

BAYES ERROR

(1)

$$(x, y) \in \mathbb{R}^d \times \{0, 1\} \sim (\mu, \eta). \mu(A) = \mathbb{P}(x \in A)$$

$\eta(x) = \mathbb{P}\{y=1 | x\} = \mathbb{E}[y|x]$. This is enough to describe $\mathbb{P}_{x,y}$ since if $C \subseteq \mathbb{R}^d \times \{0, 1\}$, then

$$\begin{aligned} C &= (C_0 \cap \mathbb{R}^d \times \{0\}) \cup (C_1 \cap \mathbb{R}^d \times \{1\}) \\ &= C_0 \times \{0\} \cup C_1 \times \{1\} \end{aligned}$$

$$\begin{aligned} \text{Then } \mathbb{P}[(x, y) \in C] &= \mathbb{P}[x \in C_0, y=0] + \mathbb{P}[x \in C_1, y=1] \text{ (disjoint)} \\ &= \mathbb{P}[y=0 | x \in C_0] \mathbb{P}[x \in C_0] + \mathbb{P}[y=1 | x \in C_1] \mathbb{P}[x \in C_1] \\ &= \int_{C_0} (1 - \eta(x)) \mu(dx) + \int_{C_1} \eta(x) \mu(dx) \end{aligned}$$

classifier $\rightarrow g: \mathbb{R}^d \rightarrow \{0, 1\}$. Error is $L(g) = \mathbb{P}[g(x) \neq y]$

Bayes classifier $\rightarrow g^*(x) = \begin{cases} 1 & \text{if } \eta(x) > 1/2 \\ 0 & \text{otherwise} \end{cases}$

Theorem g^* is optimal, i.e. for any g we have that

$$\begin{aligned} \mathbb{P}[g^*(x) \neq y] &\leq \mathbb{P}[g(x) \neq y] \\ L(g^*) &\leq L(g) \end{aligned}$$

$$\begin{aligned} \text{Proof. } \mathbb{P}[g(x) \neq y | x] &= 1 - \mathbb{P}[g(x)=y | x] \\ &= 1 - (\mathbb{P}[g(x)=0, y=0 | x] + \mathbb{P}[g(x)=1, y=1 | x]) \\ &= 1 - (\mathbb{E}[g(x)=0] \mathbb{P}[y=0 | x] + \mathbb{E}[g(x)=1] \mathbb{P}[y=1 | x]) \\ &= 1 - (\mathbb{E}[g(x)=0](1 - \eta(x)) + \mathbb{E}[g(x)=1]\eta(x)) \end{aligned}$$

Notice that $\mathbb{E}[g(x)=0] = 1 - \mathbb{E}[g(x)=1]$, thus

$$\begin{aligned} \mathbb{P}[g(x) \neq y | x] &= 1 - ((1 - \eta(x)) + \mathbb{E}[g(x)=1](2\eta(x) - 1)) \\ &= \eta(x) - (2\eta(x) - 1)\mathbb{E}[g(x)=1] \end{aligned}$$

$$P(g(x) \neq y | x) - P(g^*(x) \neq y | x) = (2\eta(x) - 1)(\mathbb{1}(g(x)=1) - \mathbb{1}(g^*(x)=1)) \quad (2)$$

If $\eta > \frac{1}{2} \Rightarrow \underbrace{(2(\frac{1}{2} + \delta) - 1)}_{>0} \underbrace{(-\mathbb{1}(g(x)=1))}_{>0} \geq 0$

If $\eta \leq \frac{1}{2} \Rightarrow \underbrace{(2(\frac{1}{2} - \delta) - 1)}_{\leq 0} \underbrace{(-\mathbb{1}(g(x)=1))}_{\leq 0} \geq 0$

thus in either case $P(g(x) \neq y | x) - P(g^*(x) \neq y | x) \geq 0$.

