



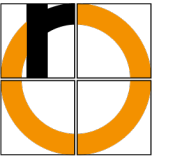
Deep Learning

Feedforward neural networks: Multi Layer Perceptron

Technische Hochschule Rosenheim
Sommer 2023
Prof. Dr. Jochen Schmidt

Many of the slides presented here are based on the Deep Learning Slides Summer Semester 2020, courtesy of **A. Maier, V. Christlein, K. Breininger, F. Denzinger, F. Thamm**, Pattern Recognition Lab, Friedrich-Alexander-University Erlangen-Nürnberg.
<https://lme.tf.fau.de/>

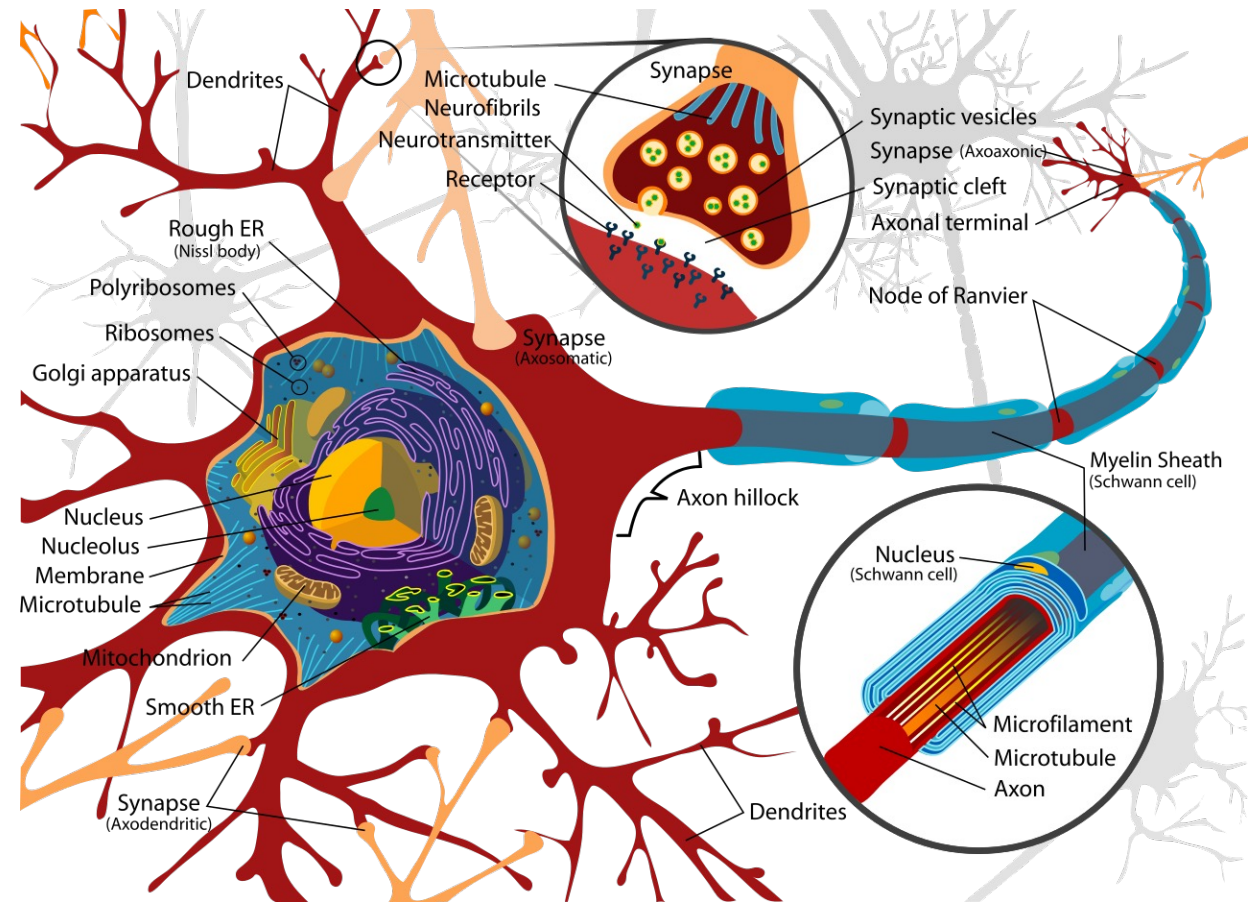
- Perceptrons
- MLP-Training Basics
- MLP – Layer Abstraction
- Universal Approximation Theorem



Perceptrons

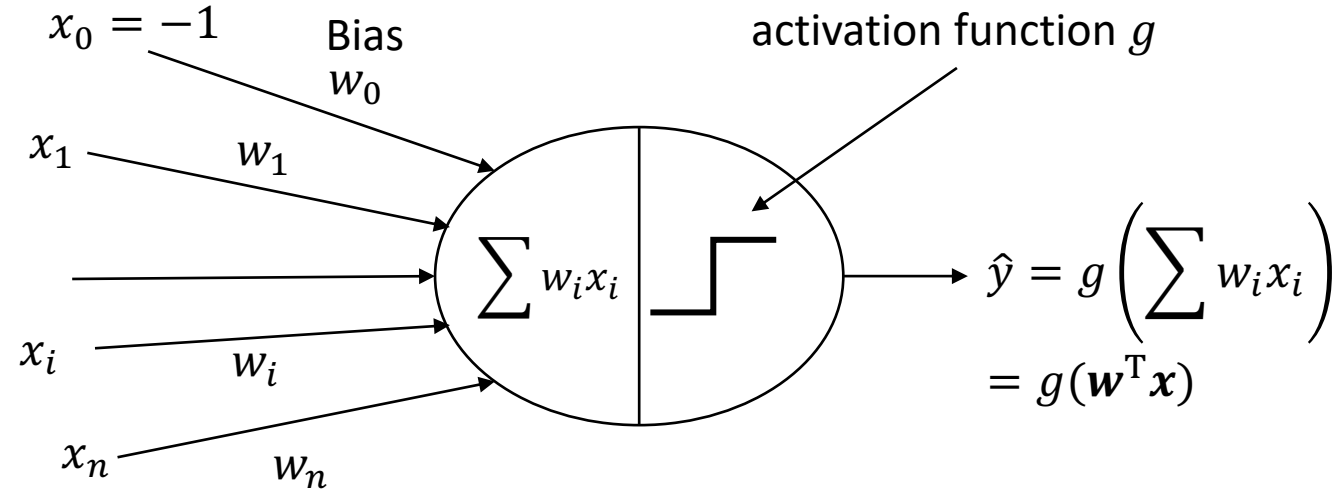
Human brain:

- 10^{11} Neurons
- more than 20 different types
- 10^{14} synapses (connections)
- Cycle time: 1-10 msec

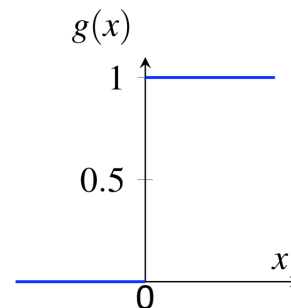




- simplified model of a real neuron
- by Frank Rosenblatt 1957
- binary classifier $y \in \{0, 1\}$
- $\mathbf{w} = (w_0, w_1, \dots, w_n)$ set of weights (w_0 : bias)
- $\mathbf{x} = (-1, x_1, \dots, x_n)$ input feature vector
- Output: linear function of input + (non-linear) threshold



$$\hat{y} = \begin{cases} 1 & \mathbf{w}^T \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases}$$



Decision threshold is always zero?

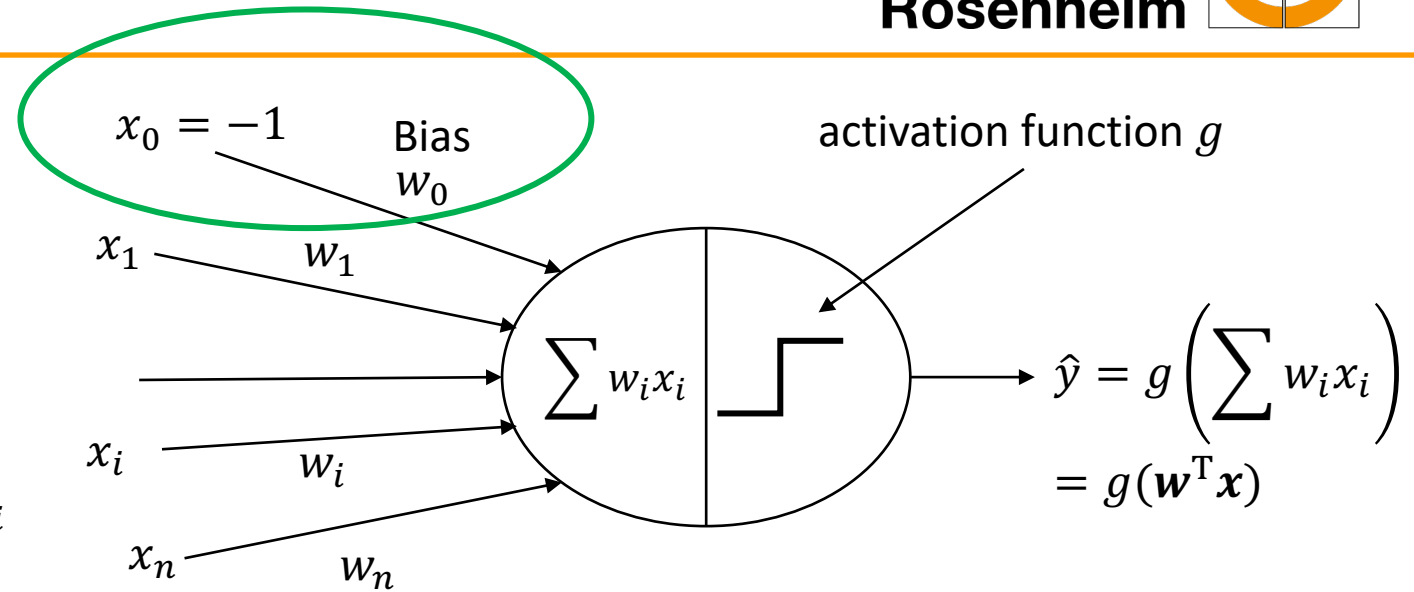
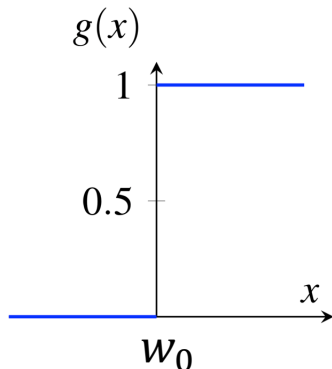
Can we set it to an arbitrary value?

YES: Bias weight = threshold position!

$$\hat{y} = \begin{cases} 1 & \mathbf{w}^T \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbf{w}^T \mathbf{x} = \sum_{i=0}^n w_i x_i = w_0 x_0 + \sum_{i=1}^n w_i x_i = -w_0 + \sum_{i=1}^n w_i x_i$$

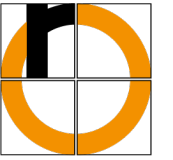
$$\sum_{i=0}^n w_i x_i > 0 \quad \rightarrow \quad \sum_{i=1}^n w_i x_i > w_0$$



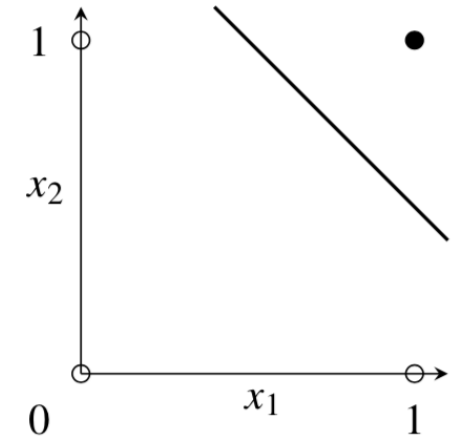
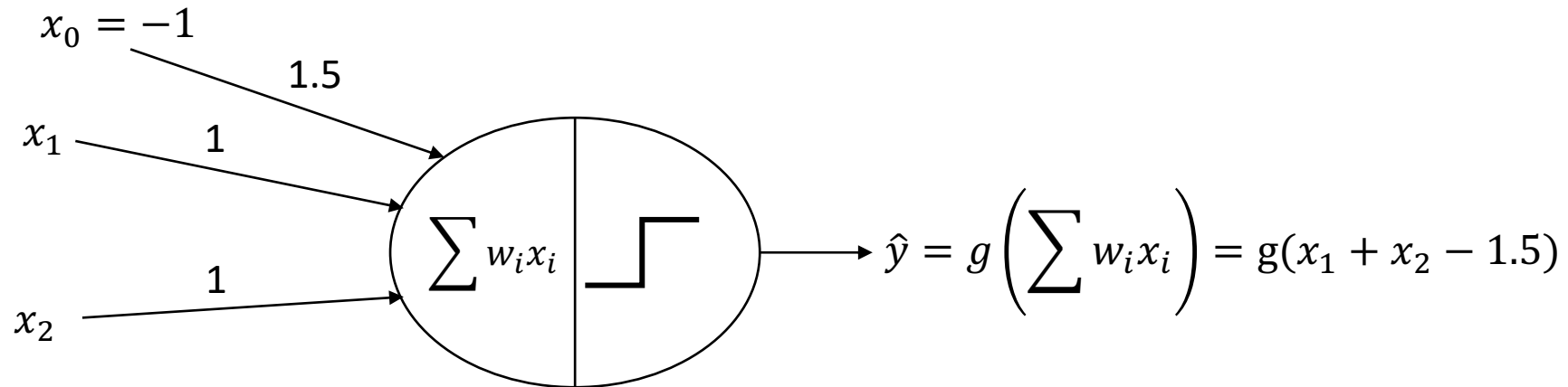
Why?

