

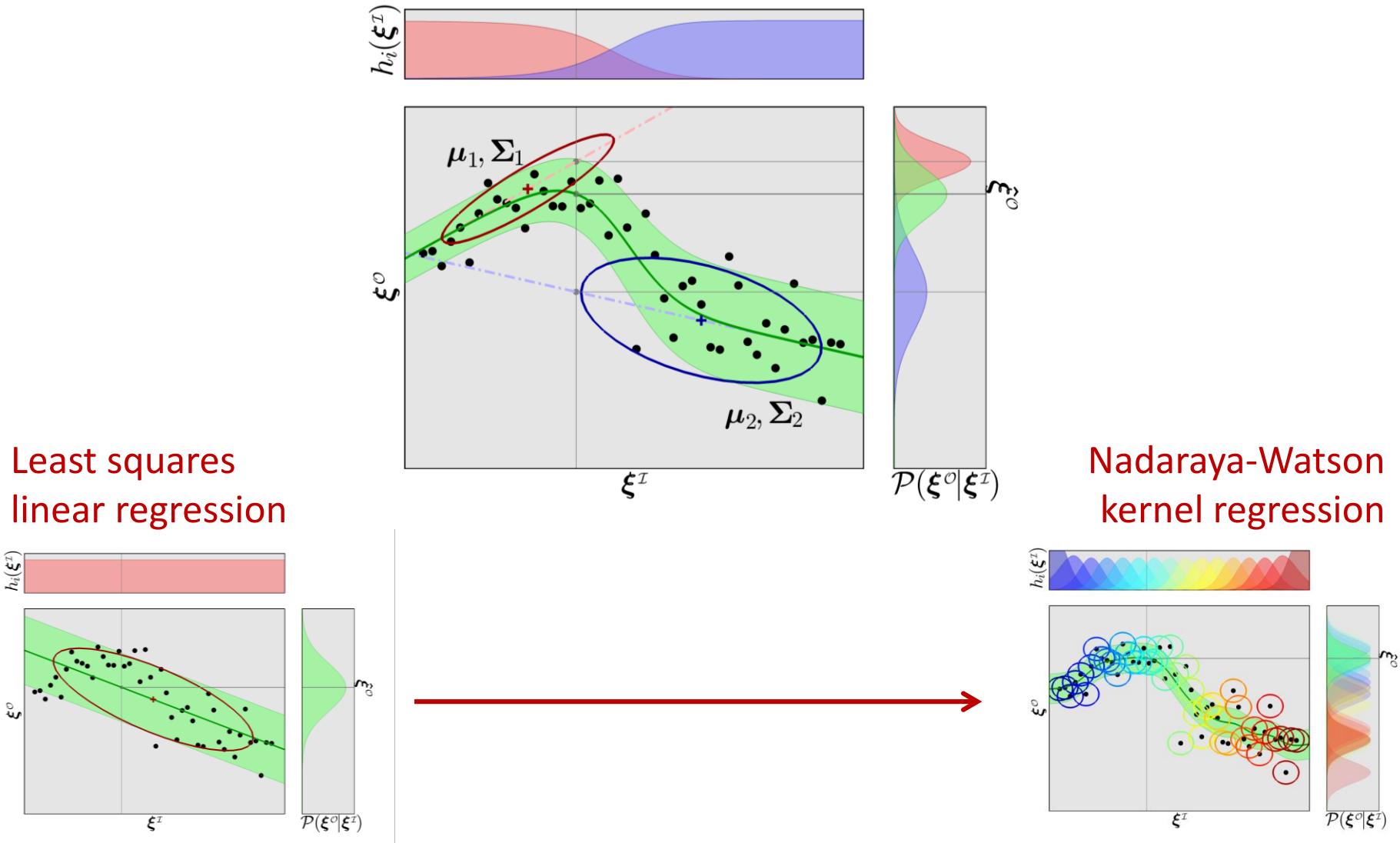
**EE613**  
**Machine Learning for Engineers**

**NONLINEAR REGRESSION II**

**Sylvain Calinon**  
**Robot Learning & Interaction Group**  
**Idiap Research Institute**  
**Dec. 19, 2019**

# Gaussian mixture regression (GMR)

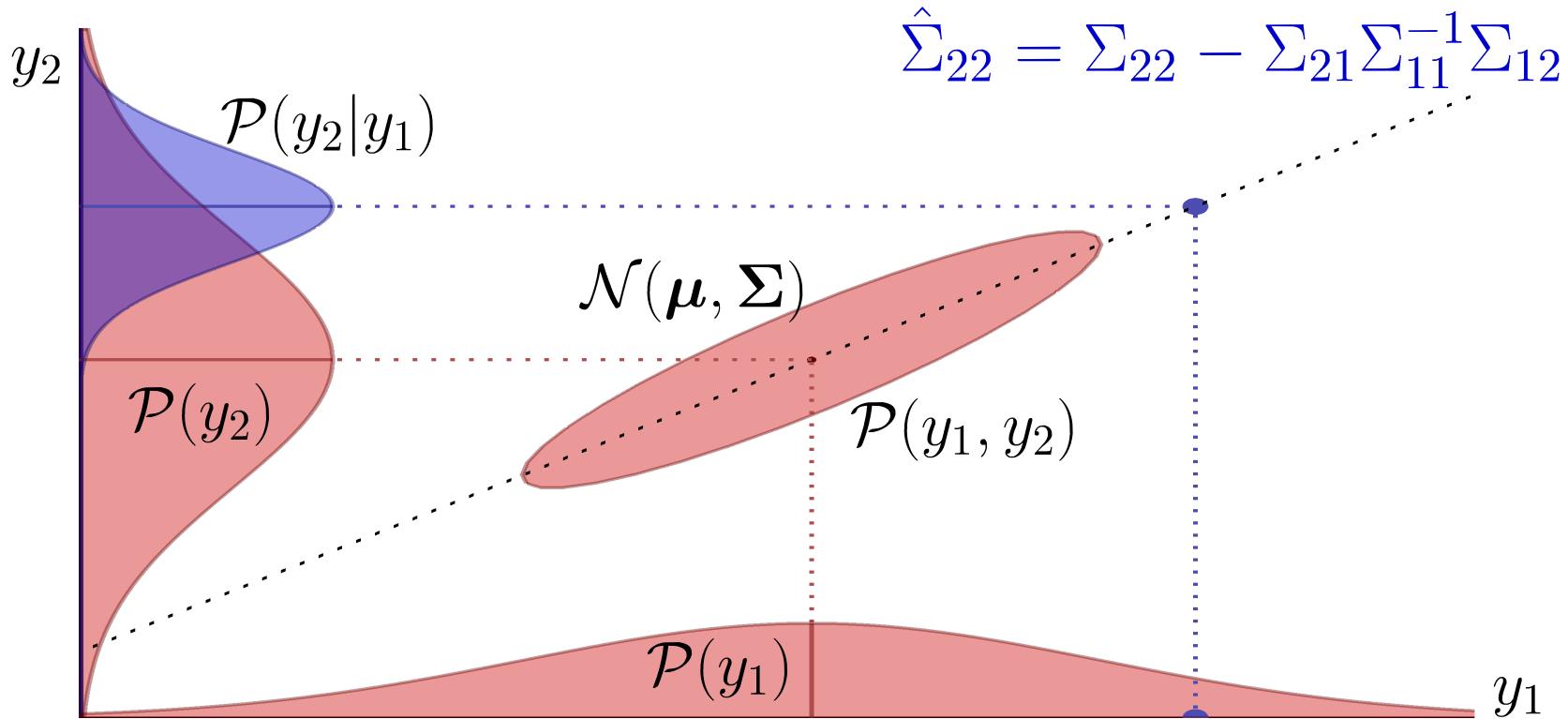
Nonlinear regression I



# Gaussian process (GP)

# Gaussian process - Informal interpretation

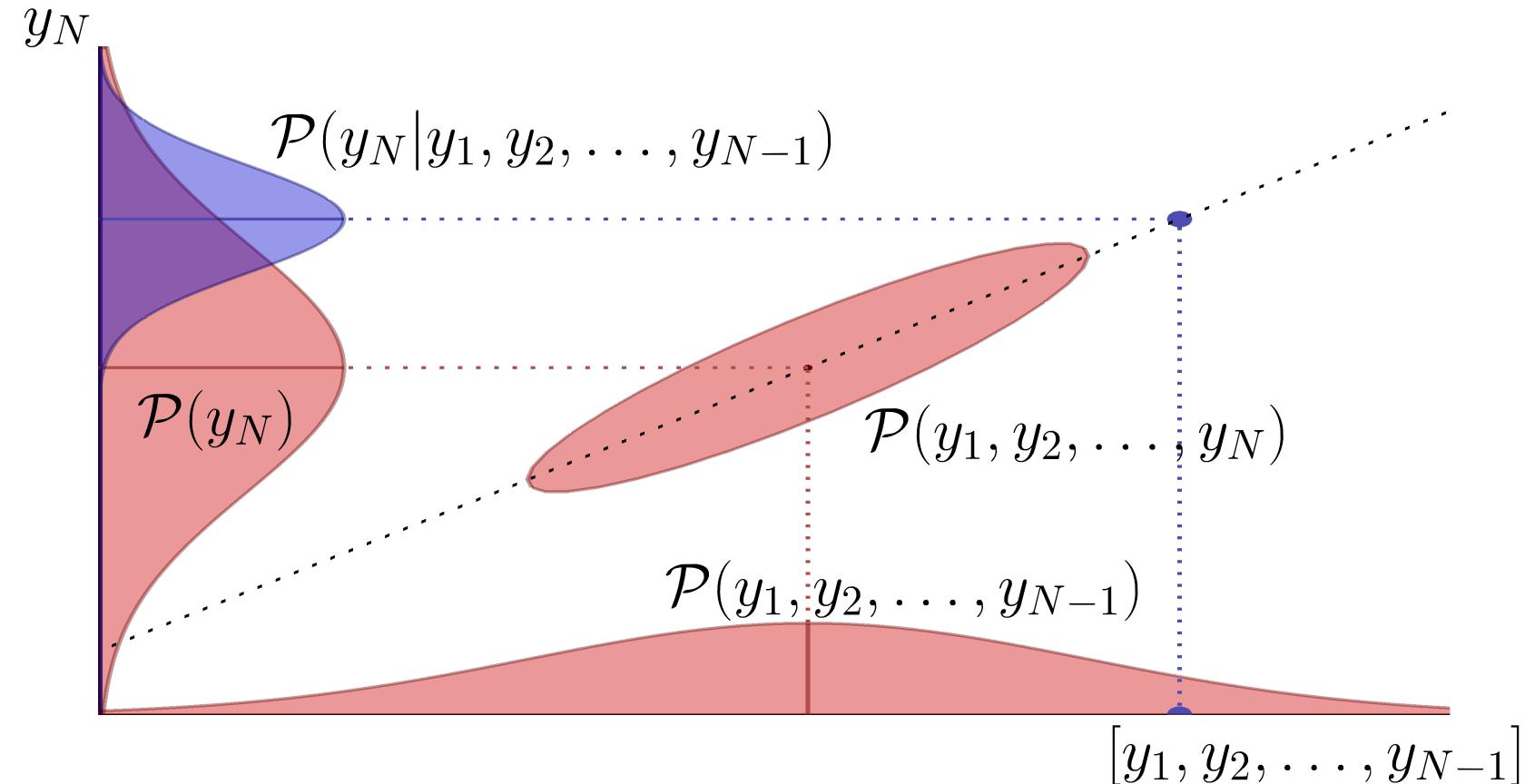
- A joint distribution represented by a bivariate Gaussian forms marginal distributions  $P(y_1)$  and  $P(y_2)$  that are unidimensional.
- Observing  $y_1$  changes our belief about  $y_2$ , giving rise to a **conditional distribution**.  
$$\hat{y}_2 = \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(y_1 - \mu_1)$$
- Knowledge of the covariance lets us shrink uncertainty in one variable based on the observation of the other.



