

Differentiating Gaussian Processes

Andrew McHutchon

April 17, 2013

1 First Order Derivative of the Posterior Mean

The posterior mean of a GP is given by,

$$\bar{f}_* = \mathbf{k}(\mathbf{x}_*, X) K(X, X)^{-1} \mathbf{y} \triangleq \mathbf{k}(\mathbf{x}_*, X) \boldsymbol{\alpha} \quad (1)$$

Only the $\mathbf{k}(\mathbf{x}_*, X)$ term depends on the test point \mathbf{x}_* , therefore to calculate the slope of the posterior mean we just need to differentiate the kernel. For the squared exponential covariance function the derivative of the kernel between \mathbf{x}_* and a training point \mathbf{x}_i is,

$$\begin{aligned} \frac{\partial k(\mathbf{x}_*, \mathbf{x}_i)}{\partial \mathbf{x}_*} &= \frac{\partial}{\partial \mathbf{x}_*} \left\{ \sigma_f^2 \exp \left(-\frac{1}{2} (\mathbf{x}_* - \mathbf{x}_i)^T \Lambda^{-1} (\mathbf{x}_* - \mathbf{x}_i) \right) \right\} \\ &= \frac{\partial}{\partial \mathbf{x}_*} \left\{ -\frac{1}{2} (\mathbf{x}_* - \mathbf{x}_i)^T \Lambda^{-1} (\mathbf{x}_* - \mathbf{x}_i) \right\} k(\mathbf{x}_*, \mathbf{x}_i) \\ &= -\Lambda^{-1} (\mathbf{x}_* - \mathbf{x}_i) k(\mathbf{x}_*, \mathbf{x}_i) \end{aligned} \quad (2)$$

which is a $D \times 1$ vector. To compute the derivative of the posterior mean we need to concatenate this derivative for each of the training points. It is helpful to define, $\tilde{X}_* = [\mathbf{x}_* - \mathbf{x}_1, \dots, \mathbf{x}_* - \mathbf{x}_N]^T$, which is an $N \times D$ matrix.

$$\begin{aligned} \frac{\partial \bar{f}_*}{\partial \mathbf{x}_*} &= \frac{\partial \mathbf{k}(\mathbf{x}_*, X)}{\partial \mathbf{x}_*} \boldsymbol{\alpha} \\ &= -\Lambda^{-1} \tilde{X}_*^T (\mathbf{k}(\mathbf{x}_*, X)^T \odot \boldsymbol{\alpha}) \end{aligned} \quad (3)$$

which is a $D \times 1$ vector. \odot represents an element-wise product.

2 Distribution over First Order Derivatives of Posterior Functions

In the previous section we found the derivative of the posterior mean of a GP. However, it is possible to find the distribution over derivatives of functions drawn from the GP posterior. Consider the random GP function values at two test point locations,

$$\begin{aligned} f(\mathbf{x}_*) &= \bar{f}(\mathbf{x}_*) + z_* \\ f(\mathbf{x}_* + \boldsymbol{\delta}) &= \bar{f}(\mathbf{x}_* + \boldsymbol{\delta}) + z_\delta \end{aligned} \quad (4)$$

where,

$$P(z_*, z_\delta) = \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} k_{**} - \mathbf{k}_*^T K^{-1} \mathbf{k}_* & k_{*\delta} - \mathbf{k}_*^T K^{-1} \mathbf{k}_\delta \\ k_{\delta*} - \mathbf{k}_\delta^T K^{-1} \mathbf{k}_* & k_{\delta\delta} - \mathbf{k}_\delta^T K^{-1} \mathbf{k}_\delta \end{bmatrix} \right) \quad (5)$$

The derivative is,

$$\begin{aligned}
\frac{\partial f_*}{\partial \mathbf{x}_*} &= \lim_{\delta \rightarrow 0} \frac{f(\mathbf{x}_* + \delta) - f(\mathbf{x}_*)}{\mathbf{x}_* + \delta - \mathbf{x}_*} \\
&= \lim_{\delta \rightarrow 0} \frac{\bar{f}(\mathbf{x}_* + \delta) + z_\delta - \bar{f}(\mathbf{x}_*) - z_*}{\delta} \\
&= \lim_{\delta \rightarrow 0} \frac{\bar{f}(\mathbf{x}_* + \delta) - \bar{f}(\mathbf{x}_*)}{\delta} + \lim_{\delta \rightarrow 0} \frac{z_\delta - z_*}{\delta} \\
&= \frac{\partial \bar{f}_*}{\partial \mathbf{x}_*} + \lim_{\delta \rightarrow 0} \frac{z_\delta - z_*}{\delta}
\end{aligned} \tag{6}$$

This is a random variable, the mean of which is given by the first term, and the variance comes from the second. The variance of the second term is found as follows,

