

Modelado Probabilístico Generativo Gaussiano de 2 Clases con Estimación de Parámetros por Máxima Verosimilitud (Marzo 2012)

Iván López Espejo

Deducción de la estimación por máxima verosimilitud de los parámetros correspondientes a las distribuciones de probabilidad de un modelo probabilístico generativo de clasificación gaussiano de 2 clases e iguales matrices de covarianza con el fin de calcular posteriormente el hiperplano de decisión para la clasificación de muestras entrantes al sistema.

I. INTRODUCCIÓN

EN ESTE trabajo se deduce la estimación por máxima verosimilitud de los parámetros correspondientes a las distribuciones de probabilidad de un modelo probabilístico generativo de clasificación gaussiano de $K = 2$ clases e iguales matrices de covarianza. Una vez estimados estos, es posible calcular el hiperplano de decisión en base a la función logística sigmoideal. Calculado este, es posible clasificar una nueva muestra de entrada al sistema en una de las dos clases posibles sin más que estudiar la probabilidad de pertenencia a cada una de las clases a partir de la anterior función logística.

Para completar el trabajo, a partir de los resultados obtenidos, se realiza una implementación en MatLab con la que llevar a cabo una experimentación práctica.

II. DESARROLLO

Partimos de poseer un conjunto de N muestras agrupables en dos clases según sendas distribuciones gaussianas. Para cada muestra, además, se tiene un valor de etiqueta, t_n , que indica a priori la pertenencia de dicha muestra a una clase, es decir,

$$\{\mathbf{x}_n, t_n\}, \quad n = 1, 2, \dots, N,$$

$$t_n = \begin{cases} 0 & \Leftrightarrow \mathbf{x}_n \in C_0 \\ 1 & \Leftrightarrow \mathbf{x}_n \in C_1' \end{cases}$$

donde, como se ve, t_n toma el valor 0 si la muestra \mathbf{x}_n pertenece a la clase 0 (C_0) o el valor 1 si pertenece a la clase 1 (C_1). La probabilidad a priori de la primera clase es $P(C_1) = \Pi$, de tal forma que, como sólo se poseen dos clases, la probabilidad a priori de la segunda es su complementario, es decir, $P(C_0) = (1 - \Pi)$. Dado que hemos dicho que la densidad de probabilidad condicional de las muestras dada la clase se modela como una distribución gaussiana, las probabilidades conjuntas de las muestras y las clases se expresan, a partir de la Regla de Bayes, como

$$P(\mathbf{x}_n, C_0) = P(\mathbf{x}_n|C_0)P(C_0) = \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_0, \boldsymbol{\Sigma})(1 - \Pi),$$

$$P(\mathbf{x}_n, C_1) = P(\mathbf{x}_n|C_1)P(C_1) = \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})\Pi.$$

Sea $\mathbf{t} = (t_1, t_2, \dots, t_N)^T$ el vector de etiquetas de pertenencia a una clase para las N muestras, la función de verosimilitud se expresa como

$$p(\mathbf{t}|\Pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})\Pi]^{t_n} [\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_0, \boldsymbol{\Sigma})(1 - \Pi)]^{(1-t_n)},$$

donde, como se puede observar, sólo las probabilidades conjuntas de las muestras que pertenecen a su clase pesan en el anterior baremo. Los parámetros que deseamos estimar a continuación por máxima verosimilitud son las probabilidades a priori de cada clase, las medias de cada una de ellas y su matriz de covarianza. Para simplificar el cálculo, optimizamos sobre el logaritmo de la función de verosimilitud, pues su monotonía no cambia tras este tipo de composición:

$$\log p(\mathbf{t}|\Pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}) = \sum_{n=1}^N t_n \log(\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})\Pi) + \sum_{n=1}^N (1 - t_n) \log(\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_0, \boldsymbol{\Sigma})(1 - \Pi)),$$

donde la distribución normal multivariada de dos clases se expresa como

$$\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{2\pi|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})\right\},$$

y, en términos logarítmicos como

$$\log \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\log(2\pi|\boldsymbol{\Sigma}|^{1/2}) - \frac{1}{2}(\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}).$$

