



Malte Bonart

Testing for Structural Breaks in Factor Copula Models - Implementation and Application in Social Media Topic Analysis

Master thesis

Supervisor: Prof.Dr. Dominik Wied

Submitted for the Master Examination in Economics at
the Faculty of Management, Economics and Social Sci-
ences of the University of Cologne in June 2018.

This is my abstract.

Contents

List of Figures	iii
List of Tables	iii
List of Abbreviations	iv
1. Introduction	1
2. Theoretical foundation	2
2.1. Copula theory	2
2.2. Copula models for multivariate time series	4
2.3. Factor copulas	5
2.4. Simulated methods of moments estimation for factor copulas	8
2.5. Structural break test for factor copulas	9
3. <i>factorcopula</i> - an R package for simulation and estimation of factor copulas	10
3.1. Overview and usage	11
3.2. Simulation study	13
4. Modelling topic dependencies over time with factor copulas	17
4.1. The <i>btw17</i> social media dataset	17
4.2. Data processing and descriptive analysis	19
4.3. Results	20
5. Discussion	24
A. Appendix	25
B. References	28
C. Statutory Declaration	30

List of Figures

2.1. Illustration of different equidependence factor copula models with $N = 2$, $\beta = 1.5$, standard normal distributed marginals and different distributions for the latent variable and the error term.	7
3.1. Approximated density of $\hat{\theta}$ for an equidependence skew t - t factor copula model with $\beta = 1.5$, $S = 25000$ and standard normal distributed marginals. Each simulation is based on 500 Monte-Carlo replications.	13
3.2. Structural break test for a bloc-equidependence model with $N = 21$ and 3 groups of equal size. The theoretical breakpoint is at $t = 1000$ and is modelled as a change of the intra- and interdependenvy of the first group from 0 to 1.5.	16
4.1. Number of active accounts and number of posts per party and month.	19
4.2. Relative weighted monthly frequency of posts matching the regular expression <code>flucht fluecht</code> . The vertical line indicates the breakpoint detected by the moments based test.	22
A.1. Scree-plot of ranked eigenvalues based on the pairwise rank-correlation matrix.	25
A.2.	26
A.3. Pairwise scatterplot of the estimated residuals.	27

List of Tables

4.1. Number of posts, accounts, likes and shares over the observation period from "2014-01-01" - "2017-12-31"	18
4.2. Estimation results for different one-factor copula specifications.	23
4.3. Estimation results for different one-factor copula specifications before and after the breakpoint detected by the moments based test.	24

List of Abbreviations

cdf	Cumulative distribution function
iid	Identically and independent distributed
SMM	Simulated methods of moments
DGP	Data generating process
btw17	Bundestag election 2017
API	Application programming interface

1. Introduction

Models with copula functions became increasingly popular since the 1990th (Nelsen 1999, p. 1). The concept was first introduced in the work by Sklar (1959). They are mainly used in two ways: First, to model the dependence structure of multivariate distributions independent of their underlying marginal distributions and second, to construct bivariate or multivariate distributions (Sempi 2011, p. 302). This paper focuses on the first application.

Sklar's theorem can be used to construct multivariate models first, by specifying the marginal distributions of the random variables involved and second, by specifying the dependence structure among the variables via a copula function Sempi 2011. By doing so, one allows for non-parametrized or semi-parametrized estimation of the marginal distributions together with a parametrized copula. High dimensional problems become traceable since the number of parameters can be drastically reduced (Patton 2009, p. 777).

For time series data copula theory can be used in two ways: First, to describe the cross sectional dependence structure by estimating the conditional copula function of the conditional joint distribution $F(\mathbf{y}|\mathcal{F}_{t-1})$ with $\mathbf{Y}_t = [Y_{1t}, \dots, Y_{nt}]'$ and past information \mathcal{F}_{t-1} . To obtain a valid distribution, the information set must be the same for both the copula and the marginal distributions (Patton 2009, p. 771).

Second, copulas can be used to describe the dependence between observations of a univariate time series $[Y_t, Y_{t+1}, \dots, Y_{t+n}]'$. This is related to the study of Markov processes. (Patton 2009, p. 774 ff). This paper focuses on the first application.

Applications for copula modeling can be found in various disciplines but they became increasingly popular in the field of finance, actuarial science and hydrology Sempi 2011.

Correlation or covariance matrices can be used to model linear dependence especially for multivariate normal or t-distributions. But they lack the ability to model the dependence e.g. in the presence of heavy tails or outliers (Kumar 2011).

2. Theoretical foundation

Rank correlation matrices such as *Spearman's rho* are invariant under monotonic transformations but they are not moment-based.

Research question: What do I want to analyze?

How similar is the dependency structure of the political communication on social media channels compared to financial markets?

Can we detect structural breaks in the dependency structure of the political communication on social media channels? This could be an indicator of political eruptions such as elections, scandals or political events.

extreme dependence during economic crisis (-> elections)

Relevance: Why do I ask this question? Why is it relevant?

The thesis is structured in four main chapters: The first chapter lays the theoretical foundation by summarizing important aspects of copula theory and by presenting the factor copula approach, its estimation strategy via simulated methods of moments and a suitable test for time varying dependence structures. The second chapter presents implementation details of the software package *factorcopula*, written in the statistical programming language R (Bonart 2018). With the package, factor copulas can be fitted to real data and structural breaks can be detected. The validity of the package and the methods is illustrated by a small simulation study. In the last chapter, the methods are applied to a large dataset of textual social media posts from german politicians and political parties. Here, the goal is to identify temporal dependencies between different topics and to test for changes in the dependence structure due to important political events. The last chapter summarizes the findings and critically discusses the presented methods.

2. Theoretical foundation

2.1. Copula theory

A function of the type $C : [0, 1]^N \rightarrow [0, 1]$, with $N \geq 2$ is called a *copula* if

2. Theoretical foundation

1. $C(u_1, \dots, u_N) = 0$, if $\exists i \in \{1, \dots, N\} : u_i = 0$
2. $C(1, \dots, 1, x_i, 1, \dots, 1) = x_i$
3. The C -volume of every N -box is postive

The last property is also called the N -increasing property (Sempi 2011, p. 302). From this definition it follows that, in a statistical sense, a copula function is a multivariate distribution $C(u_1, \dots, u_N) = P(U_1 \leq u_1, \dots, U_N \leq u_N)$ with uniform marginals $U_i \sim U(0, 1) \forall i \in \{1, \dots, N\}$ (Joe 2015, p. 7).

Sklar (1959) showed, that every d-variate distribution $F(x_1, \dots, x_n)$ can be expressed in terms of its marginal distributions $F_1(x), \dots, F_n(x)$ and a copula function $C(u_1, \dots, u_N)$ such that $F(\mathbf{x}) = C(F_1(x_1), \dots, F_N(x_N))$.

If F is continuous with marginal quantile functions $F_1^{-1}, \dots, F_N^{-1}$ then $C(\mathbf{u})$ is uniquely determined by $C(\mathbf{u}) = F(F_1^{-1}(u_1), \dots, F_n^{-1}(u_N))$.

Some bivariate measures of dependency share the property of *scale invariance*. Thus, the measures are invariant with respect to the marginal distributions and can therefore be expressed as a function of their copula (Schmid et al. 2010, p. 210). Two widely used measures are Spearman's

$$\rho_{X_1, X_2} = 12 \int \int_{[0,1]^2} u_1 u_2 dC(u_1, u_2) - 3 \quad (2.1)$$

and Kendal's rank correlation

$$\tau_{X_1, X_2} = 4 \int \int_{[0,1]^2} C(u_1, u_2) dC(u_1, u_2) - 1 = 4E(C(U_1, U_2)) - 1. \quad (2.2)$$

Bivariate measures can be extended to the multivariate case where there measure the strength of association embedded in the N -dimensional copula of a multivariate vector X .

One can also use the average over all bivariate dependency measures.

