

HGD PROPRIETIES

7/14/2023

$$x \in \mathbb{R}^d, \mu \in \mathbb{R}^d, \Sigma \in \mathbb{S}^d \Rightarrow x \sim N(\mu, \Sigma)$$

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

WHERE OUR PARAMETERS ARE DEFINED AS FOLLOWS (EXAMPLE IN 2 DIMENSIONS)

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}$$

1. NORMALIZATION

$$\int_x P(x; \mu, \Sigma) = 1$$

2. MARGINALIZATION

$$P(x_A) = \int_{x_B} P(x_A, x_B; \mu, \Sigma) dx_B$$

$$x_A \sim N(\mu_A, \Sigma_{AA})$$

WE OBTAIN A MARGINALIZED DISTRIBUTION (WHICH IS ALSO A NGD)

ex

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \dots \\ & \ddots \\ & & \Sigma_{dd} \end{bmatrix} \xrightarrow{\text{MARGINALIZED}} \begin{cases} \bullet x \sim N(\mu, \Sigma) \\ \bullet x_1 \sim N(\mu_1, \Sigma_{11}) \\ \bullet \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \dots \\ & \ddots \\ & & \Sigma_{33} \end{bmatrix}\right) \end{cases}$$

3. CONDITIONING

$$P(x_A | x_B) = \frac{P(x_A, x_B; \mu, \Sigma)}{\int_{x_A} P(x_A, x_B; \mu, \Sigma)} = \frac{P(x_A, x_B; \mu, \Sigma)}{P(x_B; \mu, \Sigma)}$$

WE ARE NOW ABLE TO DEFINE A NEW NGD

$$x_A | x_B \sim N\left(\underbrace{\mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (x_B - \mu_B)}_{\frac{\sqrt{\sigma_1 \sigma_2}}{\sigma_1^2}}, \underbrace{\Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}}_{\frac{\sqrt{\sigma_1 \sigma_2}}{\sigma_1^2} \cdot \sqrt{\sigma_2 \sigma_2} = \sqrt{\sigma_1^2 \sigma_2^2}}\right) \quad (*)$$

ex

LET'S CONSIDER TWO DIMENSIONS

$$X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \Sigma = \begin{pmatrix} \sum_{11} & \sum_{12} \\ \sum_{21} & \sum_{22} \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \overset{\sum_{11}}{\sigma_1^2} & \overset{\sum_{12}}{\rho \sigma_1 \sigma_2} \\ \overset{\sum_{21}}{\rho \sigma_1 \sigma_2} & \overset{\sum_{22}}{\sigma_2^2} \end{pmatrix}$$

WHERE σ IS THE STD. DEVIATION

σ^2 THE VARIANCE

ρ THE CORRELATION TERM

THEN, FROM (*) WE GET

$$x_1 | x_2 \sim N \left(\mu_1 + \rho \overset{\sigma_1}{\sigma_2} \frac{(x_2 - \mu_2)}{\sigma_2}, \overset{\sigma_1^2}{\sigma_1^2} - \rho^2 \overset{\sigma_1^2}{\sigma_1^2} \right)$$

4. SUMMATION

SUPPOSE TWO INDEPENDENT RANDOM VARIABLES

$$\begin{cases} y \sim N(\mu, \Sigma) \\ y' \sim N(\mu', \Sigma') \end{cases} \longrightarrow y + y' \sim N(\mu + \mu', \Sigma + \Sigma')$$

BAYESIAN LINEAR REGRESSION

$$S = \{x^{(i)}, y^{(i)}\}_{i=1}^n, \quad y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)} \quad i = 0, \dots, n \quad \text{AND} \quad \epsilon^{(i)} \sim N(0, \sigma^2)$$

LET'S NOW ADD THE BAYESIAN ASSUMPTION

$$\theta \sim N(0, \gamma^2 I)$$

STEPS OF BAYESIAN METHOD

PRIOR DISTRIB.		POSTERIOR	
$P(\theta)$	$\xrightarrow{\text{OBSERVE } S}$	$P(\theta S) \stackrel{=}{=} \frac{P(\theta) \prod_{i=1}^n \overbrace{P(y^{(i)} x^{(i)}, \theta)}^{\text{NORMAL DISTR.}}}{\int_{\theta} P(\theta) \prod_{i=1}^n P(y^{(i)} x^{(i)}, \theta) d\theta}$	

AND WE KNOW THAT

$$\theta|S \sim N\left(\frac{1}{\sigma^2} A^{-1} X^T \vec{y}, A^{-1}\right) \quad \text{WHERE} \quad A = \frac{1}{\sigma^2} X^T X + \frac{1}{\gamma^2} I$$

$$\theta|S \sim N\left(\frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} X^T X + \frac{1}{\gamma^2} I\right)^{-1} X^T \vec{y}, A^{-1}\right) = \boxed{N\left((X^T X + \frac{\sigma^2}{\gamma^2} I)^{-1} X^T \vec{y}, A^{-1}\right) \sim \theta|S}$$

