

Clasificación II

Modelos lineales de clasificación

Jonnatan Arias Garcia
jonnatan.arias@utp.edu.co
jariasg@uniquindio.edu.co

David Cardenas peña - dcardenasp@utp.edu.co
Hernán Felipe Garcia - hernanf.garcia@udea.edu.co

Contenido

Modelos discriminativos probabilísticos

- ❑ Logístico

- ❑ Optimizadores

 - ❑ Gradiente descendiente

 - ❑ Mínimos reponderados

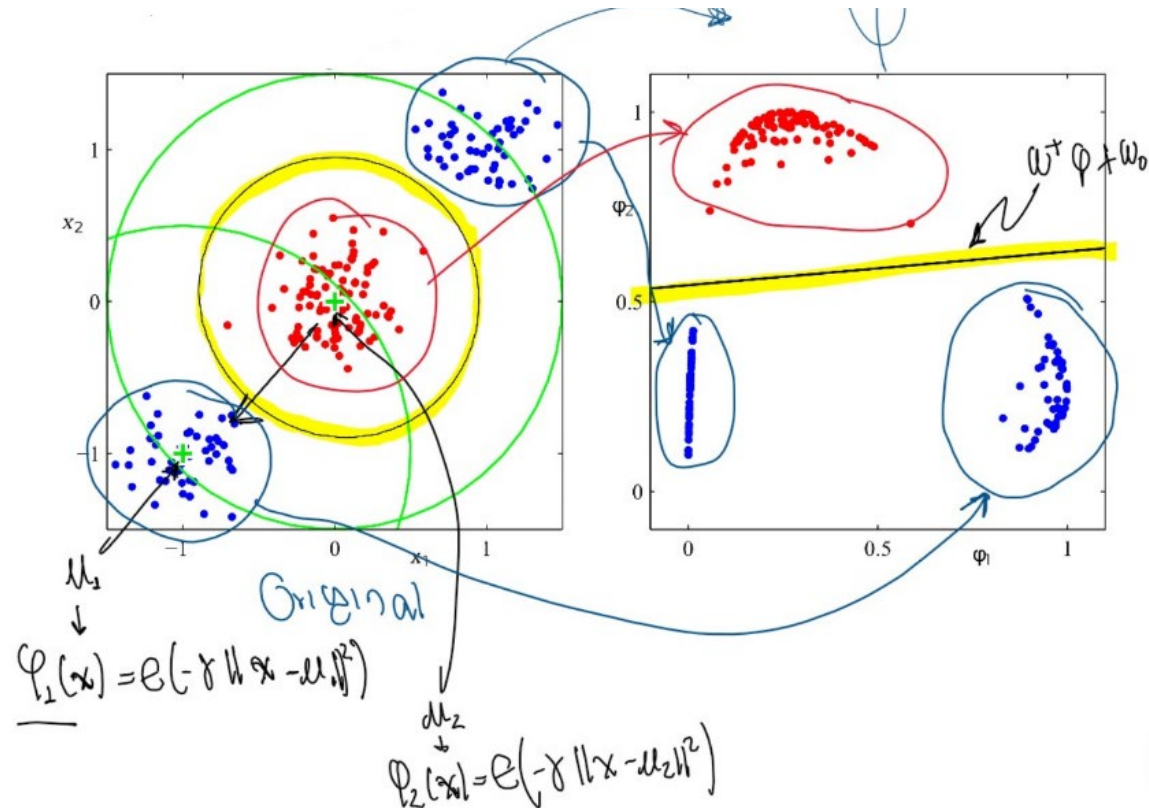
Modelos discriminativos probabilísticos

Modelos discriminativos

- ❑ Consiste en ligar modelos o distribuciones de cada clase por separado, se enfocan principalmente en la mejor manera de diferenciar entre clases.
- ❑ Proporcionan una estimación de la probabilidad de pertenencia
- ❑ Dan información de la incertidumbre debido a la probabilidad de pertenencia

Introducción

- ❑ Se modela directamente la función de probabilidad a posteriori $p(C_k|\mathbf{x})$.
- ❑ En general, se necesita determinar un número menor de parámetros.
- ❑ Se pueden introducir funciones base $\phi(\mathbf{x})$. En el espacio transformado la separación podría ser lineal.