Multivariate versions of Spearman's rank correlation can be defined in terms of an underlying copula function as ((Schmid et al. 2010, p. 215ff)):

2. Theoretical foundation

$$\rho = \frac{N+1}{2^N - (N+1)} (2^N \int_{[0,1]^N} C(\mathbf{u}) d\mathbf{u} - 1) \quad (2.3)$$

Analogously, Kendall's rank correlation is defined as

$$\tau = \frac{1}{2^{N-1} - 1} (2^N \int_{[0,1]^N} C(\mathbf{u}) d\mathbf{u} - 1). \quad (2.4)$$

Lower and upper tail dependency for two variables X and Y is defined as

$$\begin{aligned} \tau_{XY}^L &= \lim_{q \rightarrow 0} \frac{P(X \leq F_X^{-1}(q), Y \leq F_Y^{-1}(q))}{q} \\ \tau_{XY}^U &= \lim_{q \rightarrow 1} \frac{P(X > F_X^{-1}(q), Y > F_Y^{-1}(q))}{q} \end{aligned} \quad (2.5)$$

2.2. Copula models for multivariate time series

Conditional copula as presented in (Patton 2006) and (Patton 2009, p. 772)

For this work we use a semiparametric copula-based multivariate dynamic model as described in Chen and Fan (2006, p. 129 ff). The goal is to model the conditional multivariate distribution of $\mathbf{Y}_t | \mathcal{F}_{t-1}$, where the σ -algebra \mathcal{F}_{t-1} possibly contains past information and information from other exogenous variables $\{\mathbf{Y}'_{t-1}, \mathbf{Y}'_{t-2}, \dots, \mathbf{X}'_t, \mathbf{X}'_{t-1}, \dots\}$. The conditional means and variances of $\mathbf{Y}_t | \mathcal{F}_{t-1}$ are estimated parametrically. The observations are then filtered by removing serial dependence or volatility clustering such that the leftover standardized innovations are independent of past information. Finally, the innovations are modeled using a parametric copula and nonparametric rank based estimates of the marginal distributions.

If we denote the parametrized conditional mean of a single variable as $\mu_{it} = E(Y_{it} | \mathcal{F}_{t-1}; \phi)$ and the parametrized conditional standard deviation as $\sigma_{it} = \sqrt{V(Y_{it} | \mathcal{F}_{t-1}; \phi)}$ we can write the multivariate time series as:

2. Theoretical foundation

$$\mathbf{Y}_t = \boldsymbol{\mu}_t + \boldsymbol{\sigma}_t \boldsymbol{\eta}_t, \quad (2.6)$$

with $\boldsymbol{\sigma}_t = \text{diag}(\sigma_{1t}, \dots, \sigma_{Nt})$. The innovations $\boldsymbol{\eta}_t = (\eta_{1t}, \dots, \eta_{Nt})'$ are independent of past information and iid distributed according to some multivariate distribution function $F_\eta(x_1, \dots, x_N)$.

The cdf of the innovations can be expressed in terms of a copula and the marginal distributions such that $F_\eta(x_1, \dots, x_N) = C(F_{\eta_1}(x_1), \dots, F_{\eta_N}(x_N); \boldsymbol{\theta})$.

2.3. Factor copulas

(Oh and Patton 2017)

Factor copulas are a family of copulas for which the copula function $C(u_1, \dots, u_N)$ is based on a latent factor structure as defined in Oh and Patton (2017, p. 140 ff).

Consider a set of artificial variables $X_i = \sum_{k=1}^K \beta_{ik} Z_k + \epsilon_i$ with $i = 1, \dots, N$ the dimension of the observable data $\mathbf{Y} = (Y_1, \dots, Y_N)'$ and $k = 1, \dots, K$ the number of latent variables. The latent variables Z_k and the error term ϵ_i follow some parametrized distributions such that $\epsilon_i \stackrel{iid}{\sim} F_\epsilon(\gamma_\epsilon)$ and $Z_k \sim F_{Z_k}(\gamma_{Z_k})$ with $Z_i \perp Z_j \forall i \neq j$, $Z_k \perp \epsilon_i \forall i, k$ and $\gamma_\epsilon, \gamma_{Z_k}$ some distribution specific parameter vectors.

The joint probability function $F_X(x_1, \dots, x_N)$ of the artificial variables can then be expressed in terms of its marginal distributions $F_{X_i}(x)$ and a copula function $C_\theta(u_1, \dots, u_N)$ such that $F_X(x_1, \dots, x_N) = C_\theta(F_{X_1}(x_1), \dots, F_{X_N}(x_N); \boldsymbol{\theta})$.

The factor copula is therefore completely defined via the parameter vector $\boldsymbol{\theta} = (\beta_{11}, \dots, \beta_{i1}, \dots, \beta_{ik}, \gamma'_{Z_1}, \dots, \gamma'_{Z_K}, \gamma'_\epsilon)'$. The number of latent variables K and the distribution functions $F_{Z_1}, \dots, F_{Z_K}, F_\epsilon$ are hyperparameters of the model which have to be chosen prior to the estimation.¹

¹Oh and Patton (2017, p. 143ff) provide a heuristic of finding the number of latent variables by analyzing so called *scree-plots*: Ordered eigenvalues from the sample rank-correlation matrix of the data.

2. Theoretical foundation

The latent factor structure is linked to the observable data via the copula function because it holds that $F_Y(y_1, \dots, y_n) = C_\theta(F_{Y_1}(y_1), \dots, F_{Y_N}(y_N))$. The model can be summarized in the following set of equations:

$$\begin{aligned} \mathbf{Y} &= (Y_1, \dots, Y_N)' \\ F_Y &= C_f(F_{Y_1}(y_1), \dots, F_{Y_N}(y_N); \boldsymbol{\theta}) \\ \mathbf{X} &= (X_1, \dots, X_N)' = \boldsymbol{\beta}\mathbf{Z} + \boldsymbol{\epsilon} \\ F_X &= C_f(F_{X_1}(x_1), \dots, F_{X_N}(x_N); \boldsymbol{\theta}) \end{aligned} \tag{2.7}$$

It is important to note, that the artificial variables X are only used for the construction of the factor copula function $C_f(u_1, \dots, u_N)$. Once this copula function is determined, the artificial variables and its marginal distributions $F_{X_i}(x)$ are of no interest. Using the copula function together with the marginal distributions of the observable variables $F_{Y_i}(y)$ one can then determine the joint distribution of Y .

This approach allows for a two-stage estimation in which first the marginal distributions are estimated flexibly and second the factor structure for the possibly high dimensional copula function is fitted to the data. For the factor copula and the joint distribution of the artificial variables as defined in (2.7) a closed form usually does not exist. Therefore, one has to rely on simulation methods as described in section 2.4.

A lower bound for the number of parameters $P = |\boldsymbol{\theta}|$ to be estimated is given by the size of the factor matrix $\boldsymbol{\beta}$ which is $|\boldsymbol{\beta}| = N \times K$. To reduce the number of parameters Oh and Patton (2017) present two restrictions on $\boldsymbol{\beta}$: the *equidependence* and the *block-equidependence* model.

For the first model it is assumed that $K = 1$ and $\boldsymbol{\beta} = (\beta, \dots, \beta)'$. Thus, the model consists of a single latent factor and a single factor loading β which is the same for all variables. This implies that each pairwise dependency is the same for all observable variables.

2. Theoretical foundation

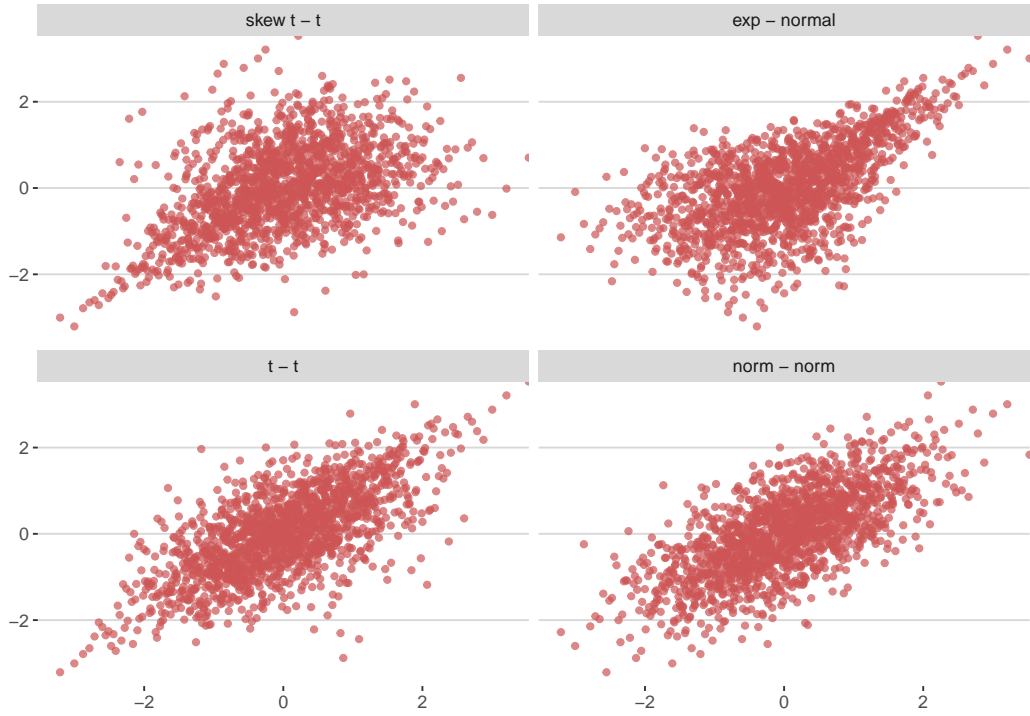


Figure 2.1: Illustration of different equidependence factor copula models with $N = 2$, $\beta = 1.5$, standard normal distributed marginals and different distributions for the latent variable and the error term.

