# Introduction to Deep Learning, PyTorch and TorchPhysics (A3)

## Part I: Introduction to Deep Learning

Daniel Otero Baguer, Janek Gödeke

Center for Industrial Mathematics (ZeTeM)
University of Bremen

Bremen
17-21.07.2023

Universität
Bremen

# Outline

Section 1

**Introduction**

Universität
Bremen

**Zentrum für**
**Technomathematik**

# Sources

- *"Deep Learning for Vision Systems"*, Mohamed Elgendy
  https://livebook.manning.com/book/
  deep-learning-for-vision-systems/welcome/v-8
- *"Convolutional Neural Networks for Visual Recognition"*
  http://cs231n.stanford.edu/

Universität
Bremen

**Zentrum für
Technomathematik**

# The learning task

Let's consider a simple regression problem where the observations are:

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), ..., (x^{(m)}, y^{(m)}) \in X \times Y \qquad (1)$$

**Task**: Infer a function $f$ from the training data, which minimizes

$$R(f) = \int_{X \times Y} c(x, y, f(x)) \mathrm{d}P(x, y) \qquad (2)$$

where $c$ is a loss function. A common choice is
$c(x, y, f(x)) = \frac{1}{2}\|f(x) - y\|^2$.

Universität
Bremen

**Zentrum für
Technomathematik**

# The learning task

Emprirical Risk

$$R(f) = \frac{1}{m} \sum_{i=1}^{m} c(x^{(i)}, y^{(i)}, f(x^{(i)})) \tag{3}$$

If we allow $f$ to be any function that maps from $X$ to $Y$, then we can minimize (5) but at the same time be very distant from the minimizer of (2).

$$f(x) = \begin{cases} y^{(i)} & x = x^{(i)} \\ 0 & otherwise \end{cases} \tag{4}$$

Universität
Bremen

# The learning task

## No Free lunch theorem

If we do not make any assumptions on the class of functions where $f$ belongs to, there is no chance to learn anything.

# Learning from data with Neural Networks (NN)

Given observations: $\{(x^{(m)}, y^{(m)})\}$ and model $\varphi_\theta : X \to Y$ (NN), minimize empirical error:

$$\arg\min_{\theta \in \Theta} = \frac{1}{m} \sum_{i=1}^{m} c(x^{(i)}, y^{(i)}, \varphi_\theta(x^{(i)})) \tag{5}$$

# Learning from data with Neural Networks (NN)

## Learning representations by back-propagating errors

David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams
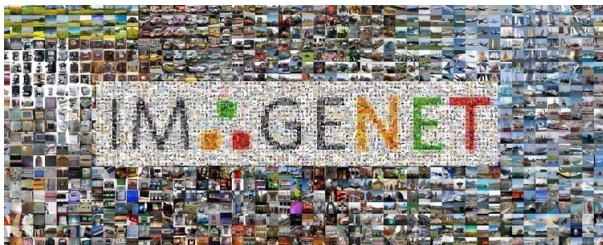
Universität
Bremen

# Why Deep learning now?

- Five decades of research in machine learning
- CPUs/GPUs/storage developed for other purposes
- Lots of data from the internet
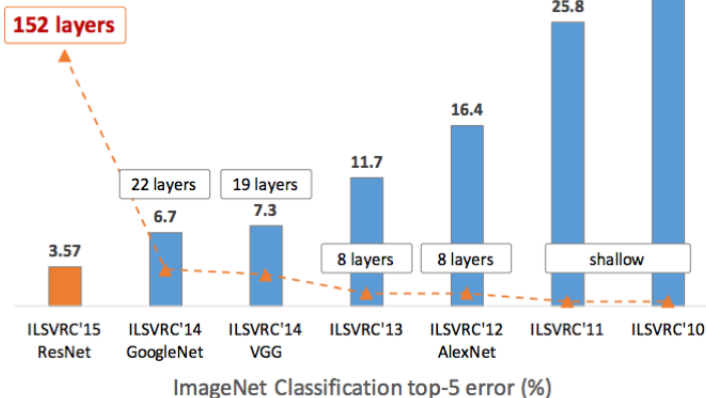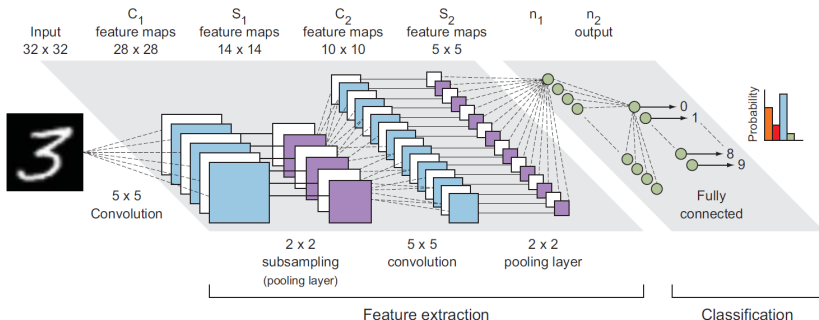- Imagenet

