

# Curs 5:

## Clasificarea datelor (III)

kNearest Neighbours și Naive Bayes

# Structura

- Clasificatori bazați pe instanțe
  - Modelul celor mai apropiați vecini (k-Nearest Neighbours)
- Clasificatori bazați pe modele probabiliste
  - Modelul Bayesian naiv (Naive Bayes)

# Clasificatori bazați pe instanțe

Ideea principală: datele similare aparțin aceleiași clase

(raționamentul bazat pe analogie este utilizat în multe domenii – exemplu: ce boală a avut un pacient cu simptome similare?)

- Modelul de clasificare constă tocmai din setul de antrenare
  - Procesul de antrenare constă doar în **stocarea datelor din set**
- Clasificarea unei noi date constă în:
  - Se calculează **similaritatea** (sau disimilaritatea) dintre noua dată și cele din setul de antrenare și se identifică **exemplarele cele mai apropiate**
  - Se **alege clasa cea mai frecvent întâlnită** în subsetul celor mai similare exemple

# Clasificatori bazați pe instanțe

Ideea principală: datele similare aparțin aceleiași clase

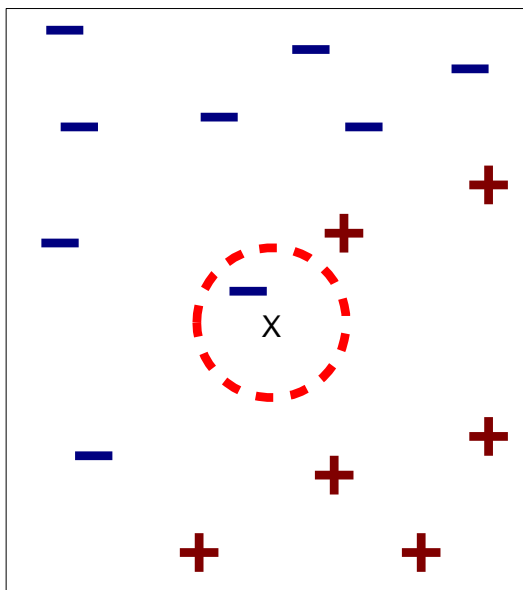
Obs:

- Astfel de clasificatori sunt considerați leneși (“lazy”) deoarece faza de antrenare nu presupune nici un efort de calcul (întregul efort este amânat pentru faza de clasificare)
- Cei mai populari clasificatori din această categorie sunt cei bazați pe principiul celui/celor mai apropiat/apropiați vecin/vecini (**k-Nearest Neighbour**)
- Aplicații:
  - sisteme de recomandare
  - diagnoza medicală

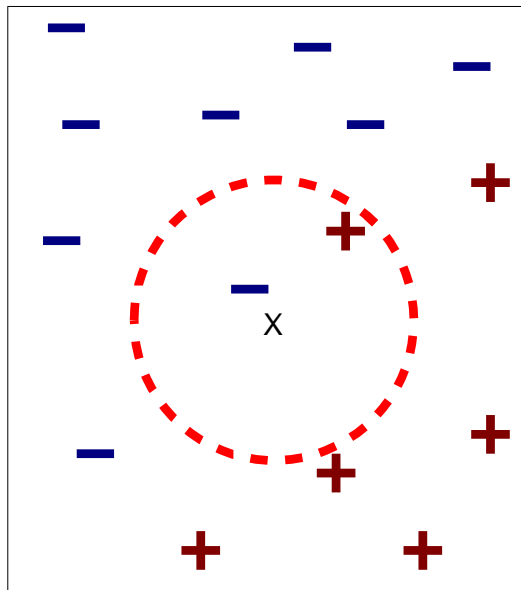
# Clasificatori bazați pe instanțe

## kNN – k Nearest Neighbour

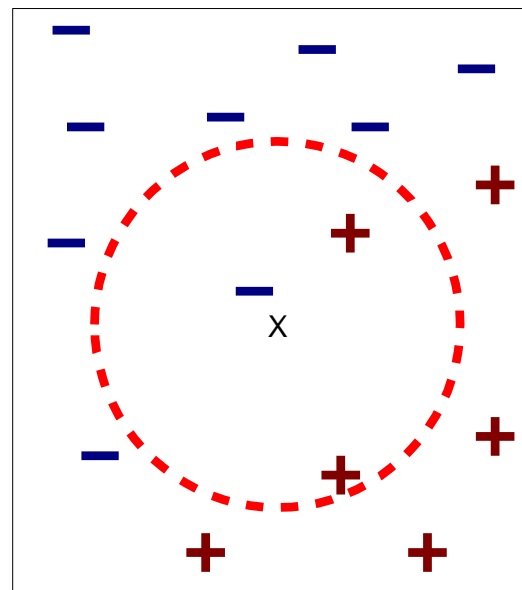
- Pt fiecare dată de clasificat:
  - Determină cele mai apropiate (mai similare) **k** exemple din setul de antrenare
  - Identifică cea mai frecventă clasă



(a) 1-nearest neighbor



(b) 2-nearest neighbor



(c) 3-nearest neighbor

# Clasificatori bazați pe instanțe

## kNN – k Nearest Neighbour

- Pt fiecare dată de clasificat:
  - Determină cele mai apropiate (mai similare) **k** exemple din setul de antrenare
  - Identifică cea mai frecventă clasă

Performanța clasificatorilor de tip kNN depinde de:

- Măsura de **similaritate/ disimilaritate**
  - Se alege în funcție de tipurile atributelor și de proprietățile problemei
- **Valoarea lui k (numărul de vecini)**
  - Cazul cel mai simplu:  $k=1$  (nu e indicat în cazul datelor afectate de zgomot)

