

## 5. Sailkatzale Bayestarrak

Irakasgaia: Konputazio Zientzien Metodo Matematikoak

Titulazioa: **Informatikan Ingeniaria**  
Konputazio Zientzia eta Adimen Artifiziala saila  
Universidad del País Vasco - Euskal Herriko Unibertsitatea

# Sailkatzaila Bayestarrak

## Sarrera

- Probabilitate teorian (**Bayes-en teorema**) oinarritutako sailkatzailaileak, gainbegiratutako sailkapenerako.
- **Abantailak:** Kasu berri bat izanik, klase-aldagaiaren iragarpenari probabilitate bat egokituko zaio. Beste sailkatzailailekin lortutako emaitzen parekoak lortzen dira
- **Eragozpenak:** Konputazionalki garestiak izan daitezke, probabilitate asko estimatu beharko direlako
- Gai honetan aztertuko ditugunak:
  - Naïve Bayes
  - Zuhaitzera zabaldutako Naïve Bayes (**Tree Augmented Naïve Bayes, TAN**), (Friedman eta abar, 1997)
  - $k$ -mendekotasuneko sailkatzaila Bayestarra ( **$k$ -Dependence Bayesian classifier,  $k$ -DB**) (Sahami, 1996)

# Bayes-en teorema

C iragarri beharreko hipotesia eta  $X$  data ezaguna izanik,

$$p(C|X) = \frac{p(C)p(X|C)}{p(X)} = \frac{p(C)p(X|C)}{\sum_{j=1}^m p(c_j)p(X|c_j)}$$

- $p(C)$ : C klase-aldagaiaren a priori probabilitatea
- $p(X)$ : X dataren a priori probabilitatea
- $p(X|C)$ : X-ren probabilitate baldintzatua C emanda ; X-ren a posteriori probabilitatea
- $p(C|X)$ : C-ren probabilitate baldintzatua, X emanda; C-ren a posteriori probabilitatea

# Eredu probabilistikotik sailkatzalea eraikitzen

Bayes-en teorema sailkapen gainbegiraturako: C klase-aldagaiaren a posteriori probabilitate maximoko balioa

## A posteriori probabilitate maximoa: MAP

$$c^* = \arg \max_{c_j \in C} p(c_j|X) = \arg \max_{c_j \in C} \frac{p(c_j)p(X|c_j)}{\sum_{j=1}^m p(c_j)p(X|c_j)}$$

- C klase-aldagaiaren (hipotesiaren) balio posiblak  $c_1, \dots, c_m$  izanik,  $c_j$  bakoitzerako  $p(c_j|X)$  kalkulatu
- $p(c_j|X)$  handieneko  $c^*$  klasea iragarriko da
- $P(X) = \sum_{j=1}^m p(c_j)p(X|c_j) \rightarrow$  hipotesiekiko independentea

# Parametroen estimazioa

D		$X_1$	...	$X_j$	...	$X_n$	C
$(x_i, c_i)$	1	$x_{11}$	...	$x_{1j}$	...	$x_{1n}$	$c_1$
		⋮	⋮	⋮	⋮	⋮	⋮
$(x_i, c_i)$	$i$	$x_{i1}$	...	$x_{ij}$	...	$x_{in}$	$c_i$
		⋮	⋮	⋮	⋮	⋮	⋮
$(x_N, c_N)$	$N$	$x_{N1}$	...	$x_{Nj}$	...	$x_{Nn}$	$c_N$
$\mathbf{x}$		$x_1$	...	$x_j$	...	$x_n$	?

$\mathbf{x} = (x_1, x_2, \dots, x_n)$  kasu berriaren sailkapenerako:

$$c^* = \arg \max_{c_j=c_1, \dots, c_m} p(c_j | x_1, x_2, \dots, x_n) = \arg \max_{c_j=c_1, \dots, c_m} p(c_j) p(x_1, x_2, \dots, x_n | c_j)$$

- $p(c_j) \rightarrow$  datubasean  $c_j$  klasearen maiztasuna
- $p(x_1, x_2, \dots, x_n | c_j) \rightarrow$  datubasea oso oso handia den kasuan bakarrik estima daiteke

Estimazio kopuru handiegia...

