

project-5

Chitra Karki

11/9/2021

Reading data

```
df.crime <- read.csv("crime.csv")
```

1. Data Preparation: Start with the data set crime.csv from the class website. We will first do some minor data preparation and exploration.

a. Remove the first five columns from the data since these are not predictive.

```
df.crime <- df.crime[, -c(1:5)]
```

b. Take a look at the missing percentage of each remaining variable. Remove those heavy missing, say, over 60%.

```
nrows <- dim(df.crime)[1]
col2drop <- which(apply(df.crime, 2, function(x) sum(is.na(x))/nrows * 100) > 60)
# columns containing more than 60 % NA values
col2drop
```

```
##      LenasSwfPerFt      LenasSwfPerPop      LenasSwfFie1dOps
##      LenasSwfFie1dPerPop      LenasTotalReq      LenasTotReqPerPop
##      PolicePerPopOffic      PolicePerPop      RacialMatchComPol
##      PctPoliceWhite      PctPoliceBlack      PctPoliceHisp
##      PctPoliceAsian      PctPoliceNor      OfficAssignPrugUnits
##      NumIndsDrugsSeiz      PoliceAveDWorked      PoliceCars
##      PoliceCoperBurg      PolicePctPolicePatr      LenasGangUnitDeploy
##      PoliceBudgePerPop
##      122
```

```
df.crime <- df.crime[, -col2drop]
# checking dimension after dropping cols with heavy NAs
dim(df.crime)
```

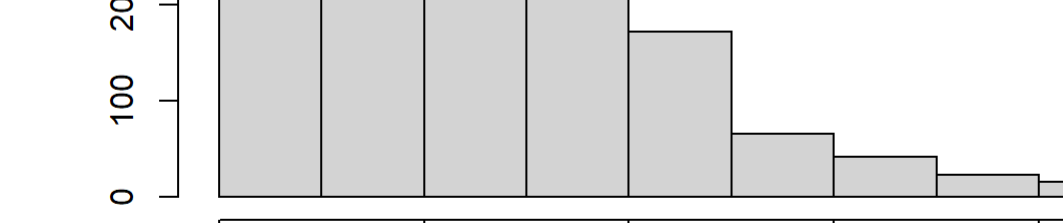
```
## [1] 1994 181
```

c. Impute or replace the remaining missing values appropriately.

```
which(apply(df.crime, 2, function(x) sum(is.na(x))/nrows * 100) > 0)
```

```
## OtherPerCap
##      26
```

```
# this shows now we are left with single NA value at the 26 i.e "otherpercap" column.
# lets replace the missing with the median value.
hist(df.crime$OtherPerCap) # looking at the distribution of data points.
```

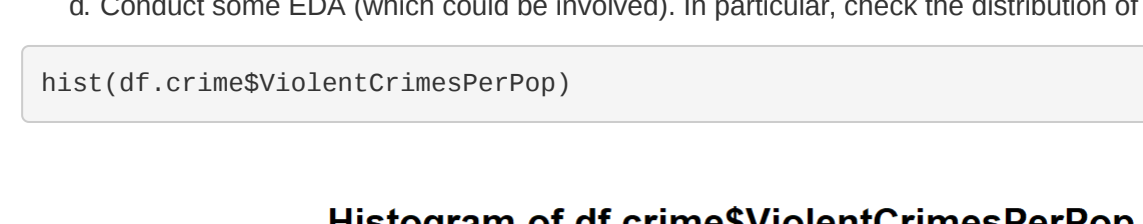


```
df.crime[which(is.na(df.crime$OtherPerCap)), "OtherPerCap"] <- median(df.crime$OtherPerCap, na.rm = T)
# checking for na
sum(is.na(df.crime)) # now the data is clean.
```

