

1) Règle de décision Bayésienne

- Dans le cas de classes équiprobables, la règle de décision Bayésienne s'écrit

$$d^*(z) = w_1 \Leftrightarrow P(w_1|z) = \frac{f(z|w_1)P(w_1)}{f(z)} \geq P(w_2|z)$$

$$d^*(z) = w_1 \Leftrightarrow f(z|w_1) \geq f(z|w_2)$$

Les densités de z conditionnellement à w_1 et w_2 s'écrivent

$$f(z|w_1) = \frac{1}{2\pi\sigma^2} \exp\left[-\frac{1}{2\sigma^2} \|z - \mu_1\|^2\right]$$

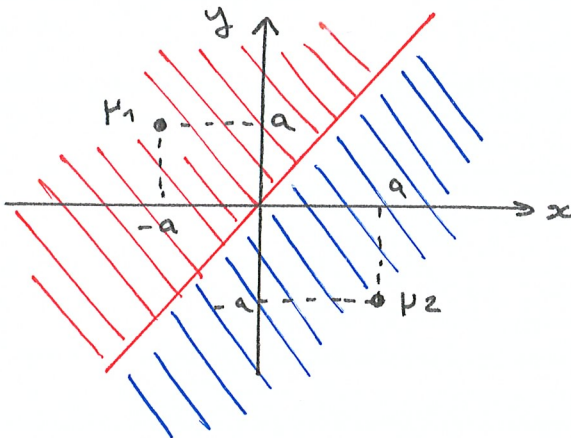
$$f(z|w_2) = \frac{1}{2\pi\sigma^2} \exp\left[-\frac{1}{2\sigma^2} \|z - \mu_2\|^2\right]$$

d'où

$$d^*(z) = w_1 \Leftrightarrow \|z - \mu_1\| \leq \|z - \mu_2\|$$

qui est la règle classique de la distance aux barycentres. Dans le cas où

$\mu_1 = \begin{pmatrix} -a \\ a \end{pmatrix}$ et $\mu_2 = \begin{pmatrix} a \\ -a \end{pmatrix}$ avec $a > 0$, on obtient les régions de décision suivantes :



$$\text{Red hatched region: } d^*(z) = w_1$$

$$\text{Blue hatched region: } d^*(z) = w_2$$

- La probabilité d'erreur associée à la règle d^* ci-dessus est définie par

$$P_e = P_1 + P_2 = \int_{R_2} f(z|w_1) P(w_1) dz + \int_{R_1} f(z|w_2) P(w_2) dz$$

où R_1 et R_2 sont les régions du plan associées aux décisions $d^*(z) = w_1$ et $d^*(z) = w_2$.

Les probabilités P_1 et P_2 se calculent comme suit

$$P_1 = \frac{1}{2} \int_{\mathbb{R}^2} \left[\int_{-\infty}^x \frac{1}{2\pi\sigma^2} \exp\left[-\frac{1}{2\sigma^2} (x+a)^2 - \frac{1}{2\sigma^2} (y-a)^2\right] dy \right] dx$$

$$= \frac{1}{2} \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left[-\frac{1}{2\sigma^2} (x+a)^2\right] \left[\int_{-\infty}^x \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left[-\frac{1}{2\sigma^2} (y-a)^2\right] dy \right] dx$$

d'où $P_1 = \int_{\mathbb{R}} \frac{1}{2} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(x+a)^2\right] \underbrace{\left[\int_{-\infty}^{\frac{x-a}{\sigma}} \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du \right]}_{\Phi\left(\frac{x-a}{\sigma}\right)} dx$

↑
chgt de variables
 $u = \frac{y-a}{\sigma}$

En faisant le changement de variables $u = \frac{x+a}{\sigma}$, on obtient alors

$$P_1 = \frac{1}{2} \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) \Phi\left(u - \frac{2a}{\sigma}\right) du$$

Par symétrie, le calcul de P_2 conduit au même résultat, d'où

$$P_e = \int_{\mathbb{R}} f(u) \Phi\left(u - \frac{2a}{\sigma}\right) du \triangleq P_e(a, \sigma)$$

Quand $a \rightarrow +\infty$ ou $\sigma \rightarrow 0$ on a $-\frac{2a}{\sigma} \rightarrow -\infty$ d'où $P_e \rightarrow 0$

• Lorsque les coûts C_{ij} sont définis par $C_{11} = C_{22} = 0$ et $C_{12} = 2C_{21} = 2$ la règle de décision Bayésienne s'écrit

