

Q2. (a) MLE:  $\hat{\Sigma} = S_0 = \begin{bmatrix} 0.0144 & 0.0117 \\ 0.0117 & 0.0146 \end{bmatrix}$ .  $\frac{n-1}{n} = \frac{41}{42} \Rightarrow \begin{bmatrix} 0.01406 & 0.01142 \\ 0.01142 & 0.01435 \end{bmatrix}$

(b)  $\frac{n-p}{p} (\bar{X} - \mu)' S_0^{-1} (\bar{X} - \mu) \sim F_{p, n-p}$ ,  $n=42$ ,  $p=2$   $F_{2,40}$ .

Find  $A = \{ \mu \in \mathbb{R}^2 : (\bar{X} - \mu)' S_0^{-1} (\bar{X} - \mu) < \frac{2}{40} F_{1-0.05; 2, 40} = 0.1615 \}$

Set  $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \Rightarrow R_1 = \{ \mu : 26.844 - 32.909\mu_1 - 58.263\mu_2 - 326.57 > \mu_1\mu_2 \}$

$R_1$  is 95% CI for  $\mu$ .  $+203.751\mu_1^2 + 201.033\mu_2^2 < 0.1615$

(c)  $\mu = \begin{pmatrix} 0.6 \\ 0.58 \end{pmatrix} \Rightarrow \text{LHS} = 0.63 > 0.1615 \therefore \mu \notin R_1$ , not in 95% confidence interval

(d)  $R_2 = \{ \mu : a^T \mu \pm \sqrt{h(\alpha)} \sqrt{a^T S_0 a} \}$ ,  $h(\alpha) = \frac{41.2}{40} F_{1-0.05; 2, 40} = 6.615$

$\mu \pm \sqrt{h(\alpha)} \sqrt{a^T S_0 a} = 0.603 \pm 0.0484 \Rightarrow R_2 = [0.555, 0.651] \Rightarrow R_2 = \begin{bmatrix} 0.555, 0.651 \\ 0.555, 0.651 \end{bmatrix}$

(e)  $\mu_0$  known,  $\Sigma$  unknown  $\Rightarrow W = n \log |S_0 + SS'| - n \log |S_0|$

$8S^T S W = 42 \cdot \begin{pmatrix} 0.555-0.6 \\ 0.651-0.58 \end{pmatrix} \begin{pmatrix} 0.0144 & 0.0117 \\ 0.0117 & 0.0146 \end{pmatrix}^{-1} \begin{pmatrix} 0.6-0.58 \\ 0.58-0.58 \end{pmatrix}^T = 26.3$

$P(F_{2,40} \geq 26.3) = 0.000 < 0.05 \therefore$  Thus we can reject  $H_0$  and accept  $\mu = \mu_0$   
p-value

Q3 (a)  $\bar{X}_n \sim N(\mu, \frac{\Sigma}{n}) \Rightarrow a_i^T \bar{X}_n \sim N(a_i^T \mu, \frac{a_i^T \Sigma a_i}{n})$   $\left( \frac{a_i^T \bar{X}_n - a_i^T \mu}{\sqrt{\frac{a_i^T \Sigma a_i}{n}}} \right) \sim N(0,1)$

$\frac{\sqrt{n} \left( \frac{a_i^T \bar{X}_n - a_i^T \mu}{\sqrt{\frac{a_i^T \Sigma a_i}{n}}} \right)}{\sqrt{\frac{a_i^T \Sigma a_i}{n}}} \sim \frac{N(0,1)}{\sqrt{\frac{\chi^2(n-1)}{n-1}}} \sim t(n-1) \Rightarrow P\left( T = \frac{a_i^T \bar{X}_n - a_i^T \mu}{\sqrt{\frac{a_i^T \Sigma a_i}{n}}} \mid T \in [-t_{\alpha/2}, t_{\alpha/2}] \right) = 1 - \alpha$

