

Introduction à l'apprentissage automatique, la science de l'intelligence artificielle

Séance 3

Théorie statistique de la décision et applications

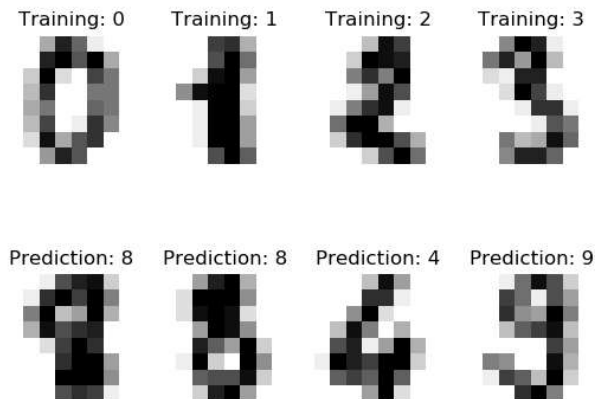
Frédéric Sur

https://members.loria.fr/FSur/enseignement/IMT_GE/

Plan

- 1 Classification et décision
 - Éléments de théorie statistique de la décision
 - Le « meilleur » classifieur : classifieur de Bayes
 - Classifieur naïf de Bayes
- 2 Mise en œuvre du classifieur de Bayes
 - Classification aux plus proches voisins
 - Régression logistique
- 3 Conclusion

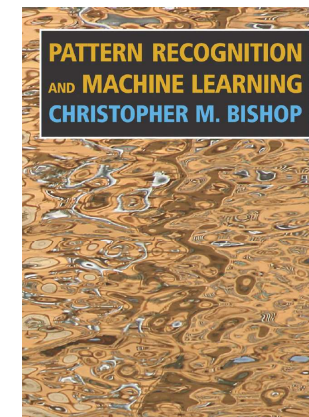
Classification supervisée



Aujourd'hui : classification supervisée

Cadre : théorie statistique de la décision

Bibliographie



Source principale (& illustrations sauf mention contraire) :
C. Bishop, *Pattern Recognition and Machine Learning*, Springer 2006
→ dans toutes les bonnes médiathèques

Notations

K classes : $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_K$

classes = partition de l'ensemble des observations possibles

Exemple : $\mathcal{C}_1=1, \mathcal{C}_2=2, \dots, \mathcal{C}_{10}=0$

N observations (base d'apprentissage) : $x_1, \dots, x_N \in \mathbb{R}^d$

Exemple : $x \in \mathbb{R}^{256}$ = vecteur des niveaux de gris d'une image 16 × 16
on connaît la classe d'appartenance de chaque x_i : $y_i (= \mathcal{C}_1, \mathcal{C}_2, \dots)$

Question du jour : prédiction de la classe d'une nouvelle observation x ?

Soit f un classifieur : $f(x) = \mathcal{C}_1$ ou $f(x) = \mathcal{C}_2$, etc.

→ une fonction sur l'ensemble des entrées possibles qui prédit la classe d'appartenance

Problème : le classifieur peut faire des erreurs

Exemple : $f(x) = \mathcal{C}_1$ alors que $y = \mathcal{C}_2$

5/22

Formalisation probabiliste

Probabilités :

- $p(\mathcal{C}_k)$ probabilité a priori (prior), $\sum_k p(\mathcal{C}_k) = 1$
→ ce qu'on suppose sans connaître d'observations
Exemple OCR : $p(a) = 0.07, p(b) = 0.01, p(c) = 0.03 \dots$
- $p(x|\mathcal{C}_k)$ proba. conditionnelle (en fait, densité - vraisemblance)
→ probabilité pour qu'une obs. tirée dans la classe \mathcal{C}_k vaille x
Exemple : $p(x|\mathcal{C}_k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1}(x-\mu_k)}$
- $p(\mathcal{C}_k|x)$ probabilité (vraisemblance) a posteriori
→ probabilité de la classe \mathcal{C}_k étant donnée l'observation x