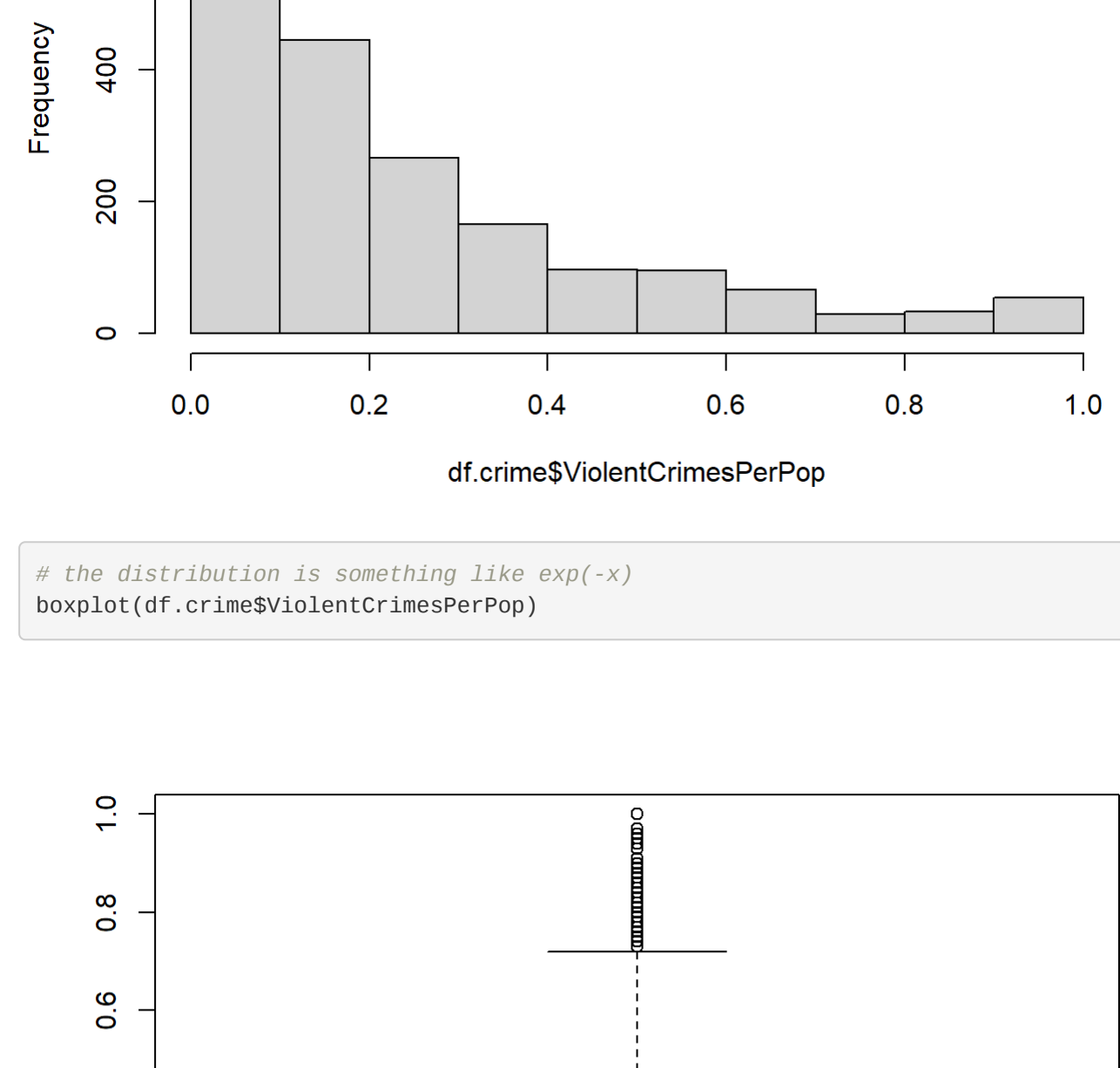
```
## [1] 0
```

d. Conduct some EDA (which could be involved). In particular, check the distribution of the target variable ViolentCrimesPerPop.

```
hist(df.crime$ViolentCrimesPerPop)
```



```
# the distribution is something like exp(-x)
boxplot(df.crime$ViolentCrimesPerPop)
```



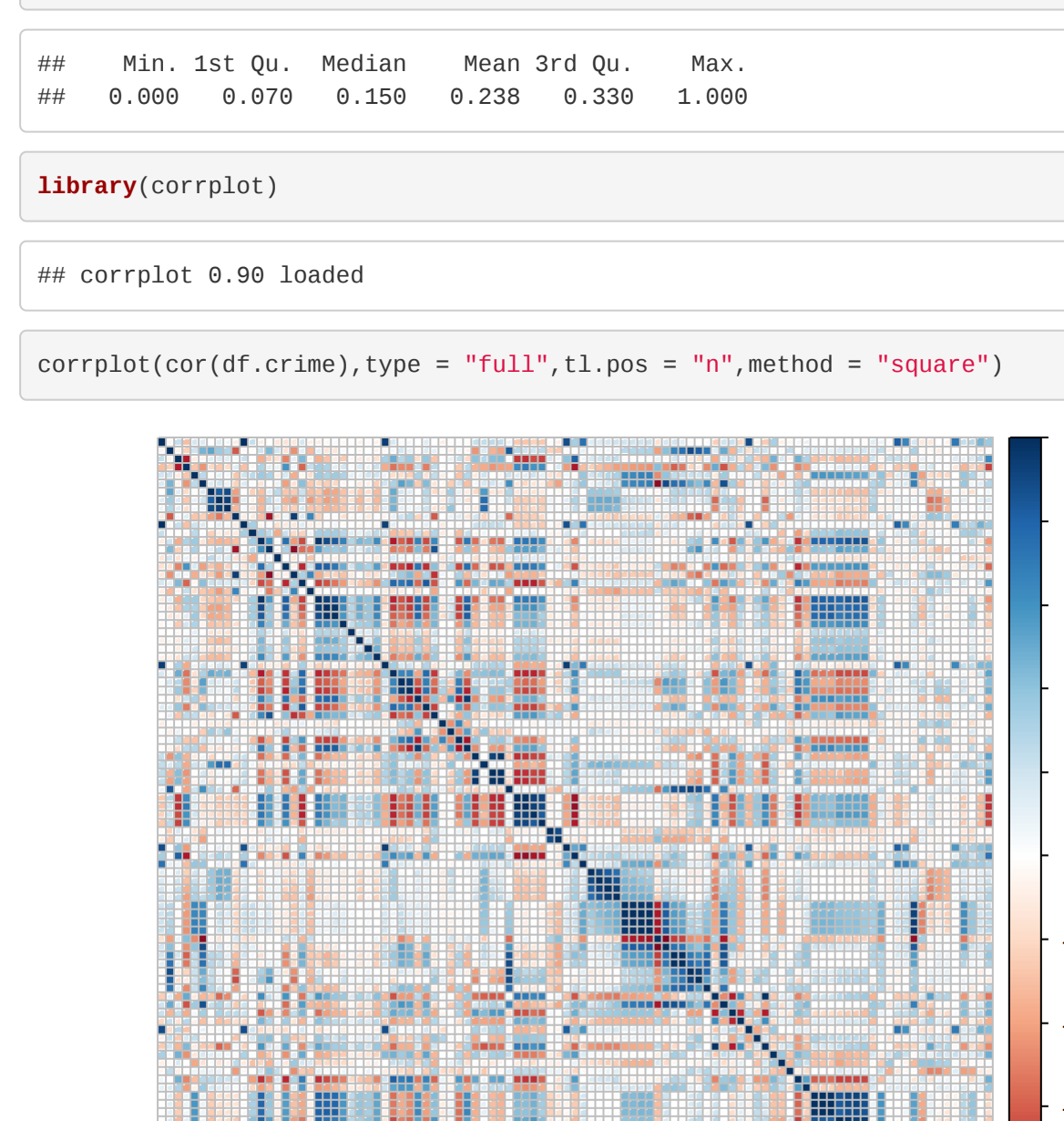
```
# there are certain outliers. More amount of data are greater than median
summary(df.crime$ViolentCrimesPerPop)
```

```
##      Min.      1st Qu.      Median      Mean      3rd Qu.      Max.
##      0.000      0.070      0.150      0.238      0.338      1.000
```

```
library(corrplot)
```

```
# corrplot 0.90 loaded
```

```
corrplot(cor(df.crime), type = "full", tl.pos = "n", method = "square")
```



