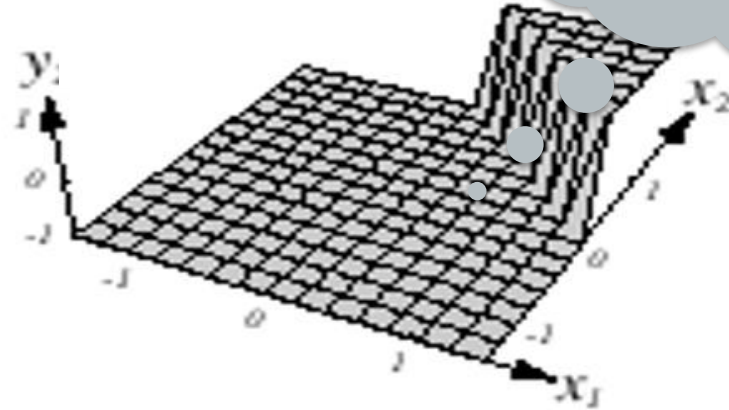
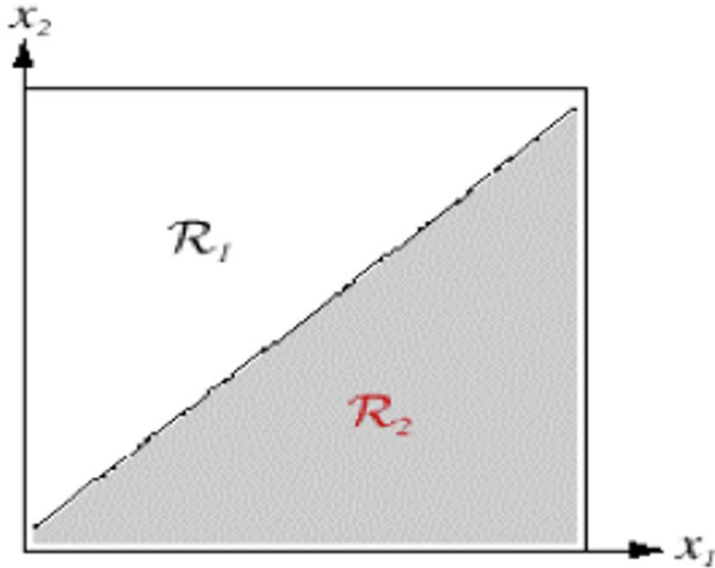


Ralf Möller, Sylvia Melzer

Training Embedding Architectures

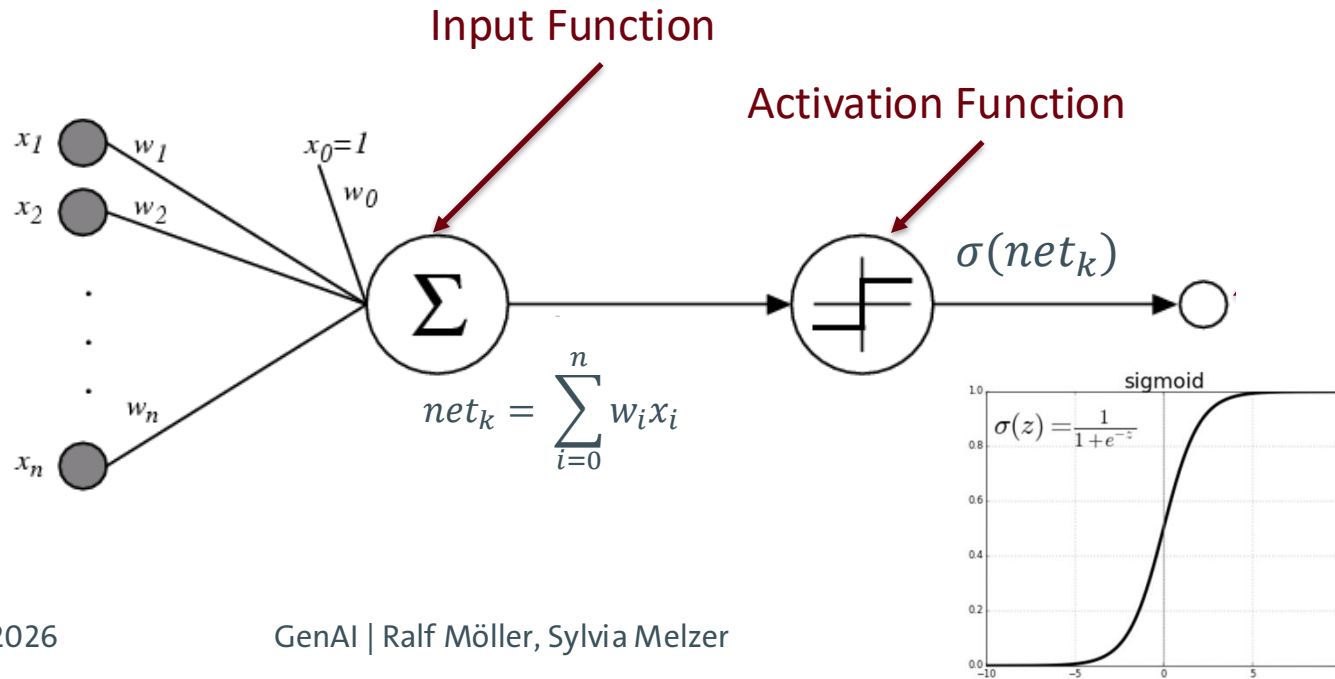
Linear Separation (Linear Models)



