

Probabilistic Artificial Intelligence

Tutorial 3: Gaussian Processes

Anastasiia Makarova

Marco Milanta

Department of Computer Science
Institute for Machine Learning

October 14, 2021

Objective

Possible assumptions on unknown f

- ▶ How to define a hypothesis set of functions to select from
 - ▶ prior knowledge
 - ▶ dependence on available features

Possible assumptions on unknown f

- ▶ How to define a hypothesis set of functions to select from
 - ▶ prior knowledge
 - ▶ dependence on available features

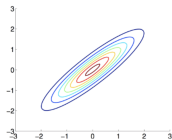
- ▶ Typical example: Lipschitz sets

$$\mathcal{F}_L = \{f : |f(\mathbf{x}) - f(\mathbf{x}')| \leq L|\mathbf{x} - \mathbf{x}'|\}$$

Possible assumptions on unknown f

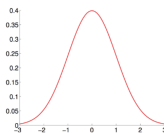
- ▶ How to define a hypothesis set of functions to select from
 - ▶ prior knowledge
 - ▶ dependence on available features
- ▶ Typical example: Lipschitz sets
$$\mathcal{F}_L = \{f : |f(\mathbf{x}) - f(\mathbf{x}')| \leq L|\mathbf{x} - \mathbf{x}'|\}$$
- ▶ How to impose regularity on a function in Bayesian case?

Gaussian are nice: let us recap why



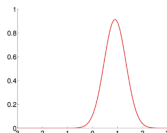
$$p(\mathbf{f}_1, \mathbf{f}_2) \sim \mathcal{N}(\mathbf{f}_1, \mathbf{f}_2 | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Joint



$$p(\mathbf{f}_1)$$

Marginal

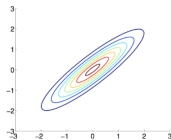


$$p(\mathbf{f}_1 | \mathbf{f}_2)$$

Conditional

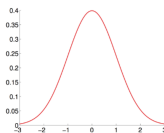
$$\begin{bmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

Gaussian are nice: let us recap why



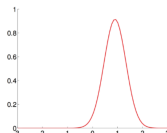
$$p(\mathbf{f}_1, \mathbf{f}_2) \sim \mathcal{N}(\mathbf{f}_1, \mathbf{f}_2 | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Joint



$$p(\mathbf{f}_1)$$

Marginal



$$p(\mathbf{f}_1 | \mathbf{f}_2)$$

Conditional

$$\begin{bmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

$$p(\mathbf{f}_1 \mid \mathbf{f}_2 = \mathbf{a}) = \mathcal{N}(\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\Sigma}})$$

$$\bar{\boldsymbol{\mu}} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{a} - \mu_2)$$

$$\bar{\boldsymbol{\Sigma}} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$

What is a GP: step from random variable to random function

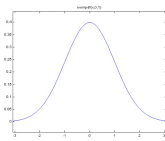
- Gaussian process = normal distribution over functions

What is a GP: step from random variable to random function

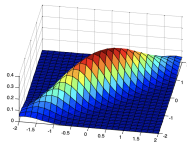
- Gaussian process = normal distribution over functions
- $GP(\mu, k)$ with mean function $\mu : \mathcal{X} \rightarrow \mathbb{R}$
and kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

What is a GP: step from random variable to random function

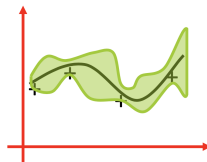
- Gaussian process = normal distribution over functions
- $GP(\mu, k)$ with mean function $\mu : \mathcal{X} \rightarrow \mathbb{R}$ and kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$



Normal dist.
(1-D Gaussian)



Multivariate normal
(n-D Gaussian)



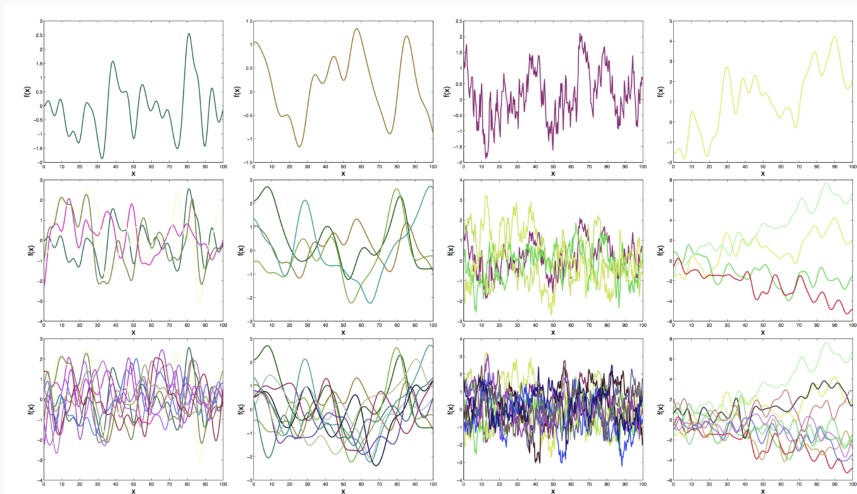
Gaussian process
(∞ -D Gaussian)

How to impose regularity on a function in Bayesian case?

