

Logistic Regression

Seungjin Choi

Department of Computer Science and Engineering
Pohang University of Science and Technology
77 Cheongam-ro, Nam-gu, Pohang 37673, Korea
seungjin@postech.ac.kr

Binary Classification

- ▶ Consider a binary classification problem, i.e., $\mathcal{C} = \{\mathcal{C}_1, \mathcal{C}_2\}$.
- ▶ Given a particular vector \mathbf{x} , we wish to assign it to one of the two classes. To this end, we use Bayes' rule and calculate the relevant posterior probability:

$$\begin{aligned} P(\mathcal{C}_1|\mathbf{x}) &= \frac{P(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1)}{P(\mathbf{x})} \\ &= \frac{P(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1)}{P(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1) + P(\mathbf{x}|\mathcal{C}_2)P(\mathcal{C}_2)} \\ &= \frac{1}{1 + \frac{P(\mathbf{x}|\mathcal{C}_2)P(\mathcal{C}_2)}{P(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1)}} \\ &= \frac{1}{1 + \exp \left\{ -\log \left[\frac{P(\mathbf{x}|\mathcal{C}_1)}{P(\mathbf{x}|\mathcal{C}_2)} \right] - \log \left[\frac{P(\mathcal{C}_1)}{P(\mathcal{C}_2)} \right] \right\}} \end{aligned}$$

- ▶ It can be written in the form of the **logistic function**:

$$P(C_1|\mathbf{x}) = \frac{1}{1 + e^{-\xi}},$$

where

$$\xi = \underbrace{\log \left[\frac{P(\mathbf{x}|C_1)}{P(\mathbf{x}|C_2)} \right]}_{\text{likelihood ratio}} + \underbrace{\log \left[\frac{P(C_1)}{P(C_2)} \right]}_{\text{prior ratio}}.$$

- ▶ We must choose a particular form for the class-conditional density function $P(\mathbf{x}|C_i)$.

Multivariate Gaussian Model

Assume that the class-conditional densities are multivariate Gaussian with identical covariance matrix Σ :

$$P(\mathbf{x}|\mathcal{C}_i) = \frac{1}{(2\pi)^{\frac{m}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right\}.$$

After a simple calculation, we have

$$P(\mathcal{C}_1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^\top \mathbf{x} + b)}},$$

where

$$\begin{aligned} \mathbf{w} &= \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \\ b &= \frac{1}{2} (\boldsymbol{\mu}_2 + \boldsymbol{\mu}_1)^\top \Sigma^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \log \left[\frac{P(\mathcal{C}_1)}{P(\mathcal{C}_2)} \right]. \end{aligned}$$

Logistic Regression

- ▶ Predict a binary output $y_t \in \{0, 1\}$ from an input \mathbf{x}_t .
- ▶ The logistic regression model has the form

$$y_t = \sigma(\mathbf{w}^\top \mathbf{x}_t) + \epsilon_t,$$

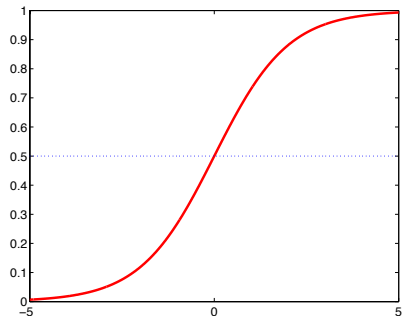
where

$$\sigma(\xi) = \frac{1}{1 + e^{-\xi}} = \frac{e^\xi}{1 + e^\xi}.$$

- ▶ Model input-output by a **conditional Bernoulli distribution**

$$P(y_t = 1 | \mathbf{x}_t) = \sigma(\mathbf{w}^\top \mathbf{x}_t).$$

Properties of Logistic Function



- ▶ $\sigma(\xi) \rightarrow 0$ as $\xi \rightarrow -\infty$.
- ▶ $\sigma(\xi) \rightarrow 1$ as $\xi \rightarrow \infty$.
- ▶ $\sigma(-\xi) = 1 - \sigma(\xi)$.
- ▶ $\frac{d}{d\xi} [\sigma(\xi)] = \sigma(\xi)\sigma(-\xi)$.

