# Ch 3.3: Even More Linear Regression Lecture 7 - CMSE 381

Prof. Guanqun Cao

Michigan State University

::

Dept of Computational Mathematics, Science & Engineering

Wed, Sep 10, 2025

### Announcements

### Last time:

• 3.2 Multiple Linear Regression

### **Announcements:**

2/29

- HW #2 Due Sunday!
- Office hours

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

## Covered in this lecture

- RSE, R<sup>2</sup>
- Confidence intervals and prediction intervals
- Qualitative predictors

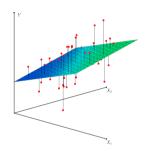
Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

### Section 1

Continued: Questions to ask of your model

Pr. Cao (MSU-CMSE) Wed, Sep 10, 2025

# Linear Regression with Multiple Variables



 $\bullet$  Predict Y on a multiple variables X

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

- Find good guesses for  $\hat{\beta}_0$ ,  $\hat{\beta}_1$ ,  $\cdots$ .
- $\bullet \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \dots + \hat{\beta}_p x_p$

- $e_i = y_i \hat{y}_i$  is the *i*th residual
- RSS =  $\sum_i e_i^2$
- RSS is minimized at least squares coefficient estimates

5/29

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

# Review: Questions to ask of your model

- Is at least one of the predictors  $X_1, \dots, X_p$  useful in predicting the response?
- ② Do all the predictors help to explain Y, or is only a subset of the predictors useful?

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

Q3

How well does the model fit the data?

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

# Assessing the accuracy of the module

Almost the same as before

## Residual standard error (RSE):

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

### R squared:

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

# $R^2$ on Advertising data

- Just TV:  $R^2 = 0.61$
- Just TV and radio:  $R^2 = 0.89719$
- All three variables:  $R^2 = 0.8972$

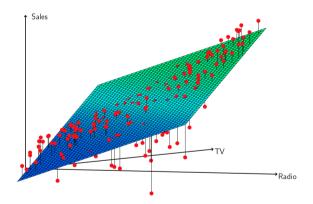
Or. Cao (MSU-CMSE) Wed, Sep 10, 2025

# RSE on Advertising Data

- Just TV: *RSE* = 3.26
- Just TV and radio: RSE = 1.681
- All three variables: RSE = 1.686

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

# If all else fails, look at the data



Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

### Q4

Given a set of predictor values, what response value should we predict, and how accurate is our prediction?

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

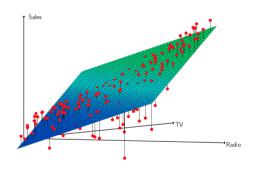
# Q4: Making predictions

Given estimates  $\hat{\beta}_0, \dots, \hat{\beta}_p$  for  $\beta_0, \dots, \beta_p$  Least squares plane:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p$$

estimate for the true population regression plane

$$f(X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$



13 / 29

Or. Cao (MSU-CMSE) Wed, Sep 10, 2025

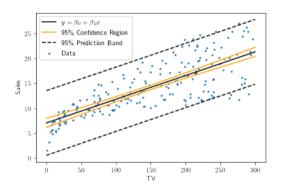
### Confidence vs Prediction Model

### **Confidence Interval**

The range likely to contain the population parameter (mean, standard deviation) of interest.

### **Prediction Interval**

The range that likely contains the value of the dependent variable for a single new observation given specific values of the independent variables.



Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025 14/29

# Specific to the Advertising Data

**Confidence interval**: quantify the uncertainty surrounding the <u>average</u> sales over a large number of cities.

### Advertising example:

If \$100K is spent on TV, and \$20K on radio, in each of *n* cities

95% CI for <u>average</u> sales: [10,985, 11,528].

**Prediction Interval:** quantify the uncertainty in sales for a particular city.

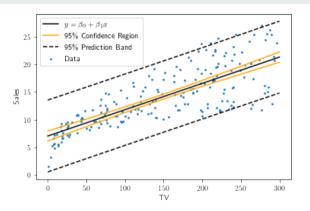
### Advertising example:

Given that \$100,000 is spent on TV advertising and \$20,000 is spent on radio advertising in **Gotham City** 

95% prediction interval for Gotham: [7,930, 14,580].

