



GEORG-AUGUST-UNIVERSITÄT  
GÖTTINGEN

# Master Thesis

Multivariate modelling of the  
dependency structure between  
article sales of a sportswear  
manufacturer

Author

**Petros Christanas**

from Nuremberg

Matriculation Number

11604278

Applied Statistics M.Sc.

Chair of Statistics and Econometrics

Supervisors

**Prof. Dr. Thomas Kneib**

**Dipl.-Vw. Quant. Fabian H. C. Raters**

Submitted March 28, 2020

Processing time of 20 weeks



## **Confidentiality Clause**



# Statutory Declaration



## Acknowledgments





# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	adidas . . . . .	1
1.2	Data Sources . . . . .	1
1.3	Motivation . . . . .	3
<b>2</b>	<b>Statistical Theory &amp; Methods</b>	<b>5</b>
2.1	Generalized Linear Models . . . . .	5
2.2	Mixed Effects Models . . . . .	5
2.3	Additive Models . . . . .	6
<b>3</b>	<b>Copulas &amp; Dependency Structures</b>	<b>9</b>
3.1	Introduction to Copulas . . . . .	9
3.2	Copula Classes . . . . .	12
3.2.1	Fundamental Copulas . . . . .	12
3.2.2	Elliptical Copulas . . . . .	13
3.2.3	Archimedean Copulas . . . . .	15
3.3	Dependence Measures . . . . .	16
3.3.1	Linear Correlation . . . . .	17
3.3.2	Rank Correlation . . . . .	17
3.3.3	Tail Dependence . . . . .	19
3.4	Conditional Copulas . . . . .	19
3.5	Vine Copulas . . . . .	19
<b>4</b>	<b>Data Exploration</b>	<b>21</b>
<b>5</b>	<b>Modelling</b>	<b>23</b>
<b>6</b>	<b>Conclusion</b>	<b>25</b>
	<b>Appendix</b>	<b>27</b>
	<b>List of Figures</b>	<b>29</b>
	<b>List of Tables</b>	<b>31</b>
	<b>List of Abbreviations</b>	<b>33</b>

**References****35**

# 1 Introduction

Write introduction here and "upper" subsections here (adidas, Motivation, etc...)

## 1.1 adidas

## 1.2 Data Sources

Throughout each season, transactional data are collected from online purchases of the sports brand's eCommerce website. Specifically, we are provided with weekly sales data for western European countries. A short description is depicted in Table 1.1.

Column	Description	Values
week_id	Calendar week of a specific year (YYYYWW)	Factors: 201648, ..., 201852
article_number	Unique article identification number (article ID)	Factors: 10669, 10, ...
min_date_of_week	Minimum date of the respective week; always a Monday (YYYY-MM-DD)	Dates: 2016-11-28, ..., 2018-12-24
art_min_price	Minimal recorded price of the article	Non-negative (integer) value
month_id	Calendar month of a specific year (YYYYMM)	Factors: 201612, ..., 201812
season	Season of year (format: SSYY) (Spring-Summer [SS]: December - May) Fall-Winter [FW]: June - November)	Factors: SS17, FW17, SS18, FW18, SS19
bf_w	Weekly "Black Friday" promotion intensity of the article	Between 0 and 1
ff_w	Weekly "Friends & Family" promotion intensity of the article	Between 0 and 1
ot_w	Weekly article promotion intensity of "Other" type	Between 0 and 1
gross_demand_quantity	Weekly amount of added articles to shopping cart	Non-negative (integer) value
base_price_locf	Retail price of the article without any discounts	Non-negative (integer) value
total_markdown_pct		
day_of_month	Day of the month	Integers: 1 - 31
month_of_year	Month of the year	Factors: January, ..., December
year	Year	Integers: 2016, 2017, 2018
week_of_year	Week of the year	Integers: 1 - 52

Table 1.1: Transactional raw data description from online purchases of western European countries

Due to legal regulations of the company, some columns had to undergo anonymization in order for the data to be released. To ensure data protection and confidentiality, numeric variables (with exception of time-indicating columns) were transformed. As a consequence for the analysis part, most integer values were converted to float numbers. This fact should be kept in mind by the reader, since the above table serves as a reminder and reference point for the data

documentation.

Another peculiarity of this setup is to be considered, too. We will often refer to the variable *gross demand quantity* as *sales*, even though it is obviously not exactly the same. In the eCommerce environment, there are several stages before the purchase is complete, e.g. addition to cart, removal from cart, proceeding to checkout & even the return of bought articles. Targeting the articles added to cart, i.e. the (gross) demand quantity, provides the optimal data extraction for analytical purposes and is the closest to adequately model the dependency structure between net sales of articles.

Besides the transactional data, *article master data*, i.e. attributes of the articles, are provided and described in Table 1.2.

