

2. K-NN sailkatzaleak

Irakasgai: Datu Meatzaritza

Espazialitatea: **Konputazioa**, hitzgaren ikasmaila
Titulazioa: **Informazioa Ingenieritzako Gradua**
Konputazio Zientzia eta Aplikazioen Fakultatea
Universidad del País Vasco - Euskal Herriko Unibertsitatea

- ▶ **K-NN:** (K-Nearest Neighbour)
- ▶ Kasu berri bat sailkatzetakoan bere hurbileneko K auzokideen klaseen artean sarrien agertzen den klasea egokituko zaio.
- ▶ Ideia simple eta intuitiboa. Implementazio erraza
- ▶ Ez dago eredu esplizitukik. Distantzien kalkuluan oinarritzen da (ikus <http://es.wikipedia.org/wiki/Distancia>)

Sarrera

Notazioa

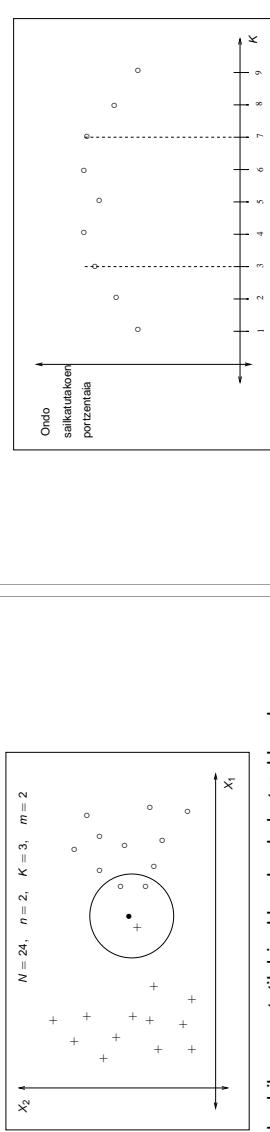
D datubasea, X_1, \dots, X_n aldagai iragarleak, C klase-aldagaiak.
 N kasu, \mathbf{x} kasu berria.

| D | \mathbf{x}_1, c_1 | $x_{11} \dots x_{ij} \dots x_{in} \dots x_n$ | C |
|-----------------------|---------------------|--|----------|
| \vdots | \vdots | \vdots | \vdots |
| (\mathbf{x}_i, c_i) | i | $x_{i1} \dots x_{ij} \dots x_{in} \dots c_i$ | \vdots |
| \vdots | \vdots | \vdots | \vdots |
| (\mathbf{x}_N, c_N) | N | $x_{N1} \dots x_{Nj} \dots x_{Nn} \dots c_N$ | \vdots |
| \mathbf{x} | | $x_1 \dots x_j \dots x_n$ | ? |

Oinarrizko K-NN algoritmoa

Adibide grafikoa: 3-NN, 24 kasu, 2 klase, 2 aldagai iragarle

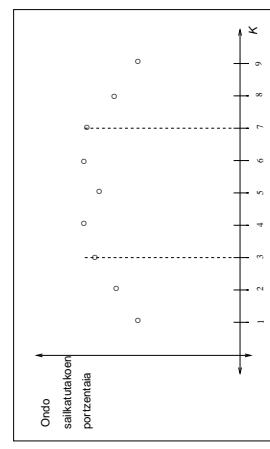
Algoritmoaren sasikodea



HASTERA
Sarrera: $D = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_N, c_N)\}$
 $\mathbf{x} = (x_1, \dots, x_n)$ sailkattutako (x_i, c_i) guztietarako
 kalkulatu $d_i = d(\mathbf{x}_i, \mathbf{x})$
 Ordenatu $d_i (i = 1, \dots, N)$ txikienetik handienera
 \mathbf{x} -tik hurbilien dauden K kasuekin geratu, $D^K_{\mathbf{x}}$
 \mathbf{x} -ri $D^K_{\mathbf{x}}$ -ko sarrieneko klasea esleitu
 AMAIERA

Oinarrizko K-NN algoritmoa

K parametroaren aukeraketa



Erlazioa ez da monotonin. K bakoitia aukeratzea komendi da
 3 auzokide hurbilien artelik bi o klasekoak, bat + klasekoak
 3-NN sailkatzalle honek o iragarriko du • kasuarentzat
 (K=1 baltz?)

