

TRABAJO TEMA 3 APR

1.- Demostrar que en cualquier problema de clasificación en C clases, la estimación de máxima verosimilitud de la probabilidad a priori de cada clase c , $1 \leq c \leq C$, es $\hat{p}_c = n_c/N$ donde $N = n_1 + \dots + n_c$ es el número total de datos observados y n_c es el número de datos de la clase c (ver el último ejemplo de aplicación de la técnica de los multiplicadores de Lagrange, transparencias 3.17 y 3.18).

Modelo:

$$P(C = c) = \hat{p}_c$$

$$\sum_{c=1}^C \hat{p}_c = 1$$

$$\boldsymbol{\theta} = (\hat{p}_1, \dots, \hat{p}_C)^t$$

Verosimilitud y logaritmo de la verosimilitud:

$$P(S|\boldsymbol{\theta}) = \prod_{i=1}^{n_1} \hat{p}_1 \dots \prod_{j=1}^{n_C} \hat{p}_C = \hat{p}_1^{n_1} \dots \hat{p}_C^{n_C}$$

$$q_S(\boldsymbol{\theta}) = L_S(\boldsymbol{\theta}) = \log P(S|\boldsymbol{\theta}) = n_1 * \log \hat{p}_1 + \dots + n_C * \log \hat{p}_C$$

Estimación de máxima verosimilitud:

$$\boldsymbol{\theta}^* = \underset{\substack{\hat{p}_1, \dots, \hat{p}_C \\ \hat{p}_1 + \dots + \hat{p}_C = 1}}{\operatorname{argmax}} (n_1 * \log \hat{p}_1 + \dots + n_C * \log \hat{p}_C)$$

Lagrangiana:

$$\Lambda(\hat{p}_1, \dots, \hat{p}_C, \beta) = n_1 * \log \hat{p}_1 + \dots + n_C * \log \hat{p}_C + \beta(1 - \hat{p}_1 - \dots - \hat{p}_C)$$

Soluciones óptimas en función del multiplicador de Lagrange:

$$\begin{aligned} \frac{\partial \Lambda}{\partial \hat{p}_1} &= \frac{n_1}{\hat{p}_1} - \beta = 0 & \hat{p}_1^*(\beta) &= \frac{n_1}{\beta} \\ &\dots & \Rightarrow & \dots \\ \frac{\partial \Lambda}{\partial \hat{p}_C} &= \frac{n_C}{\hat{p}_C} - \beta = 0 & \hat{p}_C^*(\beta) &= \frac{n_C}{\beta} \end{aligned}$$

Función dual de Lagrange

$$\begin{aligned} \Lambda_D(\beta) &= n_1 \log \frac{n_1}{\beta} + \dots + n_C \log \frac{n_C}{\beta} + \beta \left(1 - \frac{n_1}{\beta} - \dots - \frac{n_C}{\beta} \right) \\ &= \beta - N * \log \beta - N + \sum_{c=1}^C n_c \log n_c \end{aligned}$$

Valor óptimo del multiplicador de Lagrange:

$$\frac{d\Lambda_D}{d\beta} = 1 - \frac{N}{\beta} = 0 \Rightarrow \beta^* = N$$

Solución final:

$$\hat{p}_1^* = \hat{p}_1^*(\beta) = \frac{n_1}{N} \dots \hat{p}_c^* = \hat{p}_c^*(\beta) = \frac{n_c}{N}$$

2.- Aplicar la técnica de descenso por gradiente a la búsqueda del mínimo de la función:
 $q(\theta) = (\theta_1 - 1)^2 + (\theta_2 - 2)^2 + \theta_1\theta_2$ **teniendo en cuenta que $\rho_k = 1/2k$ y $\theta(1) = (-1, +1)$ y hacer una traza de las 3 primeras iteraciones.**

Algoritmo general de descenso por gradiente:

$$\theta(1) = \text{arbitrario}$$

$$\theta(k+1) = \theta(k) - \rho_k \nabla q(\theta)|_{\theta=\theta(k)}$$

Derivadas parciales:

$$\frac{dq}{d\theta_1} = 2\theta_1 + \theta_2 - 2$$

$$\frac{dq}{d\theta_2} = 2\theta_2 + \theta_1 - 4$$

Sustituimos en el algoritmo:

$$\theta(1) = \begin{pmatrix} -1 \\ +1 \end{pmatrix}$$

$$\theta(k+1) = \theta(k) - \frac{1}{2k} \begin{pmatrix} 2\theta_1 + \theta_2 - 2 \\ 2\theta_2 + \theta_1 - 4 \end{pmatrix}$$

Primera iteración:

$$\theta(2) = \begin{pmatrix} -1 \\ +1 \end{pmatrix} - \frac{1}{2 * 1} \begin{pmatrix} 2 * (-1) + 1 - 2 \\ 2 * 1 + (-1) - 4 \end{pmatrix} = \begin{pmatrix} 1/2 \\ 5/2 \end{pmatrix}$$

Segunda iteración:

$$\theta(3) = \begin{pmatrix} 1/2 \\ 5/2 \end{pmatrix} - \frac{1}{2 * 2} \begin{pmatrix} 2 * (1/2) + 5/2 - 2 \\ 2 * 5/2 + (1/2) - 4 \end{pmatrix} = \begin{pmatrix} 1/8 \\ 17/8 \end{pmatrix}$$

Tercera iteración:

$$\theta(4) = \begin{pmatrix} 1/8 \\ 17/8 \end{pmatrix} - \frac{1}{2 * 3} \begin{pmatrix} 2 * (1/8) + 17/8 - 2 \\ 2 * 17/8 + (1/8) - 4 \end{pmatrix} = \begin{pmatrix} 1/16 \\ 33/16 \end{pmatrix}$$

3.- Existe una variante de la función de Widrow-Hoff que incluye un término de regularización con el objetivo de que los pesos no se hagan demasiado grandes:

$$q_s(\theta) = \frac{1}{2} \sum_{n=1}^N (\theta^t x_n - y_n)^2 + \frac{\theta^t \theta}{2}$$

Aplicando la técnica de descenso por gradiente, obtener la correspondiente variante del algoritmo de Widrow-Hoff y la correspondiente versión muestra a muestra.

Algoritmo general de descenso por gradiente:

$$\theta(1) = \text{arbitrario}$$

$$\theta(k+1) = \theta(k) - p_k \nabla q(\theta)|_{\theta=\theta(k)}$$

Calculamos el gradiente:

$$\nabla q_s(\theta) = \nabla \left(\frac{1}{2} \sum_{n=1}^N (\theta^t x_n - y_n)^2 + \frac{\theta^t \theta}{2} \right) = \sum_{n=1}^N (\theta^t x_n - y_n) x_n + \theta$$

Versión del algoritmo:

$$\theta(1) = \text{arbitrario}$$

$$\theta(k+1) = \theta(k) - p_k \left(\sum_{n=1}^N (\theta(k)^t x_n - y_n) x_n + \theta(k) \right)$$

$$= \theta(k) + p_k \sum_{n=1}^N (y_n - \theta(k)^t x_n) x_n - p_k \theta(k)$$

Versión del algoritmo muestra a muestra:

$$\theta(1) = \text{arbitrario}$$

$$\theta(k+1) = \theta(k) + p_k (y(k) - \theta(k)^t x(k)) x(k) - p_k \theta(k)$$