

# STAT 501 Project 2

## PCA for Discrete Distribution with Gaussian Copula

Yifan Zhu

May 7, 2018

### Abstract

In the report, we extend the semiparametric scale-invariant PCA method called copula PCA proposed by Han and Liu (2014) to the case when our distribution is discrete. The senators' voting data set was used to apply this method.

## 1 Introduction

In this project, our goal is to reduce the dimension of categorical data with PCA. However, we cannot perform PCA directly on the categorical data, since the sample space only has finite points and we cannot find a proper set of basis (or principle components) for this space like what we did for continuous distribution. Hence we consider the categorical data with discrete distribution are generated from some latent variables with continuous distribution. Then instead of performing PCA on the observed data, we perform PCA on the latent variables. We also want to transform or connect the latent variables to multivariable normal r.v.s. Since for dimension reduction, PCA is an optimal dimension reduction method for multivariate normal data.

## 2 Method of copula PCA and extension to discrete case

### 2.1 nonparanormal distribution

**Definition 2.1** (nonparanormal distribution, Han and Liu (2014)). Let  $f^0 = \{f_j^0\}_{j=1}^d$  be a set of strictly increasing univariate functions. We say that a  $d$  dimensional random vector  $\mathbf{X} = (X_1, \dots, X_d)^T$  follows a nonparanormal distribution  $\text{NPN}_d(\mathbf{\Sigma}_0, f^0)$ , if

$$f^0(\mathbf{X}) := (f_1^0(X_1), \dots, f_d^0(X_d))^T \sim N_d(\mathbf{0}, \mathbf{\Sigma}^0)$$

where  $\text{diag}(\mathbf{\Sigma}^0) = \mathbf{1}$ .

From the definition, the nonparanormal distribution must be a continuous distribution. And it is actually a special case of Gaussian copula family, when the distribution is continuous.

Suppose we have  $\mathbf{X} \sim \text{NPN}_d(\mathbf{\Sigma}^0, f^0)$ , then the marginal cdfs are

$$F_j(x_j) = P(X_j \leq x_j) = P(f_j^0(X_j) \leq f_j^0(x_j)) = \Phi(f_j^0(x_j))$$

And the joint cdf is

$$\begin{aligned} F(x_1, x_2, \dots, x_d) &= P(X_1 \leq x_1, \dots, X_d \leq x_d) \\ &= P(f_1^0(X_1) \leq f_1^0(x_1), \dots, f_d^0(X_d) \leq f_d^0(x_d)) \\ &= \Phi_{\mathbf{\Sigma}^0}(f_1^0(x_1), \dots, f_d^0(x_d)) \\ &= \Phi_{\mathbf{\Sigma}^0}(\Phi^{-1}(\Phi(f_1^0(x_1))), \dots, \Phi^{-1}(\Phi(f_d^0(x_d)))) \\ &= \Phi_{\mathbf{\Sigma}^0}(\Phi^{-1}(F_1(x_1)), \dots, \Phi^{-1}(F_d(x_d))) \\ &= C_{\mathbf{\Sigma}^0}(F_1(x_1), \dots, F_d(x_d)) \end{aligned}$$

where  $\Phi_{\Sigma^0}$  is cdf of  $N_d(\mathbf{0}, \Sigma^0)$  and  $C_{\Sigma^0}$  is a Gaussian copula with correlation matrix  $\Sigma^0$ . So  $\mathbf{X}$  sin  $\text{NPN}_d(\Sigma^0, f^0)$  actually means  $(F_1(X_1), \dots, F_d(X_d))^T \sim C_{\Sigma^0}$  or  $(\Phi^{-1}(F_1(X_1)), \dots, \Phi^{-1}(F_d(X_d)))^T \sim N(\mathbf{0}, \Sigma^0)$ .