# Why Deep learning now?

- Five decades of research in machine learning
- CPUs/GPUs/storage developed for other purposes
- Lots of data from the internet
- Imagenet

**Zentrum für**
**Technomathematik**

# Why Deep learning now?

- Five decades of research in machine learning
- CPUs/GPUs/storage developed for other purposes
- Lots of data from the internet
- Imagenet

Universität
Bremen

**Zentrum für**
**Technomathematik**

# Why Deep learning now?

- Five decades of research in machine learning
- CPUs/GPUs/storage developed for other purposes
- Lots of data from the internet
- Imagenet

Universität
Bremen

Zentrum für
**Technomathematik**

# Why Deep Learning now?

Universität
Bremen

Section 2

**Convolutional Neural Networks**

Universität
Bremen

# Convolutional Neural Network

# Layer types

- Input



28 x 28
= 784 pixels

# Layer types

- Output

# Layer types

- Convolutions

Universität
Bremen

# Layer types

- Convolutions



Figure: Receptive field

# Layer types

- Convolutions



Figure: Convolutional layers are inspired on standard image filtering

# Layer types

- Pooling



Figure: Max pooling example

Universität
Bremen

Zentrum für

**Technomathematik**
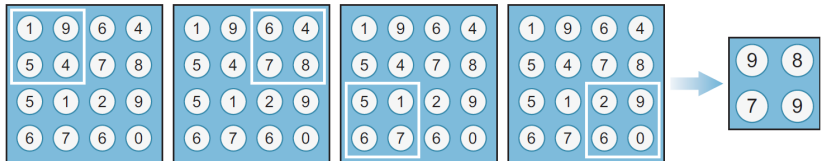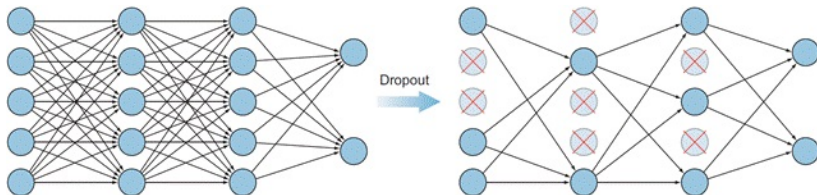
# Layer types

- Dropout



**Dropout: A Simple Way to Prevent Neural Networks from Overfitting**

*Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, Ruslan Salakhutdinov*; 15(56):1929−1958, 2014.

Universität
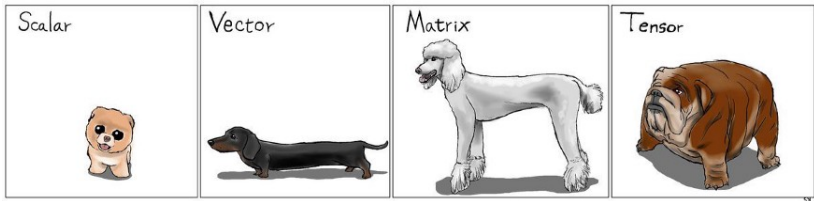Bremen

# Tensors



Figure: Tensors [1]

---

[1] Image source: https://towardsdatascience.com/understanding-pytorch-with-an-example-a-step-by-step-tutoria

