

3 Neural Networks

Leon Herrmann

Stefan Kollmannsberger

Chair of Data Engineering in Construction

Bauhaus-Universität Weimar

Paris, February 2025

*Deep Learning in Computational Mechanics – an introductory course,
Herrmann et al. 2025*



website



book



Contents

- [2 Fundamental Concepts of Machine Learning](#)
- [3.1 Fully Connected Neural Network](#)
- [3.4 Backpropagation](#)
- [3.6 Learning Algorithm](#)
- [3.8 Approximating the Sine Function](#)
- [3.9.1 Convolutional Neural Networks](#)
- [3.9.2 Graph Neural Networks](#)
- [3.9.6 Recurrent Neural Networks](#)
- [3.9.8 Physics-Inspired Architectures for Dynamics \(Hamiltonian & Lagrangian Neural Networks\)](#)
- [4 Introduction to Physics-Informed Neural Networks](#)

3 Neural Networks

- Predictions \hat{y} from an input x are performed via **forward propagation**

$$a_k^i = \sigma \left(\sum_{j=0}^n w_{kj}^i a_j^{i-1} + b_k^i \right)$$

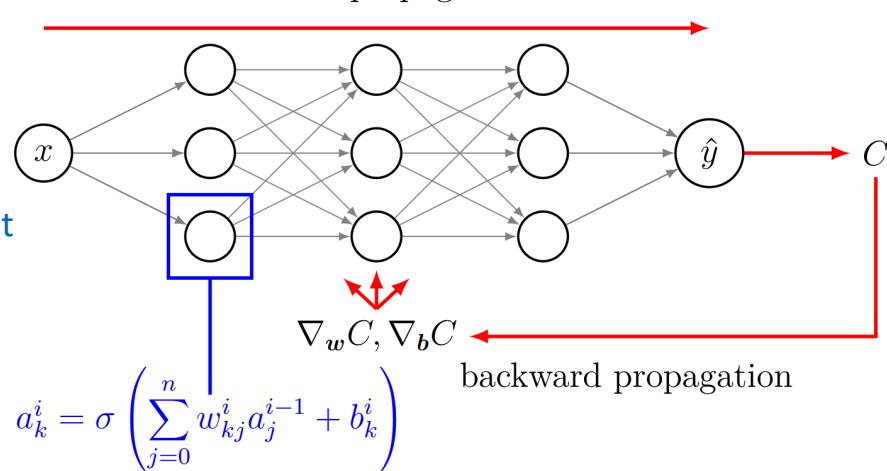
forward propagation

- Linear transformations w, b
- Nonlinear activation functions σ
- Performance assessment with **cost function** C
- Backward propagation** computes **sensitivities** $\nabla_w C, \nabla_b C$
- Learnable parameters are update through **gradient descent**

$$w \leftarrow w - \nabla_w C$$

$$b \leftarrow b - \nabla_b C$$

Backpropagation is a specialization of (reverse-mode) automatic differentiation applied to neural networks

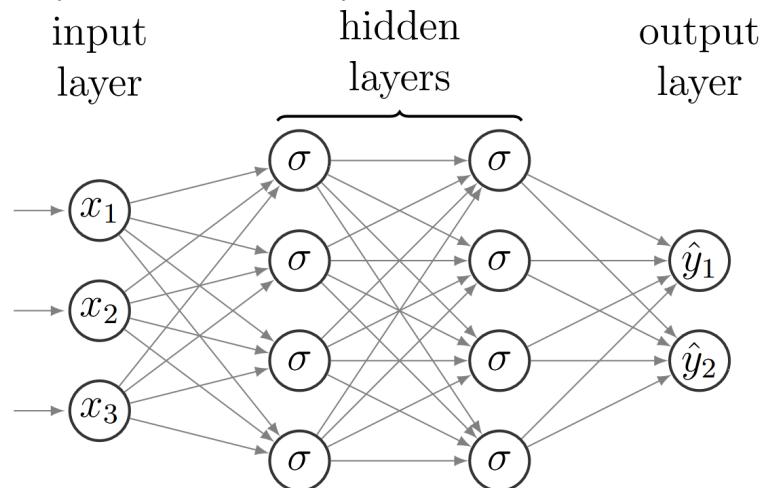


3.1 Fully Connected Neural Network

- Neural networks are a parametrized function defining a mapping $y = f_{NN}(x)$
- Neural networks $f_{NN}(x)$ are composed of nested functions $f_{NN} = f_3(f_2(f_1(x)))$
- Each nested function f_l represents a layer of the network
- The input x flows from the input layer through the hidden layers to the output layer (feed-forward)
- A network with more than one hidden layer is a **deep neural network** (otherwise **shallow**)

Components of a neural network

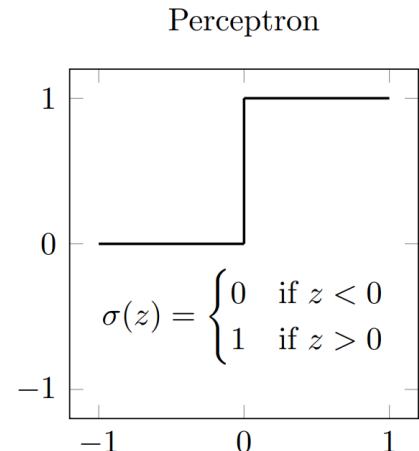
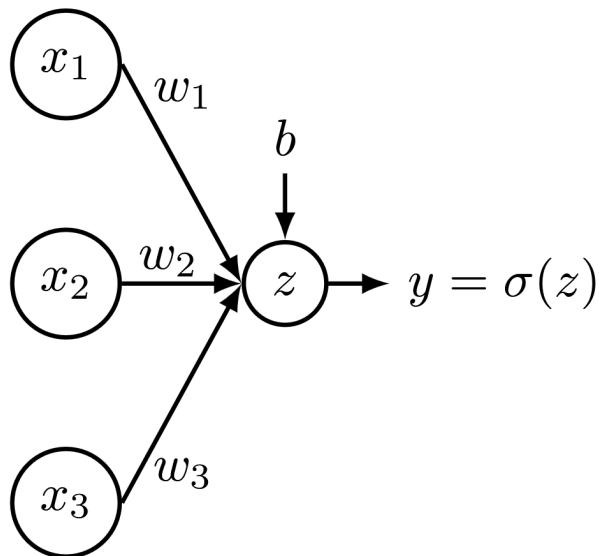
- **Neurons** (circles) represent an intermediate state in the flow and serve as input to the next neuron
- **Weights** (arrows) connect the neurons
 - If all neurons between every neighboring layer are connected, the network is **fully connected**
- **Depth:** number of hidden layers
- **Width:** number of neurons per layer



This is a fully connected feed-forward neural network. This is typically abbreviated as **fully connected neural network**.

How do neural networks predict?

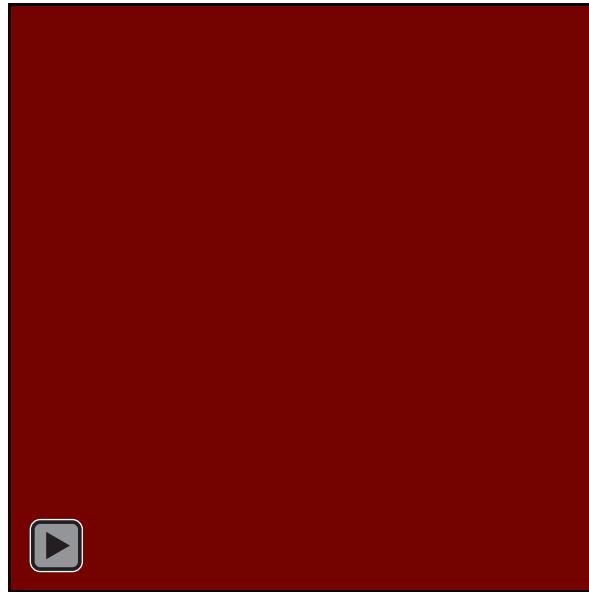
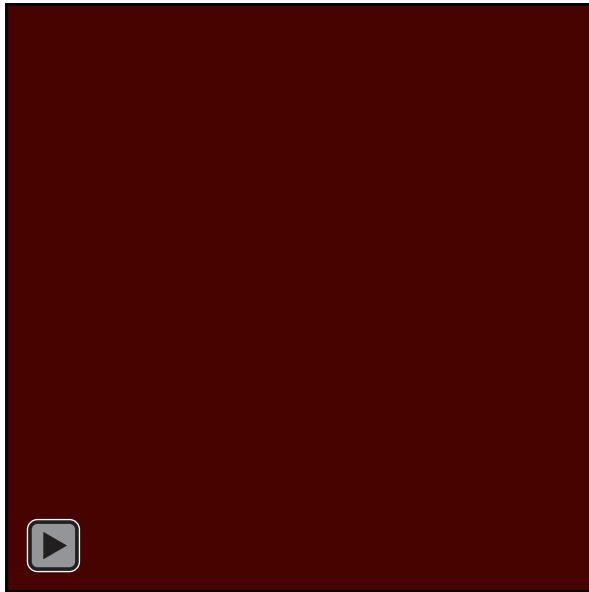
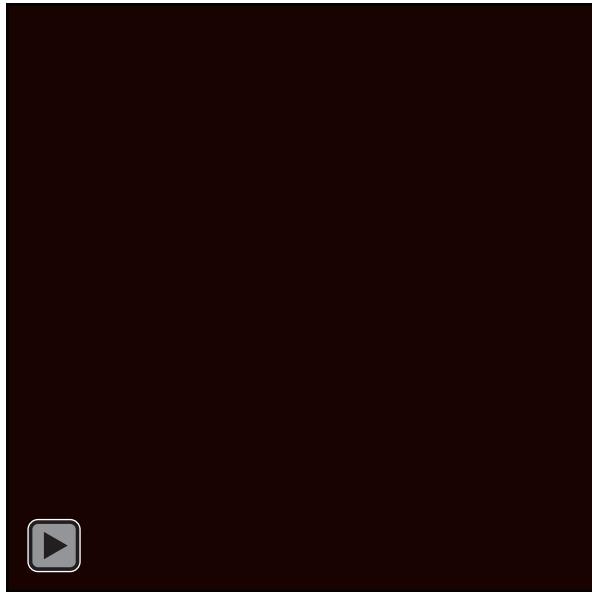
- Consider a network, that determines if you should go to a concert tonight
- Possible inputs x could be
 - x_1 = Exam tomorrow (0 or 1)
 - x_2 = Rating of the artist (float)
 - x_3 = Music quality at the concert (float)
 $y = \sigma(w_1x_1 + w_2x_2 + w_3x_3)$
- Output of the network is either 0 or 1
- The weights are chosen depending on how important each aspect is for you individually
- The bias shows the general tendency of an individual to go to a concert



A **perceptron** is an early architecture without hidden layers and a step function as activation function. Why is it problematic to train with modern techniques?

3.2 Forward Propagation

- **Universal approximation theorem:** A fully connected feed-forward neural network with one hidden layer can approximate **any continuous function** with arbitrary precision
- Output visualization of a multi-layer neural network, inputs are coordinates x, y



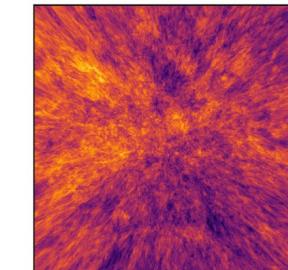
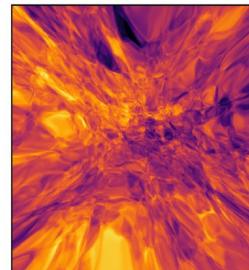
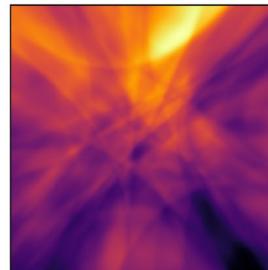
3.2 Forward Propagation

Deeper neural networks are more favorable, as fewer parameters are required to achieve greater approximation power (more nested non-linearities)

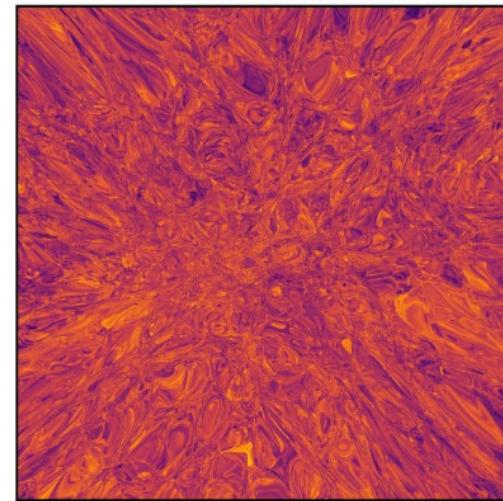
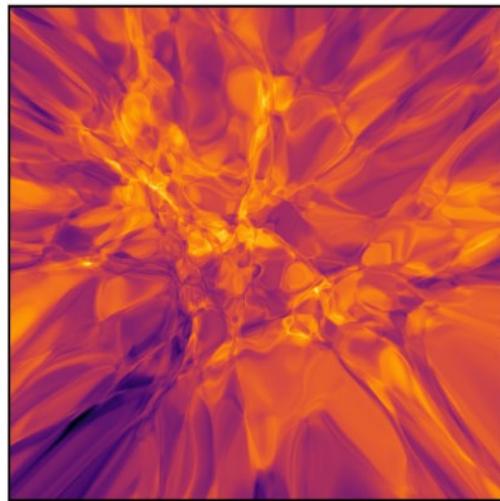
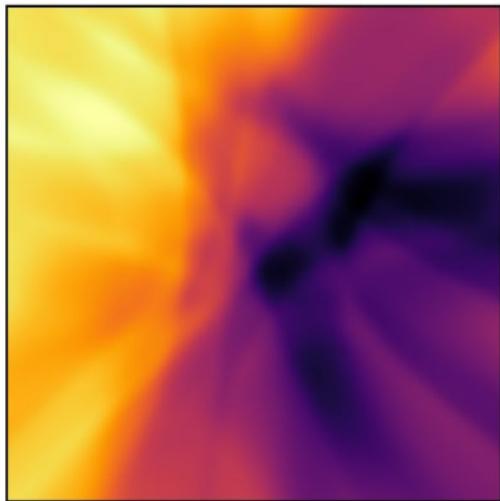
2 layers with 300 neurons: 5 layers with 150 neurons:

91'501 parameters

91'201 parameters



5 layers
with 800
neurons:
2'566'401
parameters



2 layers with 100 neurons:
10'501 parameters

5 layers with 100 neurons:
40'801 parameters

10 layers with 100 neurons:
91'301 parameters

Exercises

E.10 Neural Network Representational Capacity (C)

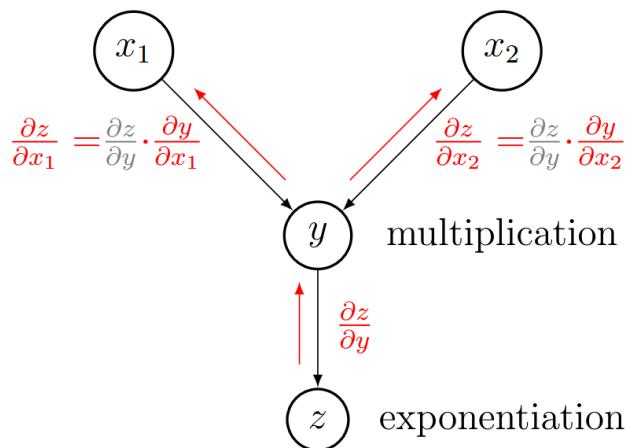
- Experience the representational capacity of neural networks by experimenting with the neural network architecture.

