

# Vorlesung Machine Learning II

Jonas Pleyer  
FDM UNI Freiburg

06.07.2022

SR 404

Math. Institut



# Theoretical Foundations

## Structure of Neural Networks (Simple)

### Definition - Training set

We have pairs  $(x_i, y_i)$  of input/output values and assume that the output  $y_i$  was generated by some function  $f^*$  depending on  $x_i$ . The set of these values is the training set  $S$

$$S = \{(x_i, y_i = f^*(x_i))\}$$

With this definition, we neglect any measurement noise since  $f^*$  is assumed to be deterministic.

### Definition - Artificial Neural Network (ANN)

An ANN is comprised of an input layer with vector input  $x_i$  an output layer with output vector  $y_j$  and  $L$  hidden layers with activation functions  $f^k$ . The layers are connected ascendingly and the input  $t^{k+1}$  of the next layer is the weighted sum of output of the previous one

$$t^{k+1} = \omega^k t^k$$

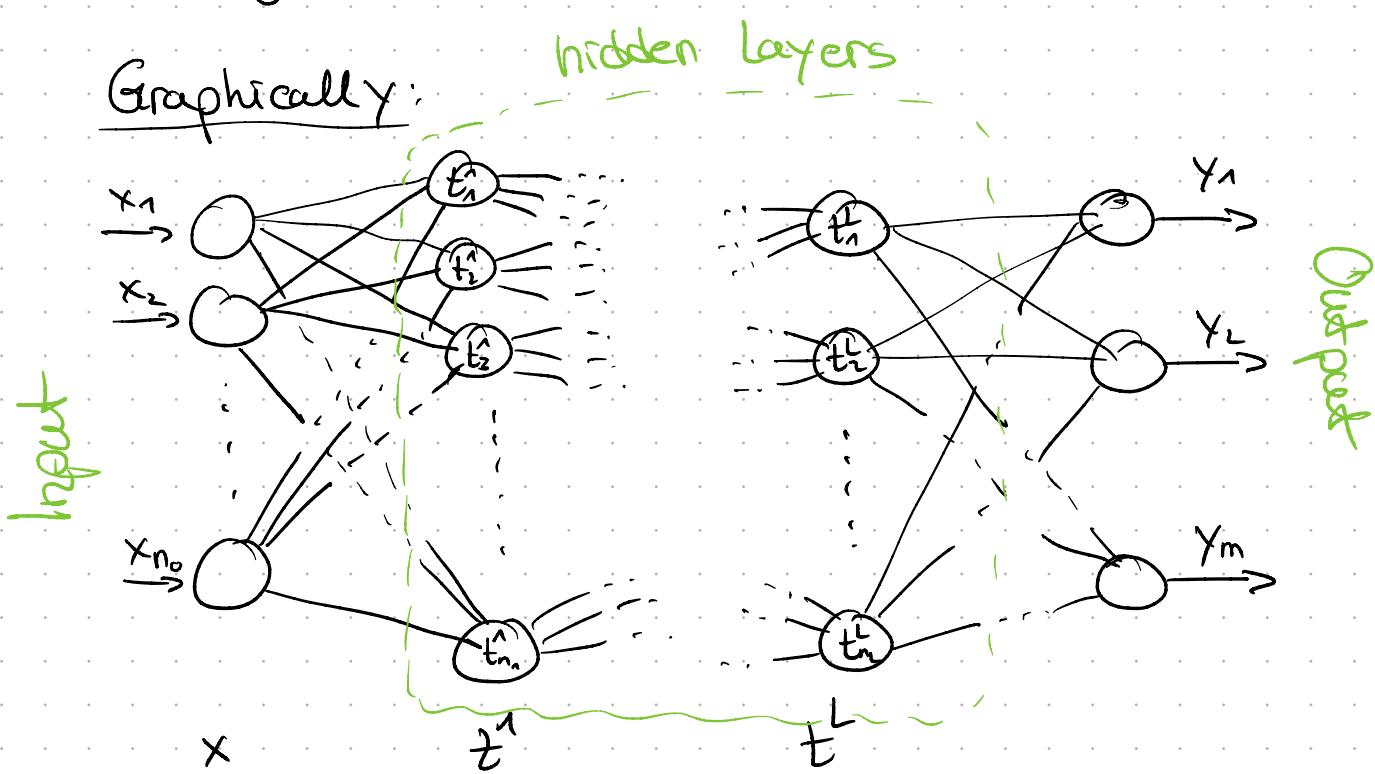
$$t_i^{k+1} = \sum_{j=1}^{n_k} w_{ij}^k t_j^k$$

