

Lecture 4: Simple Linear Regression Assumptions

Module 1: part 3

Spring 2024

Logistics

- Wrap up Module 1 today
- Module assessment due on Feb 11 11:59pm
- Module 2 will consider regression with multiple covariates
- Office hour locations: Daniel and Tathagata (Comstock 1187); Nayel in Surge B 159.

Recap

The equation for a line can be put into the following form

$$Y = b_0 + b_1X \quad (1)$$

Recap

The equation for a line can be put into the following form

$$Y = b_0 + b_1X \quad (1)$$

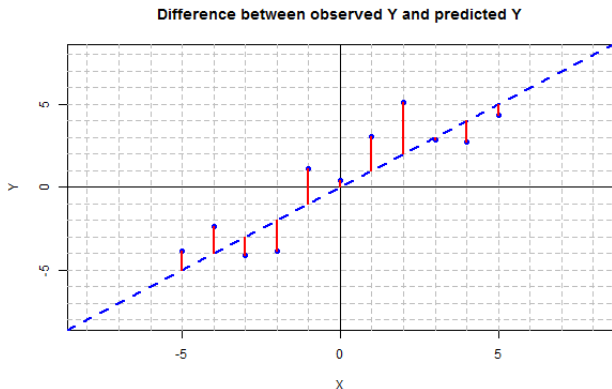
- X and Y are variables
- b_0 is the **Y-intercept**. It is the value of the Y coordinate when $X = 0$
- b_1 is the **slope**. It describes how Y changes as X changes.

Recap

Suppose we observe $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.

To select a “best line” we consider the difference between the predicted point and observed value of y_i and choose \hat{b}_0 and \hat{b}_1 to minimize the RSS:

$$RSS(\hat{b}_0, \hat{b}_1) = \sum_i (y_i - \hat{y}_i)^2 = \sum_i (y_i - (\hat{b}_0 + \hat{b}_1 x_i))^2 \quad (2)$$



Recap

Outliers:

- Points which have x values far from \bar{x} have high leverage
- Points which have high leverage may also have high influence; i.e., change the estimate when included/excluded
- When to include or exclude points with high influence?

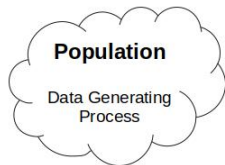
Linear Model

Interpretation

Let's take a step back and consider what we have calculated

- Still have “hat's” on \hat{b}_0 and \hat{b}_1 because they are calculated from the sample data
- We want to use the sampled data to infer something about the population

