DS 100/200: Principles and Techniques of Data Science

Date: April 10, 2020

Discussion #10

Name:

## **Bias-Variance Trade-Off**

- 1. Your team would like to train a machine learning model in order to predict the next YouTube video that a user will click on based on the videos the user has watched in the past. We extract m attributes (such as length of video, view count etc) from each video and our model will be based on the previous d videos watched by that user. Hence the number of features for each data point for the model is  $m \cdot d$ . You're not sure how many videos to consider.
  - (a) Your colleague generates the following plot, where the value d is on the x axis. However, they forgot to label the y-axis.



Which of the following could the y axis represent? Select all that apply.

- A Training Error
- **NB** Validation Error
- C. Bias
- $\Box$  D. Variance
- (b) Your colleague generates the following plot, where the value d is on the x axis. However, they forgot to label the y-axis.



Which of the following could the y axis represent? Select all that apply.

- A Training Error
   B. Validation Error
   C. Bias
   D. Variance
- We randomly sample some data (x<sub>i</sub>, y<sub>i</sub>)<sup>n</sup><sub>i=1</sub> and use it to fit a model f<sub>θ</sub>(x) according to some procedure (e.g. OLS, Ridge, LASSO). We then sample a new point that is independent from our existing points, but sampled from the same underlying truth as our data. Furthermore, assume that we have a function g(x) and some noise generation process that produces ε such that E [ε] = 0 and var(ε) = σ<sup>2</sup>. Every time we query mother nature for Y at a given a x, she gives us Y = g(x) + ε. (The true function for our data is Y = g(x) + ε.) A new ε is generated each time, independent of the last. In class, we showed that

$$\underbrace{\mathbb{E}\left[(Y - f_{\hat{\theta}}(x))^2\right]}_{\text{obs. variance}} = \underbrace{\sigma^2}_{\theta \text{ ins}^2} + \underbrace{(g(x) - \mathbb{E}\left[f_{\hat{\theta}}(x)\right])^2}_{\theta \text{ ins}^2} + \underbrace{\mathbb{E}\left[(f_{\hat{\theta}}(x) - \mathbb{E}\left[f_{\hat{\theta}}(x)\right])^2\right]}_{\text{model variance}}$$

- (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?

Y: depends on 
$$\in$$
  
 $f_{\hat{\theta}}(x)$ : random  $-$  no ise m Y (from  $\in$ )  
 $-$  generated from  $\hat{\sigma}$   
sample of Ys  
Thus an follow and emploin  $\mathbb{P}[f_{\hat{\sigma}}(x)] = 0$ 

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

True : 
$$\in$$
 and  $f_{\hat{o}}(x)$  are independent.  
For independent RVs,  $E[XY] = E[X] \in [D]$ .  
Then,  $E[ef_{\hat{o}}(x)] = E[e] \in [f_{\hat{o}}(x)] = O \cdot E[f_{\hat{o}}(x)] = O$ .

(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  $\mathbb{E}[f_{\hat{\theta}}(x)]$ 

repeat as many times as possible: \_\_\_\_\_\_ Then, take overage -sample data -fit model models

- (e) If you could collect as many data sets as you would like, how does that affect the quality of your model f<sub>θ</sub>(x)? → good estimate of E [f<sub>θ</sub>(x)]
   but if you chose a poor model to begin with, doesn't if you with, doesn't if you it
  - 3. Earlier, we posed the linear regression problem as follows: Find the  $\vec{\theta}$  value that minimizes the average squared loss. In other words, our goal is to find  $\vec{\theta}$  that satisfies the equation below:

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

Here,  $\mathbb{X}$  is a  $n \times d$  matrix,  $\vec{\theta}$  is a  $d \times 1$  vector and  $\vec{y}$  is a  $n \times 1$  vector. As we saw in lecture, the optimal  $\vec{\theta}$  is given by the closed form expression  $\vec{\theta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^t \vec{y}$ .

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

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

Recall that  $\alpha$  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:  $\vec{\theta} = (\mathbb{X}^T \mathbb{X} + n\alpha \mathbf{I})^{-1} \mathbb{X}^T \vec{y}$ . It turns out that  $\mathbb{X}^T \mathbb{X} + n\alpha \mathbf{I}$  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?

