

02. Linear Model

least-square approach minimizes  $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$

if  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ ,  $\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$ ,  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ .

coefficient accuracy:  $SE(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$ ,  $SE(\hat{\beta}_0)^2 = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$   $\Rightarrow t = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \sim t(n-2)$

overall fitness: (1)  $RSE = \sqrt{\frac{RSS}{n-2}}$

(2)  $R^2 = \frac{TSS - RSS}{TSS} \rightarrow$  fraction of variance explained (how much does the model reduce TSS).

where  $TSS = \sum (y_i - \bar{y})^2$ ,  $RSS = \sum (y_i - \hat{y}_i)^2$ .

for a simple regression model:  $r^2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$

hierarchy principle: If we include an interaction in a model, we should also include the main effects, even if the p-values associated with their coefficients are not significant.

How to allow qualitative predictors: by creating a dummy variable, i.e. create a new variable s.t. with no dummy variables  $y_i = \beta_0 + \epsilon_i$  baseline

$x_{i1} = \begin{cases} 1 & \text{if } i\text{th observation belongs to A} \\ 0 & \text{if } i\text{th observation does not belong to A} \end{cases} \Rightarrow y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots = \begin{cases} \beta_0 + \epsilon_i & \text{... A} \\ \beta_0 + \beta_1 + \epsilon_i & \dots B \\ \beta_0 + \beta_2 + \epsilon_i & \dots B \end{cases}$

Is at least one predictor useful?:  $F = \frac{(TSS - RSS)/P}{RSS/(n-P-1)} \sim F_{P, n-P-1}$

which predictors are useful?: subset selection

which model (different combination of predictors) is the best: model selection.

## W03 Linear Model Selection and Regularization

1. Linear Model Selection subset selection ( $R^2$ , RSS) + model selection (AIC, BIC, CP, adj.  $R^2$ , CV)

### 1.1 Subset Selection

#### 1.1.1 Exhausted / best subset.

method: consider all possible models; and compute least square fit for all of them  
then choose the best model based on some criterion.

steps:

(1) let  $M_0$  be the null model, which contains only the intercept  $\bar{y}$ .

(2) For  $k=1, 2, \dots, p$ ,

(a) fit all  $\binom{p}{k}$  models

(b) Pick the best  $\xrightarrow{\text{smallest RSS, or largest } R^2}$  model and let it be  $M_k$ .  $\xrightarrow{\text{model selection}}$

(3) Compare all  $p$  models and select the one with the least test error

problem: We need to consider all  $1 + \binom{p}{1} + \binom{p}{2} + \dots + \binom{p}{p} = 2^p$  models, if  $p=40$ ,  $2^p = 1$  trillion.  
Also the problem of overfitting and high variance  $\Rightarrow$  stepwise selection.

#### 1.1.2 Forward Selection.

method: start with null model, then pick the next predictor that reduces RSS the most;  
put that predictor in the model, repeat until some stopping rule is satisfied.

steps:

(1) let  $M_0$  be the null model

(2) For  $k=0, \dots, p-1$

(a) consider all  $p-k$  models that augment the predictors in  $M_k$  with one more predictor

(b) pick the best model among all  $p-k$  models and let it be  $M_{k+1}$ .

(3) Compare all  $p$  models and select the one with the least test error  $\xrightarrow{\text{Mo}}$

note: greater computing advantages since only  $p + (p-1) + \dots + 1 + 1 \approx p^2$  are considered.

problem: it's not guaranteed the forward selection would find the best combination of predictors.  
especially if there is correlation between predictors.

#### 1.1.3 Backward Selection:

method: start with full model, then remove the predictor with the largest P-value; then fit the model with  $p-1$  predictors, repeat until some stopping rule is satisfied.

steps:

(1) let  $M_p$  be the full model, which contains all  $p$  predictors.

(2) For  $k=p, p-1, \dots, 1$ ,

(a) Consider all  $k$  models that contain all but one predictors in  $M_p$ .

