

Machine Learning and Data Mining WS21/22

# “9 Neural Networks”

Dr. Zeyd Boukhers

@ZBoukhers

Institute for Web Science and Technologies  
University of Koblenz-Landau

January 5, 2022

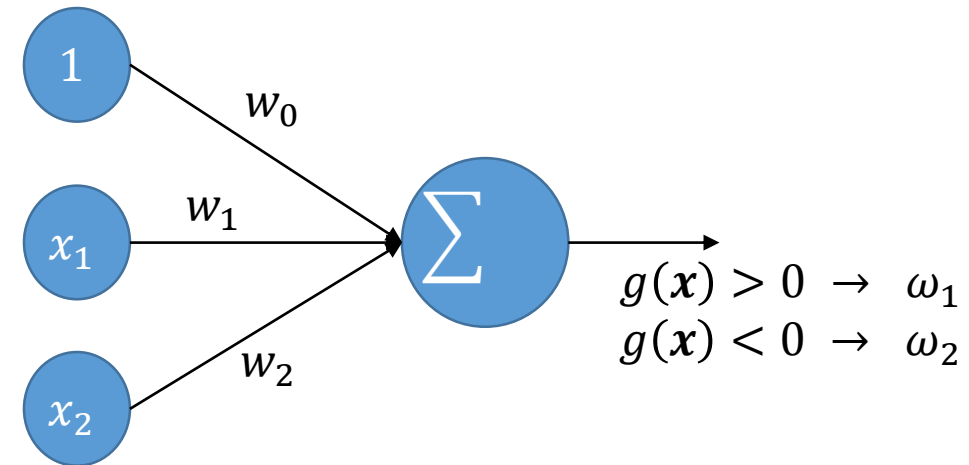
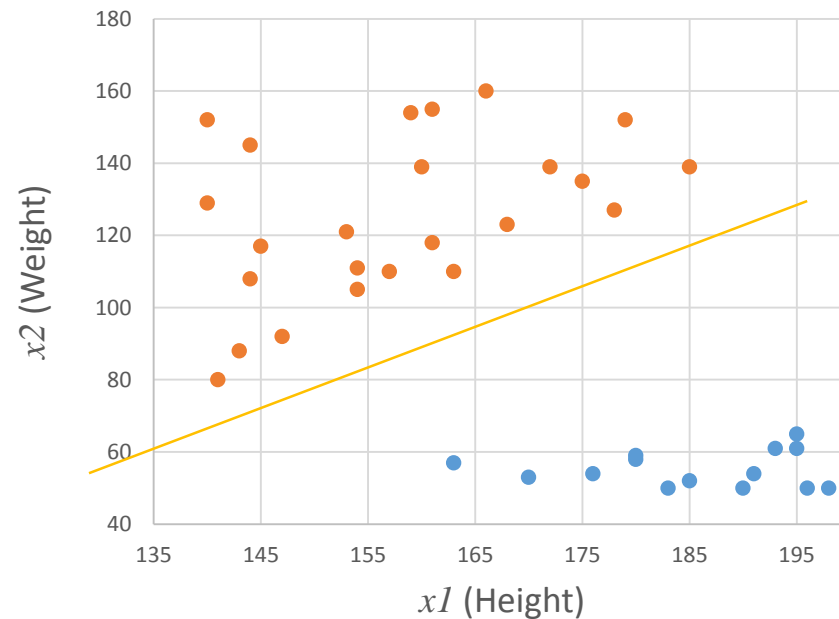


- Linear regression
  - Least squares function
  - Optimization
- Linear classification
  - Perceptron classifier
  - Support Vector Machine (SVM)
  - Optimization

- From Perceptron classifier to Multi-Layer Perceptron
  - XOR problem
  - Activation function
  - Two-Layer Perceptron
  - Gradient Decent
  - Backpropagation

# Perceptron classifier

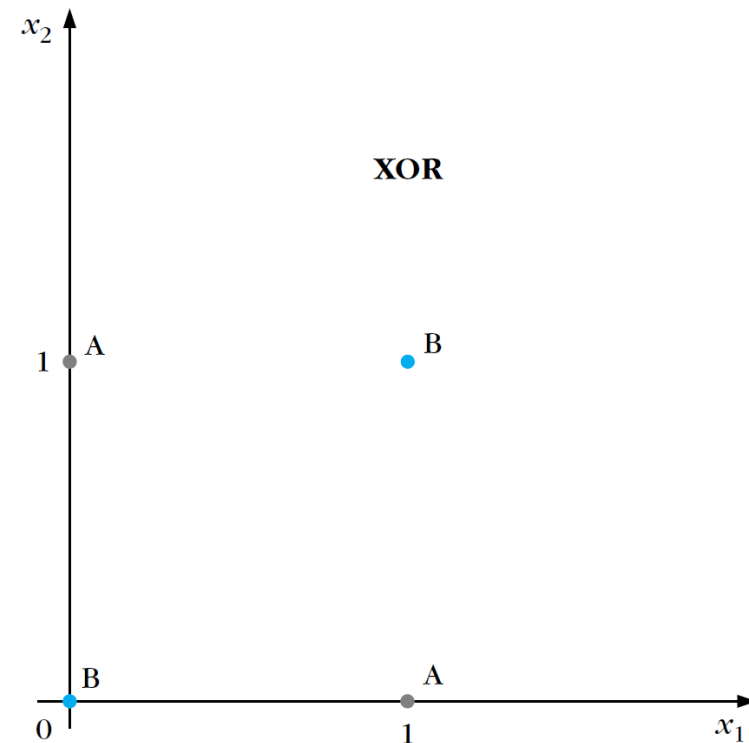
- Perceptron: the aim is to learn a hyperplane  $g(\mathbf{x}) = w_0 + \mathbf{w}^T \mathbf{x}$  that separates the classes.
- $g(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2$