Logistic Regression: MLE

- ▶ Given $\{(\mathbf{x}_t, y_t)\}_{t=1}^N$, the likelihood is given by

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \mathbf{w}) &= \prod_{t=1}^N p(y_t = 1|\mathbf{x}_t)^{y_t} (1 - p(y_t = 1|\mathbf{x}_t))^{1-y_t} \\ &= \prod_{t=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_t)^{y_t} (1 - \sigma(\mathbf{w}^\top \mathbf{x}_t))^{1-y_t}. \end{aligned}$$

- ▶ Then log-likelihood function is given by

$$\mathcal{L} = \sum_{t=1}^N \log p(y_t|\mathbf{x}_t) = \sum_{t=1}^N \{y_t \log \sigma_t + (1 - y_t) \log(1 - \sigma_t)\},$$

where $\sigma_t = \sigma(\mathbf{w}^\top \mathbf{x}_t)$.

- ▶ This is a nonlinear function of \mathbf{w} whose maximum cannot be computed in a closed form.
- ▶ **Iterative re-weighted least squares (IRLS)** is a popular algorithm, derived from Newton's method.

Mathematical Preliminaries

- ▶ Gradient
- ▶ Hessian matrix
- ▶ Gradient descent/ascent
- ▶ Newton's method

Gradient

- ▶ Let us consider a real scalar function of a real vector $\mathbf{x} \in \mathbb{R}^m$, $f(\mathbf{x})$,

$$f(\mathbf{x}) : \mathbb{R}^n \longrightarrow \mathbb{R}.$$

- ▶ The gradient of $f(\mathbf{x})$ is defined by

$$\begin{aligned}\nabla f(\mathbf{x}) &= \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \\ &= \left[\frac{\partial f}{\partial x_1} \quad \cdots \quad \frac{\partial f}{\partial x_m} \right]^\top.\end{aligned}$$

Hessian Matrix

If $f(\mathbf{x})$ belongs to the class \mathcal{C}^2 , the Hessian matrix \mathbf{H} is defined as the symmetric matrix with the (i,j) -element $\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}$,

$$\begin{aligned}\mathbf{H} &= \nabla^2 f(\mathbf{x}) \\ &= \left[\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} \right] \\ &= \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_m} \\ \vdots & & & \vdots \\ \frac{\partial^2 f}{\partial x_m \partial x_1} & \frac{\partial^2 f}{\partial x_m \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_m^2} \end{bmatrix} \\ &= \frac{\partial}{\partial \mathbf{x}} \left[\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right]^\top \\ &= \frac{\partial}{\partial \mathbf{x}} [\nabla f(\mathbf{x})]^\top\end{aligned}$$

Gradient Descent/Ascent

- ▶ The gradient descent/ascent learning is a simple (first-order) iterative method for minimization/maximization.
- ▶ **Gradient descent**: iterative minimization

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \left(\frac{\partial \mathcal{J}}{\partial \mathbf{w}} \right).$$

- ▶ **Gradient ascent**: iterative maximization

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \eta \left(\frac{\partial \mathcal{J}}{\partial \mathbf{w}} \right).$$

- ▶ **Learning rate**: $\eta > 0$

Newton's Method

The basic idea of Newton's method is to optimize the quadratic approximation of the objective function $\mathcal{J}(\mathbf{w})$ around the current point \mathbf{w}_t .

The second-order Taylor series expansion of $\mathcal{J}(\mathbf{w})$ at the current point \mathbf{w}_t gives

$$\mathcal{J}(\mathbf{w}_{t+1}) \approx \mathcal{J}(\mathbf{w}_t) + \nabla^\top \mathcal{J}(\mathbf{w}_t) \Delta \mathbf{w}_t + \frac{1}{2} [\Delta \mathbf{w}_t]^\top \nabla^2 \mathcal{J}(\mathbf{w}_t) [\Delta \mathbf{w}_t],$$

where $\nabla^2 \mathcal{J}(\mathbf{w}_t)$ is the Hessian of $\mathcal{J}(\mathbf{w})$ and $\Delta \mathbf{w}_t = \mathbf{w}_{t+1} - \mathbf{w}_t$.

Define

$$\Delta \mathcal{J}(\mathbf{w}_t) = \mathcal{J}(\mathbf{w}_{t+1}) - \mathcal{J}(\mathbf{w}_t).$$

Then, we have

$$\Delta \mathcal{J}(\mathbf{w}_t) = \nabla^\top \mathcal{J}(\mathbf{w}_t) \Delta \mathbf{w}_t + \frac{1}{2} [\Delta \mathbf{w}_t]^\top \nabla^2 \mathcal{J}(\mathbf{w}_t) [\Delta \mathbf{w}_t].$$

Differentiate this w.r.t. $\Delta \mathbf{w}_t$ and set it equal 0, which leads to

$$\nabla \mathcal{J}(\mathbf{w}_t) + \nabla^2 \mathcal{J}(\mathbf{w}_t) \Delta \mathbf{w}_t = 0,$$

