

K-plus proche Voisin (KNN/KPP)

L'apprentissage automatique (2016-2017)

UFR MIME

Université Lille 3

7 décembre 2016

Sommaire

K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation

- 1 Classification
- 2 K-plus proche voisin (KPPV)
- 3 Généralisation et Évaluation

Pourquoi Faire l'Apprentissage

K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation

Apprentissage

Au lieu de programmer les règles manuellement dans un programme, donner une ordinateur une moyenne de extraire les règles automatiquement.

Pourquoi

- Problèmes trop complexes
- Travail manuel trop coûteux
- Très grandes quantités de données
- Pour devenir plus efficace/efficace avec ces tâches
- C'est rigolo :)

Classification

K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation

Apprentissage

Au lieu de programmer les règles manuellement dans un programme, donner une ordinateur une moyenne de extraire les règles automatiquement.

Apprentissage Supervisé : Classification

- La classification consiste a prédire **une catégorie**
- On va se limiter a un sous-classe de problème :
classification binaire
- Fournir au système en entrée un ensemble de d'exemples étiquetés \mathbf{x}_i, y_i d'apprentissage.
- $\mathbf{x}_i \in \mathbb{R}^d$ et $y_i \in \{C_1, C_2, \dots, C_q\}$ avec q classes
- Pour classification binaire y_i est un catégorie soit positif / négatif dans $\{+1, -1\}$

Classification

K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation

Classification Binaire

- Tache pour nous : Retrouver une mapping (correspondance) $f_{\theta} : \mathbb{R}^d \rightarrow \{-1, +1\}$
- Algo permet ordinateur de se programmer **lui-même**.
- Ici L'algorithme du coup consiste a retrouver les paramétrés optimaux θ
- Le classifier finale vas nous aider a prédire les taches comme si :
 - **Oui(+1)/Non(-1)** Une visage apparaît dans une image ?
Entree : ensemble des pixels d'image
 - **Oui(+1)/Non(-1)** Un document parle de sport ?
Entree : Les suites des characters de texte
 - **Oui(+1)/Non(-1)** Un client risque de quitter ma banque ?
Entree : Attributs de client (age, salaire, épargne, type maison, type contrat etc)

K-plus proche Voisin (KPPV) K-Nearest neighbours (KNN)

K-plus proche
Voisin
(KNN/KPPV)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation

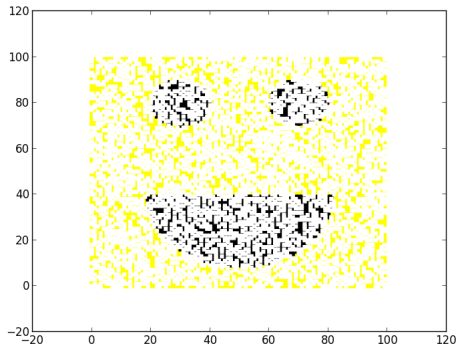


FIGURE – Comment prédire les pixels manquant ?

K-plus proche Voisin K-Nearest neighbours (KNN)

K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation

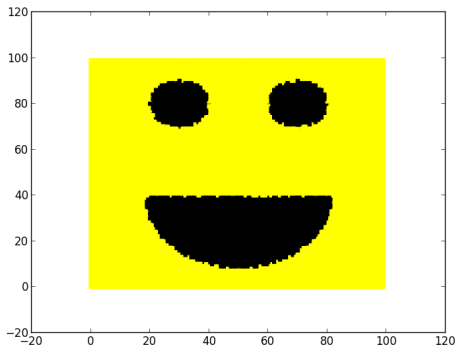


FIGURE – Utiliser les voisins de pixel a remplir pour décider la libelle (noir(+1)/jaune(-1)).

Exemple : dimension $d = 2$

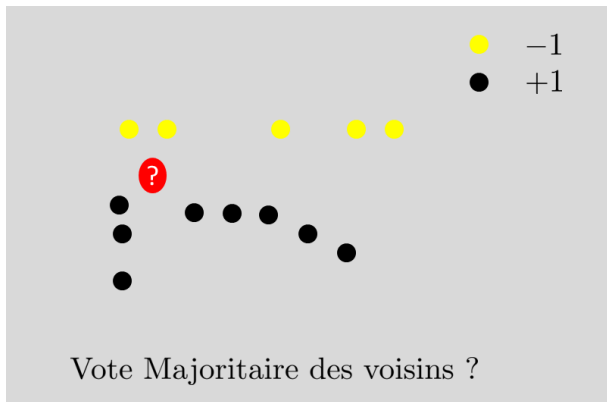
K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation



Exemple : dimension $d = 2$

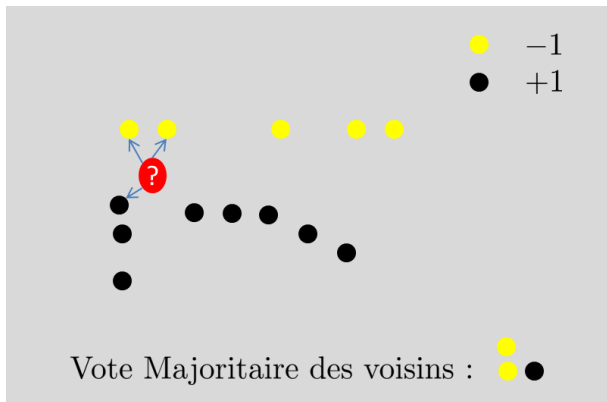
K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation



Exemple : dimension $d = 2$

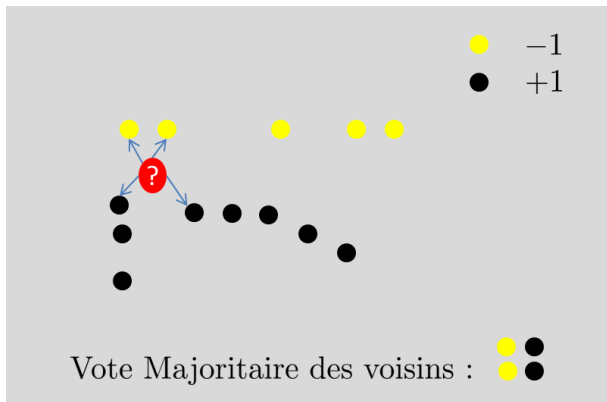
K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation



Exemple : dimension $d = 2$

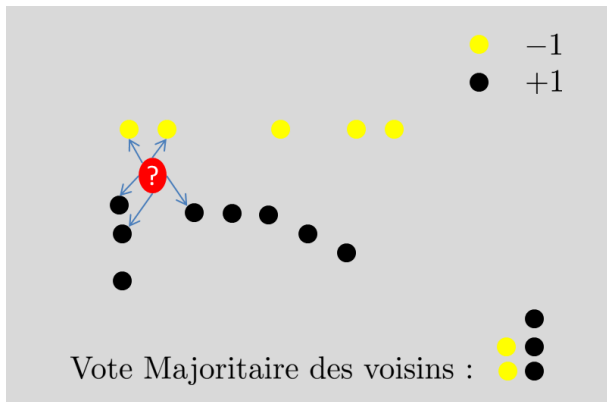
K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation



Exemple : dimension = 2

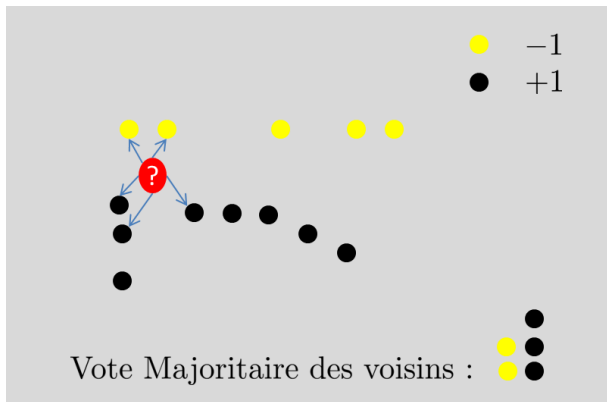
K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation



Voisinage

En réglant la **distance** → *nombre de voisins* on arrive à recouvrir plus des voisins qui change le résultat de **classification**.

Algorithme

- Ensemble d'apprentissage (ou Training set) :

$$X_{Tr} := \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$$
- $\mathbf{x}_i \in \mathbb{R}^d$ et $y_i \in \{C_q\}_{q=1}^Q$
- Pour la classification binaire $y_i \in \{-1, +1\}$
- Pour une nouvelle entrée \mathbf{z} (dans le ensemble test)
- $f_{\theta}(\mathbf{z}) = \text{VoteMajoritaire}\{y_i | i \in \text{k-plus proche voisin}(\mathbf{z})\}$
dans X_{Tr}
- Pour calculer la k-plus proche voisin de \mathbf{x} :
 - Pour calculer distance euclidien entres deux point :
 - $d(\mathbf{x}, \mathbf{x}') = \sqrt{(x_1 - x'_1)^2 + (x_2 - x'_2)^2 + \dots + (x_d - x'_d)^2}$
- $P(y = +1) = \text{Proportion des pts de classe } +1 \text{ dans le voisinage } N_k(\mathbf{z})$

