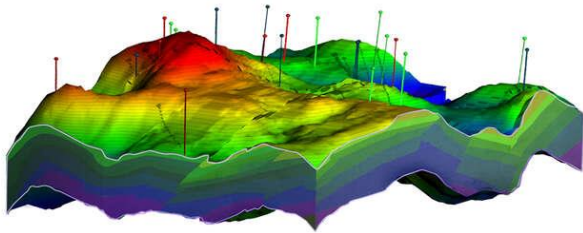DNV

# Introduction to **Gaussian Processes** for surrogate modelling

HIPERWIND PhD Summer School
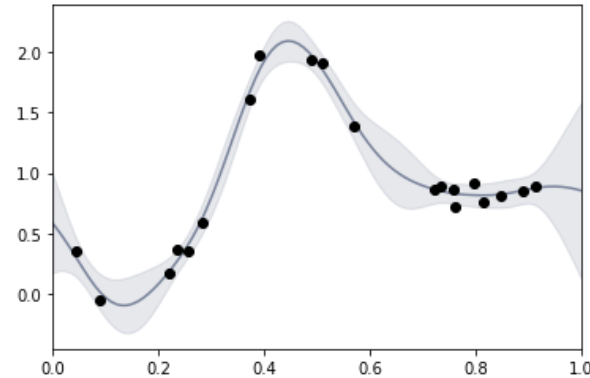
Christian Agrell

29 August 2023

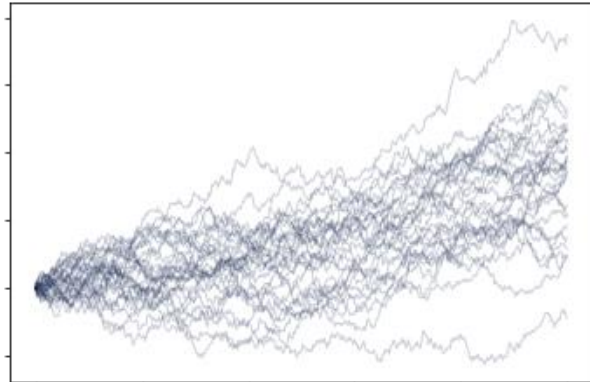# Gaussian processes appear in many different fields

**Geostatistics**
(Kriging)
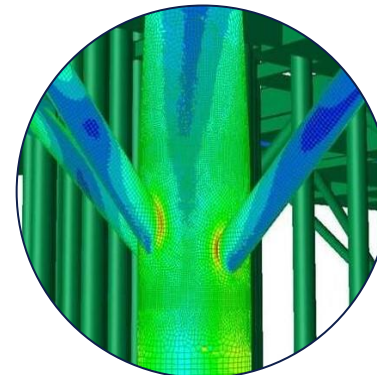
**Machine Learning**

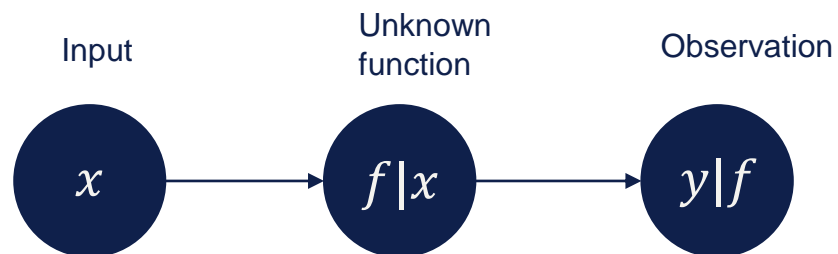**Stochastic Differential Equations**

**Uncertainty Quantification**

Computer experiment

Laboratory experiment

# Bayesian nonparametric function estimation

Input      Unknown function      Observation

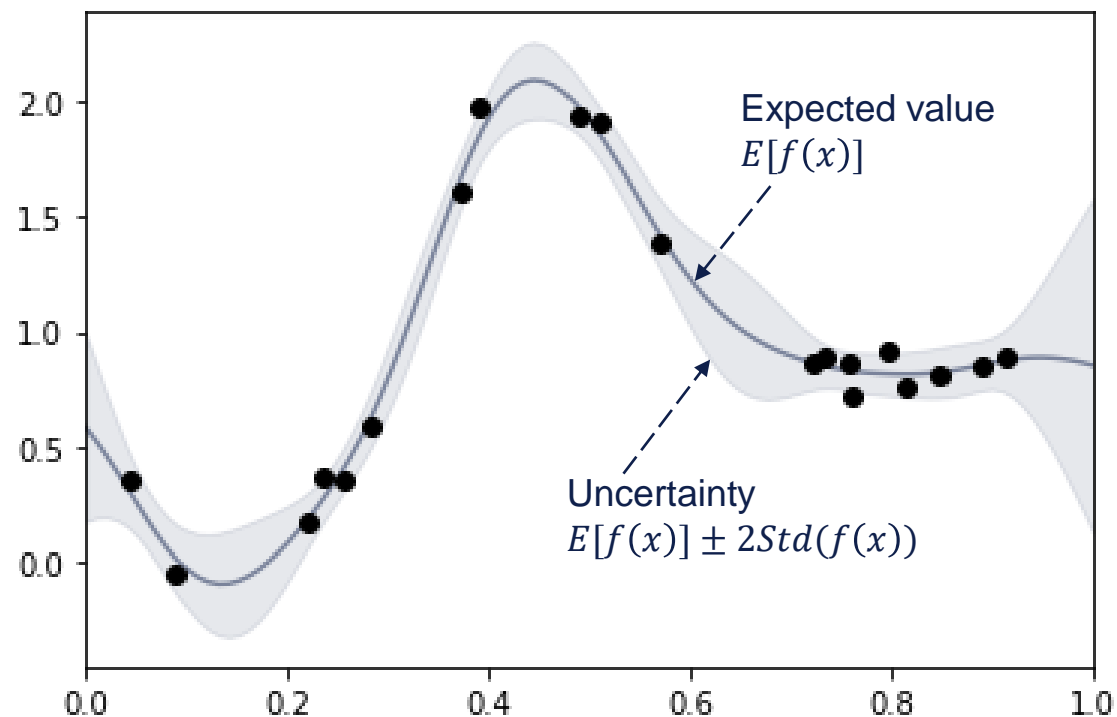$$x \rightarrow f|x \rightarrow y|f$$

**Goal**

*Infer* the function $f(x)$, given a set of observations $D = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

**Canonical case**

Input and output: $x \in \mathbb{R}^d, \ f(x) \in \mathbb{R}, \ y \in \mathbb{R}$
Observations: $y_i = f(x_i) + \varepsilon_i, \ \varepsilon_i = $ noise

Expected value
$E[f(x)]$

Uncertainty
$E[f(x)] \pm 2Std(f(x))$

# Preliminary
## The Gaussian conditional distribution

# Preliminary:
# **Multivariate Gaussian conditional distribution**

Let $X_1$ and $X_2$ be Gaussian random variables with means $\mu_1$ and $\mu_2$, variances $\sigma_1^2$ and $\sigma_2^2$, and correlation $\rho$.

- You observe $X_2 = x$.

Then $X_1 | X_2 = x$ is Gaussian with mean and variance:

$$E[X_1 | X_2 = x] = \mu_1 + \frac{\sigma_1}{\sigma_2} \rho (x - \mu_2)$$

$$Var[X_1 | X_2 = x] = (1 - \rho^2)\sigma_1^2$$

DNV

# Preliminary:
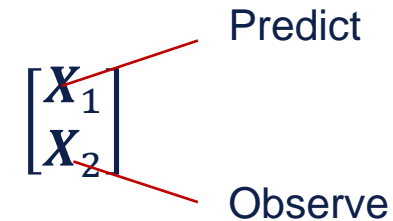# Multivariate Gaussian conditional distribution

Let $X_1$ and $X_2$ be Gaussian <u>vectors</u>, with joint mean and covariance

$$\mu = E\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$$

$$\Sigma = COV\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

$$[X_1, X_2] = [\ \underbrace{X_{11}, \dots, X_{1m}}_{N(\mu_1, \Sigma_{11})}, \ \underbrace{X_{21}, \dots, X_{2n}}_{N(\mu_2, \Sigma_{22})}\ ]$$

- You observe $X_1 = x$.

$$\begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$$

Predict

Observe

$X_1|X_2 = x$ is Gaussian with mean and variance:

$$E[X_1|X_2 = x] = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x - \mu_1)$$

$$Var[X_1|X_2 = x] = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$
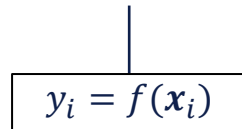
DNV

# Constructing a GP

# Constructing a Gaussian process (GP)

**Definition : Gaussian process**

$f(\cdot)$ has a Gaussian process (GP) distribution if for any $n \in \mathbb{N}$ and $x_1, \dots, x_n \in \mathcal{X}$, the joint distribution of $f(x_1), \dots, f(x_n)$ is multivariate normal

$$x_1, \dots, x_n$$

A finite collection from some domain $\mathcal{X}$ (typically $\mathbb{R}^N$)

$$y_i = f(x_i)$$

A mapping $f$

$$[y_1, \dots, y_n]$$

An $n$-dimensional Gaussian vector

# Constructing a Gaussian process (GP)

**Definition : Gaussian process**

$f(\cdot)$ has a Gaussian process (GP) distribution if for any $n \in \mathbb{N}$ and $x_1, \dots, x_n \in \mathcal{X}$, the joint distribution of $f(x_1), \dots, f(x_n)$ is multivariate normal

- A GP over functions $f : \mathcal{X} \to \mathbb{R}$ is completely specified by it's mean function $\mu$ and covariance function (kernel) $k$, where

$$\mu(x) = E[f(x)]$$
$$k(x, x') = cov(f(x), f(x'))$$

and we write $f \sim GP(\mu, k)$.

**Definition : Covariance function**

A function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ that is

- Symmetric

- Positive semi-definite

**Definition : Symmetric**
$$k(x_i, x_j) = k(x_j, x_i)$$
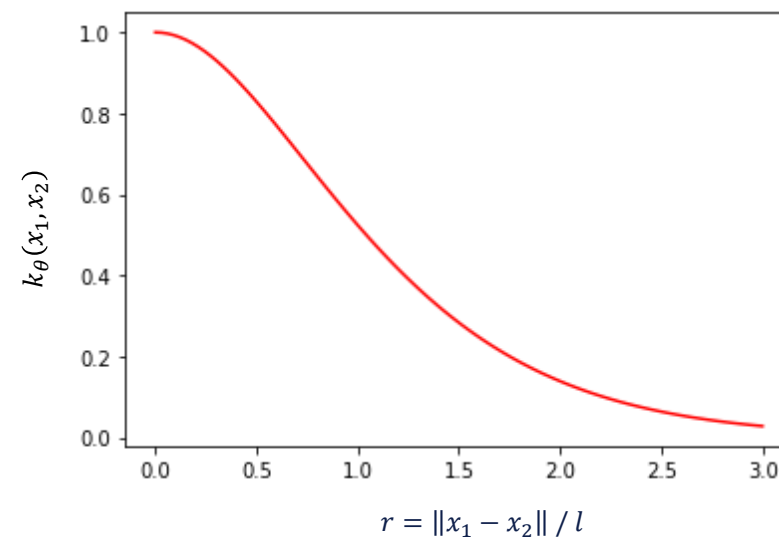
**Definition :** Positive semi-definite

$$\sum_{i,j=1}^{n} \lambda_i \lambda_j \, k(x_i, x_j) \geq 0$$

For all $x_1, \dots, x_n \in \mathcal{X}$ and $\lambda_1, \dots, \lambda_n \in \mathbb{R}$

# The covariance function (kernel)

**Recall:** $f \sim GP(\mu, k)$

- Assume $\mu = 0$
  (Or that we work with $f - \mu$)

- Assume $k$ is stationary: $k(x_1, x_2)$ can be written as $k(x_1 - x_2)$
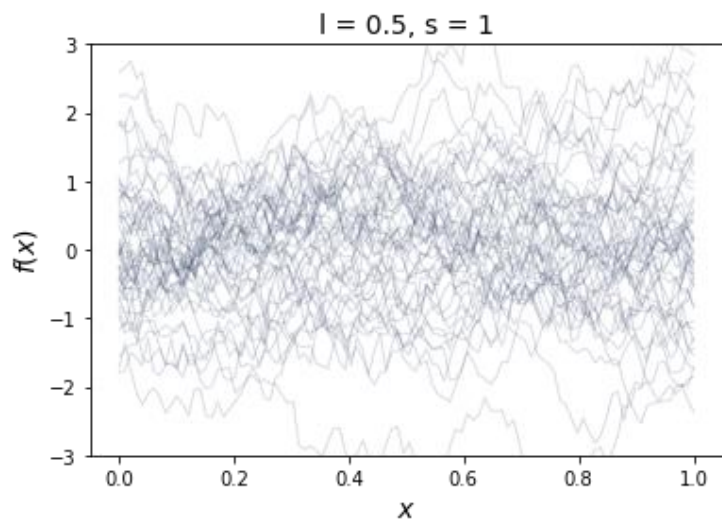  (This is often used in practice)

# The covariance function (kernel)

Different covariance functions give different «types» of functions that the GP can represent

- Differentiable
- Periodicity
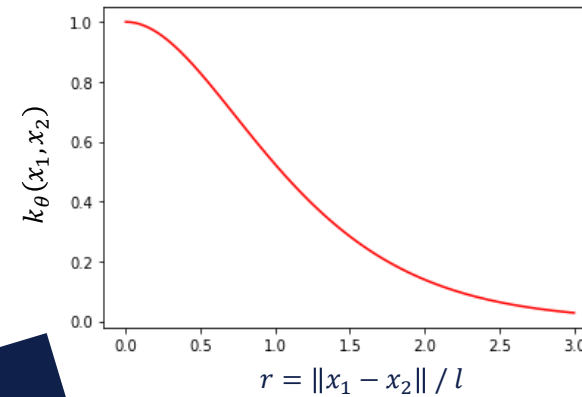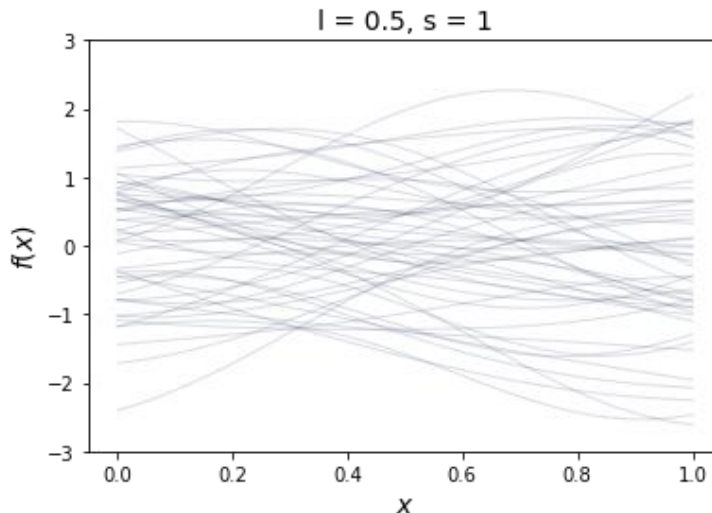- Stationary / non-stationary
- Linear trend



