

P1

GR5241

HW1

Haoyi Li

HW1

$$1. P(x|u, v) = P(x|u, v) = \left(\frac{v}{u}\right)^u \frac{x^{u-1}}{\Gamma(u)} \exp\left(-\frac{vx}{u}\right)$$

$$L(\theta) = \prod_{i=1}^n P(x_i|\theta) = \left(\frac{v}{u}\right)^{nv} \frac{\pi x_i^{u-1}}{\Gamma^{(n)}(u)} \exp\left(-\frac{v}{u} \sum_{i=1}^n x_i\right)$$

suppose $l(\theta)$ to be log-likelihood.

$$l(\theta) = \log L(\theta) = nv(\log v - \log u) - n \log \Gamma(u) + \sum_{i=1}^n (u-1) \log(x_i) - \frac{v}{u} \sum_{i=1}^n x_i$$

$$2. \textcircled{1} \frac{\partial l(\theta)}{\partial u} = -\frac{1}{u} nv + \frac{v}{u^2} \sum_{i=1}^n x_i = 0 \Rightarrow \frac{\sum_{i=1}^n x_i}{u} = n \quad \hat{u} = \frac{\sum x_i}{n}$$

$$3. \textcircled{2} \frac{\partial l(\theta)}{\partial v} = n(\log v - \log u) + n - n \frac{\Gamma'(u)}{\Gamma(u)} + \sum_{i=1}^n \log(x_i) - \frac{\sum_{i=1}^n x_i}{u} = 0$$

$$\Rightarrow n(\log v - \log u) + \sum \log x_i = \sum \log\left(\frac{x_i v}{u}\right)$$

$$-n \frac{\Gamma'(u)}{\Gamma(u)} = -\sum \phi(v)$$

$$\text{so } \textcircled{2} = \sum \log\left(\frac{x_i v}{u}\right) - \frac{\sum x_i}{u} + n - \sum \phi(v) = \sum \left(\log\left(\frac{x_i v}{u}\right) - \left(\frac{x_i}{u} - 1\right) - \phi(v) \right) = 0$$

P2. our target is to get $\min R(f)$, due to monotonicity of integral,

$\min R(f)$ is to get $\min R(f|x)$ ($R(f) = \int R(f(x)) p(x) dx$)

So we need to find f to get $\min R(f|x) = \sum_{y \in [K]} L^{0-1}(y, f(x)) p(y|x)$

$$f(x) = \arg \min_{y \in [K]} \sum_{y \in [K]} L^{0-1}(y, f(x)) p(y|x)$$

$$= \arg \min_{y \in [K]} (1 - p(y|x)) \quad (\text{explain: mismatch loss is 1 and match loss 0})$$

$$= 1 - (1 - p(\text{match})) = 1 - p(y|x)$$

So $f(x) = \arg \max_{y \in \{K\}} P(y|x)$

this is exactly $f_0(x) = \arg \max_{y \in \{K\}} P(y|x)$

So $f_0(x)$ is the classifier with minimum $R(f|x)$, i.e. $\min R(f)$.

