

Multivariate Regression und Klassifikation mit Anwendungen aus der Chemometrie

Peter Filzmoser

**Institut für Statistik und Wahrscheinlichkeitstheorie
Technische Universität Wien**

Herbstseminar des WBS, Wien

21. November 2013



Vienna University of Technology

Overview

321 + xiii pages; around 100 figures; appeared in 2009.

Written by a **chemometrician** and a **statistician**, reflects both the **practical** approach to chemometrics and the more **formally oriented** one of statistics.

Chapter 1: Introduction

Chapter 2: Multivariate Data

Chapter 3: Principal Component Analysis

Chapter 4: Calibration

Chapter 5: Classification

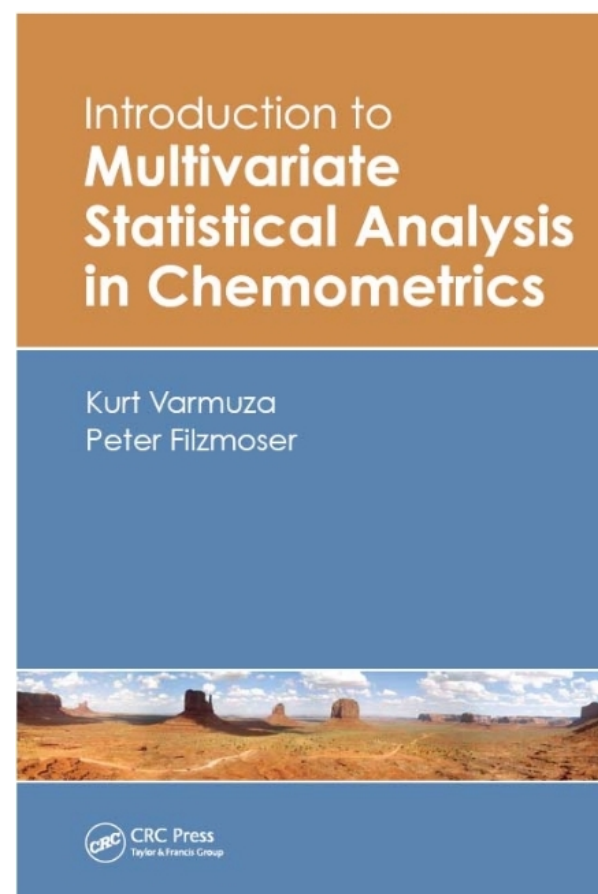
Chapter 6: Cluster Analysis

Chapter 7: Preprocessing

Appendix 1: Symbols and Abbreviations

Appendix 2: Matrix Algebra

Appendix 3: Introduction to R



Consider the **multiple linear regression model**

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i \quad \text{for } i = 1, \dots, n,$$

with

$$\mathbf{y} = (y_1, \dots, y_n)^\top$$

$$\mathbf{x}_i^\top = (1, x_{i1}, \dots, x_{ip}), \text{ forming the rows of } \mathbf{X},$$

$$\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^\top.$$

The **ordinary least-squares solution**

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

cannot be computed if $\mathbf{X}^\top \mathbf{X}$ is (near) singular. This happens typically in case of **multicollinearity** or $n < p$.

Way out: Dimension reduction in the X space.

Assume that the $n \times p$ matrix X is mean-centered.

Representation with the first k **principal components (PCs)**:

$$X = ZV^{\top} + E$$

with Z the $n \times k$ matrix of the first k PCs

V the $p \times k$ matrix with the first k PC loadings

E the error matrix (“reconstruction error”)

Principal Component Regression:

$$y = Z\theta + \varepsilon$$

with new regression coefficients θ and errors ε .

Important: “tuning” parameter k

PCR is a “technical” solution to the problem of multicollinearity:

- The PCs Z contain the “most important” information of X . But is this also the most important part for explaining y ?
- Why take the first k PCs, and not a “variable selection” of PCs?
- The real problem is to find that information of X that is best suited for explaining y .

Solution: Partial Least Squares (PLS) regression

Use again a “**latent variable model**” to explain X :

$$X = TA^T + E$$

with *scores* T ($n \times k$)
and *loadings* A ($p \times k$)

The columns of A are derived by

$$a_j = \underset{a}{\operatorname{argmax}} \operatorname{Cov}(y, Xa)$$

($j = 1, \dots, k$).

Then we have the **new regression model**:

$$\begin{aligned} y &= T\gamma + \varepsilon = (XA)\gamma + \varepsilon \\ &= X \underbrace{(A\gamma)}_{\beta} + \varepsilon = X\beta + \varepsilon \end{aligned}$$

Important for PLS: selection of the “optimal” number k of PLS components.

- k too small: underfit; poor prediction ability of the model
- k too large: overfit, and again poor prediction ability of the model

Model evaluation requires **training**, **validation** and **test** data.

Since we have **only one data set** available, we have to split appropriately.

Popular in chemometrics: **Standard Error of Prediction (SEP):**

$$\text{SEP} = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (y_i - \hat{y}_i - \text{bias})^2}$$

with

$$\text{bias} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)$$

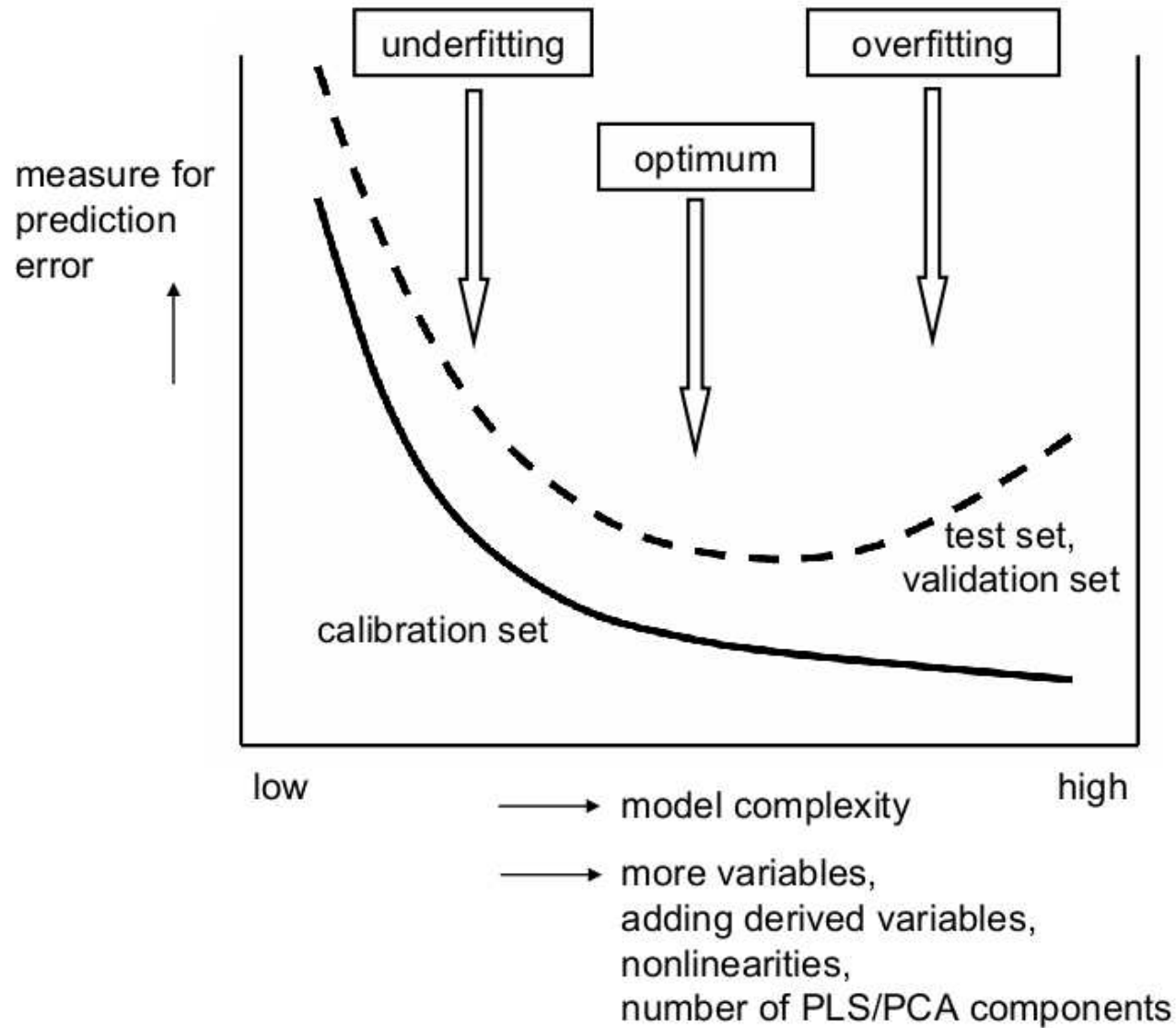
N is the number of test set observations.

Popular in many fields: **Mean Squared Error of Prediction (MSEP):**

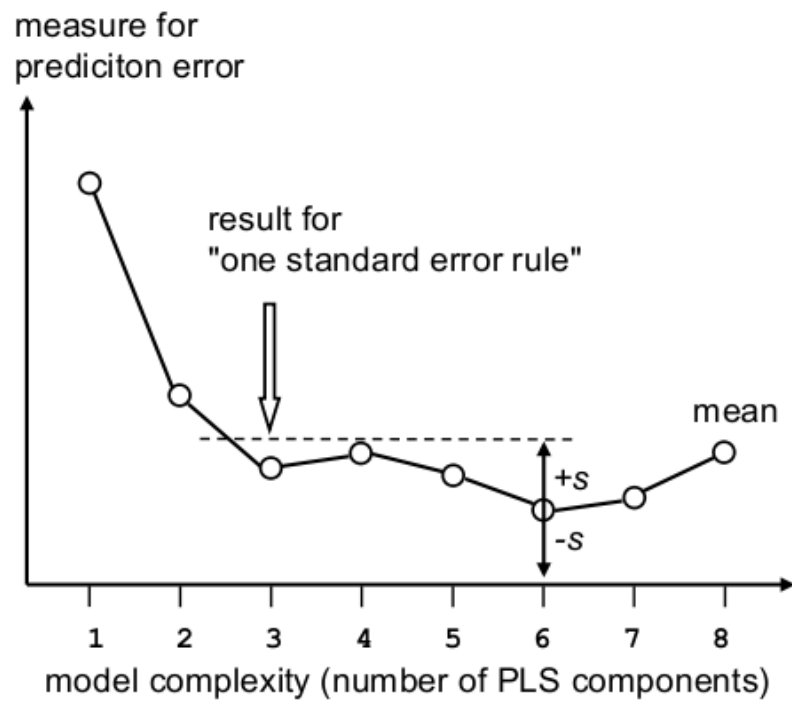
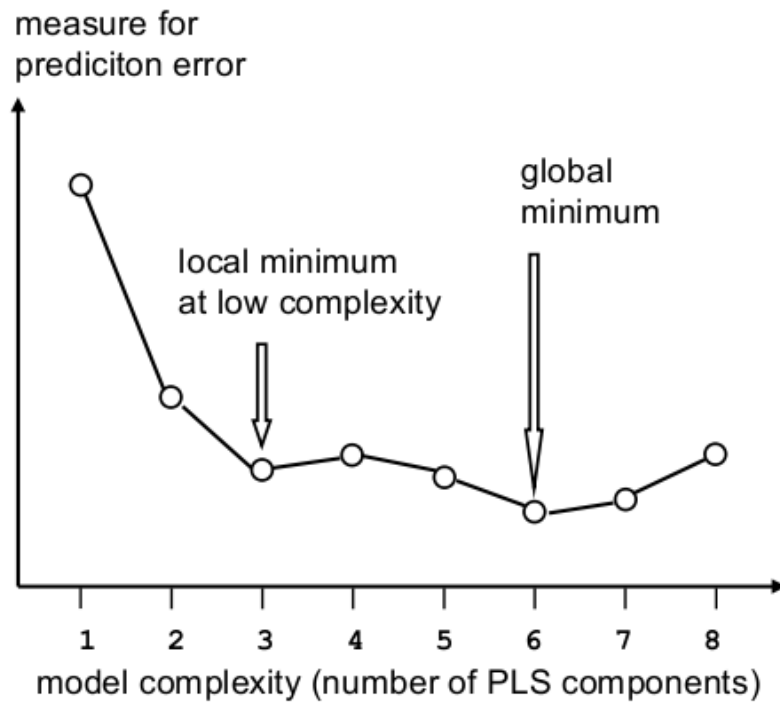
$$\text{MSEP} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

