

# Advanced Regression Techniques

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

Summer term 2025



Find all slides on [GitHub](#) (DaWe1992/Applied\_ML\_Fundamentals)

# Lecture Overview

- |   |  |
|---|--|
| <b>I</b> Machine Learning Introduction      | <b>IX</b> Evaluation                     |
| <b>II</b> Optimization Techniques           | <b>X</b> Decision Trees                  |
| <b>III</b> Bayesian Decision Theory         | <b>XI</b> Support Vector Machines        |
| <b>IV</b> Non-parametric Density Estimation | <b>XII</b> Clustering                    |
| <b>V</b> Probabilistic Graphical Models     | <b>XIII</b> Principal Component Analysis |
| <b>VI</b> Linear Regression                 | <b>XIV</b> Reinforcement Learning        |
| <b>VII</b> Logistic Regression              | ● <b>XV</b> Advanced Regression          |
| <b>VIII</b> Deep Learning                   |  |

# Agenda for this Unit

① Operations on GAUSSian Distributions

② BAYESian Linear Regression

③ Kernel Ridge Regression

④ GAUSSian Process Regression

⑤ Wrap-Up

Section:  
**Operations on GAUSSIAN Distributions**

- GAUSSian Distribution
- Conditional GAUSSian Distributions
- Marginal GAUSSian Distributions
- BAYES' Theorem for GAUSSian Variables

# Univariate GAUSSIAN Distribution

## Univariate GAUSSIAN distribution:

$$\mathcal{N}(x; \mu, \sigma^2) := \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \quad (1)$$

- A univariate GAUSSIAN random variable is uniquely identified by its **mean**  $\mu \in \mathbb{R}$  and **variance**  $\sigma^2 \in \mathbb{R}^+ := \{x \in \mathbb{R} : x \geq 0\}$
- GAUSSIAN distributions have nice **analytic properties** which we shall exploit

# Multivariate GAUSSIAN Distribution

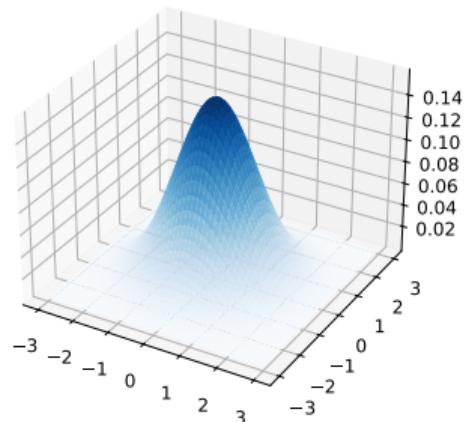
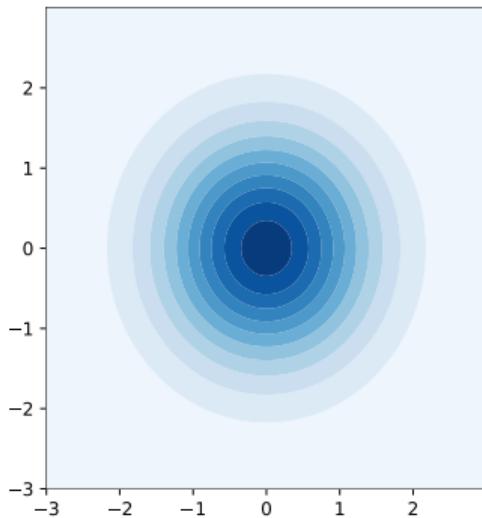
## Multivariate GAUSSIAN distribution:

$$\mathcal{N}_D(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) := \frac{1}{\sqrt{(2\pi)^D \det(\boldsymbol{\Sigma})}} \cdot \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right) \quad (2)$$

- $D \in \mathbb{N}$  is the dimension
- For  $D = 1$  definition (2) reduces to the univariate GAUSSIAN distribution (1)
- A multidimensional GAUSSIAN random variable is uniquely identified by its **mean**  $\boldsymbol{\mu} \in \mathbb{R}^D$  and its **covariance matrix**  $\boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$

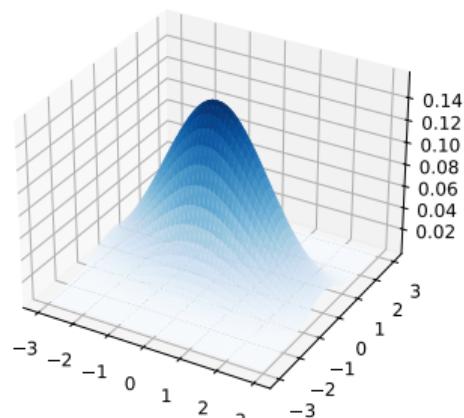
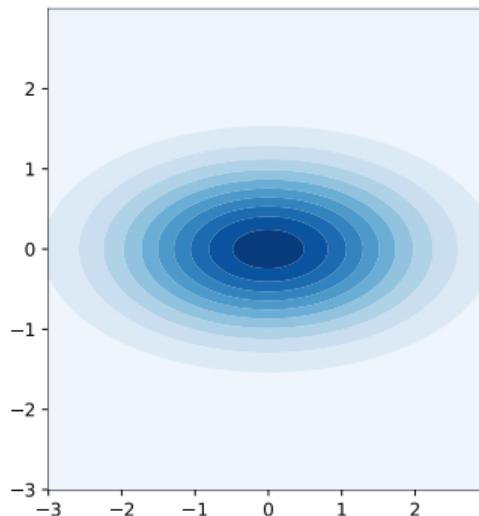
# Visualization of a multivariate GAUSSIAN Distribution

$$\Sigma = \begin{pmatrix} 1.0 & 0.0 \\ 0.0 & 1.0 \end{pmatrix}$$



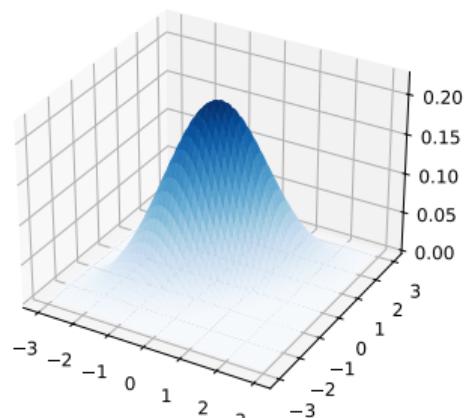
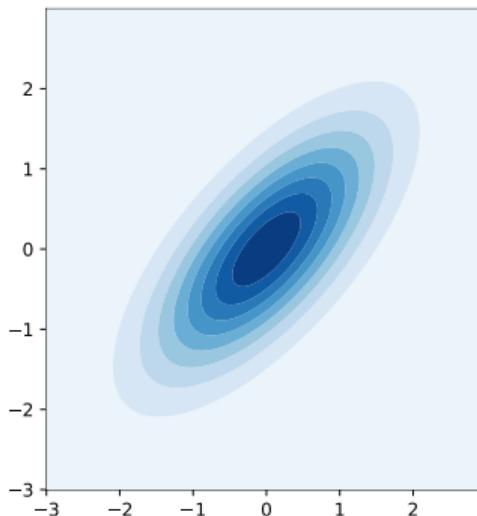
# Visualization of a multivariate GAUSSIAN Distribution (Ctd.)

$$\Sigma = \begin{pmatrix} 2.0 & 0.0 \\ 0.0 & 0.5 \end{pmatrix}$$



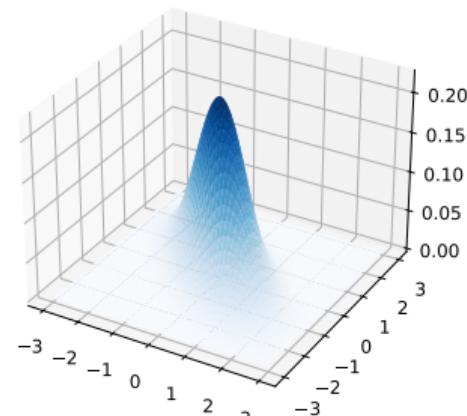
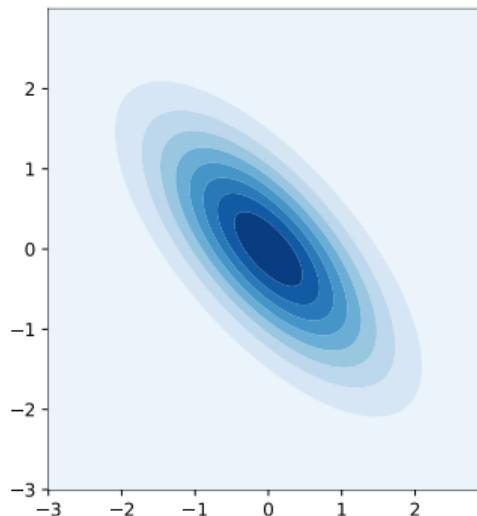
# Visualization of a multivariate GAUSSIAN Distribution (Ctd.)

$$\Sigma = \begin{pmatrix} 1.0 & 0.7 \\ 0.7 & 1.0 \end{pmatrix}$$



# Visualization of a multivariate GAUSSIAN Distribution (Ctd.)

$$\Sigma = \begin{pmatrix} 1.0 & -0.7 \\ -0.7 & 1.0 \end{pmatrix}$$



# Partitioned GAUSSIAN Distribution

- In the following we introduce **marginalization** and **conditioning** of multivariate GAUSSIAN random variables
- Let  $\mathbf{x} \in \mathbb{R}^D$  be a GAUSSIAN random variable which we partition into two disjoint sets  $\mathbf{x}_a$  and  $\mathbf{x}_b$ :

$$p(\mathbf{x}) = p\left(\begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix}\right) \quad (3)$$

- The variables  $\mathbf{x}_a \in \mathbb{R}^K$  and  $\mathbf{x}_b \in \mathbb{R}^{D-K}$  are **jointly GAUSSIAN**
- $\boldsymbol{\Sigma}_{aa}$  and  $\boldsymbol{\Sigma}_{bb}$  are symmetric matrices and we have  $\boldsymbol{\Sigma}_{ba} = \boldsymbol{\Sigma}_{ab}^\top$

# Conditional GAUSSIAN Distribution

- Suppose we observe the variables  $\mathbf{x}_b$  and want to compute the **conditional probability**  $p(\mathbf{x}_a|\mathbf{x}_b)$
- This probability is distributed GAUSSIAN with mean  $\boldsymbol{\mu}_{a|b}$  and covariance  $\boldsymbol{\Sigma}_{a|b}$ :

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a; \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b}) \quad (4)$$

- The parameters  $\boldsymbol{\mu}_{a|b}$  and  $\boldsymbol{\Sigma}_{a|b}$  can be computed in closed form:

$$\boldsymbol{\mu}_{a|b} := \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (5)$$

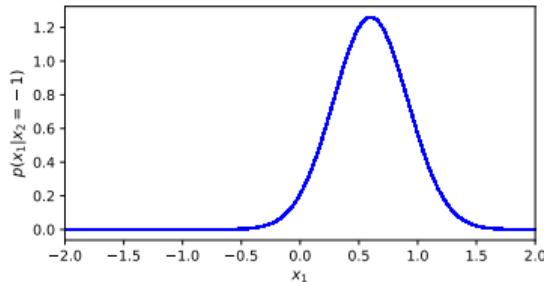
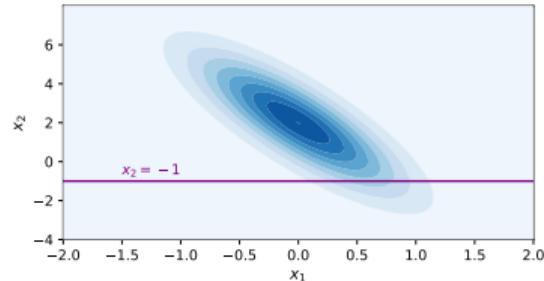
$$\boldsymbol{\Sigma}_{a|b} := \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba} \quad (6)$$

