

Artificial Neural Networks

Artificial Intelligence for Digital Health (AID)

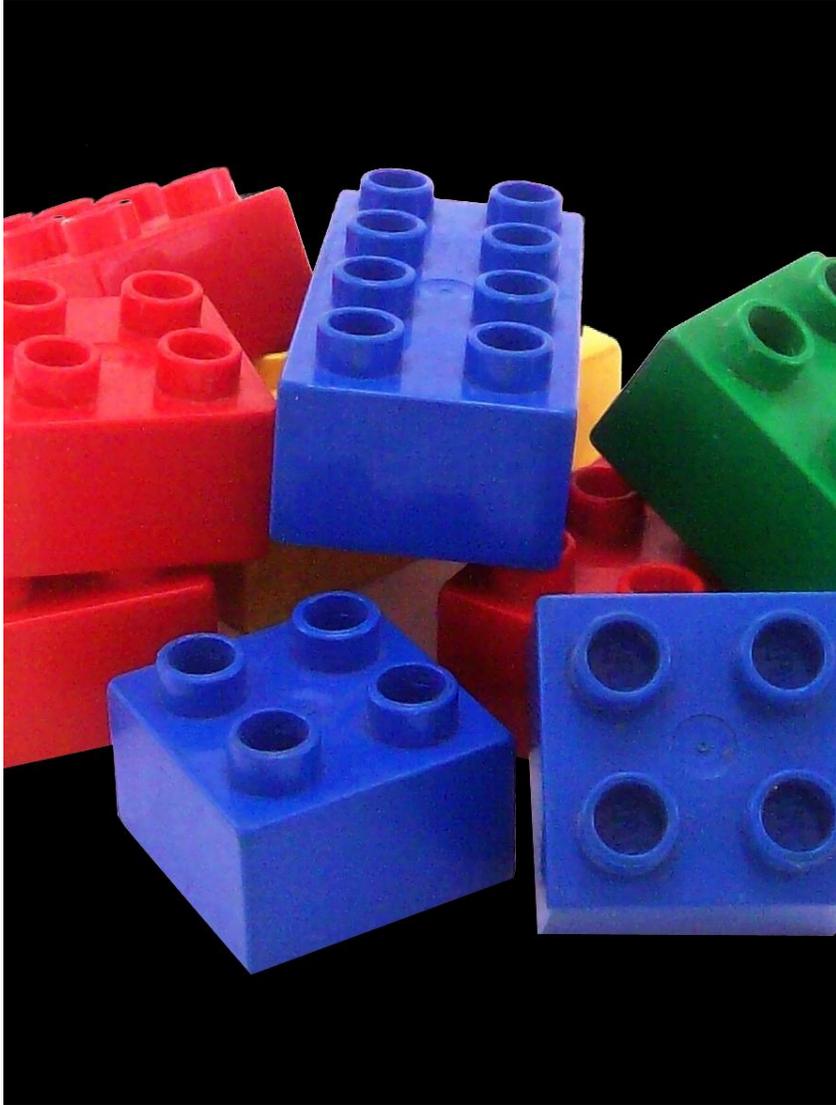
M.Sc. in Digital Health – University of Pisa

Davide Bacciu (davide.bacciu@unipi.it)



Lecture Outline

- Modeling the artificial neuron
- Artificial neural networks and the multilayer perceptron
 - Layered structure
 - Activation functions
 - Outputs and losses
- Training Artificial neural networks
 - Backpropagation algorithm
 - Loss optimization
 - Some basic tricks



Why Using Neural Networks?

- They are universal function approximators
- Non-parametric
- Scalable
- Heavily supported at SW and HW
- Fast solutions by Lego bricking

Neural Networks in Healthcare

Learn **complex patterns** from noisy and highly diverse data

Handle **vast amounts of data** in a scalable way, possibly in real time

Provide highly **accurate predictions**

When to be careful in using them?

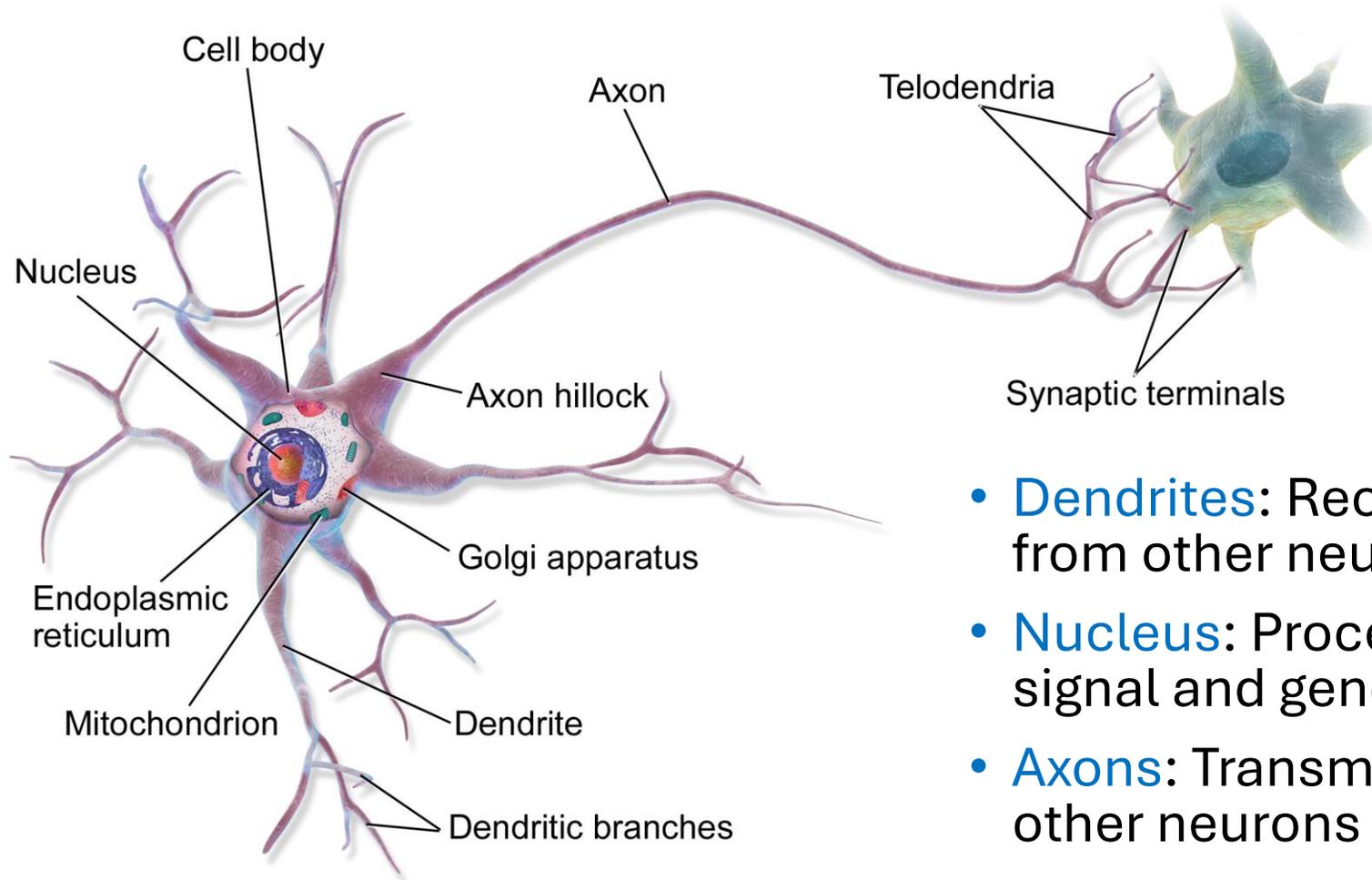
- Stringent interpretability and safety/security requirements
- Need to incorporate background/prior knowledge available
- Little-data (as opposed to Big-Data)



Essentially in our biomedical setting!

