

# PH C240D/STAT C245D: Assignment 2

Rachael Phillips

10/19/2017

## Regularized Regression

### Question 1. Ridge and LASSO regression: Orthonormal covariates.

Consider the linear regression model for which  $E[\mathbf{Y}_n|\mathbf{X}_n] = \mathbf{X}_n\beta$  and  $Cov[\mathbf{Y}_n|\mathbf{X}_n] = \sigma^2\mathbf{I}_n$ . Derive closed-form expressions for the ordinary least squares (OLS), ridge, and LASSO estimators of the regression coefficients  $\beta$  in the special case of *orthonormal covariates*, i.e.,  $\mathbf{X}_n^\top \mathbf{X}_n = \mathbf{I}_J$ . Provide the effective degrees of freedom, bias, and covariance matrix of the ridge regression estimator.

#### Solution:

Note that for orthonormal covariates,  $\mathbf{X}_n^\top \mathbf{X}_n = \mathbf{I}_J$ ,  $\det(\mathbf{X}_n)^2 = 1$ . In particular,  $\det(\mathbf{X}_n) \neq 0$  so  $\mathbf{X}_n$  is non-singular and has full rank.

#### OLS

We will begin with derivation of a closed-form expression for the ordinary least squares (OLS) estimator of the regression coefficients  $\beta$ . According to equation (10) on the ‘Regularized Regression’ lecture slides, for a design matrix of full column rank,  $\hat{\beta}_n^{\text{OLS}} = (\mathbf{X}_n^\top \mathbf{X}_n)^{-1} \mathbf{X}_n^\top \mathbf{Y}_n$ . Because our covariates are orthonormal (i.e. linearly independent) we have that  $\hat{\beta}_n^{\text{OLS}} = \mathbf{X}_n^\top \mathbf{Y}_n$ . For the parameter  $\beta = (\beta_j : j = 1, \dots, J) \in \mathbf{R}^J$ , a J-dimensional column vector of regression coefficients with  $\beta_j$  the  $j$ th regression coefficient and  $X = (X_j : j = 1, \dots, J) \in \mathbf{R}^J$ , a J-dimensional row vector of covariates, with  $X_j$  the  $j$ th covariate, we have that the OLS estimator of the  $j$ th regression coefficient  $\beta_{n,j}$  is  $\hat{\beta}_{n,j}^{\text{OLS}} = \sum_{i=1}^n X_{i,j} Y_i$ .

#### LASSO

Next, we consider the derivation of closed-form expression for the LASSO regression estimator of the regression coefficients  $\beta$ . According to equation (27) on the ‘Regularized Regression’ lecture slides,  $\hat{\beta}_n^{\text{LASSO}} \equiv \arg \min_{\beta \in \mathbf{R}^J} \|\mathbf{Y}_n - \mathbf{X}_n\beta\|_2^2 + \lambda \|\beta\|_1 = \arg \min_{\beta \in \mathbf{R}^J} \sum_{i=1}^n (Y_i - \sum_{j=1}^J \beta_j X_{i,j})^2 + \lambda \sum_{i=1}^J |\beta_j|$

#### Ridge

Finally, we consider the derivation of closed-form expression for the ridge estimator of the regression coefficients  $\beta$ . We also provide the effective degrees of freedom, bias, and covariance matrix of the ridge regression estimator. According to equation (20) on the ‘Regularized Regression’ lecture slides,  $\hat{\beta}_n^{\text{ridge}} = (\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1} \mathbf{X}_n^\top \mathbf{Y}_n$ . Since  $\mathbf{X}_n^\top \mathbf{X}_n = \mathbf{I}_J$  we have that  $\hat{\beta}_n^{\text{ridge}} = (1 + \lambda)^{-1} \mathbf{X}_n^\top \mathbf{Y}_n = \frac{1}{1+\lambda} \mathbf{X}_n^\top \mathbf{Y}_n = \frac{1}{1+\lambda} \hat{\beta}_n^{\text{OLS}}$ . Thus, for the parameter  $\beta = (\beta_j : j = 1, \dots, J) \in \mathbf{R}^J$ , a J-dimensional column vector of regression coefficients with  $\beta_j$  the  $j$ th regression coefficient and  $X = (X_j : j = 1, \dots, J) \in \mathbf{R}^J$ , a J-dimensional row vector of covariates, with  $X_j$  the  $j$ th covariate, we have that the ridge regression estimator of the  $j$ th regression coefficient  $\beta_{n,j}$  is  $\hat{\beta}_{n,j}^{\text{ridge}} = \frac{1}{1+\lambda} \sum_{i=1}^n X_{i,j} Y_i$

According to equation (21) on the ‘Regularized Regression’ lecture slides,  $E[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] = (\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1} \mathbf{X}_n^\top \mathbf{X}_n \beta = \frac{1}{1+\lambda} \beta$ . The bias is as follows,  $\text{Bias}[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] = E[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] - E[\mathbf{Y}_n|\mathbf{X}_n] = E[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] - \beta = \frac{1}{1+\lambda} \beta - \beta = (\frac{1}{1+\lambda} - 1) \beta = -\frac{\lambda}{1+\lambda} \beta$ .

According to equation (22) on the ‘Regularized Regression’ lecture slides, the covariance matrix of the ridge regression estimator,  $Cov[\hat{\beta}_n^{\text{ridge}} | \mathbf{X}_n] = \sigma^2 (\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1} \mathbf{X}_n^\top \mathbf{X}_n (\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1} = \sigma^2 \frac{1}{1+\lambda} \mathbf{I}_J \frac{1}{1+\lambda} = \frac{\sigma^2 \mathbf{I}_J}{(1+\lambda)^2}$ .

According to equation (26) on the ‘Regularized Regression’ lecture slides, for ridge regression, the effective degrees of freedom,  $df^{\text{ridge}}(\lambda) = \sum_{j=1}^J \frac{L_j^2}{L_j^2 + \lambda}$  where  $L_j$  are the singular values of  $\mathbf{X}_n$ . Thus, for our orthonormal covariates,  $df^{\text{ridge}}(\lambda) = \sum_{j=1}^J \frac{1}{1+\lambda} = \frac{J}{1+\lambda}$ .

## Question 2. Ridge and LASSO regression: Bayesian interpretation.

Ridge and LASSO estimators also arise within a *Bayesian* framework, as the *mode of a posterior distribution* for the regression coefficients  $\beta$ , with a suitably-chosen *prior distribution*. Consider the Gaussian linear regression model  $Y_n | X_n \sim N(X_n \beta, \sigma^2 I_n)$ , with  $\sigma^2$  known.

### a) Ridge regression.

Propose a prior distribution for  $\beta$  for which the ridge regression estimator of  $\beta$  is the mode of the posterior distribution for  $\beta$ . Comment on how the prior parameters control shrinkage.

**Solution:**

### b) LASSO regression.

Propose a prior distribution for  $\beta$  for which the LASSO regression estimator of  $\beta$  is the mode of the posterior distribution for  $\beta$ . Comment on how the prior parameters control shrinkage.

**Solution:**

## Question 3. Elastic net: Simulation study.

### a) Simulation model.

Consider the data structure  $(X, Y) \sim P$ , where  $Y \in \mathbb{R}$  is a scalar outcome and  $X = (X_j : j = 1, \dots, J) \in \mathbb{R}^J$  a  $J$ -dimensional vector of covariates. Assume the following Gaussian linear regression model

$$Y | X \sim N(X^T \beta, \sigma^2) \text{ and } X \sim N(0_{J \times 1}, \Gamma), \quad (1)$$

where  $0_{J \times 1}$  is a  $J$ -dimensional column vector of zeros and the covariance matrix  $\Gamma = (\gamma_{j,j'} : j, j' = 1, \dots, J)$  of the covariates has an autocorrelation of order 1, i.e., AR(1), structure,

$$\gamma_{j,j'} = \rho^{|j-j'|},$$

for  $\rho \in (-1, 1)$ . Set the parameter values to  $J = 10$ ,  $\beta = (-J/2 + 1, \dots, -2, -1, 0, 0, 1, 2, \dots, J/2 - 1)/10$ ,  $\sigma = 2$ , and  $\rho = 0.5$ .

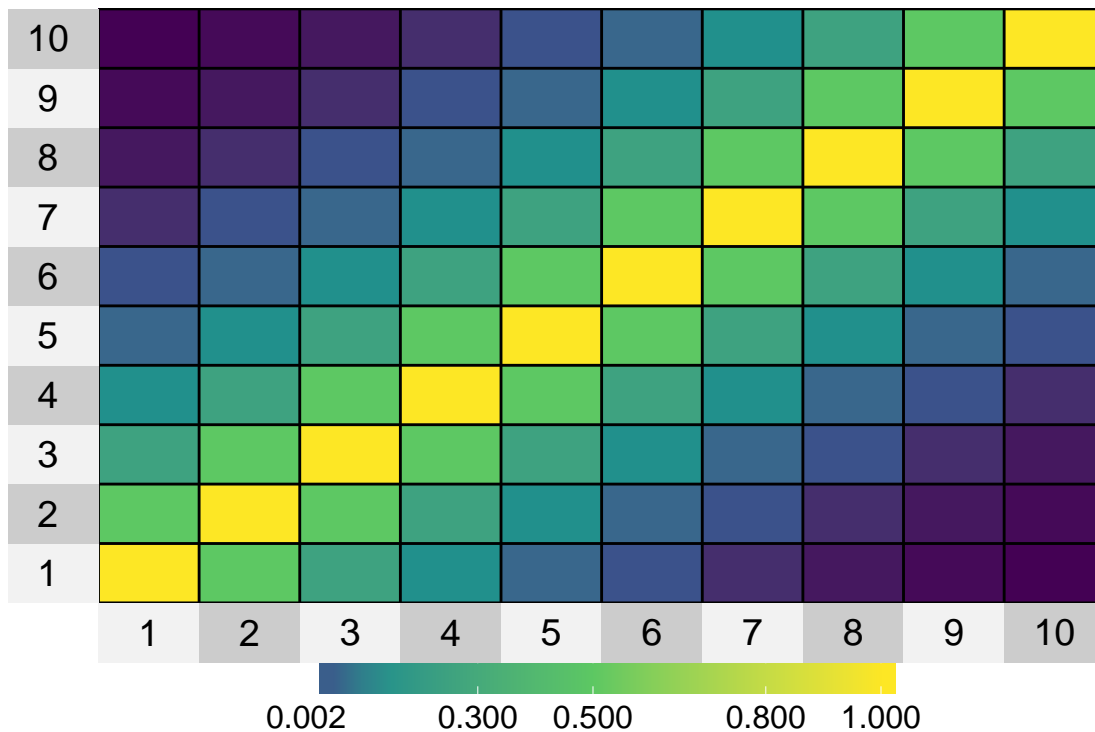
Simulate a learning set  $\mathcal{L}_n = (X_i, Y_i) : i = 1, \dots, n$  of  $n = 100$  independent and identically distributed (IID) random variables  $(X_i, Y_i) \sim P, i = 1, \dots, n$ . Also simulate an independent test set  $\mathcal{T}_{n_{TS}} = \{(X_i, Y_i) : i = 1, \dots, n_{TS}\}$  of  $n_{TS} = 1,000$  IID  $(X_i, Y_i) \sim P, i = 1, \dots, n_{TS}$ .

Provide numerical and graphical summaries of the simulation model and of the learning set.

```
#####
# Given Information
#####

# Thank you, Kelly!
J <- 10
rho <- .5
sigma <- 2

# Variance of the Covariates, gamma
gamma <- sapply(1:10, function(i){
  sapply(1:10, function(j){
    rho^(abs(i-j))
  })
})
# visualize this 10 X 10 symmetric matrix
# with i columns
superheat(gamma)
```

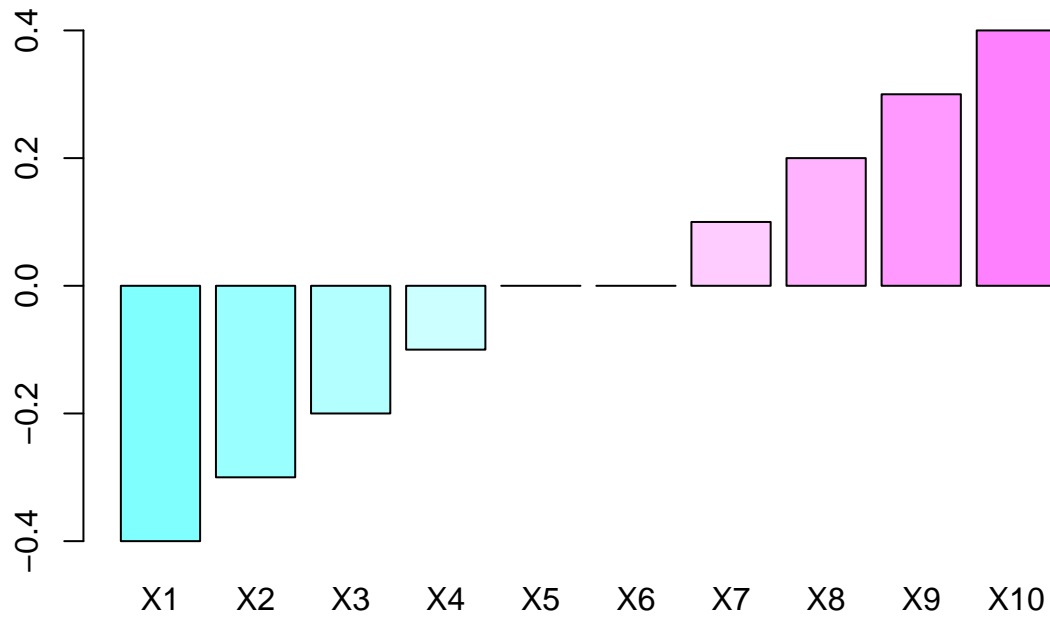


```
# regression coefficients, beta
Beta <- c((-J/2+1):0,0:(J/2-1))/10
Beta

## [1] -0.4 -0.3 -0.2 -0.1  0.0  0.0  0.1  0.2  0.3  0.4

barplot(Beta, main = expression("Beta Coefficients for Each Covariate"),
col=cm.colors(10), names.arg = c("X1", "X2", "X3", "X4", "X5", "X6", "X7",
                                "X8", "X9", "X10"))
```