where  $t^k = (t_1^k, \dots, t_{n_k}^k)$  is a vector and

$$\omega^k = \begin{pmatrix} \omega_{11}^k & \omega_{12}^k & \omega_{13}^k & \dots \\ \omega_{21}^k & \omega_{22}^k & \dots & \dots \\ \vdots & & & \end{pmatrix}$$

is the matrix containing the weights  $\omega_{ij}^k$  between layer  $k$  and  $k+1$  and node  $i$  of layer  $k+1$  and node  $j$  of layer  $k$ . The fall out  $g$  of the network  $g$  can be computed via

$$g(x) = f^L(\omega^L(f^{L-1}(\omega^{L-1} \dots f^1(\omega^1(x)) \dots)))$$



## Definition — Activation Function

The term "function" is used loosely here.

These can contain intermediate state and have non-trivial feedbacks. An activation function  $f$  takes a scalar value and produces a scalar value.

## Examples

1) Step-function

$$f(x) = \begin{cases} 1 & \text{if } x > \text{threshold} \\ 0 & \text{if } x \leq \text{threshold} \end{cases}$$

2) Linear function

$$f(x) = ax$$

3) Sigmoid

$$f(x) = \frac{1}{1+e^{-x}}$$

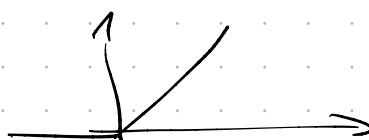


4) Tanh

$$\begin{aligned} f(x) &= 2 \operatorname{sigmoid}(x) - 1 \\ &= \frac{2}{1+e^{-x}} - \frac{1+e^{-x}}{1+e^{-x}} \\ &= \frac{1-e^{-x}}{1+e^{-x}} \end{aligned}$$

5) ReLU (Rectified Linear Unit)

$$f(x) = \max(0, x)$$



## Definition - Loss / Cost function

The loss-function  $C$  takes two output values (vectors) and compares them against each other. It returns a non-negative scalar value and is zero if and only if

$$C(y, z) = 0 \Leftrightarrow y = z$$

### Examples:

let  $y, z \in \mathbb{R}^m$  be vectors.

- 1) Choose  $C$  such that

$$C(y, z) = \frac{1}{m} \sum_{i=1}^m (y_i - z_i)^2$$

Comment:

General purpose cost-function is most often chosen.

the mean-squared error (MSE).

- 2) Assume the variables  $y$  and  $z$  are parameters for polynomials on a interval  $I$ . Then we define

$$C(x, y) = \max_{x \in I} \left| \sum_{i=1}^m y_i x^i - z_i x^i \right|$$

the supremum norm.

# Backpropagation

## Questions:

- i) How do Neural Networks learn?
- ii) How do we adjust parameters?
- iii) Are there different techniques?

## Definition - Gradient descent

Consider a differentiable function  $F$ .

The gradient of  $F$  at point  $x_0$  points into the direction of steepest ascend while its negative value points to the largest direction of descent. Starting from the initial  $x_0$ , we can obtain the next point  $x_1$  by

$$x_1 = x_0 - \gamma \nabla F(x_0)$$

and thus  $F(x_1) \leq F(x_0)$ . Iteratively, we obtain

$$x_{n+1} = x_n - \gamma \nabla F(x_n)$$

and hope for convergence as  $n \rightarrow \infty$ .

The parameter  $\gamma$  controls the speed of learning and is thus called the learning rate.

To use the gradient-descent method in ANNs, we need to compute the gradient of the cost-function with respect to the weights in the ANN.