# Example: Conditional GAUSSIAN Distribution

- Consider the GAUSSIAN

$$p(x_1, x_2) = \mathcal{N} \left( \begin{bmatrix} 0 \\ 2 \end{bmatrix}, \begin{bmatrix} 0.3 & -1 \\ -1 & 5 \end{bmatrix} \right)$$

- Question:** What is  $p(x_1|x_2 = -1)$ ?
- According to equations (5) and (6) we obtain  $\mu_{x_1|x_2} = 0.6$  and  $\sigma_{x_1|x_2}^2 = 0.1$
- Thus,  $p(x_1|x_2 = -1) = \mathcal{N}(0.6, 0.1)$



# Marginal GAUSSIAN Distribution

- Sometimes we may be interested in the **marginal distribution**  $p(\mathbf{x}_a)$
- The marginal is distributed GAUSSIAN and is given by:

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}) \quad (7)$$

- Intuitively, we ignore everything we are not interested in, i. e. we integrate out  $\mathbf{x}_b$  from the joint distribution  $p(\mathbf{x}_a, \mathbf{x}_b)$  given by equation (3)
- The corresponding results hold for  $p(\mathbf{x}_b)$  which is obtained by marginalization of  $p(\mathbf{x}_a, \mathbf{x}_b)$  over  $\mathbf{x}_a$

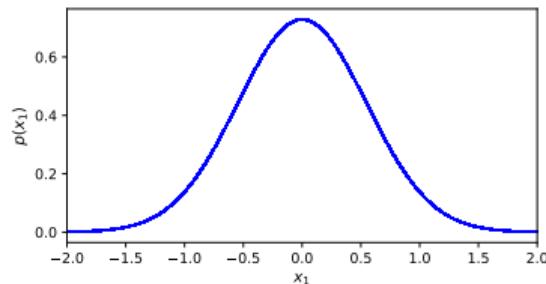
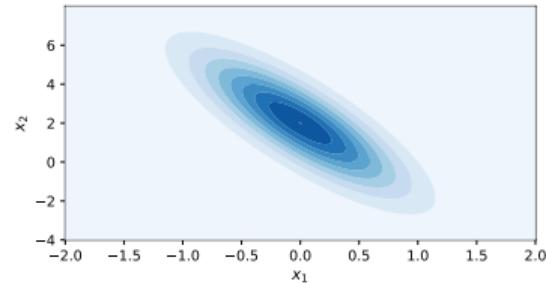
## Example: Marginal GAUSSIAN Distribution

- Consider again

$$p(x_1, x_2) = \mathcal{N} \left( \begin{bmatrix} 0 \\ 2 \end{bmatrix}, \begin{bmatrix} 0.3 & -1 \\ -1 & 5 \end{bmatrix} \right)$$

- According to equation (7), the marginal distribution  $p(x_1)$  is given by:

$$p(x_1) = \mathcal{N}(0, 0.3)$$



# BAYES' Theorem for GAUSSIAN Variables

- Let a marginal GAUSSIAN distribution for  $\mathbf{x}$  and a conditional GAUSSIAN distribution for  $\mathbf{y}$  given  $\mathbf{x}$  be given:

$$p(\mathbf{x}) := \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \text{and} \quad p(\mathbf{y}|\mathbf{x}) := \mathcal{N}(\mathbf{y}; \mathbf{A}\mathbf{x}, \boldsymbol{\Lambda}) \quad (8)$$

- Then we obtain for the GAUSSIAN distributions  $p(\mathbf{y})$  and  $p(\mathbf{x}|\mathbf{y})$ :

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top + \boldsymbol{\Lambda}) \quad (9)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\Gamma}[\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \mathbf{A}^\top \boldsymbol{\Lambda}^{-1}\mathbf{y}], \boldsymbol{\Gamma}) \quad (10)$$

$$\boldsymbol{\Gamma} := (\boldsymbol{\Sigma}^{-1} + \mathbf{A}^\top \boldsymbol{\Lambda}^{-1} \mathbf{A})^{-1}$$

# Summary: GAUSSIAN Distributions

- GAUSSIAN distributions play a major role in machine learning
- Complex operations have **analytical solutions**
- Important operations on GAUSSIANS (i. e. **conditioning, marginalization, convolution**) result again in **GAUSSIAN distributions**
- Remember:

Once a **GAUSSIAN**, always a **GAUSSIAN!**

## Section: **BAYESian Linear Regression**

- Maximum Likelihood Estimation (MLE)
- Maximum A Posteriori (MAP) Estimation
- Full BAYESian Linear Regression

# Remember: A probabilistic View on Regression

- **Assumption 1:** The target values  $y$  and the inputs  $\mathbf{x}$  are related via the equation

$$y = h_{\theta}(\mathbf{x}) + \varepsilon = \boldsymbol{\theta}^{\top} \mathbf{x} + \varepsilon, \quad (11)$$

where  $\varepsilon$  is an error term which captures unmodeled effects or noise in the data

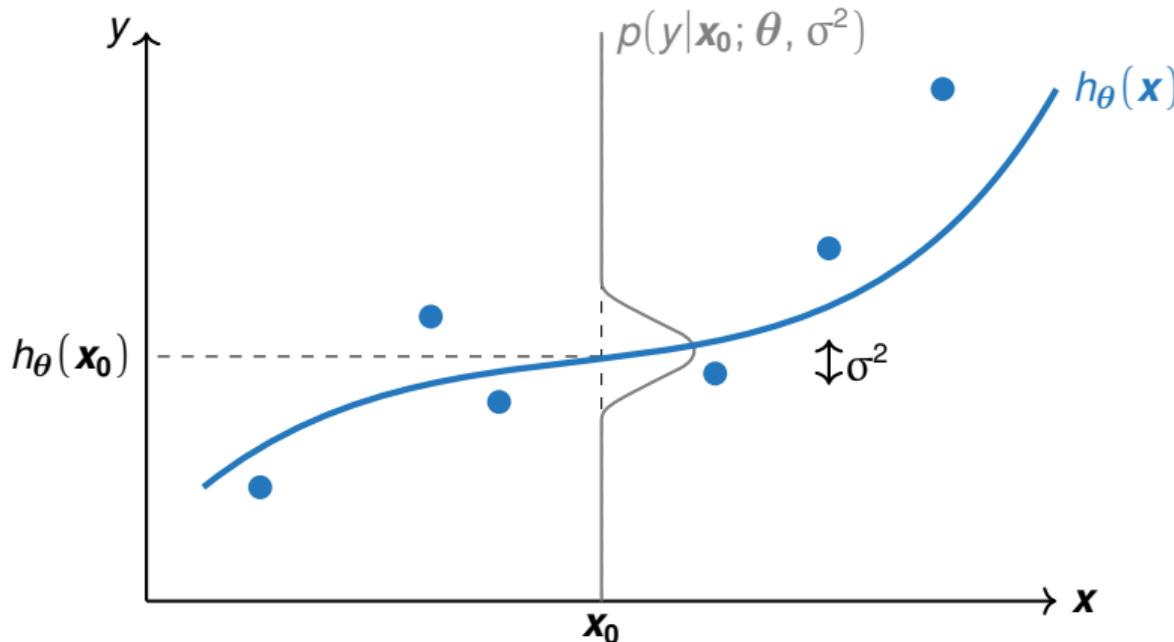
- **Assumption 2:** The noise  $\varepsilon$  is a zero mean GAUSSIAN random variable

$$\varepsilon \sim \mathcal{N}(0, \sigma^2) \quad (12)$$

- We consider  $y$  a random variable which is distributed according to

$$p(y|\mathbf{x}; \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(y; h_{\theta}(\mathbf{x}), \sigma^2) \quad (13)$$

## Remember: A probabilistic View on Regression (Ctd.)



# Remember: Likelihood Function for Regression

- We are given a training dataset  $\mathcal{D} := \{(\mathbf{x}^n, y_n)\}_{n=1}^N$
- The **(conditional) likelihood** of our regression model is given by:

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}, \sigma^2) &= \prod_{n=1}^N \mathcal{N}(y_n; h_{\boldsymbol{\theta}}(\mathbf{x}^n), \sigma^2) = \prod_{n=1}^N \mathcal{N}(y_n; \boldsymbol{\theta}^\top \boldsymbol{\varphi}(\mathbf{x}^n), \sigma^2) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{1}{2\sigma^2}(y_n - \boldsymbol{\theta}^\top \boldsymbol{\varphi}(\mathbf{x}^n))^2\right) \end{aligned} \quad (14)$$

## Remember: Likelihood Function for Regression (Ctd.)

- The **log-likelihood** is then given by (*we have computed this earlier already*):

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &:= \log p(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}, \sigma^2) \\ &= -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{\sigma^2} \cdot \underbrace{\frac{1}{2} \sum_{n=1}^N (y_n - \boldsymbol{\theta}^\top \varphi(\mathbf{x}^n))^2}_{\text{least squares error}}\end{aligned}\tag{15}$$

- We have to **minimize the least squares error** to **maximize the likelihood!**

When minimizing the squared error we implicitly assume Gaussian noise!

## Remember: The Maximum Likelihood Solution to Regression

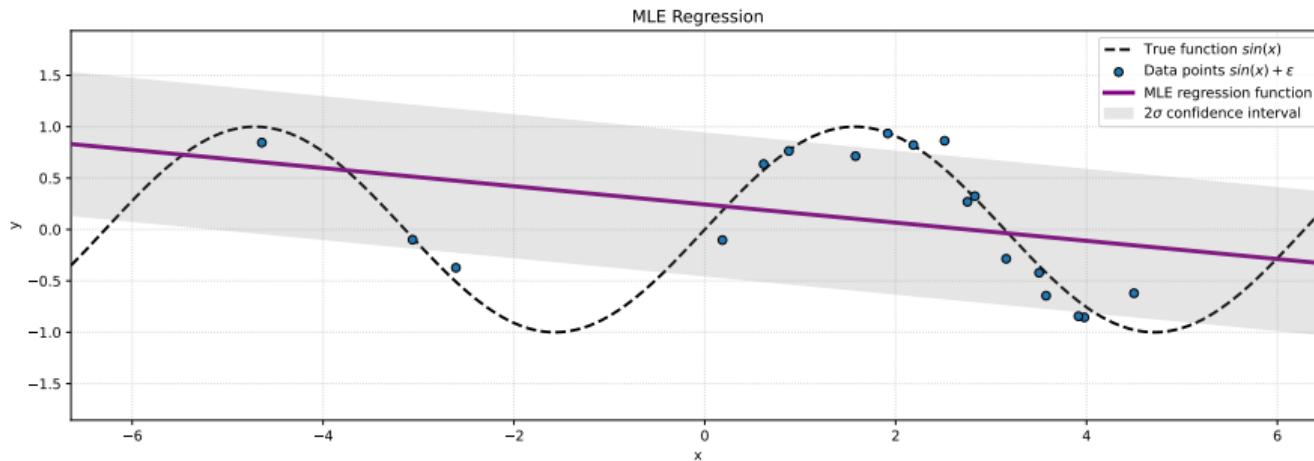
- Optimization of the log-likelihood function gives:

$$\boldsymbol{\theta}^{\text{ML}} := (\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{y} \quad (16)$$

$$(\sigma^2)^{\text{ML}} := \frac{1}{N} \sum_{n=1}^N (y_n - (\boldsymbol{\theta}^{\text{ML}})^\top \boldsymbol{\varphi}(\mathbf{x}^n))^2 \quad (17)$$