Figure 2.1 shows four different simulations from a one factor equidependence factor copula model. The marginal distributions are standard normal while the distributions of the latent variable and the error term differ.

The block-equidependence model is less restrictive and is especially suitable for variables which can be naturally partitioned into different groups.² The model assumes a common factor for all groups and a group specific factor for each group. Thus, each variable is only affected by two factors. For the factor matrix, it is further assumed that all variables in the same group have the same factor loading while variables in different groups can have different loadings. This implies that the pairwise intra-group dependencies are equal while the pairwise inter-group dependencies can vary between the groups.

Formally, consider a partition of $X = (X_1, \dots, X_N)'$ into D groups X_j^i ,

²E.g. this could be stock market prices grouped into different industry sectors.

2. Theoretical foundation

where $i = 1, \dots, D, j = 1, \dots, s_i$ and s_i the number of variables in group i . Then the model can be summarized as:

$$\begin{aligned} \mathbf{X} &= (X_1^1, \dots, X_{s_1}^1, \dots, X_1^D, \dots, X_{s_D}^D)' = \boldsymbol{\beta} \mathbf{Z} + \boldsymbol{\epsilon} \\ \mathbf{Z} &= (Z_0, Z_1, \dots, Z_D)' \\ \boldsymbol{\beta} &= \begin{pmatrix} \beta^1 & \beta^{D+1} & 0 & \dots & 0 \\ \beta^1 & \beta^{D+1} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \beta^1 & \beta^{D+1} & 0 & \dots & 0 \\ \beta^2 & 0 & \beta^{D+2} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \beta^D & 0 & 0 & \dots & \beta^{D+D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \beta^D & 0 & 0 & \dots & \beta^{D+D} \end{pmatrix}, \end{aligned} \quad (2.8)$$

where $\boldsymbol{\beta}$ is of size $N \times (D + 1)$, $N = \sum_{i=1}^D s_i$ but with only $2D$ different factor loadings.

2.4. Simulated methods of moments estimation for factor copulas

Estimation methods used for copula models depends on the degree of parametrization: For fully parametrized models for the copula and the marginal distributions maximum likelihood or multi-stage maximum likelihood is used. But one can also non parametrically estimate the marginal distributions and combine them with a parametric copula. In this case, pseudo-maximum likelihood is used. If a closed form functional relation of spearman's rho or kendall's thau to the copula parameters is available, one can also solve the system directly by using a method of moments approach. Here, the population based statistics are replaced by their sample counterparts (inversion method).

For the factor copula model a closed form one to one mapping of the copula's parameters θ to measures of dependency as defined in (2.2)

2. Theoretical foundation

- (2.5) is not available in general. If it was available, methods of moments or generalized methods of moments (if the number of moment conditions is larger than the number of parameters) could be applied (Oh and Patton 2013, p. 689f).

Instead one can use a set of scale-invariant empirical dependence measures calculated with simulations from the artificial variables X and compare them to the dependence measures obtained from the observable data Y . Minimizing the weighted squared difference of the two dependency vectors yields an estimator for θ .

Formally, the estimator is given by

$$\hat{\theta} = \arg \min Q(\theta) = \arg \min g(\theta)' \hat{W} g(\theta) \quad (2.9)$$

with

$$g(\theta) = \hat{m} - \tilde{m}(\theta), \quad (2.10)$$

where \hat{m} and \tilde{m} are the vector of dependencies computed with the observable and the simulated data respectively.

Oh and Patton (2013, p. 691ff) showed that under a set of assumptions, the SMM is weakly consistent and asymptotically normal distributed:

$$\frac{1}{\sqrt{1/T + 1/S}}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, \Sigma_0) \text{ for } T, S \rightarrow \infty \quad (2.11)$$

, with covariance matrix $\Sigma_0 = 12$

2.5. Structural break test for factor copulas

Note that we assume that the functional form of the copula is time invariant while the copula's parameters can vary over time (Patton 2006, p. 542).

The model presented in 2.2 allows for a wide variety of parametrization and copula functions. Here, we focus on the factor copula model

3. *factorcopula* - an R package for simulation and estimation of factor copulas

and the SMM estimation procedure as presented in the previous sections.

$$H_0 : \theta_1 = \theta_2 = \dots = \theta_T \quad H_1 : \theta_t \neq \theta_{t+1} \text{ for some } t = \{1, \dots, T\}$$

The test statistics is

$$P = \sup_{s \in [\epsilon, 1]} s^2 T (\boldsymbol{\theta}_{sT, S} - \boldsymbol{\theta}_{T, S})' (\boldsymbol{\theta}_{sT, S} - \boldsymbol{\theta}_{T, S}) \quad (2.12)$$

Under the null hypothesis and given some assumptions the test statistics converges in distribution to

$$P \xrightarrow{d} \sup_{s \in [\epsilon, 1]} (\mathbf{A}^*(s) - s\mathbf{A}^*(1))' (\mathbf{A}^*(s) - s\mathbf{A}^*(1)), \quad (2.13)$$

with $\mathbf{A}^*(s) = (G'WG)^{-1}G'W(A(s) - \frac{s}{\sqrt{k}}A(1))$ and $T, S \rightarrow \infty, \frac{S}{T} \rightarrow k$ or $\frac{S}{T} \rightarrow \infty$.

In the following a structural break test is presented (Manner, Stark, and Wied 2017).

(Patton 2009)

(Oh and Patton 2013)

(Manner, Stark, and Wied 2017)

3. *factorcopula* - an R package for simulation and estimation of factor copulas

In this chapter we present the implementation of the previous methods, using the programming language R. The functions are bundled in an R-package such that the methods can be easily installed and distributed. The validity of the package is tested in a simulation study.

3. *factorcopula* - an R package for simulation and estimation of factor copulas

3.1. Overview and usage

The package consists of a set of high level functions which can be used to construct, simulate and fit various factor copula models. The definition of the factor copula model is handled by the functions `config_factor`, `config_error` and `config_beta` which define the distribution of the latent variables Z , the error term ϵ and the matrix of factor loadings β .

For the specification of the distributions, the function name of any available random number generator can be used. Additional arguments for the function can be declared in a named list. Distributional parameters can either be pre-defined or passed as expressions which are not immediately evaluated. To distinguish model parameters from pre-defined distributional parameters, an additional character vector with the name of the model parameters has to be passed to `config_factor` or `config_error`.

In the following example, a factor copula with two latent variables is defined. The first factor is skewed-t distributed, the second is normal distributed. The error term is t distributed with 4 degrees of freedom. The distributional parameters of the skewed-t distribution are not pre-defined and will be included in the estimation process.

```
library(factorcopula)
Z <- config_factor(rst = list(nu = 1/dfInv, lambda = lambda),
                  rnorm = list(),
                  par = c("dfInv", "lambda"))
eps <- config_error(rt = list(df = 4))
```

For specifying the matrix of factor loadings one can either manually construct a character matrix of parameters or use the function `config_beta` which accepts an input vector k and the number of latent variables.

The vector k defines the group for each variable $i = 1, \dots, N$. Thus, an equidependence model can be specified with $k = (1, 1, \dots, 1)$, an unrestrictive model with $k = (1, \dots, N)$ and a bloc-equidependence model with $k = (1, 1, \dots, 2, 2, \dots, M, M, \dots)$, where M is the number of groups.