3.4 Backpropagation

- A **computation graph** tracks every **elementary arithmetic operation**, enabling **differentiation**
- Consider the example

$$y = x_1 \cdot x_2$$
$$z = y^2$$

- Using the **chain rule**, the derivatives of z can be obtained with respect to x_1, x_2



- In PyTorch the arithmetic operations are stored as **function handles** (`.grad_fn`)

3.4 Backpropagation

A general algorithm can be formulated by considering an arbitrary node c connected to nodes a, b, d, e

- a, b are incoming
- d, e are outgoing

Gradient with respect to c

$$\frac{\partial C}{\partial c} = \frac{\partial C}{\partial c_d} + \frac{\partial C}{\partial c_e}$$

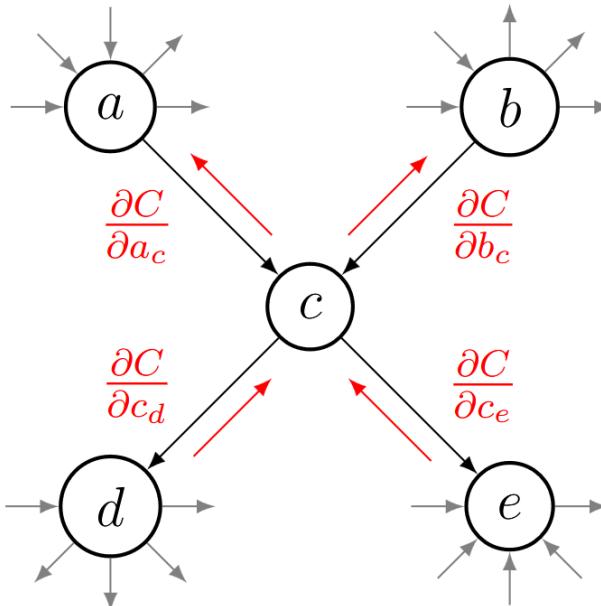
where

$$\frac{\partial C}{\partial c_d} = \frac{\partial C}{\partial d} \cdot \frac{\partial d}{\partial c_d}; \quad \frac{\partial C}{\partial c_e} = \frac{\partial C}{\partial e} \cdot \frac{\partial e}{\partial c_e}$$

(Back)propagation to incoming nodes

$$\frac{\partial C}{\partial a_c} = \frac{\partial C}{\partial c} \cdot \frac{\partial c}{\partial a_c}; \quad \frac{\partial C}{\partial b_c} = \frac{\partial C}{\partial c} \cdot \frac{\partial c}{\partial b_c}$$

which might be mixed with other gradient contributions to nodes a, b



3.6 Learning Algorithm

Prerequisites to train a neural network for a supervised learning task

- Input data \tilde{X} and corresponding targets \tilde{y} , divided into a training, validation and test set
- Network topology (more in Chapter 3.9)
- Network initialization (not covered in this lecture, see, e.g., <https://www.deeplearning.ai/ai-notes/initialization/>)
- Activation function σ
- Cost function C defining the prediction quality
- Method of computing gradients with respect to the network parameters, e.g., backpropagation
- Optimizer with hyperparameters, such as the learning rate α

Learning algorithms: 3 most common variants of one and the same idea: gradient-based optimization (as seen in Chapter 2)

- full batch
- stochastic
- mini batch

3.6 Learning Algorithm – Full-batch

Algorithm 5 Training a neural network with full-batch gradient descent. The inner loop is only displayed for a better understanding. Normally, the loop over the examples is vectorized for more efficient computations.

Require: training data $\tilde{\mathbf{X}}$, targets $\tilde{\mathbf{y}}$

define network architecture (input layer, hidden layers, output layer, activation function) set learning rate α

initialize weights \mathbf{W} and biases \mathbf{b}

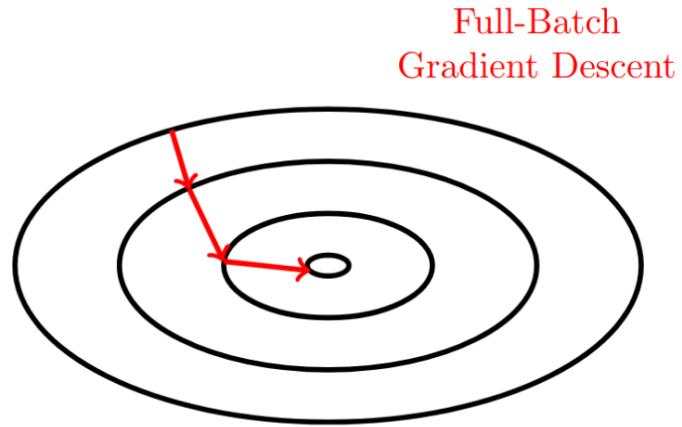
for all epochs do

for example $i \leftarrow 1$ to $m_{\mathcal{P}}$ do

apply forward propagation: $\hat{y}_i \leftarrow f_{NN}(\mathbf{x}_i; \mathbf{W}, \mathbf{b})$ ▷ cf. Section 3.2

update biases: $\mathbf{b} \leftarrow \mathbf{b} - \alpha \frac{\partial C}{\partial \mathbf{b}}$

end for



Accurate gradients at a high cost

3.6 Learning Algorithm – Stochastic

Algorithm 6 Training a neural network with stochastic gradient descent

Require: training data $\tilde{\mathbf{X}}$, targets $\tilde{\mathbf{y}}$

define network architecture (input layer, hidden layers, output layer, activation function)
set learning rate α

initialize weights \mathbf{W} and biases \mathbf{b}

for all epochs do

 for example $i \leftarrow 1$ to m_D do

 apply forward propagation $\hat{\mathbf{y}}_i \leftarrow f_{NN}(\mathbf{x}_i; \mathbf{W}, \mathbf{b})$

 compute loss: $C_i \leftarrow (\tilde{y}_i - \hat{y}_i)^2$

 apply backpropagation for gradients $\partial C_i / \partial \mathbf{W}$, $\partial C_i / \partial \mathbf{b}$

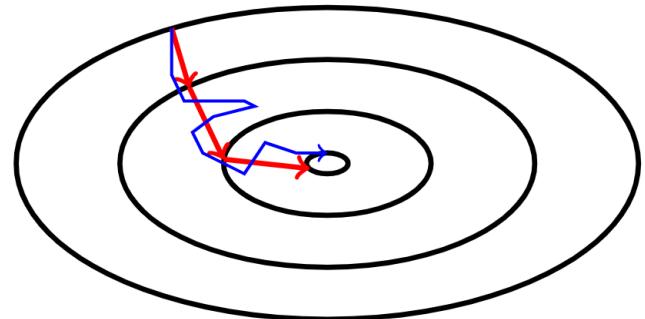
 update weights: $\mathbf{W} \leftarrow \mathbf{W} - \alpha \frac{\partial C_i}{\partial \mathbf{W}}$

 update biases: $\mathbf{b} \leftarrow \mathbf{b} - \alpha \frac{\partial C_i}{\partial \mathbf{b}}$

 end for

end for

Stochastic Gradient Descent Full-Batch Gradient Descent



▷ cf. Section 3.2

▷ cf. Section 3.4

Cheap & inaccurate approximation
of gradients (enabling escape of
local optima)

3.6 Learning Algorithm – Mini-batch

Algorithm 7 Training a neural network with mini-batch gradient descent

Require: training data $\tilde{\mathbf{X}}$, targets $\tilde{\mathbf{y}}$

define network architecture (input layer, hidden layers, output layer, activation function)
set learning rate α

initialize weights \mathbf{W} and biases \mathbf{b}

for all epochs **do**

shuffle rows of \mathbf{X} and \mathbf{y} synchronously (optional)

divide \mathbf{X} and \mathbf{y} into n batches of size k

for all batches **do**

for example $i \leftarrow 1$ to k **do**

 apply forward propagation: $\hat{\mathbf{y}}_i \leftarrow f_{NN}(\mathbf{x}_i; \mathbf{W}, \mathbf{b})$ ▷ cf. Section 3.2

 compute loss: $C_i \leftarrow (\tilde{\mathbf{y}}_i - \hat{\mathbf{y}}_i)^2$

 apply backpropagation for gradients $\partial C_i / \partial \mathbf{W}$, $\partial C_i / \partial \mathbf{b}$ ▷ cf. Section 3.4

end for

 compute mini-batch cost function: $C \leftarrow \frac{1}{k} \sum_{i=1}^k C_i$

 compute mini-batch gradient w.r.t. \mathbf{W} : $\frac{\partial C}{\partial \mathbf{W}} \leftarrow \frac{1}{k} \sum_{i=1}^k \frac{\partial C_i}{\partial \mathbf{W}}$

 compute mini-batch gradients w.r.t. \mathbf{b} : $\frac{\partial C}{\partial \mathbf{b}} \leftarrow \frac{1}{k} \sum_{i=1}^k \frac{\partial C_i}{\partial \mathbf{b}}$

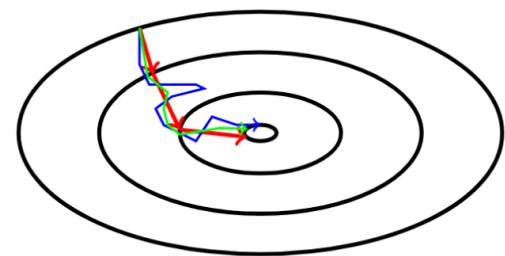
 update weights: $\mathbf{W} \leftarrow \mathbf{W} - \alpha \frac{\partial C}{\partial \mathbf{W}}$

 update biases: $\mathbf{b} \leftarrow \mathbf{b} - \alpha \frac{\partial C}{\partial \mathbf{b}}$

end for

end for

Stochastic Gradient Descent	Mini-Batch Gradient Descent	Full-Batch Gradient Descent
--------------------------------	--------------------------------	--------------------------------



Improved approximation of the
gradients at low computational cost

3.7 Regularization of Neural Networks

Noise: a Flaw in Human Judgement, Kahneman et al. 2021

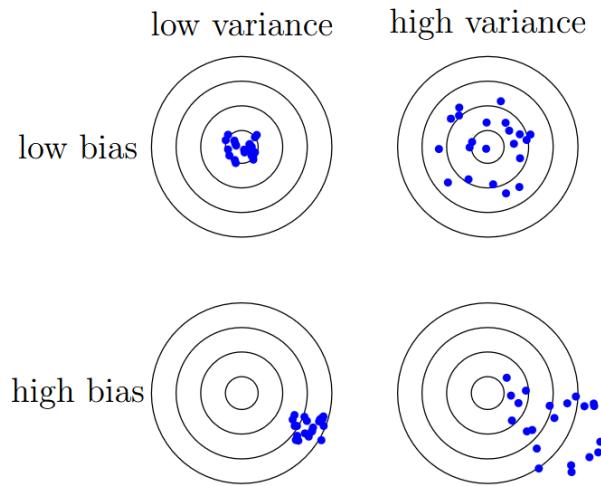
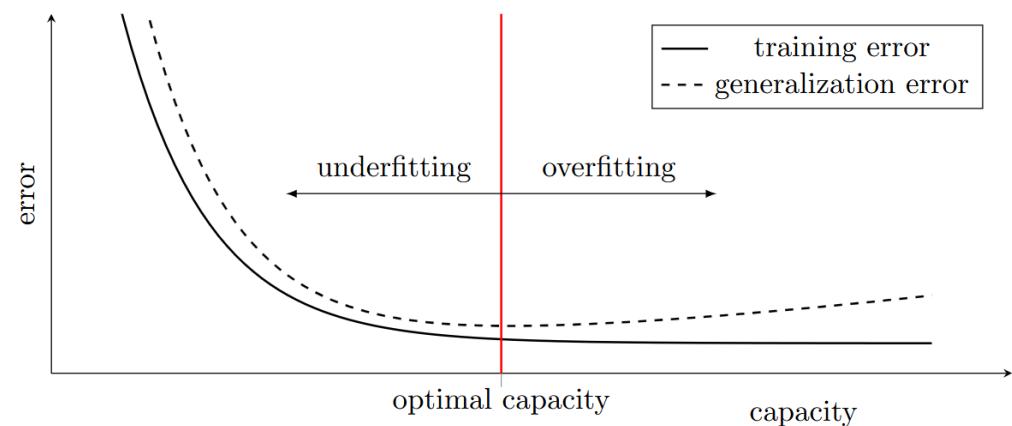
Aim is to bring **testing** and **training error** closer together (without significantly increasing training error)