# Gaussian process - Informal interpretation

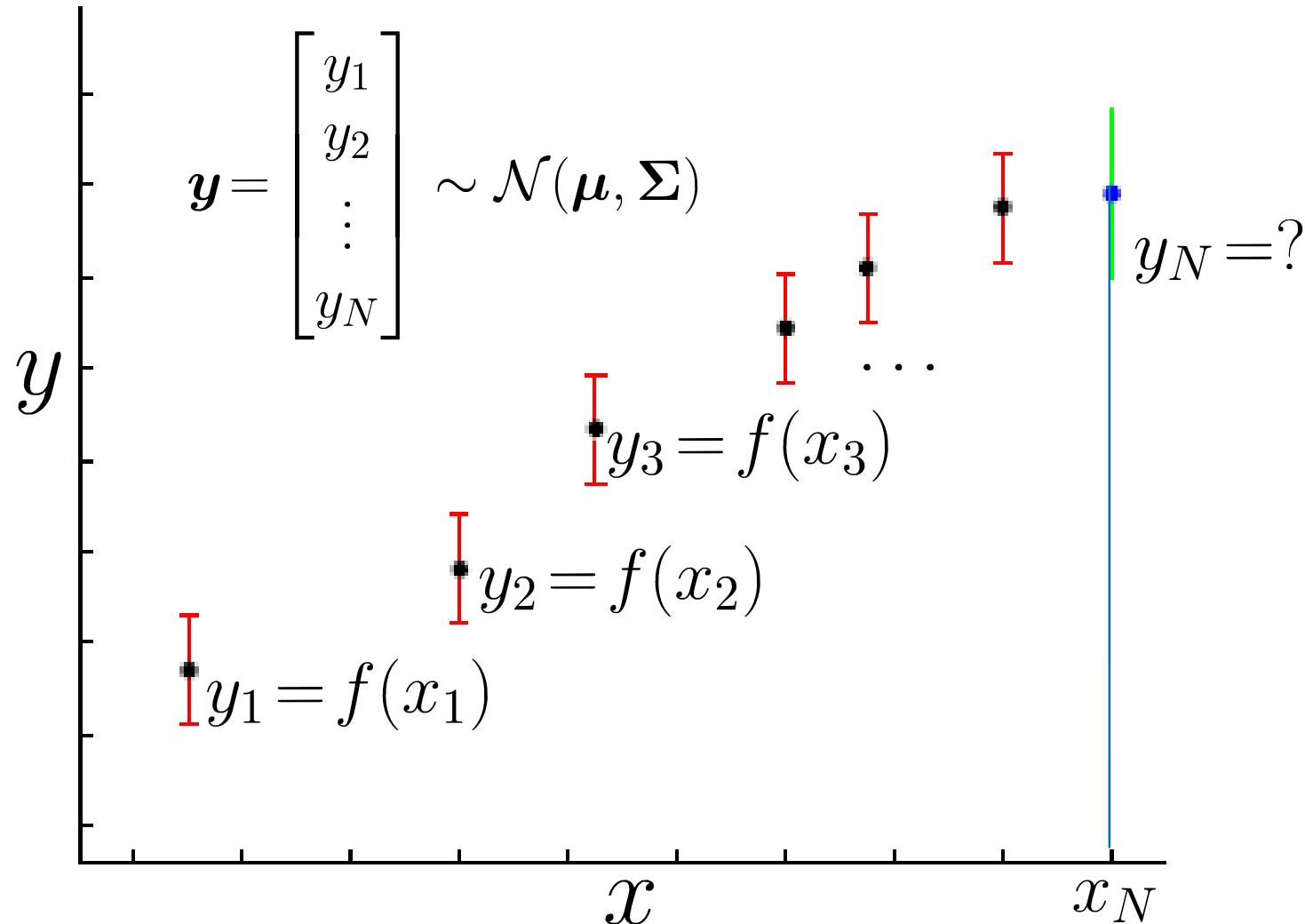
- This bivariate example can be extended to an arbitrarily large number of variables.
- Indeed, observations in an arbitrary dataset can always be imagined as a single point sampled from a multivariate Gaussian distribution.

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

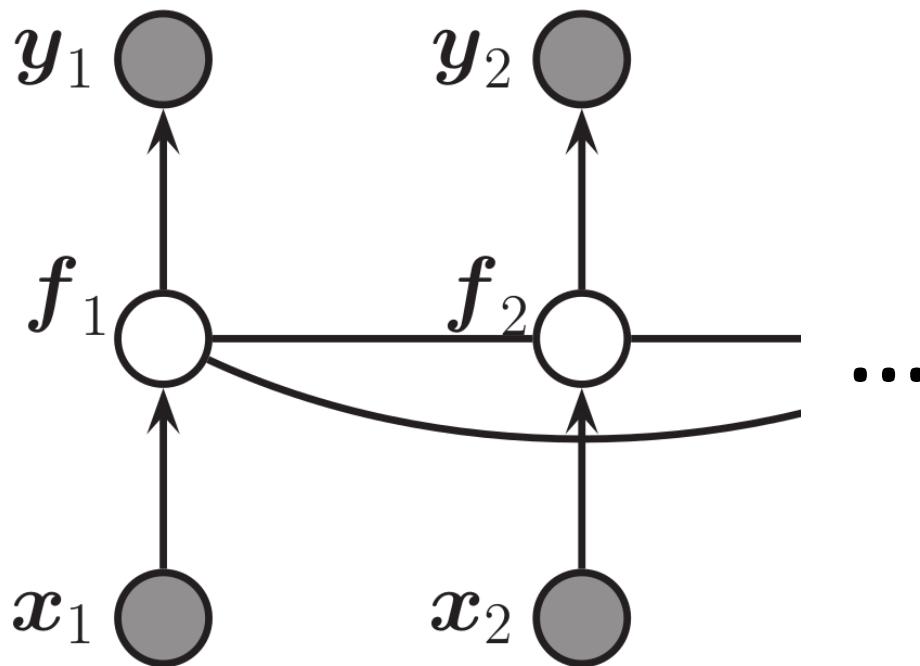


# How to construct this joint distribution in GP?

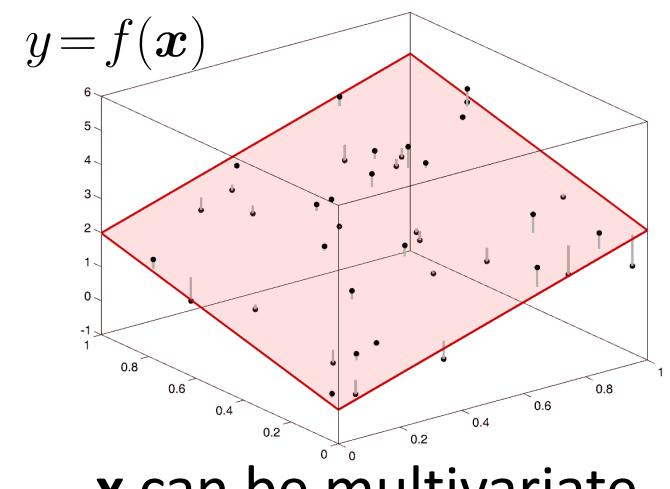
By looking at the similarities in the continuous  $\mathbf{x}$  space, representing the locations at which we evaluate  $y = f(\mathbf{x})$



# Graphical model of a Gaussian process



$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \sim \mathcal{N}(\mu, \Sigma)$$



Note that with GPs, we do not build a distribution on  $\{x_1, x_2, \dots, x_N\}$ !

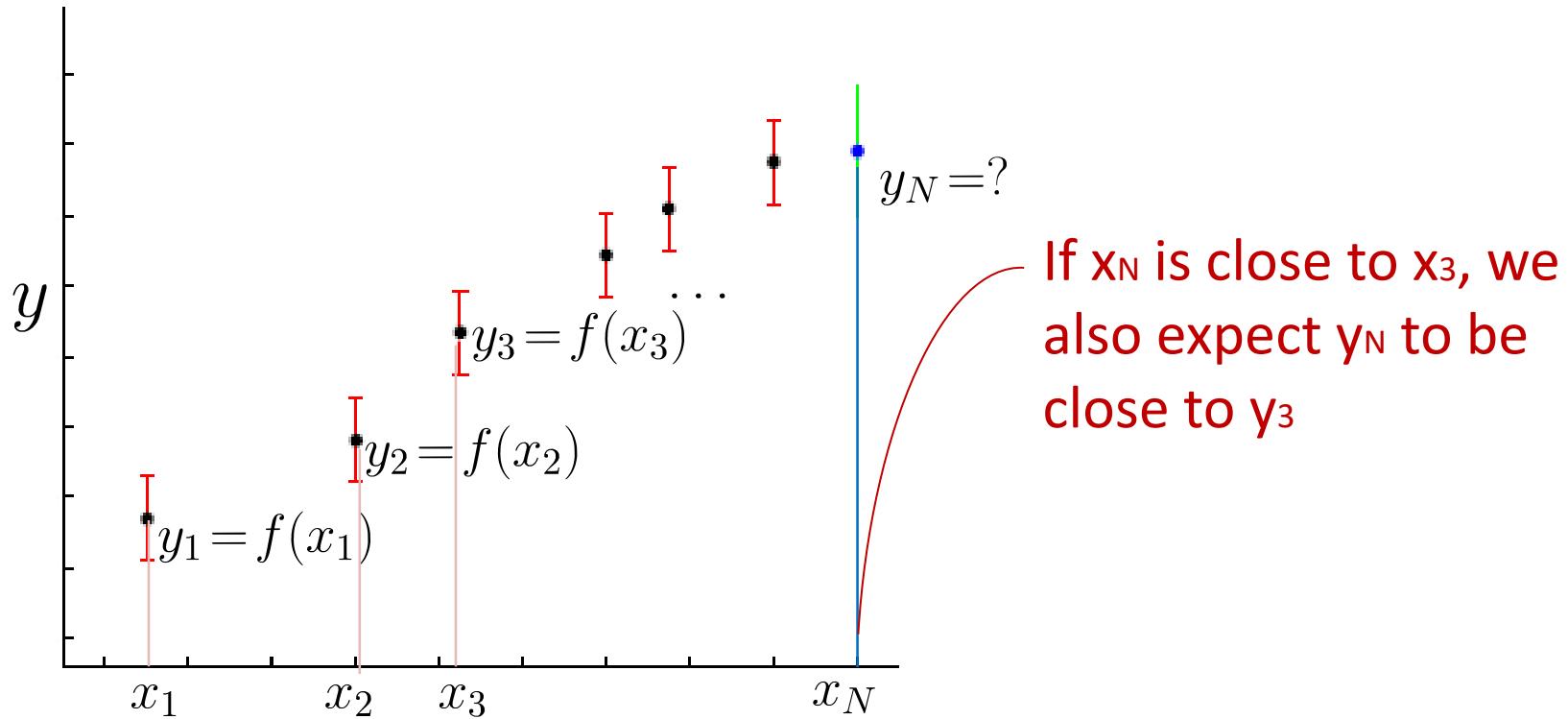
# Gaussian process (GP)

- Gaussian processes (GPs) can be seen as an infinite-dimensional generalization of multivariate normal distributions.
- The **infinite joint distribution** over all possible variables is equivalent to a **distribution over a function space**  $y = f(x)$ .
- $x$  can for be a vector or any object, but  $y$  is a scalar output.
- To understand GPs,  $N$  observations of an arbitrary data set  $y = \{y_1, \dots, y_N\}$  should be imagined as a single point sampled from an  $N$ -variate Gaussian.

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

# Gaussian process (GP)

- A covariance over an arbitrarily large set of variables can be defined through the **covariance kernel function**  $k(\mathbf{x}_i, \mathbf{x}_j)$ , providing the covariance elements between any two sample locations  $\mathbf{x}_i$  and  $\mathbf{x}_j$ .