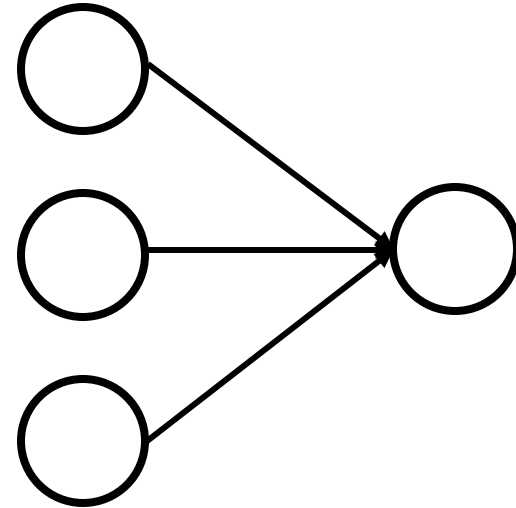
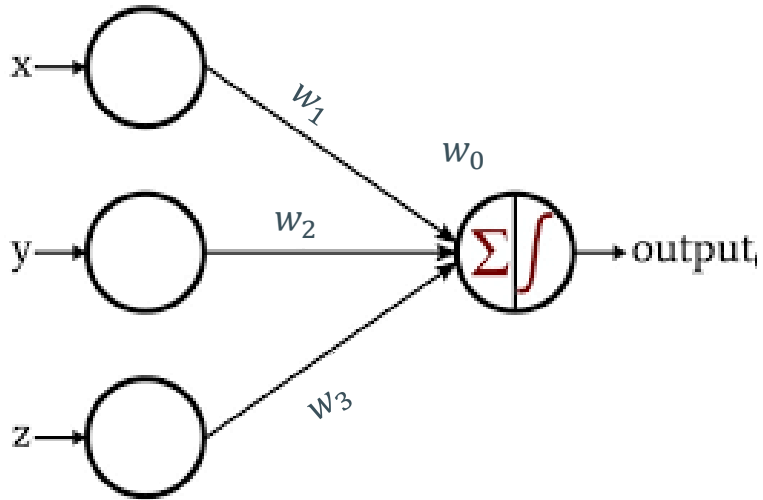
Embedding
(x_1, x_2) data
into y

Perceptron

Frank Rosenblatt, The Perceptron—a perceiving and recognizing automaton. Report 85-460-1, Cornell Aeronautical Laboratory, 1957



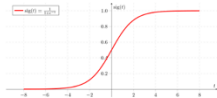
Perceptron (simplified views)



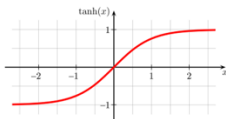
Other Examples for Activation Functions

- Linear activation: $a_k = f_{act}(net_k) = c_k \cdot net_k$
- Threshold activation: $a_k = f_{act}(net_k) = \begin{cases} 1, & \text{falls } net_k \geq \theta_k, \\ 0, & \text{sonst.} \end{cases}$
- Logistic activation: $a_k = f_{act}(net_k) = \frac{1}{1 + e^{-\frac{net_k}{T}}}$
- Hyperbolic tangent: $a_k = f_{act}(net_k) = \frac{e^{net_k} - e^{-net_k}}{e^{net_k} + e^{-net_k}} = \frac{1 + \tanh(net_k)}{2}$

Threshold (often 0)

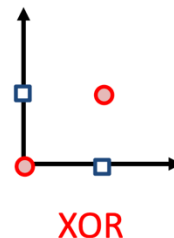
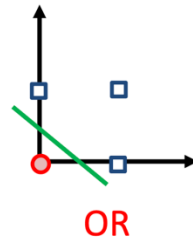
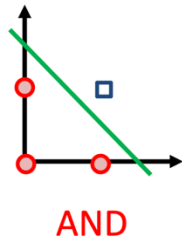


Special case Sigmoid: T=1



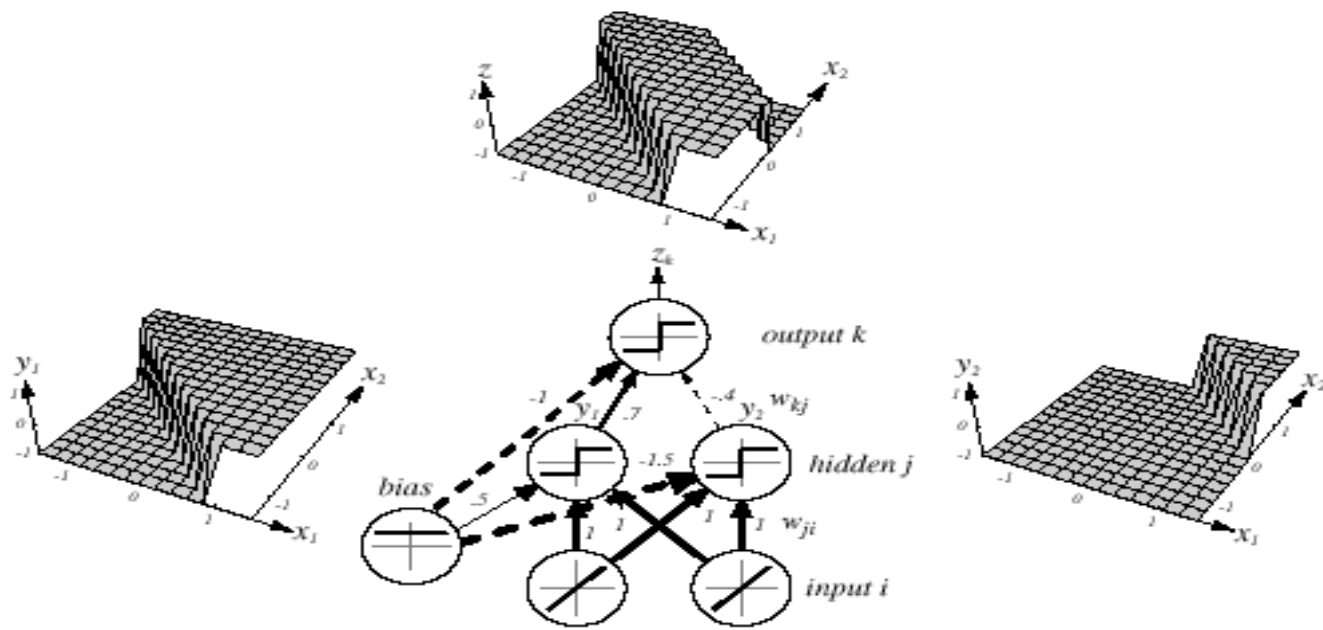
Limitations of Linear Models

Some function are linearly separable, but many are not

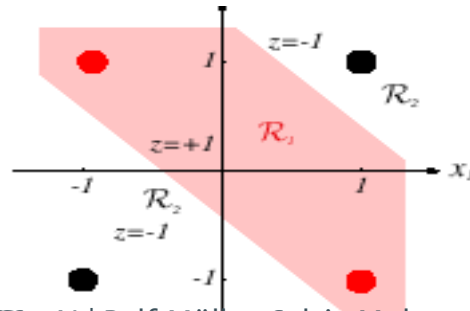


Possible solution: **composition**

$$x_1 \text{ XOR } x_2 = (x_1 \text{ OR } x_2) \text{ AND not}(x_1 \text{ AND } x_2)$$



$$Z = x_1 \text{ XOR } x_2 = (x_1 \text{ OR } x_2) \text{ AND NOT}(x_1 \text{ AND } x_2)$$



