



✓ 恭喜！您通过了！

通过条件 75% 或更高

坚持学习

成绩
100%

作業四

最新提交作业的评分

100%

1. Deterministic noise depends on \mathcal{H} , as some models approximate f better than others. Assume $\mathcal{H}' \subset \mathcal{H}$ and that f is fixed. In general (but not necessarily in all cases), if we use \mathcal{H}' instead of \mathcal{H} , how does deterministic noise behave?

10/10 分

✓ Correct

2. Consider the following hypothesis set for $\mathbf{x} \in \mathbb{R}^d$ defined by the constraint:

$$\mathcal{H}(d, d_0) = \{h \mid h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}; w_i = 0 \text{ for } i \geq d_0\},$$

which of the following statements is correct?

10/10 分

✓ Correct

3. For Questions 3-4, consider the augmented error $E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$ with some $\lambda > 0$. If we want to minimize the augmented error $E_{\text{aug}}(\mathbf{w})$ by gradient descent with η as learning rate, which of the following is a correct update rule?

10/10 分

✓ Correct

4. Let \mathbf{w}_{lin} be the optimal solution for the plain-vanilla linear regression and $\mathbf{w}_{\text{reg}}(\lambda)$ be the optimal solution for minimizing E_{aug} in Question 3, with E_{in} being the squared error for linear regression. Which of the following is correct?

10/10 分

✓ Correct

5. You are given the data points: $(-1, 0), (\rho, 1), (1, 0)$, $\rho \geq 0$, and a choice between two models:

- constant $h_0(x) = b_0$ and
- linear $h_1(x) = a_1 x + b_1$.

For which value of ρ would the two models be tied using leave-one-out cross-validation with the squared error measure?

10/10 分

✓ Correct

6. For Questions 6-7, suppose that for 5 weeks in a row, a letter arrives in the mail that predicts the outcome of the upcoming Monday night baseball game. Assume there are no tie. You keenly watch each Monday and to your surprise, the prediction is correct each time. On the day after the fifth game, a letter arrives, stating that if you wish to see next week's prediction, a payment of NTD 1000 is required. Which of the following statement is true?

10/10 分

✓ Correct

7. If the cost of printing and mailing out each letter is NTD 10. If the sender sends the minimum number of letters out, how much money can be made for the above "fraud" to succeed once? That is, one of the recipients does send him NTD 1000 to receive the prediction of the 6-th game?

10/10 分

✓ Correct

8. For Questions 8-10, please read the following story first. In our credit card example, the bank starts with some vague idea of what constitutes a good credit risk. So, as customers $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ arrive, the bank applies its vague idea to approve credit cards for some of these customers based on a formula $a(\mathbf{x})$. Then, only those who get credit cards are monitored to see if they default or not.

10/10 分

For simplicity, suppose that the first $N = 10000$ customers were given credit cards by the credit approval function $a(\mathbf{x})$. Now that the bank knows the behavior of these customers, it comes to you to improve their algorithm for approving credit. The bank gives you the data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$. Before you look at the data, you do mathematical derivations and come up with a credit approval function. You now test it on the data and, to your delight, obtain perfect prediction.

What is M , the size of your hypothesis set?

✓ Correct

9. With such an M , what does the Hoeffding bound say about the probability that the true average error rate of g is worse than 1% for $N = 10,000$?

10/10 分

✓ Correct

10. You assure the bank that you have got a system g for approving credit cards for new customers, which is nearly error-free. Your confidence is given by your answer to the previous question. The bank is thrilled and uses your g to approve credit for new customers. To their dismay, more than half their credit cards are being defaulted on. Assume that the customers that were sent to the old credit approval function and the customers that were sent to your g are indeed i.i.d. from the same distribution, and the bank is lucky enough (so the "bad luck" that "the true error of g is worse than 1%" does not happen). Which of the following claim is true?

10/10 分

✓ Correct

11. For Questions 11-12, consider linear regression with virtual examples. That is, we add K virtual examples $(\tilde{\mathbf{x}}_1, \tilde{y}_1), (\tilde{\mathbf{x}}_2, \tilde{y}_2), \dots, (\tilde{\mathbf{x}}_K, \tilde{y}_K)$ to the training data set, and solve

10/10 分

$$\min_{\mathbf{w}} \frac{1}{N+K} \left(\sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \sum_{k=1}^K (\tilde{y}_k - \mathbf{w}^T \tilde{\mathbf{x}}_k)^2 \right).$$

We will show that using some "special" virtual examples, which were claimed to be a possible way to combat overfitting in Lecture 9, is related to regularization, another possible way to combat overfitting discussed in Lecture 10. Let $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1 \tilde{\mathbf{x}}_2 \dots \tilde{\mathbf{x}}_K]^T$, and $\tilde{\mathbf{y}} = [\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_K]^T$.

