



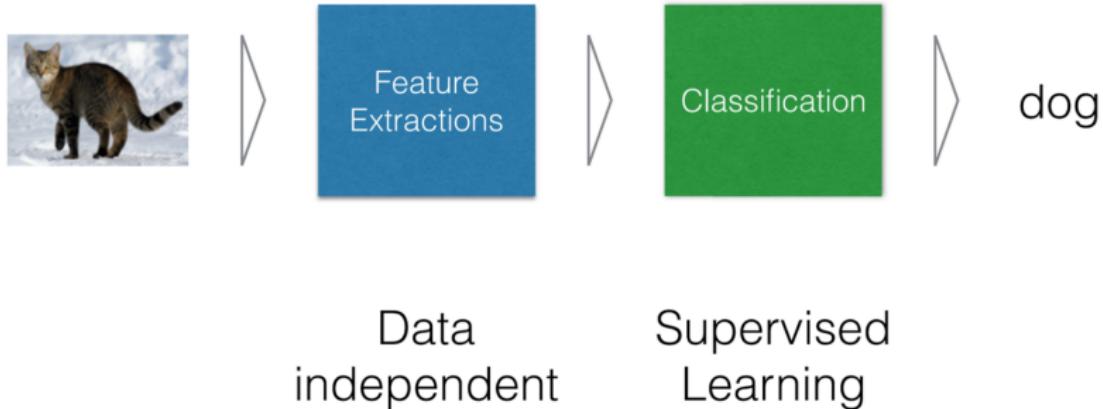
Introduction to Neural Networks

Saeed Reza Kheradpisheh

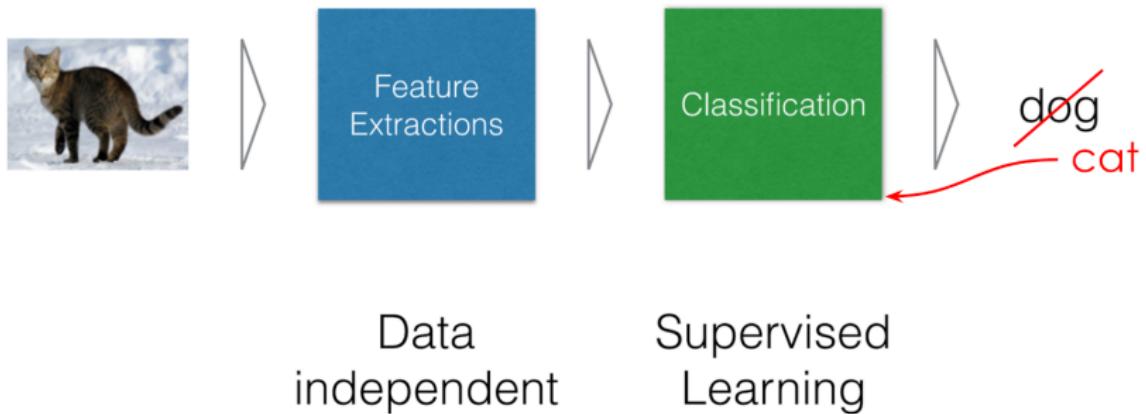
s.kheradpisheh@sbu.ac.ir

Department of Computer Science
Shahid Beheshti University
Summmer 1398

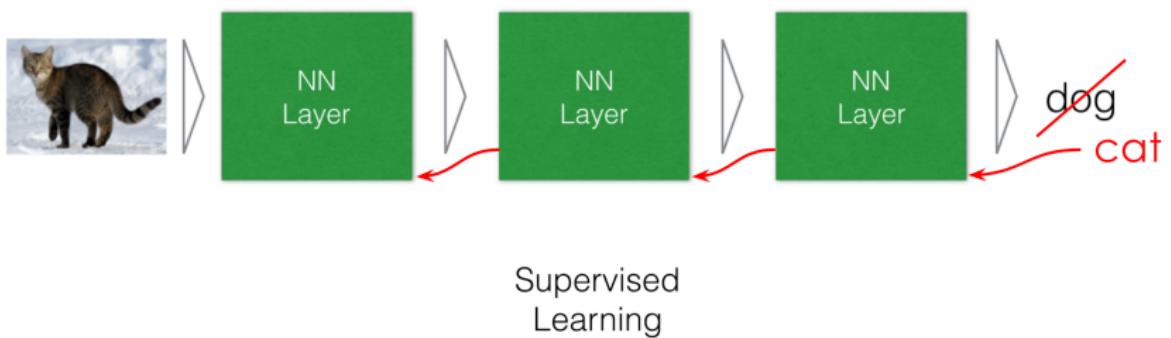
Typical ML system



Typical ML system

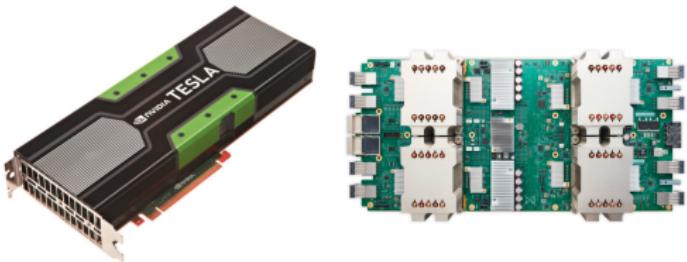


Deep Learning system



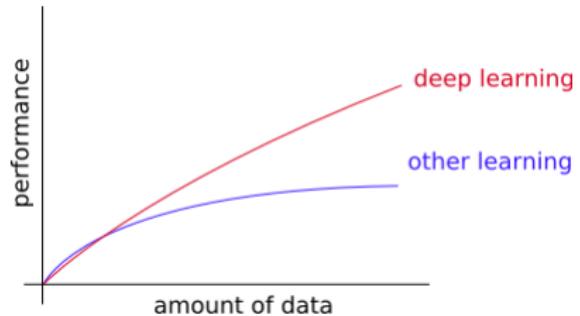
Why Deep Learning Now?

- Computing power (GPUs, TPUs, ...)



Why Deep Learning Now?

- Computing power (GPUs, TPUs, ...)
- Data with labels



Why Deep Learning Now?

- Computing power (GPUs, TPUs, ...)
- Data with labels
- Open source tools and models



theano



Microsoft
CNTK



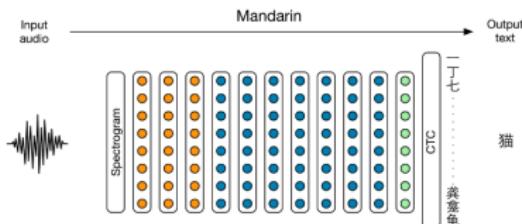
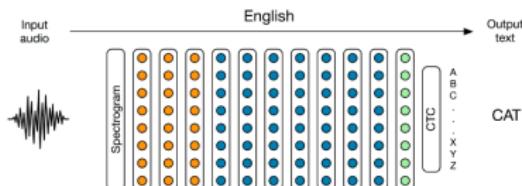
dmrc
mxnet

gensim spaCy

Why Deep Learning Now?

- Computing power (GPUs, TPUs, ...)
- Data with labels
- Open source tools and models
- Better algorithms & understanding

DL Today: Speech-to-Text



- Convolution Layer
- Recurrent Layer
- Fully Connected Layer

[Baidu 2014]

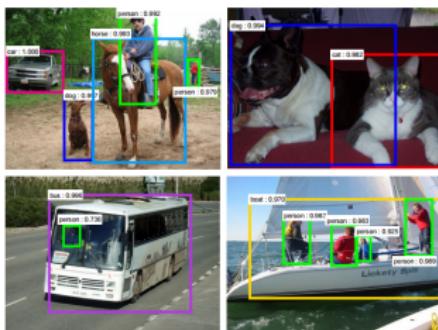
DL Today:Vision



[Krizhevsky 2012]



[Ciresan et al. 2013]



[Faster R-CNN - Ren 2015]

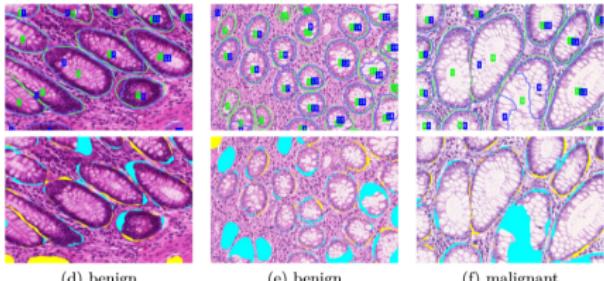


[NVIDIA dev blog]

DL Today:Vision



[Stanford 2017]



[Nvidia Dev Blog 2017]

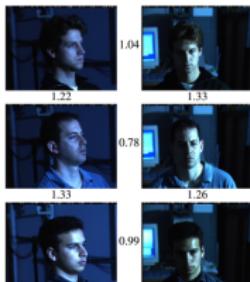
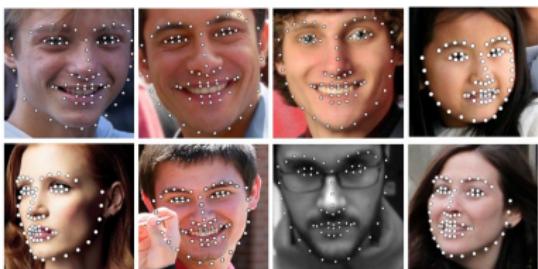


Figure 1. Illumination and Pose invariance.

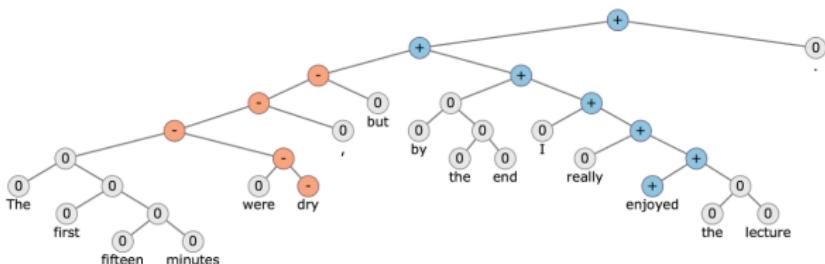
[FaceNet - Google 2015]



[Facial landmark detection CUHK 2014]

DL Today: Natural Language Processing (NLP)

The screenshot shows the Google Translate interface. On the left, the input field contains "deep learning" with a dropdown menu showing "Italian", "Chinese", "French", "English - detected", and a "Translate" button. The main area displays the French translation "l'apprentissage en profondeur" with a small info icon. Below the input field are icons for copy, paste, and history. At the bottom left, there's a "See also" section with "deep, learning". On the right, there are four colored buttons: "ENGLISH" (orange), "KOREAN" (blue), "JAPANESE" (red), and a decorative graphic of colored spheres connected by lines.



[Socher 2015]

DL Today: Natural Language Processing (NLP)

Salit Kulla
to me

11:29 AM ***

Hey, Wynton Marsalis is playing this weekend. Do you have a preference between Saturday and Sunday?

-S

I'm down for either.

Let's do Saturday.

I'm fine with whatever.



Reply

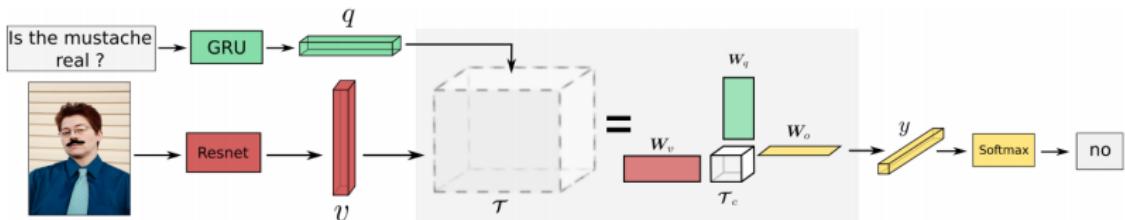


Forward

[Google Inbox Smart Reply]

[Amazon Echo / Alexa]

DL Today: Vision + NLP



[VQA - Mutan 2017]



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."

