRAM value of 0.2687. The observed value lies well beyond the 99th percentile of the null distribution, providing strong evidence of statistically significant association ( $p = 0.0011$ ).

bootstrap samples, indicating high confidence in this prediction. Similarly, for patients with short previous attack, better pain change, and present lordosis ( $X_1 = 1, X_2 = 1, X_3 = 2$ ), the prediction of category 5 is even more consistent, occurring in 87.17% of bootstrap samples.

The bootstrap analysis also reveals cases where predictions are less certain. For instance, patients with long previous attack, same pain change, and present lordosis ( $X_1 = 2, X_2 = 2, X_3 = 2$ ) show 76.92% of bootstrap samples predicting category 3 (slight improvement), while patients with long previous attack, better pain change, and present lordosis ( $X_1 = 2, X_2 = 1, X_3 = 2$ ) have more uncertainty with 52.06% of samples predicting category 5 and 36.29% predicting category 4.

The visualization of this bootstrap prediction matrix through a heatmap produced by `ccrvam` in Figure 5.6 further enhances our understanding of prediction confidence across different

combinations of predictor variables. The dotted line in the heatmap indicates the predicted category under joint independence (category 4), providing a reference point against which to compare our model predictions.

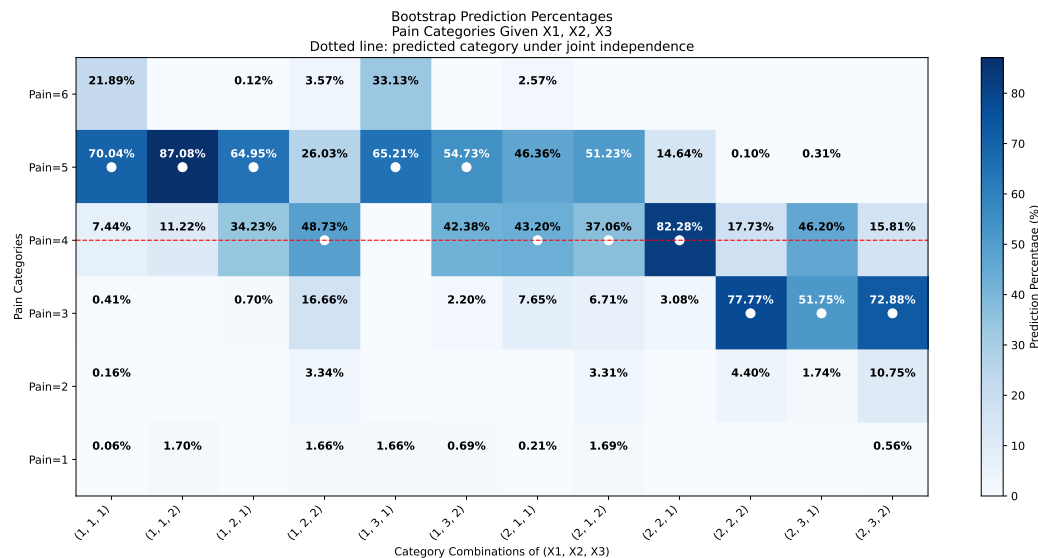


Figure 5.6: Bootstrap-based heatmap of predicted pain categories given  $(X_1, X_2, X_3)$ . Each cell shows the percentage of bootstrap samples predicting a given pain level for a specific combination of predictor values. White dots mark the most frequently predicted category, and the red dotted line indicates the expected response under joint independence.

## 5.10 Discussion and Clinical Interpretation

Our analysis of the back pain treatment dataset using the Checkerboard Copula Regression-based Visualization and Association Measure (CCRVAM) methodology reveals several clinically meaningful patterns that merit discussion. The CCRAM value of 0.2576 indicates that approximately 25.76% of the variation in back pain outcomes can be explained by the three predictor variables (Length of Previous Attack, Pain Change, and Lordosis). Similarly, the SCCRAM value of 0.2687 provides a standardized measure that accounts for the theoretical maximum association possible in this dataset structure. These values demonstrate a moderate but meaningful association between our predictors and pain outcomes, consistent with the



multifactorial nature of back pain treatment response.

