

2.3 Méthode probabiliste

2.3.1 Définitions

Avec la méthode probabiliste, on considère que le vecteur de caractéristique $\mathbf{x} = (x_1, \dots, x_N)$ a une densité de probabilité (ddp) qui dépend de sa classe.

Exemple d'une ddp pour un vecteur \mathbf{x} ?

$$P(\mathbf{x}) = \frac{1}{\left(\sqrt{2\pi}|\mathbf{\Gamma}|\right)^N} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Gamma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

Dans la suite de ce document, on utilise les notations suivantes

$$\int_{x_1} \dots \int_{x_N} g(\mathbf{x})P(\mathbf{x})dx_1\dots dx_N = \int g(\mathbf{x})P(\mathbf{x})d\mathbf{x}$$

On a en particulier

$$\begin{aligned} \int P(\mathbf{x})d\mathbf{x} &= 1 \\ \int \mathbf{x}P(\mathbf{x})d\mathbf{x} &= \boldsymbol{\mu} \\ \int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T P(\mathbf{x})d\mathbf{x} &= \mathbf{\Gamma} \end{aligned}$$

où $\boldsymbol{\mu}$ est le vecteur moyen et $\mathbf{\Gamma}$ est la matrice de covariance de \mathbf{X} .

Rappel. Un discriminateur est une application

$$\begin{aligned} \mathbb{R}^N &\mapsto \{\omega_1, \omega_2, \dots, \omega_C\} \\ \mathbf{x} &\mapsto \hat{d}(\mathbf{x}) \end{aligned}$$

Avec la méthode probabiliste, comme au moment où le discriminateur doit effectuer la discrimination, celui-ci ne connaît pas la classe ω , celle-ci est supposée être une variable aléatoire discrète.

On suppose que $\forall \mathbf{x}$, il existe une unique classe associée à \mathbf{x} . Ainsi

$$P(\omega_i \text{ et } \omega_j | \mathbf{x}) = 0 \quad \text{pour } i \neq j \quad \text{non ambiguïté}$$

$$\sum_{c=1}^C P(\omega_c | \mathbf{x}) = 1 \quad \text{exhaustivité}$$

Il y a 2 propriétés importantes à connaître et savoir manipuler :

1. Loi de Bayes

$$P(\mathbf{x}, \omega_c) = P(\mathbf{x}|\omega_c)P(\omega_c) = P(\omega_c|\mathbf{x})P(\mathbf{x})$$

2. Calcul des marginales (ou théorème des probabilités totales)

$$P(\mathbf{x}) = \sum_{c=1}^C P(\mathbf{x}, \omega_c) = \sum_{c=1}^C P(\mathbf{x}|\omega_c)P(\omega_c)$$

et

$$P(\omega_c) = \int P(\mathbf{x}, \omega_c)d\mathbf{x} = \int P(\omega_c|\mathbf{x})P(\mathbf{x})d\mathbf{x}$$

Quel est le sens des variables ci-dessous ?

- $P(\omega_c)$
- $P(\mathbf{x}|\omega_c)$
- $P(\omega_c|\mathbf{x})$
- $P(\mathbf{x})$

Pour chaque classe, on peut calculer la moyenne et la matrice de covariance intra-classe

$$\begin{aligned} \boldsymbol{\mu}_c &= \int \mathbf{x}P(\mathbf{x}|\omega_c)d\mathbf{x} \\ \mathbf{\Gamma}_c &= \int (\mathbf{x} - \boldsymbol{\mu}_c)(\mathbf{x} - \boldsymbol{\mu}_c)^T P(\mathbf{x}|\omega_c)d\mathbf{x} \end{aligned}$$

à la condition que ces intégrales existent.

2.3.2 Application de la théorie de la décision dans le cas à deux classes

Quel est le critère à utiliser pour caractériser les performances d'un discriminateur ?

$$P_{err} = P(\text{faire une erreur}) = \int_{\mathbb{R}^N} P(\mathbf{x}, \text{faire une erreur}) d\mathbf{x}$$

Dans un problème à deux classes, on a 2 erreurs possibles

réalité	Probabilité	décision
$\mathbf{x} \in \omega_1$	$P(\mathbf{x}, \omega_1)$	$\mathbf{x} \mapsto \omega_2$
$\mathbf{x} \in \omega_2$	$P(\mathbf{x}, \omega_2)$	$\mathbf{x} \mapsto \omega_1$

Ainsi

$$P_{err} = \int_{\mathcal{R}_1} P(\mathbf{x}, \omega_2) d\mathbf{x} + \int_{\mathcal{R}_2} P(\mathbf{x}, \omega_1) d\mathbf{x}$$

où \mathcal{R}_1 et \mathcal{R}_2 correspondent à la partition liée à notre discriminateur.

