

Analyse, Classification et Indexation de données

Notes

Robin Montfermé

Master 1 Informatique 2019-2020

Table des matières

Introduction	4
1 Introduction à la classification de données	5
1.1 Système et modèles	5
1.1.1 Système numérique	5
1.1.2 Modèle	5
1.2 Apprentissage	6
1.2.1 L'apprentissage supervisé	6
1.2.2 L'apprentissage non-supervisé	7
1.2.3 L'apprentissage semi-supervisé	7
1.2.4 Résumé	8
1.3 L'exemple de l'usine de tri de poisson	8
1.3.1 Problème de pertinence	8
1.3.2 Problème de dimension	9
1.3.3 Réduction de la dimension des descripteur	10
1.3.4 Sensibilité des descripteurs	11
1.3.5 Comment s'assurer de la qualité du descripteur ?	11
1.3.6 Injection de données non-numériques	14
1.3.7 Résumé et remarques	14
2 Régression linéaire	16
2.1 Définition	16
2.2 Régression linéaire	18
2.2.1 F, fonction de régression linéaire	18
2.2.2 Recherche de l'équation de la droite	19
2.2.3 Petits rappels sur les vecteurs	20
2.2.4 Forme général de F	21
2.2.5 Fonction $J(\Theta)$	22
2.2.6 Dérivé partielle et Gradient	23
2.2.7 Que cherche-t-on à résoudre ?	23
2.2.8 Descente de gradient	24

3	Descente de gradient	28
3.1	Rappels	28
3.1.1	Processus itératif	28
3.1.2	Regression linéaire	28
3.1.3	Regression linéaire (suite)	29
3.2	Dérivée de la fonction $J(\theta)$	30
3.3	Reformulation sous forme matricielle	31
3.3.1	Forme Matricielle de $J(\theta)$	31
3.3.2	Dérivé de la forme matricielle	32
3.3.3	Problèmes de la forme matricielle	32
3.3.4	À retenir	32
4	Bayesienne	33
4.0.1	Principe	33
4.1	Différentes approches	33
4.1.1	Maximum à priori	33
4.1.2	Maximum vraisemblable	34
4.1.3	Maximum à posteriori	37
4.2	Loi normal en dimension N	37
4.2.1	Fonction gaussienne dans \mathbb{R}^n	37
4.2.2	Matrice de covariance	38
4.3	$\mathcal{N}(\mu, \sigma)$ en dimension 2	39
4.3.1	Distance	39
4.3.2	Courbe d'iso probabilité	39
4.3.3	Forme de la frontière de séparation	44
5	Réduction en dimension	47
5.1	Objectif	47
5.2	Avantages et inconvénients	47
5.2.1	Inconvénients	47
5.2.2	Avantages	48
5.3	Méthodes de réduction	49
5.3.1	Méthode ACP(Analyse en composante principale)	49
5.3.2	Méthode ACI(Analyse en composantes Indépendantes)	50
5.3.3	ACP : fonctionnements et calculs	51
6	Analyse discriminante	55
6.1	Calculs	55
6.1.1	Moyenne des projetés	55
6.1.2	Variance	56
6.1.3	Nouvelle forme de la fonction à minimiser	56
6.2	Généralisation à n classes	57
6.3	Que retenir ?	57

7	Classifieurs linéaires	58
7.1	Utilisation du problème à 2 classes	58
7.1.1	Explication du problème	58
7.2	Classifieur linéaire à deux classes	59
7.2.1	Première transformation de g	59
7.2.2	Seconde transformation de g	60
7.2.3	Changement de notation	60
7.2.4	Que fait-on dans le cas où $a^T < 0$?	60
7.2.5	Système d'équation	64
7.2.6	Avec le vecteur b	65

Introduction

Chargé de Cours :

Jean-Phillipe Domenger- `jean-philippe.domengerarobaselabri.fr`

Modalité d'examen :

- CC : DS + TP (fin de semestre)
- Examen

Supports

Notes Partielles : `masterinfo.emi.u-bordeaux.fr`

Objectifs

À partir d'un ensemble de points et d'un vecteur $v \in \mathbb{R}^n$ de paramètres, il faut réussir à que les points appartiennent à une classe (*x et y ∈ classe bleu*)

Contenu

Cours magistraux : présentation des concepts.

TD : implémentations des concepts. Utilisation de MATLAB.

Bibliographie

- *Pattern classification*, R.O.Duda, chez John Wiley&Sons .

Chapitre 1

Introduction à la classification de données

1.1 Système et modèles

1.1.1 Système numérique

Un système numérique est défini par 3 éléments :

- Des capteurs, recevant des données **continues**.
- Une unité de calcul, traitant des données **discrètes**.
- Un stockage.

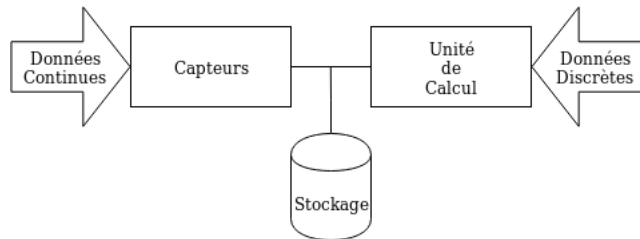


FIGURE 1.1 – Schéma d'un système numérique

1.1.2 Modèle

Comment traiter les données fournies par les capteurs au système ?

Tout algorithme est basé sur un modèle. Ce que la classification cherche à faire c'est d'intégrer à un modèle, une méthode de prise de décision, afin de créer un algorithme capable de prendre une décision.

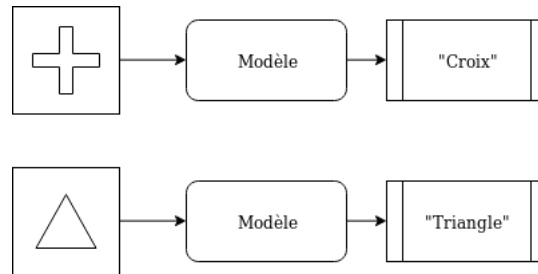


FIGURE 1.2 – Exemple de classification

1.2 Apprentissage

Pour qu'un modèle de classification soit créé, il faut passer par une phase d'apprentissage. Il y a 3 type d'apprentissage.

1.2.1 L'apprentissage supervisé

L'apprentissage supervisé consiste, à créer une observation. Une observation est un ensemble de données associées à des annotations.

Observation = (donnée 1, classe x), (donnée 2, classe y), (donnée 3, classe z)
Exemple d'observation

Une **annotation** est une indication qui identifie clairement à quelle classe appartiennent des données.

Les **données** sont des valeurs numériques obtenues via les capteurs. Elles permettent de créer des **descripteurs**.

Un **descripteur ou feature**¹ est un ensemble de données permettant d'identifier la classe à laquelle appartiennent les éléments observées.

1. Il s'agit des la notion de paramètres évoquée en introduction.

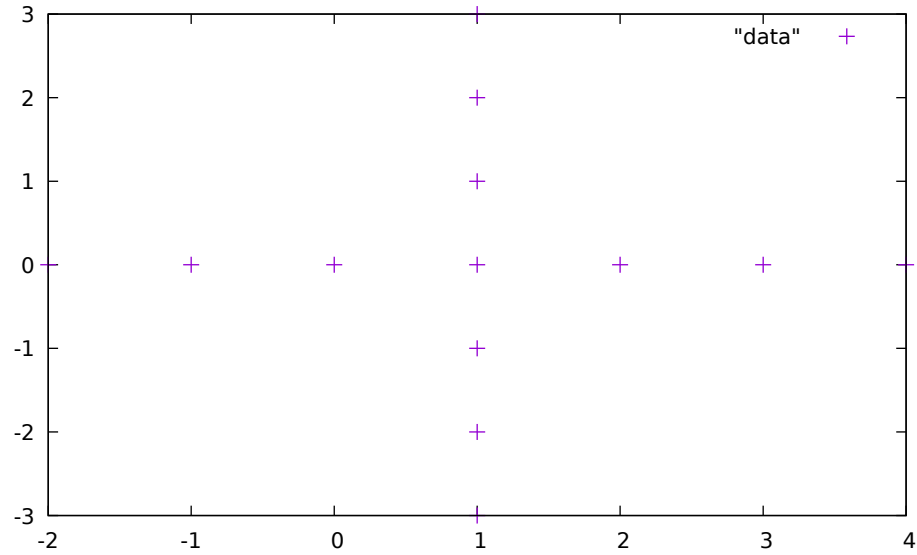


FIGURE 1.3 – Exemple de données

Les données représentées ici (1.3) peut être écrit de la façon suivante :

$$[(-1, 0), (0, 1), \dots] \in \mathbb{R}^n$$

Cet ensemble \mathbb{R}^n correspond à l'ensemble de descripteurs.

Ces ensembles de données annotées forment la **Vérité Terrain**.

En conclusion, l'apprentissage supervisé consiste à fournir au modèle une **Vérité Terrain**.

1.2.2 L'apprentissage non-supervisé

Pour cet apprentissage on ne fourni aucune annotation au modèle.

$$\{x \in \mathbb{R}^n\}$$

Le modèle n'obtient que des données

Cependant, on sait combien de classe il doit y avoir. On défini les classes à partir des données obtenues des échantillons observés.

1.2.3 L'apprentissage semi-supervisé

Il s'agit simplement d'un mélange des précédent type d'apprentissage. La **Vérité Terrain** est ici partielle

On sait :
 $(x_1, c_1), (x_2, c_2) \dots, (x_n, c_2) \dots, x_{n+1}, x_{n+2}, \dots$
 On souhaite :
 $(x_1, c_1), (x_2, c_2) \dots, (x_n, c_2) \dots, \underline{(x_{n+1}, c_2), (x_{n+2}, c_1)}, \dots$

1.2.4 Résumé

- Les **descripteurs ou features** sont des ensemble de donnée appartenant à \mathbb{R}^n .
- Les **classes** sont des données discrètes qui peuvent :
 - appartenir à \mathbb{N} si il s'agit de nombre.
 - appartenir à l'ensemble des chaînes de caractères.
- Les **observations** sont un ensemble de descripteurs et de classes.

1.3 L'exemple de l'usine de tri de poisson

Imaginons le système suivant :

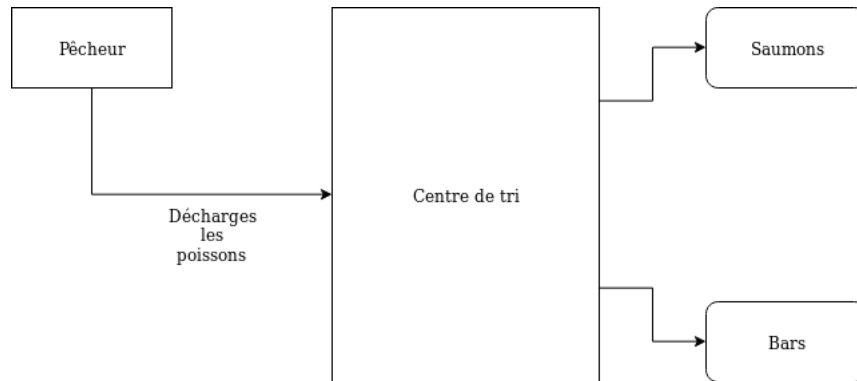


FIGURE 1.4 – Usine de tri de poisson

On veut faire en sorte que les poissons amenés par le pêcheurs soient trié dans les bacs "saumons" et "bars"². Ici on a donc identifié la classe "Saumon" et la classe "Bar".ur

1.3.1 Problème de pertinence

Il y a beaucoup de descripteurs que l'on peut obtenir à partir des poissons, le poids ou l'apparence par exemple. Mais l'important est que les descripteurs soient **pertinents**.

2. Le pêcheur n'a pêcher que ces espèces de poisson.

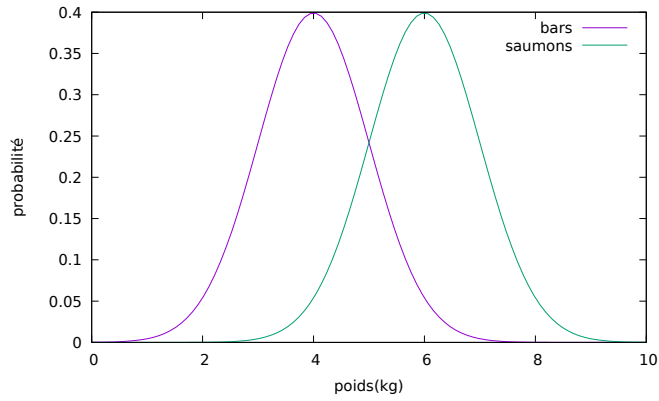


FIGURE 1.5 – Cas où le poids (seul) n’est pas un descripteur pertinent.

Sur ce graphique (1.5) on voit que les poids mesurés chez les deux espèces de poissons peuvent être très similaires, ce qui rend l’utilisation du poids en tant que descripteur n’est pas très pertinent ³.

1.3.2 Problème de dimension

On a vu que certains paramètres ne sont pas pertinents, mais il se peut que certains descripteur posent des problèmes de dimension. Si l’usine est équipée d’une caméra pour scanner les poissons, les données récoltés dépendent de la résolution de la caméra.

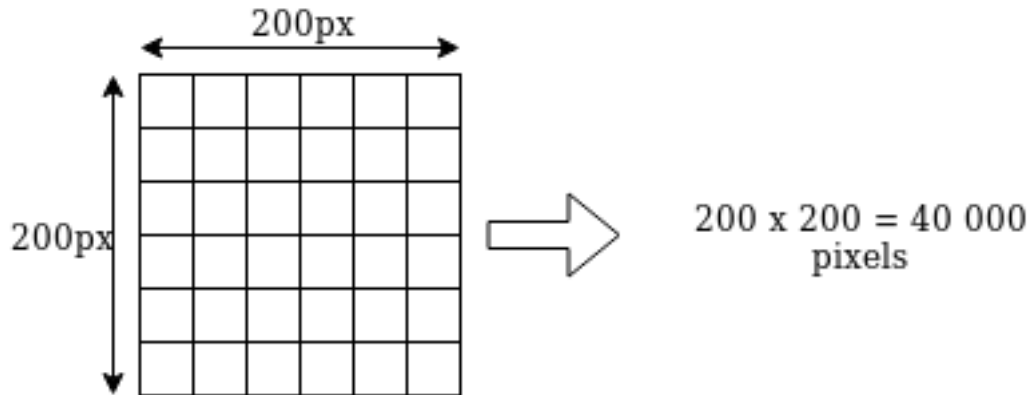


FIGURE 1.6 – Pixels de la caméra

Avec l’image obtenu de la caméra on a $200 \text{ par } 200 = 40\,000$ pixels. On rappelle qu’un descripteur est un vecteur, on a donc un descripteur $p \in \mathbb{R}^{40000}$

³. Le problème vient du fait que les courbes se croisent trop.

Cela pose un problème de **dimension**. Les données sont trop éparpillées et il se peut que lorsque en réalité, il y ait deux poissons de même classe mais étant dans une disposition différente, que ces derniers ne soient pas rassemblés sous la même classe (voir 1.7).

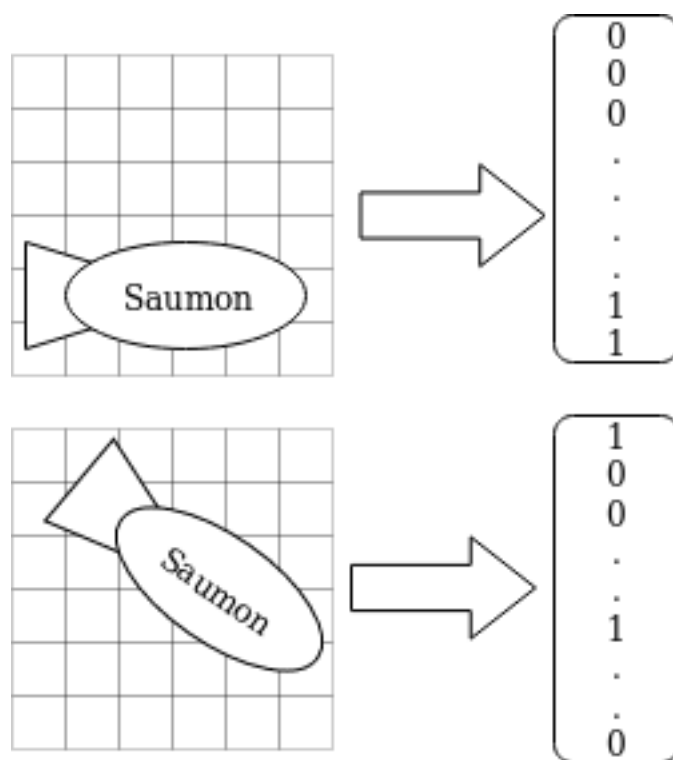


FIGURE 1.7 – Deux poissons de même classes offrant des descripteurs très différents

Il faut donc réduire la dimension de ce descripteur.