5241 HW1 Haiqi Li hl3115

Haiqi Li

27/1/2018

Problem 3

1

First I copied the chart in Wiki with all their symbols to Excel. Then I copied the symbol column and paste with transpose. Finally I get longstring here.

```
options("getSymbols.yahoo.warning"=FALSE)
options("getSymbols.warning4.0"=FALSE)
options(warn = F)
library(quantmod, warn.conflicts = F)

## Loading required package: xts
## Loading required package: zoo
##
## Attaching package: 'zoo'
## The following objects are masked from 'package:base':
##
##   as.Date, as.Date.numeric
## Loading required package: TTR
## Version 0.4-0 included new data defaults. See ?getSymbols.
longstring="AAPL  AXP BA  CAT CSCO  CVX DIS DWDP  GE  GS  HD  IBM INTC  JNJ JPM KO  MCD MMM MRK
"
DJIname <- strsplit(longstring, split = "\t")
DJIname <- DJIname[[1]] #unlist it
DJIname[30] <- "XOM" #final term has a \n
data <- getSymbols("AAPL", auto.assign = F, from = "2017-01-01", to = "2018-01-01")
#initialize data to get number of columns
data <- data[,4]
#since I only use the close price, I pick it up manually

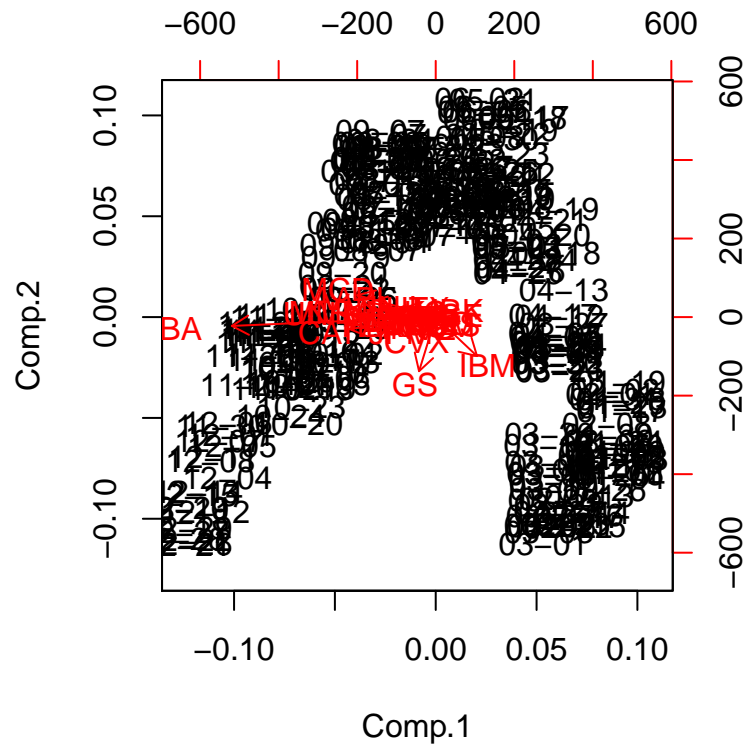
for(i in 2:30){
  datatemp <- getSymbols(DJIname[i], auto.assign = F,
    from = "2017-01-01", to = "2018-01-01")
  #since we only take use of close price, I think to get it first is better
  datatemp <- datatemp[,4]
  data <- cbind.data.frame(data, datatemp)
}
#A for-loop to get other data and put to them into one dataframe.

colnames(data) <- DJIname
for(i in 1:nrow(data)){
  tempname <- substr(rownames(data)[i], start = 6, stop=nchar(rownames(data)[i]))
  rownames(data)[i] <- tempname
}
```

```
}
#rename columns and rows
```

