

SPLEX

Statistiques pour la classification et fouille de données en génomique

Classification probabiliste - Classification Linéaire Binaire

Pierre-Henri WUILLEMIN

Decision-axe IA-LIP6
`pierre-henri.wuillemin@lip6.fr`


Rappels : Variables aléatoires et probabilités

Soient X et Y deux variables aléatoires sur Ω .

➡ Définition

- *loi marginale de X*
 $p(X = x) = p(\{X = x\} \mid \Omega)$
- *loi jointe de X et Y*
 $p(X = x, Y = y) = p(\{X = x\} \cap \{Y = y\} \mid \Omega)$
- *loi conditionnelle de X sachant Y*
 $p(X = x \mid Y = y) = p(\{X = x\} \mid \{Y = y\}, \Omega)$

Propriété

- $\sum_x p(X = x \mid Y = y) = 1$
-  $\sum_y p(X = x \mid Y = y) \neq 1$
- $\sum_x p(X = x, Y = y \mid Z = z) = p(Y = y \mid Z = z)$

Rappels : Probabilités, fonctions — théorème de Bayes

- $p(X = x|Y = y) \in [0, 1]$ est une probabilité.
- $p(X|Y = y)$ est une distribution de probabilité (une fonction).
- $p(X|Y)$ est une famille de distributions de probabilités sur des espaces différents.

Il y a une différence de nature entre $p(X | Y)$ et $p(Y | X)$. On connaît néanmoins un rapport entre ces grandeurs : $p(X, Y) = p(X | Y) \cdot p(Y) = p(Y | X) \cdot p(X)$

Théorème (de Bayes)

$$p(X | Y) = \frac{p(Y | X) \cdot p(X)}{p(Y)}$$

ou encore : $p(X | Y, Z) = \frac{p(Y|X,Z) \cdot p(X|Z)}{p(Y|Z)}$

Rappels : Variables aléatoires et indépendances

Soient X , Y et Z des variables aléatoires sur Ω .

➡ Définition (Indépendance conditionnelle)

$X \perp\!\!\!\perp Y | Z$ selon p si et seulement si $p(X|Y, Z) = p(X|Z)$

➡ Définition (indépendance marginale)

$X \perp\!\!\!\perp Y$ selon p si et seulement si $p(X|Y) = p(X)$

$p(X|Y, Z) = p(X|Z)$ peut se lire :

Augmenter notre état de connaissance Z en apprenant Y n'influence pas la distribution de probabilité qu'on attribue à X .

À noter que cette relation est **symétrique** : $X \perp\!\!\!\perp Y | Z$ si et seulement si $Y \perp\!\!\!\perp X | Z$.

On remarque que, si $X \perp\!\!\!\perp Y$ alors $p(X, Y) = p(X | Y) \cdot p(Y) = p(X) \cdot p(Y)$

De même, si $X \perp\!\!\!\perp Y | Z$ alors $p(X, Y | Z) = p(X | Y, Z) \cdot p(Y | Z) = p(X | Z) \cdot p(Y | Z)$

Approche probabiliste de la classification

Soient, à nouveau, deux v.a. X (de dimension d) discrète et Y (de dimension 1) discrète (*pas forcément binaire*).

Sur la base Π_a , on peut estimer les probabilités par des fréquences pour $P(X, Y)$.
Soit x une instantiation de X , on cherche sa classe y (valeur de Y).

1 Maximum de vraisemblance

$$y_{ML}^* = \arg \max_{y_i} P(x | y_i)$$

2 Maximum a posteriori

$$y_{MAP}^* = \arg \max_{y_i} P(y_i | x)$$

D'après la règle de Bayes, $P(Y | X) \propto P(X | Y) \cdot P(Y)$, on comprend que l'intérêt du MAP est de prendre en compte un *a priori* sur la fréquence de chaque classe.



Il peut être difficile d'obtenir ces distributions.

Particulièrement : $P(X | Y)$ peut demander beaucoup d'observation !!

Exo 1 : français ou suédois ?

On considère deux attributs pour déterminer la nationalité d'un individu.

L'attribut taille qui peut prendre les valeurs grand ou petit, l'attribut couleur des cheveux qui peut prendre les valeurs brun ou blond. Les nationalités possibles sont français et suédois.