Column	Description	Values (all Factors)
article_number	Unique article identification number (article ID)	10669, 10, ...
gender	Gender type of the article (Men, Women, Unisex)	M, W, U
age_group	Age group of the article (Adult, Infant, Junior, Kids)	A, I, J, K
product_division_descr	Product division of the article	APPAREL, FOOTWEAR, HARDWARE
product_group_descr	Product group of the article	BAGS, BALLS, FOOTWEAR ACCESSORIES, SHOES, ...
color	Consolidated color group of the article	BEIGE, BLACK, BROWN, ORANGE, PINK, RED, ...
sports_category_descr	Sports category of the article	encoded: SC_1, ..., SC_22
sales_line_descr	Sales line of the article	encoded: SL_1, ..., SL_379
business_unit_descr	The article's Business Unit membership	encoded: BU_1, ..., BU_18
business_segment_descr	The article's Business Segment membership	encoded: BS_1, ..., BS_49
sub_brand_descr	Sub-brand of the article	encoded: sub-brand_1, ..., sub-brand_4
item_type	Item type of the article	encoded: IT_1, ..., IT_171
brand_element	Brand element of the article	encoded: BE_1, ..., BE_131
product_franchise_descr	Product franchise of the article	encoded: franchise_1, ..., franchise_72
product_line_descr	Product line of the article	encoded: PL_1, ..., PL_105
franchise_bin	Franchise indicator of the article	FRANCHISE, NON-FRANCHISE
category	Category of the article	encoded: category_1, category_2

Table 1.2: Article master data

Plenty of additional information is stored in the database, but we are neglecting columns omitted in these tables, as they are redundant, already summarized, transformed or simply do not provide any value.

Overall, we will be dealing with data collected over two years, namely the years 2017 and 2018, while some transactions of late 2016 might be attached marginally. In summary, after joining the transactional observations to the article

attributes by the article ID, this translates to a dataset of over 587,000 rows and over 30 variables.

## **1.3 Motivation**



## 2 Statistical Theory & Methods

This chapter introduces various statistical methods used during the conduction of this thesis. It is assumed that basic understanding and knowledge of the reader regarding mathematical foundations of statistics (like linear algebra, probability theory, etc) already exists.

### 2.1 Generalized Linear Models

*Generalized Linear Models (GLMs)* are an extension of the classical *Linear Regression Model (LM)*

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i, \quad i = 1, \dots, n$$

which in matrix notation can be written as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where the response variable  $y_i$  can take values from several probability distributions (e.g. Poisson, Binomial, Gamma and others), which are members of the exponential family [Fahrmeir et al., 2003]. The linear predictor

$$\eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i = \mathbf{x}_i' \boldsymbol{\beta} \quad (2.1)$$

is passed through a *response function*  $h$  (a one-to-one, twice differentiable transformation), such that

$$E(y_i) = h(\eta_i)$$

i.e. the expected value of the response variable belongs to the value range of  $h$ . The inverse of the response function,

$$g = h^{-1}$$

is called the *link function*.

### 2.2 Mixed Effects Models

*Linear Mixed Models (LMMs)* are powerful tools when dealing with clustered data or data with a longitudinal structure (repeated measurements of individuals). As in the classical LM, there are population-specific effects, namely the parameter

vector of *fixed effects*  $\beta$ , as well as the cluster- or individual-specific effects of such models called *random effects* [Fahrmeir et al., 2003]. In the following, we will refer to our clusters or individuals as "groups" for briefness. Mathematically speaking, the linear predictor  $\eta_{ij} = \mathbf{x}'_{ij}\beta$  is extended to

$$\eta_{ij} = \mathbf{x}'_{ij}\beta + \mathbf{u}'_{ij}\gamma_i, \quad j = 1, \dots, m, \quad i = 1, \dots, n_i, \quad (2.2)$$

where

- $i$  is the number of groups
- $j$  is the number of observations per group
- $\beta$  is the vector of fixed effects
- $\gamma_i$  is the vector of random effects
- $\mathbf{x}'_{ij}$  is the vector of covariates and
- $\mathbf{u}'_{ij}$  is a subvector of  $\mathbf{x}'_{ij}$ .

$\mathbf{x}'_{ij} = (1, x_{ij1}, \dots, x_{ijk})$  and  $\mathbf{u}'_{ij} = (1, u_{ij1}, \dots, u_{ijk})$  are therefore the design vectors and  $\varepsilon_{ij}$  are the error terms of the *measurement model*

$$y_{ij} = \mathbf{x}'_{ij}\beta + \mathbf{u}'_{ij}\gamma_i + \varepsilon_{ij}, \quad \varepsilon_{ij} \stackrel{i.i.d.}{\sim} N(0, \sigma^2) \quad (2.3)$$

or in matrix notation

$$\mathbf{y}_i = \mathbf{X}_i\beta + \mathbf{U}_i\gamma_i + \boldsymbol{\varepsilon}_i \quad (2.4)$$

for group  $i = 1, \dots, m$  with  $E(\boldsymbol{\varepsilon}_i) = \mathbf{0}$ .