$$\begin{aligned}
\mathbb{V} \left[\lim_{\delta \rightarrow 0} \frac{z_\delta - z_*}{\delta} \right] &= \lim_{\delta \rightarrow 0} \frac{1}{\delta^2} (\mathbb{V}[z_\delta] + \mathbb{V}[z_*] - \mathbb{C}[z_\delta, z_*] - \mathbb{C}[z_*, z_\delta]) \\
&= \lim_{\delta \rightarrow 0} \frac{1}{\delta^2} \left(k_{\delta\delta} - \mathbf{k}_\delta^T K^{-1} \mathbf{k}_\delta + k_{**} - \mathbf{k}_*^T K^{-1} \mathbf{k}_* - (k_{*\delta} - \mathbf{k}_*^T K^{-1} \mathbf{k}_\delta) - (k_{\delta*} - \mathbf{k}_\delta^T K^{-1} \mathbf{k}_*) \right) \\
&= \lim_{\delta \rightarrow 0} \frac{1}{\delta^2} (k_{\delta\delta} - k_{*\delta} - k_{\delta*} + k_{**} - (\mathbf{k}_\delta - \mathbf{k}_*)^T K^{-1} (\mathbf{k}_\delta - \mathbf{k}_*)) \\
&= \frac{\partial^2 k(\mathbf{x}_1^*, \mathbf{x}_2^*)}{\partial \mathbf{x}_1^* \partial \mathbf{x}_2^*} - \frac{\partial \mathbf{k}(\mathbf{x}_1^*, X)}{\partial \mathbf{x}_1^*} K^{-1} \frac{\partial \mathbf{k}(X, \mathbf{x}_2^*)}{\partial \mathbf{x}_2^*}
\end{aligned} \tag{7}$$

which is a $D \times D$ matrix - the variances and covariances of the derivatives w.r.t. each dimension in \mathbf{x}_* . Thus,

$$P \left(\frac{\partial f_*}{\partial \mathbf{x}_*} \right) = \mathcal{N} \left(\frac{\partial \bar{f}_*}{\partial \mathbf{x}_*}, \frac{\partial^2 k(\mathbf{x}_1^*, \mathbf{x}_2^*)}{\partial \mathbf{x}_1^* \partial \mathbf{x}_2^*} - \frac{\partial \mathbf{k}(\mathbf{x}_*, X)}{\partial \mathbf{x}_*} K^{-1} \frac{\partial \mathbf{k}(X, \mathbf{x}_*)}{\partial \mathbf{x}_*} \right) \tag{8}$$

We see that the mean of the distribution of derivatives is the derivative of the posterior mean. This is to be expected as both differentiation and expectation are linear operations and so are commutative.

3 Expected Squared Derivative

When propagating variances through first order Taylor series models, one uses the square of the first order derivative. We could also take the square inside the expectation which might lead to a better model.

The square of the expected derivative is given by,

$$\mathbb{E} \left[\frac{\partial f_*}{\partial \mathbf{x}_*} \right]^2 = \frac{\partial \bar{f}_*}{\partial \mathbf{x}_*} \frac{\partial \bar{f}_*}{\partial \mathbf{x}_*}^T \tag{9}$$

We can find the expected squared derivative as follows,

$$\begin{aligned}
\mathbb{E} \left[\left(\frac{\partial f_*}{\partial \mathbf{x}_*} \right)^2 \right] &= \mathbb{V} \left[\frac{\partial f_*}{\partial \mathbf{x}_*} \right] + \mathbb{E} \left[\frac{\partial f_*}{\partial \mathbf{x}_*} \right]^2 \\
&= \frac{\partial^2 k(\mathbf{x}_1^*, \mathbf{x}_2^*)}{\partial \mathbf{x}_1^* \partial \mathbf{x}_2^*} - \frac{\partial \mathbf{k}(\mathbf{x}_*, X)}{\partial \mathbf{x}_*} K^{-1} \frac{\partial \mathbf{k}(X, \mathbf{x}_*)}{\partial \mathbf{x}_*} + \frac{\partial \bar{f}_*}{\partial \mathbf{x}_*} \frac{\partial \bar{f}_*}{\partial \mathbf{x}_*}^T
\end{aligned} \tag{10}$$

Compared to equation 9 the expected squared derivative is inflated by the variance of the derivative

4 Derivatives with Uncertain Inputs

We can also ask what the distribution over the derivatives is when the input location is Gaussian distributed, i.e.,

$$\mathbf{x}_* \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma) \tag{11}$$

4.1 The mean

We can use the rule of iterated expectations to find the mean derivative,

$$\begin{aligned}
\mathbb{E} \left[\frac{\partial f^*}{\partial \mathbf{x}^*} \right] &= \mathbb{E}_{\mathbf{x}^*} \left[\mathbb{E}_{f^*} \left[\frac{\partial f^*}{\partial \mathbf{x}^*} \right] \right] \\
&= \mathbb{E}_{\mathbf{x}^*} \left[\frac{\partial \bar{f}^*}{\partial \mathbf{x}^*} \right] \\
&= \mathbb{E}_{\mathbf{x}^*} \left[-\Lambda^{-1} \tilde{X}_* (\mathbf{k}(\mathbf{x}_*, X)^T \odot \boldsymbol{\alpha}) \right] \\
&= -\Lambda^{-1} \mathbb{E}_{\mathbf{x}^*} \left[\sum_{i=1}^N \alpha_i (\mathbf{x}_* - \mathbf{x}_i) k(\mathbf{x}_*, \mathbf{x}_i) \right] \\
&= -\Lambda^{-1} \sum_{i=1}^N \alpha_i \mathbb{E}_{\mathbf{x}^*} [(\mathbf{x}_* - \mathbf{x}_i) k(\mathbf{x}_*, \mathbf{x}_i)] \\
&= -\Lambda^{-1} \sum_{i=1}^N \alpha_i \mathbb{E}_{\mathbf{x}^*} [\mathbf{x}_* k(\mathbf{x}_*, \mathbf{x}_i)] - \alpha_i \mathbf{x}_i \mathbb{E}_{\mathbf{x}^*} [k(\mathbf{x}_*, \mathbf{x}_i)]
\end{aligned} \tag{12}$$

