

# Artificial Neural Networks 2

Morten Nielsen  
Department of Health Technology,  
DTU

---

# Outline

---

- Optimization procedures
    - Gradient decent (this you already know)
  - Network training
    - back propagation
    - cross-validation
    - Over-fitting
    - Examples
    - Deep learning
-

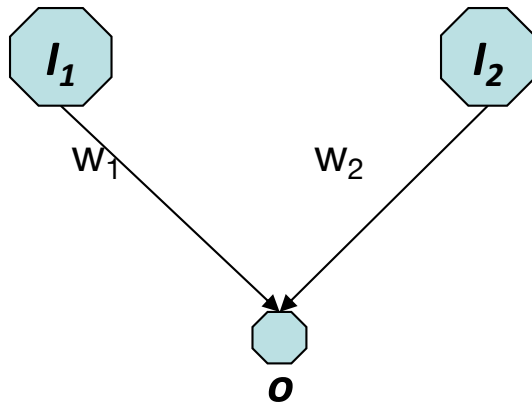
# Neural network. Error estimate

---

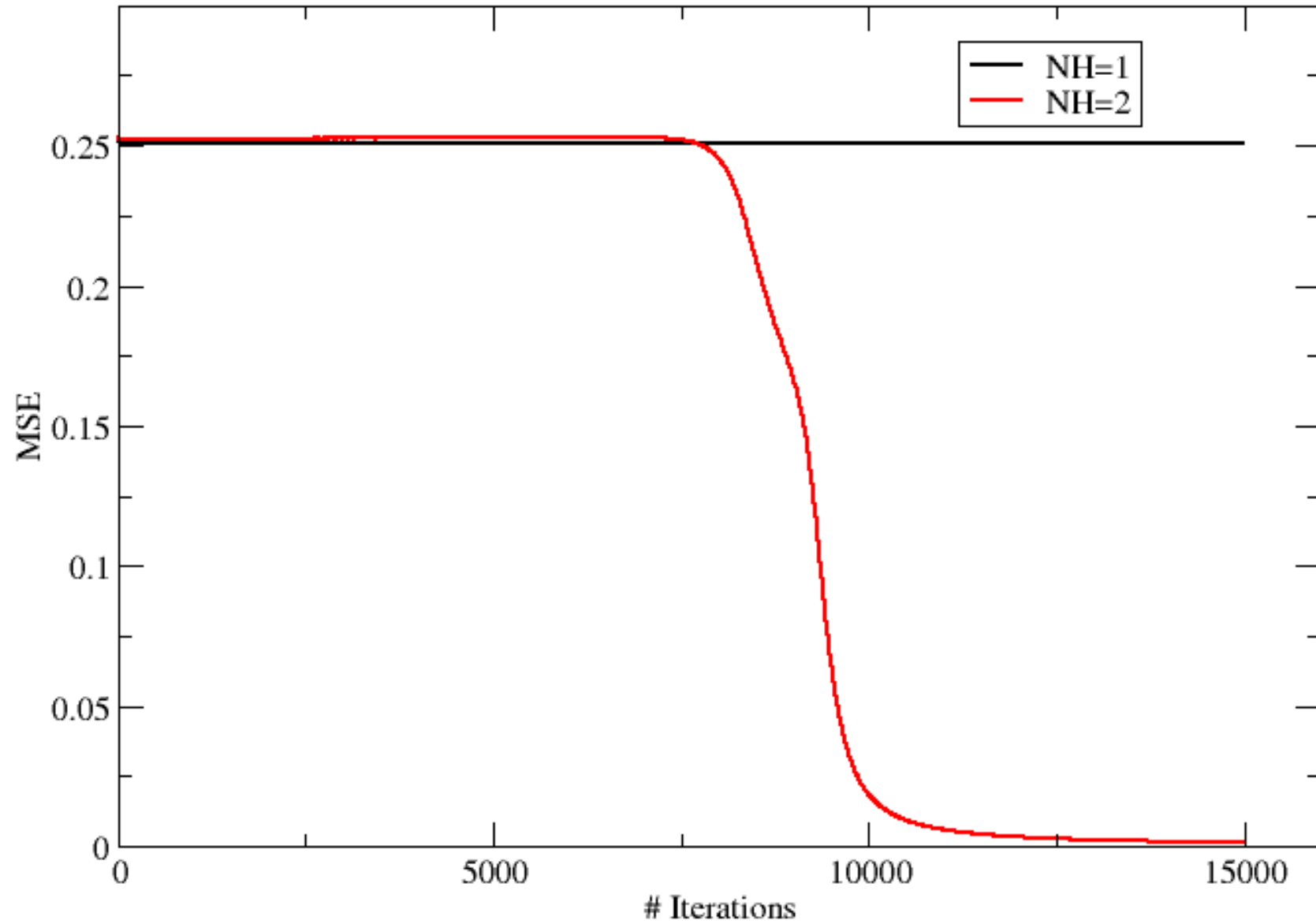
Linear function

$$o = I_1 \cdot w_1 + I_2 \cdot w_2$$

$$E = \frac{1}{2} \cdot (o - t)^2$$



# Neural networks

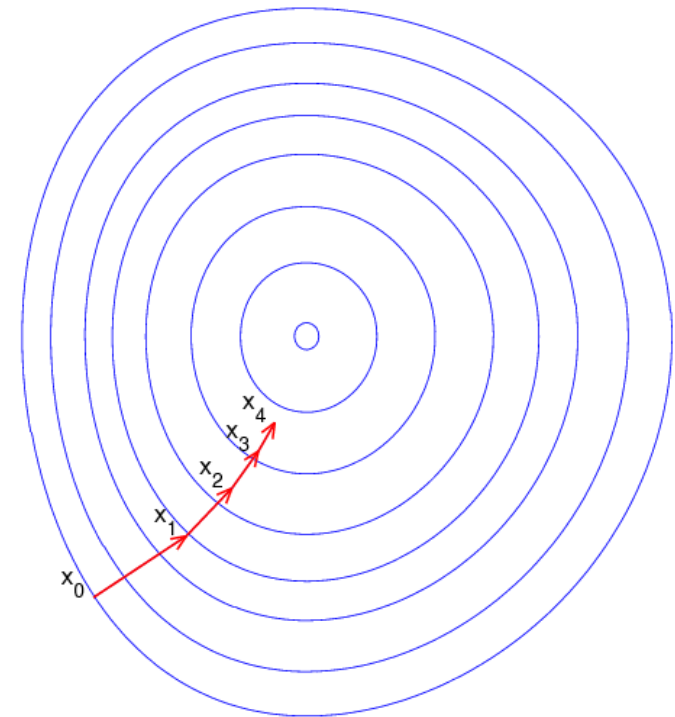


# Gradient descent (from wikipedia)

Gradient descent is based on the observation that if the real-valued function  $F(x)$  is defined and differentiable in a neighborhood of a point  $a$ , then  $F(x)$  decreases fastest if one goes from  $a$  in the direction of the negative gradient of  $F$  at  $a$ . It follows that, if

$$b = a - \varepsilon \cdot \nabla F(a)$$

for  $\varepsilon > 0$  a small enough number, then  $F(b) < F(a)$



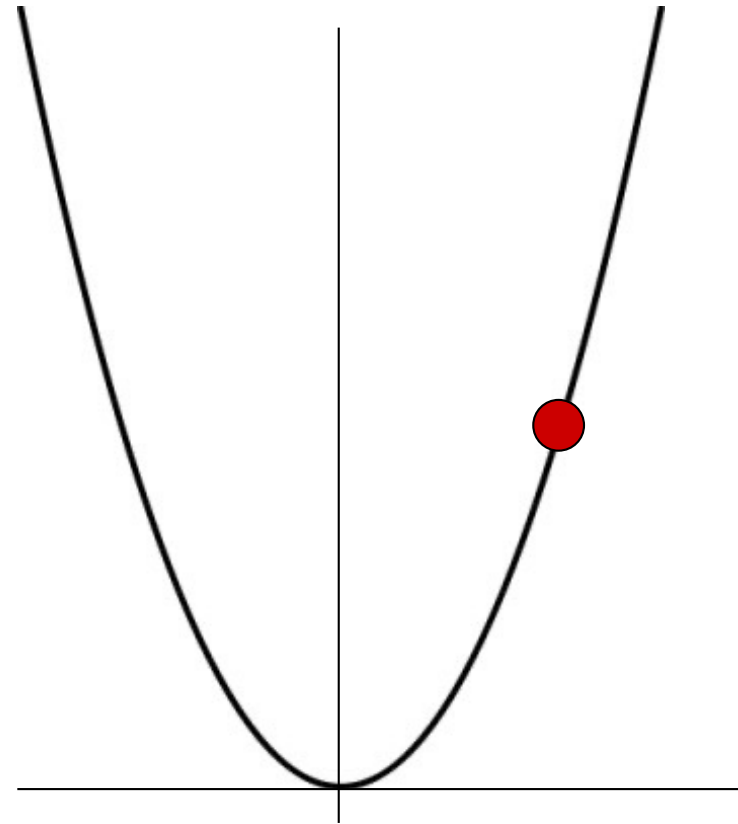
# Gradient decent (example)

---

$$F(x) = x^2$$

$$a = 2$$

$$F(a) = 4$$



# Gradient decent (example)

---

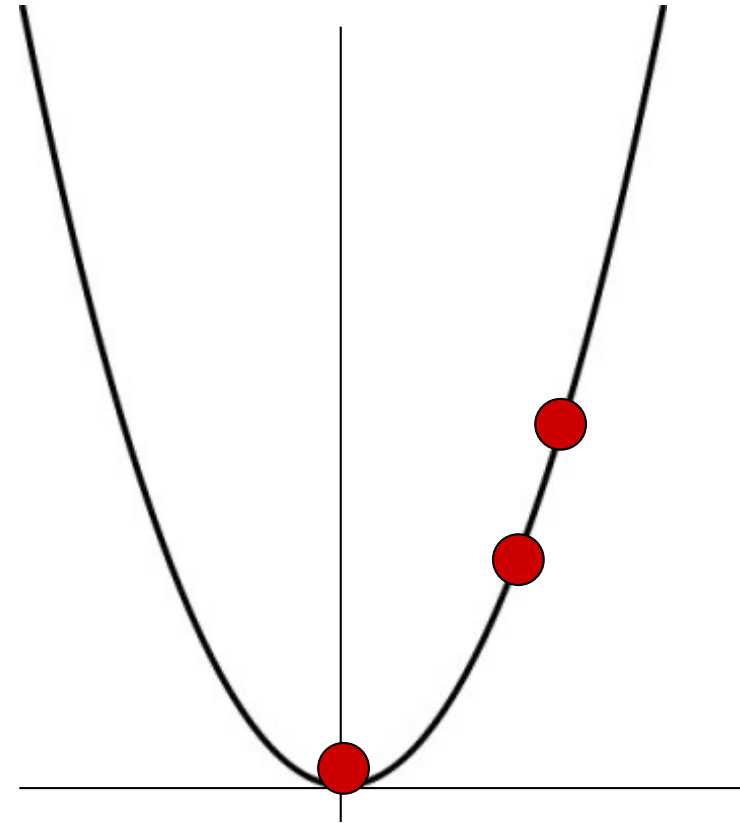
$$F(x) = x^2$$

$$a = 2$$

$$F(a) = 4$$

$$\frac{\partial F}{\partial x} = 2 \cdot x = 4$$

$$b = a - \varepsilon \cdot \nabla F(a) = 2 - 0.1 \cdot 4 = 1.6$$



# Gradient decent. Example

Weights are changed in the opposite direction of the gradient of the error

$$w_i' = w_i + \Delta w_i$$

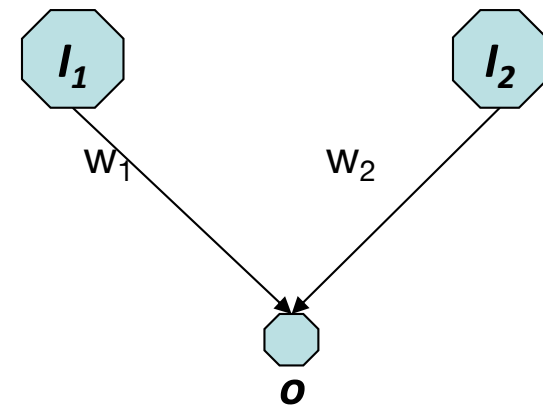
$$E = \frac{1}{2} \cdot (O - t)^2$$

$$O = \sum_i w_i \cdot I_i$$

$$\Delta w_i = -\varepsilon \cdot \frac{\partial E}{\partial w_i} = -\varepsilon \cdot \frac{\partial E}{\partial O} \cdot \frac{\partial O}{\partial w_i} = -\varepsilon \cdot (O - t) \cdot I_i$$

