

Today we have a blackboard lecture :)

Gaussian process regression

Model $\left[\begin{array}{c} y = \underset{\substack{\uparrow \\ \text{input}}}{\vec{x}}^T \underset{\substack{\uparrow \\ \text{parameters}}}{\vec{\theta}} + \underset{\substack{\uparrow \\ \text{noise}}}{\varepsilon}, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2) - \text{i.i.d.} \end{array} \right.$

Prior $\left[\begin{array}{c} \vec{\theta} \sim \mathcal{N}(\vec{0}, \beta \mathbf{I}) - \text{prior distribution} \end{array} \right.$

Bayes formula $\left[\begin{array}{c} p(\vec{\theta} | \mathcal{D}) = \frac{p(\mathcal{D} | \vec{\theta}) p(\vec{\theta})}{p(\mathcal{D})} - \text{Bayes formula} \end{array} \right.$

$\underline{p(\mathcal{D} | \vec{\theta})} = \prod_{i=1}^n p(y_i | \vec{x}_i, \vec{\theta})$
 $\propto p(\mathcal{D} | \vec{\theta}) p(\vec{\theta}) \rightarrow \max_{\vec{\theta}} p(\mathcal{D} | \vec{\theta}) = \mathcal{L}(\vec{\theta})$
 $\rightarrow \max_{\vec{\theta}} \mathcal{L}(\vec{\theta})$
 Maximum a posterior estimate:

$$\vec{\theta}_{\text{MAP}} = (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\beta} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$p(y_i | -) \sim \mathcal{N}(\vec{x}_i^T \vec{\theta}, \sigma^2)$

Posterior distribution $p(\vec{\theta} | \mathcal{D}) = \mathcal{N}(\vec{\theta}_{\text{MAP}}, \sigma^2 (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\beta} \mathbf{I})^{-1})$

dist. function

$$P(\theta) \rightarrow P(\theta | D)$$

$$f_{\theta}(\vec{x}) = \vec{\theta}^T \vec{x}$$

$$\vec{\theta}_{MAP}$$

GP_y

$$\vec{\theta} \sim N(0, \sigma^2)$$

$$\vec{\theta}_1, \vec{\theta}_2, \dots$$

$$f_{\theta_2} = \vec{\theta}_2^T \vec{x}$$

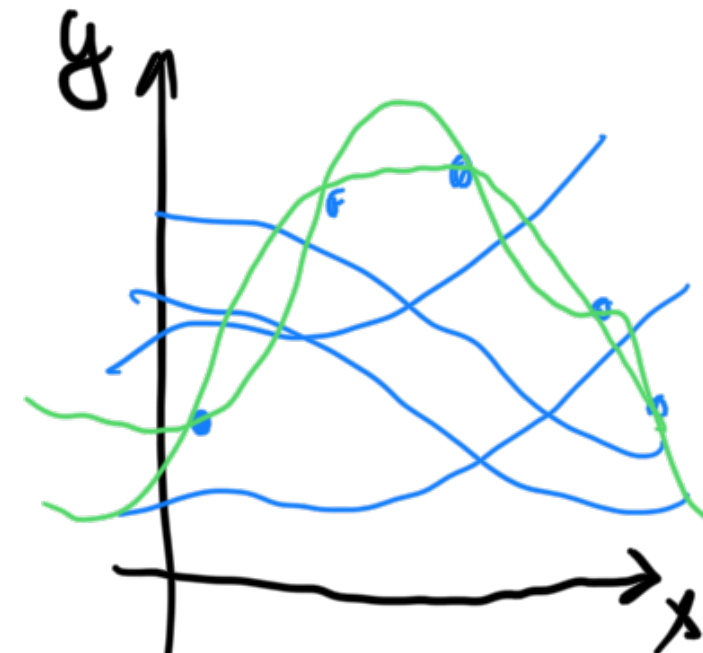
$$\vec{\theta} \sim p(\vec{\theta} | D)$$