there are several variables correlated with each other. Also some other variables are positively and negatively correlated with the response variable.

2. Partitioning Data: Randomly partition your data into two sets: the training set D1 and the test set D2 with a ratio of 2:1. In order for your results to be reproducible, report the random seed that you use in the partitioning.

```
set.seed(1) # seed 123 for reproducibility
train.index <- sample(1:nrows, 2/3*nrows, replace = F)
test.index <- train.index
d1 <- df.crime[train.index,]
d2 <- df.crime[test.index,]
```

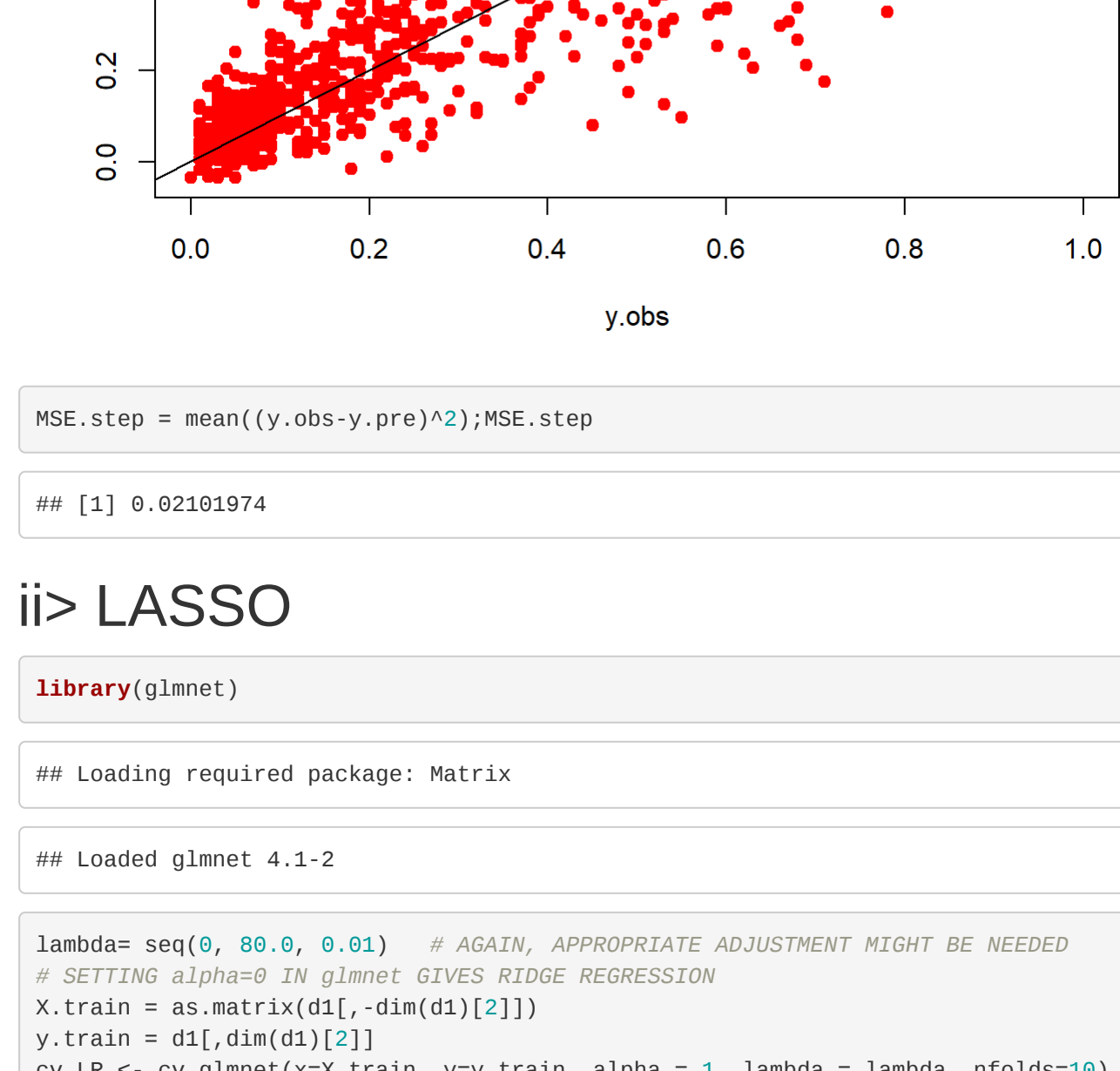
3. Predictive Modeling: Referring to the sample R code R09.R from the class website, fit at least five models of your own choice from the following models using the training set D1:

- Linear regression with stepwise selection;
- LASSO
- Ridge Regression (RR)
- Principal Components Regression (PCR)
- Partial least squares regression (PLSR)
- Weighted orthogonal components regression (WOCR)
- Total least squares regression (TLSR)
- Stagewise regression
- Least angle regression (LAR) You might want to describe how you select the tuning parameter if such decision needs to be made in each approach. Then apply the fitted model to the test set D2. Report the mean square error of prediction (MSEP) and compare.