Patients with short previous attacks ( $X_1 = 1$ ) generally experienced better outcomes (categories 4-5: moderate to marked improvement) compared to those with long previous attacks ( $X_1 = 2$ ). This finding suggests that the chronicity of pain prior to treatment may be an important prognostic factor, with early intervention potentially yielding better results. The influence of pain change ( $X_2$ ) appears to interact with the length of previous attack. For patients with short previous attacks, both better ( $X_2 = 1$ ) and worse ( $X_2 = 3$ ) pain change categories often led to marked improvement (category 5), while for patients with long previous attacks, worse pain change ( $X_2 = 3$ ) generally predicted poorer outcomes (category 3: slight improvement). This interaction effect highlights the complex nature of pain response trajectories. The presence or absence of lordosis ( $X_3$ ) appears to have a more subtle influence on outcomes compared to the other predictors, often modifying the effects of the primary predictors rather than driving outcomes independently.

The permutation test results provide strong statistical evidence that the observed associations are not due to random variation. With p-values of 0.0016 for CCRAM and 0.0011 for SCCRAM, we can confidently reject the null hypothesis of no association between our predictor variables and back pain outcomes. The observed values exceed the 99th percentile of their respective null distributions, further strengthening the significance of our findings.

The bootstrap analysis of our predictions reveals varying levels of confidence across different predictor combinations. For instance, for patients with short previous attack, better pain change, and absent lordosis ( $X_1 = 1, X_2 = 1, X_3 = 1$ ), the prediction of marked improvement (category 5) occurs in 70.90% of bootstrap samples. For patients with short previous attack, better pain change, and present lordosis ( $X_1 = 1, X_2 = 1, X_3 = 2$ ), this prediction is even

more consistent, occurring in 87.17% of bootstrap samples. In contrast, patients with long previous attack, better pain change, and present lordosis ( $X_1 = 2, X_2 = 1, X_3 = 2$ ) show more uncertainty with 52.06% of samples predicting category 5 and 36.29% predicting category 4. These confidence metrics provide valuable context for clinical decision-making, indicating where predictions are most reliable and where greater caution may be warranted.

These findings have several implications for clinical practice. Clinicians may use these results to provide more informed prognostic guidance to patients based on their specific combination of risk factors. The identification of patients with long previous attacks and worsening pain as having poorer outcomes may suggest the need for more aggressive or multimodal intervention approaches for this subgroup. The generally better outcomes observed in patients with shorter previous attacks reinforce the importance of early treatment initiation for back pain. Additionally, the bootstrap prediction matrices can help calibrate expectations for both clinicians and patients, providing nuanced probability estimates rather than deterministic predictions.

This analysis demonstrates several key advantages of the CCRVAM methodology for analyzing categorical and ordinal data in clinical research. The approach respects the ordinal nature of both predictors and outcomes without imposing arbitrary numerical scoring. CCRVAM naturally captures complex interactions among predictors without requiring explicit interaction terms. The visualization tools provided by the `ccrvam` package facilitate intuitive interpretation of multidimensional patterns. The integration of bootstrap and permutation methods provides comprehensive uncertainty assessment without parametric assumptions.

Despite these strengths, several limitations should be acknowledged. The dataset includes 101 observations distributed across 72 possible combinations of predictor and outcome categories,

resulting in sparse data for some combinations. As with any observational study, unmeasured confounders may influence the observed associations. The cross-sectional nature of the data limits our ability to assess temporal relationships and treatment dynamics.

The application of CCRVAM methodology to this back pain treatment dataset using our `ccrvam` package has yielded clinically meaningful insights while demonstrating the utility of this approach for analyzing complex categorical data in medical research. The findings suggest that patient characteristics, particularly the length of previous pain episodes and early treatment response, significantly influence treatment outcomes. These results may inform more personalized approaches to back pain management and highlight the value of sophisticated methodological tools for extracting meaningful patterns from categorical clinical data.

