

Bayesian Logistic Classification

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

February 25, 2022

Abstract

Abstract goes here

1 Introduction

2 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 :

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

In many cases, the integral $\int f(\mathbf{x})d\mathbf{x}$ is intractable and 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.

Initially, consider the truncated Taylor expansion of $\ln f(\mathbf{x})$ around a local maximum \mathbf{x}_0 :

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

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

$$\begin{aligned} \ln f(\mathbf{x}) &\approx \ln f(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \nabla^2 \ln f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}(\mathbf{x} - \mathbf{x}_0) \\ f(\mathbf{x}) &\approx f(\mathbf{x}_0) \exp\left(\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \nabla^2 \ln f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}(\mathbf{x} - \mathbf{x}_0)\right) \end{aligned} \quad (3)$$

Equation 3 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 3 to get an approximation for $p(\mathbf{x})$:

$$\begin{aligned} p(\mathbf{x}) &\approx \frac{1}{(2\pi)^{\frac{N}{2}} \det \mathbf{P}^{-\frac{1}{2}}} \exp\left(\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \mathbf{P}(\mathbf{x} - \mathbf{x}_0)\right) \\ &= \mathcal{N}(\mathbf{x}; \mathbf{x}_0, \mathbf{P}^{-1}) \end{aligned} \quad (4)$$

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 3 into the definition of K :

$$\begin{aligned} K &= \int f(\mathbf{x})d\mathbf{x} \approx f(\mathbf{x}_0) \int \exp\left(\frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T \mathbf{P}(\mathbf{x} - \mathbf{x}_0)\right) d\mathbf{x} \\ &= \frac{(2\pi)^{\frac{N}{2}}}{\det \mathbf{P}^{\frac{1}{2}}} f(\mathbf{x}_0) \end{aligned} \quad (5)$$

3 Recovering Bayesian logistic regression

3.1 Posterior distribution

In the previous report, the maximum-likelihood estimate for the model weights was used in classification¹. It is desirable to have a full posterior distribution of the weights:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{X}|\mathbf{w}, \mathbf{y})p(\mathbf{w})}{p(\mathbf{X}|\mathbf{y})} \quad (6)$$

The problem is that calculating Equation 6 requires normalising a product of logistic functions, which is difficult. The Laplace approximation can be used to calculate the normalising constant $p(\mathbf{X}|\mathbf{y})$, also called the model evidence, and calculate an approximate posterior distribution.

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

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{m}_0, \mathbf{C}_0) \quad (7)$$

In order to apply the Laplace approximation, we need the location of a maximum in the posterior. Using the results for the log-likelihood of the logistic model and the log of Equation 7, the log-posterior is:

$$\begin{aligned} \ln p(\mathbf{w}|\mathbf{X}, \mathbf{y}) &= \sum_{n=1}^N y_n \log \sigma(\mathbf{w}^T \tilde{\mathbf{x}}_n) + (1 - y_n) \log \sigma(-\mathbf{w}^T \tilde{\mathbf{x}}_n) \\ &\quad + \frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^T \mathbf{C}_0^{-1}(\mathbf{w} - \mathbf{w}_0) \\ &\quad + \text{const} \end{aligned} \quad (8)$$

The value of \mathbf{w} at the maximum, \mathbf{w}_{MAP} , can be found by setting the derivative of Equation 8 to zero. As outlined in Section 2 \mathbf{w}_{MAP} will be the mean of $q(\mathbf{w})$.

Using Equation 4, the covariance matrix of $q(\mathbf{w})$ can be found:

$$\begin{aligned} \mathbf{C}_N^{-1} &= -\nabla^2 \ln p(\mathbf{w}|\mathbf{X}, \mathbf{y})|_{\mathbf{w}=\mathbf{w}_{\text{MAP}}} \\ &= \mathbf{C}_0^{-1} + \sum_{n=1}^N \sigma(\mathbf{w}^T \tilde{\mathbf{x}}_n) \sigma(-\mathbf{w}^T \tilde{\mathbf{x}}_n) \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T \end{aligned} \quad (9)$$

Hence, using Equation 4 and Equation 5:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) \approx \mathcal{N}(\mathbf{w}; \mathbf{w}_{\text{MAP}}, \mathbf{C}_N) \quad (10)$$

$$p(\mathbf{X}|\mathbf{y}) = \frac{(2\pi)^{\frac{N}{2}}}{\det \mathbf{C}_N^{-\frac{1}{2}}} p(\mathbf{X}|\mathbf{w} = \mathbf{w}_{\text{MAP}}, \mathbf{y}) p(\mathbf{w} = \mathbf{w}_{\text{MAP}}) \quad (11)$$

3.2 Predictive distribution

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