## 2.2 Estimation of $\Sigma^0$

For nonparanormal distribution, we already know  $F_j(x_j) = \Phi(f_j^0(x_j))$ , then one natural way is to first estimate  $f_j^0(x_j) = \Phi^{-1}(F_j(x_j))$ , and then use the estimated  $\hat{f}_j^0$  to transform the data to multivariate normal and estimate  $\Sigma^0$  as if we have data from a multivariate normal distribution. Suppose we have our data  $\mathbf{x}_1, \dots, \mathbf{x}_n$  from nonparanormal distribution, and  $\mathbf{x}_i = (x_{i1}, \dots, x_{id})^T$ , then Liu et al. (2012) gave an estimation of  $f_j^0$  by

$$\hat{f}_j^0(t) = \Phi^{-1} \left( T_{\delta_n}[\hat{F}_j(t)] \right)$$

where  $T_{\delta_n}$  is a Winsorization (or truncation) operator defined as  $T_{\delta_n}(x) = \delta_n I(x < \delta_n) + x I(\delta_n \leq x \leq 1 - \delta_n) + (1 - \delta_n) I(x > 1 - \delta_n)$  with  $\delta_n = 1/(4n^{1/4} \sqrt{\pi \log n})$ . And  $\hat{F}_j(t) = \frac{1}{n+1} \sum_{i=1}^n I(x_{ij} \leq t)$  is the scaled empirical cdf of  $X_j$ . Then

$$\hat{\Sigma}_{jk}^0 = \frac{\frac{1}{n} \sum_{i=1}^n \hat{f}_j^0(x_{ij}) \hat{f}_k^0(x_{ik})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{f}_j^0(x_{ij}))^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{f}_k^0(x_{ik}))^2}}, \hat{\Sigma}^0 = [\hat{\Sigma}_{jk}^0]$$

Another way to estimate  $\Sigma^0$  directly without estimating  $f_j^0$  by Liu et al. (2012) uses Spearman's  $\rho$  and Kendall's  $\tau$  statistics. Let  $r_{ij}$  be the rank of  $x_{ij}$  among  $x_{1j}, \dots, x_{nj}$  and  $\bar{r}_j = \frac{1}{n} \sum_{i=1}^n r_{ij}$ , then

$$\begin{aligned} \hat{\rho}_{jk} &= \frac{\sum_{i=1}^n (r_{ij} - \bar{r}_j)(r_{ik} - \bar{r}_k)}{\sqrt{\sum_{i=1}^n (r_{ij} - \bar{r}_j)^2} \cdot \sqrt{\sum_{i=1}^n (r_{ik} - \bar{r}_k)^2}} \\ \hat{\tau}_{jk} &= \frac{2}{n(n-1)} \sum_{1 \leq i < i' \leq n} \text{sign}(x_{ij} - x_{i'j})(x_{ik} - x_{i'k}) \end{aligned}$$

The population version of Spearman's  $\rho$  and Kendall's  $\tau$  are

$$\rho_{jk} = \text{Corr}(F_j(X_j), F_k(X_k)), \tau_{jk} = \text{Corr}(\text{sign}(X_j - \tilde{X}_j), \text{sign}(X_k - \tilde{X}_k))$$

where  $\tilde{X}_j$  and  $\tilde{X}_k$  are two independent copies of  $X_j$  and  $X_k$ .

**Lemma 2.1** (Liu et al. (2012)). *Assuming  $\mathbf{X} \sim \text{NPN}(\Sigma^0, f^0)$ , we have*

$$\Sigma_{jk}^0 = 2 \sin \left( \frac{\pi}{6} \rho_{jk} \right) = \sin \left( \frac{\pi}{2} \tau_{jk} \right)$$

By Lemma 2.1, we then estimate  $\Sigma^0$  by

$$\hat{\Sigma}_{jk}^0 = 2 \sin \left( \frac{\pi}{6} \hat{\rho}_{jk} \right) = \sin \left( \frac{\pi}{2} \hat{\tau}_{jk} \right)$$

