

Remark: For all the 3 examples ECM (or TPM) of the ECM minimizing classification rules can be calculated once we calculate $P(1|2)$ & $P(2|1)$

e.g.: Consider "Example 1" of $\Pi_1: N_p(\underline{\mu}_1, \Sigma_1)$

$$\Pi_2: N_p(\underline{\mu}_2, \Sigma_2)$$

$$P(1|2) = P_{\Pi_2} \left(\tilde{x} \in R_1^* \right).$$

$$= P_{\Pi_2} \left((\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \tilde{x} \geq \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) + \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \right)$$

Note that under Π_2 :

$$\begin{aligned} z_1 &= (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \tilde{x} \sim N_1 \left((\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \underline{\mu}_2, \right. \\ &\quad \left. (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} (\underline{\mu}_1 - \underline{\mu}_2) \right) \\ &= \Delta \leftarrow \text{Eq. of Mahalanobis distance} \\ \Rightarrow P(1|2) &= P_{\Pi_2} \left(\frac{z_1 - (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \underline{\mu}_2}{\Delta} \geq \left[\frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \right. \right. \\ &\quad \left. \left. - (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \underline{\mu}_2 + \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) \right] \Delta^{-1} \right) \\ \text{let } y &= \frac{z_1 - (\underline{\mu}_1 - \underline{\mu}_2)' \tilde{\Sigma}^{-1} \underline{\mu}_2}{\Delta} \end{aligned}$$

$$P(1|2) = P_{\Pi_2} \left(y \geq \frac{1}{2} \Delta + \frac{1}{\Delta} \log \left(\frac{p_2 C(1|2)}{p_1 C(2|1)} \right) \right)$$

$$\Rightarrow P(1|2) = 1 - \Phi \left(\frac{\Delta}{2} + \frac{1}{\Delta} \log \left(\frac{P_2 C(1|2)}{P_1 C(2|1)} \right) \right)$$

∴ $P(2|1)$ can be calculated

Note: TPM / ECM can be estimated by estimating Δ from learning set L .

Remark: Performance measure for classifier

Apparent error rate (APER)

- Calculate the "confusion matrix"

		Predicted class		
		π_1	π_2	
Actual class	π_1	n_{1C}	n_{1M}	$(n_{1C} + n_{1M} = n_1)$
	π_2	n_{2M}	n_{2C}	$(n_{2M} + n_{2C} = n_2)$

n_{1M} : # of cases from π_1 misclassified by

n_{1C} : # of - - - π_1 correctly classified

$$\text{APER} = \frac{n_{1M} + n_{2M}}{n_1 + n_2}$$

Classification in multiclass problem

Let $\pi_1, \pi_2, \dots, \pi_c$ be C classes with prior probabilities $p(\pi_1), p(\pi_2), \dots, p(\pi_c)$, respectively.

Bayes classifier

Assign \tilde{x} to π_j If the posterior prob

of the class π_j given the obsn \tilde{x} , i.e.,

$p(\pi_j | \tilde{x})$ is the highest over all classes

π_1, \dots, π_c

i.e. Assign \tilde{x} to π_j If

$$p(\pi_j | \tilde{x}) \geq p(\pi_k | \tilde{x}) ; k=1, \dots, c \\ k \neq j$$

$$\Rightarrow \frac{p(\pi_j) p(\tilde{x} | \pi_j)}{\sum_{i=1}^c p(\pi_i) p(\tilde{x} | \pi_i)} \geq \frac{p(\pi_k) p(\tilde{x} | \pi_k)}{\sum_{i=1}^c p(\pi_i) p(\tilde{x} | \pi_i)}$$

$$\Rightarrow p(\pi_j) p(\tilde{x} | \pi_j) \geq p(\pi_k) p(\tilde{x} | \pi_k)$$

\Rightarrow as in the 2-class problem (done in class), Bayes classifier for a $C (\geq 2)$ -class problem can be expressed in terms of prior ~~and~~ prob & class conditional densities