We see that $\text{SEP}^2 = \text{MSEP} - \text{bias}^2$

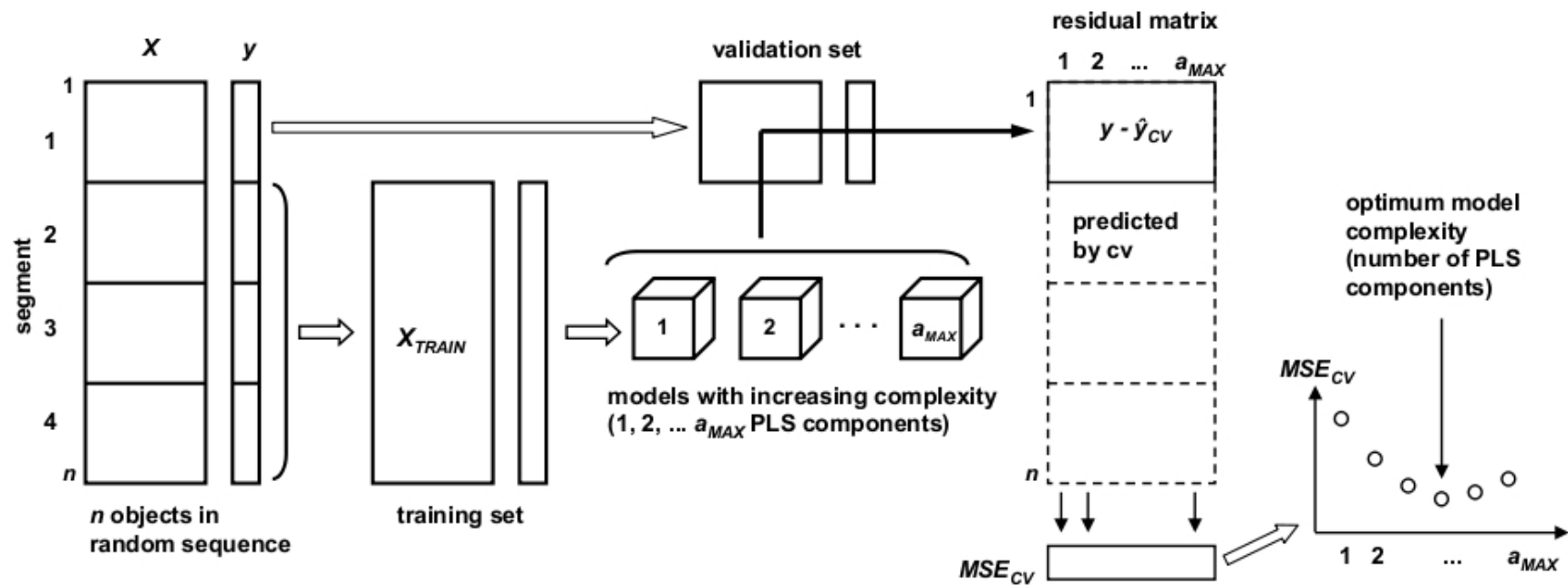
Optimal model complexity



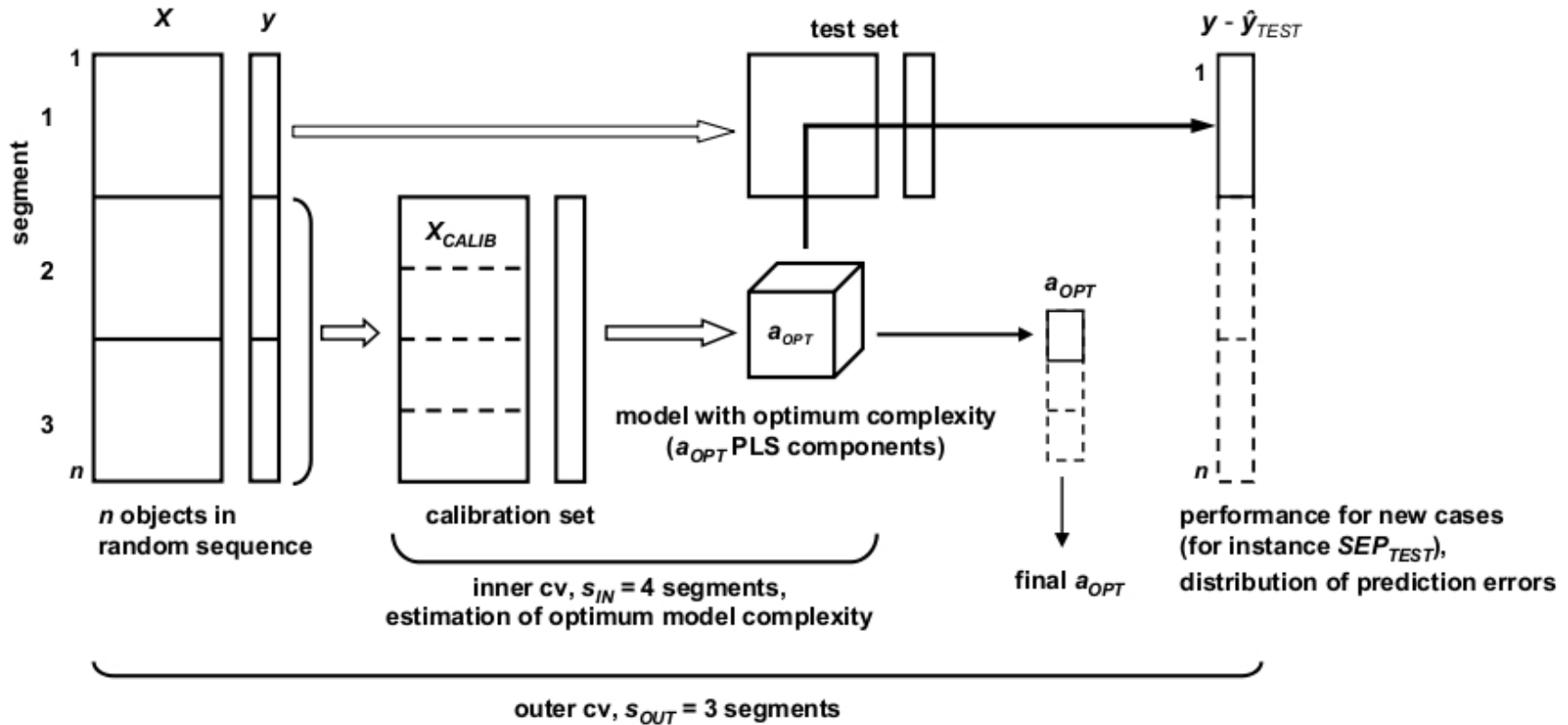
Optimal model complexity



Cross Validation



Repeated Double Cross Validation



Ridge regression (Hoerl and Kennard, 1970) is a solution to the problem of multicollinearity. The Ridge estimator is

$$\hat{\beta}_{\text{Ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

and it solves the problem

$$\hat{\beta}_{\text{Ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \right\} \quad \text{with} \quad \sum_{j=1}^p \beta_j^2 \leq \text{const.}$$

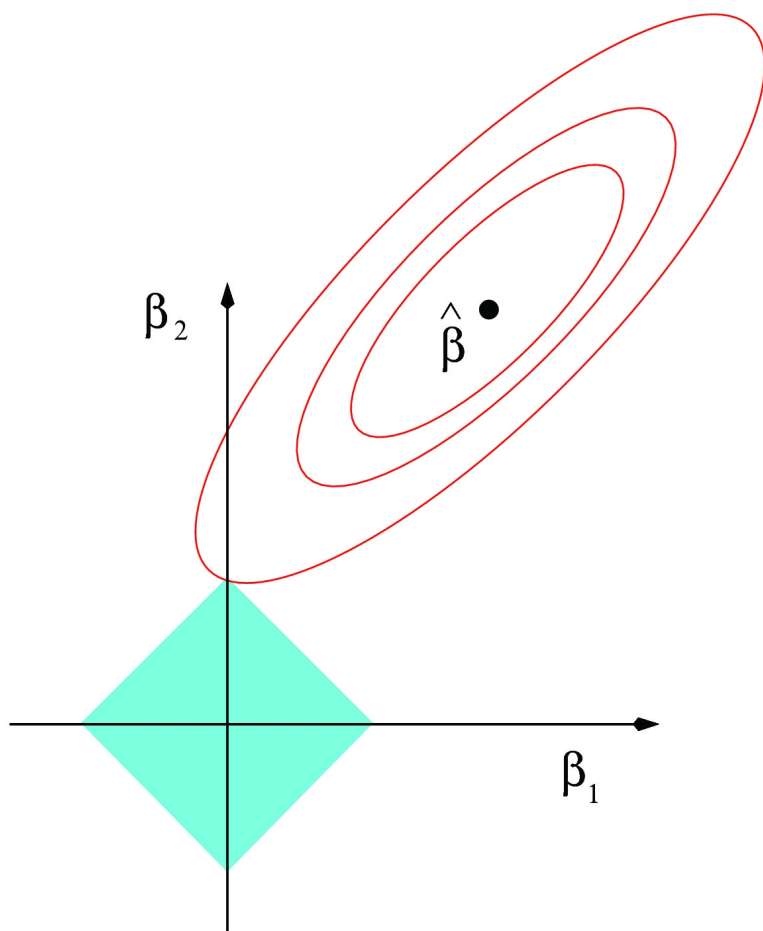
The Ridge parameter λ **shrinks the regression coefficients** (regularization).

Lasso regression (Tibshirani, 1996) solves the problem

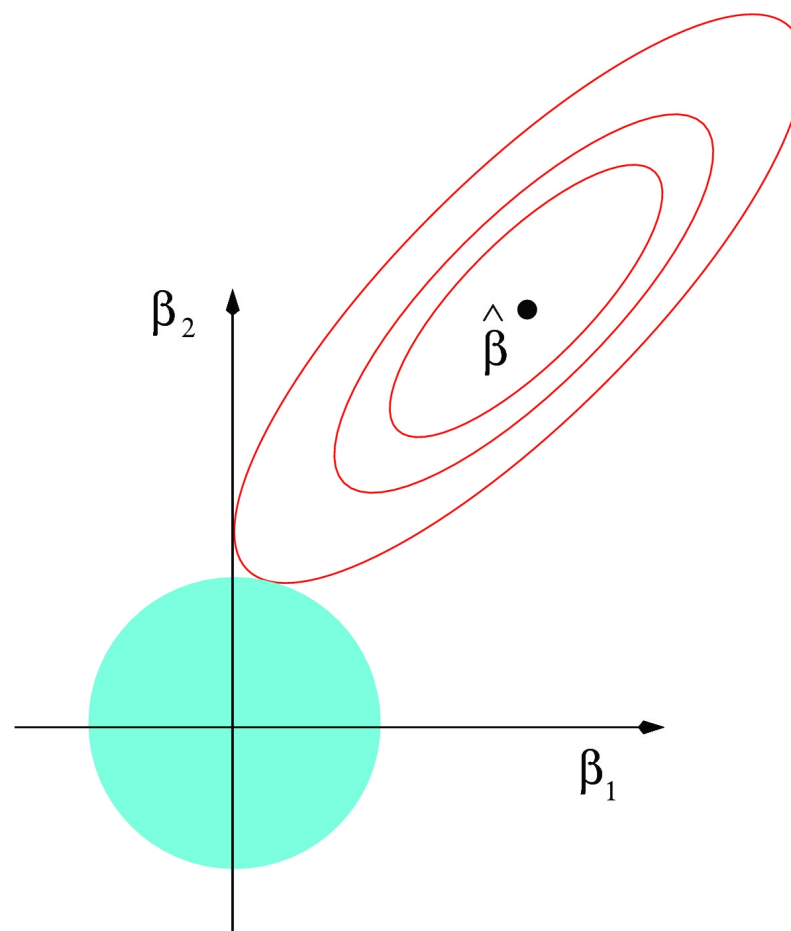
$$\hat{\beta}_{\text{Lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \right\} \quad \text{with} \quad \sum_{j=1}^p |\beta_j| \leq \text{const.}$$

No explicit formula for the solution, only numerical optimization.

Lasso: $\sum_{j=1}^p |\beta_j| \leq \text{const}$



Ridge: $\sum_{j=1}^p \beta_j^2 \leq \text{const}$



Freely available in R: <http://www.r-project.org>

Main contributions of the “chemometrics” package:

- Robustness (e.g. robust PLS)
- Model evaluation (e.g. repeated double cross-validation)
- Unified evaluation tools for parameter selection
- Diagnostic tools (e.g. for choice of number of components, visualizing effect of outliers)
- Example data sets

Example data:

Gas chromatographic retention indices of polycyclic aromatic compounds:

We consider $n = 209$ polycyclic aromatic compounds (PAC):

y -vector: GC retention index;

X -matrix: $p = 467$ descriptors of the molecular structure (Corina, Dragon).



```
> library(chemometrics)
> data(PAC)
> str(PAC)
```

List of 2

```
$ y: num [1:209] 197 197 197 200 201 ...
$ X: num [1:209, 1:467] 6.51 6.51 6.01 7.12 8.95 6.62 7.32 7.6 7.6 6.77 ...
..- attr(*, "dimnames")=List of 2
.. ..$ : chr [1:209] "1" "2" "3" "4" ...
.. ..$ : chr [1:467] "AMW" "Me" "Mp" "Ms" ...
```

Idea: reduce the number of regressor variables to a few components (PCR: using only the X -data; PLS: using both X and y data).



