

An Introduction to Gaussian Process Regression

P.L.Green
Institute for Risk and Uncertainty
School of Engineering
University of Liverpool
Liverpool, L69 7ZF
United Kingdom

p.l.green@liverpool.ac.uk
<https://www.liverpool.ac.uk/engineering/staff/peter-green/>

September 3, 2018

1 Introduction

This document is supposed to provide a tutorial-style introduction to Gaussian Process (GP) regression. It comes with some example Gaussian Process code, written in Python, that is free and should be relatively easy to use. The code is by no-means optimised, rather, it is supposed to be easy to read and easy to understand for someone who is new to GP regression.

2 The Problem

Say we believe that a system's input, \mathbf{x} , is related to its output via a relationship $f(\mathbf{x})$. The function, f , is currently unknown. To help us find f , we conduct a series of experiments where, for N different \mathbf{x} values, we observe the system's output (denoted y). Our experiment therefore yields the data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$. As with any experiment, our observations are corrupted by some errors. In the following, we treat these errors as being Gaussian and we assume the following 'noise model':

$$y_i = f(\mathbf{x}_i) + \epsilon, \quad \epsilon \sim \mathcal{N}(\epsilon; 0, \sigma^2) \quad (1)$$

that is, we assume that each observation y_i is equal to the function, f , evaluated at input, \mathbf{x}_i , but corrupted with Gaussian noise of variance σ^2 .

3 Bayesian Linear Regression

Once we have our data, $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$, we can use it to try and infer the function f . Say we assume that f takes the form

$$f(\mathbf{x}_i) = \boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}_i) \quad (2)$$

where $\boldsymbol{\theta}$ is a vector of parameters which we need to estimate and $\boldsymbol{\phi}$ is a 'basis function'. For example, if we are considering univariate inputs and set

$$\boldsymbol{\theta} = (\theta_0, \theta_1, \theta_2)^T, \quad \boldsymbol{\phi}(x_i) = (1, x_i, x_i^2)^T \quad (3)$$

then we would be trying to fit a quadratic function to our data. Equation (2) is convenient because it is a linear function of the parameters we need to estimate - this usually makes it easier to find closed-form expression for the ‘optimum choice’ of $\boldsymbol{\theta}$.

Defining

$$\mathbf{f} = (f_1, f_2, \dots, f_N)^T, \quad \text{where} \quad f_i \equiv f(\mathbf{x}_i) \quad (4)$$

we can write

$$f_i = \sum_k \theta_k \phi_k(\mathbf{x}_i). \quad (5)$$

Furthermore, defining the matrix

$$\boldsymbol{\Phi} = \begin{bmatrix} \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \dots \\ \phi_1(\mathbf{x}_2) & \phi_2(\mathbf{x}_2) & \\ \vdots & & \ddots \end{bmatrix} \quad (6)$$

lets us write $\mathbf{f} = \boldsymbol{\Phi} \boldsymbol{\theta}$.

At this point, a Bayesian approach would typically involve defining a prior probability distribution, $p(\boldsymbol{\theta})$, that is supposed to represent our prior belief in $\boldsymbol{\theta}$. We would then write down Bayes theorem:

$$p(\boldsymbol{\theta} | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) \quad (7)$$

where $\mathbf{y} = (y_1, \dots, y_N)$ is our observation data and the likelihood, $p(\mathbf{y} | \boldsymbol{\theta})$, is formed based on our noise model (equation (1) in this case). This approach gives us an expression for the posterior distribution, $p(\boldsymbol{\theta} | \mathbf{y})$, which allows us to analyse the probability that our regression parameters can take particular values, conditional on our observation data.

If we adopt such a Bayesian approach and happen choose the prior over our parameters to be

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \mathbf{0}, \alpha^{-1} \mathbf{I}) \quad (8)$$

then, after ‘feeding it through’ our model ($\mathbf{f} = \boldsymbol{\Phi} \boldsymbol{\theta}$), it follows that the covariance matrix over \mathbf{f} is

$$\text{Cov}[\mathbf{f}] = \text{E}[\mathbf{f} \mathbf{f}^T] = \boldsymbol{\Phi} \text{E}[\boldsymbol{\theta} \boldsymbol{\theta}^T] \boldsymbol{\Phi}^T = \alpha^{-1} \boldsymbol{\Phi} \boldsymbol{\Phi}^T \quad (9)$$

(which we’ll write just as $\text{Cov}[\mathbf{y}] = \mathbf{K}$ from now on). This result means that, regardless of the inputs used in our experiments, saying that our prior over $\boldsymbol{\theta}$ is Gaussian means we automatically assume a Gaussian prior over \mathbf{f} :

$$p(\mathbf{f}) \propto \exp\left(-\frac{1}{2} \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f}\right). \quad (10)$$

where the individual elements of \mathbf{K} are given by

$$K_{ij} = \alpha^{-1} \sum_i \phi_i(\mathbf{x}_i) \phi_i(\mathbf{x}_j) \quad (11)$$