1.3.3 Réduction de la dimension des descripteur

Tout comme il y a de nombreux descripteurs, il y a de nombreux moyens de les réduire. Dans notre exemple, on va faire un traitement de l'image afin d'identifier ce qui est un poisson et ce qui ne l'est pas, puis créer la **Boîte Englobante Orientée** de poisson.

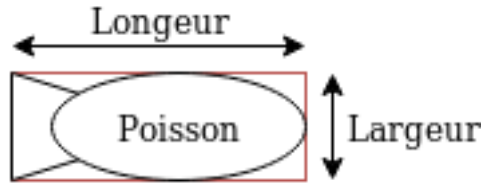


FIGURE 1.8 – Boîte englobante orienté pour les poissons

La boîte englobante orienté permet plusieurs choses :

- On obtient la longueur et la largeur du poisson, on passe de \mathbb{R}^{40000} à \mathbb{R}^2 , il n'y a plus de problème de dimension.
- Peut importe la position du poisson sur l'image on peut extraire des descripteur **pertinents**.

1.3.4 Sensibilité des descripteurs

Un autre point important est de savoir si les descripteurs sont sensibles ou non à certaines transformations.

Exemple : l'image du poisson

- Est sensible : au zoom (changement des dimensions mesurées)
- N'est pas sensible : à la rotation (Boîte Englobante Orientée)

1.3.5 Comment s'assurer de la qualité du descripteur ?

Verification graphique

Si on représente les données sur un graphes on va vouloir :

- que les points soient regroupe de manière dense
- que l'on puisse séparer les deux classes

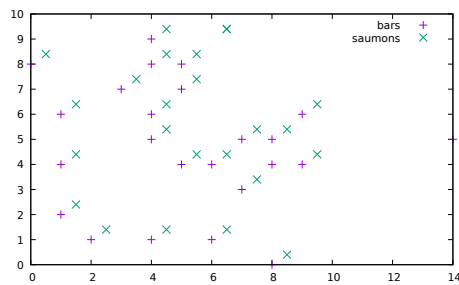


FIGURE 1.9 – Exemple de mauvais descripteurs (à gauche)

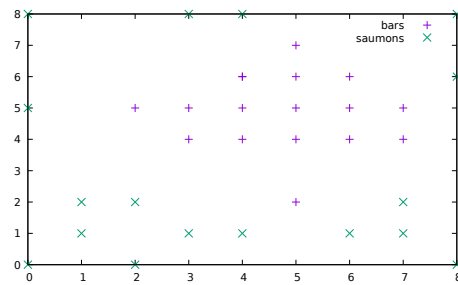


FIGURE 1.10 – Exemple de bon descripteur (à droite)

Variation de descripteur

On veut aussi que le descripteur ai peu de variation :

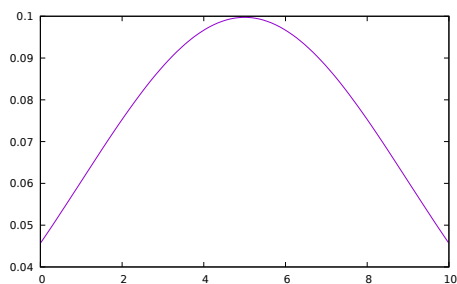


FIGURE 1.11 – Exemple de mauvaise variation (à gauche)

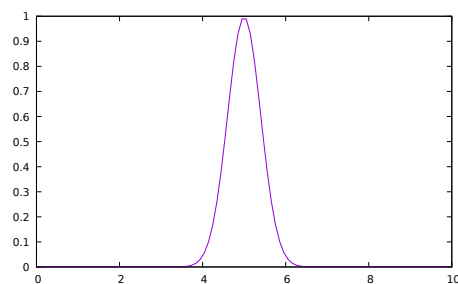


FIGURE 1.12 – Exemple de bonne variation (à droite)

Rappel sur la pertinence

On rappelle aussi qu'il faut que le descripteur doit être pertinent⁴ :

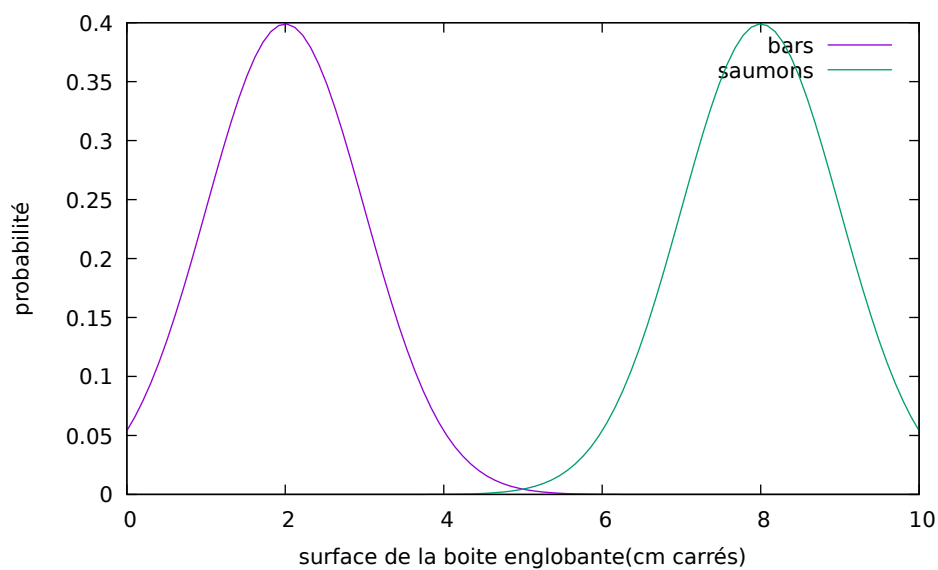


FIGURE 1.13 – La surface est plus pertinente que le poids(voir fig 1.5).

4. On remarquera au passage que la dimension de la surface de la boîte englobante est \mathbb{R}^1

Ajout de paramètre

On peut rajouter des paramètre. Par exemple, la luminance de l'image. Cela va augmenter la dimension du descripteur. Pour bien faire il faut que ce nouveau paramètre soit **indépendant** des anciens. La luminance et la surface sont des paramètres indépendant l'un de l'autre.

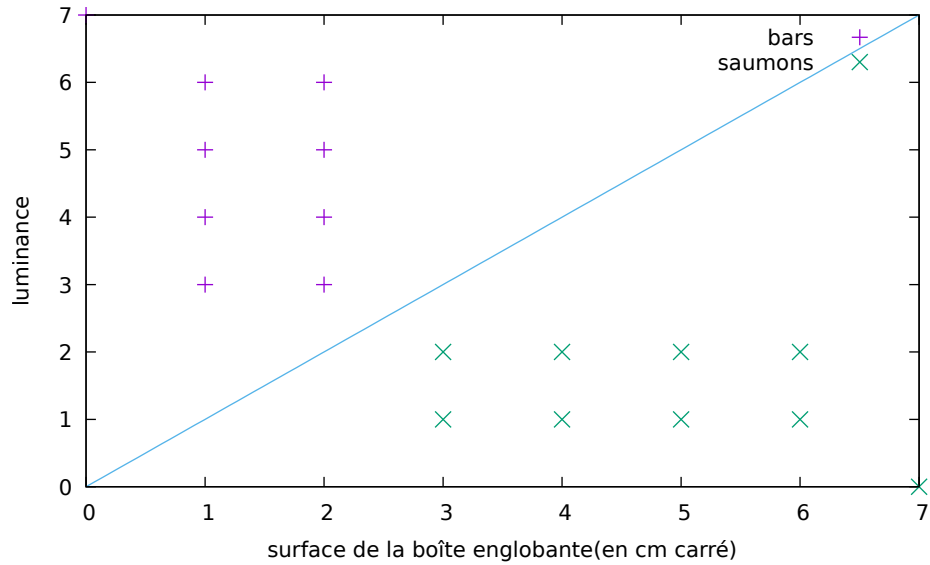


FIGURE 1.14 – L'ajout de la luminance rend le descripteur très performant

On voit bien que sur 1.14, la montée en dimension aide à la séparation.

1.3.6 Injection de données non-numériques

Il se peut que pour certaines raisons il y ait besoin d'injecter des données, non quantifiable dans le système. Par exemple, si le classifieur de l'usine doit être entraîné pendant la saison des amours du Saumon, alors il faudra peut être appliquer un seuil de décision différent qu'en temps normal.

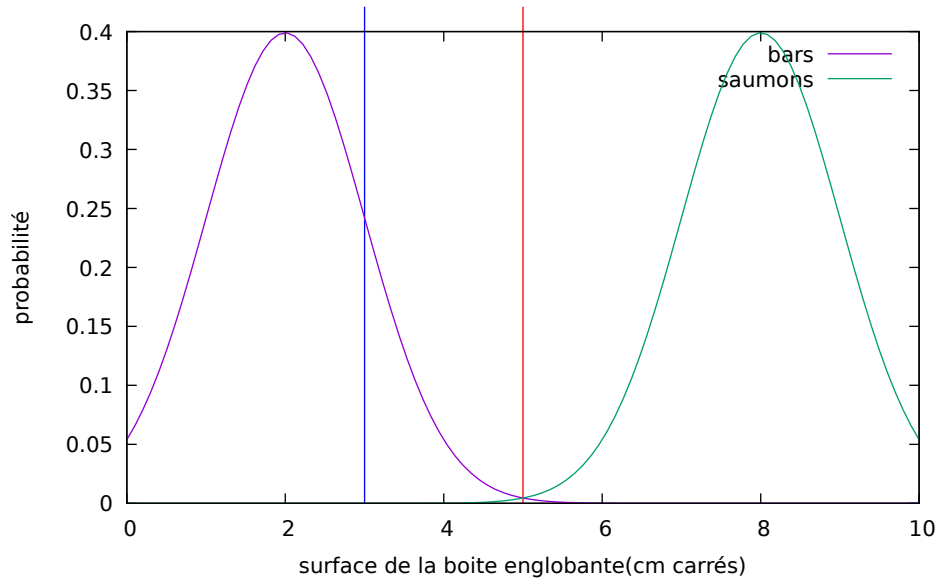


FIGURE 1.15 – En rouge seuil en temps normal, en bleu seuil pendant la saison des amours du saumon

On favorise ici un poisson plus que l'autre

1.3.7 Résumé et remarques

Résumé

- Il ne faut pas de descripteur de trop grande dimension
- Rajouter des descripteur aide à la séparation
- Les descripteurs doivent être pertinents
- Il ne doit pas avoir de corrélation entre descripteurs
- Les descripteur peuvent être mis en relation sans corrélation

Remarques

Si un classifieur obtient 0% de marge d'erreur⁵ alors il est trop adapté à la vérité terrain. Ce n'est pas une bonne chose car la vérité terrain peut être

⁵. La marge d'erreur s'évalue en fonction du pourcentage de la donnée A classé en tant que B.

amené à changer.
Le classifieur doit être adapté à la vérité terrain.

Chapitre 2

Régression linéaire

2.1 Définition

Rappel sur l'observation

Une observation est le fait d'associer des descripteurs/features à une classe/annotation. On rappelle que $Feature \in \mathbb{R}^n$ et que $classe \in \mathbb{N} | classe \in \mathbb{R}$. Entraîner un classifieur a pour but d'associer des features à des classes afin de trouver une **fonction F** qui pourra *deviner* la classe correspondant à une feature.

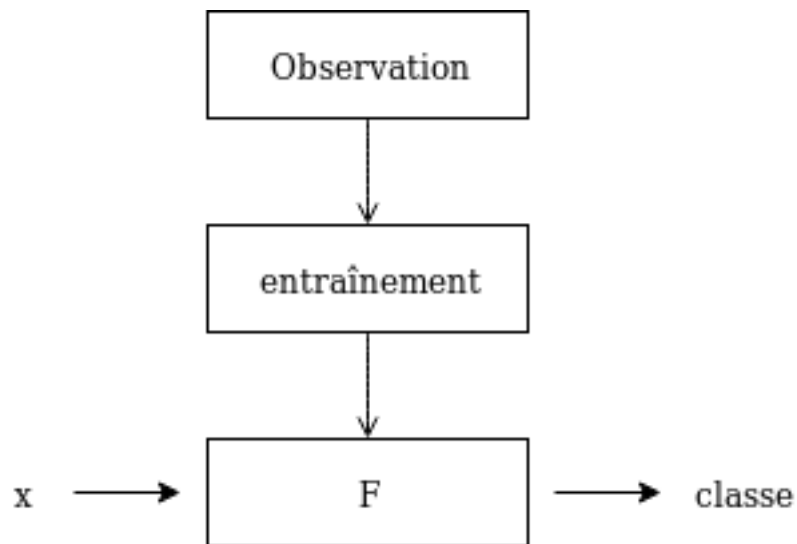


FIGURE 2.1 – L'entraînement produit une fonction F

La fonction $F : \mathbb{R}^n \rightarrow \mathbb{R}$ obtenue par entraînement du classifieur va prendre

en entrée une valeur x et va prédire une valeur y et ce à partir de données non-présentes dans l'observation¹. **Régression : Prendre quelque chose en entrée et prédire la sortie.**

Notions liés à la Regression

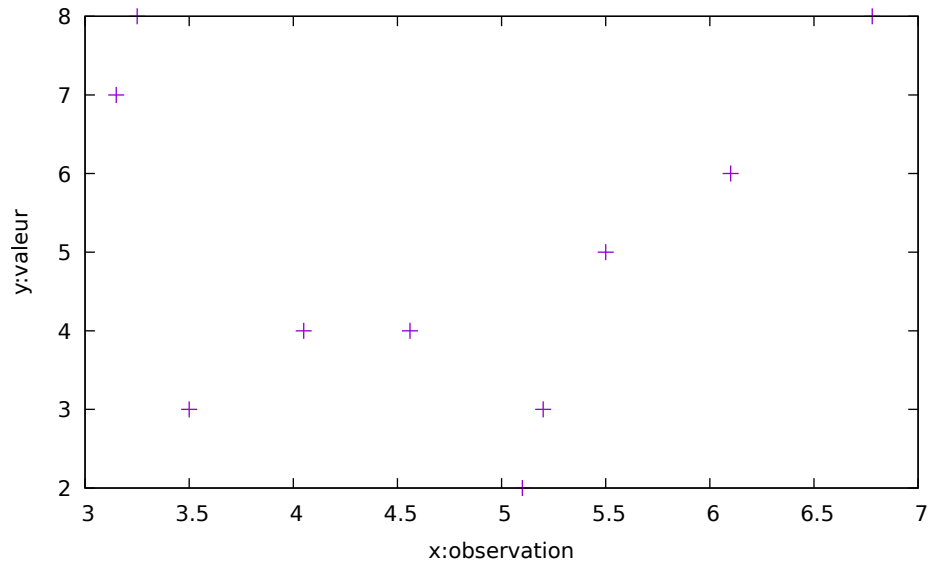


FIGURE 2.2 – Exemple d'observation

On définira m comme étant le nombre de couples (x, y) dans l'observation et n la dimension du descripteur (ici (2.2) $n = 1$).

1. non présente dans le vecteur de descripteurs.