TPM minimizing classifier

TPM for a general c-class problem can be defined as

$$\text{TPM} = \sum_{i=1}^c p(\pi_i) P(\text{error} | \pi_i)$$

↑
 misclassification
 i.e. $P(\text{error} | \pi_i)$: prob of
 misclassifying a pattern
 from π_i

Let $\{R_1, R_2, \dots, R_c\}$ be the classification partition, then

$$P(\text{error} | \pi_i) = \int_{R_i^c} f(\underline{x} | \pi_i) d\underline{x}$$

$$\Rightarrow \text{TPM} = \sum_{i=1}^c \int_{\mathcal{X} - R_i} f(\underline{x} | \pi_i) p(\pi_i) d\underline{x}$$

$$= \sum_{i=1}^c p(\pi_i) \int_{\mathcal{X} - R_i} f(\underline{x} | \pi_i) d\underline{x}$$

$$= \sum_{i=1}^c p(\pi_i) \left(\int_{\mathcal{X}} f(\underline{x} | \pi_i) d\underline{x} - \int_{R_i} f(\underline{x} | \pi_i) d\underline{x} \right)$$

$$= \sum_{i=1}^c p(\pi_i) \left(1 - \int_{R_i} f(\underline{x} | \pi_i) d\underline{x} \right)$$

i.e.

$$\text{TPM} = 1 - \sum_{i=1}^C \int_{R_i} p(\pi_i) f(\underline{x} | \pi_i) d\underline{x}$$

\Rightarrow Minimizing the TPM w.r.t. partition $\{R_1, \dots, R_C\}$ is equivalent to maximizing

$$\sum_{i=1}^C \int_{R_i} p(\pi_i) f(\underline{x} | \pi_i) d\underline{x} \text{ w.r.t. } \{R_1, \dots, R_C\}$$

\longleftrightarrow

This can be viewed as prob of correct classification

Maximization of the above is achieved by selecting R_i to be the region for which

$p(\pi_i) f(\underline{x} | \pi_i)$ is the largest among all classes

(i.e. among all $i = 1 \dots C$)

i.e. the TPM minimizing classification rule

is

Assign \underline{x} to π_i if

$$p(\pi_i) f(\underline{x} | \pi_i) \geq p(\pi_k) f(\underline{x} | \pi_k) \quad \forall k \neq i$$

Remark: As in the 2-class problem, the TPM minimizing rule is equivalent to the one obtained by maximizing the "posterior prob", i.e. the Bayes classifier

Minimum ECM classification rule

Let $p(\pi_i) = p_i \quad i = 1(1)c$

class conditional densities $f(\tilde{x} | \pi_i); i = 1(1)c$

$c(k|i)$: cost of misclassifying an item/pattern
from π_i to $\pi_k; k, i = 1(1)c$

$c(i|i) = 0 \quad i = 1(1)c$

$\{R_1, \dots, R_c\}$: classification partition

$$P(k|i) = P(\text{Classifying an item to } \pi_k | \pi_i) \\ = \int_{R_k} f(\tilde{x} | \pi_i) d\tilde{x}$$

$$P(i|i) = \int_{R_i} f(\tilde{x} | \pi_i) d\tilde{x} = \int_{\tilde{X} - \cup_{k \neq i} R_k} f(\tilde{x} | \pi_i) d\tilde{x} \\ = 1 - \sum_{\substack{k=1 \\ k \neq i}}^c P(k|i)$$

Note that the conditional ECM of \tilde{x} from π_1 ,
into π_2 or π_3 or ... or π_c is

$$\text{ECM}_1 = P(2|1)C(2|1) + P(3|1)C(3|1) + \dots + P(c|1)C(c|1) \\ = \sum_{k=2}^c P(k|1)C(k|1)$$

The above ECM_i is with prob $p(\pi_i) = p_i$

Thus for other conditional ECM's ECM_j

We can write similar expressions w.p. p_j
(as below)

