

Explaining SPGP (4)

Bishop	(article) SPGP	What
x_n	x_n	input values
y_n	f_n	GP function output
t_n	y_n	observation / target value

By GP def: $p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{m}, \mathbf{K})$ where $K_{ij} = K(x_i, x_j)$ is always the case.

• No prior knowledge of $f_n \Rightarrow$ can take $\mathbf{m} = \mathbf{0}$. (bishop p 305) $\Rightarrow p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$ p 305

Note in GP for regression (SPGP p 2, Bishop p 306):

Assume a Gaussian noise model: $y_n = f_n + \epsilon_n$ where $\epsilon_n = \mathcal{N}(\epsilon_n | 0, G^2)$
GP function values
observation / tgt value

$$p(y|\mathbf{f}) = \mathcal{N}(y|\mathbf{f}, G^2 \mathbf{I})$$

$$\text{then } p(y|\mathbf{X}, G^2) = \int_{\mathbf{f}} p(y|\mathbf{f}) p(\mathbf{f}|\mathbf{X}) d\mathbf{f} = \{\text{Bishop 2.115 p 93}\}$$

$$= \mathcal{N}(\mathbf{I}\mathbf{0} + \mathbf{0}, G^2 \mathbf{I} + \mathbf{I}\mathbf{K}\mathbf{I}^T) = \mathcal{N}(\mathbf{0}, \mathbf{K} + G^2 \mathbf{I})$$

Let's consider pseudo examples only, $\bar{\mathbf{x}}, \bar{\mathbf{f}}$ generated from \mathbf{X}, \mathbf{f} . Since pseudo inputs are constructed from the input data (they are 'artificial' training examples) we can simply impose for the pseudo inputs that their observations \bar{y} are fully explained by the GP (no noise addition from GP output $\bar{\mathbf{f}}$ to 'observed' value \bar{y} for the pseudo data)

$$\left. \begin{aligned} p(y|\bar{\mathbf{x}}) & \left\{ p(\bar{y}|\bar{\mathbf{f}}) = \mathcal{N}(\bar{y}|\bar{\mathbf{f}}, 0 \cdot \mathbf{I}) \right\} \text{ where } \\ p(\bar{\mathbf{f}}|\bar{\mathbf{x}}) & \left\{ p(\bar{\mathbf{f}}|\bar{\mathbf{x}}) = \mathcal{N}(\bar{\mathbf{f}}|\mathbf{0}, \mathbf{K}_M) \right\} \end{aligned} \right\} \Rightarrow \{\text{from Bishop 2.115 p 93}\}$$

$$\Rightarrow p(\bar{y}|\bar{\mathbf{x}}) = \int_{\bar{\mathbf{f}}} p(\bar{y}|\bar{\mathbf{f}}) p(\bar{\mathbf{f}}|\bar{\mathbf{x}}) d\bar{\mathbf{f}} = \mathcal{N}(\mathbf{I}\mathbf{0} + \mathbf{0}, \mathbf{0}\mathbf{I} + \mathbf{I}\mathbf{K}_M\mathbf{I}^T) = \mathcal{N}(\mathbf{0}, \mathbf{K}_M)$$

predictive distribution

Now we observe the input values of one of the real training examples \mathbf{x} . We want to find $p(y)$ conditioned on the pseudo inputs $\bar{\mathbf{x}}, \bar{\mathbf{f}}$ and \mathbf{x} (and any hyperparameters θ), $p(y|\mathbf{x}, \bar{\mathbf{x}}, \bar{\mathbf{f}})$. Note that the new example is noisy $y = f + \epsilon \Rightarrow p(f|\mathbf{x}) = \mathcal{N}(0, k(\mathbf{x}, \mathbf{x}))$, $p(y|\mathbf{x}) = \mathcal{N}(f, G^2) \Rightarrow p(y|\mathbf{x}) = \mathcal{N}(0, G^2 + k(\mathbf{x}, \mathbf{x}))$
 $\mathbf{I} = \bar{\mathbf{y}}$