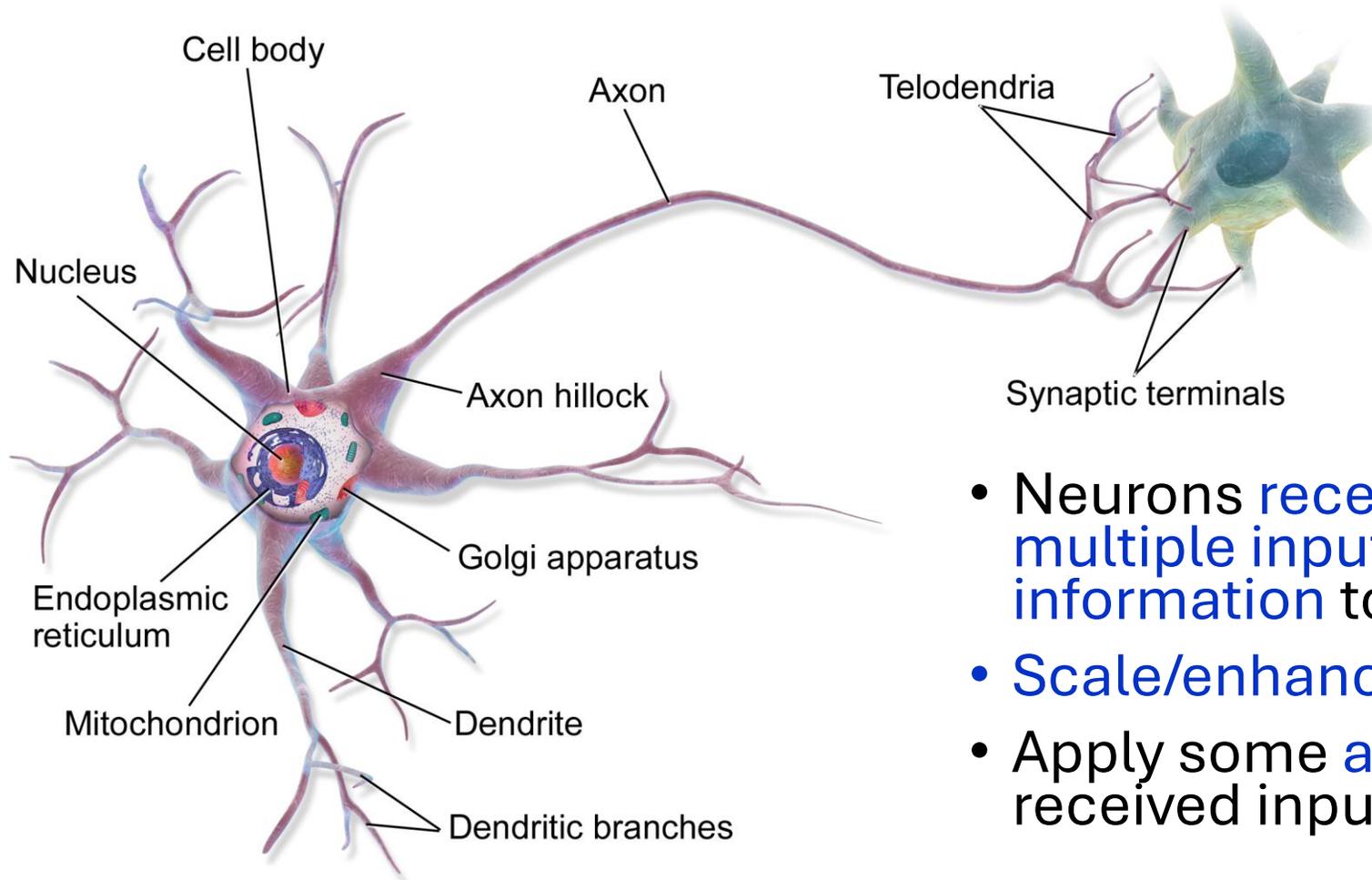
Modeling the artificial neuron

The Biological Neuron



- **Dendrites:** Receive electrical signals from other neurons
- **Nucleus:** Processes the electrical signal and generates an output signal
- **Axons:** Transmit the output signal to other neurons

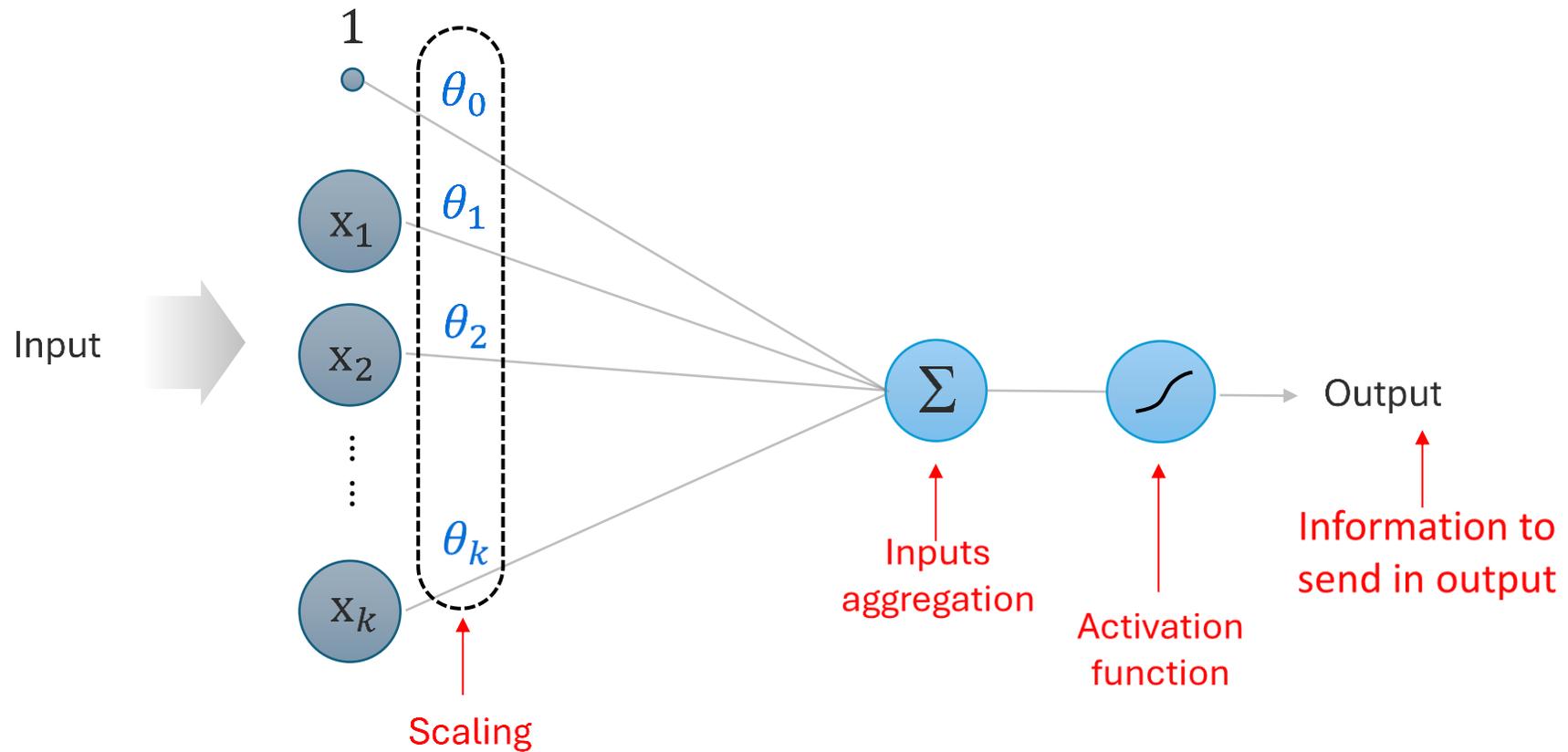
The Neuron Metaphor

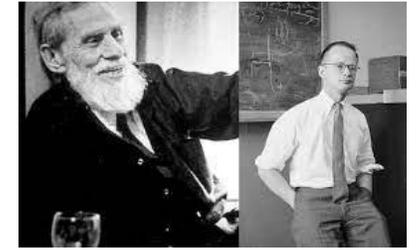


Looks like something we have already seen!

- Neurons receive information from multiple inputs and transmit information to other neurons
- Scale/enhance inputs
- Apply some activation function to received input information