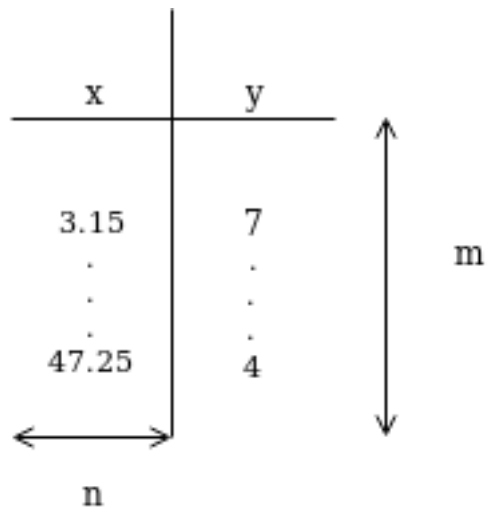


FIGURE 2.3 – Relation entre n, m, x et y

2.2 Régression linéaire

2.2.1 F , fonction de régression linéaire

A priori, F est une fonction linéaire dans un repère ou un plan².

2. On peut généraliser en parlant d'hyperplan.

2.2.2 Recherche de l'équation de la droite

Il est assez difficile de trouver une droite représentative de manière graphique.

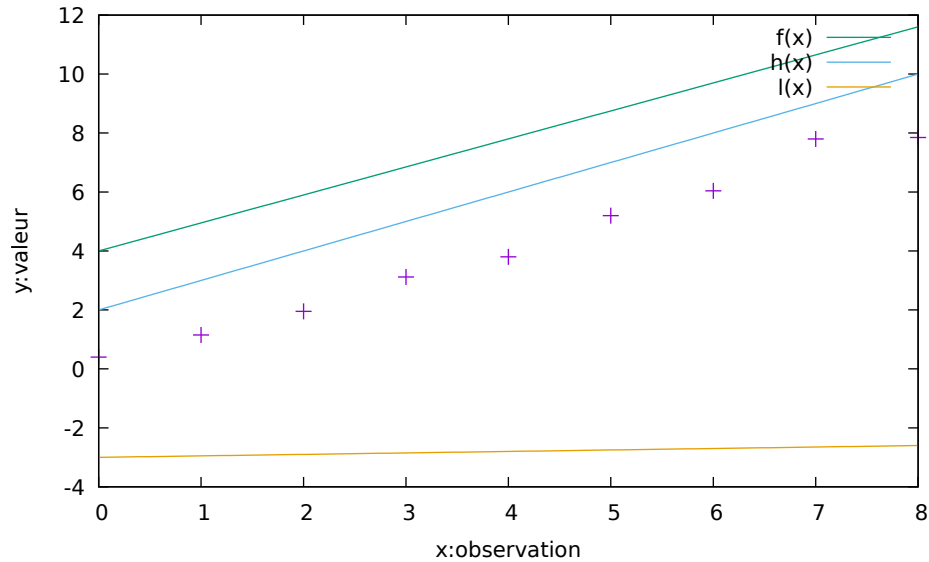


FIGURE 2.4 – Plusieurs fonction de regression peu satisfaisantes.

Il faut trouver la fonction F de manière algébrique :

$$F : y = a \cdot \vec{x} + b \quad y, b \in \mathbb{R}, \vec{x} \in \mathbb{R}^n, a \cdot \vec{x} \in \mathbb{R}$$

On remarque que $a \cdot \vec{x} \in \mathbb{R}$, on en déduit que $a \in \mathbb{R}^n$ et que il y a un produit scalaire entre a et \vec{x} . Il y a cependant un³ différence entre a et \vec{x} pour que ce produit soit possible : a est un vecteur ligne et \vec{x} un vecteur ligne.

$$a = (a_0 \quad . \quad . \quad . \quad a_{n-1}) \quad \vec{x} = \begin{pmatrix} x_0 \\ . \\ . \\ . \\ x_{n-1} \end{pmatrix}$$

a est en fait la norme de la droite représentant F , la norme étant un vecteur qui est perpendiculaire à une droite. De ce fait a est un vecteur colonne transposé et on doit noter l'équation de la droite de la manière suivante :

$$Y : y = a^T \cdot \vec{x} + b$$

On verra à quoi correspond a^T par la suite.

³. on verra en fait que a est un vecteur colonne transposé.

2.2.3 Petits rappels sur les vecteurs

Produit scalaire

On définit le produit scalaire de la manière suivante :

$$a.\vec{x} = \sum_{i=0}^n a_i.x_i$$

Transposition

La transposition d'un vecteur ligne en fait un vecteur colonne et vice versa. On note v^T la transposé du vecteur v . Ainsi :

$$\vec{c}^T = (c_0 \quad . \quad . \quad . \quad c_{n-1})$$

$$l^T = \begin{pmatrix} a_0 \\ . \\ . \\ . \\ l_{n-1} \end{pmatrix}$$

c un vecteur colonne et l un vecteur ligne.

On pourra remarquer une propriété intéressante :

$$v^{T^T} = (v^T)^T = v$$

Du coup on peut écrire :

$$(a^T.\vec{x})^T = \vec{x}^T.a^{T^T} = \vec{x}^T.a$$

On rappelle aussi la chose suivante :

$$v.v^T = \|v\|^2$$

Le produit scalaire d'un vecteur et de sa transposé est égal à la norme⁴ de ce vecteur au carré.

4. longueur

2.2.4 Forme général de F

On donne comme forme générale de F :

$$F : \Theta_1 \cdot \vec{x} + \Theta_0$$

On va chercher à trouver Θ_1 et Θ_0 sachant que :

- Θ_0 est la valeur à l'origine de F.
- Θ_1 est le vecteur normal à la droite F.

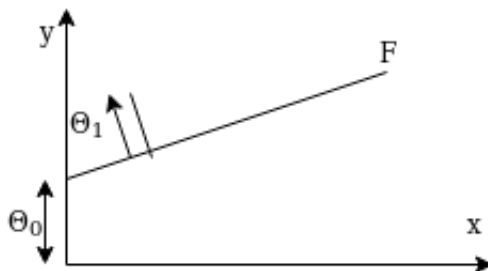


FIGURE 2.5 – Θ_0 et Θ_1 par rapport à F

Pour se faciliter la tâche on va réunir Θ_0 et Θ_1 dans un vecteur $\Theta \cdot \vec{x}'$.

$$F : \Theta \cdot \vec{x}'$$

Comment passer de x à x' ?

Prenons des valeurs en exemple :

$$x_1 = 3.3 \rightarrow x'_1 = \begin{pmatrix} 1 \\ 3.3 \end{pmatrix}$$

$$x_2 = 4.4 \rightarrow x'_2 = \begin{pmatrix} 1 \\ 4.4 \end{pmatrix}$$

$$x_n = 7.4 \rightarrow x'_n = \begin{pmatrix} 1 \\ 7.4 \end{pmatrix}$$

On a donc :

$$x' = \begin{pmatrix} 1 \\ x \end{pmatrix}$$

et

$$\Theta = (\Theta_0 \quad \Theta_1)$$

Pour calculer/apprendre Θ il est crucial de **prendre en compte l'observation**, afin de ne pas se retrouver dans la situation de la figure 2.4.

2.2.5 Fonction $J(\Theta)$

Prendre en compte les observations

On sait que $F(x_i) \rightarrow y_i$ est une fonction prédisant les valeurs de y . Pour prendre en compte les observation on va définir la fonction suivante :

$$J(\Theta) = \sum_m^{i=1} (F(x_i) - y_i)^2$$

que l'on pourra aussi écrire :

$$J(\Theta) = \sum_m^{i=1} (\Theta x'_i - y_i)^2$$

On notera que si $x \in \mathbb{R}^n$ alors $\Theta \in \mathbb{R}^{n+1}$

Analysons cette fonction :

- m correspond au nombre d'échantillons observés.
- Les valeurs que l'on cherche doivent être positives, d'où l'élevation au carré.
- $F(x_i)$ correspond aux valeurs prédites.
- y_i correspond aux valeurs observées.

Il va falloir travailler sur cette fonction $J(\Theta)$

Minimisation de $J(\Theta)$

Reprenons le processus de régression depuis le début :

1. On part d'un observation : $\{(features, classe)\}$
2. On connaît la forme de la fonction droite.
3. $\Theta_0, \Theta_1 \rightarrow \Theta$ ($\Theta_0 \quad \Theta_1$) et $x \rightarrow x' \begin{pmatrix} 1 \\ x \end{pmatrix}$.
4. On définit la fonction que l'on cherche minimiser :

$$J(\Theta) = \sum_m^{i=1} (\Theta x'_i - y_i)^2$$

Minimiser une fonction revient à trouver le minimum de la fonction. Donc minimiser une fonction $f(x)$ revient à résoudre $f(x)=0$.

Pour $J(\Theta)$ on trouve une solution unique quand $m = n + 1$, mais la plupart du temps il va falloir approximer la solution.

2.2.6 Dérivée partielle et Gradient

Dérivée partielle

Pour trouver le minimum il va falloir utiliser la notion de **dérivée partielle de fonction**. Par exemple la dérivée partielle de $X.x_0 + Y.x_1$ en x_0 puis en x_1 .

$$\frac{\partial(X.x_0 + Y.x_1)}{\partial x_0} = X \quad \frac{\partial(X.x_0 + Y.x_1)}{\partial x_1} = Y$$

Pour notre fonction $J(\Theta)$:

$$\frac{\partial J(\Theta)}{\partial \Theta_0} \quad \frac{\partial J(\Theta)}{\partial \Theta_1}$$

Gradient

On appelle gradient d'un fonction $f(x) = X.x_0 + Y.x_1$ la chose suivante :

$$\vec{\nabla} f(x) = \begin{pmatrix} \frac{\partial(X.x_0 + Y.x_1)}{\partial x_0} \\ \frac{\partial(X.x_0 + Y.x_1)}{\partial x_1} \end{pmatrix}$$

Retranscrit à la fonction $J(\Theta)$ on obtient :

$$\vec{\nabla} J(\Theta) = \begin{pmatrix} \frac{\partial J(\Theta)}{\partial \Theta_0} \\ \frac{\partial J(\Theta)}{\partial \Theta_1} \end{pmatrix}$$

2.2.7 Que cherche-t-on à résoudre ?

Le but de ce calcul est de trouver 1 ou plusieurs Θ tel que $\vec{\nabla} J(\Theta) = \vec{0}$.
Il y a deux solution :

La solution analytique

Il faut donc résoudre le système suivant :

$$\frac{\partial J(\Theta)}{\partial \Theta_0} = 0 \quad \frac{\partial J(\Theta)}{\partial \Theta_1} = 0$$

Cependant réussir à calculer la dérivée de cette fonction est compliquer rendant cette solution souvent impossible à appliquer.

La solution algorithmique la descente de gradient

$$\frac{\partial J(\Theta)}{\partial \Theta_0} \quad \frac{\partial J(\Theta)}{\partial \Theta_1}$$

Il faut de prime abord comprendre que cette méthode ne fonctionne pas toujours pour plusieurs raisons que nous allons identifier.

Le principe de cet algorithme est le suivant : **Aller dans le sens de la pente.**

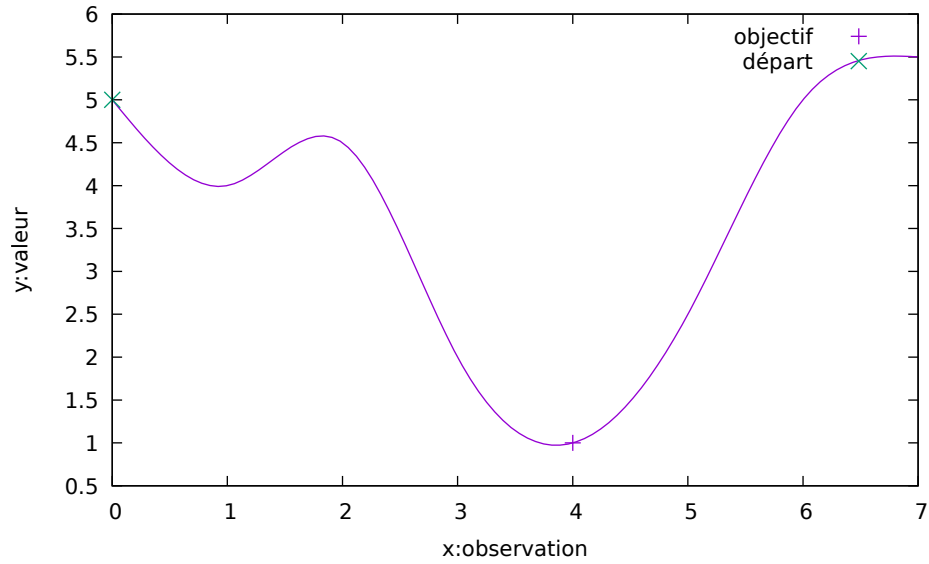


FIGURE 2.6 – Principe de la descente de gradient

L'aide de la dérivé de la fonction, on va aller vers la gauche si f est décroissante ou aller vers la droite si f est croissante.

On se déplace dans le sens inverse du signe de la dérivée. On peut tout de suite remarquer qu'il peut y avoir des "faux" minimums et que si une fonction est convexe sur un intervalle on peut se retrouver à osciller indéfiniment.

2.2.8 Descente de gradient

Trouver la formule pour la solution

La solution se calcul ainsi : Solution à $n+1$ = solution à n - $\vec{\nabla} J()$ à l'étape n . On peut donc donner cette formule :

$$\Theta_{n+1} = \Theta_n - \vec{\nabla} J(\Theta_n)$$

On ne s'arrête que lorsque $\vec{\nabla} J(\Theta) \approx \vec{0}$, c'est à dire tant que $\|J(\Theta)^{n+1}\| \geq \epsilon$. Cependant il y a de nombreux problèmes :

Faux minimum

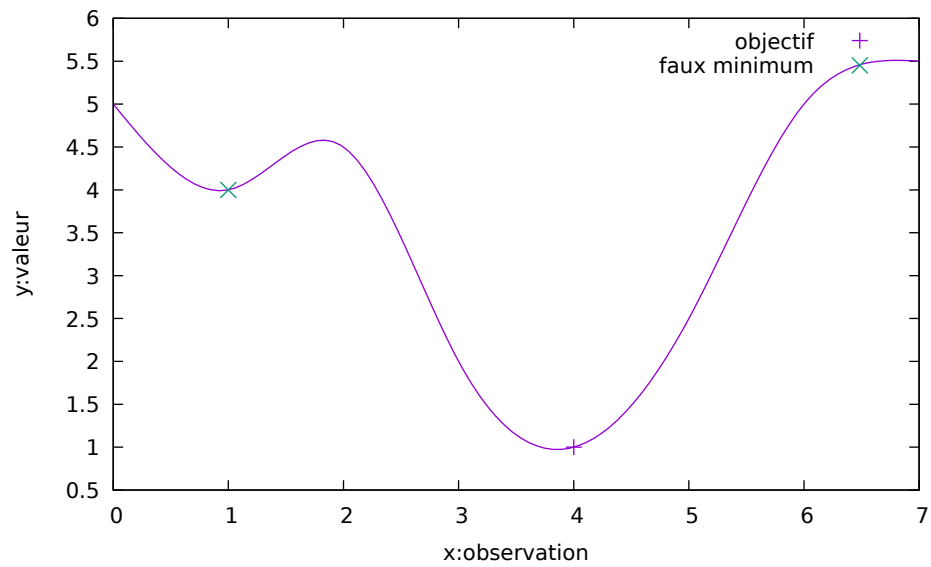


FIGURE 2.7 – Exemple de faux minimum

Le minimum réel n'est jamais atteint.