K-NN algoritmoaren aldaera

1. Aldaera: K-NN bermegabeak baztertuz

- 1. K-NN bermegabeak baztertuz
- 2. K-NN batazbesteko distantziarekin
- 3. K-NN distantzia minimoarekin
- 4. K-NN aukeratutako auzokideen pisaketarekin
- 5. K-NN aldagaien pisaketarekin

7/17

9/17

▷ Batazbesteko distantziarik txikieneiro klasoa esleitu

Batazbesteko distantziarik txikieneiro klasoa esleitu
Aldaketa. 7 auzokide hurbilena hartuz

K-NN bermegabeak baztertuz

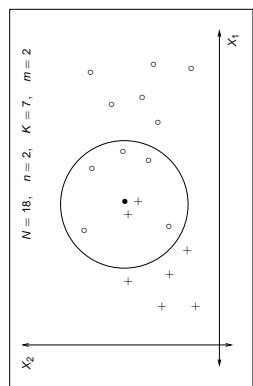
- ▶ Kasu berri bat sailkatzeko nahiko berme izan behar da
- ▶ Ez izatekotan, **kasua sailkatu gabe utz daiteteke**
- ▶ Aurrelik **atxikia bat** ezartzen da
- ▶ Gehiengo absolutua

• kasu berriari + klasoa esleitu zaio

8/17

9/17

2. Aldaera: K-NN batazbesteko distantziarekin



N = 18, n = 2, K = 7, m = 2

Batazbesteko distantziarik txikieneiro klasoa esleitu
Aldaketa. 7 auzokide hurbilena hartuz

3. Aldaera: K-NN distantzia minimoarekin

- ▶ Klase bakoitzeko **ordezkari bat** aukeratu (adb. klasoko banzerutik hurbilien dagoen kasua)
- ▶ Gorde beharreko kasuen fitxategiaien dimentsioa N-tik m-ra murritzuko da
- ▶ 1-NN algoritmoa aplikatu murritzutako fixategiari
- ▶ Eraginkortasuna klase barruko homogenotasunaren mende dago
- ▶ Klase bateko kasuak multzo edo cluster batean baino gehiagotan banatuta badau, aldaera ez da egokia geratzzen, ordezkarria ez delako klasaren adierazgarri egokia izango

10/17

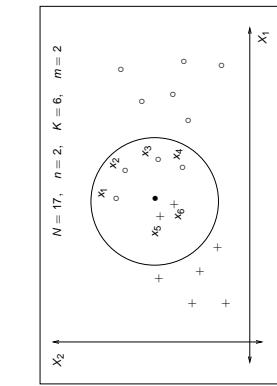
▷ Klasoa esleitu

11/17

K-NN auzokideen pisaketarekin. Adibidea

- ▶ Pisua distantziaren alderantzikro proporcionala
- ▶ Adibidea

| Auzokide hurbilena | Klasoa C | Distantzia $d(\mathbf{x}_i, \mathbf{x})$ | Pisua w_i |
|--------------------|----------|--|---------------------|
| \mathbf{x}_1 | O | 2 | $\frac{1}{2} = 0,5$ |
| \mathbf{x}_2 | O | 2 | $\frac{1}{2} = 0,5$ |
| \mathbf{x}_3 | O | 2 | $\frac{1}{2} = 0,5$ |
| \mathbf{x}_4 | O | 2 | $\frac{1}{2} = 0,5$ |
| \mathbf{x}_5 | + | 0,7 | $1/0,7 = 1,42$ |
| \mathbf{x}_6 | + | 0,8 | $1/0,8 = 1,25$ |



• klasoa esleitu + klasoa esleitu zaio

12/17

4. Aldaera: K-NN auzokideen pisaketarekin

- ▶ Kasu guziek ez dute garanzia bera. **Auzokide hurbilenei pisua bat egokitu** (distantziaren mendekoa, adibidez)
- ▶ Pisuen baturan oinarrituz emango da klasaren iragarpena

