

Module 6: Model Selection and Regularization

Lecture 14
Feb 20th, 2023
Ch 6.1- 6.2



SPARTANS WILL.

You can leave when and if you need to

Nothing on the syllabus is as important as **your well-beings**.

It is important to **ask for help** in this class and beyond

We should prioritize flexibility, grace, and care for each other

Bootstrap vs CV

The goal: quantify the uncertainty associated with a given estimator (model) or statistical learning method.

$$Y = \hat{\beta}_0 + \hat{\beta}_1 X + \epsilon$$

↑ ↑ ↑

x is really associated with y ?

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \epsilon.$$

- The linear model has advantages in terms of its **interpretability** and often shows good **predictive performance**.
- How to improve linear model? Beyond Ordinary Least Square(OLS)



Why Consider Alternatives to

- Prediction Accuracy: especially when $p > n$. Need to control variance
- Model Interpretability: Small number of features are easier to understand and design experiment.

Handwritten diagram illustrating a model structure. On the left, a vector Y is shown. To its right is an equals sign. Further right is a vertical vector. The top of this vector is labeled g_i and the bottom is labeled g_{n+1} .

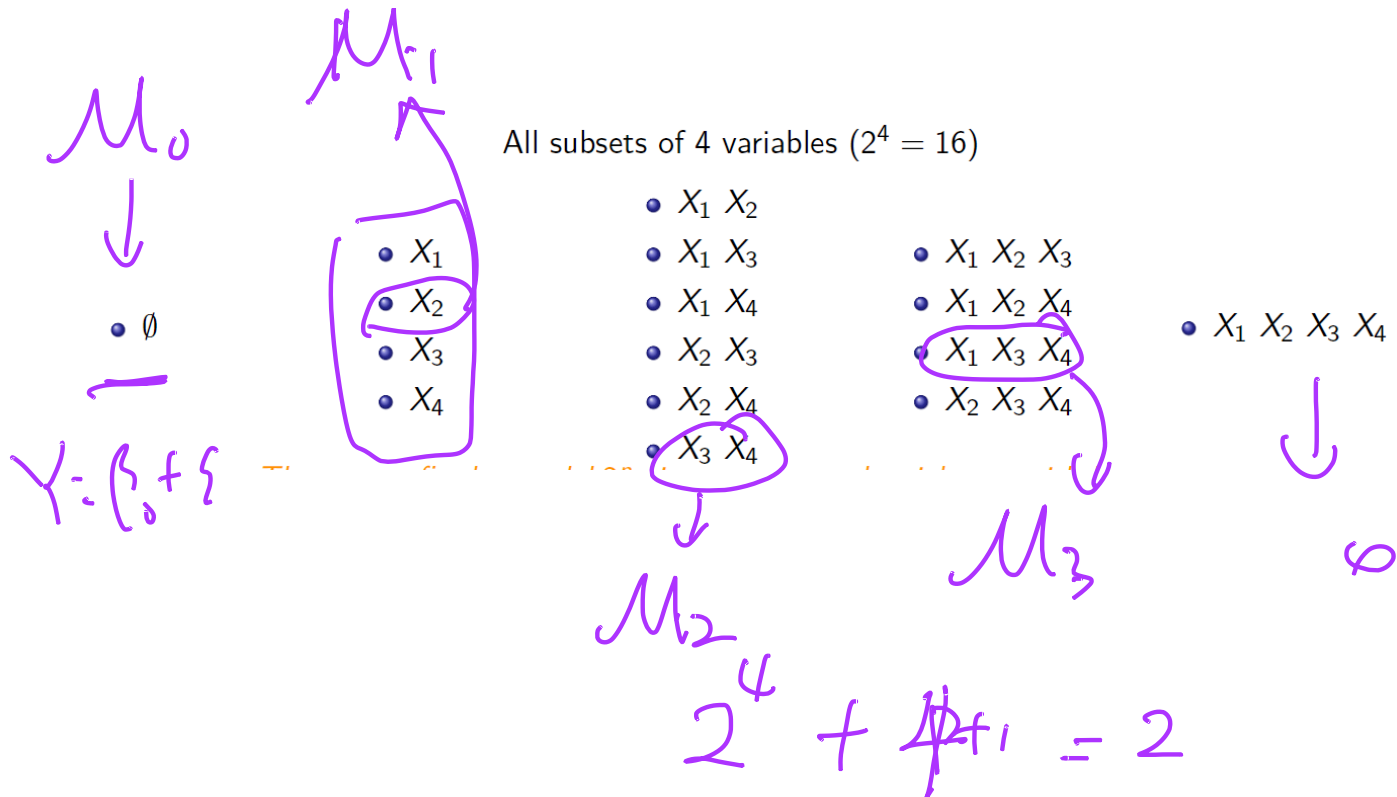
Three Classes of Methods



- **Subset Selection.** We identify a subset of the p predictors then fit a model using least squares on the reduced set of variables.
- **Shrinkage.** We fit a model involving all p predictors, but the estimated coefficients are shrunk towards zero relative to the OLS estimates. This shrinkage (also known as regularization) has the effect of **reducing** variance and can also perform variable selection.
- **Dimension Reduction.**

- Best subset selection
- Stepwise selection

Too many possible models



