

Deep Learning

A course about theory & practice

Deep Neural Networks

Marco Piastra

Feed-Forward Neural Networks (recap)

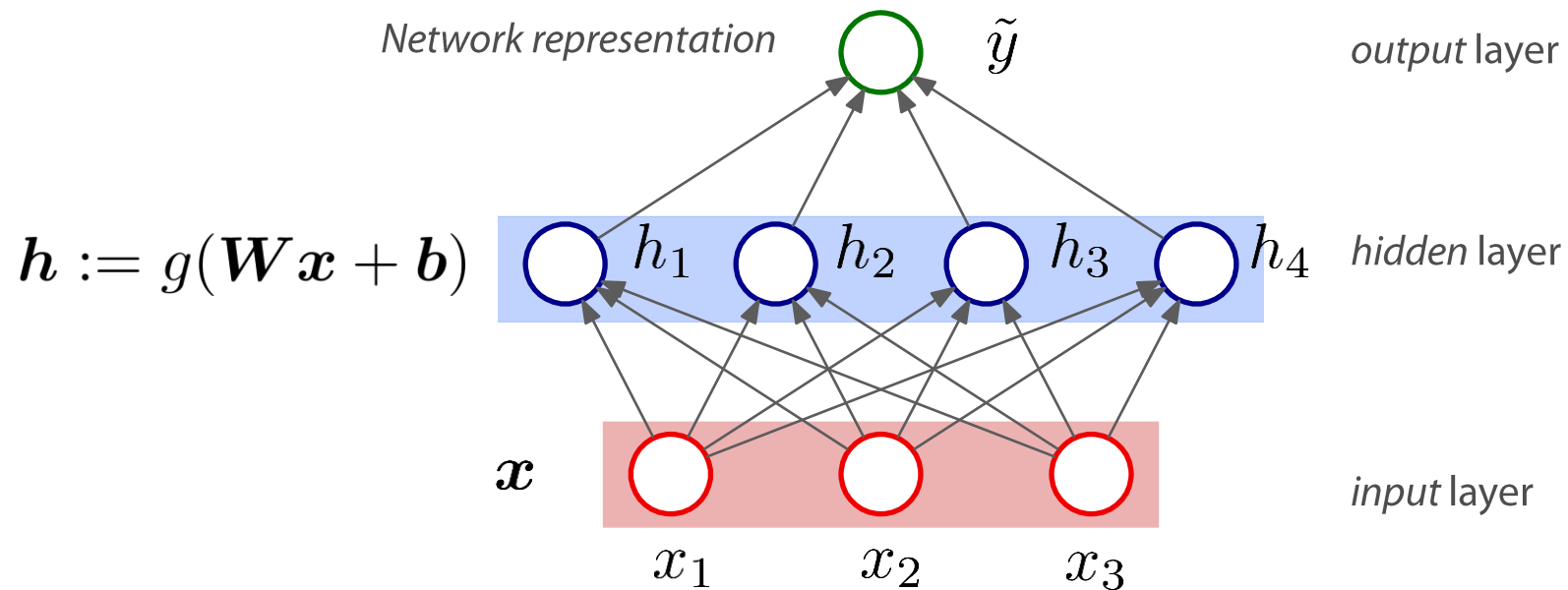
Feed-Forward Neural Network

- Approximating a target function

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Universal approximator: **feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b, \quad \mathbf{W} \in \mathbb{R}^{h \times d}, \quad \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, b \in \mathbb{R}$$



Feed-Forward Neural Network

- Approximating a target function

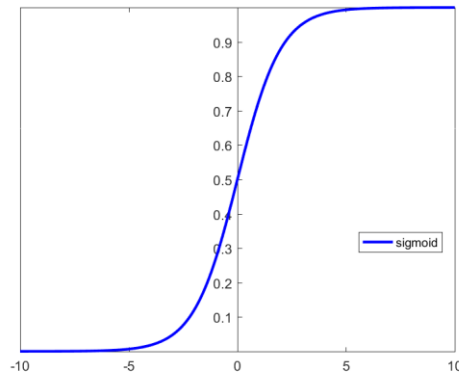
$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Universal approximator: **feed-forward neural network**

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b, \quad \mathbf{W} \in \mathbb{R}^{h \times d}, \mathbf{w}, \mathbf{b} \in \mathbb{R}^h, b \in \mathbb{R}$$

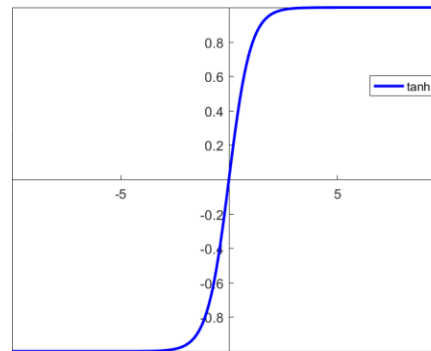
Popular choices for the non-linear function:

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



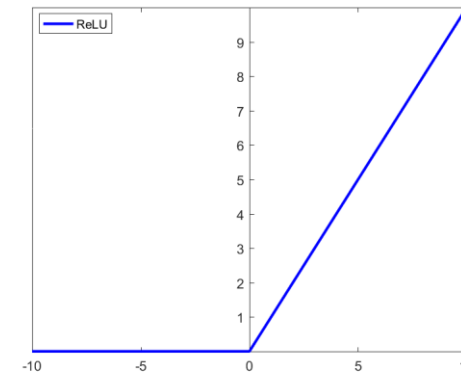
Sigmoid

$$g(x) = \tanh(x)$$



Hyperbolic Tangent

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



ReLU

Training Feed-Forward Neural Networks

■ Stochastic Gradient Descent (SGD)

1. Assign initial values to the four parameters $\mathbf{W}^{(0)}$, $\mathbf{b}^{(0)}$, $\mathbf{w}^{(0)}$, $b^{(0)}$
2. Pick up a data item $(\mathbf{x}^{(i)}, y^{(i)})$ from D with uniform probability and update the four parameters (with $\eta \ll 1.0$, $\eta \rightarrow 0$ as iterations progress)

$$\Delta \mathbf{W} = -\eta \frac{\partial}{\partial \mathbf{W}} L(\tilde{y}^{(i)}, y^{(i)})$$

$$\Delta \mathbf{b} = -\eta \frac{\partial}{\partial \mathbf{b}} L(\tilde{y}^{(i)}, y^{(i)})$$

$$\Delta \mathbf{w} = -\eta \frac{\partial}{\partial \mathbf{w}} L(\tilde{y}^{(i)}, y^{(i)})$$

$$\Delta b = -\eta \frac{\partial}{\partial b} L(\tilde{y}^{(i)}, y^{(i)})$$

3. Unless complete, return to step 2.

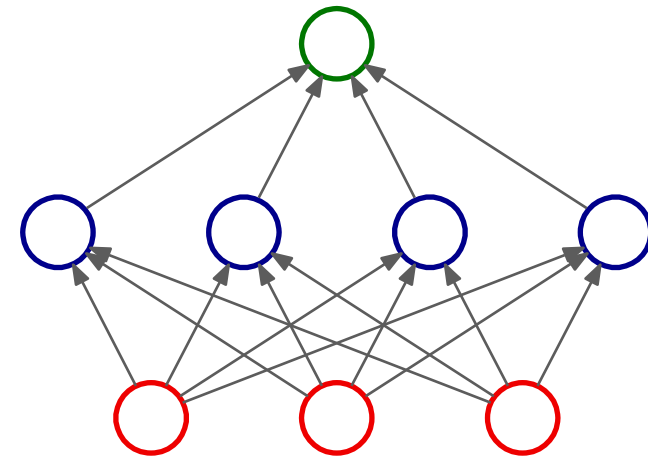
The Quest for Deeper Networks

Shallow vs. Deep Feed-Forward Neural Networks

- **Increasing network depth**

A feed-forward neural network with one hidden layer

$$\tilde{y} = \boldsymbol{w} \cdot g(\boldsymbol{W}^{[1]} \boldsymbol{x} + \boldsymbol{b}^{[1]}) + b$$

