

Gaussian Process

Eugenio Fella, Theivan Pasupathipillai, Matteo Pedrazzi,
Gaetano Ricucci, Carlo Sgorlon Gaiatto

June 28, 2023

1 Bayesian linear regression

1.1 Linear regression model

Let's consider a linear regression model:

$$\mathbf{y} = \mathbf{X}\mathbf{w} \quad (1)$$

where

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_N \end{pmatrix} \quad \text{with} \quad \mathbf{x}_i, \mathbf{w} \in \mathbb{R}^{d+1} \quad (2)$$

The posterior probability of the regression coefficients is:

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{y})} \quad (3)$$

1.2 Likelihood

Let's assume that each observation y_i is drawn from a Gaussian distribution with unknown mean $\sum_j X_{ij}w_j$ and fixed variance σ_i^2 . Therefore the likelihood is:

$$p(\mathbf{y} | \mathbf{w}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \mathbf{\Sigma}) \quad (4)$$

where

$$\mathbf{\Sigma} = \begin{pmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_N^2 \end{pmatrix} \quad (5)$$

Result 1.1. *The product of N independent Gaussian distributions is a multivariate Gaussian distribution.*

Proof.

$$\begin{aligned}
\log p(\mathbf{y} | \mathbf{w}) &= \log \prod_{i=1}^N \mathcal{N} \left(y_i | \sum_j X_{ij} w_j, \sigma_i^2 \right) \\
&= \sum_{i=1}^N \log \mathcal{N} \left(y_i | \sum_j X_{ij} w_j, \sigma_i^2 \right) \\
&= -\frac{1}{2} \left[\frac{\left(y_1 - \sum_j X_{1j} w_j \right)^2}{\sigma_1^2} + \frac{\left(y_2 - \sum_j X_{2j} w_j \right)^2}{\sigma_2^2} + \dots + \frac{\left(y_N - \sum_j X_{Nj} w_j \right)^2}{\sigma_N^2} \right] + \text{const} \\
&= -\frac{1}{2} \begin{pmatrix} y_1 - \sum_j X_{1j} w_j \\ y_2 - \sum_j X_{2j} w_j \\ \vdots \\ y_N - \sum_j X_{Nj} w_j \end{pmatrix}^T \begin{pmatrix} \frac{1}{\sigma_1^2} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2^2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\sigma_N^2} \end{pmatrix} \begin{pmatrix} y_1 - \sum_j X_{1j} w_j \\ y_2 - \sum_j X_{2j} w_j \\ \vdots \\ y_N - \sum_j X_{Nj} w_j \end{pmatrix} + \text{const} \\
&= -\frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T \mathbf{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{w}) + \text{const} \\
&= \log \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \mathbf{\Sigma})
\end{aligned}$$

□

1.3 Prior

The prior is chosen to be a conjugate prior, namely it is also Gaussian:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mu, \mathbf{V}_0) \quad (6)$$

1.4 Joint probability

The joint probability is the product of two multivariate Gaussian distributions that can be shown to be again Gaussian:

$$p(\mathbf{w}, \mathbf{y}) = p(\mathbf{y} | \mathbf{w}) p(\mathbf{w}) = \mathcal{N}(\mathbf{w}, \mathbf{y} | \mu', \mathbf{\Sigma}') \quad (7)$$

where

$$\mu' = \begin{pmatrix} \mu \\ \mathbf{X}\mu \end{pmatrix} \quad \mathbf{\Sigma}' = \begin{pmatrix} \mathbf{V}_0 & \mathbf{V}_0 \mathbf{X}^T \\ \mathbf{X} \mathbf{V}_0 & \mathbf{\Sigma} + \mathbf{X} \mathbf{V}_0 \mathbf{X}^T \end{pmatrix} \quad (8)$$

Result 1.2. *The product of two multivariate Gaussian distributions is again Gaussian.*

Proof.

$$\log p(\mathbf{w}, \mathbf{y}) = \log p(\mathbf{y} | \mathbf{w}) p(\mathbf{w}) \quad (9)$$

$$= -\frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T \mathbf{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{w}) - \frac{1}{2} (\mathbf{w} - \mu)^T \mathbf{V}_0^{-1} (\mathbf{w} - \mu) + \text{const} \quad (10)$$

$$= -\frac{1}{2} \left[\mathbf{y}^T \mathbf{\Sigma}^{-1} \mathbf{y} - \mathbf{y}^T \mathbf{\Sigma}^{-1} \mathbf{X}\mathbf{w} - (\mathbf{X}\mathbf{w})^T \mathbf{\Sigma}^{-1} \mathbf{y} + (\mathbf{X}\mathbf{w})^T \mathbf{\Sigma}^{-1} \mathbf{X}\mathbf{w} \right] \quad (11)$$

$$- \frac{1}{2} \left[\mathbf{w}^T \mathbf{V}_0^{-1} \mathbf{w} - \mathbf{w}^T \mathbf{V}_0^{-1} \mu - \mu^T \mathbf{V}_0^{-1} \mathbf{w} + \mu^T \mathbf{V}_0^{-1} \mu \right] + \text{const} \quad (12)$$

