

## 1. Sarrera

### 7. Aldagai Aukeraketa

Espesialitatea: Konputazioa, hirugarren ikasmarria  
Titulazioa: **Informazioa Ingelesezko Gradua**  
Konputazio Zientzia eta Aulmen Artifiziala Saia  
Euskal Herriko Unibertsitatea

#### Aldagaien azpimultzoaren aukeraketa

- ▶ Feature Subset Selection (**FSS**): Aldagaien azpimultzoaren aukeraketa.
- ▶ Eredu sailkatzaleak induzitzeko aldagai iragarleen azpimultzo egokia aukeratzea beharrezko eta egokia geratzen da askotan.
- ▶ Aukeraketa egiteko **bi filosofia** nagusi daude:
  - ▶ Filter edo iragazki modukoak
  - ▶ Wrapper edo bilgari modukoak

## 2. Zergatik egin aldagai iragarleen aukeraketa?

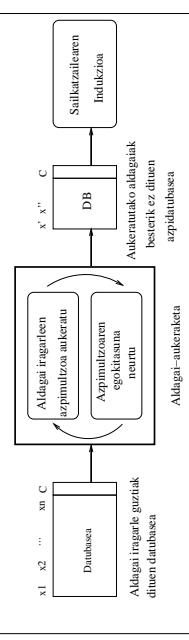
- Ez da egia aldagai iragarrie gehiago erabilizeagatik, inditzitutako sailkatzalea hobea izango denik**
- ▶ Aldagai iragarie batzuk ez dira esanguratsuak
    - ▶ Aldagai iragarie baten balioa ezagutzeak ez badu C klase-aldagaiaren iragarpenean lagunzen, esanguratsua ez dela esaten da
    - ▶ Aldagai iragarie batzuk beren artean **erredundanteak** dira badateke, erredundantea dela esaten da
      - ▶ Aldagai iragarie bat beste aldagai iragarbetatik ondoriozta diren eta erredundanteak dieren aldagaiak identifikatza
  - Aldagai-aureraketa egitearen helburua: esanguratsuak ez diren eta erredundanteak dieren aldagaiak identifikatza
  - Datuak ondo adierazten dituzten bi ereduren artean, simpleena aukeratu

## 3. Aldagai-aureraketa egitearen abantailak

### Filter Metodoa

- Filter Metodoa: iragazki moduko metodoa**
- ▶ Aldagai iraganleek C klase-aldagaiaren iragarpenerako ematen duten **informazioa neurtuko** da: elkarrekiko informazioa, irabazi-ratiora, chi karratu,...
  - ▶ Neurri horren arabera **aldagai iragariek ordenatu**ko dira, esanguratsuna denetik hasi, eta informazio gutxien ematen duenea
  - ▶ Aldagai iragarie **esanguratsuenak aukeratuko** dira, eta aukeratutako aldagai iragarie horiekin egindo da **eredu sailkatzalearen indukzioa**
  - ▶ Aldagai iraganleak ordenatzeko erabili den neurriak ez du kontuan hartzen, eta **deseegokia** izan daiteke. Gainera, aldagaiaik banaka aztertzen badira, **ez dira aldagaien arteko elkarrekintzak kontuan hartzen**. Binaka aztertu? Hirunaka?

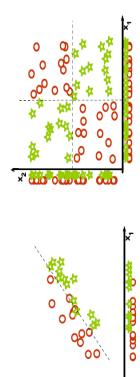
## 4. Aldagai-aureraketako metodoak



- ▶ Abantailak: Konputazionalki **azkarra**. Gainera, eredu sailkatzaleagandik independentea izateagatik, aldagai-aureraketa **behin bakarrik** egindo da, eta aldagai horietan oinarrituz, nahi adina eredu sailkatzalle indizitutiko dira desabantalak: aldagai-aureraketa eredu sailkatzaleagandik independente egiteagatik, ez da beraien arteko elkarrekintza kontuan hartzen, eta **deseegokia** izan daiteke. Gainera, aldagaiaik banaka aztertzen badira, **ez dira aldagaien arteko elkarrekintzak kontuan hartzen**. Binaka aztertu? Hirunaka?