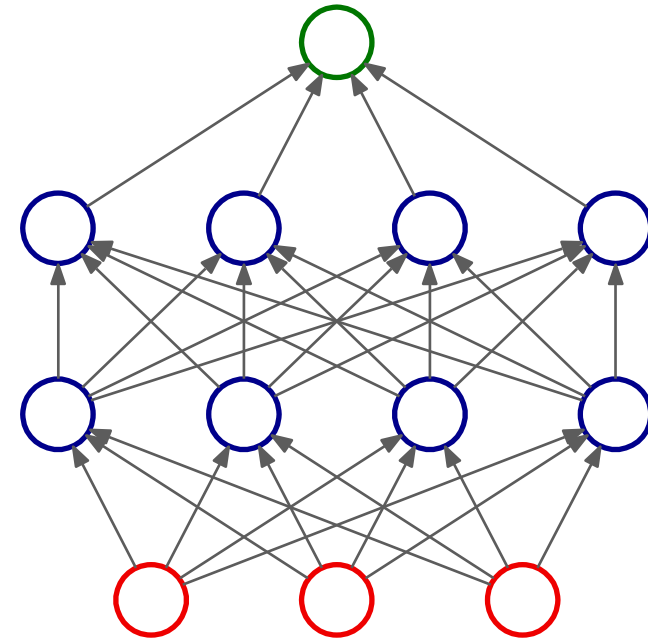


Shallow vs. Deep Feed-Forward Neural Networks

- **Increasing network depth**

A feed-forward neural network with two hidden layers

$$\tilde{y} = \boldsymbol{w} \cdot g(\boldsymbol{W}^{[2]}g(\boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}) + \boldsymbol{b}^{[2]}) + b$$

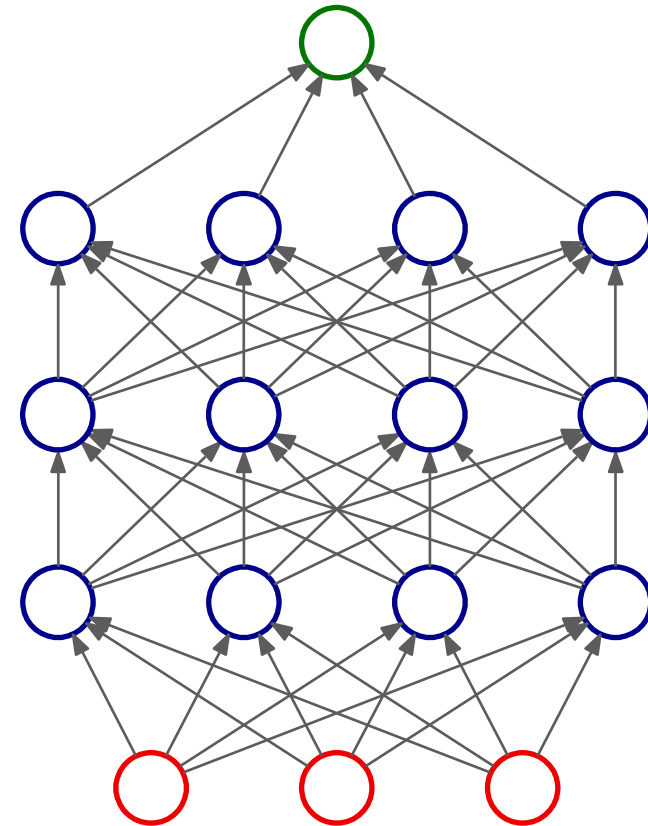


Shallow vs. Deep Feed-Forward Neural Networks

- **Increasing network depth**

A feed-forward neural network with three hidden layers

$$\tilde{y} = \boldsymbol{w} \cdot g(\boldsymbol{W}^{[3]}g(\boldsymbol{W}^{[2]}g(\boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}) + \boldsymbol{b}^{[2]}) + \boldsymbol{b}^{[3]}) + b$$



Shallow vs. Deep Feed-Forward Neural Networks

■ Increasing network depth

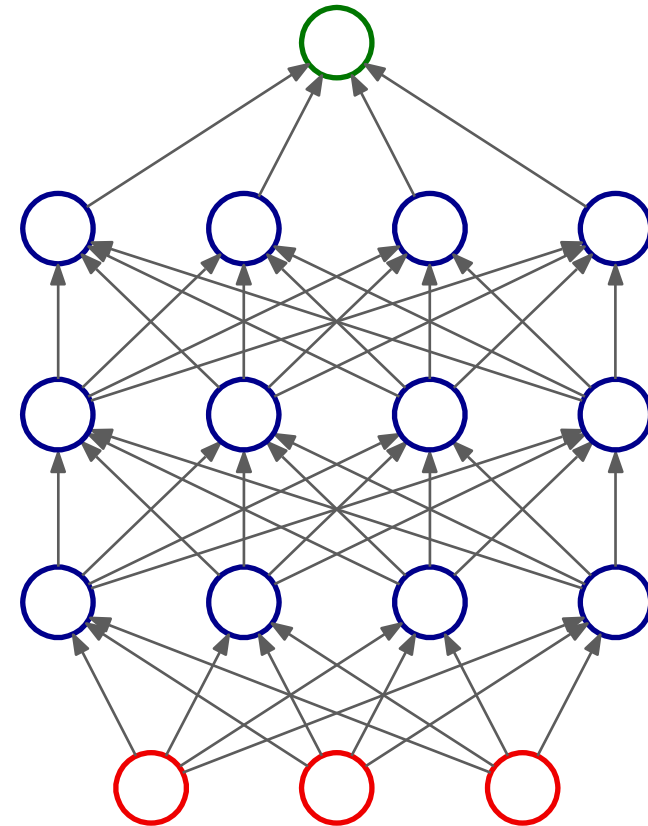
A feed-forward neural network with three hidden layers

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}^{[3]}g(\mathbf{W}^{[2]}g(\mathbf{W}^{[1]}\mathbf{x} + \mathbf{b}^{[1]}) + \mathbf{b}^{[2]}) + \mathbf{b}^{[3]}) + b$$

OK, but what is there to gain from such increase in depth?

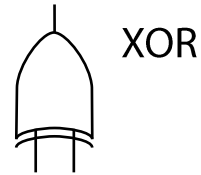
After all, the universal approximation theorem says that one layer is enough...

...and each layer brings in some extra complexity and further parameters.



Parity Circuits

A logical circuit whose output is 1 whenever the number of 1s in input is odd



XOR

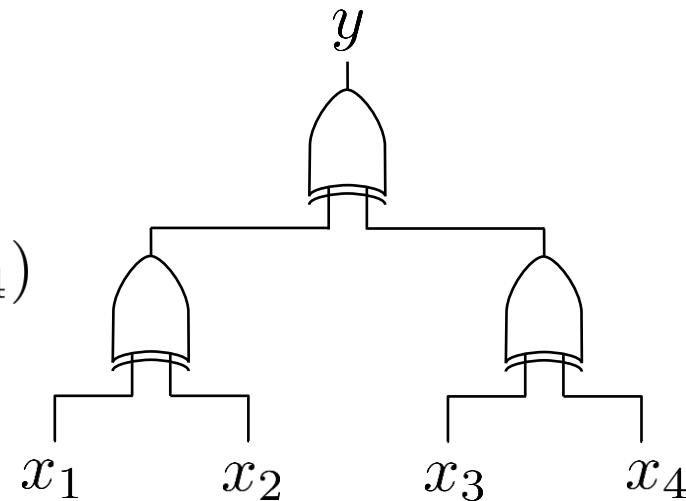
x_1	x_2	$x_1 \oplus x_2$
0	0	0
0	1	1
1	0	1
1	1	0

For instance:

$$\mathbf{x} = [0, 1, 1, 0] \rightarrow y = 0$$

$$\mathbf{x} = [1, 1, 0, 1] \rightarrow y = 1$$

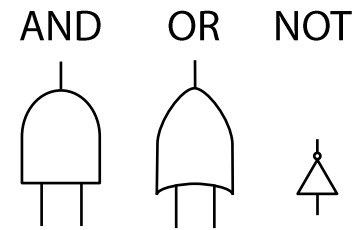
$$(x_1 \oplus x_2) \oplus (x_3 \oplus x_4)$$