Algorithm 6.1 *Best subset selection*

1. Let \mathcal{M}_0 denote the *null model*, which contains no predictors. This model simply predicts the sample mean for each observation.

2. For $k = \underline{1}, 2, \dots, \overset{p}{p}$:

(a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.

(b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here *best* is defined as having the smallest RSS, or equivalently largest R^2 .

3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .

from exam
 $\mathcal{M}_0, \mathcal{M}_1, \mathcal{M}_2, \dots$

\mathcal{M}_p

We train a model using four variables, X_1, X_2, X_3, X_4 . We're interested in getting a subset of the variables to use. The following table shows the mean squared error and the R^2 value computed for the model learned using each possible subset of variables.

	Training MSE ($\times 10^4$)	k-fold CV Testing Error
Null model	8.76	10.08
X1	8.63	9.98
X2	7.42	8.01
X3	8.16	8.3
X4	8.33	9.06
X1,X2	4.33	7.47
X1,X3	5.82	5.22
X1,X4	3.17	4.23
X2,X3	4.07	3.78
X2,X4	3.31	4.01
X3,X4	3.06	4.16
X1,X2,X3	3.08	5.49
X1,X2,X4	3.55	4.02
X1,X3,X4	2.97	4.23
X2,X3,X4	2.98	3.17
X1,X2,X3,X4	2.16	4.39

- 1 What subset of variables is found for each of the sets $\mathcal{M}_0, \mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4$ when using best subset selection?
- 2 What subset of variables is returned using best subset selection?

Group work:

M_k

$2^p \cdot k$

	Training MSE ($\times 10^7$)	k-fold CV Testing Error
Null model	8.76	10.08
X1	8.63	9.98
X2	7.42	8.01
X3	8.16	8.3
X4	8.33	9.06
X1,X2	4.33	7.47
X1,X3	5.82	5.22
X1,X4	3.17	4.23
X2,X3	4.07	3.78
X2,X4	3.31	4.01
X3,X4	3.06	4.16
X1,X2,X3	3.08	5.49
X1,X2,X4	3.55	4.02
X1,X3,X4	2.97	4.23
X2,X3,X4	2.98	3.17
X1,X2,X3,X4	2.16	4.39

μ_0 (Null model)
 μ_1 (X1)
 μ_2 (X2, X3, X4)
 μ_3 (X1, X3, X4)
 μ_4 (X2, X3, X4)

