

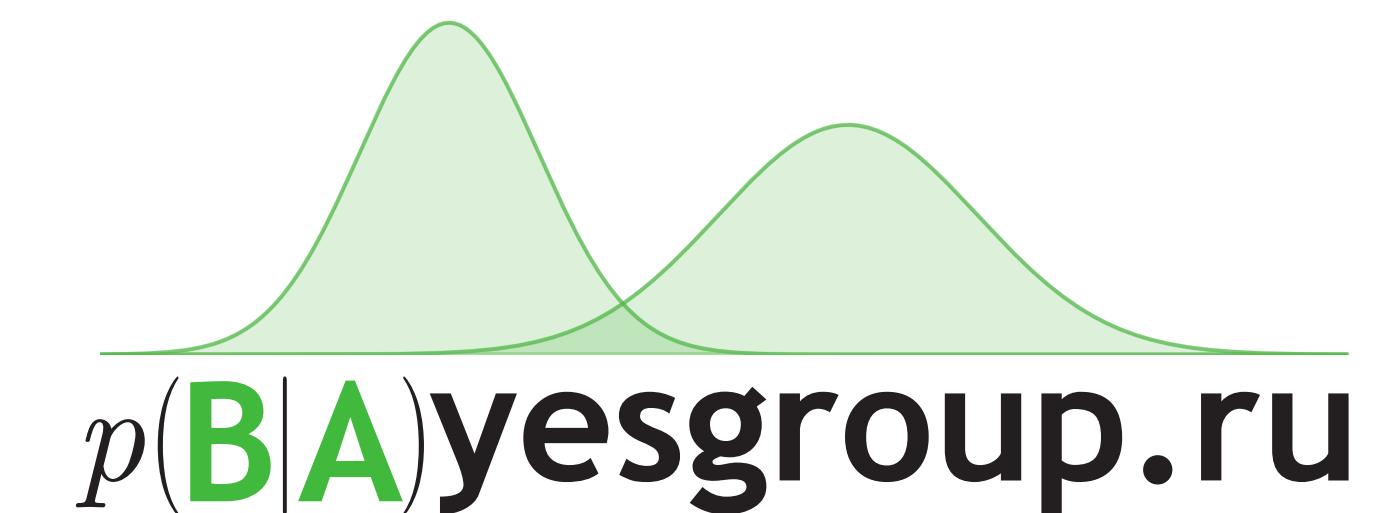
Gaussian processes

Nadia Chirkova

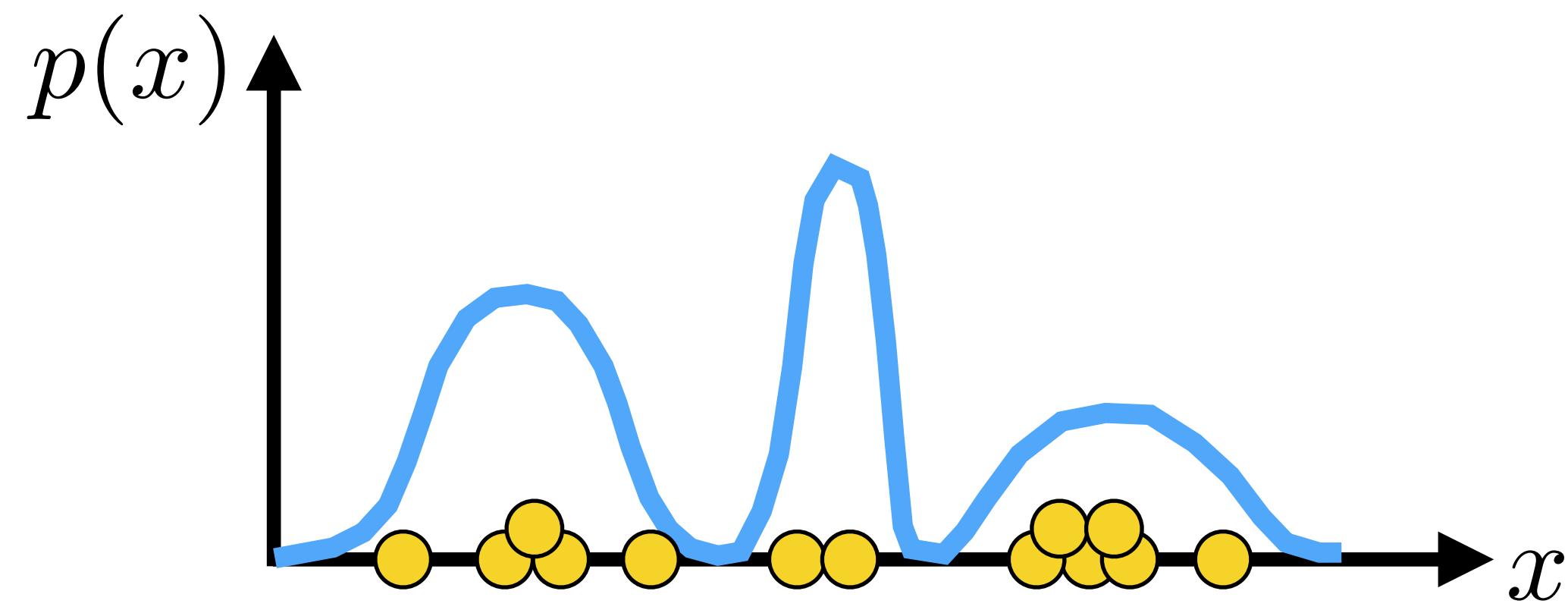
Higher School of Economics, Samsung-HSE Laboratory
Moscow, Russia



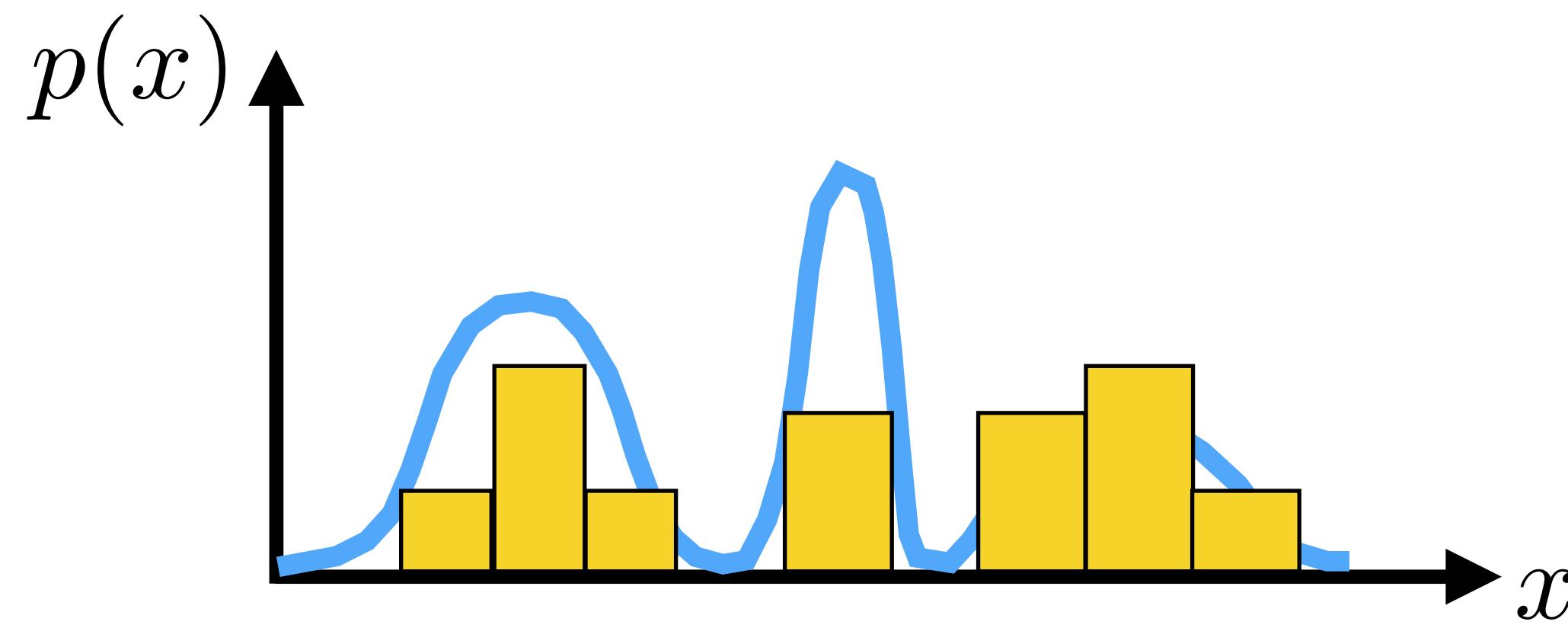
SAMSUNG
Research



Sampling points from a distribution

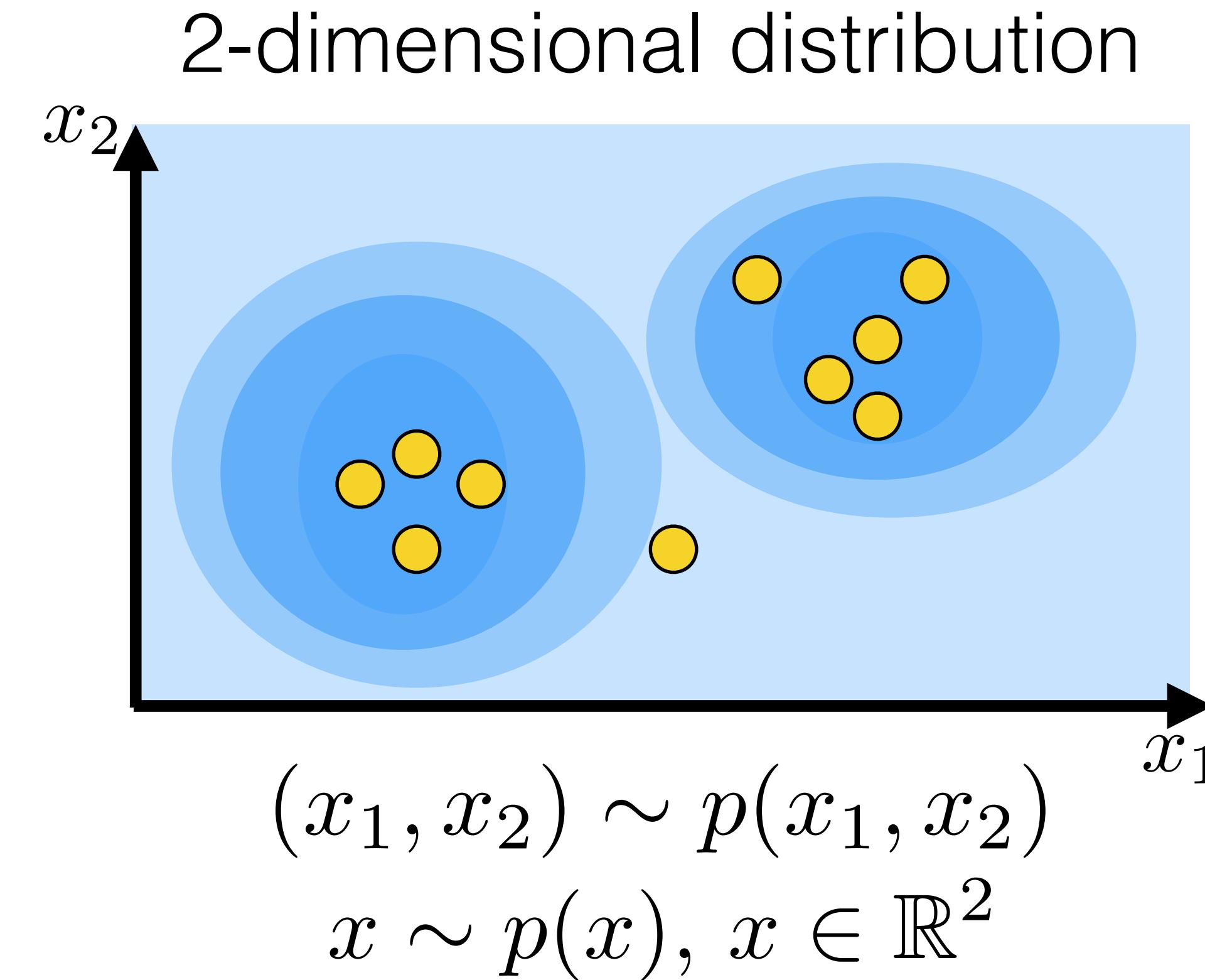
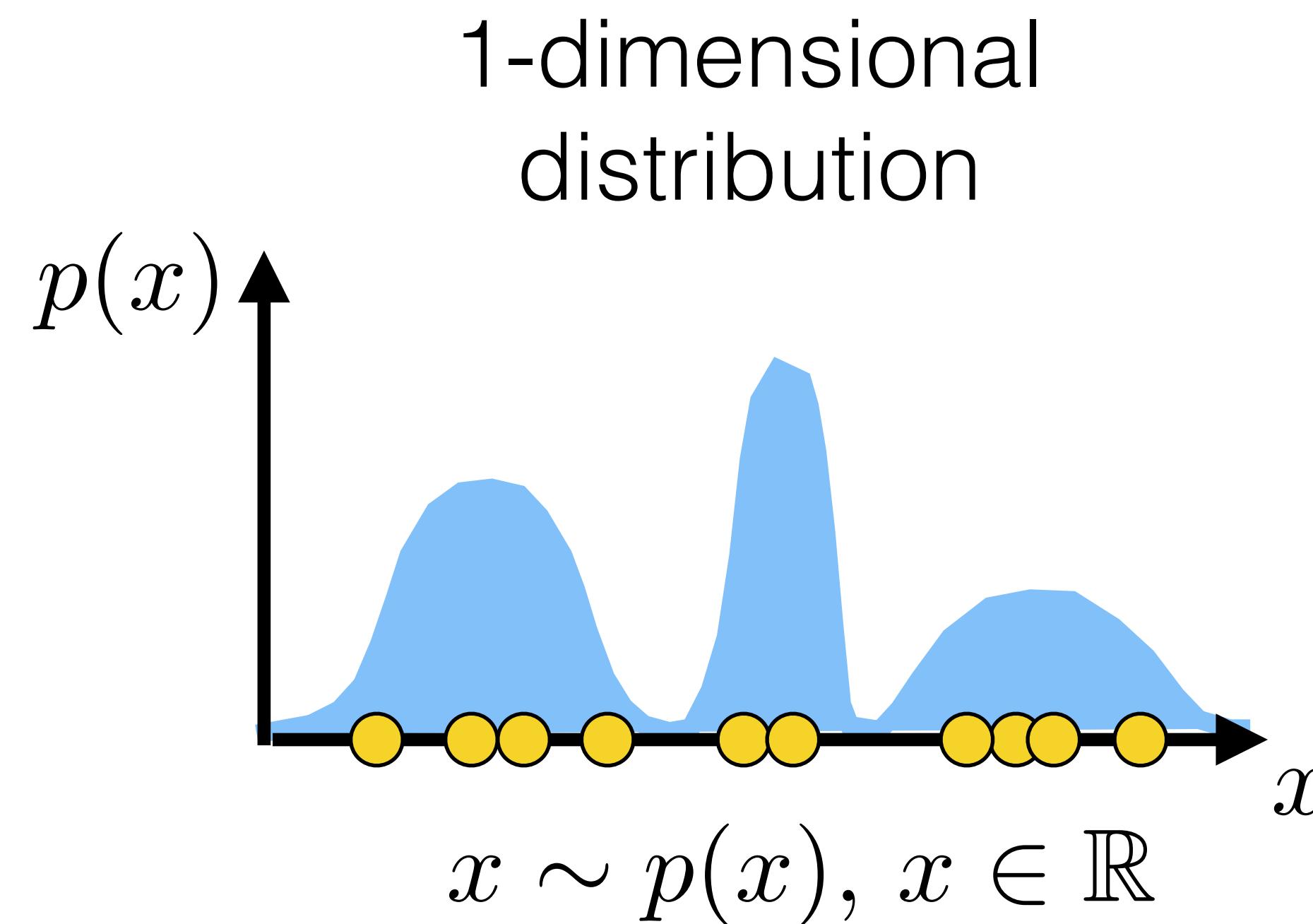


$x \sim p(x), x \in \mathbb{R}$



histogram approaches $p(x)$
when the number of points grows

Sampling points from a multivariate distribution



... n-dimensional distribution:

$$(x_1, \dots, x_n) \sim p(x_1, \dots, x_n)$$
$$x \sim p(x), x \in \mathbb{R}^n$$

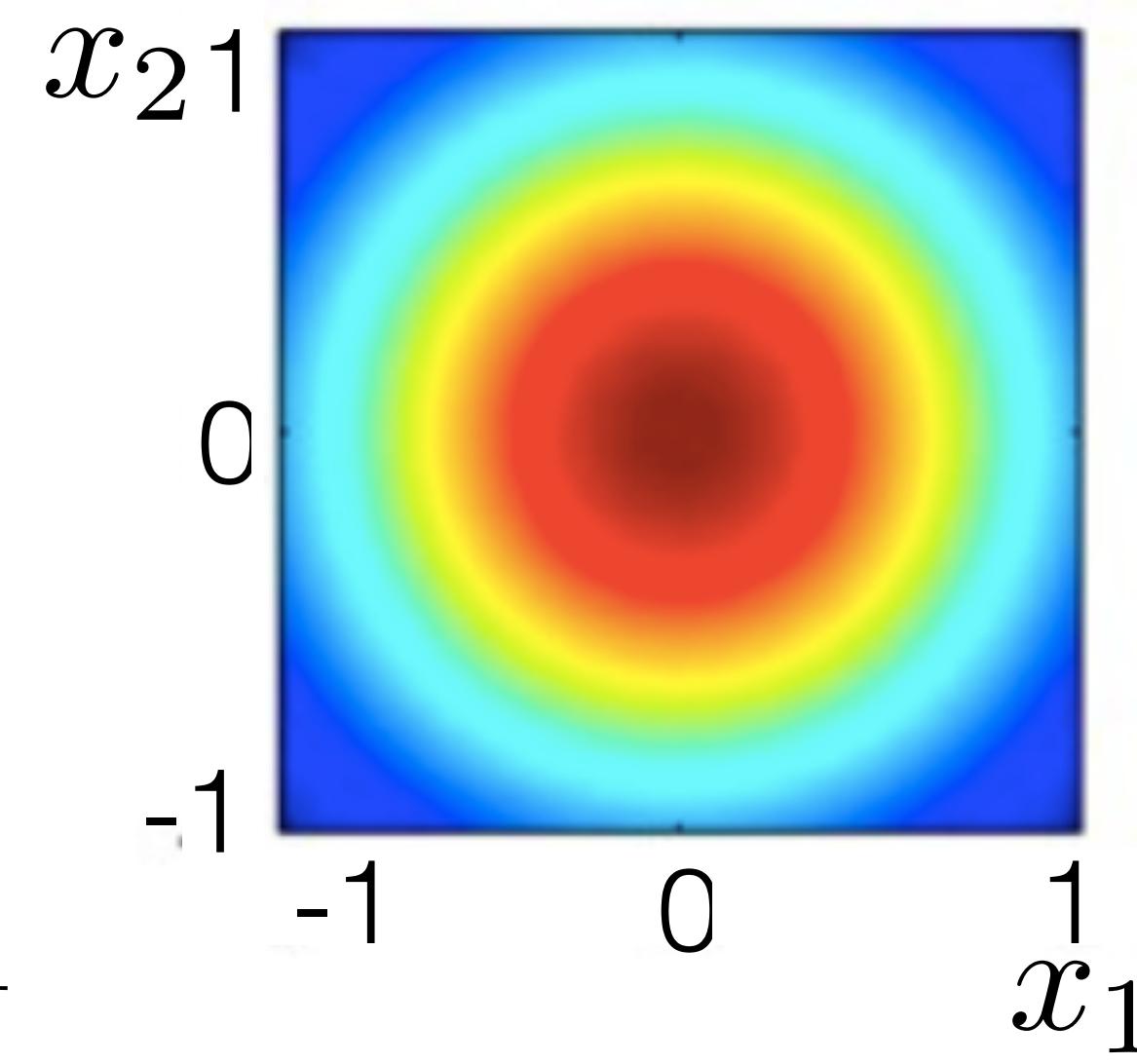
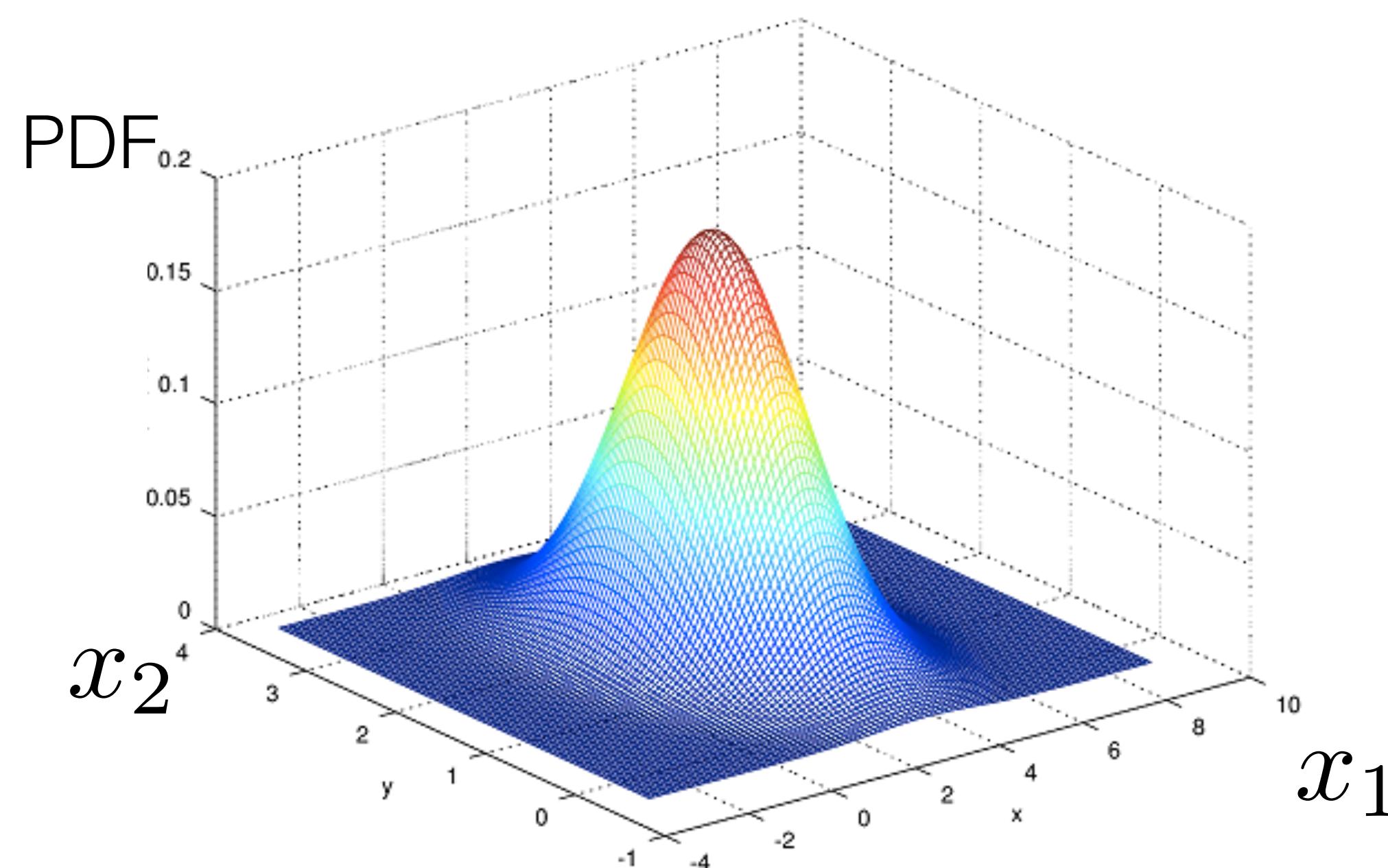
Multivariate normal (Gaussian) distribution

$$\mathcal{N}(\mu, \Sigma) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right),$$

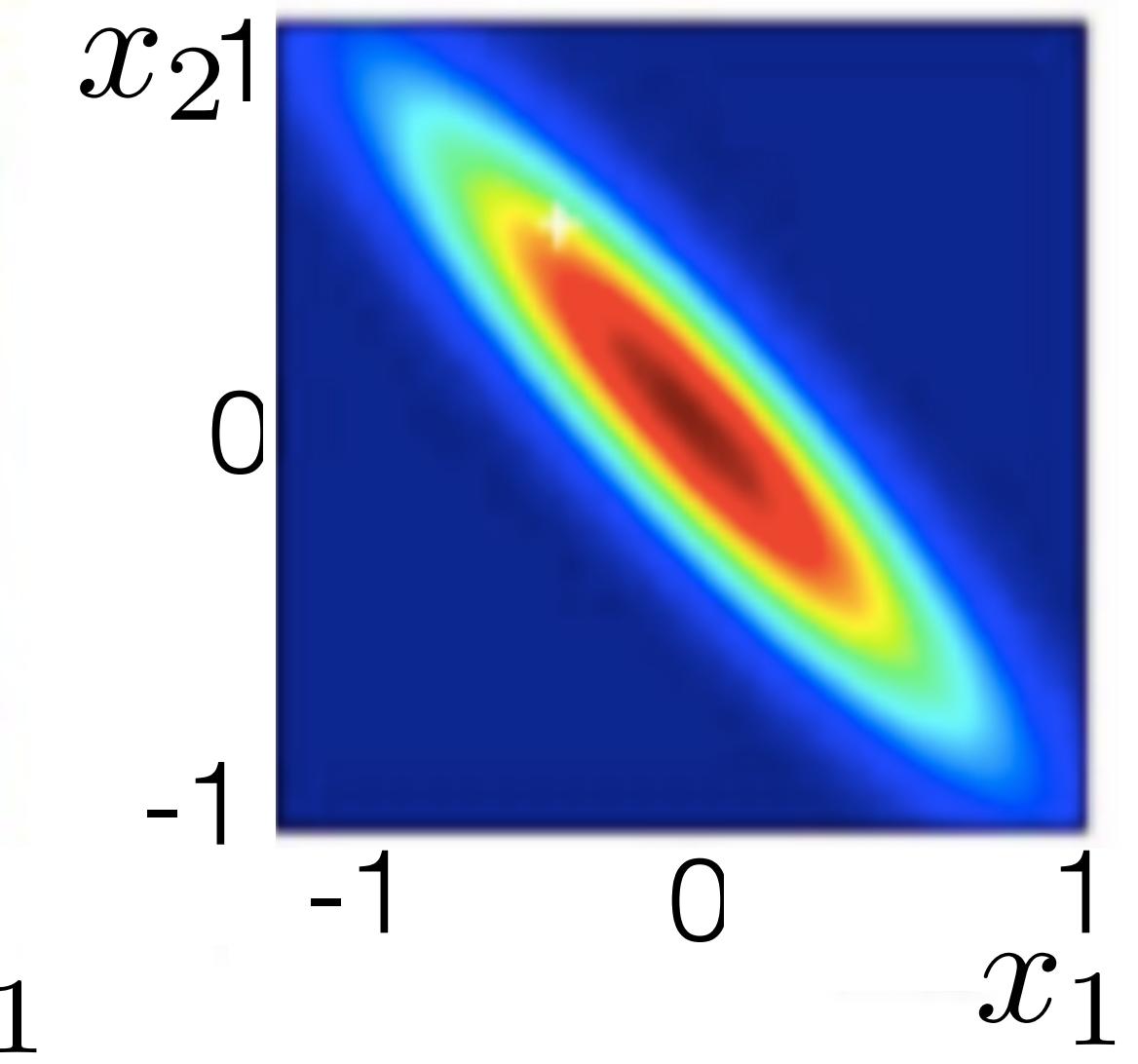
$$x \in \mathbb{R}^d$$

$$\mu \in \mathbb{R}^d$$

$$\Sigma \in \mathbb{R}^{d \times d}$$



diagonal Σ



non-diagonal Σ

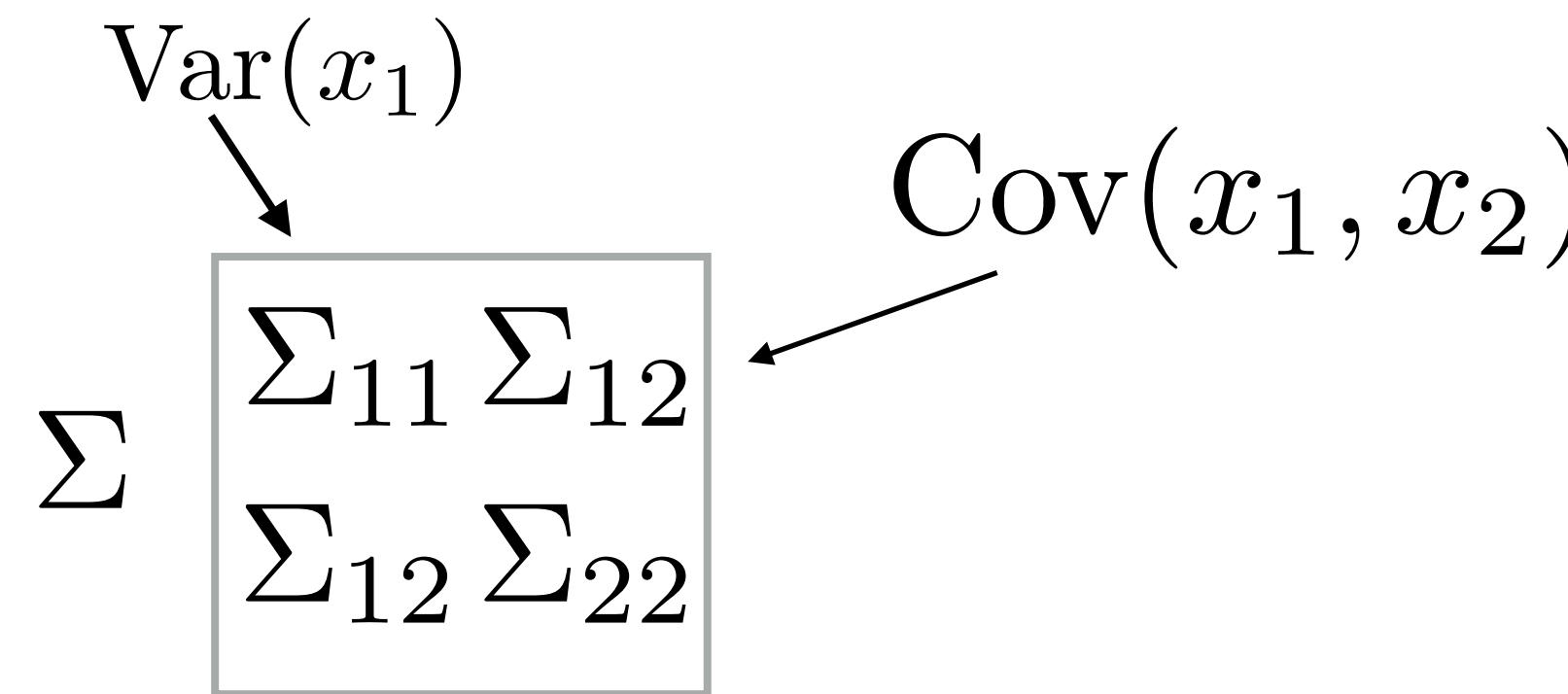
Multivariate normal (Gaussian) distribution

$$\mathcal{N}(\mu, \Sigma) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right),$$