Using $P(g(x) \neq y) = \int P(g(x) \neq y | x) h(dx)$ we have the integral of a positive integrand over a positive range, so $L(g) - L(g^*) \geq 0$. \square

$$\begin{aligned} \text{Note that } P(g(x) \neq y) &= \int P(g(x) \neq y | x) h(dx) \\ &= \int h(dx) \{1 - \eta(x) \mathbb{1}(g(x)=1) - (1-\eta(x)) \mathbb{1}(g(x)=0)\} \\ &= 1 - \mathbb{E}_x [\eta(x) \mathbb{1}(g(x)=1) + (1-\eta(x)) \mathbb{1}(g(x)=0)] \end{aligned}$$

Also $\eta(x) = P(y=1 | x) = \mathbb{E}[y | x]$

Suppose $y = f(x)$ for some predictor $f: \mathbb{R}^d \rightarrow \mathbb{R}$ and consider $E[(f(x) - y)^2]$. We have $E[(f(x) - y)^2 | X] = E[(f - \eta + \eta - y)^2 | X]$, $E[(f(x) - y)^2 | X] = E[(f - \eta)^2 + 2(f - \eta)(\eta - y) + (\eta - y)^2 | X]$, $= (f - \eta)^2 + \underbrace{2(f - \eta)(\eta - y)}_{\geq 0} + \underbrace{(\eta - y)^2}_{=0} | X$

$$E[(f(x) - y)^2 | X] \geq E[(\eta(x) - y)^2 | X] \text{ or}$$

$$E[(f(x) - y)^2] \geq E[(\eta(x) - y)^2] \quad \text{Therefore, } \eta \text{ is a minimizer of the squared error function.}$$

Let us show other formulas.

(3)

Since g^* is the best classifier, $L^* = \inf_g P[g(x) \neq y]$

• Another formula is $L^* = E\{\min(\gamma(x), 1-\gamma(x))\}$

Proof Consider $P[g(x) \neq y | x] = P[g=0, y=1 | x] + P[g=1, y=0 | x]$

$$= \mathbb{1}(g=0) P[y=1 | x] + \mathbb{1}(g=1) P[y=0 | x]$$

$$= \mathbb{1}(g=0) \gamma(x) + \mathbb{1}(g=1) (1-\gamma(x))$$

Now $L(g) = P[g(x) \neq y] = \int P[g(x) \neq y | x] \mu(dx) = E[P[g(x) \neq y | x]]$

so $L^* = E[\mathbb{1}(\gamma \leq 1/2) \gamma(x) + \mathbb{1}(\gamma > 1/2) (1-\gamma(x))]$

Now observe that } if $\gamma \leq 1/2 \Rightarrow 1-\gamma \geq 1/2 \geq \gamma$
 if $\gamma > 1/2 \Rightarrow 1-\gamma < 1/2 < \gamma$

So if $\gamma \leq 1-\gamma$ we pick the first term $L^* = E[\gamma]$, and if
 $1-\gamma < \gamma$ we pick the second $L^* = E[1-\gamma]$. In other
 words, we pick the minimum between the two. \square

• Another formula is $L^* = \frac{1}{2} - \frac{1}{2} E[|2\gamma(x)-1|]$

Proof we showed $L^* = E\{\min(\gamma, 1-\gamma)\}$.

Suppose $\gamma > 1/2$. Then $L^* = E\{1-\gamma\} = E\left\{\frac{1}{2} - \frac{1}{2}(2\gamma-1)\right\}$

$$= E\left\{\frac{1}{2} - \frac{1}{2}|2\gamma-1|\right\}$$

$$= \frac{1}{2} - \frac{1}{2} E\{|2\gamma-1|\}.$$

Suppose $\gamma \leq 1/2$. Then $L^* = E\{\gamma\} = E\left\{\frac{1}{2} - \frac{1}{2}(-2\gamma+1)\right\}$

$$= E\left\{\frac{1}{2} - \frac{1}{2}|2\gamma-1|\right\}$$

$$= \frac{1}{2} - \frac{1}{2} E\{|2\gamma-1|\} \quad \square$$

If X has density f we have

(4)

$$L^* = \int \min(\gamma(x), 1-\gamma(x)) f(x) dx$$

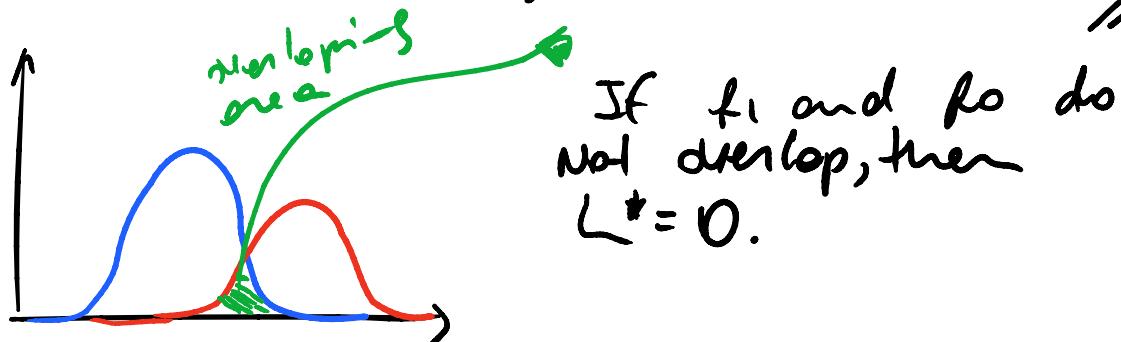
Now note that $\gamma(x) = P(Y=1|X) = \frac{p(x|Y=1)P(Y=1)}{f(x)}$

