

Algoritmi de clasificare

1. Clasificarea

1. Arbori de decizie

2. Învățarea bazată pe instanțe

- Metoda celor mai apropiați k vecini

Clasificarea bazată pe instanțe

Lazy learners

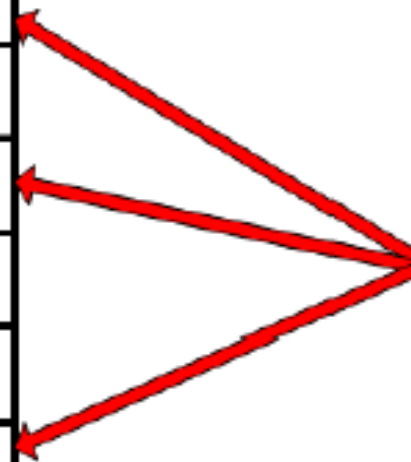
Instanțele memorate

Atr1	AtrN	Class
			A
			B
			B
			C
			A
			C
			B

- Se memorează efectiv instanțele de antrenare și se folosesc pentru a prezice clasele instanțelor noi

Instanță nouă

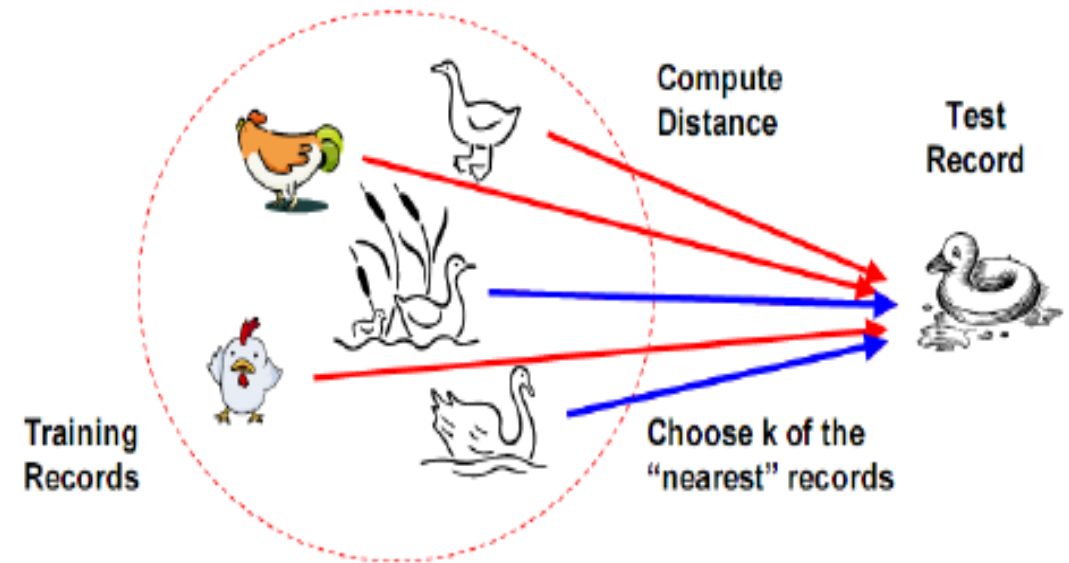
Atr1	AtrN



Metoda celor mai apropiați k vecini

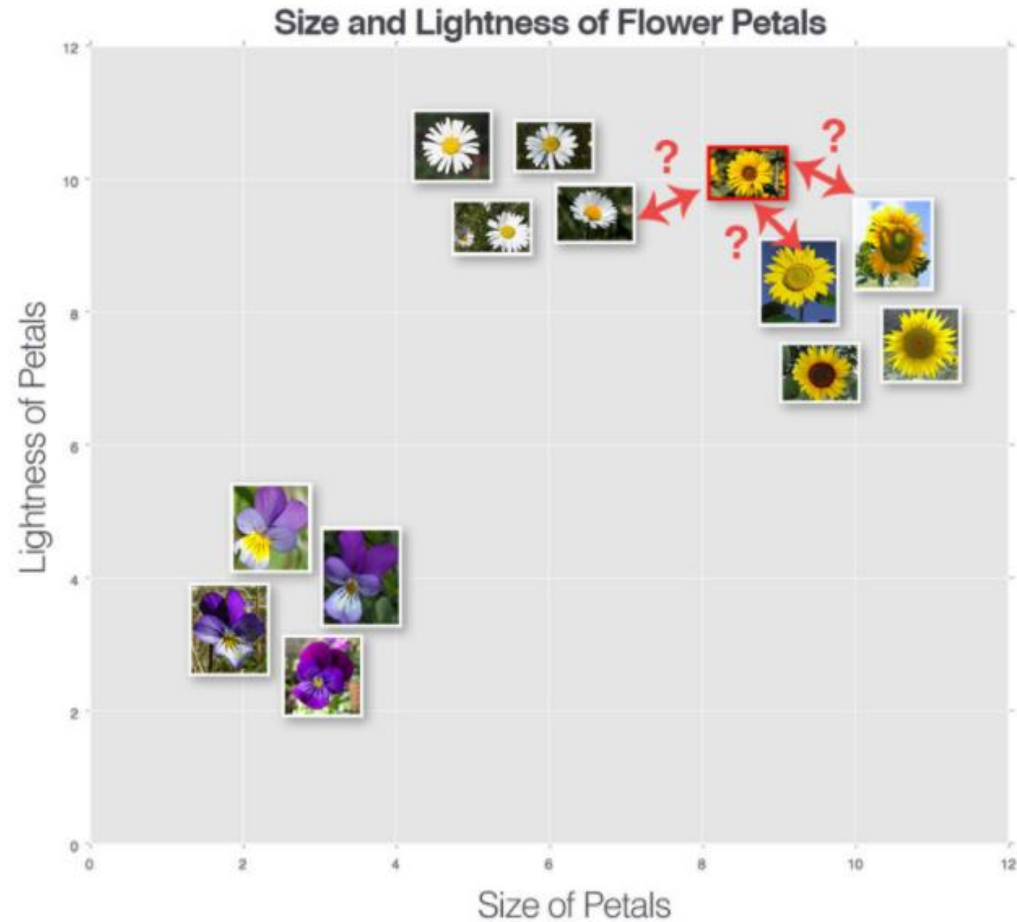
- *k-Nearest Neighbors*, prescurtat: k-NN
 - Se folosesc cele mai apropiate k instanțe pentru a realiza clasificarea

1. Cum se reprezintă instanțele
2. Cum se determină similitudinile
3. Cum se alege numărul de vecini
4. Cum se realizează clasificarea