(b) pick the best model among all  $k$  models and let it be  $M_k$ .

(3) Compare all  $p$  models and select the one with the least test error

note: backwards selection requires  $n > p$ , whereas forward selection can be used even when  $p > n$ .

problem: same as forward.

## 1.2 Model Selection.

To choose an 'optimal' member in the path of models  
(i.e. # predictors & which predictors).

Why test error? - The full model will always have the lowest training error, i.e. overfitting

- Training error is not a represent of how good the model predict on new data.

how? ① by adjusting training error to account for overfitting; ② by separating test/train data.

1. Mallon's Cp. ( $C_p$ )

$$C_p = \frac{1}{n} (RSS + 2d\hat{\sigma}^2)$$

→ # predictors, e.g.  $M_3$  has 4 predictors including intercept.  
 → estimate of  $\text{Var}(E)$ , i.e. variance of error.

problem:  $C_p$  must have  $n > p$ ; if  $n \leq p$ ,  $M_p$  is not defined and has  $C_p = \infty$ .

2. Akaike Information Criterion (AIC).

$$AIC = -2\log L + 2d$$

note: For linear model with Gaussian error,  $C_p = AIC$  because  $-2\log L = \frac{RSS}{\hat{\sigma}^2}$ .  
 i.e.  $AIC = \frac{RSS}{\hat{\sigma}^2} + 2d \Leftrightarrow C_p = \frac{1}{n} (RSS + 2d\hat{\sigma}^2)$   
 i.e.  $\frac{nC_p}{\hat{\sigma}^2} = AIC$ .

proof: for linear model with Gaussian errors.  $\Rightarrow E_i \sim N(0, \sigma^2) \Rightarrow Y_i \sim N(\hat{Y}_i, \sigma^2)$ .  
 $\Rightarrow L = \exp\left\{-\frac{(Y_i - \hat{Y}_i)^2}{2\sigma^2}\right\} = \exp\left\{-\frac{RSS}{2\sigma^2}\right\} \Rightarrow \log L = -\frac{RSS}{2\sigma^2} \Rightarrow -2\log L = \frac{RSS}{\sigma^2}$

problem: also requires  $n > p$ .

3. Bayesian Information Criterion (BIC).

$$BIC = \frac{1}{n} (RSS + \log(n)d\hat{\sigma}^2)$$

note: BIC places a heavier penalty on size of model than  $C_p$  or AIC

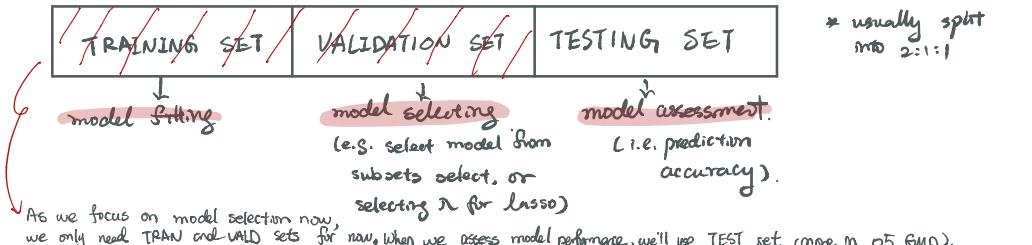
adjusted  $R^2 \leftarrow$  4. adjusted  $R^2$

$$\text{Adjusted } R^2 = 1 - \frac{RSS/(n-d-1)}{TSS/(n-1)} \quad [R^2 = \frac{TSS-RSS}{TSS} \text{ doesn't account for size of model}]$$

note: • effective part  $RSS/(n-d-1) \Rightarrow$  place a penalty on including unnecessary variables since as  $d \uparrow$ ,  $RSS/(n-d-1)$  may ↑ or ↓, and so does adjusted  $R^2$ .  
 • adjusted  $R^2$  doesn't require  $n > p$ .

separate test & train data. ← 5. Cross-Validation (CV)

Ideally, we would split\* the data into 3 sections for model selection and assessment.