4 Gaussian Process Regression

In the previous section we defined a function f that was linear in our regression parameters, $\boldsymbol{\theta}$, before defining a Gaussian prior over $\boldsymbol{\theta}$. This led to a Gaussian prior over \mathbf{f} , whose covariance matrix is given by equation (11). Key to the construction of GPs is the idea that, as we know that the aforementioned Bayesian approach will always lead to a Gaussian distribution over \mathbf{f} , we may as well define $p(\mathbf{f})$ directly and skip the part where we define a basis function and a prior over $\boldsymbol{\theta}$. Specifically, we simply write

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}) \quad (12)$$

where the matrix \mathbf{K} can be anything we like (so long as it is a valid covariance matrix). For this reason, GPs are sometimes described as a method that ‘assigns a prior over the function values, \mathbf{f} ’. The matrix \mathbf{K} (the ‘gram matrix’) is usually defined using a *kernel function*, k , such that $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$. A common choice of kernel function, for example, is

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{1}{2l^2}(\mathbf{x}_i - \mathbf{x}_j)^T(\mathbf{x}_i - \mathbf{x}_j)\right) \quad (13)$$

where we have introduced l , a ‘length scale’ parameter. So how do the kernel function (equation (13)) and length scale parameter alter our prior belief about the function f ? Consider the case where we have a vector of 100 1D inputs, that are evenly spaced between 0 and 10. We choose the kernel described by equation (13) and plot the matrix \mathbf{K} for $l = 10, 1$ and 0.1 . The resulting gram matrices are shown in Figure 1.

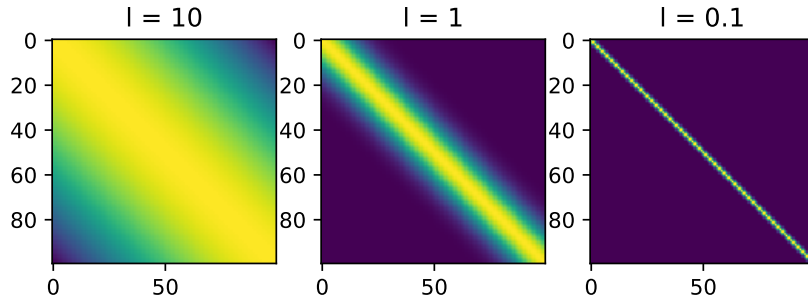


Figure 1: Plotting the gram matrix, \mathbf{K} , for different values of length scale, l , using the kernel function described by equation (13).

Figure 1 shows us that, by using larger values of l , we are introducing greater correlation between our inputs. So what does this signify? Generating 100 samples of \mathbf{f} , from $\mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K})$ for different values of l , we realise the results shown in Figure 2.

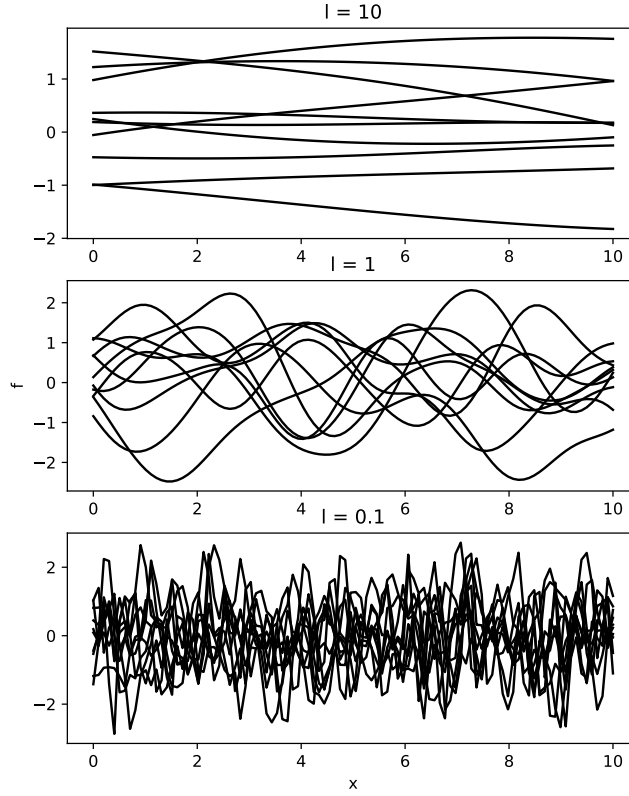


Figure 2: 100 samples of $\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \mathbf{0}, \mathbf{K})$, for different values of length scale, l .

Figure 2 shows us that, by increasing the correlation between our inputs, we are inducing *smoothness* in our prior function f . As we shall see shortly, ‘smoothness’ is quite an intuitive and useful quantity when it comes to performing regression.

