

# Probabilidades y Estadística - Clasificación

Federico Yulita

Primer cuatrimestre, 2019.

## Índice

1. Introducción	1
2. <i>Generalized Linear Model</i>	1
3. Perceptrón	2
4. Método Probabilista	2

## 1. Introducción

El problema de clasificación es un ejemplo de aprendizaje supervisado. El problema consiste en clasificar ciertos datos  $\mathbf{x}$  con distintos *labels* (o clases)  $t$  mediante una función lineal para luego poder predecir la clasificación de nuevos datos. La elección de la codificación para las clases - es decir el valor de  $t$  - tiene un gran impacto en la optimización del problema. En la Figura 1 se puede ver un ejemplo de clasificación. En el ejemplo vamos a clasificar datos que caigan por debajo de recta naranja como de clase roja y a los que caigan por encima de clase azul.

Existen dos tipos de métodos para resolver este problema: El **Método Directo** y el **Método Probabilista**. En el método directo se va a buscar optimizar algo sin modelos y en el probabilista se usa un modelo estadístico de los datos.

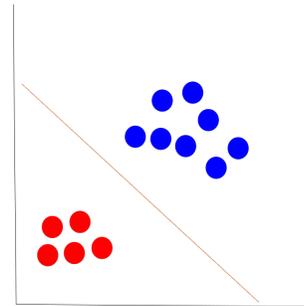


Figura 1: Ejemplo de clasificación.

## 2. *Generalized Linear Model*

En este modelo vamos a usar una función  $f$  de activación con parámetros  $\mathbf{w}$  que vamos a intentar optimizar para que la función clasifique los datos. Consideremos:

$$y_k(x) = \mathbf{w}_k^\dagger \cdot \mathbf{x},$$

donde a  $\mathbf{x} = (1, \mathbf{x}^\dagger)^\dagger$ . Es decir, le agregamos un 1 a los datos para que el primer término proveniente del producto interno es una ordenada al origen. El subíndice  $k$  indica la clase. Entonces:

$$y(\mathbf{x}) = \bar{\mathbf{W}}^\dagger \cdot \mathbf{x}. \quad (1)$$

Definimos la **Función de Error** como:

$$\begin{aligned}
Err(\bar{W}) &= \frac{1}{2} Tr \left( \left( \bar{X} \cdot \bar{W} - \bar{T} \right)^\dagger \cdot \left( \bar{X} \cdot \bar{W} - \bar{T} \right) \right) \\
&= \frac{1}{2} \sum_{i,k} \left( \mathbf{x}_i \cdot \mathbf{w}_k^\dagger - t_i \right)^2.
\end{aligned} \tag{2}$$

Minimizando esta función de error (2) podemos encontrar los parámetros óptimos y una vez hallados podemos clasificar nuevos datos usando la función (1). Estaría bueno que esta última ecuación esté definida en el rango  $[0, 1]$  para usarla como una probabilidad.

Resulta que sale de la minimización analítica que los parámetros óptimos son:

$$\bar{W} = \left( \bar{X}^\dagger \cdot \bar{X} \right)^{-1} \cdot \bar{X}^\dagger \cdot \bar{T}. \tag{3}$$

El problema de este método para clasificar es que los *outliers* tienen un gran peso sobre los parámetros hallados. Para solucionar esto lo que se hace es usar el **Discriminante de Fisher** para optimizar. Esto es, se toma el promedio de los datos de cada clase y entonces se trabajan con varios promedios de cada clase para cada set de datos. Luego se define:

$$s_k^2 = \sum_{n \in C_k} |y_k - \mathbf{w}^\dagger \cdot \mathbf{m}_k|^2,$$

donde  $C_k$  es la clase  $k$  y  $\mathbf{m}_k$  es el promedio de la clase  $k$ . Entonces,  $s_k$  es como una medida de la distancia máxima entre los promedios para cada clase. Entonces se define el discriminante de Fisher (para dos clases) como:

$$J(\mathbf{w}) = \frac{(\mathbf{w}^\dagger \cdot (\mathbf{m}_1 - \mathbf{m}_2))^2}{s_1^2 - s_2^2}. \tag{4}$$

Esta solución es mucho mejor que la otra ya que los outliers ya no pesan sobre la solución. Lo único malo es que no es una probabilidad.

