

Supervised Dimension Reduction

S. Tonini, F. Chiaromonte (special thanks to J. Di Iorio, L. Insolia and L. Testa)

March 13th 2024

Contents

Introduction	1
Libraries	1
Data	1
PCA vs LDA	3
PCA	3
LDA	3
SIR	5

Introduction

Libraries

We are going to use:

```
library(tidyverse) # data manipulation and visualization
library(plotly)    # plots in 3D
library(ggplot2)   # plots in 2D
library(ggpubr)     # to combine multiple ggplot objects (ggarrange)
library(mvtnorm)    # to generate multivariate normal distribution
library(dr)         # SIR
library(factoextra) # PCA-related functions
```

Data

Let's first define a function to generate Gaussian data. This function takes four arguments:

- n: number of observations;
- center: the mean vector
- sigma: the covariance matrix
- label: the cluster label

```
generateGaussianData <- function(n, center, sigma, label) {
  data = rmvnorm(n, center, sigma)
  data = data.frame(data)
  names(data) = c("x", "y", "z")
  data = data %>% mutate(class=factor(label))
  data
}
```

Now let's simulate a dataset.

```
covmat <- matrix(c(1,0.88,0.88,0.88, 1,0.88,0.88,0.88, 1),
  nrow = 3, byrow=T)
```

```
# cluster 1
n = 200
center = c(2, 8, 6)
sigma = covmat
group1 = generateGaussianData(n, center, sigma, 1)
```

```
# cluster 2
n = 200
center = c(4, 8, 6)
sigma = covmat
group2 = generateGaussianData(n, center, sigma, 2)
```

```
# cluster 3
n = 200
center = c(6, 8, 6)
sigma = covmat
group3 = generateGaussianData(n, center, sigma, 3)
```

```
# all data
df = bind_rows(group1, group2, group3)
```

```
head(df)
```

```
##           x           y           z class
## 1 2.2435646 8.410348 6.158682      1
## 2 2.0082015 7.583865 6.554732      1
## 3 2.4795629 8.939585 6.602703      1
## 4 1.3318241 7.032033 4.436121      1
## 5 0.6852441 6.884413 5.064675      1
## 6 2.5598488 8.091716 6.658554      1
```

```
summary(df)
```

```
##           x           y           z      class
## Min.      :-1.326   Min.      : 4.658   Min.      :2.620   1:200
## 1st Qu.: 2.398     1st Qu.: 7.260     1st Qu.:5.224   2:200
## Median : 3.897     Median : 8.025     Median :5.945   3:200
## Mean    : 3.939     Mean    : 7.964     Mean     :5.932
## 3rd Qu.: 5.487     3rd Qu.: 8.667     3rd Qu.:6.645
## Max.     : 8.647     Max.     :10.929    Max.      :8.852
```

And plot our simulated data.

```
fig <- plot_ly(df, x = ~x, y = ~y, z = ~z,
  color = ~class, colors = c('#b3e378', '#81e5f0', '#ed5391'))
fig <- fig %>% add_markers()
fig <- fig %>% layout(scene = list(xaxis = list(title = 'x'),
  yaxis = list(title = 'y'),
  zaxis = list(title = 'z'))))
fig
```