For  $\hat{\Sigma}^0 = [2 \sin(\frac{\pi}{6} \hat{\rho}_{jk})]$ , Han and Liu (2014) also showed when  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \stackrel{\text{iid}}{\sim} \text{NPN}_d(\Sigma^0, f^0)$ , for any  $n > 21/\log d + 2$ ,

$$P \left( \|\hat{\Sigma}^0 - \Sigma^0\|_{\max} \leq 8\pi \sqrt{\frac{\log d}{n}} \right) \geq 1 - \frac{2}{d^2}$$

where  $\|\mathbf{M}\|_{\max} = \max\{|\mathbf{M}_{ij}|\}$ .

Then we can use either estimation of  $\Sigma^0$  to perform our PCA. In Han and Liu (2014) and Han and Liu (2013), sparse assumption about the eigenvectors are imposed, and some theoretical properties were discussed about the sparse PCA solution.

### 2.3 Generalized distributional transform and continuous latent variables

Now we want to extend the above method to catagorical case. The main idea is to define a latent continuous variable and apply the estimation of  $\Sigma^0$  with Spearman's  $\rho$  and Kendall's  $\tau$  (These two estimations both assume data are from continuous distribution).

**Definition 2.2** (generalized distributional transformation, Rüschendorf (2013), Chapter 1). Let  $Y$  be a real random variable with distribution function  $F$  and let  $V$  be a random variable independent of  $Y$ , such that  $V \sim \text{Uniform}(0, 1)$ . The generalized distributiona function is defined by

$$F(x, \lambda) = P(X < x) + \lambda P(Y = x)$$

and we call

$$U = F(Y, V)$$

the generalized distributional transform of  $Y$ .

**Theorem 2.2** (Rüschendorf (2013), Chapter 1). *Let  $U = F(Y, V)$  be the distributional transform of  $Y$  as defined above, then*

$$U \sim \text{Uniform}(0, 1) \text{ and } Y = F^{-1}(U) \text{ a.s.}$$

where

$$F^{-1}(t) = \inf\{x \in \mathbb{R} : F(x) \geq t\}$$

is the generalized inverse of  $F$ , or the quantile transform of  $F$ .

Now suppose our  $\mathbf{X}$  is from a discrete distribution of Gaussian copula family, then

$$F(x_1, \dots, x_d) = C_{\Sigma^0}(F_1(x_1), \dots, F_d(x_d)),$$

which means

$$(U_1, \dots, U_d)^T = (F_1(X_1, V_1), \dots, F_d(X_d, V_d))^T \sim N(\mathbf{0}, \Sigma^0)$$

by Skalar's Theorem (see Rüschendorf (2013)).

Thus we can consider  $U_j = F_j(X_j, V_j)$ ,  $j = 1, \dots, d$  to be our latent variables, where  $F_j(x_j, 1) = P(X_j \leq x_j)$  is the marginal cdf of  $\mathbf{X}$ . And our data is generated by a quantile transformation of  $U_j$ , i.e.  $X_j = F_j^{-1}(U_j)$

To transform the catagorical data to the latent continuous one, we will use the empirical cdf and empirical pmf. Suppose for  $X_j$  there are  $m$  possible values which are  $c_1 < c_2, \dots, c_m$ , then the empirical pmf is

$$\hat{P}_j(c_l) = \frac{1}{n} \sum_{i=1}^n I(x_{ij} = c_l)$$

and the empirical cdf is

$$\hat{F}_j(t) = \sum_{i=1}^n I(x_{ij} \leq t)$$

and

$$\hat{F}_j(c_l) = \sum_{k=1}^l \hat{P}_j(c_k)$$

Then our transformation would be

$$\hat{F}_j(c_l, V) = \hat{F}_j(c_{l-1}) + \hat{P}_j(c_l)V$$

when  $l = 1$ , we denote  $\hat{F}_j(c_0) = 0$ .

