

Cours 2

Performance d'un classifieur

Méthode supervisée vs non supervisée

Méthode générative vs discriminative

Classification par les kppv

Arbre de décision et forêt aléatoire

Catherine ACHARD
Institut des Systèmes Intelligents et de Robotique

catherine.achard@upmc.fr

Classifier = associer une **classe** C à un **vecteur de caractéristiques** x de dimension n

Vecteur de caractéristique x = forme + variabilité + bruit de mesure

Connaissances disponibles

- Informations fournies par un « expert »
- Modèles explicites (méthode structurelle)
- Cas le plus général : **base de données étiquetées ou non**

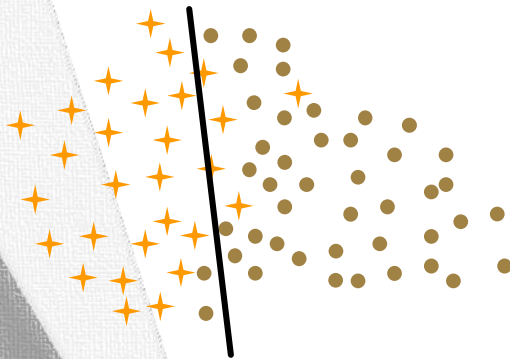
Base de données divisée en 3

- **Base d'apprentissage** : pour apprendre le modèle
- **Base de validation** : pour aider à définir certains paramètres du modèle
- **Base de test** : pour tester le modèle sur de nouvelles données jamais rencontrées

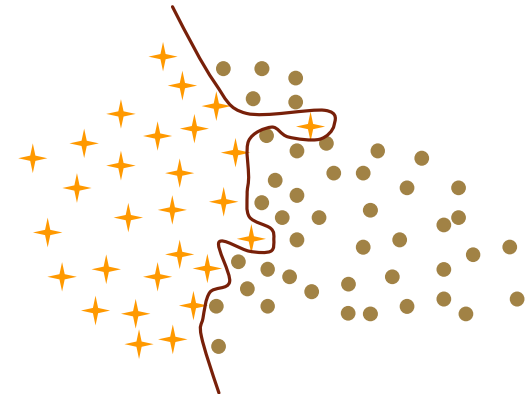
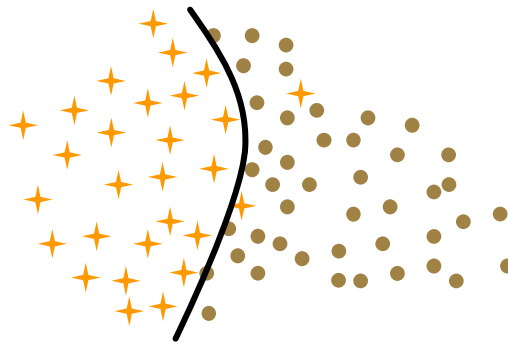
Un bon apprentissage devra

- Avoir une faible erreur sur la base d'apprentissage
- Avoir un écart faible entre l'erreur d'apprentissage et l'erreur de test

Lorsque c'est le cas, on parle de **bonne généralisation** : le système est capable de reconnaître des données jamais vues



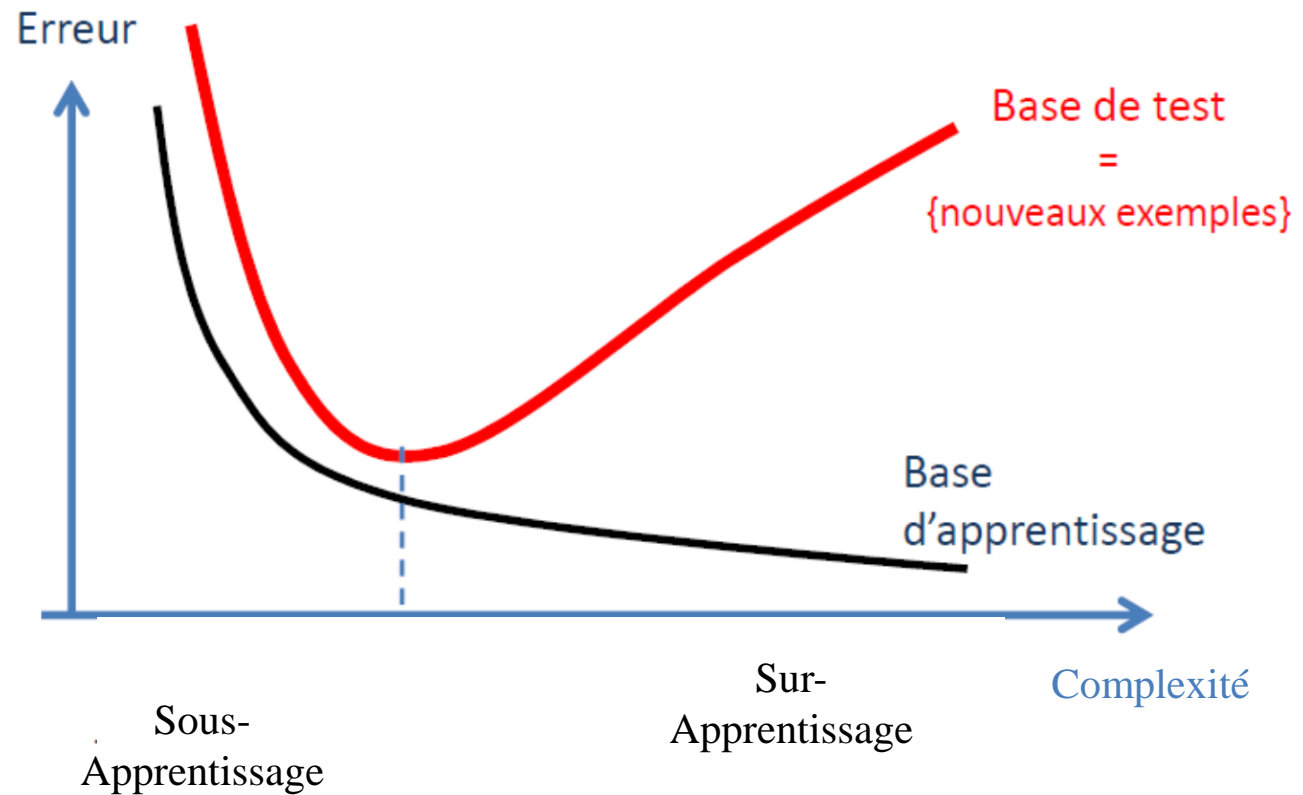
Sous-apprentissage



Sur-apprentissage

Généralisation

- Bonne généralisation, où est la frontière ?



Généralisation

Sur-apprentissage (over-fitting)

On se fie trop aux données

- petit changement de la base d'apprentissage = forte variation de la sortie
- **Variance** élevée

Sous-apprentissage (under-fitting)

On se fie moins aux données

- on s'éloigne par endroit de la vraie frontière
- on introduit un **biais**

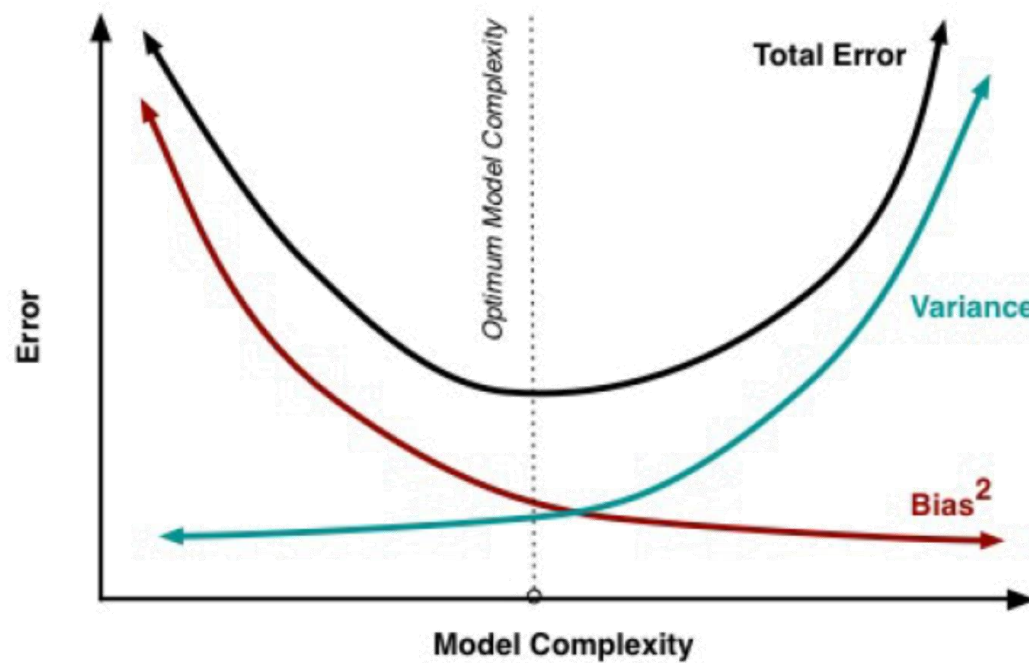
Si le modèle est trop simple → biais élevé et variance faible.

Si le modèle est trop complexe → biais faible et une variance élevée.

Compromis biais/variance

Généralisation

Erreur de prédiction = $\text{Biais}^2 + \text{Variance} + \text{Erreur Irréductible}$



<http://scott.fortmann-roe.com/docs/BiasVariance.html>

Généralisation

Comment avoir un bon compromis biais/variance ?

Comment avoir une bonne généralisation

- En diminuant la dimension des exemples
- En augmentant le nombre d'exemples
- Sélectionner un bon modèle ou en le contraignant
- Méthodes ensemblistes

Performances d'un classifieur

En apprentissage statistique, il n'y a **pas de méthode idéale**, celle-ci varie en fonction des données, de leur dimension, du problème posé,...

- Le choix de la méthode la plus adaptée à un cas précis repose sur une **estimation de l'erreur**
- Cette estimation devra être la plus rigoureuse possible

On utilise 3 bases :

- **Une base de référence ou d'apprentissage** utilisée pour mettre au point le classificateur
- **Une base de validation** pour déterminer les paramètres du classifieur
- **Une base de test** : exemples jamais vus au préalable pour évaluer le classificateur

Petite base de données

Lorsque l'on possède peu de données, il est difficile d'estimer les performances de manière fiable.