Exemple. Voir Figure 10.

Frontière qui minimise la probabilité d'erreur quand les classes sont équiprobables $P(\omega_1) = P(\omega_2)$

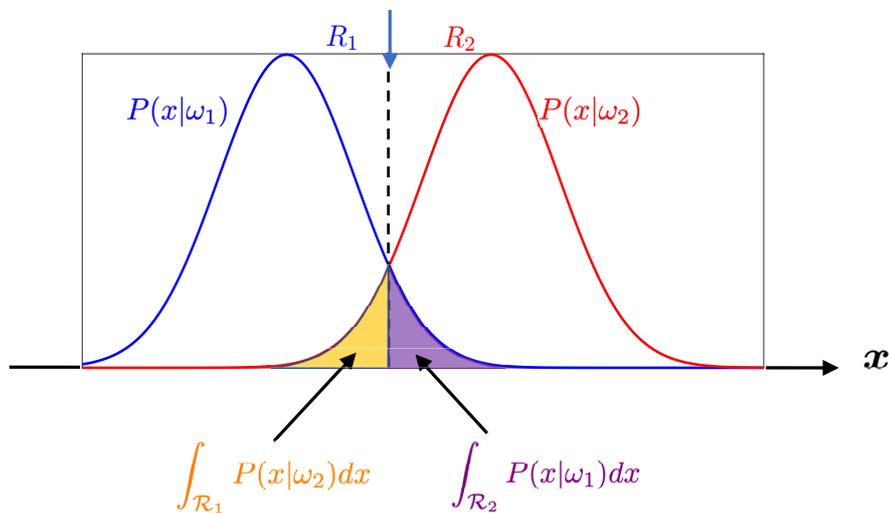


FIGURE 10 – Illustration dans le cas avec deux classes gaussiennes équiprobables.

On peut montrer que la probabilité d'erreur minimale est donnée par

$$P_{err}^{opt} = \int \min(P(\mathbf{x}, \omega_1), P(\mathbf{x}, \omega_2)) d\mathbf{x}$$

Ceci s'illustre en 1D sur la Figure 10 dans le cas équiprobable.

Ainsi la règle de décision qui minimise la probabilité d'erreur P_{err} est

$$P(\mathbf{x}, \omega_1) \underset{\omega_2}{\overset{\omega_1}{\geq}} P(\mathbf{x}, \omega_2)$$

ce qui est équivalent à

$$P(\omega_1|\mathbf{x}) \underset{\omega_2}{\overset{\omega_1}{\geq}} P(\omega_2|\mathbf{x})$$

mais ce qui est aussi équivalent à

$$P(\mathbf{x}|\omega_1)P(\omega_1) \underset{\omega_2}{\overset{\omega_1}{\geq}} P(\mathbf{x}|\omega_2)P(\omega_2)$$

Que signifie le fait que cette règle de décision est optimale quand on choisit comme critère probabilité d'erreur ?

Que faire si on connaît $P(\mathbf{x}|\omega_i)$, mais on ne connaît pas $P(\omega_i)$?

2.3.3 Cas gaussien avec 2 classes équiprobables et paramètres connus

Si on suppose

$$P(\mathbf{x}|\omega_c) = \frac{1}{\sqrt{2\pi}^N \sqrt{|\mathbf{\Gamma}_c|}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \mathbf{\Gamma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \right]$$

Alors le discriminateur qui minimise la probabilité d'erreur

$$\ln \left(\frac{P(\mathbf{x}|\omega_1)}{P(\mathbf{x}|\omega_2)} \right) \underset{\omega_2}{\overset{\omega_1}{\geq}} 0$$

peut s'écrire

$$\frac{1}{2} \ln \frac{|\mathbf{\Gamma}_2|}{|\mathbf{\Gamma}_1|} + \frac{1}{2} (Q_2(\mathbf{x}) - Q_1(\mathbf{x})) \underset{\omega_2}{\overset{\omega_1}{\geq}} 0$$

où

$$Q_c(\mathbf{x}) = (\mathbf{x} - \boldsymbol{\mu}_c)^T \mathbf{\Gamma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c)$$

Que représente $Q_c(\mathbf{x})$?