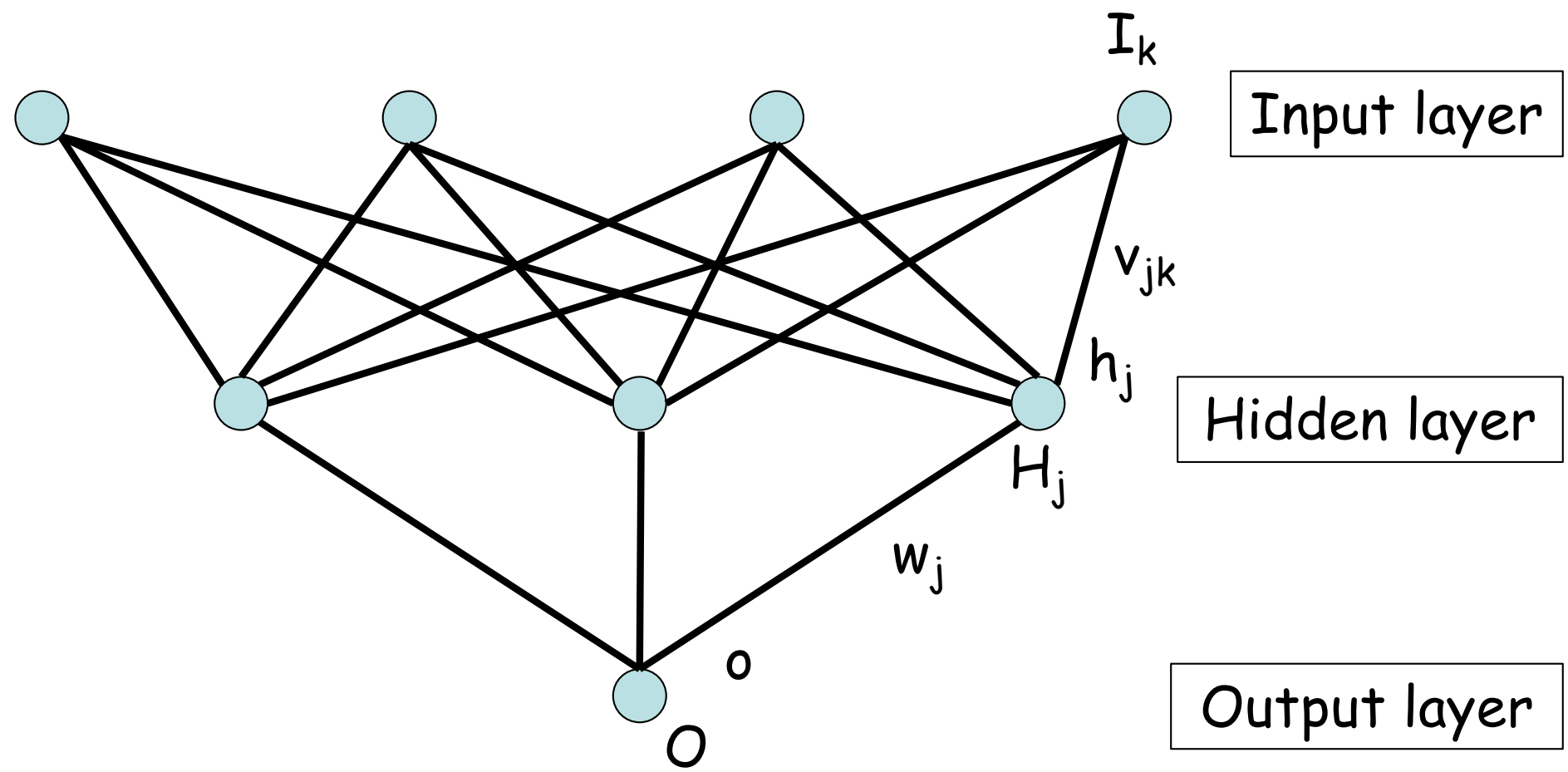
Linear function

$$O = I_1 \cdot w_1 + I_2 \cdot w_2$$





# Network architecture



# What about the hidden layer?

$$\Delta w_i = -\varepsilon \cdot \frac{\partial E}{\partial w_i}$$

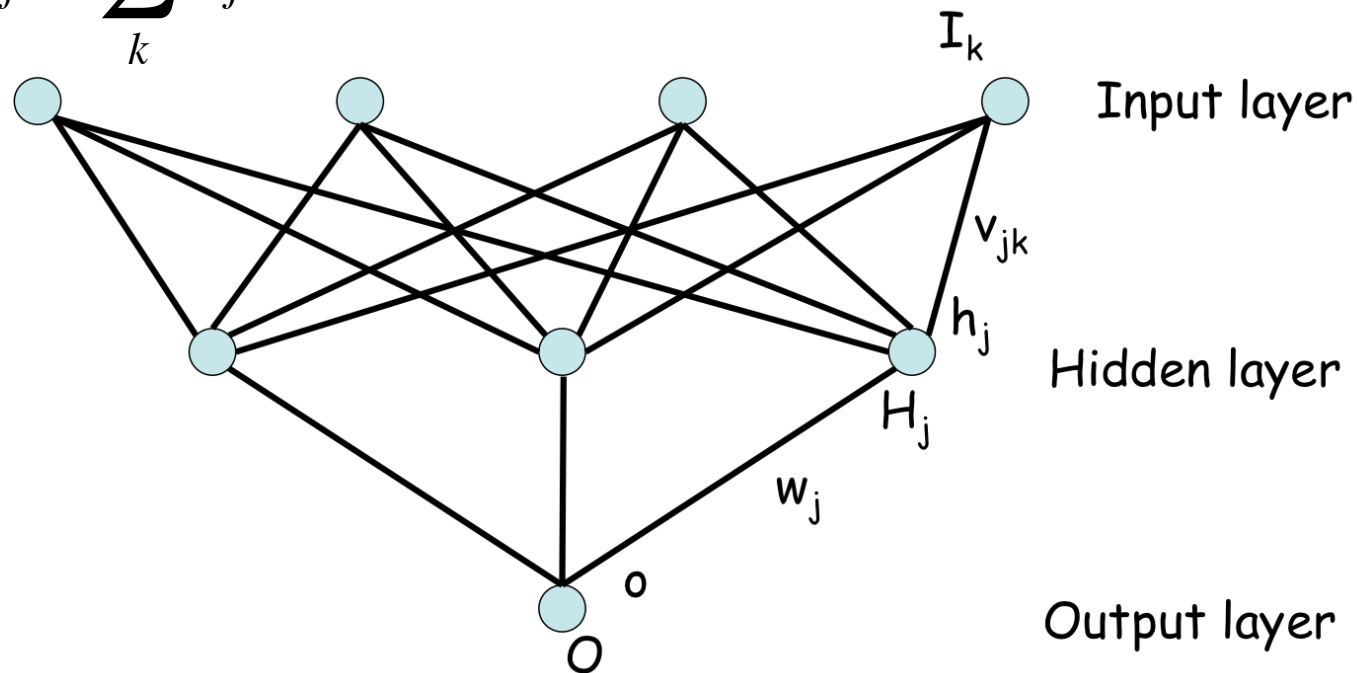
$$E = \frac{1}{2} \cdot (O - t)^2$$

$$o = \sum_j w_j \cdot H_j$$

$$O = g(o), H = g(h)$$

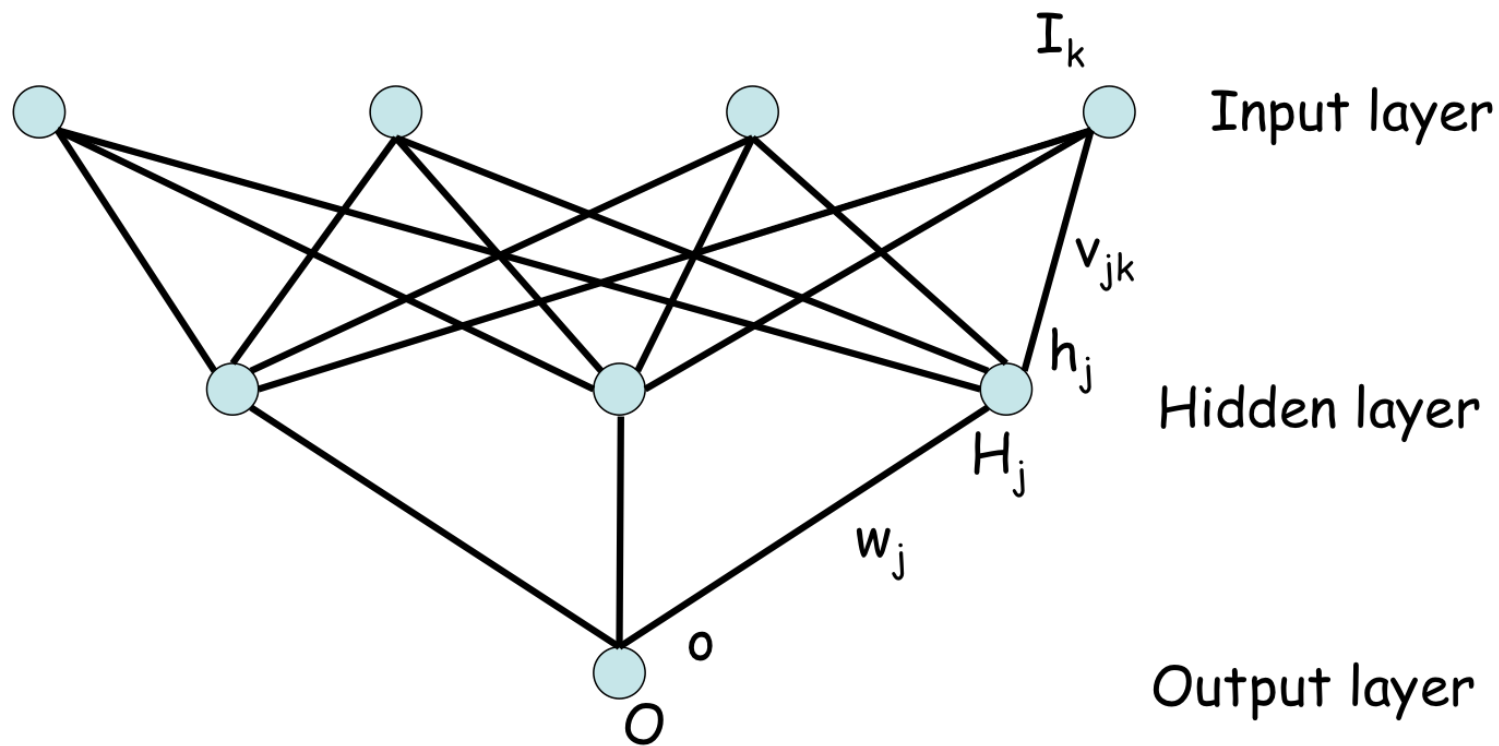
$$h_j = \sum_k v_{jk} \cdot I_k$$

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



# Hidden to output layer

$$\frac{\partial E}{\partial w_j} = \frac{\partial E(O(o(w_j)))}{\partial w_j} = \frac{\partial E}{\partial O} \cdot \frac{\partial O}{\partial o} \cdot \frac{\partial o}{\partial w_j}$$



# Hidden to output layer

$$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial O} \cdot \frac{\partial O}{\partial o} \cdot \frac{\partial o}{\partial w_j} =$$

$$\frac{\partial E}{\partial O} = (O - t)$$

$$\frac{\partial O}{\partial o} = \frac{\partial g}{\partial o} = ?$$

$$\frac{\partial o}{\partial w_j} =$$

$$O = g(o)$$

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

$$g'(x) = \frac{-1}{(1 + e^{-x})^2} \cdot (-e^{-x})$$

$$= (1 - g(x)) \cdot g(x)$$

# Hidden to output layer

$$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial O} \cdot \frac{\partial O}{\partial o} \cdot \frac{\partial o}{\partial w_j} =$$

$$\frac{\partial E}{\partial O} = (O - t)$$

$$\frac{\partial O}{\partial o} = \frac{\partial g}{\partial o} = ?$$

$$\frac{\partial o}{\partial w_j} = \frac{1}{\partial w_j} \sum_l w_l \cdot H_l = H_j$$

$$O = g(o)$$

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

$$g'(x) = \frac{-1}{(1 + e^{-x})^2} \cdot (-e^{-x})$$

$$= (1 - g(x)) \cdot g(x)$$

# Hidden to output layer

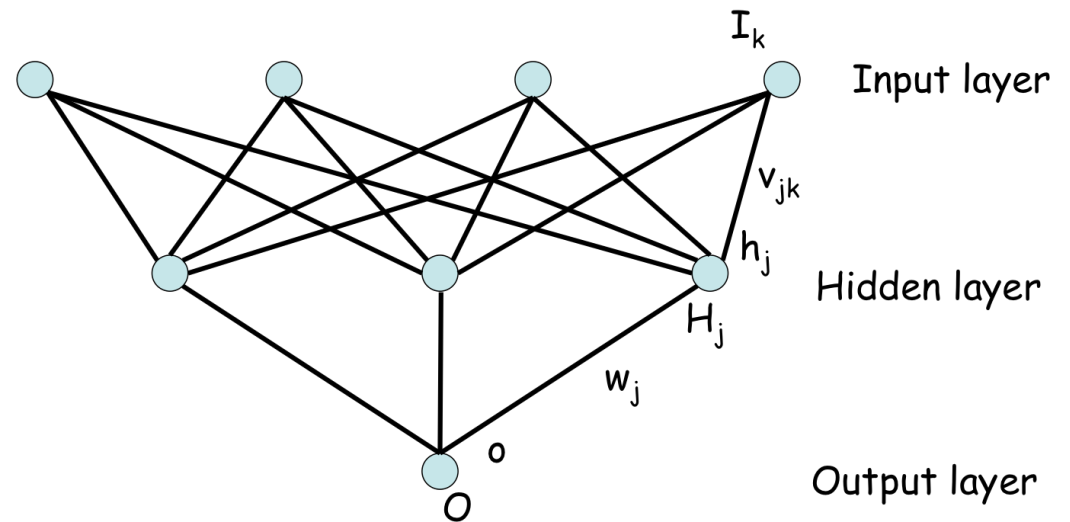
$$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial O} \cdot \frac{\partial O}{\partial o} \cdot \frac{\partial o}{\partial w_j} = (O - t) \cdot g'(o) \cdot H_j$$

$$= (O - t) \cdot (1 - O) \cdot O \cdot H_j$$

$$O = g(o)$$

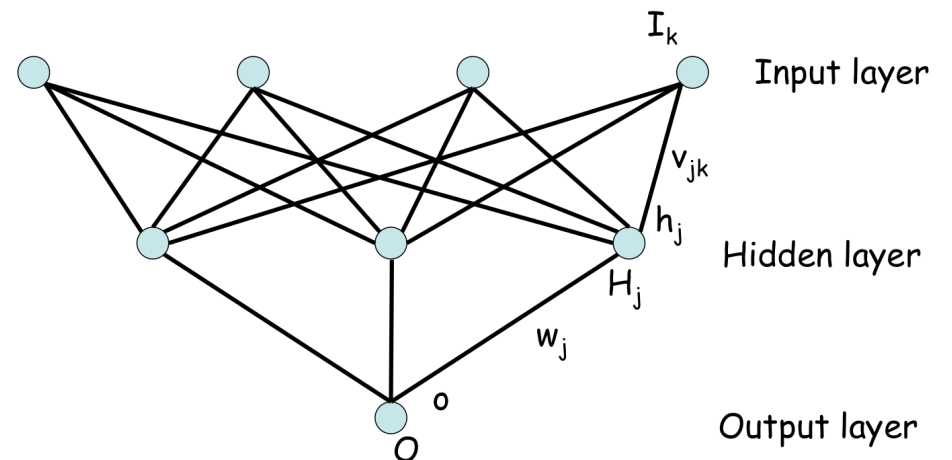
$$g'(o) = (1 - g(o)) \cdot g(o)$$

$$= (1 - O) \cdot O$$



# Input to hidden layer

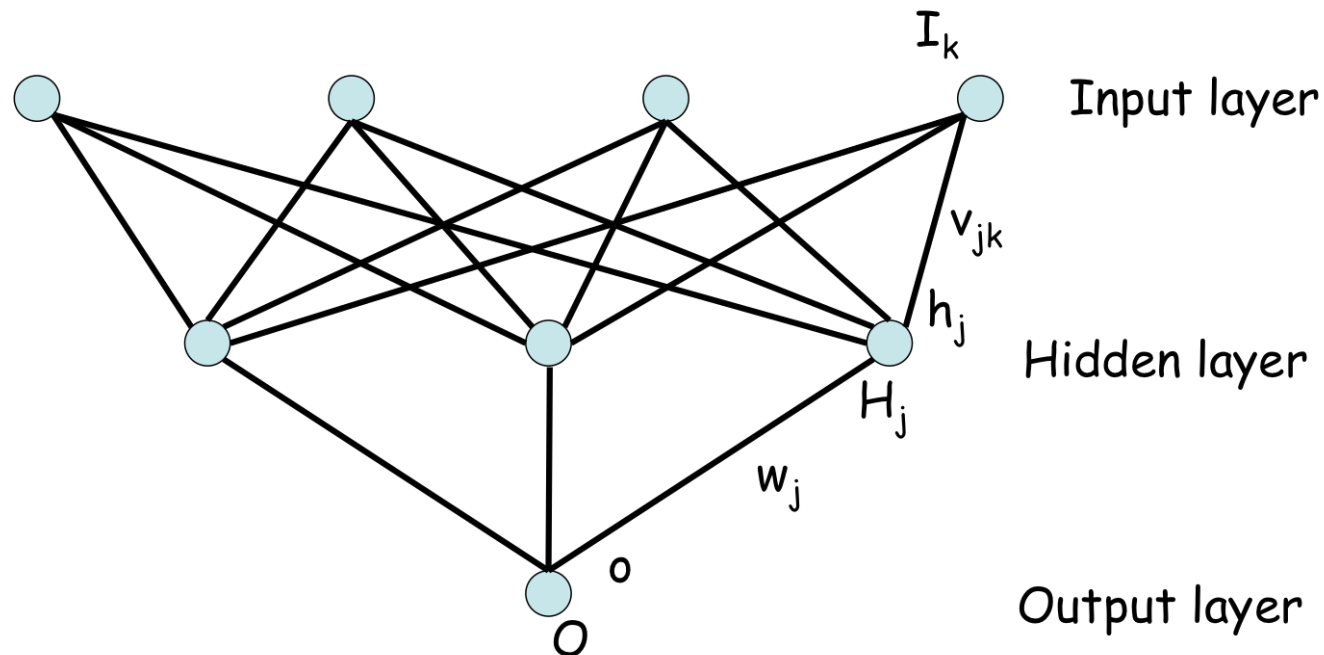
$$\begin{aligned} \frac{\partial E}{\partial v_{jk}} &= \frac{\partial E(O(o(H_j(h_j(v_{jk}))))}{\partial v_{jk}} \\ &= \frac{\partial E}{\partial O} \cdot \frac{\partial O}{\partial o} \cdot \frac{\partial o}{\partial H_j} \cdot \frac{\partial H_j}{\partial h_j} \cdot \frac{\partial h_j}{\partial v_{jk}} \\ &= (O - t) \cdot g'(o) \cdot w_j \cdot g'(h_j) \cdot I_k \end{aligned}$$



# Summary

$$\frac{\partial E}{\partial w_j} = (O - t) \cdot g'(o) \cdot H_j$$

$$\frac{\partial E}{\partial v_{jk}} = (O - t) \cdot g'(o) \cdot w_j \cdot g'(h_j) \cdot I_k$$



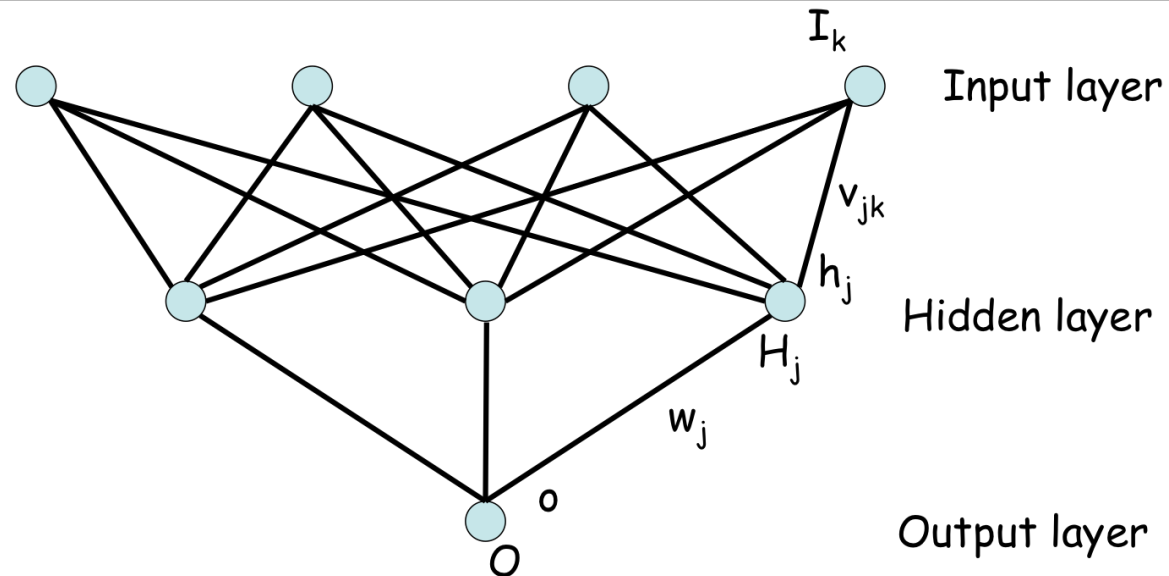


Or

$$\frac{\partial E}{\partial w_j} = (O - t) \cdot g'(o) \cdot H_j = \delta \cdot H_j$$

$$\frac{\partial E}{\partial v_{jk}} = (O - t) \cdot g'(o) \cdot w_j \cdot g'(h_j) \cdot I_k = \delta \cdot w_j \cdot g'(h_j) \cdot I_k$$

