



**FAU**

FRIEDRICH-ALEXANDER-  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG  
SCHOOL OF ENGINEERING

# Feed Forward Neural Networks

**A. Maier, V. Christlein, K. Breininger, Z. Yang, L. Rist, M. Nau, S. Jaganathan, C. Liu, N. Maul, L. Folle,  
K. Packhäuser, M. Zinnen**

Pattern Recognition Lab, Friedrich-Alexander-Universität Erlangen-Nürnberg

April 24, 2023



## Outline

### Model

Perceptrons as Universal Function Approximators

From Activations to Classifications: Softmax Function

Optimization

Layer Abstraction



**FAU**

FRIEDRICH-ALEXANDER-  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG  
SCHOOL OF ENGINEERING

# Model



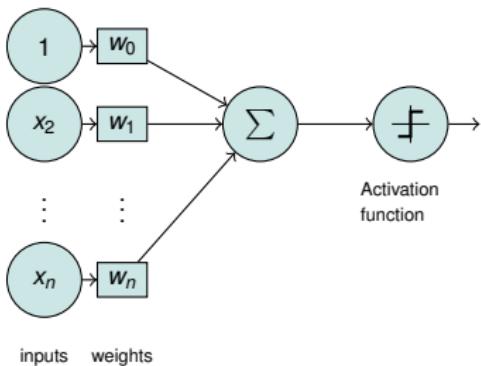
## Recap: Perceptrons

- Perceptron's decision rule:

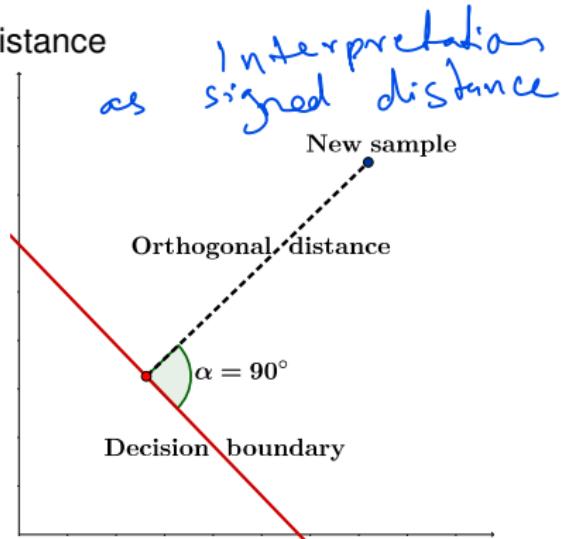
$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$

High dimensional vector to inner product

- Classification only depends on sign of distance



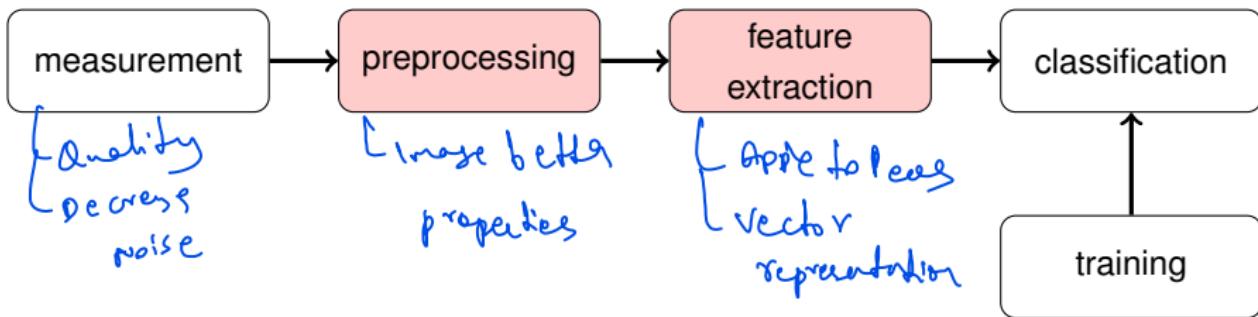
Perceptron



Linear decision boundary

# Classical PR pipeline

## Recap: Pattern Recognition Pipeline



MLP

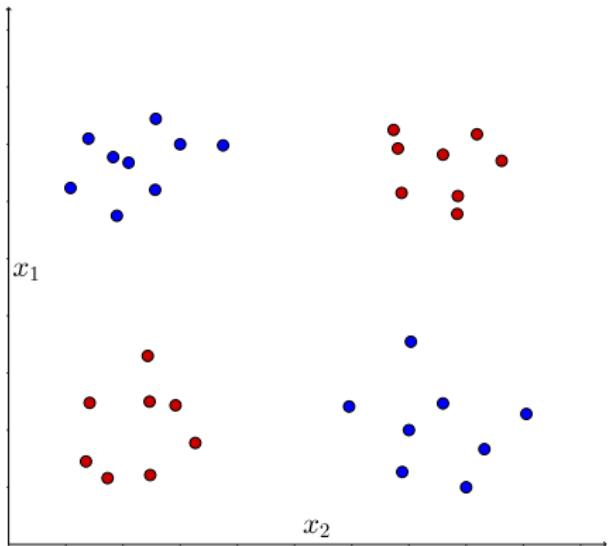
- (Multi-layer) perceptron (today's lecture) still uses predefined features
- “Hand-crafted” feature design is replaced by data driven feature learning in state-of-the-art architectures (upcoming lectures)
- Most concepts are important across architectures!

## XOR Problem

$$\begin{array}{r} \text{XOR} \\ - 0 \ 1 \ 1 \ 0 \\ A \oplus B \end{array}$$

single perceptron  $\rightarrow$  NOT possible

multiple  $\Rightarrow$  Yes



Samples from a XOR problem

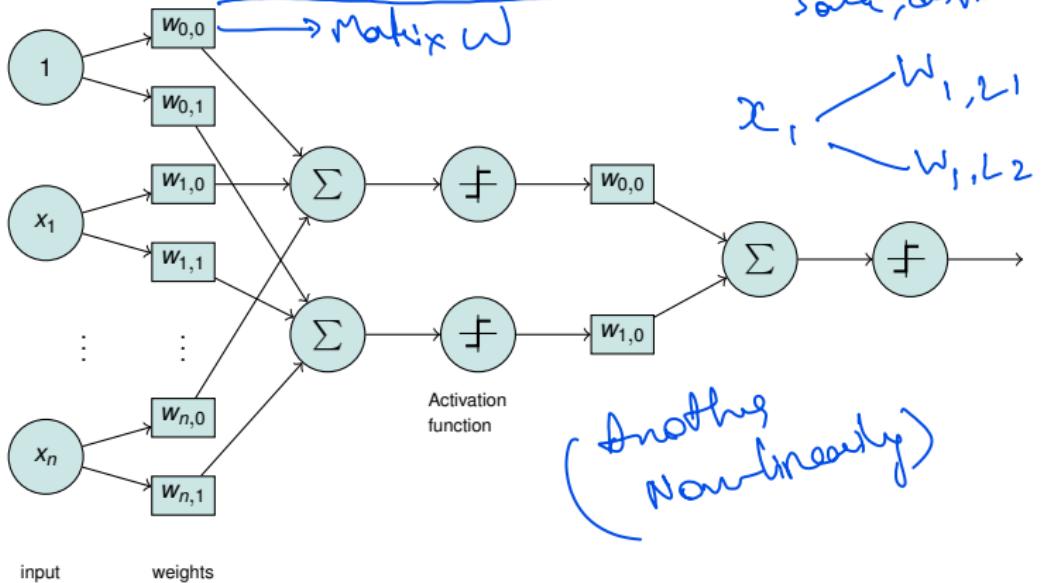
- XOR can't be solved with a line
- 1969: "Perceptrons" described limitations of neural networks
- AI funding was cut heavily
- This became known as "AI winter"

Idea: Use two lines  
or non-linear

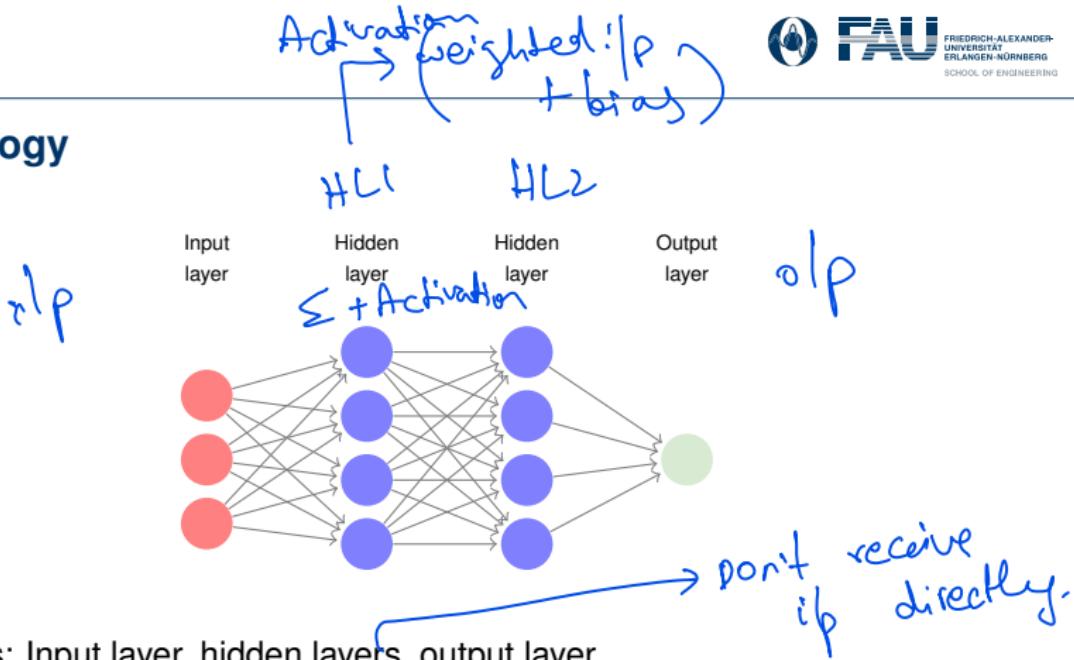
# (CS) (WSP)

## Multi-Layer Perceptron

- A perceptron resembles a single neuron
- Idea: Use multiple neurons!
- This enables non-linear decision boundaries



## Terminology



- Terms: Input layer, hidden layers, output layer
- A single hidden layer (of arbitrary width) can already be shown to be a **universal function approximator**
- Non-linear functions
  - are called activation functions in hidden layers
  - predict in the output layer and are used for the loss function



**FAU**

FRIEDRICH-ALEXANDER-  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG  
SCHOOL OF ENGINEERING

# Perceptrons as Universal Function Approximators



## Universal Approximation Theorem

- 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 exist an integer  $N$ , real constants  $v_i, b_i \in \mathbb{R}$  and real vectors  $w_i \in \mathbb{R}^m$  where  $i = 1, \dots, N$ , such that

$$F(\mathbf{x}) = \sum_{i=1}^N v_i \varphi(\mathbf{w}_i^T \mathbf{x} + b_i) \quad \text{with}$$
$$|F(\mathbf{x}) - f(\mathbf{x})| < \varepsilon$$

## Universal Approximation Theorem

- 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 exist an integer  $N$ , real constants  $v_i, b_i \in \mathbb{R}$  and real vectors  $w_i \in \mathbb{R}^m$  where  $i = 1, \dots, N$ , such that

Activation Function

$$F(\mathbf{x}) = \sum_{i=1}^N v_i \varphi(w_i^T \mathbf{x} + b_i) \quad \text{with} \quad |F(\mathbf{x}) - f(\mathbf{x})| < \varepsilon$$

$x_i$   
 $w_i$   
 $\Sigma + \varphi$   
 $F(x) = \text{linear combination of non-linearities}$

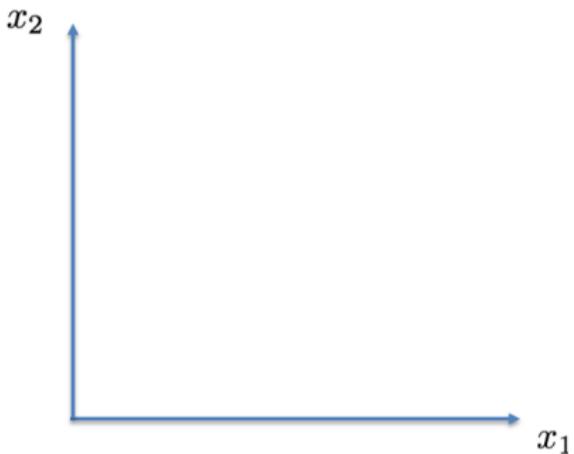
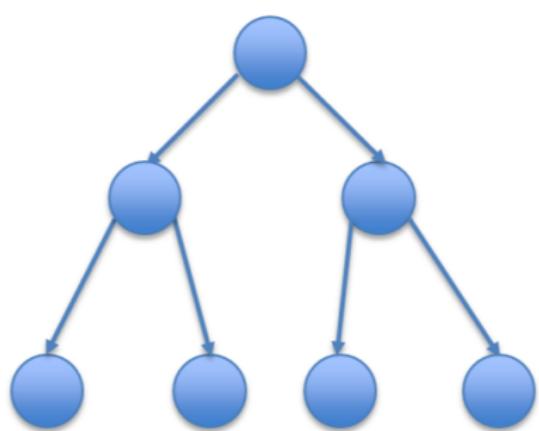
$\mathbb{R}^m \xrightarrow{w_i}$   
 $\mathbb{R}^n \xrightarrow{\varepsilon}$   
 Could be large.  
 or p

Continous  
 N → ∞  
 ε → 0

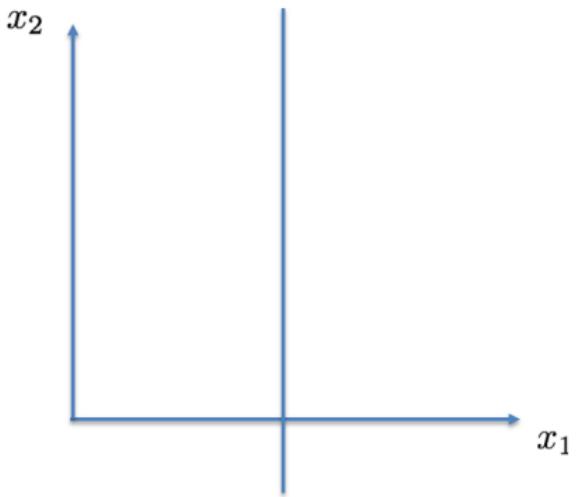
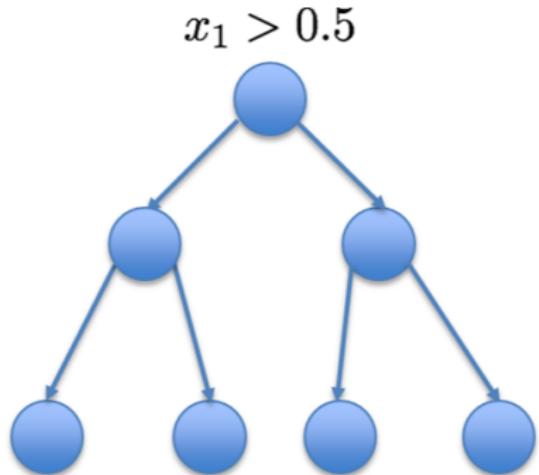
- We can approximate any function with just one hidden layer with a sensible activation function.