To find the two expectations it is useful to note the squared exponential kernel is closely related to the Gaussian p.d.f.,

$$k(\mathbf{x}_*, \mathbf{x}_i) = \sigma_f^2 \exp \left(-\frac{1}{2} (\mathbf{x}_* - \mathbf{x}_i)^T \Lambda^{-1} (\mathbf{x}_* - \mathbf{x}_i) \right) \tag{13}$$

$$= \sigma_f^2 (2\pi)^{D/2} |\Lambda|^{1/2} \mathcal{N}(\mathbf{x}_*; \mathbf{x}_i, \Lambda) \tag{14}$$

and also to quote the area under a product of two Gaussians,

$$\int \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma_1) \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \Sigma_2) d\mathbf{x} = (2\pi)^{-D/2} |\Sigma_1 + \Sigma_2|^{-1/2} \exp \left(-\frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T (\Sigma_1 + \Sigma_2)^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right) \tag{15}$$

$$\triangleq Z \tag{16}$$

We start with the simpler of the two expectations,

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}^*} [k(\mathbf{x}_*, \mathbf{x}_i)] &= \int k(\mathbf{x}_*, \mathbf{x}_i) p(\mathbf{x}_*) d\mathbf{x}_* \\
&= \sigma_f^2 (2\pi)^{D/2} |\Lambda|^{1/2} \int \mathcal{N}(\mathbf{x}_*; \mathbf{x}_i, \Lambda) \mathcal{N}(\mathbf{x}_*; \boldsymbol{\mu}, \Sigma) d\mathbf{x}_* \\
&= \sigma_f^2 (2\pi)^{D/2} |\Lambda|^{1/2} (2\pi)^{-D/2} |\Lambda + \Sigma|^{-1/2} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^T (\Lambda + \Sigma)^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) \\
&= \sigma_f^2 |\Lambda|^{1/2} |\Lambda + \Sigma|^{-1/2} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^T (\Lambda + \Sigma)^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) \\
&= |\Sigma \Lambda^{-1} + I|^{-1/2} k(\mathbf{x}_i, \boldsymbol{\mu}, \Lambda + \Sigma)
\end{aligned} \tag{17}$$

where the third argument to the covariance function specifies the lengthscales.

The second expectation,

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}^*} [\mathbf{x}_* k(\mathbf{x}_*, \mathbf{x}_i)] &= \int \mathbf{x}_* k(\mathbf{x}_*, \mathbf{x}_i) p(\mathbf{x}_*) d\mathbf{x}_* \\
&= \sigma_f^2 (2\pi)^{D/2} |\Lambda|^{1/2} \int \mathbf{x}_* \mathcal{N}(\mathbf{x}_*; \mathbf{x}_i, \Lambda) \mathcal{N}(\mathbf{x}_*; \boldsymbol{\mu}, \Sigma) d\mathbf{x}_*
\end{aligned} \tag{18}$$

which is the mean of the product of two Gaussian distributions (times a constant). Therefore,

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}_*} [\mathbf{x}_* k(\mathbf{x}_*, \mathbf{x}_i)] &= \sigma_f^2 (2\pi)^{D/2} |\Lambda|^{1/2} Z (\Lambda^{-1} + \Sigma^{-1})^{-1} (\Lambda^{-1} \mathbf{x}_i + \Sigma^{-1} \boldsymbol{\mu}) \\
&= \sigma_f^2 |\Lambda|^{1/2} |\Sigma + \Lambda|^{-1/2} \Lambda \exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu})^T (\Sigma + \Lambda)^{-1} (\mathbf{x}_i - \boldsymbol{\mu})\right) (\Sigma + \Lambda)^{-1} \Sigma (\Lambda^{-1} \mathbf{x}_i + \Sigma^{-1} \boldsymbol{\mu}) \\
&= |\Sigma \Lambda^{-1} + I|^{-1/2} k(\mathbf{x}_i, \boldsymbol{\mu}, \Lambda + \Sigma) \Lambda (\Sigma + \Lambda)^{-1} (\Sigma \Lambda^{-1} \mathbf{x}_i + \boldsymbol{\mu}) \\
&= \mathbb{E}_{\mathbf{x}_*} [k(\mathbf{x}_*, \mathbf{x}_i)] \Lambda (\Sigma + \Lambda)^{-1} (\Sigma \Lambda^{-1} \mathbf{x}_i + \boldsymbol{\mu})
\end{aligned} \tag{19}$$

Putting equations 17 and 19 into equation 12 gives,

$$\begin{aligned}
\mathbb{E} \left[\frac{\partial f_*}{\partial \mathbf{x}_*} \right] &= -\Lambda^{-1} \sum_{i=1}^N \alpha_i (\mathbb{E}_{\mathbf{x}_*} [\mathbf{x}_* k(\mathbf{x}_*, \mathbf{x}_i)] - \mathbf{x}_i \mathbb{E}_{\mathbf{x}_*} [k(\mathbf{x}_*, \mathbf{x}_i)]) \\
&= -\Lambda^{-1} \sum_{i=1}^N \left(\Lambda (\Sigma + \Lambda)^{-1} (\Sigma \Lambda^{-1} \mathbf{x}_i + \boldsymbol{\mu}) - \mathbf{x}_i \right) \alpha_i \mathbb{E}_{\mathbf{x}_*} [k(\mathbf{x}_*, \mathbf{x}_i)] \\
&= \sum_{i=1}^N \left((\Lambda^{-1} - (\Sigma + \Lambda)^{-1} \Sigma \Lambda^{-1}) \mathbf{x}_i - (\Sigma + \Lambda)^{-1} \boldsymbol{\mu} \right) \alpha_i \mathbb{E}_{\mathbf{x}_*} [k(\mathbf{x}_*, \mathbf{x}_i)] \\
&= (\Sigma + \Lambda)^{-1} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}) \alpha_i \mathbb{E}_{\mathbf{x}_*} [k(\mathbf{x}_*, \mathbf{x}_i)] \\
&= |\Sigma \Lambda^{-1} + I|^{-1/2} (\Sigma + \Lambda)^{-1} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}) \alpha_i k(\mathbf{x}_i, \boldsymbol{\mu}, \Lambda + \Sigma) \\
&= |\Sigma \Lambda^{-1} + I|^{-1/2} (\Sigma + \Lambda)^{-1} \tilde{X}_*^T (\boldsymbol{\alpha} \odot \mathbf{k}(X, \boldsymbol{\mu}, \Lambda + \Sigma))
\end{aligned} \tag{20}$$