Arrêt de descente

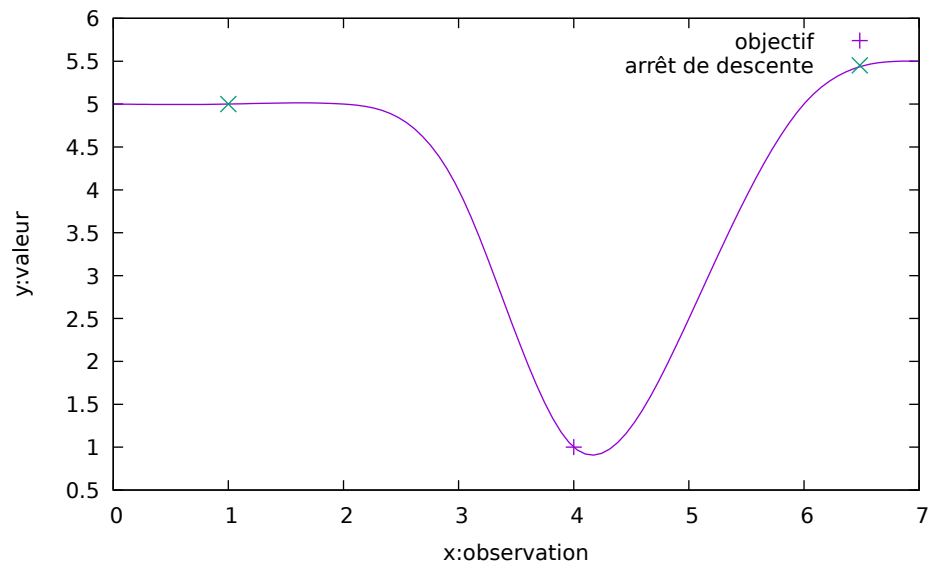


FIGURE 2.8 – Exemple d'arrêt de descente

Le minimum n'est pas garanti si la courbe est localement constante.

Pente trop douce

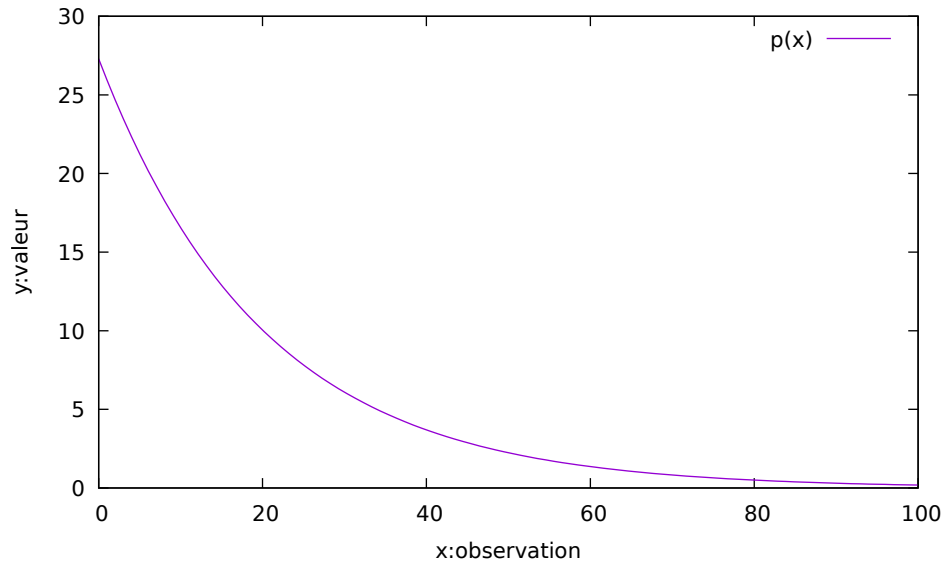


FIGURE 2.9 – Exemple de pente trop douce

Ici le problème vient du fait que plus on avance sur la pente plus $\vec{\nabla} J(\Theta)$ est trop petit, la descente s'arrête à quand ϵ sera atteint et ce bien avant le réel résultat.

Solution

Afin de résoudre ce problème on fait appel à **un coefficient d'amortissement** noté $\eta(n)$.

$\eta(n) = 1$ Descente constante

$\eta(n) > 1$ Descente accélérée

$\eta(n) < 1$ Descente amortie

$\eta(n) = 0$ Pas de descente

Ainsi on cherche le **Gradient amorti** :

$$\Theta_{n+1} = \Theta_n - \eta(n) \cdot \vec{\nabla} J(\Theta_n)$$

Remarques

Θ_0 est initialisé aléatoirement ou de manière choisie Le calcul de gradient est un calcul long, on va donc vouloir se concentrer dès le départ sur une petite portion choisie du gradient.

Chapitre 3

Descente de gradient

3.1 Rappels

3.1.1 Processus itératif

La descente de gradient s'effectue sur une fonction $J(\theta)$, une fonction objective permettant de prédire une valeur en fonction des descripteurs fournis. On va donc chercher à trouver θ tel que $J(\theta)$ soit **minimal**.

Pour cela on utilise le procédé itératif suivant :

Tant que $(\theta_{k+1} - \theta_k) \geq \epsilon$

$$\Theta_{k+1} = \Theta_k - \eta(k) \cdot \vec{\nabla} J(\Theta_k)$$

On rappelle que $\eta(k)$ est un coefficient d'amortissement/d'accélération(en fonction de sa valeur) et que $\vec{\nabla} J(\theta)$ est le gradient de la fonction $J(\theta)$ tel que :

$$\vec{\nabla} J(\theta) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_0} \\ \dots \\ \frac{\partial J(\theta)}{\partial \theta_i} \end{pmatrix}$$

La descente va dans le sens inverse du signe de la dérivée.

3.1.2 Regression linéaire

On réalise une régression à partir des observations :

observation = (x_i, y_i) , x_i est un descripteur et y_i la valeur observée *Exemple* :
 $((3, 2), 4)$, $x = (3, 2)$ et $y = 4$.

On peut définir une régression linéaire et ainsi définir la forme générale de la fonction de régression $h(x)$:

$$h(x) = \theta x + \theta_0$$

θ est de dimensions \mathbb{R}^n .

Remarque sur les hyperplans

La définition d'un hyperplan est la suivante :

Un espace vectoriel de dimension $n - 1$ si la dimension est de n

En clair si on est dans espace de dimension :

- 2 (x, y) : alors l'hyperplan sera une droite
- 3 (x, y, z) : alors l'hyperplan sera un plan
- etc..

3.1.3 Regression linéaire (suite)

On cherche ensuite à simplifier la fonction $h(x)$. Pour cela on définit :

$$\theta'' = \begin{pmatrix} \theta_0 \\ \theta \end{pmatrix}$$

De ce fait on re-définit x en x'' :

$$x_i'' = (1x_i)$$

On a donc maintenant pour fonction :

$$h(x'') = \theta'' x''$$

Et on définit la fonction $J(\theta'')$:

$$\sum_{i=1}^m (\theta'' x_i'' - y_i)^2, \quad m \text{ étant le nombre d'observation et } y \text{ la valeur observée}$$

C'est à partir de cette écriture qu'on a pu définir $\vec{\nabla} J(\Theta)$:

$$\vec{\nabla} J(\Theta_k) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_0} \\ \dots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{pmatrix} \quad n \text{ étant la dimension des descripteur}$$

3.2 Dérivée de la fonction $J(\theta)$

On va maintenant voir comment calculer et noter la dérivée de la fonction $J(\theta)$:

$$\begin{aligned} \frac{\partial \sum_{i=1}^m (\theta'' x_i'' - y_i)^2}{\partial \theta_i} &= \frac{\partial \sum_{i=1}^m ((\theta_0 \quad \dots \quad \theta_m) \begin{pmatrix} 1 \\ x_{i,1} \\ \dots \\ x_{i,m} \end{pmatrix} - y_i)^2}{\partial \theta_i} \\ &= \frac{\partial \sum_{i=1}^m ((\theta_0 + \theta_1 x_{i,1} + \dots + \theta_n x_{i,n}) - y_i)^2}{\partial \theta_i} \end{aligned}$$

On notera la chose suivante :

$$x_{i,k} = \frac{\partial (\theta_0 + \dots + \theta_n x_{i,n}) - y_i}{\partial \theta_i}$$

Ainsi on peut donner la formule suivante :

$$\frac{\partial J(\theta)}{\partial \theta_k} = 2 \sum_{i=1}^m x_{i,k} (\theta'' x_i'' - y_i)$$

Comme on cherche à minimiser la formule peut être réécrite ainsi :

$$\frac{\partial J(\theta)}{\partial \theta_k} = \sum_{i=1}^m x_{i,k} (\theta'' x_i'' - y_i)$$

en effet le coefficient 2 influe peu sur le résultat. On peut donc donner $\vec{\nabla} J(\Theta_k)$ sous cette forme

$$\vec{\nabla} J(\Theta) = \begin{pmatrix} \sum_{i=1}^m x_{i,0} (\theta'' x_i'' - y_i) \\ \dots \\ \sum_{i=1}^m x_{i,k} (\theta'' x_i'' - y_i) \\ \dots \end{pmatrix} \quad \text{avec } \forall j, x_{0,j} = 0$$

Exemple

On possède les observations suivantes :

$$((2, 3), 1) \quad ((4, 5), 3)$$

On calcule le $\vec{\nabla} J(\Theta)$ à partir des ces observations :

$$\vec{\nabla} J(\Theta) = \begin{pmatrix} 1 \times \theta \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} - 1 + 1 \times \theta \begin{pmatrix} 1 \\ 4 \\ 5 \end{pmatrix} - 3 \\ 2 \times \theta \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} - 1 + 4 \times \theta \begin{pmatrix} 1 \\ 4 \\ 5 \end{pmatrix} - 3 \\ 3 \times \theta \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} - 1 + 5 \times \theta \begin{pmatrix} 1 \\ 4 \\ 5 \end{pmatrix} - 3 \end{pmatrix}$$

Maintenant que l'on sait calculer le gradient on peut appliquer le processus itératif pour trouver un minimum de $J(\theta)$

Remarque du professeur

Savoir appliquer la descente de gradient ainsi sur un exemple est demandé à l'examen

3.3 Reformulation sous forme matricielle

Exemple de mise en forme

Soit l'observation suivante :

$$Obs = \{((2, 3), 1), ((4, 5), 3), ((5, 4), 8)\}$$

3.3.1 Forme Matricielle de $J(\theta)$

On définira les matrices et vecteurs suivants :

$$X = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 4 & 5 \\ 1 & 5 & 4 \end{bmatrix} \quad \vec{y} = \begin{bmatrix} 1 \\ 3 \\ 8 \end{bmatrix} \quad \Theta'' = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

On peut faire le lien suivant :

$$J(\theta'') = \sum_{i=1}^m (\theta'' x_i'' - y_i)^2 = (X\Theta'' - \vec{y})(X\Theta'' - \vec{y})^T$$

$(X\Theta'' - \vec{y})$ est un vecteur colonne

$(X\Theta'' - \vec{y})^T$ est un vecteur ligne

$J(\theta'')$ est donc un scalaire ($\in \mathbb{R}$)

3.3.2 Dérivé de la forme matricielle

On va très rapidement donner les étapes de la dérivation de la forme matricielle :

$$\begin{aligned}\frac{\partial J(\Theta'')}{\partial \Theta''} &= (X\Theta'' - \vec{y})(X\Theta'' - \vec{y})^T \\ &= \frac{\partial (X\Theta'')^T X\Theta'' - (X\Theta'')^T \vec{y} - \vec{y}^T X\Theta'' + \vec{y}^T \vec{y}}{\partial \Theta''} \\ &= \frac{\partial (X\Theta'')^T X\Theta'' - 2\vec{y}^T (X\Theta'')}{\partial \Theta''} \\ \frac{\partial J(\Theta'')}{\partial \Theta''} &= XX^T \Theta'' - X^T \vec{y}\end{aligned}$$

Ainsi lorsque l'on veut résoudre

$$\frac{\partial J(\Theta'')}{\partial \Theta''} = 0$$

Il suffit de résoudre :

$$XX^T \Theta = X^T \vec{y}$$

Remarque : le résultat d'un produit matriciel entre une matrice et sa transposé est une matrice carrée. XX^T est **une matrice carrée** Si la matrice XX^T est inversible alors on peut calculer $(XX^T)^{-1}$ et donc résoudre l'équation suivante¹ :

$$\Theta'' = (XX^T)^{-1} X^T \vec{y}$$

3.3.3 Problèmes de la forme matricielle

Est-ce qu'il y a forcément une solution à l'équation suivante ? :

$$X\Theta'' - \vec{y} = \vec{0} \quad (X \in (\mathbb{R}^n \times \mathbb{R}^m), \Theta'' \in \mathbb{R}^n, \vec{y} \in \mathbb{R}^m)$$

On note que qu'il y aura 10^6 avec seulement 3 inconnues, et qu'une inversion de matrice a une complexité en temps de $O(n^3)$.

La réponse est non. La solution va uniquement dépendre de l'ordre de grandeur des informations prises. Une solution est de choisir des sous-échantillons.

3.3.4 À retenir

- D'une observation on obtient θ .
- Méthode directe : $\Theta = (XX^T)^{-1} X^T \vec{y}$
- Méthode indirecte : descente de gradient par processus itératif.

1. on peut obtenir ce résultat sur matlab en écrivant $\Theta'' = y \backslash X$

Chapitre 4

Bayesienne

4.0.1 Principe

On connaît la forme de la distribution. On ne **connait pas les paramètres** mais on connaît la **la loi de probabilité**

Exemple

Problèmes à deux classes : "Bar" et "Saumon".
On possède la vérité terrain suivante :

$$VT = \begin{pmatrix} (Feature_1, "Bar") \\ (Feature_2, "Saumon") \\ \dots \\ (Feature_n, "Bar") \end{pmatrix}$$

On doit prendre une décision en fonction de la Vérité Terrain et avec l'aide d'une loi de distribution de \underline{X} (une loi de probabilité), en fonction de "Bar" ou "Saumon"

4.1 Différentes approches

4.1.1 Maximum à priori

La bonne démarche est la suivante :

$$VT = \begin{pmatrix} (\emptyset, "w_1") \\ \dots \\ (\emptyset, "w_2") \end{pmatrix}$$

On connaît le nombre de "Bar": n_{w_1} et le nombre de "Saumon" n_{w_2}
On donne les probabilités suivante :

$$P(w_1) = \frac{n_{w_1}}{|VT|}$$

$$P(w_2) = \frac{n_{w_2}}{|VT|}$$

Pour vérifier que il s'agisse bien d'une loi de probabilité on s'assure que :

$$\frac{n_{w_1} + n_{w_2}}{|VT|} = 1$$

On va donc pouvoir dire si $P(w_1) > P(w_2)$ alors " w_1 ", sinon " w_2 ".
Ce classifieur ne prend pas en compte les features.

4.1.2 Maximum vraisemblable

Pour ce classifieur on va étudier les features en fonction de l'étiquette de la classe.

Notre Vérité Terrain devient :

$$VT = \begin{pmatrix} (f_1, "w_1") \\ \dots \\ (f_n, "w_2") \end{pmatrix}$$

Il nous faudra donc définir : $P(x = f|w)$ la probabilité que f est la valeur x sachant qu'il appartient à la classe w .