3. *factorcopula* - an R package for simulation and estimation of factor copulas

In this example an equidependence model, with one β -parameter for each latent variable is constructed:

```
k <- c(1, 1, 1)
beta <- config_beta(k, Z = 2)
```

The two functions `fc_create` and `fc_fit` are used to either simulate values from a factor copula or to fit a model to a real dataset.

The function `fc_create` returns itself a function which can be used to simulate values from the copula model. The simulation function accepts a *named* vector θ of parameters, the number of simulations S and an optional random seed. If a random seed is provided the functions always uses the same underlying random numbers.

The simulation of new values is a time consuming part. To avoid unnecessary calls to the random number generators, the function remembers the state of the Z and ϵ matrix over the function calls and only updates the values if necessary. Thus if neither the seed nor the distributional parameters change, the function directly proceeds to the calculation of $\beta Z + \epsilon$.

During optimization, the random seed is kept fixed to avoid numerical instabilities (app. to Oh and Patton 2013, p. 12f). Therefore, by using the described memory functionality one can greatly save computational costs. This improves the performance of the optimization process massively, especially when only beta parameters are about to be optimized. In this case, the random number generators are only called once.

```
copfun <- fc_create(Z, eps, beta)
theta <- c(beta1 = 2, beta2 = 1, dfInv = 0.25, lambda = -0.8)
X <- copfun(theta, S = 10)
```

As its core, the function `fc_fit` uses the `sbplx` function from the `nloptr` package (Johnson n.d.). The function is a re-implementation of the Subplex algorithm by (Rowan 1990) with support for bound constraints. It is based on the well known Nelder-Mead simplex method and solves the objective function on subspaces of the input space. For parallel computations, the function `fc_fit` accepts a cluster object created by the `snow` or `parallel` package. The function uses several

3. *factorcopula* - an R package for simulation and estimation of factor copulas

randomly selected starting values. The number of different starting values can also be specified with the `trials` argument.

3.2. Simulation study

To illustrate the discussed methods and the validity and performance of the package two simulation studies were performed. First, an equidependence factor copula model with varying dimensions was estimated repeatedly to show the consistency of the SMM procedure. Second, both the moments and copula based structural break test was performed for a bloc-equidependence factor copula model.

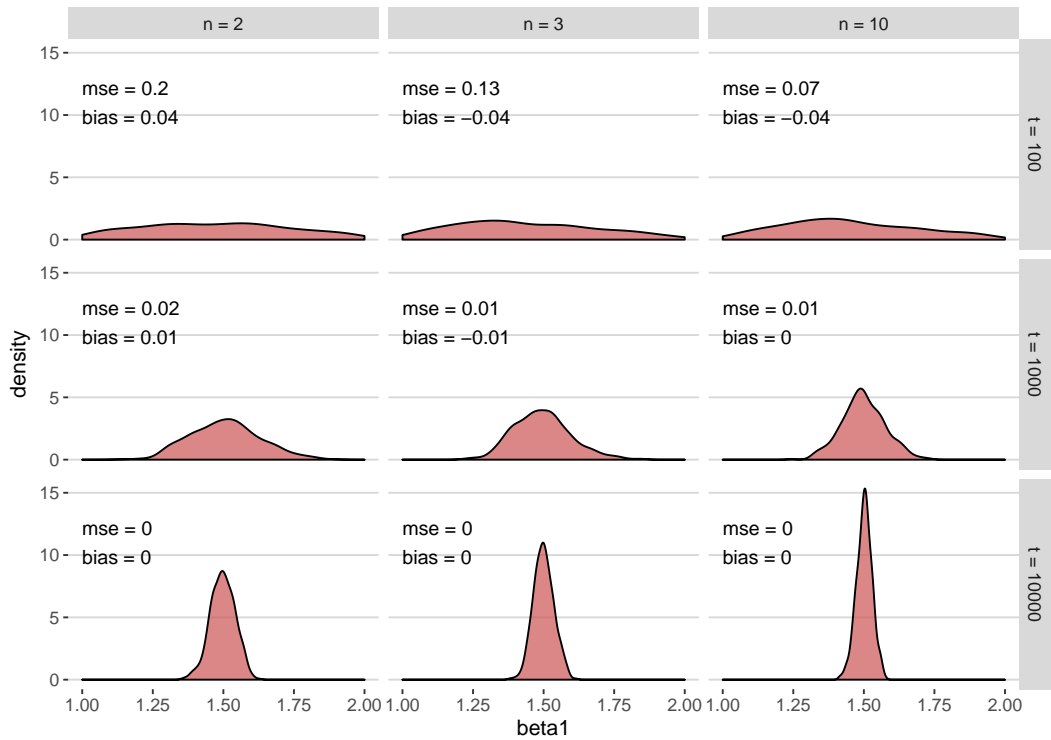


Figure 3.1: Approximated density of $\hat{\theta}$ for an equidependence skew t - t factor copula model with $\beta = 1.5$, $S = 25000$ and standard normal distributed marginals. Each simulation is based on 500 Monte-Carlo replications.

Figure 3.1 shows the results for the first study. The DGP was based on a simple equidependence model with one skew- t distributed latent variable and a single factor loading $\beta = 1.5$. The error term is t -distributed

3. *factorcopula* - an R package for simulation and estimation of factor copulas

and the marginal distribution of the observed values is iid standard normal distributed:

$$\begin{aligned} \mathbf{Y} &= (Y_1, Y_2)' \sim F_{\mathbf{Y}} = C(F_{Y_1}, F_{Y_2}) \\ (X_1, X_2)' &= (\beta, \beta)'Z + \epsilon \\ Y_1, Y_2 &\sim N(0, 1), Z \sim \text{skew} - t(4, -0.8), \epsilon \sim t(4). \end{aligned} \tag{3.1}$$

For all variations of N and T , $S = 25000$ was chosen. Each simulation was repeated 500 times. The bias and mean squared error of $\hat{\beta}$ was approximated by using the empirical average and standard deviation of all 500 simulations.

One can clearly see that for larger T the SMM estimator converges to the true parameter. For $t = 10000$ the bias and mean squared error is virtually zero. For larger N one can also note a drop in the mean squared error.

For the second study, a more sophisticated model was chosen to illustrate the effectiveness of the approach even for high dimensional problems and complicated dependence structures. Analogous to the empirical examples in Manner, Stark, and Wied (2017) and Oh and Patton (2017) a *bloc-equidependence* model as described in section 2.3 was chosen. The model was based on the equations:

$$\begin{aligned} \mathbf{Y} &= (Y_1, \dots, Y_{21})' \sim F_{\mathbf{Y}} = C(F_{Y_1}, \dots, F_{Y_{21}}) \\ (X_1, \dots, X_{21})' &= \beta \mathbf{Z} + \epsilon \\ Y_i &\sim N(0, 1), Z_0 \sim \text{skew} - t(4, -0.8), Z_j \sim t(4), \epsilon \sim t(4) \\ i &= 1, \dots, 21, j = 1, \dots, 3 \\ k_1 &= k_2 = k_3 = 7 \\ \theta &= (\beta_1, \dots, \beta_6)' \end{aligned} \tag{3.2}$$

