

# Gaussian Processes with Derivative Observations

Xiaoke Yang (das.xiaoke@hotmail.com)  
School of Automation Science and Electrical Engineering,  
Beihang University

**Abstract:** This article summarises necessary derivations for Gaussian processes (GP) with derivative observations. Please refer to ‘Solak et al, Derivative observations in Gaussian Process models of dynamic systems.’ for more details on derivative observations. Uncertainties over derivative observations are not considered in this article.

## 1 Preliminary

Assume the underlying system is modelled by a multiple-input single-output function with noisy measurements of the output, i.e.

$$y = f(\mathbf{x}) + v, \quad (1)$$

where  $\mathbf{x} \in \mathbb{R}^D$ ,  $y \in \mathbb{R}$ , and  $v$  is an additive Gaussian noise, i.e.  $v \sim \mathcal{N}(0, v)$ . The underlying function  $f(\mathbf{x})$  is modelled by a GP, i.e.

$$f \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}^m, \mathbf{x}^n)). \quad (2)$$

A zero mean function  $m(\mathbf{x}) = 0$  is used for the GP and the following covariance function is used

$$\text{cov}[y^m, y^n] = k(\mathbf{x}^m, \mathbf{x}^n) = \alpha \exp\left(-\frac{1}{2}\|\mathbf{x}^m - \mathbf{x}^n\|_{\mathbf{\Gamma}}^2\right) + v\delta_{m,n} \quad (3)$$

where  $\mathbf{x}^m, \mathbf{x}^n \in \mathbb{R}^D$  are two input points,  $\mathbf{\Gamma} = \text{diag}([\gamma_1 \ \gamma_2 \ \cdots \ \gamma_D])$ , notation  $\|\cdot\|_{\mathbf{\Gamma}}$  is defined as  $\|\mathbf{x}\|_{\mathbf{\Gamma}}^2 \triangleq \mathbf{x}^\top \mathbf{\Gamma} \mathbf{x}$ , and the Kronecker delta is defined as

$$\delta_{m,n} = \begin{cases} 1 & m = n \\ 0 & m \neq n \end{cases}. \quad (4)$$

We also define another covariance function for convenience,

$$k_f(\mathbf{x}^m, \mathbf{x}^n) = \alpha \exp\left(-\frac{1}{2}\|\mathbf{x}^m - \mathbf{x}^n\|_{\mathbf{\Gamma}}^2\right). \quad (5)$$

Basically,  $k_f(\mathbf{x}^m, \mathbf{x}^n)$  calculates the covariance between the function values, instead of the measurements of the function outputs. The vector of hyper-parameters is defined as

$$\boldsymbol{\theta} = [\gamma_1 \ \gamma_2 \ \cdots \ \gamma_D \ \alpha \ v]. \quad (6)$$

## 2 Derivative Observations

With derivative observations, data for the GP become

Table 1: Input-output data for GP		
Type	Derivative observation	Function observation
input	$\mathbf{x}_d \in \mathbb{R}^D$	$\mathbf{x} \in \mathbb{R}^D$
output	$\mathbf{y}_d \in \mathbb{R}^D$	$y \in \mathbb{R}$

If we define the input output data block as

$$\mathbf{X} = [\mathbf{x}_d^1 \quad \mathbf{x}_d^2 \quad \cdots \quad \mathbf{x}_d^M \quad \mathbf{x}^1 \quad \mathbf{x}^2 \quad \cdots \quad \mathbf{x}^N] \quad (7)$$

$$\mathbf{Y} = [(\mathbf{y}_d^1)^\top \quad (\mathbf{y}_d^2)^\top \quad \cdots \quad (\mathbf{y}_d^M)^\top \quad y^1 \quad y^2 \quad \cdots \quad y^N] \quad (8)$$

Then the output covariance matrix is

$$\mathbf{K} = \begin{bmatrix} \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n) & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{\substack{m \in \{1, \dots, M\} \\ n \in \{1, \dots, M\}}} & \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n) & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{\substack{m \in \{1, \dots, M\} \\ n \in \{1, \dots, N\}}} \\ \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & k_{xd}(\mathbf{x}^m, \mathbf{x}_d^n) & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{\substack{m \in \{1, \dots, N\} \\ n \in \{1, \dots, M\}}} & \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & k_{xx}(\mathbf{x}^m, \mathbf{x}^n) & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{\substack{m \in \{1, \dots, N\} \\ n \in \{1, \dots, N\}}} \end{bmatrix} \quad (9)$$