So, to recap, we have our prior

$$p(\mathbf{f}) = \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}) \quad (14)$$

and we have our noise model, equation (1). As we have assumed that the noise affecting our observations is independent for each data point, it follows that

$$p(\mathbf{y} | \mathbf{f}) = \mathcal{N}(\mathbf{y}; \mathbf{f}, \mathbf{I}\sigma^2) \quad (15)$$

We can then marginalise over \mathbf{f} , to find the probability of witnessing \mathbf{y} :

$$p(\mathbf{y}) = \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f}) d\mathbf{f} \quad (16)$$

Fortunately, we can find a closed-form expression for $p(\mathbf{y})$ using some generic expressions for Gaussian distributions (see Appendix A). After some manipulation, we find that

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{C}) \quad (17)$$

where

$$\mathbf{C} = \mathbf{I}\sigma^2 + \mathbf{K} \implies C_{ij} = \sigma^2\delta_{ij} + k(\mathbf{x}_i, \mathbf{x}_j). \quad (18)$$

Now, what we really want to do is analyse the probability that, in response to a new input, \mathbf{x}_* , we would observe the response y_* . Using the relationship we derived in equation (18), we can write

$$p(\mathbf{y}, y_*) = \mathcal{N}\left(\begin{pmatrix} \mathbf{y} \\ y_* \end{pmatrix}; \begin{pmatrix} \mathbf{0} \\ 0 \end{pmatrix}, \begin{bmatrix} \mathbf{C} & \mathbf{k} \\ \mathbf{k}^T & c \end{bmatrix}\right) \quad (19)$$

where the vector \mathbf{k} is defined as

$$\mathbf{k}_n = k(\mathbf{x}_n, \mathbf{x}_*), \quad n = 1, \dots, N \quad (20)$$

and

$$c = \sigma^2 + k(\mathbf{x}_*, \mathbf{x}_*). \quad (21)$$

By once again exploiting our generic expressions for Gaussian distributions, we are then able to find that

$$p(y_* | \mathbf{y}) = \mathcal{N}(y_*; m_*, \sigma_*^2) \quad (22)$$

where

$$m_* = \mathbf{k}^T \mathbf{C}^{-1} \mathbf{y}, \quad \sigma_*^2 = c - \mathbf{k}^T \mathbf{C}^{-1} \mathbf{k}. \quad (23)$$

5 Example: 1D Interpolation

Here we have generated some ‘observations’ according to:

$$y_n = \sin x_n + \epsilon, \quad \epsilon \sim \mathcal{N}(\epsilon; 0, \sigma^2) \quad (24)$$

where $\sigma = 0.2$. We are going to use a Gaussian process to ‘fill the gaps’ between the data that we have gathered. Using the kernel function in equation (13), Figure 3 shows what happens when we set $l = 0.5, 1$ and 5 respectively.

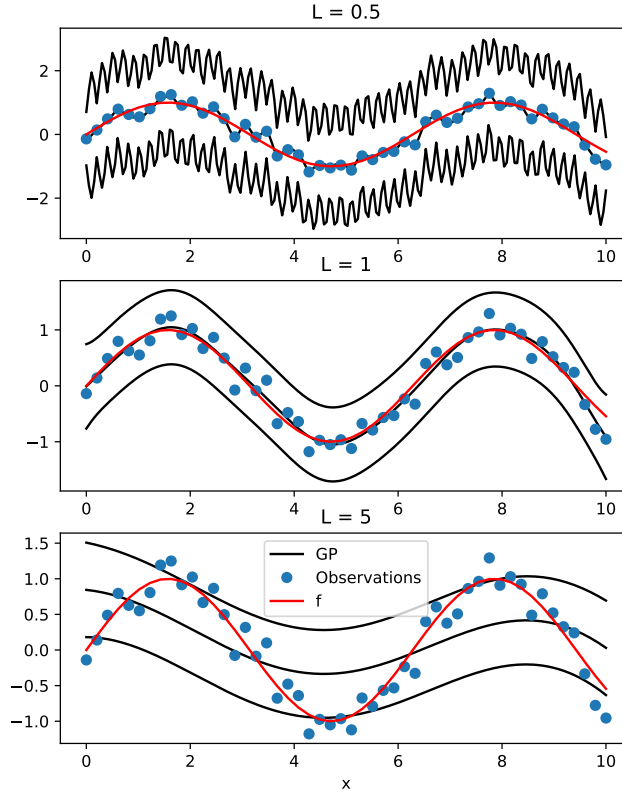


Figure 3: Performing Gaussian Process regression for different values of the length scale, l . Gaussian Process confidence bounds are 3 standard deviations from the mean.

