

Modelling, Uncertainty and Data for Engineers (MUDE)

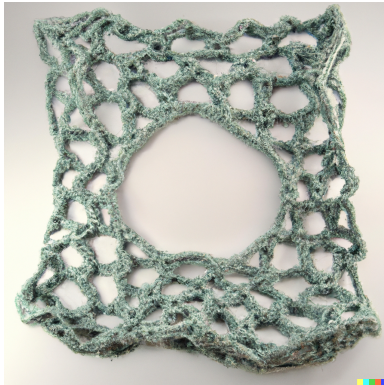
Machine Learning

Iuri Rocha, Riccardo Taormina

and the SLIMM lab and Aidrolab PhDs

A long, long, long time ago (2 years)

A Finite Element model knitted out of wool



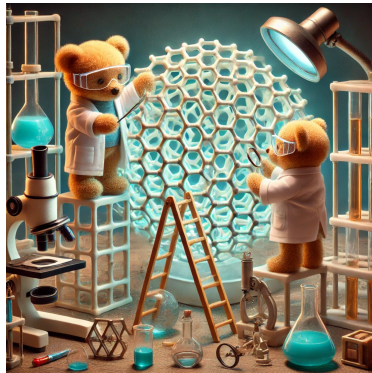
Two teddy bears discovering a new metamaterial



Fairly long ago (2 days)




Here is the creative depiction of a FEM mesh knitted out of wool! It combines the structural precision of FEM with the softness and humor of knitted material, making it an engaging and fun way to teach



Here is a delightful image of two teddy bears hard at work constructing a new metamaterial in a whimsical lab! Their tiny lab coats and safety goggles add an adorable touch while blending humor with a scientific theme.

Nowadays machine learning can do a lot more



Q Search

Ctrl K

Q1 Topics

1. Modelling Concepts

2. Propagation of Uncertainty

3. Observation theory

4. Numerical Modelling

5. Univariate Continuous Distributions

6. Multivariate Distributions

Q2 Topics

1. Finite Volume Method

2. Finite Element Method

2.1. Strong form of the 1D Poisson equation

2.2. From strong to weak form

2.3. From weak to discrete form

2.4. Finite element implementation

2.5. Elements and shape functions

2.6. Numerical integration

2.7. Poisson equation in 2D

2.8. Isoparametric mapping

3. Signal Processing

4. Time Series Analysis

5. Optimization

6. Machine Learning

Programming

Week 1: Poisson Problem

2.1. Strong form of the 1D Poisson equation

The *strong form* of the problem is a Partial Differential Equation (PDE), or a set of PDEs, describing the physics of the phenomena at hand, supplemented with appropriate boundary and initial conditions. Solving the strong form of the problem leads to the exact solution of the continuous system. However, for many engineering problems, it is not possible to obtain the exact solution, and numerical methods need to be employed to find an approximate solution. The finite element method is one approach to find approximate solutions that is applicable to a wide range of different problems. In this chapter we will focus on problems for which the strong form is the Poisson equation.

The story told below and up to the [discrete form](#) section is also presented in this video:

M

Formulating the finite element method

The Poisson equation in 1D

Approximate u as u^h with $u^h(x) = \sum_{i=1}^N N_i(x) u_i = N u$

How to find the values u_i ?

Watch on YouTube

1D Poisson equation

At this point, let's take as an example the one dimensional Poisson equation to illustrate at first the so-called strong form of the problem:

$$-\frac{\partial}{\partial x} \left(\nu(x) \frac{\partial u(x)}{\partial x} \right) = f(x) \quad (2.1)$$

where $\nu(x)$ is a material property, $u(x)$ the solution and $f(x)$ a source, all of which can depend on the coordinate x . To simplify notation we will omit the x -dependence and write ν , u and f in what follows.

4/43

Nowadays machine learning can do a lot more

The screenshot shows a course interface. On the left is a sidebar with a search bar and a table of contents. The main content area displays the title '2.1. Strong form of the 1D Poisson equation' and a paragraph explaining that it is a Partial Differential Equation (PDE). Below the text is a video player with a thumbnail titled 'Formulating the finite element method'. A red arrow points from the video player to a separate box containing the equation $-\frac{\partial}{\partial x} \left(\nu(x) \frac{\partial u(x)}{\partial x} \right) = f(x)$. Below the video player, the text '1D Poisson equation' is followed by another instance of the same equation, labeled (2.1).

Q1 Topics

- 1. Modelling Concepts
- 2. Propagation of Uncertainty
- 3. Observation theory
- 4. Numerical Modelling
- 5. Univariate Continuous Distributions
- 6. Multivariate Distributions

Q2 Topics

- 1. Finite Volume Method
- 2. Finite Element Method
 - 2.1. Strong form of the 1D Poisson equation
 - 2.2. From strong to weak form
 - 2.3. From weak to discrete form
 - 2.4. Finite element implementation
 - 2.5. Elements and shape functions
 - 2.6. Numerical integration
 - 2.7. Poisson equation in 2D
 - 2.8. Isoparametric mapping
- 3. Signal Processing
- 4. Time Series Analysis
- 5. Optimization
- 6. Machine Learning

Programming

2.1. Strong form of the 1D Poisson equation

The *strong form* of the problem is a Partial Differential Equation (PDE), or a set of PDEs, describing the physics of the phenomena at hand, supplemented with appropriate boundary and initial conditions. Solving the strong form of the problem leads to the exact solution of the continuous system. However, for many engineering problems, it is not possible to obtain the exact solution, and numerical methods need to be employed to find an approximate solution. The finite element method is one approach to find approximate solutions that is applicable to a wide range of different problems. In this chapter we will focus on problems for which the strong form is the Poisson equation.

The story told below and up to the [discrete form](#) section is also presented in this video:

Watch on

1D Poisson equation

At this point, let's take as an example the one dimensional Poisson equation to illustrate at first the so-called strong form of the problem:

$$-\frac{\partial}{\partial x} \left(\nu(x) \frac{\partial u(x)}{\partial x} \right) = f(x) \quad (2.1)$$

where $\nu(x)$ is a material property, $u(x)$ the solution and $f(x)$ a source, all of which can depend on the coordinate x . To simplify notation we will omit the x -dependence and write ν , u and f in what follows.

$$-\frac{\partial}{\partial x} \left(\nu(x) \frac{\partial u(x)}{\partial x} \right) = f(x)$$

(screenshot)

Nowadays machine learning can do a lot more

$$-\frac{\partial}{\partial x} \left(\nu(x) \frac{\partial u(x)}{\partial x} \right) = f(x) \quad (2.1)$$

What is this equation?



This equation is a one-dimensional **diffusion equation** with a spatially varying diffusion coefficient $\nu(x)$, often encountered in physics and engineering to describe processes like heat conduction, mass transport, or flow in porous media.

Components:

1. $\nu(x)$:

- A spatially dependent coefficient, representing the diffusivity or conductivity of the medium.

2. $u(x)$:

- The variable of interest, such as temperature, concentration, or velocity field.

3. $f(x)$:

- A source term, representing any external inputs or sinks in the system.

4. $-\frac{\partial}{\partial x} \left(\nu(x) \frac{\partial u(x)}{\partial x} \right)$:

- Represents the flux due to diffusion and the spatial variation of diffusivity.

Nowadays machine learning can do a lot more

$$-\frac{\partial}{\partial x} \left(\nu(x) \frac{\partial u(x)}{\partial x} \right) = f(x) \quad (2.1)$$

What is this equation?



This equation is a one-dimensional **diffusion equation** with a spatially varying diffusion coefficient $\nu(x)$, often encountered in physics and engineering to describe processes like heat conduction, mass transport, or flow in porous media.

Components:

1. $\nu(x)$:
 - A spatially dependent coefficient, representing the diffusivity or conductivity of the medium.
2. $u(x)$:
 - The variable of interest, such as temperature, concentration, or velocity field.
3. $f(x)$:
 - A source term, representing any external inputs or sinks in the system.
4. $-\frac{\partial}{\partial x} \left(\nu(x) \frac{\partial u(x)}{\partial x} \right)$:
 - Represents the flux due to diffusion and the spatial variation of diffusivity.



ChatGPT 4o

This form is particularly useful when the material properties vary with position, requiring a spatially dependent $\nu(x)$.

Can you write me a simple python code to solve this equation with FEM in 1D?



Below is a simple Python code to solve the given 1D diffusion equation using the **Finite Element Method (FEM)**. The example assumes Dirichlet boundary conditions for simplicity.

Fem 1d Diffusion

Answer in chat instead >

This code solves the equation $-\frac{\partial}{\partial x} \left(\nu(x) \frac{\partial u}{\partial x} \right) = f(x)$ using a simple FEM approach in 1D. The solution is plotted at the end. You can modify $\nu(x)$ and $f(x)$ to represent more complex problems. Let me know if you need explanations or enhancements!

🔗 📄 📱 🔄

Message ChatGPT



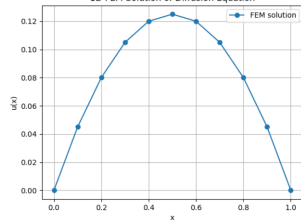
Fem 1d Diffusion

```
29
30 # Element stiffness matrix (2x2)
31 ke = (nu(x[n1]) + nu(x[n2])) / 2 / h * np.array([[1, -1],
32                                                    [-1, 1]])
33
34 # Element load vector (2x1)
35 fe = h / 2 * np.array([f(x[n1]), f(x[n2])])
36
37 # Assemble into the global system
38 K[n1:n2+1, n1:n2+1] += ke
39 F[n1:n2+1] += fe
40
41 # Apply Dirichlet boundary conditions (u(0) = 0, u(L) = 0)
42 bc_nodes = [0, n_nodes - 1] # Boundary nodes
43 for bc_node in bc_nodes:
44     K[bc_node, :] = 0
45     K[:, bc_node] = 0
46     F[bc_node] = 0
```

Console

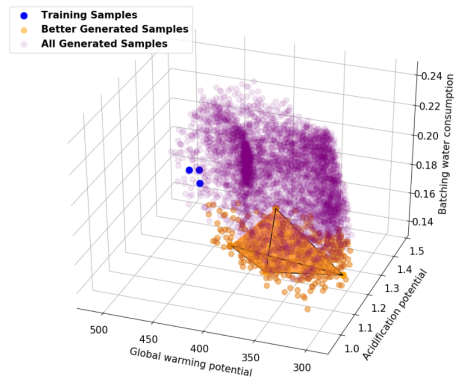
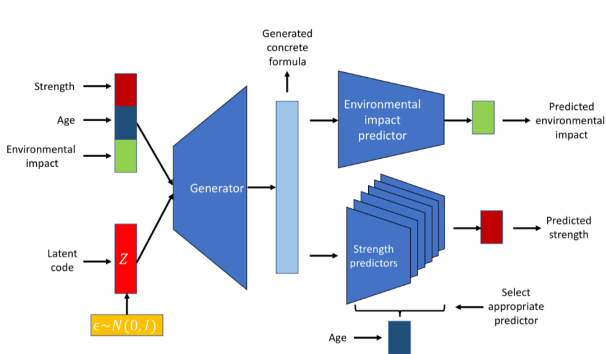
Run

1D FEM Solution of Diffusion Equation



It cannot come up with new ideas, but it can help you do it

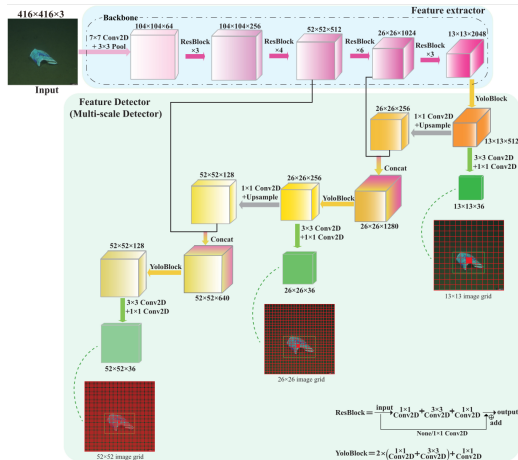
Generating new, more sustainable concrete mixtures:



[Ge et al (2019), arXiv:1905.08222]

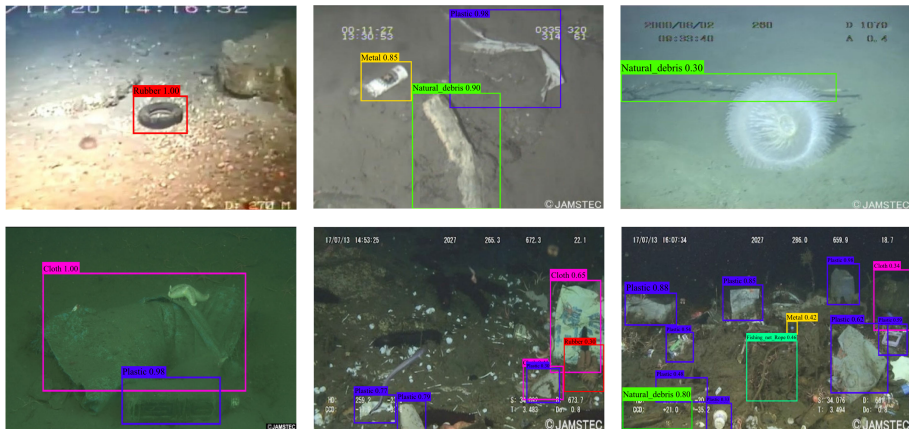
It cannot come up with new ideas, but it can help you do it

Detecting plastic in deep sea with autonomous underwater vehicles:



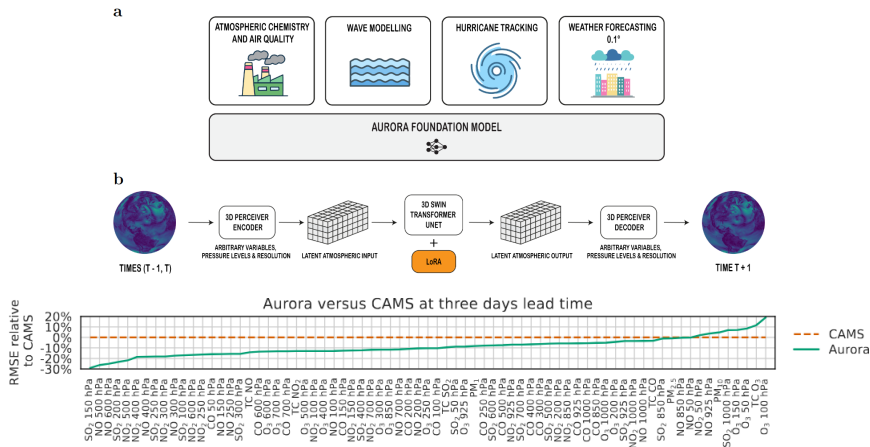
It cannot come up with new ideas, but it can help you do it

Detecting plastic in deep sea with autonomous underwater vehicles:



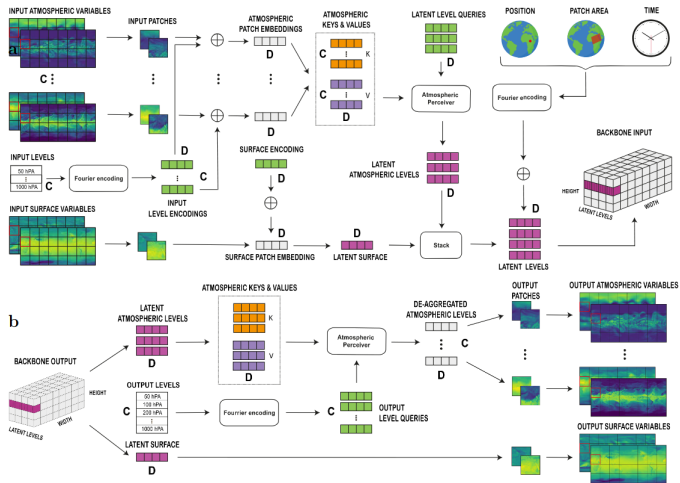
It cannot come up with new ideas, but it can help you do it

A 1.3 billion parameter foundation model for the Earth atmosphere and ocean waves:



It cannot come up with new ideas, but it can help you do it

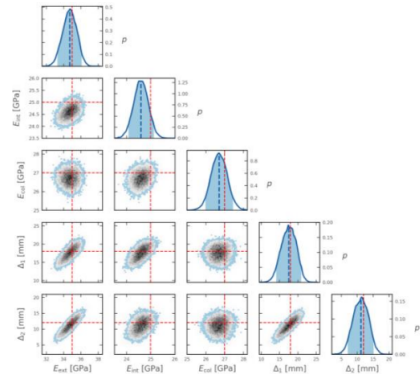
A 1.3 billion parameter foundation model for the Earth atmosphere and ocean waves:



[Bodnar et al (2024), arXiv:2405.13063]

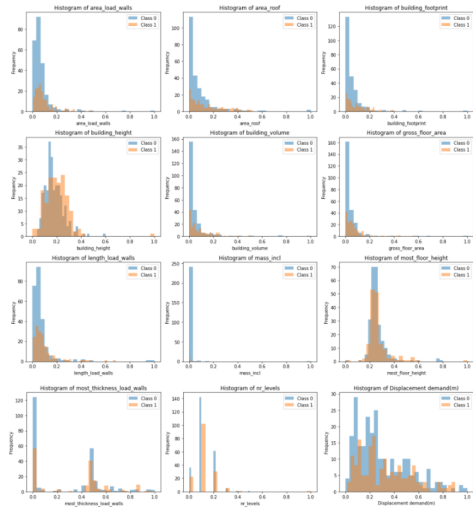
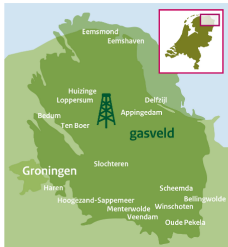
A few examples closer to home

Inverse identification of bridge health through Bayesian machine learning:



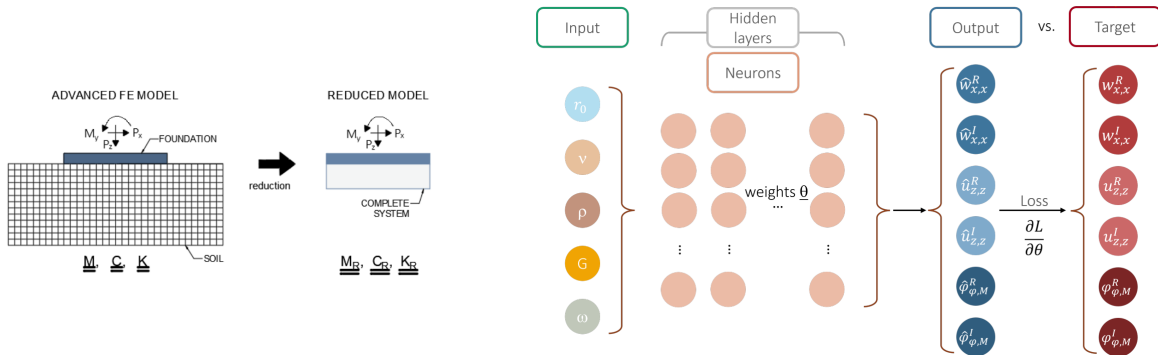
A few examples closer to home

Classifier model for predicting need for reinforcement against earthquakes in Groningen:



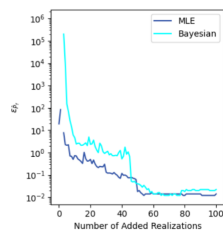
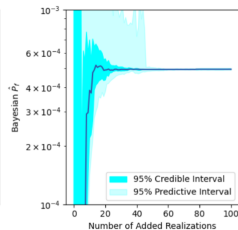
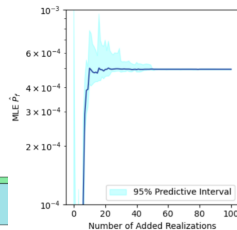
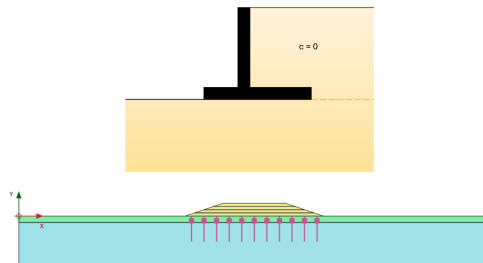
A few examples closer to home

Neural network-based mechanical analogues for soil-structure interaction:



A few examples closer to home

Bayesian machine learning for reliability analysis of geotechnical engineering problems:



Some definitions and terminology

Narrow versus General AI:

- Narrow AI can only perform one specific task \Leftarrow ML techniques live here
- General AI can perform a multitude of tasks and program itself \Leftarrow just a dream (for now)

Some definitions and terminology

Narrow versus General AI:

- Narrow AI can only perform one specific task \Leftarrow ML techniques live here
- General AI can perform a multitude of tasks and program itself \Leftarrow just a dream (for now)

Supervised Learning: Tasks with known target outcomes, requires labeled data:

- Regression: Map input features to noisy observations of continuous targets \Leftarrow this course
- Classification: Map input features to discrete class labels

Some definitions and terminology

Narrow versus General AI:

- Narrow AI can only perform one specific task \Leftarrow ML techniques live here
- General AI can perform a multitude of tasks and program itself \Leftarrow just a dream (for now)

Supervised Learning: Tasks with known target outcomes, requires labeled data:

- Regression: Map input features to noisy observations of continuous targets \Leftarrow this course
- Classification: Map input features to discrete class labels

Unsupervised Learning: Explain patterns in unlabeled data with latent (hidden) variables:

- Clustering: Split data into groups explained by discrete latents
- Dimensionality reduction: Explain the data with a manifold described by continuous latents

Some definitions and terminology

Narrow versus General AI:

- Narrow AI can only perform one specific task \Leftarrow ML techniques live here
- General AI can perform a multitude of tasks and program itself \Leftarrow just a dream (for now)

Supervised Learning: Tasks with known target outcomes, requires labeled data:

- Regression: Map input features to noisy observations of continuous targets \Leftarrow this course
- Classification: Map input features to discrete class labels

Unsupervised Learning: Explain patterns in unlabeled data with latent (hidden) variables:

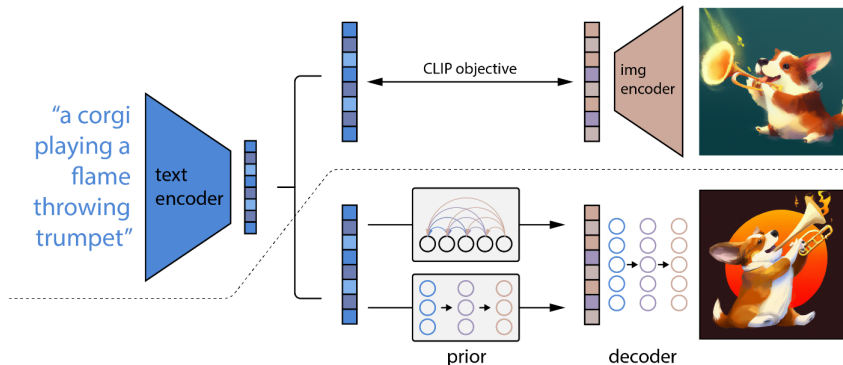
- Clustering: Split data into groups explained by discrete latents
- Dimensionality reduction: Explain the data with a manifold described by continuous latents

Reinforcement Learning: Learn a task through reward/punishment mechanisms:

- Agent(s) interacting with an environment, evolving interaction policy

Finding patterns and making good use of them

The core goal of ML is not fitting data, but finding useful representations and exploiting them



Week 2.6 – Theory

You can find all the material in the book, as always:

- Pages for each important concept, starting from scratch
- It is beneficial to go through them in order

Interactive plots:

- You will play with some interactive plots today
- Keep going in the book, good way to build ML intuition

Videos:

- Each page comes with a short video, highly recommended to watch
- In total about 40 minutes of videos for regression modeling

Quiz questions:

- In each page with hidden answer blocks
- At the end of the book section on a dedicated page – **these are exam-like questions**

Week 2.6 – PA and WS

Programming assignment:

- Split a dataset into training/validation/test blocks
- Crucial operation when training ML models!