Things we need to compute include  $k_{xx}(\mathbf{x}^m, \mathbf{x}^n) \in \mathbb{R}$ ,  $k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n) \in \mathbb{R}^{D \times 1}$ , and  $k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n) \in \mathbb{R}^{D \times D}$ .  $k_{xx}(\mathbf{x}^m, \mathbf{x}^n) = k(\mathbf{x}^m, \mathbf{x}^n)$  as in (??). For the two derivatives The following results are from reference .

$$k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n)_i = -\alpha \gamma_i (x_{d,i}^m - x_i^n) \exp \left( -\frac{1}{2} \|\mathbf{x}_d^m - \mathbf{x}^n\|_{\Gamma}^2 \right), \quad (10)$$

$$= -\gamma_i (x_{d,i}^m - x_i^n) k_f(\mathbf{x}_d^m, \mathbf{x}^n), \quad (11)$$

$$k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n)_{i,j} = \alpha \gamma_i (\delta_{i,j} - \gamma_j (x_{d,i}^m - x_{d,i}^n)(x_{d,j}^m - x_{d,j}^n)) \exp \left( -\frac{1}{2} \|\mathbf{x}_d^m - \mathbf{x}_d^n\|_{\Gamma}^2 \right), \quad (12)$$

$$= \gamma_i (\delta_{i,j} - \gamma_j (x_{d,i}^m - x_{d,i}^n)(x_{d,j}^m - x_{d,j}^n)) k_f(\mathbf{x}_d^m, \mathbf{x}_d^n), \quad (13)$$

$$(14)$$

where the subscripts  $i$  and  $j$  corresponds to the indices within the array or matrix, respectively. and

$$k_f(\mathbf{x}_d^m, \mathbf{x}^n) = \exp \left( -\frac{1}{2} \|\mathbf{x}_d^m - \mathbf{x}^n\|_{\Gamma}^2 \right) = k_f(\mathbf{x}^n, \mathbf{x}_d^m) \quad (15)$$

Then, in vector and matrix form, the above equations can be written as

$$k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n) = -k_f(\mathbf{x}_d^m, \mathbf{x}^n) \Gamma (\mathbf{x}_d^m - \mathbf{x}^n), \quad (16)$$

$$k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n) = k_f(\mathbf{x}_d^m, \mathbf{x}_d^n) \left( \Gamma - (\Gamma (\mathbf{x}_d^m - \mathbf{x}_d^n)) (\Gamma (\mathbf{x}_d^m - \mathbf{x}_d^n))^\top \right), \quad (17)$$

$$= k_f(\mathbf{x}_d^m, \mathbf{x}_d^n) \left( \Gamma - \Gamma (\mathbf{x}_d^m - \mathbf{x}_d^n) (\mathbf{x}_d^m - \mathbf{x}_d^n)^\top \Gamma \right), \quad (18)$$

The GP predictions are

$$p(y^* | \mathbf{X}, \mathbf{Y}, \mathbf{x}^*) \sim \mathcal{N}(\mu, \Sigma) \quad (19)$$

$$(20)$$

where

$$\mu = k(\mathbf{x}^*, \mathbf{X}) \mathbf{K}^{-1} \mathbf{Y} \quad (21)$$

$$\Sigma = k(\mathbf{x}^*, \mathbf{x}^*) - k(\mathbf{x}^*, \mathbf{X}) \mathbf{K}^{-1} k(\mathbf{X}, \mathbf{x}^*) \quad (22)$$

where

$$k(\mathbf{x}^*, \mathbf{X}) = [[k_{xd}(\mathbf{x}^*, \mathbf{x}_d^m)]_{m \in \{1, \dots, M\}} \quad [k_{xx}(\mathbf{x}^*, \mathbf{x}^n)]_{n \in \{1, \dots, N\}}] \quad (23)$$

$$k(\mathbf{X}, \mathbf{x}^*) = k(\mathbf{x}^*, \mathbf{X})^\top \quad (24)$$

and the marginal likelihood is

$$L = \log p(\mathbf{Y}|\mathbf{X}) = -\frac{1}{2}\mathbf{Y}^\top \mathbf{K}^{-1}\mathbf{Y} - \frac{1}{2}\log |\mathbf{K}| - \frac{MD+N}{2}\log 2\pi \quad (25)$$

