# 2. Gaussian Processes.

## 2.1 Introduction

GPs are stochastic processes identified by a mean function $m$ and a kernel $k$, $f \sim GP(m, k)$.

$\forall x \in M$, $f(x)$ is a random variable (stoch. process).

For a GP, any set $\{f(x_1) \cdots f(x_M)\}$ is distributed according to an $M$-dimensional Gaussian, s.t.

$$\begin{cases} E[f(x_i)] = m(x_i) = m_i \\ Cov[f(x_i), f(x_j)] = k(x_i, x_j) = K_{ij} > 0 \text{ def.} \end{cases}$$

$m : M \to \mathbb{R}$

$k : M \times M \to \mathbb{R}$

## 2.2 Bayesian Approach to Inverse Problems.

$$y_z = \int_M dx \, C_z(x) f(x)$$

- $f$ is promoted to a GP.

- choose a prior for $f$, $p(f)$, where $f_i = f(x_i)$

   all prior knowledge is encoded in $p(f)$,

   $\hookrightarrow$ independent of the data.

- Bayes thm. $\Rightarrow$ posterior distribution

$$\bar{p}(f) = p(f|y) = \frac{p(y|f) \, p(f)}{p(y)}$$

   $p(y|f)$ likelihood

Knowledge about the solution is encoded in $\tilde{p}(f)$. e.g.

central value : $E_{\tilde{p}}[f]$

covariance : $Cov_{\tilde{p}}[f_i, f_j]$

## 2.3 Setting the problem

$x = \{x_i, i = 1 \ldots N\}$ $\qquad \rightarrow \qquad$ $f \in R^N$, $f_i = f(x_i)$

$x^* = \{x_i^*, i = 1 \ldots M\}$ $\qquad\qquad$ $f^* \in R^M$, $f_i^* = f(x_i^*)$

Joint prior, depends on a set of hyperparameters $\theta$.

$$p(f, f^* | \theta) = \frac{1}{\sqrt{\det(2\pi K)}} \exp\left\{-\frac{1}{2}\left((f-m)^T, (f^*-m^*)^T\right) K^{-1} \begin{pmatrix} f-m \\ f^*-m^* \end{pmatrix}\right\}$$

where $\quad m_i = m(x_i)$, $m_i^* = m(x_i^*)$

$\qquad\quad K_{ij} = k(x_i, x_j)$

$$K = \begin{pmatrix} K_{xx} & K_{xx^*} \\ K_{x^*x} & K_{x^*x^*} \end{pmatrix}, \quad (N+M) \times (N+M) \text{ sym., } >0 \text{ matrix.}$$

e.g. $m(x) = 0$, $\forall x \in M$

$$k(x,x') = \sigma^2 \sqrt{\frac{2\,l(x)\,l(x')}{l(x)^2 + l(x')^2}} \exp\left\{-\frac{(x-x')^2}{l(x)^2 + l(x')^2}\right\}.$$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \hookrightarrow$ Gibbs kernel

$l(x) = l_0 (x + \delta)$

$\sigma$ and $l_0$ are the hyperparameters.

Other choices are possible, depending on the info that we want to encode.

In all cases, the prior is explicit, which is good.

## 2.4 Data & Theory Predictions

data central values $y = \{y_I, I = 1 \ldots N_{dat}\}$

fluctuations $\varepsilon \sim \mathcal{N}(0, C_y)$

theory prediction for the $I^{th}$ datapoint

$$T_I = \int_M dx \; C_I(x) \, f(x) \simeq \sum_{i=1}^{M} (FK)_{Ii} \, f_i$$

only $f_i$ involved in the theory prediction, not $f_i^*$.

$f$ is a GP $\Rightarrow$ $T_I$ are Gaussian variables.

$$\begin{cases} E_p[T_I] = (FK)_{Ii} \; E_p[f_i] = (FK)_{Ii} \, u_i \\ \text{Cov}_p[T_I, T_J] = (FK)_{Ii} \, K_{ij} \, (FK)^T_{jJ} \end{cases}$$

## 2.5 Posterior Distribution

$$\tilde{p}(f, f^*) = p(f, f^* | y)$$

$$= \int d\theta \; p(f, f^*, \theta | y) \qquad \text{marginalize wrt. } \theta.$$

$$\tilde{p}(f) = \int df^* \; \tilde{p}(f, f^*) \qquad \text{and} \qquad \tilde{p}(f^*) = \int df \; \tilde{p}(f, f^*)$$

$\tilde{p}(f^*)$ posterior distribution f- values of $f$ that do $\underline{not}$ enter in the theory prediction.

