

Introduction to Statistical learning

Omid Safarzadeh

January 19, 2022

Table of contents

- 1 Simple logistic regression
 - MLE for simple logistic regression
- 2 Multiple logistic regression
 - MLE for multiple logistic regression
- 3 Multiple logistic regression-extension
 - l_1 & l_2 regularized logistic regression
 - Logistic regression for $K > 2$ classes
- 4 Reference

***Acknowledgement:** This slide is prepared based on Murphy, 2012 and James et al., 2013

Logistic regression

- $Y \in \{0, 1\}$. Ex: 0 = *ebola*, 1 = no ebola
- $X \in \mathcal{R}$

$$\pi(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}} \text{ (logistic function)}$$

- $\lim_{X \rightarrow -\infty} \pi(X)$? $\lim_{X \rightarrow +\infty} \pi(X)$?
- $\pi(X)$ models $Pr(Y = 1|X)$
- Odds:

$$\frac{\pi(X)}{1 - \pi(X)} = e^{\beta_0 + \beta_1 X}$$

- Log-odds (logit):

$$\text{logit}(\pi(X)) = \log\left(\frac{\pi(X)}{1 - \pi(X)}\right) = \beta_0 + \beta_1 X$$

- logit is linear in X !

MLE for simple logistic regression

- Data: $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$
- Model: Y_1, \dots, Y_n are independent. $Y_i \sim \text{Bernoulli}(\pi(x_i))$

Likelihood

$$L(\beta_0, \beta_1) = p(\mathcal{D} | \beta_1, \beta_0) = \prod_{i: y_i=1} \pi(x_i) \prod_{i': y_{i'}=0} (1 - \pi(x_{i'}))$$

Log-likelihood

$$l(\beta_0, \beta_1) = \log p(\mathcal{D} | \beta_1, \beta_0) = \sum_{i=1}^n [y_i(\beta_0 + \beta_1 x_i) - \log(1 + e^{\beta_0 + \beta_1 x_i})]$$

MLE

$$(\hat{\beta}_0^{MLE}, \hat{\beta}_1^{MLE}) = \arg_{\beta_0, \beta_1} \max L(\beta_0, \beta_1) = \arg_{\beta_0, \beta_1} \max l(\beta_0, \beta_1)$$

MLE for simple logistic regression

- No closed form solution for $(\hat{\beta}_0^{MLE}, \hat{\beta}_1^{MLE})$
- MLE can be found by **Newton-Raphson method**

Multiple logistic regression

- Response: $Y \in \{0, 1\}$
- Predictors: $\mathbf{X} = [1, X_1, \dots, X_p]^T$
- Parameters: $\beta = [\beta_0, \dots, \beta_p]^T$
- Logistic function:

$$\pi(\mathbf{X}; \beta) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}} = \frac{e^{\beta^T \mathbf{X}}}{1 + e^{\beta^T \mathbf{X}}}$$

- $\pi(\mathbf{X}; \beta)$ models $Pr(Y = 1 | X_1, \dots, X_p; \beta)$
- Odds:

$$\frac{\pi(\mathbf{X}; \beta)}{1 - \pi(\mathbf{X}; \beta)} = e^{\beta^T \mathbf{X}}$$

- Log-odds (logit):

$$\text{logit}(\pi(\mathbf{X}; \beta)) = \log\left(\frac{\pi(\mathbf{X}; \beta)}{1 - \pi(\mathbf{X}; \beta)}\right) = \beta^T \mathbf{X}$$

- logit is linear in \mathbf{X} !

MLE for multiple logistic regression

- Data: $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, $\mathbf{x}_i = [1, x_{i1}, \dots, x_{ip}]^T$
- Model: Y_1, \dots, Y_n are independent.

$$Y_i \sim \text{Bernoulli}(\pi(\mathbf{X}_i))$$

- Log-likelihood

$$l(\beta) = \log p(\mathcal{D}|\beta) = \sum_{i=1}^n [y_i \beta^T \mathbf{x}_i - \log(1 + e^{\beta^T \mathbf{x}_i})]$$

- MLE

$$\hat{\beta}^{MLE} = \arg \max_{\beta \in \mathcal{R}^{p+1}} l(\beta)$$

l_1 & l_2 regularized logistic regression

Optimization problem:

- l_1 regularized logistic regression:

$$(\hat{\beta}_0^{MLE}, \hat{\beta}^{MLE}) = \arg \max_{\beta_0, \beta} \underbrace{\left(\sum_i^n [y_i \beta^T \mathbf{x}_i - \log(1 + e^{\beta^T \mathbf{x}_i})] \right)}_{\text{log-likelihood}} - \lambda \sum_{j=1}^p |\beta_j|$$

- l_2 regularized logistic regression:

$$(\hat{\beta}_0^{MLE}, \hat{\beta}^{MLE}) = \arg \max_{\beta_0, \beta} \underbrace{\left(\sum_i^n [y_i \beta^T \mathbf{x}_i - \log(1 + e^{\beta^T \mathbf{x}_i})] \right)}_{\text{log-likelihood}} - \lambda \sum_{j=1}^p \beta_j^2$$

- What happens as $\lambda \rightarrow 0?$, $\lambda \rightarrow \infty?$

Logistic regression for $K > 2$ classes

- Response: $Y \in \mathcal{C}, \mathcal{C} = \{1, \dots, K\}$.
- Predictors: $\mathbf{X} = [1, X_1, \dots, X_p]^T$
- Parameters: $\beta_j = [\beta_{j0}, \dots, \beta_{jp}]^T, j = 1, \dots, K$

Model:

$$\log \frac{\Pr(Y = 1 | \mathbf{X} = \mathbf{x})}{\Pr(Y = K | \mathbf{X} = \mathbf{x})} = \beta_1^T \mathbf{x}$$

$$\vdots$$

$$\log \frac{\Pr(Y = 1 | \mathbf{X} = \mathbf{x})}{\Pr(Y = K | \mathbf{X} = \mathbf{x})} = \beta_{K-1}^T \mathbf{x}$$

- Parameters vector $\theta = \{\beta_1^T, \dots, \beta_{K-1}^T\}$

Logistic regression for $K > 2$ classes

- Above equations can be solved for each

$$\pi_K(\mathbf{x}; \theta) = \Pr(Y = k | \mathbf{X} = \mathbf{x})$$

$$\pi_j(\mathbf{x}; \theta) = \frac{\exp(\beta_j^T \mathbf{x})}{1 + \sum_{l=1}^{K-1} \exp(\beta_l^T \mathbf{x})}, \quad j = 1, \dots, K-1$$

$$\pi_K(\mathbf{x}; \theta) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp(\beta_l^T \mathbf{x})}$$

- Is $\pi_j(\mathbf{x}; \theta) \in [0, 1]$ for all $j = 1, \dots, K$?
- What is $\sum_{j=1}^K \pi_j(\mathbf{x}; \theta)$?

How to solve logistic regression for $K > 2$ classes

- Similar but not exactly the same as the case $K = 2$
- Apply Newton method.

References

- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning: With applications in r*. Springer New York.
https://books.google.fr/books?id=qcl%5C_AAAAQBAJ
- Murphy, K. (2012). *Machine learning: A probabilistic perspective*. MIT Press.
<https://books.google.fr/books?id=NZP6AQAAQBAJ>