

## Discussion #6

ft soon

Note: Your TA will probably not cover all the problems. This is totally fine, the discussion worksheets are not designed to be finished within an hour. They are deliberately made slightly longer so they can serve as resources you can use to practice, reinforce, and build upon concepts discussed in lectures, labs, and homework.

## Driving with a Constant Model

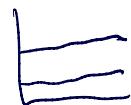
- Lillian is trying to use modeling to drive her car autonomously. To do this, she collects a lot of data from driving around her neighborhood and stores it in `drive`. She wants your help to design a model that can drive on her behalf in the future using the outputs of the models you design. First, she wants to tackle two aspects of this autonomous car modeling framework: going forward and turning.

Some statistics from the collected dataset are shown below using `drive.describe()`, which returns the mean, standard deviation, quartiles, minimum, and maximum for the two columns in the dataset: `target_speed` and `degree_turn`.

	target_speed	degree_turn
<b>count</b>	531.00000	531.00000
<b>mean</b>	32.923408	143.721153
<b>std</b>	46.678744	153.641504
<b>min</b>	0.231601	0.000000
<b>25%</b>	12.350025	6.916210
<b>50%</b>	25.820689	45.490086
<b>75%</b>	39.788716	323.197168
<b>max</b>	379.919965	359.430309

 $L_2$  MSE : $L_1$  MAE

$$y = \theta_0$$



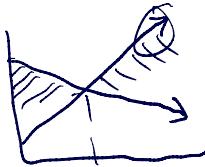
$$L_2 \text{ loss: } y = \hat{y} \approx 32.9$$

$$L_1 \text{ loss: } y = \text{median}(y) \approx 25.8$$

## Discussion #6

- (a) Suppose the first part of the model predicts the target speed of the car. Using constant models trained on the speeds of the collected data shown above with  $L_1$  and  $L_2$  loss functions, which of the following is true?

- A. The model trained with the  $L_1$  loss will always drive slower than the model trained with  $L_2$  loss.
- B. The model trained with the  $L_2$  loss will always drive slower than the model trained with  $L_1$  loss.
- C. The model trained with the  $L_1$  loss will sometimes drive slower than the model trained with  $L_2$  loss.
- D. The model trained with the  $L_2$  loss will sometimes drive slower than the model trained with  $L_1$  loss.



- (b) Finding that the model trained with the  $L_2$  loss drives too slowly, Lillian changes the loss function for the constant model where the loss is penalized **more** if the true speed is higher. That way, in order to minimize loss, the model would have to output predictions closer to the true value, particularly as speeds get faster, the end result being a higher constant speed. Lillian writes this as  $L(y, \hat{y}) = y(y - \hat{y})^2$ .

Find the optimal  $\hat{\theta}_0$  for the constant model using the new empirical risk function  $R(\theta_0)$  below:

$$\frac{d}{d\theta_0} R(\theta_0) = \frac{1}{n} \sum_i y_i (y_i - \theta_0)^2$$

$$\frac{d}{d\theta_0} \frac{1}{n} \sum_i y_i (y_i - \theta_0)^2 = \frac{1}{n} \sum_i \frac{d}{d\theta_0} y_i (y_i - \theta_0)^2$$

$$= \frac{1}{n} \sum_i y_i \cdot 2(y_i - \theta_0) \cdot -1 = -\frac{2}{n} \sum_i y_i (y_i - \theta_0)$$

$$\boxed{-\frac{2}{n} \sum_i y_i (y_i - \theta_0)} = 0 \quad \sum_i (y_i^2 - y_i \theta_0) = 0$$

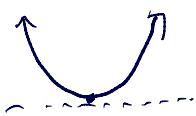
$$\downarrow \quad \sum_i y_i^2 - \sum_i y_i \theta_0 = 0$$

$$\sum_i y_i^2 = \sum_i y_i \theta_0 \rightarrow \sum_i y_i^2 = \theta_0 \sum_i y_i$$

$$\theta_0 = \frac{\sum_i y_i^2}{\sum_i y_i}$$

$$\frac{d}{d\theta_0} \left( -\frac{2}{n} \sum_i y_i^2 - y_i \theta_0 \right)$$

$$= -\frac{2}{n} \sum_i \frac{d}{d\theta_0} (y_i^2 - y_i \theta_0) = -\frac{2}{n} \sum_i -y_i = \frac{2}{n} \sum_i y_i > 0$$



- (c) Lillian's friend, Yash, also begins working on a model that predicts the degree of turning at a particular time between 0 and 359 degrees using the data in the `degree_turn` column. Explain why a constant model is likely inappropriate in this use case.

