

Nadia Chirkova



Gaussian processes

2021



Yandex

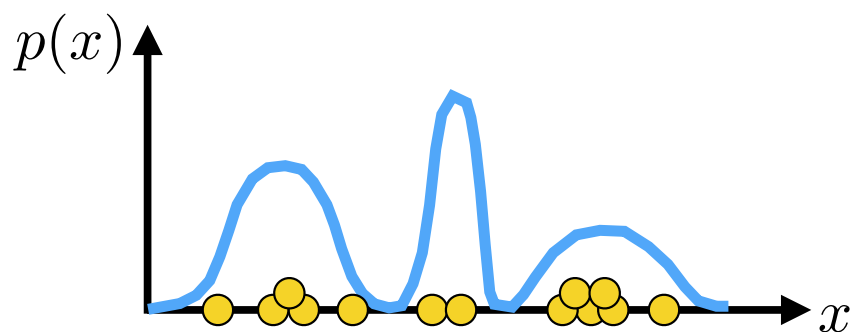


EPFL

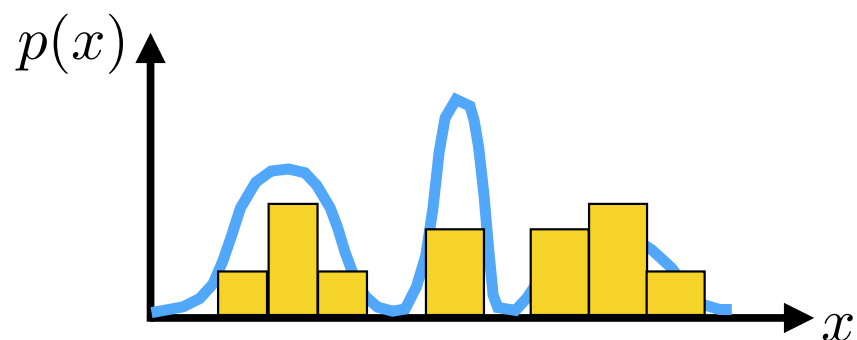


Slides are partially based on lectures of Dmitry Vetrov, Dmitry Kropotov and Kirill Struminsky, deepbayes.ru/2018

Sampling points from a distribution



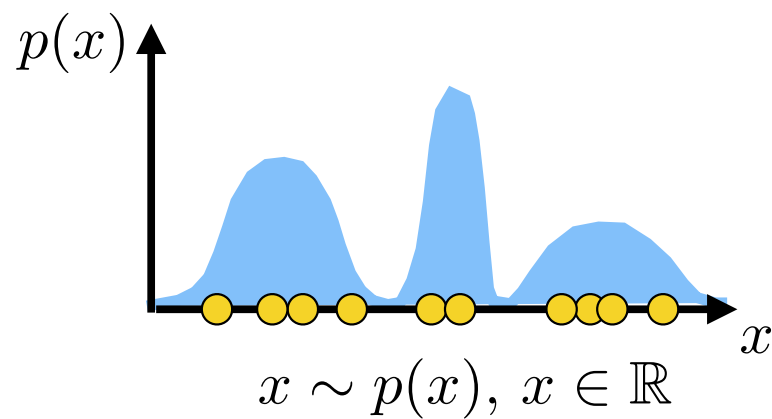
$$x \sim p(x), x \in \mathbb{R}$$



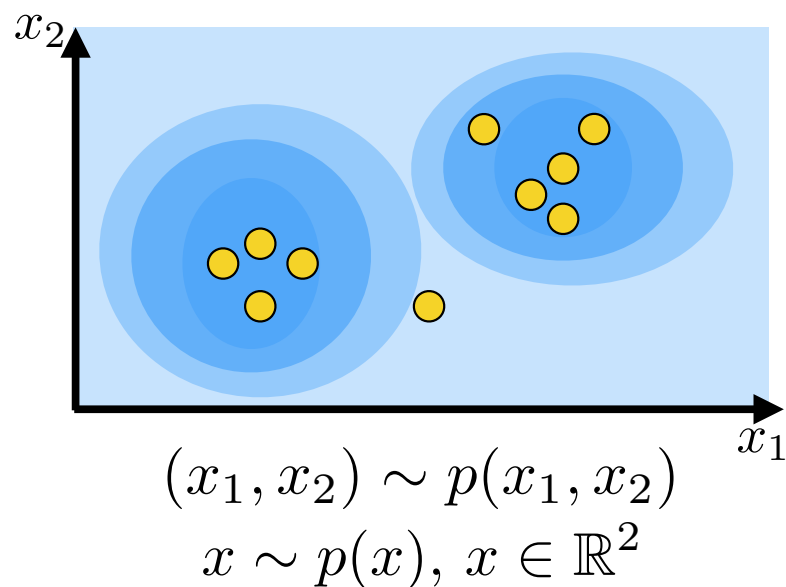
histogram approaches $p(x)$
when the number of points grows

Sampling points from a multivariate distribution

1-dimensional
distribution



2-dimensional distribution



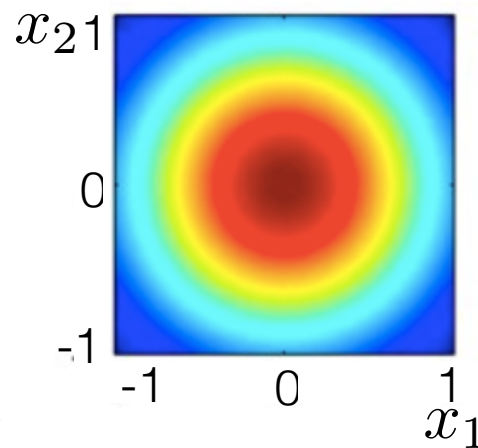
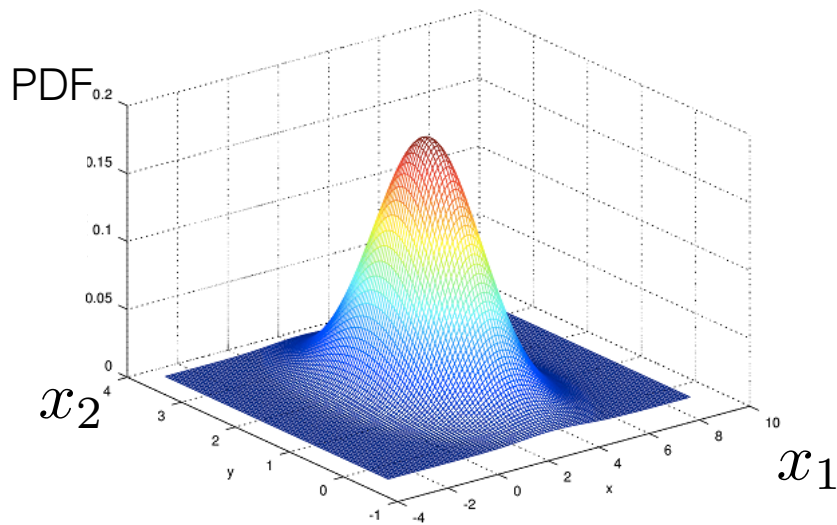
... n-dimensional distribution:

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

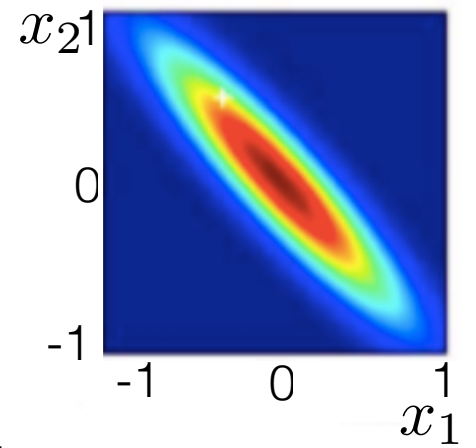
Multivariate normal (Gaussian) distribution

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

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



diagonal Σ



non-diagonal Σ

Multivariate normal (Gaussian) distribution

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

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

Diagram illustrating the covariance matrix Σ for two variables x_1 and x_2 :

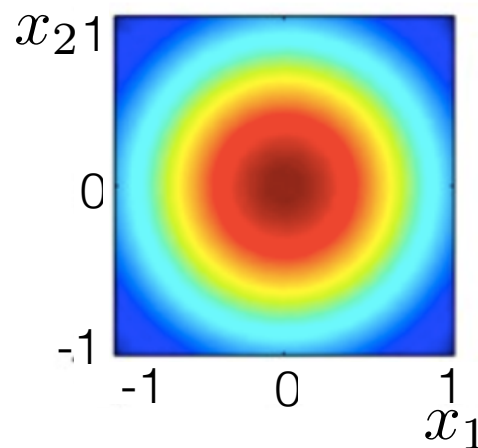
$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12} & \Sigma_{22} \end{bmatrix}$$

Annotations:

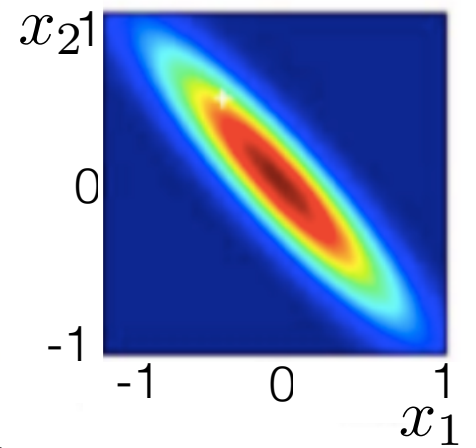
- $\text{Var}(x_1)$ points to Σ_{11} .
- $\text{Cov}(x_1, x_2)$ points to Σ_{12} .

