

# Sparse Gaussian

January 5, 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,  $[\mathbf{K}_N]_{ij} \equiv K_{\mathbf{x}^{(i)}\mathbf{x}^{(j)}} = K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ : Notice that this is the case that we have same number of  $\mathbf{x}^{(i)}, \mathbf{x}^{(j)}$ . In case of different sizes, we use  $\mathbf{K}_{NM}$ , i.e.  $N$  rows for the first input matrix,  $M$  rows for the second input matrix.

$$K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = c \exp\left[-\frac{1}{2} \sum_{k=1}^D b_k (x_k^{(i)} - x_k^{(j)})^2\right], \quad \boldsymbol{\theta} \equiv \{c, \mathbf{b}\}, \quad (1)$$

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

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

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

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

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

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

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

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

where  $\boldsymbol{\Lambda} = \text{diag}(\lambda)$ ,  $\lambda_i = K_{\mathbf{x}^{(i)}\mathbf{x}^{(i)}} - \mathbf{K}_{\mathbf{x}^{(i)}M}\mathbf{K}_M^{-1}\mathbf{K}_{\mathbf{x}^{(i)}M}^T$ , is a  $N \times N$  diagonal matrix, and  $[\mathbf{K}_{NM}]_{ij} = K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ . Together with a Gaussian prior,  $p(\bar{\mathbf{f}}|\bar{\mathbf{X}}) = \mathcal{N}(\bar{\mathbf{f}}|\mathbf{0}, \mathbf{K}_M)$ , integrate over Eq.6 we have the SPGP marginal likelihood over pseudo inputs

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \bar{\mathbf{X}}) &= \int p(\mathbf{y}|\mathbf{X}, \bar{\mathbf{f}}, \bar{\mathbf{X}}) p(\bar{\mathbf{f}}|\bar{\mathbf{X}}) d\bar{\mathbf{f}} \\ &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_{NM}\mathbf{K}_M^{-1}\mathbf{K}_{MN} + \boldsymbol{\Lambda} + \sigma^2\mathbf{I}). \end{aligned} \quad (7)$$

Same as we have done from Eq.3 to Eq.4, we first write the joint probability of  $y^*, \mathbf{y}$

$$\begin{aligned} p(y^*, \mathbf{y}|\mathbf{x}^*, \mathbf{X}, \bar{\mathbf{X}}) \\ = \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} K_{\mathbf{x}^*\mathbf{x}^*} + \sigma^2 & \mathbf{K}_{\mathbf{x}^*M}\mathbf{K}_M^{-1}\mathbf{K}_{MN} \\ (\mathbf{K}_{\mathbf{x}^*M}\mathbf{K}_M^{-1}\mathbf{K}_{MN})^T & \mathbf{K}_{NM}\mathbf{K}_M^{-1}\mathbf{K}_{MN} + \boldsymbol{\Lambda} + \sigma^2\mathbf{I} \end{pmatrix}\right), \end{aligned} \quad (8)$$

where  $\mathbf{K}_{\mathbf{x}^*M} = (K(\mathbf{x}^*, \bar{\mathbf{x}}^{(1)}), \dots, K(\mathbf{x}^*, \bar{\mathbf{x}}^{(M)}))$ , i.e.  $[\mathbf{K}_{\mathbf{x}^*M}]_i = K(\mathbf{x}^*, \bar{\mathbf{x}}^{(i)})$ . From now on we let

$$\mathbf{Q}_{\mathbf{X}, \mathbf{X}'} \equiv \mathbf{Q}(\mathbf{X}, \mathbf{X}') = \mathbf{K}_{\mathbf{X}M}\mathbf{K}_M^{-1}\mathbf{K}_{MX'} \quad (9)$$

$$\mathbf{Q}_N = \mathbf{K}_{NM}\mathbf{K}_M^{-1}\mathbf{K}_{MN}, \quad (10)$$

Also, remember that here  $N$  and  $M$  represents input and pseudo input data set, matrices, as input matrices of  $\mathbf{K}$ , respectively. And after conditioning on  $\mathbf{y}$ , we have the SPGP predictive distribution

$$p(y^*|\mathbf{y}, \mathbf{x}^*, \mathbf{X}, \bar{\mathbf{X}}) = \mathcal{N}(\mu^*, \sigma^{*2}) \quad (11)$$

$$\begin{aligned} \mu^* &= \mathbf{Q}_{\mathbf{x}^*N}(\mathbf{Q}_N + \boldsymbol{\Lambda} + \sigma^2\mathbf{I})^{-1}\mathbf{y} \\ \sigma^{*2} &= K_{\mathbf{x}^*\mathbf{x}^*} - \mathbf{Q}_{\mathbf{x}^*N}(\mathbf{Q}_N + \boldsymbol{\Lambda} + \sigma^2\mathbf{I})^{-1}\mathbf{Q}_{N\mathbf{x}^*} + \sigma^2. \end{aligned} \quad (12)$$

The pseudo input  $\bar{\mathbf{C}}$  and hyperparameters  $\boldsymbol{\Theta} = \{\boldsymbol{\theta}, \sigma^2\}$ , this can be done by maximizing Eq.7.