On suppose que les populations françaises et suédoises se répartissent selon le tableau suivant :

	petit, brun	petit, blond	grand, brun	grand, blond
Suédois	10	20	30	40
Français	25	25	25	25

- ❶ Dans une assemblée comprenant 60% de suédois et 40% de français, décrire
 - ❶ la règle de décision majoritaire
 - ❷ la règle du maximum de vraisemblance
 - ❸ la règle de Bayes
- ❷ Calculez les probabilités d'erreur de chacune de ces règles
- ❸ On suppose maintenant que l'on ne connaît plus les proportions respectives des suédois et des français. On note p la proportion des suédois ($p \in [0, 1]$). Décrire, selon les valeurs possibles de p , les règles de Bayes correspondantes.

Classifieur Bayésien Naïf

Il peut être difficile d'obtenir $P(Y)$, $P(X | Y)$, $P(Y | X)$.

Particulièrement : $P(X | Y)$ peut demander beaucoup d'observation !!

Hypothèse du classifieur bayésien naïf

On supposera que, $\forall k \neq l, X^k \perp\!\!\!\perp X^l | Y$

Cette hypothèse est très forte. Elle a peu de chance de s'avérer exacte dans un cas réel. Néanmoins cette approximation donne des résultats souvent satisfaisants.

Alors, le calcul du MAP s'écrit :

3 Maximum a posteriori

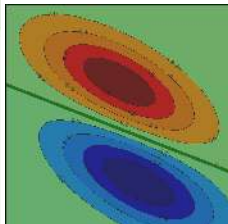
$$y = \arg \max_{y_i} \left(P(y_i) \cdot \prod_{k=1}^d P(x^k | y_i) \right)$$

Cette hypothèse permet donc de simplifier fortement les calculs nécessaires pour estimer le MAP.

Cas gaussien

Cadre

- Modèle : $\hat{C}(x) = \sigma(g(x)) = \sigma(g_{\oplus}(x) - g_{\ominus}(x))$
- Régions de décision :
 $\forall c \in \{\ominus, \oplus\}, R_c = \{x \in \mathcal{D}, \hat{C}(x) = c\}$
- Frontière de décision : $F = \{x \in \mathcal{D}, \hat{C}(x) = 0\}$
- **Multinormalité** : $\forall c \in \{\ominus, \oplus\}, g_c(x) = P(x | c) \sim \mathcal{N}(\mu_c, \Sigma_c)$



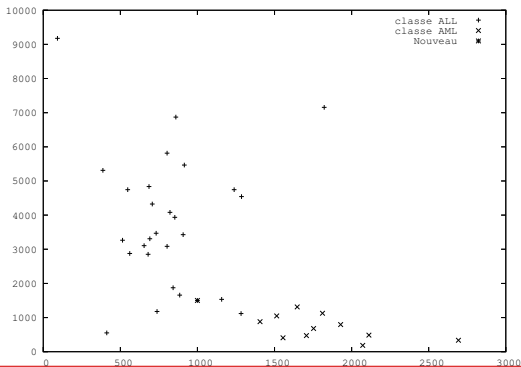
➡ Définition (Densité normale)

$$p(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^t \Sigma^{-1}(x-\mu)} \quad \text{où}$$

- μ le vecteur des moyennes,
- Σ la matrice de covariance :
 - $\forall i, j \in \{1, \dots, d\}, \Sigma_{[i,j]} = \text{cov}(X_i, X_j)$ et $\Sigma_{[i,i]} = \text{var}(X_i)$
 - Σ semi-définie positive ($\forall x, x^t \Sigma x \geq 0$)
 - $|\Sigma| \geq 0$ et Σ^{-1} existe

Un exemple

Les données –les couples (gène 1, gène 2) – suivent un certain modèle : deux distributions de probabilité sur D_X (une pour chaque classe).



Le classifieur Bayésien naïf - 1

Données :

- l'ensemble d'apprentissage $X = (x_1, x_2, \dots, x_n)$ où $x_i = (v_{i1}, v_{i2}, \dots, v_{ip})$
- un nouvel élément $x = (v_1, \dots, v_p)$

Problème : trouver la classe de x .

Approche probabiliste :

$$c_{MAP} = \operatorname{argmax}_c P(c|x) = \operatorname{argmax}_c \frac{P(x|c)P(c)}{P(x)} = \operatorname{argmax}_c P(x|c)$$

Il faut estimer les probabilités $P(c)$ et $P(x|c)$ à partir de X .