Sample data vs population distribution

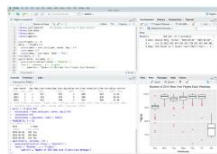


Data

id	name	age	sex	height	weight	hair
204	John	25	M	183	80.5	Black
304	John	25	M	183	80.5	Black
404	John	25	M	183	80.5	Black
504	John	25	M	183	80.5	Black
604	John	25	M	183	80.5	Black
704	John	25	M	183	80.5	Black
804	John	25	M	183	80.5	Black
904	John	25	M	183	80.5	Black
1004	John	25	M	183	80.5	Black
1104	John	25	M	183	80.5	Black
1204	John	25	M	183	80.5	Black
1304	John	25	M	183	80.5	Black
1404	John	25	M	183	80.5	Black
1504	John	25	M	183	80.5	Black
1604	John	25	M	183	80.5	Black
1704	John	25	M	183	80.5	Black
1804	John	25	M	183	80.5	Black
1904	John	25	M	183	80.5	Black
2004	John	25	M	183	80.5	Black
2104	John	25	M	183	80.5	Black
2204	John	25	M	183	80.5	Black
2304	John	25	M	183	80.5	Black
2404	John	25	M	183	80.5	Black
2504	John	25	M	183	80.5	Black
2604	John	25	M	183	80.5	Black
2704	John	25	M	183	80.5	Black
2804	John	25	M	183	80.5	Black
2904	John	25	M	183	80.5	Black
3004	John	25	M	183	80.5	Black
3104	John	25	M	183	80.5	Black
3204	John	25	M	183	80.5	Black
3304	John	25	M	183	80.5	Black
3404	John	25	M	183	80.5	Black
3504	John	25	M	183	80.5	Black
3604	John	25	M	183	80.5	Black
3704	John	25	M	183	80.5	Black
3804	John	25	M	183	80.5	Black
3904	John	25	M	183	80.5	Black
4004	John	25	M	183	80.5	Black
4104	John	25	M	183	80.5	Black
4204	John	25	M	183	80.5	Black
4304	John	25	M	183	80.5	Black
4404	John	25	M	183	80.5	Black
4504	John	25	M	183	80.5	Black
4604	John	25	M	183	80.5	Black
4704	John	25	M	183	80.5	Black
4804	John	25	M	183	80.5	Black
4904	John	25	M	183	80.5	Black
5004	John	25	M	183	80.5	Black
5104	John	25	M	183	80.5	Black
5204	John	25	M	183	80.5	Black
5304	John	25	M	183	80.5	Black
5404	John	25	M	183	80.5	Black
5504	John	25	M	183	80.5	Black
5604	John	25	M	183	80.5	Black
5704	John	25	M	183	80.5	Black
5804	John	25	M	183	80.5	Black
5904	John	25	M	183	80.5	Black
6004	John	25	M	183	80.5	Black
6104	John	25	M	183	80.5	Black
6204	John	25	M	183	80.5	Black
6304	John	25	M	183	80.5	Black
6404	John	25	M	183	80.5	Black
6504	John	25	M	183	80.5	Black
6604	John	25	M	183	80.5	Black
6704	John	25	M	183	80.5	Black
6804	John	25	M	183	80.5	Black
6904	John	25	M	183	80.5	Black
7004	John	25	M	183	80.5	Black
7104	John	25	M	183	80.5	Black
7204	John	25	M	183	80.5	Black
7304	John	25	M	183	80.5	Black
7404	John	25	M	183	80.5	Black
7504	John	25	M	183	80.5	Black
7604	John	25	M	183	80.5	Black
7704	John	25	M	183	80.5	Black
7804	John	25	M	183	80.5	Black
7904	John	25	M	183	80.5	Black
8004	John	25	M	183	80.5	Black
8104	John	25	M	183	80.5	Black
8204	John	25	M	183	80.5	Black
8304	John	25	M	183	80.5	Black
8404	John	25	M	183	80.5	Black
8504	John	25	M	183	80.5	Black
8604	John	25	M	183	80.5	Black
8704	John	25	M	183	80.5	Black
8804	John	25	M	183	80.5	Black
8904	John	25	M	183	80.5	Black
9004	John	25	M	183	80.5	Black
9104	John	25	M	183	80.5	Black
9204	John	25	M	183	80.5	Black
9304	John	25	M	183	80.5	Black
9404	John	25	M	183	80.5	Black
9504	John	25	M	183	80.5	Black
9604	John	25	M	183	80.5	Black
9704	John	25	M	183	80.5	Black
9804	John	25	M	183	80.5	Black
9904	John	25	M	183	80.5	Black
10004	John	25	M	183	80.5	Black



Statistic



Linear Models

Much of what we've talked about so far involves calculating coefficients which describe a specific set of data

- Given a sample of data $(x_1, y_1), (x_2, y_2) \dots, (x_n, y_n)$, calculate line which minimizes RSS
- Sample is all we have, but most often we are interested in quantities which describe a population
- Given a new sample (potentially repeating the experiment) will give different estimates of \hat{b}_0 and \hat{b}_1
- What can we say about \hat{b}_0, \hat{b}_1 and the “true” population process?