We have:

$$p(f, f^*, \theta | y) = \underbrace{p(f, f^* | \theta, y)}_{(a)} \; \underbrace{p(\theta | y)}_{(b)}$$

The two factors, (a) and (b), can be computed separately.

(a) $p(f, f^* | \theta, y) \propto \exp\left\{ -\frac{1}{2} \left( (f-u)^T, (f^*-u^*)^T \right) K^{-1} \begin{pmatrix} f-u \\ f^*-u^* \end{pmatrix} \right\} \times$

$$\times \exp\left\{ -\frac{1}{2} \left( (FK)f - y \right)^T C_y^{-1} \left( (FK)f - y \right) \right\}.$$

Integrate out $f^*$

$$\int df^* \, p(f, f^* | \theta, y) \propto \left[ \int df^* \exp\left\{ -\frac{1}{2} \begin{pmatrix} f-m \\ f^*-m^* \end{pmatrix}^T \begin{pmatrix} k_{xx} & k_{xx^*} \\ k_{x^*x} & k_{x^*x^*} \end{pmatrix}^{-1} \begin{pmatrix} f-m \\ f^*-m^* \end{pmatrix} \right\} \right] \times$$

$$\times \exp\left\{ -\frac{1}{2} \left( (FK)f - y \right)^T C_y^{-1} \left( (FK)f - y \right) \right\} .$$

$$= \exp\left\{ -\frac{1}{2} (f-m)^T (k_{xx})^{-1} (f-m) \right\} \exp\left\{ -\frac{1}{2} \left( (FK)f - y \right)^T C_y^{-1} \left( (FK)f - y \right) \right\} .$$

$\hookrightarrow$ quadratic form in $f$.

For linear data, $\tilde{p}(f)$ is a Gaussian

$$\tilde{p}(f | \theta, y) = \mathcal{N}(f ; \tilde{m}, \tilde{K}_{xx}) .$$

where $\begin{cases} \tilde{m} = m + K_{xx} (FK)^T C_{yT}^{-1} \left[ y - (FK)m \right] \\[4pt] \tilde{K}_{xx} = K_{xx} - K_{xx} (FK)^T C_{yT}^{-1} (FK) K_{xx} \\[4pt] C_{yT} = (FK) K_{xx} (FK)^T + C_y \end{cases}$

**N.B.** : $(\tilde{K}_{xx})_{ii} = (\tilde{\Delta} f_i)^2 \leq (K_{xx})_{ii}^2 = (\Delta f_i)^2$

since $C_{yT}$ is $> 0$ def.

Posterior dist. yields reduced errors on $f$:

Integrate out $f$. slightly trickier.

$$\tilde{p}(f^* | \theta, y) = \mathcal{N}(f^* ; \tilde{m}^*, \tilde{K}_{x^*x^*}) .$$

$$\tilde{m}^* = m^* + K_{x^*x} (FK)^T C_{yT}^{-1} \left[ y - (FK)m \right] .$$

$\hookrightarrow$ updated central value of $E_{\tilde{p}}[f^*]$ because of the correlations introduced by the prior, $K_{x^*x}$.

Note that the data is independent of $f^*$.

(b) $p(\theta|y) \propto p(y|\theta)p(\theta)$.

$$p(y|\theta) = \frac{1}{\sqrt{\det(2\pi C_{yT})}} \exp\left\{ -\frac{1}{2}\left[y - (FK)m\right]^T C_{yT}^{-1}\left[y - (FK)m\right] \right\}.$$

$$y = (FK)f + \varepsilon \quad \left| \quad \begin{array}{l} (FK)f \sim \mathcal{N}\left((FK)m, (FK)K_{xx}(FK)^T\right) \\ \\ \varepsilon \sim \mathcal{N}(0, C_y) \end{array} \right.$$

Hyperparameters $\theta$ appear in $C_{yT}$, the posterior distr. can only be sampled by MCMC.