Théorème de Bayes :

$$\forall 1 \leq k \leq K, p(\mathcal{C}_k|x) = \frac{p(x|\mathcal{C}_k)p(\mathcal{C}_k)}{p(x)}$$

$$\text{et } p(x) = \sum_{k=1}^K p(x, \mathcal{C}_k) = \sum_{k=1}^K p(x|\mathcal{C}_k)p(\mathcal{C}_k)$$

6/22

Erreur de classification

Pour simplifier, $K = 2$ dans la suite

Le classifieur f définit une partition de l'ensemble des observations possibles en régions \mathcal{R}_i telles que :

$$\mathcal{R}_i = \{x, f(x) = \mathcal{C}_i\}$$

Problème : les partitions (\mathcal{C}_i) et (\mathcal{R}_i) ne coïncident pas
(à cause des erreurs de classification)

Proposition : calcul de la proportion moyenne d'erreur théorique

$$\begin{aligned} E_{\text{err}} &= E_{X,Y} (1_{f(X) \neq Y}) = \iint 1_{f(x) \neq y} p(x, y) dx dy \\ &= \int 1_{f(x) \neq \mathcal{C}_1} p(x, \mathcal{C}_1) dx + \int 1_{f(x) \neq \mathcal{C}_2} p(x, \mathcal{C}_2) dx \\ &= \int_{\mathcal{R}_2} p(x, \mathcal{C}_1) dx + \int_{\mathcal{R}_1} p(x, \mathcal{C}_2) dx \end{aligned}$$

Question : existe-t-il un classifieur f minimisant E_{err} ?

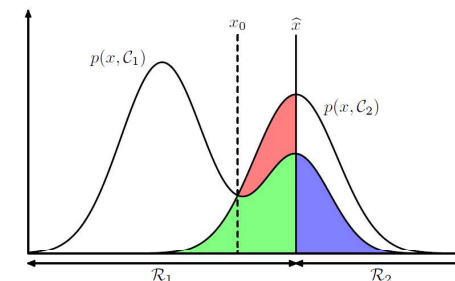
7/22

Minimisation de l'erreur moyenne

$$E_{\text{err}} = \int_{\mathcal{R}_2} p(x, \mathcal{C}_1) dx + \int_{\mathcal{R}_1} p(x, \mathcal{C}_2) dx$$

Illustration : $x \in \mathbb{R}, \mathcal{R}_1 = \{x \in \mathbb{R}, x < \hat{x}\}, \mathcal{R}_2 = \{x \in \mathbb{R}, x > \hat{x}\}$

Question : comment fixer \hat{x} de manière à minimiser E_{err} ?



$E_{\text{err}} = \text{rouge} + \text{vert} + \text{bleu}$

Or $\text{vert} + \text{bleu} = \text{Cste}$

lorsque \hat{x} varie

→ minimum atteint pour $\hat{x} = x_0$.

8/22

Classifieur de Bayes

(ce raisonnement se généralise à $x \in \mathbb{R}^d$, et $\mathcal{R}_i \neq$ intervalles)

Conséquence : la règle de classification minimisant l'erreur moyenne est

$$f(x) = \operatorname{argmax}_{C_k} p(x, C_k) = \operatorname{argmax}_{C_k} p(x) p(C_k|x)$$

Classifieur de Bayes – maximum a posteriori (MAP)

$$f(x) = \operatorname{argmax}_{C_k} p(C_k|x) = \operatorname{argmax}_{C_k} p(C_k) p(x|C_k)$$

ou, dans le cas biclasse :

$$f(x) = C_1 \text{ ssi } \frac{p(x|C_1)}{p(x|C_2)} > \frac{p(C_2)}{p(C_1)}$$

(cf test du rapport des vraisemblances)

9/22

En pratique ? (1)

Problème : comment estimer $p(C_k)$ et $p(x|C_k)$ à partir du jeu de données disponible (ensemble d'observations classifiées) ?

