

Bayesian Inference

Shantanu Kodgirwar

February 13, 2021

1 Gaussian Distribution

In the case of a single variable x , the Gaussian/normal distribution can be written in the form

$$\mathcal{N}(x | \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\} \quad (1)$$

where, μ is the mean and σ^2 is the variance.

1.1 Maximum likelihood estimation

Input vectors are given as $\mathbf{x} = (x_1, \dots, x_N)^T$ and the output/target variables as $\mathbf{t} = (t_1, \dots, t_N)^T$ and the polynomial coefficients as $\mathbf{w} = (w_1, \dots, w_M)^T$.

$$t_i = \sum_{j=1}^N y(x_i, \mathbf{w}) \quad (2)$$

$$= \sum_{i=1}^N \sum_{k=1}^M y_k(x_i) w_k \quad (3)$$

$$= (\mathbf{X}\mathbf{w})_i + n_i \quad (4)$$

Assumed distribution of n_i :

$$n_i \sim \mathcal{N}(0, \sigma^2) \quad (5)$$

The values of t , given the values of x follows a Gaussian distribution.

$$t_i \sim \mathcal{N}(y(x_i, \mathbf{w}), \sigma^2) \quad (6)$$

The likelihood function is defined as follows.

$$p(t_i | \mathbf{X}, \mathbf{w}, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp \left[-\frac{1}{2\sigma^2} (t_i - (\mathbf{X}\mathbf{w})_i)^2 \right] \quad (7)$$

A precision parameter β is defined which is given as $\beta^{-1} = \sigma^2$. Thus, we have the following modified likelihood function as follows.

$$p(t_i | \mathbf{X}, \mathbf{w}, \beta) = \sqrt{\frac{\beta}{2\pi}} \exp \left[-\frac{\beta}{2} (t_i - (\mathbf{X}\mathbf{w})_i)^2 \right] \quad (8)$$

Assuming the data is drawn independently, the likelihood is the joint probability given as the product of individual marginal probabilities. It is also assumed the value of β is known or assumed.

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N p(t_i | \mathbf{X}, \mathbf{w}, \beta) \quad (9)$$

The log likelihood of equation 9 is given as

$$\ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \sum_{i=1}^N \ln p(t_i | \mathbf{X}, \mathbf{w}, \beta) \quad (10)$$

$$= \sum_{i=1}^N \ln \left\{ \sqrt{\frac{\beta}{2\pi}} \exp \left[-\frac{\beta}{2} (t_i - (\mathbf{X}\mathbf{w})_i)^2 \right] \right\} \quad (11)$$

$$= \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \frac{\beta}{2} \sum_{i=1}^N (t_i - (\mathbf{X}\mathbf{w})_i)^2 \quad (12)$$

$$= \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 \quad (13)$$

Therefore,

$$\ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) \propto -\frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 \quad (14)$$

The posterior probability to determine the parameters is given as the product of likelihood function and prior.

$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \beta) \propto p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}) \quad (15)$$

The value of the prior $p(\mathbf{w}) = 1$.

By maximizing the negative likelihood (or posterior distribution with prior as one) with respect to \mathbf{w} ,

$$\frac{\partial}{\partial \mathbf{w}} \{-\ln p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \beta)\} \stackrel{!}{=} 0 \quad (16)$$

$$\frac{\partial}{\partial \mathbf{w}} \{-\ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta)\} \stackrel{!}{=} 0 \quad (17)$$

$$\frac{\partial}{\partial \mathbf{w}} \left\{ \frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 \right\} \stackrel{!}{=} 0 \quad (18)$$

$$\frac{\partial}{\partial \mathbf{w}} \left\{ \frac{\beta}{2} (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w}) \right\} \stackrel{!}{=} 0 \quad (19)$$

$$\beta (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{t}) \stackrel{!}{=} 0 \quad (20)$$

Therefore, \mathbf{w}_{ML} is evaluated.

$$\mathbf{w}_{\text{ML}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \quad (21)$$

It can be seen that the maximum likelihood results into least square estimator. We can similarly estimate β_{ML} by maximizing the posterior with respect to β . The known value of \mathbf{w}_{ML} can now be utilized here.

