

2 Illustration en science des données : la discrimination

Introduction

En science des données, à partir d'une observation \mathbf{x}^1 bruitée par un bruit noté \mathbf{b} , il peut y avoir plusieurs tâches possibles à réaliser :

- débruitage : $\mathbf{x} = \mathbf{s} + \mathbf{b}$. Dans ce cas, on cherche une fonction f telle que $f(\mathbf{x}) \simeq \mathbf{s}$.
- estimation : $\mathbf{x} \sim P_{\theta}$. Dans ce cas, on cherche une fonction f telle que $f(\mathbf{x}) \simeq \theta$.
- détection : $H_0 : \mathbf{x} = \mathbf{b}$ ou $H_1 : \mathbf{x} = \mathbf{s} + \mathbf{b}$. Dans ce cas, on cherche une fonction f et un seuil ρ tels que $f(\mathbf{x}) \underset{H_1}{\overset{H_0}{\geq}} \rho$ permette de discriminer correctement entre l'hypothèse H_0 et H_1 .

Ce cours porte sur la discrimination

Cela signifie qu'à partir de la connaissance de \mathbf{x} , il faut retrouver à quelle classe le vecteur \mathbf{x} appartient.

On peut donc représenter un discriminateur comme une application

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

où pour simplifier on a supposé que $\mathbf{x} \in \mathbb{R}^N$, où les ω_i sont les différentes classes possibles et $\hat{d}(\mathbf{x})$ est la classe estimée. Par exemple, on trouvera des situations pour lesquelles on a

$$\hat{d}(\mathbf{x}) = \arg \max_{c=1,2,\dots,C} f(\mathbf{x}; c)$$

où f est une fonction qui reste à déterminer et qui dépend de \mathbf{x} et de c .

Remarque. En pratique, le vecteur \mathbf{x} peut correspondre à toute sorte de données : image, spectre, son audio, ADN, cours de la bourse, vidéo, ...

Exemple. On considère un problème où une photo représente soit un chien, soit un chat. On souhaite faire un système capable de distinguer à partir d'une photo s'il s'agit d'un chien ou d'un chat.

Que signifie résoudre ce problème de discrimination par apprentissage ?

Cela signifie que pour trouver l'application \hat{d} , on utilise une base de données pour laquelle le problème a déjà été résolu. C'est à dire qu'on dispose d'un ensemble

$$\mathcal{B} = \{(\mathbf{x}^{(p)}, d^{(p)}); p \in \{1, \dots, P\}\}$$

où $d^{(p)}$ représente la classe de $\mathbf{x}^{(p)}$ (on parle aussi souvent d'étiquette) et P représente la taille de la base de données.

On appelle \mathcal{B} la base d'apprentissage car celle-ci est utilisée pour **apprendre** la fonction \hat{d} .

Est-il possible d'effectuer de la discrimination en aveugle (i.e. sans base d'apprentissage) ?

Dans ce cours, on s'intéressera uniquement à des algorithmes de discrimination avec base d'apprentissage (on dit aussi classification supervisée).

Quels sont les objectifs de ce cours ?

Donner aux élèves les éléments méthodologiques leur permettant d'analyser les comportements de discriminateurs, et en particulier d'analyser leurs performances.

Pourquoi se focaliser sur l'étude des performances ?

Parce qu'en tant qu'ingénieur généraliste, les élèves seront utilisateur de ces techniques, et pas forcément développeur. Sans être des spécialistes, ils doivent savoir en analyser les performances.

1. Dans la partie 2 du cours, on choisit la convention suivante : les vecteurs sont notés en gras (exemple \mathbf{x}), les matrices sont notées en gras majuscules (exemple \mathbf{X}).

2.1 Discrimination linéaire

2.1.1 Définitions

Soit $\mathbf{x} \in \mathcal{S} \subset \mathbb{R}^N$, un vecteur de caractéristiques.