*Extra:* If you've studied some physics, you may recognize the behavior of our constant model!



- (d) Suppose we finally expand our modeling framework to use simple linear regression (i.e.  $f_\theta(x) = \theta_{w,0} + \theta_{w,1}x$ ). For our first simple linear regression model, we predict the turn angle ( $y$ ) using target speed ( $x$ ). Our optimal parameters are:  $\hat{\theta}_{w,1} = 0.019$  and  $\hat{\theta}_{w,0} = 143.1$ .

However, we realize that we actually want a model that predicts target speed (our new  $y$ ) using turn angle, our new  $x$  (instead of the other way around)! What are our new optimal parameters for this new model?

$$y = \theta_0 + \theta_1 x \quad \theta_1 = r \frac{\sigma_y}{\sigma_x} \quad \theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\text{angle} = \theta_{w,0} + \theta_{w,1} \text{ speed}$$

$$\theta_{w,1} = r \frac{\sigma_{\text{angle}}}{\sigma_{\text{speed}}} \quad \cancel{\text{Speed}} = r \cdot \frac{\sigma_{\text{speed}}}{\sigma_{\text{angle}}} = \theta_1$$

$$\text{Speed} = \theta_0 + \theta_1 \text{ angle}$$

$$\theta_1 = r \cdot \frac{\sigma_{\text{speed}}}{\sigma_{\text{angle}}} = \theta_{w,1} \cdot \frac{\sigma_{\text{speed}}^2}{\sigma_{\text{angle}}^2} = 0.019 \cdot \frac{46.678^2}{153.642^2}$$

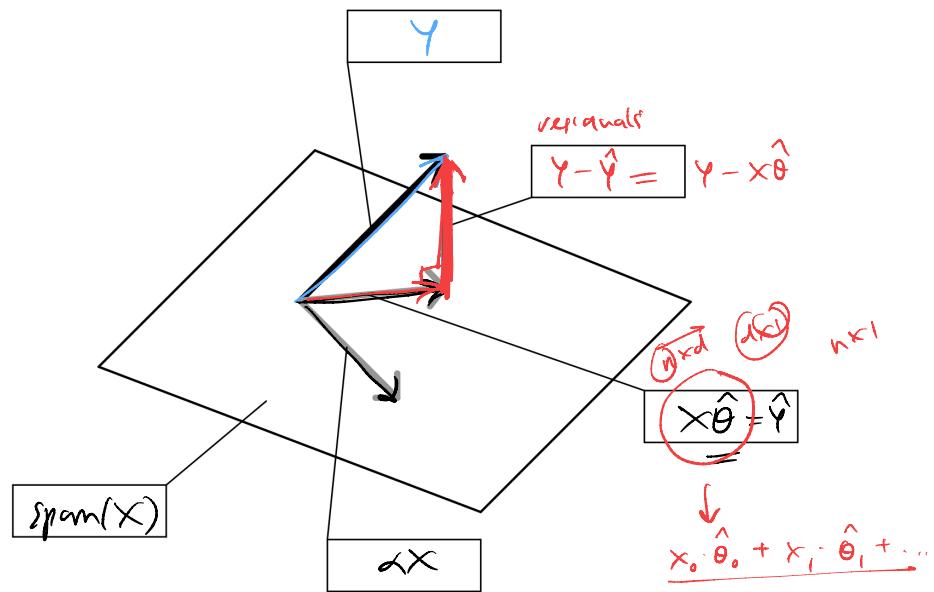
$$\theta_0 = \text{speed} - \theta_1 \text{ angle}$$

$$= 72.923 - 0.00179 \cdot 143.72 = 32.67 = \theta_0$$

## Geometry of Least Squares

2. Suppose we have a dataset represented with the design matrix  $\text{span}(X)$  and response vector  $\mathbb{Y}$ . We use linear regression to solve for this and obtain optimal weights as  $\hat{\theta}$ . Label the following terms on the geometric interpretation of ordinary least squares:

- $X$  (i.e.,  $\text{span}(X)$ )
- The prediction vector  $\mathbb{Y}\hat{\theta}$  (using optimal parameters)
- The response vector  $\mathbb{Y}$
- A prediction vector  $\mathbb{X}\alpha$  (using an arbitrary vector  $\alpha$ ).
- The residual vector  $\mathbb{Y} - \mathbb{X}\hat{\theta}$