Beta Coefficients for Each Covariate



```
#####
# Learning Set Simulation
#####

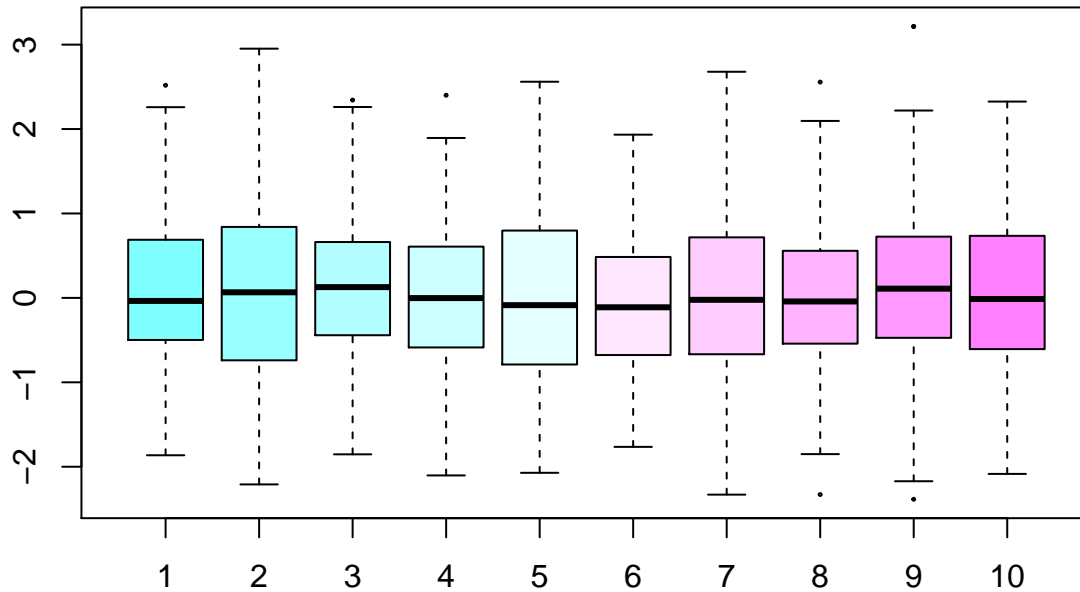
n <- 100

# Covariates
X_L <- mvrnorm(n = n, mu = rep(0,J), Sigma = gamma)
dim(X_L)

## [1] 100 10

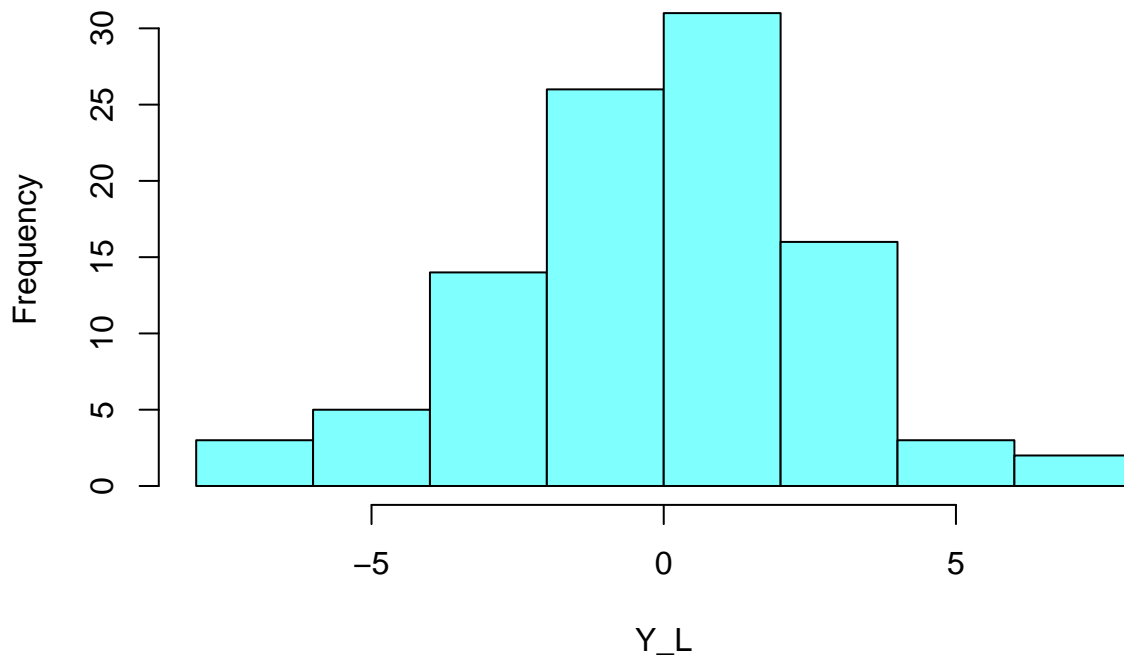
boxplot(X_L, cex=.2, col=cm.colors(10),
        main = "Boxplots of Learning Set Covariates")
```

## Boxplots of Learning Set Covariates



```
# Outcome
Y_L <- rnorm(n = n, mean = X_L %*% Beta, sd = sigma)
hist(Y_L,col=cm.colors(1), main = "Histogram of Learning Set Outcome")
```

## Histogram of Learning Set Outcome



```
summary(Y_L)
```

```
##      Min.   1st Qu.   Median     Mean   3rd Qu.     Max.
## -7.94926 -1.67424  0.27349  0.01373  1.64789  6.40388
```

```

# Combine Outcome and Covariates to make Learning Set
# we know there should be negative correlation here
Lset_X <- X_L[(1:100),]
Lset_Y <- Y_L[(1:100)]
Lset <- cbind(Lset_X,Lset_Y)
colnames(Lset) <- c("X1", "X2", "X3", "X4", "X5", "X6", "X7",
                    "X8", "X9", "X10", "Y")
summary(Lset)

```

```

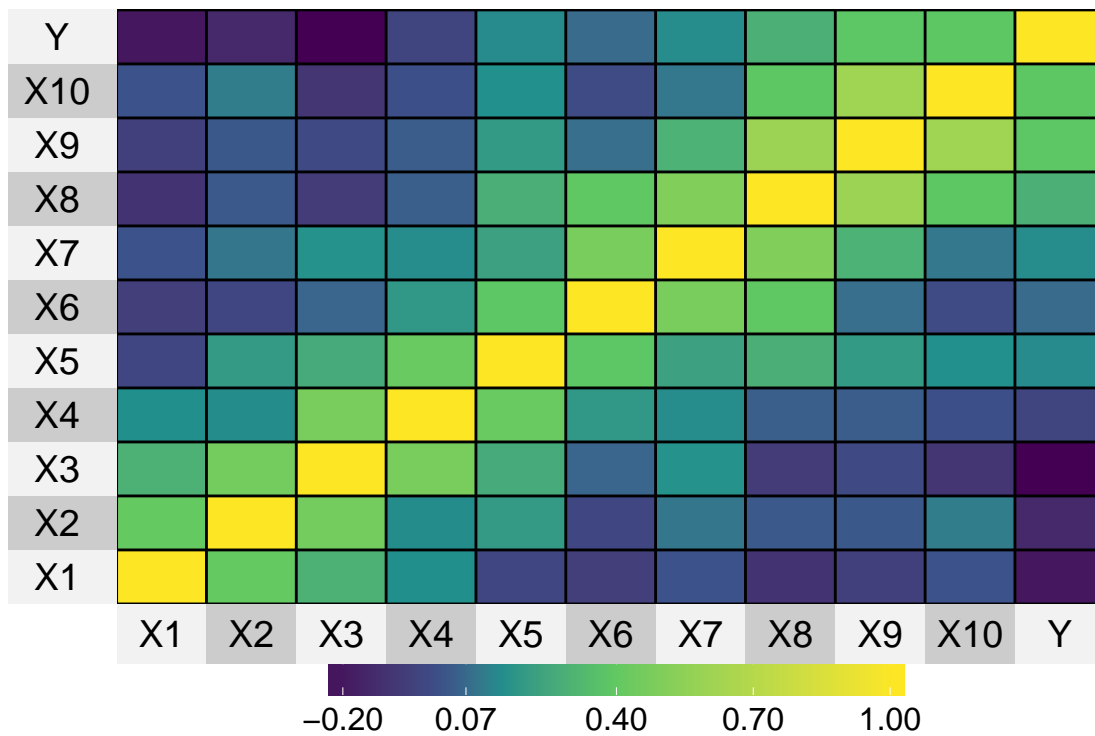
##           X1                X2                X3
## Min.      :-1.86414   Min.      :-2.20956   Min.      :-1.85300
## 1st Qu.: -0.49665   1st Qu.: -0.71494   1st Qu.: -0.44068
## Median : -0.03600   Median :  0.06624   Median :  0.12743
## Mean      : 0.09542   Mean       : 0.12134   Mean       : 0.08312
## 3rd Qu.:  0.67229   3rd Qu.:  0.83893   3rd Qu.:  0.65795
## Max.      : 2.51815   Max.       : 2.95208   Max.       : 2.34330
##           X4                X5                X6
## Min.      :-2.103144  Min.      :-2.07192   Min.      :-1.76473
## 1st Qu.: -0.584375   1st Qu.: -0.78750   1st Qu.: -0.67404
## Median : -0.002174   Median : -0.08614   Median : -0.11134
## Mean      : 0.021903   Mean       : -0.03050   Mean       : -0.09834
## 3rd Qu.:  0.582663   3rd Qu.:  0.79609   3rd Qu.:  0.47544
## Max.      : 2.400836   Max.       : 2.56048   Max.       : 1.93359
##           X7                X8                X9
## Min.      :-2.33042   Min.      :-2.3296297  Min.      :-2.3853
## 1st Qu.: -0.66005   1st Qu.: -0.5262495   1st Qu.: -0.4631
## Median : -0.02271   Median : -0.0421091   Median :  0.1091
## Mean      : 0.02171   Mean       : -0.0006869  Mean       : 0.1118
## 3rd Qu.:  0.71533   3rd Qu.:  0.5513693   3rd Qu.:  0.7167
## Max.      : 2.67876   Max.       : 2.5563781   Max.       : 3.2165
##           X10              Y
## Min.      :-2.08487   Min.      :-7.94926
## 1st Qu.: -0.60553   1st Qu.: -1.67424
## Median : -0.01341   Median :  0.27349
## Mean      : 0.08216   Mean       : 0.01373
## 3rd Qu.:  0.72089   3rd Qu.:  1.64789
## Max.      : 2.32551   Max.       : 6.40388

```

```

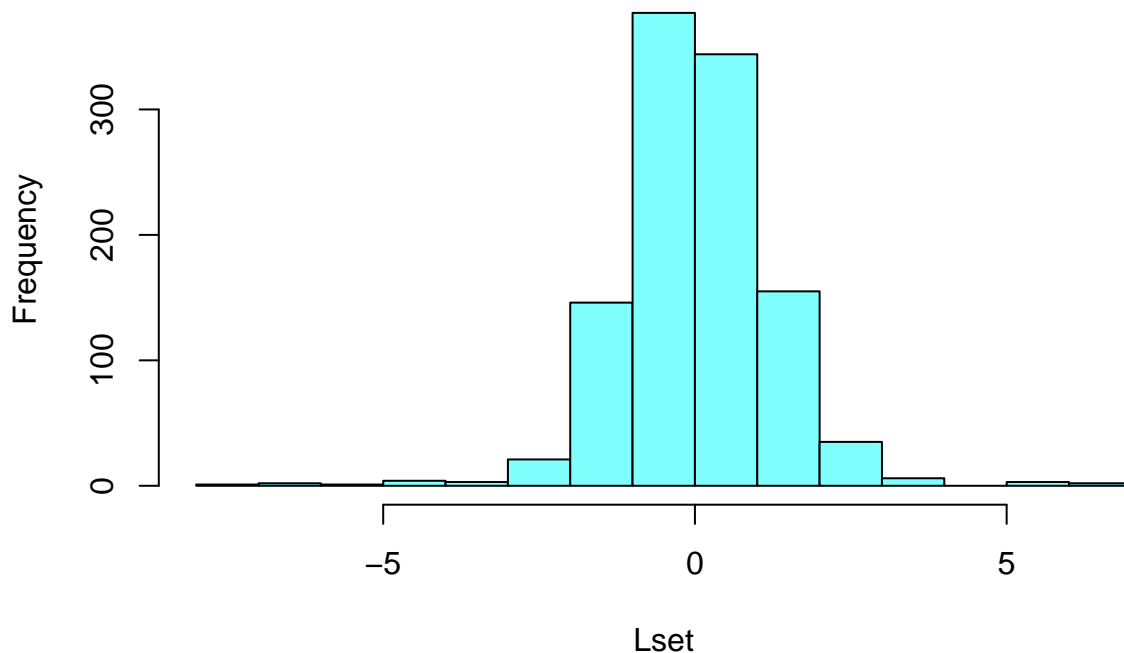
# we can visualize the correlations in a heatmap
superheat(cor(Lset))

```



```
hist(Lset, col=cm.colors(1), main = "Histogram of Learning Set")
```

**Histogram of Learning Set**



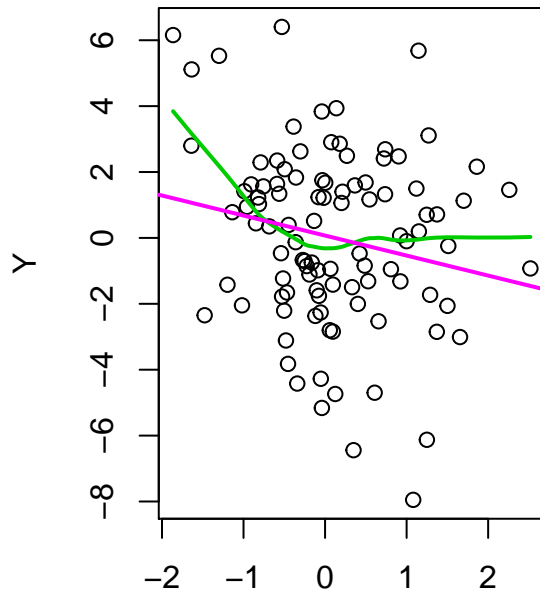
```
# Adapted from Sandrine's 'Regularized Regression:Example' code
par(mfrow=c(1,2))
for(J in 1:10){
  plot(Lset[,J], Lset_Y, xlab = colnames(Lset)[J],
       ylab = "Y", main = paste("Correlation = ", round(cor(Lset[,c(J,11)))[1,2], 2), sep = ""))
}
```

```

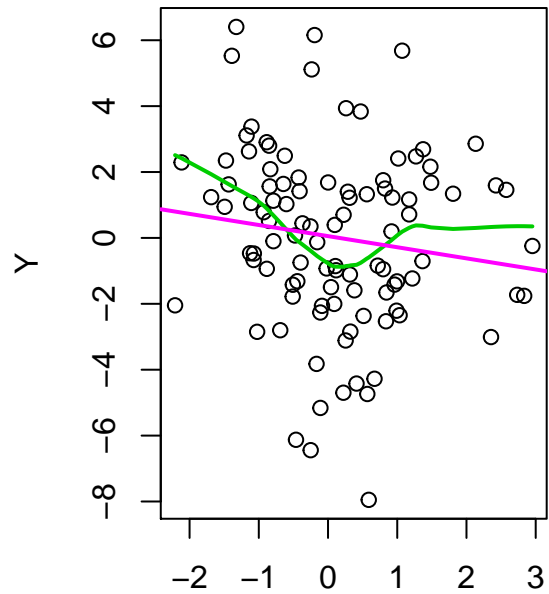
lines(lowess(Y_L ~ X_L[,J]), col=123, lwd =2)
abline(lm(Y_L ~ X_L[,J])$coef, col = 54, lwd=2)
}