And the quantile transform would be

$$\hat{F}_j^{-1}(u) = \sum_{k=1}^m c_k I(\hat{F}_j(c_{k-1}) < u \leq \hat{F}_j(c_k))$$

We perform the copula PCA with  $\Sigma^0$  estimated by Spearman's  $\rho$  or Kendall's  $\tau$  from the data we obtained with generalized distributional transformation. Then we can do dimension reduction for  $(\Phi^{-1}(U_1), \dots, \Phi^{-1}(U_d))^T$ . We can also transform the data projected to a lower dimension back to the categorical data and see the difference.

Also, we need to note that the estimated  $\hat{\Sigma}^0$  using this method might not be a correlation matrix (might not be positive and with 1's in the diagonal). Thus in practice we find the nearest correlation matrix of  $\hat{\Sigma}^0$  and use this correlation matrix to perform our PCA. Method used here is by Higham (2002) to find the nearest correlation matrix in terms of  $\infty$  norm.

### 3 Senators' voting data set

We use senators' voting data to try this method. In this data set, there are 100 observations, each observation has 542 components. Possible values are  $-1, 0, 1$ . We first transform the data to continuous one with empirical marginal distribution and then calculate the estimation of  $\Sigma^0$ . We can also test the multivariate normality, since a basic assumption we made is that the data is from Gaussian copula family.

We test the multivariate normality of the marginal transformed data by projecting the data to one dimension of some arbitrary directions, e.g. 10000 different directions. The p-value is **0.03988561**, which means our assumption that the data are from a Gaussian copula family might not hold. But we will still continue in spite of this.

We used both Spearman's  $\rho$  and Kendall's  $\tau$  estimations. For Spearman's  $\rho$ , the first 38 and 58 PCs can account for 80% and 90% of the total variance. And for Kendall's  $\tau$ , the first 39 PCs and 61 PCs can account for 80% and 90% of the total variance. The cumulative proportion of total variance against number of PCs is shown in Figure 1.

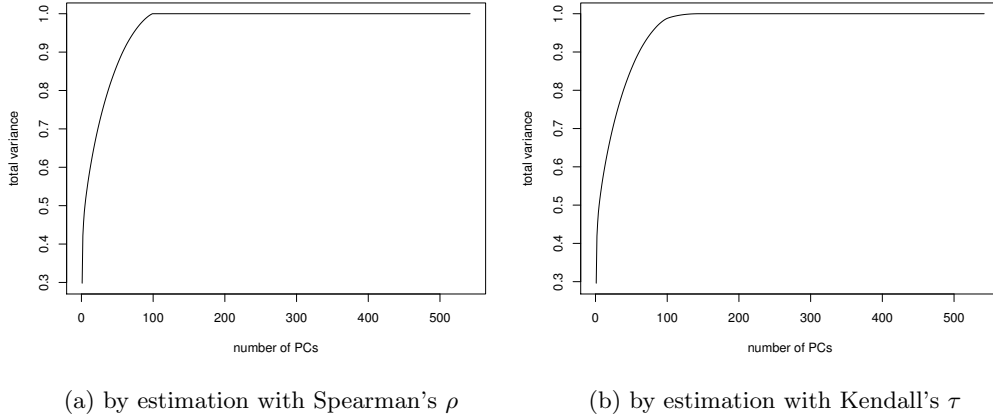


Figure 1: Cumulative proportion of variance against number of PCs with two kinds of estimations of correlation matrix

Then we reconstruct the senators' voting data by these few PCs (number of PCs accounting for 80% and 90% of variance with estimation from Spearman's  $\rho$  and Kendall's  $\tau$ ) with the help of quantile transformation, and compare how much of our categorical data matches with the original

one. The percent of matched data is shown in Table 1. We can see we can reconstruct the data pretty well although our assumption is not really met here.