- The probabilistic interpretation of linear regression allows us **to quantify the (global) uncertainty** of the model

# Example



## Adding a Prior Distribution

- We have just seen that **maximum likelihood estimation is prone to overfitting**
- We now mitigate this effect by placing a **GAUSSIAN prior distribution**

$$p(\boldsymbol{\theta}; b^2) := \mathcal{N}(\mathbf{0}, b^2 \mathbf{I}) \quad (18)$$

on the parameters  $\boldsymbol{\theta}$

- The prior distribution encodes what parameter values are plausible (**prior knowledge**) before having seen any data
- Instead of maximizing the likelihood, we seek parameters  $\boldsymbol{\theta}^{\text{MAP}}$  which **maximize the posterior distribution**  $p(\boldsymbol{\theta}|\mathcal{D})$  given a dataset  $\mathcal{D}$

## Adding a Prior Distribution (Ctd.)

**Question:** Why do we use GAUSSIANS to model the prior?

- We set  $p(\theta; b^2) := \mathcal{N}(\mathbf{0}, b^2 \mathbf{I})$
- The most important argument is **conjugacy**: Both, the prior and the likelihood are modelled by GAUSSIAN distributions. Therefore, the product of both – the posterior distribution – will also be GAUSSIAN (*only few distributions have that property*)
- There are **analytical solutions to complex operations** on GAUSSIANS which again result in GAUSSIAN distributions (*see first section for details*)

# Recall BAYES' Theorem

In the following we shall often refer to **BAYES' theorem**:

## BAYES' Theorem:

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{\overbrace{p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}; \sigma^2)}^{\text{likelihood}} \cdot \overbrace{p(\boldsymbol{\theta}; b^2)}^{\text{prior}}}{p(\mathbf{y}|\mathbf{X})} \propto p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}; \sigma^2) \cdot p(\boldsymbol{\theta}; b^2) \quad (19)$$

## Remarks:

- Until now we have maximized the likelihood without considering  $p(\boldsymbol{\theta}; b^2)$
- $\boldsymbol{\theta}$  is now a random variable!

# Log-Posterior Distribution

- We begin by computing the **log-posterior**:

$$\log p(\boldsymbol{\theta}|\mathcal{D}) = \underbrace{\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}; \sigma^2)}_{\text{log-likelihood}} + \underbrace{\log p(\boldsymbol{\theta}; b^2)}_{\text{log-prior}} + \text{const} \quad (20)$$

- The constant ‘const’ contains all terms independent of  $\boldsymbol{\theta}$
- The MAP (*maximum a posteriori*) estimate will be a **compromise** between the parameter prior and the data-dependent likelihood
- As usual, we minimize the **negative log-posterior**

$$\boldsymbol{\theta}^{\text{MAP}} := \arg \min_{\boldsymbol{\theta}} \{-\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}; \sigma^2) - \log p(\boldsymbol{\theta}; b^2)\} \quad (21)$$



# Rewriting the Log-Likelihood Function for Regression

The log-likelihood in equation (15) can be rewritten in matrix-vector notation:

- The sum in the second term can be rewritten as:

$$\sum_{n=1}^N (y_n - \boldsymbol{\theta}^\top \varphi(\mathbf{x}^n))^2 \iff (\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta})^\top (\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta}) \quad (22)$$

- Also, we notice that the first term  $-\frac{N}{2} \log(2\pi\sigma^2)$  is constant with respect to  $\boldsymbol{\theta}$
- Thus,

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}; \sigma^2) = -\frac{1}{2\sigma^2} (\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta})^\top (\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta}) + \text{const} \quad (23)$$



# Rewriting the negative Log-Posterior

**Negative log-posterior:**

$$\begin{aligned} -\log p(\boldsymbol{\theta}|\mathcal{D}) &= -\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}; \sigma^2) - \log p(\boldsymbol{\theta}; b^2) + \text{const} \\ &= \frac{1}{2\sigma^2} (\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta})^\top (\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta}) + \frac{1}{2b^2} \boldsymbol{\theta}^\top \boldsymbol{\theta} + \text{const} \\ &= \frac{1}{2\sigma^2} (\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \boldsymbol{\Phi}\boldsymbol{\theta} - (\boldsymbol{\Phi}\boldsymbol{\theta})^\top \mathbf{y} + (\boldsymbol{\Phi}\boldsymbol{\theta})^\top \boldsymbol{\Phi}\boldsymbol{\theta}) + \frac{1}{2b^2} \boldsymbol{\theta}^\top \boldsymbol{\theta} + \text{const} \\ &= \frac{1}{2\sigma^2} (\boldsymbol{\theta}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi}\boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \boldsymbol{\Phi}^\top \mathbf{y}) + \frac{1}{2b^2} \boldsymbol{\theta}^\top \boldsymbol{\theta} + \text{const} \end{aligned} \tag{24}$$

# Remember: Transposition Rules and Vector Derivatives

Remember the following rules:

## Matrix transposition rules:

$$(\mathbf{A}^T)^T = \mathbf{A} \quad (25)$$

$$(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T \quad (26)$$

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T \quad (27)$$

## Vector derivatives:

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{Ax} = 2\mathbf{Ax} \quad (28)$$

[Equation (28) only holds if  $\mathbf{A}$  is symmetric]

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{a}^T \mathbf{x} = \mathbf{a} \quad (29)$$



# Computation of the Gradient

**Gradient of the negative log-posterior:**

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\theta}} - \log p(\boldsymbol{\theta} | \mathcal{D}) &= \frac{\partial}{\partial \boldsymbol{\theta}} \frac{1}{2\sigma^2} (\boldsymbol{\theta}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \boldsymbol{\Phi}^\top \mathbf{y}) + \frac{\partial}{\partial \boldsymbol{\theta}} \frac{1}{2b^2} \boldsymbol{\theta}^\top \boldsymbol{\theta} \\ &= \frac{1}{\sigma^2} (\boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\theta} - \boldsymbol{\Phi}^\top \mathbf{y}) + \frac{1}{b^2} \boldsymbol{\theta}\end{aligned}\tag{30}$$



# Computation of the MAP Estimate

**Set gradient to zero:**

$$\begin{aligned} \frac{1}{\sigma^2}(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\theta} - \boldsymbol{\Phi}^\top \mathbf{y}) + \frac{1}{b^2} \boldsymbol{\theta} &\stackrel{!}{=} \mathbf{0} \\ \iff \quad \left( \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{\sigma^2}{b^2} \mathbf{I} \right) \boldsymbol{\theta} &= \boldsymbol{\Phi}^\top \mathbf{y} \\ \iff \quad \boldsymbol{\theta}^{\text{MAP}} &= \left( \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \boldsymbol{\Phi}^\top \mathbf{y} \end{aligned} \quad (31)$$

**Remark:** The matrix  $\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{\sigma^2}{b^2} \mathbf{I}$  is positive definite



# Matrix Invertibility

**Lemma:** The matrix  $\Phi^\top \Phi + \frac{\sigma^2}{b^2} I$  is positive definite (and therefore invertible).

**Proof:** Let  $\mathbf{z} \in \mathbb{R}^M \setminus \{\mathbf{0}\}$  be a non-zero vector. We obtain

$$\mathbf{z}^\top \left( \Phi^\top \Phi + \frac{\sigma^2}{b^2} I \right) \mathbf{z} = \mathbf{z}^\top \Phi^\top \Phi \mathbf{z} + \frac{\sigma^2}{b^2} \mathbf{z}^\top I \mathbf{z} > 0 \quad (32)$$

due to the positive semidefiniteness of  $\Phi^\top \Phi$ , the positive definiteness of  $I$ , and the positivity of  $\sigma^2$  and  $b^2$ . ■

# Comparison of MLE and MAP Estimates

## MLE Estimate:

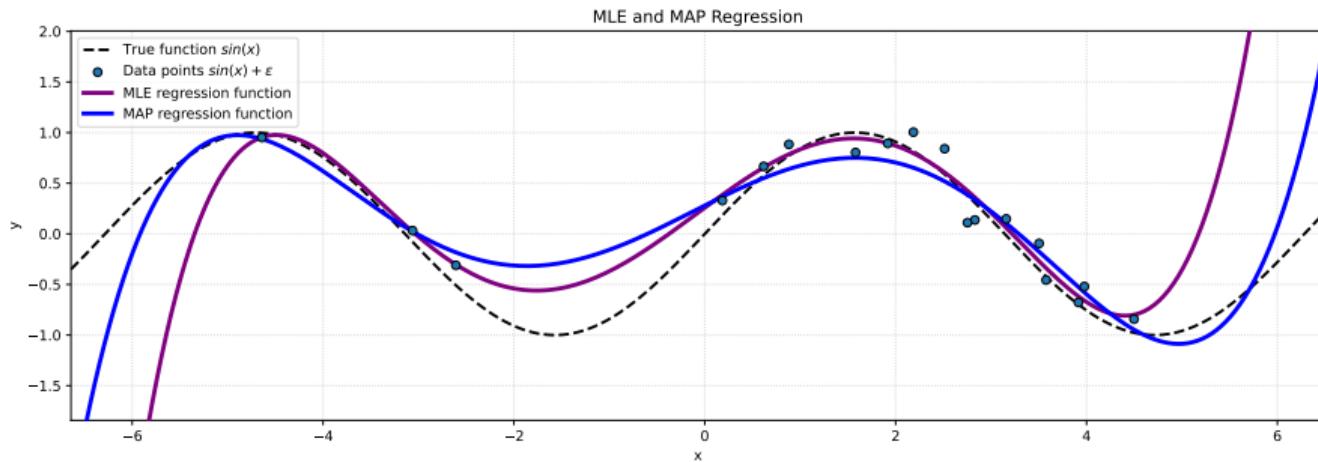
$$\boldsymbol{\theta}^{\text{ML}} := (\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{y} \quad (33)$$

## MAP Estimate:

$$\boldsymbol{\theta}^{\text{MAP}} := \left( \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \boldsymbol{\Phi}^\top \mathbf{y} \quad (34)$$

For  $\lambda := \frac{\sigma^2}{b^2}$  the MAP estimate is identical to ridge regression!