Pour x donné il faut décider s'il appartient à " w_1 " ou " w_2 " :

$$P(x = f|w_1)$$

$$P(x = f|w_2)$$

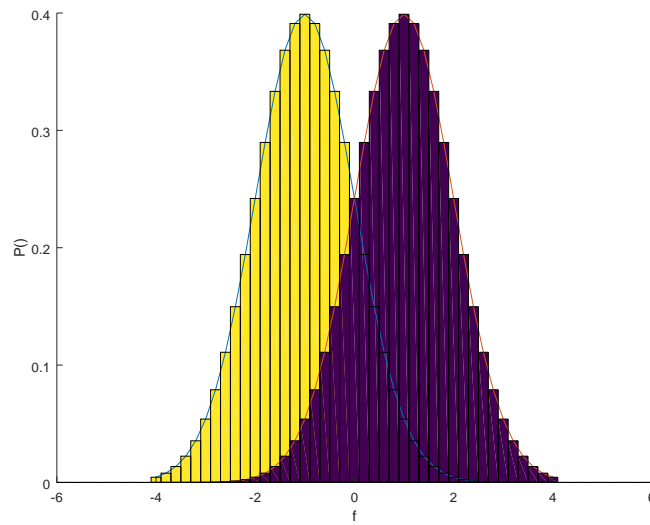


FIGURE 4.1 – Maximum vraisemblable

On passe de données discrètes à des données continue suivant une certaine loi de probabilité. Ici on suit une loi normale $\mathcal{N}(\mu_1, \sigma_1)$ et $\mathcal{N}(\mu_2, \sigma_2)$

Récapitulatif 1

1. On part d'une Vérité Terrain et on identifie 2 classes : w_1, w_2
2. On identifie $(\mu_1, \sigma_1), (\mu_2, \sigma_2)$
3. On calcule $P(x = f|w_1)$ et $P(x = f|w_2)$

$$P(x = f|w) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$P(x = f|w_1) = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2}(\frac{x-\mu_1}{\sigma_1})^2}$$

$$P(x = f|w_2) = \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2}(\frac{x-\mu_2}{\sigma_2})^2}$$

Fonctionnement

Ce classifieur prend x et répond w_1 ou w_2 .
On va donc comparer $P(x = f|w_1)$ $P(x = f|w_2)$:

Si $P(x = f|w_1) > P(x = f|w_2)$ alors w_1 sinon w_2

Probabilité d'erreur et seuil d'erreur

$$P(\text{erreur}) = \max[P(x = f|w_1), P(x = f|w_2)], \int_{f_1}^{f_n} P(\text{erreur})dx = 1$$

Pour trouver le seuil d'erreur on va résoudre l'équation suivante :

$$P(x = f|w_1) = P(x = f|w_2)$$

$$\frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2}(\frac{x-\mu_1}{\sigma_1})^2} = \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2}(\frac{x-\mu_2}{\sigma_2})^2}$$

$$\log\left(\frac{1}{\sigma_1} e^{-\frac{1}{2}(\frac{x-\mu_1}{\sigma_1})^2}\right) = \log\left(\frac{1}{\sigma_2} e^{-\frac{1}{2}(\frac{x-\mu_2}{\sigma_2})^2}\right)$$

$$-\frac{1}{2}\left(\frac{x-\mu_1}{\sigma_1}\right)^2 + \frac{1}{2}\left(\frac{x-\mu_2}{\sigma_2}\right)^2 = -\log \sigma_1 + \log \sigma_2$$

L'erreur correspondra à la surface créée par le croisement des deux courbes gaussiennes.

On obtient une équation du second degré qui est résoluble assez facilement.
Ensuite on pourra réaliser la comparaison suivante :

Si $x > Z$ alors w_1 sinon w_2

On notera qu'il y a plusieurs seuils comparables, un par descripteur. On devra donc réaliser la suite de comparaison suivante :

si $x < z_1$ alors w_1
sinon si $x < z_2$ alors w_2
...
sinon w_1

Récapitulatif 1

1. On part d'une Vérité Terrain et on identifie 2 classes : w_1, w_2
2. On identifie $(\mu_1, \sigma_1), (\mu_2, \sigma_2)$
3. On calcule $P(x = f|w_1)$ et $P(x = f|w_2)$

Récapitulatif 2

1. On part d'une Vérité Terrain et on identifie 2 classes : w_1, w_2
2. On identifie $(\mu_1, \sigma_1), (\mu_2, \sigma_2)$
3. On calcule $P(x = f|w_1)$ et $P(x = f|w_2)$
4. On résout l'équation $P(x = f|w_1) = P(x = f|w_2)$
5. On fait une comparaison sur chaque z_i de chaque descripteur ;

4.1.3 Maximum à posteriori

Pour ce type de classifieur on utilisera la formule de Bayes :

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Dans le cadre du classifieur on va chercher à donner la formule $P(w|x = f)$;

$$P(w|x = f) = \frac{P(x = f|w)P(w)}{P(x = f)}$$

Si l'on reprend le précédent classifieur, on se sert de la distribution des features dans la classe.

Ici, on se sert des features pour trouver la classe. On pourra réaliser la comparaison suivante :

$$\frac{P(w_1|f = x)P(w_1)}{P(f = x)} > \frac{P(w_2|f = x)P(w_2)}{P(f = x)}$$

$P(f = x)$ n'influençant pas la comparaison, on peut la simplifier .

$$P(f = x|w_1)P(w_1) > P(f = x|w_2)P(w_2)$$

Ainsi :

Si $P(f = x|w_1)P(w_1) > P(f = x|w_2)P(w_2)$ alors w_1 sinon w_2

4.2 Loi normal en dimension N

On s'intéresse désormais aux descripteurs de dimension \mathbb{R}^n

4.2.1 Fonction gaussienne dans \mathbb{R}^n

$$P(x) = \frac{1}{(2\pi)^{1/2}|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

Avec :

- Σ une matrice de covariance $\in \mathbb{R}^n$
- Σ^{-1} , l'inverse Σ
- μ le vecteur moyen de l'échantillon $\in \mathbb{R}^n$
- $|\Sigma|$ le déterminant de la matrice $\Sigma, \in \mathbb{R}$

4.2.2 Matrice de covariance

La matrice de covariance est une matrice carrée, symétrique, positive ayant la forme suivante :

$$\begin{bmatrix} \sigma_{1,1}^2 & \sigma_{1,2} & \cdots & \sigma_{1,n} \\ \sigma_{2,1} & \sigma_{2,2}^2 & & \\ \cdots & & \cdots & \\ \sigma_{n,1} & & & \sigma_{n,n}^2 \end{bmatrix}$$

On a les formules suivantes¹ :

$$\sigma_{i,i}^2 = \frac{1}{k} \sum_{l=1}^{l=k} (x_{l,i} - \mu_i)^2$$

et :

$$\sigma_{i,j}^2 = \frac{1}{k} \sum_{l=1}^{l=k} (x_{l,i} - \mu_i)(x_{l,j} - \mu_j)$$

On a aussi :

$x_{l,i}$, l : le l -ième descripteur, i : le i -ième élément du descripteur mail

$$\mu = \begin{bmatrix} \mu_1 \\ \cdots \\ \mu_j \\ \cdots \\ \mu_n \end{bmatrix} = \begin{bmatrix} \frac{1}{n} \sum_1^l x_{l,1} \\ \cdots \\ \frac{1}{n} \sum_1^l x_{l,j} \\ \cdots \\ \frac{1}{n} \sum_1^l x_{l,n} \end{bmatrix}$$

D'ailleurs on notera que $\sigma_{i,j} = \sigma_{j,i}$.

Exemple

On possède les descripteurs suivants :

$$(1, 4), (2, 5), (3, 9)$$

On aura la matrice suivante :

$$\begin{bmatrix} \sigma_{1,1}^2 & \sigma_{1,2} \\ \sigma_{2,1} & \sigma_{2,2}^2 \end{bmatrix}$$

On calcule donc les valeurs suivantes :

$$\mu_1 = \frac{1 + 2 + 3}{3} = 2$$

$$\mu_2 = \frac{4 + 5 + 9}{3} = 6$$

1. $\sigma_{i,i}^2$ est une notation pour signifier que ce σ précis est sur la diagonale.

$$\begin{aligned}\sigma_{1,1}^2 &= \frac{(1-2)^2 + (2-2)^2 + (3-2)^2}{3} = \frac{2}{3} \\ \sigma_{2,2}^2 &= \frac{(4-6)^2 + (5-6)^2 + (9-6)^2}{3} = \frac{14}{3} \\ \sigma_{1,2}^2 &= \frac{(1-2)(4-6) + (2-2)(5-6) + (3-2)(9-6)}{3} = \frac{5}{3}\end{aligned}$$

Ainsi on obtient :

$$\begin{bmatrix} \frac{2}{3} & \frac{5}{3} \\ \frac{5}{3} & \frac{14}{3} \end{bmatrix}$$

4.3 $\mathcal{N}(\mu, \sigma)$ en dimension 2

4.3.1 Distance

Si l'on regarde la fonction gaussienne on voit apparaître :

$$\frac{x-\mu}{\sigma}, \text{ une formule de distance}$$

En dimension 2 on peut considérer les identités suivantes :

$$\begin{aligned}\Sigma &= \sigma^2 \\ \Sigma^{-1} &= \frac{1}{\sigma^2}\end{aligned}$$

Formule de la distance de Mahalanobis

On donne la formule suivante :

$$\sqrt{(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

Cette formule permet de remplacer la distance euclidienne car elle prend en compte la matrice Σ^{-1}

4.3.2 Courbe d'iso probabilité

Il est possible de tracer des courbes d'isoprobabilité² à partir des matrices de covariance Σ . On distinguera plusieurs cas particuliers :

2. ellipses concentriques. Chaque ellipse représente une probabilité et tout point appartenant à la même ellipse ont la même probabilité.

Courbes d'iso-probabilité : cercles

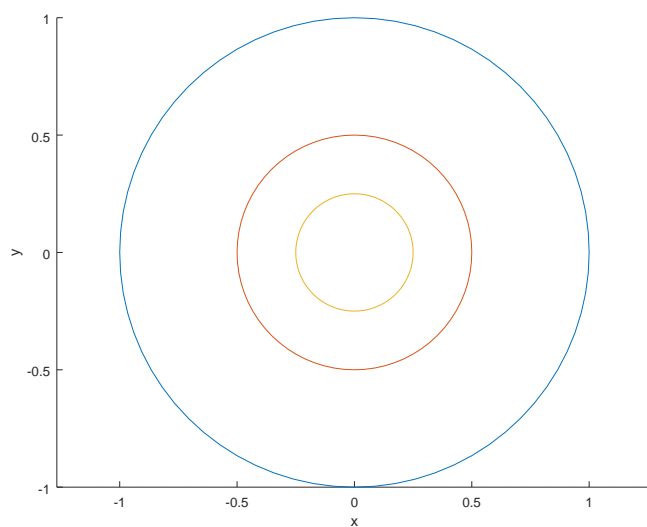


FIGURE 4.2 – Courbes d'isprobabilités

Dans ce cas particulier la matrice Σ est de la forme suivante :

$$\Sigma = \begin{bmatrix} \alpha & 0 \\ 0 & \alpha \end{bmatrix} = \alpha * Identity$$

On remarquera que si Σ est inversible alors :

$$\Sigma^{-1} = \frac{1}{\alpha} * Identity$$

Courbes d'iso-probabilité : ellipses parallèles aux axes

Dans ce cas la matrice Σ est la suivante :

$$\Sigma = \begin{bmatrix} \alpha & 0 \\ 0 & \beta \end{bmatrix}$$

La distinction va se jouer entre α et β .

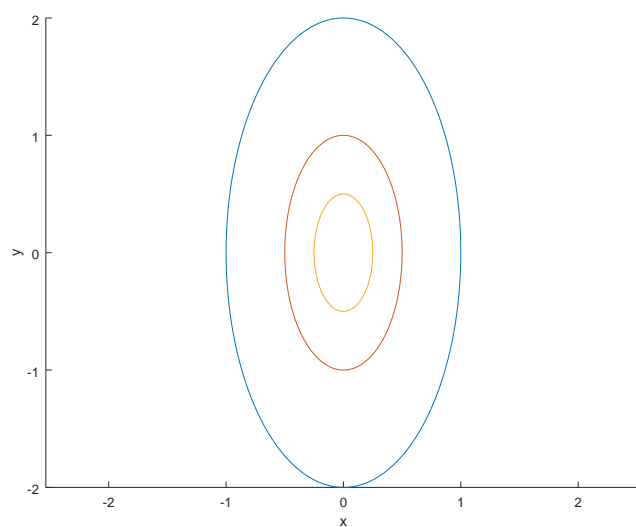


FIGURE 4.3 – Cas ou $\alpha < \beta$

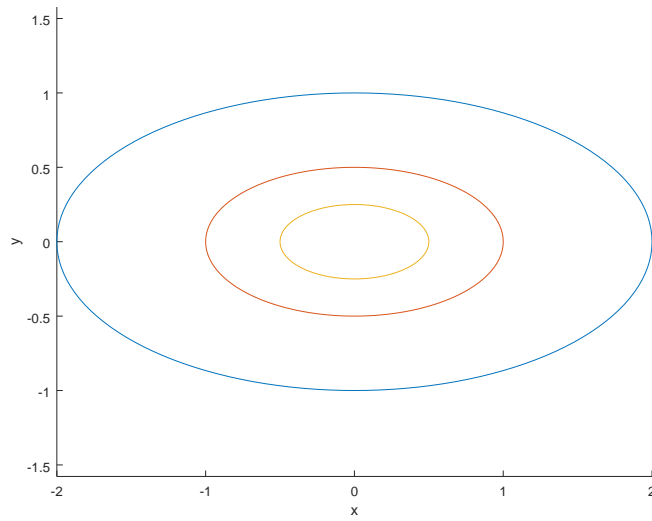


FIGURE 4.4 – Cas ou $\beta < \alpha$

Coubes d'iso-probabilité : ellipses non parallèles

Dans ce cas la matrice Σ sera similaire au cas précédent, si ce n'est que l'on introduit une variable γ :

$$\Sigma = \begin{bmatrix} \alpha & \gamma \\ \gamma & \beta \end{bmatrix}$$

Ainsi l'orientation ne dépendra plus que de α et β mais aussi de la valeur de γ .

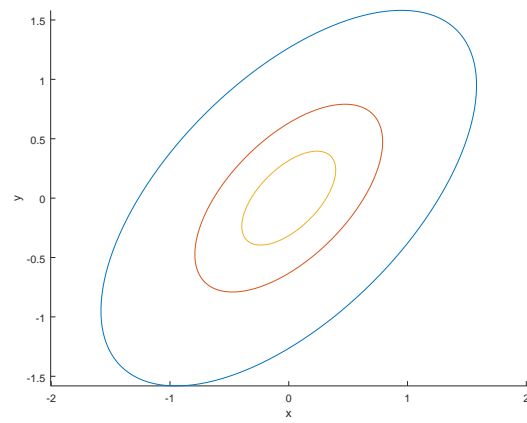


FIGURE 4.5 – $\gamma > 0$ et $\beta > \alpha$

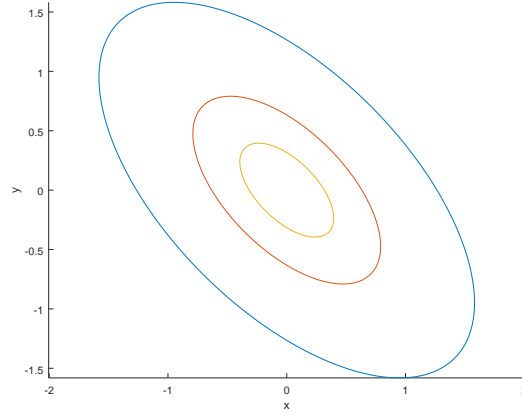


FIGURE 4.6 – $\gamma < 0$ et $\alpha > \beta$

4.3.3 Forme de la frontière de séparation

Si l'on a plusieurs observation nous obtiendrons plusieurs courbes d'iso-probabilités liées. Il va donc falloir pouvoir délimiter ces courbes entre elles.

Cas $\Sigma_1 = \Sigma_2 = \alpha * Id$

Dans le cas où $\Sigma_1 = \Sigma_2 = \alpha * Identity$, deux lois normales de matrices de covariance $\Sigma = \alpha * Identity$, on cherche x tel que $P_1(x) = P_2(x)$. Les courbes d'iso-probabilités seront 2 cercles.