• Validation:

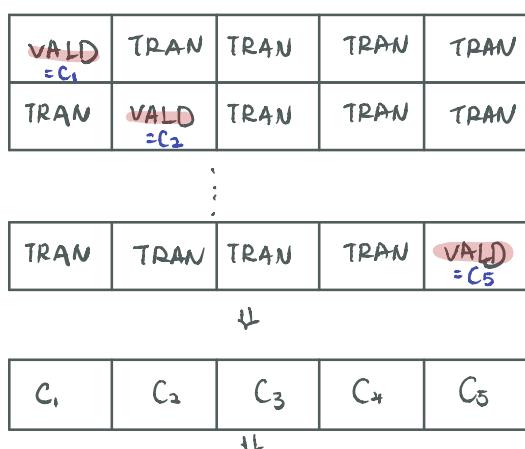
TRAINING SET	VALIDATION SET
fitting the model.	MSE / MCR.



problem:

- even though the validation estimated test error tells us which model has the lowest MSE, the validation estimated test error itself varies a lot.  $\Rightarrow$  no actual level of test error.
- a fraction of data is not used to train the model  $\Rightarrow$  waste lots of data.

### Cross-validation:



\* 5-fold, i.e.  $k=5$ .

steps:

- (1) randomly split the data into  $K$  parts.
- (2) each part will be the validation part in turn, and the rest  $K-1$  pieces are the training part
- (3) fit the model using the training data and compute the error using the validation part. repeat  $K$  times
- (4) compute the average error

$$CV \text{ error} = CV_{(k)} = \sum_{k=1}^K \frac{n_k}{n} MSE_k \quad \text{where } MSE_k = \sum_{i \in C_k} \frac{(y_i - \hat{y}_i)^2}{n_k},$$

or  $\sum_{k=1}^K \frac{n_k}{n} Err_k,$

$$Err_k = \sum_{i \in C_k} \frac{|y_i - \hat{y}_i|}{n_k},$$

$n_k = \# \text{ observ. in } k^{\text{th}} \text{ part}$

**note:** - letting  $K=n$ , yields  $n$ -fold or leave-one-out cross validation (LOOCV)  
i.e. every data point gets to be training data as well as validation data.

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{y}_i}{h_i} \right)^2, \quad \text{where } h_i \text{ is the leverage (i.e. diagonal of the "hat" matrix).}$$

- Generally  $K$ -fold preferable on LOOCV.

• Computation time shorter.

• LOOCV has high variance because each training data is highly similar/correlated

Thus, even though LOOCV has low bias, the tradeoff between bias and variance is not worth it

**Important:** subset selection should be done within CV process! cannot select the predictors and then do CV!

\* One-standard-error rule: If test error are very close for some model, we first calculate the standard error of the test error for each model, then select the smallest model for which the test error is within one standard error of the lowest point on the curve.

## 2. Shrinkage/Regularization (shifting coefficients toward or to exactly 0)

### 2.1 Ridge Regression

$$\underset{\lambda}{\text{minimize}} \quad \underset{\text{RSS}}{\underbrace{\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2}_{\downarrow}} + \underset{\text{shrinkage penalty/L2 penalty}}{\underbrace{\lambda \sum_{j=1}^p \beta_j^2}_{\downarrow}}$$

shrinkage penalty/L2 penalty

• RSS term ensure fitness of data;

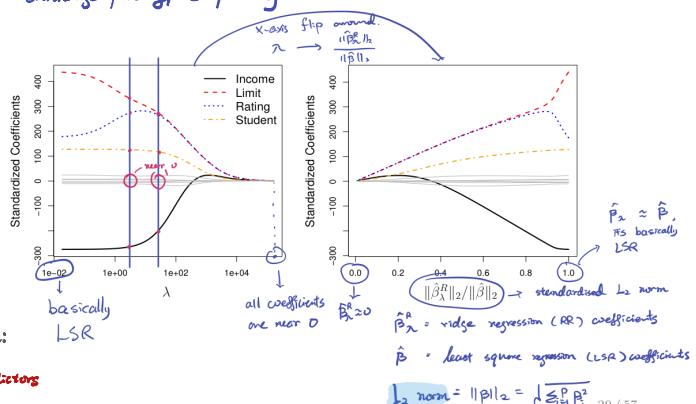
shrinkage penalty has the effect of shrinking estimated coefficients towards 0. (but  $\neq 0$ ).