# Comparison of MLE and MAP Estimates (Ctd.)

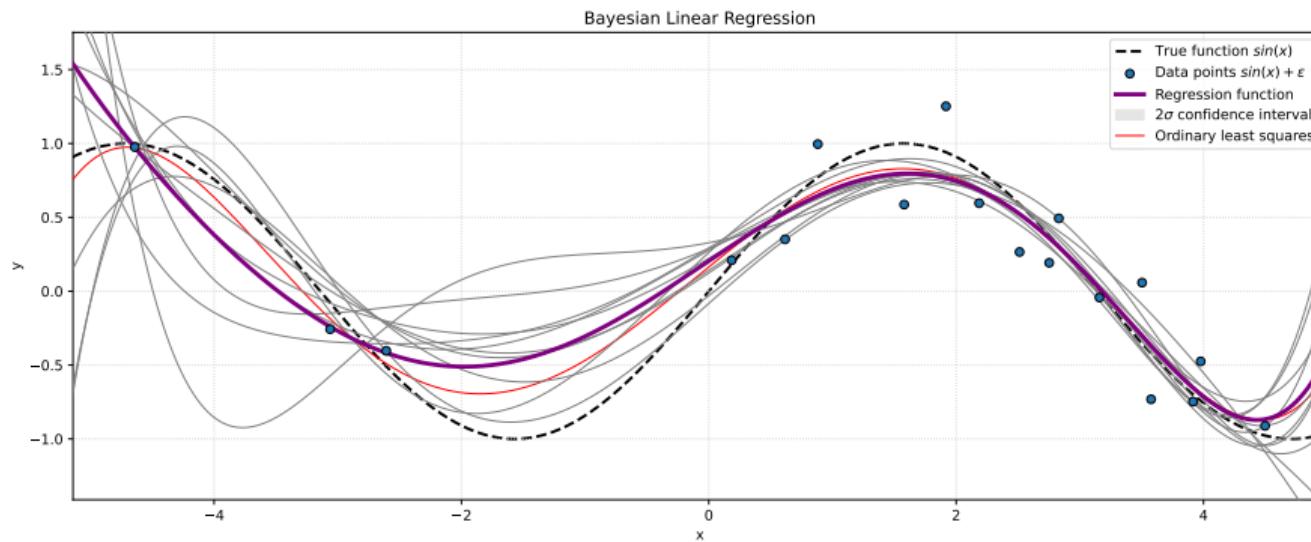


$$\sigma^2 := 0.2 \text{ and } b^2 := 0.1$$

# Full BAYESIAN Regression

- So far we have seen that MLE estimates tend to overfit the training data
- In MAP estimation we place a prior on the parameters  $\theta$  which plays the role of a **regularizer**
- **BAYESIAN linear regression** pushes this idea to the next level:
  - Again we consider a prior distribution  $p(\theta)$  of the parameters
  - But: We **do not** compute a point estimate of the parameters
  - Instead, we take the **full posterior distribution** into account when making predictions
  - **This means we do not fit any parameters, but we compute a mean over all plausible parameter settings (according to the posterior)**

# What we want to achieve:



# BAYESian Regression Model

In BAYESian linear regression we consider the following model:

$$\text{Prior} \quad p(\boldsymbol{\theta}) := \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad (35)$$

$$\text{Likelihood} \quad p(y|\mathbf{x}', \boldsymbol{\theta}) := \mathcal{N}(y; \boldsymbol{\theta}^\top \boldsymbol{\varphi}(\mathbf{x}'), \sigma^2) \quad (36)$$

$\mathbf{x}'$  is a test data point,  $\boldsymbol{\mu}_0$  is the mean of the prior distribution (usually we set  $\boldsymbol{\mu}_0 := \mathbf{0}$ ), and  $\boldsymbol{\Sigma}_0$  its covariance matrix

## Prior Predictions

- In practice, we are (usually) not interested in the parameter values themselves
- Instead, we want to **predict labels of unseen instances**
- Let  $\mathbf{x}'$  be a test data point whose label  $y'$  we want to predict
- The **prior predictive distribution** is given by

$$p(y|\mathbf{x}') = \int_{\theta} p(y|\mathbf{x}', \theta)p(\theta) d\theta = \mathbb{E}_{\theta}\{p(y|\mathbf{x}', \theta)\} \quad (37)$$

- We can interpret equation (37) as the **average prediction** for all plausible parameter values according to the prior distribution

## Prior Predictions (Ctd.)

### Parameters of the (prior) predictive distribution:

Due to the **conjugate GAUSSIAN prior**  $p(\theta) = \mathcal{N}(\mu_0, \Sigma_0)$ , we can compute the parameters of the predictive distribution in **closed form**:

$$p(y|\mathbf{x}') = \mathcal{N}(y; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (38)$$

$$\text{with} \quad \boldsymbol{\mu} := \boldsymbol{\mu}_0^\top \boldsymbol{\varphi}(\mathbf{x}') \quad (39)$$

$$\boldsymbol{\Sigma} := \boldsymbol{\varphi}(\mathbf{x}')^\top \boldsymbol{\Sigma}_0 \boldsymbol{\varphi}(\mathbf{x}') + \sigma^2 \quad (40)$$

**Remark:** We have used equation (9) here

## Prior Predictions (Ctd.)

- In equation (40), the term  $\varphi(\mathbf{x}')^\top \Sigma_0 \varphi(\mathbf{x}')$  accounts for the **uncertainty** associated with the parameters  $\theta$
- $\sigma^2$  is the uncertainty contribution due to measurement noise
- Every  $\theta^i$  we sample from the prior distribution gives rise to a function

$$f_i(\cdot) := (\theta^i)^\top \varphi(\cdot) \quad (41)$$

- The parameter distribution  $p(\theta)$  induces a distribution  $p(f(\cdot))$  over functions
- The prior predictions are made **without taking the training data into account**, therefore these functions do not fit the data well

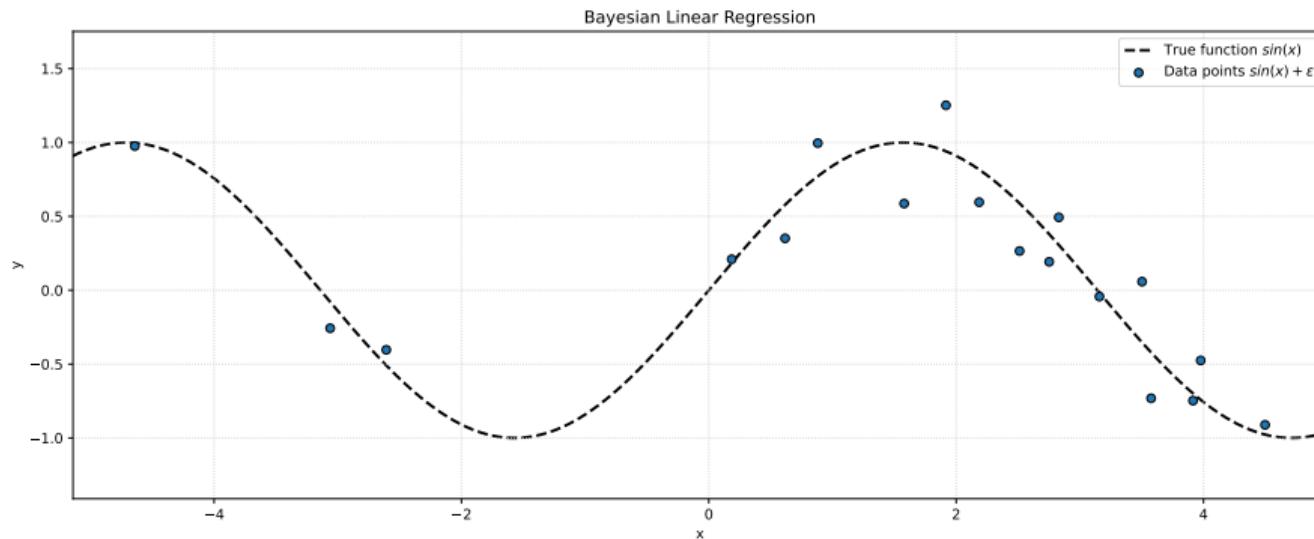
## Example: Setup

- Let us have a look at an example
- Suppose the labels are generated by the **true function**  $f(x) = \sin(x)$  plus some additive noise  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$
- A training data point  $(x_n, y_n)$  then takes the form

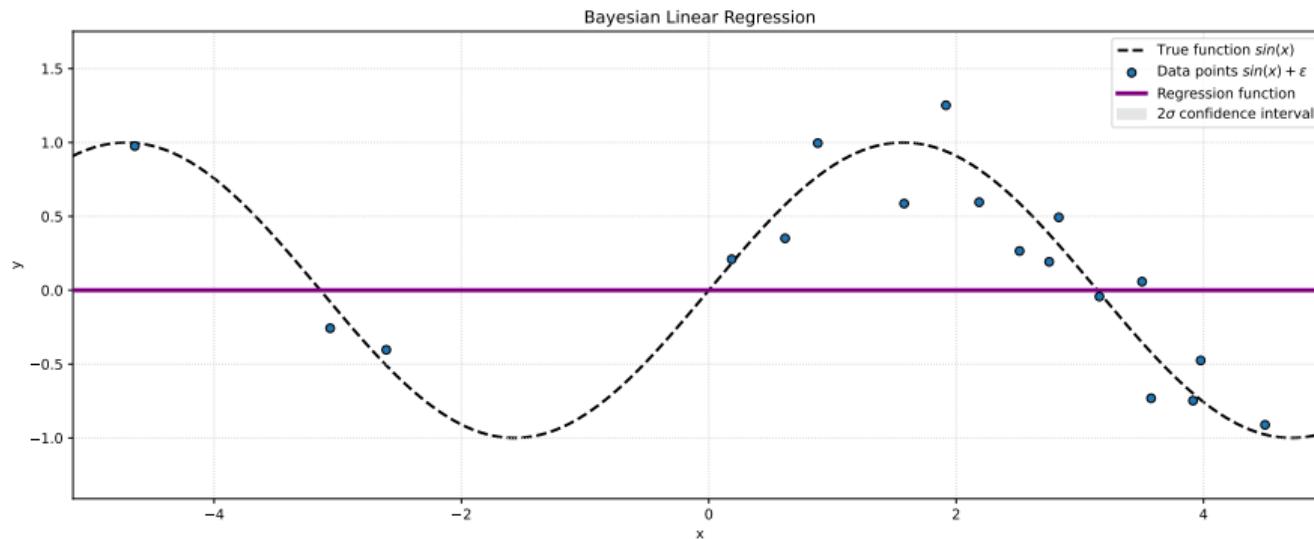
$$(x_n, y_n) = (x_n, \sin(x_n) + \varepsilon_n) \quad n = 1, 2, \dots, N \quad (42)$$

- We initialize the **parameter prior distribution** to be  $p(\theta) = \mathcal{N}(\theta | \mathbf{0}, I)$ , i. e. zero mean and **isotropic** (*i. e. rotation invariant*) covariance

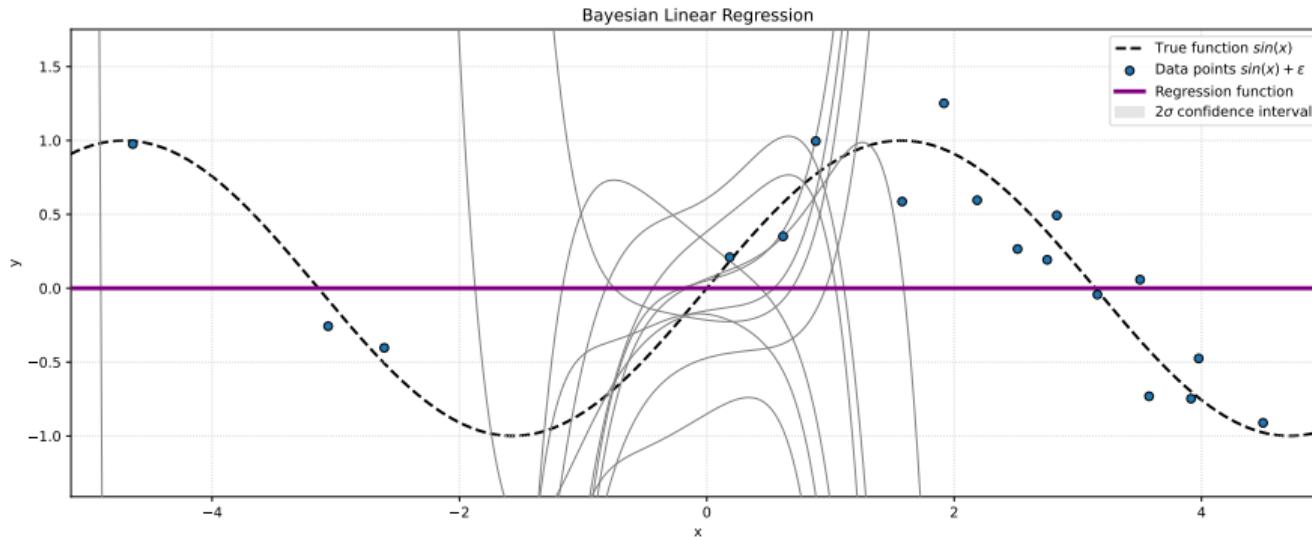
# Example: Dataset



# Example: Prior Predictions



# Example: Sampled Functions



# Incorporation of the Training Data

- So far we have used the **prior distribution**  $p(\theta)$  when computing predictions
- In the following we will replace  $p(\theta)$  with the **posterior distribution**  $p(\theta|\mathcal{D})$  of  $\theta$  given the training data  $\mathcal{D}$
- We compute the posterior distribution according to BAYES' theorem analogously to equation (19):

$$p(\theta|\mathcal{D}) = \frac{\underbrace{p(\mathbf{y}|\mathbf{X}, \theta)}_{likelihood} \cdot \underbrace{p(\theta)}_{prior}}{p(\mathbf{y}|\mathbf{X})} \quad (43)$$