## Distribution over functions in GPs

For a set of spatial or temporal locations  $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ , a positive semidefinite covariance matrix (also known as the Gram matrix) is defined as

$$\mathbf{K}(\mathbf{x}, \mathbf{x}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}.$$

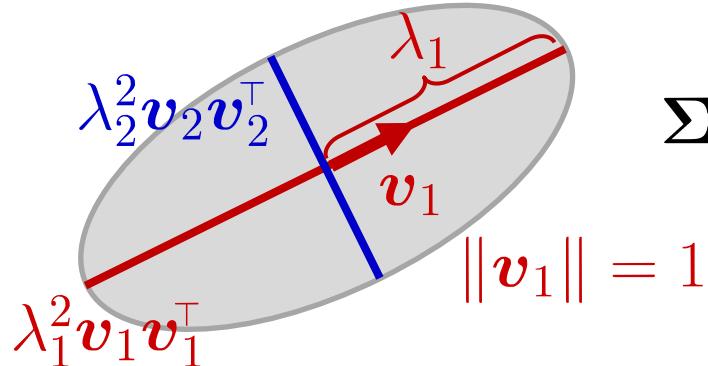
The entire function evaluation  $y_n = f(\mathbf{x}_n)$  associated with the set of inputs  $\mathbf{x}_n$  is a draw from a multivariate Gaussian distribution

$$\mathbf{y} \sim \mathcal{N}\left(\boldsymbol{\mu}(\mathbf{x}), \mathbf{K}(\mathbf{x}, \mathbf{x})\right),$$

specifying a **distribution over functions**.

# Stochastic sampling with Gaussians

The eigendecomposition of  $\Sigma$  is expressed in a matrix form as

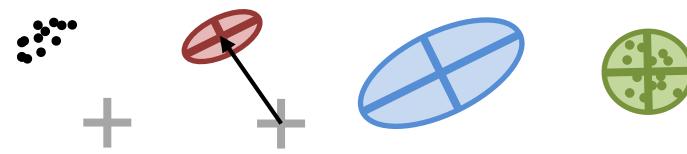


$$\Sigma = \mathbf{V} \mathbf{D} \mathbf{V}^\top = \sum_{j=1}^D \lambda_j^2 \mathbf{v}_j \mathbf{v}_j^\top$$

with  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_D]$

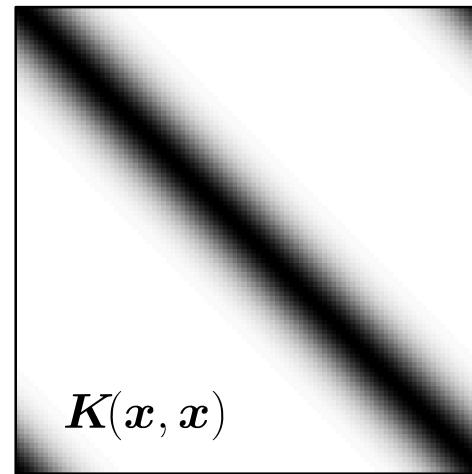
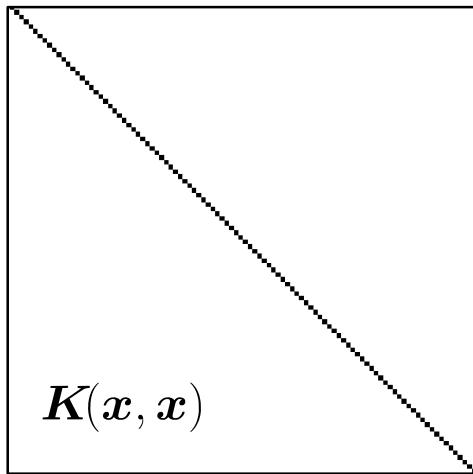
$$\mathbf{D} = \begin{bmatrix} \lambda_1^2 & 0 & \cdots & 0 \\ 0 & \lambda_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_D^2 \end{bmatrix}$$

By using this notation, datapoints can be stochastically generated with

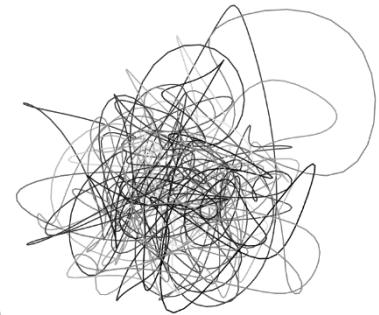
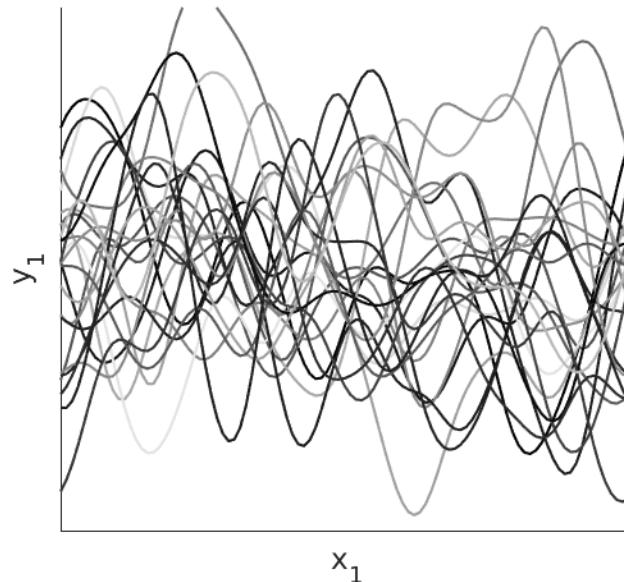
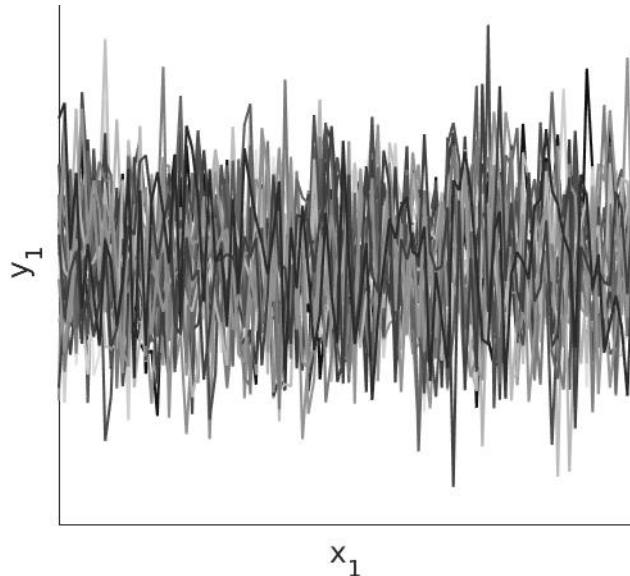


$$\xi \sim \mathcal{N}(\mu, \Sigma) \iff \xi \sim \mu + \mathbf{V} \mathbf{D}^{\frac{1}{2}} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

# Distribution over functions in GPs - Sampling



$$\mathbf{y} \sim \mathcal{N}\left(\mu(\mathbf{x}), \mathbf{K}(\mathbf{x}, \mathbf{x})\right)$$



# How to choose $k(\mathbf{x}_i, \mathbf{x}_j)$ ?

- We may know that our observations are samples from a process that is **smooth**, that is **continuous**, that has **typical amplitude**, or that the variations in the function take place within a **typical dynamic range**.
- These models require hyperparameters to be inferred, but **these hyperparameters define characteristics that are more generic** (such as the scale of a distribution) rather than acting explicitly on the structure or functional form of the signals.
- The notion of similarity will depend on the application: some of the basic aspects that can be defined through the covariance function  $k(\mathbf{x}_i, \mathbf{x}_j)$  are the process **stationarity, isotropy, smoothness or periodicity**.
- With continuous time series, past observations will be informative about current data as a function of how long ago they were observed.
- This corresponds to a **stationary** covariance, dependent on the Euclidean distance  $|\mathbf{x}_i - \mathbf{x}_j|$ .
- This process is also considered as **isotropic** if it does not depend on directions between  $\mathbf{x}_i$  and  $\mathbf{x}_j$ .
- A process that is both stationary and isotropic is **homogeneous**.

## $k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential covariance

A popular homogeneous covariance function is the **squared exponential kernel**, also known as **radial basis function**

$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp\left(-\frac{1}{\Theta_2^{\text{GP}}}(\mathbf{x}_i - \mathbf{x}_j)^\top(\mathbf{x}_i - \mathbf{x}_j)\right),$$

with two hyperparameters  $\Theta_1^{\text{GP}}$  and  $\Theta_2^{\text{GP}}$  corresponding respectively to **output and input scales** of the problem.

The radial basis function is widely employed when it is expected that nearby inputs  $\mathbf{x}_i$  and  $\mathbf{x}_j$  will have their corresponding outputs  $\mathbf{y}_i$  and  $\mathbf{y}_j$  also nearby (**assumption of continuity**).

# $k(x_i, x_j)$ as squared exponential covariance

$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp\left(-\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^{\top} \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j)\right)$$

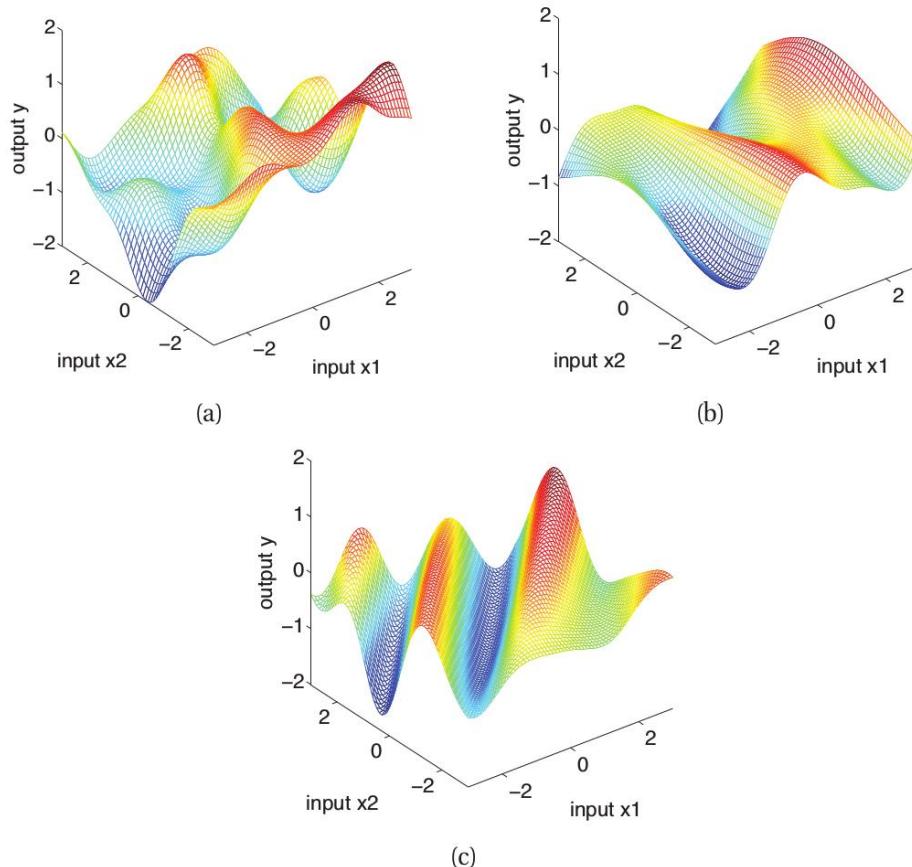
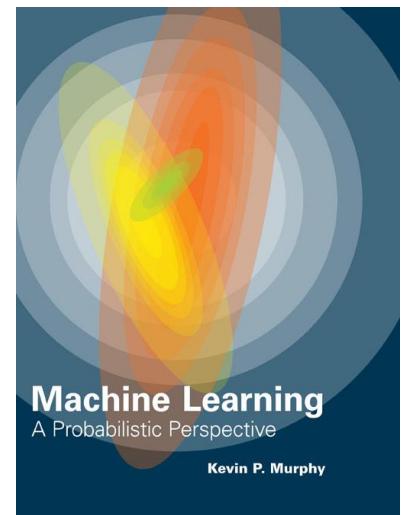


Figure from:



**Figure 15.4** Some 2d functions sampled from a GP with an SE kernel but different hyper-parameters. The kernel has the form in Equation 15.20 where (a)  $\mathbf{M} = \mathbf{I}$ , (b)  $\mathbf{M} = \text{diag}(1, 3)^{-2}$ , (c)  $\mathbf{M} = (1, -1; -1, 1) + \text{diag}(6, 6)^{-2}$ . Based on Figure 5.1 of (Rasmussen and Williams 2006). Figure generated by `gprDemoArd`, written by Carl Rasmussen.