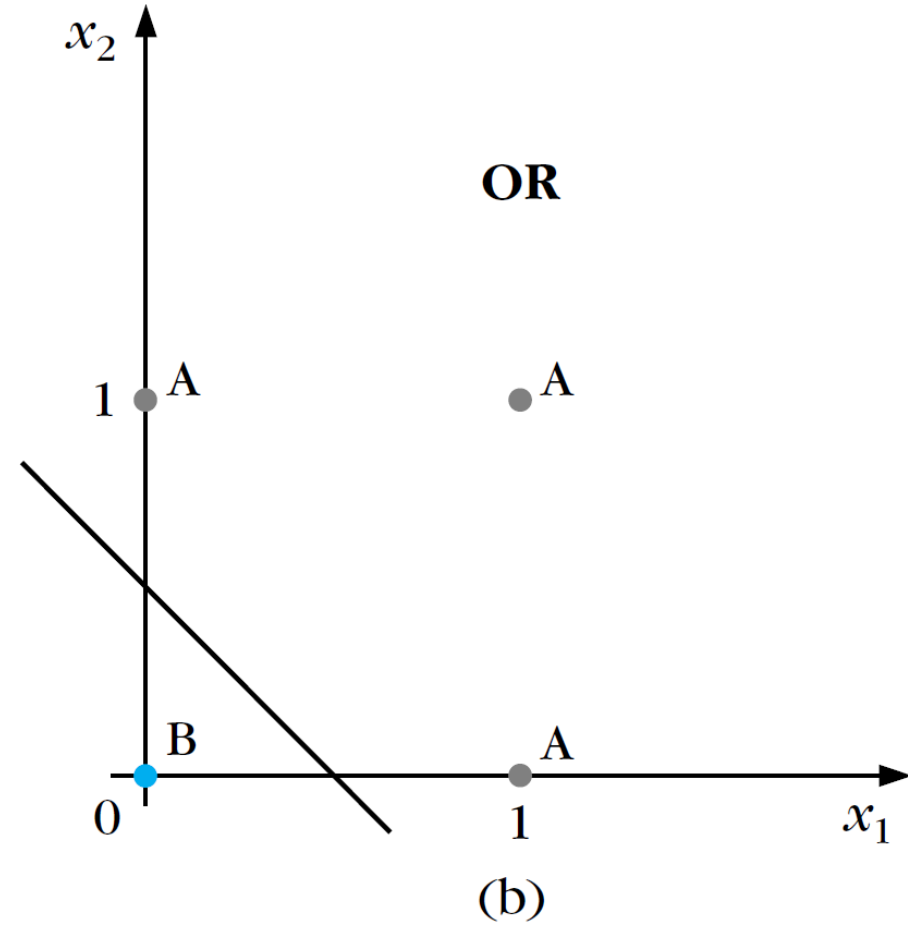
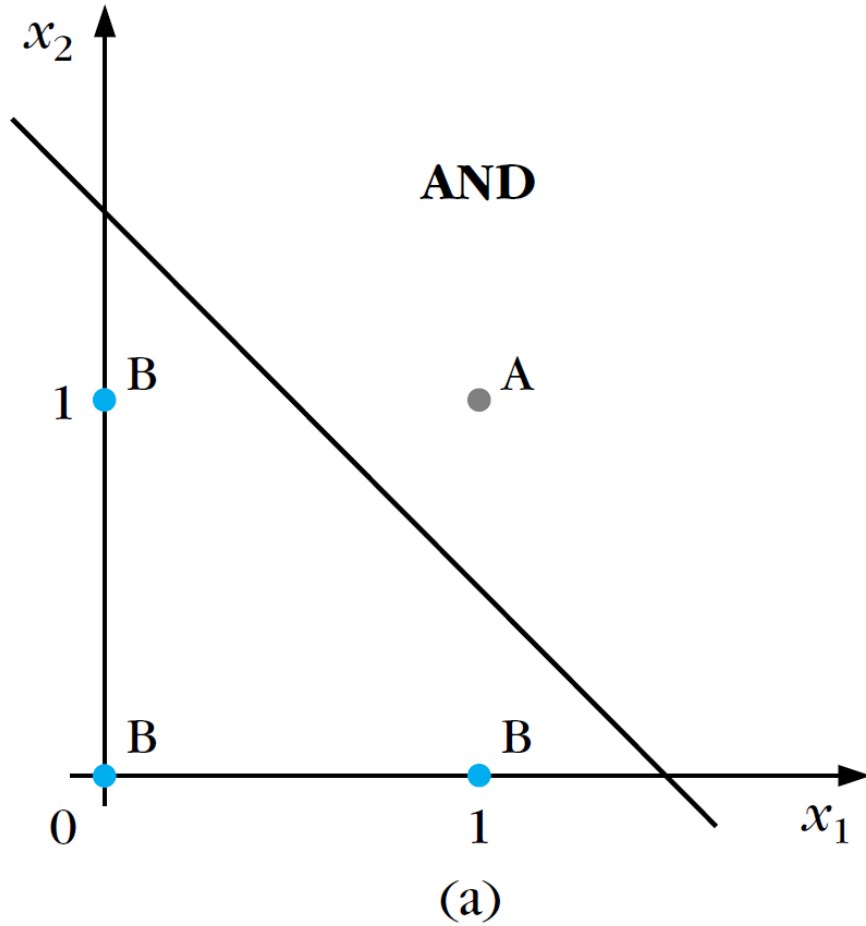
# XOR problem

- The Exclusive OR (XOR) Boolean function is a typical example of such a problem.
- Given the values of the input binary data  $\mathbf{x} = [x_1, x_2]^T$ , the output of the function is either 0 or 1, corresponding to the classes  $\omega_1$  and  $\omega_2$ , respectively.

$x_1$	$x_2$	XOR	Class
0	0	0	$\omega_1(B)$
0	1	1	$\omega_2(A)$
1	0	1	$\omega_2(A)$
1	1	0	$\omega_1(B)$



