

Let us define the  $N$  ~~test~~<sup>training</sup> points as

$$\{(x_i^o, y_i^o)\}_{i=1,2,\dots,N}, \text{ where } y_i^o = f(x_i^o)$$

and the  $M$  test points as

$$\{(x_j^o, f(x_j^o))\}_{j=1,2,\dots,M}$$

The fundamental GP assumption is:

$$[y_f] \sim N\left(\begin{bmatrix} 0 \end{bmatrix}, \begin{bmatrix} K_{xx} & K_{xx^*} \\ K_{x^*x} & K_{x^*x^*} \end{bmatrix}\right)$$

where  $x$  represents the training data  
and  $x^*$  represents the test data.

To get a posterior predictive distribution,  
we need inference in jointly Gaussian distributions.

Lemma: Let  $M$  be a general partitioned matrix,

then  $M = \begin{bmatrix} E & F \\ G & H \end{bmatrix}$

$$M^{-1} = \begin{bmatrix} (M/H)^{-1} & - (M/H)^{-1} F H^{-1} \\ -H^{-1} G (M/H)^{-1} & H^{-1} + H^{-1} G (M/H)^{-1} F H^{-1} \end{bmatrix}$$

where

$$M/H = E - FH^{-1}G$$

Proof:  $\begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} E & F \\ G & H \end{bmatrix} = \begin{bmatrix} E - FH^{-1}G & 0 \\ G & H \end{bmatrix} - 0$

and

$$\begin{bmatrix} E - FH^{-1}G & 0 \\ G & H \end{bmatrix} \begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix} = \begin{bmatrix} E - FH^{-1}G & 0 \\ 0 & H \end{bmatrix} - ②$$

Considering ① and ②

$$\underbrace{\begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix}}_X \underbrace{\begin{bmatrix} E & F \\ G & H \end{bmatrix}}_M \underbrace{\begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix}}_Z = \begin{bmatrix} E - FH^{-1}G & 0 \\ 0 & H \end{bmatrix}$$

Taking inverse on both sides,

$$Z^{-1} M^{-1} X^{-1} = W^{-1}$$

hence

$$M^{-1} = ZW^{-1}X$$

$$\begin{aligned}
 \Rightarrow \begin{bmatrix} E & F \\ G & H \end{bmatrix}^{-1} &= \begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix} \begin{bmatrix} (M/H)^{-1} & 0 \\ 0 & H^{-1} \end{bmatrix} \begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix} \\
 &= \begin{bmatrix} (M/H)^{-1} & 0 \\ -H^{-1}G(M/H)^{-1}H^{-1} & I \end{bmatrix} \begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix} \\
 &= \begin{bmatrix} (M/H)^{-1} & -(M/H)^{-1}FH^{-1} \\ -H^{-1}G(M/H)^{-1} & H^{-1} + H^{-1}G(M/H)^{-1}FH^{-1} \end{bmatrix}
 \end{aligned}$$

Now, writing the standard bivariate distribution using matrix approach,

$$P(f, y) = e^{-\frac{1}{2} \begin{bmatrix} f - \mu_1 \\ y - \mu_2 \end{bmatrix} \begin{bmatrix} K_{X^*, X^*} & K_{X^*, X} \\ K_{X, X^*} & K_{X, X} \end{bmatrix}^{-1} \begin{bmatrix} f - \mu_1 \\ y - \mu_2 \end{bmatrix}}$$

Now, using the previous equation, we get

$$\begin{aligned}
 P(f, y) &= e^{-\frac{1}{2} \begin{bmatrix} f - \mu_1 \\ y - \mu_2 \end{bmatrix} \begin{bmatrix} I & 0 \\ -k_{xx}^{-1} & I \end{bmatrix} \begin{bmatrix} (K/K_{xx})^{-1} & 0 \\ 0 & K_{xx}^{-1} \end{bmatrix} \begin{bmatrix} f - \mu_1 \\ y - \mu_2 \end{bmatrix}} \\
 &\quad \cdot \begin{bmatrix} I & -K_{x^*, X} K_{xx}^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} f - \mu_1 \\ y - \mu_2 \end{bmatrix}
 \end{aligned}$$

