

# Les données numériques: toutes des nombres!

## Algorithme kNN

Diane Lingrand



2022 - 2023

1 kNN algorithm

2 Improvements

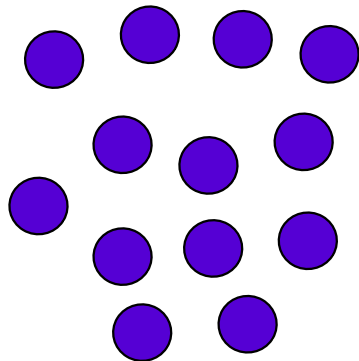
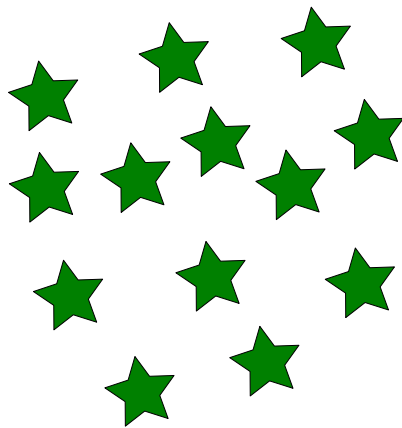
3 Experiments

1 kNN algorithm

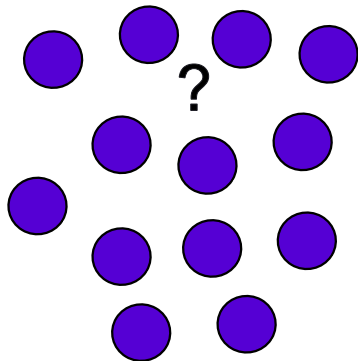
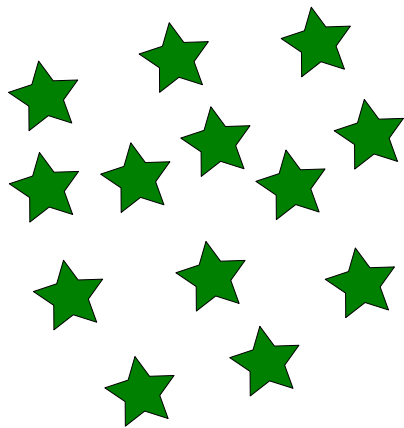
2 Improvements

3 Experiments

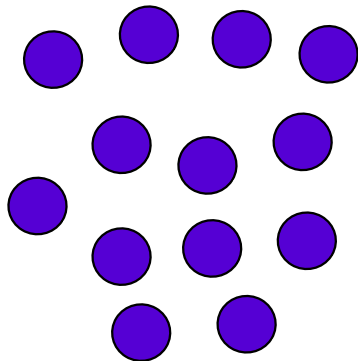
- training dataset
  - data  $x_i \in \mathbb{R}^d$  with  $0 \leq i \leq n$
  - labels or class index  $y_i \in \mathbb{R}$  with  $0 \leq i \leq n$ 
    - binary classification : positives and negatives
    - multi-class classification
- the goal is to be able to guess the class given the data
- metrics are computed on a test dataset



We observe that nearby objects belong to the same class.

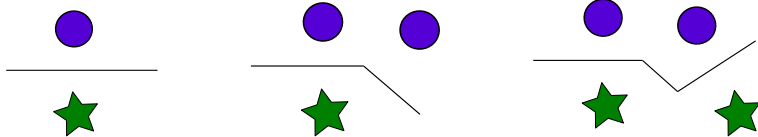


This new object will be assign to the class of purple rounds.



This new object will be assign to the class of green stars.

- intuition : nearest point
- Voronoi tessellation
  - points at same distance from two different points



