

Clasificación estadística de patrones

Clasificador gaussiano

César Martínez

cmartinez_at_sinc.unl.edu.ar

Tópicos Selectos en Aprendizaje Maquinal
Doctorado en Ingeniería, FICH-UNL

12 de setiembre de 2016



- Supervisado:

- ↳ patrones de entrenamiento *etiquetados*.

- ↳ búsqueda de regiones de decisión.

- Aproximación paramétrica: fronteras de decisión definidas por las funciones de densidad de probabilidad de las clases. Método: clasificador gaussiano.

- Aproximación no paramétrica: fronteras de decisión definidas por los prototipos. Métodos: núcleos, vecinos, árboles de decisión, otros.

- No supervisado:

- ↳ patrones de entrenamiento de *clase incierta*.

- ↳ agrupamiento de patrones sin etiquetar.

- Métodos: con cantidad de clases conocida (K -medias) o desconocida, grafos, etc.

Se desconocen $P(\omega_i)$ y $P(\mathbf{x}|\omega_i)$, pero se disponen de N muestras de etiqueta conocida (*prototipos*):

$$(x_1, c_1), \dots, (x_N, c_N) \text{ i.i.d según } P(\mathbf{x})$$

- Estimación de $P(\omega_i)$: $\hat{P}(\omega_i) = \frac{N_i}{N}$, con N_i : número de muestras de ω_i .
- Estimación de $P(\mathbf{x}|\omega_i)$: $\hat{P}(\mathbf{x}|\omega_i)$
 - Métodos paramétricos: suponen una f.d.p. de cierta forma paramétrica conocida (p.ej., gaussiana), de parámetros desconocidos.
 - Métodos no paramétricos: sin suposición sobre $P(\mathbf{x}|\omega_i)$.

Supuesto de independencia funcional: densidades condicionales independientes entre sí $\Rightarrow c$ problemas de estimación independientes.

- Función unidimensional: variable aleatoria con *fdp*

$$N(\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x - \mu}{\sigma}\right)^2\right)$$

- Densidad normal multidimensional:

$$N_d(\boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

- Clasificador de Bayes cuyas densidades condicionales son gaussianas. Si $P(\mathbf{x}|\omega_i) \approx N_d(\boldsymbol{\mu}_i, \Sigma_i)$, entonces la discriminante asociada a ω_i es:

$$\begin{aligned}g_i(\mathbf{x}) &= \log P(\mathbf{x}|\omega_i) + \log P(\omega_i) \\&= -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^t \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \log P(\omega_i) \\&= -\frac{1}{2} \mathbf{x}^t \Sigma_i^{-1} \mathbf{x} \\&\quad + \boldsymbol{\mu}_i^t \Sigma_i^{-1} \mathbf{x} \\&\quad - \frac{1}{2} (\boldsymbol{\mu}_i^t \Sigma_i^{-1} \boldsymbol{\mu}_i + d \log 2\pi + \log |\Sigma_i|) + \log P(\omega_i)\end{aligned}$$

- Caso 1: $\Sigma_i = \sigma^2 I$.
 - Matrices de covarianza cuadradas e iguales para todas las clases.
 - Variables no correlacionadas.
 - Las muestras forman nubes hiperesféricas del mismo tamaño.

Como: $|\Sigma_i| = \sigma^{2d}$ y $\Sigma_i^{-1} = \frac{1}{\sigma^2} I$:

$$\begin{aligned}g_i(\mathbf{x}) &= -\frac{1}{2\sigma^d} \mathbf{x}^t \mathbf{x} + \frac{1}{\sigma^2} \boldsymbol{\mu}_i^t \mathbf{x} - \frac{1}{2} \left(\frac{1}{\sigma^2} \boldsymbol{\mu}_i^t \boldsymbol{\mu}_i + d \log 2\pi + d \log \sigma^2 \right) + \log P(\omega_i) \\ &\equiv \frac{1}{\sigma^2} \boldsymbol{\mu}_i^t \mathbf{x} - \frac{1}{2\sigma^2} \boldsymbol{\mu}_i^t \boldsymbol{\mu}_i + \log P(\omega_i) \\ &= \mathbf{w}_i^t \mathbf{x} + w_{i0},\end{aligned}$$