$$\tag{21}$$

4.2 The variance

We can use the rule of total variance to find the variance of the derivative,

$$\begin{aligned}
\mathbb{V} \left[\frac{\partial f_*}{\partial \mathbf{x}_*} \right] &= \mathbb{E}_{\mathbf{x}_*} \left[\mathbb{V}_{f_*} \left[\frac{\partial f_*}{\partial \mathbf{x}_*} \right] \right] + \mathbb{V}_{\mathbf{x}_*} \left[\mathbb{E}_{f_*} \left[\frac{\partial f_*}{\partial \mathbf{x}_*} \right] \right] \\
&= \mathbb{E}_{\mathbf{x}_*} \left[\frac{\partial^2 k(\mathbf{x}_1^*, \mathbf{x}_2^*)}{\partial \mathbf{x}_1^* \partial \mathbf{x}_2^*} - \frac{\partial \mathbf{k}(\mathbf{x}_*, X)}{\partial \mathbf{x}_*} K^{-1} \frac{\partial \mathbf{k}(X, \mathbf{x}_*)}{\partial \mathbf{x}_*} \right] + \mathbb{V}_{\mathbf{x}_*} \left[\frac{\partial f_*}{\partial \mathbf{x}_*} \right] \\
&= \mathbb{E}_{\mathbf{x}_*} \left[\frac{\partial^2 k(\mathbf{x}_1^*, \mathbf{x}_2^*)}{\partial \mathbf{x}_1^* \partial \mathbf{x}_2^*} - \frac{\partial \mathbf{k}(\mathbf{x}_*, X)}{\partial \mathbf{x}_*} K^{-1} \frac{\partial \mathbf{k}(X, \mathbf{x}_*)}{\partial \mathbf{x}_*} \right] + \mathbb{V}_{\mathbf{x}_*} \left[-\Lambda^{-1} \tilde{X}_*^T (\mathbf{k}(\mathbf{x}_*, X)^T \odot \boldsymbol{\alpha}) \right]
\end{aligned} \tag{22}$$

We will calculate these expectations separately. Firstly the expectation of the second derivative of the kernel (see section 5 for derivation of the second derivative),

$$\begin{aligned}
\mathbb{E} \left[\frac{\partial^2 k(\mathbf{x}_1^*, \mathbf{x}_2^*)}{\partial \mathbf{x}_1^* \partial \mathbf{x}_2^*} \right] &= \mathbb{E} [\Lambda^{-1} k(\mathbf{x}_*, \mathbf{x}_*)] \\
&= \Lambda^{-1} \sigma_f^2
\end{aligned} \tag{23}$$

$$\begin{aligned}
\mathbb{E} \left[\frac{\partial \mathbf{k}(\mathbf{x}_*, X)}{\partial \mathbf{x}_i^*} K^{-1} \frac{\partial \mathbf{k}(X, \mathbf{x}_*)}{\partial \mathbf{x}_j^*} \right] &= \Lambda_i^{-1} \mathbb{E} \left[[\tilde{X}_*^{(i)T} \odot \mathbf{k}(\mathbf{x}_*, X)] K^{-1} [\mathbf{k}(X, \mathbf{x}_*) \odot \tilde{X}_*^{(j)}] \right] \Lambda_j^{-1} \\
&= \Lambda_i^{-1} \left(\mathbb{E}[\tilde{X}_*^{(i)T} \odot \mathbf{k}(\mathbf{x}_*, X)] K^{-1} \mathbb{E}[\mathbf{k}(X, \mathbf{x}_*) \odot \tilde{X}_*^{(j)}] \right. \\
&\quad \left. + \text{tr} \left(K^{-1} \mathbb{C}[\tilde{X}_*^{(i)T} \odot \mathbf{k}(\mathbf{x}_*, X), \tilde{X}_*^{(j)T} \odot \mathbf{k}(\mathbf{x}_*, X)] \right) \right) \Lambda_j^{-1}
\end{aligned} \tag{24}$$

where,

$$\tilde{X}_*^{(i)} = [x_1^{(i)} - x_*^{(i)}, \dots, x_N^{(i)} - x_*^{(i)}]^T \quad (25)$$

which is a $N \times 1$ vector.