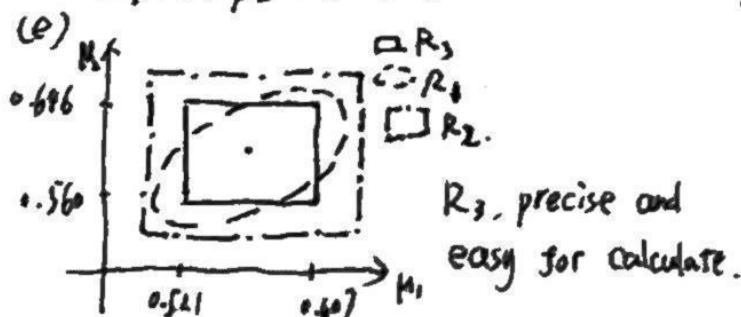
(b)  $m=3$   $P(a_i^T \mu \in C_i, i=1,2,3) = 1 - P(\text{at least 1 } i, a_i^T \mu \notin C_i) = 1 - P(a_i^T \mu \notin C_i)$   
 $= 1 - P(\cup \{a_i^T \mu \notin C_i\}) \geq 1 - \sum P(a_i^T \mu \notin C_i) = 1 - (\alpha_1 + \alpha_2 + \alpha_3) = 1 - \alpha$

(c) By loss (b), we have  $P(a_i^T \mu \in C_i(\frac{\alpha}{m}), \forall i) \geq 1 - (\frac{\alpha}{m} + \frac{\alpha}{m} + \dots) = 1 - \alpha$

(d)  $m=2$ ,  $I_{\mu_1} = 0.555 \pm (1,0) \begin{pmatrix} 0.564 \\ 0.603 \end{pmatrix} \pm t \left( \frac{0.05}{2 \times 2} \right) \cdot \left( \sqrt{\frac{(1.0) S(0)}{n}} \right) = 0.564 \pm 0.043$   
 $I_{\mu_1} = (0.521, 0.607)$

$I_{\mu_2} = (0,1) \begin{pmatrix} 0.564 \\ 0.603 \end{pmatrix} \pm 2.306 \cdot \sqrt{\frac{(1.0) S(0)}{42}} = 0.603 \pm 0.043 = (0.560, 0.646)$

$I_{\mu_1}, I_{\mu_2}$  are confidence interval of  $\mu$   $R_3 = \begin{pmatrix} I_{\mu_1} \\ I_{\mu_2} \end{pmatrix}$



Q4. ~~Regression~~

$SSR = \sum \left\| \tilde{X}_i - \frac{u' \tilde{X}_i}{u' u} u \right\|^2 = \sum \left( \tilde{X}_i^T \tilde{X}_i - \tilde{X}_i^T \frac{u' \tilde{X}_i}{u' u} u - \frac{u' \tilde{X}_i}{u' u} u^T \tilde{X}_i + \frac{u' \tilde{X}_i u u' \tilde{X}_i}{u' u u' u} \right)$   
 $= \sum \left( \tilde{X}_i^T \tilde{X}_i - 2 \frac{u' \tilde{X}_i \tilde{X}_i' u}{u' u} + \frac{u' \tilde{X}_i \tilde{X}_i' u}{u' u} \right)$   
 $= \sum (\tilde{X}_i^T \tilde{X}_i) - \sum \left( \frac{u' \tilde{X}_i \tilde{X}_i' u}{u' u} \right)$

$\max_{\min} SSR \Leftrightarrow \max \sum \left( \frac{u' \tilde{X}_i \tilde{X}_i' u}{u' u} \right)$ , i.e. max Sample Variance:  $\max \frac{1}{n-1} \sum (u' \tilde{X}_i - \frac{\sum u' \tilde{X}_i}{n})^2$

$\therefore \max \text{ Sample Var} \Leftrightarrow \min \text{ RSS}$   
 $= \frac{1}{n-1} \sum \left( u' \tilde{X}_i - \frac{u' \sum \tilde{X}_i}{n} \right)^2 = \frac{1}{n-1} \sum (u' \tilde{X}_i)^2$

# STAT 2221 Multivariate: HW 3

Ji Changcheng

Spring 2024

(a) Write down the equation for the logistic regression model of LI on remission in cancer patients (using the parameters  $\beta_0$  and  $\beta_1$ ).