$$\vec{\theta}_1, \vec{\theta}_2, \dots$$

$$f_{\theta_2'}(\vec{x}) = \vec{\theta}_2'^T \vec{x}$$

Gaussian random process

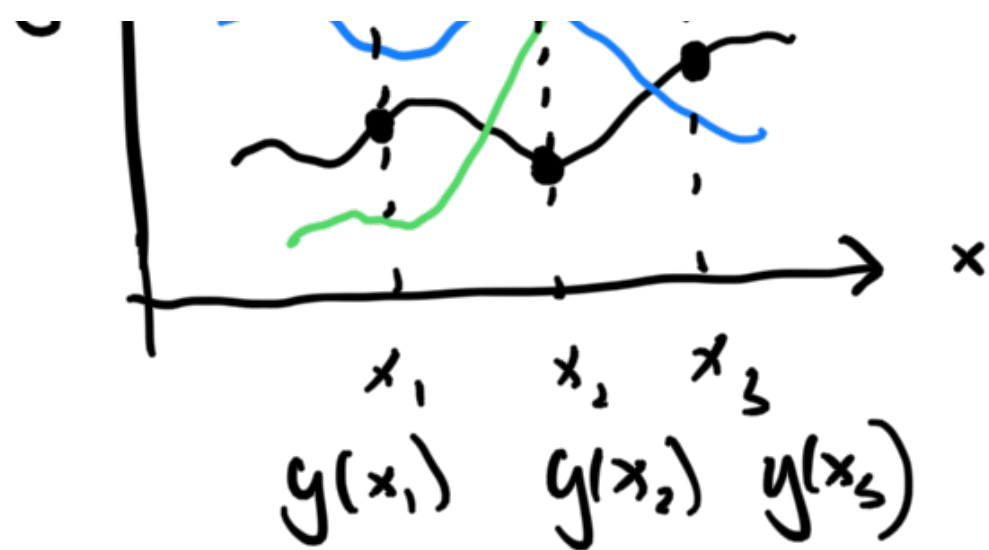
Bayesian linear regression



$$GP(0, k(x, x'))$$

$$GP(\mu, k(x, x'))$$

$$y(x_1)$$



$$\begin{pmatrix} y(x_2) \\ y(x_3) \end{pmatrix} \sim \mathcal{N}(\underline{\bar{\mu}}, \underline{\Sigma})$$

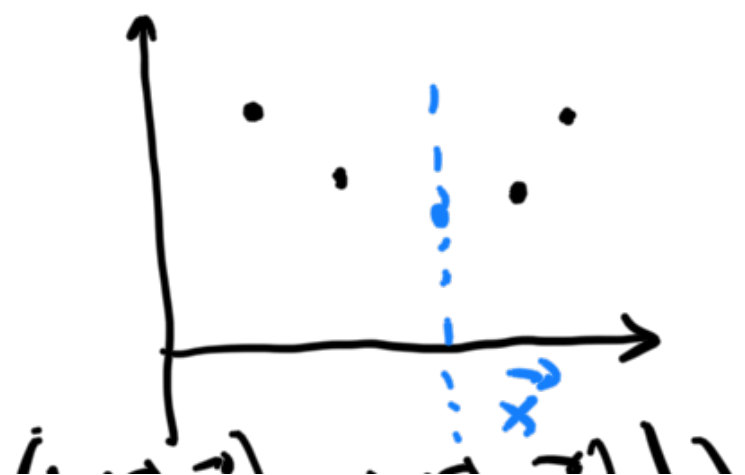
Def. Rand. process is Gaussian, if and only if
for all possible $\forall k \{x_1, \dots, x_k\}$:
 $(y(x_1), \dots, y(x_k))$ - Gaussian random vector

1. $\mathcal{N}(\underbrace{\mu(\vec{x})}_{\text{mean}}, \underbrace{k(\vec{x}, \vec{x}')}_{\text{covariance}})$

$$\begin{matrix} y(\vec{x}) & , & y(\vec{x}') \\ \mathbb{E}(y(\vec{x}) - \mathbb{E}y(\vec{x})) & & \\ (y(\vec{x}') - \mathbb{E}y(\vec{x}')) & & \end{matrix}$$

2. $\mathcal{D} = \{(\vec{x}_i, y(\vec{x}_i))\}_{i=1}^n$

3. $p(\underline{y|\vec{x}} | \mathcal{D})$



$$\begin{pmatrix} y(\vec{x}) \\ y(\vec{x}_1) \\ \vdots \\ y(\vec{x}_n) \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu(\vec{x}) \\ \vdots \\ \mu(\vec{x}_n) \end{pmatrix}, \begin{pmatrix} k(\vec{x}, \vec{x}) & k(\vec{x}, \vec{x}_1) \\ \vdots & \ddots \\ k(\vec{x}_1, \vec{x}_n) & k(\vec{x}_1, \vec{x}_1) \end{pmatrix} \right)$$

$$P(y(\vec{x}) | y(\vec{x}_1), \dots, y(\vec{x}_n)) = \mathcal{N}(\hat{y}(\vec{x}), \sigma^2(\vec{x}))$$

$$\hat{y}(\vec{x}) = \vec{k}^T K^{-1} \vec{y}_D = f(\vec{x}) = \sum_{i=1}^n \alpha_i k(\vec{x}, \vec{x}_i)$$

