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

Prof. Elizabeth Munch

Michigan State University
::
Dept of Computational Mathematics, Science & Engineering

Fri, Sep 16, 2022

**Last time:**

- 3.2 Multiple Linear Regression

**Announcements:**
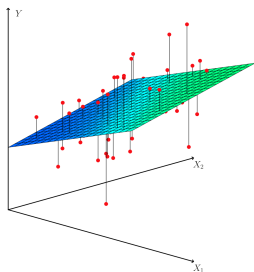
- Third homework posted on github
- Office hours

# Covered in this lecture

- Qualitative predictors
- Extending the linear model with interaction terms
- Hierarchy principle
- Polynomial regression

## Section 1

Review from last time

# Linear Regression with Multiple Variables



- 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 $\hat{\beta}_0, \hat{\beta}_1, \cdots$.
- $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \cdots + \hat{\beta}_p x_p$

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

## Questions to ask of your model

1. Is at least one of the predictors $X_1, \cdots, X_p$ useful in predicting the response?

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

3. How well does the model fit the data?

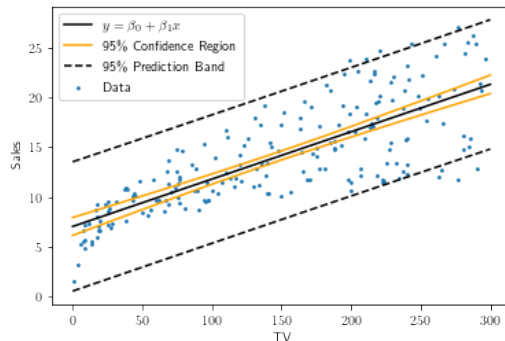4. Given a set of predictor values, what response value should we predict, and how accurate is our prediction?

# 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.

# Specific to the Advertising Data

**Confidence interval**: quantify the uncertainty surrounding the average 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 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**

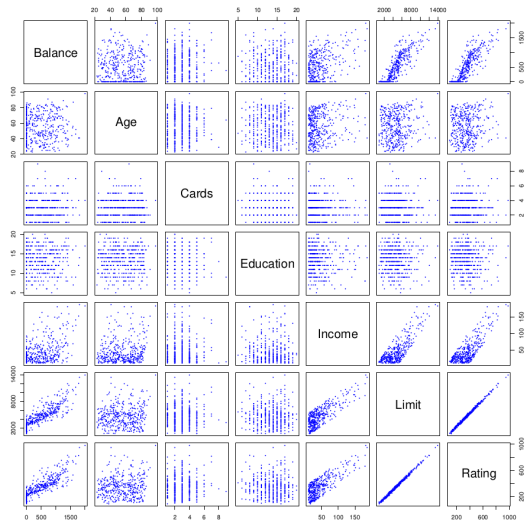95% prediction interval for Gotham:
[7,930, 14,580].

# Section 2

## Qualitative Predictors

# 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

## 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.

# 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:

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

# 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:

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

**Alternative version: $\pm 1$**

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

Model:

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

## Qualitiative Predictor with More than Two Levels

Region:

Create spare dummy variables:

| | $x_{i1}$ | $x_{i2}$ |
|---|---|---|
| South | | |
| West | | |
| East | | |

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

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

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

# More on multiple levels

| | Coefficient | Std. error | $t$-statistic | $p$-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 |

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

# Section 3

## Extending the linear model

$$\hat{Y}_{sales} = \beta_0 + \beta_1 \cdot X_{TV} + \beta_2 \cdot X_{radio} + \beta_3 \cdot X_{newspaper}$$

Assumed (implicitly) that the effect on sales by increasing one medium is independent of the others.

What if spending money on radio advertising increases the effectiveness of TV advertising? How do we model it?

## Interaction Term

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \varepsilon$$

$$Y_{sales} = \beta_0 + \beta_1 X_{TV} + \beta_2 X_{radio} + \beta_3 X_{radio} X_{TV} + \varepsilon$$
$$= \beta_0 + (\beta_1 + \beta_3 X_{radio}) X_{TV} + \beta_2 X_{radio} + \varepsilon$$

# Interaction term

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \varepsilon$$

|  | Coefficient | Std. error | $t$-statistic | $p$-value |
|---|---|---|---|---|
| Intercept | 6.7502 | 0.248 | 27.23 | < 0.0001 |
| TV | 0.0191 | 0.002 | 12.70 | < 0.0001 |
| radio | 0.0289 | 0.009 | 3.24 | 0.0014 |
| TV×radio | 0.0011 | 0.000 | 20.73 | < 0.0001 |

$$\begin{aligned} Y_{sales} &= \beta_0 + \beta_1 X_{TV} + \beta_2 X_{radio} + \beta_3 X_{radio} X_{TV} + \varepsilon \\ &= \beta_0 + (\beta_1 + \beta_3 X_{radio}) X_{TV} + \beta_2 X_{radio} + \varepsilon \end{aligned}$$

# Interpretation

|          | Coefficient | Std. error | $t$-statistic | $p$-value |
|----------|-------------|------------|---------------|-----------|
| Intercept | 6.7502     | 0.248      | 27.23         | < 0.0001  |
| TV        | 0.0191     | 0.002      | 12.70         | < 0.0001  |
| radio     | 0.0289     | 0.009      | 3.24          | 0.0014    |
| TV×radio  | 0.0011     | 0.000      | 20.73         | < 0.0001  |

