

4.3.5 Probit regression

Remain within the framework of generalized linear models (two-class classification)

$$p(t = 1 | a) = f(a)$$

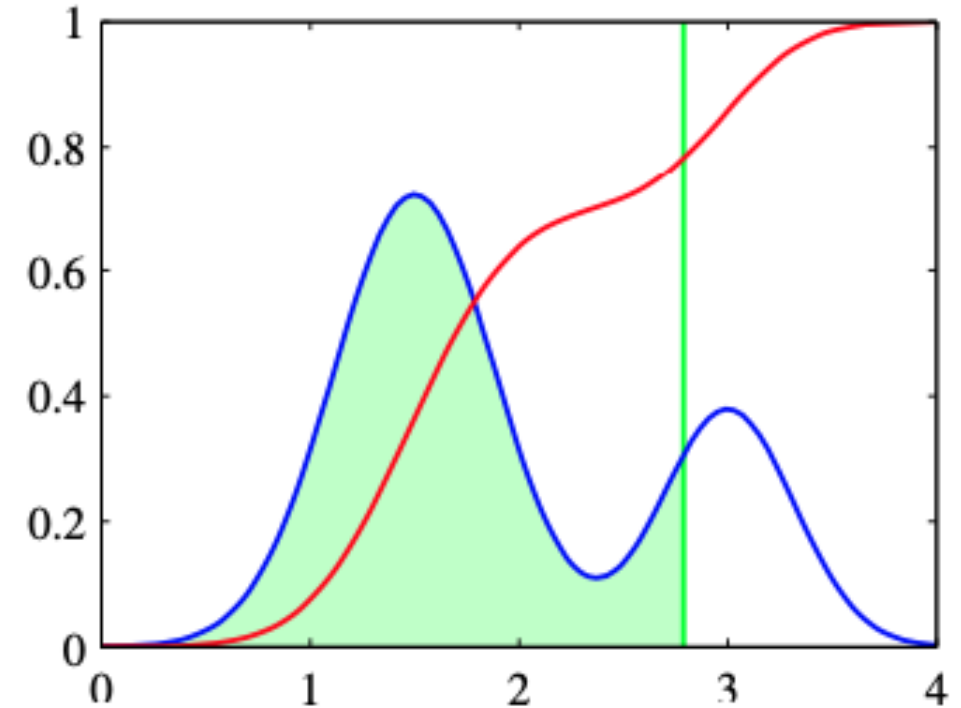
$a = \mathbf{w}^T \phi$ and $f(\cdot)$ is activation function.

$$\begin{cases} t_n = 1 & \text{if } a_n \geq \theta \\ t_n = 0 & \text{otherwise.} \end{cases}$$

■ $p(\theta)$: probability density

■ $\Phi(a)$: cumulative distribution function

- $\Phi(a) = \int_{-\infty}^a p(\theta) d\theta$
- when $p(\theta) \sim \mathcal{N}(0, 1) \Rightarrow$ probit function
- cumulative distribution function is equivalent to activation function.



Exercise 4.21

Show that the probit function (4.114) and the erf function (4.115) are related by (4.116).

$$\begin{aligned}\Phi(a) &= \int_{-\infty}^a \mathcal{N}(\theta \mid 0, 1) d\theta \\ &= \frac{1}{2} + \int_0^a \mathcal{N}(\theta \mid 0, 1) d\theta \\ &= \frac{1}{2} + \frac{1}{\sqrt{2\pi}} \int_0^a \exp(-\theta^2/2) d\theta\end{aligned}$$

Replace $\theta = \sqrt{2}\theta'$

$$\begin{aligned}&= \frac{1}{2} + \frac{1}{\sqrt{\pi}} \int_0^{a/\sqrt{2}} \exp(-\theta'^2) d\theta' \\ &= \frac{1}{2} \left\{ 1 + \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \right\}\end{aligned}$$

$$\operatorname{erf}(a) = \frac{2}{\sqrt{\pi}} \int_0^a \exp(-\theta^2) d\theta$$

