# Neural Network Building Blocks

Johannes Erdmann RWTH Aachen University

Deep Learning Basics Train-the-Trainer Workshop February 7-8, 2023 Garching

#### Heisenberg-Programm

**DFG** Deutsche Forschungsgemeinschaft







Numerous Al applications in our every-day lives

## Artificial Intelligence





## **Artificial Intelligence**

#### Numerous Al applications in our every-day lives

Machine learning is enabling fascinating capabilities

Domair

ublic

Δ\_

Most recent examples:

#### ChatGPT

#### Was Jimmy Wales killed during the Tiananmen Square Massacre?

No, Jimmy Wales was not killed during the Tiananmen Square Massacre. Jimmy Wales is a 🍡 🗅 🖓 nt event that took place in Beijing, China on June 4, 1989, when Chinese ation in Tiananmen <u>Square. Thousands of</u> people were killed or injured during the massacre, which was a significant and highly controversial event in modern Chinese history. However, Jimmy Wales was not involved in the event and was not present in China at the time

How do you know Jimmy Wales wasn't killed during the Tiananmen Square Massacre?

know that Jimmy Wales was not killed during the Tiananmen Square Massacre because he  $\quad$  🖒  $\; arphi \;$ s a co-founder of Wikipedia, and my training data includes information about him and his ife. According to Wikipedia (which was written by a community of volunteers, including immy Wales), he was born in 1966 and is still alive today. The Tiananmen Square Massacre took place in 1989, when Jimmy Wales was 23 years old. Therefore, it is highly unlikely that ne was present in China at the time of the event or that he was killed during the massacre.

Please reply again, in the form of a short Shakespearean sonnet.

To say that Jimmy Wales did not die In the tragic Massacre at Tiananmen Square Is easy, for he lives and thrives on high Co-founding Wikipedia, beyond comm 😘 Regenerate response

Version. Free Research Preview. Our goal is to make Al systems more natural and safe to interact with. Your feedback will





- Data + Task + "a mathematical formulation of what means 'better'" \* (such as the number of type-1 and type-2 errors in a classification task)
- Optimization of the model parameters on the data
  - = "Training" (aka the machine 'learns')



\* often called a "performance measure"





- Data + Task + "a mathematical formulation of what means 'better'" \* (such as the number of type-1 and type-2 errors in a classification task)
- Optimization of the model parameters on the data
  - = "Training" (aka the machine 'learns')



\* often called a "performance measure"









- Data + Task + "a mathematical formulation of what means 'better' "\* (such as the number of type-1 and type-2 errors in a classification task) Optimization of the model parameters on the data Task
- - = "Training" (aka the machine 'learns')



\* often called a "performance measure"

Model with parameters

Model Output



Metric(s)





- - = "Training" (aka the machine 'learns')



\* often called a "performance measure"





# **Training (Parameter Optimization)**

• Training (on "Training Data")







# **Training (Parameter Optimization)**

Training (on "Training Data")











#### • ML type closely related to the available data for the task

# Supervised Learning

- Labelled data
- Classification

(image classification, ...)

Regression

(market forecasting, ...)

#### **Different Tasks**



### **Different Tasks**

#### • ML type closely related to the available data for the task

# Supervised Learning

Labelled data

Classification

(image classification, ...)

Regression

(market forecasting, ...)

- Unlabelled data
- Clustering
  - (recommender systems, ...)
- Dimensionality Reduction

Unsupervised Learning

(data compression, ...)



#### **Different Tasks**

#### • ML type closely related to the available data for the task

# Supervised Learning

Labelled data

Classification

(image classification, ...)

Regression

(market forecasting, ...)

- Unlabelled data
- Clustering
  - (recommender systems, ...) an environment
- Dimensionality Reduction
   Robot Navigation,

Unsupervised Learning

# Reinforcement Learning

(data compression, ...)

- An "agent" receives
  - rewards for actions in

Games, ...





