

Bayesian Logistic Classification

3F8: Inference Coursework
Theo Brown
Selwyn College, University of Cambridge

March 12, 2022

Abstract

Abstract goes here

1 Introduction

2 Laplace approximation for logistic regression

The Laplace approximation finds a Gaussian, $q(\mathbf{x})$, that provides a good approximation to a probability distribution $p(\mathbf{x})$ near a local maximum of $p(\mathbf{x})$. This is useful when the desired distribution $p(\mathbf{x})$ is intractable and difficult to manipulate.

In the previous report, the maximum-likelihood estimate for the model weights was used in classification¹. For Bayesian classification, a full posterior distribution of the weights is required:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{X}|\mathbf{w}, \mathbf{y})p(\mathbf{w})}{p(\mathbf{X}|\mathbf{y})} \quad (1)$$

Calculating the normalising constant (also called the evidence) $p(\mathbf{X}|\mathbf{y})$ requires integrating a product of logistic functions, which is intractable. The Laplace approximation can be used to find an approximate posterior distribution and an approximate predictive distribution for the logistic classifier, while avoiding the difficult integral. The derivation of these results is found in Appendix A.

In this application, the prior on the model weights is chosen to be Gaussian:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; 0, \mathbf{I}\sigma_0^2) \quad (2)$$

This gives the Laplace approximation of the posterior as:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) \approx \mathcal{N}(\mathbf{w}; \mathbf{w}_{\text{MAP}}, \mathbf{C}_N) \quad (3)$$

And the predictive distribution as:

$$p(y^* = 1|\mathbf{x}^*, \mathbf{y}, \mathbf{X}) \approx \sigma \left(\frac{\mu_p}{\sqrt{1 + \sigma_p^2 \lambda^2}} \right) \quad (4)$$

where

$$\mathbf{C}_N = \left(\frac{1}{\sigma_0^2} \mathbf{I} + \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \sigma(1 - \sigma) \right)^{-1} \quad (5)$$

$$\sigma = \sigma(\tilde{\mathbf{X}}^T \mathbf{w}) \quad (6)$$

$$\mu_p = \mathbf{x}^{*T} \mathbf{w}_{\text{MAP}} \quad (7)$$

$$\sigma_p^2 = \mathbf{x}^{*T} \mathbf{C}_N \mathbf{x}^* \quad (8)$$

$$\lambda = \sqrt{\frac{\pi}{8}} \quad (9)$$

¹For the model definition, and definitions of \mathbf{X} , \mathbf{y} , $\tilde{\mathbf{x}}$, etc, see the previous report.

A Deriving Bayesian Logistic Classification using the Laplace Approximation

A.1 Laplace approximation

To define a probability distribution function (PDF) $p(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^N$ it is necessary to integrate a function to find the normalising constant K :

$$p(\mathbf{x}) := \frac{f(\mathbf{x})}{\int f(\mathbf{x})d\mathbf{x}} = \frac{1}{K}f(\mathbf{x}) \quad (10)$$

In many cases, the integral $\int f(\mathbf{x})d\mathbf{x}$ is intractable or does not have a closed-form solution.

The Laplace approximation finds a Gaussian, $q(\mathbf{x})$, that provides a good approximation to $p(\mathbf{x})$ near a local maximum of $p(\mathbf{x})$. As $q(\mathbf{x})$ is Gaussian, its normalising constant is easy to find, so the problem of solving $\int f(\mathbf{x})d\mathbf{x}$ is avoided. To find $q(\mathbf{x})$, we start with the truncated Taylor expansion of $\ln f(\mathbf{x})$ around a local maximum \mathbf{x}_0 :

$$\ln f(\mathbf{x}) \approx \ln f(\mathbf{x}_0) + \nabla \ln f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \nabla^2 \ln f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}(\mathbf{x} - \mathbf{x}_0)$$