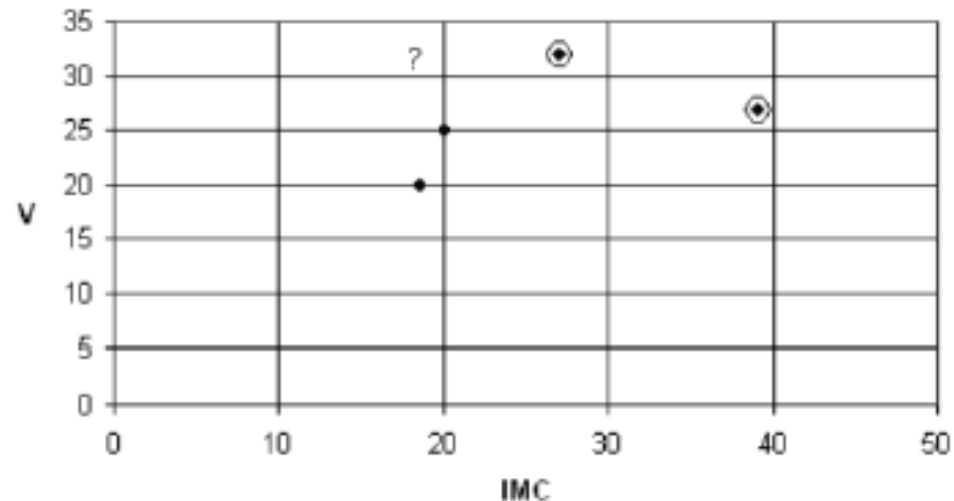
Metoda celor mai apropiați k vecini

Spune-mi cine îți sunt vecinii ca să-ți spun cine ești



1. Reprezentarea instanțelor

- Pentru n atribute, instanțele pot fi văzute ca puncte într-un spațiu n -dimensional
- De exemplu, clasificarea riscului unor pacienți
- Atribute:
 - Indicele masei corporale $IMC (= G / I^2)$
 - Vârsta V
- Instanțe:
 - $IMC = 18.5, V = 20$
 - $IMC = 27, V = 32$
 - $IMC = 39, V = 27$
 - $IMC = 20, V = 25$



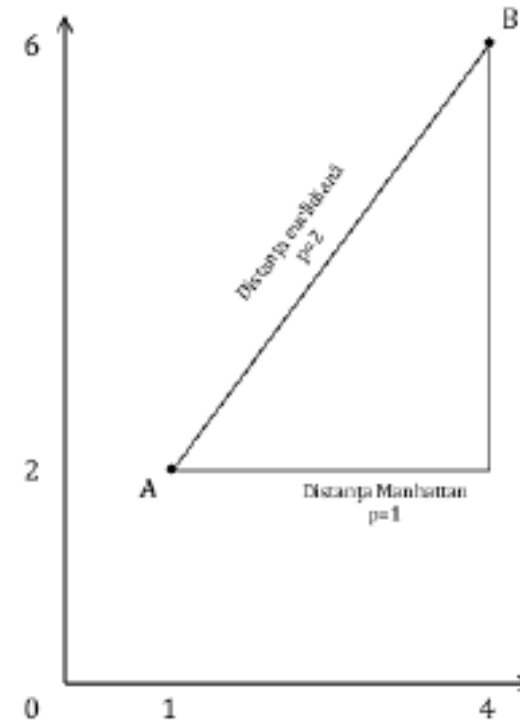
2. Calcularea similitudinii/distanței

Metrici pentru calcularea distanței

- Se folosesc în general particularizări ale distanței Minkowski

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$$

- Cele mai folosite metrici sunt:
 - Distanța euclidiană: $p = 2$
 - Distanța Manhattan: $p = 1$



2. Cum se calculeaza similitudinile

Distanța euclidiană

$$x = (x_1, \dots, x_n), \quad y = (y_1, \dots, y_n),$$

$$d_E(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Definim *distanța euclidiană* dintre $x, y \in \mathbb{R}^k$, astfel:

$$d(x, y) = \|x - y\| (= \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_k - y_k)^2}).$$

Observație. Pentru $k = 1$, $d(x, y) = \sqrt{(x - y)^2} = |x - y|$;

Pentru $k = 2$, $d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$;

Pentru $k = 3$, $d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + (x_3 - y_3)^2}$.

2. Calcularea similitudinii

Scalarea

Problemă legată de ordinul de mărime al datelor:

- avem două atribute: înălțimea și greutatea unor persoane
- înălțimea e măsurată în metri; intervalul poate fi de ex [1.50 m, 2.00 m], deci cu o diferență de maxim 0.5
- greutatea se măsoară în kilograme; intervalul poate fi [50 kg, 200 kg]
- diferențele de greutate domină pe cele în înălțime; o diferență de 1 kg este mai mare decât orice diferență de înălțime, contribuind deci prea mult la calculul distanței
- Soluție: scalarea mărimilor

2. Calcularea similitudinii

Scalarea

- Se recomandă scalarea atributelor pentru a preveni dominarea măsurii de distanță de către un anumit atribut
- De exemplu:
 - Înălțimea unei persoane $\in [1.5, 2.1]$ m
 - Greutatea unei persoane $\in [50, 120]$ kg
 - Venitul unei persoane $\in [20000, 1000000]$ lei/an
- Valorile atributelor sunt normalizate:

$$x'_i = \frac{(x_i - \min_i)}{(\max_i - \min_i)} \in [0,1]$$

2. Calcularea similitudinii

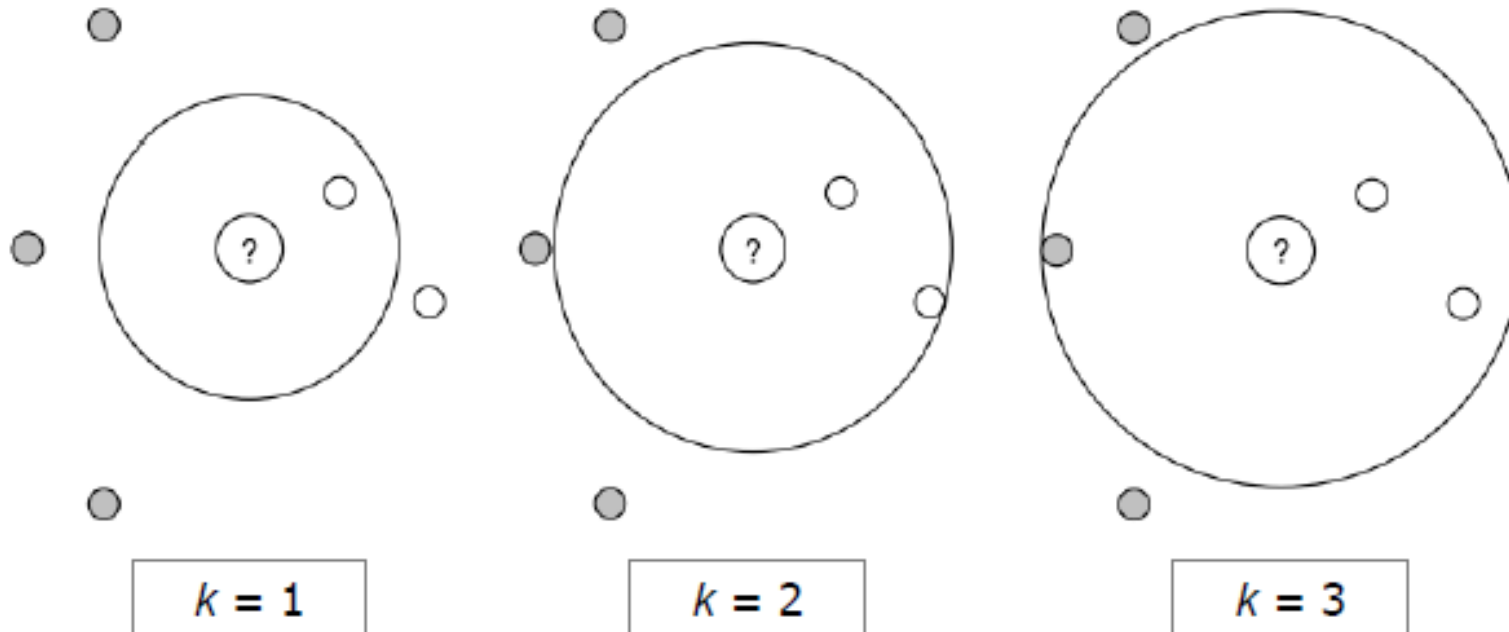
Metrici pentru calcularea distanței.

Atribute nominale

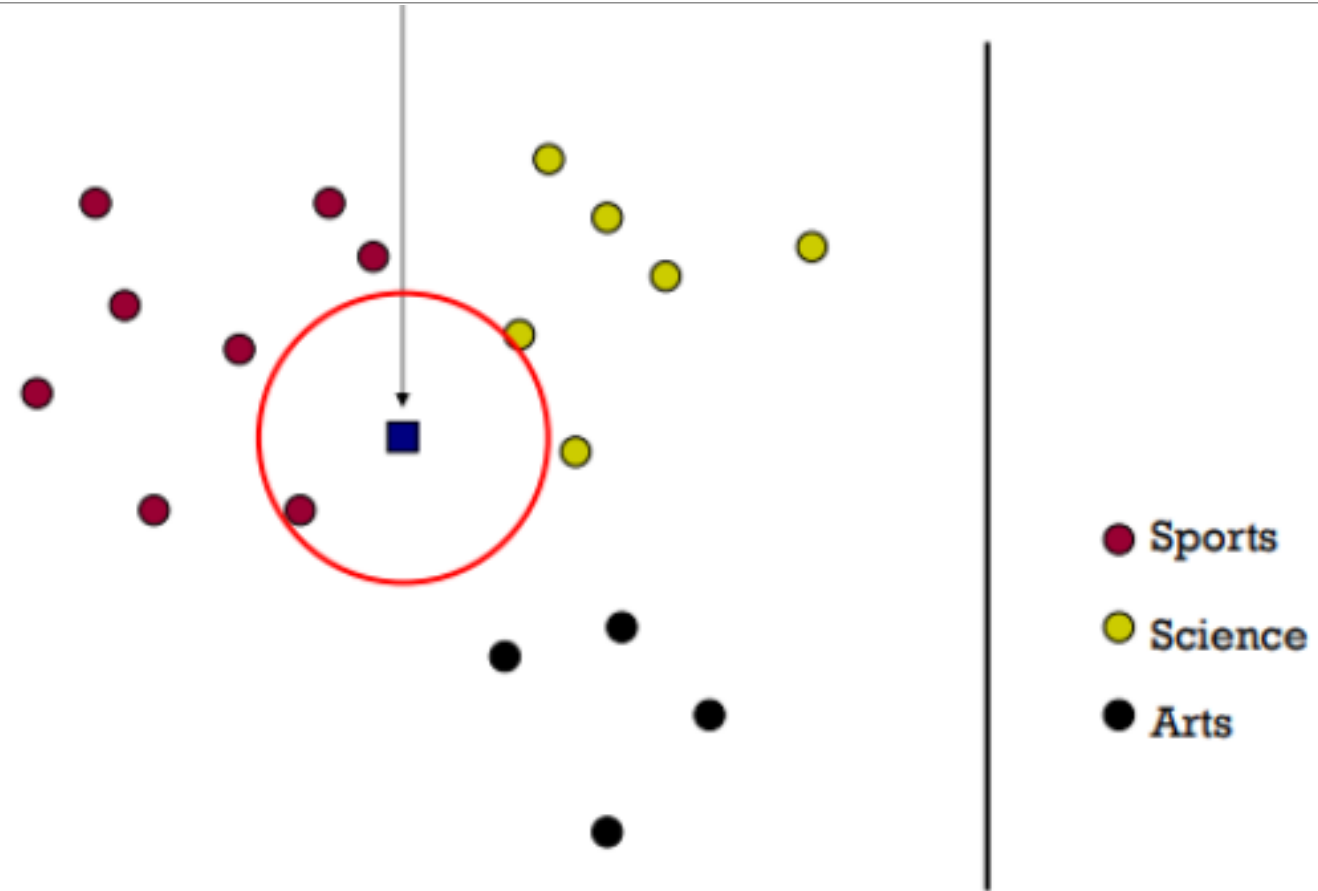
- ❑ Este necesară gasirea unei “distanțe” între valorile diferite ale atributelor nominale
 - Ex: distanța dintre valorile: roșu, galben și verde
- ❑ De obicei, se consideră distanța zero pentru valori identice și unu în caz contrar
 - Ex. având mai multe culori se poate utiliza o măsură metrică a nuanțelor din spațiul culorilor, punând galbenul mai apropiat de portocaliu decât verde.
- ❑ Unele atribute au o importanță diferită care este reflectată în distanța metrică cu ajutorul anumitor ponderi

