

Chapter ML:IV (continued)

IV. Neural Networks

- Perceptron Learning
- Multilayer Perceptron Basics
- Multilayer Perceptron with Two Layers
- Multilayer Perceptron at Arbitrary Depth
- Advanced MLPs
- Automatic Gradient Computation

Multilayer Perceptron Basics

Definition 1 (Linear Separability)

Two sets of feature vectors, X_0, X_1 , sampled from a p -dimensional feature space \mathbf{X} , are called linearly separable if $p+1$ real numbers, w_0, w_1, \dots, w_p , exist such that the following conditions holds:

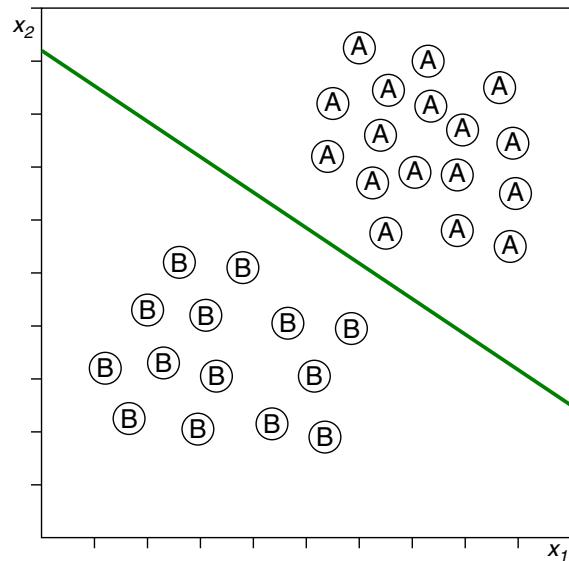
1. $\forall \mathbf{x} \in X_0: \sum_{j=0}^p w_j x_j < 0$
2. $\forall \mathbf{x} \in X_1: \sum_{j=0}^p w_j x_j \geq 0$

Multilayer Perceptron Basics

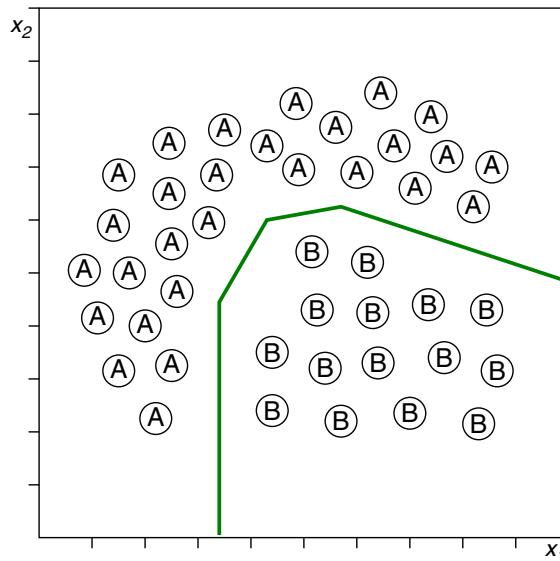
Definition 1 (Linear Separability)

Two sets of feature vectors, X_0, X_1 , sampled from a p -dimensional feature space \mathbf{X} , are called linearly separable if $p+1$ real numbers, w_0, w_1, \dots, w_p , exist such that the following conditions holds:

1. $\forall \mathbf{x} \in X_0: \sum_{j=0}^p w_j x_j < 0$
2. $\forall \mathbf{x} \in X_1: \sum_{j=0}^p w_j x_j \geq 0$



linearly separable



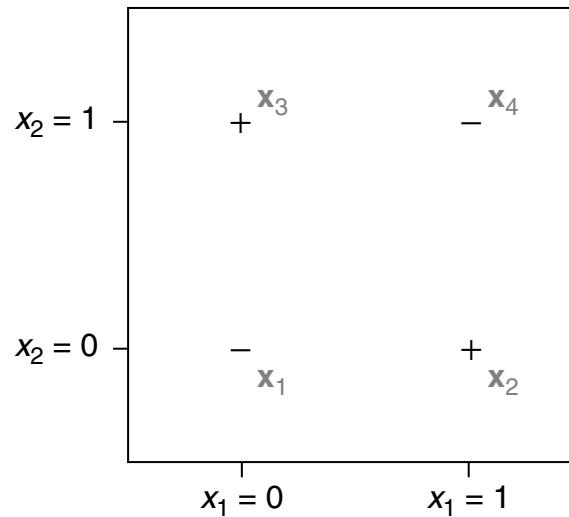
not linearly separable

Multilayer Perceptron Basics

Linear Separability (continued)

The *XOR* function defines the smallest example for two not linearly separable sets:

	x_1	x_2	XOR	c
\mathbf{x}_1	0	0	0	-
\mathbf{x}_2	1	0	1	+
\mathbf{x}_3	0	1	1	+
\mathbf{x}_4	1	1	0	-

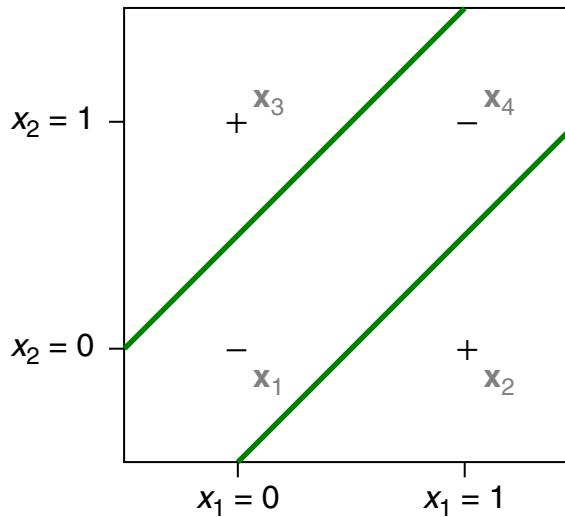


Multilayer Perceptron Basics

Linear Separability (continued)

The *XOR* function defines the smallest example for two not linearly separable sets:

	x_1	x_2	XOR	c
\mathbf{x}_1	0	0	0	-
\mathbf{x}_2	1	0	1	+
\mathbf{x}_3	0	1	1	+
\mathbf{x}_4	1	1	0	-

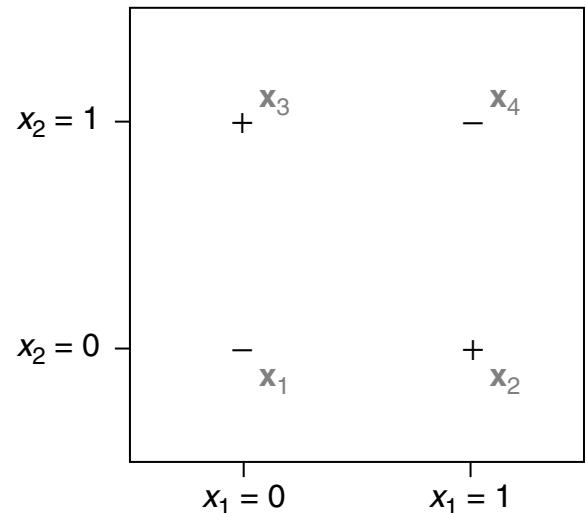
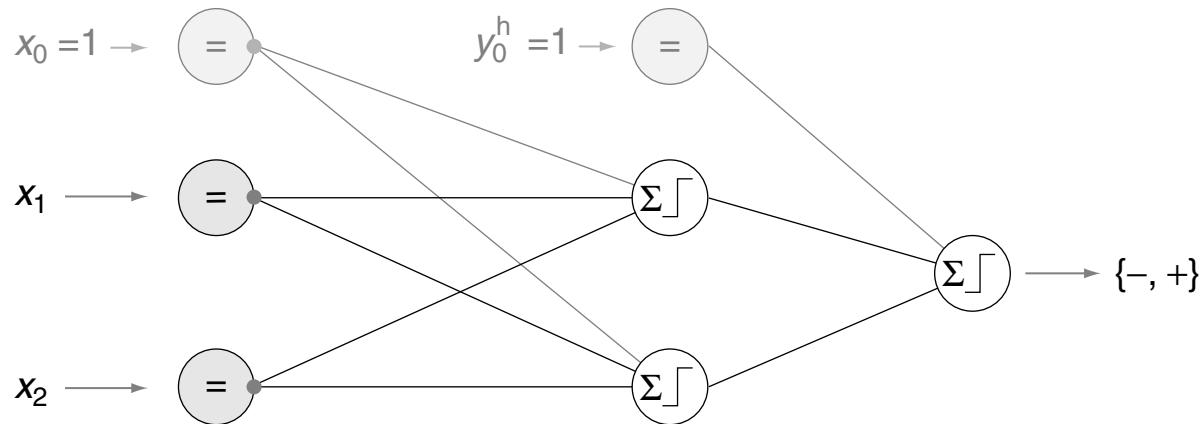


- Specification of several hyperplanes.
- Layered combination of several perceptrons: the multilayer perceptron.

Multilayer Perceptron Basics

Overcoming the Linear Separability Restriction

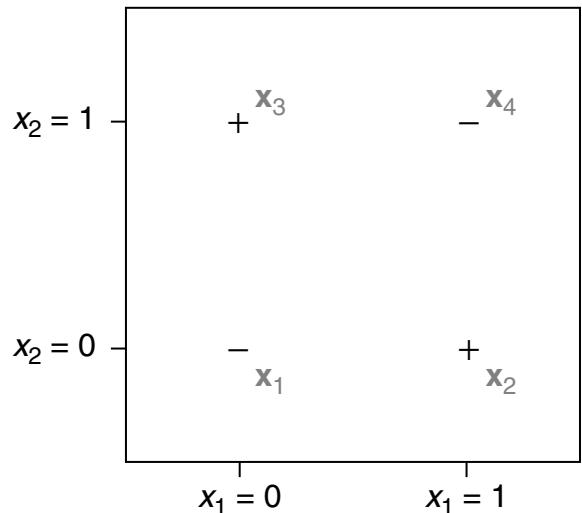
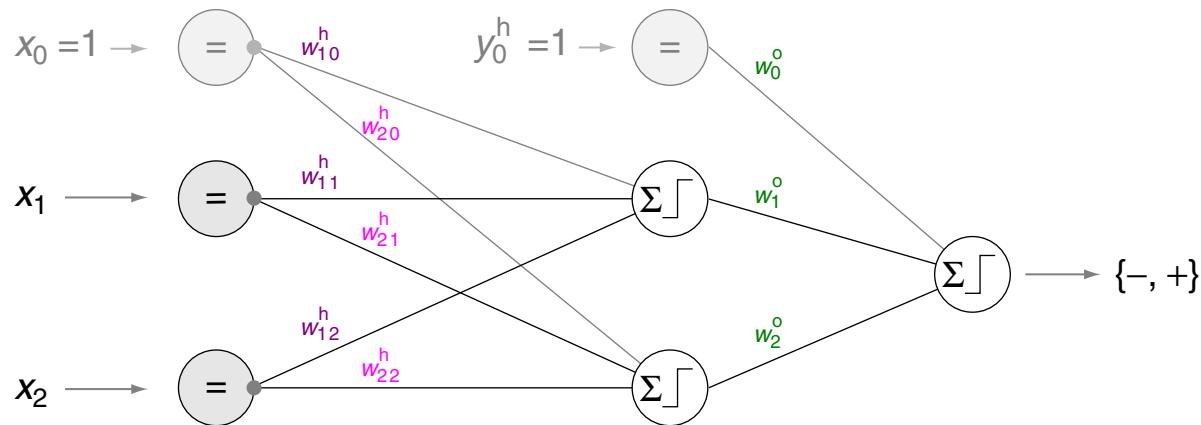
A minimum multilayer perceptron $y(\mathbf{x})$ that can handle the *XOR* problem:



Multilayer Perceptron Basics

Overcoming the Linear Separability Restriction (continued)

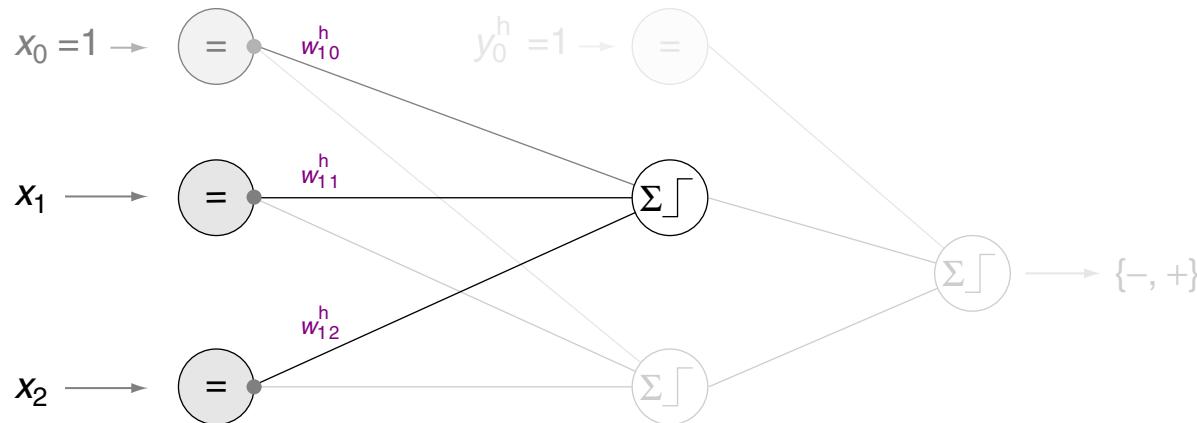
A minimum multilayer perceptron $y(\mathbf{x})$ that can handle the *XOR* problem:



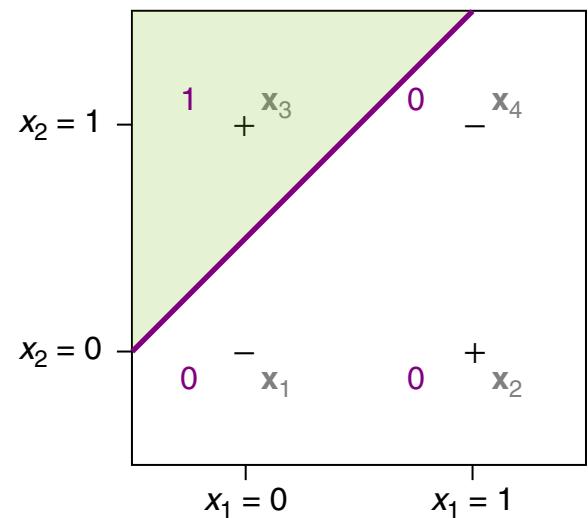
Multilayer Perceptron Basics

Overcoming the Linear Separability Restriction (continued)

A minimum multilayer perceptron $y(\mathbf{x})$ that can handle the *XOR* problem:



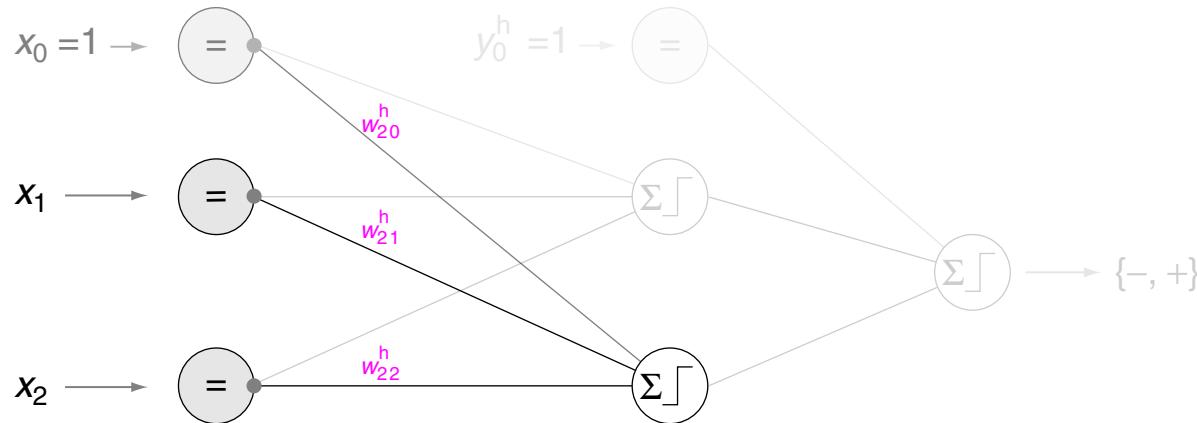
$$W^h = \begin{bmatrix} -0.5 & -1 & 1 \\ 0.5 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$



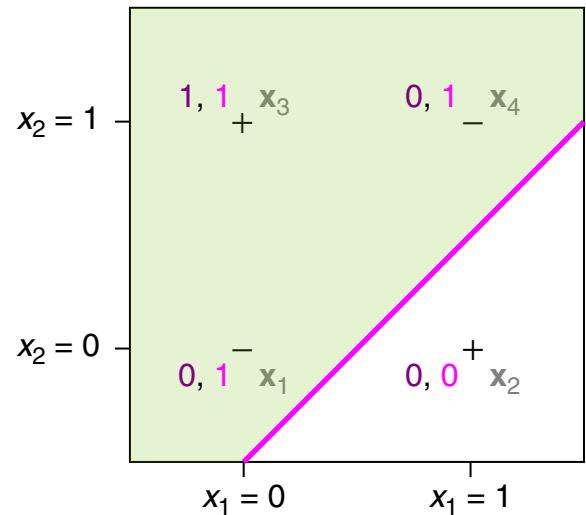
Multilayer Perceptron Basics

Overcoming the Linear Separability Restriction (continued)

A minimum multilayer perceptron $y(\mathbf{x})$ that can handle the *XOR* problem:



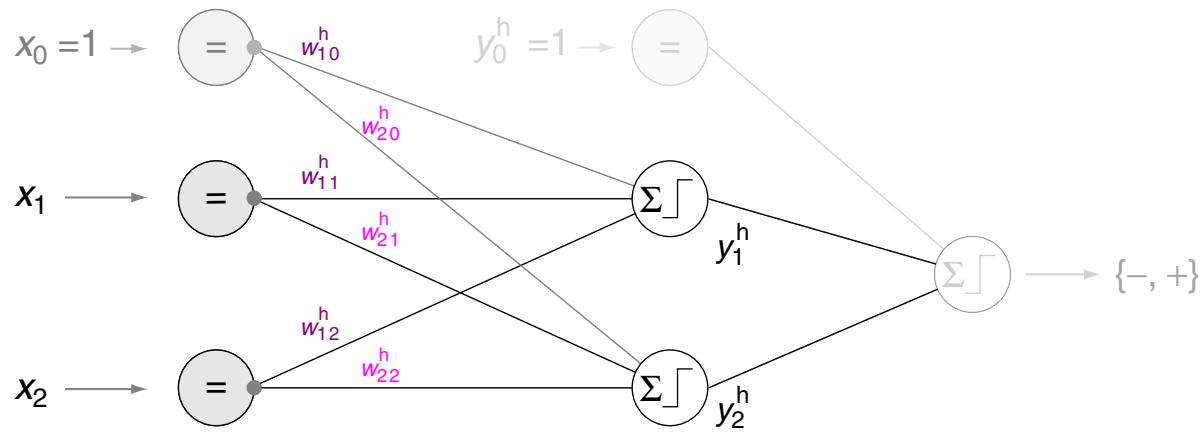
$$W^h = \begin{bmatrix} -0.5 & -1 & 1 \\ 0.5 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$



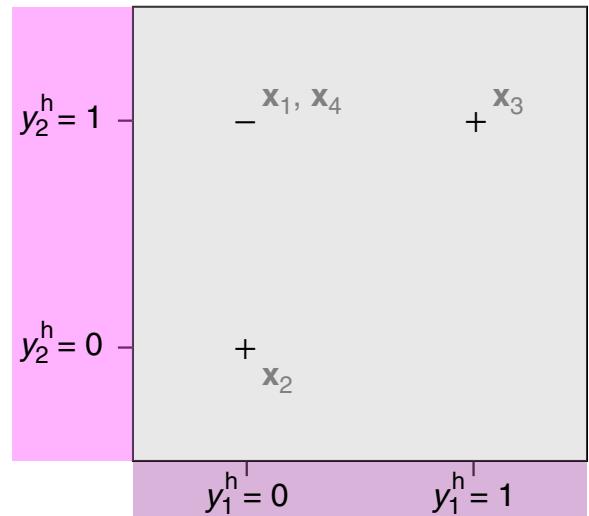
Multilayer Perceptron Basics

Overcoming the Linear Separability Restriction (continued)

A minimum multilayer perceptron $y(\mathbf{x})$ that can handle the *XOR* problem:



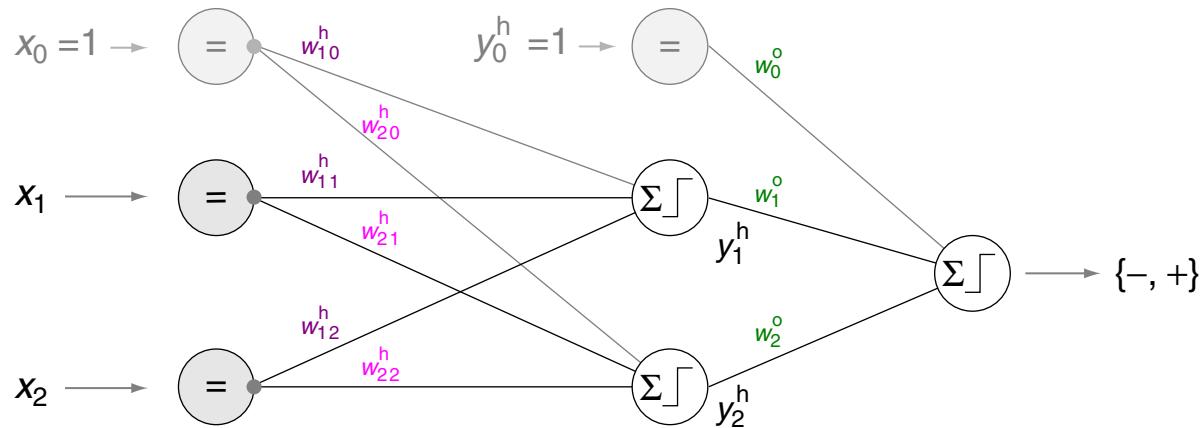
$$W^h = \begin{bmatrix} -0.5 & -1 & 1 \\ 0.5 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$



Multilayer Perceptron Basics

Overcoming the Linear Separability Restriction (continued)

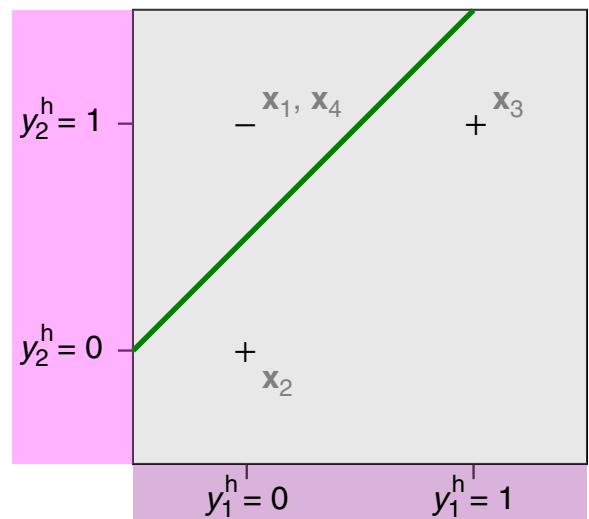
A minimum multilayer perceptron $y(\mathbf{x})$ that can handle the *XOR* problem:



$$y(\mathbf{x}) = \text{heaviside} \left(W^o \left(\text{Heaviside} \left(W^h \mathbf{x} \right) \right) \right)$$

$$W^h = \begin{bmatrix} -0.5 & -1 & 1 \\ 0.5 & -1 & 1 \end{bmatrix} \quad \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

$$W^o = [0.5 \ 1 \ -1]$$

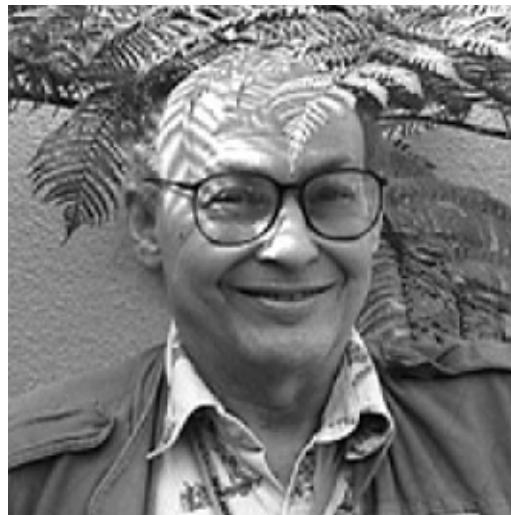


Remarks:

- The first, second, and third layer of the shown multilayer perceptron are called input, hidden, and output layer respectively. Here, in the example, the input layer is comprised of $p+1=3$ units, the hidden layer contains $l+1=3$ units, and the output layer consists of $k=1$ unit.
- Each input unit is connected via a weighted edge to all hidden units (except to the topmost hidden unit, which has a constant input $y_0^h = 1$), resulting in six weights, organized as 2×3 -matrix W^h . Each hidden unit is connected via a weighted edge to the output unit, resulting in three weights, organized as 1×3 -matrix W^o .
- The input units perform no computation but only distribute the values x_0, x_1, x_2 to the next layer. The hidden units (again except the topmost unit) and the output unit apply the *heaviside* function to the sum of their weighted inputs and propagate the result.
I.e., the nine weights $w = (w_{10}^h, \dots, w_{22}^h, w_1^o, w_2^o, w_3^o)$, organized as W^h and W^o , specify the multilayer perceptron (model function) $y(x)$ completely: $y(x) = \text{heaviside}(W^o \begin{pmatrix} 1 \\ \text{Heaviside}(W^h x) \end{pmatrix})$
- The function *Heaviside* denotes the extension of the scalar *heaviside* function to vectors. For $z \in \mathbf{R}^d$ the function *Heaviside*(z) is defined as $(\text{heaviside}(z_1), \dots, \text{heaviside}(z_d))^T$.

Remarks (history) :

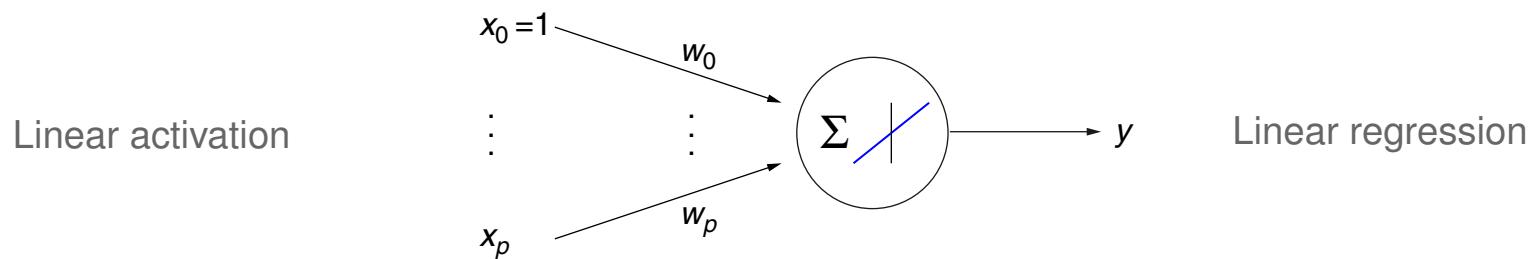
- The multilayer perceptron was presented by Rumelhart and McClelland in 1986. Earlier, but unnoticed, was a similar research work of Werbos and Parker [1974, 1982].
- Compared to a single perceptron, the multilayer perceptron poses a significantly more challenging training (= learning) problem, which requires continuous (and non-linear) threshold functions along with sophisticated learning strategies.
- Marvin Minsky and Seymour Papert showed 1969 with the *XOR* problem the limitations of single perceptrons. Moreover, they assumed that extensions of the perceptron architecture (such as the multilayer perceptron) would be similarly limited as a single perceptron. A fatal mistake. In fact, they brought the research in this field to a halt that lasted 17 years. [[Berkeley](#)]



[Marvin Minsky: [MIT Media Lab](#), [Wikipedia](#)]

Multilayer Perceptron Basics

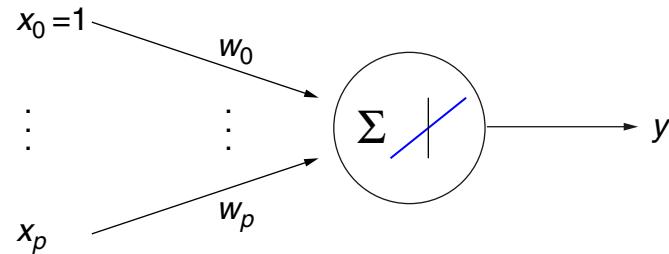
Overcoming the Non-Differentiability Restriction



Multilayer Perceptron Basics

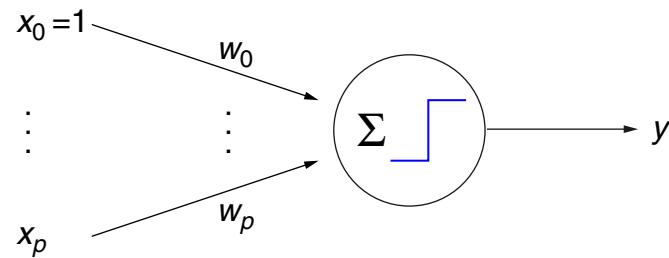
Overcoming the Non-Differentiability Restriction (continued)

Linear activation



Linear regression

Heaviside activation

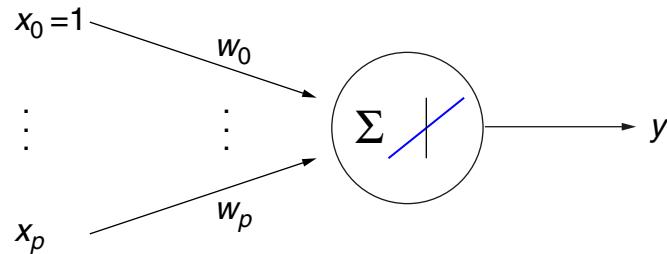


Perceptron algorithm

Multilayer Perceptron Basics

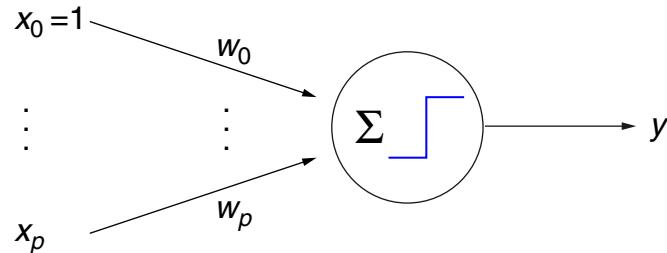
Overcoming the Non-Differentiability Restriction (continued)

Linear activation



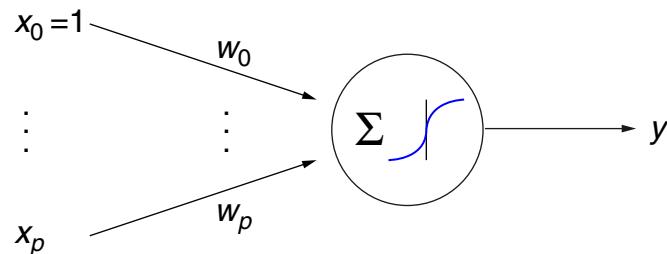
Linear regression

Heaviside activation



Perceptron algorithm

Sigmoid activation

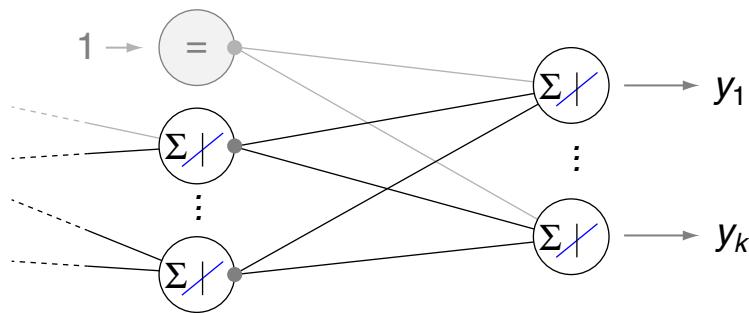


Logistic regression

Multilayer Perceptron Basics

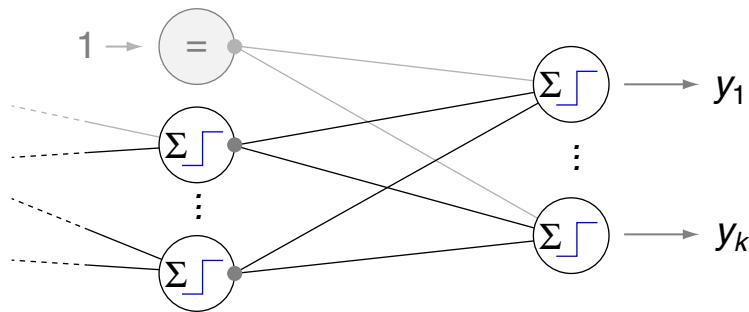
Overcoming the Non-Differentiability Restriction (continued)

Network with linear units



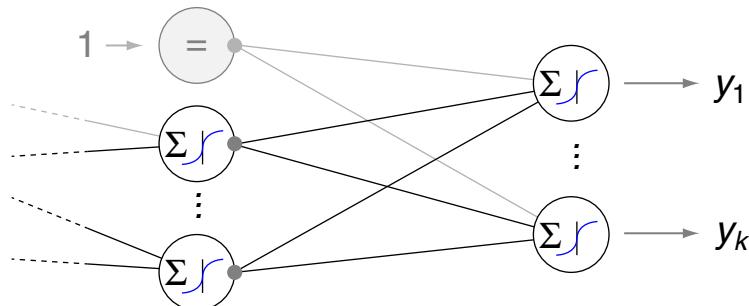
No decision power
beyond a single
hyperplane

Network with heaviside units



Nonlinear decision
boundaries but
no gradient information

Network with sigmoid units



Nonlinear decision
boundaries and
gradient information

Multilayer Perceptron Basics

Unrestricted Classification Problems

Setting:

- X is a multiset of feature vectors from an inner product space \mathbf{X} , $\mathbf{X} \subseteq \mathbf{R}^p$.
- $C = \{0, 1\}^k$ is the set of all multiclass labelings for k classes.
- $D = \{(\mathbf{x}_1, \mathbf{c}_1), \dots, (\mathbf{x}_n, \mathbf{c}_n)\} \subseteq X \times C$ is a multiset of examples.

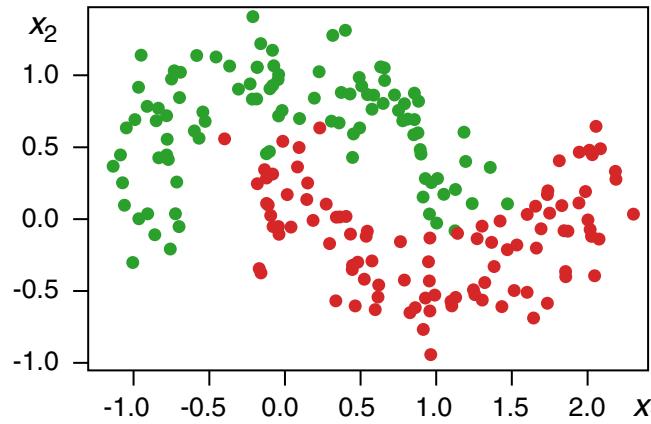
Learning task:

- Fit the examples in D with the multilayer perceptron.

