

ETC3250: Regularisation and Shrinkage

Week 9, class 2

Professor Di Cook, Econometrics and
Business Statistics

- When the number of variables (p) is large, estimating a model is problematic.
- Particularly, when it is larger than the sample size ($p \gg n$), the variance of an estimate could be ∞ .
- Constraining, or shrinking the estimates, can substantially decrease the variance, while minimally affecting the bias.

- Subset selection: Fit models to best subset
- Dimension reduction: Use combinations of variables, e.g. PCs, and feed these into your model

- Modified least squares

$$\sum_{i=1}^n (y_i - b_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \beta_j^2$$

- where λ is a tuning parameter
- Minimizing this quantity trades off error with small β 's, at least forcing some of them to be small

- More recent alternative to ridge regression

$$\sum_{i=1}^n (y_i - b_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- the change using an l_1 error, really forces some of the coefficients to be 0.

Simulation example



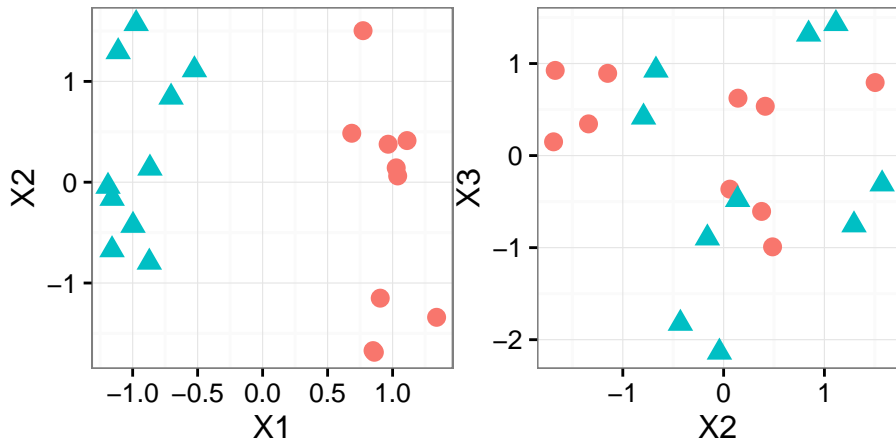
```
x<-matrix(rnorm(20*100),ncol=100)
x[1:10,1]<-x[1:10,1]+5
x<-scale(x)
x<-data.frame(x, cl=c(rep("A",10),rep("B",10)))
library(ggplot2)
qplot(X1,X2,data=x,colour=cl, size=I(3), shape=cl) +
  theme_bw() + theme(legend.position="None", aspect.ratio=1)
```

```
qplot(X2,X3,data=x,colour=cl, size=I(3), shape=cl) +
  theme_bw() + theme(legend.position="None", aspect.ratio=1)
```

Generate test data

```
x.t<-matrix(rnorm(10*100),ncol=100)
x.t[1:5,1]<-x.t[1:5,1]+5
x.t<-scale(x.t)
x.t<-data.frame(x.t, cl=c(rep("A",5),rep("B",5)))
```

Simulation example



```
## Call:
## lda(c1 ~ ., data = x[, c(1:2, 101)], prior = c(0.5, 0.5))
##
## Prior probabilities of groups:
##      A      B
## 0.5 0.5
##
## Group means:
##           X1           X2
## A  0.9555333 -0.2864578
## B -0.9555333  0.2864578
##
## Coefficients of linear discriminants:
##           LD1
## X1 -4.9308578
## X2  0.1098234
```