### 3. Perceptrón

Consideremos un vector  $\mathbf{x}$  con *features*  $\Phi(\mathbf{x})$ . definimos  $y(\mathbf{x}) = f(\mathbf{w}^\dagger \cdot \mathbf{x})$ , donde  $f(a) = \begin{cases} 1 & a \geq 0 \\ 0 & a < 0 \end{cases}$ . Vamos

a usar una codificación para las clases de tipo  $t_i \in \{-1, 1\}$ . Decimos que  $\begin{cases} \mathbf{w}^\dagger \cdot \Phi(x_n) < 0 \implies t_n = -1 \\ \mathbf{w}^\dagger \cdot \Phi(x_n) > 0 \implies t_n = 1 \end{cases}$ . El objetivo entonces es hallar los parámetros que minimicen:

$$Err(\mathbf{w}) = - \sum_n \mathbf{w}^\dagger \cdot \Phi(x_n) t_n. \tag{5}$$

Para esto puede usarse el algoritmo de **Gradient Descent**. Este algoritmo consiste en iterar  $\mathbf{w}^\dagger \stackrel{\text{def}}{=} \mathbf{w}^\dagger - \eta \nabla Err(\mathbf{w})$  hasta converger a un mínimo, donde  $\eta$  es algún *step-size*. El problema de este algoritmo es que dependiendo de tus valores iniciales del parámetro puedes llegar a distintas soluciones de los parámetros óptimos. En la Figura 2 puede verse la red neuronal utilizada en este caso. Es de una capa ya que tiene una sola columna de neuronas (en rosa) donde cada una tiene un peso, en el caso de esta función  $f$  solo una neurona tiene un valor no nulo y depende del input dado.

### 4. Método Probabilista

Consideremos la fórmula de Bayes:

$$\mathbb{P}(C_1|x) = \frac{\mathbb{P}(x|C_1)\mathbb{P}(C_1)}{\mathbb{P}(x|C_1)\mathbb{P}(C_1) + \mathbb{P}(x|C_2)\mathbb{P}(C_2)} \quad (6)$$

Notemos que esto es igual a  $\sigma(A)$ , donde  $\sigma$  es la **Función Sigmoidal**:

$$\sigma(x) = \frac{1}{1 + \exp(-x)}, \quad (7)$$

y  $A = \ln\left(\frac{\mathbb{P}(x|C_1)\mathbb{P}(C_1)}{\mathbb{P}(x|C_2)\mathbb{P}(C_2)}\right)$ . Si utilizamos el **Modelo Gaussiano** tenemos que:

$$\mathbb{P}(x|C_k) = \frac{1}{(2\pi)^d \left|\bar{\Sigma}\right|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^\dagger \cdot \bar{\Sigma}^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}_k)\right),$$

donde  $\bar{\Sigma}$  es la matriz de covarianza. Entonces, usando el método de máxima verosimilitud podemos estimar  $\boldsymbol{\mu}_k$ ,  $\bar{\Sigma}$  y  $\mathbb{P}(C_k)$  y usando Bayes para predecir las clases. Vamos a usar que  $\mathbb{P}(C_1|\mathbf{x}) = \sigma(\mathbf{w}^\dagger \cdot \mathbf{x} + w_o)$  y entonces hallamos que:

$$\begin{aligned} \mathbf{w} &= \bar{\Sigma}^{-1} \cdot (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ w_0 &= -\frac{1}{2}\boldsymbol{\mu}_1^\dagger \cdot \bar{\Sigma}^{-1} \cdot \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^\dagger \cdot \bar{\Sigma}^{-1} \cdot \boldsymbol{\mu}_2 + \ln\left(\frac{\mathbb{P}(C_1)}{\mathbb{P}(C_2)}\right). \end{aligned}$$

## Glosario

### D

Discriminante de Fisher, 2

### F

Función de Error, 1

Función Sigmoidal, 3

### G

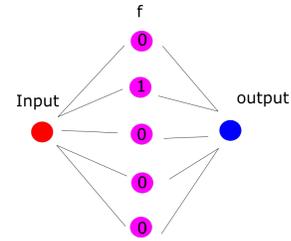
*Gradient Descent*, 2

### M

Método Directo, 1

Método Probabilista, 1

Modelo Gaussiano, 3



**Figura 2:** Red neuronal.