

Logistic regression

Sep 2022

Murphy chap8, Bishop chap4

Maximum entropy property

$$-\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K P(y_i = k|x_i) \log P(y_i = k|x_i).$$

The fixed-mean constraints,

$$\sum_{i=1}^N P(y_i = k|x_i) x_{ij} = \sum_{i=1}^N \mu_k(y_i) x_{ij}, \quad \forall k = 1, \dots, K, \quad j = 1, \dots, p.$$

The normalization constraint,

$$\sum_{k=1}^K P(y_i = k|x_i) = 1, \quad i = 1, \dots, N.$$

Maximum entropy property

6.3. Binary logistic regression. For binary classification, we encode the class labels using $\{0, 1\}$. Note $P(y = 1|x) = 1 - P(y = 0|x)$. We may write

$$\ln P(y = 0|x) = \beta_0^T x - \ln Z,$$

$$\ln(1 - P(y = 0|x)) = \beta_1^T x - \ln Z.$$

Since we have incorporated the normalization condition, the coefficients β_0^T, β_1^T must be compatible in some sense. One way to include the compatibility is to write the model as (by subtracting the second equation from the first),

$$\ln \frac{P(y = 0|x)}{1 - P(y = 0|x)} = \beta^T x, \quad \beta \in \mathbb{R}^p.$$

The LHS as a function of $P(y = 0|x)$ is called the logit or log odds function

$$\text{logit}(p) = \ln \frac{p}{1 - p},$$

which is the inverse of the logistic function

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

Therefore

$$P(y = 0|x) = \sigma(\beta^T x) = \frac{1}{1 + e^{-\beta^T x}}.$$

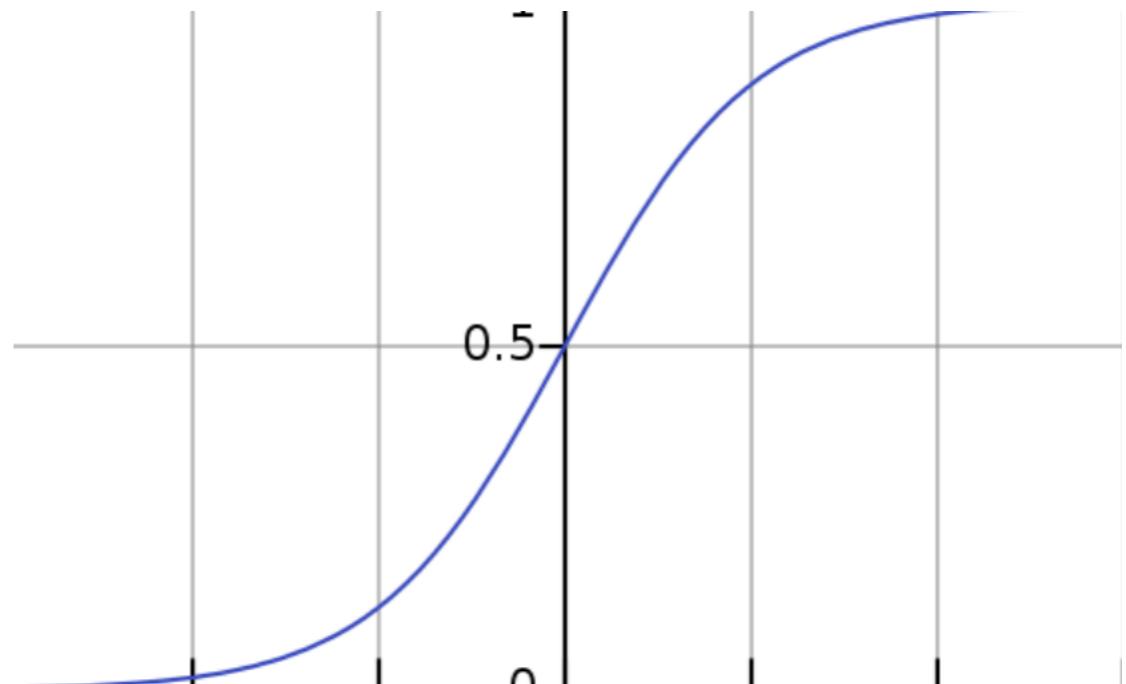
Logistic function

For later reference: the logistic function is

$$\sigma(t) = \frac{1}{1 + e^{-t}} = \frac{1}{2} + \frac{1}{2} \tanh\left(\frac{t}{2}\right)$$