[Karpathy 2015]

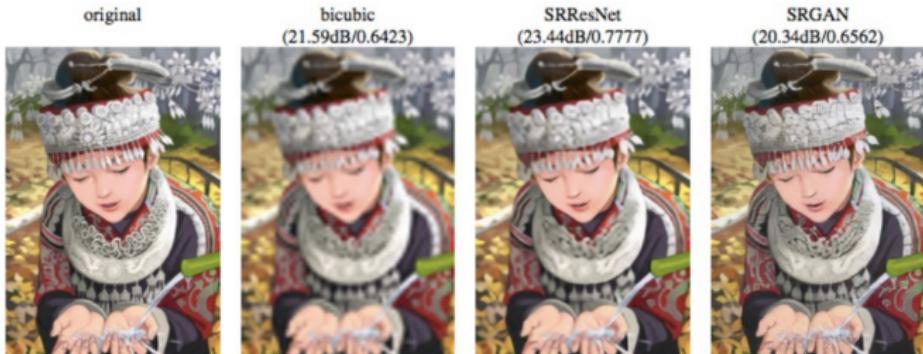
DL Today: Image translation



[DeepDream 2015]



[Gatys 2015]



[Ledig 2016]

DL Today: Generative models

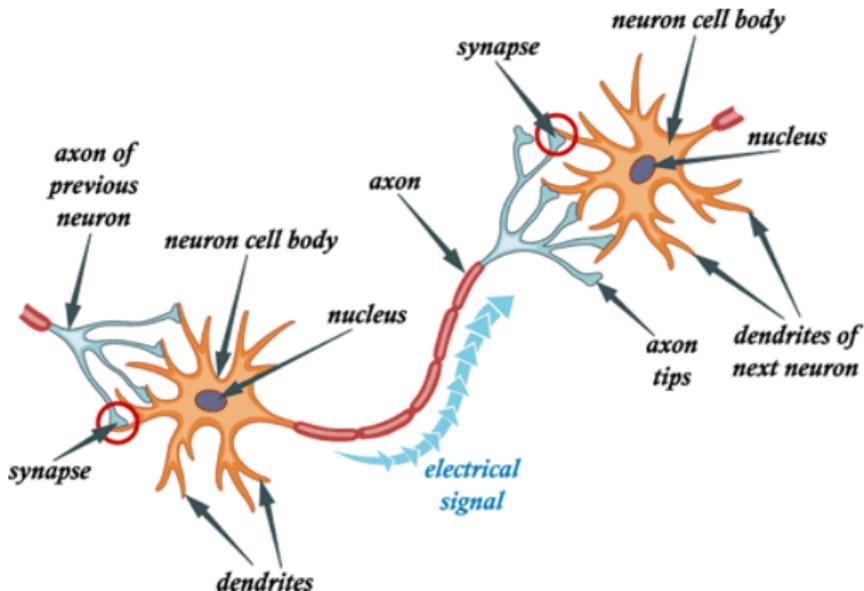
Sampled celebrities [Nvidia 2017]



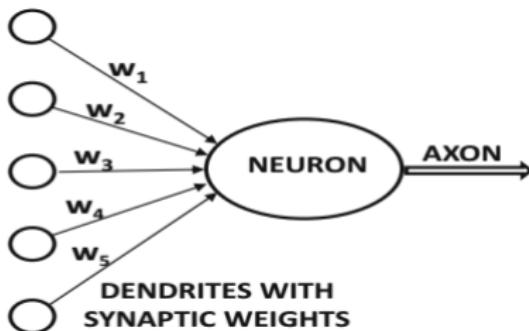
StackGAN v2 [Zhang 2017]

Text description	This bird is blue with white and has a very short beak	This bird has wings that are brown and has a yellow belly	A white bird with a black crown and yellow beak	This bird is white, black, and brown in color, with a brown beak	The bird has small beak, with reddish brown crown and gray belly	This is a small, black bird with a white breast and white on the wingbars.	This bird is white black and yellow in color, with a short black beak
Stage-I images	A row of four small, low-resolution images of birds, likely generated by a Stage-I GAN. They appear somewhat blurry and lack fine detail.				A row of seven larger, more detailed images of birds, likely generated by a Stage-II GAN. The images show more realistic features and better color representation compared to the Stage-I images.		
Stage-II images	A row of four small, low-resolution images of birds, likely generated by a Stage-I GAN. They appear somewhat blurry and lack fine detail.				A row of seven larger, more detailed images of birds, likely generated by a Stage-II GAN. The images show more realistic features and better color representation compared to the Stage-I images.		

Biological Neuron

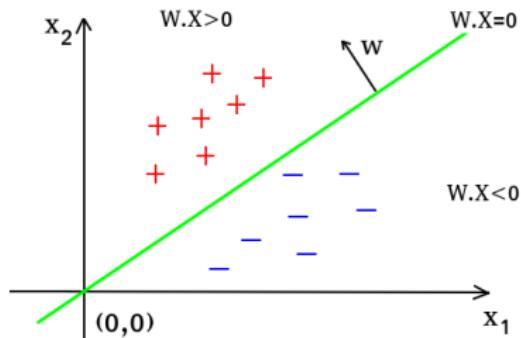
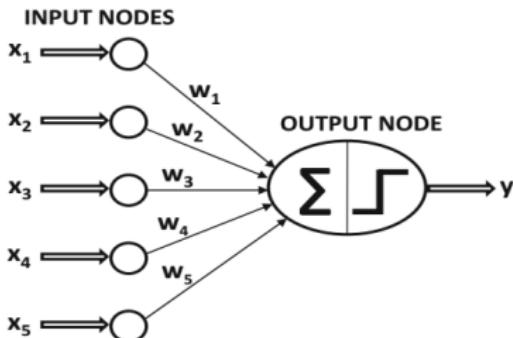


Artificial Neuron



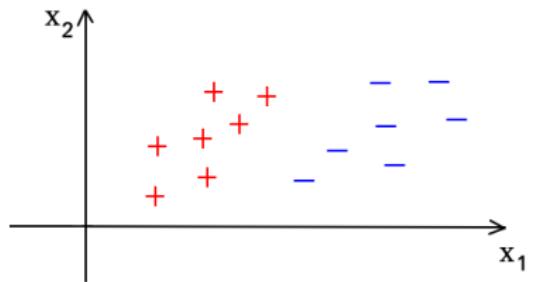
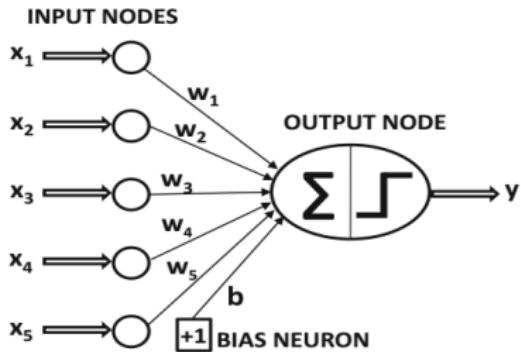
- Neuron output is a function of the inputs.
- Learning occurs by changing the weights to map inputs to outputs.
- Neural networks gain their power by putting together many of basic computing units.

Perceptron

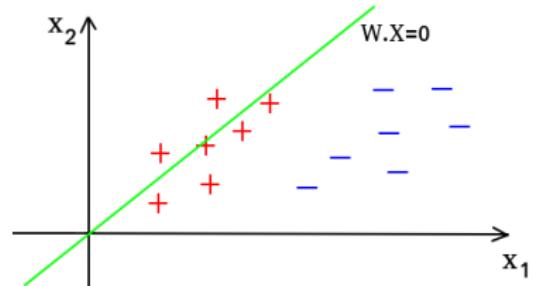
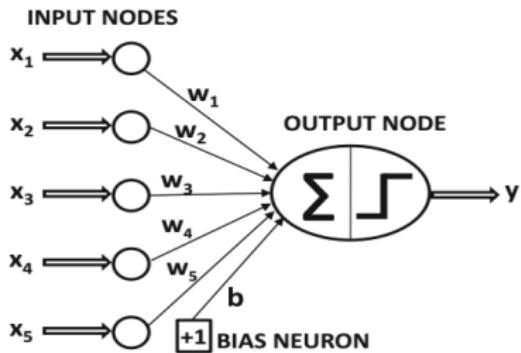


- Input vector: $X = [x_1, \dots, x_d]$
- Target output: $Y \in \{-1, +1\}$
- Input weights: $W = [w_1, \dots, w_d]$
- Predicted output: $y = \text{sign}\{W \cdot X\} = \text{sign}\{\sum_{i=1}^d w_i x_i\}$

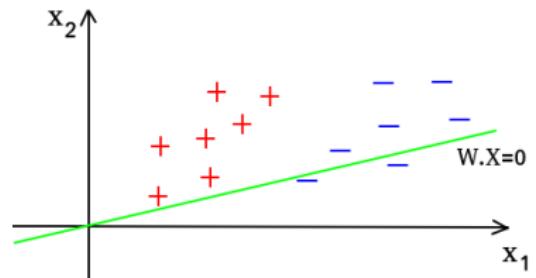
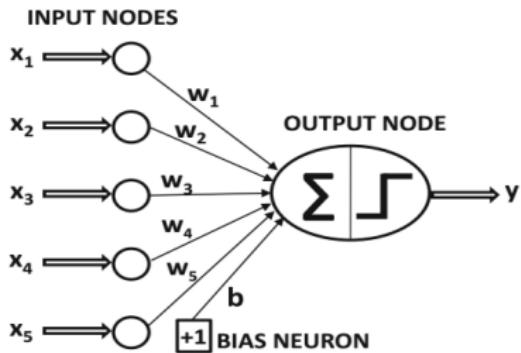
Perceptron with bias



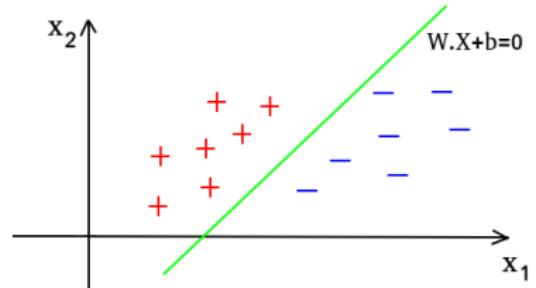
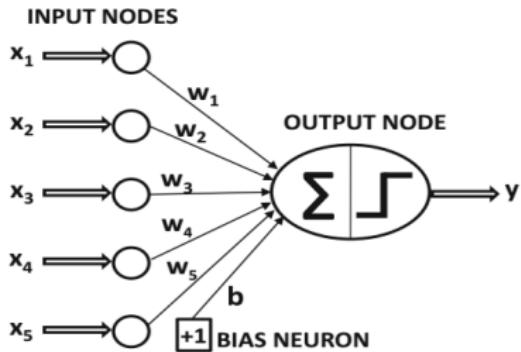
Perceptron with bias



Perceptron with bias



Perceptron with bias



Learning in Perceptrons

Consider a d -dimensional binary classification problem:

- Training set: $D = \{(X_i, Y_i) | i = 1 : N\}$
- Training sample: $X_i = [X_{i1}, \dots, X_{id}]$, $Y_i \in \{-1, +1\}$
- Perceptron predicts: $y_i = sign\{W.X_i\}$
- *Loss Function:*

$$L = \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2 = \sum_{(X_i, Y_i) \in D} (Y_i - sign\{W.X_i\})^2$$

Learning in Perceptrons

Consider a d -dimensional binary classification problem:

- Training set: $D = \{(X_i, Y_i) | i = 1 : N\}$
- Training sample: $X_i = [X_{i1}, \dots, X_{id}]$, $Y_i \in \{-1, +1\}$
- Perceptron predicts: $y_i = sign\{W.X_i\}$
- *Loss Function:*

$$L = \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2 = \sum_{(X_i, Y_i) \in D} (Y_i - sign\{W.X_i\})^2$$

- Loss function depends on W and D .