At a maximum of $f(\mathbf{x})$, $\nabla \ln f(\mathbf{x}) = 0$ as the logarithm is a monotonic function. Hence, close to \mathbf{x}_0 :

$$\begin{aligned} \ln f(\mathbf{x}) &\approx \ln f(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \nabla^2 \ln f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}(\mathbf{x} - \mathbf{x}_0) \\ f(\mathbf{x}) &\approx f(\mathbf{x}_0) \exp\left(\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \nabla^2 \ln f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}(\mathbf{x} - \mathbf{x}_0)\right) \end{aligned} \quad (11)$$

Equation 11 is of the form of an un-normalised Gaussian. Let $\mathbf{P} = -\nabla^2 \ln f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}$, and normalise Equation 11 to get an approximation for $p(\mathbf{x})$:

$$\begin{aligned} p(\mathbf{x}) &\approx \frac{1}{(2\pi)^{\frac{M}{2}} \det \mathbf{P}^{-\frac{1}{2}}} \exp\left(\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \mathbf{P}(\mathbf{x} - \mathbf{x}_0)\right) \\ &= \mathcal{N}(\mathbf{x}; \mathbf{x}_0, \mathbf{P}^{-1}) \end{aligned} \quad (12)$$

For this to hold, \mathbf{P} must be positive definite, i.e. \mathbf{x}_0 must be a maximum.

This method can be used instead to find an approximate value of the normalising constant K . Substituting Equation 11 into the definition of K :

$$\begin{aligned} K = \int f(\mathbf{x})d\mathbf{x} &\approx f(\mathbf{x}_0) \int \exp\left(\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \mathbf{P}(\mathbf{x} - \mathbf{x}_0)\right) d\mathbf{x} \\ &= \frac{(2\pi)^{\frac{N}{2}}}{\det \mathbf{P}^{\frac{1}{2}}} f(\mathbf{x}_0) \end{aligned} \quad (13)$$

A.2 Bayesian logistic regression

A.2.1 Posterior distribution

Given that the posterior Laplace approximation will be Gaussian, it is sensible to choose a Gaussian prior for the model weights:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{m}_0, \mathbf{C}_0) \quad (14)$$

In order to apply the Laplace approximation, we need the location of a maximum in the posterior. Using the results for the likelihood of the logistic model in the previous report and the log of Equation 14, the log-posterior of the model weights is:

$$\begin{aligned} \ln p(\mathbf{w}|\mathbf{X}, \mathbf{y}) &= \sum_{n=1}^N y_n \log \sigma(\mathbf{w}^T \tilde{\mathbf{x}}_n) + (1 - y_n) \log \sigma(-\mathbf{w}^T \tilde{\mathbf{x}}_n) \\ &\quad + \frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^T \mathbf{C}_0^{-1}(\mathbf{w} - \mathbf{w}_0) \\ &\quad + \text{const} \end{aligned} \quad (15)$$

The value of \mathbf{w} at the maximum, \mathbf{w}_{MAP} , can be found by setting the derivative of Equation 15 to zero. The mean of $q(\mathbf{w})$ will be set to \mathbf{w}_{MAP} .

Using Equation 12, the covariance matrix of $q(\mathbf{w})$ can be found:

$$\begin{aligned} \mathbf{C}_N^{-1} &= -\nabla^2 \ln p(\mathbf{w}|\mathbf{X}, \mathbf{y})|_{\mathbf{w}=\mathbf{w}_{\text{MAP}}} \\ &= \mathbf{C}_0^{-1} + \sum_{n=1}^N \sigma(\mathbf{w}^T \tilde{\mathbf{x}}_n) \sigma(-\mathbf{w}^T \tilde{\mathbf{x}}_n) \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T \end{aligned} \quad (16)$$

Defining $\boldsymbol{\sigma} = \sigma(\tilde{\mathbf{X}}^T \mathbf{w})$, this can be written in vector form as:

$$\mathbf{C}_N^{-1} = \mathbf{C}_0^{-1} + \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \boldsymbol{\sigma} (1 - \boldsymbol{\sigma}) \quad (17)$$