On sait que :

- $P(x|w_1)P(w_1)$ est une loi normale, multi-varié, avec $\Sigma = \alpha * Id$ et $\mu = \mu_1$
- $P(x|w_2)P(w_2)$ est une loi normal, multi-varié, avec $\Sigma = \alpha * Id$ et $\mu = \mu_2$

On souhaiterait résoudre :

$$P(x|w_1)P(w_1) = P(x|w_2)P(w_2)$$

Avec quelques calculs on peut arriver à une fonction de la forme :

$$w^T x + x_0 = 0$$

Avec :

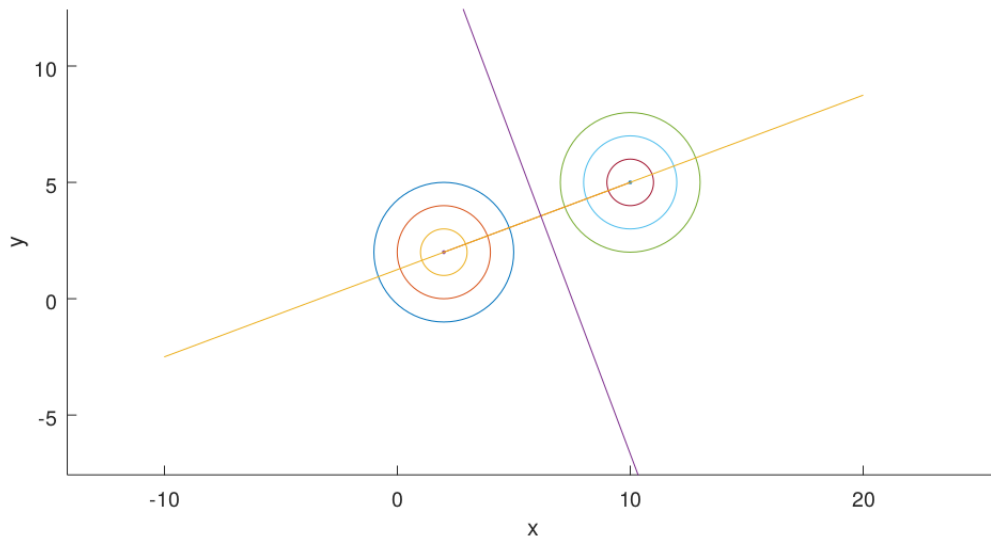
$$x_0 = -\log\left(\frac{P(w_2)}{P(w_1)}\right) + \frac{1}{2}\mu_1^T \mu_1 - \mu_2^T \mu_2$$

$$w^T = -\frac{1}{\alpha}(\mu_2^T - \mu_1^T)$$

On notera aussi :

$$w = -\frac{1}{\alpha}(\mu_2 - \mu_1)$$

La fonction est une droite orthogonale à la droite passant par le milieu (des cercles).

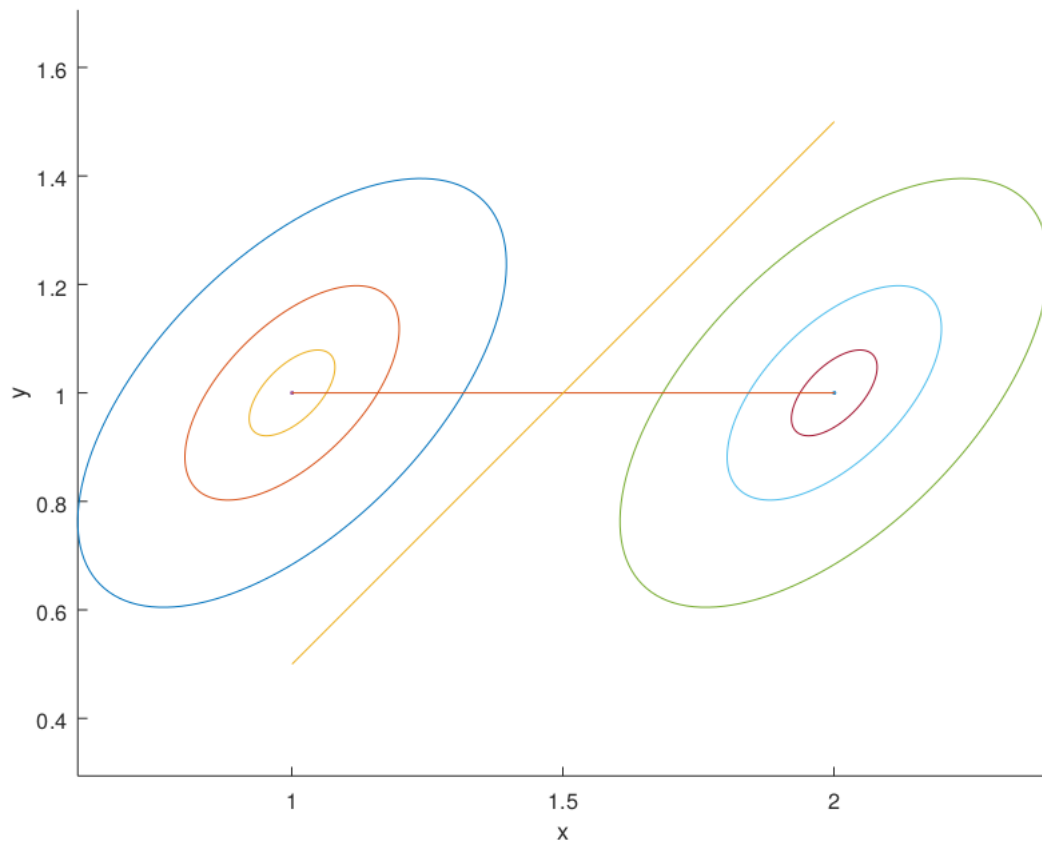


Cas $\Sigma_1 = \Sigma_2 \neq \alpha * Id$

Dans ce cas précis, les courbes sont deux ellipses parallèles par rapport au grand axe. La droite de séparation ne sera plus orthogonale à la droite passant par les milieux. La droite de séparation sera parallèle au grands axes des ellipses. On fait à nous ici face à une notion de distance (celle de Mahalanobis) :

$$(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)$$

$$(x - \mu_2)^T \Sigma^{-1} (x - \mu_2)$$



Cas $\Sigma_1 \neq \Sigma_2$

Dans ce cas précis, la fonction de la droite de séparation sera une fonction quadratique de la forme suivante :

$$x^T w x + w^T x + x_0$$

Remarque importante

Il est important de pouvoir expliquer l'origine de la forme de ces droites de séparation.

Chapitre 5

Réduction en dimension

5.1 Objectif

Le but de la réduction en dimension est de passer d'une dimension \mathbb{R}^d à une dimension \mathbb{R}^k tel que $k < d$.

On va donc passer faire la transition d'un vecteur \vec{x} de descripteur de dimension d à un vecteur \vec{y} de descripteur de dimension k . Pour ce faire il nous faudra trouver une matrice de transformation W .

$$\begin{bmatrix} x_{1,1} \\ \dots \\ x_{1,d} \end{bmatrix} \xrightarrow{W} \begin{bmatrix} y_{1,1} \\ \dots \\ y_{1,k} \end{bmatrix}$$

Sachant que l'opération à réaliser est la suivante :

$$\begin{bmatrix} x_{1,1} \\ \dots \\ x_{1,d} \end{bmatrix} = W \begin{bmatrix} y_{1,1} \\ \dots \\ y_{1,k} \end{bmatrix} \quad \text{On peut déterminer que } W \text{ est de la forme } W(k, d)$$

Ainsi :

$$W = \begin{bmatrix} w_{1,1} & \dots & w_{1,d} \\ \dots & \dots & \dots \\ w_{k,1} & \dots & w_{k,d} \end{bmatrix}$$

5.2 Avantages et inconvénients

La réduction en dimension va avoir plusieurs effets sur les observations qui se traduisent principalement par une compression de l'information.

5.2.1 Inconvénients

La compression d'information peut générer une perte d'information, en effet si on observe l'exemple suivant, la réduction de \mathbb{R}^2 à \mathbb{R}^1 provoque une perte notable d'information :

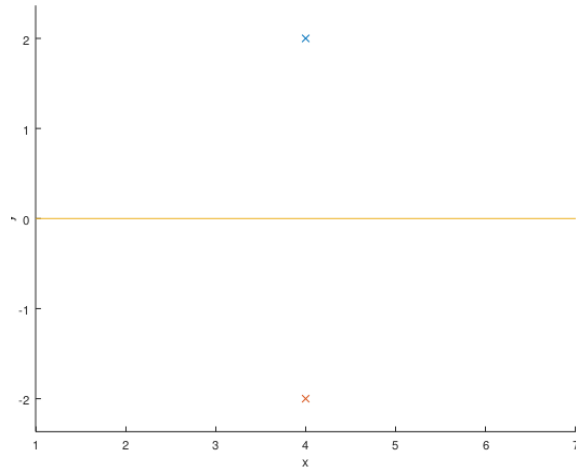


FIGURE 5.1 – Avant perte d'informations

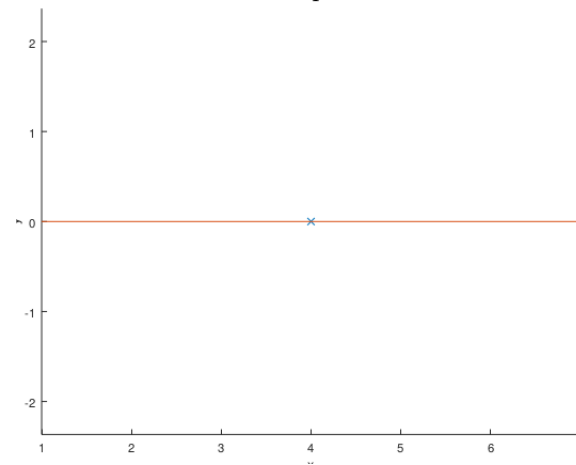


FIGURE 5.2 – Après perte d'informations

5.2.2 Avantages

La réduction va permettre de **densifier** l'information et donc provoquer une augmentation de la distance entre les échantillons (ce qui permet de trouver plus facilement une frontière). Par exemple si nous avons à faire à une observation de dimension 60 on va se poser les questions suivantes :

- Est-ce que tous ces paramètres ont la même importance ?
- Si non, quels sont les plus importants ?

Prenons dans ce cas :

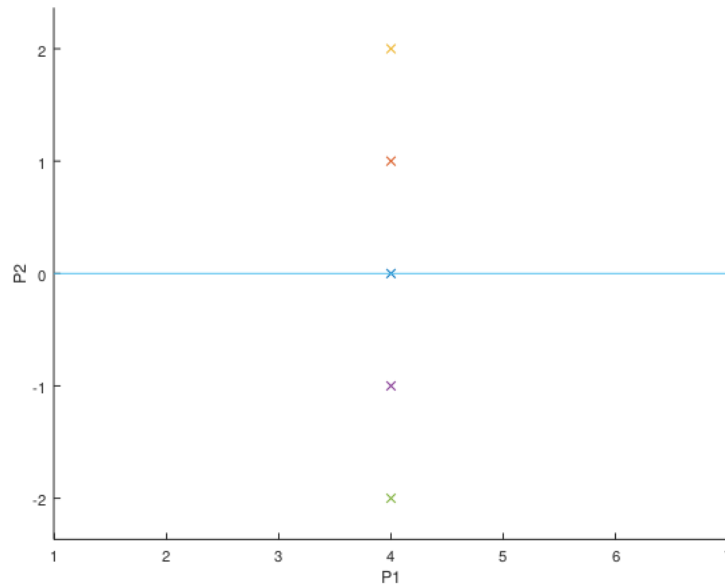


FIGURE 5.3 – Exemple de paramètres inutiles

Dans cet exemple il est clair que les paramètres P1 sont inutiles pour établir un choix : **P1 est identique pour tous les échantillons.**

Certaines méthodes ou algorithmes permettent une très grande réduction. Si nous plaçons dans \mathbb{R}^{60} avec 300 échantillons une méthode de réduction serait par exemple d'utiliser le **DeepLearning**. Cette méthode va permettre de réduire jusqu'à \mathbb{R}^3 en offrant un bon processus de sélection ¹.

5.3 Méthodes de réduction

5.3.1 Méthode ACP(Analyse en composante principale)

Il existe deux méthodes pour réaliser une réduction en dimension. La première est l'**Analyse en composante principale**. Le but est réduire l'observation de manière globale en essayant de conserver le maximum d'information. On oublierait la classe et on essaie de réduire en conservant au mieux la **variance**.

1. tout en risquant une perte d'informations.

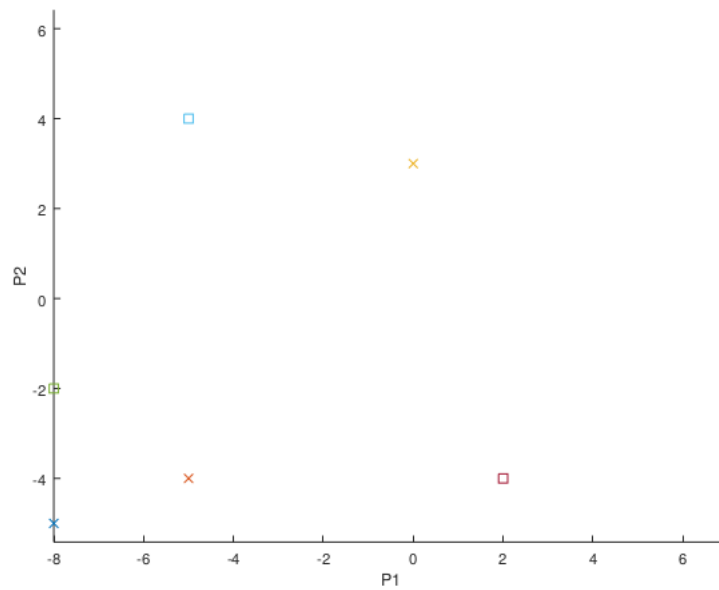


FIGURE 5.4 – Avant l'ACP

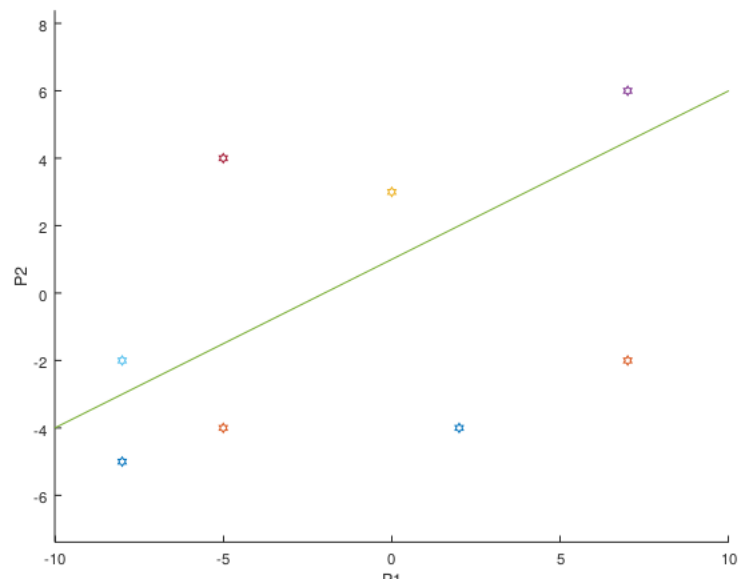


FIGURE 5.5 – Après l'ACP

5.3.2 Méthode ACI(Analyse en composantes Indépendantes)

Cette fois, la réduction va se faire en tenant compte de l'observation et en essayant de trouver un projection qui sépare au mieux les classes.

En bref, on conserve les classes et on cherche une projections des données qui les sépare au mieux.

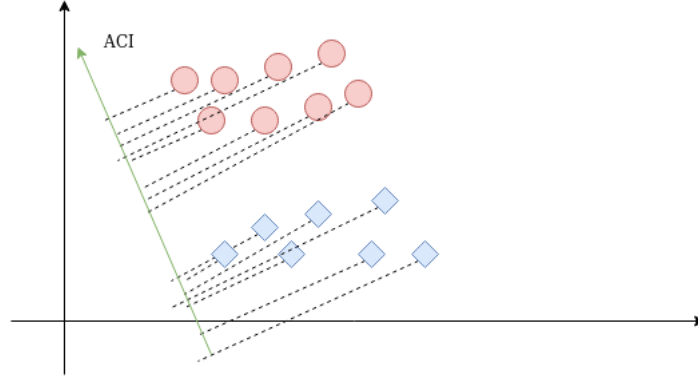


FIGURE 5.6 – Exemple d'ACI

À savoir : la méthode ACI permet une réduction de dimension entre 1 et le nombre de classe -1 . Ainsi il est possible de passer de \mathbb{R}^{60} à \mathbb{R}

5.3.3 ACP : fonctionnements et calculs

Dans cette section on verra comment on passe de $\vec{p} = \{p_1, \dots, p_d\}$, le vecteur de l'observation à $\vec{y} = \{y_1, \dots, y_k\}$ un vecteur de points projetés.

Recentrage des données

On va recentrer les points de l'observation en passant de \vec{p} à \vec{c} tel que :

$$c_i = p_i - \mu$$

Avec :

$$\mu = \frac{1}{N} \sum_{i=1}^N p_i, N = d$$

Travail à réaliser

Le but de la project peut être formalisé ainsi :

$$\vec{y} = W^T \vec{c}$$

$$\begin{bmatrix} y_1 \\ \dots \\ y_k \end{bmatrix} = \begin{bmatrix} e_{1,1} & \dots & e_{1,d} \\ \dots & \dots & \dots \\ e_{k,1} & \dots & e_{k,d} \end{bmatrix} \begin{bmatrix} c_1 \\ \dots \\ c_d \end{bmatrix}$$

On notera que l'on va passer d'une dimensions \mathbb{R}^d à une dimensions \mathbb{R}^k .