PCA vs LDA

PCA

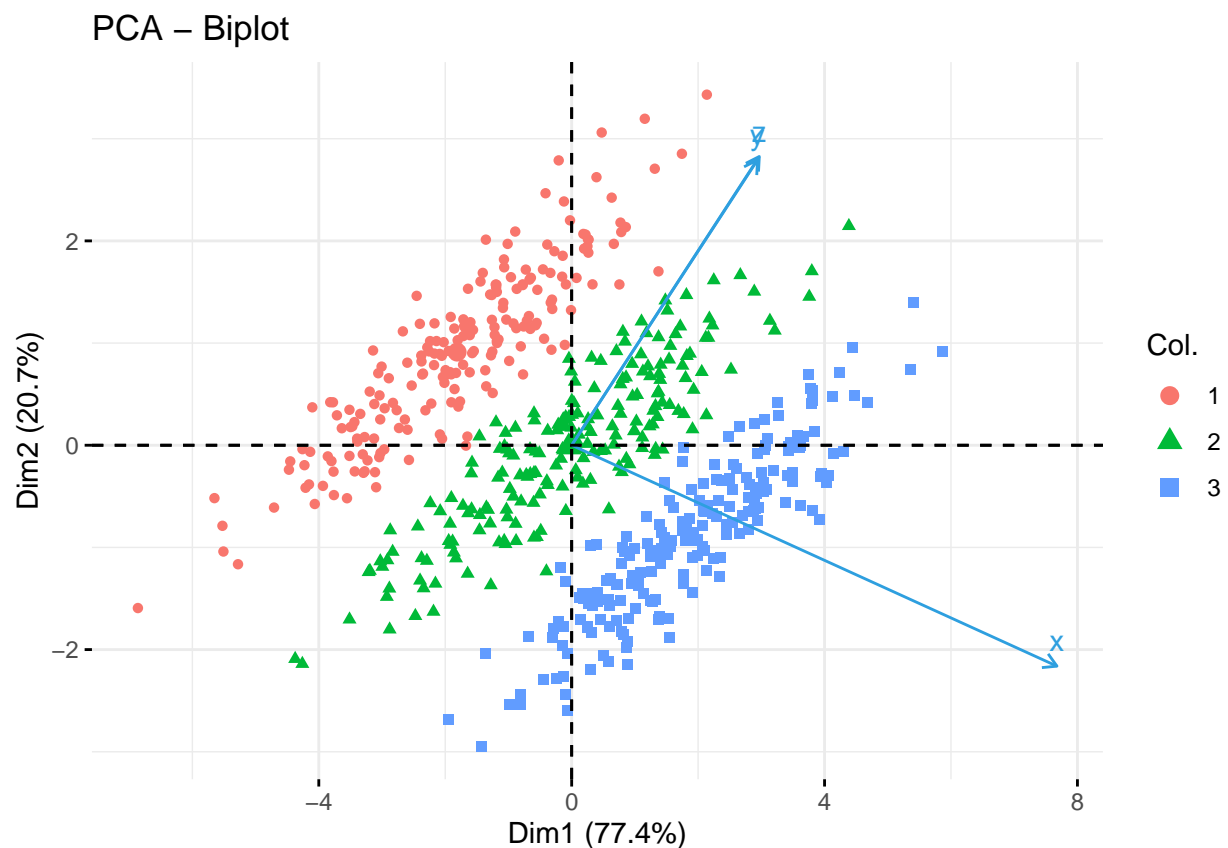
Now let us perform PCA.

```
pc <- prcomp(df[,c(1,2,3)])  
get_eig(pc)
```

```
##      eigenvalue variance.percent cumulative.variance.percent  
## Dim.1  4.5469471       77.399980           77.39998  
## Dim.2  1.2180285       20.733776           98.13376  
## Dim.3  0.1096345        1.866244          100.00000
```

This is the corresponding biplot.

```
fviz_pca_biplot(pc, col.var= "#2E9FDF", col.ind= df$class, label="var")
```



Note that just considering the first principal component it is impossible to notice differences within the three groups (all groups are overlapping).

LDA

Let's perform LDA:

```
lda.df <- lda(factor(class) ~ x + y + z, data = df)  
lda.df
```

```
## Call:  
## lda(factor(class) ~ x + y + z, data = df)
```

```
##
## Prior probabilities of groups:
##      1      2      3
## 0.3333333 0.3333333 0.3333333
##
## Group means:
##      x      y      z
## 1 1.910877 7.919929 5.883410
## 2 3.906692 7.942222 5.901631
## 3 5.998491 8.029468 6.010456
##
## Coefficients of linear discriminants:
##      LD1      LD2
## x  2.4840633 -0.02966442
## y -1.2855751 -0.77831989
## z -0.9877422  1.62387594
##
## Proportion of trace:
## LD1 LD2
##   1   0
```