This is an implementation using XOR components

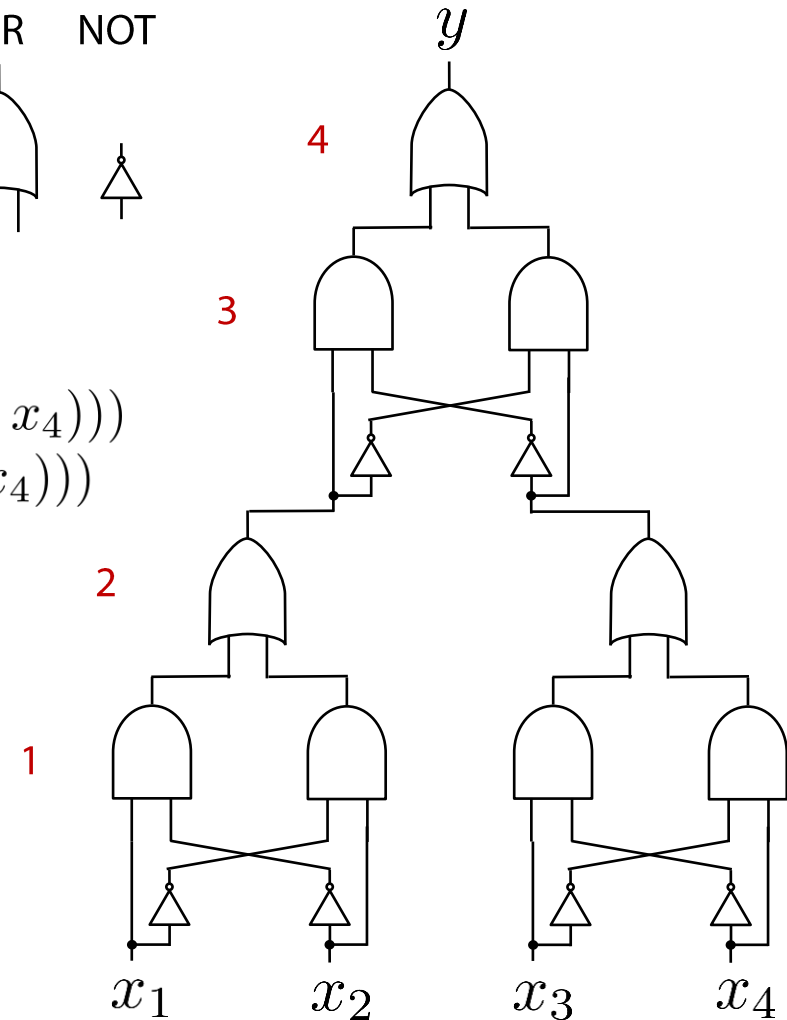
Parity Circuits

An implementation of the same parity circuit using AND, OR and NOT



$$\begin{aligned} &(((x_1 \wedge \neg x_2) \vee (\neg x_1 \wedge x_2)) \wedge \neg((x_3 \wedge \neg x_4) \vee (\neg x_3 \wedge x_4))) \\ &\vee (\neg((x_1 \wedge \neg x_2) \vee (\neg x_1 \wedge x_2)) \wedge ((x_3 \wedge \neg x_4) \vee (\neg x_3 \wedge x_4))) \end{aligned}$$

Note that, discounting NOTs, the depth of this circuit is 4



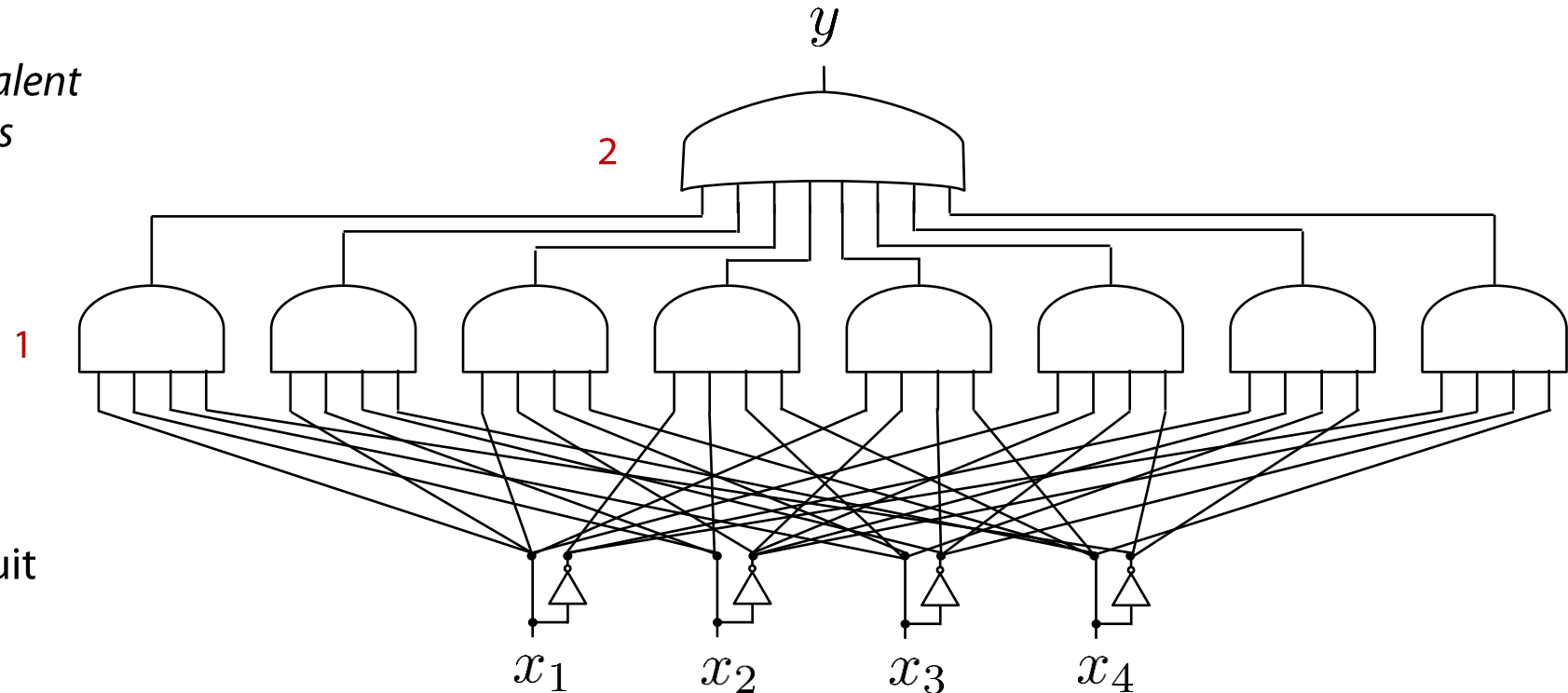
Parity Circuits

■ Disjunctive Normal Form (DNF)

Any logical formula
can be expressed
as an OR of ANDs
of the inputs
and their negations

$$\begin{aligned} & (x_1 \wedge x_2 \wedge x_3 \wedge \neg x_4) \vee (x_1 \wedge \neg x_2 \wedge x_3 \wedge x_4) \\ & \vee (x_1 \wedge x_2 \wedge \neg x_3 \wedge x_4) \vee (\neg x_1 \wedge x_2 \wedge x_3 \wedge x_4) \\ & \vee (x_1 \wedge \neg x_2 \wedge \neg x_3 \wedge \neg x_4) \vee (\neg x_1 \wedge x_2 \wedge \neg x_3 \wedge \neg x_4) \\ & \vee (\neg x_1 \wedge \neg x_2 \wedge x_3 \wedge \neg x_4) \vee (\neg x_1 \wedge \neg x_2 \wedge \neg x_3 \wedge x_4) \end{aligned}$$