4.3.5 Probit regression

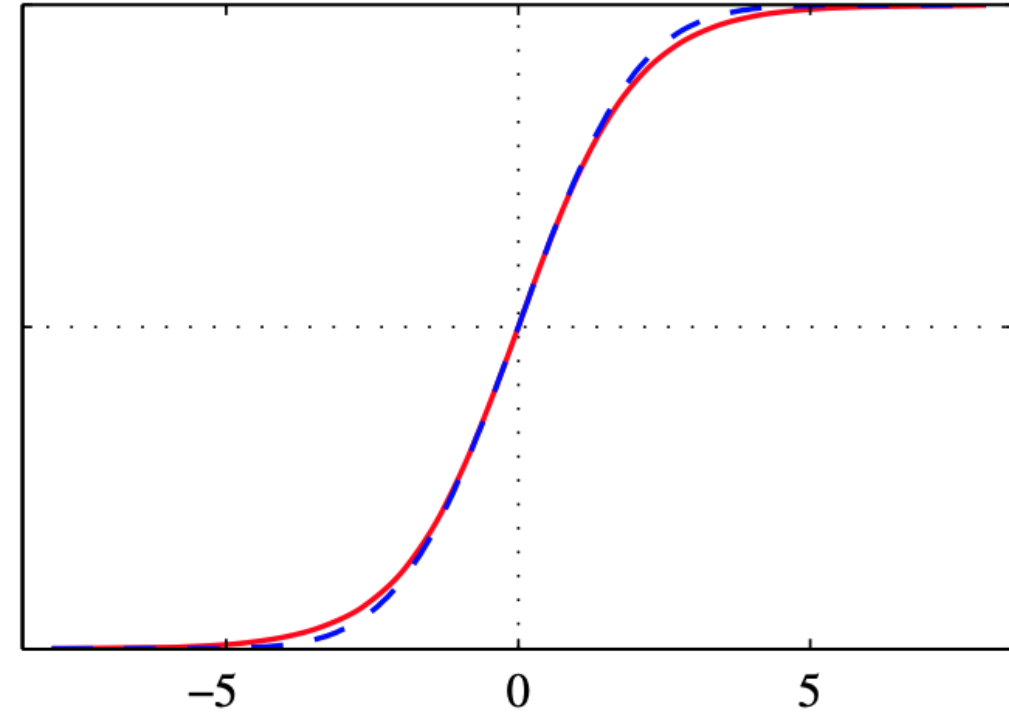
■ $\sigma(a)$: logistic function

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

■ $\Phi(\lambda a)$: probit function (scaling factor $\lambda = \frac{\pi^2}{8}$)

$$\Phi(a) = \frac{1}{2} \left\{ 1 + \frac{1}{\sqrt{2}} \operatorname{erf}(a) \right\}$$

- logistic regression: decay like $\exp(-x)$
- probit regression: decay like $\exp(-x^2)$
 - => more sensitive to outlier



4.3.6 Canonical link functions

Assumption of exponential family distribution to the target variable t

$$p(t \mid \eta, s) = \frac{1}{s} h\left(\frac{t}{s}\right) g(\eta) \exp\left\{\frac{\eta t}{s}\right\}$$

Log likelihood

$$\ln p(\mathbf{t} \mid \eta, s) = \sum_{n=1}^N \ln p(t_n \mid \eta, s) = \sum_{n=1}^N \left\{ \ln g(\eta_n) + \frac{\eta_n t_n}{s} \right\} + \text{const}$$

Derivative of the log likelihood with respect to the model parameters \mathbf{w}

$$\begin{aligned} \nabla_{\mathbf{w}} \ln p(\mathbf{t} \mid \eta, s) &= \sum_{n=1}^N \left\{ \frac{d}{d\eta_n} \ln g(\eta_n) + \frac{t_n}{s} \right\} \frac{d\eta_n}{dy_n} \frac{dy_n}{da_n} \nabla a_n \\ &= \sum_{n=1}^N \frac{1}{s} \{t_n - y_n\} \psi'(y_n) f'(a_n) \phi_n \end{aligned}$$

what is canonical link functions $y = f(\mathbf{w}^T \phi)$, $f^{-1}(y) = \psi(y)$

$$\nabla \ln E(\mathbf{w}) = \frac{1}{s} \sum_{n=1}^N \{y_n - t_n\} \phi_n$$

4.4 The Laplace Approximation

Approximate the posterior distribution of logistic regression with Gaussian $p(z) = \frac{f(z)}{Z}$
(Z = normalize coefficient)

=> Integrate the posterior probability with parameter w , we can obtain the predictive distribution(discussed in section(3.3))