Solution 1 : K-fold validation :

- Division de la base en K parties.
- $(K-1)$ parties pour l'apprentissage et la dernière en test puis estimation de l'erreur
- On fait tourner la partie test sur les K parties
- L'erreur finale est la moyenne des erreurs

Solution 2 : Leave-One-Out Cross-Validation (LOOCV)

- C'est le cas particuliers où $K = \text{nombre d'exemples}$
- $N-1$ exemples pour l'apprentissage et le dernier exemple en test
- On fait tourner l'exemple de test sur tous les exemples de la base

Solution 3 : bootstrap

Tirage aléatoire pour déterminer les exemples d'apprentissage et de test et on calcule l'erreur
Plusieurs itérations puis erreur moyenne

Matrice de confusion :

Classe réelle
↓

Classe trouvée →

| | 1 | 2 | 3 |
|---|--|--|--|
| 1 | e_{11} Nb d'exemples 1 étiquetés 1 | e_{12} Nb d'exemples 1 étiquetés 2 | e_{13} Nb d'exemples 1 étiquetés 3 |
| 2 | e_{21} Nb d'exemples 2 étiquetés 1 | e_{22} Nb d'exemples 2 étiquetés 2 | e_{23} Nb d'exemples 2 étiquetés 3 |
| 3 | e_{31} Nb d'exemples 3 étiquetés 1 | e_{32} Nb d'exemples 3 étiquetés 2 | e_{33} Nb d'exemples 3 étiquetés 3 |

Exercice

Matrice de confusion :

Classe réelle

↓

Classe trouvée →

| | 1 | 2 | 3 |
|---|----|----|-----|
| 1 | 90 | 6 | 5 |
| 2 | 20 | 70 | 4 |
| 3 | 2 | 1 | 104 |

Quel est le nombre d'exemples de la base de test ?

Que représente le chiffre 20 ?

Que représente le chiffre 104 ?

Quel est le taux de bonne reconnaissance ?

Exercice

Matrice de confusion :

Classe réelle
↓

Classe trouvée →

| | 1 | 2 | 3 |
|---|----|----|-----|
| 1 | 90 | 6 | 5 |
| 2 | 20 | 70 | 4 |
| 3 | 2 | 1 | 104 |

Quel est le nombre d'exemples de la base de test ?

La somme des éléments de la matrice = 302

Que représente le chiffre 20 ?

Le nombre d'exemples appartenant à la classe 2 étiquetés 1

Que représente le chiffre 104 ?

Le nombre d'exemples appartenant à la classe 3 étiquetés 3

Quel est le taux de bonne reconnaissance ?

La somme des éléments de la diagonale divisée par la somme des éléments
 $264/302 \rightarrow 87,4 \%$

Taux de bonne classification

Sans rejet

$$\text{Taux de bonne classification } T_b = \frac{\text{Nb exemples bien classés}}{\text{Nb exemples}}$$

$$\text{Taux d'erreur } T_e = 1 - T_b$$

Avec rejet

$$\text{Taux de rejet } T_r = \frac{\text{Nb exemples rejeté}}{\text{Nb exemples}}$$

$$\text{Taux de bonne classification } T_b = \frac{\text{Nb exemples bien classés}}{\text{Nb exemples}}$$

$$\text{Taux d'erreur } T_e = 1 - T_b - T_r$$

Taux de bonne classification

Problème :

Il s'agit d'une mesure faible qui **ne tient pas compte de la distribution des classes**

➔ Quand les classes sont très disproportionnées, prédire systématiquement la classe majoritaire amène à un bon taux de classification

Exercice :

En diagnostic médical, très peu de personnes sont malades (5%?). Quel est le taux de classification si on prédit toujours que la personne est saine ?

Exemple sur 100 personnes

| | malade | sain |
|--------|--------|------|
| malade | 0 | 5 |
| sain | 0 | 95 |

Tb=95%

Taux de bonne classification

Solution :

On tient compte de la répartition des classes et on construit une **matrice de confusion normalisée**

| | | |
|--|--------------|--------------|
| | | |
| | e_{11}/N_1 | e_{12}/N_1 |
| | e_{21}/N_2 | e_{22}/N_2 |

N_1 : Nombre d'exemples de la classe 1

N_2 : Nombre d'exemples de la classe 2

Le nouveau taux de bonne classification devient

$$tb = \frac{1}{N_c} \sum_{i=1}^{N_c} \frac{e_{ii}}{N_i} \text{ où } N_c \text{ est le nombre de classes}$$

Exercice :

reprendre l'exemple précédant avec cette normalisation

On a maintenant

| | malade | sain |
|--------|--------|------|
| malade | 0 | 1 |
| sain | 0 | 1 |

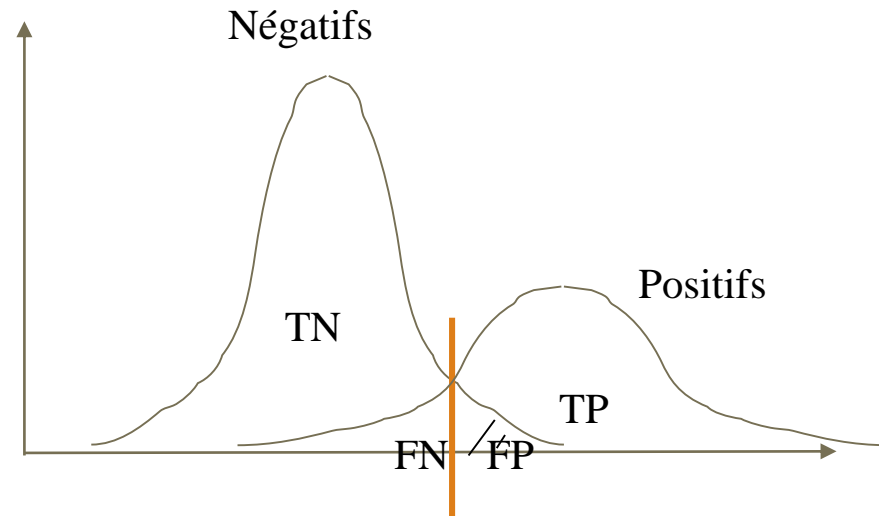
Et $tb = 50\%$

Courbe ROC (Receiver Operating Characteristic)

Permet de comparer plusieurs classificateurs indépendamment d'un seuil pour un problèmes à 2 classes (**positif et négatif**)

On définit les :

- Vrai Positif (True Positive)
- Vrai Négatif (True Négatif)
- Faux Négatif (False Négatif)
- Faux Positif (False Positif)



Classe réelle



Classe trouvée



| | + | - |
|---|----|----|
| + | TP | FN |
| - | FP | TN |

Courbe ROC (Receiver Operating Characteristic)

$$\text{Sensibilité} = \frac{TP}{TP+FN} = \frac{\text{Nb positifs bien classés}}{\text{Nb positifs}}$$

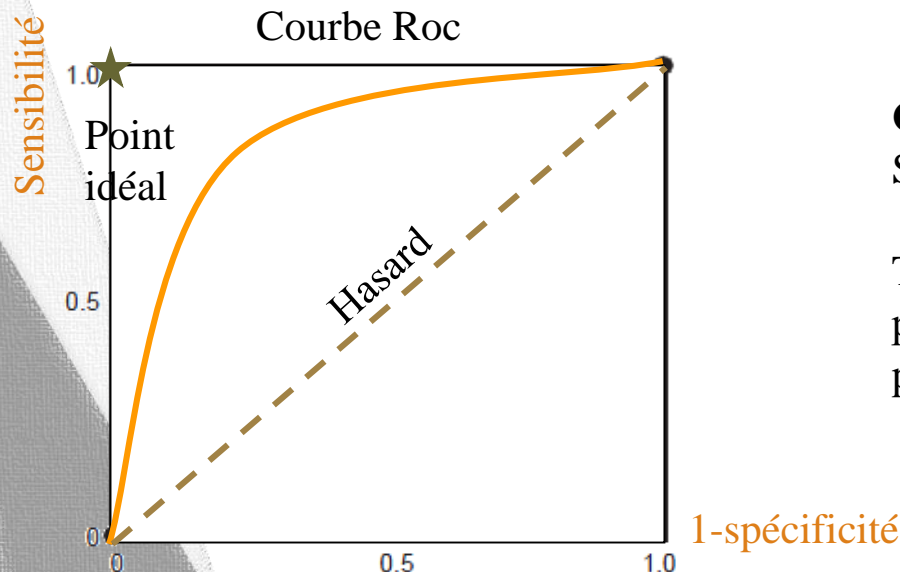
$$\text{Spécificité} = \frac{TN}{FP+TN} = \frac{\text{Nb négatifs bien classés}}{\text{Nb négatifs}}$$

| | + | - |
|---|----|----|
| + | TP | FN |
| - | FP | TN |

Un bon classificateur devra être

- sensible : détecter les positifs = pourcentage de vrais positifs détectés
- spécifique : ne pas tout détecter comme positif = pourcentage de vrais négatifs détectés

Généralement, plus un classificateur est sensible, moins il est spécifique et vice versa



Courbe ROC

$$\text{Sensibilité} = f(1\text{-spécificité})$$

Toutes les courbes ROC passent par l'origine et par le point (1,1)

Précision/Rappel

On trouve aussi parfois les scores exprimés en termes de Précision / Rappel

Précision = $\frac{TP}{TP+FP} = \frac{\text{Nb positifs bien classés}}{\text{Nb classés positifs}} = \text{pourcentage de vrais positifs parmi les positifs trouvés}$

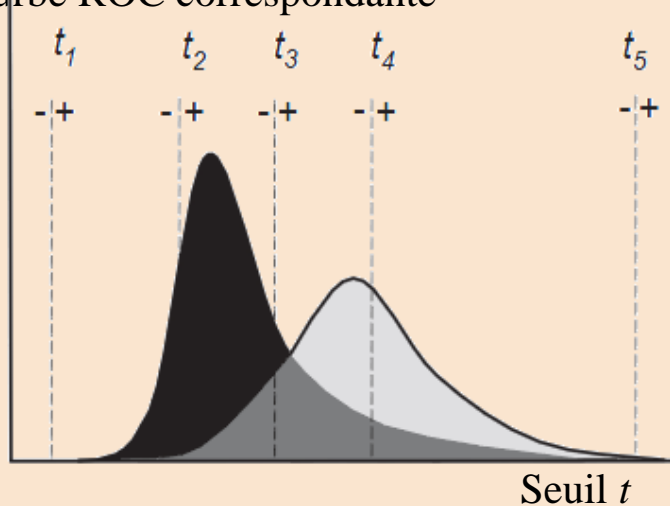
Rappel $\frac{TP}{TP+FN} = \frac{\text{Nb positifs bien classés}}{\text{Nb positifs}} = \text{pourcentage de vrais positifs détectés (même chose que sensibilité)}$

| | + | - |
|---|----|----|
| + | TP | FN |
| - | FP | TN |

Exercice :

On souhaite réaliser une classification binaire. Plusieurs algorithmes dépendant d'un seuil doivent être comparés.

