# Chapter 10: Shrinkage Methods

## Outline

- Principal Components Analysis

- Ridge regression

- Lasso

- Partial Least Squares
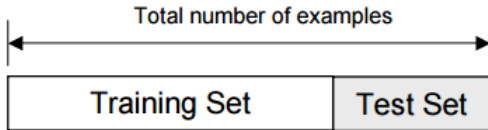
## Training and Test Data Set

Motivation

- Model Selections and Shrinkage Methods provide good models

- Hard to choose only one optimal model.

- Choose an optimal model based on its performance.
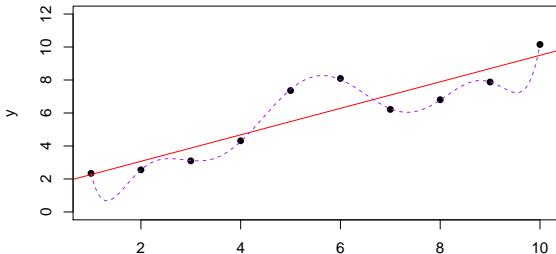
# Training and Test Data Set

Method

- Split dataset into two groups .
- Training Set: Used to train the model.
- Test Set: Used to measure the performance of the trained model.
- Validation Set: Used to train tuning parameters. (Extra).

## Training and Test Data Set

Why Test Set?

- Overfitting Issue: It refers to a model that models the training data too well.

## Training and Test Data Set

Why Test Set?

- Overfitting Issue

- Complex models tend to have a good performance in training data.

- A good model in training dataset may not be a good model in new dataset.

**Training and Test Data Set**

How to determine Test Data Set?

- Absolute Random selection: Choose $10 \sim 20\%$ of data set.

- Choose $10 \sim 20\%$ of data set with same proportion of success
  in both Training and Test data set. (To prevent the test set
  only contains success or same value of predictors. )

## Training and Test Data Set

Weakness

- Cannot be performed when sample size is small
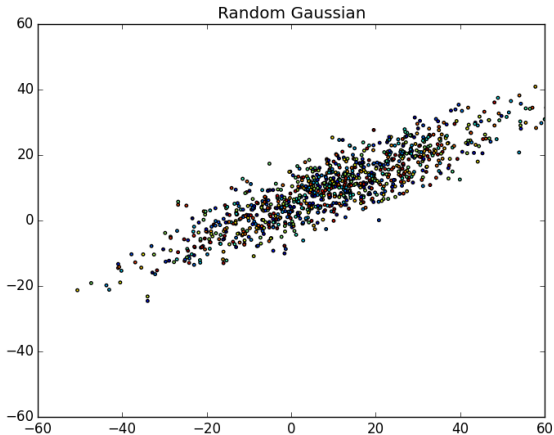
- Cross-Validation (Later)

### Principal Components Analysis (PCA)

- Special transformation on predictors

- Useful for high dimensional data

- Solve collinearity issue

## Principal Components Analysis (PCA)

- Find the $u_1$ such that $\text{var}(u_1^T X)$ is maximized subject to $u_1^T u_1 = 1$.
- Find the $u_2$ such that $\text{var}(u_2^T X)$ is maximized subject to $u_2^T u_1 = 0$ and $u_2^T u_2 = 1$.
- Keep finding directions of greatest variation orthogonal to those directions we have already found.

# Principal Components Analysis (PCA)



Random Gaussian

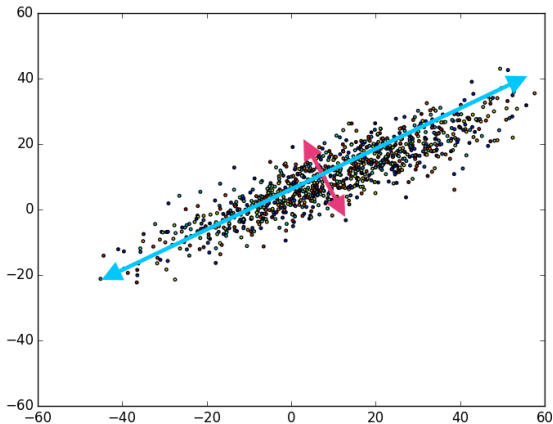# Principal Components Analysis (PCA)



Figure: PCA: $u_1$ and $u_2$

## Simulation Study

```
> prcomp(x)

Rotation:
     PC1    PC2
x1 -0.892 -0.451
x2 -0.451  0.892
```

$$Z1 = -0.892X_1 - 0.451X_2$$

$$Z2 = -0.451X_1 + 0.892X_2$$

## Simulation Study

```
> plot(x, pch = 16, xlim = c(-4, 4), ylim = c(-3,3))
> abline(a = 0, b = 0.451/-0.892, col = "red")
> abline(a = 0, b = 0.451/0.892, col = "blue")
```

