

## Assignment #4 (Neural Networks)

Instructor: Beilun Wang

Name: Haorui li, ID: 61518407

## 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  $\varepsilon$  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  $\alpha$ , 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

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

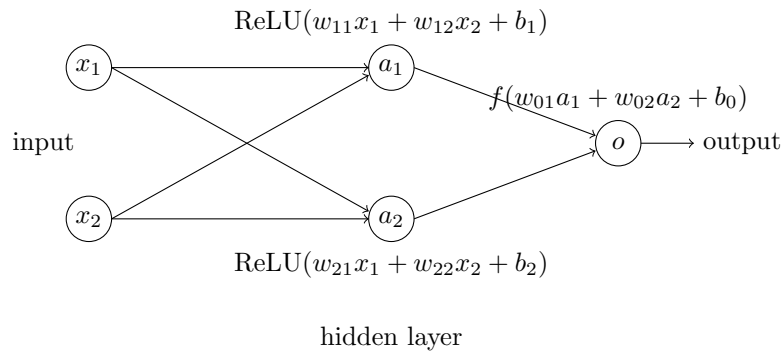
the cross-validated residual can be written as

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

(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)

$$E(\mathcal{I}\{x_i \leq t\}) = P(x_i \leq t)$$

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

$$E(F_n^2) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n E(\mathcal{I}\{x_i \leq t\} \mathcal{I}\{x_j \leq t\}) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n P(x_i \leq t, x_j \leq t)$$

$$\text{Var}(F_n(t)) = E(F_n(t)^2) - E(F_n(t))^2 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n P(x_i \leq t, x_j \leq t) - \left(\frac{1}{n} \sum_{i=1}^n F_i(t)\right)^2$$

$$\text{Let } \mu = E(\mathcal{L}(X_i \leq t)) = F(t)$$

$$\lim_{n \rightarrow \infty} P\left\{\left|\frac{1}{n} \sum_{i=1}^n \mathcal{L}(X_i \leq t) - \mu\right| < \varepsilon\right\} = 1$$

Therefore:

$$\lim_{n \rightarrow \infty} P\{|\hat{F}_n - F(t)| < \varepsilon\} = 1$$

(2)

(a)

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

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

(b)

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

(c)

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

(3)

$$S_{ii} = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X} + \lambda \Omega)^{-1} \mathbf{x}_i$$

$$\hat{\mathbf{f}}(\mathbf{x}_i) = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X} + \lambda \Omega)^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{\mathbf{f}}^{-i}(x_i) = x_i^T (X^T X - x_i x_i^T + \lambda \Omega)^{-1} (X^T y - x_i y_i) \quad (1)$$

$$= x_i^T (A - x_i x_i^T)^{-1} (X^T y - x_i y_i) \quad (2)$$

$$= y_i - \frac{y_i - \hat{f}(x_i)}{1 - S_{ii}} \quad (3)$$

### Problem 2: Neural Networks

(1)

$$w_{11} = 1; w_{12} = -1; w_{21} = -1; w_{22} = 1; w_{01} = w_{02} = 1; b_1 = b_2 = b_0 = 0$$

$$f = \begin{cases} 1, & x > 0 \\ 0, & \text{other} \end{cases}$$

(2) Any boolean functions can be represented by AND, NOT and OR.

And these logic can be represented by neural networks as follow:

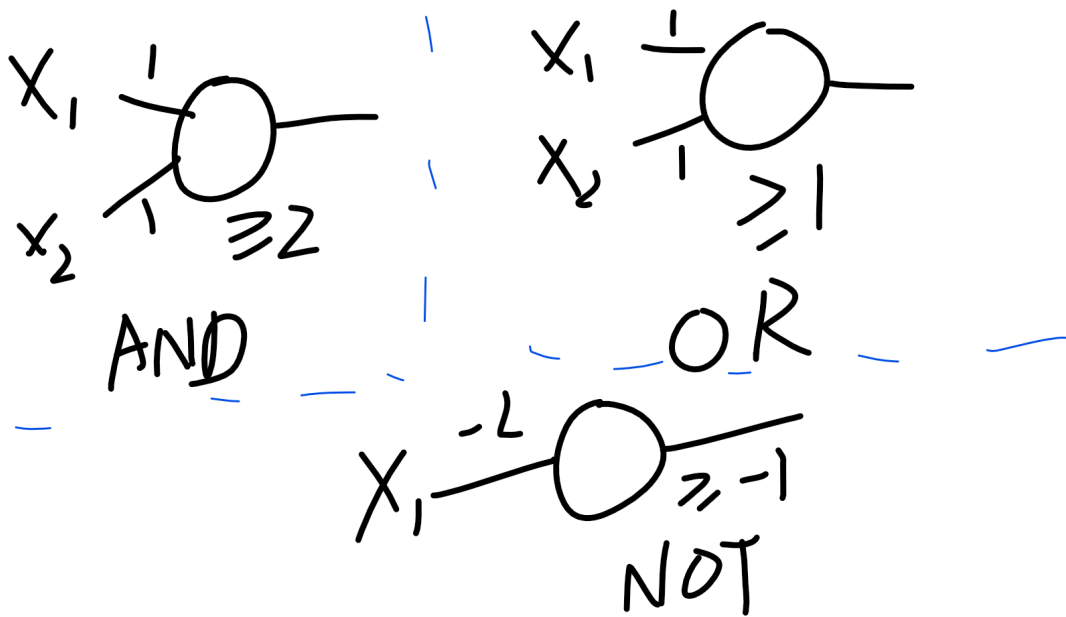


Figure 1: Network logic

So, NN can represent arbitrary binary functions.