• tuning parameter  $\lambda$  controls the relative impact of these two terms.

• choosing a good tuning parameter  $\lambda$  critical; can use CV to determine.

• Since  $\sum_{j=1}^p \beta_j^2$  is multiplied by a constant  $\lambda$ , the Ridge coefficient estimates can change quite dramatically. Thus, it's best to standardize first:

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_{ij})^2}} \rightarrow \text{standardized predictors}$$

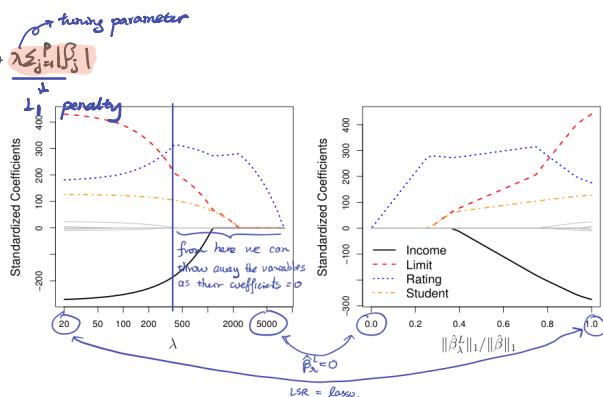


## 2.2 Lasso.

$$\text{minimize}_{\beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j|$$

↓  
RSS

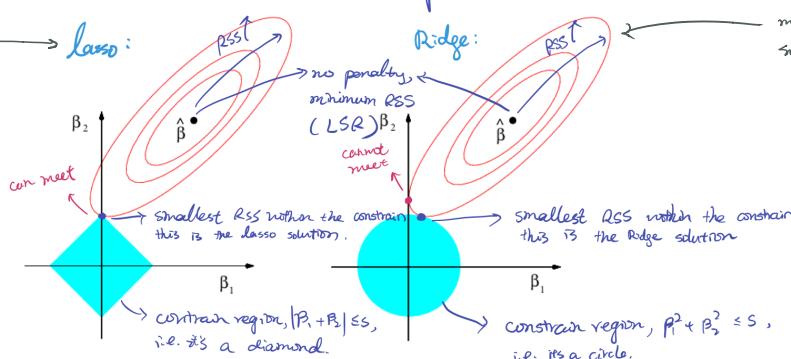
- L1 penalty has the effect of forcing some of the coefficient estimates to be exactly 0 when  $n$  is large enough
- Thus Lasso can perform **predictor selection** like subset selection, i.e. it yields **sparse model**
- CV again is used to choose  $\lambda$
- Why does Lasso can force coefficient estimates to be 0 but not Ridge:



## The Lasso Picture (2 parameters example)

minimize RSS  
subject to  $\|\beta\|_1 \leq S$

→ **Lasso:**



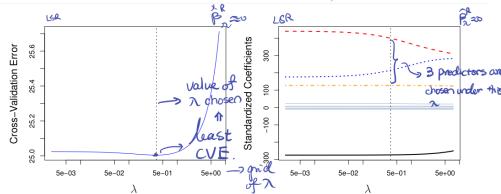
→ RSS  
minimize to  $\|\beta\|_1 \leq S$   
subject

\* the diamond has corners (for higher dimensions, we have edges and corners), so one parameter is on y-axis and it's easier for RSS circle to hit it.

37/57

### \* How to choose $\lambda$ using CV:

- choose a grid of  $\lambda$
- compute CVE for each  $\lambda$
- select the  $\lambda$  that has the lowest CVE.