Sustituyendo la anterior expansión en el desarrollo logarítmico de la función de verosimilitud, esta resulta

$$\begin{aligned} \log p(\mathbf{t}|\Pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}) &= \sum_{n=1}^N t_n \log \Pi - \sum_{n=1}^N t_n \log(2\pi|\boldsymbol{\Sigma}|^{1/2}) - \\ &- \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_1) + \sum_{n=1}^N (1 - t_n) \log(1 - \Pi) - \\ &- \sum_{n=1}^N (1 - t_n) \log(2\pi|\boldsymbol{\Sigma}|^{1/2}) - \frac{1}{2} \sum_{n=1}^N (1 - t_n) (\mathbf{x}_n - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_0). \end{aligned}$$

Derivando con respecto a Π la anterior expresión e igualando a 0 optimizamos para calcular por máxima verosimilitud la probabilidad a priori de la primera clase y, por ende, la de la segunda:

$$\frac{\partial \log p(\mathbf{t}|\Pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_0, \boldsymbol{\Sigma})}{\partial \Pi} = 0 \Rightarrow \frac{1}{\Pi} \sum_{n=1}^N t_n = \frac{1}{1-\Pi} \sum_{n=1}^N (1-t_n).$$

La anterior derivada es trivial, pues sólo depende de dos términos logarítmicos en la sumatoria de la función log-verosimilitud. Vamos a continuación a renombrar las sumatorias en ambos miembros de la anterior derivada. La suma de las etiquetas t_n nos proporciona la cantidad de muestras pertenecientes a la primera clase, es decir:

$$\sum_{n=1}^N t_n = N_1.$$

De otro lado, la suma del complementario, $(1-t_n)$, nos proporciona la cantidad de muestras pertenecientes a la clase 0, es decir:

$$\sum_{n=1}^N (1-t_n) = N_0.$$

La suma de ambos términos nos proporciona la cantidad total de muestras, N , de la forma:

$$N = N_0 + N_1.$$

Sustituyendo lo anterior en el resultado de la derivada parcial de la función log-verosimilitud con respecto a Π , llegamos a que la estimación ML de Π es

$$\begin{aligned} \frac{N_1}{\Pi} = \frac{N_0}{1-\Pi} &\Rightarrow \frac{1}{\Pi} - 1 = \frac{N_0}{N_1} \Rightarrow \Pi^{-1} = \frac{N_0}{N_1} + 1 \Rightarrow \\ &\Rightarrow \Pi = \frac{N_1}{N_0 + N_1} = \frac{N_1}{N}. \end{aligned}$$

Como era de esperar, la probabilidad a priori de C_1 es el número de muestras pertenecientes a dicha clase sobre el total, siendo la probabilidad a priori de la clase C_0 su complementario, es decir,

$$P(C_0) = (1-\Pi) = 1 - \frac{N_1}{N} = \frac{N - N_1}{N} = \frac{N_0}{N},$$

o, en otras palabras, el número de muestras pertenecientes a C_0 sobre el total.

A continuación, llevamos a cabo un procedimiento análogo al anterior para obtener la estimación de la media de la clase 1. Comenzamos derivando la función log-verosimilitud con respecto a $\boldsymbol{\mu}_1$:

$$\begin{aligned} \frac{\partial \log p(\mathbf{t}|\Pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_0, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}_1} = 0 &\Rightarrow \\ \Rightarrow -\frac{\partial}{\partial \boldsymbol{\mu}_1} \left(\frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \right) &= 0. \end{aligned}$$

Como vemos, sólo un término de la sumatoria (el del exponente de la gaussiana correspondiente) depende de la media de la clase 1. Para mayor facilidad de cálculo, expandimos su argumento:

$$\begin{aligned} (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) &= \\ = \mathbf{x}_n^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_n - \mathbf{x}_n^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_n + \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1. \end{aligned}$$

