

Assignment #4 (Neural Networks)

Instructor: Beilun Wang

Name: Qipeng Zhu, ID: 58119304

Problem Description:

Problem 1: Model selection and learning theory

(1) Let $X_i, i = 1, 2, \dots, n$ be n i.i.d observations from CDF $F(t) = P(X \leq t)$. If we estimate the true CDF by empirical CDF which is

$$\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathcal{I}(X_i \leq t)$$

where $\mathcal{I}(p) = 1$ if statement p is true, otherwise 0.

Write down the expectation and variance of $\hat{F}_n(t)$. Then use Khinchin's law (i.e., the weak law of large numbers) to show that $\hat{F}_n(t) \xrightarrow{P} F(t)$.

(2) Suppose we have a target variable y and a vector of inputs \mathbf{x} , and the true model is $f(\mathbf{x})$. If we assume that $y = f(\mathbf{x}) + \varepsilon$ where ε is a Gaussian noise with $E(\varepsilon) = 0$ and $Var(\varepsilon) = \sigma_\varepsilon^2$, we can derive an expression for the expected prediction error of a regression fit $\hat{f}(\mathbf{x})$ at an point $\mathbf{x} = \mathbf{x}_0$, using squared-error loss:

$$EPE(\mathbf{x}_0) = \sigma_\varepsilon^2 + Bias^2(\hat{f}(\mathbf{x}_0)) + Var(\hat{f}(\mathbf{x}_0)).$$

Now give training set \mathbf{X} and the corresponding labels \mathbf{y} .

(a) For a linear model fit $\hat{f}_p(\mathbf{x}) = \mathbf{x}^\top \hat{\boldsymbol{\beta}}$, where the parameter vector $\hat{\boldsymbol{\beta}}$ with p components is fit by least squares, write down the closed form solution of $\hat{\boldsymbol{\beta}}$. (Assume that $\mathbf{X}^\top \mathbf{X}$ is invertible.)

If $\hat{\boldsymbol{\beta}}$ is fit by least squares using ridge regression with regularization parameter α , we can get a model $\hat{f}_\alpha(\mathbf{x})$. Also write down the closed form solution of $\hat{\boldsymbol{\beta}}_\alpha$.

(b) For the linear model $\hat{f}_p(\mathbf{x}) = \mathbf{x}^\top \hat{\boldsymbol{\beta}}$, write down the expression of EPE at $\mathbf{x} = \mathbf{x}_0$ using the solution $\hat{\boldsymbol{\beta}}$ in (a). (Assume only y is random variable).

(c) For the model $\hat{f}_\alpha(\mathbf{x}) = \mathbf{x}^\top \hat{\boldsymbol{\beta}}_\alpha$, write down the expression of EPE at $\mathbf{x} = \mathbf{x}_0$ using the solution $\hat{\boldsymbol{\beta}}_\alpha$ in (a). (Assume only y is random variable).

(3) Suppose we stack the outcomes y_1, y_2, \dots, y_N into a vector \mathbf{y} , and similarly for the predictions $\hat{\mathbf{y}}$. Then a linear fitting method is one for which we can write

$$\hat{\mathbf{y}} = \mathbf{S}\mathbf{y},$$

where \mathbf{S} is an $N \times N$ matrix depending on the input vectors \mathbf{x}_i but not on the y_i .

Let $\hat{\mathbf{f}} = \mathbf{S}\mathbf{y}$ be a linear fitting of \mathbf{y} , and let $\hat{\mathbf{f}}^{-i}$ be the fitted function computed with (\mathbf{x}_i, y_i) removed. If S_{ii} is the i th diagonal element of \mathbf{S} , show that for

$$\mathbf{S} = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top = \mathbf{X} \mathbf{A}^{-1} \mathbf{X}^\top \quad (\mathbf{A} \text{ is PSD})$$

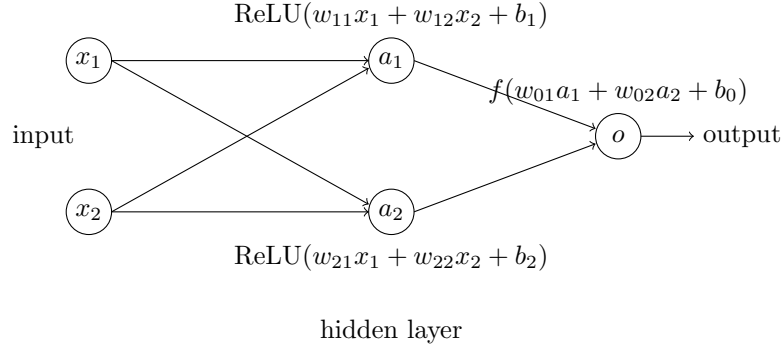
the cross-validated residual can be written as

$$y_i - \hat{f}^{-i}(\mathbf{x}_i) = \frac{y_i - \hat{f}(\mathbf{x}_i)}{1 - S_{ii}}$$

