

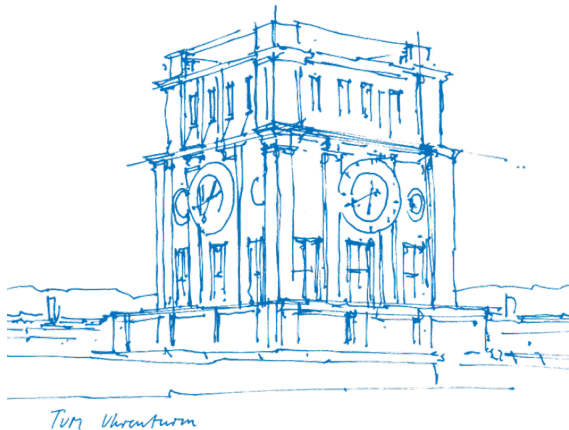
# Seminar: Gaussian Processes for Machine Learning

Gaussian processes for regression

**Davit Papikyan**

Mathematics in Data Science  
Department of Mathematics  
Technical University of Munich

February 3<sup>rd</sup>, 2022

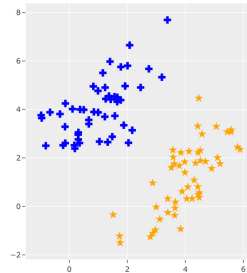
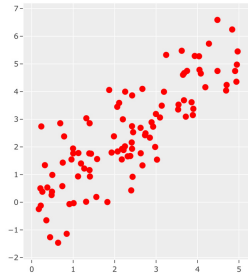


- 1 Introduction
- 2 Bayesian Linear Regression
- 3 Gaussian Processes

# Introduction

## Supervised Learning: Regression vs. Classification

- Regression: **continuous** variables (price, salary, etc.)
- Classification: **discrete** variables (spam/not spam, male/female, etc.)



1 Introduction

**2 Bayesian Linear Regression**

3 Gaussian Processes

# Bayesian Linear Regression

## Problem statement

Training set:  $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, n\}$ , where

$$\mathbf{x}_i \in \mathbb{R}^d \quad i = 1, 2, \dots, n$$

$$y_i \in \mathbb{R}$$

Using matrix notations:  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$ , where

$$\mathbf{X} \in \mathbb{R}^{d \times n} - \text{design matrix}$$

$$\mathbf{y} \in \mathbb{R}^n - \text{target vector}$$

# Bayesian Linear Regression

## Problem statement

Training set:  $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, n\}$ , where

$$\mathbf{x}_i \in \mathbb{R}^d \quad i = 1, 2, \dots, n$$

$$y_i \in \mathbb{R}$$

Using matrix notations:  $\mathcal{D} = (\mathbf{X}, \mathbf{y})$ , where

$$\mathbf{X} \in \mathbb{R}^{d \times n} - \text{design matrix}$$

$$\mathbf{y} \in \mathbb{R}^n - \text{target vector}$$

**Goal:** Make inferences about the relationship between  $\mathbf{X}$  and  $\mathbf{y}$ , i.e. the conditional distribution of the targets given inputs (without modeling the input distribution itself).

# Bayesian Linear Regression

## Problem statement

Model:  $f(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w}$ ,  $y_i = f(\mathbf{x}_i) + \epsilon$  where  $\epsilon \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$ .  
 $\mathbf{w} \in \mathbb{R}^d$  – weight vector

# Bayesian Linear Regression

## Likelihood

Model:  $f(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{w}$ ,  $y_i = f(\mathbf{x}_i) + \epsilon$  where  $\epsilon \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$   
 $\mathbf{w} \in \mathbb{R}^d$  – weight vector

Likelihood: Probability density of the observations given the parameters.

$$\begin{aligned} p(\mathbf{y} | \mathbf{X}, \mathbf{w}) &= \prod_{i=1}^n p(y_i | \mathbf{x}_i, \mathbf{w}) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{(y_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2} \right] \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left[ -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}^T \mathbf{w}\|^2 \right] = \mathcal{N}(\mathbf{X}^T \mathbf{w}, \sigma^2 \mathbf{I}). \end{aligned}$$



# Bayesian Linear Regression

## Prior

Prior: Probability density of parameters expressing our beliefs about them before having a look at the data

$$\boldsymbol{w} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_w) .$$

# Bayesian Linear Regression

## Posterior

Posterior: Probability density of parameters given the observations (data).