From Figure 3 it should be obvious that, when $l = 0.5$, our interpolating function isn't smooth enough. In fact, it passes through every observation almost exactly - by mistakenly fitting to the noise that is in the data, it has formed a poor representation of the true function, f . This phenomenon is sometimes referred to as 'overfitting'. When $l = 1$, the fit appears sensible while, when $l = 5$, our fit seems overly smooth (in other words, it is not able to react quickly enough to changes in the input, x).

6 Gaussian Process Learning

A useful aspect of GPs is that, unlike our linear regression approach where we had to begin by assuming a functional form for f (equation (2)), GPs are not so tightly constrained. They are sometimes referred to as a form of *non-parametric regression* as they are not limited to only consider one parametric family of regression functions. The GP does, however, have some parameters which require tuning (l and σ is our case). These, perhaps a little confusingly, are often referred to as *hyperparameters*.

Writing our hyperparameters together as $\boldsymbol{\theta}_{\text{hyp}}$ such that $\boldsymbol{\theta}_{\text{hyp}} = \{l, \sigma\}$ then, from equation (17) (derived in the previous section), we know that the probability of witnessing our observations conditional on $\boldsymbol{\theta}_{\text{hyp}}$ is given by

$$p(\mathbf{y} \mid \boldsymbol{\theta}_{\text{hyp}}) = \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{C}) = ((2\pi)^N |\mathbf{C}|)^{-1/2} \exp\left(-\frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y}\right) \quad (25)$$

such that:

$$\log p(\mathbf{y} | \boldsymbol{\theta}_{\text{hyp}}) = -\frac{1}{2} (N \log 2\pi + \log |\mathbf{C}| + \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y}). \quad (26)$$

Finding the maximum of equation (26), we can therefore find the maximum likelihood estimates of $\boldsymbol{\theta}_{\text{hyp}}$. This is typically what happens during the ‘learning process’ for a GP. Conveniently, to help our optimisation routine, we are usually able to evaluate the gradient of the log-likelihood with respect to our hyperparameters using closed-form expressions. For example, noting that:

$$\frac{\partial \log |\mathbf{C}|}{\partial \theta_{\text{hyp},i}} = \text{Tr} \left(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_{\text{hyp},i}} \right), \quad \frac{\partial \mathbf{C}^{-1}}{\partial \theta_{\text{hyp},i}} = -\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_{\text{hyp},i}} \mathbf{C}^{-1} \quad (27)$$

then it can be shown that

$$\frac{\partial}{\partial \theta_{\text{hyp},i}} \log p(\mathbf{y} | \boldsymbol{\theta}_{\text{hyp}}) = -\frac{1}{2} \left(\text{Tr} \left(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_{\text{hyp},i}} \right) - \mathbf{y}^T \mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_{\text{hyp},i}} \mathbf{C}^{-1} \mathbf{y} \right). \quad (28)$$

7 Example Python Code

Code accompanying this tutorial is written in Python (version 2), is heavily commented and is designed to be readable. The module ‘LIRU_GP’ contains all of the functions needed to train a GP for regression. The scripts ‘1D_regression_example’ and ‘2D_regression_example’, as the names imply, use LIRU_GP to perform regression on a 1D and 2D problem respectively. ‘1D_regression_example’ will plot results similar to those shown in Figure 4 while ‘2D_regression_example’ will produce results similar to those shown in Figure 5.

Note that, for more optimised and generic code, a variety of GP packages are offered though Python (see <https://sheffielddml.github.io/GPy/> for example). As stated previously, our code is intended to be readable, rather than generic.

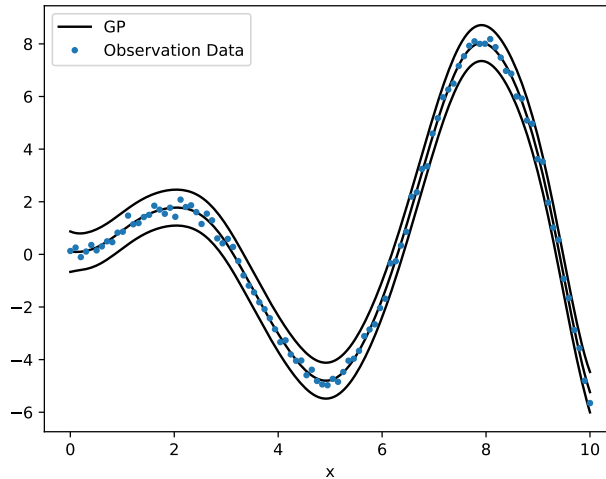


Figure 4: Example output of ‘1D_regression_example’.