Le classifieur Bayésien naïf - 2

- Estimation de $P(c)$

$$P(0) = n_0/n = 27/38$$

$$P(1) = n_1/n = 11/38$$

- Le problème est plus compliqué pour $P(x|c)$

Approche naïve Bayes : les attributs sont indépendants étant donnée la classe :

$$P(x|c) = P(g_1|c)P(g_2|c)$$

Il faut estimer $P(g_1|c)$ et $p(g_2|c)$ à partir de X .

Le classifieur Bayésien naïf - 3

Hypothèse de normalité :

- $P(g_1|c) \rightsquigarrow \mathcal{N}(g_1, \mu_{g_1,c}, \sigma_{g_1,c}^2)$
- $P(g_2|c) \rightsquigarrow \mathcal{N}(v_p, \mu_{g_2,c}, \sigma_{g_2,c}^2)$

On n'a plus qu'à estimer les moyennes et écart-types à partir de X .

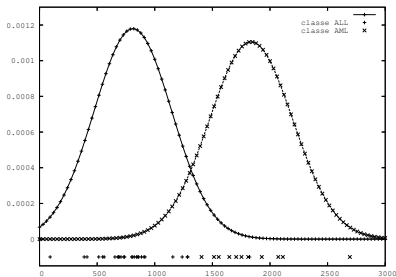
$$\mu_{g_1,0} \approx 810.29 \qquad \sigma_{g_1,0} \approx 338.17$$

$$\mu_{g_1,1} \approx 1836.27 \qquad \sigma_{g_1,1} \approx 360.88$$

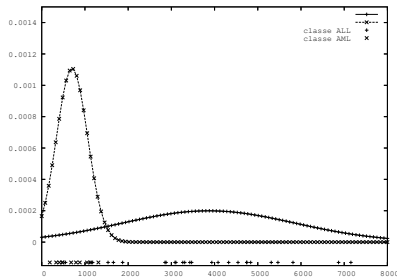
$$\mu_{g_2,0} \approx 3863.29 \qquad \sigma_{g_2,0} \approx 1999.78$$

$$\mu_{g_2,1} \approx 702.36 \qquad \sigma_{g_2,1} \approx 360.34$$

Le classifieur Bayésien naïf - 4



Fumarylacetoacetate



C-myb

Le classifieur Bayésien naïf - 5

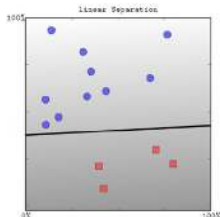
On veut déterminer la classe de $x = (1000, 1500)$:

$$\begin{aligned} P(0|x) &= P(0) \mathcal{N}(1500, \mu_{g_{1,0}}, \sigma_{g_{1,0}}^2) \mathcal{N}(4000, \mu_{g_{2,0}}, \sigma_{g_{2,0}}^2) \\ &= 0.71053 \times 1.01 \cdot 10^{-3} \times 9.92 \cdot 10^{-5} = 7.11 \cdot 10^{-8} \end{aligned}$$

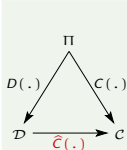
$$\begin{aligned} P(1|x) &= P(1) \mathcal{N}(1500, \mu_{g_{1,1}}, \sigma_{g_{1,1}}^2) \mathcal{N}(4000, \mu_{g_{2,1}}, \sigma_{g_{2,1}}^2) \\ &= 0.28947 \times 7.54 \cdot 10^{-5} \times 9.55 \cdot 10^{-5} = 2.09 \cdot 10^{-9} \end{aligned}$$

$$\Rightarrow c_{MAP} = 0$$

Classification linéaire binaire (CLB)



➡ Définition (CLB)



- $\mathcal{C} = \{\ominus, \oplus\}$

- $\exists w \in \mathbb{R}^d, w_0 \in \mathbb{R}, \exists f : \mathbb{R} \rightarrow \mathcal{C},$

$$\forall x \in \mathbb{R}^d, \hat{C}(x) = f\left(\sum_{i=1}^d w_i \cdot x_i + w_0\right)$$

Le problème d'apprentissage : trouver w , w_0 et f .

Modèles génératifs, modèles discriminants

- **Modèles génératifs** : classification grâce à une estimation de $P(x, y)$ à partir de Π_a et des connaissances *a priori*.
 - Classifieur bayésien (ML, MAP)
 - Classifieur bayésien naïf
 - Discriminant linéaire de Fisher
- **Modèles discriminants** : estimation directe des w , w_0 à partir de Π_a .
 - Régression logistique
 - Perceptron
 - SVM

Le classifieur bayésien naïf binaire est un CLB ?