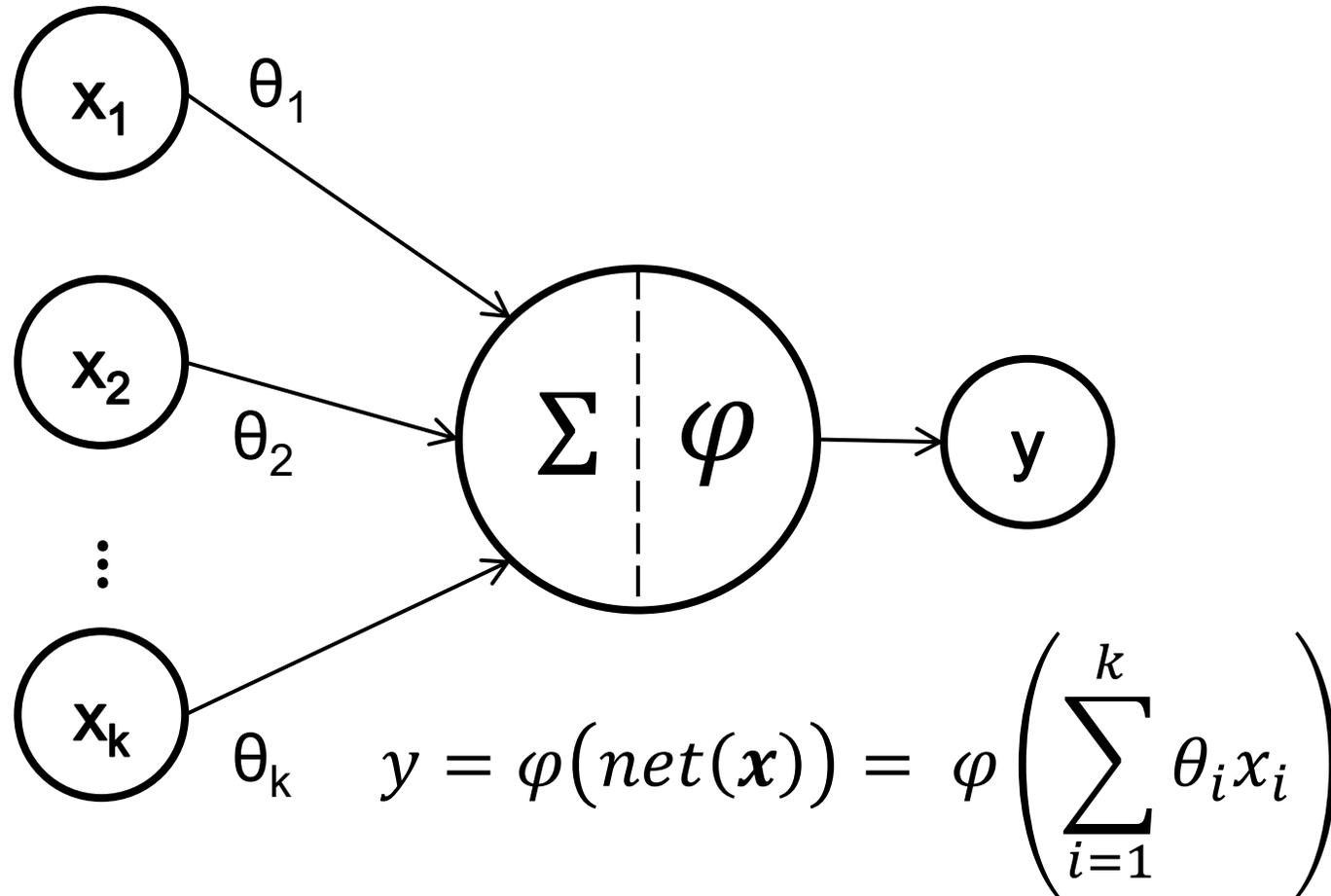
The Return of the Logistic Regression





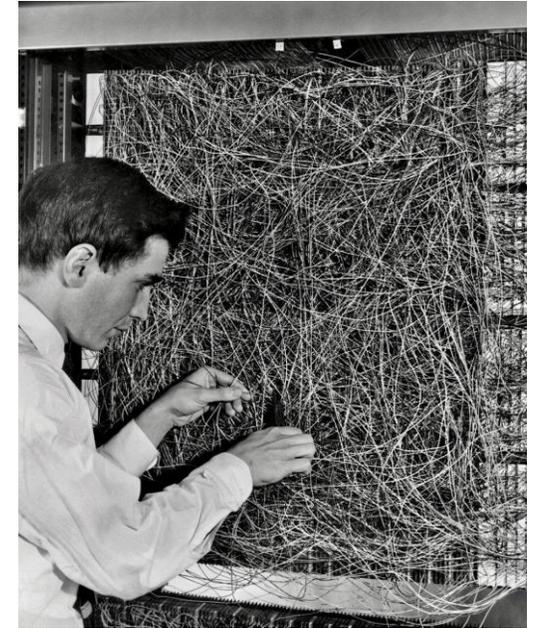
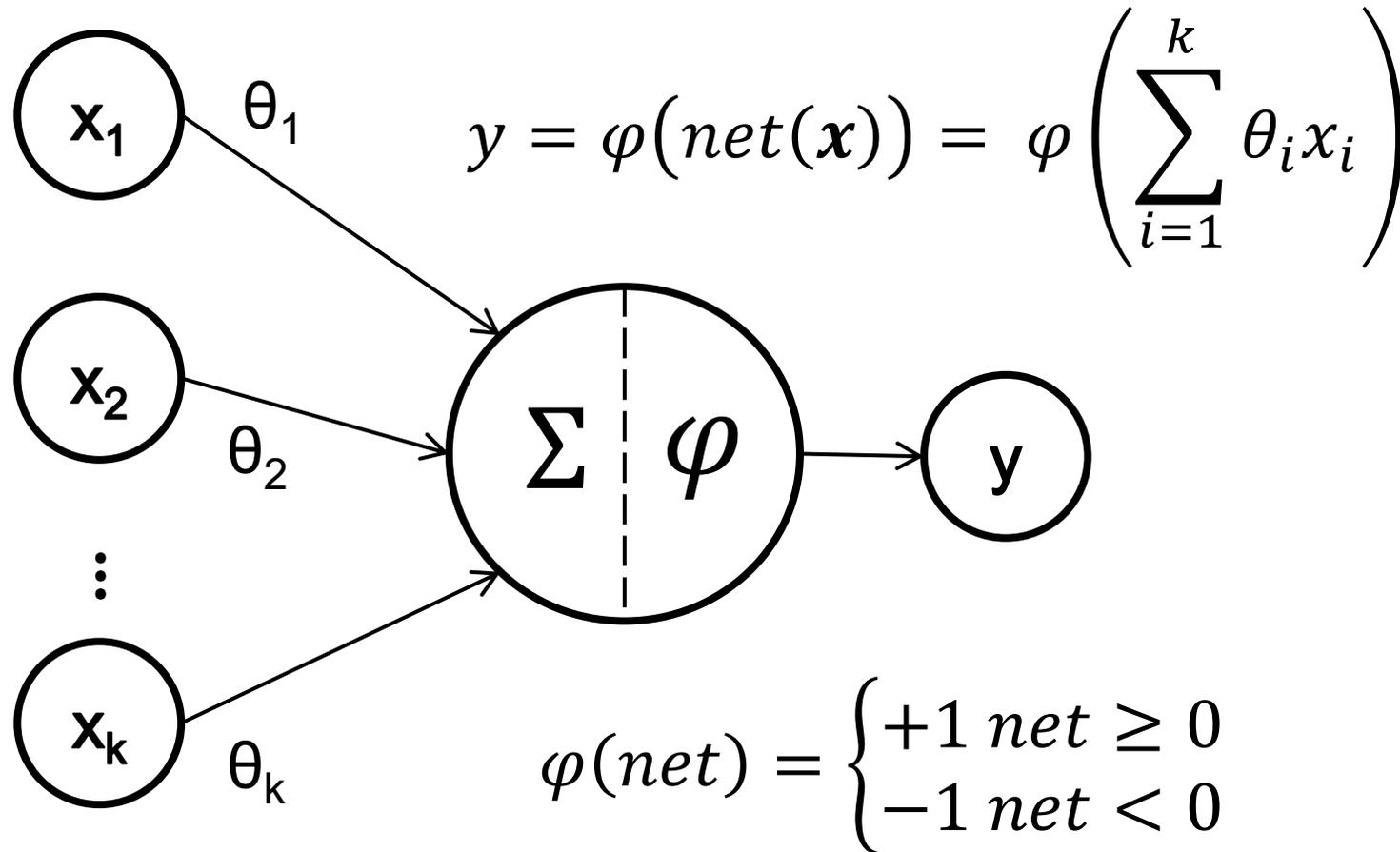
McCulloch & Pitts,
1943

The Artificial Neuron



- Input \mathbf{x}
 - Synaptic weights $\boldsymbol{\theta}$
 - Local potential
- $$\text{net}(\mathbf{x}) = \sum_{i=1}^k \theta_i x_i$$
- Activation function φ
 - Output \mathbf{y}

The Perceptron



Rosenblatt, 1953

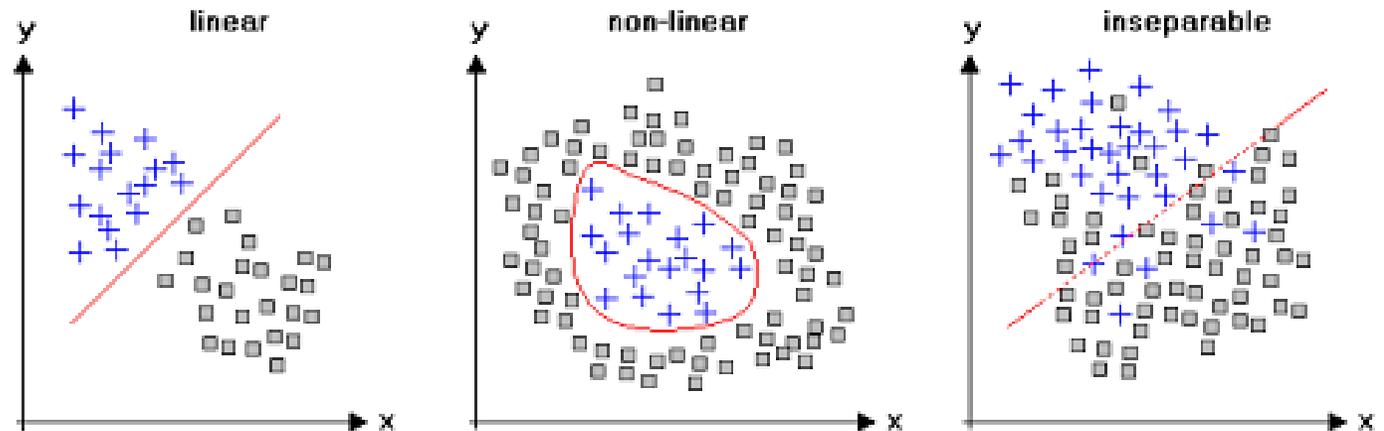
The first **learnable neuron** (the delta-rule)

The Perceptron → The Issue

Still too similar to a logistic regression

$$y = \varphi(\text{net}(\mathbf{x})) = \sigma \left(\sum_{i=1}^k \theta_i x_i \right)$$

With its limitations



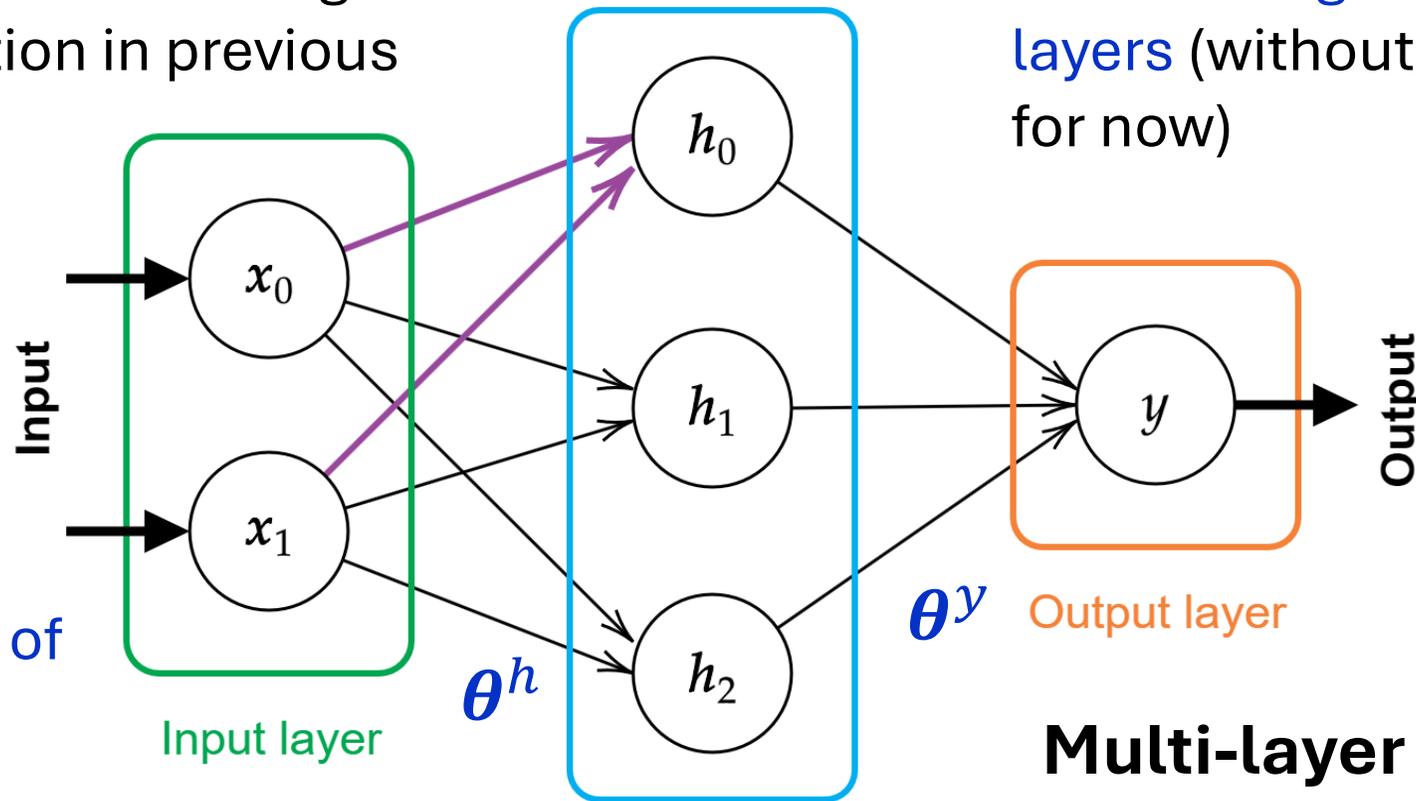
Artificial Neural Networks

Artificial neural networks – Interconnected neurons

Each neuron following the equation in previous slides

Neurons **organized in layers** (without cycles, for now)