Kernel function k encodes assumptions about correlation:

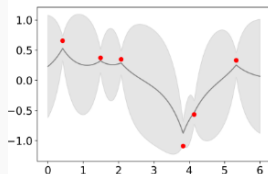
How to impose regularity on a function in Bayesian case?

Kernel function k encodes assumptions about correlation:

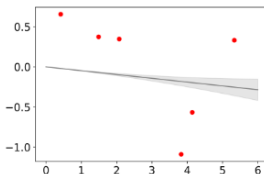


Quiz: Open Eduapp

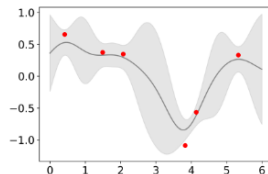
Quiz: Open Eduapp



(a) Laplace kernel



(b) Linear kernel



(c) Gaussian kernel

Bayesian Linear Regression and GPs

- ▶ Dataset $\{(x_1, y_1), \dots, (x_n, y_n)\}$ of points $x \in \mathcal{X}, y \in \mathbb{R}$
- ▶ Feature map $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$
- ▶ Our model

$$y = \phi(x)^T \mathbf{w} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

- ▶ Prior knowledge:

$$\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_d \sigma_p^2)$$

- ▶ Objective: predict y distribution in a new point x^*

Bayesian Linear Regression approach

Bayesian Linear Regression approach

Gaussian Process approach

Gaussian Process approach

GP vs BLR

► Bayesian Linear Regression

$$\mathbb{E}[y^*] = \phi(x^*)^T \frac{1}{\sigma_n^2} \left(\frac{1}{\sigma_n^2} \Phi^T \Phi + \frac{1}{\sigma_p^2} \mathbf{I}_d \right)^{-1} \Phi^T y_{1:m}$$

$$\sigma^2(y^*) = \phi(x^*)^T \left(\frac{1}{\sigma_n^2} \Phi^T \Phi + \frac{1}{\sigma_p^2} \mathbf{I}_d \right)^{-1} \phi(x^*) + \sigma_n^2$$

► Gaussian Process

$$\mathbb{E}[y^*] = \phi(x^*)^T \Phi^T \left(\Phi \Phi^T + \frac{\sigma_\varepsilon^2}{\sigma_p^2} \mathbf{I}_m \right)^{-1} y_{1:m}$$

$$\sigma^2(y^*) = \sigma_p^2 \left(\phi(x^*)^T \phi(x^*) - \phi(x^*)^T \Phi^T \left(\Phi \Phi^T + \frac{\sigma_\varepsilon^2}{\sigma_p^2} \mathbf{I}_m \right)^{-1} \Phi \phi(x^*) \right).$$

GP advantages

- ▶ Kernel formula:

$$y^* \mid y_{1:n} \sim \mathcal{N}(\tilde{\mu}, \tilde{\Sigma})$$

$$\tilde{\mu} = \mathbf{k}_{Ax^*}^T (\mathbf{K}_{AA} + \sigma_\varepsilon^2 \mathbf{I}_m)^{-1} y_{1:m}$$

$$\tilde{\Sigma} = k(x^*, x^*) - \mathbf{k}_{Ax^*}^t (\mathbf{K}_{AA} + \sigma_\varepsilon^2 \mathbf{I}_m)^{-1} \mathbf{k}_{Ax^*}$$

- ▶ High dimensionality regression

$$\phi \rightarrow \mathbb{R}^d$$

$$d = 100, d = 100000, d = \infty?$$

GP advantages

- ▶ Kernel formula:

$$y^* \mid y_{1:n} \sim \mathcal{N}(\tilde{\mu}, \tilde{\Sigma})$$

$$\tilde{\mu} = \mathbf{k}_{Ax^*}^T (\mathbf{K}_{AA} + \sigma_\varepsilon^2 \mathbf{I}_m)^{-1} y_{1:m}$$

$$\tilde{\Sigma} = k(x^*, x^*) - \mathbf{k}_{Ax^*}^t (\mathbf{K}_{AA} + \sigma_\varepsilon^2 \mathbf{I}_m)^{-1} \mathbf{k}_{Ax^*}$$

- ▶ High dimensionality regression

$$\phi \rightarrow \mathbb{R}^d$$

$d = 100$, $d = 100000$, $d = \infty$? If it converges yes

$$k(x, x') = \sigma_p^2 \sum_{i=1}^{\infty} \phi(x)_i \phi(x')_i < \infty$$

- ▶ Gaussian Kernel / Laplace Kernel / Matérn kernel

GP regression in practice

demo Here

Kalman filters?

Recap of Kalman filter and reference to hw task

Links

- ▶ Derivations of BLR and GPs: [Link](#)
- ▶ GP visualization notebook: [Link](#)