**Exponential**

$$k_\theta(x_1, x_2) = s^2 e^{-r}$$

**Gaussian**

$$k_\theta(x_1, x_2) = s^2 e^{-\frac{1}{2}r^2}$$

**Matérn 5/2**

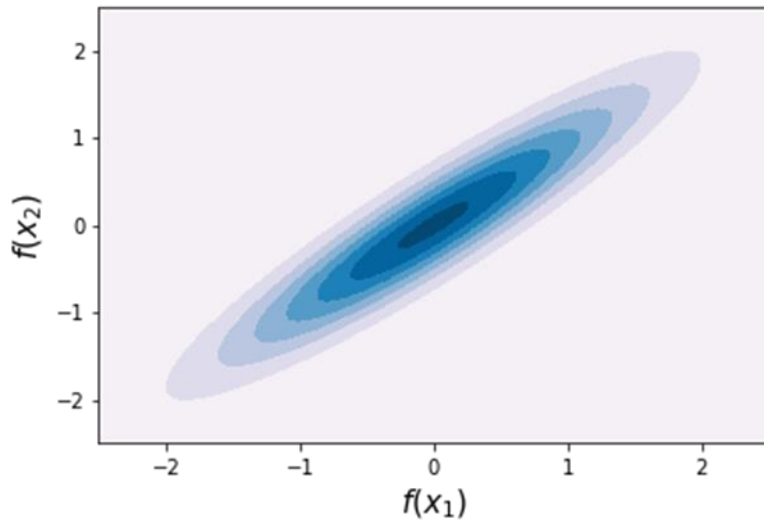$$k_\theta(x_1, x_2) = s^2 \left(1 + \sqrt{5}r + \frac{5}{3}r^2\right)e^{-\sqrt{5}r}$$

# 3 ways to think about GPs

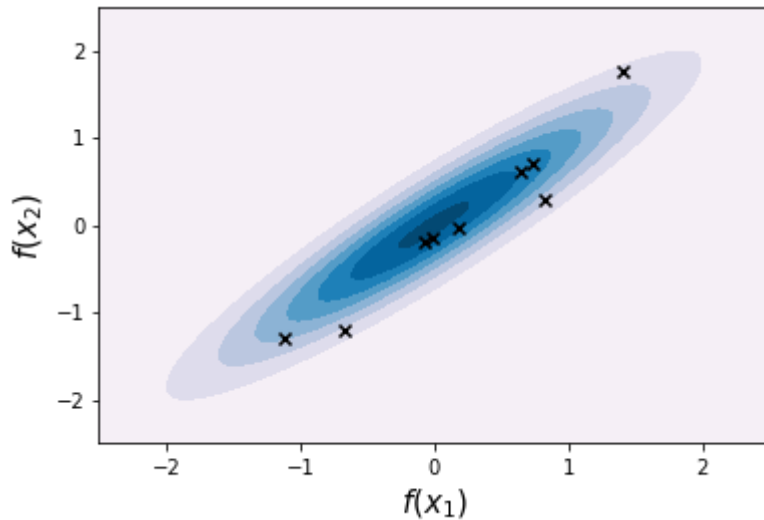# 3 Ways to think about GPs
## 1) A large vector

- Assume $X$ is finite. $X = \{x_1, \ldots, x_N\}$

- Then a GP is just a mapping from $X$
  to components of the vector
  $[f(x_1), \ldots, f(x_N)]$

DNV

# 3 Ways to think about GPs
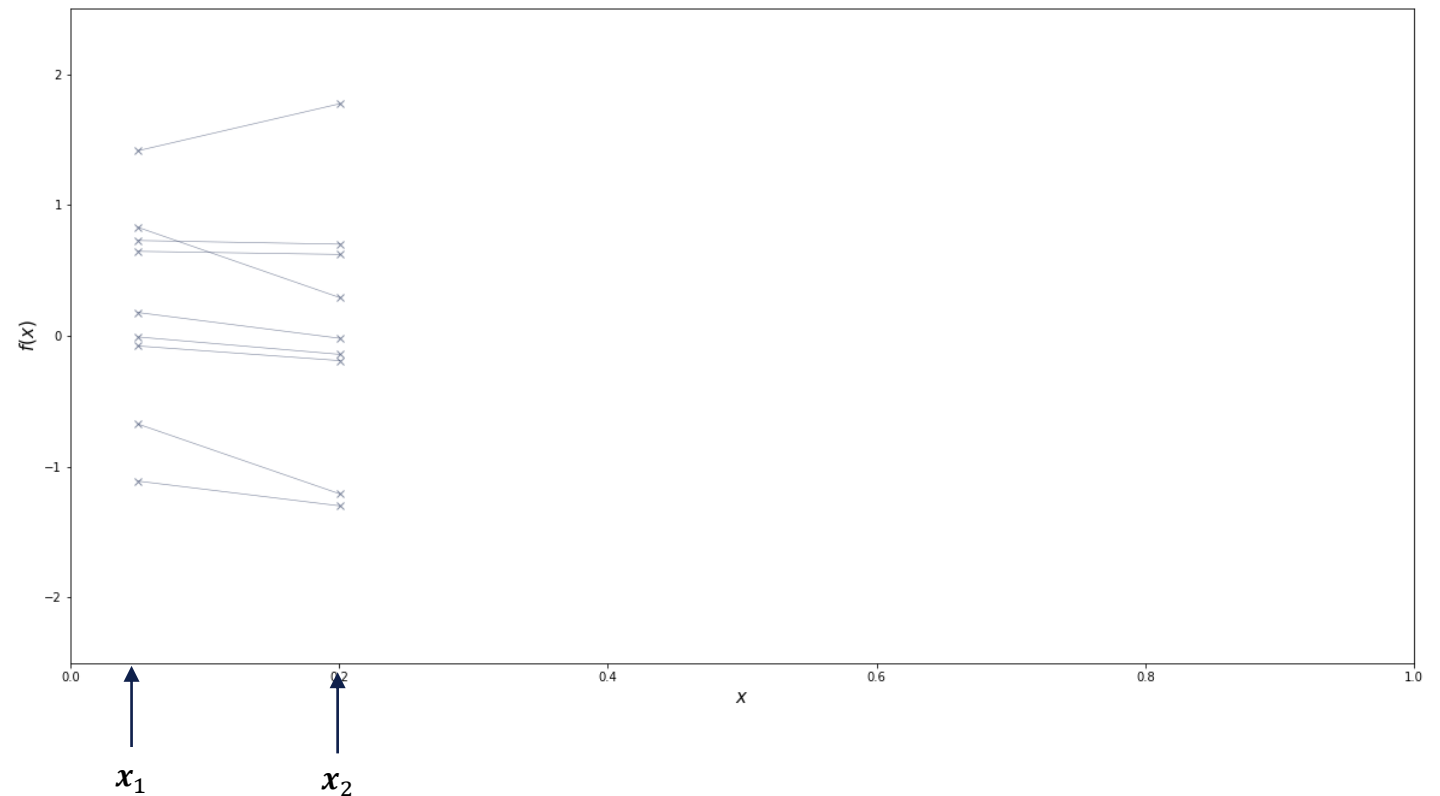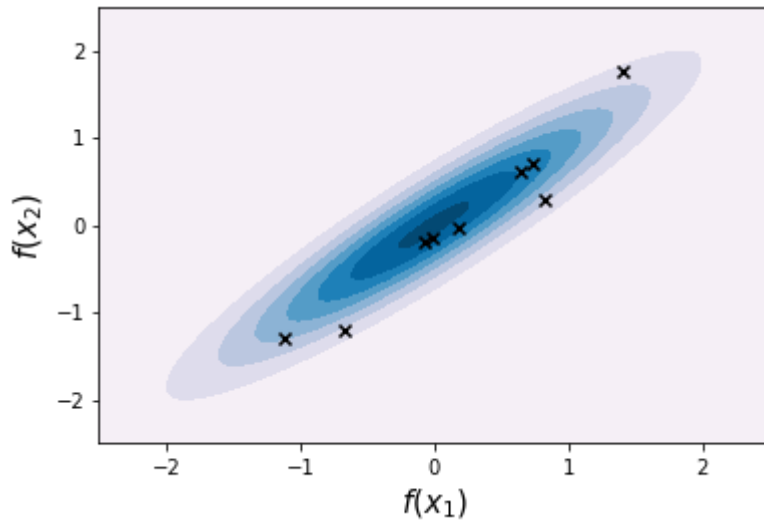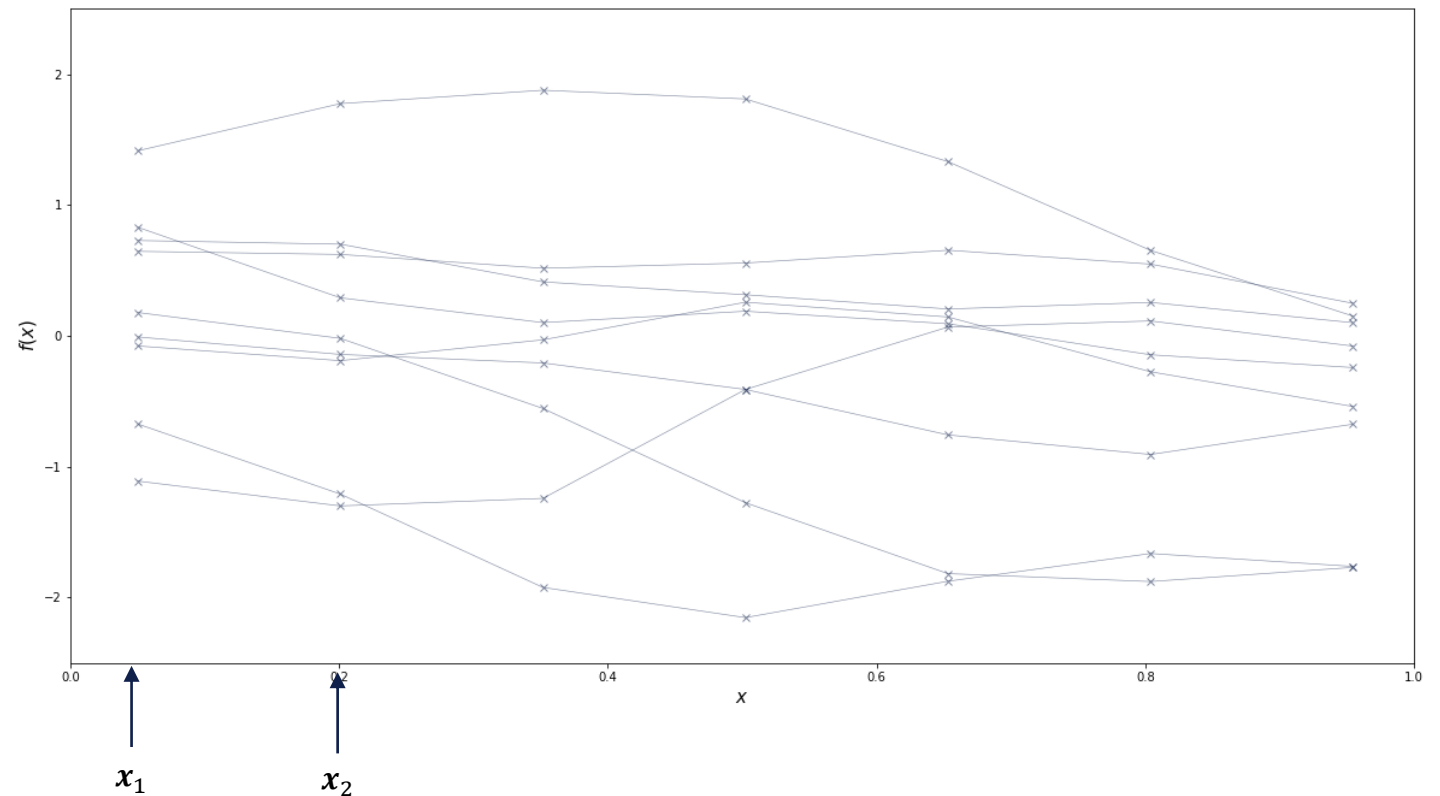## 1) A large vector

- Assume $\mathcal{X}$ is finite. $\mathcal{X} = \{\boldsymbol{x}_1, ..., \boldsymbol{x}_N\}$

- Then a GP is just a mapping from $\mathcal{X}$
  to components of the vector
  $[f(\boldsymbol{x}_1), ..., f(\boldsymbol{x}_N)]$