Universität
Bremen

Section 3

**Architectures**

Universität
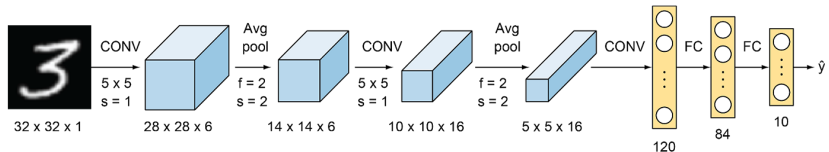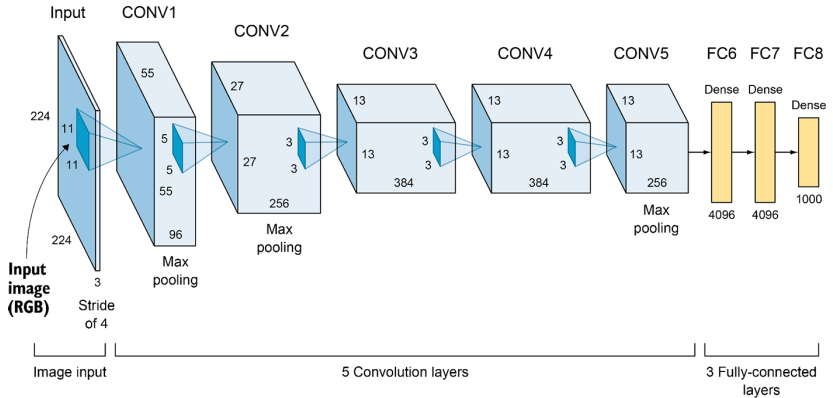Bremen

# LeNet

- Introduced in 1998, by LeCun et al. in their paper "Gradient-Based Learning Applied to Document Recognition"
- 5 layers: 3 conv + 2 fully-connected
- 61706 parameters

Universität Bremen

## AlexNet

- Winner of the ILSVRC image classification competition in 2012.
- Introduced by Alex Krizhevsky, Geoffrey Hinton and Ilya Sutskever in their paper "ImageNet Classification with Deep Convolutional Neural Networks"
- 8 layers: 5 conv + 3 fully-connected
- 60 million parameters

# AlexNet

# How deep can we go?

- Very deep networks are able to represent very complex functions
- The network can learn features at many different levels of abstraction

## Vanishing gradients

- By the chain rule, the derivatives of each layer are multiplied down the network
- Gradient decreases exponentially as we propagate down to the initial layers
- First hidden layers are learning much slower than later hidden layers
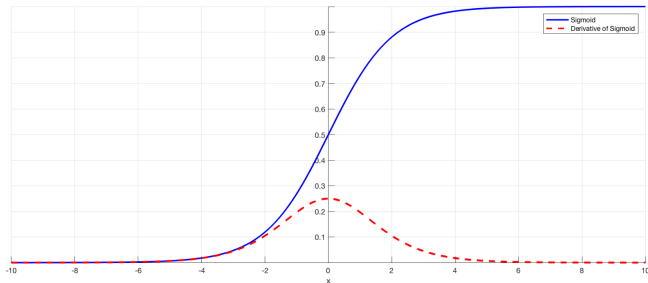
# How deep can we go?



Figure: Sigmoid activation

Universität
Bremen

# How deep can we go?



Figure: Vanishing gradient effect. First layers train much slower.

Universität
Bremen

# How deep can we go?

**Solutions:**

- Sigmoid with restricted inputs
- Use ReLu activations
- Normalization layers
- Residual Networks

**Zentrum für
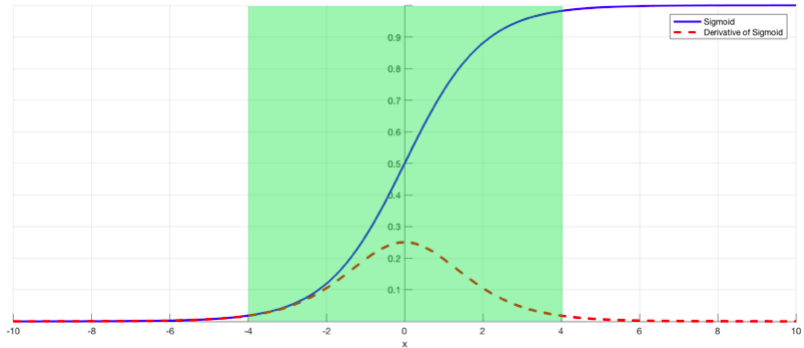Technomathematik**

# How deep can we go?



Figure: Sigmoid with restricted inputs

Universität
Bremen

# How deep can we go?



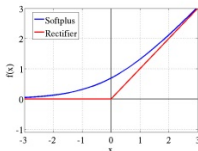Figure: ReLU (Rectifier linear unit)