Taylor expansion of $\ln f(z)$ centred on the mode z_0

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

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

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

Normalized distribution

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

4.4 The Laplace Approximation

Advantages

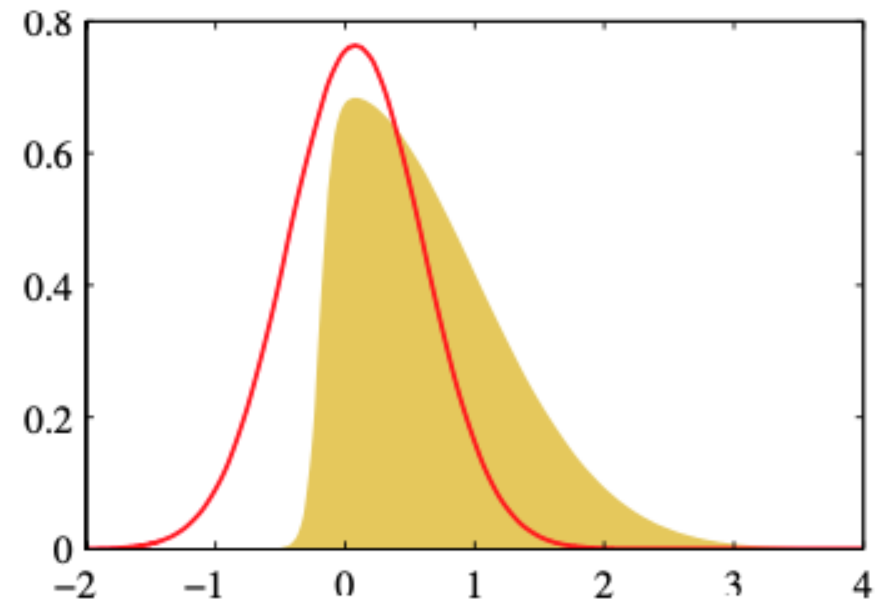
- Simple to calculate
- Approximation by Gaussian distribution will be more accurate if there is more data (the central limit theorem)

Disadvantages

- Only applicable to real variables
- Difficult to choose which mode to use when multimodal distribution

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

■ : Laplace approximation of $p(z)$



4.4.1 Model comparison and BIC

Using laplace approximation to the normalization constant Z

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

Consider a data set D and a set of models $\{M_i\}$ having parameters $\{\theta_i\}$
Model evidence $p(D | M_i)$ (omit M_i)

$$p(D) = \int p(D | \theta) p(\theta) d\theta$$

$$\ln p(D) \simeq \ln p(D | \theta_{\text{MAP}}) + \ln p(\theta_{\text{MAP}}) + \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}|$$

$$\mathbf{A} = -\nabla \nabla \ln p(D | \theta_{\text{MAP}}) p(\theta_{\text{MAP}}) = -\nabla \nabla \ln p(\theta_{\text{MAP}} | D)$$

Assume that the Gaussian prior distribution over parameters is broad, and that the Hessian has full rank

BIC(Bayesian Information Criterion)

$$\ln p(D) \simeq \ln p(D | \theta_{\text{MAP}}) - \frac{1}{2} M \ln N$$

- BIC penalizes model complexity more heavily than AIC($\ln p(D | \theta_{\text{MAP}}) - M$)

Exercise 4.22

Using the result (4.135), derive the expression (4.137) for the log model evidence under the Laplace approximation.

$$\begin{aligned} p(D) &= \int p(D \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \int f(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= f(\boldsymbol{\theta}_{MAP}) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}} \\ &= p(D \mid \boldsymbol{\theta}_{MAP}) p(\boldsymbol{\theta}_{MAP}) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}} \end{aligned}$$

$$\ln p(D) = \ln p(D \mid \boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) + \frac{M}{2} \ln 2\pi - \frac{1}{2} \ln |\mathbf{A}|$$

- M : Dimension of parameter θ
- θ_{MAP} : mode of $f(\theta)$

Exercise 4.23

4.23 (★★) **www** In this exercise, we derive the BIC result (4.139) starting from the Laplace approximation to the model evidence given by (4.137). Show that if the prior over parameters is Gaussian of the form $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \mathbf{V}_0)$, the log model evidence under the Laplace approximation takes the form

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) - \frac{1}{2}(\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m})^T \mathbf{V}_0^{-1}(\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m}) - \frac{1}{2} \ln |\mathbf{H}| + \text{const}$$

where \mathbf{H} is the matrix of second derivatives of the log likelihood $\ln p(\mathcal{D}|\boldsymbol{\theta})$ evaluated at $\boldsymbol{\theta}_{\text{MAP}}$. Now assume that the prior is broad so that \mathbf{V}_0^{-1} is small and the second term on the right-hand side above can be neglected. Furthermore, consider the case of independent, identically distributed data so that \mathbf{H} is the sum of terms one for each data point. Show that the log model evidence can then be written approximately in the form of the BIC expression (4.139).