Taking the log likelihood in equation 10, the following could be shown.

$$\frac{\partial}{\partial \beta} \{-\ln p(\mathbf{w}_{\text{ML}} | \mathbf{t}, \mathbf{X}, \beta)\} \stackrel{!}{=} 0 \quad (22)$$

$$\frac{\partial}{\partial \beta} \{-\ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}_{\text{ML}}, \beta)\} \stackrel{!}{=} 0 \quad (23)$$

$$\frac{\partial}{\partial \beta} \left\{ -\frac{N}{2} \ln \beta + \frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{ML}}\|^2 \right\} \stackrel{!}{=} 0 \quad (24)$$

$$\|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{ML}}\|^2 - \frac{N}{\beta} \stackrel{!}{=} 0 \quad (25)$$

Therefore, the value β_{ML} is determined to be

$$\beta_{\text{ML}} = \frac{N}{\|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{ML}}\|^2} \quad (26)$$

1.2 Maximum a posteriori estimation

In the case of maximum a posteriori (MAP) estimation, the distribution of prior over parameters is known.

1.2.1 Conjugate Gaussian prior

The prior distribution is given as follows

$$p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{0}, \mathbf{I}/\alpha) = \left(\frac{\alpha}{2\pi}\right)^{M/2} \exp\left\{-\frac{\alpha}{2} \|\mathbf{w}\|^2\right\} \quad (27)$$

The posterior distribution is shown as follows.

$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha) \quad (28)$$

where β as defined earlier is the precision parameter of the likelihood, α is the hyperparameter (also a precision parameter of the prior distribution) which controls the distribution of model parameters. It is assumed that the value of α and β is known.

The log of the posterior is given as follows

$$\ln p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) = \ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) + \ln p(\mathbf{w} | \alpha) \quad (29)$$

The log likelihood is known from equation 10. Therefore,

$$\ln p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 + \frac{M}{2} \ln \left(\frac{\alpha}{2\pi}\right) - \frac{\alpha}{2} \|\mathbf{w}\|^2 \quad (30)$$

Therefore,

$$\ln p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto -\frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 - \frac{\alpha}{2} \|\mathbf{w}\|^2 \quad (31)$$

Maximizing the negative log of posterior with respect to \mathbf{w} .

$$\frac{\partial}{\partial \mathbf{w}} \{-\ln p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta)\} \stackrel{!}{=} 0 \quad (32)$$

$$\frac{\partial}{\partial \mathbf{w}} \left\{ \frac{\beta}{2} [(\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w})] + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\} \stackrel{!}{=} 0 \quad (33)$$

$$\beta (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{t}) + \alpha \mathbf{w} \stackrel{!}{=} 0 \quad (34)$$

Let us assign a regularization parameter $\lambda = \alpha/\beta$. Therefore, the value of parameter using maximum a posteriori estimation \mathbf{w}_{MAP} is given as

$$\mathbf{w}_{\text{MAP}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t} \quad (35)$$

1.2.2 Jeffreys prior

Jeffreys prior is given as

$$p(\sigma) = \frac{1}{\sigma} \quad (36)$$

Parameter transformation: $\sigma \rightarrow \beta = 1/\sigma^2$

$$\begin{aligned} p_\beta(\beta) &= p_\sigma(\sigma) \left| \frac{d\sigma}{d\beta} \right| \\ &\propto \sqrt{\beta} \frac{1}{1/\sigma^3} \Big|_{\sigma=1/\sqrt{\beta}} \\ &= \sqrt{\beta} \sigma^3 \Big|_{\sigma=1/\sqrt{\beta}} \\ &= \sqrt{\beta} \beta^{-3/2} = 1/\beta \end{aligned}$$

Therefore,

$$p(\beta) = \frac{1}{\beta} \quad (37)$$

Evaluation for precision parameter β

The posterior distribution is given as

$$p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \beta) \propto p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) p(\beta) \quad (38)$$

The log likelihood of the posterior distribution is given as

$$\ln p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \beta) = \ln p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) + \ln p(\beta) \quad (39)$$

Using equation 10 to give the log likelihood, the above equation would be as follows

$$\ln p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 - \ln \beta \quad (40)$$

For MAP estimation on \mathbf{w} , $\mathbf{w}_{\text{MAP}} = \mathbf{w}_{\text{ML}}$ (refer equation 21) because only a Gaussian likelihood is used for parameter distribution.