\* Geometric way of finding OLS optimal parameter  $\hat{\theta}$

① since  $(\mathbb{Y} - \mathbb{X}\hat{\theta}) \perp \text{span}(X)$ ; dot product between  $(\mathbb{Y} - \mathbb{X}\hat{\theta})$  and  $X$  is 0:

$$\underline{x^T(\mathbb{Y} - \mathbb{X}\hat{\theta}) = 0}$$

↳ minimize norm (length) of residual vector  $\mathbb{Y} - \mathbb{X}\hat{\theta}$

↳ minimize MLE for linear model  $\hat{\mathbb{Y}} = \mathbb{X}\hat{\theta}$

②  $\cancel{x^T \mathbb{Y} = x^T \mathbb{X}\hat{\theta}}$  normal equation

↳ OLS

③  $\hat{\theta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}$

3. Using the geometry of least squares, let's answer a few questions about Ordinary Least Squares (OLS)!

- (a) Which of the following are true about the optimal solution  $\hat{\theta}$  to OLS? Recall that the least squares estimate  $\hat{\theta}$  solves the normal equation  $(\mathbf{X}^T \mathbf{X})\hat{\theta} = \mathbf{X}^T \mathbf{Y}$ .

$$\text{derivation: } \hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$\mathbf{d} \times \mathbf{d}$        $\mathbf{d} \times \mathbf{n}$        $\mathbf{n} \times \mathbf{1}$        $\mathbf{d} \times \mathbf{1}$        $\mathbf{d} \times \mathbf{n}$        $\mathbf{n} \times \mathbf{1}$

- A. Using the normal equation, we can derive an optimal solution for simple linear regression with an  $L_2$  loss.  $\rightarrow \text{MSE}$
- B. Using the normal equation, we can derive an optimal solution for simple linear regression with an  $L_1$  loss.
- C. Using the normal equation, we can derive an optimal solution for a constant model with an  $L_2$  loss.
- D. Using the normal equation, we can derive an optimal solution for a constant model with an  $L_1$  loss.
- E. Using the normal equation, we can derive an optimal solution for the model  $\hat{y} = \theta_1 x + \theta_2 \sin(x^2)$ .  $\rightarrow$   $\begin{bmatrix} \sin(x) \\ x^2 \end{bmatrix}$   $\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \rightarrow \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$
- F. Using the normal equation, we can derive an optimal solution for the model  $\hat{y} = \theta_1 x^2 + \theta_2 x^3$ .  $\text{not linear in } \theta$

- (b) Which of the following conditions are required for the least squares estimate in the previous subpart?

- A.  $\mathbf{X}$  must be full column rank.
- B.  $\mathbf{Y}$  must be full column rank.
- C.  $\mathbf{X}$  must be invertible.
- D.  $\mathbf{X}^T$  must be invertible.

$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$\mathbf{d} \times \mathbf{d}$

- for matrix  $A$  to be invertible
- ① square
  - ② full column rank (row)
  - ③ nonzero det.

- (c) What is always true about the residuals in the least squares regression? Select all that apply.

- A. They are orthogonal to the column space of the design matrix.
- B. They represent the errors of the predictions.  $\hat{y} - y$
- C. Their sum is equal to the mean squared error.
- D. Their sum is equal to zero.  $\text{not always!}$ , if  $X$  has linearly dependent columns
- E. None of the above.

SLR:  $y_i = \underline{\alpha} \underline{x}_i + \underline{\beta}$

com. dependent on  
 $y = \underline{x}^T \underline{\theta}$   
if all bias column to  $x$

$$\left[ \begin{array}{c|ccccc} x & x_{1,1} & \dots & & & \\ \hline x_{1,1} & & \dots & & & \\ x_{1,1} & & \dots & & & \\ x_{1,1} & & \dots & & & \end{array} \right]^T \left[ \begin{array}{c} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{array} \right] = \left[ \begin{array}{c} 1 \\ x_{1,1} \\ \vdots \\ x_{1,1} \end{array} \right]$$

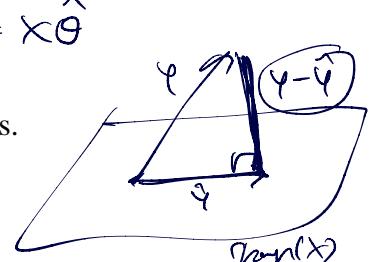
constant model

$$y_i = c \quad v_i$$

$$\underline{y} = \left[ \begin{array}{c|ccccc} 1 & x_{1,1} & \dots & & & \\ \hline x_{1,1} & & \dots & & & \\ x_{1,1} & & \dots & & & \\ x_{1,1} & & \dots & & & \end{array} \right] \left[ \begin{array}{c} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{array} \right] = \left[ \begin{array}{c} 1 \\ x_{1,1} \\ \vdots \\ x_{1,1} \end{array} \right] \left[ \begin{array}{c} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{array} \right] = \theta_0 + \theta_1 x_{1,1} + \dots$$

(d) Which of the following are true about the predictions made by OLS? Select all that apply.