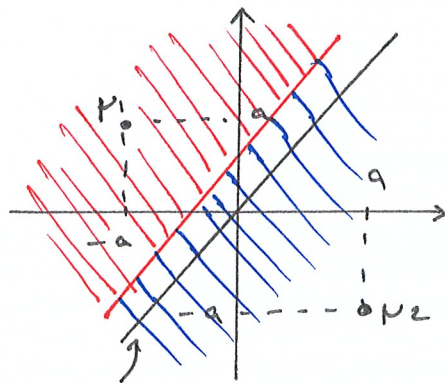
$$d^*(x) = w_1 \Leftrightarrow \frac{f(x|w_1)}{f(x|w_2)} \geq \frac{C_{12} - C_{22}}{C_{21} - C_{11}} \frac{P(w_2)}{P(w_1)} = 2$$

$$d^*(x) = w_1 \Leftrightarrow f(x|w_1) \geq 2 f(x|w_2)$$

C_{12} est le coût de prendre la décision a_1 alors que x est élément de w_2 .
Comme ce coût est plus élevé que C_{21} , on va prendre la décision a_1 moins souvent que dans le cas précédent où $d^*(x) = w_1 \Leftrightarrow f(x|w_1) \geq f(x|w_2)$
On a alors après quelques calculs élémentaires

$$d^*(x) = w_1 \Leftrightarrow \|x - \mu_1\|^2 \leq \|x - \mu_2\|^2 - 2\sigma^2 \ln 2$$

Par rapport à la question précédente, on a déplacé la droite de séparation entre les deux régions, comme illustré sur la figure ci-dessous



▨ $d^*(z) = w_1$
▨ $d^*(z) = w_2$

frontière pour la question précédente

2) Apprentissage Bayésien

On désire estimer μ_1 et σ_1^2 à l'aide des vecteurs v_1, \dots, v_{n_1} et de manière similaire estimer μ_2 et σ_2^2 à l'aide des vecteurs w_1, \dots, w_{n_2} .

La vraisemblance associée aux données v_1, \dots, v_{n_1} s'écrit en supposant l'indépendance entre ces différents vecteurs:

$$\begin{aligned}
 f(v; \mu_1, \sigma_1^2) &= f(v_1, \dots, v_{n_1}; \mu_1, \sigma_1^2) \\
 &= \prod_{i=1}^{n_1} \frac{1}{(2\pi\sigma_1^2)^{1/2}} \exp\left[-\frac{1}{2\sigma_1^2} \|z_i - \mu_1\|^2\right] \\
 &= \frac{1}{(2\pi\sigma_1^2)^{n_1}} \exp\left[-\frac{1}{2\sigma_1^2} \sum_{i=1}^{n_1} \|z_i - \mu_1\|^2\right]
 \end{aligned}$$

on en déduit la log-vraisemblance

$$\ln f(v; \mu_1; \sigma_1^2) = -n_1 \ln(2\pi) - \frac{n_1}{2} \ln(\sigma_1^2) - \frac{1}{2\sigma_1^2} \sum_{i=1}^{n_1} \|z_i - \mu_1\|^2$$

que l'on peut dériver par rapport à μ_1 et σ_1^2 pour obtenir

$$\frac{\partial \ln f}{\partial \mu_1} = 0 \Rightarrow \boxed{\hat{\mu}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} z_i}$$

$$\frac{\partial \ln f}{\partial \sigma_1^2} = 0 \Rightarrow -\frac{n_1}{\sigma_1^2} + \frac{1}{2\sigma_1^4} \sum_{i=1}^{n_1} \|z_i - \mu_1\|^2 = 0 \Rightarrow \boxed{\hat{\sigma}_1^2 = \frac{1}{2n_1} \sum_{i=1}^{n_1} \|z_i - \mu_1\|^2}$$