$$\delta = (O - t) \cdot g'(o)$$

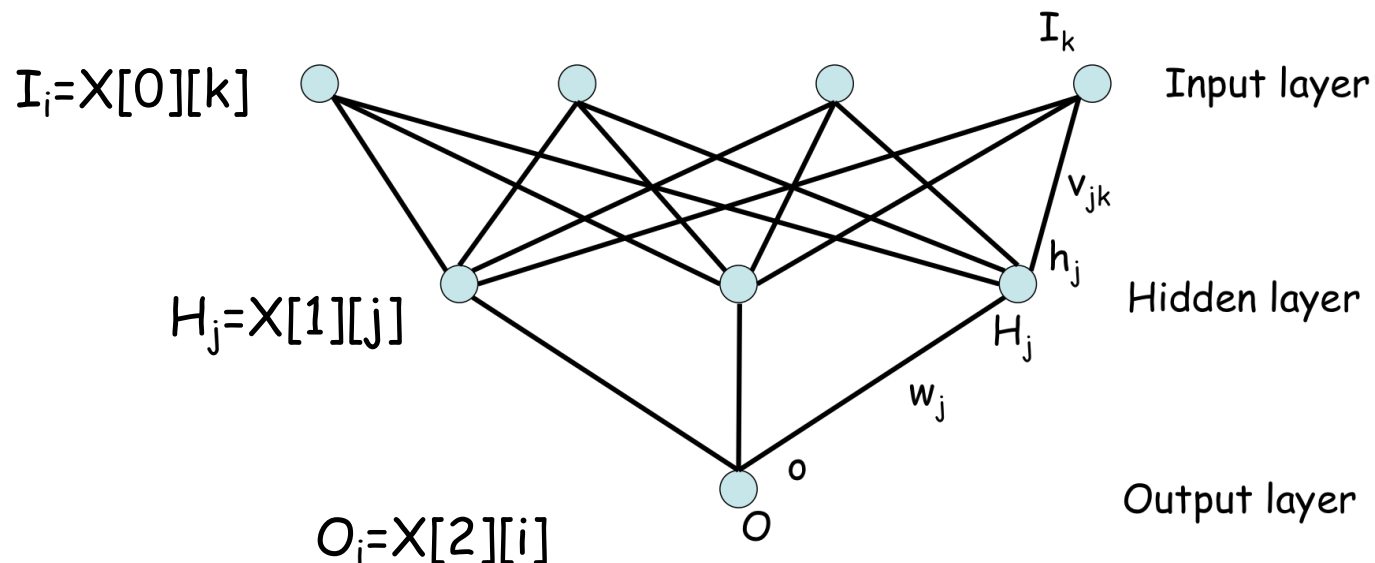


Or

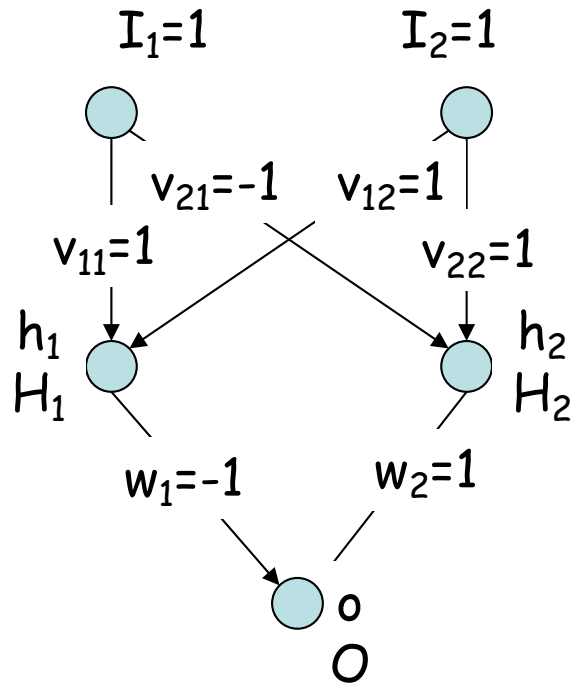
$$\frac{\partial E}{\partial w_j} = \delta \cdot H_j = \delta \cdot x[1][j]$$

$$\frac{\partial E}{\partial v_{jk}} = \delta \cdot w_j \cdot g'(h_j) \cdot I_k = \delta \cdot w_j \cdot x[1][j] \cdot (1 - x[1][j]) \cdot I_k$$

$$\delta = (O - t) \cdot g'(o) = (x[2][i] - t_i) \cdot x[2][i] \cdot (1 - x[2][i])$$



# Can you do it your self?



$$\Delta w_j = -\varepsilon \cdot \frac{\partial E}{\partial w_j}; \Delta v_{jk} = -\varepsilon \cdot \frac{\partial E}{\partial v_{jk}}$$

$$\frac{\partial E}{\partial w_j} = (O - t) \cdot g'(o) \cdot H_j$$

$$\frac{\partial E}{\partial v_{jk}} = (O - t) \cdot g'(o) \cdot w_j \cdot g'(h_j) \cdot I_k$$

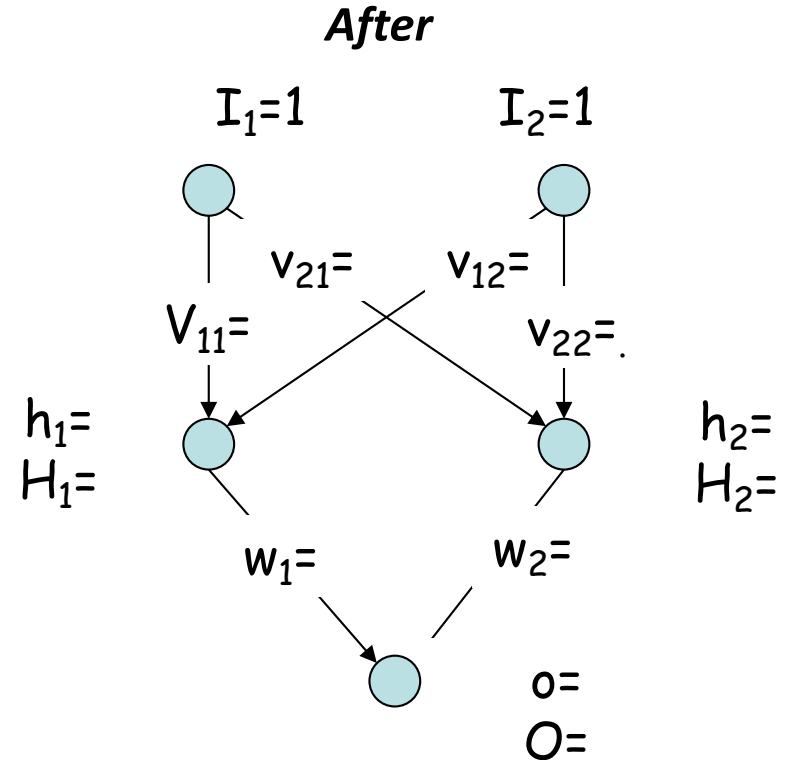
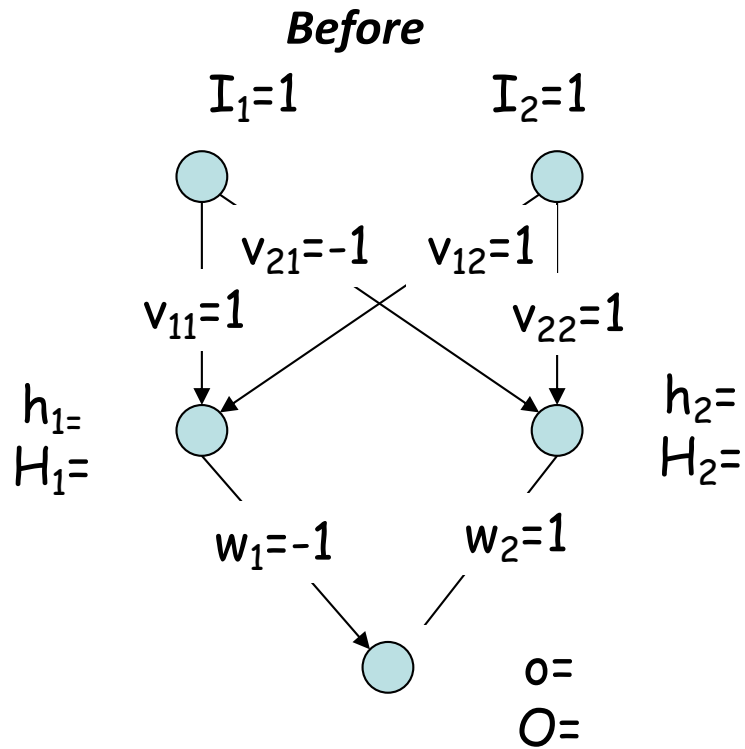
$$g'(x) = (1 - g(x)) \cdot g(x)$$

$$O = g(o)$$

What is the output ( $O$ ) from the network?  
 What are the  $\Delta w_{ij}$  and  $\Delta v_{jk}$  values if the target value is 0 and  $\varepsilon=0.5$ ?

# Can you do it your self ( $\epsilon=0.5$ ).

Has the error decreased?



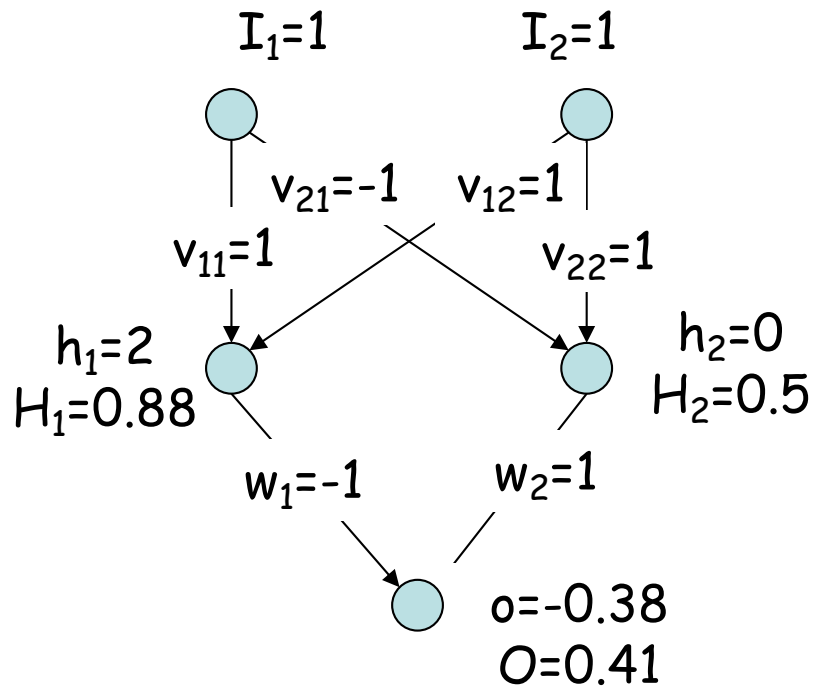
$\Delta w_1 = ??$
$\Delta w_2 = ??$

$\Delta v_{11} = ??$
$\Delta v_{12} = ??$
$\Delta v_{21} = ??$
$\Delta v_{22} = ??$

# Can you do it your self ( $\epsilon=0.5$ ).

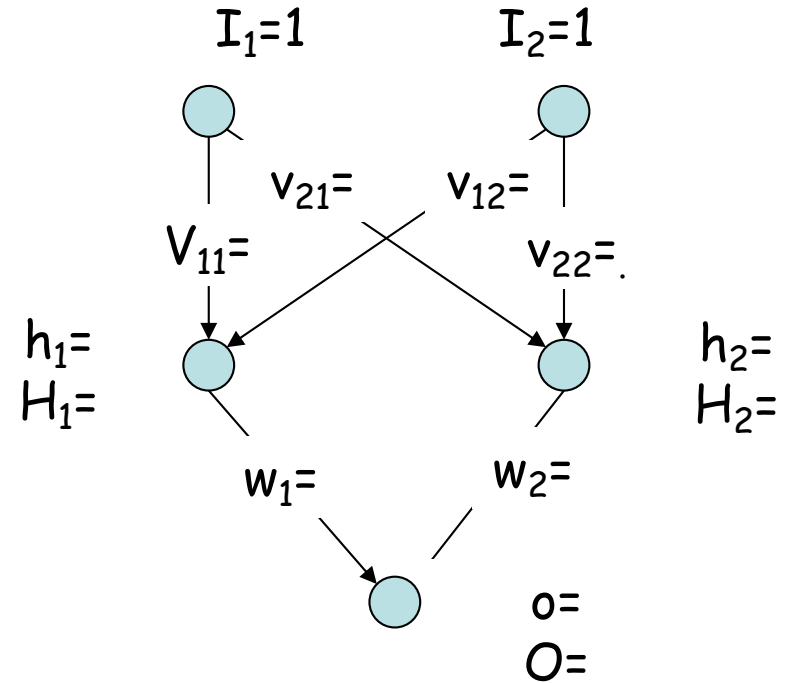
Has the error decreased?

Before



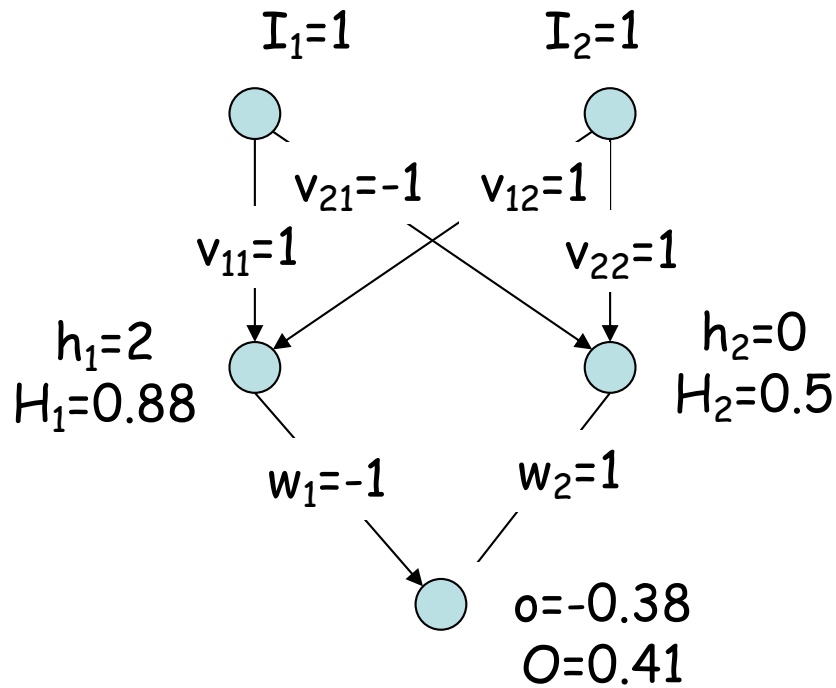
$\Delta w_1 = ??$   
 $\Delta w_2 = ??$

After



$\Delta v_{11} = ??$   
 $\Delta v_{12} = ??$   
 $\Delta v_{21} = ??$   
 $\Delta v_{22} = ??$

# Can you do it your self?



- What is the output ( $O$ ) from the network?
- What are the  $\Delta w_{ij}$  and  $\Delta v_{jk}$  values if the target value is 0?

$$\Delta w_j = -\varepsilon \cdot \frac{\partial E}{\partial w_j}; \Delta v_{jk} = -\varepsilon \cdot \frac{\partial E}{\partial v_{jk}}$$

$$\frac{\partial E}{\partial w_j} = (O - t) \cdot g'(o) \cdot H_j = \delta \cdot H_j$$

$$\delta = (O - t) \cdot g'(o) = 0.41 \cdot 0.41 \cdot (1 - 0.41) = 0.099$$

$$\Delta w_1 = -\varepsilon \cdot \delta \cdot 0.88 = -\varepsilon \cdot 0.087$$

$$\Delta w_2 = -\varepsilon \cdot \delta \cdot 0.5 = -\varepsilon \cdot 0.050$$

$$\frac{\partial E}{\partial v_{jk}} = g'(h_j) \cdot I_k \cdot (O - t) \cdot g'(o) \cdot w_j$$

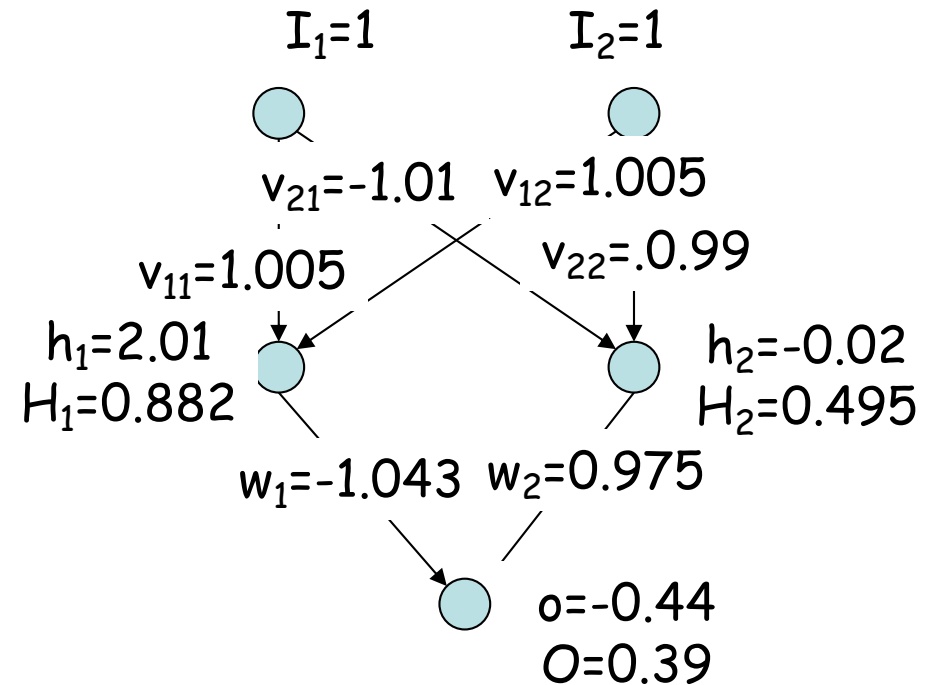
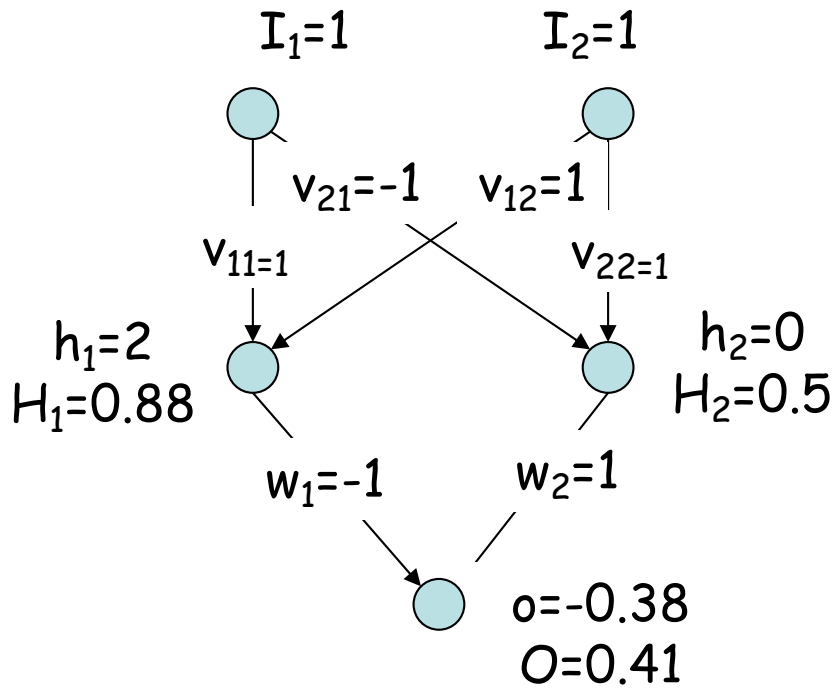
$$\Delta v_{11} = -\varepsilon \cdot H_1 \cdot (1 - H_1) \cdot 1 \cdot \delta \cdot (-1) = \varepsilon \cdot 0.01$$

$$\Delta v_{12} = \Delta v_{11}$$

$$\Delta v_{21} = -\varepsilon \cdot H_2 \cdot (1 - H_2) \cdot 1 \cdot \delta \cdot 1 = -\varepsilon \cdot 0.02$$

$$\Delta v_{22} = \Delta v_{21}$$

Can you do it your self ( $\varepsilon=0.5$ ).  
Has the error decreased?