Matrice de distance

K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}
x_2	1.5																	
x_3	1.4	1.6																
x_4	1.6	1.4	1.3															
x_5	1.7	1.4	1.5	1.5														
x_6	1.3	1.4	1.4	1.5	1.4													
x_7	1.6	1.3	1.4	1.4	1.5	1.8												
x_8	1.5	1.4	1.6	1.3	1.7	1.6	1.4											
x_9	1.4	1.3	1.4	1.5	1.2	1.4	1.3	1.5										
x_{10}	2.3	2.4	2.5	2.3	2.6	2.7	2.8	2.7	3.1									
x_{11}	2.9	2.8	2.9	3.0	2.9	3.1	2.9	3.1	3.0	1.5								
x_{12}	3.2	3.3	3.2	3.1	3.3	3.4	3.3	3.4	3.5	3.3	1.6							
x_{13}	3.3	3.4	3.2	3.2	3.3	3.4	3.2	3.3	3.5	3.6	1.4	1.7						
x_{14}	3.4	3.2	3.5	3.4	3.7	3.5	3.6	3.3	3.5	3.6	1.5	1.8	0.5					
x_{15}	4.2	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	1.7	1.6	0.3	0.5				
x_{16}	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1	1.6	1.5	0.4	0.5	0.4			
x_{17}	5.9	6.2	6.2	5.8	6.1	6.0	6.1	5.9	5.8	6.0	2.3	2.3	2.5	2.3	2.4	2.5		
x_{18}	6.1	6.3	6.2	5.8	6.1	6.0	6.1	5.9	5.8	6.0	3.1	2.7	2.6	2.3	2.5	2.6	3.0	
x_{19}	6.0	6.1	6.2	5.8	6.1	6.0	6.1	5.9	5.8	6.0	3.0	2.9	2.7	2.4	2.5	2.8	3.1	0.4

FIGURE – Pour $k = 4$ on cherche les plus proche voisin pour le vecteur x_{11} . On retrouve x_{10} (classe bleu), x_{12} , x_{13} , x_{14} (classe rouge)

Jeu de données synthétique avec deux classes

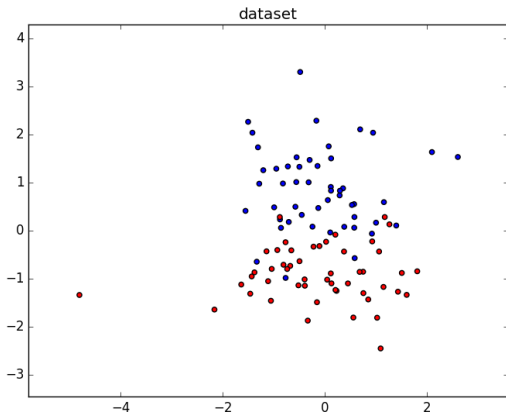
K-plus proche
Voisin
(KNN/KPP)

Lille 3

Classification

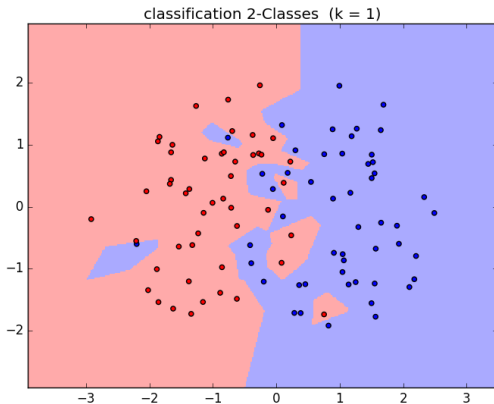
K-plus proche
voisin (KPPV)

Généralisation
et Évaluation



Exemple avec frontières de décision

$k=1$, trop compliqué



Exemple avec frontières de décision

K-plus proche
Voisin
(KNN/KPP)

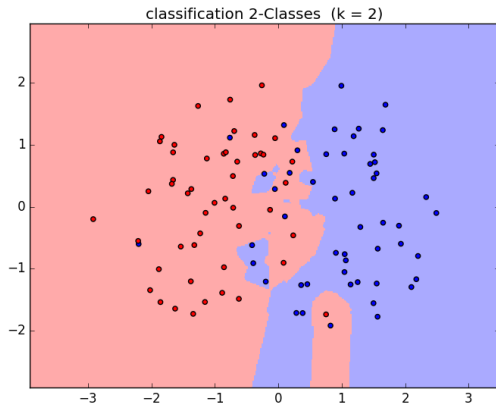
Lille 3

Classification

K-plus proche
voisin (KPPV)