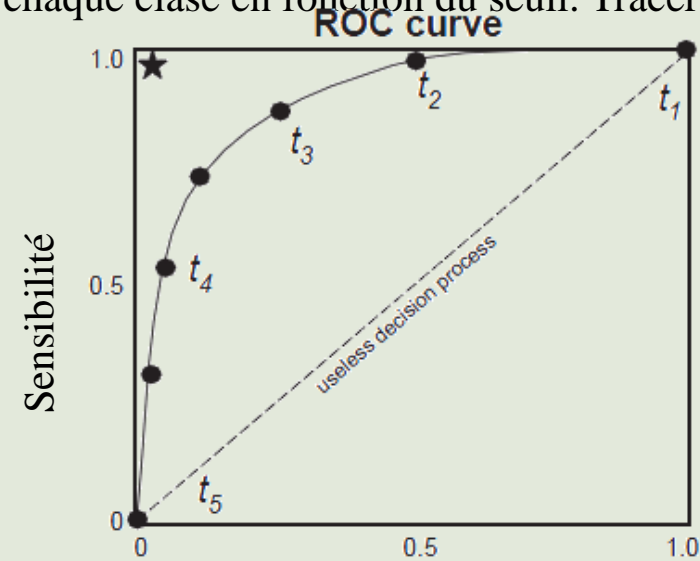
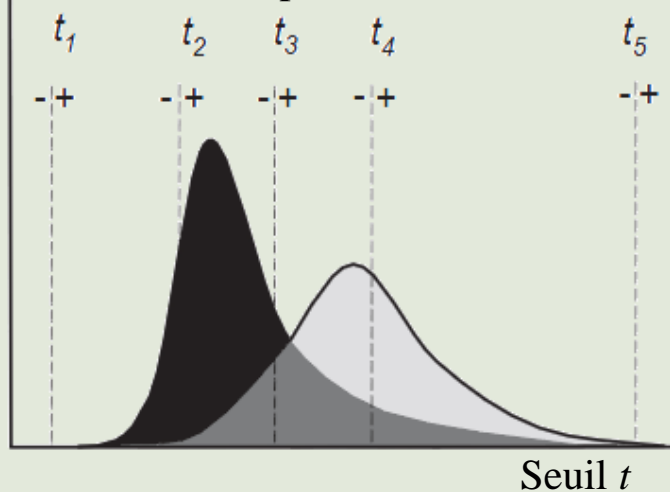
Un algorithme donné donne les densité de probabilité de chaque classe en fonction du seuil. Tracer la courbe ROC correspondante



Exercice :

On souhaite réaliser une classification binaire. Plusieurs algorithmes dépendant d'un seuil doivent être comparés.

Un algorithme donné donne les densité de probabilité de chaque classe en fonction du seuil. Tracer la courbe ROC correspondante



$$\text{Sensibilité} = \frac{TP}{TP + FN}$$

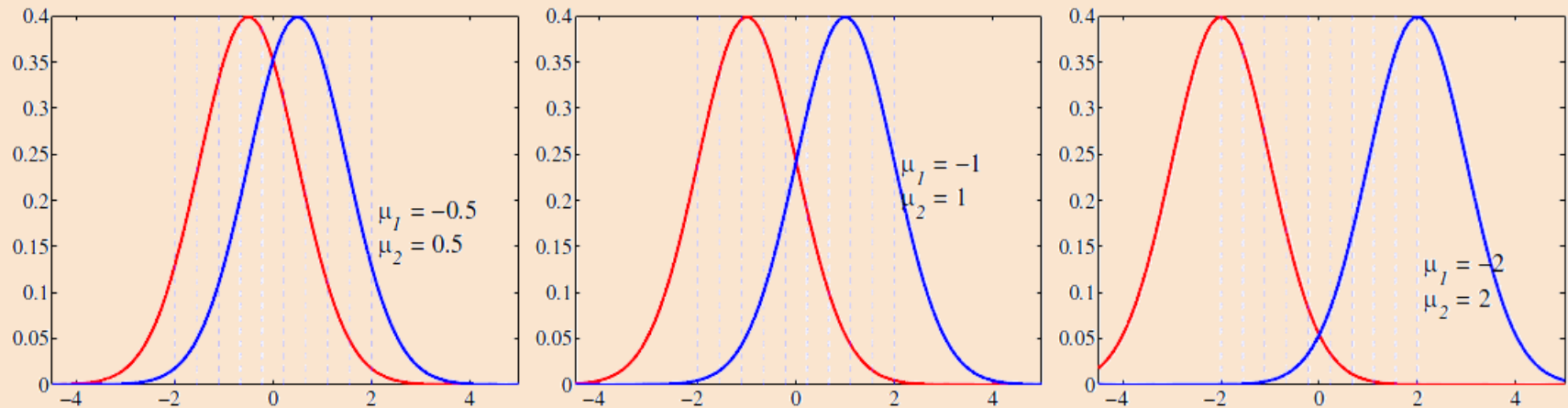
$$\text{Spécificité} = \frac{TN}{TN + FP}$$

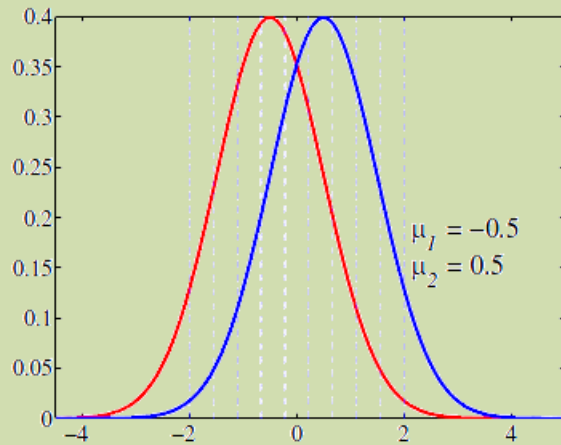
Remarque:

$$\begin{cases} TP + FN = cte \\ TN + FP = cte \end{cases}$$

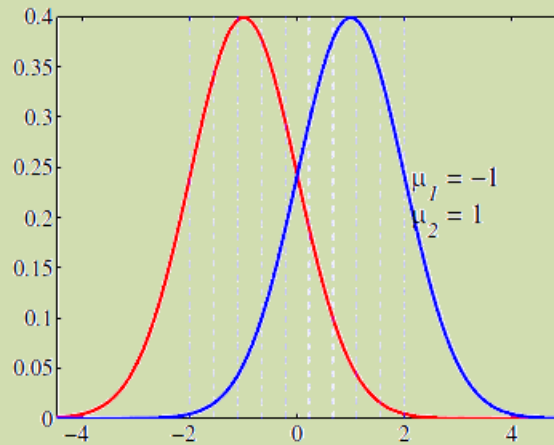
Exercice

Tracer l'allure des courbes ROC des 3 classifieurs ci-dessous (mêmes figures que pour l'exercice précédent))

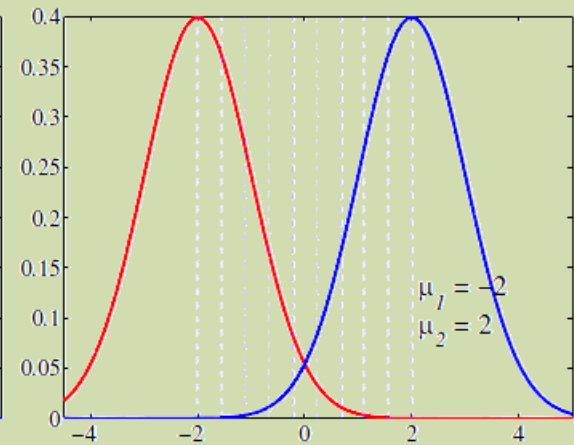




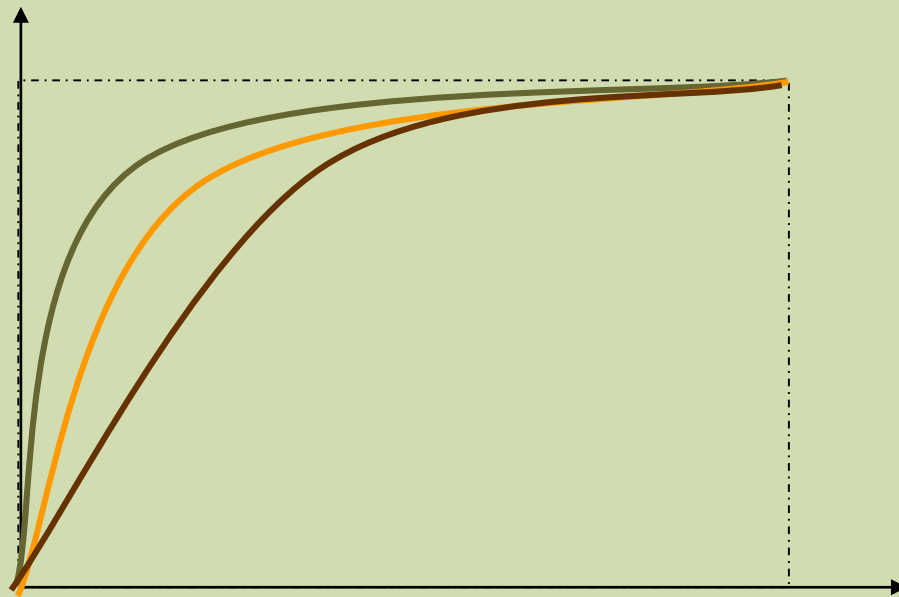
Classifieur 1



Classifieur 2



Classifieur 3



Performances d'un classifieur

Exercice

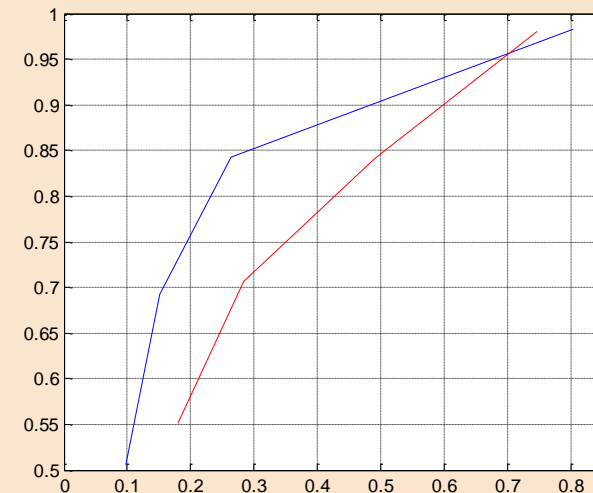
Nous souhaitons détecter des objets en mouvement (positif) dans des séquences d'images.

→ 2 algorithmes que l'on veut comparer indépendamment du seuil de détection

| | | | | | |
|---|------------------|--------|--------|--------|-------|
| A = Nombre de pixels détectés en mouvement B = Nombre de pixels détectés non mouvement C = Nombre de pixels réellement en mouvement et détecté en mouvement D = Nombre de pixels réellement en mouvement = 57587 | | | A | B | C |
| | Première méthode | Seuil1 | 372417 | 77583 | 56629 |
| | | Seuil2 | 151899 | 298101 | 48526 |
| | | Seuil3 | 98997 | 351003 | 39847 |
| | | Seuil4 | 67181 | 382819 | 29135 |
| | Deuxième méthode | Seuil1 | 349806 | 100194 | 56479 |
| | | Seuil2 | 241185 | 208815 | 48472 |
| | | Seuil3 | 151902 | 298098 | 40654 |
| | | Seuil4 | 102264 | 347736 | 31739 |

1. Donner l'expression analytique du nombre de vrais positifs (TP), vrais négatifs (TN), faux positifs (FP) et faux négatifs (FN) en fonction de A, B, C, D.
2. Dans la base de test, quel est le nombre d'exemples positifs et d'exemples négatifs ? Donner l'expression analytique et le résultat numérique.
3. Donner l'expression analytique de la sensibilité et de la spécificité en fonction de A, B, C et D.
4. Les deux courbes ROC sont représentées ci-contre. Quelle est la courbe ROC de la première méthode ? Justifier votre réponse. Quelle méthode donne les meilleurs résultats ?