In order to train the GP, the derivative  $\frac{\partial L}{\partial \boldsymbol{\theta}}$  needs to be computed, and are listed as follows.

$$\frac{\partial L}{\partial \theta_l} = \text{vec} \left( \frac{\partial L}{\partial \mathbf{K}} \right)^\top \text{vec} \left( \frac{\partial \mathbf{K}}{\partial \theta_l} \right), \quad l \in \{1, 2, \dots, D+2\}, \quad (26)$$

where  $\circ$  denotes the matrix Hadamard product or element-wise product. The first half of this product is a  $(MD+N) \times (MD+N)$  matrix, whose  $i, j^{\text{th}}$  element is defined as  $\frac{\partial L}{\partial K_{i,j}}$

$$\frac{\partial L}{\partial \mathbf{K}} = \begin{bmatrix} \cdots & \cdots & \cdots \\ \cdots & \frac{\partial L}{\partial K_{i,j}} & \cdots \\ \cdots & \cdots & \cdots \end{bmatrix}_{\substack{i \in \{1, \dots, MD+N\} \\ j \in \{1, \dots, MD+N\}}} \quad (27)$$

from matrix cookbook

$$\frac{\partial L}{\partial \mathbf{K}} = \frac{1}{2}\mathbf{K}^{-\top} \mathbf{Y} \mathbf{Y}^\top \mathbf{K}^{-\top} - \frac{1}{2}\mathbf{K}^{-\top} \quad (28)$$

Then we need to compute the right half of the product, i.e.  $\frac{\partial \mathbf{K}}{\partial \theta_l}$ , we still divide this into 4 blocks, in accordance with the definition of  $\mathbf{K}$ .

$$\frac{\partial k_{xx}(\mathbf{x}^m, \mathbf{x}^n)}{\partial \theta_l} = \frac{\partial k_f(\mathbf{x}^m, \mathbf{x}^n)}{\partial \theta_l} = \alpha \exp \left( -\frac{1}{2}\|\mathbf{x}^m - \mathbf{x}^n\|_\Gamma^2 \right) \frac{\partial \left( -\frac{1}{2}(\mathbf{x}^m - \mathbf{x}^n)^\top \Gamma (\mathbf{x}^m - \mathbf{x}^n) \right)}{\partial \theta_l}, \quad (29)$$

$$= k_f(\mathbf{x}^m, \mathbf{x}^n) \frac{\partial \left( -\frac{1}{2} \sum_{i=1}^D \gamma_i (x_i^m - x_i^n)^2 \right)}{\partial \theta_l}, \quad (30)$$

$$= -\frac{1}{2} k_f(\mathbf{x}^m, \mathbf{x}^n) (x_l^m - x_l^n)^2, \quad l \in \{1, \dots, D\} \quad (31)$$

$$\frac{\partial k_{xx}(\mathbf{x}^m, \mathbf{x}^n)}{\partial \theta_l} = \frac{\partial k_f(\mathbf{x}^m, \mathbf{x}^n)}{\partial \theta_l} = \exp \left( -\frac{1}{2}\|\mathbf{x}^m - \mathbf{x}^n\|_\Gamma^2 \right), \quad l = D+1 \quad (32)$$

$$\frac{\partial k_{xx}(\mathbf{x}^m, \mathbf{x}^n)}{\partial \theta_l} = \delta_{m,n}, \quad l = D+2 \quad (33)$$

In the following,  $i \in \{1, \dots, D\}$ .