A continuación tenemos en cuenta las siguientes tres identidades de derivación matricial para resolver el cálculo:

$$\frac{\partial (\mathbf{A}\mathbf{x})}{\partial \mathbf{x}} = \mathbf{A}^T,$$

$$\frac{\partial (\mathbf{x}^T \mathbf{A})}{\partial \mathbf{x}} = \mathbf{A},$$

$$\frac{\partial (\mathbf{x}^T \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = \mathbf{A} \mathbf{x} + \mathbf{A}^T \mathbf{x}.$$

Aplicándolas, llegamos a que la anterior derivada finalmente puede ser expresada como

$$\sum_{n=1}^N t_n [-(\mathbf{x}_n^T \boldsymbol{\Sigma}^{-1})^T - \boldsymbol{\Sigma}^{-1} \mathbf{x}_n + \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + (\boldsymbol{\Sigma}^{-1})^T \boldsymbol{\mu}_1] = 0.$$

Puesto que la matriz de covarianza es una matriz simétrica ($\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^T$), siendo su inversa también simétrica, $\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + (\boldsymbol{\Sigma}^{-1})^T \boldsymbol{\mu}_1 = 2\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1$. Además, por la regla de trasposición del producto matricial, tenemos que $(\mathbf{x}_n^T \boldsymbol{\Sigma}^{-1})^T = (\boldsymbol{\Sigma}^{-1})^T \mathbf{x}_n$. Aplicando además de nuevo la propiedad de simetría, el primer par de monomios del corchete puede simplificarse a su vez como $-(\mathbf{x}_n^T \boldsymbol{\Sigma}^{-1})^T - \boldsymbol{\Sigma}^{-1} \mathbf{x}_n = -2\boldsymbol{\Sigma}^{-1} \mathbf{x}_n$. Teniendo esto en cuenta, la estimación ML de $\boldsymbol{\mu}_1$ resulta finalmente

$$\begin{aligned} \sum_{n=1}^N t_n 2\boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \mathbf{x}_n) = 0 &\Rightarrow \boldsymbol{\mu}_1 \sum_{n=1}^N t_n = \sum_{n=1}^N t_n \mathbf{x}_n \Rightarrow \\ \Rightarrow \boldsymbol{\mu}_1 &= \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n. \end{aligned}$$

Como era de esperar, resulta la media muestral del conjunto de muestras de la clase 1.

Procedemos a continuación de forma análoga con el fin de estimar la media de la clase 0:

$$\begin{aligned} \frac{\partial \log p(\mathbf{t}|\Pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_0, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}_0} = 0 &\Rightarrow \\ \Rightarrow -\frac{\partial}{\partial \boldsymbol{\mu}_0} \left(\frac{1}{2} \sum_{n=1}^N (1-t_n) (\mathbf{x}_n - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_0) \right) &= 0. \end{aligned}$$

Aplicando las mismas identidades de derivación matricial que en el caso anterior, llegamos a

$$\sum_{n=1}^N (1-t_n) [-(\mathbf{x}_n^T \boldsymbol{\Sigma}^{-1})^T - \boldsymbol{\Sigma}^{-1} \mathbf{x}_n + \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 + (\boldsymbol{\Sigma}^{-1})^T \boldsymbol{\mu}_0] = 0.$$

Simplificando a partir de las mismas consideraciones, llegamos a que la estimación ML de la media de la clase 0 es

$$\begin{aligned} \sum_{n=1}^N (1-t_n) 2\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_0 - \mathbf{x}_n) = 0 &\Rightarrow \boldsymbol{\mu}_0 \sum_{n=1}^N (1-t_n) = \sum_{n=1}^N (1-t_n) \mathbf{x}_n \Rightarrow \\ &\Rightarrow \boldsymbol{\mu}_0 = \frac{1}{N_0} \sum_{n=1}^N (1-t_n) \mathbf{x}_n, \end{aligned}$$

es decir, de nuevo, la media muestral del conjunto de muestras pertenecientes a C_0 .

