# Lecture 2: Learning with neural networks

Deep Learning @ UvA

# Lecture Overview

- Machine Learning Paradigm for Neural Networks

- The Backpropagation algorithm for learning with a neural network

- Neural Networks as modular architectures

- Various Neural Network modules

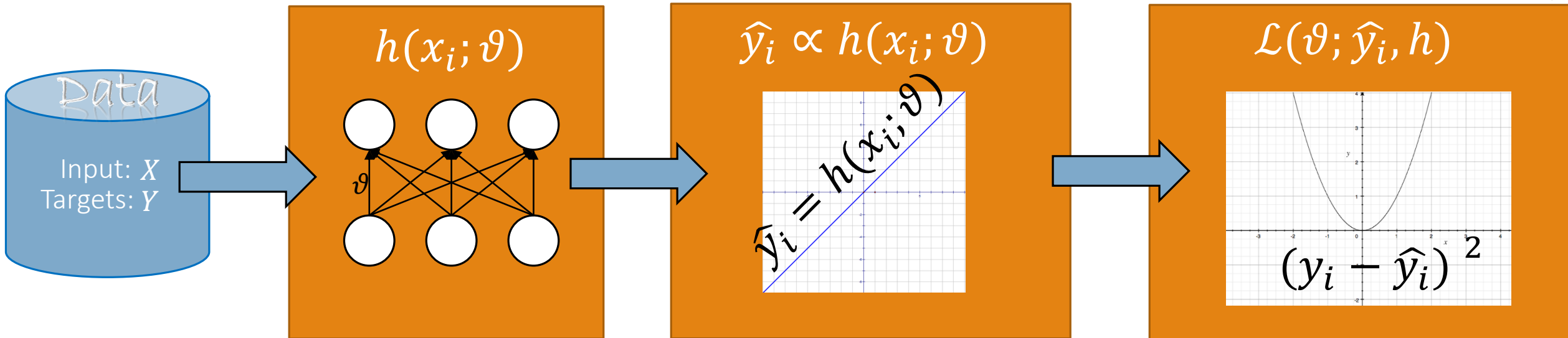- How to implement and check your very own module

# The Machine Learning Paradigm

# Forward computations
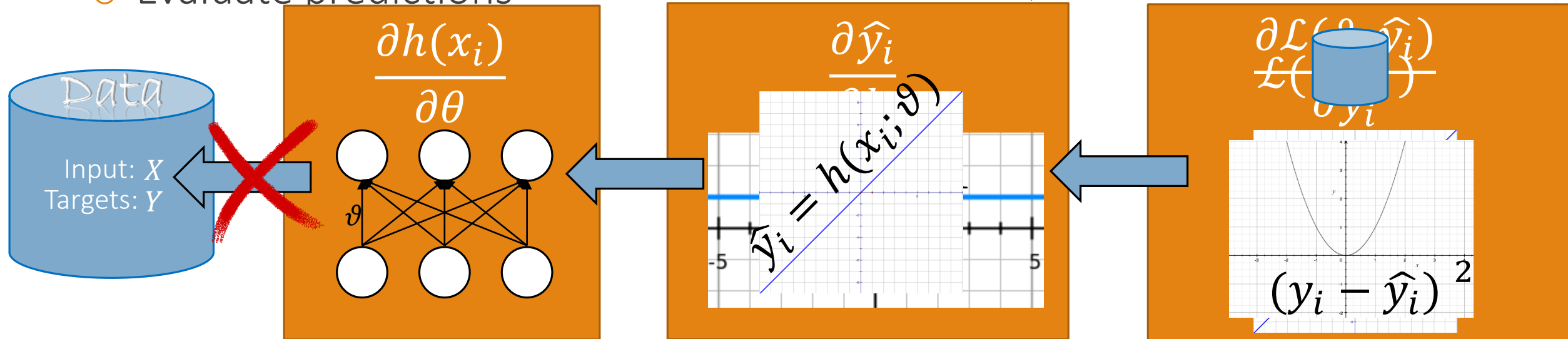
- Collect annotated data

- Define model and initialize randomly

- Predict based on current model
  - In neural network jargon **"forward propagation"**

- Evaluate predictions



Model

$$h(x_i; \vartheta)$$

Score/Prediction/Output

$$\widehat{y}_i \propto h(x_i; \vartheta)$$

Objective/Loss/Cost/Energy

$$\mathcal{L}(\vartheta; \widehat{y}_i, h)$$

$$\widehat{y}_i = h(x_i; \vartheta)$$

$$(y_i - \widehat{y}_i)^2$$

Data

Input: $X$
Targets: $Y$

$\vartheta$

# Backward computations

○ Collect gradient data

○ Define model and initialize randomly

○ Predict based on current model

  ◦ In neural network jargon **"backpropagation"**

○ Evaluate predictions

Model — Score/Prediction/Output — Objective/Loss/Cost/Energy

$$\frac{\partial h(x_i)}{\partial \theta}$$

$$\frac{\partial \widehat{y}_i}{\partial \vartheta}$$

$$\frac{\partial \mathcal{L}(y, \widehat{y}_i)}{\partial \widehat{y}_i}$$

$$\widehat{y}_i = h(x_i; \vartheta)$$

$$(y_i - \widehat{y}_i)^2$$

Input: $X$
Targets: $Y$

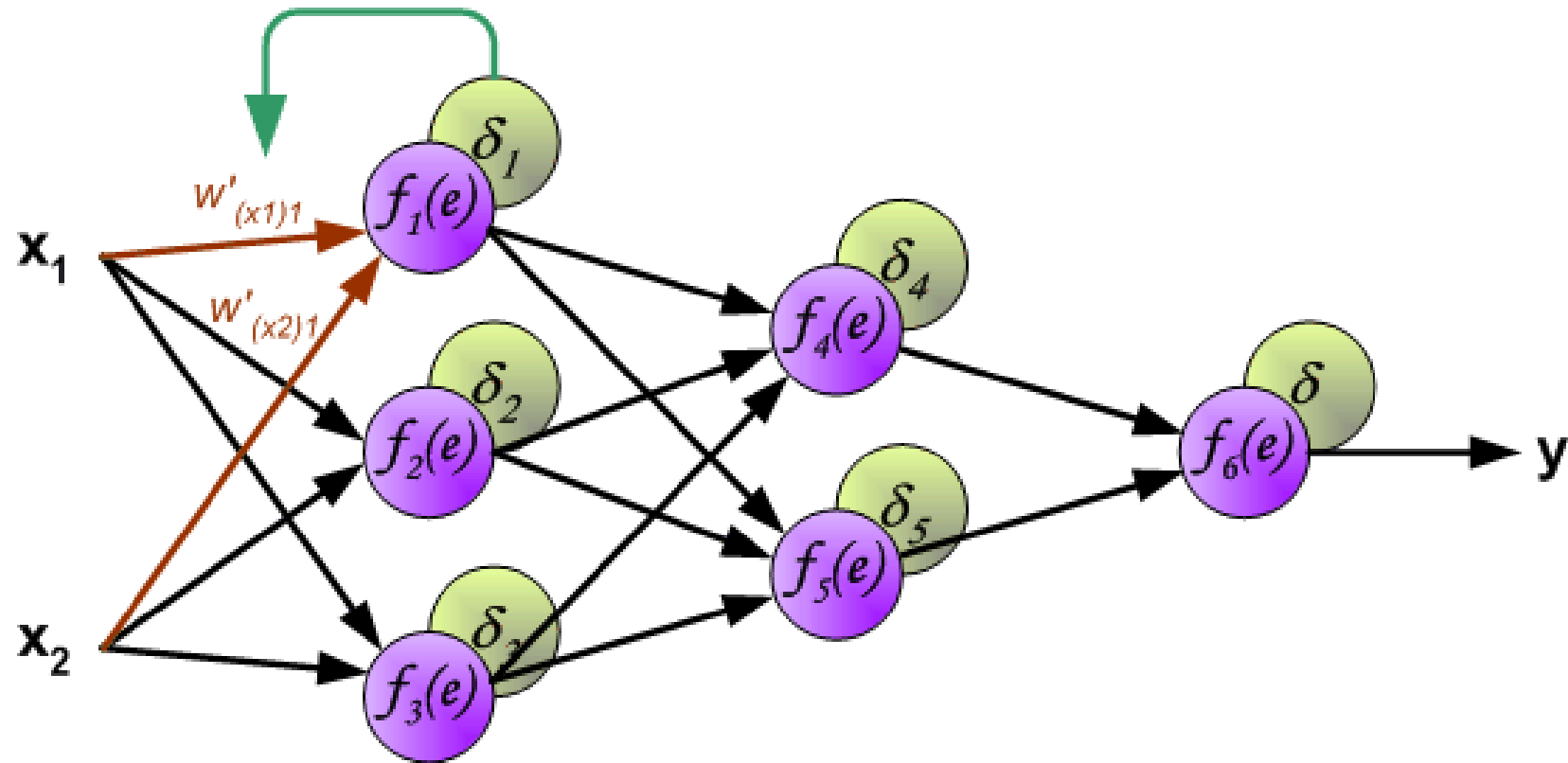Data

$\vartheta$

£( , )

