

2.11 Show that if  $(x, y) \mapsto (T(x), y)$  then  $L^*(T(x), y) \geq L^*(x, y)$ .

Recall Th. 1 which states that for any classifier  $g$  we have  $L(g^*) = \text{IP}[g^*(x) \neq y] \leq L(g) = \text{IP}[g(x) \neq y]$ , i.e  $g^*$  is optimal.

Let  $T: X \mapsto Z = T(X)$  be a transformation. In the "Z-space" the Bayes classifier is defined by

$$h^*(z) = \begin{cases} 1 & \text{if } \sigma(z) > 1/2 \\ 0 & \text{otherwise} \end{cases}$$

where  $\sigma(z) = \text{IP}[y=1 | z]$ , and its error is  $L(h^*) = \text{IP}[h^*(z) \neq y]$ . Now suppose that  $L(h^*) < L(g)$ . Then we would be able to construct the following classifier in the original "X-Space":

$$g(x) = \begin{cases} 1 & \text{if } \gamma_0 T(x) = \gamma(T(x)) > 1/2 \\ 0 & \text{otherwise} \end{cases}$$

where  $\gamma = \text{IP}[y=1 | x]$  and  $L(g) < L(g^*)$  by assumption. But this contradicts the above mentioned Th. 1, therefore we must have  $L(h^*) \leq L(g)$ .

In other words, we can incorporate the transformation in a classifier but this cannot improve Bayes error since it is optimal.

(2)

2.2 Let  $x'$  be independent of  $(x, y)$ . Show that  $L_{(x, x')}^*(g) = L_x^*$ .

Recall  $L_x(g^*) = \mathbb{P}[g(x) \neq y]$  and  $g^*(x) = \begin{cases} 1, & \gamma(x) > \frac{1}{2} \\ 0, & \text{otherwise} \end{cases}$  where  $\gamma(x) = \mathbb{P}[y=1|x]$ .

Analogously,  $L_{(x, x')}(h^*(x, x') \neq y)$  with  $h^* = \begin{cases} 1, & \delta(x, x') > \frac{1}{2} \\ 0, & \text{otherwise} \end{cases}$  where  $\delta(x, x') = \mathbb{P}[y=1|x, x']$ .

Now if  $x' \perp (x, y)$ , then  $\mathbb{P}[y|x, x'] = \mathbb{P}[y|X]$ , thus  $h^*(x, x') = g^*(x)$  implying  $L_{(x, x')}^* = L_x^*$ .

2.5 For  $c \in (0, \frac{1}{2})$  define

$$g_c(x) = \begin{cases} 1 & \text{if } \gamma(x) > \frac{1}{2} + c \\ 0 & \text{if } \gamma(x) \leq \frac{1}{2} - c \\ \text{"reject"} & \text{otherwise} \end{cases}$$

Show that, for any  $g$ , if

$$\mathbb{P}[g(x) = \text{"reject"}] \leq \mathbb{P}[g_c(x) = \text{"reject"}]$$

then

$$\mathbb{P}[g(x) \neq y | g(x) \neq \text{"reject"}] \geq \mathbb{P}[g_c(x) \neq y | g_c(x) \neq \text{"reject"}].$$

Thus decisions under  $g_c$  are optimal.

Notice that

$$\begin{aligned} \mathbb{P}[g(x) \neq y | g(x) \neq R] &= 1 - \mathbb{P}[g(x) = y | g(x) \neq R] \\ &= 1 - \frac{\mathbb{P}[g=y, g \neq R]}{\mathbb{P}[g \neq R]} \end{aligned}$$

]

③

$$\begin{aligned} \text{IP}[g \neq y | g \neq R] &= 1 - \text{IP}[g = y | g \neq R] \\ &= 1 - \frac{\text{IP}[g = y, g \neq R]}{\text{IP}[g \neq R]} \\ &= 1 - \frac{\mathbb{I}(g \neq R) \text{IP}[g = y]}{\text{IP}[g \neq R]} \end{aligned}$$

thus

$$\begin{aligned} \text{IP}[g \neq y | g \neq R] - \text{IP}[g_c \neq y | g_c \neq R] &= \frac{\mathbb{I}(g_c \neq R) \text{IP}[g_c = y]}{\text{IP}[g_c \neq R]} \\ &\quad - \frac{\mathbb{I}(g \neq R) \text{IP}[g = y]}{\text{IP}[g \neq R]} \end{aligned}$$