If the distribution is sufficiently narrow, $\theta$ can be fixed to the mode of the posterior distr.

## 2.6 Closure Test - 1

Consider synthetic data generated w. a given $f_0$.

$$y = (FK)f_0 + \eta, \quad \eta \sim \mathcal{N}(0, C_y).$$

Study the case of vanishing noise: $C_y = 0$

$\Rightarrow \quad \tilde{m} = R_{xx}^{(6)} f_0$

$$R_{xx}^{(6)} = K_{xx}(FK)^T \left[(FK)K_{xx}(FK)^T\right]^{-1}(FK).$$

$\hookrightarrow$ smearing kernel.

Note the correspondence w. BG solution.

$$\tilde{m} = a_I y_I \quad \left| \quad \begin{array}{l} a_I = K_{xx}(FK)^T \left[(FK)K_{xx}(FK)^T\right]^{-1} \\ \\ y_I = (FK)f_0 \end{array} \right.$$

cfr. w. BG solution

$$w_k(x_0) = \int dx \; c_k(x) \, k(x, v_o) \quad \longrightarrow \quad (FK)_{ki} (K_{xx})_{i i_o}$$

$$\hat{W}_{kJ} = \int dx \; c_k(x) \, k(x, x') \, c_J(x') \quad \longrightarrow \quad (FK) K_{xx} (FK)^T$$

$$a_I = \left(\hat{W}^{-1}\right)_{IJ} w_J = \left(w^T\right)_J \left(\hat{W}^{-1}\right)_{JI}$$

$\uparrow$ symmetric

Identical sol.n if $k(x, x')$ in the BG metric

$$= k(kx') \text{ for the GP.}$$

$\tilde{m} \neq f_0$ even in the absence of stat. fluctuations in the data, there is a "reconstruction" error in the space of functions $f$.

## 2.7 Bias & Variance

In data space

$$B = \sum_I \left(\tilde{T}_I - y_I\right) = \sum_I (FK)_{Ii} \left(\tilde{m}_i - f_{0i}\right)$$

$$= \sum (FK) \left[ R_{xx}^{(0)} - 1 \right] f_0 = 0$$

data is reproduced exactly !

$$V = tr\left[ (FK) \tilde{K} (FK)^T \right]$$

$$= tr\left[ (FK) \left(1 - R_{xx}^{(0)}\right) K_{xx} (FK)^T \right] = 0 \;!$$

$V = 0 \;\Rightarrow\; \text{exact reconstruction of data.}$

## 2.8 Closure Test - 2

Adding exp. errors :   $y_I = (FK)_{Ii} f_{0i} + \eta_I$

$$R_{xx} = K_{xx}(FK)^T \left[ (FK) K_{xx} (FK)^T + C_\eta \right]^{-1} (FK).$$

$$\tilde{m} = R_{xx} f_0 + a_{xx}^T \eta$$

$$\tilde{K}_{xx} = (1 - R_{xx}) K_{xx} (1 - R_{xx})^T + a_{xx}^T C_\eta a_{xx}$$

$$a_{xx}^T = K_{xx}(FK)^T C_{\eta T}^{-1}$$

$\Rightarrow R_{xx} = a_{xx}^T (FK).$

cfr.   $f_a(x) = \sum_I a_I C_I(x)$   for BG.      $\Big\}$  connection w. BG.

In data space, we can compute bias & variance

$$B = (FK)\left[ R_{xx} - 1 \right] f_0 + (FK) a_{xx}^T \eta$$

$$V = (FK)(1 - R_{xx}) K_{xx} (1 - R_{xx})^T (FK)^T + (FK) a_{xx}^T C_\eta a_{xx} (FK)^T$$

A bit of algebra yields, for $C_\eta \to 0$

$$B = - C_T^{-1} C_\eta \left( (FK) f_0 + \eta \right)$$

$$= - G^{-1} C_\eta y$$

$$V = (FK) \left[ (1 - R_{xx}) K_{xx} (1 - R_{xx})^T + a_{xx}^T C_\eta a_{xx} \right] (FK)^T$$

Explicit dependence on the prior, $K_{xx}$.

All assumptions are exposed !

No minimization needed. Only sample $p(\theta|y)$ by MCMC.