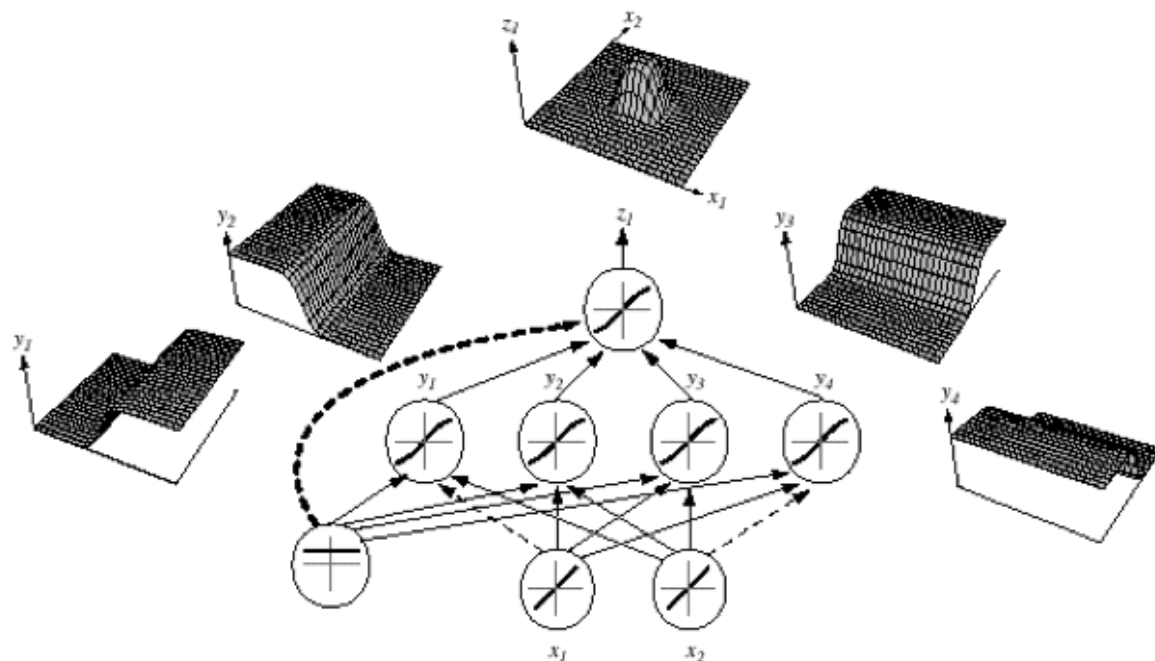


FIGURE 6.2. A 2-4-1 network (with bias) along with the response functions at different units; each hidden output unit has sigmoidal activation function $f(\cdot)$. In the case shown, the hidden unit outputs are paired in opposition thereby producing a “bump” at the output unit. Given a sufficiently large number of hidden units, any continuous function from input to output can be approximated arbitrarily well by such a network. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

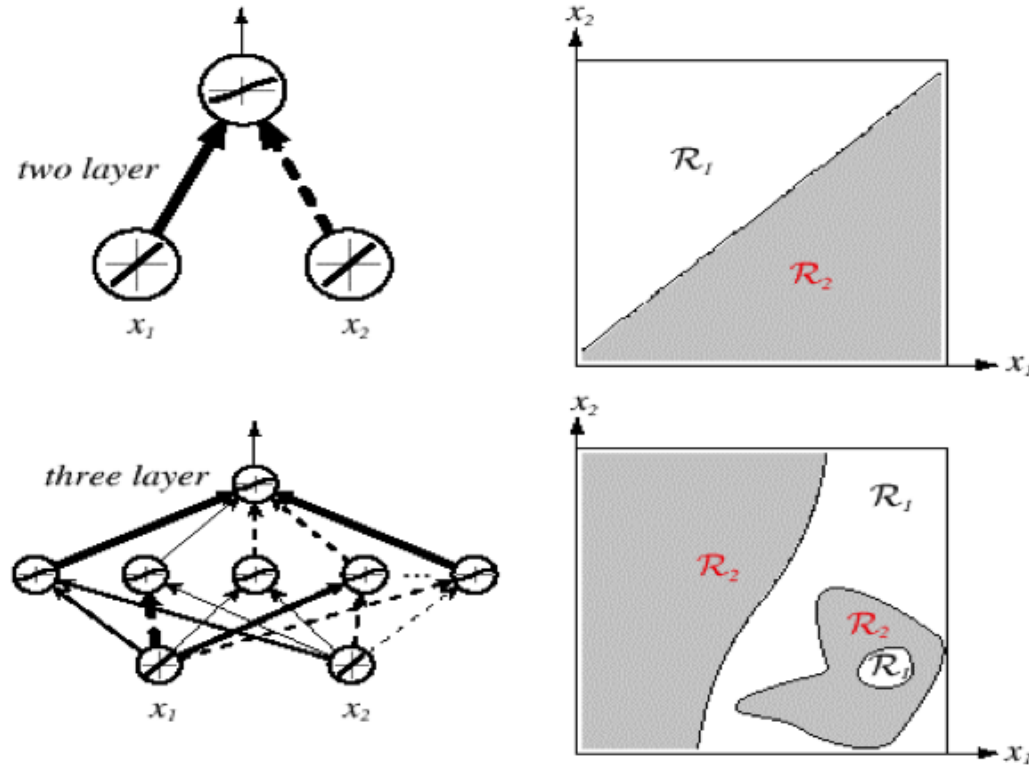
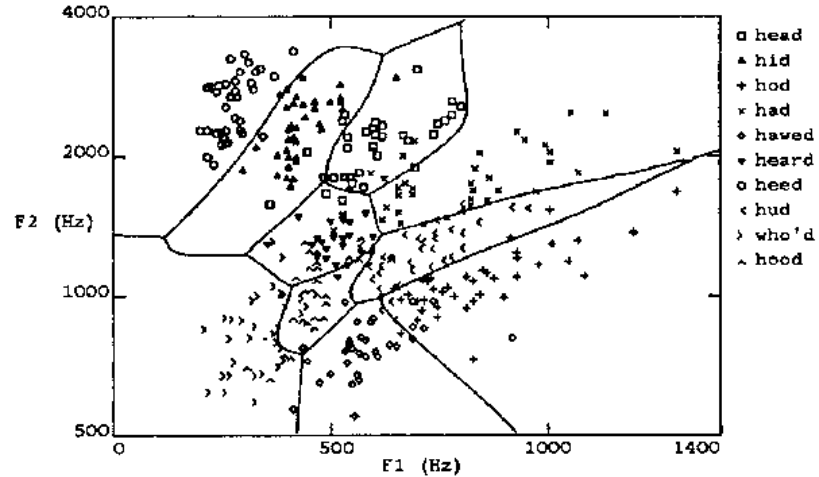
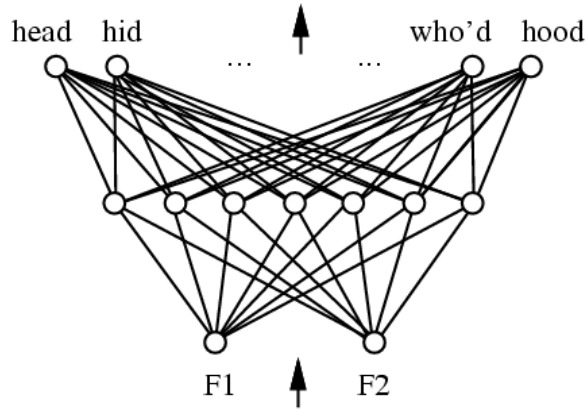