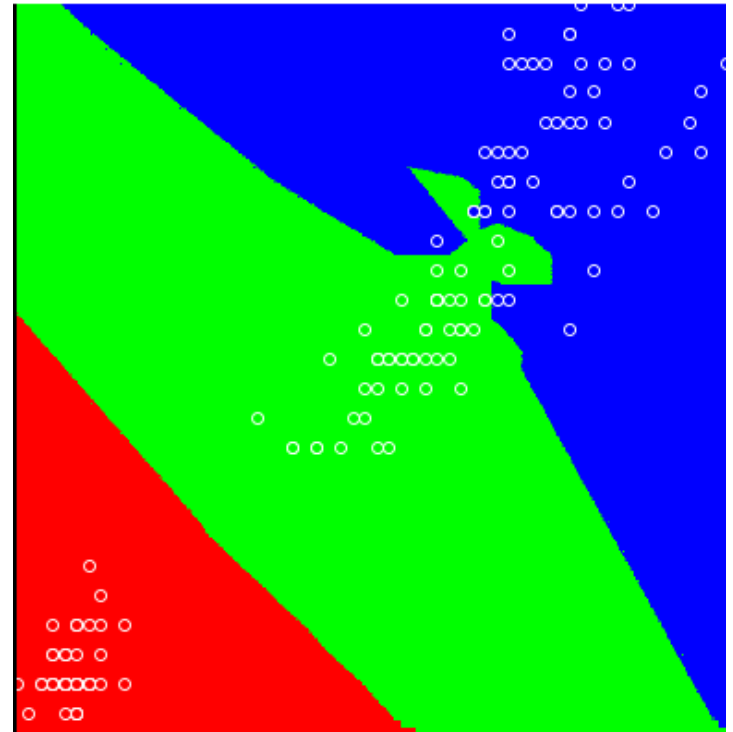
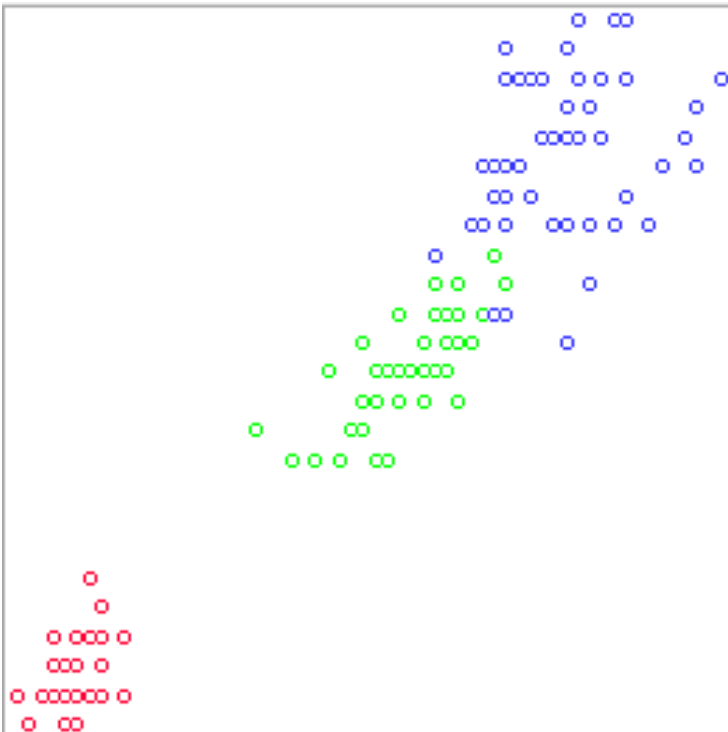
**Obs:** kNN induce o partiționare în regiuni a spațiului datelor; regiunile nu sunt calculate explicit ci sunt implicit determinate de măsura de similaritate (precum și de valoarea lui k)

# Clasificatori bazați pe instanțe

1NN = Nearest Neighbor bazat pe cel mai apropiat vecin (și **distanța euclidiană**)

Ilustrarea regiunilor. Dataset: iris2D (“petal length” and “petal width”).

Plot: [Weka->Visualization->BoundaryVisualizer](#)



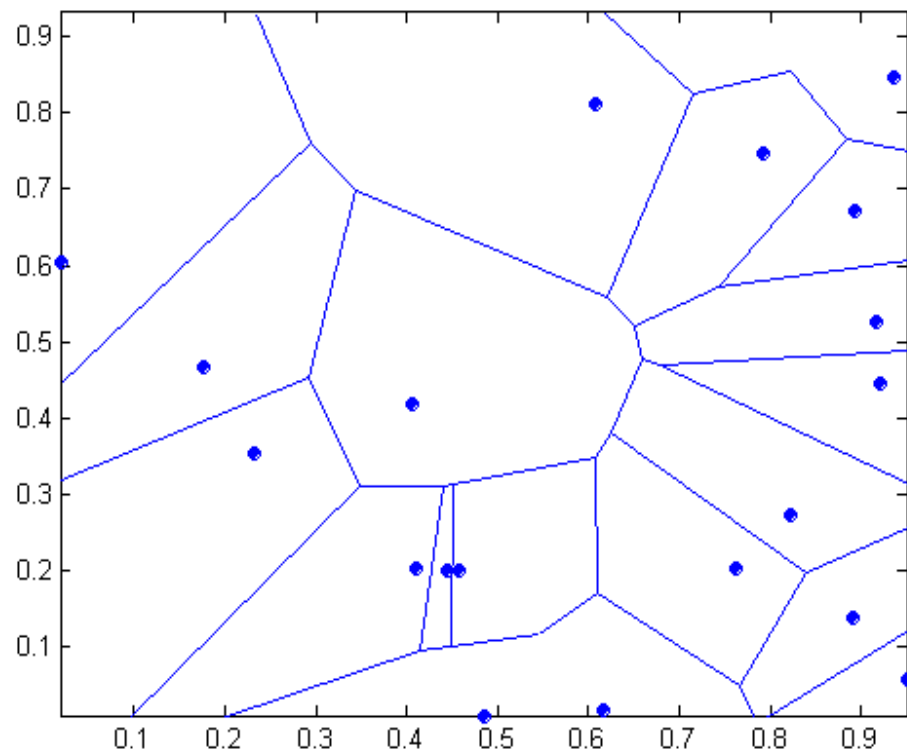
# Clasificatori bazați pe instanțe

1NN = Nearest Neighbor bazat pe cel mai apropiat vecin (și distanța euclidiană)

1NN induce o partiționare a spațiului datelor (e.g. în 2D aceasta corespunde unei diagrame Voronoi)

Obs:

Fiecare instanță din setul de antrenare (punctele din imagine) corespunde unei regiuni care cuprinde datele aflate în vecinătatea acelei instanțe



[Tan, Steinbach, Kumar; Introduction to Data Mining, slides, 2004]



# Măsuri de similaritate/ disimilaritate

Considerăm două entități (e.g. vectori de date, serii de timp etc) A and B

- O măsură de **similaritate**,  $S$ , asociază perechii (A,B) un număr,  $S(A,B)$ , care este cu atât mai mare cu cât A și B sunt mai similare
- O măsură de **disimilaritate**,  $D$ , asociază perechii (A,B) un număr,  $D(A,B)$ , care este cu atât mai mare cu cât A și B sunt mai diferite

(obs: o măsură de disimilaritate nu satisface neapărat toate proprietățile matematice ale unei metrice, de exemplu poate să nu satisfacă inegalitatea triunghiului)

Alegerea măsurii depinde de:

- Tipul atributelor
- Numărul de attribute
- Distribuția datelor
- Particularitățile problemei

# Măsuri de similaritate/ disimilaritate

## Atribute numerice

Cele mai populare măsuri de disimilaritate:

- Distanța euclidiană
- Distanța Manhattan

Obs:

- Distanța euclidiană este invariantă în raport cu rotații
- Dacă nu toate atributele au aceeași importanță (sau dacă nu sunt scalate adecvat) atunci e indicat să se folosească varianta ponderată

( $w_i(a_i - b_i)^2$  în loc de  $(a_i - b_i)^2$ )

$$d_p(A, B) = \sqrt[p]{\sum_{i=1}^n (a_i - b_i)^p} \quad (\text{Minkowski, } L_p)$$
$$d_E(A, B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2} \quad (\text{Euclidean, } p = 2)$$
$$d_M(A, B) = \sum_{i=1}^n |a_i - b_i| \quad (\text{Manhattan, } p = 1)$$
$$d_\infty(A, B) = \max_{i=1, \dots, n} |a_i - b_i| \quad (p = \infty)$$

Obs.