##

A B

A 10 0

B 0 10

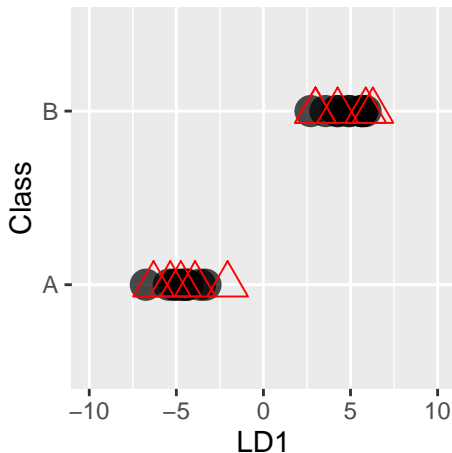
##

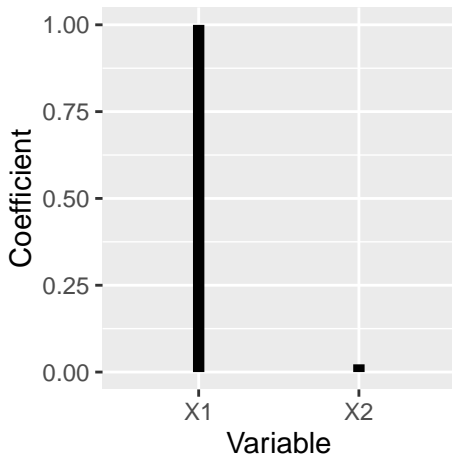
A B

A 5 0

B 0 5

Plot training data and test in discriminant space





- The next few slides repeat the results just shown for increasing number of variables
- None of the additional variables contribute to the separation between classes
- Additional variables are purely noise

$p=5$

##

A B

A 10 0

B 0 10

##

A B

A 5 0

B 0 5

$p=8$

##

A B

A 10 0

B 0 10

##

A B

A 5 0

B 0 5

$p=11$

##

A B

A 10 0

B 0 10

##

A B

A 5 0

B 0 5

$p=12$

##

A B

A 10 0

B 0 10

##

A B

A 5 0

B 0 5

$p=13$

##

A B

A 10 0

B 0 10

##

A B

A 4 1

B 0 5

$p=14$

##

A B

A 10 0

B 0 10

##

A B

A 4 1

B 0 5

$p=15$

##

A B

A 10 0

B 0 10

##

A B

A 4 1

B 1 4

$p=16$

##

A B

A 10 0

B 0 10

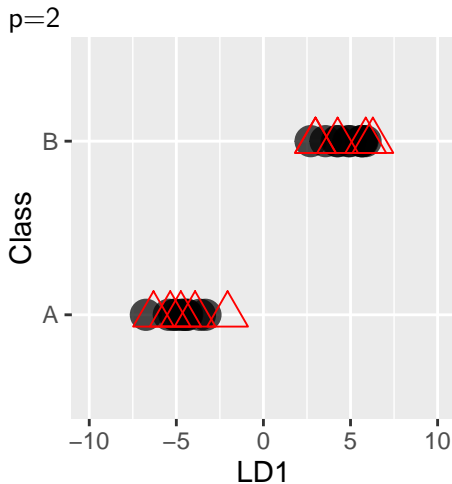
##

A B

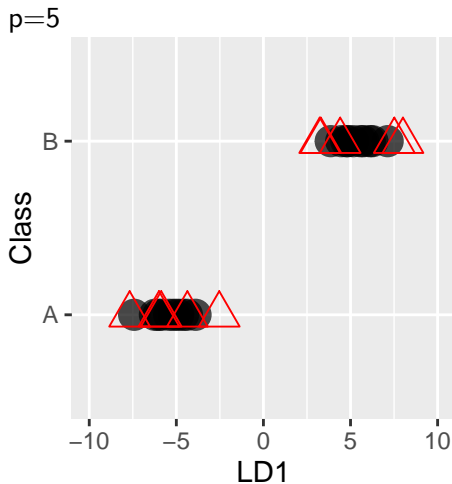
A 4 1

B 0 5

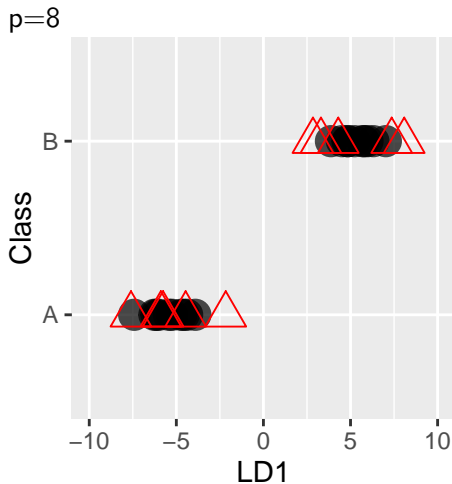
Plot training data and test in discriminant space



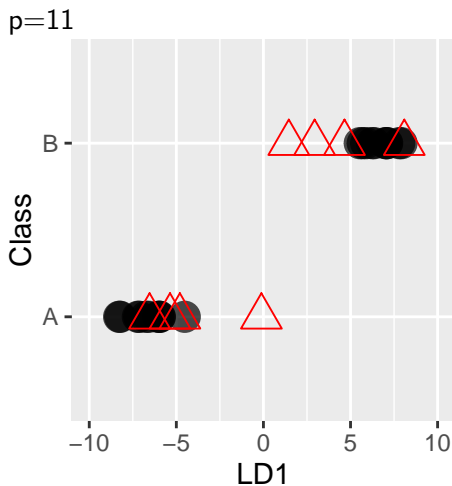
Plot training data and test in discriminant space