Quel est le point de fonctionnement avec le meilleur taux de reconnaissance ? Quel est-il ?



Performances d'un classifieur

1. Donner l'expression analytique du nombre de vrais positifs (TP), vrais négatifs (TN), faux positifs (FP) et faux négatifs (FN) en fonction de A, B, C, D.

| décision \ étiquette | + | - |
|----------------------|----|----|
| + | TP | FN |
| - | FP | TN |

$$\begin{aligned}
 A &= TP + FP \\
 B &= TN + FN \\
 C &= TP \\
 D &= TP + FN
 \end{aligned}
 \Rightarrow
 \begin{cases}
 TP = C \\
 FP = A - C \\
 FN = D - C \\
 TN = B - D + C
 \end{cases}$$

2. Dans la base de test, quel est le nombre d'exemples positifs et d'exemples négatifs ?

Donner l'expression analytique et le résultat numérique.

$$\text{Nb ex} > 0 = TP + FN = C + D - C = D = 57587$$

$$\text{Nb Ex} < 0 = TN + FP = A - C + B - D + C = A + B - D = 392413$$

3. Donner l'expression analytique de la sensibilité et de la spécificité en fonction de A, B, C et D.

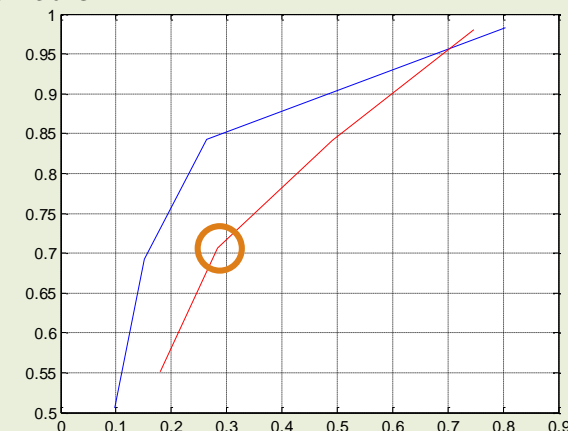
$$\text{sens} = \frac{C}{D} = \frac{C}{57587} \quad \text{spe} = \frac{B - D + C}{A + B - D} = \frac{B - D + C}{392413}$$

En prenant le seuil 3 de la méthode 2, on a :

$$\text{Sens} = 0.7 \text{ et } \text{Spe} = 0.71 \rightarrow 1 - \text{spe} = 0.28$$

→ c'est la courbe du bas

→ La courbe du haut est la première méthode et celle du bas la seconde



Méthodes supervisées / Non supervisées

Méthode supervisée

On dispose d'un ensemble d'exemples étiquetés en apprentissage

→ On souhaite savoir classer un nouvel exemple

Le nombre de classes est connu *a priori*

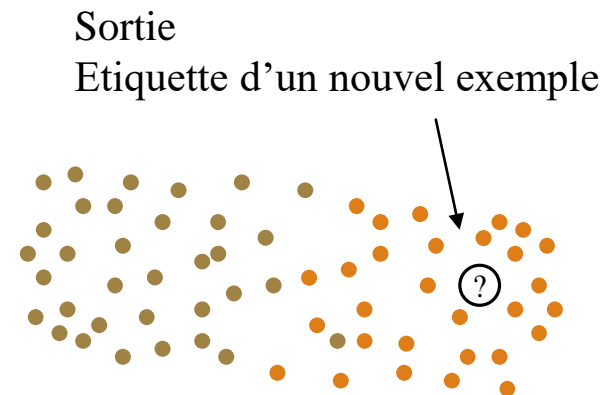
Exemple

On dispose d'un ensemble de signaux audio enregistré à partir de 50 personnes. L'identité de la personne est connue pour chaque enregistrement

Pour un nouveau signal, on souhaite déterminer l'identité de la personne



Entrée
Exemples étiquetés



Méthode non supervisée

On dispose d'un ensemble d'exemples non étiquetés en apprentissage

→ On souhaite partitionner cet ensemble en sous-ensembles homogènes

Le nombre de classes (sous-ensembles) n'est pas connu *a priori*

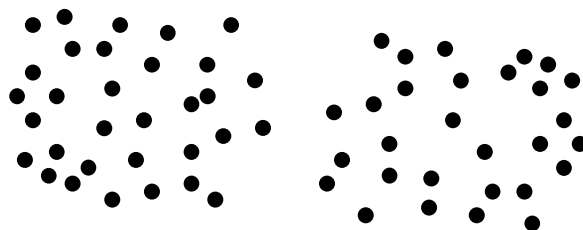
Exemple

On dispose des achats de livres de clients sur amazon

On souhaite catégoriser les clients en fonction de leur goûts afin de leur proposer des nouveaux achats pertinents

→ On peut utiliser ces méthodes pour explorer et comprendre les données

Entrée : exemples non étiquetés



Sortie: sous-ensembles homogènes



Méthodes génératives / discriminatives

Soit \mathbf{x} les exemples de dimension n et y leur classe telle que $y \in [1, \dots, K]$

But de la classification : déterminer $p(y|\mathbf{x})$

Approche discriminative

Détermine directement $p(y|\mathbf{x})$

Approche générative / Classification bayésienne :

Détermine pour chaque classe $p(\mathbf{x}|y)$ et $P(y)$ puis utilise le théorème de Bayes :

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})}$$

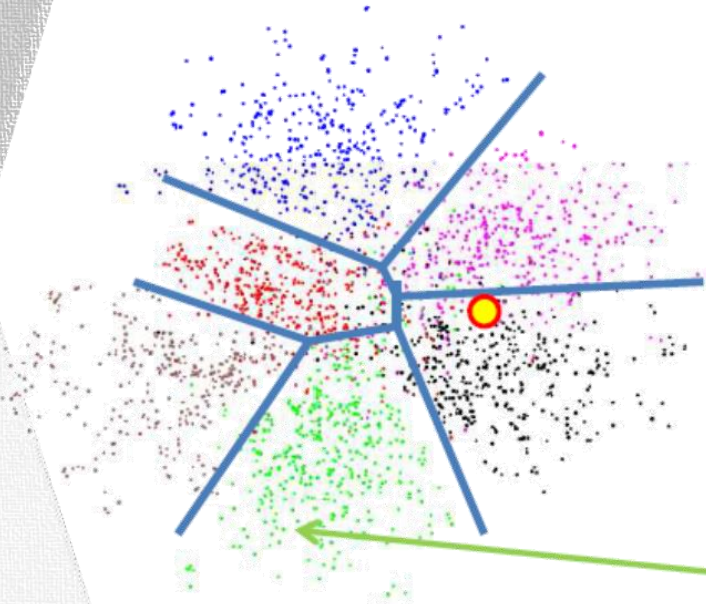
Où le dénominateur est un terme de normalisation :

$$p(\mathbf{x}) = \sum_y p(\mathbf{x}|y)p(y)$$

Cette approche est dite générative car, connaissant $p(\mathbf{x}|y)$, il est facile de générer des données dans l'espace des paramètres

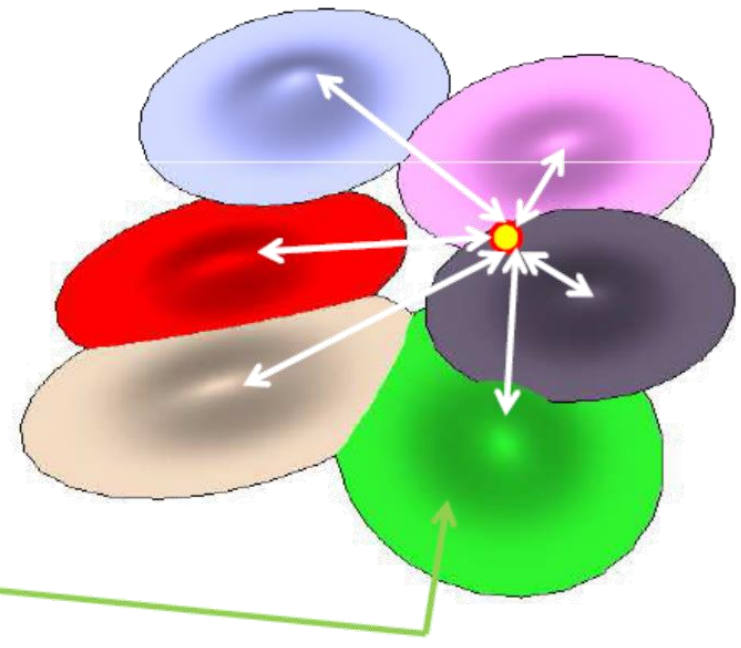
Approche discriminante

On apprend les frontières entre les classes



Approche générative

On modélise chaque classe





Exemple : on souhaite déterminer la langue parlée par une personne

Approche générative : on apprend chaque langage puis on détermine à quel langage la parole appartient (peut fonctionner avec une seule langue pour savoir si la personne parle français ou non)



Approche discriminative: on apprend les différences linguistiques entre les langages, sans apprendre le langage. Beaucoup plus simple !

Avantage/inconvénient

Approche générative

- 
- $p(\mathbf{x}|y)$ est estimée. On peut considérer $p(\mathbf{x}|y)$ comme la probabilité que \mathbf{x} soit bien modélisé par le modèle. Ceci permet de faire du rejet ou de combiner des classifieurs.
 - $p(\mathbf{x}|y)$ peut être utilisé pour générer des données
 - Permet à un système d'utiliser une seule classe. Ex : la teinte chaire
- 
- Trouver $p(\mathbf{x}|y)$ pour chaque classe est très coûteux en temps de calcul, surtout quand \mathbf{x} est de grande dimension
 - Nécessite une grande base de données, surtout quand \mathbf{x} est de grande dimension

Approche discriminative

- 
- Il est beaucoup plus rapide de déterminer $p(y|\mathbf{x})$ car la dimension de y est bien plus faible que celle de \mathbf{x}
- 
- On ne peut pas générer de données
 - On ne peut pas faire de rejet
 - Difficile de combiner des classifieurs

| Méthodes génératives | Méthodes discriminatives |
|----------------------------|----------------------------------|
| Classification bayésienne | K plus proches voisins |
| HMM (Hidden Markov Model) | Arbres de décision |
| Réseaux bayésiens | SVM (Support Vector Machine) |
| MRF (Markov Random Fields) | RVM (Relevance Vector Machine) |
| | Réseaux de neurones |
| | CRF (Conditional Random fields) |

Exemple : classification binaire

exemple issu de computer vision: models, learning and inference, Simon J.D. Prince 2012

On souhaite classifier des pixels en teinte chaire ou non chaire à partir de la quantité de rouge

x est une variable continue (quantité de rouge)

2 classes : teinte chaire $y=1$ ou non chaire $y=0$

Rappel sur la loi de Bernoulli

Pour x une variable binaire, $Bern(x = 1) = \mu$ et $Bern(x = 0) = 1 - \mu$ et de manière générale $Bern(x) = \mu^x (1 - \mu)^{1-x}$