(Hint: use a lemma $(\mathbf{A} - \mathbf{x}\mathbf{x}^\top)^{-1} = \mathbf{A}^{-1} + \frac{\mathbf{A}^{-1}\mathbf{x}\mathbf{x}^\top\mathbf{A}^{-1}}{1 - \mathbf{x}^\top\mathbf{A}^{-1}\mathbf{x}}$)

Problem 2: Neural Networks

Suppose that we apply neural networks on a problem which has boolean inputs $\mathbf{x} \in \{0, 1\}^p$ and boolean output $y \in \{0, 1\}$. The network structure example is showed as below. In this example we set $p = 2$, single hidden layer with 2 neurons, activation function $\text{ReLU}(u) = u$ if $u > 0$ otherwise 0, and an additional threshold function (e.g., $f(v) = 1$ if $v > 0$, otherwise $f(v) = 0$) for output layer.



(1) Using the structure and settings of neural network above, show that such a simple neural network could output the function $x_1 \text{ XOR } x_2$ (equals to 0 if $x_1 = x_2$ and otherwise 1), which is impossible for linear models. State the values of parameters (i.e., w_{ij} and b_i) you found.

(2) Now we allow the number of neurons in the hidden layer to be more than 2 but finite. Retain the structure and other settings. Show that such a neural network with single hidden layer could output an arbitrary binary function $h : \{0, 1\}^p \mapsto \{0, 1\}$. You can apply threshold function after each neuron in the hidden layer.

Answer:**Problem 1: Model selection and learning theory**

(1)

$$E(\hat{F}_n(t)) = \frac{1}{n} \sum_{i=1}^n E(\mathcal{I}(X_i \leq t)) = \frac{1}{n} \sum_{i=1}^n P(\mathbf{X}_i \leq \mathbf{t})$$

$\because P(\mathbf{X}_i \leq t)$ is i.i.d. \therefore according Chebyshev Theorem: we can get:

$$\hat{F}_n(t) \xrightarrow{P} \frac{1}{n} \sum_{i=1}^n P(\mathbf{X}_i \leq \mathbf{t}) \xrightarrow{P} \frac{1}{n} \times n P(\mathbf{X}_i \leq \mathbf{t}) = F(t)$$

(2) Suppose

$$J(\hat{\beta}) = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2$$

First we can simplify the optimization:

$$J(\hat{\beta}) = (\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta}) = \hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} - 2\hat{\beta}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$$

so its gradient is:

$$\nabla_{\hat{\beta}} J(\hat{\beta}) = 2\mathbf{X}^T \mathbf{X} \hat{\beta} - 2\mathbf{X}^T \mathbf{y}$$

let

$$\nabla_{\hat{\beta}} J(\hat{\beta}) = 0$$

we can get

$$\hat{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

If we use Ridge Regression: Suppose

$$J(\hat{\beta}_\alpha) = \|\mathbf{y} - \mathbf{X}\hat{\beta}_\alpha\|^2 + \alpha \|\hat{\beta}_\alpha\|^2$$

First we can simplify the optimization:

$$J(\hat{\beta}_\alpha) = (\mathbf{y} - \mathbf{X}\hat{\beta}_\alpha)^T (\mathbf{y} - \mathbf{X}\hat{\beta}_\alpha) + \alpha \hat{\beta}_\alpha^T \hat{\beta}_\alpha = \hat{\beta}_\alpha^T \mathbf{X}^T \mathbf{X} \hat{\beta}_\alpha - 2\hat{\beta}_\alpha^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} + \alpha \hat{\beta}_\alpha^T \hat{\beta}_\alpha$$

so its gradient is:

$$\nabla_{\hat{\beta}_\alpha} J(\hat{\beta}_\alpha) = 2\mathbf{X}^T \mathbf{X} \hat{\beta}_\alpha - 2\mathbf{X}^T \mathbf{y} + 2\alpha \hat{\beta}_\alpha$$

let

$$\nabla_{\hat{\beta}_\alpha} J(\hat{\beta}_\alpha) = 0$$

we can get

$$\hat{\beta}_\alpha^* = (\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

(b)

$$\text{Bias}^2(\hat{f}(\mathbf{x}_0)) = \left[f(\mathbf{x}_0) - E(\hat{f}(\mathbf{x}_0)) \right]^2$$

$$= (f(\mathbf{x}_0) - E(\mathbf{X}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}))^2$$

$$\because E\mathbf{y} = E(f(\mathbf{x}_0) + \varepsilon) = f(\mathbf{x}_0)$$

$$\therefore \text{Bias}^2(\hat{f}(\mathbf{x}_0)) = \left(1 - \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right)^2 f^2(\mathbf{x}_0)$$

$$\text{Var}(\hat{f}(\mathbf{x}_0)) = \hat{f}(\mathbf{x}_0) - E(\hat{f}(\mathbf{x}_0))^2$$

$$\begin{aligned}
&= (\mathbf{x}_0^T (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (y - E y))^2 \\
&= (\mathbf{x}_0^T (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \varepsilon)^2 \\
\therefore EPE(x_0) &= \sigma_\varepsilon^2 + \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0)) \\
&= \sigma_\varepsilon^2 + \left(1 - \mathbf{x}_0^T (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top\right)^2 f^2(x_0) + (\mathbf{x}_0^T (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \varepsilon)^2
\end{aligned}$$