- The parameter posterior is **again Gaussian**

# Posterior Parameter Distribution

## Parameter posterior distribution:

The parameter posterior distribution is given by

$$p(\boldsymbol{\theta}|\mathcal{D}) = \mathcal{N}(\boldsymbol{\theta}; \boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N) \quad (44)$$

with  $\boldsymbol{\Sigma}_N := (\boldsymbol{\Sigma}_0^{-1} + \sigma^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}$  (45)

$$\boldsymbol{\mu}_N := \boldsymbol{\Sigma}_N (\boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \sigma^{-2} \boldsymbol{\Phi}^\top \mathbf{y}) \quad (46)$$

**Remark:** We have used equation (10) here

## Posterior Predictions (Ctd.)

### Parameters of the (posterior) predictive distribution:

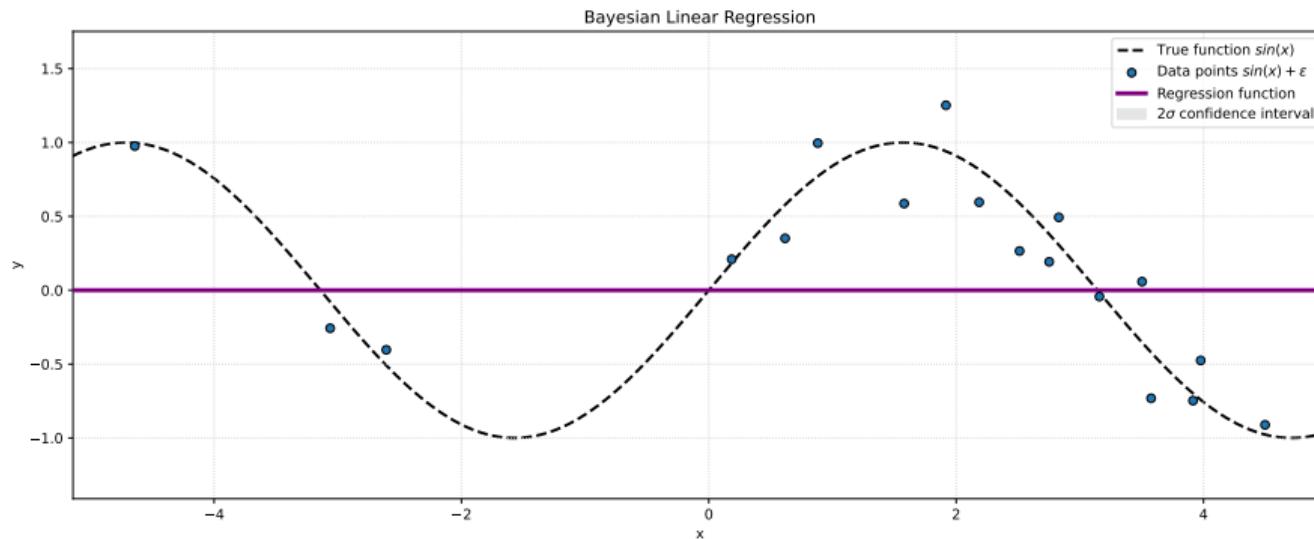
The parameters of the posterior predictive distribution are computed analogously to those of the prior predictive distribution. We replace  $\mu_0$  and  $\Sigma_0$  with  $\mu_N$  and  $\Sigma_N$ , respectively:

$$p(y|\mathbf{x}') = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (47)$$

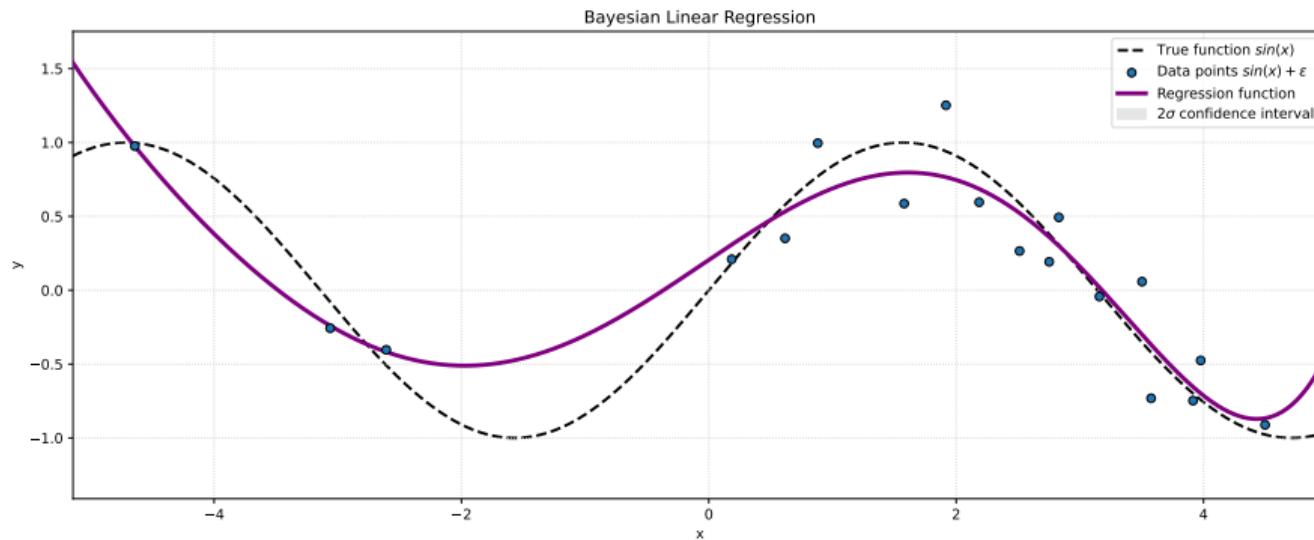
with 
$$\boldsymbol{\mu} := \boldsymbol{\mu}_N^\top \boldsymbol{\varphi}(\mathbf{x}') \quad (48)$$

$$\boldsymbol{\Sigma} := \boldsymbol{\varphi}(\mathbf{x}')^\top \boldsymbol{\Sigma}_N \boldsymbol{\varphi}(\mathbf{x}') + \sigma^2 \quad (49)$$

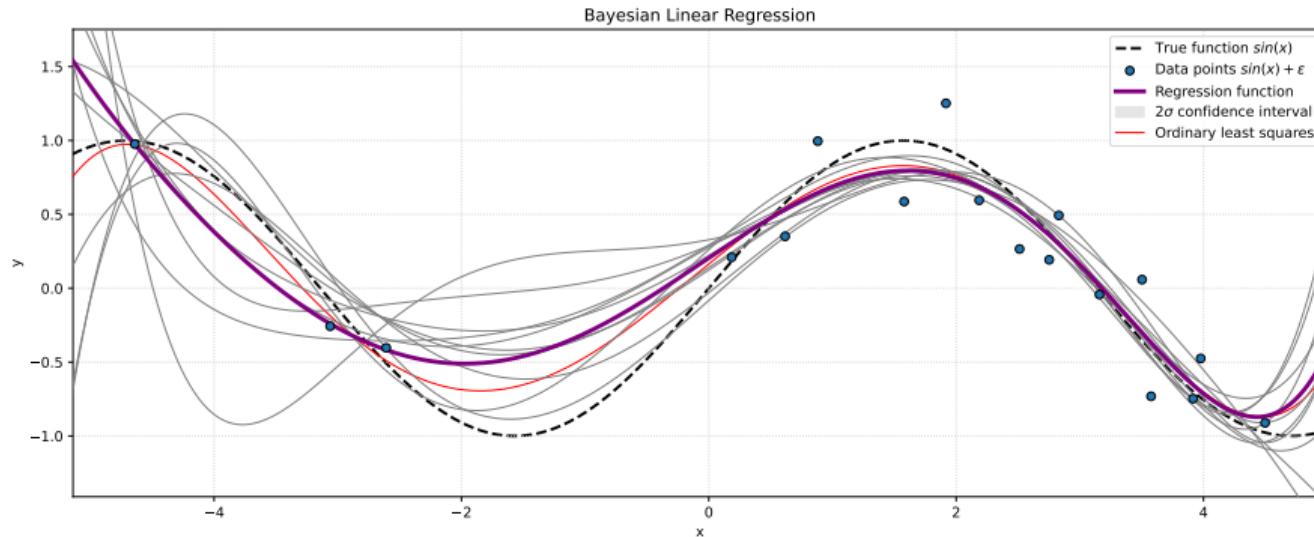
# Example: From Prior Predictions to Posterior Predictions



# Example: From Prior Predictions to Posterior Predictions (Ctd.)



# Example: From Prior Predictions to Posterior Predictions (Ctd.)

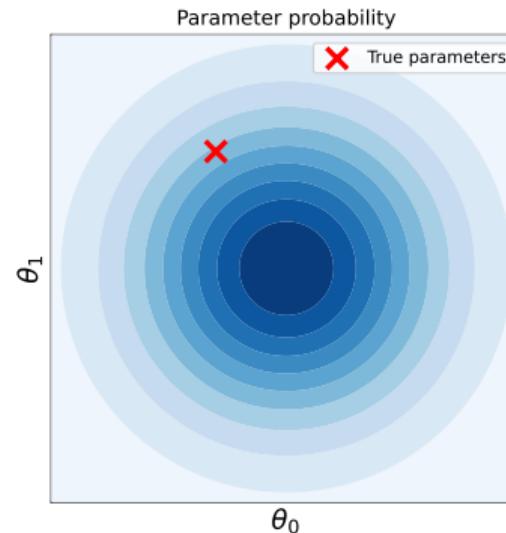


# Illustration of BAYESIAN Learning

- We sample data points from the function

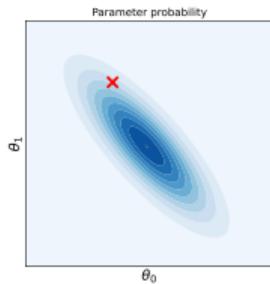
$$f(x) = 0.5x - 0.3 + \varepsilon$$

- $\varepsilon$  is additive noise with mean  $\mu_\varepsilon = 0$  and variance  $\sigma_\varepsilon^2 = 1$
- We place a **prior** on the parameters:  
 $p(\boldsymbol{\theta}) := \mathcal{N}(\mathbf{0}, \mathbf{I}_2)$