If only one HL enough? Why go Deep

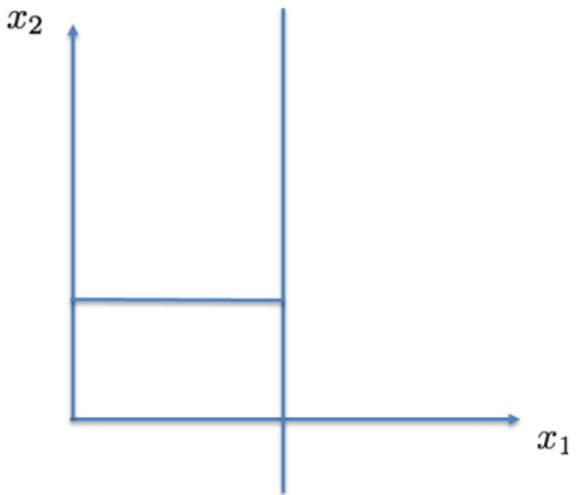
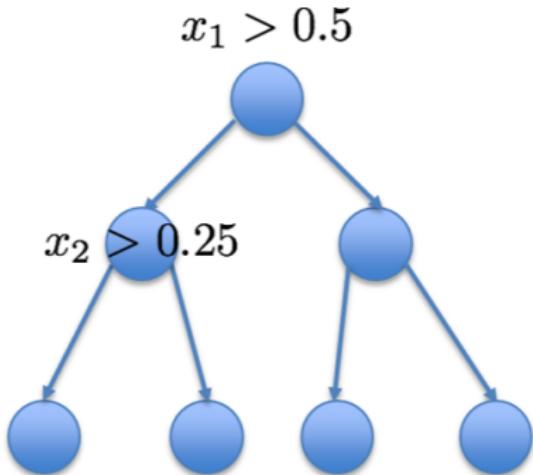
## Classification Trees: Why do we need a universal approximator?



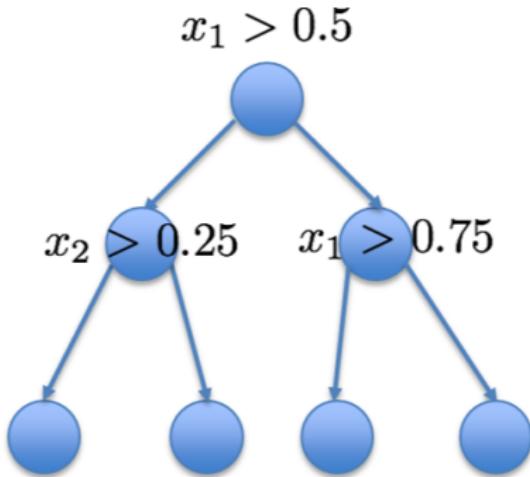
## Classification Trees: Why do we need a universal approximator?



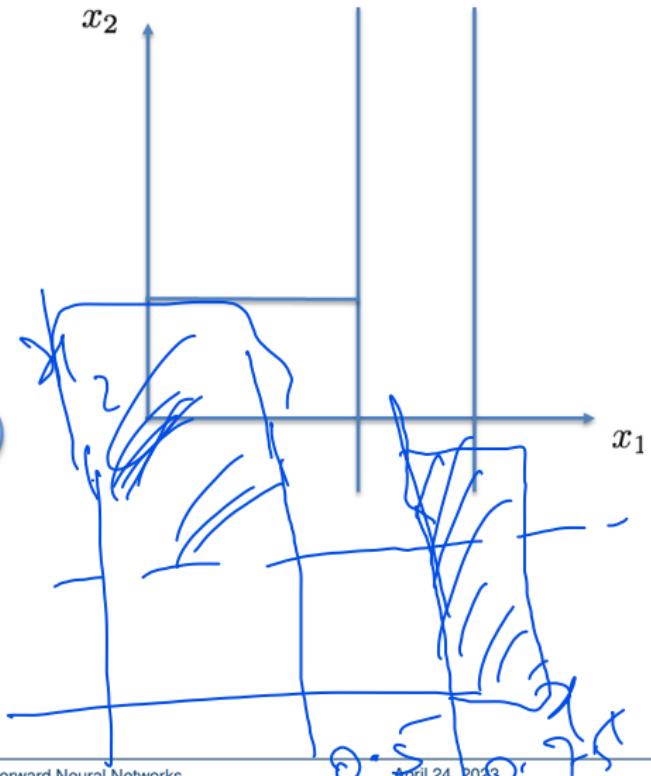
## Classification Trees: Why do we need a universal approximator?



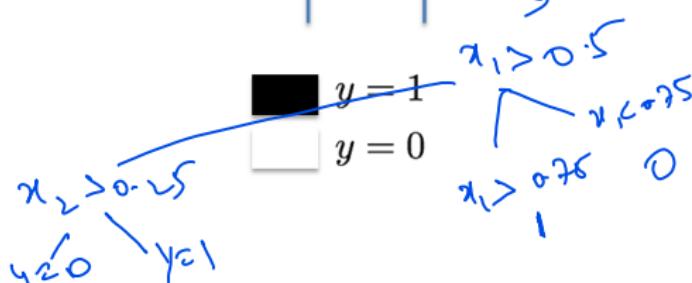
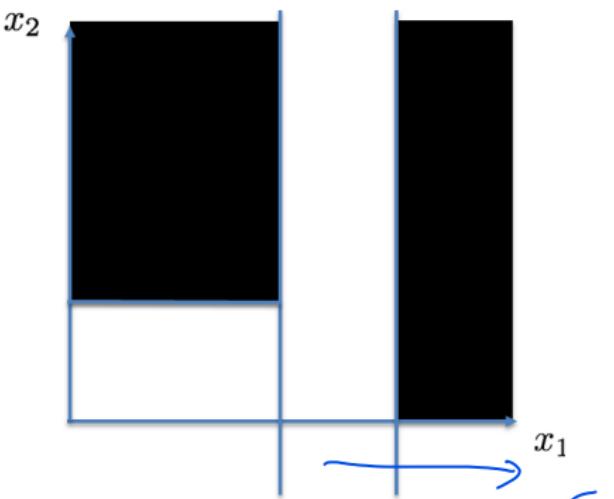
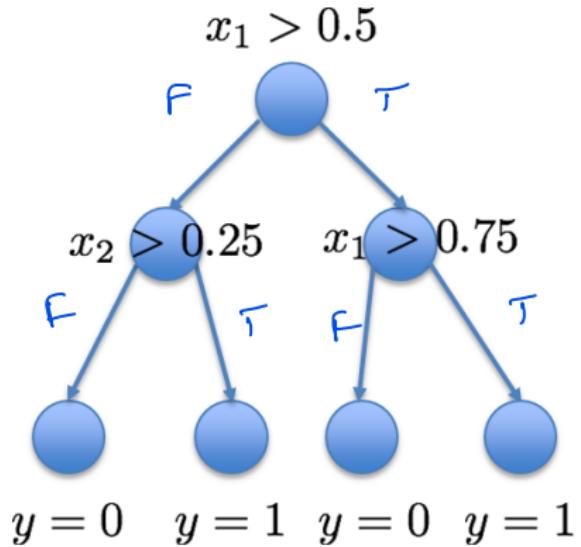
# Classification Trees: Why do we need a universal approximator?



$n_2 =$



## Classification Trees: Why do we need a universal approximator?

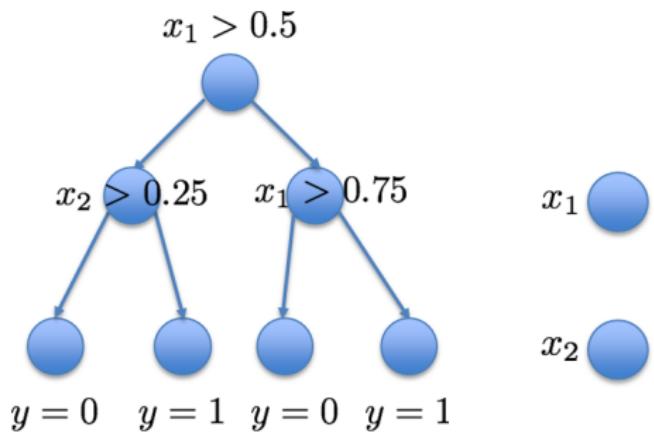


## Classification Trees: Why do we need a universal approximator?

We can transform this into a network!

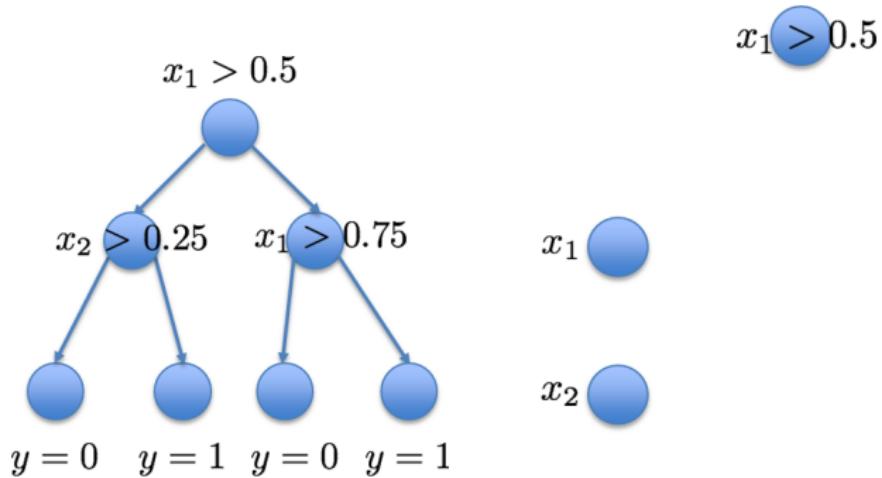
Classification  
Tree       $\Rightarrow$  Network

# Classification Trees: Why do we need a universal approximator?

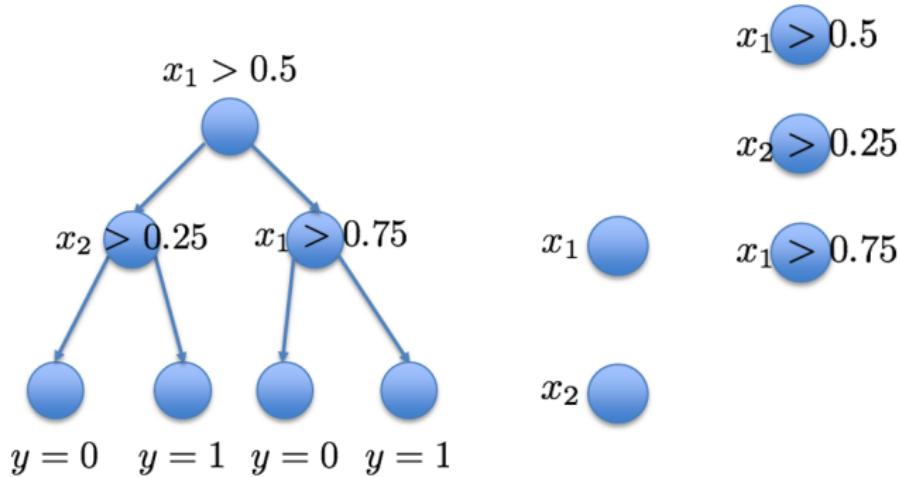


$x_1$    
 $x_2$  

# Classification Trees: Why do we need a universal approximator?

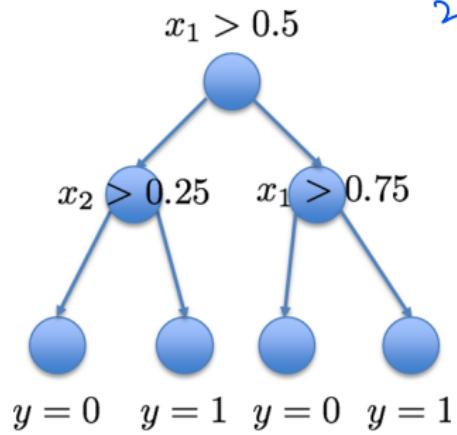


# Classification Trees: Why do we need a universal approximator?

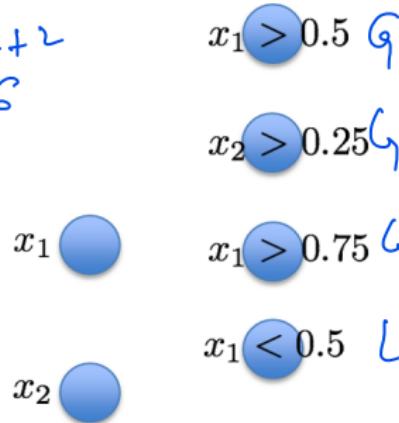


## Classification Trees: Why do we need a universal approximator?

3 nodes < Two possibilities



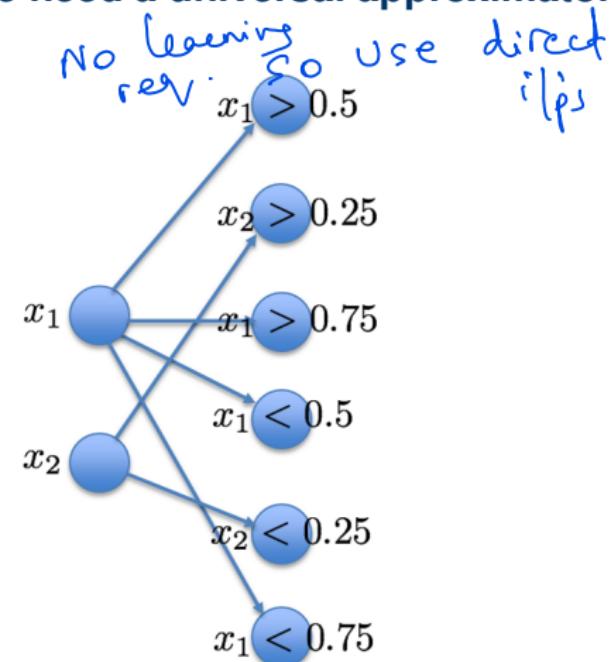
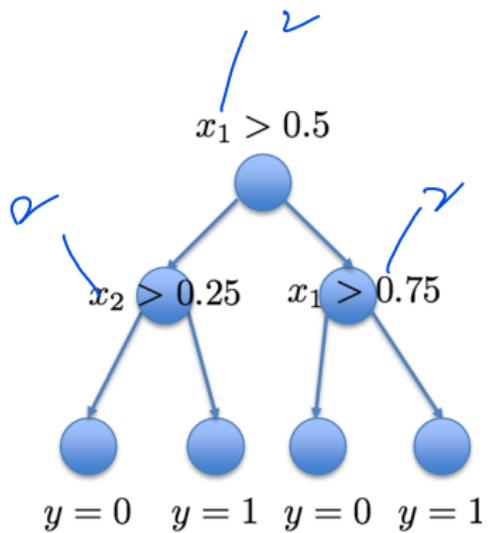
$$2 + 2 + 2 \approx 6$$