$$\begin{aligned} \mathbb{E}_{x_*} [(x_1^{(i)} - x_*^{(i)}) k(\mathbf{x}_1, \mathbf{x}_*)] &= x_1^{(i)} \mathbb{E}[k(\mathbf{x}_1, \mathbf{x}_*)] - \mathbb{E}[x_*^{(i)} k(\mathbf{x}_1, \mathbf{x}_*)] \\ &= x_1^{(i)} \bar{k}(\boldsymbol{\mu}, \Sigma, \mathbf{x}_1, \Lambda) - \bar{k}(\boldsymbol{\mu}, \Sigma, \mathbf{x}_1, \Lambda) (\Sigma_{ii} \Lambda_i^{-1} + 1)^{-1} (\Sigma_{ii} \Lambda_i^{-1} x_1^{(i)} + \mu_i) \\ &= \left(x_1^{(i)} - (\Sigma_{ii} \Lambda_i^{-1} + 1)^{-1} (\Sigma_{ii} \Lambda_i^{-1} x_1^{(i)} + \mu_i) \right) \bar{k}(\boldsymbol{\mu}, \Sigma, \mathbf{x}_1, \Lambda) \\ &= \bar{k}(\boldsymbol{\mu}, \Sigma, \mathbf{x}_1, \Lambda) (x_1^{(i)} - \mu_i) (\Sigma_{ii} \Lambda_i^{-1} + 1)^{-1} \end{aligned} \quad (26)$$

Defining \mathbf{U} to be a $N \times 1$ vector with elements, $U_k = x_k^{(i)} - \mu_i$,

$$\mathbb{E}_{x_*} [\tilde{X}_*^{(i)} \odot k(X, \mathbf{x}_*)] = \bar{k}(X, \boldsymbol{\mu}, \Sigma, \Lambda) \odot \mathbf{U} (\Sigma_{ii} \Lambda_i^{-1} + 1)^{-1} \quad (27)$$

which is a $N \times 1$ vector.

$$\begin{aligned} \mathbb{C}_{x_*} [(x_1^{(i)} - x_*^{(i)}) k(\mathbf{x}_1, \mathbf{x}_*), (x_2^{(j)} - x_*^{(j)}) k(\mathbf{x}_2, \mathbf{x}_*)] \\ &= x_1^{(i)} x_2^{(j)} \mathbb{C}_{x_*} [k(\mathbf{x}_1, \mathbf{x}_*), k(\mathbf{x}_2, \mathbf{x}_*)] - x_1^{(i)} \mathbb{C}_{x_*} [k(\mathbf{x}_1, \mathbf{x}_*), x_*^{(j)} k(\mathbf{x}_2, \mathbf{x}_*)] \\ &\quad - x_2^{(j)} \mathbb{C}_{x_*} [x_*^{(i)} k(\mathbf{x}_1, \mathbf{x}_*), k(\mathbf{x}_2, \mathbf{x}_*)] + \mathbb{C}_{x_*} [x_*^{(i)} k(\mathbf{x}_1, \mathbf{x}_*), x_*^{(j)} k(\mathbf{x}_2, \mathbf{x}_*)] \\ &= x_1^{(i)} x_2^{(j)} C_{kk}(\boldsymbol{\mu}, \mathbf{x}_1, \mathbf{x}_2, \Sigma) - x_1^{(i)} C_{xkk}^{(j)}(\boldsymbol{\mu}, \mathbf{x}_2, \mathbf{x}_1, \Sigma) - x_2^{(j)} C_{xkk}^{(i)}(\boldsymbol{\mu}, \mathbf{x}_1, \mathbf{x}_2, \Sigma) + C_{xkxk}^{(ij)}(\boldsymbol{\mu}, \mathbf{x}_1, \mathbf{x}_2, \Sigma) \end{aligned} \quad (28)$$

The C terms are derived and defined in section 6. Therefore,

$$\begin{aligned} \mathbb{C}[\tilde{X}_*^{(i)T} \odot \mathbf{k}(\mathbf{x}_*, X), \tilde{X}_*^{(j)T} \odot \mathbf{k}(\mathbf{x}_*, X)] \\ &= (X^{(i)} X^{(j)T}) \odot C_{kk}(\boldsymbol{\mu}, X, X, \Sigma) - X^{(i)} \odot C_{xkk}^{(j)}(\boldsymbol{\mu}, X, X, \Sigma) - X^{(j)T} \odot C_{xkk}^{(i)}(\boldsymbol{\mu}, X, X, \Sigma) + C_{xkxk}^{(ij)}(\boldsymbol{\mu}, X, X, \Sigma) \end{aligned} \quad (29)$$

which is a $N \times N$ matrix.

Finally we need, $\mathbb{V}_{x_*} [\Lambda^{-1} \tilde{X}_*^T (\mathbf{k}(X, \mathbf{x}_*) \odot \boldsymbol{\alpha})]$, which we will break up and compute as,

$$\begin{aligned} \mathbb{C}_{x_*} [\Lambda_i^{-1} \tilde{X}_*^{(i)T} (\mathbf{k}(X, \mathbf{x}_*) \odot \boldsymbol{\alpha}), \Lambda_j^{-1} \tilde{X}_*^{(j)T} (\mathbf{k}(X, \mathbf{x}_*) \odot \boldsymbol{\alpha})] \\ &= \Lambda_i^{-1} \Lambda_j^{-1} \mathbb{C}_{x_*} \left[\sum_k (x_k^{(i)} - x_*^{(i)}) k(\mathbf{x}_k, \mathbf{x}_*) \alpha_k, \sum_l (x_l^{(j)} - x_*^{(j)}) k(\mathbf{x}_l, \mathbf{x}_*) \alpha_l \right] \\ &= \Lambda_i^{-1} \Lambda_j^{-1} \sum_k \sum_l \alpha_k \alpha_l \left(x_k^{(i)} x_l^{(j)} C_{kk} - x_k^{(i)} C_{xkk}^{(j)} - x_l^{(j)} C_{xkk}^{(i)} + C_{xkxk}^{(ij)} \right) \\ &= \Lambda_i^{-1} \Lambda_j^{-1} \boldsymbol{\alpha}_k^T \left(X^{(i)} X^{(j)T} \odot C_{kk} - X^{(i)} C_{xkk}^{(j)} - X^{(j)T} C_{xkk}^{(i)} + C_{xkxk}^{(ij)} \right) \boldsymbol{\alpha}_l \end{aligned} \quad (30)$$