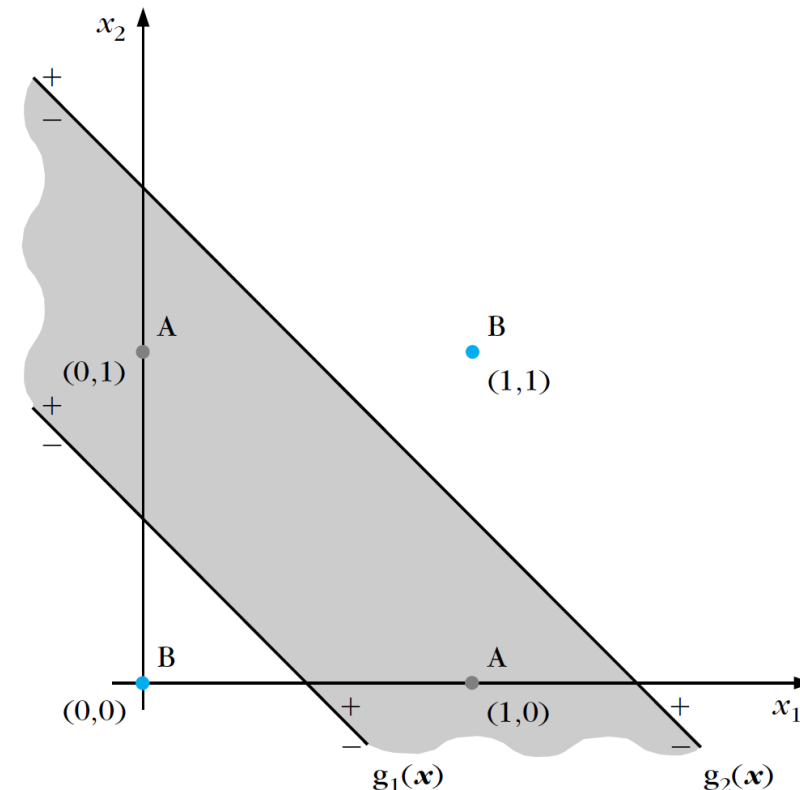
# XOR problem



# XOR problem

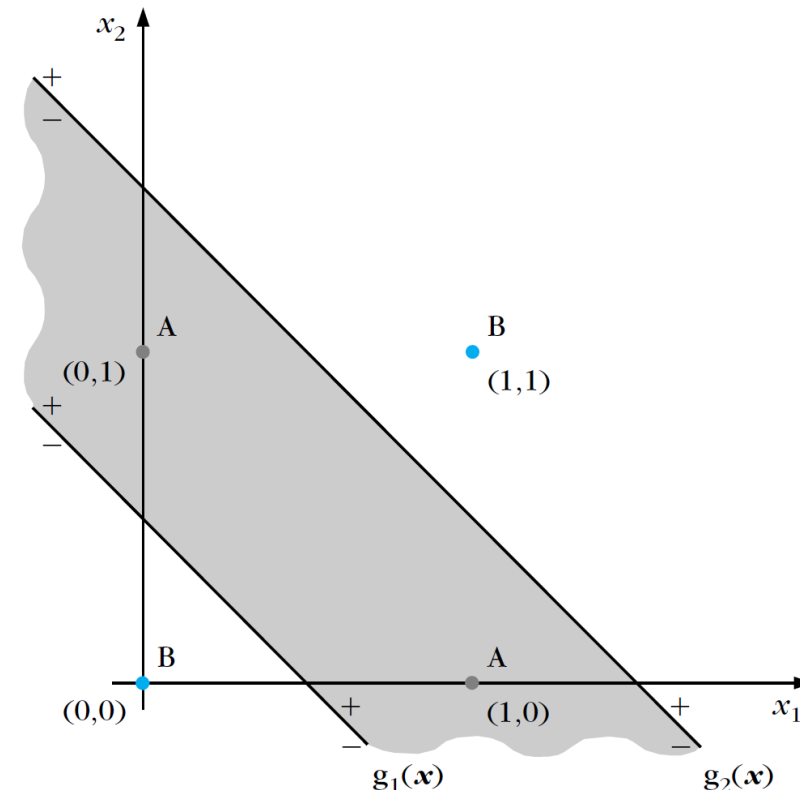
- Now the classes can be separated with the two lines:  $g_1(\mathbf{x}) = g_2(\mathbf{x}) = 0$ 
  - $A$  is to the right (+) of  $g_1(\mathbf{x})$  and to the left (−) of  $g_2(\mathbf{x})$
  - $B$  is to the right of both lines.
  - or to the left of both lines.

➤ The problem has been attacked  
In *two successive* phases.



# XOR problem

- In the first phase, the position of  $\mathbf{x}$  is computed w.r.t each line.
- In the second phase, the results of the previous phase are combined to find the position of  $\mathbf{x}$  w.r.t both lines.
  - outside or inside the shaded area.

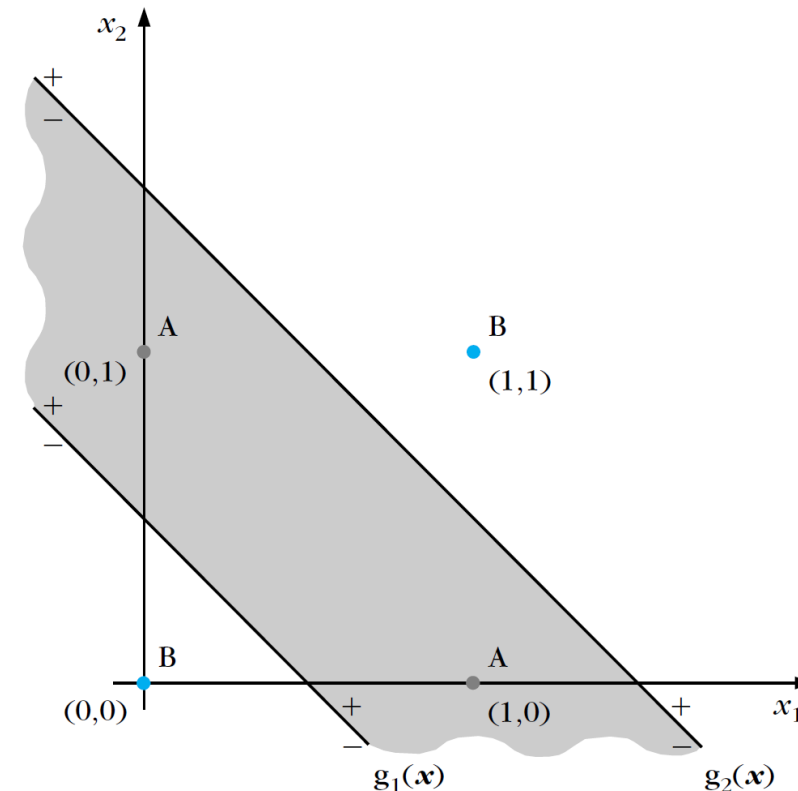




# XOR problem

- In the first phase, two perceptrons are adopted with input  $x_1$  and  $x_2$  and appropriate synaptic weights.
- $y_u = f(g_u(\mathbf{x})), u = 1, 2$ 
  - $f(\cdot)$  is the activation function.
  - $f(\mathbf{x}) = \begin{cases} 1 & \text{if } g(\mathbf{x}) \geq 0 \\ 0 & \text{if } g(\mathbf{x}) < 0 \end{cases}$