Deep Sparse Rectifier Neural Networks

**Xavier Glorot**
DIRO, Université de Montréal
Montréal, QC, Canada
glorotxa@iro.umontreal.ca

**Antoine Bordes**
Heudiasyc, UMR CNRS 6599
UTC, Compiègne, France
and
DIRO, Université de Montréal
Montréal, QC, Canada
antoine.bordes@hds.utc.fr

**Yoshua Bengio**
DIRO, Université de Montréal
Montréal, QC, Canada
bengioy@iro.umontreal.ca

# How deep can we go?

Residual networks (ResNet)

- Introduces shortcut (skip-connection) that allows the gradient to directly back-propagate to earlier layers
- Allows the layer to learn an identity function (will perform at least as good as the previous layer)



Figure: ResNet's skip connection

Zentrum für
Technomathematik
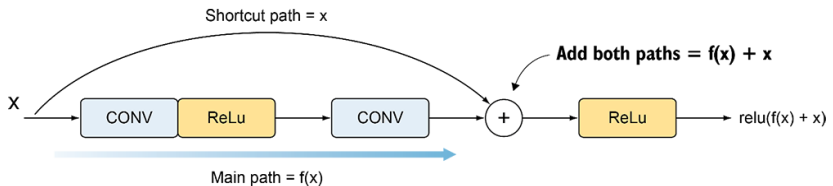
# How deep can we go?

- From 18 to 152 layers
- With 152 layers: $\approx$ 60 million parameters



Figure: Residual block

Universität
Bremen

# How deep can we go?



Figure: ResNet

Section 4

**Data augmentation**

Universität
Bremen

**Zentrum für**
**Technomathematik**

# Data augmentation

**Main idea**: Generate additional
images based on existing ones

- Random flip
- Random erasing
- Random sized crop
- Random perspective
- Color jitter



Figure: Data augmentation

Universität
Bremen

# Data augmentation

### Insights

- The epoch size does not change
- We get different randomly transformed samples every epoch (dynamic dataset)
- Helps to avoid overfitting

### Advise

- Applying heavy augmentations unnecessarily can result in poor accuracy
- Do not use data augmentation in validation or test set

# Data augmentation



Figure: Non-augmented model vs augmented model [2].

---

Universität
Bremen

Section 5

**Transfer learning**

Universität
Bremen

**Zentrum für**
**Technomathematik**

## Transfer learning

Common phenomena on CNNs trained on natural images:



Figure: First layer learn features similar to Gabor filters and color blobs. Each of the 96 filters shown here is of size $[11 \times 11 \times 3]$,

**Claim:** Such a layer is not specific to a particular dataset or task

Universität
Bremen

# Transfer learning

**Main idea:** Use a pre-trained network
**Use cases**

- Fixed feature extractor: Take a pre-trained network and train only the last fully-connected layer
- Fine-tuning: Fine-tune the weights of the pre-trained network by continuing the back-propagation
    - Use small learning rate
    - Keep some of the earlier layers fixed (due to overfitting concerns)

**Zentrum für**
**Technomathematik**

# When and how?

|                | Small dataset                         | Large dataset                      |
| -------------- | ------------------------------------- | ---------------------------------- |
| Similar        | Fixed feature extractor + Linear Classifier | Fine-tuning                  |
| Very different | Keep early layers fixed + Linear classifier | Pre-trained as initialization |

Universität
Bremen

# Practical advice

- Constraints from pre-trained models
- Small learning rate for fine-tuning (avoid distorting pre-trained parameters too quickly or too much)

Universität
Bremen

# Turing award

Zentrum für
**Technomathematik**

Thanks!

Universität
Bremen