con

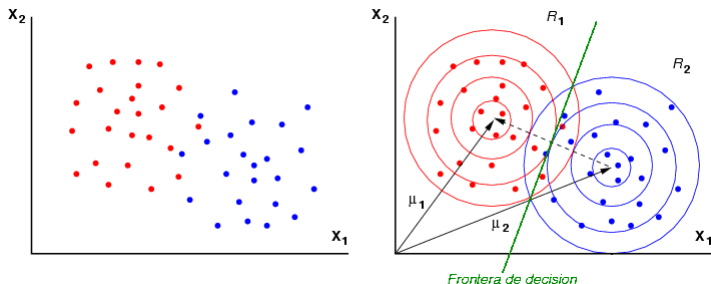
$$\begin{aligned}\mathbf{w}_i &= \frac{1}{\sigma^2} \boldsymbol{\mu}_i \\ w_{i0} &= -\frac{1}{2\sigma^2} \boldsymbol{\mu}_i^t \boldsymbol{\mu}_i + \log P(\omega_i)\end{aligned}$$

Clasificador gaussiano

- Caso 1: $\Sigma_i = \sigma^2 I$.
 - Clasificador lineal por mínima distancia Euclídea, con discriminantes

$$g_i(\mathbf{x}) = -\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2\sigma^2} + \log P(\omega_i)$$

- Fronteras de decisión: hiperplanos.



- Caso 2: $\Sigma_i = \Sigma$.
 - Matrices de covarianza arbitrarias e iguales para todas las clases.
 - Variables correlacionadas.
 - Las muestras forman nubes hiperelípticas del mismo tamaño y forma.

$$\begin{aligned}g_i(\mathbf{x}) &= \\ & -\frac{1}{2}\mathbf{x}^t\Sigma^{-1}\mathbf{x} + \boldsymbol{\mu}_i^t\Sigma^{-1}\mathbf{x} - \frac{1}{2}\left(\boldsymbol{\mu}_i^t\Sigma^{-1}\boldsymbol{\mu}_i + d\log 2\pi + \log|\Sigma|\right) + \log P(\omega_i) \\ & \equiv \boldsymbol{\mu}_i^t\Sigma^{-1}\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_i^t\Sigma^{-1}\boldsymbol{\mu}_i + \log P(\omega_i) \\ & = \mathbf{w}_i^t\mathbf{x} + w_{i0},\end{aligned}$$

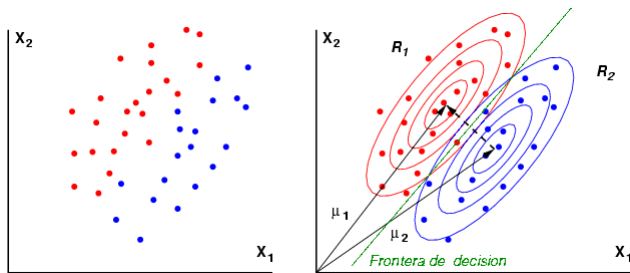
con

$$\begin{aligned}\mathbf{w}_i &= \Sigma^{-1}\boldsymbol{\mu}_i \\ w_{i0} &= -\frac{1}{2}\boldsymbol{\mu}_i^t\Sigma^{-1}\boldsymbol{\mu}_i + \log P(\omega_i)\end{aligned}$$

- Caso 2: $\Sigma_i = \Sigma$.
 - Clasificador lineal por mínima distancia de Mahalanobis:

$$g_i(\mathbf{x}) = -\frac{1}{2}\boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i + \boldsymbol{\mu}_i^T \Sigma^{-1} \mathbf{x} + \log P(\omega_i)$$

- Fronteras de decisión: hiperplanos.