i> linear regression with stepwise selections

```
library(MASS)
fit.full <- lm(d1$ViolentCrimesPerPop ~., data = d1)
fit.step <- stepAIC(fit.full, direction="both", k=log(nrow(d1)))
```

```
# fit: stepAICova
best.model <- summary(fit.step)
best.formula <- as.formula(best.model$call)
best.fit <- lm(best.formula, data = d1)
old.dat <- as.data.frame(best.matrix(best.fit))
new.dat <- d2[, names(old.dat)[-1]]
y.pre <- predict(best.fit, new.dat)
y.obs <- d2[, dim(d2)[2]]
plot(y.obs, y.pre, pch=19, col="red", main = "stepwise Regression")
abline(a=0, b=1, col="black")
```



```
MSE.step <- mean((y.obs - y.pre)^2); MSE.step
```

```
## [1] 0.02181974
```

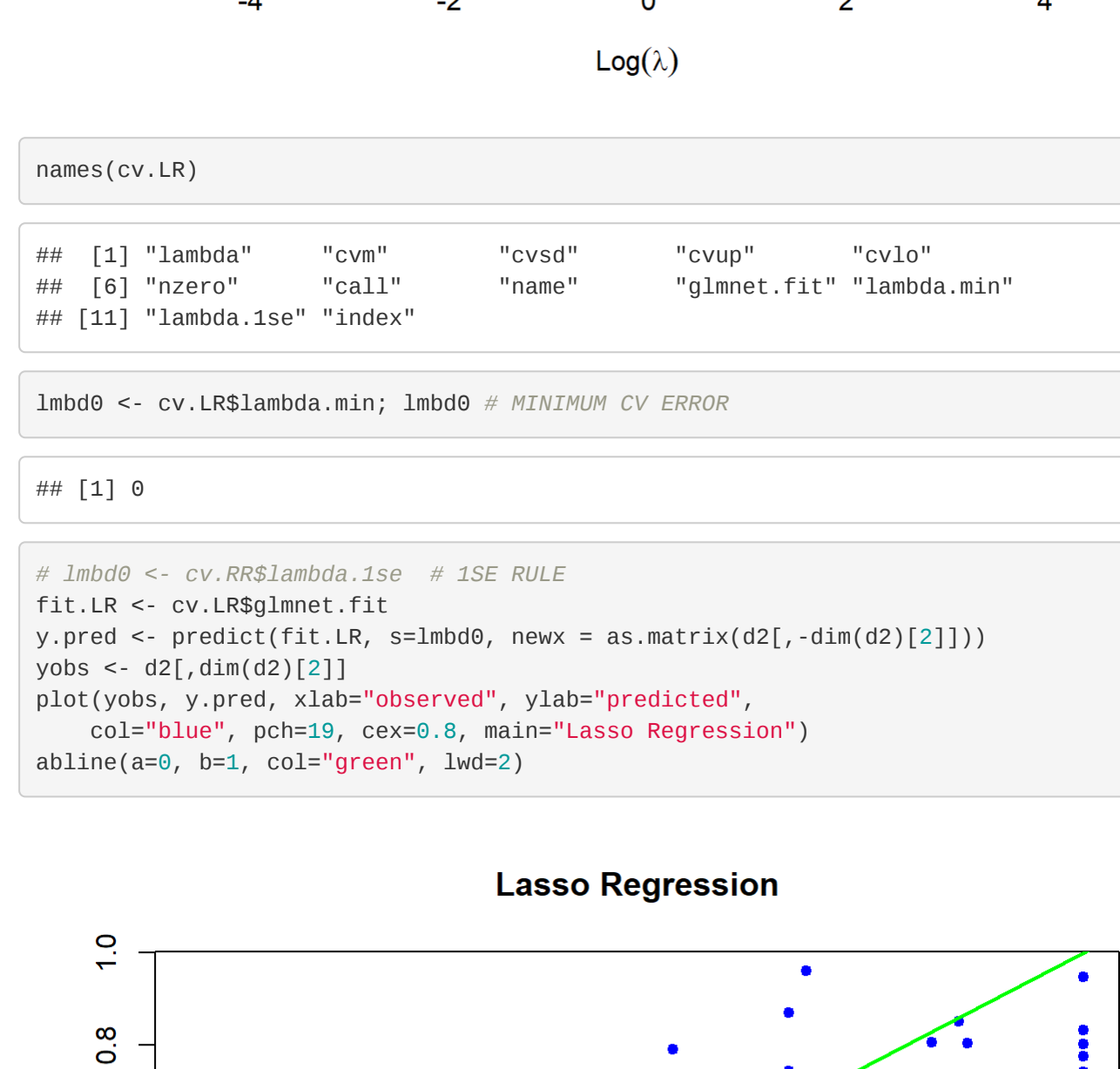
ii> LASSO

```
library(glmnet)
```

```
## Loading required package: Matrix
```

```
## Loaded glmnet 4.1-2
```

```
lambda <- seq(0, 80, 0.01) # AGAIN, APPROPRIATE ADJUSTMENT MIGHT BE NEEDED
# SETTING lambda IN glmnet GIVES RIDGE REGRESSION
X.train <- as.matrix(d1[, -dim(d1)[2]])
y.train <- d1[, dim(d1)[2]]
cv.lr <- cv.glmnet(X=train, y=y.train, alpha = 1, lambda = lambda, nfolds=10) # 10-FOLD CV BY DEFAULT
plot(cv.lr)
```



```
names(cv.lr)
## [1] "lambda"      "cvm"          "cvstd"        "cvup"         "cvlo"
## [6] "nzero"        "call"        "name"        "glmnet.fit"   "lambda.min"
## [11] "lambda.1se"   "index"
```

```
lmbd0 <- cv.lr$lambda.min; lmbd0 # MINIMUM CV ERROR
```

```
## [1] 0
```

```
# lmbd0 <- cv.lr$lambda.1se # 1SE RULE
fit.lr <- cv.glmnet.fit
y.pred <- predict(fit.lr, s=lmbd0, newx = as.matrix(d2[, -dim(d2)[2]]))
yobs <- d2[, dim(d2)[2]]
plot(yobs, y.pred, xlab="observed", ylab="predicted",
     col="blue", pch=19, cex=0.8, main="Lasso Regression")
abline(a=0, b=1, col="green", lwd=2)
```