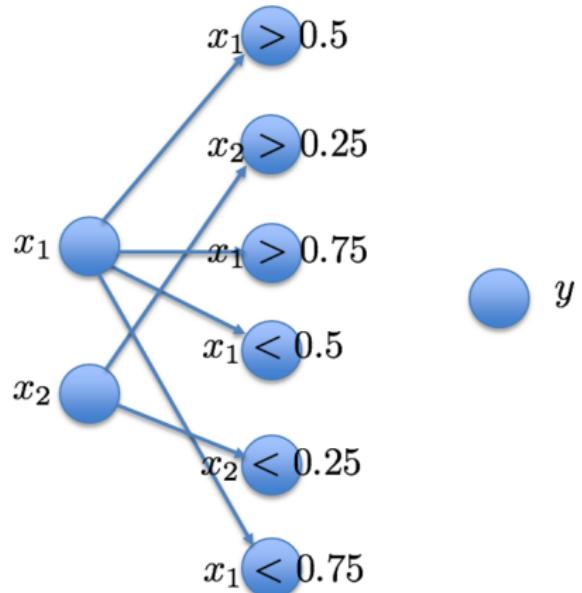
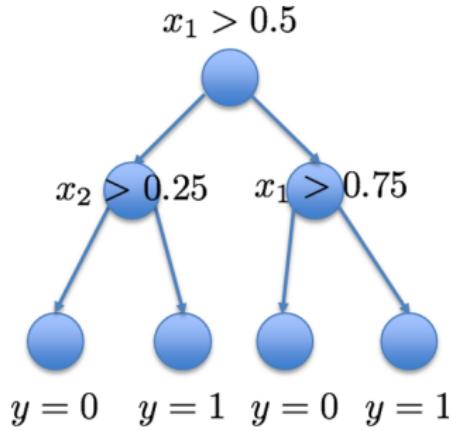
$x_1 > 0.5$   
 $x_2 > 0.25$   
 $x_1 > 0.75$   
 $x_1 < 0.5$   
 $x_2 < 0.25$   
 $x_1 < 0.75$

Sigmoid  
as  
Activation  
Use -  
Inverse  
as well

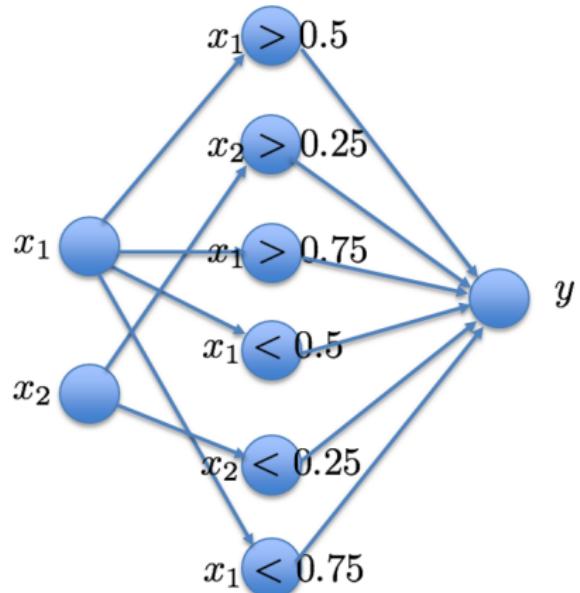
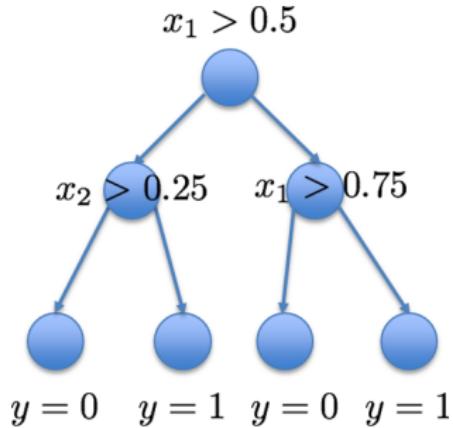
## Classification Trees: Why do we need a universal approximator?



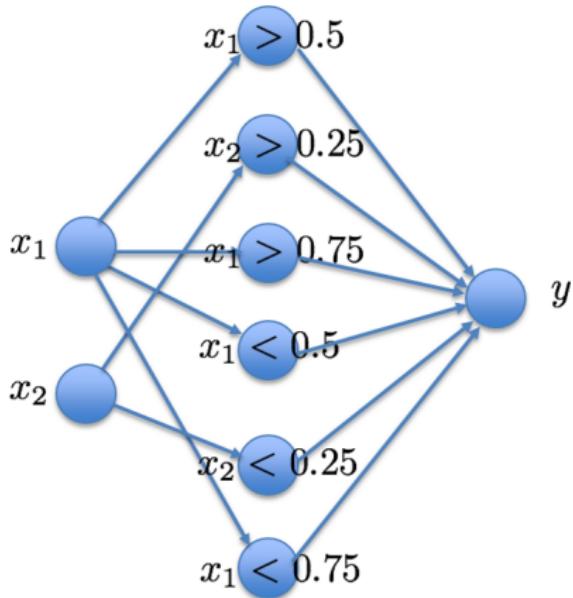
## Classification Trees: Why do we need a universal approximator?



## Classification Trees: Why do we need a universal approximator?

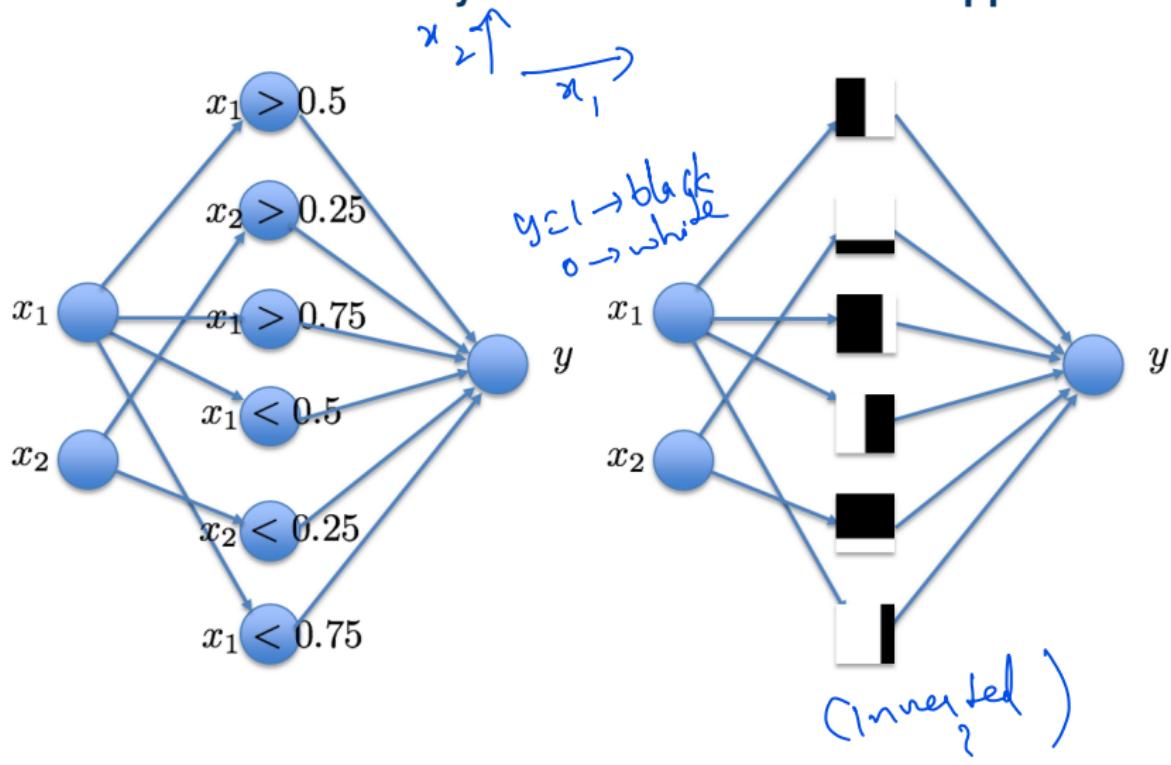


# Classification Trees: Why do we need a universal approximator?

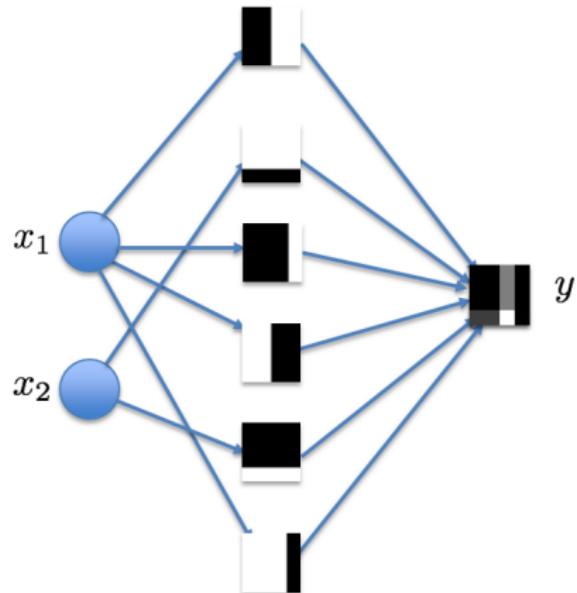
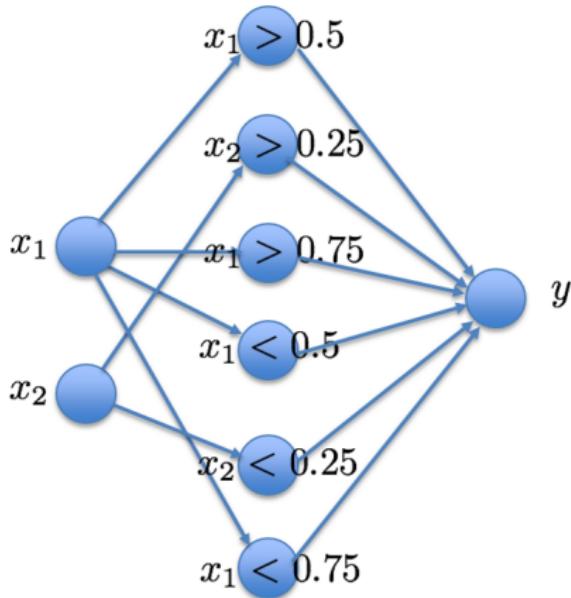


IF  $y = [0, 1]$

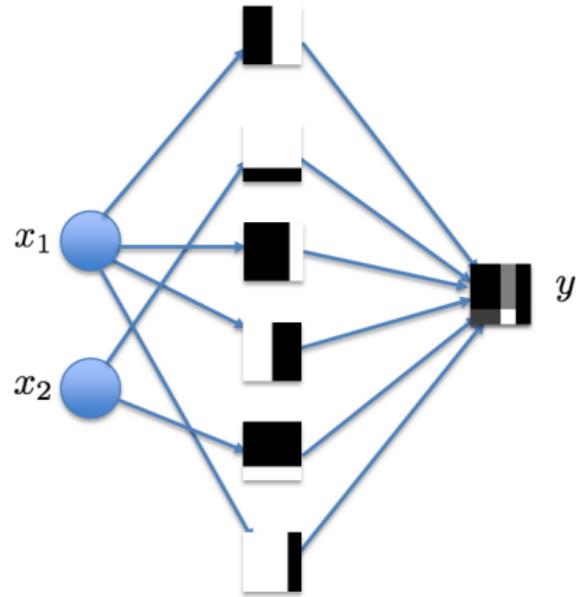
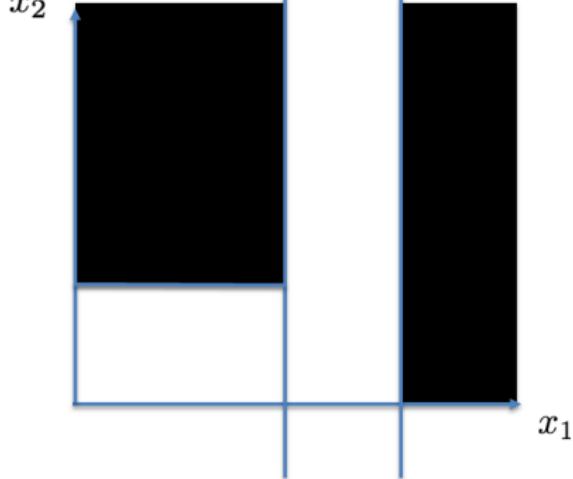
## Classification Trees: Why do we need a universal approximator?



# Classification Trees: Why do we need a universal approximator?

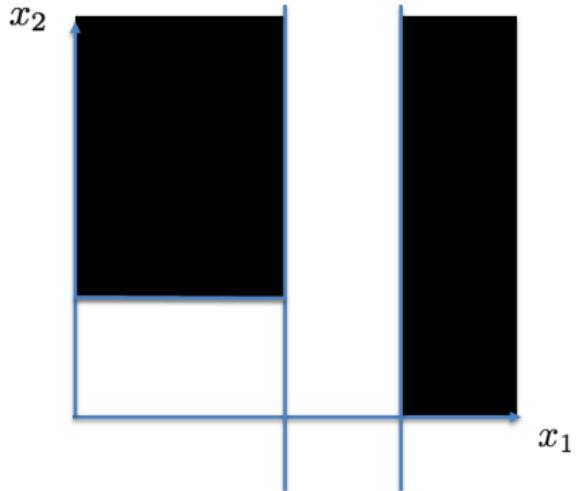


# Classification Trees: Why do we need a universal approximator?



# Classification Trees: Why do we need a universal approximator?

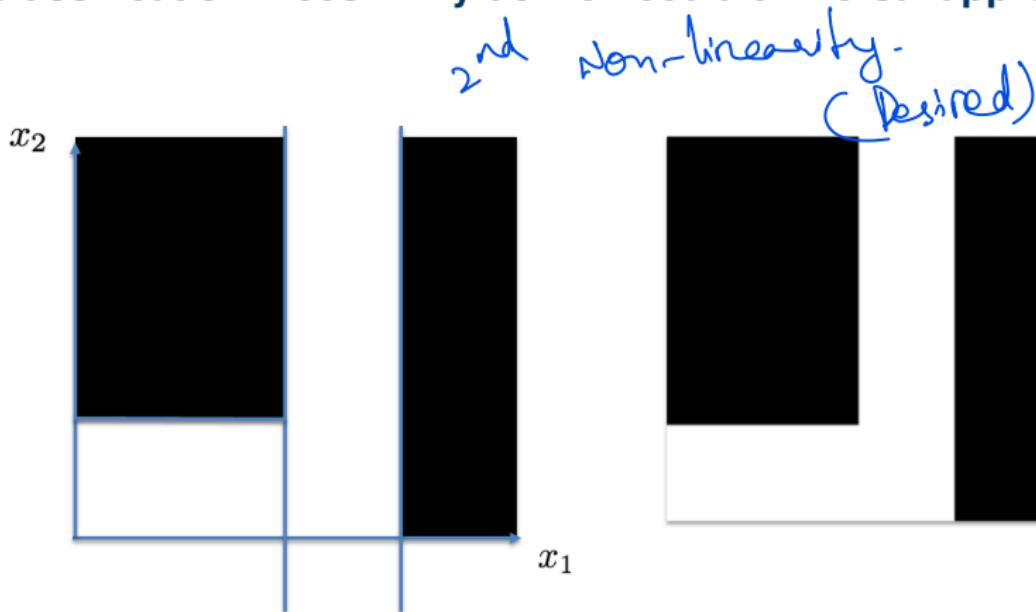
(Least Square  $\hookrightarrow$  Approximation)



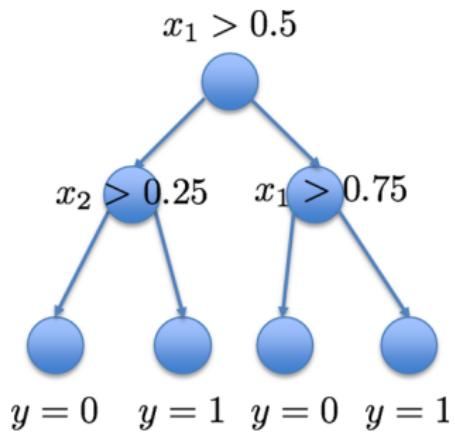
Doesn't matter  
 How many nodes  
 $n \rightarrow \infty$