$\pi(x) = \exp(\beta_0 + \beta_1 x) / (1 + \exp(\beta_0 + \beta_1 x))$  The parameters are obtained from R code result  
 $\beta_0 = -3.77714, \beta_1 = 0.14486$

## Question 1

```
df <- read.csv("homeless_smallldata.csv", header=T, row.names=1)
```

(a)

$H_0 : \mu_1 = \mu_2$   
 $H_\alpha : \mu_1 \neq \mu_2$

(b) Determine the dimension p and the sample sizes n1, n2.

```
cat("Dimension p:", ncol(df)-1, "\n")
```

```
## Dimension p: 3
```

```
cat("Sample size nonhomeless n1 (homeless=0):", sum(df$homeless == 0), "\n")
```

```
## Sample size nonhomeless n1 (homeless=0): 244
```

```
cat("Sample size homeless n2 (homeless=1):", sum(df$homeless == 1), "\n")
```

```
## Sample size homeless n2 (homeless=1): 209
```

(c) Conduct the test.

```
non_homeless_data <- df[df$homeless == 0, ]
homeless_data <- df[df$homeless == 1, ]
p <- 3
n1 <- nrow(non_homeless_data)
n2 <- nrow(homeless_data)

mean1 <- colMeans(non_homeless_data[, c(1, 2, 3)])
mean2 <- colMeans(homeless_data[, c(1, 2, 3)])
covmat <- cov(df[, c(1, 2, 3)])

cat("Means of non homeless: ", mean1, "\n")
```

```
## Means of non homeless:  49.00083 32.48683 31.84016
```

```
cat("Means of homeless: ", mean2, "\n")
```

```
## Means of homeless:  46.93678 30.73085 34.02392
```

```
cat("Covariance of vars: \n")
```

```
## Covariance of vars:
```

```
covmat
```

```
##           pcs           mcs           cesd
## pcs  116.30766   15.29467  -39.50384
## mcs   15.29467  164.84858 -109.56974
## cesd -39.50384 -109.56974  156.61170
```

```
tval <- (n1+n2-p-1)*(n1*n2)/((n1+n2)*p*(n1+n2-2))* t(mean1-mean2) %%% solve(covmat) %%% (mean1-mean2)
cat("T statistic = ",tval,"\n")
```

```
## T statistic = 2.012177
```

```
p_value <- 1 - pf(tval, p, n1+n2)
cat("p value = ",p_value,"\n")
```

```
## p value = 0.1114513
```

p-value > .05, thus we can not reject the H0 that they have the same means.

#### (d) Scatter points

```
library(dplyr)
```

```
## Warning:  程辑包‘dplyr’是用R版本4.3.2 来建造的
```

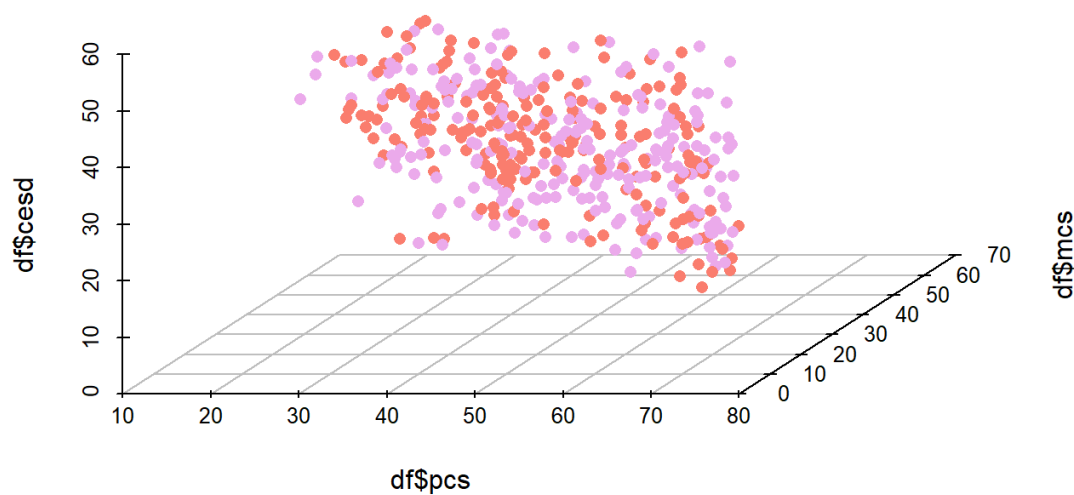
```
##
## 载入程辑包：‘dplyr’
```

```
## The following objects are masked from ‘package:stats’:
##
##      filter, lag
```

```
## The following objects are masked from ‘package:base’:
##
##      intersect, setdiff, setequal, union
```

```
library(scatterplot3d)
colors <- c("plum2", "salmon")
colors <- colors[as.factor(df$homeless)]
scatterplot3d(df$pcs, df$mcs, df$cesd,pch = 16, main = "3D Scatter Plot", angle = 45, scale.y = 0.7, pty = "s", box = FALSE,
              color=colors)
```