Plot training data and test in discriminant space



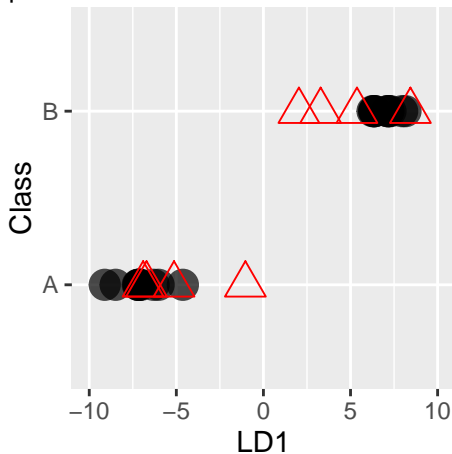
Plot training data and test in discriminant space



Plot training data and test in discriminant space



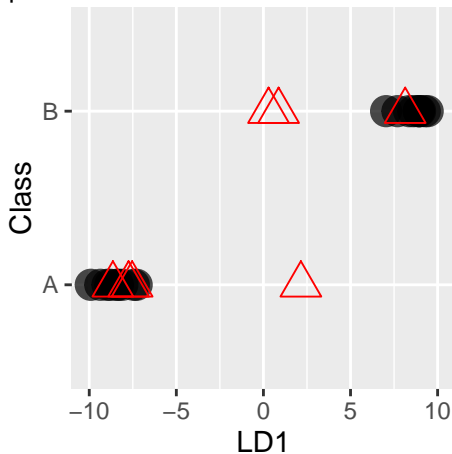
$p=12$



Plot training data and test in discriminant space



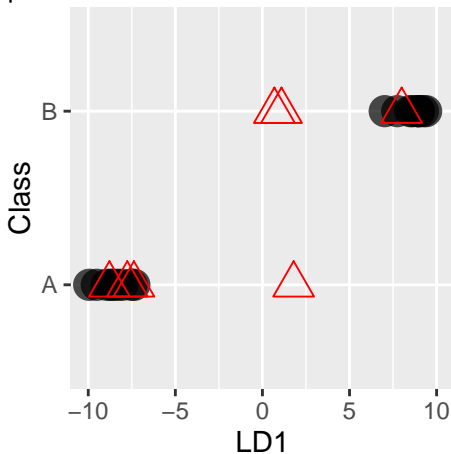
$p=13$



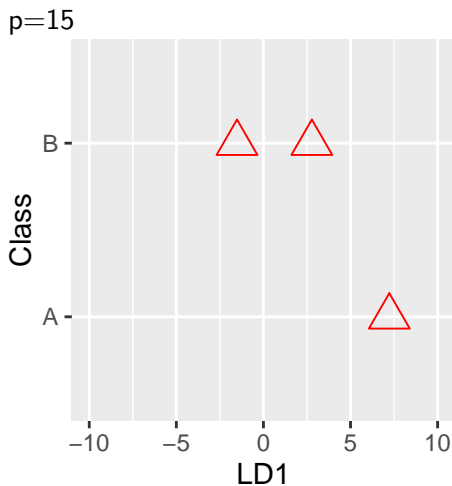
Plot training data and test in discriminant space