$$\tanh(t) = \frac{e^t - e^{-t}}{e^t + e^{-t}} = \frac{e^{2t} - 1}{e^{2t} + 1}$$

$$\frac{\partial \sigma}{\partial t} = \sigma(t)(1 - \sigma(t))$$



Model selection

Entropy, AIC, BIC?

AIC:

Historically various ‘information criteria’ have been proposed that attempt to correct for the bias of maximum likelihood by the addition of a penalty term to compensate for the over-fitting of more complex models. For example, the *Akaike information criterion*, or AIC (Akaike, 1974), chooses the model for which the quantity

$$\ln p(\mathcal{D}|\mathbf{w}_{\text{ML}}) - M \tag{1.73}$$

is largest. Here $p(\mathcal{D}|\mathbf{w}_{\text{ML}})$ is the best-fit log likelihood, and M is the number of adjustable parameters in the model. A variant of this quantity, called the *Bayesian*

Laplace Approximation

Here we introduce a simple, but widely used, framework called the Laplace approximation, that aims to find a Gaussian approximation to a probability density defined over a set of continuous variables. Consider first the case of a single continuous variable z , and suppose the distribution $p(z)$ is defined by

$$p(z) = \frac{1}{Z} f(z) \tag{4.125}$$

where $Z = \int f(z) dz$ is the normalization coefficient. We shall suppose that the value of Z is unknown. In the Laplace method the goal is to find a Gaussian approximation $q(z)$ which is centred on a mode of the distribution $p(z)$. The first step is to find a mode of $p(z)$, in other words a point z_0 such that $p'(z_0) = 0$, or equivalently