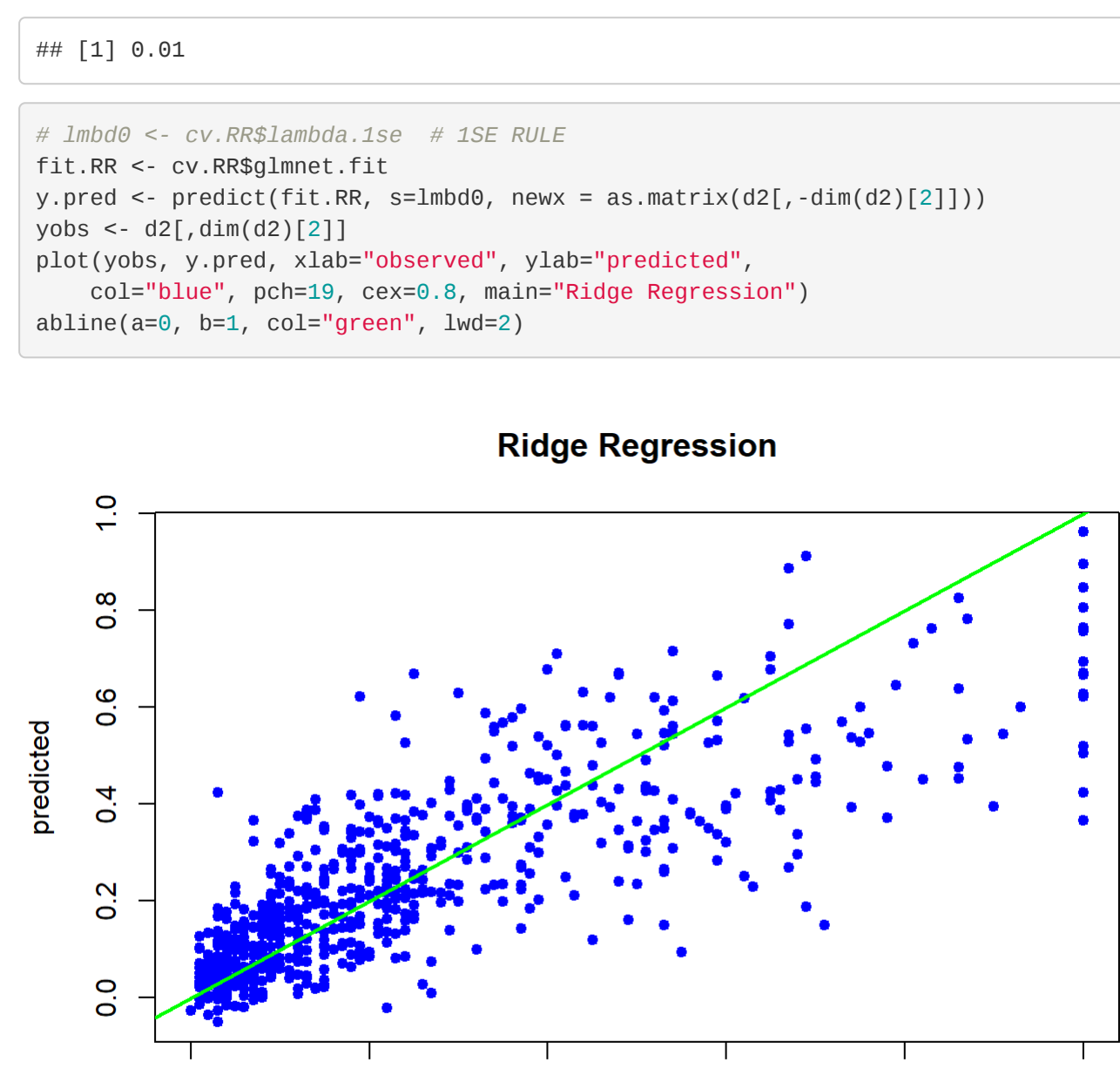
```
MSEP.lr <- mean((yobs - y.pred)^2); MSEP.lr
```

```
## [1] 0.02041893
```

iii> Ridge Regression

```
library(glmnet)
```

```
lambda <- seq(0, 80, 0.01) # AGAIN, APPROPRIATE ADJUSTMENT MIGHT BE NEEDED
# SETTING lambda IN glmnet GIVES RIDGE REGRESSION
X.train <- as.matrix(d1[, -dim(d1)[2]])
y.train <- d1[, dim(d1)[2]]
cv.rr <- cv.glmnet(X=train, y=y.train, alpha = 0, lambda = lambda, nfolds=10) # 10-FOLD CV BY DEFAULT
plot(cv.rr)
```



```
names(cv.rr)
## [1] "lambda"      "cvm"          "cvstd"        "cvup"         "cvlo"
## [6] "nzero"        "call"        "name"        "glmnet.fit"   "lambda.min"
## [11] "lambda.1se"   "index"
```

```
lmbd0 <- cv.rr$lambda.min; lmbd0 # MINIMUM CV ERROR
```

```
## [1] 0.01
```

```
# lmbd0 <- cv.rr$lambda.1se # 1SE RULE
fit.rr <- cv.glmnet.fit
y.pred <- predict(fit.rr, s=lmbd0, newx = as.matrix(d2[, -dim(d2)[2]]))
yobs <- d2[, dim(d2)[2]]
plot(yobs, y.pred, xlab="observed", ylab="predicted",
     col="blue", pch=19, cex=0.8, main="Ridge Regression")
abline(a=0, b=1, col="green", lwd=2)
```



```
MSEP.rr <- mean((yobs - y.pred)^2); MSEP.rr
```

```
## [1] 0.01970178
```