DNV

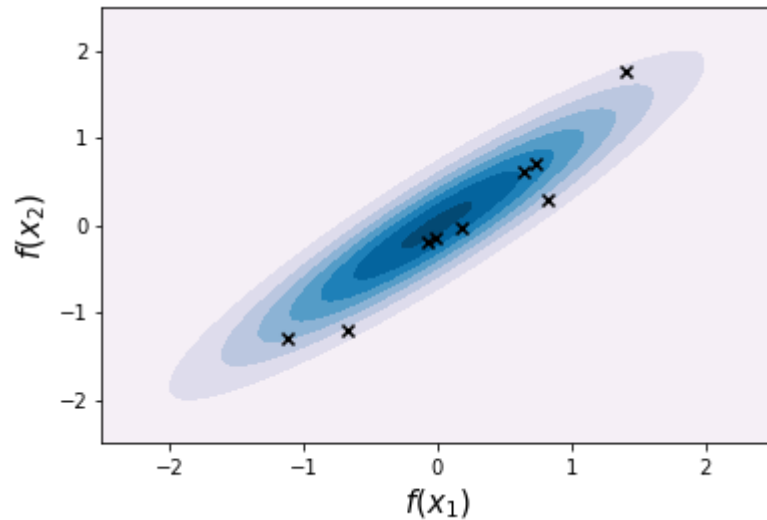# 3 Ways to think about GPs
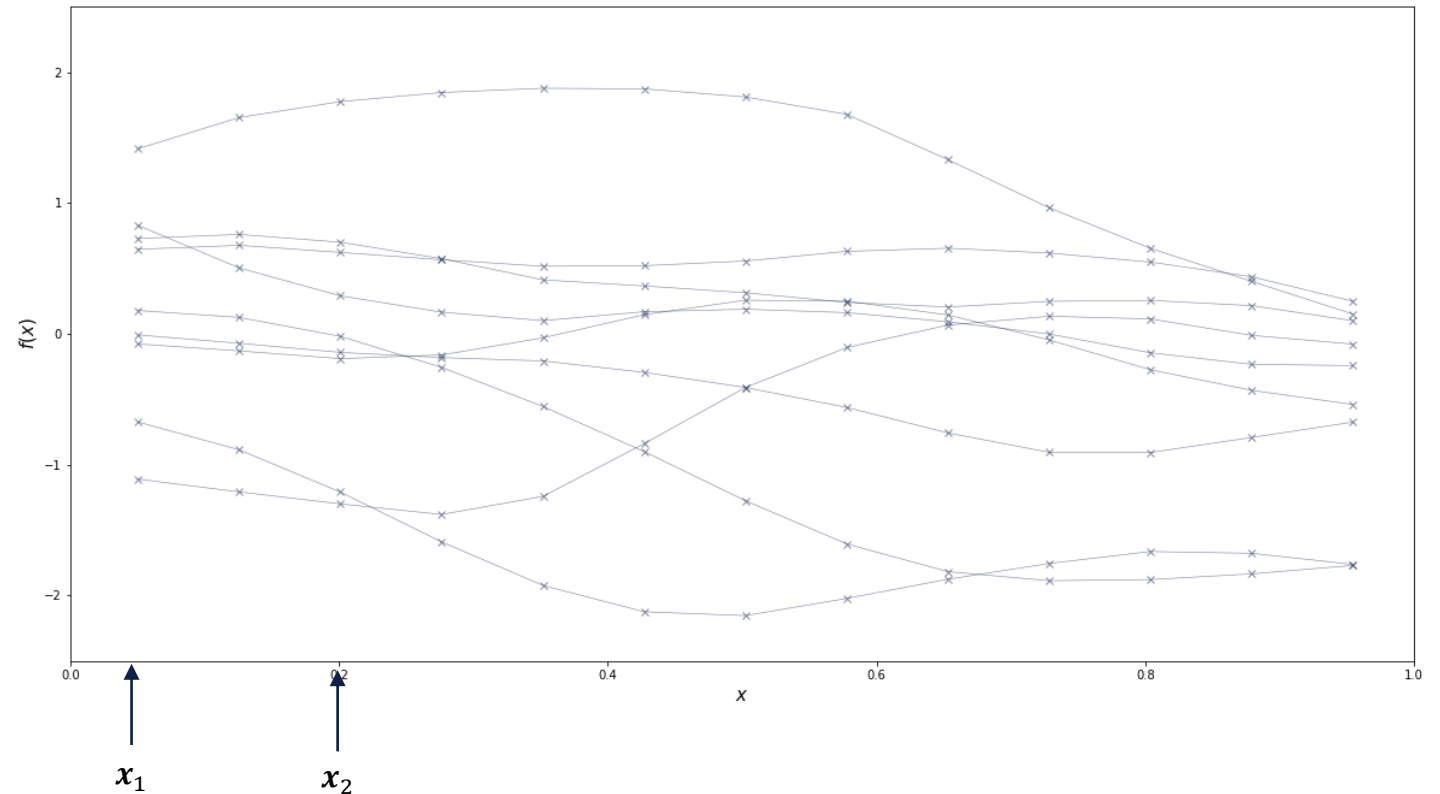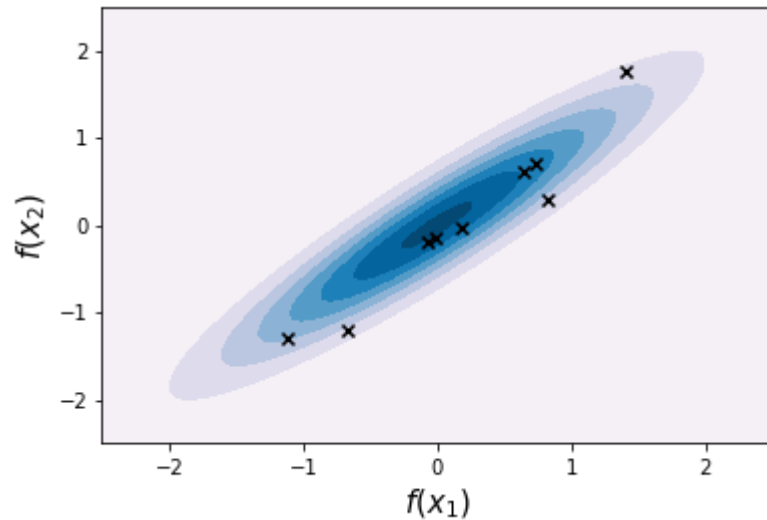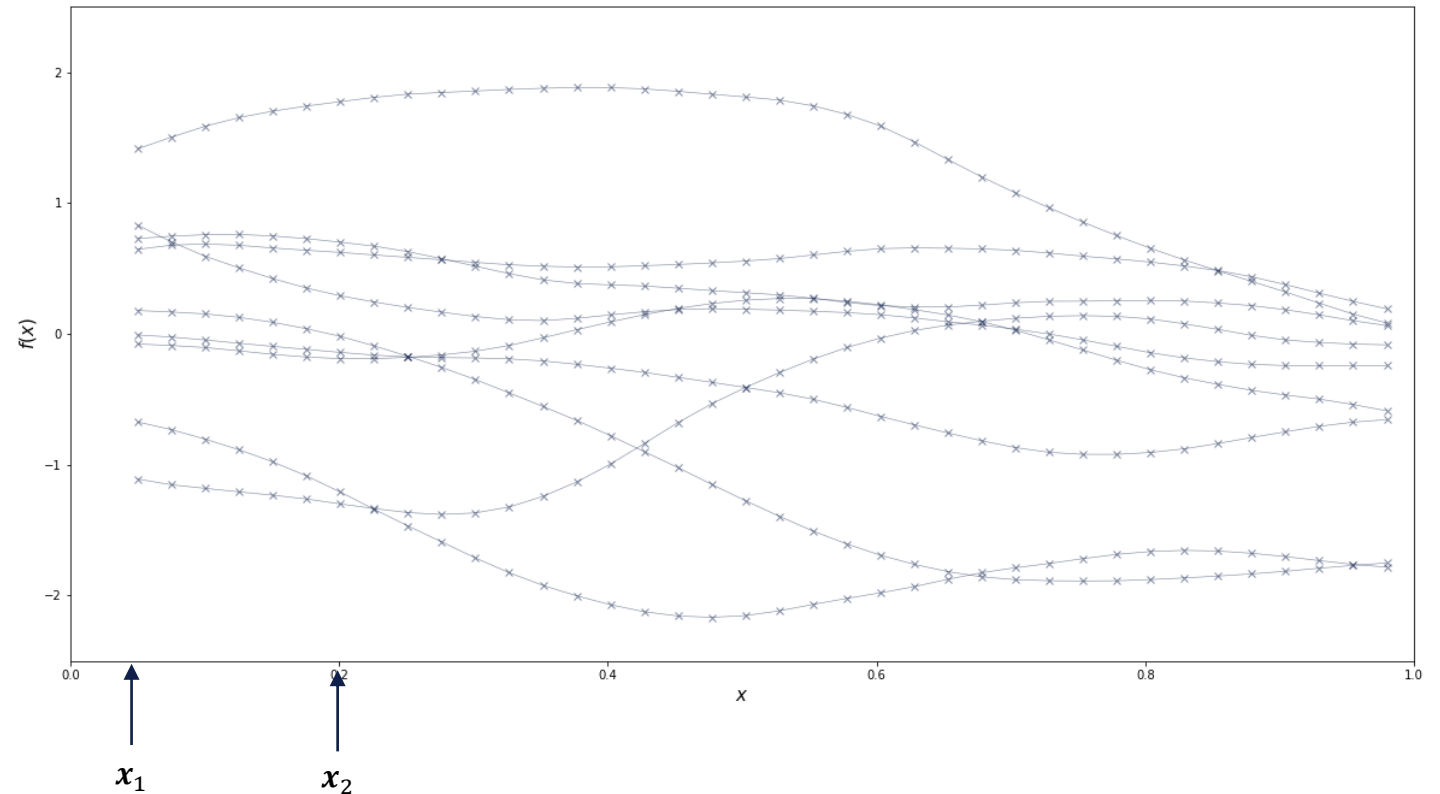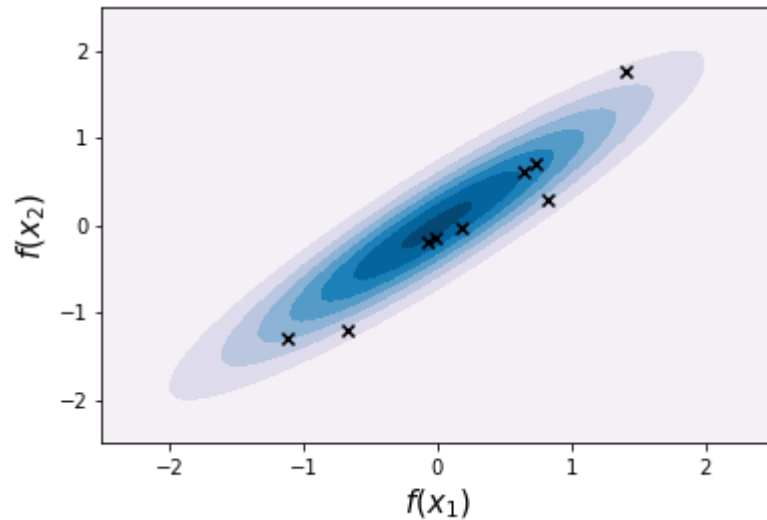## 1) A large vector

- Assume $\mathcal{X}$ is finite. $\mathcal{X} = \{x_1, \ldots, x_N\}$

- Then a GP is just a mapping from $\mathcal{X}$ to components of the vector $[f(x_1), \ldots, f(x_N)]$

# 3 Ways to think about GPs
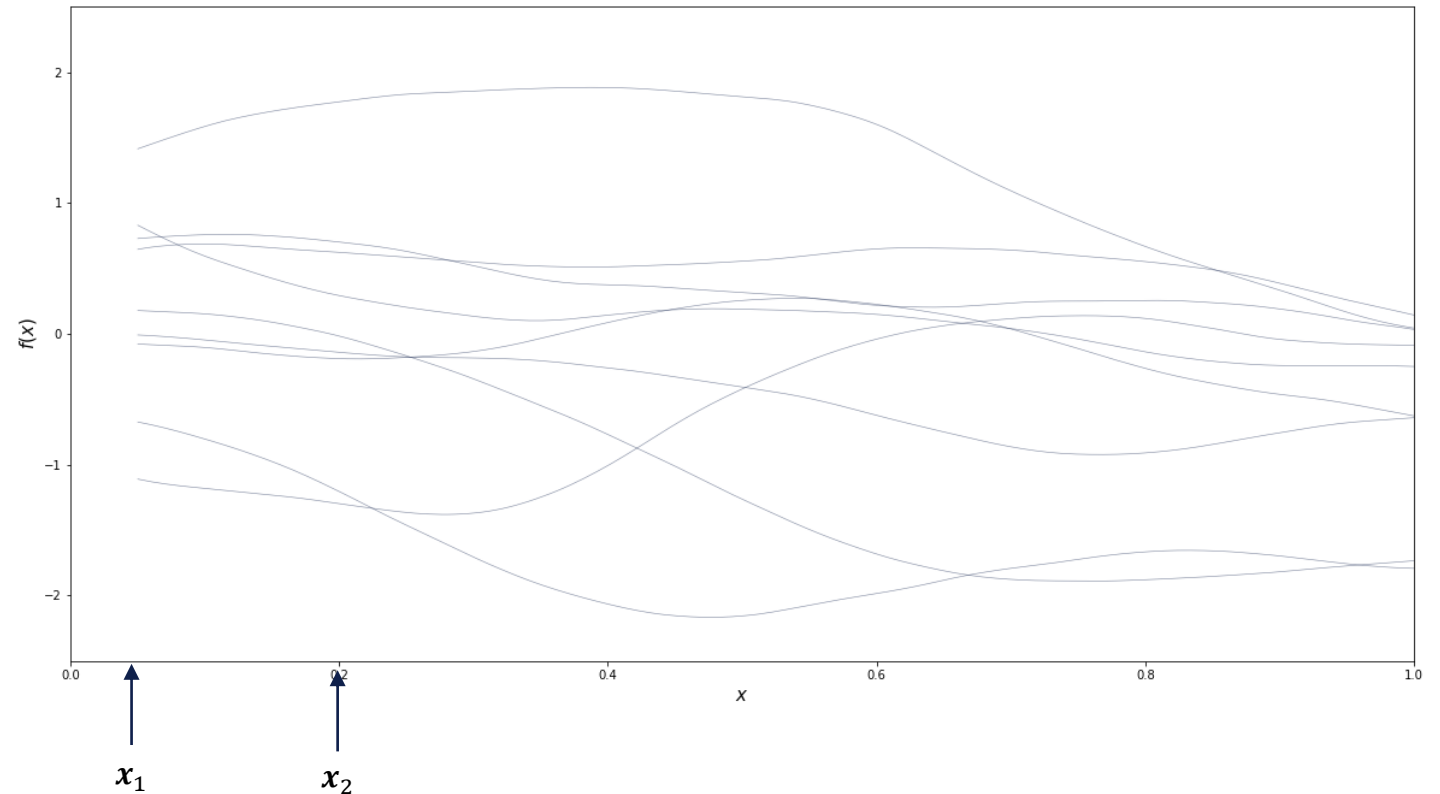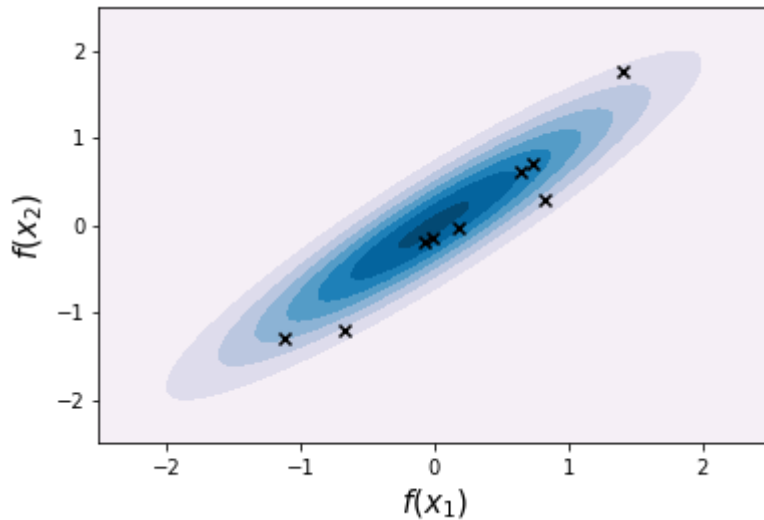## 1) A large vector

- Assume $\mathcal{X}$ is finite. $\mathcal{X} = \{\boldsymbol{x}_1, \dots, \boldsymbol{x}_N\}$

- Then a GP is just a mapping from $\mathcal{X}$ to components of the vector $[f(\boldsymbol{x}_1), \dots, f(\boldsymbol{x}_N)]$

# 3 Ways to think about GPs
## 1) A large vector

- Assume $\mathcal{X}$ is finite. $\mathcal{X} = \{\boldsymbol{x}_1, \dots, \boldsymbol{x}_N\}$

- Then a GP is just a mapping from $\mathcal{X}$
  to components of the vector
  $[f(\boldsymbol{x}_1), \dots, f(\boldsymbol{x}_N)]$

DNV

# 3 Ways to think about GPs
## 1) A large vector

- Assume $\mathcal{X}$ is finite. $\mathcal{X} = \{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\}$

- Then a GP is just a mapping from $\mathcal{X}$ to components of the vector $[f(\boldsymbol{x}_1), \ldots, f(\boldsymbol{x}_N)]$

DNV

# 3 Ways to think about GPs
## 1) A large vector

- Assume $\mathcal{X}$ is finite. $\mathcal{X} = \{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\}$

- Then a GP is just a mapping from $\mathcal{X}$ to components of the vector $[f(\boldsymbol{x}_1), \ldots, f(\boldsymbol{x}_N)]$

DNV

# 3 Ways to think about GPs
## 1) A large vector

- Assume $\mathcal{X}$ is finite. $\mathcal{X} = \{\boldsymbol{x}_1, ..., \boldsymbol{x}_N\}$

- Then a GP is just a mapping from $\mathcal{X}$
  to components of the vector
  $[f(\boldsymbol{x}_1), ..., f(\boldsymbol{x}_N)]$

# 3 Ways to think about GPs
## 1) A large vector

- Assume $\mathcal{X}$ is finite. $\mathcal{X} = \{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\}$

- Then a GP is just a mapping from $\mathcal{X}$
  to components of the vector
  $[f(\boldsymbol{x}_1), \ldots, f(\boldsymbol{x}_N)]$

DNV

# 3 Ways to think about GPs
## 2) A distribution over functions

- The GP is a *function-valued random variable*

- It is a distribution over functions

Reproducing Kernel
Hilbert Space

$$f \sim GP(\mu(\boldsymbol{x}), k(\boldsymbol{x}, \boldsymbol{x}'))$$

There is a space of functions associated with the kernel of the GP

This is called the **Reproducing Kernel Hilbert Space**

This connection is very useful for theoretical work!

DNV

# 3 Ways to think about GPs
## 3) An infinite-parameter model

- Let

$$f(\boldsymbol{x}) = \sum_{i=0}^{N} \lambda_i \xi_i \, \phi_i(\boldsymbol{x}) \qquad (1)$$

**Karhunen-Loève expansion:**

Any GP can be written as (1) with $N = \infty$

where

$\lambda_i$ = constant

$\phi(\boldsymbol{x})_i$ = deterministic function

$\xi_i$ = Standard normal variable (pairwise independent)

Then $f$ is a GP.