$$\Delta w_1 = -\varepsilon \cdot \delta \cdot 0.88 = -\varepsilon \cdot 0.087$$

$$\Delta w_2 = -\varepsilon \cdot \delta \cdot 0.5 = -\varepsilon \cdot 0.050$$

$$\Delta v_{11} = -\varepsilon \cdot H_1 \cdot (1 - H_1) \cdot 1 \cdot \delta \cdot (-1) = \varepsilon \cdot 0.01$$

$$\Delta v_{12} = \Delta v_{11}$$

$$\Delta v_{21} = -\varepsilon \cdot H_2 \cdot (1 - H_2) \cdot 1 \cdot \delta \cdot 1 = -\varepsilon \cdot 0.02$$

$$\Delta v_{22} = \Delta v_{21}$$

# Sequence encoding

---

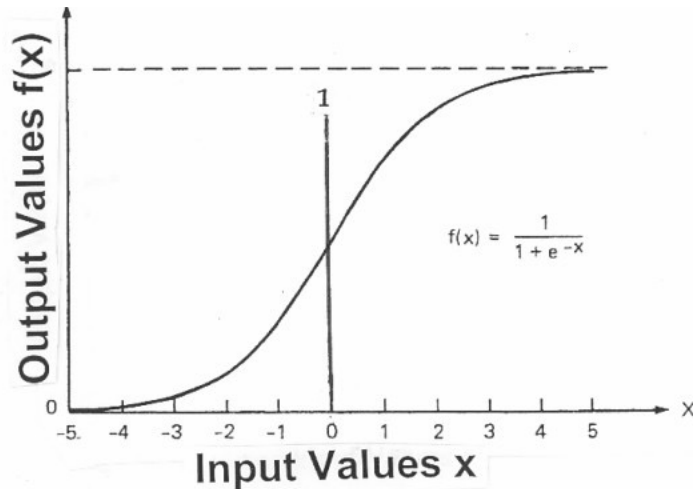
- Change in weight is linearly dependent on input value
- “True” sparse encoding (i.e 1/0) is therefore highly inefficient
- Sparse is most often encoded as
  - +1/-1 or 0.9/0.05

$$\frac{\partial E}{\partial v_{jk}} = \delta \cdot w_j \cdot g'(h_j) \cdot I_k = \delta \cdot w_j \cdot x[1][j] \cdot (1 - x[1][j]) \cdot I_k$$



# Sequence encoding - rescaling

- Rescaling the input values



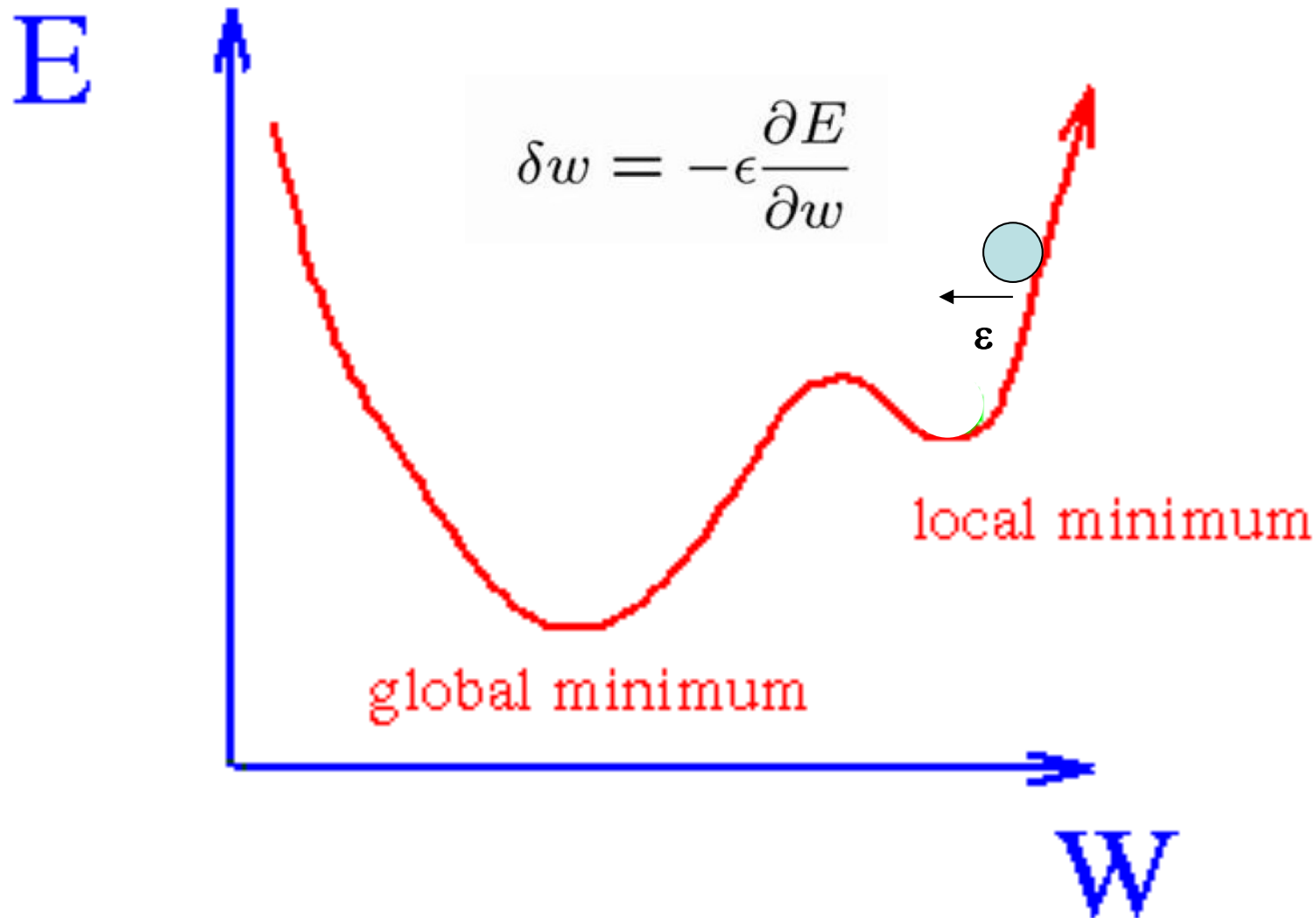
If the input ( $o$  or  $h$ ) is too large or too small,  $g'$  is zero and the weights are not changed. Optimal performance is when  $o, h$  are close to 0.5

$$\frac{\partial E}{\partial w_j} = (O - t) \cdot g'(o) \cdot H_j = \delta \cdot H_j$$

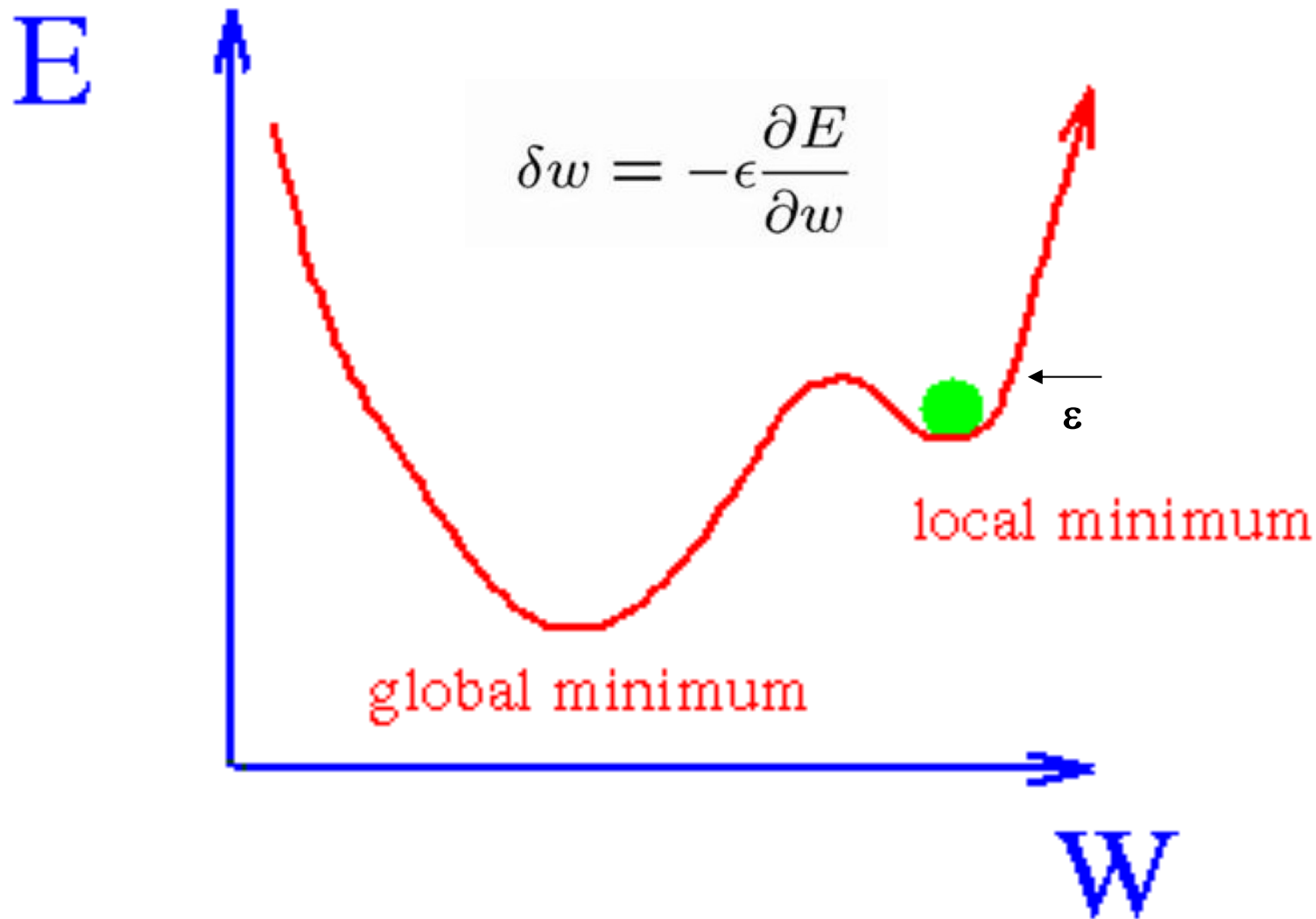
$$\frac{\partial E}{\partial v_{jk}} = g'(h_j) \cdot I_k \cdot (O - t) \cdot g'(o) \cdot w_j = g'(h_j) \cdot I_k \cdot \delta \cdot w_j$$

$$\delta = (O - t) \cdot g'(o)$$

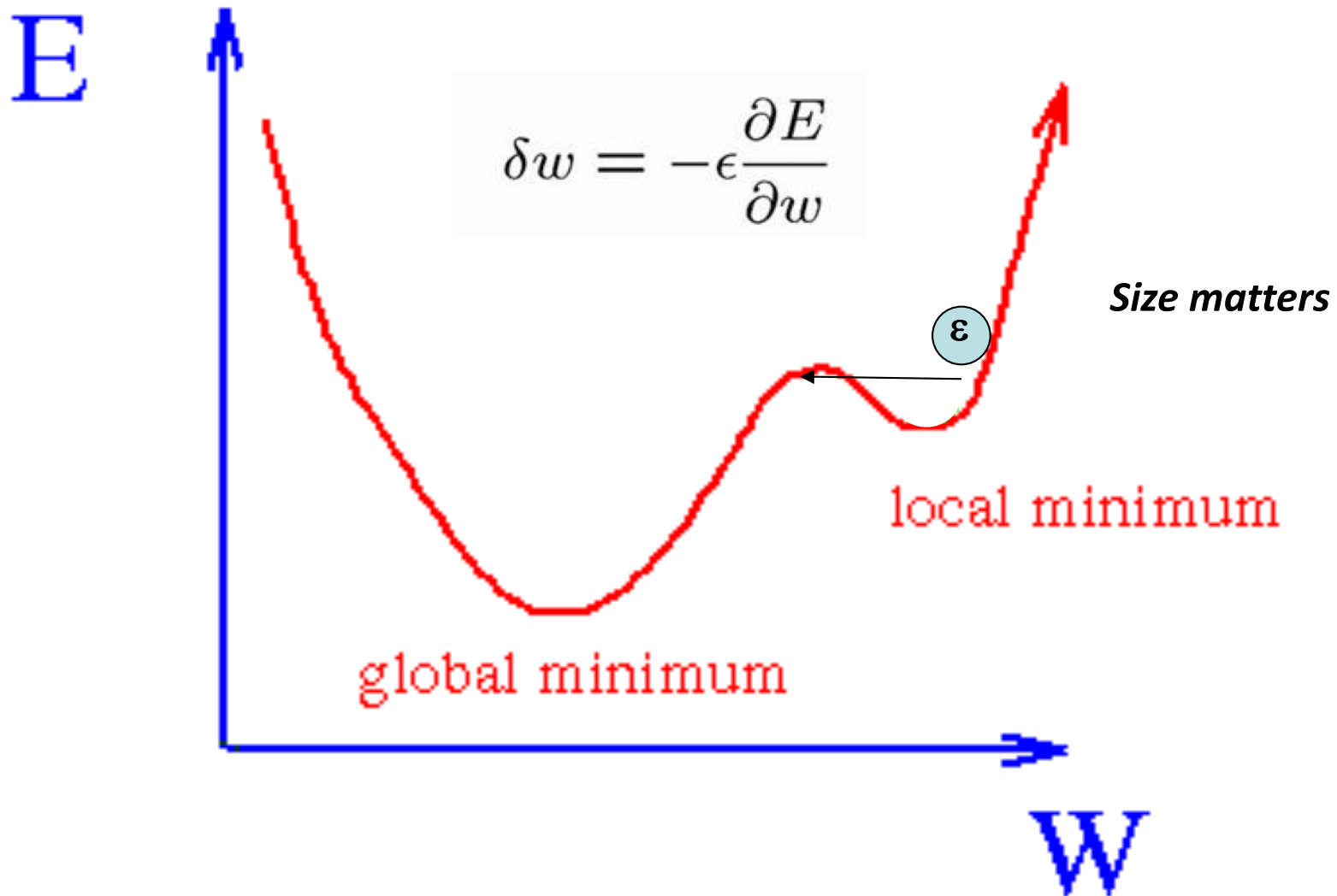
# Training and error reduction



# Training and error reduction



# Training and error reduction



# Demo

---

- <http://playground.tensorflow.org/>
-

# Do hidden neurons matter?

- The environment matters

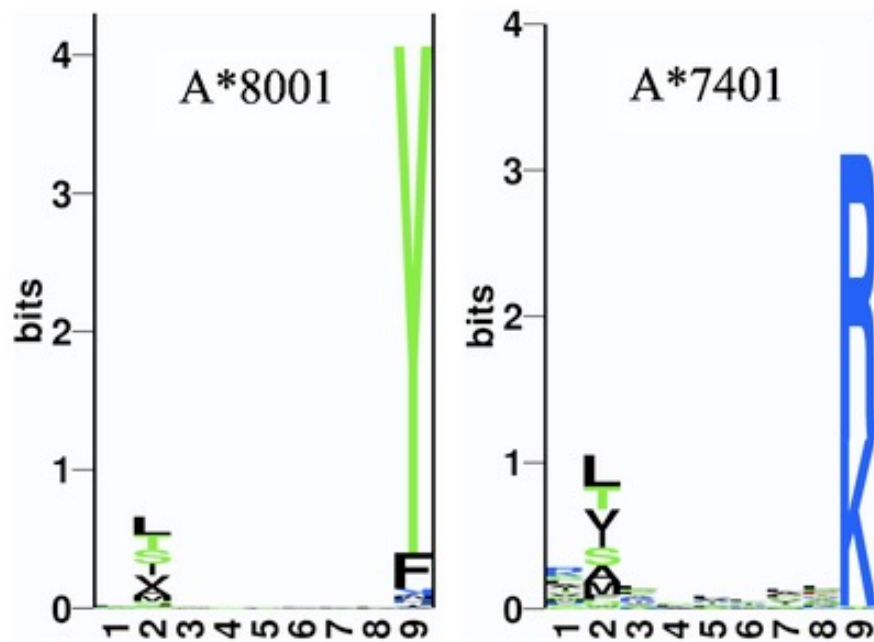
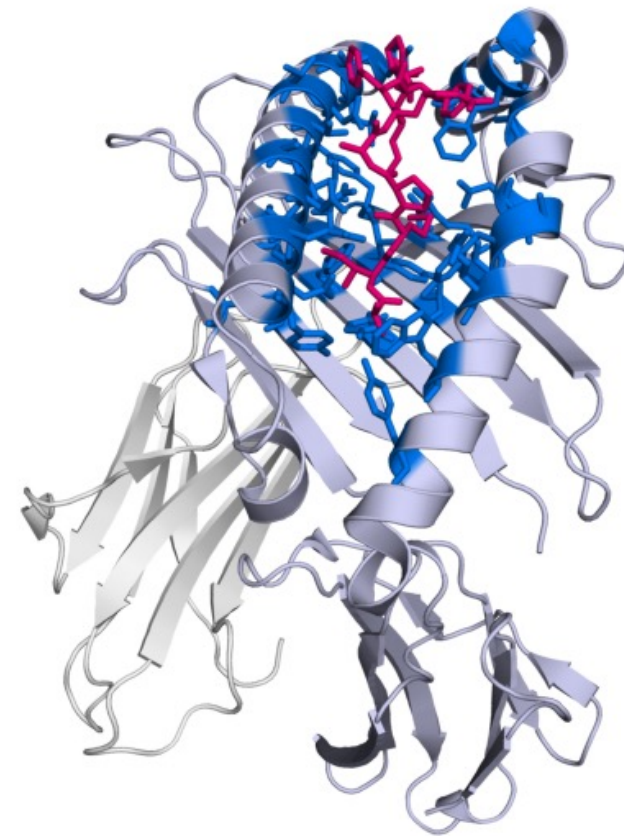


Figure 1. Prospective validation using hitherto uncharacterized HLA molecules.

