



UNIVERSITAT
POLITÈCNICA
DE VALÈNCIA

Escola Tècnica Superior d'Enginyeria Informàtica



Tema 5. Distribución de Bernoulli

Percepción (PER)

Curso 2017/2018

Departamento de Sistemas Informáticos y Computación

Índice

- 1 Introducción y motivación ▷ 3
- 2 Definición de la distribución de Bernoulli ▷ 7
- 3 Clasificador Bernoulli ▷ 10
- 4 Entrenamiento por máxima verosimilitud (MV) ▷ 14
- 5 Suavizado ▷ 16

Índice

- 1 *Introducción y motivación* ▷ 3
- 2 Definición de la distribución de Bernoulli ▷ 7
- 3 Clasificador Bernoulli ▷ 10
- 4 Entrenamiento por máxima verosimilitud (MV) ▷ 14
- 5 Suavizado ▷ 16

Introducción

Clasificador de Bayes:

$$\begin{aligned} c^*(x) &= \operatorname{argmax}_{c=1,\dots,C} P(c \mid x) = \operatorname{argmax}_{c=1,\dots,C} \frac{P(c) p(x \mid c)}{p(x)} \\ &= \operatorname{argmax}_{c=1,\dots,C} P(c) p(x \mid c) = \operatorname{argmax}_{c=1,\dots,C} \log P(c) + \log p(x \mid c) \end{aligned}$$

- $P(c)$: probabilidad *a priori*
- $p(x|c)$: función de densidad (f.d.) de probabilidad condicional para clase c

En la práctica, se usan **estimaciones** de $P(c)$ y $p(x|c)$:

$$c^*(x) \approx \operatorname{argmax}_{c=1,\dots,C} \log \hat{P}(c) + \log \hat{p}(x \mid c)$$

Introducción

$\hat{P}(c)$ y $\hat{p}(x | c)$ se estiman a partir de N muestras etiquetadas, $(x_1, c_1), \dots, (x_N, c_N)$, extraídas aleatoriamente de $p(x, c)$

Estimación de la probabilidad *a priori*:

$$\hat{P}(c) = \frac{N_c}{N} \qquad N_c = \sum_{n : c_n = c} 1$$

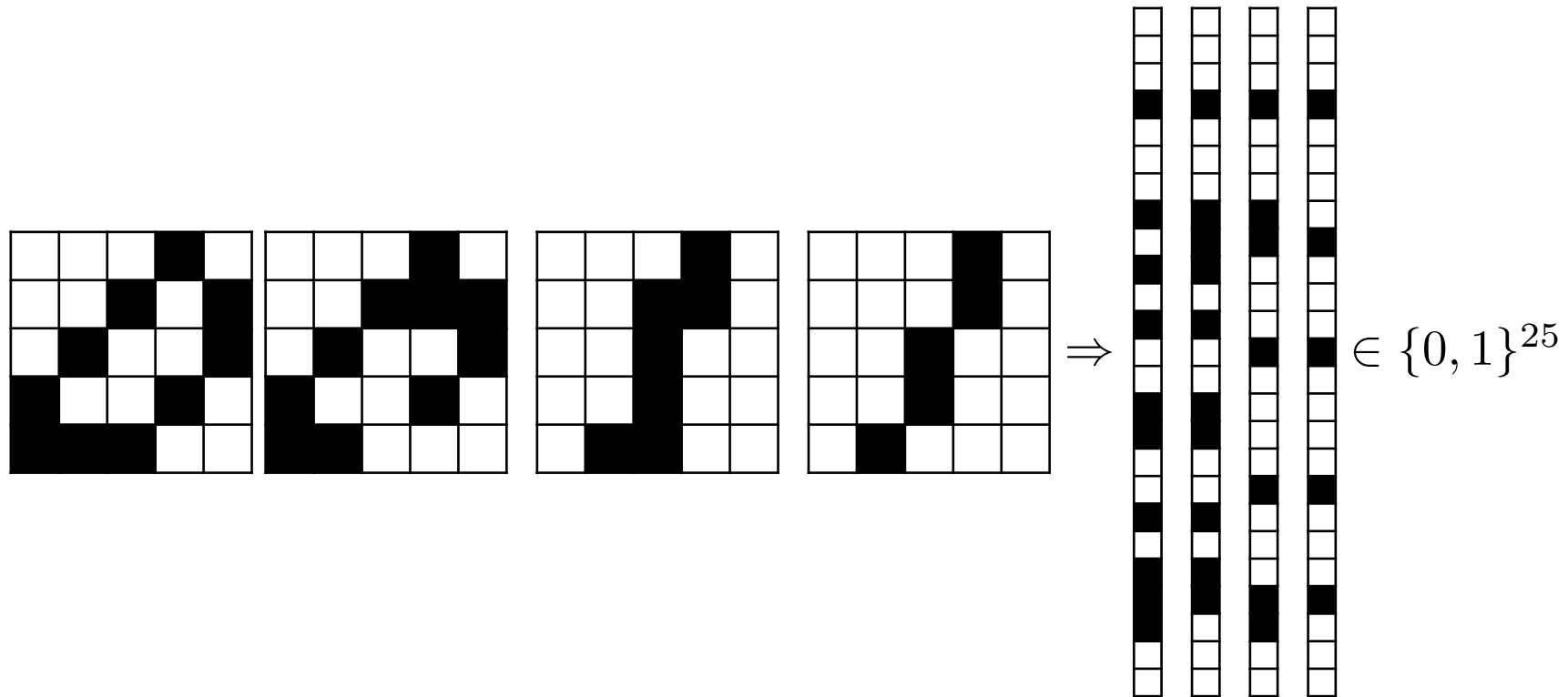
Estimación de la condicional $\hat{p}(x|c)$: a partir de las muestras x_n con $c_n = c$