$\epsilon > 0 \rightarrow$  not close to 0

## Classification Trees: Why do we need a universal approximator?



## Classification Trees: Why do we need a universal approximator?

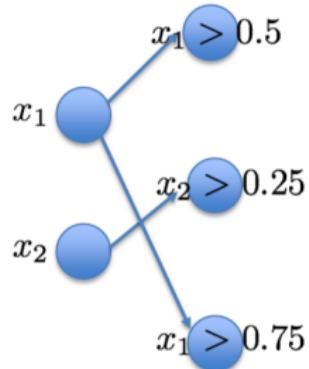
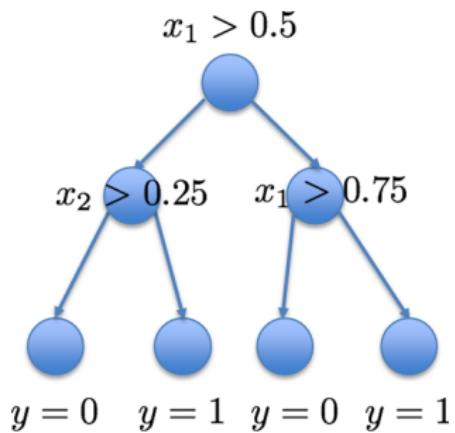


Algorithm to map  
Any Decision  
tree  
to N.N

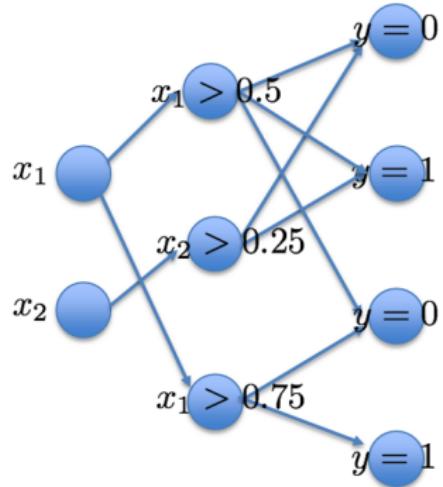
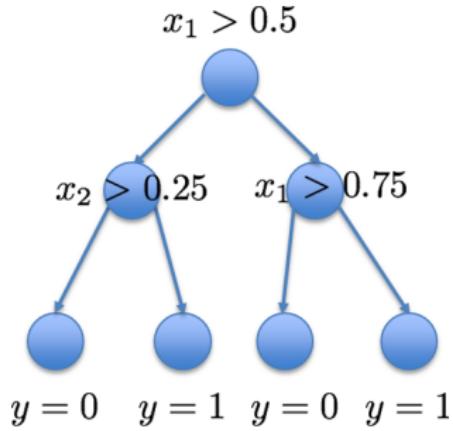


Take all  
inner  
nodes

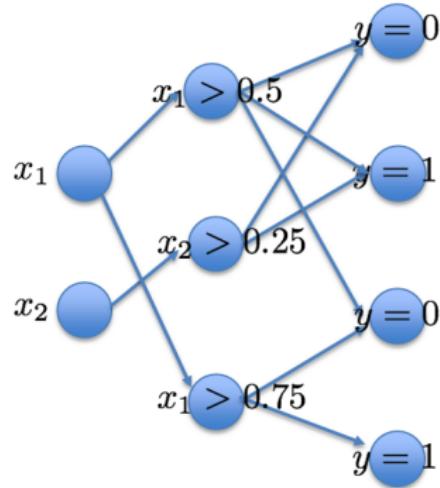
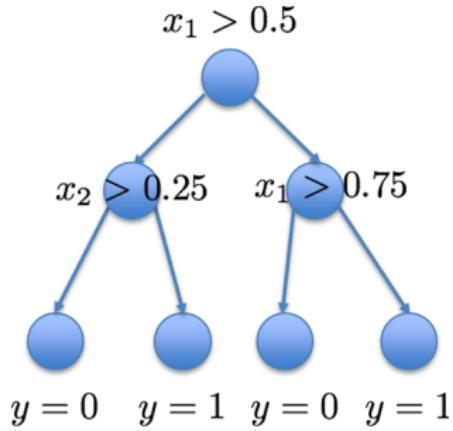
# Classification Trees: Why do we need a universal approximator?



## Classification Trees: Why do we need a universal approximator?

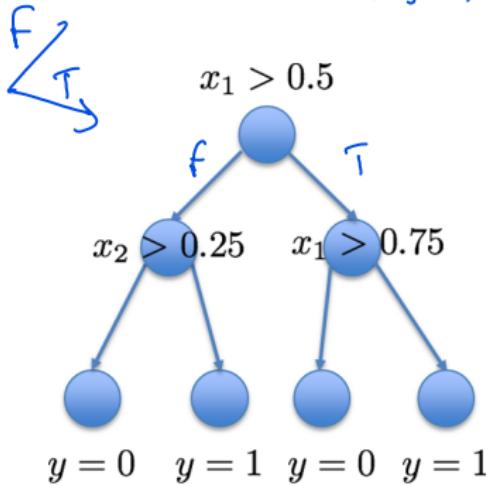


# Classification Trees: Why do we need a universal approximator?

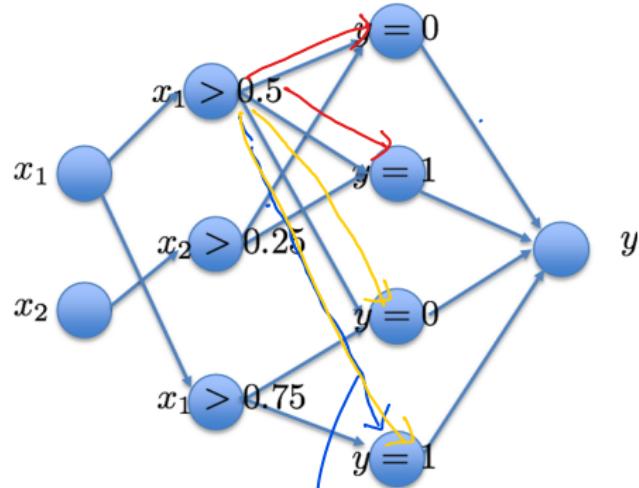


## Classification Trees: Why do we need a universal approximator?

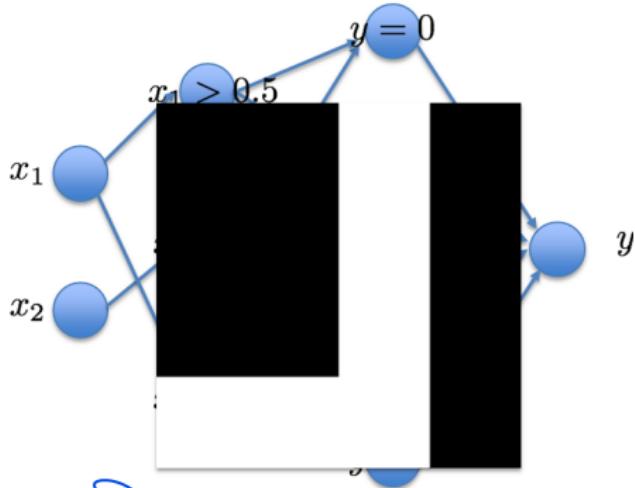
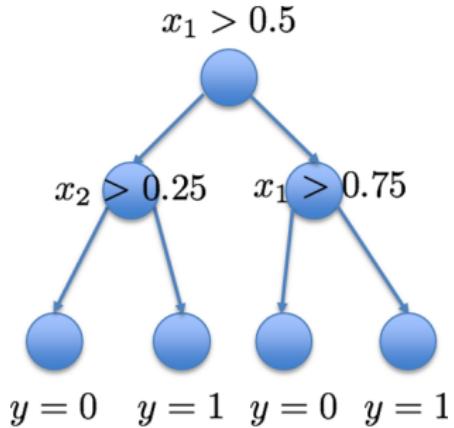
$x \cos \xi \rightarrow f_0, y_1$



(7 Nodes  
in sum  
problem solved)  
not drawn



# Classification Trees: Why do we need a universal approximator?



$(\text{more layers} \rightarrow \text{more steps} \rightarrow \text{Decomposition problem})$   
 $\text{One layer} \rightarrow \text{Solvable}$   
 But - Decomposition  $\rightarrow$  makes sense

DL = learning programs with more than 1 step

## Universal Approximation Theorem

↳ Exists

- 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 exist an integer  $N$ , real constants  $v_i, b_i \in \mathbb{R}$  and real vectors  $w_i \in \mathbb{R}^m$  where  $i = 1, \dots, N$ , such that we can define:

This is why  
we go to  
Deep learning  
DL

$$F(\mathbf{x}) = \sum_{i=1}^N v_i \varphi(\mathbf{w}_i^T \mathbf{x} + b_i) \quad \text{with} \\ |F(\mathbf{x}) - f(\mathbf{x})| < \varepsilon$$

- We can approximate *any function with just one hidden layer* with a sensible activation function.
- **We have no idea how: how many nodes, how to train, ...**

There definitely exists; How does it look like we don't know

**NEXT TIME  
ON DEEP LEARNING**

# Feed Forward Neural Networks - Part 2

A. Maier, V. Christlein, K. Breininger, Z. Yang, L. Rist, M. Nau, S. Jaganathan, C. Liu, N. Maul, L. Folle,  
K. Packhäuser, M. Zinnen

Pattern Recognition Lab, Friedrich-Alexander-Universität Erlangen-Nürnberg

April 24, 2023





**FAU**

FRIEDRICH-ALEXANDER-  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG  
SCHOOL OF ENGINEERING

# From Activations to Classifications: Softmax Function

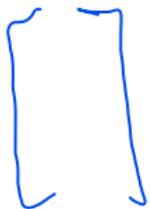


## Terminology

- So far: ground truth/estimated label is described by  $y/\hat{y} \in \{-1, 1\}$ .
- Instead, we can use a vector  $\mathbf{y} = (y_1, \dots, y_K)^T$  where  $K = \#\text{classes}$ .
- For exclusive classes,  $\mathbf{y}$  looks as follows:

vector  $\mathbf{y}$   
 scalar  $y_k = \begin{cases} 1 & \text{if } k \text{ is the index of the true class,} \\ 0 & \text{otherwise} \end{cases}$

- Called **one-hot encoding**: Only one element is  $\neq 0$ .
- Classifier output  $\hat{\mathbf{y}}$  can represent class probabilities.
- Better descriptor, especially for multi-class problems!



## Softmax activation function

- The softmax function rescales a vector  $\mathbf{x}$  using:

$$\hat{\mathbf{y}} = (\hat{y}_0, \hat{y}_1, \dots, \hat{y}_K)$$

$$\hat{y}_k = \frac{\exp(x_k)}{\sum_{j=1}^K \exp(x_j)}$$

Probabilistic o/p.

for each classes

- $\hat{\mathbf{y}}$  has two properties:

$$1. \sum_{k=1}^K \hat{y}_k = 1$$

$$2. \hat{y}_k \geq 0 \quad \forall \hat{y}_k \in \hat{\mathbf{y}}$$

- prediction vector

$$\geq 0 \quad \sum \hat{y}_k = 1$$

$$\sum \hat{y}_k = \sum e^{x_k} = 1$$

$$= \frac{e^{x_1}}{\sum e^{x_j}} > 0$$

- These are two of Kolmogorov's axioms for a probability distribution.

- This allows to treat the output as normalized probabilities.

- The softmax function is also known as the normalized exponential function.

S.F. = NEF

## Softmax activation function

- The softmax function rescales a vector  $\mathbf{x}$  using:

$$\hat{y}_k = \frac{\exp(x_k)}{\sum_{j=1}^K \exp(x_j)}$$

- Example: three-class problem



Label	$x_k$	$\exp(x_k)$	$\hat{y}_k$
Tiger			
Airplane			
Boat			

## Softmax activation function

- The softmax function rescales a vector  $\mathbf{x}$  using:

$$\hat{y}_k = \frac{\exp(x_k)}{\sum_{j=1}^K \exp(x_j)}$$

- Example: threefour-class problem



Label	$x_k$	$\exp(x_k)$	$\hat{y}_k$
Tiger			
Airplane			
Boat			
Heavy Metal			

## Softmax activation function

- The softmax function rescales a vector  $\mathbf{x}$  using:

$$\hat{y}_k = \frac{\exp(x_k)}{\sum_{j=1}^K \exp(x_j)}$$

- Example: threefour-class problem



Label	$x_k$	$\exp(x_k)$	$\hat{y}_k$
Tiger	-3.44		
Airplane	1.16		
Boat	-0.81		
Heavy Metal	3.91		

## Softmax activation function

- The softmax function rescales a vector  $\mathbf{x}$  using:

$$\hat{y}_k = \frac{\exp(x_k)}{\sum_{j=1}^K \exp(x_j)}$$

- Example: threefour-class problem



Label	$x_k$	$\exp(x_k)$	$\hat{y}_k$
Tiger	-3.44	0.03	
Airplane	1.16	3.19	
Boat	-0.81	0.44	
Heavy Metal	3.91	49.90	

## Softmax activation function

- The softmax function rescales a vector  $\mathbf{x}$  using:

$$\hat{y}_k = \frac{\exp(x_k)}{\sum_{j=1}^K \exp(x_j)}$$

- Example: threefour-class problem



T B A HM

Label	$x_k$	$\exp(x_k)$	$\hat{y}_k$
Tiger	-3.44	0.03	0.0006
Airplane	1.16	3.19	0.0596
Boat	-0.81	0.44	0.0083
Heavy Metal	3.91	49.90	0.9315

## Loss functions

→ How Good the prediction of a network is

- The cross entropy  $H$  of probability distributions  $\mathbf{p}$  and  $\mathbf{q}$  → Estimated.

$$y = [y_1 \dots y_K]$$

$$\hat{y} = (\hat{y}_1, \dots, \hat{y}_K)$$

$$H(\mathbf{p}, \mathbf{q}) = - \sum_{k=1}^K p_k \log(q_k)$$

$$L_{\text{CE}} = - \sum_{k=1}^K p_k \log(\hat{y}_k)$$

- Based on  $H$ , we formulate a loss function  $L$ :

$$(Vectors) \quad L(\mathbf{y}, \hat{\mathbf{y}}) = - \log(\hat{y}_k) \Big|_{y_k=1}$$

$$(Scalars) \quad L(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

$$L(\mathbf{y}, \hat{\mathbf{y}}) = - \log(\hat{y}_k) \Big|_{y_k=1}$$

$y_k = 1$ ;  $\hat{y}_k = 1$   
o otherwise

class  $k$  for which it is non-zero (true-class)

- We will talk more about this during the next session!