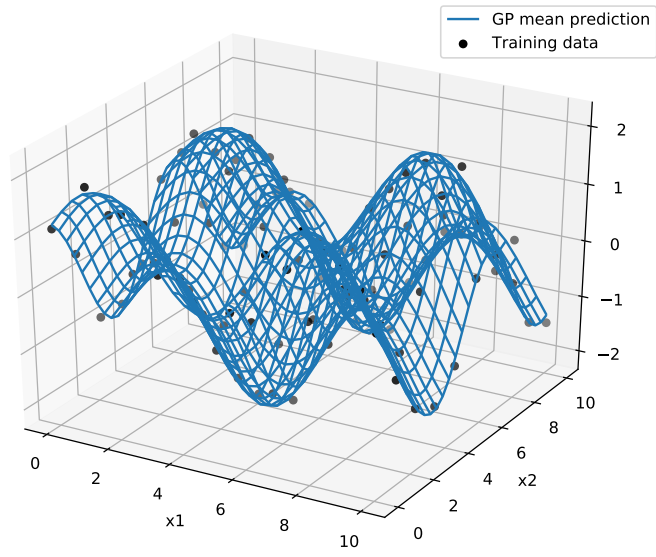


Figure 5: Example output of ‘2D_regression_example’.

A Properties of Gaussians

Before we start we should note some useful properties. All this stuff is proved in Appendix B (which is just a slightly longer version of what you can find in Bishop [1]).

Say we have a partitioned Gaussian $p(\mathbf{x}) \propto \exp(-\frac{1}{2}\phi)$ where

$$\phi = ((\mathbf{x}_a - \boldsymbol{\mu}_a)^T, (\mathbf{x}_b - \boldsymbol{\mu}_b)^T) \begin{bmatrix} \mathbf{A}_{aa} & \mathbf{A}_{ab} \\ \mathbf{A}_{ba} & \mathbf{A}_{bb} \end{bmatrix} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix}. \quad (29)$$

$$= ((\mathbf{x}_a - \boldsymbol{\mu}_a)^T, (\mathbf{x}_b - \boldsymbol{\mu}_b)^T) \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix}^{-1} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix}. \quad (30)$$

Written in terms of precision matrices we have the properties that:

Distribution	Precision	Mean
$p(\mathbf{x}_a \mathbf{x}_b)$	\mathbf{A}_{aa}	$\mathbf{A}_{aa}^{-1}(\mathbf{A}_{aa}\boldsymbol{\mu}_a - \mathbf{A}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b))$
$p(\mathbf{x}_b \mathbf{x}_a)$	\mathbf{A}_{bb}	$\mathbf{A}_{bb}^{-1}(\mathbf{A}_{bb}\boldsymbol{\mu}_b - \mathbf{A}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a))$
$p(\mathbf{x}_a)$	$\mathbf{A}_{aa} - \mathbf{A}_{ab}\mathbf{A}_{bb}^{-1}\mathbf{A}_{ba}$	$\boldsymbol{\mu}_a$
$p(\mathbf{x}_b)$	$\mathbf{A}_{bb} - \mathbf{A}_{ba}\mathbf{A}_{aa}^{-1}\mathbf{A}_{ab}$	$\boldsymbol{\mu}_b$

Written in terms of covariance matrices these can also be written as:

Distribution	Covariance	Mean
$p(\mathbf{x}_a \mathbf{x}_b)$	$\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}$	$\boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)$
$p(\mathbf{x}_b \mathbf{x}_a)$	$\boldsymbol{\Sigma}_{bb} - \boldsymbol{\Sigma}_{ba}\boldsymbol{\Sigma}_{aa}^{-1}\boldsymbol{\Sigma}_{ab}$	$\boldsymbol{\mu}_b + \boldsymbol{\Sigma}_{ba}\boldsymbol{\Sigma}_{aa}^{-1}(\mathbf{x}_a - \boldsymbol{\mu}_a)$
$p(\mathbf{x}_a)$	$\boldsymbol{\Sigma}_{aa}$	$\boldsymbol{\mu}_a$
$p(\mathbf{x}_b)$	$\boldsymbol{\Sigma}_{bb}$	$\boldsymbol{\mu}_b$

Finally, for the case where

$$p(\mathbf{x}_a) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{A}^{-1}), \quad p(\mathbf{x}_b|\mathbf{x}_a) = \mathcal{N}(\mathbf{H}\mathbf{x}_a + \mathbf{b}, \mathbf{L}^{-1}) \quad (31)$$

then we have

Distribution	Covariance	Mean
$p(\mathbf{x}_a \mathbf{x}_b)$	$(\mathbf{A} + \mathbf{H}^T\mathbf{L}\mathbf{H})^{-1}$	$(\mathbf{A} + \mathbf{H}^T\mathbf{L}\mathbf{H})^{-1}(\mathbf{A}\boldsymbol{\mu} + \mathbf{H}^T\mathbf{L}(\mathbf{x}_b - \mathbf{b}))$
$p(\mathbf{x}_b)$	$\mathbf{L}^{-1} + \mathbf{H}\mathbf{A}^{-1}\mathbf{H}^T$	$\mathbf{H}\boldsymbol{\mu} + \mathbf{b}$

These expressions are generic - they will simplify down a lot for the examples shown here.