Similar to GLMs, *Generalized Linear Mixed Models (GLMMs)* relate the linear mixed predictor 2.2 to the conditional mean  $\mu_{ij} = E(y_{ij}|\gamma_i)$  via a suitable response function  $h$ , such that  $\mu_{ij} = h(\eta_{ij})$  and thus the conditional density of  $y_{ij}$  belongs to the exponential family.

## 2.3 Additive Models

*Additive Models* expand models with just a linear predictor

$$\eta_i^{lin} = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}$$



(such as the LM) to

$$y_i = \eta_i^{add} + \varepsilon_i, \quad (2.5)$$

where

$$\eta_i^{add} = f_1(z_{i1}) + \dots + f_q(z_{iq}) + \eta_i^{lin}, \quad i = 1, \dots, n. \quad (2.6)$$

The functions  $f_1(z_1), \dots, f_q(z_q)$  are non-linear univariate *smooth effects* of the *continuous* covariates  $z_1, \dots, z_q$  and are defined as

$$f_j(z_j) = \sum_{l=1}^{d_j} \gamma_{jl} B_l(z_j) \quad (2.7)$$

with  $B_l(z_j)$  being *basis functions* for  $j = 1, \dots, q$  and  $d_j$  the number of basis functions for covariate  $z_j$ . The regression coefficients of the basis functions  $B_l(z_j)$  are labeled as  $\gamma_{jl}$ . There is a wide variety of basis functions which can be used to flexibly model the data in a non- or semiparametric manner. For more content on basis functions we refer to Wood [2017] and Fahrmeir et al. [2003]. The basis functions evaluated at the observed covariate values are summarized in the design matrices  $\mathbf{Z}_1, \dots, \mathbf{Z}_q$  and the additive model 2.5 can be written in matrix notation as

$$\mathbf{y} = \mathbf{Z}_1 \boldsymbol{\gamma}_1 + \dots + \mathbf{Z}_q \boldsymbol{\gamma}_q + \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

Accordingly, the vector of function values evaluated at the observed covariate values  $z_{1j}, \dots, z_{nj}$  is denoted by  $\mathbf{f}_j = (f_j(z_{1j}), \dots, f_j(z_{nj}))'$  and therefore  $\mathbf{f}_j = \mathbf{Z}_j \boldsymbol{\gamma}_j$ . To ensure identifiability of the additive model, the functions  $f_j(z_j)$  are centered around zero, such that

$$\sum_{i=1}^n f_1(z_{i1}) = \dots = \sum_{i=1}^n f_q(z_{iq}) = 0.$$

Extensions of additive models to non-normal responses are consequently called *Generalized Additive Models (GAMs)*.

MAYBE INCLUDE MODELS WITH INTERACTIONS; MIXED EFFECTS



### 3 Copulas & Dependency Structures

Multivariate distributions consist of the marginal distributions and the dependence structure between those marginals. These components can be specified separately in a single framework with the help of copula functions. This chapter introduces the concept of modelling such dependency structures with copulas, which is the main focus of this thesis.

#### 3.1 Introduction to Copulas

A  $d$ -dimensional function  $C : [0, 1]^d \rightarrow [0, 1]$  is called a *copula*, if it is a Cumulative Distribution Function (CDF) with uniform margins, i.e.

$$P(U_1 \leq u_1, \dots, U_d \leq u_d) = C(u_1, \dots, u_d)$$

where  $U_i, i = 1, \dots, d$  are uniformly distributed Random Variables (RVs) in  $[0, 1]$ .

Since  $C$  is a CDF, following properties emerge:

- $C(\mathbf{u}) = C(u_1, \dots, u_d)$  is increasing in each component  $u_i, i = 1, \dots, d$ .
- The  $i^{th}$  marginal distribution is obtained by setting  $u_j = 1$  for  $j \neq i$  and it has to be uniformly distributed

$$C(1, \dots, 1, u_i, 1, \dots, 1) = u_i$$

- For  $a_i \leq b_i$ , the probability  $P(U_1 \in [a_1, b_1], \dots, U_d \in [a_d, b_d])$  must be non-negative, so we obtain the *rectangle inequality*

$$\sum_{i_1=1}^2 \dots \sum_{i_d=1}^2 (-1)^{i_1+\dots+i_d} C(u_{1,i_1}, \dots, u_{d,i_d}) \geq 0, \quad (3.1)$$

where  $u_{j,1} = a_j$  and  $u_{j,2} = b_j$ .

The reverse is also true, i.e. any function  $C$  that satisfies the above properties is a copula. Furthermore,  $C(1, u_1, \dots, u_{d-1})$  is also a  $(d-1)$ -dimensional copula and thus all  $k$ -dimensional marginals with  $2 < k < d$  are copulas.

#### Generalized Inverse

For a CDF, the *generalized inverse* is defined by

$$F^{\leftarrow}(y) := \inf\{x : F(x) \geq y\}$$

(similar to the definition of a *quantile function*).