Note the **multiple set of weights**



Multi-layer Perceptron

Multi-Layer Perceptron (MLP)

Three main components (**layers**):

- **Input layer**: feeds the input \mathbf{x}
- **Hidden layer**: transforms the input \mathbf{x} into a new (vectorial) representation $h(\mathbf{x})$ (size of the number of hidden neurons) depending on the hidden layer parameters θ^h
- **Output layer**: generates the prediction \mathbf{y} (output) by combining the intermediate representation of the input provided by the hidden layer and the output layer parameters θ^y

Two key aspects

- Can learn **any decision boundary** (even non-linear ones)
- Adaptable to **any type of task** by changing the output layer

MLP – Input Layer

- The input layer is simply a replica of the input data
- However, **preprocessing is crucial** for good performance
- **Feature scaling** (normalization or standardization)
 - If features have different scales, normalization is needed to ensure stable training.
- Handling **categorical features**
 - Convert categorical variables into numerical form using **one-hot encoding**
 - For large categories, **embedding layers** can be used instead

Categorical Features

- A categorical variable is a variable that can belong to **one of k discrete categories**
- Categorical variables are **usually encoded using 1-out-of-k coding (one hot)**
 - E.g. for three colors: **red** = (1 0 0), **green** = (0 1 0), **Blue** = (0 0 1)
 - I.e. the MLP above will have 3 input
- If we used red = 1, green = 2, blue = 3, then this type of encoding imposes a representational bias which is not semantically supported
- More recently **dense embeddings** have taken over (especially for natural language)

Numerical Features

- A numerical variable (**continuous, ordinal**) can be directly fed to a neural network
- However, it is good practice to normalize data so that the dynamic range of inputs is limited

- **Min-Max Scaling** - Scales values to $[0,1]$ or $[-1,1]$

$$x' = \frac{x - x_{min}}{x_{max} - x_{min}} \rightarrow x' \in [0,1] \text{ or } x' = 2 \frac{x - x_{min}}{x_{max} - x_{min}} - 1 \rightarrow x' \in [-1,1]$$

- **Z-score Standardization**- Centers data around mean 0 and std 1

$$x' = \frac{x - \mu}{\sigma}$$

- Need also care in choosing w.r.t. whom you are normalizing **population** or **individual** (especially in biomedical applications)
- Normalization is applied also in **timeseries and graphs**

MLP – Hidden Layer

The hidden layer transforms the input \mathbf{x} into another vector \mathbf{h} of arbitrary size

- Size of \mathbf{h}
 - \mathbf{h} can be smaller or larger than \mathbf{x} , depending on the complexity of the task
 - More neurons \rightarrow More capacity to learn complex patterns (but risk of overfitting)
 - Fewer neurons \rightarrow Simpler models, better generalization
- Helps solve the task by **learning representations** that make classification, regression, or other **tasks easier**
- The **transformation from \mathbf{x} to \mathbf{h}** is governed by trainable parameters θ^h fitted during training

$$\mathbf{h} = g_{\theta^h}(\mathbf{x}) \text{ with } g_{\theta^h}(\cdot) \text{ nonlinear}$$

Hidden layer step-by-step (I)

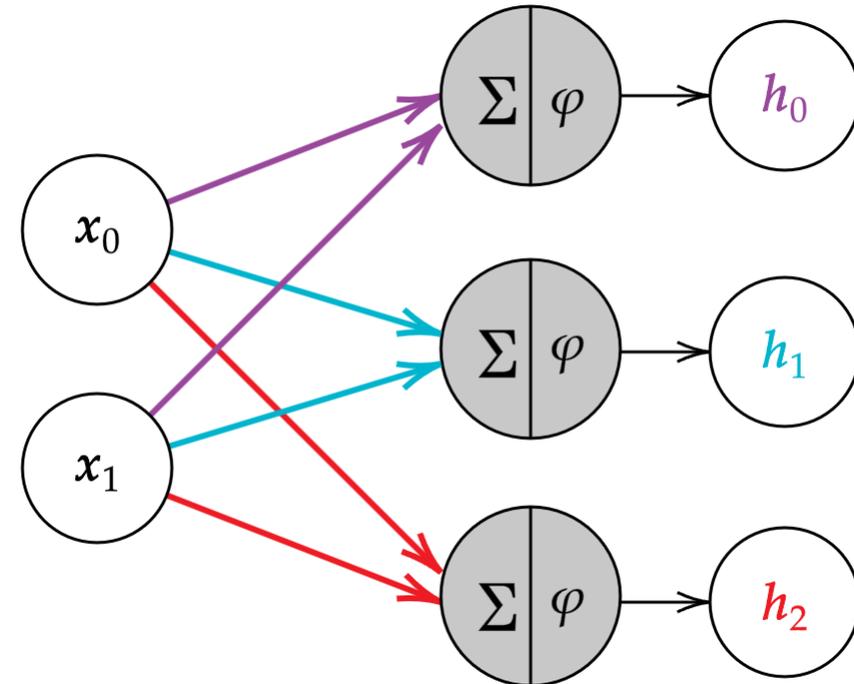
- **Input:** A vector \mathbf{x} with $k = 2$ features $\rightarrow \mathbf{x} \in \mathbb{R}^2$
- **Output:** A vector \mathbf{h} with $s = 3$ features $\rightarrow \mathbf{h} \in \mathbb{R}^3$

$$h_0 = \varphi(\boldsymbol{\theta}_0 \mathbf{x}) = \varphi(\theta_{00}x_0 + \theta_{10}x_1)$$

$$h_1 = \varphi(\boldsymbol{\theta}_1 \mathbf{x}) = \varphi(\theta_{01}x_0 + \theta_{11}x_1)$$

$$h_2 = \varphi(\boldsymbol{\theta}_2 \mathbf{x}) = \varphi(\theta_{02}x_0 + \theta_{12}x_1)$$

We transformed a 2D input vector into a 3D hidden representation by a linear combination of \mathbf{x} features + an activation function φ



Hidden layer step-by-step (II)

- θ is a **parameter matrix** with $k=2$ rows and $s=3$ columns

$$\theta = \begin{bmatrix} 2 & 3 & -1 \\ 2 & 6 & -8 \end{bmatrix}^T$$

- $\mathbf{x} = [-1 \quad 1]^T$ is a **sample input vector**

$$h_0 = \varphi \left(\begin{bmatrix} 2 & 2 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right) = \varphi(0)$$

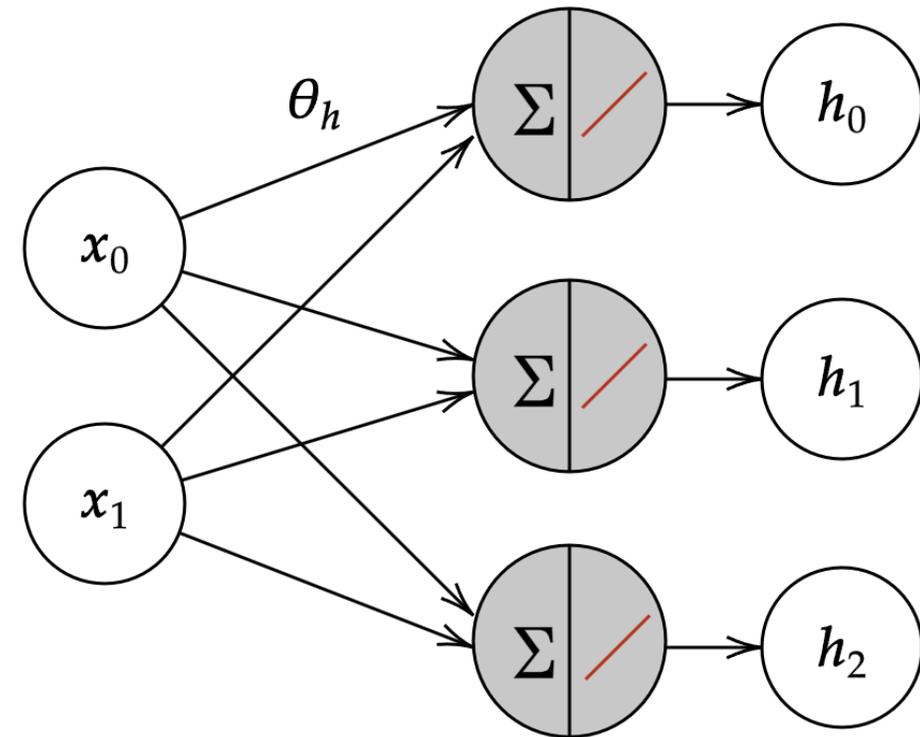
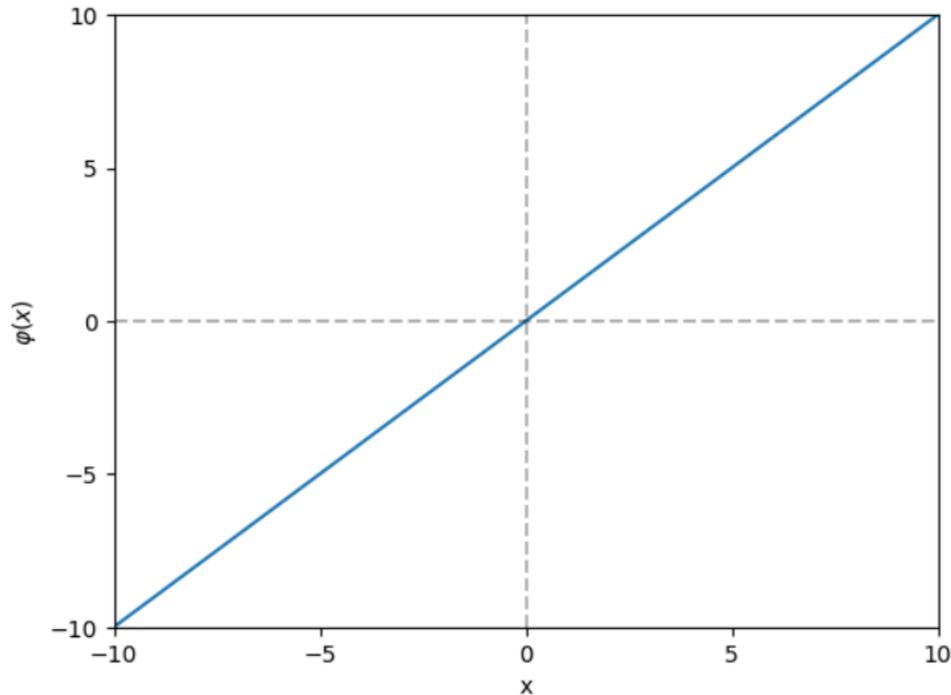
$$h_1 = \varphi \left(\begin{bmatrix} 6 & 3 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right) = \varphi(-3)$$