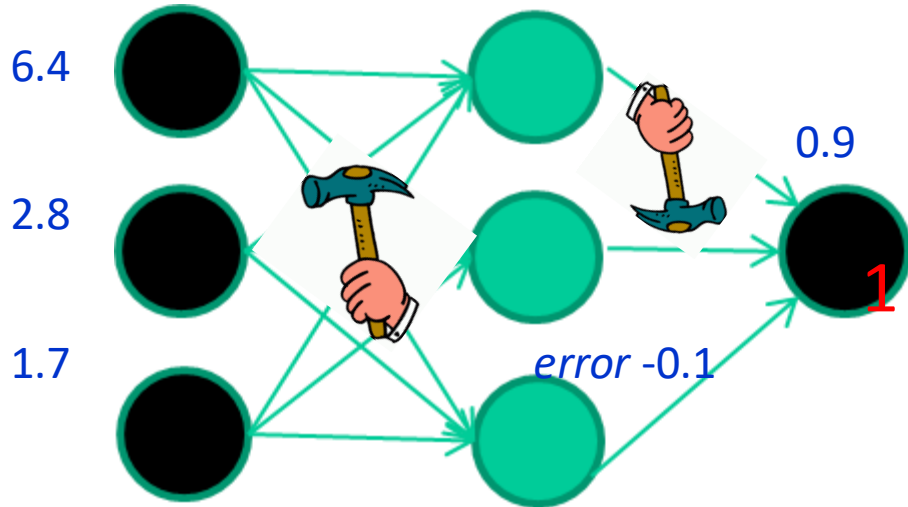
FIGURE 6.3. Whereas a two-layer network classifier can only implement a linear decision boundary, given an adequate number of hidden units, three-, four- and higher-layer networks can implement arbitrary decision boundaries. The decision regions need not be convex or simply connected. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Multi-level networks of sigmoid units



Adaptation of mapping parameters (“weights”)

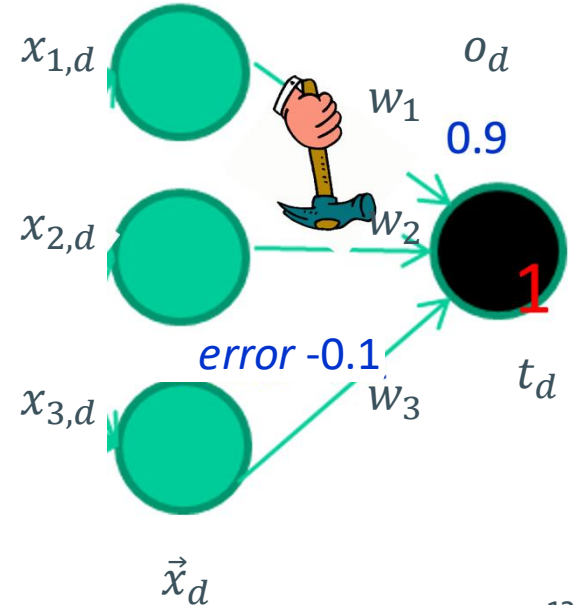
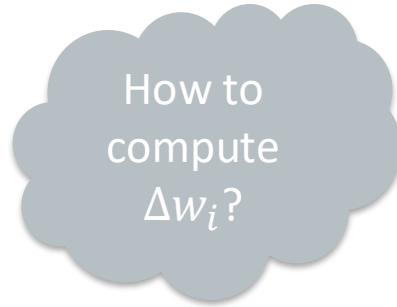
<i>Fields</i>	<i>class</i>
1.4 2.7 1.9	0
3.8 3.4 3.2	0
6.4 2.8 1.7	1
4.1 0.1 0.2	0
etc ...	



Simple case $o = Wx$

Delta Learning Rule for Perceptrons

- Given: Set T of training data $\{(\vec{x}_1, t_1), \dots, (\vec{x}_N, t_N)\}$
- Select training datum $d \in T$:
 - Input \vec{x}_d
 - Target t_d
- Output o_d
- $w_i^{k+1} = \Delta w_i + w_i^k$



Derivation of a suitable update value Δw_i

- Idea: minimize quadratic error
 -
 - t_d value for $d \in T$ an arbitrarily selected (“sampled”) datum $d = (\vec{x}_d, t_d)$ from dataset T
 - o_d output for \vec{x}_d

$$L[\vec{w}] = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

- Determine minimum with 1st derivative

$$L[\vec{w}] = \frac{1}{2} (t_d - o_d)^2$$

Gradient Descent

- Gradient

$$\Delta L[\vec{w}]$$

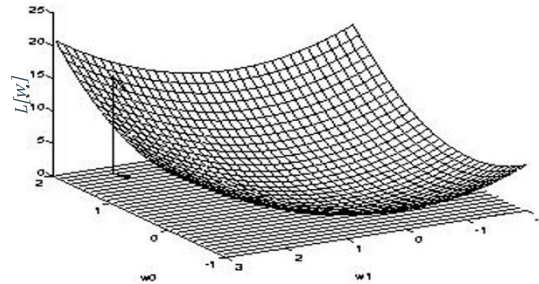
$$\equiv \left[\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_n} \right]$$

- Learning rule

$$\Delta \vec{w} = -\eta \Delta L[\vec{w}]$$

- i.e.