# Optimization through Gradient Descent

o As with many model, we optimize our neural network with Gradient Descent

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

o The most important component in this formulation is the gradient

o The backward computations return the gradients

o How are the backward computations done in a neural network?
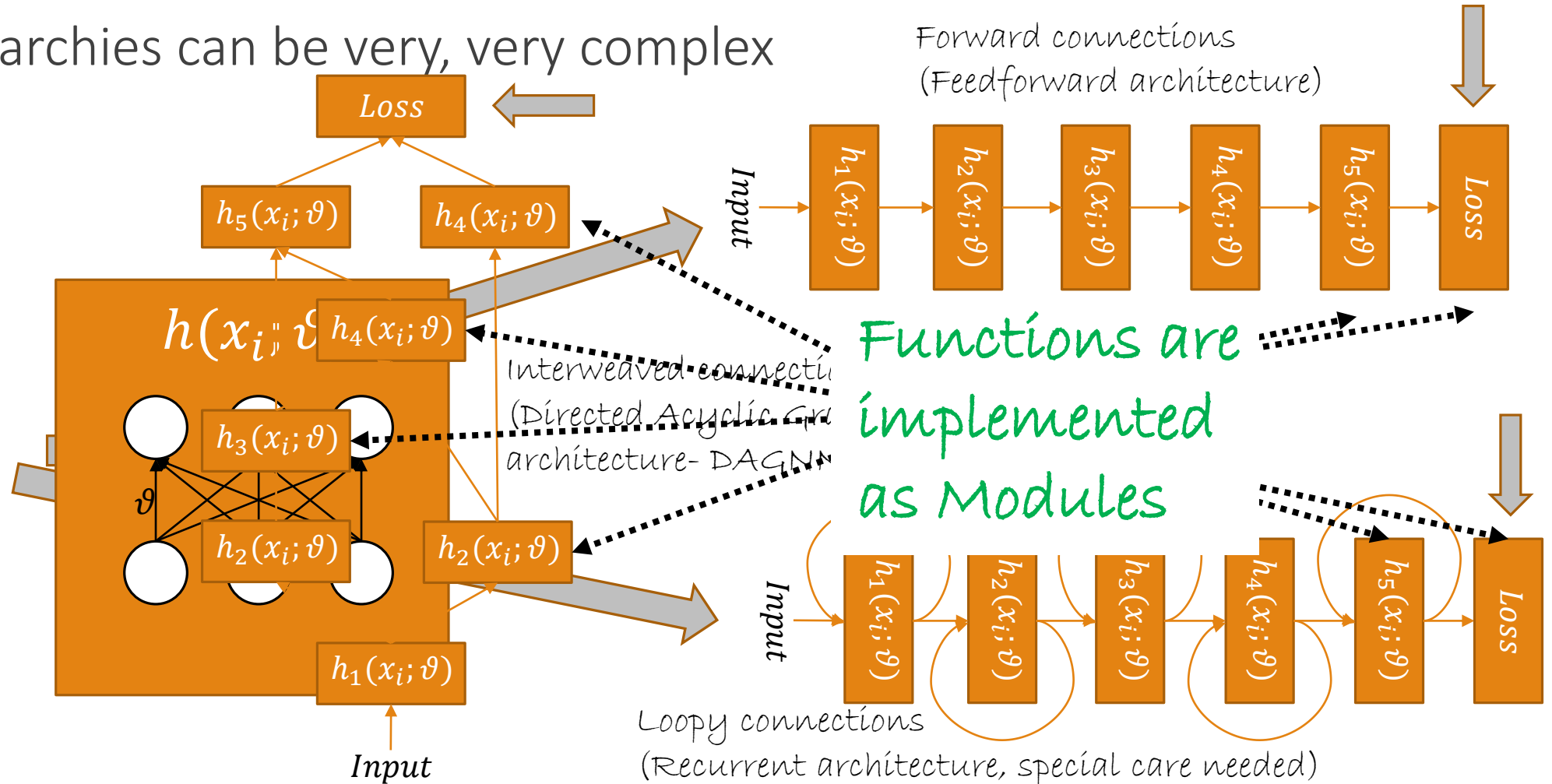
# Backpropagation

# What is a neural network again?

○ A family of parametric, non-linear and hierarchical representation learning functions, which are massively optimized with stochastic gradient descent to encode domain knowledge, i.e. domain invariances, stationarity.

○ $a_L\left(x; \theta_{1,\dots,L}\right) = h_L\left(h_{L-1}(\dots h_1(x, \theta_1), \theta_{L-1}), \theta_L\right)$

    ◦ $x$:input, $\theta_l$: parameters for layer l, $a_l = h_l(x, \theta_l)$: (non-)linear function

○ Given training corpus $\{X, Y\}$ find optimal parameters

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \ell(y, a_L\left(x; \theta_{1,\dots,L}\right))$$
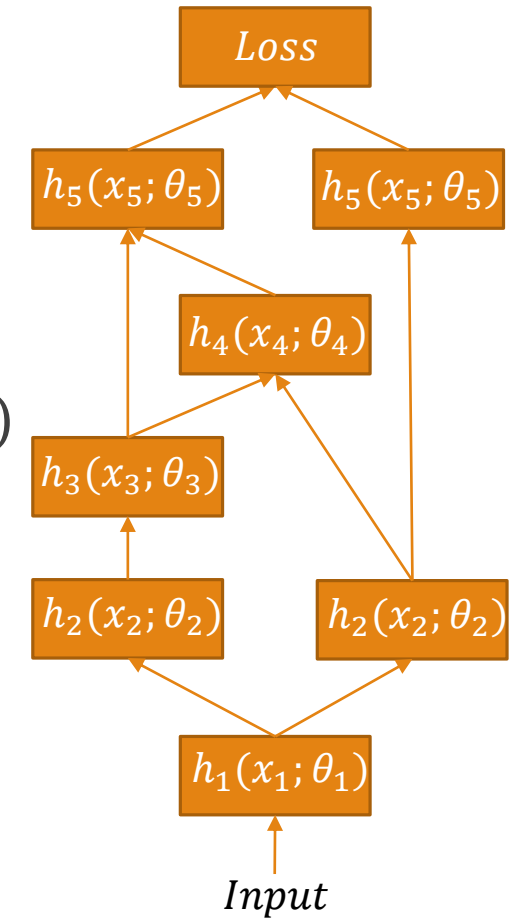
# Neural network models

o A neural network model is a series of hierarchically connected functions

o This hierarchies can be very, very complex



Forward connections
(Feedforward architecture)

$Loss$

$h_5(x_i; \vartheta)$  $h_4(x_i; \vartheta)$

$h(x_i; \vartheta)$  $h_4(x_i; \vartheta)$

$h_3(x_i; \vartheta)$

$\vartheta$

$h_2(x_i; \vartheta)$  $h_2(x_i; \vartheta)$

$h_1(x_i; \vartheta)$

Input

Input  $h_1(x_i; \vartheta)$  $h_2(x_i; \vartheta)$  $h_3(x_i; \vartheta)$  $h_4(x_i; \vartheta)$  $h_5(x_i; \vartheta)$  Loss

Interweaved connections
(Directed Acyclic Graph
architecture- DAGNN)

**Functions are implemented as Modules**

Input  $h_1(x_i; \vartheta)$  $h_2(x_i; \vartheta)$  $h_3(x_i; \vartheta)$  $h_4(x_i; \vartheta)$  $h_5(x_i; \vartheta)$  Loss

Loopy connections
(Recurrent architecture, special care needed)

# What is a module?

o A module is a building block for our network

o Each module is an object/function $a = h(x; \theta)$ that
  ◦ Contains trainable parameters ($\theta$)
  ◦ Receives as an argument an input $x$
  ◦ And returns an output $a$ based on the activation function $h(\dots)$

o The activation function should be (at least)
  **first order differentiable (almost) everywhere**

o For easier/more efficient backpropagation, the
  output of a module should be stored

# Anything goes or do special constraints exist?

- A neural network is a composition of modules (building blocks)

- Any architecture works

- If the architecture is a feedforward cascade, no special care

- If acyclic, there is right order of computing the forward computations

- If there are loops, these form **recurrent** connections (revisited later)

# Forward computations for neural networks

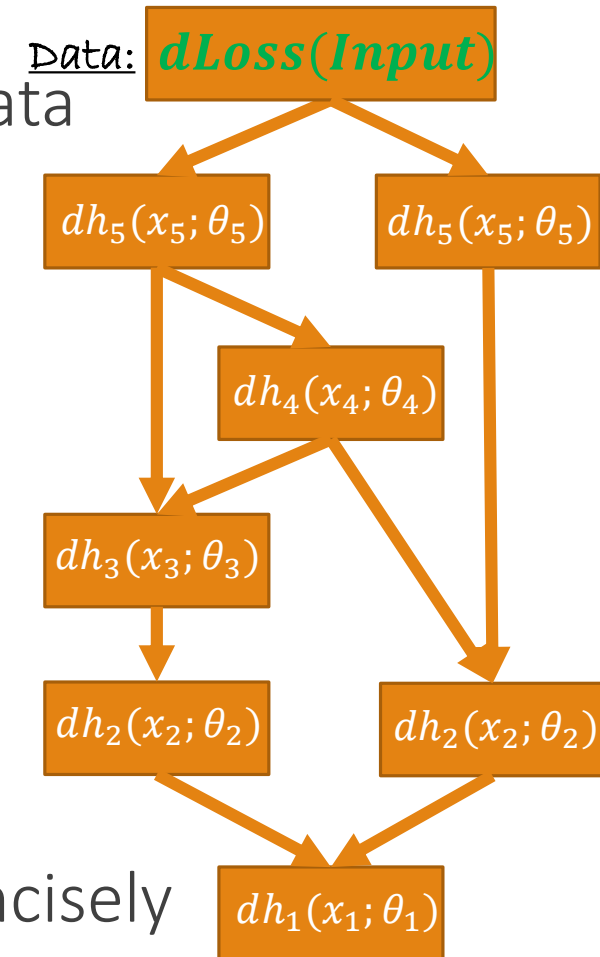o Simply compute the activation of each module in the network

$$a_l = h_l(x_l; \vartheta), \text{ where } a_l = x_{l+1} (\text{or } x_l = a_{l-1})$$

o We need to know the precise function behind each module $h_l(\dots)$

o We start from the data input, e.g. a few images

o Then, we need to compute its module's input
  ◦ It could be that the input is defined from other modules in quite different parts of the network

o So, we compute modules activations **with the right order**
  ◦ Make sure that all the inputs are computed at the right time
  ◦ Then everything goes smoothly

# Backward computations for neural networks

o Simply compute the gradients of each module for our data
  ◦ We need to know the gradient formulation of each module $\partial h_l(x_l; \theta_l)$ w.r.t. their inputs $x_l$ and parameters $\theta_l$

o We need the **forward computations first**
  ◦ Their result is the sum of losses for our input data

o Then take the reverse network (reverse connections) and traverse it backwards

o Instead of using the activation functions, we use their gradients

o The whole process can be described very neatly and concisely with the **backpropagation algorithm**

Data: $\boldsymbol{dLoss(Input)}$

$dh_5(x_5; \theta_5)$   $dh_5(x_5; \theta_5)$

$dh_4(x_4; \theta_4)$

$dh_3(x_3; \theta_3)$

$dh_2(x_2; \theta_2)$   $dh_2(x_2; \theta_2)$

$dh_1(x_1; \theta_1)$

# Again, what is a neural network again?

- $a_L(x; \theta_{1,...,L}) = h_L(h_{L-1}(...h_1(x, \theta_1), \theta_{L-1}), \theta_L)$
  - $x$:input, $\theta_l$: parameters for layer l, $a_l = h_l(x, \theta_l)$: (non-)linear function

- Given training corpus $\{X, Y\}$ find optimal parameters

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \ell(y, a_L(x; \theta_{1,...,L}))$$

- To use any gradient descent based optimization $(\theta^{(t+1)} = \theta^{(t+1)} - \eta_t \frac{\partial \mathcal{L}}{\partial \theta^{(t)}})$ we need the gradients

$$\frac{\partial \mathcal{L}}{\partial \theta_l}, l = 1, ..., L$$

- How to compute the gradients for such a complicated function enclosing other functions, like $a_L(...)$?