- The threshold becomes just another weight
- The weights are learned during training
- Thus, the threshold can be trained automatically

Example: Boolean AND



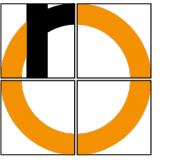
A perceptron can compute the Boolean AND-function:



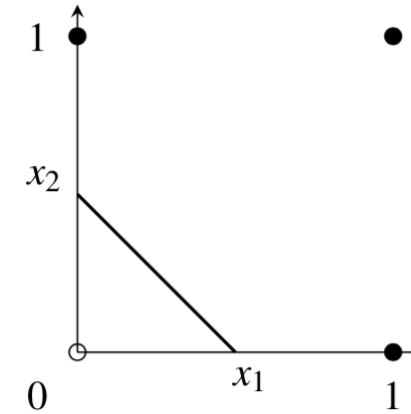
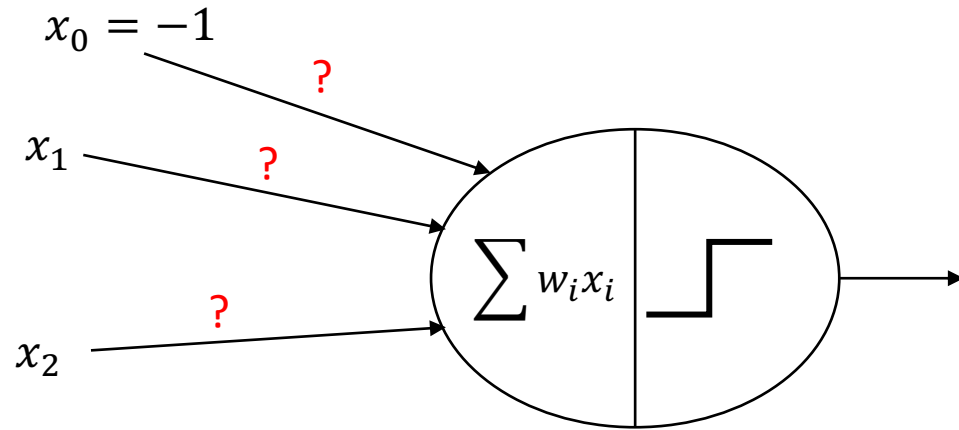
Could we use other weights?

x_1	x_2	$\sum w_i x_i$	$g\left(\sum w_i x_i\right)$
0	0	-1.5	0
0	1	-0.5	0
1	0	-0.5	0
1	1	0.5	1

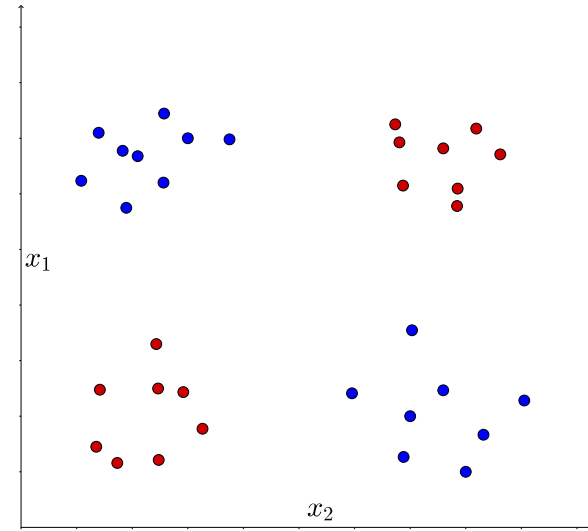
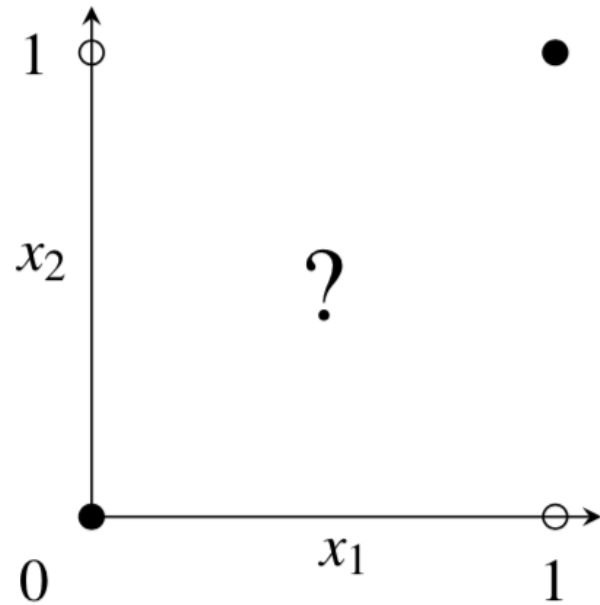
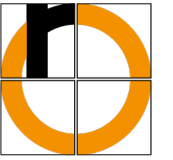
Exercise: Boolean OR



What would the weights look like for the Boolean OR-function?



Boolean XOR?



Observation: A single perceptron computes a linear separation of the feature space



$f(\mathbf{x}) = -w_0 + \sum_{i=1}^n w_i x_i = \mathbf{w}'^T \mathbf{x}' - w_0$ is the equation of an n -dimensional hyperplane

$\mathbf{w}' = (w_1, \dots, w_n)$ set of weights without bias w_0

$\mathbf{x}' = (x_1, \dots, x_n)$ input feature vector

- normal vector: \mathbf{w}'

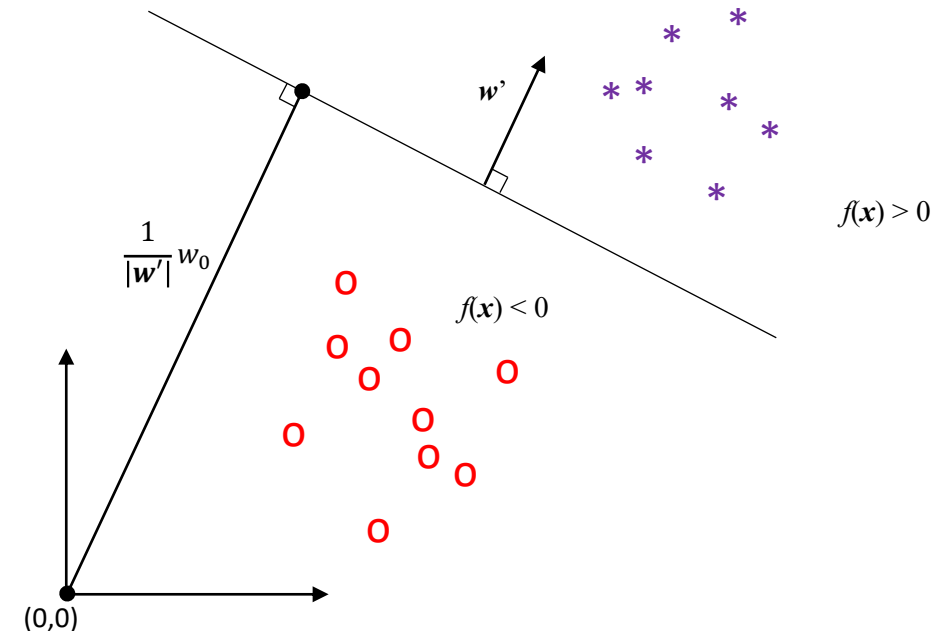
- point on plane: $f(\mathbf{x}') = 0$
- point above plane: $f(\mathbf{x}') > 0$
- point below plane: $f(\mathbf{x}') < 0$
- "above" = in direction of plane normal

- Signed distance d of a point to hyperplane

- normalize \mathbf{w}' , such that $|\mathbf{w}'| = 1$:

$$d = f'(\mathbf{x}') = \frac{1}{|\mathbf{w}'|} f(\mathbf{x}') = \frac{1}{|\mathbf{w}'|} \mathbf{w}'^T \mathbf{x}' - \frac{1}{|\mathbf{w}'|} w_0$$

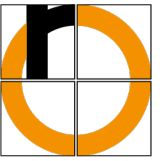
- distance of plane from origin: $\frac{1}{|\mathbf{w}'|} w_0$



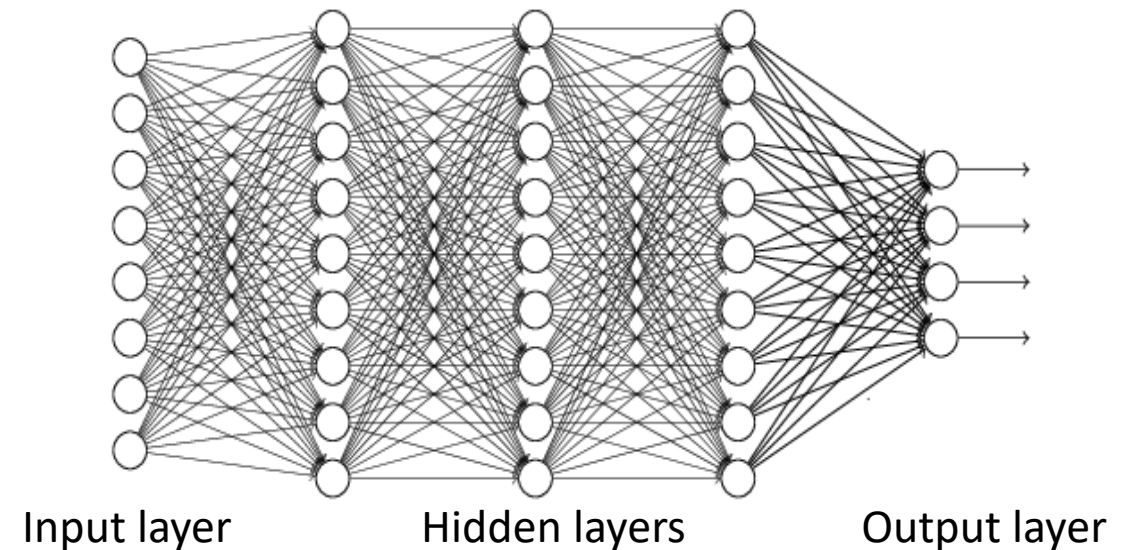
P.S.: The perceptron decision rule looks exactly like the one for an SVM (Support Vector Machine)

Can we learn the weights automatically?

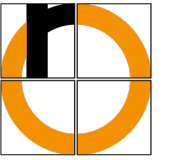
- We need a training data set: $\mathcal{S} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$
- Weights $\mathbf{w} = (w_0, w_1, \dots, w_n)$ initialized randomly
- Repeat until convergence or for a defined number of iterations:
 - for each training sample $(\mathbf{x}_e, y_e) \in \mathcal{S}$
 - compute output $\hat{y}_e = g(\mathbf{w}^T \mathbf{x}_e)$
 - compute residual error: $\varepsilon = y_e - \hat{y}_e$
 - update weights: $\mathbf{w} := \mathbf{w} + \alpha \varepsilon \mathbf{x}_e$
- α : **Learning Rate**, e.g., $\alpha = 0.1$
- this algorithm updates weights after each sample
- other strategies:
 - process all samples, sum changes, update weights after each iteration
 - process small batches, sum changes, then update



- How to obtain nonlinear decision boundaries? Use more than one neuron.
- Neurons are arranged in layers
- Layer i is (fully) connected with layer $i+1$
 - no connections within layer
 - no connections to any other layers
 - no feedback
- Information flow from one layer to the next:
feed-forward
- network has no internal state



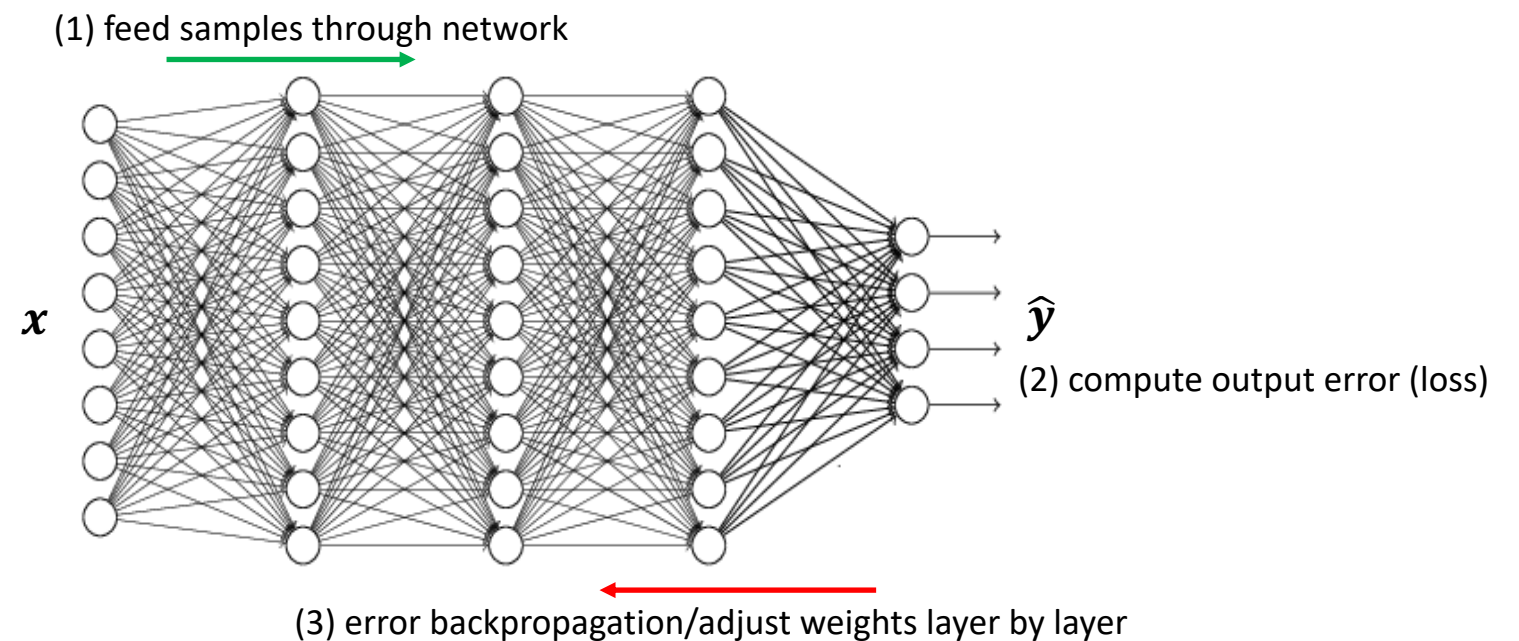
- number of input/output units is problem dependent
- number of layers/hidden units determined by developer
 - ... we will discuss this in more detail later