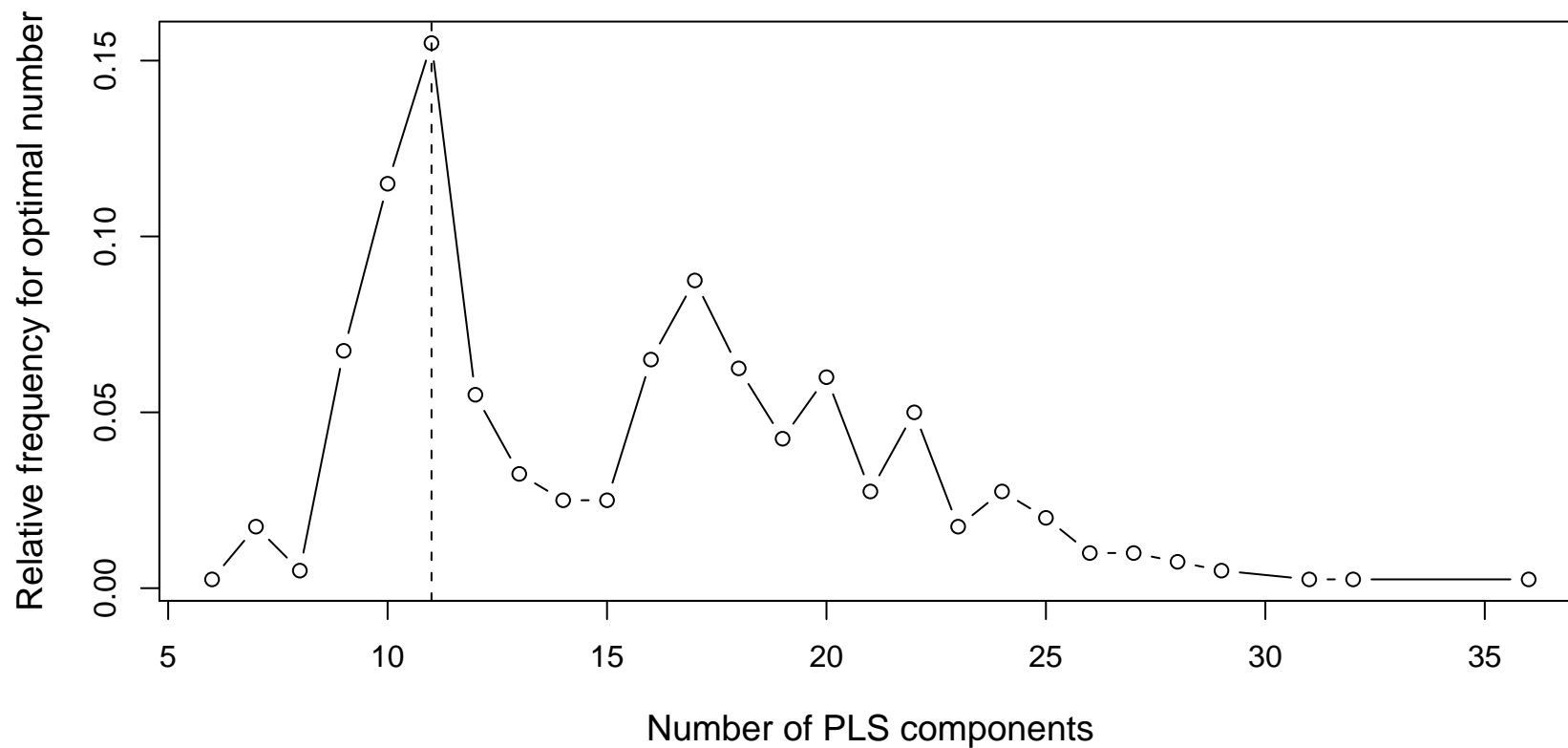
```
> pls_dcv <- mvr_dcv(y~X,ncomp=50,data=PAC,method="simpls")
# PLS with repeated double cross validation
# Default: 100 repetitions, 4 outer and 10 inner segments
# for PCR use: method="svdpc"
```

Output:

```
$ resopt : num [1:209, 1, 1:100]
$ predopt : num [1:209, 1, 1:100]
$ optcomp : int [1:4, 1:100]
$ pred : num [1:209, 1, 1:50, 1:100]
$ SEPopt : num 12
$ sIQROpt : num 8.4
$ sMADOpt : num 8.39
$ MSEPOpt : num 144
$ afinal : num 11
$ SEPfinal: Named num [1:50]
```

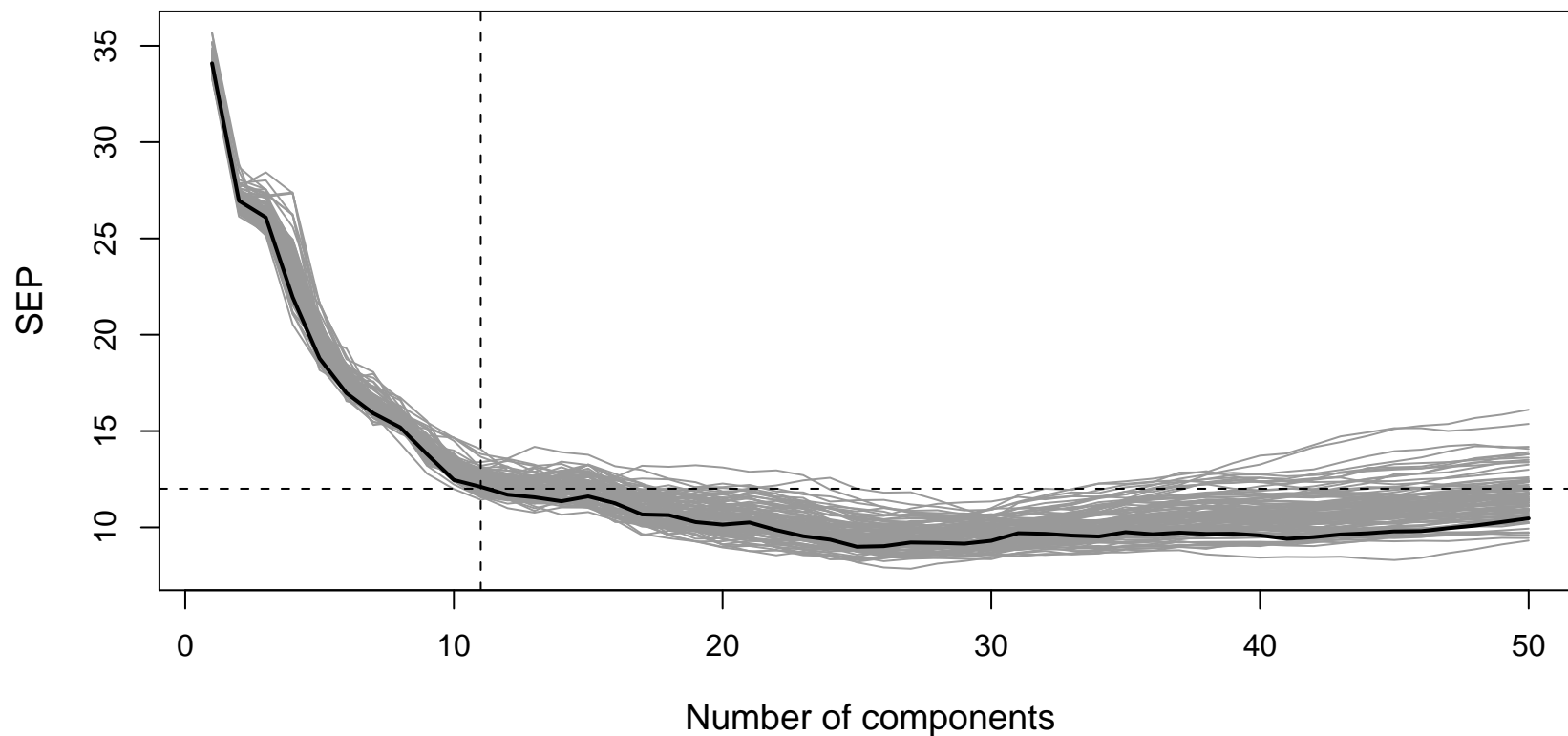
Diagnostic plots: optimal number of components (based on 4×100 values)

```
> plotcompmvr(pls_dcv)
```



Diagnostic plots: SEP for 1, ..., 50 components

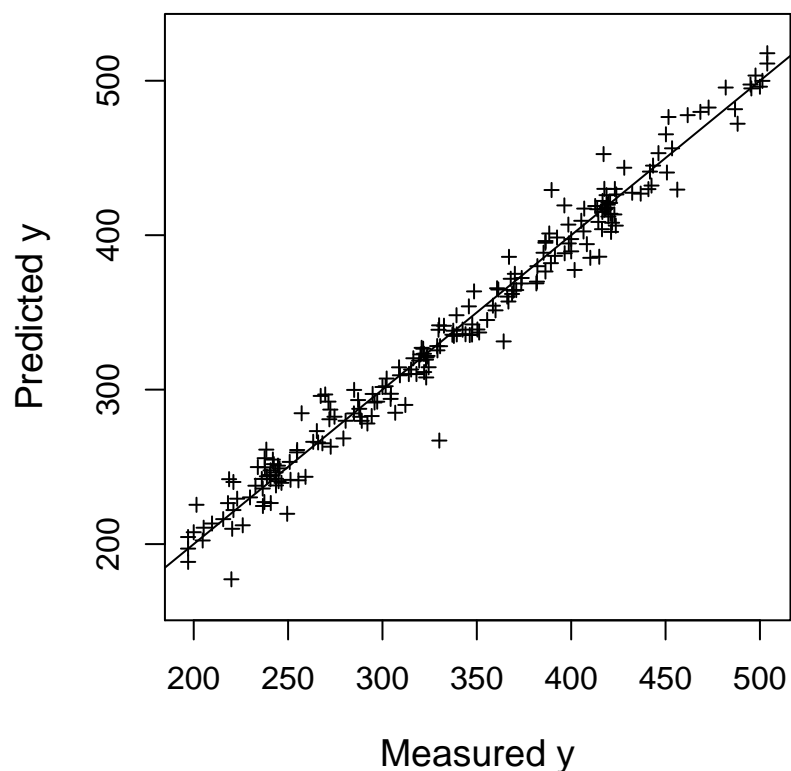
```
> plotSEPMvr(pls_dcv, optcomp=11, y=PAC$y, X=PAC$X, method="simpls")
```



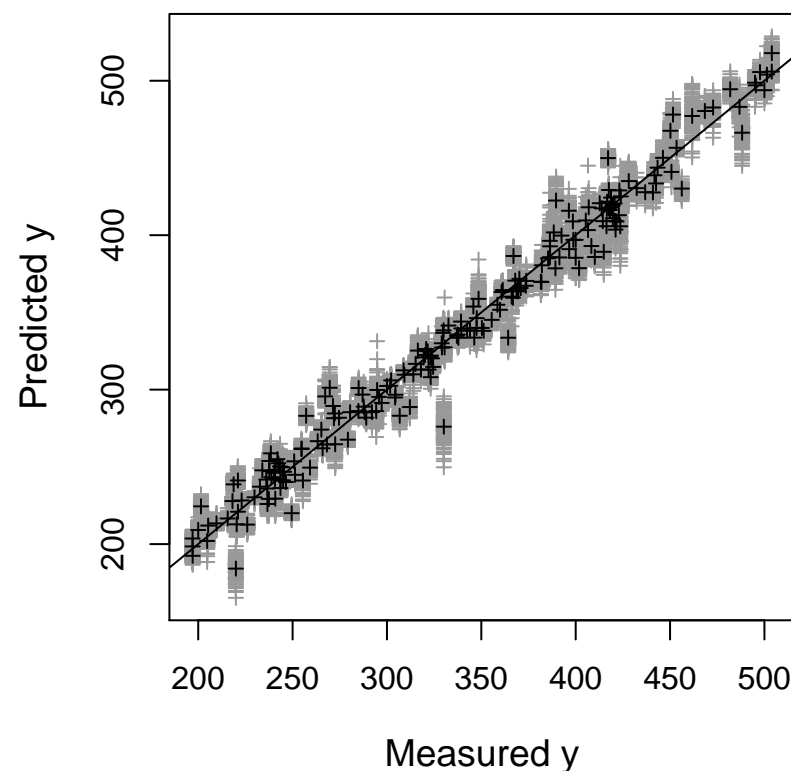
Diagnostic plots: predicted values from 100 models with 11 components

```
> plotpredmvr(pls_dcv, optcomp=11, y=PAC$y, X=PAC$X, method="simpls")
```