# Backpropagation ⟺ Chain rule!!!

- The function $\mathcal{L}(y, a_L)$ depends on $a_L$, which depends on $a_{L-1}$, which depends on $a_{L-2}, \ldots$, which depends on $a_l, \ldots$, which depends on $a_2$

- Chain rule for parameters of layer l

$$\frac{\partial \mathcal{L}(y, a_L)}{\partial \theta_l}$$

- In shorter, we can rewrite this as

*Gradient w.r.t. the module parameters*

$$\frac{\partial \mathcal{L}(y, a_L)}{\partial \theta_l} = \frac{\partial \mathcal{L}}{\partial a_l} \cdot \left(\frac{\partial a_l}{\partial \theta_l}\right)^T$$

$$a_L(x; \theta_{1,\ldots,L}) = h_L(h_{L-1}(\ldots h_1(x, \theta_1), \theta_{L-1}), \quad \theta_L)$$

# Chain rule in practice

- $\dfrac{\partial f}{\partial x} = \dfrac{\partial \, \sin(0.5x^2)}{\partial x} = \dfrac{\partial \, \mathrm{f}(g(x))}{\partial x} =,$ where $0.5x^2$

- $\dfrac{\partial f}{\partial x} = \dfrac{\partial f}{\partial g}\dfrac{\partial g}{\partial x} = x \cdot \cos(0.5x^2)$

# Backpropagation ⟺ Chain rule!!!

○ In $\frac{\partial \mathcal{L}(y, a_L)}{\partial \theta_l} = \frac{\partial \mathcal{L}}{\partial a_l} \cdot \frac{\partial a_l}{\partial \theta_l}$, we need to also easily compute $\frac{\partial \mathcal{L}}{\partial a_l}$. How?

○ Chain rule again

$$\frac{\partial \mathcal{L}}{\partial a_l} = \frac{\partial \mathcal{L}}{\partial a_L} \cdot \frac{\partial a_L}{\partial a_{L-1}} \cdot \frac{\partial a_{L-1}}{\partial a_{L-2}} \cdot \ldots \cdot \frac{\partial a_{l+1}}{\partial a_l}$$

$a_{l+1} = h_{l+1}(x_{l+1}; \theta_{l+1})$

$x_{l+1} = a_l$

$a_l = h_l(x_l; \theta_l)$

○ Remember, the output of a module is the input for the next one: $a_l = x_{l+1}$

○ In shorter, we can rewrite this as

$$\frac{\partial \mathcal{L}}{\partial a_l} = \frac{\partial \mathcal{L}}{\partial a_{l+1}} \cdot \frac{\partial a_{l+1}}{\partial a_l} = (\frac{\partial \mathcal{L}}{\partial a_{l+1}})^T \cdot \frac{\partial a_{l+1}}{\partial x_{l+1}}$$

*Recursive rule (<u>good for us</u>)!!!*

*Gradient w.r.t. the module input*

# Backpropagation for multivariate functions $f(\boldsymbol{x})$

Plenty of functions are c

$$a(x) = \exp(\boldsymbol{x}) = \exp\left(\begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \end{bmatrix}\right) = \begin{bmatrix} \exp(x^{(1)}) \\ \exp(x^{(2)}) \\ \exp(x^{(3)}) \end{bmatrix} = \begin{bmatrix} a(x^{(1)}) \\ a(x^{(2)}) \\ a(x^{(3)}) \end{bmatrix}$$

# Backpropagation for multivariate functions $f(x)$

Plenty of functions are c

$$a(x) = \exp(\boldsymbol{x}) = \exp\left(\begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \end{bmatrix}\right) = \begin{bmatrix} \exp(x^{(1)}) \\ \exp(x^{(2)}) \\ \exp(x^{(3)}) \end{bmatrix} = \begin{bmatrix} a(x^{(1)}) \\ a(x^{(2)}) \\ a(x^{(3)}) \end{bmatrix}$$

o Some functions, however, depend on multiple input variables
  ◦ Softmax!
  ◦ Each output dimension depends on multiple input dimensions

$$a^{(j)} = \frac{e^{x^{(j)}}}{e^{x^{(1)}} + e^{x^{(2)}} + e^{x^{(3)}}}$$

# Backpropagation for multivariate functions $f(x)$

Plenty of functions are c

$$(or\ \partial a l\ \partial \theta l\ )a\text{ we compute the Jacobian matrix}$$
$$l\ l\ a\ l\ \partial \theta l$$

$$a(x) = \exp(x) = \exp\left(\begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \end{bmatrix}\right) = \begin{bmatrix} \exp(x^{(1)}) \\ \exp(x^{(2)}) \\ \exp(x^{(3)}) \end{bmatrix} = \begin{bmatrix} a(x^{(1)}) \\ a(x^{(2)}) \\ a(x^{(3)}) \end{bmatrix}$$

o Some functions, however, depend on multiple input variables

$$a^{(j)} = \frac{e^{x^{(j)}}}{e^{x^{(1)}} + e^{x^{(2)}} + e^{x^{(3)}}}$$

 ◦ Softmax!

 ◦ Each output dimension depends on multiple input dimensions

o For these cases for the
$$\partial a l\ \partial x l\quad a l\ \partial x l\quad l l l a l\ \partial x l\quad l l l a l\ \partial x l\quad \left(or\ \frac{\partial a_l}{\partial \theta_l}\right)\ \text{we compute}$$
the Jacobian matrix

# The Jacobian

○ When $a(x)$ is $2-\mathbf{d}$ and depends on 3 variables, $x^{(1)}, x^{(2)}, x^{(3)}$

$$J\big(a(x)\big) = \begin{bmatrix} \dfrac{\partial a^{(1)}}{\partial x^{(1)}} & \dfrac{\partial a^{(1)}}{\partial x^{(2)}} & \dfrac{\partial a^{(1)}}{\partial x^{(3)}} \\[2em] \dfrac{\partial a^{(2)}}{\partial x^{(1)}} & \dfrac{\partial a^{(2)}}{\partial x^{(2)}} & \dfrac{\partial a^{(2)}}{\partial x^{(3)}} \end{bmatrix}$$

# Backpropagation for multivariate functions $f(\boldsymbol{x})$

o Plenty of functions are computed element-wise

   ◦ $\sigma(x), \tanh(x), \exp(x)$

   ◦ Each output dimension depends **only** on the respective input dimension

$$a(x) = \exp(\boldsymbol{x}) = \exp\left(\begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \end{bmatrix}\right) = \begin{bmatrix} \exp(x^{(1)}) \\ \exp(x^{(2)}) \\ \exp(x^{(3)}) \end{bmatrix} = \begin{bmatrix} a(x^{(1)}) \\ a(x^{(2)}) \\ a(x^{(3)}) \end{bmatrix}$$

o Some functions, however, depend on multiple input variables

   ◦ Softmax!

   ◦ Each output dimension depends on multiple input dimensions

$$a^{(j)} = \frac{e^{x^{(j)}}}{e^{x^{(1)}} + e^{x^{(2)}} + e^{x^{(3)}}}$$

o For these cases for the $\frac{\partial a_l}{\partial x_l}$ (or $\frac{\partial a_l}{\partial \theta_l}$)  we compute the Jacobian matrix

o Then, $\frac{\partial \mathcal{L}}{\partial a_l} = \left(\frac{\partial \mathcal{L}}{\partial a_{l+1}}\right)^T \cdot \frac{\partial a_{l+1}}{\partial x_{l+1}}$

# Dimension analysis

o To make sure everything is done correctly ➔ "Dimension analysis"

o The dimensions of the gradient w.r.t. $\theta_l$ must be equal to the dimensions of the respective weight $\theta_l$

$$\dim\left(\frac{\partial \mathcal{L}}{\partial a_l}\right) = \dim(a_l) \text{ and } \dim\left(\frac{\partial \mathcal{L}}{\partial \theta_l}\right) = \dim(\theta_l)$$

o E.g. for $\frac{\partial \mathcal{L}}{\partial a_l} = (\frac{\partial \mathcal{L}}{\partial a_{l+1}})^T \cdot \frac{\partial a_{l+1}}{\partial x_{l+1}}$, if $\dim(a_l) = d_l$, then it should be

$$[d_l \times 1] = [1 \times d_{l+1}] \cdot [d_{l+1} \times d_l]$$

o E.g. for $\frac{\partial \mathcal{L}}{\partial \theta_l} = \frac{\partial \mathcal{L}}{\partial a_l} \cdot (\frac{\partial a_l}{\partial \theta_l})^T$, if $\dim(\theta_l) = d_l \times d_{l-1}$, then it should be

$$[d_l \times d_{l-1}] = [d_l \times 1] \cdot [1 \times d_{l-1}]$$

# Backpropagation again

- o **Step 1.** Compute forward propagations for all layers, starting from the first layer until the last loss layer

$$a_l = h_l(x_l) \text{ and } x_{l+1} = a_l$$

- o **Step 2.** Once done with forward propagation, follow the reverse path. Start from the last layer and for each new layer compute the gradients

$$\frac{\partial \mathcal{L}}{\partial a_l} = \left(\frac{\partial \mathcal{L}}{\partial a_{l+1}}\right)^T \cdot \frac{\partial a_{l+1}}{\partial x_{l+1}} \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial \theta_l} = \frac{\partial \mathcal{L}}{\partial a_l} \cdot \left(\frac{\partial a_l}{\partial \theta_l}\right)^T$$