□

### Probability Transformation

If a RV  $Y$  has a continuous CDF  $F$ , then

$$F(Y) \sim U[0, 1]. \quad (3.2)$$

□

The reverse of the *probability transformation* is the *quantile transformation*.

### Quantile Transformation

If  $U \sim U[0, 1]$  and  $F$  be a CDF, then

$$P(F^{\leftarrow}(U) \leq x) = F(x) \quad (3.3)$$

□

The above two transformations allow us to move back and forth between  $\mathbb{R}^d$  and  $[0, 1]^d$  and are the primary building blocks regarding copulas. Against this backdrop, we introduce *Sklar's theorem* which is considered the foundation of all copula related applications.

### Sklar's Theorem Sklar [1959]

Let  $F$  be a  $d$ -dimensional CDF with marginal distributions  $F_i$ ,  $i = 1, \dots, d$ . Then there exists a copula  $C$  such that

$$F(x_1, \dots, x_d) = C(F_1(x_1), \dots, F_d(x_d)) \quad (3.4)$$

for all  $x_i \in \mathbb{R}$ ,  $i = 1, \dots, d$ .

The copula  $C$  is unique, if  $\forall i = 1, \dots, d$ ,  $F_i$  is continuous. Otherwise  $C$  is uniquely determined only on  $\text{Ran}(F_1) \times \dots \times \text{Ran}(F_d)$ , where  $\text{Ran}(F_i)$  is the range of  $F_i$ .

Conversely, if  $C$  is a  $d$ -dimensional copula and  $F_1, \dots, F_d$  are univariate CDF's, then  $F$  as defined in Equation 3.4 is a  $d$ -dimensional CDF.

□

If the copula has a Probability Density Function (PDF), then the *copula density* is defined as

$$c(\mathbf{u}) = \frac{\partial^d C(u_1, \dots, u_d)}{\partial u_1 \cdots \partial u_d} \quad (3.5)$$

for a differentiable copula function  $C$  and the realization of a random vector  $\mathbf{u} = (u_1, \dots, u_d)$ .

By virtue of Equation 3.4 in Sklar's theorem and given that

$$C(\mathbf{u}) = F(F_1^{\leftarrow}(u_1), \dots, F_d^{\leftarrow}(u_d)), \quad (3.6)$$

i.e. invertible CDFs  $F_i$ ,  $i = 1, \dots, d$ , we can rewrite the copula density to

$$c(u_1, \dots, u_d) = \frac{f(F_1^{\leftarrow}(u_1), \dots, F_d^{\leftarrow}(u_d))}{\prod_{i=1}^d f_i(F_i^{\leftarrow}(u_i))} \quad (3.7)$$

for densities  $f$  of  $F$  and  $f_1, \dots, f_d$  of the corresponding marginals.

### Invariance Principal

Suppose the RVs  $X_1, \dots, X_d$  have continuous marginals and copula  $C$ . For strictly increasing functions  $T_i : \mathbb{R} \rightarrow \mathbb{R}$ ,  $i = 1, \dots, d$ , the RVs  $T_1(X_1), \dots, T_d(X_d)$  also have copula  $C$ .

□

### Fréchet-Hoeffding Bounds

Let  $C(\mathbf{u}) = C(u_1, \dots, u_d)$  be any  $d$ -dimensional copula.

Then, for

$$W(\mathbf{u}) = \max \left\{ \sum_{i=1}^d u_i - d + 1, 0 \right\} \quad (3.8)$$

as well as

$$M(\mathbf{u}) = \min_{1 \leq i \leq d} \{u_i\}, \quad (3.9)$$

it holds that

$$W(\mathbf{u}) \leq C(\mathbf{u}) \leq M(\mathbf{u}), \quad \mathbf{u} \in [0, 1]^d. \quad (3.10)$$

We call  $W$  the *lower Fréchet-Hoeffding bound* and  $M$  the *upper Fréchet-Hoeffding*

*bound*.

Note that  $W$  is a copula if and only if  $d = 2$ , whereas  $M$  is a copula for all  $d \geq 2$  (more on this later in Section 3.2.1).

□

## MORE ON COPULA THEORY (NOTES)

### 3.2 Copula Classes

In this section we will take a look at three very popular *copula classes*, namely *fundamental*, *elliptical* and *archimedean copulas*. For each class, we will present a few (parametric) *copula families* which are widely used.

#### 3.2.1 Fundamental Copulas

Fundamental copulas are a basic class of copulas, which emerge directly from the copula framework and do not depend on any parametric components.

