

陶思都

519021910906

ECE 4710J, Spring 2022

Homework #4

## Properties of Simple Linear Regression

1. (10 points) In lecture, we spent a great deal of time talking about simple linear regression. To briefly summarize, the simple linear regression model assumes that given a single observation  $x$ , our predicted response for this observation is  $\hat{y} = \theta_0 + \theta_1 x$ . (Note: In this problem we write  $(\theta_0, \theta_1)$  instead of  $(a, b)$  to more closely mirror the multiple linear regression model notation.)

We saw that the  $\theta_0 = \hat{\theta}_0$  and  $\theta_1 = \hat{\theta}_1$  that minimize the average  $L_2$  loss for the simple linear regression model are:

$$\begin{aligned}\hat{\theta}_0 &= \bar{y} - \hat{\theta}_1 \bar{x} \\ \hat{\theta}_1 &= r \frac{\sigma_y}{\sigma_x}\end{aligned}$$

Or, rearranging terms, our predictions  $\hat{y}$  are:

$$\hat{y} = \bar{y} + r \sigma_y \frac{x - \bar{x}}{\sigma_x}$$

- (a) (4 points) As we saw in lecture, a residual  $e_i$  is defined to be the difference between a true response  $y_i$  and predicted response  $\hat{y}_i$ . Specifically,  $e_i = y_i - \hat{y}_i$ . Note that there are  $n$  data points, and each data point is denoted by  $(x_i, y_i)$ .

Prove, using the equation for  $\hat{y}$  above, that  $\sum_{i=1}^n e_i = 0$  (meaning the sum of the residuals is zero).

Proof: Since  $\hat{y} = \bar{y} + r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$ ,  $\sum_{i=1}^n y_i = n\bar{y}$ ,  $\sum_{i=1}^n x_i = n\bar{x}$ , we have

$$\begin{aligned}\sum_{i=1}^n e_i &= \sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - \bar{y} - r \frac{\sigma_y}{\sigma_x} x_i + r \frac{\sigma_y}{\sigma_x} \bar{x}) \\ &= \sum_{i=1}^n y_i - n\bar{y} - r \frac{\sigma_y}{\sigma_x} (\sum_{i=1}^n x_i - n\bar{x}) \\ &= 0 - r \frac{\sigma_y}{\sigma_x} \cdot 0 \\ &= 0\end{aligned}$$

Therefore, we have  $\sum_{i=1}^n e_i = 0$

(b) (3 points) Using your result from part (a), prove that  $\bar{y} = \hat{y}$ .

Proof: Since  $\sum_{i=1}^n \hat{y}_i = n\bar{y}$ ,  $\sum_{i=1}^n x_i = n\bar{x}$ , we have

$$\begin{aligned}\bar{y} &= \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \frac{1}{n} \sum_{i=1}^n (\bar{y} + r \frac{\partial y}{\partial x} (x_i - \bar{x})) \\ &= \frac{1}{n} \sum_{i=1}^n \bar{y} + \frac{1}{n} \cdot r \frac{\partial y}{\partial x} (\sum_{i=1}^n x_i - n\bar{x}) \\ &= \frac{1}{n} \cdot n\bar{y} + \frac{1}{n} \cdot r \frac{\partial y}{\partial x} \cdot 0 \\ &= \bar{y}\end{aligned}$$

Therefore, we have  $\bar{y} = \hat{y}$

(c) (3 points) Prove that  $(\bar{x}, \bar{y})$  is on the simple linear regression line.

Proof: Let  $(\bar{x}, \bar{y}_1)$  be on the simple linear regression line, then

$$\begin{aligned}y_1 &= \hat{\theta}_0 + \hat{\theta}_1 \bar{x} = \bar{y} + r \frac{\partial y}{\partial x} (\bar{x} - \bar{x}) \\ &= \bar{y} + r \frac{\partial y}{\partial x} \cdot 0 \\ &= \bar{y}\end{aligned}$$

$$\Rightarrow y_1 = \bar{y}$$

Therefore,  $(\bar{x}, \bar{y})$  is on the simple linear regression model.