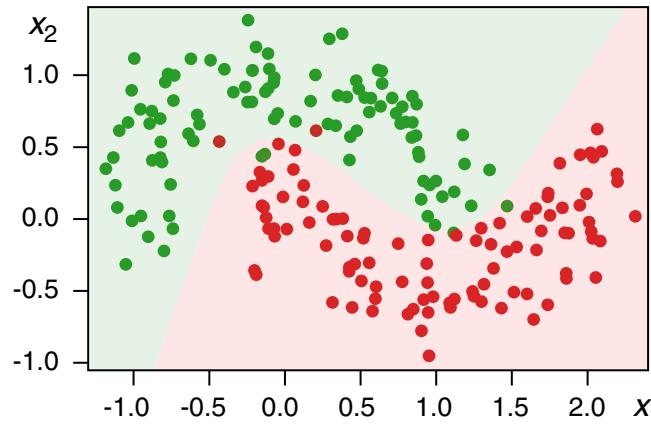
Multilayer Perceptron Basics

Unrestricted Classification Problems: Example

Two-class classification problem:



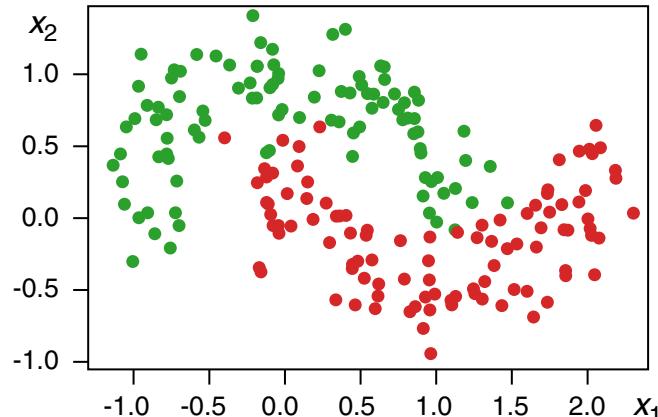
Separated classes:



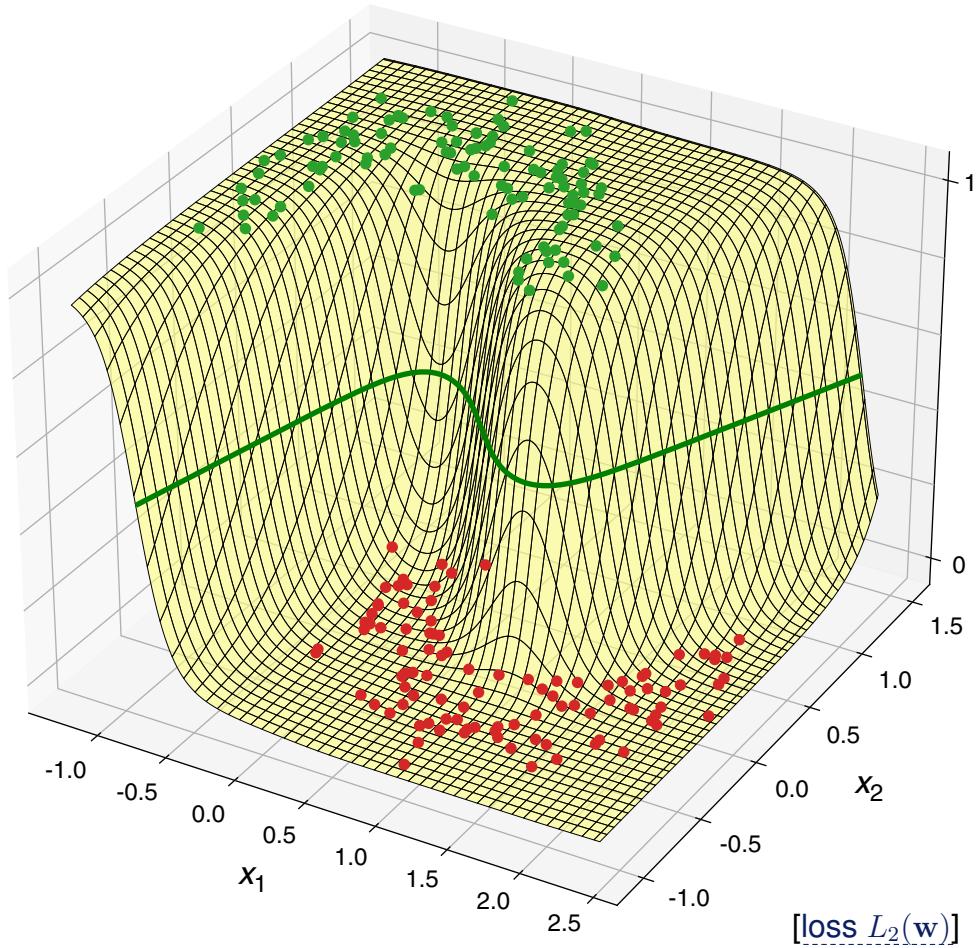
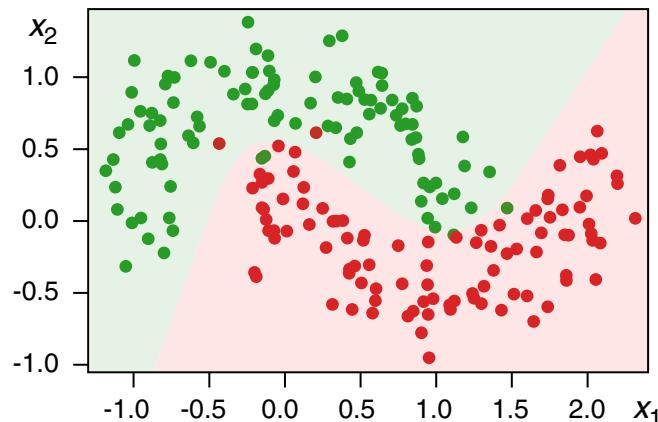
Multilayer Perceptron Basics

Unrestricted Classification Problems: Example (continued)

Two-class classification problem:



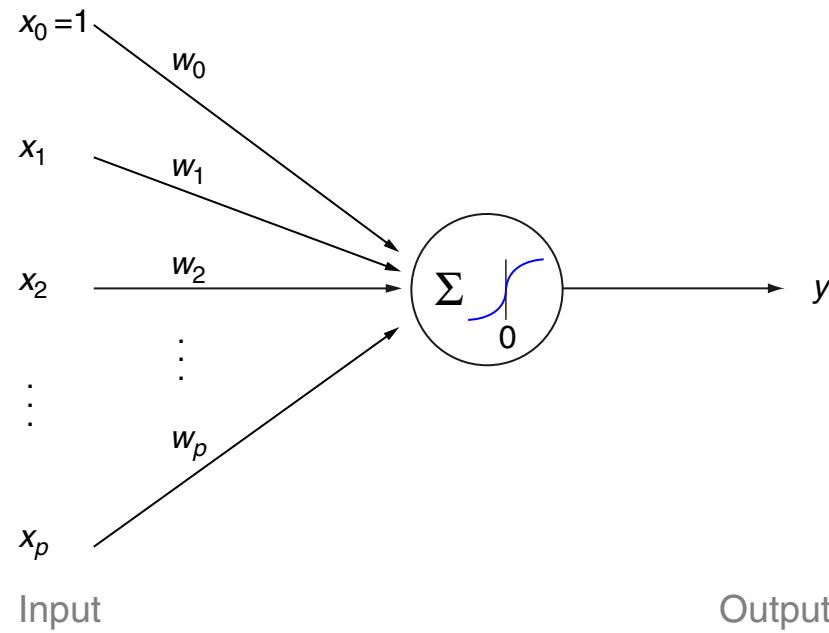
Separated classes:



Multilayer Perceptron Basics

Sigmoid Function [Heaviside]

A perceptron with a continuous *and* non-linear threshold function:



The sigmoid function $\sigma(z)$ as threshold function:

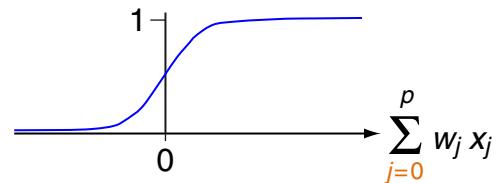
$$\sigma(z) = \frac{1}{1 + e^{-z}}, \quad \frac{d\sigma(z)}{dz} = \sigma(z) \cdot (1 - \sigma(z))$$

Multilayer Perceptron Basics

Sigmoid Function (continued)

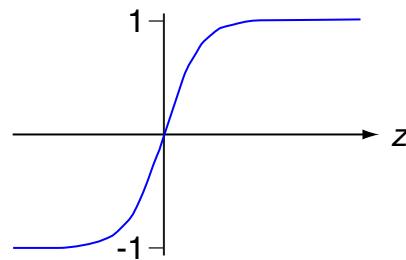
Computation of the perceptron output $y(\mathbf{x})$ via the sigmoid function σ :

$$y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$



An alternative to the sigmoid function is the tanh function:

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \frac{e^{2z} - 1}{e^{2z} + 1}$$



Remarks (derivation of $(\sigma(z))'$):

$$\begin{aligned}\square \quad \frac{d \sigma(z)}{dz} &= \frac{d}{dz} \frac{1}{1+e^{-z}} = \frac{d}{dz} (1+e^{-z})^{-1} \\ &= -1 \cdot (1+e^{-z})^{-2} \cdot (-1) \cdot e^{-z} \\ &= \sigma(z) \cdot \sigma(z) \cdot e^{-z} \\ &= \sigma(z) \cdot \sigma(z) \cdot (1+e^{-z}-1) \\ &= \sigma(z) \cdot \sigma(z) \cdot (\sigma(z)^{-1} - 1) \\ &= \sigma(z) \cdot (1 - \sigma(z))\end{aligned}$$

Remarks (limitation of linear thresholds):

- Employing a nonlinear function as threshold function in the perceptron, such as sigmoid or *heaviside*, is necessary to synthesize complex nonlinear functions via layered composition.
- A “multilayer” perceptron with linear threshold functions can be expressed as a single linear function and hence is equivalent to the power of a single perceptron only.
- Consider the following exemplary composition of three linear functions as a “multilayer” perceptron with p input units, two hidden units, and one output unit: $y(\mathbf{x}) = W^o [W^h \mathbf{x}]$
- The respective weight matrices are as follows:

$$W^h = \begin{bmatrix} w_{11}^h & \dots & w_{1p}^h \\ w_{21}^h & \dots & w_{1p}^h \end{bmatrix}, \quad W^o = \begin{bmatrix} w_1^o & w_2^o \end{bmatrix}$$

Obviously holds:

$$\begin{aligned} y(\mathbf{x}) = W^o [W^h \mathbf{x}] &= W^o \begin{bmatrix} w_{11}^h x_1 + \dots + w_{1p}^h x_p \\ w_{21}^h x_1 + \dots + w_{1p}^h x_p \end{bmatrix} \\ &= w_1^o w_{11}^h x_1 + \dots + w_1^o w_{1p}^h x_p + w_2^o w_{21}^h x_1 + \dots + w_2^o w_{1p}^h x_p \\ &= (w_1^o w_{11}^h + w_2^o w_{21}^h) x_1 + \dots + (w_1^o w_{1p}^h + w_2^o w_{1p}^h) x_p \\ &= w_1 x_1 + \dots + w_p x_p = \mathbf{w}^T \mathbf{x} \end{aligned}$$

Chapter ML:IV (continued)

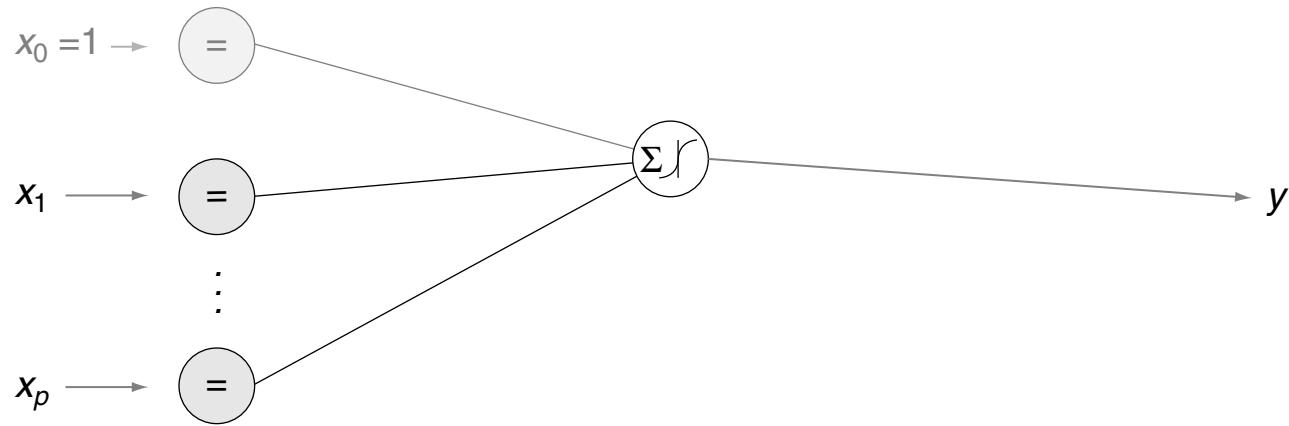
IV. Neural Networks

- Perceptron Learning
- Multilayer Perceptron Basics
- Multilayer Perceptron with Two Layers
- Multilayer Perceptron at Arbitrary Depth
- Advanced MLPs
- Automatic Gradient Computation

Multilayer Perceptron with Two Layers

Network Architecture

A single perceptron $y(\mathbf{x})$:

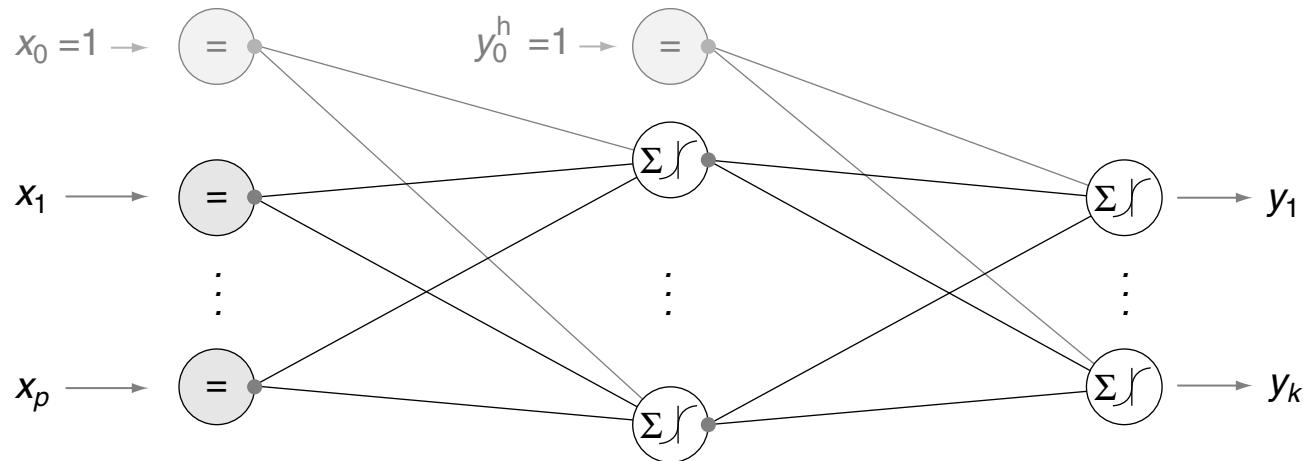


[arbitrary depth]

Multilayer Perceptron with Two Layers

Network Architecture (continued)

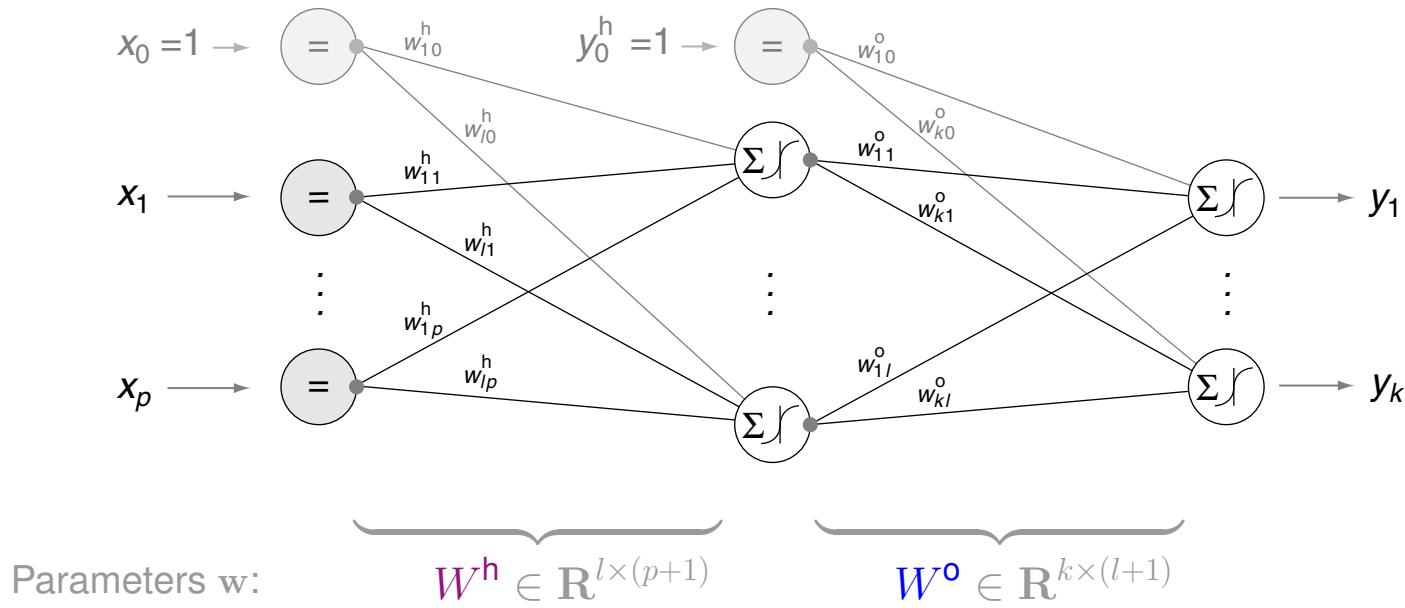
Multilayer perceptron $y(x)$ with a hidden layer and k -dimensional output layer:



Multilayer Perceptron with Two Layers

Network Architecture (continued)

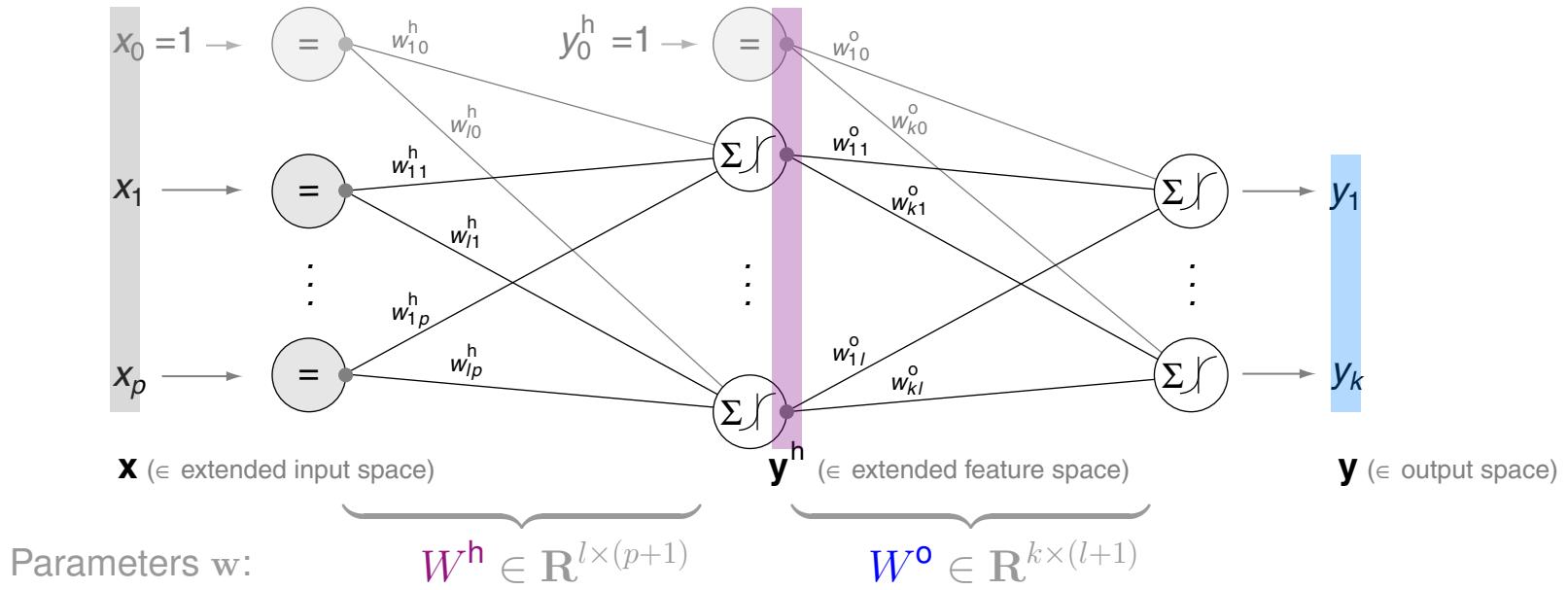
Multilayer perceptron $y(x)$ with a hidden layer and k -dimensional output layer:



Multilayer Perceptron with Two Layers

Network Architecture (continued)

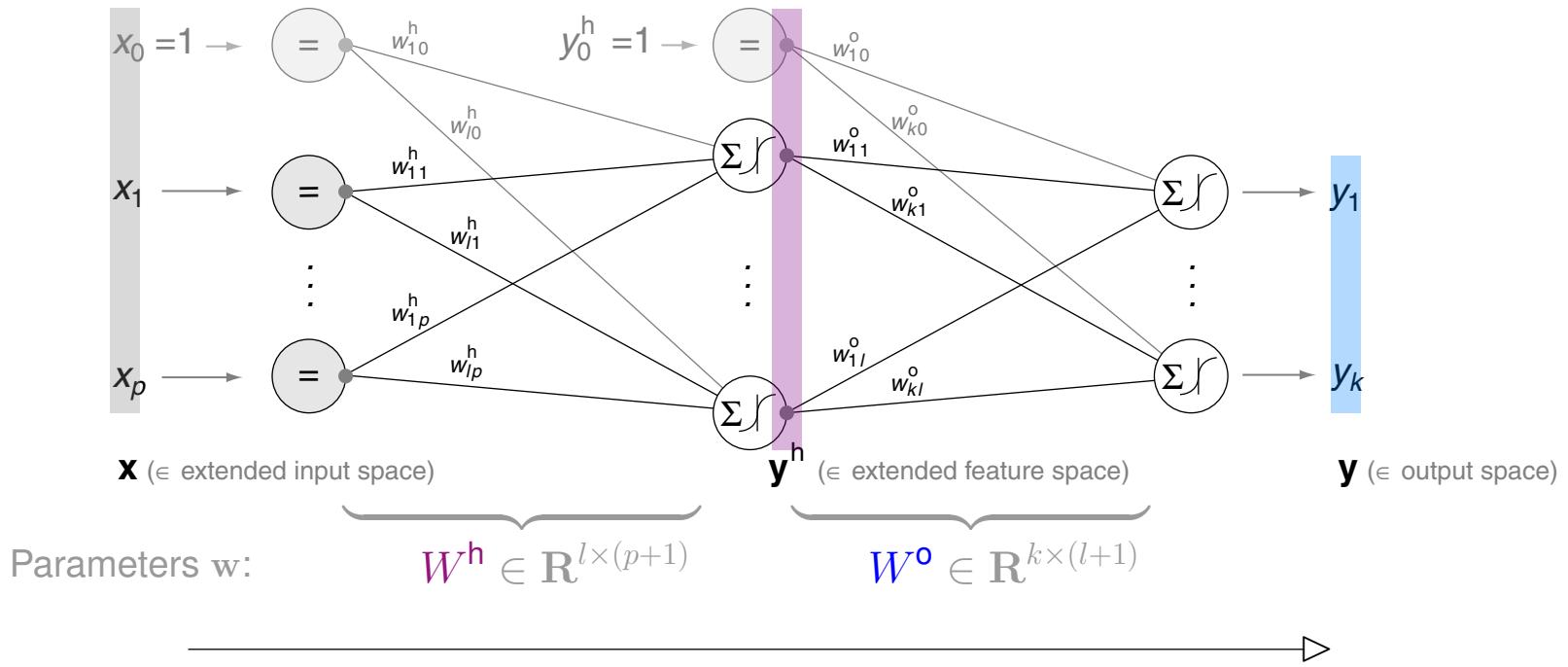
Multilayer perceptron $\mathbf{y}(\mathbf{x})$ with a hidden layer and k -dimensional output layer:



Multilayer Perceptron with Two Layers

(1) Forward Propagation [arbitrary depth]

Multilayer perceptron $\mathbf{y}(\mathbf{x})$ with a hidden layer and k -dimensional output layer:



Model function evaluation (= forward propagation):

$$\mathbf{y}(\mathbf{x}) = \sigma(W^o \mathbf{y}^h(\mathbf{x})) = \sigma\left(W^o \left(\frac{1}{\sigma}(W^h \mathbf{x}) \right)\right)$$

Remarks:

- Each input unit is connected to the hidden units $1, \dots, l$, resulting in $l \cdot (p+1)$ weights, organized as matrix $W^h \in \mathbf{R}^{l \times (p+1)}$. Each hidden unit is connected to the output units $1, \dots, k$, resulting in $k \cdot (l+1)$ weights, organized as matrix $W^o \in \mathbf{R}^{k \times (l+1)}$.
- The hidden units and the output unit(s) apply the (vectorial) sigmoid function, σ , to the sum of their weighted inputs and propagate the result as y^h and y respectively. For $z \in \mathbf{R}^d$ the vectorial sigmoid function $\sigma(z)$ is defined as $(\sigma(z_1), \dots, \sigma(z_d))^T$.
The parameter vector $w = (w_{10}^h, \dots, w_{lp}^h, w_{10}^o, \dots, w_{kl}^o)$, organized as matrices W^h and W^o , specify the multilayer perceptron (model function) $y(x)$ completely: $y(x) = \sigma(W^o (\sigma^{-1}_{(W^h x)}))$
- The shown architecture with k output units allows for the distinction of k classes, either within an exclusive class assignment setting or within a multi-label setting. In the former setting a so-called “softmax layer” can be added subsequent to the output layer to directly return the class label $1, \dots, k$.
- The non-linear characteristic of the sigmoid function allows for networks that approximate every (computable) function. For this capability only three “active” layers are required, i.e., two layers with hidden units and one layer with output units. Keyword: universal approximator
[Kolmogorov theorem, 1957]
- Multilayer perceptrons are also called multilayer networks or (artificial) neural networks, ANN for short.

Multilayer Perceptron with Two Layers

(1) Forward Propagation (continued) [network architecture]

(a) Propagate \mathbf{x} from input to hidden layer: (IGD_{MLP₂} algorithm, Line 5)

$$\mathcal{W}^h \in \mathbf{R}^{l \times (p+1)} \quad \mathbf{x} \in \mathbf{R}^{p+1}$$

$$\sigma \left(\begin{bmatrix} w_{10}^h & \dots & w_{1p}^h \\ \vdots & & \vdots \\ w_{l0}^h & \dots & w_{lp}^h \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{bmatrix} \right) = \begin{bmatrix} y_1^h \\ \vdots \\ y_l^h \end{bmatrix}$$

Multilayer Perceptron with Two Layers

(1) Forward Propagation (continued) [network architecture]

(a) Propagate \mathbf{x} from input to hidden layer: (IGD_{MLP₂} algorithm, Line 5)

$$\mathbf{W}^{\text{h}} \in \mathbf{R}^{l \times (p+1)} \quad \mathbf{x} \in \mathbf{R}^{p+1}$$

$$\sigma \left(\begin{bmatrix} w_{10}^{\text{h}} & \dots & w_{1p}^{\text{h}} \\ \vdots & & \vdots \\ w_{l0}^{\text{h}} & \dots & w_{lp}^{\text{h}} \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{bmatrix} \right) = \begin{bmatrix} y_1^{\text{h}} \\ \vdots \\ y_l^{\text{h}} \end{bmatrix}$$

(b) Propagate \mathbf{y}^{h} from hidden to output layer: (IGD_{MLP₂} algorithm, Line 5)

$$\mathbf{W}^{\text{o}} \in \mathbf{R}^{k \times (l+1)} \quad \mathbf{y}^{\text{h}} \in \mathbf{R}^{l+1} \quad \mathbf{y} \in \mathbf{R}^k$$

$$\sigma \left(\begin{bmatrix} w_{10}^{\text{o}} & \dots & w_{1l}^{\text{o}} \\ \vdots & & \vdots \\ w_{k0}^{\text{o}} & \dots & w_{kl}^{\text{o}} \end{bmatrix} \begin{bmatrix} 1 \\ y_1^{\text{h}} \\ \vdots \\ y_l^{\text{h}} \end{bmatrix} \right) = \begin{bmatrix} y_1 \\ \vdots \\ y_k \end{bmatrix}$$

Multilayer Perceptron with Two Layers

(1) Forward Propagation: Batch Mode [network architecture]

(a) Propagate x from input to hidden layer: (IGD_{MLP₂} algorithm, Line 5)

$$W^h \in \mathbf{R}^{l \times (p+1)} \quad X \subset \mathbf{R}^{p+1}$$

$$\sigma \left(\begin{bmatrix} w_{10}^h & \dots & w_{1p}^h \\ \vdots & & \vdots \\ w_{l0}^h & \dots & w_{lp}^h \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \\ x_{11} & \dots & x_{1n} \\ \vdots & & \vdots \\ x_{p1} & \dots & x_{pn} \end{bmatrix} \right) = \begin{bmatrix} y_{11}^h & \dots & y_{1n}^h \\ \vdots & & \vdots \\ y_{l1}^h & \dots & y_{ln}^h \end{bmatrix}$$

(b) Propagate y^h from hidden to output layer: (IGD_{MLP₂} algorithm, Line 5)

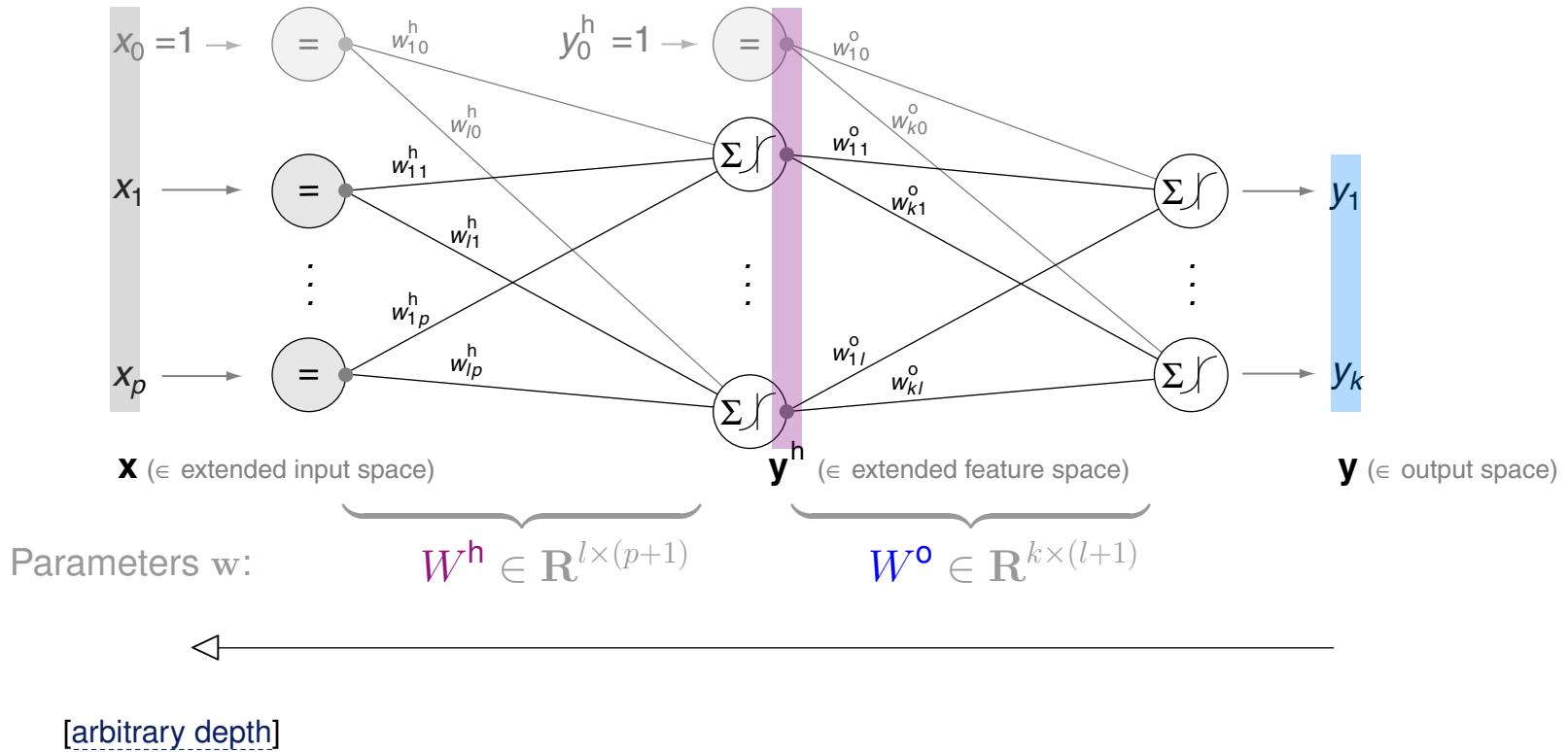
$$W^o \in \mathbf{R}^{k \times (l+1)}$$

$$\sigma \left(\begin{bmatrix} w_{10}^o & \dots & w_{1l}^o \\ \vdots & & \vdots \\ w_{k0}^o & \dots & w_{kl}^o \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \\ y_{11}^h & \dots & y_{1n}^h \\ \vdots & & \vdots \\ y_{l1}^h & \dots & y_{ln}^h \end{bmatrix} \right) = \begin{bmatrix} y_{11} & \dots & y_{1n} \\ \vdots & & \vdots \\ y_{k1} & \dots & y_{kn} \end{bmatrix}$$