MLP-Training Basics

Training = Determine weights

- define loss function L of output neurons
- minimize loss



Objective: Find optimal weights \mathbf{w} for all layers

- Abstraction: Consider the whole neural network as a function

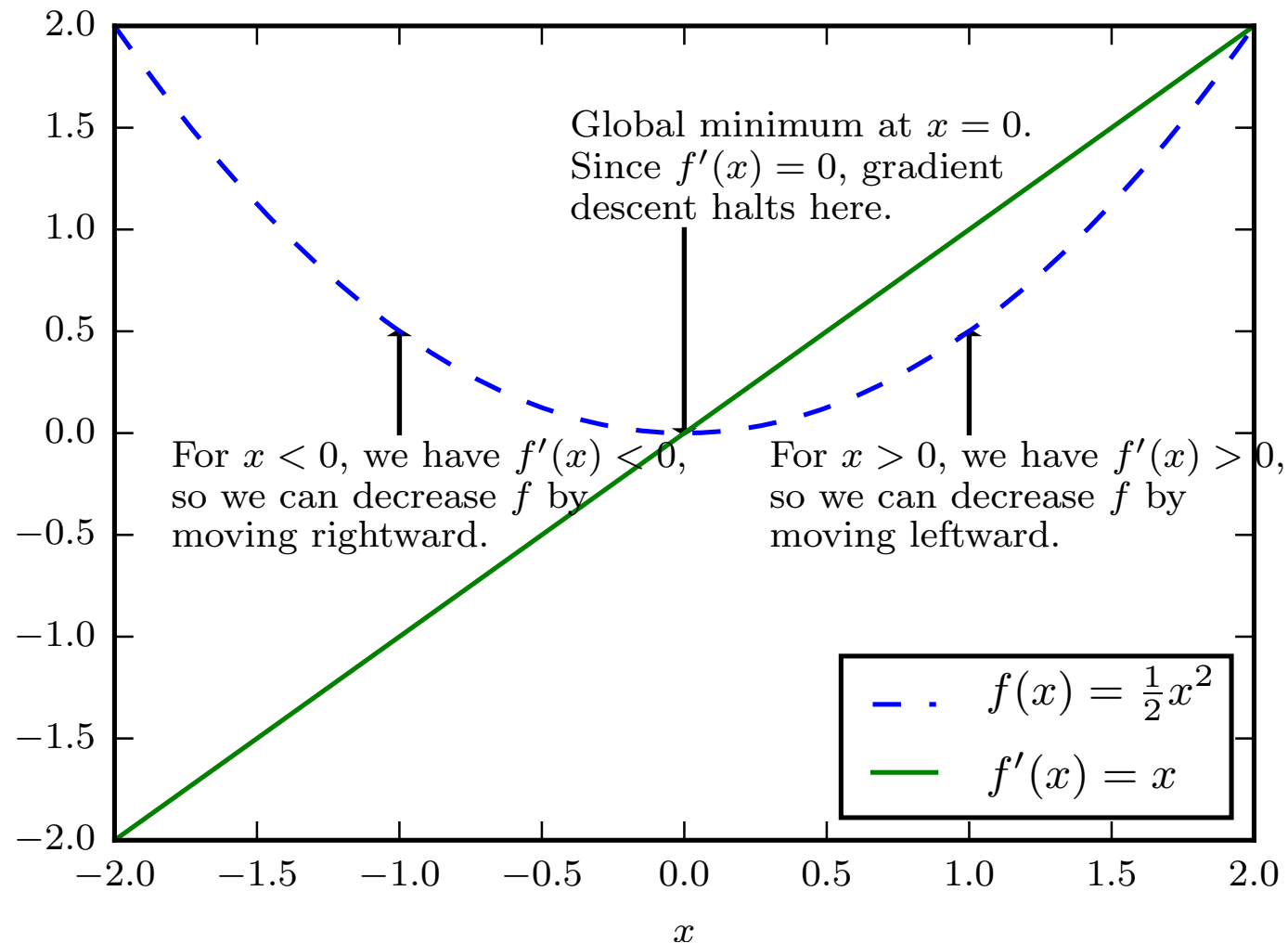
$$L(\mathbf{w}, \mathbf{x}, \mathbf{y})$$

- Consider all M training samples

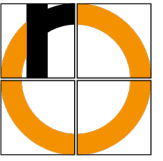
$$\frac{1}{M} \sum_{m=1}^M L(\mathbf{w}, \mathbf{x}_m, \mathbf{y}_m)$$

- Minimize (with unknowns \mathbf{w})
 - computing the first derivative and solving for zero crossings will not work (highly nonlinear)
 - method of choice: Gradient Descent

Gradient Descent Illustration



© Goodfellow, I., Bengio, Y., Courville, A.: *Deep Learning*, MIT Press, 2017.



$$\frac{1}{M} \sum_{m=1}^M L(\mathbf{w}, \mathbf{x}_m, \mathbf{y}_m)$$

1. Initialize weights \mathbf{w}
2. Iterate until convergence (or for a defined number of iterations)

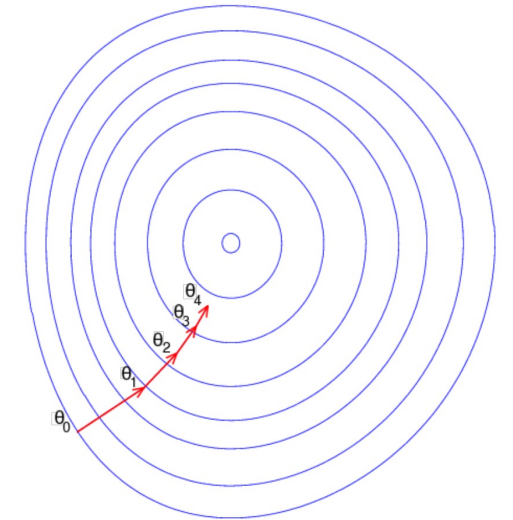
$$\mathbf{w}^{(k+1)} := \mathbf{w}^{(k)} - \alpha \nabla_{\mathbf{w}} \frac{1}{M} \sum_{m=1}^M L(\mathbf{w}, \mathbf{x}_m, \mathbf{y}_m)$$

α : Learning Rate

$$\nabla_x f(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

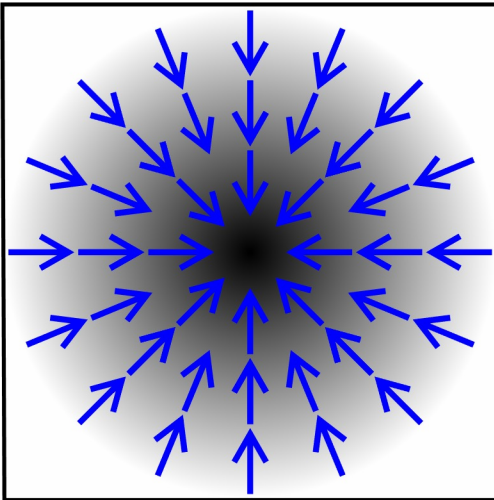
Gradient of f

- vector direction: direction of steepest increase
- vector magnitude: rate/strength of increase

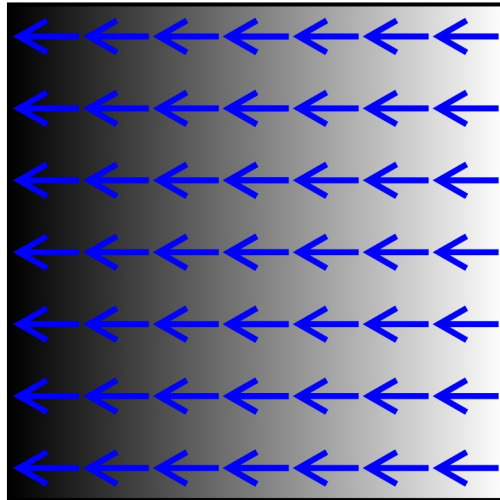


Gradient of f

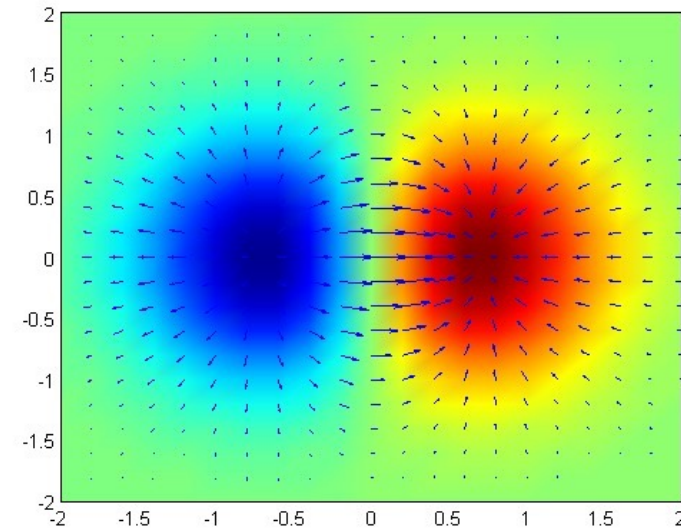
- vector direction: direction of steepest increase
- vector magnitude: rate/strength of increase



© JoKalliauer, Gradient2, CC BY-SA 2.5



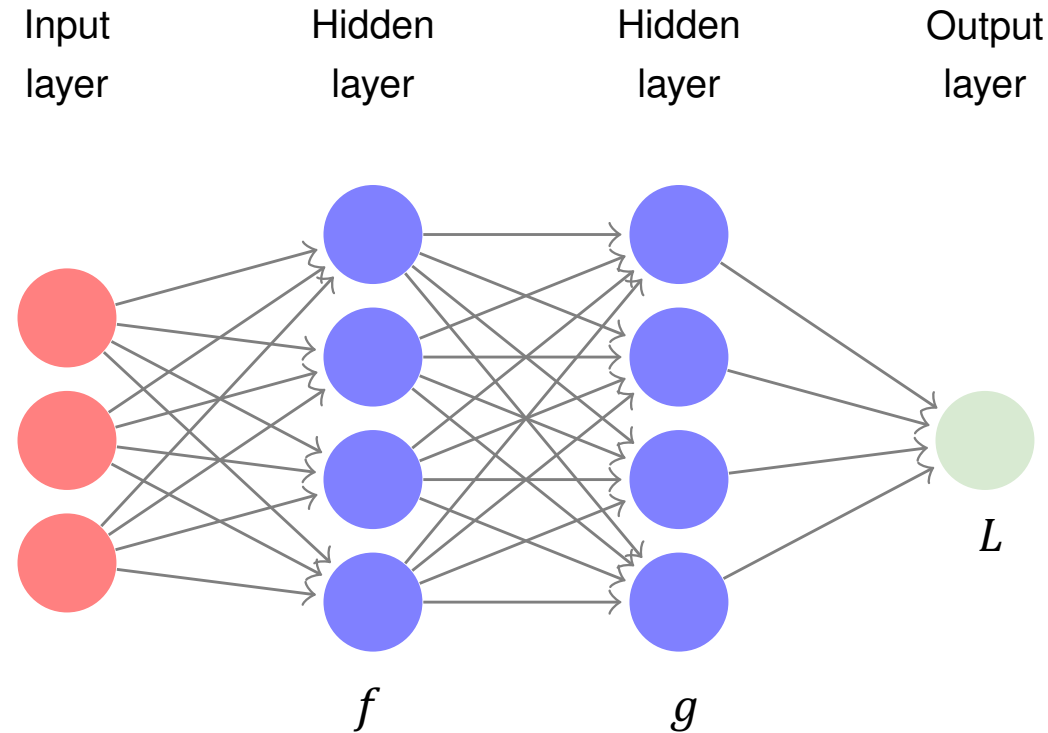
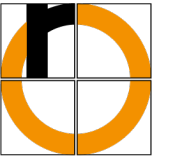
Function values in grayscale:
white = low values
dark = high values



© Vivekj78, Gradient of a Function, CC BY-SA 3.0

Gradient of $f(x, y) = \frac{x}{e^{x^2 + y^2}}$ as pseudo color plot
blue/cold = low values
red/hot = high values

What does this look like in a Neural Network?



A complex network can be looked at as composed functions:

$$L(\mathbf{w}, \mathbf{x}, \mathbf{y}) = L(g(f(\mathbf{x}, \mathbf{w}_f), \mathbf{w}_g), \mathbf{y})$$

How do we Calculate Derivatives in a Neural Network?