```

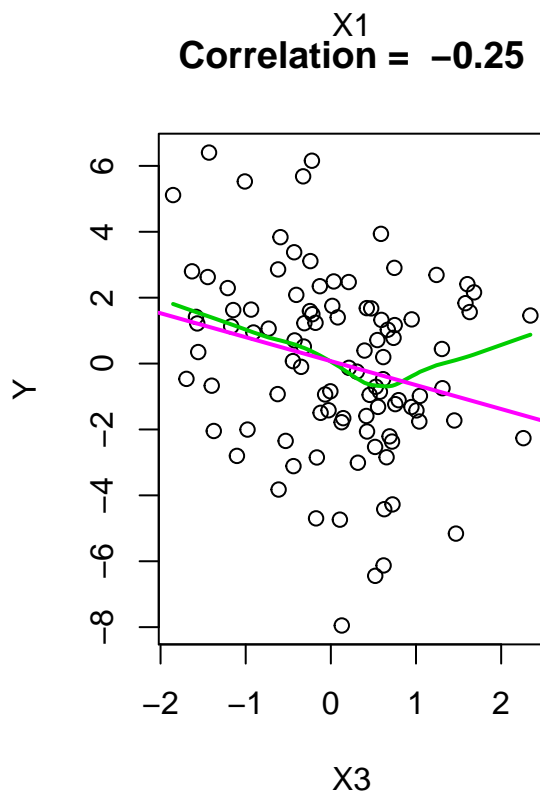
**Correlation = -0.2**



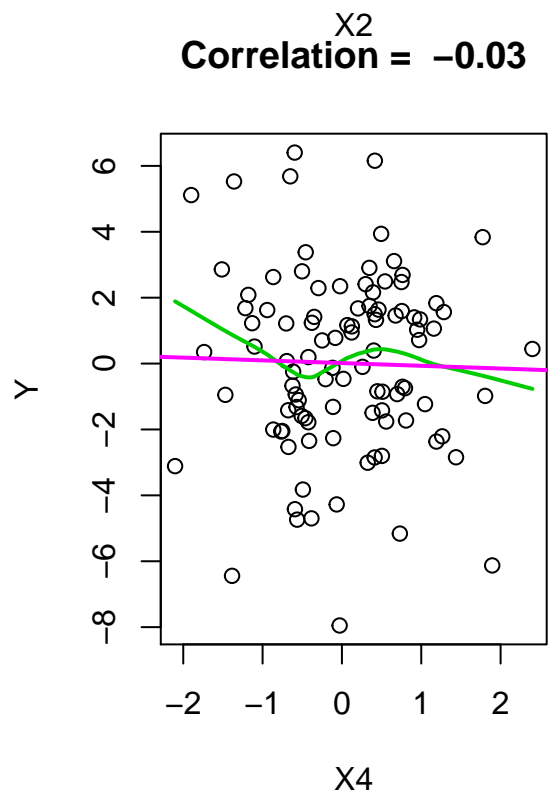
**Correlation = -0.14**



**Correlation = -0.25**

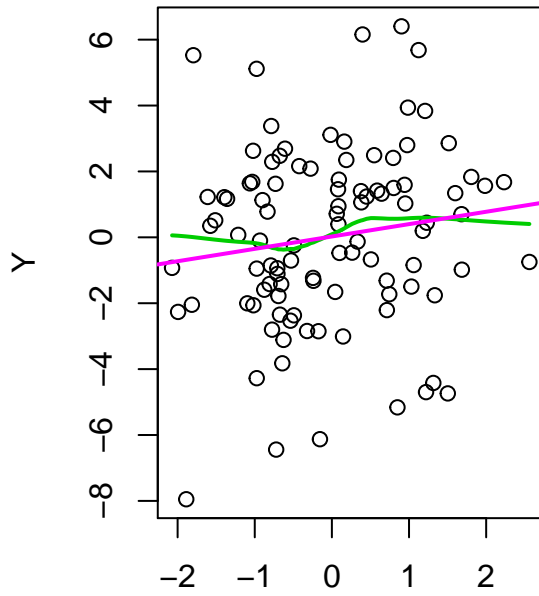


**Correlation = -0.03**

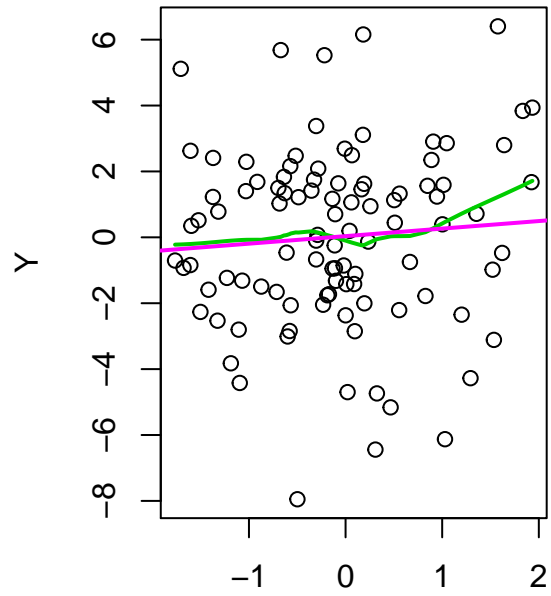




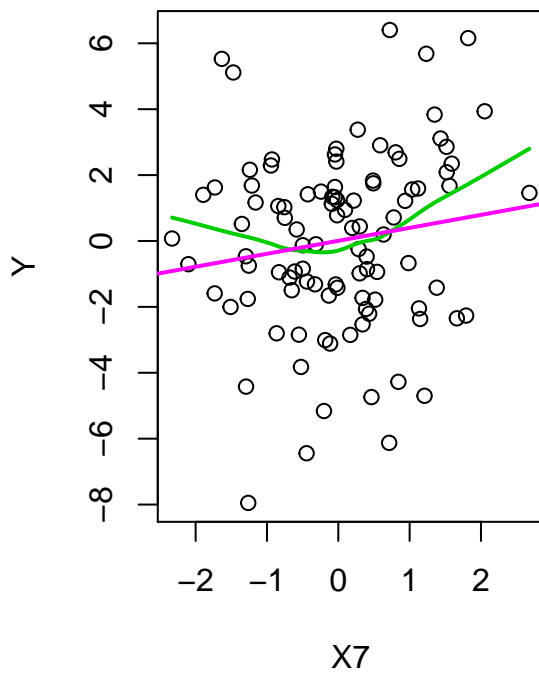
**Correlation = 0.14**



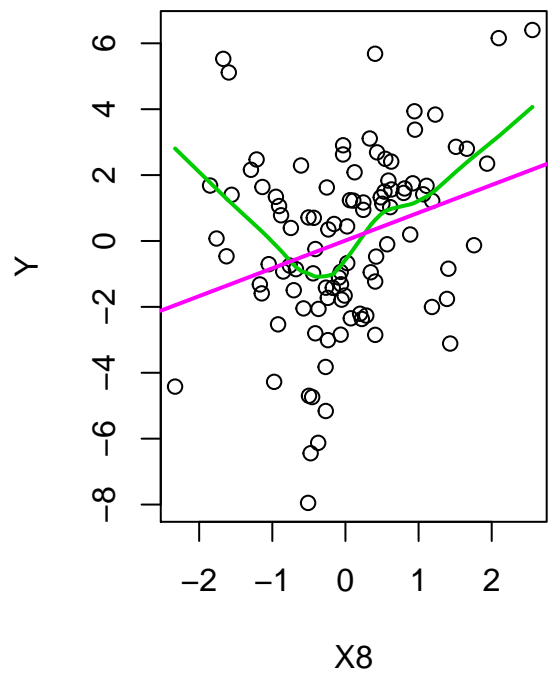
**Correlation = 0.08**



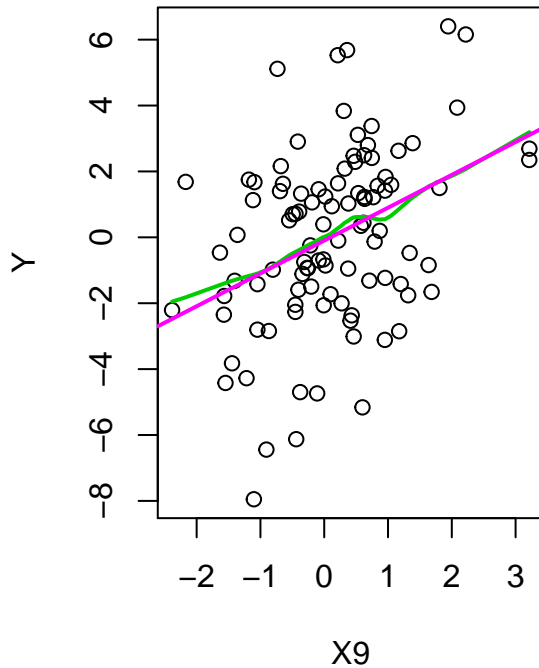
**X5  
Correlation = 0.15**



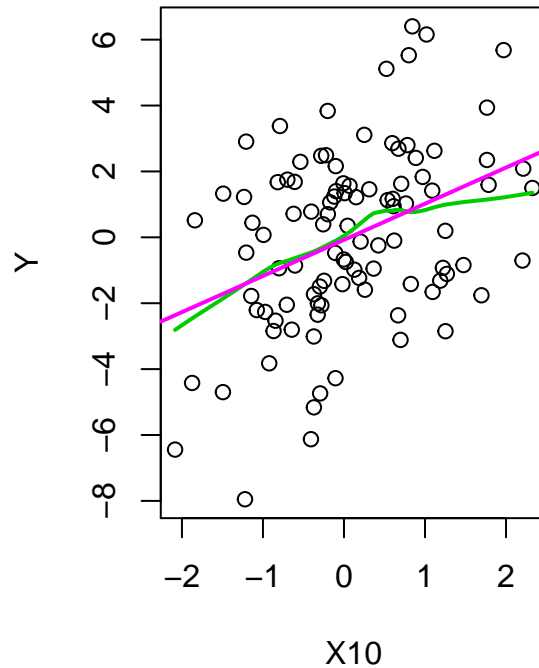
**X6  
Correlation = 0.29**



**Correlation = 0.38**



**Correlation = 0.39**



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

#####
# Test Set Simulation
#####

n <- 1000

# Covariates
X_T <- mvrnorm(n = n, mu = rep(0,J), Sigma = gamma)
dim(X_T)

## [1] 1000 10

# Outcome
Y_T <- rnorm(n = n, mean = X_T %*% Beta, sd = sigma)

# Combine Outcome and Covariates to make Test Set
Tset_X <- X_T[(1:1000),]
Tset_Y <- Y_T[(1:1000)]
Tset <- cbind(Tset_X,Tset_Y)
colnames(Tset) <- c("X1", "X2", "X3", "X4", "X5", "X6", "X7",
                   "X8", "X9", "X10", "Y")
```

#### b) Elastic net regression on learning set.

The *elastic net* estimator of the regression coefficients  $\beta$  is defined as

$$\hat{\beta}_n^{\text{enet}} \equiv \arg \min_{\beta \in \mathbb{R}^J} \|Y_n - X_n \beta\|_2^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2$$

$$= \arg \min_{\beta \in \mathbb{R}^J} \sum_{i=1}^n \left( Y_i - \sum_{j=1}^J \beta_j X_{i,j} \right)^2 + \lambda_1 \sum_{j=1}^J |\beta_j| + \lambda_2 \sum_{j=1}^J \beta_j^2$$

where the shrinkage parameters  $\lambda_1 \geq 0$  and  $\lambda_2 \geq 0$  are tuning parameters that control the strength of the penalty terms, i.e., the complexity or shrinking of the coefficients towards zero.

Obtain ridge ( $\lambda_1 = 0, \lambda_2 = \lambda$ ), LASSO ( $\lambda_1 = \lambda, \lambda_2 = 0$ ), and elastic net ( $\lambda_1 = \lambda_2 = \lambda/2$ ) estimators of the regression coefficients  $\beta$ , for  $\lambda \in \{0, 1, \dots, 100\}$ , based on the learning set simulated in a).

In particular, for each type of estimator, provide and comment on plots of the effective degrees of freedom versus the shrinkage parameter  $\lambda$  and plots of the estimated regression coefficients versus the shrinkage parameter.

For each type of estimator, obtain the learning set risk for the squared error loss function, i.e., the mean squared error (MSE),

$$MSE(\hat{\beta}_n; \mathcal{L}_n) = \frac{1}{n} \|Y_n - X_n \hat{\beta}_n\|_2^2.$$

Provide and comment on plots of the MSE versus the shrinkage parameter and report which values of the shrinkage parameter minimize risk.

**Hint.** You may use the `glmnet` function from the `glmnet` package, but be mindful of centering and scaling, of the handling of the intercept, and of the parameterization of the elastic net penalty.