Trouver une technique de classification revient à partitionner l'espace \mathcal{S} . On considère des partitionnements avec les propriétés ci-dessous :

$$\begin{aligned} \cup_n \mathcal{R}_n &= \mathcal{S} && \text{exhaustivité} \\ \mathcal{R}_n \cap \mathcal{R}_{n'} &= \emptyset && \text{non ambiguïté} \end{aligned}$$

et on définit la règle de décision suivante :

$$\text{Si } \mathbf{x} \in \mathcal{R}_n, \text{ alors on attribue la classe } \omega_n$$

Dans le reste de ce chapitre, on considère le cas à 2 classes :

$$\begin{aligned} \mathcal{R}_1 \cup \mathcal{R}_2 &= \mathcal{S} && \text{exhaustivité} \\ \mathcal{R}_1 \cap \mathcal{R}_2 &= \emptyset && \text{non ambiguïté} \end{aligned}$$

Définition. Un discriminateur linéaire signifie que l'espace \mathcal{S} est partitionné avec des frontières décrites par des hyper-plans.

Remarque. La discrimination linéaire est souvent le premier outil à utiliser car il est simple à mettre en œuvre et est très facile à interpréter. On commence donc par celui-ci. Ensuite, on verra des méthodes plus compliquées comme l'approche probabiliste ou les réseaux de neurones.

Mathématiquement la règle de décision associée à une frontière linéaire peut être décrite par

$$x_1 w_1 + \dots + x_N w_N \underset{\omega_2}{\overset{\omega_1}{\geq}} b$$

où $\mathbf{x} = (x_1, \dots, x_N)^T$ est le vecteur de caractéristiques et

$$x_1 w_1 + \dots + x_N w_N = b$$

est l'équation de l'hyper-plan qui joue le rôle de frontière. Le vecteur $\mathbf{w} = (w_1, \dots, w_N)^T$ est orthogonal à cet hyperplan.

Plutôt que faire apparaître le seuil b , en général on simplifie les notations avec les transformations :

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \dots \\ x_N \end{pmatrix} \mapsto \mathbf{x}' = \begin{pmatrix} 1 \\ x_1 \\ \dots \\ x_N \end{pmatrix}$$

et

$$\mathbf{w} = \begin{pmatrix} w_1 \\ \dots \\ w_N \end{pmatrix} \mapsto \mathbf{w}' = \begin{pmatrix} w_0 \\ w_1 \\ \dots \\ w_N \end{pmatrix}$$

où $b = -w_0$. Ainsi, l'équation de la frontière

$$x_1 w_1 + \dots + x_n w_n = b$$

qui s'écrit

$$\mathbf{w}^T \mathbf{x} = b$$

devient simplement

$$(\mathbf{w}')^T \mathbf{x}' = 0$$

Dans la suite, on ne notera plus les '.

Synthèse : à partir d'une base d'apprentissage

$$\mathcal{B}_{app} = \{(\mathbf{x}^{(p)}, d^{(p)}); p \in \{1, \dots, P_{app}\}\}$$

où P_{app} est la taille de la base d'apprentissage, un discriminateur linéaire est simplement défini par un vecteur \mathbf{w} qui est une fonction de \mathcal{B}_{app} et la règle de décision s'écrit simplement

$$\mathbf{w}^T \mathbf{x} \underset{\omega_2}{\overset{\omega_1}{\geq}} 0 \quad \forall \mathbf{x}$$

Dans la suite de ce chapitre, on considère le cas à 2 classes en choisissant la convention suivante :

Si $\mathbf{x}^{(p)}$ appartient à la classe ω_1 , alors $d^{(p)} = 1$
 Si $\mathbf{x}^{(p)}$ appartient à la classe ω_2 , alors $d^{(p)} = -1$

2.1.2 Méthode de Hebb

Cette méthode établie par M. Hebb en 1949 est très simple

$$\mathbf{w}_{\text{Hebb}} = \sum_{p=1}^{P_{app}} d^{(p)} \mathbf{x}^{(p)}$$

où la somme est effectuée sur tous les vecteurs de la base d'apprentissage. Ainsi

$$\forall \mathbf{x} \quad \hat{d}_{\text{Hebb}}(\mathbf{x}) = \text{sign}(\mathbf{w}_{\text{Hebb}}^T \mathbf{x})$$

Et ainsi

$$\mathbf{w}_{\text{Hebb}}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}_{\text{Hebb}} = \sum_{p=1}^{P_{app}} d^{(p)} \mathbf{x}^T \mathbf{x}^{(p)}$$

Comment interpréter cette relation ?