3. Cum se alege numărul de vecini Cei mai apropiați vecini

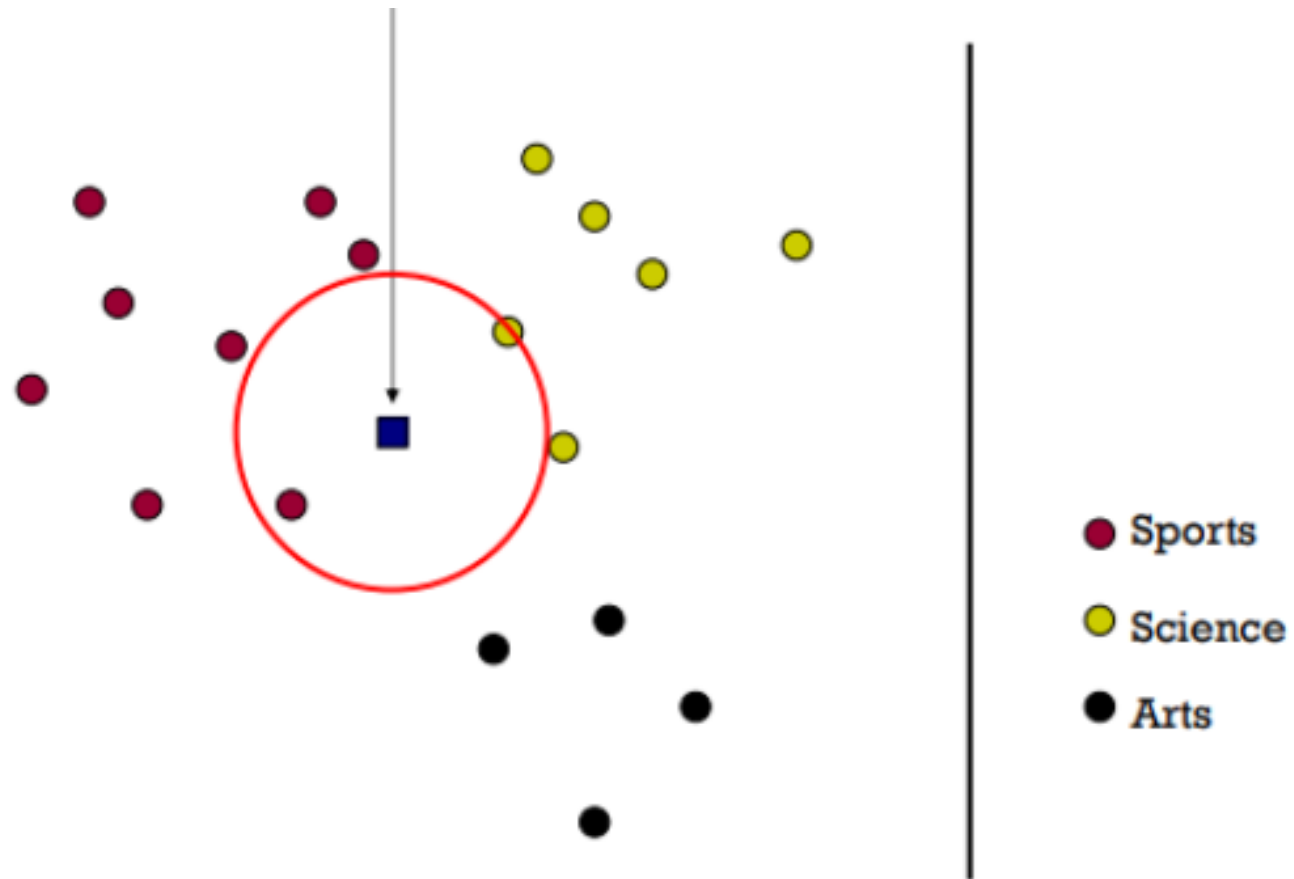
Cei mai apropiați k vecini ai unei instanțe x sunt punctele cu distanțele cele mai mici față de x