5 Differentiating the Squared Exponential kernel

The squared exponential kernel is given by,

$$k(\mathbf{x}_1, \mathbf{x}_2) = \sigma_f^2 \exp \left(-\frac{1}{2} (\mathbf{x}_1 - \mathbf{x}_2)^T \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) \right) \quad (31)$$

where the hyperparameters are the signal variance σ_f^2 and a characteristic length-scale for each dimension, $\{l_i\}_{i=1}^D$. The squared length-scales are collected into a $D \times D$, diagonal matrix Λ ,

$$\Lambda = \begin{bmatrix} l_1^2 & 0 & 0 \\ \vdots & \ddots & \vdots \\ 0 & 0 & l_D^2 \end{bmatrix} \quad (32)$$

The derivative of this kernel with respect to the first argument is,

$$\begin{aligned}
\frac{\partial k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_1} &= \frac{\partial}{\partial \mathbf{x}_1} \left\{ \sigma_f^2 \exp \left(-\frac{1}{2} (\mathbf{x}_1 - \mathbf{x}_2)^T \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) \right) \right\} \\
&= \frac{\partial}{\partial \mathbf{x}_1} \left\{ -\frac{1}{2} (\mathbf{x}_1 - \mathbf{x}_2)^T \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) \right\} k(\mathbf{x}_1, \mathbf{x}_2) \\
&= -\frac{1}{2} \frac{\partial}{\partial \mathbf{x}_1} \{ \mathbf{x}_1^T \Lambda^{-1} \mathbf{x}_1 - \mathbf{x}_2^T \Lambda^{-1} \mathbf{x}_1 - \mathbf{x}_1^T \Lambda^{-1} \mathbf{x}_2 \} k(\mathbf{x}_1, \mathbf{x}_2) \\
&= -\frac{1}{2} (2\Lambda^{-1} \mathbf{x}_1 - 2\Lambda^{-1} \mathbf{x}_2) k(\mathbf{x}_1, \mathbf{x}_2) \\
&= -\Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) k(\mathbf{x}_1, \mathbf{x}_2)
\end{aligned} \tag{33}$$

which is a $D \times 1$ vector. We can also take the derivative w.r.t. the second argument,

$$\begin{aligned}
\frac{\partial k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_2} &= \frac{\partial}{\partial \mathbf{x}_2} \left\{ \sigma_f^2 \exp \left(-\frac{1}{2} (\mathbf{x}_1 - \mathbf{x}_2)^T \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) \right) \right\} \\
&= \frac{\partial}{\partial \mathbf{x}_2} \left\{ -\frac{1}{2} (\mathbf{x}_1 - \mathbf{x}_2)^T \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) \right\} k(\mathbf{x}_1, \mathbf{x}_2) \\
&= -\frac{1}{2} \frac{\partial}{\partial \mathbf{x}_2} \{ -\mathbf{x}_2^T \Lambda^{-1} \mathbf{x}_1 - \mathbf{x}_1^T \Lambda^{-1} \mathbf{x}_2 + \mathbf{x}_2^T \Lambda^{-1} \mathbf{x}_2 \} k(\mathbf{x}_1, \mathbf{x}_2) \\
&= -\frac{1}{2} (-2\Lambda^{-1} \mathbf{x}_1 + 2\Lambda^{-1} \mathbf{x}_2) k(\mathbf{x}_1, \mathbf{x}_2) \\
&= \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) k(\mathbf{x}_1, \mathbf{x}_2)
\end{aligned} \tag{34}$$

Note that the only difference between these two derivatives (equations 33 and 34) is the minus sign in equation 33. This comes about because the distance between \mathbf{x}_1 and \mathbf{x}_2 is calculated as $\mathbf{x}_1 - \mathbf{x}_2$ and hence increasing x_1 increases the separation and so decreases the covariance; the opposite is true for \mathbf{x}_2 .

It is trivial to extend these results for the case when one of the inputs is a collection of points, such as a $N \times D$ training matrix X ,

$$X = [\mathbf{x}_1, \dots, \mathbf{x}_2]^T \tag{35}$$

$$\frac{\partial \mathbf{k}(X, \mathbf{x}_*)}{\partial \mathbf{x}_*} = \mathbf{k}(X, \mathbf{x}_*) \odot \tilde{X}_* \Lambda^{-1} \tag{36}$$

where we define \tilde{X}_* to be the $N \times D$ matrix, $[\mathbf{x}_1 - \mathbf{x}_*, \dots, \mathbf{x}_N - \mathbf{x}_*]^T$. Note that $\mathbf{k}(X, \mathbf{x}_*)$ is a $N \times 1$ column vector.