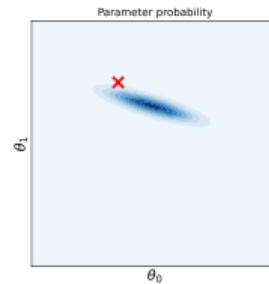


# Illustration of BAYESIAN Learning (Ctd.)

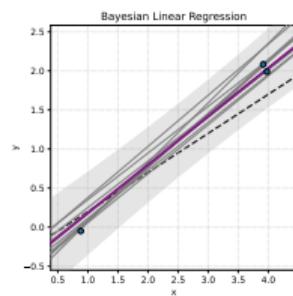
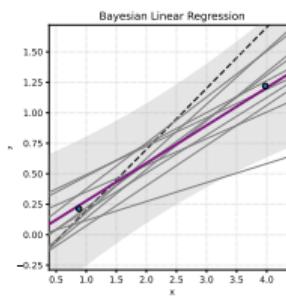
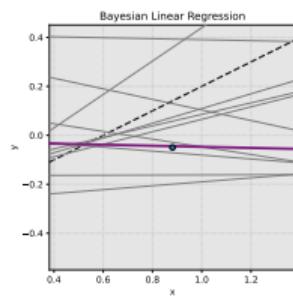
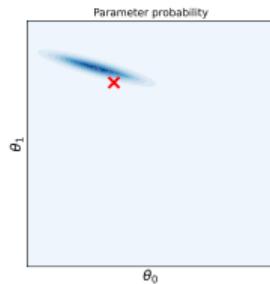
$N = 1$ :



$N = 2$ :

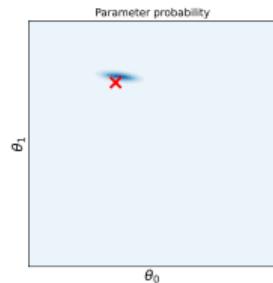


$N = 3$ :

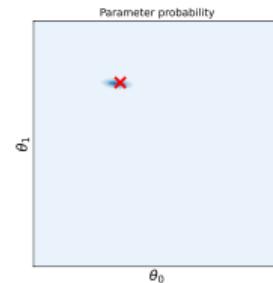


# Illustration of BAYESIAN Learning (Ctd.)

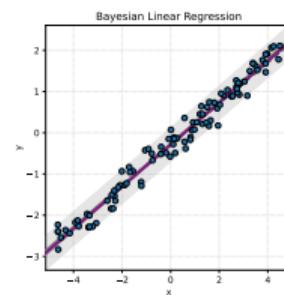
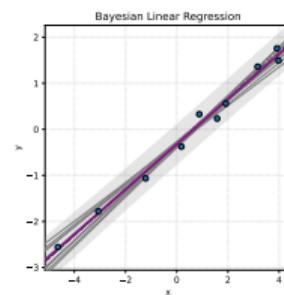
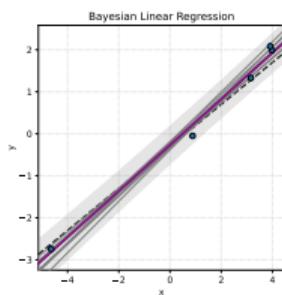
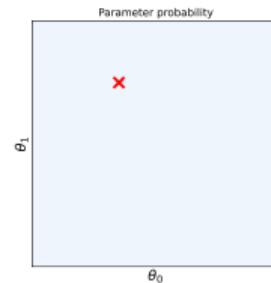
$N = 5$ :



$N = 10$ :



$N = 100$ :



## Section: Kernel Ridge Regression

- Introduction
- Woodbury Matrix Identity
- Derivation of the Algorithm
- Example

# Introduction of Kernels for Regression

- We now turn to a **kernel approach** to regression.
- Remember the **ridge regression** formula:

$$\boldsymbol{\theta} = (\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^\top \mathbf{y}, \quad (50)$$

where  $\lambda \geqslant 0$  is the regularization parameter

- Recall: In order to apply kernels, we have to rephrase this equation in terms of **dot products of the input features**



# A useful Matrix Identity

## WOODBURY matrix identity / ‘push-through’ identity:

The following identity holds for square matrices  $P$  and  $Q$ :

$$(PQ + I)^{-1}P = P(QP + I)^{-1} \quad (51)$$

**Proof:** Multiply the equation by  $(PQ + I)$  from the left, and by  $(QP + I)$  from the right. The resulting equation is  $P(QP + I) = (PQ + I)P$ . This equation is straightforward to verify by performing the multiplications according to the distributive law. ■

# Derivation of Kernel Ridge Regression

- We apply equation (51) to the matrix  $(\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top$
- For this we set:

$$\mathbf{P} := \Phi^\top \quad \mathbf{Q} := \Phi$$

- We receive:

$$(\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top = \Phi^\top (\Phi \Phi^\top + \lambda I)^{-1} \quad (52)$$

- We plug this result into equation (50) to obtain

$$\boldsymbol{\theta} = (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top \mathbf{y} = \Phi^\top (\Phi \Phi^\top + \lambda I)^{-1} \mathbf{y} \quad (53)$$

## Derivation of Kernel Ridge Regression (Ctd.)

To predict the label  $y'$  of an unseen instance  $\mathbf{x}'$  we compute:

$$\begin{aligned}y' &= \boldsymbol{\theta}^\top \boldsymbol{\varphi}(\mathbf{x}') \\&= (\boldsymbol{\Phi}^\top (\boldsymbol{\Phi}\boldsymbol{\Phi}^\top + \lambda\mathbf{I})^{-1} \mathbf{y})^\top \boldsymbol{\varphi}(\mathbf{x}') \\&= \mathbf{y}^\top (\boldsymbol{\Phi}\boldsymbol{\Phi}^\top + \lambda\mathbf{I})^{-1} \boldsymbol{\Phi} \boldsymbol{\varphi}(\mathbf{x}')\end{aligned}$$

[The features appear exclusively in terms of inner products – **we can apply kernels!**]

$$= \boxed{\mathbf{y}^\top (\boldsymbol{\kappa} + \lambda\mathbf{I})^{-1} \boldsymbol{\kappa}(\mathbf{x}')} \quad (54)$$

# Derivation of Kernel Ridge Regression (Ctd.)

In equation (54) we have set:

$$\Phi \varphi(\mathbf{x}') = \begin{pmatrix} \varphi(\mathbf{x}^1)^\top \\ \varphi(\mathbf{x}^2)^\top \\ \vdots \\ \varphi(\mathbf{x}^N)^\top \end{pmatrix} \varphi(\mathbf{x}') = \begin{pmatrix} \varphi(\mathbf{x}^1)^\top \varphi(\mathbf{x}') \\ \varphi(\mathbf{x}^2)^\top \varphi(\mathbf{x}') \\ \vdots \\ \varphi(\mathbf{x}^N)^\top \varphi(\mathbf{x}') \end{pmatrix} = \begin{pmatrix} k(\mathbf{x}^1, \mathbf{x}') \\ k(\mathbf{x}^2, \mathbf{x}') \\ \vdots \\ k(\mathbf{x}^N, \mathbf{x}') \end{pmatrix} =: \kappa(\mathbf{x}')$$

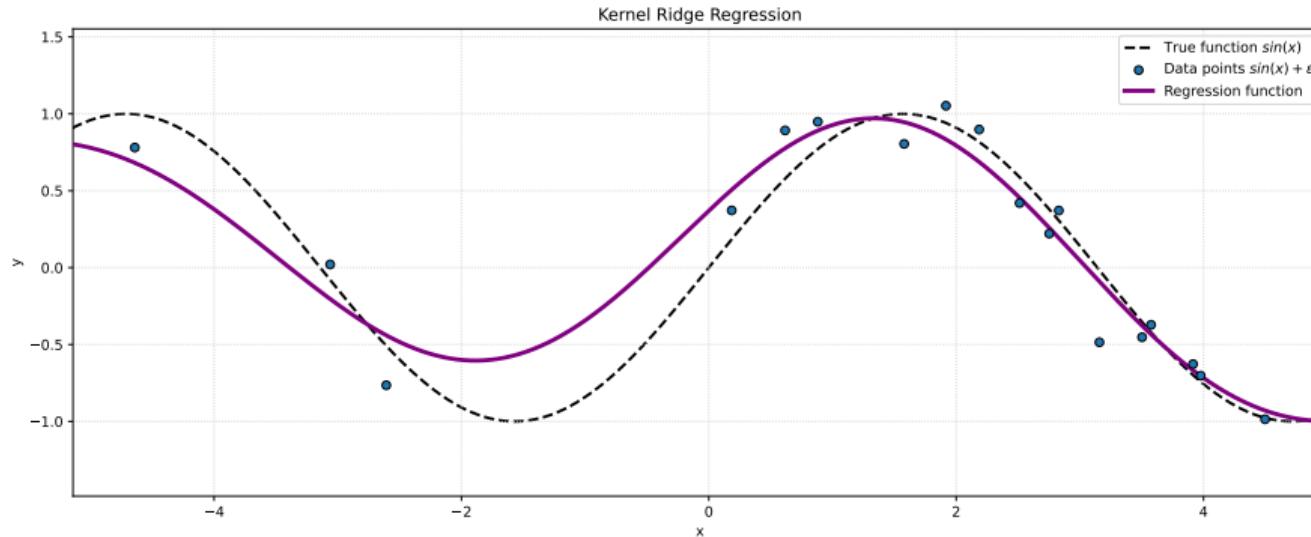
The vector  $\kappa(\mathbf{x}')$  contains the pairwise similarities of  $\mathbf{x}'$  and all training data points,  
i. e. **it must be computed from scratch each time we want to do inference!**

## Derivation of Kernel Ridge Regression (Ctd.)

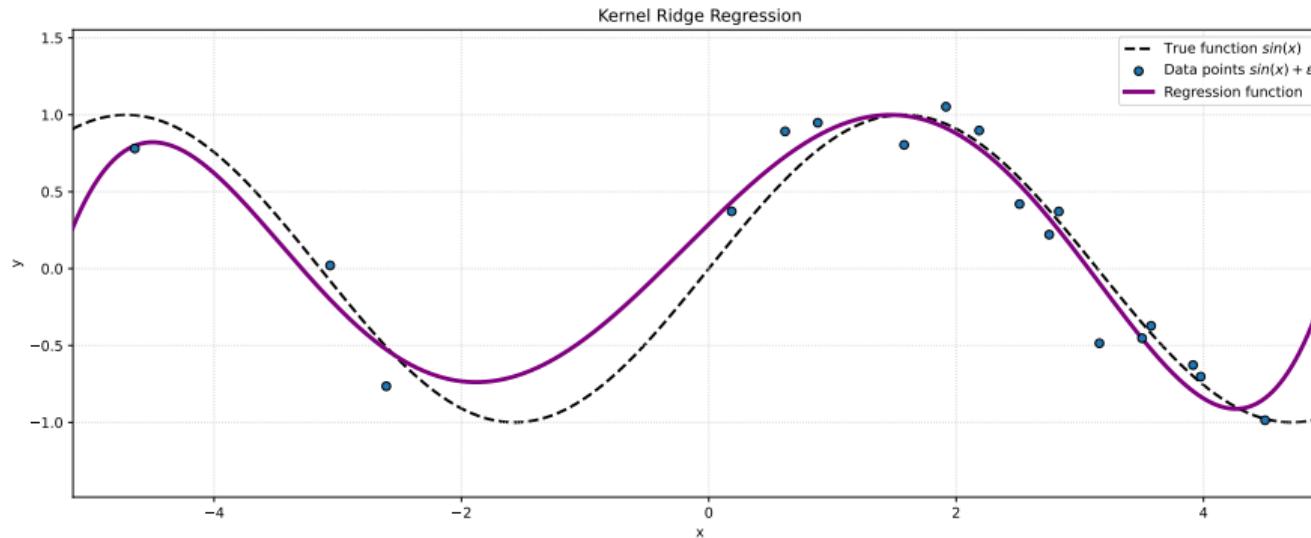
$$\begin{aligned}\boldsymbol{\Phi} \boldsymbol{\Phi}^T &= \begin{pmatrix} \varphi(\mathbf{x}^1)^T \\ \vdots \\ \varphi(\mathbf{x}^N)^T \end{pmatrix} \begin{pmatrix} \varphi(\mathbf{x}^1)^T \\ \vdots \\ \varphi(\mathbf{x}^N)^T \end{pmatrix}^T = \begin{pmatrix} \varphi(\mathbf{x}^1)^T \varphi(\mathbf{x}^1) & \dots & \varphi(\mathbf{x}^1)^T \varphi(\mathbf{x}^N) \\ \vdots & \ddots & \vdots \\ \varphi(\mathbf{x}^N)^T \varphi(\mathbf{x}^1) & \dots & \varphi(\mathbf{x}^N)^T \varphi(\mathbf{x}^N) \end{pmatrix} \\ &= \begin{pmatrix} k(\mathbf{x}^1, \mathbf{x}^1) & \dots & k(\mathbf{x}^1, \mathbf{x}^N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}^N, \mathbf{x}^1) & \dots & k(\mathbf{x}^N, \mathbf{x}^N) \end{pmatrix} =: \mathbf{K} \quad (55)\end{aligned}$$

Only the training data is needed to compute  $\mathbf{K}$ . We can do it once and reuse it!

