

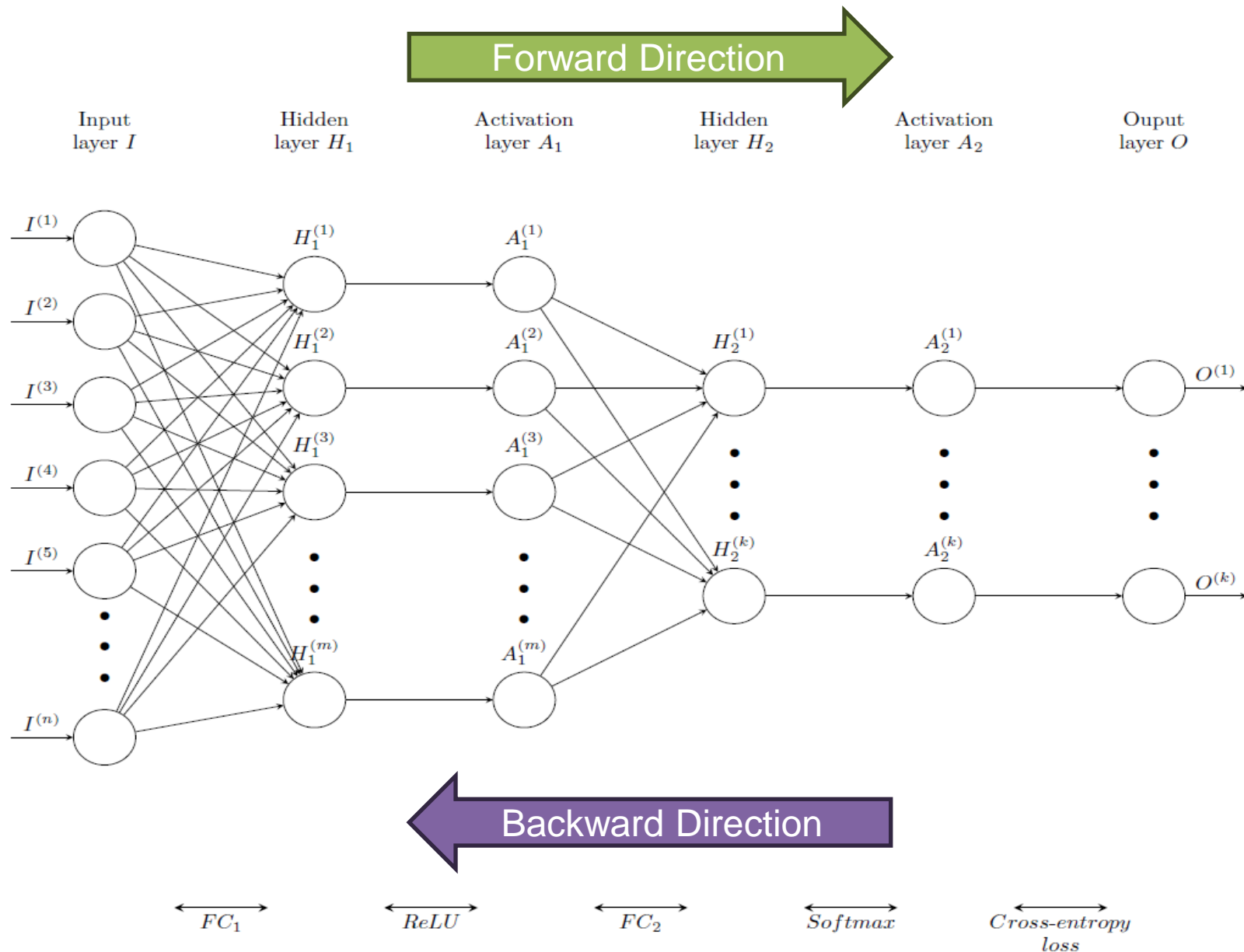
AdvPT Project



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

TECHNISCHE FAKULTÄT

Proposed FCNN Network



General Program Flow

Program workflow can be split into three parts:

- 1. Initialization**
- 2. Training**
- 3. Evaluation**

Each part will be explained in the following, but also take a look at the PyTorch reference for more insights!