Example

- function $\hat{y} = f(\mathbf{x}) = f(x_1, x_2) = (2x_1 + 3x_2)^2 + 3$
- evaluate $\frac{\partial}{\partial x_1} f(1, 3)$

Two methods:

- numerical derivative using finite differences
- analytic derivative

Definition of derivative: $f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$

This is asymmetric; due to finite precision, we prefer the symmetric definition:

$$f'(x) = \lim_{h \rightarrow 0} \frac{f\left(x + \frac{h}{2}\right) - f\left(x - \frac{h}{2}\right)}{h}$$

Thus, we get the derivative at the exact position of x

- function $\hat{y} = f(\mathbf{x}) = f(x_1, x_2) = (2x_1 + 3x_2)^2 + 3$
- evaluate $\frac{\partial}{\partial x_1} f(1, 3)$

$$f'(x) = \lim_{h \rightarrow 0} \frac{f\left(x + \frac{h}{2}\right) - f\left(x - \frac{h}{2}\right)}{h}$$

Set $h = 2 \cdot 10^{-2}$. We get

$$\frac{\partial}{\partial x_1} f(1, 3) = \frac{(2(1 + 10^{-2}) + 9)^2 + 3 - ((2(1 - 10^{-2}) + 9)^2 + 3)}{2 \cdot 10^{-2}} = \frac{124.4404 - 123.5604}{2 \cdot 10^{-2}} = 43.999$$

(the exact value is 44)

- for practical purposes it often suffices to use $h = 1 \cdot 10^{-5}$
- for a more accurate derivative use $h = \sqrt[3]{\varepsilon_f} \cdot x$
 - ε_f is the fractional accuracy with which f is computed
 - for a simple function this is approximately the machine accuracy ($\varepsilon_m \approx 10^{-7}$ for 32 Bit floats)
 - for complex functions this may be considerably larger than machine accuracy
 - x is the position, where f is being evaluated
 - prevent division by zero when $x = 0$
- see [Pre07, pp.186-189] for more details
- easy to use
- we only need to be able to evaluate functions
 - can be computed for any function, even if we do not have the derivative
- however:
 - computationally inefficient
 - not as accurate as using the analytical derivative
- frequently used for validating implementations of analytical derivatives

- function $\hat{y} = f(\mathbf{x}) = f(x_1, x_2) = (2x_1 + 3x_2)^2 + 3$
- evaluate $\frac{\partial}{\partial x_1} f(1, 3)$

$$\begin{aligned}\frac{\partial}{\partial x_1} f(x_1, x_2) &= \frac{\partial}{\partial x_1} ((2x_1 + 3x_2)^2 + 3) \\ &= \frac{\partial}{\partial x_1} (2x_1 + 3x_2)^2 + \frac{\partial}{\partial x_1} 3 \\ &= \frac{\partial}{\partial x_1} (2x_1 + 3x_2)^2 \\ &= \frac{\partial}{\partial z} z^2 \frac{\partial}{\partial x_1} (2x_1 + 3x_2) \\ &= 2z \cdot 2 = 2 \cdot (2x_1 + 3x_2) \cdot 2 = (8x_1 + 12x_2)\end{aligned}\tag{2}$$

$$\frac{\partial}{\partial x_1} f(1, 3) = (8 \cdot 1 + 12 \cdot 3) = 44$$

Required Rules:

1. $\frac{d}{dx} \text{const} = 0$
2. $\frac{d}{dx}$ is a linear operation
3. Monomials: $\frac{d}{dx} x^n = n \cdot x^{n-1}$
4. Chain rule: $\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \cdot \frac{d}{dx} g(x)$

(1)

(4) and $(2x_1 + 3x_2) = z$

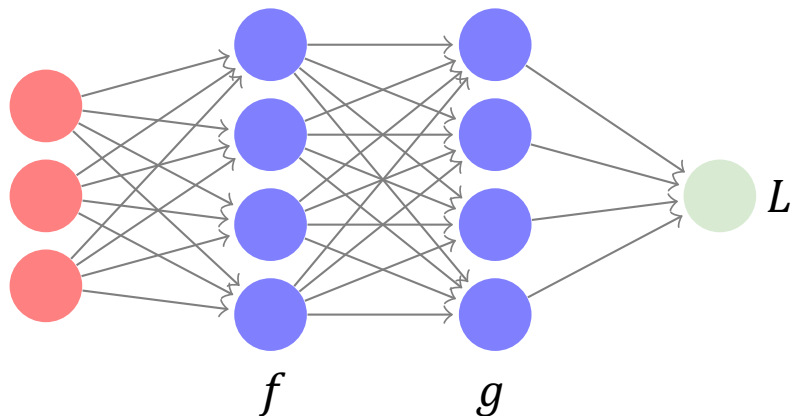
(1) + (2)



- **chain rule** and **linearity** enable us to **decompose** complex functions
- analytic derivatives have to be calculated manually
- computationally more efficient and more precise than finite differences

Can we compute analytic gradients of the network loss function automatically?

YES: Provided we have the analytic derivative of the activation function – just use chain rule and linearity

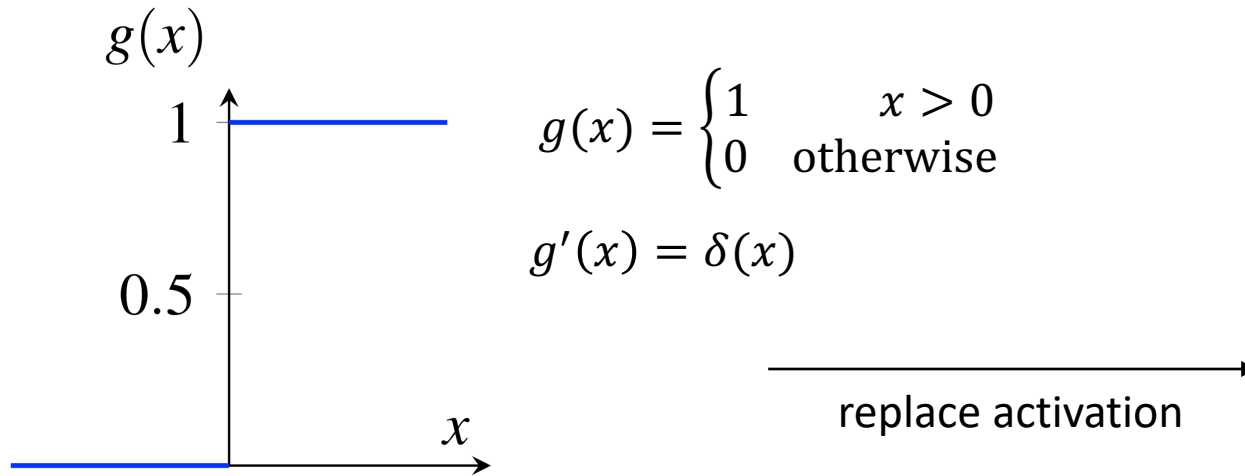


$$L(\mathbf{w}, \mathbf{x}, \mathbf{y}) = L(g(f(\mathbf{x}, \mathbf{w}_f), \mathbf{w}_g), \mathbf{y})$$

Ooops – Activation Derivatives Required!



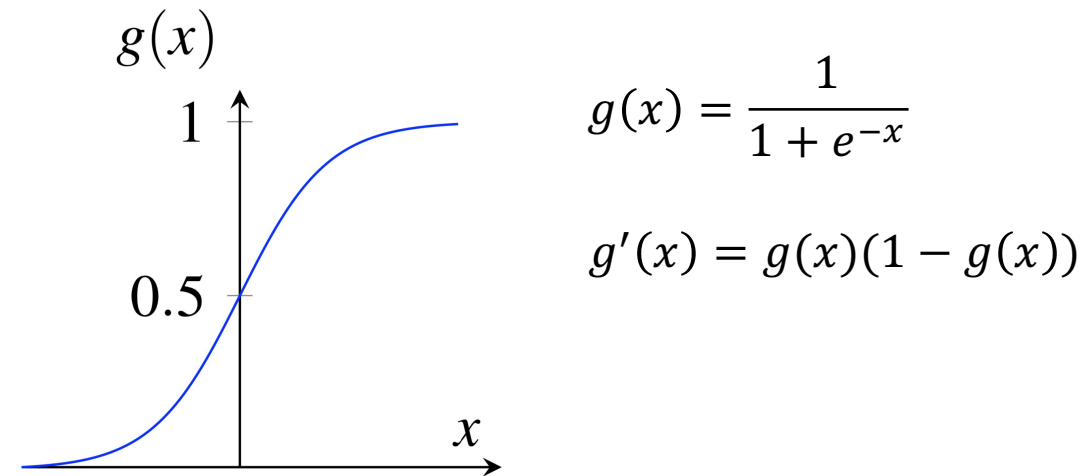
As activation we used the Heaviside step function:



- not differentiable in the usual sense
- gradient vanishes almost everywhere

Unsuitable for gradient descent!

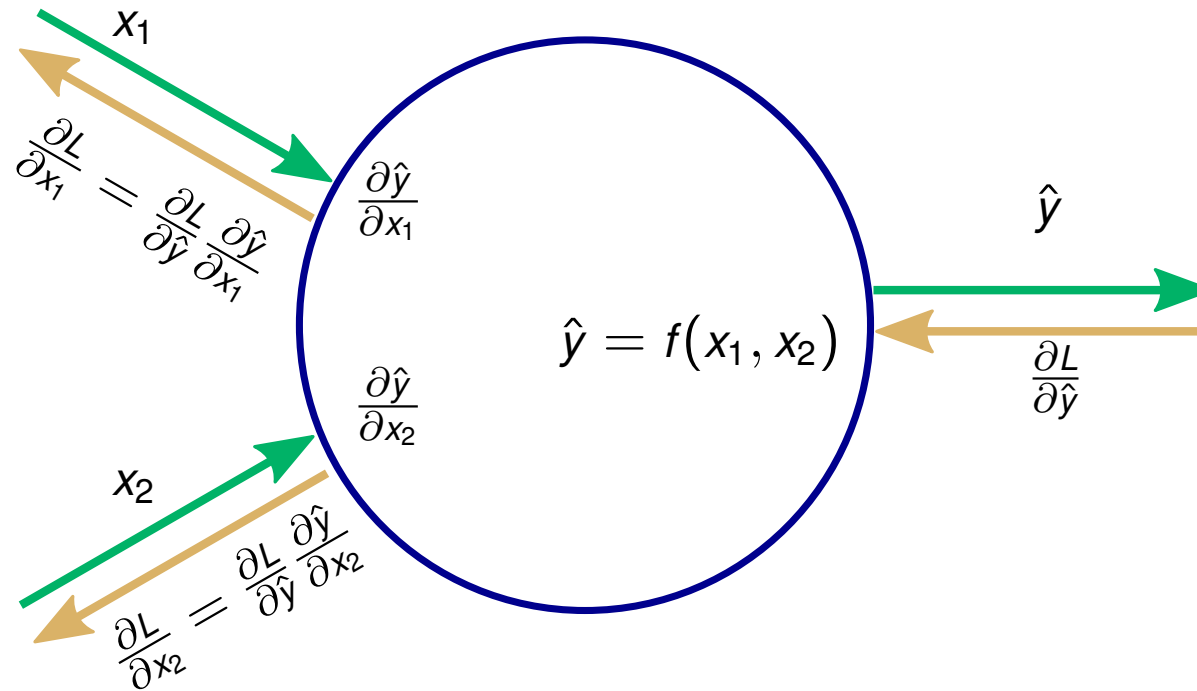
Sigmoid (also: logistic function)



- smooth approximation of step function
- efficient computation of derivative
- gradient still vanishes eventually...

we will discuss other activation functions later

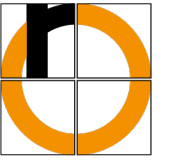
Backpropagation algorithm



- does not show weights
- to train weights, we also require the derivatives w.r.t. the weights (for update in gradient descent)

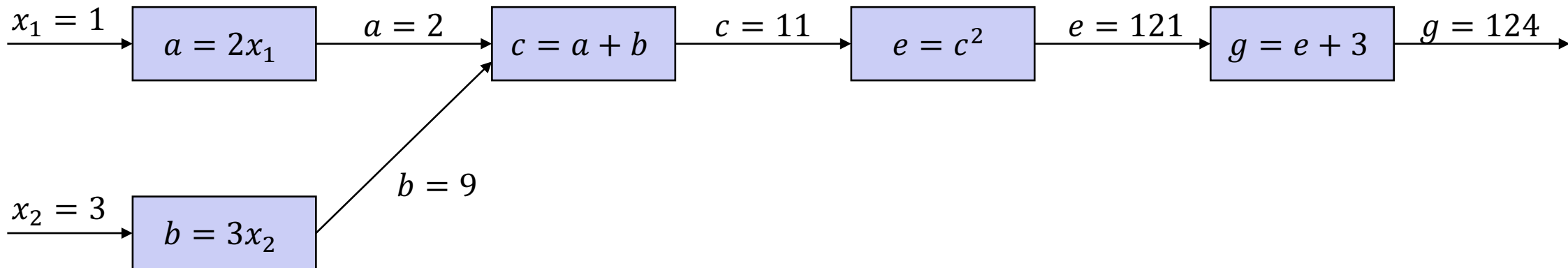
1. Forward pass: Compute activations
2. Backward pass: Recursively apply chain rule

Backpropagation – Example / Forward Pass

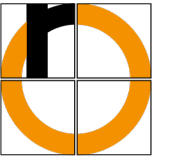


- function $\hat{y} = f(\mathbf{x}) = f(x_1, x_2) = (2x_1 + 3x_2)^2 + 3$
- evaluate $\frac{\partial}{\partial x_1} f(1, 3)$

$$\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \cdot \frac{d}{dx} g(x)$$