Finalmente, llevamos a cabo la estimación de la matriz de covarianza de ambas clases. De nuevo, derivamos la función log-verosimilitud, en este caso, respecto de $\boldsymbol{\Sigma}$. Teniendo en cuenta únicamente los términos dependientes de la matriz de covarianza, el problema se reduce a resolver la siguiente ecuación:

$$\begin{aligned} \frac{\partial \log p(t|\Pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_0, \boldsymbol{\Sigma})}{\partial \boldsymbol{\Sigma}} = 0 &\Rightarrow \\ \Rightarrow \frac{\partial}{\partial \boldsymbol{\Sigma}} \left(-\frac{1}{2} \sum_{n=1}^N t_n \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) - \right. \\ &- \frac{1}{2} \sum_{n=1}^N (1-t_n) \log |\boldsymbol{\Sigma}| - \\ &\left. - \frac{1}{2} \sum_{n=1}^N (1-t_n) (\mathbf{x}_n - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_0) \right) = 0. \end{aligned}$$

Para resolver la derivada del logaritmo del determinante de $\boldsymbol{\Sigma}$ hacemos uso de la siguiente identidad de derivación matricial:

$$\frac{\partial \log |\mathbf{A}|}{\partial \mathbf{A}} = (\mathbf{A}^{-1})^T.$$

Análogamente, aplicamos la siguiente identidad para resolver la derivada en los términos del argumento de la exponencial de las gaussianas:

$$\frac{\partial \mathbf{a}^T \mathbf{A}^{-1} \mathbf{b}}{\partial \mathbf{A}} = -(\mathbf{A}^{-1})^T \mathbf{a} \mathbf{b}^T (\mathbf{A}^{-1})^T.$$

Sustituyendo, la derivada resulta finalmente:

$$\begin{aligned} -\frac{1}{2} \sum_{n=1}^N t_n \boldsymbol{\Sigma}^{-1} + \frac{1}{2} \sum_{n=1}^N t_n \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} - \frac{1}{2} \sum_{n=1}^N (1-t_n) \boldsymbol{\Sigma}^{-1} + \\ + \frac{1}{2} \sum_{n=1}^N (1-t_n) \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_0) (\mathbf{x}_n - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} = 0. \end{aligned}$$

Para simplificar, multiplicamos por la derecha por la matriz de covarianza, resultando en

$$-\frac{N_1}{2} + \frac{\boldsymbol{\Sigma}^{-1}}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1) (\mathbf{x}_n - \boldsymbol{\mu}_1)^T - \frac{N_0}{2} +$$

$$+ \frac{\boldsymbol{\Sigma}^{-1}}{2} \sum_{n=1}^N (1-t_n) (\mathbf{x}_n - \boldsymbol{\mu}_0) (\mathbf{x}_n - \boldsymbol{\mu}_0)^T = 0.$$

Multiplicando ambos términos de la ecuación por 2 y sacando factor común la matriz de covarianza inversa,

$$\boldsymbol{\Sigma}^{-1} \left[\sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1) (\mathbf{x}_n - \boldsymbol{\mu}_1)^T + \sum_{n=1}^N (1-t_n) (\mathbf{x}_n - \boldsymbol{\mu}_0) (\mathbf{x}_n - \boldsymbol{\mu}_0)^T \right] = N_0 + N_1 = N.$$

Finalmente, la estimación ML de la matriz de covarianza resulta

$$\boldsymbol{\Sigma} = \frac{1}{N} \left[\sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1) (\mathbf{x}_n - \boldsymbol{\mu}_1)^T + \sum_{n=1}^N (1-t_n) (\mathbf{x}_n - \boldsymbol{\mu}_0) (\mathbf{x}_n - \boldsymbol{\mu}_0)^T \right],$$

que de nuevo se traduce, como era de esperar, en la covarianza muestral.