Initialization

Tasks at hand are:

- Parse input configuration file
- Load and preprocess dataset
 - Normalize pixel input data via linear mapping (cf. 1a)
 - One-hot encoding of label data (cf. 1b)
- Set up model and layers with
 - Fixed input size $n = 28^2 = 784$ and output size $k = 10$
 - Hidden layer size m specified in input configuration file
- Initialize weights and biases of layers randomly
 - Fix the random seed for reproducibility
 - **Avoid zero-value initialization** (*Problem*: exploding gradients)
 - Use other initialization techniques for better network convergence
 - Uniform distribution with negative values, e.g. in range $[-1/n, 1/n]$ (*Problem*: vanishing gradients)
 - Better: Kaiming or Xavier initialization

Repeatedly process input data (repetitions are called epochs):

- Feed input data in forward direction to obtain prediction
- Calculate errors through a loss function
- Traverse network in backward direction and compute error tensors used for optimizing the trainable parameters

Pseudo-Code

- Repeat for a number of epochs
 - For each batch/sample in *training* dataset
 1. Perform forward pass
 2. Compute loss (and potentially log it to console)
 3. Perform backward pass
 4. Optimization: update trainable parameters

Trained model is tested on unseen data to measure generalizability to new inputs

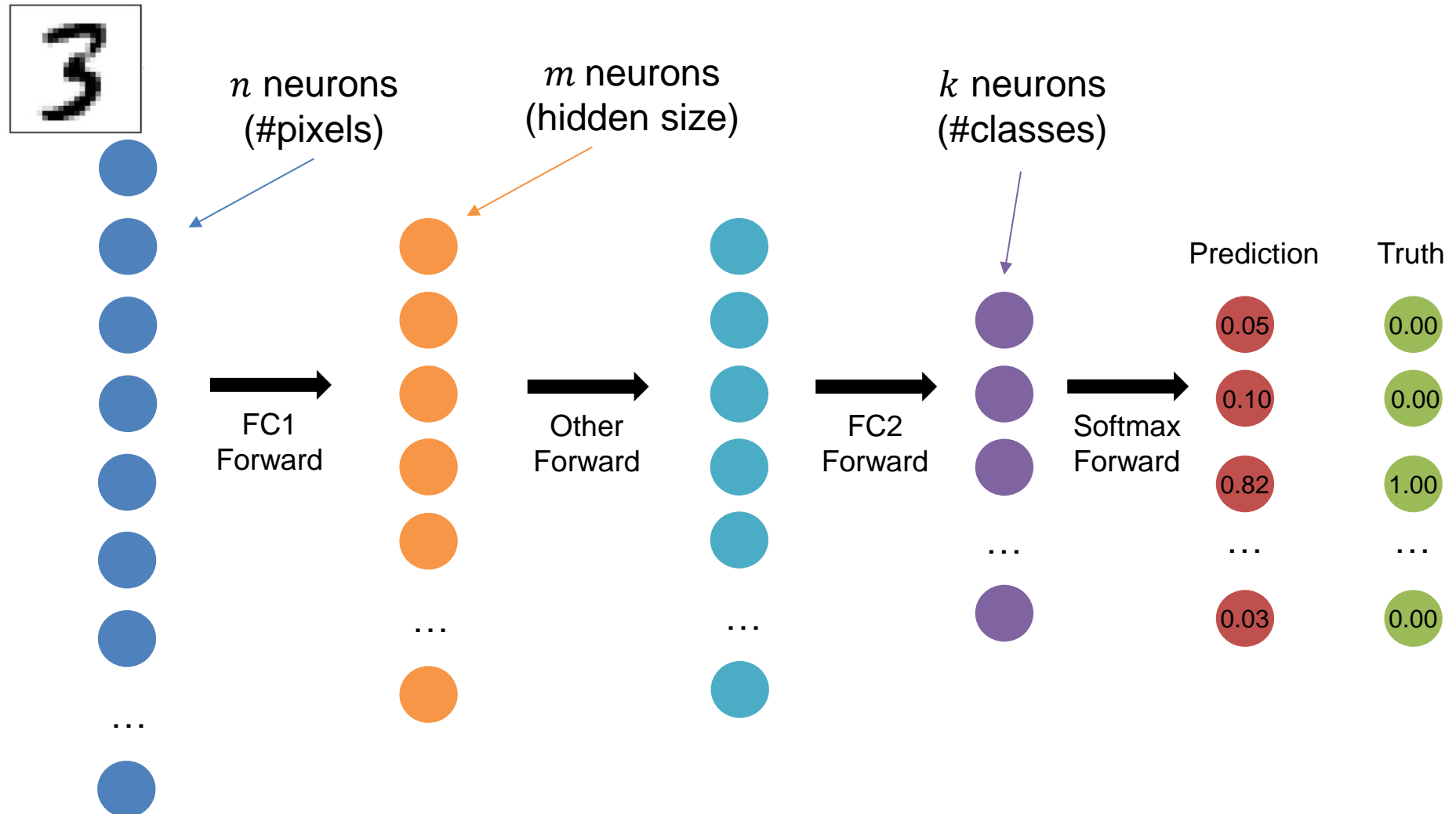
- Feed input data in forward direction to obtain prediction
- Compare predictions to ground truth (true label value) and print using the prediction logger for review
- Determine accuracy metric (i.e. percentage of correct predictions)