$$\text{Form: } \mathbf{y}_{M+1} = [y \ y' \ y'' \ \dots \ y^M]^T$$

$$\mathbf{K}_{M+1} = \begin{bmatrix} k(\mathbf{x}, \mathbf{x}) + G^2 & k(\mathbf{x}, \mathbf{x}^{1:M})^T \\ k(\mathbf{x}^{1:M}, \mathbf{x}) & \mathbf{K}_M \end{bmatrix} = \left\{ \begin{aligned} & \text{using SPGP notation } K_{xx} = k(\mathbf{x}, \mathbf{x}) \\ & \text{and noting } k(\mathbf{x}, \mathbf{x}^{1:M}) = k(\mathbf{x}^{1:M}, \mathbf{x}) \\ & \text{(since cov matrix must be symmetric)} \end{aligned} \right\} = \begin{bmatrix} K_{xx} + G^2 & k_x^T \\ k_x & K_M \end{bmatrix}$$

$\begin{matrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{matrix}$

then from Bishop p 85-87:

$$p(y|\mathbf{x}, \bar{\mathbf{x}}, \bar{\mathbf{f}}) = \mathcal{N}(\mu_{alb}, \Sigma_{alb}) \text{ where } \mu_{alb} = 0 + k_x^T K_M^{-1} (\bar{\mathbf{f}} - 0)$$

$$\Sigma_{alb} = K_{xx} + G^2 - k_x^T K_M^{-1} k_x$$

$$\Rightarrow p(y|\mathbf{x}, \bar{\mathbf{x}}, \bar{\mathbf{f}}) = \mathcal{N}(y|k_x^T K_M^{-1} \bar{\mathbf{f}}, K_{xx} - k_x^T K_M^{-1} k_x + G^2) \text{ i.e. SPGP (4)}$$

$$p(y|x, \bar{X}, \bar{f}) = \mathcal{N}(k_x^T K_M^{-1} \bar{f}, K_{xx} - k_x^T K_M^{-1} k_x + \sigma^2) \quad (4)$$

The idea now is to express the joint probability of all observed values y_n for all training data $n=1, \dots, N$, conditioned on their resp input values x_n and the pseudo data points \bar{X}, \bar{f} . If we can maximize the joint probability of all training data observations wrt the pseudo points, then we know that the pseudo points selected are the optimal for this particular training data X, f (and n.o. pseudo data points M).

If example data is i.i.d: (where we now denote pseudo inputs $\bar{x}_m, m=1, \dots, M$. Training data inputs $x_n, n=1, \dots, N$):

$$p(y|X, \bar{X}, \bar{f}) = \prod_{n=1}^N p(y_n|x_n, \bar{X}, \bar{f}) = \{\text{expressing as multivariate Gaussian}\} = p(y|\mu_y, \Sigma_y) \text{ where}$$

$$\mu_y = \begin{bmatrix} k_{x_1}^T K_M^{-1} \bar{f} \\ k_{x_2}^T K_M^{-1} \bar{f} \\ \vdots \\ k_{x_N}^T K_M^{-1} \bar{f} \end{bmatrix} = \begin{bmatrix} -k_{x_1}^T - \\ -k_{x_2}^T - \\ \vdots \\ -k_{x_N}^T - \end{bmatrix} \begin{matrix} \in \mathbb{R}^{N \times M} \\ \downarrow \\ \in \mathbb{R}^{M \times 1} \end{matrix} K_M^{-1} \bar{f} = \underbrace{\begin{bmatrix} k(x_1, \bar{x}_1) & k(x_1, \bar{x}_2) & \dots & k(x_1, \bar{x}_M) \\ k(x_2, \bar{x}_1) & k(x_2, \bar{x}_2) & \dots & k(x_2, \bar{x}_M) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_N, \bar{x}_1) & k(x_N, \bar{x}_2) & \dots & k(x_N, \bar{x}_M) \end{bmatrix}}_{K_{NM}} K_M^{-1} \bar{f} = K_{NM} K_M^{-1} \bar{f} \quad \Rightarrow$$

