

Logistic Regression

Statistical Machine Learning (ENGG*6600*08)

School of Engineering,
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghogh
Fall 2023

Logistic Regression

- Logistic regression is popular in bio-statistics and bio-informatics.
- Let $\mathbf{x} \in \mathbb{R}^d$ be data and $y \in \mathbb{R}$ be class label. Baye's rule:

$$\mathbb{P}(y|\mathbf{x}) = \frac{\mathbb{P}(\mathbf{x}|y)\mathbb{P}(y)}{\mathbb{P}(\mathbf{x})}, \quad (1)$$

where $\mathbb{P}(y|\mathbf{x})$ and $\mathbb{P}(\mathbf{x}|y)$ are the posterior and likelihood, respectively, and $\mathbb{P}(\mathbf{x})$ and $\mathbb{P}(y)$ are the priors.

- In contrast to Linear Discriminant Analysis (LDA), logistic regression works on the posterior $\mathbb{P}(y|\mathbf{x})$ directly rather than working on likelihood $\mathbb{P}(\mathbf{x}|y)$ and prior $\mathbb{P}(y)$.

Logistic Regression

- Logistic regression is a binary classifier where it assigns probability between zero and one for belonging to one of the classes.
- The logistic function, used in logistic regression, was initially proposed in 1845 for modeling the population growth [1]. It was further improved in the 20th century [2]. See [3] for the history of logistic regression.
- It considers the classification problem as a regression problem where it regresses (predicts) the probability of belonging to a class. It first considers a linear regression $\beta^\top \mathbf{x} + \beta_0$. However, in order to not have the bias, it assumes that \mathbf{x} is $d + 1$ dimensional with an additional element of 1 for bias, i.e., $\mathbf{x} = [x_1, \dots, x_d, 1]^\top$. The $\beta \in \mathbb{R}^{d+1}$ is the learnable parameter of the logistic regression model. As a result, the linear regression becomes $\beta^\top \mathbf{x}$.
- However, there is no bound on this regression while logistic regression desires the output to be in the range $[0, 1]$ to behave like a probability. Therefore, Logistic regression models the posterior using a **logistic function**, also called the **sigmoid function**, to make this regression between zero and one.

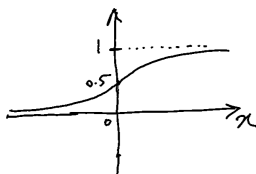
Logistic Regression

- Assume we have two classes $y \in \{0, 1\}$.
- Logistic regression models the posterior using a **logistic function**, also called the **sigmoid function**:

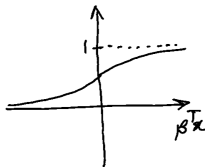
$$\mathbb{P}(y = 1|X = x) = \frac{e^{\beta^\top x}}{1 + e^{\beta^\top x}}, \quad (2)$$

$$\mathbb{P}(y = 0|X = x) = 1 - \mathbb{P}(y = 1|X = x) = \frac{1}{1 + e^{\beta^\top x}}, \quad (3)$$

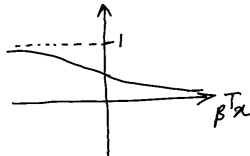
where $\beta \in \mathbb{R}^d$ is the learnable parameter of the logistic regression model.



$$\frac{1}{1 + e^{-x}} \text{ (sigmoid)}$$



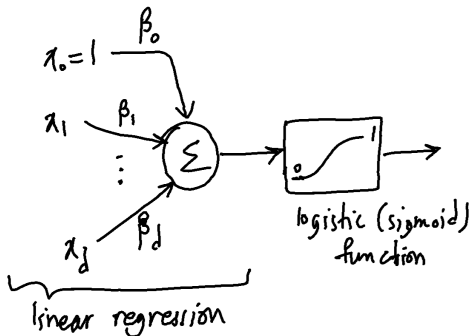
$$\frac{e^{\beta^\top x}}{1 + e^{\beta^\top x}}$$



$$\frac{1}{1 + e^{\beta^\top x}}$$

Logistic Regression as a Neural Network

- Logistic regression can be seen as a neural network with one neuron where the activation function is the nonlinear sigmoid (logistic) function.



Logistic Regression

- Consider n data points $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ in the dataset. Assuming that they are independent and identically distributed (i.i.d), the posterior over all data points is:

$$\mathbb{P}(y|X) = \prod_{i=1}^n \left(\mathbb{P}(y_i = 1|X = \mathbf{x}_i)\mathbb{I}(y_i = 1) + \mathbb{P}(y_i = 0|X = \mathbf{x}_i)\mathbb{I}(y_i = 0) \right), \quad (4)$$

where $\mathbb{I}(\cdot)$ is the indicator function which is one if its condition is satisfied and is zero otherwise.

