

La Regression logistique

February 12, 2024

La régression logistique est une technique de modélisation statistique utilisée pour prédire la probabilité qu'une variable descriptive binaire prenne l'une des deux valeurs possibles (0 ou 1), on peut noter Y cette variable (elle appartient à $\{0, 1\}^n$) en fonction d'un ensemble de variables explicatives que l'on note

$X = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} \in \mathbb{R}^{n \times p}$, avec n observations et p variables. C'est une méthode

couramment utilisée en apprentissage automatique et en statistiques.

La régression logistique a une interprétation probabiliste, elle permet de modéliser $P(Y = 1|\mathbf{x})$, où $\mathbf{x} \in \mathbb{R}^p$.

En utilisant la loi de Bayes et le fait que :

$$P(\mathbf{x}) = P(\mathbf{x}|Y = 1)P(Y = 1) + P(\mathbf{x}|Y = 0)P(Y = 0)$$

nous avons :

$$\begin{aligned} P(Y = 1|\mathbf{x}) &= \frac{P(\mathbf{x}|Y = 1)P(Y = 1)}{P(\mathbf{x}|Y = 1)P(Y = 1) + P(\mathbf{x}|Y = 0)P(Y = 0)} \\ &= \frac{1}{1 + \frac{P(\mathbf{x}|Y=0)P(Y=0)}{P(\mathbf{x}|Y=1)P(Y=1)}} \\ &= \frac{1}{1 + \frac{P(Y=0|\mathbf{x})}{P(Y=1|\mathbf{x})}} \end{aligned}$$

On note $f(\mathbf{x}) := \log \left(\frac{P(Y=1|\mathbf{x})}{P(Y=0|\mathbf{x})} \right)$.

On a ainsi $P(Y = 1|\mathbf{x}) =: \sigma(f(\mathbf{x}))$ avec $\sigma(z) = \frac{1}{1+e^{-z}}$.

La fonction σ , appelée fonction logistique, satisfait les propriétés suivantes :

$$\begin{aligned} \sigma(-z) &= 1 - \sigma(z) \\ \frac{d\sigma(z)}{dz} &= \sigma(z)\sigma(-z). \end{aligned}$$

L'intérêt de la fonction logistique réside dans sa capacité à transformer une fonction f à valeurs dans \mathbb{R} en une probabilité comprise entre 0 et 1.

La régression logistique revient en fait à supposer que f est linéaire de la forme $f : \mathbf{x} \mapsto \boldsymbol{\theta}^\top \mathbf{x}$ avec $\boldsymbol{\theta} \in \mathbb{R}^p$.

Sous cette hypothèse, la règle de classification est simplement :

$$\begin{cases} \text{si } \boldsymbol{\theta}^\top \mathbf{x} \leq 0, \text{ on étiquette 0 au point } \mathbf{x} \\ \text{si } \boldsymbol{\theta}^\top \mathbf{x} > 0, \text{ on étiquette 1 au point } \mathbf{x} \end{cases}$$

On obtient donc :

$$\begin{aligned} P(Y = 1|\mathbf{x}) &= \sigma(\boldsymbol{\theta}^\top \mathbf{x}) \\ P(Y = 0|\mathbf{x}) &= 1 - \sigma(\boldsymbol{\theta}^\top \mathbf{x}) = \sigma(-\boldsymbol{\theta}^\top \mathbf{x}) \end{aligned}$$

Le but maintenant est d'estimer $\boldsymbol{\theta}$. Nous avons $(x_i, y_i)_{1 \leq i \leq n}$ où $x_i \in \mathbb{R}^p$ et $y_i \in \{0, 1\}$ qui constitue un échantillon de taille n . On a alors :

$$P(y = y_i | \mathbf{x} = x_i) = \sigma(\boldsymbol{\theta}^\top x_i)^{y_i} \sigma(-\boldsymbol{\theta}^\top x_i)^{1-y_i}.$$

La log-vraisemblance s'exprime alors de cette manière :

$$\begin{aligned} L(\boldsymbol{\theta}) &= \sum_{i=1}^n \log (\sigma(\boldsymbol{\theta}^\top x_i)^{y_i} \sigma(-\boldsymbol{\theta}^\top x_i)^{1-y_i}) \\ &= \sum_{i=1}^n l(\boldsymbol{\theta}^\top x_i, y_i). \end{aligned}$$

où l est définie comme étant la fonction de perte logistique. Il faut ensuite avoir recours à des algorithmes itératifs (descente de gradient, méthode de Newton,...) pour trouver $\hat{\boldsymbol{\theta}}$.

On peut passer du cadre binaire au cadre multi-classes avec K classes, par exemple, c'est-à-dire en ayant Y appartenant à $\llbracket 1, K \rrbracket^n$. De nouveau, on modélise les probabilités conditionnelles des classes, ou plutôt leur log-ratio, par des quantités linéaires : $\log \left(\frac{P(Y=k|\mathbf{x})}{P(Y=K|\mathbf{x})} \right) = \boldsymbol{\theta}_k^\top \mathbf{x}_i$ pour $k \in \llbracket 1, K-1 \rrbracket$ et $\boldsymbol{\theta}_k \in \mathbb{R}^p$

On a alors pour paramètre globale: $\boldsymbol{\theta} \in [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K] \in \mathbb{R}^{p \times K}$ et pour $k \in \llbracket 1, K \rrbracket$:

$$\begin{aligned} P(y = k|\mathbf{x}) &= \frac{\exp(\langle \boldsymbol{\theta}_k, \mathbf{x} \rangle)}{\sum_{l=1}^K \exp(\langle \boldsymbol{\theta}_l, \mathbf{x} \rangle)}, \\ &= \frac{\exp(\boldsymbol{\theta}_k^\top \mathbf{x})}{\sum_{l=1}^K \exp(\boldsymbol{\theta}_l^\top \mathbf{x})} \end{aligned}$$

On peut écrire cette égalité sous forme vectorielle en utilisant la notation softmax, on a alors : $(P(y = k|\mathbf{x}))_{k=1, \dots, K} = \text{softmax}(\boldsymbol{\theta}_1^\top \mathbf{x}, \dots, \boldsymbol{\theta}_K^\top \mathbf{x})$

Pour la régression softmax, la log-vraisemblance peut être exprimée comme suit :

$$L(\boldsymbol{\theta}) = \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}(y_i = k) \log \left(\frac{e^{\boldsymbol{\theta}_k^\top x_i}}{\sum_{l=1}^K e^{\boldsymbol{\theta}_l^\top x_i}} \right)$$

On utilise ensuite des méthodes algorithmiques pour trouver $\hat{\boldsymbol{\theta}}$.