Classifieur bayésien naïf

$$y = \arg \max_{y_i} \left(P(y_i) \cdot \prod_{k=1}^d P(x^k | y_i) \right)$$

- Ici, $y_0 = \ominus$ et $y_1 = \oplus$.
- Soit $R(x) = \frac{P(\oplus) \cdot \prod_{k=1}^d P(x^k | \oplus)}{P(\ominus) \cdot \prod_{k=1}^d P(x^k | \ominus)}$
- Si $R(x) > 1$ alors $\hat{C}(x) = \hat{\oplus}$ sinon $\hat{C}(x) = \hat{\ominus}$
- Donc $\hat{C}(x) = \sigma(\log R(x))$ où $\sigma(u) = \begin{cases} -1 & \text{si } u < 0 \\ 0 & \text{si } u = 0 \\ +1 & \text{sinon} \end{cases}$
- Il vient alors

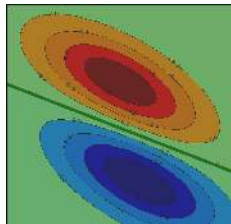
$$\hat{C}(x) = \sigma \left(\log \frac{P(\oplus)}{P(\ominus)} + \sum_{k=1}^d \log \frac{P(x^k | \oplus)}{P(x^k | \ominus)} \right)$$

Suite évidente dans le cas binomial ($\mathcal{D} = \{\ominus, \oplus\}^d$).

Discrimination linéaire - cas gaussien

Cadre gaussien

- Modèle : $\hat{C}(x) = \sigma(g(x)) = \sigma(g_{\oplus}(x) - g_{\ominus}(x))$
- Régions de décision :
 $\forall c \in \{\ominus, \oplus\}, R_c = \{x \in \mathcal{D}, \hat{C}(x) = c\}$
- Frontière de décision : $F = \{x \in \mathcal{D}, \hat{C}(x) = 0\}$
- Multinormalité : $\forall c \in \{\ominus, \oplus\}, P(x | c) \sim \mathcal{N}(\mu_c, \Sigma_c)$



CLB

Si **homoscédasticité** : $\forall c, \Sigma_c = \Sigma$ alors, la fonction discriminante devient linéaire :

$$g(x) = (\mu_{\oplus} - \mu_{\ominus})^t \Sigma^{-1} (x - x_0)$$

$$\text{avec } x_0 = \frac{1}{2} (\mu_{\oplus} + \mu_{\ominus}) + \left(\frac{1}{(\mu_{\oplus} - \mu_{\ominus})^t \Sigma^{-1} (\mu_{\oplus} - \mu_{\ominus})} \log \frac{P(\oplus)}{P(\ominus)} \right) (\mu_{\oplus} - \mu_{\ominus})$$

Rappels de géométrie

Soit $y(x) = \sum_{i=1}^d w_i \cdot x_i + w_0 \Rightarrow \hat{C}(x) = f(y(x))$, on peut également écrire :

$$y(x) = w' \cdot x + w_0 \quad \text{avec } y(x) = 0 \text{ l'équation d'un hyperplan } H$$

$$\forall a, b \in H, y(a) = y(b) = 0 \Rightarrow y(a) - y(b) = w' \cdot (a - b) = 0$$

w est un vecteur normal à H .

Soit $x \notin H$ et x_H sa projection perpendiculaire sur H , $x - x_H$ est donc colinéaire à w ,
Soit $r \in \mathbb{R}$, $x - x_H = r \cdot \frac{w}{\|w\|}$ où r est la **distance de x à H** .