# Naïve Bayes sailkatzalea

## Eredu probabilistikoaren sinplifikazioa

C klase-aldagaiaren baldintzapean **aldagai iragarleak independente** direla suposatzen badugu,

$$p(x_1, x_2, \dots, x_n | c_j) = \prod_{i=1}^n p(x_i | c_j)$$

## Naïve Bayes sailkatzalea

$$c^* = \arg \max_{c_j=c_1, \dots, c_m} p(c_j | x_1, x_2, \dots, x_n) = \arg \max_{c_j=c_1, \dots, c_m} p(c_j) \prod_{i=1}^n p(x_i | c_j)$$

# Naïve Bayes sailkatzalea

## Oharrak

- Estimatu beharreko parametro kopurua izugarri jaitsi da
- Eta, C klase-aldagaiaren baldintzapean aldagai iragarleak ez badira independente?

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | c_j) \neq \prod_{i=1}^n p(\mathbf{x}_i | c_j)$$

Hala ere, Naïve Bayes-ek emaitza onak ematen ditu...

- Eta,  $p(\mathbf{x}_i | c_j)$  probabilitatearen estimazioa kalkulatzerakoan, datubasean  $c_j$  klaseko kasuen artean  $X_i = x_i$  kasurik ez badago?

$$\hat{p}(\mathbf{x}_i | c_j) = 0 \rightarrow \hat{p}(c_j) \prod_{i=1}^n \hat{p}(\mathbf{x}_i | c_j) = 0$$

# Naïve Bayes sailkatzalea

## Oharrak

$\hat{p}(x_i|c_j) = 0$  arazoa ekiditeko, estimazioa horrela kalkulatu:

$$\hat{p}(x_i|c_j) = \frac{N(X_i = x_i, C = c_j) + 1}{N(C = c_j) + k}$$

non

- $N(C = c_j)$ :  $c_j$  klaseko kasu kopurua
- $N(X_i = x_i, C = c_j)$ :  $c_j$  klaseko kasuen artean zenbatek duten  $X_i = x_i$
- $k$ :  $X_i$  aldagai iragarlearen balio kopurua

# Naïve Bayes sailkatzalea

## Oharrak

- Eta, **aldagaiak jarraiak** badira? Honakoa bilatzen da:

$$c^* = \arg \max_c p(C=c) \prod_{i=1}^n f_{X_i|C=c}(x_i|c)$$

non  $f_{X_i|C=c}(x_i|c)$  probabilitate-banaketa normala izan ohi den:

$$f_{X_i|C=c}(x_i|c) \rightsquigarrow \mathcal{N}(x_i; \mu_i, \sigma_i)$$

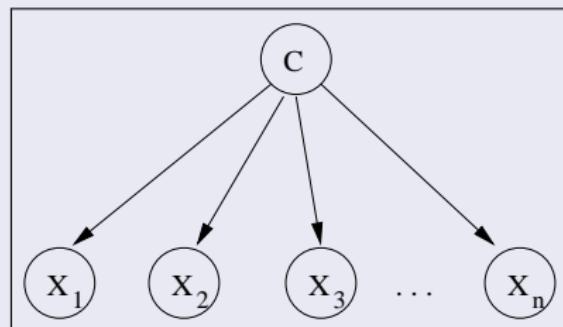
$\mu_i$  itxaropen matematikoa;  $\sigma_i$  desbideratze estandarra.

$$c^* = \arg \max_c p(C=c) \prod_{i=1}^n \frac{1}{\sigma_i \sqrt{2\pi}} e^{\frac{-1}{2} \left( \frac{x_i - \mu_i}{\sigma_i} \right)^2}$$

# Naïve Bayes sailkatzailaren egitura

## Adierazpen grafikoa

C klase-aldagaiaren baldintzaapean **aldagai iragarleak independente** direla suposatzen denez,



$$p(c_j | \textcolor{red}{X}_1, \dots, \textcolor{red}{X}_n) \propto p(c_j) \prod_{i=1}^n p(\textcolor{red}{X}_i | c_j)$$