Functions corresponding to different samples of the vector $[\xi_1, \xi_2, \dots]$

DNV

# Conditioning on data

# Data

$(x_1, y_1), (x_2, y_2), \ldots$



# Model

## Some different Gaussian process priors

DNV

# Data

$(x_1, y_1), (x_2, y_2), \dots$



# Model

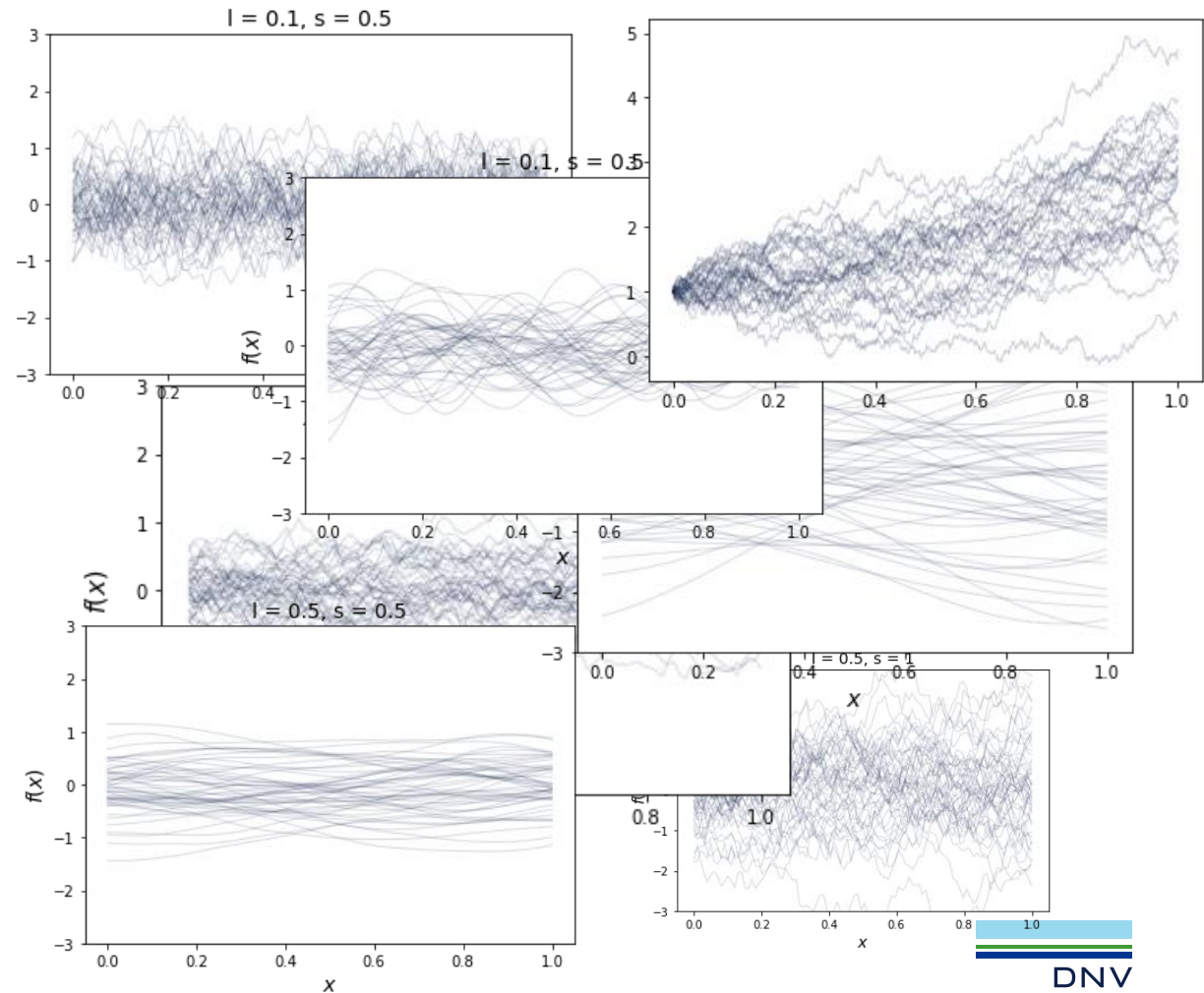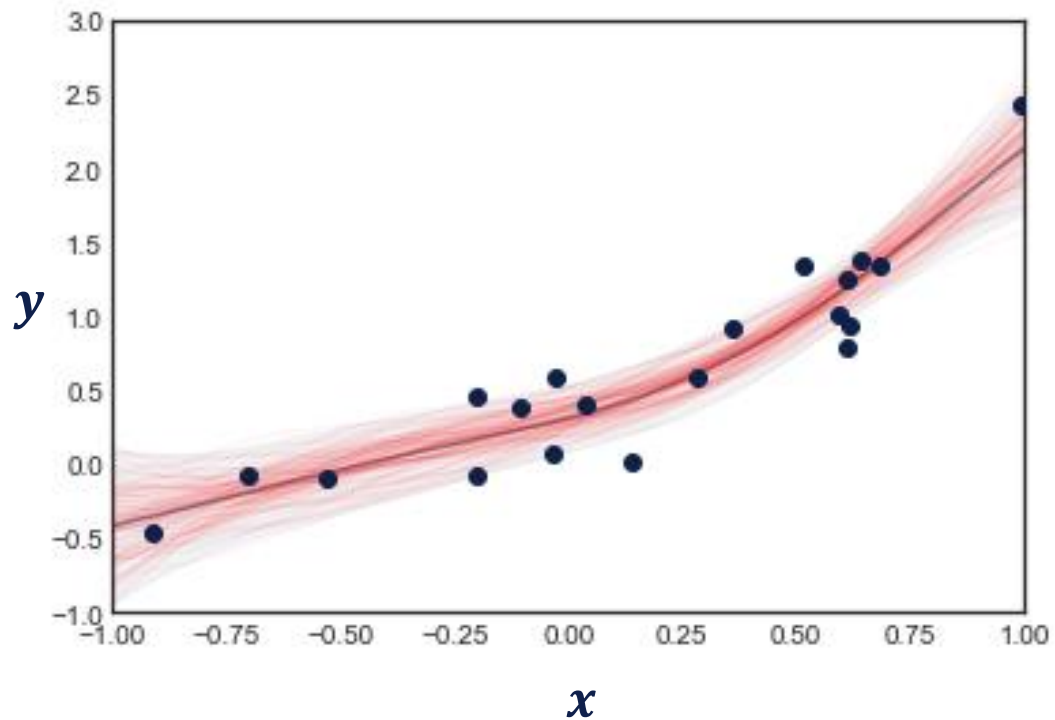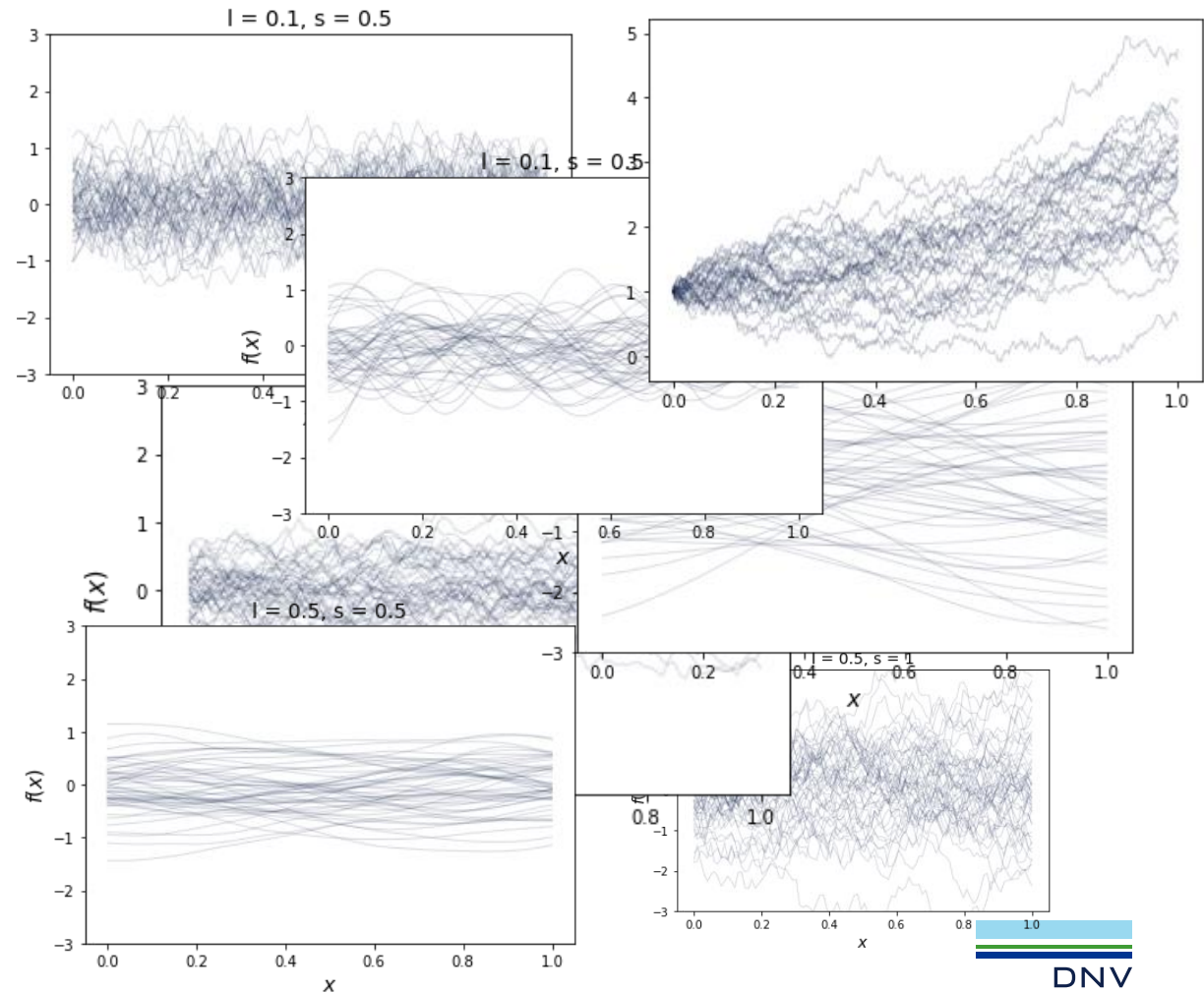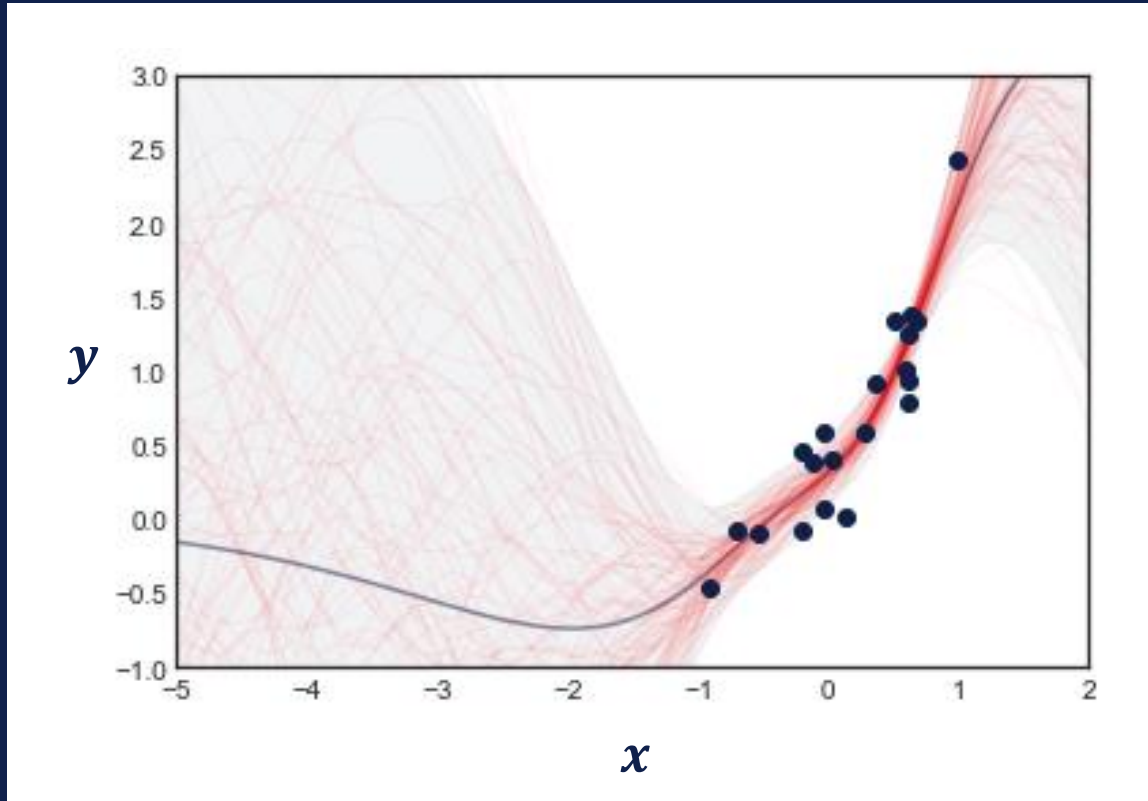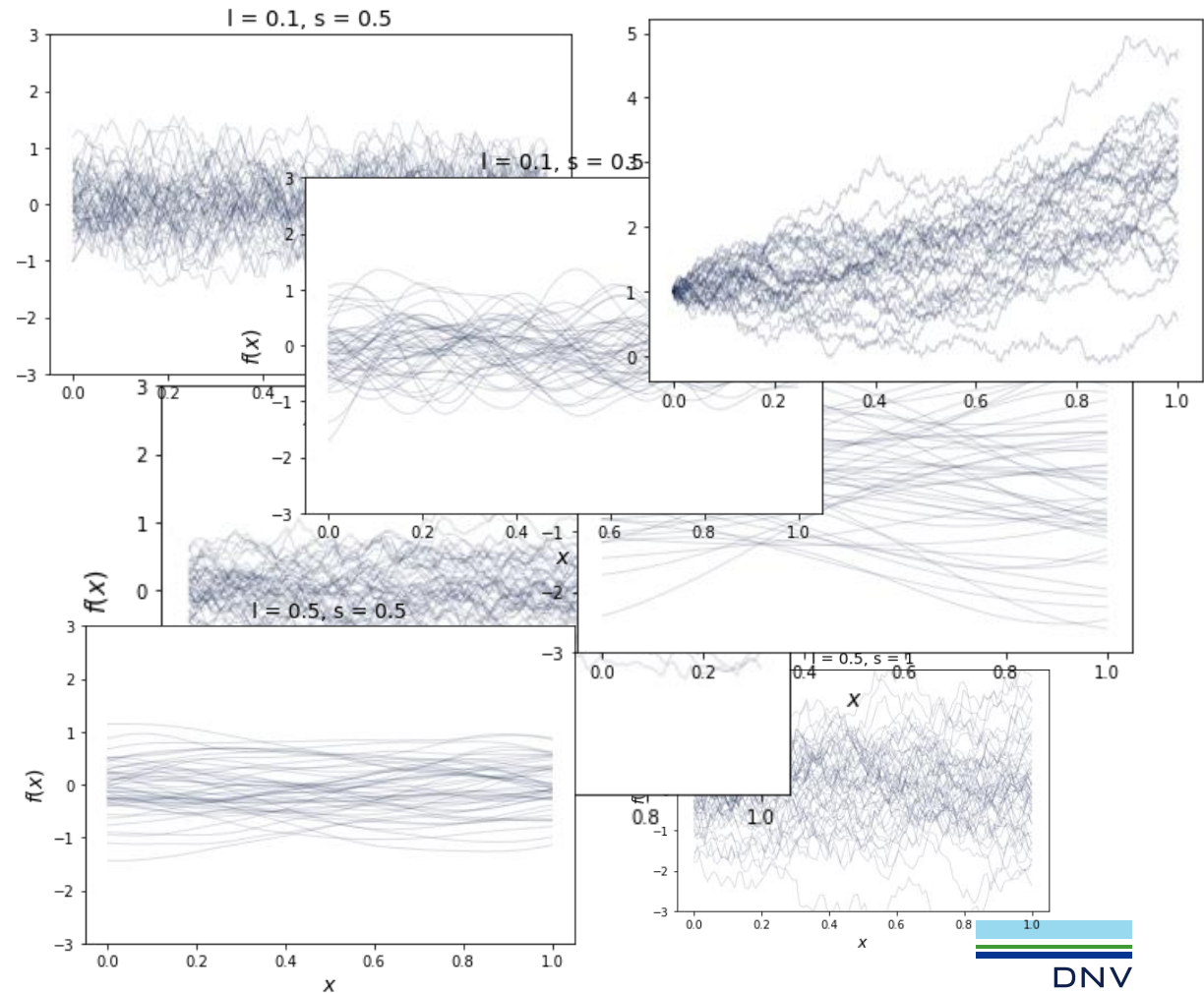## Some different Gaussian process priors