Let us plot the projections on LD1 and LD2

```
# prediction on df to get projections
predmodel.lda = predict(lda.df, data=df)

# projections with LDA classes
estclass <- as.factor(apply(predmodel.lda$posterior, 1, which.max))
newdata2 <- data.frame(type = estclass, lda = predmodel.lda$x)
p1 <- ggplot(newdata2) +
  geom_point(aes(lda.LD1, lda.LD2, colour = type), size = 2.5) +
  ggtitle("projections with LDA classes")

# projections with true classes
newdata <- data.frame(type = df$class, lda = predmodel.lda$x)
p2 <- ggplot(newdata) +
  geom_point(aes(lda.LD1, lda.LD2, colour = type), size = 2.5) +
  ggtitle("projections with true classes")

ggarrange(p1,p2,
  nrow=2)
```



SIR

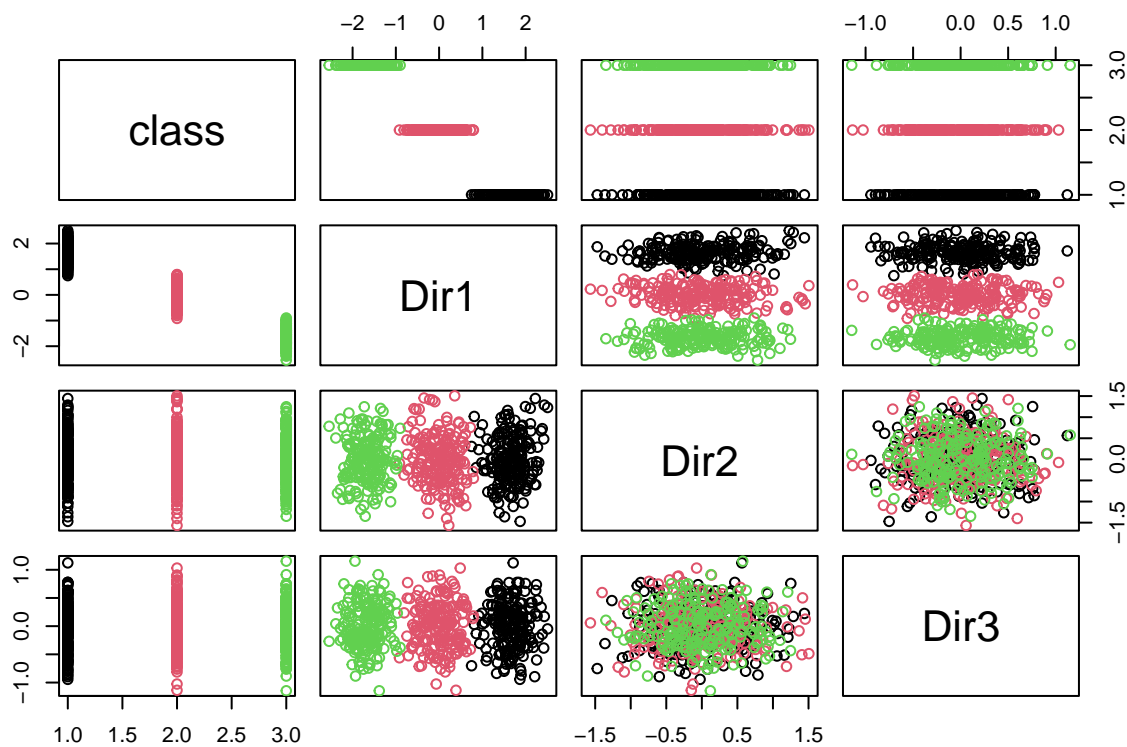
Now we use the SIR (Sliced Inversion Regression) in the `dr` package

```
# default fitting method is "sir"
help(dr)

dr_res <- dr(class ~ x + y + z, data = df, method='sir')

dr_res

##
## dr(formula = class ~ x + y + z, data = df, method = "sir")
## Estimated Basis Vectors for Central Subspace:
##      Dir1      Dir2      Dir3
## x -0.8374295 -0.0164710  0.00380392
## y  0.4333941 -0.4321576 -0.81462206
## z  0.3329885  0.9016477  0.57997968
## Eigenvalues:
## [1] 9.424499e-01 4.902462e-04 2.616200e-17
plot(dr_res, col=df$class)
```