```
#####
```

```
# Setting up
```

```
#####
```

```
# Directly from Sandrine's 'Regularized Regression:Example' code
```

```
## Elastic net
```

```
## N.B. alpha = lambda1/(lambda1+2*lambda2), lambda = (lambda1+2*lambda2)/(2*n)
```

```
myGlmnet <- function(x,y,x.new=NULL,intercept=TRUE,scale=TRUE,alpha=0,lambda=0,thresh=1e-12)
```

```
{
```

```
  n <- nrow(x)
```

```
  J <- ncol(x)
```

```
  xx <- scale(x,center=TRUE,scale=scale)
```

```
  beta0.hat <- mean(y)
```

```
  y.new <- NULL
```

```
  res <- glmnet(xx,y/sd(y),alpha=alpha,lambda=lambda,intercept=FALSE,standardize=FALSE,thresh=thresh)
```

```
  if(alpha == 0)
```

```
    df <- sapply(lambda*n, function(l) sum(diag(xx)%%solve(crossprod(xx)+l*diag(J))%t(xx))) + intercept
```

```
  else
```

```
    df <- rev(res$df) + intercept
```

```
  beta.hat <- as.matrix(t(coef(res)[-1,length(lambda):1])*sd(y))
```

```
  rownames(beta.hat) <- NULL
```

```
  y.hat <- t(predict(res,newx=xx,s=lambda)*sd(y))
```

```
  if(intercept)
```

```
  {
```

```
    beta.hat <- cbind(rep(beta0.hat,length(lambda)),beta.hat)
```

```
    y.hat <- y.hat + beta0.hat
```

```
  }
```

```

e <- scale(y.hat,center=y,scale=FALSE)
mse <- rowMeans(e^2)

if(!is.null(x.new))
  y.new <- t(predict(res,newx=scale(x.new,center=TRUE,scale=scale),s=lambda))*sd(y)+beta0.hat*intercept

res <- list(df=df,beta.hat=beta.hat,mse=mse,y.hat=y.hat,e=e,y.new=y.new)
res
}

# Beta vs. lambda
myPlotBeta <- function(x,beta,type="l",lwd=2,lty=1,col=1:ncol(beta),
  xlab=expression(lambda),ylab="",labels=paste(1:ncol(beta)),
  zero=TRUE,right=FALSE,main="",...)
{
  matplot(x,beta,type=type,lwd=lwd,lty=lty,col=col,xlab=xlab,ylab=ylab,main=main,...)
  if(right)
    text(x[length(x)],beta[length(x),],labels=labels,col=col)
  if(!right)
    text(x[1],beta[1,],labels=labels,col=col)
  if(zero)
    abline(h=0,lty=2)
}

#####
# Estimators
#####
lambda <- seq(0,100,by=1)

ridge <- myGlmnet(x = Lset_X, y = Lset_Y, x.new = Tset_X, alpha = 0, lambda = lambda/(nrow(Lset_X)))
lasso <- myGlmnet(x = Lset_X, y = Lset_Y, x.new = Tset_X, alpha = 1, lambda = lambda/(2*nrow(Lset_X)))
enet <- myGlmnet(x = Lset_X, y = Lset_Y, x.new = Tset_X, alpha = 1/2, lambda = (3*lambda)/(4*nrow(Lset_X)))

# Compare to lm
lm <- lm(Lset_Y~ scale(Lset_X), center = TRUE, scale = FALSE)
summary(lm)

##
## Call:
## lm(formula = Lset_Y ~ scale(Lset_X), center = TRUE, scale = FALSE)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.2172 -1.3311 -0.1405  1.5753  4.8433
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.01373    0.24330   0.056  0.9551
## scale(Lset_X)1 -0.28850    0.27648  -1.043  0.2996
## scale(Lset_X)2 -0.17900    0.30213  -0.592  0.5551
## scale(Lset_X)3 -0.59845    0.32136  -1.862  0.0659 .
## scale(Lset_X)4  0.09710    0.29982   0.324  0.7468
## scale(Lset_X)5  0.30618    0.30238   1.013  0.3140

```

```
## scale(Lset_X)6 -0.06087 0.30986 -0.196 0.8447
## scale(Lset_X)7 0.26075 0.30958 0.842 0.4019
## scale(Lset_X)8 -0.02070 0.36007 -0.057 0.9543
## scale(Lset_X)9 0.50090 0.36305 1.380 0.1711
## scale(Lset_X)10 0.64872 0.31872 2.035 0.0448 *
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.433 on 89 degrees of freedom
## Multiple R-squared: 0.2736, Adjusted R-squared: 0.192
## F-statistic: 3.352 on 10 and 89 DF, p-value: 0.0009658

ridge$beta.hat[1,]

##          V1          V2          V3          V4          V5
## 0.01373304 -0.28849736 -0.17899475 -0.59845312 0.09709517 0.30617854
##          V6          V7          V8          V9          V10
## -0.06087212 0.26074607 -0.02069877 0.50089470 0.64871828

lasso$beta.hat[1,]

##          V1          V2          V3          V4          V5
## 0.01373304 -0.28849729 -0.17899442 -0.59845274 0.09709529 0.30617798
##          V6          V7          V8          V9          V10
## -0.06087073 0.26074463 -0.02070036 0.50089761 0.64871737

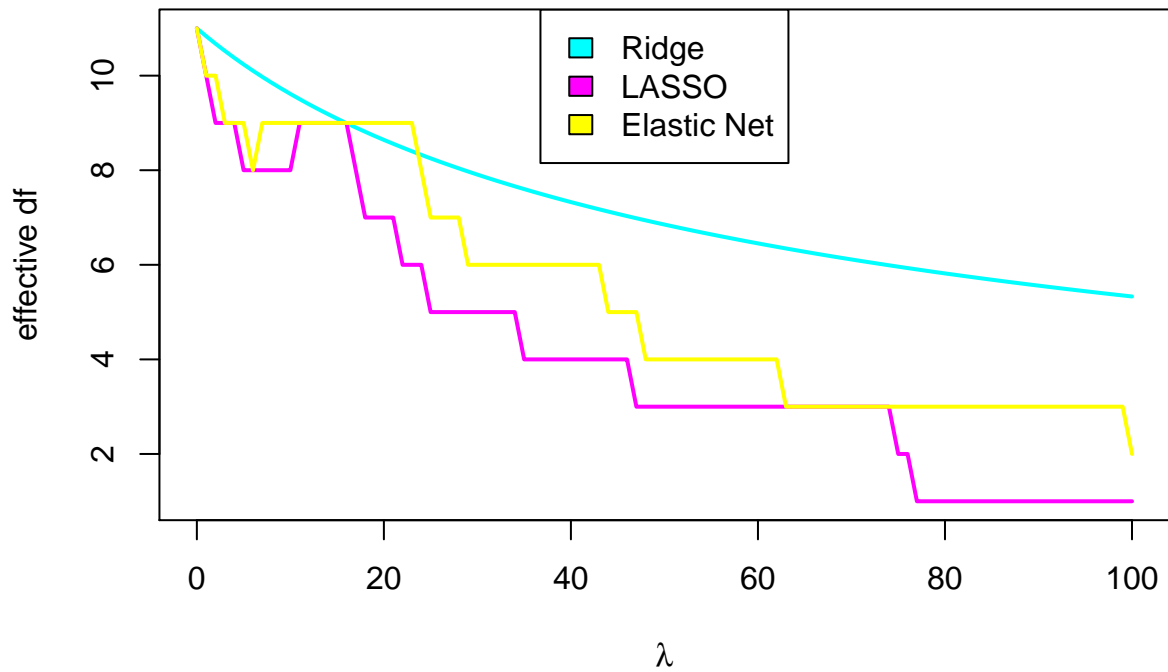
enet$beta.hat[1,]

##          V1          V2          V3          V4          V5
## 0.01373304 -0.28849744 -0.17899446 -0.59845272 0.09709544 0.30617779
##          V6          V7          V8          V9          V10
## -0.06087066 0.26074509 -0.02070059 0.50089616 0.64871833

#####
# Plots
#####

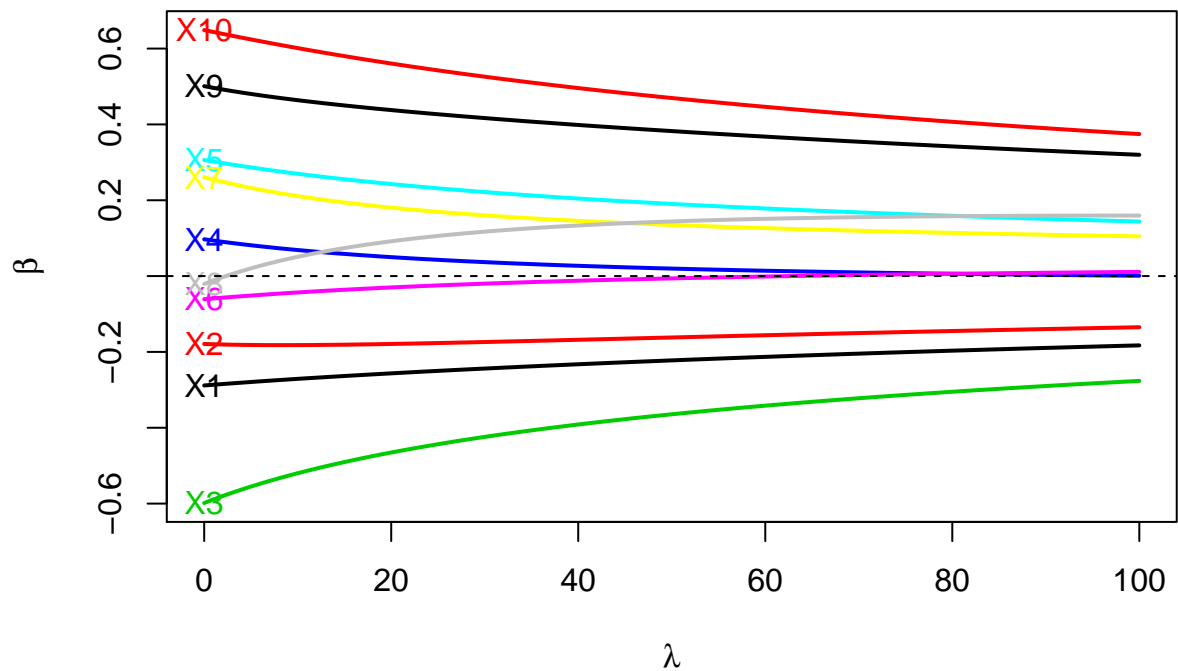
# Effective df vs. lambda
matplot(lambda, cbind(ridge$df,lasso$df, enet$df),
        type="l", lwd=2, lty=1, col=5:7,
        xlab=expression(lambda), ylab = "effective df",
        main="Learning Set: Ridge, LASSO, and Elastic Net")
legend("top", c("Ridge", "LASSO", "Elastic Net"), fill=5:7)
```

## Learning Set: Ridge, LASSO, and Elastic Net



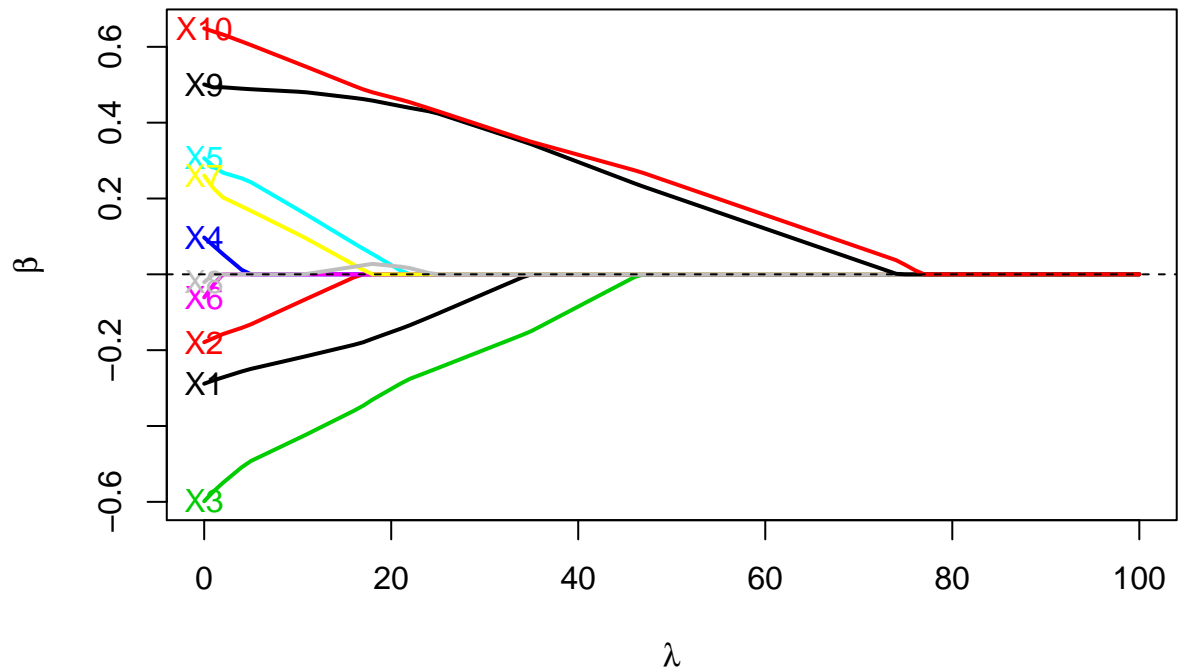
```
myPlotBeta(lambda, ridge$beta.hat[,-1],type="l",lwd=2, ylab=expression(hat(beta)),
  labels=colnames(Lset[,1:10]), main="Learning Set: Ridge")
```

## Learning Set: Ridge



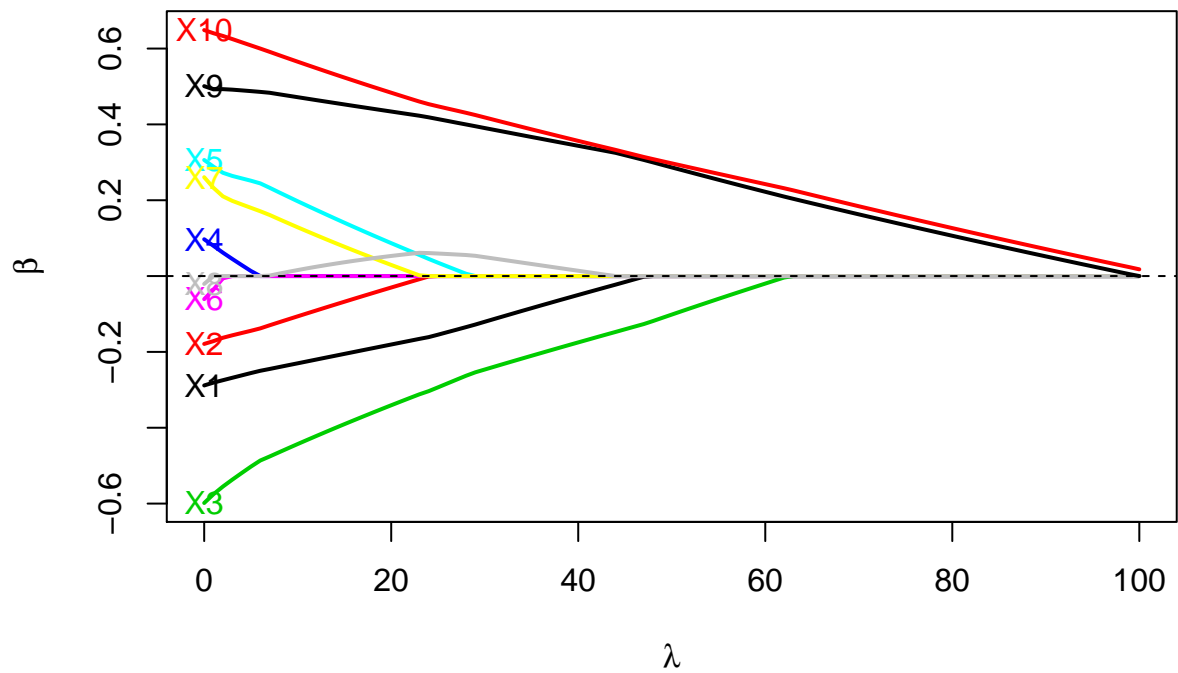
```
myPlotBeta(lambda, lasso$beta.hat[,-1],type="l",lwd=2, ylab=expression(hat(beta)),
  labels=colnames(Lset[,1:10]), main="Learning Set: LASSO")
```

### Learning Set: LASSO



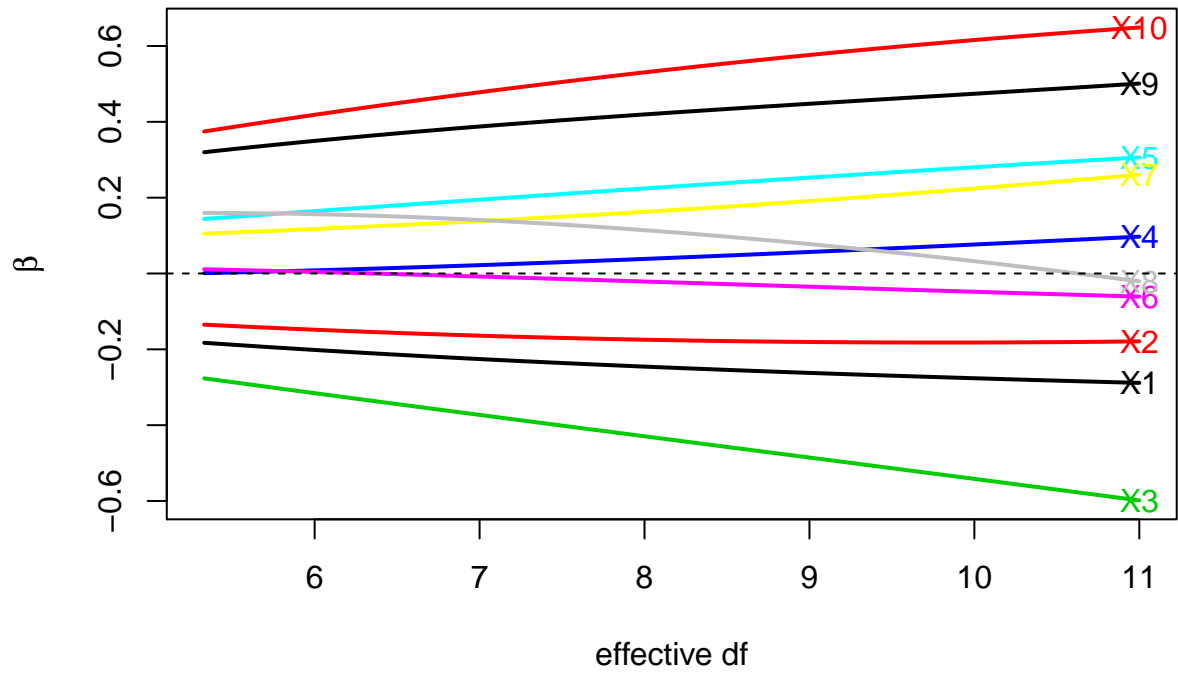
```
myPlotBeta(lambda, enet$beta.hat[,-1],type="l",lwd=2, ylab=expression(hat(beta)),
  labels=colnames(Lset[,1:10]), main="Learning Set: Elastic Net")
```

### Learning Set: Elastic Net



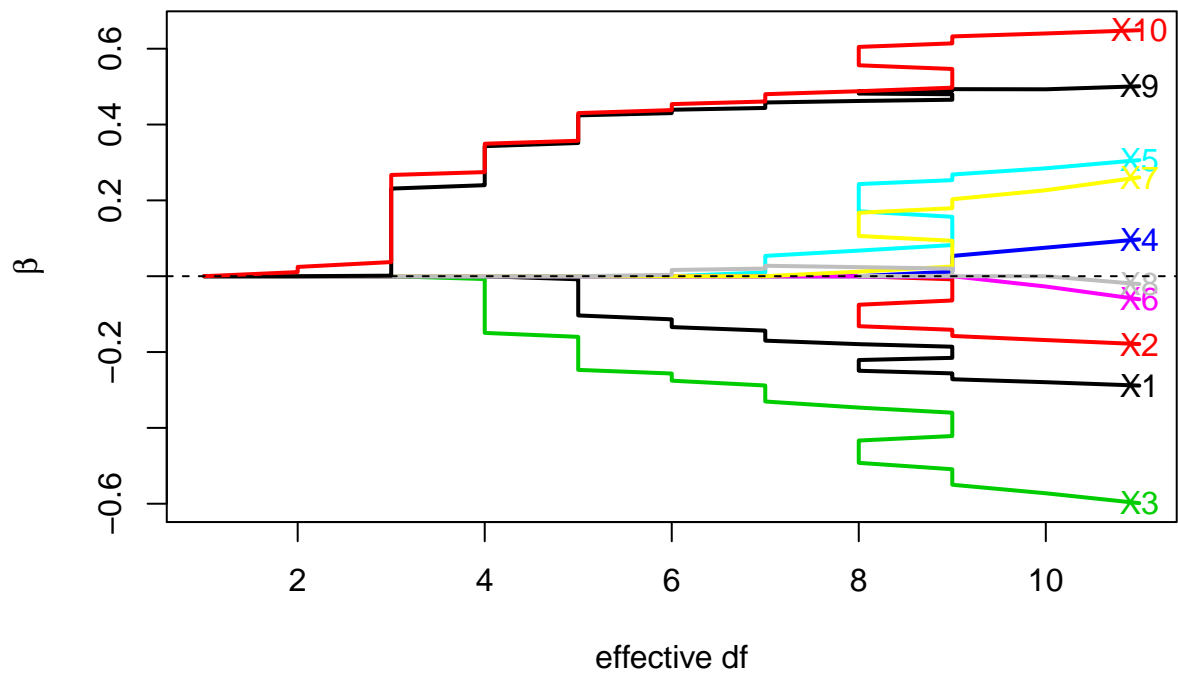
```
# Beta vs. effective df
myPlotBeta(ridge$df,ridge$beta.hat[,-1],type="l",lwd=2,xlab="effective df",
  ylab=expression(hat(beta)),labels=colnames(Lset[,1:10]),main="Learning Set: Ridge")
```

### Learning Set: Ridge