### Calculation

$$\frac{\partial C}{\partial w_{ij}^L} = ???$$

We use chain-rule applied to

$$C(y, g(x)) = C(y, f^L(w^L(f^{L-1}(w^{L-1} \dots f^1(w^1(x)) \dots))) )$$

We denoted

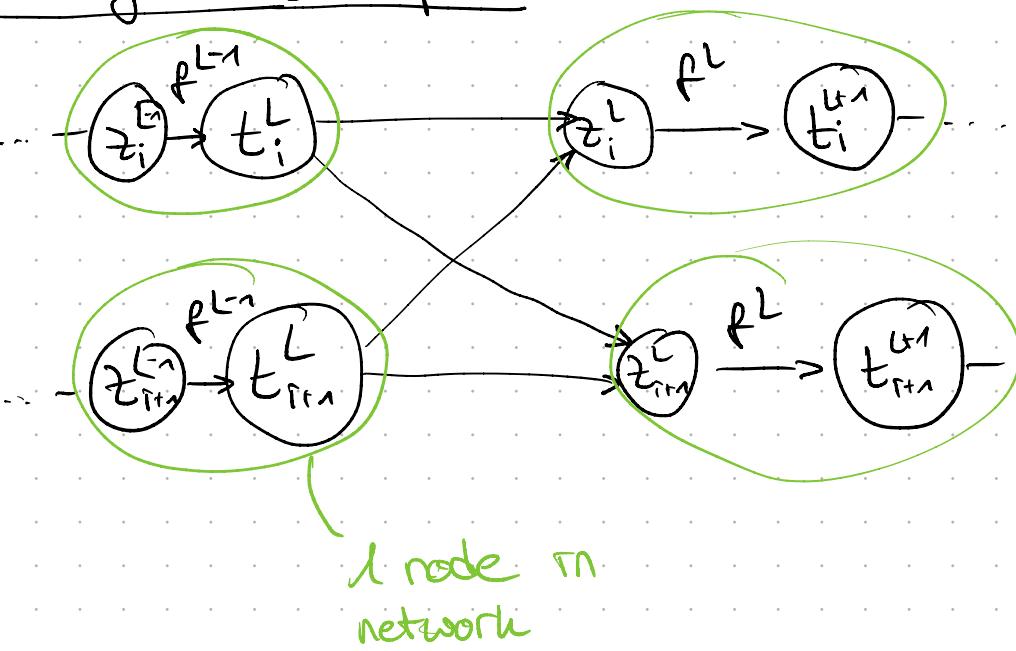
$$\begin{aligned} t^{L+1} &= f^L(w^L t^L) \\ &= f^L(w^L(f^{L-1}(w^{L-1} t^{L-1}))) \\ &\vdots \\ &= f^L(w^L(\dots f^1(w^1 x) \dots)) \end{aligned}$$

as the output of the layer L.

Define the weighted input of each layer

$$z^L = w^L t^L$$

## Diagram (Example)



$$z^{L+1} = f^L(z_i^L) = f^L(\omega^L t^L)$$

$$z^L = \omega^L t^L$$

Now we can rewrite the cost-function

$$\begin{aligned} & \nabla_x C(y, g(x)) \\ &= \underbrace{\nabla_x g(x)}_{\substack{\text{vector} \\ \text{matrix}}} \cdot \underbrace{\nabla_z C(y, g(x))}_{\substack{\text{vector} \\ \text{vector}}} \end{aligned}$$

scalar

$$= (\nabla_x z^L \cdot \nabla_{z^L} f^L) \quad \nabla_{t^L} C$$

$$= ((\nabla_x t^L \underbrace{\nabla_{t^L} z^L}_{\omega^L}) \nabla_{z^L} f^L) \quad \nabla_{t^L} C$$

$$= (((\nabla_x z^{L-1} \nabla_{z^{L-1}} f^{L-1}) \omega^L) \nabla_{z^L} f^L) \quad \nabla_{t^L} C$$

= ...

$$= \omega^1 \nabla_{z^1} f^1 \omega^2 \nabla_{z^2} f^2 \dots \omega^L \nabla_{z^L} f^L \nabla_t C$$

