

Bayesian Logistic Regression

Assaf Dvora

May 1, 2022

1 Model Definition

We observe a set of pairs $\mathcal{D} = \{(\mathbf{x}_i, y_i) | i = 1, \dots, N\}$, with $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in [0, 1]$ (a binary classification response). we assume that the class-conditional probability of belonging to the “1” class is given by a nonlinear transformation of a linear function of \mathbf{x} :

$$P(y = 1 | \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w}). \quad (1)$$

The most commonly used function σ is the logistic function:

$$\sigma(z) = \frac{1}{1 + \exp(-z)}. \quad (2)$$

The class-conditional probability can be formulated as Bernoulli distribution with probability of success $p = \sigma(\mathbf{x}^T \mathbf{w})$:

$$p(y | \mathbf{x}, \mathbf{w}) = \left(\frac{1}{1 + \exp(-\mathbf{x}^T \mathbf{w})} \right)^{y_i} \left(\frac{\exp(-\mathbf{x}^T \mathbf{w})}{1 + \exp(-\mathbf{x}^T \mathbf{w})} \right)^{1-y_i} \quad (3)$$

We further assume that that predictions for different \mathbf{x} are independent given \mathbf{w} , then the likelihood can be written

$$p(\mathcal{Y} | \mathcal{X}, \mathbf{w}) = \prod_{i=1}^N p(y_i | \mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^N \left(\frac{1}{1 + \exp(-\mathbf{x}_i^T \mathbf{w})} \right)^{y_i} \left(\frac{\exp(-\mathbf{x}_i^T \mathbf{w})}{1 + \exp(-\mathbf{x}_i^T \mathbf{w})} \right)^{1-y_i} \quad (4)$$

Since we are constructing a Bayesian model, we must assign a prior distribution on the unknown variables in the model. We choose zero mean normal priors with variance s^2 for the weights \mathbf{w} which corresponds to weak information regarding the true parameters values. I.e.,

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_p). \quad (5)$$

Using Bayes theorem, the posterior distribution is given by

$$p(\mathbf{w}|\mathcal{D}) = p(\mathbf{w}|\mathcal{Y}, \mathcal{X}) = \frac{p(\mathcal{Y}|\mathcal{X}, \mathbf{w})p(\mathbf{w})}{p(\mathcal{Y}|\mathcal{X})} = \frac{p(\mathcal{Y}|\mathcal{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathcal{Y}|\mathcal{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}}. \quad (6)$$

To make predictions based the training data \mathcal{D} for a test point \mathbf{x}_* we have

$$p(y_*|\mathbf{x}_*, \mathcal{D}) = \int p(y_*|\mathbf{w}, \mathbf{x}_*, \mathcal{D})p(\mathbf{w}|\mathcal{D}) \quad (7)$$

We note that the posterior distribution, $p(\mathbf{w}|\mathcal{D})$, does not belong to a nice parametric family. Furthermore, the integral $p(y_*|\mathbf{x}_*, \mathcal{D})$ is intractable as well.

2 MAP Estimation

The MAP estimates, $\hat{\mathbf{w}}$ is defined as:

$$\begin{aligned} \hat{\mathbf{w}} &= \underset{\mathbf{w}}{\operatorname{argmax}}[p(\mathbf{w}|\mathcal{D})] = \underset{\mathbf{w}}{\operatorname{argmax}}[p(\mathcal{Y}|\mathcal{X}, \mathbf{w})p(\mathbf{w})] = \\ &= \underset{\mathbf{w}}{\operatorname{argmax}}\left\{\sum_{i=1}^N y_i \log[\sigma(\mathbf{x}_i^T \mathbf{w})] + \sum_{i=1}^N (1 - y_i) \log[1 - \sigma(\mathbf{x}_i^T \mathbf{w})] - 0.5 \mathbf{w}^T \Sigma_p \mathbf{w}\right\} \end{aligned} \quad (8)$$

we note that $\hat{\mathbf{w}}$ has no simple analytic form. However, it is easy to show that for some sigmoid functions, such as logistic and cumulative Gaussian, the log likelihood (4) is a concave function of \mathbf{w} for fixed \mathcal{D} . As the quadratic term, $0.5 \mathbf{w}^T \Sigma_p \mathbf{w}$, is also concave then the log posterior is a concave function, which means that it is relatively easy to find its unique maximum.