*This circuit is equivalent
to the previous ones*



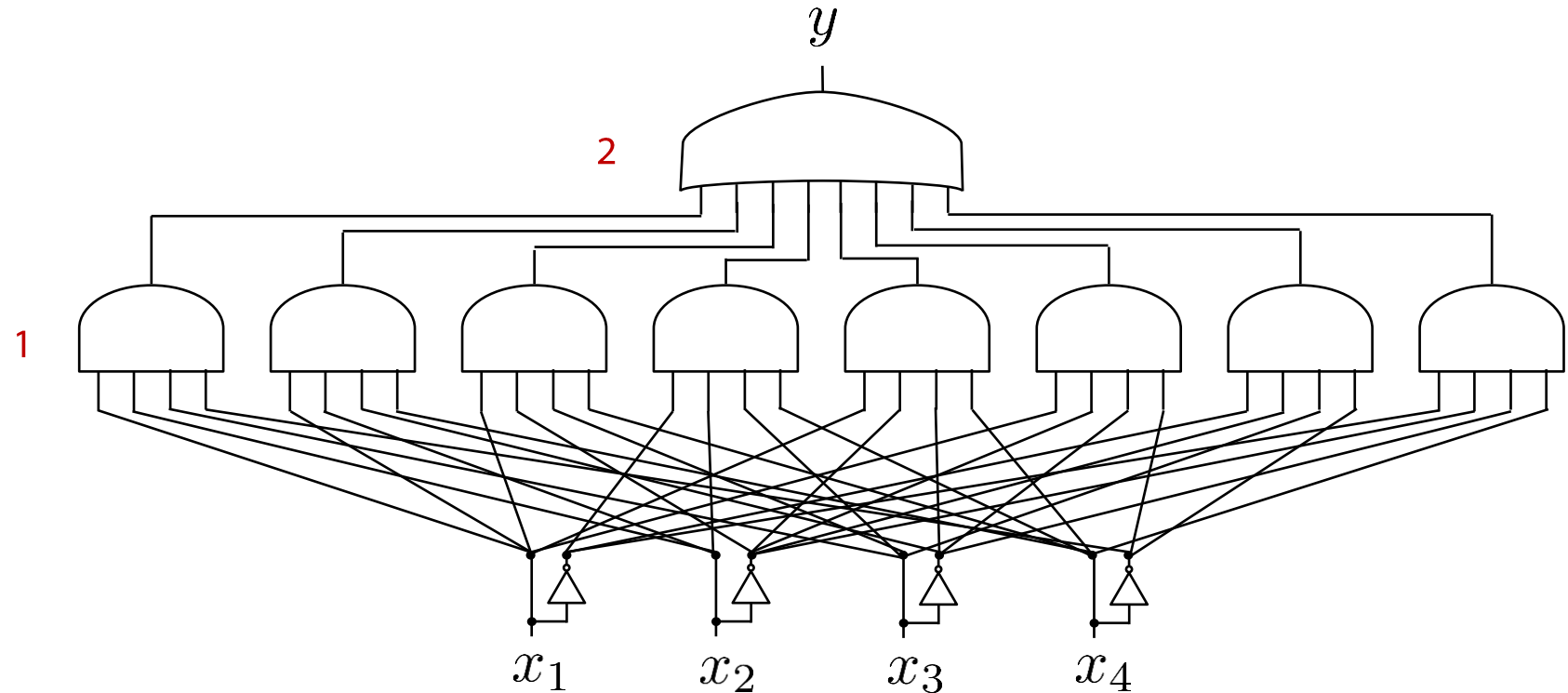
Note that this circuit
has depth 2

Parity Circuits

Any logical circuit can be re-implemented in *shallow* mode (i.e. with depth 2)

■ Question

Which way is better?
(*deep* vs. *shallow*)



Parity Circuits

Any logical circuit can be re-implemented in *shallow* mode (i.e., with depth 2)

- Lower Bound (Hastad, 1986)

For the implementation of *parity circuits*
the number of AND, OR components required is

$$\Omega \left(\exp \left(d^{\frac{1}{k-1}} \right) \right)$$

d is the number of bits in input
 k is the maximum depth allowed

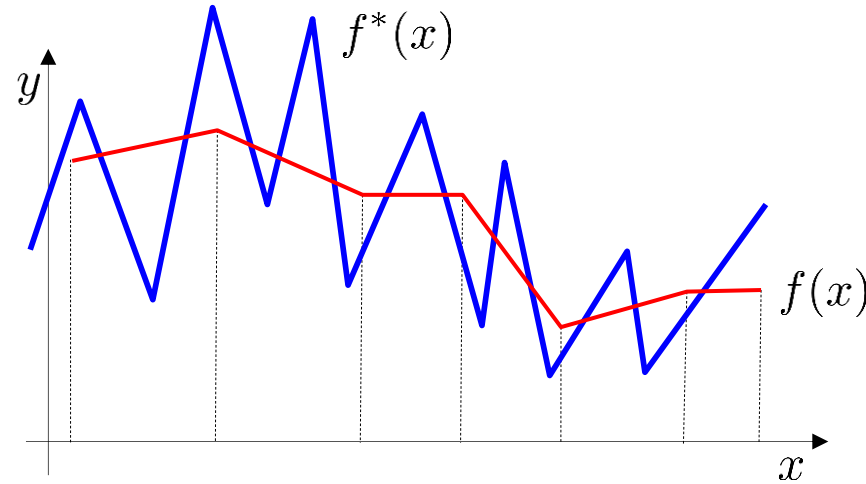
The above quantity becomes polynomial for

$$k = \frac{\log(d-1)}{\log \log(d-1) + \mathcal{O}(1)}$$

*In English: there exists a threshold $k_{\min}(d)$ beyond which
an exponential number of components w.r.t. d is no longer required*

Depth and piecewise linear functions

Example: a zig-zag target function:



Intuitively, the accuracy of the approximation depends on input space partitioning: unless we have a sufficient number of 'pieces' (i.e. regions in the partition) the approximation will be inaccurate

Assume we want to use a deep neural network with ReLU

$$\tilde{y} = \mathbf{w} \cdot \max(0, \mathbf{W}^{[k]} \dots \max(0, \mathbf{W}^{[1]}x + \mathbf{b}^{[1]}) \dots + \mathbf{b}^{[k]}) + b$$

Depth and piecewise linear functions

Construct two scalar functions using ReLU and parameters

$$\tilde{h}^{[k]} := \mathbf{w}^{[k]} \cdot \max(0, \mathbf{h}^{[k]} x)$$

$$\mathbf{h}^{[k]} := \max(0, \mathbf{W}^{[k]} x + \mathbf{b}^{[k]})$$

$$\mathbf{h}^{[1]} := [h_1^{[1]}, h_2^{[1]}]$$

$$h_1^{[1]} := \max(0, x)$$

$$h_2^{[1]} := \max(0, 2(x - 1))$$

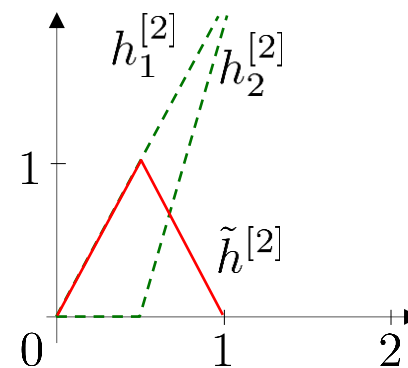
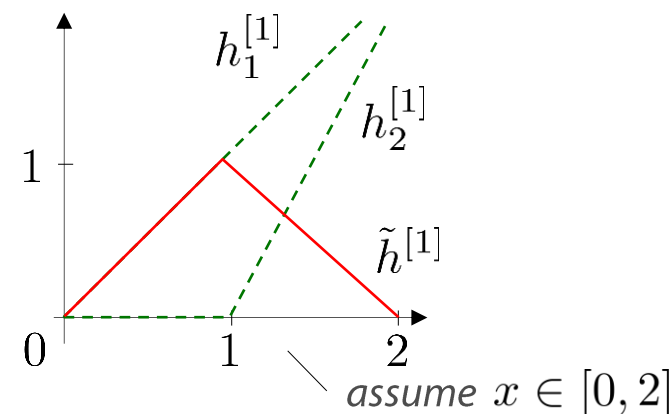
$$\tilde{h}^{[1]} := \max(0, x) - \max(0, 2(x - 1))$$