$$x \in \mathbb{R}^d$$

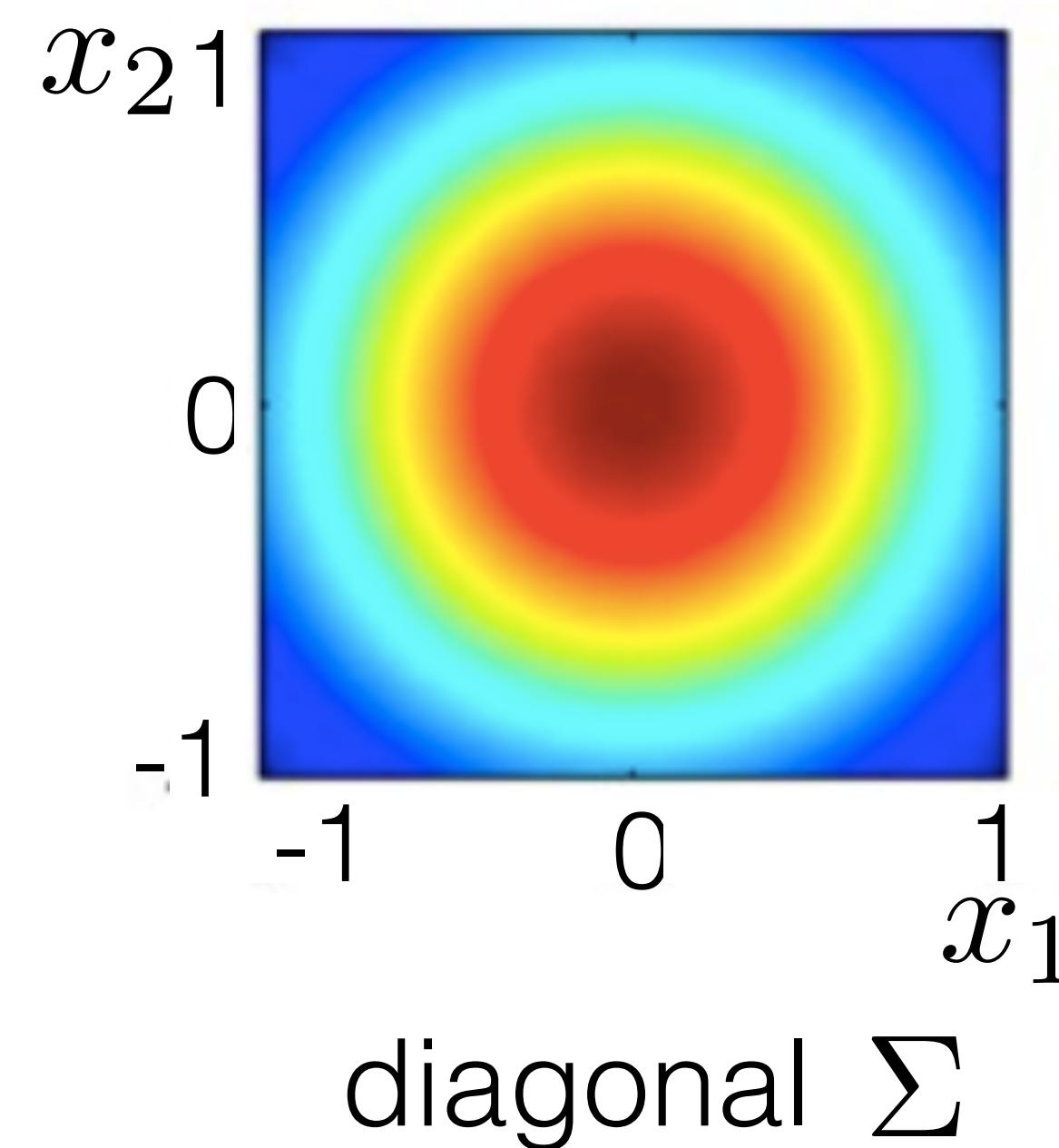
$$\mu \in \mathbb{R}^d$$

$$\Sigma \in \mathbb{R}^{d \times d}$$

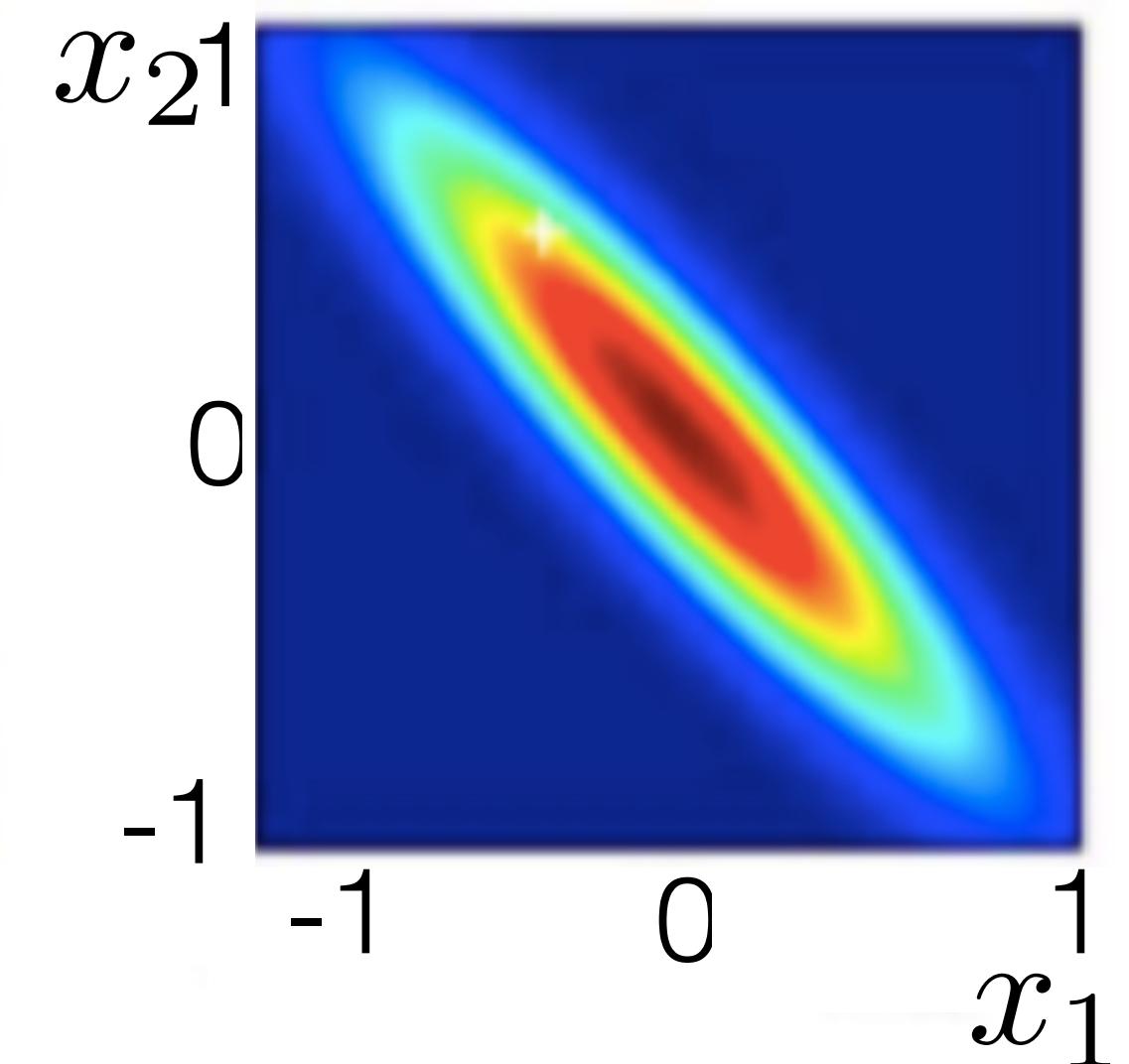


diagonal $\Sigma \Rightarrow$

independent x_1, \dots, x_d



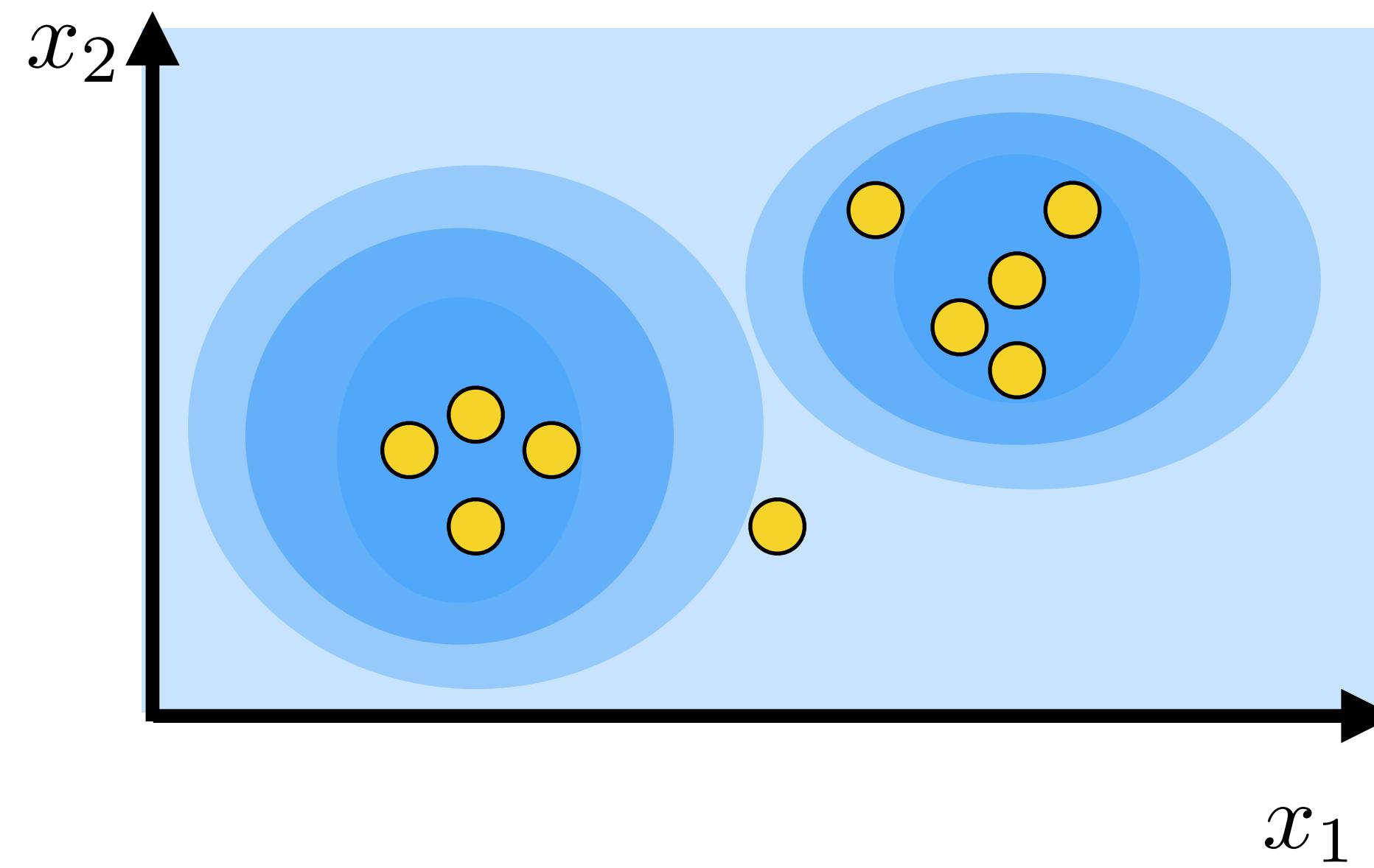
diagonal Σ



non-diagonal Σ

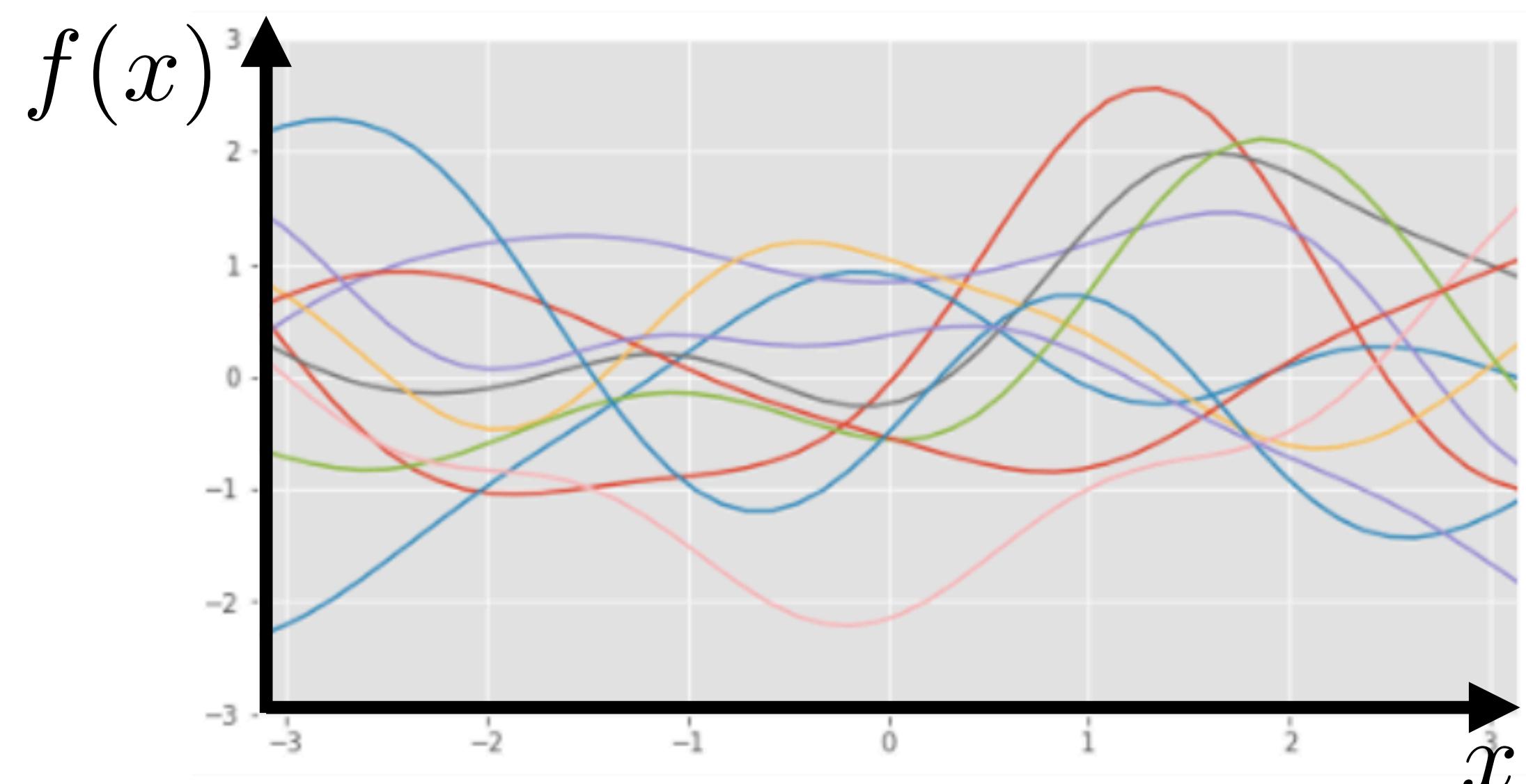
Sampling functions from a process

Multivariate distribution
— sample **points**



$$x \sim p(x)$$

Process
— sample **functions**

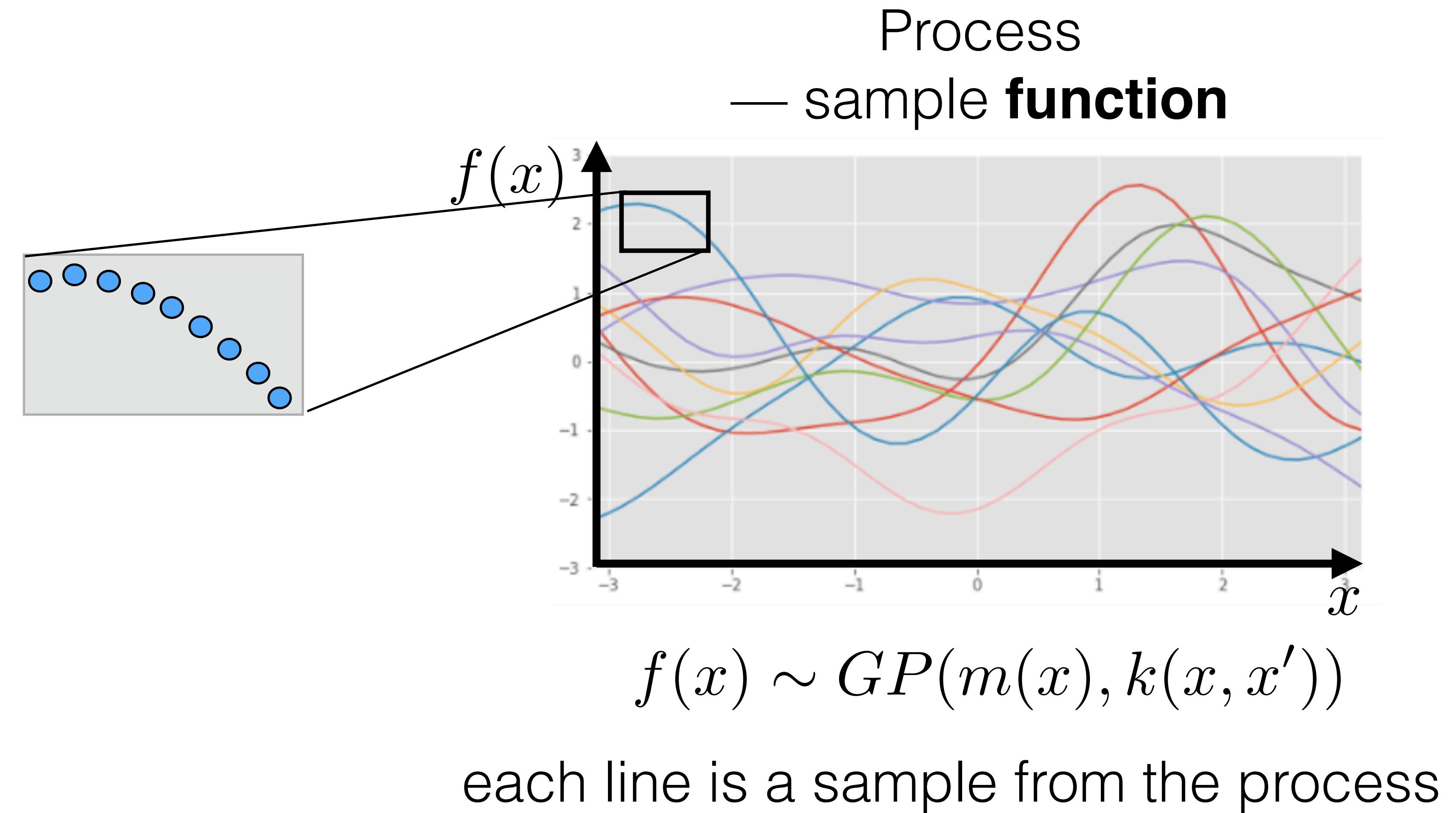


$$f(x) \sim GP(m(x), k(x, x'))$$

each line is a sample from the process

Sampling functions from a process

When we plot a function in python, we define a function as a sequence of points



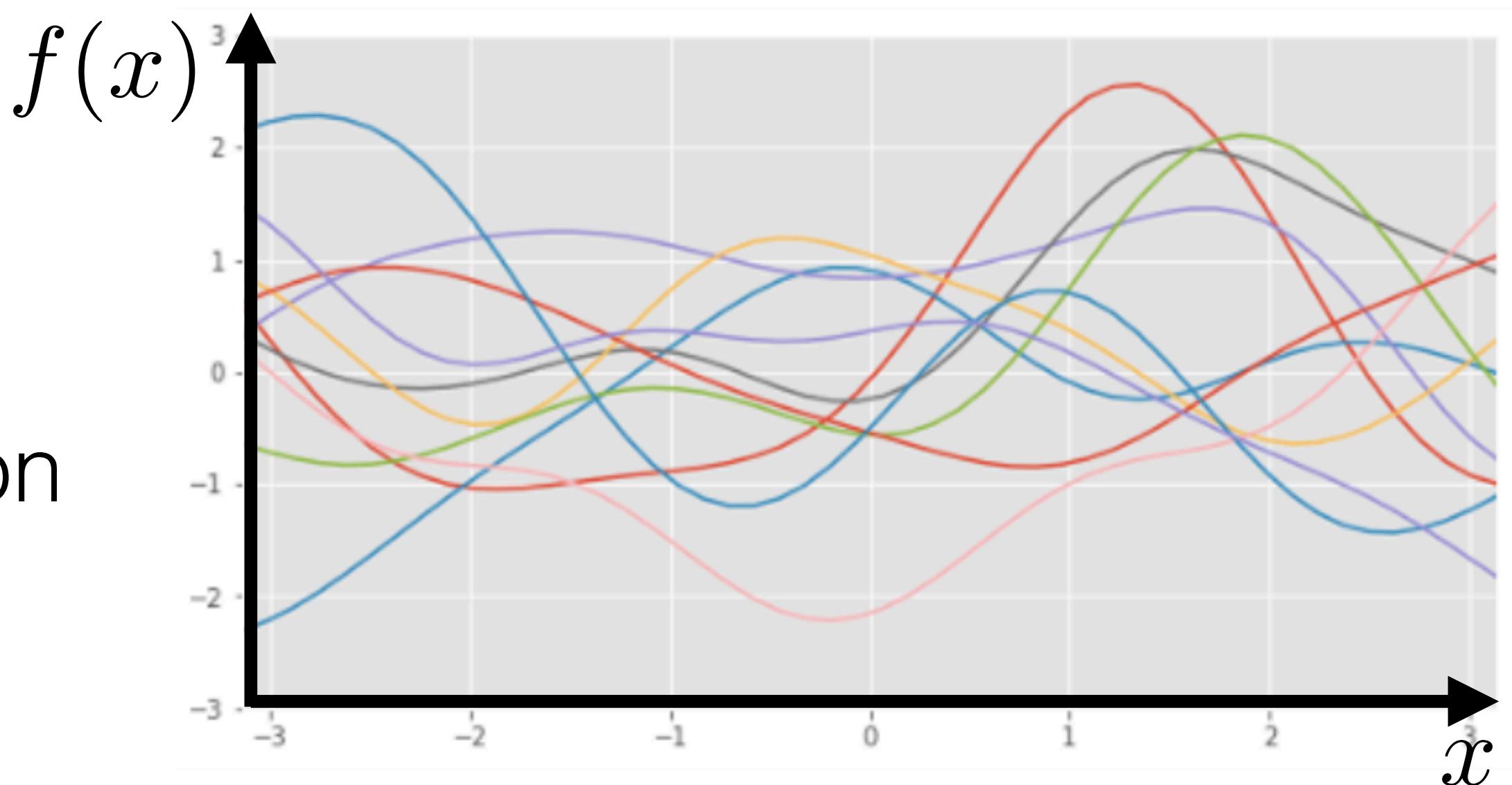
Gaussian process

$$f(x) \sim GP(m(x), k(x, x'))$$

$m(x)$ — mean function

$k(x, x')$ — covariance (or kernel) function

(x may be a vector in general case)



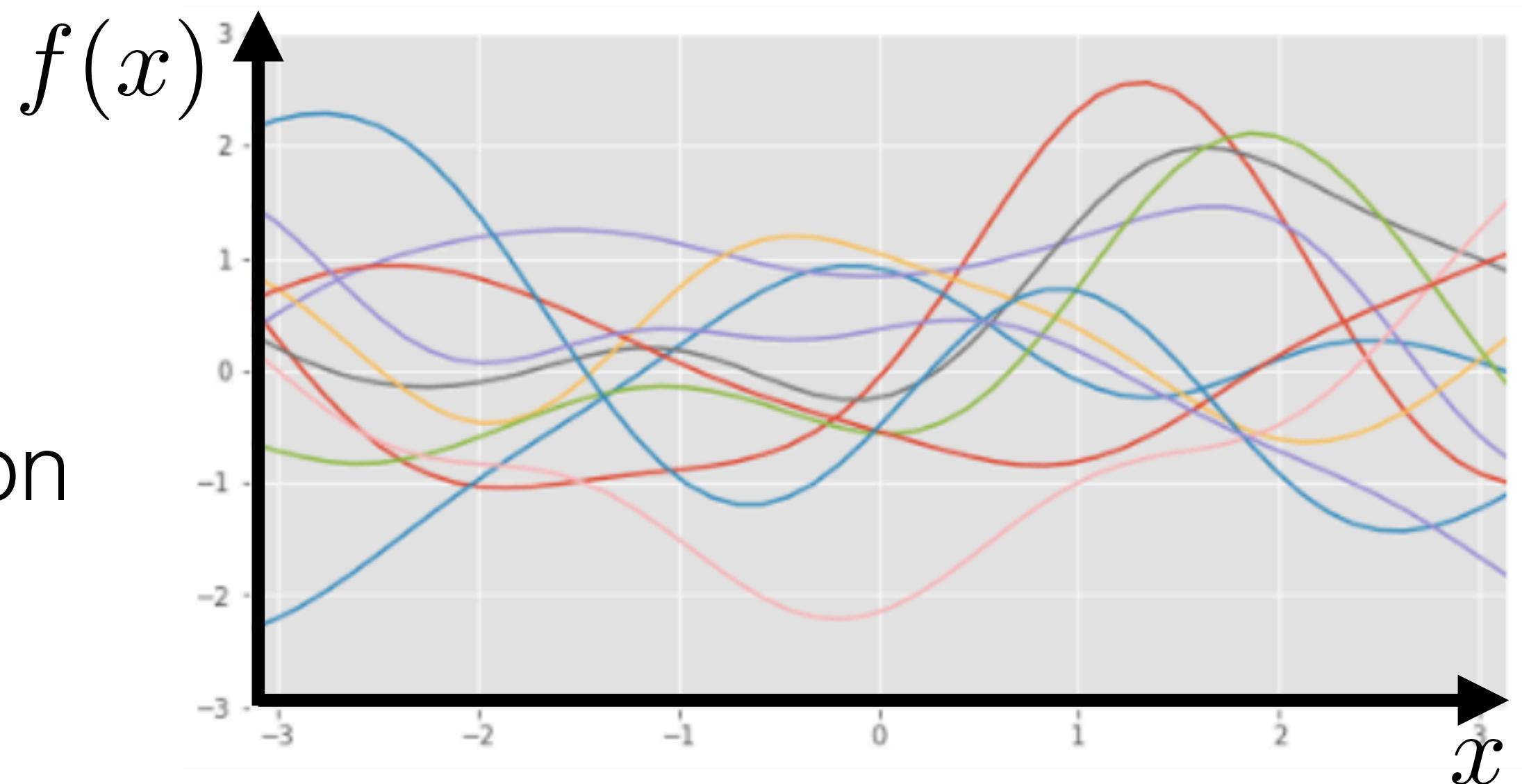
Gaussian process

$$f(x) \sim GP(m(x), k(x, x'))$$

$m(x)$ — mean function

$k(x, x')$ — covariance (or kernel) function

(x may be a vector in general case)



Definition of Gaussian process:

every finite set of function values has a multivariate normal distribution

$$\forall n \quad \forall (x_1, \dots, x_n) \quad (f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = \{m(x_i)\}_{i=1}^n \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$

Example 1: white noise

$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

$$k(x, x') = \sigma^2[x = x']$$

[condition] = 1 if condition is True else 0

Example 1: white noise

$$f(x) \sim GP(m(x), k(x, x'))$$

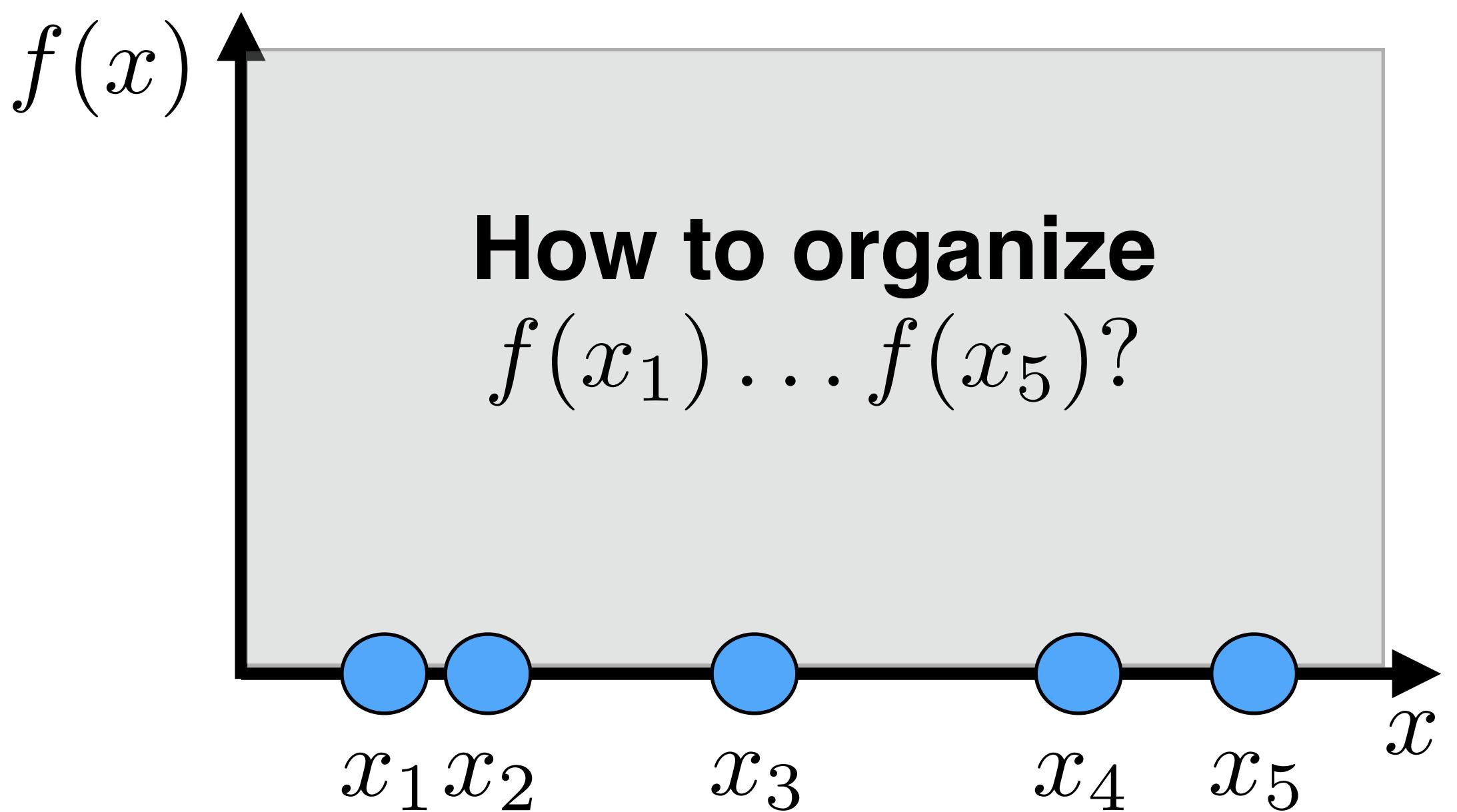
$$m(x) = 0$$

$$k(x, x') = \sigma^2[x = x']$$

[condition] = 1 if condition is True else 0

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = 0 \quad \Sigma = \sigma^2 I$$



Example 1: white noise

$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

$$k(x, x') = \sigma^2[x = x']$$

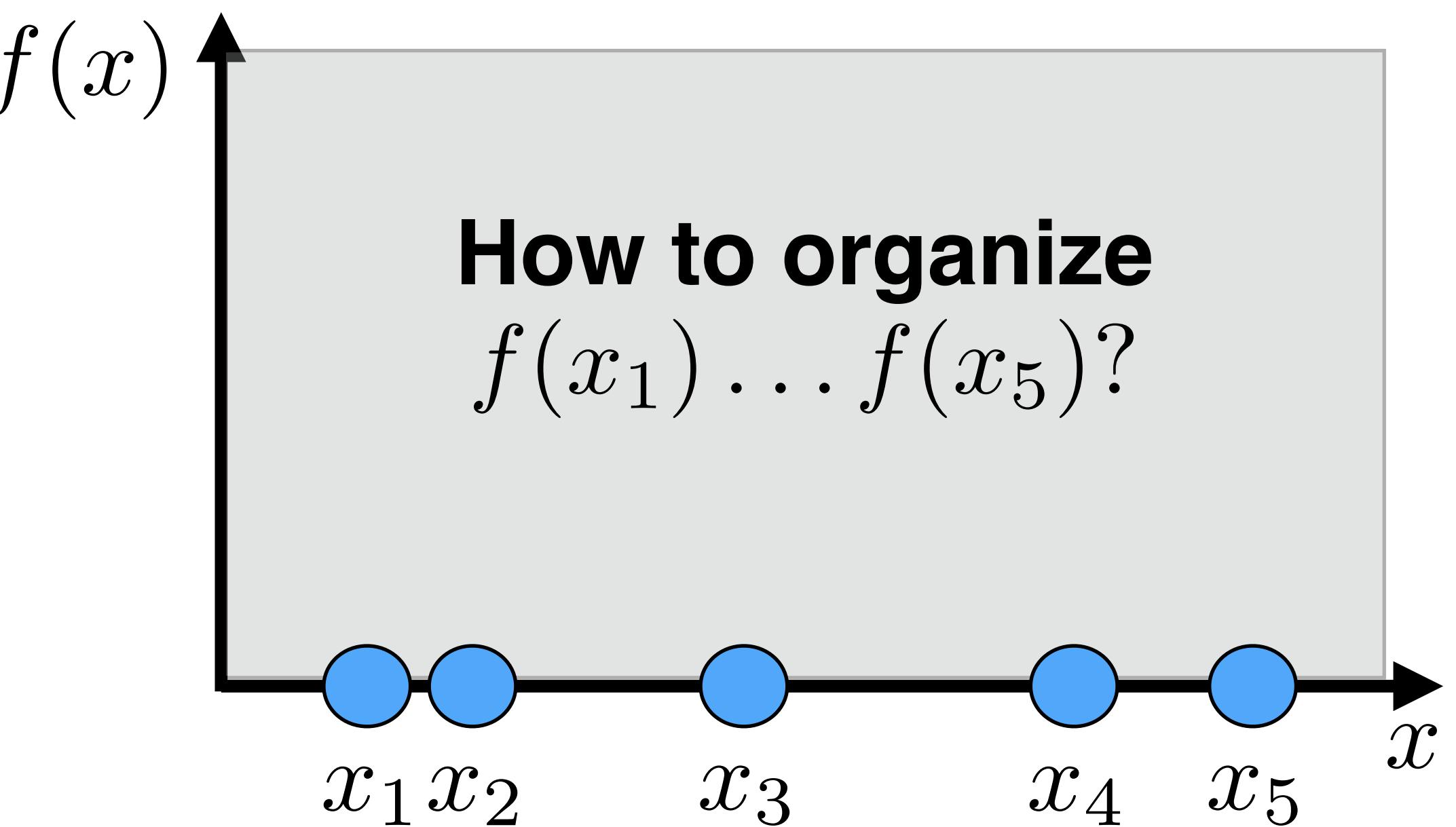
[condition] = 1 if condition is True else 0