- A. They are projections of the observations onto the column space of the design matrix.
- B. They are linear combinations of the features.  $\hat{Y} = X\hat{\theta}$
- C. They are orthogonal to the residuals.
- D. They are orthogonal to the column space of the features.
- E. None of the above.



(e) Which of the following is true of the mystery quantity  $\vec{v} = (I - \mathbb{X}(\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T) \mathbb{Y}$ ?

- A. The vector  $\vec{v}$  represents the residuals for any linear model.
- B. If the  $\mathbb{X}$  matrix contains the  $\vec{1}$  vector, then the sum of the elements in vector  $\vec{v}$  is 0 (i.e.  $\sum_i v_i = 0$ ).
- C. All the column vectors  $x_i$  of  $\mathbb{X}$  are orthogonal to  $\vec{v}$ .
- D. If  $\mathbb{X}$  is of shape  $n$  by  $p$ , there are  $p$  elements in vector  $\vec{v}$ .
- E. For any  $\vec{\alpha}$ ,  $\mathbb{X}\vec{\alpha}$  is orthogonal to  $\vec{v}$ .

Objective of OLS is L2 loss

$$R(\theta) = \frac{1}{n} \|y - \hat{y}\|_2^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

(dot product  $\rightarrow$  scalar)

\* Algebraic way of finding OLS' optimal parameter

OLS minimizer MSE of  $\hat{y} = X\hat{\theta}$ .

$$\begin{aligned} R(\theta) &= \frac{1}{n} \|y - X\theta\|_2^2 = \frac{1}{n} (y - X\theta)^T (y - X\theta) \\ &= \frac{1}{n} (y^T - X^T \theta^T)(y - X\theta) \quad \leftarrow \text{transposed} \\ &= \frac{1}{n} \left[ y^T y - y^T X \theta - (X \theta)^T y - (X \theta)^T X \theta \right] \quad \leftarrow \text{diff. w.r.t. } \theta \end{aligned}$$

① Find derivative w.r.t. parameter  $\theta$

$$\frac{d}{d\theta} R(\theta) = \frac{1}{n} (-2X^T y - 2X^T X \theta)$$

$\approx$  dx<sub>1</sub> vector, so result is dx<sub>1</sub> gradient

② Set derivative to 0 to find optimal parameter  $\hat{\theta}$

$$\frac{d}{d\theta} R(\theta) = 0$$

$$-\frac{2}{n} (X^T y - X^T X \hat{\theta}) = 0$$

$$X^T y = X^T X \hat{\theta}$$

$$\hat{\theta} = (X^T X)^{-1} X^T y. \quad \square$$

## \* Sum of residuals

Sum of residuals is 0 if there is a bias column in  $X$

Say we have a  $2 \times 1$  parameter  $\theta = [\theta_0 \ \theta_1]$  and  $X = [1 \ x_{i,1} \dots \ x_{i,2}]$   
 $R(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - (\theta_1 x_i + \theta_0))^2$  |  $\frac{d}{d\theta}$   
2x1  
dimensionality  
of model  
n x 2

$$\frac{d}{d\theta_0} R(\theta) = \frac{1}{n} \sum_{i=1}^n 2(y_i - \theta_0 - \theta_1 x_i) \cdot (-1)$$

$$0 = -\cancel{\frac{2}{n}} \sum_{i=1}^n (y_i - \hat{\theta}_0 - \theta_1 x_i)$$

$$\sum_{i=1}^n (y_i - \theta_1 x_i) - n \hat{\theta}_0 = 0$$

$$\hat{\theta}_0 = \frac{1}{n} \sum_{i=1}^n (y_i - \theta_1 x_i)$$



$$\sum \epsilon_i = \sum (y_i - \theta_0 - \theta_1 x_i) = 0.$$

where  $\epsilon$  is residuals ("error") vector