Some simplification for matrix inversion. First from matrix inversion lemma

$$(\mathbf{A} + \mathbf{U}\mathbf{B}\mathbf{U}^T)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{B}^{-1} + \mathbf{U}^T\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{U}^T\mathbf{A}^{-1} \quad (13)$$

$$\det(\mathbf{A} + \mathbf{U}\mathbf{B}\mathbf{U}^T) = \det(\mathbf{A})\det(\mathbf{B})\det(\mathbf{B}^{-1} + \mathbf{U}^T\mathbf{A}^{-1}\mathbf{U}), \quad (14)$$

we can rewrite following

$$(\mathbf{K}_{NM}\mathbf{K}_M^{-1}\mathbf{K}_{MN} + \mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1} \quad (15)$$

$$= (\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1} - (\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1}\mathbf{K}_{NM}\mathbf{B}^{-1}\mathbf{K}_{MN}(\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1}, \quad (16)$$

where  $\mathbf{B} = \mathbf{K}_M + \mathbf{K}_{MN}(\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1}\mathbf{K}_{NM}$ . Now matrix inversion only happens to  $(\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1}$  which is  $\mathcal{O}(N)$  as it is diagonal. Now Eq.12 become

$$\begin{aligned} \mu^* &= \mathbf{K}_{\mathbf{x}^*M}\mathbf{B}^{-1}\mathbf{K}_{MN}(\mathbf{\Lambda} + \sigma^2\mathbf{I})^{-1}\mathbf{y} \\ \sigma^{*2} &= \mathbf{K}_{\mathbf{x}^*\mathbf{x}^*} - \mathbf{K}_{\mathbf{x}^*M}(\mathbf{K}_M^{-1} - \mathbf{B}^{-1})\mathbf{K}_{M\mathbf{x}^*} + \sigma^2. \end{aligned} \quad (17)$$

### 3 implementation

Rewrite

$$\sigma^2\mathbf{\Gamma} = \mathbf{\Lambda} + \sigma^2\mathbf{I}, \quad (18)$$

and suppressing data dependency of Eq.7, we have

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{Q}_N + \sigma^2\mathbf{\Gamma}). \quad (19)$$

We maximize  $\mathcal{L} = -\log p(\mathbf{y})$

$$\mathcal{L} = \frac{1}{2}(\log(\det(\mathbf{Q}_N + \sigma^2\mathbf{\Gamma})) + \mathbf{y}(\mathbf{Q}_N + \sigma^2\mathbf{\Gamma})^{-1}\mathbf{y}^T + N \log(2\pi)). \quad (20)$$

Let

$$\mathcal{L}_1 = \log(\det(\mathbf{Q}_N + \sigma^2\mathbf{\Gamma})) \quad (21)$$

$$\mathcal{L}_2 = \mathbf{y}(\mathbf{Q}_N + \sigma^2\mathbf{\Gamma})^{-1}\mathbf{y}^T. \quad (22)$$

Use matrix inversion lemma again, we have

$$\begin{aligned} \mathcal{L}_1 &= \log(\det(\mathbf{K}_M + \sigma^{-2}\mathbf{K}_{MN}\mathbf{\Gamma}^{-1}\mathbf{K}_{NM}) \det(\mathbf{K}_M^{-1}) \det(\sigma^2\mathbf{\Gamma})) \\ &= \log(\det(\mathbf{A})) - \log(\det(\mathbf{K}_M)) + \log(\det(\mathbf{\Gamma})) + (N - M) \log(\sigma^2) \end{aligned} \quad (23)$$

$$\mathcal{L}_2 = \sigma^{-2}\mathbf{y}(\mathbf{\Gamma}^{-1} - \mathbf{\Gamma}^{-1}\mathbf{K}_{NM}\mathbf{A}^{-1}\mathbf{K}_{MN}\mathbf{\Gamma}^{-1})\mathbf{y}^T \quad (24)$$

$$= \sigma^{-2}(\|\mathbf{\Gamma}^{-\frac{1}{2}}\mathbf{y}^T\|^2 - \|\mathbf{A}^{-\frac{1}{2}}(\mathbf{\Gamma}^{-\frac{1}{2}}\mathbf{K}_{NM})^T\mathbf{\Gamma}^{-\frac{1}{2}}\mathbf{y}^T\|^2) \quad (25)$$

where  $\mathbf{A} = \sigma^2\mathbf{K}_M + \mathbf{K}_{MN}\mathbf{\Gamma}^{-1}\mathbf{K}_{NM}$ . The final negative log marginal likelihood is

$$\mathcal{L} = \frac{1}{2}(\mathcal{L}_1 + \mathcal{L}_2 + N \log(2\pi)). \quad (26)$$

#### 3.1 matrix derivatives

Let  $\mathbf{A}$  be a matrix with underlining parameter  $\theta$ . The derivative of the inverse matrix w.r.t  $\theta$  is

$$\frac{\partial}{\partial \theta}\mathbf{A}^{-1} = -\mathbf{A}^{-1}\frac{\partial \mathbf{A}}{\partial \theta}\mathbf{A}^{-1}, \quad (27)$$

where the partial derivative takes elementwise. If  $\mathbf{A}$  is positive definite symmetric, the derivative of the log determinant is

$$\frac{\partial}{\partial \theta} \log(\det(\mathbf{A})) = \text{tr}(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \theta}) \quad (28)$$

First ignoring the noise variance  $\sigma^2$ , do partial derivative on  $\boldsymbol{\theta}$ , we have

$$\dot{\mathcal{L}}_1 = \text{tr}(\mathbf{A}^{-\frac{1}{2}} \dot{\mathbf{A}} \mathbf{A}^{-\frac{T}{2}}) - \text{tr}(\mathbf{K}_M^{-\frac{1}{2}} \dot{\mathbf{K}}_M \mathbf{K}_M^{-\frac{T}{2}}) + \text{tr}(\boldsymbol{\Gamma}^{-\frac{1}{2}} \dot{\boldsymbol{\Gamma}} \boldsymbol{\Gamma}^{-\frac{1}{2}}) \quad (29)$$