## Aldagaien arteko elkarrekintzak baztertuz gero,

*Univariate selection may fail*



Oyvind-Eriksen et al., JMLR 2004; Springer 2006

7/23

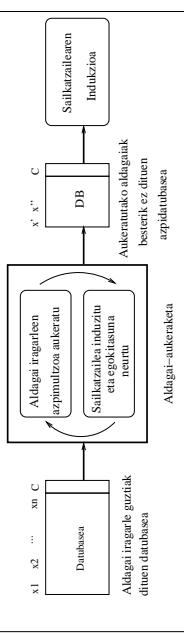
## 4. Aldagai-aukeraketarako metodoak

### Wrapper Metodoa: bilgarri moduko metodoa

- ▶ Aldagai iragarleen azpimultzoak, **azpimultzo horretan onarritzutako eredu sailkatzalearen bidez ebaluatuak** izango dira
- ▶ n aldagai iragarle izanik,  $2^n$  **azpimultzo** sor daitezke, eta horietako onena zein den erabaki beharko da
- ▶ Azpimultzo kopurua, aldagai iragarle kopurarekin batera, exponentzialki hazten da. **Ezinezkoa** geratzen da aldagaien azpimultzo guziletarako sailkatzalea entrenatu eta testeatzea, guztiuen arteko sailkatzalea entrenatzeko
- ▶ Aldagaien **azpimultzo optimoaren bilaketa** zenbait algoritmo heuristikorik erabiliz egin ohi da, halako, algoritmo genetikoak, etab

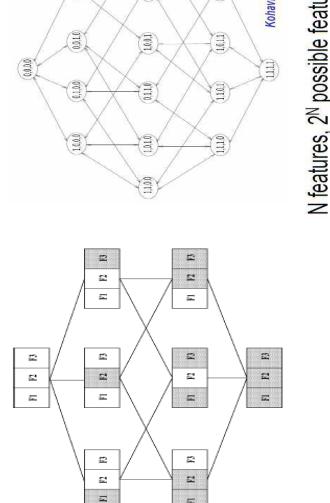
8/23

### Wrapper Metodoa



9/23

## Aldagaien azpimultzo optimoaren bilaketa

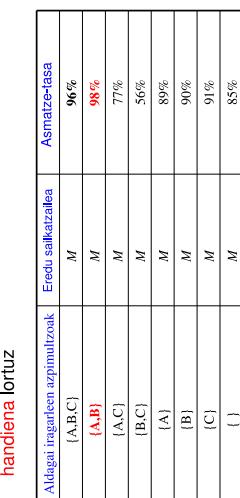


Kohavi-John, 1997  
N features,  $2^N$  possible feature subsets!

10/23

## Aldagaien azpimultzo optimoaren bilaketa

Aldagai iragarle kopuru handietarako, ez da posible azpimultzo guztiak azterearaztutako optimoaren bilaketa. Algoritmo heuristikoren bat erabili beharko da

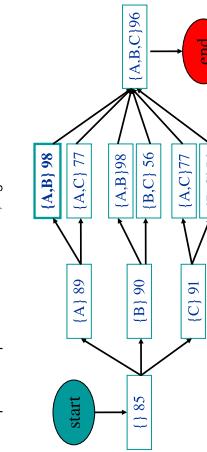


Bilaketa heuristikorekin, agian ez dugu azpimultzo optima aukerik

11/23

## Aldagaien azpimultzo optimoaren bilaketa

Aldagai iragarle kopuru handietarako, ez da posible azpimultzo guztiak azterearaztutako optimoaren bilaketa. Algoritmo heuristikoren bat erabili beharko da



