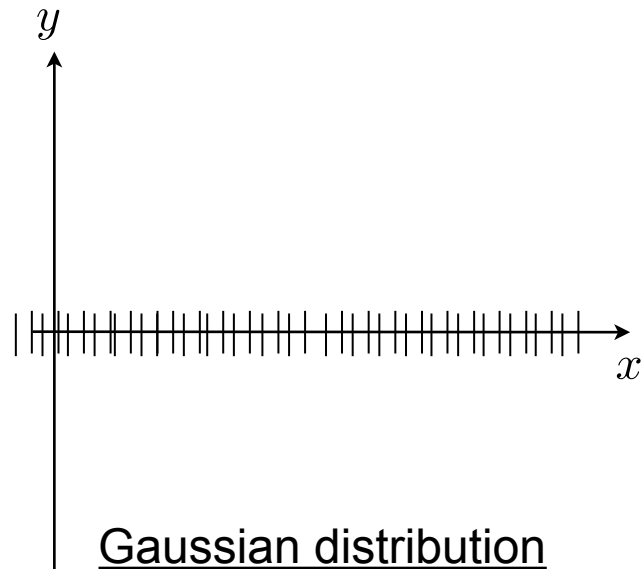


# What is Gaussian process?

# From distribution to process

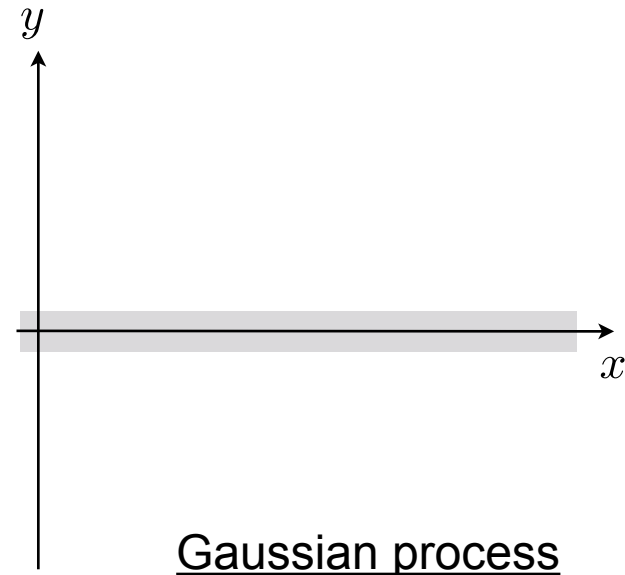
Distribution: probability defined in a **finite** dimensional space.

Process: probability defined in an **infinite** dimensional space (e.g., **function**).



$$p(\mathbf{y}) = \text{Gauss}_{100000}(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\begin{aligned}\mathbf{y} &\in \mathbb{R}^{100000}, \\ \boldsymbol{\mu} &\in \mathbb{R}^{100000}, \\ \boldsymbol{\Sigma} &\in \mathbb{R}^{100000 \times 100000},\end{aligned}$$



$$p(y) = \text{GP}(y; \nu, k)$$

$$\begin{aligned}y(x) &\in \mathbb{R} \text{ for } x \in \mathbb{R}, \\ \nu(x) &\in \mathbb{R} \text{ for } x \in \mathbb{R}, \\ k(x, x') &\in \mathbb{R} \text{ for } (x, x') \in \mathbb{R}^2,\end{aligned}$$

# Formal definition

A Gaussian process is a collection of random variables,  
any of finite number of which have a joint Gaussian distribution.

[Rasmussen&Williams:2016]

# Practical property 1 (Marginalization property)

Random functions  $f(\cdot) : \mathbb{R}^L \mapsto \mathbb{R}$  follow the Gaussian process  $f \sim \text{GP}(\nu(\cdot), k(\cdot, \cdot))$  with a mean function  $\nu(\cdot)$  and a kernel function  $k(\cdot, \cdot)$ .



For any given input set  $\{\mathbf{x}^{(n)} \in \mathbb{R}^L\}_{n=1}^N$ , it holds that

$$p(\mathbf{f} | \boldsymbol{\nu}, \mathbf{K}) \propto \exp \left( - \frac{(\mathbf{f} - \boldsymbol{\nu})^\top \mathbf{K}^{-1} (\mathbf{f} - \boldsymbol{\nu})}{2} \right),$$

where  $\mathbf{f} = \left( f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(N)}) \right)^\top,$