GIVEN THE POSTERIOR PREDICTIVE DISTRIBUTION, WE CAN ESTIMATE $\hat{\Theta}$ AS FOLLOWS

$$\hat{\Theta} = (X^T X)^{-1} X^T \vec{y}$$

AND DEFINE ITS DISTRIBUTION AS

$$\Theta \sim N \left((X^T X + \frac{\sigma^2}{\beta_2} I)^{-1} X^T \vec{y}, (X^T X + \frac{\sigma^2}{\beta_2} I)^{-1} \right)$$

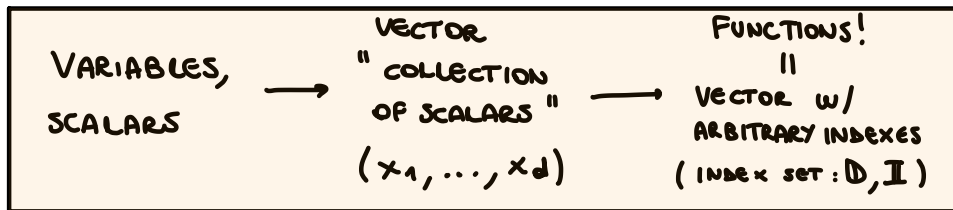
AND THE PREDICTION \hat{y} WILL BE AS FOLLOWS

$$\hat{y} = \hat{\Theta}^T x_*$$

PRED EST. Θ NEW EX

GAUSSIAN PROCESS

LET'S MAKE THIS ASSUMPTION



IF WE APPLY TO A MNGD A FUNCTION INSTEAD OF VECTORS, WE GET A GAUSSIAN PROCESS

$$\begin{array}{c} \mathbb{R}^d \\ \cup \\ x \sim N(\mu, \Sigma) \end{array} \quad \text{(vs)} \quad \begin{array}{c} \mathbb{R}^d \\ \cup \\ \vec{f} \sim N(\mu, \Sigma) \end{array} \quad \begin{array}{l} f_0 = f(0) \sim N(\mu_0, \Sigma_{00}) \\ \vdots \\ f_d = f(d) \sim N(\mu_d, \Sigma_{dd}) \end{array}$$

WE HAVE A TRAINING SET $\{ \underbrace{x^{(i)}}_{\mathbb{R}^d}, y^{(i)} \}_{i=1}^n$, NOT DEFINING A MGD OVER $x^{(i)}$ 'S. $x^{(i)}$ 'S WILL BE THE INDEXES

$x^{(i)} \in X$
ABSTRACT SPACE

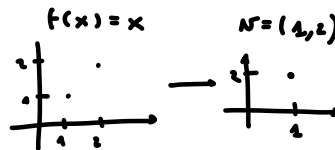
$$\begin{bmatrix} f_{x^{(1)}} \\ f_{x^{(2)}} \\ f_{x^{(3)}} \\ \vdots \\ f_{x^{(n)}} \end{bmatrix} = \begin{bmatrix} f_{x^{(1)}} \\ f_{x^{(2)}} \\ f_{x^{(3)}} \\ \vdots \\ f_{x^{(n)}} \end{bmatrix}$$

THIS ENTITY INDEXED w/ ANYTHING $\in X$ (∞ DIMENSIONAL)

$$f_X(\cdot) \sim GP(m(\cdot), K(\cdot, \cdot)) \quad \mathbb{D}: x^{(1)}, \dots, x^{(n)}; X \mapsto \mathbb{R}^d$$

$$z \sim N(\mu, \Sigma) \in \mathbb{R}^d \quad \mathbb{D}: d, \dots, d$$

INDEX: d, \dots, d (0 DIMENSIONAL)



HOW LIKELY A DETERMINISTIC FUNCTION IS

GP \rightarrow DISTRIBUTION OVER FUNCTIONS, A SAMPLE IS A FUNCTION

N \rightarrow DISTRIBUTION OVER VECTORS, A SAMPLE IS A VECTOR

$$z : [1, d] \mapsto \mathbb{R}, \mu : [1, d] \mapsto \mathbb{R}, \Sigma : [1, d] \times [1, d] \mapsto \mathbb{R}$$

PSD

$$f : x \mapsto \mathbb{R}, \mu : x \mapsto \mathbb{R}, K : x \times x \mapsto \mathbb{R}$$

KERNEL (COVARIANCE FUNCTIONS)

$$z^T M z, \dots, z, 0$$

\downarrow

$$\int \int_z f(z) K(z, x) f(x) dx dz$$

VS

$$a^T b$$

$$\sum_i a_i b_i$$

$$\int_z a(z) b(z) dz$$

ex

$$\begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_8 \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_0 \\ \mu_1 \\ \mu_2 \\ \vdots \\ \mu_8 \end{bmatrix}, \begin{bmatrix} \quad \end{bmatrix} \right)$$