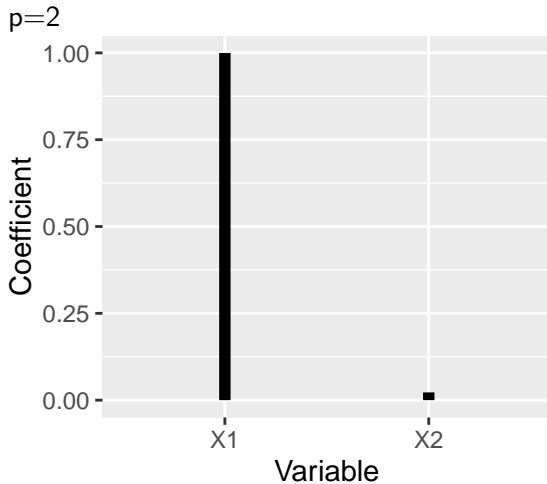


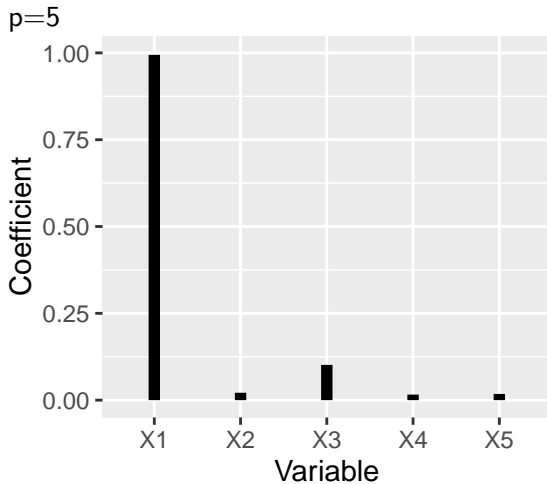
$p=14$

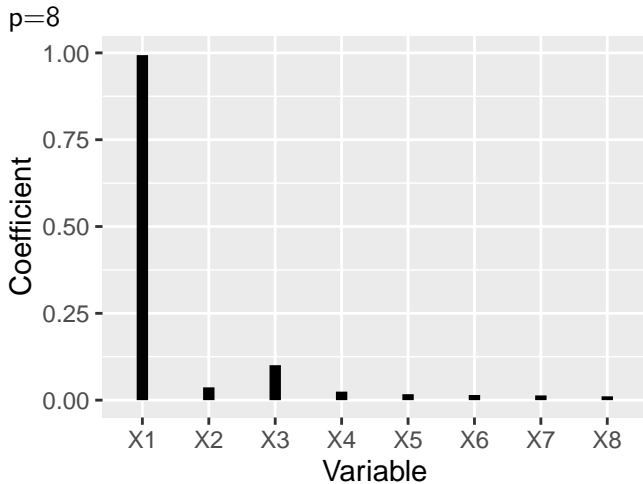


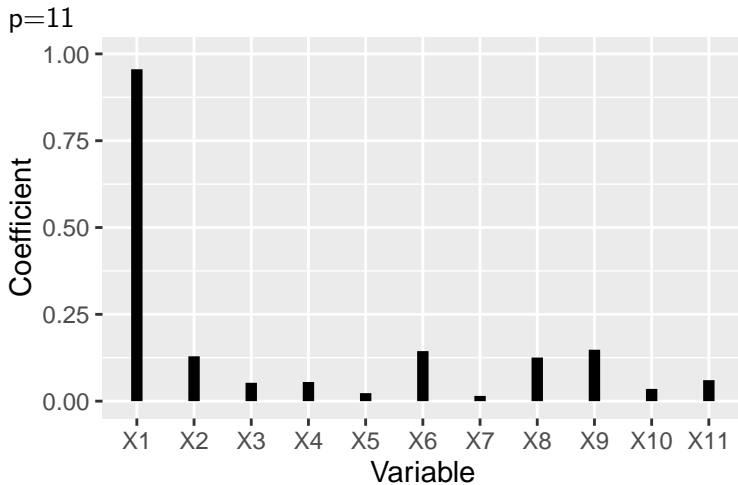
Plot training data and test in discriminant space

