

## Discussion #7

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  $f_{\hat{\theta}}(x)$  according to some procedure (e.g. OLS, Ridge, LASSO). In class, we showed that

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

- (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.
- (b) What is random in the equation above? Where does the randomness come from?
- (c) True or false and explain.  $\mathbb{E}[\epsilon f_{\hat{\theta}}(x)] = 0$
- (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  $f_{\theta}$  to your data e.g. linear regression, describe a procedure to get good estimates of  $\mathbb{E}[f_{\hat{\theta}}(x)]$  (technical point: you may assume this expectation exists).
- (e) If you could collect as many data sets as you would like, how does that affect the quality of your model  $f_{\theta}(x)$ ?

2. We find the optimal  $\theta$  that minimizes squared loss with  $L_1$  regularization:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \|\mathbf{y} - \Phi\theta\|_2^2 + \lambda \|\theta\|_1 = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - \Phi_{i,\cdot}^T \theta)^2 + \lambda \sum_{j=1}^d |\theta_j|.$$

You receive all your (nice, clean, numerical) data in a single `DataFrame` called `CompleteData`. The first column contains the responses and the remaining  $d$  columns hold the features. You want to use 60% of your data in the training set and implement part of 5-fold cross-validation with the following pseudocode:

```
Phi_train, Y_train, Phi_test, Y_test = \
    make_train_test_split(CompleteData, 0.60)
lambdas = make_lambdas(from=0.1, to=0.4, by=0.1)

n = count_rows(...)
fold_size = n / k
idx = range(n)
randomly shuffle the ordering of idx
folds = [idx[i * fold_size : (i+1) * fold_size] for i in range(k)]

for i, fold in enumerate(folds):
    for j, lam in enumerate(lambdas):
        mse[i, j] = calculate_mse_lasso(Phi__, Y__, fold, lam)
```

(a) What should the ... be in `count_rows` above? Your choices are `CompleteData`, `Phi_train`, and `Phi_test`.

(b) What should the blanks be in `calculate_mse_lasso` above? Your choices are `train` and `test`.

(c) Describe an algorithm for `calculate_mse_lasso`.

- (d) 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.

## Ridge Regression

3. Ridge regression is a variant of least squares that involves regularization. The problem is stated as follows:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} L(\theta) = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \|\mathbf{y} - \Phi\theta\|_2^2 + \lambda \|\theta\|_2^2 = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - \Phi_{i,\cdot}^T \theta)^2 + \lambda \sum_{j=1}^d \theta_j^2$$

Here,  $\lambda$  is a hyperparameter that determines the impact of the regularization term.  $\Phi$  is a  $n \times d$  matrix,  $\theta$  is a  $d \times 1$  vector and  $\mathbf{y}$  is a  $n \times 1$  vector. The optimal choice is  $\hat{\theta} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y}$ .

- As model complexity increases, what happens to the bias and variance of the model?
- In terms of bias and variance, how does the a regularized regression estimator compare to ordinary least squares regression?
- In ridge regression, what happens if we set  $\lambda = 0$ ? What happens as  $\lambda$  approaches  $\infty$ ?
- How does model complexity compare between ridge regression and ordinary least squares regression? How does this change for large and small values of  $\lambda$ ?

- (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?
- (f) What are the benefits of using ridge regression?
- (g) On last week's discussion, we discussed possible situations where the matrix  $\Phi^T \Phi$  was not invertible, such as the presence of linearly dependent columns or an insufficient number of observations. In this question, we will demonstrate that the  $L_2$  regularization penalty always ensures that the matrix  $(\Phi^T \Phi + \lambda \mathbf{I})^{-1}$  is invertible, resulting in a unique solution.
- i. A symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is **positive semi-definite** if for every non-zero vector  $\mathbf{v} \in \mathbb{R}^n$ , we have  $\mathbf{v}^T \mathbf{A} \mathbf{v} \geq 0$ . Given a matrix  $\Phi \in \mathbb{R}^{n \times d}$  (think our feature matrix), show that  $\Phi^T \Phi$  is positive semi-definite.
  - ii. A symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is **positive definite** if for every non-zero vector  $\mathbf{v} \in \mathbb{R}^n$ , we have  $\mathbf{v}^T \mathbf{A} \mathbf{v} > 0$ . Notice that the inequality is now strict. Given a matrix  $\Phi \in \mathbb{R}^{n \times d}$  (think our feature matrix) and  $\lambda > 0$  (our regularization hyperparameter), show that  $\Phi^T \Phi + \lambda \mathbf{I}$  is positive definite.
  - iii. Prove that a positive-definite matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is always invertible by showing that its null space only contains the zero vector i.e.

$$N(\mathbf{A}) = \{\mathbf{0}\}$$