It is enough to show that

$$\frac{\text{IP}[g_c = y | X]}{\text{IP}[g_c \neq R]} - \frac{\text{IP}[g = y | X]}{\text{IP}[g \neq R]} > 0 \quad (\#)$$

$$\begin{aligned} \text{Consider } \text{IP}[g = y | X] &= \text{IP}[g = 1, y = 1 | X] + \text{IP}[g = 0, y = 0 | X] \\ &= \mathbb{I}(g = 1) \gamma(x) + \mathbb{I}(g = 0) (1 - \gamma(x)) \end{aligned}$$

Thus we have

$$\frac{\mathbb{I}(g_c = 1) \gamma + \mathbb{I}(g_c = 0) (1 - \gamma)}{\text{IP}[g_c \neq R]} - \frac{\mathbb{I}(g = 1) \gamma + \mathbb{I}(g = 0) (1 - \gamma)}{\text{IP}[g \neq R]} \quad (*)$$

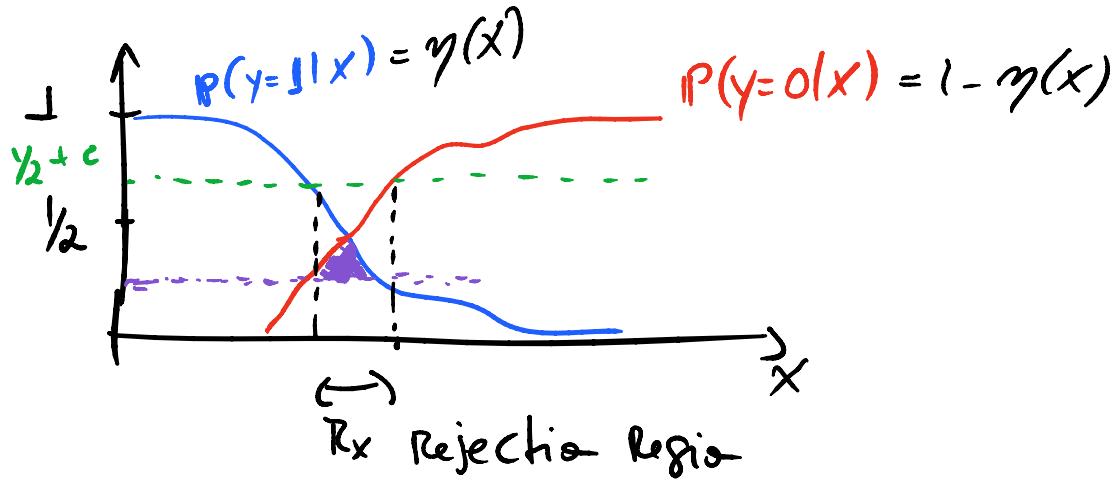
$$\text{If } g = g_c \text{ thus gives } \left( \frac{1}{\text{IP}[g_c \neq R]} - \frac{1}{\text{IP}[g \neq R]} \right) \cdot c > 0,$$

$$\text{where } c = \gamma \text{ or } c = 1 - \gamma \text{ (both are } > 0 \text{). If } g_c = 1 \text{ and } g = 0 \text{ we get } \frac{\gamma}{\text{IP}[g_c \neq R]} - \frac{(1 - \gamma)}{\text{IP}[g \neq R]} > 0 \text{ since}$$

$\gamma > 1 - \gamma$  and  $\text{IP}[g_c \neq R] < \text{IP}[g \neq R]$ . Analogous argument works for  $g_c = 0$  and  $g = 1$ . Therefore, in any case,  $(*) > 0$ , thus  $(\#) > 0$ . Integrating over  $X$  gives the claim.]

The intuition behind this is the following.

⑨



We assign  $x$  to class  $\{1, 0\}$  according to  $\max\{\gamma(x), 1-\gamma(x)\} > y_2 + c$ , and reject any decision otherwise. We only incur an error if we make a wrong assignment when we do not reject. The error comes when  $P[y=1|x]$  and  $P[y=0|x]$  have similar values and overlap, as indicated in the PURPLE shaded area. Thus if we lower the rejection threshold,  $P[g \neq \text{Reject}]$ , more points  $x$  will have non-reject decisions in the shaded area which increases the misclassification error.