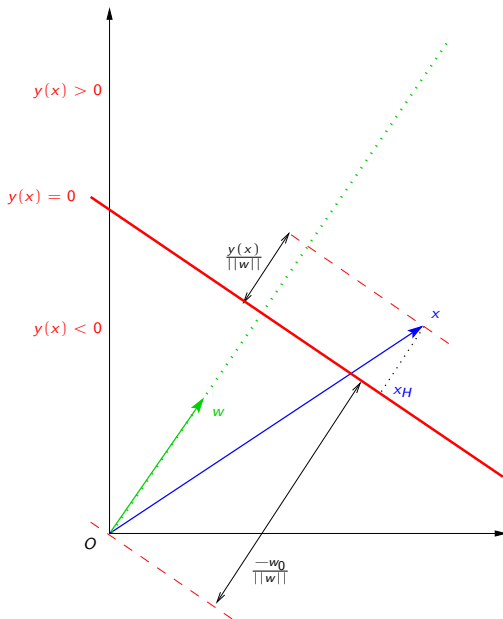
$$x = x_H + r \cdot \frac{w}{\|w\|}$$

$$w' \cdot x = w' \cdot x_H + r \cdot \frac{w' \cdot w}{\|w\|} = w' \cdot x_H + r \cdot \frac{\|w\|^2}{\|w\|} = w' \cdot x_H + r \cdot \|w\|$$

$$y(x) = w' \cdot x + w_0 = w' \cdot x_H + w_0 + r \cdot \|w\| = y(x_H) + r \cdot \|w\| = r \cdot \|w\|$$

$$\text{distance de } x \text{ à } H : r = \frac{y(x)}{\|w\|}$$

Rappels de géométrie



exemple : Hyper-plan séparateurs

La frontière entre les deux classes est donnée par $\sum_{i=1}^d w_i \cdot x_i + w_0 = 0$ qui est l'équation d'un hyper-plan.
Comment choisir cet hyper-plan ?

Exemple : CLB par régression linéaire

- Ajuster un modèle linéaire \hat{l}_k pour chaque fonction indicatrice d'une classe k :

$$\forall k \in \{\oplus, \ominus\}, \hat{l}_k(x) = \begin{cases} 1 & \text{si } x \text{ est de classe } k \\ 0 & \text{sinon.} \end{cases}$$

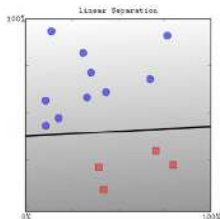
$$\hat{l}_{\oplus}(x) = \beta_{\oplus 0} + \beta'_{\oplus} \cdot x \text{ et } \hat{l}_{\ominus}(x) = \beta_{\ominus 0} + \beta'_{\ominus} \cdot x$$

- Soit un x à classifier : $\hat{C}(x) = \arg \max_k \hat{l}_k(x) = \sigma(\hat{l}_{\oplus} - \hat{l}_{\ominus})$

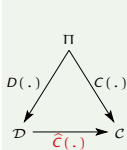
Frontière de décision :

$$\hat{f}_{\oplus}(x) = \hat{f}_{\ominus}(x) \Rightarrow \text{hyperplan : } \begin{cases} w = \beta_{\oplus} - \beta_{\ominus} \\ \text{et} \\ w_0 = \beta_{\oplus 0} - \beta_{\ominus 0} \end{cases}$$

Séparabilité



➡ Définition (CLB)



- $\mathcal{C} = \{\ominus, \oplus\}$

- $\exists w \in \mathbb{R}^d, w_0 \in \mathbb{R}, \exists f : \mathbb{R} \rightarrow \mathcal{C},$

$$\forall x \in \mathbb{R}^d, \hat{C}(x) = f\left(\sum_{i=1}^d w_i \cdot x_i + w_0\right)$$

Le problème d'apprentissage : trouver w, w_0 (**W**) et f (souvent σ).

Séparabilité sur Π_a

Soit une base de données $\Pi_a = (\mathbf{x}_i, y_i)_{i \leq N}$ où y_i est la classe de \mathbf{x}_i ($\in -1, +1$). Π_a est linéairement séparable si il existe un hyperplan d'équation $y(x) = w' \cdot x + w_0 = 0$ tel que

$$\forall i \in \{1, \dots, N\}, y(\mathbf{x}_i) \cdot y_i > 0 \text{ i.e. } (\mathbf{X} \cdot \mathbf{W}) \times \mathbf{Y} > 0$$

Optimisation de \mathbf{W} : moindres carrés



Carl Friedrich Gauss

- $\mathbf{X} \cdot \mathbf{W} - \mathbf{Y}$ est le vecteur des erreurs effectuées en classant Π_a à l'aide de \mathbf{W} .
- L'erreur quadratique obtenue sur Π_a se calcule donc comme :

$$e^2(\mathbf{W}) = (\mathbf{X} \cdot \mathbf{W} - \mathbf{Y})' \cdot (\mathbf{X} \cdot \mathbf{W} - \mathbf{Y})$$

- Minimiser cette erreur en annulant le gradient donne :

$$\mathbf{W}^* = (\mathbf{X}' \cdot \mathbf{X})^{-1} \cdot \mathbf{X}' \cdot \mathbf{Y} = \mathbf{X}^\dagger \cdot \mathbf{Y}$$

$\mathbf{X}^\dagger = (\mathbf{X}' \cdot \mathbf{X})^{-1} \cdot \mathbf{X}'$ est la pseudo-inverse de \mathbf{X} .

