

Bayesian linear regression with unknown noise precision

January 2, 2023

For Bayesian linear regression with fixed noise precision, we relied on a theorem for linear Gaussian systems to go from the prior and likelihood to the posterior. Here, we first state and prove a new theorem that's very similar to that theorem, which will be crucially for helping us go from the prior and likelihood to the posterior for Bayesian linear regression with unknown noise precision.

Theorem 1. *Given the following prior and likelihood*

$$\begin{aligned} p(\sigma^2) &= \text{IG}(a_\sigma, b_\sigma) \\ p(\mathbf{x}) &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_x, \sigma^2 \boldsymbol{\Sigma}_x), \\ p(\mathbf{y} | \mathbf{x}, \sigma^2) &= \mathcal{N}(\mathbf{x} | \mathbf{A}\mathbf{x} + \mathbf{b}, \boldsymbol{\Sigma}_y) \end{aligned}$$

the posterior is given by

$$\begin{aligned} p(\mathbf{x}, \sigma^2 | \mathbf{y}) &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{x|y}, \sigma^2 \boldsymbol{\Sigma}_{x|y}) \text{IG}(\sigma^2 | a_{\sigma|y}, b_{\sigma|y}) \\ \boldsymbol{\mu}_{x|y} &= \boldsymbol{\Sigma}_{x|y} (\boldsymbol{\Sigma}_x^{-1} \boldsymbol{\mu}_x + \mathbf{A}^T \mathbf{y}) \\ \boldsymbol{\Sigma}_{x|y} &= (\boldsymbol{\Sigma}_x^{-1} + \mathbf{A}^T \mathbf{A})^{-1} \\ a_{\sigma|y} &= a_\sigma + D_y / 2 \\ b_{\sigma|y} &= b_\sigma + \frac{1}{2} (\boldsymbol{\mu}_x^T \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\mu}_x + \mathbf{y}^T \mathbf{y} - \boldsymbol{\mu}_{x|y}^T \boldsymbol{\Sigma}_{x|y}^{-1} \boldsymbol{\mu}_{x|y}) \end{aligned}$$

Proof. The proof should be analagous to the proof of the Bayes rule for linear Gaussian systems presented in Section 4.4.3 of Murphy. \square

Likelihood.

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{i=1}^N \mathcal{N}(y_i | \mu_i, \sigma^2) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$$

with $\boldsymbol{\mu} = \mathbf{X}\mathbf{w}$.

Prior.

$$p(\mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{w} | \mathbf{w}_0, \sigma^2 \mathbf{V}_0) \text{IG}(\sigma^2 | a_0, b_0)$$

Posterior. Applying Theorem 1, we obtain:

$$\begin{aligned} p(\mathbf{w}, \sigma^2 | \mathbf{y}, \mathbf{X}) &= \mathcal{N}(\mathbf{w} | \mathbf{w}_N, \sigma^2 \mathbf{V}_N) \text{IG}(\sigma^2 | a_N, b_N) \\ \mathbf{w}_N &= \mathbf{V}_N (\mathbf{V}_0^{-1} \mathbf{w}_0 + \mathbf{X}^T \mathbf{y}) \\ \mathbf{V}_N &= (\mathbf{V}_0^{-1} + \mathbf{X}^T \mathbf{X})^{-1} \\ a_N &= a_0 + N / 2 \\ b_N &= b_0 + \frac{1}{2} (\mathbf{w}_0^T \mathbf{V}_0^{-1} \mathbf{w}_0 + \mathbf{y}^T \mathbf{y} - \mathbf{w}_N^T \mathbf{V}_N^{-1} \mathbf{w}_N), \end{aligned}$$

which is exactly the same as Equation 7.69 to 7.73 from Murphy.

Marginal posterior.

$$\begin{aligned}
p(\mathbf{w} | \mathbf{y}, \mathbf{X}) &= \int p(\mathbf{w}, \sigma^2 | \mathbf{y}, \mathbf{X}) d\sigma^2 \\
&= \int \mathcal{N}(\mathbf{w} | \mathbf{w}_N, \sigma^2 \mathbf{V}_N) \text{IG}(\sigma^2 | a_N, b_N) d\sigma^2 \\
&= \int \mathcal{N}(\mathbf{w} | \mathbf{w}_N, (1/\sigma^2) \mathbf{V}_N) \text{Ga}(1/\sigma^2 | a_N, b_N) d\sigma^2 \\
&= \mathcal{T}(\mathbf{w} | \mathbf{w}_N, (b_N/a_N) \mathbf{V}_N, 2a_N) \quad (\text{by Eq 11.61})
\end{aligned}$$

$$\begin{aligned}
p(\sigma^2 | \mathbf{y}, \mathbf{X}) &= \int p(\mathbf{w}, \sigma^2 | \mathbf{y}, \mathbf{X}) d\mathbf{w} \\
&= \int \mathcal{N}(\mathbf{w} | \mathbf{w}_N, \sigma^2 \mathbf{V}_N) \text{IG}(\sigma^2 | a_N, b_N) d\mathbf{w} \\
&= \int \mathcal{N}(\mathbf{w} | \mathbf{w}_N, \sigma^2 \mathbf{V}_N) d\mathbf{w} \text{IG}(\sigma^2 | a_N, b_N) \\
&= \text{IG}(\sigma^2 | a_N, b_N)
\end{aligned}$$

Posterior predictive.

$$\begin{aligned}
p(\tilde{\mathbf{y}} | \tilde{\mathbf{X}}, \mathbf{X}, \mathbf{y}) &= \iint \mathcal{N}(\tilde{\mathbf{y}} | \tilde{\mathbf{X}}^T \mathbf{w}, \sigma^2 \mathbf{I}_m) \text{NIW}(\mathbf{w}, \sigma^2 | \mathbf{w}_N, \mathbf{V}_N, a_N, b_N) d\mathbf{w} d\sigma^2 \\
&= \int \left[\int \mathcal{N}(\tilde{\mathbf{y}} | \tilde{\mathbf{X}}^T \mathbf{w}, \sigma^2 \mathbf{I}_m) \mathcal{N}(\mathbf{w} | \mathbf{w}_N, \sigma^2 \mathbf{V}_N) d\mathbf{w} \right] \text{IG}(\sigma^2 | a_N, b_N) d\sigma^2 \\
&= \int \mathcal{N}(\tilde{\mathbf{y}} | \tilde{\mathbf{X}}^T \mathbf{w}_N, \sigma^2 (\mathbf{I}_m + \tilde{\mathbf{X}}^T \mathbf{V}_N \tilde{\mathbf{X}})) \text{IG}(\sigma^2 | a_N, b_N) d\sigma^2 \quad (\text{by Eq 4.126}) \\
&= \int \mathcal{N}(\tilde{\mathbf{y}} | \tilde{\mathbf{X}}^T \mathbf{w}_N, (1/\sigma^2) (\mathbf{I}_m + \tilde{\mathbf{X}}^T \mathbf{V}_N \tilde{\mathbf{X}})) \text{Ga}(1/\sigma^2 | a_N, b_N) d\sigma^2 \\
&= \mathcal{T}\left(\tilde{\mathbf{y}} \middle| \tilde{\mathbf{X}}^T \mathbf{w}_N, \frac{b_N}{a_N} (\mathbf{I}_m + \tilde{\mathbf{X}}^T \mathbf{V}_N \tilde{\mathbf{X}}), 2a_N\right) \quad (\text{by Eq 11.61})
\end{aligned}$$