

GP Write Up

September 26, 2018

Convolutions on the Kernel matrix

1. Introduction

Predictive mean for a GP is:

$$\mu_L = \mathbf{K}_* (\mathbf{K}_X + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} \quad (1)$$

the parameters for the covariance function $k(\cdot, \cdot)$ are optimized by maximizing the log marginal likelihood given by:

$$\log p(\mathbf{y} | X) = -\frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K} + \sigma_n^2 \mathbf{I}| - \frac{n}{2} \log 2\pi \quad (2)$$

Equation (1) involves computing inverses and this is computationally expensive and we need to circumvent this problem subject to some assumptions on the training data.

We will see one of the ways of avoiding the inverse when the training data is on a grid and the kernel is stationary.

Say:

$$\mathbf{v} = (\mathbf{K}_X + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} \quad (3)$$

We can avoid calculating the inverse by solving the linear system:

$$(\mathbf{K}_X + \sigma_n^2 \mathbf{I}) \mathbf{v} = \mathbf{y} \quad (4)$$

for \mathbf{v} .

2. Convolutions on a grid

When the training data $X \in \mathbb{R}^{n \times n}$, say $n = 2$, and with a stationary kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = f'(\mathbf{x}_i - \mathbf{x}_j) = f(x_i^{(0)} - x_j^{(0)}, x_i^{(1)} - x_j^{(1)}) \quad (5)$$

equation (4) can be written as:

$$\mathbf{y} = \begin{bmatrix} f(0,0) & f(0,-1) & f(-1,0) & f(-1,-1) \\ f(0,1) & f(0,0) & f(-1,1) & f(-1,0) \\ f(1,0) & f(1,-1) & f(0,0) & f(0,-1) \\ f(1,1) & f(1,0) & f(0,1) & f(0,0) \end{bmatrix} \begin{bmatrix} v_{0,0} \\ v_{0,1} \\ v_{1,0} \\ v_{1,1} \end{bmatrix} + \sigma_n^2 \mathbf{I} \cdot \mathbf{v} \quad (6)$$

This matrix multiplication can be represented as a convolution:

$$\mathbf{y} = \mathbf{K}' * \mathbf{v}' + \sigma_n^2 \mathbf{I} \cdot \mathbf{v} \quad (7)$$

where:

$$\mathbf{K}' = \begin{bmatrix} f(-1, -1) & f(-1, 0) & f(-1, 1) \\ f(0, -1) & f(0, 0) & f(0, 1) \\ f(1, -1) & f(1, 0) & f(1, 1) \end{bmatrix} \quad (8)$$

and

$$\mathbf{v}' = \begin{bmatrix} v_{1,1} & v_{1,0} \\ v_{0,1} & v_{0,0} \end{bmatrix} \quad (9)$$

Lets see convolutions on a grid where $n = 3$:

like in equation (8)

$$\mathbf{K}' = \begin{bmatrix} f(-2, -2) & f(-2, -1) & f(-2, 0) & f(-2, 1) & f(-2, 2) \\ f(-1, -2) & f(-1, -1) & f(-1, 0) & f(-1, 1) & f(-1, 2) \\ f(0, -2) & f(0, -1) & f(0, 0) & f(0, 1) & f(0, 2) \\ f(1, -2) & f(1, -1) & f(1, 0) & f(1, 1) & f(1, 2) \\ f(2, -2) & f(2, -1) & f(2, 0) & f(2, 1) & f(2, 2) \end{bmatrix} \quad (10)$$

and like in equation (9)

$$\mathbf{v}' = \begin{bmatrix} v_{2,2} & v_{2,1} & v_{2,0} \\ v_{1,2} & v_{1,1} & v_{1,0} \\ v_{0,2} & v_{0,1} & v_{0,0} \end{bmatrix} \quad (11)$$

The figure displays three hand-drawn 5x5 grids, each representing a matrix decomposition. The columns are labeled -2, -1, 0, 1, 2 and the rows are labeled -2, -1, 0, 1, 2. The 3x3 blocks are highlighted with red boxes.

Top-Left Grid: The 3x3 block in the top-left corner (rows -2 to 0, columns -2 to 0) contains 0s. The rest of the matrix contains the values $k(i, j)$.

Top-Right Grid: The 3x3 block in the top-right corner (rows -2 to 0, columns 1 to 2) contains 0s. The rest of the matrix contains the values $k(i, j)$.

Bottom Grid: The 3x3 block in the bottom-right corner (rows 1 to 2, columns 1 to 2) contains 0s. The rest of the matrix contains the values $k(i, j)$.

Consider the scenario where some observations are missing on the grid and we fill the missing observations with samples from $\mathcal{N}(0, \sigma_w^2 I)$ where $\sigma_w \rightarrow \infty$. Lets say we have M observations and W missing observations for a total of N observations.

from equation (3):

$$(\mathbf{K}_N + \mathbf{D}_N)^{-1} \mathbf{y} = \mathbf{v} \quad (14)$$

we can use the block matrix inversion algorithm to see what this inverse leads to:

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} (A - BD^{-1}C)^{-1} & -A^{-1}B(I - D^{-1}CA^{-1}B)^{-1}D^{-1} \\ -D^{-1}C(A - BD^{-1}C)^{-1} & (I - D^{-1}CA^{-1}B)^{-1}D^{-1} \end{bmatrix} \quad (15)$$

where:

$$A = K_M + \sigma_M^2 I_M \quad (16)$$

$$B = K_{MW} \quad (17)$$

$$C = K_{MW}^T \quad (18)$$

$$D = K_W + \sigma_W^2 I_W = \sigma_W^2 (\sigma_W^{-2} K_W + I_W) \quad (19)$$

$$\lim_{\sigma_W \rightarrow \infty} D^{-1} = 0 \quad (20)$$

Therefore equation (3) becomes

$$(\mathbf{K}_M + \sigma_M^2 I_M)^{-1} \mathbf{y}_M = \mathbf{v}_M \quad (21)$$

and

$$\mathbf{y}_M = (\mathbf{K}_M + \sigma_M^2 I_M) \mathbf{v}_M \quad (22)$$

from equation (1) the predictive mean for the entire grid will be:

$$\mu_L = \mathbf{K}_{NM} \mathbf{v}_M \quad (23)$$

4. Loss Function

For a function:

$$f(x) = \frac{1}{2} x^T A x - x^T b + c \quad (24)$$

if A is symmetric the stationary point is at:

$$Ax = b \quad (25)$$

and if A is positive definite then the solution to equation (26) is the minimum. So the solution to equation (26) can be obtained by minimizing equation (24).

So the optimization problem is:

$$\min_v \frac{1}{2} \mathbf{v}^T \mathbf{K} \mathbf{v} - \mathbf{v}^T \mathbf{y} \quad (26)$$

which simplifies to:

$$\frac{1}{2} \mathbf{y}_m' (\mathbf{K}_M + \mathbf{D}_M)^{-1} \mathbf{y}_m' - \mathbf{y}_m' (\mathbf{K}_M + \mathbf{D}_M)^{-1} \mathbf{y}_m \quad (27)$$