Prediction from CV



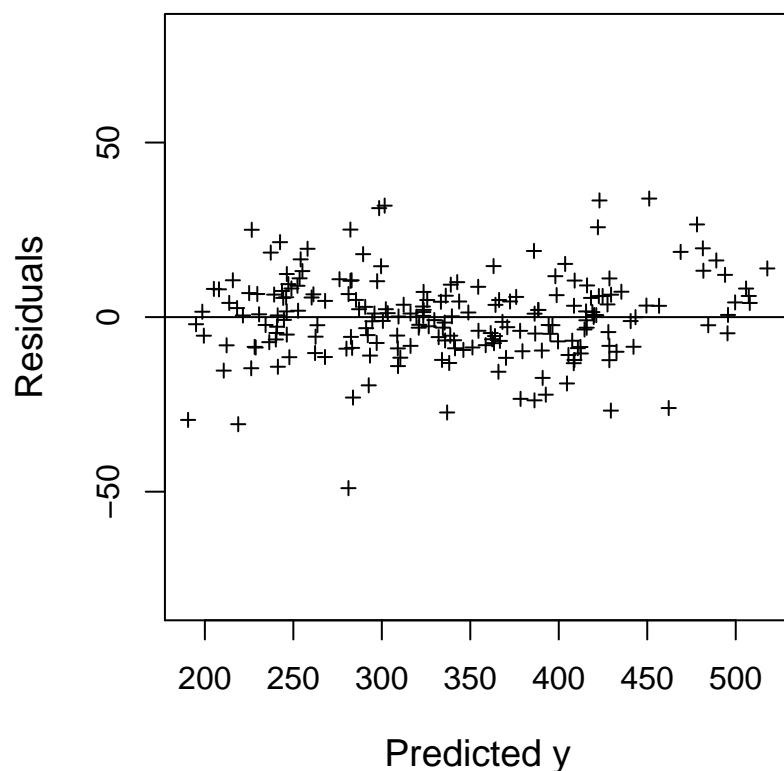
Prediction from Repeated Double-CV



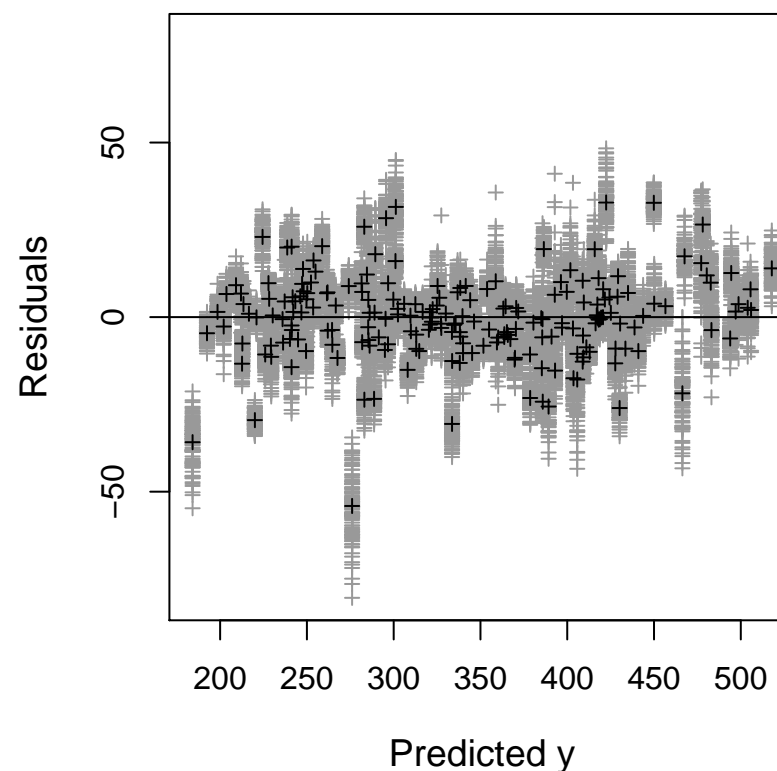
Diagnostic plots: residuals from 100 models with 11 components

```
> plotresmvr(pls_dcv, optcomp=11, y=PAC$y, X=PAC$X, method="simpls")
```

Results from CV

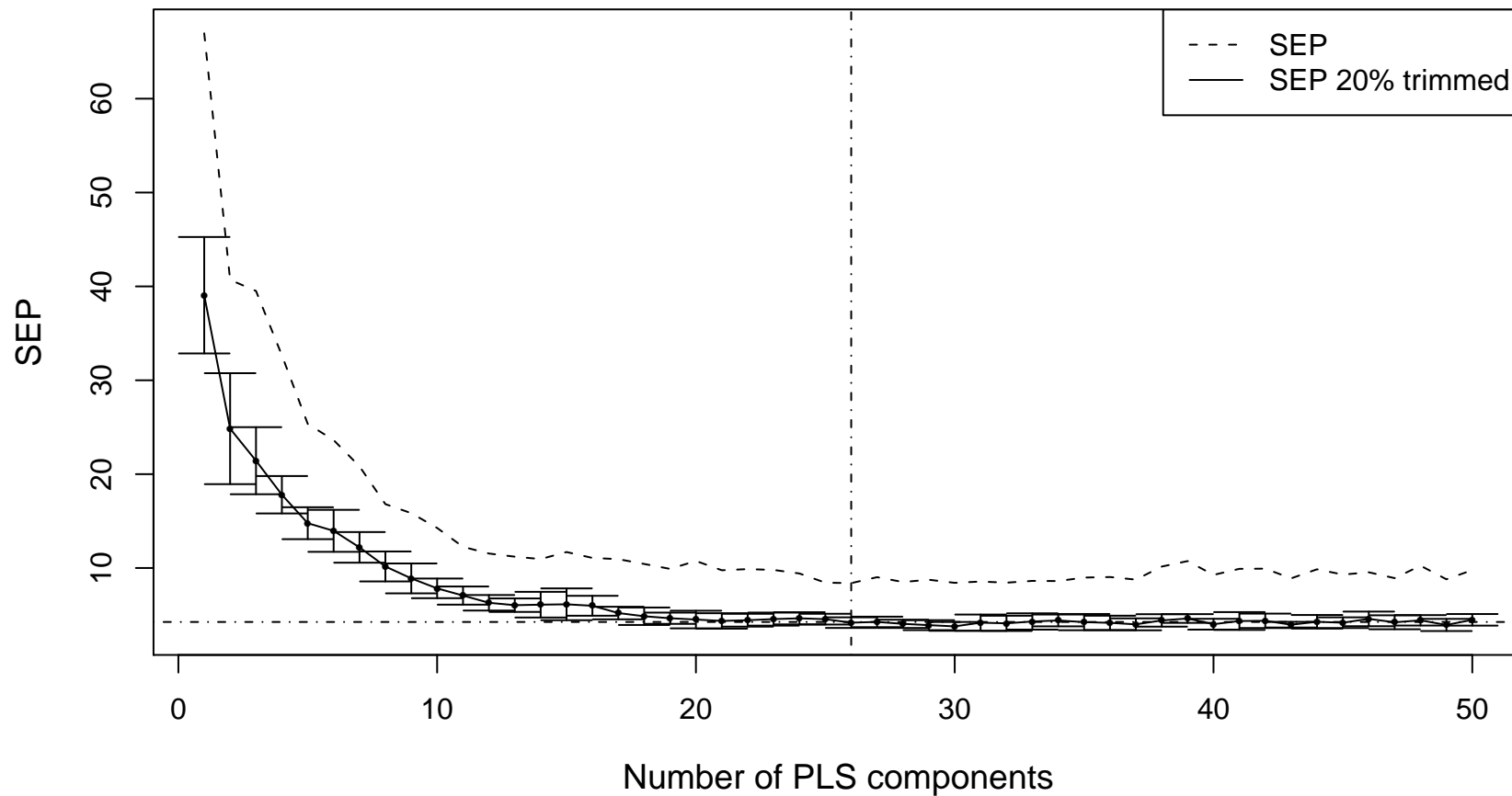


Results from Repeated Double-CV



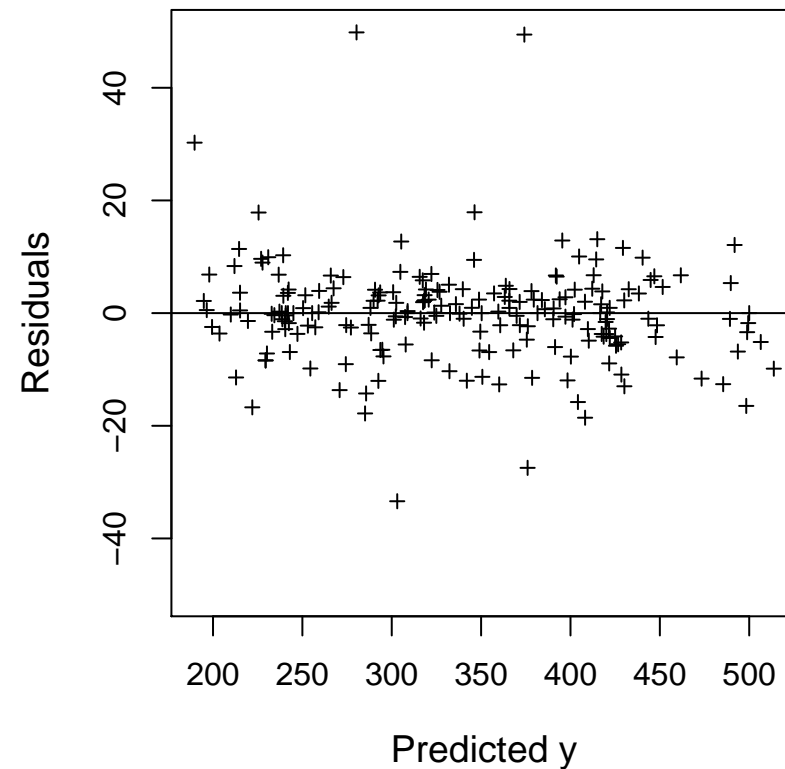
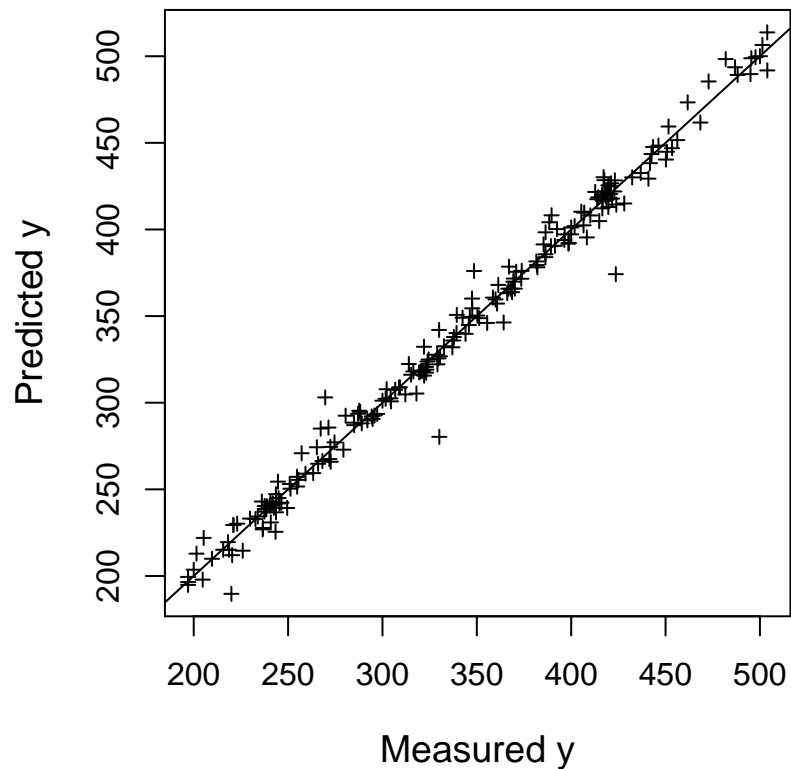
PRM: Serneels, Croux, Filzmoser, Van Espen (ChemoLab, 2005)

```
> prm_cv(PAC$X,PAC$y,a=50,trim=0.2,plot.opt=TRUE)
```