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = 0 \quad \Sigma = \sigma^2 I$$

$$p(f(x_1), \dots, f(x_n)) = \prod_{i=1}^n \mathcal{N}(0, \sigma^2)$$

(all x are independent on each other)



Example 1: white noise

$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

$$k(x, x') = \sigma^2[x = x']$$

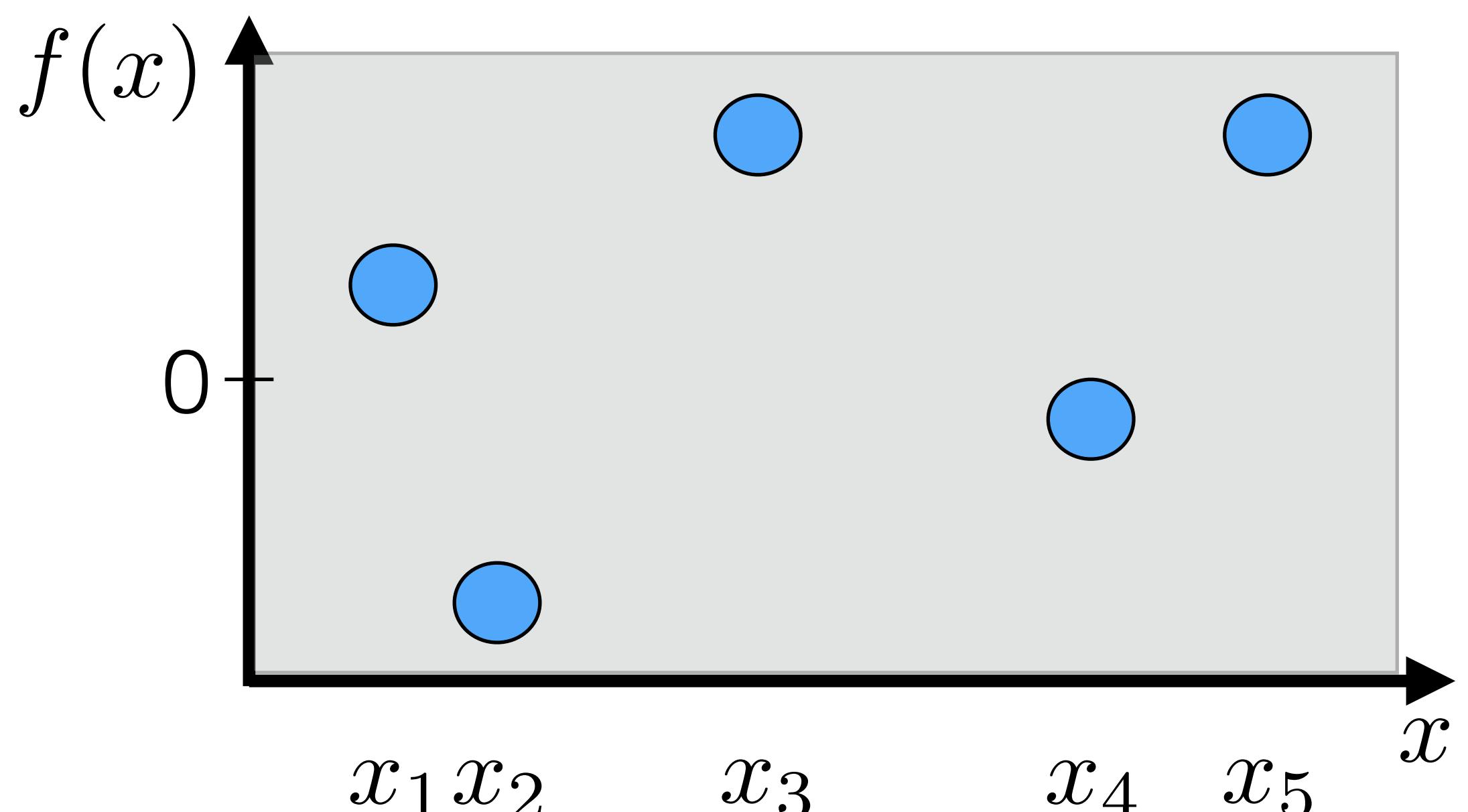
[condition] = 1 if condition is True else 0

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = 0 \quad \Sigma = \sigma^2 I$$

$$p(f(x_1), \dots, f(x_n)) = \prod_{i=1}^n \mathcal{N}(0, \sigma^2)$$

(all x are independent on each other)



for any x
 $f(x)$ is sampled
independently

Example 2: constant function

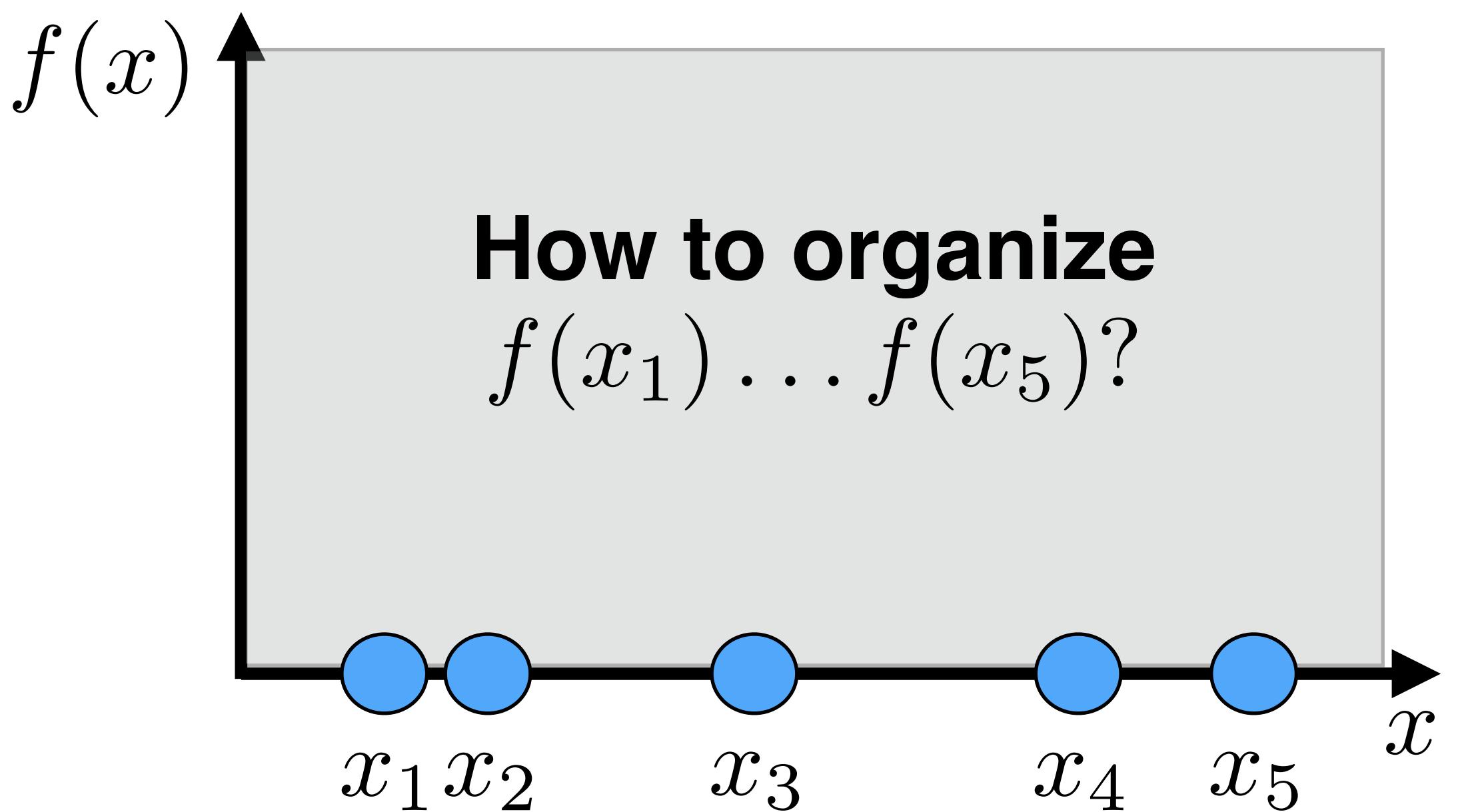
$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

$$k(x, x') = C$$

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = 0 \quad \Sigma = \{C\}_{i,j=1}^{n,n}$$



Example 2: constant function

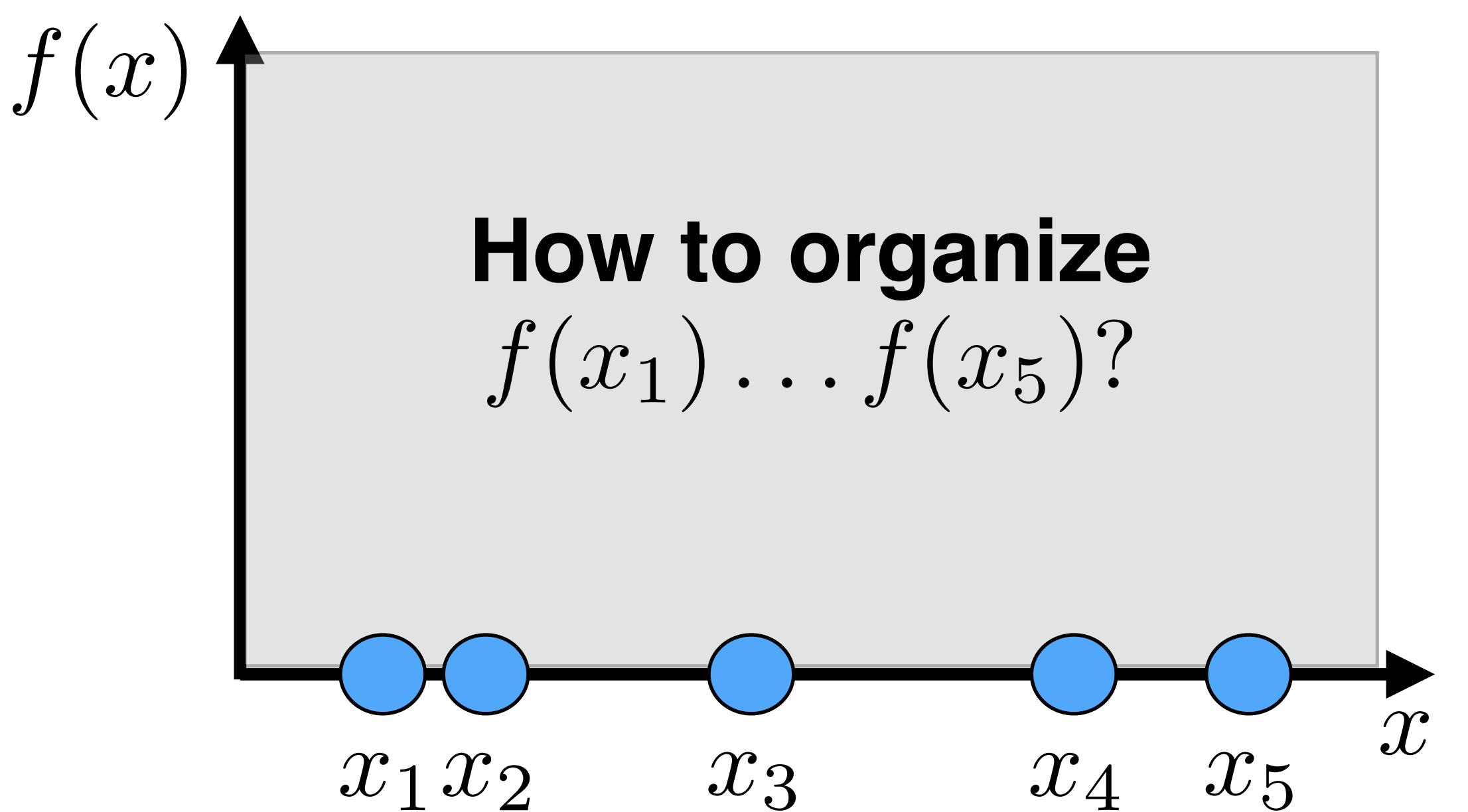
$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

$$k(x, x') = C$$

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = 0 \quad \Sigma = \{C\}_{i,j=1}^{n,n}$$



$$\forall i \neq j$$

$$\text{Corr}(f(x_i), f(x_j)) = \frac{\text{Cov}(f(x_i), f(x_j))}{\sqrt{\text{Var}(f(x_i))\text{Var}(f(x_j))}} = \frac{C}{\sqrt{C^2}} = 1 \quad \boxed{\Rightarrow f(x_i) = f(x_j)}$$

$$\text{Var}(f(x_i)) = \text{Var}(f(x_j)) = C, \quad \mathbb{E}f(x_i) = \mathbb{E}f(x_j) = 0$$

Example 2: constant function

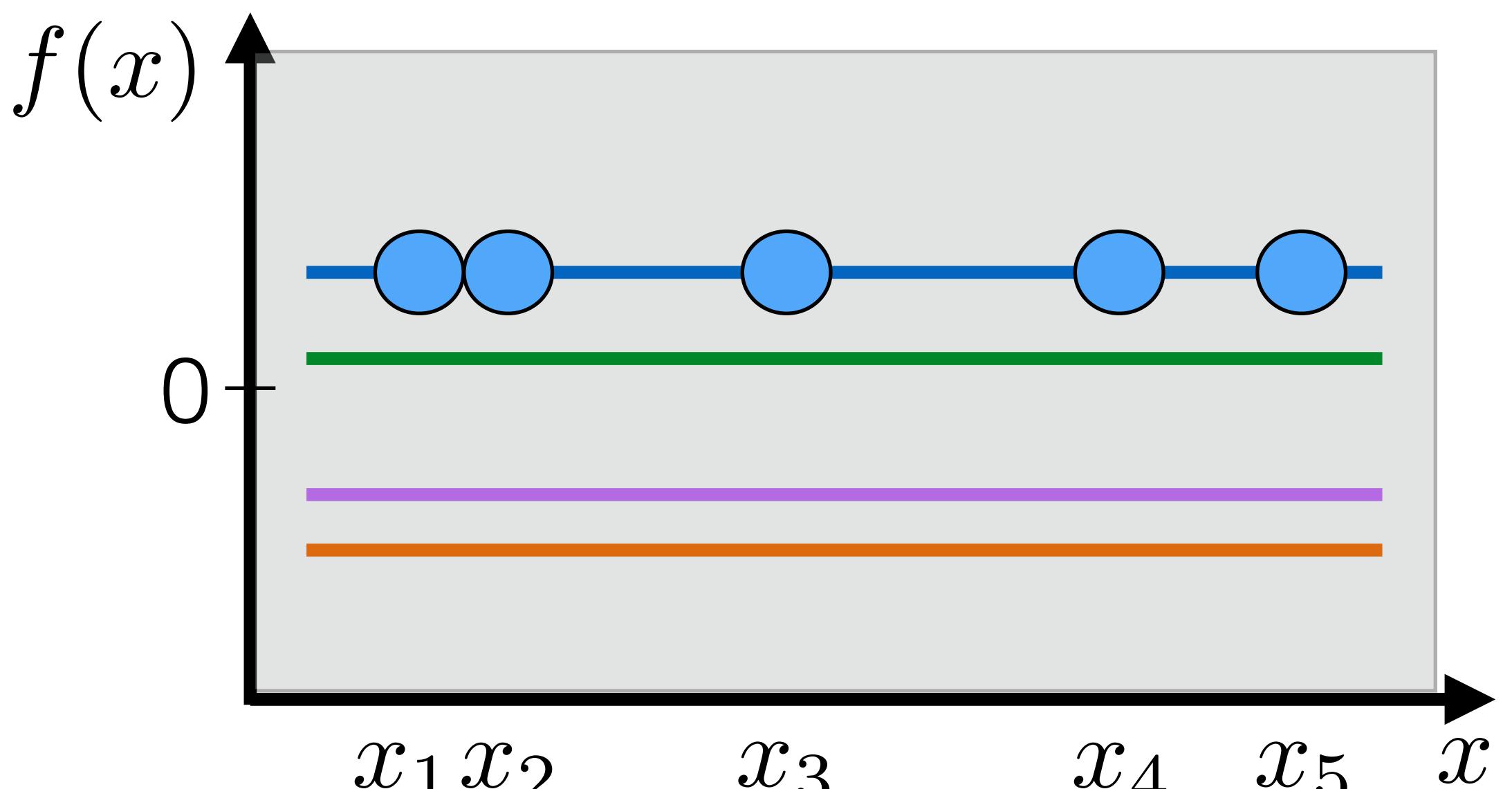
$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

$$k(x, x') = C$$

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = 0 \quad \Sigma = \{C\}_{i,j=1}^{n,n}$$



$$\forall i \neq j$$

$$\text{Corr}(f(x_i), f(x_j)) = \frac{\text{Cov}(f(x_i), f(x_j))}{\sqrt{\text{Var}(f(x_i))\text{Var}(f(x_j))}} = \frac{C}{\sqrt{C^2}} = 1 \quad \left. \right] \Rightarrow f(x_i) = f(x_j)$$

$$\text{Var}(f(x_i)) = \text{Var}(f(x_j)) = C, \quad \mathbb{E}f(x_i) = \mathbb{E}f(x_j) = 0 \quad \left. \right]$$

Example 3: RBF-kernel

$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$$

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = 0 \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$

Example 3: RBF-kernel

$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$$

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = 0 \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$

$$\text{if } \|x_i - x_j\| \approx 0 \quad \Rightarrow \quad \Sigma_{ij} \approx \sigma^2 = \Sigma_{ii} = \Sigma_{jj} \quad \Rightarrow \quad f(x_i) \approx f(x_j)$$

$$\text{if } \|x_i - x_j\| \gg 0 \quad \Rightarrow \quad \Sigma_{ij} \approx 0, \quad f(x_i) \text{ and } f(x_j) \text{ are not correlated}$$

Example 3: RBF-kernel

$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

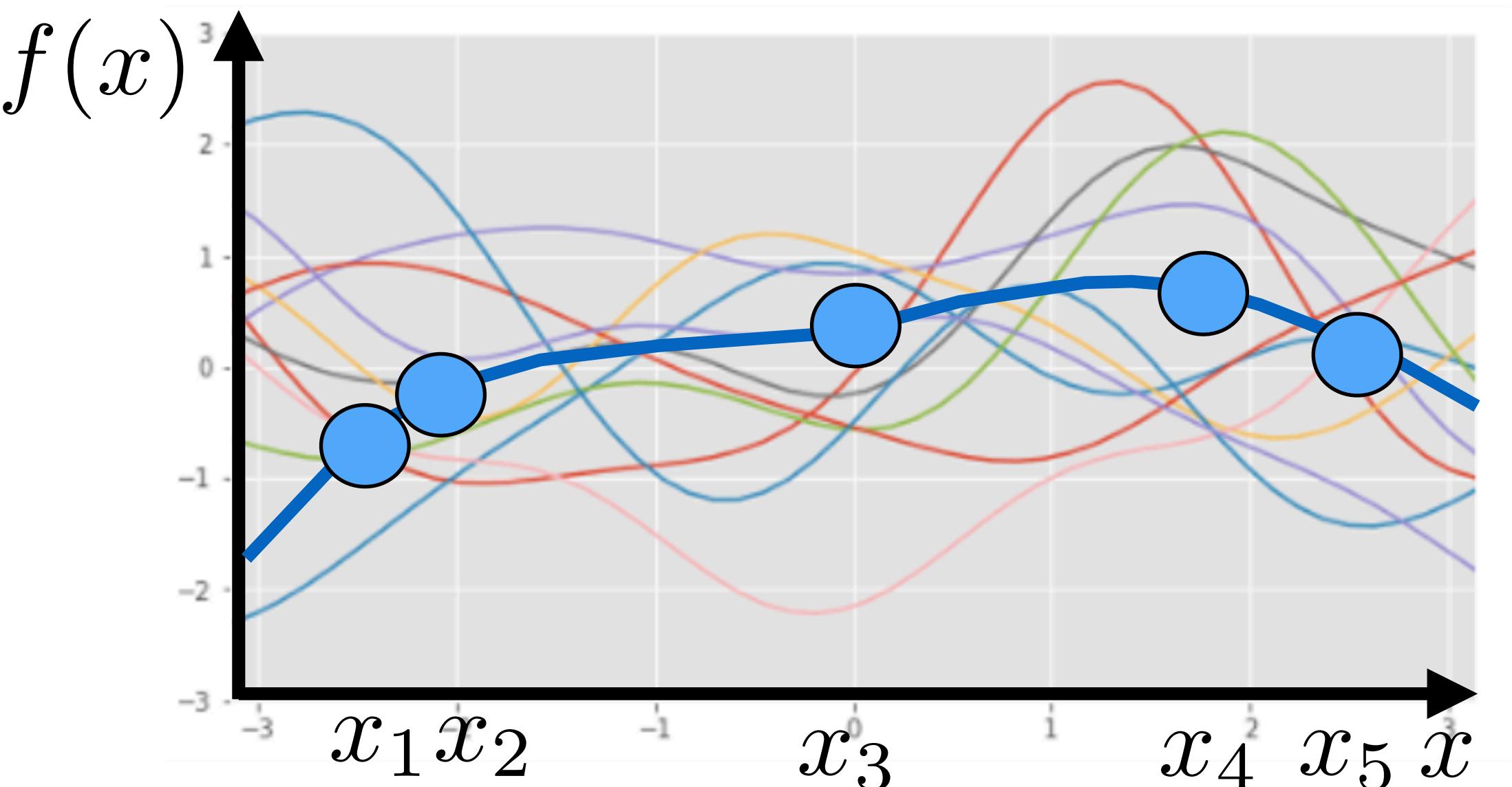
$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$$

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

$$\mu = 0 \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$

if $\|x_i - x_j\| \approx 0 \Rightarrow \Sigma_{ij} \approx \sigma^2 = \Sigma_{ii} = \Sigma_{jj} \Rightarrow f(x_i) \approx f(x_j)$

if $\|x_i - x_j\| \gg 0 \Rightarrow \Sigma_{ij} \approx 0, f(x_i)$ and $f(x_j)$ are not correlated



Example 3: RBF-kernel

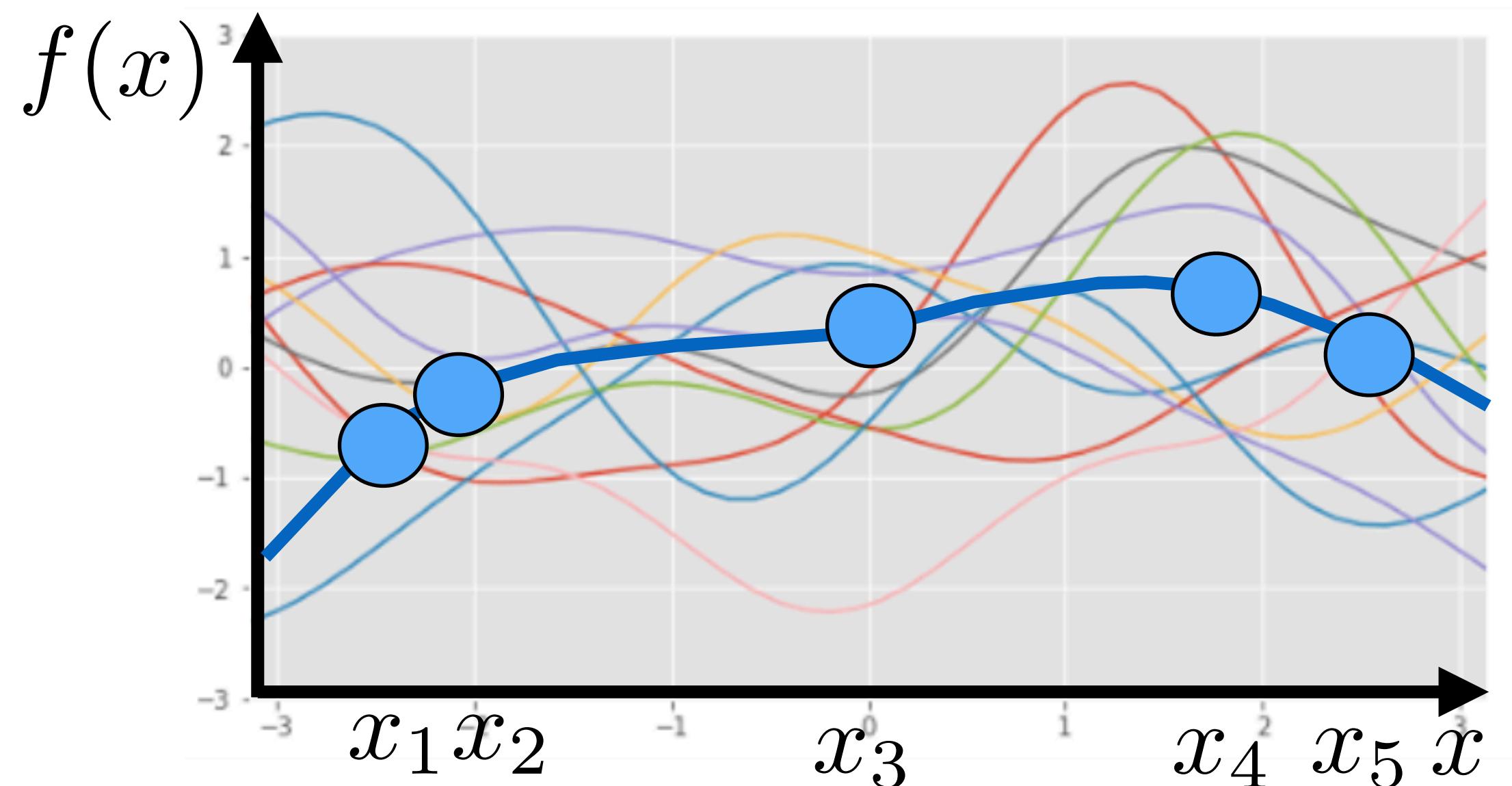
$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$$