Maximizing the log posterior with respect to β , equation 40 is used.

$$\frac{\partial}{\partial \beta} \{-\ln p(\mathbf{w}_{\text{MAP}} \mid \mathbf{t}, \mathbf{X}, \beta)\} \stackrel{!}{=} 0 \quad (41)$$

$$\frac{\partial}{\partial \beta} \left\{ \frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{MAP}}\|^2 \right\} - \frac{\partial}{\partial \beta} \left\{ \frac{N}{2} \ln \beta \right\} + \frac{\partial}{\partial \beta} \{\ln \beta\} \stackrel{!}{=} 0 \quad (42)$$

$$\frac{1}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{MAP}}\|^2 - \frac{N}{2\beta} + \frac{1}{\beta} \stackrel{!}{=} 0 \quad (43)$$

Therefore,

$$\beta_{\text{MAP}} = \frac{N-2}{\|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{MAP}}\|^2} \quad (44)$$

1.3 Jeffreys hyperprior with conjugate Gaussian prior

Jeffreys prior for the hyperparameter α is given as follows,

$$p(\alpha) = \frac{1}{\alpha} \quad (45)$$

The joint distribution $p(\mathbf{w}, \alpha)$ is given as follows

$$p(\mathbf{w}, \alpha) = p(\mathbf{w} | \alpha)p(\alpha) \quad (46)$$

The distribution $p(\mathbf{w} | \alpha)$ is known from equation 27. Therefore, the joint probability $p(\mathbf{w}, \alpha)$, i.e, Jeffreys hyperprior over hyperparameter α and Gaussian prior over parameter \mathbf{w} is evaluated as follows.

$$p(\mathbf{w}, \alpha) = \left(\frac{\alpha}{2\pi}\right)^{M/2} \exp\left\{-\frac{\alpha}{2} \|\mathbf{w}\|^2\right\} \times \frac{1}{\alpha} \quad (47)$$

$$p(\mathbf{w}, \alpha) = \left(\frac{\alpha}{2\pi}\right)^{(M-2)/2} \exp\left\{-\frac{\alpha}{2} \|\mathbf{w}\|^2\right\} \quad (48)$$

Evaluation for hyperparameter α , precision parameter β , parameter \mathbf{w}

The posterior distribution (Conjugate Gaussian prior and Jeffreys prior combined) is given as follows.

$$p(\mathbf{w}, \alpha, \beta | \mathbf{t}, \mathbf{X}) \propto p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \alpha, \beta)p(\mathbf{w}, \alpha)p(\beta) \quad (49)$$

The distribution of $p(\beta)$ is known from equation 37 to be $p(\beta) = 1/\beta$. Therefore, the log of the posterior distribution is

$$\ln p(\mathbf{w}, \alpha, \beta | \mathbf{t}, \mathbf{X}) = \ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \alpha, \beta) + \ln p(\mathbf{w}, \alpha) + \ln p(\beta) \quad (50)$$

Evaluating the above equation results to the following.

$$\ln p(\mathbf{w}, \alpha, \beta | \mathbf{t}, \mathbf{X}) = \frac{N}{2} \ln \beta - \frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 + \frac{(M-2)}{2} \ln \alpha - \frac{\alpha}{2} \|\mathbf{w}\|^2 - \ln \beta \quad (51)$$

Initially, maximizing the posterior (refer equation 51) with respect to \mathbf{w} , we get

$$\frac{\partial}{\partial \mathbf{w}} \{-\ln p(\mathbf{w}, \alpha, \beta \mid \mathbf{t}, \mathbf{X})\} \stackrel{!}{=} 0 \quad (52)$$

$$\frac{\partial}{\partial \mathbf{w}} \left\{ \frac{\beta}{2} [(\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w})] + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\} \stackrel{!}{=} 0 \quad (53)$$

$$\beta (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{t}) + \alpha \mathbf{w} \stackrel{!}{=} 0 \quad (54)$$

Given the regularization parameter $\lambda = \alpha/\beta$, \mathbf{w}_{MAP} is the same as per the equation 35.

$$\mathbf{w}_{\text{MAP}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t} \quad (55)$$