## Modeling noise in the observed $\mathbf{y}_n$

If we assume there is noise associated with the observed function values  $\mathbf{y}_n = f(\mathbf{x}_n) + \boldsymbol{\eta}$ , this noise term can also be modeled in the covariance.

This noise is most often assumed to be uncorrelated from sample to sample, meaning that the noise term is only added to the diagonal elements of  $\mathbf{K}$ , giving a modified covariance for noisy observations of the form

$$\tilde{\mathbf{K}}(\mathbf{x}, \mathbf{x}) = \mathbf{K}(\mathbf{x}, \mathbf{x}) + \Theta_3^{\text{GP}} \mathbf{I}$$

where  $\mathbf{I}$  is the identity matrix and  $\Theta_3^{\text{GP}}$  is a Gaussian process hyperparameter representing the noise variance.

## Modeling noise in the observed $\mathbf{y}_n$

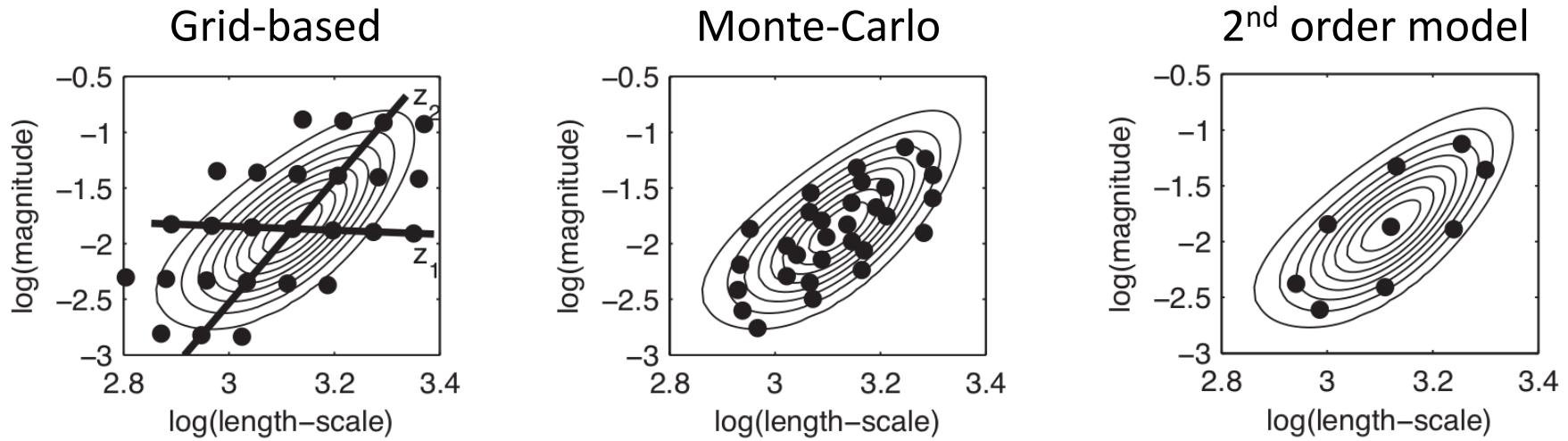
In the example of a covariance defined as squared exponential function, if noisy observations  $\mathbf{y}$  are assumed, the kernel can be directly defined as

$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp\left(-\frac{1}{\Theta_2^{\text{GP}}}(\mathbf{x}_i - \mathbf{x}_j)^\top(\mathbf{x}_i - \mathbf{x}_j)\right) + \Theta_3^{\text{GP}} \delta_{i,j},$$

where  $\delta_{i,j} = \mathbb{I}(i = j)$  is equal to one only when  $i = j$  and is zero otherwise, resulting in a covariance matrix  $\mathbf{K}(\mathbf{x}, \mathbf{x})$  with noise related to observations only present in the diagonal (noise uncorrelated from sample to sample).

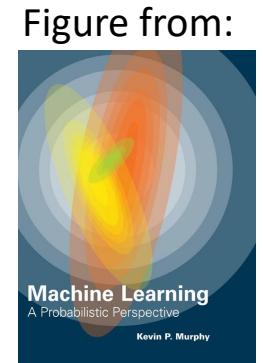
# Learning the kernel function parameters

Several approaches exist to estimate the hyperparameters of the covariance function: Maximum Likelihood Estimation (MLE), cross-validation (CV), Bayesian approaches involving sampling algorithms such as MCMC, etc.



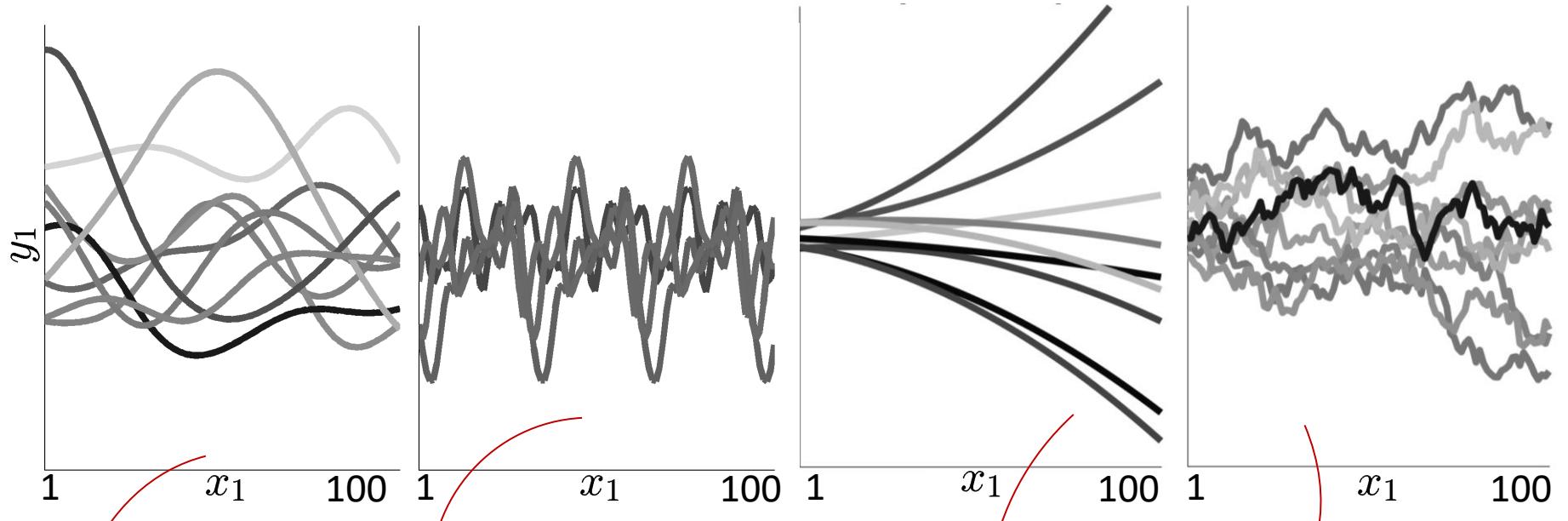
For example, given an expression for the log marginal likelihood and its derivative, we can estimate the kernel parameters using standard gradient-based optimizer.

Since the objective is not convex, local minima can still be a problem.



# Stochastic sampling from covariance matrix

$$\mathbf{y} \sim \mathcal{N}(0, K(\mathbf{x}, \mathbf{x})) \quad \mathbf{x} = [1, 2, \dots, 100]^\top$$



$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp\left(-\frac{1}{\Theta_2^{\text{GP}}} \sin^2(\Theta_4 \|\mathbf{x}_i - \mathbf{x}_j\|)\right)$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp\left(-\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j)\right)$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + \Theta_1^{\text{GP}})^2$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \min(\mathbf{x}_i, \mathbf{x}_j) + \Theta_1^{\text{GP}}$$

# Gaussian process regression (GPR)

*a.k.a.*

## Kriging

**Python notebook:**  
**demo\_GPR.ipynb**

**Matlab code:**  
**demo\_GPR01.m**

# Gaussian process regression (GPR)

We are interested in the **posterior distribution** of  $\mathbf{y}^*$  to be computed at some location(s)  $\mathbf{x}^*$ .

The **joint distribution** of the already observed  $\mathbf{y}$  (at location  $\mathbf{x}$ ) augmented by  $\mathbf{y}^*$  (at location  $\mathbf{x}^*$ ) is

$$\mathbf{K}(\mathbf{x}, \mathbf{x}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{y}^* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}(\mathbf{x}) \\ \boldsymbol{\mu}(\mathbf{x}^*) \end{bmatrix}, \begin{bmatrix} \mathbf{K}(\mathbf{x}, \mathbf{x}) & \mathbf{K}(\mathbf{x}, \mathbf{x}^*) \\ \mathbf{K}(\mathbf{x}^*, \mathbf{x}) & \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*) \end{bmatrix} \right)$$

We can use the conditional probability property of Gaussians to evaluate the posterior distribution over  $\mathbf{y}^*$ , yielding a Gaussian

$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

$$\text{with } \boldsymbol{\mu}^* = \boldsymbol{\mu}(\mathbf{x}^*) + \mathbf{K}(\mathbf{x}^*, \mathbf{x}) \mathbf{K}(\mathbf{x}, \mathbf{x})^{-1} (\mathbf{y} - \boldsymbol{\mu}(\mathbf{x}))$$

$$\boldsymbol{\Sigma}^* = \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{K}(\mathbf{x}^*, \mathbf{x}) \mathbf{K}(\mathbf{x}, \mathbf{x})^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x}^*)$$

# GPR in practice

$$\mathbf{K}(\mathbf{x}, \mathbf{x}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

It is also often assumed in practice that  $\begin{bmatrix} \boldsymbol{\mu}(\mathbf{x}) \\ \boldsymbol{\mu}(\mathbf{x}^*) \end{bmatrix} = \mathbf{0}$ .

In this case, Gaussian processes can be completely defined by second-order statistics, where the Gram matrix  $\mathbf{K}$  is a positive semi-definite covariance.