Σ_y is a diagonal $N \times N$ matrix $= \Lambda + \sigma^2 I$ where $\Lambda_{nn} = K_{x_n x_n} - k_{x_n}^T K_M^{-1} k_{x_n}$

$$p(y|X, \bar{X}, \bar{f}) = \mathcal{N}(y | K_{NM} K_M^{-1} \bar{f}, \Lambda + \sigma^2 I) \quad \text{ie SPGP (5)}$$

If we put a prior on $\bar{\mathbf{f}}$ it turns out (which will be shown) that we can integrate out $\bar{\mathbf{f}}$, thereby reducing the n.o. parameters for which to optimize the joint probability.

$$p(\bar{\mathbf{f}}|\bar{\mathbf{x}}) = \mathcal{N}(\bar{\mathbf{f}}|\mathbf{0}, \mathbf{K}_M) \quad (6)$$

In fact, this was already used in the derivation of (4). For motivation, see corresponding text.

Summarizing, we have:

$$\begin{aligned} p(\bar{\mathbf{f}}|\bar{\mathbf{x}}) &= \mathcal{N}(\bar{\mathbf{f}}|\mathbf{0}, \mathbf{K}_M) \\ p(\mathbf{y}|\bar{\mathbf{f}}, \mathbf{x}, \bar{\mathbf{x}}) &= \mathcal{N}(\mathbf{y}|\underbrace{\mathbf{K}_{NM} \mathbf{K}_M^{-1} \bar{\mathbf{f}}}_{\mathbf{A}^T \mathbf{b}}, \underbrace{\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I}}_{\mathbf{L}^{-1}}) \end{aligned}$$

Bishop notation

We can use Bayes' rule to find $p(\bar{\mathbf{f}}|\mathbf{y}, \mathbf{x}, \bar{\mathbf{x}}) = p(\bar{\mathbf{f}}|\mathbf{0}, \bar{\mathbf{x}})$. Note $p(\bar{\mathbf{f}}|\mathbf{y}, \mathbf{x}, \bar{\mathbf{x}}) = \frac{p(\mathbf{y}|\bar{\mathbf{f}}, \mathbf{x}, \bar{\mathbf{x}}) p(\bar{\mathbf{f}}|\bar{\mathbf{x}}, \mathbf{x})}{p(\mathbf{y}|\bar{\mathbf{x}}, \mathbf{x})}$

$$= \frac{p(\mathbf{y}|\bar{\mathbf{f}}, \mathbf{x}, \bar{\mathbf{x}}) p(\bar{\mathbf{f}}|\bar{\mathbf{x}})}{Z} \quad \text{where } Z \text{ is a constant normalization factor (does not depend on } \bar{\mathbf{f}} \text{)}.$$

Using Bishop 2.116 p93:

$$\underbrace{\Sigma}_{\text{posterior covariance}} = \underbrace{\left(\mathbf{K}_M^{-1} + \underbrace{(\mathbf{K}_{NM} \mathbf{K}_M^{-1})^T}_{\mathbf{A}^T} \right)}_{\mathbb{R}^{M \times M}} \underbrace{\left(\underbrace{\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I}}_{\mathbf{L}} \right)^{-1}}_{\mathbb{R}^{N \times N}} \underbrace{\left(\mathbf{K}_{NM} \mathbf{K}_M^{-1} \right)}_{\mathbf{A}}^{-1} = \underbrace{\left(\mathbf{K}_M^{-1} + \mathbf{K}_M^{-1} \mathbf{K}_{MN} (\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I})^{-1} \mathbf{K}_{NM} \mathbf{K}_M^{-1} \right)^{-1}}_{\mathbb{R}^{M \times M}}$$

Bishop notation

Note \mathbf{K}_M is symmetric $\Rightarrow \mathbf{K}_M^{-1}$ is symmetric.