- Ponderile pot fi determinate folosind tehnici de preprocesare
- Pentru a evita utilizarea ponderilor datele pot fi în prealabil scalate sau standardizate

# Măsuri de similaritate/ disimilaritate

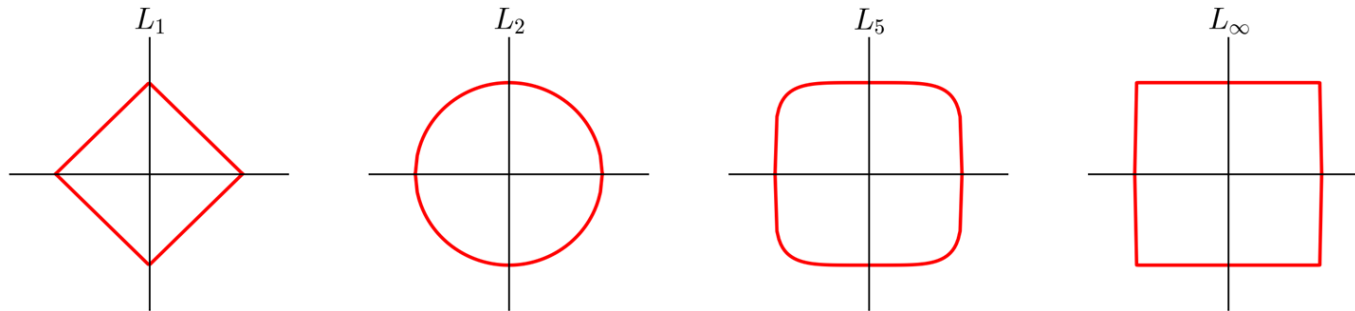
## Atribute numerice

Cele mai populare măsuri de disimilaritate:

- Distanța euclidiană
- Distanța Manhattan

Diferențe între măsurile de disimilaritate  
(vizualizarea punctelor egal distanțate  
de originea sistemului de axe de  
coordonate)

$$d_p(A, B) = \sqrt[p]{\sum_{i=1}^n (a_i - b_i)^p} \quad (\text{Minkowski, } L_p)$$
$$d_E(A, B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2} \quad (\text{Euclidean, } p = 2)$$
$$d_M(A, B) = \sum_{i=1}^n |a_i - b_i| \quad (\text{Manhattan, } p = 1)$$
$$d_\infty(A, B) = \max_{i=1, \dots, n} |a_i - b_i| \quad (p = \infty)$$



Source: Lecture 19 - Data Science, Steven Skiena, Stony Brook University

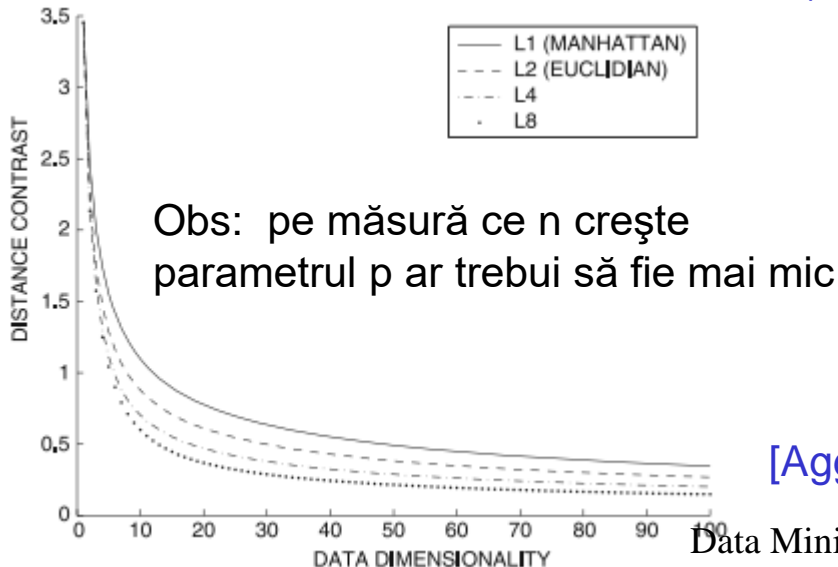
# Măsuri de similaritate/ disimilaritate

Aspecte practice – problema dimensiunii (dimensionality curse):

- Puterea de discriminare a acestor distanțe scade pe măsură ce nr de attribute (n) crește →

pt date cu multe attribute clasificatorii bazați pe distanțe devin inefectivi

- Se recomandă reducerea dimensionalității (de ex. prin PCA sau transformări neliniare – tSNE, UMAP)



$$d_p(A, B) = \sqrt[p]{\sum_{i=1}^n (a_i - b_i)^p} \quad (\text{Minkowski, } L_p)$$
$$d_E(A, B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2} \quad (\text{Euclidean, } p = 2)$$
$$d_M(A, B) = \sum_{i=1}^n |a_i - b_i| \quad (\text{Manhattan, } p = 1)$$
$$d_\infty(A, B) = \max_{i=1, \dots, n} |a_i - b_i| \quad (p = \infty)$$

$$\text{Distance contrast} = \frac{d_{\max} - d_{\min}}{\sigma}$$

$d_{\max}, d_{\min}$  = cea mai mare și cea mai mică dintre distanțe

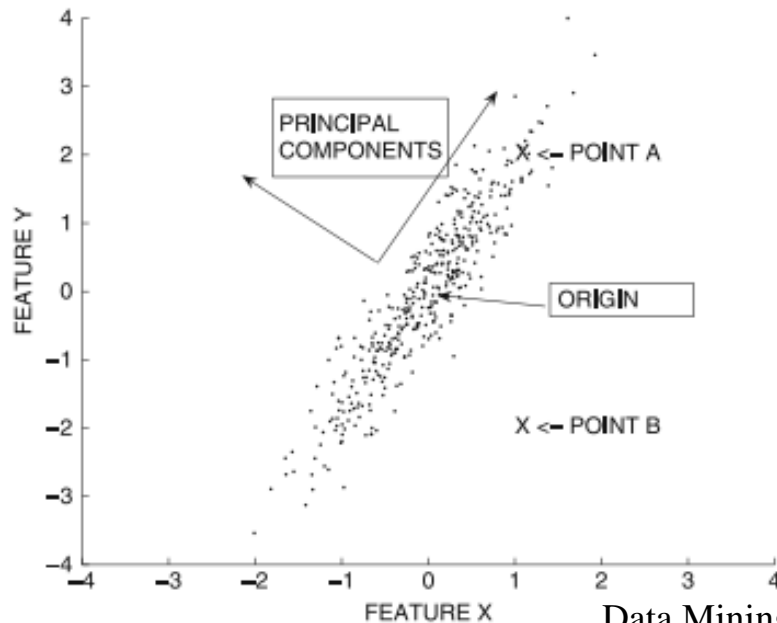
$\sigma$  = abaterea standard a distanțelor

