Ch 6.1: Subset Selection

Prof. Guanqun Cao

Michigan State University

:

Dept of Computational Mathematics, Science & Engineering

Mon, Oct 6, 2025

Dr. Cao (MSU-CMSE) Mon, Oct 6, 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)

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	

Dr. Cao (MSU-CMSE) Mon, Oct 6, 2025 2 / 23

Section 1

Previously on linear regression ...

Pr. Cao (MSU-CMSE) Mon, Oct 6, 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

Dr. Cao (MSU-CMSE) Mon, Oct 6, 2025

Go through each combo of variables exhaustively (exhausting?)

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

•
$$X_1 X_2 X_3$$

Dr. Cao (MSU-CMSE)

• Ø

Mon, Oct 6, 2025

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 .

Dr. Cao (MSU-CMSE) Mon, Oct 6, 2025

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?

Dr. Cao (MSU-CMSE) Mon, Oct 6, 2025 9 / 23

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

• 0

- X₁ X₂
- X₁ X₁ X₃
- $\bullet X_2 \bullet X_1 X_4$
- X₃ X₂ X₃
- X₄ X₂ X₄
 - X₃ X₄

- $X_1 X_2 X_3$
- X₁ X₂ X₄
- $\bullet \ X_1 \ X_3 \ X_4$
- $X_2 X_3 X_4$

• X₁ X₂ X₃ X₄

Code to do this

Section 3

Forward Selection

Dr. Cao (MSU-CMSE) Mon, Oct 6, 2025

What's the problem with best subset selection?

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

Dr. Cao (MSU-CMSE) Mon, Oct 6, 2025 14/23

An example for Forward Stepwise Selection

• Ø

- X₁
- X₂
- X₃
- X₄

- \bullet X_1 X_2
- $X_1 X_3$
- X₁ X₄
- \bullet X_2 X_3
- \bullet X_2 X_4
- $X_3 X_4$

- X₁ X₂ X₃
- X₁ X₂ X₄
- X₁ X₃ X₄
- \bullet X_2 X_3 X_4

X₁ X₂ X₃ X₄

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	k-fold CV
	MSE (x10^7)	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?

Dr. Cao (MSU-CMSE) Mon, Oct 6, 2025 16 / 23

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

X₁ X₂
X₁ X₃

• X₃

- X₂ X₁ X₄
 - X₂ X₃
- X₄ X₂ X₄
 - X₃ X₄

- X₁ X₂ X₃
- X₁ X₂ X₄
- X₁ X₃ X₄
- X₂ X₃ X₄

• X₁ X₂ X₃ X₄

Pros and Cons of Forward Stepwise

Pros: Cons:

Section 4

Backward Selection

Dr. Cao (MSU-CMSE) Mon, Oct 6, 2025

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 .

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

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 M₀,..., M_p using crossvalidated prediction error, C_p (AIC), BIC, or adjusted R².

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

22 / 23

Dr. Cao (MSU-CMSE) Mon, Oct 6, 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	

Dr. Cao (MSU-CMSE)

Mon, Oct 6, 2025