\emptyset
 X_1
 X_2
 X_3
 X_4

$X_1 X_2$
 $X_1 X_3$
 $X_1 X_4$
 $X_2 X_3$
 $X_2 X_4$
 $X_3 X_4$

$X_1 X_2 X_3$
 $X_1 X_2 X_4$
 $X_1 X_3 X_4$
 $X_2 X_3 X_4$

$X_1 X_2 X_3 X_4$

Challenges for Best Subset selection

$$2^p$$

$$p=40$$

$$2^{20,000}$$

Forward Stepwise Selection

Algorithm 6.2 *Forward stepwise selection*

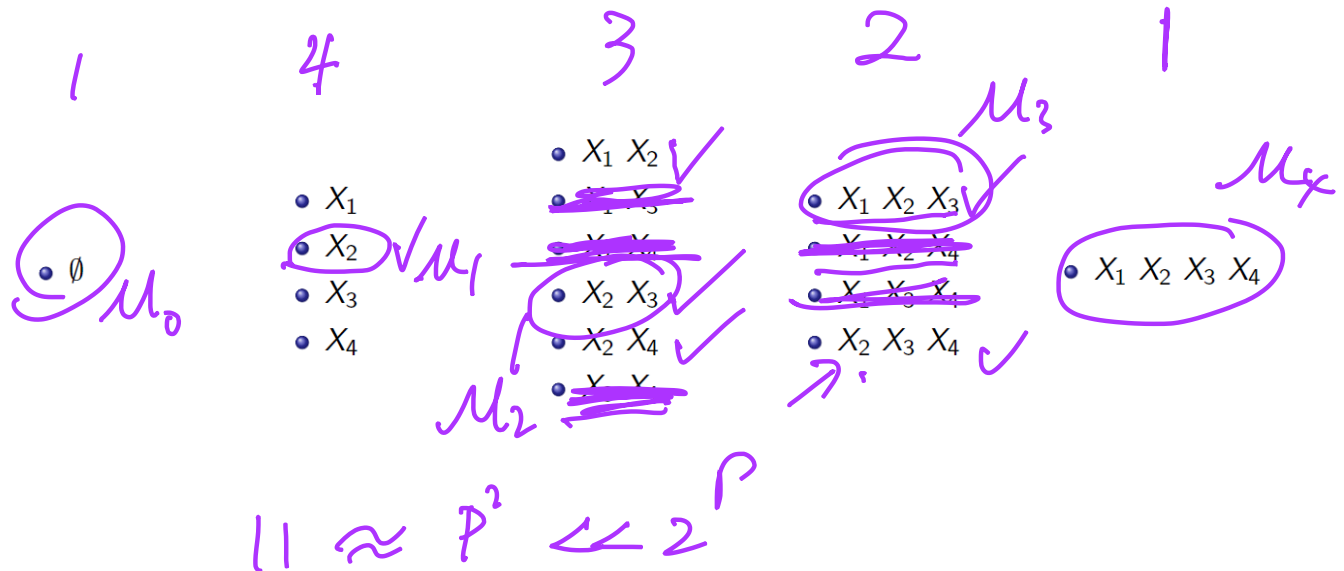
1. Let \mathcal{M}_0 denote the *null* model, which contains no predictors.
 2. For $k = 0, \dots, p-1$:

$\boxed{X_2}^+$ $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$

 - (a) Consider all $p-k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the *best* among these $p-k$ models, and call it \mathcal{M}_{k+1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

$\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p \rightarrow CV$

An Example for Forward Stepwise

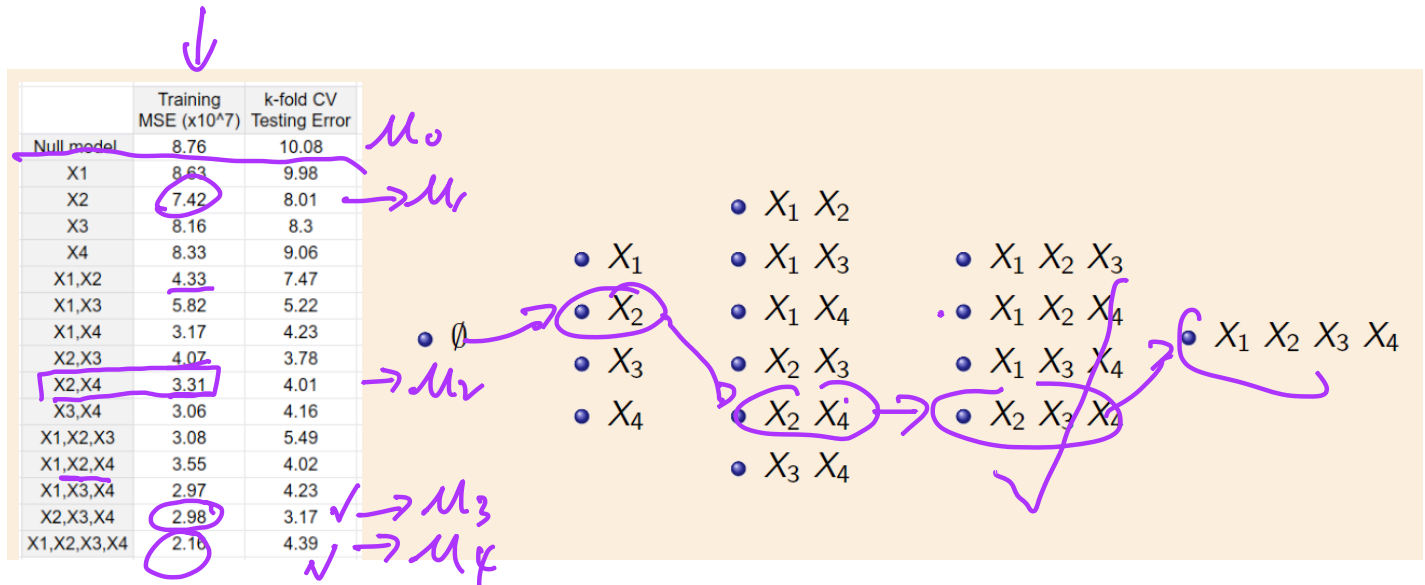


We train a model using four variables, X_1, X_2, X_3, X_4 . We're interested in getting a subset of the variables to use. The following table shows the mean squared error and the R^2 value computed for the model learned using each possible subset of variables.