La frontière qui sépare les deux régions est une fonction quadratique. Elle peut décrire plusieurs formes :

— 1 ellipse avec une équation du type

$$\frac{x_1^2}{a_1^2} + \frac{x_2^2}{a_2^2} = 1$$

— 1 hyperbole avec une équation du type

$$\frac{x_1^2}{a_1^2} - \frac{x_2^2}{a_2^2} = 1$$

— 1 parabole avec une équation du type

$$x_2^2 = ax_1$$

Dans le cas particulier où $\mathbf{\Gamma}_1 = \mathbf{\Gamma}_2 = \mathbf{\Gamma}$, le test se simplifie en

$$Q_2(\mathbf{x}) \underset{\omega_2}{\overset{\omega_1}{\geq}} Q_1(\mathbf{x})$$

ainsi

$$(\mathbf{x} - \boldsymbol{\mu}_2)^T \mathbf{\Gamma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \underset{\omega_2}{\overset{\omega_1}{\geq}} (\mathbf{x} - \boldsymbol{\mu}_1)^T \mathbf{\Gamma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)$$

En développant, on obtient

$$-\boldsymbol{\mu}_2^T \mathbf{\Gamma}^{-1} \mathbf{x} - \mathbf{x}^T \mathbf{\Gamma}^{-1} \boldsymbol{\mu}_2 + \boldsymbol{\mu}_2^T \mathbf{\Gamma}^{-1} \boldsymbol{\mu}_2 \underset{\omega_2}{\overset{\omega_1}{\geq}} -\boldsymbol{\mu}_1^T \mathbf{\Gamma}^{-1} \mathbf{x} - \mathbf{x}^T \mathbf{\Gamma}^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\mu}_1^T \mathbf{\Gamma}^{-1} \boldsymbol{\mu}_1$$

et finalement

$$\boldsymbol{\mu}_2^T \mathbf{\Gamma}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^T \mathbf{\Gamma}^{-1} \boldsymbol{\mu}_1 \underset{\omega_2}{\overset{\omega_1}{\geq}} 2 (\boldsymbol{\mu}_2^T - \boldsymbol{\mu}_1^T) \mathbf{\Gamma}^{-1} \mathbf{x}$$

De quel type de frontière s'agit-il ?

Quels sont les risques de cette approche ?

Remarque. Il est possible de chercher une frontière **linéaire**, sans avoir besoin de supposer que les classes sont gaussiennes. Pour cela, on peut repartir de la règle de décision

$$\ln \left(\frac{P(\omega_1|\mathbf{x})}{P(\omega_2|\mathbf{x})} \right) \underset{\omega_2}{\overset{\omega_1}{\geq}} 0$$

Si on suppose qu'il existe une frontière linéaire sur la statistique $\ln \left(\frac{P(\omega_1|\mathbf{x})}{P(\omega_2|\mathbf{x})} \right)$, alors il existe w_0 et \mathbf{w} tels que

$$\ln \left(\frac{P(\omega_1|\mathbf{x})}{P(\omega_2|\mathbf{x})} \right) = w_0 + \mathbf{w}^T \mathbf{x}$$

Or, comme

$$P(\omega_1|\mathbf{x}) + P(\omega_2|\mathbf{x}) = 1$$

on a donc

$$\ln\left(\frac{P(\omega_1|\mathbf{x})}{1 - P(\omega_1|\mathbf{x})}\right) = w_0 + \mathbf{w}^T \mathbf{x}$$

et donc

$$P(\omega_1|\mathbf{x}) = \frac{e^{w_0 + \mathbf{w}^T \mathbf{x}}}{1 + e^{w_0 + \mathbf{w}^T \mathbf{x}}}$$

et

$$P(\omega_2|\mathbf{x}) = 1 - P(\omega_1|\mathbf{x}) = \frac{1}{1 + e^{w_0 + \mathbf{w}^T \mathbf{x}}}$$

C'est le discriminateur *logistique*.