iv> Principal Components Regression (PCR)

```
install.packages("pls")
```

```
library(pls)
```

```
## Attaching package: 'pls'
```

```
## The following object is masked from 'package:corrplot':
```

```
##      corrplot
```

```
## The following object is masked from 'package:stats':
```

```
##      loadings
```

```
##      ncol(d1)
```

```
## [1] 181
```

```
fit.pcr <- pcr(d1$ViolentCrimesPerPop ~., ncomp=100, datad1, method = pls.options()$pcra1g,
  validation = "cv", segments = 10, segmet.type = "random", scale=TRUE)
summary(fit.pcr)
names(fit.pcr);
cv <- fit.pcr$Validation; names(cv)
```

```
## [1] "method"      "pred"          "coefficients" "gammas"      "PRESS0"
## [6] "PRESS"        "adj"           "segments"     "ncomp"
```

```
par(mfrow=c(2,1), mar=c(4,4))
plot(fit.pcr$comp, CVPRESS, xlab="number of PCs", ylab="PRESS", type="b", col="blue", lwd=2)
# plot(fit.pcr)
ncomp.best <- which.min(CVPRESS); ncomp.best
```

```
## [1] 73
```

```
# FIT THE BEST PCR MODEL WITHOUT V-FOLD CV (SETTING V=1 WOULD DO)
fit.pcr.best <- pcr(d1$ViolentCrimesPerPop ~., ncomp=ncomp.best, data=d1, method = pls.options()$pcra1g,
  segments = 1, scale=TRUE)
summary(fit.pcr.best)
```

