

Ex 3.3 (a)

$$\hat{\theta}_{LS} = a^T \hat{\beta} = a^T (X^T X)^{-1} X^T y \quad E[a^T \hat{\beta}] = a^T \beta$$

$$\tilde{\theta} = c^T y \quad \text{unbiased} \quad c^T = a^T (X^T X)^{-1} X^T + \delta^T$$

$$\begin{aligned} E[\tilde{\theta}] &= E[c^T y] \\ &= E[a^T (X^T X)^{-1} X^T y + \delta^T y] \\ &= a^T (\cancel{X^T X})^{-1} \cancel{X^T} X \beta + \delta^T X \beta \\ &= \underline{a^T \beta} + \delta^T X \beta \Rightarrow \underline{\delta^T X = 0} \end{aligned}$$

$$\begin{aligned} \text{Var}(\tilde{\theta}) &= \text{Var}(c^T y) \\ &= c^T \sigma^2 c \\ &= \sigma^2 (a^T (X^T X)^{-1} X^T + \delta^T) (a^T (X^T X)^{-1} X^T + \delta^T)^T \\ &= \sigma^2 (a^T (X^T X)^{-1} X^T + \delta^T) (X (X^T X)^{-1} a + \delta) \\ &= \sigma^2 \left( \underbrace{a^T (X^T X)^{-1} X^T X (X^T X)^{-1} a}_{=0} + \underbrace{a^T (X^T X)^{-1} X^T \delta}_{=0} + \underbrace{\delta^T X (X^T X)^{-1} a}_{=0} + \delta^T \delta \right) \\ &= \sigma^2 a^T (X^T X)^{-1} a + \sigma^2 \delta^T \delta \\ &\quad \parallel \quad \text{VI} \\ &\quad \text{Var}(\hat{\theta}_{LS}) \quad 0 \end{aligned}$$

## Variable selection

IDEA: remove irrelevant variables from the model

↑  
not useful to predict/explain the response

- less variance (despite small increase of size)
- better interpretability
- better portability

When a variable is considered irrelevant

- p-value of a test  $> \alpha$ , usually  $\alpha = 0.05$
- its inclusion increases a information criterion

## INFORMATION CRITERIA

IDEA: instead of  $\hat{\theta} = \arg \min_{\theta} L(\theta)$ ,

we find  $\hat{\theta}_{ic} = \arg \min_{\theta} \{L(\theta) + \lambda J(\theta)\}$

$$\hat{\theta}_{ic} = X\beta$$

$$\downarrow$$
  

$$-2\ell(\theta)$$

$$\downarrow$$
  

$$\sum_{j=1}^p \mathbb{1}[\theta_j \neq 0]$$

GOAL: penalize larger models

$$L(\theta) = -2\ell(\theta)$$

$$J(\theta) = \sum_{j=1}^p \mathbb{1}[\theta_j \neq 0] \quad \# \text{ variables in the model}$$

$$\lambda = \begin{cases} 2 & \text{AIC} & \text{Akaike information criterion} \\ \log(n) & \text{BIC} & \text{Bayesian " " " "} \end{cases}$$

$n$  is the sample size

NB: using AIC for model selection is like using  $\alpha = 0.157$

if two explanatory variables are strongly correlated  $\rightarrow$  collinearity

extreme case: variables linearly dependent  
 $\rightarrow$  super-collinearity

in the case of super-collinearity  $(X^T X)^{-1}$  is not invertible  
 (not full rank)

Hoerl & Kennard (1970)  $X^T X \rightarrow X^T X + \lambda I_p$

$$I_p = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}$$

$$\hat{\beta}_{\text{Bridge}}^{(\lambda)} = (X^T X + \lambda I_p)^{-1} X^T y \quad \lambda \in [0; \infty)$$

when  $\lambda \in (0; \infty)$   $(X^T X + \lambda I_p)^{-1}$  exists

RIDGE REGRESSION AS A SHRINKAGE METHOD

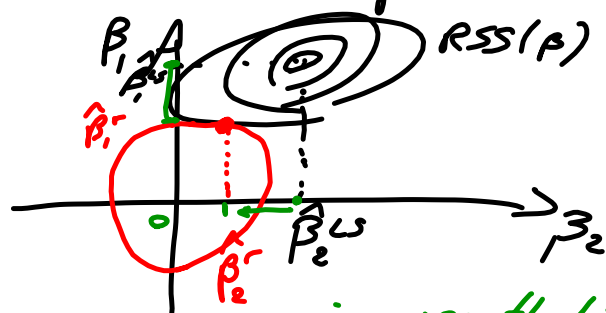
$$\bullet \hat{\beta}_{\text{Bridge}}^{(\lambda)} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

alternative formulation

$$\bullet \hat{\beta}_{\text{Bridge}}^{(\lambda)} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 \right\}$$