## Geometric Perspective of Least Squares

2. (10 points) We also viewed both the simple linear regression model and the multiple linear regression model through linear algebra. The key geometric insight was that if we train a model on some design matrix  $\mathbb{X}$  and true response vector  $\mathbb{Y}$ , our predicted response  $\hat{\mathbb{Y}} = \mathbb{X}\hat{\theta}$  is the vector in  $\text{span}(\mathbb{X})$  that is closest to  $\mathbb{Y}$  ( $\hat{\mathbb{Y}}$  is the orthogonal projection of  $\mathbb{Y}$  onto the  $\text{span}(\mathbb{X})$ ).

In the simple linear regression case, our optimal vector  $\theta$  is  $\hat{\theta} = [\hat{\theta}_0, \hat{\theta}_1]^T$ , and our design matrix is

$$\mathbb{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} = \begin{bmatrix} | & | \\ \mathbb{1} & \vec{x} \\ | & | \end{bmatrix}$$

This means we can write our predicted response vector as  $\hat{\mathbb{Y}} = \mathbb{X} \begin{bmatrix} \hat{\theta}_0 \\ \hat{\theta}_1 \end{bmatrix} = \hat{\theta}_0 \mathbb{1} + \hat{\theta}_1 \vec{x}$ .

Note, in this problem,  $\vec{x}$  refers to the  $n$ -length vector  $[x_1, x_2, \dots, x_n]^T$ . In other words, it is a feature, not an observation.

For this problem, assume we are working with the **simple linear regression model**, though the properties we establish here hold for any linear regression model that contains an intercept term.

- (a) (4 points) Using the geometric properties from lecture, prove that  $\sum_{i=1}^n e_i = 0$ .

*Hint:* Recall, we define the residual vector as  $e = \mathbb{Y} - \hat{\mathbb{Y}}$ , and  $e = [e_1, e_2, \dots, e_n]^T$ .

*Proof:* Since the model has intercept,  $\mathbb{1}$ , which is the first column of  $\mathbb{X}$ ,

$$\Rightarrow \mathbb{1} = \mathbb{X} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \mathbb{X} \vec{c}, \text{ where } \vec{c} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

then we have:

$$\begin{aligned} \sum_{i=1}^n e_i &= \mathbb{1}^T \vec{e} = \mathbb{1}^T (\mathbb{Y} - \hat{\mathbb{Y}}) \\ &= \vec{c}^T \mathbb{X}^T (\mathbb{Y} - \mathbb{X} \hat{\theta}) \\ &= \vec{c}^T \mathbb{X}^T (\mathbb{Y} - \mathbb{X} (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}) \\ &= \vec{c}^T (\mathbb{X}^T \mathbb{Y} - \mathbb{X}^T \mathbb{X} (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}) \\ &= \vec{c}^T (\mathbb{X}^T \mathbb{Y} - \mathbb{X}^T \mathbb{Y}) \\ &= 0 \end{aligned}$$

Therefore, we have  $\sum_{i=1}^n \vec{e}_i = 0$ .

- (b) (3 points) Explain why the vector  $\vec{x}$  (as defined in the problem) and the residual vector  $e$  are orthogonal. Hint: Two vectors are orthogonal if their dot product is 0.

To obtain a good model, we need to minimize the average loss with geometry, that is, to minimize the  $L_2$  norm of the residual vector. To achieve this goal, we need the vector in  $\text{span}(\vec{x})$  that is closest to  $\vec{Y}$  is the orthogonal projection of  $\vec{Y}$  onto  $\text{span}(\vec{x})$ , i.e. the vector  $\vec{x}$  is orthogonal to the residual vector  $\vec{e}$ .