Approche générative

- On modélise $p(x|y = 0)$ et $p(x|y = 1)$ par des gaussiennes $(\mu_0, \sigma_0, \mu_1, \sigma_1)$
- On modélise $p(y)$ par une loi de Bernoulli de paramètre μ . On utilise les données d'apprentissage (x_i, y_i) pour estimer les paramètres $(\mu_0, \sigma_0, \mu_1, \sigma_1, \mu)$
- On estime $p(y|x)$ en utilisant Bayes

Approche discriminative

- On modélise $p(y|x)$ par une loi de Bernoulli dont le paramètre μ dépend de x . Comme $0 < \mu < 1$, on pose $\mu = \frac{1}{1 + \exp(-\Phi_0 - \Phi_1 x)}$. On utilise les données d'apprentissage (x_i, y_i) pour estimer les paramètres (Φ_0, Φ_1) de $p(y|x)$ (2 paramètres)

Algorithmes des kppv

Il s'agit d'une **méthode supervisée**

Données de départ

Ensemble d'exemples de dimensions n étiquetés $y \in [1, \dots, K]$

But

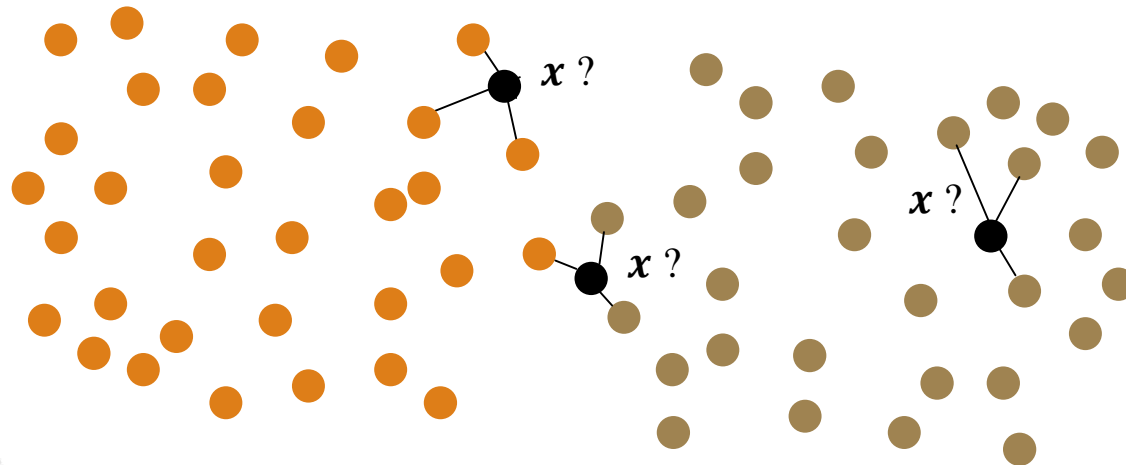
Trouver la classe d'un nouvel exemple \mathbf{x}

Méthode

- Calculer la distance entre \mathbf{x} et tous les exemples de la base de référence
- Déterminer les k exemples les plus proches.
- Affecter à \mathbf{x} la classe majoritaire parmi les k plus proches voisins

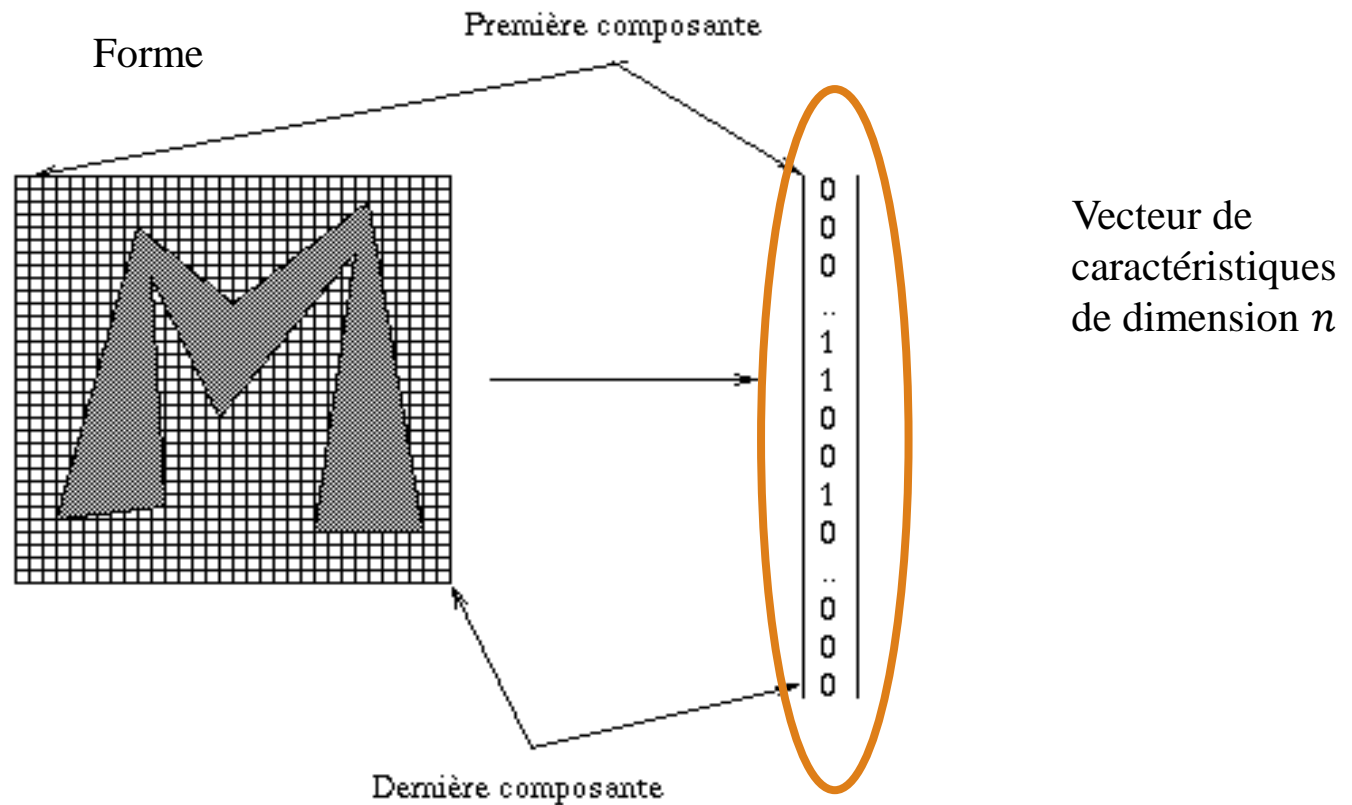
Exemple

En dimension 2, chaque exemple est caractérisé par un vecteur de dimension 2 et est donc représenté dans le plan :

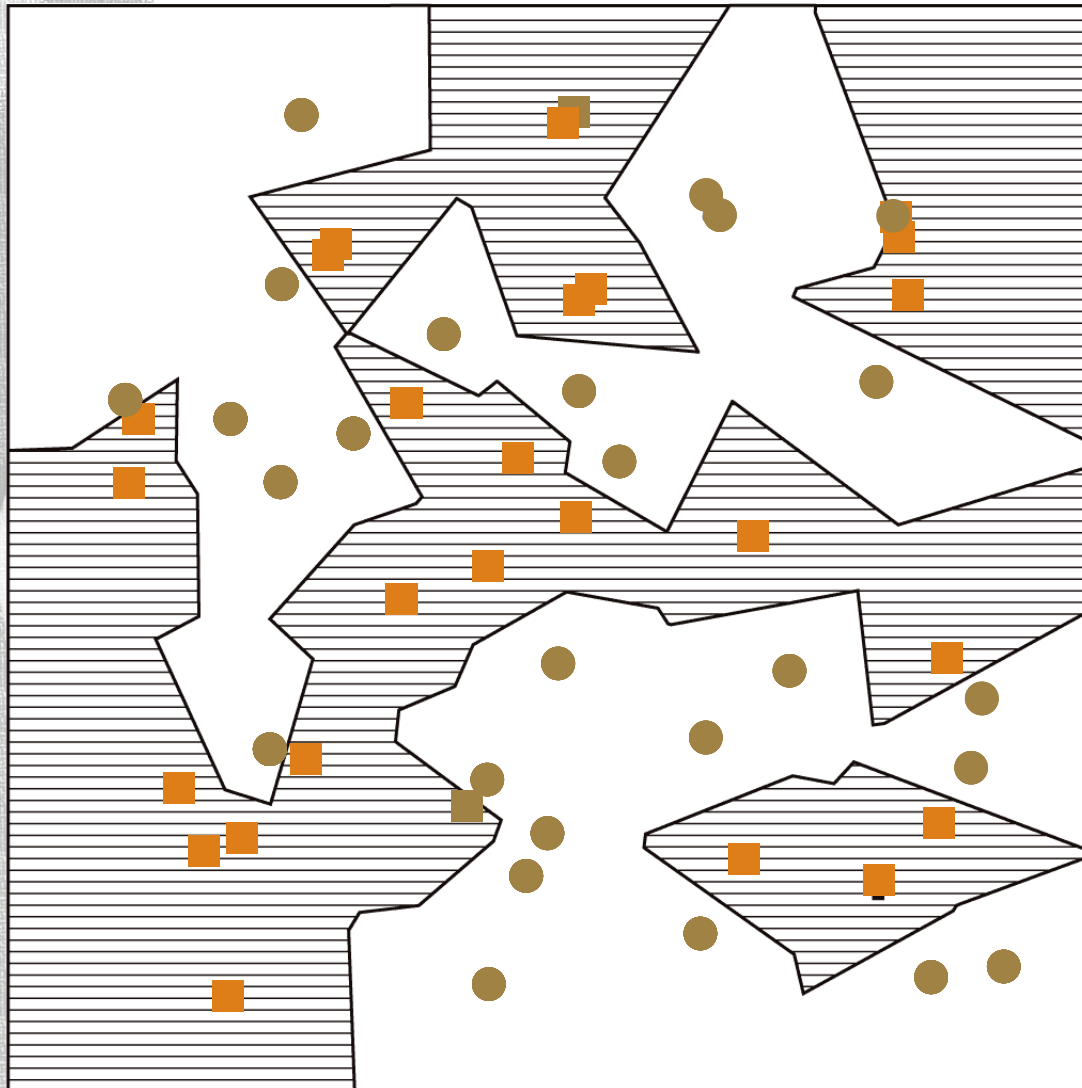


Exemple

En dimension n



Signification géométrique 1ppv

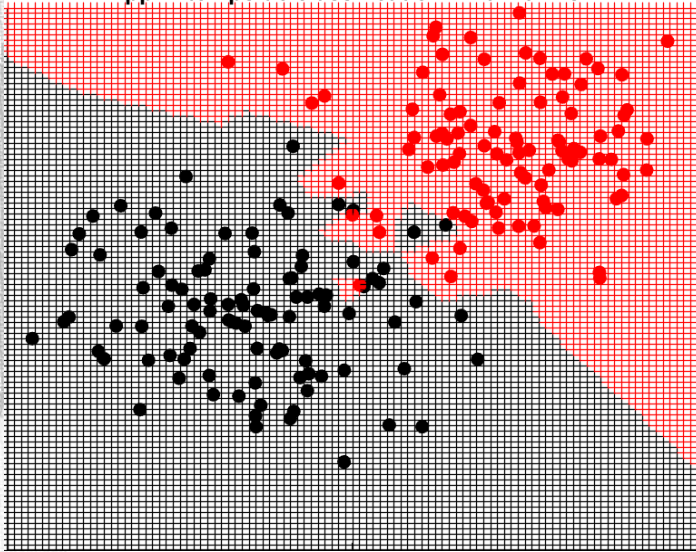


Les classes sont définies par la réunion des domaines d'influence des exemples de références