Learning in Perceptrons

Consider a d -dimensional binary classification problem:

- Training set: $D = \{(X_i, Y_i) | i = 1 : N\}$
- Training sample: $X_i = [X_{i1}, \dots, X_{id}]$, $Y_i \in \{-1, +1\}$
- Perceptron predicts: $y_i = sign\{W.X_i\}$
- *Loss Function:*

$$L = \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2 = \sum_{(X_i, Y_i) \in D} (Y_i - sign\{W.X_i\})^2$$

- Loss function depends on W and D .
- As D is given, hence, learning is to find W^* minimizing the loss function:

$$W^* = \operatorname{argmin}_W \sum_{(X_i, Y_i) \in D} (Y_i - sign\{W.X_i\})^2$$

How to find optimum weights?

- At any point $x = [x_1, \dots, x_d]$ of a function f , the gradient vector, $\nabla f(x)$, shows the direction of steepest ascend.

How to find optimum weights?

- At any point $x = [x_1, \dots, x_d]$ of a function f , the gradient vector, $\nabla f(x)$, shows the direction of steepest ascend.
- For an arbitrary weight vector, $-\nabla_w L$ shows the direction of the steepest descent of the loss function.

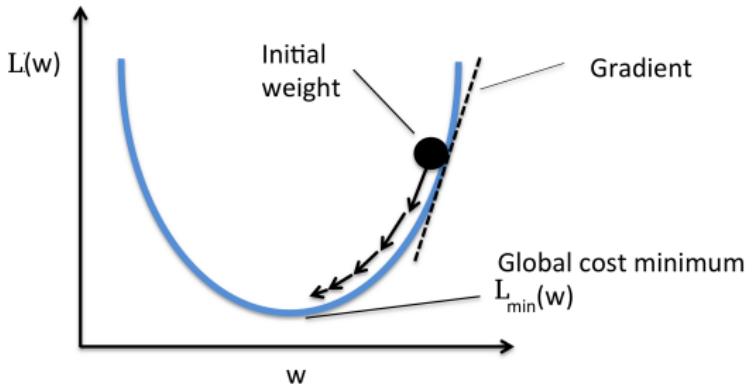
$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_d} \right]$$

How to find optimum weights?

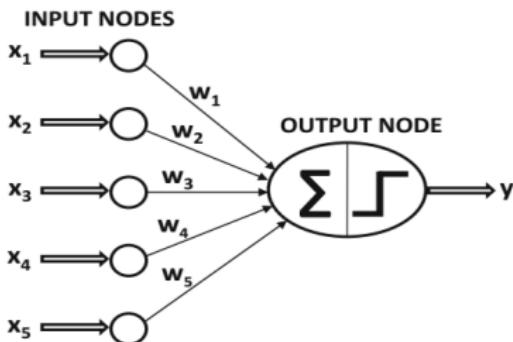
- At any point $x = [x_1, \dots, x_d]$ of a function f , the gradient vector, $\nabla f(x)$, shows the direction of steepest ascend.
- For an arbitrary weight vector, $-\nabla_w L$ shows the direction of the steepest descent of the loss function.

$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_d} \right]$$

- Find W^* by starting from a random weight vector and an iterative use of gradient:



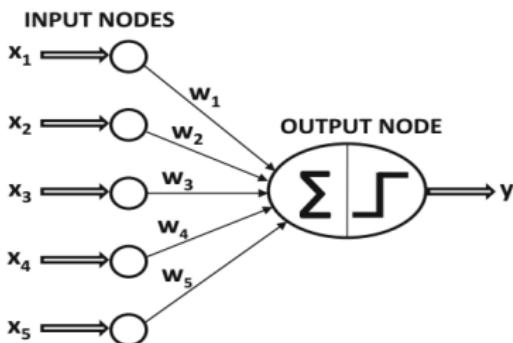
Computing gradients for Perceptron



$$L = \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - \text{sign}\{W \cdot X_i\})^2$$

Computing gradients for Perceptron

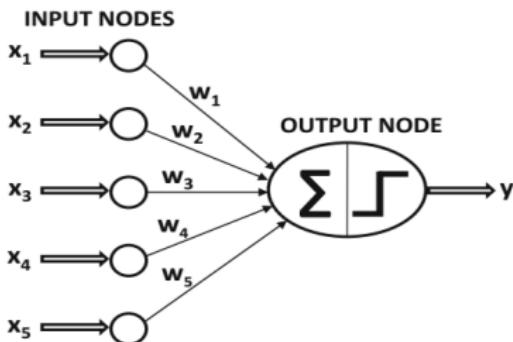


$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \dots, \frac{\partial f}{\partial w_d} \right]$$

$$L = \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - \text{sign}\{W \cdot X_i\})^2$$

Computing gradients for Perceptron



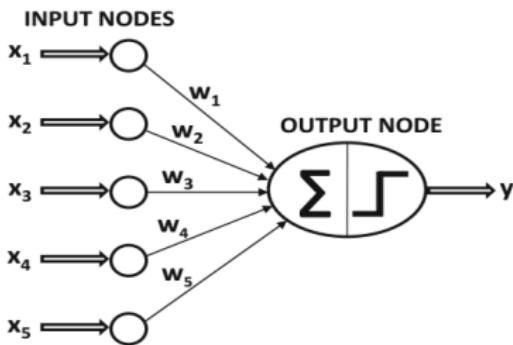
$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_d} \right]$$

$$\frac{\partial L}{\partial w_j} = \frac{\partial \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2}{\partial w_j}$$

$$L = \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - \text{sign}\{W \cdot X_i\})^2$$

Computing gradients for Perceptron



$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_d} \right]$$

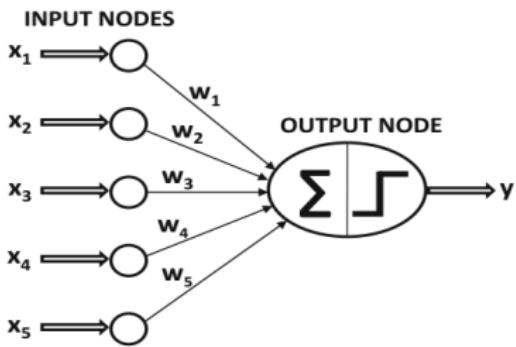
$$\frac{\partial L}{\partial w_j} = \frac{\partial \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2}{\partial w_j}$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} \frac{\partial (Y_i - y_i)^2}{\partial w_j}$$

$$L = \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - \text{sign}\{W \cdot X_i\})^2$$

Computing gradients for Perceptron



$$L = \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - \text{sign}\{W \cdot X_i\})^2$$

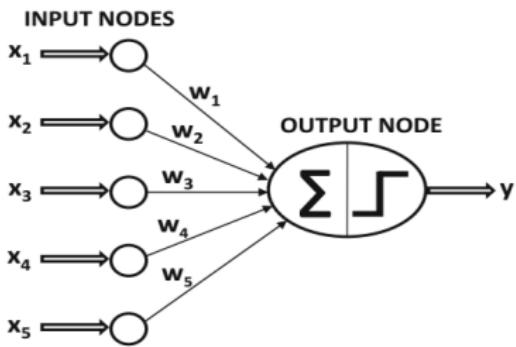
$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_d} \right]$$

$$\frac{\partial L}{\partial w_j} = \frac{\partial \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2}{\partial w_j}$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} \frac{\partial (Y_i - y_i)^2}{\partial w_j}$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} \frac{\partial (Y_i - y_i)^2}{\partial y_i} \frac{\partial y_i}{\partial w_j}$$

Computing gradients for Perceptron



$$L = \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - \text{sign}\{W \cdot X_i\})^2$$

$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_d} \right]$$

$$\frac{\partial L}{\partial w_j} = \frac{\partial \frac{1}{2} \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2}{\partial w_j}$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} \frac{\partial (Y_i - y_i)^2}{\partial w_j}$$

$$= \frac{1}{2} \sum_{(X_i, Y_i) \in D} \frac{\partial (Y_i - y_i)^2}{\partial y_i} \frac{\partial y_i}{\partial w_j}$$

$$= - \sum_{(X_i, Y_i) \in D} (Y_i - y_i) \frac{\partial y_i}{\partial w_j}$$

Computing gradients for Perceptron

The derivative of the Perceptron's predicted output is zero everywhere and is undefined at zero:

$$y_i = \text{sign}\{W.X_i\} = \text{sign}\left\{\sum_{j=1}^d w_j x_{ij}\right\}$$

$$\frac{\partial y_i}{\partial w_j} = \frac{\partial \text{sign}\{W.X_i\}}{\partial W.X_i} \frac{\partial W.X_i}{\partial w_j} = \frac{\partial \text{sign}\{W.X_i\}}{\partial W.X_i} x_{ij}$$

Computing gradients for Perceptron

The derivative of the Perceptron's predicted output is zero everywhere and is undefined at zero:

$$y_i = \text{sign}\{W.X_i\} = \text{sign}\left\{\sum_{j=1}^d w_j x_{ij}\right\}$$

$$\frac{\partial y_i}{\partial w_j} = \frac{\partial \text{sign}\{W.X_i\}}{\partial W.X_i} \frac{\partial W.X_i}{\partial w_j} = \frac{\partial \text{sign}\{W.X_i\}}{\partial W.X_i} x_{ij}$$

Hence, a surrogate gradient is used:

$$\frac{\partial y_i}{\partial w_j} = x_{ij}$$

Thus we have:

$$\frac{\partial L}{\partial w_j} = - \sum_{(X_i, Y_i) \in D} (Y_i - y_i) x_{ij}$$

Gradient Descent Learning for Perceptron

To find optimum weights (W^*):

- Start from a random initial weight vector, W^0 .

Gradient Descent Learning for Perceptron

To find optimum weights (W^*):

- Start from a random initial weight vector, W^0 .
- Through an iterative manner, use gradients and update the weights:

$$W^{t+1} = W^t + \eta \sum_{(X_i, Y_i) \in D} (Y_i - y_i) X_i$$

Gradient Descent Learning for Perceptron

To find optimum weights (W^*):

- Start from a random initial weight vector, W^0 .
- Through an iterative manner, use gradients and update the weights:

$$W^{t+1} = W^t + \eta \sum_{(X_i, Y_i) \in D} (Y_i - y_i) X_i$$

- The learning rate is controlled by η .

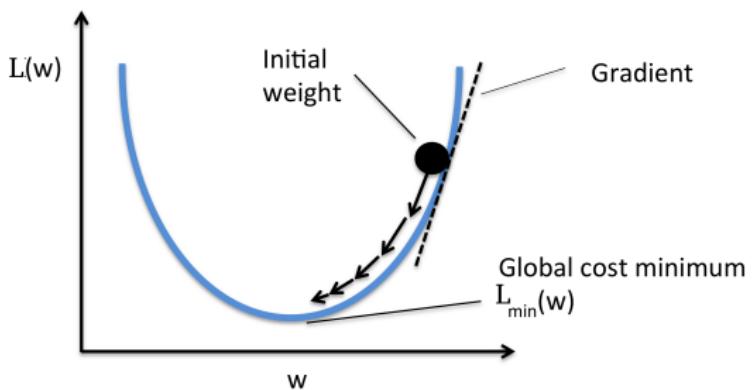
Gradient Descent Learning for Perceptron

To find optimum weights (W^*):

- Start from a random initial weight vector, W^0 .
- Through an iterative manner, use gradients and update the weights:

$$W^{t+1} = W^t + \eta \sum_{(X_i, Y_i) \in D} (Y_i - y_i) X_i$$

- The learning rate is controlled by η .