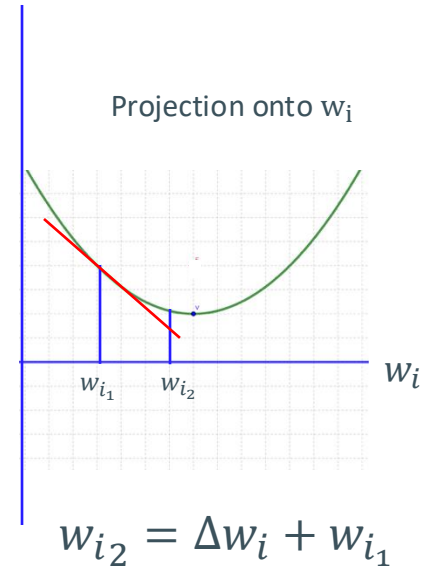
$$\Delta \vec{w}_i = -\eta \frac{\partial L}{\partial w_i}(\vec{w})$$



Visualization of L for $\vec{w} = (w_0, w_1)$

How to compute the derivative

$L(w_i)$



What if only linear mappings were considered?

$$\begin{aligned}\frac{\partial L}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} (t_d - o_d)^2 \\ &= \frac{1}{2} \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \\ &= (t_d - o_d) (-x_{i,d}) = -(t_d - o_d)x_{i,d}\end{aligned}$$

Multivariate
linear
regression

$$\Delta \vec{w}_i = -\eta \frac{\partial L}{\partial w_i} (\vec{w})$$

Iteration for sampled data points d
 $\Delta w_i = \eta (t_d - o_d) x_{i,d}$

$$L[\vec{w}] = \frac{1}{2} (t_d - o_d)^2$$

With sigmoid

Multivariate
logistic
regression

$$\frac{\partial L}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} (t_d - o_d)^2$$

$$= \frac{1}{2} \frac{\partial}{\partial w_i} (t_d - o_d)^2$$

$$= \frac{1}{2} 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d)$$

$$= (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \sigma(\vec{w} \cdot \vec{x}_d))$$

$$= (t_d - o_d) (-x_{i,d} \sigma'(\vec{w} \cdot \vec{x}_d))$$

$$= -(t_d - o_d) x_{i,d} \sigma'(\vec{w} \cdot \vec{x}_d)$$

$$= -(t_d - o_d) x_{i,d} \sigma(\vec{w} \cdot \vec{x}_d) (1 - \sigma(\vec{w} \cdot \vec{x}_d))$$

$$= -(t_d - o_d) x_{i,d} o_d (1 - o_d)$$

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

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Stochastic Gradient Descent (for sampled data point)

- Gradient

$$\Delta L[\vec{w}] \equiv \left[\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_n} \right]$$

- Learning rule

$$\Delta \vec{w} = -\eta \Delta L[\vec{w}]$$

- i.e.

$$\Delta \vec{w}_i = -\eta \frac{\partial L}{\partial w_i}(\vec{w})$$

$$\frac{\partial L}{\partial w_i} = -(t_d - o_d)x_{i,d}o_d(1 - o_d)$$

Iteration for sampled data points d

$$\Delta w_i = \eta(t_d - o_d)x_{i,d}o_d(1 - o_d)$$

Stochastic Gradient Descent for batch of data

Given **batch of data** D sampled from from training set T :

$$D = \{d_1, \dots, d_N\} \subseteq T:$$

$$\Delta w_i = \eta \frac{1}{N} \sum_{d \in D} (t_d - o_d) o(1 - o_d) x_{i,d}$$

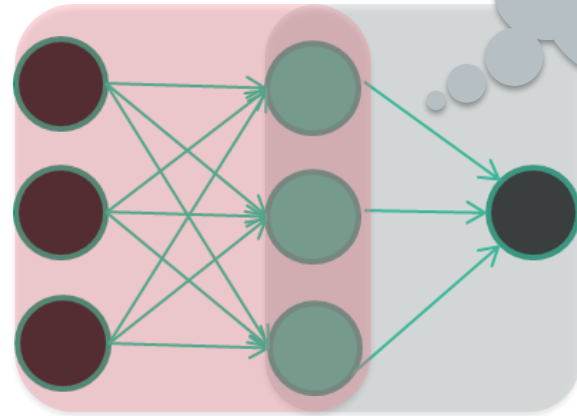
$$\Delta w_i = \eta \frac{1}{N} \sum_{d \in D} (t_d - o_d) o(1 - o_d) x_{i,d}$$

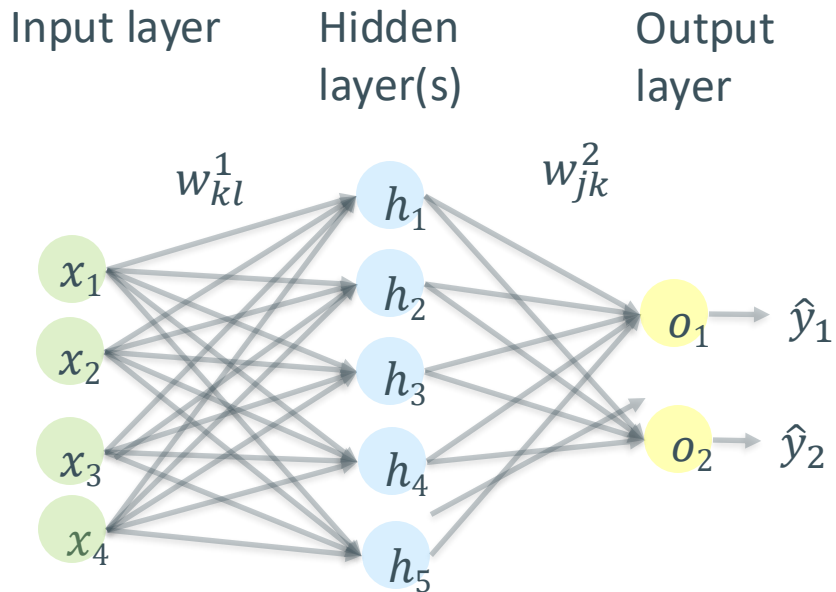
Learning Parameters (“Hyper parameters”)

- Learning rate η
- Batch size
 - Number of samples processed before model is updated
- Epochs
 - Number of complete passes through the training dataset