## Chapter 6

# Conclusion and Future Work

Exploratory Data Analysis of multivariate categorical data with ordinal responses presents unique challenges that traditional statistical methodologies often struggle to address effectively. This thesis has presented a comprehensive framework that synthesizes copula theory, regression dependence concepts, and modern computational approaches to overcome these limitations. By starting with an introduction to dependence, exploring continuous copula, bridging to discrete cases, building upon the methodological innovations proposed by Wei and Kim (2021) and extending their practical implementation, we have developed a robust analytical workflow encapsulated in a novel Python package, providing researchers with accessible tools for rigorous analysis of complex categorical data structures. The primary contribution of this work lies in bridging the theoretical foundations of copula theory with practical exploratory data analysis, resulting in a unified framework specifically tailored for categorical data analysis. The CCRVAM Python package represents the first end-to-end implementation of Checkerboard Copula Regression methods, enabling researchers to apply these sophisticated techniques without requiring extensive knowledge of the underlying

mathematical complexities. Furthermore, through rigorous testing on both simulated and real-world datasets, we have demonstrated the robustness and utility of the CCR approach and (S)CCRAM measures for quantifying regression dependence in multivariate categorical data.

Unlike traditional association measures that often treat variables symmetrically, CCR approach explicitly accommodates regression dependence, preserving the crucial distinction between response and explanatory variables. This orientation toward regression dependence makes the methods particularly valuable in fields such as healthcare, social sciences, and economics, where understanding causal and predictive relationships is paramount. The model-free nature of CCR and (S)CCRAM offers significant advantages over parametric approaches by avoiding potentially restrictive assumptions about functional forms, allowing researchers to explore dependencies in data without imposing predetermined structures. This flexibility potentially reveals patterns that might be missed by conventional modeling approaches. Additionally, the development of novel visualization methods enhances the interpretability of dependence structures in categorical data, going beyond traditional contingency table representations to provide intuitive visual insights into complex multivariate relationships.

Despite these advances, several limitations warrant acknowledgment in the current implementation. As the dimensionality of the contingency table increases, the computational requirements for CCR analysis grow substantially, potentially challenging the analysis of extremely high-dimensional datasets despite our implemented optimizations. While (S)CCRAM provides a quantitative measure of association strength, interpreting the practical significance of specific values remains more art than science, lacking the direct translation to marginal effects that traditional regression coefficients offer. Additionally, as with many non-parametric methods, CCR approaches require adequate sample sizes to reliably estimate dependence

structures, particularly for complex multivariate scenarios. The current implementation also requires complete data in contingency tables, with no built-in mechanisms for addressing missing values beyond preprocessing steps.

The integration of bootstrap resampling and permutation testing within our framework provides much-needed statistical inference capabilities previously unavailable for model-free approaches to categorical data analysis. These techniques enable researchers to quantify uncertainty in both predictions and association measures, offering confidence intervals and formal hypothesis tests that enhance the rigor of categorical data analysis. The CCRVAM package thus democratizes access to advanced statistical methods by abstracting away complex mathematical details while preserving methodological rigor, enabling researchers across disciplines to incorporate these techniques into their analytical workflows without specialized statistical expertise.

This work opens numerous avenues for future research and methodological extensions. Developing methodology for conditional (S)CCRAM as theorized in Wei et al. (2023) would significantly extend these measures' utility by quantifying partial dependence after controlling for confounding variables. Extending the framework to accommodate simultaneous analysis of continuous and categorical variables would broaden its applicability across diverse research contexts, while adapting CCR methods for longitudinal or time series categorical data would enable exploration of temporal dependencies in ordinal outcomes. From a computational perspective, research into more efficient algorithmic approaches could reduce the burden for high-dimensional scenarios, potentially leveraging sparse matrix representations or approximation methods. Implementations utilizing parallel computing frameworks and GPU acceleration would enhance scalability for large datasets, particularly for computationally intensive bootstrap and permutation procedures.

Further advances in visualization and explainability present another promising direction, with potential for developing interactive visualization components that allow researchers to dynamically explore dependence structures. Research into decomposing (S)CCRAM values into interpretable components could enhance the explanatory power of the measure. Domain-specific implementations tailored for healthcare, finance, or social sciences could enhance adoption by incorporating field-specific knowledge and constraints, while developing bridges to common statistical environments would facilitate incorporation into established research pipelines.

The exploration of regression dependence in multivariate categorical data remains a rich and challenging area of statistical research. As data collection continues to proliferate across fields, the need for sophisticated yet accessible methods for categorical data analysis will only increase. This thesis has contributed to advancing this field by reviewing theoretical insights and engineering novel practical tools through the lens of checkerboard copulas and associated regression methods. By enabling model-free exploration of complex dependence structures, we hope to enhance researchers' ability to extract meaningful insights from categorical data across diverse domains. The framework presented in this thesis represents one step toward addressing this need, but much work remains to fully realize the potential of copula-based approaches for exploratory data analysis of categorical information. We hope that both the methodological advances and software tools presented here will serve as valuable resources for researchers and inspire continued innovation.

