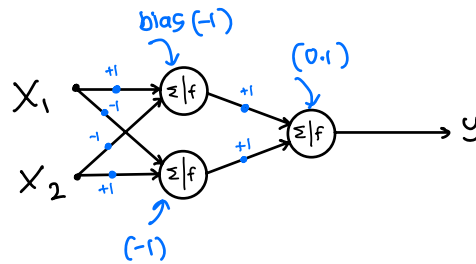


Cannot be separated with one line
Some problems need more than 1 line.

0/1

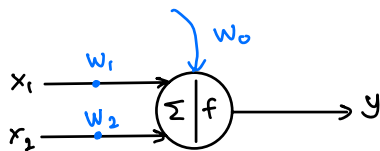
0/1



x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

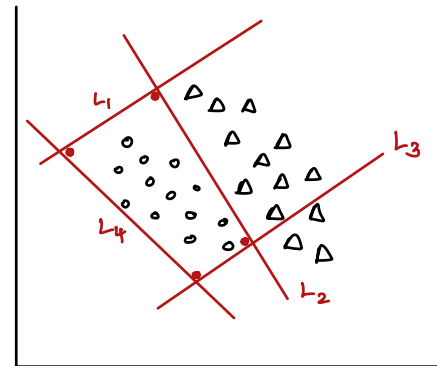
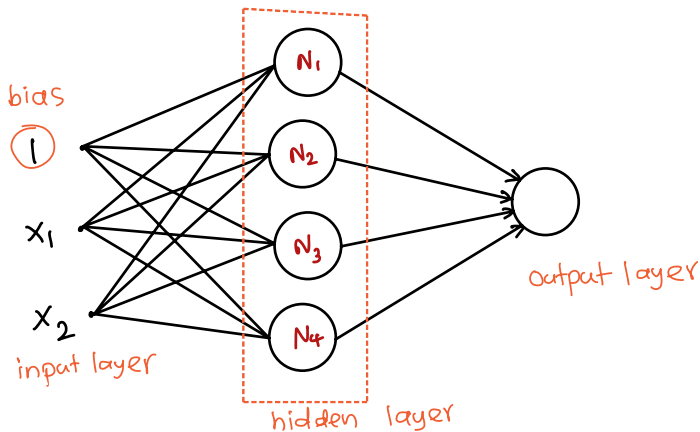
This works if used a "sign" logistic function $f(\cdot)$

This will enable a non linear separation via 2 lines. Neurons use the sign function due to its simplicity.