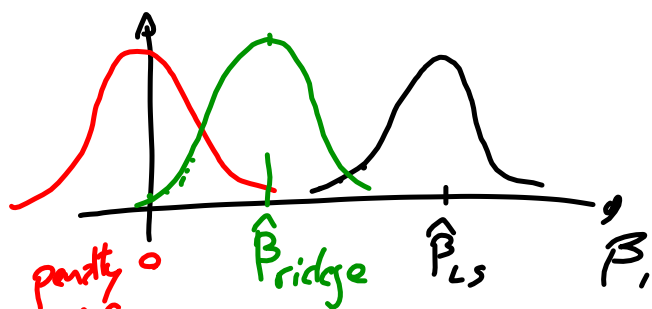
$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t$$

- one to one correspondence between  $\lambda$  and  $t$



$\rightarrow$  increase the bias, reduce the variance

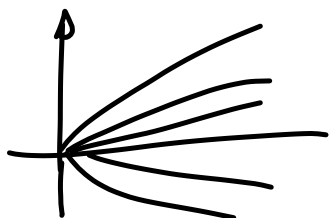
- from a Bayesian point of view, ridge estimator is the posterior mean/mode



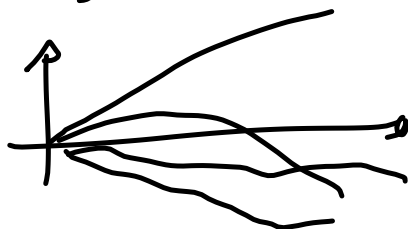
larger  $\lambda$ , more influence for the prior  
 smaller  $\lambda$ , less " " " "

---

if  $X_s$  are uncorrelated



if  $X_s$  are correlated



IMPORTANT! We need to standardize our explanatory variables before applying the ridge regression

$$E[X_j] = 0$$

$$\text{Var}[X_j] = 1$$

Expected value of ridge estimator

$$\begin{aligned} E[\hat{\beta}_{\text{ridge}}^{(\lambda)}] &= E[(X^T X + \lambda I_p)^{-1} X^T y] \\ &= E[(I_p + \lambda (X^T X)^{-1})^{-1} \underbrace{(X^T X)^{-1} X^T y}_{\hat{\beta}_{LS}}] \\ &= \underbrace{(I_p + \lambda (X^T X)^{-1})^{-1}}_{W_\lambda} E[\hat{\beta}_{LS}] \\ &= W_\lambda \beta \quad E[\hat{\beta}_{\text{ridge}}^{(\lambda)}] \neq \beta \end{aligned}$$

$$\lambda \rightarrow 0 \quad E[\hat{\beta}_{\text{ridge}}^{(\lambda)}] = \beta$$

$$\lambda \rightarrow \infty \quad E[\hat{\beta}_{\text{ridge}}^{(\lambda)}] = 0_{p \times 1} \quad \text{without intercept } \beta_0 = 0$$

important  $\lambda_a > \lambda_b \Rightarrow |\hat{\beta}_j(\lambda_a)| < |\hat{\beta}_j(\lambda_b)|$   
due to correlation

Variance of the ridge estimator

$$\begin{aligned}
 \text{Var}[\hat{\beta}_{\text{ridge}}^{(\lambda)}] &= \text{Var}[W_2 \hat{\beta}_{LS}] \\
 &= W_2 \text{Var}(\hat{\beta}_{LS}) W_2^T \\
 &= \sigma^2 W_2 (X^T X)^{-1} W_2^T
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(\hat{\beta}_{LS}) - \text{Var}(\hat{\beta}_{\text{ridge}}^{(\lambda)}) &= \sigma^2 \left[ (X^T X)^{-1} - W_2 (X^T X)^{-1} W_2^T \right] \\
 &= \sigma^2 W_2 \left[ \underline{(I + \lambda(X^T X)^{-1})} \underline{(X^T X)^{-1}} \underline{(I + \lambda(X^T X)^{-1})} - (X^T X)^{-1} \right] W_2^T \\
 &= \sigma^2 W_2 \left[ \cancel{(X^T X)^{-1}} + \lambda(X^T X)^{-2} + \lambda(X^T X)^{-2} + \lambda^2(X^T X)^{-3} - \cancel{(X^T X)^{-1}} \right] W_2^T \\
 &= \sigma^2 W_2 \left[ 2\lambda(X^T X)^{-2} + \lambda^2(X^T X)^{-3} \right] W_2^T > 0 \\
 &\Rightarrow \text{Var}(\hat{\beta}_{\text{ridge}}^{(\lambda)}) \leq \text{Var}(\hat{\beta}_{LS})
 \end{aligned}$$