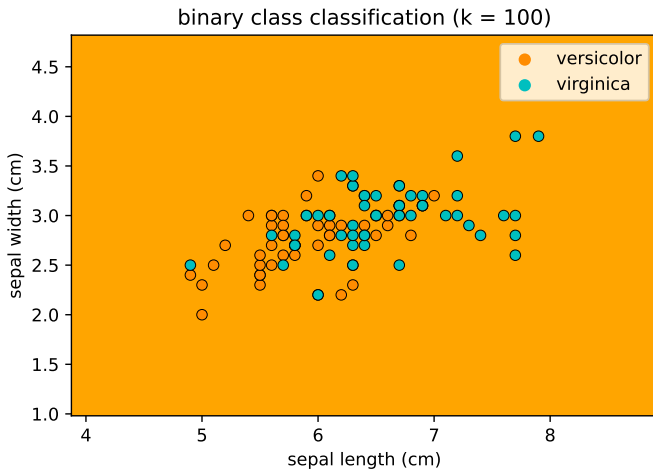
- complex decision boundary



- 1NN is very sensitive to outliers
  - idea : use more than one nearest neighbor to make decision
    - use  $k$  nearest neighbors
- algorithm
  - compute distance to every training sample  $x_i$
  - select  $k$  closest instances
  - output the class that is most frequent

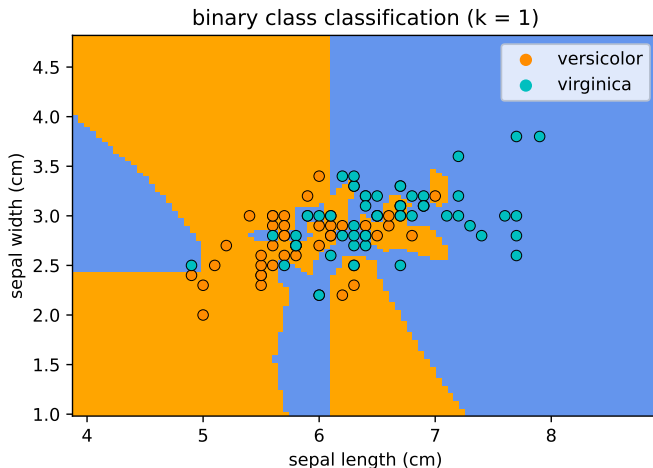
# How to choose value $k$ ?

- large value : everything classified as the most probable class



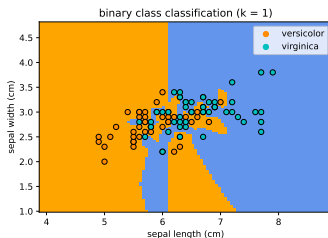
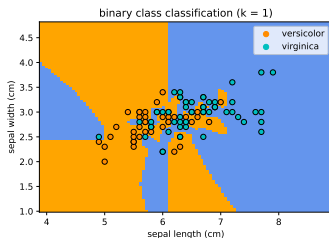
# How to choose value $k$ ?

- large value : everything classified as the most probable class
- small value : highly variable, unstable decision boundaries



# How to choose value $k$ ?

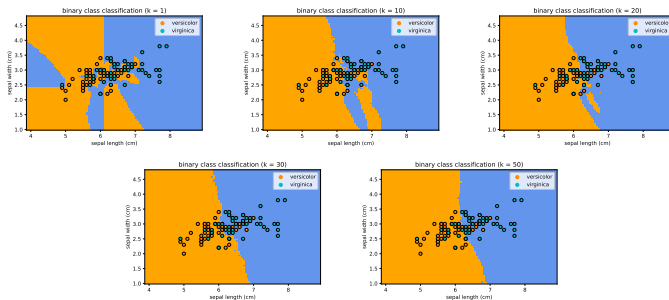
- large value : everything classified as the most probable class
- small value : highly variable, unstable decision boundaries



- small changes in the training set imply large changes in classification

# How to choose value $k$ ?

- large value : everything classified as the most probable class
- small value : highly variable, unstable decision boundaries
  - small changes in the training set imply large changes in classification
- affects smoothness of the boundary



# How to choose value $k$ ?

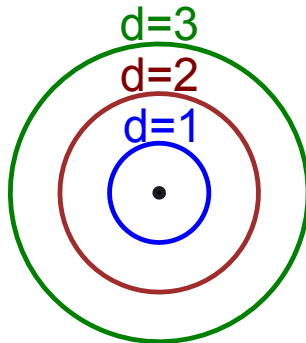
- large value : everything classified as the most probable class
- small value : highly variable, unstable decision boundaries
  - small changes in the training set imply large changes in classification
- affects smoothness of the boundary
- select value of  $k$  using validation dataset
  - it is critical to use the validation dataset
    - if we use the training dataset, 1NN is the best (the nearest example is itself and has the correct class/value)
- we thus need training / validation / test dataset

# Distance functions : key component

- Euclidean :

$$d_2(i, j) = \sqrt{\sum_{k=1}^m |x_{ik} - x_{jk}|^2}$$

- nice properties : symmetric, isotropic, ..

