

## Discussion #6 Solutions

Name:

**Bias-Variance Tradeoff**

1. Let  $X$  be a random variable with mean  $\mu = \mathbb{E}[X]$ . Using the definition  $\text{Var}(X) = \mathbb{E}[(X - \mu)^2]$ , show that for any constant  $c$ ,

$$\mathbb{E}[(X - c)^2] = (\mu - c)^2 + \text{Var}(X).$$

**Solution:** One way to show this is to write  $X - c = X - \mu + \mu - c$ . Squaring both sides,

$$\mathbb{E}[(X - c)^2] = \mathbb{E}[(X - \mu)^2 + (\mu - c)^2 + 2(X - \mu)(\mu - c)]$$

Now using linearity of expectation and pulling out the constants,

$$\begin{aligned}\mathbb{E}[(X - c)^2] &= \mathbb{E}[(X - \mu)^2] + (\mu - c)^2 + 2 \underbrace{\mathbb{E}[X - \mu]}_{=0}(\mu - c) \\ &= \text{Var}(X) + (\mu - c)^2.\end{aligned}$$

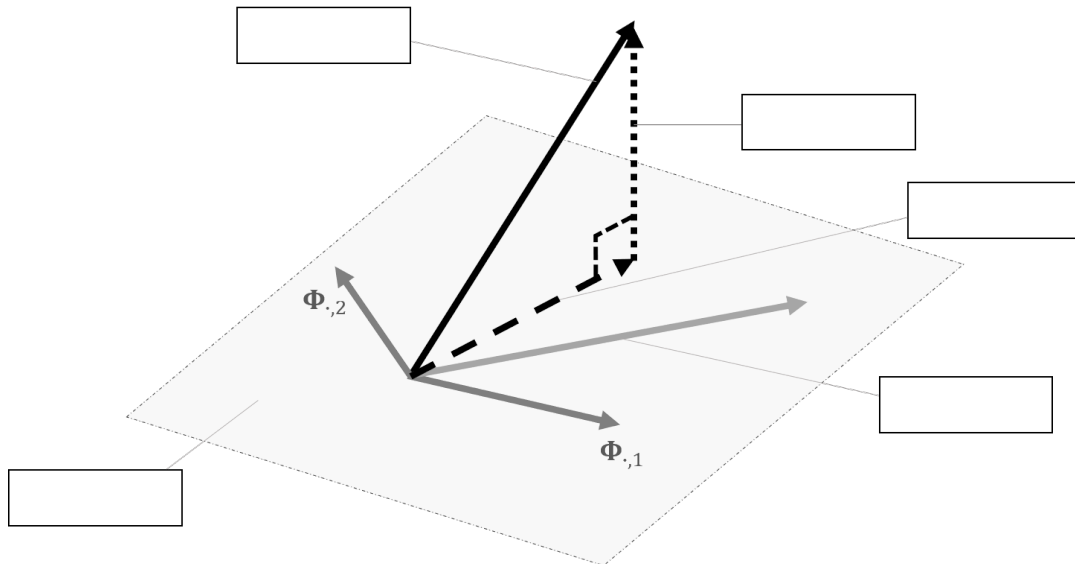
2. Use the above result to prove that

- $\text{Var}(X) \leq \mathbb{E}[(X - c)^2]$  for any  $c$
- $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$