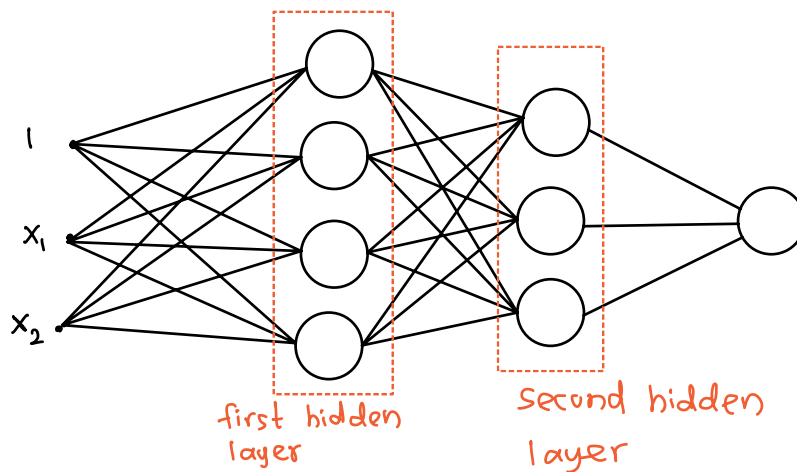


$$\frac{w_0 + w_1 x_1 + w_2 x_2 > 0}{w_0 + w_1 x_1 + w_2 x_2 < 0}$$

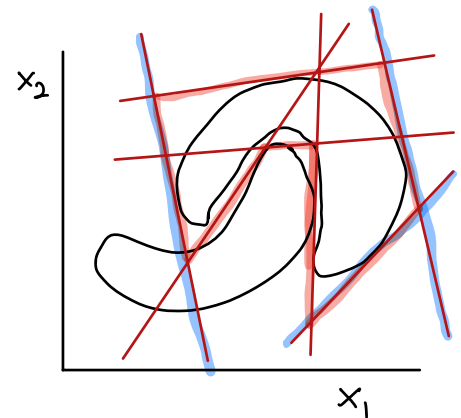
One hidden layer



Two hidden layer

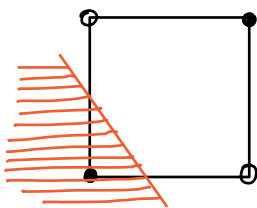


less freedom to locate lines
with 2 layers

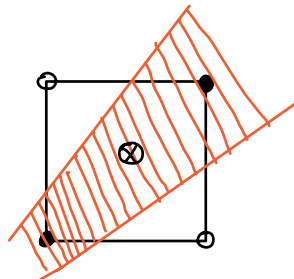


Any shape can be
bounded by # of
neurons

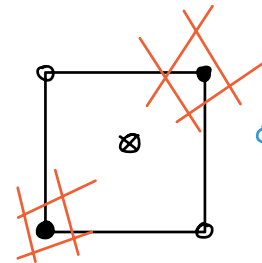
misclassified



one neuron
cannot solve the
problem



2 neurons solves
the problem



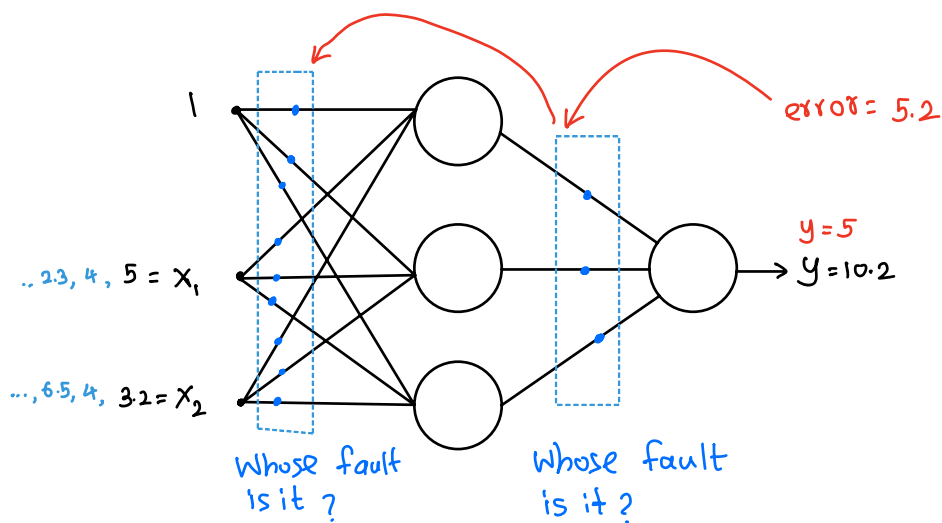
overfitted

2 layers solves the
problem too and more
accurately, but do you
need that?

Learning Algorithms

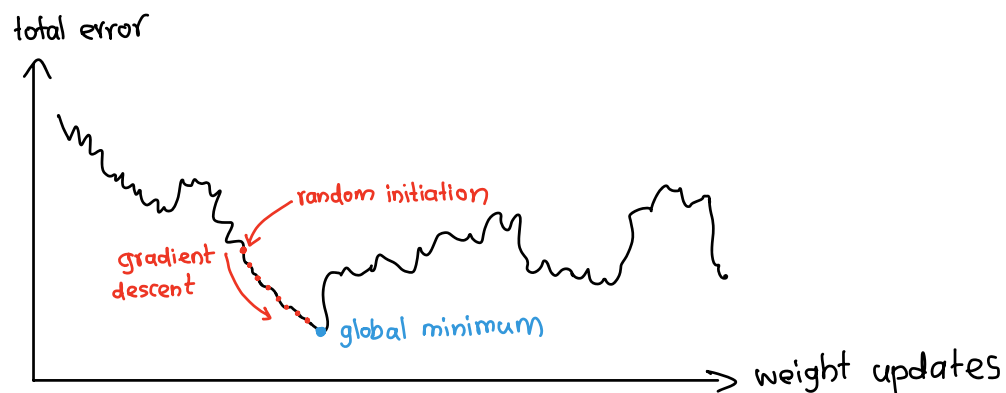
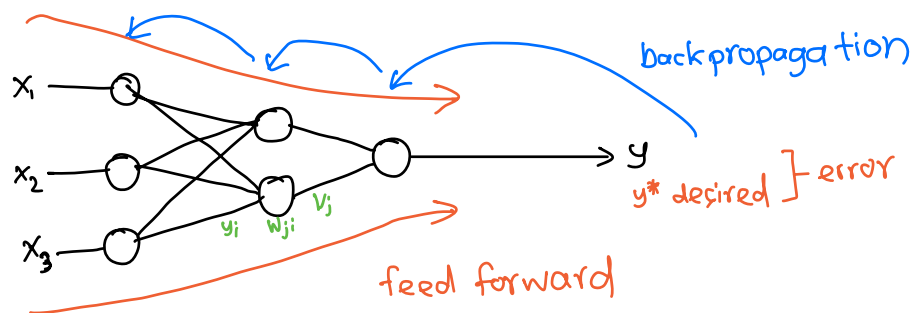
For MLPs (multi layer perceptrons) that we also called
feed forward MLPs.