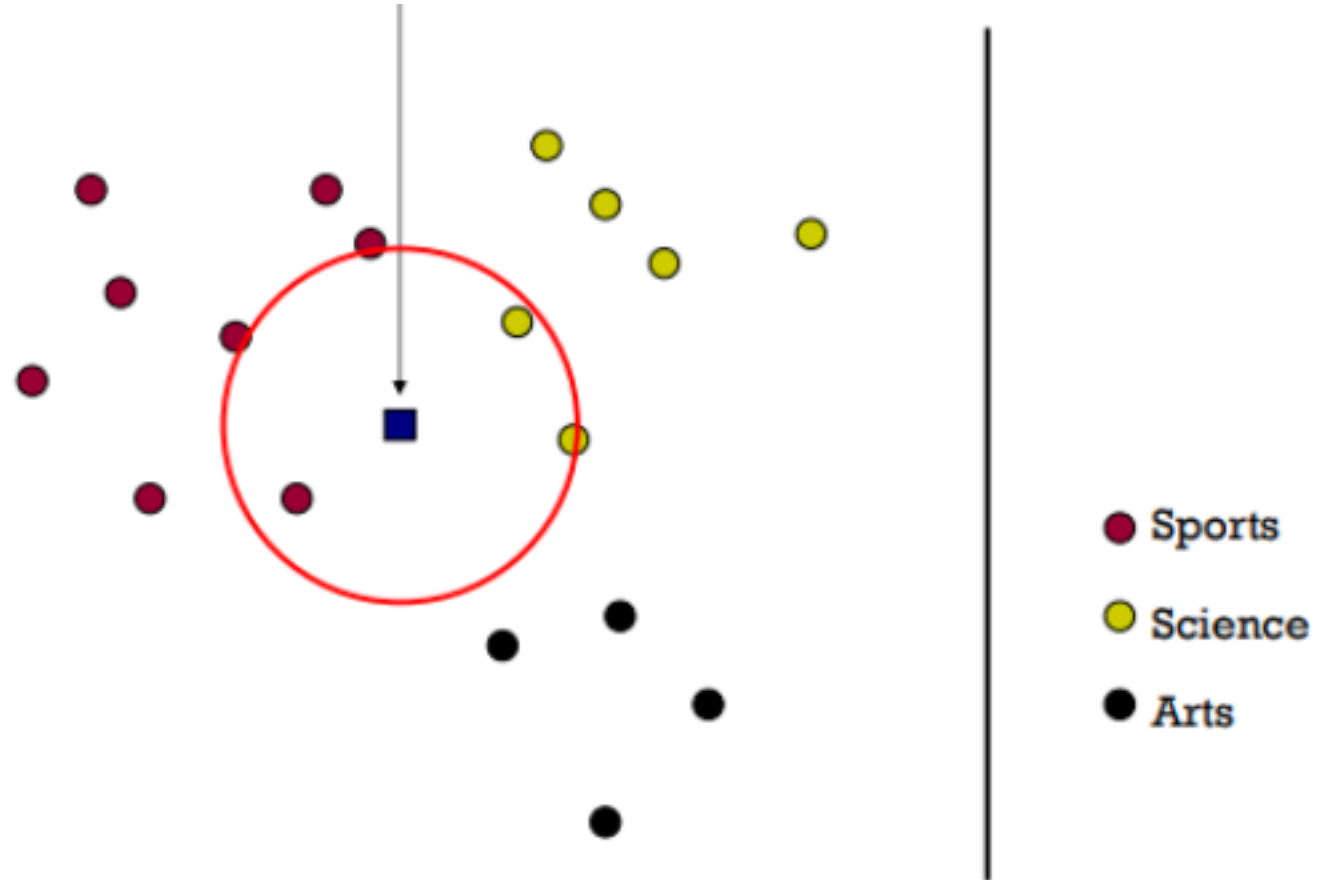
3. Cum se alege numărul de vecini 1-NN



3. Cum se alege numărul de vecini 2-NN



3. Cum se alege numărul de vecini 3-NN

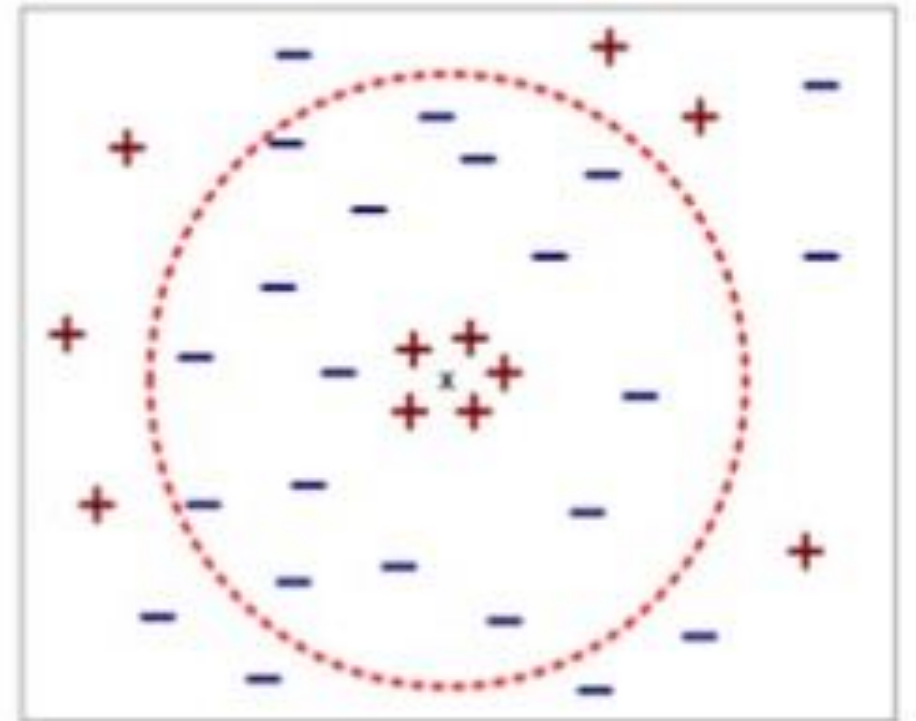


3. Cum se alege numărul de vecini

Numărul de vecini

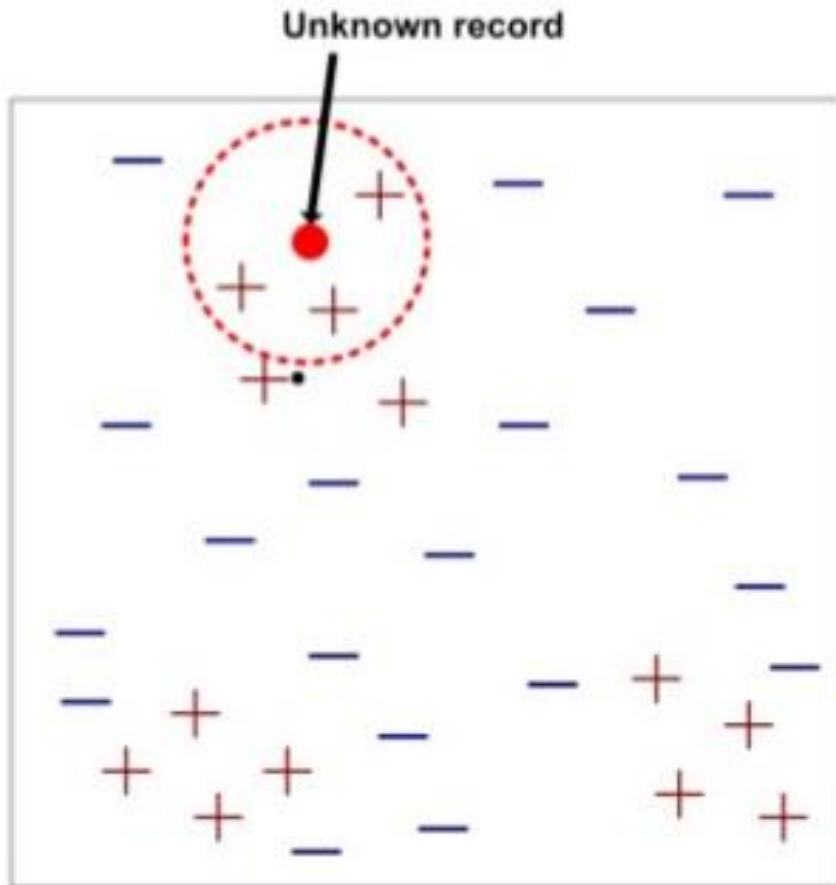
Valoarea lui k este importantă:

- ❑ dacă e prea mic, atunci clasificatorul poate fi suspectat de overfitting, pentru că devine prea sensibil la zgomotul din datele de intrare (zgomot \Rightarrow date eronate); clasificarea poate fi afectată de zgomot
- ❑ dacă e prea mare, atunci s-ar putea ca prea mulți dintre cei k vecini considerați să fie depărtați de punctul curent și deci irelevanți pentru clasificarea curentă; vecinătatea poate include puncte din alte clase

