

Handmade Logistic Regression

Álvaro Orgaz Expósito

THEORY

The aim of the *logistic regression* is to predict the probability of obtaining a positive value in the binary variable y (then it is a *classification* problem) knowing the explanatory variables $x = \{x_1, \dots, x_p\}$.

Firstly, you need to understand how the *logistic regression* use the *linear predictor* $z = w_0 + w_1x_1 + \dots + w_px_p$ to predict $P(y = 1|x) = \pi$. Since $-\infty \leq z \leq \infty$ but we want to predict $0 \leq \pi \leq 1$, we need a *link function* $g(z) = \pi$. Using the log-odds $-\infty \leq \log(\frac{\pi}{1-\pi}) \leq \infty$ we can isolate π from the expression

$$\log\left(\frac{\pi}{1-\pi}\right) = z = w_0 + w_1x_1 + \dots + w_px_p$$

getting the popular *sigmoid link function*

$$\frac{\pi}{1-\pi} = e^z \rightarrow \pi = e^z - \pi e^z \rightarrow \pi + \pi e^z = e^z \rightarrow \pi = \frac{e^z}{1+e^z} = \frac{1}{1+e^{-z}} = \text{sigmoid}(z)$$

Secondly, you need to understand the *likelihood function* for the distribution of the target variable y which is a collection of n *Bernoulli* independent samples $y = \{y_1, \dots, y_n\}$ where $y_i = \{0, 1\}$ with respective probability of positive value $P(y_i = 1|x_i) = \pi_i$.

$$L = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$$

The *likelihood function* L will be used as the loss function in the *gradient descent* optimization. The logic under this product of probabilities is: when $y_i = 1$ we should predict $\pi_i \approx 1$ and when $y_i = 0$ we should predict $\pi_i \approx 0$. The perfect prediction makes $L = 1$, but if the prediction is poor L will decrease to 0.

Thirdly, let's see the mathematical formula for the *gradient descent* used in the optimization code. Since we will predict a binary variable, the loss function J to optimize (minimize) will be the negative logarithm (just to simplify the gradient calculation) of the *likelihood function* defined above

$$J = -\log(L) = -\sum_{i=1}^n \log(\pi_i^{y_i}) + \log((1 - \pi_i)^{1-y_i}) = -\sum_{i=1}^n y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i)$$

where as defined above

$$\pi_i = P(y_i = 1|x_i) = \text{sigmoid}(z_i) = \frac{1}{1 + e^{-z_i}} = \frac{1}{1 + e^{-(w_0 + w_1x_{i1} + \dots + w_px_{ip})}}$$

and n is the number of samples or observations with p explanatory features.

Then, assuming $\pi_i = \hat{y}_i$ for simplicity, the formula for the gradient of J is

$$\frac{\partial J}{\partial w} = \begin{pmatrix} \frac{\partial J}{\partial w_0} \\ \dots \\ \frac{\partial J}{\partial w_p} \end{pmatrix} = \begin{pmatrix} -\sum_{i=1}^n y_i \frac{1}{\hat{y}_i} \frac{\partial \hat{y}_i}{\partial w_0} + (1 - y_i) \frac{1}{1 - \hat{y}_i} \frac{\partial (1 - \hat{y}_i)}{\partial w_0} \\ \dots \\ -\sum_{i=1}^n y_i \frac{1}{\hat{y}_i} \frac{\partial \hat{y}_i}{\partial w_p} + (1 - y_i) \frac{1}{1 - \hat{y}_i} \frac{\partial (1 - \hat{y}_i)}{\partial w_p} \end{pmatrix}$$

where using the *chain rule* for derivatives

$$\frac{\partial \hat{y}_i}{\partial w_p} = \frac{\partial \text{sigmoid}(z_i)}{\partial z_i} \frac{\partial z_i}{\partial w_p} = \frac{e^{-z_i}}{(1 + e^{-z_i})^2} x_{ip}$$

and

$$\frac{\partial(1 - \hat{y}_i)}{\partial w_p} = -\frac{\partial \hat{y}_i}{\partial w_p}$$