[Aggarwal, Data Mining Textbook, 2015]

# Măsuri de similaritate/ disimilaritate

Aspecte practice – impactul distribuției datelor

Intrebare: Care punct e mai aproape de origine? A sau B?



[Aggarwal, Data Mining Textbook, 2015]

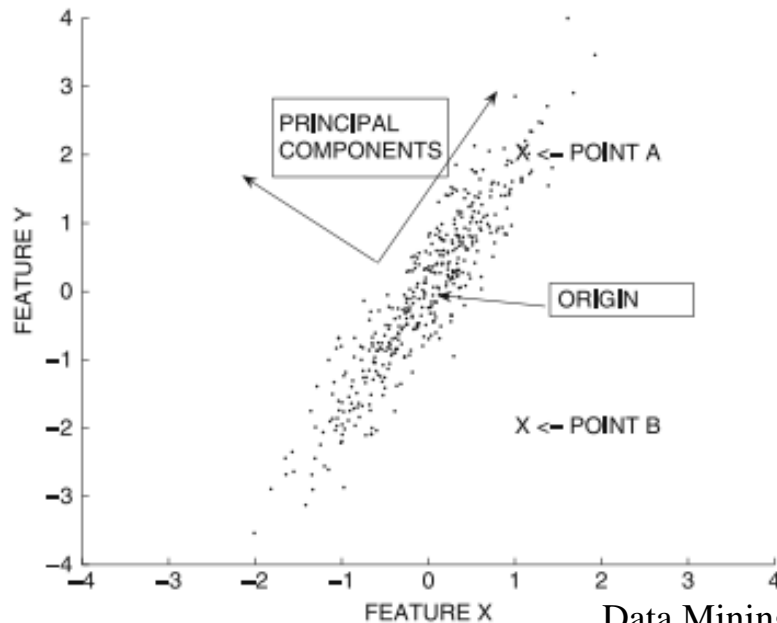
# Măsuri de similaritate/ disimilaritate

## Aspecte practice – impactul distribuției datelor

Intrebare: Care punct e mai aproape de origine? A sau B?

R:  $d(O,A) = d(O,B)$  (distanțe euclidiene egale). Luând în considerare distribuția datelor: A este mai apropiat de O decât B

**Altă întrebare:** cum poate fi inclusă distribuția datelor în calculul distanței?



Distanța Mahalanobis

$$d_{Mah}(A, B) = \sqrt{(A - B)^T C^{-1} (A - B)}$$

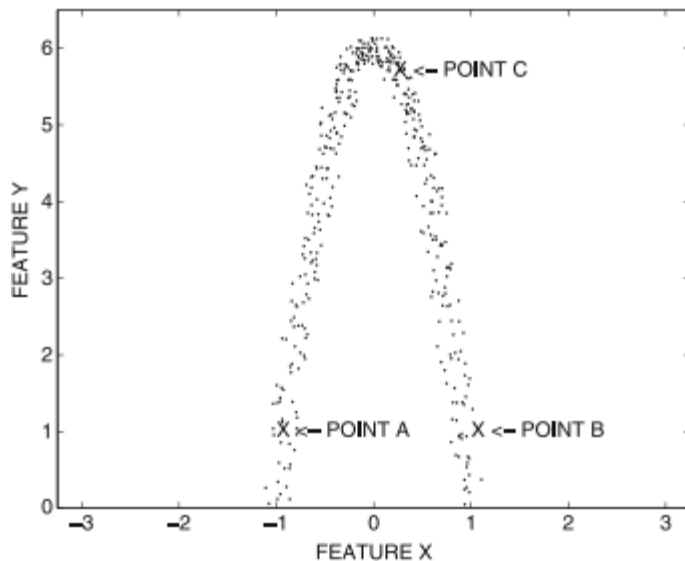
$C^{-1}$  = inversa matricii de covarianța

[Aggarwal, Data Mining Textbook, 2015]

# Măsuri de similaritate/ disimilaritate

Aspecte practice – impactul distribuției datelor

Intrebare: este distanța dintre A și B mai mică decât distanța dintre B și C?



[Aggarwal, Data Mining Textbook, 2015]

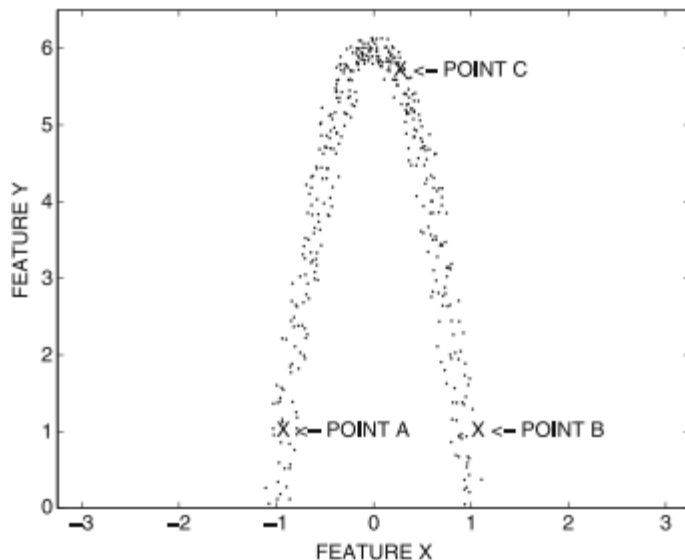
# Măsuri de similaritate/ disimilaritate

## Aspecte practice – impactul distribuției datelor

Intrebare: este distanța dintre A și B mai mică decât distanța dintre B și C?