Base

Les colonnes de la matrice W , sont de la forme² :

$$\vec{e}_i = \begin{pmatrix} e_{1,i} \\ e_{2,i} \\ \dots \\ e_{d,i} \end{pmatrix}$$

Ces vecteurs forment une base, du sous-espace vectoriel de dimension k . On considérera que la base est orthonormée, ce qui donne les propriétés suivantes :

- $\forall i \neq j, \vec{e}_i \cdot \vec{e}_j = 0$
- $\forall i, \vec{e}_i \cdot \vec{e}_i = 1$

Fonction à minimiser

On donne la fonction suivante :

$$J(\vec{e}_1 \dots \vec{e}_k, \alpha_{1,1} \dots \alpha_{n,k}) = \sum_{i=1}^N \|c_i - \tilde{c}_i\|_2$$

avec :

$$\tilde{c}_i = \sum_{j=1}^k \alpha_{i,j} \vec{e}_j : \text{le projeté de } c_i \text{ dans le sous-espace défini par la base}$$

$$\|c_i - \tilde{c}_i\|_2 : \text{l'erreur commise sur un échantillon}$$

On peut finalement donner la définition de la fonction précédente :

$$J(\vec{e}_1 \dots \vec{e}_k, \alpha_{1,1} \dots \alpha_{n,k}) = \sum_{i=1}^N \|c_i - \sum_{j=1}^k \alpha_{i,j} \vec{e}_j\|_2$$

Dérivé de J sur α

On va partiellement dériver la fonction J sur $\alpha_{i,j}$ afin d'en connaître la valeur :

$$\frac{\partial J(\vec{e}_1 \dots \vec{e}_k, \alpha_{1,1} \dots \alpha_{n,k})}{\partial \alpha_{i,j}} = 0$$

On passera les calculs pour obtenir la formule suivante :

$$\alpha_{i,j} = c_i^T \vec{e}_j$$

En faisant abstraction de calculs supplémentaires on arrive à devoir résoudre le système suivant :

$$S \vec{e}_i = \lambda_i \vec{e}_i$$

2. Attention ! On parle bien de W et non pas W^T .

avec :

$$S = \sum_{i=1}^N c_i c_i^T = \sum_{i=1}^N (p_i - \mu)(p_i - \mu)^T : S \text{ est appelée Scatter Matrix}$$

Les solutions de ce système permettront de trouver une minimisation de la fonction $J(\vec{e}_1 \dots \vec{e}_k)$. Ces solutions sont les **vecteurs propres** de la matrice S . On remarquera que S est de rang d , c'est à dire qu'elle possède d **valeurs propres**.

On formera des couples (\vec{e}_i, λ_i) que l'on classera en fonction de la valeur de λ_i

Exemple avec une base $B = (\vec{e}_1)$

On va ici les c_i sur \mathbb{R} :

$$\text{Perte d'info : } \frac{1 - \lambda_1}{\sum_{i=1}^d \lambda_i}$$

$$\text{Info préservée : } \frac{\lambda_1}{\sum_{i=1}^d \lambda_i}$$

Exemple avec une base $B = (\vec{e}_1, \vec{e}_2)$

On va ici les c_i sur \mathbb{R}^2 :

$$\text{Perte d'info : } \frac{1 - \lambda_1 + \lambda_2}{\sum_{i=1}^d \lambda_i}$$

$$\text{Info préservée : } \frac{\lambda_1 + \lambda_2}{\sum_{i=1}^d \lambda_i}$$

Exemple avec une base $B = (\vec{e}_1, \dots, \vec{e}_d)$

On va ici les c_i sur \mathbb{R}^d :

Perte d'info : 0

Info préservée : 1

Marche à suivre résumée

Pour $P = \{p_1, \dots, p_N\}$ un ensemble de points :

1. calculer $\mu = \frac{1}{N} \sum_{j=1}^N p_j$
2. recentrer P en $C = \{c_1, \dots, c_N\}, c_i = (p_i - \mu)$
3. calculer $S = \sum_{j=1}^N c_j c_j^T$
4. calculer les couples de S (λ_i, \vec{e}_i)

5. retenir les k vecteurs propres ($\vec{e}_1, \dots, \vec{e}_k$) des plus grandes valeurs propres, avec $(\lambda_1 \leq \dots \leq \lambda_k)$. On trie les vecteurs propres, en fonctions des valeurs des valeurs propres
6. former $W^T = \begin{bmatrix} e_{1,1} & \dots & e_{1,d} \\ \dots & \dots & \dots \\ e_{k,1} & \dots & e_{k,d} \end{bmatrix}$
7. faire la projection des c_i

Chapitre 6

Analyse discriminante

Cette fois ci, on va chercher la projection qui sépare *au mieux* les classes. On garde donc l'information de la classe dans le processus de réduction.

6.1 Calculs

6.1.1 Moyenne des projetés

Pour ce faire on choisira comme critères "*au mieux*" de faire attention à la moyenne de la projection des classes¹ :

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2$$

Cependant, il faut en plus, s'assurer de minimiser la somme des variances, notées : \hat{S}_1^2 et \hat{S}_2^2 tel que :

$$\frac{1}{\hat{S}_1^2 + \hat{S}_2^2}$$

Ainsi on obtient la fonction à minimiser suivante :

$$J(W) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\hat{S}_1^2 + \hat{S}_2^2}$$

W est la matrice de projection. On donne les définition des $\tilde{\mu}_i$:

$$\begin{aligned}\tilde{\mu}_i &= \frac{\sum_{j=1}^{|C_i|} W^T c_j}{|C_i|} \\ &= \frac{1}{|C_i|} W^T \sum_{j=1}^{|C_i|} c_j\end{aligned}$$

1. on travaillera ici avec 2 classes.

$$= W^T \frac{1}{|C_i|} \sum_{j=1}^{|C_i|} c_j$$

$$\tilde{\mu}_i = W^T \mu_i$$

On peut maintenant travailler sur la formule suivante :

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = W^T (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T W$$

simplifiable par, $W^T S_b W$, $S_b = (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T$

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = W^T S_b W$$

6.1.2 Variance

S_k avec $k \in \{1, 2\}$.

$$S_k = \frac{1}{|C_k|} \sum_{i=1}^{|C_k|} (c_i - \mu_k) (c_i - \mu_k)^T$$

On donne donc :

$$\begin{aligned} \hat{S}_k^2 &= \frac{1}{|C_k|} \sum_{i=1}^{|C_k|} (\tilde{y}_i - \tilde{\mu}_k)^2 \\ &= \frac{1}{|C_k|} \sum_{i=1}^{|C_k|} (W^T c_i - \tilde{\mu}_k)^2 \\ &= \frac{1}{|C_k|} \sum_{i=1}^{|C_k|} (W^T c_i - \tilde{\mu}_k)^2 \\ \hat{S}_k^2 &= \frac{1}{|C_k|} \sum_{i=1}^{|C_k|} W^T (c_i - \tilde{\mu}_k)^T (c_i - \tilde{\mu}_k) W \\ \hat{S}_k^2 &= \frac{1}{|C_k|} \sum_{i=1}^{|C_k|} W^T S_K W \end{aligned}$$

6.1.3 Nouvelle forme de la fonction à minimiser

$$J(W) = \frac{W^T S_b W}{W^T S_w W}, S_w = \sum S_k$$

On dérive partiellement $J(W)$ sur W :

$$\frac{\partial J(W)}{\partial W} = 0$$

$$S_b W - \alpha S_w W = 0$$

$$S_b W = \alpha S_w W$$

$$S_w^{-1} S_b W = \alpha W$$

$$S_w^{-1} S_b W = S_w^{-1} (\alpha (\mu_1 - \mu_2))$$

On sait grâce à cela que :

$$W = S_w^{-1} (\mu_1 - \mu_2)^T$$

2

6.2 Généralisation à n classes

$$S_b = \sum_{i=1}^n (\tilde{\mu}_i - \mu)(\tilde{\mu}_i - \mu)^T, \mu = \frac{\sum_{i=1}^N \tilde{\mu}_i}{N}$$

$$S_w = \sum_{i=1}^n \hat{S}_i$$

$$\hat{S}_k = \frac{1}{|C_k|} \sum_{i=1}^k (c_i - \mu_k)(c_i - \mu_k)^T$$

Ces formules nous amènent à trouver :

$$S_b W = \lambda S_w W$$

$$S_w^{-1} S_b W = \lambda W$$

On précise que dans ce cas $S_w^{-1} S_b$ est de rang $N - 1$

6.3 Que retenir ?

- Calculer S_b
- Calculer S_w
- $J(W) = \frac{W^T S_b W}{W^T S_w W} = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{S}_1^2 + \tilde{S}_2^2}$, dans le cas à deux classes.
- Dans le cas N classe on peut réduire de 1 à N-1

2. ajout de la transposé en TD

Chapitre 7

Classifieurs linéaires

On utilise un classifieur linéaire lorsque l'on travaille sur une vérité terrain à n classes. On va chercher à trouver un ensemble d'**hyperplan** qui sépare les classes.

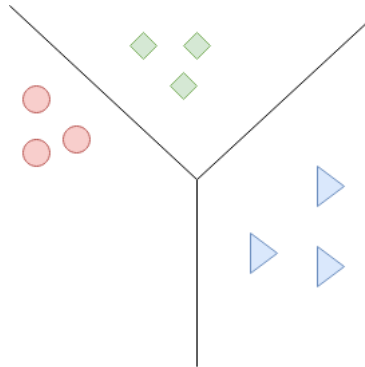


FIGURE 7.1 – Exemple de cas pouvant nécessiter un classifieur linéaire

7.1 Utilisation du problème à 2 classes

7.1.1 Explication du problème

On verra que l'on peut se servir du problème à deux classe pour le problème à n classes, cependant il existe une **marge** d'erreur.

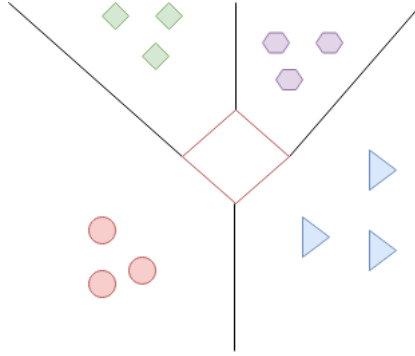


FIGURE 7.2 – Exemple de risques,lié à l'utilisation du problème à 2 classes pour n classes

7.2 Classifieur linéaire à deux classes

Dans ce classifieur linéaire on peut obtenir une fonction $g(x)$:

$$g(x) = W^T x + W_0$$

W^T est un vecteur normal à l'hyperplan, et W_0 un décalage à l'origine. On défini donc un hyperplan par g qui sépare l'espace en 2 palier :

$$g(x) > 0 \rightarrow x \in C_1$$

$$g(x) < 0 \rightarrow x \in C_2$$

$$g(x) = 0, \text{ On se sait pas}$$

Un problème ne sera linéairement séparable seulement si il existe (W, W_0) tel que :

- $\forall x_i \in C_1 : g(x_i) = W^T x_i + W_0 > 0$
- $\forall x_i \in C_2 : g(x_i) = W^T x_i + W_0 < 0$

7.2.1 Première transformation de g

On va maintenant transformer $g(x)$ en $\tilde{g}(\tilde{x})$:

$$\tilde{g}(\tilde{x}) = \tilde{W} \tilde{x}, \tilde{W} = \begin{bmatrix} W_0 \\ W \end{bmatrix}$$

et

$$\tilde{x} = \begin{bmatrix} 1 \\ x \end{bmatrix}$$

Il faut donc aussi réaliser une extension de la vérité terrain : $VT \rightarrow \tilde{V}T$. On réalise cela en rajoutant 1 devant tous les éléments de l'observation :

Exemple : $(a, b) \rightarrow (1, a, b)$.

7.2.2 Seconde transformation de g

On transforme $\tilde{g}(\tilde{x})$ en $\tilde{\tilde{g}}(\tilde{\tilde{x}})$, tel que $\forall \tilde{x} \in \tilde{\tilde{C}}_1 \cup \tilde{\tilde{C}}_2$.

Si $\tilde{\tilde{g}}(\tilde{\tilde{x}}) = \tilde{g}(\tilde{x})$, $\forall \tilde{x} \in \tilde{C}_1$, alors $\tilde{\tilde{C}}_1 = \tilde{C}_1$.

Comment faire pour $\tilde{\tilde{C}}_2$?

Pour ce cas si on va transformer $\tilde{x} \in \tilde{C}_2$ tel que $-\tilde{g}(x) = \tilde{\tilde{g}}(\tilde{x})$

Ainsi on obtient : $\tilde{g}(-\tilde{x}) = \tilde{\tilde{g}}(\tilde{x})$ et $\tilde{\tilde{C}}_2 = (-1)\tilde{C}_2$ On vient de faire les transformations suivantes :

$$VT \rightarrow \tilde{V}\tilde{T} \rightarrow \tilde{\tilde{V}}\tilde{\tilde{T}}$$

7.2.3 Changement de notation

$$\tilde{\tilde{W}} = a$$

$$\tilde{\tilde{W}}^T = a^T$$

et

$$\tilde{\tilde{x}} = y$$

Et donc on redéfinit le problème en cherchant :

$$\forall y \in \text{Observation transformé} = \{\tilde{\tilde{C}}_1 \cup \tilde{\tilde{C}}_2\}$$

$$a^T y \geq 0$$

7.2.4 Que fait-on dans le cas où $a^T < 0$?

Si $a^T y < 0$ alors on peut dire que a n'est pas solution du problème et que y est mal classé par a . On définira alors :

$$Y_m(a) = \{y \in \text{obs}\}, a^T y < 0$$

Il s'agit là de l'ensemble des y tel que $a^T y < 0$.

Si a est solution alors $Y_m(a) = \emptyset$. On doit trouver une fonction $J(a)$ dont a est la solution du problème¹ quand $J(a)$ est minimal.

On donne alors $J(a) = |Y - m(a)|$, le nombre de mal classés. Si $|Y - m(a)| = 0$ alors a est solution.

Le problème de cette fonction est qu'elle est constante par morceau :

1. séparateur linéaire

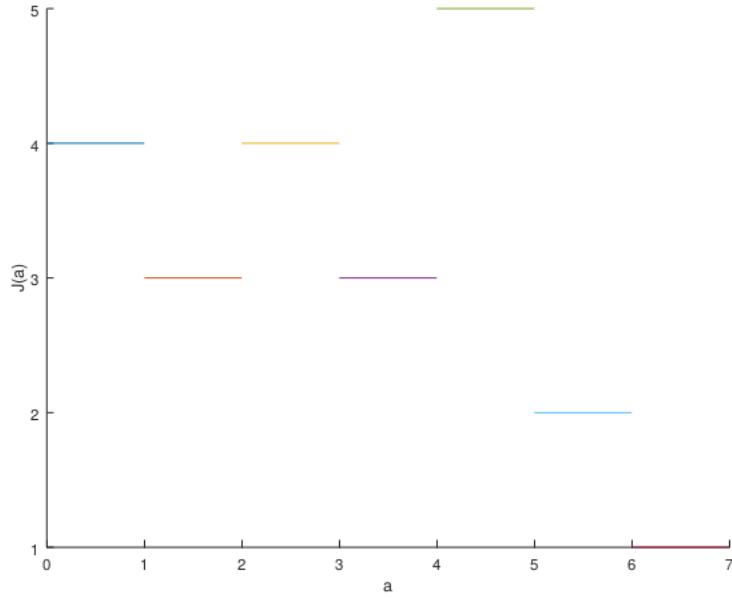


FIGURE 7.3 – $J(a)$ est constante par morceau

Comme cette fonction est constante, on ne peut pas faire de descente de gradient.