Although finding the MAP is a fast and easy way of obtaining estimates of the unknown model parameters, it is limited because there is no associated estimate of uncertainty produced with the MAP estimates.

3 Laplace Approximation

The idea of Laplace approximation is to approximate the posterior density with a Gaussian:

$$p(\mathbf{w}|\mathcal{D}) \approx \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (9)$$

The Laplace approximation is based on a second order Taylor expansion of the negative log of the unnormalized posterior

$$\Psi(\mathbf{w}) = -\log p(\mathcal{D}|\mathbf{w}) - \log p(\mathbf{w}) \quad (10)$$

around the MAP estimate $\hat{\mathbf{w}}$. This results in the following Gaussian distribution

$$p(\mathbf{w}|\mathcal{D}) \approx \mathcal{N}(\mathbf{w}|\hat{\mathbf{w}}, H^{-1}) \quad (11)$$

where H is the Hessian of $\Psi(\mathbf{w})$ evaluated at $\hat{\mathbf{w}}$:

$$H = \nabla \nabla \Psi(\mathbf{w})|_{\mathbf{w}=\hat{\mathbf{w}}} \quad (12)$$

In the case of logistic regression, the Hessian is given by

$$H = \Sigma_p^{-1} + \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \sigma(\mathbf{x}_i^T \hat{\mathbf{w}}) [1 - \sigma(\mathbf{x}_i^T \hat{\mathbf{w}})] \quad (13)$$

To make a prediction, one needs the predictive distribution (7). However, the integral can not be computed analytically even with the Gaussian approximation of the posterior (11). A numerical approximation can however be easily obtained by Monte Carlo sampling (SLLN)

$$p(y_* | \mathbf{x}_*, \mathcal{D}) \approx \frac{1}{S} \sum_{i=1}^S \sigma(\mathbf{x}_{new}^T \mathbf{w}_s)$$

where \mathbf{w}_s are independently sampled from $\mathcal{N}(\mathbf{w} | \hat{\mathbf{w}}, H^{-1})$

4 Multiclass logistic regression

In the multiclass logistic regression formulation, the responses are considered to be the set of 1-of- K encoded random vectors \mathbf{y} of dimension K having the property that exactly one element has the value 1 and the others have the value 0. In this formulation, the response vector, \mathbf{y} , can be modeled as categorical variable,

$$P(\mathbf{y} | \mathbf{x}, \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{k=1}^K p_k^{y_k}$$

where p_k are the class conditional densities given by the softmax function,

$$p_k = P(C_k | \mathbf{x}, \mathbf{w}_1, \dots, \mathbf{w}_K) = \frac{\exp(\mathbf{w}_k^T \mathbf{x})}{\sum_{k'} \exp(\mathbf{w}_{k'}^T \mathbf{x})}$$

Assuming the data samples \mathbf{x}_n , $n = 1, \dots, N$ are statistically independent given , the likelihood can be written

$$p(\mathcal{Y} | \mathcal{X}, \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N p(\mathbf{y}_n | \mathbf{x}_n, \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K \left(\frac{\exp(\mathbf{w}_k^T \mathbf{x}_n)}{\sum_{k'} \exp(\mathbf{w}_{k'}^T \mathbf{x}_n)} \right)^{y_{nk}}$$

since we are constructing a Bayesian model, we must assign a prior distribution on the unknown variables in the model. We model each weights vector \mathbf{w}_k with multivariate normal prior with zero mean vector and covariance matrix variance $\mathbf{I}_p s^2$ which corresponds to weak information regarding the true parameters values. I.e.,

$$p(\mathbf{w}_k) = \mathcal{N}(\mathbf{0}, \Sigma_p) = \mathcal{N}(\mathbf{0}, \mathbf{I}_p s^2)$$

Since $\mathbf{w}_k, k = 1, \dots, K$ are assumed statistically independent we obtain

$$p(\mathbf{W}) = p(\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{k=1}^K p(\mathbf{w}_k) = \prod_{k=1}^K \mathcal{N}(\mathbf{0}, \mathbf{I}_p s^2)$$