Degrees of freedom

$$\hat{Y}_{LS} = \underbrace{X(X^T X)^{-1} X^T}_{H} y$$

$$df = \text{tr}(H) = p$$

$$\hat{Y}_{Ridge} = \underbrace{X(X^T X + \lambda I)^{-1} X^T}_{H_\lambda} y$$

$$\text{effective } df = \text{tr}(H_\lambda)$$

$$\text{tr}(H_\lambda) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$

$$\lambda = 0 \rightarrow df = p$$

$$\lambda \rightarrow \infty \rightarrow df \rightarrow 0$$

$$AIC : L(\beta) + 2J(\theta)$$

$$\text{least square} : -2\ell(\beta) + 2p$$

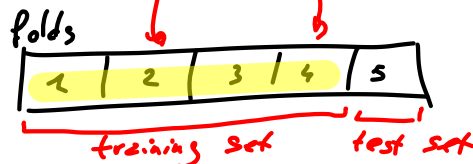
$$\text{ridge} : -2\ell(\beta) + 2 \sum \frac{d_j^2}{d_j^2 + \lambda}$$

AIC( $\lambda$ ) we can find  $\hat{\lambda}$  as  $\arg\min_{\lambda} AIC(\lambda)$

Usually  $\lambda$  is computed by cross-validation

sample:  $x_1, x_2, \dots, x_n$

split in  $K$  folds ( $K=5$  or  $K=10$ )



$$\text{err}_5 = \sum_{i=1}^n y_i - x_{(5)} \hat{\beta}_{(5)}^2$$

train the model on  $K-1$  folds  $\rightarrow \hat{\beta}_{(5)}$   
 evaluate its performance on the remaining fold  $\text{err}_5$

repeat the procedure for  $K$  times, always using a different fold as test set

$\text{err}_1$  : train set  $\{2, 3, 4, 5\}$   
           test set  $\{1\}$

$\text{err}_2$  : train set  $\{1, 3, 4, 5\}$   
           test set  $\{2\}$

$\vdots$

$\text{err}_5$

$\hat{\beta}$  computed by using all observation not belonging to fold 5

$$\sum_{k=1}^K \text{err}_k(\lambda)$$

$\lambda$  that minimizes

More about shrinkage

$$X_{n \times p} = U \Delta V^T$$

$U = n \times n$  orthogonal matrix, whose columns span the column space of  $X$

$V = p \times n$  orthogonal matrix, whose columns span the row space of  $X$

$\Delta = n \times n$  diagonal matrix,  $d_{ii}$  are single value of  $X_i$

$$U^T U = I_n = \underline{V V^T}$$

LS estimator  $\hat{\beta}_{LS}$

$$\begin{aligned} \hat{\beta}_{LS} &= (X^T X)^{-1} X^T y \\ &= (V \Delta U^T U \Delta V^T)^{-1} V \Delta U^T y \\ &= (V \Delta^2 V^T)^{-1} V \Delta U^T y \\ &= V \Delta^2 V^T V \Delta U^T y \\ &= \underline{V \Delta^2 \Delta} U^T y \end{aligned}$$

$$\begin{aligned} \hat{y}_{LS} &= X \hat{\beta}_{LS} \\ &= U \Delta V^T V \Delta^{-2} \Delta U^T y \\ &= U \Delta^2 \Delta^{-2} U^T y \\ &= U U^T y \end{aligned}$$

$$\sum_{j=1}^p \frac{d_j}{d_j^2}$$

$$\begin{aligned} \hat{\beta}_{ridge}^{(\lambda)} &= (X^T X + \lambda I)^{-1} X^T y \\ &= (V \Delta U^T U \Delta V^T + \lambda I_p)^{-1} V \Delta U^T y \\ &= (V \Delta^2 V^T + \lambda V V^T) V \Delta U^T y \\ &= V (\Delta^2 + \lambda I_p)^{-1} V^T V \Delta U^T y \\ &= V (\Delta^2 + \lambda I_p)^{-1} \Delta U^T y \end{aligned}$$