Building on equations 33 and 34 we can now find the second derivative (cross term),

$$\begin{aligned}
\frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_1 \partial \mathbf{x}_2} &= \frac{\partial}{\partial \mathbf{x}_1} \{ \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) k(\mathbf{x}_1, \mathbf{x}_2) \} \\
&= \Lambda^{-1} (I - (\mathbf{x}_1 - \mathbf{x}_2) (\mathbf{x}_1 - \mathbf{x}_2)^T \Lambda^{-1}) k(\mathbf{x}_1, \mathbf{x}_2)
\end{aligned} \tag{37}$$

We can see the following relationship,

$$\frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_1^2} = \frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_2^2} = -\frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_1 \partial \mathbf{x}_2} = -\frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_2 \partial \mathbf{x}_1} \tag{38}$$

We can summarise the derivatives as follows,

$$\frac{\partial k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_2} = \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) k(\mathbf{x}_1, \mathbf{x}_2) \tag{39}$$

$$\frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_2^2} = -\Lambda^{-1} k(\mathbf{x}_1, \mathbf{x}_2) + \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2) \frac{\partial k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_2}^T \tag{40}$$

To find higher derivatives we need to switch to writing down elements of the derivative. First rephrase the first two derivatives,

$$\frac{\partial k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_2^{(i)}} = \Lambda_{ii}^{-1} (x_1^{(i)} - x_2^{(i)}) k(\mathbf{x}_1, \mathbf{x}_2) \quad (41)$$

$$\frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(j)} \partial x_2^{(i)}} = -\Lambda_{ij}^{-1} k(\mathbf{x}_1, \mathbf{x}_2) + \Lambda_{ii}^{-1} (x_1^{(i)} - x_2^{(i)}) \frac{\partial k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(j)}} \quad (42)$$

$$\frac{\partial^3 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(k)} \partial x_2^{(j)} \partial x_2^{(i)}} = -\Lambda_{ij}^{-1} \frac{\partial k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(k)}} - \Lambda_{ik}^{-1} \frac{\partial k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(j)}} + \Lambda_{ii}^{-1} (x_1^{(i)} - x_2^{(i)}) \frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(k)} \partial x_2^{(j)}} \quad (43)$$

$$\frac{\partial^4 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(l)} \partial x_2^{(k)} \partial x_2^{(j)} \partial x_2^{(i)}} = -\Lambda_{ij}^{-1} \frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(l)} \partial x_2^{(k)}} - \Lambda_{ik}^{-1} \frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(l)} \partial x_2^{(j)}} - \Lambda_{il}^{-1} \frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(k)} \partial x_2^{(j)}} + \Lambda_{ii}^{-1} (x_1^{(i)} - x_2^{(i)}) \frac{\partial^3 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(l)} \partial x_2^{(k)} \partial x_2^{(j)}} \quad (44)$$

We sometimes need to evaluate these derivatives for $\mathbf{x}_1 = \mathbf{x}_2$,

$$k(\mathbf{x}_1, \mathbf{x}_2)|_{\mathbf{x}_1 = \mathbf{x}_2} = \sigma_f^2 \quad (45)$$

$$\left. \frac{\partial k(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_2^{(i)}} \right|_{\mathbf{x}_1 = \mathbf{x}_2} = 0 \quad (46)$$

$$\left. \frac{\partial^2 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(j)} \partial x_2^{(i)}} \right|_{\mathbf{x}_1 = \mathbf{x}_2} = -\Lambda_{ij}^{-1} \sigma_f^2 \quad (47)$$

$$\left. \frac{\partial^3 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(k)} \partial x_2^{(j)} \partial x_2^{(i)}} \right|_{\mathbf{x}_1 = \mathbf{x}_2} = 0 \quad (48)$$

$$\left. \frac{\partial^4 k(\mathbf{x}_1, \mathbf{x}_2)}{\partial x_2^{(l)} \partial x_2^{(k)} \partial x_2^{(j)} \partial x_2^{(i)}} \right|_{\mathbf{x}_1 = \mathbf{x}_2} = \Lambda_{ij}^{-1} \Lambda_{kl}^{-1} \sigma_f^2 + \Lambda_{ik}^{-1} \Lambda_{jl}^{-1} \sigma_f^2 + \Lambda_{il}^{-1} \Lambda_{jk}^{-1} \sigma_f^2 \quad (49)$$

6 Squared Exponential Kernel Moments

The squared exponential kernel,

$$k(\mathbf{x}_1, \mathbf{x}_2) = \sigma_f^2 \exp\left(-\frac{1}{2}(\mathbf{x}_1 - \mathbf{x}_2)^T \Lambda^{-1} (\mathbf{x}_1 - \mathbf{x}_2)\right) \quad (50)$$

$$= \sigma_f^2 (2\pi)^{D/2} |\Lambda|^{1/2} \mathcal{N}(\mathbf{x}_1; \mathbf{x}_2, \Lambda) \quad (51)$$

Its mean,

$$\begin{aligned} \mathbb{E}_{\mathbf{x}}[k(\mathbf{x}, \mathbf{x}_1)] &= \sigma_f^2 |\Lambda|^{1/2} |\Lambda + \Sigma|^{-1/2} \exp\left(-\frac{1}{2}(\boldsymbol{\mu} - \mathbf{x}_1)^T (\Lambda + \Sigma)^{-1} (\boldsymbol{\mu} - \mathbf{x}_1)\right) \\ &= |\Sigma \Lambda^{-1} + I|^{-1/2} k(\boldsymbol{\mu}, \mathbf{x}_1, \Lambda + \Sigma) \end{aligned} \quad (52)$$

$$\triangleq \bar{k}(\boldsymbol{\mu}, \mathbf{x}_1, \Sigma, \Lambda) \quad (1 \times 1) \quad (53)$$

