

# Stat 154 Problem Set One

Jinze Gu SID:24968967

February 20, 2014

```
# Problem One Data Cleaning:
setwd("/Volumes/æĬLëČ;ăĜžæšq/Stat 154/Stat 154")
meta <- read.csv("stock.csv")

## Warning: æŦăæşTæŁŞăijĂæŨĞăżű'stock.csv': No such file or directory
## Error: æŦăæşTæŁŞăijĂéŞĭçZŞ

# To calculate daily returns, I only need price, so I extract the data
# with price, date and company name.
stock <- data.frame(date = meta$date, COMP = meta$COMNAM, PRC = meta$PRC)

## Error: æŁĭăyDăĹřărźésă'meta'

# Just remove the raw data since it takes too much memory.
rm(meta)

## Warning: æŁĭăyDăĹřărźésă'meta'

gc()

##          used (Mb) gc trigger (Mb) max used (Mb)
## Ncells 262592 14.1      467875 25.0   407500 21.8
## Vcells 629202  4.9     1031040  7.9    786432  6.0

compname <- levels(factor(stock$COMP))

## Error: æŁĭăyDăĹřărźésă'stock'

# f is a function that extract the price of the same company with 1342
# price records. It turns out that there are 422 companies with 1342
# sotck price records.
f <- function(x) {
  if (length(stock$PRC[stock$COMP == compname[x]]) == 1342) {
    return(stock$PRC[stock$COMP == compname[x]])
  }
}

stock_l <- sapply(c(1:length(compname)), f)

## Error: æŁĭăyDăĹřărźésă'compname'

# g is a function that extract the name of the company with 1342 price
# records.
g <- function(x) {
  if (length(stock$PRC[stock$COMP == compname[x]]) == 1342) {
    return(compname[x])
  }
}

names <- sapply(c(1:length(compname)), g)
```



```

min(prin_stock$rotation[, 2])
## Error: æĲäŸDăĹřăřžëšă'prin_stock'
hist(prin_stock$rotation[, 3], main = "Third Principal Direction")
## Error: æĲäŸDăĹřăřžëšă'prin_stock'
max(prin_stock$rotation[, 3])
## Error: æĲäŸDăĹřăřžëšă'prin_stock'
min(prin_stock$rotation[, 3])
## Error: æĲäŸDăĹřăřžëšă'prin_stock'

# Comment: Based on the principal loadings, in the first three principal
# component loadings, most of the variables(companies in this case) have
# similar variance. It means that each variable are equally important in
# accounting for the variability in the PC.
# Besides, I plotted the projected data on the first and second principal
# component direction. I don't think there is an obvious clustering.
plot(prin_stock$x[, 1], prin_stock$x[, 2], xlab = "First Principal Component",
     ylab = "Second Principal Component")
## Error: æĲäŸDăĹřăřžëšă'prin_stock'

# Compar the plot that I only use the projected data in the first and
# second PC direction rather than using all of the eigenvectors to
# project the data.
prin_stock <- prcomp(t(daily_return), rtex = TRUE, scale = TRUE)
## Error: æĲäŸDăĹřăřžëšă'daily_return'

stock_proj <- t(prin_stock$rotation[, 1:2]) %*% scale(daily_return)
## Error: æĲäŸDăĹřăřžëšă'prin_stock'

plot(t(stock_proj), main = "Projected data in the first and second PCs")
## Error: æĲäŸDăĹřăřžëšă'stock_proj'

# Comment: I did not see any obviously great clustering.
# Problem One 2). If I use all 422 companies to do hieararchical
# clustering, then I have the following graph.
dist_stock <- dist(t(daily_return), method = "manhattan")
## Error: æĲäŸDăĹřăřžëšă'daily_return'

hc_stock <- hclust(dist_stock, method = "complete")
## Error: æĲäŸDăĹřăřžëšă'dist_stock'

plclust(hc_stock, labels = FALSE)
## Error: æĲäŸDăĹřăřžëšă'hc_stock'

# Comment: I chose manhattan metric since I believe daily return should
# be equally weighted for everyday and it is reasonable to calculate the
# absolute difference between daily return rather than the euclidean
# distance. From the plot, I can see that those companies can be
# classified as different groups based on their variance of daily return.
# I used only 10 companies which would generate a better plot.
sample <- c(1, 40, 60, 114, 210, 275, 89, 320, 170, 413)
dist_stock_sub <- dist(t(daily_return[, sample]), method = "manhattan")

```

```
## Error: æĹ;äÿDâĹřâřžèšq'daily_return'
hc_stock_sub <- hclust(dist_stock_sub, method = "complete")
## Error: æĹ;äÿDâĹřâřžèšq'dist_stock_sub'
plclust(hc_stock_sub, labels = NULL)
## Error: æĹ;äÿDâĹřâřžèšq'hc_stock_sub'
```

Problem One 1): Interpretation of first few vectors of PC loadings: Comment: Based on the principal loadings, in the first three principal component loadings, most of the variables(companies in this case) have similar variance. It means that each variable are equally important in accounting for the variability in the PC.

Problem One 2): Comment: I chose manhattan metric to do hclustering since I believe daily return should be equally weighted for everyday and it is reasonable to calculate the absolute difference between daily return rather than the euclidean distance. From the plot, I can see that those companies can be classified as different groups based on their variance of daily return.