Multilayer Perceptron

- General layer: Parallel composition of perceptrons
- Multilayer perceptron: series of parallel compositions
- How to train based on forward pass?
- We can adjust the mapping parameters of the last layer
- How to adjust parameters of previous layers?
 - Problem: no targets t for previous layer given
 - Idea: Propagate change from o_d to $o_{d,new}$ to input value changes (green circles) to obtain a target t for previous layer (→ Backpropagation)





$$\mathbb{R}^4 \xrightarrow{g} \mathbb{R}^5 \xrightarrow{f} \mathbb{R}^2$$

$$\mathbf{x} \mapsto \mathbf{h} \mapsto \hat{\mathbf{y}}$$

$$\hat{\mathbf{y}} = \mathbf{f}(\mathbf{g}(\mathbf{x}; \mathbf{W}^1); \mathbf{W}^2)$$

k : Layer k
 \mathbf{W}^k : Weight matrix for k



$$= \sigma(\mathbf{W}^2 \sigma(\mathbf{W}^1 \mathbf{x}))$$

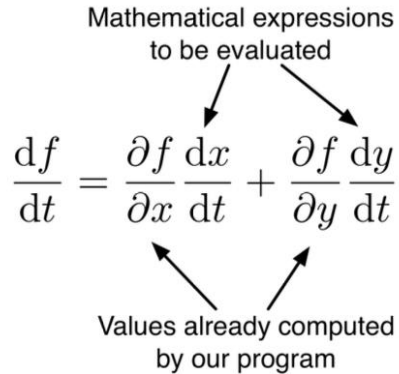
Recap: Multivariate chain rule $f(x, y)$

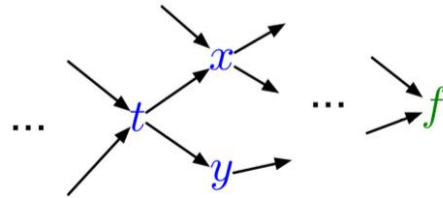
- In the context of backpropagation:

Mathematical expressions to be evaluated

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

Values already computed by our program



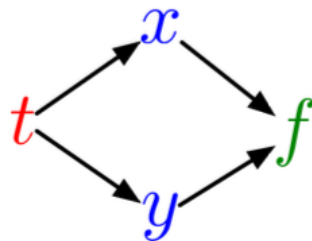


- In our notation:

$$\bar{t} = \bar{x} \frac{dx}{dt} + \bar{y} \frac{dy}{dt}$$

- Suppose we have a function $f(x, y)$ and functions $x(t)$ and $y(t)$. (All the variables here are scalar-valued.) Then

$$\frac{d}{dt} f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$



- Example:

$$f(x, y) = y + e^{xy}$$

$$x(t) = \cos t$$

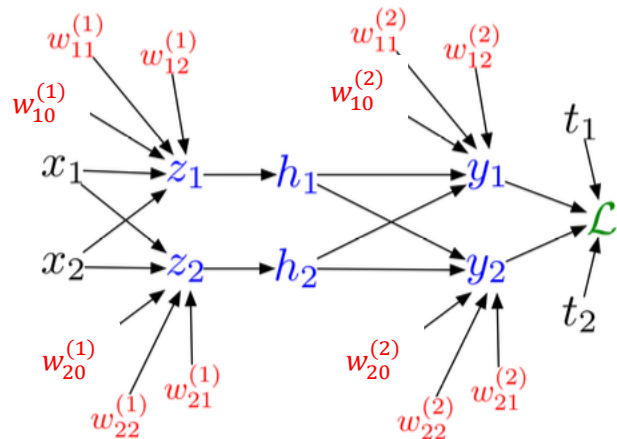
$$y(t) = t^2$$

- Plug in to Chain Rule:

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

$$= (ye^{xy}) \cdot (-\sin t) + (1 + xe^{xy}) \cdot 2t$$

Multilayer Perceptron (multiple outputs):



Backward pass:

$$\bar{\mathcal{L}} = 1$$

$$\bar{y}_k = \bar{\mathcal{L}} (y_k - t_k)$$

$$\bar{w}_{ki}^{(2)} = \bar{y}_k h_i$$

$$\bar{h}_i = \sum_k \bar{y}_k w_{ki}^{(2)}$$

$$\bar{z}_i = \bar{h}_i \sigma'(z_i)$$

$$\bar{w}_{ij}^{(1)} = \bar{z}_i x_j$$

Forward pass:

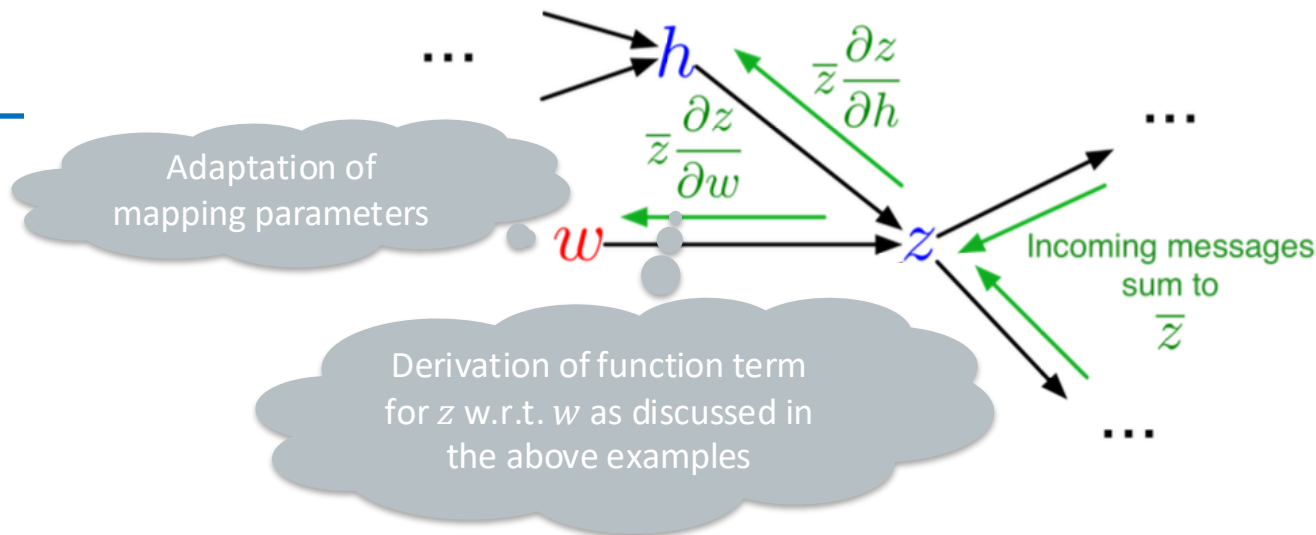
$$z_i = \sum_j w_{ij}^{(1)} x_j$$

$$h_i = \sigma(z_i)$$

$$y_k = \sum_i w_{ki}^{(2)} h_i$$

$$\mathcal{L} = \frac{1}{2} \sum_k (y_k - t_k)^2$$

Backprop as message passing:



- Each node receives a bunch of messages from its children, which it aggregates to get its error signal. It then passes messages to its parents.
- This provides modularity, since each node only has to know how to compute derivatives with respect to its arguments, and doesn't have to know anything about the rest of the graph.

