

Bayesian Linear Regression

Jan Drugowitsch

Apr 2008, last update: May 2010

The Model

The model assumes a linear relation between D -dimensional inputs \mathbf{x} and outputs y and constant-variance Gaussian noise, such that the data likelihood is given by

$$p(y|\mathbf{x}, \mathbf{w}, \tau) = \mathcal{N}(y|\mathbf{w}^T \mathbf{x}, \tau^{-1}) = \left(\frac{\tau}{2\pi}\right)^{1/2} \exp\left(-\frac{\tau}{2}(y - \mathbf{w}^T \mathbf{x})^2\right). \quad (1)$$

Given all data $\mathcal{D} = \{\mathbf{X}, \mathbf{Y}\}$, with $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and $\mathbf{Y} = \{y_1, \dots, y_N\}$, the data likelihood is

$$p(\mathbf{Y}|\mathbf{X}, \mathbf{w}, \tau) = \prod_n p(y_n|\mathbf{x}_n, \mathbf{w}, \tau). \quad (2)$$

The prior on \mathbf{w} and τ is conjugate normal inverse-gamma

$$\begin{aligned} p(\mathbf{w}, \tau|\alpha) &= \mathcal{N}(\mathbf{w}|0, (\tau\alpha)^{-1}\mathbf{I})\text{Gam}(\tau|a_0, b_0) \\ &= \left(\frac{\alpha}{2\pi}\right)^{D/2} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{D/2+a_0-1} \exp\left(-\frac{\tau}{2}(\alpha\mathbf{w}^T \mathbf{w} + 2b_0)\right), \end{aligned} \quad (3)$$

parametrised by α . This hyper-parameter is assigned the hyper-prior

$$p(\alpha) = \text{Gam}(\alpha|c_0, d_0) = \frac{1}{\Gamma(c_0)} d_0^{c_0} \alpha^{c_0-1} \exp(-d_0\alpha). \quad (4)$$

Due to the hyper-prior, there is no analytic solution to the posteriors and variational Bayesian inference will be applied.

Variational Bayesian Inference

The variational posteriors are found by maximising the variational bound

$$\mathcal{L}(q) = \iiint q(\mathbf{w}, \tau, \alpha) \ln \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{w}, \tau)p(\mathbf{w}, \tau|\alpha)p(\alpha)}{q(\mathbf{w}, \tau, \alpha)} d\mathbf{w} d\tau d\alpha \leq \ln p(\mathcal{D}), \quad (5)$$

where $p(\mathcal{D})$ is the model evidence, and under the assumption that the variational distribution $q(\mathbf{w}, \tau, \alpha)$, which approximates the posterior $p(\mathbf{w}, \tau, \alpha|\mathcal{D})$, factors into $q(\mathbf{w}, \tau)q(\alpha)$.

The variational posterior for \mathbf{w}, τ is given by

$$\ln q^*(\mathbf{w}, \tau) = \ln p(\mathbf{Y}|\mathbf{X}, \mathbf{w}, \tau) + \mathbb{E}_\alpha(\ln p(\mathbf{w}, \tau|\alpha)) + \text{const.} \quad (6)$$

$$\begin{aligned} &= \left(\frac{D}{2} + a_0 - 1 + \frac{N}{2}\right) \ln \tau \\ &\quad - \frac{\tau}{2} \left(\mathbf{w}^T \left(\mathbb{E}_\alpha(\alpha)\mathbf{I} + \sum_n \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} + \sum_n y_n^2 - 2\mathbf{w}^T \sum_n \mathbf{x}_n y_n + 2b_0 \right) + \text{const.} \\ &= \ln \mathcal{N}(\mathbf{w}|\mathbf{w}_N, \tau^{-1}\mathbf{V}_N) \text{Gam}(\tau|a_N, b_N), \end{aligned} \quad (8)$$

with

$$\mathbf{V}_N^{-1} = \mathbb{E}_\alpha(\alpha) \mathbf{I} + \sum_n \mathbf{x}_n \mathbf{x}_n^T, \quad (9)$$

$$\mathbf{w}_N = \mathbf{V}_N \sum_n \mathbf{x}_n y_n, \quad (10)$$

$$a_N = a_0 + \frac{N}{2}, \quad (11)$$

$$\begin{aligned} b_N &= b_0 + \frac{1}{2} \left(\sum_n y_n^2 - \mathbf{w}_N^T \mathbf{V}_N^{-1} \mathbf{w}_N \right) \\ &= b_0 + \frac{1}{2} \left(\sum_n (y_n - \mathbf{w}_N^T \mathbf{x}_n)^2 + \mathbb{E}_\alpha(\alpha) \mathbf{w}_N^T \mathbf{w}_N \right). \end{aligned} \quad (12)$$