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0. \tag{4.126}$$

Laplace Approximation

z_0 so that

$$\ln f(z) \simeq \ln f(z_0) - \frac{1}{2} A(z - z_0)^2 \quad (4.127)$$

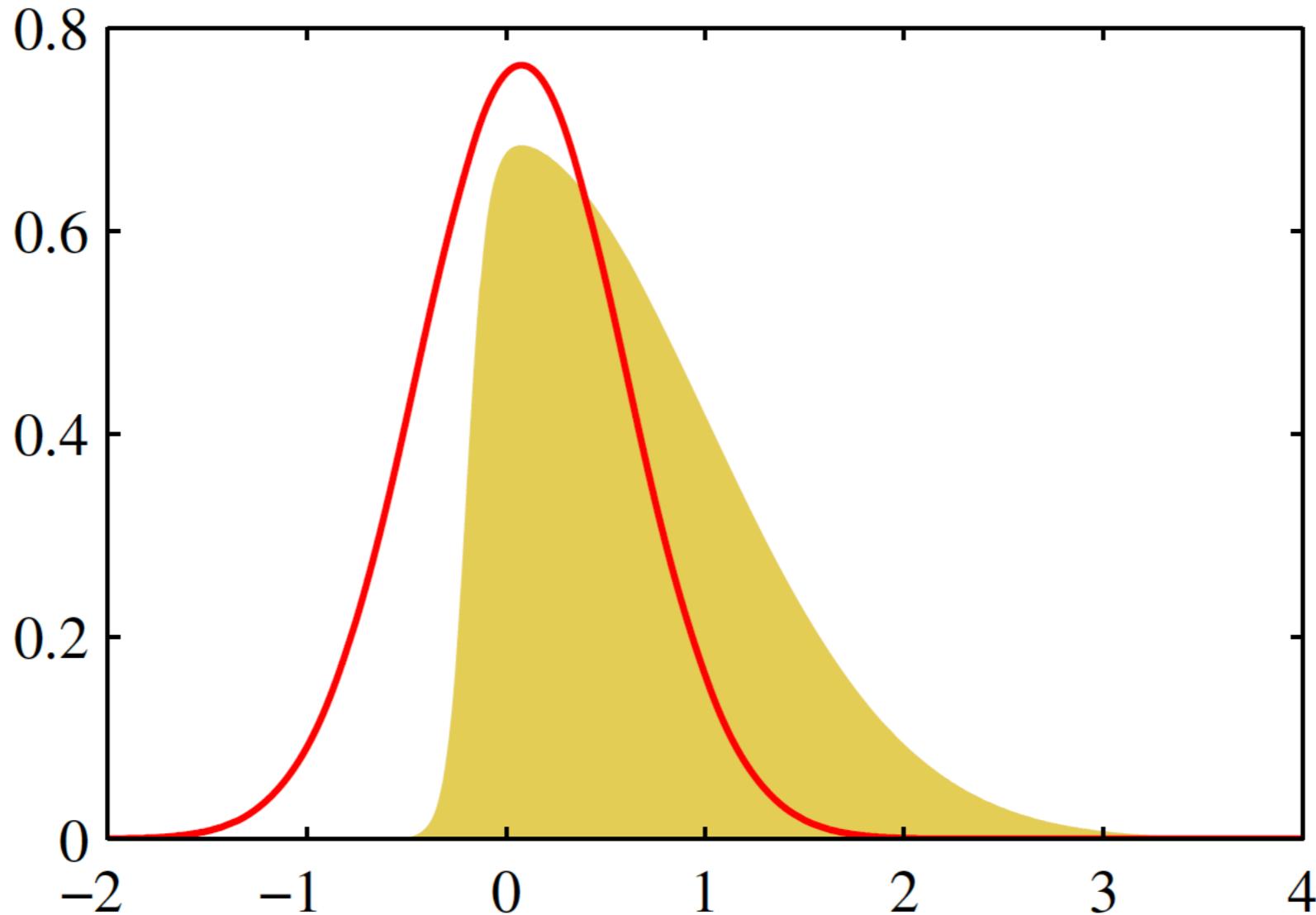
$$A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}. \quad (4.128)$$

$$f(z) \simeq f(z_0) \exp \left\{ -\frac{A}{2}(z - z_0)^2 \right\}. \quad (4.129)$$

After normalization, the approximation $q(x)$ of $p(x)$ is:

$$q(z) = \left(\frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2}(z - z_0)^2 \right\}. \quad (4.130)$$

Laplace Approximation



$$p(z) \propto \exp(-z^2/2)\sigma(20z + 4)$$

where $\sigma(z)$ is the logistic sigmoid function defined by $\sigma(z) = (1 + e^{-z})^{-1}$

Laplace Approximation

We can extend the Laplace method to approximate a distribution $p(\mathbf{z}) = f(\mathbf{z})/Z$ defined over an M -dimensional space \mathbf{z} . At a stationary point \mathbf{z}_0 the gradient $\nabla f(\mathbf{z})$ will vanish. Expanding around this stationary point we have

$$\ln f(\mathbf{z}) \simeq \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \quad (4.131)$$

where the $M \times M$ Hessian matrix \mathbf{A} is defined by

$$\mathbf{A} = -\nabla\nabla \ln f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0} \quad (4.132)$$

and ∇ is the gradient operator. Taking the exponential of both sides we obtain

$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\}. \quad (4.133)$$

Laplace Approximation

From (4.133):

$$\begin{aligned} Z &= \int f(\mathbf{z}) d\mathbf{z} \\ &\simeq f(\mathbf{z}_0) \int \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{z}_0)^T \mathbf{A} (\mathbf{z} - \mathbf{z}_0) \right\} d\mathbf{z} \\ &= f(\mathbf{z}_0) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}} \end{aligned} \tag{4.135}$$

Use (4.135) to approximate model evidence !