$$h_2 = \varphi \left(\begin{bmatrix} -8 & -1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right) = \varphi(7)$$

$$\mathbf{h} = \varphi([0 \quad -3 \quad 7])$$

Linear activation function

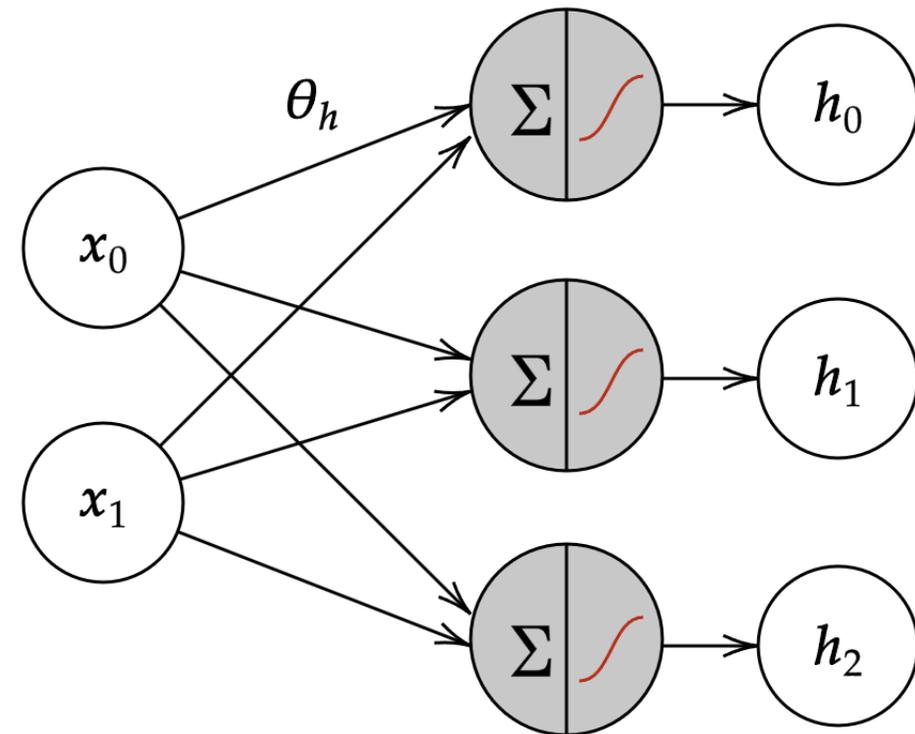
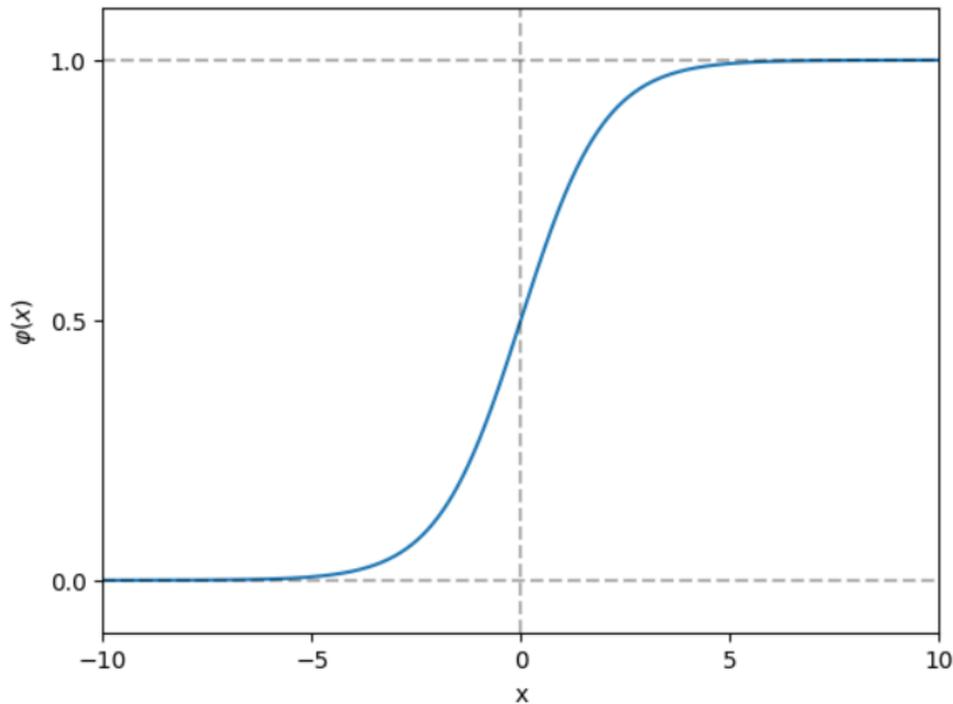
$$\varphi(\theta x) = \theta x$$



It returns the input without changing it $\mathbf{h} = \varphi([0 \ -3 \ 7]) = [0 \ -3 \ 7]$

Sigmoid activation function

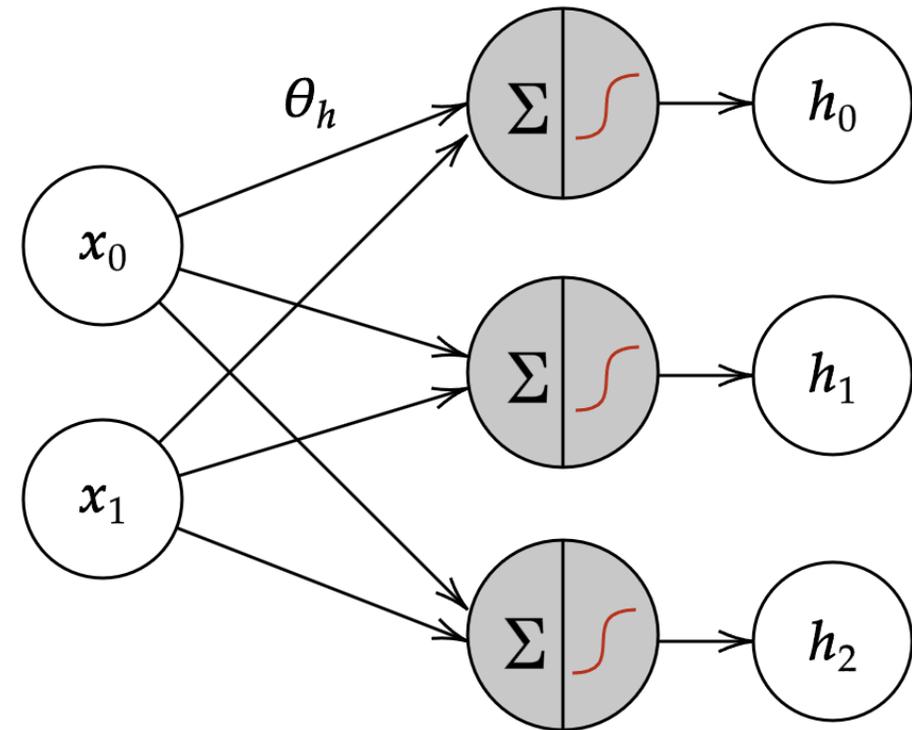
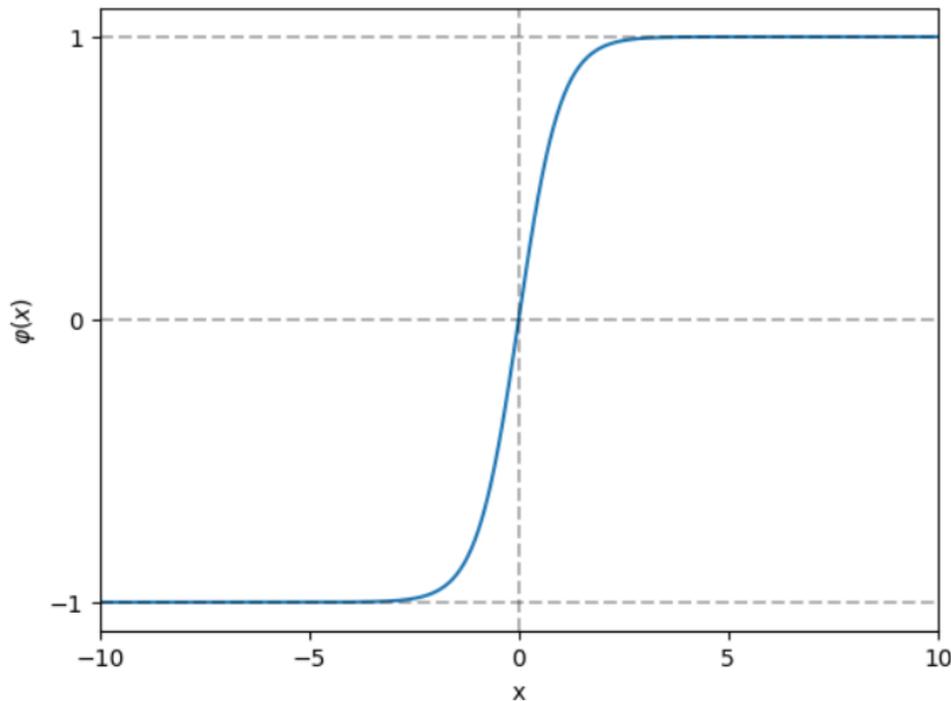
$$\varphi(\theta x) = \frac{1}{1 + e^{-\theta x}}$$