	Training MSE ($\times 10^8$)	k-fold CV Testing Error
Null model	8.76	10.08
X1	8.63	9.98
X2	7.42	8.01
X3	8.16	8.3
X4	8.33	9.06
X1,X2	4.33	7.47
X1,X3	5.82	5.22
X1,X4	3.17	4.23
X2,X3	4.07	3.78
X2,X4	3.31	4.01
X3,X4	3.06	4.16
X1,X2,X3	3.08	5.49
X1,X2,X4	3.55	4.02
X1,X3,X4	2.97	4.23
X2,X3,X4	2.98	3.17
X1,X2,X3,X4	2.16	4.39

- 1 What subset of variables is found for each of the sets $\mathcal{M}_0, \mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4$ when using forward selection?
- 2 What subset of variables is returned using forward subset selection?

Group Work



Pros and Cons of Forward Selection

Pros

less models to be tested
Computationally feasible.

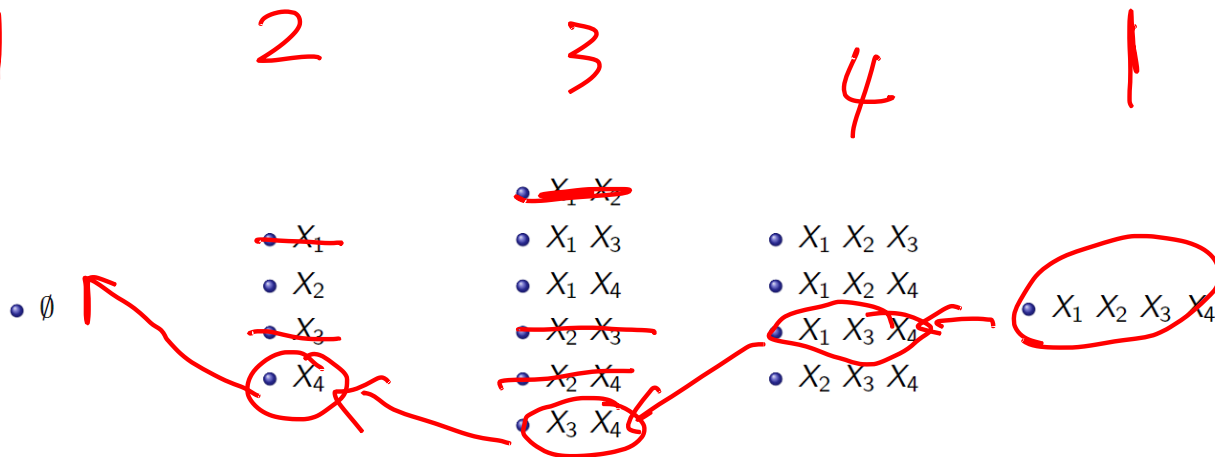
Cons

stuck on
local solution

Algorithm 6.3 *Backward stepwise selection*

1. Let \mathcal{M}_p denote 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 of the predictors in \mathcal{M}_k , for a total of $k - 1$ predictors.
 - (b) Choose the *best* among these k models, and call it \mathcal{M}_{k-1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

An example for Backward Stepwise



11 models

We train a model using four variables, X_1, X_2, X_3, X_4 . We're interested in getting a subset of the variables to use. The following table shows the mean squared error and the R^2 value computed for the model learned using each possible subset of variables.

	Training MSE ($\times 10^4$)	k-fold CV Testing Error
Null model	8.76	10.08
X1	8.63	9.98
X2	7.42	8.01
X3	8.16	8.3
X4	8.33	9.06
X1,X2	4.33	7.47
X1,X3	5.82	5.22
X1,X4	3.17	4.23
X2,X3	4.07	3.78
X2,X4	3.31	4.01
X3,X4	3.06	4.16
X1,X2,X3	3.08	5.49
X1,X2,X4	3.55	4.02
X1,X3,X4	2.97	4.23
X2,X3,X4	2.98	3.17
X1,X2,X3,X4	2.16	4.39

- 1 What subset of variables is found for each of the sets $\mathcal{M}_0, \mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4$ when using forward selection?
- 2 What subset of variables is returned using forward subset selection?