Multilayer Perceptron with Two Layers

(2) Backpropagation

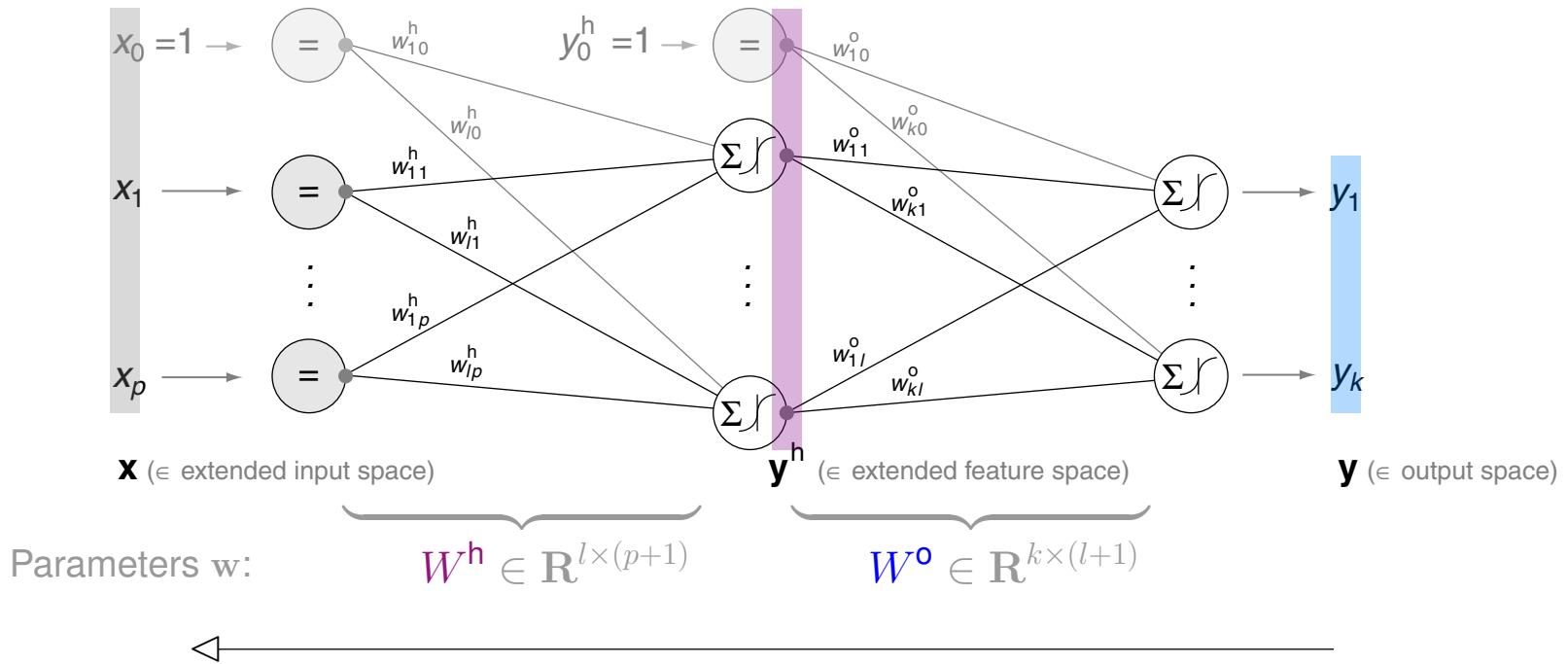
The considered multilayer perceptron $\mathbf{y}(\mathbf{x})$:



Multilayer Perceptron with Two Layers

(2) Backpropagation (continued)

The considered multilayer perceptron $y(\mathbf{x})$:



Calculation of derivatives (= backpropagation) wrt. the global squared loss:

$$L_2(\mathbf{w}) = \frac{1}{2} \cdot \text{RSS}(\mathbf{w}) = \frac{1}{2} \cdot \sum_{(\mathbf{x}, \mathbf{c}) \in D} \sum_{u=1}^k (c_u - y_u(\mathbf{x}))^2$$

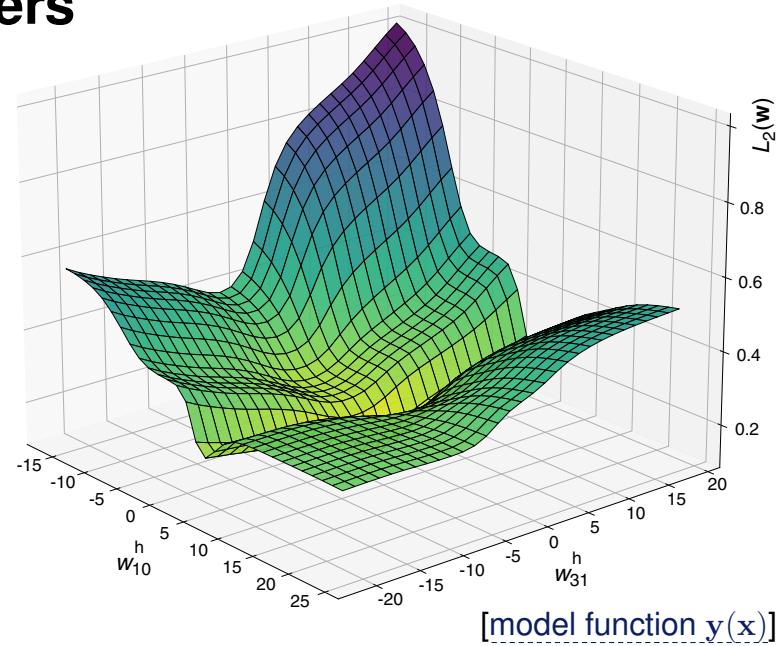
Multilayer Perceptron with Two Layers

(2) Backpropagation (continued)

$L_2(\mathbf{w})$ usually contains various local minima:

$$\mathbf{y}(\mathbf{x}) = \sigma \left(W^o \left(\sigma_{(W^h \mathbf{x})}^1 \right) \right)$$

$$L_2(\mathbf{w}) = \frac{1}{2} \cdot \sum_{(\mathbf{x}, \mathbf{c}) \in D} \sum_{u=1}^k (c_u - y_u(\mathbf{x}))^2$$



Multilayer Perceptron with Two Layers

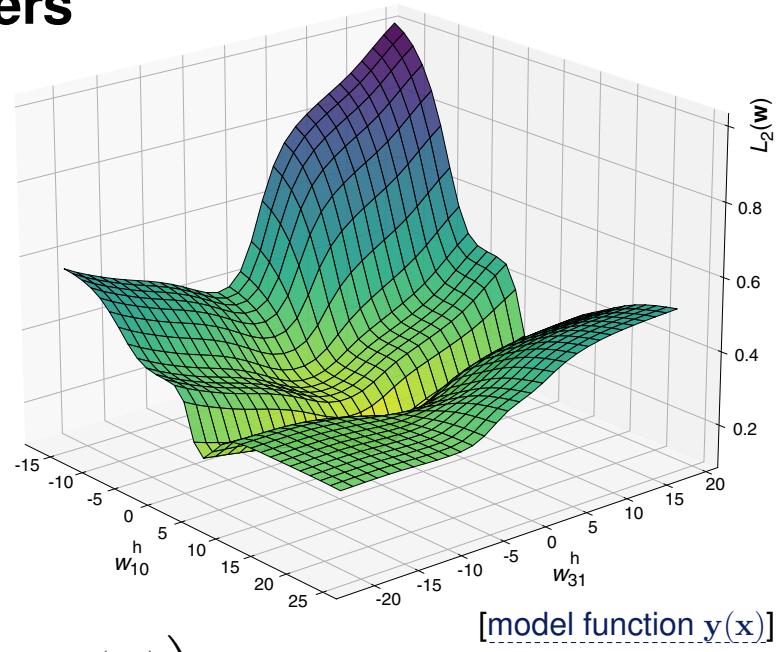
(2) Backpropagation (continued)

$L_2(\mathbf{w})$ usually contains various local minima:

$$\mathbf{y}(\mathbf{x}) = \sigma \left(W^o \left(\sigma^1_{(W^h \mathbf{x})} \right) \right)$$

$$L_2(\mathbf{w}) = \frac{1}{2} \cdot \sum_{(\mathbf{x}, \mathbf{c}) \in D} \sum_{u=1}^k (c_u - y_u(\mathbf{x}))^2$$

$$\nabla L_2(\mathbf{w}) = \left(\frac{\partial L_2(\mathbf{w})}{\partial w_{10}^o}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{kl}^o}, \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^h}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{lp}^h} \right)$$



[model function $\mathbf{y}(\mathbf{x})$]

Multilayer Perceptron with Two Layers

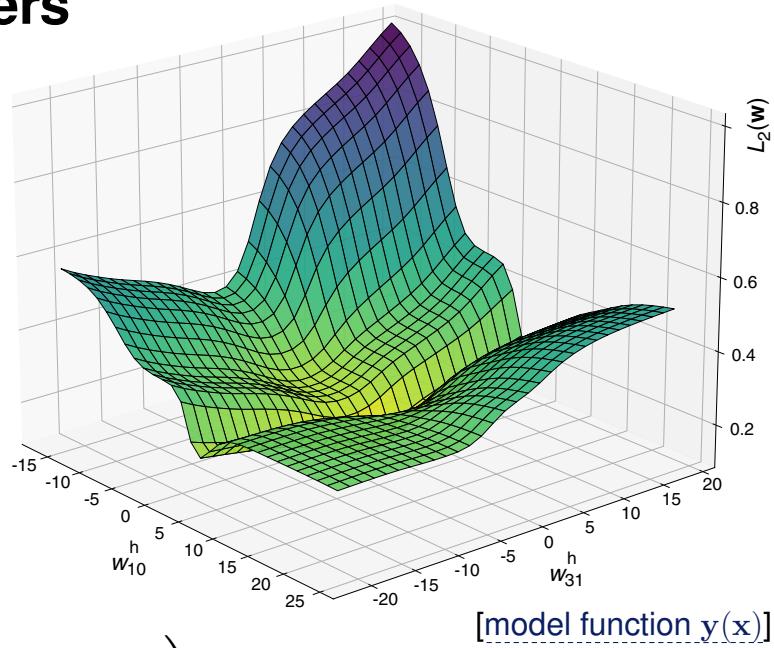
(2) Backpropagation (continued)

$L_2(\mathbf{w})$ usually contains various local minima:

$$\mathbf{y}(\mathbf{x}) = \sigma \left(W^o \left(\sigma_{(W^h \mathbf{x})}^1 \right) \right)$$

$$L_2(\mathbf{w}) = \frac{1}{2} \cdot \sum_{(\mathbf{x}, \mathbf{c}) \in D} \sum_{u=1}^k (c_u - y_u(\mathbf{x}))^2$$

$$\nabla L_2(\mathbf{w}) = \left(\frac{\partial L_2(\mathbf{w})}{\partial w_{10}^o}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{kl}^o}, \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^h}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{lp}^h} \right)$$



[model function $\mathbf{y}(\mathbf{x})$]

(a) Gradient in direction of W^o , written as matrix:

$$\begin{bmatrix} \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^o} & \dots & \frac{\partial L_2(\mathbf{w})}{\partial w_{1l}^o} \\ \vdots & \ddots & \vdots \\ \frac{\partial L_2(\mathbf{w})}{\partial w_{k0}^o} & \dots & \frac{\partial L_2(\mathbf{w})}{\partial w_{kl}^o} \end{bmatrix} \equiv \nabla^o L_2(\mathbf{w})$$

(b) Gradient in direction of W^h :

$$\begin{bmatrix} \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^h} & \dots & \frac{\partial L_2(\mathbf{w})}{\partial w_{1p}^h} \\ \vdots & \ddots & \vdots \\ \frac{\partial L_2(\mathbf{w})}{\partial w_{l0}^h} & \dots & \frac{\partial L_2(\mathbf{w})}{\partial w_{lp}^h} \end{bmatrix} \equiv \nabla^h L_2(\mathbf{w})$$

Remarks:

- “Backpropagation” is short for “backward propagation of errors”.
- Basically, the computation of the gradient $\nabla L_2(\mathbf{w})$ is independent of the organization of the weights in matrices W^h and W^o of a network (model function) $y(x)$. Adopt the following view instead:

To calculate $\nabla L_2(\mathbf{w})$ one has to calculate each of its components $\partial L_2(\mathbf{w})/\partial w$, $w \in \mathbf{w}$, since each weight (parameter) has a certain impact on the global loss $L_2(\mathbf{w})$ of the network. This impact—as well as the computation of this impact—is different for different weights, but it is canonical for all weights of the same layer though: observe that each weight w influences “only” its direct and indirect successor nodes, and that the structure of the influenced successor graph (in fact a tree) is identical for all weights of the same layer.

Hence it is convenient, but not necessary, to process the components of the gradient layer-wise (matrix-wise), as $\nabla^o L_2(\mathbf{w})$ and $\nabla^h L_2(\mathbf{w})$ respectively. Even more, due to the network structure of the model function $y(x)$ only two cases need to be distinguished when deriving the partial derivative $\partial L_2(\mathbf{w})/\partial w$ of an arbitrary weight $w \in \mathbf{w}$: (a) w belongs to the output layer, or (b) w belongs to some hidden layer.

- The derivation of the gradient for the two-layer MLP (and hence the weight update processed in the IGD algorithm) is given in the following, as special case of the derivation of the gradient for MLPs at arbitrary depth.

Multilayer Perceptron with Two Layers

(2) Backpropagation (continued) [arbitrary depth]

(a) Update of weight matrix W^o : (IGD_{MLP₂} algorithm, Lines 7+8)

$$W^o = W^o + \Delta W^o,$$

using the ∇^o -gradient of the loss function $L_2(\mathbf{w})$ to take the steepest descent:

$$\Delta W^o = -\eta \cdot \nabla^o L_2(\mathbf{w})$$

Multilayer Perceptron with Two Layers

(2) Backpropagation (continued) [arbitrary depth]

(a) Update of weight matrix W^o : (IGD_{MLP₂} algorithm, Lines 7+8)

$$W^o = W^o + \Delta W^o,$$

using the ∇^o -gradient of the loss function $L_2(\mathbf{w})$ to take the steepest descent:

$$\Delta W^o = -\eta \cdot \nabla^o L_2(\mathbf{w})$$

$$= -\eta \cdot \begin{bmatrix} \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^o} & \dots & \frac{\partial L_2(\mathbf{w})}{\partial w_{1l}^o} \\ \vdots & & \vdots \\ \frac{\partial L_2(\mathbf{w})}{\partial w_{k0}^o} & \dots & \frac{\partial L_2(\mathbf{w})}{\partial w_{kl}^o} \end{bmatrix}$$

⋮

$$= \eta \cdot \sum_D \underbrace{[(\mathbf{c} - \mathbf{y}(\mathbf{x})) \odot \mathbf{y}(\mathbf{x}) \odot (1 - \mathbf{y}(\mathbf{x}))]}_{\delta^o} \otimes \mathbf{y}^h$$

Multilayer Perceptron with Two Layers

(2) Backpropagation (continued) [arbitrary depth]

(b) Update of weight matrix W^h : (IGD_{MLP₂} algorithm, Lines 7+8)

$$W^h = W^h + \Delta W^h,$$

using the ∇^h -gradient of the loss function $L_2(\mathbf{w})$ to take the steepest descent:

$$\Delta W^h = -\eta \cdot \nabla^h L_2(\mathbf{w})$$

$$= -\eta \cdot \begin{bmatrix} \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^h} & \dots & \frac{\partial L_2(\mathbf{w})}{\partial w_{1p}^h} \\ \vdots & & \vdots \\ \frac{\partial L_2(\mathbf{w})}{\partial w_{l0}^h} & \dots & \frac{\partial L_2(\mathbf{w})}{\partial w_{lp}^h} \end{bmatrix}$$

⋮

$$= \eta \cdot \sum_D \underbrace{\left[((W^o)^T \delta^o) \odot \mathbf{y}^h(\mathbf{x}) \odot (1 - \mathbf{y}^h(\mathbf{x})) \right]_{1,\dots,l}}_{\delta^h} \otimes \mathbf{x}$$

Multilayer Perceptron with Two Layers

The IGD Algorithm

Algorithm: IGD_{MLP₂} Incremental Gradient Descent for the two-layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.

Output: W^h, W^o Weights of $l \cdot (p+1)$ hidden and $k \cdot (l+1)$ output layer units. (= hypothesis)

1. *initialize_random_weights*(W^h, W^o), $t = 0$

2. **REPEAT**

3. $t = t + 1$

4. **FOREACH** $(\mathbf{x}, \mathbf{c}) \in D$ DO

5.

6.

7a.

7b.

8.

9. **ENDDO**

10. **UNTIL**(*convergence*($D, \mathbf{y}(\cdot), t$))

11. **return**(W^h, W^o)

[[Python code](#)]

Multilayer Perceptron with Two Layers

The IGD Algorithm (continued) [arbitrary depth]

Algorithm: IGD_{MLP₂} Incremental Gradient Descent for the two-layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.

η Learning rate, a small positive constant.

Output: W^h, W^o Weights of $l \cdot (p+1)$ hidden and $k \cdot (l+1)$ output layer units. (= hypothesis)

1. *initialize_random_weights*(W^h, W^o), $t = 0$
2. **REPEAT**
3. $t = t + 1$
4. **FOREACH** $(\mathbf{x}, \mathbf{c}) \in D$ **DO**
5. $\mathbf{y}^h(\mathbf{x}) = (\sigma_{(W^h \mathbf{x})}^1)$ // forward propagation; \mathbf{x} is extended by $x_0 = 1$
6. $\mathbf{y}(\mathbf{x}) = \sigma(W^o \mathbf{y}^h(\mathbf{x}))$
- 7a.
- 7b.
- 8.
9. **ENDDO**
10. **UNTIL**(*convergence*($D, \mathbf{y}(\cdot), t$))
11. **return**(W^h, W^o)

[Python code]

Multilayer Perceptron with Two Layers

The IGD Algorithm (continued) [arbitrary depth]

Algorithm: IGD_{MLP₂} Incremental Gradient Descent for the two-layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.