We need solve the credit Assignment Problem.



Intelligence is updating the weights (= solving the credit assignment problem)

→ Backpropagating the error into the network.



Take a step in the direction resulting in a maximum decrease of the network error E .

This direction is the opposite of gradient of E.

Intelligence is updating the weights.

$$W_{ji}^{(n+1)} = W_{ji}^{(n)} + \Delta W_{ji}^{(n)}$$

$$\Delta W_{ji} = -\eta \frac{\partial E}{\partial W_{ji}} \quad (\text{opposite of gradient of error at local neuron})$$

$$\eta \in (0,1)$$

The input of the jth neuron: (In hidden layer internal neurons)

$$V_j = \sum_{i=1,2,\dots,m} W_{ji} \cdot y_i$$

Using the chain rule:

$$\frac{\partial E}{\partial W_{ji}} = \frac{\partial E}{\partial V_j} \cdot \frac{\partial V_j}{\partial W_{ji}}$$

Local gradient of the jth neuron:

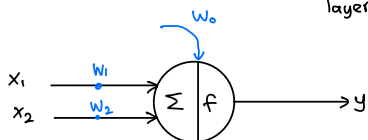
$$\delta_j = -\frac{\partial E}{\partial V_j}$$

Then from $\frac{\partial V_j}{\partial W_{ji}} = y_i$ we get

$$\Delta W_{ji} = \eta \cdot \delta_j \cdot y_i \quad \leftarrow \text{Delta Rule}$$

logistic function

$$\delta_j = \begin{cases} f'(V_j) (y_j^* - y_j) & \text{if } j \text{ is an output neuron} \\ f'(V_j) \sum_{k \text{ of next layer}} \delta_k W_{jk} & \text{if } j \text{ is an hidden neuron} \end{cases}$$



$$y = f(\sum x_i w_i + w_0)$$

$f(x)$ is the logistic function,

for instance, $f(x) = \frac{1}{1 + e^{-ax}}$

$$\begin{aligned} \frac{\partial}{\partial x} f(x) &= \frac{\partial}{\partial x} (1 + e^{-ax})^{-1} \\ &= -1 \cdot (1 + e^{-ax})^{-2} \cdot (-a e^{-ax}) \\ &= \frac{a \cdot e^{-ax}}{(1 + e^{-ax})^2} \end{aligned}$$

a is just a factor

$$= a \cdot \frac{1}{(1 + e^{-ax})} \cdot \frac{(1 + e^{-ax}) - 1}{(1 + e^{-ax})}$$

$$= \frac{a}{(1 + e^{-ax})} \cdot \left(1 - \frac{1}{1 + e^{-ax}} \right)$$

$$= a \cdot f(x) \cdot [1 - f(x)]$$

$$f'(x) = a \cdot f(x) \cdot [1 - f(x)]$$

$$f'(y_j) = a y_j (1 - y_j)$$

Backpropagation

$$n = 1$$

Initialize weights randomly $w(n)$

while (stopping criterion not satisfied)

 for each example (x, \underline{y}^*) supervised learning

```

    Run network with x and get y (feed f)
    Update weights in backpropagation.
  end for loop
  n = n + 1
end while

```

Backpropagation in Batch mode

Update weights only after all examples have been pushed forward through the network.

$$W_{ji}^{(n+1)} = W_{ji}^{(n)} + \sum_{x=\text{training sample}} \Delta W_{ji}^x$$

Training is epoch-by-epoch

Stopping criteria

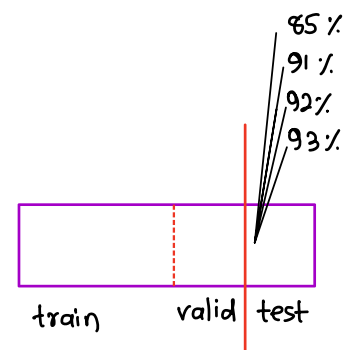
① Look at the MSE change

Network converged ($E \approx 0$) if the absolute rate of change in the average squared error per epoch is sufficiently small i.e., $[0.1, 0.01, \dots]$

② Generalization based method

Test for generalization after each epoch.