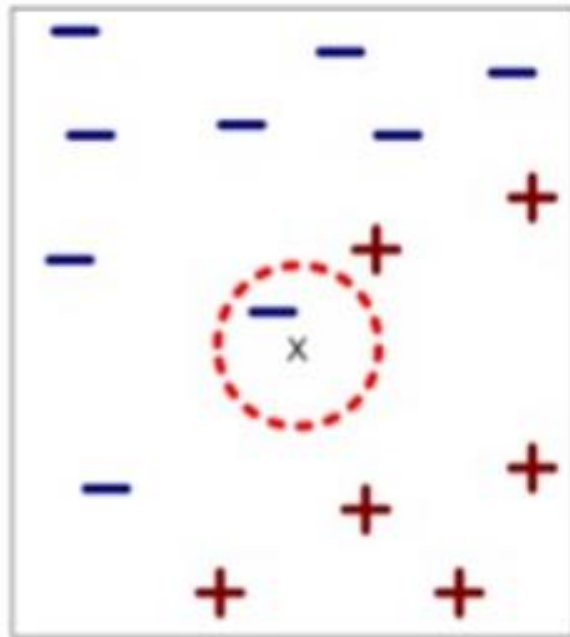


3. Cum se alege numărul de vecini

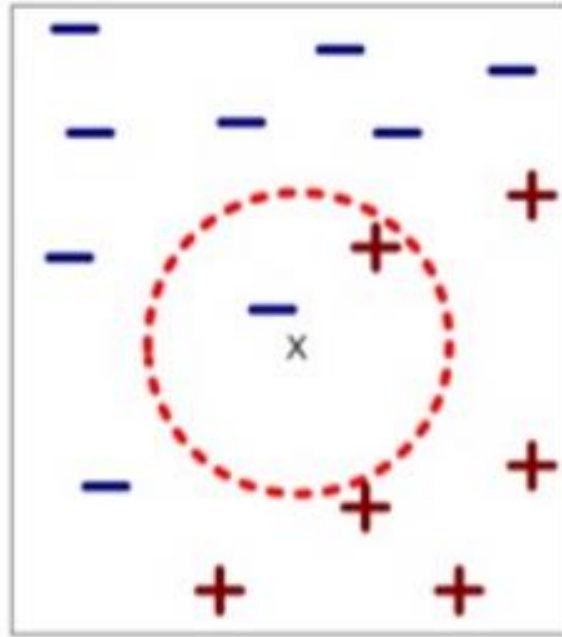
Exemplu: se consideră 3 vecini



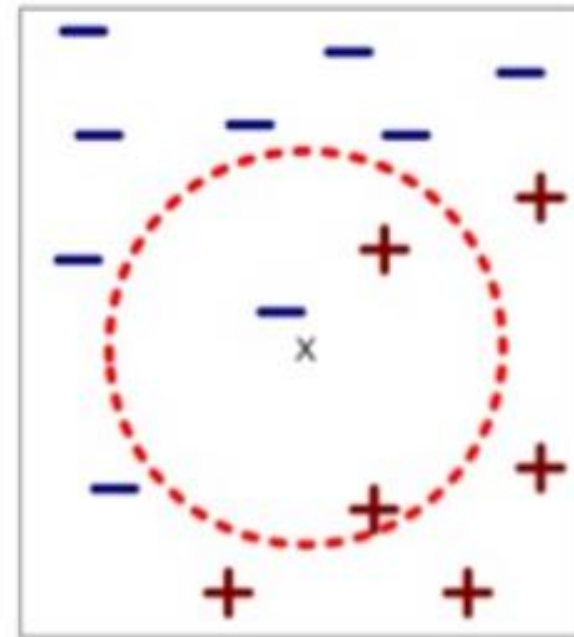
3. Cum se alege numărul de vecini Metoda celor mai apropiați k vecini