- As the labels are either zero or one, i.e., $y_i \in \{0, 1\}$, this equation can be restated as:

$$\mathbb{P}(y|X) = \prod_{i=1}^n \left(\mathbb{P}(y_i = 1|X = \mathbf{x}_i) \right)^{y_i} \left(\mathbb{P}(y_i = 0|X = \mathbf{x}_i) \right)^{1-y_i}. \quad (5)$$

- Substituting Eqs. (2) and (3) in this equation gives:

$$\mathbb{P}(y|X) = \prod_{i=1}^n \left(\frac{e^{\boldsymbol{\beta}^\top \mathbf{x}_i}}{1 + e^{\boldsymbol{\beta}^\top \mathbf{x}_i}} \right)^{y_i} \left(\frac{1}{1 + e^{\boldsymbol{\beta}^\top \mathbf{x}_i}} \right)^{1-y_i}. \quad (6)$$

Logistic Regression

- The log posterior is:

$$\begin{aligned}\ell(\beta) &:= \mathbb{P}(y|X = x) = \log \prod_{i=1}^n \left(\frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}} \right)^{y_i} \left(\frac{1}{1 + e^{\beta^\top x_i}} \right)^{1-y_i} \\&= \sum_{i=1}^n \left(\log \left(\frac{e^{\beta^\top x_i}}{1 + e^{\beta^\top x_i}} \right)^{y_i} + \log \left(\frac{1}{1 + e^{\beta^\top x_i}} \right)^{1-y_i} \right) \\&= \sum_{i=1}^n \left(y_i \log(e^{\beta^\top x_i}) - y_i \log(1 + e^{\beta^\top x_i}) - (1 - y_i) \log(1 + e^{\beta^\top x_i}) \right) \\&= \sum_{i=1}^n \left(y_i \beta^\top x_i - y_i \log(1 + e^{\beta^\top x_i}) - \log(1 + e^{\beta^\top x_i}) + y_i \log(1 + e^{\beta^\top x_i}) \right) \\&= \sum_{i=1}^n \left(y_i \beta^\top x_i - \log(1 + e^{\beta^\top x_i}) \right).\end{aligned}$$

Logistic Regression

- The log posterior is:

$$\ell(\beta) = \sum_{i=1}^n \left(y_i \beta^\top \mathbf{x}_i - \log(1 + e^{\beta^\top \mathbf{x}_i}) \right).$$

- Newton's method can be used to find the optimum β . The first derivative, or the gradient, is:

$$\frac{\partial \ell(\beta)}{\partial \beta} = \sum_{i=1}^n \left(y_i \mathbf{x}_i - \frac{1}{1 + e^{\beta^\top \mathbf{x}_i}} e^{\beta^\top \mathbf{x}_i} \mathbf{x}_i \right) = \sum_{i=1}^n \left(y_i - \frac{e^{\beta^\top \mathbf{x}_i}}{1 + e^{\beta^\top \mathbf{x}_i}} \right) \mathbf{x}_i. \quad (7)$$

Its transpose is:

$$\frac{\partial \ell(\beta)}{\partial \beta^\top} = \sum_{i=1}^n \left(y_i - \frac{e^{\beta^\top \mathbf{x}_i}}{1 + e^{\beta^\top \mathbf{x}_i}} \right) \mathbf{x}_i^\top.$$

Logistic Regression

- The second derivative is:

$$\begin{aligned}\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^\top} &= \frac{\partial}{\partial \beta} \left(\frac{\partial \ell(\beta)}{\partial \beta^\top} \right) = \frac{\partial}{\partial \beta} \left(\sum_{i=1}^n \left(y_i - \frac{e^{\beta^\top \mathbf{x}_i}}{1 + e^{\beta^\top \mathbf{x}_i}} \right) \mathbf{x}_i^\top \right) \\ &= \sum_{i=1}^n \left(- \frac{\partial}{\partial \beta} \left(\frac{e^{\beta^\top \mathbf{x}_i}}{1 + e^{\beta^\top \mathbf{x}_i}} \right) \right) \mathbf{x}_i^\top.\end{aligned}$$

- We define:

$$\mathbb{P}(\mathbf{x}_i | \beta) := \frac{e^{\beta^\top \mathbf{x}_i}}{1 + e^{\beta^\top \mathbf{x}_i}}. \quad (8)$$

Therefore:

$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^\top} = - \sum_{i=1}^n \left(\frac{\partial}{\partial \beta} (\mathbb{P}(\mathbf{x}_i | \beta)) \right) \mathbf{x}_i^\top. \quad (9)$$

Logistic Regression

- We have:

$$\begin{aligned}\frac{\partial}{\partial \beta} (\mathbb{P}(\mathbf{x}_i | \beta)) &= \frac{\partial}{\partial \beta} \left(\frac{e^{\beta^\top \mathbf{x}_i}}{1 + e^{\beta^\top \mathbf{x}_i}} \right) \\&= \frac{1}{(1 + e^{\beta^\top \mathbf{x}_i})^2} \left(e^{\beta^\top \mathbf{x}_i} \mathbf{x}_i (1 + e^{\beta^\top \mathbf{x}_i}) - e^{\beta^\top \mathbf{x}_i} (e^{\beta^\top \mathbf{x}_i} \mathbf{x}_i) \right) \\&= \frac{e^{\beta^\top \mathbf{x}_i}}{(1 + e^{\beta^\top \mathbf{x}_i})^2} \left(1 + e^{\beta^\top \mathbf{x}_i} - e^{\beta^\top \mathbf{x}_i} \right) \mathbf{x}_i = \frac{e^{\beta^\top \mathbf{x}_i}}{(1 + e^{\beta^\top \mathbf{x}_i})^2} \mathbf{x}_i \\&= \frac{e^{\beta^\top \mathbf{x}_i}}{(1 + e^{\beta^\top \mathbf{x}_i})} \frac{1}{(1 + e^{\beta^\top \mathbf{x}_i})} \mathbf{x}_i \stackrel{(8)}{=} \mathbb{P}(\mathbf{x}_i | \beta) (1 - \mathbb{P}(\mathbf{x}_i | \beta)) \mathbf{x}_i\end{aligned}$$

- Substituting it in Eq. (9) gives the second derivative, i.e., the Hessian matrix:

$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^\top} = - \sum_{i=1}^n \left(\mathbb{P}(\mathbf{x}_i | \beta) (1 - \mathbb{P}(\mathbf{x}_i | \beta)) \mathbf{x}_i \right) \mathbf{x}_i^\top. \quad (10)$$

Logistic Regression

- It is possible to write the Newton's method in matrix form. We define:

$$\mathbb{R}^{(d+1) \times n} \ni \mathbf{X} := \begin{bmatrix} 1 & 1 & \dots & 1 \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \end{bmatrix},$$

$$\mathbb{R}^{n \times n} \ni \mathbf{W} := \text{diag}\left(\mathbb{P}(\mathbf{x}_i|\boldsymbol{\beta})(1 - \mathbb{P}(\mathbf{x}_i|\boldsymbol{\beta}))\right),$$

$$\mathbb{R}^n \ni \mathbf{y} := [y_1, \dots, y_n]^\top,$$

$$\mathbb{R}^n \ni \mathbf{p} := \left[\frac{e^{\boldsymbol{\beta}^\top \mathbf{x}_1}}{1 + e^{\boldsymbol{\beta}^\top \mathbf{x}_1}}, \dots, \frac{e^{\boldsymbol{\beta}^\top \mathbf{x}_n}}{1 + e^{\boldsymbol{\beta}^\top \mathbf{x}_n}} \right]^\top.$$

- The Eqs. (7) and (10) can be restated as:

$$\mathbb{R}^{(d+1)} \ni \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \mathbf{X}(\mathbf{y} - \mathbf{p}), \quad (11)$$

$$\mathbb{R}^{(d+1) \times (d+1)} \ni \frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} = -\mathbf{X} \mathbf{W} \mathbf{X}^\top. \quad (12)$$

- Using Newton's method for maximization of the log posterior is:

$$\begin{aligned} \boldsymbol{\beta}^{(\tau+1)} &:= \boldsymbol{\beta}^{(\tau)} + \left(\frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} \right)^{-1} \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \implies \\ \boldsymbol{\beta}^{(\tau+1)} &:= \boldsymbol{\beta}^{(\tau)} - (\mathbf{X} \mathbf{W} \mathbf{X}^\top)^{-1} \mathbf{X}(\mathbf{y} - \mathbf{p}), \end{aligned} \quad (13)$$

where τ is the iteration index. It is repeated until convergence of $\boldsymbol{\beta}$.

Logistic Regression

- In the test phase, the class of a point \mathbf{x} is determined as:

$$y = \begin{cases} 1 & \text{if } \frac{e^{\beta^\top \mathbf{x}}}{1 + e^{\beta^\top \mathbf{x}}} \geq 0.5, \\ 0 & \text{Otherwise.} \end{cases} \quad (14)$$

- Comparison to LDA:

- ▶ Logistic regression estimates $(d + 1)$ parameters in β , but LDA estimates many more parameters:
 - ★ prior of each class: 1. We have two classes: $2 \times 1 = 2$.
 - ★ mean of each class: d . We have two classes: $2 \times d = 2d$.
 - ★ covariance matrix of each class: $d(d + 1)/2$. We have two classes: $2 \times (d(d + 1)/2) = d(d + 1)$.
 - ★ so, in total: $2 + 2d + d(d + 1) = d^2 + 2d + 2$.
- ▶ LDA assumes the distribution of each class is Gaussian which may not be true. However, logistic regression does not assume anything about the distribution of data.

Acknowledgment

- Some slides of this slide deck were inspired by teachings of Prof. Ali Ghodsi (at University of Waterloo, Department of Statistics).

References

- [1] P. F. Verhulst, “Resherches mathematiques sur la loi d’accroissement de la population,” *Nouveaux memoires de l’academie royale des sciences*, vol. 18, pp. 1–41, 1845.
- [2] S. H. Walker and D. B. Duncan, “Estimation of the probability of an event as a function of several independent variables,” *Biometrika*, vol. 54, no. 1-2, pp. 167–179, 1967.
- [3] J. S. Cramer, “The origins of logistic regression,” 2002.