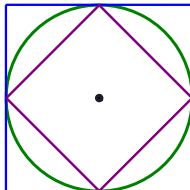


# Distance functions : key component

- Minkowski (generalisation)

$$d_q(i, j) = \sqrt[q]{\sum_{k=1}^m |x_{ik} - x_{jk}|^q}$$

- Manhattan ( $q = 1$ ) :  $d_1(i, j) = \sum_{k=1}^m |x_{ik} - x_{jk}|$
- Euclidean ( $q = 2$ ) :  $d_2(i, j) = \sqrt{\sum_{k=1}^m |x_{ik} - x_{jk}|^2}$
- max distance ( $q = \infty$ ) :  $d_\infty(i, j) = \max_k |x_{ik} - x_{jk}|$





- no majority class : equal number of neighbours for at least 2 classes
  - in case of binary classification : use odd  $k$
  - random between the equal classes
  - prior : pick class with greater prior
  - nearest : let 1-NN decide

- Parzen windows
  - instead of choosing the number of neighbours, choose the size of the neighborhood
  - or use all the points with a decay function for distances (close to kernel that converts distances to numbers)

- assumption : nearby (defined by distance fn) regions of space concern the same class
- almost no learning (except  $k$ )
- sensitive to class outliers
- sensitive to lots of irrelevant attributes
- computationnally expensive :
  - space : need to store all training samples
  - time : need to compute distances to all training samples
  - expensive at testing, no training time (which is bad)

1 kNN algorithm

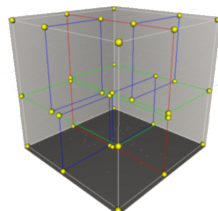
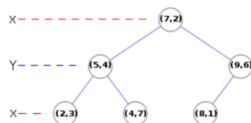
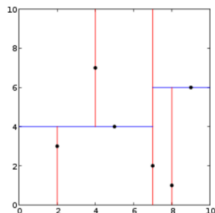
2 Improvements

3 Experiments

- time :  $O(nd)$  where  $n$  is the nb of samples and  $d$  the number of features
- reduce  $d$  (feature selection)
- reduce  $n$  :
  - idea : quickly find  $m \ll n$  potential near neighbors
  - compare only to those and pick  $k$  nearest neighbors  $O(md)$  time
  - Kd tree : low dim, real valued data
    - $O(d \log_2(n))$ , only if  $d \ll n$ , can miss neighbors
  - others methods we won't study in this class :
    - LSH (locality-sensitive hashing) for high dimensions
    - inverted lists
    - ball trees
    - ...

# Kd-trees (Friedman et al 1977)

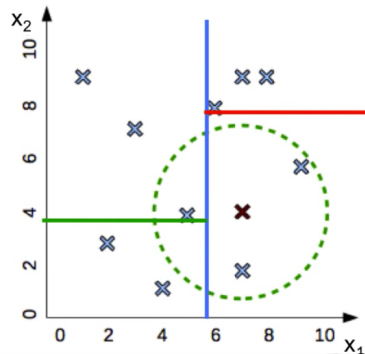
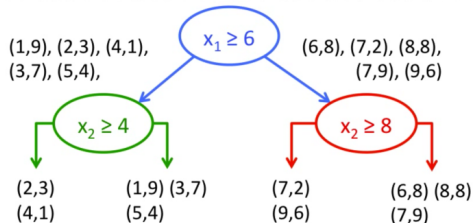
- $k$  corresponds to the dimension of the data ( $d$  in this course)
- kd-trees are binary trees
- designed to handle spatial data in a simple way



# Kd-tree : construction

- Repeat :
  - Pick random dimension
  - Find median element
  - Split data

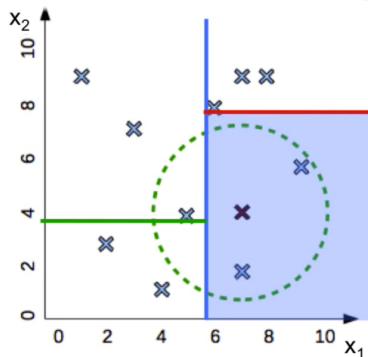
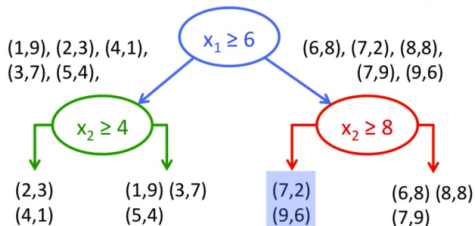
(1,9), (2,3), (4,1), (3,7), (5,4), (6,8), (7,2), (8,8), (7,9), (9,6)