Do the section on "Interaction Terms"

Sometimes $p$-value for interaction term is very small, but associated main effects are not.
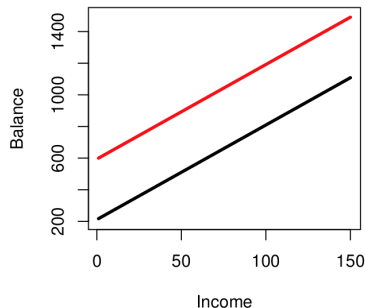
The hierarchy principle:

# Interaction term for qualitative variables
Without interaction term

For credit data set:
Predict balance using income (quantitative) and student (qualitative)

$$\texttt{balance}_i \approx \beta_0 + \beta_1 \cdot \texttt{income}_i + \begin{cases} \beta_2 & \text{if student} \\ 0 & \text{if not} \end{cases}$$

$$\approx \beta_1 \cdot \texttt{income}_i + \begin{cases} \beta_0 + \beta_2 & \text{if student} \\ \beta_0 & \text{if not} \end{cases}$$
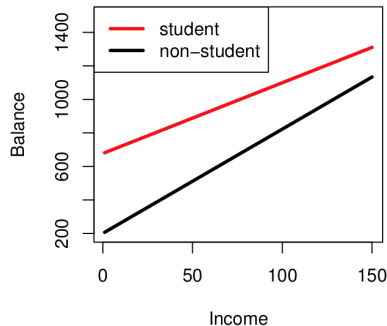
# Interaction term for qualitative variables
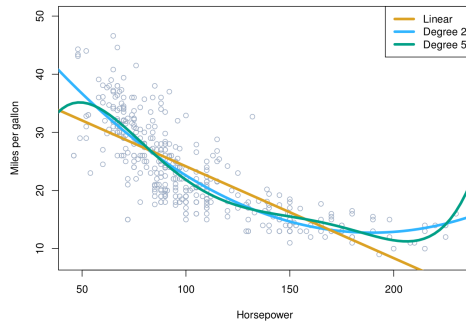With interaction term

For credit data set:
Predict balance using income (quantitative) and student (qualitative)

$$\text{balance}_i \approx \beta_0 + \beta_1 \cdot \text{income}_i + \begin{cases} \beta_2 + \beta_3 \cdot \text{income}_i & \text{if student} \\ 0 & \text{if not} \end{cases}$$

$$\approx \begin{cases} (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \cdot \text{income}_i & \text{if student} \\ \beta_0 + \beta_1 \cdot \text{income}_i & \text{if not} \end{cases}$$

# Nonlinear relationships

$$\text{mpg} = \beta_0 + \beta_1 \cdot \text{horsepower} + \beta_2 \cdot \text{horsepower}^2 + \varepsilon$$



|  | Coefficient | Std. error | $t$-statistic | $p$-value |
|---|---|---|---|---|
| Intercept | 56.9001 | 1.8004 | 31.6 | < 0.0001 |
| horsepower | −0.4662 | 0.0311 | −15.0 | < 0.0001 |
| horsepower$^2$ | 0.0012 | 0.0001 | 10.1 | < 0.0001 |

# Next time

| Lec # | | Date | Topic | Reading | Homeworks |
|---|---|---|---|---|---|
| 1 | W | Aug 31 | Intro / First day stuff / Python Review Pt 1 | 1 | |
| 2 | F | Sep 2 | What is statistical learning? | 2.1 | |
| | M | Sep 5 | No class - Labor day | | |
| 3 | W | Sep 7 | Assessing Model Accuracy | 2.2.1, 2.2.2 | HW #1 Due |
| 4 | F | Sep 9 | Linear Regression | 3.1 | |
| 5 | M | Sep 12 | More Linear Regression | 3.1/3.2 | |
| 6 | W | Sep 14 | Even more linear regression | 3.2.2 | HW #2 Due |
| 7 | F | Sep 16 | Probably more linear regression | 3.3 | |
| 8 | M | Sep 19 | Intro to classification, Logisitic Regression | 2.2.3, 4.1, 4.2, 4.3 | |
| 9 | W | Sep 21 | More logistic regression | | HW #3 Due |
| 10 | F | Sep 23 | Review | | |
| 11 | M | Sep 26 | *Midterm #1* | | |
| 12 | W | Sep 28 | [No class, Dr Munch out of town] | | |
| 13 | F | Sep 30 | [No class, Dr Munch out of town] | | |
| 14 | M | Oct 3 | Leave one out CV | 5.1.1, 5.1.2 | |
| 15 | W | Oct 5 | k-fold CV | 5.1.3 | |
| 16 | F | Oct 7 | More k-fold CV | 5.1.4 | |
| 17 | M | Oct 10 | CV for classification | 5.1.5 | HW #4 Due |
| 18 | W | Oct 12 | Resampling methods: Bootstrap | 5.2 | |
| 19 | F | Oct 14 | Subset selection | 6.1 | |
| 20 | M | Oct 17 | Shrinkage: Ridge | 6.2.1 | HW #5 Due |
| 21 | W | Oct 19 | Shrinkage: Lasso | 6.2.2 | |
| 22 | F | Oct 21 | Dimension Reduction | 6.3 | |