diagonal $\Sigma \Rightarrow$

independent x_1, \dots, x_d



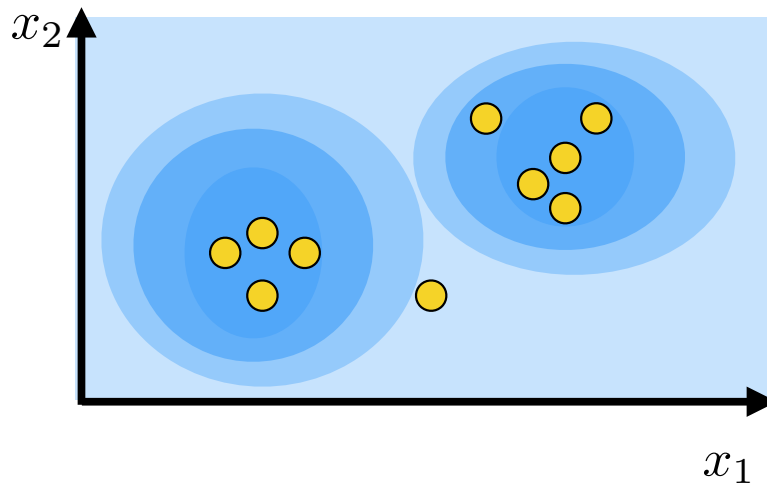
diagonal Σ



non-diagonal Σ

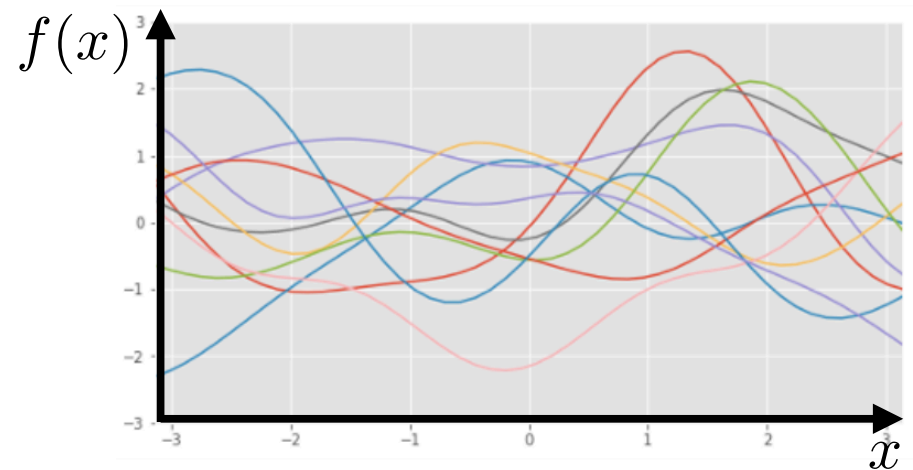
Sampling functions from a process

Multivariate distribution
— sample **points**



$$x \sim p(x)$$

Process
— sample **functions**

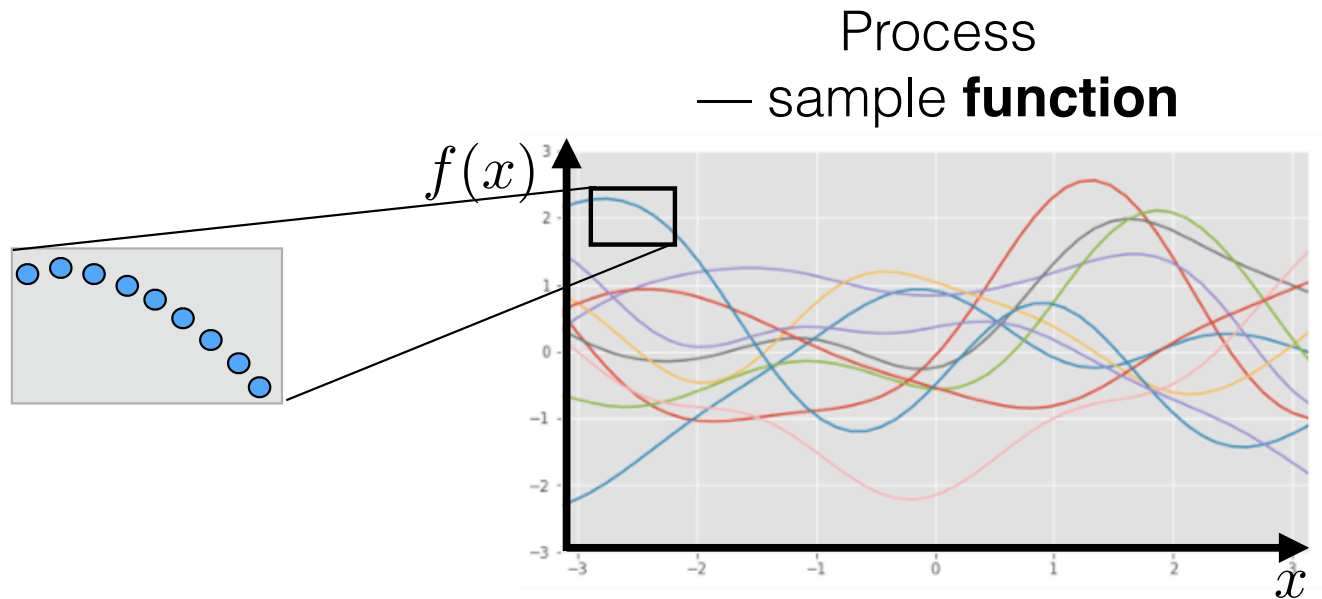


$$f(x) \sim GP(m(x), k(x, x'))$$

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



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

each line is a sample from the process

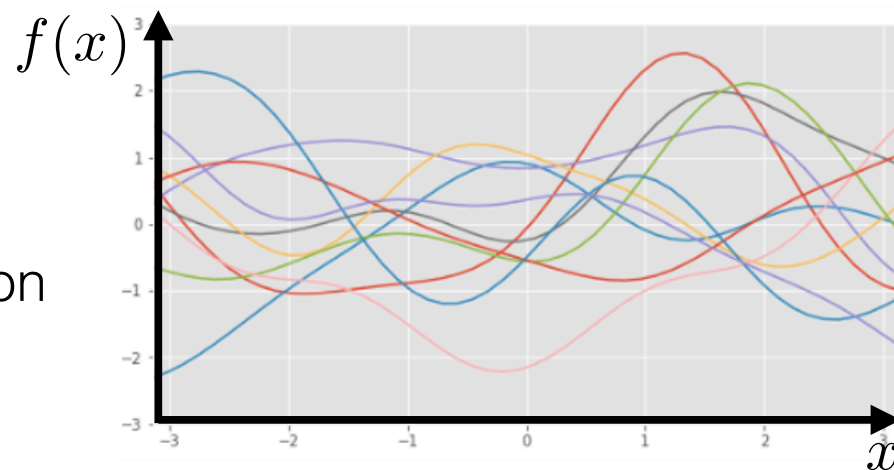
Gaussian process

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

$m(x)$ — mean function

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

(x may be a vector in general case)