Group work

↓

	Training MSE ($\times 10^7$)	k-fold CV Testing Error
Null model	8.76	10.08
X1	8.63	9.98
X2	7.42	8.01
X3	8.16	8.3
X4	8.33	9.06
X1, X2	4.33	7.47
X1, X3	5.82	5.22
X1, X4	3.17	4.23
X2, X3	4.07	3.78
X2, X4	3.31	4.01
X3, X4	3.06	4.16
X1, X2, X3	3.08	5.49
X1, X2, X4	3.55	4.02
<u>X1, X3, X4</u>	2.97	4.23
X2, X3, X4	2.98	3.17
X1, X2, X3, X4	2.16	4.39

$\rightarrow \mu_0$
 $\rightarrow \mu_1$
 $\bullet \emptyset$
 $\rightarrow \mu_2 \checkmark$
 $\rightarrow \mu_3$
 $\leftarrow \mu_4$

$\bullet X_1$
 $\bullet X_2$
 $\bullet X_3$
 $\bullet X_4$

$\bullet X_1 X_2$
 $\bullet X_1 X_3$
 $\bullet X_1 X_4$
 $\bullet X_2 X_3$
 $\bullet X_2 X_4$
 $\bullet X_3 X_4$

$\bullet X_1 X_2 X_3$
 $\bullet X_1 X_2 X_4$
 $\bullet X_1 X_3 X_4$
 $\bullet X_2 X_3 X_4$

$\bullet X_1 X_2 X_3 X_4$

$$Y = \beta_0 + \beta_1 X_3 + \beta_2 X_4 + \epsilon$$

Pros and Cons of Backward Selection

$p > n$
(i) computational feasible Cons
(ii) local solution

Forward Vs Backward $p > n$
✓ ✗

p^2

Alternatives for Approximating Test Error

Algorithm 6.3 *Backward stepwise selection*

1. Let \mathcal{M}_p denote 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 of the predictors in \mathcal{M}_k , for a total of $k - 1$ predictors.
 - (b) Choose the *best* among these k models, and call it \mathcal{M}_{k-1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

- **CV** is good but very time consuming.

- We can indirectly estimate test error by making an **adjustment** to the training error to account for the bias due overfitting

Estimating Test Error

Algorithm 6.1 Best subset selection

1. Let \mathcal{M}_0 denote the *null model*, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 1, 2, \dots, p$:
 - (a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - (b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here *best* is defined as having the smallest RSS, or equivalently largest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .

Algorithm 6.2 Forward stepwise selection

1. Let \mathcal{M}_0 denote the *null model*, which contains no predictors.
2. For $k = 0, \dots, p - 1$:
 - (a) Consider all $p - k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the *best* among these $p - k$ models, and call it \mathcal{M}_{k+1} . Here *best* is defined as having smallest RSS or highest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .

Algorithm 6.3 Backward stepwise selection

1. Let \mathcal{M}_p denote 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 of the predictors in \mathcal{M}_k , for a total of $k - 1$ predictors.
 - (b) Choose the *best* among these k models, and call it \mathcal{M}_k . Here *best* is defined as having smallest RSS or highest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .

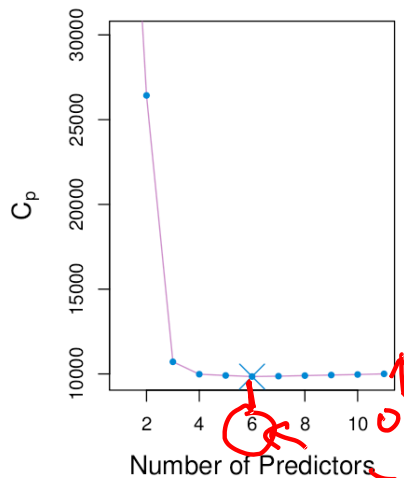
$$C_p = AIC$$

- Mallow's C_p :

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

training ↓
penalized complexity

where d is the total # of parameters used and $\hat{\sigma}^2$ is an estimate of the variance of the error ϵ associated with each response measurement.



$$2\hat{\sigma}^2 \cdot [d]$$

$$Y = f(x) + \begin{pmatrix} \epsilon \\ 1 \end{pmatrix}$$

$$\epsilon \sim N(0, \sigma^2)$$

- The *AIC* criterion is defined for a large class of models fit by maximum likelihood:

$$\text{AIC} = \underbrace{-2 \log L}_{\sim R^2} + \underbrace{2 \cdot d}_{\text{penalty}}$$

where L is the maximized value of the likelihood function for the estimated model.