# Zuhaitzera zabaldutako Naïve Bayes (TAN)

## Tree Augmented Naïve Bayes (TAN)

- Naïve Bayes sailkatzailaren hedapen bat da.
- C klase-aldagaiaren baldintzapean **aldagai iragarleen arteko mendekotasuna** egon daitekeela aintzakotzat hartzen da; elkarrekiko informazioa erabiltzen da

$$I(X, Y|C) = \sum_{i=1}^n \sum_{j=1}^m \sum_{k=1}^r p(x_i, y_j, c_k) \log_2 \frac{p(x_i, y_j|c_k)}{p(x_i|c_k) \cdot p(y_j|c_k)}$$

- Mendekotasunak **zuhaitz egitura** batera mugatzen dira; simplea izaten jarraitzen du

# Zuhaitzera zabaldutako Naïve Bayes (TAN)

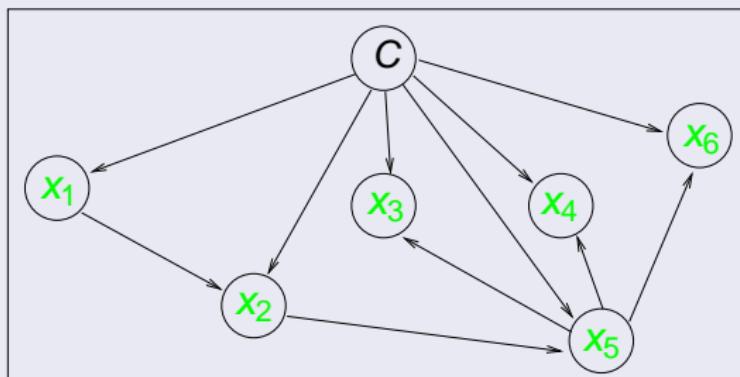
## TAN eredu sortzeko algoritmoa

- 1. urratsa.**  $I(X_i, X_j \mid C)$  kalkulatu  $i < j$  aldagaietarako  $i, j = 1, \dots, n$ .
- 2. urratsa.** Grafo bat definitu. Erpinak  $X_1, \dots, X_n$  aldagai iragarleak izango dira.
- 3. urratsa.**  $I(X_i, X_j \mid C)$  elkarrekiko informazioak handienetik txikienera ordenatu
- 4. urratsa.**  $I(X_i, X_j \mid C)$  handiena aukeratu eta konektatu  $X_i$  eta  $X_j$  erpinak ertz ez-zuzendu baten bidez.
- 5. urratsa.** Hurrengo  $I(X_i, X_j \mid C)$  handiena aukeratu.  $X_i$  eta  $X_j$  erpinak konektatzean zikloa sortzen bada, ez konektatu ertzak.
- 6. urratsa.** Errepikatu 5. urratsa  $n - 1$  ertzen bidez grafo konektatua lortuko den arte (zuhaitza)
- 7. urratsa.** Lortu berri den zuhaitz ez-zuzendua zuhaitz zuzendu bihurtu, erpin aukeratuz zuhaitzaren erro izateko eta horren arabera ertzei noranzkoa jarriz.
- 8. urratsa.** TAN eredu eraiki, C etiketadun erpina gehituz eta C-tik  $X_i$  erpin bakoitzera ertz bat gehituz.

# TAN ereduaren egitura

**Adibidea. Aldagai iragarleak:**  $X_1, X_2, X_3, X_4, X_5, X_6$

$$I(X_1, X_2 | C) > I(X_2, X_5 | C) > I(X_1, X_5 | C) > I(X_3, X_5 | C) > I(X_2, X_3 | C) > \\ I(X_4, X_5 | C) > I(X_1, X_4 | C) > I(X_5, X_6 | C) > \dots$$



$$p(c_j | x_1, \dots, x_6) \propto p(c_j) p(x_1 | c_j) p(x_2 | x_1, c_j) p(x_3 | x_5, c_j) p(x_4 | x_5, c_j) p(x_5 | x_2, c_j) p(x_6 | x_5, c_j)$$

