

## Discussion # 10 Solutions

Name:

## Bias-Variance Trade-off

1. Assume that we have a function  $h(x)$  and some noise generation process that produces  $\epsilon$  such that  $\mathbb{E}[\epsilon] = 0$  and  $\text{var}(\epsilon) = \sigma^2$ . Every time we query mother nature for  $Y$  at a given  $x$ , she gives us  $Y = h(x) + \epsilon$ . A new  $\epsilon$  is generated each time, independent of the last. We randomly sample some data  $(x_i, y_i)_{i=1}^n$  and use it to fit a model  $f_{\hat{\beta}}(x)$  according to some procedure (e.g. OLS, Ridge, LASSO). In class, we showed that

$$\underbrace{\mathbb{E}[(Y - f_{\hat{\beta}}(x))^2]}_{\text{model risk}} = \underbrace{\sigma^2}_{\text{observation variance}} + \underbrace{(h(x) - \mathbb{E}[f_{\hat{\beta}}(x)])^2}_{\text{model bias}^2} + \underbrace{\mathbb{E}[(\mathbb{E}[f_{\hat{\beta}}(x)] - f_{\hat{\beta}}(x))^2]}_{\text{model variance}}.$$

- (a) Label each of the terms above. Word bank: observation variance, model variance, observation bias<sup>2</sup>, model bias<sup>2</sup>, model risk, empirical mean square error.

**Solution:**

$$\underbrace{\mathbb{E}[(Y - f_{\hat{\beta}}(x))^2]}_{\text{model risk}} = \underbrace{\sigma^2}_{\text{observation variance}} + \underbrace{(h(x) - \mathbb{E}[f_{\hat{\beta}}(x)])^2}_{\text{model bias}^2} + \underbrace{\mathbb{E}[(\mathbb{E}[f_{\hat{\beta}}(x)] - f_{\hat{\beta}}(x))^2]}_{\text{model variance}}$$

- (b) What is random in the equation above? Where does the randomness come from?

**Solution:**  $Y$  - this is the new observation at  $x$ . Its randomness comes from the noise generation process.  $f_{\hat{\beta}}$  - this is the model fitted from the data. Its randomness comes from sampling and the noise generation process.

- (c) True or false and explain.  $\mathbb{E}[\epsilon f_{\hat{\beta}}(x)] = 0$

**Solution:** True. Since  $\epsilon$  and  $\hat{\beta}$  are independent,

$$\mathbb{E}[\epsilon f_{\hat{\beta}}(x)] = \mathbb{E}[\epsilon] \mathbb{E}[f_{\hat{\beta}}(x)] = 0$$

- (d) Suppose you lived in a world where you could collect as many data sets you would like. Given a fixed algorithm to fit a model  $f_\beta$  to your data e.g. linear regression, describe a procedure to get good estimates of  $\mathbb{E} [f_{\hat{\beta}}(x)]$  (technical point: you may assume this expectation exists).

**Solution:**

- Pick an  $x$
- Gather a data set  $\mathcal{D}_i$
- Fit a model  $f_{\hat{\beta}_i}$  to that data set
- Calculate  $f_{\hat{\beta}_i}(x)$
- Repeat many times
- Average over all the  $f_{\hat{\beta}_i}(x)$

- (e) If you could collect as many data sets as you would like, how does that affect the quality of your model  $f_\beta(x)$ ?

**Solution:** By collecting many data sets, we have an unbiased estimate of the “average” model, but this does not mean our model will have unbiased prediction.

## Ridge and LASSO Regression

2. Earlier, we posed the linear regression problem as follows: Find the  $\vec{\beta}$  value that minimizes the average squared loss. In other words, our goal is to find  $\vec{\hat{\beta}}$  that satisfies the equation below:

$$\vec{\hat{\beta}} = \underset{\vec{\beta}}{\operatorname{argmin}} L(\vec{\beta}) = \underset{\vec{\beta}}{\operatorname{argmin}} \frac{1}{n} \|\vec{y} - \mathbb{X}\vec{\beta}\|_2^2$$

Here,  $\mathbb{X}$  is a  $n \times d$  matrix,  $\vec{\beta}$  is a  $d \times 1$  vector and  $\vec{y}$  is a  $n \times 1$  vector. As we saw in lecture and in last week’s discussion, the optimal  $\vec{\hat{\beta}}$  is given by the closed form expression  $\vec{\hat{\beta}} = (\mathbb{X}^t \mathbb{X})^{-1} \mathbb{X}^t \vec{y}$ .