Regresión Logística (i)

- En forma general

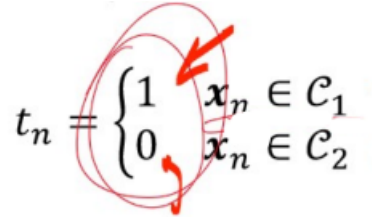
$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^\top \phi),$$

donde $\sigma(\cdot)$ es la función logística sigmoidal.

- En estadística este modelo se conoce como *regresión logística*.
- Sea un conjunto de datos $\{\phi_n, t_n\}_{n=1}^N$, con $\phi_n = \phi(\mathbf{x}_n)$ y $t_n \in \{0, 1\}$.
- La función de verosimilitud se define como

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n},$$

donde $\mathbf{t} = [t_1 \cdots t_N]^\top$, y $y_n = p(C_1|\phi_n)$.


$$t_n = \begin{cases} 1 & \mathbf{x}_n \in C_1 \\ 0 & \mathbf{x}_n \in C_2 \end{cases}$$
$$\phi_n = \phi(\mathbf{x}_n)$$

Regresión Logística (i. A demostración)

- Usando la variable booleana t_n como interruptor, la verosimilitud se puede escribir de forma resumida:

$$\begin{aligned} p(\mathbf{t}|\mathbf{w}) &= \prod_{n=1}^N p(c_1|\varphi_n)^{t_n} p(c_2|\varphi_n)^{1-t_n} \\ p(\mathbf{t}|\mathbf{w}) &= \prod_{n=1}^N p(c_1|\varphi_n)^{t_n} [1 - p(c_1|\varphi_n)]^{1-t_n} \\ p(\mathbf{t}|\mathbf{w}) &= \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n} \end{aligned}$$

$$\begin{aligned} p(c_1|\varphi_n) + p(c_2|\varphi_n) &= 1 \\ p(c_2|\varphi_n) &= 1 - p(c_1|\varphi_n) \end{aligned}$$

Regresión Logística (i. B análisis)

- Se define una función de costo tomando el logaritmo negativo de la función de verosimilitud:

costo $E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n)$

min w costo

$t_{nk} = \begin{cases} 1 & : x_n \in G_k \\ 0 & : x_n \notin G_k \end{cases}$

Regresión Logística (ii)

- Se define una función de error tomando el logaritmo negativo de la función de verosimilitud

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n).$$

donde $\sigma(\cdot)$ es la función logística sigmoideal.

- Se tiene en cuenta la relación $\frac{d\sigma}{da} = \sigma(1 - \sigma)$.
- El gradiente de la función de error con respecto a \mathbf{w} sigue la forma

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \mathbf{\Phi}^\top (\mathbf{y} - \mathbf{t}),$$

donde $\mathbf{y} = [y_1 \cdots y_N]^\top$.

Optimización por Gradiente

Descendiente (i)

□ La minimización de la función de costo se realiza mediante gradiente descendiente, así que:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \frac{de_n(y_n)}{dy_n} \frac{dy_n(\mathbf{w})}{d\mathbf{w}}$$

donde:

$$e_n = -t_n \ln y_n - (1 - t_n) \ln(1 - y_n)$$

$$y_n = \sigma(a_n)$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \frac{de_n(y_n)}{dy_n} \frac{dy_n(a_n)}{da_n} \frac{da_n}{d\mathbf{w}}$$

$$a_n = \mathbf{w}^\top \boldsymbol{\varphi}_n$$

Optimización por Gradiente

Descendiente (i. Dem. II)

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \frac{de_n(y_n)}{dy_n} \frac{dy_n(a_n)}{da_n} \frac{da_n}{d\mathbf{w}}$$

$$\frac{dy_n(a_n)}{da_n} = \nabla(a) [1 - \nabla(a)]$$

$$\frac{da_n}{d\mathbf{w}} = \varphi$$

$$\frac{de_n(y_n)}{dy_n} = -\frac{t_n}{y_n} + \frac{(1-t_n)}{1-y_n}$$

$$\nabla'(a) = \frac{d}{da} \left\{ \frac{1}{1+e(-a)} \right\} = -(1+e(-a))^{-2} \cdot e(-a)(-1)$$