##### Independence Copula

It is well known that the joint probability of a finite set of events  $E_i, i = 1, \dots, k$ , is equal to the product of the marginals if and only if the Events  $E_i$  are independent.

$$P\left(\bigcap_{i=1}^k E_i\right) = \prod_{i=1}^k P(E_i).$$

Analogously, the same concept applies when we talk about the *independence copula*

$$\Pi(\mathbf{u}) = \prod_{i=1}^d u_i. \quad (3.11)$$

As a result of Sklar's theorem the RVs  $u_i$  are independent if and only if their copula is the independence copula, i.e.

$$C(\mathbf{u}) = \Pi(\mathbf{u})$$

and thus the copula density would be

$$c(\mathbf{u}) = 1, \quad \mathbf{u} \in [0, 1]^d.$$

□

From Equation 3.10, it is obvious that the Fréchet-Hoeffding bounds correspond to the extreme cases of perfect dependence between the RVs  $X_i, i = 1, \dots, d$ .

### Comonotonicity Copula

Consider the RVs  $X_1, \dots, X_d$  and strictly increasing transformations  $T_1, \dots, T_d$  and  $X_i = T(X_i)$  for  $i = 2, \dots, d$ . Making use of the *invariance principal*, it can be shown that these RVs have as copula the upper Fréchet-Hoeffding bound

$$M(\mathbf{u}) = \min\{u_1, \dots, u_d\}.$$

Since there is perfect positive dependence between those RVs, we call  $M$  the *comonotonicity copula*. The number of dimensions  $d$  can be any finite number greater than or equal to 2 for  $M$  to be a copula, as the minimum remains well defined.

□

### Countermonotonicity Copula

Similar to the comonotonic case, it can be shown that if two RVs  $X_1$  and  $X_2$  are perfectly negatively dependent, their copula is the lower Fréchet-Hoeffding bound

$$W(\mathbf{u}) = \max\left\{\sum_{i=1}^d u_i - d + 1, 0\right\}.$$

Therefore,  $W$  is known as the *countermonotonicity copula*. Because of the fact that countermonotonicity is not valid for a dimension greater than 2, we end up with the restriction  $d = 2$  for  $W$  to be indeed a copula.

□

### 3.2.2 Elliptical Copulas

Copulas which can be derived from known multivariate distributions like for example the *Multivariate Normal (or Gaussian) Distribution* or the *Multivariate Student's t-Distribution* are called *implicit copulas*. *Elliptical copulas* are implicit copulas which arise via Sklar's theorem from elliptical distributions like the mentioned examples.

### Gaussian Copula

Without loss of generality (w.l.o.g.), for a random vector  $\mathbf{X} \sim \mathcal{N}_d(\mathbf{0}, \mathbf{P})$  and correlation matrix  $\mathbf{P}$ , the *Gaussian copula (family)* is given by

$$C_{\mathbf{P}}^{Ga}(\mathbf{u}) = \Phi_{\mathbf{P}}(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_d)), \quad (3.12)$$

where  $\Phi$  is the CDF of  $\mathcal{N}(0, \sigma^2)$  and  $\Phi_{\mathbf{P}}$  is the CDF of  $\mathcal{N}_d(\mathbf{0}, \mathbf{P})$ .

There are special cases to this copula family, namely for  $d = 2$  and correlation  $\rho$ , the *bivariate Gaussian copula*  $C_{\rho}^{Ga}$  is equivalent to

- the independence copula  $\Pi$  if  $\rho = 0$ ,
- the comonotonicity copula  $M$  if  $\rho = 1$  and
- the countermonotonicity copula  $W$  if  $\rho = -1$

The density of the Gaussian copula is given by

$$c_{\mathbf{P}}^{Ga}(\mathbf{u}) = \frac{1}{\sqrt{\det \mathbf{P}}} \exp\left(-\frac{1}{2} \mathbf{x}' (\mathbf{P}^{-1} - \mathbf{I}_d) \mathbf{x}\right), \quad (3.13)$$

where  $\mathbf{x} = (\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_d))$ .

□

### t-Copula

Consider w.l.o.g.  $\mathbf{X} \sim t_d(\nu, \mathbf{0}, \mathbf{P})$  (multivariate Student's t-distribution) with  $\nu$  Degrees of Freedom (d.o.f.) and  $\mathbf{P}$  a correlation matrix, then the *t-copula (family)* is given by

$$C_{\nu, \mathbf{P}}^t(\mathbf{u}) = t_{\nu, \mathbf{P}}(t_{\nu}^{-1}(u_1), \dots, t_{\nu}^{-1}(u_d)), \quad (3.14)$$

where  $t_{\nu}$  is the CDF of the univariate Student's t-distribution and  $t_{\nu, \mathbf{P}}$  is the CDF of the multivariate Student's t-distribution (both with  $\nu$  d.o.f.).

For the *bivariate t-copula* ( $d = 2$ ), the special cases are the same as for the Gaussian copula except that  $d = 0$  does not yield the independence copula (unless  $\nu \rightarrow \infty$  in which case  $C_{\nu, \rho}^t = C_{\rho}^{Ga}$ ).