$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

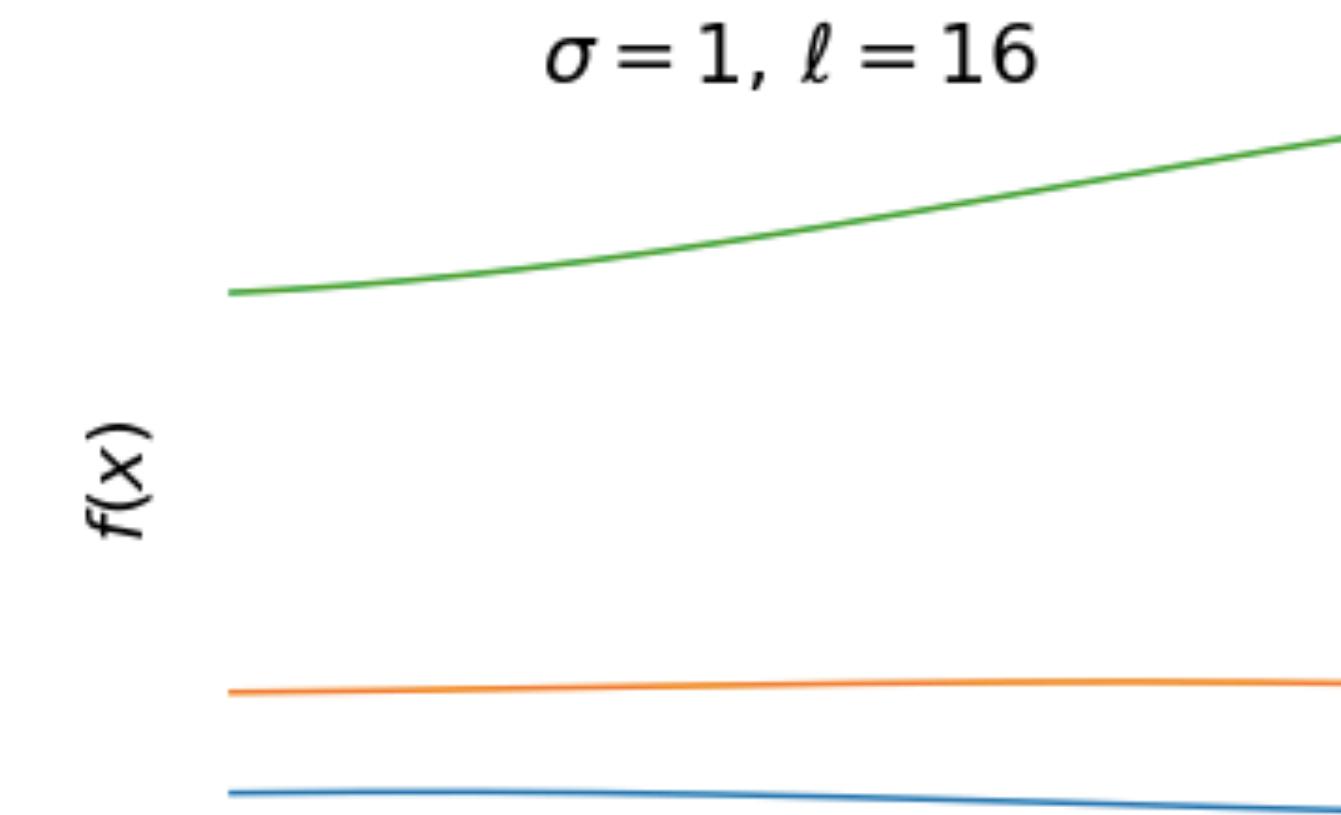
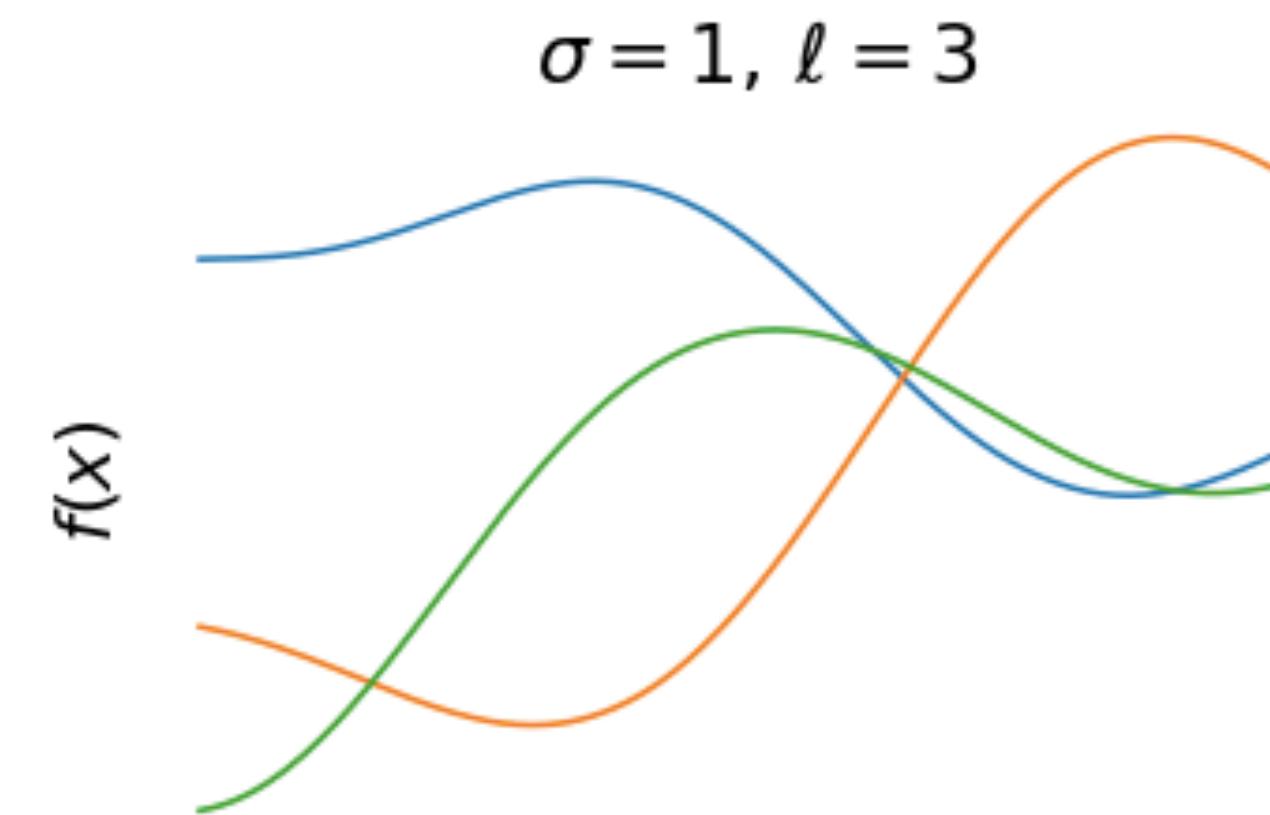
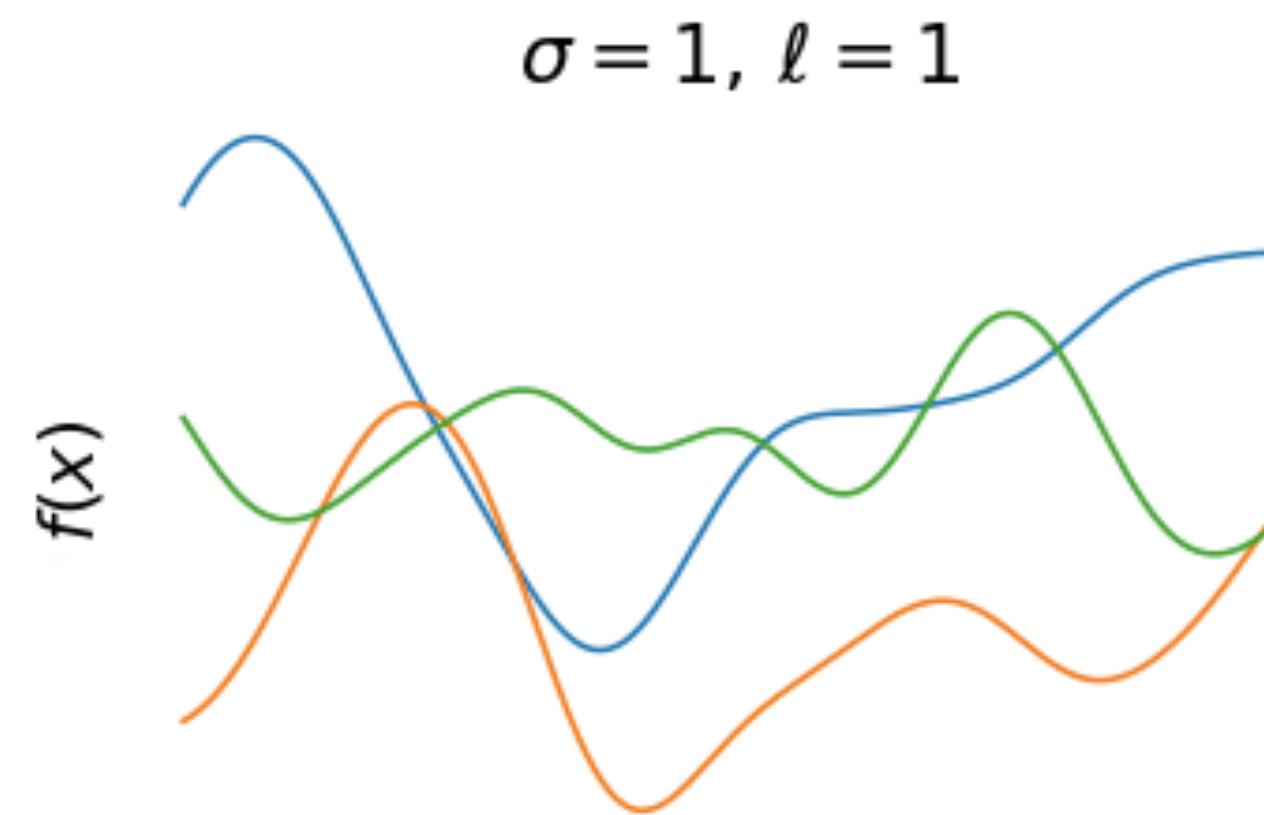
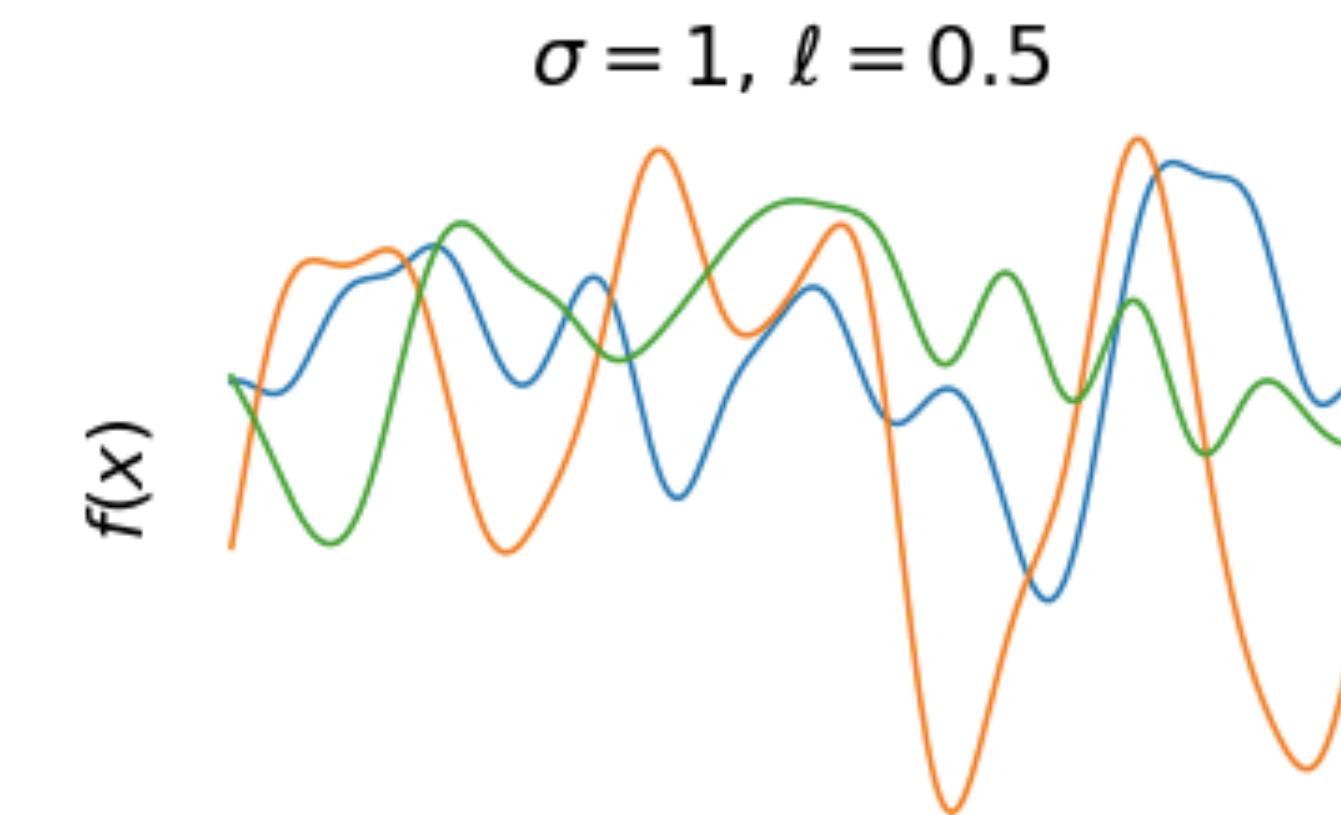
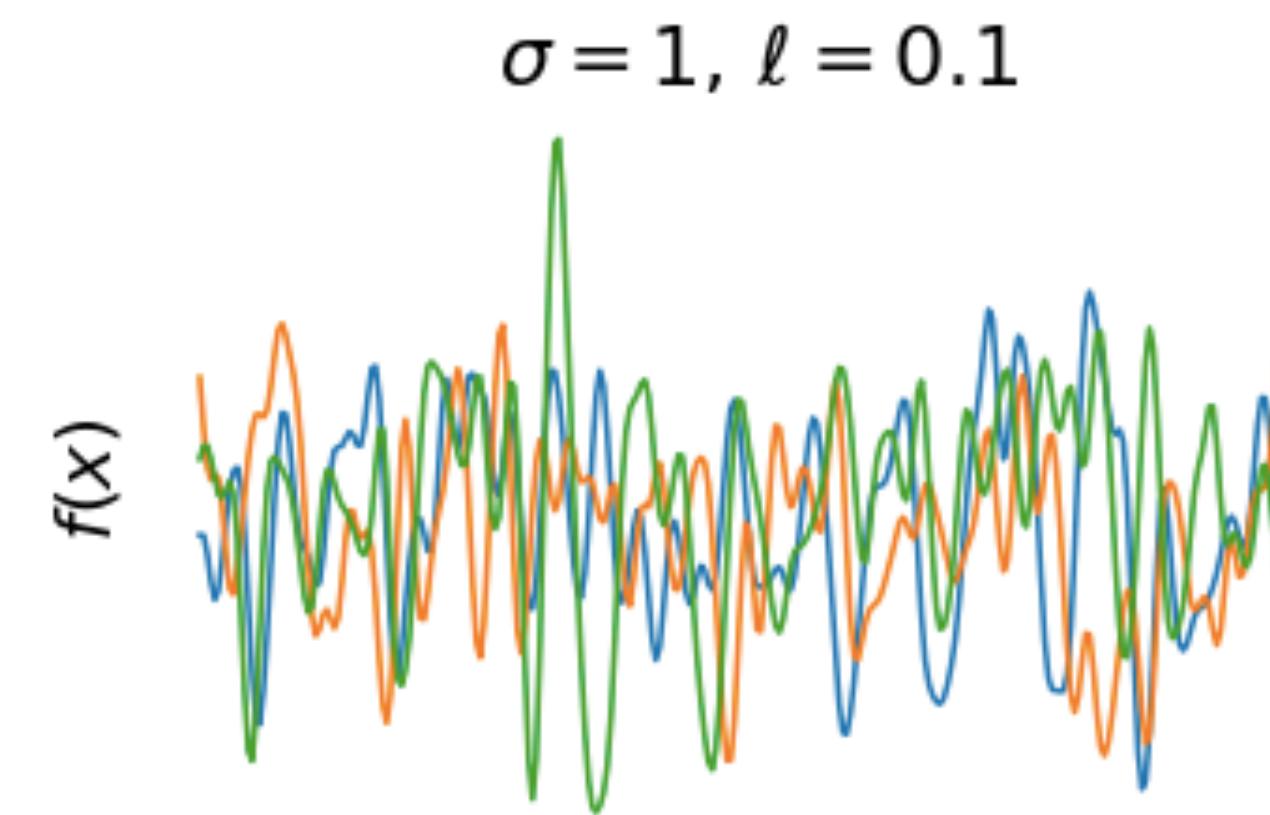
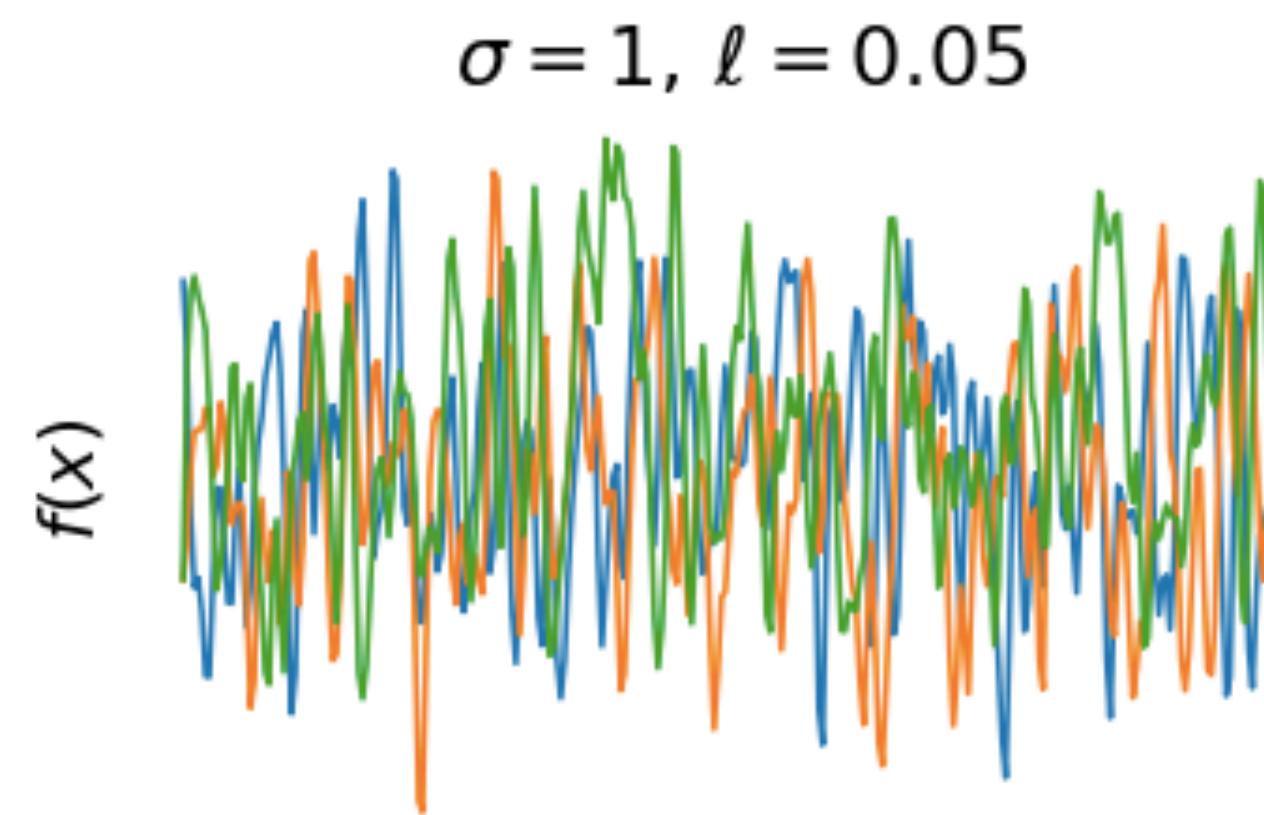
$$\mu = 0 \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$



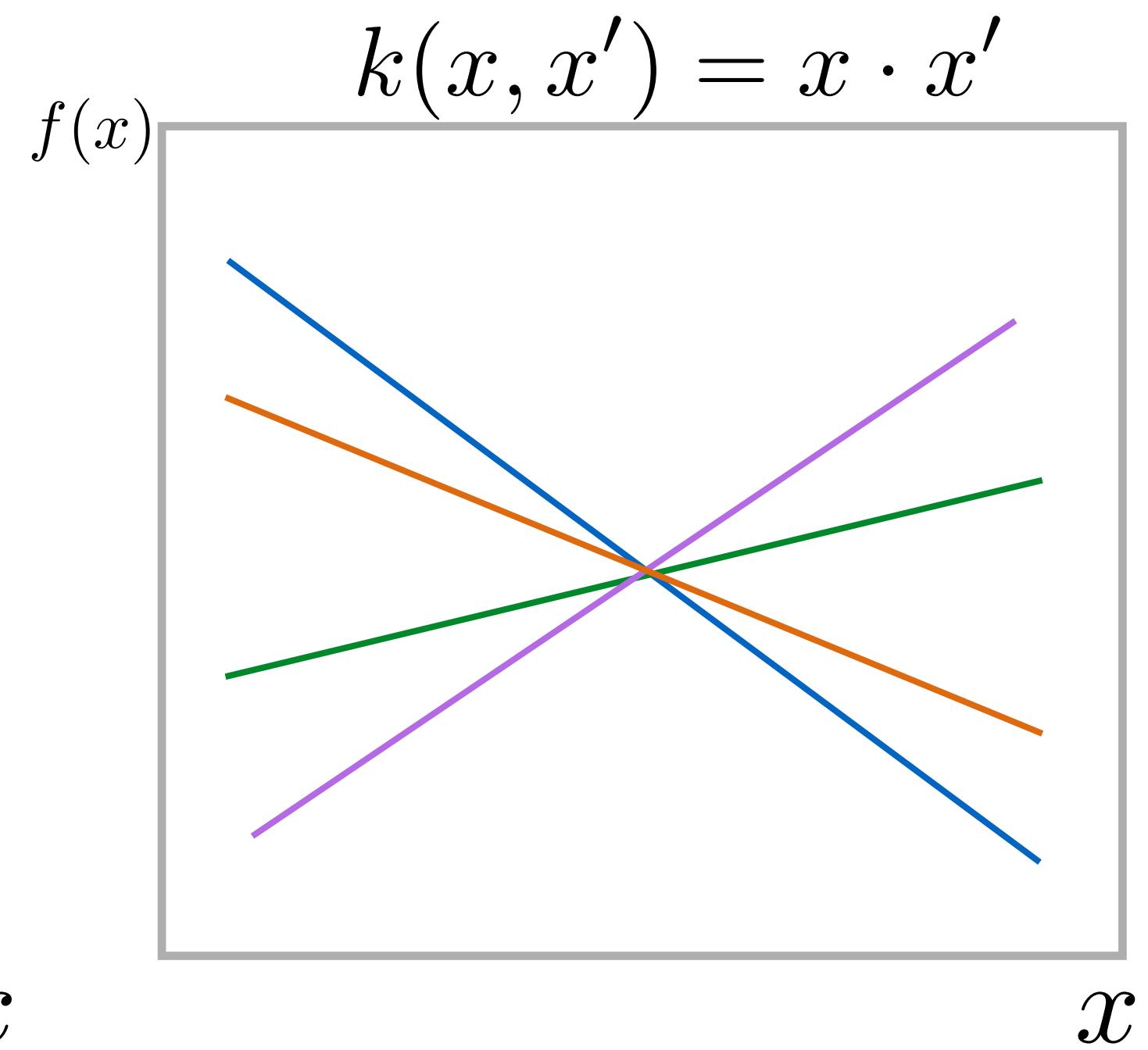
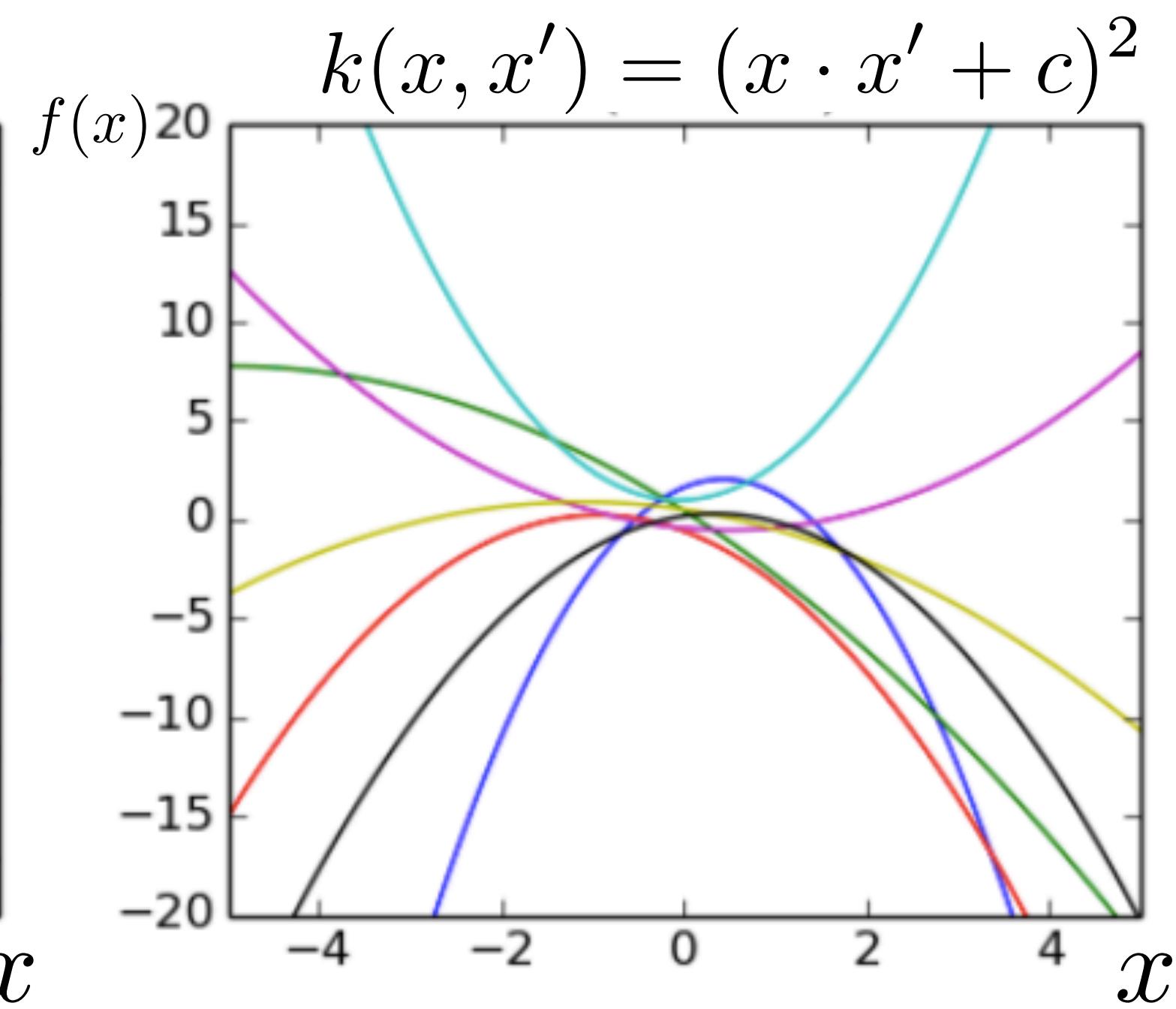
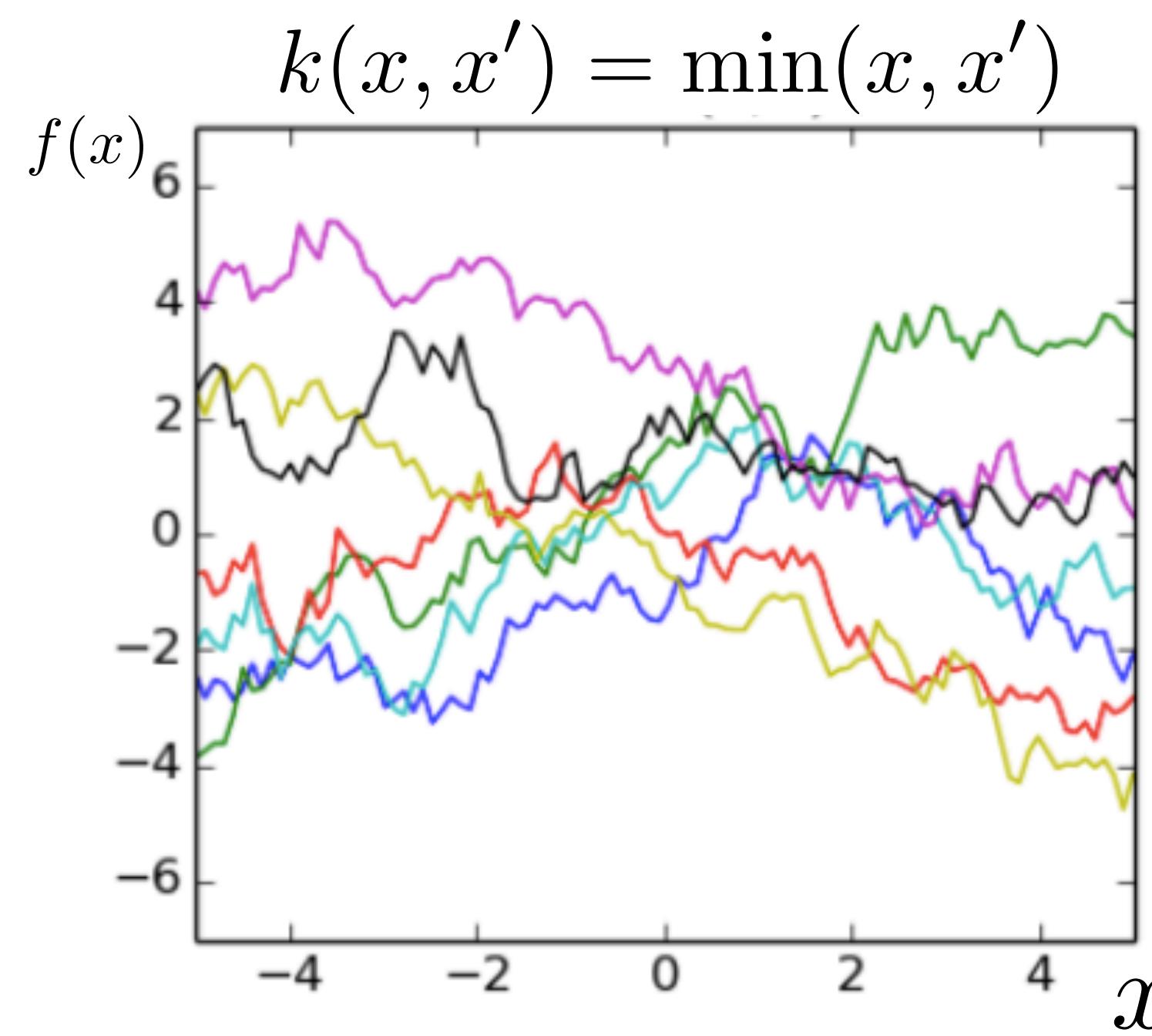
σ^2 defines the “height” of the function
 ℓ^2 defines the frequency of fluctuations

Example 3: RBF-kernel

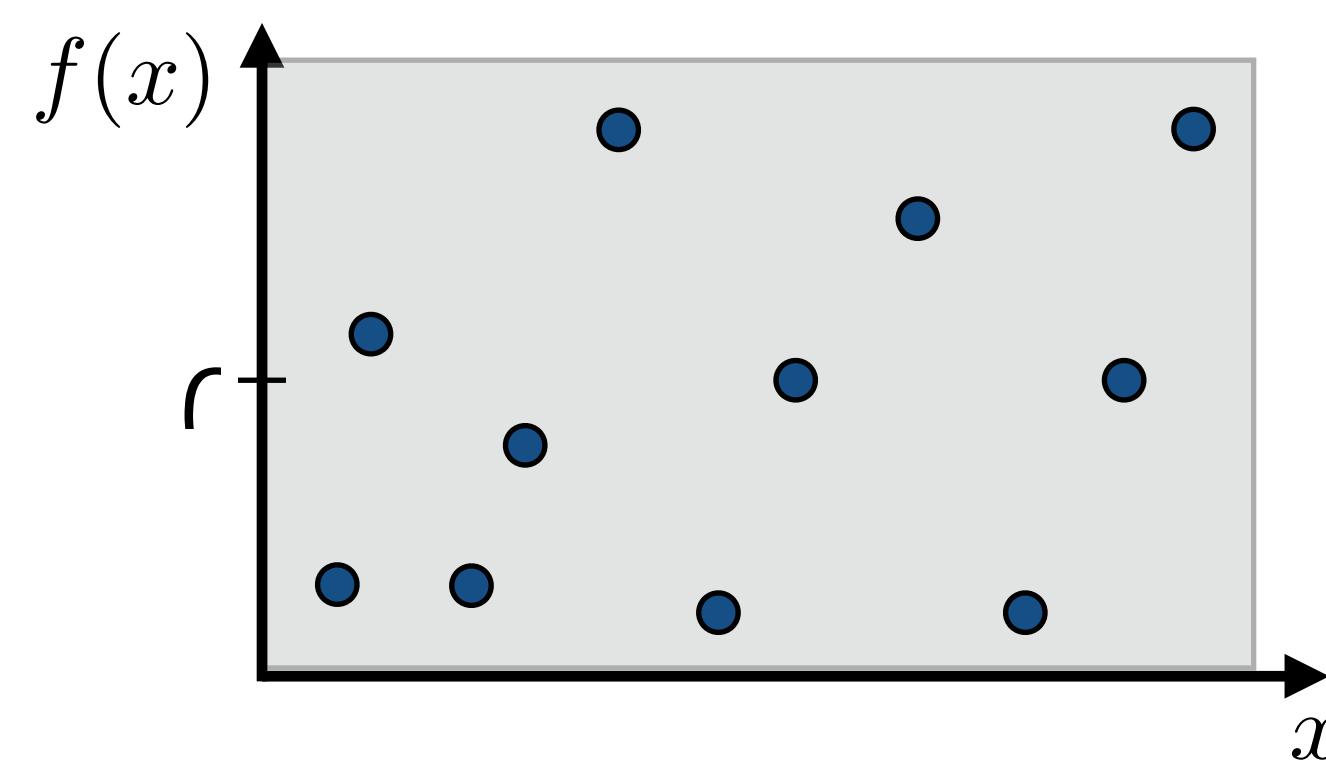
ℓ^2 defines the frequency of fluctuations