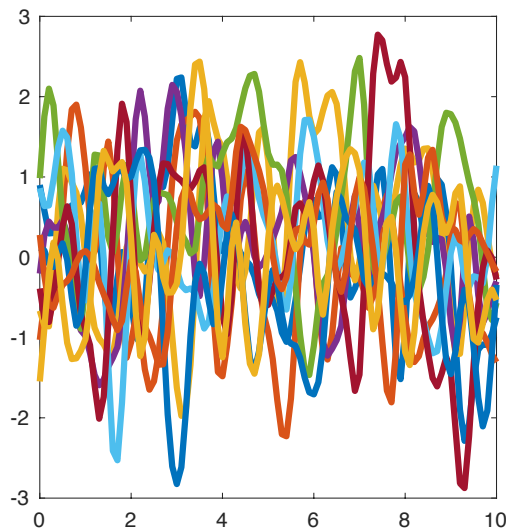
$$\boldsymbol{\nu} = \left( \nu(\mathbf{x}^{(1)}), \dots, \nu(\mathbf{x}^{(N)}) \right)^\top,$$

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ \vdots & & \vdots \\ k(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{pmatrix}.$$

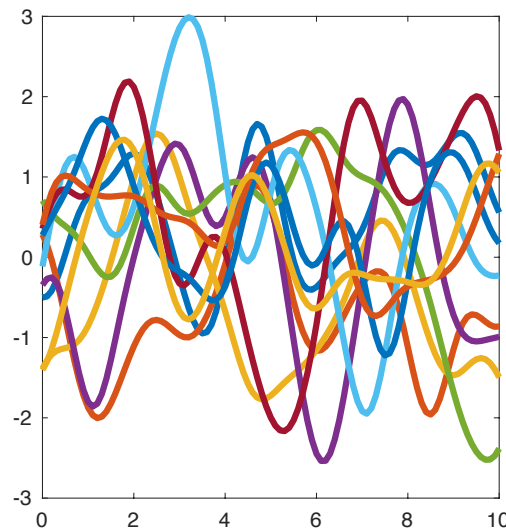
Assume infinite dimension, work in finite dimension.

# Examples

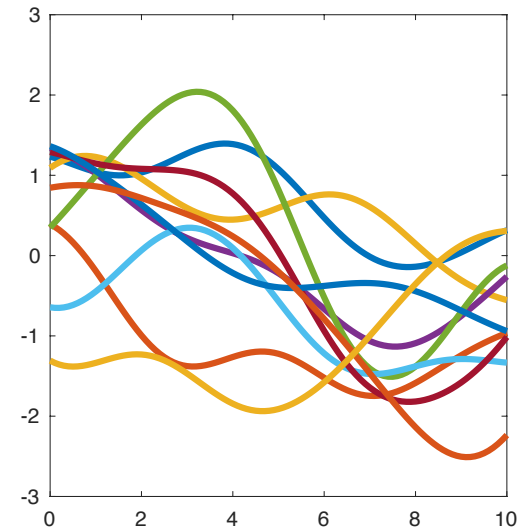
$$f \sim \text{GP}(\nu, k), \quad \text{where } \nu(\mathbf{x}) \equiv 0, \quad k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\gamma^2}\right).$$



$$\gamma^2 = 0.1$$



$$\gamma^2 = 1$$



$$\gamma^2 = 10$$

# Practical property 2

Given another input set  $\{\mathbf{x}^{*(n)} \in \mathbb{R}^L\}_{n=1}^{N^*}$ , the joint distribution is

$$p(\mathbf{f}, \mathbf{f}^*) \propto \exp \left( - \frac{\left( \begin{pmatrix} \mathbf{f} \\ \mathbf{f}^* \end{pmatrix} - \begin{pmatrix} \boldsymbol{\nu} \\ \boldsymbol{\nu}^* \end{pmatrix} \right)^\top \begin{pmatrix} \mathbf{K} & \mathbf{K}^* \\ \mathbf{K}^{*\top} & \mathbf{K}^{**} \end{pmatrix}^{-1} \left( \begin{pmatrix} \mathbf{f} \\ \mathbf{f}^* \end{pmatrix} - \begin{pmatrix} \boldsymbol{\nu} \\ \boldsymbol{\nu}^* \end{pmatrix} \right)}{2} \right),$$

where  $\mathbf{K}^*$  is the kernel between the two sets,  $\{\mathbf{x}^{(n)} \in \mathbb{R}^L\}_{n=1}^N$  and  $\{\mathbf{x}^{*(n)} \in \mathbb{R}^L\}_{n=1}^{N^*}$ , and  $\mathbf{K}^{**}$  is the kernel within the second set,  $\{\mathbf{x}^{*(n)} \in \mathbb{R}^L\}_{n=1}^{N^*}$ .