The variational posterior for α is

$$\ln q^*(\alpha) = \mathbb{E}_{\mathbf{w}, \tau}(\ln p(\mathbf{w}, \tau | \alpha)) + \ln p(\alpha) + \text{const.} \quad (13)$$

$$= \left(c_0 - 1 + \frac{D}{2} \right) \ln \alpha - \alpha \left(d_0 + \frac{1}{2} \mathbb{E}_{\mathbf{w}, \tau}(\tau \mathbf{w}^T \mathbf{w}) \right) + \text{const.} \quad (14)$$

$$= \ln \text{Gam}(\alpha | c_N, d_N), \quad (15)$$

with

$$c_N = c_0 + \frac{D}{2}, \quad (16)$$

$$d_N = d_0 + \frac{1}{2} \mathbb{E}_{\mathbf{w}, \tau}(\tau \mathbf{w}^T \mathbf{w}). \quad (17)$$

The expectations are evaluated with respect to the variational distribution and are given by

$$\mathbb{E}_{\mathbf{w}, \tau}(\tau \mathbf{w}^T \mathbf{w}) = \frac{a_N}{b_N} \mathbf{w}_N^T \mathbf{w}_N + \text{Tr}(\mathbf{V}_N), \quad (18)$$

$$\mathbb{E}_\alpha(\alpha) = \frac{c_N}{d_N}. \quad (19)$$

The variational bound itself consists of

$$\begin{aligned} \mathcal{L}(q) &= \mathbb{E}_{\mathbf{w}, \tau}(\ln p(\mathbf{Y} | \mathbf{X}, \mathbf{w}, \tau)) + \mathbb{E}_{\mathbf{w}, \tau, \alpha}(\ln p(\mathbf{w}, \tau | \alpha)) + \mathbb{E}_\alpha(\ln p(\alpha)) \\ &\quad - \mathbb{E}_{\mathbf{w}, \tau}(\ln p(\mathbf{w}, \tau)) - \mathbb{E}_\alpha(\ln p(\alpha)), \end{aligned} \quad (20)$$

$$\begin{aligned} \mathbb{E}_{\mathbf{w}, \tau}(\ln p(\mathbf{Y} | \mathbf{X}, \mathbf{w}, \tau)) &= \frac{N}{2} (\psi(a_N) - \ln b_N - \ln 2\pi) \\ &\quad - \frac{1}{2} \sum_n \left(\frac{a_N}{b_N} (y_n - \mathbf{w}_N^T \mathbf{x}_n)^2 + \mathbf{x}_n^T \mathbf{V}_N \mathbf{x}_n \right), \end{aligned} \quad (21)$$

$$\begin{aligned} \mathbb{E}_{\mathbf{w}, \tau, \alpha}(\ln p(\mathbf{w}, \tau | \alpha)) &= \frac{D}{2} (\psi(a_N) - \ln b_N + \psi(c_N) - \ln d_N - \ln 2\pi) \\ &\quad - \frac{1}{2} \frac{c_N}{d_N} \left(\frac{a_N}{b_N} \mathbf{w}_N^T \mathbf{w}_N + \text{Tr}(\mathbf{V}_N) \right) \\ &\quad - \ln \Gamma(a_0) + a_0 \ln b_0 + (a_0 - 1)(\psi(a_N) - \ln b_N) - b_0 \frac{a_N}{b_N}, \end{aligned} \quad (22)$$

$$\mathbb{E}_\alpha(\ln p(\alpha)) = -\ln \Gamma(c_0) + d_0 \ln c_0 + (c_0 - 1)(\psi(c_N) - \ln d_N) - d_0 \frac{c_N}{d_N}, \quad (23)$$

$$\begin{aligned} \mathbb{E}_{\mathbf{w}, \tau}(\ln q(\mathbf{w}, \tau)) &= \frac{D}{2} (\psi(a_N) - \ln b_N - \ln 2\pi - 1) - \frac{1}{2} \ln |\mathbf{V}_N| \\ &\quad - \ln \Gamma(a_N) + a_N \ln b_N + (a_N - 1)(\psi(a_N) - \ln b_N) - a_N, \end{aligned} \quad (24)$$

$$\mathbb{E}_\alpha(\ln q(\alpha)) = -\ln \Gamma(c_N) + (c_N - 1)\psi(c_N) + \ln d_N - c_N. \quad (25)$$