Let's separate second order terms with respect to the linear ones and collect everything that does not depend on \mathbf{w} or \mathbf{y} into *const*. Moreover note that:

$$(\mathbf{X}\mathbf{w})^T = \mathbf{w}^T \mathbf{X}^T \quad (13)$$

$$\mu^T \mathbf{V}_0^{-1} \mathbf{w} = (\mathbf{w}^T \mathbf{V}_0^{-1} \mu)^T = \mathbf{w}^T \mathbf{V}_0^{-1} \mu \quad (14)$$

Therefore we have:

$$\log p(\mathbf{w}, \mathbf{y}) = -\frac{1}{2} \mathbf{w}^T (\mathbf{V}_0^{-1} + \mathbf{X}^T \Sigma^{-1} \mathbf{X}) \mathbf{w} + \frac{1}{2} \mathbf{w}^T \mathbf{X}^T \Sigma^{-1} \mathbf{y} + \frac{1}{2} \mathbf{y}^T \Sigma^{-1} \mathbf{X} \mathbf{w} - \frac{1}{2} \mathbf{y}^T \Sigma^{-1} \mathbf{y} + \mathbf{w}^T \mathbf{V}_0^{-1} \mu + \text{const} \quad (15)$$

$$= -\frac{1}{2} \begin{pmatrix} \mathbf{w} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \mathbf{V}_0^{-1} + \mathbf{X}^T \Sigma^{-1} \mathbf{X} & -\mathbf{X}^T \Sigma^{-1} \\ -\Sigma^{-1} \mathbf{X} & \Sigma^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ \mathbf{y} \end{pmatrix} + \begin{pmatrix} \mathbf{w} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \mathbf{V}_0^{-1} \mu \\ \mathbf{0} \end{pmatrix} + \text{const} \quad (16)$$

To determine the covariance matrix, let's define:

$$\Sigma'^{-1} = \begin{pmatrix} \mathbf{V}_0^{-1} + \mathbf{X}^T \Sigma^{-1} \mathbf{X} & -\mathbf{X}^T \Sigma^{-1} \\ -\Sigma^{-1} \mathbf{X} & \Sigma^{-1} \end{pmatrix} \quad (17)$$

To get Σ' we use the following result for the inverse of a partitioned matrix:

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix} \quad (18)$$

where

$$\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} \quad (19)$$

Doing the (simple) math we get:

$$\Sigma' = \begin{pmatrix} \mathbf{V}_0 & \mathbf{V}_0 \mathbf{X}^T \\ \mathbf{X} \mathbf{V}_0 & \Sigma + \mathbf{X} \mathbf{V}_0 \mathbf{X}^T \end{pmatrix} \quad (20)$$

To determine the mean, note that we can write the exponent of a multivariate Gaussian distribution in a general form involving a quadratic term, a linear term and a constant term:

$$-\frac{1}{2}(\mathbf{z} - \mu')^T \Sigma'^{-1}(\mathbf{z} - \mu') = -\frac{1}{2} \mathbf{z}^T \Sigma'^{-1} \mathbf{z} + \mathbf{z}^T \Sigma'^{-1} \mu' + \text{const} \quad (21)$$

Thus, comparing with our result we can get the mean μ' :

$$\Sigma'^{-1} \mu' = \begin{pmatrix} \mathbf{V}_0^{-1} \mu \\ \mathbf{0} \end{pmatrix} \quad (22)$$

multiplying by Σ'

$$\mu' = \Sigma' \begin{pmatrix} \mathbf{V}_0^{-1} \mu \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mu \\ \mathbf{X} \mu \end{pmatrix} \quad (23)$$

□

1.5 Posterior

The posterior can be found by starting with the joint probability and showing that the conditional probability is also Gaussian:

$$p(\mathbf{w} \mid \mathbf{y}) = \mathcal{N}(\mathbf{w} \mid \mu_{\mathbf{w}|\mathbf{y}}, \Sigma_{\mathbf{w}|\mathbf{y}}) \quad (24)$$

where

$$\mu_{\mathbf{w}|\mathbf{y}} = \mu + \mathbf{V}_0 \mathbf{X}^T (\Sigma + \mathbf{X} \mathbf{V}_0 \mathbf{X}^T)^{-1} (\mathbf{y} - \mathbf{X} \mu) \quad (25)$$

$$\Sigma_{\mathbf{w}|\mathbf{y}} = \mathbf{V}_0 - \mathbf{V}_0 \mathbf{X}^T (\Sigma + \mathbf{X} \mathbf{V}_0 \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{V}_0 \quad (26)$$