DNV

# Data

$(x_1, y_1), (x_2, y_2), \ldots$



# Model

Some different Gaussian process priors

DNV

# Data

$(x_1, y_1), (x_2, y_2), \ldots$


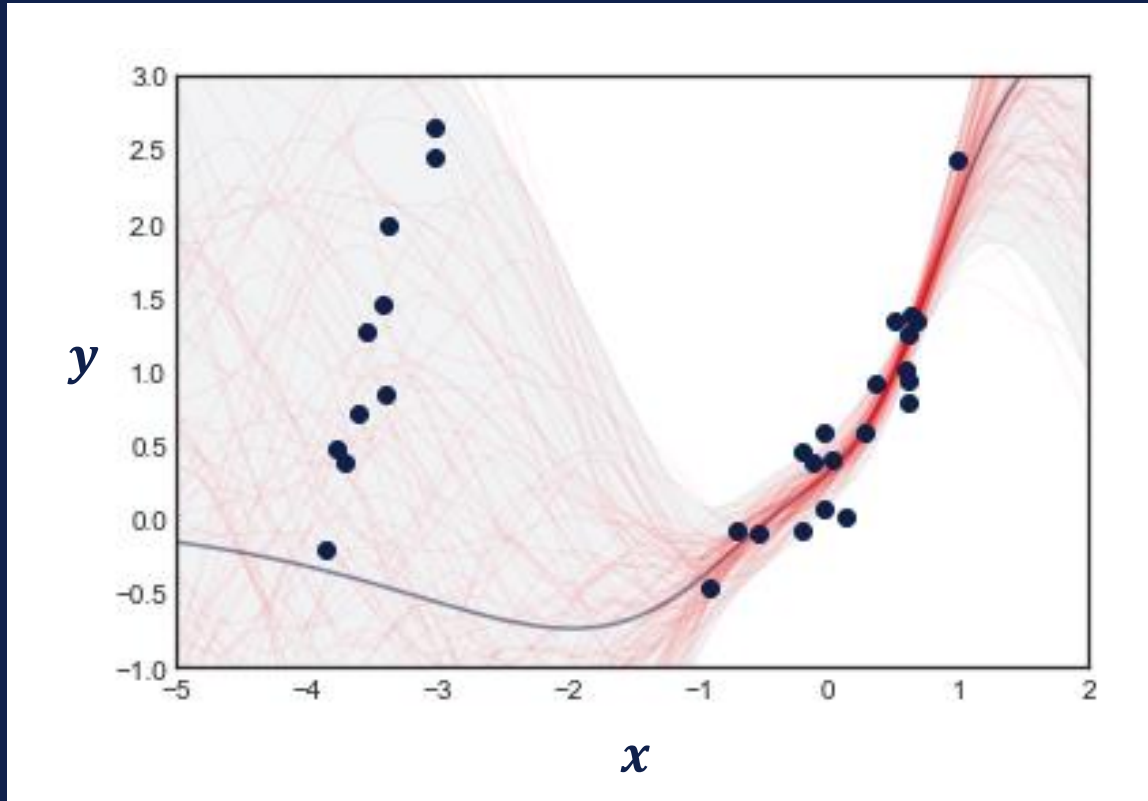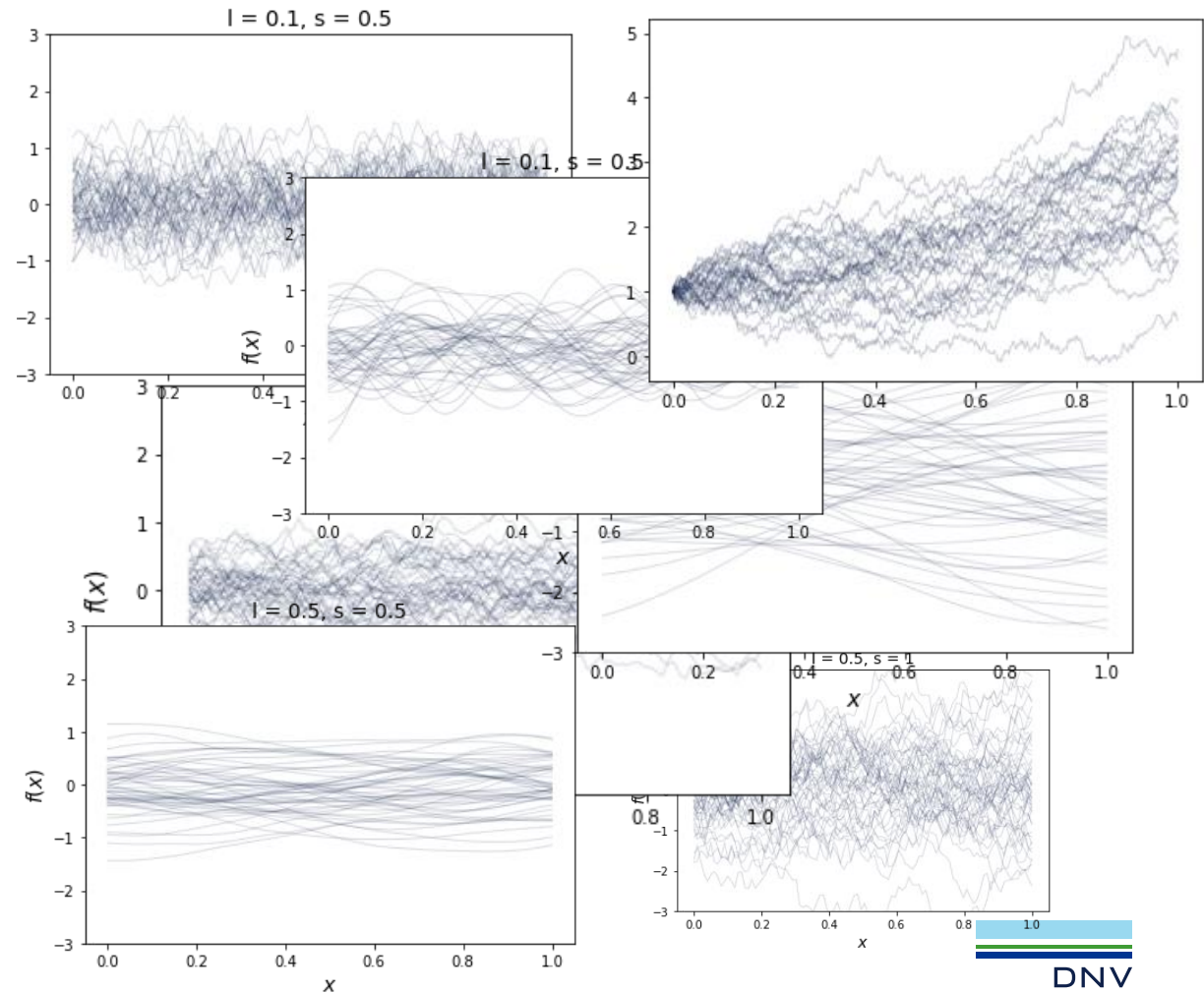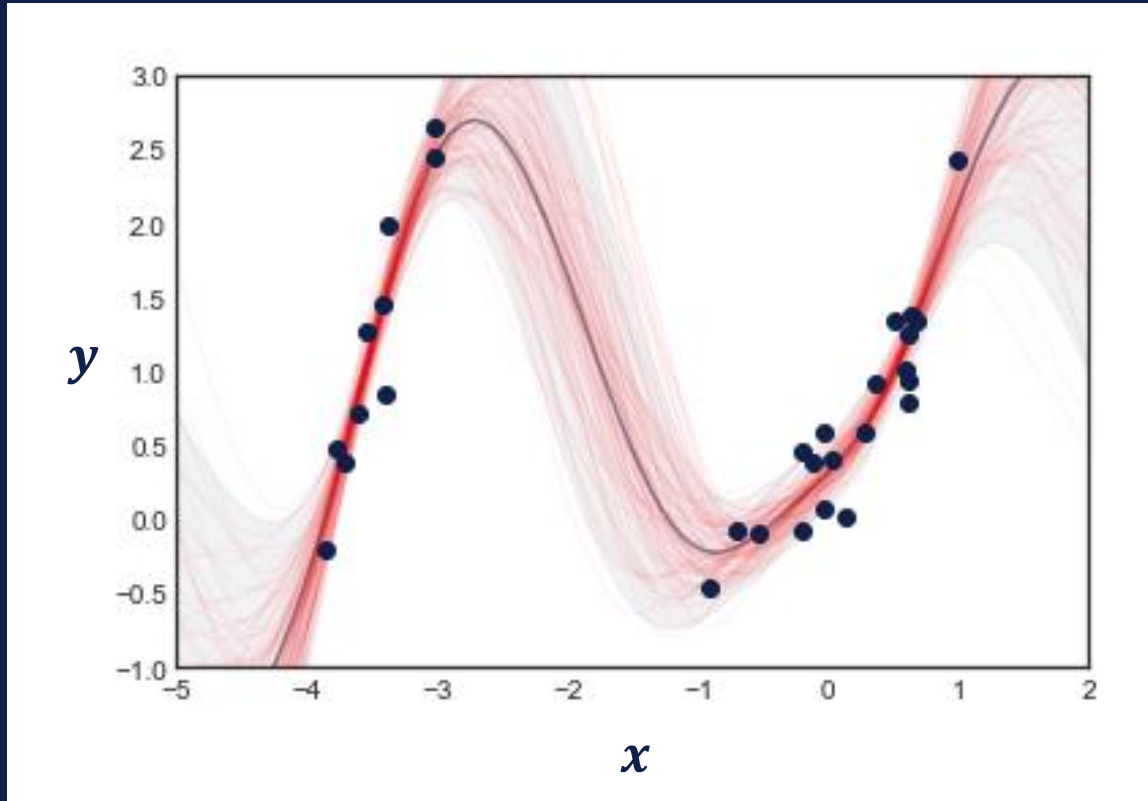
# Model

## Some different Gaussian process priors

DNV

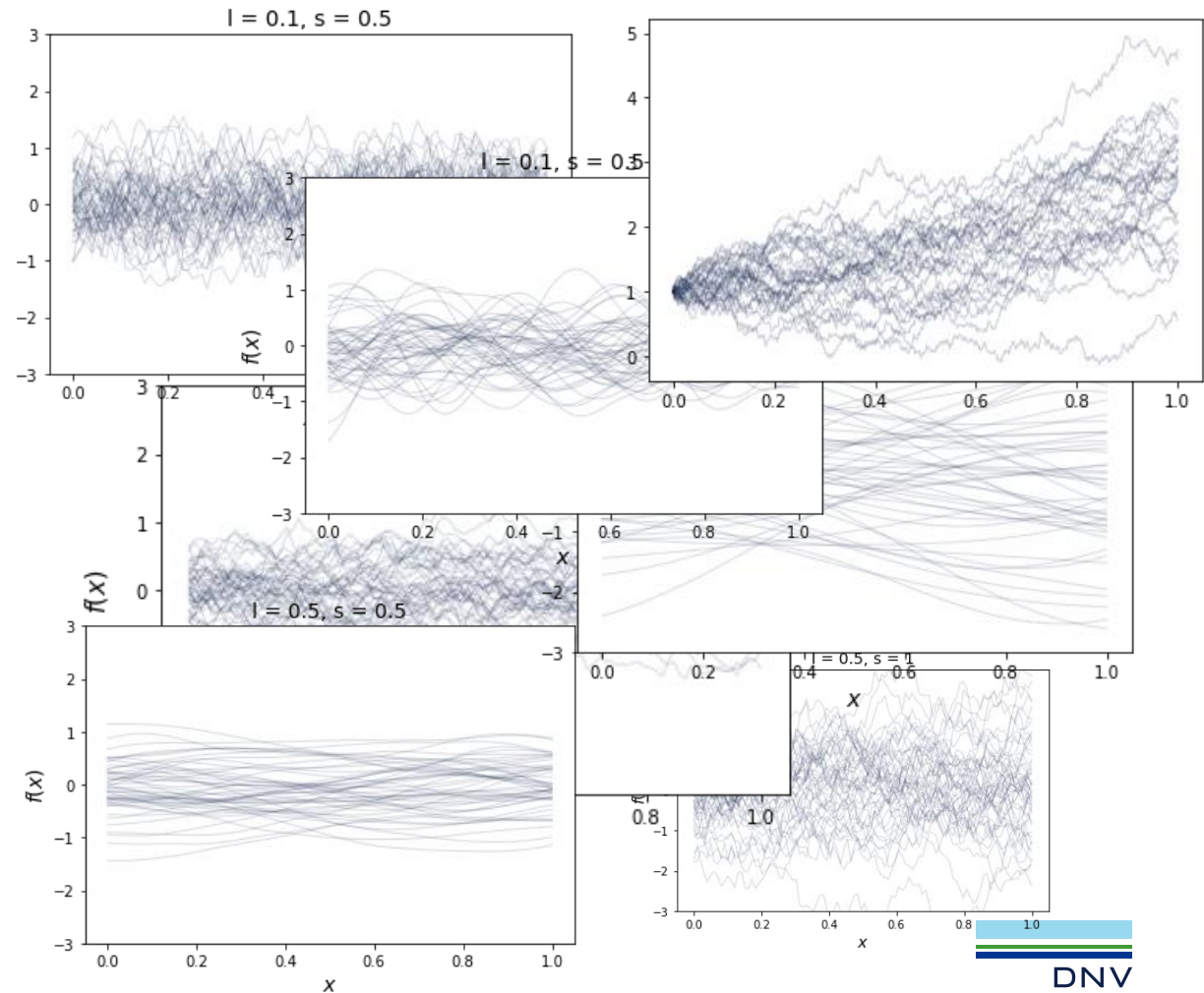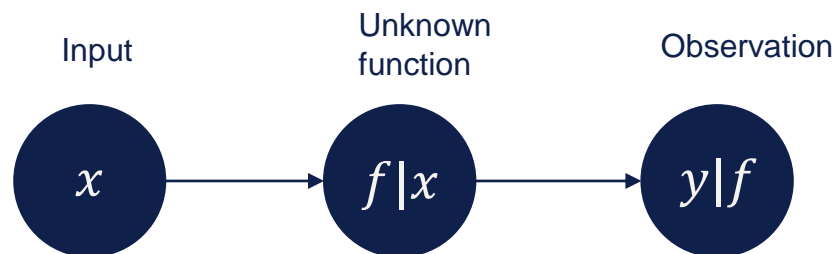# Data

$(x_1, y_1), (x_2, y_2), \dots$



# Model

## Some different Gaussian process priors

DNV

# Conditioning on a set of observation

Input    Unknown function    Observation

$x$ → $f|x$ → $y|f$

**Goal**

*Infer* the function $f(x)$, given a set of observations D=$\{(x_1, y_1), \dots, (x_n, y_n)\}$

**Canonical case**

Input and output: $x \in \mathbb{R}^d$,   $f(x) \in \mathbb{R}$, $y \in \mathbb{R}$
Observations: $y_i = f(x_i) + \varepsilon_i$, $\varepsilon_i$ = noise

**Standard GP regression**

- Assume the noise terms are i.i.d. $\varepsilon_i \sim N(0, \sigma^2)$.