## Simulation Study

```
> prX = prcomp(x)
> summary(prX)

Importance of components:
                          PC1     PC2
Standard deviation       1.210  0.1918
Proportion of Variance   0.976  0.0245
Cumulative Proportion    0.976  1.0000
```

## Simulation Study

```
> round(prX$rot[,1],3)
x1     x2
-0.892 -0.451
```

$$Z_1 = -0.892X_1 - 0.451X_2$$

$Z_1$ explains 97.6% of the both $X_1$ and $X_2$

## Simulation Study

```
> lm0 = lm(Y ~ X1 + X2)
> summary(lm0)

Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.1164    0.1091  -1.067    0.289
x1           0.8597    0.2020   4.255 4.82e-05 ***
x2           1.2688    0.2219   5.717 1.20e-07 ***
---
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1   1

Residual standard error: 1.086 on 97 degrees of freedom
Multiple R-squared:  0.7988,Adjusted R-squared:  0.7946
F-statistic: 192.5 on 2 and 97 DF,  p-value: < 2.2e-16
```

# Simulation Study

```
> lm1 = lm(Y ~ Z)
> summary(lm1)

Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept)  -0.1092     0.1094  -0.999      0.32
z            -1.4875     0.0762 -19.520    <2e-16 ***
---
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1   1

Residual standard error: 1.09 on 98 degrees of freedom
Multiple R-squared:  0.7954,Adjusted R-squared:  0.7933
F-statistic:    381 on 1 and 98 DF,  p-value: < 2.2e-16
```

# Fat Example

- Response: fat
- Predictors: neck, chest, abdom, hip, thigh, knee, ankle, biceps, forearm, wrist

```
> cfat = fat[,9:18]
> prfat = prcomp(cfat)
> dim(prfat$rot)
[1] 10 10
> dim(prfat$x)
[1] 252  10
```

# Fat Example

```
> summary(prfat)

Importance of components:
                        PC1    PC2    PC3    PC4     PC5      PC6      PC7      PC8      PC9     PC10
Standard deviation    15.990 4.0658 2.9660 2.0004 1.69408 1.49881 1.30322 1.25478 1.10955 0.52737
Proportion of Variance 0.867 0.0561 0.0298 0.0136 0.00973 0.00762 0.00576 0.00534 0.00417 0.00094
Cumulative Proportion  0.867 0.9230 0.9529 0.9664 0.97617 0.98378 0.98954 0.99488 0.99906 1.00000
```

## Fat Example

```
> round(prfat$rot[,1],2)
 neck   chest   abdom    hip   thigh    knee
 0.12    0.50    0.66   0.42    0.28    0.12

ankle  biceps forearm  wrist
 0.06    0.15    0.07   0.04
```

# Fat Example

```
> prfatc = prcomp(cfat, scale = TRUE)
> summary(prfatc)
Importance of components:
                          PC1    PC2    PC3    PC4    PC5    PC6    PC7    PC8     PC9    PC10
Standard deviation      2.650 0.8530 0.8191 0.7011 0.5471 0.5283 0.4520 0.4054 0.27827 0.2530
Proportion of Variance  0.702 0.0728 0.0671 0.0492 0.0299 0.0279 0.0204 0.0164 0.00774 0.0064
Cumulative Proportion   0.702 0.7749 0.8420 0.8911 0.9211 0.9490 0.9694 0.9859 0.99360 1.0000
```

## Fat Example

```
> round(prfatc$rot[,1],3)

 neck   chest  abdom    hip  thigh   knee  ankle
0.327   0.339  0.334  0.348  0.333  0.329  0.247

 biceps forearm  wrist
 0.322   0.270  0.299
```

## Fat Example

```
> round(prfatc$rot[,2],3)

 neck   chest  abdom     hip  thigh   knee  ankle
-0.003 -0.273 -0.398  -0.255 -0.191  0.022  0.625

 biceps forearm wrist
 0.022   0.363  0.377
```

## Fat Example

```
> lmoda = lm(fat$brozek ~., data = cfat)
> summary(lmoda)
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept)  7.228749   6.214309   1.163  0.24588
neck        -0.581947   0.208580  -2.790  0.00569 **
chest       -0.090847   0.085430  -1.063  0.28866
abdom        0.960229   0.071582  13.414  < 2e-16 ***
hip         -0.391355   0.112686  -3.473  0.00061 ***
thigh        0.133708   0.124922   1.070  0.28554
knee        -0.094055   0.212394  -0.443  0.65828
ankle        0.004222   0.203175   0.021  0.98344
biceps       0.111196   0.159118   0.699  0.48533
forearm      0.344536   0.185511   1.857  0.06450 .
wrist       -1.353472   0.471410  -2.871  0.00445 **

Residual standard error: 4.071 on 241 degrees of freedom
Multiple R-squared:  0.7351,Adjusted R-squared:  0.7241
F-statistic: 66.87 on 10 and 241 DF,  p-value: < 2.2e-16
```