Squashes the input in $[0,1] \rightarrow \varphi([0 \ -3 \ 7]) = [0.5 \ 0.047 \ 0.999]$

Hyperbolic tangent activation function

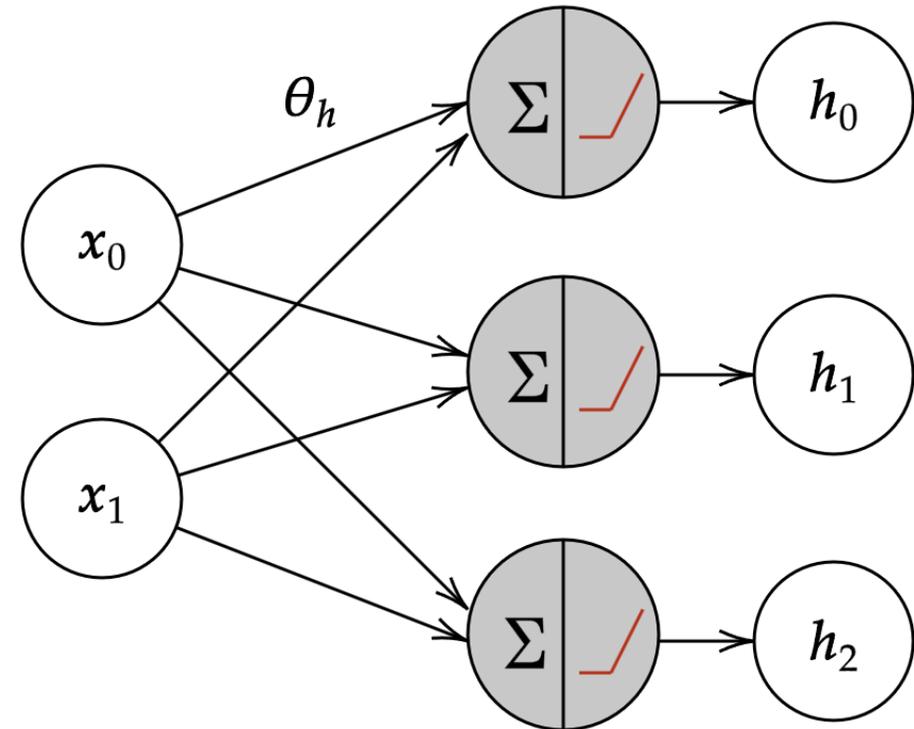
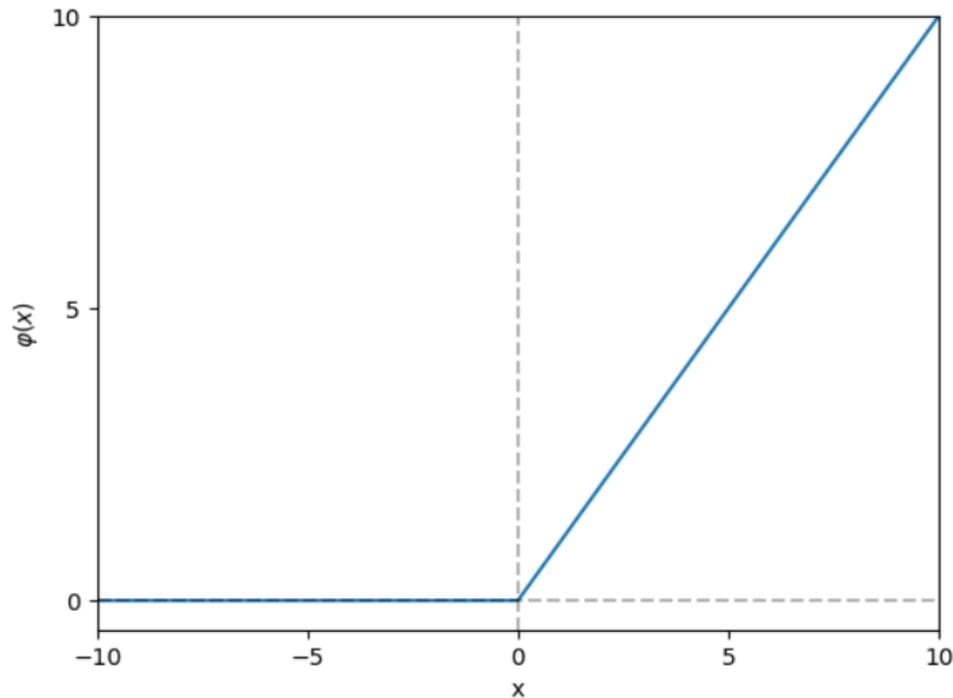
$$\varphi(\theta x) = \frac{e^{\theta x} - e^{-\theta x}}{e^{\theta x} + e^{-\theta x}}$$



Squashes the input in $[-1, 1] \rightarrow \varphi([0 \ -3 \ 7]) = [0 \ -0.995 \ 0.999]$

Rectified Linear Unit (ReLU)

$$\varphi(\theta x) = \max(0, \theta x)$$



Zeros the negative components of the input vector, leaving the rest unchanged

$$\rightarrow \varphi([0 \ -3 \ 7]) = [0 \ 0 \ 7]$$

MLP – Output Layer

The output layer transforms the hidden layer output \mathbf{h} into the prediction y

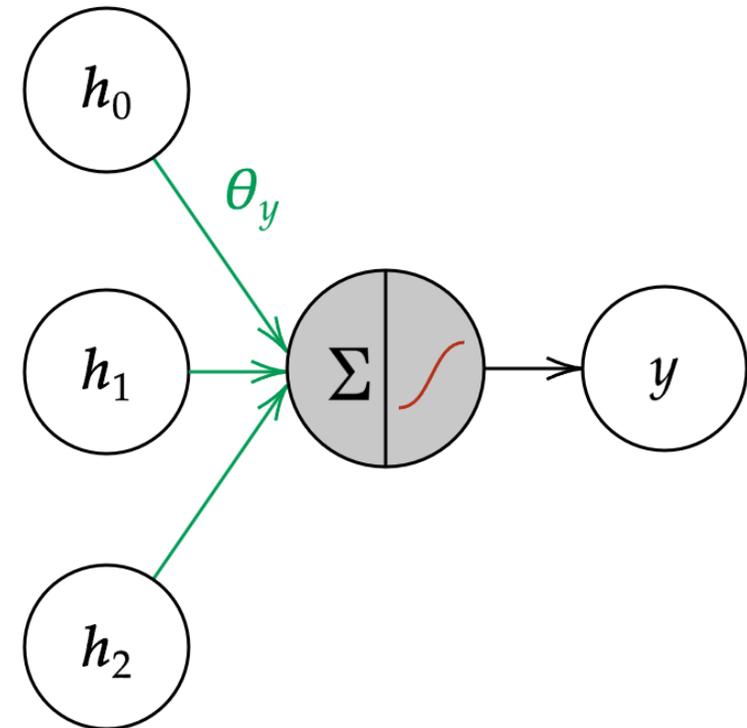
- Provides the final prediction of the MLP
- The transformation from \mathbf{h} to y is parameterized by θ^y which are learned during training
- It implements a parameterized function $y = g_{\theta^y}(\mathbf{h})$
- Computes similarly to the hidden layer, but it is **associated to activation functions which are task-specific** and are linked to a loss function

Output layer for binary classification

Goal: classify a sample \mathbf{x} into either class $y \in \{0,1\}$

- Linear combination of the inputs
- Followed by sigmoid activation

$$\hat{y} = \sigma(\boldsymbol{\theta}^y \mathbf{h})$$



Binary cross entropy loss

$$BCE(\hat{y}, y) = -y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

Output layer for multiple binary classification

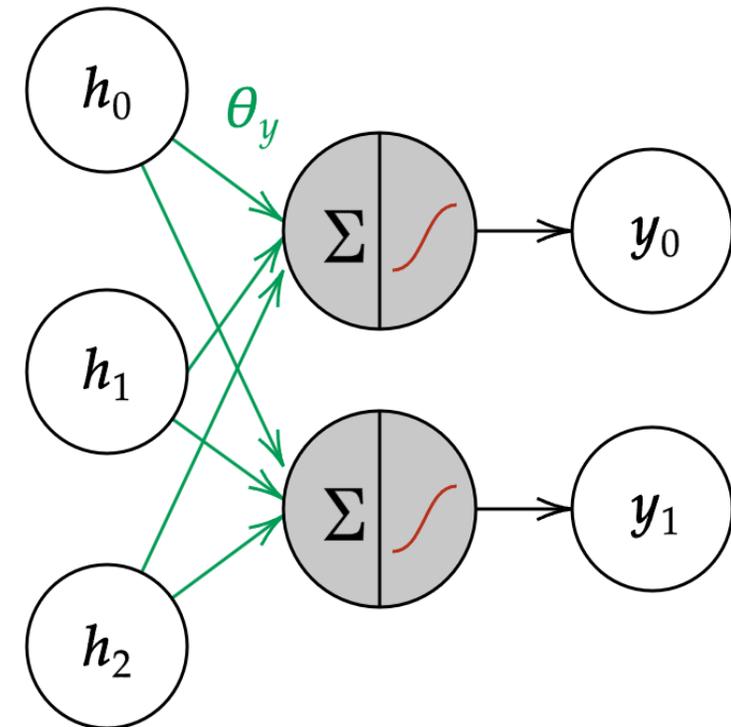
Goal: predict a d -dimensional vector \mathbf{y} of 0 and 1 for a sample \mathbf{x}