Finally, maximizing posterior (refer equation 51) with respect to β , we get the following

$$\frac{\partial}{\partial \beta} \{-\ln p(\mathbf{w}_{\text{MAP}}, \alpha, \beta \mid \mathbf{t}, \mathbf{X})\} \stackrel{!}{=} 0 \quad (56)$$

$$\frac{\partial}{\partial \beta} \left\{ \frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{MAP}}\|^2 \right\} - \frac{\partial}{\partial \beta} \left\{ \frac{N}{2} \ln \beta \right\} + \frac{\partial}{\partial \beta} \{\ln \beta\} \stackrel{!}{=} 0 \quad (57)$$

$$\frac{1}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{MAP}}\|^2 - \frac{N}{2\beta} + \frac{1}{\beta} \stackrel{!}{=} 0 \quad (58)$$

Therefore, β_{MAP} is as follows.

$$\beta_{\text{MAP}} = \frac{N - 2}{\|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{MAP}}\|^2} \quad (59)$$

It can also be seen that $\beta_{\text{MAP}} = \beta_{\text{ML}}$ (refer equation 26).

Maximizing the posterior (refer equation 51) with respect to α , we get the following.

$$\frac{\partial}{\partial \alpha} \{-\ln p(\mathbf{w}_{\text{MAP}}, \alpha, \beta_{\text{MAP}} \mid \mathbf{t}, \mathbf{X})\} \stackrel{!}{=} 0 \quad (60)$$

$$\|\mathbf{w}_{\text{MAP}}\|^2 - \frac{M - 2}{\alpha} \stackrel{!}{=} 0 \quad (61)$$

Therefore, α_{MAP} is given as follows.

$$\alpha_{\text{MAP}} = \frac{M - 2}{\|\mathbf{w}_{\text{MAP}}\|^2} \quad (62)$$

2 Multivariate Gaussian Distribution

For a N-dimensional vector $\mathbf{x} = (x_1, \dots, x_N)^T$, the multivariate Gaussian distribution takes the following form

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{N/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (63)$$

where $\boldsymbol{\mu}$ is an N -dimensional mean vector, $\boldsymbol{\Sigma}$ is an $N \times N$ covariance matrix and $|\boldsymbol{\Sigma}|$ is the determinant of $\boldsymbol{\Sigma}$.

2.1 Conjugate Bayesian analysis of the Gaussian distribution

In the conjugate Bayesian analysis, the prior and likelihood distribution is the same. It can be shown that the resulting posterior is a multivariate Gaussian when the prior and likelihood are a Gaussian distribution.

Rewriting the posterior distribution as proportional to likelihood and prior.

$$p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} \mid \alpha) \quad (64)$$

Therefore, it can be shown that the $\ln p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \alpha, \beta)$ for a Gaussian distribution as

$$\ln p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto -\frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 - \frac{\alpha}{2} \|\mathbf{w}\|^2 \quad (65)$$

Evaluating this further,

$$\ln p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto -\frac{\beta}{2} (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w}) - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \quad (66)$$

$$\ln p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto -\frac{1}{2} \mathbf{w}^T \underbrace{(\beta \mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})}_{\boldsymbol{\Sigma}^{-1}} \mathbf{w} + \mathbf{w}^T \underbrace{\beta \mathbf{X}^T \mathbf{t}}_{\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}} \quad (67)$$

Therefore, $\boldsymbol{\Sigma}$ and $\boldsymbol{\mu}$ is given as

$$\boldsymbol{\Sigma} = (\beta \mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})^{-1} \quad (68)$$

$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \beta \mathbf{X}^T \mathbf{t} \quad (69)$$

$$\boldsymbol{\mu} = (\beta \mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})^{-1} \beta \mathbf{X}^T \mathbf{t} \quad (70)$$

$$\boldsymbol{\mu} = (\beta \mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})^{-1} \beta \mathbf{X}^T \mathbf{t} \quad (71)$$

$$\boldsymbol{\mu} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t} \quad (72)$$

$$\boldsymbol{\mu} \equiv \mathbf{w}_{\text{MAP}} \quad (73)$$

The posterior $\ln p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta)$ is further given as follows

$$\ln p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto -\frac{1}{2} \mathbf{w}^T \boldsymbol{\Sigma}^{-1} \mathbf{w} + \mathbf{w}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \quad (74)$$

$$\ln p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto -\frac{1}{2} \mathbf{w}^T \boldsymbol{\Sigma}^{-1} \mathbf{w} + \mathbf{w}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \frac{1}{2} \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \frac{1}{2} \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \quad (75)$$

$$\ln p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto -\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{w} - \boldsymbol{\mu}) + \frac{1}{2} \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \quad (76)$$

$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) \propto \exp \left\{ -\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right\} \quad (77)$$

The posterior over the regression parameters $\mathbf{w} = (w_1, w_2, \dots, w_M)^T$ is a multivariate Gaussian.

