



# Tema 6. Distribución Multinomial

Percepción (PER)

Curso 2017/2018

Departamento de Sistemas Informáticos y Computación

- 1 Introducción y motivación ⊳ 3
- 2 Definición de la distribución multinomial > 5
- 3 Clasificador multinomial ▷ 9
- 4 Entrenamiento por máxima verosimilitud (MV) ▷ 12
- 5 Suavizado ▷ 14





- 1 Introducción y motivación ▷ 3
  - 2 Definición de la distribución multinomial > 5
  - 3 Clasificador multinomial ▷ 9
  - 4 Entrenamiento por máxima verosimilitud (MV) ▷ 12
  - 5 Suavizado ▷ 14

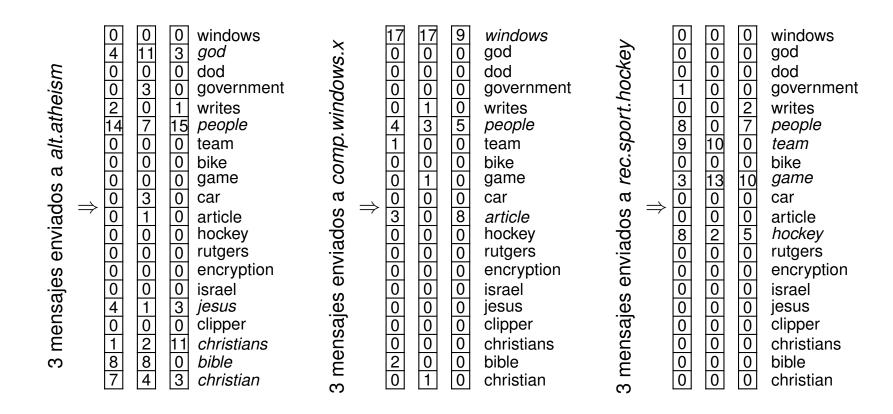




#### Distribución multinomial: motivación

Algunas tareas de RF representan objetos como vectores de cuentas

**Ejemplo:** Texto representado como *bag-of-words* 



**Idea:** usar la *distribución multinomial* para modelizar la condicional  $p({m x}|c)$ 





- 1 Introducción y motivación ⊳ 3
- 2 Definición de la distribución multinomial ▷ 5
  - 3 Clasificador multinomial ▷ 9
  - 4 Entrenamiento por máxima verosimilitud (MV) ▷ 12
  - 5 Suavizado ▷ 14





### Definición: distribución multinomial

Sea una población  $\mathcal{Y} = \{y_1, \dots, y_n\}$  con  $y_i \in \{1, \dots, D\}$ ,

Sean las proporciones  $p_d$  de los tipos de elemento  $\{1,\ldots,D\}$  dadas por:

$$p = (p_1, \dots, p_D)^t \in [0, 1]^D$$
 con  $\sum_{d=1}^D p_d = 1$ 

Sea una secuencia de N elementos formada por extracción aleatoria con reemplazo desde  $\mathcal Y$ 

$$w_1^N = w_1 \, w_2 \, \cdots \, w_N$$

Número de secuencias distintas de longitud N:

$$VR_{D,N} = D^N$$





### Definición: distribución multinomial

Asumiendo independencia entre elementos:

$$p(w_1^N) = p_{w_1} p_{w_2} \cdots p_{w_N}$$

No depende del orden de los elementos, sino de su número de ocurrencias:

- $x_d$ : el número de ocurrencias del elemento d en  $w_1^N$
- $x = (x_1, \dots, x_D)^t$ : vector de ocurrencias (número de ocurrencias de cada elemento en  $w_1^N$ )

$$p(w_1^N) = p_1^{x_1} \cdots p_D^{x_D} = \prod_{d=1}^D p_d^{x_d}$$

El número de secuencias diferentes con el mismo vector de ocurrencias es un coeficiente multinomial:

$$\binom{N}{\boldsymbol{x}} = \binom{N}{x_1, \dots, x_D} = \frac{N!}{x_1! \cdots x_D!}$$





### Definición: distribución multinomial

**Distribución multinomial**: se define sobre el espacio de vectores de ocurrencias

La probabilidad de x es la suma de probabilidades de todas las secuencias con vector de ocurrencias x:

$$p(\boldsymbol{x}) = {N \choose \boldsymbol{x}} \prod_{d=1}^{D} p_d^{x_d}$$

p(x) es una f.d. multinomial:

- *D*-dimensional
- Longitud  $N = \sum_{d=1}^{D} x_d$
- Prototipo p

De ahora en adelante, usaremos  $x_+ = N$ .