The density of  $C_{\nu, \mathbf{P}}^t$  is given by

$$c_{\nu, \mathbf{P}}^t(\mathbf{u}) = \frac{\Gamma((\nu + d)/2)}{\Gamma(\nu/2)\sqrt{\det \mathbf{P}}} \left( \frac{\Gamma(\nu/2)}{\Gamma((\nu + 1)/2)} \right)^d \frac{(1 + \mathbf{x}' \mathbf{P}^{-1} \mathbf{x}/\nu)^{-(\nu+d)/2}}{\prod_{j=1}^d (1 + x_j^2/\nu)^{-(\nu+1)/2}}, \quad (3.15)$$

where  $\mathbf{x} = (t_{\nu}^{-1}(u_1), \dots, t_{\nu}^{-1}(u_d))$ .

□



### 3.2.3 Archimedean Copulas

Unlike implicit copulas, *explicit copulas* can be specified directly by taking into account certain constructional principles. The most important aspects of a such explicit copulas, in particular *archimedean copulas*, are showcased in this subsection. Archimedean copulas are of the general form

$$C(\mathbf{u}) = \phi^{-1}(\phi(u_1) + \cdots + \phi(u_d)), \quad (3.16)$$

where the function  $\phi : [0, 1] \rightarrow [0, \infty)$  is the (*archimedean*) *generator* and satisfies the following properties:

- $\phi$  is strictly decreasing in the entire domain  $[0, 1]$
- We set  $\phi(1) = 0$
- If  $\phi(0) = \lim_{u \rightarrow 0^-} \phi(u) = \infty$ , then  $\phi$  is called *strict*.

Based on Equation 3.16 and according to the form of the generator, we can construct several copula families. Three of the most popular ones are the *Gumbel*, the *Clayton* and the *Frank copula*, which will be discussed.<sup>1</sup> The advantage of such copulas lies in the fact that they interpolate between certain fundamental dependency structures.

#### Clayton Copula

If the generator takes on the form

$$\phi_{Cl}(u) = \frac{1}{\theta} (u^{-\theta} - 1) \quad (3.17)$$

then we obtain the *Clayton copula* given by

$$C_{\theta}^{Cl}(u_1, u_2) = (\max\{u_1^{-\theta} + u_2^{-\theta} - 1, 0\})^{-\frac{1}{\theta}}, \quad (3.18)$$

where  $\theta \in [-1, \infty) \setminus \{0\}$ .

For  $\theta > 0$  the generator of the Clayton copula is strict and we arrive at

$$C_{\theta}^{Cl}(u_1, u_2) = (u_1^{-\theta} + u_2^{-\theta} - 1)^{-\frac{1}{\theta}}. \quad (3.19)$$

Note that for  $\theta = -1$ , we obtain the lower Fréchet-Hoeffding bound, whereas for the limits  $\theta \rightarrow 0$  and  $\theta \rightarrow \infty$  we arrive at the independence copula and the comonotonicity copula respectively.

<sup>1</sup>We will look into these copulas for the bivariate case ( $d = 2$ ) only.

□

### Gumbel Copula

If the generator takes on the form

$$\phi_{Gu}(u) = (-\ln u)^\theta, \quad \theta \in [1, \infty), \quad (3.20)$$

then we arrive at the *Gumbel copula* given by

$$C_\theta^{Gu}(u_1, u_2) = \exp \left[ - \left( (-\ln u_1)^\theta + (-\ln u_2)^\theta \right)^{\frac{1}{\theta}} \right]. \quad (3.21)$$

Note that for  $\theta = 1$ , we obtain the independence copula, while for  $\theta \rightarrow \infty$  the Gumbel copula converges to the comonotonicity copula.

□

### Frank Copula

If the generator takes on the form

$$\phi_{Fr}(u) = \ln(e^{-\theta} - 1) - \ln(e^{-\theta u} - 1), \quad \theta \in \mathbb{R} \setminus \{0\}, \quad (3.22)$$

we obtain the *Frank copula* given by

$$C_\theta^{Fr}(u_1, u_2) = -\frac{1}{\theta} \ln \left( 1 + \frac{(e^{-\theta u_1} - 1) \cdot (e^{-\theta u_2} - 1)}{e^{-\theta} - 1} \right) \quad (3.23)$$

□

MAYBE MORE ON THESE WITH SOME PRETTY PLOTS

## 3.3 Dependence Measures

*Dependence measures* allow us to summarize a particular kind of dependence into a single number.<sup>2</sup> Recall the Fréchet-Hoeffding bounds (Equation 3.8 and Equation 3.9). They are an example of such kind of dependence measures. After all, they represent perfect negative or positive dependence. In this section, we will take a closer look into three classes of dependence measures along with appropriate association metrics.

---

<sup>2</sup>In the bivariate case

### 3.3.1 Linear Correlation

Undoubtedly, the most famous association metric for two RVs  $X_1$  and  $X_2$  is the *Linear or Pearson's correlation coefficient*

$$\rho(X_1, X_2) = \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1)}\sqrt{\text{Var}(X_2)}} \in [-1, 1]. \quad (3.24)$$

Note that  $E(X_1) < \infty$  and  $E(X_2) < \infty$  have to hold, i.e. the first two moments have to exist for  $\rho$  to be defined.