La règle de classification qui découle de cette procédure d'apprentissage s'écrit

$$d^*(z) = w_1 \Leftrightarrow \hat{f}(z|w_1) \geq \hat{f}(z|w_2)$$

$$\Leftrightarrow \frac{1}{\hat{\sigma}_1^2} \exp\left[-\frac{1}{2\hat{\sigma}_1^2} \|z - \hat{\mu}_1\|^2\right] \geq \frac{1}{\hat{\sigma}_2^2} \exp\left[-\frac{1}{2\hat{\sigma}_2^2} \|z - \hat{\mu}_2\|^2\right]$$

$$\Leftrightarrow \ln\left(\frac{\hat{\sigma}_2^2}{\hat{\sigma}_1^2}\right) - \frac{1}{2\hat{\sigma}_1^2} \|z - \hat{\mu}_1\|^2 \geq -\frac{1}{2\hat{\sigma}_2^2} \|z - \hat{\mu}_2\|^2$$

c'est à dire

(4)

$$d^*(z) = w_1 \Leftrightarrow \|z - \hat{\mu}_1\|^2 \leq \frac{\hat{\sigma}_1^2}{\hat{\sigma}_2^2} \|z - \hat{\mu}_2\|^2 + 2 \frac{\hat{\sigma}_1^2}{\hat{\sigma}_2^2} \ln \left(\frac{\hat{\sigma}_2^2}{\hat{\sigma}_1^2} \right)$$

- Dans le cas de données multidimensionnelles, on peut estimer une densité à l'aide de la méthode des fenêtres de Parzen ou la méthode des k plus proches voisins. Considérons par exemple la méthode des fenêtres de Parzen. Un estimateur de densité s'écrit pour $x \in \mathbb{R}^p$

$$\hat{f}(x) = \frac{1}{nh^p} \sum_{i=1}^n \phi\left(\frac{x - z_i}{h}\right)$$

où z_i est un vecteur de la base d'apprentissage de densité f , ϕ est un noyau choisi a priori (par exemple noyau Gaussien) et h est un paramètre (appelé parfois bande passante du noyau) qu'il faut régler de façon à avoir une estimation peu bruitée avec une résolution suffisante.

On parle d'estimation de densité non paramétrique car ces méthodes ne supposent pas que la loi est définie par un vecteur paramètre θ que l'on cherche à estimer comme pour la méthode d'apprentissage basé sur l'estimateur du maximum de vraisemblance.

3) Perception Multi-couches

- Cherchons tout d'abord le vecteur $w_1 = (w_{11}, w_{21}, a)^T$ qui minimise $E[e_1^2(n)] = E[(d_1(n) - y_1(n))^2]$

Puisque $f(t) = t$, on a $y_1(n) = w_{11} x_1(n) + w_{21} x_2(n) - a$ et donc

$$\begin{cases} \frac{\partial E[e_1^2(n)]}{\partial w_{11}} = 0 \\ \frac{\partial E[e_1^2(n)]}{\partial w_{21}} = 0 \\ \frac{\partial E[e_1^2(n)]}{\partial a} = 0 \end{cases} \Rightarrow \begin{cases} E\left[e_1(n) \frac{\partial y_1(n)}{\partial w_{11}} \right] = 0 \\ E\left[e_1(n) \frac{\partial y_1(n)}{\partial w_{21}} \right] = 0 \\ E\left[e_1(n) \frac{\partial y_1(n)}{\partial a} \right] = 0 \end{cases}$$

c'est-à-dire

$$\begin{cases} E[e_1(n)x_1(n)] = 0 \\ E[e_1(n)x_2(n)] = 0 \\ E[e_1(n)] = 0 \end{cases} \Leftrightarrow \begin{cases} E[d_1(n)x_1(n)] = w_{11} E[x_1^2(n)] + w_{21} E[x_1(n)x_2(n)] - a E[x_1(n)] \\ E[d_1(n)x_2(n)] = w_{11} E[x_1(n)x_2(n)] + w_{21} E[x_2^2(n)] - a E[x_2(n)] \\ E[d_1(n)] = w_{11} E[x_1(n)] + w_{21} E[x_2(n)] - a \end{cases}$$

On a donc un système de 3 équations à 3 inconnues qui nous permet de trouver w_{11} , w_{21} et a . Ceci suppose bien sûr qu'on connaisse toutes les espérances intervenant dans ce système

On fait de même pour estimer w_2 en minimisant $E[e_2^2(n)]$ et pour estimer w_3 en minimisant $E[e_3^2(n)]$.