Un intérêt de cette solution est que pour identifier \mathbf{w} , on peut utiliser comme critère

$$J(\mathbf{w}) = \sum_{\{p|d^{(p)}=1\}} \log \frac{e^{\mathbf{w}^T \mathbf{x}^{(p)}}}{1 + e^{\mathbf{w}^T \mathbf{x}^{(p)}}} + \sum_{\{p|d^{(p)}=-1\}} \log \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x}^{(p)}}}$$

qui revient à maximiser une sorte de vraisemblance sur la base d'apprentissage.

On peut montrer que cette optimisation est bien plus facile à réaliser que celle du gradient par exemple.

Synthèse sur le linéaire. On a plusieurs solutions pour faire du linéaire.

- Hebb, pas de critère mathématique, il s'agit d'une solution plutôt empirique.
- pseudo-inverse (PI), le critère est

$$J = \sum_p \|d^{(p)} - \mathbf{w}^T \mathbf{x}^{(p)}\|^2$$

- RA et gradient utilisent le même critère, mais ils régularisent l'inversion de matrice.
- le cas gaussien avec $\mathbf{\Gamma}_1 = \mathbf{\Gamma}_2$,
- le logistique avec le critère ci-dessus.

On a donc 6 discriminateurs qui mènent à une frontière linéaire.

Est-ce que les solutions donnent des résultats identiques ?

2.3.4 Cas gaussien avec 2 classes équiprobables et paramètres inconnus

Si en pratique, on ne connaît pas les paramètres associés à chacune des lois $P(\mathbf{x}|\omega_c)$, il faut mettre en œuvre des estimateurs.

Quels estimateurs vous connaissez ?

Dans le cas gaussien où l'on ne suppose pas que les covariances sont égales, si on a N_c vecteur de la classe c dans la base d'apprentissage, alors les estimateurs s'écrivent $\forall c$

$$\hat{\boldsymbol{\mu}}_c = \frac{1}{N_c} \sum_{p=1}^{N_c} \mathbf{x}^{(p)}$$

et

$$\hat{\boldsymbol{\Gamma}}_c = \frac{1}{N_c} \sum_{p=1}^{N_c} (\mathbf{x}^{(p)} - \hat{\boldsymbol{\mu}}_c) (\mathbf{x}^{(p)} - \hat{\boldsymbol{\mu}}_c)^T$$

Si on sait que les covariances sont égales, alors on peut montrer que l'estimateur de la matrice de covariance $\hat{\boldsymbol{\Gamma}}$ s'écrit alors

$$\hat{\boldsymbol{\Gamma}} = \frac{1}{2} (\hat{\boldsymbol{\Gamma}}_1 + \hat{\boldsymbol{\Gamma}}_2)$$

Qu'est-ce qu'on perd quand les paramètres sont inconnus ?

Qu'est-ce qu'il convient de faire dans ce cas ?

2.3.5 Généralisation dans le cas à C classes

Dans le cas à C classes, il faut sommer toutes les erreurs possibles

$$P_{err} = \sum_{i=1}^C \int_{\mathcal{R}_i} \sum_{j \text{ avec } j \neq i} P(\mathbf{x}, \omega_j) d\mathbf{x}$$

ce qui peut s'écrire

$$P_{err} = \sum_{i=1}^C \int_{\mathcal{R}_i} \sum_j P(\mathbf{x}, \omega_j) d\mathbf{x} - \sum_{i=1}^C \int_{\mathcal{R}_i} P(\mathbf{x}, \omega_i) d\mathbf{x}$$

Or $\sum_j P(\mathbf{x}, \omega_j) d\mathbf{x} = P(\mathbf{x})$ et ainsi

$$P_{err} = \int_{\mathbb{R}} P(\mathbf{x}) d\mathbf{x} - \sum_{i=1}^C \int_{\mathcal{R}_i} P(\mathbf{x}, \omega_i) d\mathbf{x} = 1 - \sum_{i=1}^C \int_{\mathcal{R}_i} P(\mathbf{x}, \omega_i) d\mathbf{x}$$

Pour minimiser la probabilité d'erreur, il faut donc choisir comme partitionnement \mathcal{R}_i

$$\mathcal{R}_i = \{\mathbf{x} | P(\mathbf{x}, \omega_i) \geq P(\mathbf{x}, \omega_j) \quad \forall j\}$$

Si toutes les classes sont équiprobables, alors quelle la règle qui minimise la probabilité d'erreur ?