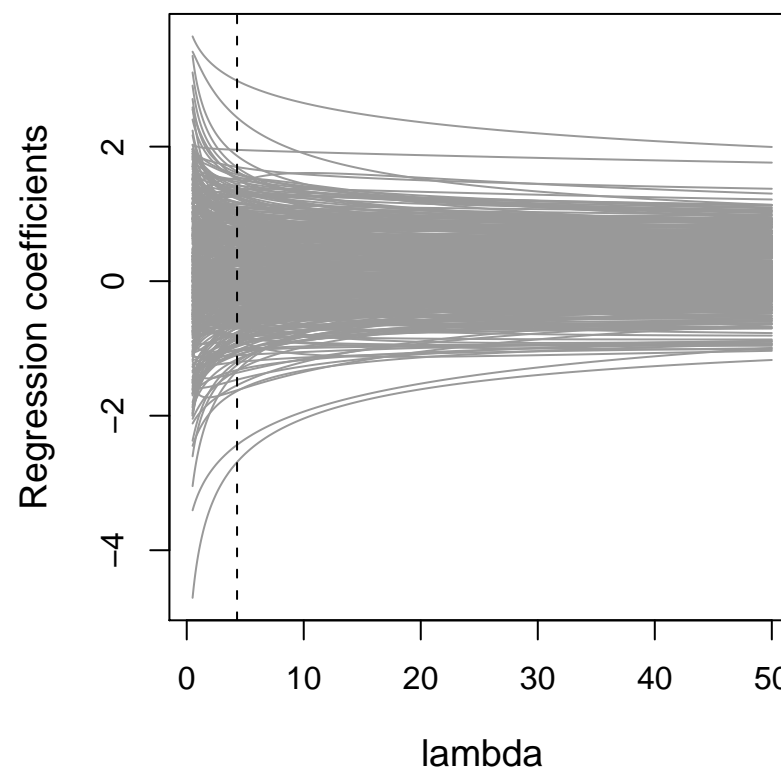
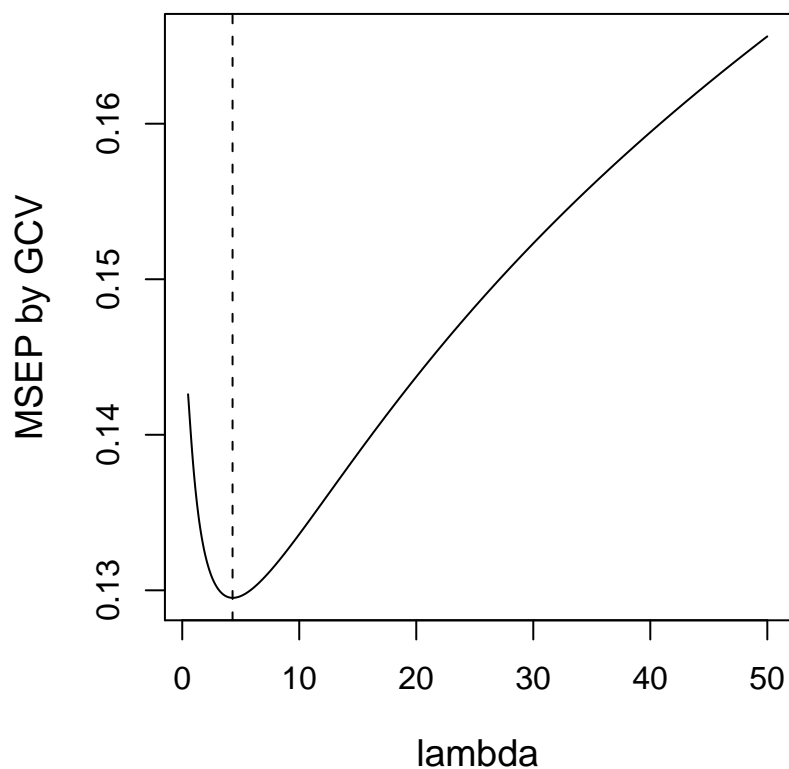
Diagnostic plots: predicted values and residuals using 26 components

```
> plotprm(resprmcv,PAC$y)
```



Diagnostic plot: choice of ridge parameter

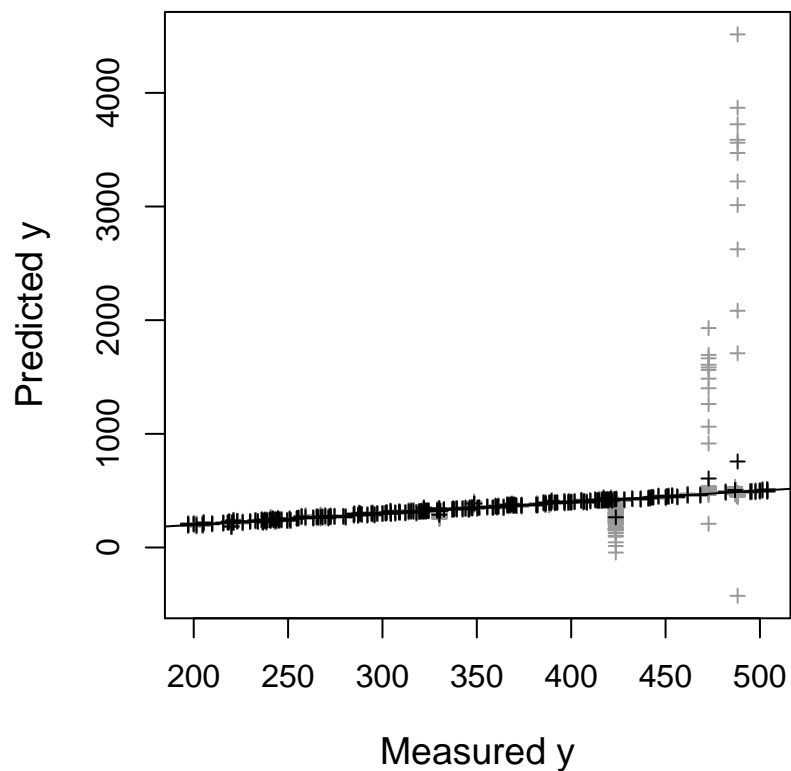
```
> resR <- plotRidge(y~X,data=PAC,lambda=seq(0.5,50,by=0.05))
```



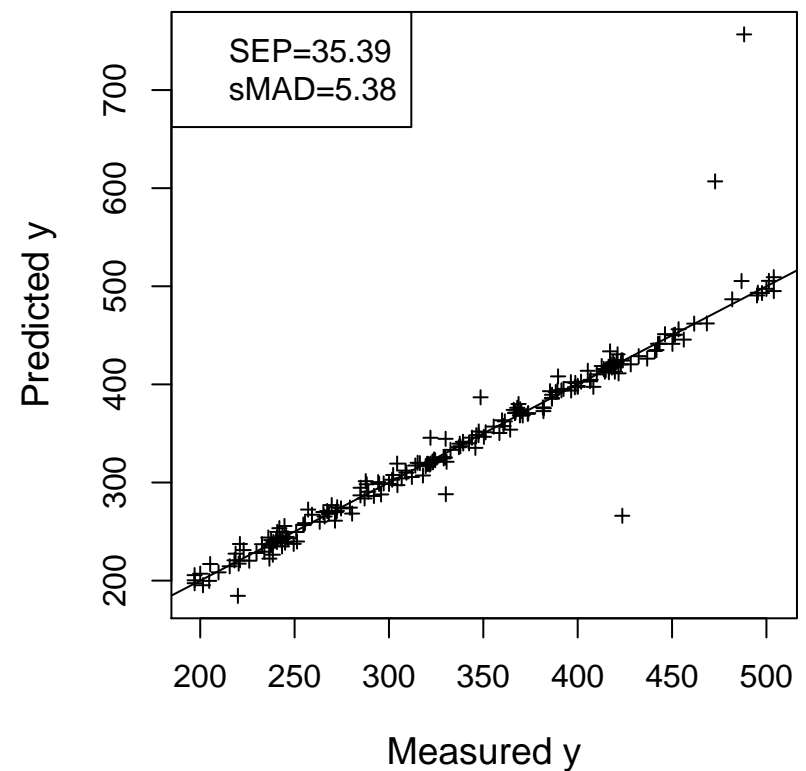
Diagnostic plots: from repeated cross validation

```
> resRcv <- ridgeCV(y~X,data=PAC,repl=100,lambda=resR$lambdaopt)
```

Predictions from Repeated CV

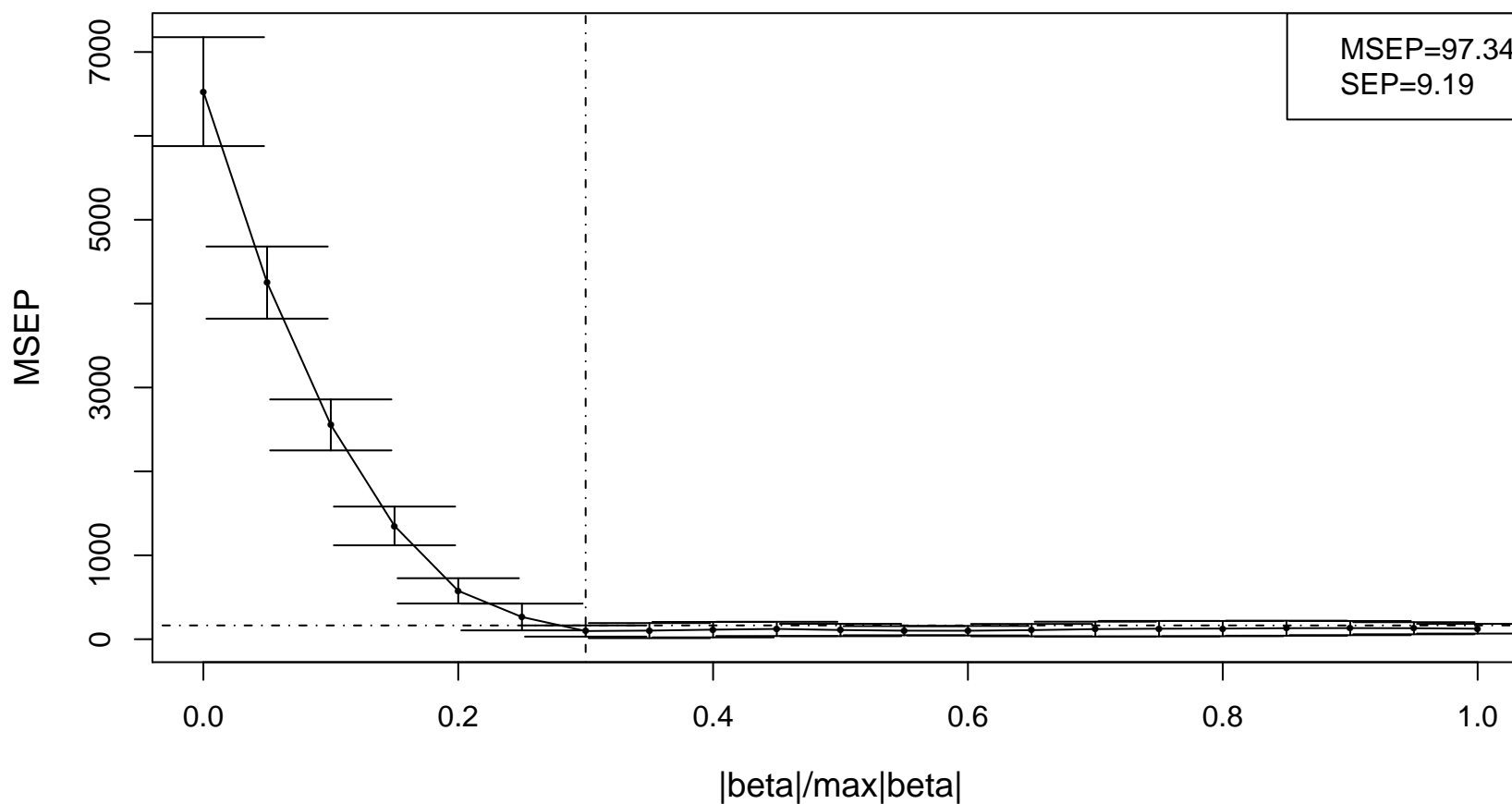


Average of predictions



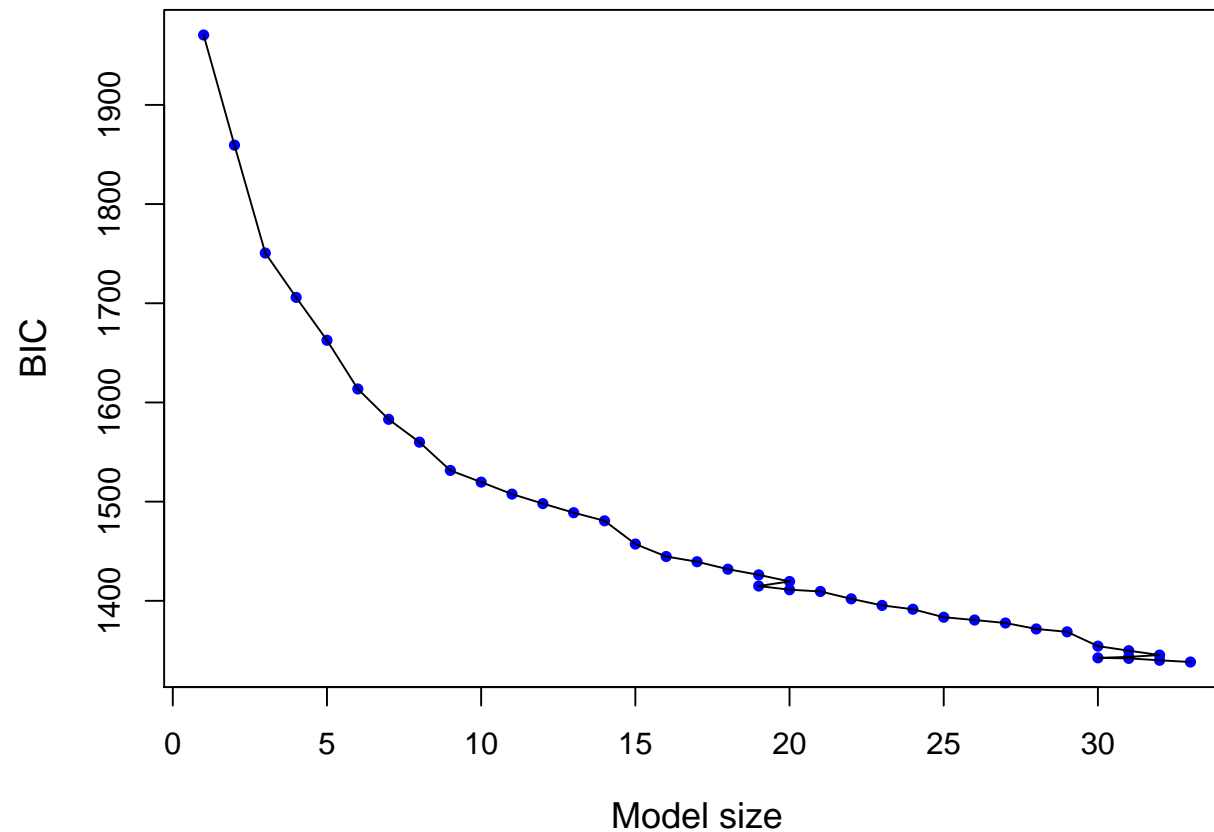
Diagnostic plot: choice of Lasso parameter

```
> resL <- lassoCV(y~X,data=PAC,K=10,fraction=seq(0,1,by=0.05))
```