- Trade **reduction of variance** for a slightly **increased bias**

Common approaches to reduce the variance

- **Constraints** on model parameters (**regularization**)
- **Simpler models** for better generalization

In deep learning, we typically use larger models with appropriate regularization mechanisms



3.8 Approximating the Sine Function

Goal is to approximate

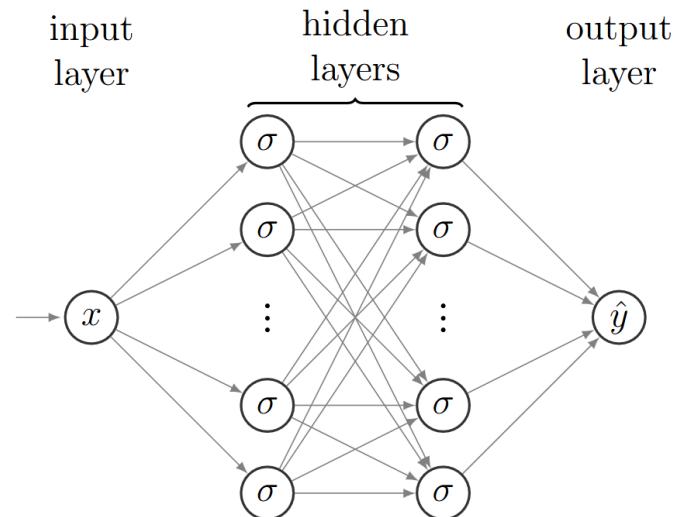
$$f(x) = \sin(2\pi x), x \in [-1, 1]$$

Neural network architecture

- 2 hidden layers
- 50 neurons each
- Activation function $\sigma(z_i) = \tanh(z_i)$
- $f_{NN} = w_3^T \sigma(W_2(\sigma(w_1^T x + b_1)) + b_2) + b_3 = \hat{y}$

Data

- Training set contains 40 samples
- Validation set contains 40 samples
- Generated with $y = \sin(2\pi x) + \epsilon, \epsilon = 0.1 \cdot U(-1,1)$
where $U(-1,1)$ is a uniform random distribution in the interval $[-1,1]$
- Test set is the analytical solution



3.8 Approximating the Sine Function

- Goal is to approximate

$$f(x) = \sin(2\pi x), x \in [-1, 1]$$

- Network architecture

$$f_{NN} = w_3^T \sigma \left(W_2 \left(\sigma(w_1^T x + b_1) \right) + b_2 \right) + b_3 = \hat{y}$$

- Cost function used for training

$$\mathcal{L}_D^{(\text{train})} = \frac{1}{m_D^{(\text{train})}} \sum_{i=1}^{m_D^{(\text{train})}} \left(\tilde{y}_i^{(\text{train})} - \hat{y}_i^{(\text{train})} \right)^2$$

- Neural network parameter initialization with a uniform distribution with special bounds
- Optimization with Adam for 10'000 epochs

3.8 Approximating the Sine Function – PyTorch

Neural network definition

```
modules = []
modules.append(torch.nn.Linear(1, 50))
modules.append(torch.nn.Tanh())
modules.append(torch.nn.Linear(50, 50))
modules.append(torch.nn.Tanh())
modules.append(torch.nn.Linear(50, 1))

model = torch.nn.Sequential(*modules)
```

Training data

```
x = torch.rand((40, 1)) * 2 - 1
noise = torch.rand(x.shape) * 0.2 - 0.1
y = torch.sin(2 * torch.pi * x) + noise
```

3.8 Approximating the Sine Function – PyTorch

Cost function definition

```
def costFunction(y, yPred):  
    return torch.mean((y - yPred) ** 2)
```

Prediction and cost function evaluation

```
yPred = model(x)  
cost = costFunction(y, yPred)  
cost.backward()
```

Optimizer definition

```
epochs = 10000  
optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
```

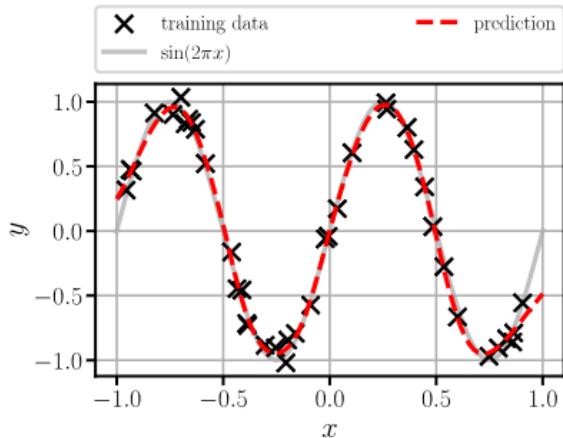
3.8 Approximating the Sine Function – PyTorch

Training loop

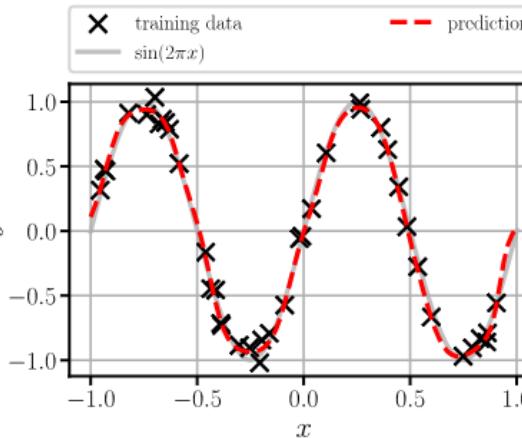
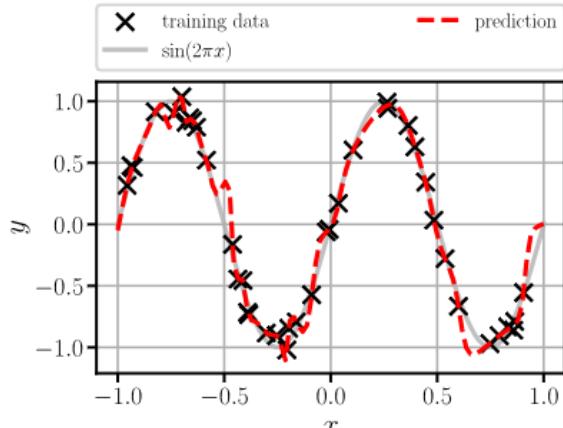
```
for epoch in range(epochs):
    optimizer.zero_grad()
    yPred = model(x)
    cost = costFunction(y, yPred)
    cost.backward()
    optimizer.step()
```

