

Regresión Logística

Ya conocemos el método de regresión lineal para predecir valores numéricos, pero ¿qué pasa si queremos hacer clasificación?

Uno pensaría que podemos ajustar una regresión lineal y definir un threshold. Así, los valores mayores al threshold son de una clase, mientras que los menores son de otra.

Sin embargo esta opción no nos gusta porque los valores que entrega la regresión no están acotados.

Afortunadamente podemos usar el siguiente truco!

Supongamos una observación con k features x y un vector de coeficientes β . Un modelo útil para clasificar es:

$$\hat{p} = \sigma(x^T \beta) \quad \text{con} \quad \sigma(t) = \frac{1}{1 + e^{-t}} = \frac{e^t}{e^t + 1}$$

↙

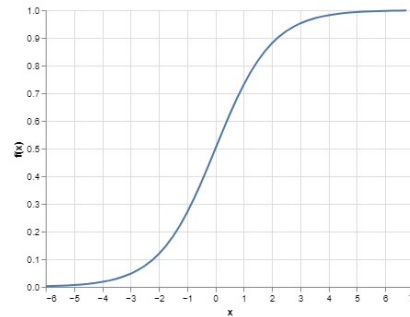
$$\hat{p} = \frac{1}{1 + e^{-(x^T \beta)}}$$

Esto es, intuitivamente, peser el resultado de la regresión lineal por la función logística!

La función logística
va de $]-\infty, \infty[\mapsto [0, 1]$

Luego para predecir:

$$\hat{y} = \begin{cases} 0 & \text{si } \hat{p} < 0.5 \\ 1 & \text{si } \hat{p} \geq 0.5 \end{cases}$$



Y como $\sigma(0) = 0.5$, el modelo predice 1 si $x^T \beta$ es mayor o igual que 0.

Ahora, ¿cómo encontramos los valores de β ? Vamos a usar el método de máxima verosimilitud (Maximum likelihood).

Máxima verosimilitud

El método de máxima verosimilitud nos sirve para encontrar los parámetros de una distribución dado que conocemos ciertas observaciones (datos).

Sean X_1, \dots, X_n valores IID con función de densidad $f(x; \theta)$, queremos encontrar θ !

→ Normalmente uno conoce θ (el parámetro) y pregunta la probabilidad de ver X_i .

Ej. Si sabemos que la estatura distribuye normal de parámetros μ, σ en una población, preguntamos la probabilidad de encontrar a alguien de 1.7 metros

En nuestro caso tenemos observaciones X_1, \dots, X_n que vendrían a ser "estaturas" en el ejemplo anterior, y queremos aprender μ y σ (Ojo: aquí sabemos la distribución, nos faltan sus parámetros).

Ahora, la función de verosimilitud es:

$$\mathcal{L}(\theta) = \prod_{i=1}^n f(x_i; \theta)$$

Y el estimador de máxima verosimilitud $\hat{\theta}$ es el valor de θ que maximiza $\mathcal{L}(\theta)$

Así encontramos el valor de θ que mejor se ajusta a nuestros datos. Pero, ¿cuál es la intuición detrás?

Dado que tenemos las observaciones X_1, \dots, X_n , queremos su distribución conjunta:

$$f(x_1, x_2, \dots, x_n | \theta) = f(x_1 | \theta) \cdot f(x_2 | \theta) \cdot \dots \cdot f(x_n | \theta)$$

Donde sabemos el tipo de distribución (por ejemplo, normal) pero queremos aprender θ , que serían los parámetros de la distribución. Además, recordemos que las observaciones son IID, por eso la distribución conjunta es la multiplicación.

Queremos que el valor de $f(x_1, \dots, x_n | \theta)$ sea lo más grande posible porque "vimos" x_1, \dots, x_n .

Entrenando una regresión logística

Recordemos que $\hat{p} = \sigma(x^T \beta)$ y queremos aprender β . Como esto es clasificación binaria, modelamos esto como una distribución Bernoulli:

$$\left| \begin{array}{l} X \sim \text{Be}(p) \\ f(x) = p^x (1-p)^{1-x} \\ \text{con } x \in \{0, 1\} \end{array} \right|$$

$$P(\underbrace{Y=y}_{\text{La respuesta}} \mid \underbrace{X=x}_{\text{La instancia con sus features}}) = \sigma(x^T \beta) \cdot (1 - \sigma(x^T \beta))^{1-y}$$

$$\log \left(\mathcal{L}(\beta) = \prod_{i=1}^n [\sigma(x_i^T \beta)^{y_i} (1 - \sigma(x_i^T \beta))^{1-y_i}] \right)$$

$$\log(\mathcal{L}(\beta)) = \sum_{i=1}^n [y_i \log(\sigma(x_i^T \beta)) + (1-y_i) \log(1 - \sigma(x_i^T \beta))]$$

Así que lo que hacemos es maximizar $\log(\mathcal{L}(\beta))$ que es equivalente a maximizar $\mathcal{L}(\beta)$ o minimizar $-\log(\mathcal{L}(\beta))$.

También encontraremos en la literatura que busquemos minimizar:

$$-\frac{1}{n} \log(\mathcal{L}(\beta)) = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\sigma(x_i^T \beta)) + (1-y_i) \log(1 - \sigma(x_i^T \beta))]$$

Pero lamentablemente no hay una fórmula explícita para encontrar el óptimo. Sin embargo, al ser una función convexa, podemos usar Gradient Descent:

$$\frac{\partial}{\partial \beta_j} -\frac{1}{n} \log(\mathcal{L}(\beta)) = \frac{1}{n} \sum_{i=1}^n (\sigma(x_i^T \beta) - y_i) x_{ij}$$

Feature j
de la
observación
 i

Así, podemos encontrar el óptimo para nuestra función objetivo. El valor de β en el óptimo son los parámetros que mejor se aproximan para los datos que tenemos.

Odds y logit

Recordemos que nuestro modelo se ve:

$$\hat{p} = \frac{1}{1 + e^{-x^T \beta}}$$

Para entender mejor nuestro modelo (ej. cómo varían las respuestas cuando cambia X) queremos sacar la función e^x . Para esto trabajemos con la razón $\frac{p}{1-p}$ (Odds en inglés).

$$\text{Odds} = \frac{p}{1-p} \longrightarrow p = \frac{\text{Odds}}{1 + \text{Odds}}$$

$$\text{Odds} = e^{x^T \beta}$$

$\log(\text{Odds}) = x^T \beta \rightarrow$ Así podemos conocer cómo cambia nuestra predicción (odds) ante una variación de x .

$$\log(\text{Odds}) = \text{logit}(P) = \log\left(\frac{P}{1-P}\right)$$

Así vemos que estamos usando un modelo lineal para predecir una probabilidad.