R: da, dacă ignorăm distribuția datelor și folosim distanța euclidiană

Totuși, distribuția datelor nu poate fi ignorată întrucât este cea care furnizează contextul problemei, iar în acest context  $d(A,B) > d(B,C)$



### Distanța geodesică:

- Se construiește un graf ce are în noduri punctele iar muchiile unesc nodurile vecine (ex: cei mai apropiați k vecini)
- Calculează distanța dintre două puncte ca fiind cea mai scurtă cale în graf

[Aggarwal, Data Mining Textbook, 2015]



# Măsuri de similaritate/ disimilaritate

## Atribute numerice – măsură de similaritate

- Măsura cosinus:  $\text{sim}(A,B)=A^T B/(||A|| \ ||B||)$  (produsul scalar dintre A și B împărțit la produsul normelor)

### Remarcă:

- In cazul vectorilor normalizați ( $||A||=||B||=1$ ) similaritatea e maximă când distanța euclidiană este minimă:

$$\begin{aligned} d_E^2(A, B) &= (A - B)^T (A - B) = A^T A - 2A^T B + B^T B \\ &= 2(1 - A^T B) = 2(1 - \text{sim}(A, B)) \end{aligned}$$

# Măsuri de similaritate/ disimilaritate

## Atribute nominale

**Abordare 1:** Transformarea atributelor nominale în atribute numerice (prin binarizare = one hot encoding) și utilizarea măsurilor de similaritate/disimilaritate pentru vectori binari:

- Disimilaritate: distanța **Hamming** = distanța Manhattan:  
 $d_H(A, B) = d_M(A, B)$

- **Jaccard similarity:**

$$J(A, B) = \frac{\sum_{i=1}^n a_i b_i}{\sum_{i=1}^n (a_i^2 + b_i^2 - a_i b_i)} = \frac{\text{card}(S_A \cap S_B)}{\text{card}(S_A \cup S_B)}$$

**Obs:**  $S_A$  și  $S_B$  sunt submulțimi ale mulțimii globale cu  $n$  atribute care corespund vectorilor de apartenență  $A$  și  $B$ .

# Măsuri de similaritate/ disimilaritate

## Atribute nominale

**Abordare 2:** Utilizează măsuri locale de similaritate (între valorile atributelor)

$$S(A, B) = \sum_{i=1}^n S(a_i, b_i)$$
$$S(a_i, b_i) = \begin{cases} 1 & \text{if } a_i = b_i \\ 0 & \text{if } a_i \neq b_i \end{cases}$$

**Obs:** similaritățile mai puțin frecvente pot fi considerate mai relevante decât cele frecvente

$$S(a_i, b_i) = \begin{cases} 1/f^2(a_i) & \text{if } a_i = b_i \\ 0 & \text{if } a_i \neq b_i \end{cases}$$

$f(a_i)$  = frecvența valorii  $a_i$  în setul de date (pt atributul  $i$ )

# Măsuri de similaritate/ disimilaritate

**Atribute mixte:** se combină măsurile corespunzătoare celor două tipuri de atribute (utilizând ponderi specifice)

$$S(A, B) = \lambda S_{numerical}(A, B) + (1 - \lambda) S_{nominal}(A, B)$$

**Alte tipuri de date:**

- **Siruri** (e.g. text sau secvențe biologice) – se utilizează **distanța de editare (distanța Levenshtein)**
- **Concepte** (e.g. noduri într-o ontologie) – distanțe bazate pe **cele mai scurte căi în grafuri sau arbori**
- **Grafuri** (e.g. rețele sociale sau biologice) – ponderea structurilor (tiparelor) similare în cele două structuri

# kNN: alegerea lui k

Performanța clasificatorilor de tip kNN depinde de numărul de vecini

Cazuri extreme:

- $k=1$  - clasificatorul nu este robust (erorile din setul de date influențează răspunsul clasificatorului)
- $k=N$  - e echivalent cu ZeroR fiind bazat doar pe modul de distribuire a datelor în clase

Cum se alege k?

- Abordare de tip trial-and-error: se încearcă diferite valori și se alege valoarea care maximizează performanța

## kNN: cost

Clasificarea unei noi instanțe necesită calculul a  $N$  distanțe (sau măsuri de similaritate pt un set de date cu  $N$  elemente care au  $n$  attribute precum și selecția celor mai mici  $k$  distanțe  $\rightarrow O(Nn+kN)$  (costul de calcul a similarității / disimilarității poate fi diferit de  $Nn$  – depinde de structura datelor)

Dacă  $N$  e mare această prelucrare poate fi costisitoare (întrucât trebuie efectuată pentru fiecare instanță care trebuie clasificată)

### Abordări posibile:

- Crearea unei structuri de indexare a datelor din setul de antrenare care permite identificarea celor mai apropiați vecini într-un mod eficient (arbori k-d = generalizare a arborilor binari de căutare pentru date k-dimensionale; obs: în notațiile noastre  $k=n$  și nu are legătură cu numărul de vecini)
- Reducerea numărului de date din setul de antrenare prin gruparea lor în clustere și înlocuirea fiecărui cluster cu un singur prototip
- Selecția unor prototipuri din set

# Modele probabiliste de clasificare

**Exemplu:** Presupunem că ne interesează să estimăm probabilitatea ca un pacient care are simptomul S să aibă boala T

- Probabilitatea de estimat:  $P(T|S)$  = probabilitatea evenimentului T condiționată de evenimentul S
- Presupunem că se cunosc:
  - $P(S)$  – probabilitatea ca simptomul sa fie prezent (**evidence**)
  - $P(T)$  – prevalența bolii T estimată pe baza unor studii populaționale (e o măsură a frecvenței de apariție a bolii) (**prior**)
  - $P(S|T)$  – se estimează pe baza cunoștințelor medicale (cât de frecvent este simptomul S în cazul bolii T) (**likelihood**)
- Regula de calcul (regula probabilității condiționate – formula lui Bayes):

$$P(T|S) = P(S|T)P(T)/P(S) = \text{likelihood} * \text{prior} / \text{evidence}$$

- Cum se analizează cazul în care nu e un singur simptom S, ci mai multe simptome  $S_1, S_2, \dots, S_n$ ?