3.8 Approximating the Sine Function – Results

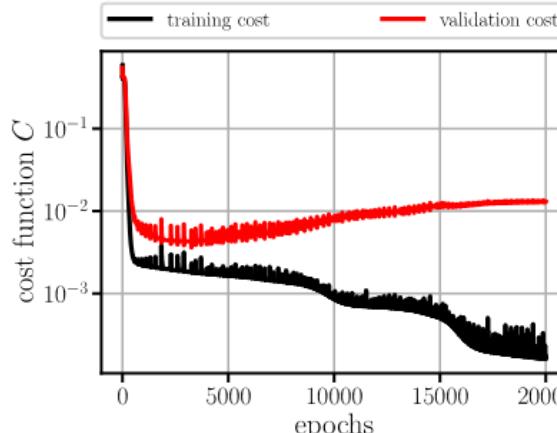
500 epochs



20'000 epochs



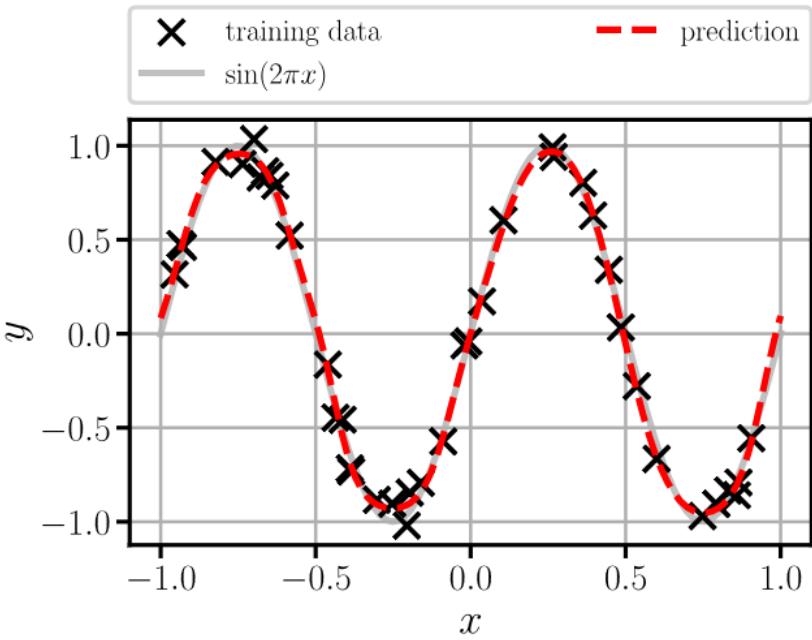
5'000 epochs



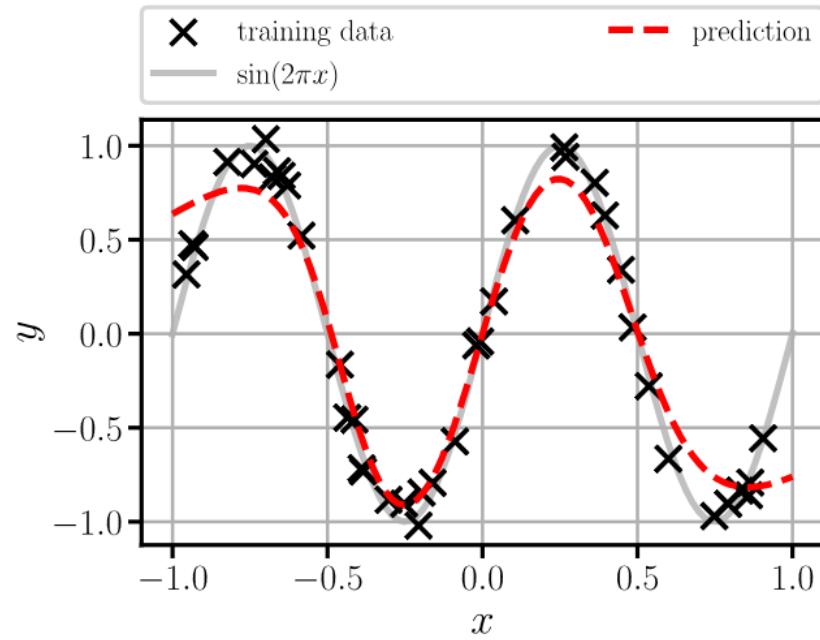
Learning history

3.8 Approximating the Sine Function – Results

20'000 epochs with regularization

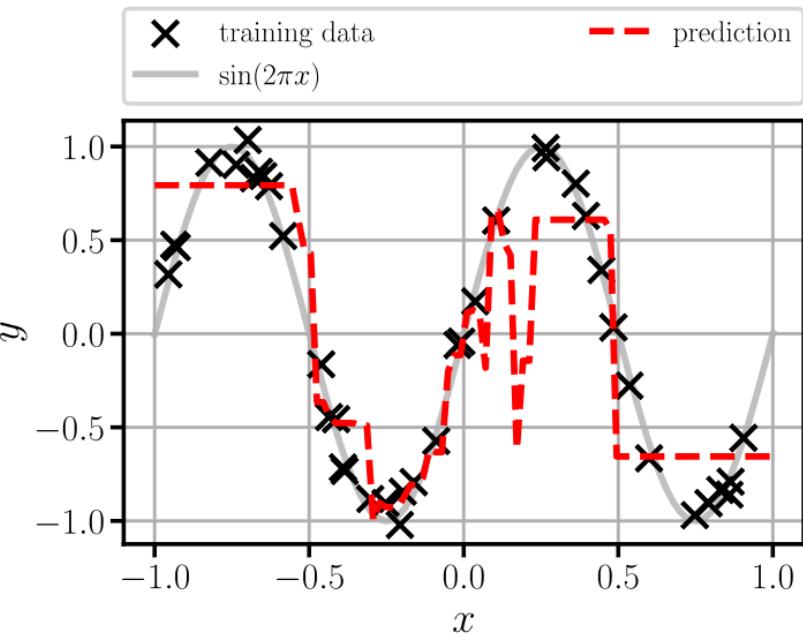


20'000 epochs without regularization

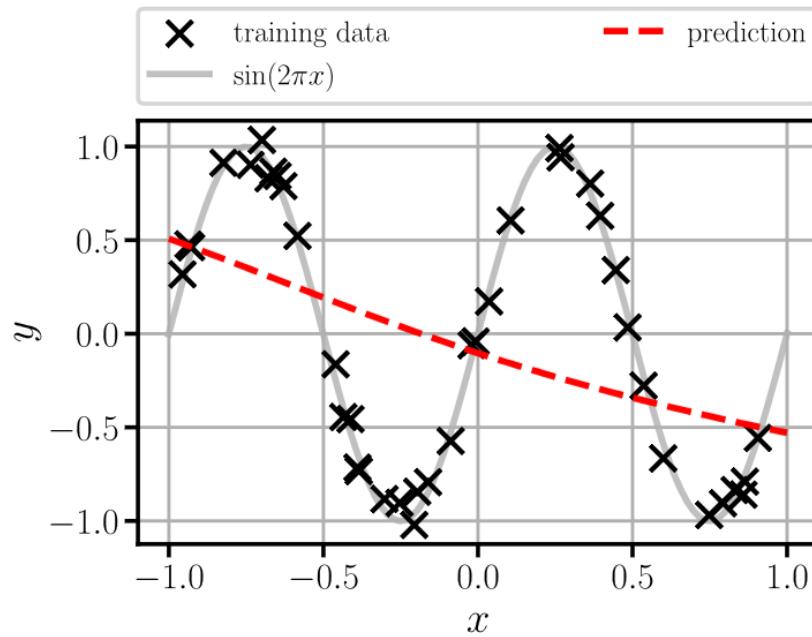


3.8 Approximating the Sine Function – Results

5'000 epochs with Adam at $\alpha = 0.1$



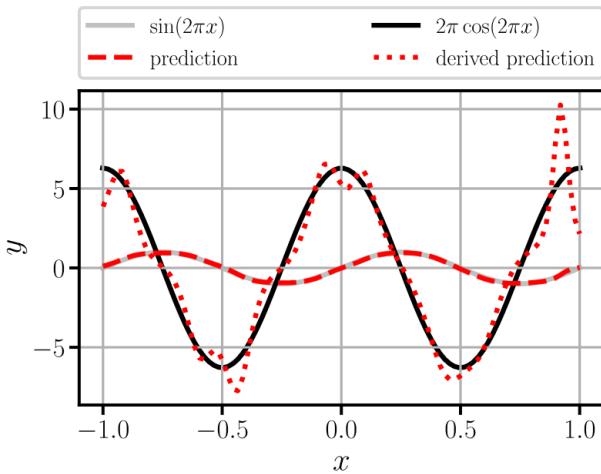
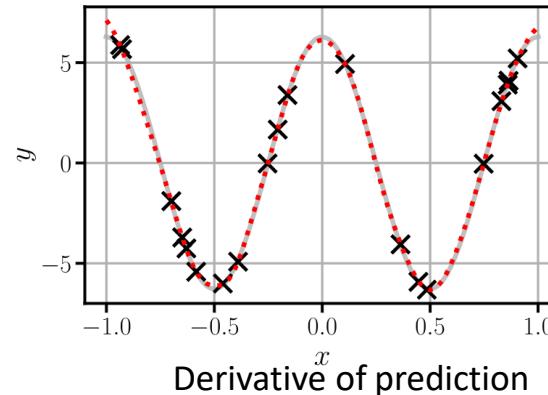
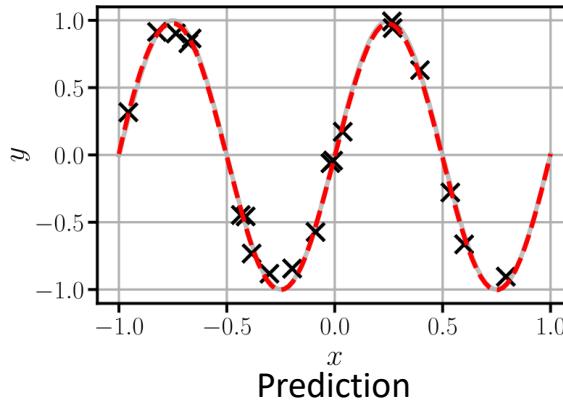
5'000 epochs with SGD



3.8 Approximating the Sine Function – Sobolev

- Derivatives of predictions are typically inaccurate
- Sobolev training** increases accuracy by incorporating derivatives in training data (higher order derivatives also improve)

```
def costFunction(y, yPred, dy, dyPred):
    lossy = torch.mean((y - yPred) ** 2)
    lossdy = torch.mean((dy - dyPred) ** 2)
return lossy + lossdy
```



Exercises

E.7 Approximating the Sine Function (C)

- In PyTorch, build a fully connected neural network, that approximates the sine function.

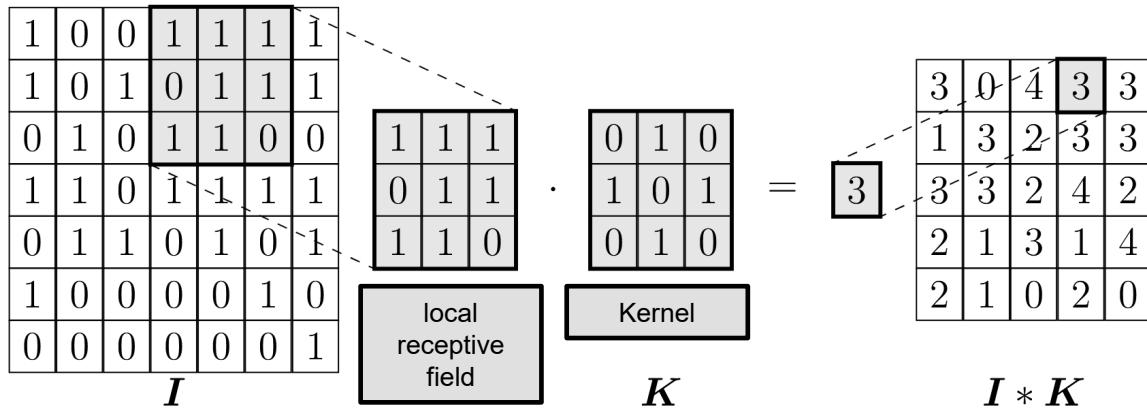
3.9.1 Convolutional Neural Networks

- Originally designed for image processing
- Identify relative positions of features
- Operate with convolutional layers

Convolutional layers

consist of **filters/kernels** $K(w)$ that are applied to an input field I and defined by

- Kernel size
- Number of kernels
- **Stride length**: how large is the shift of the kernel per step
- **Padding**: do you embed the input field?



In neural networks: The values in the kernel are not preset. They are learnt by the network and represent the weights.

3.9.1 Convolutional Neural Networks

1	0	0	1	1	1	1
1	0	1	0	1	1	1
0	1	0	1	1	0	0
1	1	0	1	1	1	1
0	1	1	0	1	0	1
1	0	0	0	0	1	0
0	0	0	0	0	0	1

$$V_{ij} \cdot K_{ij} = O_{kl}$$

3	0	4	3	3
1	3	2	3	3
3	3	2	4	2
2	1	3	1	4
2	1	0	2	0

The convolution $*$ is computed as a combination of the multidimensional dot product, also known as [tensordot](#) product:

$$o = V_{ij} \cdot K_{ij}$$

Example (2D [convolutional filter](#) applied with a [stride](#) length of $s = 1$ and no padding)

$$1 \cdot 0 + 1 \cdot 1 + 1 \cdot 0 + 0 \cdot 1 + 1 \cdot 0 + 1 \cdot 1 + 1 \cdot 0 + 1 \cdot 1 + 0 \cdot 0 = 3$$

This operation is then repeated by shifting the [view](#) V_{ij} with the stride s

$$O_{kl} = \sum_{i,j} V_{(i+k)(j+l)} \cdot K_{ij} = V * K$$

3.9.1 Convolutional Neural Networks – Filters

Example filters

Identity filter



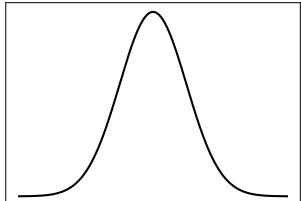
$$\mathbf{K} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}, s = 1$$

Gaussian filter

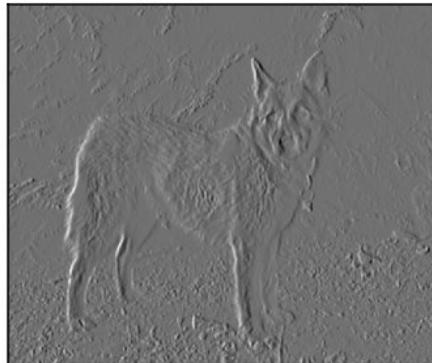


$$\mathbf{K} = \frac{1}{16} \begin{pmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{pmatrix}, s = 1$$

Gaussian Distribution



Sobel filter in x -direction



$$\mathbf{K} = \begin{pmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{pmatrix}, s = 1$$

Averaging filter



$$\mathbf{K} = \frac{1}{9} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}, s = 3$$

For other filters for gradient computation, see Chapter 9.

finite difference approximation convolutional filter

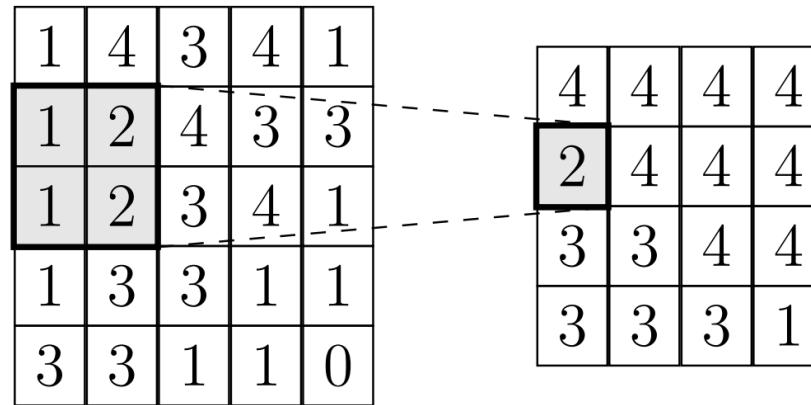
$$f'(x) = \frac{f(x+\Delta x) - f(x)}{\Delta x} \quad \frac{1}{\Delta x} [-1, 1]$$

$$f''(x) = \frac{f(x+\Delta x) - 2f(x) + f(x-\Delta x)}{\Delta x^2} \quad \frac{1}{\Delta x^2} [1, -2, 1]$$

3.9.1 Convolutional Neural Networks – Pooling

Pooling layers reduce the spatial size by extracting relevant features (filtering!)

- **Pooling types**
 - Max Pooling
 - Average Pooling
- **Pooling layers are defined by**
 - Kernel size
 - Stride length
 - Padding



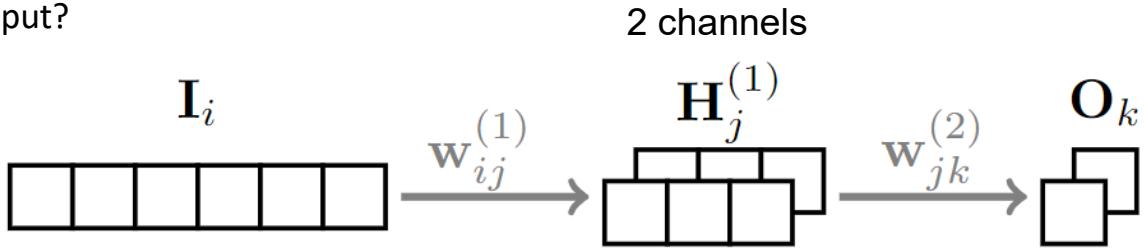
- Example “max pooling”(2D convolutional filter applied with a stride length of $s = 1$ and no padding)
 $\max\{1, 2, 1, 2\} = 2$

3.9.1 Convolutional Neural Networks – Example

What if we want two features O from one input?

Example with 1D convolutional filters

- Stride length $s = 1$
- Without a bias
- Input (given)



$$\mathbf{I} = [0, 1, 3, 2, 5]$$

- Weights from input layer \mathbf{I} to first hidden layer $\mathbf{H}^{(1)}$ (given)

$$\mathbf{w}_{11}^{(1)} = [0, 2, 1], \mathbf{w}_{12}^{(1)} = [1, 0, 3]$$

- Forward pass to first hidden layer $\mathbf{H}_j^{(1)} = \sum_i I_i * \mathbf{w}_{ij}^{(1)}$

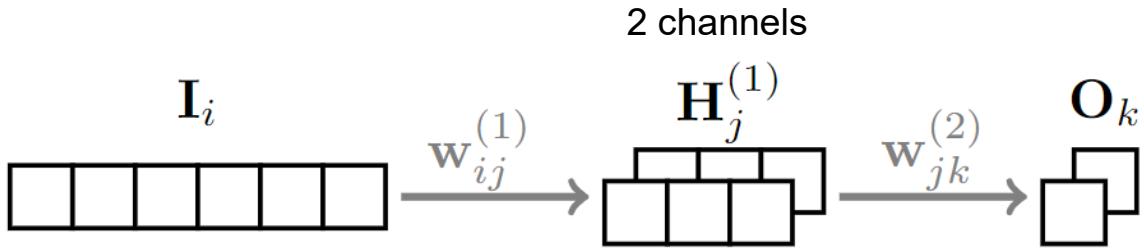
$$\mathbf{H}_1^{(1)} = [0 \cdot 0 + 1 \cdot 2 + 3 \cdot 1, 1 \cdot 0 + 3 \cdot 2 + 2 \cdot 1, 3 \cdot 0 + 2 \cdot 2 + 5 \cdot 1] = [5, 8, 9]$$

$$\mathbf{H}_2^{(1)} = [0 \cdot 1 + 1 \cdot 0 + 3 \cdot 3, 1 \cdot 1 + 3 \cdot 0 + 2 \cdot 3, 3 \cdot 1 + 2 \cdot 0 + 5 \cdot 3] = [9, 7, 18]$$

3.9.1 Convolutional Neural Networks – Example

Example with 1D convolutional filters

- Stride length $s = 1$
- Without a bias
- Hidden layer $\mathbf{H}^{(1)}$ (given)



$$\mathbf{H}_1^{(1)} = [5, 8, 9]$$

$$\mathbf{H}_2^{(1)} = [9, 7, 18]$$

- Weights from first hidden layer $\mathbf{H}^{(1)}$ to output layer $\mathbf{O}_k = \sum_j \mathbf{H}_j^{(1)} * \mathbf{w}_{jk}^{(2)}$ (given)

$$\mathbf{w}_{11}^{(2)} = [1, 0, 0], \mathbf{w}_{21}^{(2)} = [2, 1, 0]$$

$$\mathbf{w}_{12}^{(2)} = [0, 0, 3], \mathbf{w}_{22}^{(2)} = [0, 1, 1]$$

- Forward pass to output layer

$$\mathbf{O}_1 = [(5 \cdot 1 + 8 \cdot 0 + 9 \cdot 0) + (9 \cdot 2 + 7 \cdot 1 + 18 \cdot 0)] = 30$$

$$\mathbf{O}_2 = [(5 \cdot 0 + 8 \cdot 0 + 9 \cdot 3) + (9 \cdot 0 + 7 \cdot 1 + 18 \cdot 1)] = 52$$

3.9.1 Convolutional Neural Networks – Depth

ImageNet Classification with Deep Convolutional Neural Networks, (Alex) Krizhevsky et al. 2012