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}} \qquad p(\mathbf{w} \mid \mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) \times p(\mathbf{w})}{p(\mathbf{y} \mid \mathbf{X})}$$

where marginal likelihood is independent of  $\mathbf{w}$ , plays a role of normalizing constant and is given by

$$p(\mathbf{y} \mid \mathbf{X}) = \int p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w} .$$

# Bayesian Linear Regression

## Posterior

$$\begin{aligned} p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) &\propto \exp \left[ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}^T \mathbf{w})^T (\mathbf{y} - \mathbf{X}^T \mathbf{w}) \right] \exp \left[ \frac{1}{2} \mathbf{w}^T \boldsymbol{\Sigma}_w \mathbf{w} \right] \\ &\propto \exp \left[ -\frac{1}{2} (\mathbf{w} - \tilde{\mathbf{w}})^T \underbrace{\left( \frac{1}{\sigma^2} \mathbf{X} \mathbf{X}^T + \boldsymbol{\Sigma}_w^{-1} \right)}_{\mathbf{A}} (\mathbf{w} - \tilde{\mathbf{w}}) \right] \end{aligned}$$

$$\text{where } \tilde{\mathbf{w}} = \frac{1}{\sigma^2} \mathbf{A} \mathbf{X} \mathbf{y} \quad \text{and} \quad \mathbf{A} = \frac{1}{\sigma^2} \mathbf{X} \mathbf{X}^T + \boldsymbol{\Sigma}_w^{-1}.$$

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) \sim \mathcal{N}(\tilde{\mathbf{w}}, \mathbf{A}^{-1}).$$

# Bayesian Linear Regression

## Posterior predictive distribution

Assume we have a new data point  $\mathbf{x}_*$  and goal is to predict  $y_*$ .

$$p(y_* \mid \mathbf{x}_*, \mathbf{X}, \mathbf{y}) =$$

# Bayesian Linear Regression

## Posterior predictive distribution

Assume we have a new data point  $\mathbf{x}_*$  and goal is to predict  $y_*$ .

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_*, \mathbf{w} | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) d\mathbf{w}$$

# Bayesian Linear Regression

## Posterior predictive distribution

Assume we have a new data point  $\mathbf{x}_*$  and goal is to predict  $y_*$ .

$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_*, \mathbf{w} | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) d\mathbf{w} \\ &= \int p(y_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w} \\ &= \mathcal{N}\left(\frac{1}{\sigma^2} \mathbf{x}_*^T \mathbf{A}^{-1} \mathbf{X} \mathbf{y}, \mathbf{x}_*^T \mathbf{A}^{-1} \mathbf{x}_*\right). \end{aligned}$$

$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y})$  is known as posterior predictive distribution or just predictive distribution.

Note that we can now estimate the uncertainty of our prediction.

# Bayesian Linear Regression

## Addressing nonlinear dependency

What if the dependency between  $y$  and  $x$  is not linear?

# Bayesian Linear Regression

## Addressing nonlinear dependency

What if the dependency between  $y$  and  $x$  is not linear?

**Idea:** Project data into some high dimensional space and then apply the linear model in that space.



# Bayesian Linear Regression

## Addressing nonlinear dependency

$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^k$  with  $\Phi := \Phi(\mathbf{X}) \in \mathbb{R}^{k \times n}$  be the aggregation of columns  $\phi(x)$  for all samples in training set.

Now the model becomes  $f(x) = \phi(x)^T \mathbf{w}$  where  $\mathbf{w} \in \mathbb{R}^k$ .

# Bayesian Linear Regression

## Addressing nonlinear dependency

$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^k$  with  $\Phi := \Phi(\mathbf{X}) \in \mathbb{R}^{k \times n}$  be the aggregation of columns  $\phi(x)$  for all samples in training set.

