

DSO530 Statistical Learning Methods

Lecture 6 part I: Linear Model Selection

Dr. Xin Tong

Department of Data Sciences and Operations

Marshall School of Business

University of Southern California

xint@marshall.usc.edu

Outline

When $n > 100$ only then least squares is useful.

sample size \nearrow When sample size limited, you cannot train effectively, more often in scientific fields.

of predictors \nearrow random error

- Previously, we studied the linear regression model

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, \sigma^2).$$

- One typically fits this model using least squares (go back to Lecture 2. if you do not recall the concept). *↳ minimizing sum of squared residuals.*
- We might want to use another fitting procedure instead of least squares on all p features to yield better prediction accuracy and model interpretability.
- In this lecture, we discuss one alternative approach: subset selection.
- Subset selection methods include best subset selection, stepwise selection.
- This set of slides only covers the subset selection methods ideas. A Python tutorial will be released later.

Best Subset Selection

$\{x_1, \dots, x_p\}$ How many subsets are there?

$$M_0 = \{\emptyset\} = 1 = {}^P C_0$$

$$M_1 = \{x_1, \emptyset, \dots, \{x_p\}\} = {}^P C_1$$

$$M_2 = \{x_1, x_2\}, \dots, \{x_{p-1}, x_p\} = {}^P C_2$$

:

:

:

$$M_p = \{x_1, \dots, x_p\} = 1 = {}^P C_p$$

$${}^P C + {}^P C_1 + \dots + {}^P C_p = 2^P$$

This is just enumeration, not selection.
1 option: choose best set at each level w/ highest R^2 or smallest RSS
Since models are of same size, in sample R^2 can be used.

a) $P!$ b) 2^P c) Neither
 \downarrow
 $2 \times 2 \dots \times 2$

e.g. All subsets: $\{\emptyset, x_1, x_2\}$
 $\{\emptyset\}, \{x_1\}, \{x_2\}, \{x_1, x_2\}$

If $P=2^0$

$2^P = 2^{2^0}$ subsets

$${n \choose m} = \frac{n!}{m!(n-m)!}$$

$${}^P C_p = {}^P C_{p-r}$$

Best Subset Selection

- To perform *best subset selection*, we fit a separate least squares regression for each possible combination of the p predictors
- Potential problem: selecting the best model from among too many possibilities considered by best subset selection is not trivial.
- Best subset selection is *usually* implemented by:

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 .

If we stop the best subset at k^* ($\leq p$) then how many models?

- a) 2^k b) $k^*!$ c) Neither A or B

A few questions on best subset selection

$$p=2 \quad k^* = 1 \quad \{x_1, x_2\} \quad \text{Subsets : } \{\}, \{x_1\}, \{x_2\}$$

$$\text{Ans}] P_{C_0} + P_{C_1} + \dots + P_{C_{k^*}}$$

models of same size

- In step 2b), why can we just look at the RSS on training error?
- In step 3), why don't we just look at the RSS on training error?
- How many models do we search through in best subset selection?

different size

$$P = 20$$

$$\begin{matrix} 20 \\ 2 \\ 2 \end{matrix} > 1,000,000$$

$$P^2 = 400$$

Forward Stepwise Selection

\Rightarrow A subset selection can be slow sometimes

date subset
depends on previous
level.

this is
where forward

Selection & subset method
differ. e.g. x_1 is always in
there is stepwise in:

$$\mathcal{M}_0: \{\emptyset\}$$

$$\mathcal{M}_1: \{\{x_1\}\}, \{\{x_2\}\}, \dots, \{\{x_p\}\}$$

$$\mathcal{M}_2: \{\{x_1, x_2\}\}, \dots, \{\{x_1, x_p\}\}$$

$p-1$ models

$p-2$

Algorithm 6.2 Forward stepwise selection

- Let \mathcal{M}_0 denote the *null model*, which contains no predictors.
- For $k = 0, \dots, p-1$: \therefore total # models: $1 + p + (p-1) + \dots + 1 = 1 + \frac{p(p+1)}{2} = \mathcal{O}(p^2)$
 - Consider all $p-k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - 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 .
- 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 .

$$P = 20 \\ 2^P > 1,000,000 \\ P^2 = 400.$$

potential problem: some best subsets may be left out
Q. If we stop fss at $k^*(p)$ then how many models?

Backward Stepwise Selection

$$\mathcal{M}_p : \{x_1, \dots, x_p\}^{\text{total } 1}$$

$$\mathcal{M}_{p-1} : \{x_1, x_2, \dots, x_p\}, \{x_1, x_3, \dots, x_p\}, \dots, \{x_1, \dots, x_p\}$$

$$\mathcal{M}_{p-2} : \{x_2, \dots, x_p\}, \{x_1, x_3, \dots, x_p\}, \dots, \{x_1, x_2, \dots, x_p\}$$

$$\vdots$$

$$\mathcal{M}_0 = \{\}^{\text{total } 1}$$

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 .
-

- In either forward or backward selection, we search through $1 + p(p + 1)/2$ models. This is a huge saving compared with best subset selection
- However, forward stepwise selection and backward stepwise selection might miss the optimal subset of features. This is a price we have to pay for computational advantages.

Choose the model in Step 3)

In best subset, forward and backward selection algorithms, the step 3)'s are the same. There are essentially two ideas in this step

- directly estimate the test error, using either a validation set approach or a cross-validation approach, as discussed in Lecture 5.
- indirectly estimate test error by making an adjustment to the training error (e.g., adjusted R^2 , AIC, BIC and C_p). *# of predictors available in select mod.*
-

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

i.e. $\{x_1, x_2, x_3\}$ has $d = 3$

Q: how is adjusted R^2 different from R^2 ?

C_p , AIC and BIC (for linear regression)

- Mallow's C_p

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

- Akaike information criterion (AIC)

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

- Bayesian information criterion(BIC)

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

- Note that $\hat{\sigma}^2$ is the estimated variance of random error term ϵ using the full model (i.e., p predictors).
- The definitions of AIC and BIC ignored some constant

trade off b/w model fitting
& complexity

A few questions

B

- For linear regression, will C_p and AIC give the same ranking of models? Yes.
- When n is bigger than 7, $\log n > 2$. This means BIC penalizes larger models heavier compared to AIC. So which criterion encourages smaller models? JC

A few questions

- for $p > n$, backward X.
 - we can use best subset & forward w/ early stopping at k^* (2^n)
 - CV is best approach.

- In the definitions of adjusted R^2 , AIC, BIC and C_p , do you see the trade-offs between fitting on training data and model complexity?
- When $p > n$, estimating $\hat{\sigma}^2$ is a big problem. Then which method do you prefer for model selection?
- Among adjusted R^2 , AIC, C_p , BIC and cross-validation, which one is the most easily generalizable beyond least squares linear regression?