BIC

Consider a data set \mathcal{D} and a set of models $\{\mathcal{M}_i\}$ having parameters $\{\boldsymbol{\theta}_i\}$. For each model we define a likelihood function $p(\mathcal{D}|\boldsymbol{\theta}_i, \mathcal{M}_i)$. If we introduce a prior $p(\boldsymbol{\theta}_i|\mathcal{M}_i)$ over the parameters, then we are interested in computing the model evidence $p(\mathcal{D}|\mathcal{M}_i)$ for the various models. From now on we omit the conditioning on \mathcal{M}_i to keep the notation uncluttered. From Bayes' theorem the model evidence is given by

$$p(\mathcal{D}) = \int p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (4.136)$$

Identifying $f(\boldsymbol{\theta}) = p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$ and $Z = p(\mathcal{D})$, and applying the result (4.135), we obtain

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) + \underbrace{\ln p(\boldsymbol{\theta}_{\text{MAP}}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}|}_{\text{Occam factor}} \quad (4.137)$$

where $\boldsymbol{\theta}_{\text{MAP}}$ is the value of $\boldsymbol{\theta}$ at the mode of the posterior distribution, and \mathbf{A} is the *Hessian* matrix of second derivatives of the negative log posterior

$$\mathbf{A} = -\nabla\nabla \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}})p(\boldsymbol{\theta}_{\text{MAP}}) = -\nabla\nabla \ln p(\boldsymbol{\theta}_{\text{MAP}}|\mathcal{D}). \quad (4.138)$$

BIC

If we assume that the Gaussian prior distribution over parameters is broad, and that the Hessian has full rank, then we can approximate (4.137) very roughly using

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) - \frac{1}{2}M \ln N \quad (4.139)$$

where N is the number of data points, M is the number of parameters in $\boldsymbol{\theta}$ and we have omitted additive constants. This is known as the *Bayesian Information Criterion* (BIC) or the *Schwarz criterion* (Schwarz, 1978). Note that, compared to AIC given by (1.73), this penalizes model complexity more heavily.

Bayesian logistic regression

4.5. Bayesian Logistic Regression

We now turn to a Bayesian treatment of logistic regression.

begin with a Gaussian prior, which we write in the general form

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) \quad (4.140)$$

where \mathbf{m}_0 and \mathbf{S}_0 are fixed hyperparameters. The posterior distribution over \mathbf{w} is given by

$$p(\mathbf{w} | \mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t} | \mathbf{w}) \quad (4.141)$$

Bayesian logistic regression

where $\mathbf{t} = (t_1, \dots, t_N)^T$. Taking the log of both sides, and substituting for the prior distribution using (4.140), and for the likelihood function using (4.89), we obtain

$$\begin{aligned}\ln p(\mathbf{w}|\mathbf{t}) &= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) \\ &\quad + \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} + \text{const}\end{aligned}\quad (4.142)$$

where $y_n = \sigma(\mathbf{w}^T \boldsymbol{\phi}_n)$. To obtain a Gaussian approximation to the posterior distribution, we first maximize the posterior distribution to give the MAP (maximum posterior) solution \mathbf{w}_{MAP} , which defines the mean of the Gaussian. The covariance is then given by the inverse of the matrix of second derivatives of the negative log likelihood, which takes the form

$$\mathbf{S}_N = -\nabla \nabla \ln p(\mathbf{w}|\mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n(1 - y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T. \quad (4.143)$$

[next page](#)

Bayesian logistic regression

$$\begin{aligned} & \frac{\partial(tlny + (1 - t)ln(1 - y))}{\partial w} \\ &= t \frac{y(1 - y)}{y} \phi + (1 - t) \frac{(-y)(1 - y)}{1 - y} \phi \\ &= t\phi - ty\phi + ty\phi - y\phi \\ &= (t - y)\phi \end{aligned}$$

$$\begin{aligned} & \frac{\partial^2(tlny + (1 - t)ln(1 - y))}{\partial w^2} \\ &= -y(1 - y)\phi\phi^T \end{aligned}$$

Bayesian logistic regression

The Gaussian approximation to the posterior distribution therefore takes the form

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{\text{MAP}}, \mathbf{S}_N). \quad (4.144)$$

Bayesian logistic regression

Having obtained a Gaussian approximation to the posterior distribution, there remains the task of marginalizing with respect to this distribution in order to make predictions.