- On sait que la dérivée de la fonction sigmoïde vérifie

$$f'(x) = \alpha f(x) [1 - f(x)]$$

La mise à jour des poids w_{11} , w_{21} et a à l'aide de la règle du gradient s'écrit sous la forme

$$w_{11}(n+1) = w_{11}(n) - \mu \left. \frac{\partial e_1^2(n)}{\partial w_{11}} \right|_{w_{11}=w_{11}(n)}$$

$$w_{21}(n+1) = w_{21}(n) - \mu \left. \frac{\partial e_1^2(n)}{\partial w_{21}} \right|_{w_{21}=w_{21}(n)}$$

$$a(n+1) = a(n) - \mu \left. \frac{\partial e_1^2(n)}{\partial a} \right|_{a=a(n)}$$

c'est-à-dire

$$\begin{cases} w_{11}(n+1) = w_{11}(n) - 2\mu e_1(n) \frac{\partial e_1(n)}{\partial w_{11}} = w_{11}(n) + 2\mu e_1(n) \frac{\partial y_1(n)}{\partial w_{11}} \\ w_{21}(n+1) = w_{21}(n) - 2\mu e_1(n) \frac{\partial e_1(n)}{\partial w_{21}} = w_{21}(n) + 2\mu e_1(n) \frac{\partial y_1(n)}{\partial w_{21}} \\ a(n+1) = a(n) - 2\mu e_1(n) \frac{\partial e_1(n)}{\partial a} = a(n) + 2\mu e_1(n) \frac{\partial y_1(n)}{\partial a} \end{cases}$$

Les dérivées de $y_1(n) = f[w_{11}x_1(n) + w_{21}x_2(n) - a]$ par rapport à w_{11} , w_{21} et a se calculent facilement - Par exemple

$$\frac{\partial y_1(n)}{\partial w_{11}} = \alpha y_1(n) [1 - y_1(n)] x_1(n)$$

On en déduit les règles de mise à jour de w_{11} , w_{21} et a :

$$\begin{aligned}
 w_{11}(n+1) &= w_{11}(n) + \delta e_1(n) y_1(n) [1 - y_1(n)] x_1(n) \\
 w_{21}(n+1) &= w_{21}(n) + \delta e_1(n) y_1(n) [1 - y_1(n)] x_2(n) \\
 a(n+1) &= a(n) + \delta e_1(n) y_1(n) [1 - y_1(n)] (-1)
 \end{aligned}$$

4) Règle DCMS

Le coût de prendre la décision "classe c" prend les valeurs

- s_{1c} avec proba p_{1c}
- s_{2c} avec proba p_{2c}
- \vdots
- s_{cc} avec proba p_{cc}

La moyenne de ce coût s'écrit donc $E[\phi_c] = \sum_{i=1}^c s_{ic} p_{ic}$

Puisque le coût associé à une décision correcte, noté s_{cc} , est nul, on a

$$E[\phi_c] = \sum_{\substack{i=1 \\ i \neq c}}^c s_{ic} p_{ic} \quad (*)$$

Le coût de ne pas prendre la décision "classe c" alors que les données proviennent de la classe c, prend les valeurs

- s_{c1} avec la proba p_{c1}
- s_{c2} avec la proba p_{c2}
- \vdots
- s_{cc} avec la proba p_{cc}

La valeur moyenne est donc

$$E[V_c] = \sum_{\substack{i=1 \\ i \neq c}}^c s_{ci} p_{ci} \quad (**)$$

La question la plus difficile de l'examen car il fallait avoir compris l'article ! On utilise un tiers des données pour déterminer les vecteurs w_1, w_2 et w_3 en appliquant la règle du gradient. Une fois que ces vecteurs ont été déterminés, on estime les probabilités p_{ij} à l'aide de la relation (1) et du second tiers des données de la base d'apprentissage. Ensuite, à l'aide des relations (*) et (**) ci-dessus, on peut déterminer $E[\phi_c]$ et $E[V_c]$.

On en déduit alors la matrice de coûts

$$\begin{bmatrix} 0 & E[V_c] \\ E[\phi_c] & 0 \end{bmatrix}$$

à laquelle on peut appliquer une règle de Bayes standard.

• Distance de Mahalanobis

Pour une chaque classe, on estime les vecteurs moyennes m_i et les matrices de covariance Σ_i à l'aide de la base d'apprentissage. On notera \hat{m}_i et $\hat{\Sigma}_i$ les estimations obtenues. Le classifieur construit à partir de la règle de la distance aux barycentres s'écrit

$$d(\mathbf{z}) = w_j \iff d(\mathbf{z}, \hat{\mu}_j) \leq d(\mathbf{z}, \hat{\mu}_k) \quad \forall k$$

avec $d(\mathbf{z}, \hat{\mu}_j) = (\mathbf{z} - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (\mathbf{z} - \hat{\mu}_j)$.

Classifieur KNN

Pour classifier \mathbf{z} , on cherche ses K plus proches voisins parmi la base d'apprentissage et on affecte \mathbf{z} à la classe majoritairement représentée par ces voisins.