Table 1: Percent of matching data between dimension reduced and original

	Spearman's $\rho$	Kendall's $\tau$
80%	0.9196863	0.9167159
90%	0.9295387	0.9244834

## 4 Discussion

We can see this method is pretty simple and performs well in the senators' voting data set. Although the assumption is not really met we still get good result: the dimension is reduced a lot and not a lot of information is lost. One reason might be our estimation is based on the rank information of data, and that provided some robustness when deviating from the assumption. We could investigate about this in the future if possible.

## References

- Fang Han and Han Liu. Principal component analysis on non-gaussian dependent data. In *International Conference on Machine Learning*, pages 240–248, 2013.
- Fang Han and Han Liu. High dimensional semiparametric scale-invariant principal component analysis. *IEEE transactions on pattern analysis and machine intelligence*, 36(10):2016–2032, 2014.
- Nicholas J Higham. Computing the nearest correlation matrix problem from finance. *IMA journal of Numerical Analysis*, 22(3):329–343, 2002.
- Han Liu, Fang Han, Ming Yuan, John Lafferty, and Larry Wasserman. The nonparanormal skeptic. *arXiv preprint arXiv:1206.6488*, 2012.
- Ludger Rüschendorf. Mathematical risk analysis. *Springer Ser. Oper. Res. Financ. Eng. Springer, Heidelberg*, 2013.