Forma habitual: se asume un ***tipo de distribución*** sobre los datos de la clase y ***se estiman sus parámetros***

Motivación

Algunas tareas de RF conllevan objetos representados como un *vector de bits*.

Ejemplo: imágenes binarias de $5 \times 5 \rightarrow$ vectores de bits de 25 dimensiones



Idea: distribución de Bernoulli para modelar la condicional $p(x|c)$

Índice

- 1 Introducción y motivación ▷ 3
- 2 *Definición de la distribución de Bernoulli* ▷ 7
- 3 Clasificador Bernoulli ▷ 10
- 4 Entrenamiento por máxima verosimilitud (MV) ▷ 14
- 5 Suavizado ▷ 16

Definición: Bernoulli unidimensional

Sea $p \in [0, 1]$ y $q = 1 - p$.

Sea x una variable aleatoria que sigue una distribución de Bernoulli de parámetro p ($x \sim Be(p)$)

La f.d. de x es:

$$p(x) = \begin{cases} p & \text{si } x = 1 \\ q & \text{si } x = 0 \end{cases} = p x + q (1 - x) = p^x q^{1-x}$$

Nota: $0^0 = 1$ y $0 \log 0 = 0$

Definición: Bernoulli multidimensional

Sean $x_1 \sim Be(p_1), \dots, x_D \sim Be(p_D)$ independientes

En ese caso, $\mathbf{x} = (x_1, \dots, x_D)^t$ sigue una Bernoulli D -dimensional de parámetro $\mathbf{p} = (p_1, \dots, p_D)^t$

La f.d. de \mathbf{x} es:

$$p(\mathbf{x}) = \prod_{d=1}^D p(x_d) = \prod_{d=1}^D p_d x_d + q_d (1 - x_d) = \prod_{d=1}^D p_d^{x_d} q_d^{(1-x_d)}$$

Índice

- 1 Introducción y motivación ▷ 3
- 2 Definición de la distribución de Bernoulli ▷ 7
- 3 *Clasificador Bernoulli* ▷ 10
- 4 Entrenamiento por máxima verosimilitud (MV) ▷ 14
- 5 Suavizado ▷ 16

Clasificador Bernoulli

Clasificador Bernoulli: clasificador de Bayes en el caso particular en que la f.d. condicional $p(\mathbf{x}|c)$ es una Bernoulli:

$$p(\mathbf{x} | c) \sim Be_D(\mathbf{p}_c), \quad c = 1, \dots, C.$$

Por tanto:

$$\begin{aligned} c^*(\mathbf{x}) &= \operatorname{argmax}_{c=1, \dots, C} \log P(c) + \log p(\mathbf{x} | c) \\ &= \operatorname{argmax}_{c=1, \dots, C} \log P(c) + \log \prod_{d=1}^D p_{cd}^{x_d} (1 - p_{cd})^{(1-x_d)} \\ &= \operatorname{argmax}_{c=1, \dots, C} \log P(c) + \sum_{d=1}^D x_d \log p_{cd} + (1 - x_d) \log(1 - p_{cd}) \end{aligned}$$