*Verbes*

$$L(y, x) = -\log(\hat{y}_k) \Big|_{y_k=1}$$

## "Softmax loss"

- Cross-entropy and the Softmax function typically appear together

$$L(y, x) = -\log \left( \frac{\exp(x_k)}{\sum_{j=1}^K \exp(x_j)} \right) \Big|_{y_k=1}$$

- One-hot encoding very convenient → represents a histogram
- Naturally handles multiple class problems

$$\hat{y}_k = \frac{\exp(x_k)}{\sum_{j=1}^K \exp(x_j)}$$



**FAU**

FRIEDRICH-ALEXANDER-  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG  
SCHOOL OF ENGINEERING

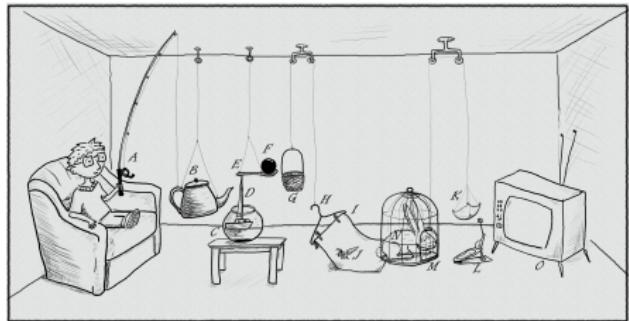
# Optimization



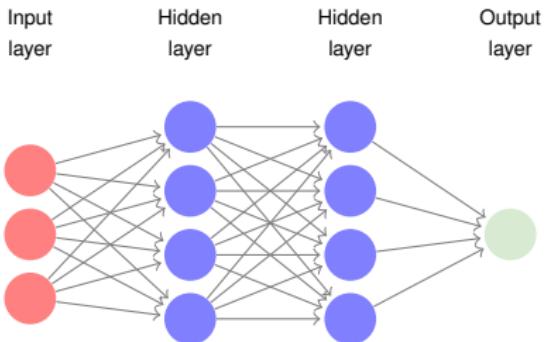
# How do learn those networks

## Credit Assignment Problem

- What do those two images have in common?



Source: <https://krypt3ia.files.wordpress.com/2011/11/rube.jpg>

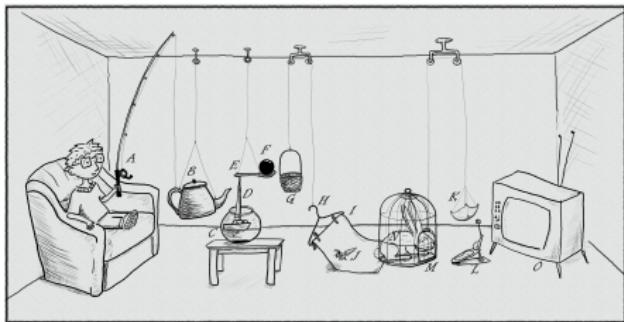


change anything in system  
→ effects chain of operations.

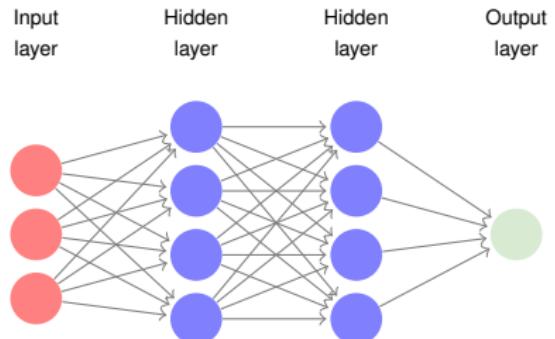
Be very careful

## Credit Assignment Problem

- What do those two images have in common?



Source: <https://krypt3ia.files.wordpress.com/2011/11/rube.jpg>



- If it doesn't work it's hard to know which parts to adjust.

## Formalization as Optimization Problem

**Goal:** Find optimal weights  $w$  for all layers:

- Abstract the whole network as a function:

$$L(w, x, y)$$

↑ *vectors*  
 ↳ *How Good is Our Fit*

- Consider all  $M$  training samples:

*Expected loss value*

$$\mathbb{E}_{x,y \sim \hat{p}_{\text{data}}(x,y)} [L(w, x, y)] = \frac{1}{M} \sum_{m=1}^M L(w, x, y)$$

- We now know what to do:

$$\underset{w}{\text{minimize}} \quad \{L(w, x, y)\}$$

## Gradient Descent

Minimize loss for all samples  
↳ on learning samples

one way  
Random  
Select best fit

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \frac{1}{M} \sum_{m=1}^M L(\mathbf{w}, \mathbf{x}, \mathbf{y}) \right\}$$

- Method of choice: Gradient Descent

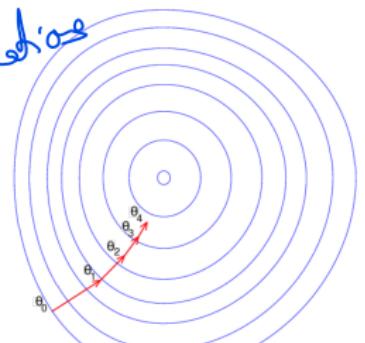
- Initialize  $\mathbf{w}$  (Random)
- Iterate until convergence:

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

- $\eta$  is commonly referred to as the **learning rate**

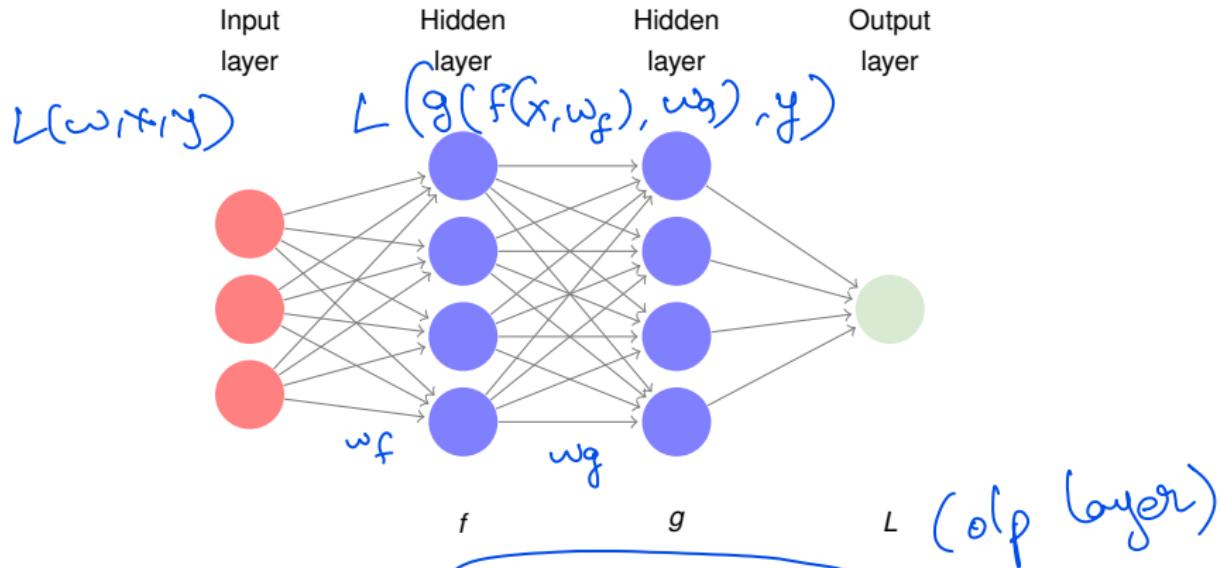
(new weight)  
iteration  $w^{k+1} = w^k - \eta \nabla_w L(w^k)$

Gradient  
(old) of loss function



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

## What is this $L$ we are trying to optimize?



- Complex network can be seen as composed functions:

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

*vectors*

**NEXT TIME  
ON DEEP LEARNING**

# Feed Forward Neural Networks - Part 3

A. Maier, V. Christlein, K. Breininger, Z. Yang, L. Rist, M. Nau, S. Jaganathan, C. Liu, N. Maul, L. Folle,  
K. Packhäuser, M. Zinnen

Pattern Recognition Lab, Friedrich-Alexander-Universität Erlangen-Nürnberg

April 24, 2023



# How to Calculate Derivatives in Complex Neural Networks?

## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

## Two algorithms:

- Finite differences
- Analytic derivative

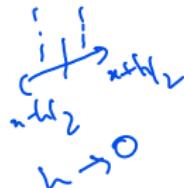
# Finite Differences

## Definition of derivative:

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

Due to finite precision the symmetric definition is preferred:

$$f'(x) = \lim_{h \rightarrow 0} \frac{f\left(x + \frac{1}{2}h\right) - f\left(x - \frac{1}{2}h\right)}{h}$$



## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x + \frac{1}{2}h) - f(x - \frac{1}{2}h)}{h}$$

Let's calculate it:

- Set  $h$  to  $2 \cdot 10^{-2}$

$$\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix} = \frac{\left( (2(1 + 10^{-2}) + 9)^2 + 3 \right) - \left( (2(1 - 10^{-2}) + 9)^2 + 3 \right)}{2 \cdot 10^{-2}}$$

## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x + \frac{1}{2}h) - f(x - \frac{1}{2}h)}{h}$$

Let's calculate it:

- Set  $h$  to  $2 \cdot 10^{-2}$

$$\begin{aligned} \frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix} &= \frac{\left( (2(1 + 10^{-2}) + 9)^2 + 3 \right) - \left( (2(1 - 10^{-2}) + 9)^2 + 3 \right)}{2 \cdot 10^{-2}} \\ &= \frac{(124.4404 - 123.5604)}{2 \cdot 10^{-2}} \end{aligned}$$

## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x + \frac{1}{2}h) - f(x - \frac{1}{2}h)}{h}$$

Let's calculate it:

- Set  $h$  to  $2 \cdot 10^{-2}$

$$\begin{aligned} \frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix} &= \frac{\left( (2(1 + 10^{-2}) + 9)^2 + 3 \right) - \left( (2(1 - 10^{-2}) + 9)^2 + 3 \right)}{2 \cdot 10^{-2}} \\ &= \frac{(124.4404 - 123.5604)}{2 \cdot 10^{-2}} \\ &= 43.9999 \end{aligned}$$

*here 43.9999  
there 44*

## Finite Differences Summed up

- For practical use it often suffices to use  $h = 1 \cdot 10^{-5}$
- For a more accurate derivative [7] use:  $h = \epsilon_f^{\frac{1}{3}} \cdot x_c$ 
  - Where  $\epsilon_f \approx 10^{-7}$
  - The characteristic scale is approximated as  $x_c = x$
  - Prevent division by zero at  $x = 0$

## Conclusion

- **Easy** to use
- We only need to be able to **evaluate** functions
- Computationally **inefficient**
- **Frequently used** to check implementations

Gradient =  
 set { All partial  
 Derivatives  
 $\nabla(f(100))$  variable  
 is not good  
 Just Verification

[CLMC]

## Analytic gradient

### Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

### Four analytic rules:

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

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

[CLMC]

## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

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

Let's calculate it:

## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

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

Let's calculate it:

$$\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix} = \frac{\partial}{\partial x_1} (2x_1 + 9)^2$$

Rules 1 and 2

## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

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

Let's calculate it:

$$\begin{aligned}\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix} &= \frac{\partial}{\partial x_1} (2x_1 + 9)^2 \\ &= \frac{\partial}{\partial z} (z)^2 \frac{\partial}{\partial x_1} (2x_1 + 9)\end{aligned}$$

Rules 1 and 2

Rule 4 and  $2x_1 + 9 = z$

## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

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

Let's calculate it:

$$\begin{aligned}
 \frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix} &= \frac{\partial}{\partial x_1} (2x_1 + 9)^2 && \text{Rules 1 and 2} \\
 &= \frac{\partial}{\partial z} (z)^2 \frac{\partial}{\partial x_1} (2x_1 + 9) && \text{Rule 4 and } 2x_1 + 9 = z \\
 &= 2(2x_1 + 9) \frac{\partial}{\partial x_1} (2x_1 + 9) && \text{Rule 3}
 \end{aligned}$$

## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

- $\frac{d}{dx} \text{const} = 0$
- $\frac{d}{dx}$  is linear
- $\frac{d}{dx} x^n = n \cdot x^{n-1}$
- $\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \cdot \frac{d}{dx} g(x)$

Let's calculate it:

$$\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix} = \frac{\partial}{\partial x_1} (2x_1 + 9)^2$$

Rules 1 and 2

$$= \frac{\partial}{\partial z} (z)^2 \frac{\partial}{\partial x_1} (2x_1 + 9)$$

Rule 4 and  $2x_1 + 9 = z$

$$= 2(2x_1 + 9) \frac{\partial}{\partial x_1} (2x_1 + 9)$$

Rule 3

$$= 2(2x_1 + 9) \cdot 2$$

Rules 1 and 2 and  $x_1 = 1$

44

## Analytic Gradient Summed up

- **Chain rule** and **Linearity** enable to **decompose** complex functions
- Analytic formulas have to be calculated manually
- Computationally more **efficient** than finite differences

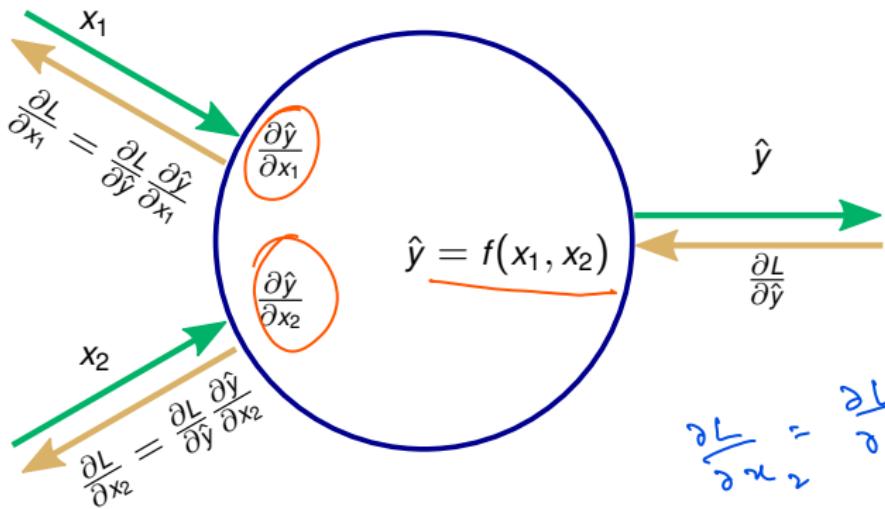
Calculate once; use for all

Can we compute analytic gradients automatically?