# References

- Anderson, J. A. (1984), “[Regression and ordered categorical variables](#),” *Journal of the Royal Statistical Society. Series B (Methodological)*, [Royal Statistical Society, Oxford University Press], 46, 1–30.
- Anscombe, F. J. (1973), “[Graphs in statistical analysis](#),” *The American Statistician*, American Statistical Association, Taylor & Francis, Ltd., 27, 17–21.
- Batchelder, N., and Coverage.py, C. to (2025), “Coverage.py: The code coverage tool for python,” <https://coverage.readthedocs.io/>.
- Buja, A., Cook, D., Hofmann, H., Lawrence, M., Lee, E.-K., Swayne, D. F., and Wickham, H. (2009), “Statistical inference for exploratory data analysis and model diagnostics,” *Philosophical Transactions of the Royal Society of London*, 367, 4361–4383.
- Davison, A. C., and Hinkley, D. V. (1997), *Bootstrap methods and their application*, Cambridge series in statistical and probabilistic mathematics, Cambridge: Cambridge University Press.
- Denuit, M., and Lambert, P. (2005), “Constraints on concordance measures in bivariate discrete data,” *Journal of Multivariate Analysis*, 93, 40–57.
- Developers, N. (2025), “NumPy 2.2.4: Fundamental package for array computing in python,” <https://pypi.org/project/numpy/>.



- Donoho, D. L. (2017), “50 years of data science,” *Journal of Computational and Graphical Statistics*, 26, 745–766.
- Efron, B. (1987), “Better bootstrap confidence intervals,” *Journal of the American Statistical Association*, 82, 171–185. <https://doi.org/10.2307/2289144>.
- 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>.
- Gelman, A., and Vehtari, A. (2021), “What are the most important statistical ideas of the past 50 years?” *Journal of the American Statistical Association*, 116, 2087–2097.
- Genest, C., Nešlehová, J. G., and Rémillard, B. (2014), “On the empirical multilinear copula process for count data,” *Bernoulli*, 20, 1344–1371.
- Genest, C., Nešlehová, J., and Remillard, B. (2017), “Asymptotic behavior of the empirical multilinear copula process under broad conditions,” *Journal of Multivariate Analysis*, 159. <https://doi.org/10.1016/j.jmva.2017.04.002>.
- Hesterberg, T. (2014), “What teachers should know about the bootstrap: Resampling in the undergraduate statistics curriculum.”
- Hofert, M., Kojadinovic, I., Maechler, M., and Yan, J. (2018), *Elements of Copula Modeling with r*, Springer Use R! Series.
- Hunter, J. D. (2007), “Matplotlib: A 2D graphics environment,” *Computing in Science & Engineering*, IEEE COMPUTER SOC, 9, 90–95. <https://doi.org/10.1109/MCSE.2007.55>.
- Joe, H. (2014), *Dependence modeling with copulas*, Chapman; Hall/CRC, New York.

- Krekel, H. (2025), “Pytest 8.3.5: Simple powerful testing with python,” <https://pypi.org/project/pytest/>.
- Liao, S.-M., Wang, L., and Kim, D. (2024), “Visualization of dependence in multidimensional contingency tables with an ordinal dependent variable via copula regression,” in *Dependent data in social sciences research: Forms, issues, and methods of analysis*, eds. M. Stemmler, W. Wiedermann, and F. L. Huang, Cham: Springer International Publishing, pp. 517–538. [https://doi.org/10.1007/978-3-031-56318-8\\_21](https://doi.org/10.1007/978-3-031-56318-8_21).
- Mckinney, W. (2011), “Pandas: A foundational python library for data analysis and statistics,” *Python High Performance Science Computer*.
- Nelsen, R. B. (2006), *An introduction to copulas*, Springer Science & business media.
- Nešlehová, J. (2007), “On rank correlation measures for non-continuous random variables,” *Journal of Multivariate Analysis*, 98, 544–567.
- 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>.
- Schweizer, B., and Sklar, A. (1974), “Operations on distribution functions not derivable from operations on random variables,” *Studia Mathematica*, 52, 43–52.
- Shearer, C. (2000), “The CRISP-DM model: The new blueprint for data mining,” *Journal of Data Warehousing*, 5, 13–22.
- Sklar, M. (1959), “Fonctions de repartition an dimensions et leurs marges,” *Publ. inst. statist. univ. Paris*, 8, 229–231.
- Sveidqvist, K., and Mermaid, C. to (2014), *Mermaid: Generate diagrams from markdown-like text*, <https://mermaid.js.org/>.
- Tufte, E. R. (1983), *The visual display of quantitative information*, Cheshire, CT: Graphics