Using Bayes theorem, the posterior distribution is given by

$$\begin{aligned} p(\mathbf{W}|D) &= p(\mathbf{w}_1, \dots, \mathbf{w}_K|\mathcal{D}) = p(\mathbf{w}_1, \dots, \mathbf{w}_K|\mathcal{Y}, \mathcal{X}) = \\ &= \frac{p(\mathcal{Y}|\mathcal{X}, \mathbf{w}_1, \dots, \mathbf{w}_K)p(\mathbf{w}_1, \dots, \mathbf{w}_K)}{p(\mathcal{Y}|\mathcal{X})} = \\ &= \frac{p(\mathcal{Y}|\mathcal{X}, \mathbf{w}_1, \dots, \mathbf{w}_K)p(\mathbf{w}_1, \dots, \mathbf{w}_K)}{\int p(\mathcal{Y}|\mathcal{X}, \mathbf{w}_1, \dots, \mathbf{w}_K)p(\mathbf{w}_1, \dots, \mathbf{w}_K)d\mathbf{w}_1 \dots d\mathbf{w}_K} \end{aligned}$$

To make predictions based on the training data \mathcal{D} for a test point \mathbf{x}_* we have

$$p(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}) = \int p(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}, \mathbf{W})p(\mathbf{W}|D)$$

The MAP estimates, $\hat{\mathbf{w}}_k$ is defined as:

$$\begin{aligned} \hat{\mathbf{w}}_k &= \arg \max_{\mathbf{w}_k} p(\mathbf{w}_1, \dots, \mathbf{w}_K|\mathcal{D}) = \arg \max_{\mathbf{w}_k} \log p(\mathbf{w}_1, \dots, \mathbf{w}_K|\mathcal{D}) = \\ &= \arg \max_{\mathbf{w}_k} [\log p(\mathcal{Y}|\mathcal{X}, \mathbf{w}_1, \dots, \mathbf{w}_K) + \log p(\mathbf{w}_1, \dots, \mathbf{w}_K)] \\ &= \arg \max_{\mathbf{w}_k} \left[\sum_{n=1}^N \sum_{k=1}^K y_{nk} \log \frac{\exp(\mathbf{w}_k^T \mathbf{x}_n)}{\sum_{k'} \exp(\mathbf{w}_{k'}^T \mathbf{x}_n)} - 0.5s^{-2} \sum_{k=1}^K \mathbf{w}_k^T \mathbf{w}_k \right] \end{aligned}$$

The unnormalized log posteriror is given by

$$\Psi(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N \sum_{k=1}^K y_{nk} \log \frac{\exp(\mathbf{w}_k^T \mathbf{x}_n)}{\sum_{k'} \exp(\mathbf{w}_{k'}^T \mathbf{x}_n)} + 0.5s^{-2} \sum_{k=1}^K \mathbf{w}_k^T \mathbf{w}_k$$

The Hessian matrix is given by (see Section 5)

$$H = \mathcal{X}D\mathcal{X}^T + \Sigma_p^{-1}$$

where

$$\begin{aligned} \mathcal{D} &\in \mathbb{R}^{N \times N}, \mathcal{D} = \text{diag}(d_1, d_2, \dots, d_N), d_i = [1 - \sigma(\mathbf{x}_i^T \mathbf{w})]\sigma(\mathbf{x}_i^T \mathbf{w}) \\ \mathcal{X} &\in \mathbb{R}^{d \times N}, \mathcal{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \end{aligned}$$

5 Appendix A - Derivation of the Hessian Matrix

The Hessian Matrix is defined as follows

$$H = \nabla \nabla \Psi(\mathbf{w}) = \frac{\partial \Psi(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} \quad (14)$$

for selecting $\Psi(\mathbf{w})$ to be the unnormalized log posterior for the logistic regression model we obtain

$$H = \frac{\partial \Psi(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = \frac{\partial}{\partial \mathbf{w} \partial \mathbf{w}^T} \left\{ - \sum_{i=1}^N y_i \log[\sigma(\mathbf{x}_i^T \mathbf{w})] - \sum_{i=1}^N (1 - y_i) \log[1 - \sigma(\mathbf{x}_i^T \mathbf{w})] + 0.5 \mathbf{w}^T \Sigma_p^{-1} \mathbf{w} \right\} \quad (15)$$

We note the following properties of logistic function

$$1 - \sigma(z) = 1 - \frac{1}{1 + \exp(-z)} = \frac{\exp(-z)}{1 + \exp(-z)} = \frac{1}{1 + \exp(z)} = \sigma(-z) \quad (16)$$

$$\frac{\partial \sigma(z)}{\partial z} = \frac{\exp(-z)}{[1 + \exp(-z)]^2} = \frac{\exp(-z)}{1 + \exp(-z)} \frac{1}{1 + \exp(-z)} = [1 - \sigma(z)]\sigma(z) \quad (17)$$