Now the model becomes  $f(x) = \phi(x)^T \mathbf{w}$  where  $\mathbf{w} \in \mathbb{R}^k$ .

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}\left(\frac{1}{\sigma^2} \phi(\mathbf{x}_*)^T \mathbf{A}^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^T \mathbf{A}^{-1} \phi(\mathbf{x}_*)\right)$$

where  $\mathbf{A} = \frac{1}{\sigma^2} \Phi \Phi^T + \Sigma_w^{-1}$ .

# Bayesian Linear Regression

## Addressing nonlinear dependency

$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^k$  with  $\Phi := \Phi(\mathbf{X}) \in \mathbb{R}^{k \times n}$  be the aggregation of columns  $\phi(x)$  for all samples in training set.

Now the model becomes  $f(x) = \phi(x)^T \mathbf{w}$  where  $\mathbf{w} \in \mathbb{R}^k$ .

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}\left(\frac{1}{\sigma^2} \phi(\mathbf{x}_*)^T \mathbf{A}^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^T \mathbf{A}^{-1} \phi(\mathbf{x}_*)\right)$$

where  $\mathbf{A} = \frac{1}{\sigma^2} \Phi \Phi^T + \Sigma_w^{-1}$ .

Problem:  $\mathbf{A}^{-1}$  is not tractable for large  $k$  !

# Bayesian Linear Regression

## Addressing nonlinear dependency

Denote  $K := \Phi^T \Sigma_w \Phi$  and  $\phi_* := \phi(x_*)$ . Let's rewrite the mean:

# Bayesian Linear Regression

## Addressing nonlinear dependency

Denote  $K := \Phi^T \Sigma_w \Phi$  and  $\phi_* := \phi(\mathbf{x}_*)$ . Let's rewrite the mean:

$$A \Sigma_w \Phi = \frac{1}{\sigma^2} \Phi \Phi^T \Sigma_w \Phi + \Phi = \frac{1}{\sigma^2} \Phi (\Phi^T \Sigma_w \Phi + \sigma^2 I) = \frac{1}{\sigma^2} \Phi (K + \sigma^2 I)$$

# Bayesian Linear Regression

## Addressing nonlinear dependency

Denote  $\mathbf{K} := \Phi^T \Sigma_w \Phi$  and  $\phi_* := \phi(\mathbf{x}_*)$ . Let's rewrite the mean:

$$\mathbf{A} \Sigma_w \Phi = \frac{1}{\sigma^2} \Phi \Phi^T \Sigma_w \Phi + \Phi = \frac{1}{\sigma^2} \Phi (\Phi^T \Sigma_w \Phi + \sigma^2 \mathbf{I}) = \frac{1}{\sigma^2} \Phi (\mathbf{K} + \sigma^2 \mathbf{I}).$$

Multiplying both sides by  $\mathbf{A}^{-1}$  from left and  $(\mathbf{K} + \sigma^2 \mathbf{I})^{-1}$  from right we will get

$$\Sigma_w \Phi (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} = \frac{1}{\sigma^2} \mathbf{A}^{-1} \Phi,$$

by using which we can rewrite the mean of posterior predictive distribution as

$$\frac{1}{\sigma^2} \phi_*^T \mathbf{A}^{-1} \Phi \mathbf{y} = \phi_*^T \Sigma_w \Phi (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}.$$

# Bayesian Linear Regression

## Addressing nonlinear dependency

Now let's rewrite the variance.

**Lemma: Matrix inversion lemma or Woodbury matrix identity**

$$(B + UCV)^{-1} = B^{-1} - B^{-1}U(C^{-1} + VB^{-1}U)^{-1}VB^{-1}.$$

$$A^{-1} = \left( \frac{1}{\sigma^2} \Phi \Phi^T + \Sigma_w^{-1} \right)^{-1} = \Sigma_w - \Sigma_w \Phi \underbrace{(\sigma^2 I + \Phi^T \Sigma_w \Phi)^{-1}}_K \Phi^T \Sigma_w$$

# Bayesian Linear Regression

## Addressing nonlinear dependency

Now let's rewrite the variance.