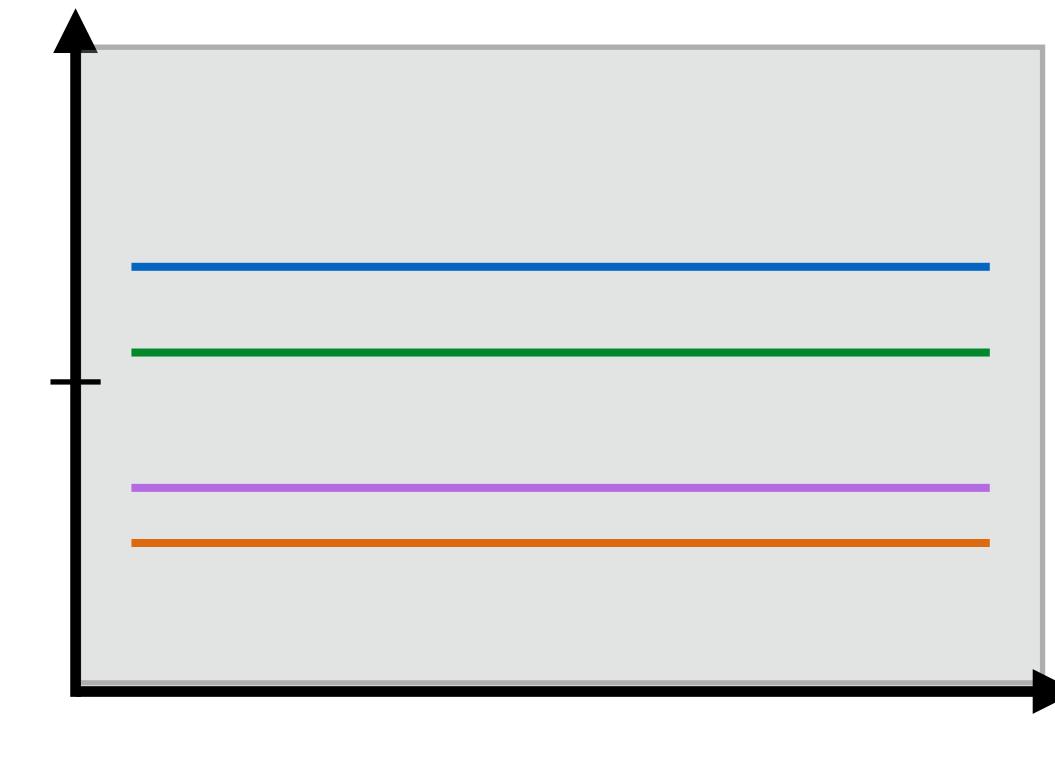
More kernels



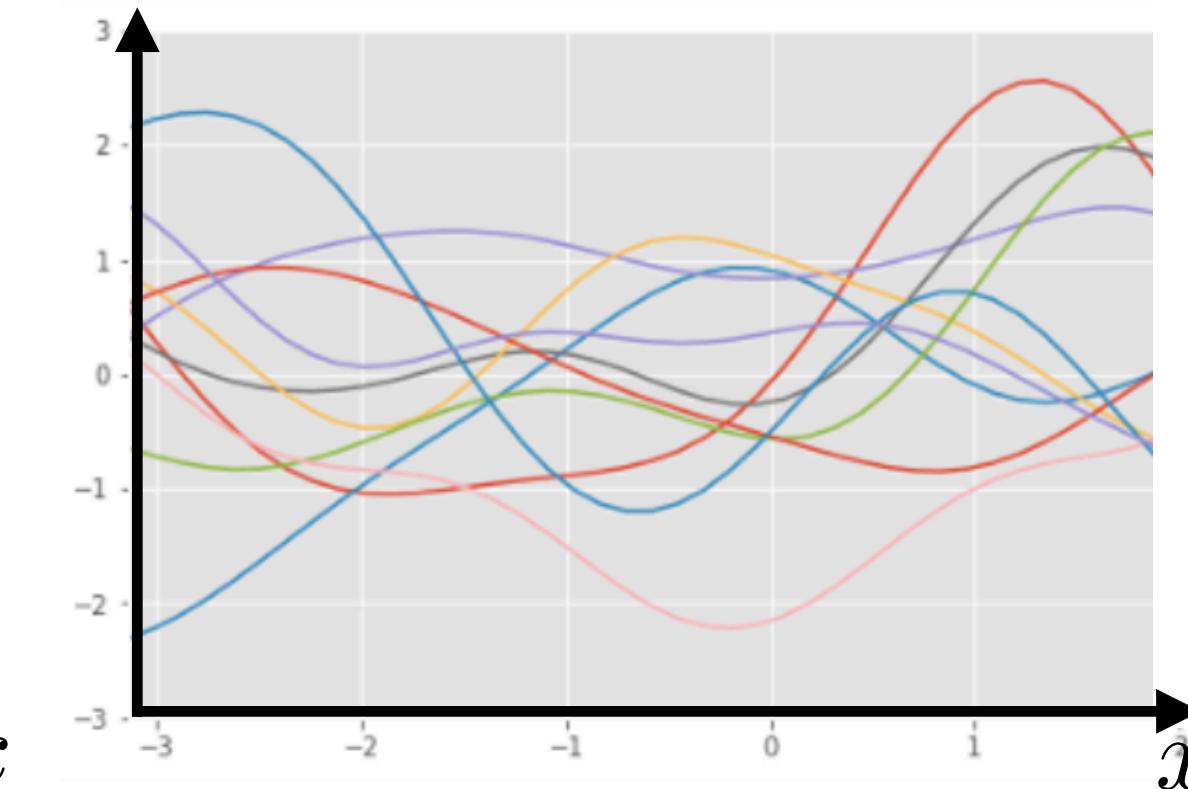
Sum-kernel



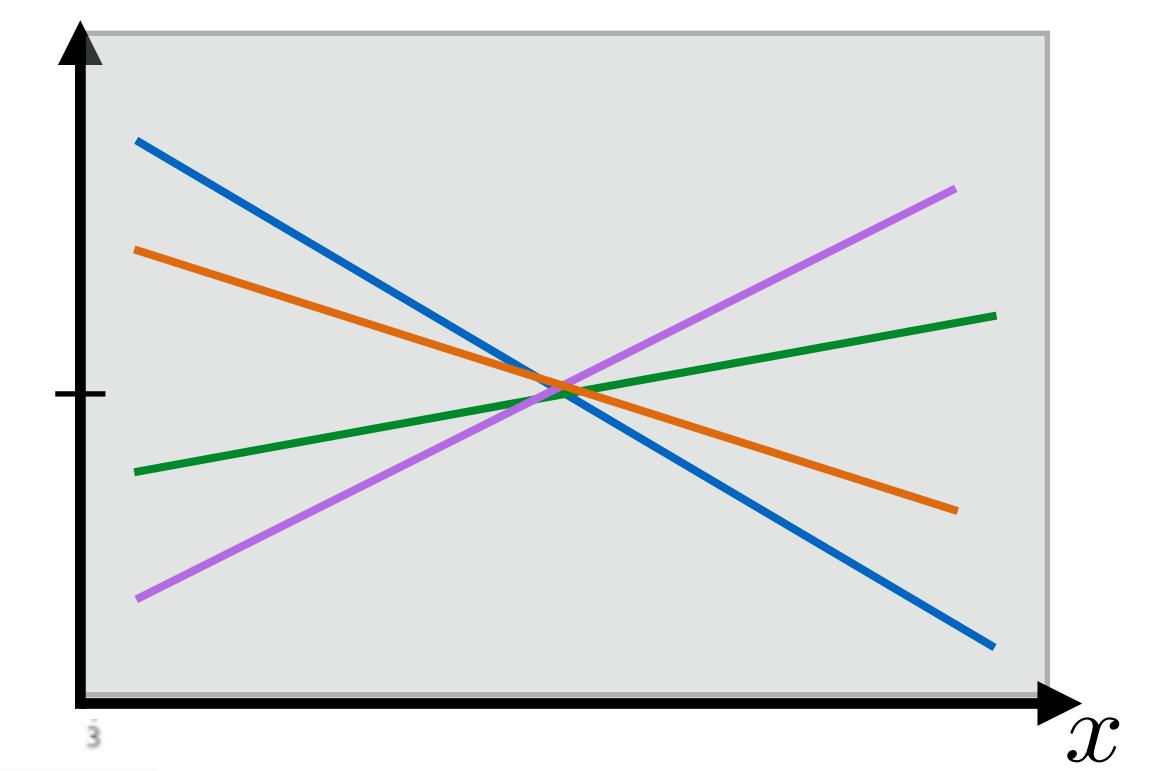
$$k(x, x') = \sigma^2[x = x']$$



$$k(x, x') = C$$



$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2r^2}\right)$$

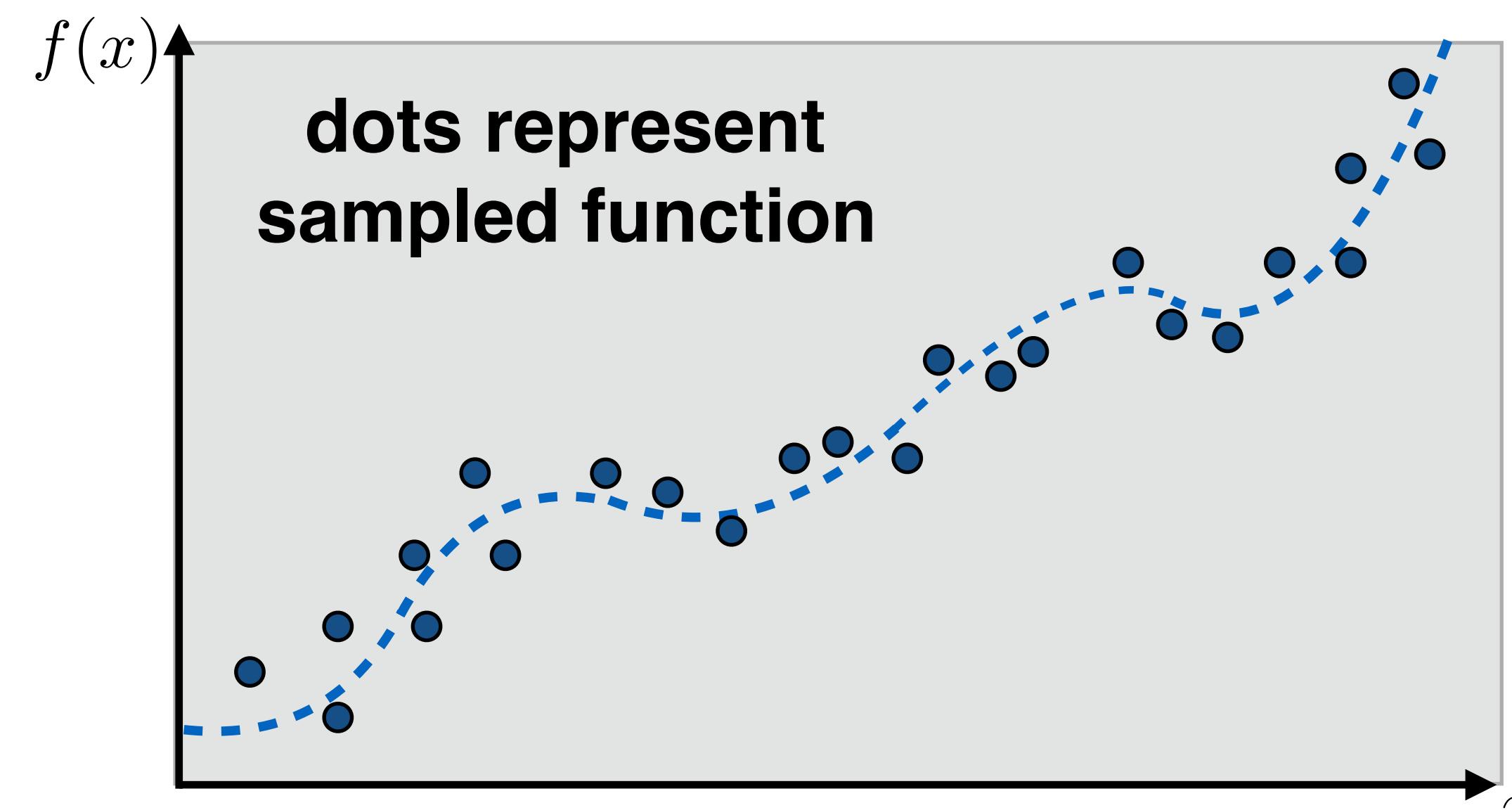


$$k(x, x') = x \cdot x'$$

Sum-kernel (multidimensional case):

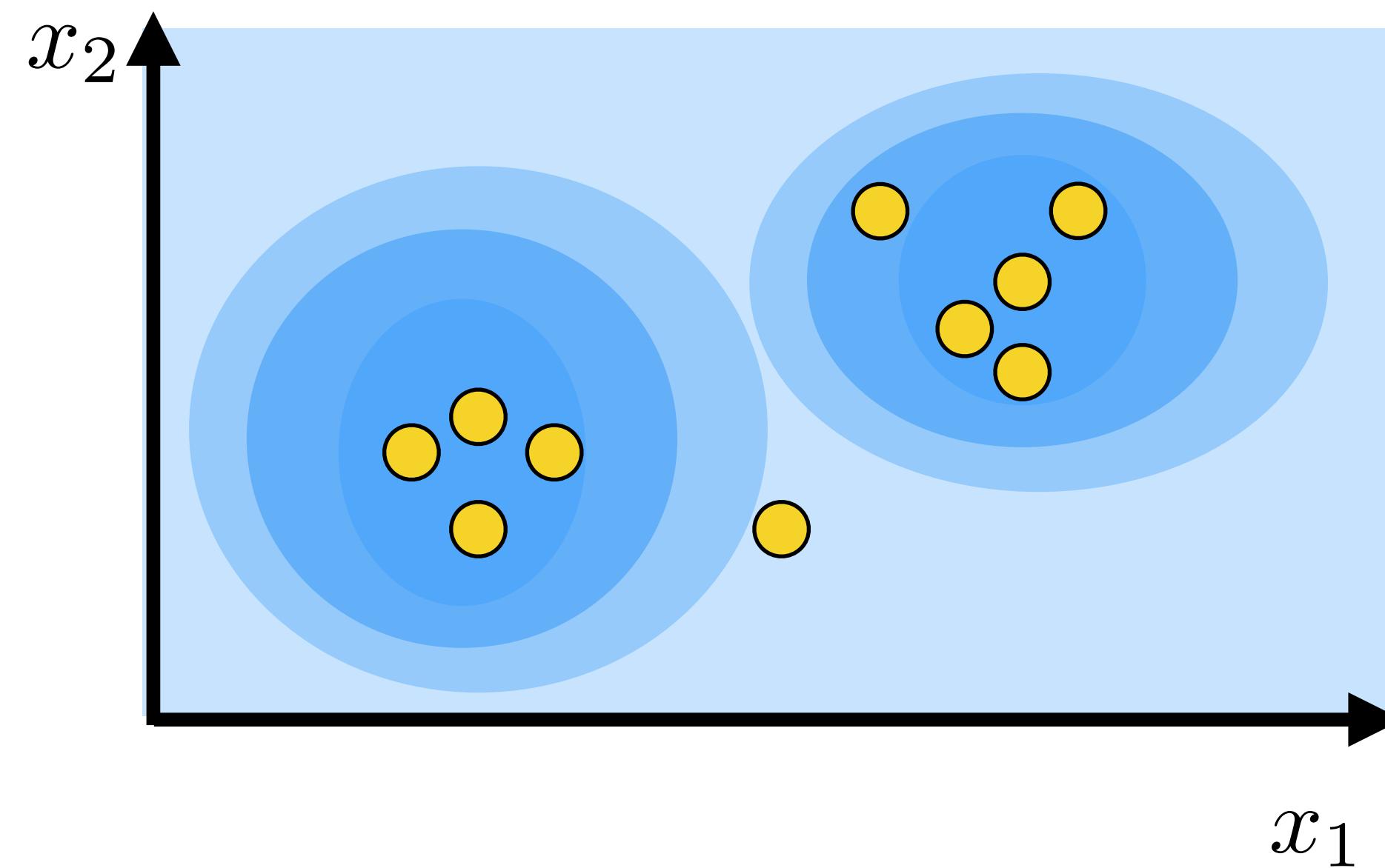
$$k(x, x') = x^T x' + \sigma_1^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right) + \sigma_2^2[x = x'] + \sigma_3^2$$

$x, x' \in \mathbb{R}^d$, d – number of features



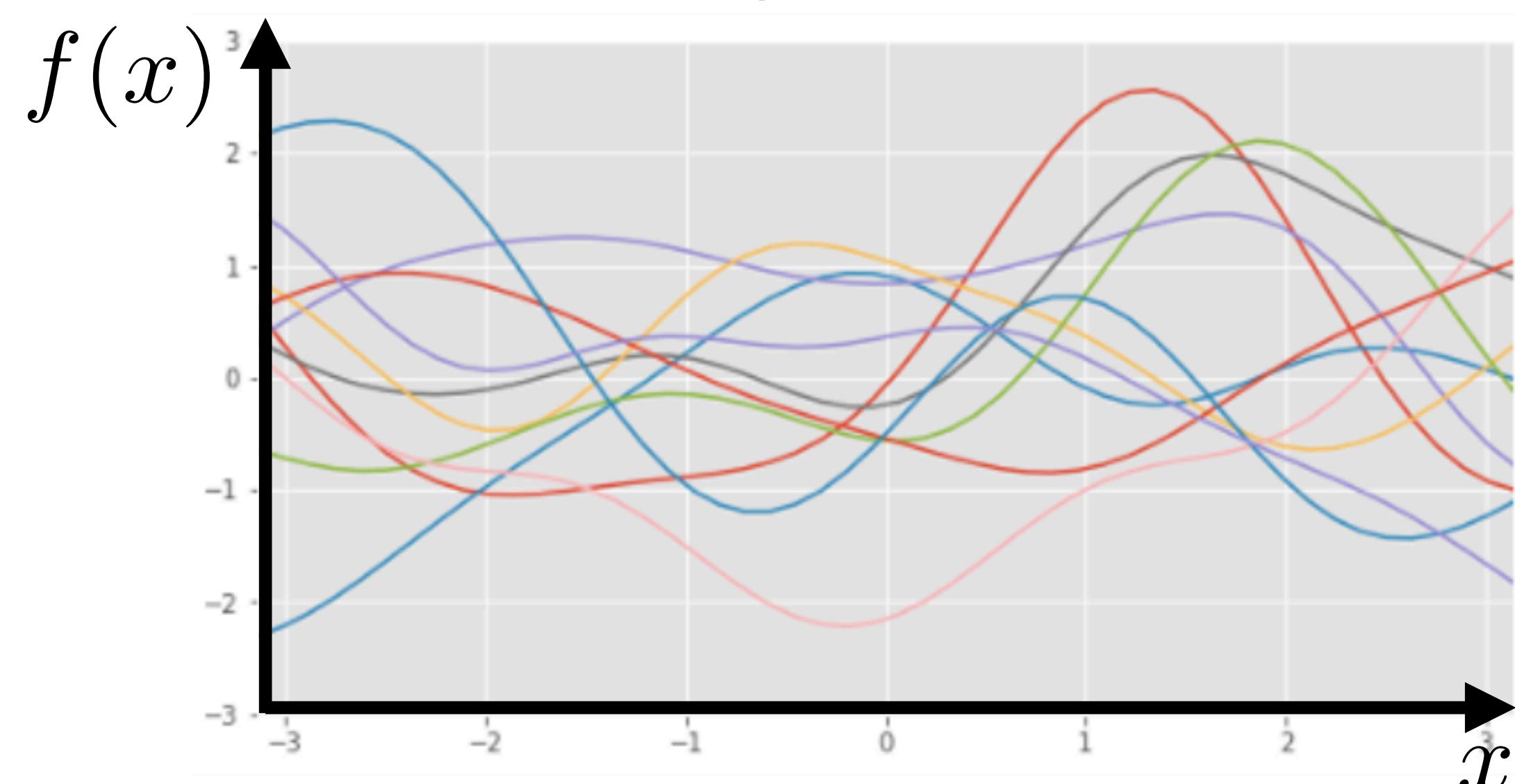
Sampling functions from a process

Multivariate distribution
— sample **points**



$$x \sim p(x)$$

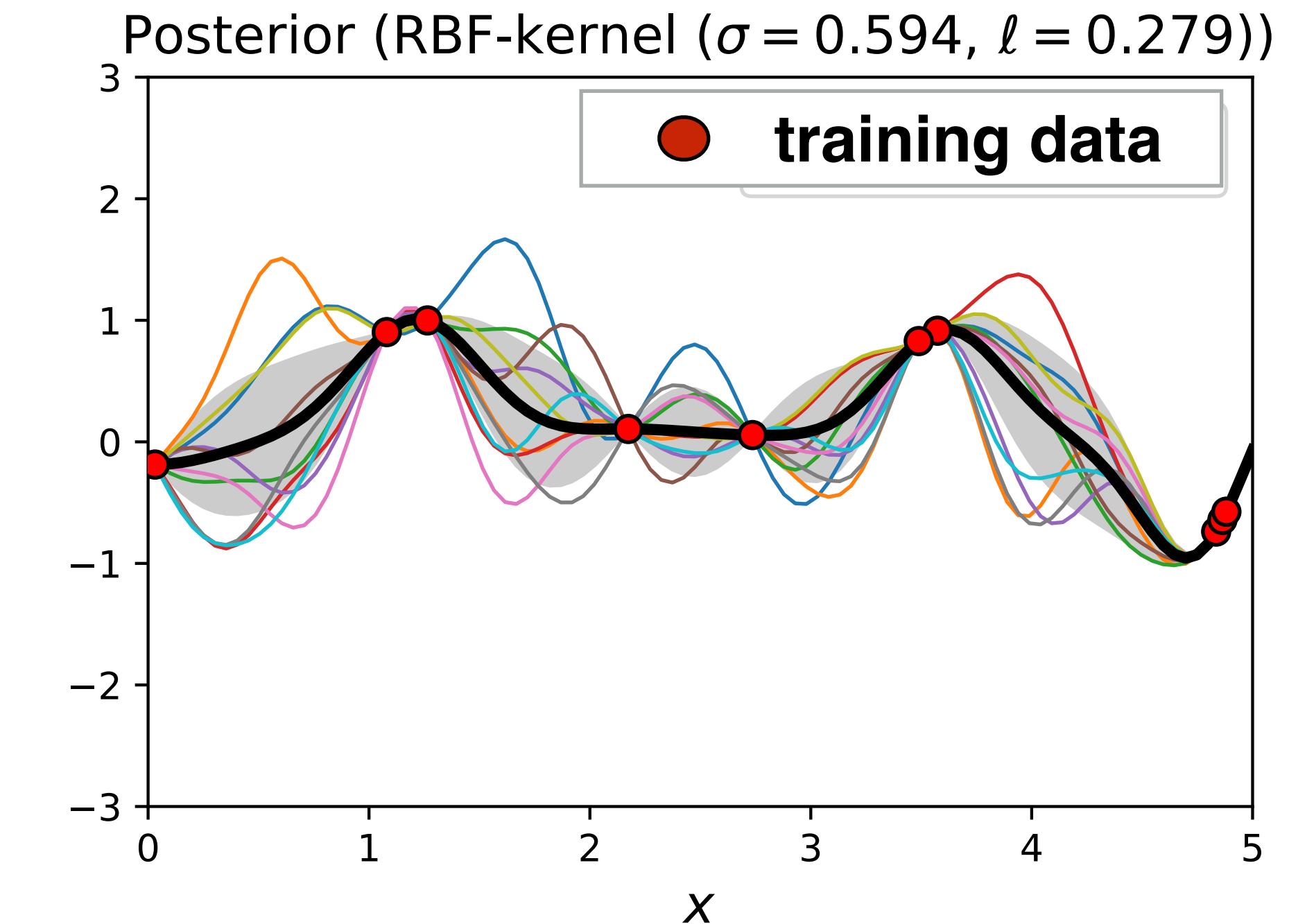
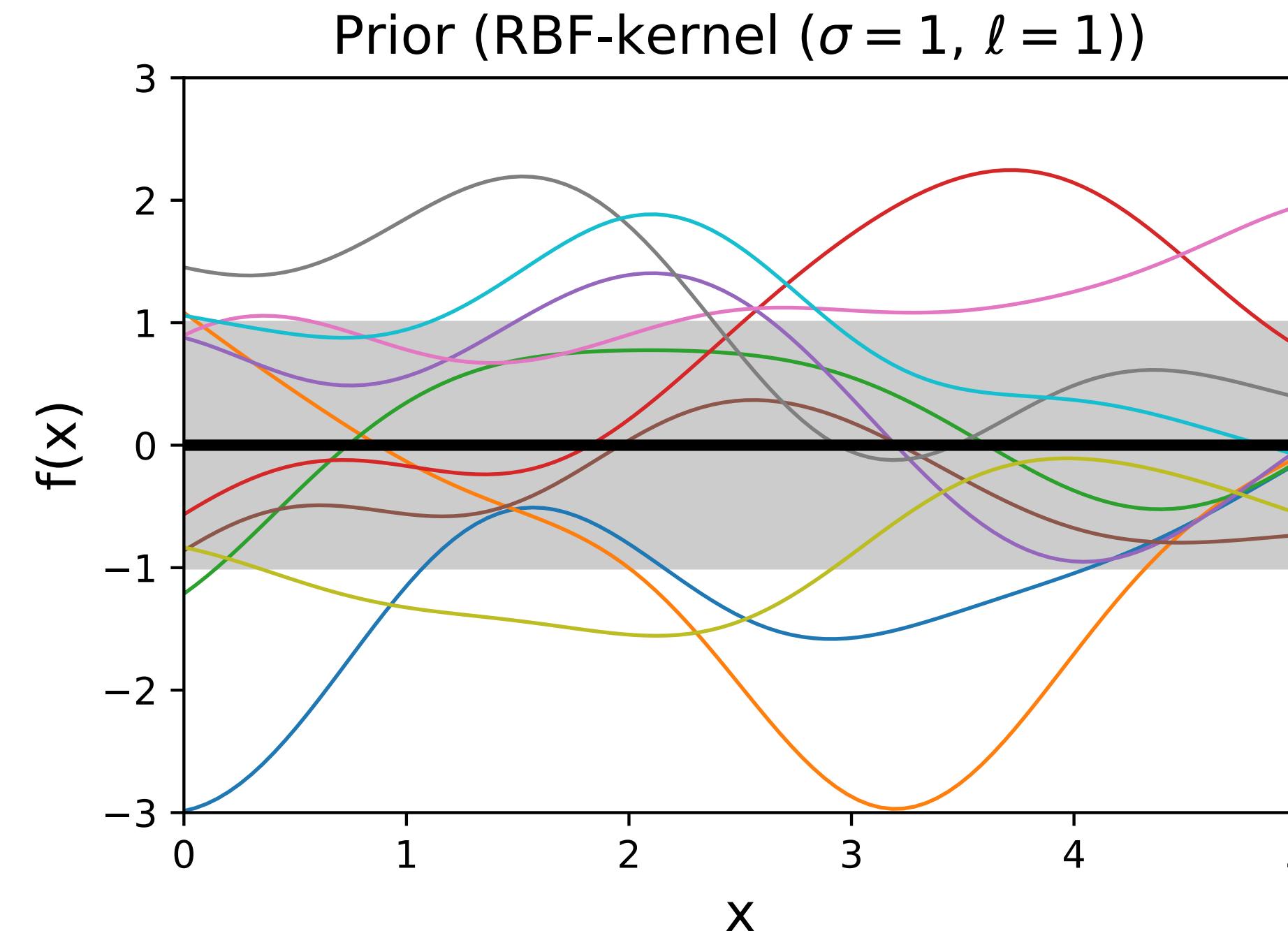
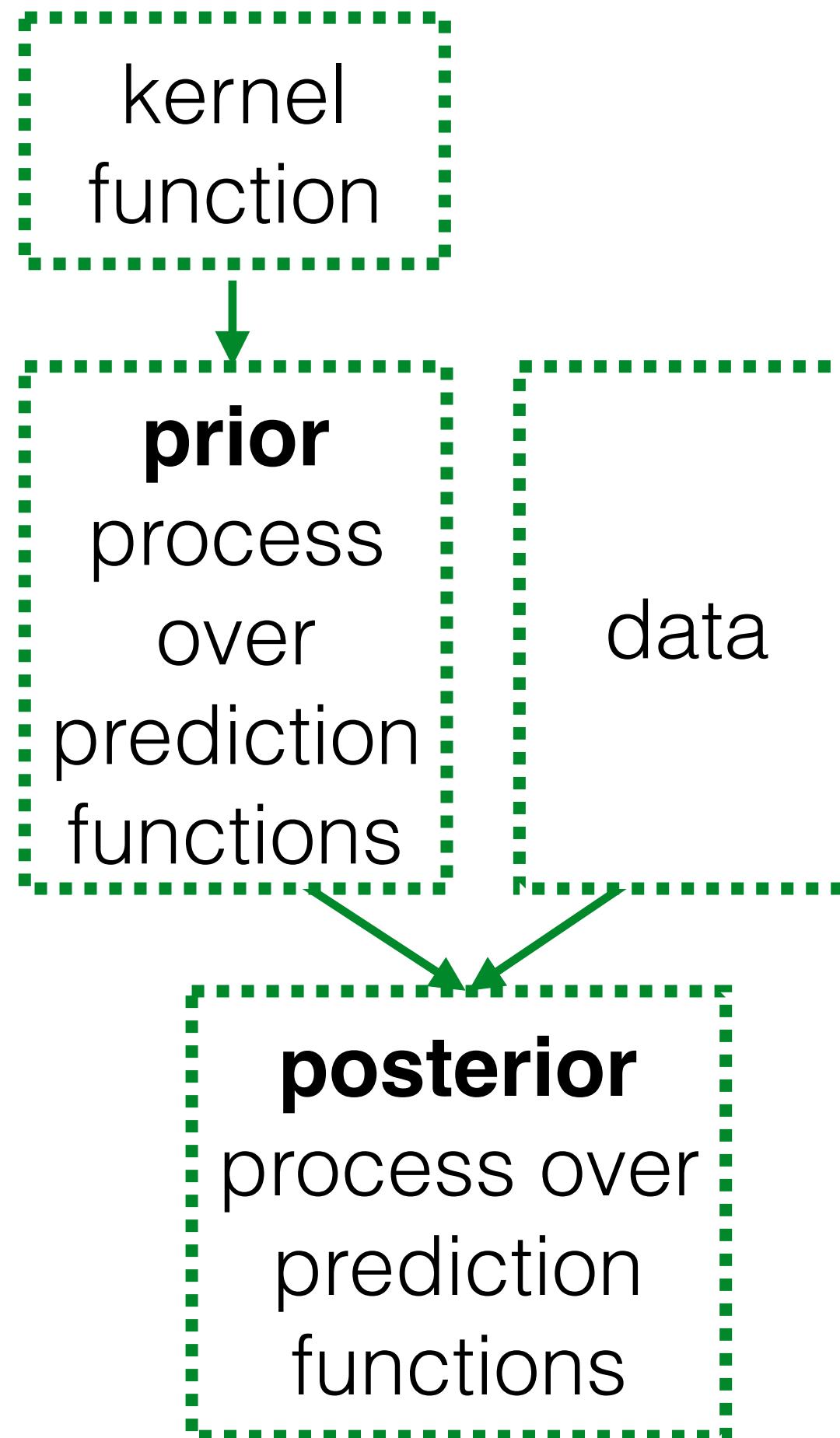
Process
— sample **functions**



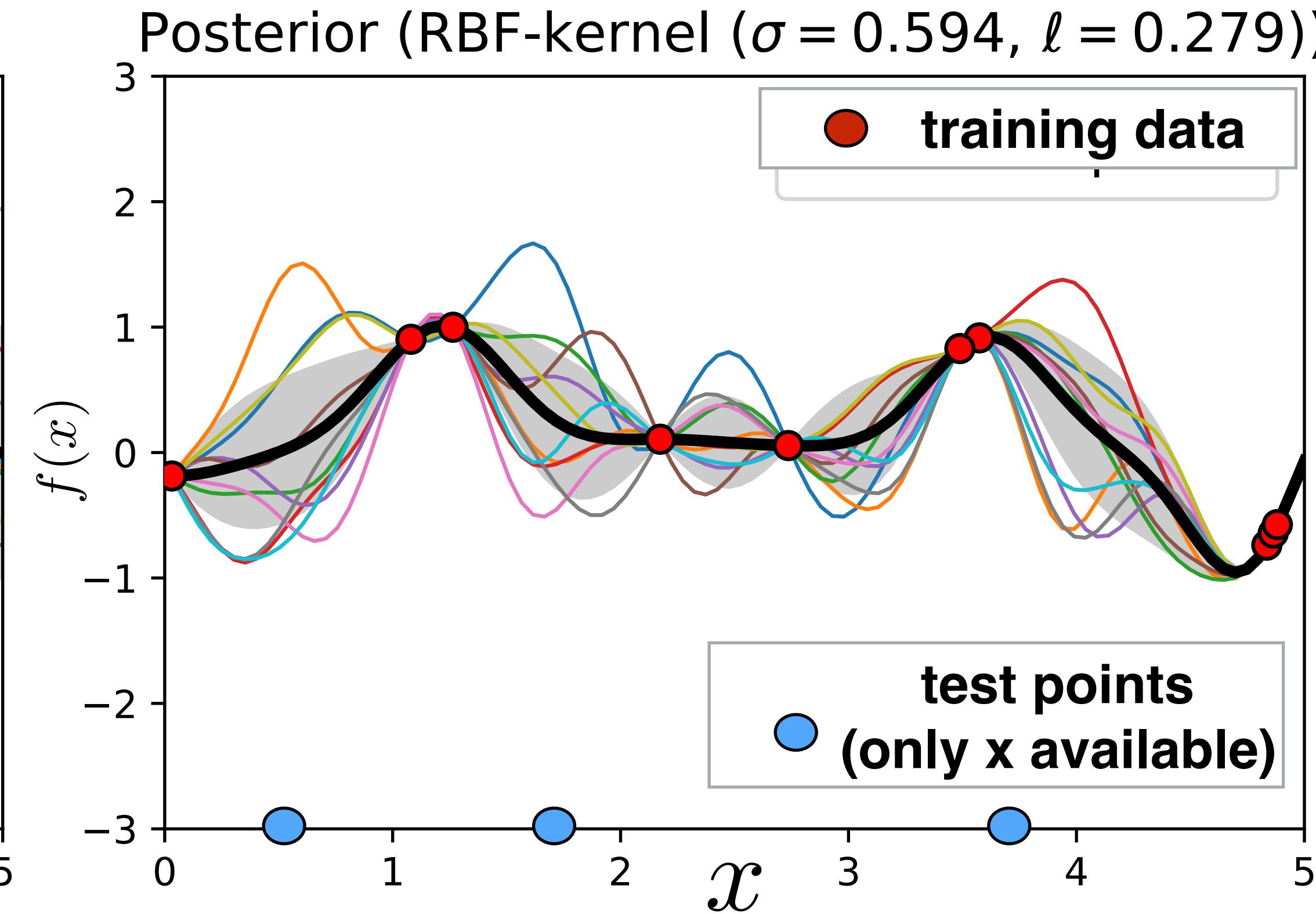
$$f(x) \sim GP(m(x), k(x, x'))$$

each line is a sample from the process

Gaussian processes for regression



Gaussian processes for regression



Given: (1) training data and
(2) prior Gaussian process over
prediction functions $a(x)$