- ◦ Cache computations when possible to avoid redundant operations

- o **Step 3.** Use the gradients $\frac{\partial \mathcal{L}}{\partial \theta_l}$ with Stochastic Gradient Descend to train your network

Vector with dimensions $[d_{l+1} \times 1]$

Vector with dimensions $[d_l \times 1]$

Matrix with dimensions $[d_l \times d_{l-1}]$

Vector with dimensions $[1 \times d_{l-1}]$
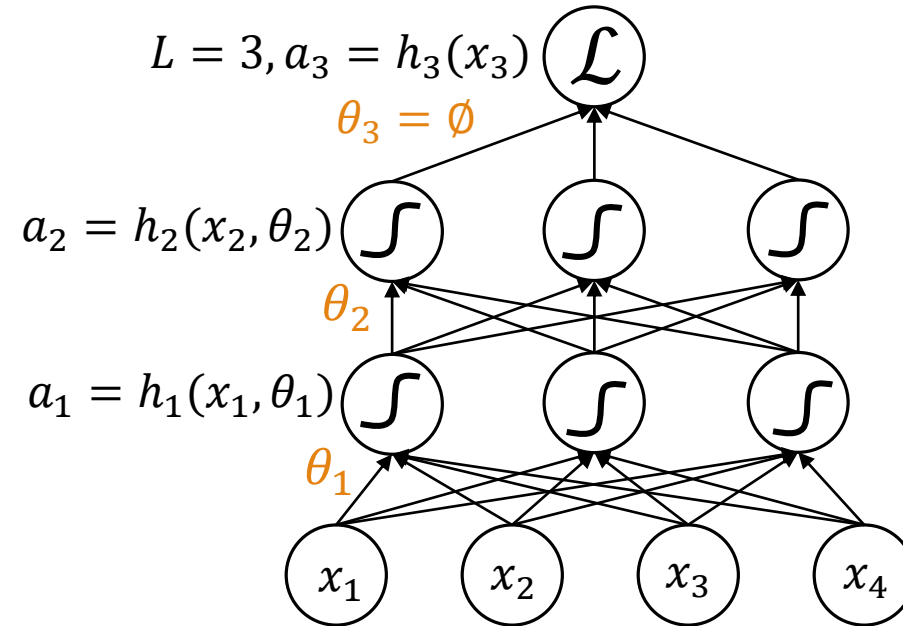
Vector with dimensions $[d_l \times 1]$

Jacobian matrix with dimensions $[d_{l+1} \times d_l]$

# Practical example and dimensionality analysis

o Layer $l-1$ has 15 neurons ($d_{l-1} = 15$), $l$ has 10 neurons ($d_l = 10$) and $l+1$ has 5 neurons ($d_{l+1} = 5$)

o My activation functions are $a_l = w_l x_l$ and $a_{l+1} = w_{l+1} x_{l+1}$

o The dimensionalities are *(remember $x_l = a_{l-1}$)*
  ◦ $a_{l-1} \to [15 \times 1], \ a_l \quad \to [10 \times 1], \ a_{l+1} \to [5 \times 1]$
  ◦ $x_l \quad \to [15 \times 1], \ x_{l+1} \quad \to [10 \times 1]$
  ◦ $\theta_l \quad \to [10 \times 15], w_{l+1} \to [5 \times 10]$

o The gradients are
  ◦ $\dfrac{\partial \mathcal{L}}{\partial a_l} \quad \to [1 \times 5] \cdot [5 \times 10] = [1 \times 10]$
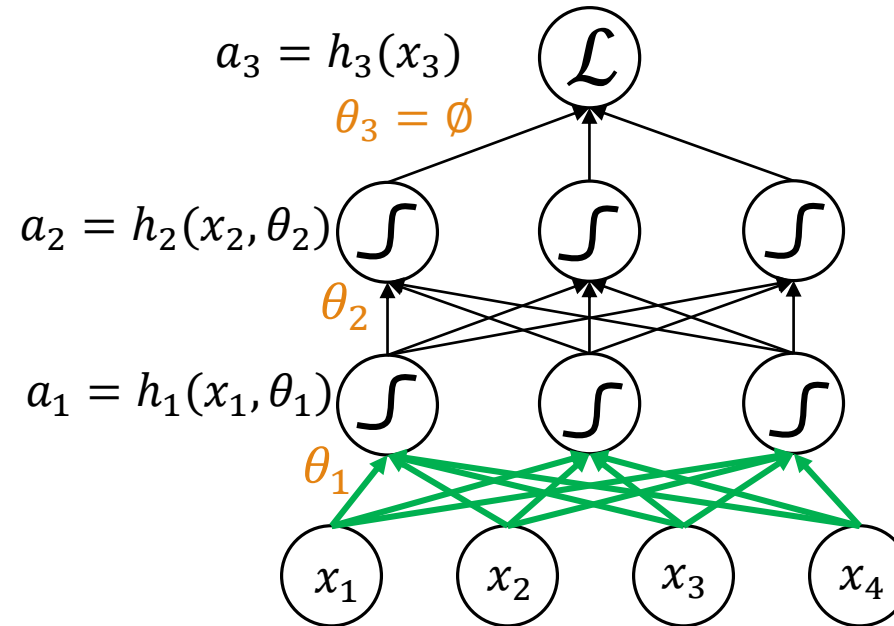  ◦ $\dfrac{\partial \mathcal{L}}{\partial \theta_l} \quad \to [10 \times 1] \cdot [1 \times 15] = [10 \times 15]$

# Backpropagation visualization



$L = 3, a_3 = h_3(x_3)$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$

$\theta_1$

$x_1$ $x_2$ $x_3$ $x_4$

# Backpropagation visualization at epoch ($t$)

Forward propagations

Compute and store $a_1 = h_1(x_1)$

$a_3 = h_3(x_3)$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$

$\theta_1$

$x_1$  $x_2$  $x_3$  $x_4$

Example

$a_1 = \sigma(\theta_1 x_1)$

Store!!!

Forward propagations

Compute and store $a_2 = h_2(x_2)$

$L = 3, a_3 = h_3(x_3)$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$

$\theta_1$

$x_1 \quad x_2 \quad x_3 \quad x_4$

Example

$a_1 = \sigma(\theta_1 x_1)$

$a_2 = \sigma(\theta_2 x_2)$

Store!!!

Forward propagations

Compute and store $a_3 = h_3(x_3)$

$L = 3, a_3 = h_3(x_3)$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$

$\theta_1$

$x_1$  $x_2$  $x_3$  $x_4$

Example

$a_1 = \sigma(\theta_1 x_1)$

$a_2 = \sigma(\theta_2 x_2)$

$a_3 = \|y - x_3\|^2$

Store!!!

# Backpropagation visualization at epoch $(t)$

Backpropagation

$$\frac{\partial \mathcal{L}}{\partial a_3} = \ldots \leftarrow \text{Direct computation}$$

~~$\frac{\partial \mathcal{L}}{\partial \theta_3}$~~

$a_3 = h_3(x_3)$ $\mathcal{L}$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$ $\int$ $\int$ $\int$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$ $\int$ $\int$ $\int$

$\theta_1$

$x_1$ $x_2$ $x_3$ $x_4$

Example
$$a_3 = \mathcal{L}(y, x_3) = h_3(x_3) = 0.5 \|y - x_3\|^2$$
$$\frac{\partial \mathcal{L}}{\partial x_3} = -(y - x_3)$$

# Backpropagation visualization at epoch $(t)$

# Backpropagation visualization at epoch ($t$)

Backpropagation

$$\frac{\partial \mathcal{L}}{\partial a_1} = \frac{\partial \mathcal{L}}{\partial a_2} \cdot \frac{\partial a_2}{\partial a_1}$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{\partial \mathcal{L}}{\partial a_1} \cdot \frac{\partial a_1}{\partial \theta_1}$$

$a_3 = h_3(x_3)$ $\mathcal{L}$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$ $\int$ $\int$ $\int$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$ $\int$ $\int$ $\int$

$\theta_1$

$x_1$ $x_2$ $x_3$ $x_4$

Computed from the exact previous
backpropagation step (Remember, recursive rule)

Example

$$\mathcal{L}(y, a_3) = 0.5 \|y - a_3\|^2$$
$$a_2 = \sigma(\theta_2 x_2)$$
$$x_2 = a_1$$
$$a_1 = \sigma(\theta_1 x_1)$$
$$\frac{\partial a_2}{\partial a_1} = \frac{\partial a_2}{\partial x_2} = \theta_2 a_2 (1 - a_2)$$
$$\frac{\partial a_1}{\partial \theta_1} = x_1 a_1 (1 - a_1)$$