- Let $f \sim GP(\mu, k)$.

For a new set of input locations, $x_1^*, \dots, x_M^*$, let $f^*|D$ denote the posterior process evaluated at each new input, $f^*|D = [f(x_1^*), \dots, f(x_M^*)] | D$.

- Then $f^*|D \sim N(\mu_{f^*|D}, \Sigma_{f^*|D})$ with

$$\mu_{f^*|D} = \mu^* + K^*(K + \sigma^2 I)^{-1}(Y - \mu)$$

$$\Sigma_{f^*|D} = K^{**} - K^*(K + \sigma^2 I)^{-1}(K^*)^T$$
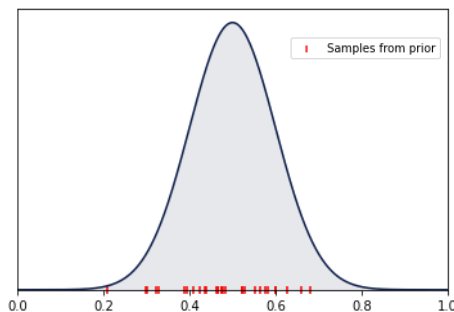
$O(n^3)$ computation
$O(n^2)$ memory

where $(\mu)_i = \mu(x_i)$, $(K)_{i,j} = k(x_i, x_j)$, $(K^*)_{i,j} = k(x_i^*, x_j)$ etc.
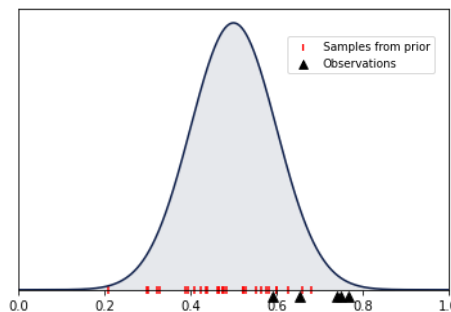
# GP as a prior over functions
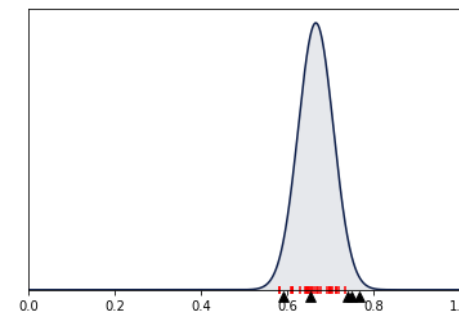
**Prior**

Prior distribution
$p(x|\theta)$

**Observations**

$y_i = x_{true} + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$

**Posterior**

Bayesian inference

(find $x_{true}$)

**Prior**

Prior process
$GP(\mu(x|\theta), k(x, x'|\theta))$

**Observations**

$y_i = f_{true}(x_i) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$

**Posterior**

Bayesian inference over functions

(find $f_{true}$)

# GP as a prior over functions



Prior distribution
$p(x|\theta)$

**Prior**

**Observations**
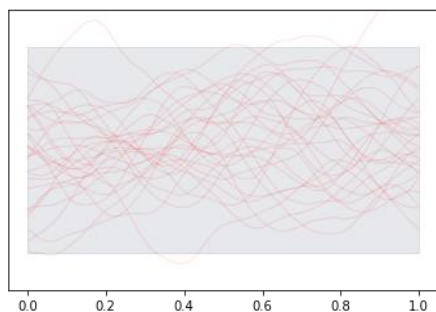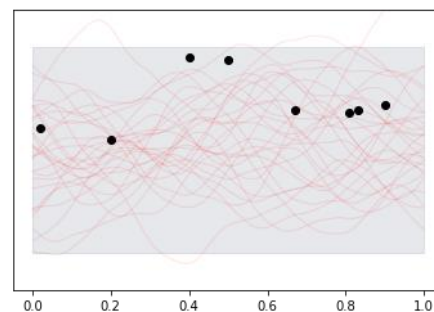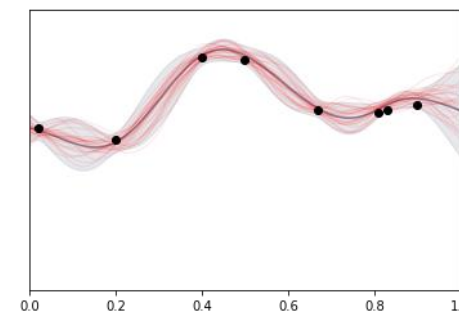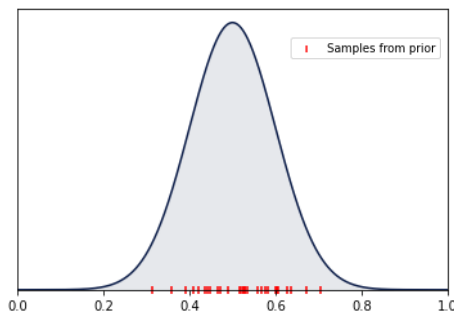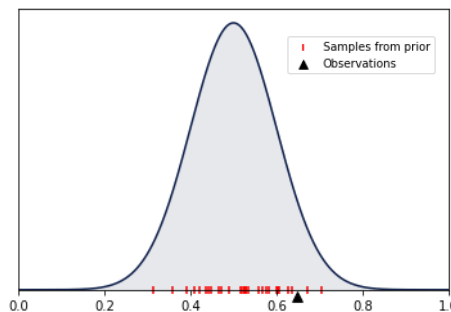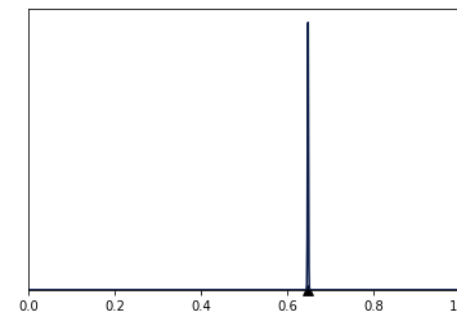$y_i = x_{true} + \varepsilon_i, \quad \varepsilon_i = 0$
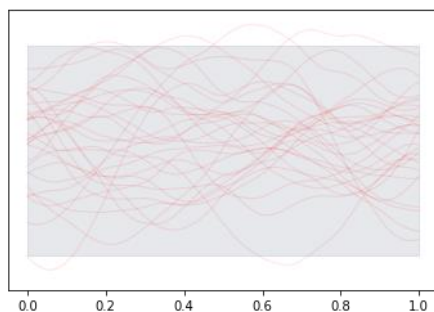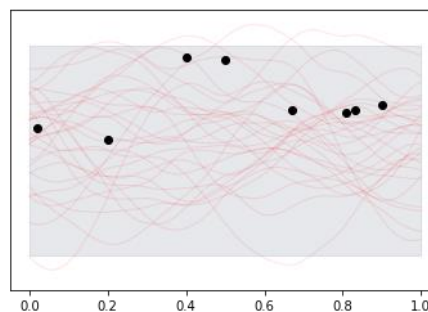
**Posterior**

Bayesian
inference

(find $x_{true}$)

Prior process
$GP(\mu(x|\theta), k(x, x'|\theta))$

**Prior**

**Observations**
$y_i = f_{true}(x_i) + \varepsilon_i, \ \varepsilon_i = 0$

**Posterior**

Bayesian
inference over
functions

(find $f_{true}$)

# GP as a prior over functions

**Prior distribution**
$$p(x|\boldsymbol{\theta})$$

The prior depends on a parameter $\theta$

**Prior process**
$$GP(\mu(x|\boldsymbol{\theta}), k(x, x'|\boldsymbol{\theta}))$$

**Prior**



**Observations**
$$y_i = x_{true} + \varepsilon_i, \quad \varepsilon_i = 0$$



**Posterior**



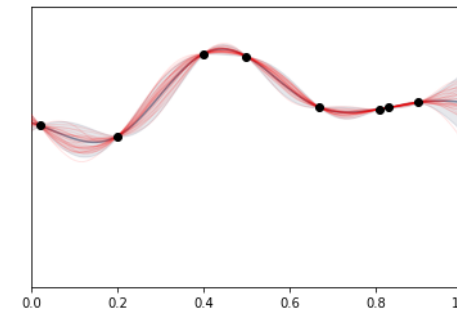Bayesian inference

(find $x_{true}$)

**Prior**



**Observations**
$$y_i = f_{true}(x_i) + \varepsilon_i, \ \varepsilon_i = 0$$



**Posterior**
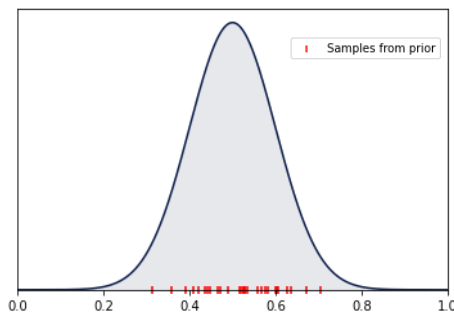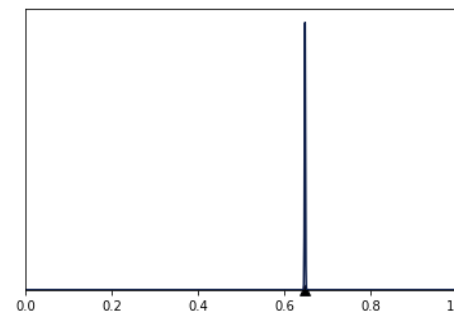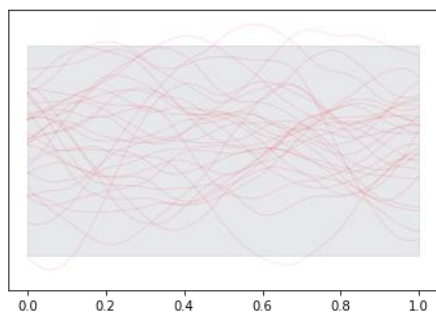


Bayesian inference over functions

(find $f_{true}$)

# Hyperparameter estimation

# The covariance function

**Recall:** $f \sim GP(\mu, k)$

- Assume $\mu = 0$
  (Or that we work with $f - \mu$. This is not a restrictive assumption)

- Assume $k$ is stationary: $k(x_1, x_2)$ can be written as $k(x_1 - x_2)$
  (Needed for theoretical analysis. This is often used in practice)

**How do we choose an appropriate covariance function $k$ ?**

- Let $\{k_\theta \mid \theta \in \Theta\}$ be a set of covariance functions parameterised by $\theta$

- Below are some examples with $\theta = (s, l)$ and $r = \|x_1 - x_2\| / l$

### Exponential

$$k_\theta(x_1, x_2) = s^2 e^{-r}$$

### Gaussian

$$k_\theta(x_1, x_2) = s^2 e^{-\frac{1}{2}r^2}$$

### Matérn 5/2

$$k_\theta(x_1, x_2) = s^2 \left(1 + \sqrt{5}r + \frac{5}{3}r^2\right) e^{-\sqrt{5}r}$$

# The covariance function

**Recall:** $f \sim GP(\mu, k)$

- Assume $\mu = 0$
  (Or that we work with $f - \mu$. This is not a restrictive assumption)
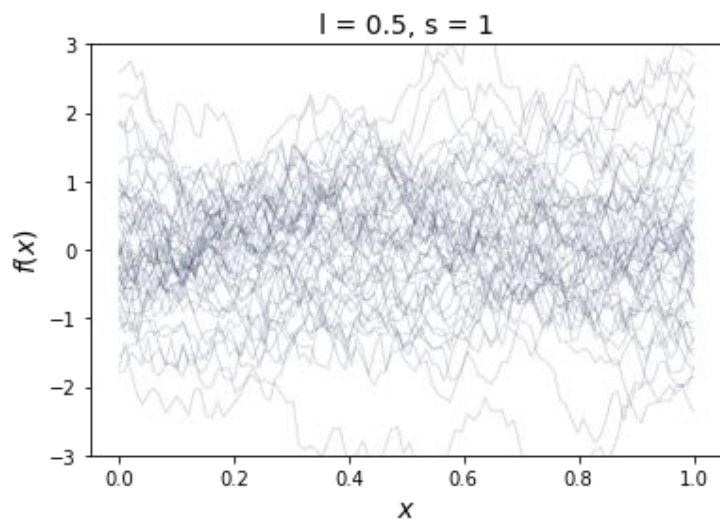
- Assume $k$ is stationary: $k(x_1, x_2)$ can be written as $k(x_1 - x_2)$
  (Needed for theoretical analysis. This is often used in practice)

**How do we choose an appropriate covariance function $k$ ?**

- Let $\{k_\theta \mid \theta \in \Theta\}$ be a set of covariance functions parameterised by $\theta$

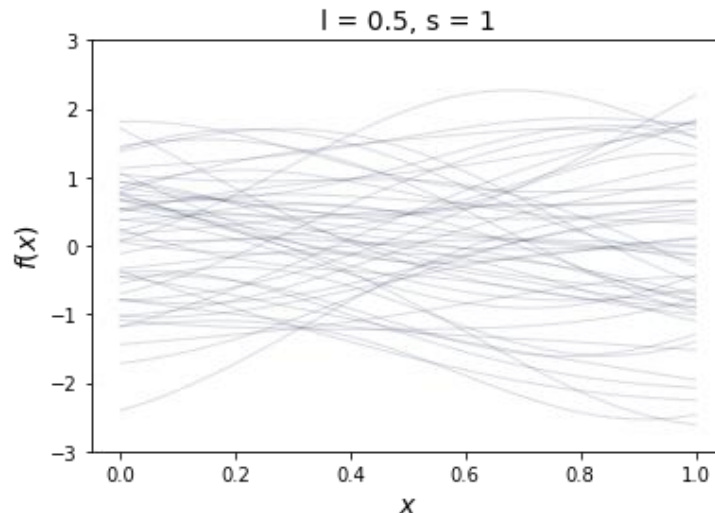- Below are some examples with $\theta = (s, l)$ and $r = \|x_1 - x_2\| / l$

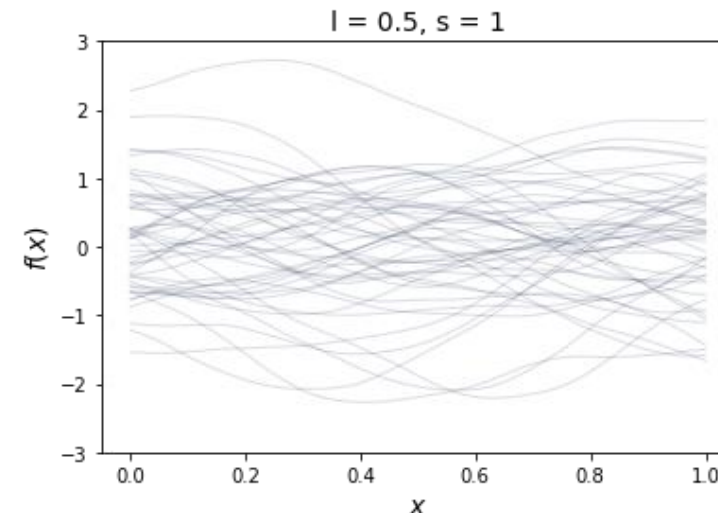**Exponential**
$$k_\theta(x_1, x_2) = s^2 e^{-r}$$