# Fat Example

```
> lmodpcr = lm(fat$brozek ~ prfatc$x[,1:2])
> summary(lmodpcr)
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept)          18.9385     0.3291  57.542   <2e-16 ***
prfatc$x[, 1:2]PC1    1.8420     0.1245  14.800   <2e-16 ***
prfatc$x[, 1:2]PC2   -3.5505     0.3866  -9.184   <2e-16 ***
---
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1   1

Residual standard error: 5.225 on 249 degrees of freedom
Multiple R-squared:  0.5492,Adjusted R-squared:  0.5456
F-statistic: 151.7 on 2 and 249 DF,  p-value: < 2.2e-16
```

# Fat Example

```
> lmodpcr2 = lm(fat$brozek ~ prfatc$x[,1:10])
> summary(lmodpcr2)

Estimate Std. Error t value Pr(>|t|)
(Intercept)             18.93849    0.25647  73.843  < 2e-16 ***
prfatc$x[, 1:10]PC1      1.84198    0.09698  18.993  < 2e-16 ***
prfatc$x[, 1:10]PC2     -3.55053    0.30126 -11.785  < 2e-16 ***
prfatc$x[, 1:10]PC3      0.25669    0.31374   0.818 0.414067
prfatc$x[, 1:10]PC4      0.54094    0.36652   1.476 0.141273
prfatc$x[, 1:10]PC5      3.72632    0.46973   7.933 8.03e-14 ***
prfatc$x[, 1:10]PC6     -1.48784    0.48642  -3.059 0.002474 **
prfatc$x[, 1:10]PC7      1.94878    0.56859   3.427 0.000716 ***
prfatc$x[, 1:10]PC8     -0.12247    0.63390  -0.193 0.846967
prfatc$x[, 1:10]PC9     -1.71366    0.92351  -1.856 0.064731 .
prfatc$x[, 1:10]PC10    -9.01059    1.01566  -8.872  < 2e-16 ***

Residual standard error: 4.071 on 241 degrees of freedom
Multiple R-squared:  0.7351,Adjusted R-squared:  0.7241
F-statistic: 66.87 on 10 and 241 DF,  p-value: < 2.2e-16
```

## Fat Example

Benefits of PCA

- Orthogonal Predictors

- No collinearity Issue

- Sometimes easy to Interpret

## Choice of the number of variables

How to choose the optimal number of variables.

- Interpretability: It is important to examine the interpretability of the components and make sure that those providing a interpretable result are retained.
- Total variance
- Eigenvalues (Skip)

## Remarks on PCA

- Interpretation may be easy or difficult

- Sufficiently reduce the number of predictors

- Difficult to decide the number of predictors

### Choice of the number of variables

How to choose the optimal number of variables.

- Training and Test Set

- Root meas square of error (RMSE)

**Root meas square of error (RMSE)**

Definition:
$$\text{RMSE} = \sqrt{\frac{1}{n}\sum_i^n (\hat{y}_i - y_i)^2}$$

Motivation:

- Bias: $y - \mathbb{E}(\hat{y})$
- Variance: $\text{Var}(\hat{y})$

# Bias

## Variance

$y_1 = 10$

- Small variance case: 95% confidence interval for $y_1$:
  (9.99, 10.01)
  $\rightarrow$ 95% sure that $y_1$ is in (9.99, 10.01)

- :Large variance case: 95% confidence interval for $y_1$: $(-5, 25)$
  $\rightarrow$ 95% sure that $y_1$ is in $(-5, 25)$

## Root meas square of error (RMSE)

- Mean Square Error $=$ Bias$(\hat{y})^2 +$ Variance$(\hat{y})$
- RMSE is a good criterion to choose an optimal model

**Simulation Study: Choice of the number of variables**

How to choose the optimal number of variables.

- 4 predictors: $X_1, X_2, X_3, X_4$

- Training set: 40 observations

- Test set: 10 observations

## Simulation Study

```
> ran = sample(1:50, replace = F)[1:40]
> train = data[ran,]
> test = data[setdiff(1:50,ran),]
> prx = prcomp(x[ran,])
> summary(prx)

Importance of components:
                          PC1     PC2     PC3      PC4
Standard deviation      1.3194  0.5817  0.4301  0.21264
Proportion of Variance  0.7538  0.1465  0.0801  0.01958
Cumulative Proportion   0.7538  0.9003  0.9804  1.00000
```

## Simulation Study

```
> lmodpcr2 = lm(y ~ prx$x[,1:2], train)
> z1 = prx$rotation[,1] %*% t( test[,1:4] )
> z2 = prx$rotation[,2] %*% t( test[,1:4] )
> z = rbind(z1, z2)
> ypred = coef(lmodpcr2) %*% rbind(1, z)
> sqrt( mean( (test$y - ypred)^2) )
[1] 1.294111
```