Comparison to fully connected neural network

- Channels are analogous to neurons
- Filters are analogous to linear transformations (connections)
- Number of parameters in fully connected neural network:

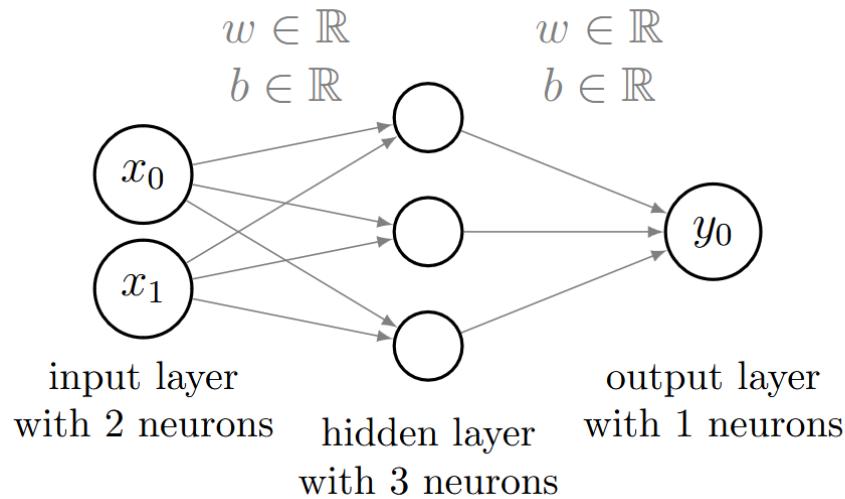
Number of parameters per weight/kernel

Number of connections (weights)

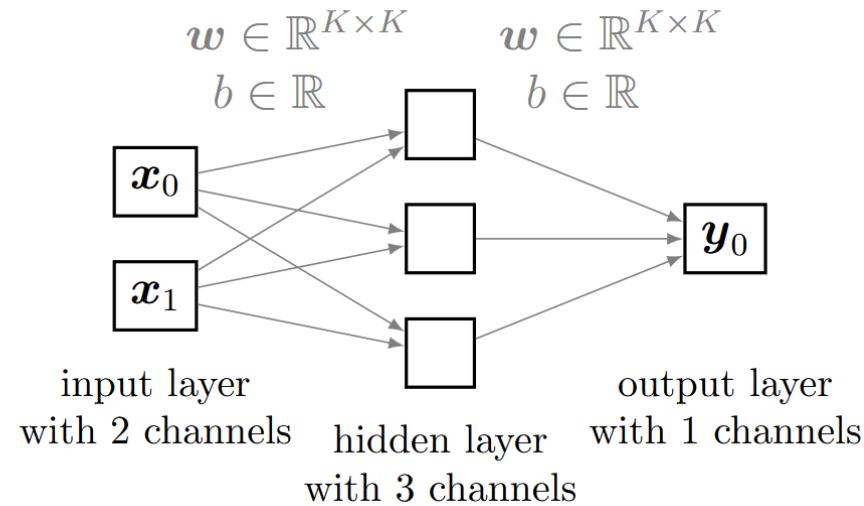
Number of biases

$$1 \cdot (2 \cdot 3 + 3 \cdot 1) + (3 + 1)$$

$$(K \cdot K) \cdot (2 \cdot 3 + 3 \cdot 1) + (3 + 1)$$



Fully connected neural network



Convolutional neural network

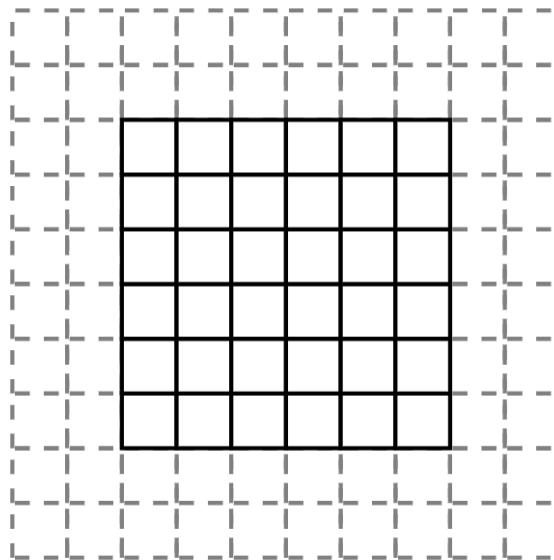
3.9.1 Convolutional Neural Networks – Padding

Padding counteracts the shrinking of the output, padding can be used

- Padding types
 - Valid padding (No padding)
 - Zero padding
 - Reflection padding
- The output dimension O after a convolutional layer

$$O = \left\lceil \frac{I - K + 2P}{s} \right\rceil + 1$$

- Input dimension I
- Kernel size K
- Padding size P
- Stride length s

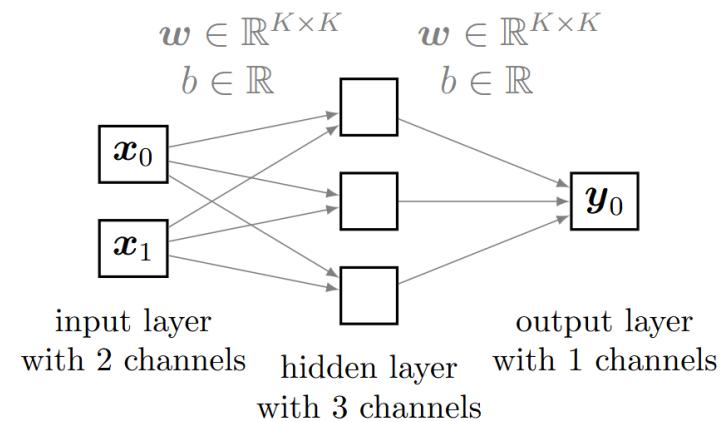


3.9.1 Convolutional Neural Networks – PyTorch

```
class CNN(torch.nn.Module):
    def __init__(self, inputImages, hiddenImages, outputImages):
        super(CNN, self).__init__()
        self.cnn1 = torch.nn.Conv2d(inputImages, hiddenImages, kernel_size=3,
                                  stride=1, padding=1)
        self.cnn2 = torch.nn.Conv2d(hiddenImages, outputImages, kernel_size=3,
                                  stride=1, padding=1)
        self.activation = torch.nn.ReLU()

    def forward(self, x):
        y = self.activation(self.cnn1(x))
        y = self.cnn2(y)
        return y

model = CNN(2, 3, 1) # input has size
                     (1, 2, imageSize, imageSize)
```



Exercises

E.11 Convolutional Neural Network (P & C)

- Familiarize yourself with convolutional neural networks through pen-and-paper forward propagation and through the application of different convolutional filters to images (using PyTorch).

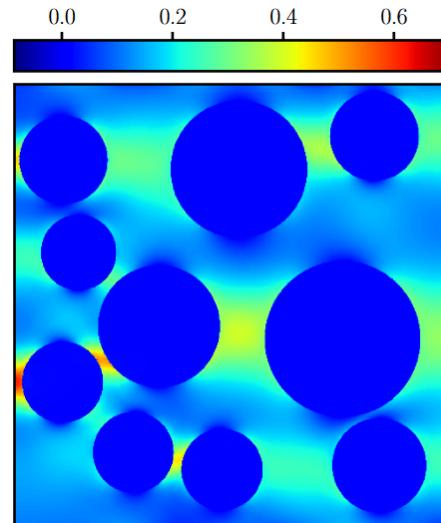
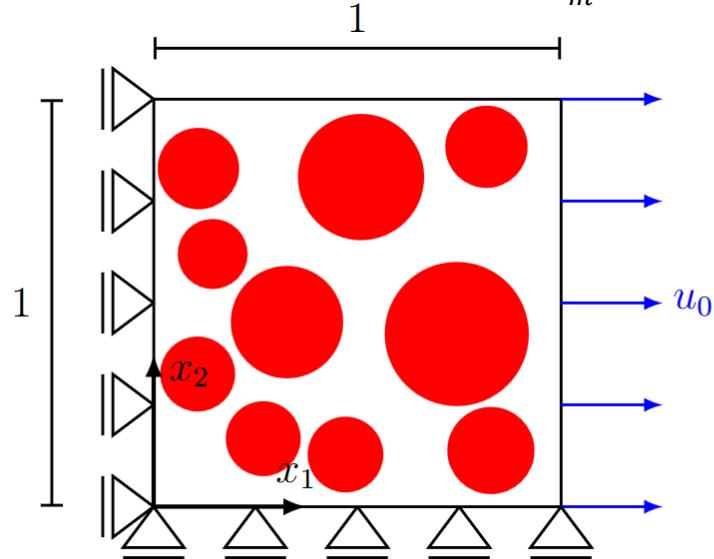
3.9.1 Convolutional Neural Networks – Example

Learning strain distributions from data

- Goal is to make the neural network predict the strain of a 2D domain under a uni-axial load
- Domain is defined by fiber material (red) and matrix material (white) with the elastic properties

$$E_f = 85'000, \nu_f = 0.22$$

$$E_m = 3'000, \nu_m = 0.4$$



3.9.1 Convolutional Neural Networks – Example

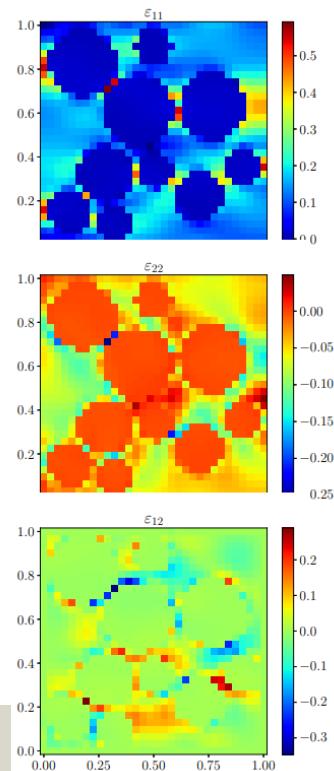
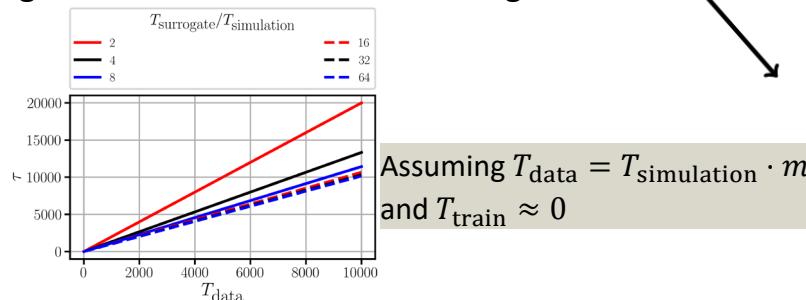
- To reduce the complexity of the task, the domain is discretized in 32×32 grid
- Strain prediction from the underlying material distribution
- As neural network architecture, a U-Net (see Chapter 8) is used

U-Net: Convolutional Networks for Biomedical Image Segmentation, Ronneberger et al. 2015

- When training a surrogate model, the breakeven threshold τ must be considered

$$\tau = \frac{T_{\text{data}} + T_{\text{train}}}{T_{\text{simulation}} + T_{\text{surrogate}}}$$

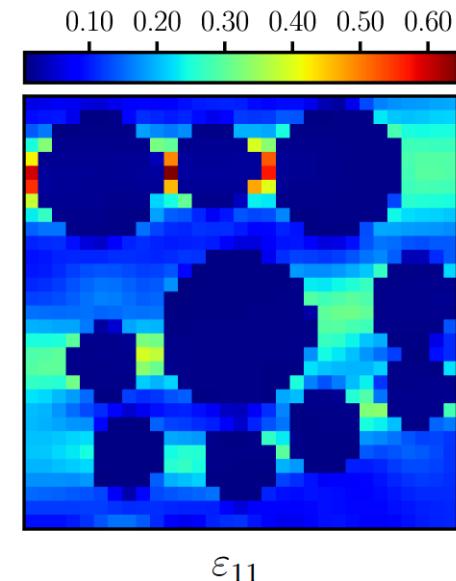
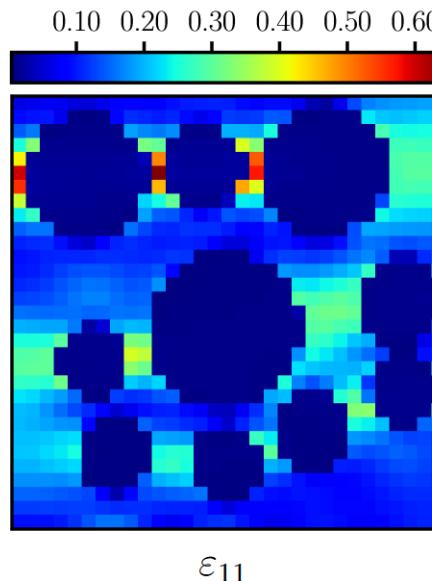
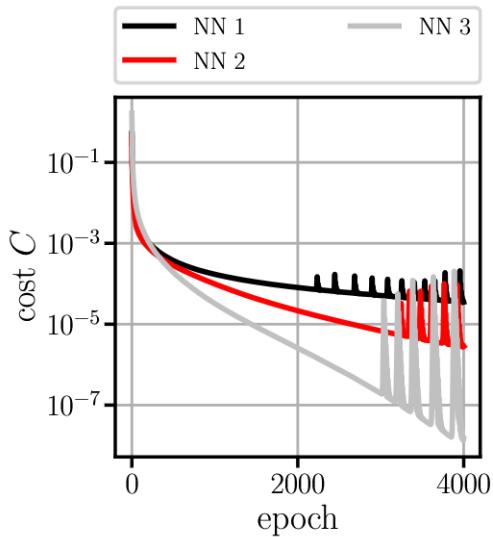
- Which represents the number of surrogate evaluations to make training worth it
- T is the time of
 - Data collection (T_{data})
 - Training (T_{train})
 - 1 Simulation ($T_{\text{simulation}}$)
 - 1 Surrogate evaluation ($T_{\text{surrogate}}$)