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025 15/29

# Comparing the two



16 / 29

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

Go take a look at the code under Q4

# Review: Questions to ask of your model

- Is at least one of the predictors  $X_1, \dots, X_p$  useful in predicting the response?
- Oo all the predictors help to explain Y, or is only a subset of the predictors useful?
- 4 How well does the model fit the data?
- Given a set of predictor values, what response value should we predict, and how accurate is our prediction?

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

# Section 2

# **Qualitative Predictors**

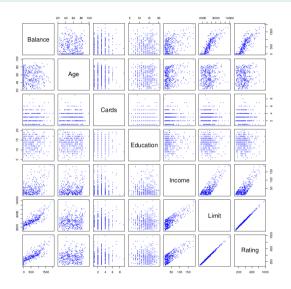
Pr. Cao (MSU-CMSE) Wed, Sep 10, 2025

Reminder: Qualitative vs Quantitative predictors

Quantitative:

Qualitative/Categorical:

### New data set! Credit card balance



own: house ownership

student: student status

• status: marital status

• region: East, West, or South

21 / 29

. Cao (MSU-CMSE) Wed, Sep 10, 2025

### What if....

- ... your variables aren't quantitative?
- Home ownership
- Student status
- Major
- Gender
- Ethnicity
- Country of origin

## Example

Investigate differences in credit card balance between people who own a house and those who don't, ignoring the other variables.

22 / 29

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

# One-hot encoding

### Create a new variable

$$x_i = \begin{cases} 1 & \text{if } i \text{th person is a student} \\ 0 & \text{if } i \text{th person is not a student} \end{cases}$$

### Model:

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_i + \varepsilon_i \\ &= \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & \text{if $i$th person is student} \\ \beta_0 + \varepsilon_i & \text{if $i$th person isn't} \end{cases} \end{aligned}$$

Dr. Cao (MSU-CMSE)

# Interpretation

coef		std err	t	P> t	[0.025	0.975]
Intercept	480.3694	23.434	20.499	0.000	434.300	526.439
Student[T.Yes]	396.4556	74.104	5.350	0.000	250.771	542.140

### Model:

$$y = 480.36 + 396.46 \cdot x_{student}$$

# Who cares about 0/1?

# Old version: 0/1

$$x_i = \begin{cases} 1 & \text{if } i \text{th person is a student} \\ 0 & \text{if } i \text{th person is not a student} \end{cases}$$

### Model:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$= \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & \text{if } i \text{th person is student} \\ \beta_0 + \varepsilon_i & \text{if } i \text{th person isn't} \end{cases}$$

### Alternative version: $\pm 1$

$$x_i = egin{cases} 1 & ext{if } i ext{th person is a student} \ -1 & ext{if } i ext{th person is not a student} \end{cases}$$

### Model:

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_i + \varepsilon_i \\ &= \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & \text{if $i$th person is student} \\ \beta_0 - \beta_1 + \varepsilon_i & \text{if $i$th person isn't} \end{cases} \end{aligned}$$

# Qualitiative Predictor with More than Two Levels

### Region:

# South West East

## Create spare dummy variables:

$$x_{i1} = \begin{cases} 1 & \text{if } i \text{th person from South} \\ 0 & \text{if } i \text{th person not from South} \end{cases}$$

$$x_{i2} = \begin{cases} 1 & \text{if } i \text{th person from West} \\ 0 & \text{if } i \text{th person not from West} \end{cases}$$

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i \\ &= \begin{cases} \beta_0 + \beta_1 x_{i1} + \varepsilon_i & \text{if $i$th person from South} \\ \beta_0 + \beta_2 x_{i2} + \varepsilon_i & \text{if $i$th person from West} \\ \beta_0 + \varepsilon_i & \text{if $i$th person from East} \end{cases} \end{aligned}$$

# More on multiple levels

	Coefficient	Std. error	t-statistic	<i>p</i> -value
Intercept	531.00	46.32	11.464	< 0.0001
region[South]	-18.69	65.02	-0.287	0.7740
region[West]	-12.50	56.68	-0.221	0.8260

27 / 29

Or. Cao (MSU-CMSE) Wed, Sep 10, 2025

Do code section on "Playing with multi-level variables"

Dr. Cao (MSU-CMSE) Wed, Sep 10, 2025

# 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	M	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. 4-1-1 007	E 4 0			

Or. Cao (MSU-CMSE) Wed, Sep 10, 2025