**Gaussian**
$$k_\theta(x_1, x_2) = s^2 e^{-\frac{1}{2}r^2}$$

**Matérn 5/2**
$$k_\theta(x_1, x_2) = s^2 \left(1 + \sqrt{5}r + \frac{5}{3}r^2\right) e^{-\sqrt{5}r}$$

# The covariance function
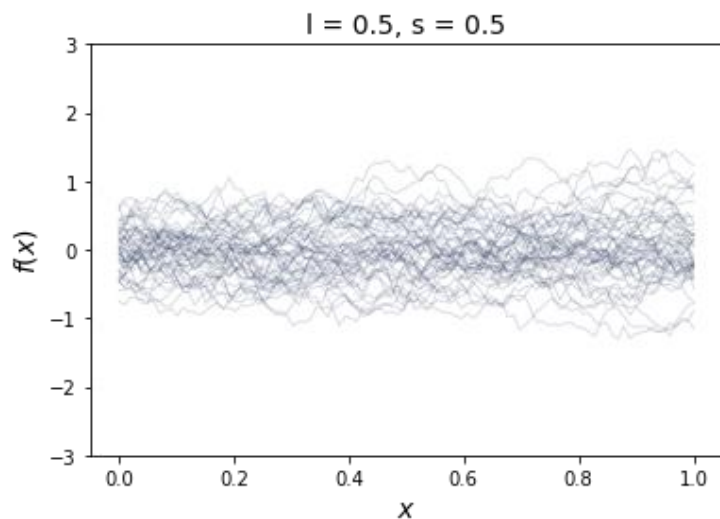
**Recall:** $f \sim GP(\mu, k)$

- Assume $\mu = 0$
  (Or that we work with $f - \mu$. This is not a restrictive assumption)

- Assume $k$ is stationary: $k(x_1, x_2)$ can be written as $k(x_1 - x_2)$
  (Needed for theoretical analysis. This is often used in practice)

**How do we choose an appropriate covariance function $k$ ?**

- Let $\{k_\theta \mid \theta \in \Theta\}$ be a set of covariance functions parameterised by $\theta$

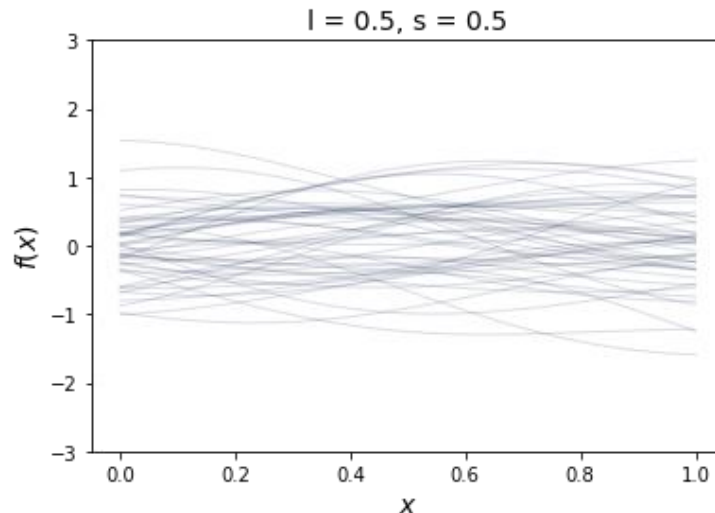- Below are some examples with $\theta = (s, l)$ and $r = \|x_1 - x_2\| / l$

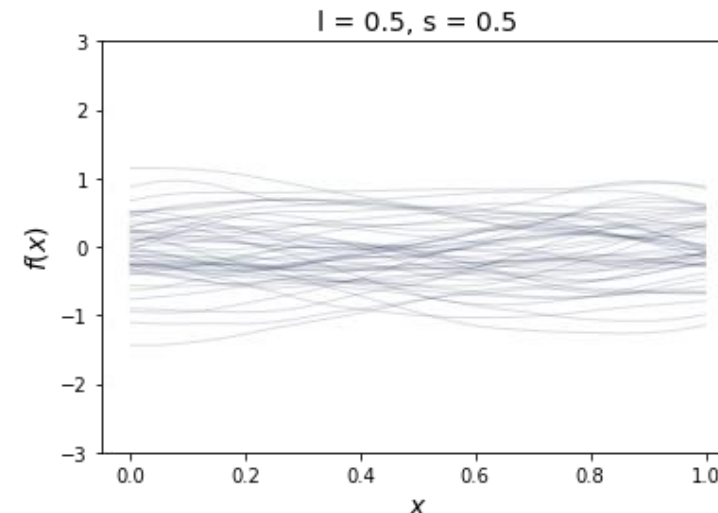### Exponential
$$k_\theta(x_1, x_2) = s^2 e^{-r}$$

### Gaussian
$$k_\theta(x_1, x_2) = s^2 e^{-\frac{1}{2}r^2}$$

### Matérn 5/2
$$k_\theta(x_1, x_2) = s^2 \left(1 + \sqrt{5}r + \frac{5}{3}r^2\right) e^{-\sqrt{5}r}$$

# Tree ways of estimating $\theta$

1. **Maximum likelihood (ML)**
   - Most common alternative

2. **Cross validation (CV)**
   - Leave-one-out predictions can be made efficient

3. **Bayesian**
   - MAP estimates
   - Full Bayesian treatment with MCMC to sample from $p(\theta|D)$[1]
   - Some use within Uncertainty Quantification[2]

**The plug-in approach**

(Also called Type-II maximum likelihood)

- Compute a fixed estimate $\hat{\theta}$

- Treat $\hat{\theta}$ as the "true" value and compute the posterior GP for $k_{\hat{\theta}}$

**Most common to use one of these and the plug-in approach**

# Maximum likelihood (ML)

**We have**

- GP prior: $f \sim GP(0, k)$

- Data $\{(x_i, y_i)\}_{i=1}^n$ where: $y_i = f(x_i)$

- This means that

$$Y \sim N(0, \mathrm{K})$$

with $\mathrm{K}_{i,j} = k(x_i, x_j)$, $Y_i = y_i$

**The covariance matrix that depends on $\theta$**

- $\mathrm{K} = \mathrm{K}_\theta$ depends on some parameter $\theta$

**Recall the Gaussian density**

$$p(Y|X, \theta) = \frac{1}{(2\pi)^{n/2} \sqrt{|\mathrm{K}_\theta|}} e^{-\frac{1}{2} Y^T \mathrm{K}_\theta^{-1} Y}$$
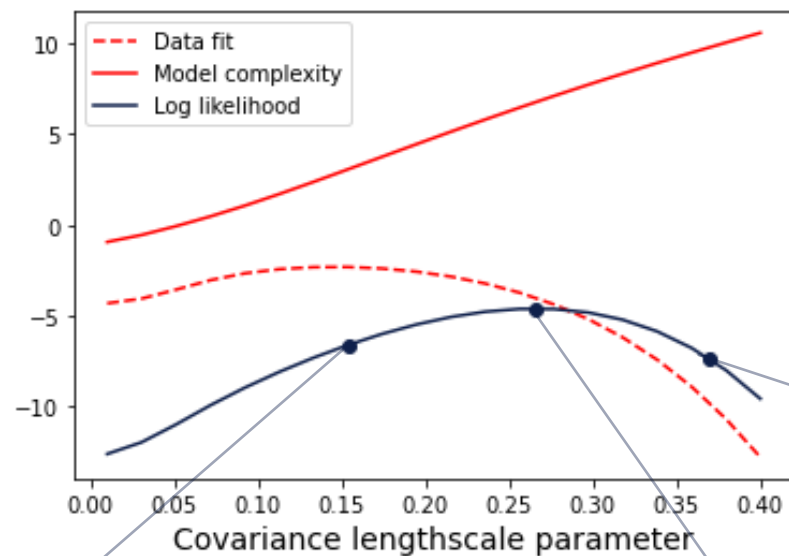
**The log likelihood:**

$$L(\theta) = -\frac{1}{2} Y^T \mathrm{K}_\theta^{-1} Y \; \underbrace{\phantom{xx}} \; - \; \frac{1}{2} \log|\mathrm{K}_\theta| \; \underbrace{\phantom{xx}} \; - \; \frac{n}{2} \log 2\pi$$

Data fit  ·  Model complexity

**Remark**

If $y_i = f(x_i) + \varepsilon_i$, $\varepsilon_i \sim N(0, \sigma^2)$, we would make us of the covariance matrix $\Sigma_\theta = (K_\theta + \sigma^2 I)$

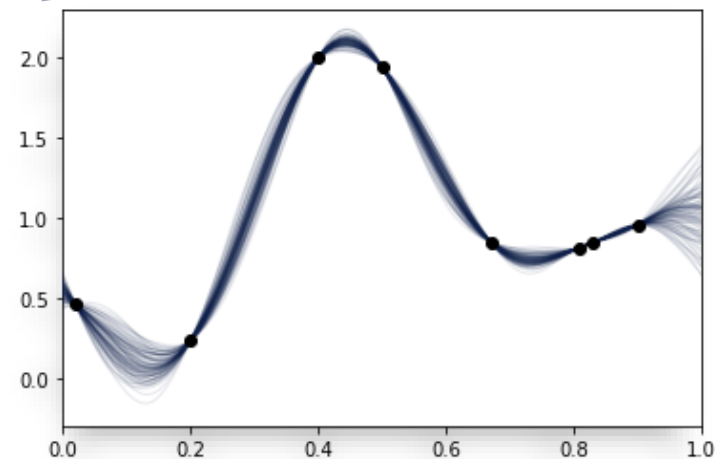- The noise variance $\sigma^2$ could also be estimated together with $\theta$

- We call $L(\theta)$ the *log marginal likelihood*

(3) C. E. Rasmussen and C. K. I. Williams. Gaussian Processes for Machine Learning. MIT Press, 2006.

# Maximum likelihood (ML)



The log likelihood:

$$L(\theta) = -\frac{1}{2}Y^T K_\theta^{-1} Y \; - \; \frac{1}{2}\log|K_\theta| \; - \; \frac{n}{2}\log 2\pi$$

Data fit                     Model complexity

# Putting it all together
- Gaussian process regression

# Gaussian process regression

**Prior**

1) Select model

2) Select likelihood

3) Optimize hyperparameters

**Posterior**

4) Condition on data

5) Make predictions

DNV

# Gaussian process regression

**Prior**

1) Select model
2) Select likelihood
3) Optimize hyperparameters

**Posterior**

4) Condition on data
5) Make predictions

Select mean $\mu(\boldsymbol{x})$ and covariance function $k(\boldsymbol{x}, \boldsymbol{x}')$
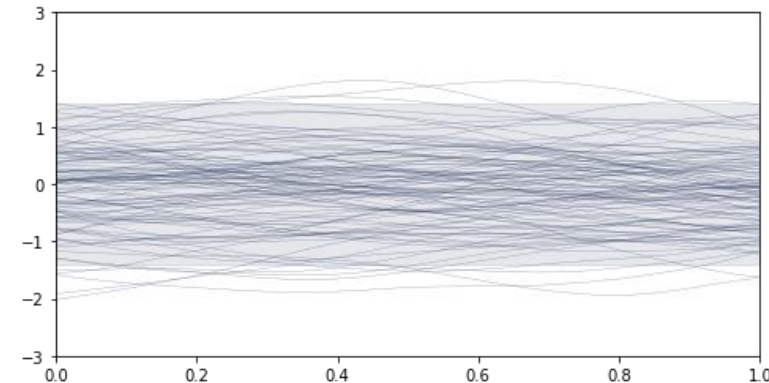
*Tip: Start with $\mu = 0$ and $k = Matérn$*

DNV

# Gaussian process regression

**Prior**

1) Select model

2) Select likelihood

3) Optimize hyperparameters

**Posterior**

4) Condition on data

5) Make predictions

Assume additive Gaussian noise: $y = f(\boldsymbol{x}) + \varepsilon$, $\varepsilon \sim N(0, \sigma^2)$.

Decide if $\sigma^2$ is fixed or unknown.

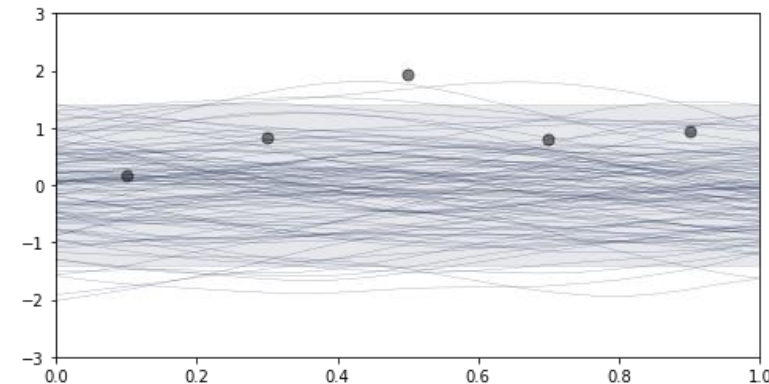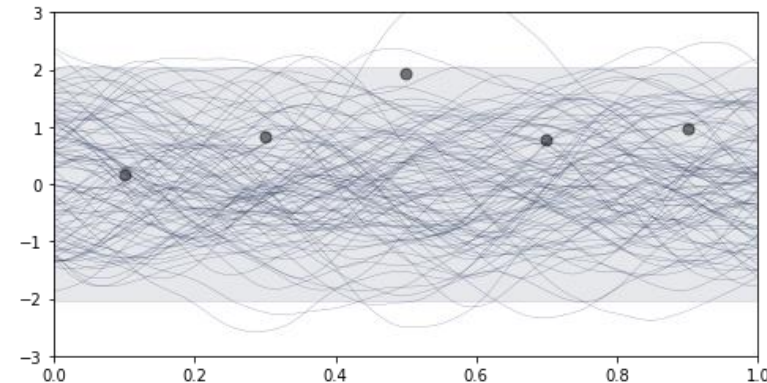*Tip: Set $\sigma^2 = 10^{-6} \approx 0$ for noiseless data.*

# Gaussian process regression

**Prior**

1) Select model

2) Select likelihood

3) Optimize hyperparameters

**Posterior**

4) Condition on data

5) Make predictions

Identify which parameters of the mean function, covariance function, and likelihood to optimize:

$$\mu(\boldsymbol{x}|\beta), \ k(\boldsymbol{x}, \boldsymbol{x}'|\theta), \ \varepsilon \sim N(0, \sigma^2).$$

Optimize using e.g. maximum likelihood

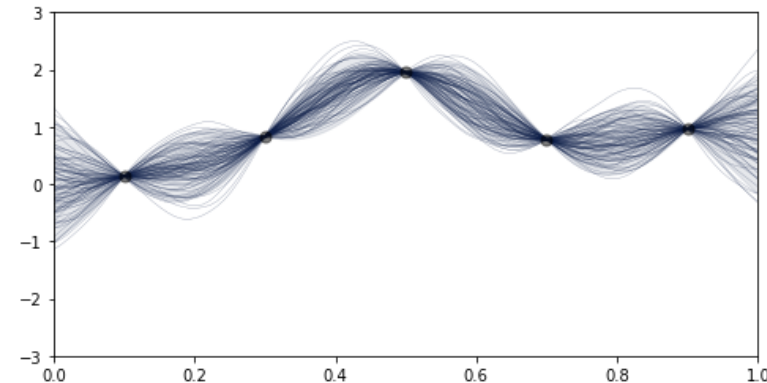$$(\beta, \theta, \sigma) \in argmax \ L(\beta, \theta, \sigma)$$

DNV

# Gaussian process regression

**Prior**

1) Select model

2) Select likelihood

3) Optimize hyperparameters

**Posterior**

4) Condition on data

5) Make predictions

The posterior GP can be computed analytically

$$\mu_{f^*|D} = \mu^* + K^*(K + \sigma^2 I)^{-1}(Y - \mu)$$

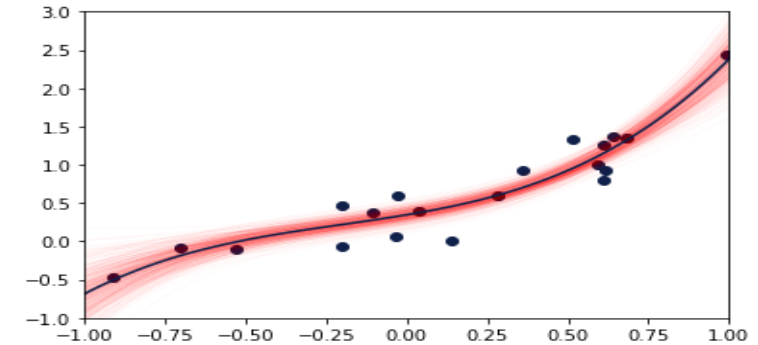$$\Sigma_{f^*|D} = K^{**} - K^*(K + \sigma^2 I)^{-1}(K^*)^T$$

DNV

# Gaussian process regression

**Prior**
1) Select model
2) Select likelihood
3) Optimize hyperparameters

**Posterior**
4) Condition on data
5) Make predictions $\longrightarrow$
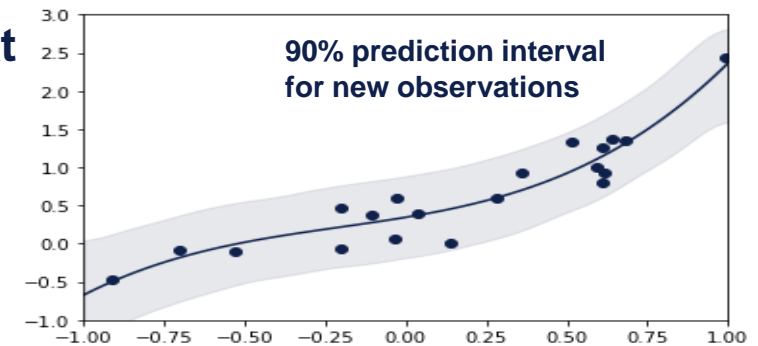
**Predict the latent function**

$$f(x)$$



**Predict the next observation**

$$f(x) + \varepsilon$$

90% prediction interval for new observations

DNV

# Gaussian process regression

**Generalisations:**

We have focused on functions $f \colon \mathbb{R} \to \mathbb{R}$ with i.i.d. noise.