## 3D Scatter Plot



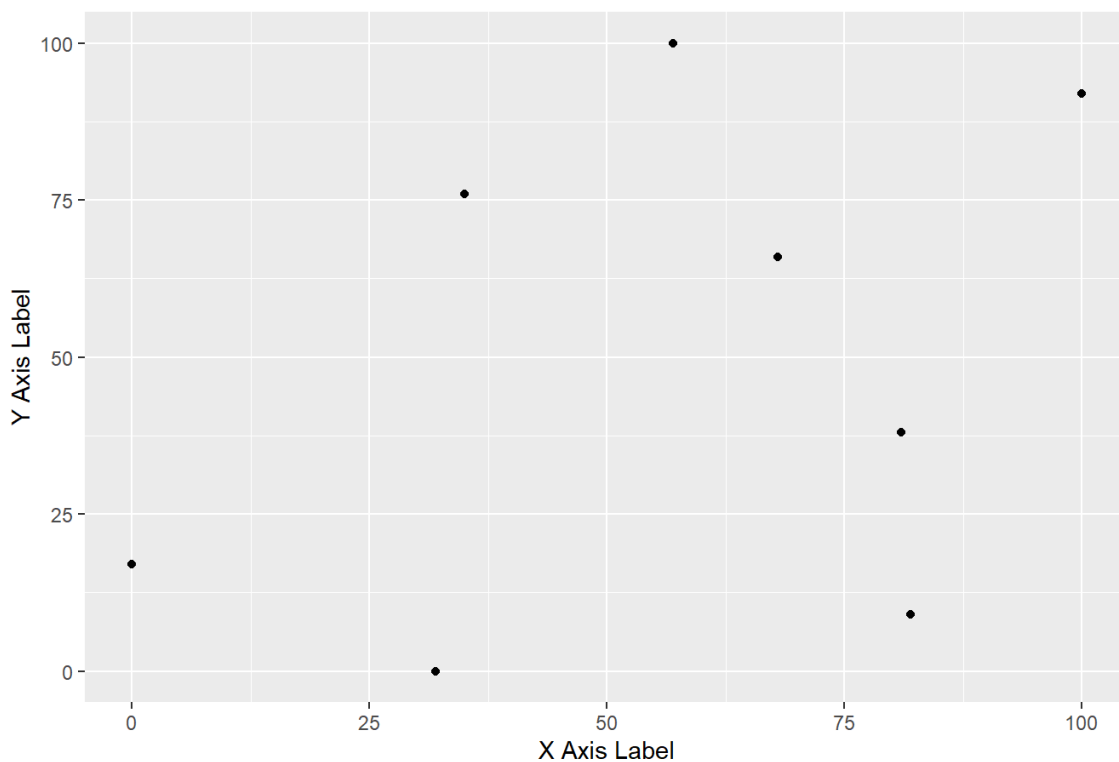
Data are generally distributed in the same way.

