# Ch 6.1: Subset Selection

Prof. Mengsen Zhang

Michigan State University

:

Dept of Computational Mathematics, Science & Engineering

Mon, Oct 6th, 2025

#### Announcements

#### Last time

• k-fold CV for Classification

#### Covered in this lecture

- response to feedback (from Quiz 5)
- Subset selection
- Forward and Backward Selection

#### **Announcements:**

- HW #4 Due Sunday (10/12)
- HW #5 due very late, to avoid Fall break weekend. I will release it early so you can work on it to prep for the exam

#### CMSE381 F2025 Schedule: Schedule

	M	9/22	Project Day & Review		
	W	9/24	Midterm #1		
12	F	9/26	Leave one out CV	5.1.1, 5.1.2	
13	М	9/29	k-fold CV	5.1.3	
14	W	10/1	More k-fold CV	5.1.4-5	
15	F	10/3	k-fold CV for classification	5.1.5	
16	М	10/6	Subset selection	6.1	
17	W	10/8	Shrinkage: Ridge	6.2.1	
18	F	10/10	Shrinkage: Lasso	6.2.2	HW #4 Due
19	М	10/13	PCA	6.3	Sun 10/12
20	W	10/15	PCR	6.3	
	F	10/17	Review		
	М	10/20	Fall Break		
	W	10/22	Midterm #2		
21	F	10/24	Polynomial & Step Functions	7.1-7.2	HW #5 Due
22	М	10/27	Step Functions; Basis functions; Start Splines	7.2-7.4	Sun 10/28
23	W	10/29	Regression Splines	7.4	

2 / 25

Dr. Zhang (MSU-CMSE)

Mon, Oct 6th, 2025

### Feedback collected in Quiz 5

- what I can do:
  - ▶ I will add overview slide at the beginning of each lecture. This will become a study guide for the exam.
  - I will provide you a practice exam with matched difficulty. Those questions will not be in the exam.
  - ► Code portfolios as bonus assignments. I will provide a template.
- what you can do:
  - Take better notes
  - Read the textbook (read before class if you have trouble following the lecture)
  - Ask more questions (in class, on slack, office hours/helproom) and respond to the review class "burning questions" survey

Or. Zhang (MSU-CMSE) Mon, Oct 6th, 2025

# What should you learn from this lecture?

- Why shouldn't you have as many predictors in your model as possible?
- What are the three basic methods could be used to select the appropriate predictors (feature selection)?
- What are the steps/procedures in the algorithm for each of these methods?
- How to implement these procedures by hand given the training and CV test error for each combo of predictors?
- What are the pros and cons for choosing each of these feature selection methods? When you can or cannot use them?
- How do you implement the feature selection algorithm in Python?

Dr. Zhang (MSU-CMSE)

Mon, Oct 6th, 2025

### Section 1

Previously on linear regression ...

rr. Zhang (MSU-CMSE) Mon, Oct 6th, 2025

# The problem of many features (p) relative to samples (n)

Up to now, we've focused on standard linear model:  $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon$  and done least squares estimation.

**Prediction accuracy** 

# The problem of many features (p) relative to samples (n)

Up to now, we've focused on standard linear model:  $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon$  and done least squares estimation.

### Model Interpretability

### Section 2

### Best Subset Selection

# Go through each combo of variables exhaustively (exhausting?)

All subsets of 4 variables ( $2^4 = 16$ )

X<sub>1</sub>

• X2

X<sub>3</sub>

X<sub>4</sub>

- $\bullet$   $X_1$   $X_3$ X<sub>1</sub> X<sub>4</sub>
- $\bullet$   $X_2$   $X_3$
- $\bullet$   $X_2$   $X_4$
- $\bullet$   $X_3$   $X_4$

- $\bullet$   $X_1$   $X_2$   $X_3$
- X<sub>1</sub> X<sub>2</sub> X<sub>4</sub>
- X<sub>1</sub> X<sub>3</sub> X<sub>4</sub>
  - $\bullet$   $X_2$   $X_3$   $X_4$

• Ø

• X<sub>1</sub> X<sub>2</sub> X<sub>3</sub> X<sub>4</sub>

# One way of breaking this up

#### 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  $\mathcal{M}_0, \ldots, \mathcal{M}_p$  using cross-validated prediction error,  $C_p$  (AIC), BIC, or adjusted  $R^2$ .

# Calculate by hand

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 MSE value computed for the model learned using each possible subset of variables.

	Training MSE (x10^7)	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

- 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?
- What subset of variables is returned using best subset selection?