Output: W^h, W^o Weights of $l \cdot (p+1)$ hidden and $k \cdot (l+1)$ output layer units. (= hypothesis)

1. *initialize_random_weights*(W^h, W^o), $t = 0$
2. **REPEAT**
3. $t = t + 1$
4. **FOREACH** $(\mathbf{x}, \mathbf{c}) \in D$ **DO**
5. $\mathbf{y}^h(\mathbf{x}) = (\sigma_{(W^h \mathbf{x})}^1)$ // forward propagation; \mathbf{x} is extended by $x_0 = 1$
 $\mathbf{y}(\mathbf{x}) = \sigma(W^o \mathbf{y}^h(\mathbf{x}))$
6. $\boldsymbol{\delta} = \mathbf{c} - \mathbf{y}(\mathbf{x})$
- 7a.
- 7b.
- 8.
9. **ENDDO**
10. **UNTIL**(*convergence*($D, \mathbf{y}(\cdot), t$))
11. **return**(W^h, W^o)

[Python code]

Multilayer Perceptron with Two Layers

The IGD Algorithm (continued) [arbitrary depth]

Algorithm: IGD_{MLP₂} Incremental Gradient Descent for the two-layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.

Output: W^h, W^o Weights of $l \cdot (p+1)$ hidden and $k \cdot (l+1)$ output layer units. (= hypothesis)

1. *initialize_random_weights*(W^h, W^o), $t = 0$
2. **REPEAT**
3. $t = t + 1$
4. **FOREACH** $(\mathbf{x}, \mathbf{c}) \in D$ **DO**
5. $\mathbf{y}^h(\mathbf{x}) = (\sigma_{(W^h \mathbf{x})}^1)$ // forward propagation; \mathbf{x} is extended by $x_0 = 1$
 $\mathbf{y}(\mathbf{x}) = \sigma(W^o \mathbf{y}^h(\mathbf{x}))$
6. $\boldsymbol{\delta} = \mathbf{c} - \mathbf{y}(\mathbf{x})$
- 7a. $\boldsymbol{\delta}^o = \boldsymbol{\delta} \odot \mathbf{y}(\mathbf{x}) \odot (1 - \mathbf{y}(\mathbf{x}))$ // backpropagation (Steps 7a+7b)
 $\boldsymbol{\delta}^h = [((W^o)^T \boldsymbol{\delta}^o) \odot \mathbf{y}^h \odot (1 - \mathbf{y}^h)]_{1, \dots, l}$
- 7b. $\Delta W^h = \eta \cdot (\boldsymbol{\delta}^h \otimes \mathbf{x})$
 $\Delta W^o = \eta \cdot (\boldsymbol{\delta}^o \otimes \mathbf{y}^h(\mathbf{x}))$
- 8.
9. **ENDDO**
10. **UNTIL**(*convergence*($D, \mathbf{y}(\cdot), t$))
11. **return**(W^h, W^o)

[Python code]

Multilayer Perceptron with Two Layers

The IGD Algorithm (continued) [arbitrary depth]

Algorithm: IGD_{MLP₂} Incremental Gradient Descent for the two-layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.

Output: W^h, W^o Weights of $l \cdot (p+1)$ hidden and $k \cdot (l+1)$ output layer units. (= hypothesis)

1. *initialize_random_weights*(W^h, W^o), $t = 0$
2. **REPEAT**
3. $t = t + 1$
4. **FOREACH** $(\mathbf{x}, \mathbf{c}) \in D$ **DO**
5. $\mathbf{y}^h(\mathbf{x}) = (\sigma_{(W^h \mathbf{x})}^1)$ // forward propagation; \mathbf{x} is extended by $x_0 = 1$
 $\mathbf{y}(\mathbf{x}) = \sigma(W^o \mathbf{y}^h(\mathbf{x}))$
6. $\boldsymbol{\delta} = \mathbf{c} - \mathbf{y}(\mathbf{x})$
- 7a. $\boldsymbol{\delta}^o = \boldsymbol{\delta} \odot \mathbf{y}(\mathbf{x}) \odot (1 - \mathbf{y}(\mathbf{x}))$ // backpropagation (Steps 7a+7b)
 $\boldsymbol{\delta}^h = [((W^o)^T \boldsymbol{\delta}^o) \odot \mathbf{y}^h \odot (1 - \mathbf{y}^h)]_{1, \dots, l}$
- 7b. $\Delta W^h = \eta \cdot (\boldsymbol{\delta}^h \otimes \mathbf{x})$
 $\Delta W^o = \eta \cdot (\boldsymbol{\delta}^o \otimes \mathbf{y}^h(\mathbf{x}))$
8. $W^h = W^h + \Delta W^h$, $W^o = W^o + \Delta W^o$
9. **ENDDO**
10. **UNTIL**(*convergence*($D, \mathbf{y}(\cdot), t$))
11. **return**(W^h, W^o)

[Python code]

Multilayer Perceptron with Two Layers

The IGD Algorithm (continued) [arbitrary depth]

Algorithm: IGD_{MLP₂} Incremental Gradient Descent for the two-layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.

Output: W^h, W^o Weights of $l \cdot (p+1)$ hidden and $k \cdot (l+1)$ output layer units. (= hypothesis)

1. *initialize_random_weights*(W^h, W^o), $t = 0$
2. **REPEAT**
3. $t = t + 1$
4. **FOREACH** $(\mathbf{x}, \mathbf{c}) \in D$ DO
5. 
Model function evaluation.
6. 
Calculation of residual vector.
- 7a. 
Calculation of derivative of the loss.
- 7b. 
Parameter vector update $\hat{=}$ one gradient step down.
8. **ENDDO**
9. **UNTIL**(*convergence*($D, \mathbf{y}(\cdot), t$))
10. **return**(W^h, W^o)

[Python code]

Remarks:

- The symbol $\gg\odot\ll$ denotes the Hadamard product, also known as the element-wise or the Schur product. It is a binary operation that takes two matrices of the same dimensions and produces another matrix of the same dimension as the operands, where each element is the product of the respective elements of the two original matrices. [\[Wikipedia\]](#)
- The symbol $\gg\otimes\ll$ denotes the dyadic product, also called outer product or tensor product. The dyadic product takes two vectors and returns a second order tensor, called a dyadic in this context: $v \otimes w \equiv vw^T$. [\[Wikipedia\]](#)
- ΔW and ${}_D\Delta W$ indicate an update of the weight matrix per batch, D , or per instance, $(x, c) \in D$, respectively.

Chapter ML:IV (continued)

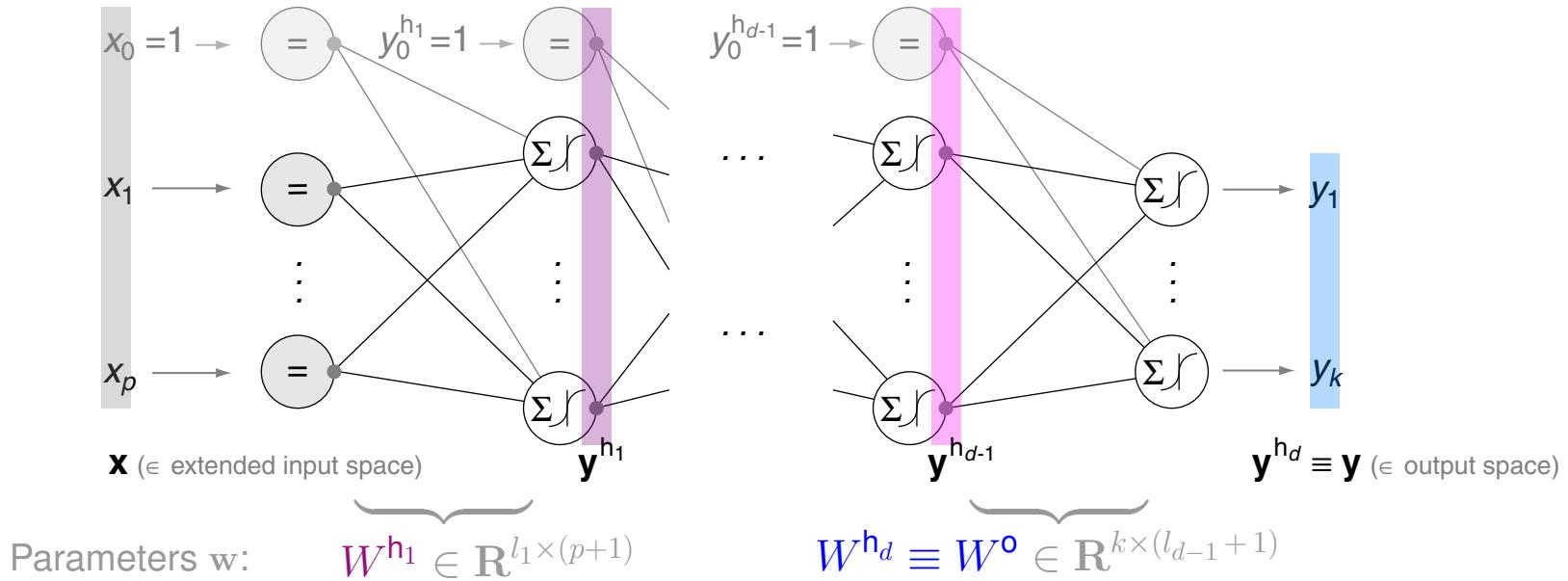
IV. Neural Networks

- Perceptron Learning
- Multilayer Perceptron Basics
- Multilayer Perceptron with Two Layers
- Multilayer Perceptron at Arbitrary Depth
- Advanced MLPs
- Automatic Gradient Computation

Multilayer Perceptron at Arbitrary Depth

Network Architecture [two layers]

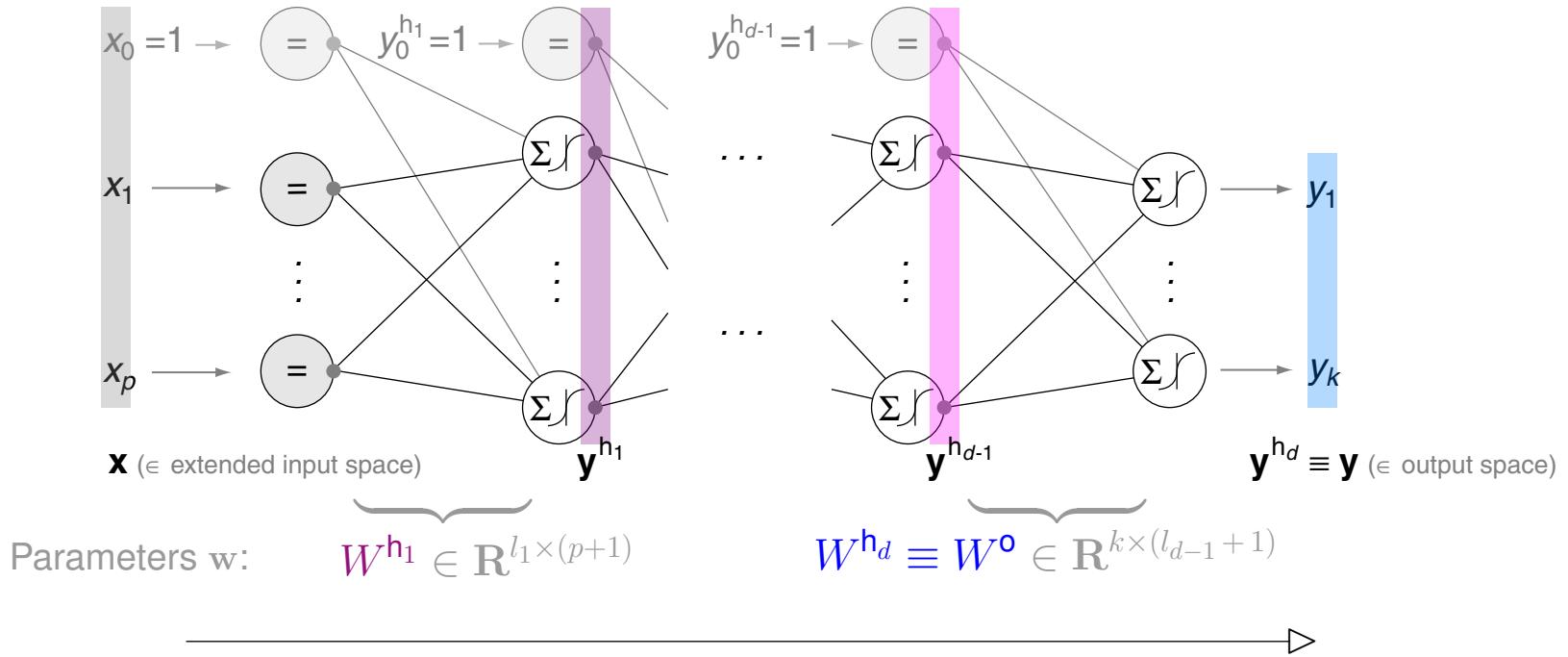
Multilayer perceptron $\mathbf{y}(\mathbf{x})$ with d layers and k -dimensional output:



Multilayer Perceptron at Arbitrary Depth

(1) Forward Propagation [two layers]

Multilayer perceptron $\mathbf{y}(\mathbf{x})$ with d layers and k -dimensional output:



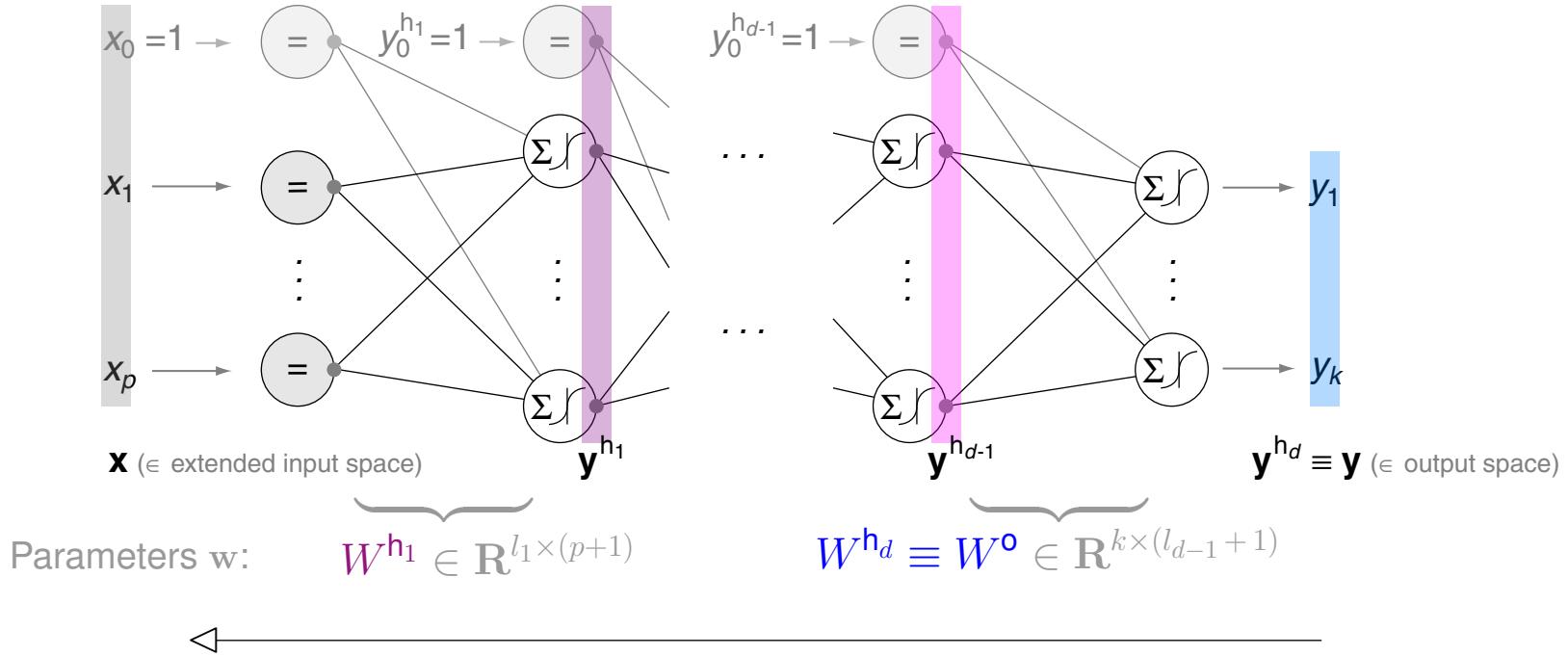
Model function evaluation (= forward propagation):

$$\mathbf{y}^{h_d}(\mathbf{x}) \equiv \mathbf{y}(\mathbf{x}) = \sigma(W^{h_d} \mathbf{y}^{h_{d-1}}(\mathbf{x})) = \dots = \sigma \left(W^{h_d} \left(\sigma \left(\dots \left(\sigma \left(W^{h_1} \mathbf{x} \right) \dots \right) \right) \right) \right)$$

Multilayer Perceptron at Arbitrary Depth

(2) Backpropagation [two layers]

The considered multilayer perceptron $y(\mathbf{x})$:



Calculation of derivatives (= backpropagation) wrt. the global squared loss:

$$L_2(\mathbf{w}) = \frac{1}{2} \cdot \text{RSS}(\mathbf{w}) = \frac{1}{2} \cdot \sum_{(\mathbf{x}, \mathbf{c}) \in D} \sum_{u=1}^k (c_u - y_u(\mathbf{x}))^2$$

Multilayer Perceptron at Arbitrary Depth

(2) Backpropagation (continued)

$$\nabla L_2(\mathbf{w}) = \left(\frac{\partial L_2(\mathbf{w})}{\partial w_{10}^{h_1}}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{l_1 p}^{h_1}}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^{h_d}}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{k l_{d-1}}^{h_d}} \right), \text{ where } l_s = \text{no_rows}(W^{h_s})$$