## Simulation Study

```
> lmodpcr3 = lm(y ~ prx$x[,1:3], train)
> z1 = prx$rotation[,1] %*% t( test[,1:4] )
> z2 = prx$rotation[,2] %*% t( test[,1:4] )
> z3 = prx$rotation[,3] %*% t( test[,1:4] )
> z = rbind(z1, z2, z3)
> ypred = coef(lmodpcr3) %*% rbind(1, z)
> sqrt( mean( (test$y - ypred)^2) )
[1] 1.29453
```

## Ridge Regression

Penalizing the square of the coefficients

$$\min_{\beta} \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2$$

Assumption:

- Regression Coefficients should not be very large (after standardization).
- A large number of predictors should be considered.
- High collinearity exists.

# Ridge Regression

Penalizing the square of the coefficients

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2$$

- The coefficients $\hat{\boldsymbol{\beta}}^{\mathrm{ridge}}$ are shrunken towards zero.

- $\lambda \geq 0$ is a tuning parameter.

- $\lambda$ controls the amount of shrinkage.

- What happens if $\lambda \to 0$?

- What happens if $\lambda \to \infty$?

**Equivalent Formulation**

$$\min_{\boldsymbol{\beta}} \quad \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2$$

$$\text{subject to} \quad \sum_{j=1}^{p} \beta_j^2 \le s$$

- Explicitly constraint the size of the coefficients.

## Equivalent Formulation

## Ridge Regression

When there are many highly correlated variables

- $\hat{\beta}^{\mathrm{ols}}$ may have a large coefficient on one variable and a similarly large negative coefficient on its correlated variable (Unstable).

- In ridge regression, the size constraint tries to avoid this phenomenon.

Often standardize the predictors first.

## Solution for Ridge Regression

- The solution is

$$\hat{\boldsymbol{\beta}}^{\mathrm{ridge}} = (\boldsymbol{X}^{\mathsf{T}}\boldsymbol{X} + \lambda\boldsymbol{I})^{-1}\boldsymbol{X}^{\mathsf{T}}\boldsymbol{y}$$

- $\hat{\boldsymbol{\beta}}$ is linear in $\boldsymbol{y}$.
- $\hat{\boldsymbol{\beta}}$ is biased.

## Comparison to LSE

$$\hat{\boldsymbol{\beta}}^{\mathrm{ridge}} = (\boldsymbol{X}^{\mathsf{T}}\boldsymbol{X} + \lambda \boldsymbol{I})^{-1}\boldsymbol{X}^{\mathsf{T}}\boldsymbol{y}$$

- Even if $\boldsymbol{X}$ is not full-rank, $(\boldsymbol{X}^{\mathsf{T}}\boldsymbol{X} + \lambda \boldsymbol{I})$ is invertible, thus solve exact collinearity issue.
- $\hat{\boldsymbol{\beta}}^{\mathrm{ridge}}$ has smaller variance than the OLS, thus may have smaller mean square error (MSE).

## Shrinkage in Ridge

Suppose orthonormal design ($\boldsymbol{X}^\mathsf{T}\boldsymbol{X} = \boldsymbol{I}$). Then $\hat{\boldsymbol{\beta}}^{\mathrm{ols}} = \boldsymbol{X}^\mathsf{T}\boldsymbol{y}$, and

$$(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})^\mathsf{T}(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}) = \mathrm{constant} + \sum_{j=1}^{p}(\beta_j - \hat{\beta}_j^{\mathrm{ols}})^2.$$

Then ridge regression minimizes

$$\sum_{j=1}^{p}(\beta_j - \hat{\beta}_j^{\mathrm{ols}})^2 + \lambda\sum_{j=1}^{p}\beta_j^2.$$

Equivalent to the component-wise minimization

$$\min_{\beta_j}(\beta_j - \hat{\beta}_j^{\mathrm{ols}})^2 + \lambda\beta_j^2 \Longrightarrow \hat{\beta}_j^{\mathrm{ridge}} = \frac{1}{1+\lambda}\hat{\beta}_j^{\mathrm{ols}}.$$

## Shrinkage in Ridge

- Shrink the estimate towards zero by a positive constant less than 1
- $\mathrm{Var}(\hat{\beta}_j^{\mathrm{ridge}}) = \frac{1}{(1+\lambda)^2}\mathrm{Var}(\hat{\beta}_j^{\mathrm{ols}})$.
- $\lambda \uparrow$, shrinkage $\uparrow$, bias $\uparrow$, variance $\downarrow$
- $\lambda \downarrow$, shrinkage $\downarrow$, bias $\downarrow$, variance $\uparrow$.