3.9.1 Convolutional Neural Networks – Example

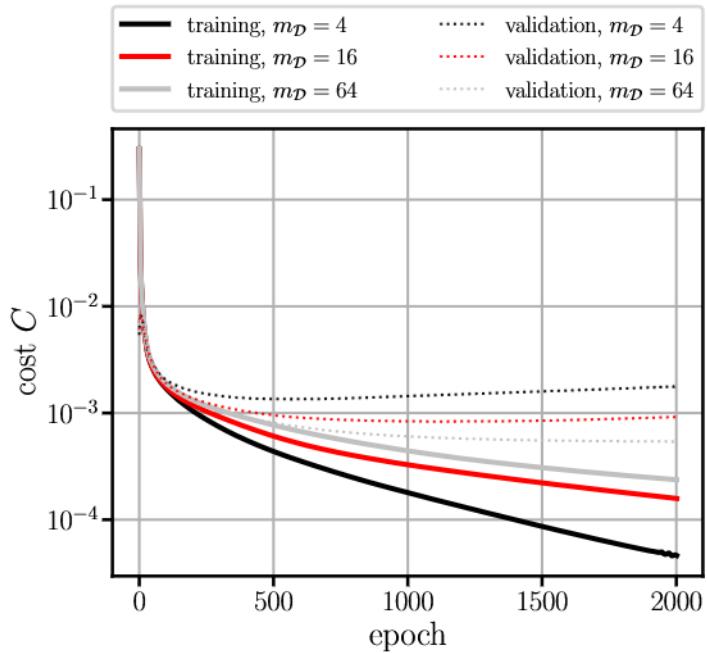
Model selection (make sure that model can learn the training data by heart)

- NN 1: U-net with 5 convolutional layers: 94'978 parameters
- NN 2: feed-forward convolutional neural network with 4 layers: 38'405 parameters
- NN 3: U-net followed by 3 convolutional layers: 154'471 parameters
- Training with 1 datapoint

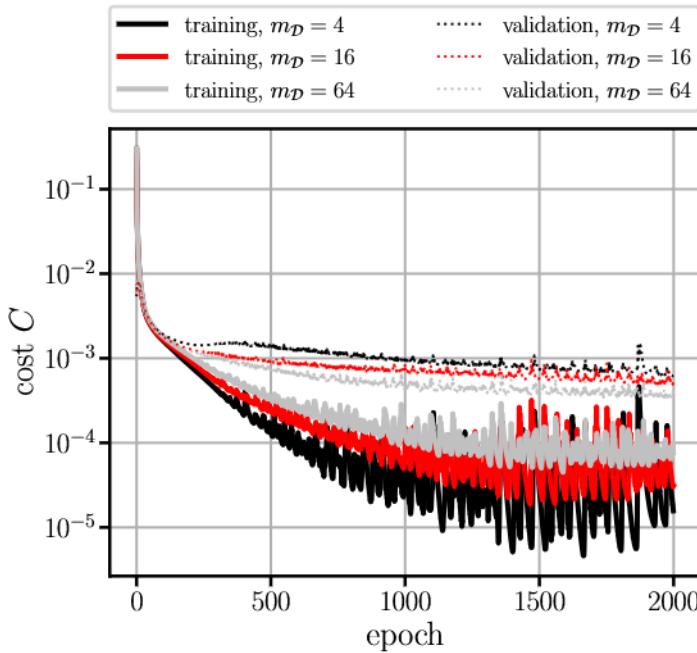


3.9.1 Convolutional Neural Networks – Example

With selected model (NN 2), increase the training data (and potentially tune the architecture & optimizer)



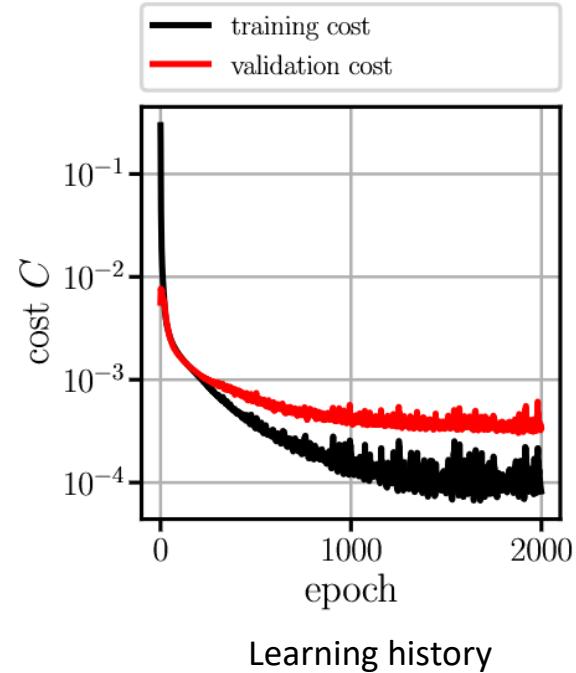
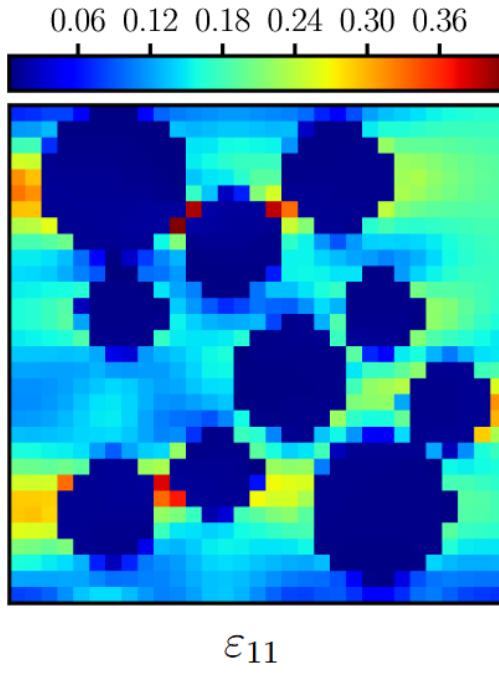
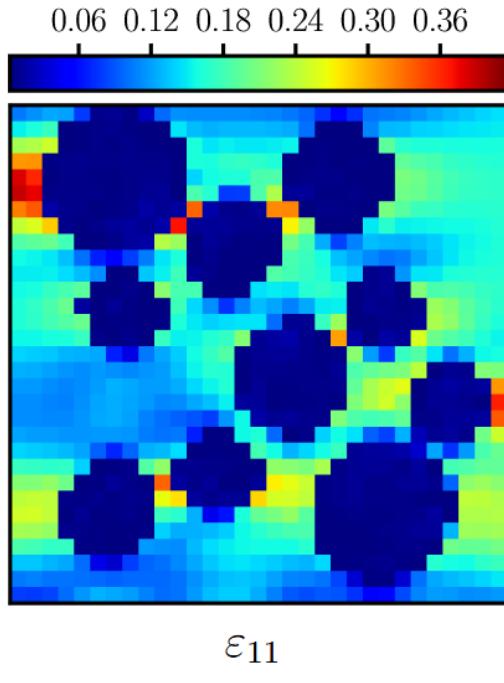
Without regularization



With regularization

3.9.1 Convolutional Neural Networks – Example

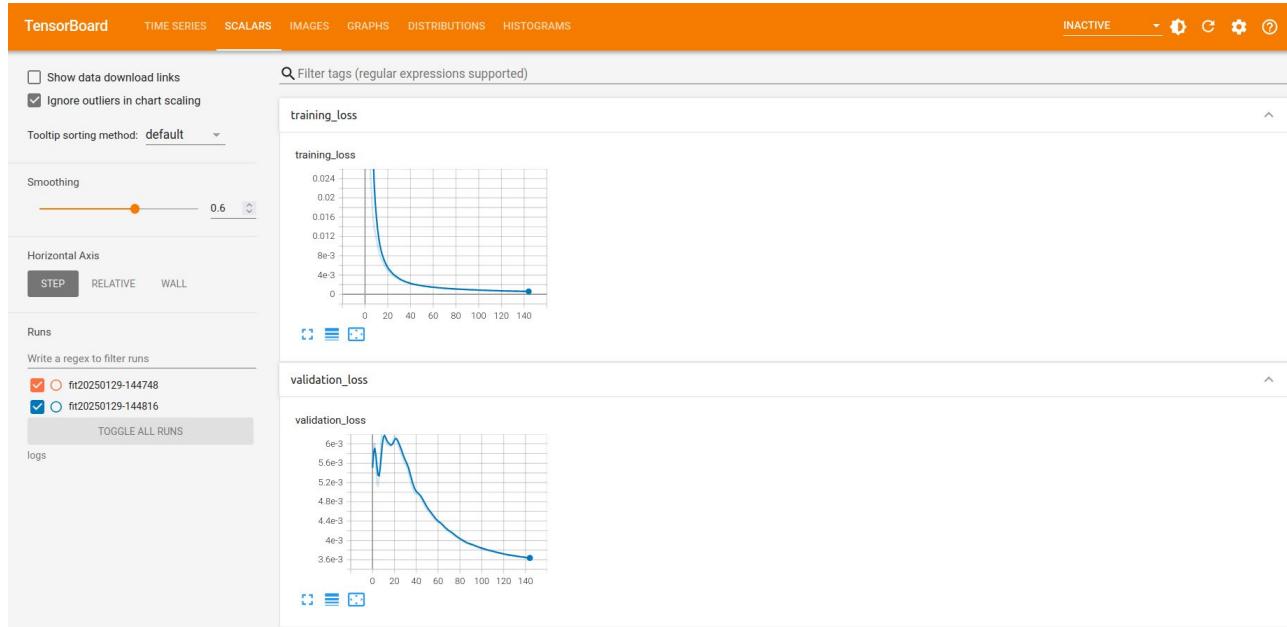
Prediction on unseen data (i.e., validation data)



3.9.1 Convolutional Neural Networks – TensorBoard

TensorBoard is a [visualization tool](#) for TensorFlow (& other machine learning libraries, e.g., PyTorch)

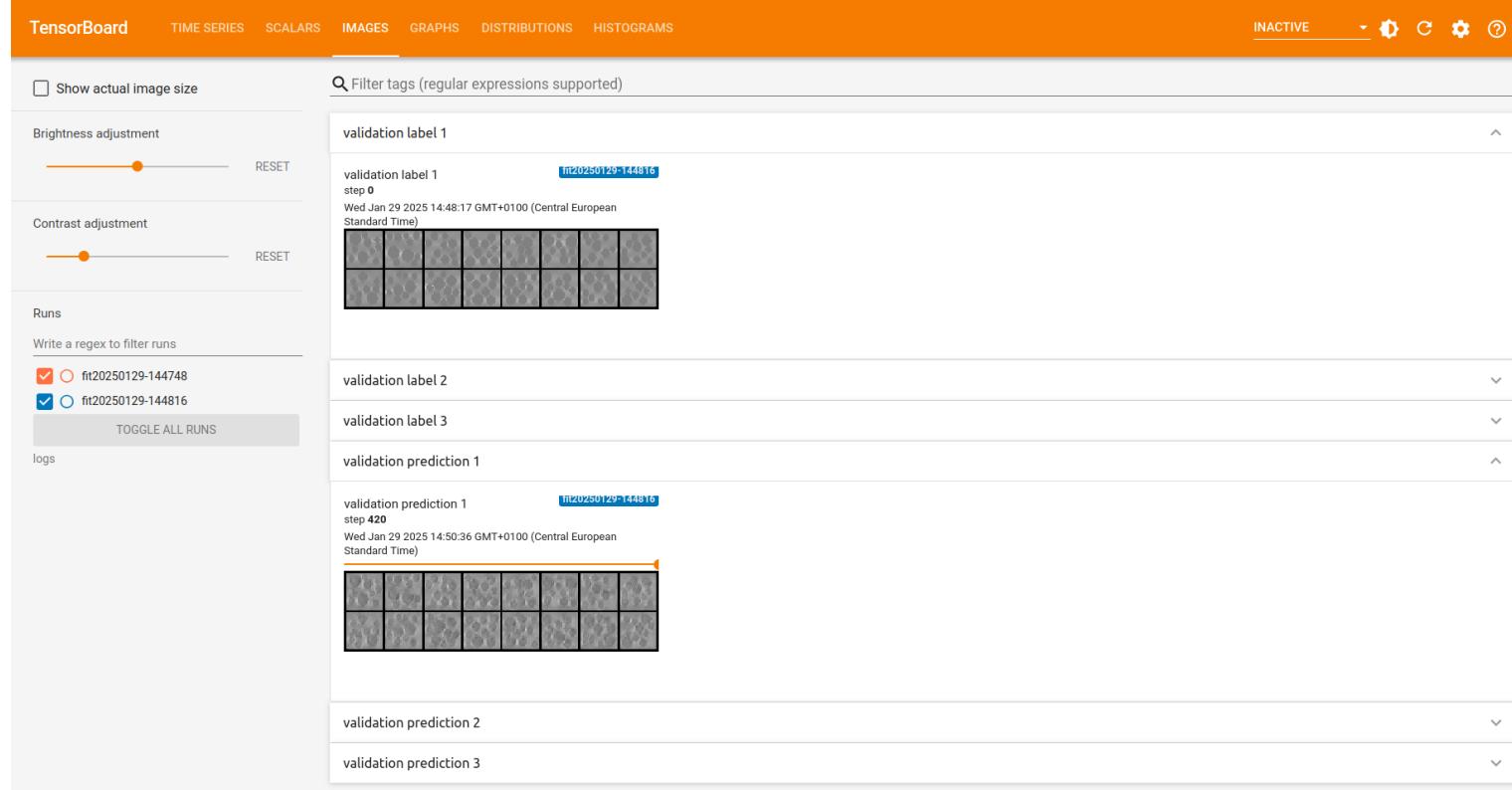
- Beneficial for [debugging](#) and [hyperparameter tuning](#)/[model selection](#)
- Includes features such as [graphs](#), [histograms](#), [output visualization](#)



