

Discussion #9

Cross Validation

1. After running 5-fold cross validation, we get the following mean squared errors for each fold and value of λ when using Ridge regularization:

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 λ ? or a specific fold- λ pair? Explain.

2. 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

Guessing at Random

3. A multiple choice test has 100 questions, each with five possible answers of which one is right. The grading scheme is as follows:
- 4 points are awarded for each right answer.
 - For each other answer (wrong, missing, etc), one point is taken off; that is, -1 points are awarded.

A student hasn't studied at all and therefore guesses each answer uniformly at random, independently of all the other answers.

Define the following random variables:

- R : the number of answers the student gets right
- W : the number of answers the student does not get right
- S : the student's score on the test

We analyze the random variable R , which denotes the number of answers the student got right.

- (a) What is the distribution of R ? Provide the name and parameters of the appropriate distribution. Explain your answer.

- (b) Find $\mathbb{E}(R)$.

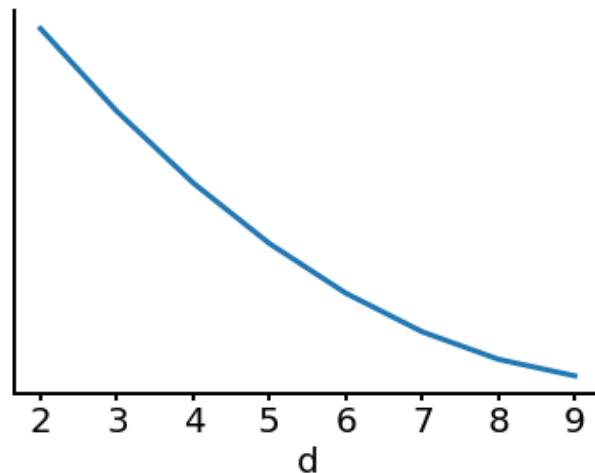
4. True or False: $\mathbb{SD}(R) = \mathbb{SD}(W)$.
5. Find $\mathbb{E}(S)$, the student's expected score on the test.
6. Find $\mathbb{SD}(S)$.

Bias-Variance Trade-Off

7. 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.

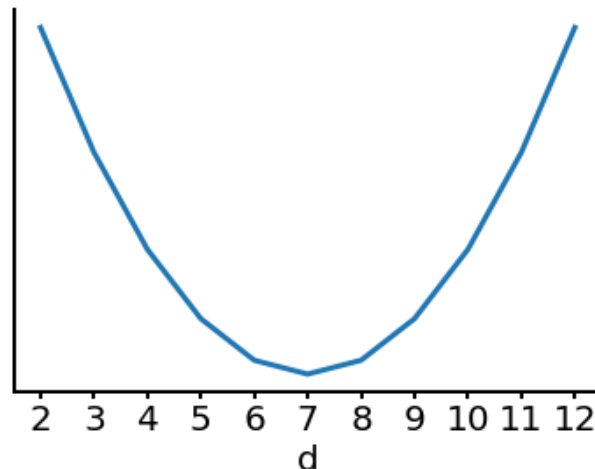
Hence the number of features for each data point for the model is $m \cdot d$. Currently, you're not sure how many videos to consider.

- (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
 - ☐ B. Validation Error
 - ☐ C. Bias
 - ☐ 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 again.



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

- ☐ A. Training Error
- ☐ B. Validation Error
- ☐ C. Bias
- ☐ D. Variance

8. 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). 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 $\mathbb{E}[\epsilon] = 0$ and $\text{var}(\epsilon) = \sigma^2$. Every time we query mother nature for Y at a given x , she gives us $Y = g(x) + \epsilon$. (The true function for our data is $Y = g(x) + \epsilon$.) A new ϵ is generated each time, independent of the last. In class, we showed that

$$\underbrace{\mathbb{E}[(Y - f_{\hat{\theta}}(x))^2]}_{\text{Empirical Mean Square Error}} = \underbrace{\sigma^2}_{\text{Observation Variance}} + \underbrace{(g(x) - \mathbb{E}[f_{\hat{\theta}}(x)])^2}_{\text{Model Bias}^2} + \underbrace{\mathbb{E}[(f_{\hat{\theta}}(x) - \mathbb{E}[f_{\hat{\theta}}(x)])^2]}_{\text{Model Variance}}$$

- (a) Label each of the terms above.

Word Bank: observation variance, model variance, observation bias², model bias², model risk, empirical mean square error.

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

- (c) Calculate the value of $\mathbb{E}[\epsilon f_{\hat{\theta}}(x)]$.