```
myPlotBeta(lasso$df,lasso$beta.hat[, -1], type="l", lwd=2, xlab="effective df",
            ylab=expression(hat(beta)), labels=colnames(Lset[, 1:10]), main="Learning Set: LASSO")
```

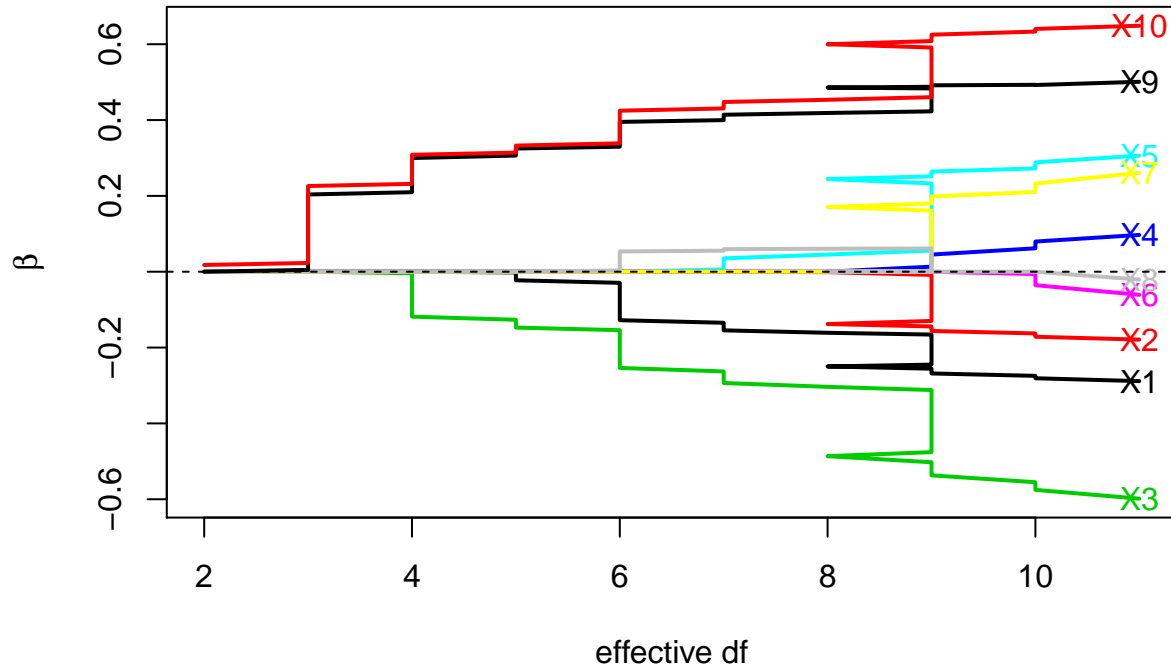
### Learning Set: LASSO





```
myPlotBeta(enet$df,enet$beta.hat[,-1],type="l",lwd=2,xlab="effective df",
  ylab=expression(hat(beta)),labels=colnames(Lset[,1:10]),main="Learning Set: Elastic Net")
```

## Learning Set: Elastic Net



```
#####
# MSE
#####

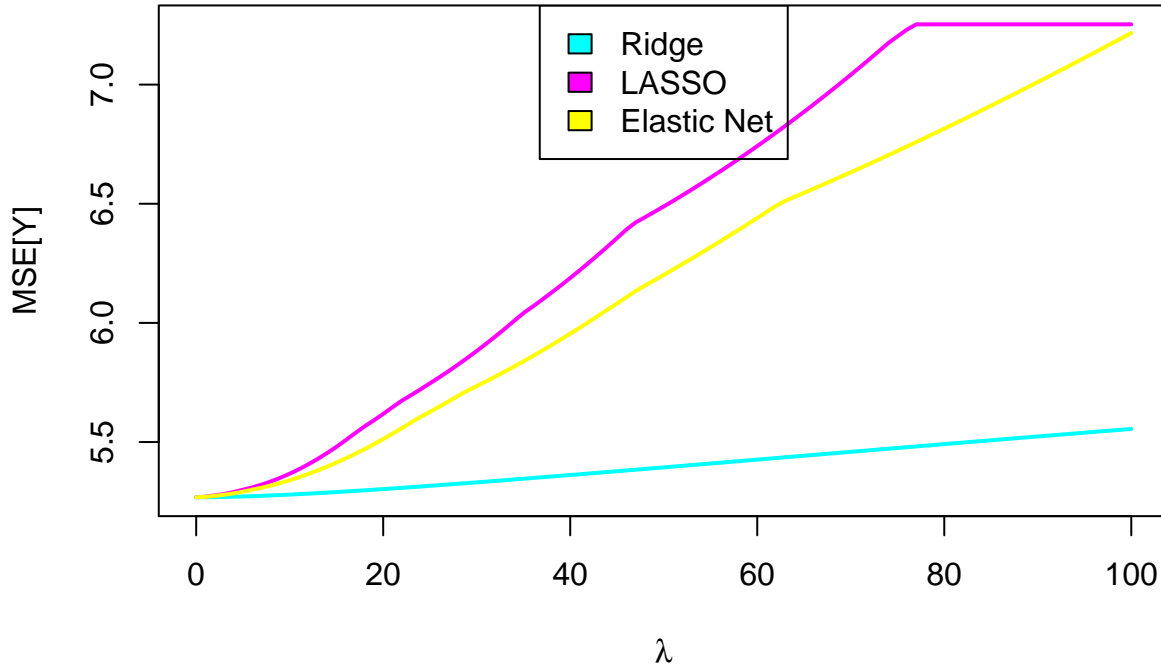
mse.Lset <- cbind("Ridge" = ridge$mse,"LASSO" = lasso$mse, "Elastic Net" = enet$mse)

summary(mse.Lset)
```

	Ridge	LASSO	Elastic Net
## Min.	:5.268	Min. :5.268	Min. :5.268
## 1st Qu.:	:5.316	1st Qu.:5.746	1st Qu.:5.627
## Median	:5.394	Median :6.489	Median :6.201
## Mean	:5.398	Mean :6.410	Mean :6.188
## 3rd Qu.:	:5.475	3rd Qu.:7.202	3rd Qu.:6.722
## Max.	:5.555	Max. :7.253	Max. :7.216

```
matplot(lambda, mse.Lset, type="l",lwd=2,lty=1,col=5:7,
  xlab=expression(lambda), ylab="MSE[Y]",
  main="Learning Set: Ridge, LASSO, and Elastic Net")
legend("top",c("Ridge", "LASSO", "Elastic Net"),fill=5:7)
```

## Learning Set: Ridge, LASSO, and Elastic Net



We notice an inverse relationship between the effective degrees of freedom and the shrinkage parameter across all three methods. For LASSO and elastic net, this happens in a stepwise manner since these are step-wise functions. Additionally, across all three methods, as the shrinkage parameter increases, the estimated regression coefficient shrinks towards zero. In fact, for large enough shrinkage parameters, the regression coefficients are set to zero. We see this occurring for the LASSO and elastic net estimators. We also note that as the effective degrees of freedom increase the estimated regression coefficient blows up, as expected. Lastly and as expected, the mean squared error (MSE) of the fitted values of the learning set increase as the shrinkage parameter increases, corresponding to the estimators become less data-adaptive. We also see that the MSE is minimized for all three types of estimators when the shrinkage parameter is set to zero.

### c) Performance assessment on test set.

For each estimator in b), obtain the test set risk  $MSE(\hat{\beta}_n; \mathcal{T}_{n_{TS}})$  for the squared error loss function (i.e., MSE). Provide and comment on plots of risk versus the shrinkage parameter and report which values of the shrinkage parameter minimize risk. Examine the corresponding three “optimal” estimators of the regression coefficients.