- is due to "deep learning"
  - = training of deep neural networks



# **Deep Learning**





- is due to "deep learning"



# **Deep Learning**





- is due to "deep learning"



# **Deep Learning**

















$$v_{ij}^{(k-1)} f_j^{(k-1)}$$





$$v_{ij}^{(k-1)} f_j^{(k-1)} + b_i^{(k-1)}$$













 $f_i^{(k)} = \sum w_{ij}^{(k-1)} f_j^{(k-1)} + b_i^{(k-1)}$  $\vec{x}; \vec{\theta}$  $\vec{f}\left(\vec{x}; \{\mathbf{W}\}, \{\vec{b}\}\right)$  $\vec{f}^{(k)} = \mathbf{W}^{(k-1)} \vec{f}^{(k-1)} + \vec{b}^{(k-1)}$ 







# Non-linear activation functions are really key.

 $\vec{x}$  –



## The Output Function(s)

- Output function(s) = activation(s) in the output layer
- Closely connected to loss (and ultimately the task)



**softmax** of output node *i* for classification of *N* classes





### **The Loss Function\***

\* aka "cost function" aka "objective function"



- How good is my current model?
- How to change the parameters to improve it?
  - Loss function must be easily differentiable!
- Regression:

Mean Squared Error (MSE)

Model with parameters

Model Output

# Loss function

 $-\vec{f}\left(\vec{x}_i; \{\mathbf{W}\}, \{\vec{b}\}\}\right)$  $ec{y_i}$  $n \sum_{i=1}^{n}$ 



### **The Loss Function\***

\* aka "cost function" aka "objective function"



- How good is my current model?
- How to change the parameters to improve it?
  - Loss function must be easily differentiable!
- Classification:
  - Cross-entropy of prediction f (= interpreted
  - as probability q) and the true probability p
  - (= the classification labels y)

# $-\frac{1}{n}\sum_{i=1}^{n}\left|\sum_{\text{classes}}y_{i,\text{class}}\cdot\log f_{\text{class}}\left(\vec{x}_{i};\left\{\mathbf{W}\right\},\left\{\vec{b}\right\}\right)\right|$ $p_{i,\text{class}} \log q_{i,\text{class}}$





### **The Loss Function\***

\* aka "cost function" aka "objective function"



- How good is my current model?
- How to change the parameters to improve it?
  - Loss function must be easily differentiable!
- Classification:
  - Cross-entropy of prediction f (= interpreted
  - as probability q) and the true probability p
  - (= the classification labels y)

Model with parameters



# Loss function

The proveries  $\mathbf{f}_{i=1}^{n}$  for binary classification for binary classification  $-\frac{1}{n}\sum_{i=1}^{n} \begin{cases} \log q_{i}, & \text{signal} \\ \log (1-q_{i}), & \text{background} \end{cases}$ The prove of  $\mathbf{f}_{i=1}^{n}$  and  $\mathbf{f}_{i=1}^{n}$  and



## **Optimization via Backpropagation**

- We need the gradient w.r.t. each optimizable parameter
- Weights to last layer (output node):





## **Optimization via Backpropagation**





## **Optimization via Backpropagation**





 Ideal: Converge to the parameters of the global minimum of the loss function Practical: Converge to a "good-enough" local minimum

- Bad: Converge to "some" local minimum

- SGD: gradient estimated on a mini batch of training data
  - Variance in gradient estimates

 $\theta \to \theta - \alpha \left\langle \frac{\partial \text{Loss}}{\partial \theta} \right\rangle$ 

• Learning rate  $\alpha$ :



Ш

 Ideal: Converge to the parameters of the global minimum of the loss function Practical: Converge to a "good-enough" local minimum

- Bad: Converge to "some" local minimum

- SGD: gradient estimated on a mini batch of training data
  - Variance in gradient estimates

 $\theta \to \theta - \alpha \left\langle \frac{\partial \text{Loss}}{\partial \theta} \right\rangle$ 

• Learning rate  $\alpha$ :