## Simulation Study: Almost Independent Predictors

3 predictors: $X_1, X_2, X_3$

```
> cor(data[,1:3])
       x1      x2      x3
x1  1.0000  0.0505 -0.142
x2  0.0505  1.0000 -0.143
x3 -0.1425 -0.1428  1.000

> lmod = lm(y ~ x1 + x2 + x3, data)
> lmod$coefficients
(Intercept)    x1        x2        x3
-0.724      1.023     1.299     1.817
```

## Simulation Study

```
> require(MASS)
> #lambda = 0
> lmrid = lm.ridge(y~x1 + x2 + x3, data, lmabda = 0)
> lmrid
(Intercept)     x1          x2          x3
-0.724        1.023       1.299       1.817
```

## Simulation Study

```
> #lambda = 0.1
> lmrid2 = lm.ridge(y~x1 + x2 + x3, data, lambda = 0.1)
> lmrid2
(Intercept)     x1          x2          x3
-0.669        1.021       1.297       1.809
>
> #lambda = 10
> lmrid3 = lm.ridge(y~x1 + x2 + x3, data, lambda = 10)
> lmrid3
(Intercept)     x1          x2          x3
3.819         0.848       1.116       1.187
```

## Simulation Study: Correlated Predictors

```
> lmod$coefficients
(Intercept)          x1           x2           x3
-0.676            1.132       -0.275        1.779

> cor(data[,1:3])
x1    x2    x3
x1 1.000 0.996 0.999
x2 0.996 1.000 0.998
x3 0.999 0.998 1.000
```

## Simulation Study: Correlated Predictors

```
> require(MASS)
> #lambda = 0
> lmrid = lm.ridge(y~x1 + x2 + x3, data, lmabda = 0)
> lmrid
(Intercept)         x1          x2          x3
-0.676            1.132       -0.275      1.779
>
> #lambda = 0.1
> lmrid2 = lm.ridge(y~x1 + x2 + x3, data, lambda = 0.1)
> lmrid2
(Intercept)         x1          x2          x3
-0.720            1.405       0.327       1.195
```

## Simulation Study: Correlated Predictors

```
> #lambda = 1
> lmrid3 = lm.ridge(y~x1 + x2 + x3, data, lambda = 1)
> lmrid3
(Intercept)         x1          x2          x3
-0.166            1.260       0.968       0.848
>
> #lambda = 10
> lmrid3 = lm.ridge(y~x1 + x2 + x3, data, lambda = 10)
> lmrid3
(Intercept)         x1          x2          x3
4.85             1.12        1.07        0.74
```

## Simulation Study: Comparison to LSE

```
# Generate Training/Test sets
> ran = sample(1:50, replace = F)[1:40]
> train = data[ran,]
> test = data[setdiff(1:50,ran),]

> lmod = lm(y ~ x1 + x2 + x3, train)
> lmrid = lm.ridge(y~x1 + x2 + x3, train, lambda = 0.1)
> lmrid2 = lm.ridge(y~x1 + x2 + x3, train, lambda = 1)
```

# Simulation Study: Comparison to LSE

```
# RMSE
> sqrt(mean((test$y - predict(lmod,test))^2))
[1] 5.4921
>
> ypred = cbind(1, as.matrix(test[,-4])) %*% coef(lmrid)
> sqrt(mean((test$y - ypred)^2))
[1] 5.1544
>
> ypred = cbind(1, as.matrix(test[,-4])) %*% coef(lmrid2)
> sqrt(mean((test$y - ypred)^2))
[1] 5.5591
```

## Simulation Study: Normalization

```
> data.scale = scale(data[,1:3])
> data.scale = data.frame(data.scale, y = data$y)
>
> train = data.scale[ran,]
> test = data.scale[setdiff(1:50,ran),]
>
> lmod = lm(y ~ x1 + x2 + x3, train)
> lmrid = lm.ridge(y~x1 + x2 + x3, train, lambda = 0.1)
> lmrid2 = lm.ridge(y~x1 + x2 + x3, train, lambda = 1)
```

## Simulation Study: Normalization

```
> sqrt(mean((test$y - predict(lmod,test))^2))
[1] 4.8518
>
> ypred = cbind(1, as.matrix(test[,-4])) %*% coef(lmrid)
> sqrt(mean((test$y - ypred)^2))
[1] 4.8484
>
> ypred = cbind(1, as.matrix(test[,-4])) %*% coef(lmrid2)
> sqrt(mean((test$y - ypred)^2))
[1] 4.8587
```

## Choice of Tuning Parameter $\lambda$

Determination of the tuning parameter $\lambda$

- Generalized Cross-Validation (GCV): Almost same concept as Training and Test Set.

- GCV
$$V(\lambda) = \frac{\frac{1}{n}\|(I - A(\lambda))y\|^2}{(\frac{1}{n}tr(I - A(\lambda)))^2},$$

  where $A(\lambda) = X(X^T X + n\lambda I)^{-1} - X^T$.