```
#####
# MSE
#####

mse.ridge.Tset <- rowMeans(scale(ridge$y.new, center=Tset[,11], scale=FALSE)^2)
mse.lasso.Tset <- rowMeans(scale(lasso$y.new, center=Tset[,11], scale=FALSE)^2)
mse.enet.Tset <- rowMeans(scale(enet$y.new, center=Tset[,11], scale=FALSE)^2)

l1 <- which.min(mse.ridge.Tset)
l2 <- which.min(mse.lasso.Tset)
l3 <- which.min(mse.enet.Tset)

mse.Tset <- cbind("Ridge" = mse.ridge.Tset,
```

```

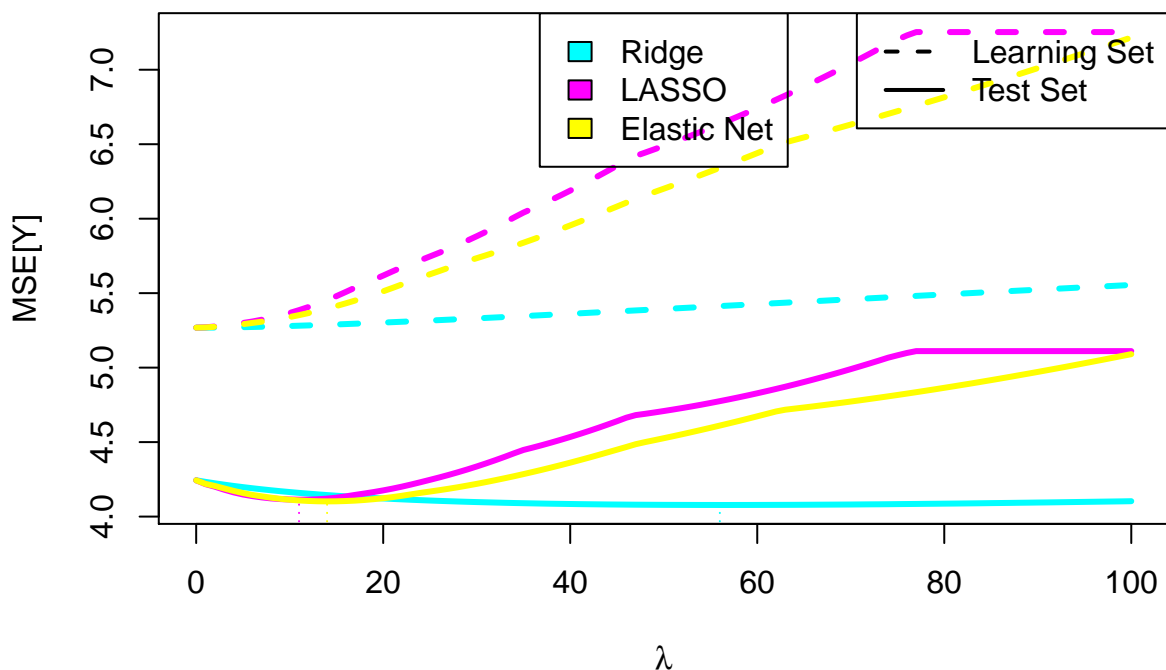
"LASSO" = mse.lasso.Tset,
"Elastic Net" = mse.enet.Tset)
summary(mse.Tset)

##      Ridge      LASSO      Elastic Net
## Min.   :4.078   Min.   :4.113   Min.   :4.102
## 1st Qu.:4.081   1st Qu.:4.247   1st Qu.:4.197
## Median :4.090   Median :4.709   Median :4.526
## Mean   :4.106   Mean   :4.658   Mean   :4.528
## 3rd Qu.:4.108   3rd Qu.:5.083   3rd Qu.:4.816
## Max.   :4.244   Max.   :5.111   Max.   :5.090

matplot(lambda, cbind(ridge$mse, mse.ridge.Tset,
                      lasso$mse, mse.lasso.Tset,
                     enet$mse, mse.enet.Tset),
        type="l", lwd=3, col=rep(5:7, each=2), lty=rep(2:1,2),
        xlab=expression(lambda), ylab="MSE[Y]",
        main="Ridge, LASSO, and Elastic Net")
legend("topright", c("Learning Set", "Test Set"), lty=2:1, lwd=2)
legend("top", c("Ridge", "LASSO", "Elastic Net"), fill=5:7)
lines(lambda[rep(11, 2)], c(0, min(mse.ridge.Tset)), col=5, lty=3)
lines(lambda[rep(12, 2)], c(0, min(mse.lasso.Tset)), col=6, lty=3)
lines(lambda[rep(13, 2)], c(0, min(mse.enet.Tset)), col=7, lty=3)

```

## Ridge, LASSO, and Elastic Net



```

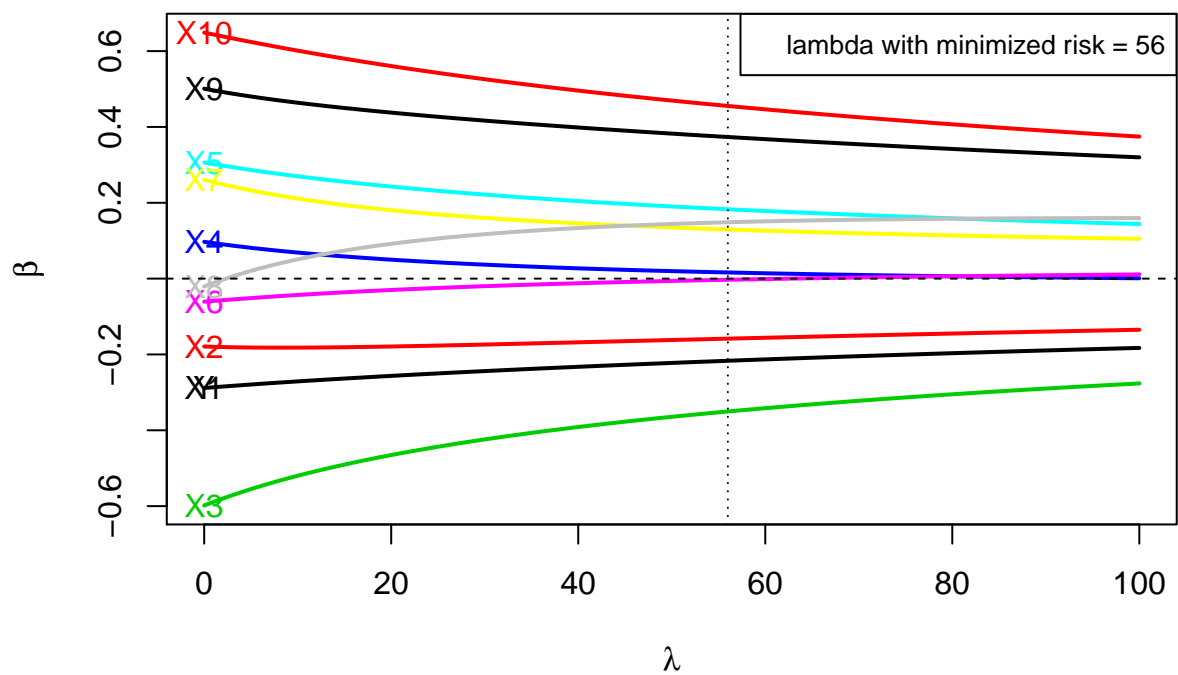
#####
# Risk v Lambda
#####

myPlotBeta(lambda, ridge$beta.hat[, -1], labels=colnames(Lset), right=FALSE,
            ylab=expression(hat(beta)), main="Ridge")

```

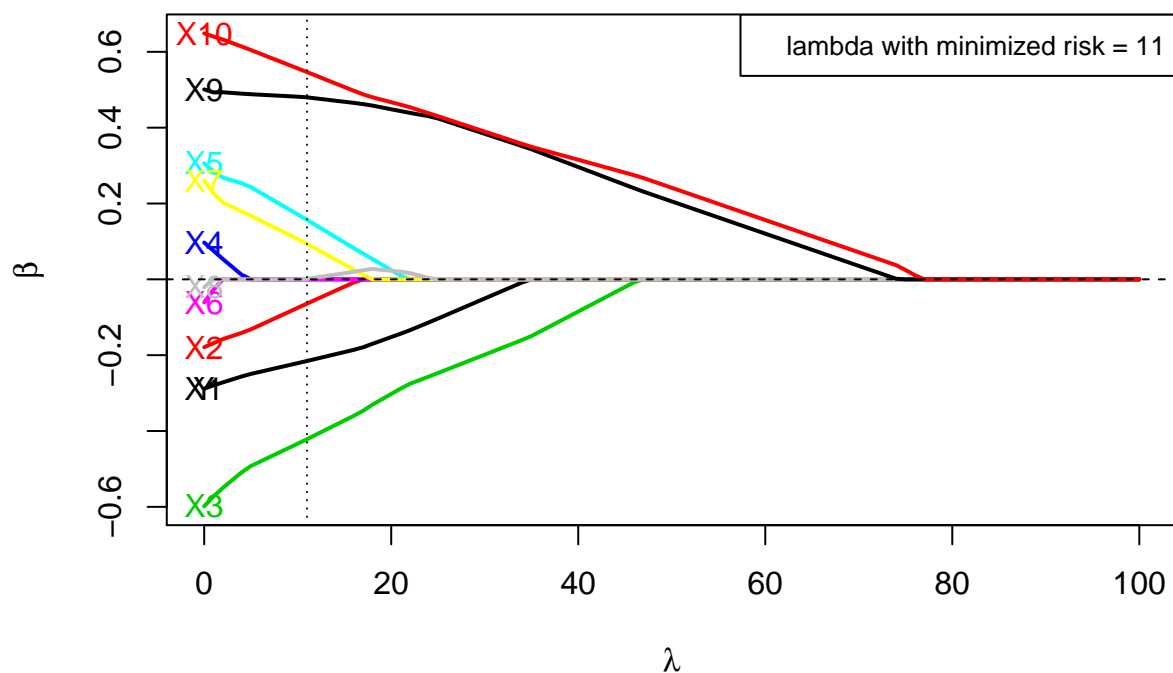
```
abline(v = lambda[11],col=1,lty=3)
legend("topright", pt.cex = 1, cex=.8,
      c(paste("lambda with minimized risk = ", lambda[11], sep = "")))
```

## Ridge



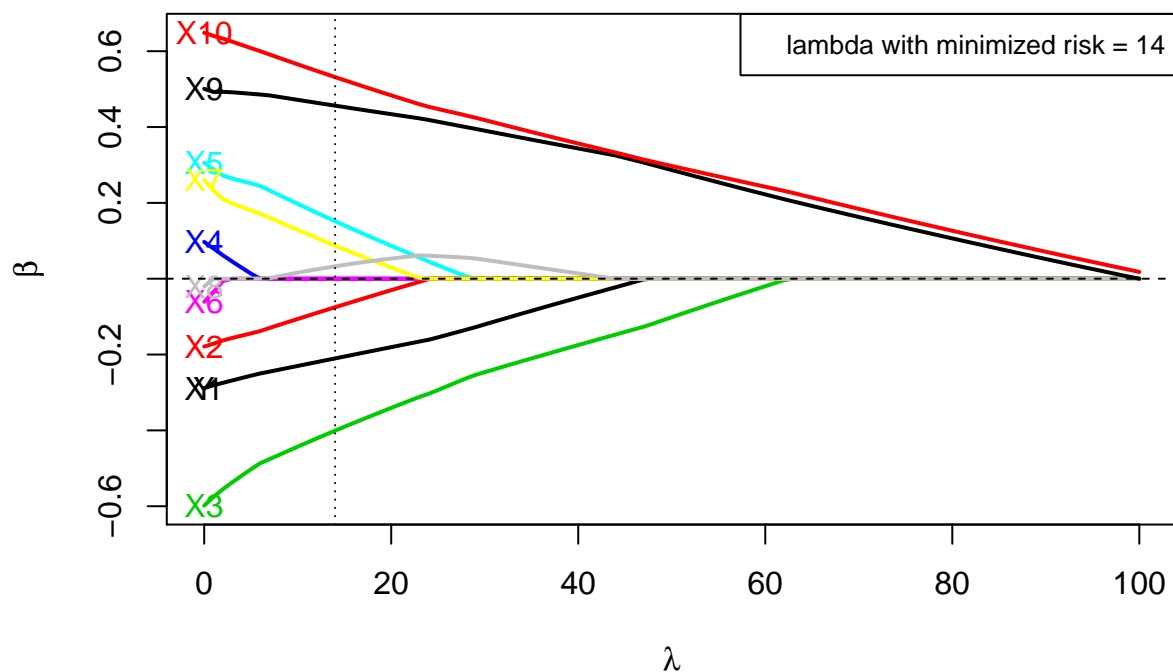
```
myPlotBeta(lambda, lasso$beta.hat[,-1], labels=colnames(Lset), right=FALSE,
            ylab=expression(hat(beta)), main="LASSO")
abline(v = lambda[12], col=1, lty=3)
legend("topright", pt.cex = 1, cex=.8,
      c(paste("lambda with minimized risk = ", lambda[12], sep = "")))
```

## LASSO



```
myPlotBeta(lambda, enet$beta.hat[,-1], labels=colnames(Lset), right=FALSE,
            ylab=expression(hat(beta)), main="Elastic Net")
abline(v = lambda[13], col=1, lty=3)
legend("topright", pt.cex = 1, cex=.8,
       c(paste("lambda with minimized risk = ", lambda[13], sep = "")))
```

## Elastic Net

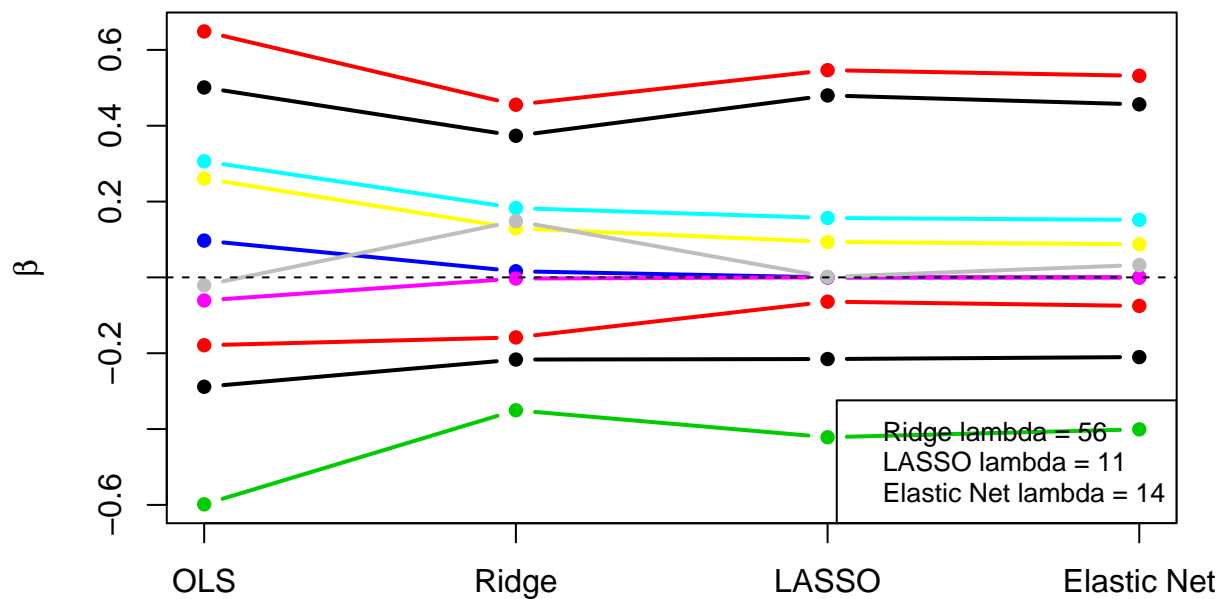