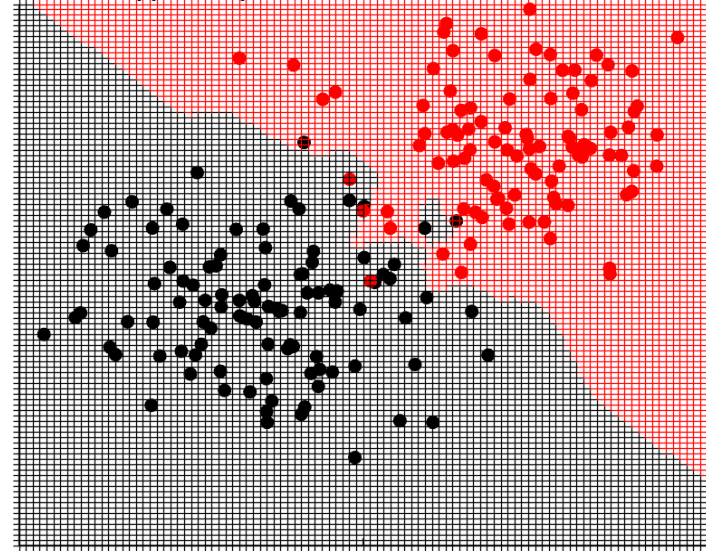
La résolution spatiale des frontières est liée au nombre d'exemples et à leur densité

Méthode discriminative : kppv

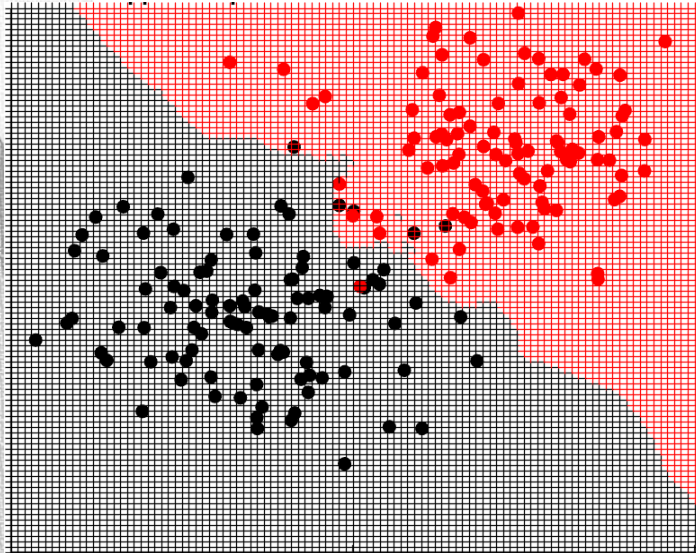
k=1



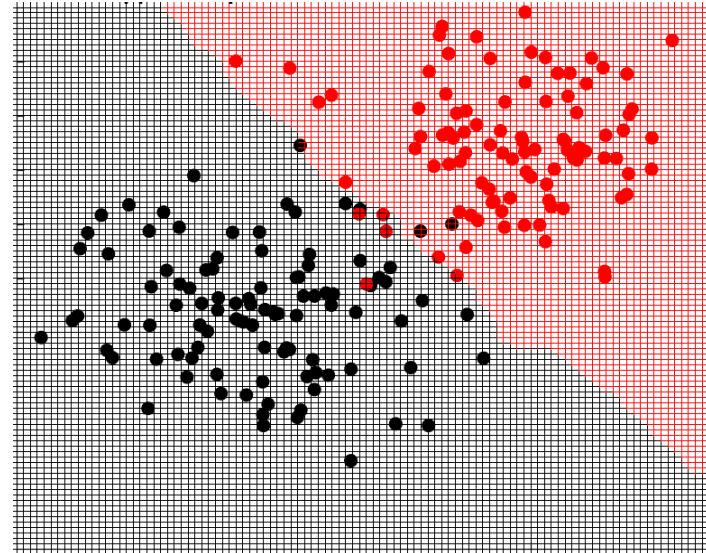
k=3



k=5



k=11



Dilemme biais/variance

k faible

- ➔ Bonne résolution des frontières entre classe
- ➔ Très sensible aux échantillons de la base de référence
- ➔ Petit biais Grande variance

k grand

- ➔ Mauvaise résolution des frontières entre classe : lissage des frontières
- ➔ Peu sensible aux échantillons de la base de référence
- ➔ Grand biais Faible variance

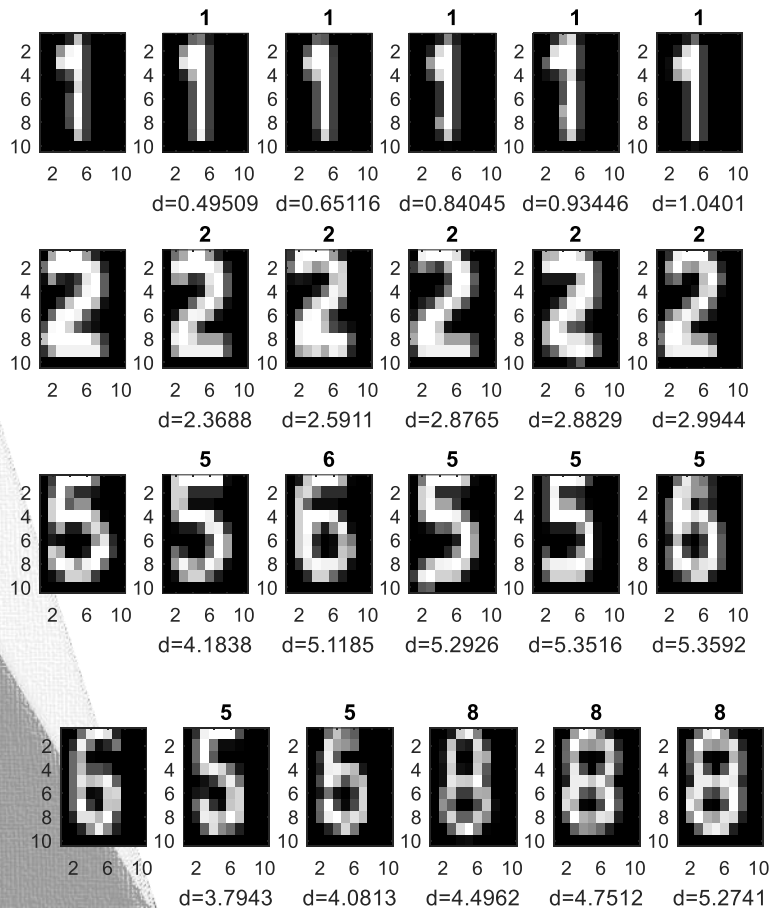
Comment choisir k ?

3 bases de données

- ➔ Base de référence où sont stockés les exemples
- ➔ Base de validation utilisée pour optimiser k
- ➔ Base de test pour évaluer les performances

Exemple en reconnaissance de caractères

- 800 exemples dans la base de références
- Chaque exemple est de dimension 100 (codage rétinien 10x10)



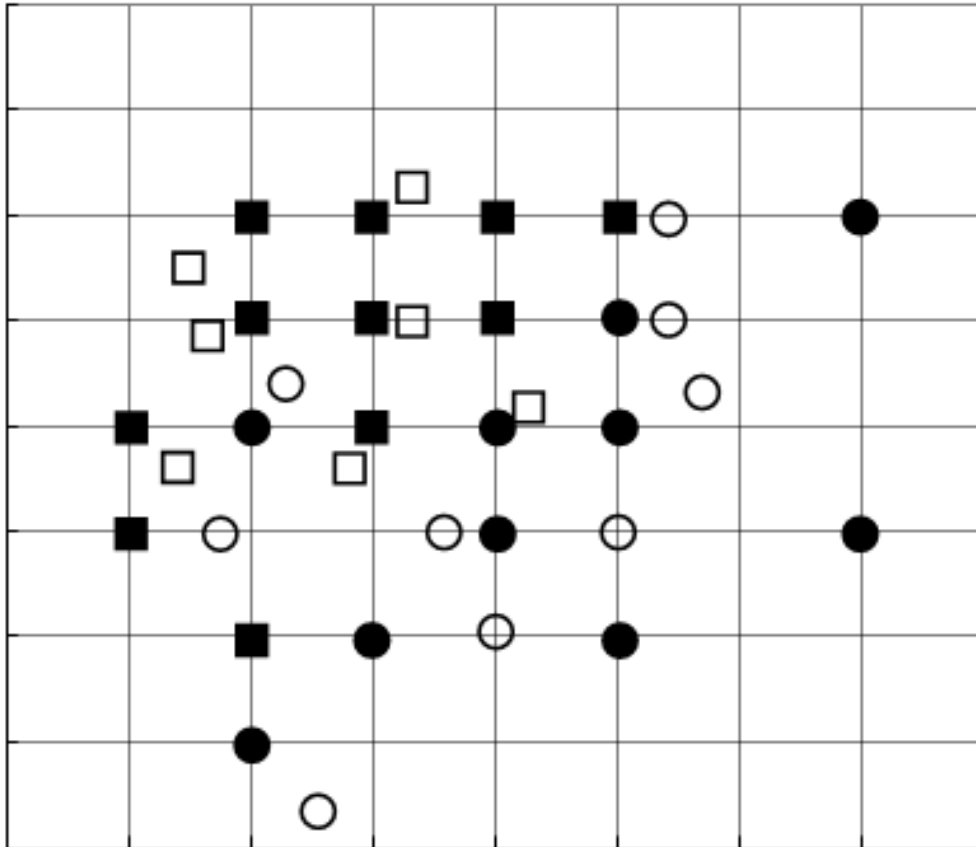
Avantages :

- Pas d'hypothèses
- Simple à mettre en œuvre
- Incrémental

Inconvénients :

- Quantité de calculs quasi-proportionnelle au nombre d'exemples
- Pas d'extraction d'information utile