BIC expression (4.139)

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

$$\begin{aligned}
\ln p(D) &\simeq \ln p(D \mid \boldsymbol{\theta}_{MAP}) + \ln p(\boldsymbol{\theta}_{MAP}) + \frac{M}{2} \ln 2\pi - \frac{1}{2} \ln |\mathbf{A}| \\
&= \ln p(D \mid \boldsymbol{\theta}_{MAP}) - \frac{M}{2} \ln 2\pi - \frac{1}{2} \ln |\mathbf{V}_0| - \frac{1}{2} (\boldsymbol{\theta}_{MAP} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta}_{MAP} - \mathbf{m}) \\
&\quad + \frac{M}{2} \ln 2\pi - \frac{1}{2} \ln |\mathbf{A}| \\
&= \ln p(D \mid \boldsymbol{\theta}_{MAP}) - \frac{1}{2} \ln |\mathbf{V}_0| - \frac{1}{2} (\boldsymbol{\theta}_{MAP} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta}_{MAP} - \mathbf{m}) - \frac{1}{2} \ln |\mathbf{A}|
\end{aligned}$$

$$\begin{aligned}
\mathbf{A} &= -\nabla \nabla \ln p(D \mid \boldsymbol{\theta}) p(\boldsymbol{\theta})|_{\theta=\theta_{MAP}} \\
&= -\nabla \nabla \ln p(D \mid \boldsymbol{\theta})|_{\theta=\theta_{MAP}} - \nabla \nabla \ln p(\boldsymbol{\theta})|_{\theta=\theta_{MAP}} \\
&= \mathbf{H} - \nabla \nabla \left\{ -\frac{1}{2} (\boldsymbol{\theta} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta} - \mathbf{m}) \right\}|_{\theta=\theta_{MAP}} & \ln \{ |\mathbf{V}_0| \cdot |\mathbf{H} + \mathbf{V}_0^{-1}| \} \\
&= \mathbf{H} + \nabla \{ \mathbf{V}_0^{-1} (\boldsymbol{\theta} - \mathbf{m}) \}|_{\theta=\theta_{MAP}} & = \ln \{ |\mathbf{V}_0 \mathbf{H} + \mathbf{I}| \} \\
&= \mathbf{H} + \mathbf{V}_0^{-1} & = \ln |\mathbf{V}_0| + \ln |\mathbf{H}|
\end{aligned}$$

$$\mathbf{H} = \sum_{n=1}^N \mathbf{H}_n = N \hat{\mathbf{H}}$$

$$\begin{aligned}
\ln p(D) &= \ln p(D \mid \boldsymbol{\theta}_{MAP}) - \frac{1}{2} (\boldsymbol{\theta}_{MAP} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta}_{MAP} - \mathbf{m}) - \frac{1}{2} \ln |\mathbf{H}| - \ln |\mathbf{V}_0| \\
&= \ln p(D \mid \boldsymbol{\theta}_{MAP}) - \frac{1}{2} (\boldsymbol{\theta}_{MAP} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta}_{MAP} - \mathbf{m}) - \frac{1}{2} \ln |N \hat{\mathbf{H}}| - \ln |\mathbf{V}_0| \\
&= \ln p(D \mid \boldsymbol{\theta}_{MAP}) - \frac{1}{2} (\boldsymbol{\theta}_{MAP} - \mathbf{m})^T \mathbf{V}_0^{-1} (\boldsymbol{\theta}_{MAP} - \mathbf{m}) - \frac{M}{2} \ln N - \frac{1}{2} \ln |\hat{\mathbf{H}}| - \ln |\mathbf{V}_0|
\end{aligned}$$

$$N \gg 1$$

$$\approx \ln p(D \mid \boldsymbol{\theta}_{MAP}) - \frac{M}{2} \ln N$$