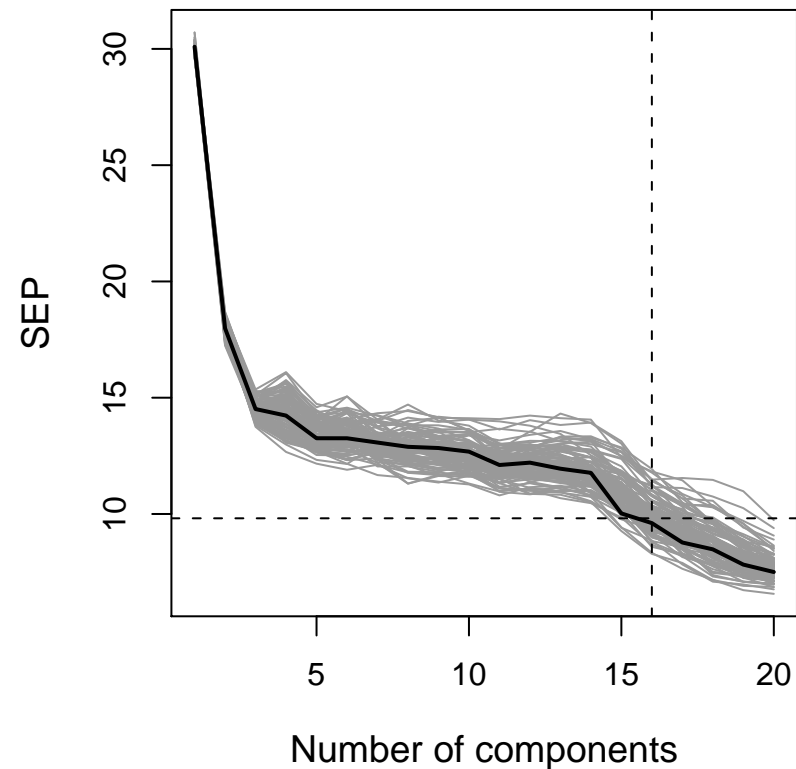
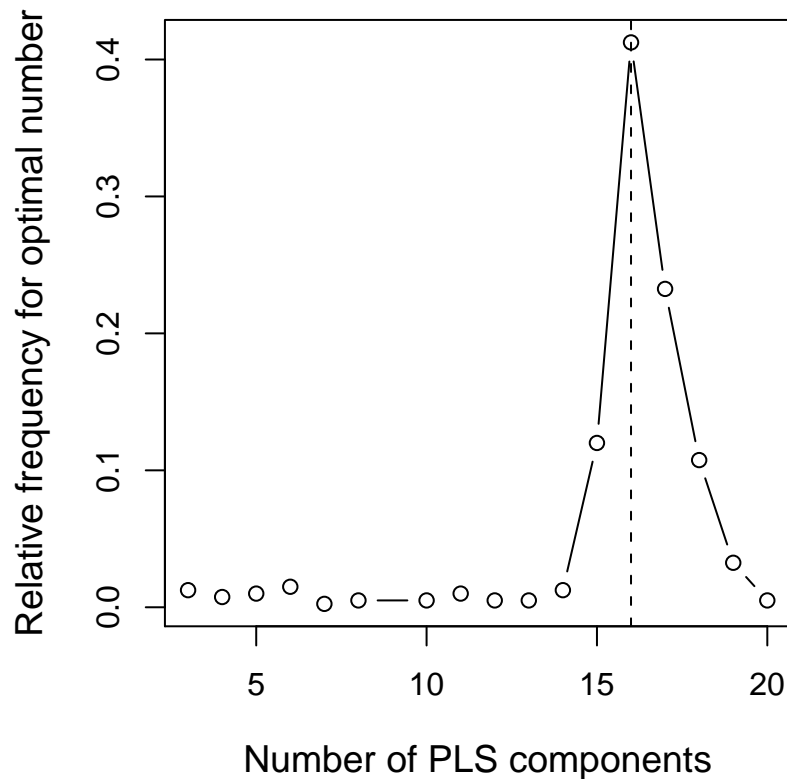
Diagnostic plot: choice of model

```
> resS <- stepwise(y~X,data=PAC)
```



Diagnostic plot: choice of number of PLS components

```
> resSdcv <- mvr_dcv(y~., ncomp=20, data=PACred, method="simpls")
```



Comparison of Results

Method	p^*	k	SEP_{Test}	SEP_{CV}	$SEP^{0.2}$
PCR	467	21	14.2	—	7.9
PLS	467	11	12.0	—	5.7
Robust PLS	467	26	—	8.9	4.0
Ridge regression	467	—	—	28.4	4.0
Lasso regression	145	—	—	7.7	5.0
Stepwise variable selection + PLS	33	16	9.6	—	4.4

Example data: Origin of glass samples:

- $n = 214$ glass samples
- 6 different glass types (e.g. windows, headlamps, tableware, containers)
- $p = 9$ variables (refractive index, mass-% of Al, Ba, Ca, Fe, K, Mg, Na, Si)



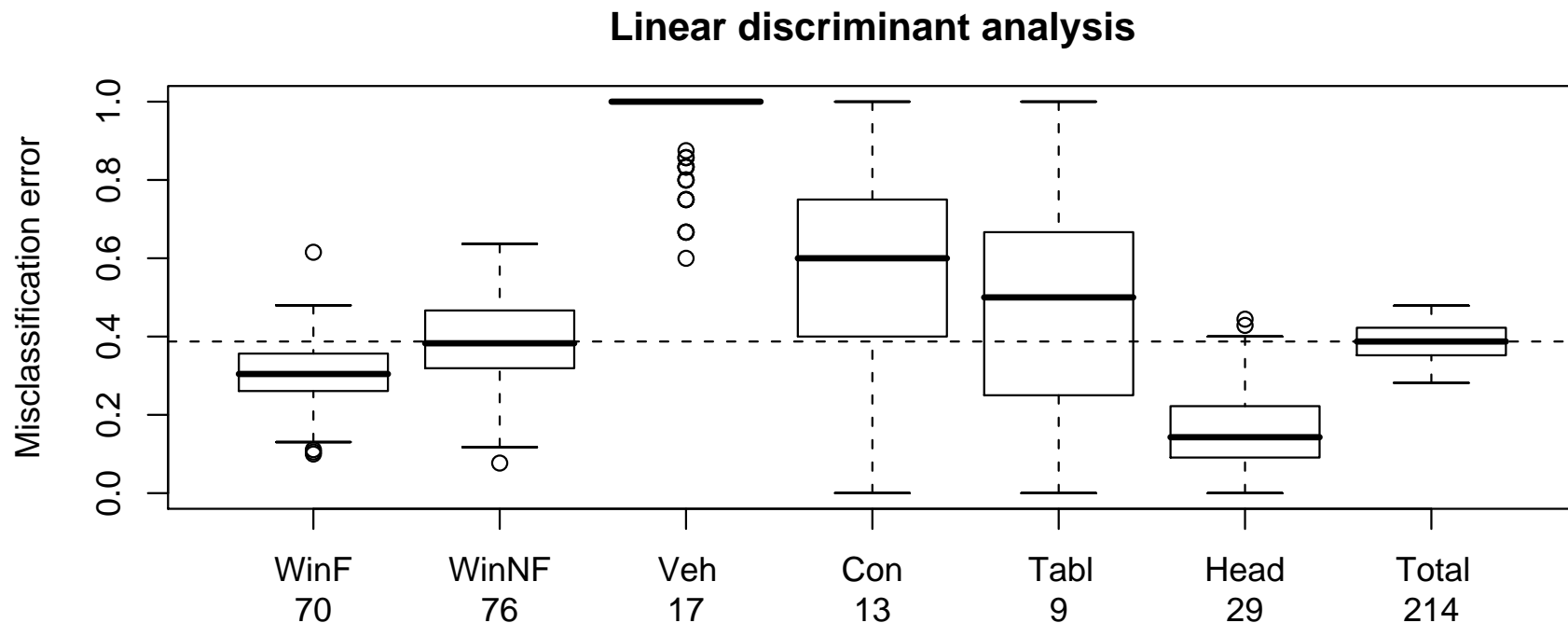
```
> library(MASS)
> data(fgl)
> grp=fgl$type
> X <- scale(fgl[,1:9])
> dim(X)
```

```
[1] 214 9
```

LDA (Linear Discriminant Analysis)

LDA: obtain LDA-rule for training data, apply to test data; repeat

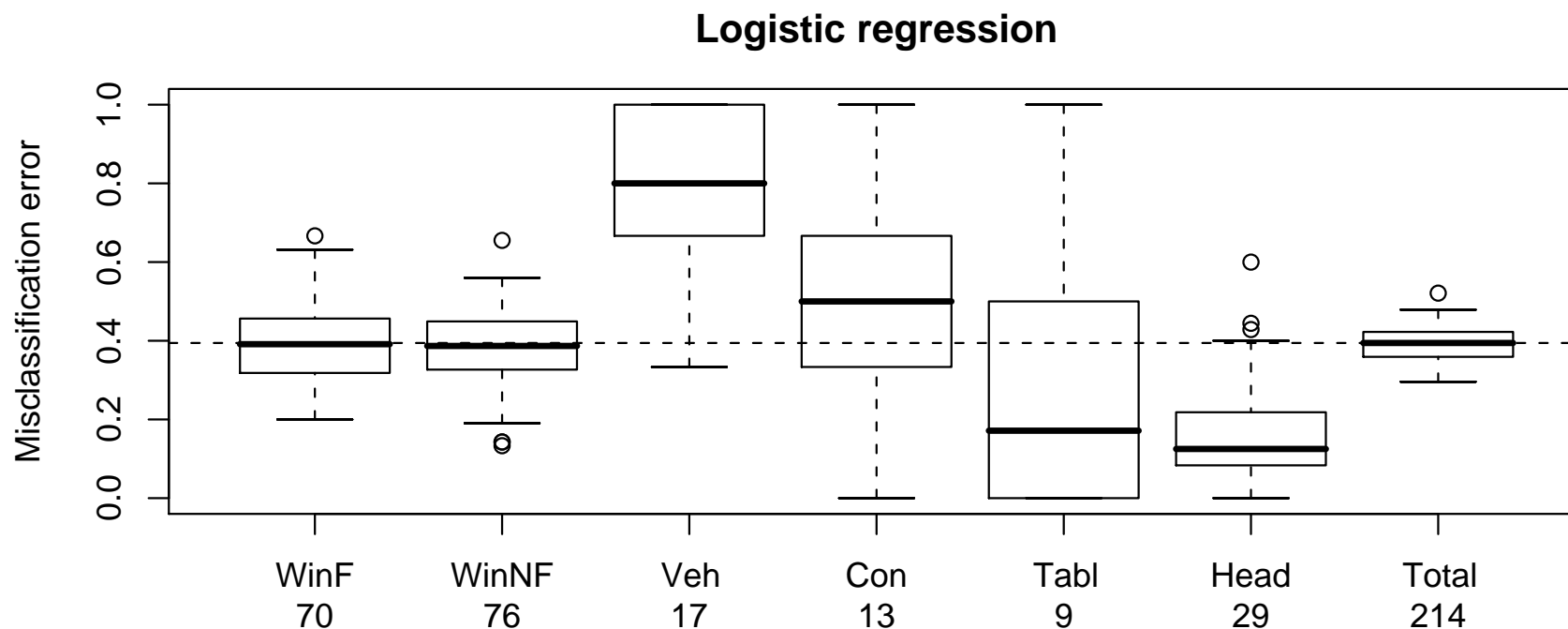
```
> train <- sample(1:n,ntrain)
> reslda <- lda(X[train,],grp[train])
> pred <- predict(reslda,newdata=X[-train,])$class
> tab <- table(grp[-train],pred)
```



LR (Logistic Regression)

