

# Sparse Gaussian

January 2, 2017

## 1 mathematical notations

We first give brief description of mathematical notations will be used throughout the project.

The original data set will be denoted as  $\mathcal{D}$  which consists of  $N$   $d$ -dimensional vectors  $\mathbf{X} = \{\mathbf{x}^{(i)} = (x_1, \dots, x_d) \mid i = 1, \dots, N\}$ . Let the new input data be  $\mathbf{x}^* = (x_1^*, \dots, x_d^*)$ . The pseudo input data set is denoted as  $\bar{\mathcal{D}}$  consists of  $\bar{\mathbf{X}} = \{\bar{\mathbf{x}}^{(i)} = (x_1, \dots, x_d) \mid i = 1, \dots, M\}$ .  $\mathbf{X}$  is paired with target  $\mathbf{Y} = (y^{(1)}, \dots, y^{(N)})$ , notice that  $y^{(i)}$  are scalars.  $\mathbf{x}^*$  is paired with new target  $y^*$ . The underlining latent function is denoted as  $\mathbf{f}(\mathbf{x}) = \mathbf{y}$  and the pseudo one is  $\bar{\mathbf{f}}$ . A Gaussian distribution is denoted as  $\mathcal{N}(\mathbf{f}|\mathbf{m}, \mathbf{V})$  with mean  $\mathbf{m}$  and variance  $\mathbf{V}$ .

## 2 sparse Gaussian process

We first give a zero mean Gaussian prior over the underlining latent function:  $p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_N)$  where  $\mathbf{K}_N$  is our kernel matrix with elements given by,  $[\mathbf{K}_N]_{ij} \equiv K_{\mathbf{x}^{(i)}\mathbf{x}^{(j)}} = K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ :

$$K(\mathbf{x}, \mathbf{x}') = c \exp\left[-\frac{1}{2} \sum_{i=1}^D b_i (x_i^{(n)} - x_i^{(n')})^2\right], \quad \boldsymbol{\theta} \equiv \{c, \mathbf{b}\}, \quad (1)$$

where  $\boldsymbol{\theta}$  is the hyperparameters. We provide noises to  $\mathbf{f}$  such that  $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I})$ . By integrating out the latent function we have the marginal likelihood

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_N + \sigma^2 \mathbf{I}) \quad (2)$$

For prediction, the new input  $\mathbf{x}^*$  conditioning on the observed data and hyperparameters. Let write the joint probability first

$$p(y^*, \mathbf{y}|\mathbf{x}^*, \mathcal{D}, \boldsymbol{\theta}) = \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \mathbf{K}_N + \sigma^2 \mathbf{I} & \mathbf{K}_{\mathbf{x}^*} \\ \mathbf{K}_{\mathbf{x}^*}^T & K_{\mathbf{x}^*\mathbf{x}^*} + \sigma^2 \end{pmatrix}\right), \quad (3)$$

where  $\mathbf{K}_{\mathbf{x}^*} = (K(\mathbf{x}^*, \mathbf{x}^{(1)}), \dots, K(\mathbf{x}^*, \mathbf{x}^{(N)}))$ , i.e.  $[\mathbf{K}_{\mathbf{x}^*}]_i = K(\mathbf{x}^*, \mathbf{x}^{(i)})$ , and  $K_{\mathbf{x}^*\mathbf{x}^*} = K(\mathbf{x}^*, \mathbf{x}^*)$ . Now we can condition on  $\mathbf{y}$  and get

$$\begin{aligned} p(y^*|\mathbf{y}, \mathbf{x}^*, \mathcal{D}, \boldsymbol{\theta}) \\ = \mathcal{N}(y^*|\mathbf{K}_{\mathbf{x}^*}^T (\mathbf{K}_N + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, K_{\mathbf{x}^*\mathbf{x}^*} + \sigma^2 - \mathbf{K}_{\mathbf{x}^*}^T (\mathbf{K}_N + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{\mathbf{x}^*}). \end{aligned} \quad (4)$$

For detailed proof, check Theorem 4.3.1 in Murphy's machine learning a probabilistic perspective.

Now we consider pseudo input  $\bar{\mathbf{X}}$ . Everything still holds except that there are no noises in it. The new input and target pair  $(\mathbf{x}^*, y^*)$  is replaced by one of the actually data set and targets pairs  $(\mathbf{x}^{(i)}, y_i)$ . We therefore just use  $\bar{\mathbf{f}}$  represents the pseudo outputs and  $\bar{\boldsymbol{\theta}}$ , and the single point likelihood is given by

$$p(y|\mathbf{x}, \bar{\mathbf{f}}, \bar{\mathcal{D}}) = \mathcal{N}(y|\mathbf{K}_{\mathbf{x}}^T \mathbf{K}_M^{-1} \bar{\mathbf{f}}, K_{\mathbf{x}\mathbf{x}} + \sigma^2 - \mathbf{K}_{\mathbf{x}}^T \mathbf{K}_M^{-1} \mathbf{K}_{\mathbf{x}}), \quad (5)$$

where  $\mathbf{K}_{\mathbf{x}} = (K(\mathbf{x}, \bar{\mathbf{x}}^{(1)}), \dots, K(\mathbf{x}, \bar{\mathbf{x}}^{(M)}))$ , i.e.  $[\mathbf{K}_{\mathbf{x}}]_i = K(\mathbf{x}, \bar{\mathbf{x}}^{(i)})$ . As the target data are i.i.d given the inputs, the complete data likelihood is given by

$$p(\mathbf{y}|\mathbf{X}, \bar{\mathbf{f}}, \bar{\mathcal{D}}) = \prod_{i=1}^N p(y_i|\mathbf{x}^{(i)}, \bar{\mathbf{f}}, \bar{\mathcal{D}}) = \mathcal{N}(\mathbf{y}|\mathbf{K}_{NM} \mathbf{K}_M^{-1} \bar{\mathbf{f}}, \boldsymbol{\Lambda} + \sigma^2 \mathbf{I}), \quad (6)$$

where  $\boldsymbol{\Lambda} = \text{diag}(\boldsymbol{\lambda})$ ,  $\lambda_i = K_{\mathbf{x}^{(i)}\mathbf{x}^{(i)}} - \mathbf{K}_{\mathbf{x}^{(i)}}^T \mathbf{K}_M^{-1} \mathbf{K}_{\mathbf{x}^{(i)}}$ , is a  $N \times N$  diagonal matrix, and  $[\mathbf{K}_{NM}]_{ij} = K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ . Together with a Gaussian prior,  $p(\bar{\mathbf{f}}|\bar{\mathbf{X}}) = \mathcal{N}(\bar{\mathbf{f}}|\mathbf{0}, \mathbf{K}_M)$