J

(5)

$$\underline{2.9} \quad g(x) = \begin{cases} 0 & \text{if } \tilde{\eta}_1(x) \leq \tilde{\eta}_0(x) \\ 1 & \text{otherwise} \end{cases}$$

$$\begin{aligned} \tilde{\eta}_1(x) &\sim \eta(x) && \text{It may be the case that} \\ \tilde{\eta}_0(x) &\sim 1 - \eta(x) && \tilde{\eta}_0 + \tilde{\eta}_1 \neq 1 \end{aligned}$$

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

PROOF

$$\begin{aligned} P[g(x) \neq y | x] &= 1 - P[g(x) = y | x] \\ &= 1 - P[g=1, y=1 | x] - P[g=0, y=0 | x] \\ &= 1 - \mathbf{1}(g=1) P[y=1 | x] - \mathbf{1}(g=0) P[y=0 | x] \\ &= 1 - \mathbf{1}(g=1) \tilde{\eta}_1(x) - \mathbf{1}(g=0) \tilde{\eta}_0(x) \end{aligned}$$

$$P[g^*(x) \neq y | x] = 1 - \mathbf{1}(g^*=1) \eta(x) - \mathbf{1}(g^*=0) (1 - \eta(x))$$

Therefore,

$$\begin{aligned} L(g) - L(g^*) &= \int \mu(dx) (\eta(x) \mathbf{1}(g^*=1) - \tilde{\eta}_1(x) \mathbf{1}(g=1)) \\ &\quad + \int \mu(dx) ((1 - \eta(x)) \mathbf{1}(g^*=0) - \tilde{\eta}_0(x) \mathbf{1}(g=0)) \end{aligned}$$

]

(6)

$$L(g) - L(g^*) = \int \mu(dx) (\eta \mathbb{1}(\eta > \frac{1}{2}) - \tilde{\eta}_1 \mathbb{1}(\tilde{\eta}_1 > \tilde{\eta}_0)) \\ + \int \mu(dx) ((1-\eta) \mathbb{1}(\eta \leq \frac{1}{2}) - \tilde{\eta}_0 \mathbb{1}(\tilde{\eta}_1 \leq \tilde{\eta}_0))$$

Now we consider each possibility.

1.  $\eta > \frac{1}{2}$ ,  $\tilde{\eta}_1 > \tilde{\eta}_0$ . The above integrand is  
 $|\eta - \tilde{\eta}_1| \leq |\eta - \tilde{\eta}_1| \leq |\eta - \tilde{\eta}_1| + |(1-\eta) - \tilde{\eta}_0|$

2.  $\eta \leq \frac{1}{2}$ ,  $\tilde{\eta}_1 \leq \tilde{\eta}_0$ . The integrand is  
 $|(1-\eta) - \tilde{\eta}_0| \leq |(1-\eta) - \tilde{\eta}_0| \leq |(1-\eta) - \tilde{\eta}_0| + |\eta - \tilde{\eta}_1|$

3.  $\eta > \frac{1}{2}$ ,  $\tilde{\eta}_1 \leq \tilde{\eta}_0$ . The integrand is  
 $|\eta - \tilde{\eta}_0| \leq |\eta - \tilde{\eta}_1| + |(1-\eta) - \tilde{\eta}_0|$

4.  $\eta \leq \frac{1}{2}$ ,  $\tilde{\eta}_1 > \tilde{\eta}_0$ . The integrand is  
 $|(1-\eta) - \tilde{\eta}_1| \leq |(1-\eta) - \tilde{\eta}_0| \leq |(1-\eta) - \tilde{\eta}_0| + |\eta - \tilde{\eta}_1|$

Thus in either case we get

$$L(g) - L(g^*) \leq \int \mu(dx) |\eta - \tilde{\eta}_1| + \int \mu(dx) |(1-\eta) - \tilde{\eta}_0|$$

as desired.

]

## TRUNK's PROBLEM

(7)

Consider a two-class problem where

$$p(x|w_1) = N(x|\mu, I)$$

$$p(x|w_2) = N(x|-\mu, I)$$

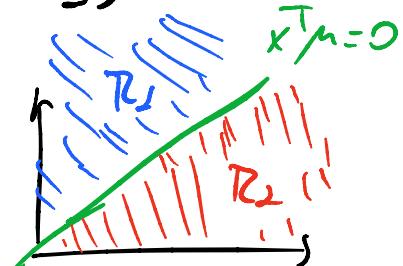
with  $p(w_1) = p(w_2) = \frac{1}{2}$  and  $\mu \in \mathbb{R}^D$  with elements given by  $\mu_i = \sqrt{i}$ ,  $i=1, \dots, D$ .