The conditional distribution of  $\mathbf{f}^*$  given  $\mathbf{f}$  is written as

$$p(\mathbf{f}^* | \mathbf{f}) \propto \exp \left( - \frac{(\mathbf{f}^* - \boldsymbol{\mu}_c)^\top \boldsymbol{\Sigma}_c^{-1} (\mathbf{f}^* - \boldsymbol{\mu}_c)}{2} \right),$$

where

$$\begin{aligned} \boldsymbol{\mu}_c &= \boldsymbol{\nu}^* + \mathbf{K}^{*\top} \mathbf{K}^{-1} (\mathbf{f} - \boldsymbol{\nu}), \\ \boldsymbol{\Sigma}_c &= \mathbf{K}^{**} - \mathbf{K}^{*\top} \mathbf{K}^{-1} \mathbf{K}^*. \end{aligned}$$

Information transfer from training locations to test locations!

# GP regression

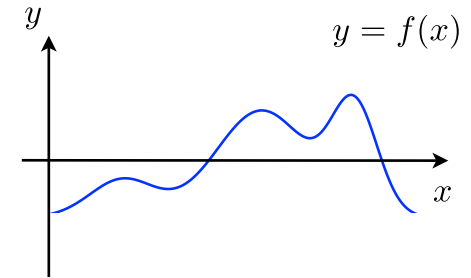
# GP regression

Observed:  $\mathbf{x} \in \mathbb{R}^L, y \in \mathbb{R}$

Parameter:  $f : \mathbb{R}^L \rightarrow \mathbb{R}$

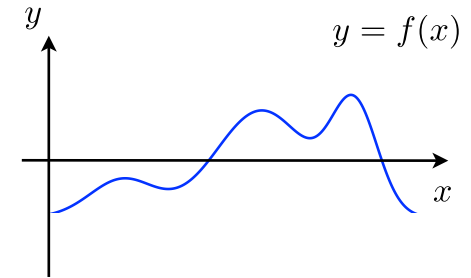
Model dist.:  $p(y|f) \propto \exp\left(-\frac{\|y - f(\mathbf{x})\|^2}{2\sigma^2}\right)$

Prior dist.:  $p(f) = \text{GP}(f; 0, k)$  with  $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\gamma^2}\right)$ .





# GP regression



Observed:  $\mathbf{x} \in \mathbb{R}^L, y \in \mathbb{R}$

Parameter:  $f : \mathbb{R}^L \rightarrow \mathbb{R}$

Model dist.:  $p(y|f) \propto \exp\left(-\frac{\|y - f(\mathbf{x})\|^2}{2\sigma^2}\right)$

Prior dist.:  $p(f) = \text{GP}(f; 0, k)$  with  $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\gamma^2}\right)$ .

Data set:  $\{\mathbf{x}^{(n)}, y^{(n)}\}_{n=1}^N$



$$p(\mathbf{y}|\mathbf{f}) \propto \exp\left(-\frac{\|\mathbf{y} - \mathbf{f}\|^2}{2\sigma^2}\right)$$

$$p(\mathbf{f}) \propto \exp\left(-\frac{\mathbf{f}^\top \mathbf{K}^{-1} \mathbf{f}}{2}\right)$$

$$\mathbf{y} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{pmatrix}$$

$$\mathbf{f} = \begin{pmatrix} f(\mathbf{x}^{(1)}) \\ \vdots \\ f(\mathbf{x}^{(N)}) \end{pmatrix}$$

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ \vdots & & \vdots \\ k(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & \dots & k(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{pmatrix}.$$

# GP regression

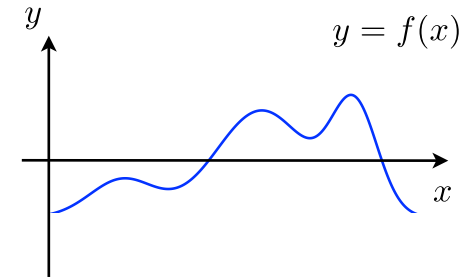
Posterior on  $\mathbf{f}$ :

$$\begin{aligned}
 p(\mathbf{f}|\mathbf{y}) &\propto p(\mathbf{y}|\mathbf{f})p(\mathbf{f}) \\
 &\propto \exp\left(-\frac{\|\mathbf{y}-\mathbf{f}\|^2}{2\sigma^2} - \frac{\mathbf{f}^\top \mathbf{K}^{-1} \mathbf{f}}{2}\right) \\
 &\propto \exp\left(-\frac{-2\sigma^{-2}\mathbf{y}^\top \mathbf{f} + \mathbf{f}^\top (\mathbf{K}^{-1} + \sigma^{-2}\mathbf{I}_N) \mathbf{f}}{2}\right) \\
 &\propto \exp\left(-\frac{(\mathbf{f}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{f}-\boldsymbol{\mu})}{2}\right),
 \end{aligned}$$

where

$$\begin{aligned}
 \boldsymbol{\mu} &= \sigma^{-2} \boldsymbol{\Sigma} \mathbf{y}, \\
 \boldsymbol{\Sigma} &= (\mathbf{K}^{-1} + \sigma^{-2} \mathbf{I}_N)^{-1}.
 \end{aligned}$$