\Rightarrow

$$\text{ECM} = p_1 \text{ECM}_1 + p_2 \text{ECM}_2 + \dots + p_c \text{ECM}_c$$

$$\text{i.e. } \text{ECM} = p_1 \left(\sum_{k=1}^c p(k|1) c(k|1) \right)$$

$$+ p_2 \left(\sum_{k=1}^c p(k|2) c(k|2) \right) + \dots$$

$$+ p_c \left(\sum_{k=1}^{c-1} p(k|c) c(k|c) \right)$$

$$\text{i.e. } \text{ECM} = \sum_{i=1}^c p_i \sum_{k=1}^c p(k|i) c(k|i)$$

$$K \neq i$$

$$= \sum_{k=1}^c \sum_{\substack{i=1 \\ i \neq k}}^c p_i p(k|i) c(k|i)$$

$$= \sum_{k=1}^c \sum_{\substack{i=1 \\ i \neq k}}^c \int_{R_K} p_i c(k|i) f(x| \pi_i) dx$$

$$\text{i.e. } \text{ECM} = \sum_{K=1}^C \int \left(\sum_{\substack{i=1 \\ i \neq K}}^C p_i c(k|i) f(\underline{x}) \pi_i \right) d\underline{x}$$

$$= \sum_{K=1}^C \int_{R_K} h_K(\underline{x}) d\underline{x}$$

Note that

$$\sum_{K=1}^C \int_{R_K} \left(h_K(\underline{x}) - \min_j h_j(\underline{x}) \right) d\underline{x} \geq 0$$

with equality only if

$$h_K(\underline{x}) = \underline{\min_j h_j(\underline{x})} \quad \underline{\text{if } \underline{x} \text{ in } R_K}$$

Thus the ECM minimizing classification rule
is

Assign \underline{x} to π_K if

$$\sum_{\substack{i=1 \\ i \neq K}}^C p_i c(k|i) f(\underline{x}) \pi_i$$

is minimum among all such C expressions

Remark: Under equal cost setup, i.e. all misclassification costs are same

We assign \underline{x} to π_K if

$$\sum_{\substack{i=1 \\ i \neq K}}^C p_i f(\underline{x}) \pi_i$$

is the smallest

$$\text{i.e. } \sum_{\substack{i=1 \\ i \neq k}}^c p_i f(\underline{x} | \pi_i) \leq \sum_{\substack{i=1 \\ i \neq j \\ i \neq k}}^c p_i f(\underline{x} | \pi_i) \quad (*)$$

Subtracting $\sum_{\substack{i=1 \\ i \neq k, j}}^c p_i f(\underline{x} | \pi_i)$ from both sides

of (*), we get

$$\begin{aligned} & \sum_{\substack{i=1 \\ i \neq k}}^c p_i f(\underline{x} | \pi_i) - \sum_{\substack{i=1 \\ i \neq k, j}}^c p_i f(\underline{x} | \pi_i) \\ & \leq \sum_{\substack{i=1 \\ i \neq j}}^c p_i f(\underline{x} | \pi_i) - \sum_{\substack{i=1 \\ i \neq k, j}}^c p_i f(\underline{x} | \pi_i) \end{aligned}$$

$$\text{i.e. } p_j f(\underline{x} | \pi_j) \leq p_k f(\underline{x} | \pi_k) \quad \forall j \neq k$$

i.e. the rule is

assign \underline{x} to π_k if

$p_k f(\underline{x} | \pi_k)$ is largest

This is the same as TPM minimizing rule
(What we expect under equal misclassification cost).

Further, the above is (under equal cost)
equivalent to Bayes classifier.