Linear Model Assumptions

Commonly used linear model where ε_i is an error term:

$$Y_i = b_0 + b_1 X_i + \varepsilon_i$$

Linear Model Assumptions

Commonly used linear model where ε_i is an error term:

$$Y_i = b_0 + b_1 X_i + \varepsilon_i$$

Assumptions of the model:

- Linear function: $E(Y_i | X_i = x) = b_0 + b_1 x$
- Independence across observations: ε_i is independent of ε_k where i and k denote different observations
- Independence of errors: ε_i is independent of X_i with mean 0 and variance σ^2

Linear Model Assumptions

Commonly used linear model where ε_i is an error term:

$$Y_i = b_0 + b_1 X_i + \varepsilon_i$$

Assumptions of the model:

- Linear function: $E(Y_i | X_i = x) = b_0 + b_1 x$
- Independence across observations: ε_i is independent of ε_k where i and k denote different observations
- Independence of errors: ε_i is independent of X_i with mean 0 and variance σ^2

Less important assumption:

- Normality: sometimes, we assume that $\varepsilon_i \sim N(0, \sigma^2)$

Model Implications

Conditional expectation: $E(Y_i | X_i = x) = b_0 + b_1x$

Model Implications

Conditional expectation: $E(Y_i | X_i = x) = b_0 + b_1x$

Interpretation

- b_0 is the expected value of Y_i when conditioning on $X_i = 0$
- b_1 is the difference of the expected value of Y_i when conditioning on values of X_i which differ by 1 unit.

$$b_1 = E(Y_i | X_i = x + 1) - E(Y_i | X_i = x)$$

Conditional Expectation

In general, the conditional expectation is not the same as “intervening” on X

Conditional Expectation

In general, the conditional expectation is not the same as “intervening” on X

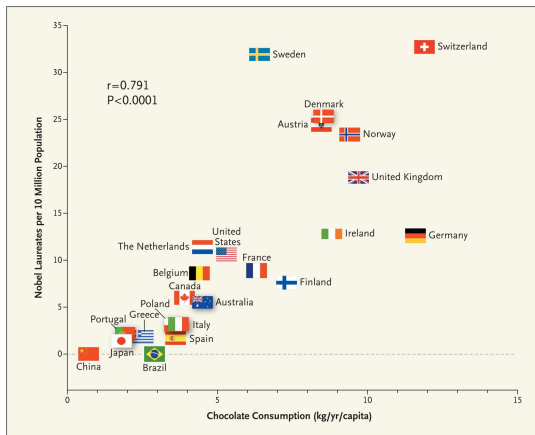


Figure: Messerli 2012, NEJM

Interpretation

Correct Interpretations

- Given two observations whose X values differ by 1 unit, we would **expect** the observation with the larger X value to have a Y value b_1 units larger than the observation with the smaller X value
- Given two observations whose X values differ by 1 unit, **on average** the observation with the larger X value will have a Y value b_1 units larger than the observation with the smaller X value

Interpretation

Correct Interpretations

- Given two observations whose X values differ by 1 unit, we would **expect** the observation with the larger X value to have a Y value b_1 units larger than the observation with the smaller X value
- Given two observations whose X values differ by 1 unit, **on average** the observation with the larger X value will have a Y value b_1 units larger than the observation with the smaller X value
- A 1 unit difference in X is associated with a b_1 unit difference in Y

Interpretation

Correct Interpretations

- Given two observations whose X values differ by 1 unit, we would **expect** the observation with the larger X value to have a Y value b_1 units larger than the observation with the smaller X value
- Given two observations whose X values differ by 1 unit, **on average** the observation with the larger X value will have a Y value b_1 units larger than the observation with the smaller X value
- A 1 unit difference in X is associated with a b_1 unit difference in Y

Incorrect Interpretations

- Increasing X by 1 unit increases Y by b_1 units
- A 1 unit increase in X causes Y to increase by b_1 units