Analytical — More accurate

→ Compute Gradients *Automatically p. derivatives w.r.t Loss*

## Backpropagation Algorithm



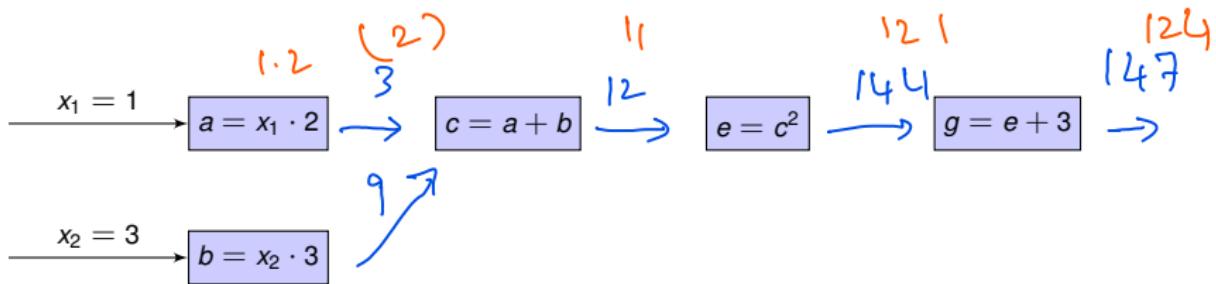
$$\frac{\partial L}{\partial x_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial x_2}$$

$$\exp(x_i) / \sum_{j=1}^n \exp x_j$$

1. Forward pass: Compute activations  $\rightarrow \hat{y}$
2. Backward pass: Recursively apply chain rule

## Example Problem

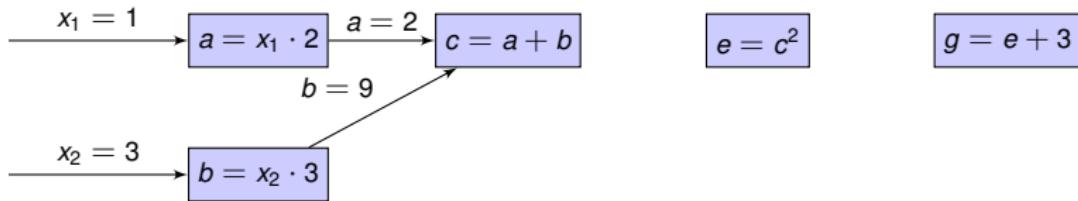
- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
  - Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$
- $$\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \frac{d}{dx} g(x)$$



## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

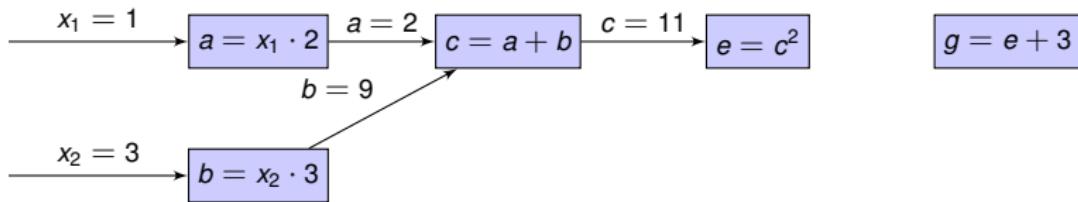
$$\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \frac{d}{dx} g(x)$$



## Example Problem

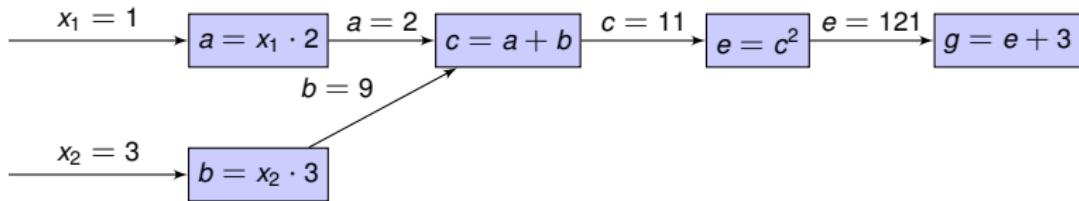
- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

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



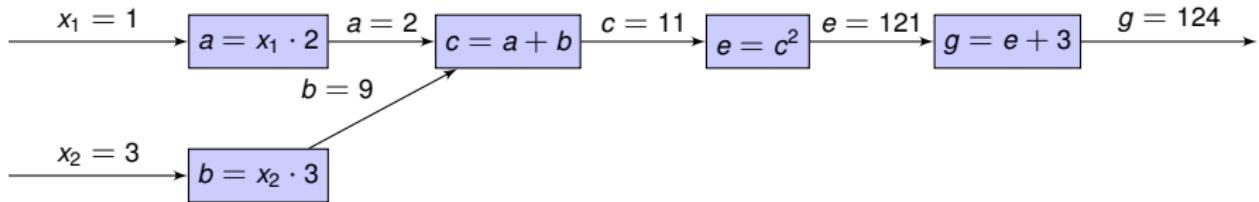
## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
  - Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$
- $$\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \frac{d}{dx} g(x)$$



## Example Problem

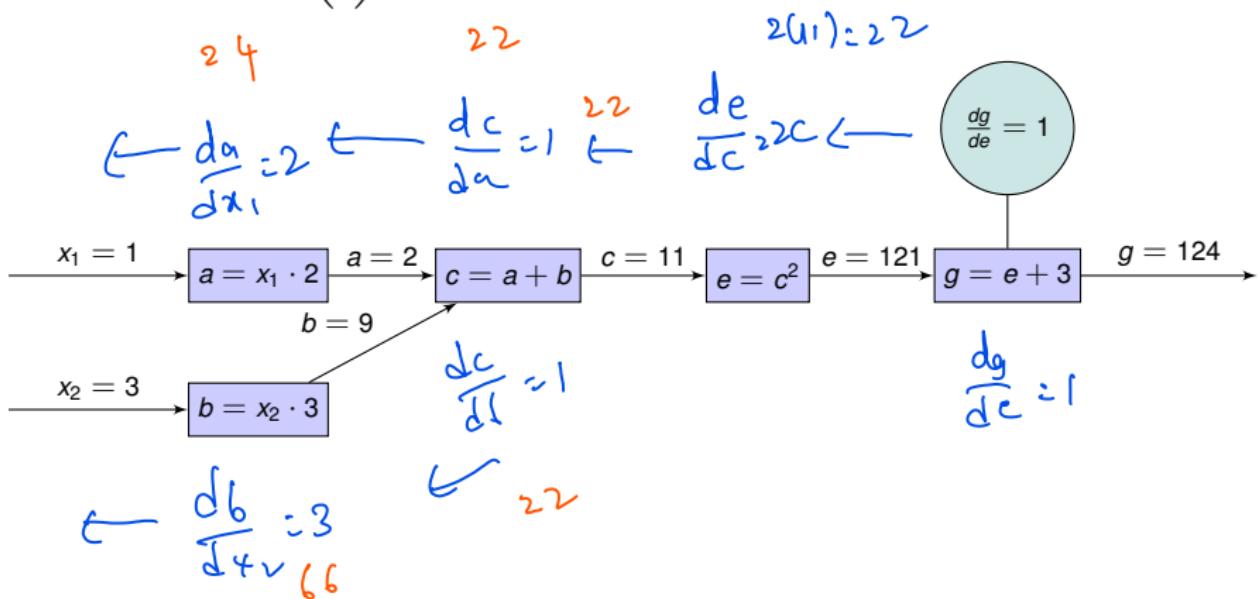
- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
  - Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$
- $$\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \frac{d}{dx} g(x)$$



## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

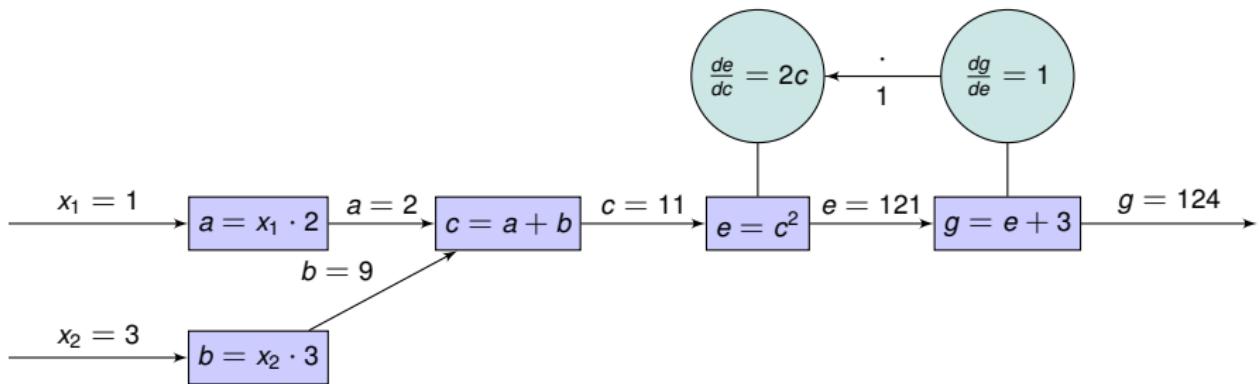
$$\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \frac{d}{dx} g(x)$$



## Example Problem

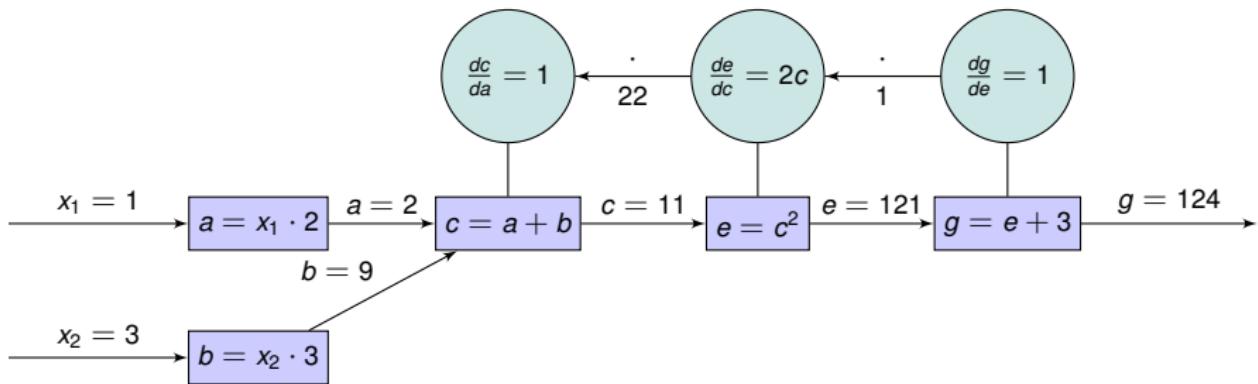
- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

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



## Example Problem

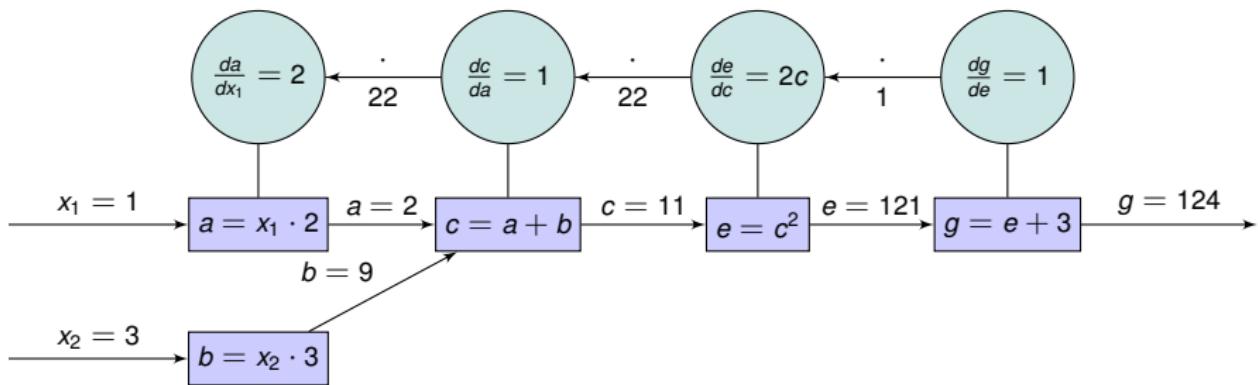
- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
  - Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$
- $$\frac{d}{dx} f(g(x)) = \frac{d}{dg} f(g) \frac{d}{dx} g(x)$$



## Example Problem

- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

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



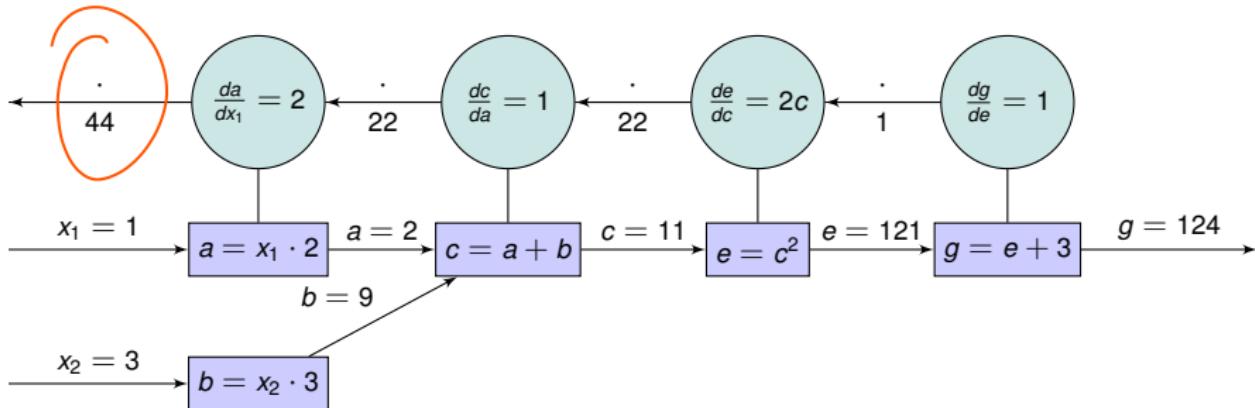
# Calculating Gradients

via BP

## Example Problem

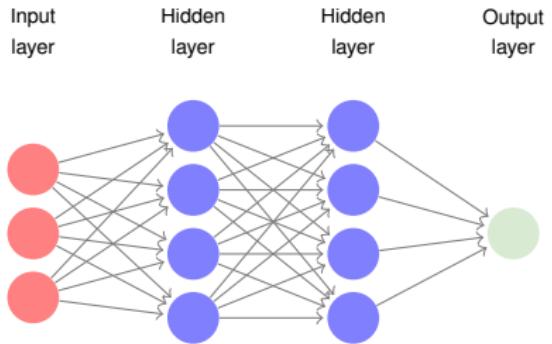
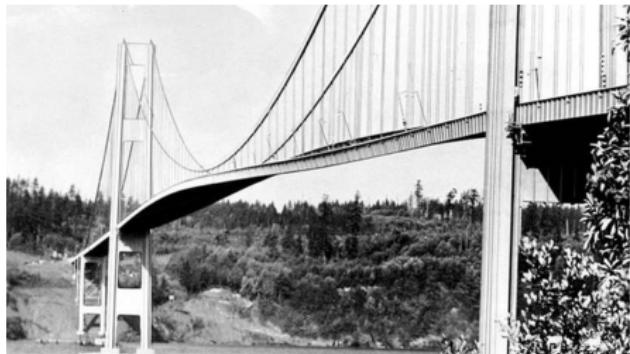
- Function:  $\hat{y} = f(\mathbf{x}) = (2x_1 + 3x_2)^2 + 3$
- Evaluate  $\frac{\partial}{\partial x_1} f \begin{pmatrix} 1 \\ 3 \end{pmatrix}$

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



## Stability Problem

- What do those two images have in common?

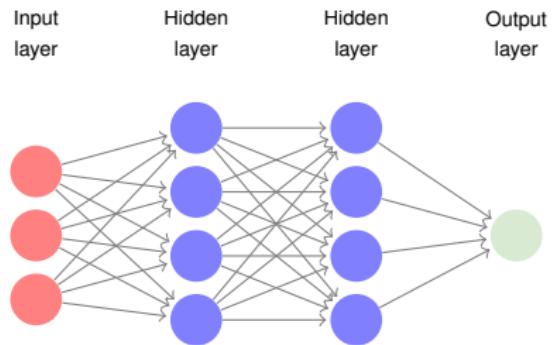
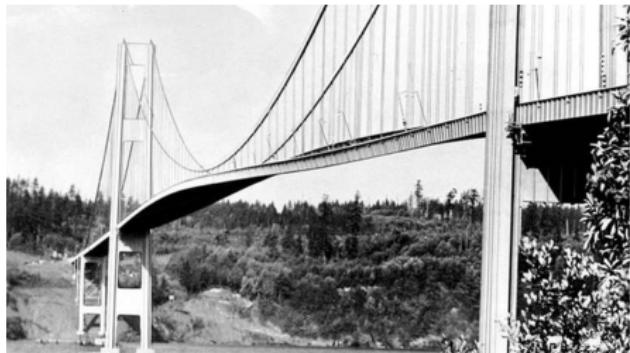


Click for video

Multiplying High & Low  
numbers  
frequently  
⇒ Disasters

## Stability Problem

- What do those two images have in common?