$O(n^2)$
 $O(n)$

$$\vec{y}_D = \begin{pmatrix} y(\vec{x}_1) \\ \vdots \\ y(\vec{x}_n) \end{pmatrix}$$

$$\vec{k} = \begin{pmatrix} k(\vec{x}, \vec{x}_1) \\ \vdots \\ k(\vec{x}, \vec{x}_n) \end{pmatrix}$$

$$\vec{\alpha} = K^{-1} \vec{y}_D$$

$$K = \begin{pmatrix} k(\vec{x}_1, \vec{x}_1) & \dots & k(\vec{x}_1, \vec{x}_n) \\ \vdots & & \vdots \end{pmatrix}$$



$$O(n^2)$$

$$O(n^3)$$

$$\mathbb{E} (y - \hat{y}(x))^2$$

$$k \dots k(x_n, x_n)$$

$$O(1)$$

$$\hat{g}^2(\bar{x}) = k(\bar{x}, \bar{x}) - \underline{\underline{k^T K^{-1} k}}$$

$$\hat{y}(x) = \sum_{i=1}^n \underline{w_i} y_i, \quad y_i = y(x_i) \quad (*)$$

$$\underline{\vec{w}} = \underline{k(\bar{x})} \underline{K^{-1}}$$

$$\hat{y}_{\text{ker}}(x) = \frac{1}{k} \sum_{i=1}^n \underline{k(x, x_i)} \underline{(y_i)} - \text{kernel regression}$$

$$k(x, x_i) = \exp(-\|\underline{x} - \underline{x_i}\|)$$

$$k = \sum_{i=1}^n k(x, x_i)$$

$$A1: \mu(\bar{x}) = 0$$

$$A2: k(\bar{x}, \bar{x}') = \exp\left(-\sum_{i=1}^d \frac{(x_i - x'_i)^2}{\Theta_i^2}\right)$$

$$\Theta = (\Theta_1, \dots, \Theta_d)$$

- how do we estimate $\vec{\Theta}$?

Maximum likelihood estimate!

MLE

$$\underline{\vec{y}} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \sim \mathcal{N}(\vec{0}, K_{\Theta}), \quad K_{\Theta} = \begin{pmatrix} k_{\Theta}(x_1, x_1) & \dots & k_{\Theta}(x_1, x_n) \\ \vdots & & \vdots \\ k_{\Theta}(x_n, x_1) & \dots & k_{\Theta}(x_n, x_n) \end{pmatrix}$$

$$L(\Theta) = \log p(\vec{y} | \Theta) \propto \underbrace{-\frac{1}{2} \log |K_{\Theta}|}_{\text{"regularization"}} - \underbrace{\frac{1}{2} \vec{y}^T K_{\Theta}^{-1} \vec{y}}_{\text{"goodness of fit"}}$$

$$L(\Theta) \rightarrow \max_{\Theta}$$

$$\frac{\partial L(\Theta)}{\partial \Theta} = -\frac{1}{2} \left(K_{\Theta}^{-1} \frac{\partial K_{\Theta}}{\partial \Theta} \right) + \frac{1}{2} \vec{y}^T K_{\Theta}^{-1} \frac{\partial K_{\Theta}}{\partial \Theta} K_{\Theta}^{-1} \vec{y}$$

$\left(\frac{\partial k_{\Theta}(x_1, x_1)}{\partial \Theta} \quad \frac{\partial k_{\Theta}(x_1, x_n)}{\partial \Theta} \right)$

k_θ - size
 $n \times n$
 \downarrow
 $\mathcal{O}(n^2)$

$$\begin{pmatrix} \frac{\partial \mathcal{L}}{\partial \theta} & \dots & \frac{\partial \mathcal{L}}{\partial \theta} \\ & \ddots & \\ & & \frac{\partial \mathcal{L}}{\partial \theta} \end{pmatrix}$$

$$\vec{\theta}' = \vec{\theta} - \alpha \frac{\partial \mathcal{L}(\vec{\theta})}{\partial \vec{\theta}} \quad \underline{\mathcal{O}(n^3)}$$

$$n \ll 10000$$

$$n: 100, 50, d: < 50$$

Interpolation

$$\begin{aligned} \hat{y}_i &= k_i^T K^{-1} \vec{y} \\ \vdots & \\ \hat{y}_n &= \cancel{K}^{-1} \vec{y} \end{aligned}$$

$$\varepsilon_1, \dots, \varepsilon_n \sim \mathcal{N}(0, \sigma^2) \text{ - i.i.d.}$$

$$y := \underbrace{v(x_1)}_{\text{...}}, \quad \rightarrow$$

$$k_y = \begin{pmatrix} k(x_1, x_1) + \delta^2 & k(x_1, x_n) \\ \vdots & \vdots \\ k(x_n, x_n) + \delta^2 \end{pmatrix}$$

$$\text{cond}(k_f) = \frac{\lambda_1}{\lambda_n}$$

min \uparrow λ_1 \downarrow max λ_n

$$\text{cond}(k_y) = \frac{\lambda_1 + \delta^2}{\lambda_n + \delta^2}$$

$$\frac{\lambda_1}{\lambda_n} < \frac{\lambda_1 + \delta^2}{\lambda_n + \delta^2}$$