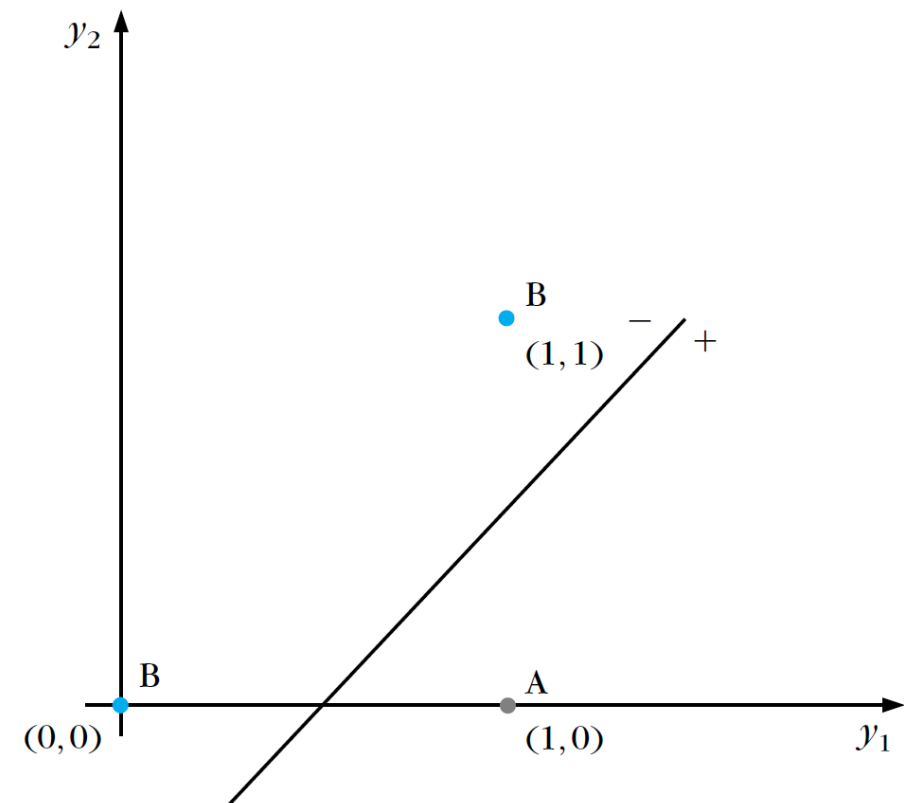
Phase 1				Phase 2	
$x_1$	$x_2$	$y_1$	$y_2$	XOR	Class
0	0	0	0	0	$\omega_1(B)$
0	1	1	0	1	$\omega_2(A)$
1	0	1	0	1	$\omega_2(A)$
1	1	1	1	0	$\omega_1(B)$



# XOR problem

- In the first phase, the input vector  $x$  is mapped to a new one  $y = [y_1, y_2]$ .
- In the second phase, the decision is made based on the transformed data.
  - $y = [0, 0] \rightarrow \omega_1$
  - $y = [1, 1] \rightarrow \omega_1$
  - $y = [1, 0] \rightarrow \omega_2$

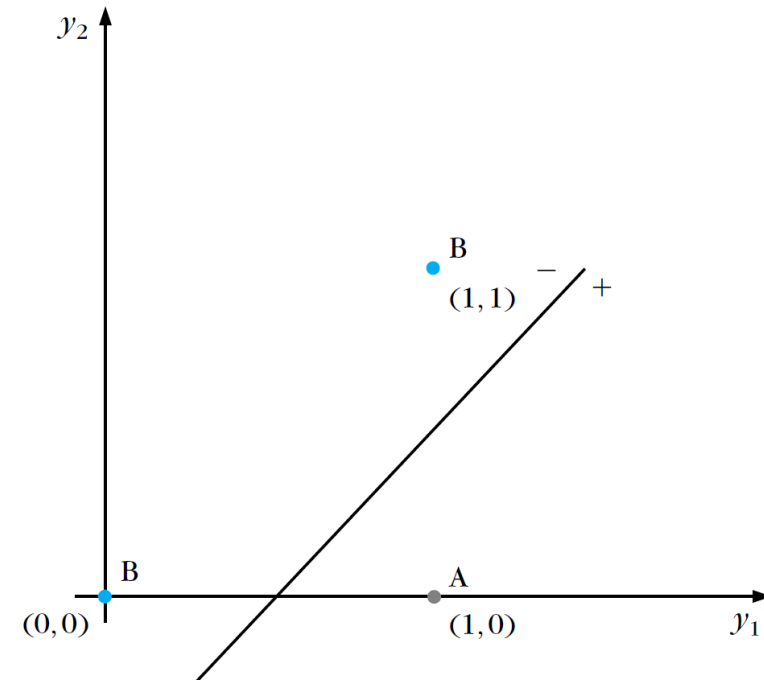
Phase 1				Phase 2	
$x_1$	$x_2$	$y_1$	$y_2$	XOR	Class
0	0	0	0	0	$\omega_1(B)$
0	1	1	0	1	$\omega_2(A)$
1	0	1	0	1	$\omega_2(A)$
1	1	1	1	0	$\omega_1(B)$



# XOR problem

- What happens?
    - The first phase transforms (by mapping  $x$  to  $y$ ) the nonlinearly separable problem to a linearly separable one.
  - This multilayer architecture is a generalization of the perceptron.
- *two-layer perceptron.*

Phase 1				Phase 2	
$x_1$	$x_2$	$y_1$	$y_2$	XOR	Class
0	0	0	0	0	$\omega_1(B)$
0	1	1	0	1	$\omega_0(A)$
1	0	1	0	1	$\omega_0(A)$
1	1	1	1	0	$\omega_1(B)$



- Why is it important?
  - Without it, the Multi-Layer Perceptron (MLP) is just a composition of successive linear functions.
    - A linear function.
- ✓ A non linear function (activation function) that introduces non-linearity into the model.

*"If the activation functions of all the hidden units in a network are taken to be linear, then for any such network we can always find an equivalent network without hidden units." [9] (Bishop 2006)*

# Two-layer perceptron

- It is also called “*two-layer feedforward neural network*”.