$$\frac{\partial k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n)_i}{\partial \theta_l} = \frac{\partial \left( -\gamma_i (x_{d,i}^m - x_i^n) k_f(\mathbf{x}_d^m, \mathbf{x}^n) \right)}{\partial \theta_l}, \quad (34)$$

$$= -\gamma_i (x_{d,i}^m - x_i^n) \frac{\partial k_f(\mathbf{x}_d^m, \mathbf{x}^n)}{\partial \theta_l} - \delta_{i,l} (x_{d,i}^m - x_i^n) k_f(\mathbf{x}_d^m, \mathbf{x}^n), \quad l \in \{1, \dots, D\} \quad (35)$$

$$\frac{\partial k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n)_i}{\partial \theta_l} = -\gamma_i (x_{d,i}^m - x_i^n) \frac{\partial k_f(\mathbf{x}_d^m, \mathbf{x}^n)}{\partial \theta_l}, \quad (36)$$

$$= -\gamma_i (x_{d,i}^m - x_i^n) \exp \left( -\frac{1}{2}\|\mathbf{x}_d^m - \mathbf{x}^n\|_\Gamma^2 \right), \quad l = D+1, \quad (37)$$

$$\frac{\partial k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n)_i}{\partial \theta_l} = 0, \quad l = D+2 \quad (38)$$

or by using matrix differentiation, we can get the following directly from xx.

$$\frac{\partial k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n)}{\partial \theta_l} = -\mathbf{\Gamma}(\mathbf{x}_d^m - \mathbf{x}^n) \frac{\partial k_f(\mathbf{x}_d^m, \mathbf{x}^n)}{\partial \theta_l} - \mathbf{\Delta}_l(\mathbf{x}_d^m - \mathbf{x}_d^n) k_f(\mathbf{x}_d^m, \mathbf{x}_d^n), \quad (39)$$

$$= \frac{1}{2} k_f(\mathbf{x}_d^m, \mathbf{x}^n) (x_{d,l}^m - x_l^n)^2 \mathbf{\Gamma}(\mathbf{x}_d^m - \mathbf{x}^n) - \mathbf{\Delta}_l(\mathbf{x}_d^m - \mathbf{x}^n) k_f(\mathbf{x}_d^m, \mathbf{x}^n), \quad l \in \{1, \dots, D\}, \quad (40)$$

$$\frac{\partial k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n)}{\partial \theta_l} = -\mathbf{\Gamma}(\mathbf{x}_d^m - \mathbf{x}^n) \exp\left(-\frac{1}{2} \|\mathbf{x}_d^m - \mathbf{x}^n\|_{\mathbf{\Gamma}}^2\right), \quad l = D+1, \quad (41)$$

$$\frac{\partial k_{dx}(\mathbf{x}_d^m, \mathbf{x}^n)}{\partial \theta_l} = \mathbf{0}_{D \times 1} \quad l = D+1, \quad (42)$$

$$(43)$$

where  $\mathbf{\Delta}_l \in \mathbb{R}^{D \times D}$  is a square matrix with the  $l^{\text{th}}$  diagonal element  $\Delta_{l,l} = 1$  and all other elements as 0.

In the following  $i \in \{1, \dots, D\}, j \in \{1, \dots, D\}$ .

$$\frac{\partial k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n)_{i,j}}{\partial \theta_l} = \frac{\partial \left[ \gamma_i (\delta_{i,j} - \gamma_j (x_{d,i}^m - x_{d,i}^n) (x_{d,j}^m - x_{d,j}^n)) k_f(\mathbf{x}_d^m, \mathbf{x}_d^n) \right]}{\partial \theta_l}, \quad (44)$$

$$= \frac{\partial \left[ \gamma_i (\delta_{i,j} - \gamma_j (x_{d,i}^m - x_{d,i}^n) (x_{d,j}^m - x_{d,j}^n)) \right]}{\partial \theta_l} k_f(\mathbf{x}_d^m, \mathbf{x}_d^n) + \gamma_i (\delta_{i,j} - \gamma_j (x_{d,i}^m - x_{d,i}^n) (x_{d,j}^m - x_{d,j}^n)) \frac{\partial k_f(\mathbf{x}_d^m, \mathbf{x}_d^n)}{\partial \theta_l}, \quad (45)$$

