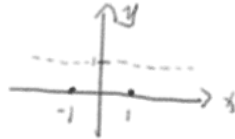


作业四

2021年12月10日 20:44

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 $h_1(x) = a_1x + b_1$.



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

$$E_{LOOCV_0} = 0.5^2 + 1^2 + 0.5^2 = 1.5$$

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

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

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?

Linear Regression Linear Regression Algorithm

Optimal Linear Regression Weights

task: find \mathbf{w}_{LIN} such that $\frac{2}{N} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) = \nabla E_{\text{in}}(\mathbf{w}) = 0$

invertible $\mathbf{X}^T \mathbf{X}$	singular $\mathbf{X}^T \mathbf{X}$
<ul style="list-style-type: none"> • easy! unique solution $\mathbf{w}_{\text{LIN}} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T}_{\text{pseudo-inverse } \mathbf{X}^\dagger} \mathbf{y}$ <ul style="list-style-type: none"> • often the case because $N \gg d+1$ 	<ul style="list-style-type: none"> • many optimal solutions • one of the solutions $\mathbf{w}_{\text{LIN}} = \mathbf{X}^\dagger \mathbf{y}$ <p>by defining \mathbf{X}^\dagger in other ways</p>

practical suggestion:
use **well-implemented** \dagger routine
instead of $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
for numerical stability when **almost singular**

17:14 / 20:03

$$\begin{aligned} \mathbf{X} &= \begin{pmatrix} \mathbf{x}_0 \\ \tilde{\mathbf{x}} \end{pmatrix} \quad \mathbf{Y} = \begin{pmatrix} \mathbf{y}_0 \\ \tilde{\mathbf{y}} \end{pmatrix} \\ \mathbf{w} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= [(\mathbf{x}_0^T \tilde{\mathbf{x}}^T) (\mathbf{x}_0^T \tilde{\mathbf{x}}^T)]^{-1} (\mathbf{x}_0^T \tilde{\mathbf{x}}^T) \begin{pmatrix} \mathbf{y}_0 \\ \tilde{\mathbf{y}} \end{pmatrix} \\ &= [\mathbf{x}_0^T \mathbf{x}_0 + \tilde{\mathbf{x}}^T \tilde{\mathbf{x}}]^{-1} [\mathbf{x}_0^T \mathbf{y}_0 + \tilde{\mathbf{x}}^T \tilde{\mathbf{y}}] \end{aligned}$$

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

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

$$\tilde{E}_{\text{aug}} = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\lambda}{N} \|\mathbf{w}\|^2$$

$$= \tilde{E}_{\text{sq}}(\mathbf{w}) + \underbrace{\frac{\lambda}{N} \mathbf{w}^T \mathbf{w}}_{\text{weight decay Regularization}}$$

☐ $\tilde{\mathbf{X}} = \mathbf{I}, \tilde{\mathbf{y}} = \mathbf{0}$

☒ $\tilde{\mathbf{X}} = \sqrt{\lambda} \mathbf{I}, \tilde{\mathbf{y}} = \mathbf{0}$

☐ $\tilde{\mathbf{X}} = \lambda \mathbf{I}, \tilde{\mathbf{y}} = \mathbf{1}$

☐ $\tilde{\mathbf{X}} = \sqrt{\lambda} \mathbf{X}, \tilde{\mathbf{y}} = \mathbf{y}$

☐ none of the other choices

Regularization Weight Decay Regularization

Augmented Error

- if oracle tells you $\lambda > 0$, then

solving $\nabla E_{\text{in}}(\mathbf{w}_{\text{REG}}) + \frac{2\lambda}{N} \mathbf{w}_{\text{REG}} = \mathbf{0}$

$$\frac{2}{N} (\mathbf{Z}^T \mathbf{Z} \mathbf{w}_{\text{REG}} - \mathbf{Z}^T \mathbf{y}) + \frac{2\lambda}{N} \mathbf{w}_{\text{REG}} = \mathbf{0}$$

- optimal solution:

$$\mathbf{w}_{\text{REG}} \leftarrow (\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^T \mathbf{y}$$

—called ridge regression in Statistics

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \lambda \mathbf{I} \quad \tilde{\mathbf{X}}^T \tilde{\mathbf{y}} = \mathbf{0}$$

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.

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? *No data snooping*

☒ 1

☐ N

☐ 2^N

☐ N^2

10. You assure the bank that you have a 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?

- ☐ By applying $a(\mathbf{x}) \text{ NOR } g(\mathbf{x})$ to approve credit for new customers, the performance of the overall credit approval system can be improved with guarantee provided by the previous problem.
- ☐ By applying $a(\mathbf{x}) \text{ NAND } g(\mathbf{x})$ to approve credit for new customers, the performance of the overall credit approval system can be improved with guarantee provided by the previous problem.
- ☐ By applying $a(\mathbf{x}) \text{ OR } g(\mathbf{x})$ to approve credit for new customers, the performance of the overall credit approval system can be improved with guarantee provided by the previous problem.
- ☒ By applying $a(\mathbf{x}) \text{ AND } g(\mathbf{x})$ to approve credit for new customers, the performance of the overall credit approval system can be improved with guarantee provided by the previous problem.
- ☐ none of the other choices

Origin Data $\xrightarrow{\text{if } a(x)=1} (x_1, y_1), \dots, (x_n, y_n) \xrightarrow{g(x)} \text{perfect prediction}$