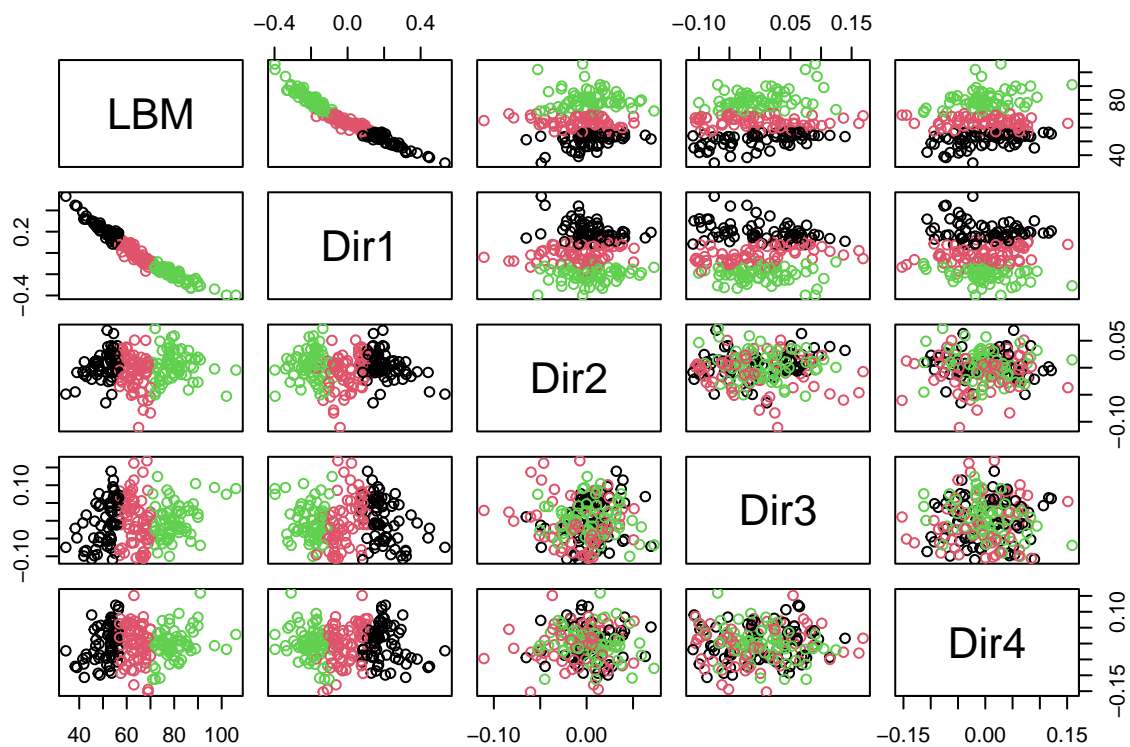


```
names(dr_res)
```

```
## [1] "x"          "y"          "weights"    "method"     "cases"
## [6] "qr"         "group"      "chi2approx" "evectors"   "evalues"
## [11] "numdir"     "raw.evectors" "M"         "slice.info" "call"
## [16] "y.name"     "terms"
```

We perform SIR on real data with continuos outcome

```
data(ais)
?ais
dr_res3 <-dr(LBM~log(SSF)+log(Wt)+log(Hg)+log(Ht)+log(WCC)+log(RCC)+
             log(Hc)+log(Ferr),data=ais, nslices=3)
plot(dr_res3, col=dr_res3$slice.info$slice.indicator)
```



```
dr_res6 <-dr(LBM~log(SSF)+log(Wt)+log(Hg)+log(Ht)+log(WCC)+log(RCC)+
             log(Hc)+log(Ferr),data=ais, nslices=6)
plot(dr_res6, col=dr_res6$slice.info$slice.indicator)
```