$$= [\delta_{i,l} \delta_{j,l} - (\delta_{i,l} \gamma_j + \delta_{j,l} \gamma_i) (x_{d,i}^m - x_{d,i}^n) (x_{d,j}^m - x_{d,j}^n)] k_f(\mathbf{x}_d^m, \mathbf{x}_d^n) + \gamma_i (\delta_{i,j} - \gamma_j (x_{d,i}^m - x_{d,i}^n) (x_{d,j}^m - x_{d,j}^n)) \frac{\partial k_f(\mathbf{x}_d^m, \mathbf{x}_d^n)}{\partial \theta_l}, \quad l \in \{1, \dots, D\} \quad (46)$$

$$\frac{\partial k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n)_{i,j}}{\partial \theta_l} = \gamma_i (\delta_{i,j} - \gamma_j (x_{d,i}^m - x_{d,i}^n) (x_{d,j}^m - x_{d,j}^n)) \frac{\partial k_f(\mathbf{x}_d^m, \mathbf{x}_d^n)}{\partial \theta_l} \quad (47)$$

$$= \gamma_i (\delta_{i,j} - \gamma_j (x_{d,i}^m - x_{d,i}^n) (x_{d,j}^m - x_{d,j}^n)) \exp\left(-\frac{1}{2} \|\mathbf{x}_d^m - \mathbf{x}_d^n\|_{\mathbf{\Gamma}}^2\right), \quad l = D+1 \quad (48)$$

$$\frac{\partial k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n)_{i,j}}{\partial \theta_l} = 0, \quad l = D+2 \quad (49)$$

or by using matrix differentiation, we can get the following directly from xx.

$$\begin{aligned} \frac{\partial k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n)}{\partial \theta_l} &= (\mathbf{\Delta}_l - \mathbf{\Gamma}(\mathbf{x}_d^m - \mathbf{x}_d^n)(\mathbf{x}_d^m - \mathbf{x}_d^n)^\top \mathbf{\Delta}_l - \mathbf{\Delta}_l(\mathbf{x}_d^m - \mathbf{x}_d^n)(\mathbf{x}_d^m - \mathbf{x}_d^n)^\top \mathbf{\Gamma}) k_f(\mathbf{x}_d^m, \mathbf{x}_d^n) \\ &\quad + \mathbf{\Gamma}(\mathbf{I} - (\mathbf{x}_d^m - \mathbf{x}_d^n)(\mathbf{x}_d^m - \mathbf{x}_d^n)^\top \mathbf{\Gamma}) \left(-\frac{1}{2} k_f(\mathbf{x}_d^m, \mathbf{x}_d^n) (x_{d,l}^m - x_{d,l}^n)^2\right), \quad l \in \{1, \dots, D\} \end{aligned} \quad (50)$$

$$\frac{\partial k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n)}{\partial \theta_l} = \mathbf{\Gamma}(\mathbf{I} - (\mathbf{x}_d^m - \mathbf{x}_d^n)(\mathbf{x}_d^m - \mathbf{x}_d^n)^\top \mathbf{\Gamma}) \exp\left(-\frac{1}{2} \|\mathbf{x}_d^m - \mathbf{x}_d^n\|_{\mathbf{\Gamma}}^2\right), \quad l = D+1 \quad (51)$$

$$\frac{\partial k_{dd}(\mathbf{x}_d^m, \mathbf{x}_d^n)}{\partial \theta_l} = \mathbf{0}_{D \times D}, \quad l = D+2 \quad (52)$$

Note, in some cases, the parameter  $\gamma_l$  take an inverse exponentiated square form of

$$\gamma_l = \frac{1}{\exp(\lambda_l)^2}, \quad l \in \{1, \dots, D\} \quad (53)$$

In this case, all the above differentiation should have one more step, i.e.

$$\frac{\partial k(\cdot, \cdot)}{\partial \lambda_l} = \frac{\partial k(\cdot, \cdot)}{\partial \gamma_l} \frac{d\gamma_l}{d\lambda_l} = -2 \exp(\lambda_l)^{-2} \frac{\partial k(\cdot, \cdot)}{\partial \gamma_l} = -2 \gamma_l \frac{\partial k(\cdot, \cdot)}{\partial \gamma_l}, \quad l \in \{1, \dots, D\}. \quad (54)$$