```
#####
# Optimal Beta
#####

matplot(t(cbind(ridge$beta.hat[1,-1], ridge$beta.hat[11,-1],
               lasso$beta.hat[12,-1], enet$beta.hat[13,-1])),
        type="b", lty=1, lwd=2, pch=16, col=1:ncol(X_T),
        ylab=expression(hat(beta)), axes=FALSE,
        main="Optimal Beta Coefficient")

box()
axis(1, at=1:4, c("OLS", "Ridge", "LASSO", "Elastic Net"))
axis(2)
abline(h=0, lty=2)
legend("bottomright", pt.cex = 1, cex=.8,
      c(paste("Ridge lambda = ", lambda[11], sep = ""),
        paste("LASSO lambda = ", lambda[12], sep = ""),
        paste("Elastic Net lambda = ", lambda[13], sep = "")))
```

### Optimal Beta Coefficient

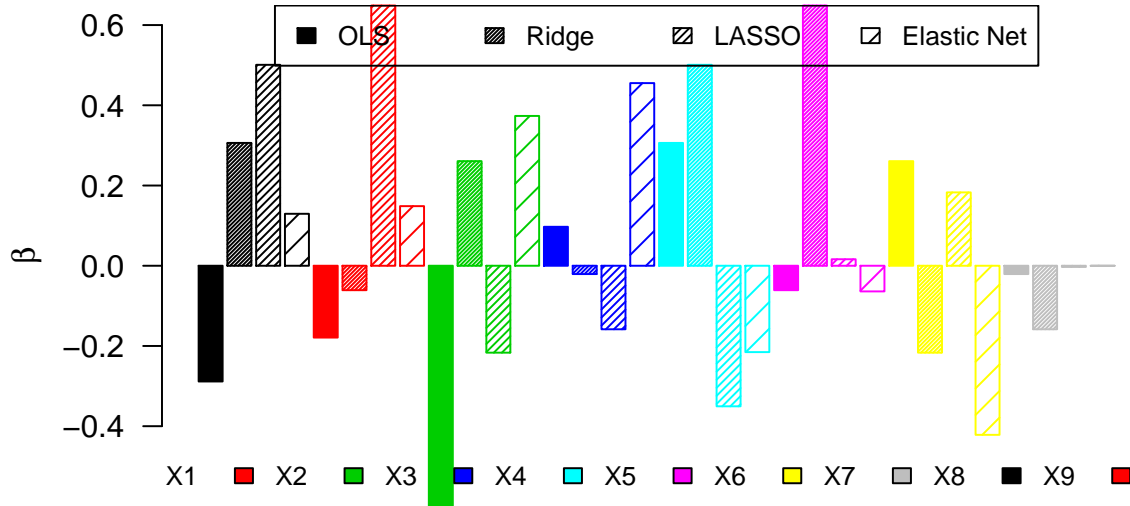


```
barplot(c(ridge$beta.hat[1,-1],ridge$beta.hat[11,-1],lasso$beta.hat[12,-1],
          enet$beta.hat[13,-1])[as.vector(sapply(1:8,function(x)c(x,x+4,x+8,x+16)))],
        col=rep(1:ncol(X_T),each=4),density=rep(c(-9,66,33,10),ncol(X_T)),
        border=TRUE, ylab=expression(hat(beta)), names.arg="",
        las=2, main="Optimal Betas for Each Covariate and Estimator")

legend("bottom", colnames(Lset)[1:10], fill=1:ncol(X_T),
      pt.cex = 1, cex=.8, bty = "n", horiz = TRUE)

legend("top", density=c(-9,66,33,10),
      c("OLS", "Ridge", "LASSO", "Elastic Net"), pt.cex = 1, cex=.8, horiz = TRUE)
```

## Optimal Betas for Each Covariate and Estimator



```
df <- data.frame("OLS" = c(ridge$beta.hat[1,-1]), "Ridge" = c(ridge$beta.hat[11,-1]),
  "LASSO" = c(ridge$beta.hat[12,-1]),
  "Elastic Net" = c(ridge$beta.hat[12,-1]))
rownames(df) <- c("X1", "X2", "X3", "X4", "X5", "X6", "X7",
  "X8", "X9", "X10")
df
```

	OLS	Ridge	LASSO	Elastic.Net
## X1	-0.28849736	-0.216668047	-0.26950020	-0.26950020
## X2	-0.17899475	-0.158348789	-0.18187158	-0.18187158
## X3	-0.59845312	-0.350397527	-0.51359879	-0.51359879
## X4	0.09709517	0.016200038	0.06645509	0.06645509
## X5	0.30617854	0.182865102	0.26673938	0.26673938
## X6	-0.06087212	-0.003025458	-0.04164557	-0.04164557
## X7	0.26074607	0.129539668	0.20725939	0.20725939
## X8	-0.02069877	0.148594421	0.05591133	0.05591133
## X9	0.50089470	0.373439180	0.46091800	0.46091800
## X10	0.64871828	0.455316910	0.59684761	0.59684761

The plots of the risk versus the shrinkage parameter show us that the risk is minimized for smaller values of the shrinkage parameter for the LASSO and Elastic Net regression estimators in comparison to the Ridge regression estimator. We examine the corresponding “optimal” estimators of the regression coefficients that we constructed as well as OLS with plots and a table. These visuals show us that the optimal estimators of the regression coefficients across all of the covariates are most similar for the Ridge, LASSO, and Elastic Net regression estimators and differ widely from the optimal OLS estimators of the regression coefficients across all of the covariates.

### d) Ridge regression: Bias, variance, and mean squared error of estimated regression coefficients.

Derive the bias, variance, and mean squared error of the ridge estimators of the regression coefficients. Be specific about assumptions and which variables you are conditioning on.

For the simulation model of a), provide and comment on graphical displays of the bias, variance, and MSE of the ridge estimators based on the learning set. For each coefficient, provide the value of the shrinkage parameter  $\lambda$  minimizing the MSE and the corresponding estimate.

### Solution:

According to equation (21) on the ‘Regularized Regression’ lecture slides,  $E[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] = (\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1} \mathbf{X}_n^\top E[\mathbf{Y}_n|\mathbf{X}_n] = (\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1} \mathbf{X}_n^\top \mathbf{X}_n \beta$ . The bias is as follows,  $\text{Bias}[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] = E[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] - \beta = (\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1} \mathbf{X}_n^\top \mathbf{X}_n \beta - \beta = ((\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1} \mathbf{X}_n^\top \mathbf{X}_n - \mathbf{I}) \beta$ .

According to equation (22) on the ‘Regularized Regression’ lecture slides, the covariance matrix of the ridge regression estimator,  $\text{Cov}[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] = \sigma^2 (\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1} \mathbf{X}_n^\top \mathbf{X}_n (\mathbf{X}_n^\top \mathbf{X}_n + \lambda \mathbf{I}_J)^{-1}$ .

Thus, for the parameter  $\beta = (\beta_j : j = 1, \dots, J) \in \mathbf{R}^J$ , a J-dimensional column vector of regression coefficients we have a J-dimensional vector of mean squared errors for each  $\beta_j$  is  $\text{MSE}[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] = \text{Var}[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n] + (\text{Bias}[\hat{\beta}_n^{\text{ridge}}|\mathbf{X}_n])^2$ .

According to the slides, and we can see here, that the ridge estimator is biased. As the shrinkage parameter increases, the bias tends to increase while variance tends to decrease. This is because we become more data-adaptive and less smooth as we increase the shrinkage parameter, highlighting the bias-variance trade-off of the ridge regression estimator. It should be noted that we assume the model in Equation (1) and the bias and covariance matrices of the ridge regression estimator are conditional on the design matrix of the learning set.

```
#####
# Setting up
#####

# Directly from Sandrine's 'Regularized Regression:Example' code

## Ridge regression: Bias, variance, and MSE
## N.B. Do not fit intercept.

myRidgePerf <- function(x,y,beta=0,sigma=1,scale=FALSE,lambda=0)
{
  n <- nrow(x)
  J <- ncol(x)
  df <- rep(NA,length(lambda))
  beta.hat <- bias <- var <- mse <- matrix(NA,length(lambda),J)
  cov <- array(NA,c(length(lambda),J,J))

  xx <- scale(x,center=TRUE,scale=scale)

  for(l in 1:length(lambda))
  {
    a <- solve(crossprod(xx)+lambda[l]*diag(J))
    df[l] <- sum(diag(xx%*%a%*%t(xx)))
    beta.hat[l,] <- a%*%crossprod(xx,y)
    bias[l,] <- a%*%t(xx)%*%x%*%beta - beta
    cov[l,,] <- sigma^2*a%*%crossprod(xx)%*%a
    var[l,] <- diag(cov[l,,])
    mse[l,] <- var[l,] + bias[l,]^2
  }

  res <- list(df=df,beta.hat=beta.hat,bias=bias,cov=cov,var=var,mse=mse)
  res
}

#####
```



```

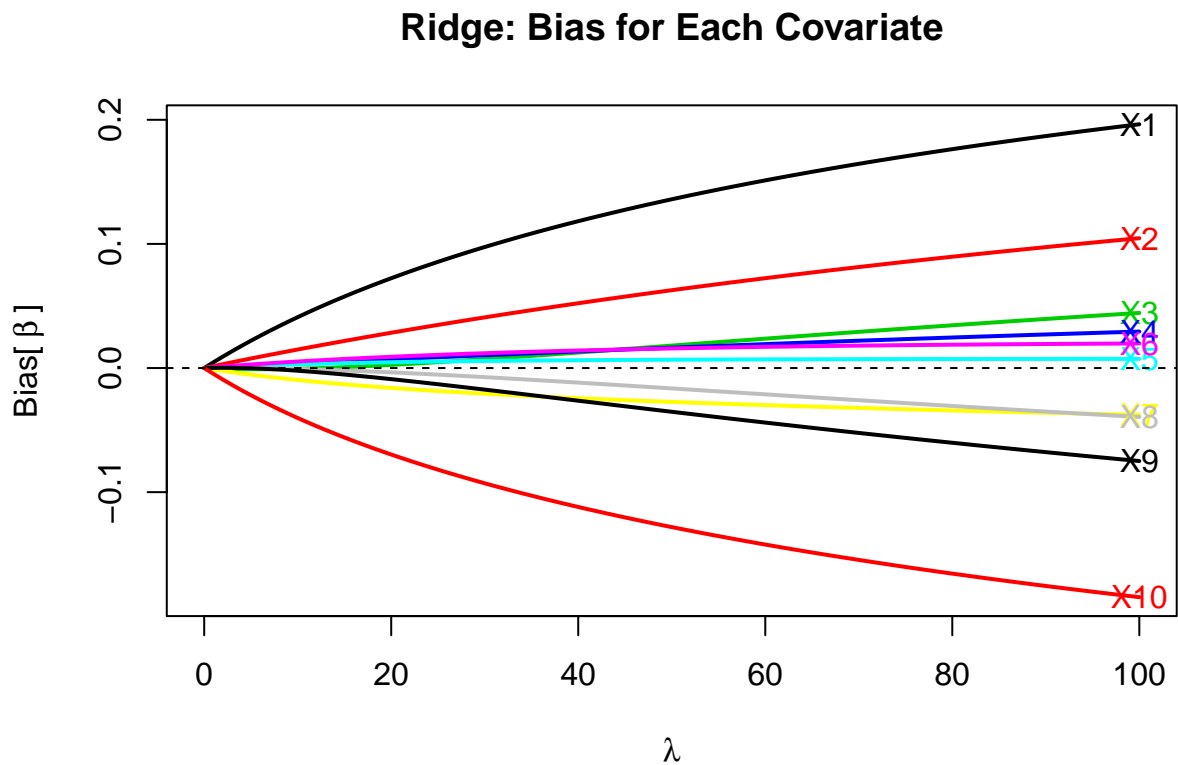
# Plots
#####

ridge_perf <- myRidgePerf(x=Lset_X, y=Lset_Y, beta=Beta, sigma=gamma, lambda=lambda)