Result 1.3. *The conditional distribution derived from a multivariate Gaussian distribution is again Gaussian.*

Proof. TODO □

2 Bayesian nonparametric linear regression

This approach allows to make predictions at new locations:

$$\mathbf{y}^* = \mathbf{X}^* \mathbf{w} \quad (27)$$

It can be shown that the distribution of predictions given the observation is again Gaussian and independent of the regression coefficients:

$$p(\mathbf{y}^* \mid \mathbf{y}) = \mathcal{N}(\mathbf{y}^* \mid \mu_{\mathbf{y}^*|\mathbf{y}}, \Sigma_{\mathbf{y}^*|\mathbf{y}}) \quad (28)$$

where

$$\mu_{\mathbf{y}^*|\mathbf{y}} = \mathbf{X}^* \mu_{\mathbf{y}^*|\mathbf{y}} = \mathbf{X}^* \mu + \mathbf{X}^* \mathbf{V}_0 \mathbf{X}^T (\Sigma + \mathbf{X} \mathbf{V}_0 \mathbf{X}^T)^{-1} (\mathbf{y} - \mathbf{X} \mu) \quad (29)$$

$$\Sigma_{\mathbf{y}^*|\mathbf{y}} = \mathbf{X}^* \Sigma_{\mathbf{y}^*|\mathbf{y}} \mathbf{X}^{*T} = \mathbf{X}^* \mathbf{V}_0 \mathbf{X}^{*T} - \mathbf{X}^* \mathbf{V}_0 \mathbf{X}^T (\Sigma + \mathbf{X} \mathbf{V}_0 \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{V}_0 \mathbf{X}^{*T} \quad (30)$$

Result 2.1. *The distribution of a linear transformation of Gaussian distributed random variable is again Gaussian*

Proof.

$$\mu_{\mathbf{y}^*|\mathbf{y}} = \mathbb{E}[\mathbf{y}^*] = \mathbb{E}[\mathbf{X}^* \mathbf{w}] = \mathbf{X}^* \mathbb{E}[\mathbf{w}] = \mathbf{X}^* \mu_{\mathbf{w}|\mathbf{y}} \quad (31)$$

$$\Sigma_{\mathbf{y}^*|\mathbf{y}} = \mathbb{E}[(\mathbf{y}^* - \mathbb{E}[\mathbf{y}^*])(\mathbf{y}^* - \mathbb{E}[\mathbf{y}^*])^T] = \mathbf{X}^* \mathbb{E}[(\mathbf{w} - \mathbb{E}[\mathbf{w}])(\mathbf{w} - \mathbb{E}[\mathbf{w}])^T] \mathbf{X}^{*T} = \mathbf{X}^* \Sigma_{\mathbf{w}|\mathbf{y}} \mathbf{X}^{*T} \quad (32)$$

□

2.1 Kernel trick

To increase the expressiveness of the model it is common to use a non linear feature mapping:

$$\Phi = \phi(\mathbf{X}) \quad (33)$$

The kernel trick emerges from the observation that the inner products $\Phi \mathbf{V}_0 \Phi^T$ can be equivalently computed by evaluating the corresponding kernel function k for all pairs to form the matrix $\mathbf{K}_{\mathbf{X}\mathbf{X}}$:

$$\mathbf{K}_{\mathbf{X}\mathbf{X}}(i, j) = k(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^T \mathbf{V}_0 \Phi(\mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle_{\mathbf{V}_0} \quad (34)$$

The kernel trick allows us to specify an intuitive similarity between pairs of points, rather than a feature map Φ , which in practice can be hard to define.

2.2 Gaussian Process

Applying the feature map and the kernel trick describe before, we obtain a Gaussian Process:

$$p(\mathbf{y}^* | \mathbf{y}) = \mathcal{N}(\mathbf{y}^* | \mu_{\mathbf{y}^*|\mathbf{y}}, \Sigma_{\mathbf{y}^*|\mathbf{y}}) \quad (35)$$

where

$$\mu_{\mathbf{y}^*|\mathbf{y}} = \Phi^* \mu + \mathbf{K}_{\mathbf{X}^*\mathbf{X}} \hat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} (\mathbf{y} - \Phi \mu) \quad (36)$$

$$\Sigma_{\mathbf{y}^*|\mathbf{y}} = \mathbf{K}_{\mathbf{X}^*\mathbf{X}^*} - \mathbf{K}_{\mathbf{X}^*\mathbf{X}} \hat{\mathbf{K}}_{\mathbf{X}\mathbf{X}}^{-1} \mathbf{K}_{\mathbf{X}\mathbf{X}^*} \quad (37)$$