$$\mathbf{h}^{[2]} := [h_1^{[2]}, h_2^{[2]}]$$

$$h_1^{[2]} := \max(0, 2x)$$

$$h_2^{[2]} := \max(0, 4(x - 1/2))$$

$$\tilde{h}^{[2]} := \max(0, 2x) - \max(0, 4(x - 1/2))$$



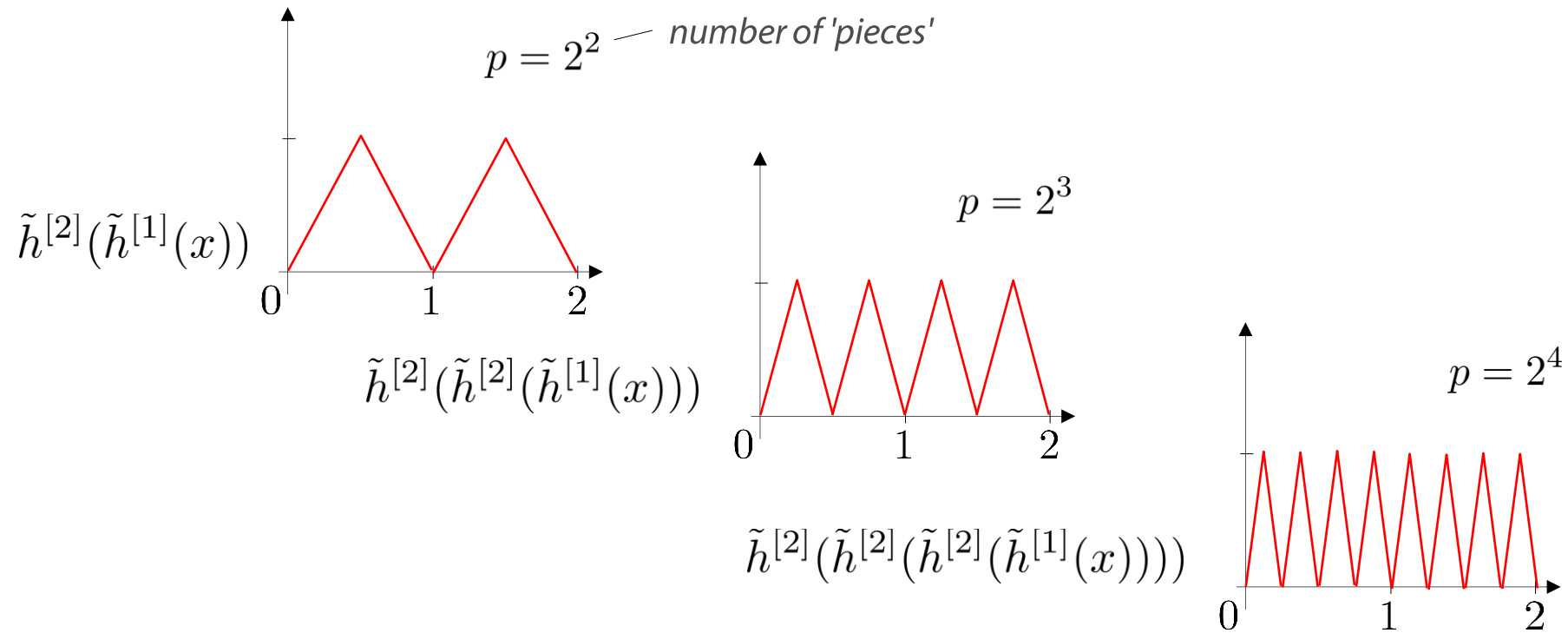
Depth and piecewise linear functions

Construct two scalar functions using ReLU plus parameters

$$\tilde{h}^{[k]} := \mathbf{w}^{[k]} \cdot \max(0, \mathbf{h}^{[k]} x)$$

$$\mathbf{h}^{[k]} := \max(0, \mathbf{W}^{[k]} x + \mathbf{b}^{[k]})$$

By nesting the two scalar functions:



Depth and piecewise linear functions

Deeper networks can make more 'pieces' with the same number of units

- **A lower bound that grows with depth** [Montufar et al. 2014]

For a network with one hidden layer of ReLU units of size h
the max number of pieces for the piecewise linear approximator is

$$p_{\max} = \sum_{i=0}^d \binom{h}{i} \leq h^d \text{ — input dimension}$$

For a network with k hidden *layers* of ReLU units, each of size h ,
the max number of such pieces is

$$p_{\max} = \mathcal{O}(2^k), \quad p_{\max} = \Omega \left(\left(\frac{h}{d} \right)^{(k-1)d} h^d \right)$$

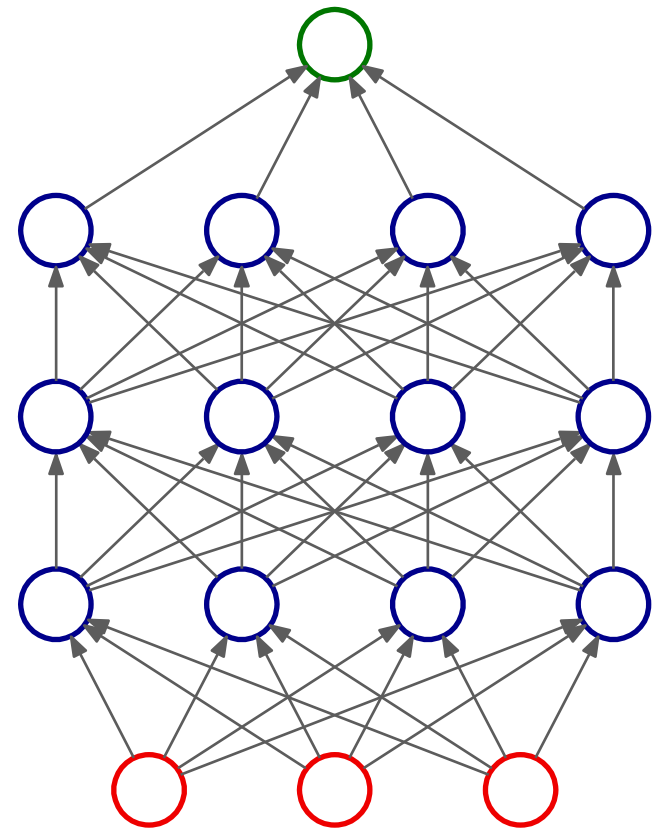
Moral: p_{\max} grows polynomially with layer size h but exponentially with depth k

Layerwise differentiation

Generalizing Deep Feed-Forward Neural Networks

- A feed-forward neural network with three hidden layers

$$\tilde{y} = w \cdot g(\mathbf{W}^{[3]}g(\mathbf{W}^{[2]}g(\mathbf{W}^{[1]}x + \mathbf{b}^{[1]}) + \mathbf{b}^{[2]}) + \mathbf{b}^{[3]}) + b$$



Generalizing Deep Feed-Forward Neural Networks

- A feed-forward neural network with three hidden layers

$$\tilde{y} = \boldsymbol{w} \cdot g(\boldsymbol{W}^{[3]}g(\boldsymbol{W}^{[2]}g(\boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}) + \boldsymbol{b}^{[2]}) + \boldsymbol{b}^{[3]}) + b$$

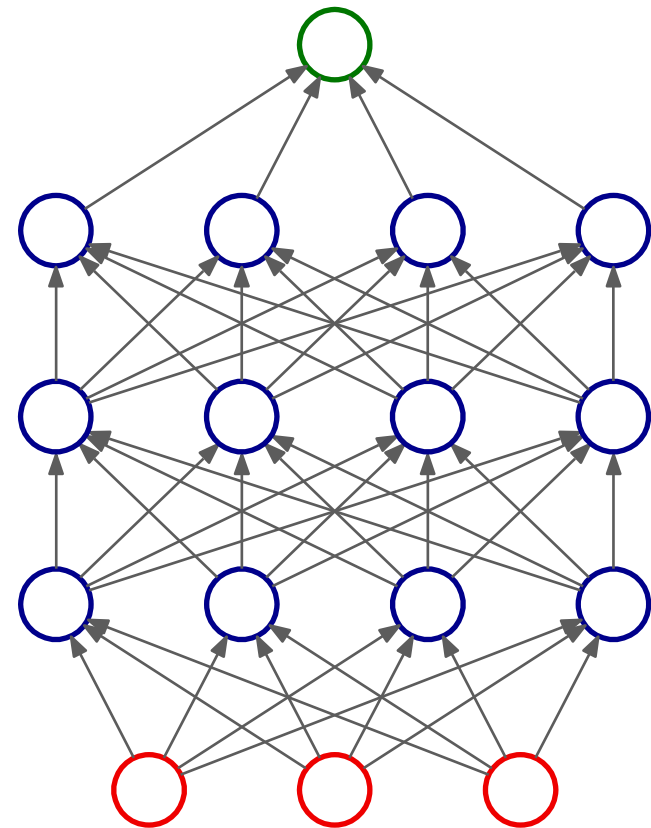
$$\tilde{y} := \boldsymbol{w} \cdot \boldsymbol{h}^{[3]} + b$$

$$\boldsymbol{h}^{[3]} := g(\boldsymbol{W}^{[3]}\boldsymbol{h}^{[2]} + \boldsymbol{b}^{[3]})$$

$$\boldsymbol{h}^{[2]} := g(\boldsymbol{W}^{[2]}\boldsymbol{h}^{[1]} + \boldsymbol{b}^{[2]})$$

$$\boldsymbol{h}^{[1]} := g(\boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]})$$