$$= \frac{e(-a)}{(1+e(-a))^2} = \underbrace{\frac{1}{1+e(-a)}}_{\nabla(a)} \cdot \frac{e(-a)}{1+e(-a)}$$

$$= \nabla(a) \frac{e(-a) + 1 - 1}{1+e(-a)} = \nabla(a) \left[\frac{1+e(-a)}{1+e(-a)} - \underbrace{\frac{1}{1+e(-a)}}_{\nabla(a)} \right]$$

$$\nabla_w E = \sum_n \left(-\frac{t_n}{y_n} + \frac{(1-t_n)}{1-y_n} \right) \nabla(a_n) [1 - \nabla(a_n)] \varphi_n$$

$$= \sum_n \left[-\frac{t_n}{y_n} + \frac{y_n(1-t_n)}{1-y_n} \right] [1-y_n] \varphi_n$$

$$= \sum_n \left[-(1-y_n)t_n + \frac{y_n(1-t_n)(1-y_n)}{1-y_n} \right] \varphi_n$$

$$= \sum_n (-t_n + y_n + y_n - y_n t_n) \varphi_n$$

$$\nabla_w E = \sum_n \underbrace{(y_n - t_n)}_{\text{"Error"}} \varphi_n \rightarrow \sum_n (\nabla(\mathbf{w}^T \varphi_n) - t_n) \varphi_n$$

Optimización por Gradiente Descendiente (ii)

- El algoritmo de gradiente descendiente queda:
 1. Inicializar el modelo con \mathbf{w}_0 y fijar una tasa de aprendizaje μ .
 2. Realizar las predicciones con los parámetros actuales $\mathbf{y} = \sigma(\Phi \mathbf{w}_k)$
 3. Calcular la función de costo total para los parámetros actuales
 $L(\mathbf{w}_k) = -\mathbf{t}^\top \ln \mathbf{y} - (1 - \mathbf{t})^\top \ln(1 - \mathbf{y})$.
 4. Actualizar los parámetros para la siguiente iteración $\mathbf{w}_{k+1} = \mathbf{w}_k - \mu \nabla E(\mathbf{w}_k)$
$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mu \Phi^\top (\mathbf{t} - \mathbf{y})$$
 5. Hasta la convergencia volver al paso 2.

Mínimos cuadrados reponderados iterativos (i)

- El gradiente de la función de costo con respecto a \mathbf{w} sigue la forma

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\varphi}_n = \boldsymbol{\Phi}^\top (\mathbf{y} - \mathbf{t})$$

donde $\mathbf{y} \in [0,1]^N$

- La función de costo puede minimizarse también usando el algoritmo de *Newton-Raphson* bajo la regla

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

donde \mathbf{H} es la matriz de segundas derivadas (Hessiana), $\mathbf{H} = \nabla \nabla E(\mathbf{w})$:

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \frac{d}{d\mathbf{w}} \{ \boldsymbol{\Phi}^\top (\mathbf{y}(\mathbf{w}) - \mathbf{t}) \}$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n(\mathbf{w})(1 - y_n(\mathbf{w})) \boldsymbol{\varphi}_n \boldsymbol{\varphi}_n^\top$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \boldsymbol{\Phi}^\top \mathbf{R}(\mathbf{w}) \boldsymbol{\Phi}$$

siendo $\mathbf{R} \in [0,1]^{N \times N}$ una matriz diagonal con elementos $R_{nn} = y_n(1 - y_n)$

Mínimos cuadrados reponderados iterativos (ii)

- Note que la solución para \mathbf{w} debe encontrarse de forma iterativa, debido a que los elementos de \mathbf{R} dependen de \mathbf{w} .
- La solución para \mathbf{w} se puede escribir como

$$\mathbf{w}_{k+1} = (\Phi^T \mathbf{R}_k \Phi)^{-1} \Phi^T \mathbf{R}_k \mathbf{z}_k$$

$$\text{donde } \mathbf{z}_k = \Phi \mathbf{w}_k - \mathbf{R}_k (\mathbf{y}_k - \mathbf{t})$$

- Teniendo en cuenta que:
 - La solución anterior es parecida a la solución de mínimos cuadrados para el problema de regresión lineal.
 - Las ecuaciones normales se deben aplicar iterativamente.

Este algoritmo se conoce como *mínimos cuadrados reponderados iterativos* (IRLS - Iterative Reweighted Least Squares).