The Pearson correlation coefficient is interpretable for RVs which have (approximately) a linear relationship, where  $\rho = -1$  indicates perfect negative linear correlation,  $\rho = 1$  indicates perfect positive linear correlation and  $\rho = 0$  indicates no correlation between  $X_1$  and  $X_2$ . However, comprehensibility of this measure comes along with some drawbacks:

- A correlation of 0 is in general not equivalent to independence. This property holds only for normally distributed RVs.<sup>3</sup>
- $\rho$  is invariant only under linear transformations, but not under transformations in general.
- Given the marginals and correlation  $\rho$ , one is able to construct a joint distribution only for the class of elliptical distributions. (MAYBE PLOT qrm)
- Given the marginals, only for elliptically distributed RVs any  $\rho \in [-1, 1]$  is attainable.

### 3.3.2 Rank Correlation

To compensate some of the drawbacks of linear correlation, we take advantage of correlation measures based on the ranks of data. *Rank correlation coefficients*, like the ones presented below, are always defined and obey to the invariance principal. This means that these coefficients only depend on the underlying copula and they can thereof be directly derive.

#### Spearman's Rho

Consider two RVs  $X_1$  and  $X_2$  with continuous CDFs  $F_1$  and  $F_2$ , then the

---

<sup>3</sup>e.g.  $X_2 = X_1^2$  implies perfect dependence, yet  $\rho(X_1, X_2) = 0$ . Conversely though, independence always yields  $\rho = 0$ .

*Spearman's rho correlation coefficient* is simply the linear correlation between the CDFs

$$\rho_S = \rho(F_1(X_1), F_2(X_2)). \quad (3.25)$$

The reason being is that by applying the CDF to data, naturally a multiple of the ranks of the data are obtained, which essentially is equivalent to

$$\rho_S = \rho(\text{Ran}(X_1), \text{Ran}(X_2)) \quad (3.26)$$

Due to the invariance principle, we also obtain Spearman's rho directly from the unique copula via

$$\rho_S = 12 \int_0^1 \int_0^1 C(u_1, u_2) du_1 du_2 - 3. \quad (3.27)$$

□

### Kendall's Tau

Let  $X_1 \sim F_1$  and  $X_2 \sim F_2$  be two RV and let  $(\tilde{X}_1, \tilde{X}_2)$  be an independent copy<sup>4</sup> of  $(X_1, X_2)$ . Then *Kendall's tau* is defined by

$$\begin{aligned} \rho_\tau &= E[\text{sign}((X_1 - X'_1)(X_2 - X'_2))] \\ &= P((X_1 - X'_1)(X_2 - X'_2) > 0) - P((X_1 - X'_1)(X_2 - X'_2) < 0). \end{aligned} \quad (3.28)$$

Similarly to Spearman's rho, using the invariance principal, we can directly derive Kendall's tau from the unique copula by

$$\rho_\tau(X_1, X_2) = 4 \int_0^1 \int_0^1 C(u_1, u_2) dC(u_1, u_2) - 1. \quad (3.29)$$

□

Both  $\rho_S, \rho_\tau \in [-1, 1]$ , where any value within this interval is attainable in contrast to the Pearson coefficient. If any of these rank correlations is  $-1$  (or  $1$ ), we are in the countermonotonic (or comonotonic) case. If  $\rho_S$  (or  $\rho_{\text{tau}}$ )  $= 0$ , this does not necessarily imply independence between  $X_1$  and  $X_2$ , although the opposite direction holds. Furthermore, they are not limited to be invariant just under linear transformations.

---

<sup>4</sup>An independent copy  $\tilde{X}$  of a RV  $X$  is a RV that inherits from the same distribution as  $X$  and is independent of  $X$ .

### 3.3.3 Tail Dependence

*Coefficients of tail dependence* express the strength of the dependency in the extremes of distributions, i.e. the joint tails. We distinguish between *lower* and *upper tail dependence* between  $X_j \sim F_j, j = 1, 2$  and provided that the below limits exist, they are given by

$$\lambda_l = \lim_{q \rightarrow 0^+} P(X_2 \leq F_2^{\leftarrow}(q) | X_1 \leq F_1^{\leftarrow}(q)) \quad (3.30)$$

and

$$\lambda_u = \lim_{q \rightarrow 1^-} P(X_2 > F_2^{\leftarrow}(q) | X_1 > F_1^{\leftarrow}(q)). \quad (3.31)$$

If  $\lambda_l$  (or  $\lambda_u$ ) = 0, then we say that  $X_1$  and  $X_2$  are *asymptotically independent* in the lower (or upper) tail,<sup>5</sup> otherwise we have lower (or upper) tail dependence.