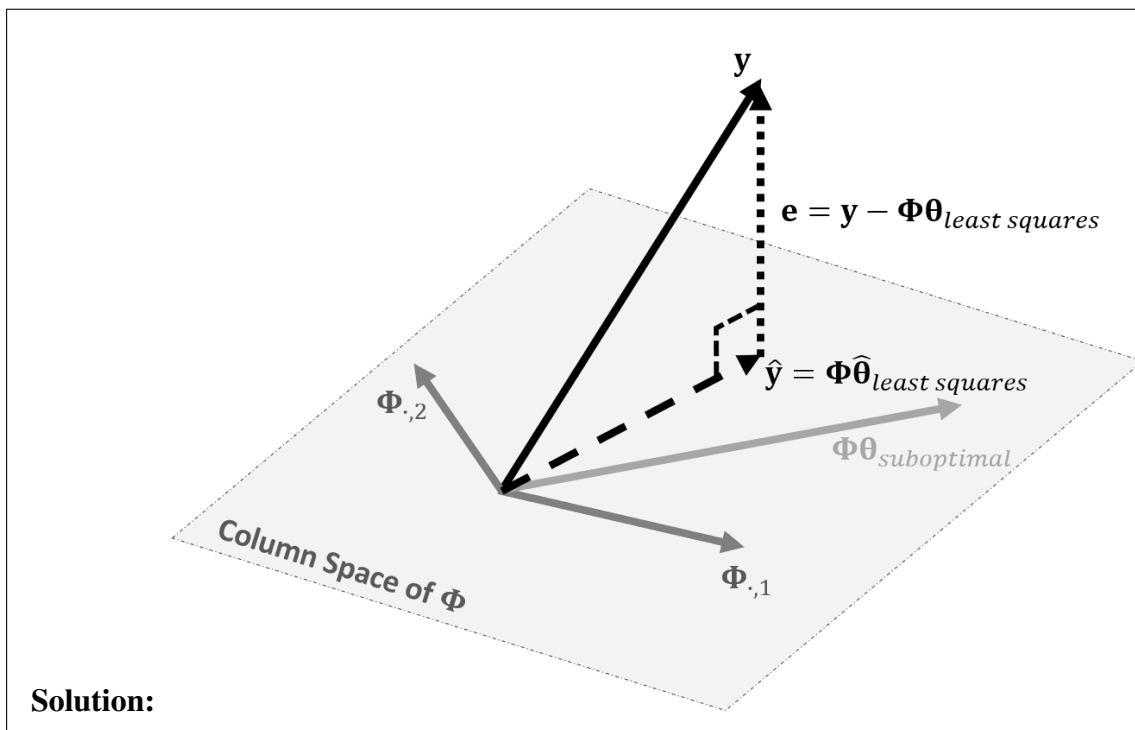
**Solution:** The first bullet follows from using  $(\mu - c)^2 \geq 0$ , and the second bullet follows from plugging in  $c = 0$ .

## Geometry of Least Squares

3. The following question will refer to the diagram below:



- (a) Fill in the diagram of the geometric interpretation of 1) the column space of the design matrix, 2) the response vector ( $y$ ), 3) the residuals and 4) the predictions



- (b) From the image above, what can we say about the residuals and the column space of  $\Phi$ ? Write this mathematically and prove this statement with a calculus-based argument and a linear-algebra-based argument.

**Solution:** We can say that the residuals are orthogonal to the column space of  $\Phi$ . Mathematically  $\Phi^T \mathbf{e} = \Phi^T (\mathbf{y} - \Phi\theta) = 0$ .

Recall that the estimator  $\hat{\theta}$  linear regression is trying to solve this problem:

$$\hat{\theta} \in \underset{\theta}{\operatorname{argmin}} (\mathbf{y} - \Phi\theta)^T (\mathbf{y} - \Phi\theta)$$

Since this is a convex function, the minimizing  $\theta$  is where the gradient is zero.

$$\begin{aligned} 0 &= \nabla \left[ (\mathbf{y} - \Phi\theta)^T (\mathbf{y} - \Phi\theta) \right] \\ &= \nabla (\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \Phi\theta - \theta^T \Phi^T \mathbf{y} + \theta^T \Phi^T \Phi\theta) \\ &= \nabla (\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \Phi\theta + \theta^T \Phi^T \Phi\theta) \\ &= 0 - 2(\mathbf{y}^T \Phi)^T + [(\Phi^T \Phi) + (\Phi^T \Phi)^T] \theta \\ &= -2\Phi^T \mathbf{y} + 2\Phi^T \Phi\theta \\ &\implies \Phi^T (\mathbf{y} - \Phi\theta) = 0 \end{aligned}$$

Alternatively, let  $V$  be a vector space equipped with an inner product  $(\cdot, \cdot)$  and  $W \subseteq V$  be a subspace of  $V$ . In the linear regression setting:  $V = \mathbb{R}^k$  for some  $k \in \mathbb{N}$ .  $W = \operatorname{colspace}(\Phi)$ . Consider a vector  $\mathbf{y} \in V$

Suppose we have  $\hat{\mathbf{y}} \in W$  such that  $\mathbf{y} - \hat{\mathbf{y}} \perp W$ . Given an arbitrary  $\mathbf{w} \in W$ , we have:

$$\|\mathbf{y} - \mathbf{w}\|^2 = \|(\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \mathbf{w})\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 + \|\hat{\mathbf{y}} - \mathbf{w}\|^2 \geq \|\mathbf{y} - \hat{\mathbf{y}}\|^2$$

Since this is true for any  $\mathbf{w}$ , we have that  $\hat{\mathbf{y}}$  is the minimizer. Furthermore, since  $\hat{\mathbf{y}} \in W$ , then there exists some  $\hat{\theta} \in \mathbb{R}^d$  such that  $\hat{\mathbf{y}} = \Phi\hat{\theta}$ .