Statistic is unbiased

Under the assumptions that ε_i is independent of X_i , we have:

$$E(\hat{b}_1) = b_1$$

$$E(\hat{b}_0) = b_0$$

so that the estimated values are “unbiased” estimators of the true values

- If you replicate the experiment many different times, you will get a different estimate, each time, but the average will be the “truth”

Potentially helpful (but not necessary) math

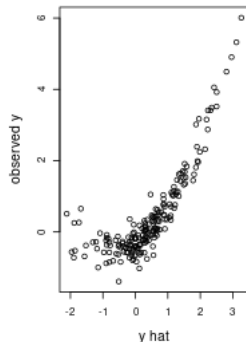
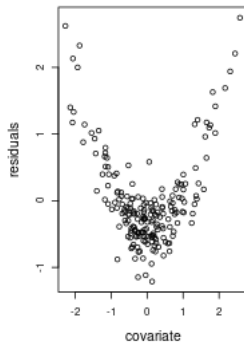
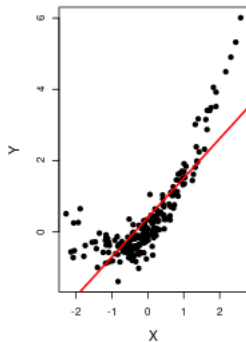
Under the assumptions, we have:

$$\bar{y} = \frac{1}{n} \sum_i (b_0 + b_1 x_i + \varepsilon_i) = b_0 + \frac{1}{n} \sum_i b_1 x_i + \frac{1}{n} \sum_i \varepsilon_i = b_0 + b_1 \bar{x} + \bar{\varepsilon}$$

$$\begin{aligned} E(\hat{b}_1 | X) &= E\left(\frac{\sum_i (y_i - \bar{y})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \mid X\right) \\ &= E\left(\frac{\sum_i (b_0 + b_1 x_i + \varepsilon_i - b_0 - b_1 \bar{x} - \bar{\varepsilon})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \mid X\right) \\ &= E\left(\frac{b_1 \sum_i (x_i - \bar{x})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \mid X\right) + \underbrace{E\left(\frac{\sum_i (\varepsilon_i - \bar{\varepsilon})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \mid X\right)}_{\text{cov}(\varepsilon_i, X_i)=0} \\ &= b_1 + 0 \end{aligned}$$

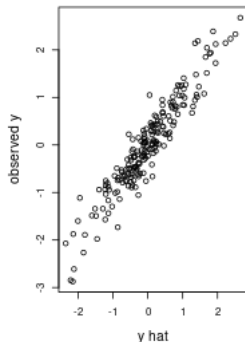
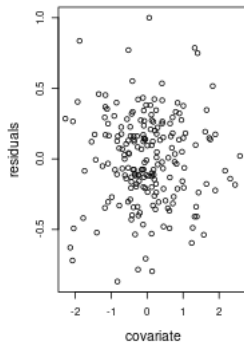
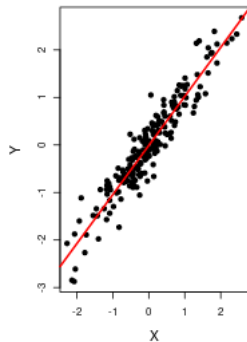
Model Assumptions: Linearity