TensorBoard

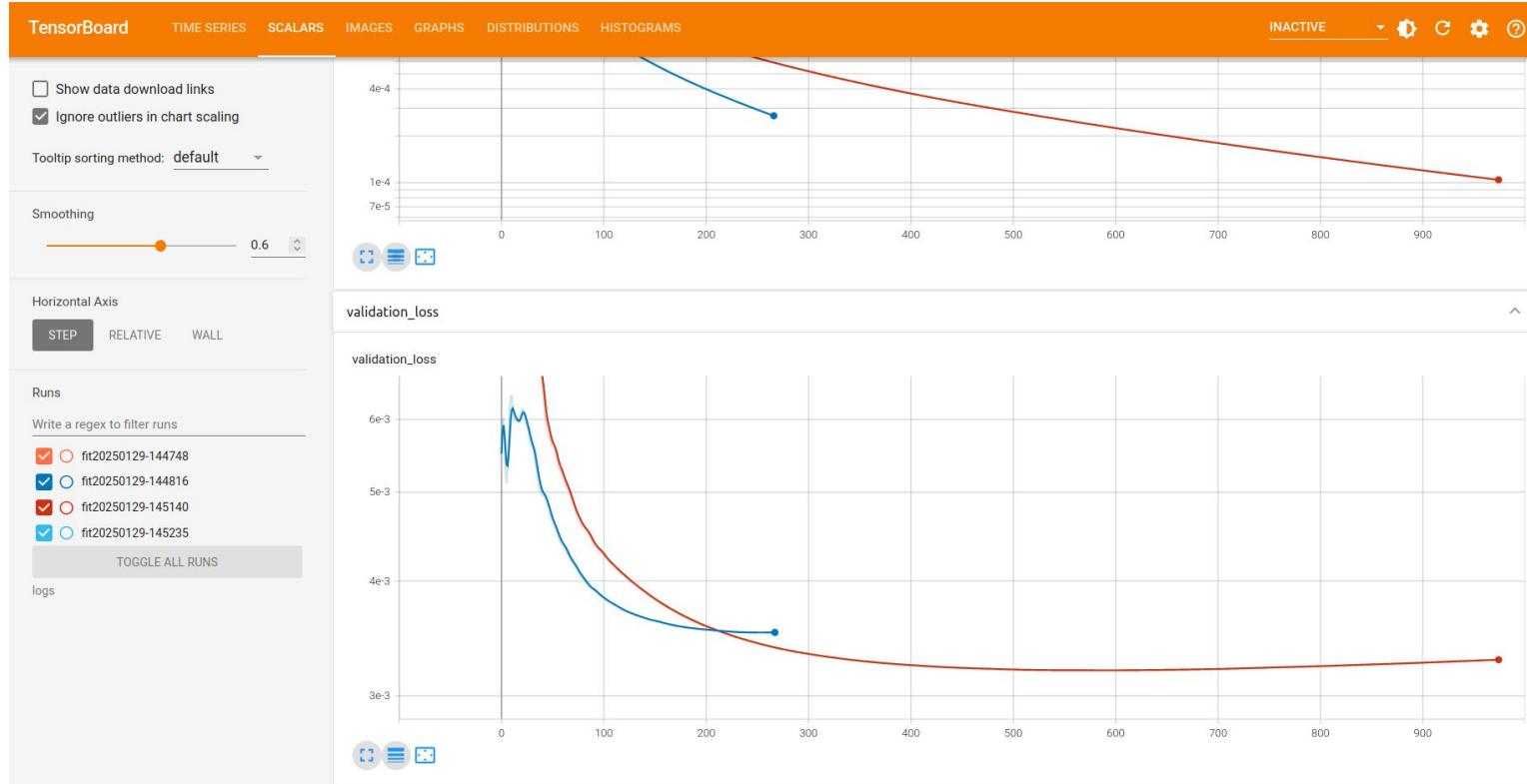
3.9.1 Convolutional Neural Networks – TensorBoard

Track the progress in prediction quality by comparing outputs with labels



3.9.1 Convolutional Neural Networks – TensorBoard

Compare different architectures and hyperparameters through different training runs



Exercises

E.12 Learning Strain Distributions (C)

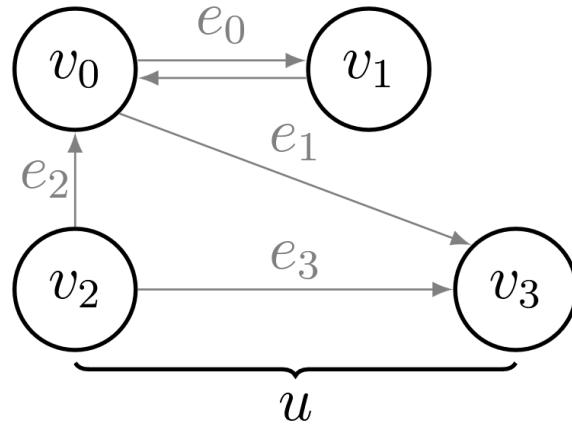
- Train a convolutional neural network as a surrogate, that maps the Young's modulus distribution to strain distributions. Experiment with different hyperparameters and different network architectures. For this, use TensorBoard.

3.9.2 Graph Neural Networks

- Convolutional neural networks are limited to data aligned on rectangular & uniform grids
- **Graph neural networks** are a generalization to general graphs

- **Graph**

- Nodes/vertices v_i
- Edges e_i
- Global attribute u



- Many different architectures exist, (message passing networks, graph convolutional networks, graph transformers...)
- We will focus on **message passing networks**, which consider the following invariant

$$v_i^s \xrightarrow{e_i} v_i^r$$

*Relational inductive biases, deep learning,
and graph networks, Battaglia et al. 2018*

- Each edge e_i has a sender v_i^s and receiver node v_i^r

3.9.2 Graph Neural Networks

- Edges are updated through the **first** fully connected neural network $e'_i = f^e(e_i, v_i^s, v_i^r, u)$
- Nodes are updated by **aggregating** connecting edges \bar{e}_i (e.g., max or sum): **second** fully connected neural network $v'_i = f^v(\bar{e}', v_i, u)$
- The global attribute is updated by **aggregating** all edges and nodes: **third** fully connected neural network $u' = f^u(\bar{e}', \bar{v}', u)$

Algorithm 8 Graph block in a message passing neural network

Require: Graph consisting of nodes \mathbf{v} , edges \mathbf{e} and a global attribute u .

```
for all edges  $e_i$  do
    Update edges:  $e'_i = f^e(e_i, v_i^s, v_i^r, u; \Theta_e)$ 
end for
```

```
for all nodes  $v_i$  do
    Find all edges connecting to node  $v_i$ :  $\bar{e}'_i$ 
    Aggregate adjacent edges:  $\bar{e}'_i = \rho^{e \rightarrow v}(\mathbf{e}'_i)$ 
    Update nodes:  $v'_i = f^v(\bar{e}'_i, v_i, u; \Theta_v)$ 
end for
```

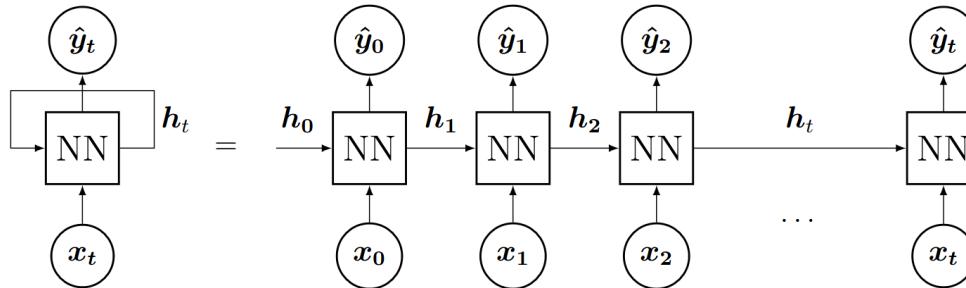
```
Aggregate all edges  $\mathbf{e}$ :  $\bar{e}' = \rho^{e \rightarrow u}(\mathbf{e})$ 
Aggregate all nodes  $\mathbf{v}$ :  $\bar{v}' = \rho^{v \rightarrow u}(\mathbf{v})$ 
Update global attribute:  $u' = f^u(\bar{e}', \bar{v}', u; \Theta_u)$ 
```

*Learning Mesh-Based Simulation with
Graph Networks, Pfaff et al. 2020*

3.9.6 Recurrent Neural Networks

Recurrent Neural Networks reuse information from previous states by

- looping over itself to generate sequences.



- Equivalent to using multiple copies of the same network with a **hidden state** h_t (unrolled recurrent neural network)

$$h_t = \sigma(W_x x_t + b_x + W_h h_{t-1} + b_h)$$

- W_x, W_h are the weights and b_x, b_h the biases
- y_t is obtained by a (learnable) mapping between h_t and y_t
- Useful for applications with **sequential data**, such as speech recognition, language modeling, translation
- In practice, only able to connect recent information with each other, e.g., "*We used too many lemons for our lemonade ... We did not like the lemonade, it was too sour*" poses a problem, as the cause for the is too far from the outcome
- Overcome by long short-term memory networks (LSTMs), gated recurrent units (GRUs), and transformers

3.9.6 Recurrent Neural Networks - PyTorch

```
class RNN(torch.nn.Module):
    def __init__(self, inputSize, hiddenStateSize, outputSize):
        super(RNN, self).__init__()
        self.rnn = torch.nn.RNN(inputSize, hiddenStateSize, nonlinearity='relu',
num_layers=2)
        self.linear = torch.nn.Linear(hiddenStateSize, outputSize)

    def forward(self, x):
        h, _ = self.rnn(x) # unrolls the RNN
        y = self.linear(h) # transforms hidden layer to output layer
        return y

inputSize = 1
hiddenStateSize = 500
outputSize = 1
model = CNN(inputSize, hiddenStateSize, outputSize)
```

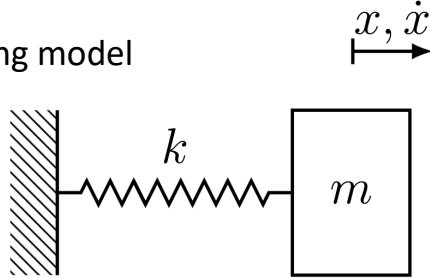
Recurrent neural networks do not only have to be combined with fully connected neural networks, convolutional neural networks are for example also possible

3.9.8 Physics-Inspired Architectures for Dynamics

- The previous architecture learn physical behavior only from data
- Governing laws can be incorporated into neural networks
 - Through training (via penalty terms in the cost function, i.e., **weak enforcement**/regularization; see Chapters 4 & 5)
 - Via the neural network architecture (by constraining the learnable space, i.e., **strong enforcement**; see Chapter 7)
- Strong enforcement of physics
 - Dynamics
 - **Hamiltonian neural networks**
 - **Lagrangian neural networks**
 - Constitutive modeling
 - Input-convex neural networks (see Chapter 7)
 - Boundary conditions
 - Strong enforcement of boundary conditions (see Chapters 4 & 5)

3.9.8 Physics-Inspired Architectures for Dynamics

Consider a one-dimensional mass-spring model



Described by the [ordinary differential equation](#)

$$m\ddot{x}(t) + kx(t) = 0$$

And the solution

$$x(t) = A\sin(\omega t + \phi)$$

Where $\omega = \sqrt{k/m}$ is the natural frequency and A, ϕ are determined by the initial conditions $x(0), \dot{x}(0)$

Alternatively the system can be described by

- [Hamiltonian Mechanics](#)
- [Euler-Lagrange Equation](#)

3.9.8.2 Hamiltonian Mechanics

In Hamiltonian mechanics, systems are described by coordinate pairs (position & momentum, i.e., $\mathbf{p}(t) = m\dot{\mathbf{x}}(t)$)
 $[\mathbf{x}(t), \mathbf{p}(t)]$

Scalar Hamiltonian $\mathcal{H}(\mathbf{x}(t), \mathbf{p}(t))$ represents the system's total energy Π_{tot} and fulfills

$$\frac{d\mathbf{x}}{dt} = \frac{\partial \mathcal{H}}{\partial \mathbf{p}},$$

$$\frac{d\mathbf{p}}{dt} = -\frac{\partial \mathcal{H}}{\partial \mathbf{x}}$$

Given the Hamiltonian, the time evolution of \mathbf{x}, \mathbf{p} can be computed via integration

$$\mathbf{x}(t + \Delta t) = \mathbf{x}(t) + \int_t^{t+\Delta t} \left. \frac{\partial \mathcal{H}}{\partial \mathbf{p}} \right|_{\tau} d\tau$$

$$\mathbf{p}(t + \Delta t) = \mathbf{p}(t) - \int_t^{t+\Delta t} \left. \frac{\partial \mathcal{H}}{\partial \mathbf{x}} \right|_{\tau} d\tau$$

Which can, e.g., be discretized with a forward Euler scheme

$$\mathbf{x}(t + \Delta t) = \mathbf{x}(t) + \left. \frac{\partial \mathcal{H}}{\partial \mathbf{p}} \right|_t \Delta t, \quad \mathbf{p}(t + \Delta t) = \mathbf{p}(t) - \left. \frac{\partial \mathcal{H}}{\partial \mathbf{x}} \right|_t \Delta t$$

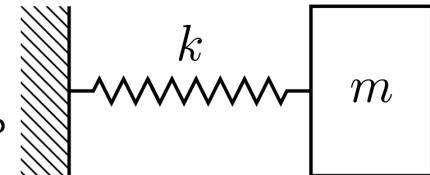
3.9.8.3 Hamiltonian Neural Networks

x, \dot{x}

For the one-dimensional mass-spring model, the **Hamiltonian** is known as

$$\mathcal{H} = \Pi_{\text{tot}} = \Pi_{\text{kin}} + \Pi_{\text{pot}} = \frac{1}{2}m\dot{x}^2 + \frac{1}{2}kx^2$$

But what if the Hamiltonian is unknown, but the system must obey Hamiltonian mechanics?