- Output vector $\mathbf{y} \in \{0,1\}^d$
- Repeated application of the sigmoid for the d -dimensions

$$\hat{\mathbf{y}} = \sigma(\boldsymbol{\theta}^y \mathbf{h})$$

Average binary cross entropy loss

$$BCE_{multi}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{d} \sum_d BCE(\hat{y}_d, y_d)$$



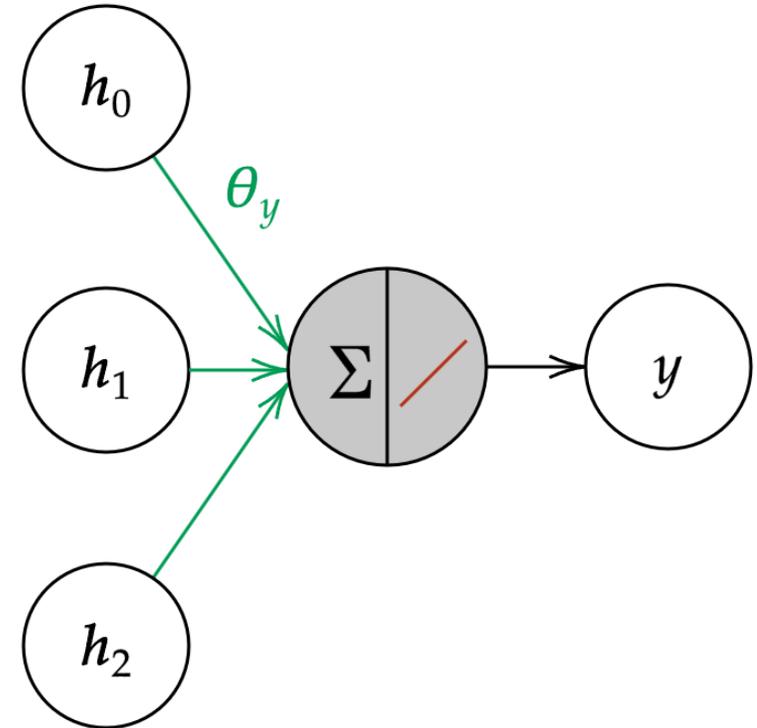
Output layer for regression

Goal: predict a real value y in response to sample \mathbf{x}

- Just a linear combination of the inputs
- Followed by linear (identity) activation
 $\hat{y} = (\boldsymbol{\theta}^y \mathbf{h})$

Mean squared error loss

$$MSE(\hat{y}, y) = (y - \hat{y})^2$$



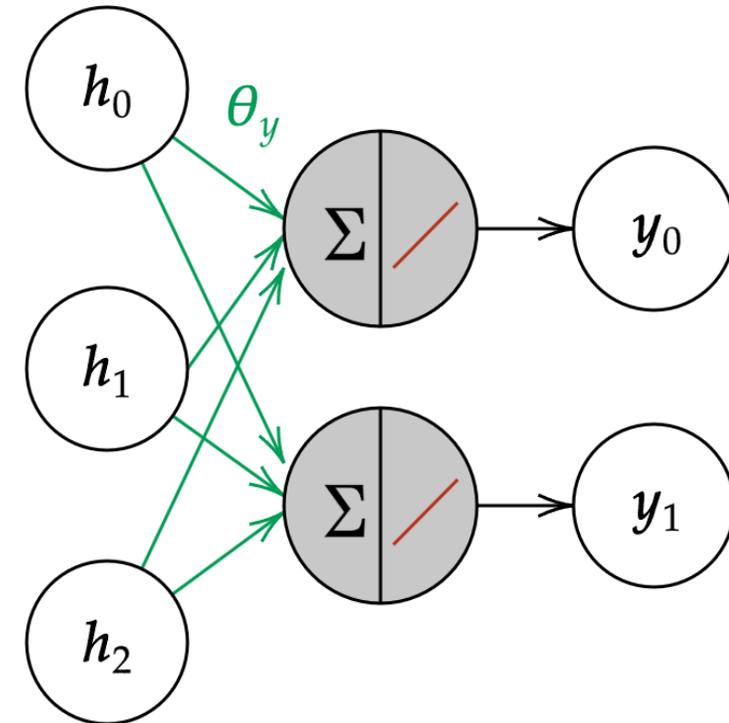
Output layer for multiple regression

Goal: predict a real valued vector $\mathbf{y} \in \mathbb{R}^d$ in response to sample \mathbf{x}

- Multiple linear combination of \mathbf{h} , one for each component of \mathbf{y}
- Followed by multiple application of identity

Mean squared error loss

$$MSE_{multi}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{d} \sum_d MSE(\hat{y}_d, y_d)$$



Output layer for multiclass classification

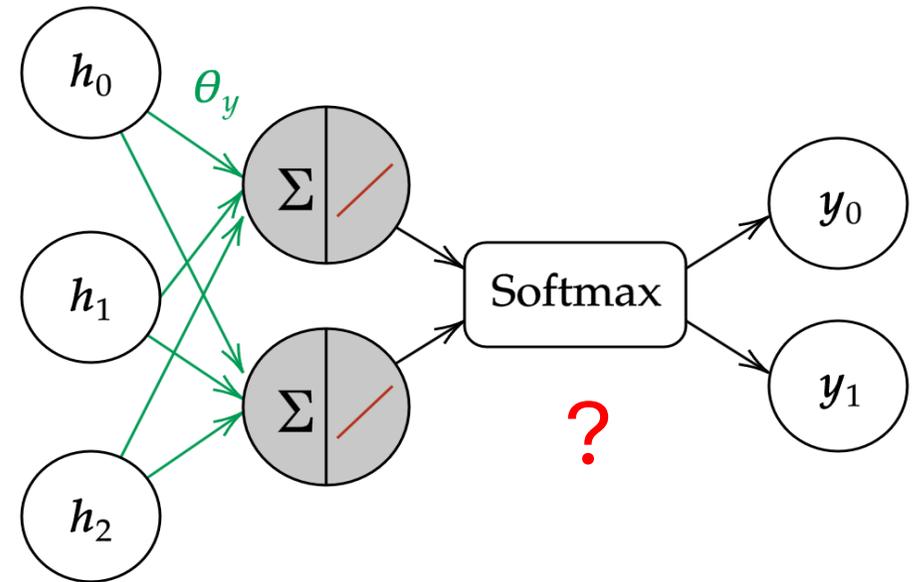
Goal: assign a sample \mathbf{x} to one of d classes $\{c_1, \dots, c_d\}$

- Output $\mathbf{y} \in \{0,1\}^d$ as a one-hot encoding of the class
- Linear combination of \mathbf{h} followed by identify
- Output generated by a softmax function

$$\hat{\mathbf{y}} = \text{softmax}(\boldsymbol{\theta}^y \mathbf{h})$$

Cross-Entropy loss

$$CE(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_d y_d \log(\hat{y}_d)$$



Softmax function

- Transforms a (dense) vector in a categorical (discrete) probability distribution
- Given a vector \mathbf{x} with n components

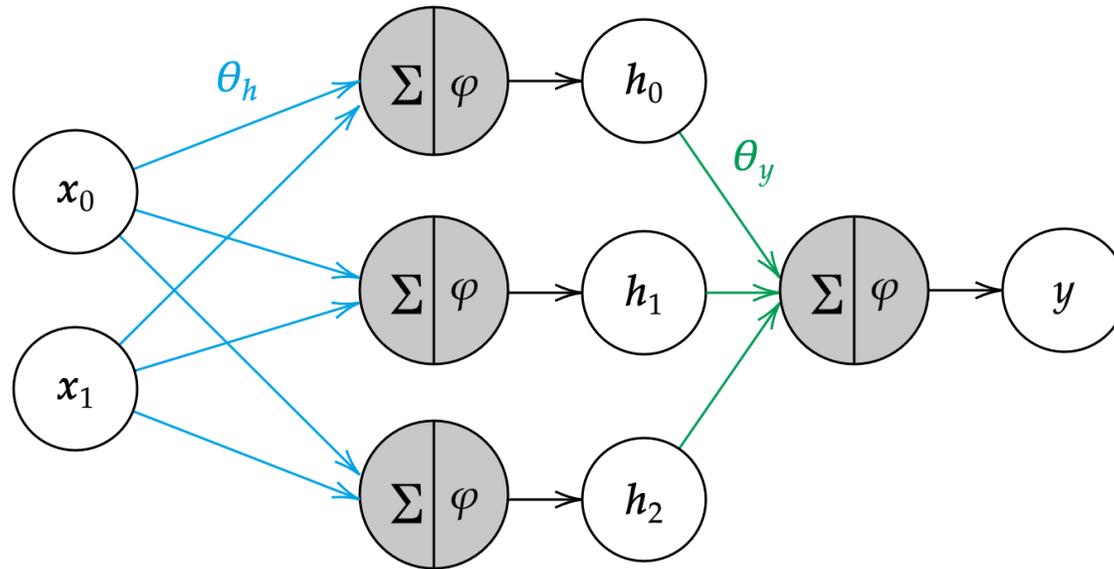
$$\text{softmax}_i(\mathbf{z}) = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

In other words

- Take the exponentiation of each component of \mathbf{z}
- Compute the sum of exponentiated components
- Divide each exponentiated component by the overall sum
- The resulting softmax vector will sum-to-1 (like a probability distribution)

$$\text{softmax}([-1 \ 0 \ 1]) = [0.09 \ 0.244 \ 0.666]$$

MLP: putting things back together



- Transform \mathbf{x} into \mathbf{h} through the hidden layer $\mathbf{h} = g_{\theta^h}(\mathbf{x}) = \varphi(\boldsymbol{\theta}^h \mathbf{x})$
- Transform \mathbf{h} into y through the output layer $y = g_{\theta^y}(\mathbf{h}) = \varphi(\boldsymbol{\theta}^y \mathbf{h})$
- The model is $MLP_{\theta}(\mathbf{x}) = g_{\theta^y}(g_{\theta^h}(\mathbf{x}))$ with $\theta = \{\theta^h, \theta^y\}$ being learned parameters

What do we need to build an MLP?