$p(C_k)$:

- information connue a priori
exemple : OCR
- ou fréquence estimée à partir de la base d'observations
exemple : $p(C_k) = \frac{\#\{x_i \in C_k, 1 \leq i \leq N\}}{N}$
où $\#$ désigne le cardinal d'un ensemble
- ou, dans le cas où les $p(C_k)$ sont égaux :
le classifieur de Bayes se simplifie en $f(x) = \operatorname{argmax}_k p(x|C_k)$
→ règle du maximum de vraisemblance (ML)

10/22

En pratique ? (2)

$p(x|C_k)$: toute une partie de l'apprentissage non-supervisé concerne l'estimation de densités de probabilité

si $x \in \mathbb{R}^d$ avec $d \ll \text{grand} \gg$: attention, *curse of dimensionality* !

Expérience :

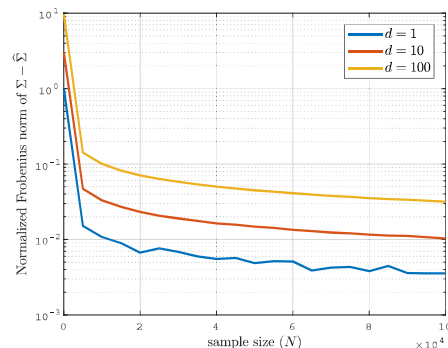
estimation matrice de covariance

N observations x_n i.i.d. $\mathcal{N}_d(0, \Sigma)$

$$\hat{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N x_n x_n^T$$

$$\|\Sigma - \hat{\Sigma}\|_{\text{F}}^2 = \frac{1}{d} \sum_{i=1}^d \lambda_i^2$$

où λ_i v.p. de $\Sigma - \hat{\Sigma}$



→ en pratique, on a intérêt à réduire la dimension d ...

11/22

Classifieur naïf de Bayes

Une manière de battre la malédiction de la dimensionnalité...

Si $x = (x^1, x^2, \dots, x^d) \in \mathbb{R}^d$, on suppose les composantes (conditionnellement) statistiquement indépendantes

$$\text{Donc : } p(x|C_k) = \prod_{i=1}^d p(x^i|C_k)$$

Gros avantage : plutôt qu'estimer la distribution $p(x|C_k)$ sur \mathbb{R}^d , on estime les d distributions $p(x^i|C_k)$ sur \mathbb{R} .

Classifieur naïf de Bayes

$$f(x) = \operatorname{argmax}_{C_k} p(C_k) \prod_{i=1}^d p(x^i|C_k)$$

Exemple : classifieur naïf Gaussien

→ on suppose les distributions $p(x^i|C_k)$ gaussiennes

→ deux paramètres : μ_k, σ_k

12/22

Plan

1 Classification et décision

- Éléments de théorie statistique de la décision
- Le « meilleur » classifieur : classifieur de Bayes
- Classifieur naïf de Bayes

2 Mise en œuvre du classifieur de Bayes

- Classification aux plus proches voisins
- Régression logistique

3 Conclusion

13/22

Estimateur de distribution aux plus proches voisins

On va chercher à estimer $\phi(x_0) = p(x_0|C_k) \dots$

Pour estimer $\phi(x_0)$, on considère une boule B_{x_0} centrée en x_0 , contenant P échantillons de \mathcal{C} parmi $M = \#\mathcal{C}$:

- $\int_B \phi(x) dx \simeq \frac{P}{M}$ (par la loi des grands nombres)
- et $\int_B \phi(x) dx \simeq \phi(x_0) \times \text{Volume}(B_{x_0})$ (si ϕ constant sur B_{x_0})

Donc : $\frac{P}{M} \simeq \int_B \phi(x) dx \simeq \phi(x_0) \times \text{Volume}(B_{x_0})$

D'où l'estimateur des P plus proches voisins :

$$\phi(x) = \frac{P}{M V_P(x)}$$

où $V_P(x)$: volume d'une boule contenant les P p.p.v.

→ hypothèse ϕ constant sur B , donc il faudrait une boule « pas trop grosse » (malédiction dimensionnalité ?)

→ en fait, utilisé pour la classification supervisée. . .

14/22

Classification aux P plus proches voisins

Base de données : N observations x_1, \dots, x_N et classes associées parmi $\mathcal{C}_1, \dots, \mathcal{C}_K$.

→ N_1 observations dans $\mathcal{C}_1, \dots, N_K$ dans \mathcal{C}_K , t.q. $\sum_k N_k = N$

Problème : étant donnée une nouvelle observation x , comment prédire sa classe ?