Look for patterns in residuals if the linearity assumption is violated



Model Assumptions: Linearity

Look for patterns in residuals if the linearity assumption is violated



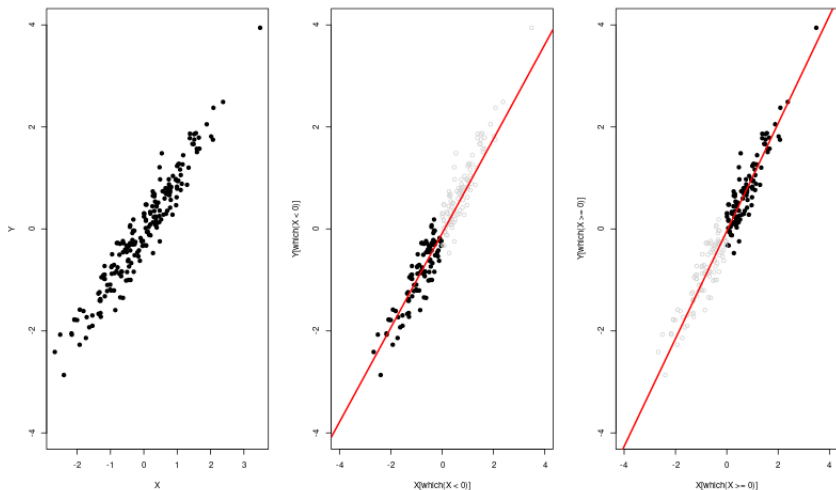
Model Assumptions: Linearity

What happens if the linearity assumption is violated?

- Consider transforming your data with a non-linear transformation
- Adding other covariates can be “helpful”
- b_1 no longer corresponds to change in conditional expectation, but the sign of coefficient can still be useful for interpretability
- Parameters are the best “linear approximation”
- Best linear approximation depends on the range of the X values

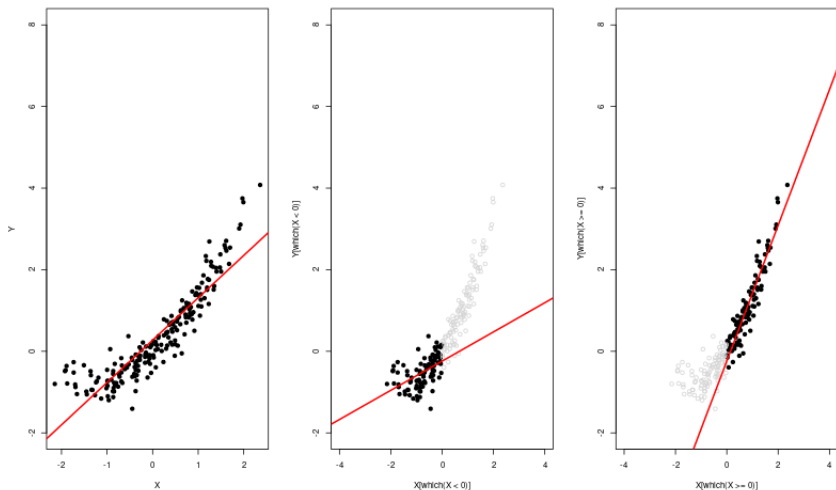
Model Assumptions: Linearity

Best linear approximation depends on the range of the X values



Model Assumptions: Linearity

Best linear approximation depends on the range of the X values



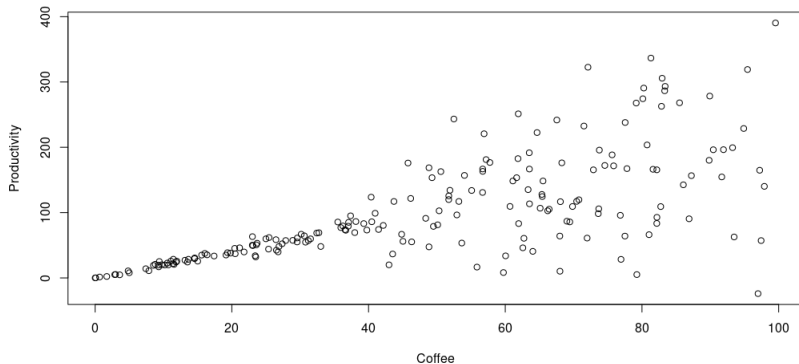
Model Assumptions: independence across observations

- Observations are independent if the value of one observation does not influence or provide information about the value of another observation.
- Ensures that the estimated coefficients and their associated statistical inferences (e.g., confidence intervals, hypothesis tests) are valid and reliable.
- Observations collected over time (e.g., stock prices, temperature readings) are often correlated with past values (autocorrelation).