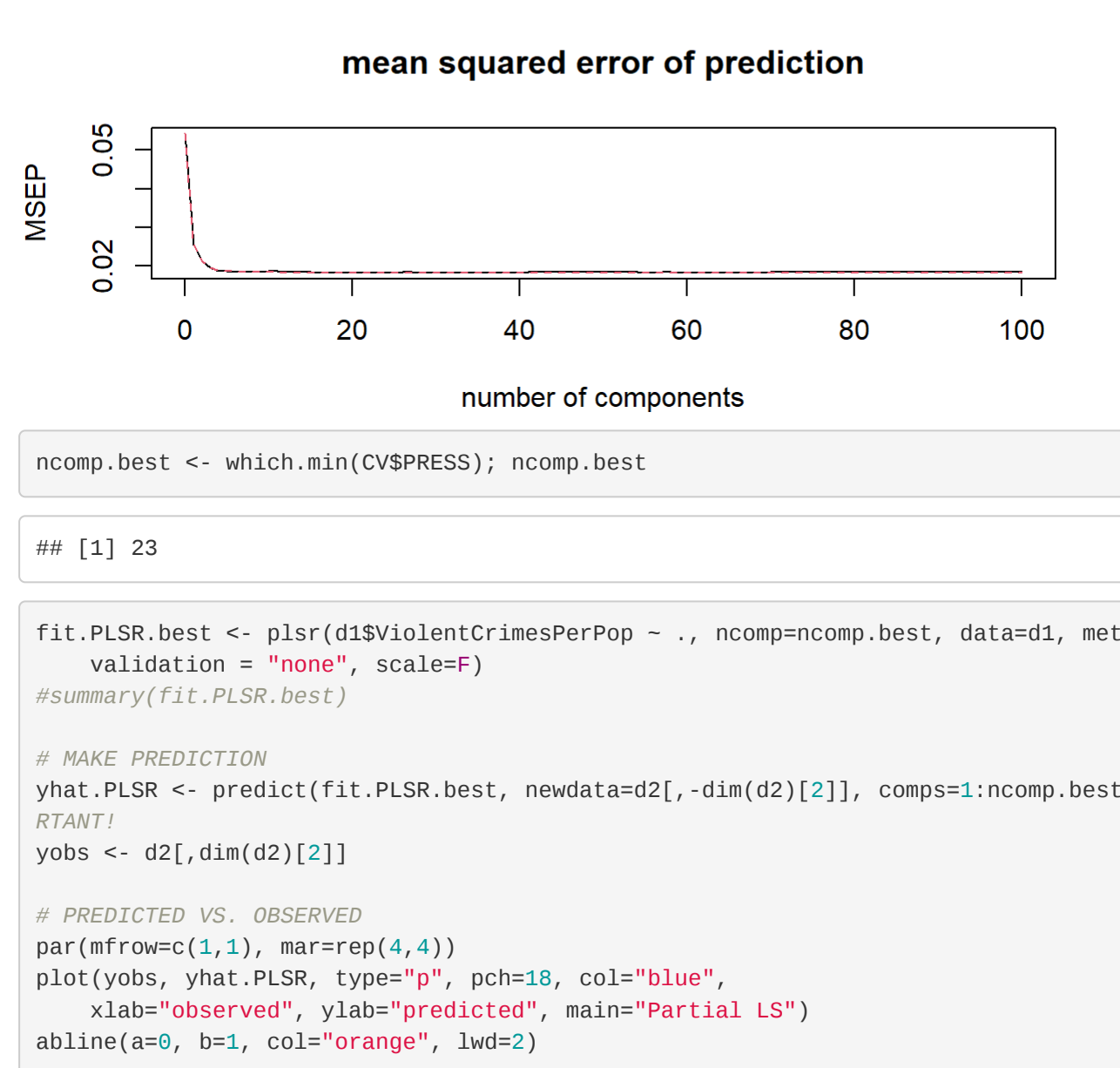
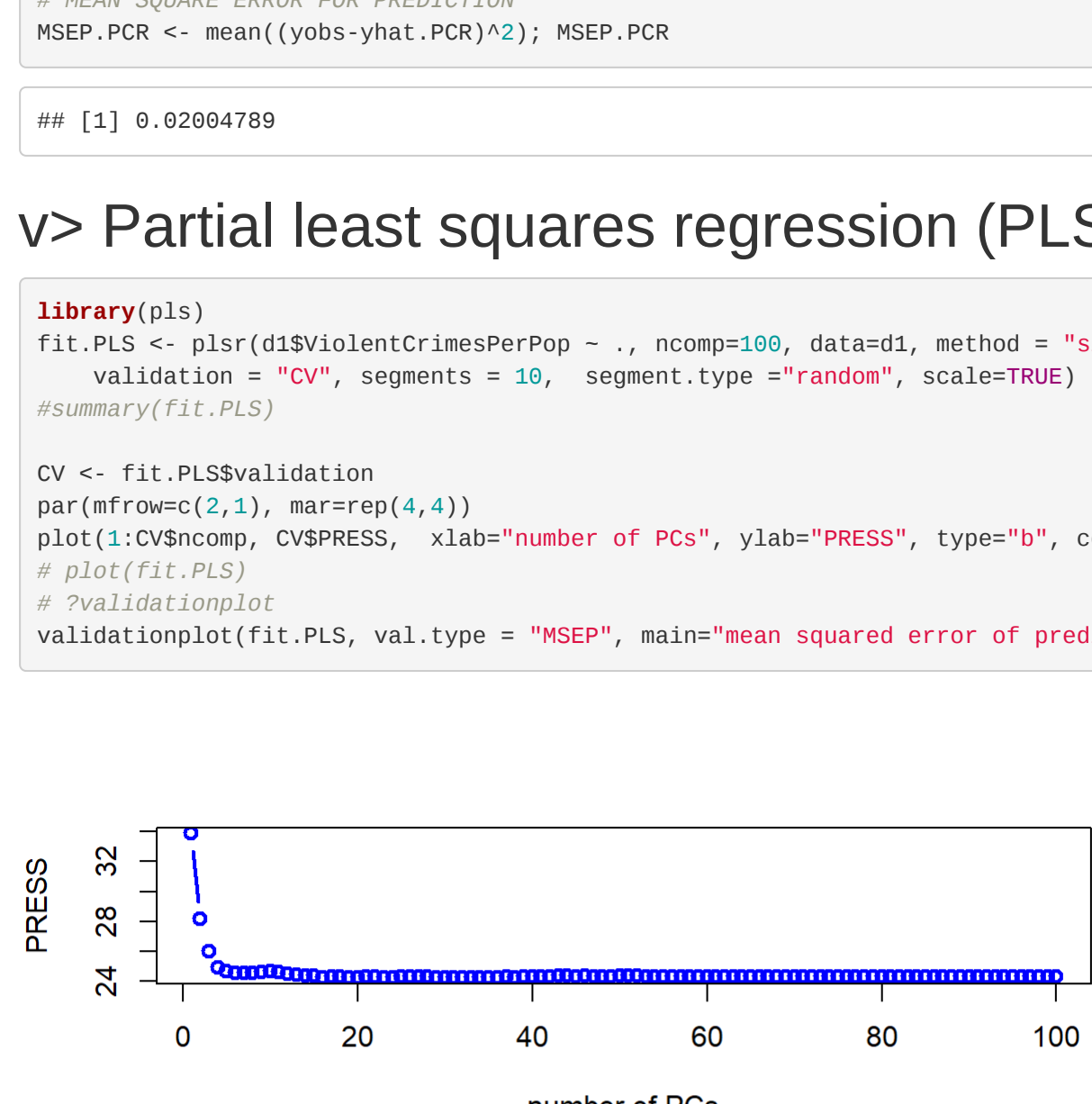
```
# PREDICTION
?predict.nvr
```

```
## starting httpd help server ...
```

```
## done
```

```
yhat.pcr <- predict(fit.pcr.best, newdata=d2[, -dim(d2)[2]], comps=1:ncomp.best) # THE ARGUMENT comps IS IMPORTANT
ANT)
yobs <- d2[, dim(d2)[2]]
```

```
# PREDICTED VS. OBSERVED
par(mfrow=c(1,1), mar=c(4,4))
plot(yobs, yhat.pcr, type="p", pch=18, col="blue",
     xlab="observed", ylab="predicted", main="PCR")
abline(a=0, b=1, col="orange", lwd=2)
```