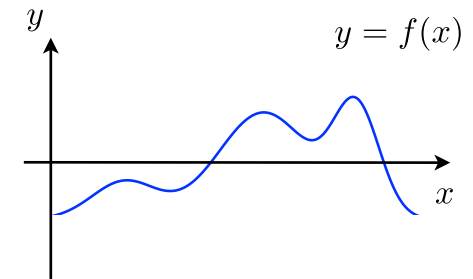
Posterior at training data points are Gaussian!



# GP regression

Predictive for test points  $\{\mathbf{x}^{*(n)}\}_{n=1}^{N^*}$ :

$$\begin{aligned}
 p(\mathbf{y}^* | \mathbf{y}) &= \int p(\mathbf{y}^* | f) p(f | \mathbf{y}) df \\
 &= \int p(\mathbf{y}^* | \mathbf{f}^*) \int p(\mathbf{f}, \mathbf{f}^* | \mathbf{y}) d\mathbf{f} d\mathbf{f}^* \\
 &= \underbrace{\int p(\mathbf{y}^* | \mathbf{f}^*)}_{\text{model dist.}} \underbrace{\int p(\mathbf{f}^* | \mathbf{f})}_{\text{conditional (transfer)}} \underbrace{p(\mathbf{f} | \mathbf{y})}_{\text{posterior on } \mathbf{f}} d\mathbf{f} d\mathbf{f}^*
 \end{aligned}$$



$$\mathbf{y}^* = \begin{pmatrix} y^{*(1)} \\ \vdots \\ y^{*(N)} \end{pmatrix}$$

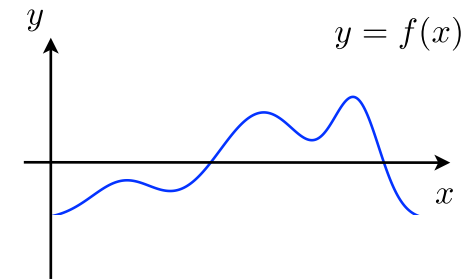
$$\mathbf{f}^* = \begin{pmatrix} f(\mathbf{x}^{*(1)}) \\ \vdots \\ f(\mathbf{x}^{*(N)}) \end{pmatrix}$$

Everything is Gaussian!

# GP regression

Predictive for test points  $\{\mathbf{x}^{*(n)}\}_{n=1}^{N^*}$ :

$$\begin{aligned}
 p(\mathbf{y}^*|\mathbf{y}) &= \int p(\mathbf{y}^*|f)p(f|\mathbf{y})df \\
 &= \int p(\mathbf{y}^*|\mathbf{f}^*) \int p(\mathbf{f}, \mathbf{f}^*|\mathbf{y})d\mathbf{f}d\mathbf{f}^* \\
 &= \int p(\mathbf{y}^*|\mathbf{f}^*) \underbrace{\int p(\mathbf{f}^*|\mathbf{f})p(\mathbf{f}|\mathbf{y})d\mathbf{f}}_{\text{posterior on } \mathbf{f}^*}d\mathbf{f}^*
 \end{aligned}$$



$$\mathbf{y}^* = \begin{pmatrix} y^{*(1)} \\ \vdots \\ y^{*(N)} \end{pmatrix}$$

$$\mathbf{f}^* = \begin{pmatrix} f(\mathbf{x}^{*(1)}) \\ \vdots \\ f(\mathbf{x}^{*(N)}) \end{pmatrix}$$