La probabilidad de la primera clase dada la muestra se puede expresar en términos de la función logística como $P(C_1|\mathbf{x}_n) = \sigma(\mathbf{w}^T \mathbf{x}_n + w_0)$, donde

$$\sigma(\mathbf{w}^T \mathbf{x}_n + w_0) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x}_n + w_0)}}.$$

Además, tenemos que $a = \mathbf{w}^T \mathbf{x}_n + w_0 = \log \frac{P(C_1|\mathbf{x}_n)}{P(C_0|\mathbf{x}_n)}$, de tal forma que en el límite, cuando las dos clases son equiprobables, $a = 0 = \mathbf{w}^T \mathbf{x}_n + w_0$, lo que constituye nuestro hiperplano de decisión. Resolviendo la anterior ecuación, llegamos a la forma explícita del hiperplano de decisión:

$$\begin{aligned} \mathbf{w}^T \mathbf{x} + w_0 = 0 &\Rightarrow (w_1 \ w_2) \begin{pmatrix} x \\ y \end{pmatrix} + w_0 = 0 \Rightarrow w_1 x + w_2 y + w_0 = 0 \Rightarrow \\ &\Rightarrow y(x) = -\frac{w_1 x + w_0}{w_2}, \end{aligned}$$

donde $\mathbf{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$ y $w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \log \frac{P(C_1)}{P(C_0)}$.

Finalmente se ha llevado a cabo una implementación en MatLab que realiza todo el proceso: generación de muestras de dos clases según distribuciones gaussianas, estimación de sus parámetros por máxima verosimilitud, cálculo del hiperplano de decisión y clasificación de una nueva muestra de entrada. La figura 1 muestra un ejemplo de resultado del funcionamiento de las rutinas programadas.

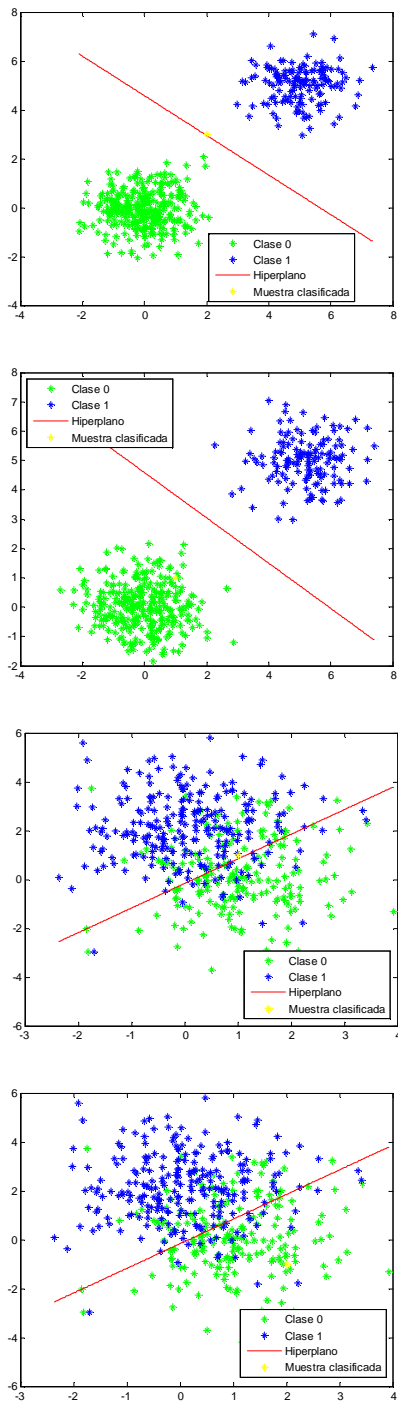


Fig. 1. De arriba a abajo: (a) dos clases con vectores de medias $(0,0)$ y $(5,5)$ y matriz de covarianzas diagonal con autovalores 0.8 y 0.6 donde se clasifica una muestra muy cercana al hiperplano de decisión, (b) igual caso que (a) pero donde se clasifica una muestra con claridad en la clase 0, (c) dos clases con vectores de medias $(1,0)$ y $(0,2)$ y matriz de covarianzas diagonal con autovalores 1 y 2 donde se clasifica una muestra muy cercana al hiperplano de decisión, (d) igual caso que (c) pero donde se clasifica una muestra más probable en clase 0.