(c)

$$\begin{aligned}
\text{Bias}^2(\hat{f}(\mathbf{x}_0)) &= \left[f(\mathbf{x}_0) - E(\hat{f}(\mathbf{x}_0))\right]^2 \\
&= (f(\mathbf{x}_0) - E(\mathbf{X}^\top (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top y))^2 \\
&= \left(1 - \mathbf{x}_0^T (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top\right)^2 f^2(\mathbf{x}_0) \\
\text{Var}(\hat{f}(\mathbf{x}_0)) &= \hat{f}(\mathbf{x}_0) - E(\hat{f}(\mathbf{x}_0))^2 \\
&= (\mathbf{x}_0^T (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top (y - E y))^2 \\
&= (\mathbf{x}_0^T (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top \varepsilon)^2 \\
\therefore EPE(x_0) &= \sigma_\varepsilon^2 + \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0)) \\
&= \sigma_\varepsilon^2 + \left(1 - \mathbf{x}_0^T (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top\right)^2 f^2(\mathbf{x}_0) + (\mathbf{x}_0^T (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top \varepsilon)^2
\end{aligned}$$

(3) Suppose we stack the outcomes y_1, y_2, \dots, y_N into a vector \mathbf{y} , and similarly for the predictions $\hat{\mathbf{y}}$. Then a linear fitting method is one for which we can write

$$\hat{\mathbf{y}} = \mathbf{S}\mathbf{y},$$

where \mathbf{S} is an $N \times N$ matrix depending on the input vectors \mathbf{x}_i but not on the y_i .

Let $\hat{\mathbf{f}} = \mathbf{S}\mathbf{y}$ be a linear fitting of \mathbf{y} , and let $\hat{\mathbf{f}}^{-i}$ be the fitted function computed with (\mathbf{x}_i, y_i) removed. If S_{ii} is the i th diagonal element of \mathbf{S} , show that for

$$\mathbf{S} = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{\Omega})^{-1} \mathbf{X}^\top = \mathbf{X} \mathbf{A}^{-1} \mathbf{X}^\top \quad (\mathbf{A} \text{ is PSD})$$

$$(\mathbf{A} - \mathbf{x}\mathbf{x}^\top)^{-1} = \mathbf{A}^{-1} + \frac{\mathbf{A}^{-1} \mathbf{x}\mathbf{x}^\top \mathbf{A}^{-1}}{1 - \mathbf{x}^\top \mathbf{A}^{-1} \mathbf{x}}$$

the cross-validated residual can be written as

$$y_i - \hat{f}^{-i}(\mathbf{x}_i) = \frac{y_i - \hat{f}(\mathbf{x}_i)}{1 - S_{ii}}$$

Problem 2: Neural Networks

(1) It's clear that we can find when $w_{11} = w_{22} = 1, w_{12} = w_{21} = -1, b_1 = b_2 = 0$ and $w_{01} = w_{02} = -1, b_0 = 1$, the neural network can output the function $x_1 \text{ XOR } x_2$. This is because:

When $x_1 = x_2 = 0$,

$$a_1 = \text{ReLU}(x_{11}x_1 + w_{12}x_2 + b_1) = 0,$$

Similarly, $a_2 = a_1 = 0$, so the output is $f(1) = 1$.

When $x_1 = 0, x_2 = 1$,

$$a_1 = \text{ReLU}(x_{11}x_1 + w_{12}x_2 + b_1) = 0,$$

$$a_2 = \text{ReLU}(x_{21}x_1 + w_{22}x_2 + b_2) = 1.$$

The output is $f(w_{01}a_1 + w_{02}a_2 + b_0) = 0$.

Similarly, when $x_1 = 1, x_2 = 0$, the result is 0.

Now when $x_1 = x_2 = 1$,

$$a_1 = a_2 = 0.$$

The output is $f(1) = 1$. So the result is the definition of function XOR.

(2) For 2 input signals x_1, x_2 , we have a total of $2^{2^2} = 16$ different binary functions h , say, h_1, h_2, \dots, h_{16} . To implement these functions, we can simply use $16 \times 2 = 32$ hidden neurons, each pair of them can represent a specific function.

For example, if function h_1 is the XOR function, then we can implement it by simply setting the first to hidden neurons and the corresponding weights to the numbers in question (1) while keeping other 30 neurons' weights as 0 (not used). Similarly, other XOR-like functions can be implemented in this way.

This is because we know that each pair of 2-neurons sub-network can implement any binary non-linear functions like XOR or XNOR as well as linear functions like AND or OR.