Cette méthode souffre de plusieurs problèmes :

- Instabilité numérique (pour des \mathbf{X} de grande taille principalement),
- Manque de robustesse pour des distributions larges de classes.

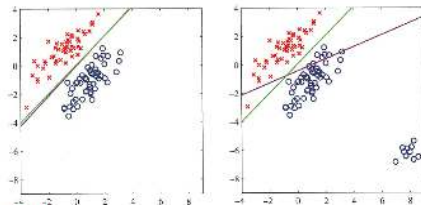


Figure 4.4 The left plot shows data from two classes, denoted by red crosses and blue circles, together with the decision boundary found by least squares (magenta curve) and also by the logistic regression model (green curve), which is discussed later in Section 4.3.2. The right-hand plot shows the corresponding 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.

From : *Pattern Recognition and Machine Learning* – C.Bishop – p186



Discriminants de Fisher : séparation entre les classes

- On note que $y = \mathbf{w}' \cdot \mathbf{x}$ correspond à la projection de \mathbf{x} (de dimension $d + 1$) sur la droite vectorielle \mathbf{w} .
- Soit $\mathbf{M}_{\oplus} = \frac{1}{N_{\oplus}} \sum_{i \in \oplus} \mathbf{X}_i$ et $\mathbf{M}_{\ominus} = \frac{1}{N_{\ominus}} \sum_{i \in \ominus} \mathbf{X}_i$
- On peut alors utiliser $\Delta_{\mathbf{w}} = \mathbf{w}' \cdot (\mathbf{M}_{\oplus} - \mathbf{M}_{\ominus})$ comme mesure de la séparation des classes selon \mathbf{w} .
Afin de supprimer l'influence sur $\Delta_{\mathbf{w}}$ de la norme de \mathbf{w} , on peut soit normaliser \mathbf{w} , soit utiliser $\frac{\Delta_{\mathbf{w}}}{\|\mathbf{w}\|}$ comme mesure.

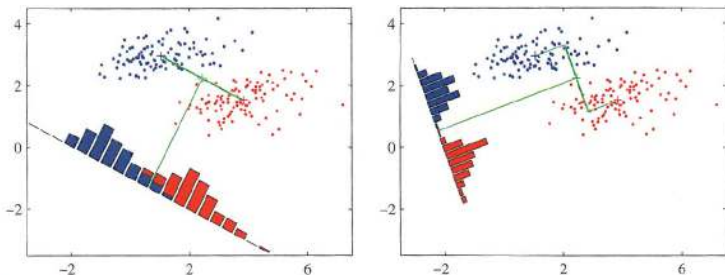


Figure 4.6 The left plot shows samples from two classes (depicted in red and blue) along with the histograms resulting from projection onto the line joining the class means. Note that there is considerable class overlap in the projected space. The right plot shows the corresponding projection based on the Fisher linear discriminant, showing the greatly improved class separation.

From : *Pattern Recognition and Machine Learning* – C.Bishop – p188

- La séparation des classes n'est intéressante qu'en fonction de la dispersion de chaque classe, i.e.
 $\forall k \in \oplus, \ominus, s_k = \sum_{i \in k} (y_i - \mathbf{w}' \cdot \mathbf{M}_k)^2$ les variances *intra-classe*.

De la régression linéaire vers la régression logistique

Régression linéaire

$$\hat{y}(x) = w' \cdot x + w_0$$

Frontière de séparation : hyperplan d'équation $y(x) = w' \cdot x + w_0 = 0$

En réutilisant MAP pour décider :

$$\hat{y} = \arg \max_{c \in \{\oplus, \ominus\}} p(c | x)$$

On ne peut pas ajuster linéairement une probabilité : une droite n'est pas bornée par $[0, 1]$.