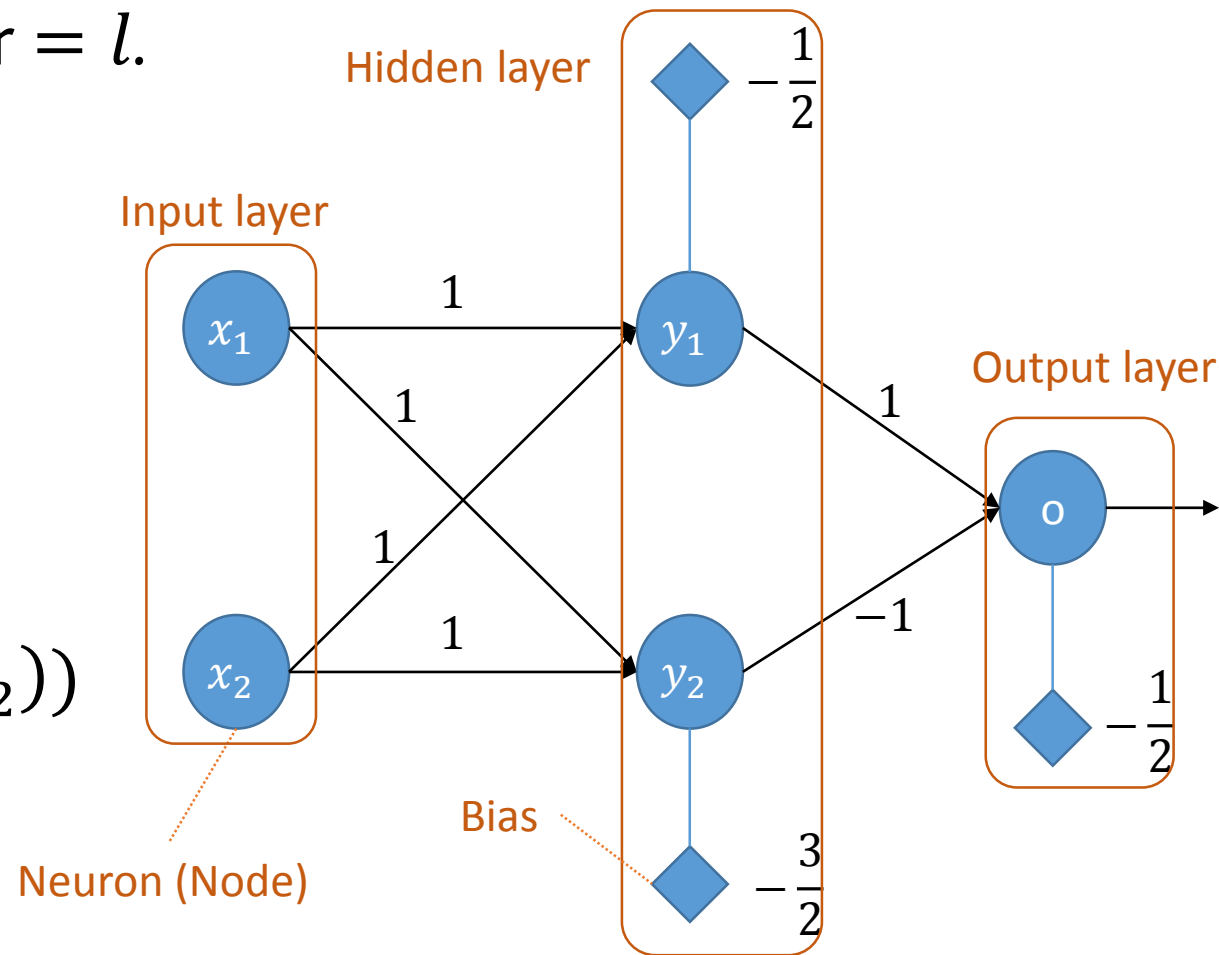
- The number of nodes in the input layer =  $l$ .

- $g_1(\mathbf{x}) = x_1 + x_2 - \frac{1}{2} = 0$

- $g_2(\mathbf{x}) = x_1 + x_2 - \frac{3}{2} = 0$

- $g(\mathbf{y}) = y_1 - y_2 - \frac{1}{2} = 0$

- $g(\mathbf{y}, \tilde{\mathbf{w}}) = f(g_1(\mathbf{x}, \tilde{\mathbf{w}}_1)) + f(g_2(\mathbf{x}, \tilde{\mathbf{w}}_2))$



# Two-layer perceptron

- $g(\mathbf{y}; \tilde{\mathbf{w}}) = f(g_1(\mathbf{x}; \tilde{\mathbf{w}}_1)) + f(g_2(\mathbf{x}; \tilde{\mathbf{w}}_2))$ 
  - Where  $\tilde{\mathbf{w}}_j = [\mathbf{w}_j^T, b_j]$ ,  $\forall j$
  - Let  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_l]$
- $g_1(\mathbf{X}; \tilde{\mathbf{w}}_1) = \mathbf{w}_1^T \mathbf{X} + b = [1, 1]^T \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix} - \frac{1}{2} = [0 \quad 1 \quad 1 \quad 2] - \frac{1}{2} =$   
 $\begin{bmatrix} -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{3}{2} \end{bmatrix}$
- $f(g_1(\mathbf{x}, \tilde{\mathbf{w}}_1)) = [0 \quad 1 \quad 1 \quad 1]$
- Let  $\tilde{\mathbf{W}} = [\tilde{\mathbf{w}}_1, \tilde{\mathbf{w}}_2]$

Compute  $f(g_2(\mathbf{x}, \tilde{\mathbf{w}}_2))$   
Compute  $f(g(\mathbf{y}, \tilde{\mathbf{w}}))$

$$f(g(\mathbf{y})) = f(\tilde{\mathbf{w}}^T f(\tilde{\mathbf{W}}^T \mathbf{x}))$$

- The architecture of *two-layer perceptron* can be generalized
  - to  $l$ -dimensional input vectors and
  - to more than two neurons in the hidden layer and
  - to more than one neuron in the output layer
  - to more than two layers

- In other words, how to find the weights?
  - Linear regression → Linear programming (linear optimization)
  - SVM → Quadratic programming
  - Perceptron → Iterative optimization (Non-linear function)
    - Gradient descent.
- MLP ?
  - Gradient descent
  - Stochastic gradient descent



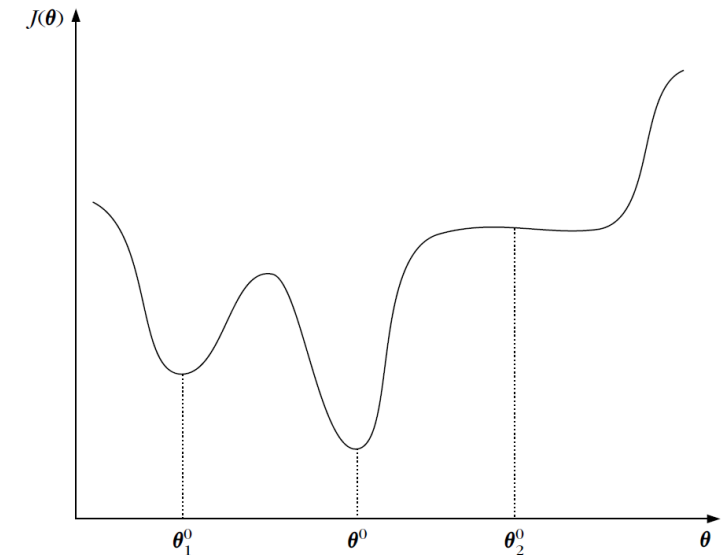
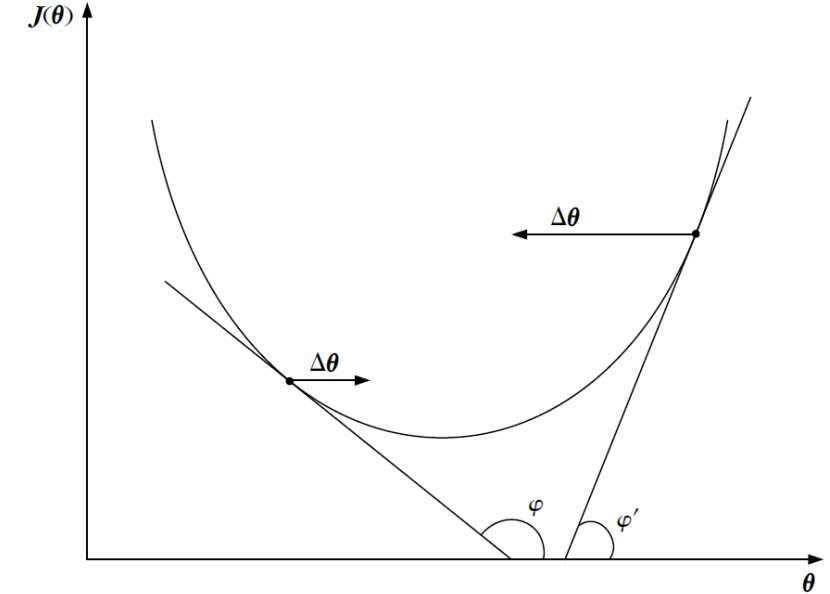
- Given  $N$  training samples, let:
  - $\theta$  be the set of parameters (to be found)
  - $p(\hat{y}_i | \mathbf{x}_i; \theta)$  be the likelihood of the true label of  $\mathbf{x}_i$