Note that  $\mathbf{K}(\mathbf{x}, \mathbf{x})^{-1}$  can be pre-computed so that the posterior distribution can be computed faster

$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

$$\text{with } \boldsymbol{\mu}^* = \boldsymbol{\mu}(\mathbf{x}^*) + \mathbf{K}(\mathbf{x}^*, \mathbf{x}) \mathbf{K}(\mathbf{x}, \mathbf{x})^{-1} (\mathbf{y} - \boldsymbol{\mu}(\mathbf{x}))$$

$$\boldsymbol{\Sigma}^* = \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{K}(\mathbf{x}^*, \mathbf{x}) \mathbf{K}(\mathbf{x}, \mathbf{x})^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x}^*)$$

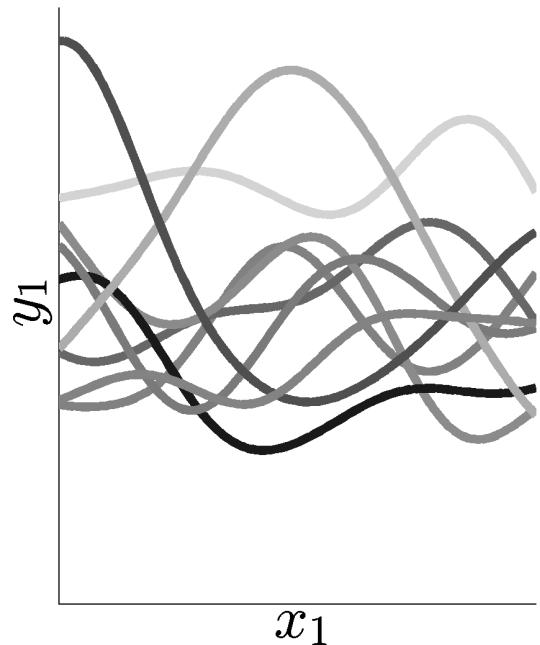
# **Examples of covariance functions**

# $k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential covariance

$$\Theta_1^{\text{GP}} = 1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0$$

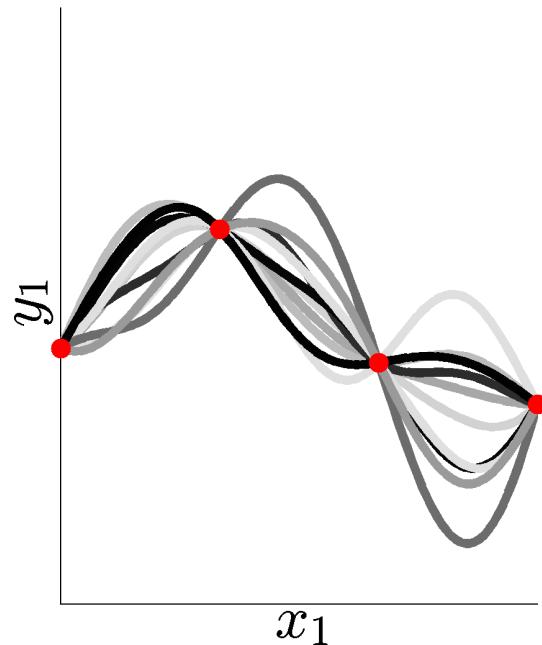
$$\mathbf{y}^* \sim \mathcal{N}(\mu(x^*), \mathbf{K}(x^*, x^*))$$

Samples from prior



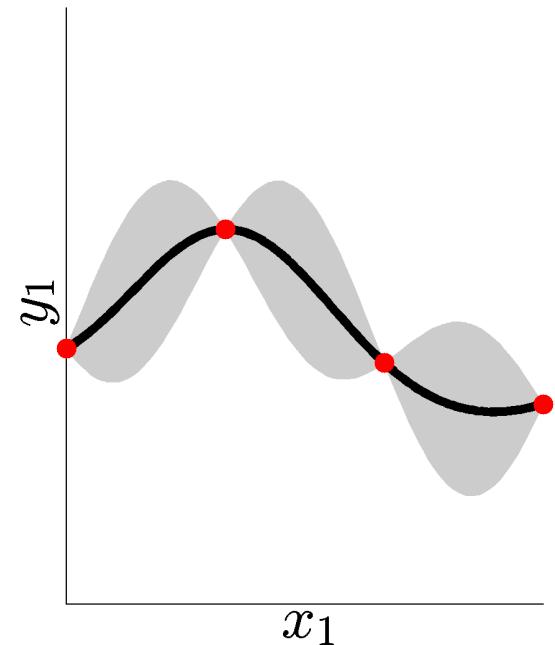
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\mu^*, \Sigma^*)$$

Samples from posterior



$$\mathcal{N}(\mu^*, \Sigma^*)$$

Trajectory distribution



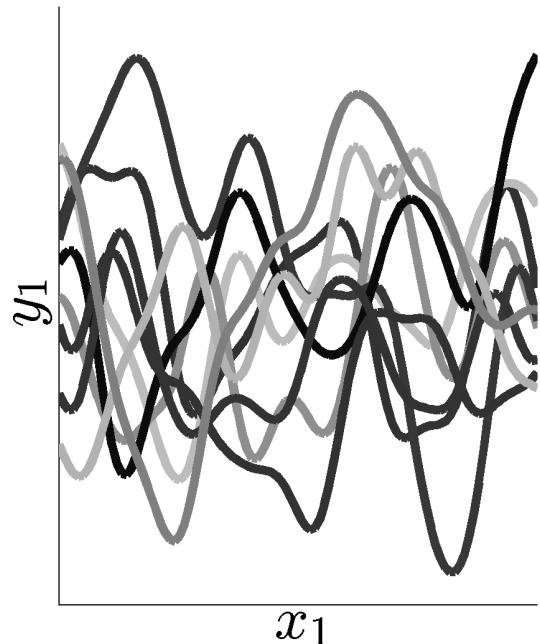
$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp \left( -\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) \right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

# $k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential covariance

$$\Theta_1^{\text{GP}} = 1, \quad \Theta_2^{\text{GP}} = 0.01, \quad \Theta_3^{\text{GP}} = 0$$

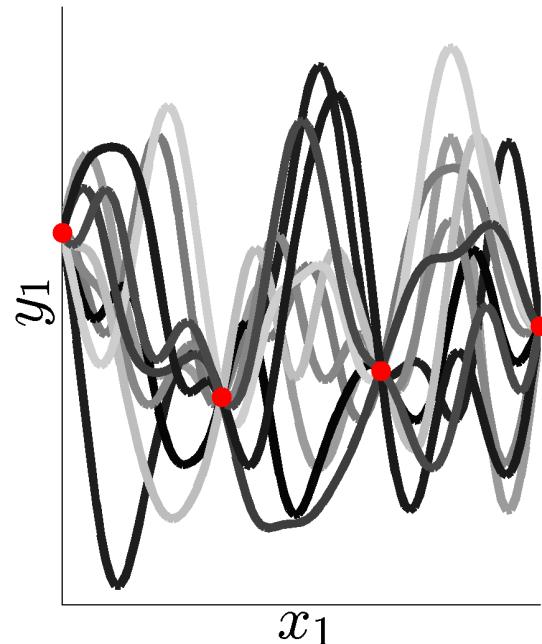
$$\mathbf{y}^* \sim \mathcal{N}(\mu(x^*), \mathbf{K}(x^*, x^*))$$

Samples from prior



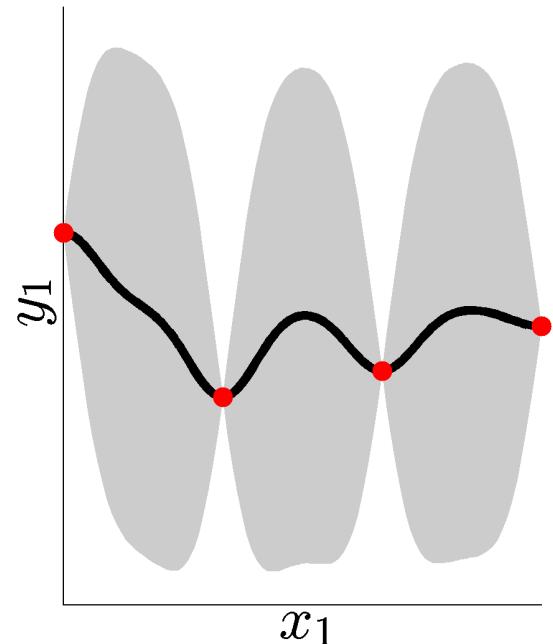
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\mu^*, \Sigma^*)$$

Samples from posterior



$$\mathcal{N}(\mu^*, \Sigma^*)$$

Trajectory distribution



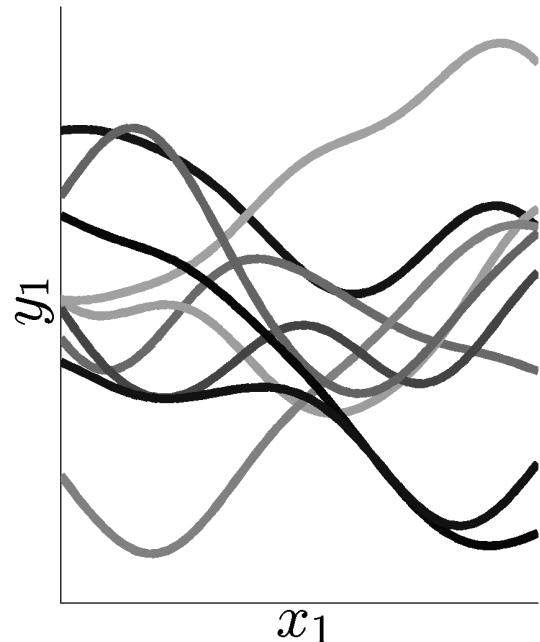
$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp \left( -\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) \right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

# $k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential covariance

$$\Theta_1^{\text{GP}} = 1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0.01$$

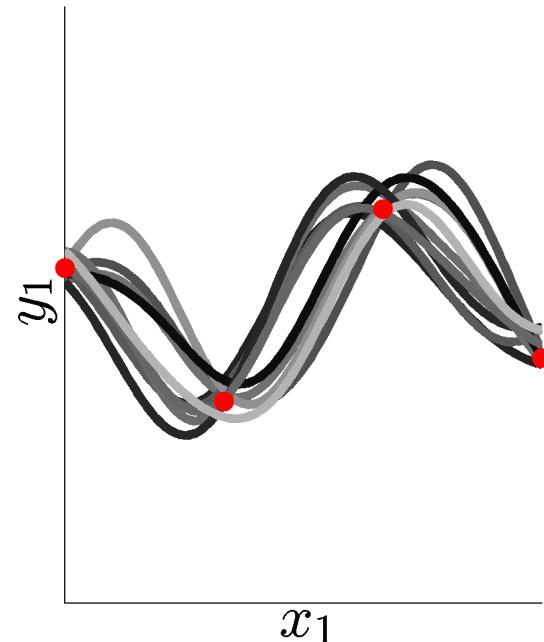
$$\mathbf{y}^* \sim \mathcal{N}(\mu(x^*), \mathbf{K}(x^*, x^*))$$

Samples from prior



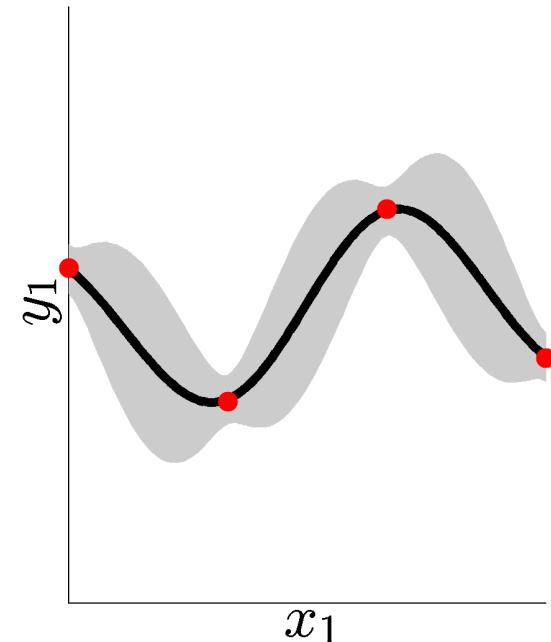
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\mu^*, \Sigma^*)$$

Samples from posterior



$$\mathcal{N}(\mu^*, \Sigma^*)$$

Trajectory distribution



$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp \left( -\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) \right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

# $k(\mathbf{x}_i, \mathbf{x}_j)$ as squared exponential covariance

$$\Theta_1^{\text{GP}} = 1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0.01$$

