

# HW 1

1. Consider the following statement: "For the ordinary least squares method to be fully valid, it is required that the distribution of Y be normal." Is this statement true or false, and why?

False, as long as the errors average out to 0, the variance is constant for all  $x_i$ s, and the errors are independent in the model (no covariance), then OLS can be utilized.

2. Read section 1.8 of the textbook. When there is a Normal distribution assumption on the error terms, we can also formulate Maximum Likelihood Estimators for  $\beta_0$ ,  $\beta_1$ , &  $\sigma^2$ . Use the likelihood function (1.26) to find the estimators  $\hat{\beta}_0$ ,  $\hat{\beta}_1$ , &  $\hat{\sigma}^2$ , and show that the estimators  $\beta_0$ ,  $\beta_1$ , are the same as the Least Square estimators. You do not need to check second derivatives to prove the maximum values.

$$L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left(-\frac{1}{2\sigma^2} (y_i - \beta_0 - \beta_1 x_i)^2\right)}$$
$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n e^{\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2\right)}$$

log likelihood

$$l = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Set derivatives to 0

$$\frac{\partial l}{\partial \beta_0} = 2 \cdot -\frac{1}{2\sigma^2} \cdot \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (-1)$$

chain rule

$$\frac{\partial L}{\partial \beta_0} = \frac{1}{\sigma^2} \cdot \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\frac{\partial L}{\partial \beta_1} = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

replace  $\sigma^2$  with  $X$

$$\frac{\partial L}{\partial \sigma^2} = -\frac{n}{2} \frac{1}{X} \cdot \frac{1}{2} \frac{1}{X^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = 0$$

$$\frac{n}{2\sigma^2} = \frac{1}{2\sigma^4} \left( \right)$$

$$\sigma^2 = \frac{\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2}{n}$$

$n$

Normal eqn

log eqn

$$\sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$n \hat{\beta}_0 = \sum y_i - \hat{\beta}_1 \sum x_i$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \checkmark$$

$$\sum y_i = n \beta_0 + \beta_1 \sum x_i$$

$$\sum x_i y_i = \beta_0 \sum x_i + \beta_1 \sum x_i^2$$

$$\beta_0 = \frac{1}{n} (\sum y_i - \beta_1 \sum x_i)$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x} \quad \checkmark$$

$$\sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \quad \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \checkmark$$

$$\sum (y_i - \bar{y} + \hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i) x_i = 0$$

$$\sum y_i x_i - \bar{y} \sum x_i + \hat{\beta}_1 \bar{x} \sum x_i - \hat{\beta}_1 \sum x_i^2 = 0$$

$$\sum y_i x_i - \bar{y} \sum x_i = \hat{\beta}_1 (\sum x_i^2 - \bar{x} \sum x_i)$$

$$\hat{\beta}_1 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \quad \checkmark$$

3. The solution for the LS estimator of the slope in simple linear regression is:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2}.$$

Prove that this expression is equivalent to the alternate formulation:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{SSXY}{SSX}.$$

$$\begin{aligned} SSXY &= \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - x_i \bar{y} - y_i \bar{x} - \bar{x} \bar{y} \\ &= \sum x_i y_i - \bar{y} \sum x_i - \bar{x} \sum y_i + n \bar{x} \bar{y} \\ &= \sum x_i y_i - \bar{y} n \bar{x} - \bar{x} n \bar{y} + n \bar{x} \bar{y} \\ &= \sum x_i y_i - n \bar{x} \bar{y} \end{aligned}$$

$$\begin{aligned} SSX &= \sum (x_i - \bar{x})^2 = \sum x_i^2 - 2 \bar{x} \sum x_i + n \bar{x}^2 \\ &= \sum x_i^2 - 2 \bar{x} \cdot n \bar{x} + n \bar{x}^2 = \sum x_i^2 - n \bar{x}^2 \end{aligned}$$

$$\frac{SSXY}{SSX} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} \quad \checkmark$$

4. Recall that the residual for the  $i^{th}$  observation is defined as  $e_i = y_i - \hat{y}_i$ . Prove that  $E(e_i) = 0$ .