Press.

Tukey, J. W. (1977), *Exploratory data analysis*, AddisonWesley; Boston, MA.

Turner, A. (2025), “Sphinx 8.2.3: Python documentation generator,” <https://pypi.org/project/Sphinx/>.

Ushey, K., and Wickham, H. (2024), *Renv: Project environments*.

VG, C. (2024), “PyPI-template 0.8.0: Template-based common/best practices for managing a python package on PyPi,” <https://pypi.org/project/pypi-template/>.

Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., Walt, S. J. van der, Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., Carey, C. J., Polat, İ., Feng, Y., Moore, E. W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H., Pedregosa, F., Mulbregt, P. van, Vijaykumar, A., Bardelli, A. P., Rothberg, A., Hilboll, A., Kloeckner, A., Scopatz, A., Lee, A., Rokem, A., Woods, C. N., Fulton, C., Masson, C., Häggström, C., Fitzgerald, C., Nicholson, D. A., Hagen, D. R., Pasechnik, D. V., Olivetti, E., Martin, E., Wieser, E., Silva, F., Lenders, F., Wilhelm, F., Young, G., Price, G. A., Ingold, G.-L., Allen, G. E., Lee, G. R., Audren, H., Probst, I., Dietrich, J. P., Silterra, J., Webber, J. T., Slavič, J., Nothman, J., Buchner, J., Kulick, J., Schönberger, J. L., Miranda Cardoso, J. V. de, Reimer, J., Harrington, J., Rodríguez, J. L. C., Nunez-Iglesias, J., Kuczynski, J., Tritz, K., Thoma, M., Newville, M., Kümmerer, M., Bolingbroke, M., Tartre, M., Pak, M., Smith, N. J., Nowaczyk, N., Shebanov, N., Pavlyk, O., Brodtkorb, P. A., Lee, P., McGibbon, R. T., Feldbauer, R., Lewis, S., Tygier, S., Sievert, S., Vigna, S., Peterson, S., More, S., Pudlik, T., Oshima, T., Pingel, T. J., Robitaille, T. P., Spura, T., Jones, T. R., Cera, T., Leslie, T., Zito, T., Krauss, T., Upadhyay, U., Halchenko, Y. O., and Vázquez-Baeza, Y. (2020), “SciPy 1.0: Fundamental algorithms

for scientific computing in python,” *Nature Methods*, Springer Science; Business Media LLC, 17, 261–272. <https://doi.org/10.1038/s41592-019-0686-2>.

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>.

Wei, Z., Wang, L., Liao, S.-M., and Kim, D. (2023), “On the exploration of regression dependence structures in multidimensional contingency tables with ordinal response variables,” *Journal of Multivariate Analysis*, 196, 105179. <https://doi.org/https://doi.org/10.1016/j.jmva.2023.105179>.

# Appendix A

## Code availability

This thesis is written using Quarto with **renv** (Ushey and Wickham 2024) to create a reproducible environment. All materials required to reproduce this document—including data sets, source files, and implementation code—are publicly available in the GitHub repository [github.com/DhyeyMavani2003/ccrvam](https://github.com/DhyeyMavani2003/ccrvam). The repository contains the complete codebase for the CCRVAM methodology, including the test suite that verifies correct implementation.

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

Contributors are welcome to submit issues or pull requests to the GitHub repository to help improve the methodology and its implementation.