Adibidea

• Pisua distantziaren alderantzikro proporcionala

13/17

K-NN auzokideen pisaketarekin. Adibidea

- ▶ Pisua distantziaren alderantzikro proporcionala
- ▶ Adibidea

| Auzokide hurbilena | Klasoa C | Distantzia $d(\mathbf{x}_i, \mathbf{x})$ | Pisua w_i |
|--------------------|----------|--|---------------------|
| \mathbf{x}_1 | O | 2 | $\frac{1}{2} = 0,5$ |
| \mathbf{x}_2 | O | 2 | $\frac{1}{2} = 0,5$ |
| \mathbf{x}_3 | O | 2 | $\frac{1}{2} = 0,5$ |
| \mathbf{x}_4 | O | 2 | $\frac{1}{2} = 0,5$ |
| \mathbf{x}_5 | + | 0,7 | $1/0,7 = 1,42$ |
| \mathbf{x}_6 | + | 0,8 | $1/0,8 = 1,25$ |

• klasoa esleitu + klasoa esleitu zaio

14/17

• klasoa esleitu + klasoa esleitu zaio

15/17

5. Aldaera: K-NN aldagaien pisaketarekin

Aldagai iragarleen pisaketa

- Aldagai guztiek C klase-aldagaiaren iragarpenerako garantzia bera ez dutenean erabilitzea komeni da
- Distantzia eta distantzia ponderatua:

$$d(\mathbf{x}, \mathbf{x}_r) = \sqrt{\sum_{j=1}^n w_j (x_j - x_{rj})^2}$$

- X_j aldagaiari w_j pisua egokitu, X_j aldagaiaren eta C klase-aldagaiaren (X_j, C) elkarrekiko informazioaren neurria erabiliz, adibidez. Informazio-teoria

Kasuen aukeraketa

- Datubaseetik kasuen aukeraketa egin daitete, sailkatzaille **azkarragoak** eta **hobeak** lortzen.
- Distantzien kalkuluan oinarritutako K-NN teknikek **kostu konputaziozial handia** izaten dute datubase handietarako.
- Datubasean gaizki etiketatuako **kasuak** badade, K-NN sailkatzaillearen **asmatze-tasa** jaitsi egindo da.
- Bi teknika: **edizio-teknikak** eta **kondentsazio-teknikak**.
- **Edizio-teknikak**: datubasean zaratza sortzen duen kasuak kenduz, erabaki-muga leunak sortzen salaitzen dira.
- **Kondentsazio-teknikak**: jatorrizko datubaseko kasuetatik indizitutako erabaki-mugak mantentzen salaitzen dira
- Datubase murriztua erabiliko da sailkatzalea sortzea

Ediziorkako teknikak: Wilson-en edizioa

- Klaseen arteko mugak garbitu, gainjarritako kasuak kenduz
- K-NN sailkatzaillek gaizki sailkatuko lituzkeen kasuak datubasetik kendu
- **Wilson-en ediziorkako sasikodea**
- Errepikatu i guztietaarako
 - Datubaseko $\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_N$ kasuen artean **K auzokide hurbilenak aurkitu**
 - \mathbf{x}_i kasuari bere auzokide hurbilenean artean **sarrien agertzen den klasea esleitu**
 - Esleitutako klasea eta klase erreala bat **ez datozen kasu guztiak datubasetik ezabatu**

Kondentsazio-teknikak: Hart-en kondentsazioa

Hart-en kondentsaziorkako sasikodea

1. Datubaseko lehenengo kasua **STORE** zerrendan gordetzea
2. Datubaseko bigarren kasuaren iragarpena egin **STORE**ko kasuetan oinarrituz. Sailkapena zuzen bada, **GRABBAG** zerrendan sartu, bestela **STORE**en sartu
3. Modu berean jarraitu datubaseko kasu guztiekin
4. Datubaseko kasu guztik aztertu ondoren, **GRABBAG** zerrendakoekin prozedura errepikatu behin eta berri. Prozesuaren amaiera bi modutara lortu daiteke:
 - **GRABBAG** zerrenda hutsa dago (kasu guztik **STORE** zerrendan daudel)
 - **GRABBAG** zerrendako kasu guztik aztertu ondoren bat berri ere **STORE** zerrendara ez da pasa

5. **GRABBAG**-eko kasuak datubasetik ezabatu

Oinarrizko bibliografia

- Liburua: **Aprendizaje Automático: conceptos básicos y avanzados**
- **Capítulo 3: Algoritmos de clasificación por vecindad**
- Koordinatzailea: Basilio Sierra Araujo
- Departamento de Ciencias de la Computación e Inteligencia Artificial, Universidad del País Vasco
- Argitaratzalea: Pearson Prentice Hall, 2004
- ISBN 10: 84-8322-318-X, ISBN 13: 978-84-8322-318-5