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} \prod_{i=1}^N p(\hat{y}_i | \mathbf{x}_i; \theta) \\ &= \sum_{i=1}^N \log(p(\hat{y}_i | \mathbf{x}_i; \theta))\end{aligned}$$

- Maximizing the log likelihood of a categorical distribution (classification) corresponds to minimising the cross-entropy between the approximated distribution and the true distribution.
  - $H(\hat{\mathbf{y}}, \mathbf{y}) = -\sum_i y_i \log(\hat{y}_i)$
  - $H(\hat{\mathbf{y}}, \mathbf{y}) \neq H(\mathbf{y}, \hat{\mathbf{y}})$
- Maximising the log of a continuous distribution (regression) corresponds to minimising the mean squared error between the approximated mean and true mean.

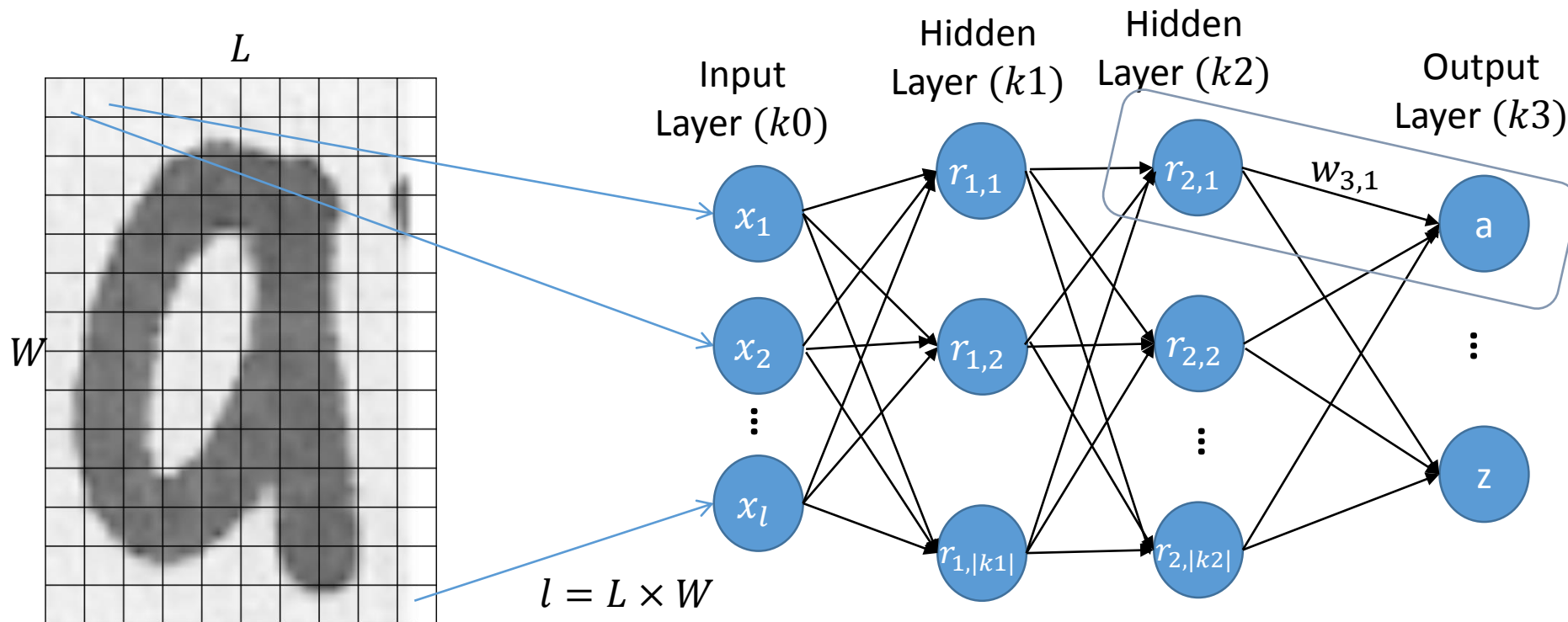
# Gradient descent

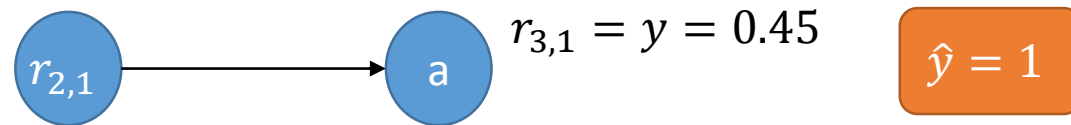
- Goal: finding  $\theta$  that minimizes  $J(\theta)$
- How: Iterative process
  - $\theta_{new} = \theta_{old} + \Delta\theta$
  - $\Delta\theta = -\rho \left. \frac{\partial J(\theta)}{\partial \theta} \right|_{\theta=\theta_{old}}$
  - $\theta_{new}$  is chosen in the direction that decreases  $J(\theta)$
  - $\rho$  is the learning rate
    - High  $\rho \rightarrow$  big step
    - Low  $\rho \rightarrow$  small step
  - Convergence: when the gradient = 0



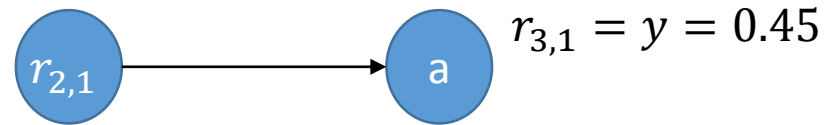
# How to build MLP?