# Example: Kernel Ridge Regression – Gaussian Kernel, $s := 2.0$



# Example: Kernel Ridge Regression – Polynomial Kernel, $p := 5$



## Section: **GAUSSIAN PROCESS REGRESSION**

- Introduction
- Predictions using noise-free Observations
- Predictions using noisy Observations
- Learning the Hyperparameters

## So far...

**So far:** In BAYESIAN linear regression we focused on **parametric representations** of the regression function  $h$ , i. e. we introduced a set of adjustable parameters  $\theta$  and assumed  $h(\mathbf{x}) = \theta^\top \mathbf{x}$  to be of a fixed parametric form.

We then inferred a posterior distribution  $p(\theta|\mathcal{D})$  over the parameters given the training data. We could use the posterior to predict the label of unseen instances.

# What is a GAUSSIAN Process?

**Now:** A **GAUSSIAN process (GP)** model extends on the **non-parametric** kernel-approach used in kernel ridge regression and performs **BAYESIAN inference over functions themselves**. This means that there are no parameters  $\theta$  involved that need to be learned from training data. Similarly to BAYESian regression, GPs produce a probabilistic output.

Formally, a GAUSSIAN process is a collection of random variables, any finite number of which has a **joint GAUSSIAN distribution**.

# GAUSSian Processes for Regression

- In this section we introduce GPs for regression
- Let the **prior** on the regression function  $h$  be a GP, denoted by

$$y = h(\mathbf{x}) \sim \text{GP}\left(m(\mathbf{x}), k(\mathbf{x}, \tilde{\mathbf{x}})\right), \quad (56)$$

where  $m(\mathbf{x})$  is the **mean function** and  $k(\mathbf{x}, \tilde{\mathbf{x}})$  the positive semi-definite kernel function (**covariance function**), i. e.

$$m(\mathbf{x}) = \mathbb{E}\{h(\mathbf{x})\} \quad (57)$$

$$k(\mathbf{x}, \tilde{\mathbf{x}}) = \mathbb{E}\left\{\left(h(\mathbf{x}) - m(\mathbf{x})\right)\left(h(\tilde{\mathbf{x}}) - m(\tilde{\mathbf{x}})\right)^{\top}\right\} \quad (58)$$

## GAUSSIAN Processes for Regression (Ctd.)

For any **finite set** of points  $\mathbf{x}^1, \dots, \mathbf{x}^N$ , this GAUSSIAN process defines a **joint GAUSSIAN** distribution:

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}, \mathbf{K}) \quad (59)$$

where

$$\boldsymbol{\mu} := (m(\mathbf{x}^1), \dots, m(\mathbf{x}^N))^{\top} \quad (60)$$

$$\mathbf{K} := [k_{ij}]_{i,j=1,\dots,N} \quad \text{with} \quad k_{ij} := k(\mathbf{x}^i, \mathbf{x}^j) \quad (61)$$

# Noise-free Observations: Training Setup

- Suppose we observe a training set

$$\mathcal{D} := \left\{ (\mathbf{x}^n, y_n) \right\}_{n=1}^N$$

where we assume  $y_n$  to be a **noise-free** observation

- Given a set of  $N_*$  test features  $\mathbf{X}_*$  we want to predict the respective labels  $\mathbf{y}_*$
- When assuming noise-free observations, a GP will act as an **interpolator** of the training data, i. e. it predicts the training data with **no uncertainty**

# Joint Distribution (noise-free)

By definition the **joint distribution** of training and test data takes the form of a **partitioned Gaussian distribution** – compare with equation (3):

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{y}_* \end{pmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^\top & \mathbf{K}_{**} \end{bmatrix} \right), \quad (62)$$

where:

$$\mathbf{K} := k(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{N \times N}$$

$$\mathbf{K}_* := k(\mathbf{X}, \mathbf{X}_*) \in \mathbb{R}^{N \times N_*}$$

$$\mathbf{K}_{**} := k(\mathbf{X}_*, \mathbf{X}_*) \in \mathbb{R}^{N_* \times N_*}$$

# GP Posterior Distribution (noise-free)

## GP posterior distribution (noise-free observations)

Using rules (5) and (6) for conditioning Gaussians, we get

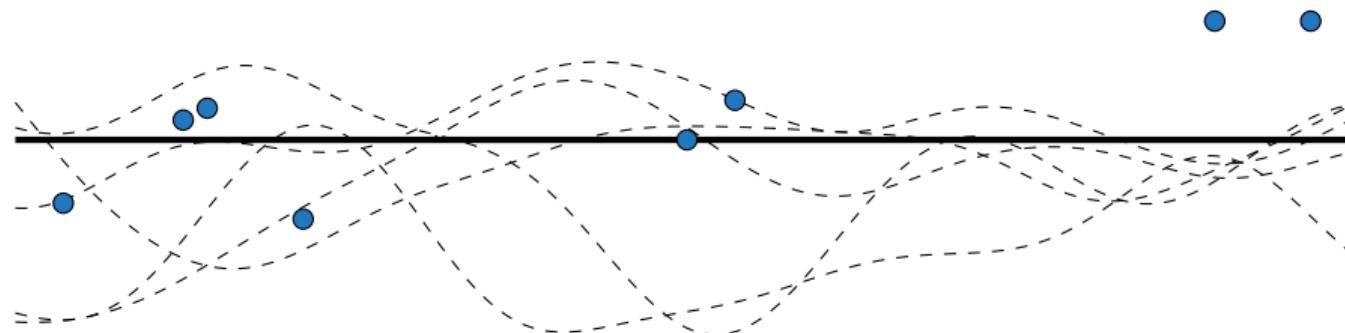
$$p(\mathbf{y}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{y}_*; \boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \quad (63)$$

$$\boldsymbol{\mu}_* := \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{K}_* \quad (64)$$

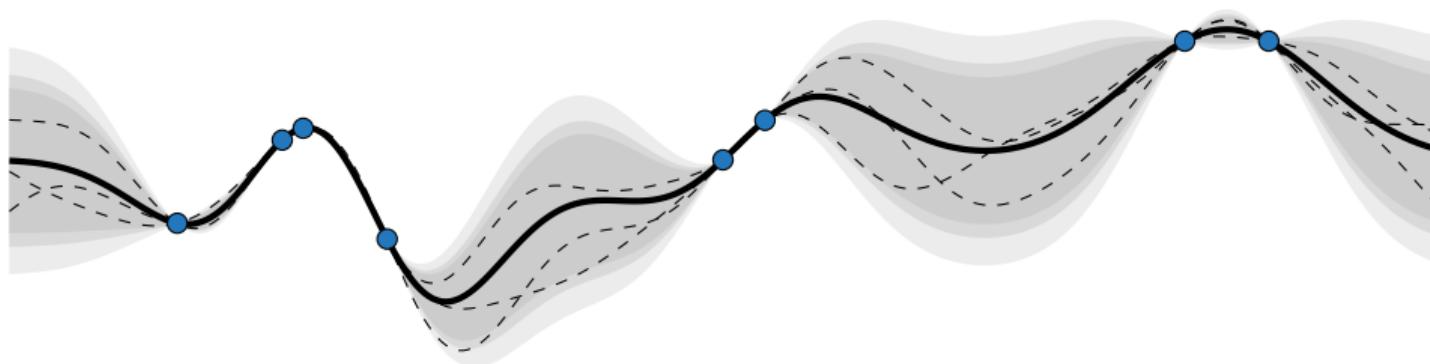
$$\boldsymbol{\Sigma}_* := \mathbf{K}_{**} - \mathbf{K}_*^\top \mathbf{K}^{-1} \mathbf{K}_* \quad (65)$$

**Remark:**  $\boldsymbol{\Sigma}_*$  is called the **SCHUR complement**

# Example: GP Prior Distribution



# Example: GP Posterior Distribution (noise-free)



## Joint Distribution (noisy)

- Now let us consider the case where what we observe is a **noisy version** of the underlying function:  $y = f(\mathbf{x}) + \varepsilon$ , where  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$
- The covariance of the observed noisy data is given by:

$$\tilde{\mathbf{K}} := \mathbf{K} + \sigma^2 \mathbf{I}_N \quad (66)$$

- Thus, the joint GAUSSIAN density becomes:

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{y}_* \end{pmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \tilde{\mathbf{K}} & \mathbf{K}_* \\ \mathbf{K}_*^\top & \mathbf{K}_{**} \end{bmatrix} \right) \quad (67)$$

# GP Posterior Distribution (noisy)

## GP posterior distribution (noisy observations)

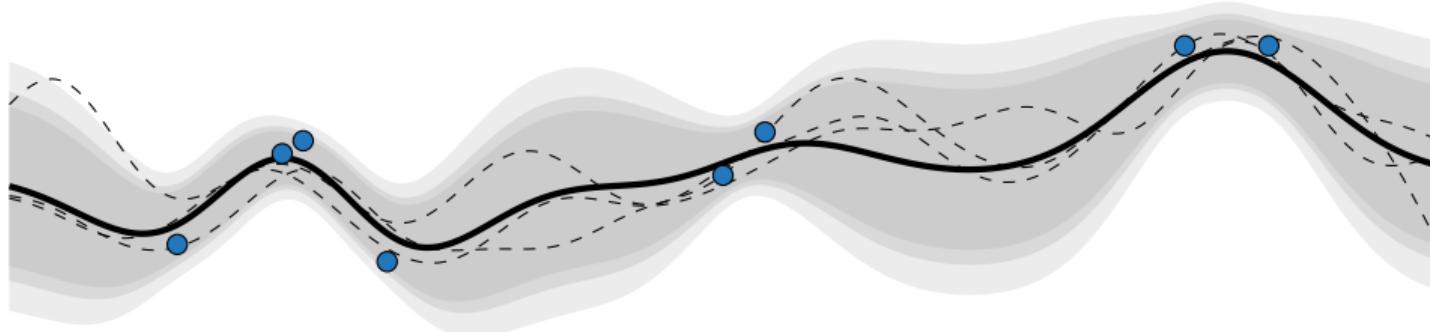
Using rules (5) and (6) for conditioning GAUSSIANS, we get

$$p(\mathbf{y}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*) \quad (68)$$

$$\boldsymbol{\mu}_* := \mathbf{y}^\top \tilde{\mathbf{K}}^{-1} \mathbf{K}_* \quad (69)$$

$$\boldsymbol{\Sigma}_* := \mathbf{K}_{**} - \mathbf{K}_*^\top \tilde{\mathbf{K}}^{-1} \mathbf{K}_* \quad (70)$$