B Deriving the properties of Gaussian distributions

B.1 Partitioned Gaussian

We have a Gaussian distribution $p(\mathbf{x}) = \exp(-\frac{1}{2}\phi(\mathbf{x}))$ where

$$\phi = (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{A} (\mathbf{x} - \boldsymbol{\mu}) \quad (32)$$

where \mathbf{A} is the precision matrix. We are going to partition this up:

$$\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)^T, \quad \boldsymbol{\mu} = (\boldsymbol{\mu}_a, \boldsymbol{\mu}_b)^T, \quad \mathbf{A} = \begin{bmatrix} \mathbf{A}_{aa} & \mathbf{A}_{ab} \\ \mathbf{A}_{ba} & \mathbf{A}_{bb} \end{bmatrix} \quad (33)$$

such that $p(\mathbf{x}) = p(\mathbf{x}_a, \mathbf{x}_b) \propto \exp(-\frac{1}{2}\phi)$ where

$$\phi = ((\mathbf{x}_a - \boldsymbol{\mu}_a)^T, (\mathbf{x}_b - \boldsymbol{\mu}_b)^T) \begin{bmatrix} \mathbf{A}_{aa} & \mathbf{A}_{ab} \\ \mathbf{A}_{ba} & \mathbf{A}_{bb} \end{bmatrix} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix} \quad (34)$$

$$\begin{aligned} &= (\mathbf{x}_a - \boldsymbol{\mu}_a)^T \mathbf{A}_{aa} (\mathbf{x}_a - \boldsymbol{\mu}_a) + (\mathbf{x}_a - \boldsymbol{\mu}_a)^T \mathbf{A}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad + (\mathbf{x}_b - \boldsymbol{\mu}_b)^T \mathbf{A}_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) + (\mathbf{x}_b - \boldsymbol{\mu}_b)^T \mathbf{A}_{bb} (\mathbf{x}_b - \boldsymbol{\mu}_b). \end{aligned} \quad (35)$$

B.2 Conditional Gaussian

We are going to take $p(\mathbf{x}_a, \mathbf{x}_b)$ and find $p(\mathbf{x}_a|\mathbf{x}_b)$. This essentially involves taking equation (35) and treating all the \mathbf{x}_b terms as being constant (which we'll lose when the distribution is normalised). Consequently, defining C as 'some constant' and noting that

$$(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \mathbf{A}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) = (\mathbf{x}_b - \boldsymbol{\mu}_b)^T \mathbf{A}_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) \quad (36)$$

(as $\mathbf{A}_{ab}^T = \mathbf{A}_{ba}$) then we find that

$$\phi = (\mathbf{x}_a - \boldsymbol{\mu}_a)^T \mathbf{A}_{aa} (\mathbf{x}_a - \boldsymbol{\mu}_a) + 2(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \mathbf{A}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) + C \quad (37)$$

$$= \mathbf{x}_a^T \mathbf{A}_{aa} \mathbf{x}_a - 2\mathbf{x}_a^T (\mathbf{A}_{aa}\boldsymbol{\mu}_a - \mathbf{A}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)) + C \quad (38)$$

and so, by completing the square, we find that $p(\mathbf{x}_a|\mathbf{x}_b)$ is Gaussian with precision matrix

$$\mathbf{A}_{a|b} = \mathbf{A}_{aa} \quad (39)$$

and mean

$$\boldsymbol{\mu}_{a|b} = \mathbf{A}_{aa}^{-1} (\mathbf{A}_{aa}\boldsymbol{\mu}_a - \mathbf{A}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)) \quad (40)$$

$$= \boldsymbol{\mu}_a - \mathbf{A}_{aa}^{-1} \mathbf{A}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b). \quad (41)$$

By inverting the precision matrix (which we won't go into here) we can also show that

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (42)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba}. \quad (43)$$

Using the same basic method you can get similar expressions for $p(\mathbf{x}_b|\mathbf{x}_a)$.

B.3 Marginalised Gaussian

Now we will find $p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$. To do this we will collect together all of the \mathbf{x}_b terms in equation (35) and complete the square. By doing this, when we integrate over \mathbf{x}_b , we will simply get the normalising constant of a Gaussian which can be absorbed into the constant C . We then just collect together all the remaining \mathbf{x}_a terms and complete the square again.