# Practical property 2

Given another input set  $\{\mathbf{x}^{*(n)} \in \mathbb{R}^L\}_{n=1}^{N^*}$ , the joint distribution is

$$p(\mathbf{f}, \mathbf{f}^*) \propto \exp \left( - \frac{\left( \begin{pmatrix} \mathbf{f} \\ \mathbf{f}^* \end{pmatrix} - \begin{pmatrix} \boldsymbol{\nu} \\ \boldsymbol{\nu}^* \end{pmatrix} \right)^\top \begin{pmatrix} \mathbf{K} & \mathbf{K}^* \\ \mathbf{K}^{*\top} & \mathbf{K}^{**} \end{pmatrix}^{-1} \left( \begin{pmatrix} \mathbf{f} \\ \mathbf{f}^* \end{pmatrix} - \begin{pmatrix} \boldsymbol{\nu} \\ \boldsymbol{\nu}^* \end{pmatrix} \right)}{2} \right),$$

where  $\mathbf{K}^*$  is the kernel between the two sets,  $\{\mathbf{x}^{(n)} \in \mathbb{R}^L\}_{n=1}^N$  and  $\{\mathbf{x}^{*(n)} \in \mathbb{R}^L\}_{n=1}^{N^*}$ , and  $\mathbf{K}^{**}$  is the kernel within the second set,  $\{\mathbf{x}^{*(n)} \in \mathbb{R}^L\}_{n=1}^{N^*}$ .

The conditional distribution of  $\mathbf{f}^*$  given  $\mathbf{f}$  is written as

$$p(\mathbf{f}^* | \mathbf{f}) \propto \exp \left( - \frac{(\mathbf{f}^* - \boldsymbol{\mu}_c)^\top \boldsymbol{\Sigma}_c^{-1} (\mathbf{f}^* - \boldsymbol{\mu}_c)}{2} \right),$$

where

$$\boldsymbol{\mu}_c = \boldsymbol{\nu}^* + \mathbf{K}^{*\top} \mathbf{K}^{-1} (\mathbf{f} - \boldsymbol{\nu}),$$

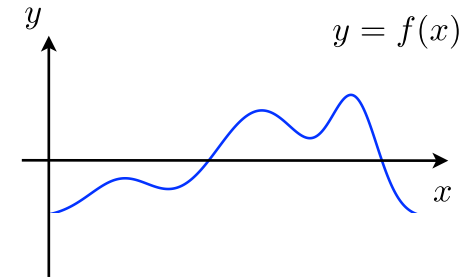
$$\boldsymbol{\Sigma}_c = \mathbf{K}^{**} - \mathbf{K}^{*\top} \mathbf{K}^{-1} \mathbf{K}^*.$$

Information transfer from the training locations to the test locations!

# GP regression

Predictive for test points  $\{\mathbf{x}^{*(n)}\}_{n=1}^{N^*}$ :

$$\begin{aligned}
 p(\mathbf{y}^* | \mathbf{y}) &= \int p(\mathbf{y}^* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} \\
 &= \int p(\mathbf{y}^* | \mathbf{f}^*) \int p(\mathbf{f}, \mathbf{f}^* | \mathbf{y}) d\mathbf{f} d\mathbf{f}^* \\
 &= \int p(\mathbf{y}^* | \mathbf{f}^*) \underbrace{\int p(\mathbf{f}^* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f}}_{\text{posterior on } \mathbf{f}^*} d\mathbf{f}^*
 \end{aligned}$$



$$\mathbf{y}^* = \begin{pmatrix} y^{*(1)} \\ \vdots \\ y^{*(N)} \end{pmatrix}$$

$$\mathbf{f}^* = \begin{pmatrix} f(\mathbf{x}^{*(1)}) \\ \vdots \\ f(\mathbf{x}^{*(N)}) \end{pmatrix}$$

After some algebra...

$$p(\mathbf{f}^* | \mathbf{y}) = \int p(\mathbf{f}^* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} \propto \exp \left( -\frac{(\mathbf{f}^* - \boldsymbol{\mu}_f)^\top \boldsymbol{\Sigma}_f^{-1} (\mathbf{f}^* - \boldsymbol{\mu}_f)}{2} \right),$$

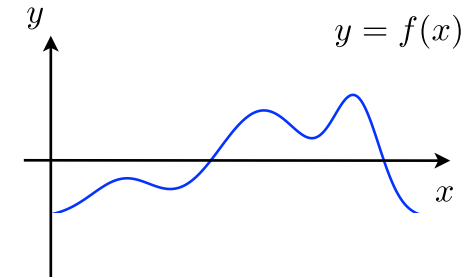
$$\text{where } \boldsymbol{\mu}_f = \mathbf{K}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y},$$

$$\boldsymbol{\Sigma}_f = \mathbf{K}^{**} - \mathbf{K}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{K}^*.$$

Posterior at test data (any) points are Gaussian!