Proof : Since  $\vec{x}$  is the second column of  $\vec{X}$ , we have

$$\Rightarrow \vec{x} = \vec{X} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \vec{X} \vec{c}_1, \text{ where } \vec{c}_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

then we have :

$$\begin{aligned} \vec{x}^T \vec{e} &= \vec{c}_1^T \vec{X}^T (\vec{Y} - \hat{\vec{Y}}) \\ &= \vec{c}_1^T (\vec{X}^T \vec{Y} - \vec{X}^T \vec{X} (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{Y}) \\ &= \vec{c}_1^T \cdot (\vec{X}^T \vec{Y} - \vec{X}^T \vec{Y}) \\ &= 0 \end{aligned}$$

Therefore, we have  $\vec{x}$  and  $\vec{e}$  are orthogonal.

- (c) (3 points) Explain why the predicted response vector  $\hat{\vec{Y}}$  and the residual vector  $e$  are orthogonal.

Since our predicted response  $\hat{\vec{Y}}$  is in  $\text{span}(\vec{x})$  by definition and  $\vec{x}^T \vec{e} = 0$ , i.e.  $\vec{x}$  is orthogonal to  $\vec{e}$ ,  $\hat{\vec{Y}}$  should also be orthogonal to  $\vec{e}$ .

Proof : Since  $\hat{\vec{Y}}$  is in  $\text{span}(\vec{x})$ , we have

$$\Rightarrow \hat{\vec{Y}} \in \text{span}(\vec{x})$$

According to the previous question (b), we know  $\vec{x}$  is orthogonal to  $\vec{e}$ . Since  $\vec{x}$  is a feature and is an arbitrary vector, we have :

$$\Rightarrow \text{span}(\vec{x}) \text{ is orthogonal to vector } \vec{e}$$

Therefore,  $\hat{\vec{Y}}$  is orthogonal to vector  $\vec{e}$ .

## Properties of a Linear Model With No Constant Term

Suppose that we don't include an intercept term in our model. That is, our model is now

$$\hat{y} = \gamma x,$$

where  $\gamma$  is the single parameter for our model that we need to optimize. (In this equation,  $x$  is a scalar, corresponding to a single observation.)

As usual, we are looking to find the value  $\hat{\gamma}$  that minimizes the average  $L_2$  loss (mean squared error) across our observed data  $\{(x_i, y_i)\}, i = 1, \dots, n$ :

$$R(\gamma) = \frac{1}{n} \sum_{i=1}^n (y_i - \gamma x_i)^2$$

The normal equations derived in lecture no longer hold. In this problem, we'll derive a solution to this simpler model. We'll see that the least squares estimate of the slope in this model differs from the simple linear regression model, and will also explore whether or not our properties from the previous problem still hold.

3. (5 points) Use calculus to find the minimizing  $\hat{\gamma}$ . That is, prove that

$$\hat{\gamma} = \frac{\sum x_i y_i}{\sum x_i^2}$$

Note: This is the slope of our regression line, analogous to  $\hat{\theta}_1$  from our simple linear regression model.

According to the description, MSE is  $R(\gamma) = \frac{1}{n} \sum_{i=1}^n (y_i - \gamma x_i)^2$ .

$$\begin{aligned} R(\gamma) &= \frac{1}{n} \sum_{i=1}^n (y_i - \gamma x_i)^2 = \frac{1}{n} \sum_{i=1}^n (y_i^2 - 2\gamma y_i x_i + \gamma^2 x_i^2) \\ &= \frac{1}{n} \sum_{i=1}^n y_i^2 - \frac{2\gamma}{n} \sum_{i=1}^n x_i y_i + \frac{\gamma^2}{n} \sum_{i=1}^n x_i^2 \\ \Rightarrow R'(\gamma) &= \frac{dR(\gamma)}{d\gamma} = -\frac{2}{n} \sum_{i=1}^n x_i y_i + \frac{2\gamma}{n} \sum_{i=1}^n x_i^2 \end{aligned}$$