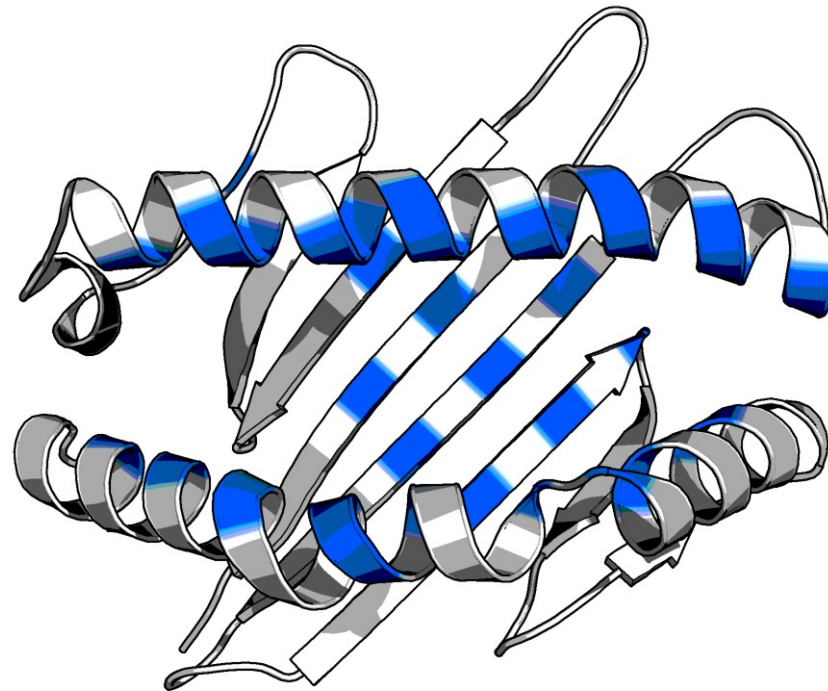


NetMHCpan

# Context matters

---

FMIDWILDA	YFAMYGE <b>KVAHTHVD</b> TLY <b>VR</b> YH <b>Y</b> YTWAV <b>L</b> A <b>Y</b> T <b>W</b> Y	0.89	A0201
FMIDWILDA	YFAMYQ <b>EN</b> MAHT <b>D</b> ANTLY <b>II</b> YR <b>D</b> YTWV <b>AR</b> V <b>Y</b> RGY	0.08	A0101
DSDGSFFLY	YFAMYGE <b>KVAHTHVD</b> TLY <b>VR</b> YH <b>Y</b> YTWAV <b>L</b> A <b>Y</b> T <b>W</b> Y	0.08	A0201
DSDGSFFLY	YFAMYQ <b>EN</b> MAHT <b>D</b> ANTLY <b>II</b> YR <b>D</b> YTWV <b>AR</b> V <b>Y</b> RGY	0.85	A0101



# Summary

---

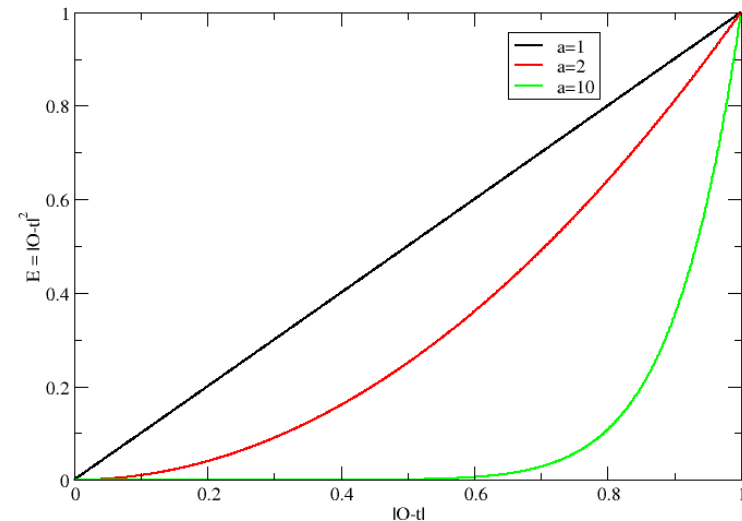
- Gradient descent is used to determine the updates for the synapses in the neural network
  - Some relatively simple math defines the gradients
    - Networks without hidden layers can be solved on the back of an envelope (SMM exercise)
    - Hidden layers are a bit more complex, but still ok
  - Always train networks using a test set to stop training
    - Be careful when reporting predictive performance
      - Use “nested” cross-validation for small data sets
  - And hidden neurons do matter (sometimes)
-



# And some more stuff for the long cold and rainy summer nights

- Can it maybe be made differently?

$$E = \frac{1}{\alpha} \cdot (O - t)^\alpha$$

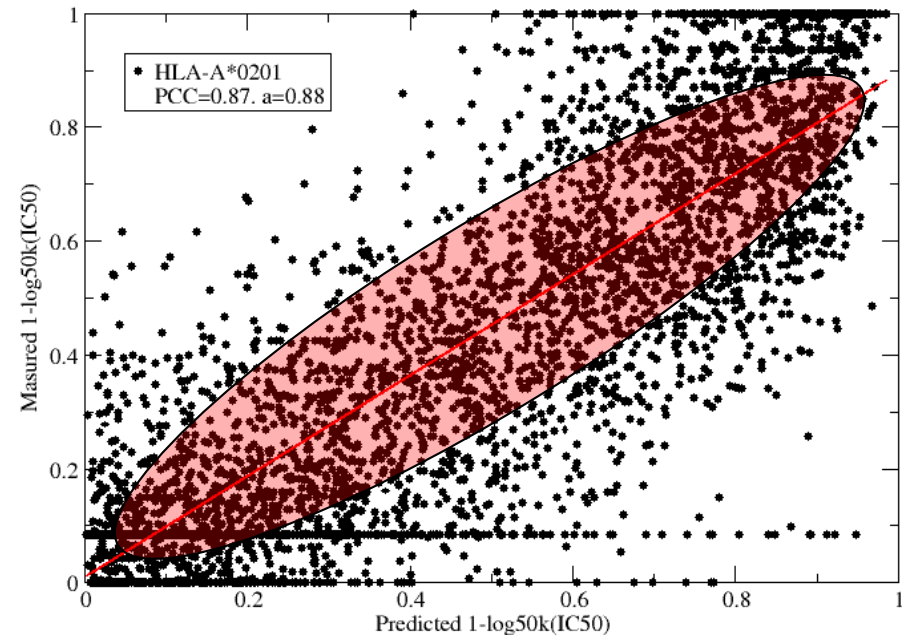
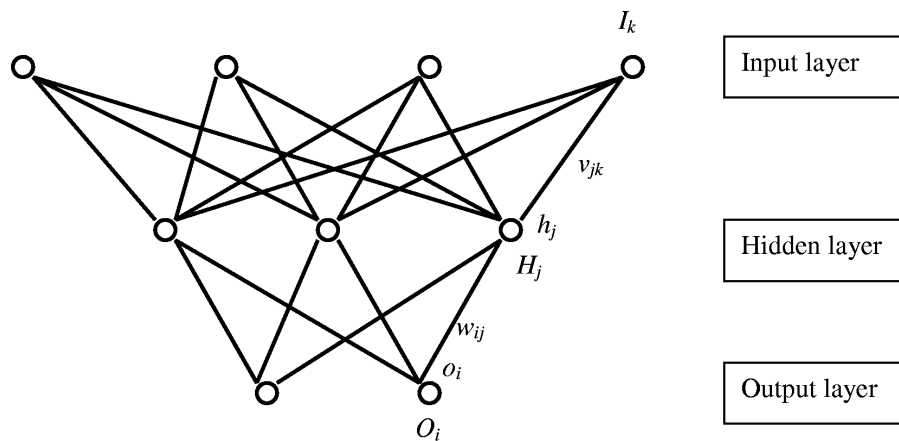


# Predicting accuracy

- Can it be made differently?

$$E = \frac{1}{2} \cdot (O_1 - t)^2 \cdot O_2 + \lambda \cdot (1 - O_2)$$

Reliability



# Making sense of ANN weights

---

- Identification of position specific receptor ligand interactions by use of artificial neural network decomposition.  
An investigation of interactions in the MHC:peptide system

Master thesis' by Frederik Otzen Bagger  
and Piotr Chmura

---

# Making sense of ANN weights

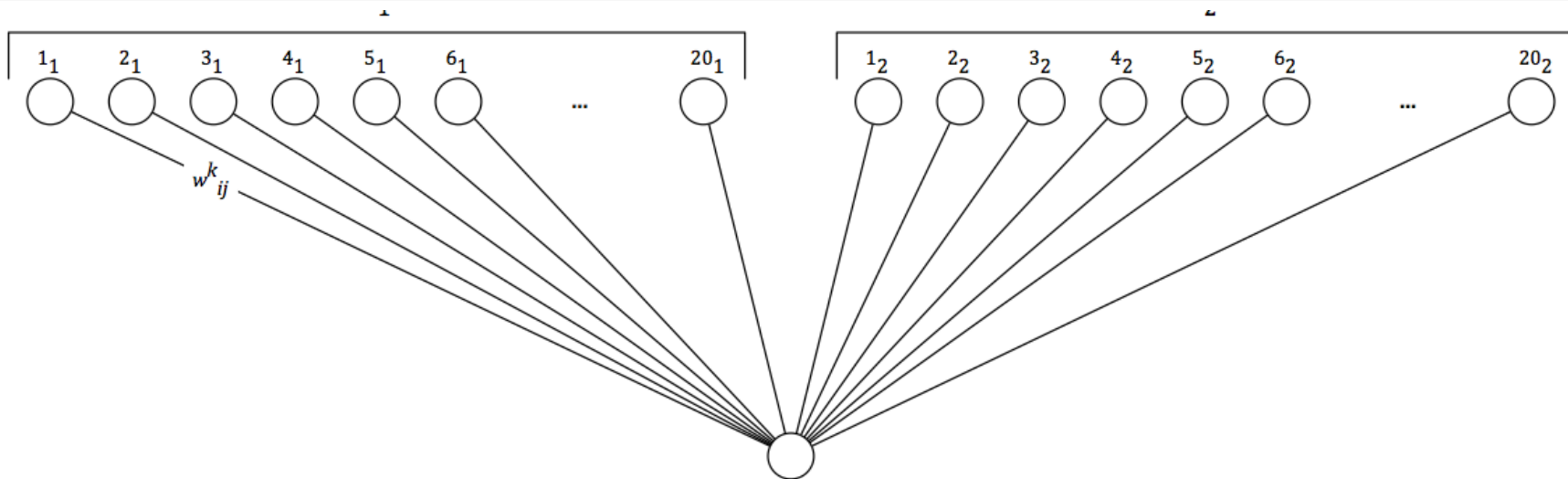
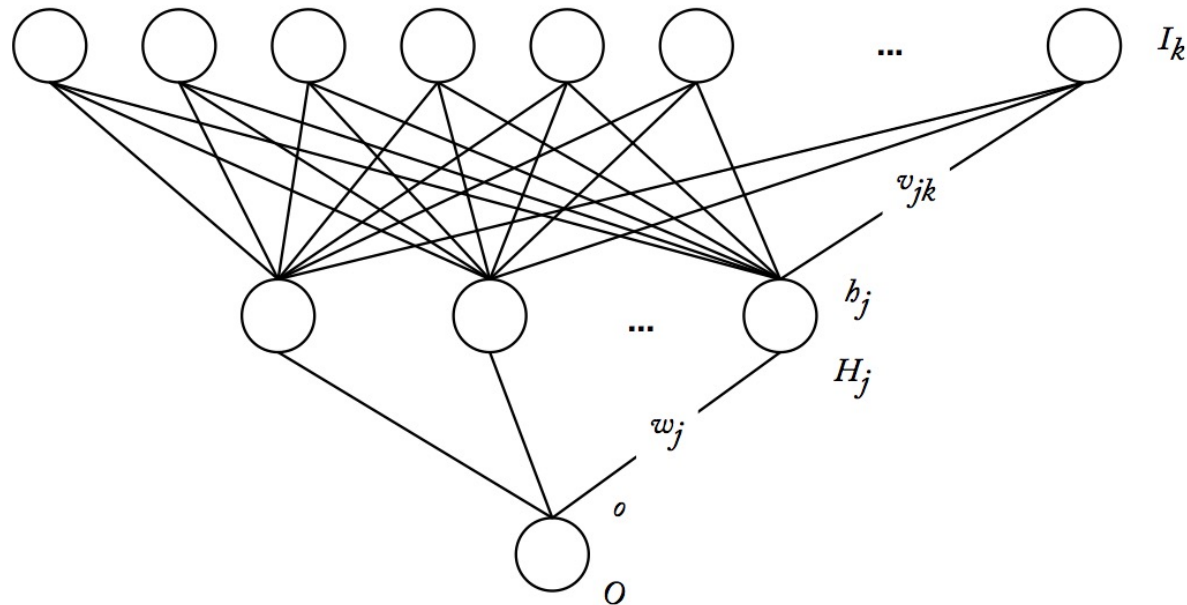


Figure 2.1. Two layer ANN with two amino acids as input,  $p = \{1,2\}$ , and one output neuron. The direction is downwards, and the graph is directed.