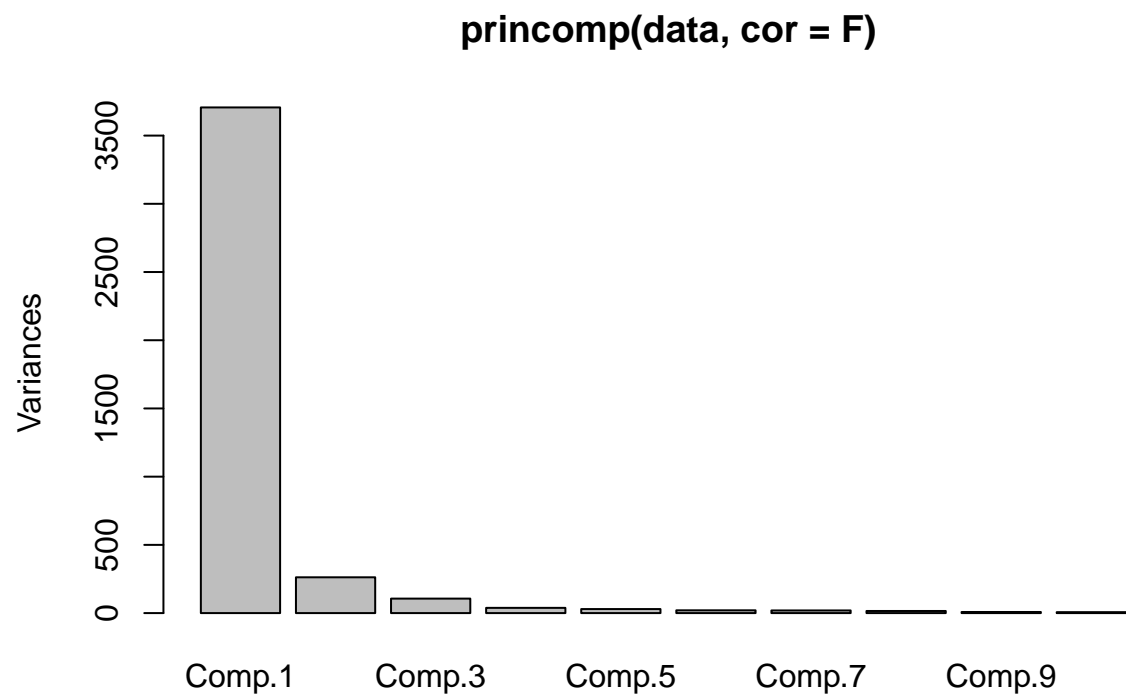
2

```
biplot(princomp(data,cor=F))
```



The biplot here is not very informative since all vectors are very condensed in the picture.

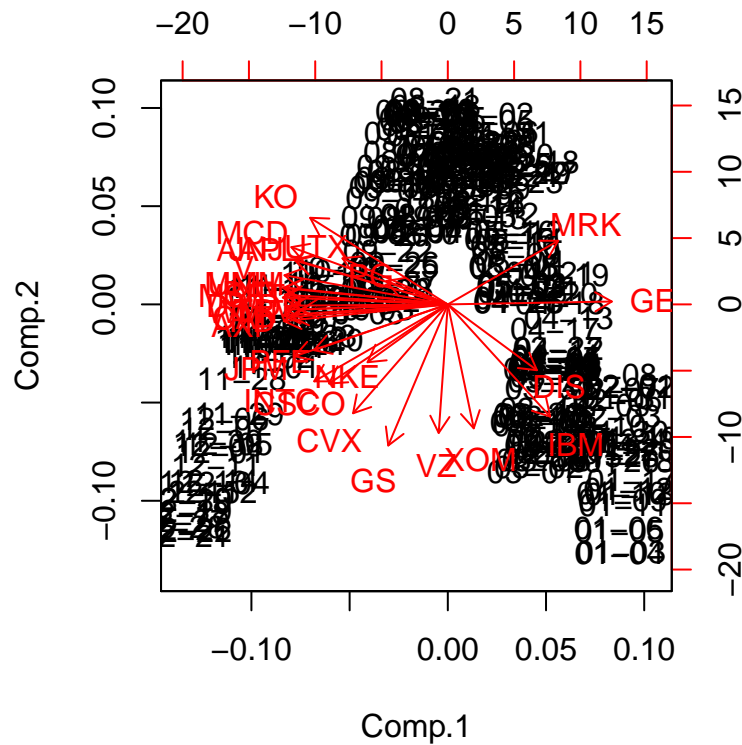
```
screepplot(princomp(data,cor=F))
```



The screeplot here shows that component 1 takes most of variance. I think only one component is really important.