To that end, all the terms involving \mathbf{x}_b in equation (35) are:

$$(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \mathbf{A}_{bb} (\mathbf{x}_b - \boldsymbol{\mu}_b) + (\mathbf{x}_b - \boldsymbol{\mu}_b)^T \mathbf{A}_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) \quad (44)$$

$$= \mathbf{x}_b^T \mathbf{A}_{bb} \mathbf{x}_b - 2\mathbf{x}_b^T (\mathbf{A}_{bb} \boldsymbol{\mu}_b - \mathbf{A}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a)) + C \quad (45)$$

$$= (\mathbf{x}_b - \mathbf{A}_{bb}^{-1} \mathbf{m})^T \mathbf{A}_{bb} (\mathbf{x}_b - \mathbf{A}_{bb}^{-1} \mathbf{m}) - \mathbf{m}^T \mathbf{A}_{bb}^{-1} \mathbf{m} + C \quad (46)$$

where

$$\mathbf{m} = \mathbf{A}_{bb} \boldsymbol{\mu}_b - \mathbf{A}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a). \quad (47)$$

Notice we've not included $\mathbf{m}^T \mathbf{A}_{bb}^{-1} \mathbf{m}$ in C as it contains some \mathbf{x}_a terms. Once we have integrated over \mathbf{x}_b we will be left with the \mathbf{x}_a terms left over from equation (35) and $-\mathbf{m}^T \mathbf{A}_{bb} \mathbf{m}$.

With

$$\mathbf{m}^T \mathbf{A}_{bb}^{-1} \mathbf{m} = (\mathbf{A}_{bb} \boldsymbol{\mu}_b - \mathbf{A}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a))^T \mathbf{A}_{bb}^{-1} (\mathbf{A}_{bb} \boldsymbol{\mu}_b - \mathbf{A}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a)) \quad (48)$$

$$\mathbf{x}_a^T \mathbf{A}_{ab} \mathbf{A}_{bb}^{-1} \mathbf{A}_{ba} \mathbf{x}_a - 2\mathbf{x}_a^T (\mathbf{A}_{ab} \boldsymbol{\mu}_b + \mathbf{A}_{ab} \mathbf{A}_{bb}^{-1} \mathbf{A}_{ba} \boldsymbol{\mu}_a) + C \quad (49)$$

and, as we have

$$\mathbf{x}_a^T \mathbf{A}_{aa} \mathbf{x}_a - 2\mathbf{x}_a^T \mathbf{A}_{aa} \boldsymbol{\mu}_a - 2\mathbf{x}_a^T \mathbf{A}_{ab} \boldsymbol{\mu}_b + C \quad (50)$$

from the left over \mathbf{x}_a terms, it follows that after integrating out \mathbf{x}_b we are left with

$$\mathbf{x}_a^T (\mathbf{A}_{aa} - \mathbf{A}_{ab} \mathbf{A}_{bb}^{-1} \mathbf{A}_{ba}) \mathbf{x}_a - 2\mathbf{x}_a^T (\mathbf{A}_{aa} - \mathbf{A}_{ab} \mathbf{A}_{bb}^{-1} \mathbf{A}_{ba}) \boldsymbol{\mu}_a + C. \quad (51)$$

Consequently, $p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$ has precision matrix

$$\mathbf{A}_a = \mathbf{A}_{aa} - \mathbf{A}_{ab} \mathbf{A}_{bb}^{-1} \mathbf{A}_{ba} \quad (52)$$

and mean

$$\mathbf{A}_a^{-1} (\mathbf{A}_{aa} - \mathbf{A}_{ab} \mathbf{A}_{bb}^{-1} \mathbf{A}_{ba}) \boldsymbol{\mu}_a = \boldsymbol{\mu}_a \quad (53)$$

which is nice. Furthermore, taking the inverse of the precision matrix (which we won't do here) we find that

$$\boldsymbol{\Sigma}_a = \boldsymbol{\Sigma}_{aa} \quad (54)$$

which is also nice.

B.4 Bayes' theorem for a linear Gaussian model

Say we have

$$p(\mathbf{x}_a) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{A}^{-1}), \quad p(\mathbf{x}_b | \mathbf{x}_a) = \mathcal{N}(\mathbf{H} \mathbf{x}_a + \mathbf{b}, \mathbf{L}^{-1}) \quad (55)$$

and we wish to evaluate Bayes' theorem:

$$p(\mathbf{x}_a | \mathbf{x}_b) \propto p(\mathbf{x}_b | \mathbf{x}_a) p(\mathbf{x}_a). \quad (56)$$

The easiest way to do is actually to write an expression for $p(\mathbf{x}_a, \mathbf{x}_b)$ before finding $p(\mathbf{x}_a | \mathbf{x}_b)$ using the stuff we derived earlier. To that end, writing $\mathbf{z}^T = (\mathbf{x}_1^T, \mathbf{x}_2^T)$ then we have $p(\mathbf{z}) \propto \exp(-\frac{1}{2} \phi)$ where

$$\phi = (\mathbf{x}_a - \boldsymbol{\mu})^T \mathbf{A} (\mathbf{x}_a - \boldsymbol{\mu}) + (\mathbf{x}_b - \mathbf{H} \mathbf{x}_a - \mathbf{b})^T \mathbf{L} (\mathbf{x}_b - \mathbf{H} \mathbf{x}_a - \mathbf{b}). \quad (57)$$