# Modele probabiliste de clasificare

**Exemplu:** Presupunem că ne interesează să estimăm probabilitatea ca un pacient care are simptomele  $S_1, S_2, \dots, S_n$  să aibă boala  $T$

- Probabilitatea de estimat:  $P(T | S_1, S_2, \dots, S_n)$
- Se folosește regula Bayes:
  - $P(T | S_1, S_2, \dots, S_n) = P(S_1, S_2, \dots, S_n | T)P(T) / P(S_1, S_2, \dots, S_n)$
- **Ipoteză simplificatoare:** simptomele  $(S_1, S_2, \dots, S_n)$  corespund unor evenimente independente în condițiile în care clasa este cunoscută (această ipoteză nu este întotdeauna adevărată însă poate fi acceptată în anumite situații practice)
- Întrucât  $P(S_1, S_2, \dots, S_n)$  nu depinde de  $T$ , dacă se dorește doar să se determine clasa ( $T$  pozitiv sau  $T$  negativ) este suficient să se considere că  $P(T | S_1, S_2, \dots, S_n)$  este proporțional cu  $P(S_1 | T) P(S_2 | T) \dots P(S_n | T)P(T)$



# Clasificatorul Naïve Bayes

## Problema de clasificare:

- Pentru o dată  $D_i=(a_{i1},a_{i2},\dots,a_{in})$  se pune problemă determinării clasei căreia îi aparține

## Ideea principală

- Estimează  $P(C_k|D_i)=P(a_{i1},a_{i2},\dots,a_{in}|C_k)P(C_k)/P(a_{i1},a_{i2},\dots,a_{in})$  pt fiecare  $k$  din  $\{1,2,\dots,K\}$  și selectează probabilitatea maximă; aceasta va indica cărei clase îi aparține, cel mai probabil, data; întrucât  $P(a_{i1},a_{i2},\dots,a_{in})$  este aceeași indiferent de clasă, probabilitatea de la numitor nu influențează decizia (se poate considera egală cu 1)
- **Ipoteză simplificatoare**: attributele sunt **independente** (acesta este motivul pentru care clasificadorul este denumit “naiv”)
- $P(C_k|D_i)=P(a_{i1}|C_k)P(a_{i2}|C_k)\dots P(a_{in}|C_k)P(C_k)$
- Estimarea probabilității de clasificare necesită cunoașterea lui  $P(a_{i1}|C_k)$ ,  $P(a_{i2}|C_k)$ , ...,  $P(a_{in}|C_k)$  și  $P(C_k)$
- Aceste probabilități pot estimate pe baza setului de antrenare (ca frecvențe relative) – această estimare corespunde procesului de învățare specific clasificadorului Naïve Bayes

# Clasificatorul Naïve Bayes

Exemplu:

Relation: weather.symbolic

| No. | outlook<br>Nominal | temperature<br>Nominal | humidity<br>Nominal | windy<br>Nominal | play<br>Nominal |
|-----|--------------------|------------------------|---------------------|------------------|-----------------|
| 1   | sunny              | hot                    | high                | FALSE            | no              |
| 2   | sunny              | hot                    | high                | TRUE             | no              |
| 3   | overcast           | hot                    | high                | FALSE            | yes             |
| 4   | rainy              | mild                   | high                | FALSE            | yes             |
| 5   | rainy              | cool                   | normal              | FALSE            | yes             |
| 6   | rainy              | cool                   | normal              | TRUE             | no              |
| 7   | overcast           | cool                   | normal              | TRUE             | yes             |
| 8   | sunny              | mild                   | high                | FALSE            | no              |
| 9   | sunny              | cool                   | normal              | FALSE            | yes             |
| 10  | rainy              | mild                   | normal              | FALSE            | yes             |
| 11  | sunny              | mild                   | normal              | TRUE             | yes             |
| 12  | overcast           | mild                   | high                | TRUE             | yes             |
| 13  | overcast           | hot                    | normal              | FALSE            | yes             |
| 14  | rainy              | mild                   | high                | TRUE             | no              |

$$P(C1)=P(\text{no})=5/14 \quad P(C2)=P(\text{yes})=9/14$$

A1: outlook

$$P(\text{sunny}|C1)=P(\text{sunny},C1)/P(C1) \\ = (3/14)/(5/14)=3/5$$

$$P(\text{sunny}|C2)=P(\text{sunny},C2)/P(C2) \\ = (2/14)/(9/14)=2/9$$

$$P(\text{overcast}|C1)=P(\text{overcast},C1)/P(C1) \\ = 0$$

$$P(\text{overcast}|C2)=P(\text{overcast},C2)/P(C2) \\ = (4/14)/(9/14)=4/9$$

$$P(\text{rainy}|C1)=P(\text{rainy},C1)/P(C1) \\ = (2/14)/(5/14)=2/5$$

$$P(\text{rainy}|C2)=P(\text{rainy},C2)/P(C2) \\ = (3/14)/(9/14)=3/9$$

# Clasificatorul Naïve Bayes

Exemplu:

| Relation: weather.symbolic |                    |                        |                     |                  |                 |
|----------------------------|--------------------|------------------------|---------------------|------------------|-----------------|
| No.                        | outlook<br>Nominal | temperature<br>Nominal | humidity<br>Nominal | windy<br>Nominal | play<br>Nominal |
| 1                          | sunny              | hot                    | high                | FALSE            | no              |
| 2                          | sunny              | hot                    | high                | TRUE             | no              |
| 3                          | overcast           | hot                    | high                | FALSE            | yes             |
| 4                          | rainy              | mild                   | high                | FALSE            | yes             |
| 5                          | rainy              | cool                   | normal              | FALSE            | yes             |
| 6                          | rainy              | cool                   | normal              | TRUE             | no              |
| 7                          | overcast           | cool                   | normal              | TRUE             | yes             |
| 8                          | sunny              | mild                   | high                | FALSE            | no              |
| 9                          | sunny              | cool                   | normal              | FALSE            | yes             |
| 10                         | rainy              | mild                   | normal              | FALSE            | yes             |
| 11                         | sunny              | mild                   | normal              | TRUE             | yes             |
| 12                         | overcast           | mild                   | high                | TRUE             | yes             |
| 13                         | overcast           | hot                    | normal              | FALSE            | yes             |
| 14                         | rainy              | mild                   | high                | TRUE             | no              |

$$P(C1)=P(\text{no})=5/14 \quad P(C2)=P(\text{yes})=9/14$$

A2: temperature

$$P(\text{hot}|C1)=P(\text{hot},C1)/P(C1)=2/5$$

$$P(\text{hot}|C2)=P(\text{hot},C2)/P(C2)=2/9$$

$$P(\text{mild}|C1)=P(\text{mild},C1)/P(C1)=2/5$$

$$P(\text{mild}|C2)=P(\text{mild},C2)/P(C2)=4/9$$

$$P(\text{cool}|C1)=P(\text{cool},C1)/P(C1) \\ = (2/14)/(5/14)=1/5$$

$$P(\text{cool}|C2)=P(\text{cool},C2)/P(C2)=2/9$$

# Clasificatorul Naïve Bayes

Exemplu:

| Relation: weather.symbolic |                    |                        |                     |                  |                 |
|----------------------------|--------------------|------------------------|---------------------|------------------|-----------------|
| No.                        | outlook<br>Nominal | temperature<br>Nominal | humidity<br>Nominal | windy<br>Nominal | play<br>Nominal |
| 1                          | sunny              | hot                    | high                | FALSE            | no              |
| 2                          | sunny              | hot                    | high                | TRUE             | no              |
| 3                          | overcast           | hot                    | high                | FALSE            | yes             |
| 4                          | rainy              | mild                   | high                | FALSE            | yes             |
| 5                          | rainy              | cool                   | normal              | FALSE            | yes             |
| 6                          | rainy              | cool                   | normal              | TRUE             | no              |
| 7                          | overcast           | cool                   | normal              | TRUE             | yes             |
| 8                          | sunny              | mild                   | high                | FALSE            | no              |
| 9                          | sunny              | cool                   | normal              | FALSE            | yes             |
| 10                         | rainy              | mild                   | normal              | FALSE            | yes             |
| 11                         | sunny              | mild                   | normal              | TRUE             | yes             |
| 12                         | overcast           | mild                   | high                | TRUE             | yes             |
| 13                         | overcast           | hot                    | normal              | FALSE            | yes             |
| 14                         | rainy              | mild                   | high                | TRUE             | no              |

$$P(C1)=P(\text{no})=5/14 \quad P(C2)=P(\text{yes})=9/14$$

A3: humidity

$$P(\text{high}|C1)=P(\text{high},C1)/P(C1)=4/5$$

$$P(\text{high}|C2)=P(\text{high},C2)/P(C2)=3/9$$

$$P(\text{normal}|C1)=P(\text{normal},C1)/P(C1)=1/5$$

$$P(\text{normal}|C2)=P(\text{normal},C2)/P(C2)=6/9$$

# Clasificatorul Naïve Bayes

Exemplu:

| Relation: weather.symbolic |                    |                        |                     |                  |                 |
|----------------------------|--------------------|------------------------|---------------------|------------------|-----------------|
| No.                        | outlook<br>Nominal | temperature<br>Nominal | humidity<br>Nominal | windy<br>Nominal | play<br>Nominal |
| 1                          | sunny              | hot                    | high                | FALSE            | no              |
| 2                          | sunny              | hot                    | high                | TRUE             | no              |
| 3                          | overcast           | hot                    | high                | FALSE            | yes             |
| 4                          | rainy              | mild                   | high                | FALSE            | yes             |
| 5                          | rainy              | cool                   | normal              | FALSE            | yes             |
| 6                          | rainy              | cool                   | normal              | TRUE             | no              |
| 7                          | overcast           | cool                   | normal              | TRUE             | yes             |
| 8                          | sunny              | mild                   | high                | FALSE            | no              |
| 9                          | sunny              | cool                   | normal              | FALSE            | yes             |
| 10                         | rainy              | mild                   | normal              | FALSE            | yes             |
| 11                         | sunny              | mild                   | normal              | TRUE             | yes             |
| 12                         | overcast           | mild                   | high                | TRUE             | yes             |
| 13                         | overcast           | hot                    | normal              | FALSE            | yes             |
| 14                         | rainy              | mild                   | high                | TRUE             | no              |

$$P(C1)=P(\text{no})=5/14 \quad P(C2)=P(\text{yes})=9/14$$

A4: windy

$$P(\text{FALSE}|C1)=P(\text{FALSE},C1)/P(C1)=2/5$$

$$P(\text{FALSE}|C2)=P(\text{FALSE},C2)/P(C2)=6/9$$

$$P(\text{TRUE}|C1)=P(\text{TRUE},C1)/P(C1)=3/5$$

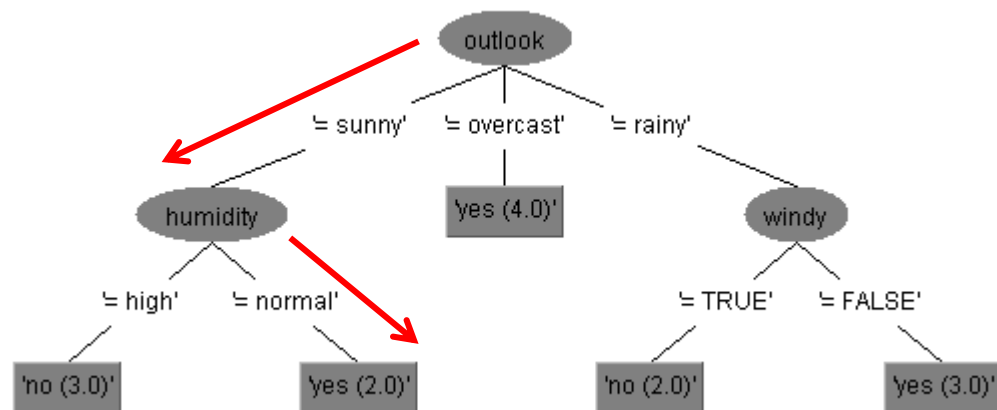
$$P(\text{TRUE}|C2)=P(\text{TRUE},C2)/P(C2)=3/9$$

# Clasificatorul Naïve Bayes

Exemplu:

Relation: weather.symbolic

| No. | outlook<br>Nominal | temperature<br>Nominal | humidity<br>Nominal | windy<br>Nominal | play<br>Nominal |
|-----|--------------------|------------------------|---------------------|------------------|-----------------|
| 1   | sunny              | hot                    | high                | FALSE            | no              |
| 2   | sunny              | hot                    | high                | TRUE             | no              |
| 3   | overcast           | hot                    | high                | FALSE            | yes             |
| 4   | rainy              | mild                   | high                | FALSE            | yes             |
| 5   | rainy              | cool                   | normal              | FALSE            | yes             |
| 6   | rainy              | cool                   | normal              | TRUE             | no              |
| 7   | overcast           | cool                   | normal              | TRUE             | yes             |
| 8   | sunny              | mild                   | high                | FALSE            | no              |
| 9   | sunny              | cool                   | normal              | FALSE            | yes             |
| 10  | rainy              | mild                   | normal              | FALSE            | yes             |
| 11  | sunny              | mild                   | normal              | TRUE             | yes             |
| 12  | overcast           | mild                   | high                | TRUE             | yes             |
| 13  | overcast           | hot                    | normal              | FALSE            | yes             |
| 14  | rainy              | mild                   | high                | TRUE             | no              |



$D=(\text{outlook}=\text{sunny}, \text{temperature}=\text{mild}, \text{humidity}=\text{normal}, \text{windy}=\text{False})$

$$P(C1|D)=P(\text{sunny}|C1)*P(\text{mild}|C1)*P(\text{normal}|C1)*P(\text{FALSE}|C1)*P(C1)/P(D)=$$

$$=3/5*2/5*1/5*2/5*5/14/P(D)=60/8750/P(D) = 0.006875/P(D)$$

$$P(C2|D)=P(\text{sunny}|C2)*P(\text{mild}|C2)*P(\text{normal}|C2)*P(\text{FALSE}|C2)*P(C2)/P(D)=$$

$$=2/9*4/9*6/9*6/9*9/14/P(D)=\mathbf{2592/91854=0.028219/P(D)} \rightarrow \mathbf{yes}$$

