# Lecture 22: Gaussian process regression

Professor Ilias Bilionis

## Gaussian process regression with measurement noise

# The likelihood of the observations

$$\varepsilon_i \sim N(0, \sigma^2)$$

$$x_{1:n} = (x_1, \cdots, x_n) \; ; \; y_{1:n} = (y_1, \cdots, y_n) \; ; \; y_i = f(x_i) + \varepsilon_i .$$

likelihood of
a single
observation

$$p(y_i \mid f(x_i)) = N(y_i \mid f(x_i), \sigma^2)$$

$$f_{1:n} = (f(x_1), \cdots, f(x_n))$$

$$p(y_{1:n} \mid f_{1:n}) = \prod_{i=1}^{n} p(y_i \mid f(x_i)) = N(y_{1:n} \mid f_{1:n}, \sigma^2 I)$$
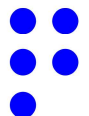
Likelihood of observed data.

# The joint probability density over observations and test points

$$\text{test points}: x_{1:n*}^* = (x_1^*, \ldots, x_n^*)$$

$$f_{1:n*}^* = (f(x_1^*), \ldots, f(x_n^*))$$

$$f(\cdot) \sim \mathcal{GP}(m(\cdot), c(\cdot, \cdot))$$

$$p(f_{1:n}, f_{1:n*}^* \mid x_{1:n}, x_{1:n*}^*) = \mathcal{N}\left(\begin{matrix} f_{1:n} \\ f_{1:n*}^* \end{matrix} \middle| \begin{pmatrix} m_{1:n} \\ m_{1:n*} \end{pmatrix}, \begin{pmatrix} C_n & B \\ B^T & C_{n*} \end{pmatrix}\right)$$

# Conditioning on observations

**Likelihood:**

$$P(y_{1:n} \mid f_{1:n}) = \prod_{i=1}^{n} P(y_i \mid f(x_i)) = N\left(y_{1:n} \mid f_{1:n}, \sigma^2 I\right)$$

**Joint:**

$$P\left(f_{1:n}, f_{1:n}^* \mid x_{1:n}, x_{1:n}^*\right) = N\left(\begin{pmatrix} f_{1:n} \\ f_{1:n}^* \end{pmatrix} \middle| \begin{pmatrix} m_{1:n} \\ m_{1:n}^* \end{pmatrix}, \begin{pmatrix} C_n & B \\ B^T & C_{n*} \end{pmatrix}\right)$$

We are after:

*this posterior*

$$P\left(f_{1:n}^* \mid x_{1:n}, y_{1:n}, x_{1:n}^*\right) \underset{\text{Rule}}{\overset{\text{Sum}}{=}} \int P\left(f_{1:n}, f_{1:n}^* \mid x_{1:n}, y_{1:n}, x_{1:n}^*\right) df_{1:n}$$

*joint posterior*

$$\underset{\text{Rule}}{\overset{\text{Bayes}}{\propto}} \int P(y_{1:n} \mid x_{1:n}, f_{1:n}) P\left(f_{1:n}, f_{1:n}^* \mid x_{1:n}, x_{1:n}^*\right) df_{1:n}$$

*likelihood*      *joint of observed & test inputs*

$$\overset{\text{Complete}}{\underset{\text{the square}}{=}} N\left(f_{1:n}^* \mid m_{1:n}^*, C_{n*}^*\right)$$

$$m_{1:n}^* = m_{1:n}^* - B^T \left[C_n + \sigma^2 I_n\right]^{-1} (y_{1:n} - m_{1:n})$$

$$C_{n*}^* = C_{n*} - B^T \left[C_n + \sigma^2 I_n\right]^{-1} B$$

*corrections*

# The posterior Gaussian process

test inputs
are arbitrary

$\Downarrow$

posterior
mean
function

posterior
covariance
function

posterior
Gaussian Process

$$f(\cdot) \mid x_{1:n}, y_{1:n} \sim \mathcal{GP}\left(m_n^*(\cdot), c_n^*(\cdot, \cdot)\right)$$
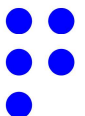
$$m_n^*(x) = m(x) - c(x, x_{1:n})\left[C_n + \sigma^2 I_n\right]^{-1}(y_{1:n} - m_{1:n})$$

$$c_n^*(x, x') = c(x, x') - c(x, x_{1:n})\left[C_n + \sigma^2 I_n\right]^{-1} c(x_{1:n}, x')$$

$1 \times n$

$\left(c(x, x_1), \cdots, c(x, x_n)\right)$

noise
component

$n \times 1$

$$\begin{pmatrix} c(x_1, x') \\ \vdots \\ c(x_n, x') \end{pmatrix}$$

✶ This summarizes everything about the functions
after you have seen the data

**PREDICTIVE SCIENCE LABORATORY**

13

# The point predictive distribution

function value
at $x$

post.

① $p(f(x) \mid x_{1:n}, y_{1:n}) \underset{GP}{=} N\left(f(x) \mid \mu_n^*(x), \sigma_n^{*2}(x)\right)$

$\hookrightarrow$ Best I can say about the function values. $C_n^*(x, x)$

Uncertainty here is epistemic:

limited number of
observations being used

$\underset{G}{\underline{\quad}} \times \underset{G}{\underline{\quad}}$

measurement

Sum

② $p(y \mid x, x_{1:n}, y_{1:n}) \underset{Rule}{=} \int \underbrace{p(y \mid f(x))}_{N(y \mid f(x), \sigma^2)} p(f(x) \mid x_{1:n}, y_{1:n}) \, df(x)$

first introduced
here

complete
$\underset{the \ square}{=} N\left(y \mid \mu_n^*(x), \underset{epistemic}{\sigma_n^{*2}(x)} + \sigma^2\right)$ ✦

aleatory

14

$\Downarrow$ 95% credible int.

① $f(x) \in \left[\mu_n^*(x) - 2\sigma_n^*(x), \mu_n^*(x) + 2\sigma_n^*(x)\right]$

② $y \in \left[\mu_n^*(x) - 2\sqrt{\sigma_n^{*2}(x) + \sigma^2}, \mu_n^*(x) + 2\sqrt{\sigma_n^{*2}(x) + \sigma^2}\right]$

# Gaussian process regression - Noisy observations

The smaller you make the noise variance, the more the GP trusts the data



Each choice of the noise corresponds to a different interpretation of the data.

# Even when there is not any noise, including it improves numerical stability

- It is common to use small noise even if there is not any in the data.

- Cholesky fails when covariance is close to being semi-positive definite.

- Adding a small noise improves numerical stability.

- It is known as the "jitter" or as the "nugget" in this case.