Pseudo-Code

- For each batch/sample in *testing* dataset
 1. Perform forward pass
 2. Compare predictions to true label values and log
 3. Compute accuracy

Forward Pass

Idea: Feed input data through network to obtain prediction that can be compared to ground truth



Quantifies difference between predicted and true class probability distributions

- Shows how well the network performs at prediction task
- Quantity to be minimized → Optimizer often make use of loss gradients

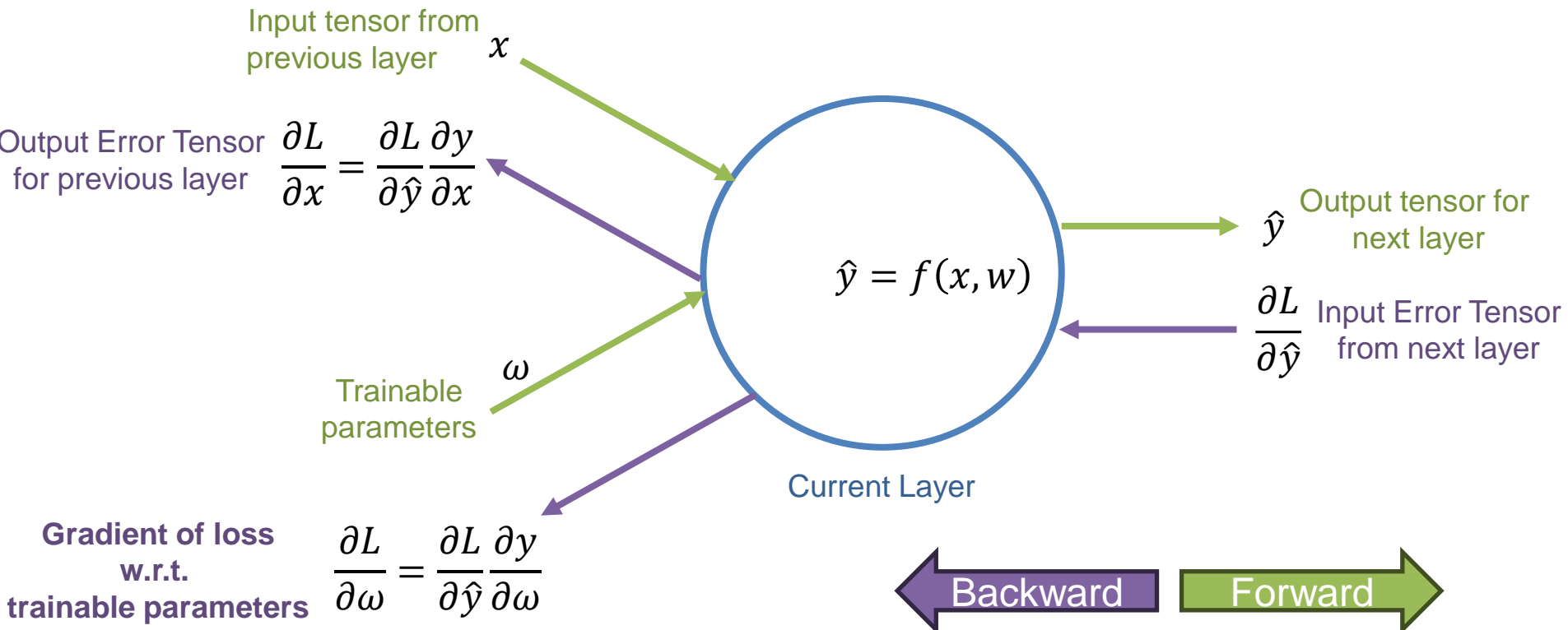
Proposed loss function: **cross-entropy loss**