$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \alpha, \beta) = \frac{1}{|2\pi|^{(M+1)/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right\} \quad (78)$$

where $\boldsymbol{\Sigma}$ is an $(M+1) \times (M+1)$ covariance matrix, given as $\boldsymbol{\Sigma} = (\beta \mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})^{-1}$, $|\boldsymbol{\Sigma}|$ is the determinant of $\boldsymbol{\Sigma}$. and $\boldsymbol{\mu}$ is an $(M+1)$ -dimensional mean vector given as $\boldsymbol{\mu} = \boldsymbol{\Sigma} \beta \mathbf{X}^T \mathbf{t}$.

3 Gamma Distribution

The likelihood function for the target values given β (precision parameter) can be written as

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(t_i | \mathbf{X}, \mathbf{w}, \beta^{-1}) \quad (79)$$

$$\propto \beta^{\frac{N}{2}} \exp \left\{ -\frac{\beta}{2} \sum_{i=1}^N (t_i - (\mathbf{X} \mathbf{w})_i)^2 \right\} \quad (80)$$

The corresponding conjugate prior¹ should be proportional to the product of the power of β and the exponential of the linear function of β .

Gamma distribution is a conjugate prior to several likelihood distributions (Gaussian, Poisson, exponential, etc). For example, gamma distribution over precision parameter β can be given as follows.

$$\text{Gam}(\beta | a, b) = \frac{b^a}{\Gamma(a)} \beta^{a-1} \exp(-b\beta) \quad (81)$$

¹posterior distribution is in the same probability distribution family as the prior distribution

where, $\Gamma(a)$ is a gamma function that ensures the gamma distribution is properly normalized, a is called as the shape parameter and b is called as the rate parameter.

The expectation $\mathbb{E}[\beta]$ is given as follows.

$$\mathbb{E}[\beta] = \frac{a}{b} \quad (82)$$

The mode, $\text{mode}[\beta]$ which is equivalent to maximizing the posterior with respect to β is given as follows.

$$\text{mode}[\beta] = \frac{a - 1}{b} \quad (83)$$

3.1 Jeffreys Priors

Let us consider the posterior distribution considering the Jeffreys prior for $p(\beta)$ as given in equation 37.

$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \beta) \propto p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) p(\beta) \quad (84)$$

$$\propto \beta^{\frac{N}{2}} \exp \left\{ -\frac{\beta}{2} \sum_{i=1}^N (t_i - (\mathbf{X}\mathbf{w})_i)^2 \right\} \times \frac{1}{\beta} \quad (85)$$

$$\propto \beta^{\frac{N}{2}-1} \exp \left\{ -\frac{\beta}{2} \sum_{i=1}^N (t_i - (\mathbf{X}\mathbf{w})_i)^2 \right\} \quad (86)$$

By comparing this to the gamma distribution in equation 81, $\text{mode}[\beta | \mathbf{t}, \mathbf{X}, \mathbf{w}]$ is as follows

$$\text{mode}[\beta | \mathbf{t}, \mathbf{X}, \mathbf{w}] = \frac{N - 2}{\|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2} \quad (87)$$

This is equivalent to β_{MAP} from the result in equation 44. Similarly, $\mathbb{E}[\beta | \mathbf{t}, \mathbf{X}, \mathbf{w}]$ is given as follows.

$$\mathbb{E}[\beta | \mathbf{t}, \mathbf{X}, \mathbf{w}] = \frac{N}{\|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2} \quad (88)$$

This can also be shown to hold true for the joint distribution $p(\mathbf{w}, \alpha)$ in equation 46.

$$p(\mathbf{w}, \alpha) \propto \alpha^{\frac{M}{2}-1} \exp \left\{ -\frac{\alpha}{2} \|\mathbf{w}\|^2 \right\} \quad (89)$$

$$\text{mode}[\alpha | \mathbf{w}] = \frac{M - 2}{\|\mathbf{w}_{\text{MAP}}\|^2} \quad (90)$$

This is same as α_{MAP} in equation 62. Similarly $\mathbb{E}[\alpha | \mathbf{w}]$ is

$$\mathbb{E}[\alpha | \mathbf{w}] = \frac{M}{\|\mathbf{w}_{\text{MAP}}\|^2} \quad (91)$$

3.2 Gamma Priors with Jeffreys Priors

Jeffreys priors are given as $p(\alpha) = \alpha^{-1}$ and $p(\beta) = \beta^{-1}$ for hyperparameter α (refer equation 45) and precision parameter β (refer equation 37) respectively.