On redéfinit alors la fonction $J(a)$:

$$J(a) = \sum_{y \in Y_m(a)} \frac{a^T y}{\|a\|}$$

On obtient un problème de minimisation que l'on va transformer en problème de maximisation.

$$J(a) = - \sum_{y \in Y_m(a)} \frac{a^T y}{\|a\|}$$

Ainsi si $Y_m(a) = \emptyset \leftarrow J(a) = 0$ et cette fonction est cette fois ci continue par morceaux :

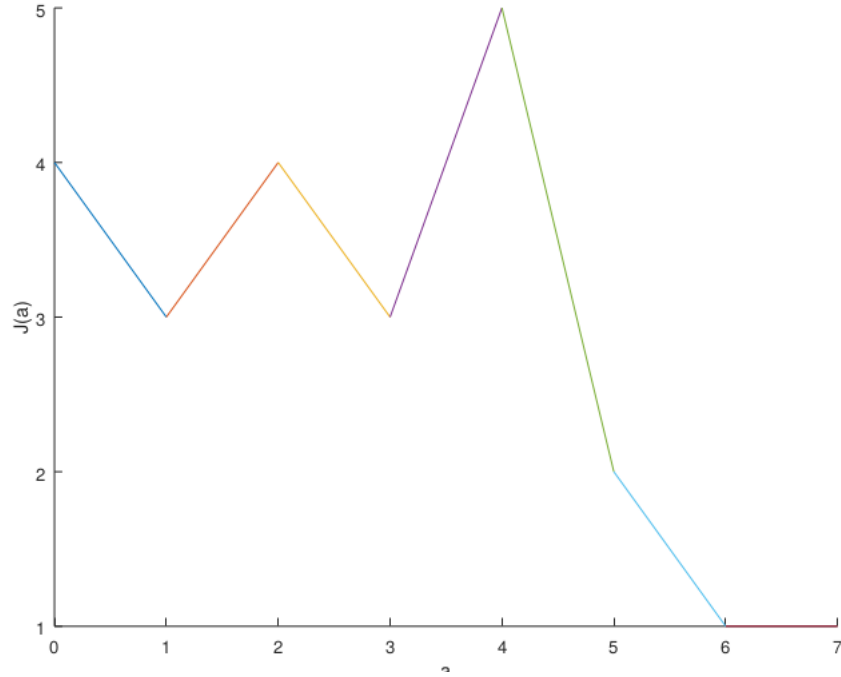


FIGURE 7.4 – $J(a)$ est continue par morceaux

On peut remarquer que l'on peut simplifier la fonction à maximiser :

$$\begin{aligned}
 J(a) &= - \sum_{y \in Y_m(a)} \frac{a^T y}{\|a\|} \iff J(a) = - \sum_{y \in Y_m(a)} a^T y \\
 &\iff J(a) = \sum_{y \in Y_m(a)} -a^T y
 \end{aligned}$$

Dérivée

Pour maximiser il nous faut la fonction dérivée de $J(a)$:

$$\frac{\partial J(a)}{\partial a} = - \sum_{y \in Y_m(a)} y$$

Le problème que l'on a désormais vient du fait que l'on a pas d'expression analytique pour $Y_m(a)$. Tel que, pour calculer $Y_m(a)$ il faut passer en revue, l'ensemble des éléments de l'observation.

On utilise donc une descente de gradient :

a_0 , à calculer ou choisir

$a_k = a_0$, à l'initialisation

$$a_{k+1} = a_k + \sum_{y \in Y_m(a)} y$$

Exemple :

On prend comme valeurs :

$$\tilde{obs} = \{(1, 1, 2), (1, 3, 2), (-1, -1, -1), (-1, -2, -3)\}$$

$$a_0 = [1, -2, 2]$$

On peut donc commencer les calculs :

$$Y_m(a_0) = ?$$

$$a_0^T y_1 = [1, -2, 2] \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix} = 3 \text{ bien classé}$$

$$a_0^T y_2 = [1, -2, 2] \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = 3 \text{ bien classé}$$

$$a_0^T y_3 = [1, -2, 2] \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} = -1 \text{ mal classé}$$

$$a_0^T y_4 = [1, -2, 2] \begin{bmatrix} -1 \\ -3 \\ -2 \end{bmatrix} = 1 \text{ bien classé}$$

$$Y_m(a_0) = \{y_3\} = \{(-1, -1, -1)\}$$

Ainsi on calcule a_1 :

$$a_1 = a_0 + \sum_{y \in Y_m(a_0)} y = \begin{bmatrix} 1 \\ -2 \\ 2 \end{bmatrix} + \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ -3 \\ 1 \end{bmatrix}$$

Puis on réitère et on obtient :

$$a = [0, -4, 3]$$

Conclusion

Ainsi pour évaluer a_k on possède deux stratégies :

1. Calculer toute l'observation, $Y_m(a_k)$
2. On prend le premier mal classé et on l'ajoute à a_k

Ainsi :

- Si le problème est linéairement séparable, l'algorithme converge vers un des solutions.
- Si le problème n'est pas linéairement séparable, alors il n'existe pas de a_k tel que $y_m(a_k) = \emptyset$ on entre dans une boucle infinie.

Pour résoudre la deuxième situation on va atténuer la contribution du gradient via un coefficient $\eta(k)$:

$$a_{k+1} = a_k + \eta(k) \sum_{y \in Y_m(a_0)} y$$

Ainsi dans le cas où le problème n'est pas linéairement séparable, la solution dépendra de a_0 et aussi de $\eta(k)$.

7.2.5 Système d'équation

On a cherché a tel que $\forall y \in \tilde{obs}, a^T y > 0$.

On donne une nouvelle approche :

Soit $b \in \mathbb{R}^{n+1}$, on cherche a tel que $Ya = b$ avec :

$$Y = \begin{bmatrix} e_i \in \tilde{C}_1 \\ e_i \in \tilde{C}_2 \\ \dots \\ e_i \in \tilde{C}_j \end{bmatrix} \quad \text{la matrice des éléments de l'observation}$$

$$b = \begin{bmatrix} b_1 > 0 \dots \\ b_{max} > 0 \end{bmatrix}$$

avec $max = |\tilde{Obs}|$.

On donne ainsi :

$$Ya = \begin{bmatrix} (Y_1 a^T) = (a^T Y_1) \\ \dots \\ (Y_k a^T) = (a^T Y_k) \end{bmatrix}$$

$$Ya = b \rightarrow \begin{bmatrix} (Y_1 a^T) = (a^T Y_1) \\ \dots \\ (Y_k a^T) = (a^T Y_k) \end{bmatrix} = \begin{bmatrix} b_1 > 0 \dots \\ b_{max} > 0 \end{bmatrix}$$

Remarque

$a^T > 0$ englobe $Ya = b$ car on fixe b dans $Ya = b$ ce qui revient à dire que $Ya = b$ est un sous problème de $a^T y > 0$

7.2.6 Avec le vecteur b

Il y a une solution si le nombre d'éléments de l'observation est égal à la dimension de l'espace qui définit a .

Ainsi, Si $x \in obs, x \in \mathbb{R}^2$ et $y \in \tilde{obs}, y \in \mathbb{R}^3$. Il faut 3 éléments dans l'observation.

Si on a pour a , 3 inconnues alors $a \in \mathbb{R}^3$. ainsi si on a 200 observations, on obtient un système à 200 équations et à 3 inconnues. On peut facilement affirmer que le système n'as pas une solution unique et que :

$$Ya \approx b$$

C'est donc la fonction suivante que l'on va minimiser :

$$J(a) = ||Ya - b||^2, \text{ avec } b \text{ choisi}$$

Comment obtenir cette fonction ?

On cherche a tel que $\forall y_i \in \tilde{obs}, a^T y_i > 0$

À partir de ceci on souhaite calculer $Ya - b$. Y est un vecteur où chaque y_i correspond à la distance entre un élément de l'observation et la projection pour a .

On associe donc à un y_i un a

$$y_i, a : a^T y_i - b_i$$

On donne ainsi :

$$\begin{aligned} J(a) &= \sum_{i=1}^m (a^T y_i - b_i)^2 \\ &= ||Ya - b||^2 \end{aligned}$$

Dérivée

La dérivée partielle de cette fonction est la suivante :

$$\frac{\partial J(a)}{\partial a} = 2y^T(ya - b)$$

Donc, si $2y^T(ya - b) = 0 \implies \nabla J(a) = 0$.

À partir de là il y a deux solutions :

1. résolution de $2y^T(ya - b)$, soit une résolution (approximée) du système d'équation.
2. une descente de gradient

Résolution du système

On souhaite résoudre :

$$2y^T(ya - b) = 0$$

Que l'on transforme en :

$$Y^T(Ya - b = 0)$$

$$Y^TYa = Y^Tb$$

On notera que Y^TY est une matrice carrée.

Si Y^TY est inversible alors $(Y^TY)^{-1}$ existe. Ainsi on dit :

$$(Y^TY)^{-1}Y^TYa = (Y^TY)^{-1}Y^Tb$$

$$a = (Y^TY)^{-1}Y^Tb$$

On obtient alors ici un a qui fait de $Ya \approx b$ une bonne approximation. Au final on obtient la relation suivante :

$$Ya \approx b \approx \begin{bmatrix} b_1 + \epsilon_1 \\ \dots \\ b_n + \epsilon_n \end{bmatrix}$$

Cependant rien ne garantit que $\epsilon_i = 0$. De ce fait on ne peut pas garantir que la a obtenue sépare linéairement le problème et ce même si le problème est linéairement séparable

On se rend compte que la solution va dépendre de la valeur de b .

$$\text{Si } b = \begin{bmatrix} 1 \\ \dots \\ 1 \end{bmatrix}, \text{ alors } \|Ya - \begin{bmatrix} 1 \\ \dots \\ 1 \end{bmatrix}\|^2 \rightarrow a.$$

$$\text{Et si } b = \begin{bmatrix} 2 \\ \dots \\ 2 \end{bmatrix}, \text{ alors } \|Ya' - \begin{bmatrix} 2 \\ \dots \\ 2 \end{bmatrix}\|^2 \rightarrow a' \text{ alors :}$$

$$a' = 2a$$

$$\|Y2a - \begin{bmatrix} 2 \\ \dots \\ 2 \end{bmatrix}\|^2 = 2\|Ya - \begin{bmatrix} 1 \\ \dots \\ 1 \end{bmatrix}\|^2$$

On voit bien que b n'intervient pas dans la détermination de a

$$\text{Si } \forall i, b_i = \text{constante alors la solution peut se ramener à } \begin{bmatrix} 1 \\ \dots \\ 1 \end{bmatrix}.$$

$$\text{Si } a \text{ est solution pour } b = \begin{bmatrix} 1 \\ \dots \\ 1 \end{bmatrix} \text{ alors constante } *a \text{ est solution de } b = \begin{bmatrix} \text{constante} \\ \dots \\ \text{constante} \end{bmatrix}.$$

Ainsi on pose :

$$Y a = \begin{bmatrix} a^T y_1 \\ \dots \\ a^T y_n \end{bmatrix} \text{ avec } a^T y_i < 0$$

Et de ce fait :

$$Y \text{ constante } a \begin{bmatrix} \text{constante } a^T y_1 \\ \dots \\ \text{constante } a^T y_n \end{bmatrix} \text{ et constante } a^T y_i < 0$$

Cas à prendre en compte : si $b = \begin{bmatrix} 1 \\ \dots \\ 1 \\ 10 \end{bmatrix}$ On voit que le dernier terme

de ce vecteur est différent : il représente un traitement particulier. Notamment, le point associé sera plus loin que les autres dans la solution.

La valeur de ce b_i pondère la distance enter la projection du i-ème échantillon et la droite obtenue

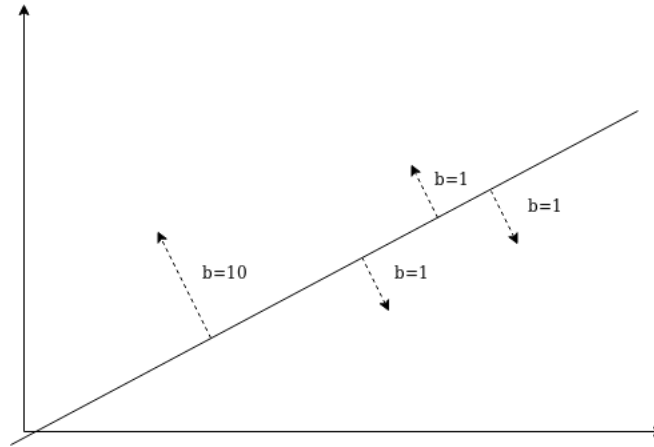


FIGURE 7.5 – Plus b est grand plus la distance est grande

Comme on le voit **plus b_i est grand, plus la distance est grande**, ce qui veut aussi dire que **moins b_i contribue par rapport aux autres**

Descente de gradient en fonction de a

On donne l'algorithme de descente de gradient suivant :

$$\text{Tant que } |a_{k+1}| \geq \epsilon$$

$$a_{k+1} = a_k - \eta(k) \nabla J(a_k)$$

On rappelle que $\nabla J(a_k) = 2Y^T(Ya - b)$, cependant il s'agit là de données difficiles à stocker. Ainsi on se propose de le faire élément par élément :

$$a_{k+1} = a_k - \eta(k) 2y_i[(a_k^T y_i) - b_i]$$

Cette descente de gradient se fait uniquement en fonction de a, b lui étant fixé. Si jamais après la descente, a n'est pas solution Ya , il faut changer la valeur de b

Descente de gradient en fonction de a et de b

On fait maintenant minimiser en fonction de a et de b :

$$J(a, b) = \|Ya - b\|^2, b \geq 0$$

On se donne un a et on approxime b . On commence par calculer une dérivée partielle :

$$\frac{\partial J(a, b)}{\partial b} = 2(Ya - b)$$

Si on essaies de résoudre :

$$2(Ya - b) = 0$$

$$Ya = b$$

Il n'y a pas de solution analytique.

On effectue donc une descente de gradient :

$$b_{k+1} = b_k - 2(Ya - b)$$

$$= b_k - \eta(k) 2(Ya - b)$$

Cependant avec cette descente de gradient il se peut que b_{k+1} prenne une valeur négative. On donne alors le nouvel algorithme suivant avec :

$$e_k = Ya - b_k$$

$$b_{k+1} = b_k + \eta(e_k + |e_k|)$$

Si $e_k > 0 \rightarrow e_k + |e_k| \implies 2e_k$. Sinon $e_k < 0 \rightarrow e_k + |e_k| \implies 0$.

Algorithme Final :

1. On se donne $b_0 > 0$
 2. On calcul a_k en utilisant :
 - la pseudo-inverse
 - la descente de gradient
 3. On bloque a_k obtenu
 4. On fait une descente gradient en utilisant a_k et on calcule b_k (on remonte à 2 si besoin)
- On s'arrêtera quand :

$$b_{k+1} = b_k + \epsilon$$

ou alors quand :

$$a_{k+1} = a_k + \epsilon$$

Problèmes

Que l'on utilise l'algorithme du Perceptron ou que l'on calcul une approximation :

On ne marise pas la position de la droite/hyperplan

La droite doit être aussi près de certains éléments de C_1 que de C_2 :

$$\min(\text{dist}(D, y \in C_1)) = \min(\text{dist}(D, y \in C_2))$$

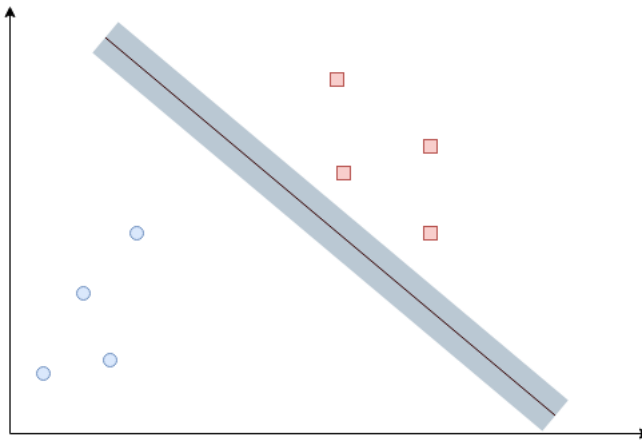


FIGURE 7.6 – la droite doit assez proches des différents éléments des classes