Full backpropagation algorithm:

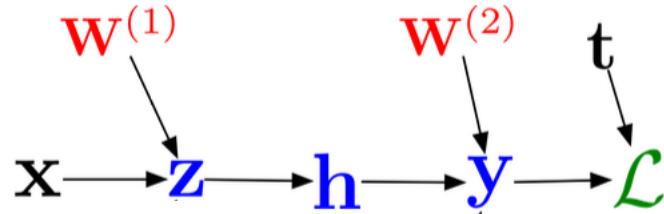
Let v_1, \dots, v_N be a **topological ordering** of the computation graph
(i.e. parents come before children.)

v_N denotes the variable we're trying to compute derivatives of (e.g. loss).

forward pass $\left[\begin{array}{l} \text{For } i = 1, \dots, N \\ \text{Compute } v_i \text{ as a function of } \text{Pa}(v_i) \end{array} \right.$

backward pass $\left[\begin{array}{l} \overline{v}_N = 1 \\ \text{For } i = N - 1, \dots, 1 \\ \overline{v}_i = \sum_{j \in \text{Ch}(v_i)} \overline{v}_j \frac{\partial v_j}{\partial v_i} \end{array} \right.$

In vectorized form:



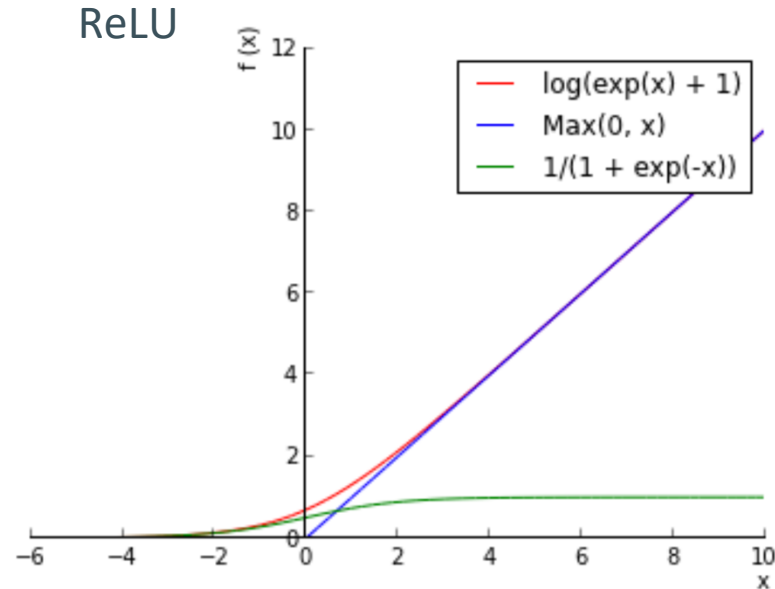
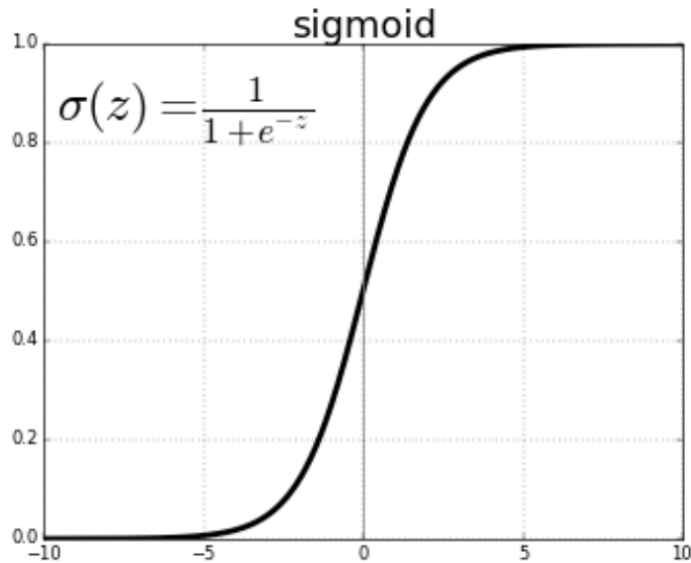
Forward pass:

$$\begin{aligned} \mathbf{z} &= \mathbf{W}^{(1)} \mathbf{x} \\ \mathbf{h} &= \sigma(\mathbf{z}) \\ \mathbf{y} &= \mathbf{W}^{(2)} \mathbf{h} \\ \mathcal{L} &= \frac{1}{2} \|\mathbf{t} - \mathbf{y}\|^2 \end{aligned}$$

Backward pass:

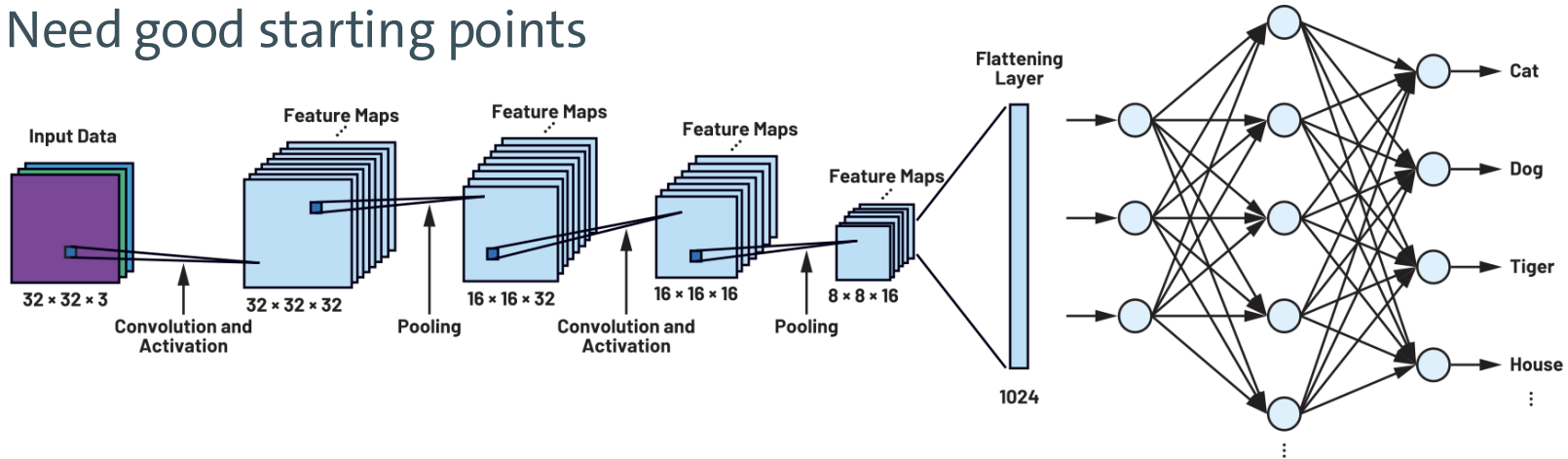
$$\begin{aligned} \bar{\mathcal{L}} &= 1 \\ \bar{\mathbf{y}} &= \bar{\mathcal{L}} (\mathbf{y} - \mathbf{t}) \\ \overline{\mathbf{W}^{(2)}} &= \bar{\mathbf{y}} \mathbf{h}^\top \\ \overline{\mathbf{b}^{(2)}} &= \bar{\mathbf{y}} \\ \bar{\mathbf{h}} &= \mathbf{W}^{(2)\top} \bar{\mathbf{y}} \\ \bar{\mathbf{z}} &= \bar{\mathbf{h}} \circ \sigma'(\mathbf{z}) \\ \overline{\mathbf{W}^{(1)}} &= \bar{\mathbf{z}} \mathbf{x}^\top \\ \overline{\mathbf{b}^{(1)}} &= \bar{\mathbf{z}} \end{aligned}$$

Sigmoid vs. ReLU (Rectified Linear Unit)

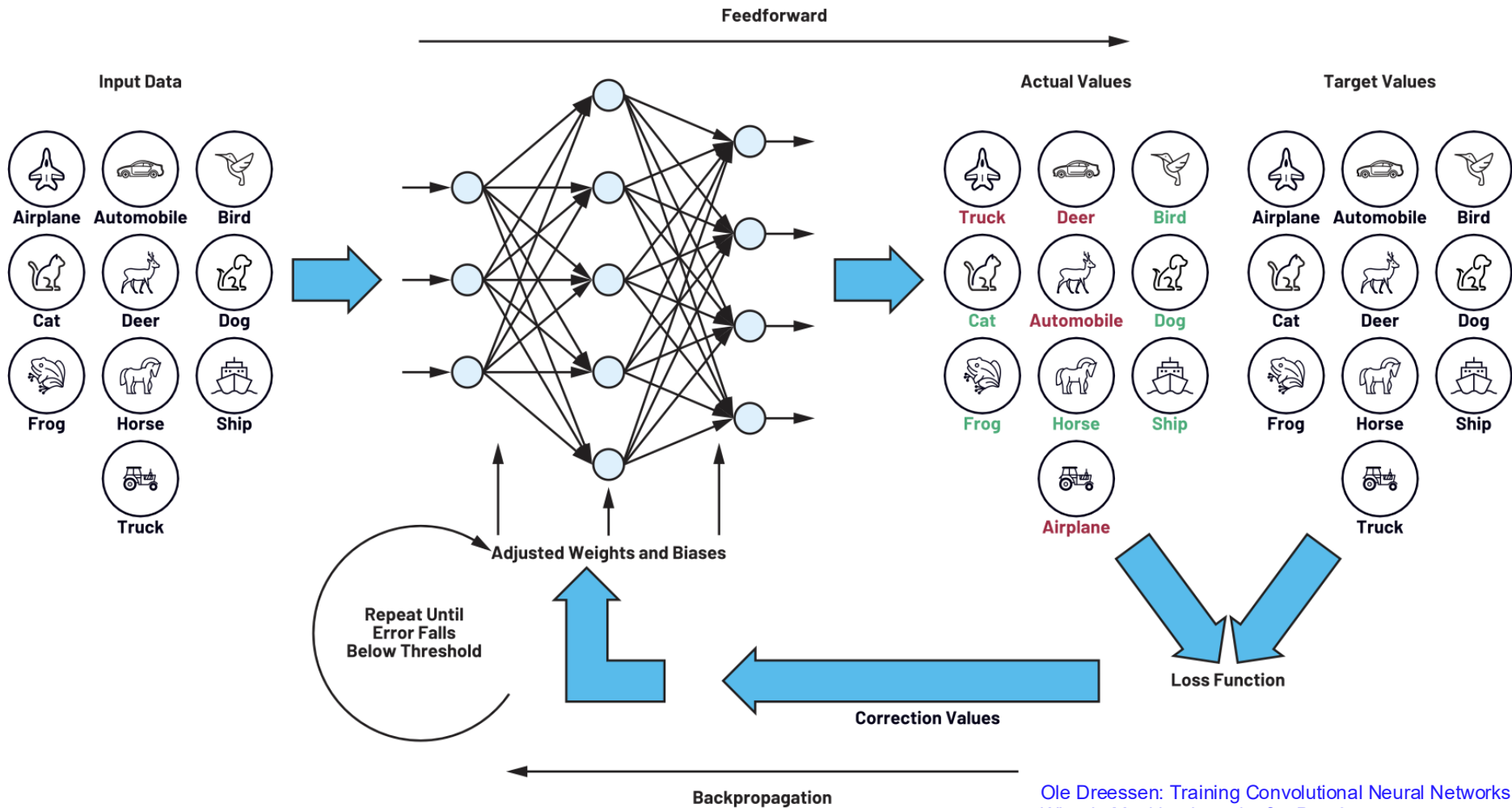


Autoencoder

- Training of convolutional networks
- Need good starting points



Ole Dreesen: Training Convolutional Neural Networks:
What Is Machine Learning?—Part 2



Softmax output