(a) 1-nearest neighbor

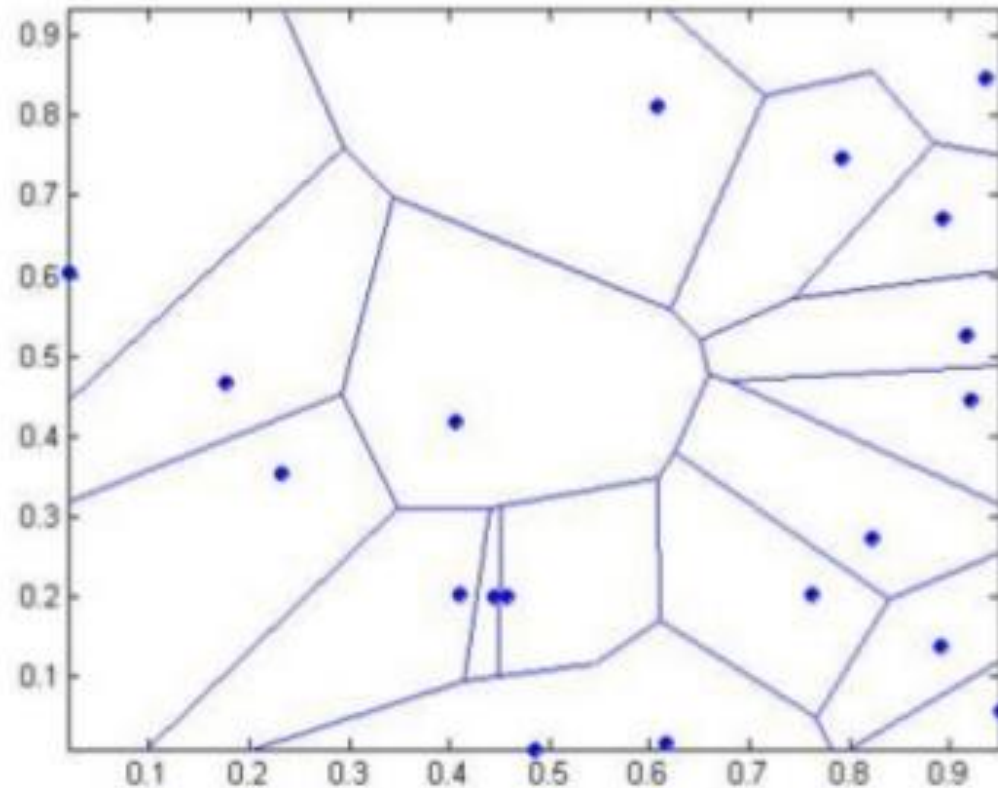


(b) 2-nearest neighbor



(c) 3-nearest neighbor

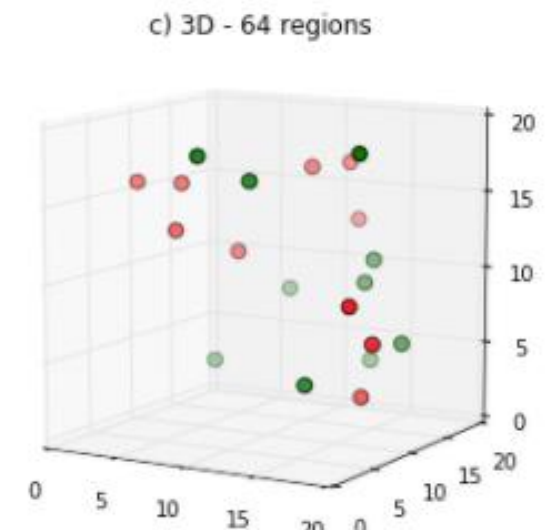
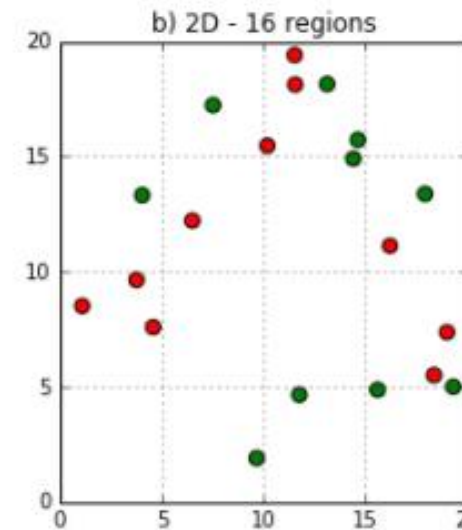
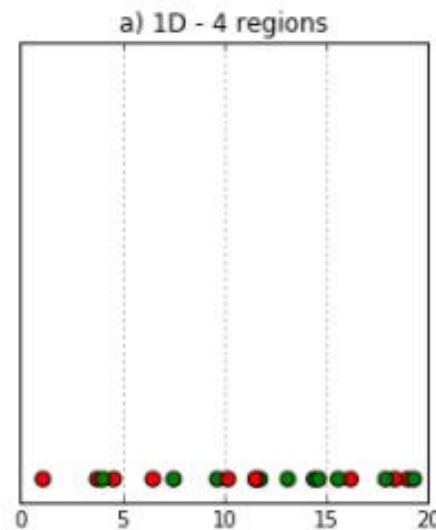
3. Cum se alege numărul de vecini 1-NN Diagrama Voronoi



- Pentru $k=1$ (Nearest neighbor) se obține diagrama Voronoi.
- În interiorul unei zone delimitate, orice punct are aceeași clasă ca punctul marcat din acea zonă.

Dimensionalitatea datelor

- ❑ Blestemul dimensionalității engl. “curse of dimensionality”
- ❑ Datele devin mai rare în spațiile multi-dimensionale
- ❑ Dacă numărul de attribute este mare, este nevoie de mai multe instanțe de antrenare pentru a forma un model corect



Algoritm

Metoda celor mai apropiați k vecini

Fie k numărul de vecini considerați și D setul de date de antrenare

Pentru fiecare exemplu de test x'

- Calculează $d(x, x')$, distanța între exemplul de test și datele (x, \cdot) din D
- Selectează $D_z \subseteq D$ setul celor mai apropiați k vecini ai lui x'
- Calculează clasa asociată lui x'

$$y' = \arg \max_v \sum_{(x_i, y_i) \in D_z} I(v = y_i)$$

Observații:

- funcția $I(\cdot)$ este funcția indicator, cu valoare 1 dacă argumentul are valoarea adevărat, 0 altfel.
- dacă există mai mulți v care maximizează partea dreaptă a expresiei de mai sus, atunci se alege arbitrar unul din aceștia
- se folosește un sistem de votare în care fiecare vecin are același impact în determinarea clasei estimate

Weighted k-NN

Pentru a reduce sensibilitatea algoritmului k-NN față de alegerea lui k se poate folosi o ponderare a vecinilor

Toți cei k vecini participă la decizia legată de clasa actuală, dar cu ponderi diferite:

- vecinii mai apropiați au pondere mai mare
- vecinii mai depărtați au pondere mai mică

Ponderea w poate fi descrescătoare cu distanța față de punctul ce se vrea a fi clasificat

$$y' = \arg \max_v \sum_{(x_i, y_i) \in D_z} I(v = y_i)$$

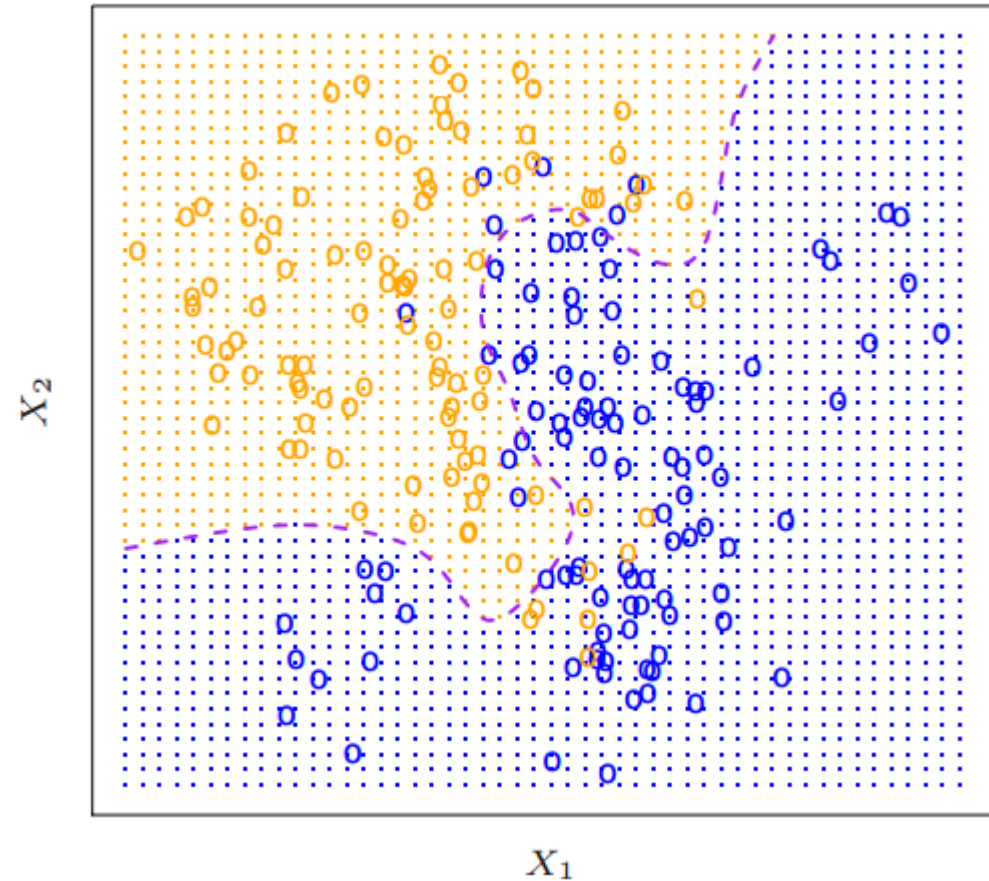
$$y' = \arg \max_v \sum_{(x_i, y_i) \in D_z} w_i \times I(v = y_i)$$