From the gamma distribution in equation 81, a prior distribution $p(\alpha | a_\alpha, b_\alpha)$ can be given as follows.

$$p(\alpha | a_\alpha, b_\alpha) = \frac{b_\alpha^{a_\alpha}}{\Gamma(a_\alpha)} \alpha^{a_\alpha-1} \exp(-b_\alpha \alpha) \quad (92)$$

Similarly for $p(\beta | a_\beta, b_\beta)$ can be given as follows.

$$p(\beta | a_\beta, b_\beta) = \frac{b_\beta^{a_\beta}}{\Gamma(a_\beta)} \beta^{a_\beta-1} \exp(-b_\beta \beta) \quad (93)$$

In the special case when $a_\alpha = 0, b_\alpha = 0, a_\beta = 0, b_\beta = 0$. The gamma distribution above results to Jeffreys priors $p(\alpha)$ and $p(\beta)$.

The posterior distribution $p(\alpha | a_\alpha, b_\alpha, \mathbf{w})$ is therefore,

$$p(\alpha | a_\alpha, b_\alpha, \mathbf{w}) = p(\mathbf{w} | \alpha) p(\alpha | a_\alpha, b_\alpha) \quad (94)$$

$$= \left(\frac{\alpha}{2\pi} \right)^{M/2} \exp \left\{ -\frac{\alpha}{2} \|\mathbf{w}\|^2 \right\} \frac{b_\alpha^{a_\alpha}}{\Gamma(a_\alpha)} \alpha^{a_\alpha-1} \exp(-b_\alpha \alpha) \quad (95)$$

Similarly, the posterior distribution $p(\beta | a_\beta, b_\beta, \mathbf{w})$

$$p(\beta | a_\beta, b_\beta, \mathbf{w}) = p(\mathbf{w} | \beta) p(\beta | a_\beta, b_\beta) \quad (96)$$

$$= \left(\frac{\beta}{2\pi} \right)^{N/2} \exp \left\{ -\frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 \right\} \frac{b_\beta^{a_\beta}}{\Gamma(a_\beta)} \beta^{a_\beta-1} \exp(-b_\beta \beta) \quad (97)$$

Initially, maximizing log posterior from equation 94 with respect to alpha.

$$\frac{\partial}{\partial \alpha} \{-\ln p(\alpha | a_\alpha, b_\alpha, \mathbf{w})\} \stackrel{!}{=} 0 \quad (98)$$

$$\frac{\alpha}{2} \|\mathbf{w}\|^2 - \frac{M}{2} \ln \alpha - (a_\alpha - 1) \ln \alpha + b_\alpha \alpha \stackrel{!}{=} 0 \quad (99)$$

Therefore, α_{MAP} is given as

$$\alpha_{MAP} = \frac{M - 2 + 2a_\alpha}{\|\mathbf{w}\|^2 + 2b_\alpha} \quad (100)$$

In the special case, when $a_\alpha = 0$ and $b_\alpha = 0$, the above equation is equal to 62.

Similarly, maximizing log posterior from equation 96 with respect to beta,

$$\frac{\partial}{\partial \beta} \{-\ln p(\beta | a_\beta, b_\beta, \mathbf{w})\} \stackrel{!}{=} 0 \quad (101)$$

$$\frac{\beta}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 - \frac{N}{2} \ln \beta - (a_\beta - 1) \ln \beta + b_\beta \beta \stackrel{!}{=} 0 \quad (102)$$

Therefore, β_{MAP} is given as follows

$$\beta_{MAP} = \frac{N - 2 + 2a_\beta}{\|\mathbf{t} - \mathbf{X}\mathbf{w}\|^2 + 2b_\beta} \quad (103)$$

In the special case, $a_\beta = 0$ and $b_\beta = 0$, the above equation is equal to 44.