What is the optimal \mathbf{w} to the optimization problem above, assuming that all the inversions exist?

✓ Correct

12. For what $\tilde{\mathbf{X}}$ and $\tilde{\mathbf{y}}$ will the solution of the linear regression problem above equal to

10/10 分

$$\mathbf{w}_{\text{reg}} = \operatorname{argmin}_{\mathbf{w}} \frac{\lambda}{N} \|\mathbf{w}\|^2 + \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2?$$

✓ Correct

✓ Correct

13. Consider regularized linear regression (also called ridge regression) for classification

10/10 分

$$\mathbf{w}_{\text{reg}} = \operatorname{argmin}_{\mathbf{w}} \left(\frac{\lambda}{N} \|\mathbf{w}\|^2 + \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \right).$$

Run the algorithm on the following data set as \mathcal{D} :

https://www.csie.ntu.edu.tw/~htlin/mooc/datasets/mlfound_algo/hw4_train.dat

and the following set for evaluating E_{out} :

https://www.csie.ntu.edu.tw/~htlin/mooc/datasets/mlfound_algo/hw4_test.dat

Because the data sets are for classification, please consider only the 0/1 error for all Questions below.

Let $\lambda = 10$, which of the followings is the corresponding E_{in} and E_{out} ?

✓ Correct

14. Following the previous Question, among $\log_{10} \lambda = \{2, 1, 0, -1, \dots, -8, -9, -10\}$. What is the λ with the minimum E_{in} ? Compute λ and its corresponding E_{in} and E_{out} then select the closest answer. Break the tie by selecting the largest λ .

10/10 分

✓ Correct

15. Following the previous Question, among $\log_{10} \lambda = \{2, 1, 0, -1, \dots, -8, -9, -10\}$. What is the λ with the minimum E_{out} ? Compute λ and the corresponding E_{in} and E_{out} then select the closest answer. Break the tie by selecting the largest λ .

10/10 分

✓ Correct

16. Now split the given training examples in \mathcal{D} to the first 120 examples for $\mathcal{D}_{\text{train}}$ and 80 for \mathcal{D}_{val} . Ideally, you should randomly do the 120/80 split. Because the given examples are already randomly permuted, however, we would use a fixed split for the purpose of this problem.

10/10 分

Run the algorithm on $\mathcal{D}_{\text{train}}$ to get g_{λ}^{-} , and validate g_{λ}^{-} with \mathcal{D}_{val} . Among $\log_{10} \lambda = \{2, 1, 0, -1, \dots, -8, -9, -10\}$. What is the λ with the minimum $E_{\text{train}}(g_{\lambda}^{-})$? Compute λ and the corresponding $E_{\text{train}}(g_{\lambda}^{-})$, $E_{\text{val}}(g_{\lambda}^{-})$ and $E_{\text{out}}(g_{\lambda}^{-})$ then select the closet answer. Break the tie by selecting the largest λ .

✓ Correct

17. Following the previous Question, among $\log_{10} \lambda = \{2, 1, 0, -1, \dots, -8, -9, -10\}$. What is the λ with the minimum $E_{\text{val}}(g_{\lambda}^{-})$? Compute λ and the corresponding $E_{\text{train}}(g_{\lambda}^{-})$, $E_{\text{val}}(g_{\lambda}^{-})$ and $E_{\text{out}}(g_{\lambda}^{-})$ then select the closet answer. Break the tie by selecting the largest λ .

10/10 分

✓ Correct

18. Run the algorithm with the optimal λ of the previous Question on the whole \mathcal{D} to get g_{λ} . Compute $E_{\text{in}}(g_{\lambda})$ and $E_{\text{out}}(g_{\lambda})$ then select the closet answer.

10/10 分

✓ Correct

19. For Questions 19-20, split the given training examples in \mathcal{D} to five folds, the first 40 being fold 1, the next 40 being fold 2, and so on. Again, we take a fixed split because the given examples are already randomly permuted.

10/10 分

Among $\log_{10} \lambda = \{2, 1, 0, -1, \dots, -8, -9, -10\}$. What is the λ with the minimum E_{cv} , where E_{cv} comes from the five folds defined above? Compute λ and the corresponding E_{cv} then select the closest answer. Break the tie by selecting the largest λ .

✓ Correct

20. Run the algorithm with the optimal λ of the previous problem on the whole \mathcal{D} to get g_λ . Compute $E_{in}(g_\lambda)$ and $E_{out}(g_\lambda)$ then select the closest answer.

10/10 分

✓ Correct