In combination, that gives

$$\begin{aligned}\mathcal{L}(q) = & -\frac{N}{2} \ln 2\pi - \frac{1}{2} \sum_n \left(\frac{a_N}{b_N} (y_n - \mathbf{w}_N^T \mathbf{x}_n)^2 + \mathbf{x}_n^T \mathbf{V}_N \mathbf{x}_n \right) + \frac{1}{2} \ln |\mathbf{V}_N| + \frac{D}{2} \\ & - \ln \Gamma(a_0) + a_0 \ln b_0 - b_0 \frac{a_N}{b_N} + \ln \Gamma(a_N) - a_N \ln b_N + a_N \\ & - \ln \Gamma(c_0) + c_0 \ln d_0 + \ln \Gamma(c_N) - c_N \ln d_N\end{aligned}\quad (26)$$

This bound is maximised by iterating over the updates for \mathbf{V}_N , \mathbf{w}_N , a_N , b_N , c_N , and d_N until $\mathcal{L}(q)$ reaches a plateau.

Predictive Density

The predictive density is evaluated by approximating the posterior $p(\mathbf{w}, \tau | \mathcal{D})$ by its variational counterpart $q(\mathbf{w}, \tau)$, to get

$$p(y | \mathbf{x}, \mathcal{D}) = \iint p(y | \mathbf{x}, \mathbf{w}, \tau) p(\mathbf{w}, \tau | \text{data}) d\mathbf{w} d\tau \quad (27)$$

$$\approx \iint p(y | \mathbf{x}, \mathbf{w}, \tau) q(\mathbf{w}, \tau) d\mathbf{w} d\tau \quad (28)$$

$$= \iint \mathcal{N}(y | \mathbf{w}^T \mathbf{x}, \tau^{-1}) \mathcal{N}(\mathbf{w} | \mathbf{w}_N, \tau^{-1} \mathbf{V}_N) \text{Gam}(\tau | a_N, b_N) d\mathbf{w} d\tau \quad (29)$$

$$= \int \mathcal{N}(y | \mathbf{w}^T \mathbf{x}, \tau^{-1} (1 + \mathbf{x}^T \mathbf{V}_N \mathbf{x})) \text{Gam}(\tau | a_N, b_N) d\tau \quad (30)$$

$$= \text{St} \left(y | \mathbf{w}^T \mathbf{x}, (1 + \mathbf{x}^T \mathbf{V}_N \mathbf{x})^{-1} \frac{a_N}{b_N}, 2a_N \right), \quad (31)$$

where standard results of convolving Gaussians with other Gaussians and Gamma distributions were used. The resulting distribution is a Student's t distribution with mean $\mathbf{w}^T \mathbf{x}$, precision $(1 + \mathbf{x}^T \mathbf{V}_N \mathbf{x})^{-1} a_N / b_N$, and $2a_N$ degrees of freedom, which has a variance of $(1 + \mathbf{x}^T \mathbf{V}_N \mathbf{x}) b_N / (a_N - 1)$.

Using Automatic Relevance Determination

Automatic Relevance Determination (ARD) determines the relevance of the elements of the input to determine the output by assigning a separate shrinkage prior to each element of the weight vector, which is in turn adjusted by a hyper-prior. While the data likelihood remains unchanged, the prior on \mathbf{w}, τ is modified to be

$$\begin{aligned}p(\mathbf{w}, \tau | \boldsymbol{\alpha}) &= \mathcal{N}(\mathbf{w} | \mathbf{0}, (\tau \mathbf{A})^{-1}) \text{Gam}(\tau | a_0, b_0) \\ &= \frac{|\mathbf{A}|^{1/2}}{\sqrt{2\pi}^D} \frac{b_0^{a_0}}{\Gamma(a_0)} \tau^{D/2 + a_0 - 1} \exp \left(-\frac{\tau}{2} (\mathbf{w}^T \mathbf{A} \mathbf{w} + 2b_0) \right),\end{aligned}\quad (32)$$

where the vector $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_D)^T$ forms the diagonal of \mathbf{A} . All of the α 's are independent, such that the hyper-prior is given by

$$p(\boldsymbol{\alpha}) = \prod_i \text{Gam}(\alpha_i | c_0, d_0) = \prod_i \frac{1}{\Gamma(c_0)} d_0^{c_0} \alpha_i^{c_0 - 1} \exp(-d_0 \alpha_i). \quad (33)$$