4.5.2 Predictive distribution

The predictive distribution for class \mathcal{C}_1 , given a new feature vector $\phi(\mathbf{x})$, is

$$p(\mathcal{C}_1|\phi, \mathbf{t}) = \int p(\mathcal{C}_1|\phi, \mathbf{w})p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \simeq \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} \quad (4.145)$$

Denoting $a = \mathbf{w}^T \phi$, we have

$$\sigma(\mathbf{w}^T \phi) = \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da \quad (4.146)$$

where $\delta(\cdot)$ is the Dirac delta function. From this we obtain

$$\int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} = \int \sigma(a) p(a) da \quad (4.147)$$

Bayesian logistic regression

$$\sigma(\mathbf{w}^T \boldsymbol{\phi}) = \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}) \sigma(a) da \quad (4.146)$$

where $\delta(\cdot)$ is the Dirac delta function. From this we obtain

$$\int \sigma(\mathbf{w}^T \boldsymbol{\phi}) q(\mathbf{w}) d\mathbf{w} = \int \sigma(a) p(a) da \quad (4.147)$$

where

$$p(a) = \int \delta(a - \mathbf{w}^T \boldsymbol{\phi}) q(\mathbf{w}) d\mathbf{w}. \quad (4.148)$$

p(a) still a Gaussian distribution !

Bayesian logistic regression

$$a = \mathbf{w}^T \boldsymbol{\phi},$$

$$\mu_a = \mathbb{E}[a] = \int p(a) a \, da = \int q(\mathbf{w}) \mathbf{w}^T \boldsymbol{\phi} \, d\mathbf{w} = \mathbf{w}_{\text{MAP}}^T \boldsymbol{\phi} \quad (4.149)$$

where we have used the result (4.144) for the variational posterior distribution $q(\mathbf{w})$. Similarly

$$\begin{aligned} \sigma_a^2 &= \text{var}[a] = \int p(a) \{a^2 - \mathbb{E}[a]^2\} \, da \\ &= \int q(\mathbf{w}) \{(\mathbf{w}^T \boldsymbol{\phi})^2 - (\mathbf{m}_N^T \boldsymbol{\phi})^2\} \, d\mathbf{w} = \boldsymbol{\phi}^T \mathbf{S}_N \boldsymbol{\phi}. \end{aligned} \quad (4.150)$$

$$p(\mathcal{C}_1 | \mathbf{t}) = \int \sigma(a) p(a) \, da = \int \sigma(a) \mathcal{N}(a | \mu_a, \sigma_a^2) \, da. \quad (4.151)$$

See 8.4.4.2 for the approximation of the integral

Monte Carlo Approximation

focusses on prediction. The posterior predictive distribution has the form

$$p(y|\mathbf{x}, \mathcal{D}) = \int p(y|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w} \quad (8.59)$$

Unfortunately this integral is intractable.

The simplest approximation is the plug-in approximation, which, in the binary case, takes the form

$$p(y = 1|\mathbf{x}, \mathcal{D}) \approx p(y = 1|\mathbf{x}, \mathbb{E}[\mathbf{w}]) \quad (8.60)$$

where $\mathbb{E}[\mathbf{w}]$ is the posterior mean. In this context, $\mathbb{E}[\mathbf{w}]$ is called the **Bayes point**. Of course, such a plug-in estimate underestimates the uncertainty. We discuss some better approximations below.

Monte Carlo Approximation

Monte Carlo approximation

A better approach is to use a Monte Carlo approximation, as follows:

$$p(y = 1 | \mathbf{x}, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S \text{sigm}((\mathbf{w}^s)^T \mathbf{x}) \quad (8.61)$$

where $\mathbf{w}^s \sim p(\mathbf{w} | \mathcal{D})$ are samples from the posterior. (This technique can be trivially extended to the multi-class case.) If we have approximated the posterior using Monte Carlo, we can reuse these samples for prediction. If we made a Gaussian approximation to the posterior, we can draw *independent* samples from the Gaussian using standard methods.