Model Assumptions: independence of error and covariate

We made a strong assumption that ε_i is mean 0 and independent of X_i

- What if the variance of ε_i depends on X_i ? i.e., model is heteroscedastic
- As long as $E(\varepsilon_i | X_i) = 0$, estimates are still unbiased $E(\hat{b}_1) = b_1$
- Will effect testing procedures!



Discussion

- What is a scientific question that you are interested in?
- Are you trying to do prediction or modeling?
- Are the assumptions we discussed today reasonable for your setting?
 - Linearity
 - Independence across observations
 - Independence of errors and covariates

Assessing explanatory power

Components of the squared error

How can we assess how useful the explanatory variable is for predicting the response variable?

$$\begin{aligned}(y_i - \bar{y}) &= (y_i - \hat{y}_i + \hat{y}_i - \bar{y}) \\ &= (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}) \\ &= \text{residual} + \text{predicted deviation from mean}\end{aligned}\tag{3}$$

Components of the squared error

How can we assess how useful the explanatory variable is for predicting the response variable?

$$\begin{aligned}(y_i - \bar{y}) &= (y_i - \hat{y}_i + \hat{y}_i - \bar{y}) \\ &= (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}) \\ &= \text{residual} + \text{predicted deviation from mean}\end{aligned}\tag{3}$$

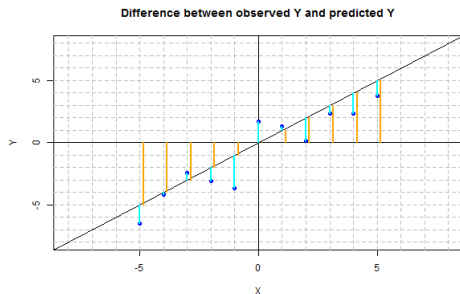
Using a bit of algebra, we can decompose the total sum of squares for Y into

$$SS_{total} = \sum_i (y_i - \bar{y})^2 = \underbrace{\sum_i (\hat{y}_i - \bar{y})^2}_{SS_{regression}} + \underbrace{\sum_i (y_i - \hat{y}_i)^2}_{SS_{error}}\tag{4}$$

Components of the squared error

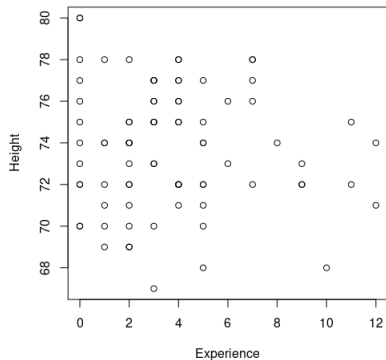
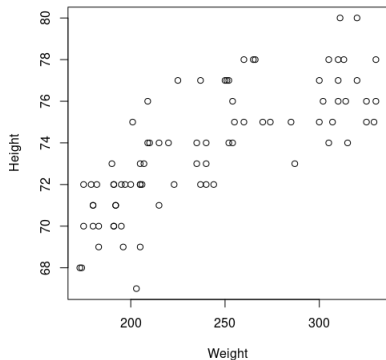
If $SS_{regression}$ is large compared to SS_{error} , then the explanatory variable is a good predictor of the response variable

$$\underbrace{1 - \frac{SS_{error}}{SS_{total}} = \frac{SS_{regression}}{SS_{total}} = \frac{\sum_i (\hat{y}_i - \bar{y})}{\sum_i (y_i - \bar{y})} = r_{XY}^2}_{\text{Referred to as } R^2} \quad (5)$$



Example: Components of the squared error

The R^2 for height and weight is .59 while the R^2 for height and experience is .01.



Wrap-up

- If we assume the true population process is a linear model, we can describe properties of the estimated regression coefficients
- Estimated slope is estimated difference in conditional expectation associated with difference in X
- If assumptions are violated, interpretation is not as straightforward
- Explanatory power of regression can be summarized by R^2 value
- Next module will consider setting with more than 1 covariate