We take  $R'(\gamma) = 0$ . Then we have:

$$\begin{aligned} \frac{2\gamma}{n} \sum_{i=1}^n x_i^2 &= \frac{2}{n} \sum_{i=1}^n x_i y_i \\ \Rightarrow \hat{\gamma} &= \frac{\frac{n}{n} \sum_{i=1}^n x_i y_i}{\frac{n}{n} \sum_{i=1}^n x_i^2} \end{aligned}$$

Since  $R'(\gamma)$  is monotonically increasing,  $R(\gamma)$  is minimum when  $R'(\gamma) = 0$ , i.e.  $\hat{\gamma} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$

Therefore, the minimizing  $\hat{\gamma}$  is  $\frac{\sum x_i y_i}{\sum x_i^2}$

4. (10 points) For our new simplified model, our design matrix  $\mathbb{X}$  is:

$$\mathbb{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} | \\ | \\ | \\ | \end{bmatrix} \vec{x}.$$

Therefore our predicted response vector  $\hat{Y}$  can be expressed as  $\hat{Y} = \hat{\gamma} \vec{x}$ . ( $\vec{x}$  here is defined the same way it was in Question 2.)

Earlier in this homework, we established several properties that held true for the simple linear regression model that contained an intercept term. For each of the following four properties, state whether or not they still hold true even when there isn't an intercept term. Be sure to justify your answer.

(a) (2 points)  $\sum_{i=1}^n e_i = 0$ .

No. When there is not an intercept term, we have

$$\begin{aligned} \sum_{i=1}^n e_i &= \mathbb{1}^T \vec{e} = \mathbb{1}^T (\mathbb{Y} - \hat{\mathbb{Y}}) = \mathbb{1}^T (\mathbb{Y} - \hat{\gamma} \vec{x}) \\ &= \mathbb{1}^T (\mathbb{Y} - (\vec{x}^T \vec{x})^{-1} \vec{x}^T \mathbb{Y} \vec{x}) \\ &= \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \cdot \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} \\ &= n\bar{y} - n\bar{x} \cdot \frac{\sum x_i y_i}{\sum x_i^2} \end{aligned}$$

Only when  $\bar{y} = \bar{x} \cdot \frac{\sum x_i y_i}{\sum x_i^2}$ ,  $\sum_{i=1}^n e_i = 0$

$\Rightarrow \sum_{i=1}^n e_i$  is not necessarily equal to 0.

(b) (3 points) The column vector  $\vec{x}$  and the residual vector  $e$  are orthogonal.

Yes. We have  $\hat{\gamma} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y} = (\vec{x}^T \vec{x})^{-1} \vec{x}^T \mathbb{Y}$ . Then

$$\begin{aligned} \vec{x}^T \vec{e} &= \vec{x}^T (\mathbb{Y} - \hat{\mathbb{Y}}) = \vec{x}^T \mathbb{Y} - \vec{x}^T \vec{\gamma} \hat{\mathbb{Y}} \\ &= \vec{x}^T \mathbb{Y} - \vec{x}^T \vec{x} (\vec{x}^T \vec{x})^{-1} \vec{x}^T \mathbb{Y} \\ &= \vec{x}^T \mathbb{Y} - \vec{x}^T \mathbb{Y} \\ &= 0 \end{aligned}$$

Therefore,  $\vec{x}$  and  $\vec{e}$  are still orthogonal.

- (c) (3 points) The predicted response vector  $\hat{Y}$  and the residual vector  $e$  are orthogonal.

Yes. We have  $\hat{Y} = \vec{x} \hat{y} \in \text{span}(\vec{x})$ . Then

$$\begin{aligned}\hat{Y}^T \cdot \vec{e} &= \hat{y} \vec{x}^T (\vec{Y} - \vec{x} \hat{y}) = \hat{y} (\vec{x}^T \vec{Y} - \vec{x}^T \vec{x} (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{Y}) \\ &= \hat{y} (\vec{x}^T \vec{Y} - \vec{x}^T \vec{Y}) \\ &= 0\end{aligned}$$

Therefore,  $\hat{Y}$  and  $\vec{e}$  are still orthogonal.