Stochastic Gradient Descent (SGD) for Perceptron

In SGD, learning is performed sample by sample:

- ➊ Shuffle the training set.

Stochastic Gradient Descent (SGD) for Perceptron

In SGD, learning is performed sample by sample:

- ① Shuffle the training set.
- ② Compute the perceptron's output y_i for the i -th sample.

Stochastic Gradient Descent (SGD) for Perceptron

In SGD, learning is performed sample by sample:

- ① Shuffle the training set.
- ② Compute the perceptron's output y_i for the i -th sample.
- ③ Update the weights using $W^{t+1} = W^t + \eta(Y_i - y_i)X_i$.

Stochastic Gradient Descent (SGD) for Perceptron

In SGD, learning is performed sample by sample:

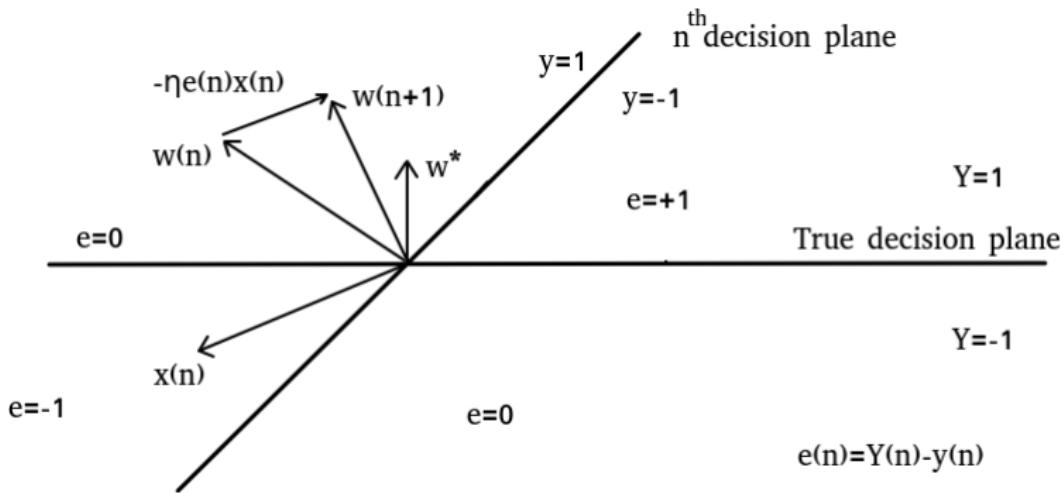
- ① Shuffle the training set.
- ② Compute the perceptron's output y_i for the i -th sample.
- ③ Update the weights using $W^{t+1} = W^t + \eta(Y_i - y_i)X_i$.
- ④ Repeat steps 2 to 3 for all training samples.

Stochastic Gradient Descent (SGD) for Perceptron

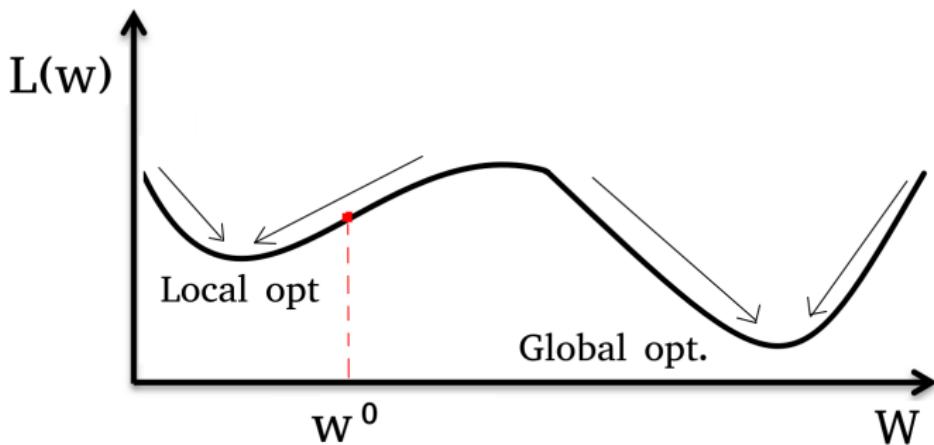
In SGD, learning is performed sample by sample:

- ① Shuffle the training set.
- ② Compute the perceptron's output y_i for the i -th sample.
- ③ Update the weights using $W^{t+1} = W^t + \eta(Y_i - y_i)X_i$.
- ④ Repeat steps 2 to 3 for all training samples.
- ⑤ Jump to step 1 if the total loss $L = \sum_{(X_i, Y_i) \in D} (Y_i - y_i)^2$ is below a certain value or the maximum number of iteration is reached

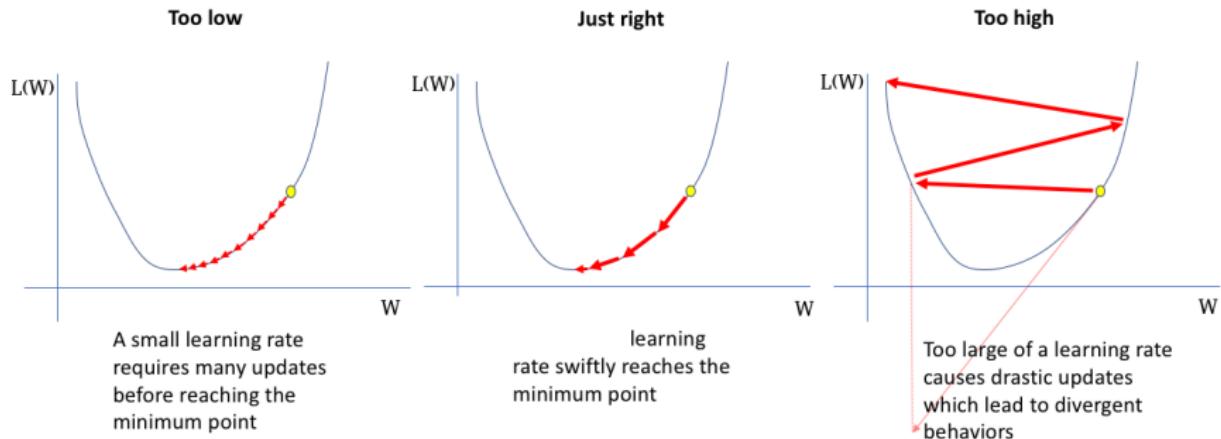
SGD



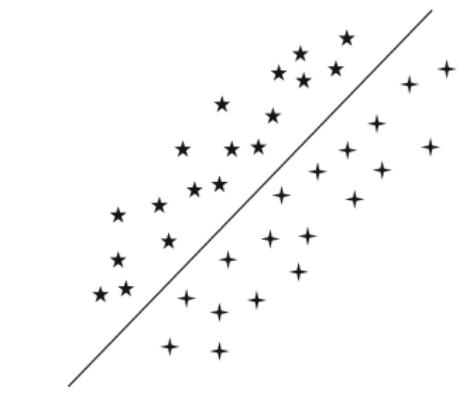
Initial weights matter



Learning rate matters



Perceptron solves linearly separable problems



$$\overline{w} \cdot \overline{x} = 0$$

LINEARLY SEPARABLE



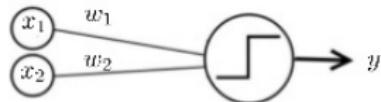
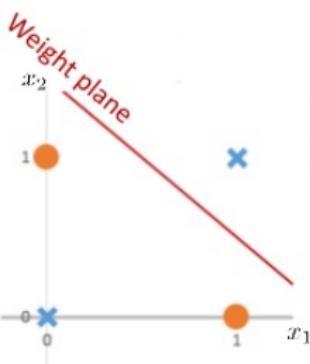
NOT LINEARLY SEPARABLE

The XOR problem

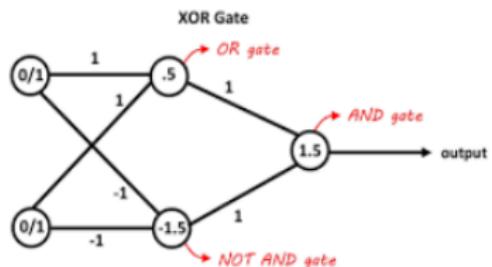
XOR Problem

Training Data

x_1	x_2	t
0	0	0
1	0	1
0	1	1
1	1	0

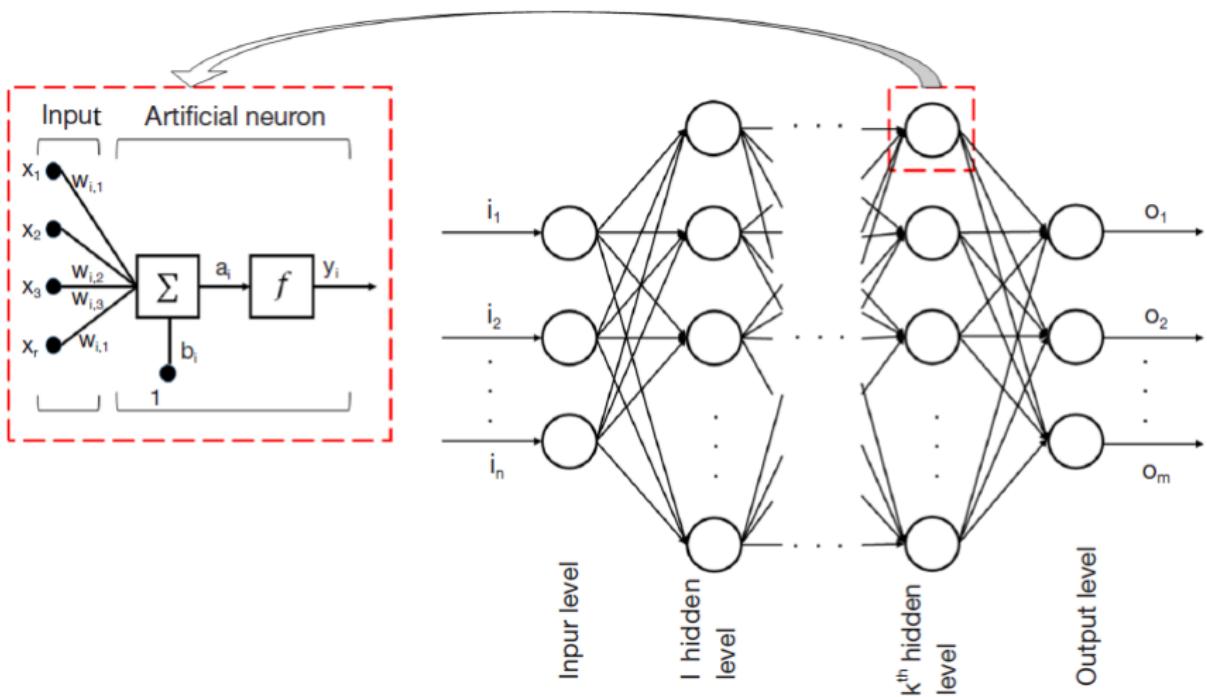


A single perceptron can only solve linear problems.