Define

$$g^L := \nabla_{z^L} f^L \omega^{L+1} \dots \omega^L \nabla_{z^L} f^L \nabla_t C$$

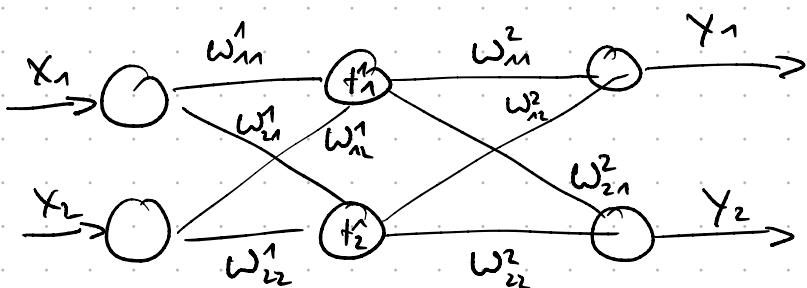
then the gradient at layer L becomes

$$\nabla_{\omega^L} C = t^{L-1} g^L$$

$$g^{L-1} = \nabla_{z^L} f^{L-1} \omega^L g^L$$

backwards-propagation

Example



suppose that  $f(x) = \alpha x$

$$C(z, y) = \frac{(z_1 - y_1)^2 + (z_2 - y_2)^2}{2}$$

$$g(x) = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

$$C(z, g(x)) = \frac{(z_1 - g_1(x))^2 + (z_2 - g_2(x))^2}{2}$$

$$\nabla_y C = \begin{pmatrix} \frac{\partial C}{\partial y_1} \\ \frac{\partial C}{\partial y_2} \end{pmatrix} = \begin{pmatrix} -(z_1 - y_1) \\ -(z_2 - y_2) \end{pmatrix}$$

$$g^2 = \underbrace{\nabla_{z^2} f}_{a} \nabla_y C$$

$$= a \nabla_y C = a \begin{pmatrix} y_1 - z_1 \\ y_2 - z_2 \end{pmatrix}$$

$$\Rightarrow g^1 = \underbrace{\nabla_{z^2} f}_{a^2} \omega^2 g^2$$

$$= a^2 \omega^2 \begin{pmatrix} y_1 - z_1 \\ y_2 - z_2 \end{pmatrix}$$

$$= a^2 \begin{pmatrix} \omega_{11}^2 & \omega_{12}^2 \\ \omega_{21}^2 & \omega_{22}^2 \end{pmatrix} \begin{pmatrix} y_1 - z_1 \\ y_2 - z_2 \end{pmatrix}$$

$$= a^2 \left( \omega_{11}^2 (y_1 - z_1) + \omega_{12}^2 (y_2 - z_2) \right. \\ \left. \omega_{21}^2 (y_1 - z_1) + \omega_{22}^2 (y_2 - z_2) \right)$$

$$\nabla_{\omega^L} = t^L g^L$$

$$\nabla_{\omega^1} = t^1 g^1$$

$$= f(\omega^1 x) g^1$$

$$= a \begin{pmatrix} \hat{\omega}_{11} & \hat{\omega}_{12} \\ \hat{\omega}_{21} & \hat{\omega}_{22} \end{pmatrix} (\dots)$$

$$= a^3 \left( \hat{\omega}_{11}^1 (\omega_{11}^2 (y_1 - z_1) + \omega_{12}^2 (y_2 - z_2)) + \hat{\omega}_{12}^1 (\dots) \right. \\ \left. \hat{\omega}_{21}^1 (\omega_{21}^2 (y_1 - z_1) + \omega_{22}^2 (y_2 - z_2)) + \hat{\omega}_{22}^1 (\dots) \right)$$

$$\nabla_{\omega^2} = t^2 S^2$$

$$= f(\omega^2 \underbrace{f(\omega^1 x)}_{\text{already calculated}}) S^2$$

$$S^2 = a \begin{pmatrix} y_1 - z_1 \\ y_2 - z_2 \end{pmatrix}$$