- (c) Derive the normal equations from the fact above.

**Solution:**  $\Phi^T (\mathbf{y} - \Phi\theta) = 0$  We can distribute and rearrange terms:  $\Phi^T \mathbf{y} = \Phi^T \Phi\theta$   
We can left multiply by  $(\Phi^T \Phi)^{-1}$  in order to get the least squares estimator for  $\theta$ :  
 $\theta = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$ .

- (d) Let  $\Phi$  be a  $n \times p$  design matrix with full column rank. In this question, we will look at properties of matrix  $H = \Phi(\Phi^T \Phi)^{-1} \Phi^T$  that appears in linear regression.
- Recall for a vector space  $V$  that a projection  $\mathbf{P} : V \rightarrow V$  is a linear transformation such that  $\mathbf{P}^2 = \mathbf{P}$ . Show that  $\mathbf{H}$  is a projection matrix.

**Solution:**

$$\begin{aligned}
 \mathbf{H}^2 &= (\Phi(\Phi^T\Phi)^{-1}\Phi^T)(\Phi(\Phi^T\Phi)^{-1}\Phi^T) \\
 &= \Phi(\Phi^T\Phi)^{-1}(\Phi^T\Phi)(\Phi^T\Phi)^{-1}\Phi^T \\
 &= \Phi\mathbf{I}(\Phi^T\Phi)^{-1}\Phi^T \\
 &= \Phi(\Phi^T\Phi)^{-1}\Phi^T \\
 &= \mathbf{H}
 \end{aligned}$$

- ii. This is often called the “hat matrix” because it puts a hat on  $\mathbf{y}$ , the observed responses used to train the linear model. Show that  $\mathbf{H}\mathbf{y} = \hat{\mathbf{y}}$

**Solution:**  $\mathbf{H}\mathbf{y} = \Phi(\Phi^T\Phi)^{-1}\Phi^T\mathbf{y} = \Phi[(\Phi^T\Phi)^{-1}\Phi^T\mathbf{y}] = \Phi\hat{\boldsymbol{\theta}} = \hat{\mathbf{y}}$

- iii. Show that  $\mathbf{M} = \mathbf{I} - \mathbf{H}$  is a projection matrix.

**Solution:**  $\mathbf{M}^2 = (\mathbf{I} - \mathbf{H})^2 = \mathbf{I}^2 - \mathbf{I}\mathbf{H} - \mathbf{H}\mathbf{I} + \mathbf{H}^2 = \mathbf{I} - \mathbf{H} - \mathbf{H} + \mathbf{H} = \mathbf{I} - \mathbf{H} = \mathbf{M}$

- iv. Show that  $\mathbf{M}\mathbf{y}$  results in the residuals of the linear model.

**Solution:**  $\mathbf{M}\mathbf{y} = (\mathbf{I} - \mathbf{H})\mathbf{y} = \mathbf{y} - \mathbf{H}\mathbf{y} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{e}$

- v. Prove that  $\mathbf{H} \perp \mathbf{M}$

**Solution:**  $\mathbf{H}\mathbf{M} = \mathbf{H}(\mathbf{I} - \mathbf{H}) = \mathbf{H} - \mathbf{H}^2 = \mathbf{H} - \mathbf{H} = \mathbf{0}$

- vi. Notice that the hat matrix is a function of our observations  $\Phi$  rather than our response variable  $\mathbf{y}$ . Intuitively, what do the values in our hat matrix represent? It might be helpful to write  $\hat{y}_i$  as a summation.

**Solution:** Using the fact that  $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ , we can write the  $i$ th prediction as  $\hat{y}_i = \sum_{j=1}^n h_{ij}y_j$ . From this, we can see that the values of the hat matrix allow us to understand how much each observation influences our prediction of  $\hat{y}_i$ . Moreover,  $h_{ii}$  is a measure of the diagonal elements of  $\mathbf{H}$  denoted as  $h_{ii}$  are called leverages. The larger a point's leverage is, the more influence the point has on the regression predictions.