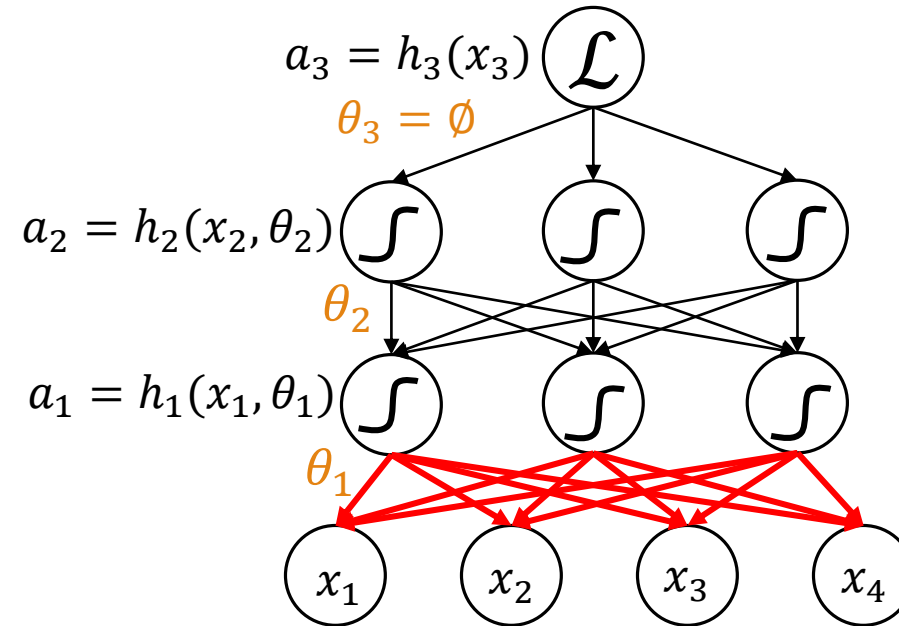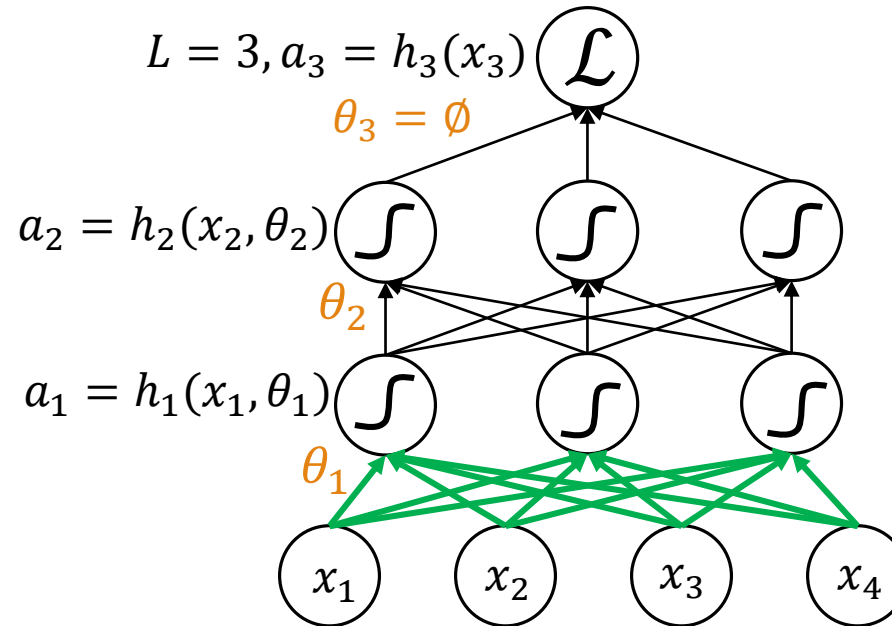$$\frac{\partial \mathcal{L}}{\partial a_1} = \frac{\partial \mathcal{L}}{\partial a_2} \theta_2 a_2 (1 - a_2)$$
$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{\partial \mathcal{L}}{\partial a_1} x_1 a_1 (1 - a_1)$$

Forward propagations

Compute and store $a_1 = h_1(x_1)$

$L = 3, a_3 = h_3(x_3)$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$

$\theta_1$

$x_1$ $x_2$ $x_3$ $x_4$

Example

$a_1 = \sigma(\theta_1 x_1)$

Store!!!

Forward propagations

Compute and store $a_2 = h_2(x_2)$

$L = 3, a_3 = h_3(x_3)$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$

$\theta_1$

$x_1$  $x_2$  $x_3$  $x_4$
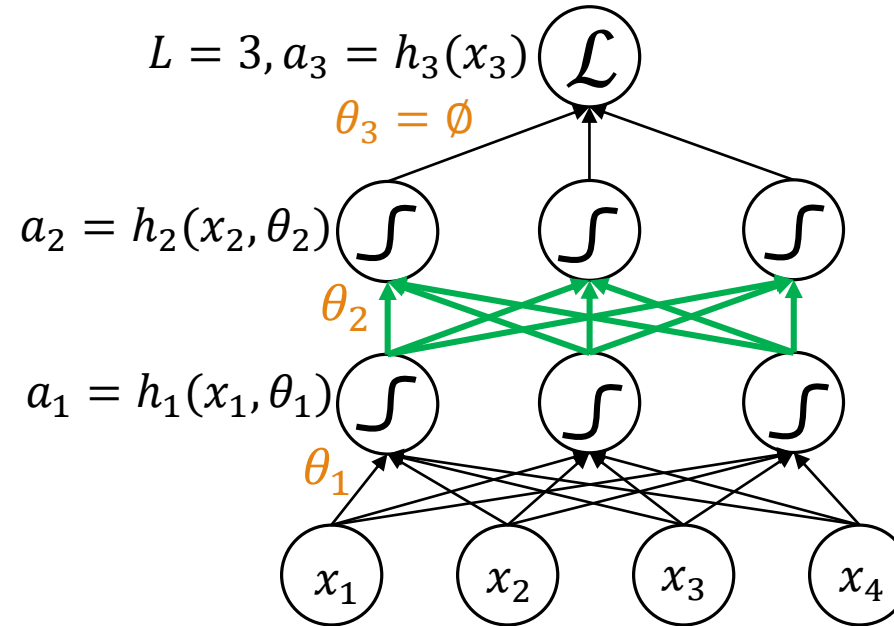
Example

$a_1 = \sigma(\theta_1 x_1)$
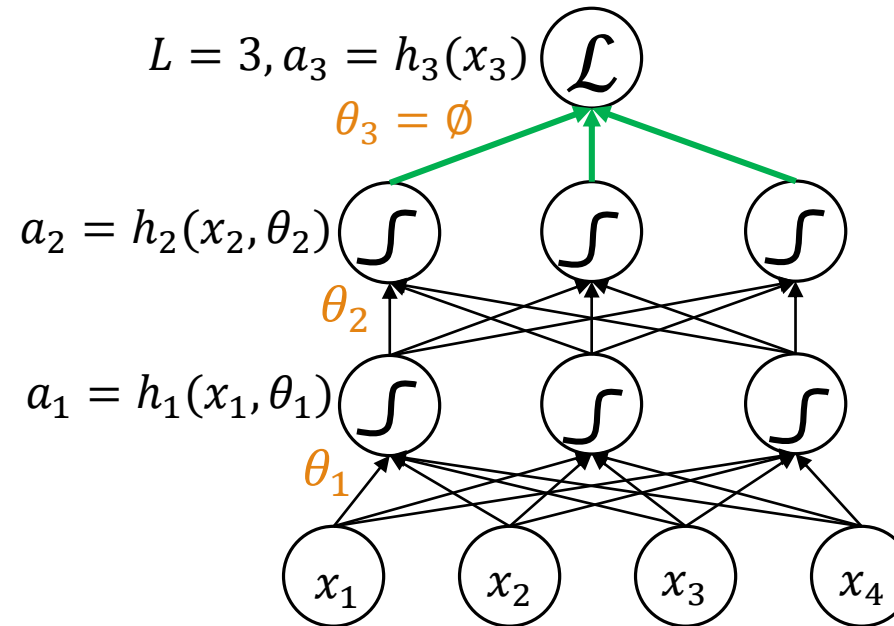$a_2 = \sigma(\theta_2 x_2)$

Store!!!

# Backpropagation visualization at epoch $(t+1)$

_Forward propagations_

_Compute and store_ $a_3 = h_3(x_3)$

$L = 3, a_3 = h_3(x_3)$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$

$\theta_1$

$\mathcal{L}$

$x_1 \quad x_2 \quad x_3 \quad x_4$

_Example_

$a_1 = \sigma(\theta_1 x_1)$

$a_2 = \sigma(\theta_2 x_2)$

$a_3 = \|y - x_3\|^2$

_Store_!!!

Backpropagation

$$\frac{\partial \mathcal{L}}{\partial a_3} = \ldots \leftarrow \text{Direct computation}$$