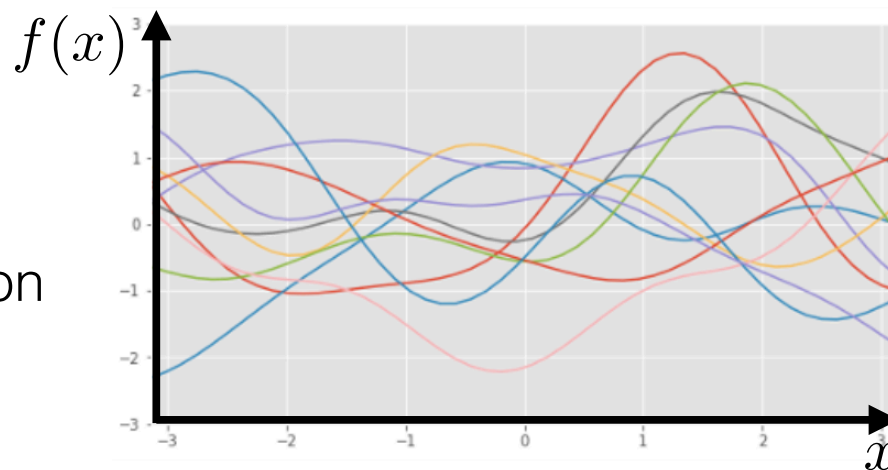
Gaussian process

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

$m(x)$ — 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

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

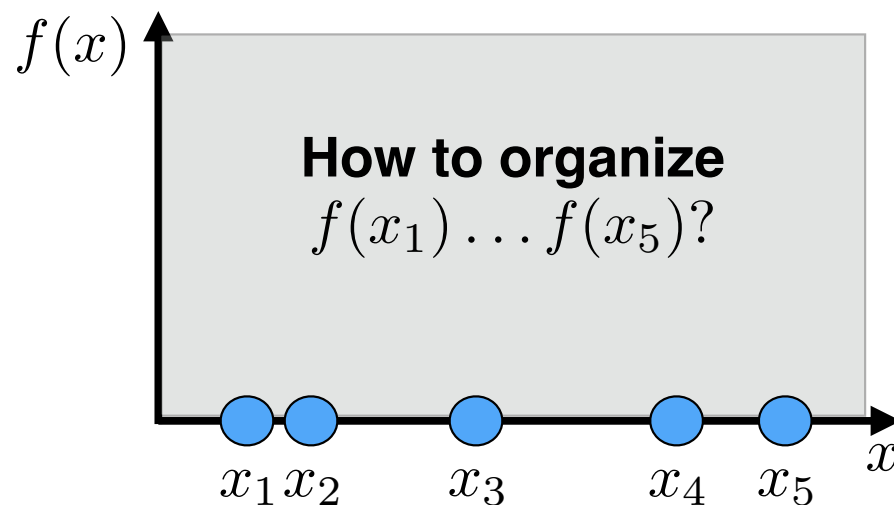
$$m(x) = 0$$

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

[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

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

$$m(x) = 0$$

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

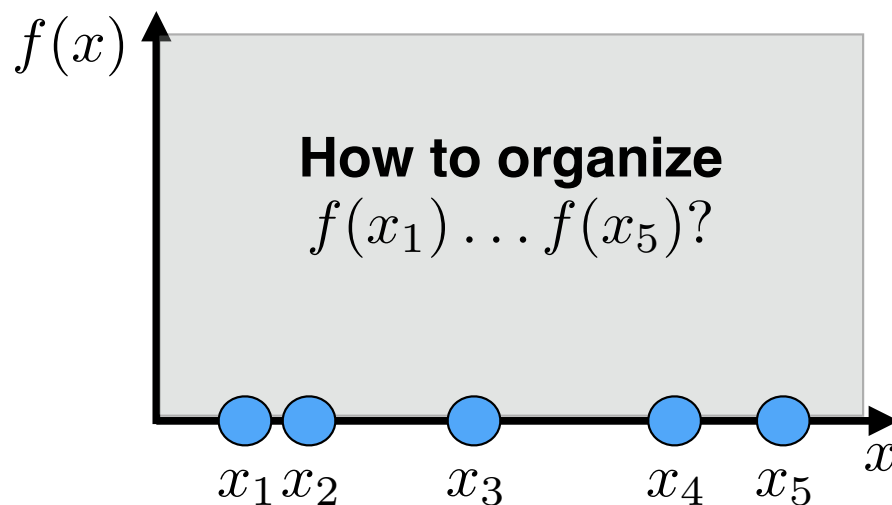
[condition] = 1 if condition is True else 0

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

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

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

(all x are independent on each other)



Example 1: white noise

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

$$m(x) = 0$$

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

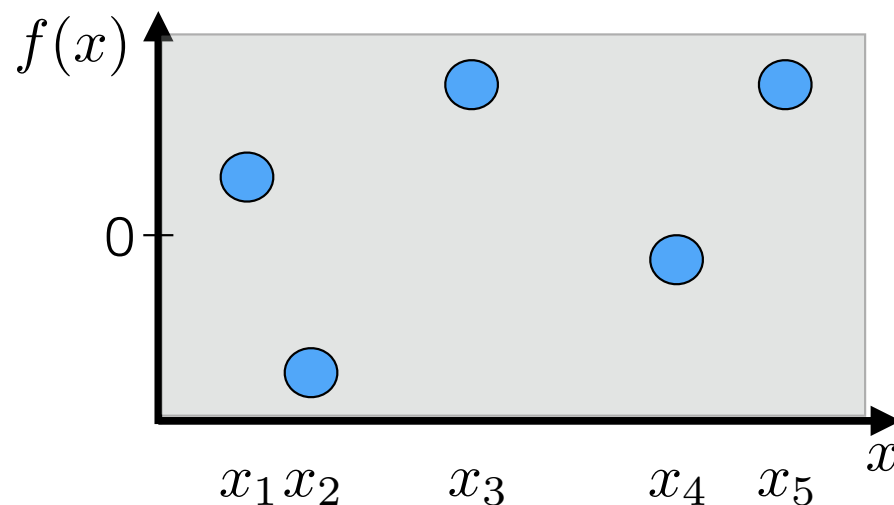
[condition] = 1 if condition is True else 0

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

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

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

(all x are independent on each other)



for any x

$f(x)$ is sampled
independently

Example 2: constant function

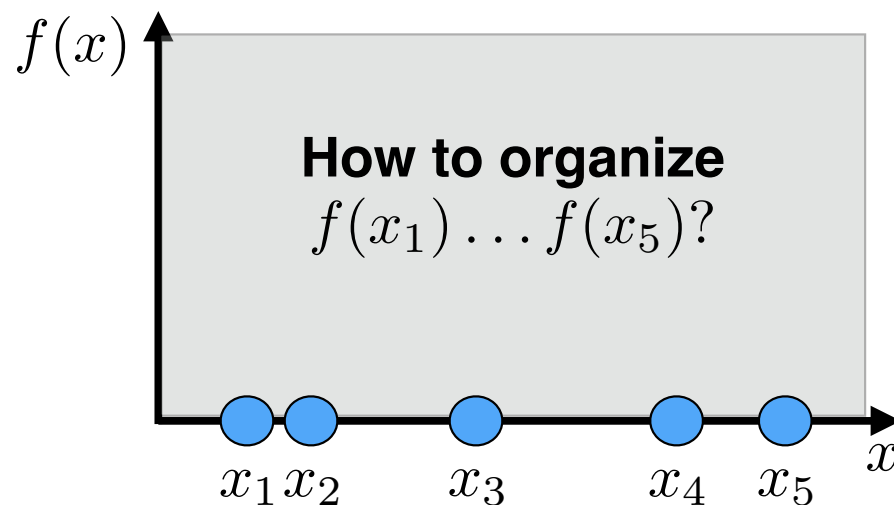
$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

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

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

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



Example 2: constant function

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

$$m(x) = 0$$

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

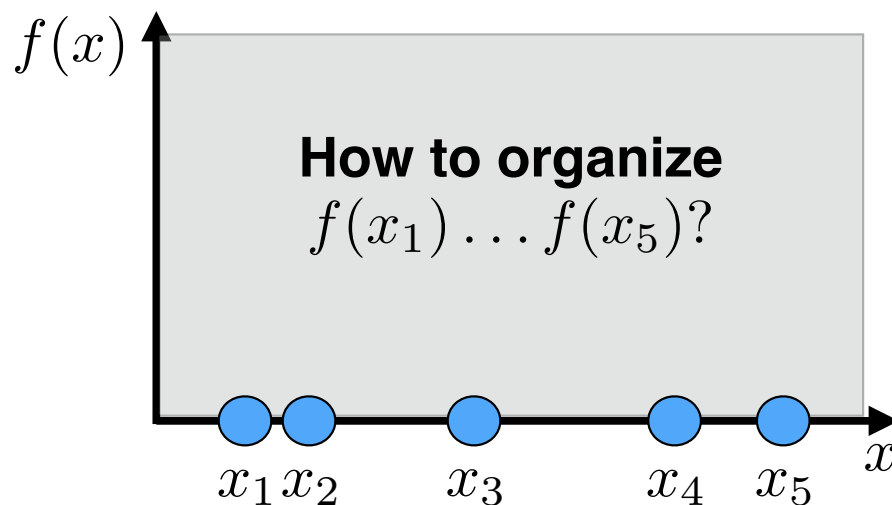
$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

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

$$\forall i \neq j$$

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

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



Example 2: constant function

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

$$m(x) = 0$$

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

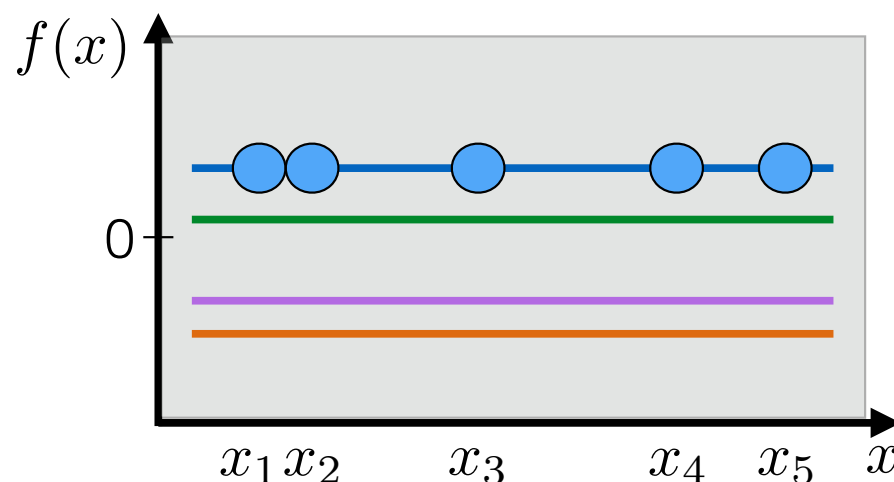
$$(f(x_1), \dots, f(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

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

$$\forall i \neq j$$

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

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



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

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

$$m(x) = 0$$

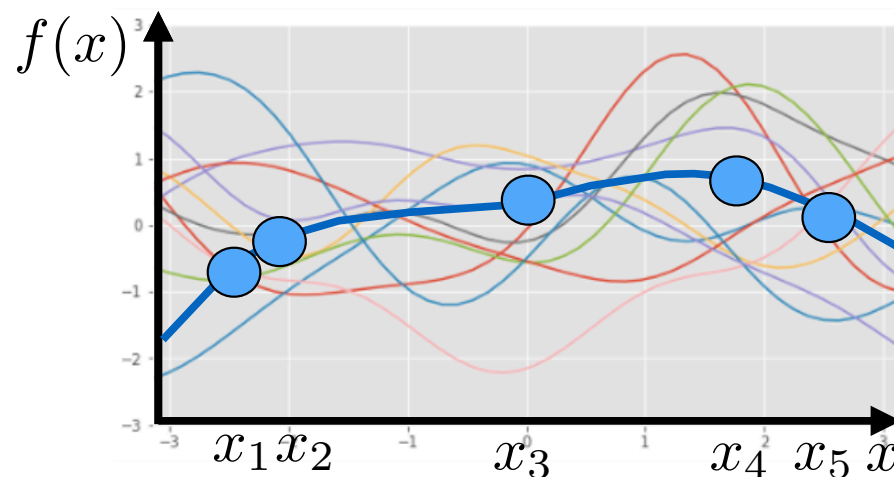
$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$$