Parmi les P p.p.v. de x : P_1 dans $\mathcal{C}_1, \dots, P_K$ dans \mathcal{C}_K ($\sum_k P_k = P$)

Par estimation aux P -p.p.v. : $p(x|C_k) = \frac{P_k}{N_k V_P(x)}$

De plus : $p(C_k) = \frac{N_k}{N}$

Classifieur MAP :

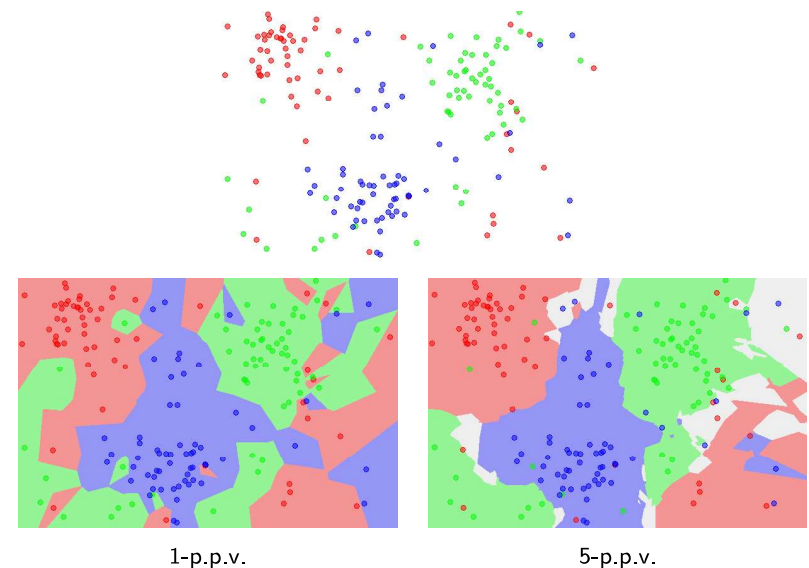
$\arg\max_k p(C_k)p(x|C_k) = \arg\max_k P_k / (N V_P(x)) = \arg\max_k P_k$

Définition : le classifieur aux P -p.p.v. affecte une nouvelle observation à la classe majoritaire parmi les P observations les plus proches

→ implémente le classifieur de Bayes sous hypothèses (très) simplificatrices

15/22

Exemple



By Agor153 - Own work, CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=24350617>

16/22

La régression logistique

Dans le cas bi-classe :

$$p(C_1|x) = \frac{p(C_1)p(x|C_1)}{p(C_1)p(x|C_1) + p(C_2)p(x|C_2)} = \frac{1}{1 + \frac{p(C_2)p(x|C_2)}{p(C_1)p(x|C_1)}}$$

Avec $f(x) = \log\left(\frac{p(x|C_2)}{p(x|C_1)}\right) + \log\left(\frac{p(C_2)}{p(C_1)}\right)$: $p(C_1|x) = \frac{1}{1 + e^{-f(x)}}$

Définition : fonction logistique (ou sigmoïde) : $\sigma(t) = \frac{1}{1+e^{-t}}$

Hypothèse simplificatrice : $f(x) = \beta_0 + \beta_1 \cdot x$

→ c'est aussi une manière de contrer la malédiction de la dimensionnalité, en réduisant le nombre de paramètres à estimer

→ si on sait estimer β_0 et β_1 , le classifieur de Bayes devient :

classifieur de la régression logistique :

x dans C_1 ssi $p(C_1|x) > 1/2 \iff f(x) > 0$

(séparation des deux classes par un hyperplan)

17/22

Régression logistique et estimation des paramètres

Hypothèse simplificatrice : $f(x) = \beta_0 + \beta_1 \cdot x$

Remarque : cas de classes gaussiennes de même variance

si $p(x|C_k) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)}$,

$$f(x) = -\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1) + \frac{1}{2}(x - \mu_2)^T \Sigma^{-1}(x - \mu_2) + \log\left(\frac{p(C_1)}{p(C_2)}\right)$$