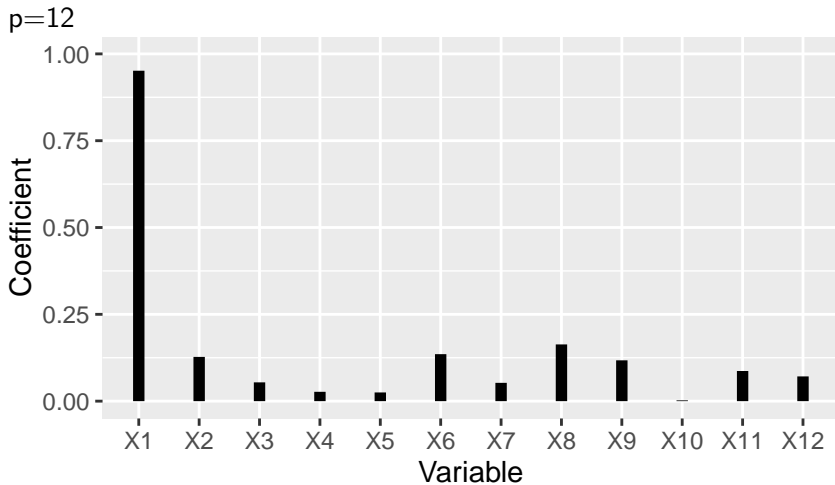


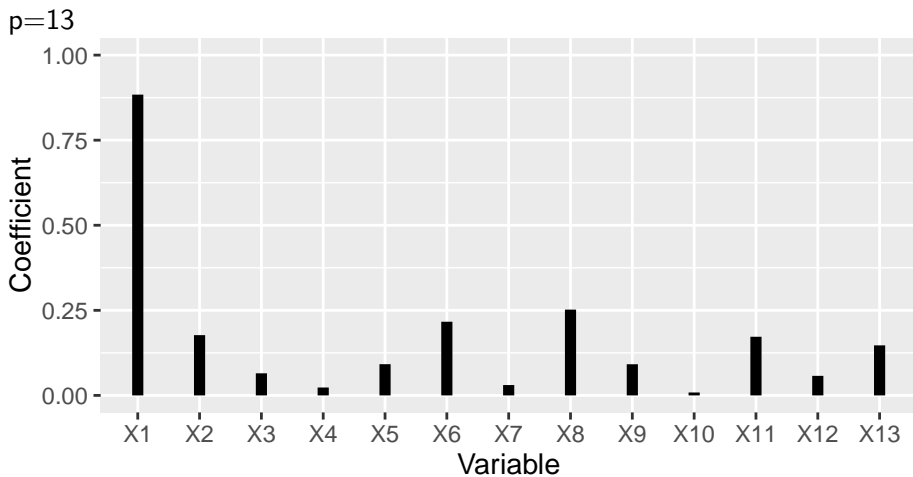


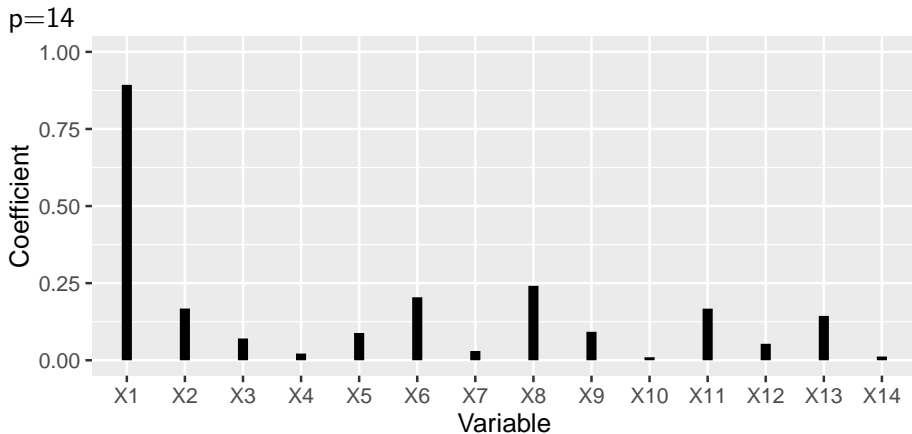


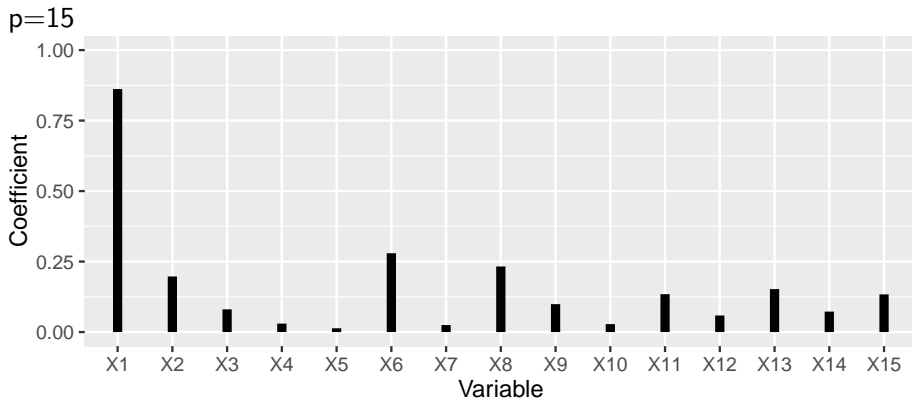


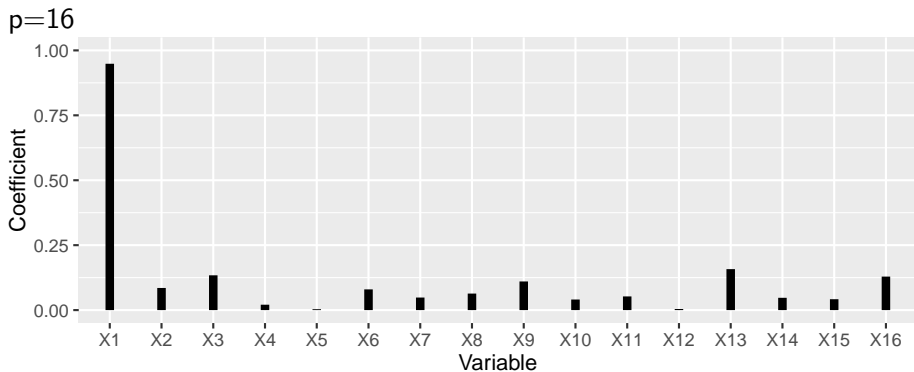












- Reference: <http://faculty.washington.edu/dwitten/Papers/JRSSBPenLDA.pdf>
- LDA does an eigendecomposition of $W^{-1}B$, which creates an estimation problem if $\ll p$
- Instead compute regularised versions of W, B

$$\text{maximize}_{\beta_k} \beta_k^T \hat{\Sigma}_b^k \beta_k - \lambda_k \sum_{j=1}^p |\hat{\sigma}_k \beta_{kj}|$$

subject to

$$\beta_k^T \hat{\Sigma}_w \beta_k \leq 1,$$

where $\hat{\Sigma}_b^k = \frac{1}{n} X^T Y (Y^T Y)^{-1/2} P_k^\perp (Y^T Y)^{-1/2} Y^T X$, $\hat{\Sigma}_w$ is a positive definite estimate of Σ_w .