- R provides a good tuning parameter $\lambda$ (Not always optimal).

## Simulation Study

```
> lmrid = lm.ridge(y~x1 + x2 + x3, train, lambda = seq(0,1, len
> lmrid$GCV
0.00    0.25    0.50    0.75    1.00
0.68954 0.66428 0.66800 0.67131 0.67449
>
> which.min(lmrid$GCV)
0.25
```

## Simulation Study

```
> lmrid = lm.ridge(y~x1 + x2 + x3, train,
                        lambda = seq(0,1, len = 100))
> which.min(lmrid$GCV)
0.161616

> lmrid_GCV = lm.ridge(y~x1 + x2 + x3, train, lambda = 0.161616)
> ypred = cbind(1, as.matrix(test[,-4])) %*% coef(lmrid_GCV)
> sqrt(mean((test$y - ypred)^2))
[1] 4.8405
```

## Simulation Study: Another Training/Test Set

```
> ran = sample(1:50, replace = F)[1:40]
> train = data.scale[ran,]
> test = data.scale[setdiff(1:50,ran),]
>
> lmrid = lm.ridge(y~x1 + x2 + x3, train,
                           lambda = seq(0,1, len = 100))
> which.min(lmrid$GCV)
0.080808
```

## Simulation Study: Another Training/Test Set

```
> ypred = cbind(1, as.matrix(test[,-4])) %*% coef(lmrid_GCV)
> lmrid_GCV = lm.ridge(y~x1 + x2 + x3, train, lambda = 0.0808)
>
> ypred = cbind(1, as.matrix(test[,-4])) %*% coef(lmrid_GCV)
> sqrt(mean((test$y - ypred)^2))
[1] 4.486
>
> lmod = lm(y ~ x1 + x2 + x3, train)
> sqrt(mean((test$y - predict(lmod,test))^2))
[1] 4.7677
```

## Simulation Study: Comparison to PCA

```
> X = train[,1:3]
> prx = prcomp(X)
> summary(prx)
Importance of components:
PC1     PC2     PC3
Standard deviation      1.755 0.06021 0.03225
Proportion of Variance 0.998 0.00117 0.00034
Cumulative Proportion  0.998 0.99966 1.00000
```

## Simulation Study: Comparison to PCA

```
> z = 0.57778 * X[,1] + 0.57848 * X[,2] + 0.57579 *X[,3]
> lmodpcr = lm(train$y ~ z)
> summary(lmodpcr)

Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept)   88.530      0.852    103.9    <2e-16 ***
z             29.304      0.491     59.6    <2e-16 ***
---
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1   1

Residual standard error: 5.39 on 38 degrees of freedom
Multiple R-squared:  0.989,Adjusted R-squared:  0.989
F-statistic: 3.56e+03 on 1 and 38 DF,  p-value: <2e-16
```

## Simulation Study: Comparison to PCA

```
> ypred = cbind(1, as.matrix(0.57778 * test[,1]
  + 0.57848 * test[,2] + 0.57579 *test[,3]))
  %*% coef(lmodpcr)
> sqrt(mean((test$y - ypred)^2))
[1] 4.1143
```
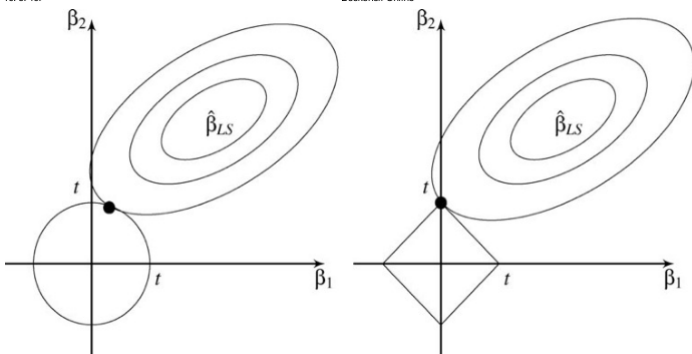
# LASSO

Least absolute shrinkage and selection operator (Chen, Donoho and Saunders 1996; Tibshirani 1996)

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$

- Shrinkage
- Sparsity: some fitted coefficients are exactly zero

Continuous variable selection

# Equivalent Formulation

**Equivalent Formulation**

$$\min_{\boldsymbol{\beta}} \quad \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2$$

$$\text{subject to} \quad \sum_{j=1}^{p} |\beta_j| \leq s$$

## Soft Thresholding

When $\boldsymbol{X}$ is orthonormal, we can minimize over $\boldsymbol{\beta}$ componentwise
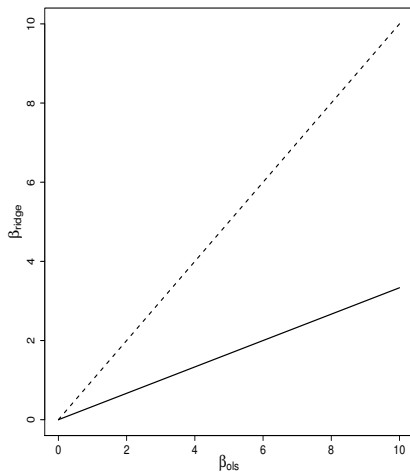
$$\min_{\beta_j} \; (\beta_j - \hat{\beta}_j^{\mathrm{ols}})^2 + \lambda|\beta_j|.$$