12/23

## 5. Bestelako metodoak

- ▶ **Embedded methods.** Eredu sailkatzalea eraikitzearekin batera egiten da aldagaien azpimultzo optimoaren bilaketa. Adibidez, sailkapen-zuhaitzak. Wrapper metodoetan bezala, aldagai-aukeraketa eredu sailkatzalearentzat bereziki egiten da, wrapper metodoak bezain garestik izan gabe
- ▶ Aldagai-aukeraketa egin beharean, **aldagai iragarleei pisu bat egoki** dakiene
- ▶ **Datu en dimentsioaren murrizketa** egin daitene aldagai-aukeraketa egin gabe, jatorrizko aldagai iragarleetik beste aldagai batzuk eratorriz. Teknika desberdinak daude. Adibidez, transformazio linealen bidez, **SVD (Singular Value Decomposition)**

## SVD. Testuen sailkapenerako adibidea

### Testuen sailkapena

- ▶ Ohiko aplikazioak: Posta elektronikoan zaborra diren mezuk antzematea, internet bidezko dokumentuen bilaketa, hitzen adieren desanbignazioa, etab
- ▶ **Bag of Words** bidezko adierazpena. **Dokumentua** hiztegiaren tamainako **bektorrearen bidez** adierazten da, dokumentuan hitz baloritzak duen agerpenarenak
- ▶ Ohikoak da hiztegiak **ehundaka milaka hitz** izatea; sarritasun txikieneko eta handieneko hitzak hasieratik kenduta, 15000 inguru jartsi daiteke hiztegiaren tamaina
- ▶ 5000-800000 dokumentuko datubaseak izaten dira erabilgarri ikerketarako

## SVD. Adibidea.

### Dokumentuak

- ▶ c1-c5: "gizaki-konputagailu arteko interakcioa"
- ▶ m1-m4: "grafoen teoria"
- ▶ **c1:** Human machine interface for Lab ABC computer applications.  
**c2:** A survey of user opinion of computer system response time.  
**c3:** The EPS user interface management system.
- ▶ **c4:** System and human system engineering testing of EPS.
- ▶ **c5:** Relation of user-perceived response time to error measurement.
- ▶ **m1:** The generation of random, binary, unordered trees.
- ▶ **m2:** The intersection graph of paths in trees.
- ▶ **m3:** Graph minors IV: Widths of trees and well-quasi-ordering.
- ▶ **m4:** Graph minors: A survey.

13/23

14/23

15/23

## SVD adibidea: Dokumentuen adierazpena

### M matrizea: terminoak x dokumentuak

$$M = \begin{pmatrix} 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 2 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix}$$

### SVD: Singular Value Decomposition

$$M=U\Sigma V^T = \begin{pmatrix} 0.22 & -0.11 & 0.29 & -0.41 & -0.11 & -0.34 & 0.52 & -0.06 & -0.41 \\ 0.20 & -0.07 & 0.14 & -0.55 & 0.28 & 0.50 & -0.07 & -0.01 & -0.11 \\ 0.24 & 0.04 & -0.16 & -0.59 & -0.11 & -0.25 & -0.30 & 0.06 & 0.49 \\ 0.40 & 0.06 & -0.34 & 0.10 & 0.33 & 0.38 & 0.00 & 0.01 & 0 \\ 0.64 & 0.17 & 0.36 & 0.33 & -0.16 & -0.21 & -0.17 & 0.03 & 0.27 \\ 0.27 & 0.11 & -0.43 & 0.07 & 0.08 & -0.17 & 0.28 & -0.02 & -0.05 \\ 0.27 & 0.11 & -0.43 & 0.07 & 0.08 & -0.17 & 0.28 & -0.02 & -0.05 \\ 0.30 & -0.14 & 0.33 & 0.19 & 0.11 & 0.27 & 0.03 & -0.17 & 0 \\ 0.21 & 0.27 & -0.18 & -0.03 & -0.54 & 0.08 & -0.47 & -0.04 & -0.58 \\ 0.01 & 0.49 & 0.23 & 0.03 & 0.59 & -0.39 & -0.29 & 0.25 & -0.23 \\ 0.04 & 0.62 & 0.22 & 0.00 & -0.07 & 0.11 & 0.16 & -0.68 & 0.23 \\ 0.03 & 0.45 & 0.14 & -0.01 & -0.30 & 0.28 & 0.34 & 0.68 & 0.18 \end{pmatrix}$$