# Making sense of ANN weights

A0201

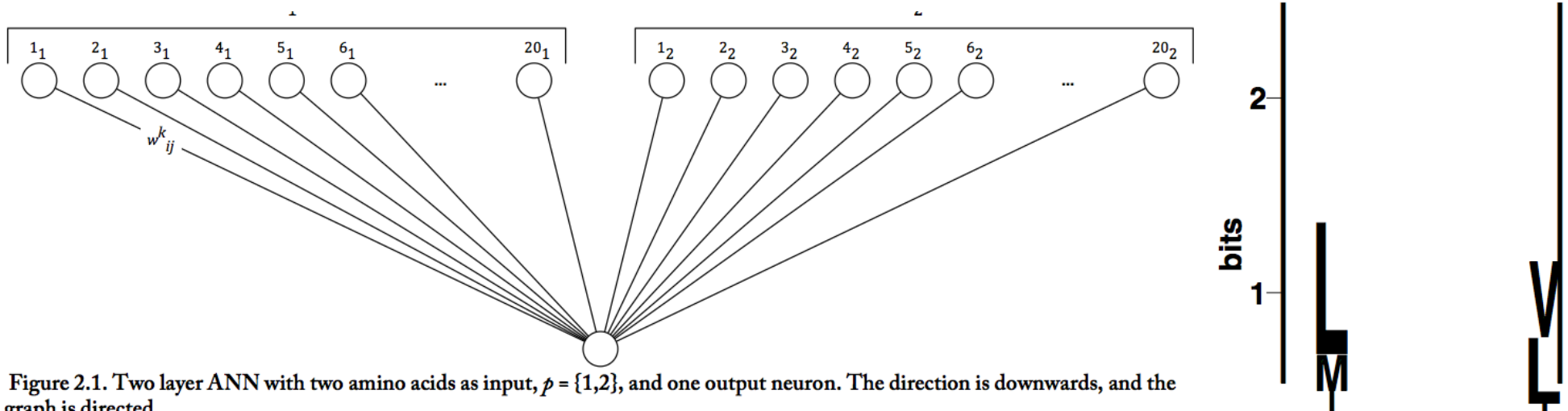
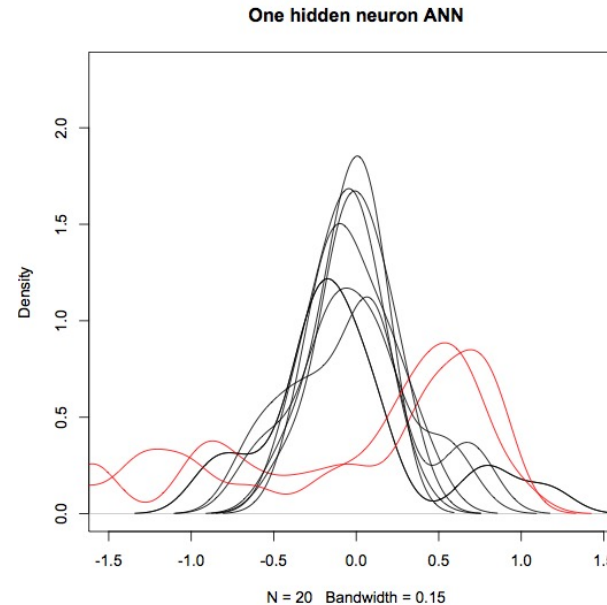
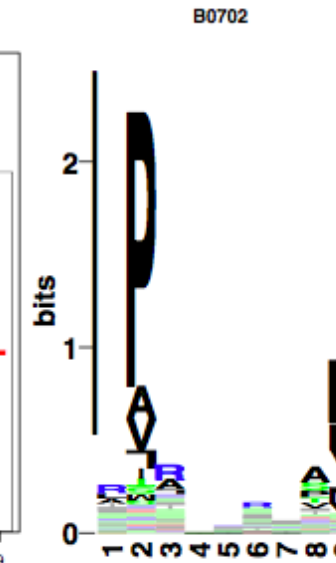
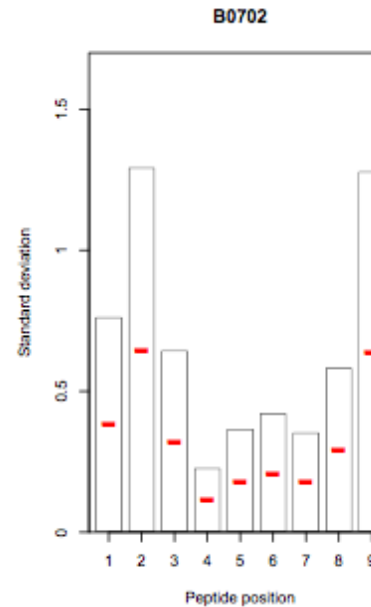
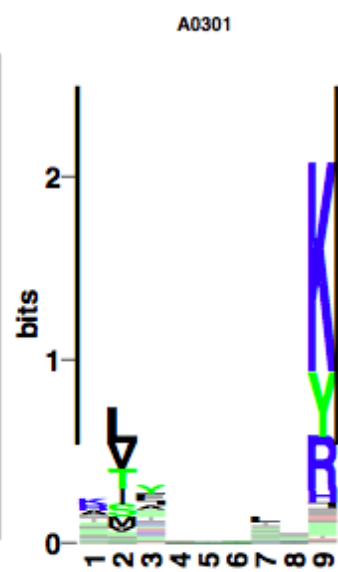
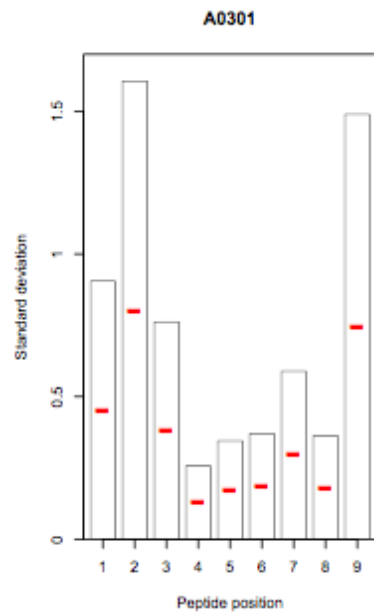
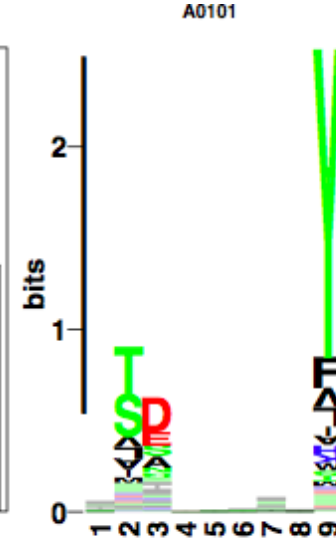
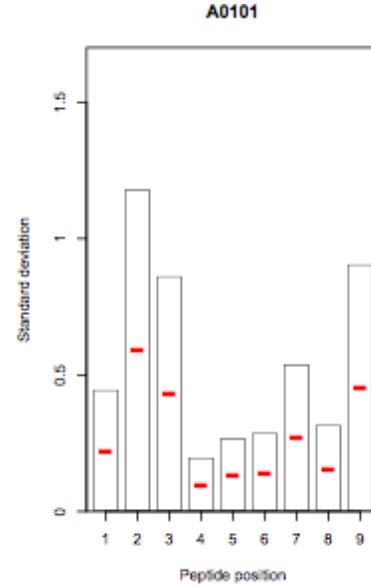
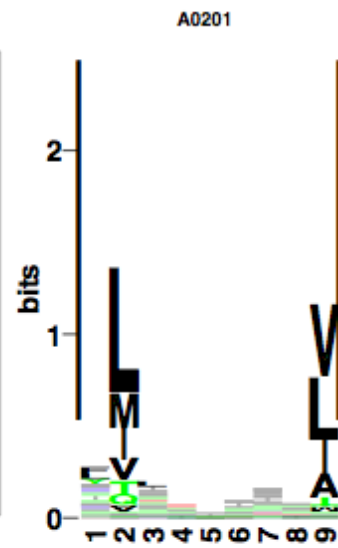
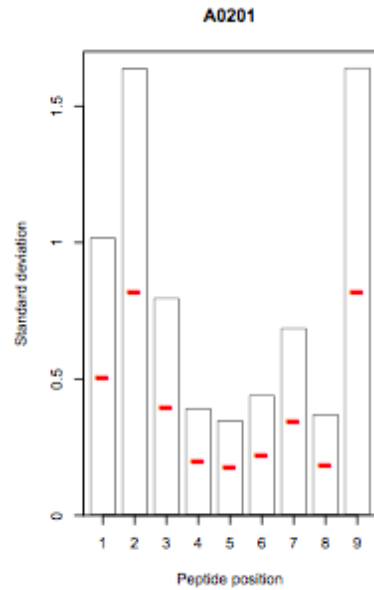


Figure 2.1. Two layer ANN with two amino acids as input,  $p = \{1,2\}$ , and one output neuron. The direction is downwards, and the graph is directed.

$$\delta_j^p = \sqrt{\frac{1}{N-1} \sum_A (v_{Aj}^p - \mu_j^p)^2}, \text{ where } \mu_j^p = \frac{1}{N} \sum_A v_{Aj}^p$$

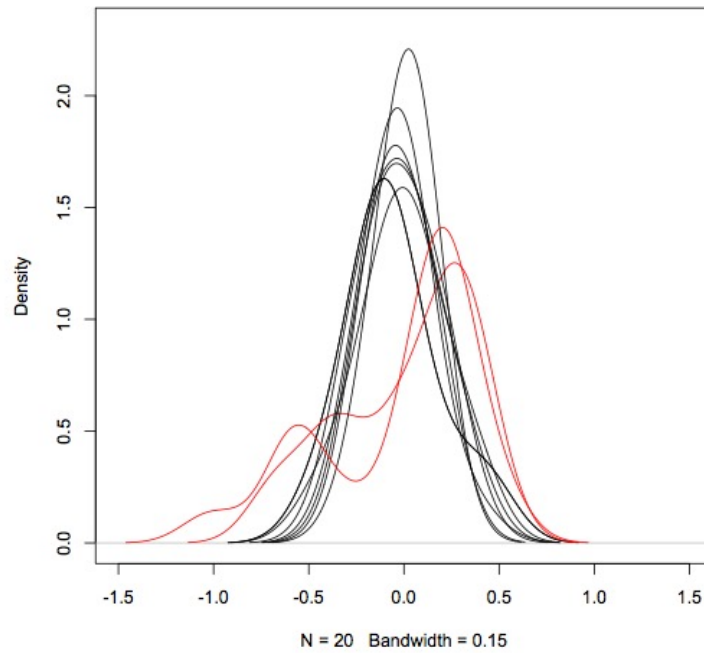


# Making sense of ANN weights

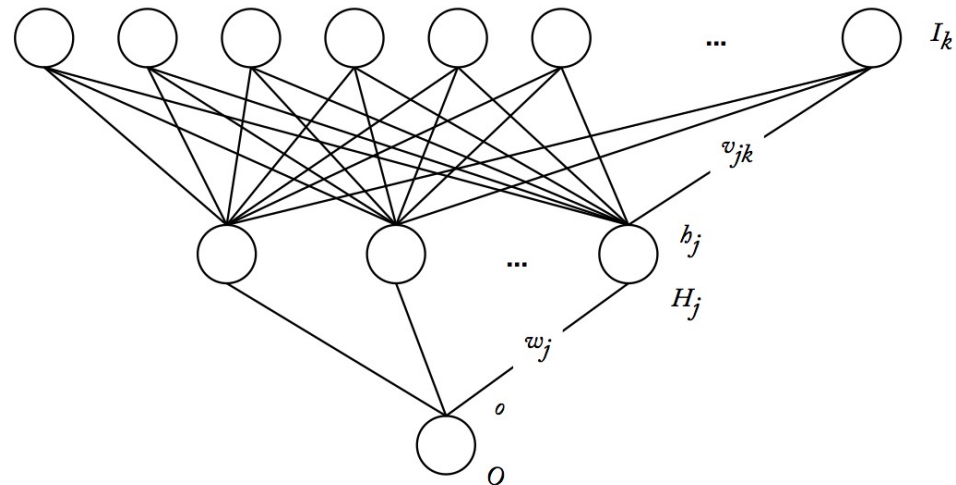
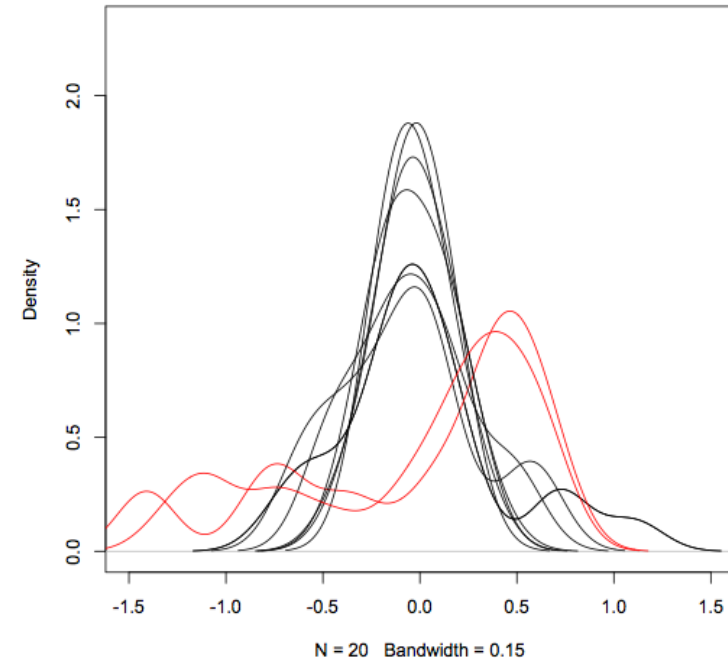


# Making sense of ANN weights

Two hidden neuron ANN, 1th hidden neuron

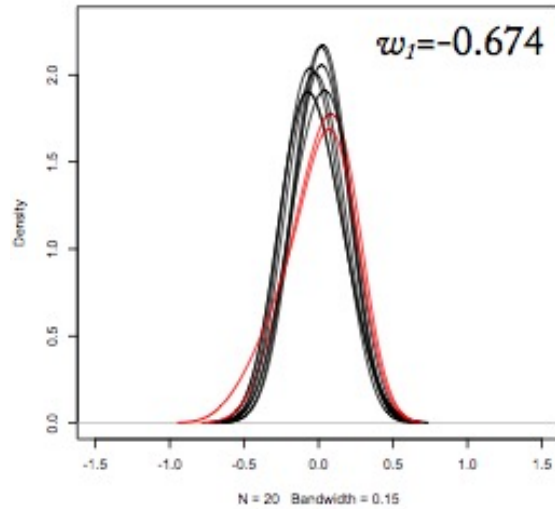


Two hidden neuron ANN, 2nd hidden neuron

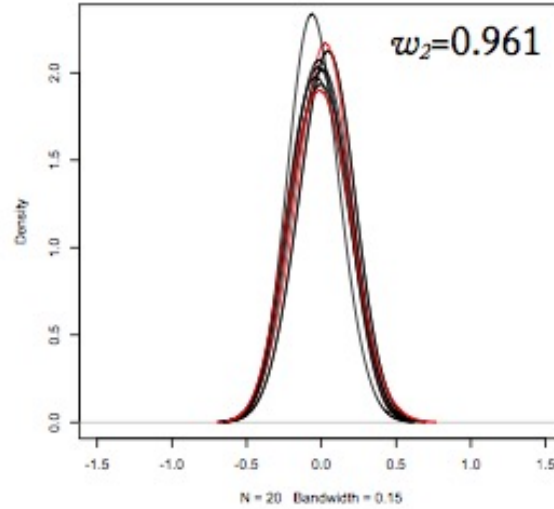


# Making sense of ANN weights

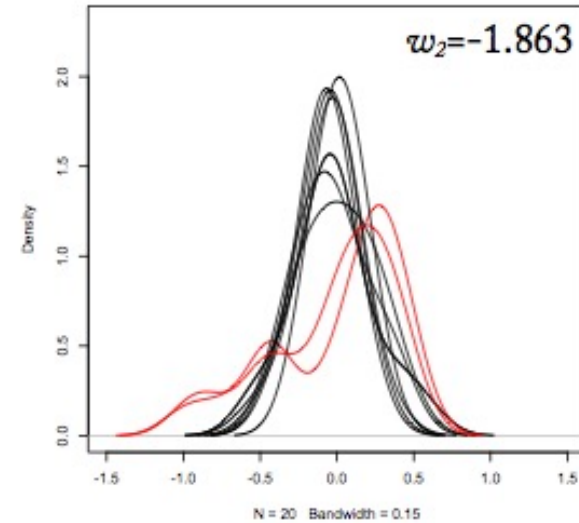
Five hidden neuron ANN, 1st hidden neuron



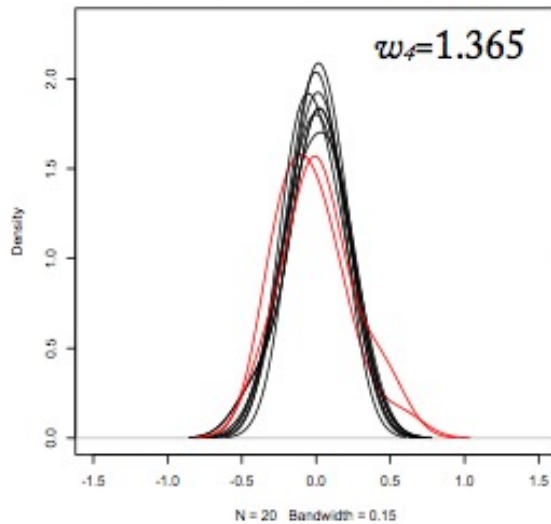
Five hidden neuron ANN, 2nd hidden neuron



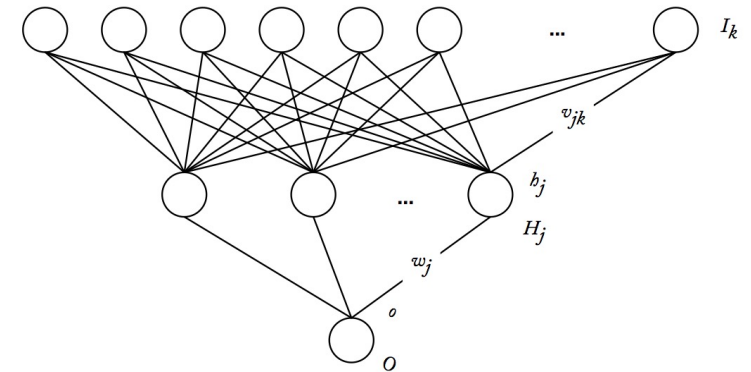
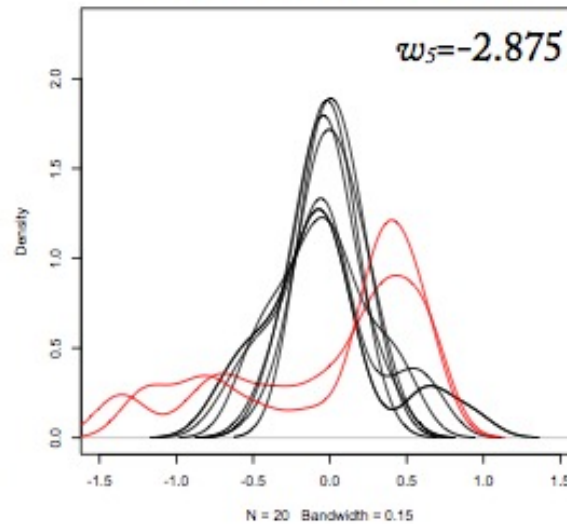
Five hidden neuron ANN, 3rd hidden neuron