Multilayer Perceptron at Arbitrary Depth

(2) Backpropagation (continued) [two layers]

$$\nabla L_2(\mathbf{w}) = \left(\frac{\partial L_2(\mathbf{w})}{\partial w_{10}^{h_1}}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{l_1 p}^{h_1}}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^{h_d}}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{k l_{d-1}}^{h_d}} \right), \text{ where } l_s = \text{no_rows}(W^{h_s})$$

Update of weight matrix W^{h_s} , $1 \leq s \leq d$: (IGD_{MLP_d} algorithm, Lines 7+8)

$$W^{h_s} = W^{h_s} + \Delta W^{h_s},$$

using the ∇^{h_s} -gradient of the loss function $L_2(\mathbf{w})$ to take the steepest descent:

$$\Delta W^{h_s} = -\eta \cdot \nabla^{h_s} L_2(\mathbf{w})$$

Multilayer Perceptron at Arbitrary Depth

(2) Backpropagation (continued) [two layers]

$$\nabla L_2(\mathbf{w}) = \left(\frac{\partial L_2(\mathbf{w})}{\partial w_{10}^{h_1}}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{l_1 p}^{h_1}}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^{h_d}}, \dots, \frac{\partial L_2(\mathbf{w})}{\partial w_{k l_{d-1}}^{h_d}} \right), \text{ where } l_s = \text{no_rows}(W^{h_s})$$

Update of weight matrix W^{h_s} , $1 \leq s \leq d$: (IGD_{MLP_d} algorithm, Lines 7+8)

$$W^{h_s} = W^{h_s} + \Delta W^{h_s},$$

using the ∇^{h_s} -gradient of the loss function $L_2(\mathbf{w})$ to take the steepest descent:

$$\Delta W^{h_s} = -\eta \cdot \nabla^{h_s} L_2(\mathbf{w})$$

$$= -\eta \cdot \begin{bmatrix} \frac{\partial L_2(\mathbf{w})}{\partial w_{10}^{h_s}} & \cdots & \frac{\partial L_2(\mathbf{w})}{\partial w_{l_{s-1}}^{h_s}} \\ \vdots & \ddots & \vdots \\ \frac{\partial L_2(\mathbf{w})}{\partial w_{l_s 0}^{h_s}} & \cdots & \frac{\partial L_2(\mathbf{w})}{\partial w_{l_s l_{s-1}}^{h_s}} \end{bmatrix}, \quad \begin{aligned} l_s &= \text{no_rows}(W^{h_s}), \\ \text{where } \mathbf{y}^{h_0} &\equiv \mathbf{x}, \\ \mathbf{y}^{h_d} &\equiv \mathbf{y} \end{aligned}$$

⋮

→ p. 110

Multilayer Perceptron at Arbitrary Depth

(2) Backpropagation (continued) [two layers]

$$\Delta W^{h_s} = \begin{cases} \eta \cdot \sum_D \underbrace{[(\mathbf{c} - \mathbf{y}(\mathbf{x})) \odot \mathbf{y}(\mathbf{x}) \odot (1 - \mathbf{y}(\mathbf{x}))]}_{\delta^{h_d} \equiv \delta^o} \otimes \mathbf{y}^{h_{d-1}}(\mathbf{x}) & \text{if } s = d \\ \eta \cdot \sum_D \underbrace{\left[\left((W^{h_{s+1}})^T \delta^{h_{s+1}} \right) \odot \mathbf{y}^{h_s}(\mathbf{x}) \odot (1 - \mathbf{y}^{h_s}(\mathbf{x})) \right]_{1, \dots, l_s}}_{\delta^{h_s}} \otimes \mathbf{y}^{h_{s-1}}(\mathbf{x}) & \text{if } 1 < s < d \\ \eta \cdot \sum_D \underbrace{\left[\left((W^{h_2})^T \delta^{h_2} \right) \odot \mathbf{y}^{h_1}(\mathbf{x}) \odot (1 - \mathbf{y}^{h_1}(\mathbf{x})) \right]_{1, \dots, l_1}}_{\delta^{h_1}} \otimes \mathbf{x} & \text{if } s = 1 \end{cases}$$

where $l_s = \text{no_rows}(W^{h_s})$

Multilayer Perceptron at Arbitrary Depth

The IGD Algorithm

Algorithm: $\text{IGD}_{\text{MLP}_d}$ Incremental Gradient Descent for the d -layer MLP
Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.
Output: W^{h_1}, \dots, W^{h_d} Weight matrices of the d layers. (= hypothesis)

1. FOR $s = 1$ TO d DO *initialize_random_weights*(W^{h_s}) ENDDO, $t = 0$

2. **REPEAT**

3. $t = t + 1$

4. FOREACH $(\mathbf{x}, \mathbf{c}) \in D$ DO

5.

6.

7a.

7b.

8.

9. ENDDO

10. **UNTIL**(*convergence*($D, \mathbf{y}(\cdot), t$))

11. *return*(W^{h_1}, \dots, W^{h_d})

[Python code]

Multilayer Perceptron at Arbitrary Depth

The IGD Algorithm (continued) [two layers]

Algorithm: $\text{IGD}_{\text{MLP}_d}$ Incremental Gradient Descent for the d -layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.

Output: W^{h_1}, \dots, W^{h_d} Weight matrices of the d layers. (= hypothesis)

```
1. FOR  $s = 1$  TO  $d$  DO initialize_random_weights( $W^{h_s}$ ) ENDDO,  $t = 0$ 
2. REPEAT
3.    $t = t + 1$ 
4.   FOREACH  $(\mathbf{x}, \mathbf{c}) \in D$  DO
5.      $\mathbf{y}^{h_1}(\mathbf{x}) = (\sigma_{(W^{h_1} \mathbf{x})}^1)$  // forward propagation;  $\mathbf{x}$  is extended by  $x_0 = 1$ 
       FOR  $s = 2$  TO  $d-1$  DO  $\mathbf{y}^{h_s}(\mathbf{x}) = (\sigma_{(W^{h_s} \mathbf{y}^{h_{s-1}}(\mathbf{x}))}^1)$  ENDDO
        $\mathbf{y}(\mathbf{x}) = \sigma(W^{h_d} \mathbf{y}^{h_{d-1}}(\mathbf{x}))$ 
6.
7a.

7b.

8.
9.   ENDDO
10.  UNTIL(convergence( $D, \mathbf{y}(\cdot), t$ ))
11.  return( $W^{h_1}, \dots, W^{h_d}$ )
```

[Python code]

Multilayer Perceptron at Arbitrary Depth

The IGD Algorithm (continued) [two layers]

Algorithm: $\text{IGD}_{\text{MLP}_d}$ Incremental Gradient Descent for the d -layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.

Output: W^{h_1}, \dots, W^{h_d} Weight matrices of the d layers. (= hypothesis)

```
1. FOR  $s = 1$  TO  $d$  DO initialize_random_weights( $W^{h_s}$ ) ENDDO,  $t = 0$ 
2. REPEAT
3.    $t = t + 1$ 
4.   FOREACH  $(\mathbf{x}, \mathbf{c}) \in D$  DO
5.      $\mathbf{y}^{h_1}(\mathbf{x}) = (\sigma_{(W^{h_1} \mathbf{x})}^1)$  // forward propagation;  $\mathbf{x}$  is extended by  $x_0 = 1$ 
       FOR  $s = 2$  TO  $d-1$  DO  $\mathbf{y}^{h_s}(\mathbf{x}) = (\sigma_{(W^{h_s} \mathbf{y}^{h_{s-1}}(\mathbf{x}))}^1)$  ENDDO
        $\mathbf{y}(\mathbf{x}) = \sigma(W^{h_d} \mathbf{y}^{h_{d-1}}(\mathbf{x}))$ 
6.    $\delta = \mathbf{c} - \mathbf{y}(\mathbf{x})$ 
7a.
7b.

8.
9.   ENDDO
10.  UNTIL(convergence( $D, \mathbf{y}(\cdot), t$ ))
11.  return( $W^{h_1}, \dots, W^{h_d}$ )
```

[Python code]

Multilayer Perceptron at Arbitrary Depth

The IGD Algorithm (continued) [two layers]

Algorithm: $\text{IGD}_{\text{MLP}_d}$ Incremental Gradient Descent for the d -layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbf{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.

η Learning rate, a small positive constant.

Output: W^{h_1}, \dots, W^{h_d} Weight matrices of the d layers. (= hypothesis)

```
1. FOR  $s = 1$  TO  $d$  DO initialize_random_weights( $W^{h_s}$ ) ENDDO,  $t = 0$ 
2. REPEAT
3.    $t = t + 1$ 
4.   FOREACH  $(\mathbf{x}, \mathbf{c}) \in D$  DO
5.      $\mathbf{y}^{h_1}(\mathbf{x}) = (\sigma_{(W^{h_1} \mathbf{x})}^1)$  // forward propagation;  $\mathbf{x}$  is extended by  $x_0 = 1$ 
       FOR  $s = 2$  TO  $d-1$  DO  $\mathbf{y}^{h_s}(\mathbf{x}) = (\sigma_{(W^{h_s} \mathbf{y}^{h_{s-1}}(\mathbf{x}))}^1)$  ENDDO
        $\mathbf{y}(\mathbf{x}) = \sigma(W^{h_d} \mathbf{y}^{h_{d-1}}(\mathbf{x}))$ 
6.      $\delta = \mathbf{c} - \mathbf{y}(\mathbf{x})$ 
7a.     $\delta^{h_d} = \delta \odot \mathbf{y}(\mathbf{x}) \odot (1 - \mathbf{y}(\mathbf{x}))$  // backpropagation (Steps 7a+7b)
          FOR  $s = d-1$  DOWNTO 1 DO  $\delta^{h_s} = [((W^{h_{s+1}})^T \delta^{h_{s+1}}) \odot \mathbf{y}^{h_s}(\mathbf{x}) \odot (1 - \mathbf{y}^{h_s}(\mathbf{x}))]_{1, \dots, l_s}$  ENDDO
7b.     $\Delta W^{h_1} = \eta \cdot (\delta^{h_1} \otimes \mathbf{x})$ 
          FOR  $s = 2$  TO  $d$  DO  $\Delta W^{h_s} = \eta \cdot (\delta^{h_s} \otimes \mathbf{y}^{h_{s-1}}(\mathbf{x}))$  ENDDO
8.
9.   ENDDO
10.  UNTIL(convergence( $D, \mathbf{y}(\cdot), t$ ))
11.  return( $W^{h_1}, \dots, W^{h_d}$ )
```

[Python code]

Multilayer Perceptron at Arbitrary Depth

The IGD Algorithm (continued) [two layers]

Algorithm: $\text{IGD}_{\text{MLP}_d}$ Incremental Gradient Descent for the d -layer MLP

Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.

Output: W^{h_1}, \dots, W^{h_d} Weight matrices of the d layers. (= hypothesis)

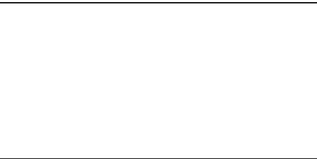
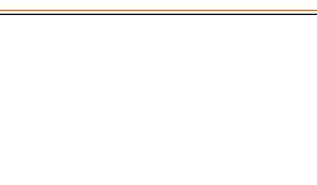
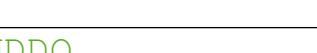
```
1. FOR  $s = 1$  TO  $d$  DO initialize_random_weights( $W^{h_s}$ ) ENDDO,  $t = 0$ 
2. REPEAT
3.    $t = t + 1$ 
4.   FOREACH  $(\mathbf{x}, \mathbf{c}) \in D$  DO
5.      $\mathbf{y}^{h_1}(\mathbf{x}) = (\sigma_{(W^{h_1} \mathbf{x})}^1)$  // forward propagation;  $\mathbf{x}$  is extended by  $x_0 = 1$ 
      FOR  $s = 2$  TO  $d-1$  DO  $\mathbf{y}^{h_s}(\mathbf{x}) = (\sigma_{(W^{h_s} \mathbf{y}^{h_{s-1}}(\mathbf{x}))}^1)$  ENDDO
       $\mathbf{y}(\mathbf{x}) = \sigma(W^{h_d} \mathbf{y}^{h_{d-1}}(\mathbf{x}))$ 
6.      $\delta = \mathbf{c} - \mathbf{y}(\mathbf{x})$ 
7a.     $\delta^{h_d} = \delta \odot \mathbf{y}(\mathbf{x}) \odot (1 - \mathbf{y}(\mathbf{x}))$  // backpropagation (Steps 7a+7b)
          FOR  $s = d-1$  DOWNTO 1 DO  $\delta^{h_s} = [((W^{h_{s+1}})^T \delta^{h_{s+1}}) \odot \mathbf{y}^{h_s}(\mathbf{x}) \odot (1 - \mathbf{y}^{h_s}(\mathbf{x}))]_{1, \dots, l_s}$  ENDDO
7b.     $\Delta W^{h_1} = \eta \cdot (\delta^{h_1} \otimes \mathbf{x})$ 
          FOR  $s = 2$  TO  $d$  DO  $\Delta W^{h_s} = \eta \cdot (\delta^{h_s} \otimes \mathbf{y}^{h_{s-1}}(\mathbf{x}))$  ENDDO
8.    FOR  $s = 1$  TO  $d$  DO  $W^{h_s} = W^{h_s} + \Delta W^{h_s}$  ENDDO
9.  ENDDO
10. UNTIL(convergence( $D, \mathbf{y}(\cdot), t$ ))
11. return( $W^{h_1}, \dots, W^{h_d}$ )
```

[Python code]

Multilayer Perceptron at Arbitrary Depth

The IGD Algorithm (continued) [two layers]

Algorithm: $\text{IGD}_{\text{MLP}_d}$ Incremental Gradient Descent for the d -layer MLP
Input: D Multiset of examples (\mathbf{x}, \mathbf{c}) with $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.
Output: W^{h_1}, \dots, W^{h_d} Weight matrices of the d layers. (= hypothesis)

1. FOR $s = 1$ TO d DO *initialize_random_weights*(W^{h_s}) ENDDO, $t = 0$
2. **REPEAT**
3. $t = t + 1$
4. FOREACH $(\mathbf{x}, \mathbf{c}) \in D$ DO
5. 
Model function evaluation.
6. 
Calculation of residual vector.
- 7a. 
Calculation of derivative of the loss.
- 7b. 
Parameter vector update $\hat{=}$ one gradient step down.
8. 
9. ENDDO
10. **UNTIL**(*convergence*($D, \mathbf{y}(\cdot), t$))
11. **return**(W^{h_1}, \dots, W^{h_d})

[Python code]

Remarks (derivation of $\nabla^{\mathbf{h}_s} L_2(\mathbf{w})$) :

- Partial derivative for a weight in a weight matrix $W^{\mathbf{h}_s}$, $1 \leq s \leq d$:

$$\begin{aligned}
\frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} L_2(\mathbf{w}) &= \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} \frac{1}{2} \cdot \sum_{(\mathbf{x}, \mathbf{c}) \in D} \sum_{u=1}^k (c_u - y_u(\mathbf{x}))^2 \\
&= \frac{1}{2} \cdot \sum_D \sum_{u=1}^k \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} (c_u - y_u(\mathbf{x}))^2 \\
&= - \sum_D \sum_{u=1}^k (c_u - y_u(\mathbf{x})) \cdot \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} y_u(\mathbf{x}) \\
&\stackrel{(1,2)}{=} - \sum_D \sum_{u=1}^k \underbrace{(c_u - y_u(\mathbf{x})) \cdot y_u(\mathbf{x}) \cdot (1 - y_u(\mathbf{x}))}_{\delta_u^{\mathbf{h}_d} \equiv \delta_u^o} \cdot \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} W_{u*}^{\mathbf{h}_d} \mathbf{y}^{\mathbf{h}_{d-1}}(\mathbf{x}) \\
&\stackrel{(3)}{=} - \sum_D \sum_{u=1}^k \delta_u^{\mathbf{h}_d} \cdot \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} \sum_{v=0}^{l_{d-1}} w_{uv}^{\mathbf{h}_d} \cdot y_v^{\mathbf{h}_{d-1}}(\mathbf{x})
\end{aligned}$$

- Partial derivative for a weight in $W^{\mathbf{h}_d}$ (output layer), i.e., $s = d$:

$$\begin{aligned}
\frac{\partial}{\partial w_{ij}^{\mathbf{h}_d}} L_2(\mathbf{w}) &= - \sum_D \sum_{u=1}^k \delta_u^{\mathbf{h}_d} \cdot \sum_{v=0}^{l_{d-1}} \frac{\partial}{\partial w_{ij}^{\mathbf{h}_d}} w_{uv}^{\mathbf{h}_d} \cdot y_v^{\mathbf{h}_{d-1}}(\mathbf{x}) \quad // \text{ Only for the term where } u = i \text{ and } v = j \text{ the partial derivative is nonzero. See the illustration.} \\
&= - \sum_D \delta_i^{\mathbf{h}_d} \cdot y_j^{\mathbf{h}_{d-1}}(\mathbf{x})
\end{aligned}$$

Remarks (derivation of $\nabla^{\mathbf{h}_s} L_2(\mathbf{w})$) : (continued)

