

Discussion #6

*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 ,

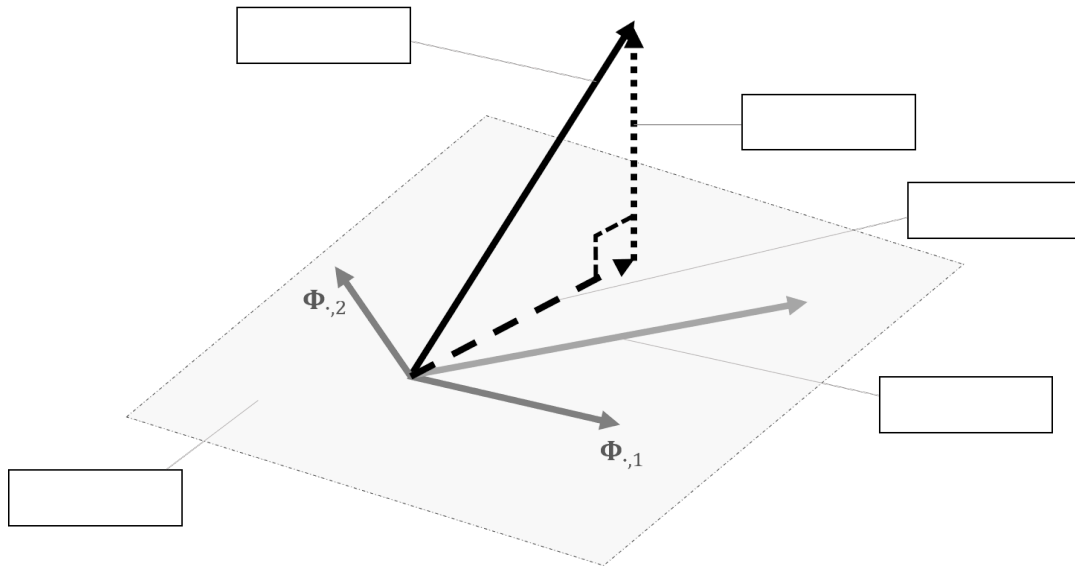
$$\mathbb{E}[(X - c)^2] = (\mu - c)^2 + \text{Var}(X).$$

2. Use the above result to prove that

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

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 (\mathbf{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 Φ ? Write this mathematically and prove this statement with a calculus-based argument and a linear-algebra-based argument.
- (c) Derive the normal equations from the fact above.

- (d) Let Φ 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.
- i. 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.
 - 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 $\mathbf{H}\mathbf{y} = \hat{\mathbf{y}}$
 - iii. Show that $\mathbf{M} = \mathbf{I} - \mathbf{H}$ is a projection matrix.
 - iv. Show that $\mathbf{M}\mathbf{y}$ results in the residuals of the linear model.
 - v. Prove that $\mathbf{H} \perp \mathbf{M}$
 - vi. Notice that the hat matrix is a function of our observations Φ 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.

(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 Φ is defined as all the vectors that get sent to 0 by Φ i.e.

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

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

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

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 Φ might not have full column rank.

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 β 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 β under this model. Use λ as the regularization parameter.

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