**Lemma: Matrix inversion lemma or Woodbury matrix identity**

$$(B + UCV)^{-1} = B^{-1} - B^{-1}U(C^{-1} + VB^{-1}U)^{-1}VB^{-1}$$

$$A^{-1} = \left( \frac{1}{\sigma^2} \Phi \Phi^T + \Sigma_w^{-1} \right)^{-1} = \Sigma_w - \Sigma_w \Phi (\sigma^2 I + \underbrace{\Phi^T \Sigma_w \Phi}_K)^{-1} \Phi^T \Sigma_w ,$$

$$\phi_*^T A^{-1} \phi_* = \phi_*^T \Sigma_w \phi_* - \phi_*^T \Sigma_w \Phi (K + \sigma^2 I)^{-1} \Phi^T \Sigma_w \phi_* .$$



# Bayesian Linear Regression

## Addressing nonlinear dependency

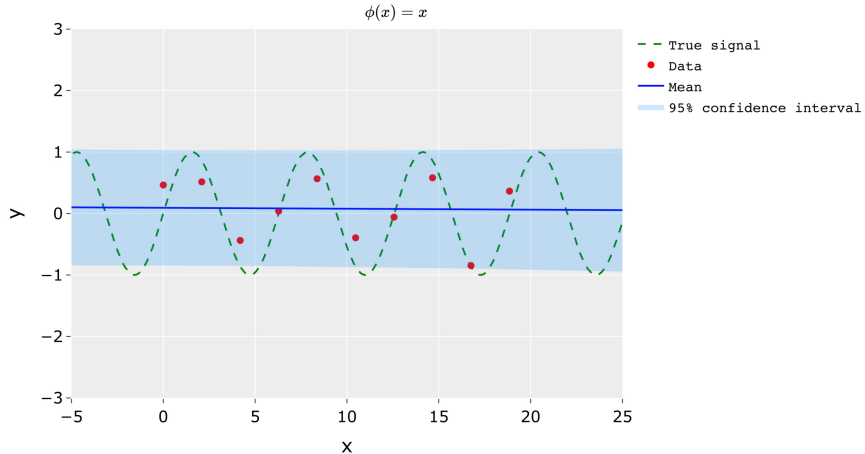
Combining the results we get:

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N} \left( \phi_*^T \Sigma_w \Phi (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \right. \\ \left. \phi_*^T \Sigma_w \phi_*^T - \phi_*^T \Sigma_w \Phi (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \Phi^T \Sigma_w \phi_* \right)$$

where we need to invert matrix of size  $n \times n$  which is more convenient when  $n < k$ .

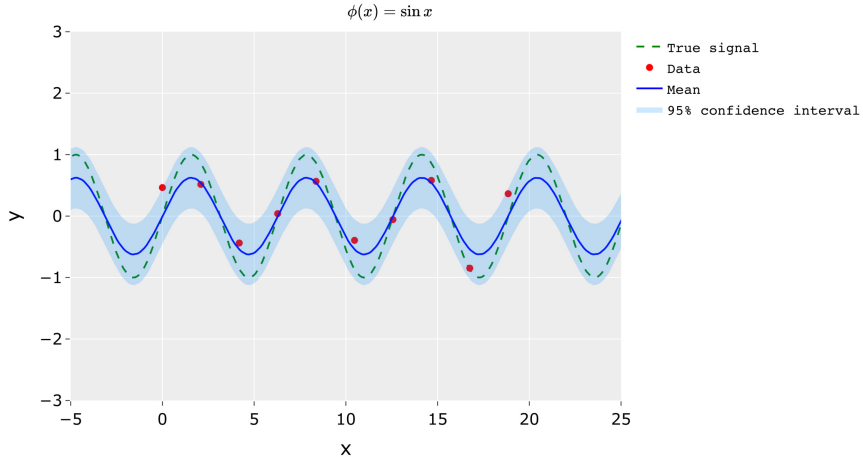
# Bayesian Linear Regression

## Addressing nonlinear dependency



# Bayesian Linear Regression

## Addressing nonlinear dependency



# Bayesian Linear Regression

## Kernel trick

Notice that feature space always enters in the forms of  $\Phi^T \Sigma_w \Phi$ ,  $\phi_*^T \Sigma_w \Phi$ ,  $\phi_*^T \Sigma_w \phi_*$ .  
Entries of these matrices are of the form  $\phi(x)^T \Sigma_w \phi(x') =: k(x, x')$ .