$$L(y, \hat{y}) = \sum_{i=1}^k y_i \log \hat{y}_i$$

- k : number of labels (here: 10)
- y_i : true class probability (one-hot encoded)
- \hat{y}_i : predicted class probability output by previous softmax layer

Backward Pass

- Idea: Compute gradient of loss w.r.t. trainable parameters ω for optimization realized by recursive application of chain rule
- **Remember to store input tensors in your layer classes**



Adjust trainable parameters (weights and biases) to minimize loss
➔ Essential for improving performance of networks during training

Proposed optimizer: **stochastic gradient descent (SGD)**

$$\omega_{new} = \omega - \eta \frac{\partial L}{\partial \omega}$$

- η : Learning rate, provided in configuration file
- $\frac{\partial L}{\partial \omega}$: Gradient of loss w.r.t. ω . Computed in backward pass



Layer Description

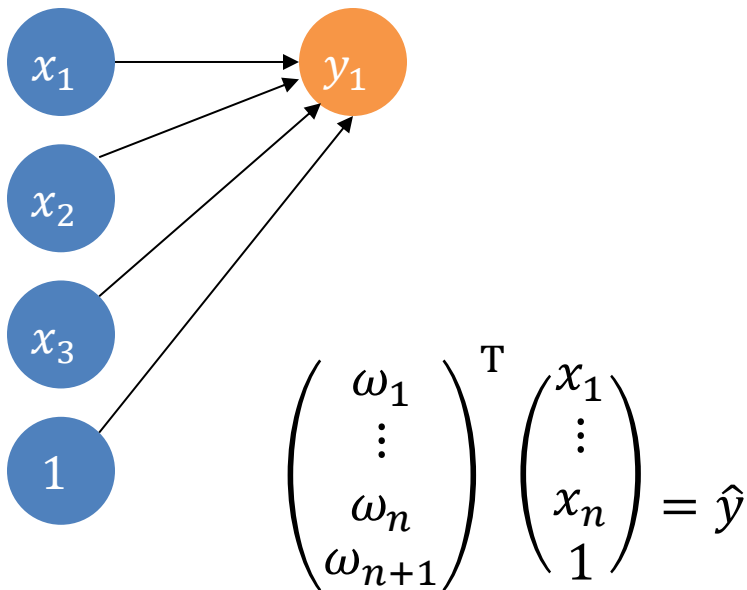


FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

TECHNISCHE FAKULTÄT

Forward pass with n input neurons and m output neurons (no batching)