The solution is

$$
\begin{aligned}
\hat{\beta}_j^{\text{lasso}} &= \begin{cases} \hat{\beta}_j^{\text{ols}} - \frac{\lambda}{2} & \text{if } \hat{\beta}_j^{\text{ols}} > \frac{\lambda}{2} \\ 0 & \text{if } |\hat{\beta}_j^{\text{ols}}| \leq \frac{\lambda}{2} \\ \hat{\beta}_j^{\text{ols}} + \frac{\lambda}{2} & \text{if } \hat{\beta}_j^{\text{ols}} < -\frac{\lambda}{2} \end{cases} \\
&= \text{sign}(\hat{\beta}_j^{\text{ols}}) \cdot \left( |\hat{\beta}_j^{\text{ols}}| - \frac{\lambda}{2} \right)_+
\end{aligned}
$$

- Lasso shrinks large coefficients by a constant.
- Lasso truncates small coefficients to zero.

# Ridge vs Lasso
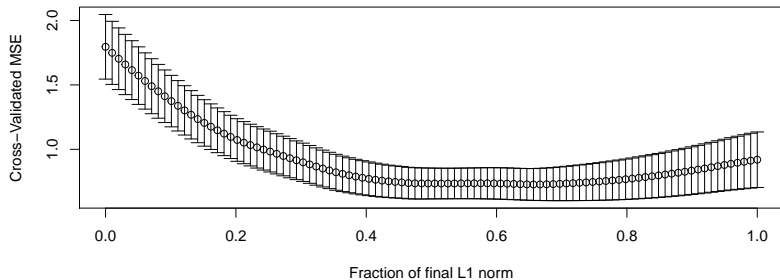
# Example: Life Expectancy

```
> plot(lmod)
> require(lars)
> data(state)
> statedata = data.frame(state.x77, row.names = state.abb)
> colnames(statedata)
[1] "Population" "Income"     "Illiteracy" "Life.Exp"   "Murder"     "HS.Grad"
[7] "Frost"      "Area"

> lmod = lars(as.matrix(statedata[,-4]), statedata$Life)
> coef(lmod)
Population   Income Illiteracy Murder HS.Grad    Frost     Area
[1,]  0.00e+00  0.00e+00    0.0000  0.000  0.0000  0.00000  0.00e+00
[2,]  0.00e+00  0.00e+00    0.0000 -0.141  0.0000  0.00000  0.00e+00
[3,]  0.00e+00  0.00e+00    0.0000 -0.203  0.0282  0.00000  0.00e+00
[4,]  1.28e-05  0.00e+00    0.0000 -0.216  0.0308  0.00000  0.00e+00
[5,]  4.90e-05  0.00e+00    0.0000 -0.298  0.0461 -0.00576  0.00e+00
[6,]  4.90e-05 -5.22e-08    0.0000 -0.298  0.0461 -0.00576  0.00e+00
[7,]  4.97e-05 -8.19e-06    0.0000 -0.298  0.0467 -0.00581 -7.79e-09
[8,]  5.18e-05 -2.18e-05    0.0338 -0.301  0.0489 -0.00574 -7.38e-08
```

## Example: Life Expectancy

```
> cvlmod = cv.lars(as.matrix(statedata[,-4]), statedata$Life.Exp)
> which.min( cvlmod$cv )
[1] 66
> cvlmod$index[66]
[1] 0.657
```

## Example: Life Expectancy

```
> predict(lmod, s=0.657, type="coef", mode="fraction")$coef

Population    Income   Illiteracy     Murder
2.35e-05    0.00e+00   0.00e+00     -2.40e-01

HS.Grad     Frost       Area
3.53e-02   -1.70e-03   0.00e+00
```

## Example: Life Expectancy

```
> g = lm(Life.Exp ~ Population + Murder + HS.Grad +Frost, stated
> coef(g)
(Intercept)  Population      Murder      HS.Grad        Frost
   7.10e+01    5.01e-05    -3.00e-01     4.66e-02    -5.94e-03
```

# Example: Life Expectancy

```
# Ridge
> require(MASS)
> g = lm.ridge(Life.Exp ~., statedata, lambda = seq(0, 4, len = 50))
> which.min(g$GCV)
2.7755
> g = lm.ridge(Life.Exp ~., statedata, lambda = 2.7755)
> g

          Population   Income   Illiteracy    Murder
7.08e+01  4.13e-05     2.32e-05  -7.89e-02     -2.64e-01

 HS.Grad    Frost       Area
 4.60e-02   -5.15e-03   -3.89e-07
```