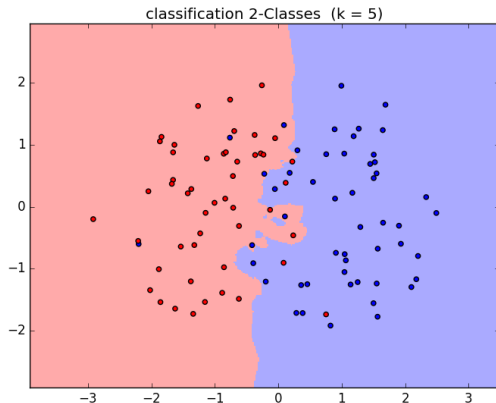
Généralisation
et Évaluation

$k=2$, compliqué



Exemple avec frontières de décision

k=5, OK !



Exemple avec frontières de décision

K-plus proche
Voisin
(KNN/KPP)

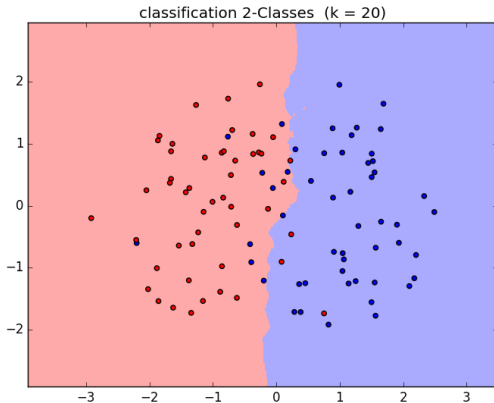
Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation

$k=20$, pas mal



Exemple avec frontières de décision

K-plus proche
Voisin
(KNN/KPP)

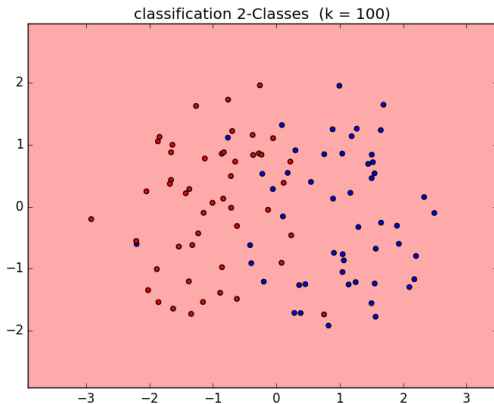
Lille 3

Classification

K-plus proche
voisin (KPPV)

Généralisation
et Évaluation

$k=100$, trop simplifié



Sur-Sous Apprentissage

- f_{θ} est trop compliqué comme frontière des décision :
sur-apprentissage
 - Il prédit très bien sur donnée d'apprentissage mais pas très bien sur les nouvelles données
- f_{θ} est trop simple comme frontière des décision :
sous-apprentissage
 - Il prédit ni bien sur donnée d'apprentissage ni sur les nouvelles données
 - Ceci est possible si k est très grandes.
 - Il prédit le vote majoritaire "globale"
- Compromis Simple-Compiqué (Biais-Variance)
 - Choix de paramétré optimale θ : ici nombre de voisins K

Performance

Comment savoir si notre classifieur va bien généraliser ? → vas bien prédire sur les nouvelles données ?

Erreur de classification

- Étant donné un scénario d'apprentissage :
 - On peut utiliser l'ensemble d'apprentissage pour l'évaluation d'erreur ?
 - Il nous faut des nouveaux exemples ? Pas déjà pris en compte par le classifieur ? Pourquoi ?
 - Quelle est l'erreur des prédictions pour $k=1$ pour KPPV sur l'ensemble d'apprentissage ?

Split/Decoupage

- $D = \text{Donnée apprentissage} \cup \text{donnée test}$
- $\text{Donnée apprentissage} \cap \text{donnée test} = \emptyset$
- Erreur est évalué sur l'ensemble test en comparant
 - Catégorie prédit par le classifieur (classification)
 - Le vrai valeur de catégorie (vérité terrain)

TP : Codez les différentes fonctions pour KPPV

- qui calcul de distance a pair pour l'ensemble d'apprentissage
- qui calcul le vote majoritaire étant données k (classification)
- fonction qui va évaluer les erreur de classification sur ensemble d'évaluation (test)