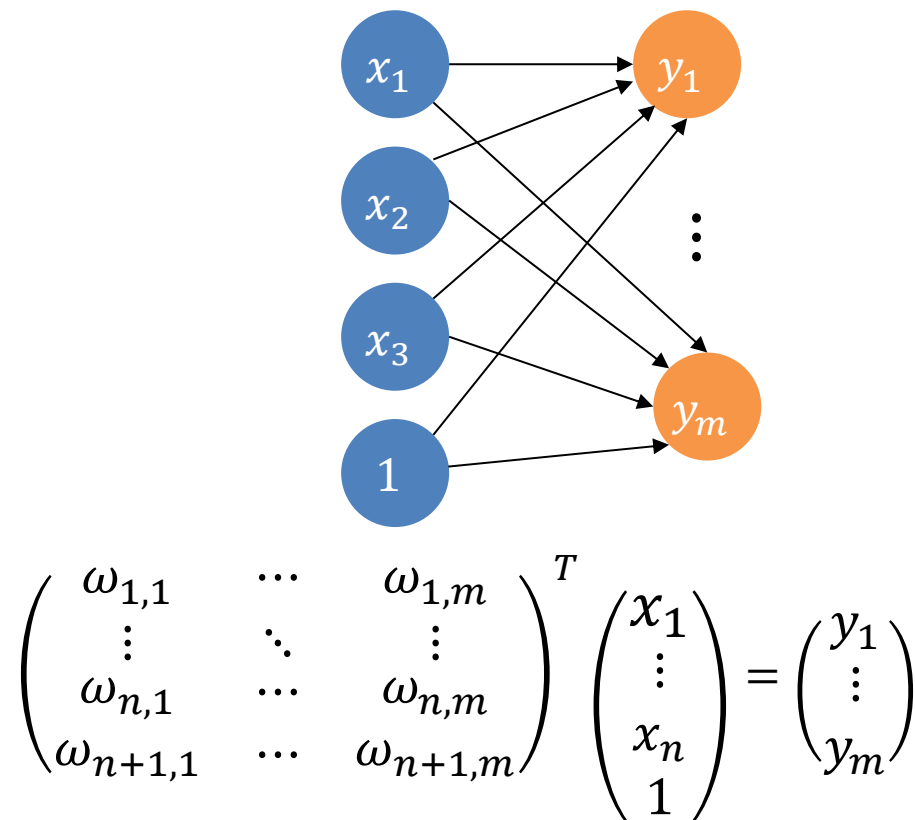
For single neuron



→ $\mathbf{w}^T \mathbf{x} = \hat{y}$

- \mathbf{w}^T : Vector with train params
- $\omega_1 \dots \omega_n$: weights
- ω_{n+1} : bias
- \mathbf{x} : input tensor

For all neurons



→ $\mathbf{W} \mathbf{x} = \hat{\mathbf{y}}$

Forward pass with batching

$$\begin{matrix} M \times (N+1) & (N+1) \times B & M \times B \\ \begin{pmatrix} \omega_{1,1} & \cdots & \omega_{1,m} \\ \vdots & \ddots & \vdots \\ \omega_{n,1} & \cdots & \omega_{n,m} \\ \omega_{n+1,1} & \cdots & \omega_{n+1,m} \end{pmatrix}^T & \begin{pmatrix} x_{1,1} & \cdots & x_{1,b} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,b} \\ 1 & \cdots & 1 \end{pmatrix} & = \begin{pmatrix} y_{1,1} & \cdots & y_{1,b} \\ \vdots & \ddots & \vdots \\ y_{m,1} & \cdots & y_{m,b} \end{pmatrix} \end{matrix}$$

$\longrightarrow \mathbf{WX} = \hat{\mathbf{Y}}$

But: inefficient memory layout causes strided data access since values are stored row-wise in C/C++

Fully Connected Layer

Backward pass:

- Return gradient with respect to (w.r.t.) X :

$$E_{n-1} = W^T E_n$$

- Update W using gradient w.r.t. W using SGD:

$$W_{new} = W - \eta E_n X^T$$

- E_n : Error tensor received from previous layer (backward direction)
- E_{n-1} : Error tensor passed to next layer (backward direction)
- η : Learning rate

Fully Connected Layer

Forward Pass *with improved memory layout* for batching caused adjustment to formulas

$$\begin{array}{ccc}
 \text{B} \times (\text{N}+1) & (\text{N}+1) \times \text{M} & \text{B} \times \text{M} \\
 \begin{pmatrix} x_{1,1} & \cdots & x_{1,b} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,b} \\ 1 & \cdots & 1 \end{pmatrix}^T & \begin{pmatrix} \omega_{1,1} & \cdots & \omega_{1,m} \\ \vdots & \ddots & \vdots \\ \omega_{n,1} & \cdots & \omega_{n,m} \\ \omega_{n+1,1} & \cdots & \omega_{n+1,m} \end{pmatrix} & = \begin{pmatrix} y_{1,1} & \cdots & y_{1,m} \\ \vdots & \ddots & \vdots \\ y_{b,1} & \cdots & y_{b,m} \end{pmatrix} \\
 & \longrightarrow & \mathbf{X}'\mathbf{W}' = \hat{\mathbf{Y}}'
 \end{array}$$