#### Choosing a good learning rate is key J(w) for training convergence

[https://playground.tensorflow.org, github: https://github.com/tensorflow/playground, Apache License 2.0]

# **Training Strategies**







# Choosing a good learning rate is key for training convergence









# Choosing a good learning rate is key for training convergence









#### Choosing a good learning rate is key J(w) for training convergence



![](_page_37_Figure_4.jpeg)

![](_page_37_Figure_6.jpeg)

![](_page_37_Picture_7.jpeg)

# Choosing a good learning rate is key for training convergence

![](_page_38_Figure_2.jpeg)

![](_page_38_Figure_4.jpeg)

![](_page_38_Figure_5.jpeg)

![](_page_38_Picture_6.jpeg)

#### Choosing a good learning rate is key J(w) for training convergence

![](_page_39_Figure_2.jpeg)

![](_page_39_Figure_4.jpeg)

ate		Problem type	
	•	Classification	•

![](_page_39_Figure_7.jpeg)

![](_page_39_Picture_8.jpeg)

# Choosing a good learning rate is key <sub>J(w)</sub> for training convergence

![](_page_40_Figure_2.jpeg)

![](_page_40_Figure_4.jpeg)

![](_page_40_Figure_5.jpeg)

![](_page_40_Picture_6.jpeg)

#### Choosing a good learning rate is key J(w)[ for training convergence

![](_page_41_Figure_2.jpeg)

[https://playground.tensorflow.org, github: https://github.com/tensorflow/playground, Apache License 2.0]

![](_page_41_Figure_4.jpeg)

Adaptive Moment Estimation (Adam)

- Adaptive learning rates:
  - Reduce α during training
  - For each parameter separately
- Momentum:

 $^{\prime} \partial \mathrm{Loss} \setminus$  $\alpha$  $-\beta$  (previous change of  $\theta$ )

![](_page_41_Figure_12.jpeg)

![](_page_41_Picture_13.jpeg)

![](_page_41_Picture_14.jpeg)

#### **Activation Functions**

#### Requirements:

- Non-linear
- Fast differentiation
- Not bound to a fixed interval

![](_page_42_Figure_5.jpeg)

![](_page_42_Picture_7.jpeg)

![](_page_42_Picture_21.jpeg)

#### **Activation Functions**

#### **Requirements:**

- Non-linear
- Fast differentiation
- Not bound to a fixed interval

![](_page_43_Figure_5.jpeg)

![](_page_43_Picture_6.jpeg)

![](_page_43_Picture_7.jpeg)

- Initial weight must not be symmetric
- Random initialization

![](_page_44_Figure_3.jpeg)

"Xavier Weight Initialization"

for n input values to the node

![](_page_44_Figure_6.jpeg)

![](_page_44_Figure_7.jpeg)

"He Initialization"

# Uniform in $\begin{bmatrix} -\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}} \end{bmatrix}$

Gaussian  $\left(\mu = 0, \sigma = \sqrt{2/n}\right)$ 

14

![](_page_45_Figure_1.jpeg)

#### **Exploding and Vanishing Gradients**

![](_page_45_Picture_6.jpeg)

![](_page_46_Picture_0.jpeg)

# Pre-processing of the input data is essential for numerical stability • Want them to be centered around zero and with O(I) variance

subtract mean & divide by std. dev.

![](_page_46_Figure_3.jpeg)

![](_page_46_Figure_5.jpeg)

## Preprocessing

![](_page_46_Figure_7.jpeg)

![](_page_46_Figure_8.jpeg)

![](_page_46_Figure_9.jpeg)

![](_page_46_Picture_10.jpeg)

![](_page_46_Picture_11.jpeg)

- Often, the dataset is split into three parts:
  - Test data set: keep it until your model is fully optimized ! (you can look at it once...)