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

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

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

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



Example 3: RBF-kernel

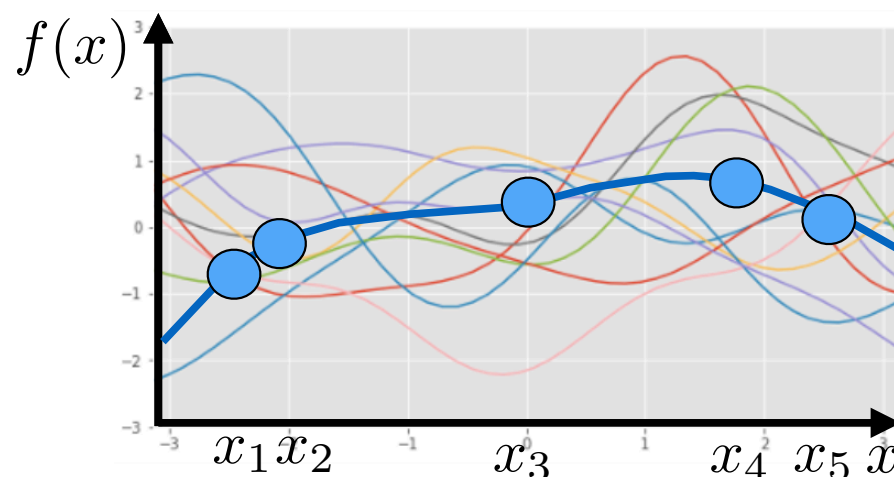
$$f(x) \sim GP(m(x), k(x, x'))$$

$$m(x) = 0$$

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

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

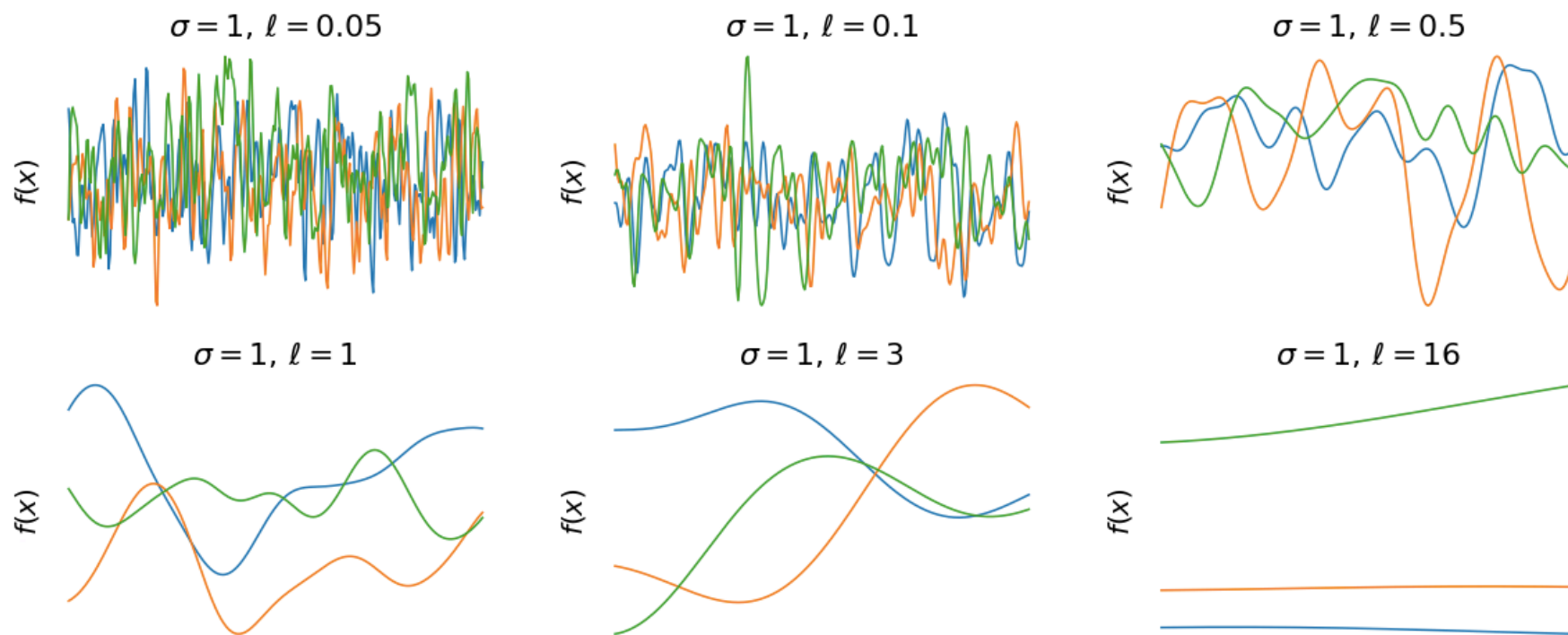
$$\mu = 0 \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$



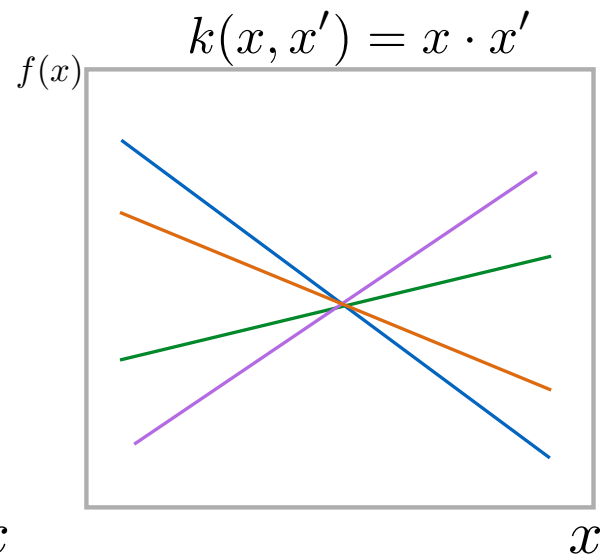
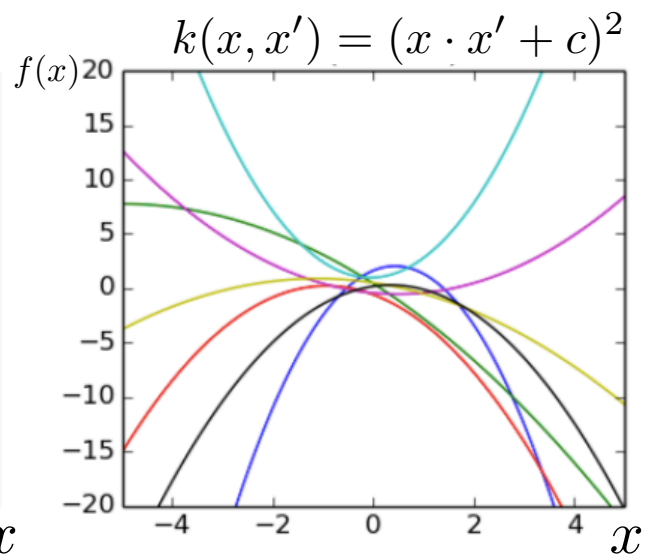
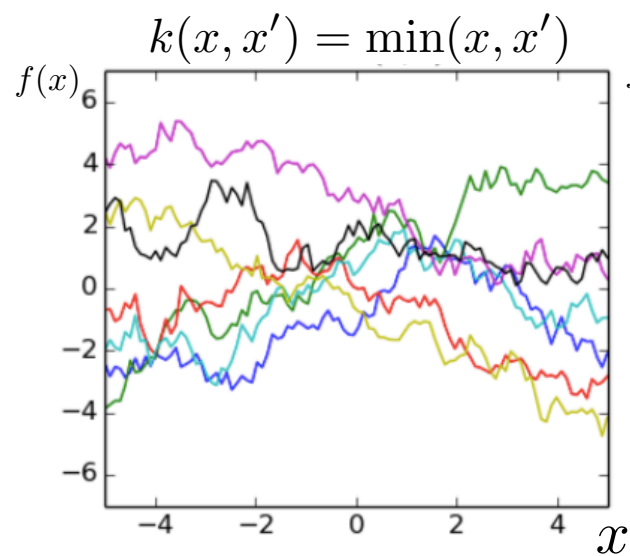
σ^2 defines the “height” of the function
 ℓ^2 defines the frequency of fluctuations