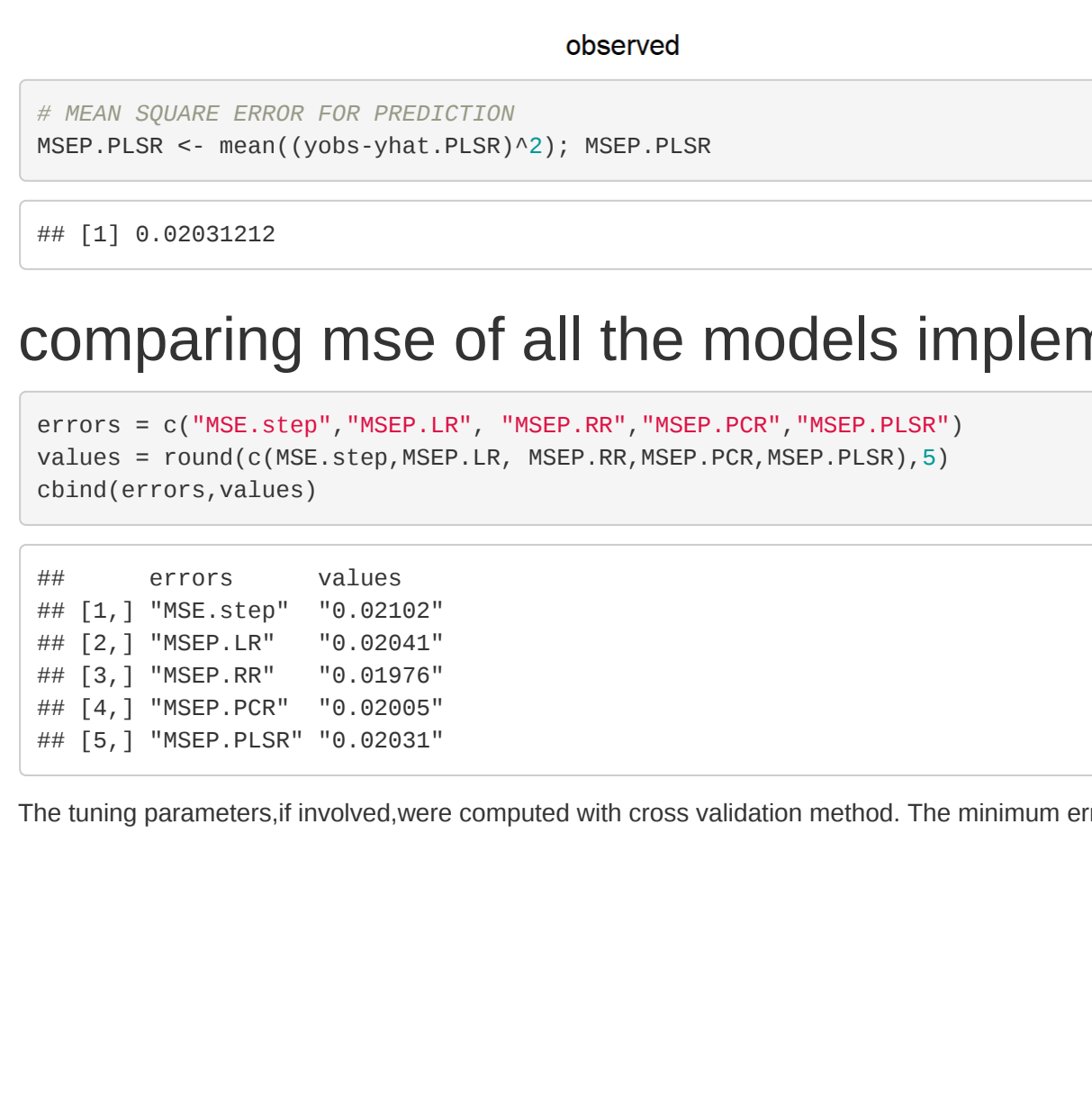
```
MSEP.pcr <- mean((yobs - yhat.pcr)^2); MSEP.pcr
```

```
## [1] 0.02084789
```

v> Partial least squares regression (PLSR)

```
library(pls)
fit.pls <- plsr(d1$ViolentCrimesPerPop ~., ncomp=100, datad1, method = "simpls",
  validation = "cv", segments = 10, segment.type = "random", scale=TRUE)
summary(fit.pls)
```

```
cv <- fit.pls$Validation
par(mfrow=c(1,1), mar=c(4,4))
plot(fit.pls$comp, CVPRESS, xlab="number of PCs", ylab="PRESS", type="b", col="blue", lwd=2)
# plot(fit.pls)
# validation summary
validationplot(fit.pls, val.type = "MSEP", main="mean squared error of prediction")
```



```
ncomp.best <- which.min(CVPRESS); ncomp.best
```

```
## [1] 23
```

```
fit.plsr.best <- plsr(d1$ViolentCrimesPerPop ~., ncomp=ncomp.best, data=d1, method = "simpls",
  validation = "none", scale=F)
summary(fit.plsr.best)
```

```
# MAKE PREDICTION
yhat.plsr <- predict(fit.plsr.best, newdata=d2[, -dim(d2)[2]], comps=1:ncomp.best) # THE ARGUMENT comps IS IMPORTANT
RANK1)
yobs <- d2[, dim(d2)[2]]
```

```
# PREDICTED VS. OBSERVED
par(mfrow=c(1,1), mar=c(4,4))
plot(yobs, yhat.plsr, type="p", pch=18, col="blue",
     xlab="observed", ylab="predicted", main="Partial LS")
abline(a=0, b=1, col="orange", lwd=2)
```



```
MSEP.plsr <- mean((yobs - yhat.plsr)^2); MSEP.plsr
```

```
## [1] 0.02031212
```

comparing mse of all the models implemented above

```
errors <- c("MSE.step", "MSEP.lr", "MSEP.rr", "MSEP.pcr", "MSEP.plsr")
errors <- round(c(MSEP.step, MSEP.lr, MSEP.rr, MSEP.pcr, MSEP.plsr), 5)
colnames(errors) <- values
```

```
##      errors      values
```

```
## [1] "MSE.step"      "0.02182"
```

```
## [2] "MSEP.lr"         "0.02042"
```

```
## [3] "MSEP.rr"         "0.01970"
```

```
## [4] "MSEP.pcr"        "0.02085"
```

```
## [5] "MSEP.plsr"       "0.02031"
```

The tuning parameters, if involved, were computed with cross validation method. The minimum error is for Ridge Regression.