\boldsymbol{x}



Generalizing Deep Feed-Forward Neural Networks

- A feed-forward neural network with three hidden layers

$$\tilde{y} = \boldsymbol{w} \cdot g(\boldsymbol{W}^{[3]}g(\boldsymbol{W}^{[2]}g(\boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]}) + \boldsymbol{b}^{[2]}) + \boldsymbol{b}^{[3]}) + b$$

$$\tilde{y}(\boldsymbol{h}^{[3]}, \boldsymbol{\vartheta}^{[\tilde{y}]})$$

$$\tilde{y} := \boldsymbol{w} \cdot \boldsymbol{h}^{[3]} + b$$

$$\boldsymbol{h}^{[3]}(\boldsymbol{h}^{[2]}, \boldsymbol{\vartheta}^{[3]})$$

$$\boldsymbol{h}^{[3]} := g(\boldsymbol{W}^{[3]}\boldsymbol{h}^{[2]} + \boldsymbol{b}^{[3]})$$

$$\boldsymbol{h}^{[2]}(\boldsymbol{h}^{[1]}, \boldsymbol{\vartheta}^{[2]})$$

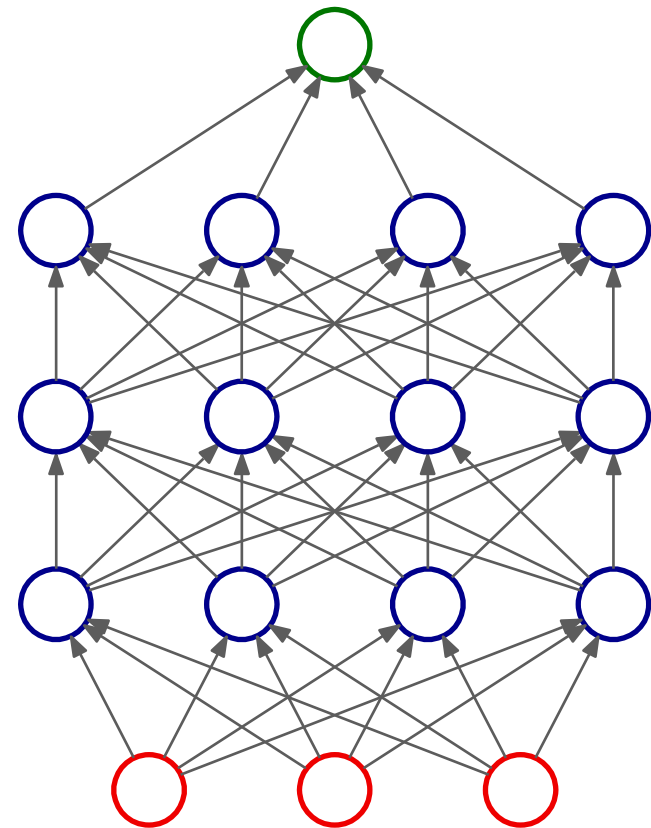
$$\boldsymbol{h}^{[2]} := g(\boldsymbol{W}^{[2]}\boldsymbol{h}^{[1]} + \boldsymbol{b}^{[2]})$$

$$\boldsymbol{h}^{[1]}(\boldsymbol{x}, \boldsymbol{\vartheta}^{[1]})$$

$$\boldsymbol{h}^{[1]} := g(\boldsymbol{W}^{[1]}\boldsymbol{x} + \boldsymbol{b}^{[1]})$$

\boldsymbol{x}

\boldsymbol{x}



Generalizing Deep Feed-Forward Neural Networks

- A feed-forward neural network with three hidden layers

$$L(\tilde{y}, y) = (\tilde{y} - y)^2$$

$$\tilde{y}(h^{[3]}, \vartheta^{[\tilde{y}]})$$

$$h^{[3]}(h^{[2]}, \vartheta^{[3]})$$

$$h^{[2]}(h^{[1]}, \vartheta^{[2]})$$

$$h^{[1]}(x, \vartheta^{[1]})$$

$$x$$

Generalizing Deep Feed-Forward Neural Networks

- **Computing gradient** (layerwise)

$$L(\tilde{y}, y) = (\tilde{y} - y)^2 \quad \frac{\partial}{\partial \boldsymbol{\vartheta}[\tilde{y}]} (\tilde{y} - y)^2 = 2(\tilde{y} - y) \frac{\partial \tilde{y}}{\partial \boldsymbol{\vartheta}[\tilde{y}]}$$

$$\tilde{y}(\boldsymbol{h}^{[3]}, \boldsymbol{\vartheta}^{[\tilde{y}]}) \quad \frac{\partial \tilde{y}}{\partial \boldsymbol{\vartheta}^{[\tilde{y}]}}$$

$$\boldsymbol{h}^{[3]}(\boldsymbol{h}^{[2]}, \boldsymbol{\vartheta}^{[3]})$$

$$\boldsymbol{h}^{[2]}(\boldsymbol{h}^{[1]}, \boldsymbol{\vartheta}^{[2]})$$

$$\boldsymbol{h}^{[1]}(\boldsymbol{x}, \boldsymbol{\vartheta}^{[1]})$$

$$\boldsymbol{x}$$

Generalizing Deep Feed-Forward Neural Networks

- **Computing gradient** (layerwise)

$$L(\tilde{y}, y) = (\tilde{y} - y)^2 \quad \frac{\partial}{\partial \boldsymbol{\vartheta}^{[3]}} (\tilde{y} - y)^2 = 2(\tilde{y} - y) \frac{\partial \tilde{y}}{\partial \boldsymbol{\vartheta}^{[3]}}$$

$$\tilde{y}(\boldsymbol{h}^{[3]}, \boldsymbol{\vartheta}^{[3]}) \quad \frac{\partial \tilde{y}}{\partial \boldsymbol{\vartheta}^{[3]}} = \frac{\partial \tilde{y}}{\partial \boldsymbol{h}^{[3]}} \frac{\partial \boldsymbol{h}^{[3]}}{\partial \boldsymbol{\vartheta}^{[3]}}$$

$$\boldsymbol{h}^{[3]}(\boldsymbol{h}^{[2]}, \boldsymbol{\vartheta}^{[3]}) \quad \frac{\partial \boldsymbol{h}^{[3]}}{\partial \boldsymbol{\vartheta}^{[3]}}$$

$$\boldsymbol{h}^{[2]}(\boldsymbol{h}^{[1]}, \boldsymbol{\vartheta}^{[2]})$$

$$\boldsymbol{h}^{[1]}(\boldsymbol{x}, \boldsymbol{\vartheta}^{[1]})$$

$$\boldsymbol{x}$$

Generalizing Deep Feed-Forward Neural Networks

- **Computing gradient** (layerwise)

$$L(\tilde{y}, y) = (\tilde{y} - y)^2 \quad \frac{\partial}{\partial \boldsymbol{\vartheta}^{[2]}} (\tilde{y} - y)^2 = 2(\tilde{y} - y) \frac{\partial \tilde{y}}{\partial \boldsymbol{\vartheta}^{[2]}}$$

$$\tilde{y}(\boldsymbol{h}^{[3]}, \boldsymbol{\vartheta}^{[\tilde{y}]}) \quad \frac{\partial \tilde{y}}{\partial \boldsymbol{\vartheta}^{[2]}} = \frac{\partial \tilde{y}}{\partial \boldsymbol{h}^{[3]}} \frac{\partial \boldsymbol{h}^{[3]}}{\partial \boldsymbol{\vartheta}^{[2]}}$$

$$\boldsymbol{h}^{[3]}(\boldsymbol{h}^{[2]}, \boldsymbol{\vartheta}^{[3]}) \quad \frac{\partial \boldsymbol{h}^{[3]}}{\partial \boldsymbol{\vartheta}^{[2]}} = \frac{\partial \boldsymbol{h}^{[3]}}{\partial \boldsymbol{h}^{[2]}} \frac{\partial \boldsymbol{h}^{[2]}}{\partial \boldsymbol{\vartheta}^{[2]}}$$

$$\boldsymbol{h}^{[2]}(\boldsymbol{h}^{[1]}, \boldsymbol{\vartheta}^{[2]}) \quad \frac{\partial \boldsymbol{h}^{[2]}}{\partial \boldsymbol{\vartheta}^{[2]}}$$

$$\boldsymbol{h}^{[1]}(\boldsymbol{x}, \boldsymbol{\vartheta}^{[1]})$$

\boldsymbol{x}

Generalizing Deep Feed-Forward Neural Networks

- **Computing gradient** (layerwise)

$$L(\tilde{y}, y) = (\tilde{y} - y)^2 \quad \frac{\partial}{\partial \boldsymbol{\vartheta}^{[j]}} (\tilde{y} - y)^2 = 2(\tilde{y} - y) \frac{\partial \tilde{y}}{\partial \boldsymbol{\vartheta}^{[j]}}$$

...

...

$$\boldsymbol{h}^{[i]}(\boldsymbol{h}^{[i-1]}, \boldsymbol{\vartheta}^{[i]}) \quad \frac{\partial \boldsymbol{h}^{[i]}}{\partial \boldsymbol{\vartheta}^{[i]}}, \quad j = i$$

$$\frac{\partial \boldsymbol{h}^{[i]}}{\partial \boldsymbol{\vartheta}^{[j]}} = \frac{\partial \boldsymbol{h}^{[i]}}{\partial \boldsymbol{h}^{[i-1]}} \frac{\partial \boldsymbol{h}^{[i-1]}}{\partial \boldsymbol{\vartheta}^{[j]}}, \quad j < i$$

...

...