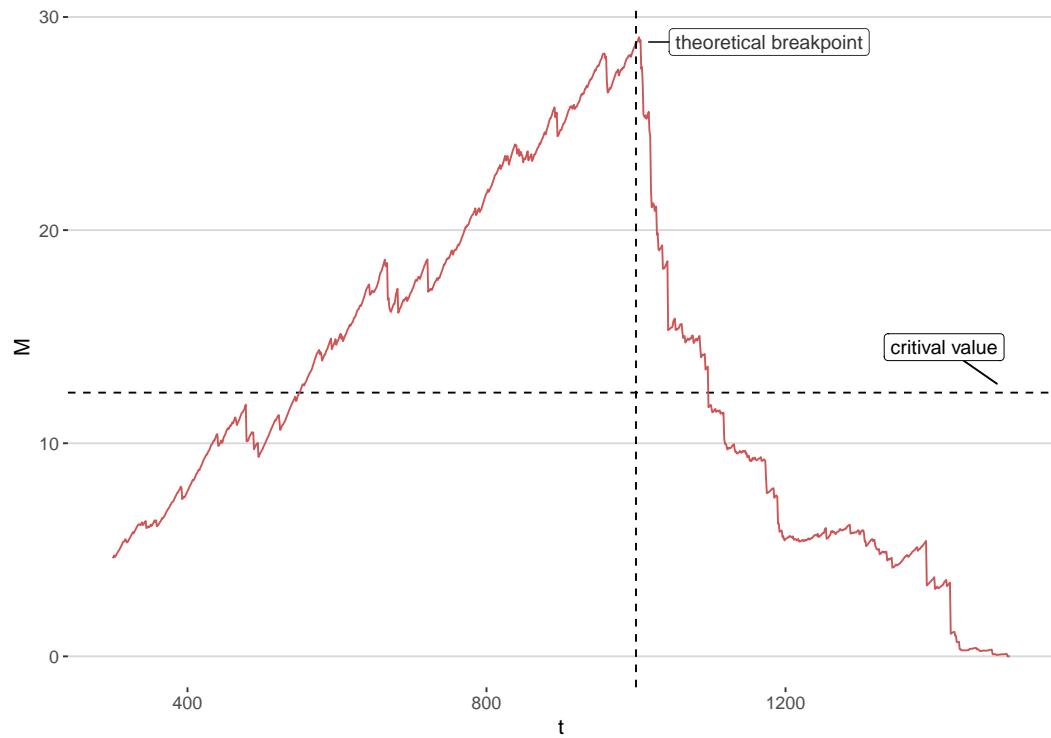
The 21 observable variables were partitioned in 3 groups of equal size. Therefore the parameters to be estimated reduced from $0.5 * N * (N - 1) = 210$ to just $2M = 6$. As in the first example the marginal distributed are iid standard normal. We chose $T = 1500$, $S = 25000$ and a breakpoint at $t = 1000$. Before the break, $\theta_0 = (0, 1, 1, 0, 1, 1)$ and after

3. *factorcopula* - an R package for simulation and estimation of factor copulas

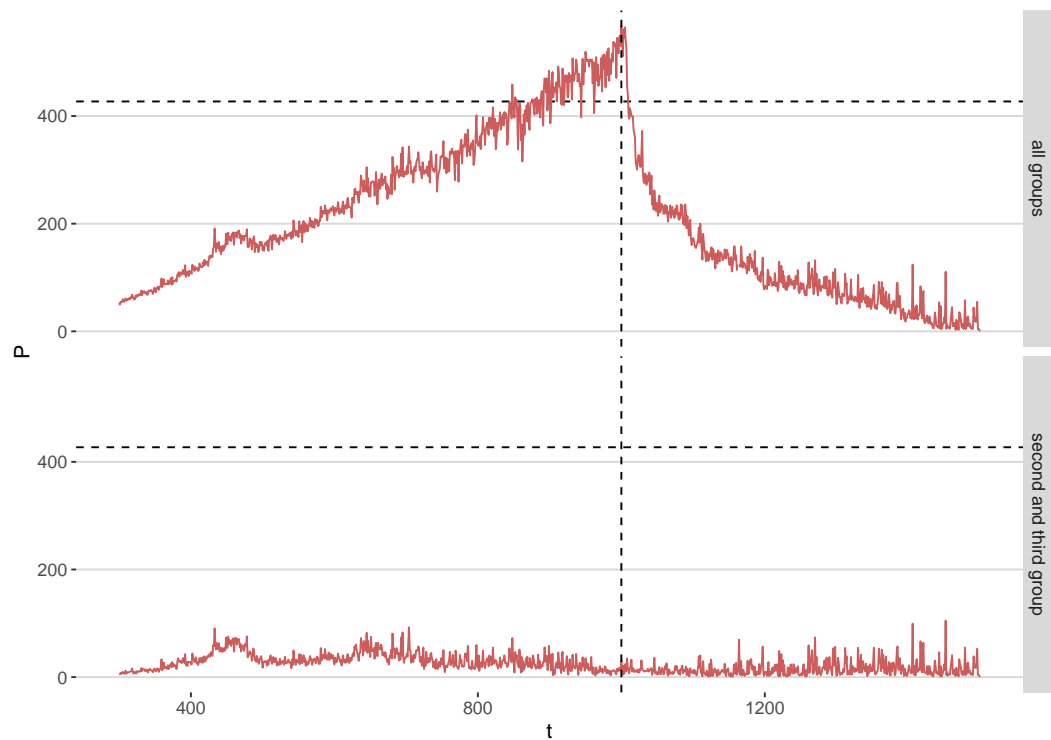
the break $\theta_1 = (1.5, 1, 1, 1.5, 1, 1)$. Thus, only the intra- and interdependence for the first group increases from 0 to 1.5.

Figure 3.2 shows the result of the break test for a single recursive run of the simulation for $t = 300, \dots, 1500$. Both the moments and copula based versions are shown together with the estimated critical value based on 1000 bootstrap replications.

3. *factorcopula* - an R package for simulation and estimation of factor copulas



(a) Moments based test.



(b) Copula based test.

Figure 3.2: Structural break test for a bloc-equipendence model with $N = 21$ and 3 groups of equal size. The theoretical breakpoint is at $t = 1000$ and is modelled as a change of the intra- and interdependenvy of the first group from 0 to 1.5.

4. *Modelling topic dependencies over time with factor copulas*

The moments based test detects a breakpoint at $t = 962$ which is close to the true value. The null hypothesis of no break is clearly rejected since the test statistic fluctuates too strong.

4. Modelling topic dependencies over time with factor copulas

In this chapter we apply the previously discussed methods to a set of relative frequencies over time. The frequencies are derived by counting the relative number of times a specific regular expression appeared in documents of a large text collection of social media posts.

First the dataset is presented. Second, we give a detailed description of the feature generation process and some descriptive overview. Third, the model setup is explained and estimation results of various factor copula models applied to the residuals of the word frequencies are presented. Finally, we apply both the moments and copula based structural break test.

4.1. The *btw17* social media dataset

The original dataset consists of social media posts published by public pages on *Facebook* between January 2014 and December 2017.

Using the official list of the candidates for the Bundestag election in 2017 (*btw17*), the public facebook page for each of the politicians was manually researched. Only candidates from the seven largest parties (CDU, CSU, SPD, Die Linke, Bündnis 90/ Die Grünen, AfD, FDP) were considered. Around 83.8% of all 2516 candidates have an account on Facebook (see Stier et al. 2018, p. 16).

Due to API and privacy restrictions only information from public pages can be accessed such that only around 52.4% of the social media accounts could be considered for the data collection. In addition, 113 official pages from political parties, both on the federal and regional level, were added.

4. Modelling topic dependencies over time with factor copulas

The data collection software used the *restfb* Java client library to call Facebook’s official *Graph API*. The posts were stored in a document orientated database on servers in Germany. In addition to the posts’s content, they are tagged with a timestamp, the user-id of the author and the number of likes and shares the post has gained upon collecting it from the API. (Facebook 2018; Allen and Bartels 2018; MongoDB Inc 2018).

For this analysis the data is restricted on textual posts only. This results in nearly 664 thousand posts. Figure 4.1 shows the monthly number of active accounts and the monthly number of posts by party.³ In early 2014 approximately 500 accounts were active. This number increased steadily to roughly 750 accounts in mid 2016. From then until the election in September 2017 the number increased rapidly to almost 1200 active accounts followed by drop after the election. One can also observe a sharp rise in the number of posts in the month of the election, followed by a decline afterwards.

Table 4.1 shows the overall number of *textual* posts, accounts, likes and shares per party. The likes and shares are based on the aggregated sum of the number of likes and shares for each post from this party. Although the right and left wing parties *AfD* and *Linke* have a relatively small number of accounts and posts they generate by far the greatest number of attention in terms of likes and shares. The social democratics party *SPD* has over twice more posts than the *AfD* but roughly generates only half of the likes and only a fifth of the shares.

party	posts	active accounts	likes (in million)	shares (in million)
AfD	74724	162	18.36	7.59
CDU/CSU	169115	267	14.72	1.83
FDP	71083	201	6.32	0.77
Grüne	67188	139	4.03	1.28
Linke	84723	158	16.03	4.23
SPD	196805	290	10.11	1.57
Sum	663638	1217	69.57	17.27

Table 4.1: Number of posts, accounts, likes and shares over the observation period from "2014-01-01" - "2017-12-31"

³An account was defined active if it has at least one post in the month considered.