Hence, using Equation 12, the posterior distribution is:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) \approx \mathcal{N}(\mathbf{w}; \mathbf{w}_{\text{MAP}}, \mathbf{C}_N) \quad (18)$$

And using Equation 13, the normalizing constant is:

$$p(\mathbf{X}|\mathbf{y}) = \frac{(2\pi)^{\frac{N}{2}}}{\det \mathbf{C}_N^{-\frac{1}{2}}} p(\mathbf{X}|\mathbf{w} = \mathbf{w}_{\text{MAP}}, \mathbf{y}) p(\mathbf{w} = \mathbf{w}_{\text{MAP}}) \quad (19)$$

A.3 Predictive distribution

The predictive distribution can also be approximated using the Laplace method:

$$\begin{aligned} p(y^* = 1|\mathbf{x}^*, \mathbf{y}, \mathbf{X}) &= \int p(y^* = 1|\mathbf{x}^*, \mathbf{w}) p(\mathbf{w}|\mathbf{y}, \mathbf{X}) d\mathbf{w} \\ &\approx \int \sigma(\mathbf{w}^T \mathbf{x}^*) q(\mathbf{w}) d\mathbf{w} \end{aligned} \quad (20)$$

Using the sifting property of the delta function:

$$\sigma(\mathbf{w}^T \mathbf{x}) = \int \delta(a - \mathbf{w}^T \mathbf{x}) \sigma(a) da \quad (21)$$

Hence:

$$\begin{aligned} p(y^* = 1|\mathbf{x}^*, \mathbf{y}, \mathbf{X}) &\approx \int \int \delta(a - \mathbf{w}^T \mathbf{x}^*) \sigma(a) q(\mathbf{w}) d\mathbf{w} da \\ &= \int \sigma(a) \int \delta(a - \mathbf{x}^{*T} \mathbf{w}) q(\mathbf{w}) d\mathbf{w} da \end{aligned}$$

The inner integral applies a linear constraint to $q(\mathbf{w})$, as the argument of the delta function is 0 unless $a = \mathbf{x}^{*T} \mathbf{w}$. Hence, the approximate predictive distribution is:

$$\begin{aligned} p(y^* = 1|\mathbf{x}^*, \mathbf{y}, \mathbf{X}) &\approx \int \sigma(a) \mathcal{N}(a; \mathbf{x}^{*T} \mathbf{w}_{\text{MAP}}, \mathbf{x}^{*T} \mathbf{C}_N \mathbf{x}^*) da \\ &= \int \sigma(a) \mathcal{N}(a; \mu_p, \sigma_p^2) da \end{aligned} \quad (22)$$

Where

$$\mu_p = \mathbf{x}^{*T} \mathbf{w}_{\text{MAP}} \quad (23)$$

$$\sigma_p^2 = \mathbf{x}^{*T} \mathbf{C}_N \mathbf{x}^* \quad (24)$$

This integral cannot be expressed analytically, so another approximation is required. The logistic function can be approximated well by a probit function scaled such that the gradient of the two functions at the origin are equal. It can be shown that this gives:

$$\sigma(x) \approx \Phi^{-1} \left(\sqrt{\frac{\pi}{8}} x \right) = \Phi^{-1}(\lambda x) \quad (25)$$

Substituting into Equation 22 and evaluating using properties of the probit function gives:

$$\begin{aligned}
p(y^* = 1 | \mathbf{x}^*, \mathbf{y}, \mathbf{X}) &\approx \int \Phi^{-1}(\lambda x) \mathcal{N}(a; \mu_p, \sigma_p) da \\
&= \Phi^{-1} \left(\frac{\mu_p}{\sqrt{\lambda^{-2} + \sigma_p^2}} \right)
\end{aligned} \tag{26}$$

Using Equation 25, this can be converted back into a logistic function:

$$p(y^* = 1 | \mathbf{x}^*, \mathbf{y}, \mathbf{X}) \approx \sigma \left(\frac{\mu_p}{\sqrt{1 + \sigma_p^2 \lambda^2}} \right) \tag{27}$$