Exercice

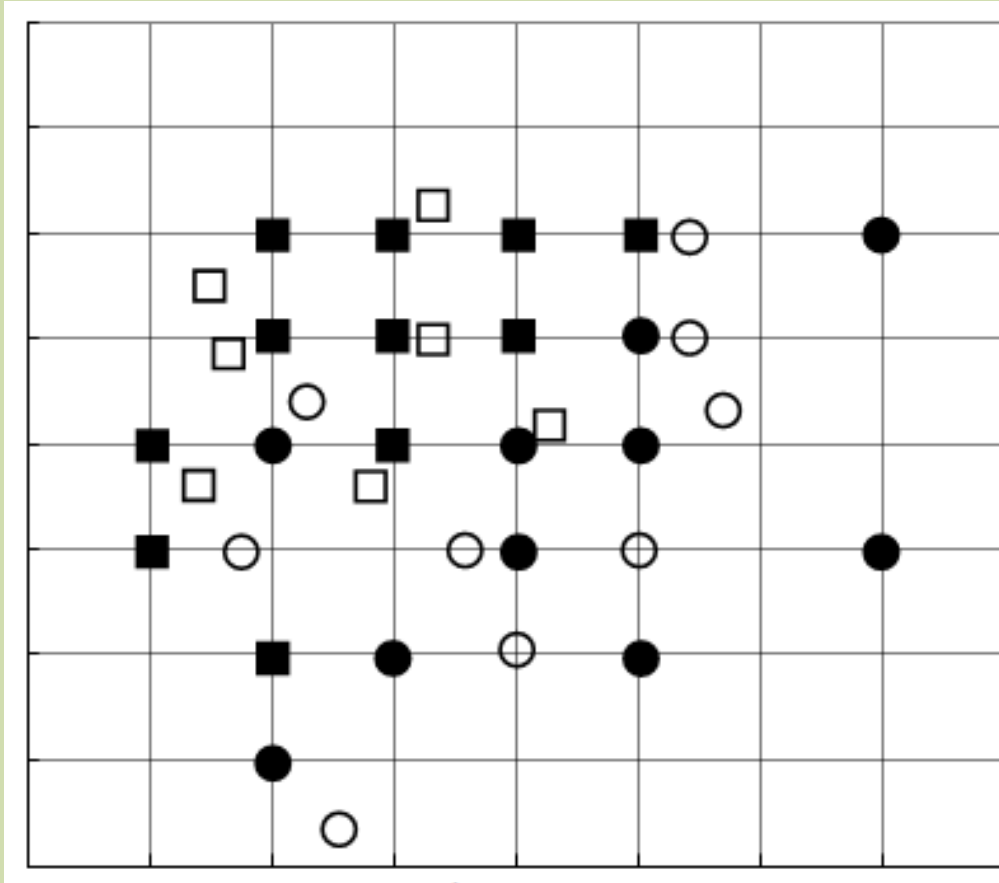


Noir : base de référence

Blanc : base de test

Donner la matrice de confusion avec
l'algorithme du 1ppv

Exercice



Noir : base de référence

Blanc : base de test

| | □ | ○ |
|---|---|---|
| □ | 6 | 1 |
| ○ | 2 | 7 |

Accélération des k-PPV

2 solutions

Réduction de la dimension de chaque exemple

- ACP
- LDA

Réduction de taille de la base de référence

- On ne représente plus chaque classe que par sa moyenne
- Génération de prototypes : LVQ



Dilemme robustesse/accélération

Représentation de chaque classe par sa moyenne μ_c

distances euclidiennes $d_E(x, \mu_c)$ entre l'exemple à classe x et les centres μ_c

$$d_E(x, c) = (x - \mu_c)^T (x - \mu_c)$$

→ L'exemple est classé à la classe de la distance la plus faible

Représentation de chaque classe par sa moyenne μ_k et sa matrice de de covariance Σ_c

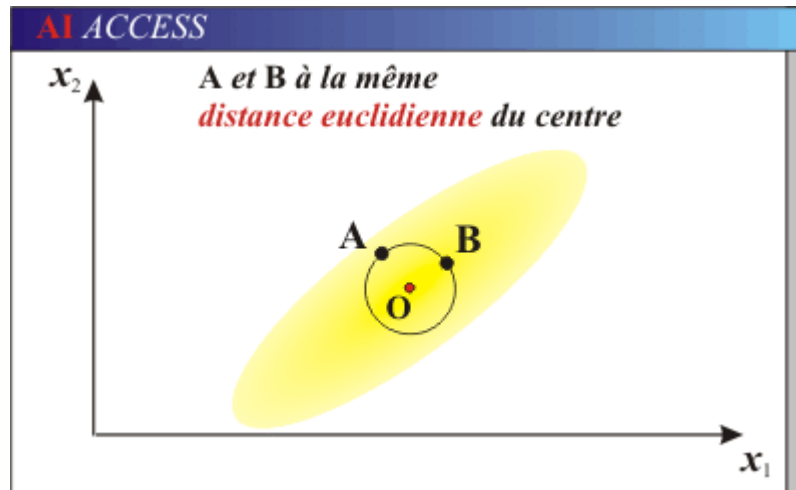
→ **distances de Mahalanobis $d_M(x, c)$ entre l'exemple à classe x et les centres μ_c**

$$d_M(x, c) = (x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c)$$

→ L'exemple est classé à la classe de la distance la plus faible

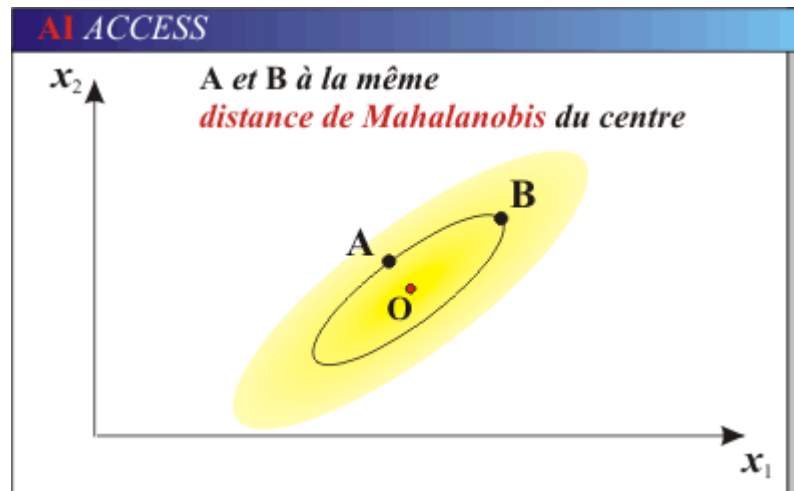
Rq : Si Σ =Identité, on retrouve la distance euclidienne

Comparaison – distance euclidienne – distance de Mahalanobis



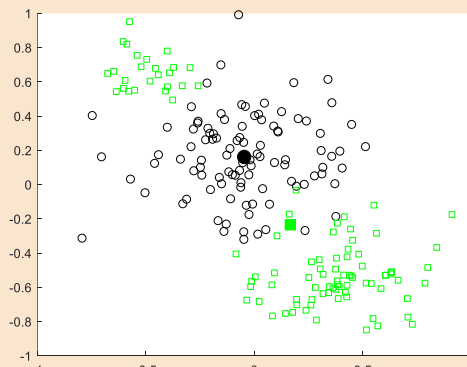
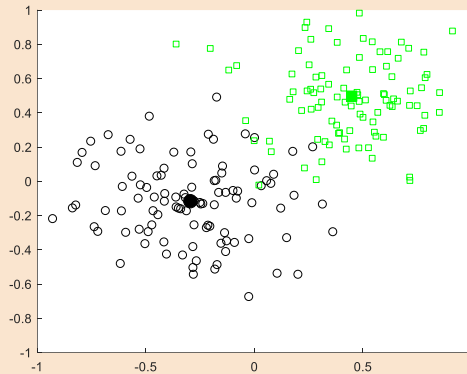
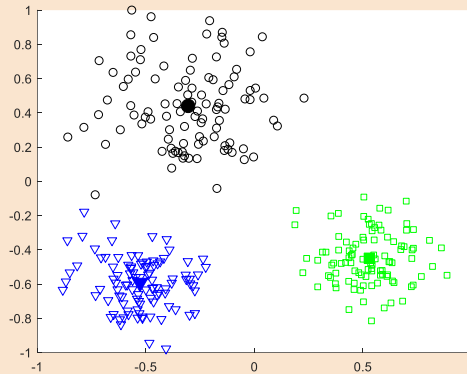
Les deux points A et B sont à la même distance euclidienne de O

→ Pas logique

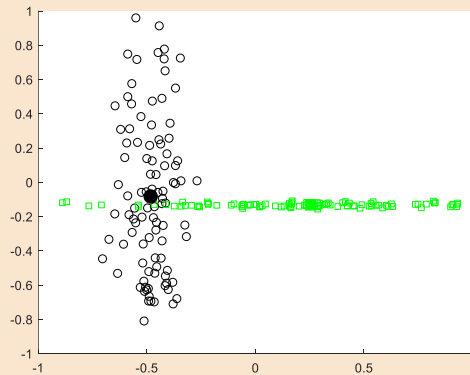
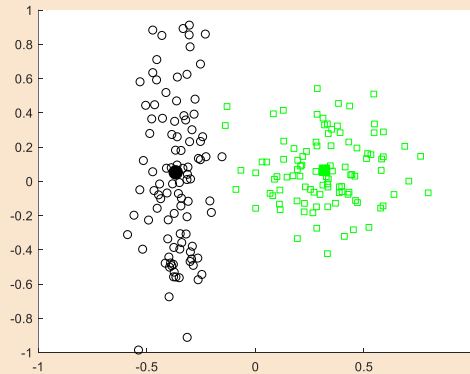