# **k-mendekotasuneko sailkatzaile Bayestarra (k-DB)**

## **k-Dependence Bayesian classifiers, k-DB**

- **Naïve Bayes sailkatzailearen beste hedapen bat** da, non aldagai iragarle bakoitzak gehienez  $k$  aldagai guraso izango dituen,  $C$  klase-aldaagaia kontatu gabe.
- **$k$  parametroarekin** finkatuko dugu ereduan adierazi nahi dugun aldagai iragarleen arteko **mendekotasun maila**
  - Naïve Bayes: 0-mendekotasuneko sailk. Bayestarra
  - TAN eredua: 1-mendekotasuneko sailkatzaile Bayestarra
- **Konputazionalki nahiko eraginkorra** da
- Aldagai iragarleen arteko mendekotasunak maila altuagoan adierazi ahal izateak **emaitza hobeak** ematen ditu zenbait problema konplexutan.
- **$k$  optimoa?**

# $k$ -DB eredua lortzeko algoritmoa

1. **urratsa.**  $I(X_i, C)$  elkarrekiko informazioak kalkulatu,  $i = 1, \dots, n$ .
2. **urratsa.**  $X_j$  eta  $X_j$  pare guzietarako  $i \neq j$ ,  $i, j = 1, \dots, n$ , kalkulatu  $I(X_i, X_j | C)$ .
3. **urratsa.** Erabilitako aldagaien zerrenda,  $\aleph$ , multzo hutsarekin hasieratu.
4. **urratsa.** Hasieratu BN egitura Bayestarra erpin bakarrarekin;  $C$ .
5. **urratsa.**  $\aleph$  multzoan aldagai guztiak egongo diren arte errepikatu:
  - 5.1.  $X_{max}$  aukeratu,  $I(X_{max}, C) = \max_{X \notin \aleph} I(X, C)$ .
  - 5.2. BN sareari erpin bat gehitu  $X_{max}$  aldagaiarentzat.
  - 5.3. BN sarean  $C$  erpinetik  $X_{max}$  erpinera ertza gehitu.
  - 5.4.  $\aleph$  multzoko  $X_j$  aldagaietatik  $I(X_{max}, X_j | C)$  balio handieneko  $m = \min(|\aleph|, k)$  aldagai aukeratu eta ertza gehitu  $X_{max}$  erpinera.
  - 5.5. Gehitu  $X_{max}$  aldagaiaren  $\aleph$  multzoari.
6. **urratsa.** Lortutako BN egiturari jarraituz, sailkapen-eredua adierazi baldintzapeko probabilitateen bidez.

# Oinarrizko bibliografia

- Introducción a la Minería de Datos, 10: **Métodos Bayesianos**, J. Hernández, M<sup>a</sup>J. Ramírez, C. Ferri, Pearson Prentice Hall, 2004
- “Datu Meatzaritza” irakasgaiaren web orriko: “**5. Diagnostikoaren Paradigma Klasikotik Naïve Bayes Sailkatzailera**”
- “An Analysis of Bayesian Classifiers”. Langley et al. 1992
- “Building Classifiers using Bayesian Networks”. Friedman and Goldszmidt. Machine Learning 29:131-163, (1997) → **TAN** ereduaurkezten dute
- “Learning Limited Dependence Bayesian Classifiers”. Sahami. Knowledge Discovery and Data Mining, KDD-96, 335-338, (1996) → **k-DB** ereduaurkezten dute
- [http://en.wikipedia.org/wiki/Bayes'\\_theorem](http://en.wikipedia.org/wiki/Bayes'_theorem)
- [http://en.wikipedia.org/wiki/Naive\\_Bayes\\_classifier](http://en.wikipedia.org/wiki/Naive_Bayes_classifier)
- [http://en.wikipedia.org/wiki/Rule\\_of\\_succession](http://en.wikipedia.org/wiki/Rule_of_succession)
- [http://eu.wikipedia.org/wiki/Banakuntza\\_normal](http://eu.wikipedia.org/wiki/Banakuntza_normal)