2.1.3 Méthode de la pseudo-inverse

Plutôt qu'une méthode empirique comme celle de Hebb, on analyse à présent une méthode "mathématique" pour trouver un discriminateur. Avoir une méthode mathématique signifie qu'elle est fondé sur la minimisation d'un critère.

Avec la pseudo-inverse, on utilise comme critère

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

Pourquoi choisir ce critère ?

Si on concatène les vecteurs de la base d'apprentissage pour créer une matrice \mathbf{X} de taille $N \times P_{app}$

$$\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(P_{app})})$$

et que l'on définit le vecteur $\mathbf{d} = (d^{(1)}, \dots, d^{(P_{app})})^T$ alors on peut montrer que

$$J = (\mathbf{d} - \mathbf{X}^T \mathbf{w})^T (\mathbf{d} - \mathbf{X}^T \mathbf{w})$$

En effet, on a

$$\mathbf{d} - \mathbf{X}^T \mathbf{w} = \begin{pmatrix} d^{(1)} - (\mathbf{x}^{(1)})^T \mathbf{w} \\ d^{(2)} - (\mathbf{x}^{(2)})^T \mathbf{w} \\ \dots \\ d^{(P_{app})} - (\mathbf{x}^{(P_{app})})^T \mathbf{w} \end{pmatrix} = \begin{pmatrix} d^{(1)} - \mathbf{w}^T \mathbf{x}^{(1)} \\ d^{(2)} - \mathbf{w}^T \mathbf{x}^{(2)} \\ \dots \\ d^{(P_{app})} - \mathbf{w}^T \mathbf{x}^{(P_{app})} \end{pmatrix}$$

Vu le critère J , il est clair que l'on cherche le vecteur \mathbf{w} tel que

$$\boxed{\mathbf{X}^T \mathbf{w} = \mathbf{d}} \quad (35)$$

\mathbf{X}^T est une matrice rectangulaire, il faut calculer $(\mathbf{X}^T)^*$ la pseudo-inverse de \mathbf{X}^T

$$\mathbf{w} = (\mathbf{X}^T)^* \mathbf{d}$$

Pour calculer la pseudo-inverse, on peut remarquer qu'en multipliant l'équation (35) à gauche par \mathbf{X} , on obtient

$$\mathbf{X} \mathbf{X}^T \mathbf{w} = \mathbf{X} \mathbf{d}$$

Si $\mathbf{X} \mathbf{X}^T$ est inversible, on trouve ainsi

$$\boxed{\mathbf{w}_{PI} = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{d}}$$

Ce qui signifie que la pseudo inverse de \mathbf{X}^T est dans ce cas

$$(\mathbf{X}^T)^* = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X}$$

Est-ce que $\mathbf{X} \mathbf{X}^T$ est toujours inversible ?

\mathbf{X} est une matrice de taille $N \times P_{app}$, la matrice $\mathbf{X} \mathbf{X}^T$ est de taille $N \times N$ et peut s'écrire

$$\mathbf{X} \mathbf{X}^T = \sum_{p=1}^{P_{app}} \mathbf{x}^{(p)} (\mathbf{x}^{(p)})^T$$

De plus, $\mathbf{X} \mathbf{X}^T$ est symétrique, donc elle est décomposable dans une base orthogonale de vecteurs propres

$$\mathbf{X} \mathbf{X}^T = \mathbf{U} \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ 0 & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & \lambda_N \end{pmatrix} \mathbf{U}^T$$

où $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$ et $\mathbf{U} \mathbf{U}^T = \mathbf{Id}$. Ainsi, on a

$$\det(\mathbf{X} \mathbf{X}^T) = \prod_{n=1}^N \lambda_n$$

A quelle condition sur λ la matrice $\mathbf{X} \mathbf{X}^T$ est inversible ?