4. Modelling topic dependencies over time with factor copulas

Accounts and posts by party over time

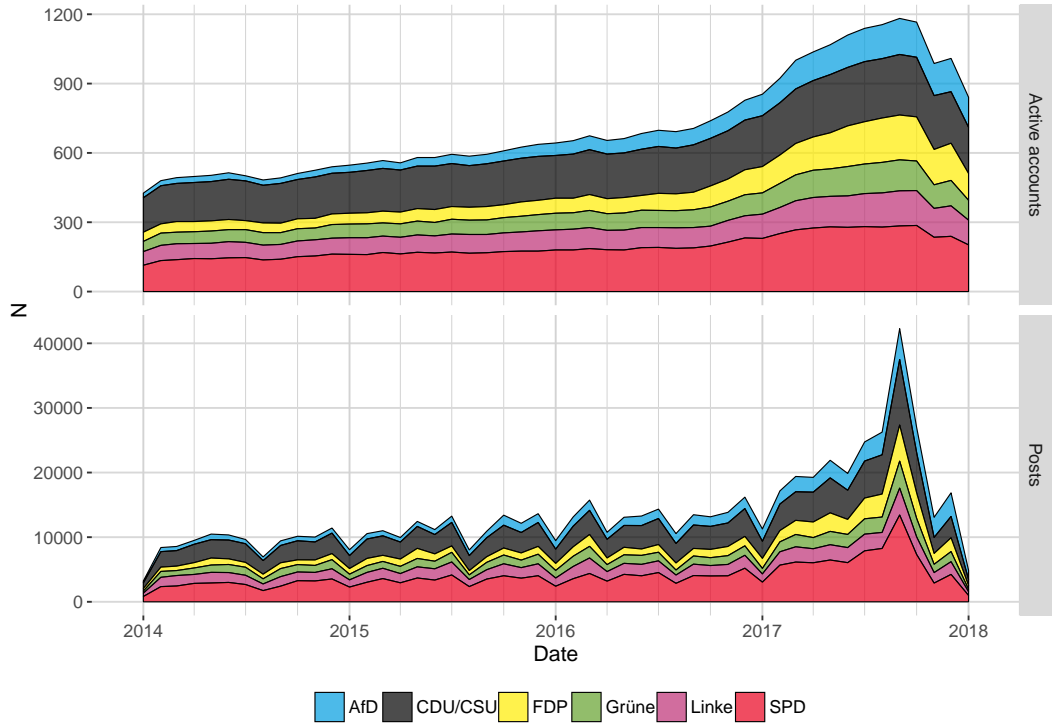


Figure 4.1: Number of active accounts and number of posts per party and month.

4.2. Data processing and descriptive analysis

For this analysis we focus on how the German refugee crisis in 2015 and 2016 is perceived by the political parties in the social media.

First, each post in the dataset is cleaned by removing links and stop-words, transforming umlauts and converting the text to lowercase letters.

Second, a post is considered to be related to the refugee crisis if it matches the regular expression `flucht|fluecht`. For example a match occurs if a post contains the words “flüchtlingskrise”, “fluchtursachen”, “flüchten” or “flüchtlingsheime”.

Third, for each party and day we count the relative number of posts related to the refugee crisis. To account for the importance of a post for the public audience we calculate a weighted average using the weights $\log(l_i + 1)$, where l_i is the number of likes of the post. The resulting

4. *Modelling topic dependencies over time with factor copulas*

dataset consists of a columns for each party, six in total. Each row represents the daily weighted relative frequency of posts matching the regular expression (compare to Figure A.2 in the Appendix). Using the full time span of four years results in $T = 1461$ observations. Thus, the aggregated dataset can be seen as an indicator for how much a party discusses and talks about the topic “refugee crisis” on each day. A positive dependence between all parties would indicate that this topic is relevant for all parties whereas we would see no dependence if the parties talk about this topic independently from each other.

To remove time dependence as discussed in ??, we estimate an ARIMA model for each of the time series. The best univariate model was chosen by running some model candidates over a range of parameters. The model candidate with the lowest Akaike information criterion was considered. All univariate models used $I = 1$, $AR \in \{0, 1\}$ and $MA \in \{0, 1\}$. In the following the estimated standardized residuals are used.

4.3. Results

Table ?? shows the overall pairwise dependency measures used also for the SMM procedure. The dependencies are very weak to weak.

4. Modelling topic dependencies over time with factor copulas

	Pairs	Rank - correlation	Quantile-dependence			
			0.05	0.1	0.90	0.95
AfD-CDU/CSU		0.05	0.14	0.22	0.18	0.10
AfD-FDP		0.06	0.16	0.21	0.16	0.18
AfD-Grüne		0.04	0.14	0.19	0.13	0.15
AfD-Linke		0.10	0.15	0.20	0.22	0.14
AfD-SPD		0.09	0.15	0.21	0.18	0.21
CDU/CSU-FDP		0.17	0.36	0.38	0.24	0.22
CDU/CSU-Grüne		0.12	0.26	0.30	0.25	0.10
CDU/CSU-Linke		0.07	0.23	0.27	0.19	0.12
CDU/CSU-SPD		0.17	0.26	0.30	0.26	0.15
FDP-Grüne		0.10	0.23	0.29	0.21	0.16
FDP-Linke		0.09	0.21	0.23	0.21	0.12
FDP-SPD		0.12	0.29	0.30	0.26	0.21
Grüne-Linke		0.10	0.16	0.26	0.23	0.18
Grüne-SPD		0.14	0.27	0.31	0.23	0.21
Linke-SPD		0.14	0.15	0.25	0.27	0.19
Average		0.10	0.21	0.26	0.21	0.16

4. Modelling topic dependencies over time with factor copulas

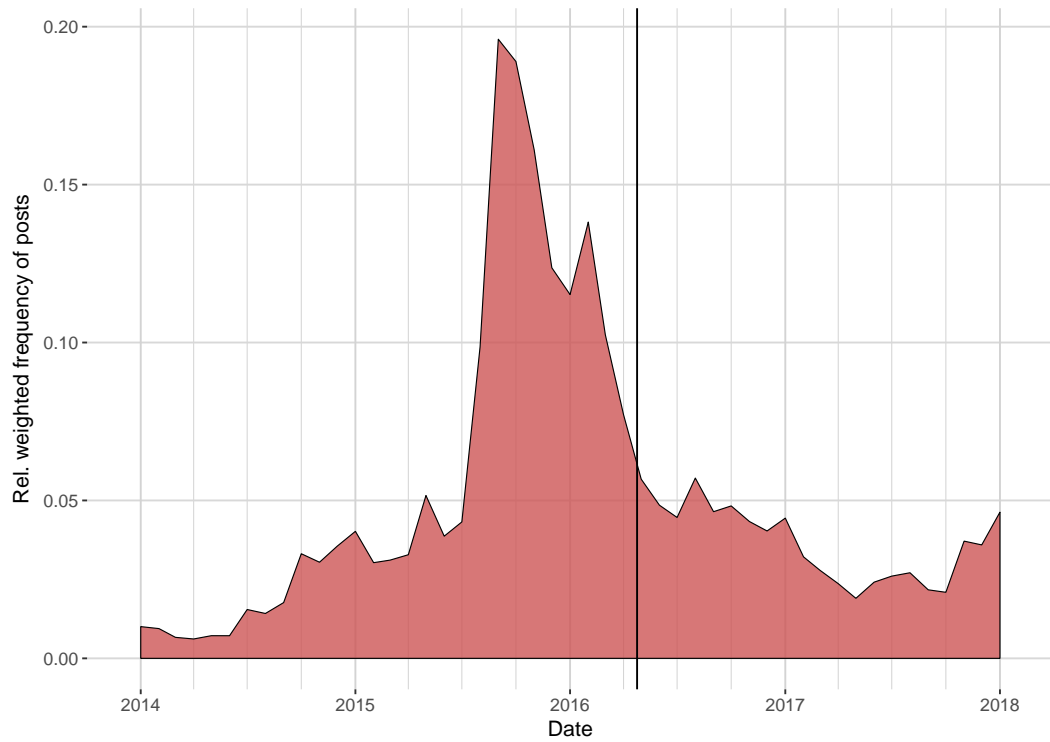


Figure 4.2: Relative weighted monthly frequency of posts matching the regular expression `flucht|fluecht`. The vertical line indicates the breakpoint detected by the moments based test.