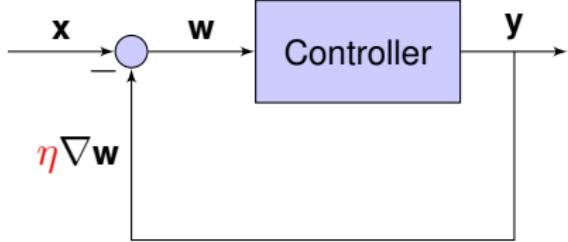


**Click for video**

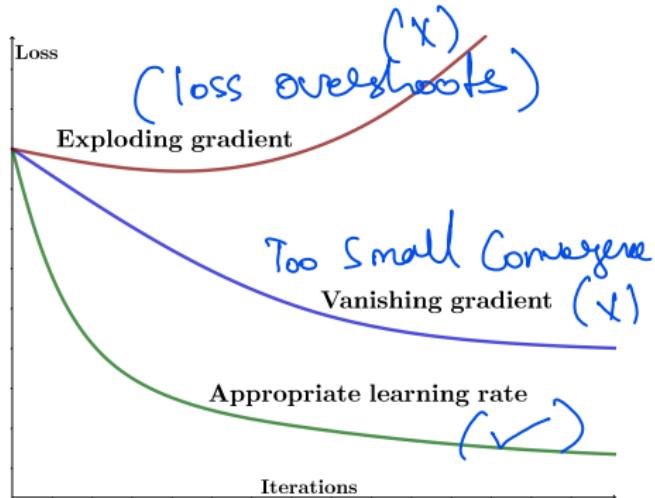
- Both suffer from positive feedback!
- This can cause disaster

(Bridge Shaking)

# Feedback loop



Analogy to control theory



- If  $\eta$  is too high → **positive feedback** → loss grows **without bounds**
- If  $\eta$  is too small → **negative feedback** → **gradient vanishes**
- Choice of  $\eta$  is **critical** for learning

## Backpropagation Summed up

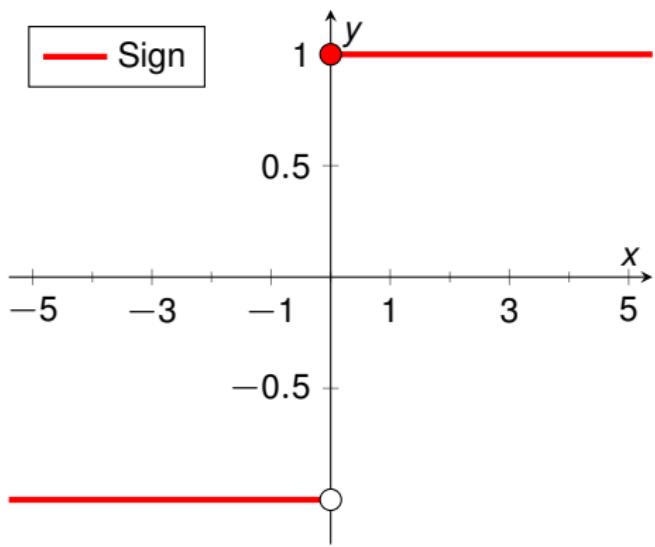
- Built around the **chain rule**
  - Uses a **forward-pass** through the function
  - Computationally very **efficient** by using a **dynamic programming** approach
  - Is **no training algorithm**, because it just computes a gradient
- $\nabla$ -calculation

### Consequences

- Product of partials  $\rightarrow$  numerical **errors multiply** low values
- Product of partials  $\rightarrow$  **vanishing** or **exploding** gradient

(Small numbers)      (High numbers)  
 very slow                  loss grows without bounds  
 (gradient vanishes)

## About the sign Activation Function



Sign

$$f(x) = \begin{cases} +1 & \text{for } x \geq 0 \\ -1 & \text{for } x < 0 \end{cases}$$

$f'(x) = 2\delta(x)$

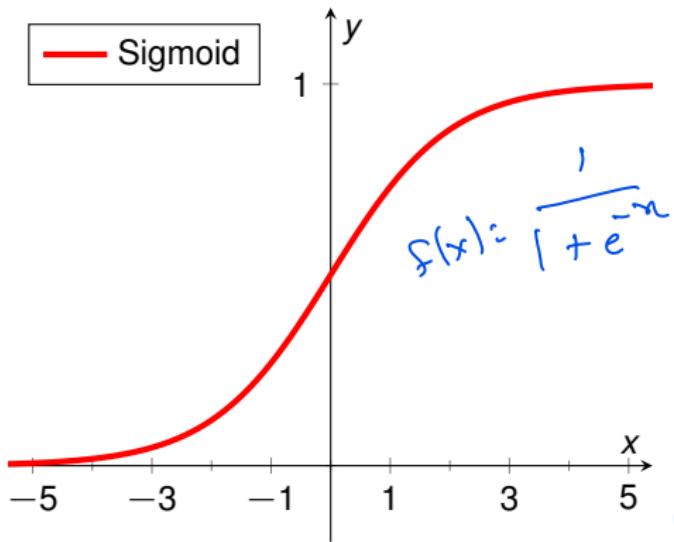
$\nabla$

Prob le mat ic

with Gradient  
Descent

- + Normalized output
- Gradient vanishes almost everywhere!

## Smooth Activation Function



Sigmoid (logistic function)

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

$$f'(x) = f(x)(1 - f(x))$$

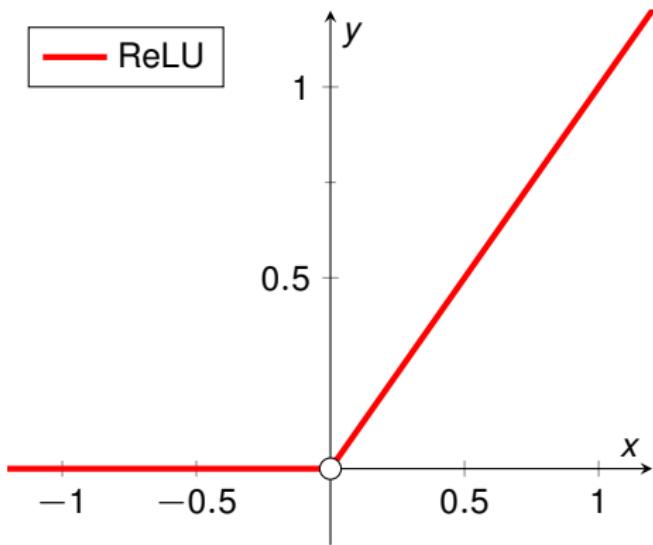
Quick efficiently.

$[-3, 73]$  — Vanishes  
 (low values)

+ Normalized output

- Gradient still eventually vanishes

## Piecewise-linear Activation Function



Rectified Linear Unit (ReLU)

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

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

+ Less vanishing gradient

Till now B.P  
at only node level.

**NEXT TIME**  
**ON DEEP LEARNING**



**FAU**

FRIEDRICH-ALEXANDER-  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG  
SCHOOL OF ENGINEERING

# Feed Forward Neural Networks - Part 4

A. Maier, V. Christlein, K. Breininger, Z. Yang, L. Rist, M. Nau, S. Jaganathan, C. Liu, N. Maul, L. Folle,  
K. Packhäuser, M. Zinnen

Pattern Recognition Lab, Friedrich-Alexander-Universität Erlangen-Nürnberg

April 24, 2023





**FAU**

FRIEDRICH-ALEXANDER-  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG  
SCHOOL OF ENGINEERING

# Layer Abstraction



## From Graphs of Nodes to Graphs of Layers

- We introduced **layers** but computed everything on individual nodes
- It is convenient to add further abstraction
- But how can we express this?

### Recall: Single neuron

- Add a bias unit to  $\mathbf{x} \in \mathbb{R}^{N-1}$  by adding a dimension with  $x_n = 1$
- This is a connection from every input element to the single output element:

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

*Saturation = Vector · vector*

$$\hat{y} = \mathbf{w}^T \cdot \mathbf{x}^T$$

$$[ ] \quad [ ]$$

## Representing the connections

- Assume we have  $M$  neurons  $\rightarrow M$  sets of weights:  $\mathbf{w}_m$  for  $m \in \{1, \dots, M\}$

$$\hat{y}_m = \mathbf{w}_m^T \mathbf{x}$$

- We rewrite this operation as matrix-vector multiplication:

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

- This is known as **fully connected layer**.
- It represents any arbitrary connection topology between layers.
- We can describe back-propagation in this more abstract view as well!

Convert from Indv.  $\rightarrow$  layers  
 by using Matrix notation

Woi

Earlier - scalar -  $\hat{y} = w^T \cdot x$

## Fully Connected Layer



$$\sigma_j = \dots$$

$$\hat{y} = w^T \cdot x$$

$$W = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$$

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

belongs to class

$[ ]$

- The forward-pass is:

$$W = \begin{pmatrix} w_1^T \\ w_2^T \end{pmatrix}$$

$$\hat{y} = Wx$$

Notice - NO bias pass

- After the forward-pass through all layers, we can compute a loss that depends on our loss function  $L$ .

- We need two gradients for the backward-pass:

- Gradient with respect to the weights:  $\frac{\partial L}{\partial w}$  for gradient descend

- Gradient with respect to the inputs:  $\frac{\partial L}{\partial x}$  for backpropagation

## Fully Connected Layer Summed up

- Can represent any connection topology
- Enables higher level view concentrating on layers instead of nodes
- Is a matrix multiplication:

$$\hat{y} = Wx$$

- Its gradient with respect to the weights:

matrix

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial W} = \frac{\partial L}{\partial \hat{y}} x^T$$

- Its gradient with respect to the input:

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial x} = W^T \frac{\partial L}{\partial \hat{y}}$$

col vector  
 $\begin{bmatrix} \cdot \\ \vdots \\ i \\ \cdot \\ n \end{bmatrix}$

col  
 $\frac{\partial y}{\partial w} = f^T$   
 $\begin{bmatrix} \cdot \\ \vdots \\ i+0 \\ \cdot \\ n \end{bmatrix}$   
 opt

layer  $\rightarrow$  nodes



## Fully Connected Layer: Simple example

- Assume we are looking at a simple network (no activation function) with the forward pass:

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

- We try to find parameters  $\mathbf{W}$  that minimize the following loss function:

*(L<sub>2</sub>-Norm)<sup>2</sup> Regression loss*

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

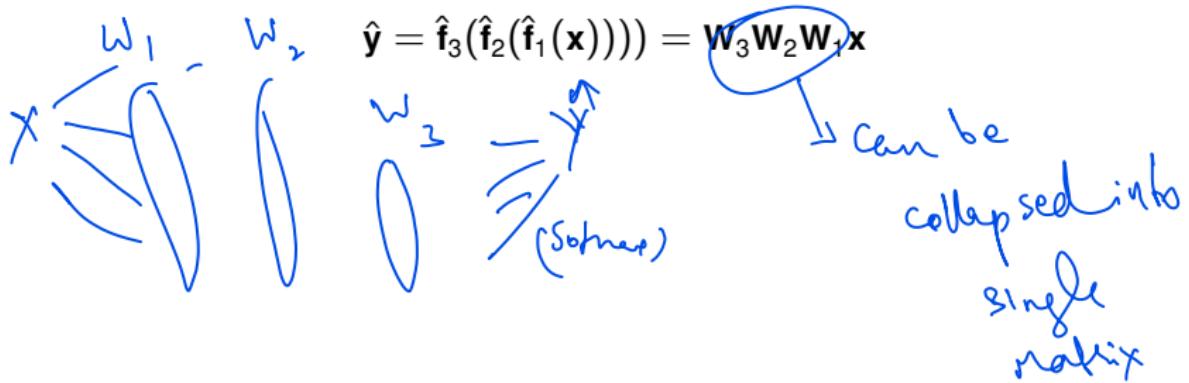
$$L(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{2} (\hat{\mathbf{y}} - \mathbf{y})^2$$

$$\approx \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \frac{1}{2} \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \hat{\mathbf{y}} + \frac{1}{2} \hat{\mathbf{y}}^T \hat{\mathbf{y}}$$

- Then simply  $\frac{\partial L}{\partial \hat{\mathbf{y}}} = \hat{\mathbf{y}} - \mathbf{y} = \mathbf{W}\mathbf{x} - \mathbf{y}$
- The gradient with respect to the weights:  $\frac{\partial L}{\partial \mathbf{W}} = (\mathbf{W}\mathbf{x} - \mathbf{y})\mathbf{x}^T$
- The gradient with respect to the inputs:  $\frac{\partial L}{\partial \mathbf{x}} = \mathbf{W}^T(\mathbf{W}\mathbf{x} - \mathbf{y})$

## Linear Network in Matrix notation

- Let's add some layers



## Linear Network in Matrix notation

- Let's add some layers

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

*Loss* :  $L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

## Linear Network in Matrix notation

- Let's add some layers

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3\mathbf{W}_2\mathbf{W}_1\mathbf{x}$$

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3\mathbf{W}_2\mathbf{W}_1\mathbf{x} - \mathbf{y}\|_2^2$$

- Gradients?

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Last layer gradient

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

$$L = \frac{1}{2} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2$$

- Last layer gradient

$$\frac{\partial L}{\partial \mathbf{W}_3} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \cdot \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{w}} = \frac{\partial L}{\partial \hat{f}_3} \cdot \frac{\partial \hat{f}_3}{\partial \mathbf{w}_3}$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Last layer gradient

$$\frac{\partial L}{\partial \mathbf{W}_3} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})}$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Last layer gradient

$$\frac{\partial L}{\partial \mathbf{W}_3} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \cdot \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_3}}_{(\mathbf{W}_2 \mathbf{W}_1 \mathbf{x})^T}$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\begin{pmatrix} w_3 \\ w_2 \\ w_1 \\ x-y \end{pmatrix}^\top \begin{pmatrix} w_3 \\ w_2 \\ w_1 \\ x \end{pmatrix}$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Last layer gradient

$$\frac{\partial L}{\partial \mathbf{W}_3} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_3}}_{\cdot (\mathbf{W}_2 \mathbf{W}_1 \mathbf{x})^\top} = (\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})(\mathbf{W}_2 \mathbf{W}_1 \mathbf{x})^\top$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Deeper gradients

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2}$$

## Linear Network in Matrix notation

- Associated loss function:

$$\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta},$$

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{y} = \hat{f}_3(\hat{f}_2(\hat{f}_1(\mathbf{x}))) = \mathbf{W}_3(\mathbf{W}_2(\mathbf{W}_1(\mathbf{x})))$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{f}_3} \frac{\partial \hat{f}_3}{\partial \mathbf{W}_2}$$

$$\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \mathbf{w}_2} \quad \hat{f}_3(z) \cdot \mathbf{w}_3(z)$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})}$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2}}_{(\mathbf{W}_3)^T}.$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2}}_{(\mathbf{W}_3)^T} \cdot \underbrace{\frac{\partial \hat{\mathbf{f}}_2}{\partial \mathbf{W}_2}}_{(\mathbf{W}_1 \mathbf{x})^T}$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2}}_{(\mathbf{W}_3)^T} \cdot \underbrace{\frac{\partial \hat{\mathbf{f}}_2}{\partial \mathbf{W}_2}}_{\cdot (\mathbf{W}_1 \mathbf{x})^T} = \boxed{\mathbf{W}_3^T (\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}) (\mathbf{W}_1 \mathbf{x})^T}$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2}}_{(\mathbf{W}_3)^T} \underbrace{\frac{\partial \hat{\mathbf{f}}_2}{\partial \mathbf{W}_2}}_{\cdot (\mathbf{W}_1 \mathbf{x})^T} = \mathbf{W}_3^T (\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}) (\mathbf{W}_1 \mathbf{x})^T$$

$$\frac{\partial L}{\partial \mathbf{W}_1}$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

$$\begin{aligned} & \frac{1}{2} \|\mathbf{f}_2(\mathbf{y}) - \mathbf{z}\|_2^2 \\ &= \frac{1}{2} \|\hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x})) - \mathbf{y})\|_2^2 \end{aligned}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2}}_{(\mathbf{W}_3)^T} \underbrace{\frac{\partial \hat{\mathbf{f}}_2}{\partial \mathbf{W}_2}}_{\cdot (\mathbf{W}_1 \mathbf{x})^T} = \mathbf{W}_3^T (\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}) (\mathbf{W}_1 \mathbf{x})^T$$

$$\frac{\partial L}{\partial \mathbf{W}_1} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_1} \quad \text{=} \quad$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2}}_{(\mathbf{W}_3)^T} \underbrace{\frac{\partial \hat{\mathbf{f}}_2}{\partial \mathbf{W}_2}}_{\cdot (\mathbf{W}_1 \mathbf{x})^T} = \mathbf{W}_3^T (\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}) (\mathbf{W}_1 \mathbf{x})^T$$

$$\frac{\partial L}{\partial \mathbf{W}_1} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_1} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2} \frac{\partial \hat{\mathbf{f}}_2}{\partial \mathbf{W}_1}$$