# Example: Life Expectancy

```
# AIC
> g = lm(Life.Exp ~., statedata)
> step(g, direction = "backward", k = 2)

Step:  AIC=-28.2
Life.Exp ~ Population + Murder + HS.Grad + Frost

Df Sum of Sq RSS   AIC
<none>                  23.3 -28.2
- Population  1   2.1 25.4 -25.9
- Frost       1   3.1 26.4 -23.9
- HS.Grad     1   5.1 28.4 -20.2
- Murder      1  34.8 58.1  15.5
```

## Example: Life Expectancy

```
Call:
lm(formula = Life.Exp ~ Population + Murder + HS.Grad + Frost,
data = statedata)

Coefficients:
(Intercept)   Population      Murder      HS.Grad        Frost
   7.10e+01     5.01e-05    -3.00e-01     4.66e-02    -5.94e-03
```

## Lasso

- Useful for high-dimensional data

- Still works when $p >> n$

- Theoretically guarantees

## Simulation Study: High-Dimensional Data

- Response: $Y$
- Predictors: $X_1, X_2, ..., X_40$
- 30 samples

```
> dim(data)
[1] 30 41
> g = lm(Y ~., data)
```

## Simulation Study: High-Dimensional Data

```
> coef(g)
(Intercept)   X1          X2          X3          X4          X5
-0.3313       1.8273      -1.3044     8.6732      -2.8432     -1.5281
X6            X7          X8          X9          X10         X11
-3.2323       -3.0793     3.0940      -0.4947     -1.5609     -1.0737
X12           X13         X14         X15         X16         X17
0.0377        3.1824      -3.0936     -0.5792     -0.2540     3.0077
X18           X19         X20         X21         X22         X23
4.3260        -1.1224     2.3879      1.7569      2.8271      -0.5614
X24           X25         X26         X27         X28         X29
2.6789        5.5562      -0.3190     -1.1525     -2.5788     1.4921
X30           X31         X32         X33         X34         X35
NA            NA          NA          NA          NA          NA
X36           X37         X38         X39         X40
NA            NA          NA          NA          NA
```

## Simulation Study: High-Dimensional Data

```
> cvlmod = cv.lars(as.matrix(data[,-1]), data$Y)
> which.min( cvlmod$cv )
[1] 24
> cvlmod$index[24]
[1] 0.232
> predict(lmod, s = 0.232, type = "coef", mode = "fraction")$coef
X1      X2      X3      X4      X5      X6      X7      X8      X9
0.6951  0.0212  0.0000  0.0000  0.0000  0.0000  0.0000  0.0276  0.0038
X10     X11     X12     X13     X14     X15     X16     X17     X18
0.0000 -0.0953  0.0000  0.0000  0.0523  0.0000  0.0000  0.0000  0.0000
X19     X20     X21     X22     X23     X24     X25     X26     X27
0.0000  0.0000  0.0000 -0.2470  0.0000  0.0000  0.0000  0.0000 -0.2140
X28     X29     X30     X31     X32     X33     X34     X35     X36
-0.3813  0.0000  0.0000  0.0000  0.1900  0.0000  0.0000  0.0000  0.0000
X37     X38     X39     X40
0.0000  0.0086  0.2222 -0.2335
```

# Simulation Study: High-Dimensional Data

```
> g = lm.ridge(Y ~., data, lambda = 1)
> g
X1        X2        X3        X4        X5        X6
1.06e-01  5.55e-01  2.84e-01  4.34e-02  2.17e-02  2.49e-02  2.53e-01
X7        X8        X9        X10       X11       X12       X13
4.12e-02  1.82e-01  2.16e-01  3.13e-02 -2.00e-01 -2.18e-01  2.80e-01
X14       X15       X16       X17       X18       X19       X20
2.90e-01  6.60e-02 -2.50e-01  1.73e-01 -1.66e-01  5.52e-02  2.26e-01
X21       X22       X23       X24       X25       X26       X27
-7.37e-02 -8.42e-02 -9.07e-05  5.40e-02  2.26e-01  1.05e-01 -2.78e-01
X28       X29       X30       X31       X32       X33       X34
-5.33e-01 -1.74e-01  1.48e-02 -2.68e-02 -5.19e-02  3.57e-01 -7.81e-02
X35       X36       X37       X38       X39       X40
-1.66e-01 -4.16e-02 -9.16e-02  4.18e-01  2.67e-01 -6.45e-01
```

## Partial Least Square Regression

Partial least squares (PLS) is a method for relating a set of input variables $X_1, X_p$ and outputs $Z_1, , Z_{p'}$. In addition, regress $Y$ over $Z_1, , Z_{p'}$.