

1 Apprentissage machine

1.1 Probabilité

Soit \mathbb{P} une fonction $\mathbb{P}: \mathcal{F} \mapsto \mathbb{R}$ sur certains¹ sous-ensembles $A \in \mathcal{F} \subseteq 2^\Omega$. L'ensemble Ω s'appelle *l'univers des événements atomiques* $x \in \Omega$. \mathbb{P} est une **probabilité** si et seulement si elle satisfait les **axiomes de Kolmogorov** comme suit :

★ la probabilité est toujours entre 0 et 1 :

$$0 \leq \mathbb{P}A \leq 1; \quad (1.1)$$

★ l'univers comprenant tous les événements est certain

$$\mathbb{P}\Omega = 1; \quad (1.2)$$

★ elle est σ -additive : pour toute suite (infinie) d'ensembles $A_i: i = 0, 1, 2, \dots$ disjoints ($A_i \cap A_j = \emptyset$ pour tout $i \neq j$),

$$\mathbb{P} \bigcup_{i=0}^{\infty} A_i = \sum_{i=0}^{\infty} \mathbb{P}A_i. \quad (1.3)$$

1.2 Variables aléatoires (v.a.)

Une variable aléatoire X est définie par sa probabilité $(\Omega, \mathcal{F}, \mathbb{P}_X)$. Ainsi, la probabilité $\mathbb{P}_X\{X \in A\} = \mathbb{P}A$ est bien définie pour tout $A \in \mathcal{F}$. On écrit $\mathbb{E}X$ pour l'**espérance** de X . Pour tout ensemble de variables aléatoires (définies sur le même univers)

$$\mathbb{E}[X_1 + X_2 + X_3 + \dots] = (\mathbb{E}X_1) + (\mathbb{E}X_2) + (\mathbb{E}X_3) + \dots \quad (1.4)$$

(Notez que l'indépendance des X_i n'est pas important ici.)

Indicateurs. Un **indicateur** pour l'ensemble A est une v.a. \mathbb{I}_A qui prend la valeur 0 ou 1 selon la mesure de A :

$$\begin{aligned} \mathbb{P}_{\mathbb{I}_A}\{\mathbb{I}_A = 1\} &= \mathbb{P}A \\ \mathbb{P}_{\mathbb{I}_A}\{\mathbb{I}_A = 0\} &= \mathbb{P}\bar{A} = 1 - \mathbb{P}A. \end{aligned}$$

¹La domaine \mathcal{F} de \mathbb{P} est σ -algèbre (ou *tribu*) ce qui veut dire que (1) si un sous-ensemble $A \in \mathcal{F}$ est mesurable, alors son complémentaire $\bar{A} = \Omega \setminus A$ l'est aussi, et que (2) si les ensembles $A_i: i = 0, 1, 2, \dots$ sont mesurables, alors leur union $\bigcup_{i=0}^{\infty} A_i$ l'est aussi. En conséquence, \mathcal{F} est stable aussi par intersection et par soustraction.

Pour économie, on écrit l'indicateur de A simplement par A . Comme cela, $\{X = 1\}$ dénote l'événement que $X = 1$, ainsi que l'indicateur X pour ce même événement. Soit X l'indicateur pour A . On a alors $\mathbb{E}X = \mathbb{P}\{X = 1\} = \mathbb{P}A$, ou autrement $\mathbb{E}A = \mathbb{P}A$.