Example 3: RBF-kernel

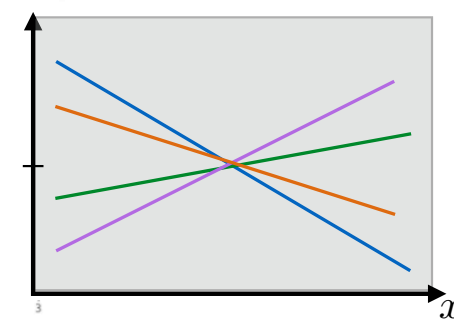
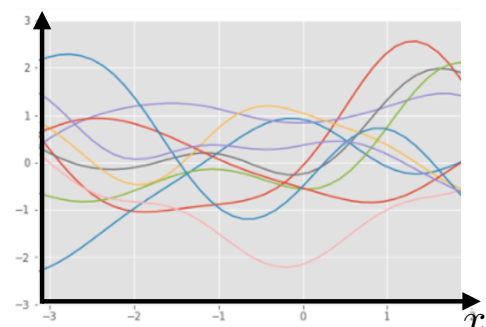
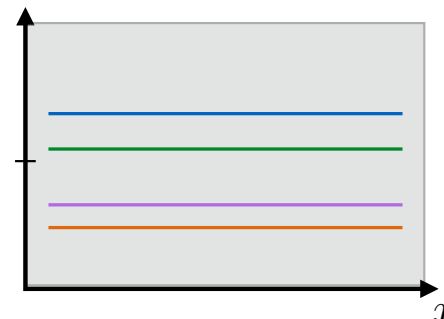
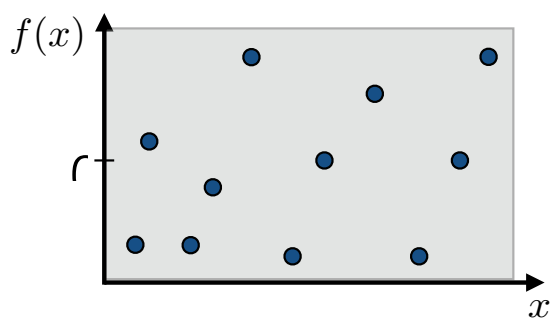
ℓ^2 defines the frequency of fluctuations



More kernels



Sum-kernel



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

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

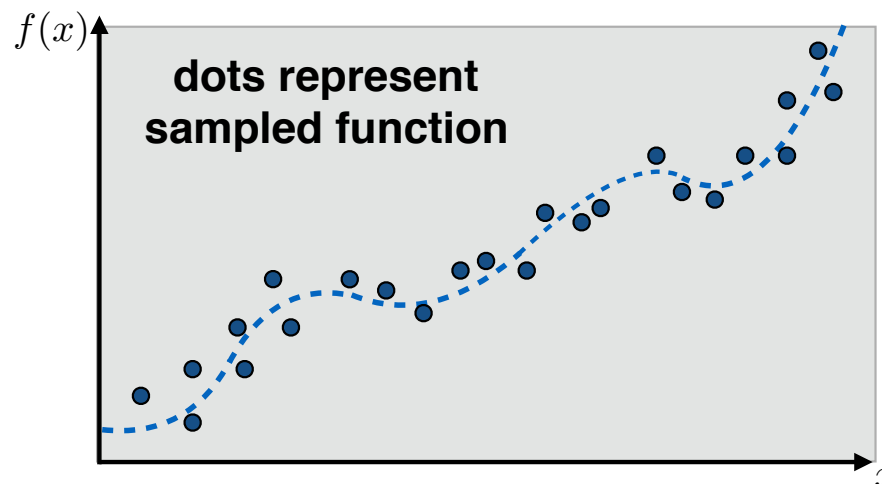
$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2r^2}\right)$$

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

Sum-kernel (multidimensional case):

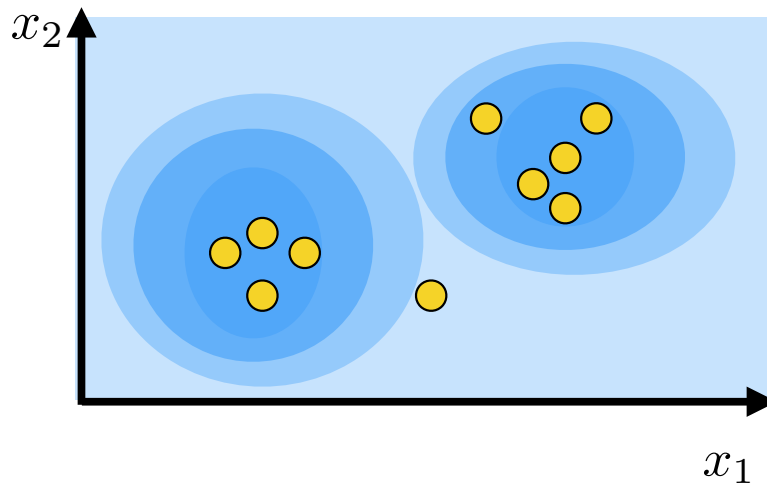
$$k(x, x') = x^T x' + \sigma_1^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right) + \sigma_2^2 [x = x'] + \sigma_3^2$$

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



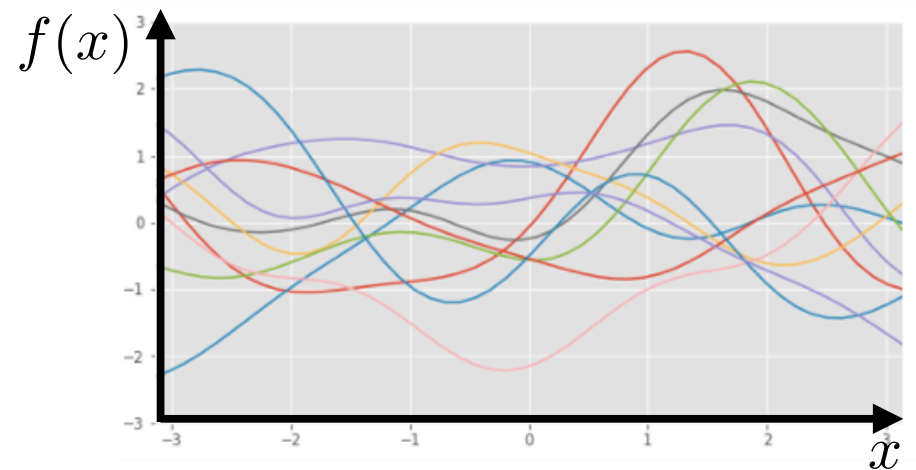
Sampling functions from a process

Multivariate distribution
— sample **points**



$$x \sim p(x)$$

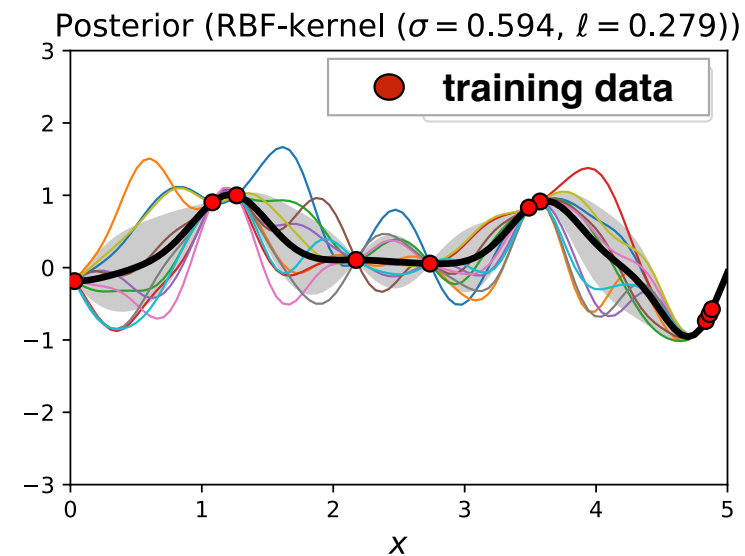
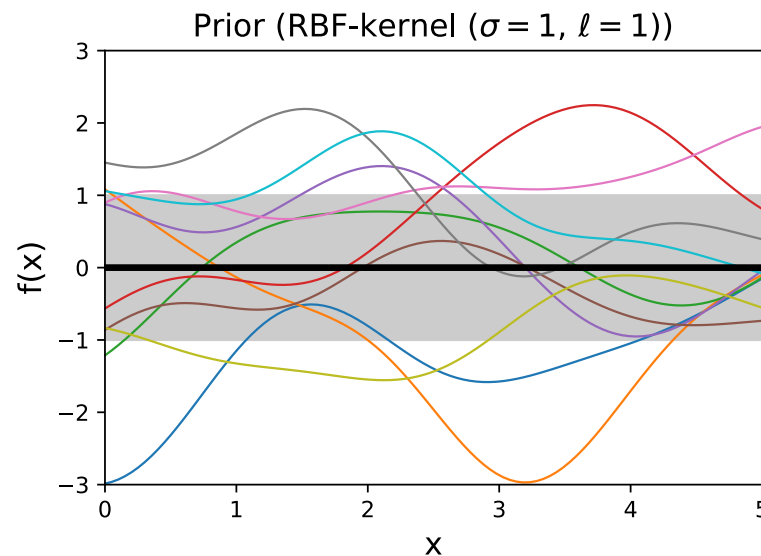
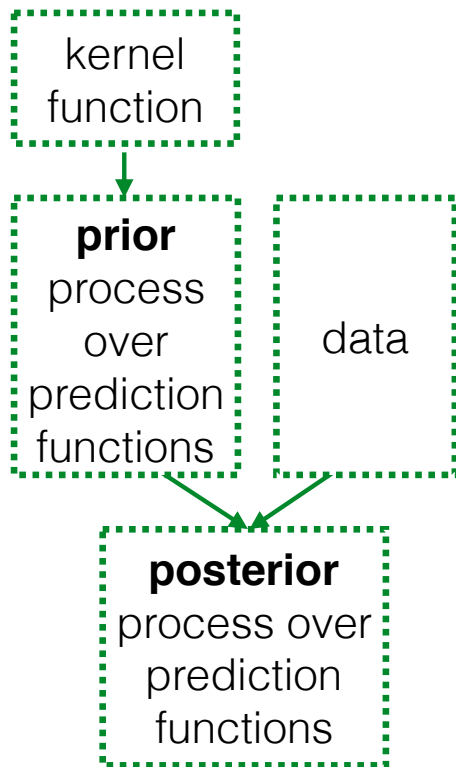
Process
— sample **functions**