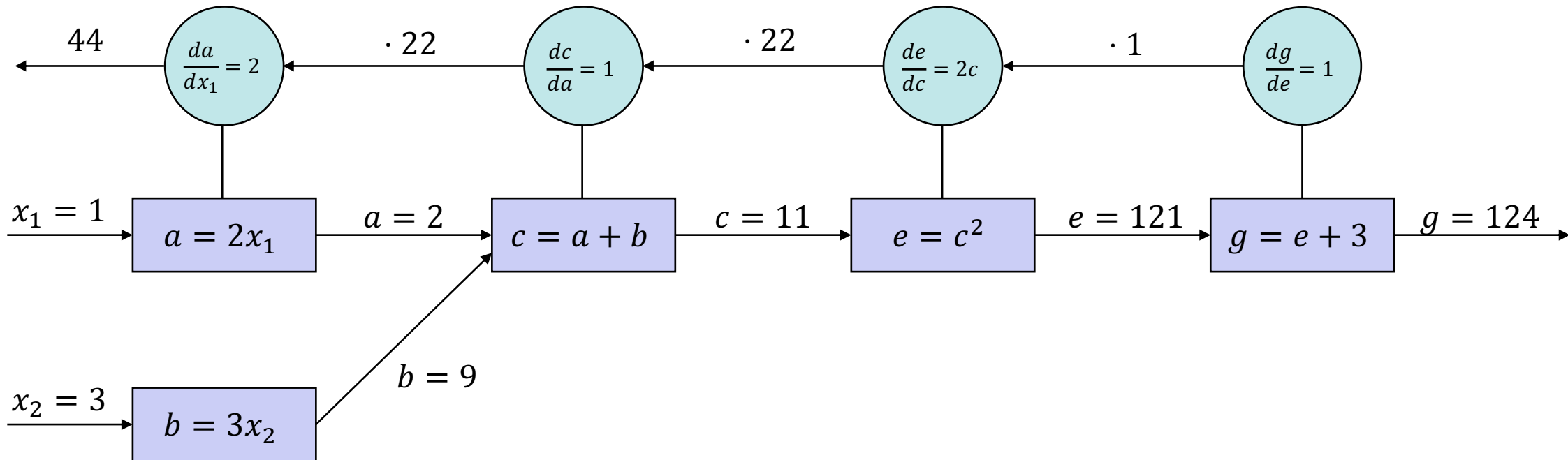


Backpropagation – Example / Backward Pass



- function $\hat{y} = f(\mathbf{x}) = f(x_1, x_2) = (2x_1 + 3x_2)^2 + 3$
- evaluate $\frac{\partial}{\partial x_1} f(1, 3)$

$$\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \cdot \frac{d}{dx} g(x)$$

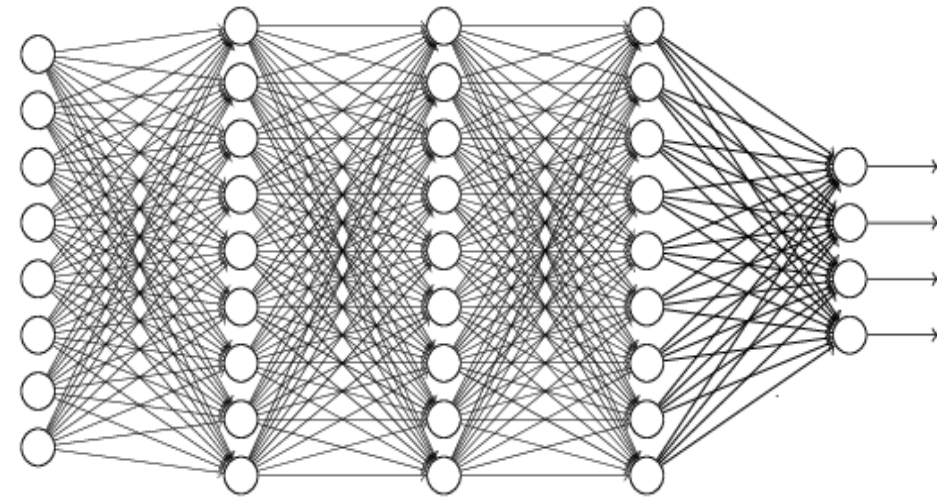


What do these two systems have in common?



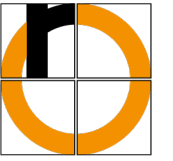
<https://youtu.be/esfpcnQW6qs>

- Both suffer from positive feedback
- This can lead to disaster

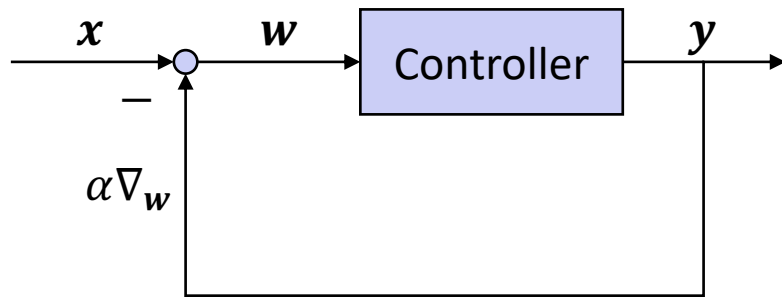


Backpropagation:

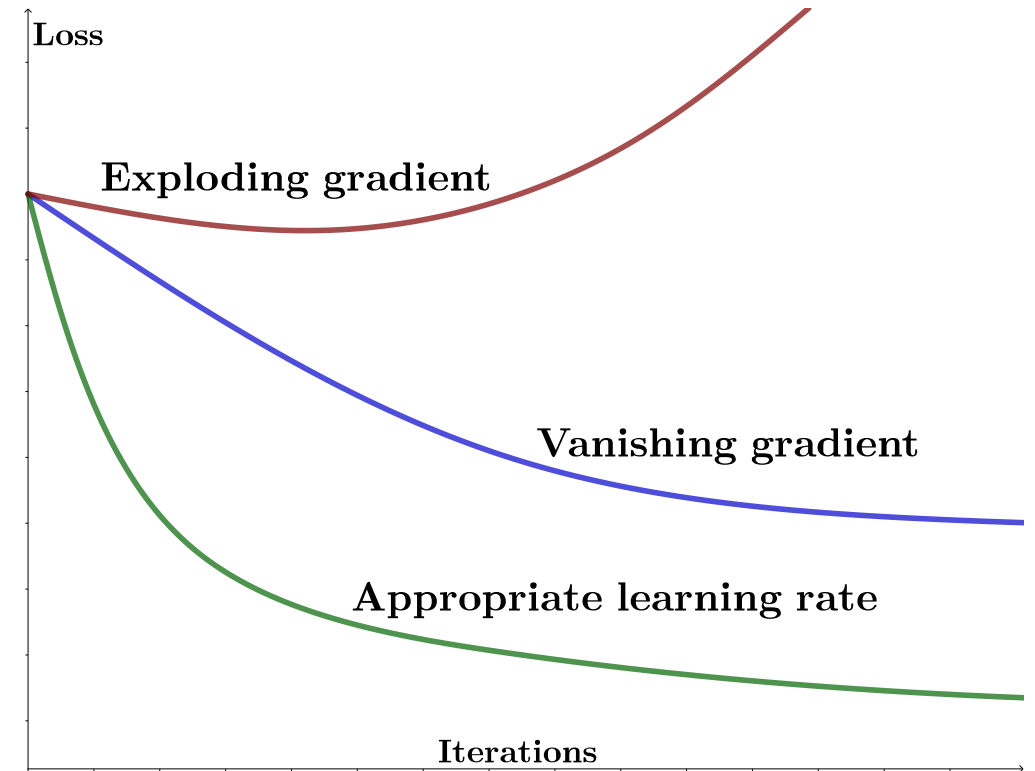
Long chains of float multiplications of potentially large and small numbers → positive feedback



Analogy to control theory

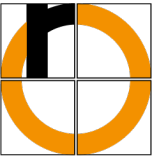


- if α is too large → **positive feedback**
- if α is too small → **negative feedback**
- choice of α is critical for learning



- loss grows without bounds
- gradient vanishes

- built around the **chain rule**
- uses a **forward pass** computing the function output
- and a **backward pass** for computing the gradient
- computationally very **efficient** by using a dynamic programming approach
 - store inputs and intermediate results during forward pass, so that they are available for backward pass
- product of partials \rightarrow numerical **errors multiply**
- product of partials \rightarrow **vanishing** or **exploding** gradient
 - can be mitigated to a certain degree by choosing other activation functions
 - \rightarrow we will discuss these later




MLP – Layer Abstraction

- we introduced layers but computed everything on individual nodes
- it is convenient to add further abstraction
 - more compact and clear notation
 - we can talk directly about gradients of entire layers
 - training is implemented on GPUs: matrix/vector notation required

- recall a single neuron: $\hat{y} = g(\sum w_i x_i) = g(\mathbf{w}^T \mathbf{x})$
- for simplicity of notation, we will use the identity mapping as activation – i.e., we omit g
- assume we have M neurons $\rightarrow M$ sets of weights \mathbf{w}_m for $m \in \{1, \dots, M\}$.
- for each neuron we get $\hat{y}_m = \mathbf{w}_m^T \mathbf{x}$
- we collect all outputs in a vector and all weights in a matrix:

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

$$\mathbf{W} = \begin{array}{c} \text{\#weights} \\ \text{\#neurons} = M \end{array} \begin{array}{c} (= \dim(\mathbf{x})) \\ \text{\#neurons} = M \end{array}$$


- this is known as a **fully connected layer**
- the non-linear activation g is then applied to each element of $\hat{\mathbf{y}}$

- we can formulate error backpropagation using matrix calculus
- the forward pass is: $\hat{\mathbf{y}} = \mathbf{W}\mathbf{x}$
- after the forward pass through all layers, we compute a loss that depends on our loss function L
- two gradients are required for the backward pass:
 - gradient with respect to **weights** for gradient descent: $\frac{\partial L}{\partial \mathbf{W}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{W}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \mathbf{x}^T$ (this is a matrix)
 - gradient with respect to **inputs** for backpropagation: $\frac{\partial L}{\partial \mathbf{x}} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{x}} \frac{\partial L}{\partial \hat{\mathbf{y}}} = \mathbf{W}^T \frac{\partial L}{\partial \hat{\mathbf{y}}}$ (this is a vector)

order to be consistent
with matrix multiplication

- assume we have simple network with no activation function. Forward pass: $\hat{\mathbf{y}} = \mathbf{W}\mathbf{x}$
- we want to find weights minimizing the following loss function (L_2 loss):

$$L(\mathbf{x}, \mathbf{W}, \mathbf{y}) = \frac{1}{2} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 = \frac{1}{2} (\hat{\mathbf{y}} - \mathbf{y})^T (\hat{\mathbf{y}} - \mathbf{y}) = \frac{1}{2} (\mathbf{W}\mathbf{x} - \mathbf{y})^T (\mathbf{W}\mathbf{x} - \mathbf{y})$$

- we get

- $\frac{\partial L}{\partial \hat{\mathbf{y}}} = \hat{\mathbf{y}} - \mathbf{y} = \mathbf{W}\mathbf{x} - \mathbf{y}$

- gradient with respect to weights:

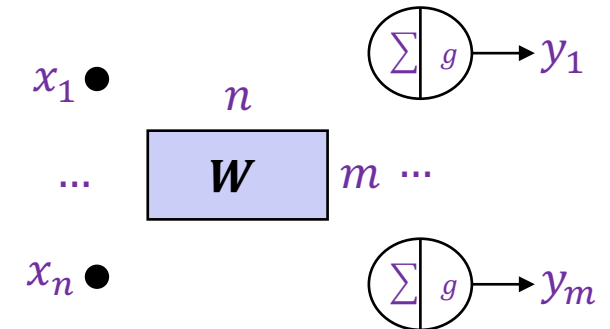
$$\frac{\partial L}{\partial \mathbf{W}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{W}} = (\mathbf{W}\mathbf{x} - \mathbf{y})\mathbf{x}^T$$

$m \times n$ $m \times 1$ $1 \times n$

- gradient with respect to inputs:

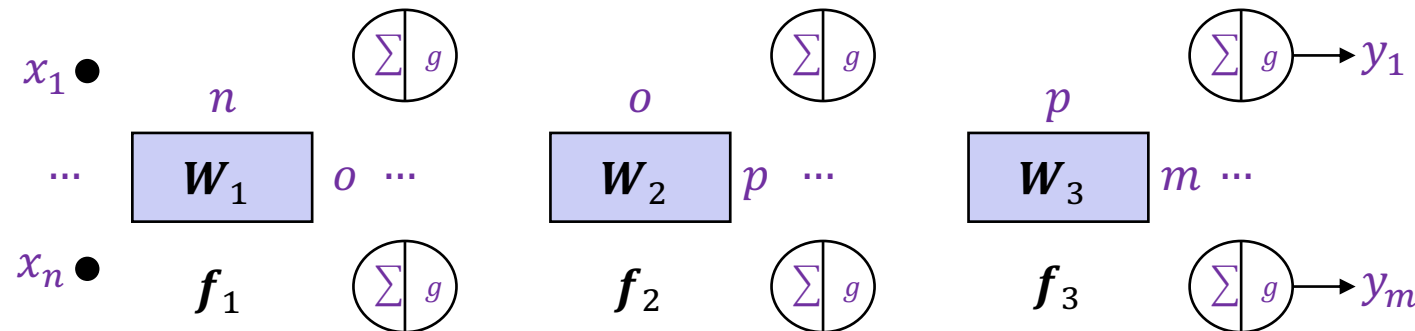
$$\frac{\partial L}{\partial \mathbf{x}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{x}} = \mathbf{W}^T (\mathbf{W}\mathbf{x} - \mathbf{y})$$

$n \times 1$ $n \times m$ $m \times 1$



- now we add two more layers: $\hat{y} = f_3(f_2(f_1(x))) = \mathbf{W}_3\mathbf{W}_2\mathbf{W}_1x$

this is why non-linear activation is required:
without, this will collapse to a single layer!