Consider $\mathbf{P} = \mathbf{K}_M (\mathbf{K}_M + \mathbf{K}_{MN} (\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I})^{-1} \mathbf{K}_{NM})^{-1} \mathbf{K}_M$ (the covariance expression in SPGP (7))

$$\text{Note } \mathbf{P}^{-1} = \mathbf{K}_M^{-1} (\mathbf{K}_M + \mathbf{K}_{MN} (\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I})^{-1} \mathbf{K}_{NM}) \mathbf{K}_M^{-1}$$

$$= \mathbf{K}_M^{-1} \mathbf{K}_M \mathbf{K}_M^{-1} + \mathbf{K}_M^{-1} \mathbf{K}_{MN} (\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I})^{-1} \mathbf{K}_{NM} \mathbf{K}_M^{-1}$$

$$= \mathbf{K}_M^{-1} + \mathbf{K}_M^{-1} \mathbf{K}_{MN} (\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I})^{-1} \mathbf{K}_{NM} \mathbf{K}_M^{-1} = \Sigma^{-1}$$

so \mathbf{P} and Σ share the same inverse $\Rightarrow \Sigma = \mathbf{P}$ (since $\mathbf{P} = \mathbf{P} \mathbf{I} = \mathbf{P} (\mathbf{P}^{-1} \Sigma) = \mathbf{I} \Sigma = \Sigma$)

posterior mean

The mean is given by (Bishop notation): $\Sigma (\mathbf{A}^T \mathbf{L} (\mathbf{y} - \mathbf{b})) = \Sigma \mathbf{K}_M^{-1} \mathbf{K}_{MN} (\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I})^{-1} \mathbf{y}$

$$\text{Let } \mathbf{Q}_M = \mathbf{K}_M + \mathbf{K}_{MN} (\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I})^{-1} \mathbf{K}_{NM} \Rightarrow \Sigma = \mathbf{K}_M \mathbf{Q}_M^{-1} \mathbf{K}_M$$

$$\text{mean} = \mathbf{K}_M \mathbf{Q}_M^{-1} \mathbf{K}_{MN} (\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I})^{-1} \mathbf{y}$$

Thus $p(\bar{\mathbf{f}}|\mathbf{y}, \mathbf{x}, \bar{\mathbf{x}}) = p(\bar{\mathbf{f}}|\mathbf{0}, \bar{\mathbf{x}}) = \mathcal{N}(\bar{\mathbf{f}}|\mathbf{K}_M \mathbf{Q}_M^{-1} \mathbf{K}_{MN} (\mathbf{\Lambda} + \mathbf{G}^2 \mathbf{I})^{-1} \mathbf{y}, \mathbf{K}_M \mathbf{Q}_M^{-1} \mathbf{K}_M)$ ie SPGP (7)

$$\text{SPGP (4)} + (7) \Rightarrow (8)$$

Finding the predictive distribution y_* for a new point x_* given the training data X, y and pseudo inputs \bar{X} .

From (4) we have the predictive distribution over the target value for a new point given the pseudo inputs \bar{X} and pseudo targets \bar{f} ($= \bar{y}$ since no noise for pseudo points)

$$(4): p(y_* | x_*, \bar{X}, \bar{f}) = \mathcal{N}(y_* | \underbrace{k_{x_*}^T K_M^{-1} \bar{f}}_{\substack{A \quad b=0 \\ L^{-1}}}, \underbrace{K_{x_* x_*} - k_{x_*}^T K_M^{-1} k_{x_*} + G^2}_{L^{-1}}) \quad \leftarrow \text{Bishop's notation 2.114}$$

From (7) we have the distribution over the pseudo targets \bar{f} given the training data X, y and pseudo inputs \bar{X}

$$(7): p(\bar{f} | \mathcal{D}, \bar{X}) = \mathcal{N}(\bar{f} | \underbrace{K_M Q_M^{-1} K_{MN} (\Lambda + G^2 I)^{-1} y}_M, \underbrace{K_M Q_M^{-1} K_M}_{\Lambda^{-1}}) \quad \leftarrow \text{Bishop's notation 2.113}$$

$$\text{Note that } p(y_* | x_*, \mathcal{D}, \bar{X}) = \int_{\bar{f}} p(y_* | x_*, \bar{f}, \mathcal{D}, \bar{X}) d\bar{f}$$

$$= \int_{\bar{f}} p(y_* | x_*, \bar{X}, \bar{f}) p(\bar{f} | \mathcal{D}, \bar{X}) d\bar{f}$$