The decision boundary is determined by:

$$\underline{p(w_1|x)} > \underline{p(w_2|x)} \quad (x \in C_1)$$

$$\frac{p(x|w_1)p(w_1)}{p(x)} > \frac{p(x|w_2)p(w_2)}{p(x)}$$

$$\frac{p(x|w_1)}{p(x|w_2)} > 1$$



$$\text{or } \log p(x|w_1) - \log p(x|w_2) > 0$$

$$-\|x - \mu\|^2 + \|x + \mu\|^2 > 0$$

$$-\|x\|^2 + 2x^T \mu - \|\mu\|^2 + \|x\|^2 + 2x^T \mu + \|\mu\|^2 > 0$$

$$\boxed{x^T \mu > 0}$$

Thus the decision rule is  $\begin{cases} x \in C_1 & \text{if } x^T \mu > 0 \\ x \in C_2 & \text{if } x^T \mu \leq 0 \end{cases}$ .

The missclassification error is

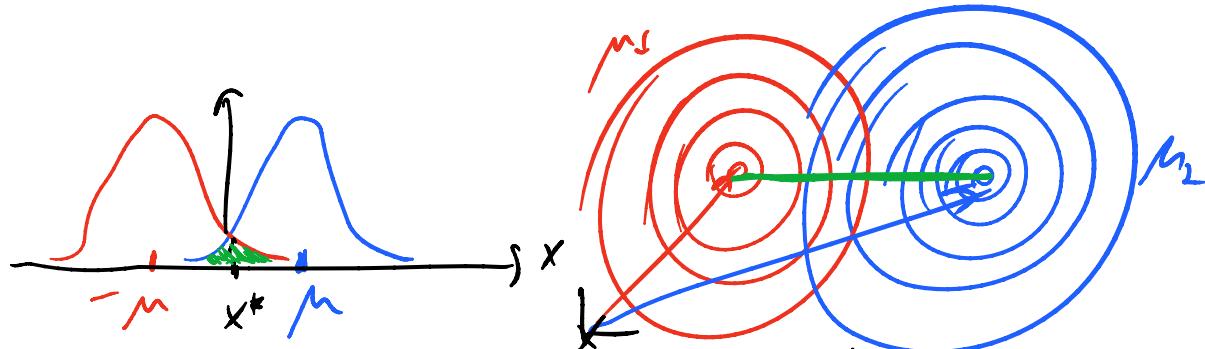
$$P_e = P(X \in R_1, x \in C_2) + P(X \in R_2, x \in C_1)$$

$$= \int_{R_1} p(x, w_2) dx + \int_{R_2} p(x, w_1) dx$$

]

$$P_e = \frac{1}{2} \int_{R_1} p(x|w_2) dx + \frac{1}{2} \int_{R_2} p(x|w_1) dx$$

$$= \frac{1}{2} \int_{R_1} \frac{1}{(2\pi)^{D/2}} e^{-\frac{1}{2}\|x - \mu_2\|^2} dx + \frac{1}{2} \int_{R_2} \frac{1}{(2\pi)^{D/2}} e^{-\frac{1}{2}\|x - \mu_1\|^2} dx$$



We can write  $x$  in a coordinate system where one of the axis is in the direction  $\mu_1 - \mu_2$ , the other ones being orthogonal, i.e.,  $\mu = (\|\mu\|, 0, \dots, 0)^T$  and  $X = (X_1, X_2, \dots, X_D)$ . Then only the first component contribute. We thus have

$$\begin{aligned} P_e &= \frac{1}{2} \int_0^\infty \frac{1}{(2\pi)^{D/2}} e^{-\frac{1}{2}(x+\mu)^2} dx + \frac{1}{2} \int_{-\infty}^0 \frac{1}{(2\pi)^{D/2}} e^{-\frac{1}{2}(x-\mu)^2} dx \\ &= \frac{1}{2} \int_0^\infty \frac{1}{\mu (2\pi)^{D/2}} e^{-\frac{1}{2}z^2} dz + \frac{1}{2} \int_{-\infty}^{\mu} \frac{1}{(2\pi)^{D/2}} e^{-\frac{1}{2}z^2} dz \\ &= \int_\mu^\infty \frac{1}{(2\pi)^{D/2}} e^{-\frac{1}{2}z^2} dz \end{aligned}$$

where  $\mu \equiv \|\mu\| = \left( \sum_{m=1}^D \frac{1}{m} \right)$ .

Thus,  $\lim_{D \rightarrow \infty} P_e = 0$ . This would not happen if the components of  $\mu$  decay faster than  $\frac{1}{\sqrt{m}}$ .