Thus we have

$$\Delta \mathbf{w}_t = - [\nabla^2 \mathcal{J}(\mathbf{w}_t)]^{-1} \nabla \mathcal{J}(\mathbf{w}_t).$$

Hence the Newton's method has the form

$$\mathbf{w}_{t+1} = \mathbf{w}_t - [\nabla^2 \mathcal{J}(\mathbf{w}_t)]^{-1} \nabla \mathcal{J}(\mathbf{w}_t).$$

Remark: The Hessian $\nabla^2 \mathcal{J}(\mathbf{w}_t)$ should be **positive definite** for all t .

Logistic Regression: Gradient Ascent

The gradient ascent learning has the form

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \eta \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right).$$

Calculate the gradient

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \sum_t \{y_t - \sigma(\mathbf{w}^\top \mathbf{x}_t)\} \mathbf{x}_t.$$

Hence, the gradient ascent update rule for \mathbf{w} is

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \eta \sum_t \left\{ y_t - \sigma \left((\mathbf{w}^{\text{old}})^\top \mathbf{x}_t \right) \right\} \mathbf{x}_t.$$

Logistic Regression: IRLS

- ▶ Newton's update has the form

$$\Delta \mathbf{w} = \underbrace{\left[\sum_{t=1}^N \sigma_t (1 - \sigma_t) \mathbf{x}_t \mathbf{x}_t^\top \right]^{-1}}_{\text{inverse of Hessian, } [\nabla^2 \mathcal{L}]^{-1}} \underbrace{\left[\sum_{t=1}^N (y_t - \sigma_t) \mathbf{x}_t \right]}_{\text{gradient, } \nabla \mathcal{L}},$$

- ▶ Newton's update reduces to **iterative re-weighted least squares (IRLS)**:

$$\Delta \mathbf{w} = \left(\mathbf{X}^\top \mathbf{V}_r \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{V}_r \mathbf{Y}^*,$$

where

$$\mathbf{V}_r = \begin{bmatrix} \sigma_1(1 - \sigma_1) & & 0 \\ \vdots & \ddots & \vdots \\ 0 & & \sigma_N(1 - \sigma_N) \end{bmatrix},$$
$$\mathbf{Y}^* = \begin{bmatrix} \frac{y_1 - \sigma_1}{\sigma_1(1 - \sigma_1)} & & 0 \\ \vdots & \ddots & \vdots \\ 0 & & \frac{y_N - \sigma_N}{\sigma_N(1 - \sigma_N)} \end{bmatrix}.$$

Logistic Regression: Generalized Linear Model

- ▶ Bernoulli distribution for a univariate binary random variable $y \in \{0, 1\}$ with mean p ,

$$P(y|p) = p^y(1-p)^{1-y}.$$

Logit transform (log-odds), $\theta = \log\left(\frac{p}{1-p}\right)$, leads to $[0, 1] \mapsto \mathbb{R}$.

- ▶ The logit transform gives $p = \frac{1}{1+e^{-\theta}} = \sigma(\theta)$.
- ▶ Thus, we have

$$\begin{aligned} P(y|\theta) &= \sigma(\theta)^y (1 - \sigma(\theta))^{1-y} \\ &= \sigma(\theta)^y \sigma(-\theta)^{1-y} \end{aligned}$$

Multiclass Logistic Regression (Multinomial Logistic Regression)

- ▶ Model input-output by a **softmax transformation** of linear functions of inputs:

$$p(y_t = k | \mathbf{x}_t) = \frac{\exp\{\mathbf{w}_k^\top \mathbf{x}_t\}}{\sum_{j=1}^K \exp\{\mathbf{w}_j^\top \mathbf{x}_t\}}.$$

- ▶ Given $\mathbf{Y} \in \mathbb{R}^{K \times N}$ (each column $\mathbf{y}_t \in \mathbb{R}^K$ follows the 1-of- K coding) and $\mathbf{X} \in \mathbb{R}^{D \times N}$, the likelihood is given by

$$p(\mathbf{Y} | \mathbf{X}, \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{t=1}^N \prod_{k=1}^K p(y_t = k | \mathbf{x}_t)^{Y_{kt}}.$$

- ▶ The log-likelihood is given by

$$\mathcal{L} = \sum_{t=1}^N \sum_{k=1}^K Y_{kt} \log [p(y_t = k | \mathbf{x}_t)].$$

- ▶ One can apply Newton's update to derive IRLS, as in logistic regression.