# Clasificatorul Naïve Bayes

## Exemplu:

| Relation: weather.symbolic |                    |                        |                     |                  |                 |
|----------------------------|--------------------|------------------------|---------------------|------------------|-----------------|
| No.                        | outlook<br>Nominal | temperature<br>Nominal | humidity<br>Nominal | windy<br>Nominal | play<br>Nominal |
| 1                          | sunny              | hot                    | high                | FALSE            | no              |
| 2                          | sunny              | hot                    | high                | TRUE             | no              |
| 3                          | overcast           | hot                    | high                | FALSE            | yes             |
| 4                          | rainy              | mild                   | high                | FALSE            | yes             |
| 5                          | rainy              | cool                   | normal              | FALSE            | yes             |
| 6                          | rainy              | cool                   | normal              | TRUE             | no              |
| 7                          | overcast           | cool                   | normal              | TRUE             | yes             |
| 8                          | sunny              | mild                   | high                | FALSE            | no              |
| 9                          | sunny              | cool                   | normal              | FALSE            | yes             |
| 10                         | rainy              | mild                   | normal              | FALSE            | yes             |
| 11                         | sunny              | mild                   | normal              | TRUE             | yes             |
| 12                         | overcast           | mild                   | high                | TRUE             | yes             |
| 13                         | overcast           | hot                    | normal              | FALSE            | yes             |
| 14                         | rainy              | mild                   | high                | TRUE             | no              |

**Obs:** dacă pt o anumită valoare de atribut ( $a_{ij}$ ) și o anumită clasă  $C_k$  nu există exemplu în setul de antrenare, atunci  $P(a_{ij} | C_k) = 0$  și (datorită ipotezei de independență) pt orice instanță având valoarea  $a_{ij}$  pt atributul  $A_i$ , probabilitatea să aparțină clasei  $C_k$  este 0.

Această situație poate să apară în special în cazul claselor mici.

Tratarea acestor situații prin regula de “netezire”:

$$P(a_{ij} | C_k) = (\text{count}(a_{ij}, C_k) + \alpha) / (\text{count}(C_k) + m_i * \alpha)$$

$\alpha$  = parametru de netezire Laplace (exemplu:  $\alpha=1$ )

$m_i$  = nr de valori distincte ale atributului  $A_i$

# Clasificatorul Naïve Bayes

## Obs:

- Acest model poate fi aplicat direct atributelor discrete și se bazează pe unul din următoarele modele probabiliste:
  - Binomial (pentru attribute binare)
  - Multinomial (pentru attribute discrete)
- În cazul atributelor numerice care iau valori într-un domeniu continuu există două abordări principale:
  - Attributele sunt **discretizate** înainte de utilizarea clasificadorului (performanța acestuia depinde de procesul de discretizare)
  - Se folosesc modele probabiliste continue (e.g. Gaussian) cu parametri estimați pe baza setului de antrenare



# Rețele Bayesiene pentru clasificare

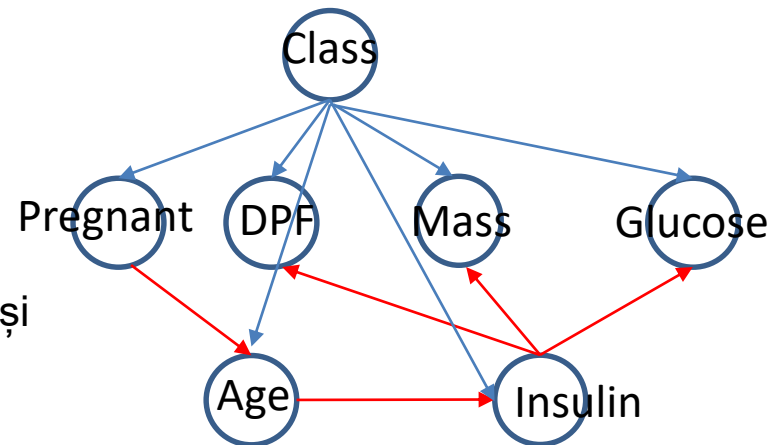
## Motivație:

- Ipoteza referitoare la independența atributelor poate să fie prea restrictivă
  - modelul Naive Bayes corespunde unei **structuri arborescente cu un singur nivel**
  - interacțiunile dintre atribute pot fi modelate utilizând **rețele Bayesiene** (grafuri orientate aciclice):
    - Nod = atribut
    - Arc = corespunde interacțiunii dintre atribute

## Exemplu (Pima dataset)

[Friedman et al, Bayesian Network Classifiers, 1997]

- Naïve Bayes:
  - săgeți albastre (interacțiunea dintre atributul de clasă și atributele predictive)
- Rețea Bayesiană:
  - săgeți albastre + săgeți roșii (interacțiuni între atribute predictive)



# Rețele Bayesiene pentru clasificare

- **Tree Augmented Naive Bayes** (TAN) – [Friedman et al, Bayesian Network Classifiers, 1997]
- **Idee**
  - Se calculează informația mutuală dintre oricare două atribute și se utilizează ca pondere pentru a defini o rețea de interacțiuni între atribute
  - Se determină arborele de acoperire de pondere maximă (maximum weight spanning tree)
  - Se stabilesc orientările muchiilor astfel încât acestea să pornească de la nodul corespunzător atributului “clasa”
- **Implementări**
  - R – bnlearn (<https://www.bnlearn.com/examples/classifiers/>)
  - Python – pyAgrum (<https://pyagrum.readthedocs.io/en/latest/skbnClassifier.html>)

# Sumar

- Clasificatori bazați pe instanțe (k-Nearest Neighbour)
  - **Avantaj:** Proces de antrenare simplu (doar stocarea exemplelor)
  - **Dezavantaj:**
    - Costul clasificării poate fi mare dacă nu sunt utilizate structuri eficiente de căutare (pt eficientizare se pot utiliza arbori k-d sau tabele de hashing care țin cont de similaritate)
    - Performanța depinde de măsura de similaritate și de valoarea lui k

# Sumar

- Clasificatori de tip Naive Bayes
  - Avantaj:
    - Ușor de construit (se bazează doar pe calcule de frecvențe)
    - Nu necesită un volum mare de date de antrenate
    - Sunt eficienți în faza de clasificare
  - Dezavantaj:
    - În varianta naivă nu țin cont de interacțiunile dintre atribute
    - Valorile estimate ale probabilităților trebuie tratate cu grijă (nu este recomandat să fie utilizate ca atare ci doar în contextul deciziei legat de clasă și care este bazată doar pe compararea valorilor)

# Curs următor

Modele funcționale (bazate pe transformări ale datelor folosind diverse funcții)

- Rețele neuronale (Feedforward Neural Networks – Multilayer Peceptrons)
- Modele bazate pe vectori suport (Support Vector Machines)