Ch 3.2: Multiple Linear Regression Lecture 6 - CMSE 381

Prof. Guanqun Cao

Michigan State University

::

Dept of Computational Mathematics, Science & Engineering

Mon, Sep 8, 2025

Announcements

Last time:

• 3.1 (Simple) linear regression

Announcements:

2/24

 Homework #2 Due Sunday on Crowdmark

Or. Cao (MSU-CMSE) Mon, Sep 8, 2025

Covered in this lecture

- Multiple linear regression
- Hypothesis test in that case
- Forward and Backward Selection

Dr. Cao (MSU-CMSE) Mon, Sep 8, 2025

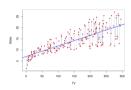
Section 1

Review from last time

Or. Cao (MSU-CMSE)

Mon, Sep 8, 2025

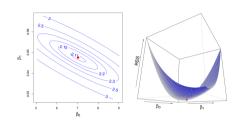
Linear Regression with One Variable



• Predict Y on a single variable X

$$Y \approx \beta_0 + \beta_1 X$$

- Find good guesses for $\hat{\beta}_0$, $\hat{\beta}_1$.
- $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$
- $e_i = y_i \hat{y}_i$ is the *i*th residual
- residual sum of squares RSS = $\sum_{i} e_{i}^{2}$



• RSS is minimized at *least* squares coefficient estimates

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \overline{x})(y_i - \overline{y})}{\sum_{i=1}^n (x_i - \overline{x})^2}$$
$$\hat{\beta}_0 = \overline{y} - \hat{\beta}_1 \overline{x}$$

Evaluating the model

- Linear regression is unbiased
- Variance of linear regression estimates:

$$SE(\hat{\beta}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\overline{x}^2}{\sum_{i=1}^n (x_i - \overline{x})^2} \right]$$
$$SE(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \overline{x})^2}$$

where
$$\sigma^2 = \operatorname{Var}(\varepsilon)$$

• Estimate σ : $\hat{\sigma}^2 = \frac{RSS}{n-2}$

• The 95% confidence interval for β_1 approximately takes the form

$$\hat{\beta}_1 \pm 2 \cdot \text{SE}(\hat{\beta}_1)$$

• Hypothesis test:

$$H_0$$
: $\beta_1 = 0$
 H_a : $\beta_1 \neq 0$

▶ Test statistic $t = \frac{\hat{eta}_1 - 0}{\operatorname{SE}(\hat{eta}_1)}$

Assessing the accuracy of the model

Residual standard error (RSE):

$$RSE = \sqrt{\frac{1}{n-2}RSS}$$

R squared:

$$R^{2} = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$$
$$TSS = \sum_{i} (y_{i} - \overline{y})^{2}$$

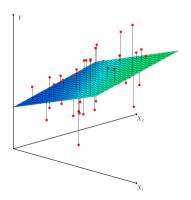
Section 2

Multiple Linear Regression

Or. Cao (MSU-CMSE) Mon, Sep 8, 2025

Setup

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \varepsilon$$



Estimation and Prediction

Given estimates $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \cdots, \hat{\beta}_p$, prediction is

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$$

Minimize the sum of squares

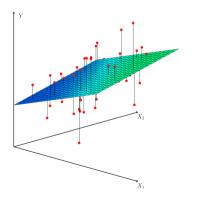
$$RSS = \sum_{i} (y_{i} - \hat{y}_{i})^{2}$$

$$= \sum_{i} (y_{i} - \hat{\beta}_{0} - \hat{\beta}_{1}x_{i1} - \dots - \hat{\beta}_{p}x_{ip})^{2}$$

Coefficients are closed form but UGLY

Advertising data set example

Sales =
$$\beta_0 + \beta_1 \cdot TV + \beta_2 \cdot radio + \beta_3 \cdot newspaper$$