- $k(x, x')$  is known as *kernel* or *covariance function*

# Bayesian Linear Regression

## Kernel trick

Notice that feature space always enters in the form of  $\Phi^T \Sigma_w \Phi$ ,  $\phi_*^T \Sigma_w \Phi$ ,  $\phi_*^T \Sigma_w \phi_*$ .  
Entries of these matrices are of the form  $\phi(x)^T \Sigma_w \phi(x') =: k(x, x')$ .

- $k(x, x')$  is known as *kernel* or *covariance function*

Define  $\psi(x) := \Sigma_w^{1/2} \phi(x)$  we can rewrite kernel as a simple dot product

$$k(x, x') = \psi(x)^T \psi(x'),$$

where  $\Sigma_w^{1/2} = U D^{1/2} U^T$  for  $U D U^T$  being the SVD of  $\Sigma_w$ .

# Bayesian Linear Regression

## Kernel trick

Notice that feature space always enters in the form of  $\Phi^T \Sigma_w \Phi$ ,  $\phi_*^T \Sigma_w \Phi$ ,  $\phi_*^T \Sigma_w \phi_*$ .  
Entries of these matrices are of the form  $\phi(x)^T \Sigma_w \phi(x') =: k(x, x')$ .

- $k(x, x')$  is known as *kernel* or *covariance function*

Define  $\psi(x) := \Sigma_w^{1/2} \phi(x)$  we can rewrite kernel as a simple dot product

$$k(x, x') = \psi(x)^T \psi(x'),$$

where  $\Sigma_w^{1/2} = U D^{1/2} U^T$  for  $U D U^T$  being the SVD of  $\Sigma_w$ .

Kernels provide an efficient way to transform data without having to compute coordinates in a new space, i.e.  $\phi(x)$ .

- 1 Introduction
- 2 Bayesian Linear Regression
- 3 Gaussian Processes**

# Gaussian Processes

## Problem statement

**Goal:** Find hypothesis function that fits the data.



# Gaussian Process Regression

## Problem statement

**Goal:** Find hypothesis function that fits the data.

**Assumption:**  $y_i = f(\mathbf{x}_i) + \epsilon$

- $\epsilon \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$
- $f$  is a function-valued r.v. drawn from a Gaussian Process conditioned on data

# Gaussian Processes

## Gaussian Process

### Definition: Gaussian Process

Let  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be a positive definite kernel and  $m : \mathbb{R}^d \rightarrow \mathbb{R}$  be a finite real-valued function. Then a random function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is said to be a Gaussian Process (GP) with mean  $m$  and covariance kernel  $k$ , denoted by  $GP(m, k)$ , if the following holds: for any finite set  $X = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \subset \mathbb{R}^d$  of any size  $n \in \mathbb{N}$ , the random vector

$$f_X = (f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n))^T \in \mathbb{R}^n$$

follows a multivariate normal distribution  $\mathcal{N}(m_X, k_{XX})$  with covariance

$k_{XX} = [k(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^n$  and mean vector  $m_X = [m(\mathbf{x}_1), m(\mathbf{x}_2), \dots, m(\mathbf{x}_n)]^T$ . It is written as

$$f \sim GP(m, k).$$

# Gaussian Processes

## Gaussian Process

Since  $m$  and  $k$  are respectively the mean and covariance functions of a Gaussian process, they can be written as

$$m(\mathbf{x}) = \mathbb{E}_{f \sim GP(m,k)}[f(\mathbf{x})],$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{f \sim GP(m,k)}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))].$$

The function values  $f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)$  are jointly Gaussian for any finite  $n$ .