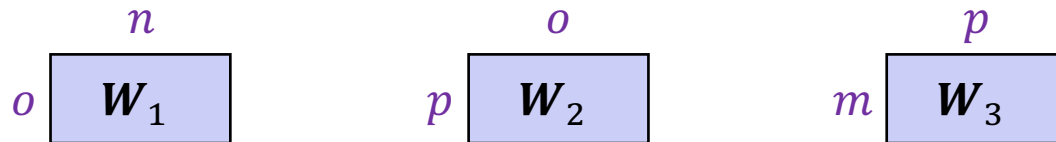
- loss function: $L(x, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3, x, y) = \frac{1}{2} \|\mathbf{W}_3\mathbf{W}_2\mathbf{W}_1x - y\|_2^2$

- now we add two more layers: $\hat{y} = f_3 \left(f_2(f_1(x)) \right) = W_3 W_2 W_1 x$
- loss function: $L(x, W_1, W_2, W_3, x, y) = \frac{1}{2} \|W_3 W_2 W_1 x - y\|_2^2$
- gradient of last layer:
$$\frac{\partial L}{\partial W_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial W_3} = (W_3 W_2 W_1 x - y)(W_2 W_1 x)^T$$


$m \times p$

$m \times 1$

$1 \times p$

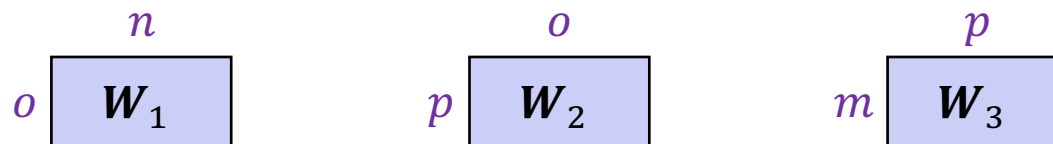


- now we add two more layers: $\hat{y} = f_3 \left(f_2(f_1(x)) \right) = W_3 W_2 W_1 x$
- loss function: $L(x, W_1, W_2, W_3, x, y) = \frac{1}{2} \|W_3 W_2 W_1 x - y\|_2^2$

- deeper gradients: $\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial W_2} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial f_2} \frac{\partial f_2}{\partial W_2} = W_3^T (W_3 W_2 W_1 x - y) (W_1 x)^T$


order to be consistent
with matrix multiplication

$p \times o$ $p \times m$ $m \times 1$ $1 \times o$



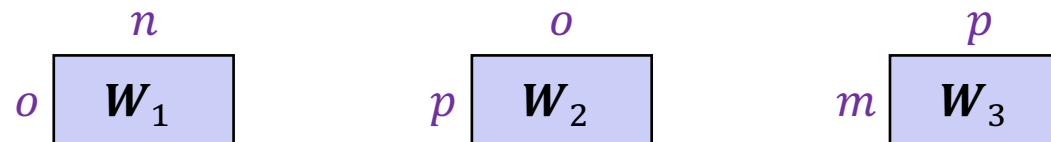
- now we add two more layers: $\hat{y} = f_3(f_2(f_1(x))) = W_3 W_2 W_1 x$
- loss function: $L(x, W_1, W_2, W_3, x, y) = \frac{1}{2} \|W_3 W_2 W_1 x - y\|_2^2$

- deeper gradients:

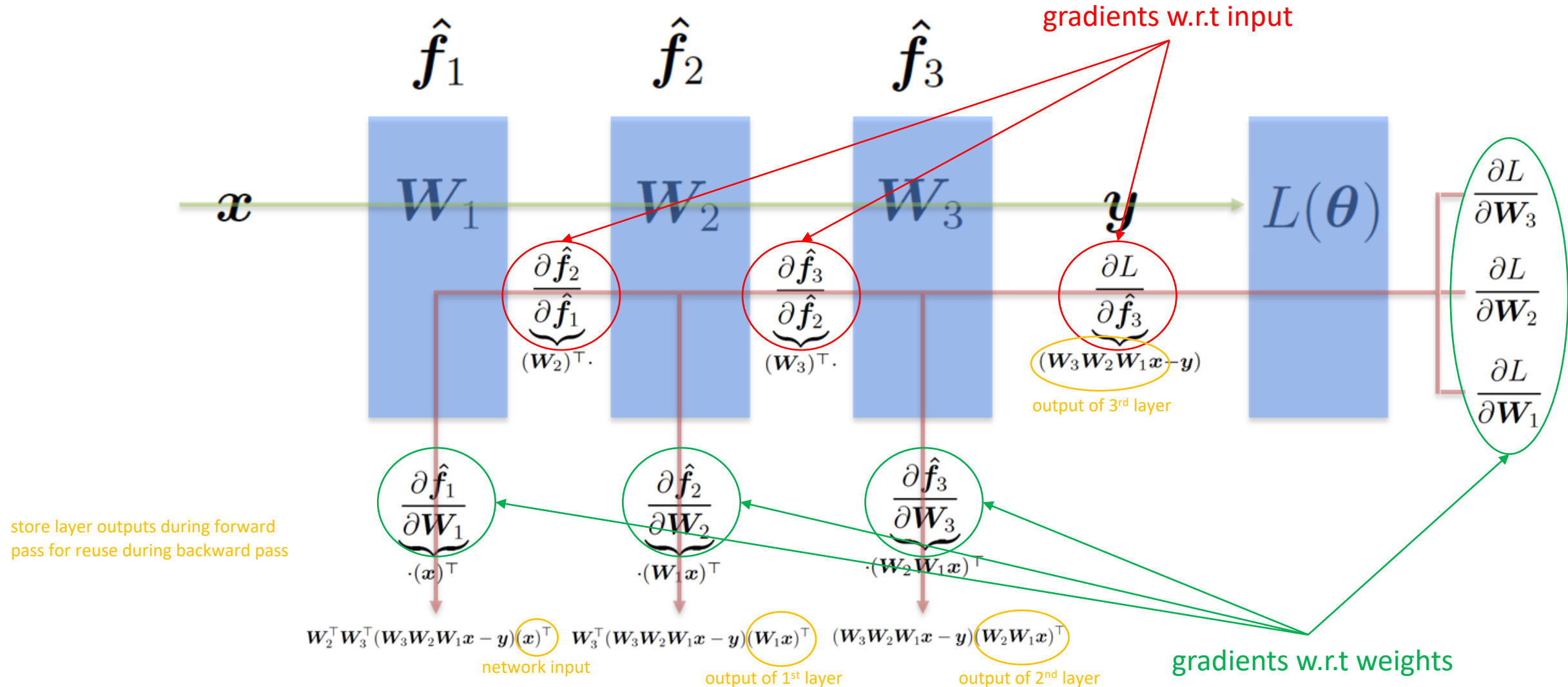
$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial W_1} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial f_2} \frac{\partial f_2}{\partial W_1} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial f_2} \frac{\partial f_2}{\partial f_1} \frac{\partial f_1}{\partial W_1} = W_2^T W_3^T (W_3 W_2 W_1 x - y) x^T$$

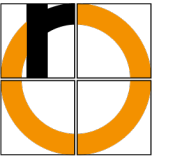
$o \times n$
 $o \times p$
 $p \times m$
 $m \times 1$
 $1 \times n$

order to be consistent with matrix multiplication



Linear Network in Matrix Notation





Universal Approximation Theorem

Universal Approximation Theorem

A **single hidden layer** is already sufficient to be a **universal function approximator**

- Proof by [Hor89, Hor91, Cyb89]; for a more recent paper see [Kra21]
- Let $\varphi(\cdot)$ be a non-constant, bounded and monotonically increasing function
- For any $\varepsilon > 0$ and any continuous function f defined on a compact subset of \mathbb{R}^m there exists an integer N , real constants v_i, b_i and real vectors $\mathbf{w}_i \in \mathbb{R}^m$ such that

#neurons in hidden layer N

activation function φ

input vector \mathbf{x}

Bias b_i

2nd layer of weights \mathbf{w}_i

first layer of weights (input – hidden)

$$F(\mathbf{x}) = \sum_{i=1}^N v_i \varphi(\mathbf{w}_i^T \mathbf{x} + b_i)$$

with $|F(\mathbf{x}) - f(\mathbf{x})| < \varepsilon$

only a single hidden layer is required!

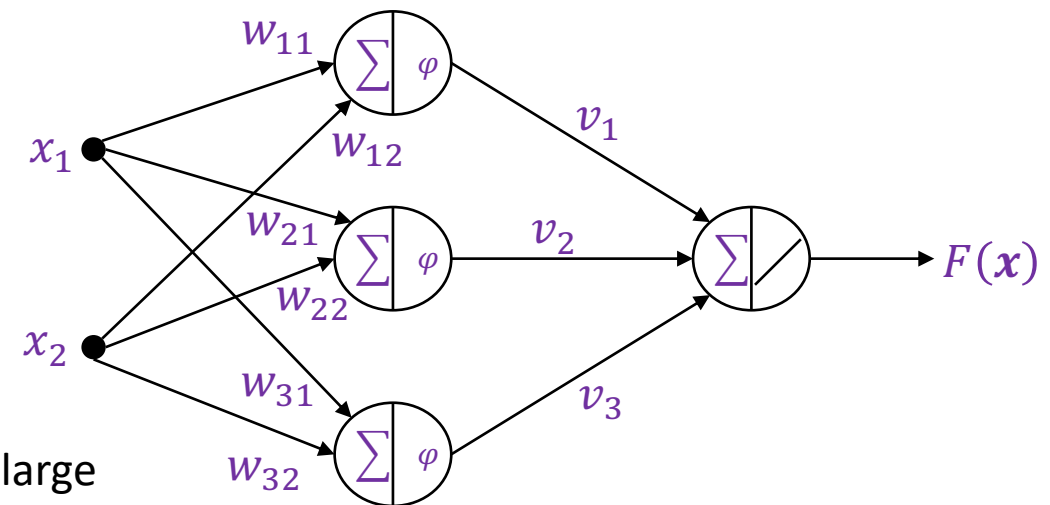
Observations:

- f can be approximated with arbitrary accuracy – just increase N
- it does not say anything about how to choose N , so ε may be quite large
- no nonlinear activation is required in the output neuron

Example:

$$\mathbf{x} = (x_1, x_2)$$

$$\mathbf{w}_1 = (w_{11}, w_{12}), \mathbf{w}_2 = (w_{21}, w_{22}), \mathbf{w}_3 = (w_{31}, w_{32})$$



A **single hidden layer** is also sufficient to train arbitrary (including non-contiguous) class boundaries [Mak89]

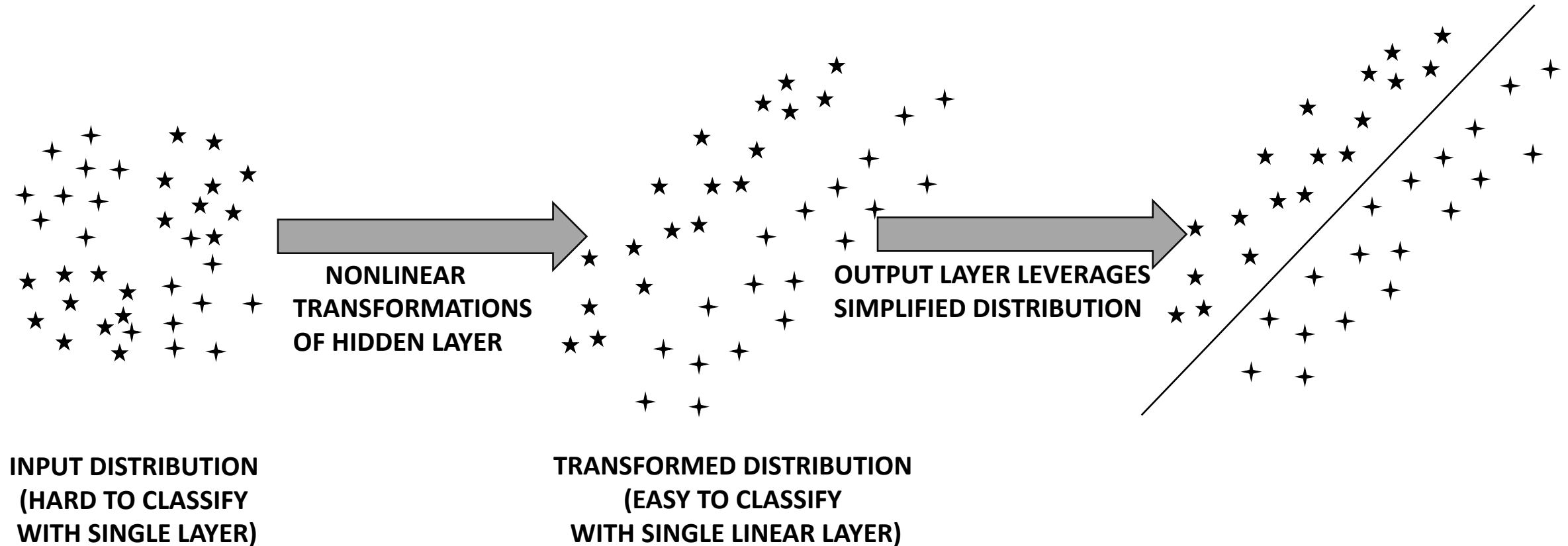


Figure from (C. Aggarwal. Neural Networks and Deep Learning, Springer 2018) Used with Permission. All rights reserved.