- Extending to $f \colon \mathbb{R}^N \to \mathbb{R}^M$ is trivial

- Extending to general Gaussian noise is trivial

**Limitations:**

- Non-Gaussian noise

- Large datasets (due to cubic time complexity)

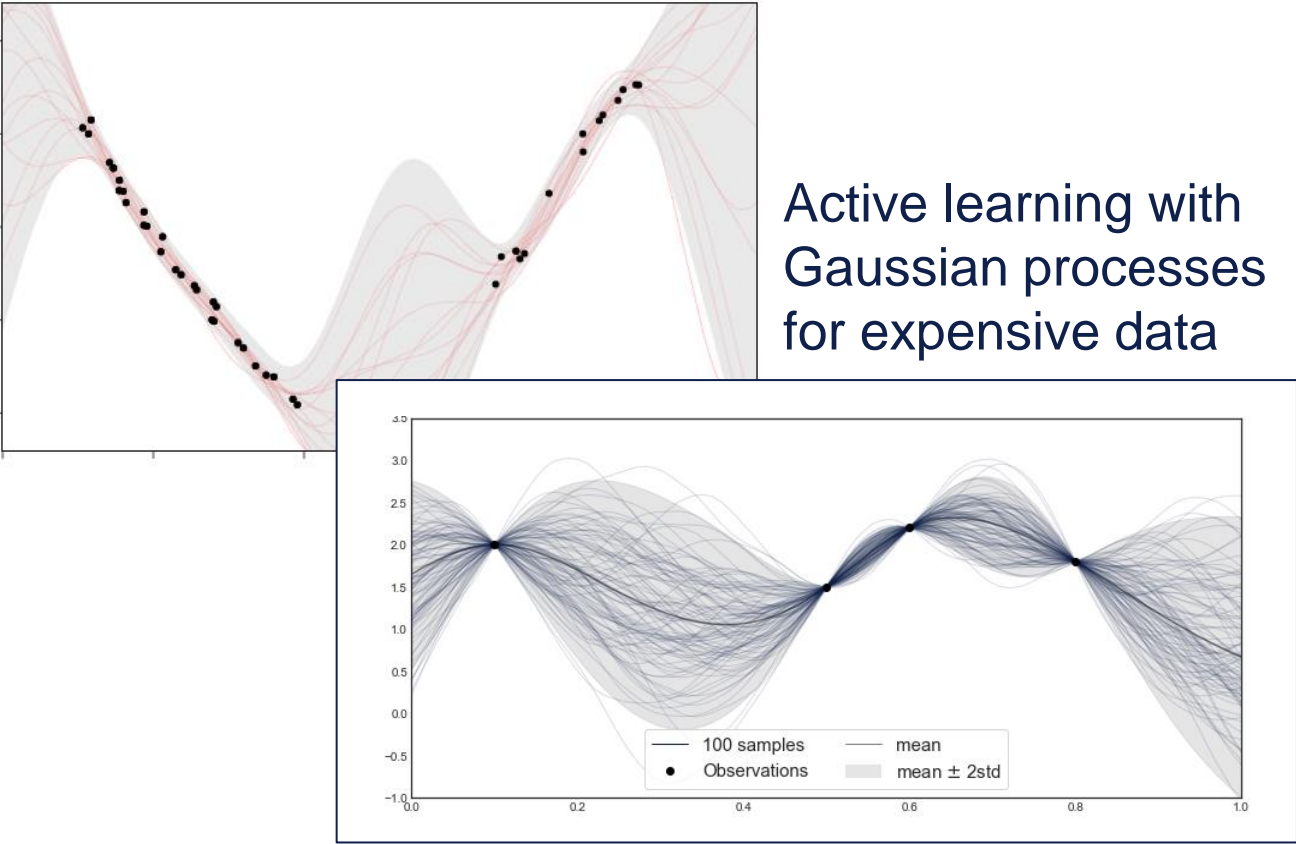For this, special methods are needed.

**Software**

There are many ML and UQ software packages for GPs
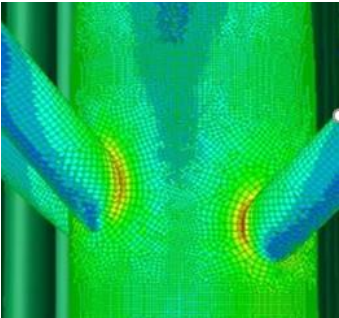
Some Python alternatives:

- scikit-learn

- GPy

- GPyTorch
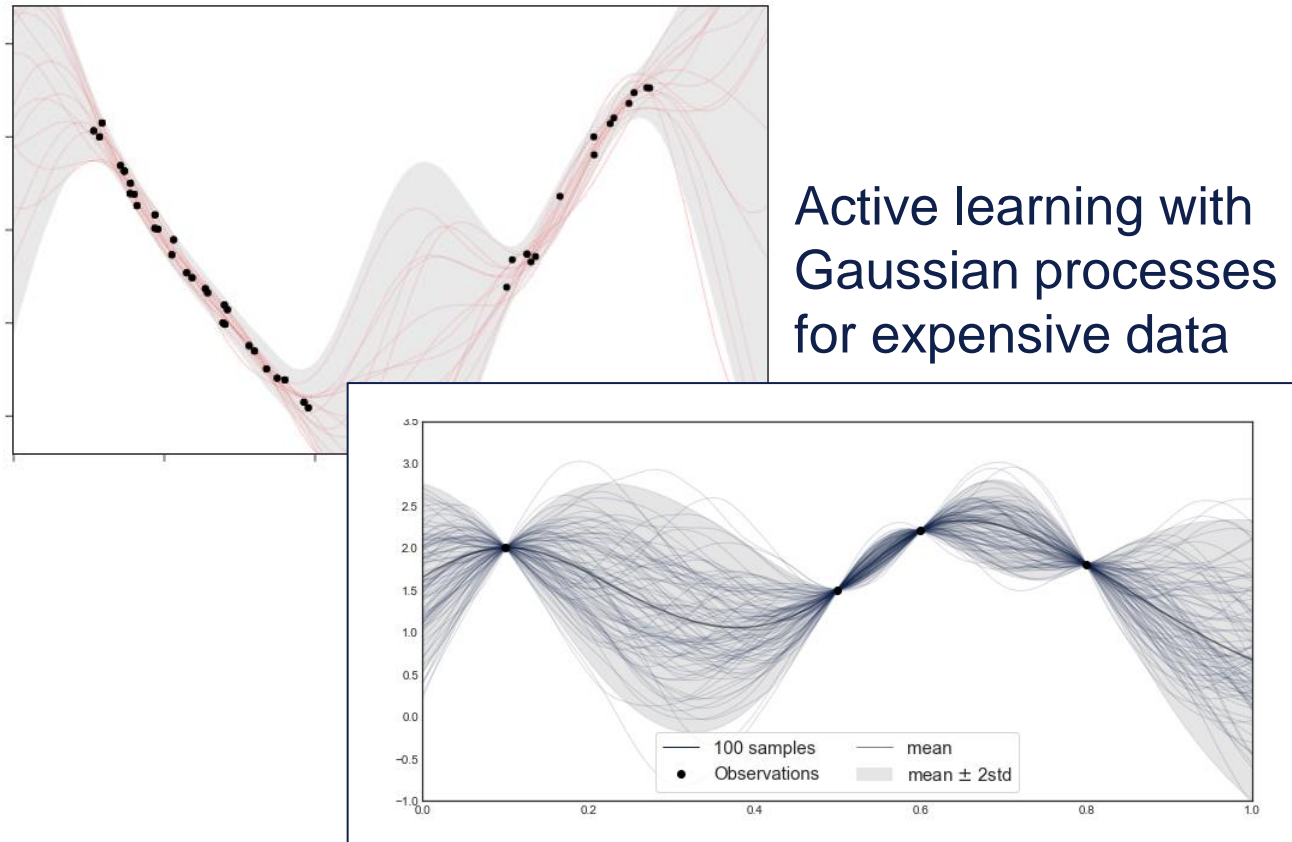
- GPflow

DNV

# DOE

# Design of experiments



Active learning with Gaussian processes for expensive data



**Computer experiment**



**Lab experiment**



**Inspection**

DNV

# Design of experiments



Active learning with Gaussian processes for expensive data

To find the maximum of the function that has generated the data, we can use e.g.

**Upper confidence band**

$$x \in \text{argmax}(E[f(x)] + \lambda \cdot Std[f(x)])$$

**Expected improvement**

$$x \in \text{argmax}\, E[\max(f(x) - f(x^+), 0]$$

Where $f(x)$ is the GP (or any other object with epistemic uncertainty that depends on $x$)

# Bayesian optimization
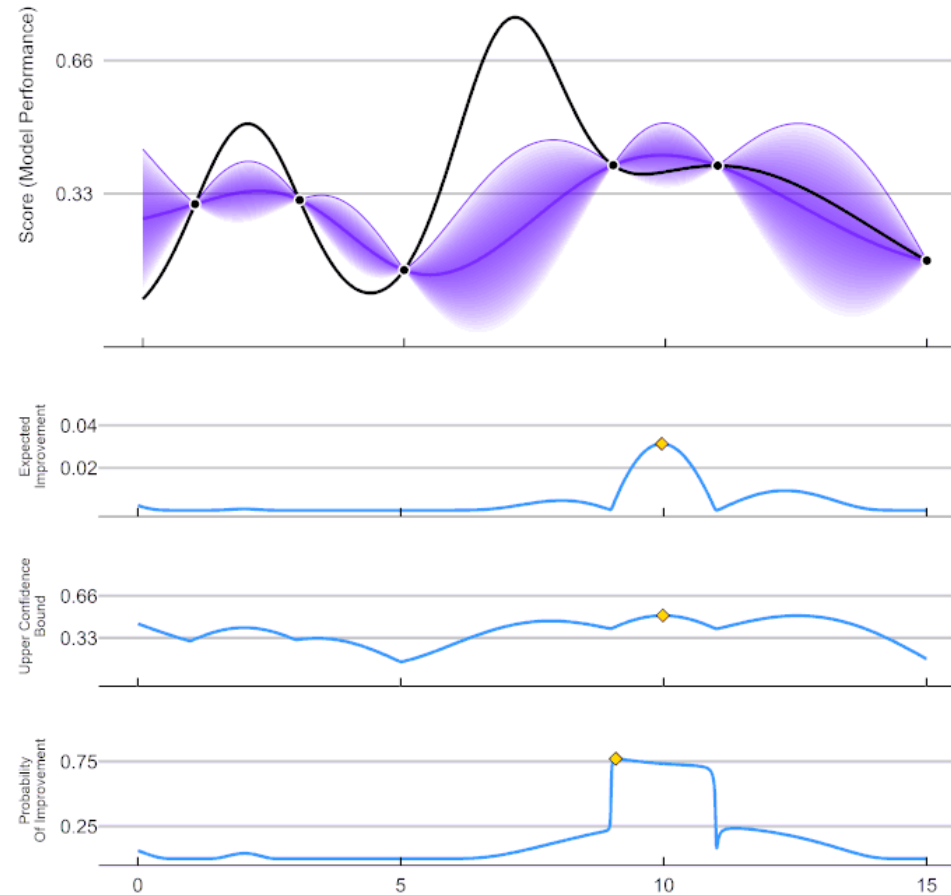


**Expected improvement**

$$E[\max(f(x) - f(x^+), 0]$$

**Upper confidence band**

$$E[f(x)] + \lambda \cdot Std[f(x)]$$

**Probability of improvement**

$$P(f(x) > f(x^+))$$

DNV
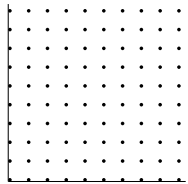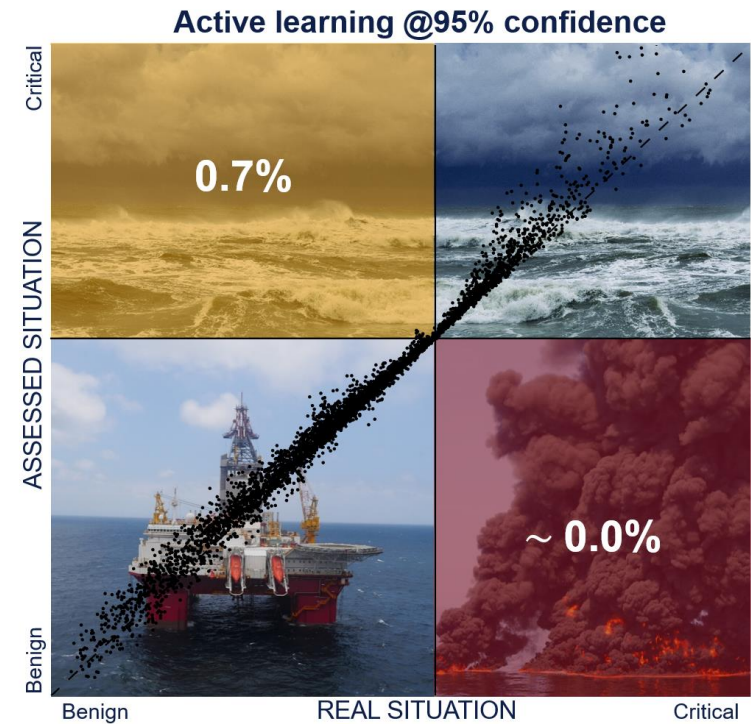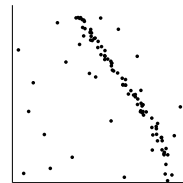
# Example - Risk-based design of experiments

Eldevik, S. and Sætre, S. (2020) *Offshore Workover Operations: Reducing Uncertainty of Critical Weather Scenarios by Optimal Use of Simulations and Probabilistic Machine Learning.* ESREL 2020.

Standard approach

Optimized for safety-critical decisions



Structured exploration (Grid)

1.3%

3.6%

ASSESSED SITUATION

Critical / Benign

REAL SITUATION

Benign / Critical



Active learning @95% confidence

0.7%

~ 0.0%

ASSESSED SITUATION

Critical / Benign

REAL SITUATION

Benign / Critical

DNV

**www.dnv.com**

DNV