## Question 4

(a)

```
library(ggplot2)
df2 <- read.table("pendigit3.txt", sep = ",")
df2 <- df2[, -17]
obs<-data.frame(x = as.numeric(df2[1, seq(1, 16, by = 2)]), y = as.numeric(df2[1, seq(2, 16, by = 2)]))
ggplot(obs, aes(x = x, y = y)) +
  geom_point() +
  ggtitle("Scatter Plot of All Points") +
  xlab("X Axis Label") +
  ylab("Y Axis Label")
```

### Scatter Plot of All Points



(b)

```
scaled_data <- scale(df2)
```

```
# Step 2: Perform PCA
```

```
pca_result <- prcomp(scaled_data)
```

```
v <- pca_result$sdev^2
```

```
variance_explained <- round(cumsum(v/sum(v)), 3)
```

```
pca_data <- data.frame(
```

```
  PC = 1:length(variance_explained),
```

```
  VarianceExplained = variance_explained
```

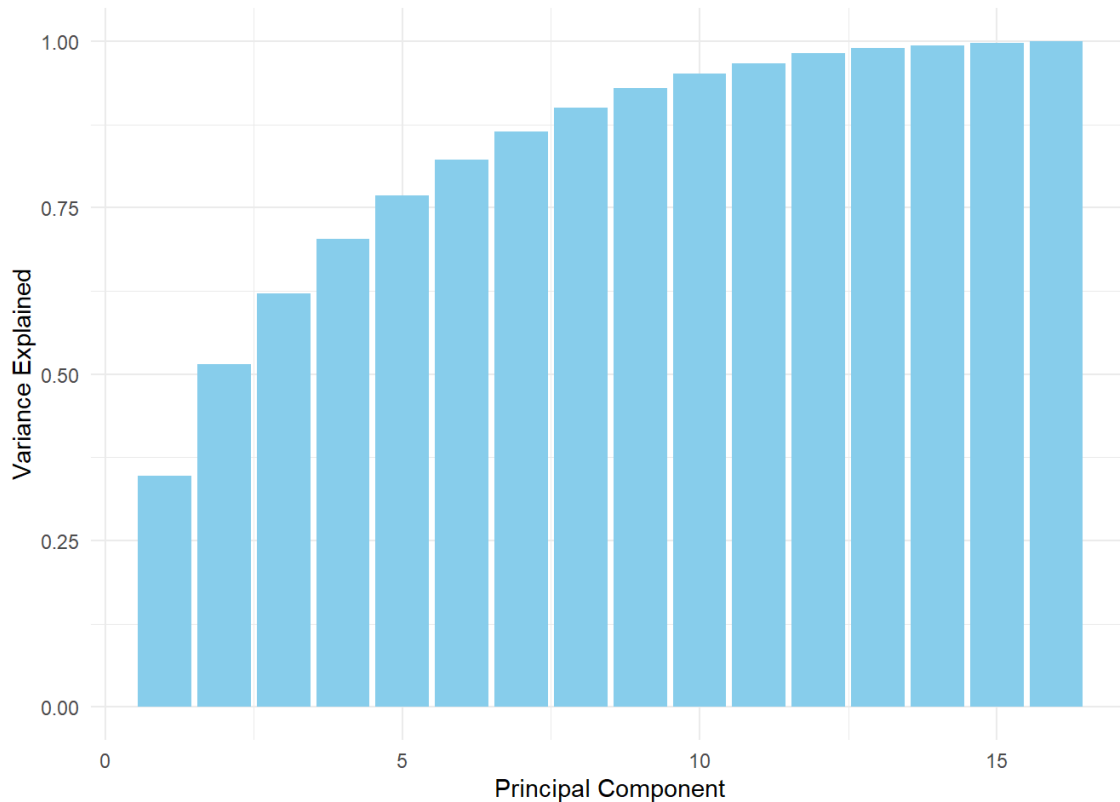
```
)
```

```
ggplot(pca_data, aes(x = PC, y = variance_explained)) +
```

```
  geom_bar(stat = "identity", fill = "skyblue") +
```

```
  labs(x = "Principal Component", y = "Variance Explained") +
```

```
  theme_minimal()
```