$$\begin{aligned} \hat{y}_{ridge} &= X \hat{\beta}_{ridge} \\ &= U \Delta V^T V (\Delta^2 + \lambda I_p)^{-1} \Delta U^T y \\ &= U \Delta^2 (\Delta^2 + \lambda I_p)^{-1} U^T y \end{aligned}$$

$$\sum_{j=1}^p \frac{d_j}{d_j^2 + \lambda}$$

$$\sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$



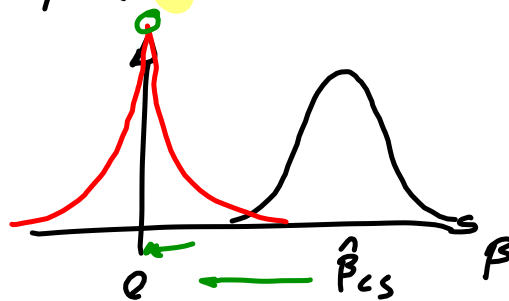
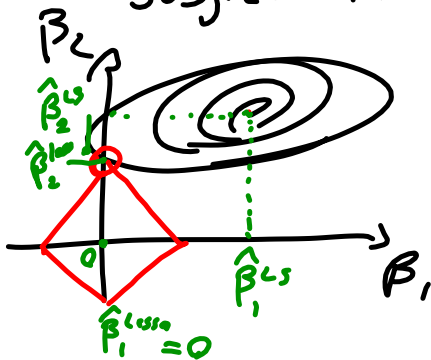
## LASSO

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

or, equivalently

$$\hat{\beta}_{\text{Lasso}} = \arg\min_{\beta} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 \right.$$

subject to  $\sum_{j=1}^p |\beta_j| \leq t$



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

$q=1 \rightarrow \text{LASSO}$

$q=2 \rightarrow \text{RIDGE REGRESSION}$

$1 < q < 2 \approx \text{average between lasso and ridge penalties}$

## ELASTIC NET

penalty :  $\lambda \left( \alpha \sum_{j=1}^p \beta_j^2 + (1-\alpha) \sum_{j=1}^p |\beta_j| \right)$



# STK4030 - Statistical Learning: Advanced Regression and Classification

Riccardo De Bin

`debin@math.uio.no`

Procedure	Logic
BE-only	Estimate full model on $x_1, \dots, x_k$ . Repeat: while the least significant term has $P \geq \alpha_2$ , remove it and re-estimate the model
BE	<p>Estimate full model on <math>x_1, \dots, x_k</math>. If the least significant term has <math>P \geq \alpha_2</math>, remove it and re-estimate; otherwise stop.</p> <p>Again: if the least significant term has <math>P \geq \alpha_2</math>, remove it and re-estimate; otherwise stop</p> <p>Repeat</p> <ul style="list-style-type: none"> <li>• if most significant excluded term has <math>P &lt; \alpha_1</math>, add it and re-estimate;</li> <li>• if least significant included term has <math>P \geq \alpha_2</math>, remove it and re-estimate;</li> </ul> <p>until neither action is possible.</p>
FS-only	Estimate null model. Repeat: while the most significant excluded term has $P < \alpha_1$ , add it and re-estimate.
FS	<p>Estimate null model. If the most significant excluded term has <math>P &lt; \alpha_1</math>, add it and re-estimate; otherwise stop.</p> <p>Again: if the most significant excluded term has <math>P &lt; \alpha_1</math>, add it and re-estimate; otherwise stop.</p> <p>Repeat</p> <ul style="list-style-type: none"> <li>• if least significant included term has <math>P \geq \alpha_2</math>, remove it and re-estimate;</li> <li>• if most significant included term has <math>P &lt; \alpha_1</math>, add it and re-estimate; until neither action is possible.</li> </ul>

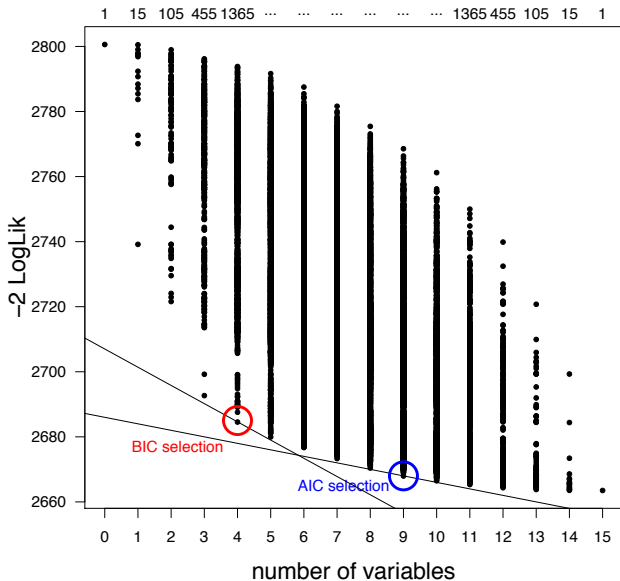
from: Royston & Sauerbrei (2008)