Multiclass classification problem

Example 1: 3-class problem

Consider the following cost-prior table

		True membership			mis class ↳ cost matrix
		π_1	π_2	π_3	
π_i	π_1	$c(1 1) = 0$	$c(1 2) = 500$	100	
	π_2	10	0	50	
	π_3	50	200	0	
prior prob		$p_1 = 0.05$	$p_2 = 0.60$	$p_3 = 0.35$	

Let \tilde{x}_0 be a new obsn \Rightarrow

$$f(\tilde{x}_0 | \pi_1) = 0.01 ; f(\tilde{x}_0 | \pi_2) = 0.85 \\ f(\tilde{x}_0 | \pi_3) = 2$$

ECM minimizing classifier

Compute $\sum_{\substack{i=1 \\ i \neq k}}^3 p_i c(k|i) f(\tilde{x}_0 | \pi_i)$ $\forall k = 1, 2, 3$

$$\underline{k=1} : \sum_{i=2}^3 p_i c(1|i) f(\tilde{x}_0 | \pi_i) \\ = p_2 c(1|2) f(\tilde{x}_0 | \pi_2) + p_3 f(\tilde{x}_0 | \pi_3) c(1|3) \\ = 325$$

$$\underline{K=2} : \sum_{\substack{i=1 \\ i \neq 2}}^3 p_i c(2|i) f(\tilde{x}_0 | \pi_i)$$

$$= p_1 c(2|1) f(\tilde{x}_0 | \pi_1) + p_3 c(2|3) f(\tilde{x}_0 | \pi_3)$$

$$= 35.06$$

$$\underline{K=3} : \sum_{i=1}^2 p_i c(3|i) f(\tilde{x}_0 | \pi_i) = 102.03$$

Since $\sum_{\substack{i=1 \\ i \neq 2}}^3 p_i c(2|i) f(\tilde{x}_0 | \pi_i)$ is the smallest

we assign \tilde{x}_0 to π_2

Note: Suppose in the same example, we take equal cost \Rightarrow TPM minimizing rule

$$\text{Calculate } p_1 f(\tilde{x}_0 | \pi_1) = 0.0005$$

$$p_2 f(\tilde{x}_0 | \pi_2) = 0.510$$

$$p_3 f(\tilde{x}_0 | \pi_3) = 0.7$$

Since $p_3 f(\tilde{x}_0 | \pi_3) > p_i f(\tilde{x}_0 | \pi_i) \quad i=1, 2$

\tilde{x}_0 is allocated to π_3

Bayes classifier would also have done the same assignment

Example 3 : Discrete populations

Consider 3 bivariate discrete populations with the following jt p.m.f.s :

		Π_1		Π_2		Π_3			
		x_1	x_2	x_1	x_2	x_1	x_2		
1	1	0.5	0.2	1	0.2	0.1	1	0.25	0.25
	2	0.1	0.2	2	0.3	0.4	2	0.25	0.25

Prior probabilities are equal, i.e. $p(\Pi_i) = \frac{1}{3}; i=1,2,3$
 $(= p_i)$

Misclassification costs are

$$c(1|i) = 1; i=2,3$$

$$c(2|i) = 2; i=1,3$$

$$c(3|i) = 3; i=1,2$$

Let us denote by $f_i(x)$ to be class conditional prob for class $\Pi_i; i=1,2,3$.

ECM classification rule is :

Assign \tilde{x} to Π_k if

$$\sum_{\substack{i=1 \\ i \neq k}}^3 p_i f_i(\tilde{x}) c(k|i) \text{ is smallest}$$

Based on the above find rule for all pairs \tilde{x}

pair (1,1)

K=1

$$\sum_{i=2}^3 p_i f_i(\underline{x}) \underset{\substack{\downarrow \\ \underline{x} = (1,1)}}{c(1|i)}$$

$$= p_2 f_2(\underline{x}) c(1|2) + p_3 f_3(\underline{x}) c(1|3)$$

$$= \frac{1}{3} (f_2(\underline{x}) + f_3(\underline{x}))$$

$$= \frac{1}{3} (0.2 + 0.25) = 0.15$$

K=2

$$\sum_{\substack{i=1 \\ i \neq 2}}^3 p_i f_i(\underline{x}) \underset{\substack{\downarrow \\ \underline{x} = (1,1)}}{c(2|i)}$$