Generalizing Deep Feed-Forward Neural Networks

- **Computing gradient** (layerwise)

$$L(\tilde{y}, y) = (\tilde{y} - y)^2 \quad \frac{\partial}{\partial \vartheta[j]} (\tilde{y} - y)^2 = 2(\tilde{y} - y) \frac{\partial \tilde{y}}{\partial \vartheta[j]}$$

...

...

$$\mathbf{h}^{[i]}(\mathbf{h}^{[i-1]}, \boldsymbol{\vartheta}^{[i]})$$

$$\boxed{\frac{\partial \mathbf{h}^{[i]}}{\partial \boldsymbol{\vartheta}^{[i]}}, j = i}$$

Each layer
'needs to know'
just these two derivatives

$$\frac{\partial \mathbf{h}^{[i]}}{\partial \boldsymbol{\vartheta}^{[j]}} = \boxed{\frac{\partial \mathbf{h}^{[i]}}{\partial \mathbf{h}^{[i-1]}}} \frac{\partial \mathbf{h}^{[i-1]}}{\partial \boldsymbol{\vartheta}^{[j]}}, j < i$$

...

...

Function approximation (a.k.a. regression) vs. classification

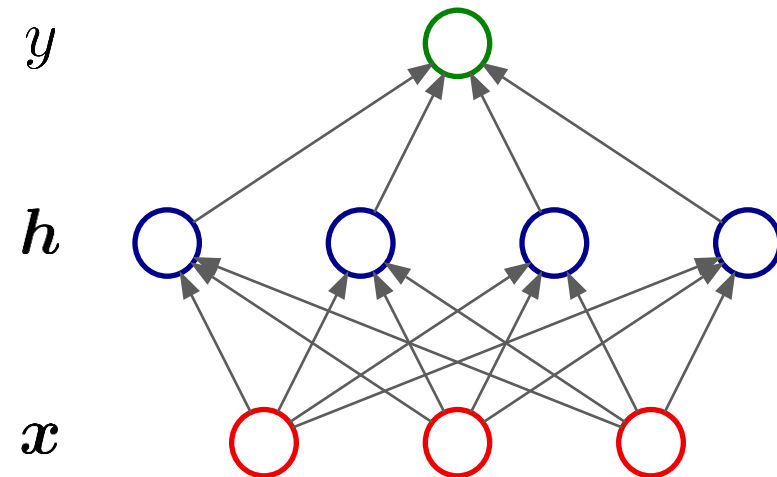
Classification: Softmax

- **Function approximation** (a.k.a. regression)

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

Feed-forward neural network

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}) + b$$



Classification: Softmax

■ Classification

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d, \quad y \in \{\text{class}_i\}_{i=1}^k$$

Feed-forward neural network with a *Softmax* layer

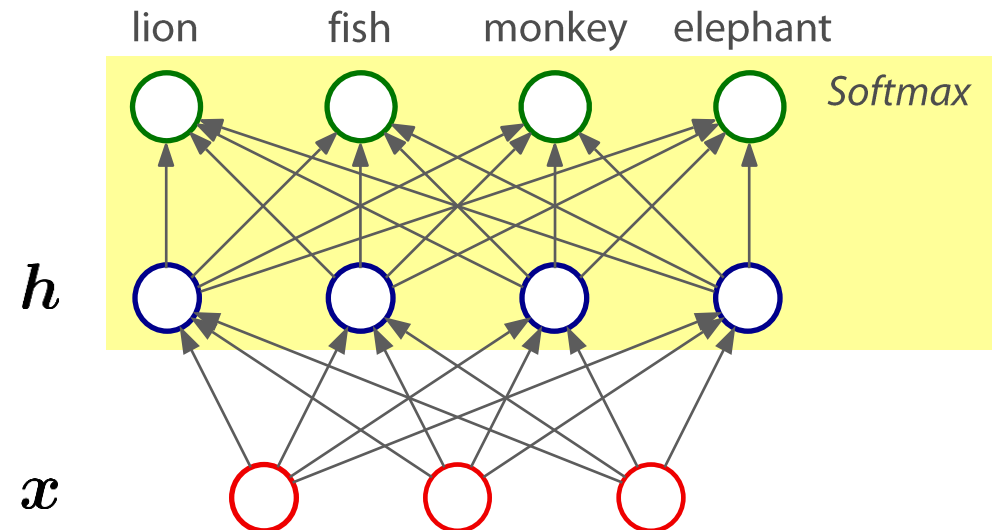
$$P(\tilde{y} = \text{class}_i \mid \mathbf{x}) := \frac{\exp(\mathbf{w}_i \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}))}{\sum_{j=1}^k \exp(\mathbf{w}_j \cdot g(\mathbf{W}\mathbf{x} + \mathbf{b}))}$$

From now on

$$P(\tilde{y} = \text{class}_i \mid \mathbf{x})$$

will be written as

$$P(\tilde{y} = i \mid \mathbf{x})$$



Classification: Softmax

■ Classification

$$y = f^*(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d, \quad y \in \{\text{class}_i\}_{i=1}^k$$

The *Softmax* layer can be rewritten as:

$$P(\tilde{y} = \text{class}_i \mid \mathbf{h}) := \frac{\exp(\mathbf{w}_i \cdot \mathbf{h} + b_i)}{\sum_{j=1}^k \exp(\mathbf{w}_j \cdot \mathbf{h} + b_j)}$$

where, in this case: $\mathbf{h} := g(\mathbf{W}\mathbf{x} + \mathbf{b})$

(yet, more in general, \mathbf{h} can be anything)

Classification: Softmax

- ***Softmax as a layer***

The entire *Softmax* layer can be rewritten as:

$$P((\tilde{y} = i)_1^k \mid \mathbf{h}) := \frac{\exp(\mathbf{W}_S \mathbf{h} + \mathbf{b}_S)}{\sum \exp(\mathbf{W}_S \mathbf{h} + \mathbf{b}_S)}$$

Probability distribution
(a vector)

Sum of all components

where:

$$\mathbf{W}_S := \begin{bmatrix} - & \mathbf{w}_1 & - \\ & \vdots & \\ - & \mathbf{w}_k & - \end{bmatrix} \quad \mathbf{b}_S := \begin{bmatrix} b_1 \\ \vdots \\ b_k \end{bmatrix}$$

The vector $\mathbf{W}_S \mathbf{h} + \mathbf{b}_S$ is sometimes referred to as the **logit**

Classification: Softmax

■ Cross-entropy in general

P and Q are probability distributions on a discrete random variable $y \in \{1, \dots, k\}$

$$H(Q, P) := - \sum_{j=1}^k Q(y = j) \log P(\tilde{y} = j)$$

■ As a loss function for Softmax

Q in this case is the 'true' classification, i.e. the one in the dataset

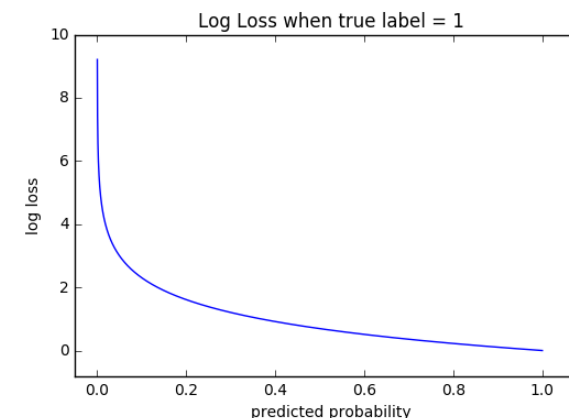
$$Q(y = j) := \delta(y = j) \quad \text{--- Kronecker delta}$$

while P is the output of the Softmax layer

$$P(\tilde{y} = j | \mathbf{h})$$

Hence, the loss is:

$$\begin{aligned} L(\mathbf{h}^{(i)}, y^{(i)}) &:= - \sum_{j=1}^k \delta(y^{(i)} = j) \log P(\tilde{y} = j | \mathbf{h}^{(i)}) \\ &= - \log P(\tilde{y} = y^{(i)} | \mathbf{h}^{(i)}) \end{aligned}$$



Classification: Softmax

■ Cross-entropy for Softmax

$$L(\mathbf{h}^{(i)}, y^{(i)}) := - \sum_{j=1}^k \delta(y^{(i)} = j) \log P(\tilde{y} = j | \mathbf{h}^{(i)})$$