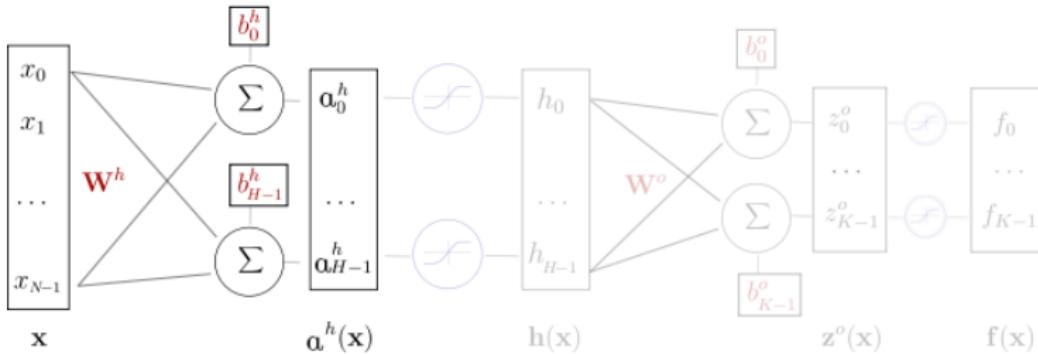


Multi-layer perceptron can solve non-linearly separable problems.

Multi Layer Perceptron (MLP)

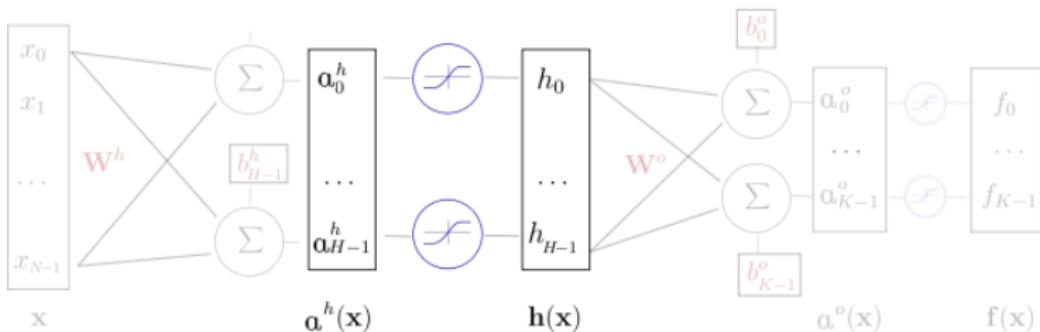


Forward propagation



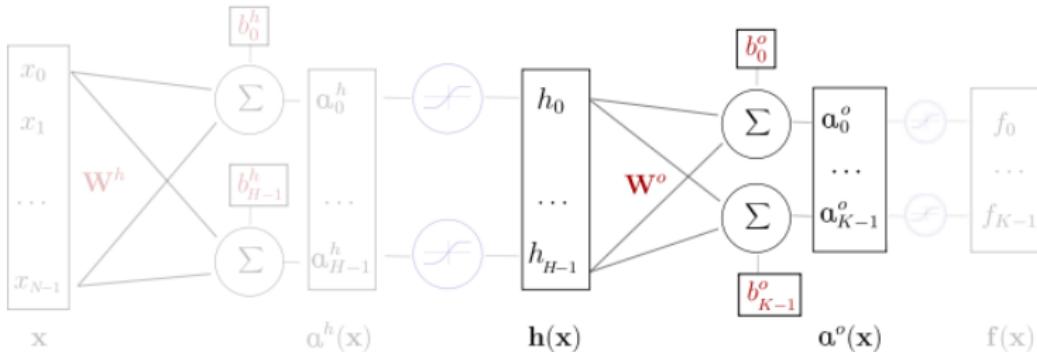
- $a^h(x) = W^h \cdot x + b^h$

Forward propagation



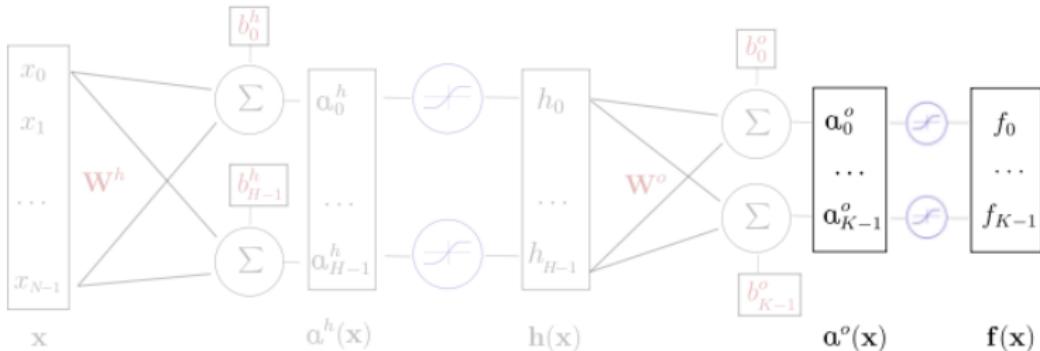
- $a^h(x) = W^h \cdot x + b^h$
- $h(x) = \Phi(a^h(x)) = \Phi(W^h \cdot x + b^h)$

Forward propagation



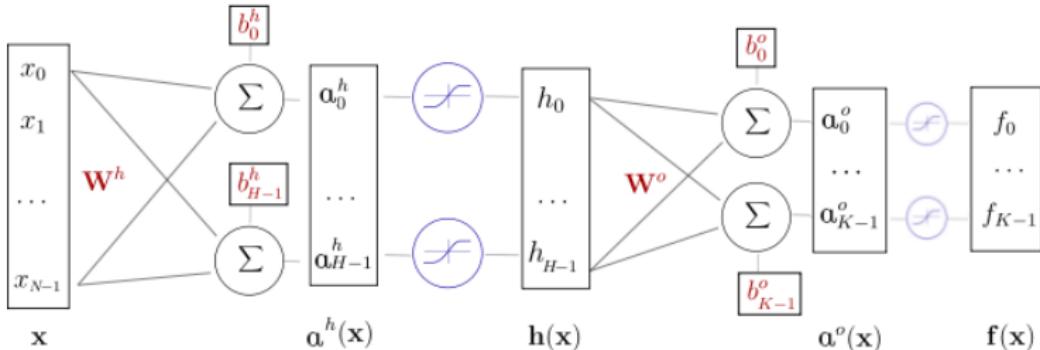
- $a^h(x) = W^h \cdot x + b^h$
- $h(x) = \Phi(a^h(x)) = \Phi(W^h \cdot x + b^h)$
- $a^o(x) = W^o \cdot h(x) + b^o$

Forward propagation



- $a^h(x) = W^h \cdot x + b^h$
- $h(x) = \Phi(a^h(x)) = \Phi(W^h \cdot x + b^h)$
- $a^o(x) = W^o \cdot h(x) + b^o$
- $f(x) = \Phi(a^o(x)) = \Phi(W^o \cdot h(x) + b^o)$

Forward propagation



- $a^h(x) = W^h \cdot x + b^h$
- $h(x) = \Phi(a^h(x)) = \Phi(W^h \cdot x + b^h)$
- $a^o(x) = W^o \cdot h(x) + b^o$
- $f(x) = \Phi(a^o(x)) = \Phi(W^o \cdot h(x) + b^o)$

Multilayer Network as a Computational Graph

- A multilayer network evaluates compositions of functions computed at individual nodes:

$$f(g_1(.), \dots, g_k(.)) \quad (1)$$

Multilayer Network as a Computational Graph

- A multilayer network evaluates compositions of functions computed at individual nodes:

$$f(g_1(.), \dots, g_k(.)) \quad (1)$$

- The use of nonlinear activation functions is the key to increase the power of multiple layers.

Multilayer Network as a Computational Graph

- A multilayer network evaluates compositions of functions computed at individual nodes:

$$f(g_1(.), \dots, g_k(.)) \quad (1)$$

- The use of nonlinear activation functions is the key to increase the power of multiple layers.
- A network with a single hidden layer of nonlinear units can compute almost any reasonable function.

Multilayer Network as a Computational Graph

- A multilayer network evaluates compositions of functions computed at individual nodes:

$$f(g_1(.), \dots, g_k(.)) \quad (1)$$

- The use of nonlinear activation functions is the key to increase the power of multiple layers.
- A network with a single hidden layer of nonlinear units can compute almost any reasonable function.
- The number of hidden units required to do so is rather large, which increases the number of parameters to be learned.

The Power of Function Composition

A multi-layer network that uses only the identity activation function in all its layers reduces to a single-layer network performing linear regression.

$$\bar{h}_1 = \Phi(W_1^T \bar{x}) = W_1^T \bar{x}$$

$$\bar{h}_{p+1} = \Phi(W_{p+1}^T \bar{h}_p) = W_{p+1}^T \bar{h}_p \quad \forall p \in \{1 \dots k-1\}$$

$$\bar{o} = \Phi(W_{k+1}^T \bar{h}_k) = W_{k+1}^T \bar{h}_k$$

$$\begin{aligned}\bar{o} &= W_{k+1}^T W_k^T \dots W_1^T \bar{x} \\ &= \underbrace{(W_1 W_2 \dots W_{k+1})^T}_{W_{xo}^T} \bar{x}\end{aligned}$$

Element-wise Activation Functions

Sign function: $\Phi(a) = sign(a)$

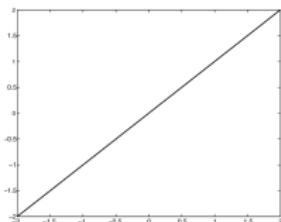
Sigmoid function: $\Phi(a) = \frac{1}{1+e^{-a}}$

Tanh function: $\Phi(a) = \frac{e^{2a}-1}{e^{2a}+1}$

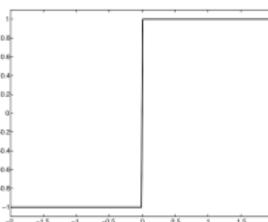
ReLU: $\Phi(a) = max\{a, 0\}$

Hard Tanh:

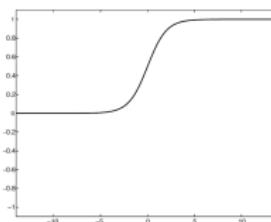
$\Phi(a) = max\{min[v, 1], -1\}$



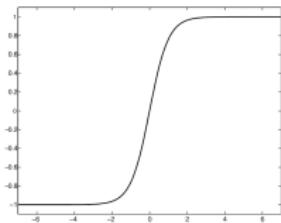
(a) Identity



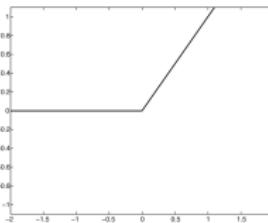
(b) Sign



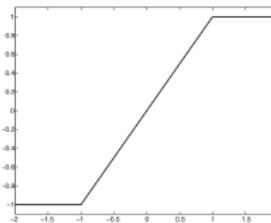
(c) Sigmoid



(d) Tanh

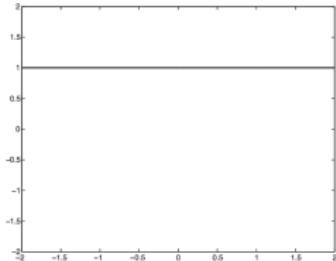


(e) ReLU

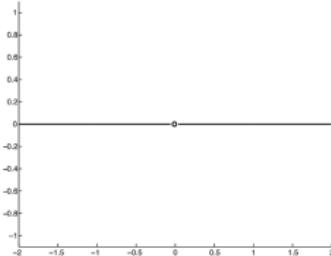


(f) Hard Tanh

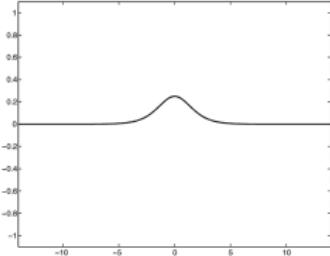
Dervation of Activation Functions



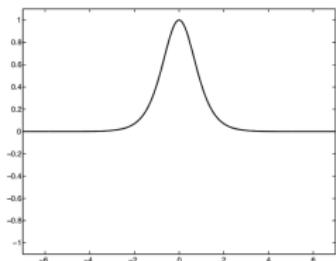
(a) Identity



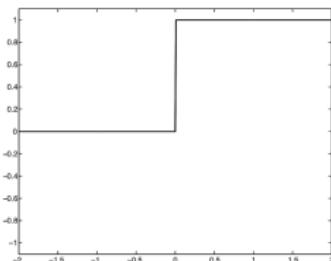
(b) Sign



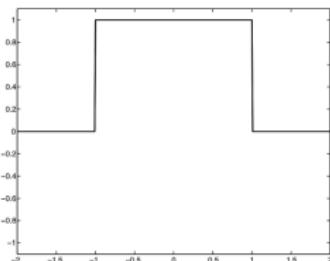
(c) Sigmoid



(d) Tanh



(e) ReLU



(f) Hard Tanh

Dervation of Activation Functions

Assume $o = \Phi(v)$ thus we have:

Sigmoid function:

$$\frac{\partial o}{\partial v} = \frac{\exp(-v)}{(1 + \exp(-v))^2}$$

$$\frac{\partial o}{\partial v} = o(1 - o)$$

Tangh function:

$$\frac{\partial o}{\partial v} = \frac{4 \cdot \exp(2v)}{(\exp(2v) + 1)^2}$$

$$\frac{\partial o}{\partial v} = 1 - o^2$$

Groupe Activation: Softmax function

$$\text{softmax}(x) = \frac{1}{\sum_{i=1}^n e^{x_i}} \cdot \begin{bmatrix} e^{x_1} \\ \vdots \\ e^{x_n} \end{bmatrix}$$

Groupe Activation: Softmax function

$$\text{softmax}(x) = \frac{1}{\sum_{i=1}^n e^{x_i}} \cdot \begin{bmatrix} e^{x_1} \\ \vdots \\ e^{x_n} \end{bmatrix}$$

- Softmax activation for each neuron is in range $[0,1]$.

Groupe Activation: Softmax function

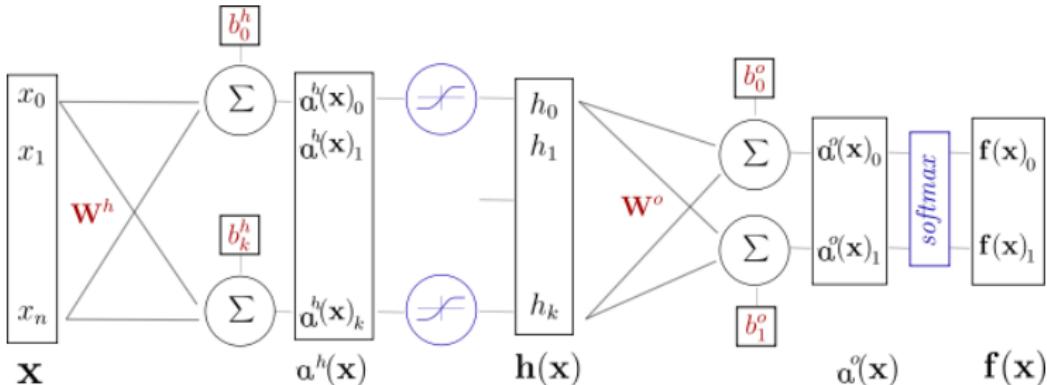
$$\text{softmax}(x) = \frac{1}{\sum_{i=1}^n e^{x_i}} \cdot \begin{bmatrix} e^{x_1} \\ \vdots \\ e^{x_n} \end{bmatrix}$$

- Softmax activation for each neuron is in range $[0,1]$.
- The summation of neurons' activation is 1.

Groupe Activation: Softmax function

$$\text{softmax}(x) = \frac{1}{\sum_{i=1}^n e^{x_i}} \cdot \begin{bmatrix} e^{x_1} \\ \vdots \\ e^{x_n} \end{bmatrix}$$

- Softmax activation for each neuron is in range $[0,1]$.
- The summation of neurons' activation is 1.
- It is usually used in the output layer.



Dervation of the Softmax Function

Assume $o_i = \text{softmax}(v_i)$ thus we have:

$$\frac{\partial o_i}{\partial v_j} = \begin{cases} o_i \cdot (1 - o_i) & i = j \\ -o_i \cdot o_j & i \neq j \end{cases}$$

Error Backpropagation

- Consider a sequence of hidden units followed by an output unit.
- To update any weight of the output layer, we use the gradient of the loss function with respect to that weight.
- Consider the Loss Function to be $L(X) = \frac{1}{2}(Y - o)^2$



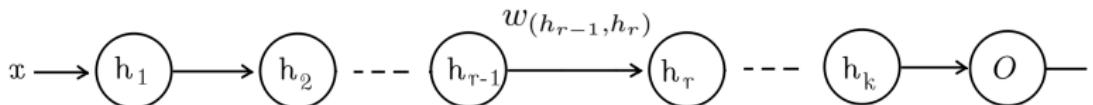
$$\frac{\partial L}{\partial w_{(h_k, o)}} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial w_{(h_k, o)}} = \Delta(o, o) \cdot h_k \cdot \Phi'(a_o)$$

$$\Delta(o, o) = \frac{\partial L}{\partial o} = -(Y - O)$$

$$\frac{\partial o}{\partial w_{(h_k, o)}} = \frac{\partial o}{\partial a_o} \frac{\partial a_o}{\partial w_{(h_k, o)}} = h_k \cdot \Phi'(a_o)$$

Error Backpropagation

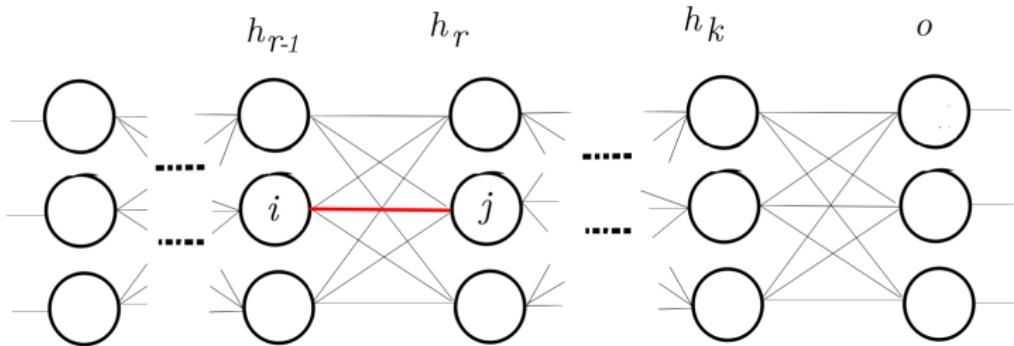
- Consider a sequence of hidden units followed by an output unit.
- To update a connection weight, we should compute the gradient of the loss function with respect to that weight.



$$\frac{\partial L}{\partial w_{(h_{r-1}, h_r)}} = \frac{\partial L}{\partial o} \cdot \left[\frac{\partial o}{\partial h_k} \prod_{i=r}^{k-1} \frac{\partial h_{i+1}}{\partial h_i} \right] \frac{\partial h_r}{\partial w_{(h_{r-1}, h_r)}}$$

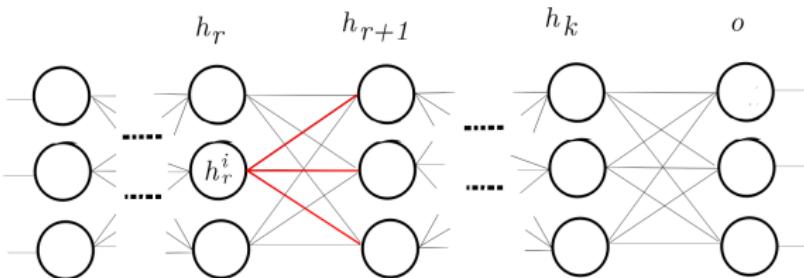
Error Backpropagation

- Now consider a sequence of hidden layers followed by an output unit (P is the set of paths from h_r to o):

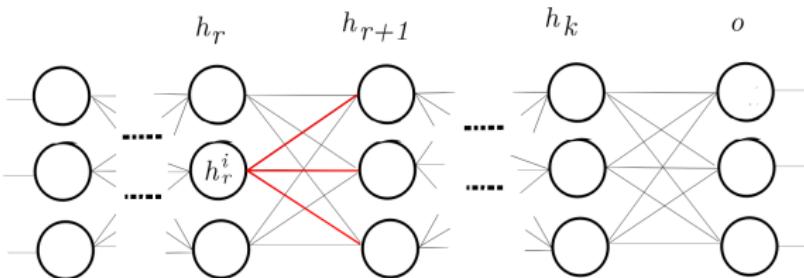


$$\frac{\partial L}{\partial w_{(h_{r-1}^i, h_r^j)}} = \underbrace{\left[\sum_{[h_r^j, h_{r+1}, \dots, h_k, o] \in \mathcal{P}} \frac{\partial L}{\partial o} \frac{\partial o}{\partial h_k} \prod_{i=r}^{k-1} \frac{\partial h_{i+1}}{\partial h_i} \right]}_{\Delta(h_r^j, o) = \frac{\partial L}{\partial h_r^j}} \frac{\partial h_r^j}{\partial w_{(h_{r-1}^i, h_r^j)}}$$

Error Backpropagation

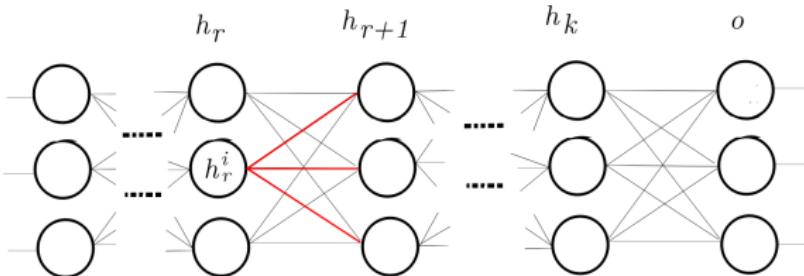


Error Backpropagation



$$\Delta(o, o) = \frac{\partial L}{\partial o}$$

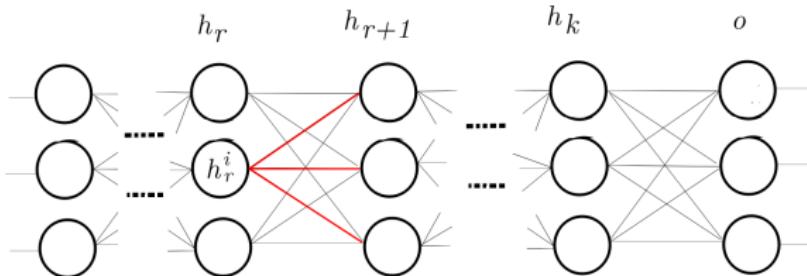
Error Backpropagation



$$\Delta(o, o) = \frac{\partial L}{\partial o}$$

$$\Delta(h_r^i, o) = \frac{\partial L}{\partial h_r^i} = \sum_{h: h_r^i \Rightarrow h_{r+1}} \frac{\partial L}{\partial h} \frac{\partial h}{\partial h_r^i} = \sum_{h: h_r^i \Rightarrow h_{r+1}} \frac{\partial h}{\partial h_r^i} \Delta(h, o)$$