Idee : La frontière de décision correspond à

$$p(\oplus | x) = p(\ominus | x) \iff \frac{p(\oplus | x)}{p(\ominus | x)} = 1 \iff \log \frac{p(\oplus | x)}{p(\ominus | x)} = 0$$

On peut renforcer l'idée que la frontière est un hyperplan (CLB) par :

Régression logistique

$$\exists w, w_0, \log \frac{p(\oplus | x)}{p(\ominus | x)} = w' \cdot x + w_0$$

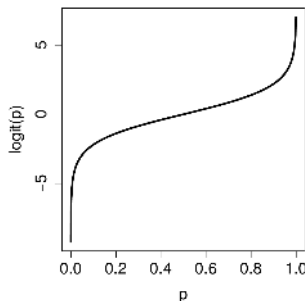
fonction logit

On peut écrire $\log \frac{p(\oplus|x)}{p(\ominus|x)} = \log \frac{p(\oplus|x)}{1-p(\oplus|x)}$

Fonction logit (*log-odds*)

$$\text{logit}(p) = \log \frac{p}{1-p}$$

La fonction logit est non bornée et donc peut être ajuster linéairement.



$$\text{logit}(p) = w' \cdot x + w_0 \iff \frac{p}{1-p} = e^{w' \cdot x + w_0} \iff p = \frac{e^{w' \cdot x + w_0}}{1 + e^{w' \cdot x + w_0}}$$

Modèle de la régression logistique

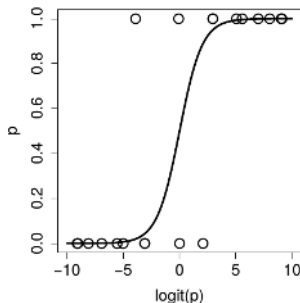
$$p(\oplus | x) = \frac{e^{w' \cdot x + w_0}}{1 + e^{w' \cdot x + w_0}} \text{ et } p(\ominus | x) = \frac{1}{1 + e^{w' \cdot x + w_0}}$$

Utilisation de la régression logistique

Soit une base $\Pi_a = (X_i, Y_i)$ avec $Y_i \in \{\oplus, \ominus\}$,

On peut calculer pour chacun $w' \cdot x + w_0$,

et donc calculer $p(\oplus | x) = \frac{e^{w' \cdot x + w_0}}{1 + e^{w' \cdot x + w_0}}$.



Estimation des paramètres w, w_0

Comment calculer les valeurs de w et w_0 de la régression logistique ?

- **Moindre carrés** ? Impossible car les erreurs ne sont pas distribuées suivant une loi normale : Elle est quasi nulle quand p proche de 0 ou 1 et plus importante quand $p \approx 0.5$.

- Utilisation du **Maximum de Vraisemblance** :

- Exprimer la vraisemblance $L(X; w, w_0)$ pour w et w_0 ,
- Essayer de maximiser la vraisemblance
- En annulant la dérivée **mais** pas de forme exacte de la dérivée.
- Utiliser une méthode approchée : *Algorithme de Newton-Raphson*.

- Soit une base de données $(X, Y)_{i \leq N}$. Avec $y_i = 1$ si \oplus et 0 si \ominus .

- $\forall i, L(x_i; w, w_0) = y_i \cdot p(x_i | \oplus) + (1 - y_i) \cdot p(x_i | \ominus)$
- Or si $\log \frac{p(\oplus|x)}{p(\ominus|x)} = w' \cdot x + w_0$ alors $\exists \beta, \beta_0, \log \frac{p(x|\oplus)}{p(x|\ominus)} = \beta' \cdot x + \beta_0$
- $p(x | \oplus) = \frac{e^{\beta' \cdot x + \beta_0}}{1 + e^{\beta' \cdot x + \beta_0}}$ et $p(x | \ominus) = \frac{1}{1 + e^{\beta' \cdot x + \beta_0}}$

Estimation des paramètres $\beta^+ = (\beta, \beta_0)$

En sommant sur toute la base la log-vraisemblance,

$$LL(\beta^+) = \sum_{i=1}^N \left[y_i \cdot (\beta^{+'} \cdot x_i^+) - \log(1 + \beta^{+'} \cdot x_i^+) \right]$$

On veut maximiser la log-vraisemblance.

$$\frac{\partial LL(\beta^+)}{\partial \beta_i^+} = \sum_{i=1}^N x_i \cdot (y_i - p(x_i; \beta^+))$$

Pas de forme simple, il faut utiliser une méthode approchée (Newton-Raphson) utilisant la dérivée seconde (le Hessian) $\frac{\partial^2 LL(\beta^+)}{\partial \beta^+ \partial \beta^{+'}}$.