# GP regression

Predictive for test points  $\{\mathbf{x}^{*(n)}\}_{n=1}^{N^*}$ :



$$\begin{aligned}
 p(\mathbf{y}^* | \mathbf{y}) &= \int p(\mathbf{y}^* | f) p(f | \mathbf{y}) df \\
 &= \int p(\mathbf{y}^* | \mathbf{f}^*) \int p(\mathbf{f}, \mathbf{f}^* | \mathbf{y}) d\mathbf{f} d\mathbf{f}^* \\
 &= \int p(\mathbf{y}^* | \mathbf{f}^*) \underbrace{\int p(\mathbf{f}^* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f}}_{\text{posterior on } \mathbf{f}^*} d\mathbf{f}^*
 \end{aligned}$$

$$\mathbf{y}^* = \begin{pmatrix} y^{*(1)} \\ \vdots \\ y^{*(N)} \end{pmatrix}$$

$$\mathbf{f}^* = \begin{pmatrix} f(\mathbf{x}^{*(1)}) \\ \vdots \\ f(\mathbf{x}^{*(N)}) \end{pmatrix}$$

After some algebra...

$$p(\mathbf{f}^* | \mathbf{y}) = \int p(\mathbf{f}^* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} \propto \exp \left( -\frac{(\mathbf{f}^* - \boldsymbol{\mu}_f)^\top \boldsymbol{\Sigma}_f^{-1} (\mathbf{f}^* - \boldsymbol{\mu}_f)}{2} \right),$$

$$\text{where } \boldsymbol{\mu}_f = \mathbf{K}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y},$$

$$\boldsymbol{\Sigma}_f = \mathbf{K}^{**} - \mathbf{K}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{K}^*.$$



Posterior is GP!  $p(f | \mathbf{y}) = \text{GP}(f; \boldsymbol{\mu}_f, \boldsymbol{\Sigma}_f)$

# GP regression

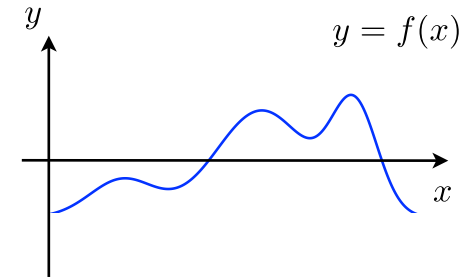
Predictive for test points  $\{\mathbf{x}^{*(n)}\}_{n=1}^{N^*}$ :

$$p(\mathbf{y}^*|\mathbf{y}) \propto \exp \left( -\frac{(\mathbf{y} - \boldsymbol{\mu}_y)^\top \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y)}{2} \right),$$

where

$$\boldsymbol{\mu}_y = \mathbf{K}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y},$$

$$\boldsymbol{\Sigma}_y = \mathbf{K}^{**} - \underbrace{\mathbf{K}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{K}^*}_{\text{blue bracket}} + \sigma^2 \mathbf{I}_{N^*}.$$

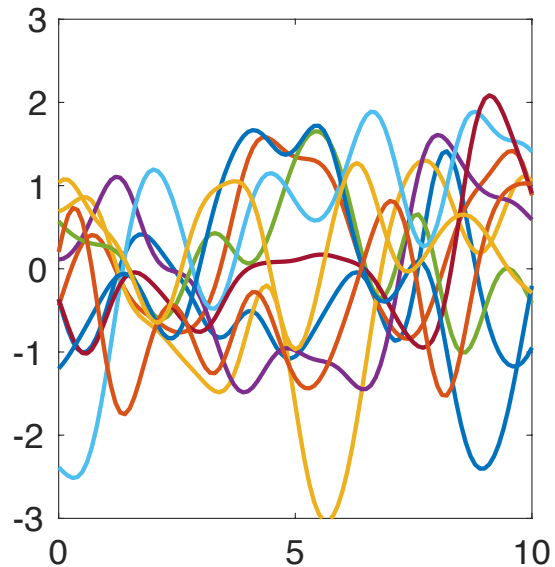


Bayesian predictive is analytically obtained without any approximation!



# GP regression

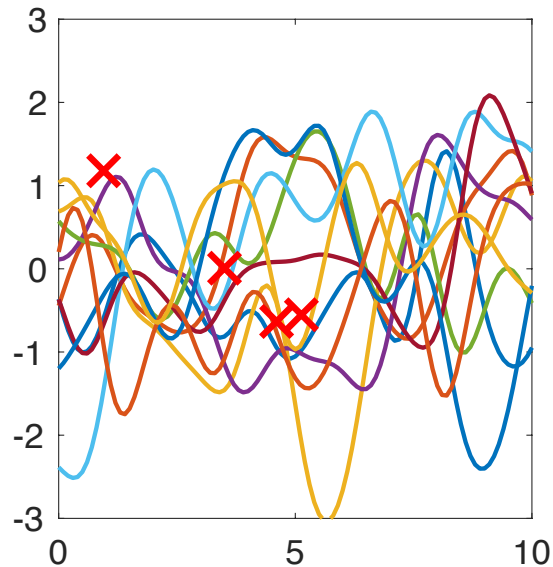
Prior dist.:  $p(f) = \text{GP}(f; 0, k)$  with  $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\gamma^2}\right)$ .