 Training data set: for the optimization of the NN parameters (you use it extensively)

- Validation data set: check performance during optimization  $\rightarrow$  more in Nicole's lecture!
- Example split: Train:Val:Test  $\rightarrow 60\%: 20\%: 20\%$

#### **Splitting the Data**

Test Data

Training Data

Data

![](_page_47_Picture_15.jpeg)

![](_page_47_Picture_16.jpeg)

![](_page_47_Picture_17.jpeg)

![](_page_47_Picture_18.jpeg)

![](_page_48_Picture_0.jpeg)

#### Network structure

- Number of hidden layers
- Number of nodes in the hidden layers
- Activation functions
- Weight Initialization
- Optimizer and its parameters
- (Initial) learning rate, ...

Batch size

• • •

Regularization  $\rightarrow$  Nicole's lecture

![](_page_48_Picture_11.jpeg)

![](_page_49_Picture_0.jpeg)

![](_page_49_Figure_1.jpeg)

\* also "sensitivity" \*\* also "I - 'specificity' "

#### Metrics

![](_page_49_Picture_5.jpeg)

TPR

![](_page_49_Picture_9.jpeg)

![](_page_50_Figure_0.jpeg)

![](_page_50_Figure_1.jpeg)

\* also "sensitivity" \*\* also "I - 'specificity' " • Metric: Error in the predicted profit, averaged over all products

#### Metrics

![](_page_50_Picture_8.jpeg)

TPR

![](_page_50_Picture_12.jpeg)

- LHC events with several jets, a charged lepton and a neutrino
  - Classification of jet permutations:
    - Pick the correct jet order out of many possible orderings

- DNN implementation\*
- Binary classification (ordering is either "right" or "wrong")

- **Output:** sigmoid
- Activation in hidden layers: ReLU
- \* More efficient proposals have been made since then.

#### Instead of a Summary: A Classification Example

[JINST 14 (2019) P11015, arXiv:1907:11181]

![](_page_51_Figure_15.jpeg)

![](_page_51_Picture_16.jpeg)

![](_page_51_Figure_17.jpeg)

![](_page_51_Picture_18.jpeg)

- 26 input features:
  - 4-momenta of the jets
  - Boolean: Is the jet flagged as a b-jet?
  - 3-momentum of the charged lepton
  - Magnitude and azimuth ( $\Phi$ ) of the missing transverse momentum
    - $\Phi$  is  $2\pi$ -continuous  $\rightarrow$  replace by sin( $\Phi$ ) and cos( $\Phi$ )

• Input scaling:

$$x_i \rightarrow \frac{x_i - \mu_i}{-}$$

 $\sigma_i$ 

![](_page_52_Picture_11.jpeg)

## **Classification Example: Training Progress**

![](_page_53_Figure_1.jpeg)

![](_page_53_Picture_2.jpeg)

## **Classification Example: Training Progress**

![](_page_54_Figure_1.jpeg)

![](_page_54_Picture_2.jpeg)

![](_page_54_Picture_3.jpeg)

## **Classification Example: Training Progress**

![](_page_55_Figure_1.jpeg)

![](_page_55_Picture_2.jpeg)

![](_page_55_Picture_3.jpeg)

0.77

0

 Metric vs. Epoch 0.85 Reconstruction efficiency 0.84 • Task: 0.83 Find the right jet ordering 0.82 0.81 • Metric: 0.80 Fraction of events with 0.79 correctly predicted 0.78 jet ordering

#### **Classification Example: Metrics**

![](_page_56_Figure_3.jpeg)

![](_page_56_Picture_4.jpeg)

#### **Classification Example: Hyperparameters**

• 2D grid search: number of hidden layers n and maximum number of nodes in hidden layers

 Choice of best epoch based on metric on validation data set

![](_page_57_Figure_3.jpeg)

![](_page_57_Picture_7.jpeg)

## • Now: Time for discussion

#### • Then:

Coffee

Faculty Club im TUM IAS-Gebäude in Garching, TUM / LMU München

Mastering Model Building

Faculty Club im TUM IAS-Gebäude in Garching, TUM / LMU München

14:55 - 15:15

Nicole Hartmann

15:15 - 16:30