# Example: GP Posterior Distribution (noisy)





# Marginal Likelihood

- To estimate the kernel parameters  $\theta$ , we could use grid search
- Here we consider an **empirical BAYES approach** – in particular we will maximize the **marginal likelihood**:

$$\log p(\mathbf{y}|\mathbf{X}) = \underbrace{-\frac{1}{2}\mathbf{y}^\top \tilde{\mathbf{K}}^{-1}\mathbf{y}}_{\textcircled{1}} - \underbrace{\frac{1}{2}\log \det \tilde{\mathbf{K}}}_{\textcircled{2}} - \underbrace{\frac{N}{2}\log 2\pi}_{\textcircled{3}} \quad (71)$$

- Term **①** is a **data fit term**, term **②** is a **model complexity term**, and the last term **③** is just a constant



# Gradient of the Marginal Likelihood

- To maximize expression (71), we have to compute its gradient
- The partial derivative with respect to the  $j$ -th kernel parameter is given by:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} \log p(\mathbf{y}|\mathbf{X}) &= \frac{1}{2} \mathbf{y}^\top \tilde{\mathbf{K}}^{-1} \frac{\partial \tilde{\mathbf{K}}}{\partial \theta_j} \tilde{\mathbf{K}}^{-1} \mathbf{y} - \frac{1}{2} \text{tr} \left( \tilde{\mathbf{K}}^{-1} \frac{\partial \tilde{\mathbf{K}}}{\partial \theta_j} \right) \\ &= \frac{1}{2} \text{tr} \left( (\boldsymbol{\alpha} \boldsymbol{\alpha}^\top - \tilde{\mathbf{K}}^{-1}) \frac{\partial \tilde{\mathbf{K}}}{\partial \theta_j} \right)\end{aligned}\tag{72}$$

- In equation (72) we have set  $\boldsymbol{\alpha} := \tilde{\mathbf{K}}^{-1} \mathbf{y}$
- The expression  $\frac{\partial \tilde{\mathbf{K}}}{\partial \theta_j}$  depends on the kernel function used



## Some useful Rules

- ① Trace of a matrix  $\mathbf{A} \in \mathbb{R}^{D \times D}$  (*sum of the elements on the main diagonal*):

$$\text{tr}(\mathbf{A}) := \sum_{i=1}^D a_{ii} \quad (73)$$

- ② Derivative of an inverse:

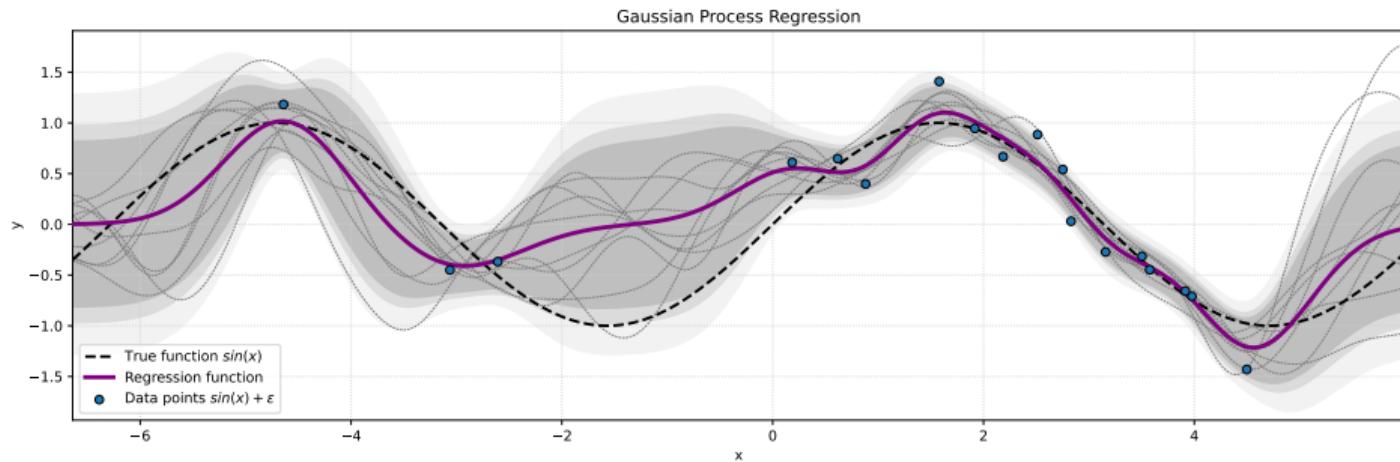
$$\frac{\partial \mathbf{A}^{-1}}{\partial x} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \quad (74)$$

- ③ Derivative of a determinant:

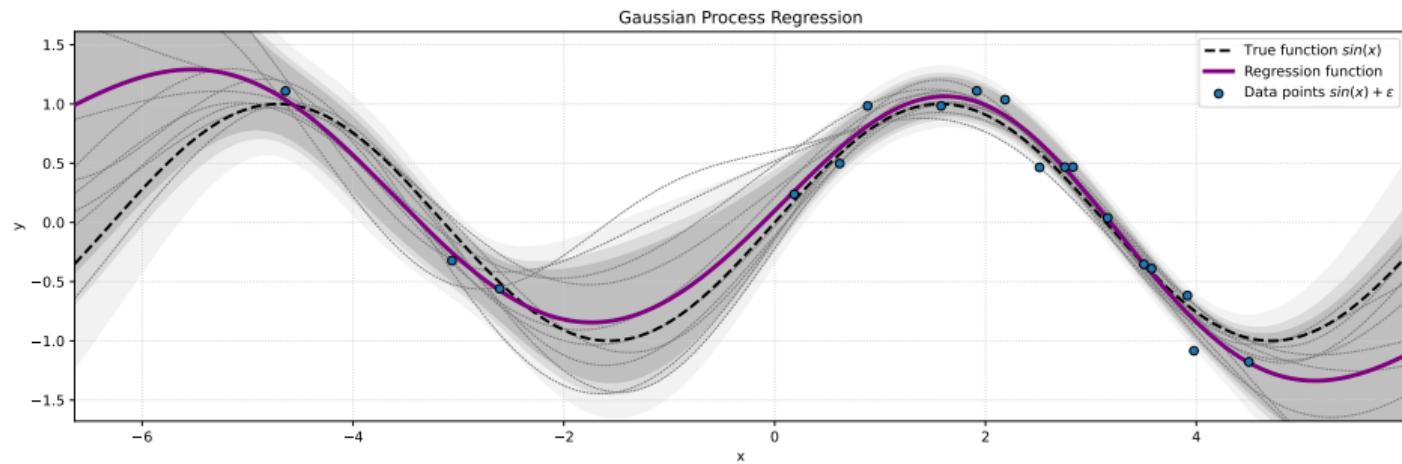
$$\frac{\partial \det \mathbf{A}}{\partial x} = \det(\mathbf{A}) \cdot \text{tr} \left( \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right) \quad (75)$$

- ④  $\mathbf{a}^\top \mathbf{a} = \text{tr}(\mathbf{a}\mathbf{a}^\top)$

# Example: GAUSSIAN Process – not optimized



# Example: GAUSSIAN Process – optimized



## Section: **Wrap-Up**

Summary

Recommended Literature

Self-Test Questions

# Summary

- The maximum likelihood (MLE) estimate is **prone to overfitting**
- We can mitigate this issue by putting a **prior on the parameters** which leads to a maximum a posteriori (MAP) estimate
- In BAYESian regression we sum over all possible models and obtain a model average
- BAYESian regression also models the **local uncertainty**
- Kernel regression uses kernels to get rid of the parameters
- A GAUSSian process is a non-parametric version of BAYESian regression

# Recommended Literature

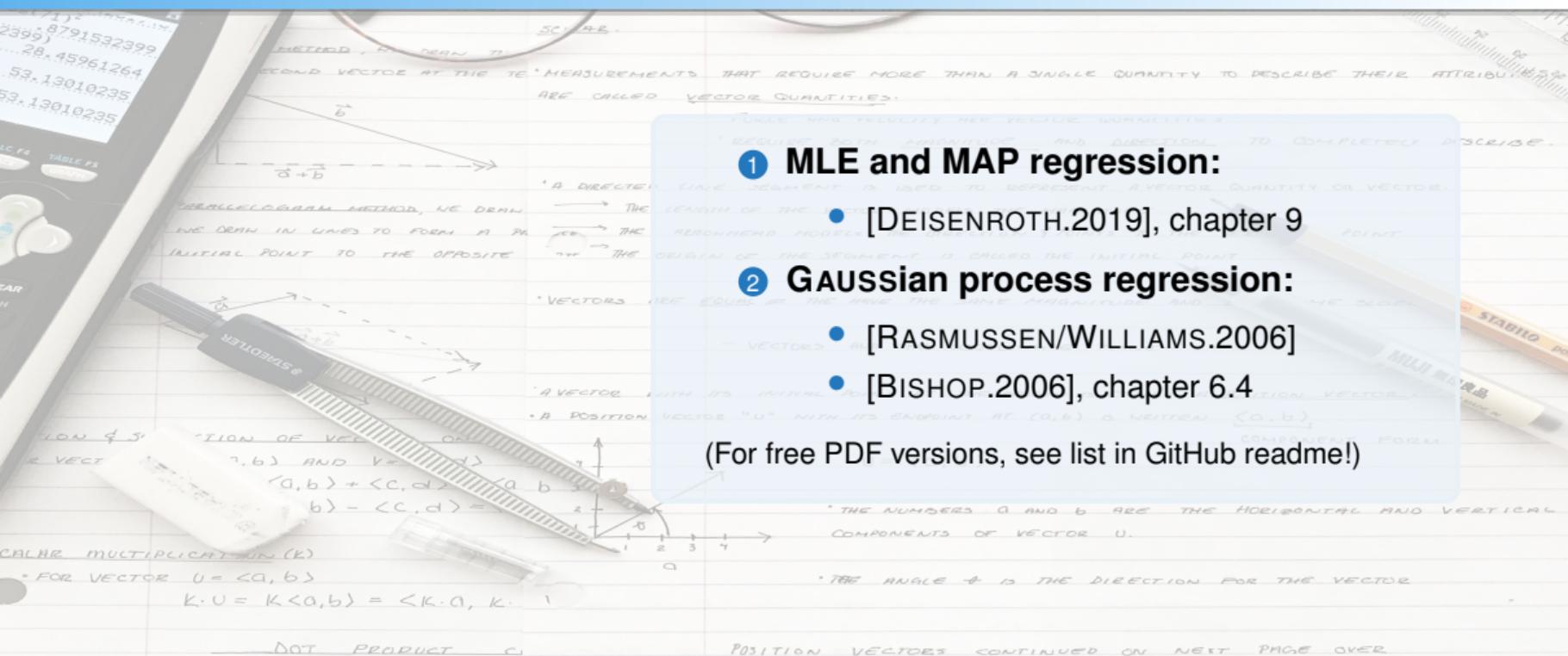
## 1 MLE and MAP regression:

- [DEISENROTH.2019], chapter 9

## 2 Gaussian process regression:

- [RASMUSSEN/WILLIAMS.2006]
- [BISHOP.2006], chapter 6.4

(For free PDF versions, see list in GitHub readme!)



# Self-Test Questions

- ① How is the least squares error related to the MLE estimate?
- ② What is the difference between the MLE estimate and the MAP estimate?
- ③ Explain what BAYESian regression does? How does it work?
- ④ What is the prerequisite for using kernels in regression?
- ⑤ What are advantages / disadvantages of the kernel method?
- ⑥ What is a GAUSSIAN process?
- ⑦ How can we learn suitable hyperparameters?

# Thank you very much for the attention!

\* \* \* Artificial Intelligence and Machine Learning \* \* \*

**Topic:** Advanced Regression Techniques

**Term:** Summer term 2025

**Contact:**

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

[daniel.wehner@sap.com](mailto:daniel.wehner@sap.com)

Do you have any questions?