Error Backpropagation



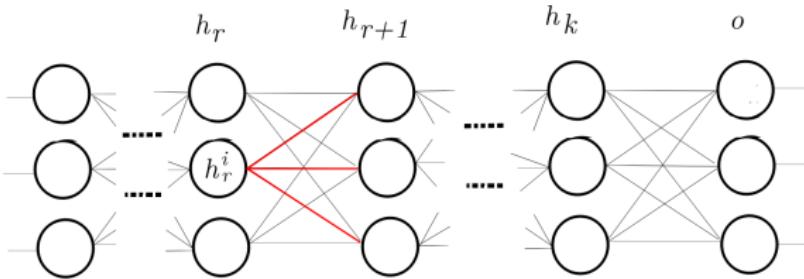
$$a_h = W_h \cdot h_r = \sum_j w_{(h_r^j, h)} h_r^j$$

$$h = \Phi(a_h) = \Phi\left(\sum_j w_{(h_r^j, h)} h_r^j\right)$$

$$\frac{\partial h}{\partial h_r^i} = \frac{\partial h}{\partial a_h} \cdot \frac{\partial a_h}{\partial h_r^i} = \frac{\partial \Phi(a_h)}{\partial a_h} \cdot w_{(h_r^i, h)} = \Phi'(a_h) \cdot w_{(h_r^i, h)}$$

$$\frac{\partial h}{\partial w_{(h_r^i, h)}} = \frac{\partial h}{\partial a_h} \cdot \frac{\partial a_h}{\partial w_{(h_r^i, h)}} = \Phi'(a_h) \cdot h_r^i$$

Error Backpropagation

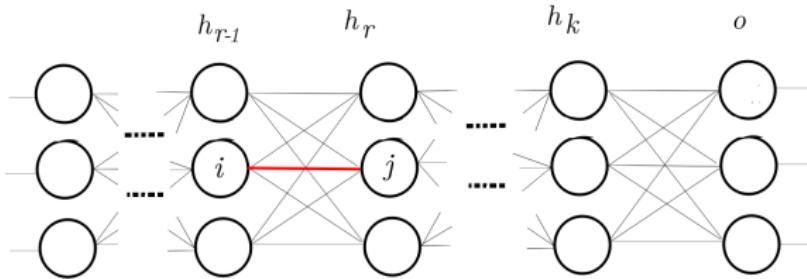


$$\Delta(o, o) = \frac{\partial L}{\partial o}$$

$$\Delta(h_r^i, o) = \frac{\partial L}{\partial h_r^i} = \sum_{h: h_r^i \Rightarrow h_{r+1}} \frac{\partial L}{\partial h} \frac{\partial h}{\partial h_r^i} = \sum_{h: h_r^i \Rightarrow h_{r+1}} \frac{\partial h}{\partial h_r^i} \Delta(h, o)$$

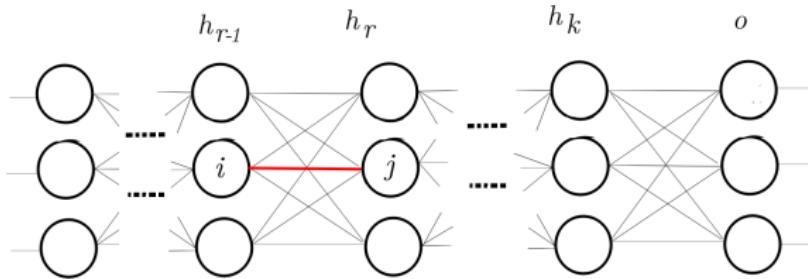
$$\Delta(h_r^i, o) = \sum_{h: h_r^i \Rightarrow h} \Phi'(a_h) \cdot w_{(h_r^i, h)} \cdot \Delta(h, o)$$

Error Backpropagation



$$\frac{\partial L}{\partial w_{(h_{r-1}^i, h_r^j)}} = \underbrace{\left[\sum_{[h_r^j, h_{r+1}, \dots, h_k, o] \in \mathcal{P}} \frac{\partial L}{\partial o} \frac{\partial o}{\partial h_k} \prod_{i=r}^{k-1} \frac{\partial h_{i+1}}{\partial h_i} \right]}_{\Delta(h_r^j, o) = \frac{\partial L}{\partial h_r^j}} \frac{\partial h_r^j}{\partial w_{(h_{r-1}^i, h_r^j)}}$$

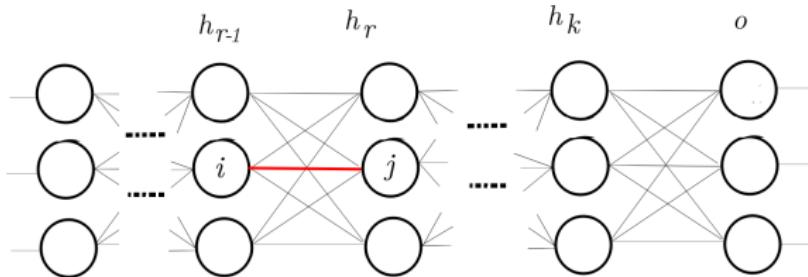
Error Backpropagation



$$\frac{\partial L}{\partial w_{(h_{r-1}^i, h_r^j)}} = \underbrace{\left[\sum_{[h_r^j, h_{r+1}, \dots, h_k, o] \in \mathcal{P}} \frac{\partial L}{\partial o} \frac{\partial o}{\partial h_k} \prod_{i=r}^{k-1} \frac{\partial h_{i+1}}{\partial h_i} \right]}_{\Delta(h_r^j, o) = \frac{\partial L}{\partial h_r^j}} \frac{\partial h_r^j}{\partial w_{(h_{r-1}^i, h_r^j)}}$$

$$\frac{\partial h_r^j}{\partial w_{(h_{r-1}^i, h_r^j)}} = h_{r-1}^i \cdot \Phi'(a_{h_r^j})$$

Error Backpropagation



$$\frac{\partial L}{\partial w_{(h_{r-1}^i, h_r^j)}} = \underbrace{\left[\sum_{[h_r^j, h_{r+1}, \dots, h_k, o] \in \mathcal{P}} \frac{\partial L}{\partial o} \frac{\partial o}{\partial h_k} \prod_{i=r}^{k-1} \frac{\partial h_{i+1}}{\partial h_i} \right]}_{\Delta(h_r^j, o) = \frac{\partial L}{\partial h_r^j}} \frac{\partial h_r^j}{\partial w_{(h_{r-1}^i, h_r^j)}}$$

$$\frac{\partial h_r^j}{\partial w_{(h_{r-1}^i, h_r^j)}} = h_{r-1}^i \cdot \Phi'(a_{h_r^j})$$

$$\frac{\partial L}{\partial w_{(h_{r-1}^i, h_r^j)}} = \Delta(h_r^i, o) \cdot h_{r-1}^i \cdot \Phi'(a_{h_r^i})$$

SGD using Backpropagation

For each training sample:

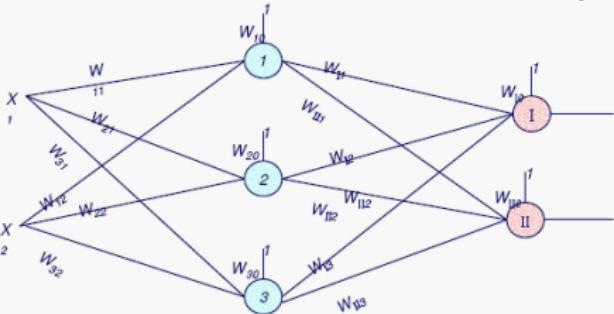
- Compute the forward path.
- Compute $\Delta(o, o)$ for each output neuron.
- Update each connecting weight of the output layer as

$$w(h_k, o) = w(h_k, o) - \eta \cdot \Delta(o, o) \cdot h_k \cdot \Phi'(a_o)$$

- For $r = k, k-1, \dots, 1$:
 - Compute $\Delta(h_r^i, o)$ for the i -th neuron at the r -th hidden layer.
 - Update each connecting weight of the i -th neuron at the r -th hidden layer as:

$$w(h_{r-1}, h_r) = w(h_{r-1}, h_r) - \eta \cdot \Delta(h_r, o) \cdot h_{r-1} \cdot \Phi'(a_{h_r})$$

MLP example



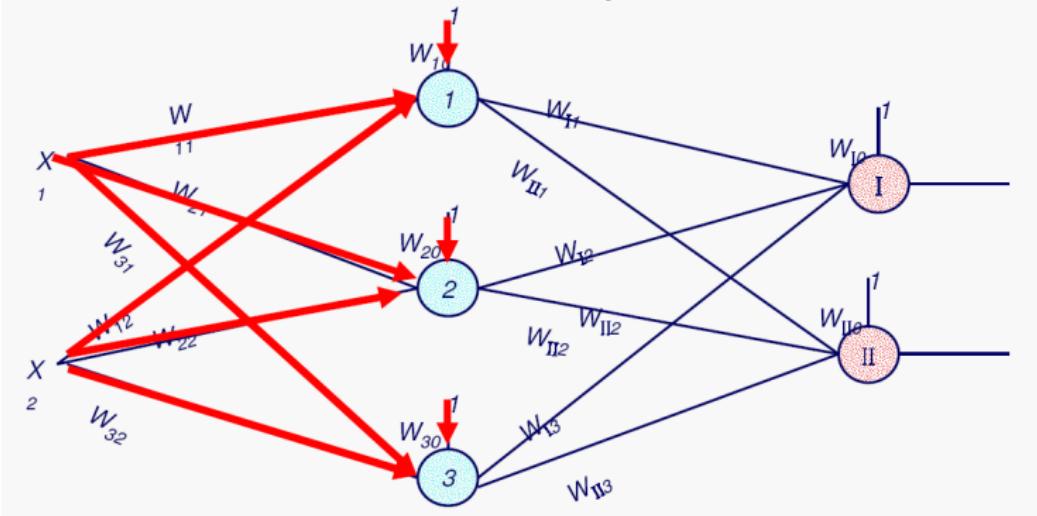
Training Data

<u>Inputs</u>	<u>Target Concept</u>
0,0	0,0
0,1	0,1
1,0	0,1
1,1	1,0

Activation Function for the hidden and output layers:

$$f(y) = \frac{1}{1+e^{-y}}$$

MLP example

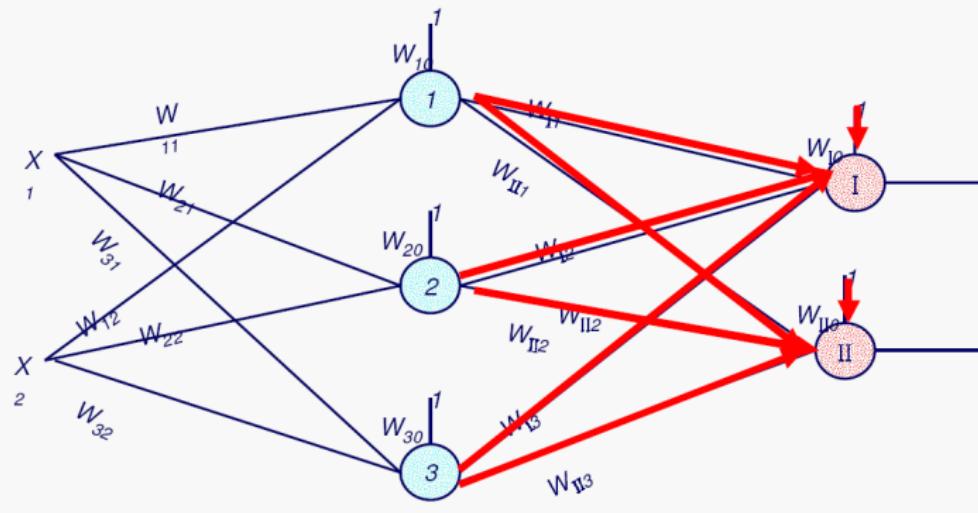


$$O_1 = O\left(\sum_{i=0} W_{1i} X_i\right) = O(W_{10} \cdot 1 + W_{11} X_1 + W_{12} X_2)$$

$$O_2 = O\left(\sum_{i=0} W_{2i} X_i\right) = O(W_{20} \cdot 1 + W_{21} X_1 + W_{22} X_2)$$

$$O_3 = O\left(\sum_{i=0} W_{3i} X_i\right) = O(W_{30} \cdot 1 + W_{31} X_1 + W_{32} X_2)$$