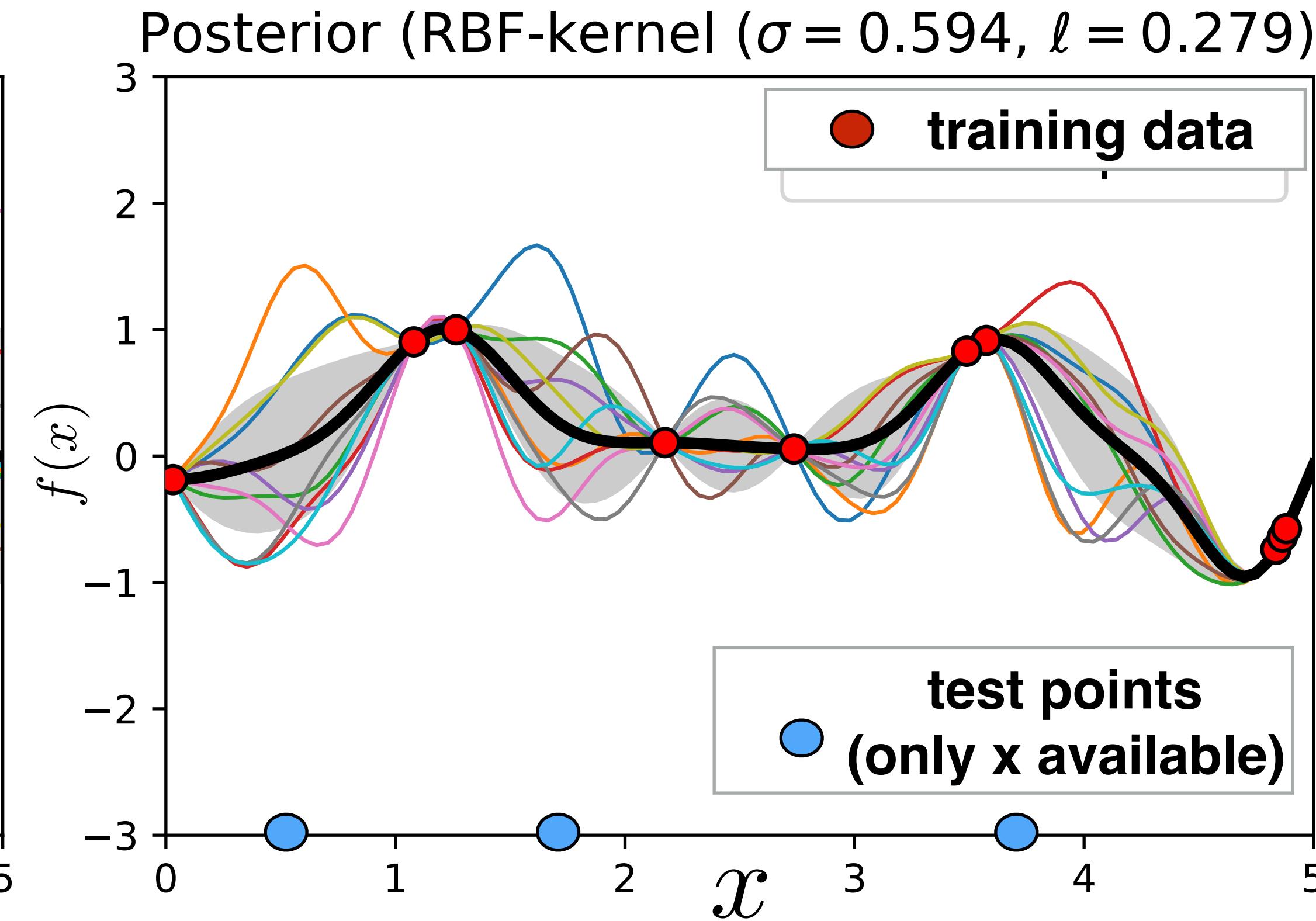
Training data ● ● ● :

$X^{tr} = \{x_i^{tr}\}_{i=1}^N, x_i^{tr} \in \mathbb{R}^d$ — input data

$Y^{tr} = \{y_i^{tr}\}_{i=1}^N, y_i^{tr} \in \mathbb{R}$ — targets

N — number of objects, d — number of features

Gaussian processes for regression



Given: (1) training data and
(2) prior Gaussian process over
prediction functions $a(x)$

Training data ● ● ● :

$X^{tr} = \{x_i^{tr}\}_{i=1}^N, x_i^{tr} \in \mathbb{R}^d$ — input data

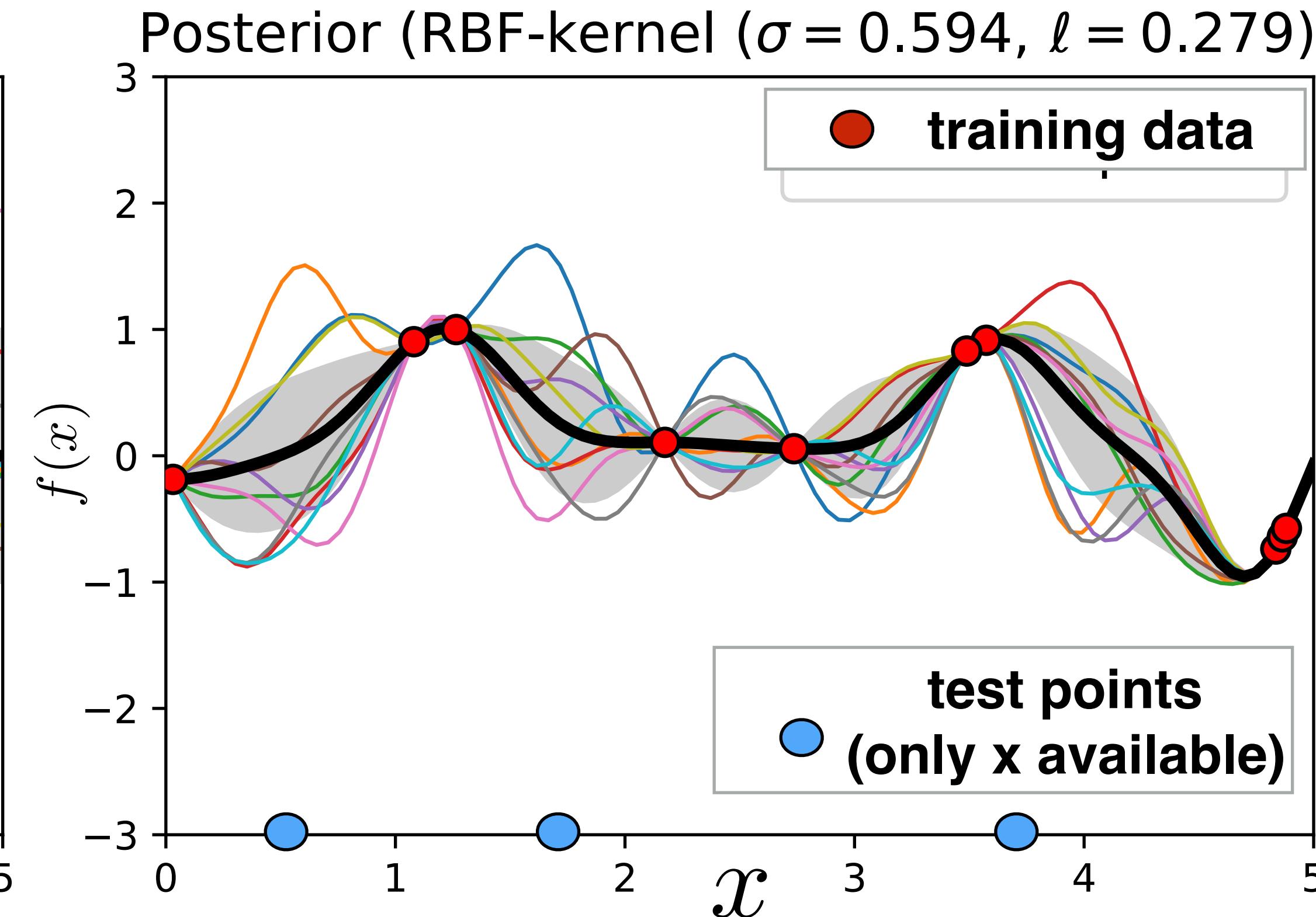
$Y^{tr} = \{y_i^{tr}\}_{i=1}^N, y_i^{tr} \in \mathbb{R}$ — targets

N — number of objects, d — number of features

Test points ● ● ● (any set of points):

$X^{te} = \{x_i^{te}\}_{i=1}^M, x_i^{te} \in \mathbb{R}^d$

Gaussian processes for regression



Given: (1) training data and
(2) prior Gaussian process over
prediction functions $a(x)$

Training data ● ● ● :

$X^{tr} = \{x_i^{tr}\}_{i=1}^N, x_i^{tr} \in \mathbb{R}^d$ — input data

$Y^{tr} = \{y_i^{tr}\}_{i=1}^N, y_i^{tr} \in \mathbb{R}$ — targets

N — number of objects, d — number of features

Test points ● ● ● (any set of points):

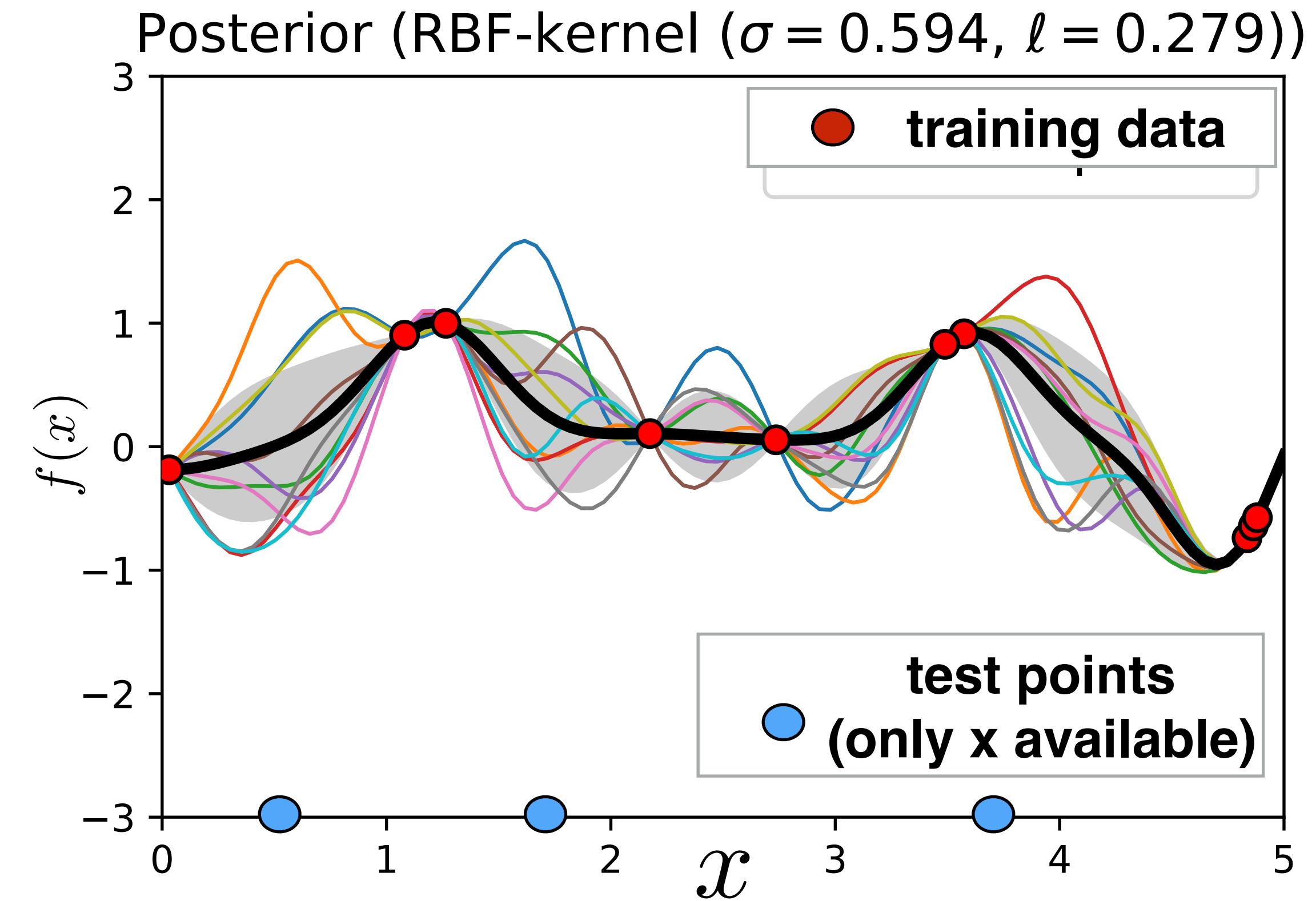
$X^{te} = \{x_i^{te}\}_{i=1}^M, x_i^{te} \in \mathbb{R}^d$

Find:

$p(a(x_1^{te}), \dots, a(x_M^{te})) - ?$

p(●●●)

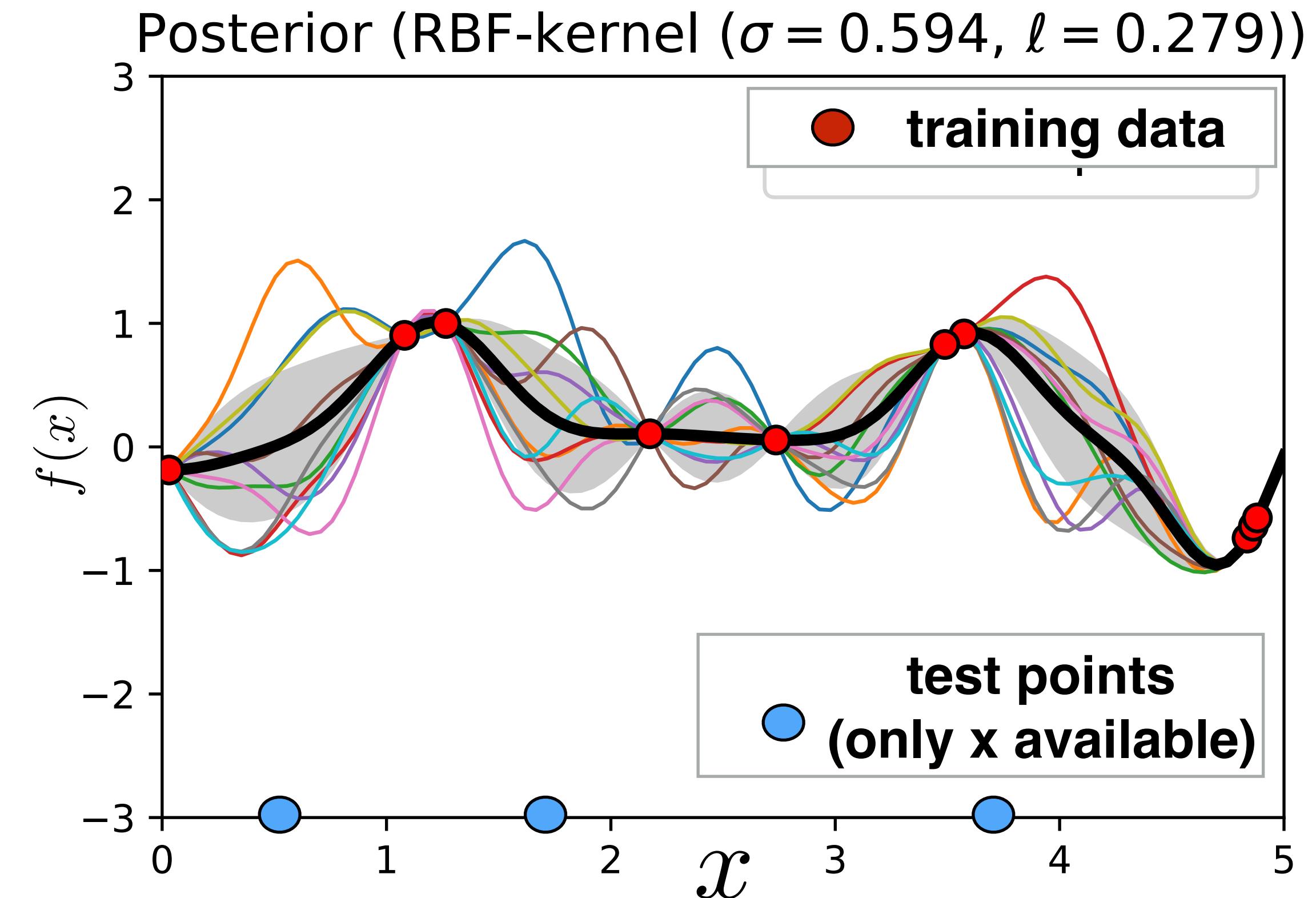
Conditioning in multivariate normal distribution



Definition of Gaussian process:
every finite set of function values
has a multivariate normal distribution

$$\forall n \quad \forall (x_1, \dots, x_n) \quad (a(x_1), \dots, a(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$
$$\mu = \{m(x_i)\}_{i=1}^n \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$

Conditioning in multivariate normal distribution

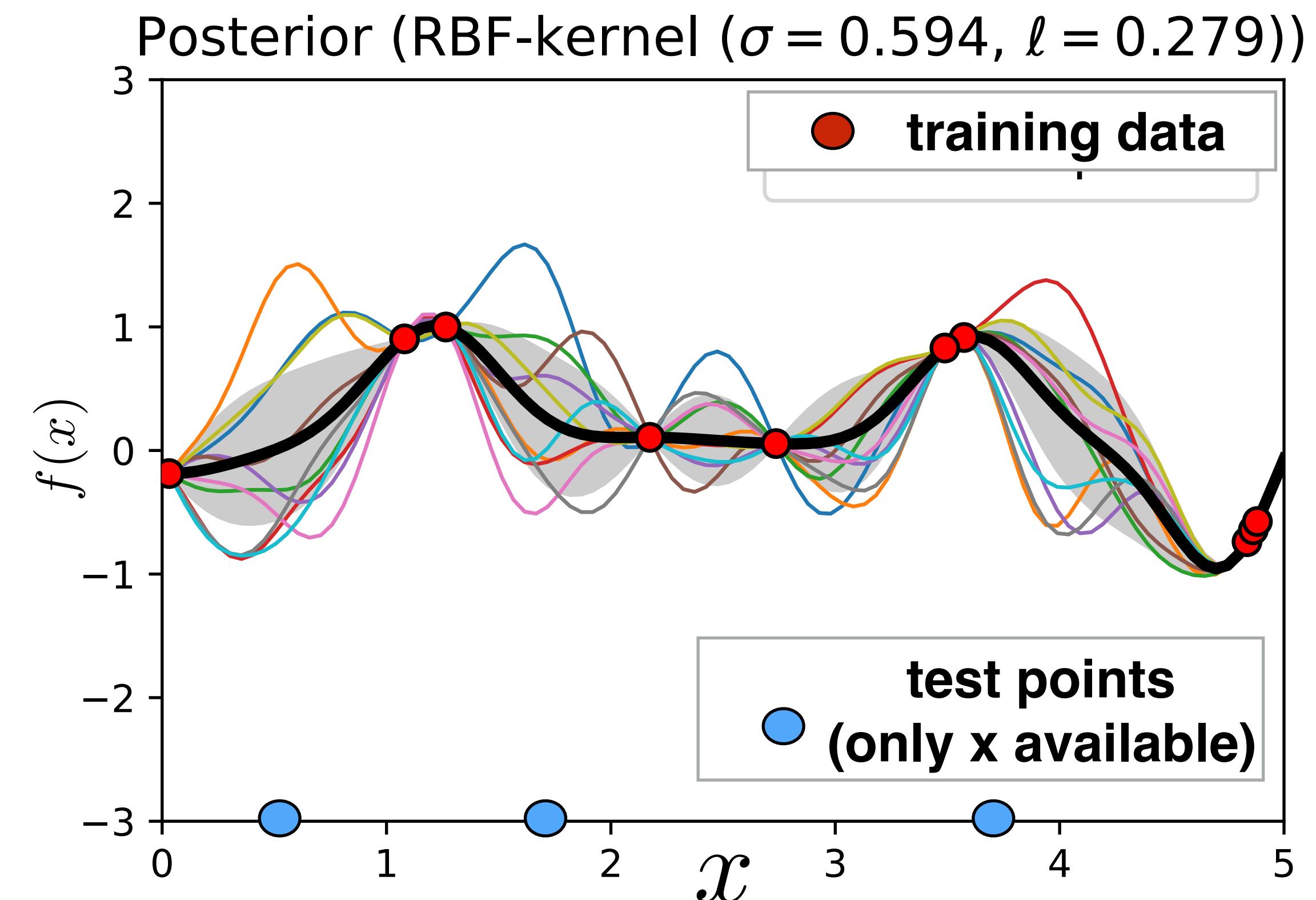


Definition of Gaussian process:
every finite set of function values
has a multivariate normal distribution

$$\forall n \quad \forall (x_1, \dots, x_n) \quad (a(x_1), \dots, a(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$
$$\mu = \{m(x_i)\}_{i=1}^n \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$

$$(\bullet \bullet \bullet \bullet \bullet \bullet) \sim \mathcal{N}\left(\begin{matrix} 0 \\ \vdots \\ 0 \end{matrix}, \begin{pmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & \ddots \end{pmatrix} \right)$$

Conditioning in multivariate normal distribution



According to properties
of normal distribution:

$$(\bullet\bullet\bullet) \sim \mathcal{N}(\begin{matrix} \text{grey box} \\ \text{red box} \end{matrix}, \begin{matrix} \text{grey box} \\ \text{red box} \end{matrix}^{-1})$$

defines a new
mean function

Definition of Gaussian process:
every finite set of function values
has a multivariate normal distribution

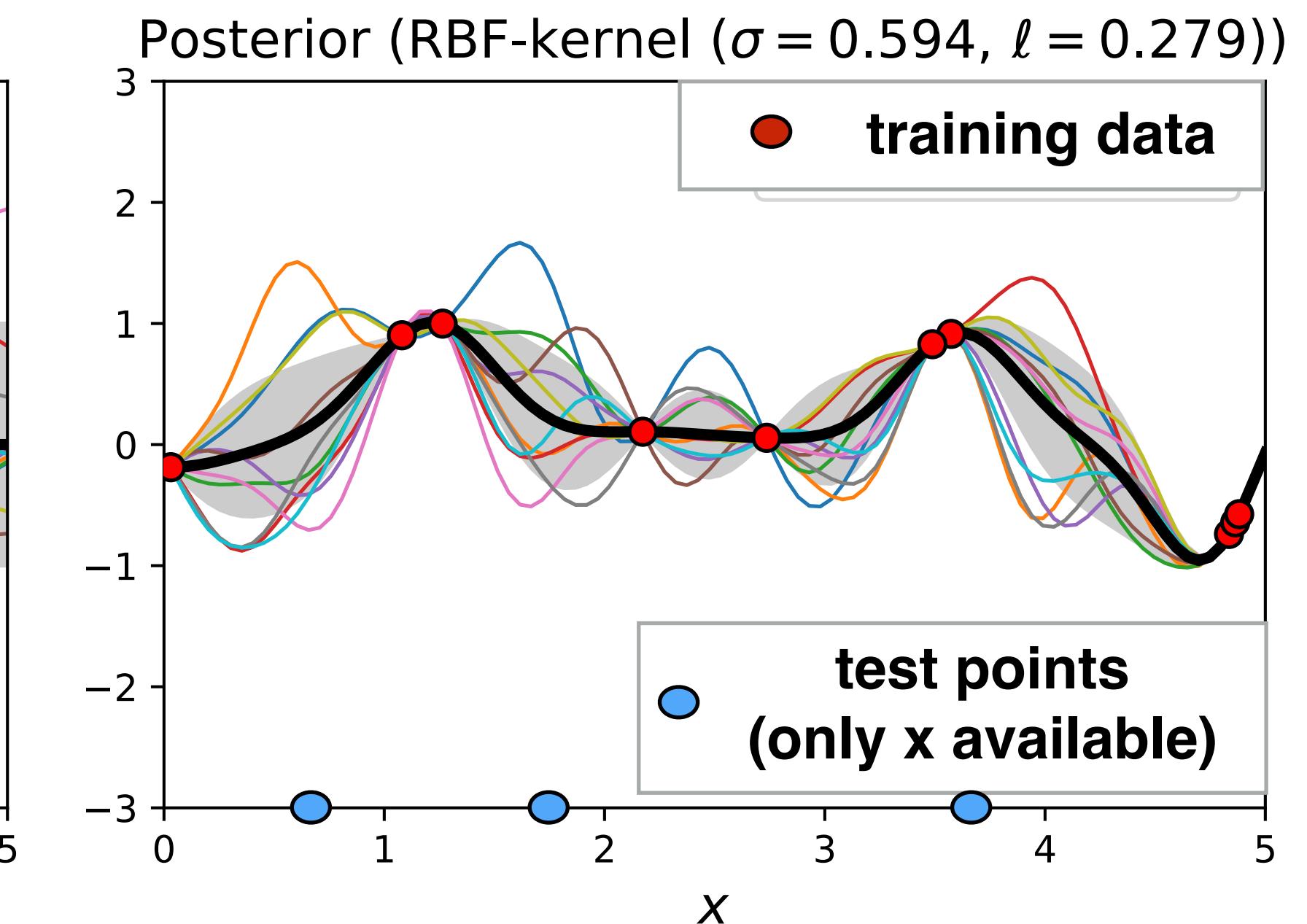
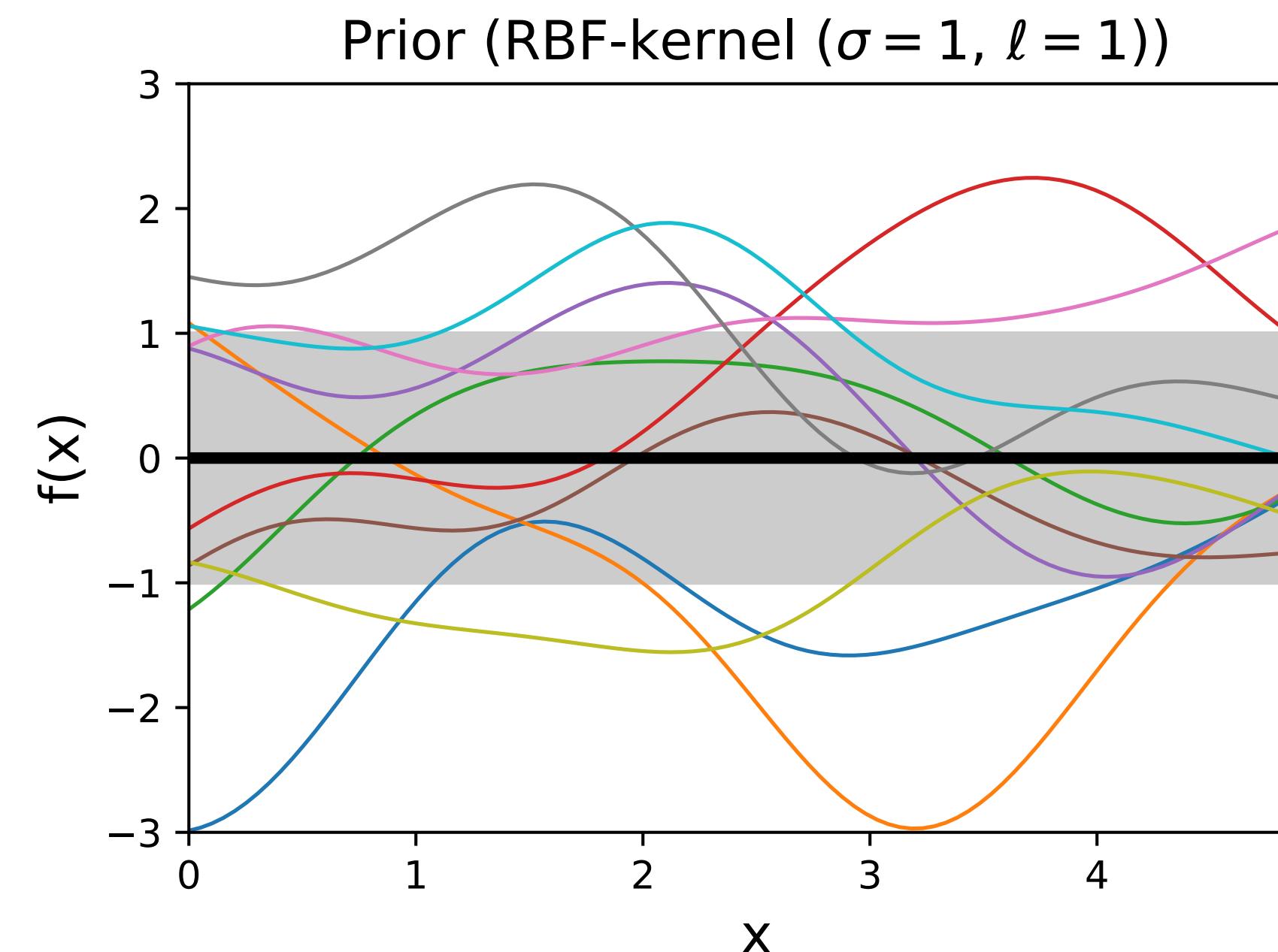
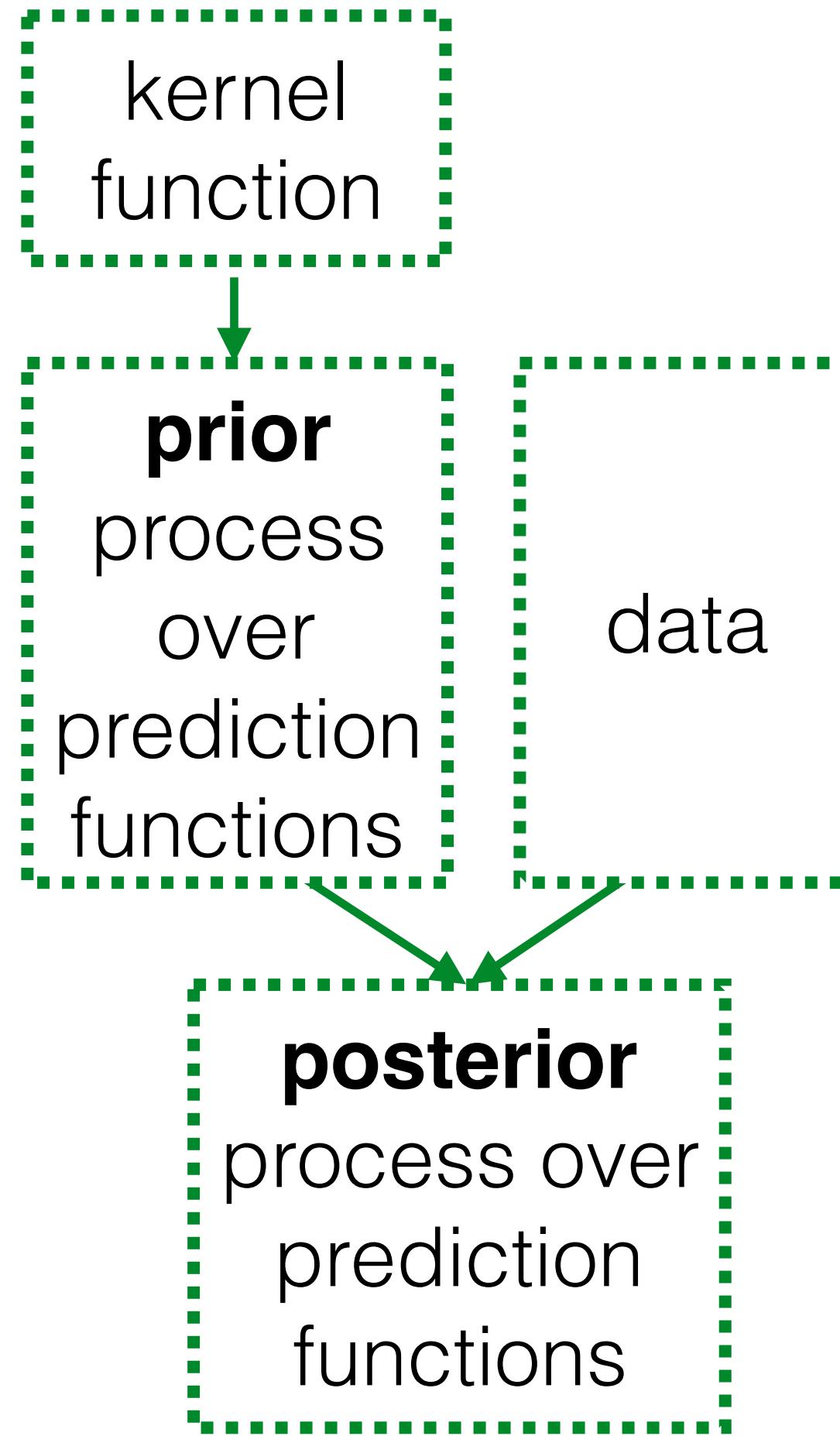
$$\forall n \quad \forall (x_1, \dots, x_n) \quad (a(x_1), \dots, a(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$
$$\mu = \{m(x_i)\}_{i=1}^n \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$

