

Data Splitting for Binary Classification

1 Keywords

- Area Under Curve

$$\text{AUC} \approx \sum_k \frac{(\text{TPR}_{\lambda_k} + \text{TPR}_{\lambda_{k+1}})}{2} (\text{FPR}_{\lambda_{k+1}} - \text{FPR}_{\lambda_k}) \quad (1)$$

- Binary Classification Problem

$$C = 2 \quad (2)$$

- Bin Width

$$\text{BW} = \frac{\max_n X_{n,f} - \min_n X_{n,f}}{K} \quad (3)$$

where K is the number of bins.

- False Negative

$$\text{FN} = \sum_n \mathbf{1}(\hat{x}_n = 0) \mathbf{1}(y_n = 1) \quad (4)$$

- False Positive

$$\text{FP} = \sum_n \mathbf{1}(\hat{x}_n = 1) \mathbf{1}(y_n = 0) \quad (5)$$

- False Positive Rate (1-specificity)

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{Precision} = \frac{\sum_n \mathbf{1}(\hat{y}_n = 1) \mathbf{1}(y_n = 0)}{\sum_n \mathbf{1}(y_n = 0)} \quad (6)$$

- Normalized Histogram (\tilde{x}_k, h_k)

$$\begin{cases} \forall k \in \{0, \dots, K\} & x_k = \min_n X_{n,f} + k\text{BW} \\ \forall k \in \{1, \dots, K\} & \tilde{x}_k = \frac{x_{k-1} + x_k}{2} \\ \forall k \in \{1, \dots, K\} & h_k = \frac{1}{N} \sum_n \mathbf{1}(x_{k-1} \leq X_{n,f} < x_k) \end{cases} \quad (7)$$

- Normalized Conditional Histogram (\tilde{x}_k, h_k^c)

$$\begin{cases} \forall k \in \{0, \dots, K\} & x_k = \min_n X_{n,f} + k\text{BW} \\ \forall k \in \{1, \dots, K\} & \tilde{x}_k = \frac{x_{k-1} + x_k}{2} \\ \forall k \in \{1, \dots, K\} & h_k^c = \frac{\sum_n \mathbf{1}(x_{k-1} \leq X_{n,f} < x_k) \mathbf{1}(Y_n = c)}{\sum_n \mathbf{1}(Y_n = c)} \end{cases} \quad (8)$$

where $c \in \{0, 1\}$.

- Precision

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{\sum_n \mathbf{1}(\hat{y}_n = 1) \mathbf{1}(y_n = 1)}{\sum_n \mathbf{1}(\hat{x}_n = 1)} \quad (9)$$

- Recall

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\sum_n \mathbf{1}(\hat{y}_n = 1) \mathbf{1}(y_n = 1)}{\sum_n \mathbf{1}(y_n = 1)} \quad (10)$$

- ROC curve:

$$(\text{FPR}_\lambda, \text{TPR}_\lambda)_\lambda \quad (11)$$

- True Negative

$$\text{TN} = \sum_n \mathbf{1}(\hat{y}_n = 0)\mathbf{1}(y_n = 0) \quad (12)$$

- True Positive Rate (sensitivity)

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \text{Recall} = \frac{\sum_n \mathbf{1}(\hat{y}_n = 1)\mathbf{1}(y_n = 1)}{\sum_n \mathbf{1}(y_n = 1)} \quad (13)$$

- True Positive

$$\text{TP} = \sum_n \mathbf{1}(\hat{y}_n = 1)\mathbf{1}(y_n = 1) \quad (14)$$