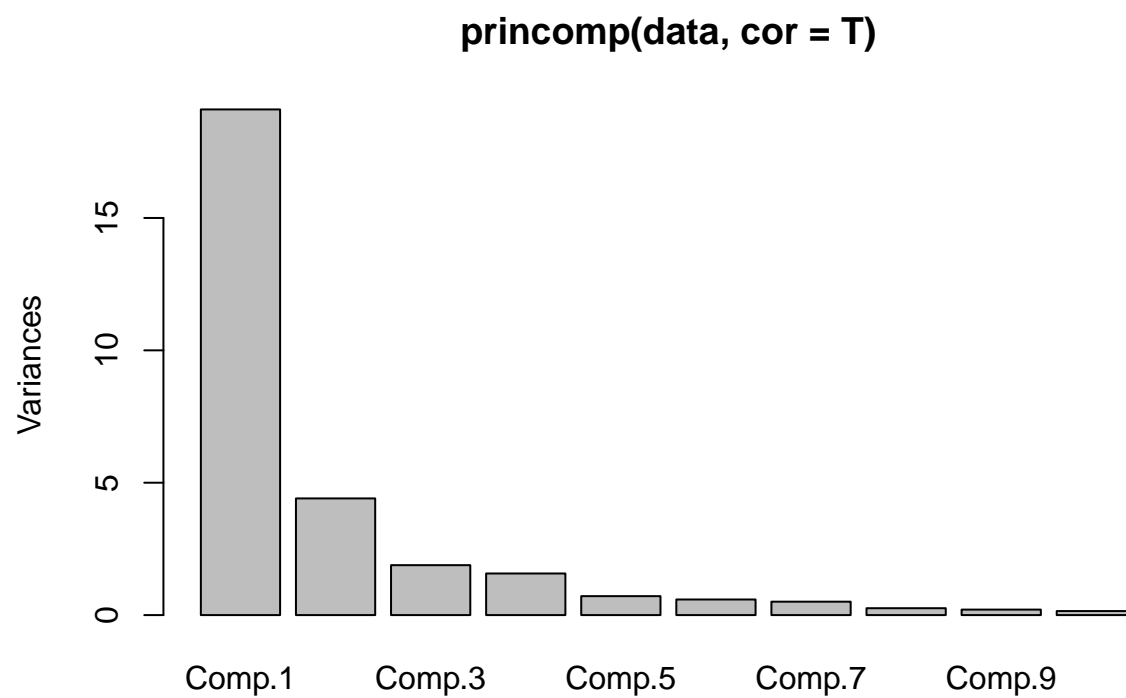
3

```
biplot(princomp(data, cor=T))
```



After modification of scale, I think this time the biplot is much more informative. I noticed that McDonald and Coca-Cola are very close to each other and they are all food companies. Also, most financial companies like JPMorgan, Travelers and Visa are of negative component 2 and are close to each other. Maybe Goldman Sachs is here.