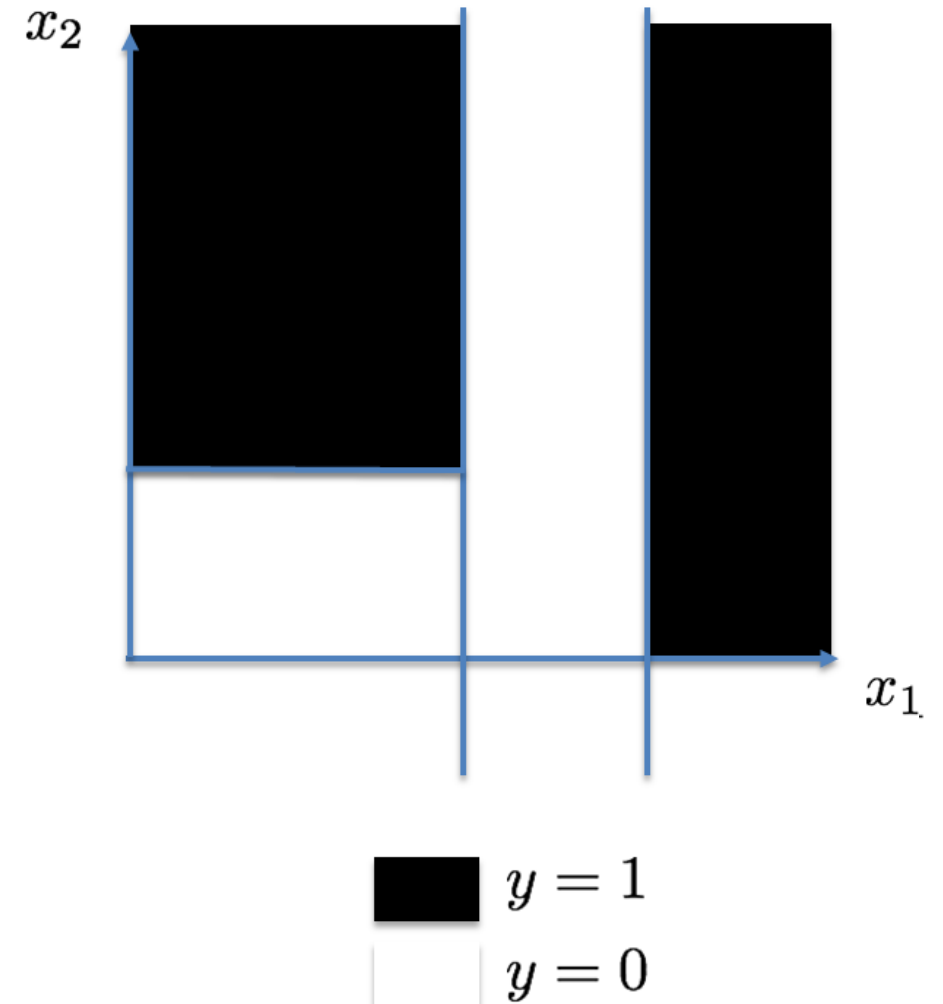
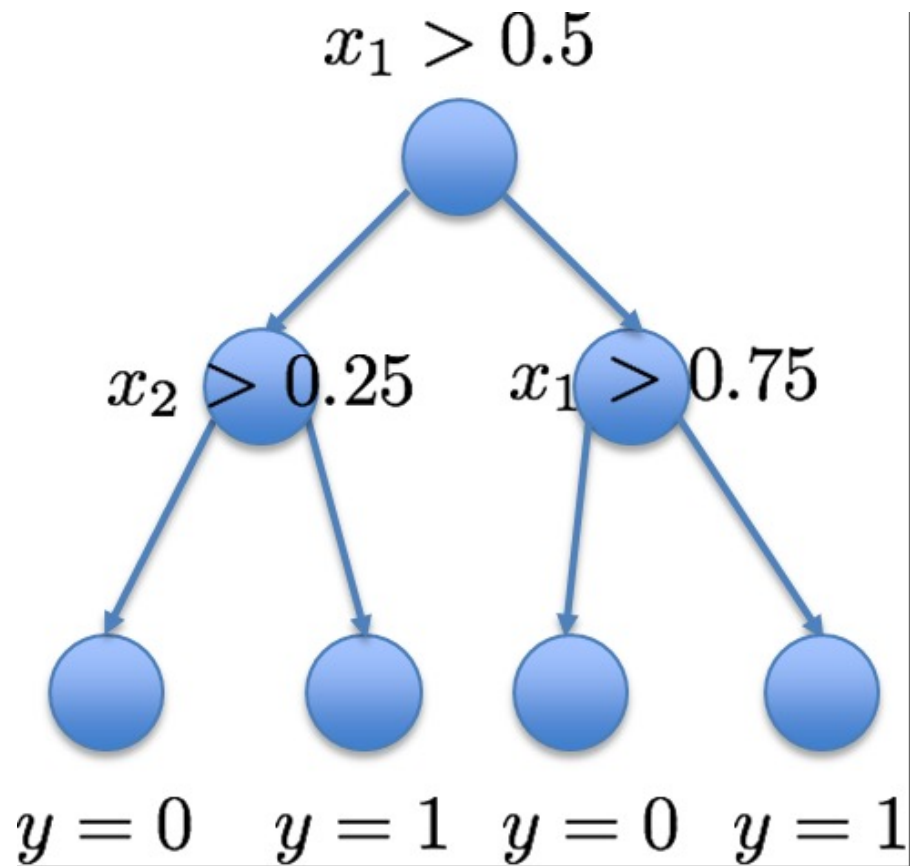
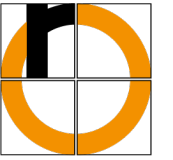
A MLP can

- compute any Boolean function (AND, OR, XOR, ...)
- approximate any non-linear function
- define arbitrary class boundaries
- do supervised learning using a sample set (Error Backpropagation)
- 2 layers of weights (= 1 hidden layer of neurons) are sufficient for
 - any compact continuous function (= on closed and bounded intervals)
 - arbitrary (including non-contiguous) class boundaries [Mak89]
 - a finite number of neurons in the hidden layer suffices; however, a lot of neurons may be required in that layer
- using 3 layers of weights, the function intervals can be unbounded

So: Why would we need Deep Learning then?

A MLP is not Turing complete!

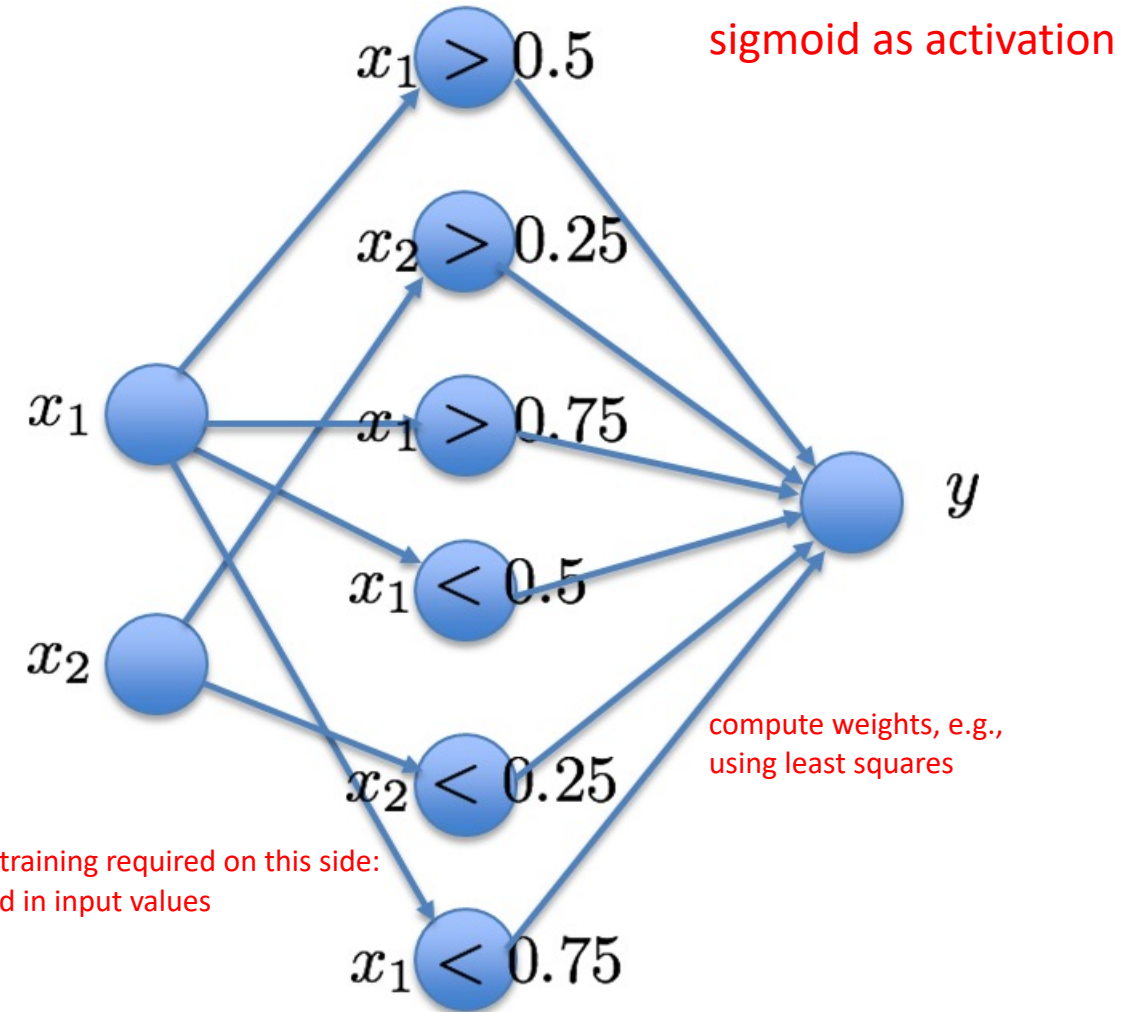
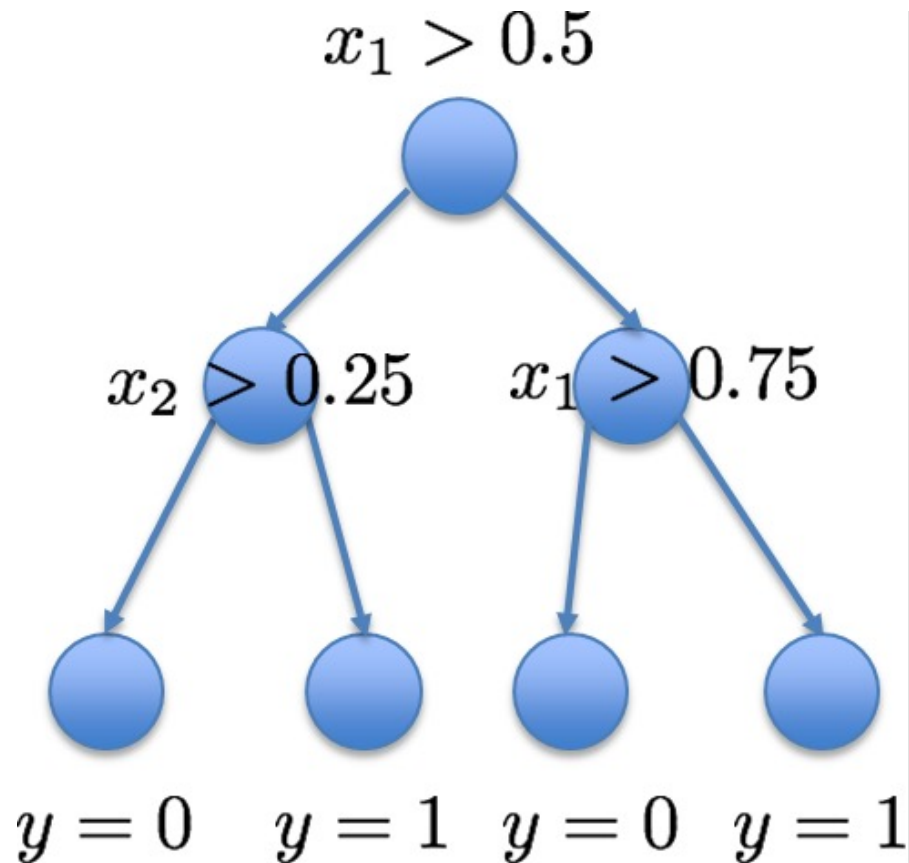
Example: Decision Tree