Caracteristici ale metodei k-NN

- ❑ Un tip particular de “instance-based learning”
- ❑ Nu produce efectiv un model; timpul de învățare este 0
- ❑ Clasificarea poate fi lentă, deoarece se face uz de tot corpusul de date din setul de instruire
- ❑ Clasificarea se face pe baza informației locale (alți clasificatori folosesc un model global)
- ❑ k-NN poate produce suprafețe de decizie arbitrar de complexe; suprafețele pot avea variabilitate mare, puternic influențate de setul de instruire
- ❑ Dacă nu se folosește preprocesare (scalarea diferitelor atribute ar trebui luată în considerare) sau o măsură de similaritate adecvată, valorile prezise pot fi greșite

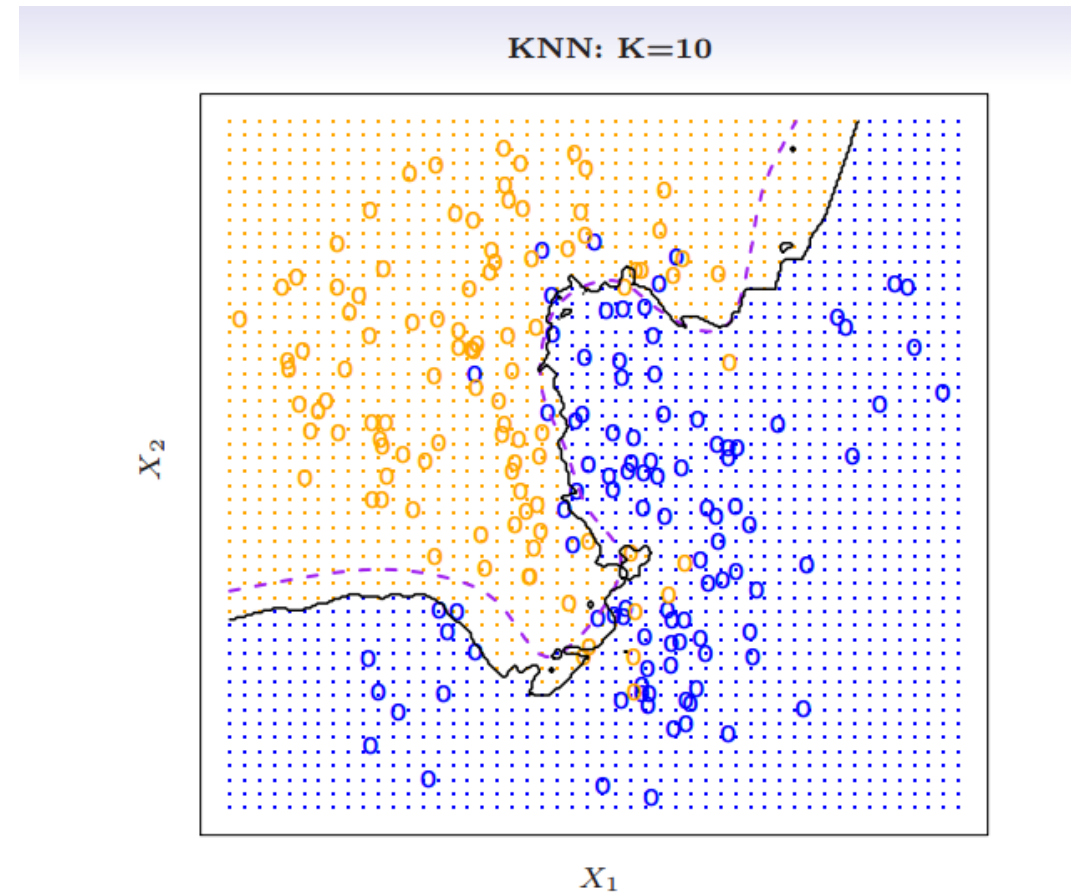
Exemplu

- 2 clase (galben si albastru)
- Decision boundary: linia punctata



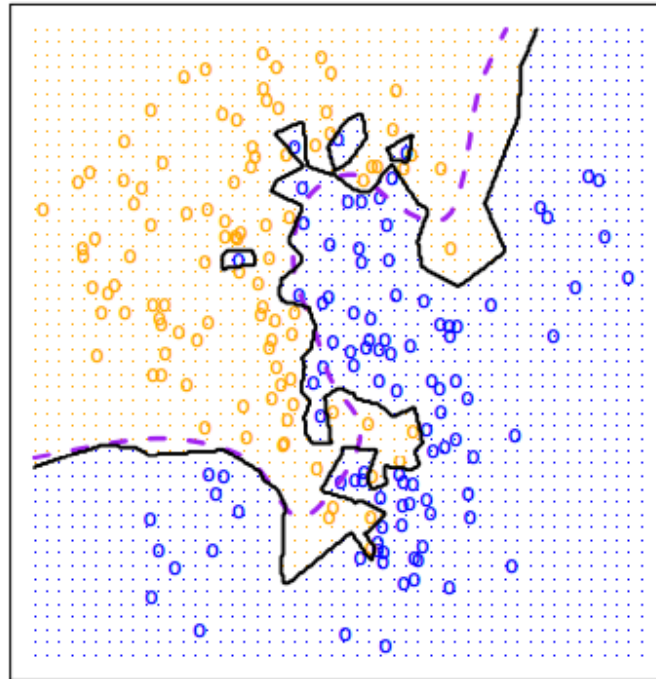
Exemplu

Decizia optima: $k=10$



Exemplu

KNN: $K=1$



KNN: $K=100$