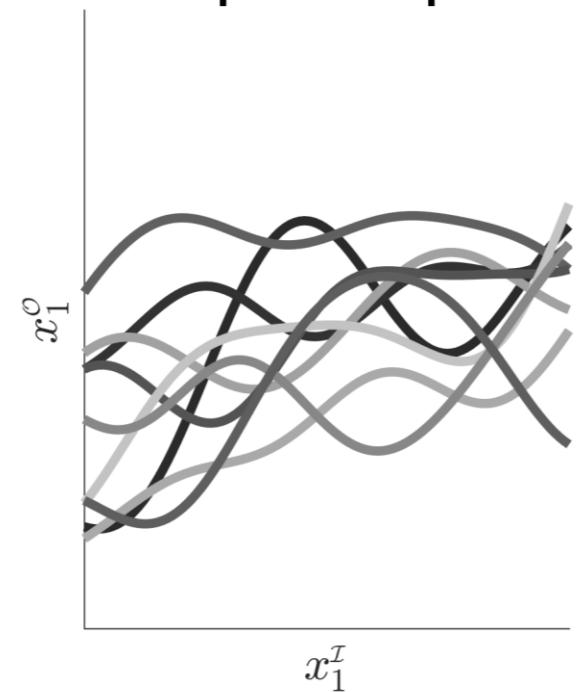
$$\mu(\mathbf{x}) = \alpha \mathbf{x}$$

$$\mathbf{y}^* \sim \mathcal{N}(\mu(x^*), \mathbf{K}(x^*, x^*))$$

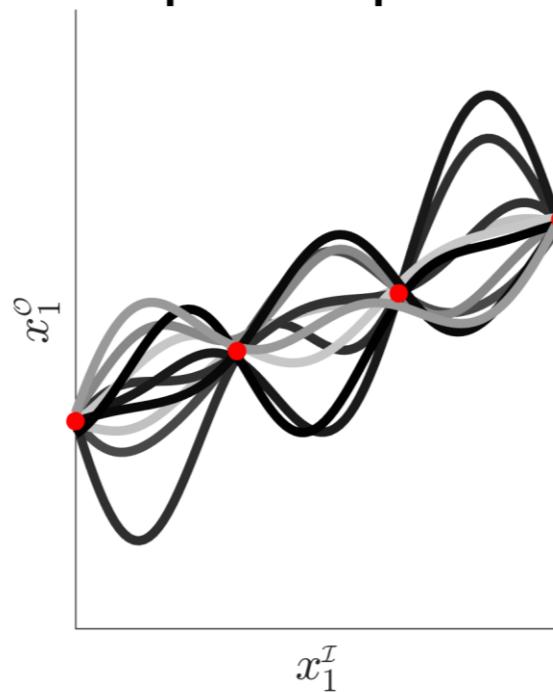
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\mu^*, \Sigma^*)$$

$$\mathcal{N}(\mu^*, \Sigma^*)$$

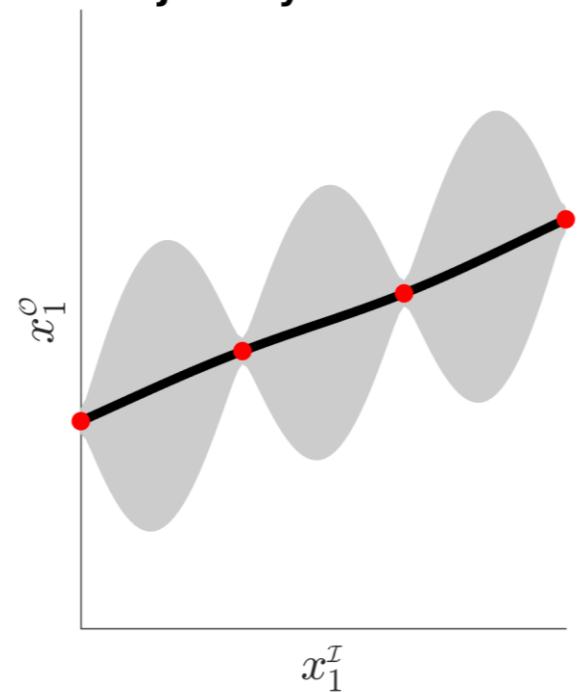
Samples from prior



Samples from posterior



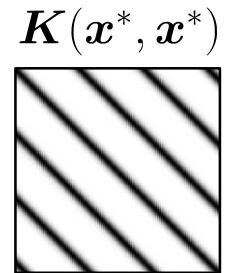
Trajectory distribution



$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp \left( -\frac{1}{\Theta_2^{\text{GP}}} (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) \right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

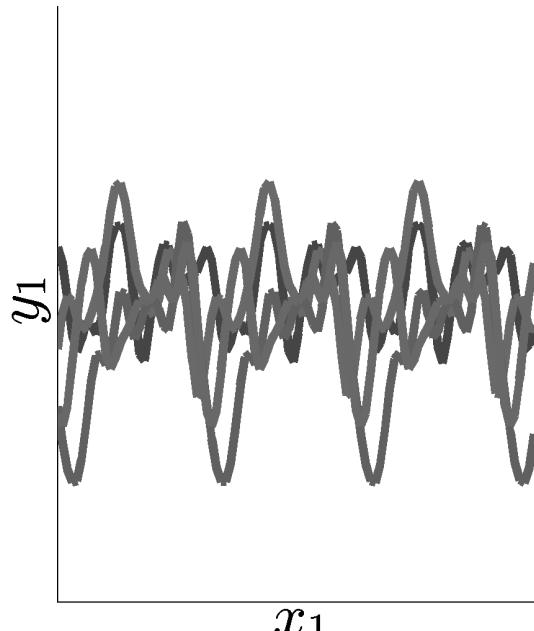
# $k(\mathbf{x}_i, \mathbf{x}_j)$ as **periodic** covariance function

$$\Theta_1^{\text{GP}} = 0.1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0, \quad \Theta_4^{\text{GP}} = 10$$



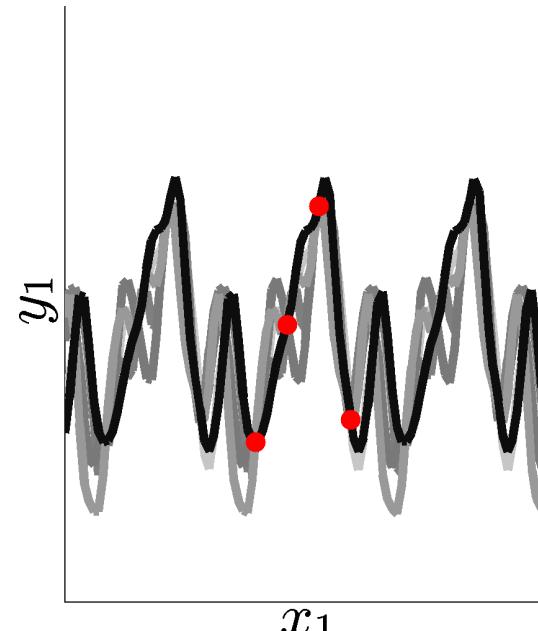
$$\mathbf{y}^* \sim \mathcal{N}(\mu(x^*), \mathbf{K}(x^*, x^*))$$

Samples from prior



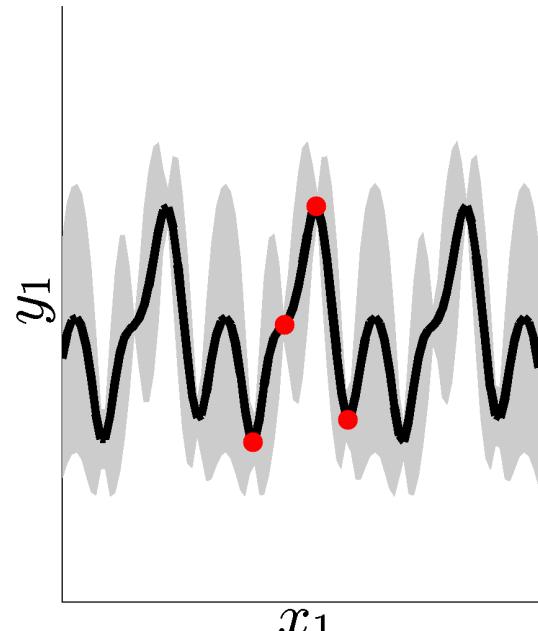
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\mu^*, \Sigma^*)$$

Samples from posterior



$$\mathcal{N}(\mu^*, \Sigma^*)$$

Trajectory distribution



$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \exp \left( -\frac{1}{\Theta_2^{\text{GP}}} \sin^2(\Theta_4^{\text{GP}} |\mathbf{x}_i - \mathbf{x}_j|) \right) + \Theta_3^{\text{GP}} \delta_{i,j}$$

## $k(x_i, x_j)$ as Matérn covariance function

Another popular covariance kernel function is the Matérn function

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \sqrt{2\nu} \frac{d}{\rho} \right)^\nu K_\nu \left( \sqrt{2\nu} \frac{d}{\rho} \right)$$

$$\text{with } d = \|\mathbf{x}_i - \mathbf{x}_j\|$$

where  $\Gamma$  is the gamma function,  $K_\nu$  is the modified Bessel function of the second kind, and  $\rho$  and  $\nu$  are non-negative parameters of the covariance.

A Gaussian process with Matérn covariance has sample paths that are  $\lfloor \nu - 1 \rfloor$  times differentiable.

## **k(x<sub>i</sub>,x<sub>j</sub>) as Matérn covariance function**

### Simplification for $\nu$ half integer

When  $\nu = p + 1/2$ ,  $p \in \mathbb{N}^+$ , the Matérn covariance can be written as a product of an exponential and a polynomial of order  $p$ :

$$C_{p+1/2}(d) = \sigma^2 \exp\left(-\frac{\sqrt{2\nu}d}{\rho}\right) \frac{\Gamma(p+1)}{\Gamma(2p+1)} \sum_{i=0}^p \frac{(p+i)!}{i!(p-i)!} \left(\frac{\sqrt{8\nu}d}{\rho}\right)^{p-i}$$

$$\text{For } \nu = 1/2 \text{ (} p = 0 \text{): } C_{1/2}(d) = \sigma^2 \exp\left(-\frac{d}{\rho}\right)$$

$$\text{For } \nu = 3/2 \text{ (} p = 1 \text{): } C_{3/2}(d) = \sigma^2 \left(1 + \frac{\sqrt{3}d}{\rho}\right) \exp\left(-\frac{\sqrt{3}d}{\rho}\right)$$

$$\text{For } \nu = 5/2 \text{ (} p = 2 \text{): } C_{5/2}(d) = \sigma^2 \left(1 + \frac{\sqrt{5}d}{\rho} + \frac{5d^2}{3\rho^2}\right) \exp\left(-\frac{\sqrt{5}d}{\rho}\right)$$

As  $\nu \rightarrow \infty$ , the Matérn covariance converges to the squared exponential covariance function.