- (d) (2 points)  $(\bar{x}, \bar{y})$  is on the regression line.

No. Let  $(\bar{x}, \bar{y})$  be on the regression line. Then,

$$y_1 = \bar{x} \hat{y} = \bar{x} \cdot (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{Y} = \bar{x} \cdot \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$$

Only when  $\bar{y} = \bar{x} \cdot \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$ , we have  $y_1 = \bar{y}$ ,  $(\bar{x}, \bar{y})$  on the regression line.

$\Rightarrow (\bar{x}, \bar{y})$  is not necessarily on the regression line.

## MSE “Minimizer”

5. (15 points) Recall from calculus that given some function  $g(x)$ , the  $x$  you get from solving  $\frac{dg(x)}{dx} = 0$  is called a *critical point* of  $g$  – this means it could be a minimizer or a maximizer for  $g$ . In this question, we will explore some basic properties and build some intuition on why, for certain loss functions such as squared  $L_2$  loss, the critical point of the empirical risk function (defined as average loss on the observed data) will always be the minimizer.

Given some linear model  $f(x) = \gamma x$  for some real scalar  $\gamma$ , we can write the empirical risk of the model  $f$  given the observed data  $\{x_i, y_i\}, i = 1, \dots, n$  as the average  $L_2$  loss, also known as mean squared error (MSE):

$$\frac{1}{n} \sum_{i=1}^n (y_i - \gamma x_i)^2.$$

- (a) (2 points) Let's break the function above into individual terms. Complete the following sentence by filling in the blanks using one of the options in the parenthesis following each of the blanks:

quadratic

The mean squared error can be viewed as a sum of  $n$  \_\_\_\_\_ (linear/quadratic/logarithmic/exponential) terms, each of which can be treated as a function of  $\underline{x}$  ( $x_i/y_i/\gamma$ ).

- (b) (4 points) Let's investigate one of the  $n$  functions in the summation in the MSE. Define  $g_i(\gamma) = \frac{1}{n}(y_i - \gamma x_i)^2$  for  $i = 1, \dots, n$ . Recall from calculus that we can use the 2nd derivative of a function to describe its curvature about a certain point (if it is facing concave up, down, or possibly a point of inflection). You can take the following as a fact: A function is convex if and only if the function's 2nd derivative is non-negative on its domain. Based on this property, verify that  $g_i$  is a **convex function**.

$$g_i(\gamma) = \frac{1}{n}(y_i - \gamma x_i)^2 = \frac{1}{n} y_i^2 - \frac{2\gamma}{n} x_i y_i + \frac{\gamma^2}{n} x_i^2$$

$$\frac{dg_i(\gamma)}{d\gamma} = -\frac{2}{n} x_i y_i + \frac{2\gamma}{n} x_i^2$$

$$\frac{d^2g_i(\gamma)}{d\gamma^2} = \frac{2}{n} x_i^2 = \frac{2x_i^2}{n} \geq 0$$

Since the 2nd derivative of  $g_i(\gamma)$  is non-negative on its domain  $\mathbb{R}$ ,  
 $\Rightarrow g_i$  is a convex function

- (c) (3 points) Briefly explain in words why given a convex function  $g(x)$ , the critical point we get by solving  $\frac{dg(x)}{dx} = 0$  minimizes  $g$ . You can assume that  $\frac{dg(x)}{dx}$  is a function of  $x$  (and not a constant).

Since  $g(x)$  is a convex function, we have  $\frac{d^2g(x)}{dx^2} \geq 0$  on its domain.