~~$$\lambda_1 \lambda_n + \lambda_1 \delta^2 \leq \lambda_1 \lambda_n + \lambda_n \delta^2$$~~

$\lambda_1 \leq \lambda_n$

Cov. functions

$$k(\underline{x}, \underline{x}') = \exp\left(-\sum_{i=1}^d \frac{(x_i - x'_i)^2}{\Theta_i^2}\right) \leftarrow \text{squared covariance} \quad \leftarrow \text{if } \rightarrow \infty$$

$$= \exp \left(- \frac{\sum_{i=1}^n |x_i - x'_i|}{\Theta_i} \right) \leftarrow \nu = 1 \quad \text{function}$$

Matern covariance

functions

$$k(x, x') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} z}{e} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu} z}{e} \right),$$

$$z = \|x - x'\|$$

$(\nu - 1)$ - differentiable

$$\nu = \frac{3}{2}, \quad \nu = \frac{5}{2}$$

Bessel functional
of ν kind

$$1. \quad k(x, x) = 1$$

$$2. \quad k(x, x') \geq 0$$

$$3. \quad k(x, x') \text{ - positive definite}$$



~~Kriging~~, kriging
 $\sqrt{s^2/2}$

← Geoscience



$$(\vec{x}_1, \dots, \vec{x}_n) \rightarrow k \geq 0$$

$$\vec{S}^T K \vec{S} \geq 0 \quad - \text{non-negative definite}$$

Ex.

$$k(\vec{x}, \vec{x}') = \vec{x}^T \vec{x}'$$

\Rightarrow Regularized linear regression

Rasmussen "Gaussian processes in Machine Learning", 2006

Comp. comp. $O(n^3)$ - prohibitive if $n > 10\,000$

K \rightsquigarrow K
Nystrom approximation

$$X_n = (x_1 \dots x_n) \sim X_m = (x_1 \dots x_m) \quad m \ll n$$

$$K_n \approx K_{nm} K_m^{-1} K_{mn} \quad \begin{matrix} n \times n & m \times m & n \times m \end{matrix}$$

$\underbrace{K_{nm}}_{O(m, n)} \underbrace{K_m^{-1}}_{O(m^2)} \underbrace{K_{mn}}_{O(m^2, n)}$

$$K_n, K_m - \text{c.m. } X_n \text{ and } X_m$$

$$K_{nm} = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_m) \\ \vdots & & \vdots \end{pmatrix} \quad \text{size } n \times m$$

$$O(m^2 n)$$

$$2) \quad k(x, x') \approx \sum_{k=1}^S \underbrace{\varphi_k(x) \varphi_k(x')}_{\text{scalar product}}$$

$$k(\|x - x'\|) = \int s(z) \exp(-iz\|x - x'\|) dz$$

$$\approx \sum$$

$$\tilde{w} = (X'X)^{-1} X' \tilde{y}$$

$\tilde{w} \in \mathbb{R}^d$ $\Theta \in \mathbb{R}^{d \times d}$ $X \in \mathbb{R}^{d \times n}$ $\tilde{y} \in \mathbb{R}^n$
 $\mathcal{O}(d^2 n)$ $\mathcal{O}(s^2 n)$

$n = 10 \cdot 10^6$

Dimension

$$d \leq 50$$

1. Projection : $\tilde{x} \mapsto \tilde{z} = W_P \tilde{x}$

2. Representation : $\tilde{x} \mapsto e(\tilde{x}) \mapsto \tilde{z}$

Bayesian optimization

$$f(\tilde{x}) \rightarrow \max_{\tilde{x}}$$

1. smooth (m.b. with noise)

- expensive: only n points

2. $x \in \mathbb{R}^d$, $d \leq 50$

$$\mathcal{D} = \begin{pmatrix} (\vec{x}_1, \dots, \vec{x}_m) \\ (y_1, \dots, y_m) \end{pmatrix} \quad \mathcal{D} \rightarrow \hat{f}(\vec{x}) \approx f(\vec{x})$$



$$x_k^* = \arg \max_{x \in \mathcal{X}} \left(\hat{f}_k(\vec{x}) + \beta \sigma_{f_k}^2(\vec{x}) \right)$$

2. $f(x_k^*)$

3. $(x_k^*, f(x_k^*)) \rightarrow \mathcal{D}_k \Rightarrow \mathcal{D}_{k+1}$

4. f_{k+1} - from \mathcal{D}_{k+1}

Exploration + Exploitation
greedy behavior

Gaussian process to construct $\hat{f}_k(x)$!

Regression to the Mean