### A.1 Chapter 2 Code

The following code was used to create Chapter 2. The motivating example in Chapter 2 demonstrates how correlation measures can vary dramatically under nonlinear transforma-

tions while copula-based measures remain invariant. The code generates bivariate normal data and transforms it via gamma and beta distributions to illustrate this phenomenon:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.stats import beta, expon, norm, gamma, binom
4 import os
5
6 # Create directory if not exists
7 fig_dir = "fig"
8 os.makedirs(fig_dir, exist_ok=True)
9
10 # Generate Data
11 np.random.seed(8990)
12 n = 10000
13 mean = [0, 0]
14 cov = [[1, 0.8], [0.8, 1]]
15 X = np.random.multivariate_normal(mean, cov, size=n)
16 X1, X2 = X[:, 0], X[:, 1]
17
18 # Transform U_X1 and U_X2 to uniform [0, 1]
19 # using the CDF of the normal distribution
20 U_X1 = norm.cdf(X1)
21 U_X2 = norm.cdf(X2)
```

```

22
23 # Transform U_X1 and U_X2 into Gamma and Beta distributions
24 Y1 = gamma.ppf(U_X1, a=3, scale=1/15)
25 Y2 = beta.ppf(U_X2, a=5, b=3)
26
27 # Calculate Pearson Correlation Coefficients
28 rho_X = np.corrcoef(X1, X2)[0, 1]
29 rho_Y = np.corrcoef(Y1, Y2)[0, 1]
30 print("Pearson correlation for (X1, X2):", rho_X)
31 print("Pearson correlation for (Y1, Y2):", rho_Y)
32
33 # Create Layout design and Set Size-Ratio
34 fig, axes = plt.subplots(1, 2, figsize=(6, 4))
35
36 # Scatter plot for (X1, X2)
37 axes[0].scatter(X1, X2, alpha=0.3, s=5)
38 axes[0].set_title("Scatter plot of (X1, X2)")
39 axes[0].set_xlabel("X1")
40 axes[0].set_ylabel("X2")
41
42 # Add marginal histograms
43 axes[0].hist(
44     X1, bins=50, density=True, alpha=0.9, color='blue',
45     orientation='vertical', histtype='step'

```

```

46 )
47 axes[0].hist(
48     X2, bins=50, density=True, alpha=0.9, color='red',
49     histtype='step', orientation='horizontal'
50 )
51
52 # Scatter plot for (Y1, Y2)
53 axes[1].scatter(Y1, Y2, alpha=0.3, s=5)
54 axes[1].set_title("Scatter plot of (Y1, Y2)")
55 axes[1].set_xlabel("Y1")
56 axes[1].set_ylabel("Y2")
57
58 # Add marginal histograms
59 axes[1].hist(
60     Y1, bins=50, density=True, alpha=0.9, color='blue',
61     orientation='vertical', histtype='step'
62 )
63 axes[1].hist(
64     Y2, bins=50, density=True, alpha=0.9, color='red',
65     histtype='step', orientation='horizontal'
66 )
67
68 # Organize into a tight layout as per matplotlib

```



```

69 plt.tight_layout()
70
71 # Save figure instead of showing it
72 fig_path = os.path.join(fig_dir, "motivating_example.png")
73 plt.savefig(fig_path, dpi=300, bbox_inches='tight')
74
75 # Close the figure to prevent rendering output
76 plt.close(fig)

```

```

1 # Set random seed for reproducibility
2 np.random.seed(8990)
3
4 # Apply probability integral transformation to all variables
5 # in order to make them uniform
6 U_Y1 = gamma.cdf(Y1, a=3, scale=1/15)
7 U_Y2 = beta.cdf(Y2, a=5, b=3)
8
9 # Calculate Pearson Correlation Coefficients
10 rho_U_X = np.corrcoef(U_X1, U_X2)[0, 1]
11 rho_U_Y = np.corrcoef(U_Y1, U_Y2)[0, 1]
12 print("Pearson correlation for (F_1(X_1)$, F_2(X_2)$):", rho_U_X)
13 print("Pearson correlation for (G_1(Y_1)$, G_2(Y_2)$):", rho_U_Y)
14
15 # Combine transformed data

```

