# References

1. Pattern Classification

   Duda, Hart, Stork

2. Probabilistic Theory
   of Pattern Recognition

   Devroye, Gyorti, Lugosi

3. Pattern recognition
   and Machine Learning

   Chris Bishop

# Recap

$X \rightarrow$ Instance space

$Y \rightarrow$ Label space

Binary Classification
$$Y = \{-1, 1\}$$

Classifier $\qquad X \subseteq \mathbb{R}^d$
$$h: X \rightarrow \{-1, 1\}$$

## Prior

$$P_1 = P(Y = 1), \qquad P_2 = P(Y = -1)$$

Class conditional Distribution

$$P(X = x | Y = 1), \qquad P(X = x | Y = -1)$$

$$[\text{Densities} / P.m.f]$$

$$\eta(x) = P(Y = 1 \mid X = x)$$

$$1 - \eta(x) = P(Y = -1 \mid X = x)$$

$$P(Y = 1 \mid X = x) = \frac{P(Y = 1) P(X = x \mid Y = 1)}{P(X = x)}$$

## Bayes Classifier

$$h(x) = \begin{cases} 1 & P(Y = 1 \mid X = x) > P(Y = -1 \mid X = x) \\ -1 & P(Y = -1 \mid X = x) \geq P(Y = 1 \mid X = x) \end{cases}$$

$$h(x) = \begin{cases} 1 & \eta(x) > 1 - \eta(x) \\ -1 & 1 - \eta(x) \geq \eta(x) \end{cases}$$

$$\overset{x}{h(x)} = \text{sign}\left(2\eta(x) - 1\right) \longrightarrow BC$$

$$\text{sign}(z) = \begin{cases} 1 & z > 0 \\ -1 & z \leq 0 \end{cases}$$

Specific case

$x \in \mathbb{R}^d, \mu \in \mathbb{R}^d$

$C \in \mathbb{R}^{d \times d}$, Symmetric positive semi-definite

$$N(x| \mu, C)$$

$$= \frac{1}{(\sqrt{2\pi})^d |C|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T C^{-1}(x-\mu)}$$

$$|C| \equiv \det(C)$$

$$P(X = x | Y = 1) = N(x|\mu_1, C_1)$$

$$P(X = x | Y = -1) = N(x|\mu_2, C_2)$$

If $C_1 = C_2 = C$

$$h^*(x) = \text{sign}\left(w^T x + b\right)$$

$$w = C^{-1}\left(\mu_1 - \mu_2\right)$$

$$b = \log \frac{p_1}{p_2} - \frac{1}{2}\left(\mu_1^T C^{-1} \mu_1 - \mu_2^T C^{-1} \mu_2\right)$$

How good is the
Bayes classifier (BC)

Thm 2.1 DG2
$$P(h(x) \neq Y) - P(h^*(x) \neq Y) \geq 0$$

$$P(h(x) \neq Y) = E_x E_{Y/x}(h(x) \neq Y))$$

$$E_{Y/x = x}\{h(x) \neq Y\}$$

$$= \eta(x) \, 1_{\{h(x) \neq 1\}} + (1-\eta(x)) 1_{\{h(x) \neq -1\}}$$

$$= \begin{cases} (1-\eta(x)) & h(x) = 1 \\ \\ \eta(x) & h(x) = -1 \end{cases}$$

If $y = -1$, but $h(x) = 1$, then there is a mistake.

If $y = 1$, but $h(x) = -1$ then again there is a mistake

If $\eta(x) > (1 - \eta(x))$ predict $h(x) = 1$

The minimum over $h$ (all possible classifiers) is attained at choosing $h(x)$ such that

$$E_{Y|X=x} 1\{h(x) = Y\} = \min\left(\eta(x), 1 - \eta(x)\right)$$

$$h^*(x) = \begin{cases} 1 & \eta(x) > 1 - \eta(x) \\ \\ -1 & 1 - \eta(x) > \eta(x) \end{cases}$$

# Bayes decision theory

We have $x \in \mathbb{R}^d$ with label $y \in \{-1, 1\}$

We predict $\hat{y} \in \{1, 1\}$

$$\ell(\hat{y}, y) : Y \times Y \Rightarrow \mathbb{R}_+$$

### Expected loss

$$R(h) = \mathbb{E}_{X,Y} \ell(h(x), Y)$$

$$\min_h R(h)$$

$$\underset{h}{\text{minimize}} \quad R(h) = R(\tilde{h})$$

$$\tilde{h}(x) = \underset{h}{\min} \, E_{Y|X=x} \, \ell\left(\underline{\tilde{h}(x)}, Y\right)$$

$R(\tilde{h})$ is called the Bayes error-rate.

For a given instance it chooses the label which yields minimum loss.

minimize $R(h)$
h

$$\tilde{h}(x) = \min_{h} E_{Y|X=x} \ell(h(x), Y)$$

## 2 class problem

$$\min_{h} E_{Y|X=x} \ell(h(x), Y) \quad Y \in \{-1, 1\}$$

$$E_{Y|X=x} \ell(1, Y) = \ell(1,1) P(Y=1 | X=x)$$
$$+ \ell(1,-1) P(Y=-1 | X=x)$$
$$= \ell(1,1) \eta(x) + \ell(1,-1)(1-\eta(x))$$

$$E_{Y|X=x} \ell(-1, Y) = \ell(-1,1) \eta(x) + \ell(-1,-1)(1-\eta(x))$$

Expected Loss

Choose $h$ such that

$$E_{Y|X=x} \, \ell(h(x), Y) \text{ should be minimum}$$

Choose class $\boxed{h(x) = 1}$

if $E_{Y|X} \, \ell(1, Y) < E_{Y|X} \, \ell(-1, Y)$

Choose class $\boxed{h(x) = -1}$

if $E_{Y|X=x} \, \ell(-1, Y) < E_{Y|X} \, \ell(1, Y)$

$$\underline{\tilde{h}(x) = 1}$$

$$\therefore \ell(1,1)\eta(x) + \ell(1,-1)(1-\eta(x))$$
$$< \ell(-1,1)\eta(x) + \ell(-1,-1)(1-\eta(x))$$

$$\left(\ell(1,1) - \ell(-1,1)\right)\eta(x)$$
$$< \left(\ell(-1,-1) - \ell(1,-1)\right)(1-\eta(x))$$

$$\tilde{h}(x) = -1$$
$$\ell(-1,1)\eta(x) + \ell(-1,-1)(1-\eta(x))$$
$$< \ell(1,1)\eta(x) + \ell(1,-1)(1-\eta(x))$$

$$\Big(\ell(-1,-1) - \ell(1,-1)\Big)\big(1 - \eta(x)\big)$$

$$< \Big(\ell(1,1) - \ell(-1,1)\Big)\eta(x)$$

.

Properties of $\ell$

① non-negative

② should penalize
            mistakes more

Choose $\ell(-1, 1) = \ell(1, -1) = 1$

$\ell(1, 1) = \ell(-1, -1) = 0$

loss should penalize mistakes more

$\tilde{h}(x) = 1$ if

$(0 - 1) \eta(x) < (0 - 1)(1 - \eta(x))$

$\Rightarrow \quad \eta(x) > 1 - \eta(x)$

$\tilde{h}(x) = -1$ if

$1 - \eta(x) > \eta(x)$.

Then $\tilde{h}(x) = h^*(x)$

$R(\tilde{h}) < R(h)$

$R(\tilde{h})$ is called the B          error rate