To prevent overfitting, we saw that we can instead minimize the sum of the average squared loss plus a regularization function  $\lambda \mathcal{S}(\vec{\beta})$ . If we use the function  $\mathcal{S}(\vec{\beta}) = \|\vec{\beta}\|_2^2$ , we have “ridge regression”. If we use the function  $\mathcal{S}(\vec{\beta}) = \|\vec{\beta}\|_1$ , we have “LASSO regression”. For example, if we choose  $\mathcal{S}(\vec{\beta}) = \|\vec{\beta}\|_2^2$ , our goal is to find  $\vec{\hat{\beta}}$  that satisfies the equation below:

$$\hat{\vec{\beta}} = \underset{\vec{\beta}}{\operatorname{argmin}} L(\vec{\beta}) = \underset{\vec{\beta}}{\operatorname{argmin}} \frac{1}{n} \|\vec{y} - \mathbb{X}\vec{\beta}\|_2^2 + \lambda \|\vec{\beta}\|_2^2 = \underset{\vec{\beta}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbb{X}_{i,\cdot}^T \vec{\beta})^2 + \lambda \sum_{j=1}^d \beta_j^2$$

Recall that  $\lambda$  is a hyperparameter that determines the impact of the regularization term. Though we did not discuss this in lecture, we can also find a closed form solution to ridge regression:  $\vec{\hat{\beta}} = (\mathbb{X}^T \mathbb{X} + n\lambda \mathbf{I})^{-1} \mathbb{X}^T \vec{y}$ . It turns out that  $\mathbb{X}^T \mathbb{X} + n\lambda \mathbf{I}$  is guaranteed to be invertible (unlike  $\mathbb{X}^T \mathbb{X}$  which might not be invertible).

- (a) As model complexity increases, what happens to the bias and variance of the model?

**Solution:** Model complexity is inversely related to the regularization parameter  $\lambda$ . As  $\lambda$  increases, Bias tends to increase and variance tends to decrease.

- (b) In terms of bias and variance, how does a regularized model compare to ordinary least squares regression?

**Solution:** Regularized regression has higher bias and lower variance relative to ordinary least squares regression.

- (c) In ridge regression, what happens if we set  $\lambda = 0$ ? What happens as  $\lambda$  approaches  $\infty$ ?

**Solution:** If we set  $\lambda = 0$  we end up with OLS. As  $\lambda$  approaches  $\infty$  then  $\vec{\beta}$  goes to 0.

- (d) How does model complexity compare between ridge regression and ordinary least squares regression? How does this change for large and small values of  $\lambda$ ?

**Solution:** Ridge regression in general will result in simpler models, as we penalize for large components in  $\vec{\beta}$ .  $\lambda$  is inversely related to model complexity, e.g. larger values of  $\lambda$  represent larger penalties, meaning even lower model complexity.

- (e) If we have a large number of features (10,000+) and we suspect that only a handful of features are useful, which type of regression (Lasso vs Ridge) would be more helpful in interpreting useful features?

**Solution:** LASSO would be better as it sets many values to 0, so it would be effectively selecting useful features and “ignoring” bad ones.

- (f) What are the benefits of using ridge regression?

**Solution:** If  $\mathbf{X}^T\mathbf{X}$  is not full rank (not invertible), then we end up with infinitely many solutions for least squares. But if we use ridge regression,  $\hat{\beta} = (\mathbf{X}^T\mathbf{X} + n\lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{Y}$ . This guarantees invertibility and a unique solution, for  $\lambda > 0$ .

## Cross Validation

3. After running 5-fold cross validation, we get the following mean squared errors for each fold and value of  $\lambda$ :

| Fold Num | $\lambda = 0.1$ | $\lambda = 0.2$ | $\lambda = 0.3$ | $\lambda = 0.4$ | Row Avg |
|----------|-----------------|-----------------|-----------------|-----------------|---------|
| 1        | 80.2            | 70.2            | 91.2            | 91.8            | 83.4    |
| 2        | 76.8            | 66.8            | 88.8            | 98.8            | 82.8    |
| 3        | 81.5            | 71.5            | 86.5            | 88.5            | 82.0    |
| 4        | 79.4            | 68.4            | 92.3            | 92.4            | 83.1    |
| 5        | 77.3            | 67.3            | 93.4            | 94.3            | 83.0    |
| Col Avg  | 79.0            | 68.8            | 90.4            | 93.2            |         |

How do we use the information above to choose our model? Do we pick a specific fold? a specific lambda? or a specific fold-lambda pair? Explain.

**Solution:** We should use  $\lambda = 0.2$  because this value has the least average MSE across all folds.

4. You build a model with two regularization hyperparameters  $\lambda$  and  $\gamma$ . You have 4 good candidate values for  $\lambda$  and 3 possible values for  $\gamma$ , and you are wondering which  $\lambda, \gamma$  pair will be the best choice. If you were to perform five-fold cross-validation, how many validation errors would you need to calculate?

**Solution:** There are  $4 \times 3 = 12$  pairs of  $\lambda, \gamma$  and each pair will have 5 validation errors, one for each fold.

5. In the typical setup of k-fold cross validation, we use a different parameter value on each fold, compute the mean squared error of each fold and choose the parameter whose fold has the lowest loss.

- ☐ A. True  
☒ B. False