# $k(\mathbf{x}_i, \mathbf{x}_j)$ as Matern covariance function

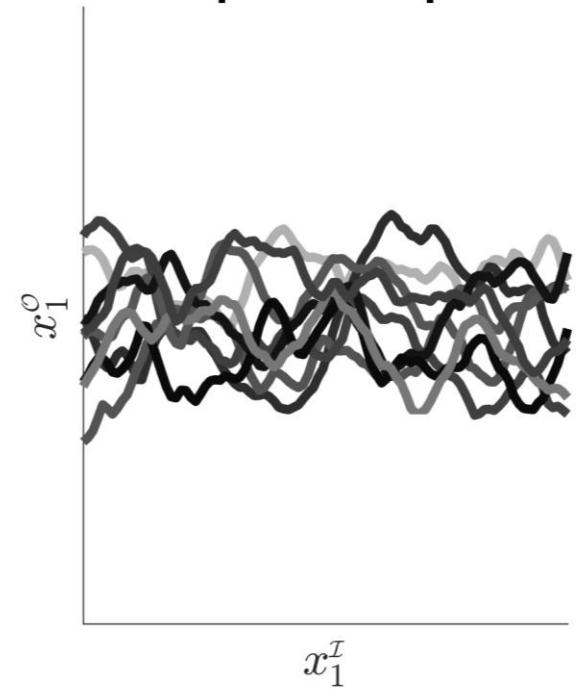
$$\Theta_1^{\text{GP}} = 0.1, \quad \Theta_2^{\text{GP}} = 0.1, \quad \Theta_3^{\text{GP}} = 0.0001$$

$$\mathbf{y}^* \sim \mathcal{N}(\mu(\mathbf{x}^*), \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*))$$

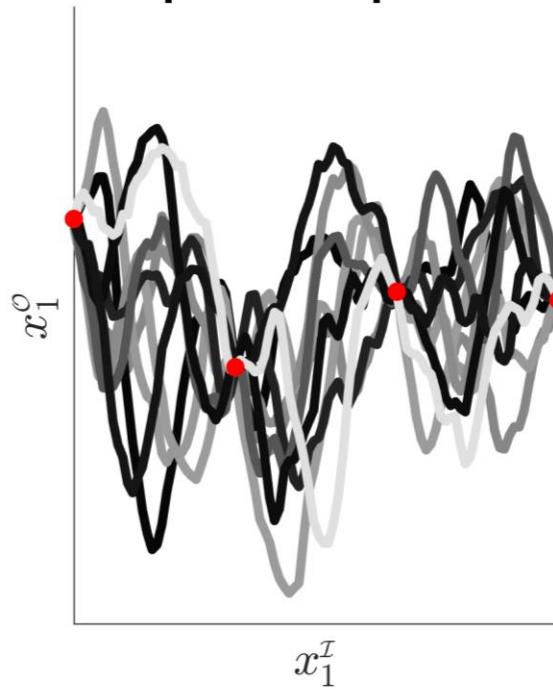
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\mu^*, \Sigma^*)$$

$$\mathcal{N}(\mu^*, \Sigma^*)$$

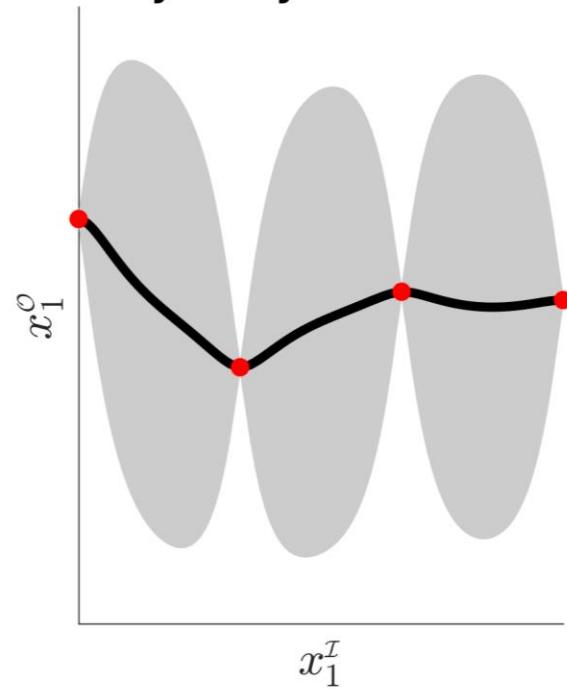
Samples from prior



Samples from posterior



Trajectory distribution



$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} \left( 1 + \frac{\sqrt{3} \|\mathbf{x}_i - \mathbf{x}_j\|}{\Theta_2^{\text{GP}}} \right) \exp \left( -\frac{\sqrt{3} \|\mathbf{x}_i - \mathbf{x}_j\|}{\Theta_2^{\text{GP}}} \right)$$

# $k(\mathbf{x}_i, \mathbf{x}_j)$ as Brownian motion covariance function

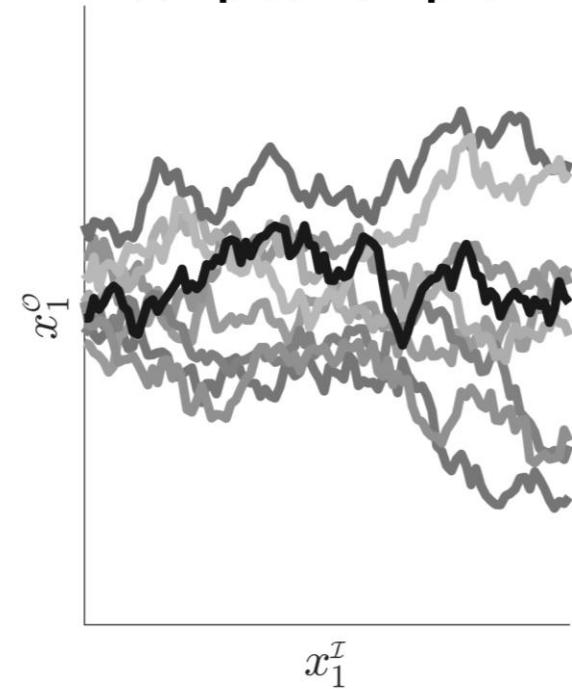
The **Wiener process** is a simple continuous-time stochastic process often put in connection to the Brownian motion.

$$\mathbf{y}^* \sim \mathcal{N}(\mu(x^*), \mathbf{K}(x^*, x^*))$$

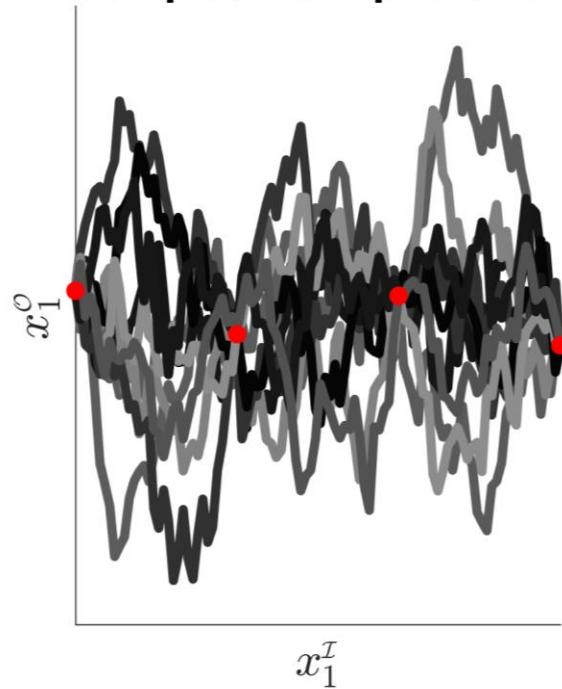
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\mu^*, \Sigma^*)$$

$$\mathcal{N}(\mu^*, \Sigma^*)$$

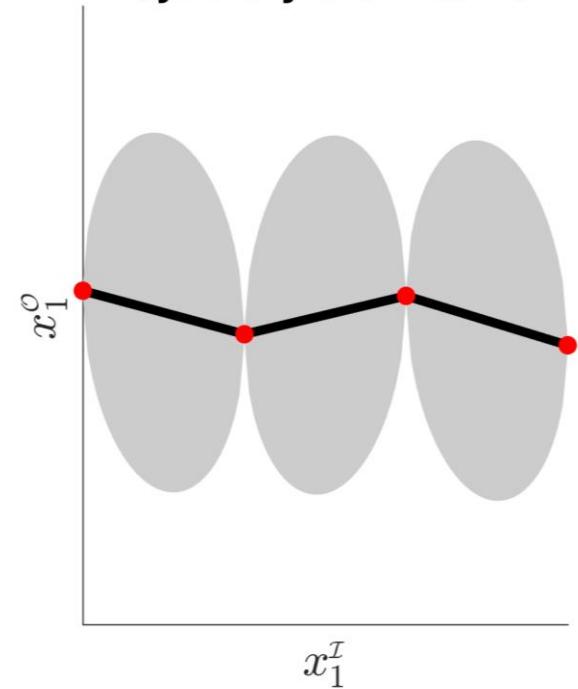
Samples from prior



Samples from posterior



Trajectory distribution

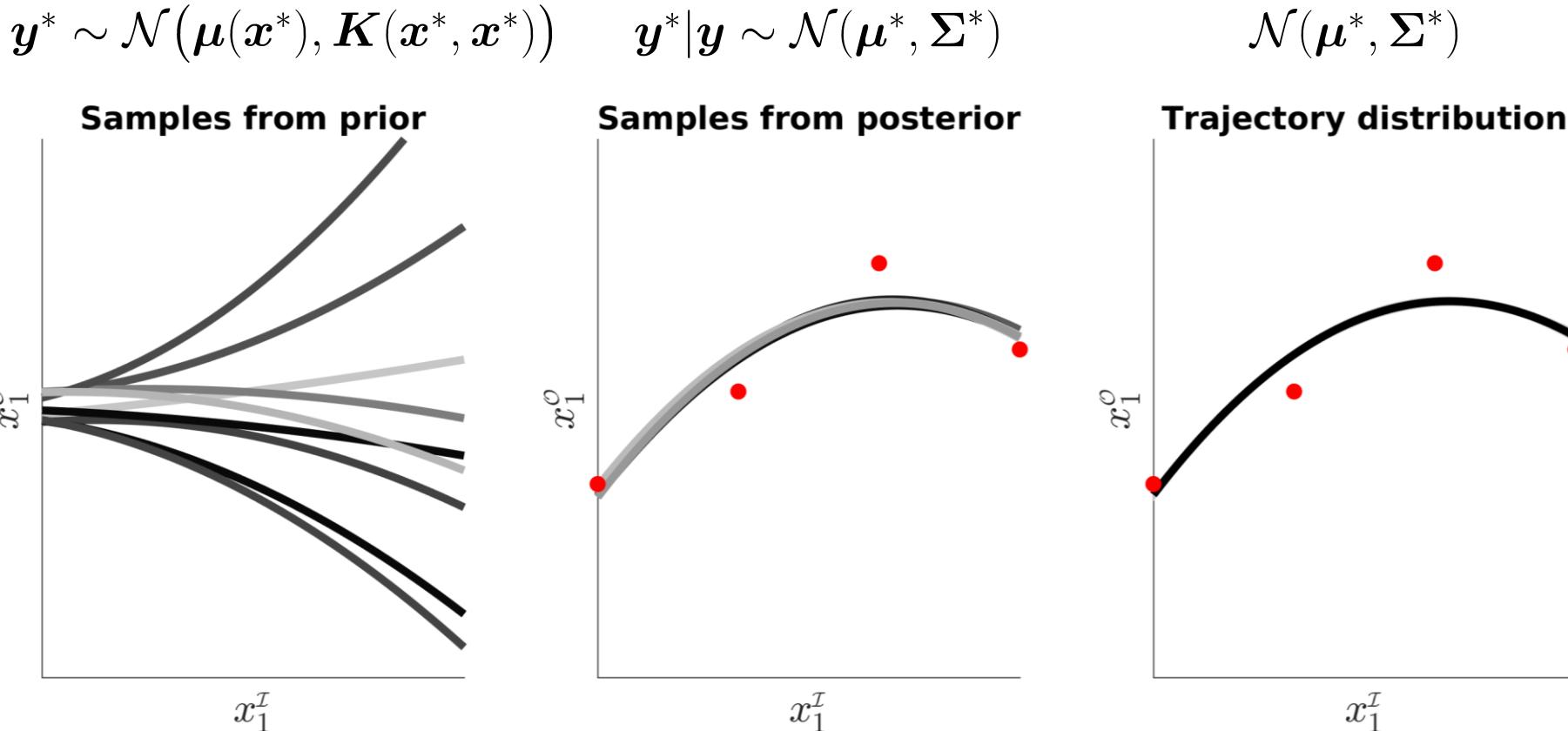


$$k(\mathbf{x}_i, \mathbf{x}_j) = \min(\mathbf{x}_i, \mathbf{x}_j) + \Theta_1^{\text{GP}}$$

$$\Theta_1^{\text{GP}} = 0.1$$

# $k(\mathbf{x}_i, \mathbf{x}_j)$ as quadratic covariance function

Bayesian linear regression is equivalent to a GP with covariance function  $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$ .