1.3 Conditionnalité

Soit $(\Omega, \mathcal{F}, \mathbb{P})$ un espace probabilisé. La probabilité conditionnée sur un événement $A \in \mathcal{F}$ avec $\mathbb{P}A > 0$ est l'espace $(A, \mathcal{F}, \mathbb{P}_{|A})$, avec $\mathbb{P}_{|A}B = \mathbb{P}\{B \mid A\}$ définie par $\mathbb{P}A \cap B = \mathbb{P}\{B \mid A\} \cdot \mathbb{P}A$ pour tout $B \in \mathcal{F}$. (Si $\mathbb{P}A = 0$, $\mathbb{P}\{B \mid A\}$ est non-déterminée.) Par Axiome (1.3),

$$\mathbb{P}B = \mathbb{P}B \cap A + \mathbb{P}B \cap \bar{A} = \mathbb{P}\{B \mid A\} \cdot \mathbb{P}A + \mathbb{P}\{B \mid \bar{A}\} \cdot \mathbb{P}\bar{A}, \quad (1.5a)$$

et, en général, pour tout $A_i: i = 0, 1, 2, \dots$ disjoints avec $A_i \cap A_j = \emptyset$ pour $i \neq j$ et $\bigcup_i A_i = \Omega$, on obtient la **formule des probabilités totales** :

$$\mathbb{P}B = \sum_i \mathbb{P}\{B \mid A_i\} \cdot \mathbb{P}A_i. \quad (1.5b)$$

L'espérance conditionnelle s'écrit par $\mathbb{E}[X \mid A]$. Si X, Y sont des v.a. réelles, $\mathbb{E}[X \mid Y]$ est une variable aléatoire : elle prend la valeur $\mathbb{E}[X \mid Y = y]$ quand $Y = y$. On a alors $\mathbb{E}[\mathbb{E}[X \mid Y]] = \mathbb{E}X$.

Théorème de Bayes. Soit $A_i: 0, 1, 2, \dots$ une famille exhaustive d'événements disjoints : $A_i \cap A_j = \emptyset$ et $\bigcup_i A_i = \Omega$. Par (1.5b),

W^(fr)

$$\mathbb{P}\{A_i \mid B\} = \frac{\mathbb{P}A_i \cap B}{\mathbb{P}B} = \frac{\mathbb{P}\{B \mid A_i\} \cdot \mathbb{P}A_i}{\mathbb{P}B} = \frac{\mathbb{P}\{B \mid A_i\} \cdot \mathbb{P}A_i}{\sum_j \mathbb{P}\{B \mid A_j\} \cdot \mathbb{P}A_j}. \quad (1.6)$$

Équation (1.6) s'appelle le théorème de Bayes.

1.4 Prédiction

On veut prédire une v.a. Y à partir d'une autre v.a. X . Dans d'autres mots, à partir d'une observation $X = x$, on prédit $\hat{y} = f(x)$, en se servant d'une fonction de prédiction f .

Classification. Si la valeur prédite prend des valeurs finies ou dénombrables, alors on parle d'un problème de classification. En particulier, on a k classes si $\mathbb{P}\{Y \in \{0, 1, \dots, k-1\}\} = 1$. Le plus commun est d'avoir $k = 2$, ou classification binaire : Y est l'indicateur d'un ensemble A . La prédiction s'évalue par la probabilité d'erreur

$$\text{erreur}(f) = \mathbb{P}\{f(X) \neq Y\}.$$

Si $k = 2$, on a

$$\begin{aligned} \text{erreur}(f) &= \mathbb{P}\{Y = 1 \mid f(X) = 0\} \cdot \mathbb{P}\{f(X) = 0\} \\ &\quad + \mathbb{P}\{Y = 0 \mid f(X) = 1\} \cdot \mathbb{P}\{f(X) = 1\} \quad (1.7) \end{aligned}$$

Erreur de Bayes. Soit $g_i(x) = \mathbb{P}\{Y = i \mid X = x\}$ pour $i = 0, 1$ et $h(x) = g_1(x) - 1/2$. Soit f^* l'indicateur

$$f^*(x) = \{h(x) > 0\} = \{g_1(x) > g_0(x)\}.$$

La fonction f^* , appelée le **prédicteur Bayes** fournit la meilleure prédiction possible. Lors d'une observation $X = x$, on devine f^* . Si $f^*(x) = 1$, $h(x) > 0$, et donc $Y = 1$ est plus probable que $Y = 0$. On se trompe avec une probabilité $\mathbb{P}\{Y = 0 \mid X = x\} = g_0(x) = 1/2 - h(x)$. Si, par contre, $f^*(x) = 0$, alors $h(x) \leq 0$ et on se trompe avec une probabilité $\mathbb{P}\{Y = 1 \mid X = x\} = g_1(x) = 1/2 + h(x)$. Dans les deux cas, l'erreur est $1/2 - |h(x)|$. Par conséquent, l'**erreur Bayes** s'écrit par