Discussion #10

 $fl \cdot c = 1$ 

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

decreased complexity

(c) In ridge regression, what happens if we set  $\alpha = 0$ ? What happens as  $\alpha$  approaches  $\infty$ ? d:o: same as OLS

```
\alpha \rightarrow \infty : A \rightarrow 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  $\alpha$ ?

## answered above

(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?

```
LASSO : encourages spensity
```

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

```
guaranteed solution, more general model
```

## **Random Variables**

4. The average response time for a question on Piazza this semester was 11 minutes. As always, the number of questions answered by each TA is highly variable, with a few TAs going above and beyond the call of duty. Below are the number of contributions for the top four TAs (out of 20,000 total Piazza contributions):

TA	# contributions
Daniel	2000
Suraj	1800
Manana	700
Allen	500

Suppose we take a sample with replacement of size n = 500 contributions from the original 20,000 contributions. We will also define some random variables:

•  $D_i = 1$  when the *i*<sup>th</sup> contribution in our sample is made by Daniel; else  $D_i = 0$ .

- $S_i = 1$  when the *i*<sup>th</sup> contribution in our sample is made by Suraj; else  $S_i = 0$ .
- $M_i = 1$  when the *i*<sup>th</sup> contribution in our sample is made by Manana; else  $M_i = 0$ .
- $A_i = 1$  when the *i*<sup>th</sup> contribution in our sample is made by Allen; else  $A_i = 0$ .
- $O_i = 1$  when the  $i^{\text{th}}$  contribution is made by anyone other than Daniel, Suraj, Manana,

or Aller; else, 
$$O_i = 0$$
  
(a) i. What is  $P(A_1 = 1)$ ?  
 $P(A_1 = 1) = \underbrace{500}_{20000}$   
ii. What is  $\mathbb{E}[S_1]$ ?  
 $\mathbb{E}[S_1] = \underbrace{100}_{20000}$   
iii. What is  $\mathbb{E}[M_{100}]$ ?  
 $\mathbb{E}[M_{100}] = \underbrace{100}_{20000}$   
iv. What is  $Var[D_{50}]$ ?  
 $Var[D_{50}] = \underbrace{10}_{0} \cdot (1 - \frac{1}{10})$   
v. What is  $Var[D_{50}]$ ?  
 $Var[D_{50}] = \underbrace{10}_{0} \cdot (1 - \frac{1}{10})$   
v. What is  $Var(D_{50}]$ ?  
 $Var[D_{50}] = \underbrace{10}_{0} \cdot (1 - \frac{1}{10})$   
v. What is  $Var(D_{50}]$ ?  
 $N_{5} = \sum_{i=1}^{500} D_{i}$   
 $N_{5} = \sum_{i=0}^{500} A_{i}$   
 $N_{6} = \sum_{i=1}^{500} \sum_$ 

20000

2

$$\operatorname{Var}(N_D + N_S + N_A + N_M + N_O) =$$

(c) Now, suppose we take a sample with replacement of 20 contributions, what is the probability that 7 were by Daniel?

distribution Probability =

- $y = \begin{pmatrix} 20 \\ 7 \end{pmatrix} \left(\frac{1}{10}\right)^{7} \left(\frac{9}{10}\right)^{13}$
- $p = \frac{1}{10}$

20000

(d) Finally, suppose we take a sample with replacement of 10 contributions. What is the probability that 3 were by Daniel, 3 were by Suraj, and 4 were by Manana? (Note: Refer to Lecture 2 to refresh your knowledge on how to calculate this type of probability)

Probability 
$$\frac{10!}{3!3!4!} (\frac{1}{10})^3 (\frac{1800}{20000}) (\frac{700}{20000})^{-1}$$
  
 $\frac{10!}{3!3!4!} = (\frac{1800}{3}) (\frac{7}{20000})^{-1}$   
 $\frac{10!}{3!3!4!} = (\frac{10}{3}) (\frac{7}{3})$   
 $\frac{10!}{7} = \frac{10}{7}$   
 $\frac{10!}{7} = \frac{10}{7}$   

6