[Note:
 k_{x_*} here is k_* in SPGP
 $K_{x_* x_*}$ " K_{**} "]

{Using Bishop 2.115 p.93,}

$$p(y_* | x_*, \mathcal{D}, \bar{X}) = \mathcal{N}(y_* | \mu_*, G_*^2) \quad \text{where}$$

Bishop's notation

$$\mu_* = A\mu + b$$

$$G_*^2 = L^{-1} + A\Lambda^{-1}A^T$$

$$= k_{x_*}^T K_M^{-1} K_{MN} Q_M^{-1} K_{MN} (\Lambda + G^2 I)^{-1} y = k_{x_*}^T Q_M^{-1} K_{MN} (\Lambda + G^2 I)^{-1} y$$

$$= K_{x_* x_*} - k_{x_*}^T K_M^{-1} k_{x_*} + G^2 + k_{x_*}^T K_M^{-1} K_{MN} Q_M^{-1} K_M K_M^{-1} k_{x_*}$$

$$= K_{x_* x_*} - k_{x_*}^T K_M^{-1} k_{x_*} + k_{x_*}^T Q_M^{-1} k_{x_*} + G^2$$

$$= \underline{K_{x_* x_*} - k_{x_*}^T (K_M^{-1} + Q_M^{-1}) k_{x_*} + G^2}$$

ie SPGP (8)

The predictive distribution is conditioned on the pseudo input locations \bar{X} (and hyper-parameters Θ). To find them we use (9).

$$SPGP (5) + (6) \Rightarrow (9)$$

We can find Maximum Likelihood estimates for $\Theta = \{c, \underbrace{b}_{\text{used in Kernel}}, G\}$ and the pseudo input locations \bar{X} .

The ML estimations are done using the marginal probability distribution of training data target values (observations) y , given training data inputs X , pseudo inputs \bar{X} and hyperparameters Θ .

$$\text{From (5): } p(y|X, \bar{X}, \bar{f}) = \mathcal{N}(y | \underbrace{K_{NM} K_M^{-1}}^A \bar{f}, \underbrace{\Lambda + G^2 I}_{b=0, L^{-1}}) \leftarrow \text{Bishop's notation 2.114}$$

$$(6): p(\bar{f}|\bar{X}) = \mathcal{N}(\bar{f} | \underbrace{\mu}_{\Lambda^{-1}}, K_M) \leftarrow \text{Bishop's notation 2.113}$$

Recall that all probability distributions so far have been implicitly conditioned on Θ

$$\text{Then } p(y|X, \bar{X}, \Theta) = \int_{\bar{f}} p(y, \bar{f} | X, \bar{X}, \Theta) d\bar{f}$$

$$= \int_{\bar{f}} p(y|X, \bar{X}, \bar{f}, \Theta) p(\bar{f}|\bar{X}, \Theta) d\bar{f}$$

{Using Bishop 2.115 p 93}

$$p(y|X, \bar{X}, \Theta) = \mathcal{N}(y | m, S) \text{ where}$$

Bishop's notation

$$m = A\mu + b = K_{NM} K_M^{-1} \mathbf{0} + \mathbf{0} = \mathbf{0}$$

$$S = L^{-1} + A\Lambda^{-1}A^T = \Lambda + G^2 I + K_{NM} K_M^{-1} K_M K_M^{-1} K_{MN} = \underline{K_{NM} K_M^{-1} K_{MN} + \Lambda + G^2 I}$$

$$\text{so } p(y|X, \bar{X}, \Theta) = \mathcal{N}(y | \mathbf{0}, K_{NM} K_M^{-1} K_{MN} + \Lambda + G^2 I) \text{ ie SPGP (9)}$$