But this expression gives a lot of computational errors in the optimization, then let's simplify it as

$$\frac{\partial \hat{y}_i}{\partial w_p} = \frac{1}{(1 + e^{-z_i})} \frac{e^{-z_i}}{(1 + e^{-z_i})} x_{ip} = \frac{1}{1 + e^{-z_i}} \left(\frac{1 + e^{-z_i}}{1 + e^{-z_i}} - \frac{1}{(1 + e^{-z_i})} \right) x_{ip} = \hat{y}_i(1 - \hat{y}_i)x_{ip}$$

Then, substituting in the previous formula, the final gradient of J is

$$\begin{aligned} \frac{\partial J}{\partial w} &= \begin{pmatrix} -\sum_{i=1}^n y_i \frac{1}{\hat{y}_i} \hat{y}_i(1 - \hat{y}_i) + (1 - y_i) \frac{1}{1 - \hat{y}_i} (-1) \hat{y}_i(1 - \hat{y}_i) & \dots \\ -\sum_{i=1}^n y_i \frac{1}{\hat{y}_i} \hat{y}_i(1 - \hat{y}_i)x_{ip} + (1 - y_i) \frac{1}{1 - \hat{y}_i} (-1) \hat{y}_i(1 - \hat{y}_i)x_{ip} \end{pmatrix} = \begin{pmatrix} -\sum_{i=1}^n y_i(1 - \hat{y}_i) + (1 - y_i)(-1)\hat{y}_i & \dots \\ -\sum_{i=1}^n y_i(1 - \hat{y}_i)x_{ip} + (1 - y_i)(-1)\hat{y}_i x_{ip} \end{pmatrix} \\ &= \begin{pmatrix} -\sum_{i=1}^n y_i - y_i \hat{y}_i - \hat{y}_i + y_i \hat{y}_i & \dots \\ -\sum_{i=1}^n y_i x_{ip} - y_i \hat{y}_i x_{ip} - \hat{y}_i x_{ip} + y_i \hat{y}_i x_{ip} \end{pmatrix} = \begin{pmatrix} -\sum_{i=1}^n (y_i - \hat{y}_i) & \dots \\ -\sum_{i=1}^n (y_i - \hat{y}_i)x_{ip} \end{pmatrix} \end{aligned}$$

and we will use the gradient of J to update the weights in each iteration of the optimization process with

$$w^{new} = \begin{pmatrix} w_0^{new} \\ \dots \\ w_p^{new} \end{pmatrix} = \begin{pmatrix} w_0 - \eta \frac{\partial J}{\partial w_0} \\ \dots \\ w_p - \eta \frac{\partial J}{\partial w_p} \end{pmatrix}$$

where η is the learning rate parameter.

CODE

The data used is the popular dataset *titanic_train* in R library *titanic*. Let's predict the binary variable *Survived* with the explanatory features: *Fare* and *Age*. *Note*: The function *complete.cases* help us to select only the observations without missing values in the needed variables.

```
library(titanic)
x1 <- titanic_train[complete.cases(titanic_train[,c("Fare", "Age", "Survived")]), "Fare"]
x2 <- titanic_train[complete.cases(titanic_train[,c("Fare", "Age", "Survived")]), "Age"]
y <- titanic_train[complete.cases(titanic_train[,c("Fare", "Age", "Survived")]), "Survived"]
```

Set initial values for the weights.

```
w0 <- 0
w1 <- 0
w2 <- 0
```

Start the weights optimization with *gradient descent*.

```
# Define the sigmoid function
sigmoid <- function(z){
  1/(1+exp(-z))
}
learning <- 0.000001
rounds <- 1000000
for(i in 1:rounds){
  # Predict the data
  z <- w0+w1*x1+w2*x2
  y_predicted <- sigmoid(z)
```

```

# Calculate the J gradient by weights
error <- y-y_predicted
w0_gradient <- -sum(error*x1)
w1_gradient <- -sum(error*x1)
w2_gradient <- -sum(error*x2)
# Update the weights
w0 <- w0-learning*w0_gradient
w1 <- w1-learning*w1_gradient
w2 <- w2-learning*w2_gradient
}

```

Let's see the optimal estimated weights.

```
c(w0,w1,w2)
```

```
## [1] -0.41705506  0.01725837 -0.01757841
```

Now, compare our optimal weights with the implemented function in R for *logistic regression* and see how it provides the same estimation for the weights.

```
glm(y~x1+x2,family="binomial")$coefficients
```

```
## (Intercept)          x1          x2
## -0.41705506  0.01725837 -0.01757841
```