Hamiltonian neural networks

- **Data triplets** $\{\tilde{x}(t_i), \tilde{\dot{x}}(t_i), \tilde{\ddot{x}}(t_i)\}_{i=1}^m$ (enable computation of $\tilde{p}(t_i), \tilde{\dot{p}}(t_i)$)
- Mapping to the **Hamiltonian** is to be learned $\hat{\mathcal{H}} = \mathcal{H}(\mathbf{x}(t), \mathbf{p}(t); \Theta)$
- From $\hat{\mathcal{H}}$, $d\hat{x}/dt$ and $d\hat{p}/dt$ can be computed $\frac{dx}{dt} = \frac{\partial \mathcal{H}}{\partial p}, \frac{dp}{dt} = -\frac{\partial \mathcal{H}}{\partial x}$
- A comparison of $d\hat{x}/dt$ and $d\hat{p}/dt$ to the corresponding data in the data triplets, enables a **cost function**

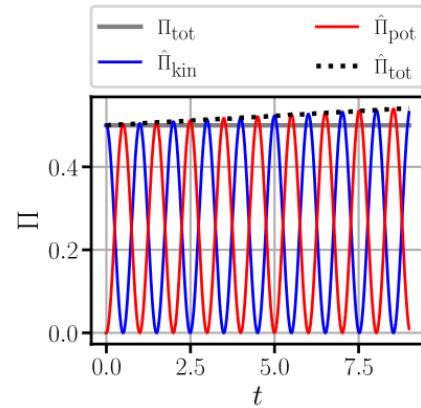
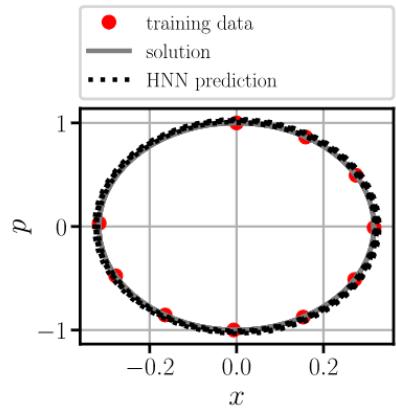
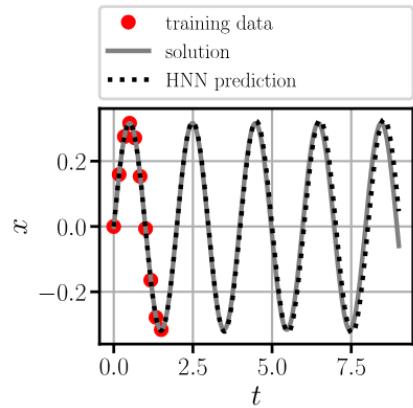
$$C = \frac{1}{m} \sum_{i=1}^m \left(\left\| \frac{d\tilde{x}(t_i)}{dt} - \frac{\partial \hat{\mathcal{H}}(\tilde{x}, \tilde{p}(t_i); \Theta)}{\partial \tilde{p}(t_i)} \right\|^2 + \left\| \frac{d\tilde{p}(t_i)}{dt} + \frac{\partial \hat{\mathcal{H}}(\tilde{x}(t_i), \tilde{p}(t_i); \Theta)}{\partial \tilde{x}(t_i)} \right\|^2 \right)$$

- Prediction of the trajectory $\hat{x}(t)$ beyond the datapoints seen in the data triplets, possible through **forward Euler**

$$x(t + \Delta t) = x(t) + \frac{\partial \mathcal{H}}{\partial p} \Big|_t \Delta t, \quad p(t + \Delta t) = p(t) - \frac{\partial \mathcal{H}}{\partial x} \Big|_t \Delta t$$

3.9.8.3 Hamiltonian Neural Networks - Results

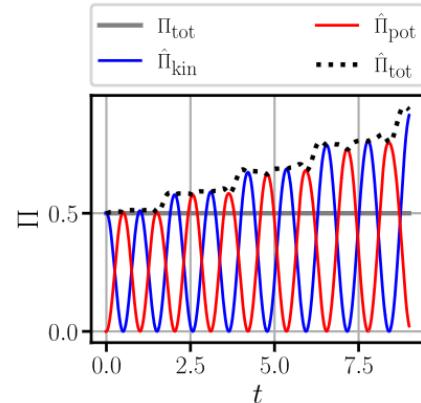
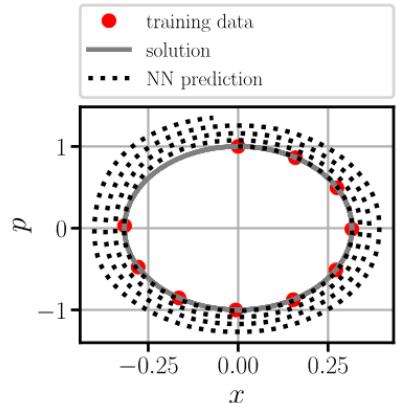
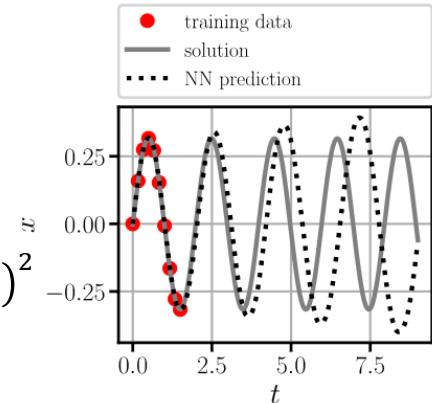
Hamiltonian neural network



Standard neural network

$$\begin{pmatrix} \hat{x} \\ \hat{p} \end{pmatrix} = f(x, p; \Theta)$$

$$C = \frac{1}{m} \sum_{i=1}^m (\tilde{x}_i - \hat{x}_i)^2 + (\tilde{p}_i - \hat{p}_i)^2$$



3.9.8.4 Euler-Lagrange Equation

The **Lagrangian framework** offers a more general framework than Hamiltonian mechanics

The **Lagrangian** of the mass-spring system is defined as

$$L = \Pi_{\text{kin}} - \Pi_{\text{pot}} = \frac{1}{2}m\dot{x}^2 - \frac{1}{2}kx^2$$

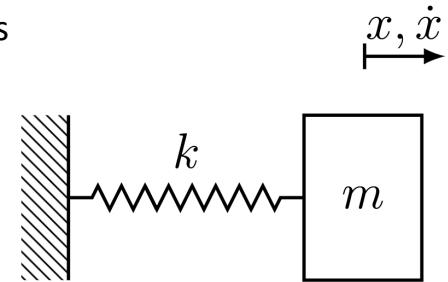
Which obeys the **Euler-Lagrange equation** (ensuring the principle of least action)

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{x}} = \frac{\partial L}{\partial x}$$

The acceleration \ddot{x} can be obtained by rewriting the Euler-Lagrange equation

$$\begin{aligned} \frac{d}{dt} \frac{\partial L}{\partial \dot{x}} &= \frac{\partial^2 L}{\partial x \partial \dot{x}} \dot{x} + \frac{\partial^2 L}{\partial \dot{x}^2} \ddot{x} = \frac{\partial L}{\partial x} \\ \ddot{x} &= \left(\frac{\partial^2 L}{\partial \dot{x}^2} \right)^{-1} \left(\frac{\partial L}{\partial x} - \frac{\partial^2 L}{\partial x \partial \dot{x}} \dot{x} \right) \end{aligned}$$

Similar as for the Hamiltonian mechanics, the evolution of $x(t)$ can be computed by **integrating** \dot{x}, \ddot{x}



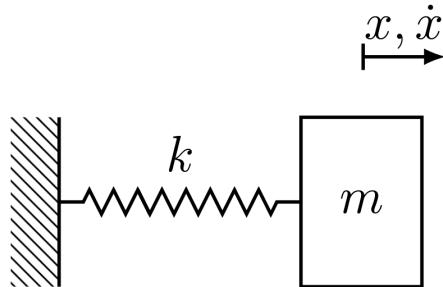
3.9.8.5 Lagrangian Neural Networks

- Data triplets $\{\tilde{x}(t_i), \dot{\tilde{x}}(t_i), \ddot{\tilde{x}}(t_i)\}_{i=1}^m$
- Mapping to the Lagrangian is to be learned $\hat{L} = L(x(t), \dot{x}(t); \Theta)$
- The acceleration \hat{x} can be computed from

$$\ddot{x} = \left(\frac{\partial^2 L}{\partial \dot{x}^2} \right)^{-1} \left(\frac{\partial L}{\partial x} - \frac{\partial^2 L}{\partial x \partial \dot{x}} \dot{x} \right)$$

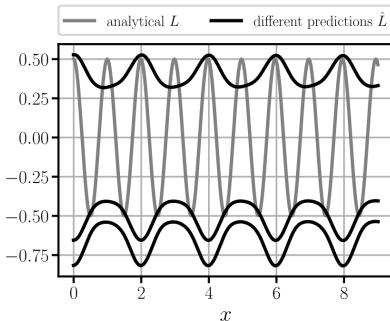
- With the predicted acceleration a cost function can be formulated as

$$C = \frac{1}{m} \sum_{i=1}^m \left\| \ddot{\tilde{x}}(t_i) - \hat{x}(\hat{L}; \tilde{x}(t_i), \dot{\tilde{x}}(t_i)) \right\|^2$$

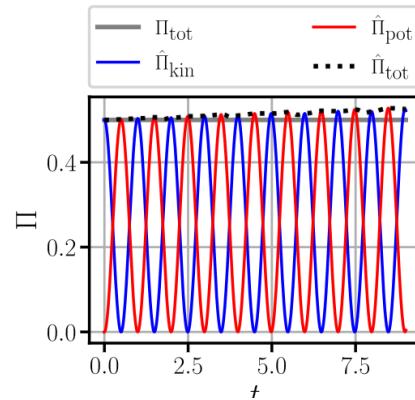
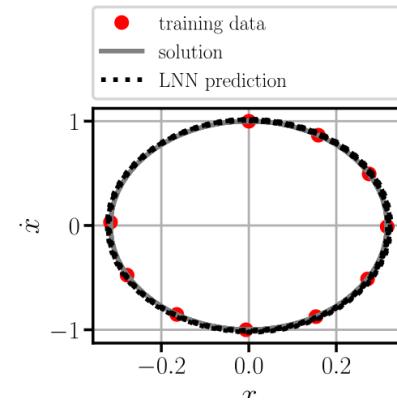
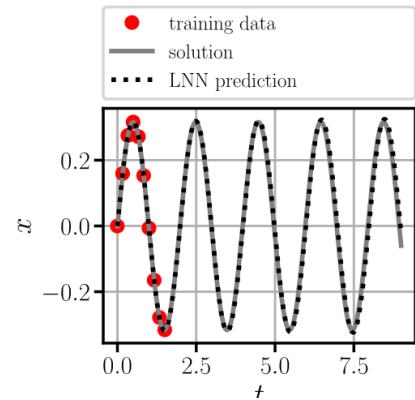


3.9.8.5 Lagrangian Neural Networks - Results

Lagrangian neural network



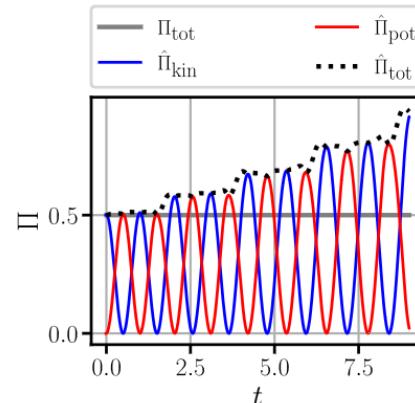
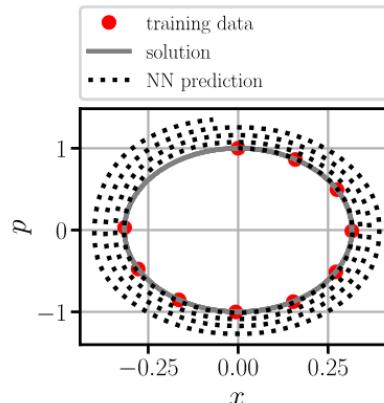
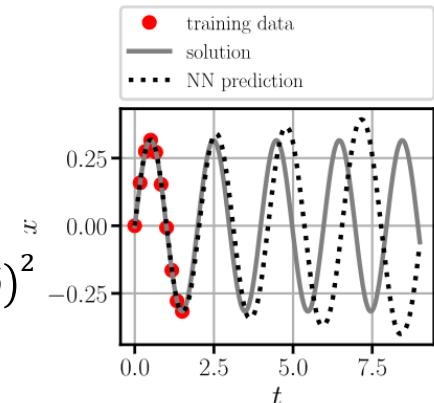
The Lagrangian is not unique



Standard neural network

$$\begin{pmatrix} \hat{x} \\ \hat{p} \end{pmatrix} = f(x, p; \Theta)$$

$$C = \frac{1}{m} \sum_{i=1}^m (\tilde{x} - \hat{x})^2 + (\tilde{p} - \hat{p})^2$$



Exercises

E.13 Hamiltonian & Lagrangian Neural Networks (C)

- Implement a Hamiltonian and a Lagrangian neural network and apply it to the mass-spring model. Compare the generalization capabilities to a standard neural network.

Contents

- [2 Fundamental Concepts of Machine Learning](#)
- [3.1 Fully Connected Neural Network](#)
- [3.4 Backpropagation](#)
- [3.6 Learning Algorithm](#)
- [3.8 Approximating the Sine Function](#)
- [3.9.1 Convolutional Neural Networks](#)
- [3.9.2 Graph Neural Networks](#)
- [3.9.6 Recurrent Neural Networks](#)
- [3.9.8 Physics-Inspired Architectures for Dynamics \(Hamiltonian & Lagrangian Neural Networks\)](#)
- [4 Introduction to Physics-Informed Neural Networks](#)

3 Neural Networks

Leon Herrmann

Stefan Kollmannsberger

Chair of Data Engineering in Construction

Bauhaus-Universität Weimar

Paris, February 2025

*Deep Learning in Computational Mechanics – an introductory course,
Herrmann et al. 2025*



website



book