## Linear Network in Matrix notation

$$\hat{y} = \frac{1}{2} \|\hat{f}_3(\hat{f}_2(\hat{f}_1(x))) - y\|_2^2$$

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{f}_3(\hat{f}_2(\hat{f}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

$$\begin{aligned} f_3(f_2(x)) \\ = w_3 w_2 x \end{aligned}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{f}_3} \frac{\partial \hat{f}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{f}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{f}_3}{\partial \hat{f}_2}}_{(\mathbf{w}_3)^T} \underbrace{\frac{\partial \hat{f}_2}{\partial \mathbf{W}_2}}_{\cdot (\mathbf{W}_1 \mathbf{x})^T} = \mathbf{W}_3^T (\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}) (\mathbf{W}_1 \mathbf{x})^T$$

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}_1} &= \frac{\partial L}{\partial \hat{f}_3} \frac{\partial \hat{f}_3}{\partial \mathbf{W}_1} = \underbrace{\frac{\partial L}{\partial \hat{f}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{f}_3}{\partial \hat{f}_2}}_{\frac{\partial \hat{f}_3}{\partial \hat{f}_2}} \underbrace{\frac{\partial \hat{f}_2}{\partial \mathbf{W}_1}}_{\cdot \mathbf{w}_2^T [\mathbf{W}_1 \mathbf{x}]} \\ &\quad \cdot \frac{\partial \hat{f}_2}{\partial \mathbf{W}_1} = \underbrace{\frac{\partial L}{\partial \hat{f}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \cdot \frac{\partial \hat{f}_3}{\partial \hat{f}_2} \cdot \mathbf{w}_2^T [\mathbf{W}_1 \mathbf{x}] \\ &\quad \cdot \frac{\partial \hat{f}_2}{\partial \mathbf{W}_1} = \underbrace{\frac{\partial L}{\partial \hat{f}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \cdot \frac{\partial \hat{f}_3}{\partial \hat{f}_2} \cdot \mathbf{w}_2^T \cdot \mathbf{w}_3^T \cdot \mathbf{w}_1^T [\mathbf{W}_1 \mathbf{x}] \end{aligned}$$

$$f_1(x) = w_1 x$$

$$f_2(x) = w_2 x$$

$$f_3(x) = w_3 x$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \| \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y} \|_2^2$$

$$\hat{\mathbf{y}} = \hat{f}_3(\hat{f}_2(\hat{f}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{f}_3} \frac{\partial \hat{f}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{f}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{f}_3}{\partial \hat{f}_2}}_{(\mathbf{W}_3)^T} \underbrace{\frac{\partial \hat{f}_2}{\partial \mathbf{W}_2}}_{\cdot (\mathbf{W}_1 \mathbf{x})^T} = \mathbf{W}_3^T (\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}) (\mathbf{W}_1 \mathbf{x})^T$$

$$\frac{\partial L}{\partial \mathbf{W}_1} = \frac{\partial L}{\partial \hat{f}_3} \frac{\partial \hat{f}_3}{\partial \mathbf{W}_1} = \frac{\partial L}{\partial \hat{f}_3} \frac{\partial \hat{f}_3}{\partial \hat{f}_2} \frac{\partial \hat{f}_2}{\partial \mathbf{W}_1} = \underbrace{\frac{\partial L}{\partial \hat{f}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{f}_3}{\partial \hat{f}_2}}_{(\mathbf{W}_3)^T} \underbrace{\frac{\partial \hat{f}_2}{\partial \mathbf{W}_1}}_{\cdot (\mathbf{W}_2)^T \cdot (\mathbf{x})^T}$$

$$\hat{f}_3 = \mathbf{w}_3 \hat{f}_2$$

## Linear Network in Matrix notation

- Associated loss function:

$$L(\theta) = \frac{1}{2} \|\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}\|_2^2$$

$$\hat{\mathbf{y}} = \hat{\mathbf{f}}_3(\hat{\mathbf{f}}_2(\hat{\mathbf{f}}_1(\mathbf{x}))) = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

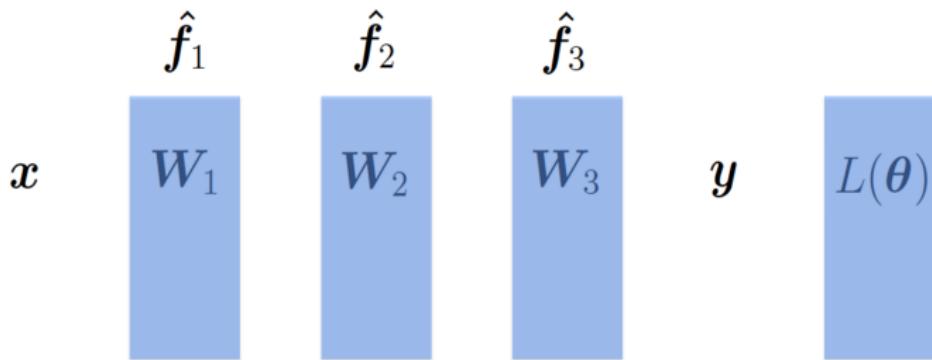
- Deeper gradients

$$\frac{\partial L}{\partial \mathbf{W}_2} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_2} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2}}_{(\mathbf{W}_3)^T} \underbrace{\frac{\partial \hat{\mathbf{f}}_2}{\partial \mathbf{W}_2}}_{\cdot (\mathbf{W}_1 \mathbf{x})^T} = \mathbf{W}_3^T (\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}) (\mathbf{W}_1 \mathbf{x})^T$$

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}_1} &= \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \mathbf{W}_1} = \frac{\partial L}{\partial \hat{\mathbf{f}}_3} \frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2} \frac{\partial \hat{\mathbf{f}}_2}{\partial \mathbf{W}_1} = \underbrace{\frac{\partial L}{\partial \hat{\mathbf{f}}_3}}_{(\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y})} \underbrace{\frac{\partial \hat{\mathbf{f}}_3}{\partial \hat{\mathbf{f}}_2}}_{(\mathbf{W}_3)^T} \underbrace{\frac{\partial \hat{\mathbf{f}}_2}{\partial \mathbf{W}_1}}_{\cdot (\mathbf{W}_2)^T \cdot (\mathbf{x})^T} \\ &= \mathbf{W}_2^T \mathbf{W}_3^T (\mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} - \mathbf{y}) (\mathbf{x})^T \end{aligned}$$

B.P similar to above

## Linear Network in Matrix notation

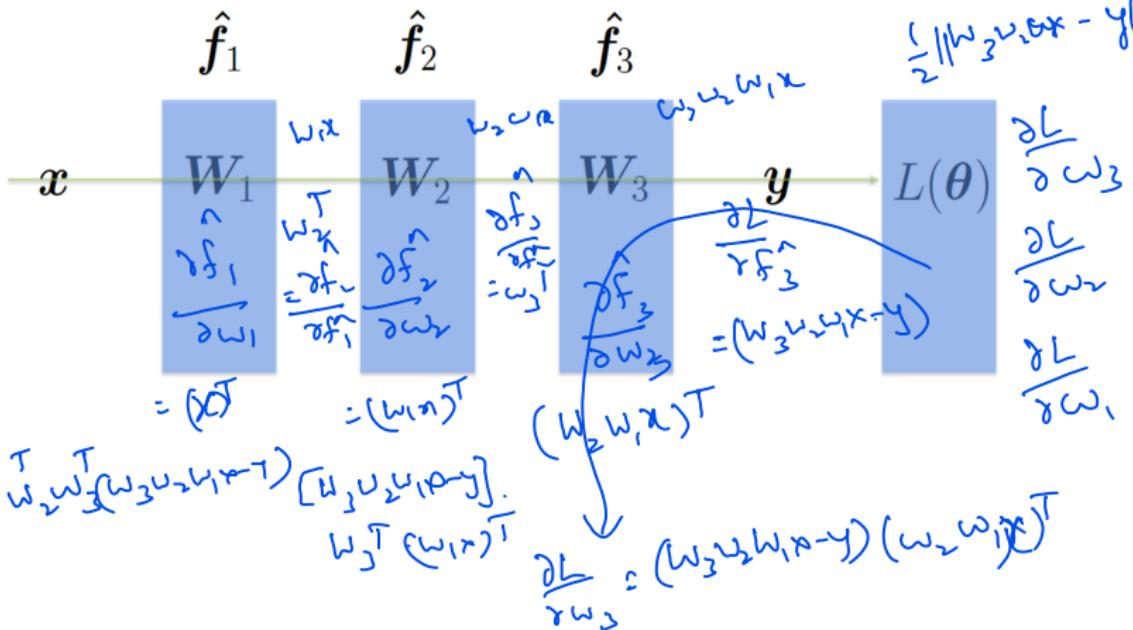


## Linear Network in Matrix notation

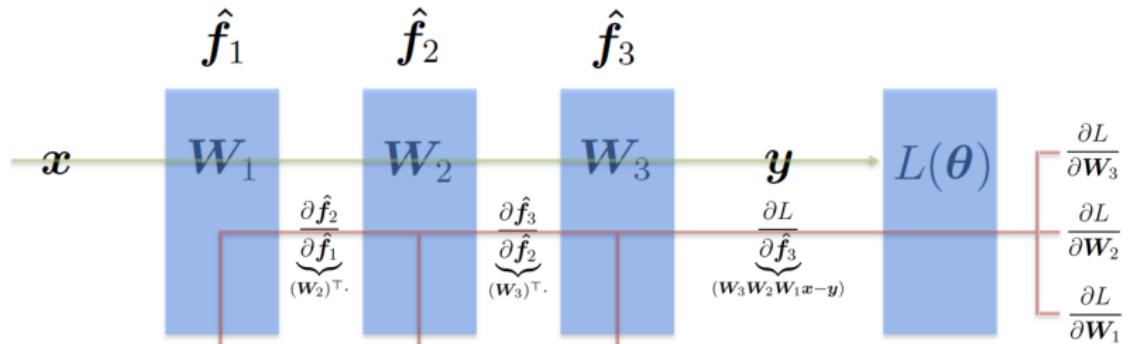
$$\hat{f}_1 = w_3 \hat{f}_2$$

$$\hat{f}_2 = w_2 \hat{f}_1$$

$$\hat{f}_1 = w_1 x$$



## Linear Network in Matrix notation



$$W_2^T W_3^T (W_3 W_2 W_1 x - y)(x)^T \quad W_3^T (W_3 W_2 W_1 x - y)(W_1 x)^T \quad (W_3 W_2 W_1 x - y)(W_2 W_1 x)^T$$

many items can be  
Hence very effective  
precomputed.

## Summary

$$y_k = \frac{e^{x_k}}{\sum_i e^{x_i}}$$

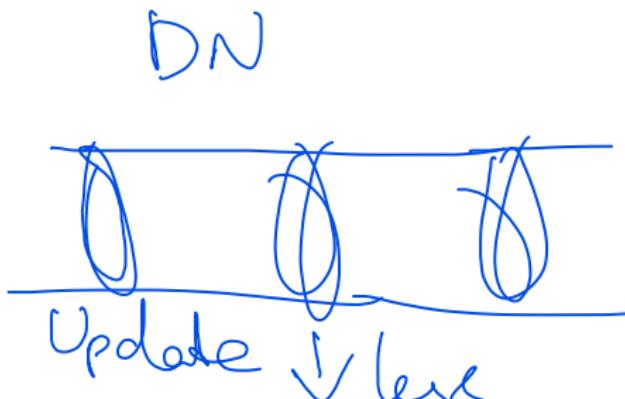
$$-\log(\hat{y}_k) \Big|_{\hat{y}_k=1}$$

- Softmax activation function with cross entropy loss mostly go together as "Softmax Loss".
  - Gradient descent is our default training algorithm in deep learning. *(local minima)*
  - We can compute gradients using finite differences to check our implementation.
  - We use the backpropagation algorithm to compute gradients efficiently.
  - The fully connected layer is the most general connectivity between layers in a feed-forward neural network.
- not training just slopes*

**NEXT TIME  
ON DEEP LEARNING**

next items

- Problem adapted loss functions
- Sophisticated optimization routines
- Optimization adapted to the needs of every single parameter
- An argument why neural networks shouldn't perform well
- Some very recent insights why they do perform well



## Comprehensive Questions

→ softmax loss (CS)

- Name a loss function for multi-class classification in deep learning.
- Explain how this loss function works.
- How can you check if the derivative implementation of a loss function is correct?
- What does backpropagation do?
- How does backpropagation work?
- Explain the exploding and vanishing gradient problems.
- Why is the signum function not used in deep learning?

→ finite differences  
 → analytical derivation  
 (negative fb)  
 loss groups without ends  
 Non-differentiable except 0  
 Explodes based on chain rule.  
 Vanishes good enough.

## Further Reading

- [Link](#) - The original paper popularizing ReLUs
- [Link](#) - The original paper popularizing backpropagation
- [Link](#) - Bishop - Mathematical compendium for machine learning
- [Link](#) - Blog article putting backpropagation in a very general context



**FAU**

FRIEDRICH-ALEXANDER-  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG  
SCHOOL OF ENGINEERING

# References



## References I

- [1] R. O. Duda, P. E. Hart, and D. G. Stork. Pattern Classification. John Wiley and Sons, inc., 2000.
- [2] Christopher M. Bishop. Pattern Recognition and Machine Learning (Information Science and Statistics). Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [3] F. Rosenblatt. "The perceptron: A probabilistic model for information storage and organization in the brain.". In: Psychological Review 65.6 (1958), pp. 386–408.
- [4] WS. McCulloch and W. Pitts. "A logical calculus of the ideas immanent in nervous activity.". In: Bulletin of mathematical biophysics 5 (1943), pp. 99–115.
- [5] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. "Learning representations by back-propagating errors.". In: Nature 323 (1986), pp. 533–536.

## References II

- [6] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. “Deep Sparse Rectifier Neural Networks”. In:  
Proceedings of the Fourteenth International Conference on Artificial Intelligence  
Vol. 15. 2011, pp. 315–323.
- [7] William H. Press, Saul A. Teukolsky, William T. Vetterling, et al.  
Numerical Recipes 3rd Edition: The Art of Scientific Computing. 3rd ed. New York, NY, USA: Cambridge University Press, 2007.