# Gaussian Processes

## Example: Prior

Define

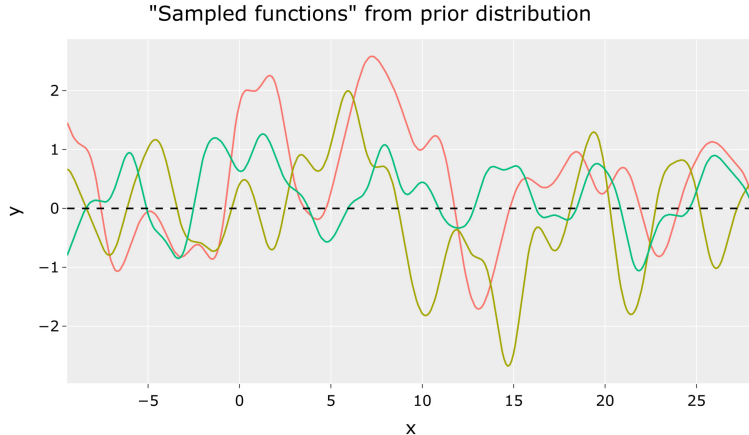
$$k(\mathbf{x}, \mathbf{x}') = \text{cov}(f(\mathbf{x}), f(\mathbf{x}')) = \exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|_2^2\right).$$

With this covariance function we can define covariance matrix  $K(\mathbf{X}_*, \mathbf{X}_*)$  between two sets of points by applying equation above elementwise for each pair of points.

Let's draw some function values  $f \sim \mathcal{N}(\mathbf{0}, K(\mathbf{X}_*, \mathbf{X}_*))$ .

# Gaussian Processes

## Example: Prior



# Gaussian Processes

## Computing the predictions

Assumption:  $y = f(\mathbf{x}) + \epsilon$  where noise is i.i.d.  $\mathcal{N}(0, \sigma^2)$ .

Then,

$$\begin{aligned}\text{cov}(y_i, y_j) &= \text{cov}(f(\mathbf{x}_i) + \epsilon_i, f(\mathbf{x}_j) + \epsilon_j) = \\ &\quad \text{cov}(f(\mathbf{x}_i), f(\mathbf{x}_j)) + \text{cov}(f(\mathbf{x}_i), \epsilon_j) + \\ &\quad \text{cov}(\epsilon_i, f(\mathbf{x}_j)) + \text{cov}(\epsilon_i, \epsilon_j) = \\ &\quad k(\mathbf{x}_i, \mathbf{x}_j) + \sigma^2 \delta_{ij}\end{aligned}$$

or

$$\text{cov}(\mathbf{y}) = k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}.$$

# Gaussian Processes

## Computing the predictions: Joint distribution of observed and test values

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & K(\mathbf{X}, \mathbf{X}_*) \\ K(\mathbf{X}_*, \mathbf{X}) & K(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix} \right)$$

# Gaussian Processes

## Computing the predictions: Posterior predictive distribution

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & K(\mathbf{X}, \mathbf{X}_*) \\ K(\mathbf{X}_*, \mathbf{X}) & K(\mathbf{X}_*, \mathbf{X}_*) \end{bmatrix}\right)$$

$$\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*))$$

where

$$\bar{\mathbf{f}}_* := \mathbb{E}[\mathbf{f}_* | \mathbf{X}, \mathbf{y}, \mathbf{X}_*] = K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(\mathbf{X}_*, \mathbf{X}_*) - K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I})^{-1} K(\mathbf{X}, \mathbf{X}_*).$$



# Gaussian Processes

## Computing the predictions: Posterior predictive distribution

Denote  $\mathbf{K} := K(\mathbf{X}, \mathbf{X})$ ,  $\mathbf{K}_* := K(\mathbf{X}, \mathbf{X}_*)$ ,

$\mathbf{k}_* := \mathbf{k}(\mathbf{x}_*) = [k(\mathbf{x}_i, \mathbf{x}_*)]_{i=1}^n$  – the covariance vector of a test point  $\mathbf{x}_*$  and  $n$  training points.