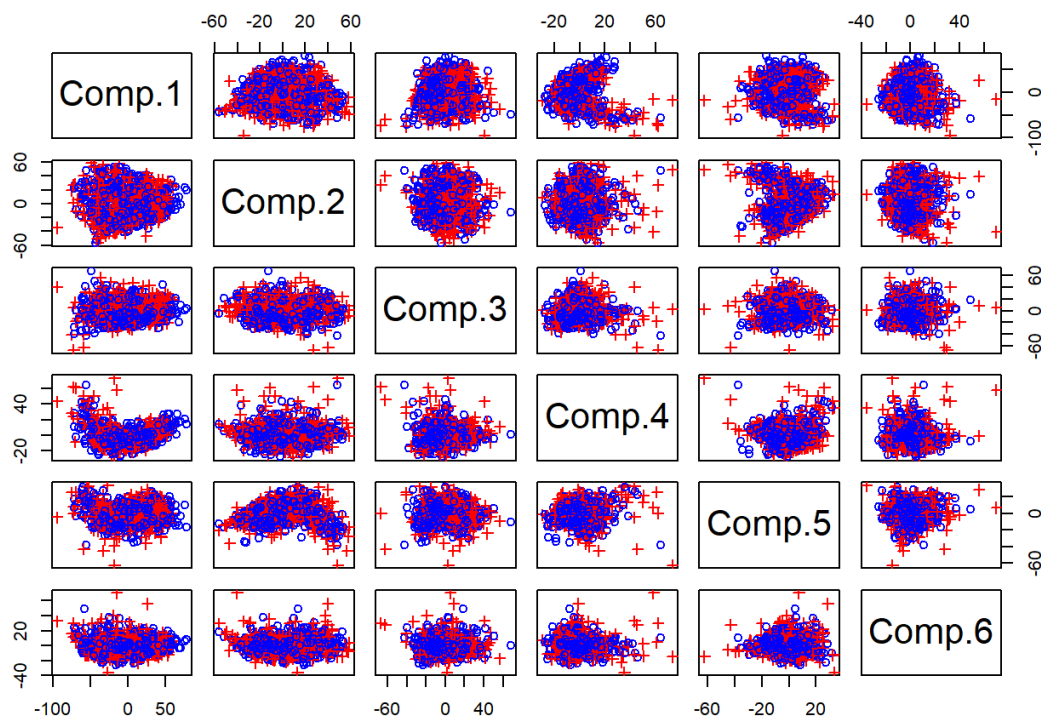


```
pca_result<-princomp(df2)
```

```
Z <-pca_result$scores
```

```
# scatterplot of principal components
```

```
pairs(Z[, c(1, 2, 3, 4, 5, 6)], pch=c(rep(1, 100), rep(3, 100)), col=c(rep("blue", 100), rep("red", 100)))
```