### SVD: Singular Value Decomposition

$$M=U\Sigma V^T = \begin{pmatrix} 3.34 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2.54 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2.35 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1.64 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1.50 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1.31 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.85 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.56 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\Sigma =$$

16/23

17/23

18/23

## SVD: Singular Value Decomposition

### 2 dimentsiotara murriztua, $M \rightarrow M_2$

$$M = U \Sigma V^T$$

$$V = \begin{pmatrix} 0.20 & -0.06 & 0.11 & -0.95 & 0.05 & -0.08 & 0.18 & -0.01 & -0.06 \\ 0.61 & 0.17 & -0.50 & -0.03 & -0.21 & -0.26 & -0.43 & 0.05 & 0.24 \\ 0.46 & -0.13 & 0.21 & 0.04 & 0.38 & -0.27 & 0.21 & 0.01 & 0.02 \\ 0.54 & -0.23 & 0.57 & 0.27 & -0.21 & -0.37 & 0.26 & -0.02 & -0.08 \\ 0.28 & 0.11 & -0.51 & 0.15 & 0.33 & 0.03 & 0.67 & -0.06 & -0.26 \\ 0.00 & 0.19 & 0.10 & 0.02 & 0.39 & -0.30 & -0.34 & 0.45 & -0.62 \\ 0.01 & 0.44 & 0.19 & 0.02 & 0.35 & -0.21 & -0.15 & -0.76 & 0.02 \\ 0.02 & 0.62 & 0.25 & 0.01 & 0.15 & 0.00 & 0.25 & 0.45 & 0.52 \\ 0.08 & 0.53 & 0.08 & -0.03 & -0.60 & 0.36 & 0.04 & -0.07 & -0.45 \end{pmatrix}$$

$$M_2 = U_2 \Sigma_2 V_2^T = \sum_{i=1}^2 \sigma_i u_i v_i^T = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T$$

19/23

20/23

## SVD: 2 dimentsioko espazio murriztua

### Dokumentuen adierazpen grafikoa planoan

d dokumentu bakoitzaren proiekzioa planoan

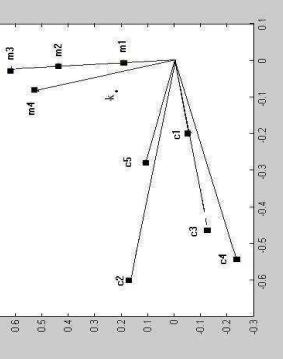
$$\mathbf{d}_2 = \mathbf{d}^T \mathbf{U}_2 \Sigma_2^{-1}$$

Adibideko 9 dokumentuen proiekzioak

$$\mathbf{D}_2 = \begin{pmatrix} -0.20 & -0.61 & -0.46 & -0.54 & -0.28 & -0.01 & -0.02 & -0.08 \\ -0.06 & 0.17 & -0.13 & -0.23 & 0.11 & 0.19 & 0.44 & 0.62 \\ \downarrow & \downarrow \\ c1 & c2 & c3 & c4 & c5 & m1 & m2 & m3 & m4 \end{pmatrix}$$

Dokumentuen proiekzioen koordenatuak planoan.

21/23



22/23

## SVD: 2 dimentsioko espazio murriztua

### Bibliografia

- ▶ [http://en.wikipedia.org/wiki/Feature\\_selection](http://en.wikipedia.org/wiki/Feature_selection)
- ▶ “An Introduction to Variable and Feature Selection”. Isabelle Guyon, André Elisseeff.
- ▶ Journal of Machine Learning Research 3, (2003), pp. 1157-1182
- ▶ “A Review of Feature Selection Techniques in Bioinformatics”. Yvan Saeys, Iñaki Inza, Pedro Larrañaga.
- ▶ Bioinformatics, Review, Vol.23, n.19, (2007), pp. 2507-2517

23/23