$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + \Theta_1^{\text{GP}})^2 \quad \Theta_1^{\text{GP}} = 0.1$$

# $k(\mathbf{x}_i, \mathbf{x}_j)$ as **polynomial** covariance function

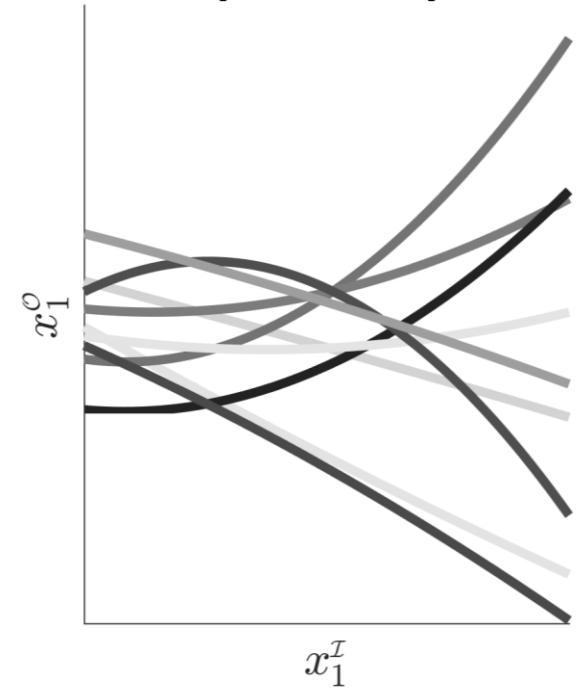
$$\Theta_1^{\text{GP}} = 0.1$$

$$\mathbf{y}^* \sim \mathcal{N}(\mu(x^*), \mathbf{K}(x^*, x^*))$$

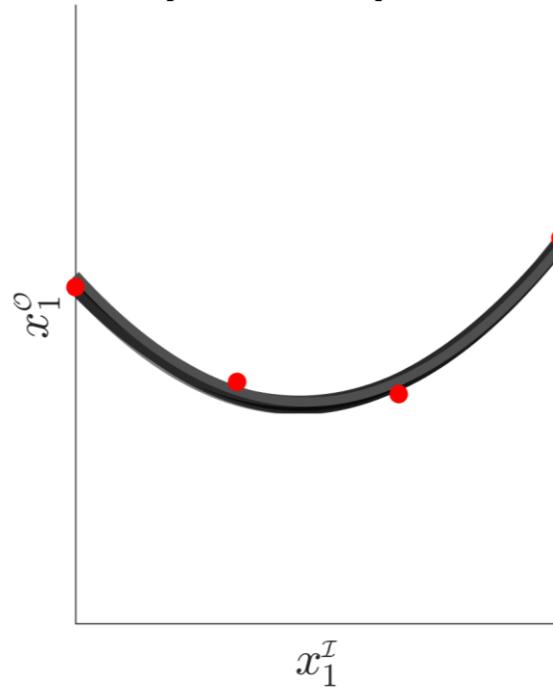
$$\mathbf{y}^* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

$$\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

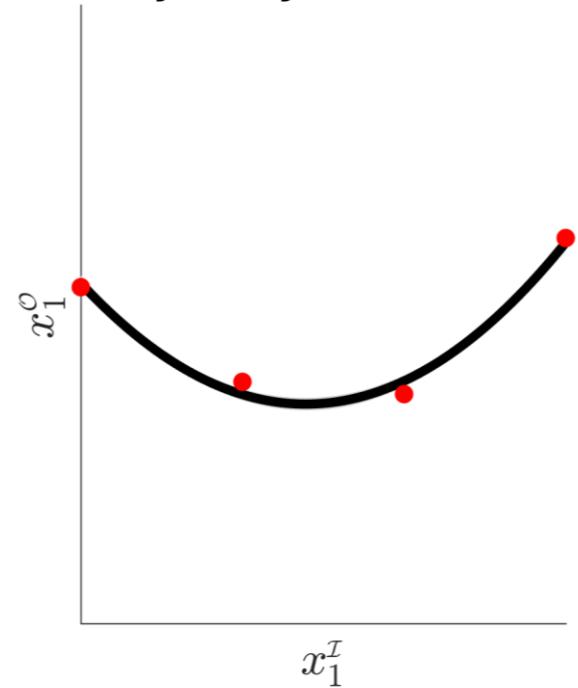
Samples from prior



Samples from posterior



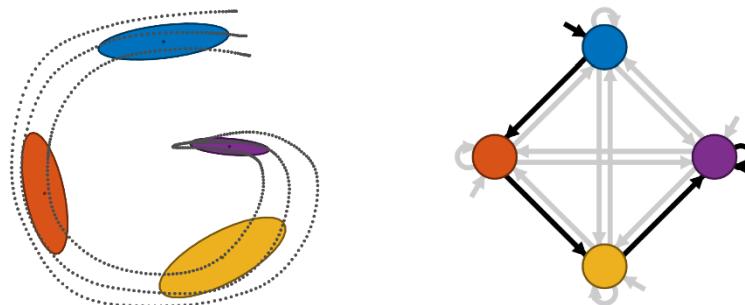
Trajectory distribution



$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j)^2 + \mathbf{x}_i^\top \mathbf{x}_j + \Theta_1^{\text{GP}}$$

## $k(x_i, x_j)$ as **probabilistic model covariance**

- Another powerful approach to the construction of kernels is to exploit probabilistic models.
- Given a generative model  $P(\mathbf{x})$ , a valid kernel can be defined as  $k(x_i, x_j) = P(x_i) P(x_j)$ , which can be interpreted as an inner product in the one-dimensional feature space defined by the mapping  $P(\mathbf{x})$ .
- Namely, two inputs  $x_i$  and  $x_j$  will be similar if they both have high probabilities to belong to the model.
- This can bring additional properties to the underlying process such as the capability of handling missing data or partial sequences of various lengths (e.g., with HMM).



# $k(\mathbf{x}_i, \mathbf{x}_j)$ as **weighted sum of kernel functions**

- A covariance function can be defined as a **linear combination of other covariance functions**, which can be exploited to incorporate different insights about the dataset.
- Such an approach can be exploited as an alternative to optimizing kernel parameters (also known as multiple kernel learning).
- The idea is to define the kernel as a **weighted sum of basis kernels**, and then to **optimize the weights instead of the kernel parameters**.

Dictionary of basis kernel functions

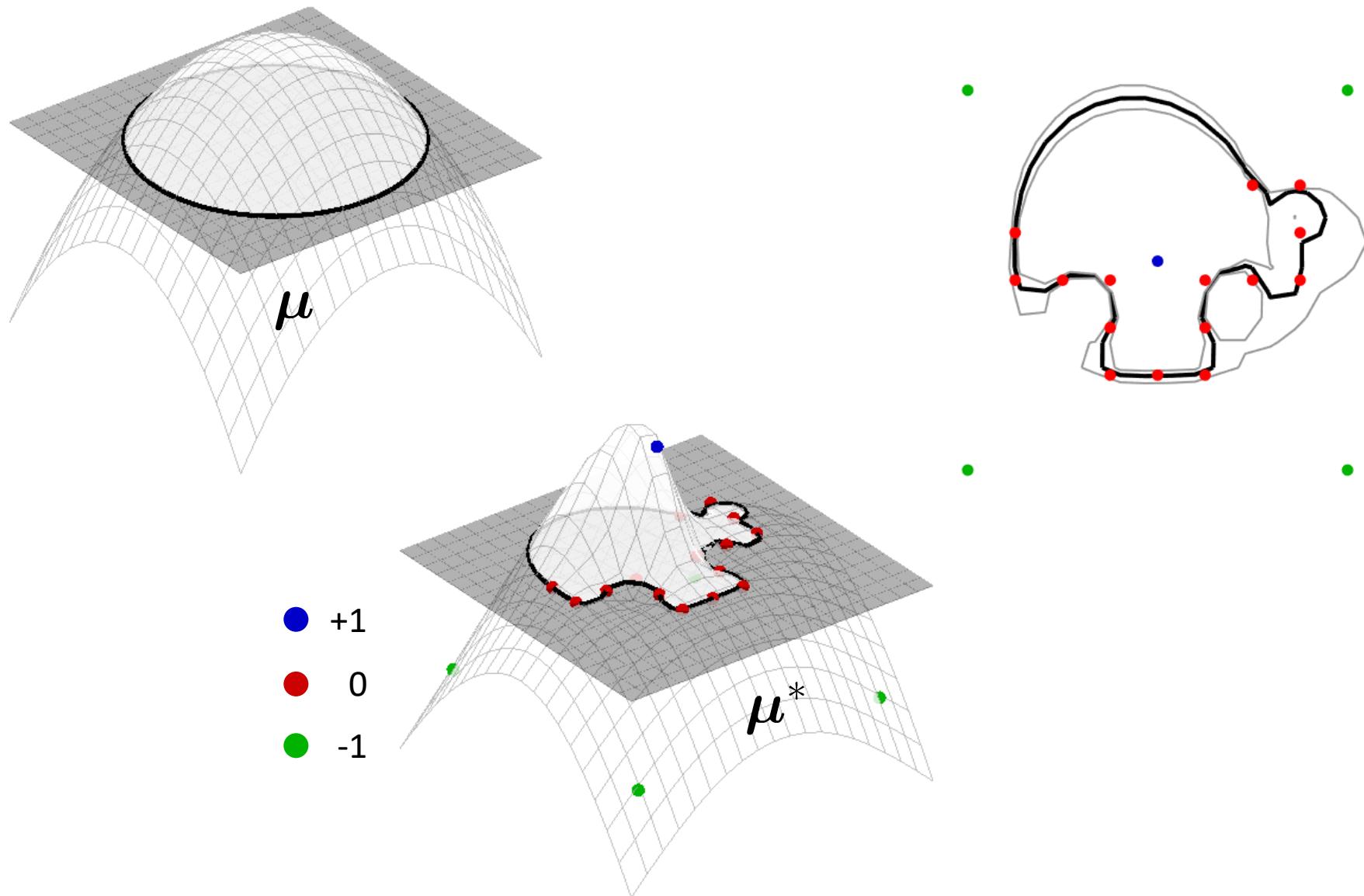
$$\left\{ \begin{array}{l} k_1(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j)^2 + \mathbf{x}_i^\top \mathbf{x}_j \\ k_2(\mathbf{x}_i, \mathbf{x}_j) = \min(\mathbf{x}_i, \mathbf{x}_j) \\ k_3(\mathbf{x}_i, \mathbf{x}_j) = \exp \left( - (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) \right) \end{array} \right.$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \Theta_1^{\text{GP}} k_1(\mathbf{x}_i, \mathbf{x}_j) + \Theta_2^{\text{GP}} k_2(\mathbf{x}_i, \mathbf{x}_j) + \Theta_3^{\text{GP}} k_3(\mathbf{x}_i, \mathbf{x}_j)$$

# Some extensions of Gaussian processes

- **Cokriging:**  
Extending GPR to multiple target variables  $y$ .
- **Sparse GP:**  
A known bottleneck in Gaussian process prediction is that the computational complexity of prediction is  $O(N^3)$   
→ not feasible for large data sets!  
Sparse Gaussian processes circumvent this issue by building a representative set for the given process  $y = f(x)$ .
- **Gaussian process latent variable models (GPLVM):**  
GPLVM is a probabilistic dimensionality reduction method that uses GPs to find a lower dimensional non-linear embedding of high dimensional data.

# Gaussian Process Implicit Surface (GPIS)



# References

## GPR

- C.K.I. Williams and C.E. Rasmussen. Gaussian processes for regression.  
In Advances in Neural Information Processing Systems (NIPS), pages 514–520, 1996
- C.E. Rasmussen and C.K.I. Williams. Gaussian processes for machine learning. MIT Press, Cambridge, MA, USA, 2006
- S. Roberts, M. Osborne, M. Ebden, S. Reece, N. Gibson, and S. Aigrain. Gaussian processes for time-series modelling. Philosophical Trans. of the Royal Society A, 371(1984):1–25, 2012

## GPLVM

- N. Lawrence. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. Journal of machine learning research, 6:1783-1816, 2005

## GPIs

- O. Williams and A. Fitzgibbon. Gaussian Process Implicit Surfaces. In Gaussian Processes in Practice, 2007