~~$$\frac{\partial \mathcal{L}}{\partial \theta_3}$$~~

$a_3 = h_3(x_3)$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$

$\theta_1$

$x_1 \quad x_2 \quad x_3 \quad x_4$

Example

$$a_3 = \mathcal{L}(y, x_3) = h_3(x_3) = 0.5\, \|y - x_3\|^2$$
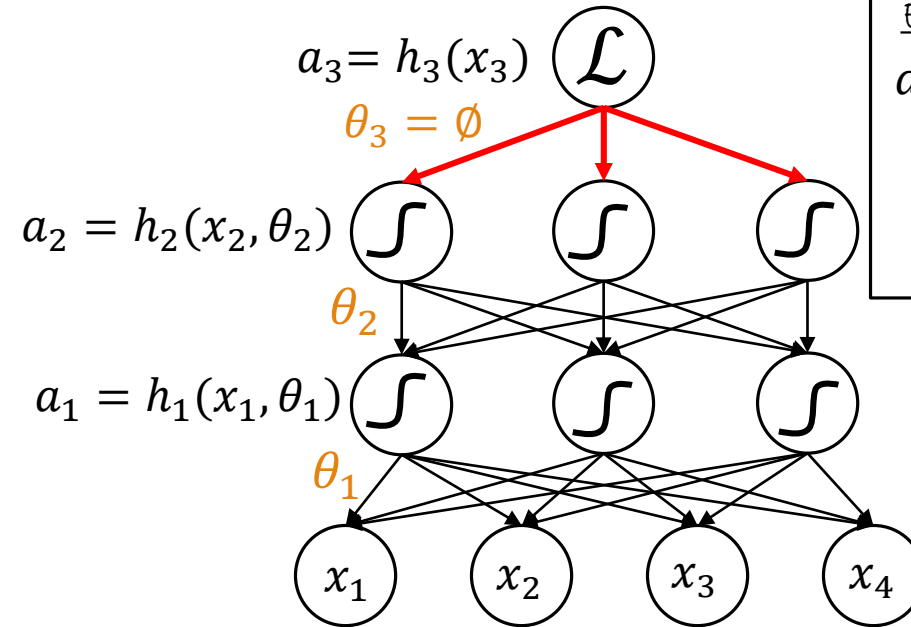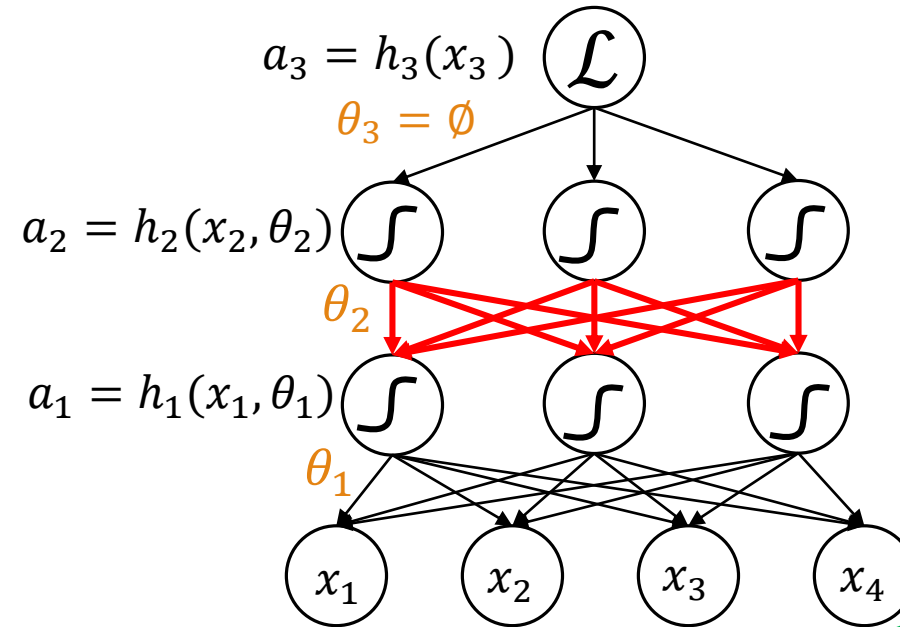
$$\frac{\partial \mathcal{L}}{\partial x_3} = -(y - x_3)$$

# Backpropagation visualization at epoch $(t+1)$

Backpropagation

$$\frac{\partial \mathcal{L}}{\partial a_2} = \frac{\partial \mathcal{L}}{\partial a_3} \cdot \frac{\partial a_3}{\partial a_2}$$

$$\frac{\partial \mathcal{L}}{\partial \theta_2} = \frac{\partial \mathcal{L}}{\partial a_2} \cdot \frac{\partial a_2}{\partial \theta_2}$$

$a_3 = h_3(x_3)$  $\mathcal{L}$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$  ∫  ∫  ∫

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$  ∫  ∫  ∫

$\theta_1$

$x_1$  $x_2$  $x_3$  $x_4$

Stored during forward computations

Example

$$\mathcal{L}(y, x_3) = 0.5 \, \|y - x_3\|^2$$

$$x_3 = a_2$$
$$a_2 = \sigma(\theta_2 x_2)$$
$$\frac{\partial \mathcal{L}}{\partial a_2} = \frac{\partial \mathcal{L}}{\partial x_3} = -(y - x_3)$$
$$\partial \sigma(x) = \sigma(x)(1 - \sigma(x))$$
$$\frac{\partial a_2}{\partial \theta_2} = x_2 \sigma(\theta_2 x_2)(1 - \sigma(\theta_2 x_2))$$
$$= x_2 a_2 (1 - a_2)$$

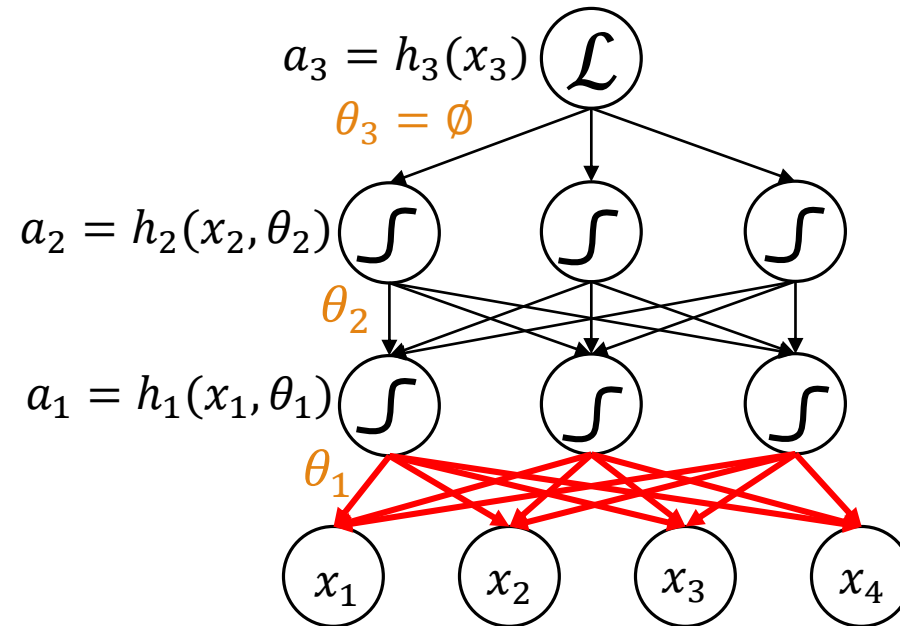$$\frac{\partial \mathcal{L}}{\partial a_2} = -(y - x_3)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_2} = \frac{\partial \mathcal{L}}{\partial a_2} x_2 a_2 (1 - a_2)$$

# Backpropagation visualization at epoch $(t + 1)$

Backpropagation

$$\frac{\partial \mathcal{L}}{\partial a_1} = \frac{\partial \mathcal{L}}{\partial a_2} \cdot \frac{\partial a_2}{\partial a_1}$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{\partial \mathcal{L}}{\partial a_1} \cdot \frac{\partial a_1}{\partial \theta_1}$$

$a_3 = h_3(x_3)$ $\mathcal{L}$

$\theta_3 = \emptyset$

$a_2 = h_2(x_2, \theta_2)$

$\theta_2$

$a_1 = h_1(x_1, \theta_1)$

$\theta_1$

$x_1$  $x_2$  $x_3$  $x_4$

Computed from the exact previous
backpropagation step (Remember, recursive rule)

Example

$$\mathcal{L}(y, a_3) = 0.5 \|y - a_3\|^2$$
$$a_2 = \sigma(\theta_2 x_2)$$
$$x_2 = a_1$$
$$a_1 = \sigma(\theta_1 x_1)$$
$$\frac{\partial a_2}{\partial a_1} = \frac{\partial a_2}{\partial x_2} = \theta_2 a_2 (1 - a_2)$$
$$\frac{\partial a_1}{\partial \theta_1} = x_1 a_1 (1 - a_1)$$

$$\frac{\partial \mathcal{L}}{\partial a_1} = \frac{\partial \mathcal{L}}{\partial a_2} \theta_2 a_2 (1 - a_2)$$
$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{\partial \mathcal{L}}{\partial a_1} x_1 a_1 (1 - a_1)$$

# Some practical tricks of the trade

o For classification use cross-entropy loss

o Use Stochastic Gradient Descent on mini-batches

o Shuffle training examples **at each** new epoch

o Normalize input variables to $(\mu, \sigma^2) = (0,1)$
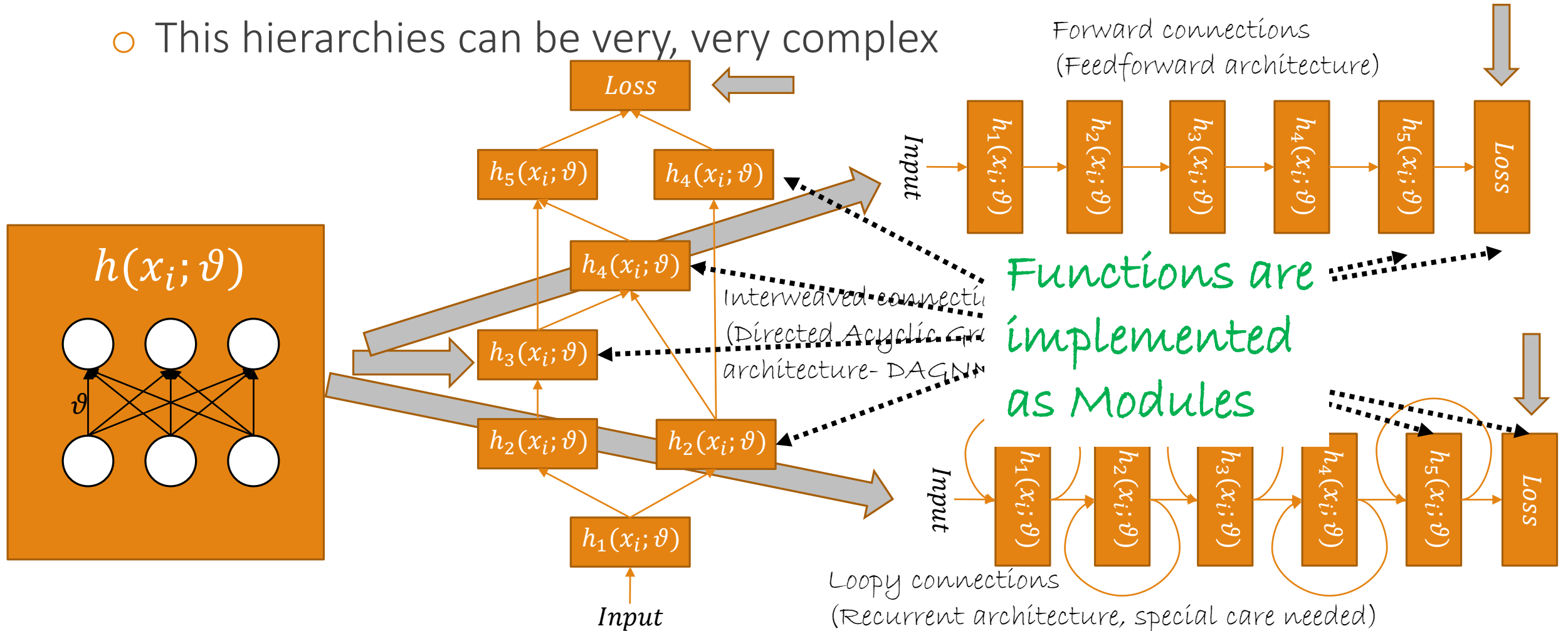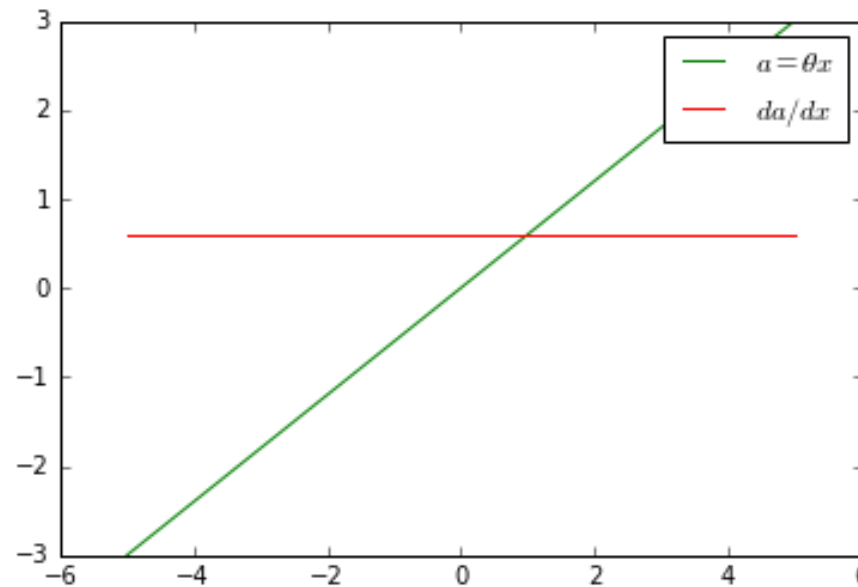
# Everything is a *module*