First various factor copula models were fitted to the complete residuals.

4. Modelling topic dependencies over time with factor copulas

	Equidependence			Unrestrictive		
	norm-norm	t-t	skewt-t	norm-norm	t-t	skewt-t
β_1	0.5526	0.4002	2.9903	0.2910	0.2195	0.9460
β_2	-	-	-	0.6499	0.5053	2.1030
β_3	-	-	-	0.6821	0.5297	2.3059
β_4	-	-	-	0.5178	0.4012	1.8858
β_5	-	-	-	0.4397	0.3377	1.4266
β_6	-	-	-	0.6674	0.5148	2.1945
df	-	0.4787	0.4900	-	0.4893	0.4703
λ	-	-	-0.1111	-	-	-0.3160
Q	0.0299	0.0124	0.0110	0.5384	0.2704	0.2702

Table 4.2: Estimation results for different one-factor copula specifications.

The unrestrictive moments based test detects a breakpoint at “2016-04-24” with a test statistics of 464. The estimated critical value for an alpha of 0.05 is 103.

Moments based test based on equidependence detect the same breakpoint with a test statistics of 27 and a critical value of 3.61.

We also estimated a recursive equidependence skewt-t factor copula model with fixed distributional parameters.

We also fitted various factor copula models before and after the break detected by the moments based test.

5. Discussion

	Equidependence				Unrestrictive			
	t-t copula		skew t-t copula		skew t-t copula		t-t copula	
	before	after	before	after	before	after	before	after
β_1	0.41	0.30	1.01	0.54	0.46	0.23	0.15	0.14
β_2	-	-	-	-	0.77	0.35	0.23	0.19
β_3	-	-	-	-	0.69	0.36	0.29	0.22
β_4	-	-	-	-	0.59	0.54	0.17	0.35
β_5	-	-	-	-	0.50	0.46	0.23	0.26
β_6	-	-	-	-	0.74	0.67	0.29	0.37
df	2.11	2.20	2.26	2.35	2.69	3.12	1.00	1.63
λ	-	-	-0.10	0.06	-0.42	0.16	-	-
Q	0.04	0.00	0.01	0.00	0.27	0.12	0.24	0.11
T	845	616	845	616	845	616	845	616

Table 4.3: Estimation results for different one-factor copula specifications before and after the breakpoint detected by the moments based test.

5. Discussion

A. Appendix

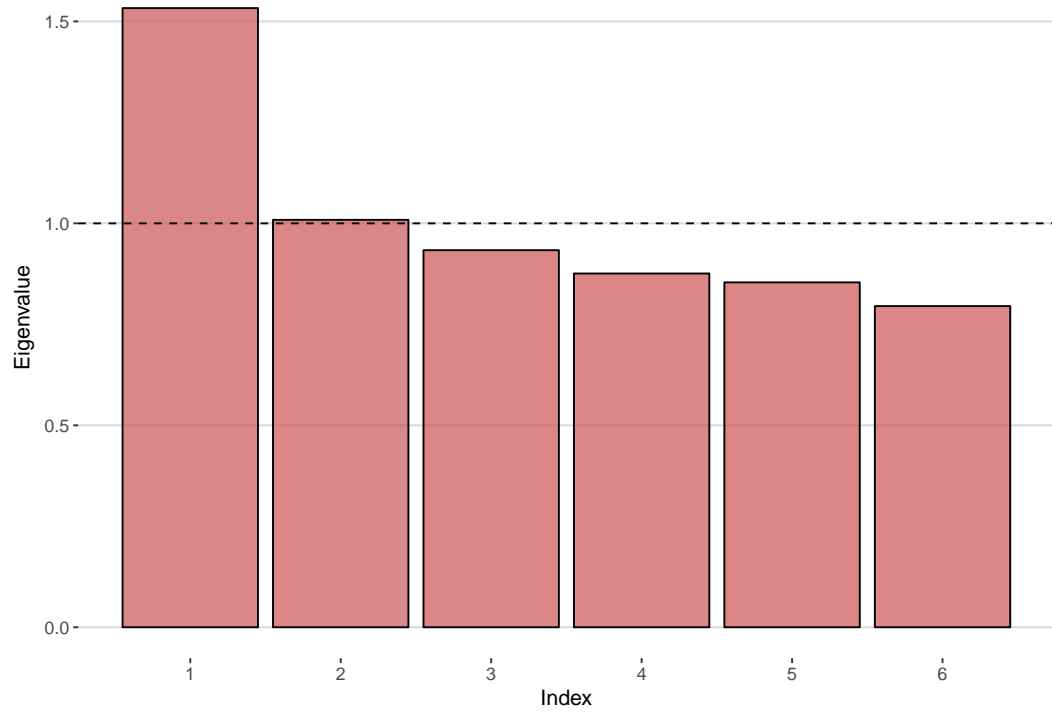
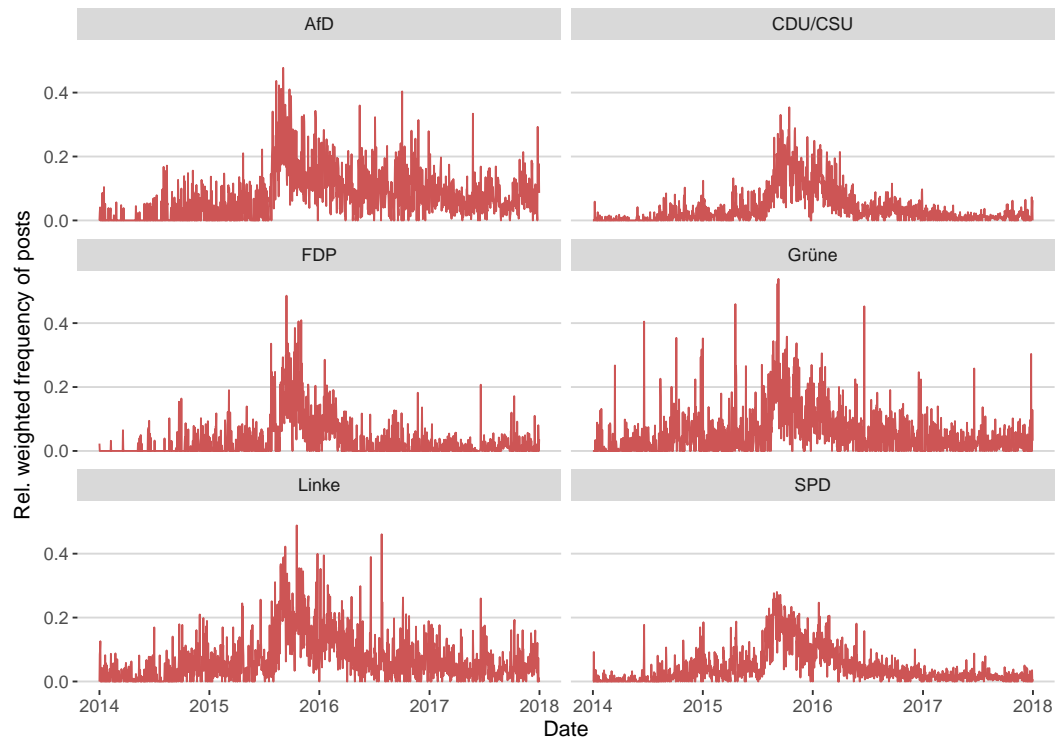
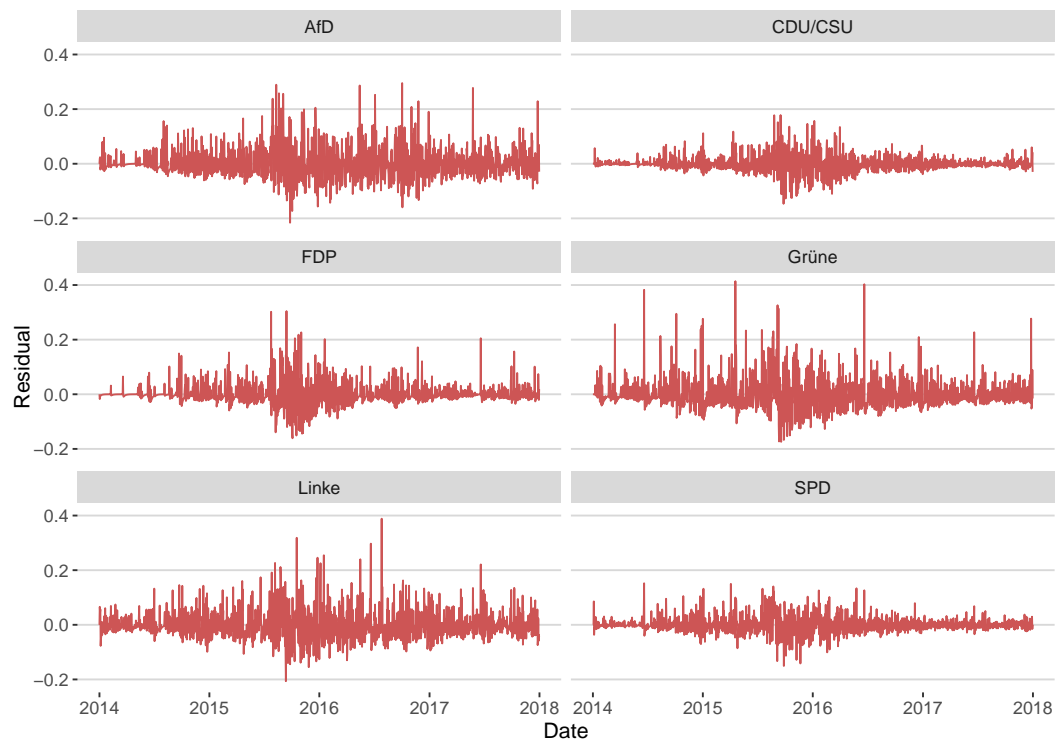