The product of two squared exponential kernels,

$$k(\mathbf{x}, \mathbf{x}_1) k(\mathbf{x}, \mathbf{x}_2) = k\left(\frac{\mathbf{x}_1}{2}, \frac{\mathbf{x}_2}{2}, \frac{\Lambda}{2}\right) k\left(\mathbf{x}, \frac{\mathbf{x}_1 + \mathbf{x}_2}{2}, \frac{\Lambda}{2}\right) \quad (54)$$

$$= \sigma_f^2 \pi^{D/2} |\Lambda|^{1/2} k\left(\frac{\mathbf{x}_1}{2}, \frac{\mathbf{x}_2}{2}, \frac{\Lambda}{2}\right) \mathcal{N}\left(\mathbf{x}; (\mathbf{x}_1 + \mathbf{x}_2)/2, \Lambda/2\right) \quad (55)$$

The mean of a product,

$$\mathbb{E}_{\mathbf{x}}[k(\mathbf{x}_1, \mathbf{x}) k(\mathbf{x}, \mathbf{x}_2)] = k(\mathbf{x}_1/2, \mathbf{x}_2/2, \Lambda/2) \bar{k}((\mathbf{x}_1 + \mathbf{x}_2)/2, \boldsymbol{\mu}_1, \Sigma, \Lambda/2) \quad (56)$$

$$= |2\Sigma\Lambda^{-1} + I|^{-1/2} k(\mathbf{x}_1/2, \mathbf{x}_2/2, \Lambda/2) k((\mathbf{x}_1 + \mathbf{x}_2)/2, \boldsymbol{\mu}_1, \Lambda/2 + \Sigma) \quad (57)$$

$$\triangleq E_{kk}(\boldsymbol{\mu}, \mathbf{x}_1, \mathbf{x}_2, \Sigma) \quad (1 \times 1) \quad (58)$$

Therefore, the covariance of two kernels,

$$\begin{aligned} \mathbb{C}_{\mathbf{x}}[k(\mathbf{x}, \mathbf{x}_1), k(\mathbf{x}, \mathbf{x}_2)] &= |2\Sigma\Lambda^{-1} + I|^{-1/2} k(\mathbf{x}_1/2, \mathbf{x}_2/2, \Lambda/2) k((\mathbf{x}_1 + \mathbf{x}_2)/2, \boldsymbol{\mu}, \Lambda/2 + \Sigma) \\ &\quad - |\Sigma\Lambda^{-1} + I|^{-1} k(\boldsymbol{\mu}, \mathbf{x}_1, \Lambda + \Sigma) k(\boldsymbol{\mu}, \mathbf{x}_2, \Lambda + \Sigma) \end{aligned} \quad (59)$$

$$\triangleq C_{kk}(\boldsymbol{\mu}, \mathbf{x}_1, \mathbf{x}_2, \Sigma) \quad (1 \times 1) \quad (60)$$

The covariance between \mathbf{x} times a covariance function with another covariance function,

$$\begin{aligned} \mathbb{C}_{\mathbf{x}}[\mathbf{x} k(\mathbf{x}, \mathbf{x}_1), k(\mathbf{x}, \mathbf{x}_2)] \\ &= k\left(\frac{\mathbf{x}_1}{2}, \frac{\mathbf{x}_2}{2}, \frac{\Lambda}{2}\right) \kappa(\boldsymbol{\mu}, \Sigma, (\mathbf{x}_1 + \mathbf{x}_2)/2, \Lambda/2) - \kappa(\boldsymbol{\mu}, \Sigma, \mathbf{x}_1, \Lambda) \bar{k}(\boldsymbol{\mu}, \mathbf{x}_2, \Sigma, \Lambda) \\ &\triangleq C_{xkk}(\boldsymbol{\mu}, \mathbf{x}_1, \mathbf{x}_2, \Sigma) \quad (D \times 1) \end{aligned} \quad (61)$$

The covariance of \mathbf{x} times a covariance function with \mathbf{x} times another covariance function,

$$\begin{aligned} \mathbb{C}_{\mathbf{x}}[\mathbf{x} k(\mathbf{x}, \mathbf{x}_1), \mathbf{x} k(\mathbf{x}, \mathbf{x}_2)] \\ &= k\left(\frac{\mathbf{x}_1}{2}, \frac{\mathbf{x}_2}{2}, \frac{\Lambda}{2}\right) \bar{k}(\boldsymbol{\mu}, \Sigma, (\mathbf{x}_1 + \mathbf{x}_2)/2, \Lambda/2) \\ &\quad [(2\Sigma\Lambda^{-1} + I)^{-1}(\Sigma\Lambda^{-1}(\mathbf{x}_1 + \mathbf{x}_2) + \boldsymbol{\mu})(\Sigma\Lambda^{-1}(\mathbf{x}_1 + \mathbf{x}_2) + \boldsymbol{\mu})^T (2\Sigma\Lambda^{-1} + I)^{-1} + \Sigma(\Sigma + \Lambda/2)^{-1}\Lambda/2] \\ &\quad - \kappa(\boldsymbol{\mu}, \mathbf{x}_1, \Sigma, \Lambda) \kappa(\boldsymbol{\mu}, \mathbf{x}_2, \Sigma, \Lambda)^T \\ &\triangleq C_{xkxk}(\boldsymbol{\mu}, \mathbf{x}_1, \mathbf{x}_2, \Sigma) \quad (D \times D) \end{aligned} \quad (62)$$