Five hidden neuron ANN, 4th hidden neuron

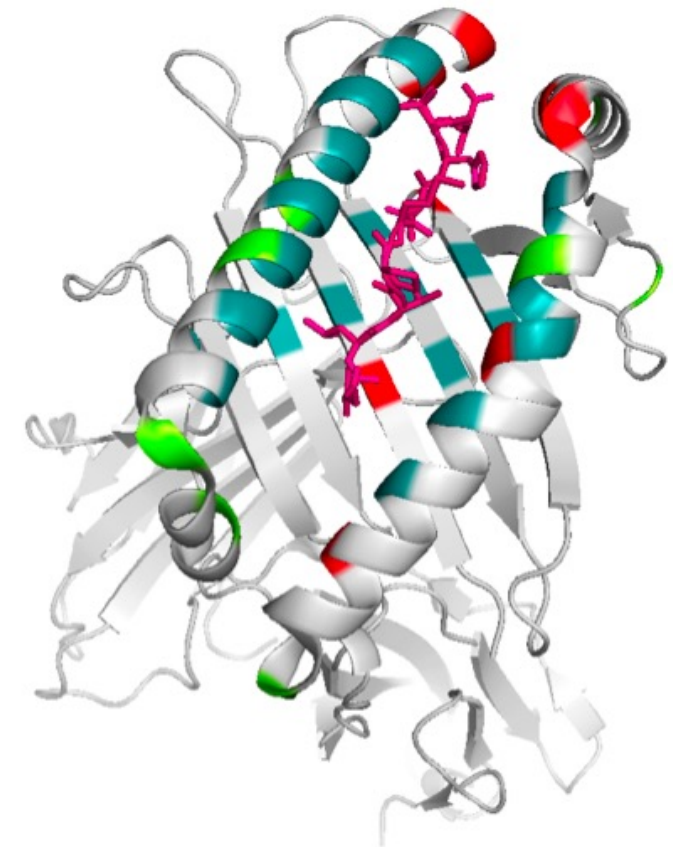
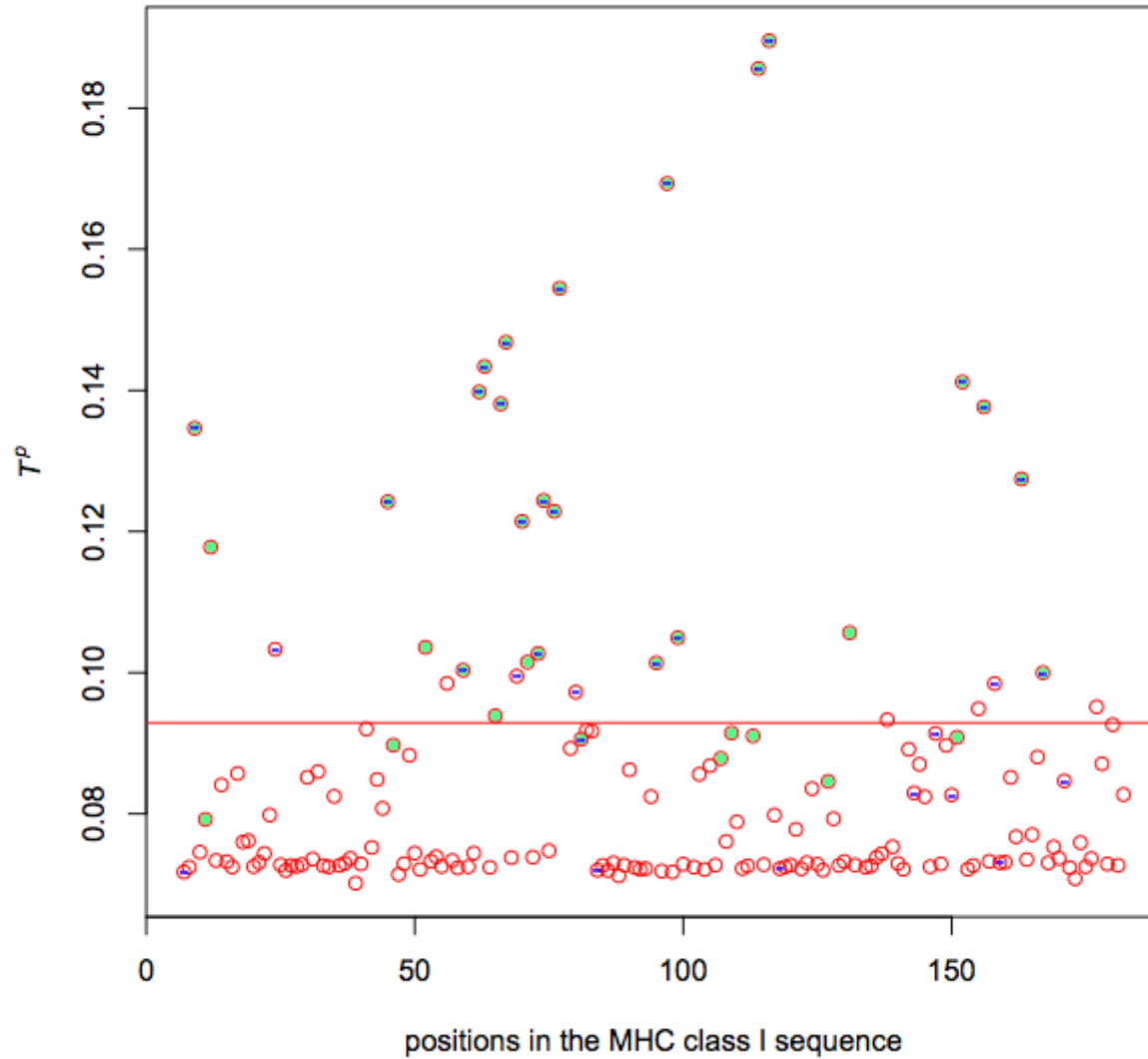


Five hidden neuron ANN, 5th hidden neuron





# Making sense of ANN weights



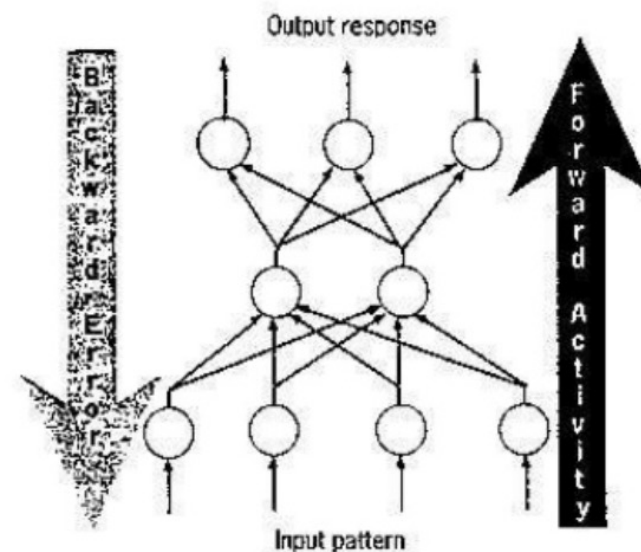
## Back Propagation

### Advantages

- Multi layer Perceptron network can be trained by the back propagation algorithm to perform any mapping between the input and the output.

### What is wrong with back-propagation?

- It requires labeled training data.  
Almost all data is unlabeled.
- The learning time does not scale well  
It is very slow in networks with multiple hidden layers.
- It can get stuck in poor local optima.



A backpropagation network trains with a two-step procedure. The activity from the input pattern flows forward through the network, and the error signal flows backward to adjust the weights.

# Deep(er) Network architecture

$$E = \frac{1}{2} \cdot (O - t)^2$$

$$O = g(o), H = g(h)$$

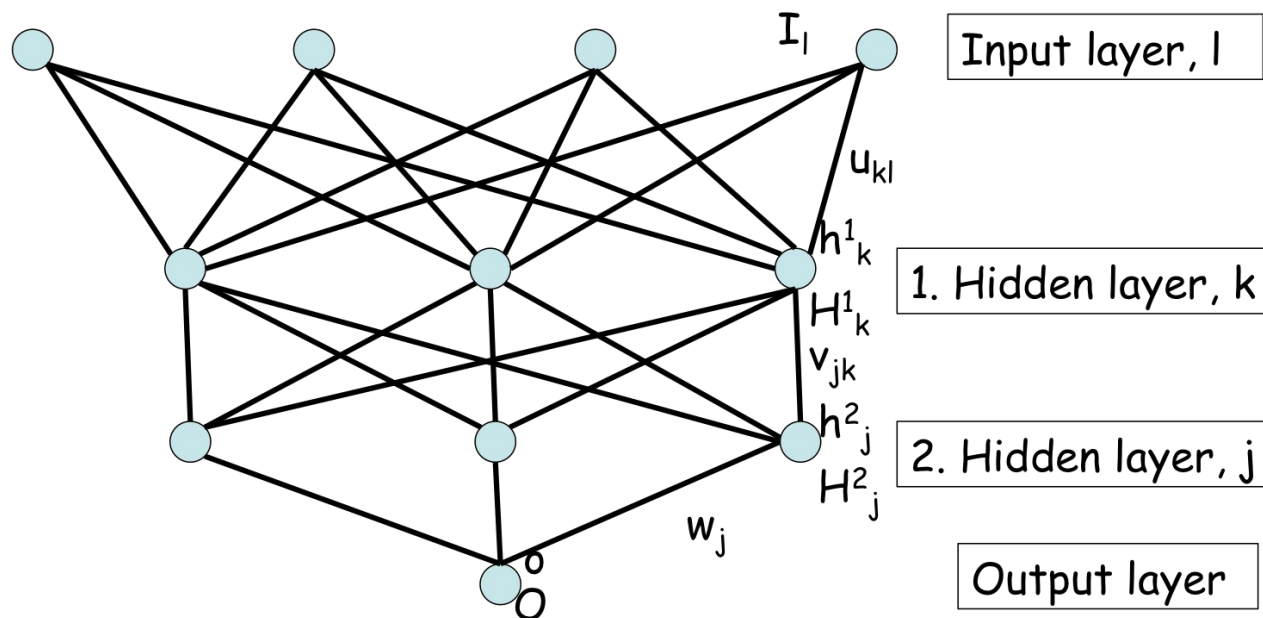
$$g(x) = \frac{1}{1 + e^{-x}}$$

$$o = \sum_j w_j \cdot H_j^2$$

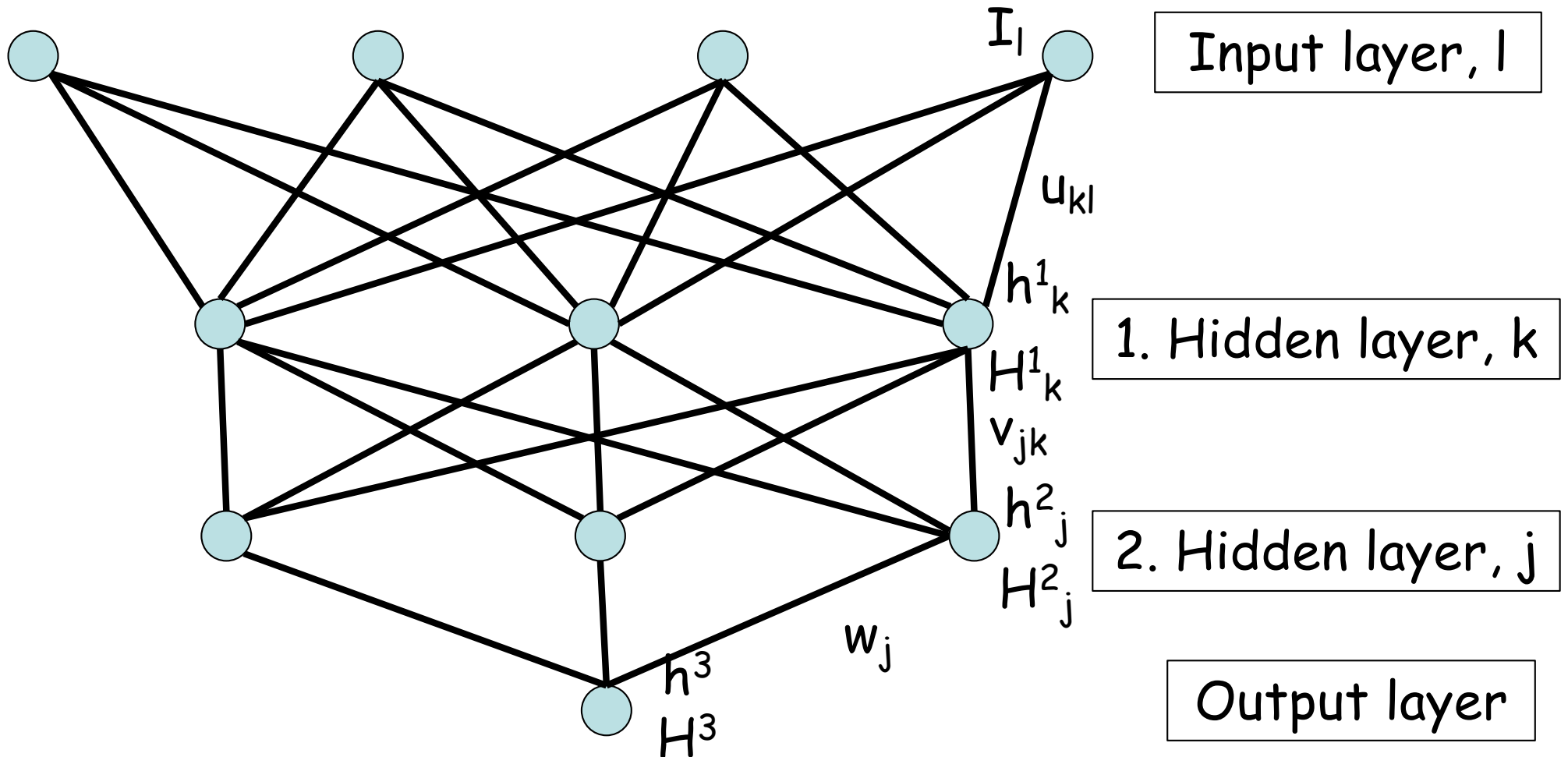
$$h_j^2 = \sum_k v_{jk} \cdot H_k^1$$

$$h_k^1 = \sum_l u_{kl} \cdot I_l$$

$$\Delta w_i = -\varepsilon \cdot \frac{\partial E}{\partial w_i}$$

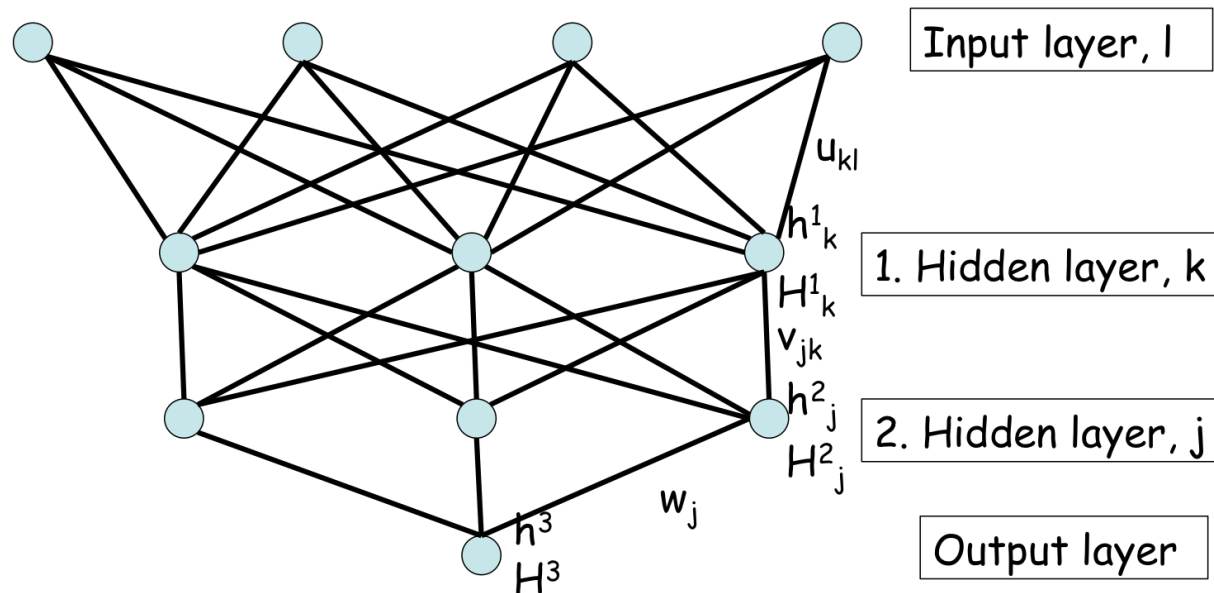


# Deeper Network architecture



$$\frac{\partial E}{\partial w_j} = \frac{\partial E(H^3(h^3(w_j)))}{\partial w_j} = \frac{\partial E}{\partial H^3} \cdot \frac{\partial H^3}{\partial h^3} \cdot \frac{\partial h^3}{\partial w_j} = (H^3 - t) \cdot g'(h^3) \cdot H_j^2$$

