

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, $[\mathbf{K}_N]_{nn'} = K(\mathbf{x}, \mathbf{x}')$:

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

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.