- Partial derivative for a weight in a weight matrix $W^{\mathbf{h}_s}$, $s \leq d-1$:

$$\begin{aligned}
 \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} L_2(\mathbf{w}) &= - \sum_D \sum_{u=1}^k \delta_u^{\mathbf{h}_d} \cdot \sum_{v=0}^{l_{d-1}} \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} w_{uv}^{\mathbf{h}_d} \cdot y_v^{\mathbf{h}_{d-1}}(\mathbf{x}) \quad // \text{Every component of } y^{\mathbf{h}_{d-1}}(\mathbf{x}) \\
 &\stackrel{(1,2)}{=} - \sum_D \sum_{u=1}^k \delta_u^{\mathbf{h}_d} \cdot \sum_{v=1}^{l_{d-1}} w_{uv}^{\mathbf{h}_d} \cdot y_v^{\mathbf{h}_{d-1}}(\mathbf{x}) \cdot (1 - y_v^{\mathbf{h}_{d-1}}(\mathbf{x})) \cdot \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} W_{v*}^{\mathbf{h}_{d-1}} \mathbf{y}^{\mathbf{h}_{d-2}}(\mathbf{x}) \\
 &\stackrel{(4)}{=} - \sum_D \sum_{v=1}^{l_{d-1}} \sum_{u=1}^k \delta_u^{\mathbf{h}_d} \cdot w_{uv}^{\mathbf{h}_d} \cdot y_v^{\mathbf{h}_{d-1}}(\mathbf{x}) \cdot (1 - y_v^{\mathbf{h}_{d-1}}(\mathbf{x})) \cdot \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} W_{v*}^{\mathbf{h}_{d-1}} \mathbf{y}^{\mathbf{h}_{d-2}}(\mathbf{x}) \\
 &\stackrel{(5)}{=} - \sum_D \sum_{v=1}^{l_{d-1}} \underbrace{\left(W_{*v}^{\mathbf{h}_d} \right)^T \boldsymbol{\delta}^{\mathbf{h}_d} \cdot y_v^{\mathbf{h}_{d-1}}(\mathbf{x}) \cdot (1 - y_v^{\mathbf{h}_{d-1}}(\mathbf{x}))}_{\delta_v^{\mathbf{h}_{d-1}}} \cdot \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} W_{v*}^{\mathbf{h}_{d-1}} \mathbf{y}^{\mathbf{h}_{d-2}}(\mathbf{x}) \\
 &\stackrel{(3)}{=} - \sum_D \sum_{v=1}^{l_{d-1}} \delta_v^{\mathbf{h}_{d-1}} \cdot \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} \sum_{w=0}^{l_{d-2}} w_{vw}^{\mathbf{h}_{d-1}} \cdot y_w^{\mathbf{h}_{d-2}}(\mathbf{x})
 \end{aligned}$$

- Partial derivative for a weight in $W^{\mathbf{h}_{d-1}}$ (next to output layer), i.e., $s = d-1$:

$$\begin{aligned}
 \frac{\partial}{\partial w_{ij}^{\mathbf{h}_{d-1}}} L_2(\mathbf{w}) &= - \sum_D \sum_{v=1}^{l_{d-1}} \delta_v^{\mathbf{h}_{d-1}} \sum_{w=0}^{l_{d-2}} \frac{\partial}{\partial w_{ij}^{\mathbf{h}_{d-1}}} w_{vw}^{\mathbf{h}_{d-1}} \cdot y_w^{\mathbf{h}_{d-2}}(\mathbf{x}) \quad // \text{Only for the term where } v = i \\
 &= - \sum_D \delta_i^{\mathbf{h}_{d-1}} \cdot y_j^{\mathbf{h}_{d-2}}(\mathbf{x})
 \end{aligned}$$

Remarks (derivation of $\nabla^{\mathbf{h}_s} L_2(\mathbf{w})$) : (continued)

- Instead of writing out the recursion further, i.e., considering a weight matrix $W^{\mathbf{h}_s}$, $s \leq d-2$, we substitute s for $d-1$ (similarly: $s+1$ for d) to derive the general backpropagation rule:

$$\begin{aligned}\frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} L_2(\mathbf{w}) &= - \sum_D \delta_i^{\mathbf{h}_s} \cdot y_j^{\mathbf{h}_{s-1}}(\mathbf{x}) \quad // \quad \delta_i^{\mathbf{h}_s} \text{ is expanded based on the definition of } \delta_v^{\mathbf{h}_{d-1}}. \\ &= - \sum_D \underbrace{(W_{*i}^{\mathbf{h}_{s+1}})^T \boldsymbol{\delta}^{\mathbf{h}_{s+1}} \cdot y_i^{\mathbf{h}_s}(\mathbf{x}) \cdot (1 - y_i^{\mathbf{h}_s}(\mathbf{x})) \cdot y_j^{\mathbf{h}_{s-1}}(\mathbf{x})}_{\delta_i^{\mathbf{h}_s}}\end{aligned}$$

- Plugging the result for $\frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} L_2(\mathbf{w})$ into $-\eta \cdot [\dots]$ yields the update formula for $\Delta W^{\mathbf{h}_s}$. In detail:

- For updating the output matrix, $W^{\mathbf{h}_d} \equiv W^{\mathbf{o}}$, we compute

$$\boldsymbol{\delta}^{\mathbf{h}_d} = (\mathbf{c} - \mathbf{y}(\mathbf{x})) \odot \mathbf{y}(\mathbf{x}) \odot (1 - \mathbf{y}(\mathbf{x}))$$

- For updating a matrix $W^{\mathbf{h}_s}$, $1 \leq s < d$, we compute

$$\boldsymbol{\delta}^{\mathbf{h}_s} = [((W^{\mathbf{h}_{s+1}})^T \boldsymbol{\delta}^{\mathbf{h}_{s+1}}) \odot \mathbf{y}^{\mathbf{h}_s}(\mathbf{x}) \odot (1 - \mathbf{y}^{\mathbf{h}_s}(\mathbf{x}))]_{1, \dots, l_s}, \quad \text{where} \quad \begin{aligned}W^{\mathbf{h}_{s+1}} &\in \mathbf{R}^{l_{s+1} \times (l_s + 1)}, \\ \boldsymbol{\delta}^{\mathbf{h}_{s+1}} &\in \mathbf{R}^{l_{s+1}}, \\ \mathbf{y}^{\mathbf{h}_s} &\in \mathbf{R}^{l_s + 1}, \\ \text{and} \quad \mathbf{y}^{\mathbf{h}_0}(\mathbf{x}) &\equiv \mathbf{x}.\end{aligned}$$

Remarks (derivation of $\nabla^{\mathbf{h}_s} L_2(\mathbf{w})$) : (continued)

□ Hints:

$$(1) \quad y_u(\mathbf{x}) = \left[\sigma \left(W^{\mathbf{h}_d} \mathbf{y}^{\mathbf{h}_{d-1}}(\mathbf{x}) \right) \right]_u = \sigma \left(W_{u*}^{\mathbf{h}_d} \mathbf{y}^{\mathbf{h}_{d-1}}(\mathbf{x}) \right)$$

(2) Chain rule with $\frac{d}{dz}\sigma(z) = \sigma(z) \cdot (1 - \sigma(z))$, where $\sigma(z) := y_u(\mathbf{x})$ and $z = W_{u*}^{\mathbf{h}_d} \mathbf{y}^{\mathbf{h}_{d-1}}(\mathbf{x})$:

$$\frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} y_u(\mathbf{x}) \equiv \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} \left(\sigma \left(W_{u*}^{\mathbf{h}_d} \mathbf{y}^{\mathbf{h}_{d-1}}(\mathbf{x}) \right) \right) \equiv \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} (\sigma(z)) = y_u(\mathbf{x}) \cdot (1 - y_u(\mathbf{x})) \cdot \frac{\partial}{\partial w_{ij}^{\mathbf{h}_s}} \left(W_{u*}^{\mathbf{h}_d} \mathbf{y}^{\mathbf{h}_{d-1}}(\mathbf{x}) \right)$$

Note that in the partial derivative expression the symbol \mathbf{x} is a constant, while $w_{ij}^{\mathbf{h}_s}$ is the variable whose effect on the change of the loss L_2 (at input \mathbf{x}) is computed.

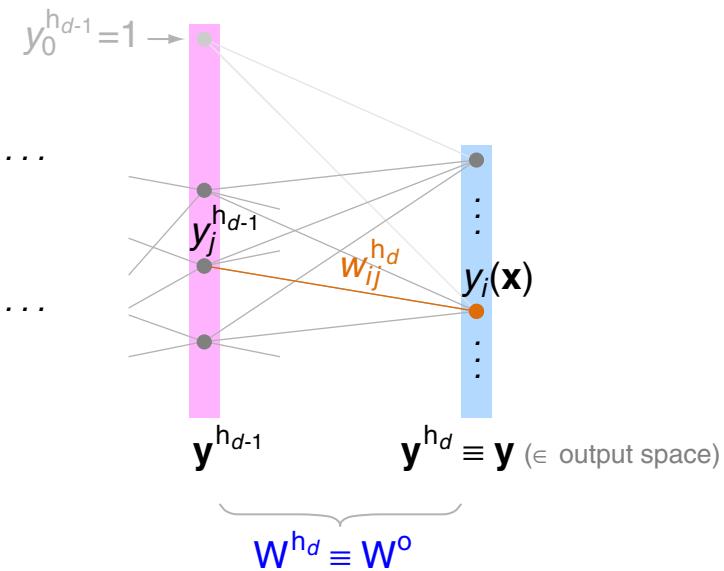
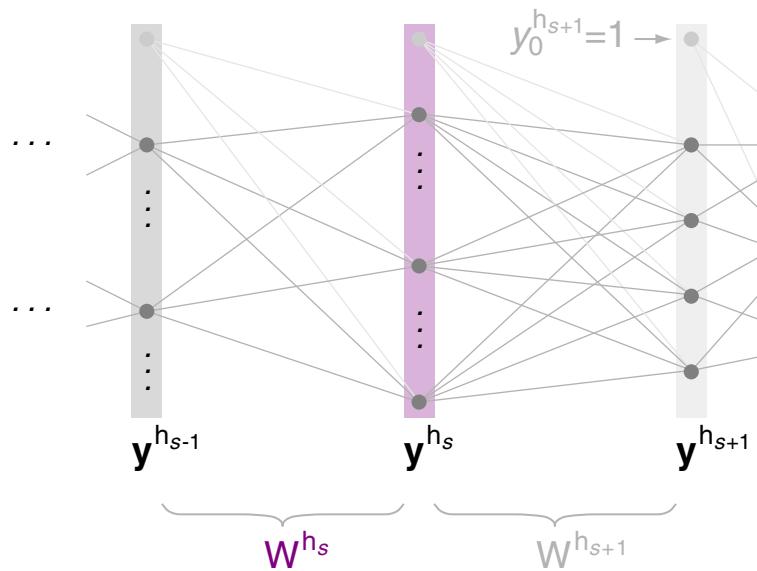
(3) $W_{u*}^{\mathbf{h}_d} \mathbf{y}^{\mathbf{h}_{d-1}}(\mathbf{x}) = w_{u0}^{\mathbf{h}_d} \cdot y_0^{\mathbf{h}_{d-1}}(\mathbf{x}) + \dots + w_{uj}^{\mathbf{h}_d} \cdot y_j^{\mathbf{h}_{d-1}}(\mathbf{x}) + \dots + w_{ul_{d-1}}^{\mathbf{h}_d} \cdot y_{l_{d-1}}^{\mathbf{h}_{d-1}}(\mathbf{x})$,
where $l_{d-1} = \text{no._rows}(W^{\mathbf{h}_{d-1}})$.

(4) Rearrange sums to reflect the nested dependencies that develop naturally from the backpropagation. We now can define $\delta_v^{\mathbf{h}_{d-1}}$ in layer $d-1$ as a function of $\delta^{\mathbf{h}_d}$ (layer d).

$$(5) \quad \sum_{u=1}^k \delta_u^{\mathbf{h}_d} \cdot w_{uv}^{\mathbf{h}_d} = (W_{*v}^{\mathbf{h}_d})^T \boldsymbol{\delta}^{\mathbf{h}_d} \quad (\text{scalar product}).$$

Remarks (derivation of $\nabla^{h_s} L_2(\mathbf{w})$) : (continued)

- The figures below show $y(x)$ as a function of some $w_{ij}^{h_s}$ exemplary in the output layer W^o and some middle layer W^{h_s} . To calculate the partial derivative of $y_u(x)$ with respect to $w_{ij}^{h_s}$, one has to determine those terms in $y_u(x)$ that depend on $w_{ij}^{h_s}$, which are shown orange here. All other terms are in the role of constants.



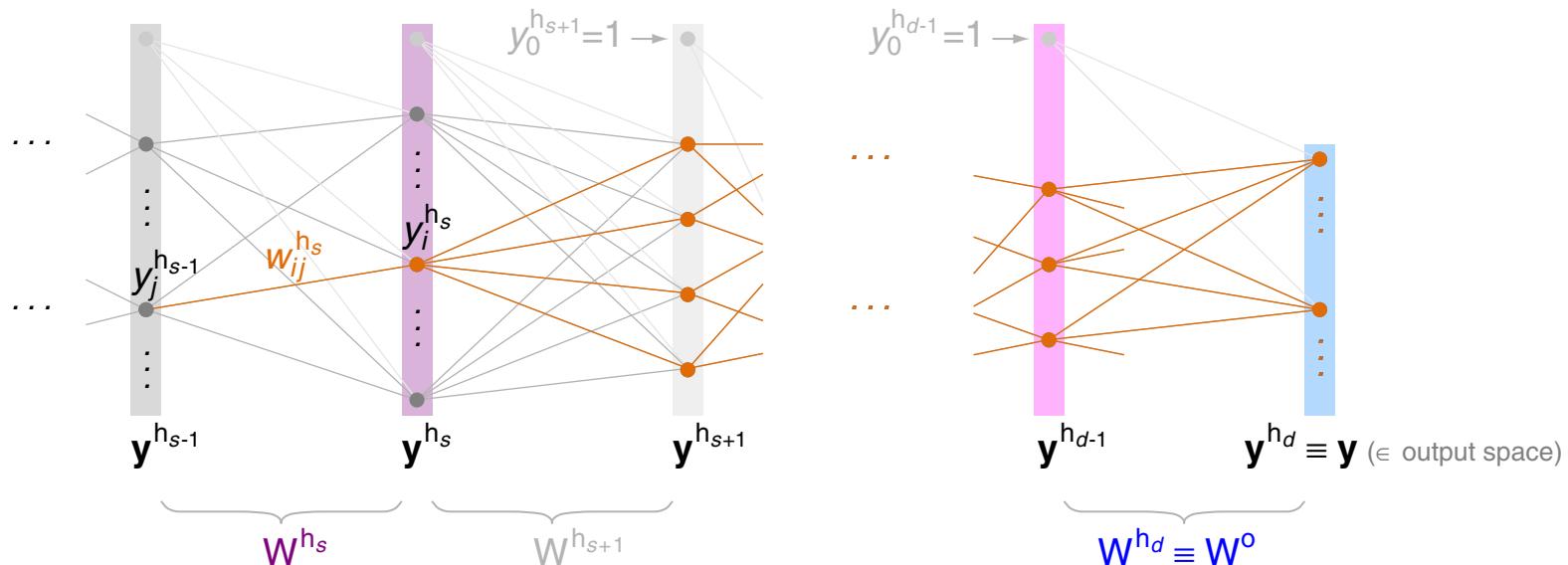
$$y_u(\mathbf{x}) = \left[\sigma \left(W^{h_d} \left(\sigma \left(\dots \left(\sigma \left(W^{h_{s+1}} \left(\sigma \left(W^{h_s} \mathbf{y}^{h_{s-1}}(\mathbf{x}) \right) \right) \right) \dots \right) \right) \right) \right)_u \right]$$

$$\begin{aligned} \sigma(\dots) &\sim \mathbf{y}^{h_d}(\mathbf{x}) \equiv \mathbf{y}(\mathbf{x}) \\ ((\dots)) &\sim \mathbf{y}^{h_{d-1}}(\mathbf{x}) \\ ((\dots)) &\sim \mathbf{y}^{h_{s+1}}(\mathbf{x}) \\ ((\dots)) &\sim \mathbf{y}^{h_s}(\mathbf{x}) \end{aligned}$$

- Compare the above illustration to the multilayer perceptron network architecture.

Remarks (derivation of $\nabla^{h_s} L_2(\mathbf{w})$) : (continued)

- The figures below show $y(x)$ as a function of some $w_{ij}^{h_s}$ exemplary in the output layer W^o and some middle layer W^{h_s} . To calculate the partial derivative of $y_u(x)$ with respect to $w_{ij}^{h_s}$, one has to determine those terms in $y_u(x)$ that depend on $w_{ij}^{h_s}$, which are shown orange here. All other terms are in the role of constants.



$$y_u(\mathbf{x}) = \left[\sigma \left(W^{h_d} \left(\sigma \left(\dots \left(\sigma \left(W^{h_{s+1}} \left(\sigma \left(W^{h_s} \mathbf{y}^{h_{s-1}}(\mathbf{x}) \right) \right) \right) \dots \right) \right) \right) \right)_u \right]$$

$$\begin{aligned} \sigma(\dots) &\sim \mathbf{y}^{h_d}(\mathbf{x}) \equiv \mathbf{y}(\mathbf{x}) \\ ((\dots)) &\sim \mathbf{y}^{h_{d-1}}(\mathbf{x}) \\ ((\dots)) &\sim \mathbf{y}^{h_{s+1}}(\mathbf{x}) \\ ((\dots)) &\sim \mathbf{y}^{h_s}(\mathbf{x}) \end{aligned}$$

- Compare the above illustration to the multilayer perceptron network architecture.

Remarks (derivation of $\nabla^o L_2(\mathbf{w})$ and $\nabla^h L_2(\mathbf{w})$ for MLP at depth one) :

- $\nabla^o L_2(\mathbf{w}) \equiv \nabla^{h_d} L_2(\mathbf{w})$, and hence $\delta^o \equiv \delta^{h_d}$.
- $\nabla^h L_2(\mathbf{w})$ is a special case of the s -layer case, and we obtain δ^h from δ^{h_s} by applying the following identities:
 $W^{h_{s+1}} = W^o$, $\delta^{h_{s+1}} = \delta^{h_d} = \delta^o$, $\mathbf{y}^{h_s} = \mathbf{y}^h$, and $l_s = l$.