La mise à jour (jusque convergence) de β^+ prend la forme :

$$\beta_{t+1}^+ = \beta_t^+ - \left(\frac{\partial^2 LL(\beta^+)}{\partial \beta^+ \partial \beta^{+'}} \right)^{-1} \cdot \frac{\partial LL(\beta^+)}{\partial \beta^+}$$

Méthode de Newton (1/2)

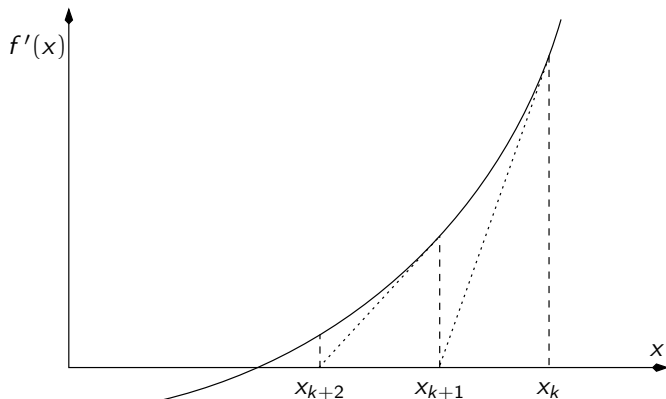
fonction de classe C^2

- $f : [a, b] \subset \mathbb{R} \mapsto \mathbb{R}$
- f : 2 fois dérivable
- f'' continue

Méthode de Newton-Raphson : recherche de 0 de la dérivée

- *principe* : engendrer une suite de points (x^k) tendant vers un point stationnaire
- point stationnaire : $f'(x^*) = 0$
- itération k : f' est remplacée par sa linéarisée en x^k :
$$l(x) = f'(x^k) + [x - x^k]f''(x^k)$$
- x^{k+1} déterminé par $l(x^{k+1}) = 0$:
$$\implies x^{k+1} = x^k - \frac{f'(x^k)}{f''(x^k)}$$

Méthode de Newton (2/2)



Un exemple (1/3)

Example

Diabetes data set

- ▶ Input X is two dimensional. X_1 and X_2 are the two principal components of the original 8 variables.
- ▶ Class 1: without diabetes; Class 2: with diabetes.
- ▶ Applying logistic regression, we obtain

$$\beta = (0.7679, -0.6816, -0.3664)^T .$$

From Jia Li (Pensylvania State University)



Un exemple (2/3)

- ▶ The posterior probabilities are:

$$\begin{aligned}Pr(G = 1 \mid X = x) &= \frac{e^{0.7679 - 0.6816X_1 - 0.3664X_2}}{1 + e^{0.7679 - 0.6816X_1 - 0.3664X_2}} \\Pr(G = 2 \mid X = x) &= \frac{1}{1 + e^{0.7679 - 0.6816X_1 - 0.3664X_2}}\end{aligned}$$

- ▶ The classification rule is:

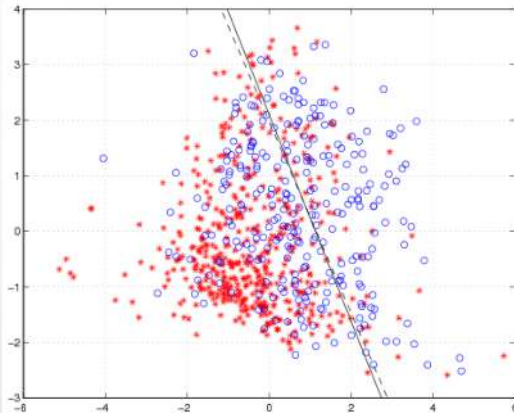
$$\hat{G}(x) = \begin{cases} 1 & 0.7679 - 0.6816X_1 - 0.3664X_2 \geq 0 \\ 2 & 0.7679 - 0.6816X_1 - 0.3664X_2 < 0 \end{cases}$$

From Jia Li (Pensylvania State University)



Un exemple (3/3)

Solid line: decision boundary obtained by logistic regression. Dash line: decision boundary obtained by LDA.



- ▶ Within training data set classification error rate: 28.12%.
- ▶ Sensitivity: 45.9%.
- ▶ Specificity: 85.8%.

From Jia Li (Pennsylvania State University)