# Neural network models

o A neural network model is a series of hierarchically connected functions

o This hierarchies can be very, very complex



Forward connections
(Feedforward architecture)

Functions are
implemented
as Modules

Interweaved connections
(Directed Acyclic Graph
architecture- DAGNN)

Loopy connections
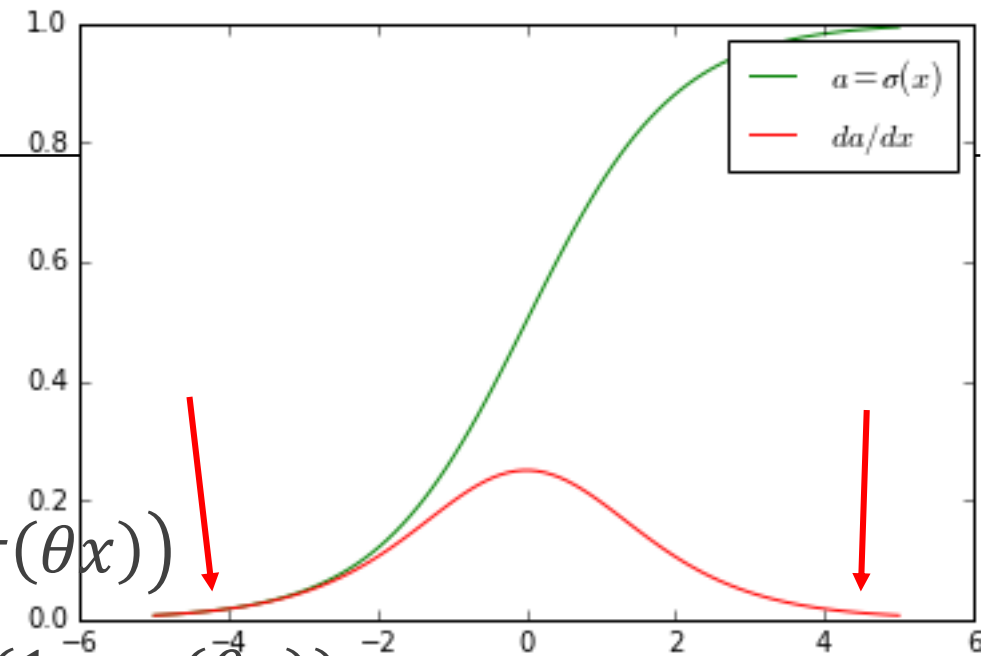(Recurrent architecture, special care needed)

# Linear module

○ Activation function $a = \theta x$

○ Gradient with respect to the input $\dfrac{\partial a}{\partial x} = \theta$

○ Gradient with respect to the parameters $\dfrac{\partial a}{\partial \theta} = x$

# Sigmoid module



o Activation function $a = \sigma(x) = \frac{1}{1+e^{-x}}$

o Gradient wrt the input $\frac{\partial a}{\partial x} = \sigma(x)(1 - \sigma(x))$

o Gradient wrt the input $\frac{\partial \sigma(\theta x)}{\partial x} = \theta \cdot \sigma(\theta x)(1 - \sigma(\theta x))$

o Gradient wrt the parameters $\frac{\partial \sigma(\theta x)}{\partial \theta} = x \cdot \sigma(\theta x)(1 - \sigma(\theta x))$

o Output can be interpreted as probability

o Always bounds the outputs between 0 and 1, so the network cannot overshoot

o Gradients can be small in deep networks because we always multiply with <1

o The gradients at the tails are flat to 0, hence no serious updates
  ◦ Overconfident, but not necessarily "correct", neurons get stuck

# Simplifying backpropagation equations

o We often want to apply a non-linearity $\sigma(\dots)$ on top of an activation $\theta x$

$$a = \sigma(\theta x)$$

o This way we end up with quite complicated backpropagation equations

o Since **everything is a module**, we can decompose this **to 2 modules**

$$\boxed{a_1 = \theta x} \rightarrow \boxed{a_2 = \sigma(a_1)}$$

o We now have to perform two backpropagation steps instead of one

o **But now our gradients are simpler**
  ◦ The complications happen when non-linear functions are parametric
  ◦ We avoid taking the extra gradients w.r.t. parameters inside a non-linearity
  ◦ This is usually how networks are implemented in Torch

# Tanh module

o Activation function $a = tanh(x) = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$

o Gradient with respect to the input $\dfrac{\partial a}{\partial x} = 1 - tanh^2(x)$

o Similar to sigmoid, but with different output range
   ◦ $[-1, +1]$ instead of $[0, +1]$
   ◦ Stronger gradients, because data is centered around 0 (not 0.5)
   ◦ Less bias to hidden layer neurons as now outputs can be both positive and negative (more likely to have zero mean in the end)

# Softmax module

o Activation function $a^{(k)} = softmax(x^{(k)}) = \frac{e^{x^{(k)}}}{\sum_j e^{x^{(j)}}}$

　◦ This activation function is mostly used for making decisions in a form of a probability

　◦ $\sum_{k=1}^{K} a^{(k)} = 1$ for $K$ classes

o Exploiting the fact that $e^{a+b} = e^a e^b$, we usually compute

$$a^{(k)} = \frac{e^{x^{(k)}-\mu}}{\sum_j e^{x^{(j)}-\mu}}, \mu = \max_k x^{(k)} \text{ as } \frac{e^{x^{(k)}-\mu}}{\sum_j e^{x^{(j)}-\mu}} = \frac{e^\mu e^{x^{(k)}}}{e^\mu \sum_j e^{x^{(j)}}} = \frac{e^{x^{(k)}}}{\sum_j e^{x^{(j)}}}$$

　◦ This provides better stability because avoids exponentianting large numbers

# Euclidean loss module

○ Activation function $a(x) = 0.5 \|y - x\|^2$

   ◦ Mostly used to measure the loss in regression tasks

○ Gradient with respect to the input $\dfrac{\partial a}{\partial x} = x - y$

# Cross-entropy loss (log-loss or log-likelihood) module

- Activation function $a(x) = -\sum_{k=1}^{K} y^{(k)} \log x^{(k)}, \quad y^{(k)} = \{0, 1\}$

- Gradient with respect to the input $\frac{\partial a}{\partial x^{(k)}} = -\frac{1}{x^{(k)}}$

- The cross-entropy loss is the most popular classification losses for classifiers that output probabilities (not SVM)

- The cross-entropy loss couples well with certain input activations, such as the softmax module or the sigmoid module
  - Often the gradients of the cross-entropy loss are computed in conjunction with the activation function from the previous layer

- Generalization of logistic regression for more than 2 outputs

# More specific modules for later

o There are many more modules that are quite often used in Deep Learning

o Convolutional filter modules

o Rectified Linear Unit (ReLU) module

o Parametric ReLU module

o Regularization modules
  ◦ Dropout

o Normalization modules
  ◦ $\ell_2$-normalization

o Loss modules
  ◦ Hinge loss

o and others, which we are going to discuss later in the course

# Make your own module

# New modules

o Everything can be a module, given some ground rules

o How to make our own module?
  ◦ Write a function that follows the ground rules

o Needs to be (at least) first-order differentiable (almost) everywhere

o Hence, we need to be able to compute the

$$\frac{\partial a(x;\theta)}{\partial x} \text{ and } \frac{\partial a(x;\theta)}{\partial \theta}$$