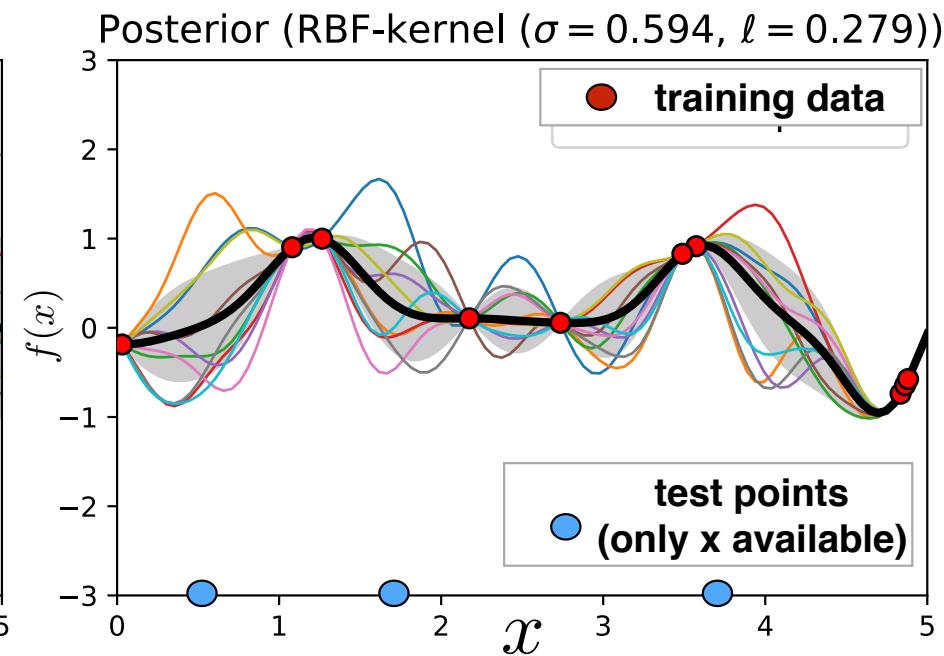
$$f(x) \sim GP(m(x), k(x, x'))$$

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 $a(x)$

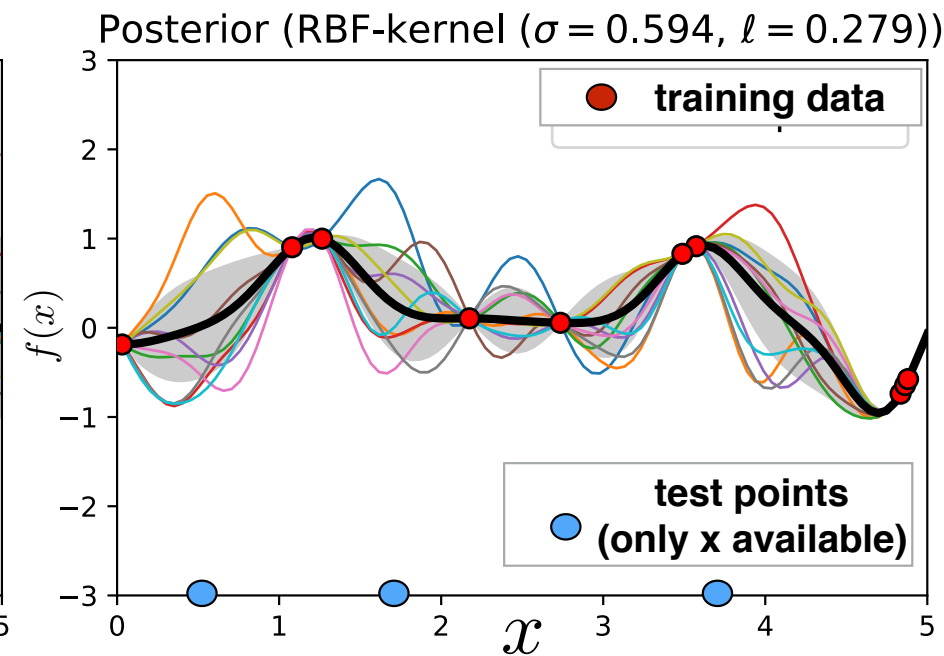
Training data ● ● ● :

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

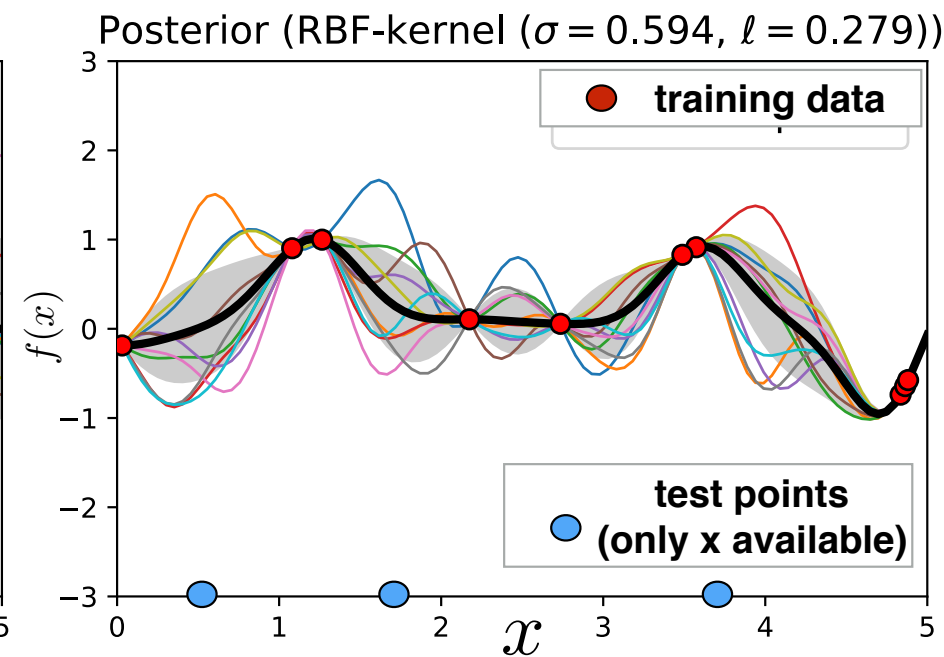
$Y^{tr} = \{y_i^{tr}\}_{i=1}^N$, $y_i^{tr} \in \mathbb{R}$ — targets

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

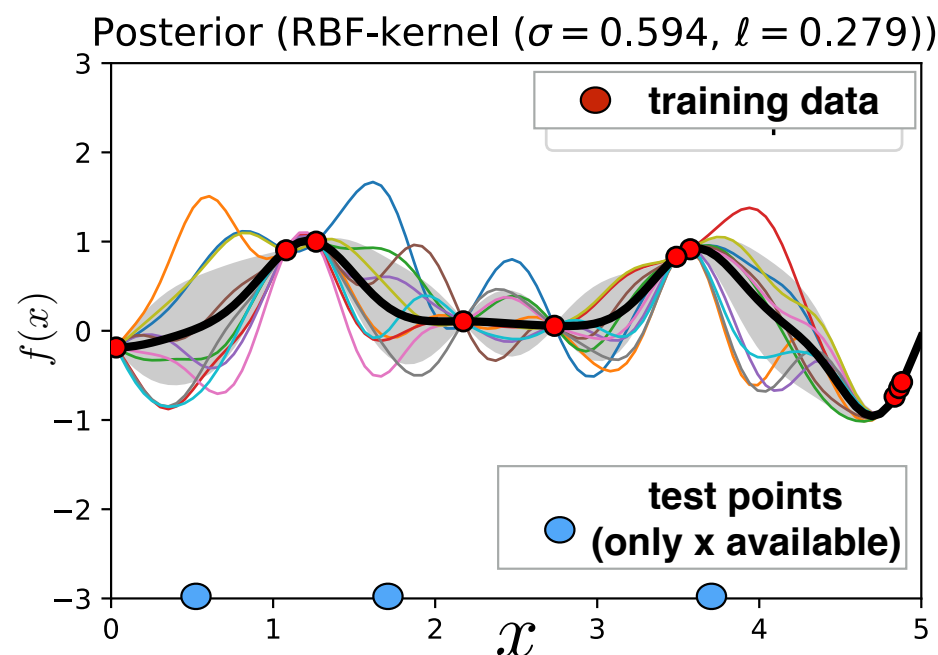
Test points ● ● ● (any set of points):

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

Find:

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

Conditioning in multivariate normal distribution



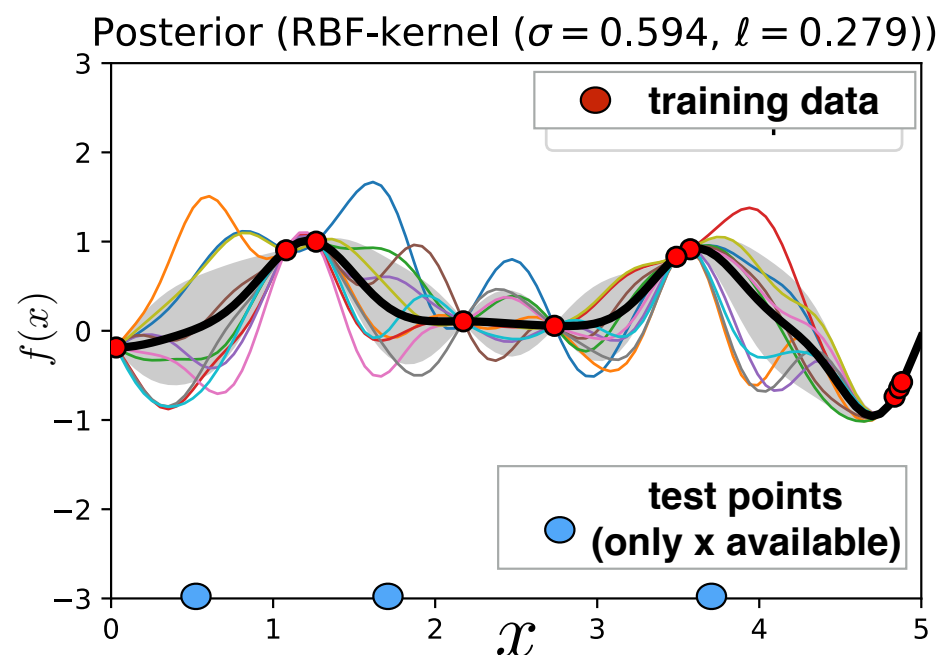
Definition of Gaussian process:

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

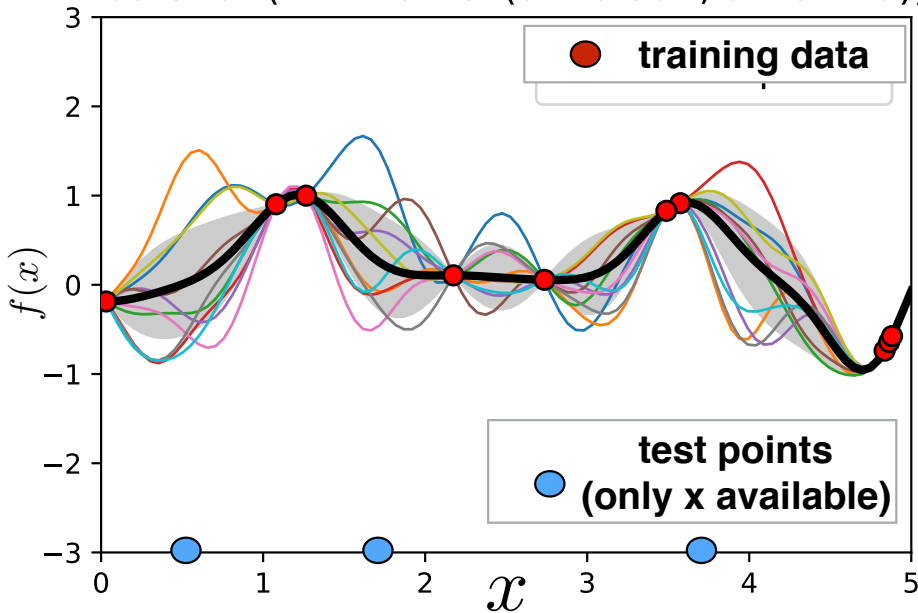
$$\forall n \quad \forall (x_1, \dots, x_n) \quad (a(x_1), \dots, a(x_n)) \sim \mathcal{N}(\mu, \Sigma)$$

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

$$(\text{red dots}, \text{blue dots}) \sim \mathcal{N}\left(\begin{bmatrix} 0 \end{bmatrix}, \begin{bmatrix} \text{red} & \text{grey} \\ \text{grey} & \text{blue} \end{bmatrix}\right)$$

Conditioning in multivariate normal distribution

Posterior (RBF-kernel ($\sigma = 0.594$, $\ell = 0.279$))



Definition of Gaussian process:

every finite set of function values
has a multivariate normal distribution

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

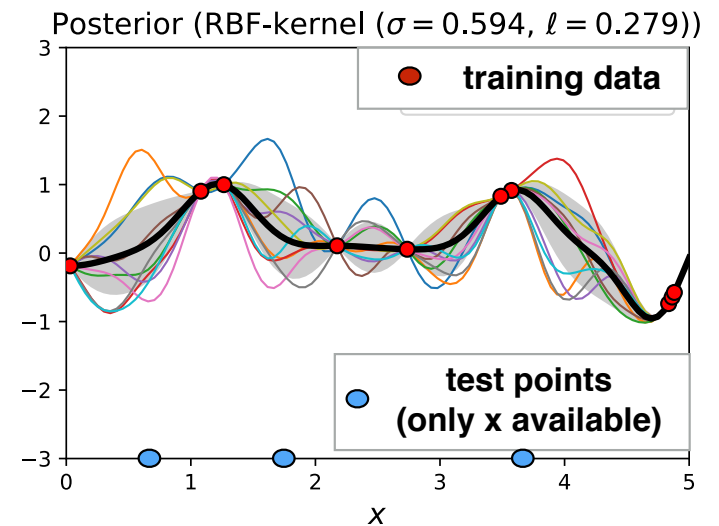
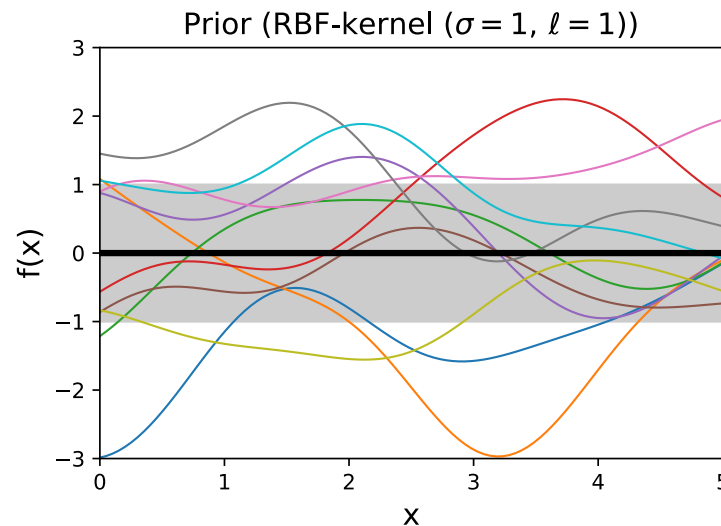
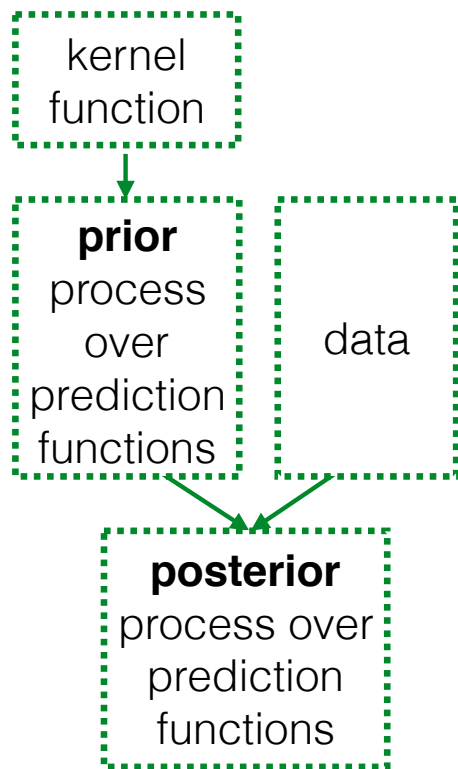
$$\mu = \{m(x_i)\}_{i=1}^n \quad \Sigma = \{k(x_i, x_j)\}_{i,j=1}^{n,n}$$

$$(\text{red dots}, \text{blue dots}) \sim \mathcal{N}(\underbrace{0}_{\text{white box}}, \begin{bmatrix} \text{red box} & \text{gray box} \\ \text{gray box} & \text{blue box} \end{bmatrix})$$

According to properties
of normal distribution:

$$(\text{blue dots}) \sim \mathcal{N}(\underbrace{\begin{bmatrix} \text{gray box} & \text{red box} \end{bmatrix} \begin{matrix} -1 \\ \vdots \\ -1 \end{matrix}}_{\text{defines a new mean function}}, \underbrace{\begin{bmatrix} \text{blue box} & - & \text{gray box} & \text{red box} & -1 \\ & & & & \end{bmatrix}}_{\text{defines a new covariance function}})$$