→ d'où $f(x) = \beta_0 + \beta_1 \cdot x$ avec :

$$\beta_0 = \frac{1}{2}\mu_2^T \Sigma^{-1}\mu_2 - \frac{1}{2}\mu_1^T \Sigma^{-1}\mu_1 + \log\left(\frac{p(C_1)}{p(C_2)}\right)$$

$$\beta_1 = \Sigma^{-1}(\mu_1 - \mu_2)$$

Estimation de (β_0, β_1) :

maximisation de la log-vraisemblance conditionnelle

$$\begin{aligned} \ell_{(x_i, y_i)_{1 \leq i \leq N}}(\beta_0, \beta_1) &= \log \prod_{i=1}^n \left((p(C_1|x_i))^{y_i} (1 - p(C_1|x_i))^{1-y_i} \right) \\ &= \sum_{i=1}^n \left(-(1 - y_i)(\beta_0 + \beta_1 x_i) - \log(1 + \exp(\beta_0 + \beta_1 x_i)) \right) \end{aligned}$$

avec **ici** : $y_i = 1 \iff x_i \in C_1$, $y_i = 0 \iff x_i \in C_2$, $p(C_2|x_i) = 1 - p(C_1|x_i)$

→ fonction **concave**, donc maximum unique

18/22

Illustration

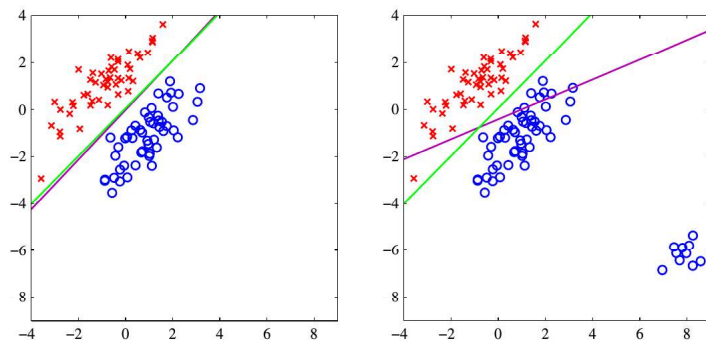


Figure 4.4 The left plot shows data from two classes, denoted by red crosses and blue circles, together with the decision boundary found by least squares (magenta curve) and also by the logistic regression model (green curve), which is discussed later in Section 4.3.2. The right-hand plot shows the corresponding results obtained when extra data points are added at the bottom left of the diagram, showing that least squares is highly sensitive to outliers, unlike logistic regression.

Illustration : C. Bishop, *Pattern Recognition and Machine Learning*, Springer 2006

→ on reviendra sur cette propriété en séance 4

19/22

Plan

- 1 Classification et décision
 - Éléments de théorie statistique de la décision
 - Le « meilleur » classifieur : classifieur de Bayes
 - Classifieur naïf de Bayes
- 2 Mise en œuvre du classifieur de Bayes
 - Classification aux plus proches voisins
 - Régression logistique
- 3 Conclusion

20/22

Fonction discriminante

Notion de **fonction discriminante** :

f_k tel que $f(x) = \operatorname{argmax}_k f_k(x)$

- MAP (classifieur de Bayes, **théoriquement** optimal) :

$$f(x) = \operatorname{argmax}_k p(C_k|x) = \operatorname{argmax}_k p(C_k)p(x|C_k)$$

Problème : on ne connaît pas les distributions de probabilité

→ on ajoute des hypothèses : classifieur naïf de Bayes, régression logistique, k -p.p.v...

- Cours suivants : autres hypothèses, autres fonctions discriminantes

Conclusion – Résumé

Théorie statistique de la décision, et mise en œuvre :

- le classifieur bayésien (MAP) minimise l'erreur moyenne de classification (classifieur **idéal**)
- simplification si *prior* uniformes : classification au maximum de vraisemblance (ML)
- simplification si composantes indépendantes pour $x \in \mathbb{R}^d$: classifieur naïf de Bayes
- simplification si $p(C_1|x) = \sigma(\beta_0 + \beta_1 \cdot x)$: régression logistique
- simplification si $p(x|C_1)$ ne varie pas trop et suffisamment d'observations : classification aux K plus proches voisins