Collecting together all the second order terms we get

$$\mathbf{x}_a^T \mathbf{A} \mathbf{x}_a + \mathbf{x}_b^T \mathbf{L} \mathbf{x}_b - 2\mathbf{x}_a^T \mathbf{H}^T \mathbf{L} \mathbf{x}_b + \mathbf{x}_a^T \mathbf{H}^T \mathbf{L} \mathbf{H} \mathbf{x}_a \quad (58)$$

$$= \begin{pmatrix} \mathbf{x}_a^T & \mathbf{x}_b^T \end{pmatrix} \begin{bmatrix} \mathbf{A} + \mathbf{H}^T \mathbf{L} \mathbf{H} & -\mathbf{H}^T \mathbf{L} \\ -\mathbf{L} \mathbf{H}^T & \mathbf{L} \end{bmatrix} \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad (59)$$

and so \mathbf{z} has precision matrix

$$\mathbf{R} = \begin{bmatrix} \mathbf{A} + \mathbf{H}^T \mathbf{L} \mathbf{H} & -\mathbf{H}^T \mathbf{L} \\ -\mathbf{L} \mathbf{H}^T & \mathbf{L} \end{bmatrix} \quad (60)$$

and, it can be shown, covariance matrix

$$\mathbf{R}^{-1} = \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{A}^{-1} \mathbf{H}^T \\ \mathbf{H} \mathbf{A}^{-1} & \mathbf{L}^{-1} + \mathbf{H} \mathbf{A}^{-1} \mathbf{H}^T \end{bmatrix}. \quad (61)$$

Collecting the first order terms:

$$-2\mathbf{x}_a^T \mathbf{A} \boldsymbol{\mu} - 2\mathbf{x}_b^T \mathbf{L} \mathbf{b} + 2\mathbf{x}_a^T \mathbf{H}^T \mathbf{L} \mathbf{b} \quad (62)$$

$$= -2 \begin{pmatrix} \mathbf{x}_a^T & \mathbf{x}_b^T \end{pmatrix} \begin{pmatrix} \mathbf{A} \boldsymbol{\mu} - \mathbf{H}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} \quad (63)$$

and so

$$\mathbb{E}[\mathbf{z}] = \mathbf{R}^{-1} \begin{pmatrix} \mathbf{A} \boldsymbol{\mu} - \mathbf{H}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} \quad (64)$$

$$= \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{A}^{-1} \mathbf{H}^T \\ \mathbf{H} \mathbf{A}^{-1} & \mathbf{L}^{-1} + \mathbf{H} \mathbf{A}^{-1} \mathbf{H}^T \end{bmatrix} \begin{pmatrix} \mathbf{A} \boldsymbol{\mu} - \mathbf{H}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{H} \boldsymbol{\mu} + \mathbf{b} \end{pmatrix}. \quad (65)$$

From our expressions for marginalised Gaussians it is clear that

$$p(\mathbf{x}_b) = \mathcal{N}(\mathbf{H} \boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{H} \mathbf{A}^{-1} \mathbf{H}^T). \quad (66)$$

Now, making the substitutions

$$\mathbf{A}_{aa} \equiv \mathbf{A} + \mathbf{H}^T \mathbf{L} \mathbf{H}, \quad \boldsymbol{\mu}_a \equiv \boldsymbol{\mu}, \quad \boldsymbol{\mu}_b \equiv \mathbf{H} \boldsymbol{\mu} + \mathbf{b}, \quad \mathbf{A}_{ab} \equiv -\mathbf{H}^T \mathbf{L} \quad (67)$$

we can use our expressions for the conditional Gaussian to find $p(\mathbf{x}_a | \mathbf{x}_b)$. Firstly, note that the posterior precision matrix is

$$\mathbf{A} + \mathbf{H}^T \mathbf{L} \mathbf{H}. \quad (68)$$

Secondly, as

$$(\mathbf{A} + \mathbf{H}^T \mathbf{L} \mathbf{H}) \boldsymbol{\mu} + \mathbf{H}^T \mathbf{L} (\mathbf{x}_b - \mathbf{H} \boldsymbol{\mu} - \mathbf{b}) = \mathbf{A} \boldsymbol{\mu} + \mathbf{H}^T \mathbf{L} (\mathbf{x}_b - \mathbf{b}) \quad (69)$$

then mean is

$$(\mathbf{A} + \mathbf{H}^T \mathbf{L} \mathbf{H})^{-1} (\mathbf{A} \boldsymbol{\mu} + \mathbf{H}^T \mathbf{L} (\mathbf{x}_b - \mathbf{b})). \quad (70)$$

References

- [1] Christopher M Bishop. *Pattern Recognition and Machine Learning*, volume 53. 2013.