# GP regression

Prior dist.:  $p(f) = \text{GP}(f; 0, k)$  with  $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\gamma^2}\right)$ .

observation!

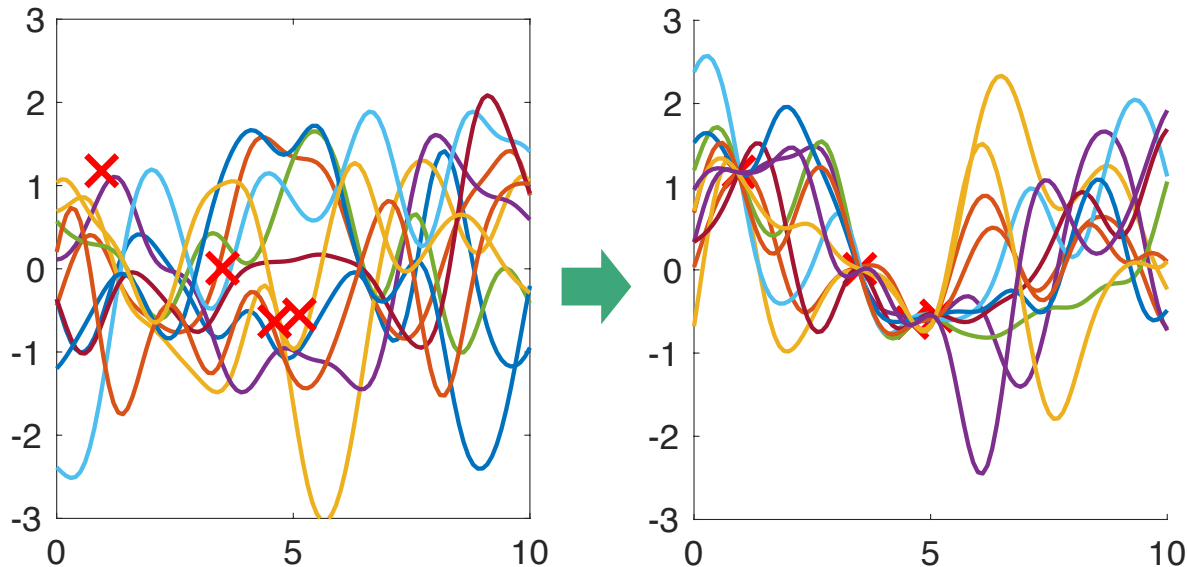


# GP regression

Prior dist.:  $p(f) = \text{GP}(f; 0, k)$  with  $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\gamma^2}\right)$ .

Posterior:  $p(f|\mathbf{y}) = \text{GP}(f; \boldsymbol{\mu}_f, \boldsymbol{\Sigma}_f)$

observation!