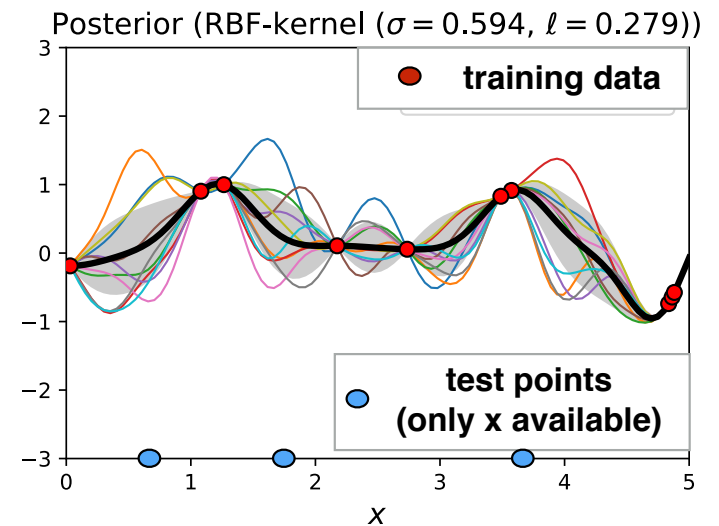
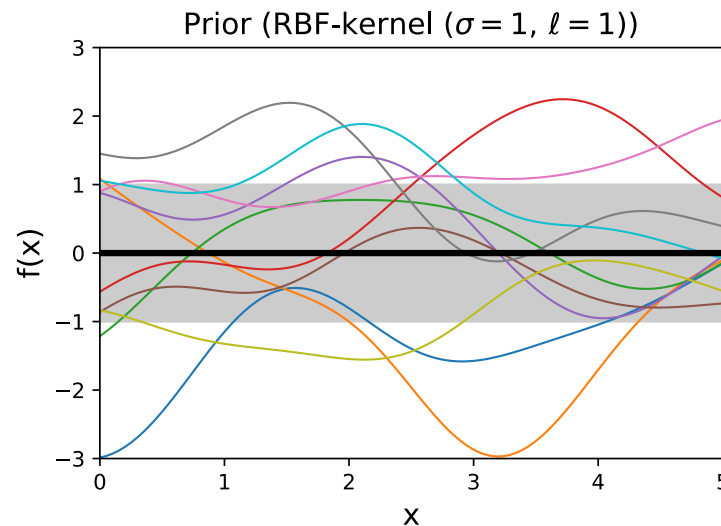
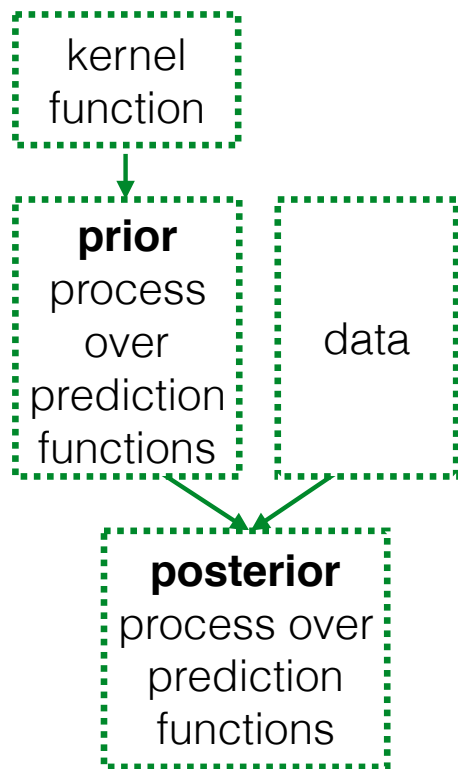
Gaussian processes for regression



$$p(\bullet\bullet\bullet) = \mathcal{N}\left(\begin{bmatrix} \text{gray} & \text{red} \end{bmatrix}^{-1} \begin{bmatrix} \bullet\bullet\bullet \\ \bullet\bullet\bullet \end{bmatrix}, \begin{bmatrix} \text{blue} & - & \text{gray} & \text{red} \end{bmatrix}^{-1} \begin{bmatrix} \text{gray} \end{bmatrix}\right)$$

$$\begin{aligned} \text{gray} &= k(X^{te}, X^{tr}) & \text{red} &= k(X^{tr}, X^{tr}) & \text{blue} &= k(X^{te}, X^{te}) \\ \bullet\bullet\bullet\bullet &= a(X^{tr}) = Y^{tr} & \bullet\bullet\bullet &= a(X^{te}) \end{aligned}$$

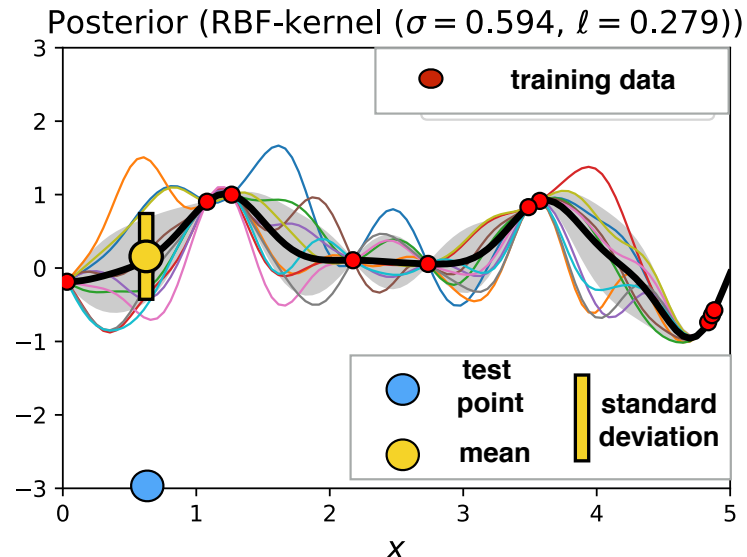
Gaussian processes for regression



$$p(\bullet \bullet \bullet) = \mathcal{N} \left(\begin{bmatrix} \text{gray} & \text{red} \end{bmatrix}^{-1} \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix}, \begin{bmatrix} \text{blue} & - & \text{gray} & \text{red} & \text{gray} \end{bmatrix}^{-1} \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix} \right)$$




Training? Prediction?



Training and prediction in GP for regression



Prediction:

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

 = $k(x_*, X^{tr})$
 = $k(X^{tr}, X^{tr})$
 = $k(x_*, x_*)$

 = $a(X^{tr}) = Y^{tr}$
 = $a(x_*)$

Training and prediction in GP for regression

Training:

$$\mathbf{p}(\bullet\bullet\bullet\bullet) = \mathcal{N}\left(\boxed{0}, \text{red square}\right) \rightarrow \max_{\sigma_1, \sigma_2, \sigma_3, \ell}$$

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

parameters of the
kernel (covariance) function

Prediction:

$$\mathbf{p}(\bullet) = \mathcal{N}\left(\underbrace{\text{grey rectangle} \text{red square}^{-1} \text{red dots}}_{\text{mean}}, \underbrace{\text{blue square} - \text{grey rectangle} \text{red square}^{-1} \text{grey rectangle}}_{\text{variance}}\right)$$

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

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:

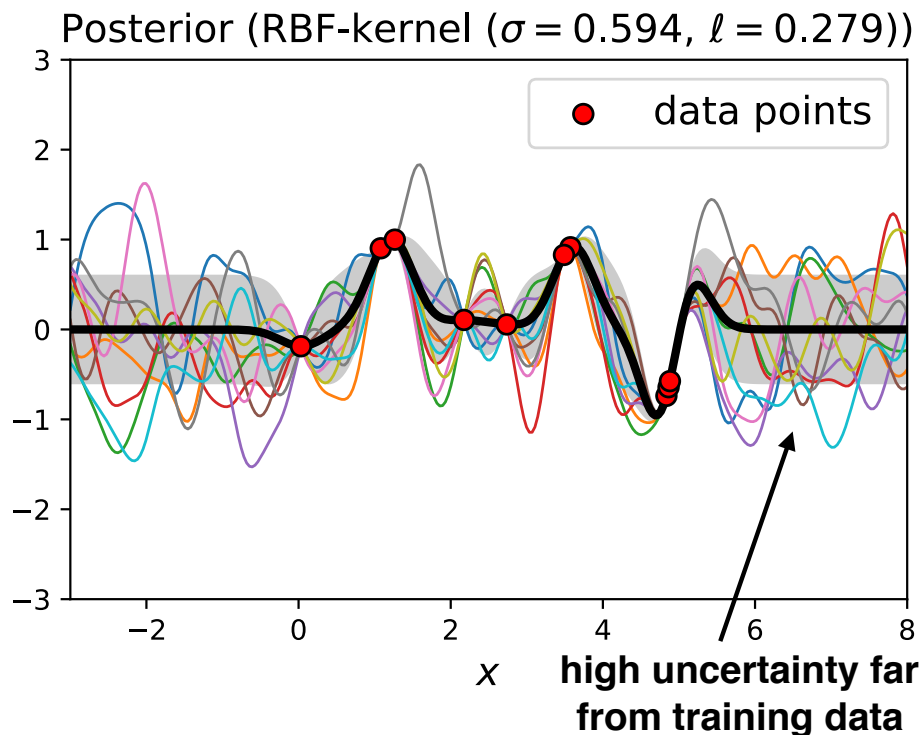
$a(x)$ – function of x
and training data

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: $O(N^3)$

$$p(\bullet\bullet\bullet\bullet) = \mathcal{N}\left(\begin{bmatrix} 0 \end{bmatrix}, \begin{bmatrix} \text{red square} \end{bmatrix}\right) \rightarrow \max_{\sigma_1, \sigma_2, \sigma_3, \ell}$$

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

$$p(\bullet) = \mathcal{N}\left(\underbrace{\begin{bmatrix} \text{grey rectangle} & \text{red square} \end{bmatrix}}_{\text{mean}} \begin{bmatrix} -1 \\ \bullet \\ -1 \end{bmatrix}, \underbrace{\begin{bmatrix} \text{blue square} & \text{grey rectangle} & \text{red square} \end{bmatrix}}_{\text{variance}} \begin{bmatrix} -1 \\ \bullet \\ -1 \end{bmatrix}\right)$$

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

Questions

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:

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

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

Exercise

x	y
-1.5	1
0.5	3
0.7	2.5

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

For $x_* = 0$:

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

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

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

For $x_* = 0$:

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

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

Formulas:

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

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