- Given the architecture below, where all neurones employ the same activation function (sigmoid)
- We assume that  $N$  training samples  $(\mathbf{x}_i, \hat{y}_i), i = 1, \dots, N$  are available.



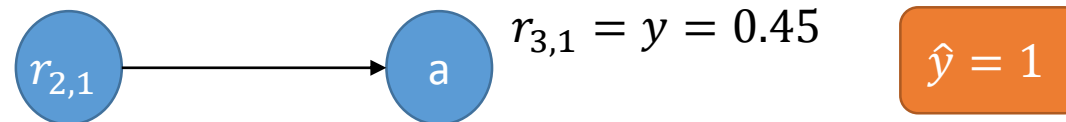


- Since the output is no longer scalar, it will be  $\mathbf{y} = [y_1, y_2, \dots, y_{|k3|}]^T$ 
  - The label will be  $\hat{\mathbf{y}} = [\hat{y}_1 = 1, \hat{y}_2 = 0, \dots, \hat{y}_{|k3|} = 0]^T$
- The goal is to minimize the cost function  $J = \sum_{i=1}^N \varepsilon(\mathbf{y}_i, \hat{\mathbf{y}}_i)$ 
  - For all training samples.
  - Let  $\varepsilon(.,.)$  be the MSE:  $\varepsilon(\mathbf{y}_i, \hat{\mathbf{y}}_i) = \sum_{j=1}^{k3} (y_{i,j} - \hat{y}_{i,j})^2$



$$\hat{y} = 1$$

- Let's consider one sample, one neurone per layer:
  - $\varepsilon(y_{i,1}, \hat{y}_{i,1})$  depends on  $y_{i,1}$  (or  $r_{3,1}$ )
  - Let get rid of  $i$
  - $y_1 = r_{3,1} = f(g(r_{2,1})) = \sigma(g(r_{2,1})) = \sigma(r_{2,1}w_{3,1} + b_{3,1})$ 
    - For simplicity,  $w_0$  is substituted with  $b$
  - $f(g(r_{2,1}))$  depends on  $g(r_{2,1})$
  - $g(r_{2,1})$  depends on:
    - $r_{2,1}$
    - $\tilde{w}_{3,1}$  (including the bias)
      - This means the weights of the first neuron in the third layer.



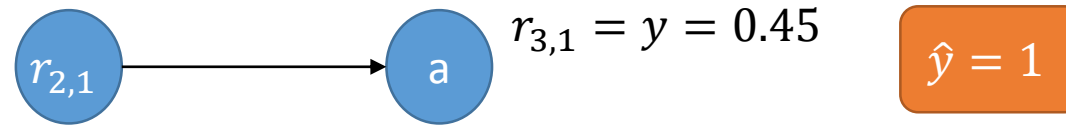
- By the chain rule:

$$\frac{\partial \varepsilon(y_1, \hat{y}_1)}{\partial w_{3,1}} = \frac{\partial \varepsilon(y_1, \hat{y}_1)}{\partial f(g(r_{2,1}))} \frac{\partial f(g(r_{2,1}))}{\partial g(r_{2,1})} \frac{\partial g(r_{2,1})}{\partial w_{3,1}}$$

- $\frac{\partial \varepsilon(y_1, \hat{y}_1)}{\partial f(g(r_{2,1}))} = 2(y_1 - \hat{y}_1)$
- $\frac{\partial f(g(r_{2,1}))}{\partial g(r_{2,1})} = \sigma'(g(r_{2,1}))$
- $\frac{\partial g(r_{2,1})}{\partial w_{3,1}} = r_{2,1}$

Remember:

$$\Delta w_{3,1} = -\rho \left. \frac{\partial \varepsilon(y_1, \hat{y}_1)}{\partial w_{3,1}} \right|_{w_{3,1}=w_{3,1}old}$$



- In the subsequent iteration:

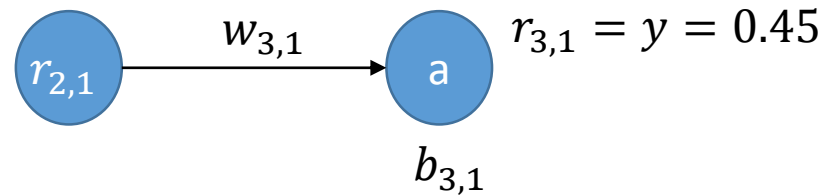
$$w_{3,1}(new) = w_{3,1}(old) - \rho 2(y_1 - \hat{y}_1)\sigma' \left( g(r_{2,1}) \right) r_{2,1}$$

- For  $N$  samples:

$$w_{3,1}(new) = w_{3,1}(old) - \rho \sum_{i=1}^N 2(y_{i,1} - \hat{y}_{i,1})\sigma' \left( g(r_{2,1}) \right) r_{2,1}$$

We can do the same thing for  $b$





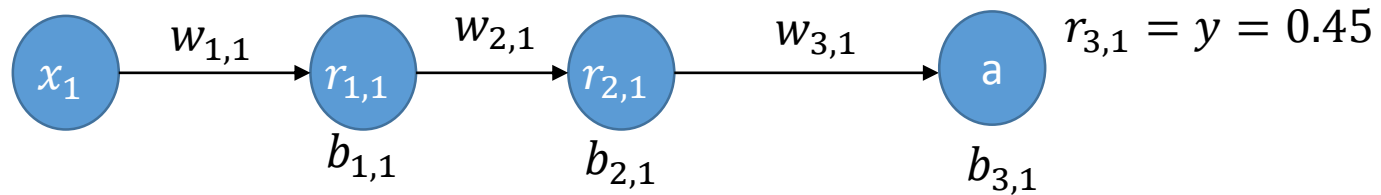
$$\hat{y} = 1$$

- Since we don't have only one sample but  $N$  samples:

$$\frac{\partial \varepsilon(y_{i=1:N,1}, \hat{y}_{i=1:N,1})}{\partial w_{3,1}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \varepsilon(y_{i,1}, \hat{y}_{i,1})}{\partial w_{3,1}}$$