Under the assumption that $a_\alpha = b_\alpha = a_\beta = b_\beta = \epsilon$, expectation, variance and mode can be given as follows.

$$\mathbb{E}[\alpha] = \frac{a_\alpha}{b_\alpha} = 1 \quad (104)$$

$$\text{var}[\alpha] = \frac{a_\alpha}{b_\alpha^2} = \frac{1}{\epsilon} \quad (105)$$

$$\text{mode}[\alpha] = \frac{a_\alpha - 1}{b_\alpha} = \frac{\epsilon - 1}{\epsilon} \quad (106)$$

The results are similar to $\mathbb{E}[\beta]$, $\text{var}[\beta]$ and $\text{mode}[\beta]$.

If the value of σ is known, a_β and b_β can be estimated.

$$\mathbb{E}[\beta] = \frac{a_\beta}{b_\beta} = \frac{1}{\sigma^2} \quad (107)$$

The $\text{var}[\beta]$ is also known. Let us assume $\text{var}[\beta] = \epsilon$.

$$\text{var}[\beta] = \mathbb{E}[\beta^2] - \mathbb{E}[\beta]^2 \quad (108)$$

$$= \frac{a_\beta}{b_\beta} = \epsilon \quad (109)$$

Therefore this results to the following.

$$a_\beta = \frac{1}{\epsilon\sigma^4} \quad (110)$$

$$b_\beta = \frac{1}{\epsilon\sigma^2} \quad (111)$$

4 Laplace distribution

The Laplace distribution is also called as the double exponential distribution because it can be thought of as 2 exponential distributions spliced together back-to-back. A

random variable has a Laplace(μ, σ) distribution if its probability density function is given as follows.

$$f(x | \mu, \sigma) = \frac{1}{2\sigma} \exp - \frac{|x - \mu|}{\sigma} \quad (112)$$

where, μ is the location parameter and σ as scale parameter.

Under the assumption that every output/target values follows a Laplace distribution. The likelihood $p(t_i | \mathbf{X}, \mathbf{w}, \sigma)$ following the Laplace distribution is given as follows.

$$p(t_i | \mathbf{X}, \mathbf{w}, \sigma) = \frac{1}{2\sigma} \exp \left\{ - \frac{|t_i - (\mathbf{X}\mathbf{w})_i|}{\sigma} \right\} \quad (113)$$

Let us define a precision parameter β which is given as follows.

$$\beta = \frac{1}{\sigma} \quad (114)$$

According to the product rule, the total likelihood is given as the product of the individual marginal probabilities as follows.

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N p(t_i | \mathbf{X}, \mathbf{w}, \beta) \quad (115)$$

The log likelihood of the above equation 115 can be given as follows.

$$\ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \sum_{i=1}^N \ln p(t_i | \mathbf{X}, \mathbf{w}, \beta) \quad (116)$$

$$= \sum_{i=1}^N \ln \left\{ \frac{\beta}{2} \exp (-\beta |t_i - (\mathbf{X}\mathbf{w})_i|) \right\} \quad (117)$$

$$= \frac{N}{2} \ln \beta - \beta \sum_{i=1}^N |t_i - (\mathbf{X}\mathbf{w})_i| \quad (118)$$

$$= \frac{N}{2} \ln \beta - \beta \|\mathbf{t} - \mathbf{X}\mathbf{w}\| \quad (119)$$

By maximizing the log likelihood with respect to \mathbf{w} , the value of \mathbf{w}_{ML} can be found. However, it cannot be treated analytically. Therefore, a suitable optimizer such as Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm could be used.

After evaluating \mathbf{w}_{ML} , β_{ML} could be found out by maximizing the log likelihood with respect to β .

$$\frac{\partial}{\partial \beta} \{-\ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}_{ML}, \beta)\} \stackrel{!}{=} 0 \quad (120)$$

$$\frac{\partial}{\partial \beta} \left\{ \beta \|\mathbf{t} - \mathbf{X}\mathbf{w}_{ML}\| - \frac{N}{2} \ln \beta \right\} \stackrel{!}{=} 0 \quad (121)$$

Therefore, β_{ML} is given as follows

$$\beta_{\text{ML}} = \frac{N}{2 \|\mathbf{t} - \mathbf{X}\mathbf{w}_{\text{ML}}\|} \quad (122)$$

If priors are considered, such as Gamma priors or Jeffreys priors, the above equation would be modified.