

# Laplace Approximation for Bayesian Logistic Regression: A Derivation

Wei Dai

March 13, 2013

This derivation is based on Bishop (2006). Given data set  $\{\phi_n, t_n\}_{n=1}^N$  where  $\phi_n$  are the feature vectors and  $t_n \in \{0, 1\}$  are the labels, we can write the likelihood function for logistic regression as

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n} \quad (1)$$

where  $\mathbf{t} = (t_1, \dots, t_N)^T$  and  $y_n = p(\mathcal{C}_1|\phi_n) = \sigma(\mathbf{w}^T \phi_n)$  and  $\sigma(s) = \frac{1}{1+e^{-s}}$ . Using Bayes rule, the posterior distribution over  $\mathbf{w}$  is

$$p(\mathbf{w}|\mathbf{t}) = \frac{p(\mathbf{w})p(\mathbf{t}|\mathbf{w})}{p(\mathbf{t})} \quad (2)$$

where  $p(\mathbf{t}) = \int p(\mathbf{w})p(\mathbf{t}|\mathbf{w})d\mathbf{w}$  involves logistic sigmoid functions and is intractable. Laplace approximation approximate this posterior with a multivariate Gaussian:

$$\begin{aligned} q(\mathbf{w}) &= \frac{1}{(2\pi)^{M/2}|\mathbf{S}_N|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{w} - \mathbf{w}_{MAP})^T \mathbf{S}_N^{-1}(\mathbf{w} - \mathbf{w}_{MAP}) \right\} \\ &= \mathcal{N}(\mathbf{w}; \mathbf{w}_{MAP}, \mathbf{S}_N) \end{aligned} \quad (3)$$

where  $\mathbf{w}_{MAP}$  is the maximum *a posteriori* and thus a mode of the posterior and  $\mathbf{S}_N^{-1} = -\nabla_{\mathbf{w}}^2 \ln p(\mathbf{w}|\mathbf{t})|_{\mathbf{w}=\mathbf{w}_{MAP}}$  is the Hessian at  $\mathbf{w}_{MAP}$ . Since we are approximating the posterior with Gaussian, it is convenient to use conjugate prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{m}_0, \mathbf{S}_0)$ . Thus we have

$$\ln p(\mathbf{w}|\mathbf{t}) = -\frac{1}{2}(\mathbf{w}-\mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w}-\mathbf{m}_0) + \sum_{n=1}^N \{t_n \ln y_n + (1-y_n) \ln(1-y_n)\} + const \quad (4)$$

Since eq. 4 is convex in  $\mathbf{w}$  (?), we can use gradient descent to find  $\mathbf{w}_{MAP} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{t}) = \arg \max_{\mathbf{w}} \ln p(\mathbf{w}|\mathbf{t})$ . Using the following facts:

$$\nabla_{\mathbf{x}} \mathbf{m}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{x} \quad (5)$$

$$\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{m} = \mathbf{A} \mathbf{x} \quad (6)$$

$$\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{A} \mathbf{x} + \mathbf{A}^T \mathbf{x} \quad (7)$$

$$\nabla_{\mathbf{x}} \mathbf{A} \mathbf{x} = \mathbf{A} \quad (8)$$

$$\frac{\partial}{\partial s} \sigma(s) = \sigma(s)(1 - \sigma(s)) \quad (9)$$

where  $\mathbf{x}, \mathbf{m} \in \mathbb{R}^d$ ,  $\mathbf{A} \in \mathbb{R}^{d \times d}$ . Also note that  $\mathbf{S}_0^{-1} = \mathbf{S}_0^{-T}$ , we arrive at

$$\nabla_{\mathbf{w}} \ln p(\mathbf{w}|\mathbf{t}) = -\mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) + \sum_{n=1}^N (t_n - y_n) \phi_n \quad (10)$$

Since we generally want small  $\|\mathbf{w}\|$  to avoid overfitting, we use  $\mathbf{m}_0 = \mathbf{0}$  and  $\mathbf{S}_0 = \sigma^2 \mathbf{I}$ . We have the following update rule:

$$\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} + \eta \left( \sum_{n=1}^N (t_n - y_{n,(t-1)}) \phi_n - \frac{1}{\sigma^2} \mathbf{w}_{t-1} \right) \quad (11)$$

where  $\eta$  is the learning rate constant. We can also get

$$\mathbf{S}_N^{-1} = -\nabla_{\mathbf{w}}^2 \ln p(\mathbf{w}|\mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n(1 - y_n) \phi_n \phi_n^T \quad (12)$$

Thus, we have the approximated posterior  $q(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{w}_{MAP}, \mathbf{S}_N)$ .