Variational Bayesian inference is performed as before, to get the variational posteriors

$$q^*(\mathbf{w}, \tau) = \mathcal{N}(\mathbf{w} | \mathbf{w}_N, \tau^{-1} \mathbf{V}_N) \text{Gam}(\tau | a_N, b_N), \quad q^*(\boldsymbol{\alpha}) = \prod_i \text{Gam}(\alpha_i | c_N, d_{Ni}), \quad (34)$$

with

$$\mathbf{V}_N^{-1} = \mathbb{E}_\alpha(\mathbf{A}) + \sum_n \mathbf{x}_n \mathbf{x}_n^T, \quad (35)$$

$$\mathbf{w}_N = \mathbf{V}_N \sum_n \mathbf{x}_n y_n, \quad (36)$$

$$a_N = a_0 + \frac{N}{2} \quad (37)$$

$$\begin{aligned} b_N &= b_0 + \frac{1}{2} \left(\sum_n y_n^2 - \mathbf{w}_N^T \mathbf{V}_N^{-1} \mathbf{w}_N \right) \\ &= b_0 + \frac{1}{2} \left(\sum_n (\mathbf{w}_N^T \mathbf{x}_n - y_n)^2 + \mathbf{w}_N^T \mathbb{E}_\alpha(\mathbf{A}) \mathbf{w}_N \right), \end{aligned} \quad (38)$$

$$c_N = c_0 + \frac{1}{2}, \quad (39)$$

$$d_{Ni} = d_0 + \frac{1}{2} \mathbb{E}_{\mathbf{w}, \alpha}(\tau w_i^2), \quad (40)$$

with expectations $\mathbb{E}_{\mathbf{w}, \alpha}(\tau w_i^2) = w_{Ni}^2 a_N / b_N + (\mathbf{V}_N)_{ii}$, and $\mathbb{E}_\alpha(\mathbf{A}) = \mathbf{A}_N$ is a diagonal matrix with elements $\mathbb{E}_\alpha(\alpha_i) = c_N / d_{Ni}$.

The variational bound changes to

$$\begin{aligned} \mathcal{L}(q) &= -\frac{N}{2} \ln 2\pi - \frac{1}{2} \sum_n \left(\frac{a_N}{b_N} (y_n - \mathbf{w}_N^T \mathbf{x}_n)^2 + \mathbf{x}_n^T \mathbf{V}_N \mathbf{x}_n \right) + \frac{1}{2} \ln |\mathbf{V}_N| + \frac{D}{2} \\ &\quad - \ln \Gamma(a_0) + a_0 \ln b_0 - b_0 \frac{a_N}{b_N} + \ln \Gamma(a_N) - a_N \ln b_N + a_N \\ &\quad + \sum_i (-\ln \Gamma(c_0) + c_0 \ln d_0 + \ln \Gamma(c_N) - c_N \ln d_{Ni}). \end{aligned} \quad (41)$$

The predictive distribution remains unchanged, as the prior does not appear in the expression for the variational posterior $p(\mathbf{w}, \tau)$.

Implementation

The scripts `bayes_linear_fit.m` and `bayes_linear_fit_ard.m` are straightforward implementations that compute the posterior parameters, without and with ARD, respectively. They operate by iteratively updating the parameters of $q^*(\mathbf{w}, \tau)$ and $q^*(\alpha)$, while monitoring $\mathcal{L}(q)$. The scripts stop as soon as either the change in $\mathcal{L}(q)$ between two consecutive iterations drops below 0.001% or the number of iterations exceeds 100.

The code is vectorised in order to speed up computation. In particular, the input \mathbf{X} is assumed to be an $N \times D$ matrix, with \mathbf{x}_n^T as its rows. \mathbf{y} is a column vector, containing all y_n 's. This allows for several vectorised operations, such as

$$\sum_n \mathbf{x}_n \mathbf{x}_n^T = \mathbf{X}' * \mathbf{X}, \quad (42)$$

$$\sum_n \mathbf{x}_n y_n = \mathbf{X}' * \mathbf{y}, \quad (43)$$

$$\mathbf{w}_N^T \mathbf{x}_n = (\mathbf{X} * \mathbf{w})_n, \quad (44)$$

$$\sum_n \mathbf{x}_n^T \mathbf{V}_N \mathbf{x}_n = \text{sum}(\text{sum}(\mathbf{X} .* (\mathbf{X} * \mathbf{V}))). \quad (45)$$

The rest of the code should be self-explanatory. The only mayor difference between `bayes_linear_fit .m` and `bayes_linear_fit_ard .m` is that in the latter version, the variables `E_a`, `dn`, and `E_t` are vectors, and all operations on these variables are vectorised.

The script `bayes_linear_post.m` computes the predictive density parameters for a set of input vectors, again given in matrix form X .