①

Now suppose we have  $m_1$  samples from  $p(x|w_1)$  and  $m_2$  samples from  $p(x|w_2)$ :

$$\{x_1, \dots, x_{m_1}\} \rightarrow \hat{\mu}_1 = \frac{1}{m_1} \sum_i x_i = \hat{\mu}$$

$$\{\tilde{x}_1, \dots, \tilde{x}_{m_2}\} \quad \hat{\mu}_2 = \frac{1}{m_2} \sum_i \tilde{x}_i = -\hat{\mu}$$

or  $\tilde{x}_i \rightarrow -x_i$  and we have  $\{x_1, \dots, x_N\}$ ,  $N = m_1 + m_2$ , elements. Then

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i \quad (\text{MLE estimator}).$$

The prob. of error is

$$\begin{aligned} P_e &= P(x \in R_1, x \in G_2) + P(x \in R_2, x \in G_1) \\ &= \frac{1}{2} P(x^T \hat{\mu} > 0 | w_2) + \frac{1}{2} P(x^T \hat{\mu} \leq 0 | w_1) \\ &= P(x^T \hat{\mu} > 0 | w_2) \\ &\stackrel{\text{iid}}{=} \end{aligned}$$

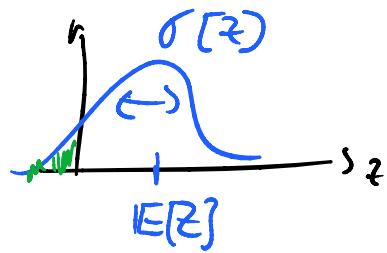
$$\begin{aligned} \overline{E[z]} &= \frac{1}{N} \sum_{i=1}^N E[x^T x_i] = \frac{1}{N} \sum_{i=1}^N E[\bar{x}^T] E[x_i] \\ &= \frac{1}{N} N \mu^T \mu = \|\mu\|^2 = \underbrace{\sum_{m=1}^D \frac{1}{m}}_{\text{iid}} \end{aligned}$$

$$\text{Var}(z) = \left(1 + \frac{1}{N}\right) \sum_{m=1}^D \frac{1}{m} + \frac{D}{N} \quad (\text{show this ???})$$

10

thus the error is

$$\hat{P}_e = \int_{-\infty}^0 \frac{1}{(2\pi)^{1/2}} \frac{e^{-\frac{1}{2}(\frac{z - E[z]}{\sigma[z]})^2}}{\sigma[z]} dz$$



Let  $\frac{z - E[z]}{\sigma[z]} \equiv y$ . Notice that  $y \sim N(0, 1)$  according to the CLT

Then,

$$\hat{P}_e = \int_{-\infty}^0 \frac{1}{(2\pi)^{1/2}} e^{-\frac{1}{2}y^2} dy = \int_{\frac{E[z]}{\sigma[z]}}^{\infty} \frac{1}{(2\pi)^{1/2}} e^{-\frac{1}{2}y^2} dy$$

Notice that  $\sigma^2(z) = O(\frac{1}{N})$ , thus  $\frac{|E(z)|}{\sigma(z)} = O\left(\frac{N^{1/2}}{D^{1/2}}\right) \sum_{m=1}^D \frac{1}{m}$ .

Now, by definition

$$\lim_{D \rightarrow \infty} \left( \sum_{m=1}^D \frac{1}{m} - \log D \right) = \gamma \quad (\text{Euler-Mascheroni constant})$$

Therefore,

$$\lim_{D \rightarrow \infty} \left( \frac{N}{D} \right)^{1/2} \sum_{m=1}^D \frac{1}{m} - \underbrace{\frac{N^{1/2} \log D}{D}}_{\rightarrow 0} = \underbrace{\left( \frac{N}{D} \right)^{1/2} \gamma}_{\rightarrow 0}$$

$\therefore \frac{|E(z)|}{\sigma(z)} \rightarrow 0 \text{ as } D \rightarrow \infty$ , which implies that

$$\boxed{\hat{P}_e \rightarrow \frac{1}{2} \text{ as } D \rightarrow \infty} !$$

1

11

## Comments:

- (1) This shows that, for this example, if  $\mu$  is known theoretically, the Bayes error  $P_e \rightarrow 0$  as  $D \rightarrow \infty$ . However, if we estimate  $\mu$  by N-size sample, the plugin error  $\hat{P}_e \rightarrow 1/2$  as  $D \rightarrow \infty$ .
- (2) The example is of course rather artificial.  
Does this happen in real data?  
Notice that if the components of  $\mu$  decay fast, for instance if  $\mu_i = \frac{1}{i^{\gamma} + \epsilon}$  with  $\epsilon > 0$ , the series would not diverge and  $P_e \not\rightarrow 0$ .
- (3) In the plugin error, we are assuming  $\frac{N}{D} \ll 0$ . If we have enough data,  $\frac{N}{D} = O(1)$ , and  $\hat{P}_e < 1/2$ . By increasing  $N$ ,  $\hat{P}_e$  should approach  $P_e$ .

1