Les deux points A et B sont à la même distance de Mahalanobis de O

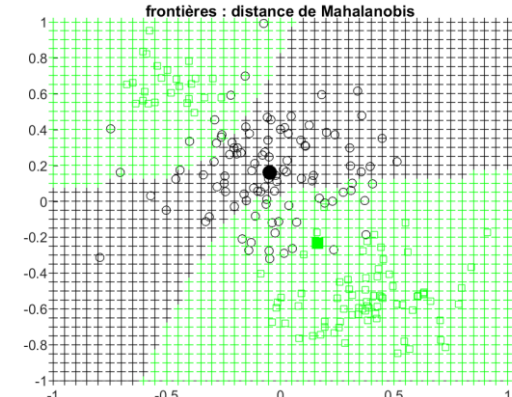
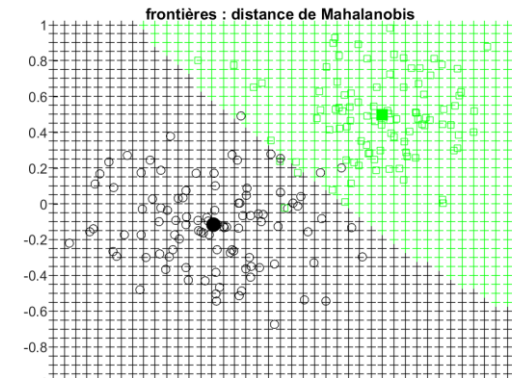
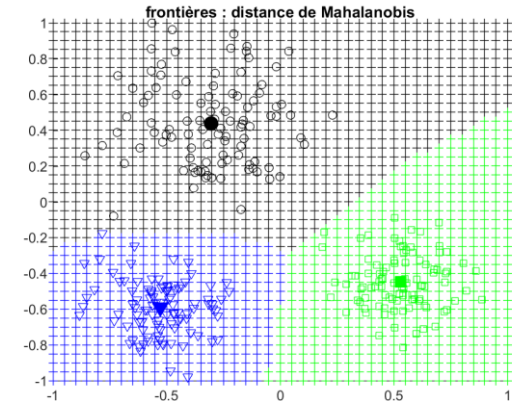
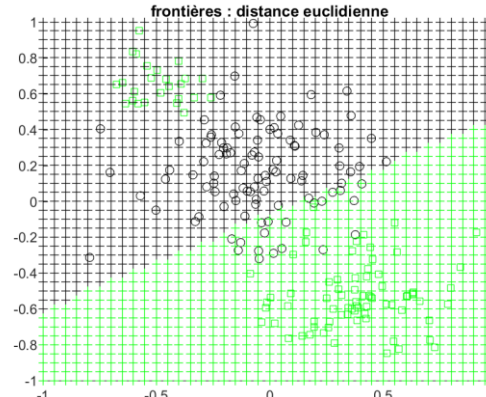
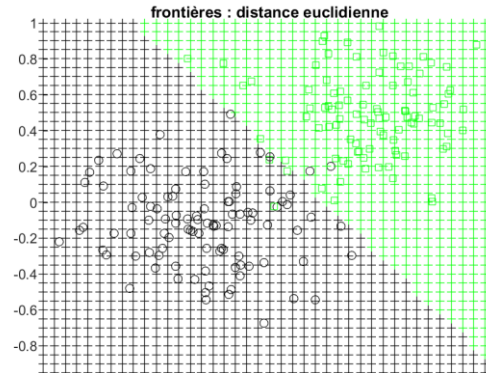
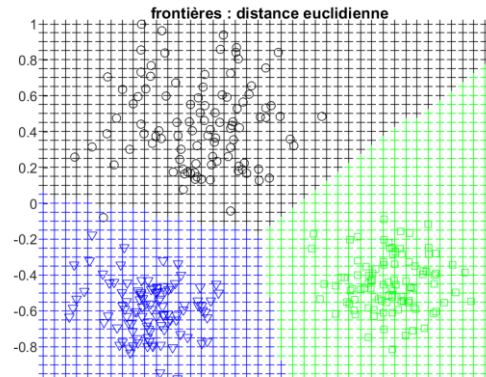
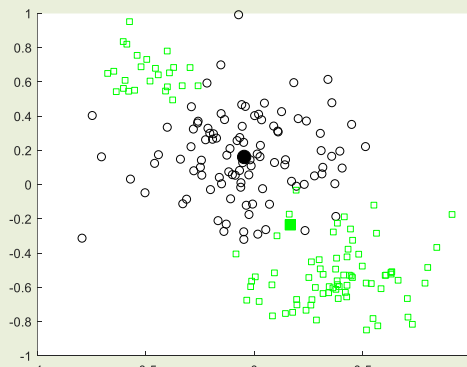
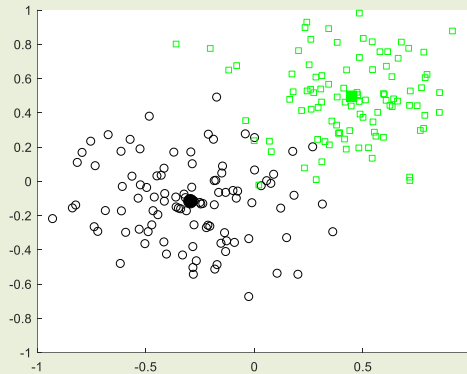
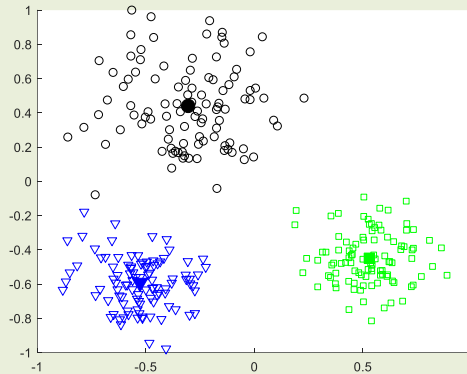
Exercice : tracer les frontières entre classes avec la distance euclidienne et la distance de Mahalanobis



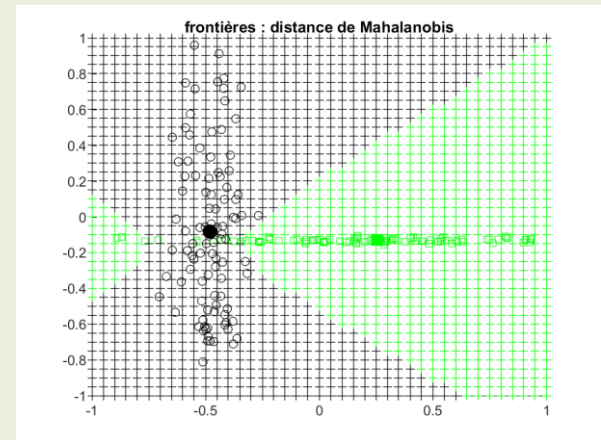
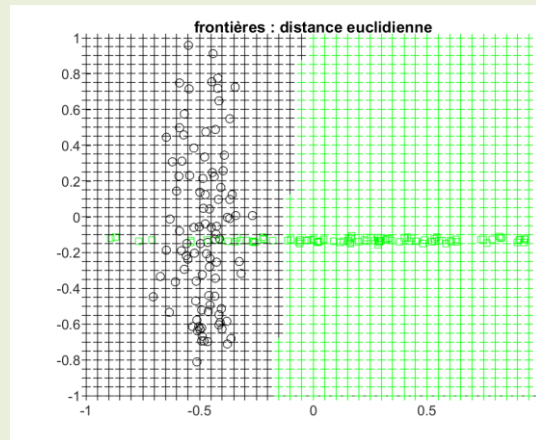
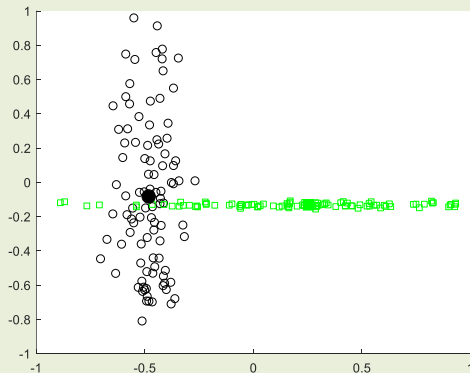
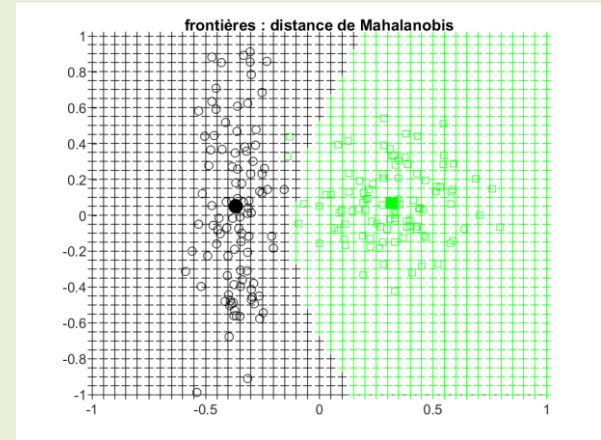
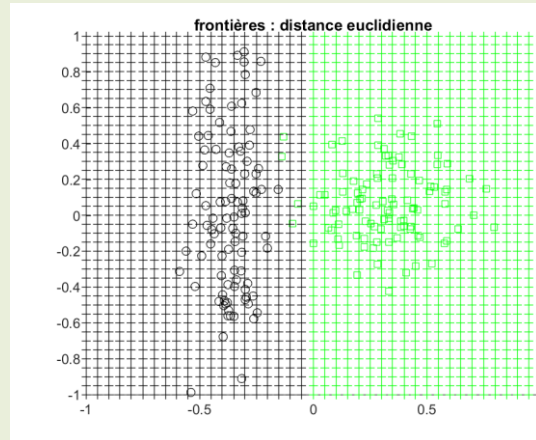
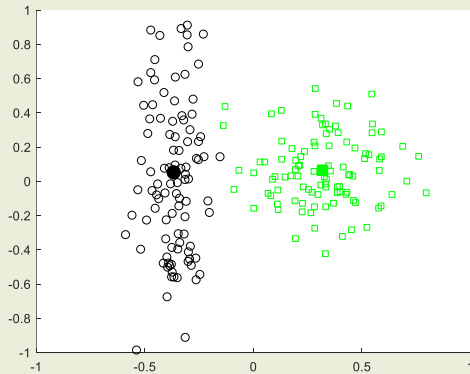
Exercice : tracer les frontières entre classes avec la distance euclidienne et la distance de Mahalanobis



Exercice : tracer les frontières entre classes avec la distance euclidienne et la distance de Mahalanobis



Exercice : tracer les frontières entre classes avec la distance euclidienne et la distance de Mahalanobis



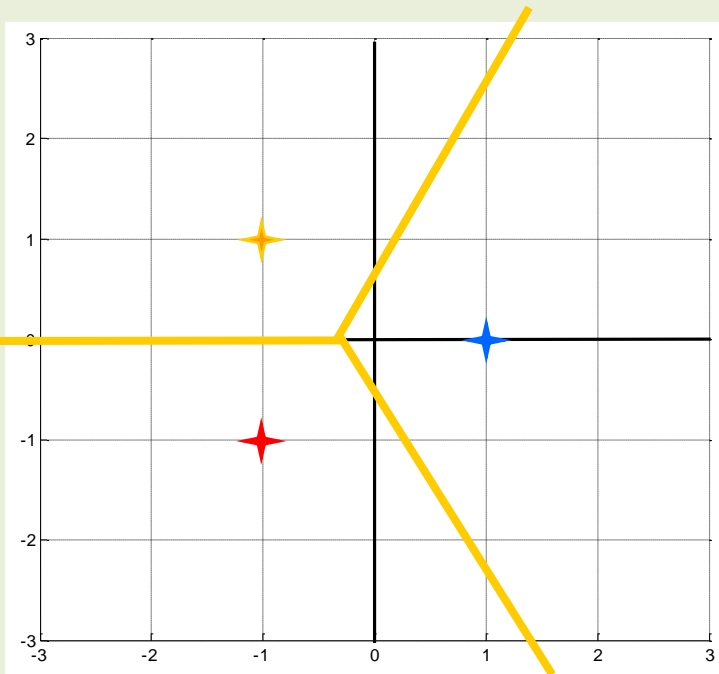
Exercice:

Supposons que l'on ait un problème à 3 classes, de dimension 2. Sur la base de référence, on a estimé:

$$\begin{aligned} m1 &= \begin{pmatrix} -1 \\ 1 \end{pmatrix} & m2 &= \begin{pmatrix} -1 \\ -1 \end{pmatrix} & m3 &= \begin{pmatrix} 1 \\ 0 \end{pmatrix} \\ \Sigma1 &= \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} & \Sigma2 &= \begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix} & \Sigma3 &= \begin{pmatrix} 4 & 0 \\ 0 & 4 \end{pmatrix} \end{aligned}$$

1. Tracer les frontières entre les classes en utilisant l'algorithme *nearest-mean* et la distance euclidienne.
2. Tracer les frontières approximatives entre les classes en utilisant l'algorithme *nearest-mean* et la distance de Mahalanobis .

Distance euclidienne



Distance de Mahalanobis