Figure A.1: Scree-plot of ranked eigenvalues based on the pairwise rank-correlation matrix.

A. Appendix



(a) Observed daily relative frequency of posts which match the regular expression.



(b) Residuals of ARIMA models applied to the time series.

Figure A.2

A. Appendix

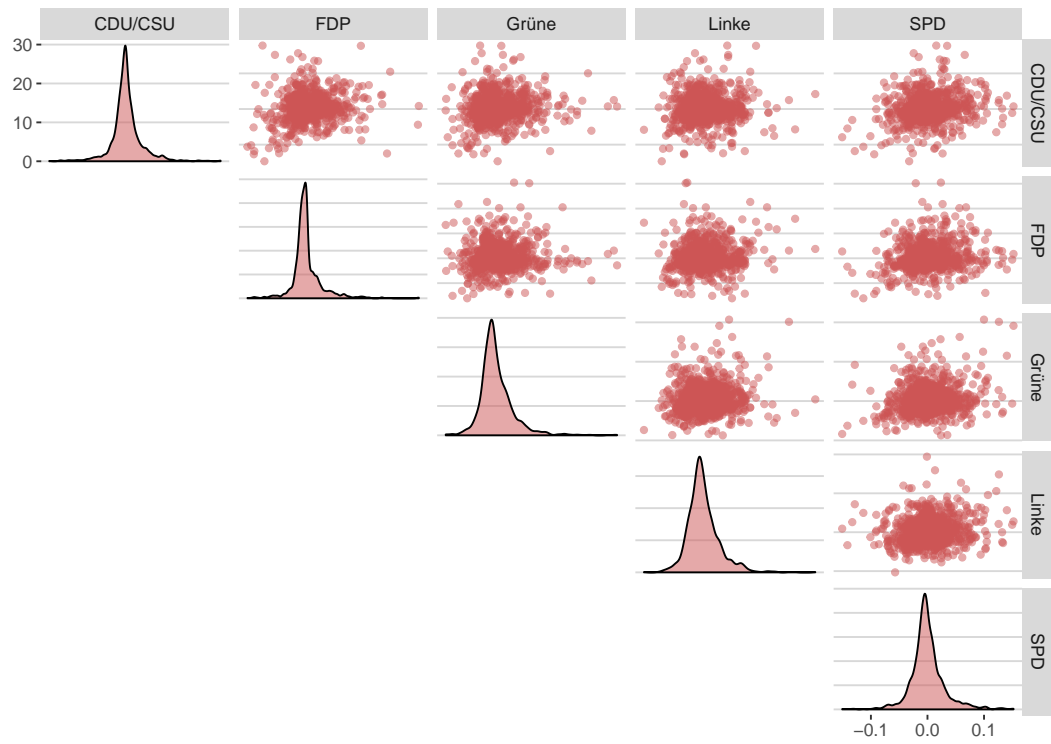


Figure A.3: Pairwise scatterplot of the estimated residuals.

B. References

- Allen, M. and N. Bartels (2018). *RestFB - a pure Java Facebook Graph API client*. <http://restfb.com>.
- Bonart, M. (2018). *factorcopula - high dimensional dependence modelling with factor copulas*. <https://github.com/bonartm/factorcopula>.
- Chen, X. and Y. Fan (2006). "Estimation and model selection of semi-parametric copula-based multivariate dynamic models under copula misspecification". In: *Journal of Econometrics* 135.1, pp. 125–154.
- Facebook (2018). *Facebook's Graph API*. <https://developers.facebook.com/docs/graph-api>.
- Joe, H. (2015). *Dependence modelling with copulas*. Taylor & Francis Group.
- Johnson, S. G. (n.d.). *The NLOpt nonlinear-optimization package*. URL: <http://ab-initio.mit.edu/nlopt>.
- Kumar, P. (2011). "Copulas: Distribution Functions and Simulation". In: *International Encyclopedia of Statistical Science*. Ed. by M. Lovric. Berlin, Heidelberg: Springer, pp. 309–312.
- Manner, H., F. Stark, and D. Wied (2017). "Testing for Structural Breaks in Factor Copula Models". Working Paper.
- MongoDB Inc (2018). *Java driver for Mongo DB*. <https://mongodb.github.io/mongo-java-driver/>.
- Nelsen, R. B. (1999). *An Introduction to Copulas*. Springer.
- Oh, D. H. and A. J. Patton (2013). "Simulated Method of Moments Estimation for Copula-Based Multivariate Models". In: *Journal of the American Statistical Association* 108.502, pp. 689–700.
- (2017). "Modeling Dependence in High Dimensions With Factor Copulas". In: *Journal of Business & Economic Statistics* 35.1, pp. 139–154.
- Patton, A. J. (2006). "MODELLING ASYMMETRIC EXCHANGE RATE DEPENDENCE". In: *International Economic Review* 47.2, pp. 527–556.
- (2009). "Copula-Based Models for Financial Time Series". In: *Handbook of Financial Time Series*. Ed. by T. Andersen, R. A. D. and J.-P. Kreiß, and T. Mikosch. Berlin, Heidelberg: Springer, pp. 767–785.
- Rowan, T. H. (1990). "Functional Stability Analysis of Numerical Algorithms". UMI Order No. GAX90-31702. PhD thesis. Austin, TX, USA.

B. References

- Schmid, F. et al. (2010). “Copula-Based Measures of Multivariate Association”. In: *Copula Theory and Its Applications*. Ed. by P. Jaworski et al. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 209–236.
- Sempi, C. (2011). “Copulas”. In: *International Encyclopedia of Statistical Science*. Ed. by M. Lovric. Berlin, Heidelberg: Springer, pp. 302–305.
- Sklar, A. (1959). “Fonctions de répartition à n dimensions et leurs marges”. In: *Publications de l’Institut de Statistique de L’Université de Paris* 8, pp. 229–231.
- Stier, S. et al. (2018). “Systematically Monitoring Social Media: the case of the German federal election 2017”. In: *GESIS papers* 2018/04, p. 25.

C. Statutory Declaration

Eidesstattliche Versicherung

Hiermit versichere ich an Eides statt, dass ich die vorliegende Arbeit selbstständig und ohne die Benutzung anderer als der angegebenen Hilfsmittel angefertigt habe. Alle Stellen, die wörtlich oder sinngemäß aus veröffentlichten und nicht veröffentlichten Schriften entnommen wurden, sind als solche kenntlich gemacht. Die Arbeit ist in gleicher oder ähnlicher Form oder auszugsweise im Rahmen einer anderen Prüfung noch nicht vorgelegt worden. Ich versichere, dass die eingereichte elektronische Fassung der eingereichten Druckfassung vollständig entspricht.

Köln, den 17. April 2018

(Malte Bonart)