```

16 uniform_data = np.vstack([U_X1, U_X2, U_Y1, U_Y2]).T
17
18 # Verify the uniformity of transformed data (Should be 0.5 in value)
19 print("U_X1 mean:", U_X1.mean(), "U_X2 mean:", U_X2.mean())
20 print("U_Y1 mean:", U_Y1.mean(), "U_Y2 mean:", U_Y2.mean())
21
22 # Create Layout design and Set Size-Ratio
23 fig, axes = plt.subplots(1, 2, figsize=(6, 4))
24
25 # Scatter plot for (U_X1, U_X2)
26 axes[0].scatter(U_X1, U_X2, alpha=0.3, s=5)
27 axes[0].set_title("Scatter plot of ( $F_1(X_1)$ ,  $F_2(X_2)$ )")
28 axes[0].set_xlabel(" $F_1(X_1)$ ")
29 axes[0].set_ylabel(" $F_2(X_2)$ ")
30
31 # Add marginal histograms
32 axes[0].hist(
33     U_X1, bins=50, density=True, alpha=0.9, color='blue',
34     orientation='vertical', histtype='step'
35 )
36 axes[0].hist(
37     U_X2, bins=50, density=True, alpha=0.9, color='red',
38     histtype='step', orientation='horizontal'

```

```

39 )
40
41 # Scatter plot for (U_Y1, U_Y2)
42 axes[1].scatter(U_Y1, U_Y2, alpha=0.3, s=5)
43 axes[1].set_title("Scatter plot of ($G_1(Y_1)$, $G_2(Y_2)$)")
44 axes[1].set_xlabel("$G_1(Y_1)$")
45 axes[1].set_ylabel("$G_2(Y_2)$")
46
47 # Add marginal histograms
48 axes[1].hist(
49     U_Y1, bins=50, density=True, alpha=0.9, color='blue',
50     orientation='vertical', histtype='step'
51 )
52 axes[1].hist(
53     U_Y2, bins=50, density=True, alpha=0.9, color='red',
54     histtype='step', orientation='horizontal'
55 )
56
57 # Organize into a tight layout as per matplotlib
58 plt.tight_layout()
59
60 # Save figure instead of showing it
61 fig_path = os.path.join(fig_dir, "transformed_motivating_example.png")

```

```

62 plt.savefig(fig_path, dpi=300, bbox_inches='tight')
63
64 # Close the figure to prevent rendering output
65 plt.close(fig)

1 # Set random seed for reproducibility
2 np.random.seed(8990)
3
4 # Transform (Y1,Y2) back to normal marginals using quantile transformation
5 F1_Y1 = norm.ppf(gamma.cdf(Y1, a=3, scale=1/15))
6 F2_Y2 = norm.ppf(beta.cdf(Y2, a=5, b=3))
7
8 # Calculate Pearson Correlation Coefficients
9 rho_F_Y = np.corrcoef(F1_Y1, F2_Y2)[0, 1]
10 print("Pearson correlation for transformed:", rho_F_Y)
11 print("Pearson correlation between X1 and X2:", rho_X)
12
13 # Plot the scatter plots with marginal histograms
14 fig, axes = plt.subplots(1, 2, figsize=(6, 4))
15
16 # Scatter plot for original normal marginals (X1, X2)
17 axes[0].scatter(X1, X2, alpha=0.3, s=10)
18 axes[0].set_title("Scatter plot ($F_1(X_1)$, $F_2(X_2)$)")
19 axes[0].set_xlabel("$F_1(X_1)$")

```

```

20 axes[0].set_ylabel("$F_2(X_2)$")
21 axes[0].hist(
22     X1, bins=50, density=True, alpha=0.6, color='blue', histtype='step'
23 )
24 axes[0].hist(
25     X2, bins=50, density=True, alpha=0.6, color='red',
26     histtype='step', orientation='horizontal'
27 )
28
29 # Scatter plot for transformed normal marginals (F1_Y1, F2_Y2)
30 axes[1].scatter(F1_Y1, F2_Y2, alpha=0.3, s=10)
31 axes[1].set_title(
32     "Scatter plot ( $F_1^{-1}(G_1(Y_1))$ ,  $F_2^{-1}(G_2(Y_2))$ )"
33 )
34 axes[1].set_xlabel(" $F_1^{-1}(G_1(Y_1))$ ")
35 axes[1].set_ylabel(" $F_2^{-1}(G_2(Y_2))$ ")
36 axes[1].hist(
37     F1_Y1, bins=50, density=True, alpha=0.6, color='blue', histtype='step'
38 )
39 axes[1].hist(
40     F2_Y2, bins=50, density=True, alpha=0.6, color='red',
41     histtype='step', orientation='horizontal'
42 )

```