Wednesday workshop:

- Build your first neural networks with `scikit-learn`
- Experiment with data normalization, overfitting, underfitting and model selection
- Revisit the road deformation project from Week 1.3, now with neural nets

Code from both assignments will be very helpful for Friday

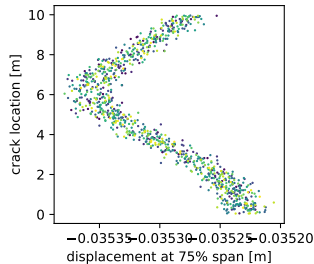
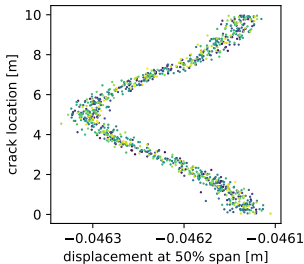
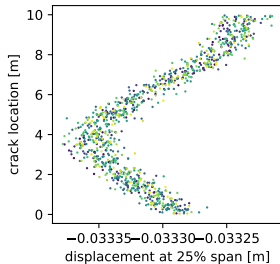
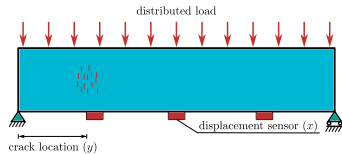
Week 2.6 – Friday project

Detecting cracks in bridges with neural networks:

- Beams with a crack somewhere along the span
- We do not know the crack location but we do know how the beam deforms
- We train a net with a dataset of 800 beams and predict 200 unseen ones

A complete ML regression workflow:

- Pre-process the data and get it ready for training
- Experiment with feature selection, starting with one sensor at midspan
- Perform a well-structured model selection procedure to pick the best architecture



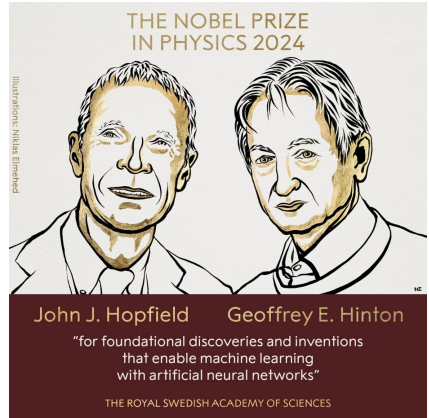
Objectives

Contents:

- Decision theory for regression, k-Nearest Neighbors estimator
- Linear regression with nonlinear basis functions
- Introduction to neural networks for regression

By the end of the week, you will be able to:

- Compare different regression modeling approaches
- Construct parsimonious regression models
- Critically assess model performance



Objectives

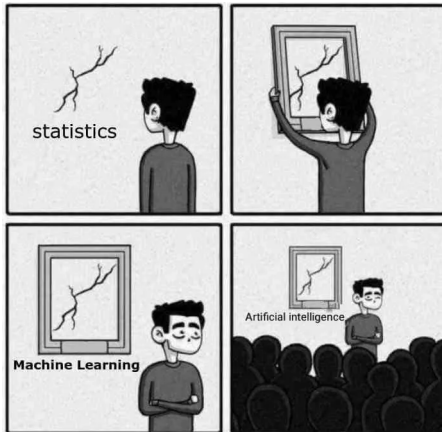
Contents:

- Decision theory for regression, k-Nearest Neighbors estimator
- Linear regression with nonlinear basis functions
- Introduction to neural networks for regression

By the end of the week, you will be able to:

- Compare different regression modeling approaches
- Construct parsimonious regression models
- Critically assess model performance

Brace yourself for some statistics:



Quick statistics recap

Continuous random variables represented by **probability densities**:

- Joint, marginal and conditional densities

$$p(x, y) \quad p(x) \quad p(y) \quad p(x|y) \quad p(y|x)$$

- For two independent variables x and y it holds

$$p(x, y) = p(x)p(y)$$

Quick statistics recap

Continuous random variables represented by **probability densities**:

- Joint, marginal and conditional densities

$$p(x, y) \quad p(x) \quad p(y) \quad p(x|y) \quad p(y|x)$$

- For two independent variables x and y it holds

$$p(x, y) = p(x)p(y)$$

Some useful integrals

- **Expectation** of a function of a random variable

$$\mathbb{E}[f(x)] = \int f(x)p(x) \, dx$$

- **Monte Carlo** approximation of the expectation, with N samples x_i from $p(x)$:

$$\mathbb{E}[f(x)] \approx \frac{1}{N} \sum_i^N f(x_i)$$

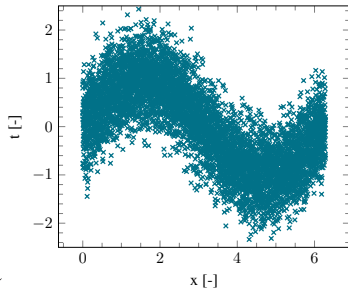
- **Variance** of a function of a random variable

$$\text{var}[f(x)] = \mathbb{E} \left[(f(x) - \mathbb{E}[f(x)])^2 \right]$$

Regression problems

The problem we would like to solve:

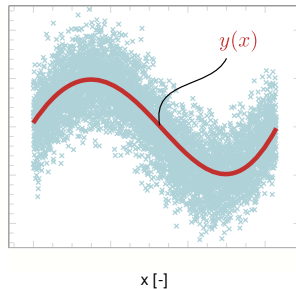
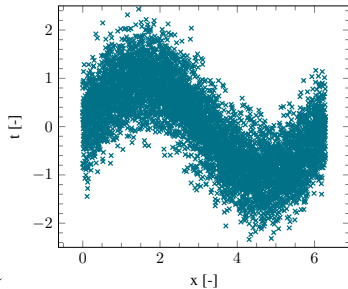
- **Given:** Some complex process $p(x, t)$, usually highly nonlinear



Regression problems

The problem we would like to solve:

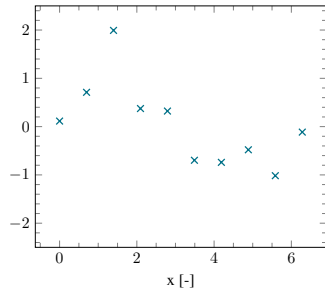
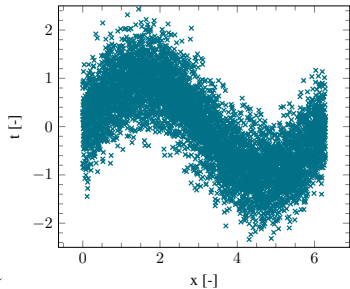
- **Given:** Some complex process $p(x, t)$, usually highly nonlinear
- **Goal:** Construct a model $y(x)$ that is **as close as possible** to t



Regression problems

The problem we would like to solve:

- **Given:** Some complex process $p(\mathbf{x}, t)$, usually highly nonlinear
- **Goal:** Construct a model $y(\mathbf{x})$ that is **as close as possible** to t
- **In practice:** We do not know $p(\mathbf{x}, t)$, but only have N observations of it

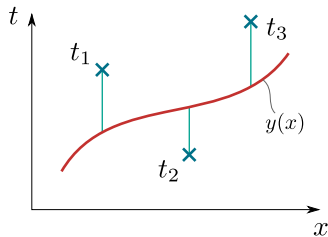
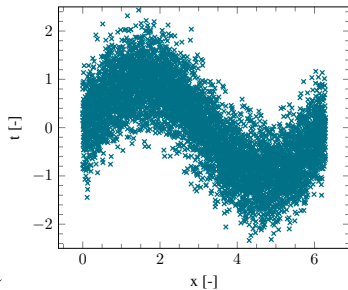


Regression problems

The problem we would like to solve:

- **Given:** Some complex process $p(\mathbf{x}, t)$, usually highly nonlinear
- **Goal:** Construct a model $y(\mathbf{x})$ that is **as close as possible** to t
- **In practice:** We do not know $p(\mathbf{x}, t)$, but only have N observations of it
- How to measure this "closeness"? The **squared loss function** is a popular choice:

$$L(t, y(\mathbf{x})) = (y(\mathbf{x}) - t)^2$$



Regression problems

The problem we would like to solve:

- **Given:** Some complex process $p(\mathbf{x}, t)$, usually highly nonlinear
- **Goal:** Construct a model $y(\mathbf{x})$ that is **as close as possible** to t
- **In practice:** We do not know $p(\mathbf{x}, t)$, but only have N observations of it
- How to measure this "closeness"? The **squared loss function** is a popular choice:

$$L(t, y(\mathbf{x})) = (y(\mathbf{x}) - t)^2$$

- Here it is natural to go for the expectation:

$$\mathbb{E}[L] = \int \int (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

Regression problems

The problem we would like to solve:

- **Given:** Some complex process $p(\mathbf{x}, t)$, usually highly nonlinear
- **Goal:** Construct a model $y(\mathbf{x})$ that is **as close as possible** to t
- **In practice:** We do not know $p(\mathbf{x}, t)$, but only have N observations of it
- How to measure this "closeness"? The **squared loss function** is a popular choice:

$$L(t, y(\mathbf{x})) = (y(\mathbf{x}) - t)^2$$

- Here it is natural to go for the expectation:

$$\mathbb{E}[L] = \int \int (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

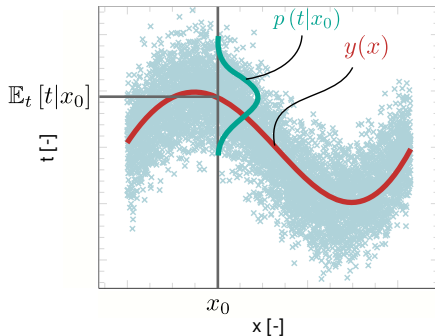
- Solving for the **regression function** $y(\mathbf{x})$ gives:

$$y(\mathbf{x}) = \int t p(t|x) \, dt = \mathbb{E}_t[t|\mathbf{x}]$$

Regression problems

The problem we would like to solve:

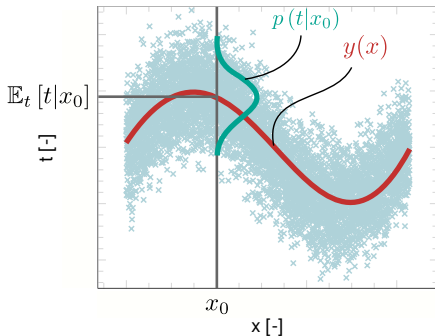
- **Given:** Some complex process $p(x, t)$, usually highly nonlinear
- **Goal:** Construct a model $y(x)$ that is **as close as possible** to t



Regression problems

The problem we would like to solve:

- **Given:** Some complex process $p(x, t)$, usually highly nonlinear
- **Goal:** Construct a model $y(x)$ that is **as close as possible** to t



- We can see this as an **"ideal model"** $h(x) = \mathbb{E}_t[t|x]$, but it requires full knowledge of $p(x, t)$ (i.e. infinite data)

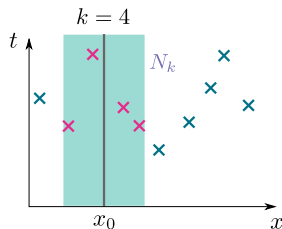
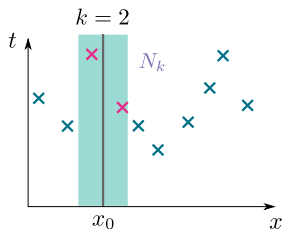
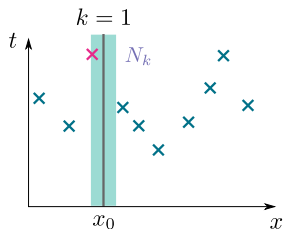
Reaching a model $y(x)$ with local approximations

Let us start with a very simple k -Nearest Neighbors (kNN) model:

Reaching a model $y(x)$ with local approximations

Let us start with a very simple **k -Nearest Neighbors (kNN) model**:

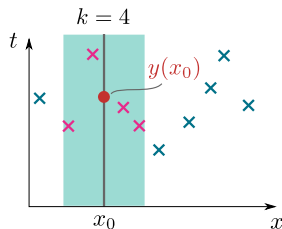
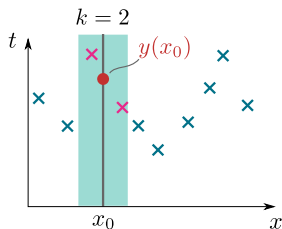
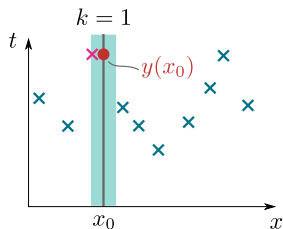
- We look at a neighborhood N_k around x_0 until we find k points \Rightarrow approximate conditioning $\mathbb{E}[t|x]$



Reaching a model $y(x)$ with local approximations

Let us start with a very simple **k -Nearest Neighbors (kNN) model**:

- We look at a neighborhood N_k around x_0 until we find k points \Rightarrow approximate conditioning $\mathbb{E}[t|x]$
- We then average these points to obtain $y(x_0) \Rightarrow$ approximate expectation $\mathbb{E}[t|x]$

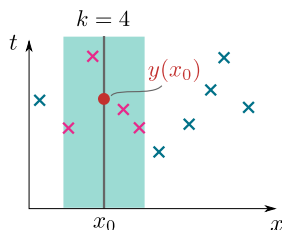
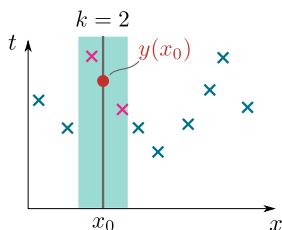
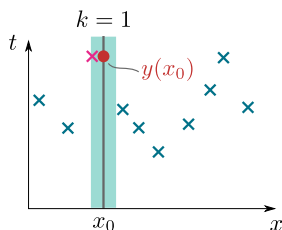


$$y(x_0) = \frac{1}{k} \sum_{x_i \in N_k} t_i$$

Reaching a model $y(x)$ with local approximations

Let us start with a very simple **k -Nearest Neighbors (kNN) model**:

- We look at a neighborhood N_k around x_0 until we find k points \Rightarrow approximate conditioning $\mathbb{E}[t|x]$
- We then average these points to obtain $y(x_0) \Rightarrow$ approximate expectation $\mathbb{E}[t|x]$



$$y(x_0) = \frac{1}{k} \sum_{x_i \in N_k} t_i$$

Now we just need to pick k by minimizing the loss:

- Since we only have N data points we minimize a Monte Carlo approximation:

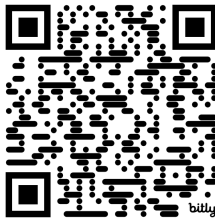
$$\int \int (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt \approx \frac{1}{N} \sum_i^N (y(\mathbf{x}_i) - t_i)^2$$

Now let us try this out

Go to bit.ly/engmechml or scan the QR code:

- Look at **the first** interactive plot
- Change the value of k until you are satisfied with the model
- Change the value of k until the training loss is as small as possible:

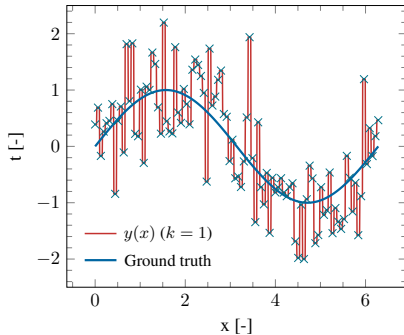
$$\mathbb{E}[L] \approx \frac{1}{N} \sum_i^N (y(\mathbf{x}_i, k) - t_i)^2$$



Overfitting and underfitting

This is the model we get if we are just trying to minimize the training loss:

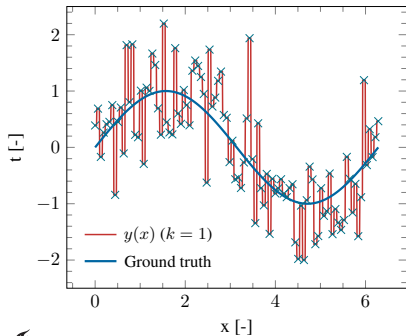
- Model fits the noise in the dataset and cannot generalize
- The error is exactly zero, but this is not a good model



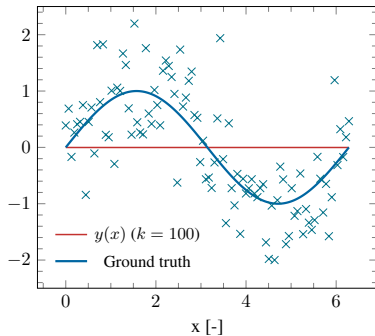
Overfitting and underfitting

This is the model we get if we are just trying to minimize the training loss:

- Model fits the noise in the dataset and cannot generalize
- The error is exactly zero, but this is not a good model
- Too much freedom? What if we increase k ?



increasing k
(a bit too much)



Model selection

Clearly, choosing k is tricky:

- Too low: we fit the noise in the data \Rightarrow **overfitting!**
- Too high: we oversmooth and lose detail \Rightarrow **underfitting!**
- The training set cannot be trusted to give us k , it will always lead to $k = 1$

Model selection

Clearly, choosing k is tricky:

- Too low: we fit the noise in the data \Rightarrow **overfitting!**
- Too high: we oversmooth and lose detail \Rightarrow **underfitting!**
- The training set cannot be trusted to give us k , it will always lead to $k = 1$

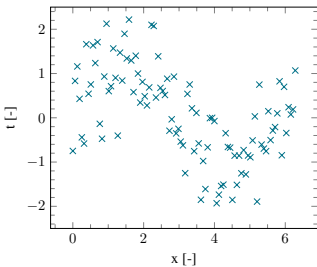
The solution is to introduce a **validation dataset**:

- A new dataset that cannot be used for training
- We can then use it to find the **hyperparameter** k :

$$k = \arg \min_{\bar{k}} \frac{1}{N_{\text{val}}} \sum_i^{N_{\text{val}}} \left(y(\mathbf{x}_i, \bar{k}) - t_i \right)^2$$

Model selection

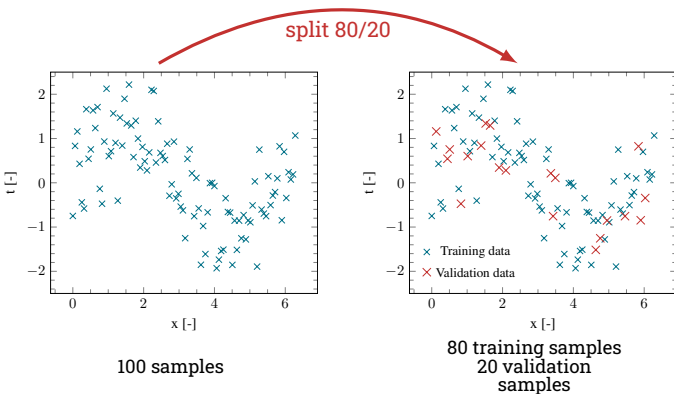
But how do we pick a validation set?



100 samples

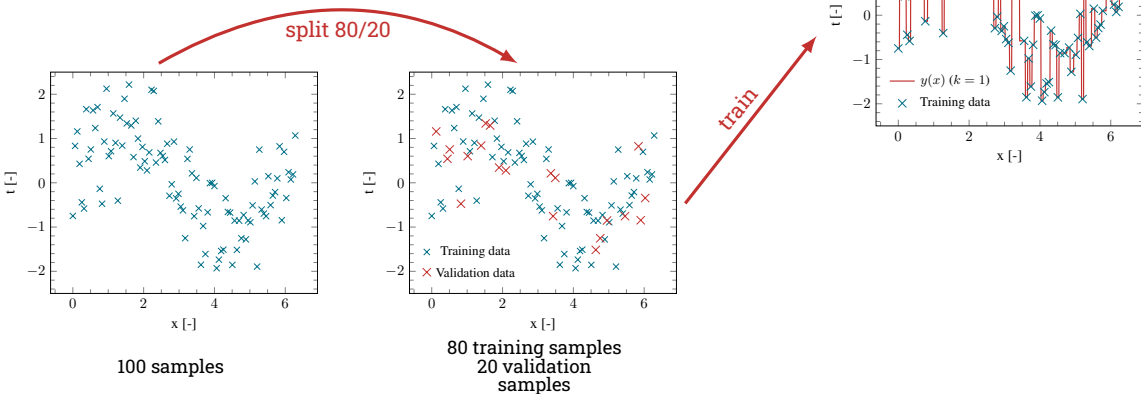
Model selection

But how do we pick a validation set?



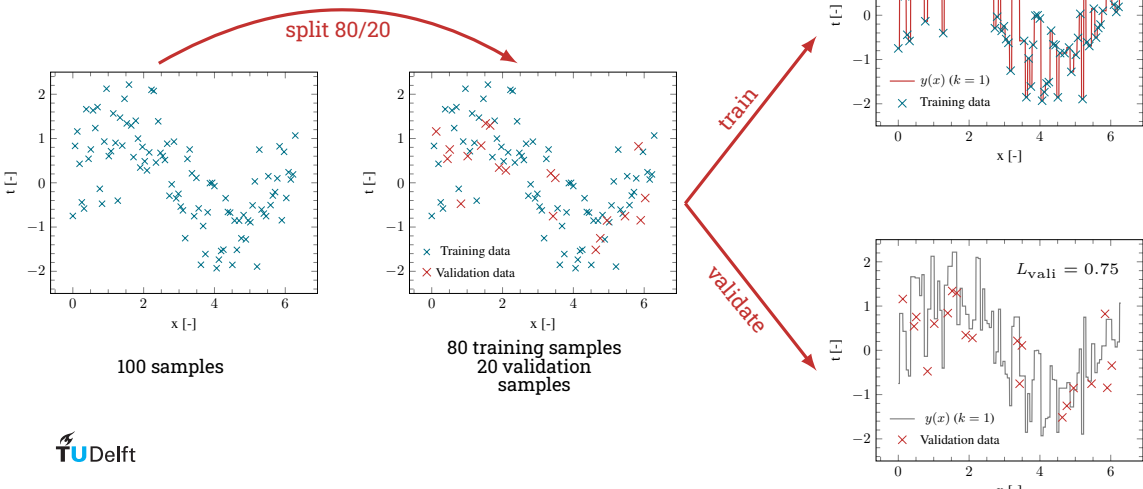
Model selection

But how do we pick a validation set?



Model selection

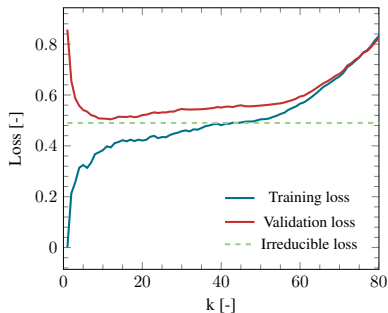
But how do we pick a validation set?