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

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

$$e_i = \beta_0 + \beta_1 x_i + \varepsilon_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$$

$$= (\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) x_i + \varepsilon_i$$

$$E(e_i) = E \left[ \begin{array}{c} \downarrow \end{array} \right]$$

$$E(e_i) = E(\beta_0 - \hat{\beta}_0) + E(\beta_1 - \hat{\beta}_1) x_i + E(\varepsilon_i)$$

if unbiased  $\hat{\beta}_0 = \beta_0$ ,  $\hat{\beta}_1 = \beta_1$ ,

and  $E(\varepsilon_i) = 0$  by property of OLS (Gauss-Markov thm.)

$$E(e_i) = E(0) + E(0) x_i + 0$$

$$E(e_i) = 0$$

5. (a) When asked to state the simple linear regression model, a student wrote:

$$E(Y_i) = \beta_0 + \beta_1 x_i + \epsilon_i.$$

Do you agree? Why or why not?

- (b) Consider the classical simple linear regression model. Suppose that the true parameter values are  $\beta_0 = 2$ ,  $\beta_1 = 4$ , and  $\sigma^2 = 9$ . State the distributions of  $Y$  at  $x = 1, 2$ , and  $4$ , and explain how you found them.

a) No, the expected value would not include the error at that term and would just be

$$E(Y_i) = \beta_0 + \beta_1 x_i$$

the actual model would just be

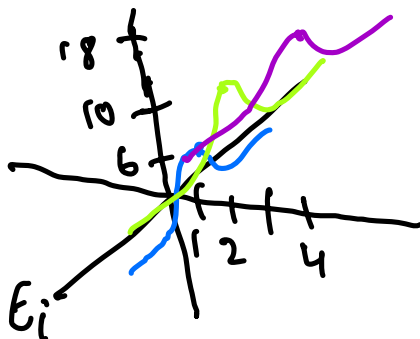
$$b) \quad Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$x=1 \quad E(Y_i) = 2 + 4(1) = 6: Y \sim N(6, 9)$$

$$x=2 \quad E(Y_i) = 2 + 4(2) = 10: Y \sim N(10, 9)$$

$$x=4 \quad E(Y_i) = 2 + 4(4) = 18: Y \sim N(18, 9)$$

Expected value with normally distributed errors



6. Consider the Rotten Tomatoes movie rating example.

- (a) Interpret the slope and the intercept in the real-life context of the problem.
- (b) Suppose the Borderlands movie is about to be released and critics have given it a score of 8 (out of 100) on Rotten Tomatoes. Using the fitted simple linear regression line, what do we predict the audience rating will be?
- (c) What does the SLR prediction in (b) suggest about who the regression line thinks is the harsher judge: audiences or critics?
- (d) Suppose "The Quiet Place: Day One" movie is about to be released and critics have given it a score of 86 (out of 100) on Rotten Tomatoes. Using the fitted simple linear regression line, what do we predict the audience rating will be?

a)  $\beta_0$ : Hypothetical predicted audience rating based on critic rating of 0

$\beta_1$ : Change in predicted audience rating at each individual point increase in critic rating (or decrease if negative)

b)  $\hat{Y} = \beta_0 + \beta_1 X = \beta_0 + \beta_1 \cdot 8$

See python notebook,  $\hat{\beta}_0 \approx 34.51$

$$\hat{\beta}_1 \approx 0.446$$

$$\hat{Y}(8) \approx 38.08$$

c) The regression model thinks  
Critics are much harsher

$$\hat{Y}(8) \gg 8$$

d)  $\hat{Y}(86) \approx 72.87$  (see python notebook)

- (e) What does the SLR prediction (d) suggest about who the regression line thinks is the harsher judge: audiences or critics?
- (f) Consider your findings in (c) and (e). Provide a reasonable explanation as to how one can reconcile these two results.
- (g) What value of critic ratings will the SLR model predict the exact same score for audience ratings? Derive a general formula for this value in terms of  $\hat{\beta}_0$  &  $\hat{\beta}_1$ .

e) Since  $\hat{Y}(86) \ll 86$  suggests  
audiences are harsher

f) It seems to suggest the ratings  
are not truly linearly related across  
the whole range.

g) 
$$X = \beta_1 X + \beta_0$$
$$X = \frac{\beta_0}{1 - \beta_1} \approx \frac{34.508}{1 - 0.446} \approx 62.29$$