Clasificador Bernoulli

Agrupando términos dependientes e independientes de x_d :

$$c^*(\mathbf{x}) = \operatorname{argmax}_{c=1,\dots,C} \left(\sum_{d=1}^D x_d (\log p_{cd} - \log(1 - p_{cd})) \right) + \left(\log P(c) + \sum_{d=1}^D \log(1 - p_{cd}) \right)$$

Reescribimos la expresión anterior como:

$$c^*(\mathbf{x}) = \operatorname{argmax}_{c=1,\dots,C} \sum_{d=1}^D w_{cd} x_d + w_{c0}$$

donde

$$w_{cd} = \log p_{cd} - \log(1 - p_{cd}) \quad w_{c0} = \log P(c) + \sum_{d=1}^D \log(1 - p_{cd})$$

Clasificador Bernoulli

Por tanto, es un *clasificador lineal* sobre x :

$$c^*(x) = \operatorname{argmax}_{c=1,\dots,C} g_c(x) = \operatorname{argmax}_{c=1,\dots,C} \sum_{d=1}^D w_{cd} x_d + w_{c0}$$

Reescribiendo la expresión anterior como un producto escalar de dos vectores:

$$c^*(x) = \operatorname{argmax}_{c=1,\dots,C} w_c^t x + w_{c0}$$

donde

$$w_c = \log p_c - \log(1 - p_c)$$

Índice

- 1 Introducción y motivación ▷ 3
- 2 Definición de la distribución de Bernoulli ▷ 7
- 3 Clasificador Bernoulli ▷ 10
- 4 *Entrenamiento por máxima verosimilitud (MV)* ▷ 14
- 5 Suavizado ▷ 16

Entrenamiento por máxima verosimilitud

Sean un conjunto de entrenamiento de N muestras independientes e idénticamente distribuidas (i.i.d.) extraídas aleatoriamente de C distribuciones Bernoulli:

$$\{(\mathbf{x}_n, c_n)\}_{n=1}^N \quad \text{i.i.d.} \quad p(\mathbf{x}, c) = P(c) p(\mathbf{x}|c), \quad p(\mathbf{x}|c) \sim Be_D(\mathbf{p}_c)$$

Conjunto de parámetros a estimar Θ :

- Probabilidades *a priori*: $P(1) \dots, P(C)$
- Parámetros de las Bernoulli para cada clase c : \mathbf{p}_c , $c = 1, \dots, C$

Por ***criterio de máxima verosimilitud*** (MV), se estima Θ como:

$$\hat{P}(c) = \frac{N_c}{N} \quad c = 1, \dots, C$$

$$\hat{\mathbf{p}}_c = \frac{1}{N_c} \sum_{n: c_n=c} \mathbf{x}_n \quad c = 1, \dots, C$$

Índice

- 1 Introducción y motivación ▷ 3
- 2 Definición de la distribución de Bernoulli ▷ 7
- 3 Clasificador Bernoulli ▷ 10
- 4 Entrenamiento por máxima verosimilitud (MV) ▷ 14
- 5 *Suavizado* ▷ 16

Suavizado de la distribución Bernoulli

Problema: muchos criterios de entrenamiento (incluido MV) pueden generar clasificadores sobreentrenados

Soluciones:

- Cambiar el criterio de aprendizaje
- ***Suavizar*** los parámetros estimados

Opciones de suavizado en Bernoulli:

- Truncamiento simple
- Muestra ficticia

Suavizado de la distribución Bernoulli

Truncamiento simple

Dado ϵ , $0 \leq \epsilon \leq 0.5$, redefinir \hat{p}_{cd} como:

$$\tilde{p}_{cd} = \begin{cases} \epsilon & \text{si } \hat{p}_{cd} < \epsilon \\ 1 - \epsilon & \text{si } \hat{p}_{cd} > 1 - \epsilon \\ \hat{p}_{cd} & \text{en otro caso} \end{cases}$$

Muestra ficticia

Añadir al conjunto de aprendizaje $(\mathbf{0}, c)$ y $(\mathbf{1}, c)$, $c = 1, \dots, C$.

Equivale a redefinir la estimación de \hat{p}_c como:

$$\tilde{p}_c = \frac{1}{N_c + 2} \left(1 + \sum_{n: c_n = c} x_n \right)$$