## Appendix A R codes

```
1 library(readxl)
2 senators<-read_xls("senate_voting_data.xls")
3
4 #b) Plot Andrews' curves
5
6 senators.names<-names(senators)[-c(1,2)]
7 rev.party.state.names<-lapply(X=strsplit(gsub(pattern=".",replacement="",x=senators.names),split="
  ↪   "),FUN = rev)
8
9 senators.party <- lapply(X = rev.party.state.names, FUN = function(x) (unlist(x)[1]))
10 senators.party <- unlist(senators.party)
11
12 senators.last.names <- lapply(X = rev.party.state.names, FUN =
  ↪   function(x) (unlist(x)[4]))
13 senators.last.names <- unlist(senators.last.names)
14
15
16 #Create new data.frame for plotting
17 senators_new <- as.data.frame(t(senators[,-c(1,2)]))
18
19 colnames(senators_new) <- NULL
20 rownames(senators_new) <- NULL
21
22 senators_new <- data.frame(senators_new, party = senators.party)
23
24
25 # Use the codes from Canvas
26 source("ggandrews.R")
27
28 # Display the Andrews' curves
29 ggandrews(senators_new, type = 2, clr = 543, linecol = c("blue", "green", "red"),
  ↪   return_value = FALSE)
30
31
32 epmf <- function(x){
33   n <- length(x)
34   return(c(sum(x==1), sum(x==0), sum(x==1))/n)
35 }
36
37 ecdf <- function(epmf){
38   return(c(0, epmf[1], epmf[1]+epmf[2], 1))
39 }
40
41 senator_epmf <- sapply(senators_new[,ncol(senators_new)], epmf)
42 senator_ecdf <- apply(senator_epmf, MARGIN = 2, ecdf)
43
44 library(plyr)
45
46 gdtrans <- function(i){
47   u1 <- mapvalues(senators_new[,i], from = c(-1,0,1), to = senator_epmf[,i])
48   u2 <- mapvalues(senators_new[,i], from = c(-1,0,1), to = senator_ecdf[1:3,i])
49   return(u2 + u1*runif(n = length(u1)))
50 }
51
52 inv_gdtrans <- function(u, i){
53   -1* (senator_ecdf[1,i] < u[,i] & senator_ecdf[2,i] >= u[,i]) + 1* (senator_ecdf[3,i] <
  ↪   u[,i] & senator_ecdf[4,i] >= u[,i])
54 }
55
56 source("testnormality.R")
57
58
59
```

```

60 senators_c <- sapply(1:542, gdtrans)
61
62 senators_normal <- qnorm(senators_c, mean = 0, sd = 1)
63
64 testnormality(senators_normal)
65
66 library(energy)
67 mvnorm.etest(x = senators_normal, R = 1000)
68
69 library(MASS)
70 library(Matrix)
71
72 # rho based estimation
73 spearmanrho <- cor(senators_c, method = "spearman")
74 sigmahat_rho <- 2*sin(pi/6*spearmanrho)
75 sigmahat_rho <- nearPD(sigmahat_rho, corr = T)
76 decomp_rho <- eigen(sigmahat_rho$mat, symmetric = T)
77
78 s <- decomp_rho$values
79
80 pvar<-s/sum(s)
81
82 # cumulative proportion of total variance explained
83 # by each component
84
85 cpvar <- cumsum(s)/sum(s)
86 plot(x = 1:length(cpvar), y = cpvar, 'n', xlab = 'number of PCs', ylab = 'total
  ↳ variance')
87 lines(x = 1:length(cpvar), y = cpvar)
88
89 npc80 <- min(which(cpvar > 0.8))
90 npc90 <- min(which(cpvar > 0.9))
91
92 senators_normal_80 <-
  ↳ senators_normal%*%decomp_rho$eigenvectors[,1:npc80]%*%t(decomp_rho$eigenvectors[,1:npc80])
93
94 senators_u_80 <- pnorm(senators_normal_80, mean = 0, sd = 1)
95
96 senators_80 <- sapply(1:542, inv_gdtrans, u = senators_u_80)
97
98 sum(senators_80==senators_new[,ncol(senators_new)])/(542*100)
99
100
101 senators_normal_90 <-
  ↳ senators_normal%*%decomp_rho$eigenvectors[,1:npc90]%*%t(decomp_rho$eigenvectors[,1:npc90])
102
103 senators_u_90 <- pnorm(senators_normal_90, mean = 0, sd = 1)
104
105 senators_90 <- sapply(1:542, inv_gdtrans, u = senators_u_90)
106
107 sum(senators_90==senators_new[,ncol(senators_new)])/(542*100)
108
109 # tau based estimation
110 kendalltau <- cor(senators_c, method = "kendall")
111 sigmahat_tau <- sin(pi/2*kendalltau)
112 sigmahat_tau <- nearPD(sigmahat_tau, corr = T)
113 decomp_tau <- eigen(sigmahat_tau$mat, symmetric = T)
114
115 s <- decomp_tau$values
116
117 pvar<-s/sum(s)
118
119 # cumulative proportion of total variance explained
120 # by each component
121
122 cpvar <- cumsum(s)/sum(s)

```

```

123 plot(x = 1:length(cpvar), y = cpvar, 'n', xlab = 'number of PCs', ylab = 'total
    ↳ variance')
124 lines(x = 1:length(cpvar), y = cpvar)
125 npc80 <- min(which(cpvar > 0.8))
126 npc90 <- min(which(cpvar > 0.9))
127
128 senators_normal_80 <-
    ↳ senators_normal%%decomp_tau$variables[,1:npc80]%%t(decomp_tau$variables[,1:npc80])
129
130 senators_u_80 <- pnorm(senators_normal_80, mean = 0, sd = 1)
131
132 senators_80 <- sapply(1:542, inv_gdtrans, u = senators_u_80)
133
134 sum(senators_80==senators_new[,ncol(senators_new)])/(542*100)
135
136
137 senators_normal_90 <-
    ↳ senators_normal%%decomp_tau$variables[,1:npc90]%%t(decomp_tau$variables[,1:npc90])
138
139 senators_u_90 <- pnorm(senators_normal_90, mean = 0, sd = 1)
140
141 senators_90 <- sapply(1:542, inv_gdtrans, u = senators_u_90)
142
143 sum(senators_90==senators_new[,ncol(senators_new)])/(542*100)

```