if adequate Generalization \longrightarrow stop



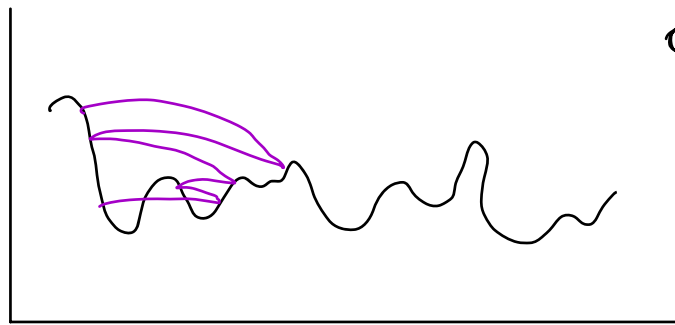
Delta Rule

$$\Delta W_{ji}^{(n)} = \eta \delta_j(n) y_i(n)$$

$\eta \rightarrow 0$: no learning

$\eta \rightarrow 1$: large changes \rightarrow unstable

(leads to weight oscillation)



$$\Delta W_{ji}(n) = \underbrace{\eta(n) \delta_j(n) y_i(n)}_{\text{begin(dominant)}} + \underbrace{\alpha(n) \Delta W_{ji}(n-1)}_{\text{end(dominant) } \alpha \in [0,1] \text{ momentum}}$$

$$\delta_j(n) = -\frac{\partial E}{\partial v_j}$$

Generalized Delta rule

Topology of the network

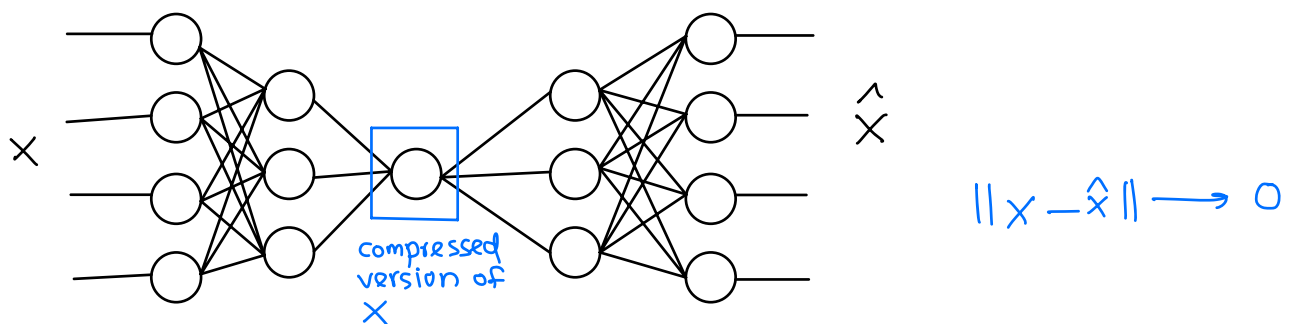
of layers
of neurons per layer } entirely task dependent
mostly done via trial and error

model size ① too small : underfitting
② too large : overfitting

Large net \rightarrow Remove neurons until performance starts to degrade \rightarrow Stop

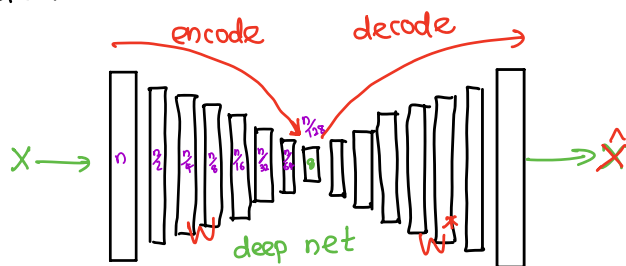
Small net \longrightarrow Add more neurons until performance is acceptable \longrightarrow stop

Autoencoders



Bottleneck concept

Force the network to reduce dimensionality of data.



* PCA is easy and more practical

$$W^* = W^T$$

* backprop doesnt work

inputs $x \in [0, 1]^d$

encoding $y \in [0, 1]^{d'}$ $d' \ll d$

decoding $z = g(w^*y + b^*)$

Error $\mathcal{L}(x, \hat{x}) = \|x - \hat{x}\|^2$

bit vector $L_H(x, \hat{x}) = - \sum_{k=1}^d [x_k \log \hat{x}_k + (1-x_k) \log (1-\hat{x}_k)]$

(error as tradition of entrophy)

cross entropy loss function

How to train deepnets

idea: Layerwise pretraining followed by greedy layerwise supervised training (with fine tuning)