$$\begin{aligned} \text{erreur}(f^*) &= \mathbb{P}\{Y = 1 \mid h(X) \leq 0\} \cdot \mathbb{P}\{h(X) \leq 0\} \\ &\quad + \mathbb{P}\{Y = 0 \mid h(X) > 0\} \cdot \mathbb{P}\{h(X) > 0\} \\ &= \mathbb{E}[1/2 - |h(x)|] \end{aligned}$$

Pour voir que f^* est la meilleure prédiction possible, soit f un prédicteur quelconque. Si $f^*(x) = 1$ et $f(x) = 0$, alors f se trompe avec une probabilité $\mathbb{P}\{Y = 1 \mid X = x\} = 1/2 + h(x)$. Si $f^*(x) = 0$ et $f(x) = 1$, alors f se trompe avec une probabilité $\mathbb{P}\{Y = 0 \mid X = x\} = 1/2 - h(x)$. Dans les deux cas, la probabilité d'erreur est $1/2 + |h(x)|$. On a donc

$$\text{erreur}(f) = \text{erreur}(f^*) + 2\mathbb{E}[|h(X)| \mid f(X) \neq f^*(X)] \cdot \mathbb{P}\{f(X) \neq f^*(X)\}.$$

Régression. Si Y est réelle ($Y \in \mathbb{R}$) on parle d'un problème de régression. L'erreur d'un prédicteur f se mesure par la distance Δ entre la prédiction et la valeur observée : $\Delta(f(X), Y)$. Typiquement, on vise à minimiser soit l'erreur L_1 : $\Delta(f(X), Y) = |f(X) - Y|$, ou l'erreur L_2 avec $\Delta(f(X), Y) = (f(X) - Y)^2$.

1.5 Apprentissage machine

En pratique, la distribution jointe de X et Y est rarement connue, et donc le prédicteur Bayes n'est pas disponible. On construit un prédicteur plutôt à partir d'un échantillon d'apprentissage (*learning sample*) D . L'échantillon comprend n paires observées : $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$. Typiquement, on assume que les paires sont indépendantes et identiquement distribuées (iid). En apprentissage machine, on cherche des algorithmes de construction de prédicteurs.

Apprentissage non-supervisé. Quand on n'a qu'un échantillon $D = \{x_1, x_2, \dots, x_n\}$ sans étiquettes, on veut «découvrir» des classes à partir des similarités observées parmi les points x_i . On cherche un *groupage* (*clustering*) des données : il s'agit d'une fonction $f: \mathcal{R}^d \mapsto \{0, 1, 2, \dots, k-1\}$ pour classifier les points x_i (de d dimensions) en k classes différentes. (Comment définir l'erreur d'un groupage ? On retournera à cette question plus tard.)

1.6 Modèles paramétriques

Les méthodes d'apprentissage machine peuvent être paramétriques ou non : tombent en une des classes suivantes

- ★ **méthodes non-paramétriques** : la prédiction se fait à partir de l'échantillon directement. Exemples : plus proche voisin, histogramme.
- ★ **méthodes paramétriques** : la prédiction se fait en choisissant le meilleur prédicteur à partir d'une classe fixée. Exemples : machine à vecteurs de support, réseau de neurones.

En bio-informatique, on se sert souvent de modèles paramétriques qui capturent les aspects spécifiques du problème biologique analysé. Un tel modèle impose des dépendances entre composantes de l'observation X et la sortie Y . En particulier, le modèle donne corps à une hypothèse, permettant sa validation et quantification. Exemple : modèles d'évolution de séquences moléculaires.