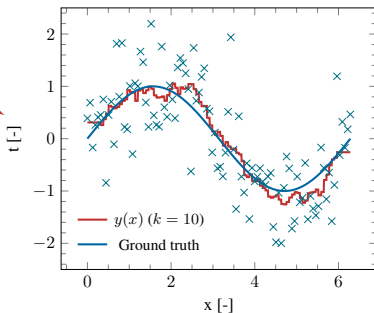
Model selection

The bias-variance tradeoff:

- Overly flexible models have **low bias** and **high variance**
- Overly rigid models have **high bias** and **low variance**
- We may accept some bias in exchange for a lower variance... but not too much



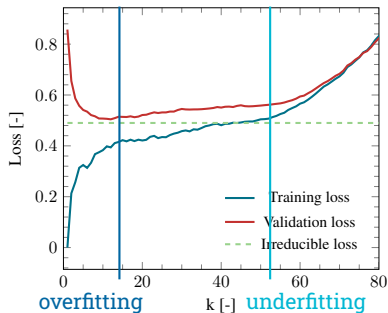
pick k
at minimum loss



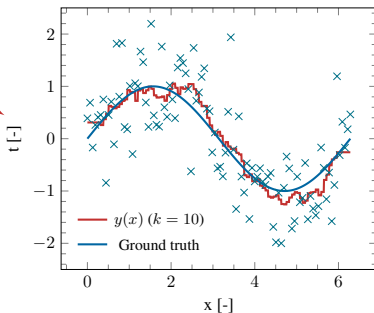
Model selection

The bias-variance tradeoff:

- Overly flexible models have **low bias** and **high variance**
- Overly rigid models have **high bias** and **low variance**
- We may accept some bias in exchange for a lower variance... but not too much



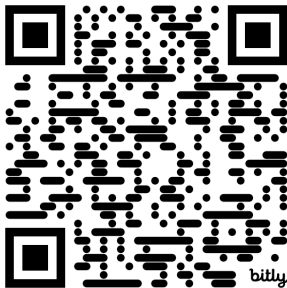
pick k
at minimum loss



Let us do it one last time

Go to bit.ly/engmechml or scan the QR code:

- Look at **the third** interactive plot
- Change the value of k until the validation loss is as low as possible

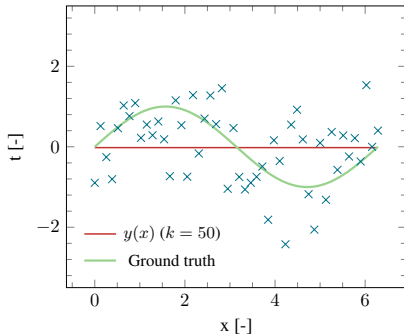
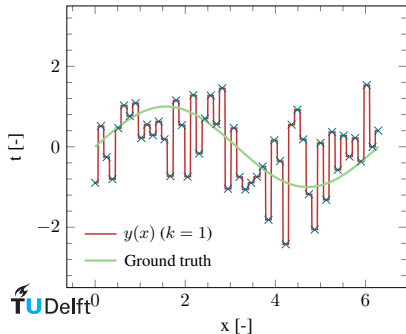


The bias-variance tradeoff

Why do we say flexible models have high variance? A closer look:

- Same example as before, but now **1000 different datasets** of $N = 50$ each
- How much does the choice of dataset affect the final model?

$$L(\mathbf{x}) = \underbrace{(\mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})] - h(\mathbf{x}))^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [(y(\mathbf{x}, \mathcal{D}) - \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})])^2]}_{\text{variance}} + \underbrace{\int (h(\mathbf{x}) - t)^2 p(t|\mathbf{x}) dt}_{\text{irreducible noise}}$$

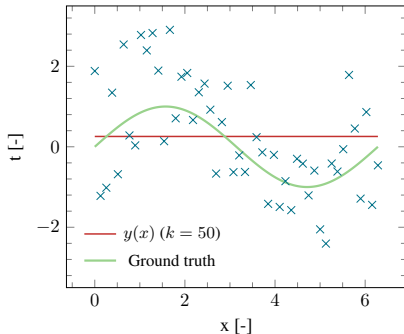
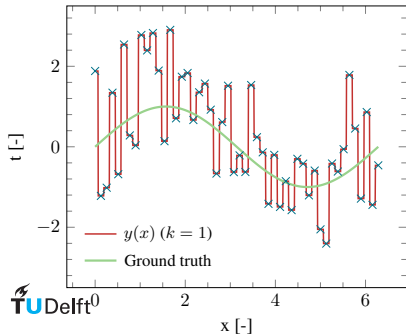


The bias-variance tradeoff

Why do we say flexible models have high variance? A closer look:

- Same example as before, but now **1000 different datasets** of $N = 50$ each
- How much does the choice of dataset affect the final model?

$$L(\mathbf{x}) = \underbrace{(\mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})] - h(\mathbf{x}))^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [(y(\mathbf{x}, \mathcal{D}) - \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})])^2]}_{\text{variance}} + \underbrace{\int (h(\mathbf{x}) - t)^2 p(t|\mathbf{x}) dt}_{\text{irreducible noise}}$$

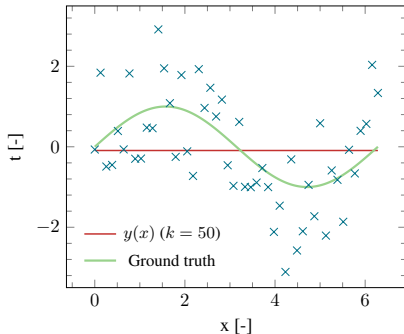
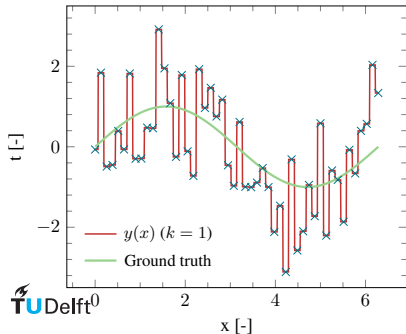


The bias-variance tradeoff

Why do we say flexible models have high variance? A closer look:

- Same example as before, but now **1000 different datasets** of $N = 50$ each
- How much does the choice of dataset affect the final model?

$$L(\mathbf{x}) = \underbrace{(\mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})] - h(\mathbf{x}))^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [(y(\mathbf{x}, \mathcal{D}) - \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})])^2]}_{\text{variance}} + \underbrace{\int (h(\mathbf{x}) - t)^2 p(t|\mathbf{x}) dt}_{\text{irreducible noise}}$$

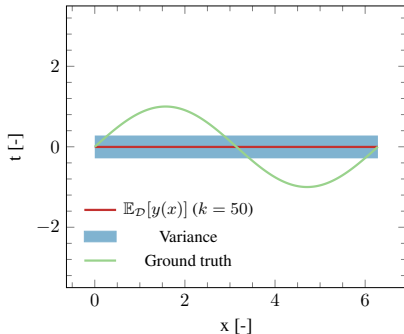
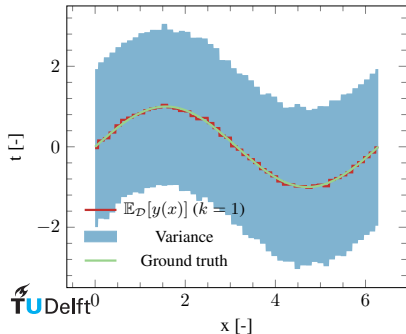


The bias-variance tradeoff

Why do we say flexible models have high variance? A closer look:

- Same example as before, but now **1000 different datasets** of $N = 50$ each
- How much does the choice of dataset affect the final model?

$$L(\mathbf{x}) = \underbrace{(\mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})] - h(\mathbf{x}))^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [(y(\mathbf{x}, \mathcal{D}) - \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})])^2]}_{\text{variance}} + \underbrace{\int (h(\mathbf{x}) - t)^2 p(t|\mathbf{x}) dt}_{\text{irreducible noise}}$$

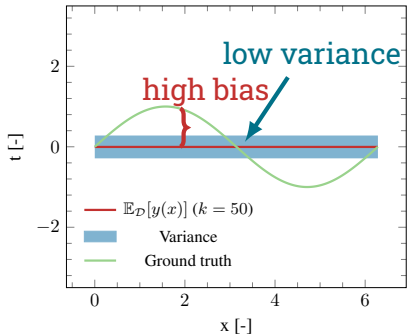
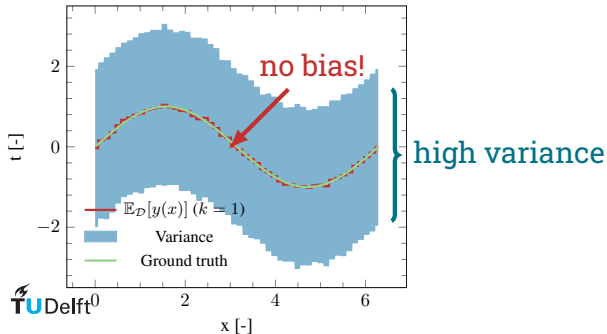


The bias-variance tradeoff

Why do we say flexible models have high variance? A closer look:

- Same example as before, but now **1000 different datasets** of $N = 50$ each
- How much does the choice of dataset affect the final model?

$$L(\mathbf{x}) = \underbrace{(\mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})] - h(\mathbf{x}))^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [(y(\mathbf{x}, \mathcal{D}) - \mathbb{E}_{\mathcal{D}} [y(\mathbf{x}, \mathcal{D})])^2]}_{\text{variance}} + \underbrace{\int (h(\mathbf{x}) - t)^2 p(t|\mathbf{x}) dt}_{\text{irreducible noise}}$$



Reaching a model $y(x)$ with a global approximation

Observation model:

- We adopt a parametric model $y(\mathbf{x}, \mathbf{w})$ and assume some additive Gaussian noise:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

Reaching a model $y(x)$ with a global approximation

Observation model:

- We adopt a parametric model $y(\mathbf{x}, \mathbf{w})$ and assume some additive Gaussian noise:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

- Under the squared loss we have seen before, the **regression function** is simply:

$$\mathbb{E}_t[t | \mathbf{x}] = \int t p(t | \mathbf{x}) dt = y(\mathbf{x}, \mathbf{w})$$

- The learning problem then boils down to fitting a Gaussian distribution

Reaching a model $y(x)$ with a global approximation

Observation model:

- We adopt a parametric model $y(\mathbf{x}, \mathbf{w})$ and assume some additive Gaussian noise:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

- Under the squared loss we have seen before, the **regression function** is simply:

$$\mathbb{E}_t[t | \mathbf{x}] = \int t p(t | \mathbf{x}) dt = y(\mathbf{x}, \mathbf{w})$$

- The learning problem then boils down to fitting a Gaussian distribution

The approximations for this flavor of models are now:

- Noise is Gaussian
- Response is unimodal (because the Gaussian has only one peak)
- The function $y(\mathbf{x}, \mathbf{w})$ might not be infinitely flexible (bias $\neq 0$)

Giving $y(x)$ some shape – linear basis function models

Simple linear regression, assuming D input features in \mathbf{x} :

- **Parametric model**, linear in its arguments:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \cdots + w_D x_D$$

Giving $y(x)$ some shape – linear basis function models

Simple linear regression, assuming D input features in \mathbf{x} :

- **Parametric model**, linear in its arguments:

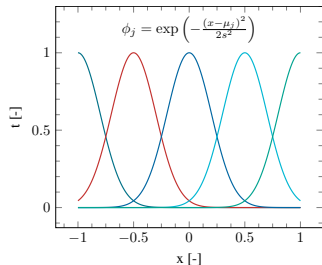
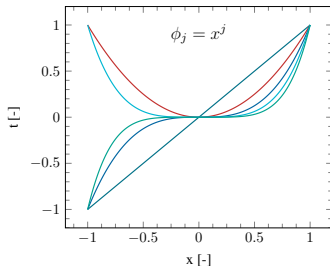
$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \cdots + w_D x_D$$

Here we make them more flexible:

- General **nonlinear** functions of \mathbf{x} as regressors:

$$y(\mathbf{x}, \mathbf{w}) = \sum_j^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

- A **bias** term $\phi_0 = 1$ is usually included in $\boldsymbol{\phi}$
- We are now unshackled from the original dimensionality D



Maximum Likelihood Estimation

Computing the likelihood of our data:

- The probability density of a given value t is:

$$p(t | \mathbf{w}) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

- Given a dataset \mathcal{D} with observations $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} / \mathbf{t} = [t_1, \dots, t_N]$,

Maximum Likelihood Estimation

Computing the likelihood of our data:

- The probability density of a given value t is:

$$p(t | \mathbf{w}) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

- Given a dataset \mathcal{D} with observations $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} / \mathbf{t} = [t_1, \dots, t_N]$,
- The likelihood of drawing our whole dataset from this Gaussian is therefore:

$$p(\mathcal{D} | \mathbf{w}) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

Maximum Likelihood Estimation

Computing the likelihood of our data:

- The probability density of a given value t is:

$$p(t | \mathbf{w}) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

- Given a dataset \mathcal{D} with observations $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} / \mathbf{t} = [t_1, \dots, t_N]$,
- The likelihood of drawing our whole dataset from this Gaussian is therefore:

$$p(\mathcal{D} | \mathbf{w}) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

- Applying the natural logarithm to both sides, we get:

$$\ln p(\mathcal{D} | \mathbf{w}) = \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta \left\{ \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 \right\}$$

Maximum Likelihood Estimation

Computing the likelihood of our data:

- The probability density of a given value t is:

$$p(t | \mathbf{w}) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

- Given a dataset \mathcal{D} with observations $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} / \mathbf{t} = [t_1, \dots, t_N]$,
- The likelihood of drawing our whole dataset from this Gaussian is therefore:

$$p(\mathcal{D} | \mathbf{w}) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

- Applying the natural logarithm to both sides, we get:

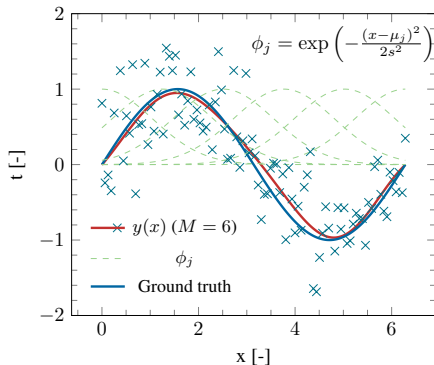
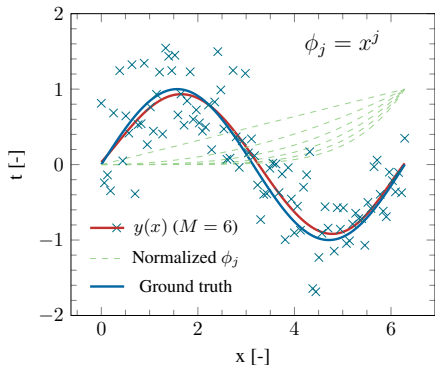
$$\ln p(\mathcal{D} | \mathbf{w}) = \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta \left\{ \frac{1}{2} \sum_{n=1}^N \left(t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) \right)^2 \right\}$$

- **Maximizing the likelihood** is therefore equivalent to **minimizing the error** in red
- This is where the usual loss function for ML regression comes from

Maximum Likelihood Estimation

How does this look like? An example:

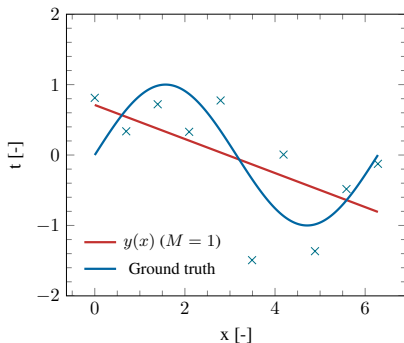
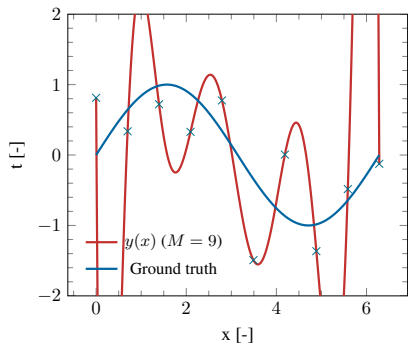
- Dataset with $N = 100$ observations, $M = 6$ basis functions (polynomials or Gaussians)



Overfitting and underfitting MLE models

Also here, flexibility is not always a good thing:

- Dataset with $N = 10$ observations, model with complete order M polynomials
- Again a tradeoff between **bias** and **variance**

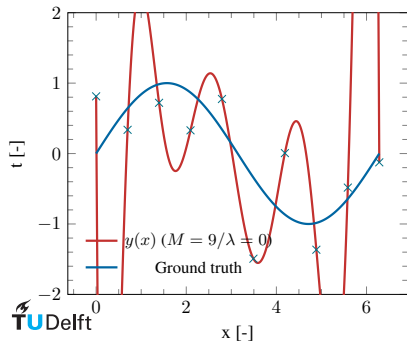


Regularized MLE models

L2 regularization, also known as Ridge Regression or Weight Decay:

- Model complexity is a bit hidden here. We can make it explicit by doing:

$$\mathbf{w}_{\text{ML}} = \arg \min_{\bar{\mathbf{w}}} \left\{ \frac{1}{2} \sum_{n=1}^N \left(t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \right\} \Rightarrow \mathbf{w}_{\text{ML}} = \left(\Phi^T \Phi + \lambda \mathbf{I} \right)^{-1} \Phi^T \mathbf{t}$$

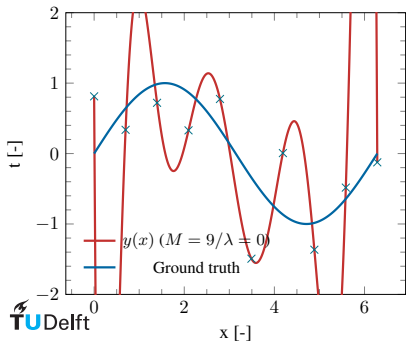


Regularized MLE models

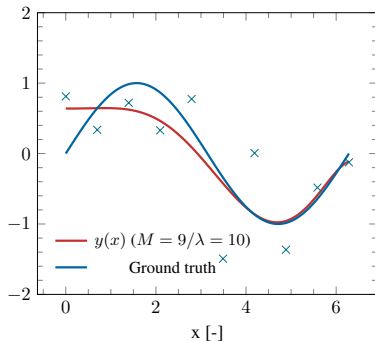
L2 regularization, also known as Ridge Regression or Weight Decay:

- Model complexity is a bit hidden here. We can make it explicit by doing:

$$\mathbf{w}_{\text{ML}} = \arg \min_{\bar{\mathbf{w}}} \left\{ \frac{1}{2} \sum_{n=1}^N \left(t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \right\} \Rightarrow \mathbf{w}_{\text{ML}} = \left(\Phi^T \Phi + \lambda \mathbf{I} \right)^{-1} \Phi^T \mathbf{t}$$



regularizing
and learning λ



Stochastic Gradient Descent

For now we have trained with the complete dataset at once:

- The error function contains **all N data points**:

$$E_{\mathcal{D}} = \frac{1}{2} \sum_{n=1}^N \left(t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right)^2$$

Situations when it is interesting (or necessary) to deviate from this:

- N is too large and computing $(\Phi^T \Phi)^{-1}$ becomes prohibitive
- The model is nonlinear (in \mathbf{w}) and \mathbf{w}_{ML} does not have a closed-form solution
- The dataset is arriving sequentially (e.g. in real time from a sensor)

Stochastic Gradient Descent

Instead of solving directly for \mathbf{w}_{ML} , we can use Gradient Descent:

- Pick a (random) subset \mathcal{B} of the dataset with $N_{\mathcal{B}}$ observations

Stochastic Gradient Descent

Instead of solving directly for \mathbf{w}_{ML} , we can use Gradient Descent:

- Pick a (**random**) subset \mathcal{B} of the dataset with $N_{\mathcal{B}}$ observations
- Update \mathbf{w} with gradients coming from \mathcal{B} and with a fixed **learning rate** η :

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_{\mathcal{B}} \quad \text{with} \quad \nabla E_{\mathcal{B}} = - \sum_{n=1}^{N_{\mathcal{B}}} \left(t_n - \mathbf{w}^{(\tau)\top} \phi(\mathbf{x}_n) \right) \phi(\mathbf{x}_n)^{\top}$$

Stochastic Gradient Descent

Instead of solving directly for \mathbf{w}_{ML} , we can use Gradient Descent:

- Pick a (random) subset \mathcal{B} of the dataset with $N_{\mathcal{B}}$ observations
- Update \mathbf{w} with gradients coming from \mathcal{B} and with a fixed learning rate η :

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_{\mathcal{B}} \quad \text{with} \quad \nabla E_{\mathcal{B}} = - \sum_{n=1}^{N_{\mathcal{B}}} \left(t_n - \mathbf{w}^{(\tau)\text{T}} \phi(\mathbf{x}_n) \right) \phi(\mathbf{x}_n)^{\text{T}}$$

- Every time the complete dataset has been seen, we say an epoch has passed

Stochastic Gradient Descent

Instead of solving directly for \mathbf{w}_{ML} , we can use Gradient Descent:

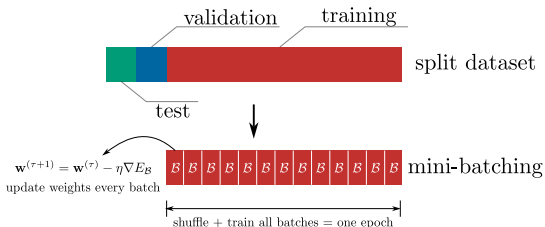
- Pick a (**random**) subset \mathcal{B} of the dataset with $N_{\mathcal{B}}$ observations
- Update \mathbf{w} with gradients coming from \mathcal{B} and with a fixed **learning rate** η :

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_{\mathcal{B}} \quad \text{with} \quad \nabla E_{\mathcal{B}} = - \sum_{n=1}^{N_{\mathcal{B}}} \left(t_n - \mathbf{w}^{(\tau)\text{T}} \phi(\mathbf{x}_n) \right) \phi(\mathbf{x}_n)^{\text{T}}$$

- Every time the complete dataset has been seen, we say **an epoch** has passed

Variations:

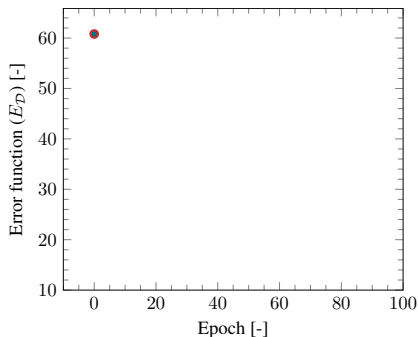
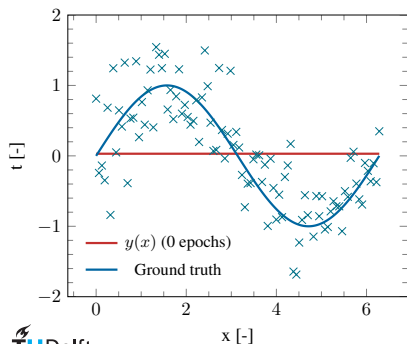
- $N_{\mathcal{B}} = 1$: Online stochastic gradient descent
- $1 < N_{\mathcal{B}} < N$: Minibatch SGD (**most popular**)
- $N_{\mathcal{B}} = N$: Full batch gradient descent



Stochastic Gradient Descent

An example:

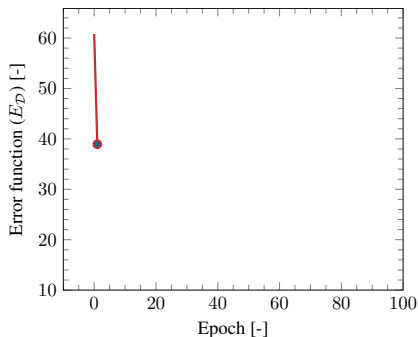
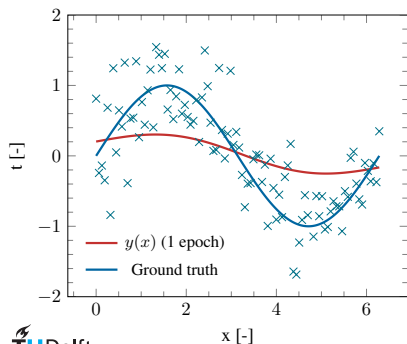
- Same example as before, with $N = 100$ and $M = 6$ polynomial basis functions
- We fix the learning rate $\eta = 0.001$ and minibatch size $N_{\mathcal{B}} = 10$



Stochastic Gradient Descent

An example:

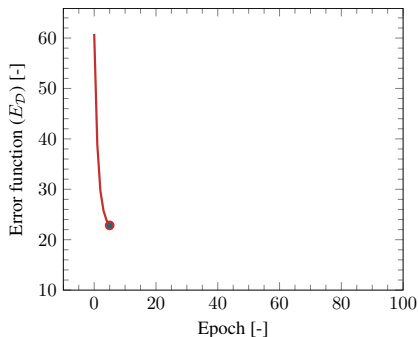
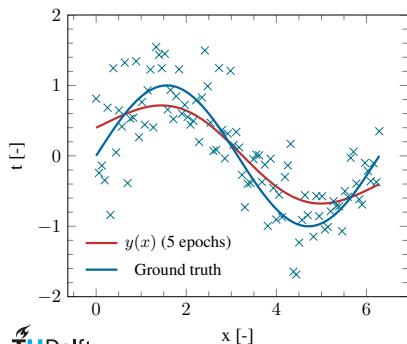
- Same example as before, with $N = 100$ and $M = 6$ polynomial basis functions
- We fix the learning rate $\eta = 0.001$ and minibatch size $N_{\mathcal{B}} = 10$



Stochastic Gradient Descent

An example:

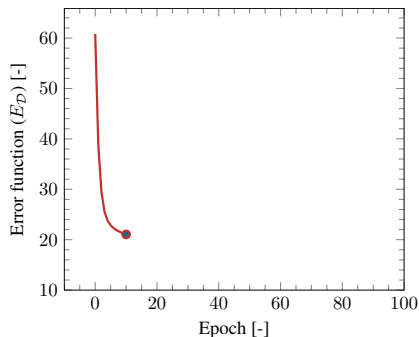
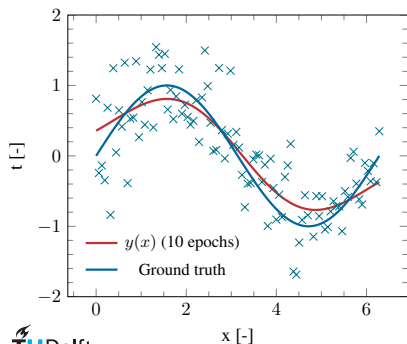
- Same example as before, with $N = 100$ and $M = 6$ polynomial basis functions
- We fix the learning rate $\eta = 0.001$ and minibatch size $N_{\mathcal{B}} = 10$



Stochastic Gradient Descent

An example:

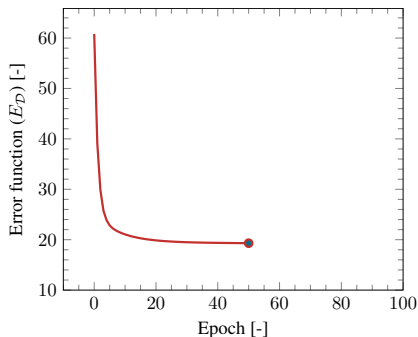
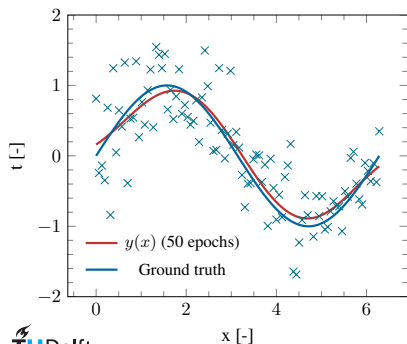
- Same example as before, with $N = 100$ and $M = 6$ polynomial basis functions
- We fix the learning rate $\eta = 0.001$ and minibatch size $N_{\mathcal{B}} = 10$



Stochastic Gradient Descent

An example:

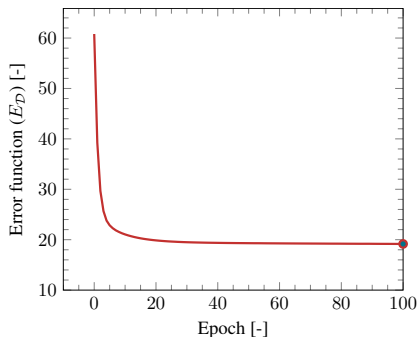
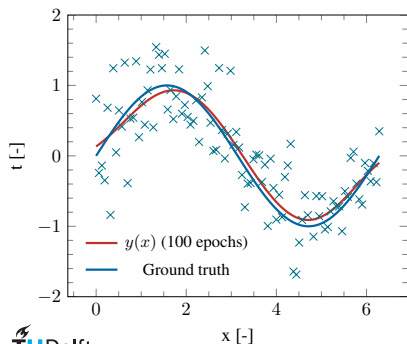
- Same example as before, with $N = 100$ and $M = 6$ polynomial basis functions
- We fix the learning rate $\eta = 0.001$ and minibatch size $N_{\mathcal{B}} = 10$



Stochastic Gradient Descent

An example:

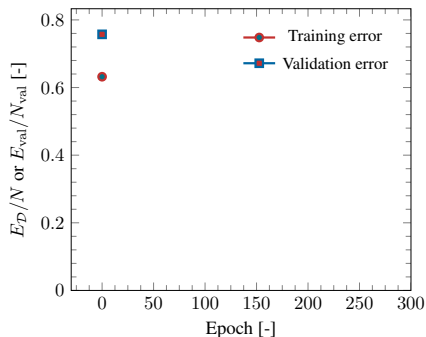
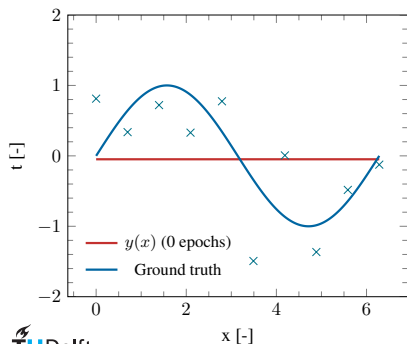
- Same example as before, with $N = 100$ and $M = 6$ polynomial basis functions
- We fix the learning rate $\eta = 0.001$ and minibatch size $N_{\mathcal{B}} = 10$



Stochastic Gradient Descent

Using SGD progress to spot signs of overfitting:

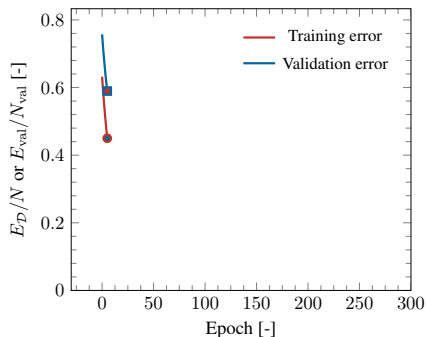
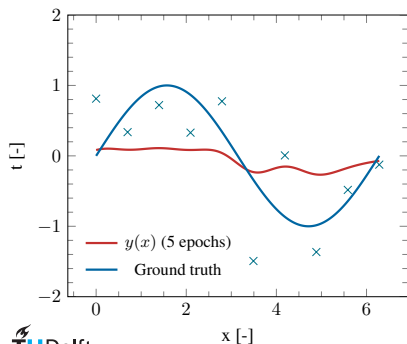
- Tracking the error on a **validation dataset** after every epoch
- This motivates the **early stopping** strategy popular in the deep learning community



Stochastic Gradient Descent

Using SGD progress to spot signs of overfitting:

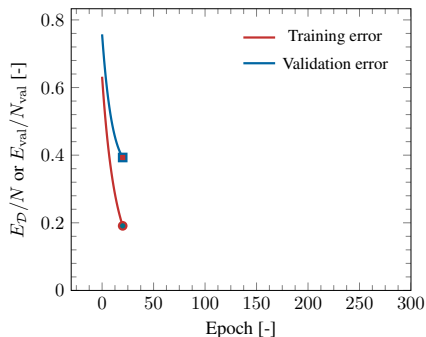
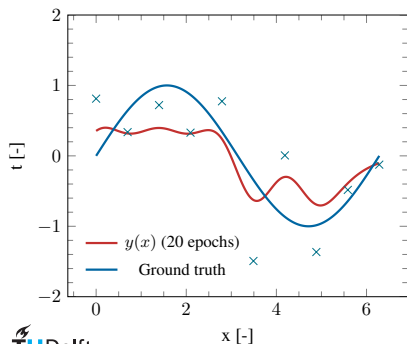
- Tracking the error on a **validation dataset** after every epoch
- This motivates the **early stopping** strategy popular in the deep learning community



Stochastic Gradient Descent

Using SGD progress to spot signs of overfitting:

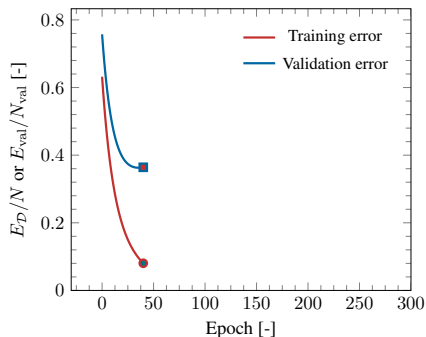
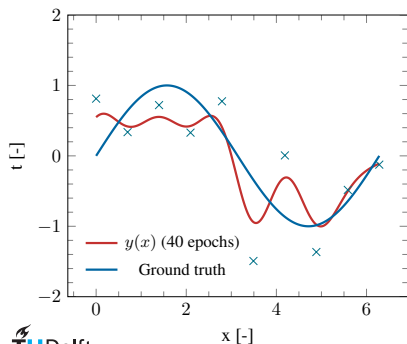
- Tracking the error on a **validation dataset** after every epoch
- This motivates the **early stopping** strategy popular in the deep learning community



Stochastic Gradient Descent

Using SGD progress to spot signs of overfitting:

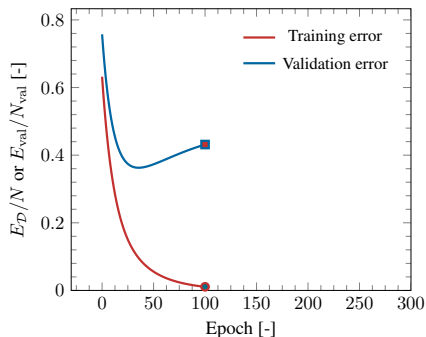
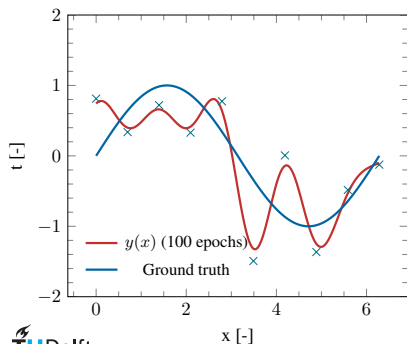
- Tracking the error on a **validation dataset** after every epoch
- This motivates the **early stopping** strategy popular in the deep learning community



Stochastic Gradient Descent

Using SGD progress to spot signs of overfitting:

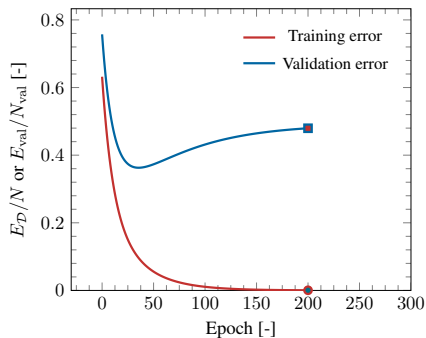
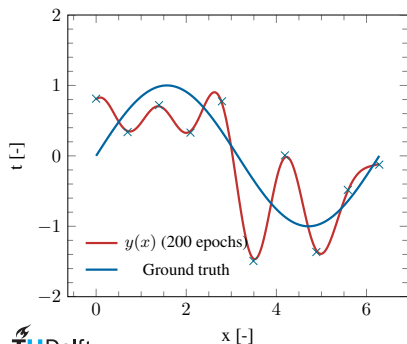
- Tracking the error on a **validation dataset** after every epoch
- This motivates the **early stopping** strategy popular in the deep learning community



Stochastic Gradient Descent

Using SGD progress to spot signs of overfitting:

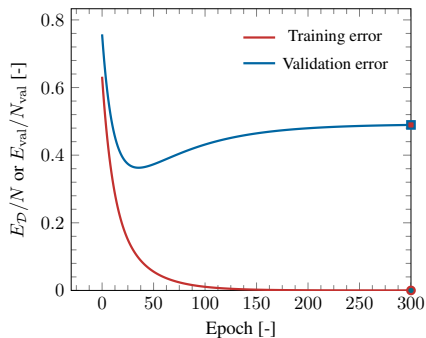
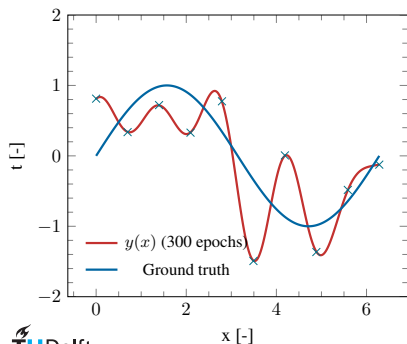
- Tracking the error on a **validation dataset** after every epoch
- This motivates the **early stopping** strategy popular in the deep learning community



Stochastic Gradient Descent

Using SGD progress to spot signs of overfitting:

- Tracking the error on a **validation dataset** after every epoch
- This motivates the **early stopping** strategy popular in the deep learning community

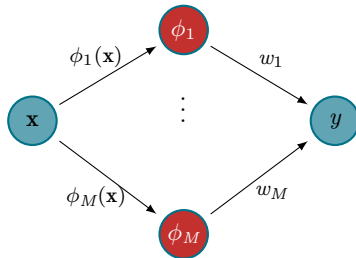


Adaptive basis functions

Up until now, the basis functions have been fixed a priori:

- Polynomials: number of terms M , polynomial degrees of each term
- Gaussians: bandwidth s , basis function centers μ_j

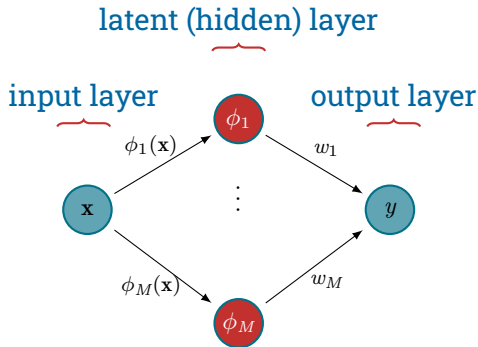
$$y = \phi_1(\mathbf{x})w_1 + \phi_2(\mathbf{x})w_2 + \cdots + \phi_M(\mathbf{x})w_M$$



Adaptive basis functions

For now, only half of the model is trainable:

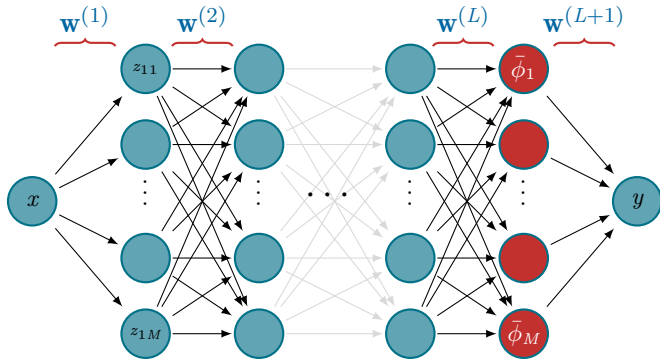
- Input to hidden **encoding** ($\phi_1 \dots \phi_M$) fixed, hidden to output **decoding** (w) trained
- What if we could also train the first half?



Artificial Neural Networks

Replacing basis functions by several layers of nonlinear transformations:

- **Neural Network:** layers of **neurons** linked by **weighted connections**
- Values at the red layer can be seen as coming from new, learned basis functions $\bar{\phi}(x)$



Neural Networks – Activation functions

For a given neuron, **forward propagation** happens in two steps:

- A **linear** combination of values from the previous layer:

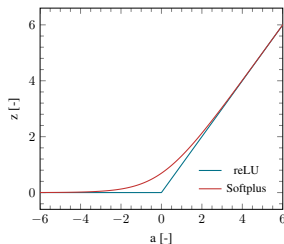
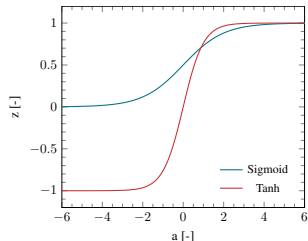
$$a_{lj} = \sum_i^D w_{ji}^{(l)} z_i^{(l-1)} + w_{j0}^{(l)}$$

- A **nonlinear** transformation with an **activation function**:

$$z_{lj} = h(a_{lj})$$

Choosing the activation function:

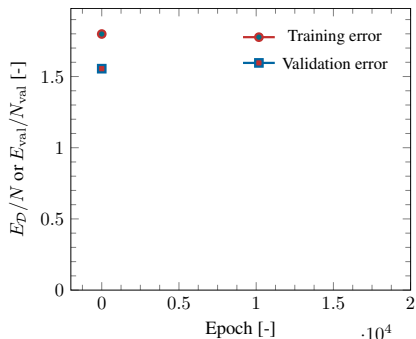
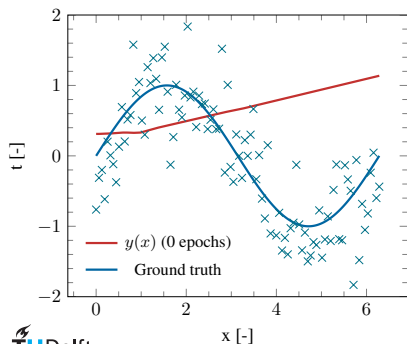
- Application dependent
- Can be seen as a **hyperparameter**



Neural Networks – Example

Same example as before, but now with a neural network:

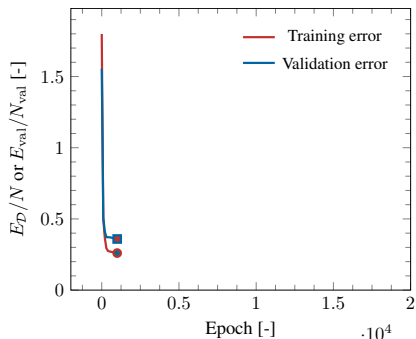
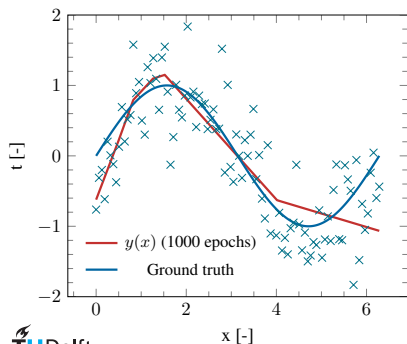
- Full batch **Adam SGD** (variable learning rate)
- Two hidden layers, 10 neurons each, **ReLU** activation



Neural Networks – Example

Same example as before, but now with a neural network:

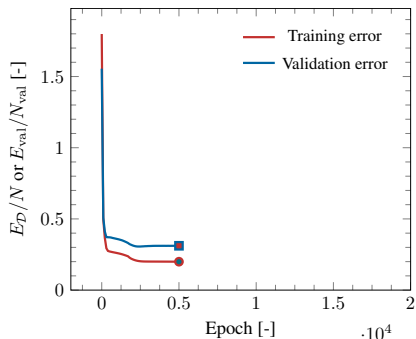
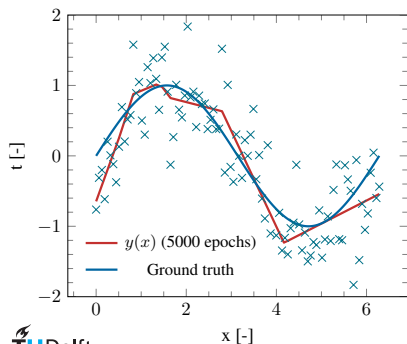
- Full batch **Adam SGD** (variable learning rate)
- Two hidden layers, 10 neurons each, **ReLU** activation



Neural Networks – Example

Same example as before, but now with a neural network:

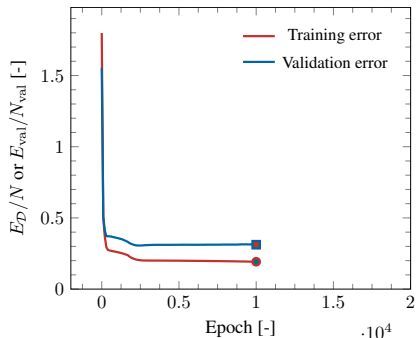
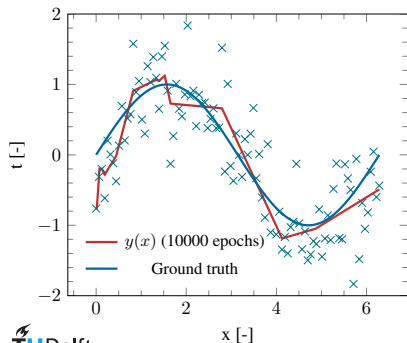
- Full batch **Adam SGD** (variable learning rate)
- Two hidden layers, 10 neurons each, **ReLU** activation



Neural Networks – Example

Same example as before, but now with a neural network:

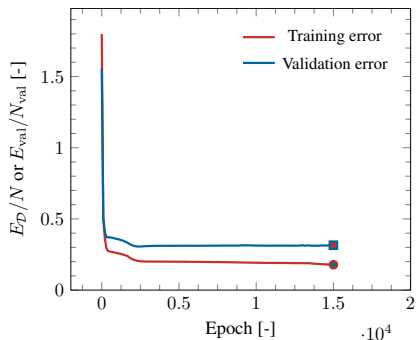
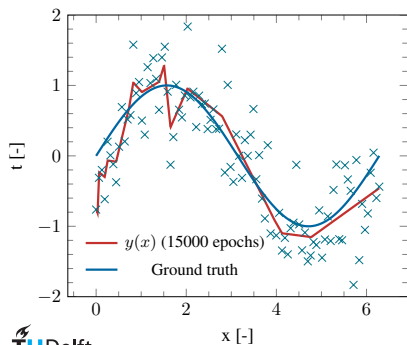
- Full batch **Adam SGD** (variable learning rate)
- Two hidden layers, 10 neurons each, **ReLU** activation



Neural Networks – Example

Same example as before, but now with a neural network:

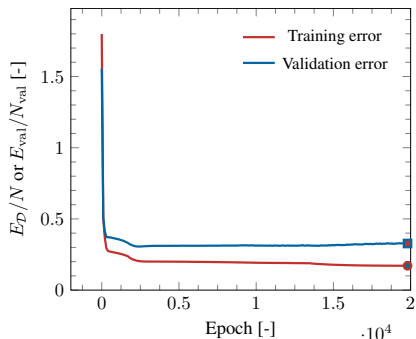
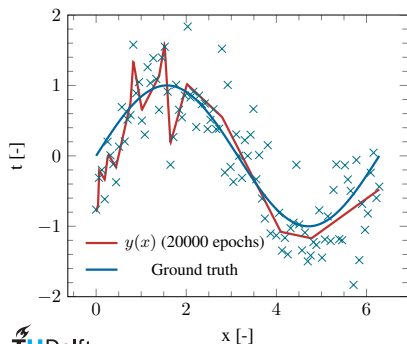
- Full batch **Adam SGD** (variable learning rate)
- Two hidden layers, 10 neurons each, **ReLU** activation



Neural Networks – Example

Same example as before, but now with a neural network:

- Full batch **Adam SGD** (variable learning rate)
- Two hidden layers, 10 neurons each, **ReLU** activation



Recap

You should now be able to:

- Understand and compare different approaches for regression modeling
- Construct parsimonious regression models for general applications
- Critically assess model performance from a probabilistic standpoint

Main takeaways:

- Model flexibility is not always beneficial, **overfitting** can be a major issue
- Selecting good models always boils down to balancing **bias and variance**
- Simply **fitting data is not learning**. ML is all about **finding hidden patterns in data**
- **Neural Networks are not magic**, can be seen as adaptive versions of simpler models
- Fancier deep learning models share this same foundation but add extra heuristics