$$f_* \mid \mathbf{X}, \mathbf{y}, \mathbf{x}_* \sim \mathcal{N}(\bar{f}_*, \mathbb{V}[f_*])$$

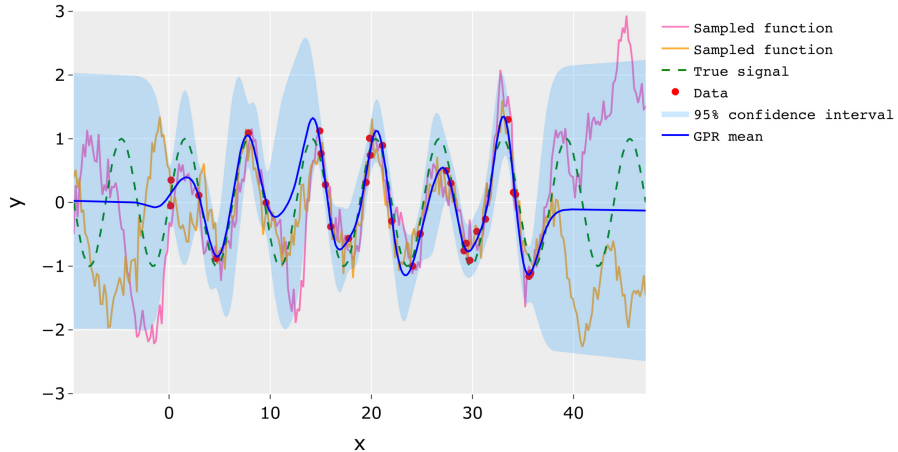
where

$$\bar{f}_* = \mathbf{k}_*^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y},$$

$$\mathbb{V}[f_*] = \mathbf{k}(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_*.$$

# Gaussian Processes

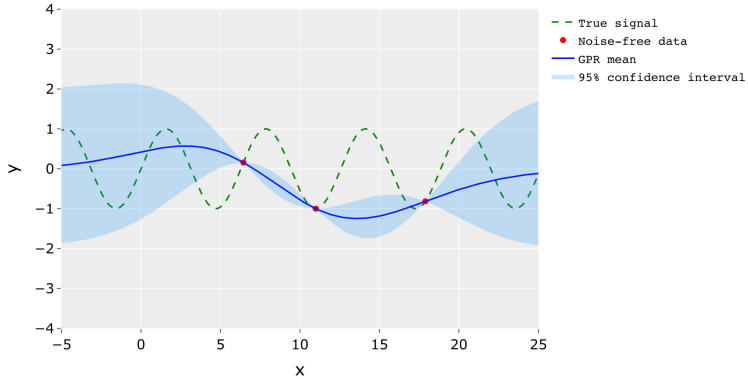
## Example



# Gaussian Processes

## Example: Noise-free case

$$\epsilon = 0 \Leftrightarrow \sigma^2 = 0$$



# Gaussian Processes

## Equivalence to Bayesian Linear Regression

$$\bar{f}(\mathbf{x}_*) = \mathbf{k}(\mathbf{x}_*)^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} = \sum_{i=1}^n \alpha_i \mathbf{y}_i$$

Corresponds to posterior predictive mean of Bayesian Linear Regression (slide 23)

$$\boldsymbol{\phi}_*^T \boldsymbol{\Sigma}_w \boldsymbol{\Phi} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

if we define  $k(\mathbf{x}_i, \mathbf{x}_j) := \boldsymbol{\phi}(\mathbf{x}_i)^T \boldsymbol{\Sigma}_w \boldsymbol{\phi}(\mathbf{x}_j)$ .