$$(\bullet\bullet\bullet \bullet\bullet\bullet) \sim \mathcal{N}(\begin{matrix} \text{grey box} \\ 0 \end{matrix}, \begin{matrix} \text{red box} & \text{grey box} \\ \text{grey box} & \text{blue box} \end{matrix})$$

$$\begin{matrix} \text{blue box} \\ - \end{matrix} \begin{matrix} \text{grey box} \\ \text{red box} \end{matrix} \begin{matrix} \text{grey box} \\ \text{blue box} \end{matrix} \sim \mathcal{N}(\begin{matrix} \text{grey box} \\ \text{red box} \end{matrix}, \begin{matrix} \text{grey box} \\ \text{red box} \end{matrix}^{-1})$$

defines a new
covariance function

Gaussian processes for regression

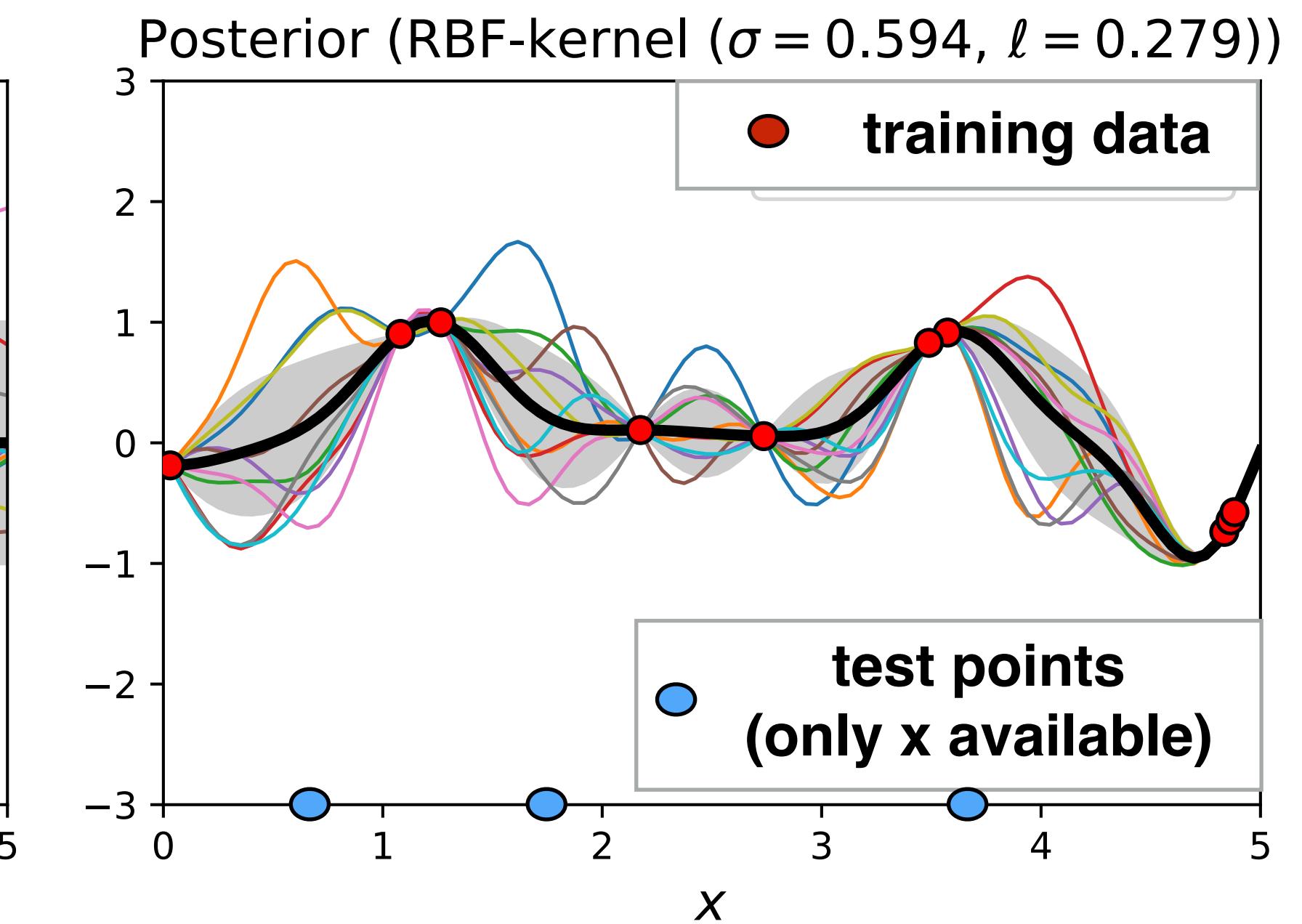
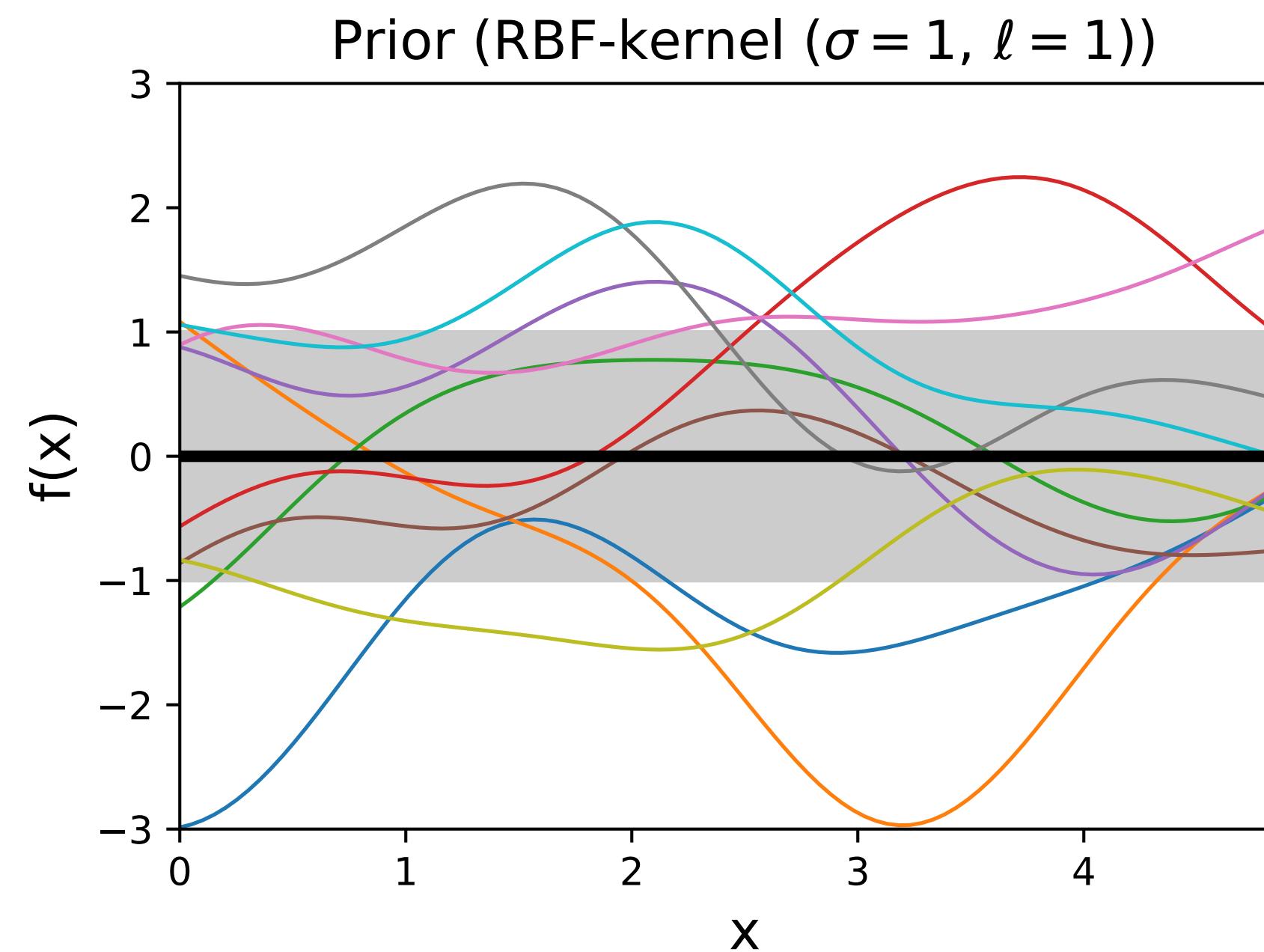
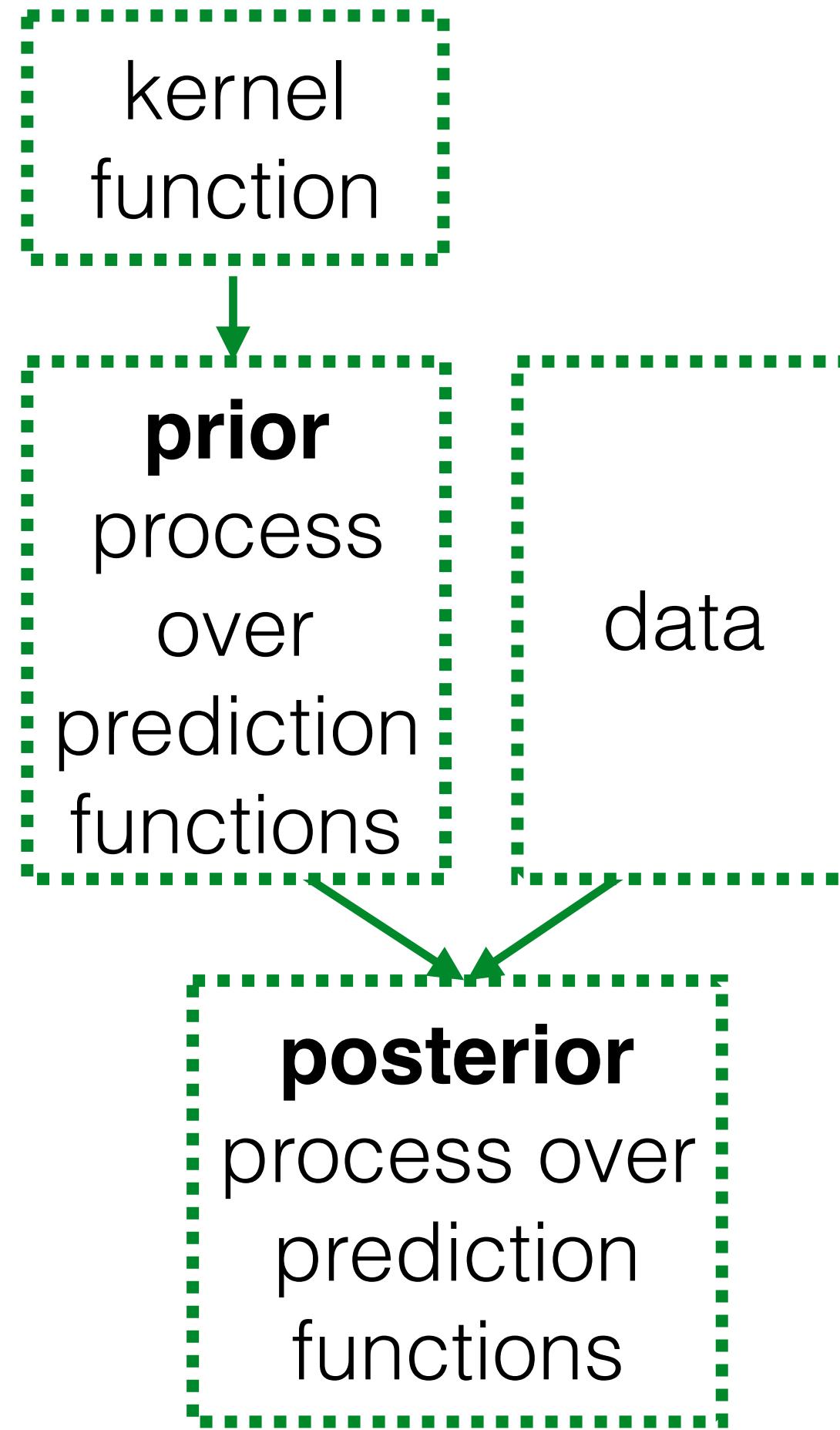


$$p(\bullet\bullet\bullet) = \mathcal{N} \left(\begin{matrix} \text{[grey box]} & \text{[red box]} \\ \text{[blue box]} & -1 \end{matrix}, \begin{matrix} \text{[grey box]} & \text{[red box]} \\ \text{[blue box]} & -1 \end{matrix} \right)$$

$\text{[grey box]} = k(X^{te}, X^{tr})$ $\text{[red box]} = k(X^{tr}, X^{tr})$ $\text{[blue box]} = k(X^{te}, X^{te})$

$\bullet\bullet\bullet = a(X^{tr}) = Y^{tr}$ $\bullet\bullet\bullet = a(X^{te})$

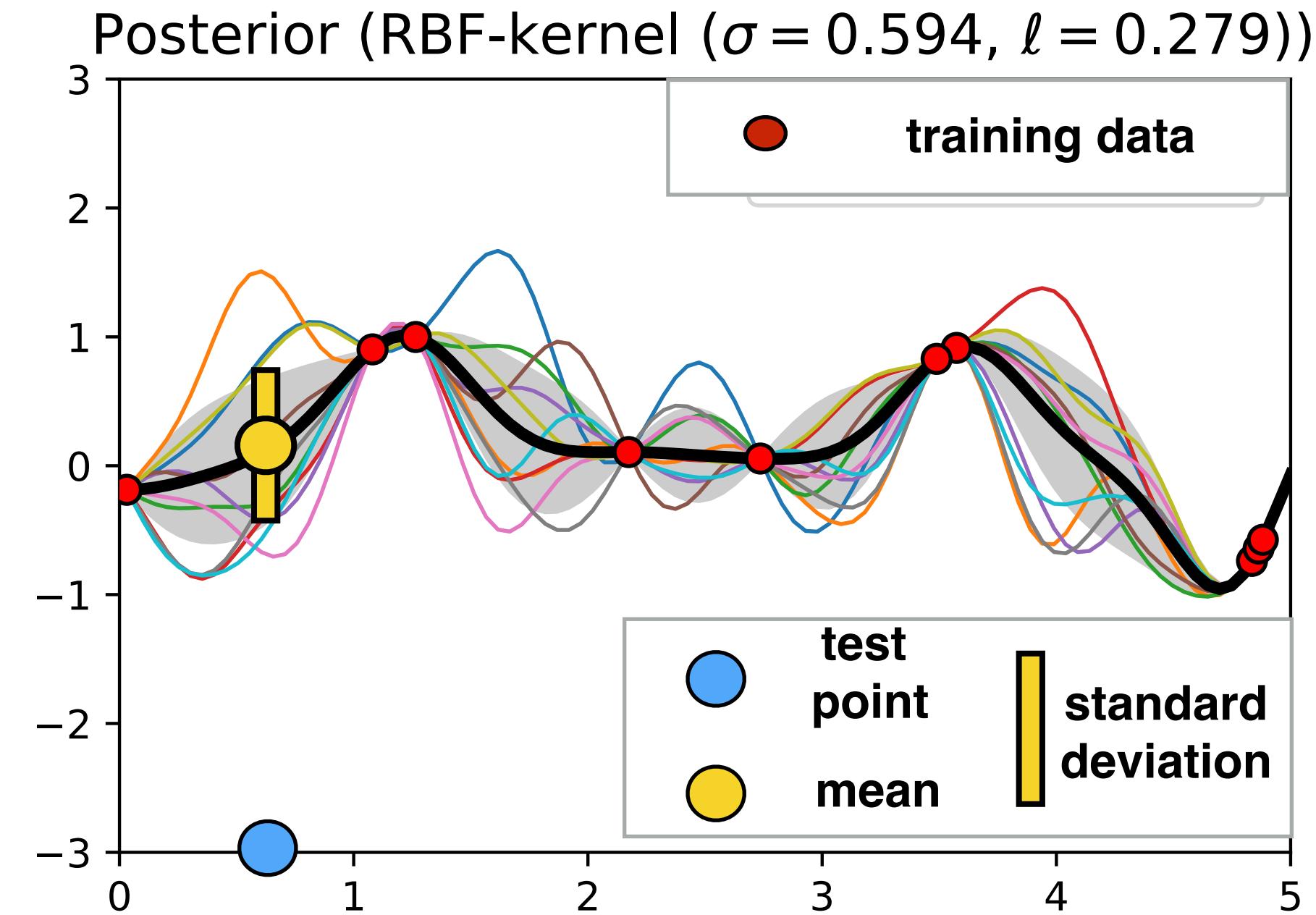
Gaussian processes for regression



$$p(\bullet \bullet \bullet) = \mathcal{N} \left(\begin{matrix} \text{grey box} \\ \text{red box} \end{matrix}, \begin{matrix} -1 \\ , \quad \text{blue box} - \begin{matrix} \text{grey box} \\ \text{red box} \end{matrix} \end{matrix} \right)$$

Training? Prediction?

Training and prediction in GP for regression



Prediction:

$$p(\bullet) = \mathcal{N}(\text{mean}, \text{variance})$$

x

-1
-1

mean variance

$\boxed{}$	$= k(x_*, X^{tr})$	$\boxed{}$	$= k(X^{tr}, X^{tr})$	$\boxed{}$	$= k(x_*, x_*)$
$\bullet \bullet \bullet$	$= a(X^{tr}) = Y^{tr}$			\bullet	$= a(x_*)$

Training and prediction in GP for regression

Training:

$$p(\bullet\bullet\bullet\bullet) = \mathcal{N}(\begin{array}{|c|}\hline 0 \\ \hline\end{array}, \begin{array}{|c|}\hline \textcolor{red}{\bullet\bullet\bullet\bullet} \\ \hline\end{array}) \rightarrow \max_{\sigma_1, \sigma_2, \sigma_3, \ell}$$

$$k(x, x') = x^T x' + \sigma_1^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right) + \sigma_2^2[x = x'] + \sigma_3^2$$

parameters of the el (covariance) function

Prediction:

$$p(\bullet) = \mathcal{N}(\text{mean}, \text{variance})$$

$$\begin{array}{ccc} \text{[Gray Box]} = k(x_*, X^{tr}) & \text{[Red Box]} = k(X^{tr}, X^{tr}) & \text{[Blue Box]} = k(x_*, x_*) \\ \bullet\bullet\bullet\bullet = a(X^{tr}) = Y^{tr} & & \bullet = a(x_*) \end{array}$$

Training and prediction in GP for regression

$$p(a(x_*)) = \mathcal{N}(m_*, \sigma_*) \quad m_* = k_*^T K^{-1} Y, \quad \sigma_*^2 = k_{**} - k_*^T K^{-1} k_* \\ k_{**} = k(x_*, x_*), \quad k_* = \{k(x_i, x_*)\}_{i=1}^N, \quad K = \{k(x_i, x_j)\}_{i,j=1}^{N,N}, \quad Y = \{y_i\}_{i=1}^N$$

Prediction:

$$p(\bullet) = \mathcal{N}\left(\begin{array}{c} \text{mean} \\ \hline \text{---} \\ \text{variance} \end{array}, \begin{array}{c} \text{---} \\ \hline \text{---} \end{array}\right)$$

$$\begin{array}{ccc} \text{---} & = k(x_*, X^{tr}) & \text{---} = k(X^{tr}, X^{tr}) & \text{---} = k(x_*, x_*) \\ \bullet \bullet \bullet \bullet & = a(X^{tr}) = Y^{tr} & \bullet & = a(x_*) \end{array}$$

Exercise

- Consider we are given the following training data (1 feature):

x	y
-1.5	1
0.5	3
0.7	2.5

We use zero mean function and RBF-kernel: $k(x, x') = 0.5 \exp\left(-\frac{(x - x')^2}{2}\right)$

- What prediction will we make for a new object $x_* = 0?$ for $x_* = 3?$

Formulas:

$$p(a(x_*)) = \mathcal{N}(m_*, \sigma_*), \quad m_* = k_*^T K^{-1} Y, \quad \sigma_*^2 = k_{**} - k_*^T K^{-1} k_*$$
$$k_{**} = k(x_*, x_*), \quad k_* = \{k(x_i, x_*)\}_{i=1}^N, \quad K = \{k(x_i, x_j)\}_{i,j=1}^{N,N}, \quad Y = \{y_i\}_{i=1}^N$$

Exercise

x	y
-1.5	1
0.5	3
0.7	2.5

$$K = 0.5 \cdot \begin{bmatrix} 1 & \exp\left(-\frac{(-1.5-0.5)^2}{2}\right) & \exp\left(-\frac{(-1.5-0.7)^2}{2}\right) \\ \exp\left(-\frac{(-1.5-0.5)^2}{2}\right) & 1 & \exp\left(-\frac{(0.5-0.7)^2}{2}\right) \\ \exp\left(-\frac{(-1.5-0.7)^2}{2}\right) & \exp\left(-\frac{(0.5-0.7)^2}{2}\right) & 1 \end{bmatrix}$$

Formulas:

$$\begin{aligned} p(a(x_*)) &= \mathcal{N}(m_*, \sigma_*) & m_* &= k_*^T K^{-1} Y, & \sigma_*^2 &= k_{**} - k_*^T K^{-1} k_* \\ k_{**} &= k(x_*, x_*), & k_* &= \{k(x_i, x_*)\}_{i=1}^N, & K &= \{k(x_i, x_j)\}_{i,j=1}^{N,N}, & Y &= \{y_i\}_{i=1}^N \end{aligned}$$

Exercise

x	y
-1.5	1
0.5	3
0.7	2.5

$$K = 0.5 \cdot \begin{bmatrix} 1 & \exp\left(-\frac{(-1.5-0.5)^2}{2}\right) & \exp\left(-\frac{(-1.5-0.7)^2}{2}\right) \\ \exp\left(-\frac{(-1.5-0.5)^2}{2}\right) & 1 & \exp\left(-\frac{(0.5-0.7)^2}{2}\right) \\ \exp\left(-\frac{(-1.5-0.7)^2}{2}\right) & \exp\left(-\frac{(0.5-0.7)^2}{2}\right) & 1 \end{bmatrix}$$

For $x_* = 0$:

$$k_* = 0.5 \cdot \begin{bmatrix} \exp\left(-\frac{(0+1.5)^2}{2}\right) & \exp\left(-\frac{(0-0.5)^2}{2}\right) & \exp\left(-\frac{(0-0.7)^2}{2}\right) \end{bmatrix} \quad k_{**} = [0.5]$$

Formulas:

$$\begin{aligned} p(a(x_*)) &= \mathcal{N}(m_*, \sigma_*) \quad m_* = k_*^T K^{-1} Y, \quad \sigma_*^2 = k_{**} - k_*^T K^{-1} k_* \\ k_{**} &= k(x_*, x_*), \quad k_* = \{k(x_i, x_*)\}_{i=1}^N, \quad K = \{k(x_i, x_j)\}_{i,j=1}^{N,N}, \quad Y = \{y_i\}_{i=1}^N \end{aligned}$$

Exercise

x	y
-1.5	1
0.5	3
0.7	2.5

$$K = 0.5 \cdot \begin{bmatrix} 1 & \exp\left(-\frac{(-1.5-0.5)^2}{2}\right) & \exp\left(-\frac{(-1.5-0.7)^2}{2}\right) \\ \exp\left(-\frac{(-1.5-0.5)^2}{2}\right) & 1 & \exp\left(-\frac{(0.5-0.7)^2}{2}\right) \\ \exp\left(-\frac{(-1.5-0.7)^2}{2}\right) & \exp\left(-\frac{(0.5-0.7)^2}{2}\right) & 1 \end{bmatrix}$$

For $x_* = 0$:

$$k_* = 0.5 \cdot \begin{bmatrix} \exp\left(-\frac{(0+1.5)^2}{2}\right) & \exp\left(-\frac{(0-0.5)^2}{2}\right) & \exp\left(-\frac{(0-0.7)^2}{2}\right) \end{bmatrix} \quad k_{**} = [0.5]$$

$$\mu_* = [0.162 \quad 0.441 \quad 0.391] \begin{bmatrix} 0.5 & 0.067 & 0.044 \\ 0.067 & 0.5 & 0.490 \\ 0.044 & 0.490 & 0.5 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 3 \\ 2.5 \end{bmatrix} \quad \sigma_*^2 = 0.5 - [0.162 \quad 0.441 \quad 0.391] \begin{bmatrix} 0.5 & 0.067 & 0.044 \\ 0.067 & 0.5 & 0.490 \\ 0.044 & 0.490 & 0.5 \end{bmatrix}^{-1} \begin{bmatrix} 0.162 \\ 0.441 \\ 0.391 \end{bmatrix}$$

Formulas:

$$p(a(x_*)) = \mathcal{N}(m_*, \sigma_*) \quad m_* = k_*^T K^{-1} Y, \quad \sigma_*^2 = k_{**} - k_*^T K^{-1} k_*$$

$$k_{**} = k(x_*, x_*), \quad k_* = \{k(x_i, x_*)\}_{i=1}^N, \quad K = \{k(x_i, x_j)\}_{i,j=1}^{N,N}, \quad Y = \{y_i\}_{i=1}^N$$

Training and prediction in GP for regression

Training:

$$p(\bullet\bullet\bullet\bullet) = \mathcal{N}(\begin{array}{|c|}\hline 0 \\ \hline\end{array}, \begin{array}{|c|}\hline \textcolor{red}{\bullet\bullet\bullet\bullet} \\ \hline\end{array}) \rightarrow \max_{\sigma_1, \sigma_2, \sigma_3, \ell}$$

$$k(x, x') = x^T x' + \sigma_1^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right) + \sigma_2^2[x = x'] + \sigma_3^2$$

parameters of the
el (covariance) function

Prediction:

$$p(\bullet) = \mathcal{N}(\text{mean}, \text{variance})$$

$$\begin{array}{ccc} \text{[Gray Box]} = k(x_*, X^{tr}) & \text{[Red Box]} = k(X^{tr}, X^{tr}) & \text{[Blue Box]} = k(x_*, x_*) \\ \bullet\bullet\bullet\bullet = a(X^{tr}) = Y^{tr} & & \text{[Blue Circle]} = a(x_*) \end{array}$$

Parametric vs non-parametric models

Parametric models:

Prediction: $a(x)$ – function of x and parameters θ

Training: finding θ based on training data

Examples: decision trees
(Bayesian) linear regression

Non-parametric models:

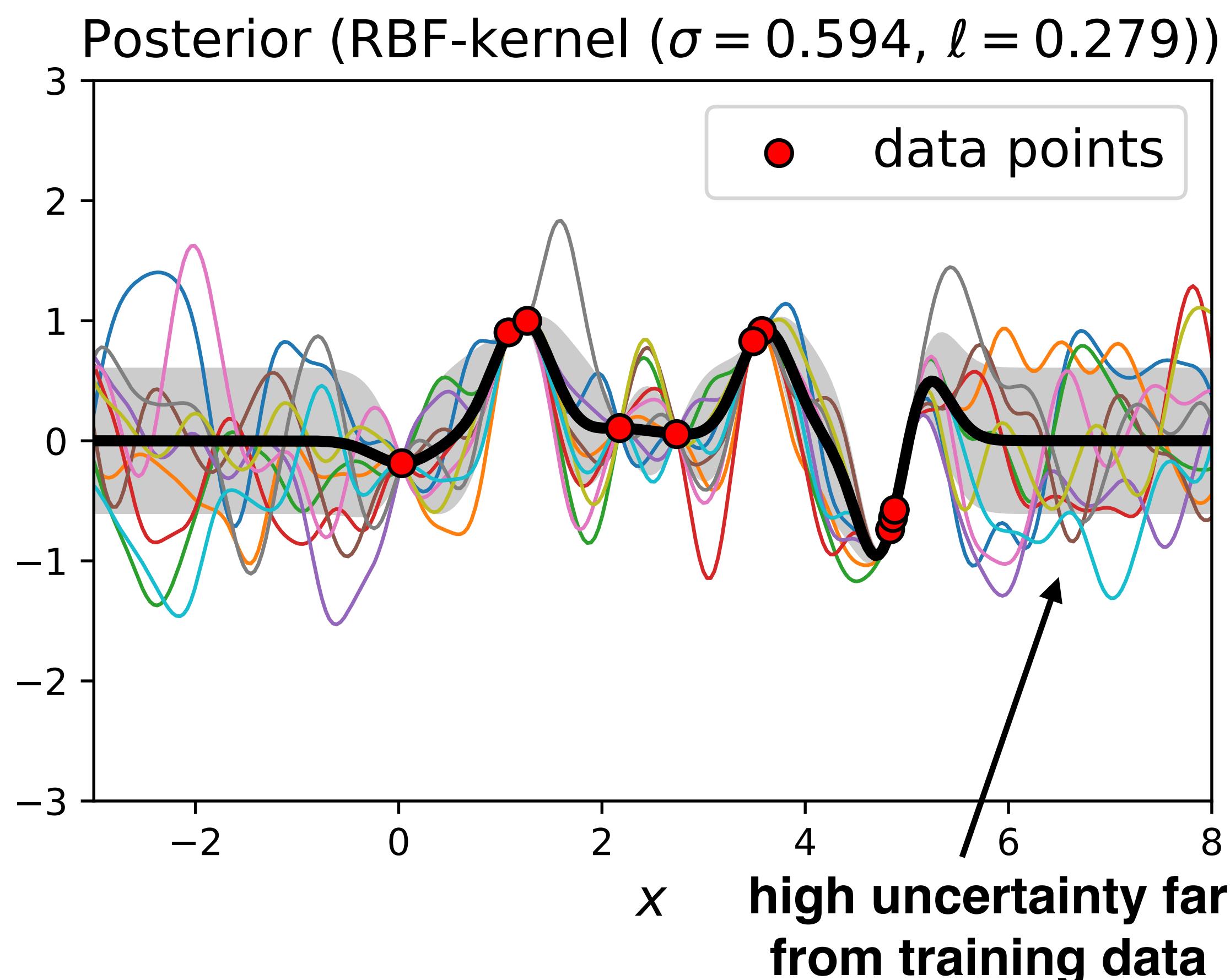
$a(x)$ – function of x and training data

none
(or tuning a small number of parameters)

kNN
Gaussian processes

Pros and cons of Gaussian processes

+ uncertainty estimation



- kernel (covariance) function?