- 1 Introducción y motivación ⊳ 3
- 2 Definición de la distribución multinomial > 5
- 3 Clasificador multinomial ▷ 9
  - 4 Entrenamiento por máxima verosimilitud (MV) ▷ 12
  - 5 Suavizado ▷ 14





#### Clasificador multinomial

**Clasificador multinomial**: clasificador de Bayes donde la f.d. condicional  $p(\boldsymbol{x}|c)$  es una multinomial

$$p(\boldsymbol{x} \mid c) \sim Mult_D(x_+, \boldsymbol{p}_c), \quad c = 1, \dots, C.$$

Por tanto:

$$c^*(\boldsymbol{x}) = \underset{c=1,...,C}{\operatorname{argmax}} \log P(c) + \log p(x \mid c)$$

$$= \underset{c=1,...,C}{\operatorname{argmax}} \log P(c) + \log \frac{x_{+}!}{x_{1}! \cdots x_{D}!} \prod_{d=1}^{D} p_{cd}^{x_{d}}$$

$$= \underset{c=1,...,C}{\operatorname{argmax}} \log P(c) + \log \frac{x_{+}!}{x_{1}! \cdots x_{D}!} + \sum_{d=1}^{D} x_{d} \log p_{cd}$$





### Clasificador multinomial

Eliminando el término independiente de c:

$$c^*(x) = \underset{c=1,...,C}{\operatorname{argmax}} \log P(c) + \sum_{d=1}^{D} x_d \log p_{cd}$$

Expresando el sumatorio en forma de producto escalar:

$$c^*(\boldsymbol{x}) = \underset{c=1,...,C}{\operatorname{argmax}} (\log \boldsymbol{p}_c)^t \boldsymbol{x} + \log P(c)$$

En forma de clasificador lineal:

$$c^*(\boldsymbol{x}) = \underset{c=1,...,C}{\operatorname{argmax}} \ g_c(\boldsymbol{x}) = \underset{c=1,...,C}{\operatorname{argmax}} \ \boldsymbol{w}_c^t \boldsymbol{x} + w_{c0}$$

Con:

$$\boldsymbol{w}_c = \log \boldsymbol{p}_c$$
  $w_{c0} = \log P(c)$ 





- 1 Introducción y motivación ⊳ 3
- 2 Definición de la distribución multinomial > 5
- 3 Clasificador multinomial ▷ 9
- 4 Entrenamiento por máxima verosimilitud (MV) ▷ 12
  - 5 Suavizado ▷ 14





### Entrenamiento por máxima verosimilitud

Sean N muestras de entrenamiento aleatoriamente extraídas de C distribuciones multinomiales independientes:

$$\{(\boldsymbol{x}_n,c_n)\}_{n=1}^N$$
 i.i.d.  $p(\boldsymbol{x},c)=P(c)\,p(\boldsymbol{x}|c), \quad p(\boldsymbol{x}|c)\sim Mult_D(x_+,\boldsymbol{p}_c)$ 

Conjunto de parámetros a estimar  $\Theta$ :

- Probabilidades a priori:  $P(1) \dots, P(C)$
- lacktriangle Prototipos de las multinomiales para cada clase c:  $m{p}_c$ ,  $c=1,\ldots,C$

Por criterio de máxima verosimilitud (MV), se estima  $\Theta$  como:

$$\hat{P}(c) = \frac{N_c}{N} \qquad \hat{\boldsymbol{p}}_c = \frac{1}{\sum_{\substack{n: c_n = c \\ d = 1}} \sum_{\substack{n: c_n = c \\ d = 1}} \boldsymbol{x}_{nd}} \sum_{\substack{n: c_n = c \\ d = 1}} \boldsymbol{x}_n \qquad c = 1, \dots, C$$





- 1 Introducción y motivación ⊳ 3
- 2 Definición de la distribución multinomial > 5
- 3 Clasificador multinomial ▷ 9
- 4 Entrenamiento por máxima verosimilitud (MV) ▷ 12
- 5 Suavizado ▷ 14





#### Suavizado de la distribución multinomial

**Laplace**: suma una constante  $\epsilon > 0$  a cada parámetro y renormaliza

### Descuento Absoluto (DA):

- 1. Descuenta una constante  $\epsilon>0$  (pequeña) a cada parámetro mayor que cero
- 2. Distribuir la probabilidad descontada según una distribución generalizada:
  - Entre todos los parámetros nulos (backing-off)
  - Entre todos los parámetros (interpolación)

