

Principles of Statistical Machine Learning

Introduction to Nearest Neighbors Learning

Fokoué Ernest, PhD
Professor of Statistics



@ErnestFokoue

To understand God's thoughts, one must study statistics, the measure of His purpose
Florence Nightingale

Exercise: 1NN Classification I

- ① Consider the k Nearest Neighbors learning machine with $k = 1$, also known as 1NN or NN learning machine. Given a sample of size n ,

$$\text{namely } \mathcal{D}_n = \left\{ (\mathbf{x}_i, y_i) \stackrel{iid}{\sim} p_{\mathbf{xy}}(\mathbf{x}, y), \mathbf{x}_i \in \mathcal{X}, y_i \in \{1, 2, \dots, G\} \right\},$$

- Find the in sample prediction $\hat{f}_{\text{NN}}(\mathbf{x}_{i_m})$ for \mathbf{x}_{i_m} , where $i_m = \lceil (n+1)/2 \rceil$. *If $\mathbf{x}_i \in \mathcal{D}_n$, then its nearest neighbor is itself \mathbf{x}_i , so that $\hat{f}_{\text{NN}}(\mathbf{x}_i) = y_i$, because*

$$\min_{\mathbf{x}_j \in \mathcal{D}_n} d(\mathbf{x}_i, \mathbf{x}_j) = d(\mathbf{x}_i, \mathbf{x}_i) = 0$$

Therefore, the nearest neighbor to \mathbf{x}_i in \mathcal{D}_n is \mathbf{x}_i

$$\mathbf{x}_i = \underset{\mathbf{x}_j \in \mathcal{D}_n}{\operatorname{argmin}} d(\mathbf{x}_i, \mathbf{x}_j).$$

Since $\mathbf{x}_{i_m} \in \mathcal{D}_n$, it follows that

$$\hat{f}_{\text{NN}}(\mathbf{x}_{i_m}) = y_{i_m}.$$

Exercise: 1NN Classification II

- Find the in-sample error rate under the zero-one loss.

Since $\hat{f}_{\text{NN}}(\mathbf{x}_i) = y_i$, it holds true that $\mathbb{1}(y_i \neq \hat{f}_{\text{NN}}(\mathbf{x}_i)) = 0$, for all $\mathbf{x}_i \in \mathcal{D}_n$. Therefore, it immediately follows that

$$\hat{R}_n(\hat{f}_{\text{NN}}) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i \neq \hat{f}_{\text{NN}}(\mathbf{x}_i)) = 0, \quad \forall \mathcal{D}_n \subset \mathcal{X} \times \mathcal{Y}.$$

- Argue on a reasonable VC Dimension for

$$\mathcal{H} = \{\mathbf{x}_i \mapsto f_{1\text{NN}}(\mathbf{x}_i) = y_i, \quad \forall (\mathbf{x}_i, y_i) \in \mathcal{D}_n\}$$

Since $\hat{R}_n(\hat{f}_{\text{NN}}) = 0, \forall n \in \mathbb{N}$, it follows naturally that

$$\text{VCdim}(\mathcal{H}) = +\infty.$$

This result has a lot of consequences on the generalizability of the nearest neighbors learning machine when $k = 1$.

Exercise: 1NN Classification III

- 2 Consider the k Nearest Neighbors learning machine with $k = 1$, also known as 1NN or NN learning machine. Given a sample of size n , namely $\mathcal{D}_n = \{(\mathbf{x}_i, y_i) \stackrel{iid}{\sim} p_{\mathbf{xy}}(\mathbf{x}, y), \mathbf{x}_i \in \mathcal{X}, y_i \in \{0, 1\}\}$, with $\pi = \Pr[Y_i = 1] = \mathbb{E}[Y_i]$ for all $i = 1, \dots, n$.

- What is the bias of this learning machine?

Now, $\forall \mathbf{x}_i \in \mathcal{D}_n$, the point-wise bias $\hat{f}_{\text{NN}}(\mathbf{x}_i)$ of \hat{f}_{NN} is

$$\text{Bias}_n(\hat{f}_{\text{NN}}(\mathbf{x}_i)) = \mathbb{E}[\hat{f}_{\text{NN}}(\mathbf{x}_i)] - f^*(\mathbf{x}_i)$$

where $f^*(\mathbf{x}_i)$ is the theoretical (true) underlying response yield by the generator of the data, corresponding to the Bayes learning machine prediction.

Exercises and Problems

- ① Consider the k Nearest Neighbors learning machine with $k = n$, also known as n NN learning machine. Given a sample of size n , namely

$$\mathcal{D}_n = \{(\mathbf{x}_i, y_i) \stackrel{iid}{\sim} p_{\mathbf{xy}}(\mathbf{x}, y), \mathbf{x}_i \in \mathcal{X}, y_i \in \mathbb{R}, i = 1, \dots, n\},$$

- Find the in sample prediction $\hat{f}_{\text{nNN}}(\mathbf{x}_{i_m})$ for \mathbf{x}_{i_m} , where $i_m = \lceil (n+1)/2 \rceil$. $\forall \mathbf{x}_i \in \mathcal{D}_n$, it is true that $\hat{f}_{\text{nNN}}(\mathbf{x}_i) = \text{constant} = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$, therefore it follows that

$$\hat{f}_{\text{nNN}}(\mathbf{x}_{i_m}) = \bar{y}, \quad \text{since } \mathbf{x}_{i_m} \in \mathcal{D}_n.$$

- Find the in-sample error rate under the squared error loss $\mathcal{L}(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$.

$$\hat{R}_n(\hat{f}_{\text{nNN}}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, \hat{f}_{\text{nNN}}(\mathbf{x}_i)) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \left(\frac{n-1}{n} \right) S_y^2$$

- What is the variance of this learning machine? Since $\hat{f}_{\text{nNN}}(\mathbf{x}_i) = \text{constant} = \bar{y}$ for all $\mathbf{x}_i \in \mathcal{D}_n$, it follows that

$$\text{Variance}(\hat{f}_{\text{nNN}}) = 0 \quad \text{smallest possible value.}$$