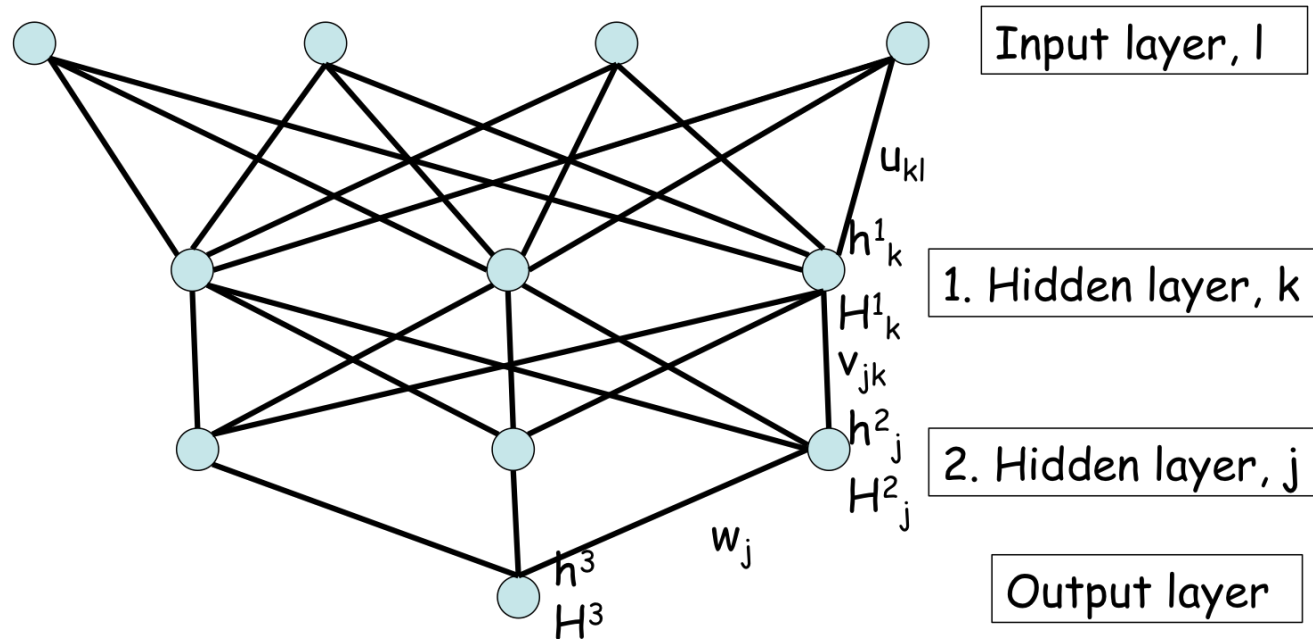
# Network architecture (hidden to hidden)



$$\frac{\partial E}{\partial v_{jk}} = \frac{\partial E}{\partial H^3} \cdot \frac{\partial H^3}{\partial h^3} \cdot \frac{\partial h^3}{\partial H^2_j} \cdot \frac{\partial H^2_j}{\partial h^2_j} \cdot \frac{\partial h^2_j}{\partial v_{jk}}$$

$$= (H^3 - t) \cdot g'(h^3) \cdot w_j \cdot g'(h^2_j) \cdot H^1_k$$

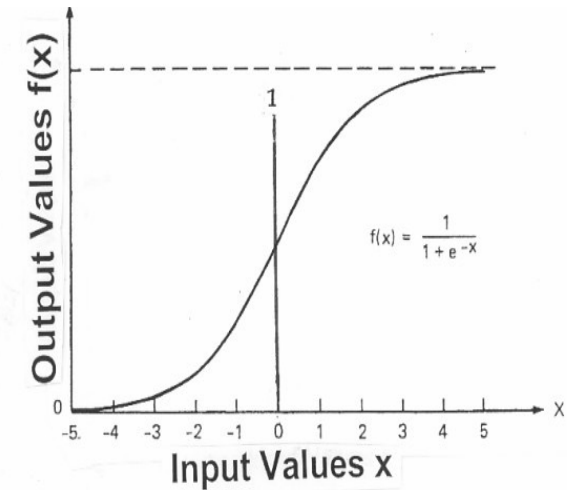
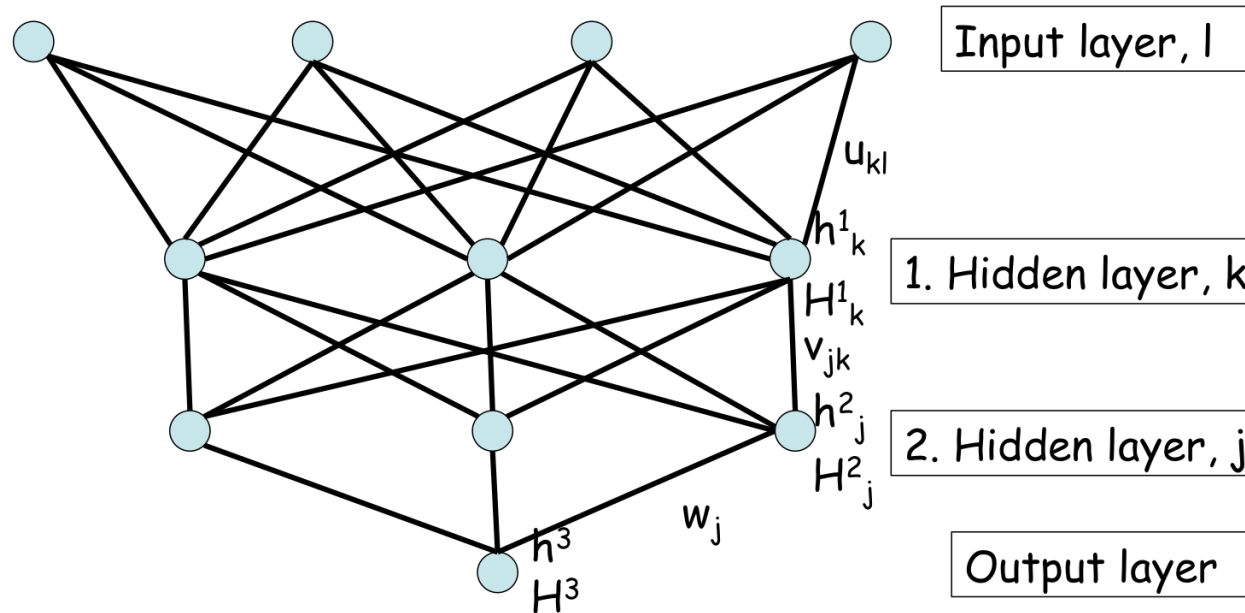
# Network architecture (input to hidden)



$$\frac{\partial E}{\partial u_{kl}} = \frac{\partial E}{\partial H^3} \cdot \frac{\partial H^3}{\partial h^3} \cdot \sum_j \frac{\partial h^3}{\partial H^2_j} \cdot \frac{\partial H^2_j}{\partial h^2_j} \cdot \frac{\partial h^2_j}{\partial H^1_k} \cdot \frac{\partial H^1_k}{\partial h^1_k} \cdot \frac{\partial h^1_k}{\partial u_{kl}}$$

$$= (H^3 - t) \cdot g'(h^3) \cdot \sum_j w_j \cdot g'(h^2_j) \cdot v_{jk} \cdot g'(h^1_k) \cdot I_l$$

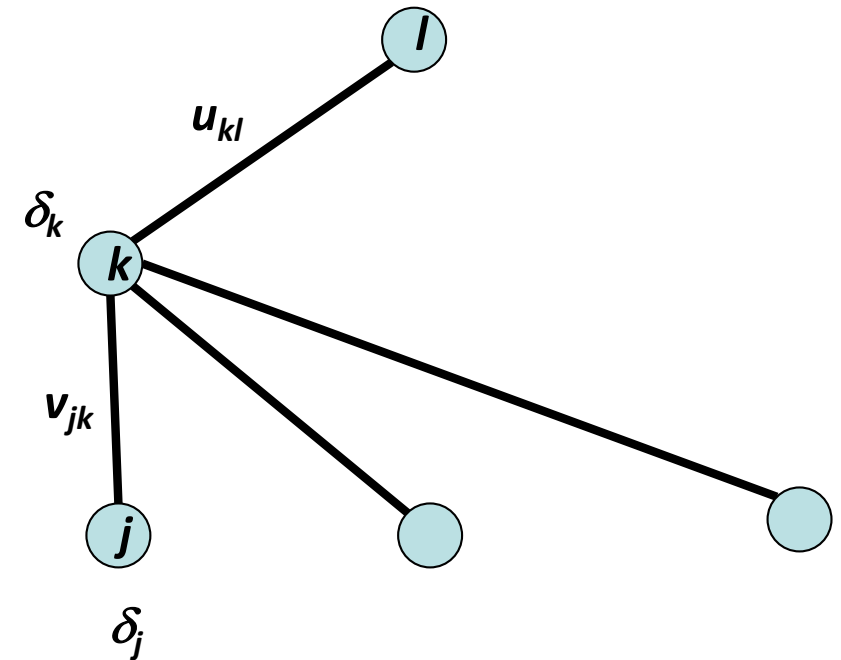
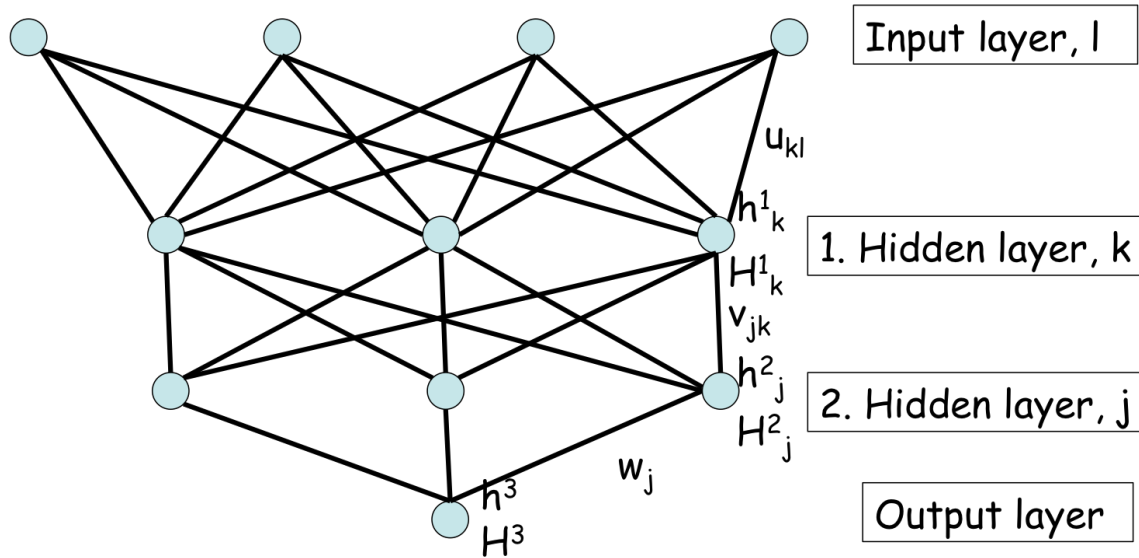
# Network architecture (input to hidden)



$$\frac{\partial E}{\partial u_{kl}} = \frac{\partial E}{\partial H^3} \cdot \frac{\partial H^3}{\partial h^3} \cdot \sum_j \frac{\partial h^3}{\partial H_j^2} \cdot \frac{\partial H_j^2}{\partial h_j^2} \cdot \frac{\partial h_j^2}{\partial H_k^1} \cdot \frac{\partial H_k^1}{\partial h_k^1} \cdot \frac{\partial h_k^1}{\partial u_{kl}}$$

$$= (H^3 - t) \cdot g'(h^3) \cdot \sum_j w_j \cdot g'(h_j^2) \cdot v_{jk} \cdot g'(h_k^1) \cdot I_l$$

# Speed. Use delta's



$$h^q_j = \sum_i w_{ji} H^{q-1}_i$$

$$H^{q-1}_i = g(h^{q-1}_i)$$

*Bishop, Christopher (1995). Neural networks for pattern recognition. Oxford: Clarendon Press. [ISBN 0-19-853864-2](https://doi.org/10.1017/CBO9780511566187).*



# Use delta's

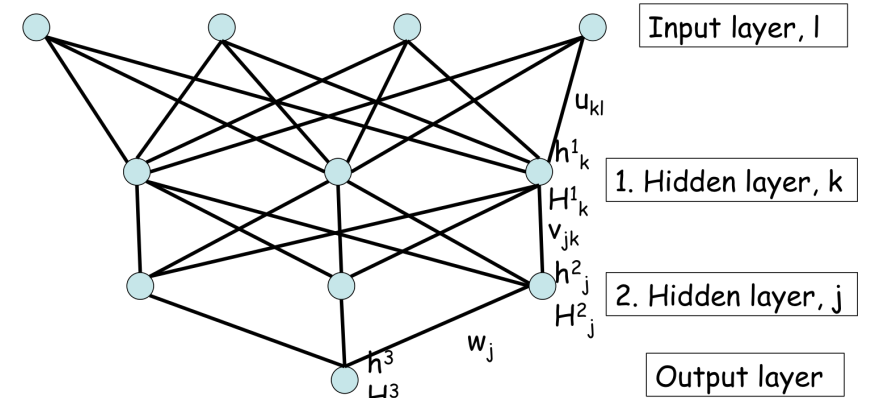
$$\frac{\partial E}{\partial w_{ji}^q} = \frac{\partial E}{\partial h_j^q} \cdot \frac{\partial h_j^q}{\partial w_{ji}^q} = \delta_j^q \cdot H_i^{q-1}$$

$$\delta_j^q = \frac{\partial E}{\partial h_j^q}$$

$$\delta^3 = \frac{\partial E}{\partial h^3} = \frac{\partial E}{\partial H^3} \cdot \frac{\partial H^3}{\partial h^3} = (H^3 - t) \cdot g'(h^3)$$

$$\delta_j^2 = \frac{\partial E}{\partial h_j^2} = \frac{\partial E}{\partial h^3} \cdot \frac{\partial h^3}{\partial h_j^2} = \frac{\partial E}{\partial h^3} \cdot \frac{\partial h^3}{\partial H_j^2} \cdot \frac{\partial H_j^2}{\partial h_j^2} = g'(h_j^2) \cdot \delta^3 \cdot v_{jk}$$

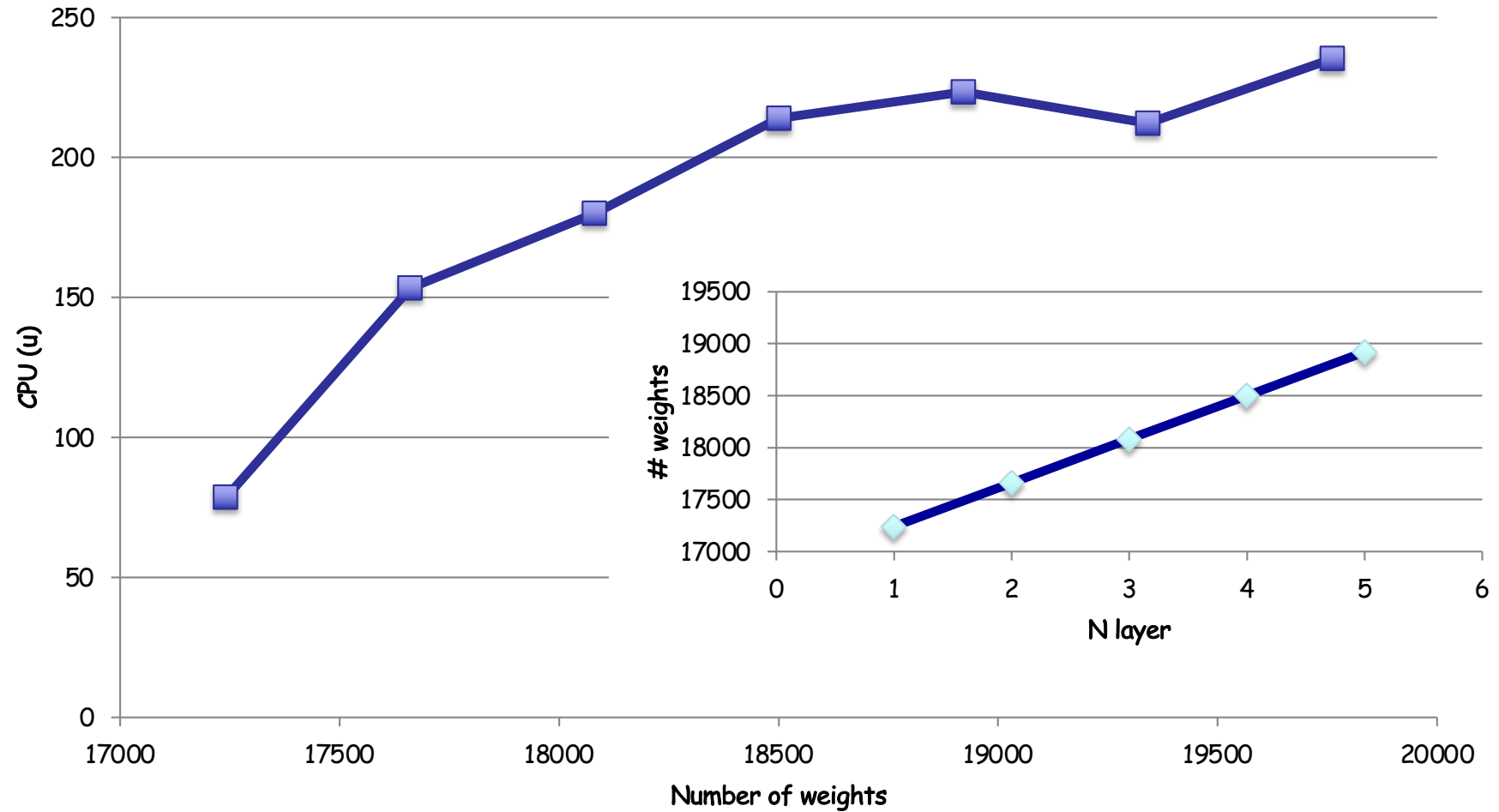
$$\delta_k^1 = \frac{\partial E}{\partial h_k^1} = \sum_j \frac{\partial E}{\partial h_j^2} \cdot \frac{\partial h_j^2}{\partial h_k^1} = \sum_j \frac{\partial E}{\partial h_j^2} \cdot \frac{\partial h_j^2}{\partial H_k^1} \cdot \frac{\partial H_k^1}{\partial h_k^1} = g'(h_k^1) \cdot \sum_j \delta_j^2 \cdot v_{jk}$$



$$h_j = \sum_i w_{ji} H_i$$

$$H_i = g(h_i)$$

# Deep learning - time is not an issue



## Deep Neural Networks

- Standard learning strategy
  - Randomly initializing the weights of the network
  - Applying gradient descent using backpropagation
- But, backpropagation does not work well (if randomly initialized)
  - Deep networks trained with back-propagation (without unsupervised pre-train) perform worse than shallow networks
  - ANN have limited to one or two layers

## Recent Deep Learning Highlights

- Google Goggles uses *Stacked Sparse Auto Encoders* (Hartmut Neven @ ICML 2011)
- The monograph or review paper *Learning Deep Architectures for AI* (Foundations & Trends in Machine Learning, 2009).
- *Exploring Strategies for Training Deep Neural Networks*, Hugo Larochelle, Yoshua Bengio, Jerome Louradour and Pascal Lamblin in: The Journal of Machine Learning Research, pages 1-40, 2009.
- The LISA publications database contains *a deep architectures category*. <http://www.iro.umontreal.ca/~lisa/twiki/bin/view.cgi/Public/ReadingOnDeepNetworks>
- Deep Machine Learning – *A New Frontier in Artificial Intelligence Research* – a survey paper by Itamar Arel, Derek C. Rose, and Thomas P. Karnowski.