# A module of modules

o   As everything can be a module, a module of modules could also be a module

   ◦ In fact, [Lin2014] proposed a Network-in-Network architecture

o   We can therefore make new building blocks as we please, if we expect them to be used frequently

o   Of course, the same rules for the eligibility of modules still apply

# Radial Basis Function (RBF) Network module

o Assume we want to build an RBF module

$$a = \sum_j u_j \exp(-\beta_j (x - w_j)^2)$$

o To avoid computing the full derivations, we can decompose this module into a cascade of modules

$$a_1 = (x - w)^2 \rightarrow a_2 = \exp(-\beta a_1) \rightarrow a_3 = u a_2 \rightarrow a_4 = plus(\ldots, a_3^{(j)}, \ldots)$$

o An RBF module is good for regression problems, in which cases it is followed by a Euclidean loss module

o The Gaussian centers $w_j$ can be initialized externally, e.g. with k-means

# An RBF visually

$$a_5 = \|y - a_4\|^2$$

RBF module

$$a_4 = plus(a_3)$$

$$\alpha_3^{(1)} = u_1 \alpha_2^{(1)} \qquad \alpha_3^{(2)} = u_2 \alpha_2^{(2)} \qquad \alpha_3^{(3)} = u_3 \alpha_2^{(3)}$$

$$\alpha_2^{(1)} = \exp(-\beta_1 \alpha_1^{(1)}) \qquad \alpha_2^{(2)} = \exp(-\beta_1 \alpha_1^{(2)}) \qquad \alpha_2^{(3)} = \exp(-\beta_1 \alpha_1^{(3)})$$

$$\alpha_1^{(1)} = (x - w_1)^2 \qquad \alpha_1^{(2)} = (x - w_1)^2 \qquad \alpha_1^{(3)} = (x - w_1)^2$$

$$a_1 = (x - w)^2 \rightarrow a_2 = \exp(-\beta a_1) \rightarrow a_3 = u a_2 \rightarrow a_4 = plus(\ldots, a_3^{(j)}, \ldots)$$

# Gradient check

Original gradient definition: $\frac{df(x)}{dx} = \lim_{h \to 0} \frac{f(x+h)}{\Delta h}$

- The most dangerous part when implementing new modules is to get your gradients right
  - The math might be wrong, the code might be wrong, …

- Check your module with gradient checks.
  - Compare your explicit gradient with computational gradient $g\left(\theta^{(i)}\right) \approx \frac{a(\theta + \varepsilon) - a(\theta - \varepsilon)}{2\varepsilon}$

$$\Delta\left(\theta^{(i)}\right) = \left\| \frac{\partial a\left(x; \theta^{(i)}\right)}{\partial \theta^{(i)}} - g\left(\theta^{(i)}\right) \right\|^2$$

  - If result is smaller than $\delta \in (10^{-4}, 10^{-7})$, then your gradients are good

- Perturb one parameter $\theta^{(i)}$ at a time, $\theta^{(i)} + \varepsilon$, then check its $\Delta\left(\theta^{(i)}\right)$
  - **Do not** perturb the whole parameter vector $\theta + \varepsilon$, it will give **wrong results**

- Good practice: check your network gradients too

# Checking your gradients in practice (for a module)

```
 1  require 'torch'
 2  require 'nn'
 3  require 'MyModules/MySin'
 4
 5  -- define inputs and module
 6  -- parameters
 7  precision = 1e-5
 8  jac = nn.Jacobian
 9
10  input = torch.Tensor():ones(2, 1)
11  module = nn.MySin(3, 2)
12
13  err = jac.testJacobian(module,input) -- test backprop, with Jacobian
14  print('==> Error: ' .. err)
15  if err<precision then
16      print('==> The module is OK')
17  else
18      print('==> The error too large, incorrect implementation')
19  end
20
21
```

*Import our module*

*Call the Jacobian module , which can check the Jacobian matrix*

*Generate an instance for our new module*

*Check the Jacobians for our new module*

# Checking your gradients in practice (for a network)

```lua
 1  local mymodel = require 'mymodel'
 2  ------------------------------------------------
 3  -- HELPER FUNCTIONS
 4  ------------------------------------------------
 5
 6  -- function that numerically checks gradient of the loss:
 7  -- f is the scalar-valued function
 8  -- g returns the true gradient (assumes input to f is a 1d tensor)
 9  -- returns difference, true gradient and estimated gradient
10  local function checkgrad(f, g, x, eps)
11
12    -- compute true gradient
13    local grad = g(x) -- this is the explicit implementation of your gradient function
14
15    -- compute numeric approximations to gradient
16    local eps = eps or 1e-5
17    local grad_est = torch.DoubleTensor(grad:size())
18    for i = 1, grad:size(1) do -- Check your gradient dimensions one at a time
19      local xorig = x[i]
20      ...
21      ...
22      ...
23      grad_est[i] = ...
24    end
25
26    -- computes (symmetric) relative error of gradient
27    local diff = ...
28    return diff, grad, grad_est
29  end
30
31  function generate_fake_data(n)
32    local data = {}
33    data.inputs = torch.randn(...)        -- random standard normal distribution for inputs
34    data.targets = torch.rand(n):mul(3):add(1):floor()  -- random integers from {1,2,3}
35    return data
36  end
37  ------------------------------------------------
```

```lua
37  ------------------------------------------------
38  -- MAIN
39  ------------------------------------------------
40
41  torch.manualSeed(1)
42  torch.setdefaulttensortype('torch.DoubleTensor')
43  precision = 1e-5
44  local data = generate_fake_data(1)
45  local model, criterion = define_model2(4, 5, 3)
46  local parameters, gradParameters = model:getParameters()
47
48  local f = function(x) -- returns the loss(params) of the network given the data and the parameters
49    ...
50  end
51
52
53  local g = function(x) -- returns dloss(params)/dparams given the data. To compute
54    -- that you first do a forward propagation, then a backpropagation step
55    ...
56  end
57
58  local err, grad, grad_est = checkgrad(f, g, parameters)
59
60  print('------------------------------------------------')
61  print('==> Error per dimension:\n')
62  print(torch.cat(grad, grad_est, 2))
63  print('------------------------------------------------')
64  print('==> Total error: ' .. err)
65  print('------------------------------------------------')
66
67  if err<precision then
68    print('==> The model is OK')
69  else
70    print('==> The error too large, something is wrong...')
71  end
72
```

*To make it faster, sample few only dimensions. Sample carefully though, when testing whole network*

*To make it faster, few only training points*
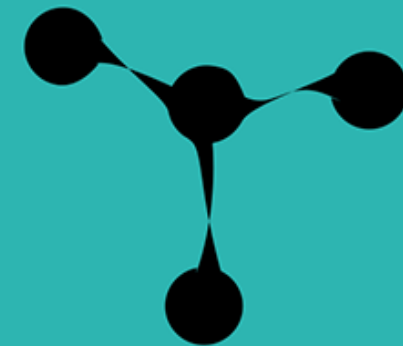
# Come up with new modules

- What about trigonometric modules

- Or polynomial modules

- Or new loss modules

- In the Lab Assignment 2 you will have the chance to think of new modules

# Implementation of basic networks and modules in Torch

UVA DEEP LEARNING COURSE
EFSTRATIOS GAVVES  & MAX WELLING

OPTIMIZING NEURAL NETWORKS IN THEORY
AND IN PRACTICE - PAGE 61

torch

Facebook AI Research

# Building a module

○ For a new module you must re-implement two functions in Torch
  ◦ One to compute the result of the forward propagation for the module `mymodule.updateOutput(…)`
  ◦ And one computing the **gradient of the loss** w.r.t. the input

  `mymodule.updateGradInput(…)` $\quad \dfrac{\partial \mathcal{L}(a_L, y)}{\partial x} = \dfrac{\partial \mathcal{L}}{\partial a_{above}} \cdot \dfrac{\partial a}{\partial x}$

  ◦ Of course you can implement other helper functions too

○ If, and only if, your module is parametric, namely has trainable parameters
  ◦ You must also implement a function for the gradient of the loss w.r.t. the parameters

  `mymodule.updateGradParameters(…)` $\quad \dfrac{\partial \mathcal{L}(a_L, y)}{\partial \theta} = \dfrac{\partial \mathcal{L}}{\partial a} \cdot \dfrac{\partial a}{\partial \theta}$

○ If your trainable parameters are boil down to a linear product $\theta \mathrm{x}$, you can simply cascade this module and avoid taking an extra gradient

$$a_1 = \theta \mathrm{x} \rightarrow a_2 = nonlinear(a_1)$$

# Make a module in Torch

```
1   local MySin, Parent = torch.class('nn.MySin', 'nn.Module')
2
3   function MySin:__init(outputsize, inputsize)
4       Parent.__init(self)
5       self.clasvar1  = ... -- Define class variables you want to use in the computations
6       self.output    = ... -- e.g. the self.output will hold the result of the forward propagation
7       self.gradInput = ... -- the gradInput will hold the gradient with respect to input, dL/dx_module
8       self.gradWeight = ... -- the gradWeight will hold the gradient with respect to params, dL/dtheta_module
9       ...
10  end
11
12  function MySin:updateOutput(input)
13      self.output = ... -- The result of forward propagation for the module
14                        -- This depends on the input of course
15      return self.output
16  end
17
18  function MySin:updateGradInput(input, gradOutput)
19      self.gradInput = ... -- The result of gradient of the module wrt input
20      return self.gradInput
21  end
22
23  -- This still needs to be understood well
24  function MySin:accGradParameters(input, gradOutput, scale)
25      self.gradWeight = ... -- If the module is parametric, you compute here the gradient wrt params
26                            -- Otherwise you do not need to reimplement the method
27  end
28
```

*Probably you will need to define some class variables*

*The forward propagation function*

*The backward propagation function wrt the input of the module*

*The backward propagation function wrt the parameters of the module*

*If module is not parametric, you don't need to implement this function*

# Summary

o We introduced how does the machine learning paradigm for neural networks

o We described the backpropagation algorithm, which is the backbone for neural network training

o We explained the neural network in terms of modular architecture and described various possible architectures

o We described different neural network modules, as well as how to implement and how to check your own module

# Next lecture

o We are going to see how to use backpropagation to optimize our neural network

o We are going to review different methods and algorithms for optimizing our neural network, especially our deep networks, better

o We are going to revisit different learning paradigms, e.g. what loss functions should be used for different machine learning tasks

o And if we have time, some more advanced modules