$$= \frac{f_1(x)P}{f(x)}$$

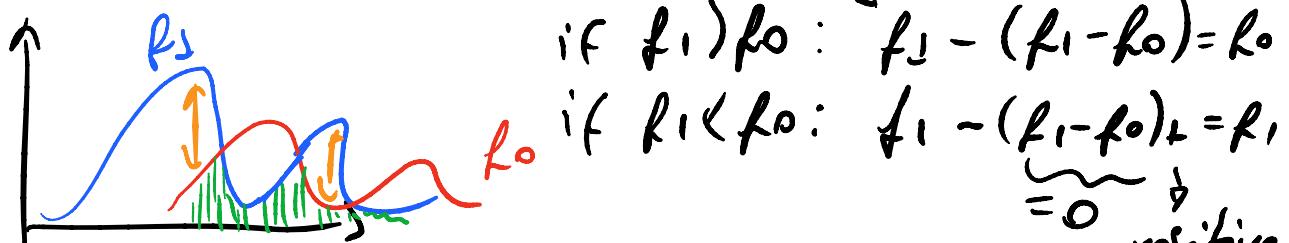
Analogously, $1-\gamma(x) = P(Y=0|X) = \frac{p(x|Y=0)P(Y=0)}{f(x)}$

$$= \frac{f_0(x)(1-p)}{f(x)}$$

Replacing this we have, $L^* = \int \min(pf_1(x), (1-p)f_0(x)) dx$

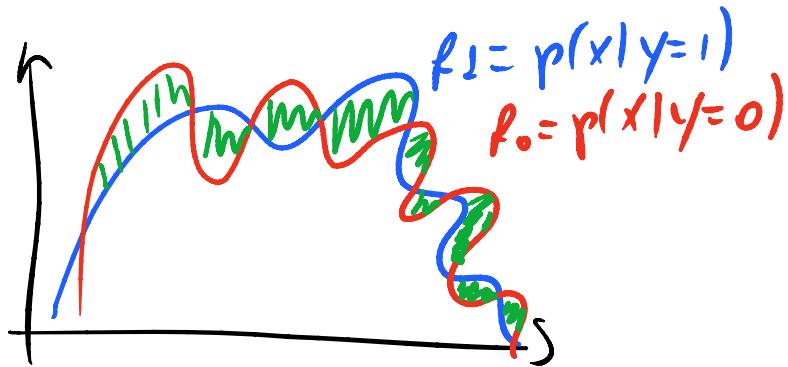


Suppose $p = 1/2$. $L^* = \frac{1}{2} \int \min(f_1, f_0) dx$



$$\begin{aligned} \min(f_1, f_0) &= f_1 - (f_1 - f_0)_+ \\ &= f_1 - \frac{1}{2}|f_1 - f_0| \end{aligned}$$

$$\begin{aligned} L^* &= \frac{1}{2} \underbrace{\int f_1 dx}_{=: J} - \frac{1}{4} \int |f_1 - f_0| dx \\ &= \frac{1}{2} J - \frac{1}{4} \int |f_1 - f_0| dx \quad L_1 \text{ distance between } \\ &\quad \text{the class densities.} \end{aligned} \tag{5}$$



Plugin Decision

$$g^*(x) = \begin{cases} 0 & \text{if } \gamma(x) \leq \frac{1}{2} \text{ (or } 1-\gamma \geq \gamma) \\ 1 & \text{otherwise.} \end{cases}$$

γ is unknown. Suppose we have an approximation $\tilde{\gamma}(x)$. Then we define

$$g(x) = \begin{cases} 0 & \text{if } \tilde{\gamma}(x) \leq \frac{1}{2} \\ 1 & \text{otherwise} \end{cases}$$

$$\text{Theo. } L(g) - L(g^*) = 2 \int_{\mathbb{R}^d} |\gamma - \frac{1}{2}| \mathbf{1}(g(x) \neq g^*(x)) \mu(dx)$$

$$\leq 2 \int_{\mathbb{R}^d} |\gamma - \tilde{\gamma}| \mu(dx) = 2 E(|\gamma - \tilde{\gamma}|)$$

So if $\tilde{\gamma}$ is close to γ in L_1 sense, L is close to L^* .