- Input preprocessing
- Size of the hidden layer s
- Hidden layer activation functions φ
- Configuration of the output layer (guided by the task)

} Model selection choices

Training an Artificial Neural Network

Let's recall logistic regression training

For a certain number of iterations (**epochs**):

For each training pair $(\mathbf{x}^i, y^i) \in D_{train}$:

1. Compute the prediction $h_{\theta}(\mathbf{x}^i) = \sigma(\mathbf{x}^i \boldsymbol{\theta})$
2. Compute the loss L of the prediction $h_{\theta}(\mathbf{x}^i)$ compared to the true label y^i
3. Compute the gradient of the loss $\nabla_{\theta} L = \mathbf{x}(y^i - h_{\theta}(\mathbf{x}^i))$
4. Update the parameters $\boldsymbol{\theta}_{new} = \boldsymbol{\theta} - \eta \nabla_{\theta} L$
5. Use the updated parameters in the next iteration $\boldsymbol{\theta} = \boldsymbol{\theta}_{new}$

MLP training

For a certain number of iterations (**epochs**):

For each training pair $(\mathbf{x}^i, y^i) \in D_{train}$:

1. Compute the prediction $MLP_{\theta}(\mathbf{x}^i)$
2. Compute the loss L of $MLP_{\theta}(\mathbf{x}^i)$ compared to the true label y^i
3. Compute the gradient of the loss $\nabla_{\theta}L$
4. Update the parameters $\theta_{new} = \theta - \eta \nabla_{\theta}L$
5. Use the updated parameters in the next iteration $\theta = \theta_{new}$

Actually, it is a bit more articulated than that

MLP training, more realistically

For a certain number of iterations (**epochs**):

Shuffle D_{train}

For each **subsample B** of pairs (\mathbf{x}^B, y^B) extracted from D_{train} :

1. Compute the prediction $MLP_{\theta}(\mathbf{x}^B)$
2. Compute the loss L of $MLP_{\theta}(\mathbf{x}^B)$ compared to the **true labels y^B**
3. Compute the gradient of the **loss $\nabla_{\theta}L$**
4. Update the parameters $\theta_{new} = \theta - \eta \nabla_{\theta}L$
5. Use the updated parameters in the next iteration $\theta = \theta_{new}$

The **gradient of the loss needs to be computed w.r.t. all parameters**, including those from the hidden layer!

Minibatching

Mini-batching is a **technique used to split the dataset into smaller subsets**, each containing B examples. Instead of updating the parameters after processing all examples or after each single example, updates are done after processing each mini-batch

1. Divide the dataset into multiple mini-batches of size B (**batch size, hyperparameter**)
2. Shuffle the dataset at the beginning of each epoch to avoid bias in training.
3. Update model parameters once per mini-batch.

- **Why is Mini-Batching Useful?** It is a compromise between:
 - Stochastic Gradient Descent (SGD) ($B = 1$): Updates after each example \rightarrow fast but noisy updates.
 - Batch Gradient Descent ($B = n$): Updates after the entire dataset \rightarrow stable but slow updates.
- Mini-batching **provides a balance**
 - Faster convergence than full batch updates
 - More stable updates than SGD
 - Efficient computation by leveraging GPU parallelism

Gradient of the loss

We have different loss functions based on the output layer, and we can compute their gradients as follows

Loss Function	Formula	Gradient
Binary Cross-Entropy (BCE)	$-[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$	$x(\hat{y} - y)$
Mean Squared Error (MSE)	$(\hat{y} - y)^2$	$x(\hat{y} - y)$
Categorical Cross-Entropy (CE)	$-\sum_{k=1}^K y_k \log(\hat{y}_k)$	$\hat{y} - y$

We can **update the output layer parameters** (θ^y) using these gradients

However, we don't yet know **how to compute the gradient for the hidden layer parameters** (θ^h)

Backpropagation at the rescue

Key Idea: We apply the Chain Rule to propagate the gradient from the output layer to the hidden layer

1. Compute the gradient for the output layer

- Example: For BCE, the gradient is

$$\nabla_{\theta^y} L = x(\hat{y} - y)$$

2. Propagate the gradient back to the hidden layer

- Using the Chain Rule, the gradient for the hidden layer parameters θ^h is:

$$\nabla_{\theta^h} L = \nabla_h L \cdot \nabla_{\theta^h} h$$

- $\nabla_h L$ is the gradient of the loss with respect to the hidden layer output \mathbf{h}
- $\nabla_{\theta^h} h$ is the gradient of \mathbf{h} with respect to θ^h

3. Compute the gradient of \mathbf{h}

- If the hidden layer transformation is:

$$\mathbf{h} = \varphi(\boldsymbol{\theta}^h \mathbf{x})$$

- Then, using the chain rule:

$$\nabla_{\theta^h} L = \nabla_h L \cdot (\varphi'(\boldsymbol{\theta}^h \mathbf{x})) \cdot \mathbf{x}$$

- Where φ' is the derivative of the activation function

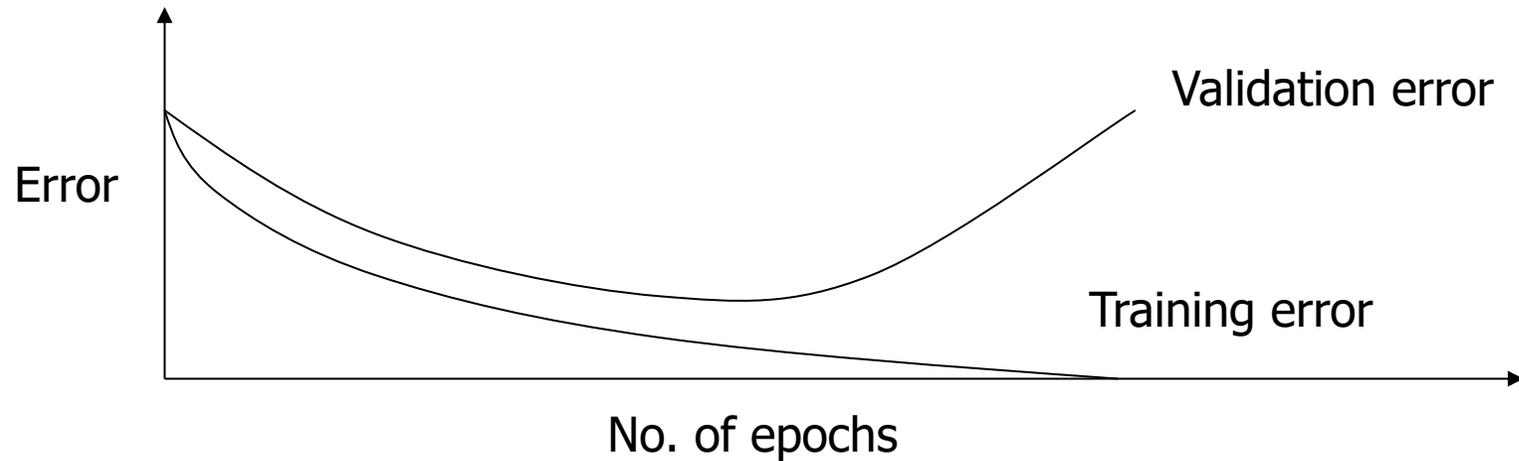
Luckily for you
backpropagation is
automatically
handled by the NN
libraries

Convergence Criteria

- Learning is obtained by iteratively supplying shuffled training data and adjusting by backpropagation
 - Typically, 1 training set presentation = 1 epoch
- We need a stopping criteria to define convergence
 - Validation for generalization performance: stop when generalization performance reaches a peak

Early Stopping

- Keep a hold-out validation set and assess accuracy after (every/some) epoch.
- Maintain weights for best performing network on the validation set and stop training when error increases beyond this
- Always let the network run for some epochs before deciding to stop (**patience parameter**), then backtrack to best result



Wrap-up

Take home lessons

- Artificial Neural Networks are **universal function approximators** that can learn complex patterns from noisy and diverse data
 - In healthcare they allow handling vast amounts of heterogenous data in a scalable way, providing highly accurate predictions
- The artificial neuron is loosely inspired by biological neurons
 - Artificial neurons receive and aggregate inputs, apply activation functions, and transmit outputs
- The MLP has a **layered structure** consisting of input, hidden, and output layers
 - Capable of learning any decision boundary and adaptable to various tasks
 - Several architectural/model selection choices need to be taken even for simple networks
- Training ANNs involves a **specialized gradient descent algorithm** (backpropagation)
 - Optimization aspects needs to be carefully curated
 - Stopping conditions are relevant as it is an iterative process
- Optimize neural networks from the start
 - Applying the right **weight initialization strategies** significantly impacts effectiveness
 - Needs to be tailored to the activation function

Next Lecture

- Next week lecture
 - [Risk scoring and stratification](#) for population health management
 - Machine learning for risk stratification
 - Identifying risk factors
 - Assessment and validation of risk predictors
 - Censoring and its impact in risk scoring