$$= p_1 f_1(\underline{x}) c(2|1) + p_3 f_3(\underline{x}) c(2|3)$$

$$= \frac{2}{3} (f_1(\underline{x}) + f_3(\underline{x})) = \frac{2}{3} (0.5 + 0.25)$$

$$= 0.5$$

K=3

$$\sum_{i=1}^2 p_i f_i(\underline{x}) c(3|i)$$

$$= p_1 f_1(\underline{x}) c(3|1) + p_2 f_2(\underline{x}) c(3|2)$$

$$= \frac{3}{3} (f_1(\underline{x}) + f_2(\underline{x})) = 0.7$$

$\sum_{\substack{i=1 \\ i \neq 1}}^3 p_i f_i(\underline{x}) c(1|i)$ is smallest

\Rightarrow Assign $\underline{x}^{(1,1)}$ to π_1
i.e. $(1,1) \in R_1$ partition

Pair (1, 2)

$$\underline{k=1} \rightarrow \frac{1}{3} (f_2(1, 2) + f_3(1, 2)) = \frac{1}{3} (0.55) \leftarrow \text{smallest}$$

$$\underline{k=2} \rightarrow \frac{2}{3} (f_1(\overset{(1, 2)}{1}, 2) + f_3(1, 2)) = \frac{2}{3} (0.35)$$

$$\underline{k=3} \rightarrow \frac{3}{3} (f_1(1, 2) + f_2(1, 2)) = 0.4$$

k=1 gives smallest

$$\Rightarrow (1, 2) \in R_1, \text{ partition}$$

Pair (2, 1)

$$\underline{k=1} \rightarrow \frac{1}{3} (f_2(2, 1) + f_3(2, 1)) = \frac{1}{3} (0.35)$$

$$\underline{k=2} \rightarrow \frac{2}{3} (f_1(2, 1) + f_3(2, 1)) = \frac{2}{3} (0.45)$$

$$\underline{k=3} \rightarrow \frac{3}{3} (f_1(2, 1) + f_2(2, 1)) = 0.3$$

k=1 gives smallest

$$\Rightarrow (2, 1) \in R_1, \text{ partition}$$

Pair (2, 2)

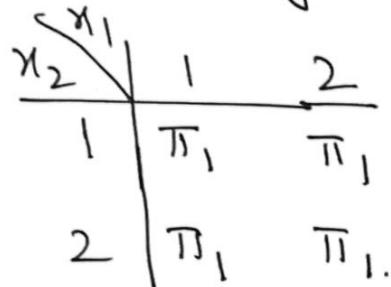
$$\underline{k=1} \rightarrow \frac{1}{3} (f_2(2, 2) + f_3(2, 2)) = \frac{1}{3} (0.65)$$

$$\underline{k=2} \rightarrow \frac{2}{3} (f_1(2, 2) + f_3(2, 2)) = \frac{2}{3} (0.45)$$

$$\underline{k=3} \rightarrow \frac{3}{3} (f_1(2, 2) + f_2(2, 2)) = 0.6$$

$$\Rightarrow (2, 2) \in R_1$$

\Rightarrow ECM minimizing classification rule is



Suppose we have a set \mathcal{D} of preclassified examples

$$\mathcal{D} = \left\{ ((1,2), \pi_1), ((1,1), \pi_1), ((1,1), \pi_1), ((1,1), \pi_3), ((2,2), \pi_2), ((2,1), \pi_3), ((1,2), \pi_3), ((2,2), \pi_3) \right\}$$

If we apply ECM rule on \mathcal{D} , we get

	Coming from		
	π_1	π_2	π_3
π_1	3	1	4
π_2	0	0	0
π_3	0	0	0

classified into

Misclassification rate would be $\frac{5}{8}$!!

↑
not good at all!