With:

- $\mathbf{X}' = \mathbf{X}^T$
- $\mathbf{W}' = \mathbf{W}^T$
- $\hat{\mathbf{Y}}' = \hat{\mathbf{Y}}^T$
- $\hat{\mathbf{Y}}^T = (\mathbf{W}\mathbf{X})^T = \mathbf{X}^T\mathbf{W}^T$

Fully Connected Layer

Backward pass *with improved memory layout*

- Return gradient with respect to (w.r.t.) X :

$$E'_{n-1} = E'_n W'^T$$

- Update W using gradient w.r.t. W using SGD:

$$W_{new} = W - \eta X'^T E'_n$$

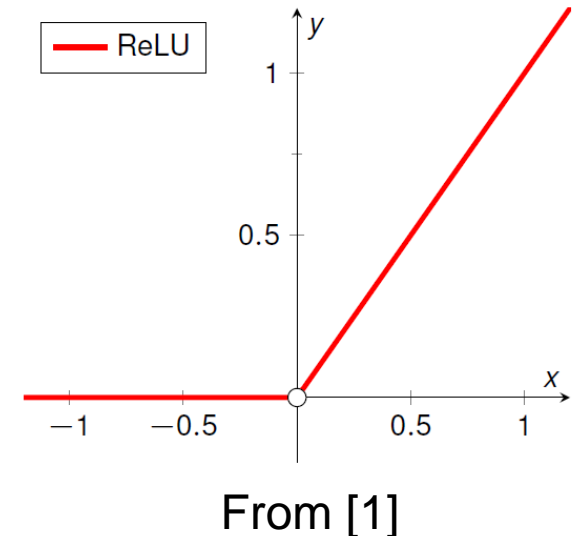
- E'_n : Error tensor received from previous layer (backward direction)
- E'_{n-1} : Error tensor passed to next layer (backward direction)
- η : Learning rate

ReLU Activation Function

Rectified Linear Unit (ReLU), introduces non-linearity to the model

Forward Pass

- Governing formula: $f(x) = \max(0, x)$



Backward Pass

- ReLU is not continuously differentiable
- Requires information about input tensor from forward pass
- Note: Activation function operate elementwise for each x (\odot operator)

$$E_{n-1} = E_n \odot \begin{cases} 0, & \text{if } x \leq 0 \\ 1, & \text{else} \end{cases}$$

Forward Pass encodes input neurons x to a probability function

$$\hat{y} = \frac{e^{x_i}}{\sum_{j=1}^k e^{x_j}}$$

Properties:

- k : Number of classes (here: 10)
- $\sum_{i=1}^k \hat{y}_i = 1$
- $\hat{y}_i \geq 0$
- If $x_i > 0$, exponent might become very large $\rightarrow x_i$ can be shifted to increase numerical stability by using:

$$\tilde{x}_i = x_i - \max(\mathbf{x})$$

Backward Pass

$$E_{n-1} = \hat{y} \left(E_n - \sum_{i=1}^k E_{n,i} \hat{y}_i \right)$$

- Note again: elementwise operations
- Remember to store \hat{y}_i

Forward Pass (cf. previous slide)

Backward Pass

$$E_n = \frac{-y}{\hat{y}}$$

- First layer in backward direction to compute error tensor
- \hat{y} : Predicted class probability
- y : True class probability

[1]: Andreas Meier, Activation Functions and Convolutional Neural Networks. Deep Learning Lecture Notes, 2020.

The top of the slide features a dark blue background with a faint, light blue image of the FAU main building and its statues. On the right side, there is a large, semi-circular seal containing the word 'ACADEMIA' and a profile of a person.

Thank you and good luck!