(6)

meat
faceless + hat

$$\begin{aligned} \mathbb{P}(g(x) \neq y|x) - \mathbb{P}(g^*(x) \neq y|x) &= (2\gamma(x) - 1)(\mathbb{I}(g^*(x)=1) - \mathbb{I}(g(x)=1)) \\ &= |2\gamma(x) - 1| \mathbb{I}(g(x) \neq g^*(x)) \end{aligned}$$

This gives the first equality upon integration.

If $g \neq g^*$ then $|\tilde{\gamma} - \gamma| \geq |\gamma - 1/2|$. Indeed:

1. $\tilde{\gamma} = 0, \gamma = 1 \Rightarrow \gamma \leq 1/2, \tilde{\gamma} > 1/2$
2. $\tilde{\gamma} = 1, \gamma = 0 \Rightarrow \gamma > 1/2, \tilde{\gamma} \leq 1/2$

Case 1: $\gamma = 1/2 - \delta, \tilde{\gamma} = 1/2 + \tilde{\delta} \therefore |\gamma - 1/2| = \delta$
 $|\tilde{\gamma} - \gamma| = \delta + \tilde{\delta} \quad]$

Case 2: $\gamma = 1/2 + \delta, \tilde{\gamma} = 1/2 - \tilde{\delta} \therefore |\gamma - 1/2| = \delta$
 $|\tilde{\gamma} - \gamma| = \delta + \tilde{\delta} \quad] \quad \square$

Suppose $\tilde{\gamma}_1(x) \approx \gamma(x), \tilde{\gamma}_0(x) \approx 1 - \gamma(x)$ (approximations).
Define $g(x) = \begin{cases} 0 & \text{if } \tilde{\gamma}_1 \leq \tilde{\gamma}_0 \\ 1 & \text{otherwise} \end{cases}$.

We don't have necessarily $\tilde{\gamma}_1 + \tilde{\gamma}_0 = 1$. Thus the previous result differs slightly:

$$\begin{aligned} \text{Theo. } L(g) - L(g^*) &\leq \int_{\mathbb{R}^d} |(1-\gamma) - \tilde{\gamma}_0| \mu(dx) \\ &\quad + \int_{\mathbb{R}^d} |\gamma - \tilde{\gamma}_1| \mu(dx) \end{aligned}$$

$$\text{Recall that } \gamma(x) = \mathbb{P}(y=1|X) = \frac{p(x|y=1) \mathbb{P}(y=1)}{p(x)} = \frac{f_1 p}{f} \quad (7)$$

$$1 - \gamma(x) = \mathbb{P}(y=0|x) = \dots = \frac{f_0(1-p)}{f}$$

Let $\tilde{f}_1 \approx f_1$, $\tilde{p}_1 \approx p$, $\tilde{f}_0 \approx f_0$, $\tilde{p}_0 \approx 1-p$.

Define $g(x) = \begin{cases} 0 & \text{if } \tilde{p}_1 \tilde{f}_1 \leq \tilde{p}_0 \tilde{f}_0 \\ 1 & \text{otherwise.} \end{cases}$

Then from the previous theorem we have

$$|(1-\gamma) - \tilde{\gamma}_0| = \left| \frac{(1-p)f_0 - \tilde{p}_0 \tilde{f}_0}{f} \right|$$

$$|\gamma - \tilde{\gamma}_1| = \left| \frac{pf_1 - \tilde{p}_1 \tilde{f}_1}{f} \right|$$

$$h(dx) = r(x) dx$$

$$\leq L - L^+ \leq \int_{\mathbb{R}^d} |(1-p)f_0 - \tilde{p}_0 \tilde{f}_0| dx + \int_{\mathbb{R}^d} |pf_1 - \tilde{p}_1 \tilde{f}_1| dx$$

Role of Dimension: If we knew $\gamma(x)$ the prob. would be 1D since we could just forget about X altogether.

However, γ is unknown. In general, every transformation $x \rightarrow T(x)$ increases the Bayes risk, since it destroys information. Nevertheless, there are transformations that leaves the Bayes error unchanged.