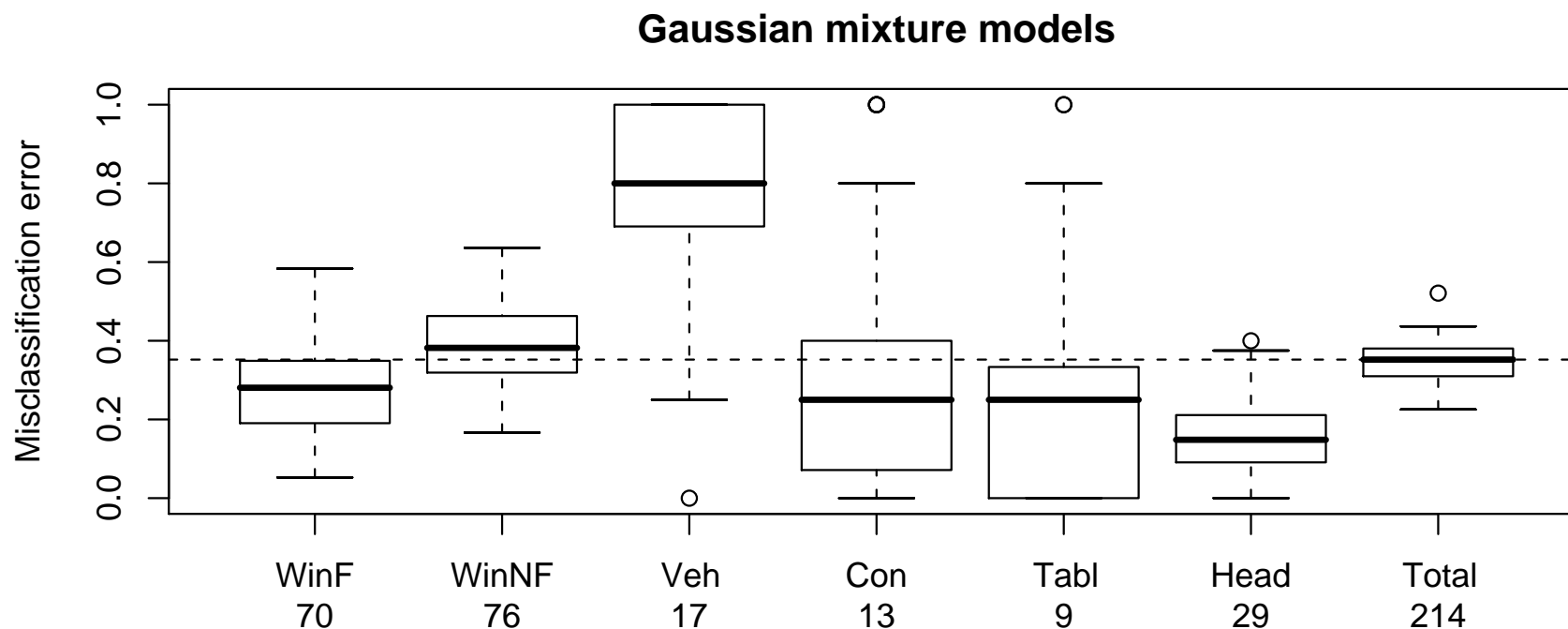
LR: obtain LR-rule for training data, apply to test data; repeat

```
> train <- sample(1:n,ntrain)
> reslr <- vglm(grp ~ .,data=dat[train,],family=multinomial)
> predmix <- predict(reslr,dat[-train,],type="response")
> predgrp <- apply(predmix,1,which.max)
> tab <- table(grp[-train],predgrp)
```



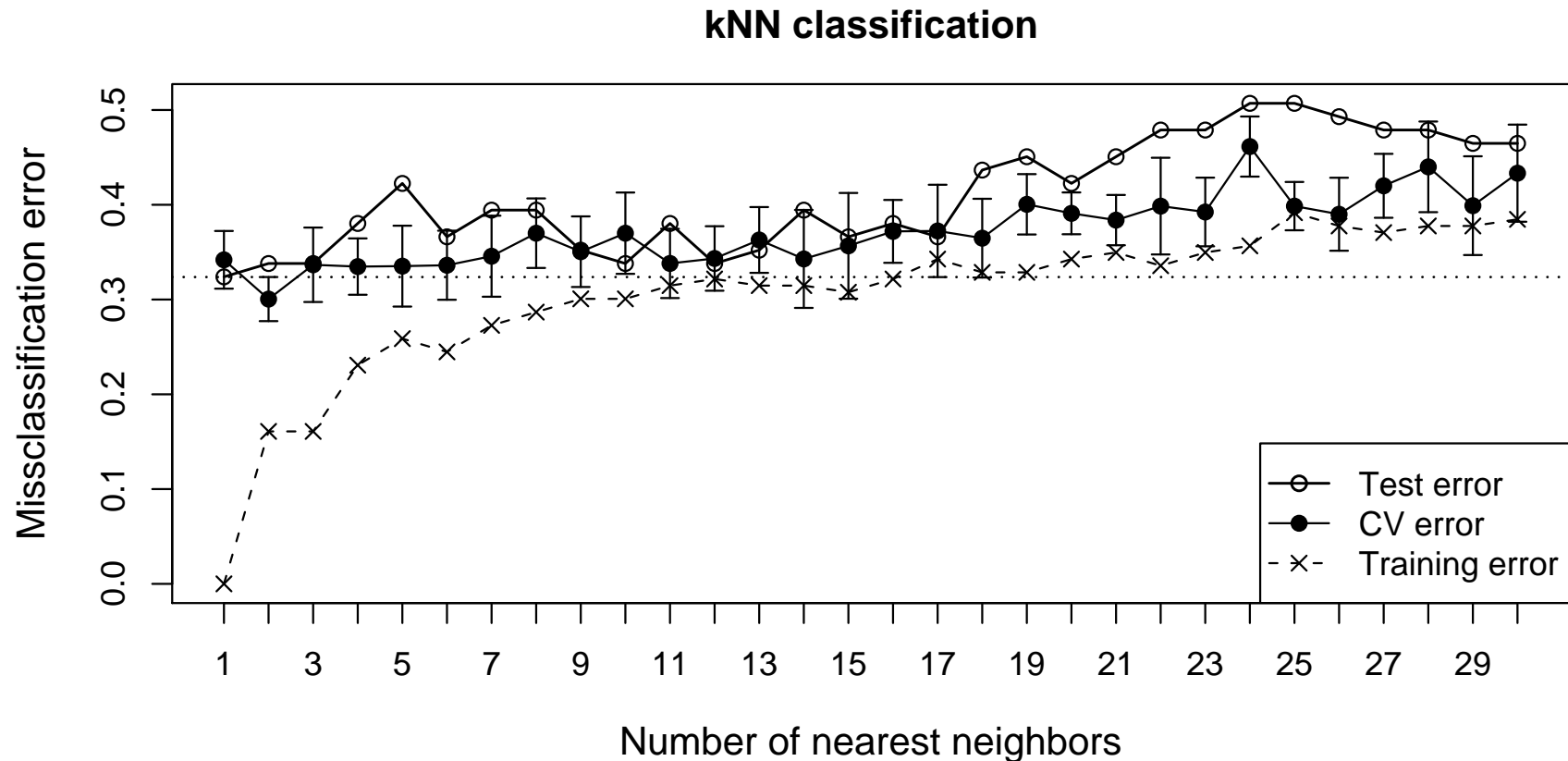
Mix: obtain models for training data, apply to test data; repeat

```
> train <- sample(1:n,ntrain)
> resgmm <- mda(grp ~ .,data=dat[train,])
> predgmm <- predict(resgmm,dat[-train,],type="post")
> predgrp <- apply(predgmm,1,which.max)
> tab <- table(grp[-train],predgrp)
```



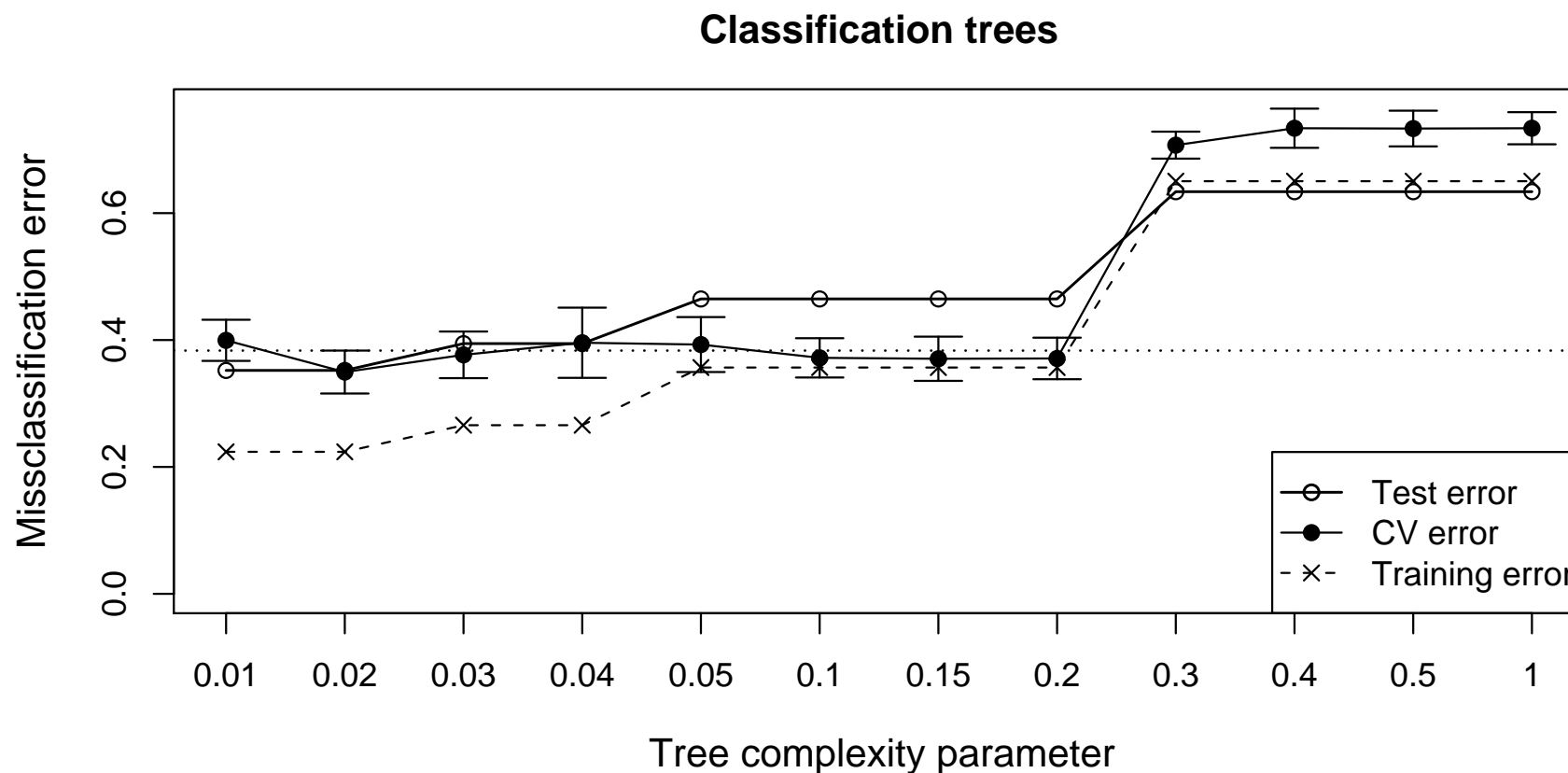
kNN: select tuning parameter “k” (number of neighbors)

```
> train <- sample(1:n,ntrain)
> resknn <- knnEval(X,grp,train,knnvec=seq(1,30,by=1))
```



Tree: select tuning parameter “cp” (tree complexity)

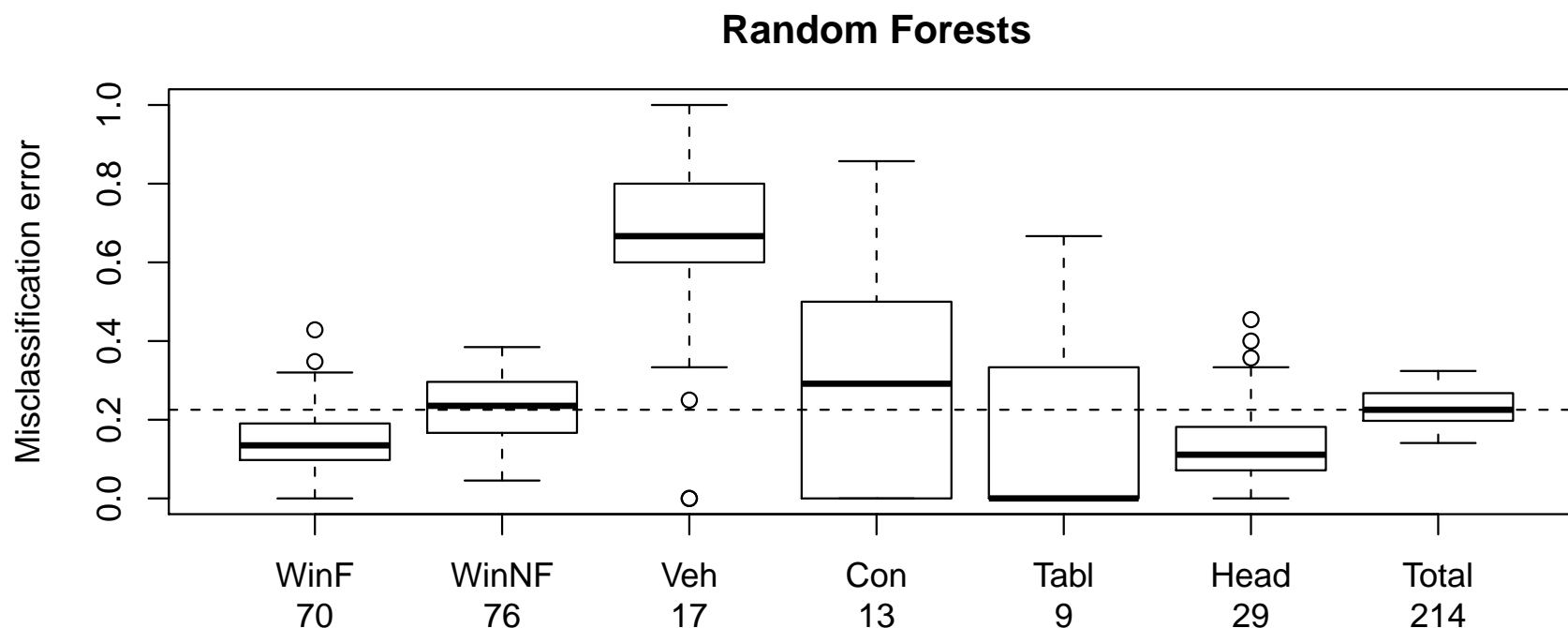
```
> train <- sample(1:n,ntrain)
> cptry <- c(0.01,0.02,0.03,0.04,0.05,0.1,0.15,0.2,0.3,0.4,0.5,1)
> restree <- treeEval(X,grp,train,cp=cptry)
```