# Kd-tree : Nearest Neighbor

- Example : find NNs for new point (7,4)
  - Find region containing (7,4)
  - Find kNN among all the points in this region

(1,9), (2,3), (4,1), (3,7), (5,4), (6,8), (7,2), (8,8), (7,9), (9,6)





1 kNN algorithm

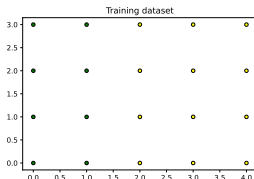
2 Improvements

3 Experiments

- Many implementations are available :
  - scikit-learn library
  - scipy library
- But today, you will write YOUR implementation, in `python`
  - brute force (original kNN), Euclidean distance
  - kd-tree

# Brute force implementation

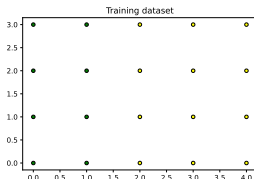
- start with a simple dataset



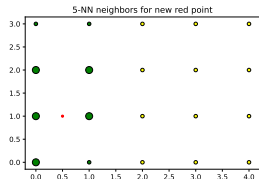
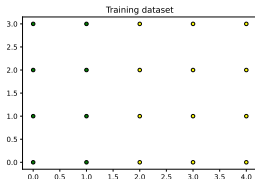
- $k$  in  $k$ -NN is a variable
- start by searching for the class of one single new point

# Brute force implementation

- start with a simple dataset

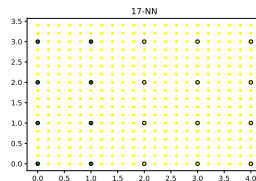
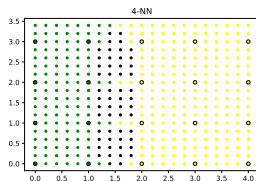
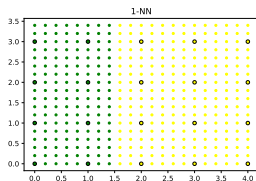


- $k$  in  $k$ -NN is a variable
- start by searching for the class of one single new point
  - list of  $k$  neighbors
  - vote for majority class
  - decide what to do in case of ambiguity



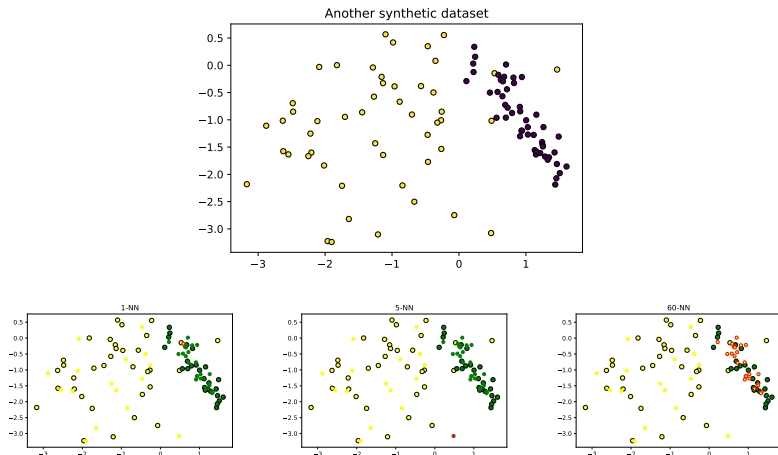
# Brute force implementation (2)

- start with a simple dataset
  - $k$  in  $k$ -NN is a variable
  - start by searching for the class of one single new point
  - **then predict class for a set of new points**
    - varying  $k$



# Brute force implementation (3)

- test on a more complex synthetic dataset



- Classification of 3 classes 'cat', 'dog' and 'bird' from the Speech Commands Dataset
  - 200 sounds per class for training, 200 other sounds per class for testing