# Gaussian Processes

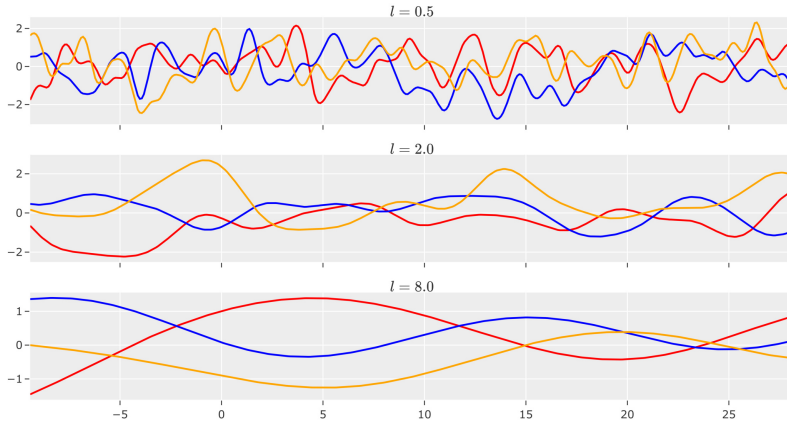
## Hyperparameter optimization: RBF kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp \left( - \frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2 l^2} \right)$$

- How to choose hyperparameters (e.g.  $l$ )?

# Gaussian Processes

## Hyperparameter optimization: RBF kernel



# Gaussian Processes

## Hyperparameter optimization

Let's introduce **marginal likelihood**:

$$p(\mathbf{y} | \mathbf{X}) = \int \overbrace{p(\mathbf{y} | \mathbf{f}, \mathbf{X})}^{\text{likelihood}} \overbrace{p(\mathbf{f} | \mathbf{X})}^{\text{prior}} d\mathbf{f}$$

# Gaussian Processes

## Hyperparameter optimization

Let's introduce **marginal likelihood**:

$$p(\mathbf{y} | \mathbf{X}) = \int \overbrace{p(\mathbf{y} | \mathbf{f}, \mathbf{X})}^{\text{likelihood}} \overbrace{p(\mathbf{f} | \mathbf{X})}^{\text{prior}} d\mathbf{f}$$

Since prior is Gaussian, i.e.  $\mathbf{f} | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$ , then

$$\log p(\mathbf{f} | \mathbf{X}) = -\frac{1}{2} \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f} - \frac{1}{2} \log |\mathbf{K}| - \frac{n}{2} \log 2\pi .$$



# Gaussian Processes

## Hyperparameter optimization

Using the fact that likelihood is also Gaussian,  $\mathbf{y} | \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$ , the final log marginal likelihood is given by

$$\log p(\mathbf{y} | \mathbf{X}) = -\frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K} + \sigma^2 \mathbf{I}| - \frac{n}{2} \log 2\pi .$$

# Gaussian Processes

## Hyperparameter optimization

Using the fact that likelihood is also Gaussian,  $\mathbf{y} | \mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I})$ , the final log marginal likelihood is given by

$$\log p(\mathbf{y} | \mathbf{X}) = -\frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K} + \sigma^2 \mathbf{I}| - \frac{n}{2} \log 2\pi .$$

Let  $\Theta$  be the set of all hyperparameters. Note that  $\mathbf{K}$  is dependent on  $\Theta$ . Then

$$\Theta^* = \underset{\Theta}{\operatorname{argmax}} \log p(\mathbf{y} | \mathbf{X}, \Theta) .$$

# Gaussian Processes

## Advantages vs. Disadvantages

Gaussian Process Regression is a function-space view of regression problem.

### ■ Advantages

- Flexibility and interpretability
- Can exactly be optimized (given the hyperparameters)

### ■ Disadvantages

- Prone to outliers (use the whole data to perform prediction)
- Not efficient for high-dimensional data

Thank you!

## References

- [1] M. Kanagawa, P. Hennig, D. Sejdinovic, and B. K. Sriperumbudur. *Gaussian Processes and Kernel Methods: A Review on Connections and Equivalences*. 2018. [arXiv: 1807.02582 \[stat.ML\]](#).
- [2] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.
- [3] C. E. Rasmussen and C. K. I. Williams. *Gaussian processes for machine learning*. Adaptive computation and machine learning. MIT Press, 2006, pp. I–XVIII, 1–248. ISBN: 026218253X.