Expressing the loss function in vector form:

$$\mathbf{y} := \begin{bmatrix} y_1 \\ \vdots \\ y_k \end{bmatrix}, \quad y_j := \delta(y = j) \quad \text{'one hot' representation} \quad \mathbf{p} := \begin{bmatrix} p_1 \\ \vdots \\ p_k \end{bmatrix}, \quad p_j := P(\tilde{y} = j | \mathbf{h})$$

$$L(\mathbf{h}^{(i)}, \mathbf{y}^{(i)}) = - \mathbf{y}^{(i)} \cdot \log(\mathbf{p}^{(i)})$$

which implies that also the dataset has to be transformed in the 'one hot' representation

$$D := \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N \quad \implies \quad D := \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$$

Classification: Softmax

■ Gradient of **Softmax** (layerwise)

$$L(D) = \sum_{i=1}^N L(\mathbf{h}^{(i)}, \mathbf{y}^{(i)}) = - \sum_{i=1}^N \mathbf{y}^{(i)} \cdot \log(\mathbf{p}^{(i)})$$

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(D) &= \frac{\partial}{\partial \boldsymbol{\vartheta}} \sum_{i=1}^N L(\mathbf{h}^{(i)}, \mathbf{y}^{(i)}) = - \frac{\partial}{\partial \boldsymbol{\vartheta}} \sum_{i=1}^N \mathbf{y}^{(i)} \cdot \log(\mathbf{p}^{(i)}) \\ &= - \sum_{i=1}^N \mathbf{y}^{(i)} \cdot \frac{\partial}{\partial \boldsymbol{\vartheta}} \log(\mathbf{p}^{(i)}) \end{aligned}$$

This is a matrix

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} \log(\mathbf{p}) = \begin{bmatrix} \frac{\partial}{\partial \vartheta_1} \log(p_1) & \dots & \frac{\partial}{\partial \vartheta_d} \log(p_1) \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial \vartheta_1} \log(p_k) & \dots & \frac{\partial}{\partial \vartheta_d} \log(p_k) \end{bmatrix} = \begin{bmatrix} - & \frac{\partial}{\partial \boldsymbol{\vartheta}} \log(p_1) & - \\ & \vdots & \\ - & \frac{\partial}{\partial \boldsymbol{\vartheta}} \log(p_k) & - \end{bmatrix}$$

Classification: Softmax

■ Gradient of **Softmax** (layerwise)

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\vartheta}} \log(p_j) &= \frac{\partial}{\partial \boldsymbol{\vartheta}} \log P(\tilde{y} = j \mid \mathbf{h}) \\&= \frac{\partial}{\partial \boldsymbol{\vartheta}} \log \frac{\exp(\mathbf{w}_j \cdot \mathbf{h} + b_j)}{\sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)} \\&= \frac{\partial}{\partial \boldsymbol{\vartheta}} \left(\log \exp(\mathbf{w}_j \cdot \mathbf{h} + b_j) - \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l) \right) \\&= \frac{\partial}{\partial \boldsymbol{\vartheta}} (\mathbf{w}_j \cdot \mathbf{h} + b_j) - \frac{\partial}{\partial \boldsymbol{\vartheta}} \log \sum_{l=1}^k \exp(\mathbf{w}_l \cdot \mathbf{h} + b_l)\end{aligned}$$

Classification: Softmax

■ Gradient of *Softmax* (layerwise)

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} \log(p_j) = \frac{\partial}{\partial \boldsymbol{\vartheta}} (\boldsymbol{w}_j \cdot \boldsymbol{h} + b_j) - \frac{\partial}{\partial \boldsymbol{\vartheta}} \log \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l)$$

Case 1: $\boldsymbol{\vartheta} = \boldsymbol{w}_r$ or $\boldsymbol{\vartheta} = b_r$

Case 2: $\boldsymbol{h}(\boldsymbol{\vartheta})$ i.e. $\boldsymbol{\vartheta}$ is a generic parameter on which \boldsymbol{h} depends

$$\frac{\partial \boldsymbol{h}^{[i]}}{\partial \boldsymbol{\vartheta}^{[i]}} \\ \frac{\partial \boldsymbol{h}^{[i]}}{\partial \boldsymbol{\vartheta}^{[j]}}, \quad j < i$$

Let's compute the two contributions separately

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} (\boldsymbol{w}_j \cdot \boldsymbol{h} + b_j) \\ \frac{\partial}{\partial \boldsymbol{\vartheta}} \log \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l)$$

Classification: Softmax

■ Gradient of **Softmax** (layerwise)

$$\frac{\partial}{\partial \vartheta} (\mathbf{w}_j \cdot \mathbf{h} + b_j)$$

Case 1: $\vartheta = w_r$ or $\vartheta = b_r$

$$\frac{\partial}{\partial w_r} (\mathbf{w}_j \cdot \mathbf{h} + b_j) = \begin{cases} 0 & \text{if } r \neq j \\ \mathbf{h} & \text{otherwise} \end{cases}$$

$$\frac{\partial}{\partial b_r} (\mathbf{w}_j \cdot \mathbf{h} + b_j) = \begin{cases} 0 & \text{if } r \neq j \\ 1 & \text{otherwise} \end{cases}$$

Case 2: $\mathbf{h}(\vartheta)$ i.e. ϑ is a generic parameter on which \mathbf{h} depends

$$\frac{\partial}{\partial \vartheta} (\mathbf{w}_j \cdot \mathbf{h} + b_j) = \mathbf{w}_j \cdot \frac{\partial}{\partial \vartheta} \mathbf{h}$$

Classification: Softmax

■ Gradient of **Softmax** (layerwise)

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} \log \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l)$$

Case 1: $\boldsymbol{\vartheta} = \boldsymbol{w}_r$ or $\boldsymbol{\vartheta} = b_r$

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{w}_r} \log \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) &= \\ &= \frac{1}{\sum_{m=1}^k \exp(\boldsymbol{w}_m \cdot \boldsymbol{h} + b_m)} \frac{\partial}{\partial \boldsymbol{w}_r} \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) \\ &= \frac{1}{\sum_{m=1}^k \exp(\boldsymbol{w}_m \cdot \boldsymbol{h} + b_m)} \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) \frac{\partial}{\partial \boldsymbol{w}_r} (\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) \\ &= \frac{\exp(\boldsymbol{w}_r \cdot \boldsymbol{h} + b_r)}{\sum_{m=1}^k \exp(\boldsymbol{w}_m \cdot \boldsymbol{h} + b_m)} \boldsymbol{h} = p_r \boldsymbol{h} \end{aligned}$$

Classification: Softmax

■ Gradient of **Softmax** (layerwise)

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} \log \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l)$$

Case 1: $\boldsymbol{\vartheta} = \boldsymbol{w}_r$ or $\boldsymbol{\vartheta} = b_r$

$$\begin{aligned} \frac{\partial}{\partial b_r} \log \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) &= \\ &= \frac{1}{\sum_{m=1}^k \exp(\boldsymbol{w}_m \cdot \boldsymbol{h} + b_m)} \frac{\partial}{\partial b_r} \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) \\ &= \frac{1}{\sum_{m=1}^k \exp(\boldsymbol{w}_m \cdot \boldsymbol{h} + b_m)} \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) \frac{\partial}{\partial b_r} (\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) \\ &= \frac{\exp(\boldsymbol{w}_r \cdot \boldsymbol{h} + b_r)}{\sum_{m=1}^k \exp(\boldsymbol{w}_m \cdot \boldsymbol{h} + b_m)} = p_r \end{aligned}$$

Classification: Softmax

■ Gradient of **Softmax** (layerwise)

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} \log \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l)$$

Case 2: $\boldsymbol{h}(\boldsymbol{\vartheta})$ i.e. $\boldsymbol{\vartheta}$ is a generic parameter on which \boldsymbol{h} depends

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\vartheta}} \log \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) &= \\ &= \frac{1}{\sum_{m=1}^k \exp(\boldsymbol{w}_m \cdot \boldsymbol{h} + b_m)} \frac{\partial}{\partial \boldsymbol{\vartheta}} \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) \\ &= \frac{1}{\sum_{m=1}^k \exp(\boldsymbol{w}_m \cdot \boldsymbol{h} + b_m)} \sum_{l=1}^k \exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) \frac{\partial}{\partial \boldsymbol{\vartheta}} (\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l) \\ &= \sum_{l=1}^k \frac{\exp(\boldsymbol{w}_l \cdot \boldsymbol{h} + b_l)}{\sum_{m=1}^k \exp(\boldsymbol{w}_m \cdot \boldsymbol{h} + b_m)} \boldsymbol{w}_l^T \frac{\partial}{\partial \boldsymbol{\vartheta}} \boldsymbol{h} = \left(\sum_{l=1}^k p_l \boldsymbol{w}_l^T \right) \frac{\partial}{\partial \boldsymbol{\vartheta}} \boldsymbol{h} \end{aligned}$$