

1 Criteria of Learning Performance of A Classifier

Suppose that we are given a data set $D = \{(x_1, a_1), (x_2, a_2), \dots, (x_m, a_m)\}$ such that $x_i \in \mathbb{R}^K$ and $a_i \in \{0, 1\}$, $i = 1, 2, \dots, m$, where we call a_i the *class of x_i* . A feature vector x_i with $a_i = 1$ (resp., 0) is called *positive* (resp., *negative*). We call a function $\eta : \mathbb{R}^K \rightarrow \{0, 1\}$ that estimates the class of a feature vector a *classifier*.

There are various criteria to evaluate the learning performance of a classifier. We summarize some of them.

1.1 Accuracy Based Criteria

Let us partition the data set D into $D = D_1 \cup D_0$, where D_1 (resp., D_0) is the set of all positive (resp., negative) feature vectors in D . For a classifier η , a feature vector $x_i \in D$ is called

- *true positive* if $a_i = 1$ and $\eta(x_i) = 1$;
- *true negative* if $a_i = 0$ and $\eta(x_i) = 0$;
- *false positive* if $a_i = 0$ and $\eta(x_i) = 1$; and
- *false negative* if $a_i = 1$ and $\eta(x_i) = 0$.

We denote by $TP(\eta; D)/TN(\eta; D)/FP(\eta; D)/FN(\eta; D)$ the sets of true positive/true negative/false positive/false negative feature vectors, respectively. We define $TPR(\eta; D)$, $TNR(\eta; D)$, $FPR(\eta; D)$ and $FNR(\eta; D)$ as follows;

$$\begin{aligned} TPR(\eta; D) &\triangleq \frac{TP(\eta; D)}{|D_1|}; & TNR(\eta; D) &\triangleq \frac{TN(\eta; D)}{|D_0|}; \\ FPR(\eta; D) &\triangleq \frac{FP(\eta; D)}{|D_0|}; & FNR(\eta; D) &\triangleq \frac{FN(\eta; D)}{|D_1|}. \end{aligned}$$

It holds that $TPR(\eta; D) + FNR(\eta; D) = TNR(\eta; D) + FPR(\eta; D) = 1$. The *accuracy* $ACC(\eta; D)$ is defined to be:

$$ACC(\eta; D) \triangleq \frac{TP(\eta; D) + TN(\eta; D)}{|D|}.$$

The *balanced accuracy* $B-ACC(\eta; D)$ is defined to be:

$$B-ACC(\eta; D) \triangleq \frac{1}{2}(TPR(\eta; D) + TNR(\eta; D)).$$

1.2 ROC Curve and AUC

Let $f : \mathbb{R}^K \rightarrow \mathbb{R}$ be a function and $\theta \in \mathbb{R}$ be a real number. We construct a classifier $\eta_{f,\theta} : \mathbb{R}^K \rightarrow \{0, 1\}$ as follows; for $x \in \mathbb{R}^K$,

$$\eta_{f,\theta}(x) \triangleq \begin{cases} 1 & \text{if } f(x) \geq \theta, \\ 0 & \text{otherwise.} \end{cases}$$

Suppose that x_1, x_2, \dots, x_m are sorted so that $f(x_1) \leq f(x_2) \leq \dots \leq f(x_m)$ holds. For $i = 1, \dots, m-1$, let us define

$$\theta_i \triangleq \frac{f(x_i) + f(x_{i+1})}{2},$$

where we let $\theta_0 \triangleq f(x_1) - \varepsilon$ and $\theta_m \triangleq f(x_m) + \varepsilon$ for a small constant $\varepsilon \in \mathbb{R}_+ \setminus \{0\}$.

We have $m+1$ classifiers $\eta_{f,\theta_0}, \eta_{f,\theta_1}, \dots, \eta_{f,\theta_m}$. Let us denote a 2D point $p_i \triangleq (\text{FPR}(\eta_{f,\theta_i}; D), \text{TPR}(\eta_{f,\theta_i}; D))$, $i = 0, 1, \dots, m$. Observe that $p_0 = (1, 1)$ holds since $\eta_{f,\theta_0}(x) = 1$ holds for all $x \in D$, and that $p_m = (0, 0)$ holds since $\eta_{f,\theta_m}(x) = 0$ holds for all $x \in D$. Also we have;

$$\text{FPR}(\eta_{f,\theta_m}; D) = 0 \leq \text{FPR}(\eta_{f,\theta_{m-1}}; D) \leq \dots \leq \text{FPR}(\eta_{f,\theta_0}; D) = 1;$$

$$\text{TPR}(\eta_{f,\theta_m}; D) = 0 \leq \text{TPR}(\eta_{f,\theta_{m-1}}; D) \leq \dots \leq \text{TPR}(\eta_{f,\theta_0}; D) = 1.$$

The *Receiver Operating Characteristic curve (ROC curve)* of f is a set of m line segments $(p_m = (0, 0), p_{m-1}), (p_{m-1}, p_{m-2}), \dots, (p_1, p_0 = (1, 1))$. The *Area Under Curve (AUC)* of f , which we denote by $\text{AUC}(f; D)$, is defined to be the area between the ROC curve and the x-axis. Hence we have $0 \leq \text{AUC}(f; D) \leq 1$.

Theorem:

When $f : \mathbb{R}^K \rightarrow \{0, 1\}$, $\text{AUC}(f; D) = \text{BACC}(\eta; D)$