- slow computation

Training: $O(N^3)$

$$p(\bullet\bullet\bullet\bullet) = \mathcal{N}(\begin{matrix} 0 \\ \vdots \end{matrix}, \begin{matrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & \ddots \end{matrix}) \xrightarrow{\max_{\sigma_1, \sigma_2, \sigma_3, \ell}}$$

Prediction: $O(N)$ – mean, $O(N^2)$ – std

$$p(\bullet) = \mathcal{N}(\begin{matrix} \text{mean} \\ \vdots \end{matrix}, \begin{matrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & \ddots \end{matrix})$$

mean

variance

N – number of training objects

Summary

- Gaussian process is a “distribution” over *functions*
- Regression with Gaussian Process generalizes kNN in a probabilistic manner
- Gaussian Processes provide reliable uncertainty estimates but require careful choice of kernel function and are slow in training and testing

Optimization of expensive to evaluate function

$$f(x) \rightarrow \min_x$$

computing $f(x)$
takes hours or more

Examples:

- choosing locations for extraction of minerals / oil / gold etc
- choosing hyperparameters of neural network (weight decay, dropout rate etc)
- generating new molecules (optimizing w. r. t. molecule embedding)

Optimization of expensive to evaluate function

$$f(x) \rightarrow \min_x$$

computing $f(x)$
takes hours or more

Examples:

- choosing locations for extraction of minerals / oil / gold etc
- choosing hyperparameters of neural network (weight decay, dropout rate etc)
- generating new molecules (optimizing w. r. t. molecule embedding)

General idea 1:

model $f(x)$ and
optimize model w.r.t. x



Optimization of expensive to evaluate function

$$f(x) \rightarrow \min_x$$

computing $f(x)$
takes hours or more

Examples:

- choosing locations for extraction of minerals / oil / gold etc
- choosing hyperparameters of neural network (weight decay, dropout rate etc)
- generating new molecules (optimizing w. r. t. molecule embedding)

General idea 1:
model $f(x)$ and
optimize model w.r.t. x

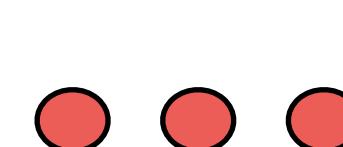


General idea 2:
model both $f(x)$ and
uncertainty in its prediction



Bayesian optimization

$$f(x) \rightarrow \min_x$$



points with known
function value

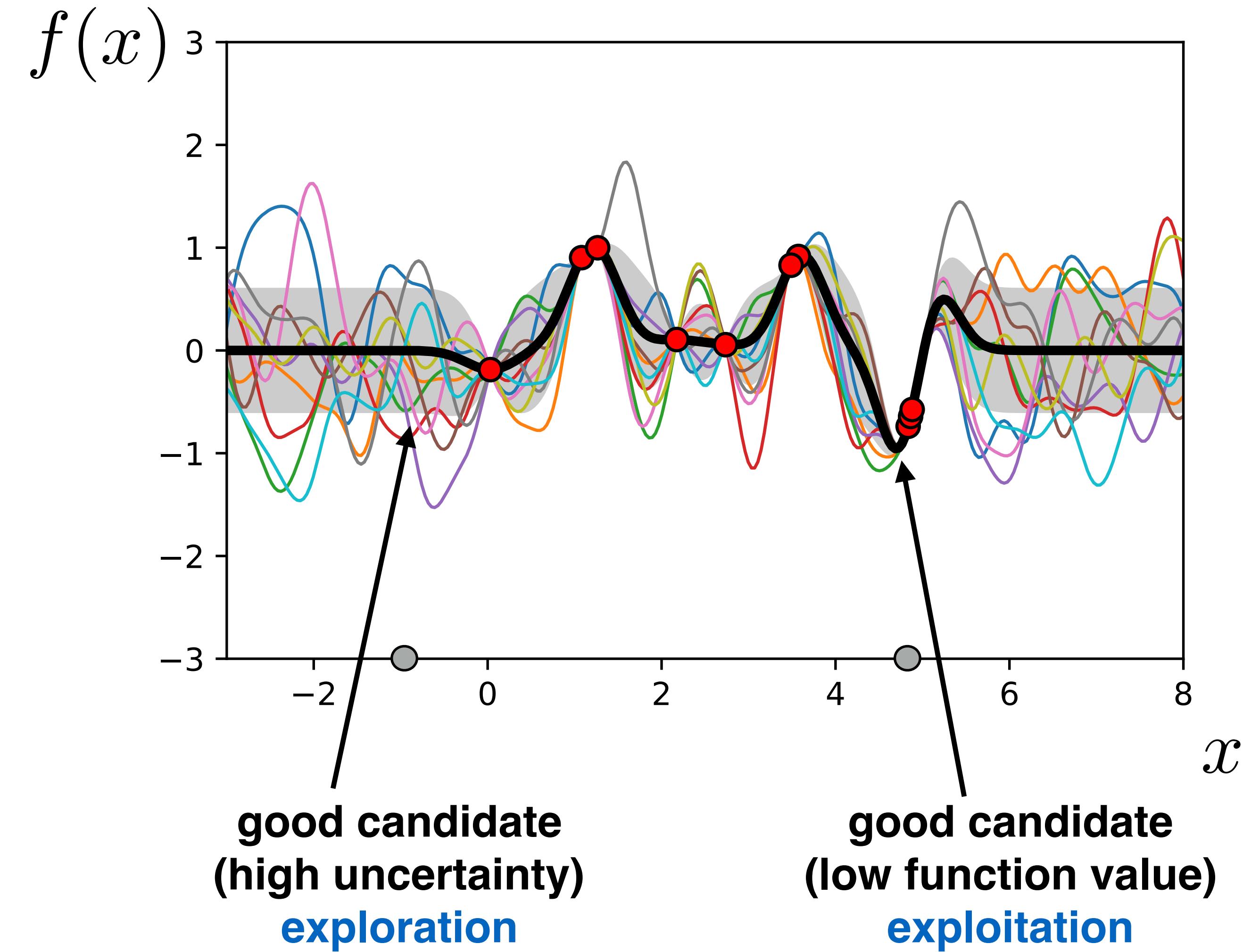


predicted
function value



uncertainty
in prediction

**Next point for computing
function value?**



Bayesian optimization

$$f(x) \rightarrow \min_x$$

computing $f(x)$
takes hours or more

Choosing a new point x_{n+1} for evaluation:

Given: $\mathcal{D}_n = \{(x_1, f(x_1)), \dots, (x_n, f(x_n))\}$ — set of currently known fun. values

1. Model $f(x)$ with a Gaussian process using \mathcal{D}_n :
ability to compute $\mu(x), \sigma(x)$ for any x
2. Acquaintance function: $\alpha(x) = \alpha(\mu(x), \sigma(x))$
Optimize acquaintance function w. r. t. x : $x_{n+1} = \operatorname{argmin}_x \alpha(x)$

Example: $\alpha(x) = \mu(x) - \beta\sigma(x)$ (β — hyperparameter)

Bayesian optimization

$$f(x) \rightarrow \min_x$$

computing $f(x)$
takes hours or more

Choosing a new point x_{n+1} for evaluation:

Given: $\mathcal{D}_n = \{(x_1, f(x_1)), \dots, (x_n, f(x_n))\}$ — set of currently known fun. values

1. Model $f(x)$ with a Gaussian process using \mathcal{D}_n :
ability to compute $\mu(x), \sigma(x)$ for any x
2. Acquaintance function: $\alpha(x) = \alpha(\mu(x), \sigma(x))$
Optimize acquaintance function w. r. t. x : $x_{n+1} = \operatorname{argmin}_x \alpha(x)$

Repeat until time limit and choose $x_* = \operatorname{argmin}_{x=x_1, \dots, x_N} f(x)$

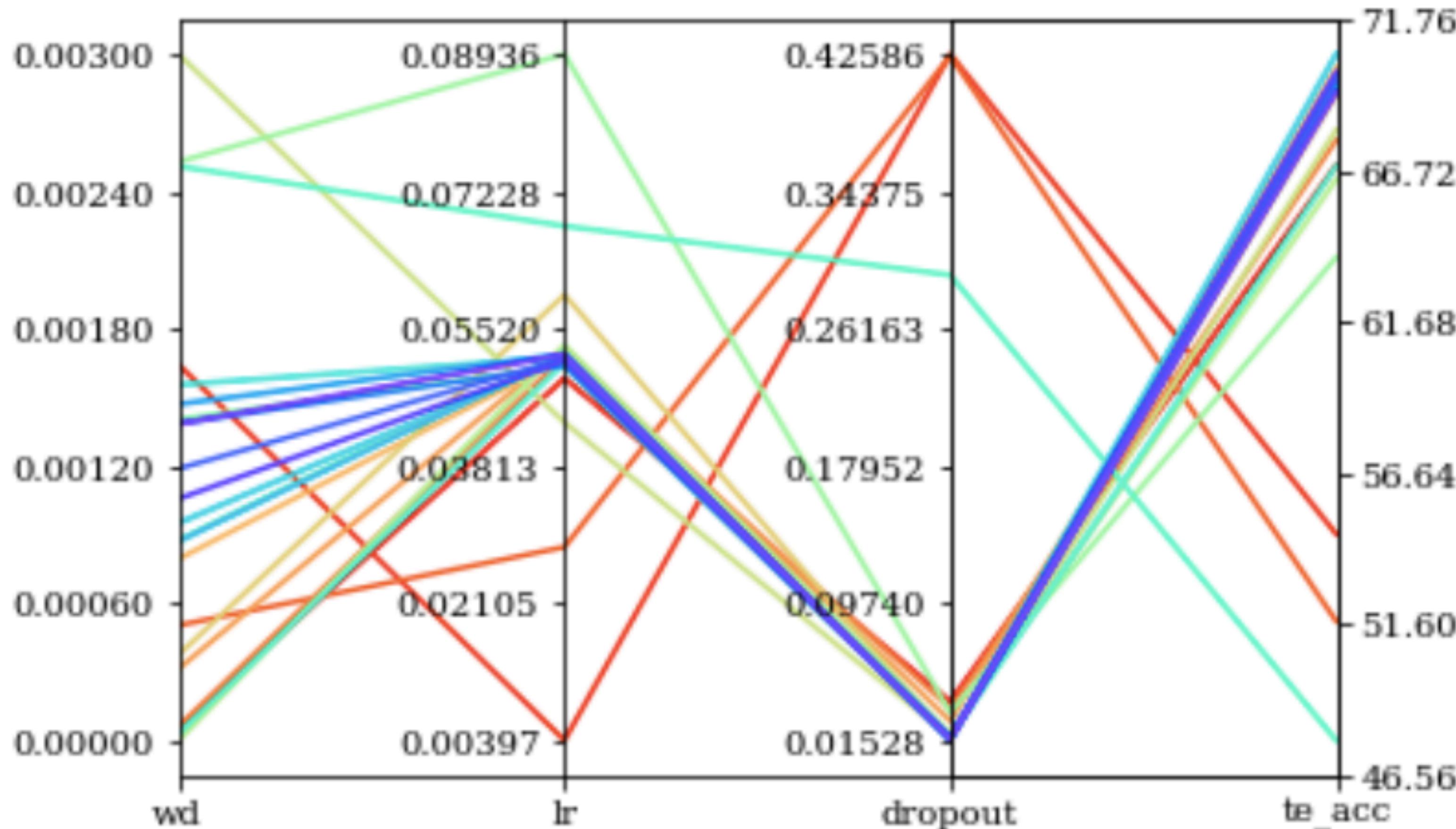
Bayesian optimization for hyperparameter search

WideResNet, CIFAR100, thin network

Choosing weight decay, learning rate, dropout

2 initial trials + 20 BO iterations

red: early iters **blue**: late iters



wd	lr	dropout	te_acc
5.614498e-05	0.049026	0.038425	67.68
1.642479e-03	0.003974	0.425863	54.14
5.029366e-04	0.028018	0.424888	50.94
6.099681e-05	0.051738	0.034468	68.60
3.129397e-04	0.052818	0.026868	70.36
7.926548e-04	0.051617	0.017005	71.26
3.743663e-04	0.059411	0.016053	70.28
3.000000e-03	0.043537	0.017029	68.98
1.000000e-07	0.053162	0.019330	67.16
2.529865e-03	0.089355	0.031902	64.30
1.404169e-03	0.050548	0.016009	70.44
8.739062e-04	0.052254	0.016042	70.60
2.509659e-03	0.067971	0.293746	46.56
2.973249e-05	0.050953	0.016744	67.66
1.554621e-03	0.051628	0.016651	70.50
9.497889e-04	0.051552	0.016384	71.76
8.709719e-04	0.051907	0.016620	70.76
1.469804e-03	0.052046	0.015283	70.88
1.384762e-03	0.050585	0.016426	70.72
1.187798e-03	0.051475	0.015822	70.82
1.055963e-03	0.051222	0.016678	71.10
1.383532e-03	0.052000	0.015923	70.36

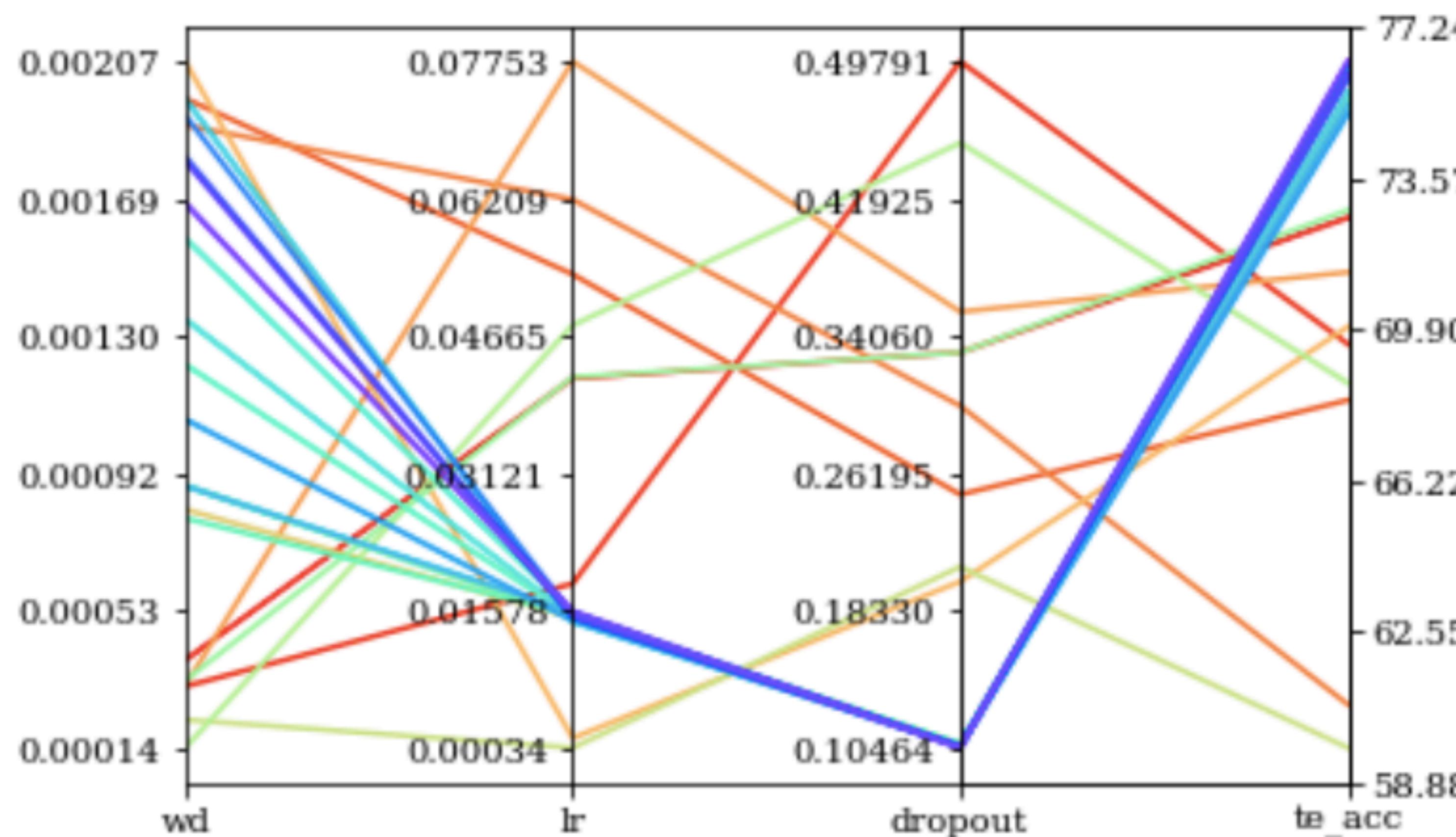
Bayesian optimization for hyperparameter search

WideResNet, CIFAR100, thin network

Choosing weight decay, learning rate, dropout

2 initial trials + 20 BO iterations

red: early iters **blue**: late iters

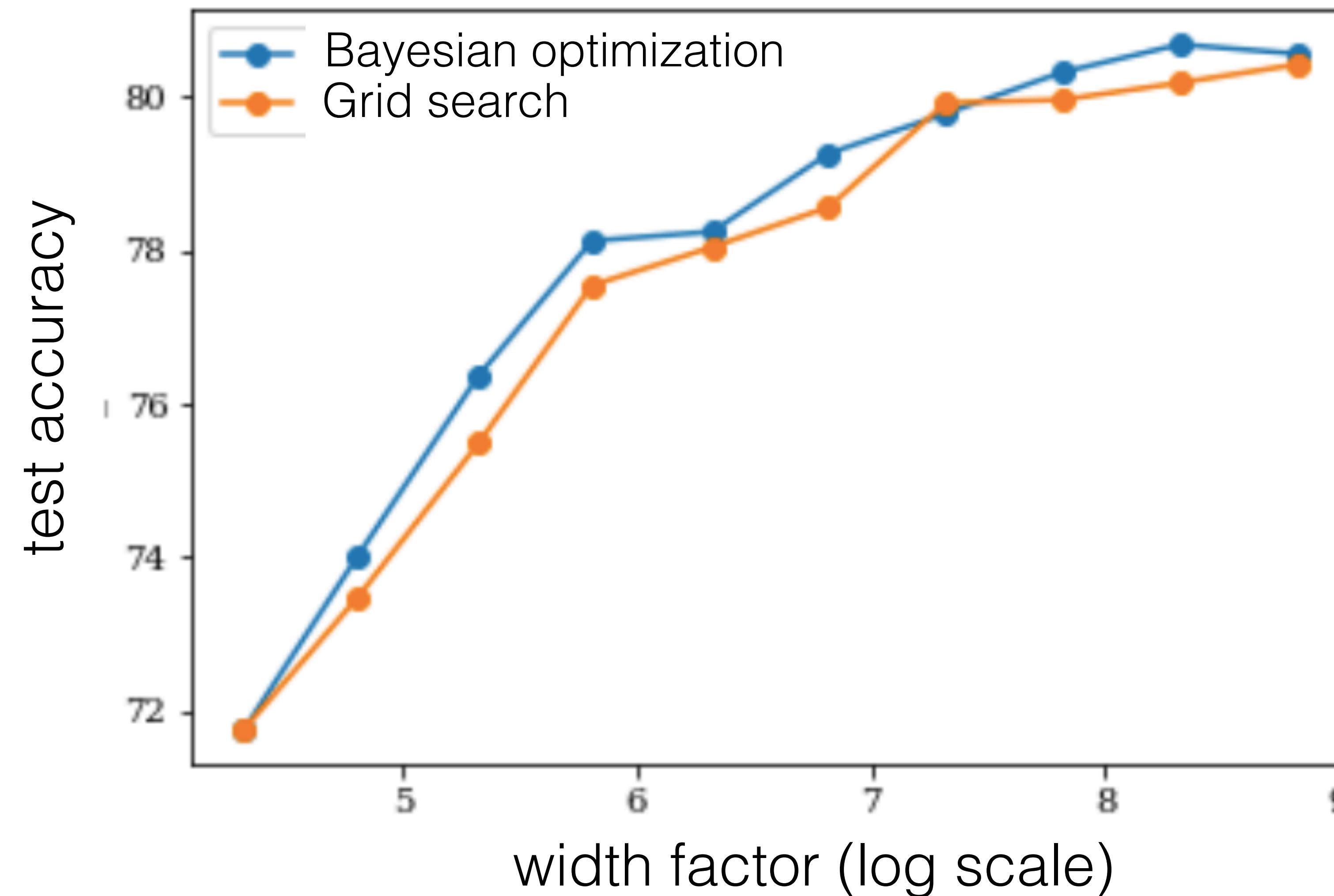


wd	lr	dropout	te_acc
0.000390	0.041890	0.331147	73.06
0.000317	0.018800	0.497906	69.68
0.001971	0.053580	0.249726	68.18
0.001895	0.061989	0.300542	60.04
0.000319	0.077525	0.354442	71.60
0.002075	0.001431	0.199903	70.14
0.000814	0.014996	0.106773	76.14
0.000223	0.000341	0.208378	58.88
0.000144	0.047858	0.451247	68.64
0.000333	0.042020	0.331092	73.26
0.000882	0.015054	0.106817	76.16
0.000790	0.015362	0.106535	76.36
0.001222	0.015142	0.106228	76.84
0.001577	0.015526	0.106035	76.80
0.001350	0.015079	0.105402	77.06
0.001970	0.014531	0.105515	76.40
0.000881	0.015591	0.105086	76.10
0.001068	0.014617	0.105467	75.84
0.001923	0.015366	0.105173	77.16
0.001808	0.015190	0.105586	76.76
0.001795	0.015050	0.104644	76.90
0.001677	0.015691	0.104903	77.24

Bayesian optimization for hyperparameter search

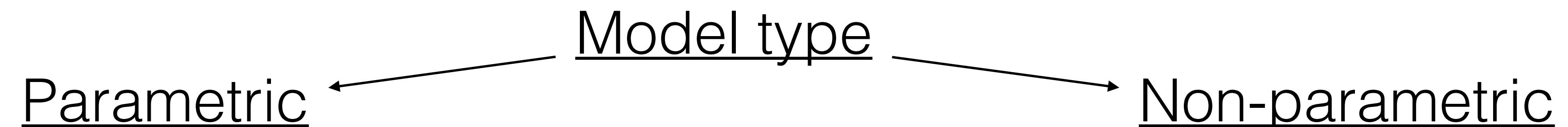
WideResNet, CIFAR100,

Choosing weight decay, learning rate, dropout for different widths



Bayesian methods for supervised learning: overview

Bayesian methods for supervised learning: overview



Model definition:

$p(y|x, \theta)$ — likelihood

$p(\theta)$ — prior distribution over model parameters

Model definition:

$p(y|f(x))$ — likelihood

$p(f(x))$ — prior process

Bayesian methods for supervised learning: overview



Model definition:

$$p(y|x, \theta) \text{ — likelihood}$$

$$p(\theta) \text{ — prior distribution over model parameters}$$

Model definition:

$$p(y|f(x)) \text{ — likelihood}$$

$$p(f(x)) \text{ — prior process}$$

Examples:

Bayesian linear regression:

$$p(y|x, w) = \mathcal{N}(y|w^T x, \beta)$$

$$p(w) = \mathcal{N}(w|0, \alpha I)$$

Sparse variational dropout (Bayesian neural net.)

$$p(y|x, w) = \text{Mult}(y|\text{NN}(x, w))$$

$$p(w) = \prod_{w_i} \frac{1}{|w_i|}$$

Gaussian process regression:

$$p(y|f(x)) = \delta(f(x))$$

$$p(f(x)) = GP(f(x)|\mu(x), k(x, x'))$$

$$\text{or } p(y|f(x)) = \mathcal{N}(y|f(x), \beta)$$

Bayesian methods for supervised learning: overview

X, Y — training data

Model type

Parametric

Non-parametric

Model definition:

$\overbrace{p(y|x, \theta)}^{\text{likelihood}}$ — likelihood

$\overbrace{p(\theta)}^{\text{prior distribution over model parameters}}$ — prior distribution over model parameters

Model definition:

$p(y|f(x))$ — likelihood

$p(f(x))$ — prior process

Finding posterior $\overbrace{p(\theta|X, Y)}^{\text{posterior}}$

Training:

None or
Finding parameters of
the prior process

$$\begin{aligned} p(y_{new}|x_{new}, X, Y) &= \\ &= \int \overbrace{p(y_{new}|x_{new}, \theta)}^{\text{likelihood}} \overbrace{p(\theta|X, Y)}^{\text{posterior}} d\theta \end{aligned}$$

Prediction:

$\overbrace{p(y_{new}|x_{new}, X, Y)}^{\text{finding}} \overbrace{\text{directly from data } X, Y}^{\text{posterior}}$

Bayesian methods for supervised learning: overview

X, Y — training data

Model type

Parametric

Non-parametric

Model definition:

$\underline{p(y|x, \theta)}$ — likelihood

$\underline{p(\theta)}$ — prior distribution over model parameters

Inference type (training)

Full Bayesian
inference

$$\underline{p(\theta|X, Y)} = \frac{\underline{p(Y|X, \theta)} \underline{p(\theta)}}{\int \underline{p(Y|X, \tilde{\theta})} \underline{p(\tilde{\theta})} d\tilde{\theta}}$$

lectures 1, 2

lecture 5

Maximum
posterior

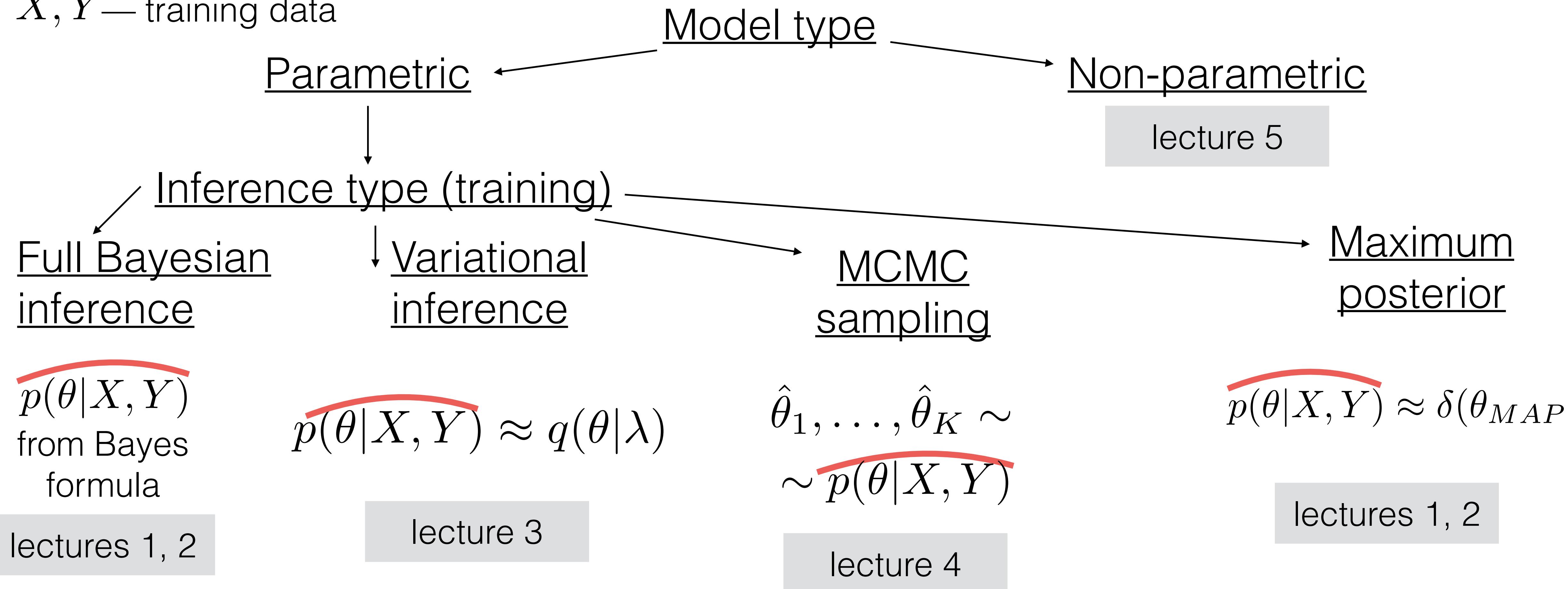
$$\begin{aligned}\theta_{MAP} &= \operatorname{argmax}_{\theta} \underline{p(\theta|X, Y)} = \\ &= \operatorname{argmax}_{\theta} \underline{p(Y|X, \theta)} \underline{p(\theta)}\end{aligned}$$

$$\underline{p(\theta|X, Y)} \approx \delta(\theta_{MAP})$$

lectures 1, 2

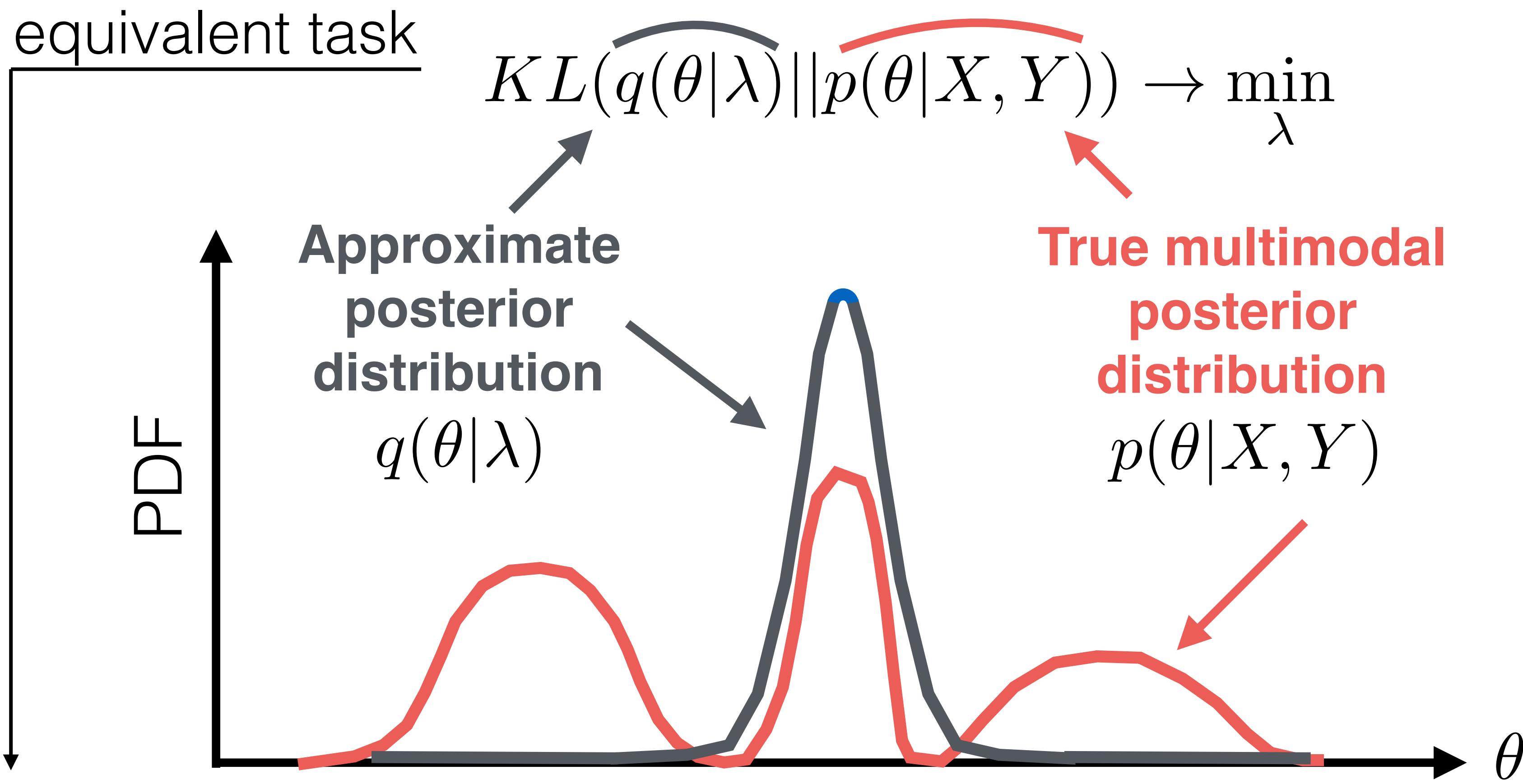
Bayesian methods for supervised learning: overview

X, Y — training data



$\overbrace{p(y|x, \theta)}^{\text{likelihood}}$ — likelihood
 $\overbrace{p(\theta)}^{\text{prior}}$ — prior

Variational inference



$$\sum_{i=1}^N \mathbb{E}_{q(\theta|\lambda)} \log \overbrace{p(y^i|x^i, \theta)}^{\text{Data term}} - \underbrace{KL(q(\theta|\lambda) || p(\theta))}_{\text{Regularizer}} \rightarrow \max_{\lambda}$$

$\overbrace{p(y|x, \theta)}$ — likelihood

$\overbrace{p(\theta)}$ — prior

MCMC Sampling

$$p(y_*|x_*, X, Y) = \mathbb{E}_{\overbrace{p(\theta|X, Y)}} \overbrace{p(y_*|x_*, \theta)} \approx \frac{1}{K} \sum_{k=1}^K \overbrace{p(y_*|x_*, \theta^k)}$$
$$\theta^k \sim \overbrace{p(\theta|X, Y)}$$

We only need to know how to **sample** from the **unnormalized** posterior!

$$p(\theta|X, Y) = \frac{\boxed{\overbrace{p(Y|X, \theta)} \overbrace{p(\theta)}}}{\int p(Y|X, \tilde{\theta}) p(\tilde{\theta}) d\tilde{\theta}}$$

unnormalized posterior

Bayesian methods for supervised learning: overview

X, Y — training data