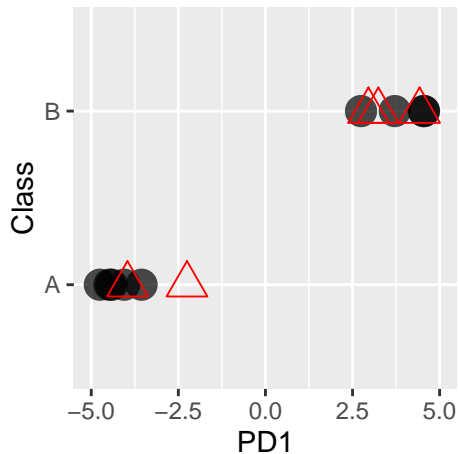
```
library(penalizedLDA)
cv.out<-PenalizedLDA.cv(as.matrix(x[,c(1:15)]), as.numeric(x$y))

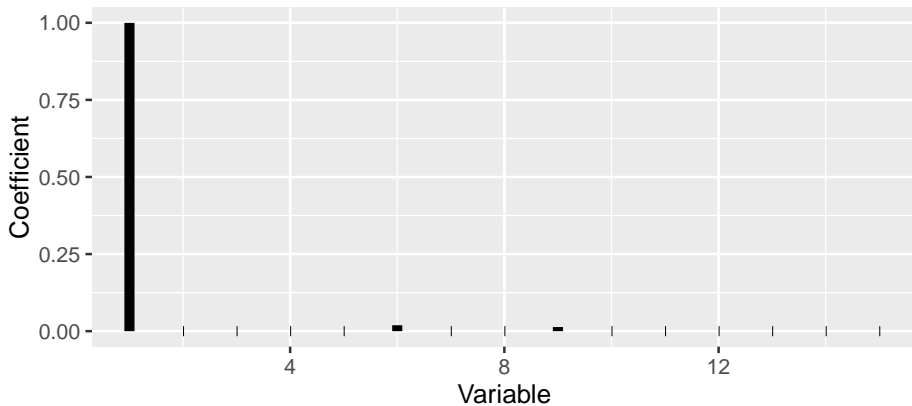
## Fold 1
## 12345Fold 2
## 12345Fold 3
## 12345Fold 4
## 12345Fold 5
## 12345Fold 6
## 12345

x.pda<-PenalizedLDA(as.matrix(x[,c(1:15)]), xte=as.matrix(x.t),
table(x.t$c1, x.pda$ypred)

##
##      1 2
##    A 5 0
```

Plot training and test





This work is licensed under the Creative Commons Attribution-Noncommercial 3.0 United States License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/3.0/us/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.