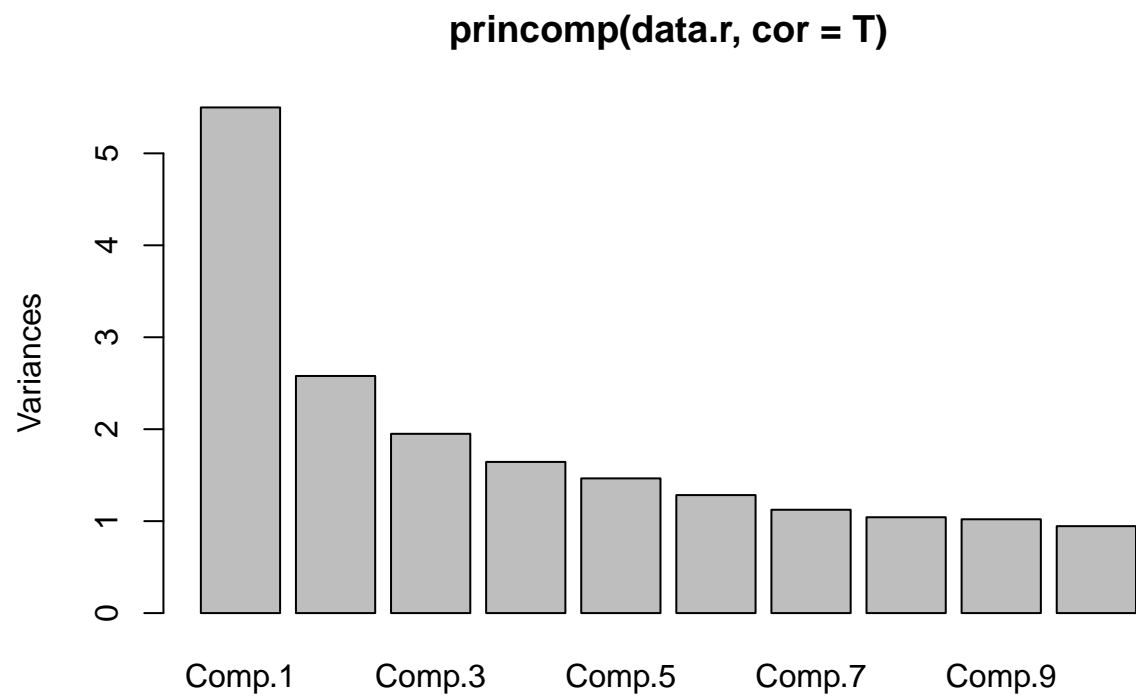
```
screeplot(princomp(data, cor=T))
```



This time, the screeplot shows that Component 2 may count for something, too.

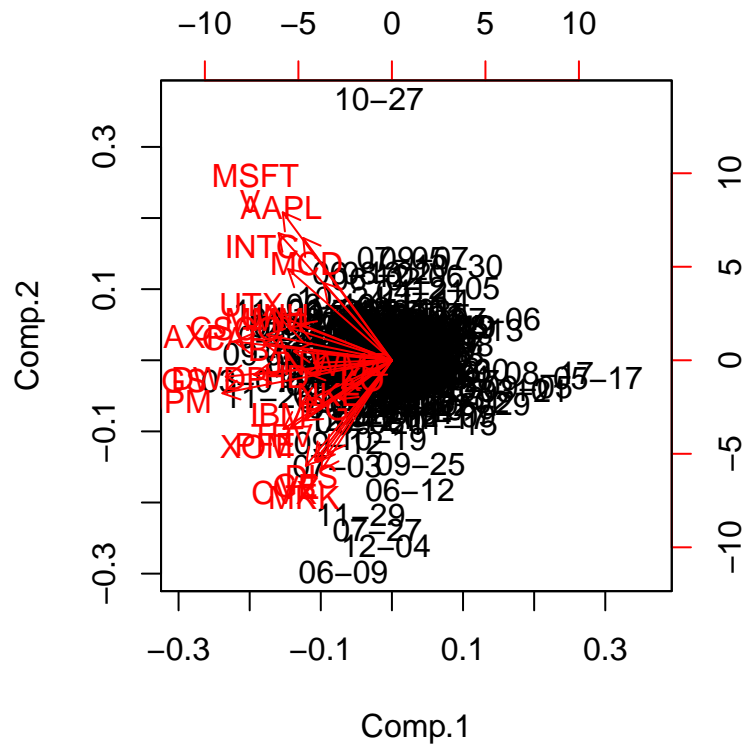
4

```
data <- as.matrix(data)
data.lag <- diff(data, lag=1, differences = 1)
data <- data[-1,]
data.r <- data.lag/data
screeplot(princomp(data.r, cor=T))
```



The screeplot here shows that only component 1 and 2 could not explain all variance. Many other components do make sense.

```
biplot(princomp(data.r, cor=T))
```

The biplot show that all stocks are of negative component 1. I think this tell us some information about the whole market trend since they are all in one direction.

If the stocks fluctuate randomly and independent to each,I think the direction of each stock vector should distribute more uniform, just like a circle.