(i) Try to find the TPM opt rule

(ii) Try to calculate ECM/TPM of the two optimum rule

Example 2: Multiclass Gaussian

$$\pi_i \equiv N_p(\mu_i, \Sigma_i)$$

Assume for simplicity that $C(k|i)$'s are all same

We know that under this equal cost, the rule is

Assign \tilde{x} to π_k if

$$p_k f(\tilde{x} | \pi_k) = \max_i p_i f(\tilde{x} | \pi_i)$$

$$\text{i.e. if } \log(p_k f(\tilde{x} | \pi_k)) = \max_i \log(p_i f(\tilde{x} | \pi_i))$$

$$\text{with } f(\tilde{x} | \pi_k) = (2\pi)^{-p/2} |\Sigma_k|^{-1/2} \exp\left(-\frac{1}{2}(\tilde{x} - \mu_k)' \Sigma_k^{-1} (\tilde{x} - \mu_k)\right)$$

Define quadratic discriminant score for the i th popn, π_i , as

$$S_i^{QDS}(\tilde{x}) = -\frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\tilde{x} - \mu_i)' \Sigma_i^{-1} (\tilde{x} - \mu_i) + \log p_i$$

Classification Rule

Assign \tilde{x} to π_k if

$$S_K^{QDS}(\tilde{x}) = \max_i S_i^{QDS}(\tilde{x})$$

124

Note: With a given learning sample L of preclassified examples, we calculate the quadratic discriminant score for each class.

$$\hat{g}_i^{QDS}(x) = -\frac{1}{2} \log |\mathcal{S}_i| - \frac{1}{2} (x - \bar{x}_i)^T S_i^{-1} (x - \bar{x}_i) + \log p_i$$

S_i : Sample variance-covariance matrix based on preclassified examples corresponding to class Π_i in \mathcal{D}

\bar{x}_i : sample mean vector based on
 π_i in d.

Remark: Sp case: suppose π_i is $N_p(\mu_i, \Sigma)$ $i=1(1)3$

$$\log p_i f_i(\underline{x}) = -\frac{1}{2} |\Sigma| - \frac{1}{2} \underline{x}' \Sigma^{-1} \underline{x} + \underline{\mu}_i' \Sigma^{-1} \underline{x} - \frac{1}{2} \underline{\mu}_i' \Sigma^{-1} \underline{\mu}_i + \log p_i$$

Same & pop's

Define the linear dimension. Same & pop's

$$f_i^{LDS}(x) = \hat{m}_i' \hat{\Sigma}^{-1} x - \frac{1}{2} \hat{m}_i' \hat{\Sigma}^{-1} \hat{m}_i + \log p_i$$

Assignment rule is:

Assign χ to π_k if

$$f_K^{LDS}(\tilde{x}) = \max_i f_i^{LDS}(\tilde{x})$$

We have to use λ to get the following estimated LDS

$$\hat{f}_i^{\text{LDS}}(\tilde{x}) = \tilde{x}' \hat{\delta}^{-1} \tilde{x} - \frac{1}{2} \tilde{x}' \hat{\delta}^{-1} \tilde{x} + \log p_i$$

\tilde{x}_i for $i=1 \dots c$ calculated from pre-classified examples of π_i class in λ .

$$\lambda \left(\sum_{i=1}^c n_i - c \right) \delta = \sum_{i=1}^c (n_i - 1) \delta_i$$

δ_i for $i=1 \dots c$ calculated from pre-classified examples of π_i class in λ .
 ↑ (divisor $(n_i - 1)$ sample covariance matrix)

Estimated classification rule is :

Assign \tilde{x} to π_K if

$$\hat{f}_K^{\text{LDS}}(\tilde{x}) = \max_i \hat{f}_i^{\text{LDS}}(\tilde{x})$$

Note: The above rule is equivalent to

Assign \tilde{x} to π_K if

$$-\frac{1}{2} (\tilde{x} - \tilde{x}_K)' \hat{\delta}^{-1} (\tilde{x} - \tilde{x}_K) + \log p_K$$

i.e assign \tilde{x} to p_K closest in terms of λ 's Mahalanobis distance penalized by $\log p_K$.

Note: We may also need to estimate p_i 's from λ .