WE ONLY CARE ABOUT THE TRAIN AND TEST SET, SO WE MARGINALIZE THE DISTRIBUTION AND OBTAIN A NEW ONE. WE CAN DEFINE THE NON-PARAMETRIC PREDICTION (LINEAR REGRESSION) AS FOLLOWS

$$y^{(i)} = f(x^{(i)}) + \epsilon^{(i)}, \quad i = 1, \dots, n$$

↓
WITH FUNCTION,
NOT DIRECTLY
THE PARAMETER

WITH PARAMETERS BEING DEFINED AS FOLLOWS

TRAIN SET

$$\underline{X} = \begin{bmatrix} - & x^{(1)} & - \\ & \vdots & \\ - & x^{(n)} & - \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \vec{f} = \begin{bmatrix} f(x^{(1)}) \\ \vdots \\ f(x^{(n)}) \end{bmatrix} \quad \vec{\epsilon} = \begin{bmatrix} \epsilon^{(1)} \\ \vdots \\ \epsilon^{(n)} \end{bmatrix} \quad \underline{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \in \mathbb{R}^n$$

X, y KNOWN

TEST
SET

$$\underline{X}_* = \begin{bmatrix} -x_*^{(1)} - \\ \vdots \\ -x_*^{(n)} - \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \vec{f} = \begin{bmatrix} f(x_*^{(1)}) \\ \vdots \\ f(x_*^{(n)}) \end{bmatrix} \quad \vec{E}_* = \begin{bmatrix} \epsilon_*^{(1)} \\ \vdots \\ \epsilon_*^{(n)} \end{bmatrix} \quad \underline{\vec{y}}_* = \begin{bmatrix} y_*^{(1)} \\ \vdots \\ y_*^{(n)} \end{bmatrix} \in \mathbb{R}$$

X KNOWN
y TO PREDICT

WE CAN OBSERVE THAT

$$f(x_*) = \mathbb{E}[y|x] \quad (\text{MEAN } 0)$$

↓
DETERMINISTIC

LET'S ASSUME THAT $f \sim N(0, K(\cdot, \cdot))$

WE NEED THE FULL DISTRIBUTION $y|x$. BY MARGINALIZING A GAUSSIAN PROCESS WE OBTAIN A M-NGD, AS FOLLOWS

$$\begin{bmatrix} \vec{f} \\ \vec{f}_* \end{bmatrix} | x, x_* \sim N \left(\vec{0}, \begin{bmatrix} K(x, x) & K(x, x_*) \\ K(x_*, x) & K(x_*, x_*) \end{bmatrix} \right)^{(*)}$$

NOW WE NEED THE ESTIMATES OF THE NOISE

$$(*) + \mathcal{E} = \begin{bmatrix} \vec{f} \\ \vec{f}_* \end{bmatrix}$$

NOW WE CAN FINALLY COME TO THE JOINT DISTRIBUTION \vec{y}, \vec{y}_* AND FINALLY OBTAIN A POSTERIOR PREDICTIVE DISTRIBUTION (NON PARAMETRIC)

$$\begin{bmatrix} \vec{y} \\ \vec{y}_* \end{bmatrix} | x, x_* = \begin{bmatrix} \vec{f} \\ \vec{f}_* \end{bmatrix} + \begin{bmatrix} \mathcal{E} \\ \mathcal{E}_* \end{bmatrix} \stackrel{\text{SUMMATION INDEPENDENCY}}{=} N \left(\vec{0} + \vec{0}, \begin{bmatrix} K(x, x) + \sigma^2 \mathbf{I} & K(x, x_*) \\ K(x_*, x) & K(x_*, x_*) + \sigma^2 \mathbf{I} \end{bmatrix} \right)$$

LET'S NOW APPLY THE CONDITIONING PROPERTY, AND WE OBTAIN THE POSTERIOR PREDICTIVE DISTRIBUTION

$$\boxed{\vec{y}_* | \vec{y}, x, x_* \sim N(\mu_*, \Sigma_*)} \xrightarrow{\text{PARAMS}} \begin{cases} \mu_* = K(x, x_*) (K(x, x) + \sigma^2 I)^{-1} \vec{y} \\ \Sigma_* = K(x_*, x) + \sigma^2 I - K(x_*, x) (K(x, x) + \sigma^2 I)^{-1} K(x, x_*) \end{cases}$$
$$= X_*^T X (X^T X + \sigma^2 I)^{-1} \vec{y}$$

LET'S RECALL SOME CONCEPTS ABOUT POSTERIOR DISTRIBUTIONS

→ IF OVER \vec{y} : PREDICTIVE

→ IF OVER OTHER PARAMS SUCH AS θ : NON-PREDICTIVE

REMEMBER THAT WE HAVEN'T MADE ANY ASSUMPTION ON \vec{f} AT ALL, THAT COULD BE "VERY NON-LINEAR"

GP: ARE PRETTY EXPENSIVE BUT VERY GOOD IN PREDICTING ($y \in \mathbb{R}$)

WE CAN HAVE NO PARAMETRIC CONSTRAINT, BUT STILL NEED TO CHOOSE A KERNEL FUNCTION (EVERYTHING DERIVES FROM THAT CONSEQUENTIALLY)