$M_1 \rightarrow AIC$

$M_2 \rightarrow BIC$

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

$n = 100$

$d \rightarrow$ model complexity

$\log(n) \hat{\sigma}^2$
 $\approx 2 \cdot \hat{\sigma}^2$

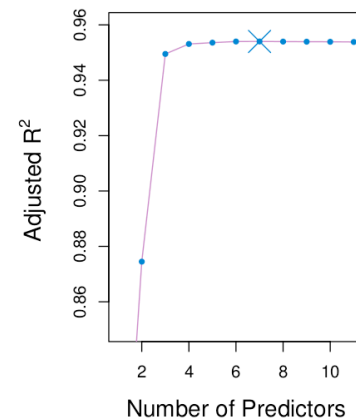
- Like C_p , the BIC will tend to take on a small value for a model with a low test error, and so generally we select the model that has the lowest BIC value.
- Notice that BIC replaces the $2d\hat{\sigma}^2$ used by C_p with a $\log(n)d\hat{\sigma}^2$ term, where n is the number of observations.
- Since $\log n > 2$ for any $n > 7$, the BIC statistic generally places a heavier penalty on models with many variables, and hence results in the selection of smaller models than

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

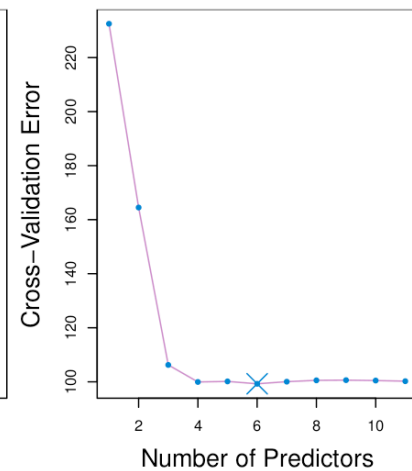
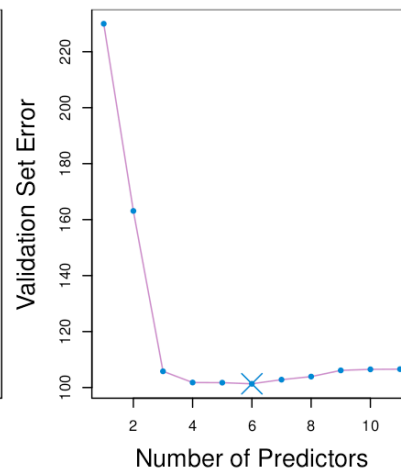
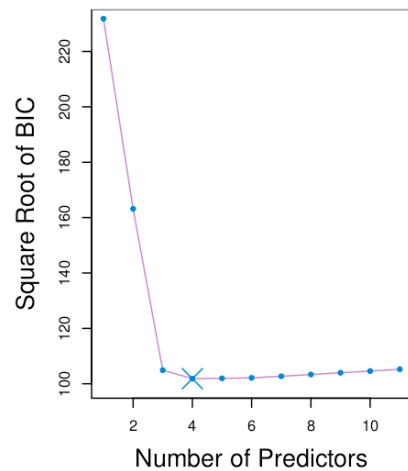
$$\text{Adjusted } R^2 = 1 - \frac{\text{RSS}/(n - d - 1)}{\text{TSS}/(n - 1)}.$$

where TSS is the total sum of squares.

- Unlike C_p , AIC, and BIC, for which a *small* value indicates a model with a low test error, a *large* value of adjusted R^2 indicates a model with a small test error.



Comparisons



Bonus Quiz 11

We have two training points ($x_1 = 1, x_2 = 1, y = 1$) and ($x_1 = 2, x_2 = 3, y = 2$). We want to fit a linear model to minimize the MSE

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

What is the minimum MSE? What is the corresponding model?

$$\begin{cases} 1 = \beta_0 + \beta_1 \cdot 1 + \beta_2 \cdot 1 \\ 2 = \beta_0 + 2 \cdot \beta_1 + 3 \cdot \beta_2 \end{cases}$$