**Table 2.2** Myeloma study (65 patients, 48 events). Results of applying variable selection strategies.<sup>a</sup>

Variable	Nominal significance level, $\alpha$						All-subsets AIC
	0.05			0.157			
	Full	BE	FS	Full	BE	FS	
$x_1$	*	✓	✓	*	✓	✓	✓
$x_2$			✓			✓	
$x_3$	*	✓		*	✓		✓
$x_4$		✓		*	✓		✓
$x_5$							
$x_6$		✓		*	✓		✓
$x_7$	*	✓		*	✓		✓
$x_8$				*	✓		✓
$x_9$							
$x_{10}$							
$x_{11}$							
$x_{12}$	*	✓		*	✓		✓
$x_{13}$	*	✓		*	✓		✓
$x_{14}$				*			
$x_{15}$							
$x_{16}$							

<sup>a</sup> BE(0.01) and FS(0.01) selected only one variable,  $x_1$ ; \* denotes that a variable is significant at the relevant  $\alpha$ -level in the full model.

from:  
Royston &  
Sauerbrei  
(2008)



**Table 2.3** Educational body-fat data. Full model and that selected by BE(0.05). Final three columns give details of the full model excluding  $x_6$ .

Variable	Full model			BE(0.05)			Full model excl. $x_6$		
	$\hat{\beta}$	SE	$\hat{\beta}/SE$	$\hat{\beta}$	SE	$\hat{\beta}/SE$	$\hat{\beta}$	SE	$\hat{\beta}/SE$
$x_1$	0.074	0.032	2.31	0.056	0.024	2.35	0.211	0.034	6.20
$x_2$	-0.019	0.067	-0.28				0.227	0.074	3.08
$x_3$	-0.249	0.191	-1.30	-0.322	0.121	-2.65	-0.915	0.212	-4.32
$x_4$	-0.394	0.234	-1.68				-0.378	0.278	-1.36
$x_5$	-0.119	0.108	-1.10				0.150	0.124	1.21
$x_6$	0.901	0.091	9.90	0.774	0.033	23.26	-	-	-
$x_7$	-0.146	0.144	-1.02				0.163	0.166	0.98
$x_8$	0.178	0.146	1.22				0.231	0.173	1.33
$x_9$	-0.041	0.245	-0.17				-0.095	0.291	-0.33
$x_{10}$	0.185	0.220	0.85				-0.053	0.259	-0.21
$x_{11}$	0.178	0.170	1.04				-0.066	0.200	-0.33
$x_{12}$	0.277	0.207	1.34				0.058	0.244	0.24
$x_{13}$	-1.830	0.529	-3.46	-1.943	0.406	-4.78	-2.692	0.620	-4.34

from: Royston & Sauerbrei (2008)

```
library(Biobase)
library(breastCancerVDX)

# ids of genes FL0T1
idFL0T1 <- which(fData(vdx)[,5] == 10211)

# ids of ERBB2
idERBB2 <- which(fData(vdx)[,5] == 2064)

# get expression levels of probes mapping to FL0T genes
X <- t(exprs(vdx)[idFL0T1,])
X <- sweep(X, 2, colMeans(X))

# get expression levels of probes mapping to FL0T genes
Y <- t(exprs(vdx)[idERBB2,])
Y <- sweep(Y, 2, colMeans(Y))

# regression analysis
summary(lm(formula = Y[,1] ~ X[,1] + X[,2] + X[,3] + X[,4]))
```

from: Van der Wieringen (2015)

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.0000	0.0633	0.0000	1.0000
X[, 1]	0.1641	0.0616	2.6637	0.0081 **
X[, 2]	0.3203	0.3773	0.8490	0.3965
X[, 3]	0.0393	0.2974	0.1321	0.8949
X[, 4]	0.1117	0.0773	1.4444	0.1496

---

Signif. codes: 0 \*\*\* 0.001 \*\* 0.01 \* 0.05 . 0.1 1

Residual standard error: 1.175 on 339 degrees of freedom  
Multiple R-squared: 0.04834, Adjusted R-squared: 0.03711  
F-statistic: 4.305 on 4 and 339 DF, p-value: 0.002072

from: Van der Wieringen (2015)