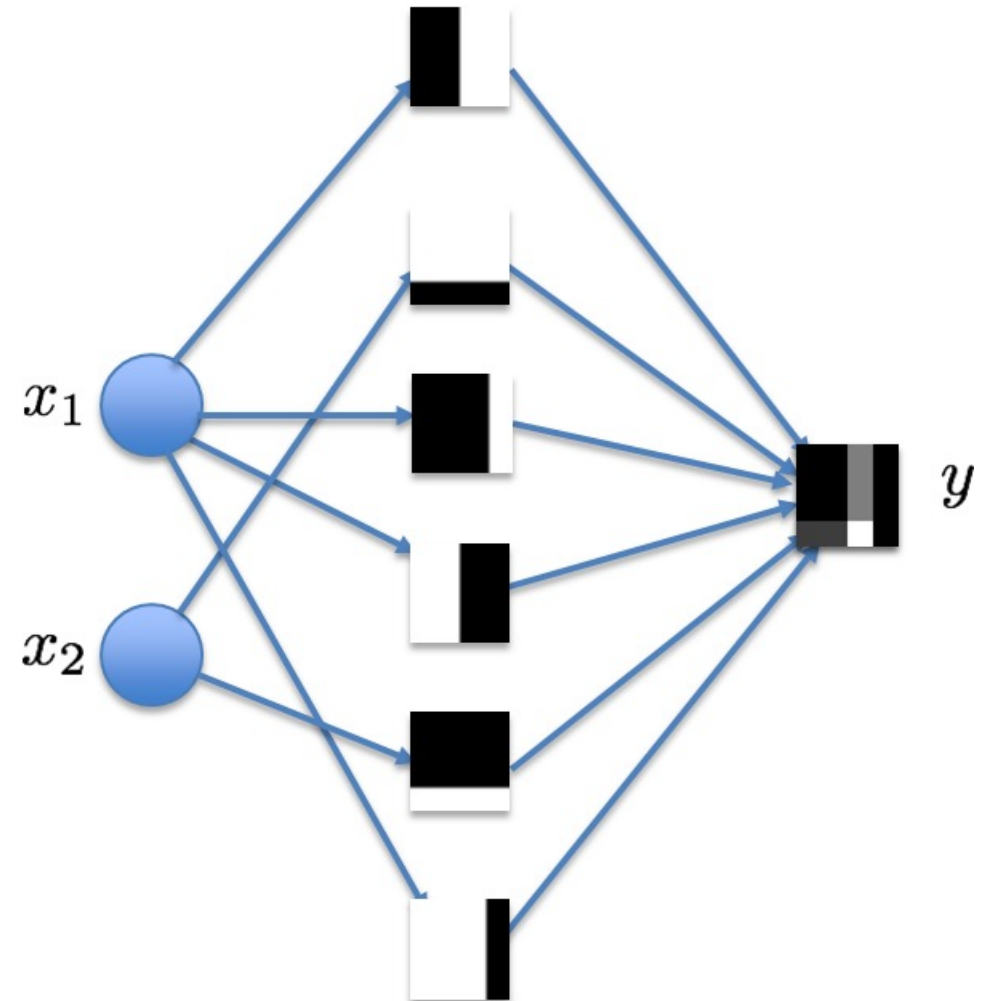
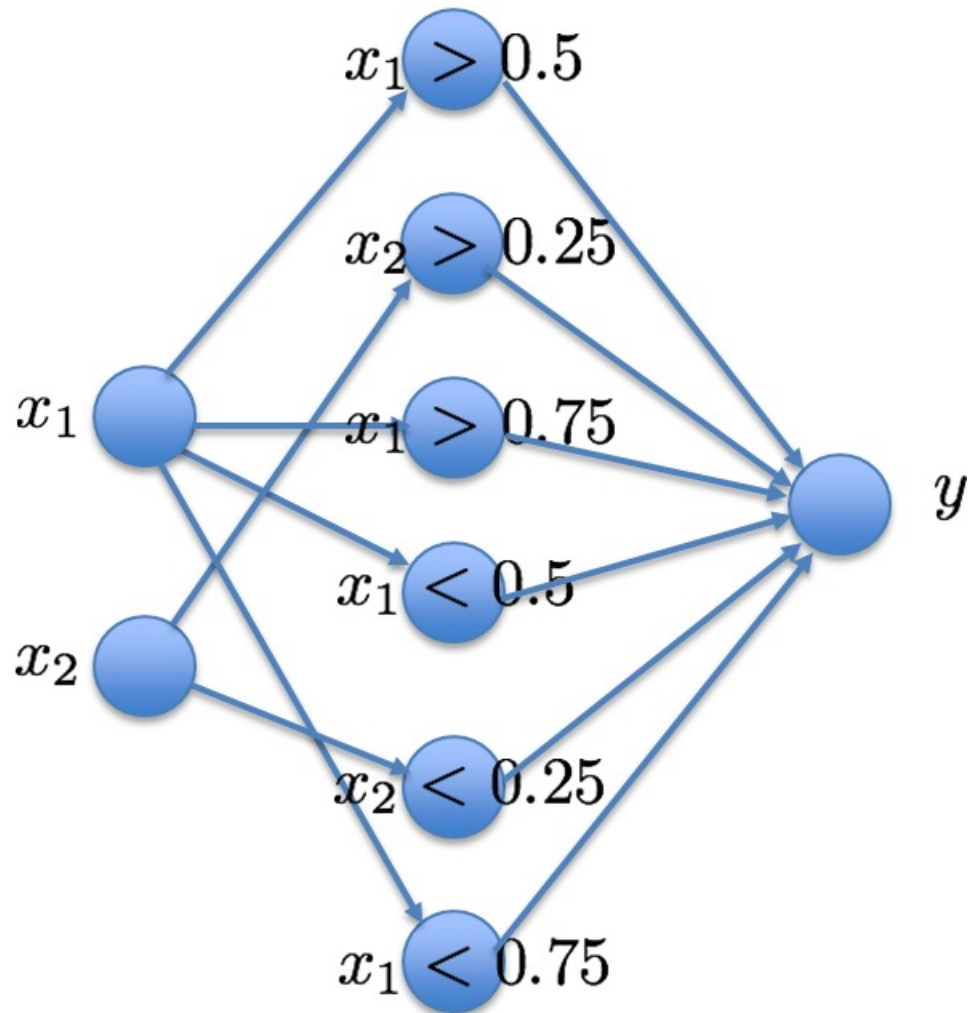
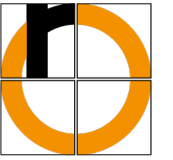
Example: Decision Tree



A neural network is an universal approximator \rightarrow
We can transform the decision tree into a network:



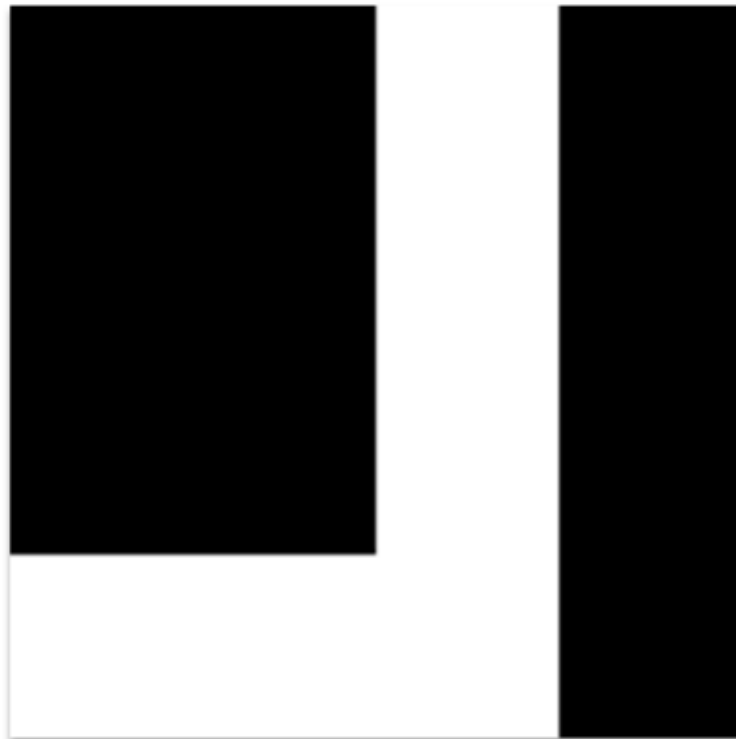
Example: Decision Tree



Example: Decision Tree



output of decision tree



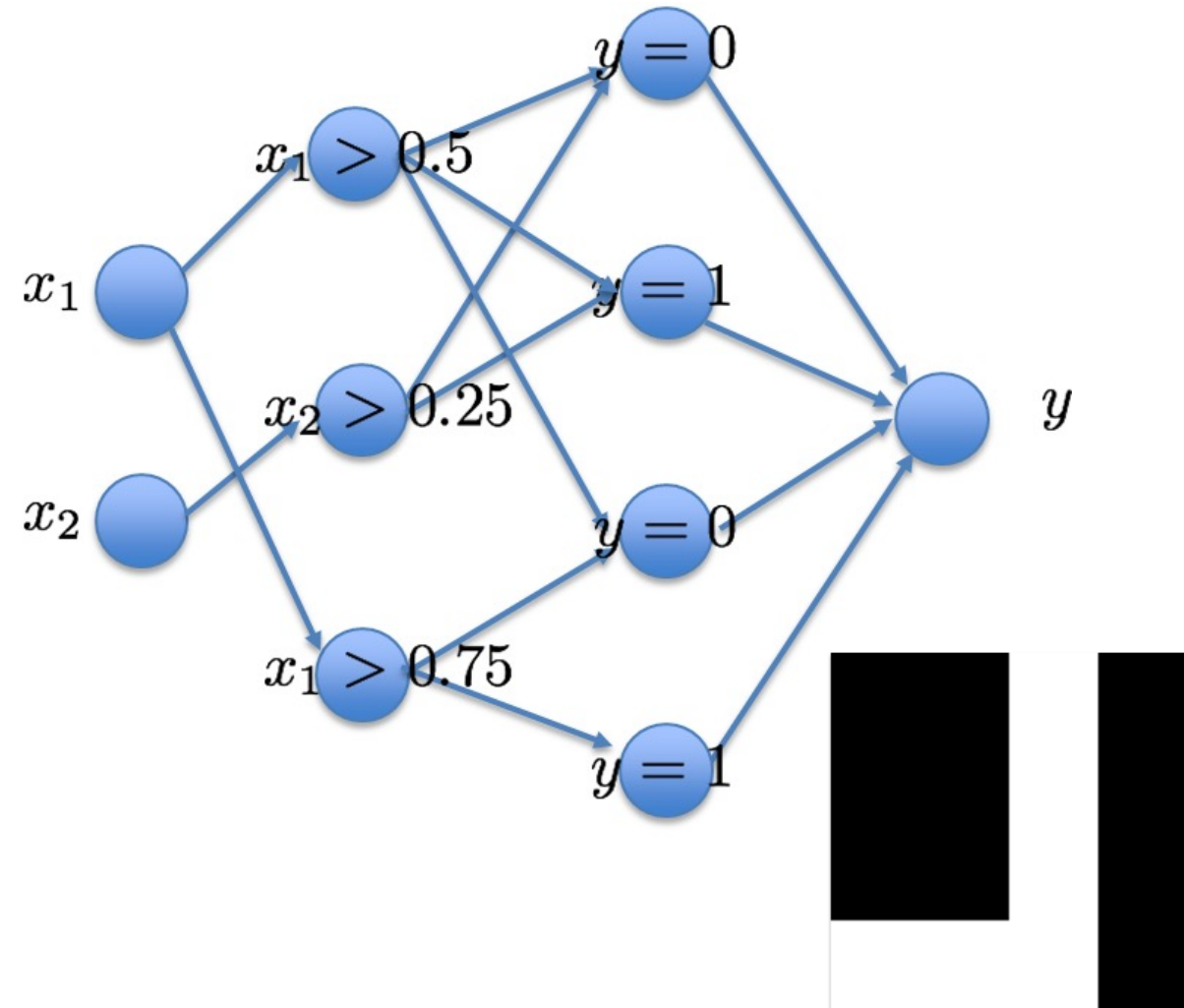
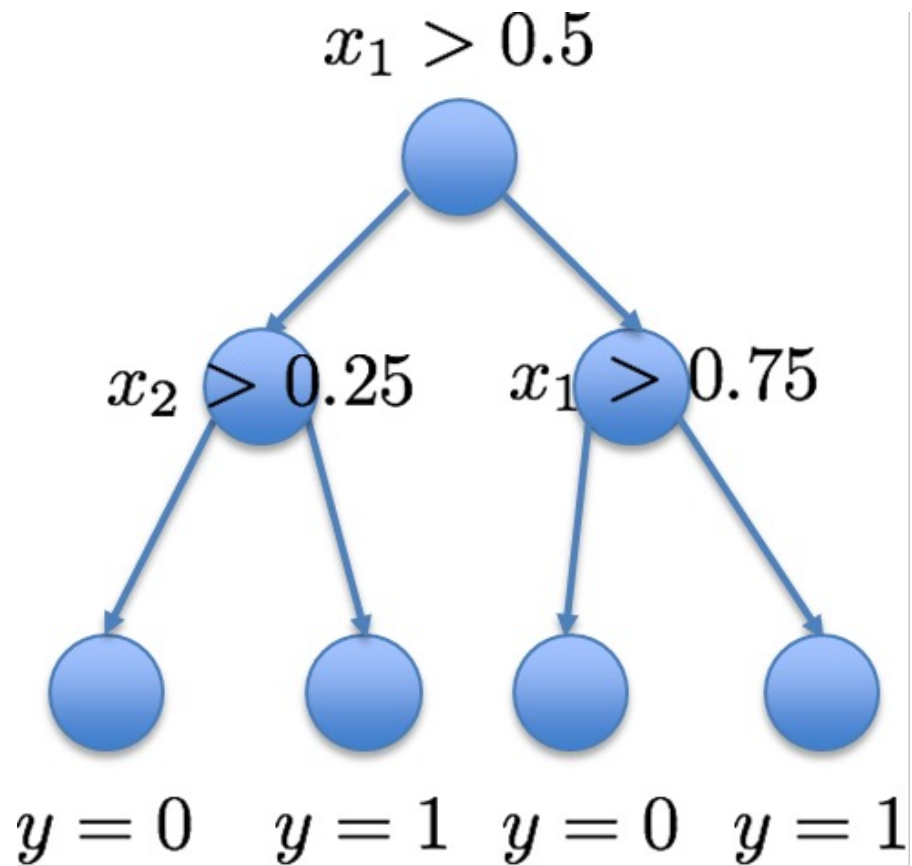
approximation of neural network with six neurons:



this is a quite hard problem – we would have to increase the number of neurons in the layer considerably to obtain a low-error approximation

Example: Decision Tree

add another layer



- newer proofs regarding deep networks with ReLU activation [Lu17, Han17]:
 - deep networks require less neurons than flat ones (for the same approximation quality)
 - they showed that
 - any continuous function with d -dimensional input and k -dimensional output can be approximated with arbitrary accuracy
 - using $d + k$ neurons in all hidden layers.
 - less than $d + k$ neurons are never sufficient, no matter how deep the network gets.
- for practical purposes
 - more neurons and deeper networks may be desirable, as the training may converge faster (or at all).
 - strictly speaking we do not have continuous functions anyway:
 - the number of training samples are always finite,
 - we have only discrete samples of the function at certain positions

- MLPs are a universal function approximator
- gradient descent is the default training algorithm for deep learning
- backpropagation allows for computing gradients efficiently

- [Han17] B. Hanin and M. Sellke. Approximating Continuous Functions by ReLU Nets of Minimal Width, 2017. [arXiv:abs/1710.11278](https://arxiv.org/abs/1710.11278).
- [Hor89] K. Hornik, M. Stinchcombe, and H. White. Multilayer Feedforward Networks are Universal Approximators. *Neural Networks*, 2(5):359 – 366, 1989.
- [Hor91] K. Hornik. Approximation Capabilities of Multilayer Feedforward Networks. *Neural Networks*, 4(2):251 – 257, 1991.
- [Cyb89] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4):303–314, 1989.
- [Kra21] A. Kratsios. The Universal Approximation Property. *Ann. Math. Artif. Intell.* (2021).
- [Lu17] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang. The Expressive Power of Neural Networks: A View from the Width. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan und R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, S. 6231–6239. Curran Associates, Inc., 2017.
- [Mai20] A. Maier, V. Christlein, K. Breininger, F. Denzinger, and F. Thamm: Deep Learning Slides Summer Semester 2020. Pattern Recognition Lab, Friedrich-Alexander-University Erlangen-Nürnberg. CC-BY 4.0.
- [Mak89] Makhoul, El-Jaroudi, and Schwartz. Formation of disconnected decision regions with a single hidden layer. In *International 1989 Joint Conference on Neural Networks*, S. 455–460 vol.1. 1989.
- [Pre07] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery. *Numerical Recipes 3rd Edition: The Art of Scientific Computing*. Cambridge University Press, New York, NY, USA, 3. Ed., 2007.