RF: Random Forests

RF: obtain rule for training data, apply to test data; repeat

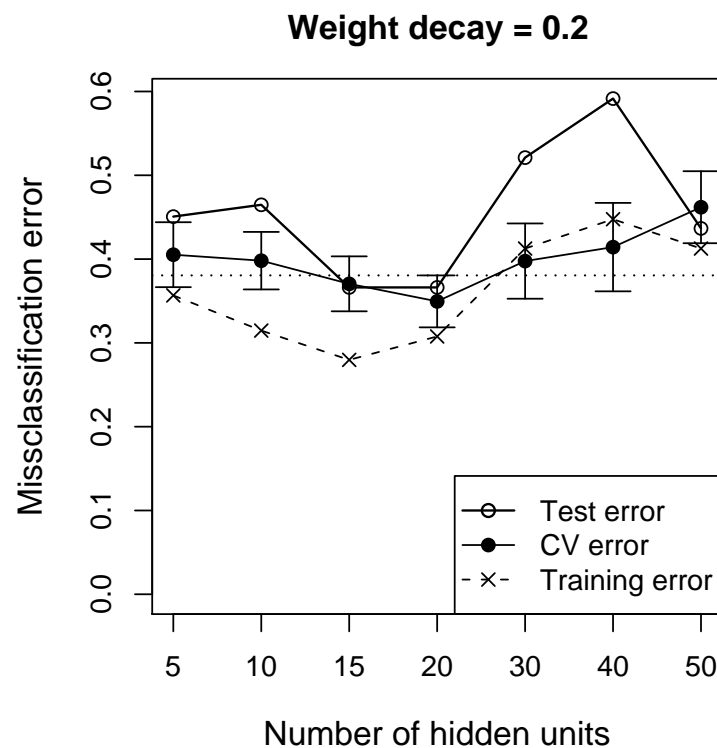
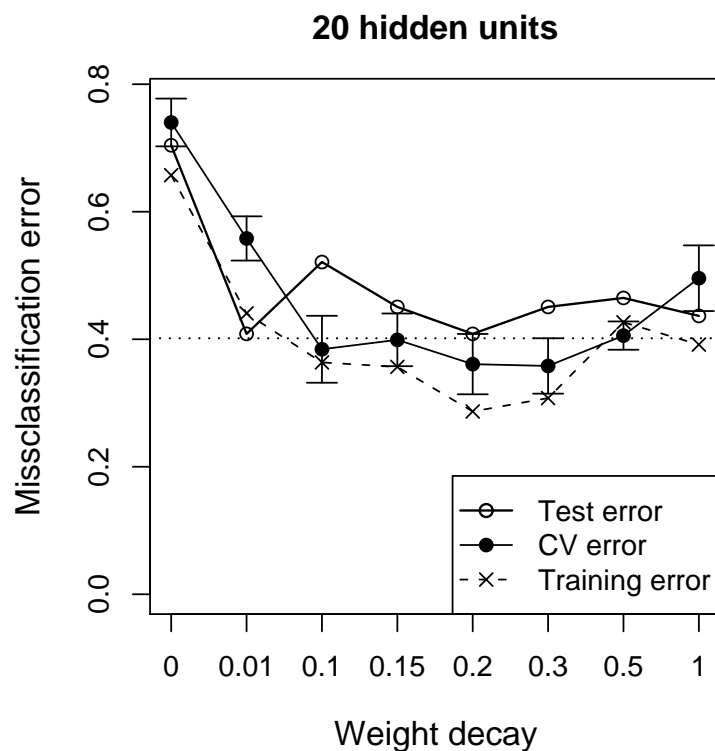
```
> train <- sample(1:n,ntrain)
> resRF <- randomForest(grp~.,data=dat,subset=train)
> predRF <- predict(resRF, dat[-train,])
> table(grp[-train],predRF)
```



ANN: Artificial Neural Networks

ANN: select tuning parameters “weight decay” and “number of hidden units”

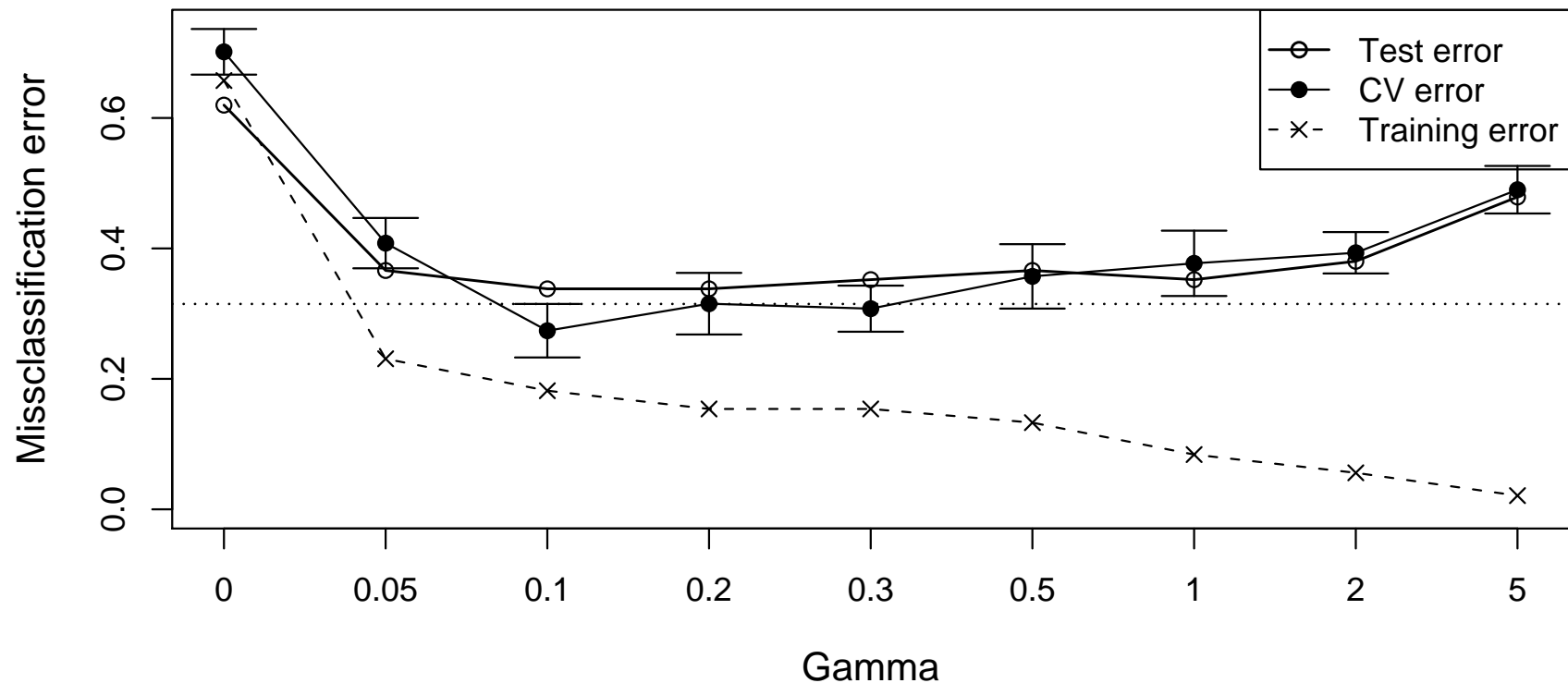
```
> train <- sample(1:n,ntrain)
> wd <- c(0,0.01,0.1,0.15,0.2,0.3,0.5,1)
> sz <- c(5,10,15,20,30,40,50)
> resnnet=nnetEval(X,grp,train,decay=wd,size=20)
> resnnet=nnetEval(X,grp,train,decay=0.2,size=sz)
```



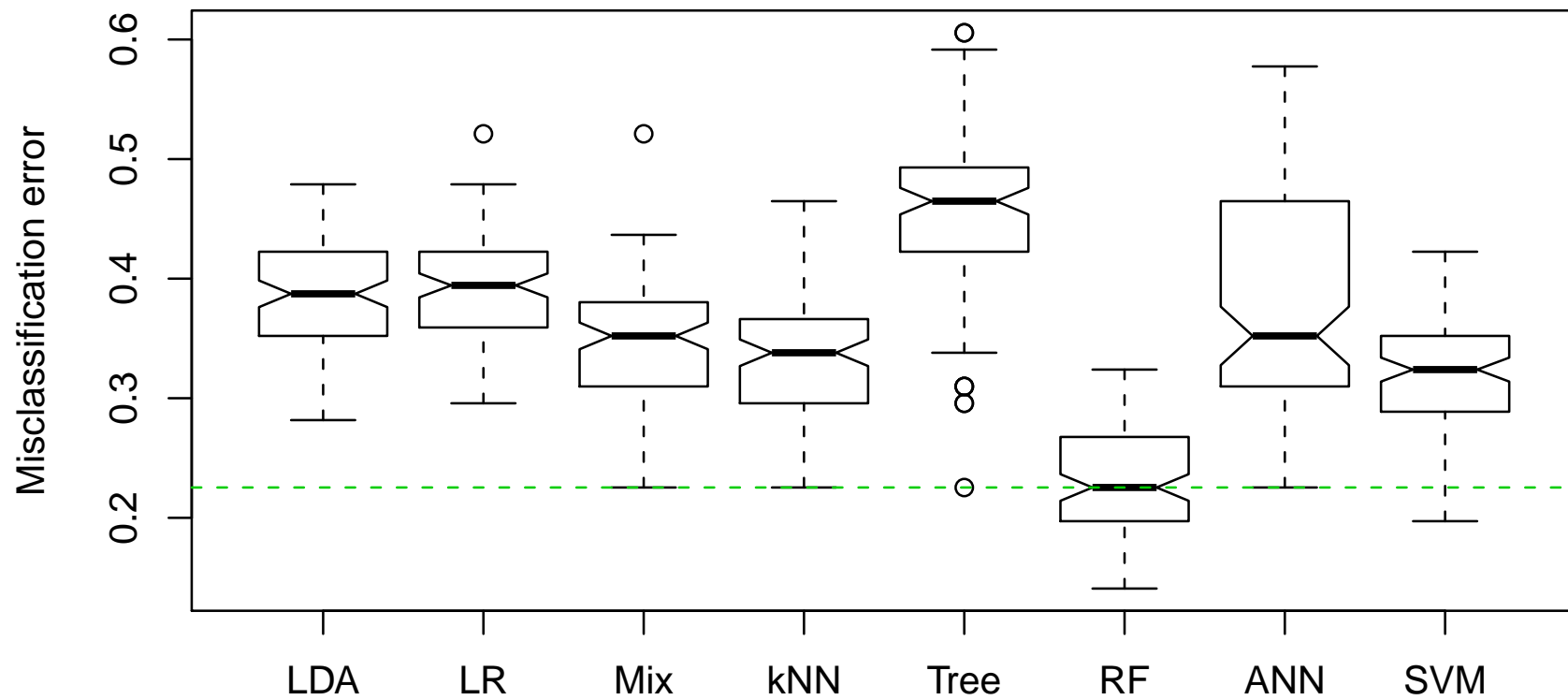
SVM: select tuning parameter “ γ ” (size constraint of slack variables)

```
> train <- sample(1:n,ntrain)
> gv <- c(0,0.05,0.1,0.2,0.3,0.5,1,2,5)
> ressvm <- svmEval(X,grp,train,gamvec=gv)
```

Support vector machines



Overall Comparison



- The book is not only suitable for chemometricians, but for all who want to learn about **multivariate statistical methods** from a theoretical and practical perspective.
- Data sets and methods treated in the book are included in the R package **chemometrics**. The different evaluation procedures (*repeated cross validation*) are implemented in a unified manner.
- We already received a lot of positive feedback to both, the book and the R package.