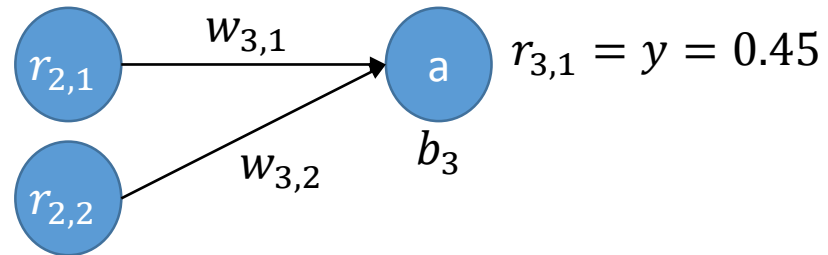
- $\varepsilon(y_1, \hat{y}_1)$  is also sensitive to  $r_{2,1}$ :

$$\begin{aligned} \blacksquare \frac{\partial \varepsilon(y_1, \hat{y}_1)}{\partial r_{2,1}} &= \frac{\partial \varepsilon(y_1, \hat{y}_1)}{\partial f(g(r_{2,1}))} \frac{\partial f(g(r_{2,1}))}{\partial g(r_{2,1})} \frac{g(r_{2,1})}{\partial r_{2,1}} \\ \blacksquare \frac{\partial g(r_{2,1})}{\partial r_{2,1}} &= w_{3,1} \end{aligned}$$



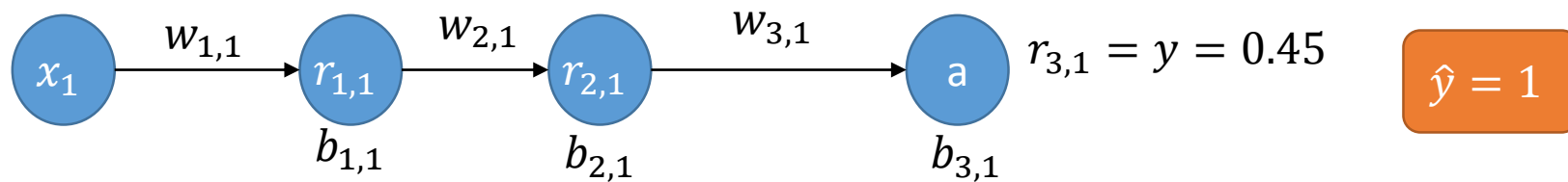
$$\hat{y} = 1$$

- But the cost function is also sensitive to previous layers:
  - $\varepsilon(y_1, \hat{y}_1)$  depends on  $y_1$  (or  $r_{3,1}$ )
  - $y_1 = r_{3,1} = f(g(r_{2,1})) = \sigma(g(r_{2,1})) = \sigma(r_{2,1}w_{3,1} + b_{3,1})$
  - $f(g(r_{2,1}))$  depends on  $g(r_{2,1})$
  - $g(r_{2,1})$  depends on  $r_{2,1}$  and  $\tilde{w}_{3,1}$
  - $r_{2,1} = f(g(r_{1,1}))$
  - $f(g(r_{1,1}))$  depends on  $g(r_{1,1})$
  - $g(r_{1,1})$  depends on  $r_{1,1}$  and  $\tilde{w}_{2,1}$



$$\hat{y} = 1$$

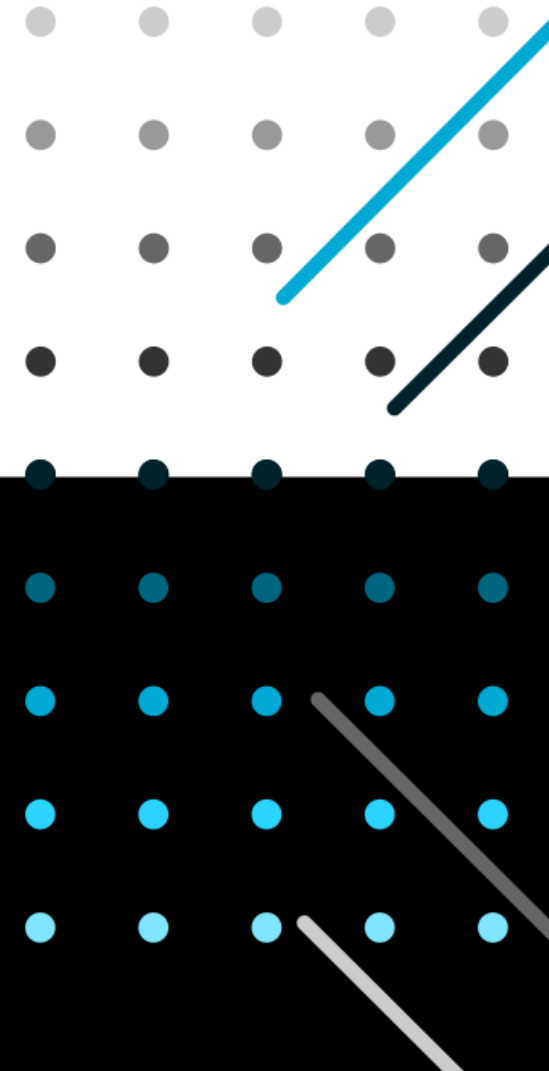
- We have two neurones in  $k=2$ :
  - $y_1 = r_{3,1} = f(g(r_{2,1}, r_{2,2})) = \sigma(\sum_{j=1}^2 (r_{2,j} w_{3,j}) + b_3)$
  - We generalize this for any number of neurones:
  - $r_{k,u} = \sigma \left( \sum_{j=1}^{|k-1|} (r_{k-1,j} w_{k,j}) + b_{k,u} \right)$



- $$\nabla \varepsilon(y_{i=1:N,1}, \hat{y}_{i=1:N,1})_{\tilde{\mathbf{w}}} = \left[ \frac{\partial \varepsilon(y_{i=1:N,1}, \hat{y}_{1:N,1})}{\partial w_{1,1}}, \frac{\partial \varepsilon(y_{1:N,1}, \hat{y}_{1:N,1})}{\partial b_{1,1}}, \dots, \frac{\partial \varepsilon(y_{1:N,1}, \hat{y}_{1:N,1})}{\partial w_{3,1}}, \frac{\partial \varepsilon(y_{1:N,1}, \hat{y}_{1:N,1})}{\partial b_{3,1}} \right]^T$$

# Summary

---



- From Perceptron classifier to Multi-Layer Perceptron
  - XOR problem
  - Activation function
  - Two-Layer Perceptron
- Learning (Optimization)
  - Gradient Decent
  - Backpropagation

# Thank you!



## Zeyd Boukhers

E-mail: [Boukhers@uni-koblenz.de](mailto:Boukhers@uni-koblenz.de)

Phone: +49 (0) 261 287-2765

Web: [Zeyd.Boukhers.com](http://Zeyd.Boukhers.com)

University of Koblenz-Landau  
Universitätsstr. 1  
56070 Koblenz