```

43
44 # Layout adjustment and save the figure
45 plt.tight_layout()
46 fig_path = os.path.join(
47     fig_dir, "quantile_transformed_motivating_example.png"
48 )
49 plt.savefig(fig_path, dpi=300, bbox_inches="tight")
50
51 # Close the figure to prevent rendering output
52 plt.close(fig)

```

## A.2 Chapter 3 Code

The following code was used to create Chapter 3. The code below creates the visual representation of a checkerboard copula density function that serves as the foundation for our methodology:

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.colors as colors
4
5 def create_copula_density_plot():
6     # Create figure
7     fig, ax = plt.subplots(figsize=(3, 3))
8

```

```

9      # Define grid divisions
10     u1_divisions = np.array([0, 2/8, 3/8, 5/8, 6/8, 1])
11     u2_divisions = np.array([0, 2/8, 4/8, 1])
12
13     # Determine the number of blocks
14     n_rows = len(u1_divisions) - 1
15     n_cols = len(u2_divisions) - 1
16
17     # Create meshgrid for the full plot
18     du1 = np.diff(u1_divisions)
19     du2 = np.diff(u2_divisions)
20
21     # Define the density values in each block
22     # 0: light blue, 2: purple, 4: magenta
23     # Reshaping the density values to match the grid structure
24     density_values = np.array([
25         [0, 0, 2], # Bottom row (U1 from 0 to 2/8)
26         [0, 4, 0], # Second row (U1 from 2/8 to 3/8)
27         [4, 0, 0], # Third row (U1 from 3/8 to 5/8)
28         [0, 4, 0], # Fourth row (U1 from 5/8 to 6/8)
29         [0, 0, 2], # Top row (U1 from 6/8 to 1)
30     ])
31

```

```

32     # Create colormap for the specific values (0, 2, 4)
33     cmap = colors.ListedColormap(['lightblue', 'violet', 'magenta'])
34     bounds = [-0.5, 0.5, 2.5, 4.5]
35     norm = colors.BoundaryNorm(bounds, cmap.N)
36
37     # Plot the piecewise constant density
38     for i in range(n_rows):
39         for j in range(n_cols):
40             value = density_values[i, j]
41             rect = plt.Rectangle(
42                 (u2_divisions[j], u1_divisions[i]),
43                 du2[j], du1[i],
44                 facecolor=cmap(norm(value)),
45                 alpha=1.0,
46                 edgecolor='black',
47                 linewidth=0.5
48             )
49             ax.add_patch(rect)
50
51     # Add text with density value
52     ax.text(
53         u2_divisions[j] + du2[j]/2,
54         u1_divisions[i] + du1[i]/2,

```



```

55         str(int(value)),
56         horizontalalignment='center',
57         verticalalignment='center',
58         fontsize=16,
59         color='black'
60     )
61
62     # Set axis labels and limits
63     ax.set_xlabel('$U_2$', fontsize=14)
64     ax.set_ylabel('$U_1$', fontsize=14)
65     ax.set_xlim(0, 1)
66     ax.set_ylim(0, 1)
67
68     # Add grid lines at each division
69     for u in u1_divisions:
70         ax.axhline(y=u, color='black', linestyle=':', linewidth=1)
71     for u in u2_divisions:
72         ax.axvline(x=u, color='black', linestyle=':', linewidth=1)
73
74     # Add colorbar
75     sm = plt.cm.ScalarMappable(cmap=cmap, norm=norm)
76     sm.set_array([])
77     cbar = plt.colorbar(sm, ax=ax, ticks=[0, 2, 4])

```

```

78     cbar.set_label('')
79
80     # Add title
81     fig.text(
82         0.5, 0.05, "Copula density  $c^{+}(u_1, u_2)$ ",
83         ha='center', fontsize=14
84     )
85
86     # Adjust layout
87     plt.tight_layout(rect=[0, 0.07, 1, 1])
88
89     return fig
90
91 # Generate the plot
92 fig = create_copula_density_plot()
93
94 # For .qmd file, you would save the figure
95 plt.savefig('fig/copula_density_plot.png', dpi=300, bbox_inches='tight')
96 plt.close(fig)

```

## 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. ...