# Extra work space if it helps

	Training MSE (x10^7)	k-fold CV Testing Erro
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

Ŋ

- X<sub>1</sub> X<sub>2</sub>
  X<sub>1</sub> X<sub>3</sub>
- $\bullet$   $X_2$   $\bullet$   $X_1$   $X_4$
- X<sub>3</sub> X<sub>2</sub> X<sub>3</sub>
- X<sub>4</sub> X<sub>2</sub> X<sub>4</sub>
  - X<sub>3</sub> X<sub>4</sub>

- X<sub>1</sub> X<sub>2</sub> X<sub>3</sub>
- X<sub>1</sub> X<sub>2</sub> X<sub>4</sub>
- $X_1 X_3 X_4$
- $X_2 X_3 X_4$

O-1 611 2025

12 / 25

• X<sub>1</sub> X<sub>2</sub> X<sub>3</sub> X<sub>4</sub>

# Code to do this

Or. Zhang (MSU-CMSE) Mon, Oct 6th, 2025

### Section 3

### Forward Selection

What's the problem with best subset selection?

Or. Zhang (MSU-CMSE) Mon, Oct 6th, 2025

# Forward Stepwise Selection

#### Algorithm 6.2 Forward stepwise selection

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

# An example for Forward Stepwise Selection

• Ø

- X<sub>1</sub>
- X<sub>2</sub>
- X<sub>3</sub>
- X<sub>4</sub>

- $\bullet$   $X_1$   $X_2$
- $X_1 X_3$
- X<sub>1</sub> X<sub>4</sub>
- $\bullet$   $X_2$   $X_3$
- $\bullet$   $X_2$   $X_4$
- $X_3 X_4$

- X<sub>1</sub> X<sub>2</sub> X<sub>3</sub>
- X<sub>1</sub> X<sub>2</sub> X<sub>4</sub>
- X<sub>1</sub> X<sub>3</sub> X<sub>4</sub>
- $\bullet$   $X_2$   $X_3$   $X_4$

X<sub>1</sub> X<sub>2</sub> X<sub>3</sub> X<sub>4</sub>

# Group work: by hand same example with forward example

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 (x10^7)	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

- 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?
- What subset of variables is returned using forward subset selection?

# Extra work space if it helps

	Training MSE (x10^7)	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

• Ø

- X<sub>1</sub> X<sub>2</sub>
  X<sub>1</sub> X<sub>3</sub>
  X<sub>2</sub> X<sub>1</sub> X<sub>4</sub>
  X<sub>3</sub> X<sub>2</sub> X<sub>3</sub>
  X<sub>4</sub> X<sub>2</sub> X<sub>4</sub>
  X<sub>3</sub> X<sub>4</sub>
- X<sub>1</sub> X<sub>2</sub> X<sub>3</sub>
  X<sub>1</sub> X<sub>2</sub> X<sub>4</sub>
- X<sub>1</sub> X<sub>3</sub> X<sub>4</sub>
- X<sub>2</sub> X<sub>3</sub> X<sub>4</sub>

• X<sub>1</sub> X<sub>2</sub> X<sub>3</sub> X<sub>4</sub>

# Pros and Cons of Forward Stepwise

Pros: Cons:

### Section 4

### **Backward Selection**

# Backward stepwise selection

#### 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  $\mathcal{M}_0, \ldots, \mathcal{M}_p$  using cross-validated prediction error,  $C_p$  (AIC), BIC, or adjusted  $R^2$ .

# Pros and Cons of Backward Stepwise

Pros: Cons:

# TL;DR

#### 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  $R^2$ .
- Select a single best model from among M<sub>0</sub>,...,M<sub>p</sub> using crossvalidated prediction error, C<sub>p</sub> (AIC), BIC, or adjusted R<sup>2</sup>.

- Modify step 2 with forward or backward selection
- Choose best model in step 3 using one of our adjusted training scores or CV

24 / 25

Test your understanding: PollEv

Or. Zhang (MSU-CMSE) Mon, Oct 6th, 2025

### Next time

#### CMSE381\_F2025\_Schedule : Schedule

	M	9/22	Project Day & Review		
	W	9/24	Midterm #1		
12	F	9/26	Leave one out CV	5.1.1, 5.1.2	
13	М	9/29	k-fold CV	5.1.3	
14	W	10/1	More k-fold CV	5.1.4-5	
15	F	10/3	k-fold CV for classification	5.1.5	
16	М	10/6	Subset selection	6.1	
17	W	10/8	Shrinkage: Ridge	6.2.1	
18	F	10/10	Shrinkage: Lasso	6.2.2	HW #4 Due
19	М	10/13	PCA	6.3	Sun 10/12
20	W	10/15	PCR	6.3	
	F	10/17	Review		
	М	10/20	Fall Break		
	W	10/22	Midterm #2		
21	F	10/24	Polynomial & Step Functions	7.1-7.2	HW #5 Due
22	М	10/27	Step Functions; Basis functions; Start Splines	7.2-7.4	Sun 10/28
23	W	10/29	Regression Splines	7.4	