De plus, en introduisant les vecteurs colonnes $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_N)$, on peut aussi écrire

$$\mathbf{X} \mathbf{X}^T = \sum_{n=1}^N \lambda_n \mathbf{u}_n \mathbf{u}_n^T$$

Si $P_{app} < N$, alors comme on a

$$\mathbf{X} \mathbf{X}^T = \sum_{p=1}^{P_{app}} \mathbf{x}^{(p)} (\mathbf{x}^{(p)})^T$$

il est clair qu'on a forcément $\lambda_n = 0$ pour $n > P_{app}$. Donc dans ce cas, $\det(\mathbf{X} \mathbf{X}^T) = 0$ et la matrice est singulière.

Comment faire dans ce cas ?

On considère la décomposition en valeurs singulières

$$\mathbf{X} = \mathbf{U} \Delta_{\delta_{1:r}} \mathbf{V}^T$$

avec

$$\Delta_{\delta_{1:r}} = \begin{pmatrix} \delta_1 & 0 & \dots & 0 & 0 & \cdot & 0 \\ 0 & \delta_2 & \dots & 0 & 0 & \cdot & 0 \\ 0 & \cdot & \cdot & 0 & 0 & \cdot & 0 \\ 0 & \cdot & \cdot & \delta_r & 0 & \cdot & 0 \\ 0 & \cdot & \cdot & 0 & 0 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & 0 & \cdot & 0 \\ 0 & \cdot & \cdot & 0 & 0 & \cdot & 0 \end{pmatrix}$$

où r est donc le rang de la matrice.
Ainsi l'équation

$$\mathbf{X} \mathbf{X}^T \mathbf{w} = \mathbf{X} \mathbf{d}$$

s'écrit

$$\mathbf{U} \Delta_{\delta_{1:r}} \Delta_{\delta_{1:r}}^T \mathbf{U}^T \mathbf{w} = \mathbf{U} \Delta_{\delta_{1:r}} \mathbf{V}^T \mathbf{d} \quad (36)$$

On définit

$$\mathbf{w} = \mathbf{U} \Delta_{1/\delta_{1:r}} \mathbf{V}^T \mathbf{d}$$

Est-ce que ce vecteur \mathbf{w} est solution de l'équation (36) ?

En pratique, pour calculer \mathbf{w} , plutôt que de calculer la décomposition en valeurs singulières de \mathbf{X} , on préfère souvent utiliser la décomposition en valeurs propres

$$\mathbf{X} \mathbf{X}^T = \mathbf{U} \Delta_{\lambda_{1:r}} \mathbf{U}^T$$

Quel est le lien entre λ_n et δ_n ?

A partir de cette décomposition, il est possible de définir

$$(\mathbf{X} \mathbf{X}^T)^* = \mathbf{U} \Delta_{1/\lambda_{1:r}} \mathbf{U}^T$$

pour finalement retrouver \mathbf{w} avec

$$\mathbf{w} = (\mathbf{X} \mathbf{X}^T)^* \mathbf{X} \mathbf{d}$$

car en effet, on a :

$$(\mathbf{X} \mathbf{X}^T)^* \mathbf{X} \mathbf{d} = \mathbf{U} \Delta_{1/\lambda_{1:r}} \mathbf{U}^T \mathbf{U} \Delta_{\delta_{1:r}} \mathbf{V}^T \mathbf{d} = \mathbf{U} \Delta_{1/\delta_{1:r}} \mathbf{V}^T \mathbf{d}$$

Ainsi dans le cas où $\mathbf{X} \mathbf{X}^T$ n'est pas inversible, la solution s'écrit simplement

$$\mathbf{w} = (\mathbf{X} \mathbf{X}^T)^* \mathbf{X} \mathbf{d}$$

Cette solution sera analysé lors du prochain TP et comparé à celle de M. Hebb.

2.1.4 Analyse des performances (avec barre d'erreur)

On dispose d'une base $\mathcal{B} = \{(\mathbf{x}^{(p)}, d^{(p)}); p \in \{1, \dots, P\}\}$.

Comment faire pour analyser les performances d'un discriminateur ?

La solution la plus courante consiste à couper la base en deux. On utilise une partie pour réaliser l'apprentissage (i.e. apprendre \mathbf{w}) et une autre partie pour réaliser l'analyse des performances de ce discriminateur.

Mathématiquement on définit une partition $\{1, \dots, P\} = I_{app} \cup I_{gen}$ avec $I_{app} \cap I_{gen} = \emptyset$.