- Caso 3: Σ_i arbitraria.
 - Matrices de covarianza arbitrarias y diferentes para cada clase, sin restricciones.
 - Las muestras forman nubes hiperelípticas cualquier tamaño y orientación.

$$\begin{aligned}g_i(\mathbf{x}) &= \\ & -\frac{1}{2}\mathbf{x}^t\Sigma_i^{-1}\mathbf{x} + \boldsymbol{\mu}_i^t\Sigma_i^{-1}\mathbf{x} - \frac{1}{2}\left(\boldsymbol{\mu}_i^t\Sigma_i^{-1}\boldsymbol{\mu}_i + d\log 2\pi + \log |\Sigma_i|\right) + \log P(\omega_i) \\ & \equiv -\frac{1}{2}\mathbf{x}^t\Sigma_i^{-1}\mathbf{x} + \boldsymbol{\mu}_i^t\Sigma_i^{-1}\mathbf{x} - \frac{1}{2}\left(\boldsymbol{\mu}_i^t\Sigma_i^{-1}\boldsymbol{\mu}_i + \log |\Sigma_i|\right) + \log P(\omega_i) \\ & = \mathbf{x}^t W_i \mathbf{x} + \mathbf{w}_i^t \mathbf{x} + w_{i0},\end{aligned}$$

con

$$W_i = -\frac{1}{2}\Sigma_i^{-1}$$

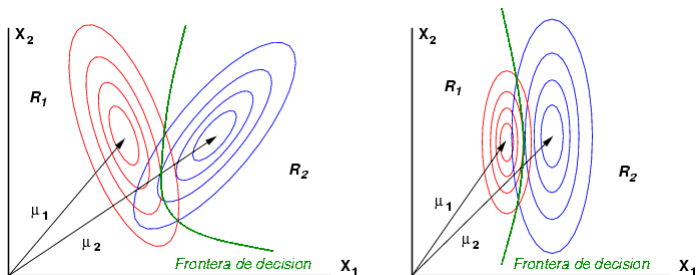
$$\mathbf{w}_i = \Sigma_i^{-1}\boldsymbol{\mu}_i$$

$$w_{i0} = -\frac{1}{2}\boldsymbol{\mu}_i^t\Sigma_i^{-1}\boldsymbol{\mu}_i - \frac{1}{2}\log |\Sigma_i| + \log P(\omega_i)$$

- Caso 3: Σ_i arbitraria.
 - Clasificador cuadrático, con discriminantes

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \frac{1}{2} \log |\Sigma_i| + \log P(\omega_i)$$

- Fronteras de decisión: hipercuádricas.



- Aprendizaje: casos

$P(\mathbf{x}|\omega_i) \approx N(\boldsymbol{\mu}_i, \Sigma_i)$, con $\boldsymbol{\mu}_i$ desconocida: $\Theta_i = \boldsymbol{\mu}_i$

$$P(\mathbf{x}|\omega_i, \Theta_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \Theta_i)^t \Sigma_i^{-1} (\mathbf{x} - \Theta_i)\right)$$

$P(\mathbf{x}|\omega_i) \approx N(\boldsymbol{\mu}_i, \Sigma_i)$, con $\boldsymbol{\mu}_i$ y Σ_i desconocidas: $\Theta_{i1} = \boldsymbol{\mu}_i$, $\Theta_{i2} = \Sigma_i$

$$P(\mathbf{x}|\omega_i, \Theta_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Theta_{i2}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \Theta_{i1})^t \Theta_{i2}^{-1} (\mathbf{x} - \Theta_{i1})\right)$$

- Estimación por máxima verosimilitud: proceso iterativo sobre los datos.

- Principio general de estimación por máxima verosimilitud:
Sean $X = \mathbf{x}_1, \dots, \mathbf{x}_N$: muestras i.i.d. según $P(\mathbf{x}|\Theta)$.
Como las muestras fueron extraídas independientemente:

$$P(X|\Theta) = \prod_{k=1}^N P(\mathbf{x}_k|\Theta),$$

función denominada *verosimilitud* de Θ respecto a X .