	Coefficient
Intercept	2.939
TV	0.046
radio	0.189
newspaper	-0.001

11 / 24

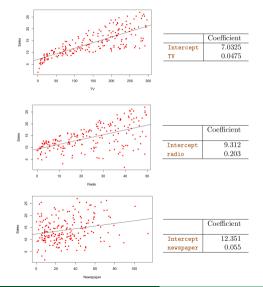
Or. Cao (MSU-CMSE) Mon, Sep 8, 2025

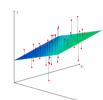
Interpretation of coefficients

$$\mathtt{Sales} = \beta_0 + \beta_1 \cdot \mathtt{TV} + \beta_2 \cdot \mathtt{radio} + \beta_3 \cdot \mathtt{newspaper}$$

	Coefficient
Intercept	2.939
TV	0.046
radio	0.189
newspaper	-0.001

Single regression vs multi-regression





	Coefficient
Intercept	2.939
TV	0.046
radio	0.189
newspaper	-0.001

Correlation matrix

	TV	radio	newspaper	sales
TV	1.0000	0.0548	0.0567	0.7822
radio		1.0000	0.3541	0.5762
newspaper			1.0000	0.2283
sales				1.0000

Coding group work

Run the section titled "Multiple Linear Regression"

Dr. Cao (MSU-CMSE) Mon, Sep 8, 2025

Section 3

Ch 3.2.2: Questions to ask of your regression

Pr. Cao (MSU-CMSE) Mon, Sep 8, 2025

Question 1

Is at least one of the predictors X_1, \dots, X_p useful in predicting the response?

Q1: Hypothesis test

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_p = 0$$

 H_a : At least one β_j is non-zero

F-statistic:

$$F = \frac{(\mathit{TSS} - \mathit{RSS})/p}{\mathit{RSS}/(n-p-1)} \sim F_{p,n-p-1}$$

The F-statistic for the hypothesis test

$$F = \frac{(TSS - RSS)/p}{RSS/(n-p-1)} \sim F_{p,n-p-1}$$

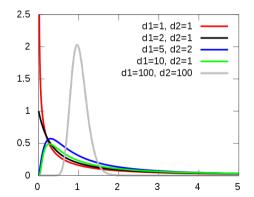


Image from wikipedia, By IkamusumeFan - Own work, CC BY-SA 4.0,

Or. Cao (MSU-CMSE) Mon, Sep 8, 2025

Do Q1 section in jupyter notebook

Q2

Do all the predictors help to explain Y, or is only a subset of the predictors useful?

Or. Cao (MSU-CMSE) Mon, Sep 8, 2025

Q2: A first idea

Great, you know at least one variable is important, so which is it?....

Dr. Cao (MSU-CMSE) Mon, Sep 8, 2025

Do Q2 section in jupyter notebook

Why is this a bad idea?

Next time

CMSE381_F2025_Schedule : Schedule

Lec #	Date		Topic	Reading	HW	
1	M	8/25	Intro / Python Review	1		
2	W	8/27	What is statistical learning	2.1		
3	F	8.29	Assessing Model Accuracy	2.2.1, 2.2.2		
	M	9/1	Labor Day - No Class			
4	W	9/3	Linear Regression	3.1		
5	F	9/5	More Linear Regression	3.1	HW #1 Due	
6	М	9/8	Multi-linear Regression	3.2	Sun 9/7	
7	W	9/10	Probably More Linear Regression	3.3		
8	F	9/12	Last of the Linear Regression		HW #2 Due	
9	М	9/15	Intro to classification, Bayes classifier, KNN classifier	2.2.3	Sun 9/14	
10	W	9/17	Logistic Regression	4.1, 4.2, 4.3.1-3		
11	F	9/19	Multiple Logistic Regression / Multinomial Logistic Regression	4.3.4-5	HW #3 Due Sun 9/21	
	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		
40		0/00	1. 6-1-1-007	F 4 0		

Dr. Cao (MSU-CMSE) Mon, Sep 8, 2025