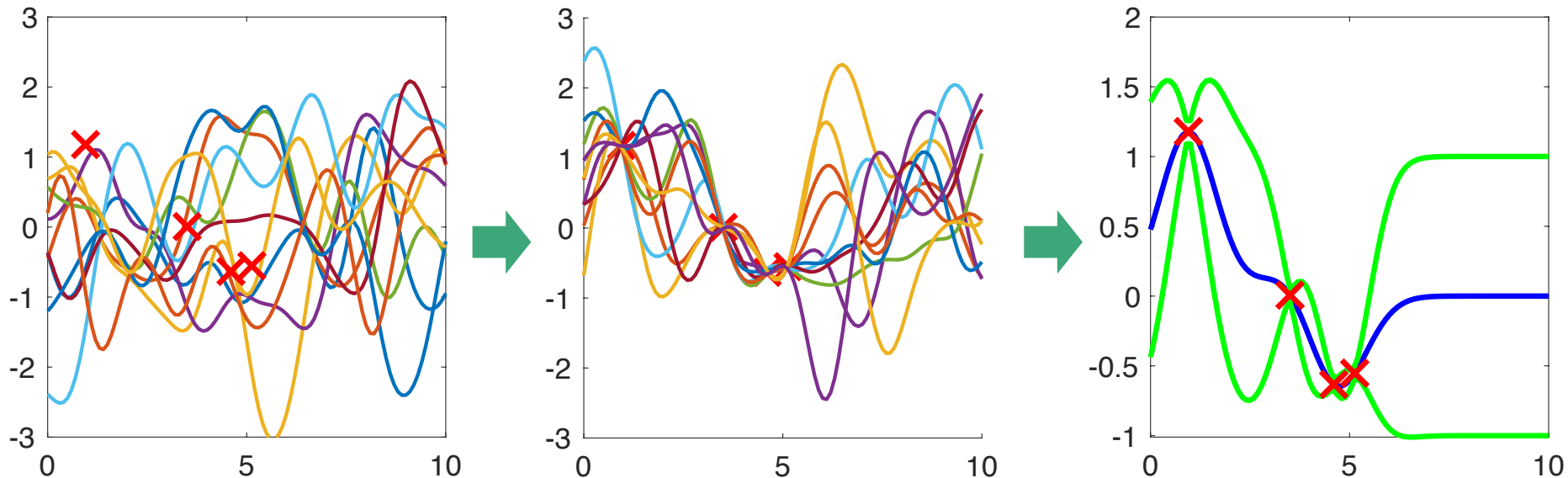
# GP regression

Prior dist.:  $p(f) = \text{GP}(f; 0, k)$  with  $k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{\gamma^2}\right)$ .

Posterior:  $p(f|y) = \text{GP}(f; \mu_f, \Sigma_f)$

Predictive:  $p(y^*|y) \propto \exp\left(-\frac{(y - \mu_y)^\top \Sigma_y^{-1} (y - \mu_y)}{2}\right),$

observation!

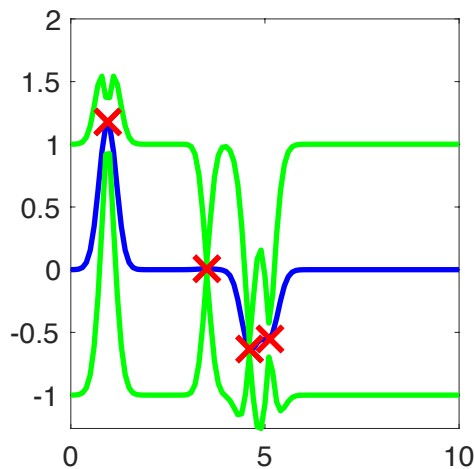


# GP regression

Prior dist.:  $p(f) = \text{GP}(f; 0, k)$  with  $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\underbrace{\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\gamma^2}}\right)$ .

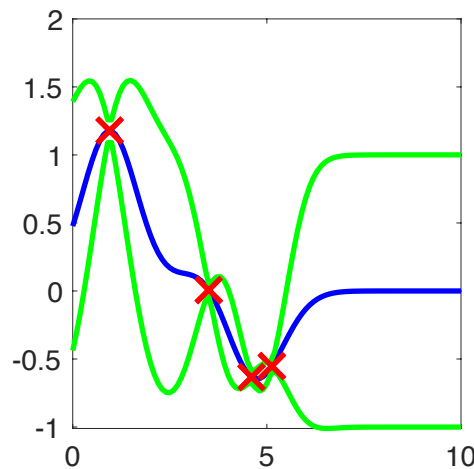
Model selection with marginal likelihood or cross validation,  
but note that criterion is task dependent!

$$p(\mathbf{y}|\mathbf{K}) = (2\pi)^{-N/2} |\sigma^2 \mathbf{I}_N + \mathbf{K}|^{-1/2} \exp\left(-\frac{\mathbf{y}^\top (\sigma^2 \mathbf{I}_N + \mathbf{K})^{-1} \mathbf{y}}{2}\right),$$



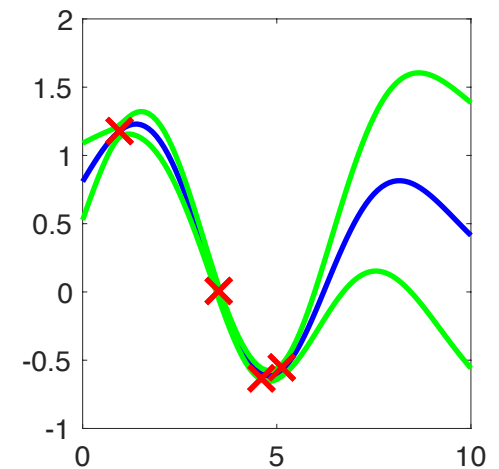
$\gamma^2 = 0.1$

$$-\log p(\mathbf{y}|\gamma^2) = 1.4544$$



$\gamma^2 = 1$

$$-\log p(\mathbf{y}|\gamma^2) = 1.4490$$



$\gamma^2 = 10$

$$-\log p(\mathbf{y}|\gamma^2) = 0.6408$$