# Bias

myPlotBeta(lambda, ridge_perf$bias, labels=colnames(Lset[,1:10]), right=TRUE,
            ylab=expression("Bias[" ~ hat(beta) ~ "]" ),
            main = "Ridge: Bias for Each Covariate")

```



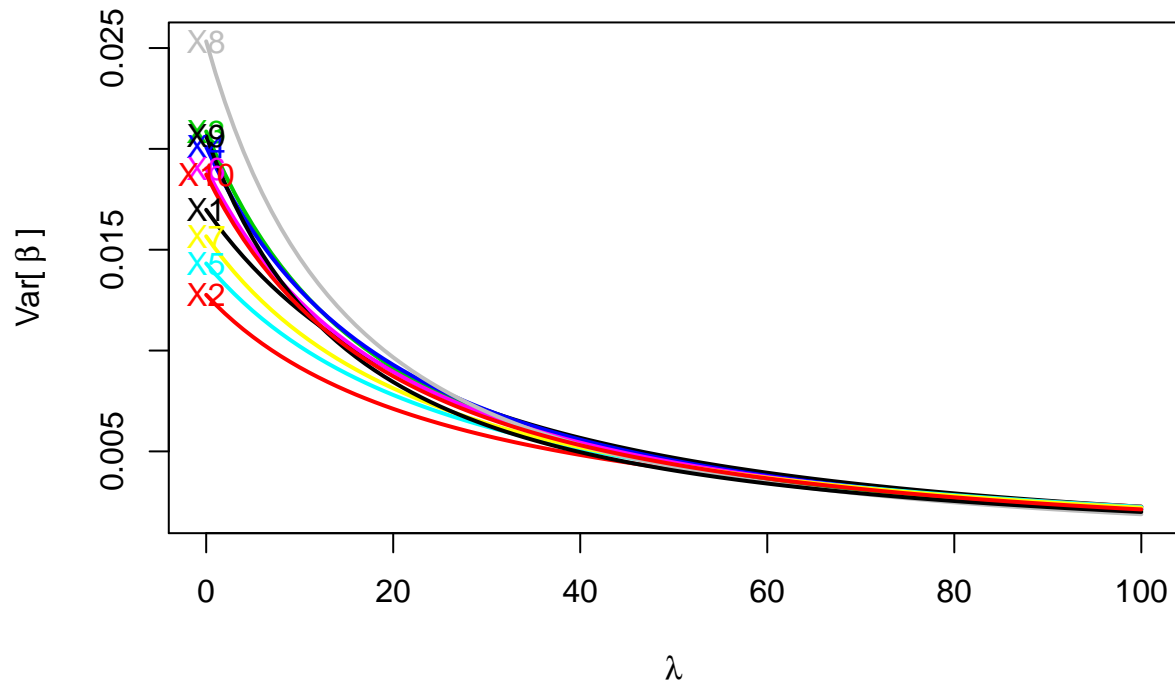
```

# Variance

myPlotBeta(lambda, ridge_perf$var, labels=colnames(Lset[,1:10]), right=FALSE,
            ylab=expression("Var[" ~ hat(beta) ~ "]" ),
            main = "Ridge: Variance for Each Covariate")

```

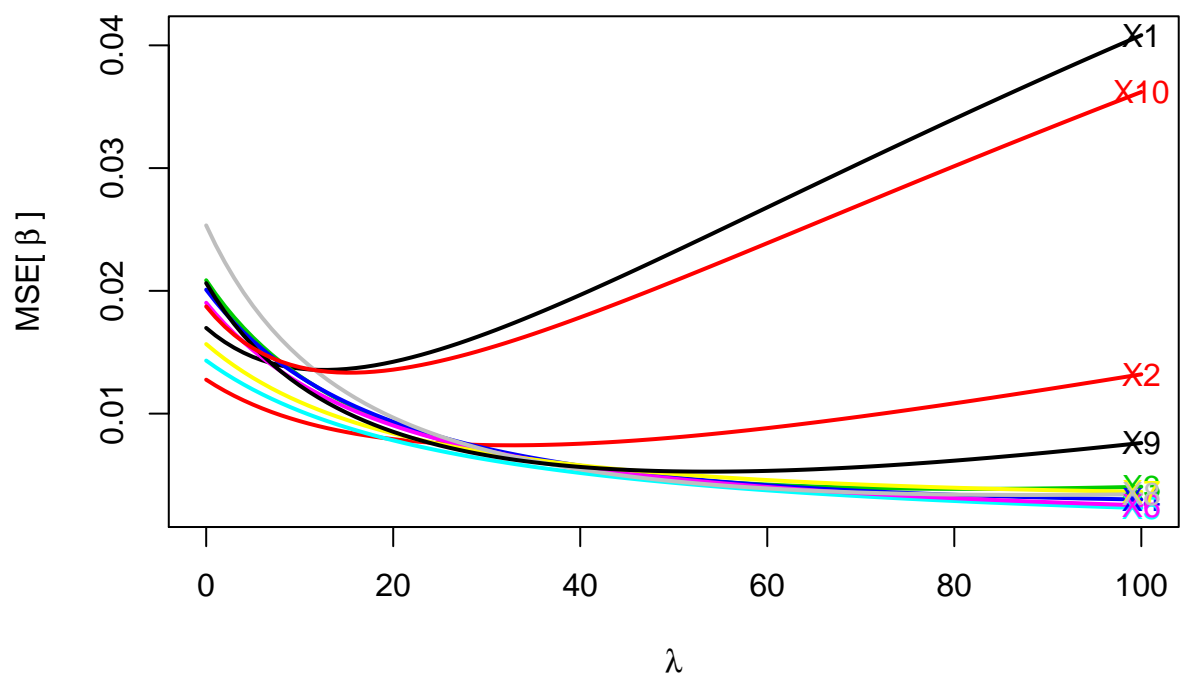
### Ridge: Variance for Each Covariate



# MSE

```
myPlotBeta(lambda, ridge_perf$mse, labels=colnames(Lset[,1:10]), right=TRUE,
  ylab=expression("MSE[" ~ hat(beta) ~ "]"),
  main = "Ridge: MSE for Each Covariate")
```

### Ridge: MSE for Each Covariate



```
#####
# Optimal lambda with beta
#####

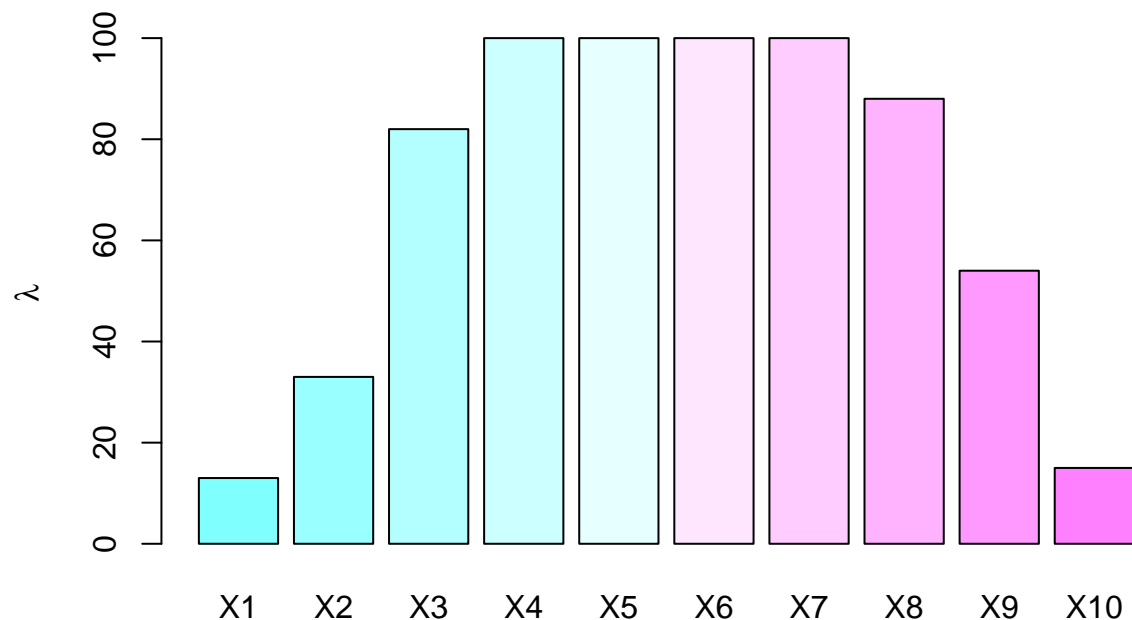
index <- apply(ridge_perf$mse, 2, which.min)
lambda_optimal <- lambda[index]

beta_hat <- rep(NA,J)
beta <- rep(NA,J)
error <- rep(NA,J)
for(j in 1:J){
  beta_hat[j] <- ridge_perf$beta.hat[index[j],j]
  beta[j] <- Beta[j]
  error[j] <- beta_hat[j] - beta[j]
}

df <- data.frame(lambda_optimal, beta_hat, beta, error)
rownames(df) <- c("X1","X2","X3","X4","X5",
                  "X6","X7","X8","X9","X10")

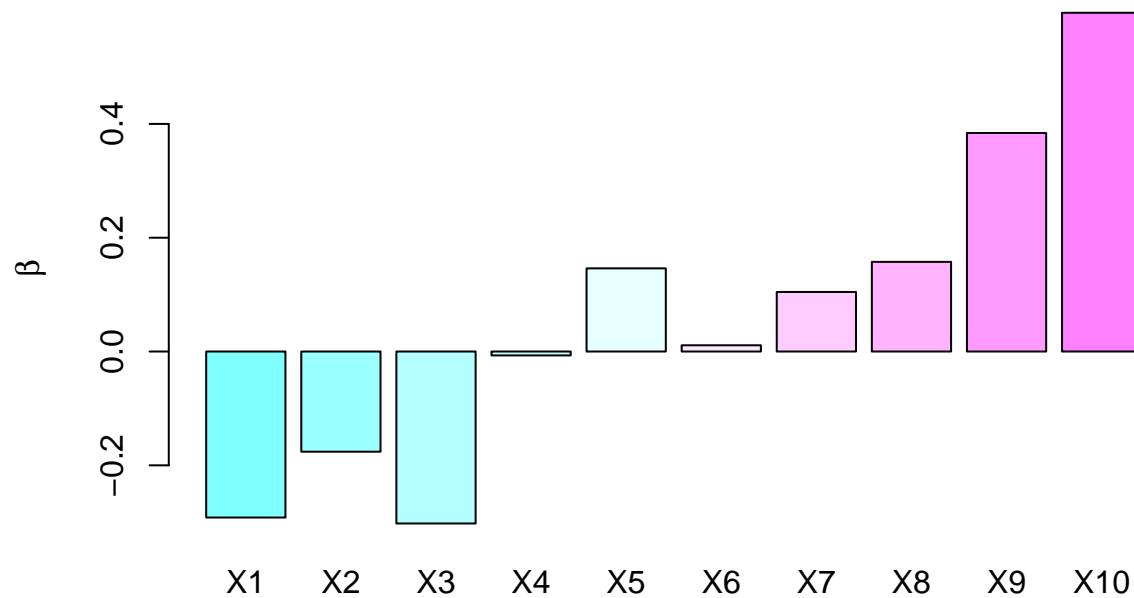
# plot of optimal lambda
barplot(df$lambda_optimal, main = expression("Ridge: Optimal shrinkage parameter for each covariate"), ylab =  $\lambda$ )
```

Ridge: Optimal shrinkage parameter for each covariate



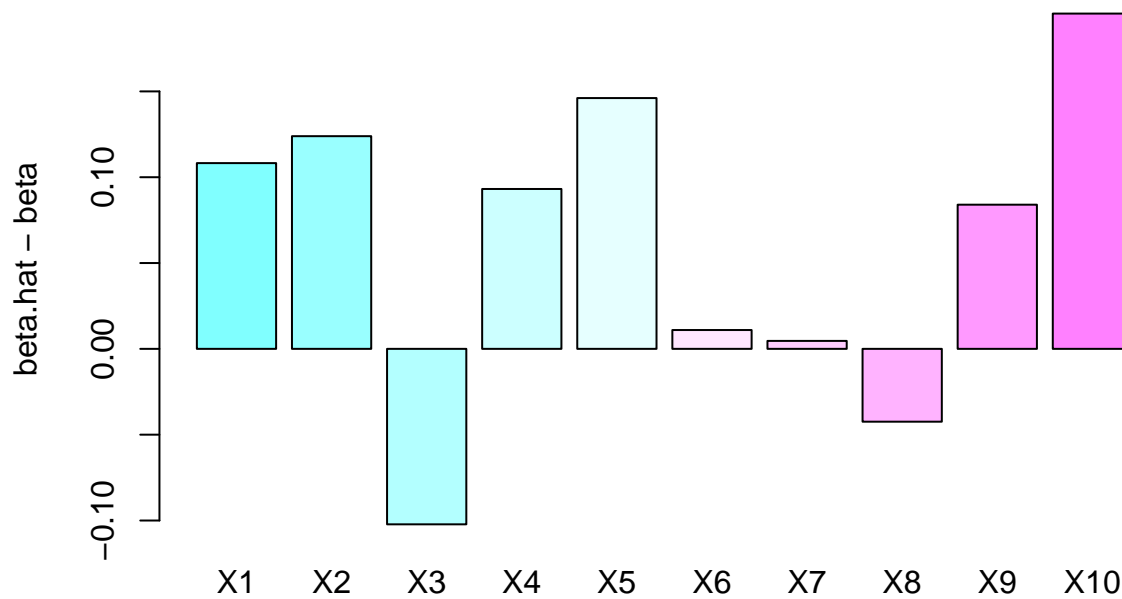
```
barplot(df$beta_hat, main = expression("Ridge: Corresponding regression coefficient for the optimal shrinkage parameter"), ylab = "beta_hat")
```

Ridge: Corresponding regression coefficient for the optimal shrinkage paran



```
par(mfrow=c(1,1))
barplot(df$error, main = expression("Ridge: Error for the Optimal shrinkage
parameter for each covariate"),
ylab = "beta.hat - beta", col=cm.colors(10), names.arg = rownames(df))
```

Ridge: Error for the Optimal shrinkage  
parameter for each covariate



```
df
##      lambda_optimal      beta_hat beta      error
## X1                13 -0.291788212 -0.4  0.108211788
```

## X2	33	-0.176091021	-0.3	0.123908979
## X3	82	-0.302227822	-0.2	-0.102227822
## X4	100	-0.006855764	-0.1	0.093144236
## X5	100	0.146114152	0.0	0.146114152
## X6	100	0.010975557	0.0	0.010975557
## X7	100	0.104634178	0.1	0.004634178
## X8	88	0.157565930	0.2	-0.042434070
## X9	54	0.383998544	0.3	0.083998544
## X10	15	0.595337160	0.4	0.195337160

The graphs nicely display the bias-variance trade-off mentioned in the beginning of this solution. We see that as we increase the shrinkage parameter the bias increases and the variance decreases. The MSE plot shows us that there surely exist optimal values of the shrinkage parameter; as we increase the shrinkage parameter the MSE decrease a bit across all covariates and then it increases if the shrinkage parameter increases too much. We find the values of the optimal shrinkage parameter (those that minimize the MSE) with the corresponding regression coefficient estimate and compare this to the true regression coefficients in the last plot. There is a noticeable amount of variability for the error (distance from the estimate to the truth) across the covariates.

**e) LASSO regression: Bias, variance, and mean squared error of estimated regression coefficients.**

For the LASSO, there are no closed-form expressions for the bias, variance, and mean squared error of the estimators of the regression coefficients.

Describe how one can estimate these quantities using the simulation model of a). In particular, provide and comment on graphical displays of the bias, variance, and MSE of the LASSO estimators based on the learning set. For each coefficient, provide the value of the shrinkage parameter  $\lambda$  minimizing the MSE and the corresponding estimate. Again, be specific about assumptions and which variables you are conditioning on.