## 3. Dimension reduction (linear combination of predictors, e.g. PCR, PLS).

method: transform the predictors into linear combinations and then fit a least squares model

step:

+ new predictor  
+ no constant

- let  $Z_m = \sum_{j=1}^p \phi_{mj} X_j$  → original predictors
- fit the model:  $y_i = \beta_0 + \sum_{m=1}^M \theta_m Z_m + \epsilon_i$ ,  $i=1, 2, \dots, n$ .

note:

- $M < p \Rightarrow$  dimension reduction
- $\sum_{m=1}^M \theta_m Z_m = \sum_{m=1}^M \theta_m \sum_{j=1}^p \phi_{mj} X_j = \sum_{j=1}^p \sum_{m=1}^M \theta_m \phi_{mj} X_j = \sum_{j=1}^p \beta_j X_j$  where  $\beta_j = \sum_{m=1}^M \theta_m \phi_{mj}$ .  
or it's a special case of OLS but with constraints  $\beta_j = \sum_{m=1}^M \theta_m \phi_{mj}$ .

How to choose linear combinations of predictors:

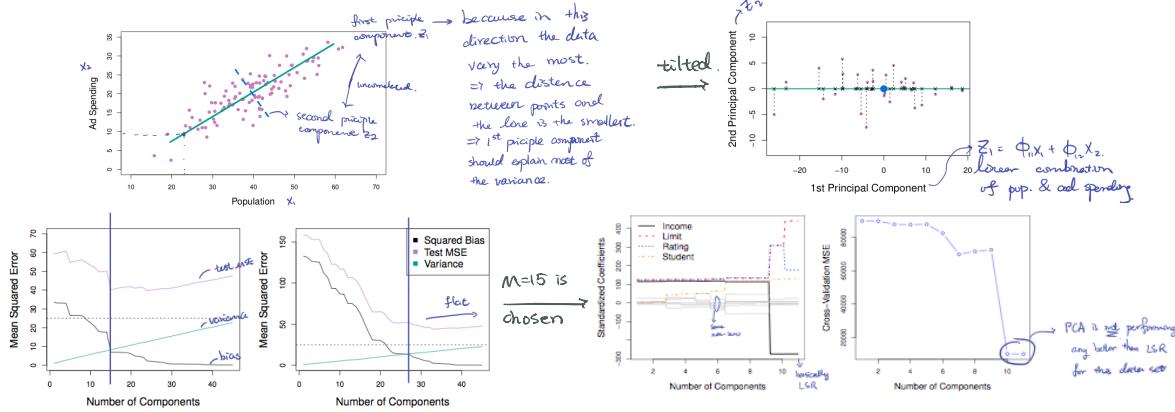
### 3.1 Principle Components Regression (PCR).

method: apply Principle Components Analysis (PCA) to define linear combinations of the predictors

1<sup>st</sup> principle component: normalized linear combination of predictors with the largest variance.

2<sup>nd</sup> ... : largest variance, subject to being uncorrelated with the 1<sup>st</sup>

N<sup>th</sup> principle component: ...  
⇒ replace the original correlated predictors with a small set of independent principle components that capture their joint variation



problem: these directions (or PCs) are identified in an unsupervised way

- ⇒ there is no guarantee that the directions that best explain the predictors will also be the best directions to use for predicting the response.
- ⇒ PLS identifies these new features in a supervised way, i.e. it makes use of the response in order to identify new features that are related to the response. ⇒ PLS explains both response and the predictors.

### 3.2 Partial Least Square (PLS).

- standardise  $p$  predictors, i.e.  $\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\sum_i (x_{ij} - \bar{x}_{ij})^2}}$
- $z_1$ : set  $\phi_{1j}$  equal to coefficients from simple linear regression of  $Y$  onto  $x_j$ . one-by-one.  
 $\Rightarrow z_1 = \bar{y}_j \phi_{1j} x_j$ , so  $\phi_{1j}$  is proportional to the correlation between  $Y$  and  $x_j$ .
- PLS places the highest weight on the predictor that are most related to  $Y$ .
- take the residuals and repeat the above