Thus,  $\frac{dg(x)}{dx}$  is monotonically increasing on its domain. Let  $x_0$  be the point such that  $\frac{dg(x)}{dx}|_{x=x_0} = 0$ , then we have

① when  $x < x_0$ ,  $\frac{dg(x)}{dx} < 0$ ,  $g(x)$  is monotonically decreasing.

② when  $x > x_0$ ,  $\frac{dg(x)}{dx} > 0$ ,  $g(x)$  is monotonically increasing.

Therefore, the critical point  $x_0$  we get by solving  $\frac{dg(x)}{dx} = 0$  minimizes  $g$ .

- (d) (4 points) Now that we have shown that each term in the summation of the MSE is a convex function, one might wonder if the entire summation is convex given that it is a sum of convex functions.

Let's look at the formal definition of a **convex function**. Algebraically speaking, a function  $g(x)$  is convex if for any two points  $(x_1, g(x_1))$  and  $(x_2, g(x_2))$  on the function,

$$g(cx_1 + (1 - c)x_2) \leq cg(x_1) + (1 - c)g(x_2)$$

for any real constant  $0 \leq c \leq 1$ .

The above definition says that, given the plot of a convex function  $g(x)$ , if you connect 2 randomly chosen points on the function, the line segment will always lie on or above  $g(x)$  (try this with the graph of  $y = x^2$ ).

- i. (2 points) Using the definition above, show that if  $g(x)$  and  $h(x)$  are both convex functions, their sum  $g(x) + h(x)$  will also be a convex function.

Let  $p(x) = g(x) + h(x)$ . Since  $g(x)$  and  $h(x)$  are convex functions,

$$g(cx_1 + (1 - c)x_2) \leq cg(x_1) + (1 - c)g(x_2)$$

$$h(cx_1 + (1 - c)x_2) \leq ch(x_1) + (1 - c)h(x_2)$$

for any two points  $x_1, x_2$  and any constant  $0 \leq c \leq 1$ . Sum them up:

$$g(cx_1 + (1 - c)x_2) + h(cx_1 + (1 - c)x_2) \leq cg(x_1) + h(x_1) + (1 - c)g(x_2) + h(x_2)$$

$$\Rightarrow p(cx_1 + (1 - c)x_2) \leq cg(x_1) + (1 - c)g(x_2) + h(x_1) + h(x_2)$$

Therefore,  $p(x) = g(x) + h(x)$  will also be a convex function.

- ii. (2 points) Based on what you have shown in the previous part, explain intuitively why the sum of  $n$  convex functions is still a convex function when  $n > 2$ .

Summing the functions will not change the features of functions.

Since the functions are independent, summing is only a process of simple superposition, the relationship between the left side and right side of the inequality will not change. Thus, the sum of  $n$  convex functions is still a convex function when  $n > 2$ .

- (e) (2 points) Finally, using the previous parts, explain why in our case that, when we solve for the critical point of the MSE by taking the gradient with respect to the parameter and setting the expression to 0, it is guaranteed that the solution we find will minimize the MSE.

According to (b), we know that  $g_i(\gamma) = \frac{1}{n} (y_i - \gamma x_i)^2$  for  $i=1, 2, \dots, n$  is convex function

According to (c), we know that for a convex function, the critical point we get by solving

$$\frac{df(x)}{dx} = 0 \text{ minimizes the function.}$$

According to (d), we know that the sum of convex functions is still a convex function.

$$\text{Since } \text{MSE} = \sum_{i=1}^n \frac{1}{n} (y_i - \gamma x_i)^2 = \sum_{i=1}^n g_i(\gamma),$$

$\Rightarrow$  it is the sum of convex functions, that it is also a convex function.

Since it is a convex function, according to (b), we have

$$\Rightarrow \text{the critical point we get by solving } \frac{d\text{MSE}}{d\gamma} = 0 \text{ minimizes MSE.}$$

Therefore, we can conclude that the solution we find in the described way will minimize the MSE.