Ainsi, on a une base d'apprentissage

$$\mathcal{B}_{app} = \{(\mathbf{x}^{(p)}, d^{(p)}); p \in I_{app}\}$$

pour apprendre à discriminer et une base de généralisation

$$\mathcal{B}_{gen} = \{(\mathbf{x}^{(p)}, d^{(p)}); p \in I_{gen}\}$$

pour analyser les performances.

Dans le reste du cours, on note P_{app} la taille de la base d'apprentissage et P_{gen} la taille de la base de généralisation.

Pourquoi on coupe la base en 2 ?

Comment analyser les performances d'un discriminateur ?

Pour n'importe quel vecteur $\mathbf{x}^{(p)}$ de la base, on peut définir le score du discriminateur \hat{d} avec

$$\eta(\mathbf{x}^{(p)}) = \begin{cases} 1 & \text{si } \hat{d}(\mathbf{x}^{(p)}) = d^{(p)} \\ 0 & \text{si } \hat{d}(\mathbf{x}^{(p)}) \neq d^{(p)} \end{cases}$$

On peut alors calculer le taux de réussite obtenu sur la base d'apprentissage

$$\tau_{app} = \frac{1}{P_{app}} \sum_{p \in I_{app}} \eta(\mathbf{x}^{(p)})$$

On appellera ce taux, le *taux d'apprentissage*.

Néanmoins, ce qui nous intéresse pour caractériser les performances, c'est plutôt le taux de réussite obtenu sur la base de généralisation

$$\tau_g = \frac{1}{P_{gen}} \sum_{p \in I_{gen}} \eta(\mathbf{x}^{(p)})$$

que l'on appelle *taux de généralisation*.

Cette fois, la somme est réalisée sur les vecteurs de la base de généralisation, et ceux-ci sont **différents** de ceux de la base d'apprentissage.

Le problème pratique est le suivant : **quand on a deux discriminateurs (par exemple en TP, vous aurez $\tau_g^{(Hebb)}$ et $\tau_g^{(PI)}$), comment savoir si les différences de performances sont significatives ?**

Pour répondre à cette question, il faut aborder des questions préliminaires.

Comment représenter $\eta(\mathbf{x}^{(p)})$?

Quelle est la loi suivie par $\eta(\mathbf{x}^{(p)})$?

Par conséquent **comment représenter τ_g ?**

Quelle est la loi suivie par τ_g ?

Enfin, on revient sur la question initiale, **quand on a deux discriminateurs (par exemple en TP, on aura $\tau_g^{(Hebb)}$ et $\tau_g^{(PI)}$), comment savoir si les différences de performances sont significatives ?**

Comment calculer les écart-types de τ_g ?

Question préliminaire : **que peut-on dire de la variance de τ_g ?**

$$\text{var}(\tau_g) = \left\langle \left(\frac{1}{P_{gen}} \sum_p \eta_p - \tau_{gen} \right)^2 \right\rangle$$

devient

$$\text{var}(\tau_g) = \frac{1}{P_{gen}^2} \left\langle \left(\sum_p (\eta_p - \tau_{gen}) \right)^2 \right\rangle$$

En supposant les η_p indépendants, on obtient

$$\text{var}(\tau_g) = \frac{1}{P_{gen}^2} \sum_p \langle (\eta_p - \tau_{gen})^2 \rangle$$

et comme les η_p sont identiquement distribués, on a

$$\text{var}(\tau_g) = \frac{P_{gen}}{P_{gen}^2} \langle (\eta_p - \tau_{gen})^2 \rangle = \frac{1}{P_{gen}} \text{var}(\eta_p)$$

Or pour une loi de Bernoulli, on a $\text{var}(\eta_p) = \tau_{gen}(1 - \tau_{gen})$, donc on obtient

$$\text{var}^{1/2}(\tau_g) = \sqrt{\frac{\tau_{gen}(1 - \tau_{gen})}{P_{gen}}}$$

Remarque. En pratique, on a pas accès à τ_{gen} . **Comment faire pour estimer $\text{var}^{1/2}(\tau_g)$?** (voir TP)