$$\text{where } K/K_{xx} = K_{x^*, X^*} - (K_{x^*, X})(K_{xx})^{-1}(K_{xx}^*)$$

$$\begin{aligned}
 &= e^{\left\{ -\frac{1}{2} (f - \mu_1 - (K_{x^*, X})(K_{xx})^{-1}(y - \mu_2))^T (K/K_{xx})^{-1} (f - \mu_1 - (K_{x^*, X})(K_{xx})^{-1}(y - \mu_2)) \right\}} \\
 &\quad \cdot e^{\left\{ -\frac{1}{2} (y - \mu_2)^T K_{xx}^{-1} (y - \mu_2) \right\}}
 \end{aligned}$$

hence

We have written  
 $p(f, y) = (\cdot) \cdot p(y)$

$$\Rightarrow (\cdot) = p(\frac{f}{y}) \quad [\text{posterior distribution}]$$

hence, we can write  
 $u_1 u_2$  as  $u_1 + (K_{x^*x})(K_{xx})^{-1}(y - u_2)$

we consider  $y$  to be distributed around  $u_2$

hence

$$u_1 u_2 = u_1 + (K_{x^*x})(K_{xx})^{-1}(y)$$

if we take  $u_1$  as zero,  
 we get

$$u_f = (K_{x^*x})(K_{xx})^{-1}(y)$$

Similarly, we get

$$K_f = K / K_{xx} = K_{x^*x} - K_{x^*x} (K_{xx})^{-1} K_{x^*x}$$

hence,  $f | x^*, D \sim N(u_f, K_f)$

$\downarrow$   
prior data

→ The kernel function needs to be selected on the basis of trends in the data.

Eg: If the data shows a periodic trend, then we can add a periodic kernel. The effect of the kernels can be stacked.

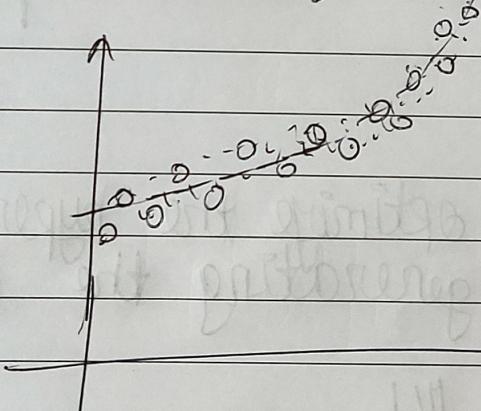
Kernels can be added :

$$K_C(x, x' | \mathcal{G}_C) = K_A(x, x' | \mathcal{G}_A) + K_B(x, x' | \mathcal{G}_B)$$

and kernels multiplication:

$$K_C(x, x' | \mathcal{G}_C) \approx K_A(x, x' | \mathcal{G}_A) \cdot K_B(x, x' | \mathcal{G}_B)$$

Hence consider the following data distribution in 2-D.



The distribution can be thought of as ~~a~~ periodic with a quadratic trend.

Hence, model kernel =  $(v_1 x \cdot x^2) \cdot (v_2 x \cdot x^2) + e^{(\frac{x^2}{\sigma^2}) \sin^2(\frac{\pi}{P}(x - x'))}$

(linear)(linear)

periodic kernel

→ The utility of GP is the fact that it does not simply give us a prediction, but rather a predictive distribution. This predictive distribution can be used to assess the confidence of the prediction. If the confidence of the model goes up as new samples are encountered, we can be sure that it is extrapolating well.

~~we can~~

→ Yes, we can adjust this system to perform classification. For example, we can use the regressor output and squash it using the sigmoid function to limit its output b/w 0 and 1 (for binary classification)

→ We ~~will~~ need to optimize the hyperparameters of the kernel generating the prior distribution.

We can use MLL

$\log(p(y|X, \gamma))$ , where  $X$  is the training data and  $\gamma$  represents the hyperparameters of the kernel.

$$\text{Now } \log(p(y|X, \gamma)) = \log(N(y|0, K_{xx}))$$

Since the prior samples are sampled from mean 0 and covariance matrix  $K_{xx}$ .

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:						

Now,  $\mathcal{J}$  is differentiable wrt RHS hence it can be optimised.