La *estimación de máxima verosimilitud* de Θ se define como:

$$\hat{\Theta} = \arg \max_{\Theta} \prod_{k=1}^N P(\mathbf{x}_k|\Theta) = \arg \max_{\Theta} \sum_{k=1}^N \log P(\mathbf{x}_k|\Theta)$$

Se calcula mediante: $\nabla_{\hat{\Theta}} \left(\sum_{k=1}^N \log P(\mathbf{x}_k|\Theta) \right) = \mathbf{0}$

- Caso 1D: μ desconocida

Tenemos que $\Theta = \mu$, entonces

$$P(x_k|\Theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_k - \Theta)^2\right)$$

$$\log P(x_k|\Theta) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(x_k - \Theta)^2$$

$$\sum_{k=1}^N \nabla_{\Theta} \log P(x_k|\hat{\Theta}) = 0 \Rightarrow \hat{\Theta} = \hat{\mu} = \frac{1}{N} \sum_{k=1}^N x_k, \text{ (media muestral)}$$

- Caso 1D: μ y σ^2 desconocidas

Tenemos $\Theta = \begin{pmatrix} \Theta_1 \\ \Theta_2 \end{pmatrix}$, con $\Theta_1 = \mu$ y $\Theta_2 = \sigma^2$, así:

$$P(x_k|\Theta) = \frac{1}{\sqrt{2\pi\Theta_2}} \exp\left(-\frac{1}{2\Theta_2}(x_k - \Theta_1)^2\right)$$

$$\log P(x_k|\Theta) = -\frac{1}{2} \log(2\pi\Theta_2) - \frac{1}{2\Theta_2}(x_k - \Theta_1)^2$$

$$\sum_{k=1}^N \nabla_{\Theta} \log P(x_k|\hat{\Theta}) = \mathbf{0} \Rightarrow \begin{cases} \hat{\Theta}_1 = \hat{\mu} = \frac{1}{N} \sum_{k=1}^N x_k \\ \hat{\Theta}_2 = \hat{\sigma}^2 = \frac{1}{N} \sum_{k=1}^N (x_k - \hat{\mu})^2 \end{cases}$$

- Caso N -D: $\boldsymbol{\mu}$ desconocido

$$\hat{\boldsymbol{\Theta}} = \hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k$$

- Caso N -D: $\boldsymbol{\mu}$ y Σ desconocidos

$$\left\{ \begin{array}{l} \hat{\boldsymbol{\Theta}}_1 = \hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k \\ \hat{\boldsymbol{\Theta}}_2 = \hat{\sigma}^2 = \frac{1}{N} \sum_{k=1}^N (\mathbf{x}_k - \hat{\boldsymbol{\mu}})(\mathbf{x}_k - \hat{\boldsymbol{\mu}})^t \end{array} \right.$$

- 1 Sea un clasificador geométrico lineal definido por:

$$g_1(\mathbf{x}) = -x_1$$

$$g_2(\mathbf{x}) = x_1 + x_2 - 1$$

$$g_3(\mathbf{x}) = x_1 - x_2 - 1$$

- Calcule y grafique las fronteras y regiones de decisión.
- Clasifique los puntos: (2,1); (2,-1); (-1,1); (-1,-1).

- 2 Sean A y B dos clases de igual probabilidad a priori definidas por:

$$P(\mathbf{x}|A) \rightsquigarrow N\left(\begin{pmatrix} 3 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right)$$

$$P(\mathbf{x}|B) \rightsquigarrow N\left(\begin{pmatrix} 0 \\ 3 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right)$$

- Construya un clasificador gaussiano y encuentre la frontera.
- Realice una representación gráfica del problema.

- 3 Implemente un clasificador gaussiano por ML para el mini-corpus de muestras proporcionado, realizando una representación gráfica de la situación.

Calcule el desempeño del sistema al clasificar todas las muestras del corpus.

- Bibliografía:

- Duda R., Hart P, Stork D., Pattern Classification, Second Edition. Wiley-Interscience, 2000. Capítulos 1, 2.1 a 2.6, 3.1 a 3.2.
- Theodoridis S., Koutroumbas K., Pattern Recognition, Fourth Edition. Elsevier, 2009. Capítulo 2 (2.1 a 2.5).
- Murphy K., Machine Learning-A Probabilistic Perspective. The MIT Press, 2012. Capítulos 2 a 4 (secciones salteadas).