For continuous CDFs and by using Bayes' theorem, these expressions can be re-written to

$$\begin{aligned} \lambda_l &= \lim_{q \rightarrow 0^+} \frac{P(X_2 \leq F_2^{\leftarrow}(q), X_1 \leq F_1^{\leftarrow}(q))}{P(X_1 \leq F_1^*(q))} \\ &= \lim_{q \rightarrow 0^+} \frac{C(q, q)}{q} \end{aligned}$$

and similarly

$$\lambda_u = 2 - \lim_{q \rightarrow 1^-} \frac{1 - C(q, q)}{1 - q}.$$

Therefore, tail dependencies can be assessed by means of the copula itself when approaching the points (0, 0) and (1, 1). In addition, for all radially symmetric copulas (e.g. the bivariate Gaussian or the t-copula) we have  $\lambda_l = \lambda_u = \lambda$ .

Some examples are:

- Clayton:  $\lambda_l = 2^{-1/\theta}, \lambda_u = 0$  (only lower tail dependence)
- Gumbel:  $\lambda_l = 0, \lambda_u = 2 - 2^{1/\theta}$  (only upper tail dependence)
- Frank:  $\lambda_l = 0, \lambda_u = 0$  (no tail dependence)

Following such guidelines, the choice of a practicable copula can be facilitated.

## 3.4 Conditional Copulas

## 3.5 Vine Copulas

---

<sup>5</sup>Not necessarily true for the other way around



## 4 Data Exploration





## 5 Modelling



## 6 Conclusion



# Appendix

Include appendix here...



## List of Figures





**List of Tables**

1.1 Transactional raw data description from online purchases of  
western European countries . . . . . 1

1.2 Article master data . . . . . 2



## List of Abbreviations

**BIC** Bayesian Information Criterion

**GLM** Generalized Linear Model

**LM** Linear Regression Model

**LMM** Linear Mixed Model

**GLMM** Generalized Linear Mixed Model

**CDF** Cumulative Distribution Function

**PDF** Probability Density Function

**RV** Random Variable

**w.l.o.g.** Without loss of generality

**d.o.f.** Degrees of Freedom

**GAM** Generalized Additive Model



## References

- [Embrechts et al. 2001] Embrechts, Paul ; Lindskog, Filip ; McNeil, Alexander: Modelling dependence with copulas. In: *Rapport technique, Département de mathématiques, Institut Fédéral de Technologie de Zurich, Zurich* 14 (2001)
- [Fahrmeir et al. 2003] Fahrmeir, L. ; Kneib, T. ; Lang, S. ; Marx, B.: *Regression; Models, Methods and Applications. 2013.* 2003
- [Hofert et al. 2019] Hofert, Marius ; Kojadinovic, Ivan ; Mächler, Martin ; Yan, Jun: *Elements of copula modeling with R.* Springer, 2019
- [Klein and Kneib 2016] Klein, Nadja ; Kneib, Thomas: Simultaneous inference in structured additive conditional copula regression models: a unifying Bayesian approach. In: *Statistics and Computing* 26 (2016), Nr. 4, S. 841–860
- [Lütkepohl 2005] Lütkepohl, Helmut: *New introduction to multiple time series analysis.* Springer Science & Business Media, 2005
- [McNeil et al. 2015] McNeil, Alexander J. ; Frey, Rüdiger ; Embrechts, Paul: *Quantitative risk management: concepts, techniques and tools-revised edition.* Princeton university press, 2015
- [Ruppert and Matteson 2015] Ruppert, David ; Matteson, David S.: *Copulas.* S. 183–215. In: *Statistics and Data Analysis for Financial Engineering: with R examples.* New York, NY : Springer New York, 2015. – URL [https://doi.org/10.1007/978-1-4939-2614-5\\_8](https://doi.org/10.1007/978-1-4939-2614-5_8). – ISBN 978-1-4939-2614-5
- [Schmidt 2007] Schmidt, Thorsten: Coping with copulas. In: *Copulas-From theory to application in finance* (2007), S. 3–34
- [Sklar 1959] Sklar, M.: Fonctions de repartition an dimensions et leurs marges. In: *Publ. inst. statist. univ. Paris* 8 (1959), S. 229–231
- [Vatter and Chavez-Demoulin 2015] Vatter, Thibault ; Chavez-Demoulin, Valérie: Generalized additive models for conditional dependence structures. In: *Journal of Multivariate Analysis* 141 (2015), S. 147–167

- [Vatter and Nagler 2018] Vatter, Thibault ; Nagler, Thomas: Generalized additive models for pair-copula constructions. In: *Journal of Computational and Graphical Statistics* 27 (2018), Nr. 4, S. 715–727
- [Vatter and Nagler 2019] Vatter, Thibault ; Nagler, Thomas: gamCopula-package: Generalized Additive Models for Bivariate Conditional... (2019)
- [Wood 2017] Wood, Simon N.: *Generalized additive models: an introduction with R*. CRC press, 2017