We start by taking the derivative of $\Psi(\mathbf{w})$ with respect to \mathbf{w}^T

$$\begin{aligned} \frac{\partial \Psi(\mathbf{w})}{\partial \mathbf{w}^T} &= \frac{\partial}{\partial \mathbf{w}^T} \left\{ - \sum_{i=1}^N y_i \log[\sigma(\mathbf{x}_i^T \mathbf{w})] - \sum_{i=1}^N (1 - y_i) \log[1 - \sigma(\mathbf{x}_i^T \mathbf{w})] + 0.5 \mathbf{w}^T \Sigma_p^{-1} \mathbf{w} \right\} = \\ &\stackrel{(16)}{=} \frac{\partial}{\partial \mathbf{w}^T} \left\{ - \sum_{i=1}^N y_i \log[\sigma(\mathbf{x}_i^T \mathbf{w})] - \sum_{i=1}^N (1 - y_i) \log[\sigma(-\mathbf{x}_i^T \mathbf{w})] + 0.5 \mathbf{w}^T \Sigma_p^{-1} \mathbf{w} \right\} \\ &= - \sum_{i=1}^N y_i \sigma(\mathbf{x}_i^T \mathbf{w})^{-1} \frac{\partial}{\partial \mathbf{w}^T} \sigma(\mathbf{x}_i^T \mathbf{w}) - \sum_{i=1}^N (1 - y_i) \sigma(-\mathbf{x}_i^T \mathbf{w})^{-1} \frac{\partial}{\partial \mathbf{w}^T} \sigma(-\mathbf{x}_i^T \mathbf{w}) + \Sigma_p^{-1} \mathbf{w} = \end{aligned} \quad (18)$$

$$\begin{aligned} &= - \sum_{i=1}^N y_i [1 - \sigma(\mathbf{x}_i^T \mathbf{w})] \mathbf{x}_i + \sum_{i=1}^N (1 - y_i) [1 - \sigma(-\mathbf{x}_i^T \mathbf{w})] \mathbf{x}_i + \Sigma_p^{-1} \mathbf{w} = \\ &= \sum_{i=1}^N \{ -y_i [1 - \sigma(\mathbf{x}_i^T \mathbf{w})] + (1 - y_i) \sigma(\mathbf{x}_i^T \mathbf{w}) \} \mathbf{x}_i + \Sigma_p^{-1} \mathbf{w} = \\ &= \sum_{i=1}^N [-y_i + \sigma(\mathbf{x}_i^T \mathbf{w})] \mathbf{x}_i + \Sigma_p^{-1} \mathbf{w} \end{aligned} \quad (19)$$

To obtain the Hessian matrix, we take the derivative of (18) with respect to \mathbf{w}

$$\begin{aligned}
H &= \frac{\partial \Psi(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = \frac{\partial}{\partial \mathbf{w}} \left\{ \sum_{i=1}^N [-y_i + \sigma(\mathbf{x}_i^T \mathbf{w})] \mathbf{x}_i + \Sigma_p^{-1} \mathbf{w} \right\} = \\
&= \sum_{i=1}^N \mathbf{x}_i \frac{\partial}{\partial \mathbf{w}} \sigma(\mathbf{x}_i^T \mathbf{w}) + \Sigma_p^{-1} = \\
&= \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T [1 - \sigma(\mathbf{x}_i^T \mathbf{w})] \sigma(\mathbf{x}_i^T \mathbf{w}) + \Sigma_p^{-1}
\end{aligned} \tag{20}$$

Denote

$$\mathcal{X} \in \mathbb{R}^{d \times N}, \mathcal{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \tag{21}$$

$$\mathcal{D} \in \mathbb{R}^{N \times N}, \mathcal{D} = \text{diag}(d_1, d_2, \dots, d_N), d_i = [1 - \sigma(\mathbf{x}_i^T \mathbf{w})] \sigma(\mathbf{x}_i^T \mathbf{w}) \tag{22}$$

Thus (20) can be expressed as

$$H = \mathcal{X} \mathcal{D} \mathcal{X}^T + \Sigma_p^{-1}$$