(e) Suppose  $\Phi \in \mathbb{R}^{n \times d}$  does not have full column rank. Then  $\Phi^T \Phi$  is not invertible. Why is that? Complete the argument below:

- i. Recall that the null space  $N(\Phi)$  of a matrix  $\Phi$  is defined as all the vectors that get sent to 0 by  $\Phi$  i.e.

$$N(\Phi) = \{\mathbf{x} \mid \Phi \mathbf{x} = \mathbf{0}\}$$

Show that the null space of  $\Phi$  is a subset of the null space of  $\Phi^T \Phi$ .

**Solution:** Let  $\mathbf{x} \in N(\Phi)$ . Then

$$\begin{aligned}\Phi \mathbf{x} &= \mathbf{0} \\ \implies \Phi^T \Phi \mathbf{x} &= \mathbf{0} \\ \implies \mathbf{x} &\in N(\Phi^T \Phi)\end{aligned}$$

- ii. Show that the reverse inclusion is also true i.e. that  $N(\Phi^T \Phi) \subseteq N(\Phi)$

**Solution:** Let  $\mathbf{x} \in N(\Phi^T \Phi)$ . Then

$$\begin{aligned}\Phi^T \Phi \mathbf{x} &= \mathbf{0} \\ \implies \mathbf{x}^T \Phi^T \Phi \mathbf{x} &= 0 \\ \implies (\Phi \mathbf{x})^T (\Phi \mathbf{x}) &= \|\Phi \mathbf{x}\|^2 = 0 \\ \implies \Phi \mathbf{x} &= \mathbf{0} \\ \implies \mathbf{x} &\in N(\Phi)\end{aligned}$$

We can then conclude that  $N(\Phi^T \Phi) = N(\Phi)$ , which implies  $\dim(N(\Phi^T \Phi)) = \dim(N(\Phi))$ . By the rank-nullity theorem,  $\text{rank}(\Phi^T \Phi) = \text{rank}(\Phi)$ . Thus if  $\text{rank}(\Phi) < d$ , then  $\text{rank}(\Phi^T \Phi) < d$ . But  $\Phi^T \Phi \in \mathbb{R}^{d \times d}$ , so there's no hope for invertibility.

- iii. List some reasons why  $\Phi$  might not have full column rank.

**Solution:**

- There are fewer observations than features i.e.  $n < d$
- Some features are linear combinations of others e.g. gross income, expenses, and net profit are all included in the design matrix

## Regularization

4. In a petri dish, yeast populations grow exponentially over time. In order to estimate the growth rate of a certain yeast, you place yeast cells in each of  $n$  petri dishes and observe the population  $y_i$  at time  $x_i$  and collect a dataset  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ . Because yeast populations are known to grow exponentially, you propose the following model:

$$\log(y_i) = \beta x_i \quad (1)$$

where  $\beta$  is the growth rate parameter (which you are trying to estimate). We will derive the  $L_2$  regularized estimator least squares estimate.

- (a) Write the *regularized least squares loss function* for  $\beta$  under this model. Use  $\lambda$  as the regularization parameter.

**Solution:**

$$L(\beta) = \frac{1}{n} \sum_{i=1}^n (\log(y_i) - \beta x_i)^2 + \lambda \beta^2 \quad (2)$$

- (b) Solve for the optimal  $\hat{\beta}$  as a function of the data and  $\lambda$ .

**Solution:** Taking the partial derivative of the regularized loss function function:

$$\begin{aligned} \frac{\partial}{\partial \beta} L(\beta) &= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \beta} (\log(y_i) - \beta x_i)^2 + \frac{\partial}{\partial \beta} \lambda \beta^2 \\ &= -\frac{2}{n} \sum_{i=1}^n (\log(y_i) - \beta x_i) x_i + 2\lambda \beta \\ &= -\frac{2}{n} \sum_{i=1}^n \log(y_i) x_i + \frac{2}{n} \sum_{i=1}^n \beta x_i^2 + 2\lambda \beta \\ &= -\frac{2}{n} \sum_{i=1}^n \log(y_i) x_i + \frac{2}{n} \beta \left( \sum_{i=1}^n x_i^2 + \lambda n \right) \end{aligned}$$

Setting the derivative equal to zero and solving for  $\beta$ :

$$0 = -\frac{2}{n} \sum_{i=1}^n \log(y_i) x_i + \frac{2\beta}{n} \left( \lambda n + \sum_{i=1}^n x_i^2 \right)$$
$$\beta \left( \lambda n + \sum_{i=1}^n x_i^2 \right) = \sum_{i=1}^n \log(y_i) x_i$$
$$\beta = \left( \lambda n + \sum_{i=1}^n x_i^2 \right)^{-1} \sum_{i=1}^n \log(y_i) x_i$$