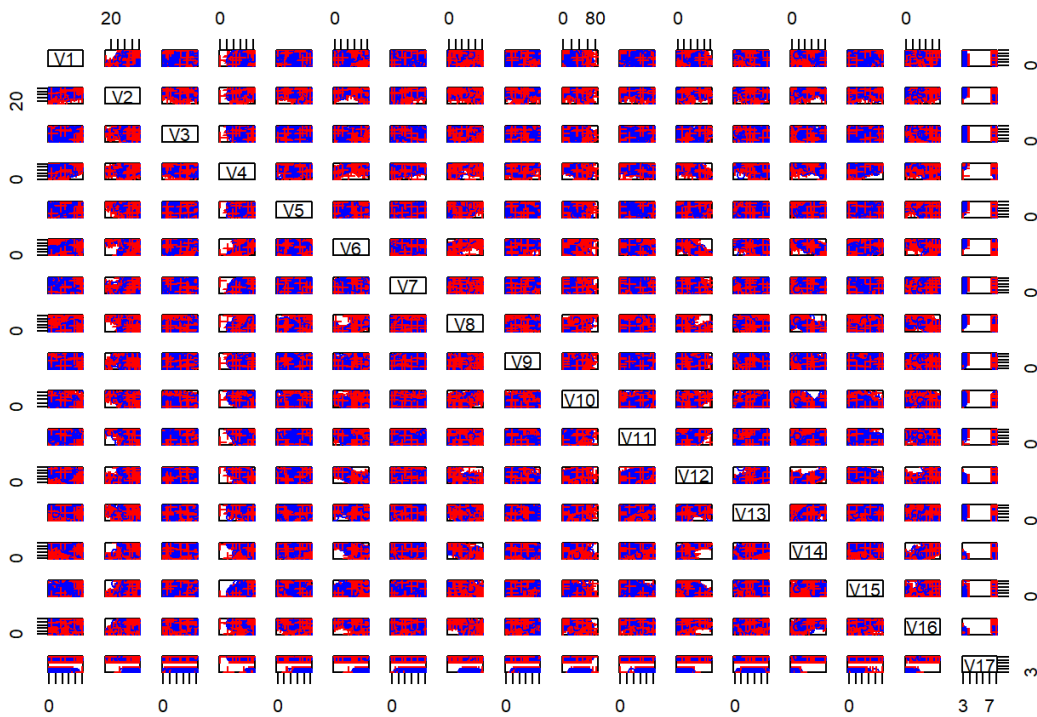


- (c) The result shows strong randomness so that it is likely to be a MVN distribution
- (d) I prefer to keep 6 as the explained variance reach 80%
- (e) PC1 is the linear combination that captures the maximum variance in the data. Along the direction of PC1, the variability of the data is maximized. PC2, orthogonal to PC1, have the next highest variance. Along the direction of PC2, the variability of the data is second highest.
- (f) From the graph, there is a nice separation between “3” and “8”.

```
df3 <- rbind(read.table("pendigit3.txt", sep = ","), read.table("pendigit8.txt", sep = ","))
```

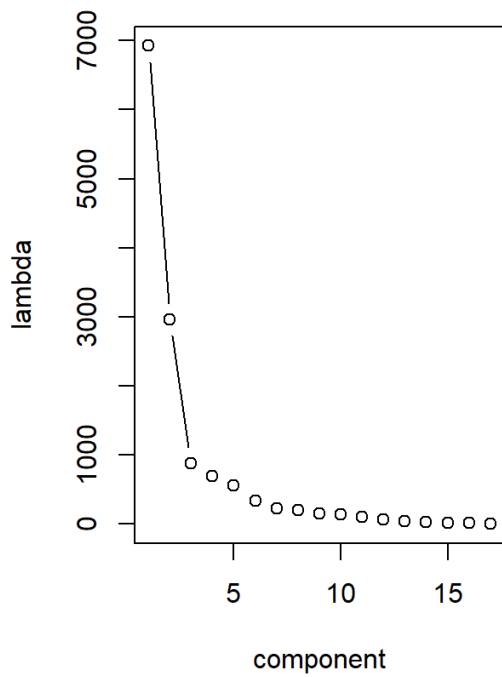
```
spr <- princomp(df3)
U<-spr$loadings
L<-(spr$sdev)^2
Z <-spr$scores

# scatterplot of principal components
pairs(df3,pch=c(rep(1,100),rep(3,100)),col=c(rep("blue",100),rep("red",100)))
```

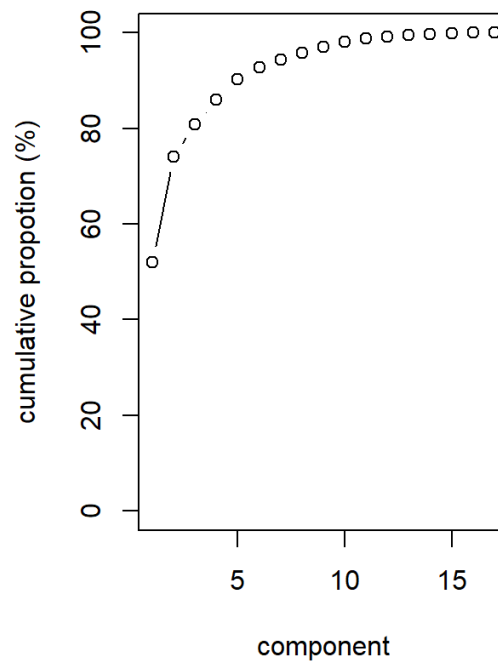


```
par(mfrow=c(1,2))
plot(L, type="b", xlab="component", ylab="lambda", main="Scree plot")
plot(cumsum(L)/sum(L)*100, ylim=c(0,100), type="b", xlab="component", ylab="cumulative propotion (%)", main="Cum. Scree plot")
```

**Scree plot**



**Cum. Scree plot**



```
# biplot
par(mfrow=c(1,1))
```