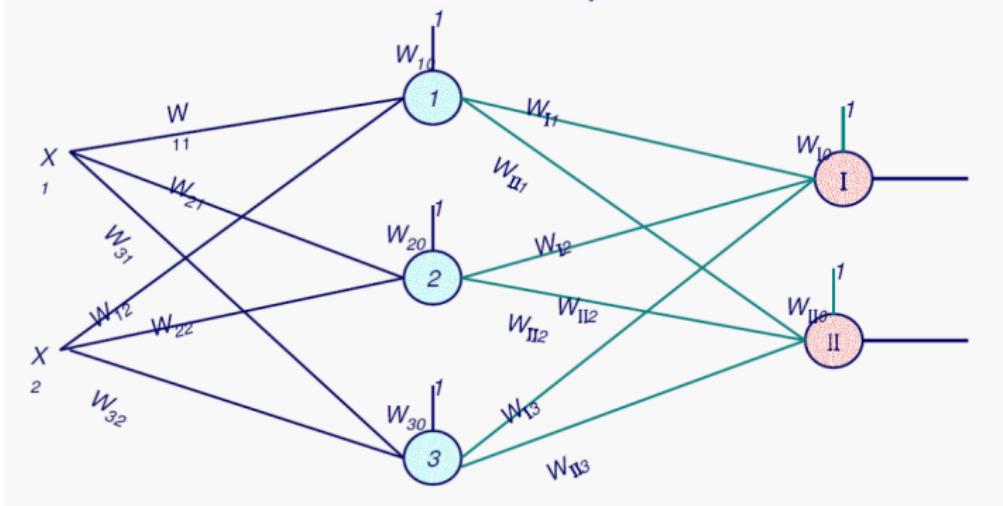
MLP example



$$O_I = O\left(\sum_{i=0} W_{Ii} O_i\right) = O(W_{I0} \cdot 1 + W_{I1} O_1 + W_{I2} O_2 + W_{I3} O_3)$$

$$O_{II} = O\left(\sum_{i=0} W_{IIi} O_i\right) = O(W_{II0} \cdot 1 + W_{II1} O_1 + W_{II2} O_2 + W_{II3} O_3)$$

MLP example

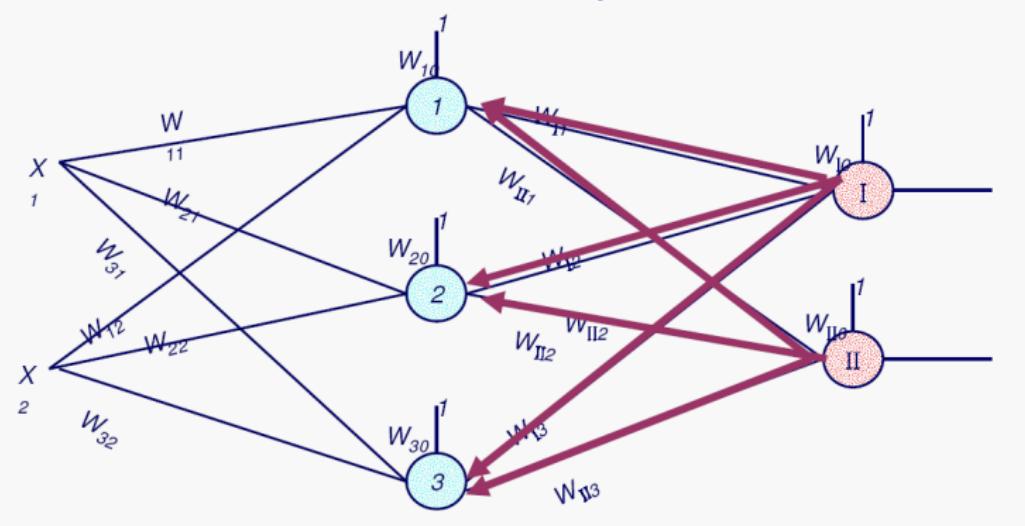


$$\begin{aligned}\Delta_I &= -(t_I - O_I) \\ \Delta_{II} &= -(t_{II} - O_{II})\end{aligned}$$

$$\begin{aligned}W_{I0} &= W_{I0} - \eta \Delta_I O_I (1 - O_I) \\ W_{I1} &= W_{I1} - \eta \Delta_I O_1 O_I (1 - O_I) \\ W_{I2} &= W_{I2} - \eta \Delta_I O_2 O_I (1 - O_I) \\ W_{I3} &= W_{I3} - \eta \Delta_I O_3 O_I (1 - O_I)\end{aligned}$$

$$\begin{aligned}W_{II0} &= W_{II0} - \eta \Delta_{II} O_{II} (1 - O_{II}) \\ W_{II1} &= W_{II1} - \eta \Delta_{II} O_1 O_{II} (1 - O_{II}) \\ W_{II2} &= W_{II2} - \eta \Delta_{II} O_2 O_{II} (1 - O_{II}) \\ W_{II3} &= W_{II3} - \eta \Delta_{II} O_3 O_{II} (1 - O_{II})\end{aligned}$$

MLP example

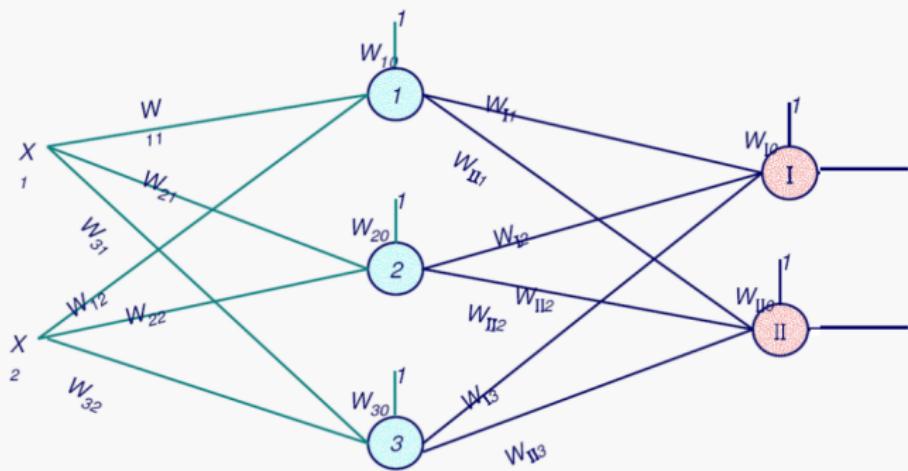


$$\Delta_1 = \sum_{k=1} W_{k1} \Delta_k O_k (1 - O_k) = (W_{I1} \Delta_I O_I (1 - O_I) + W_{II1} \Delta_{II} O_{II} (1 - O_{II}))$$

$$\Delta_2 = \sum_{k=1} W_{k2} \Delta_k O_k (1 - O_k) = (W_{I2} \Delta_I O_I (1 - O_I) + W_{II2} \Delta_{II} O_{II} (1 - O_{II}))$$

$$\Delta_3 = \sum_{k=1} W_{k3} \Delta_k O_k (1 - O_k) = (W_{I3} \Delta_I O_I (1 - O_I) + W_{II3} \Delta_{II} O_{II} (1 - O_{II}))$$

MLP example



$$W_{10} = W_{10} - \eta \Delta_1 O_1 (1 - O_1)$$

$$W_{11} = W_{11} - \eta \Delta_1 X_1 O_1 (1 - O_1)$$

$$W_{12} = W_{12} - \eta \Delta_1 X_2 O_1 (1 - O_1)$$

$$W_{20} = W_{20} - \eta \Delta_2 O_2 (1 - O_2)$$

$$W_{21} = W_{21} - \eta \Delta_2 X_1 O_2 (1 - O_2)$$

$$W_{22} = W_{22} - \eta \Delta_2 X_2 O_2 (1 - O_2)$$

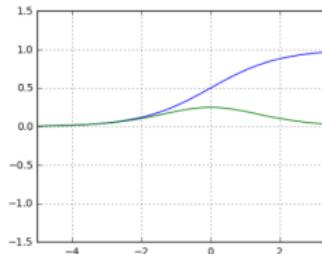
$$W_{30} = W_{30} - \eta \Delta_3 O_3 (1 - O_3)$$

$$W_{31} = W_{31} - \eta \Delta_3 X_1 O_3 (1 - O_1)$$

$$W_{32} = W_{32} - \eta \Delta_3 X_2 O_3 (1 - O_3)$$

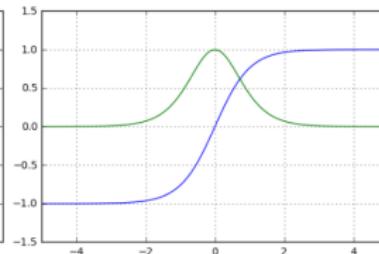
learning slowdown issue for sigmoid and tanh

The derivation of sigmoid and tanh is near zero for small and large pre-activations. This slows down the learning speed.



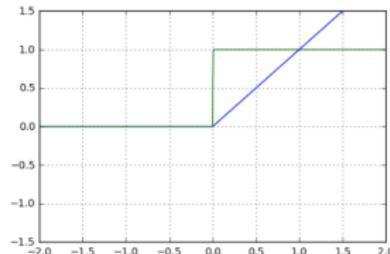
$$\text{sigm}(x) = \frac{1}{1 + e^{-x}}$$

$$\text{sigm}'(x) = \text{sigm}(x)(1 - \text{sigm}(x))$$



$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

$$\tanh'(x) = 1 - \tanh(x)^2$$



$$\text{relu}(x) = \max(0, x)$$

$$\text{relu}'(x) = 1_{x>0}$$

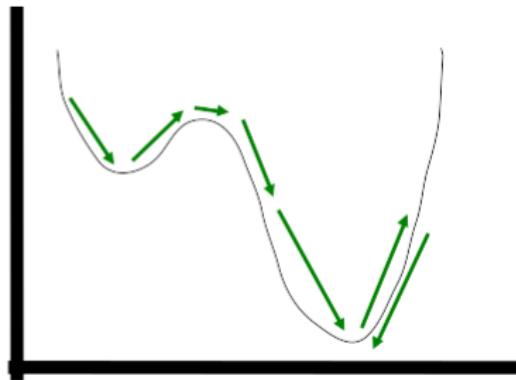
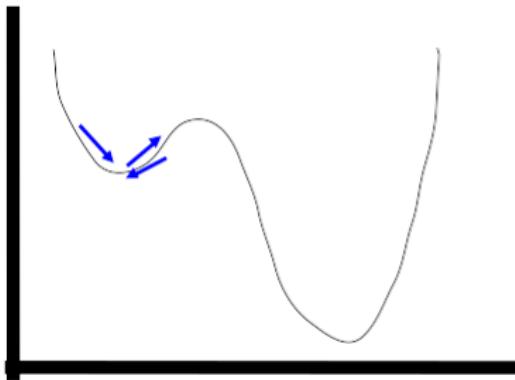
Solutions:

- Use ReLU activation function.
- Do weight and input normalization.

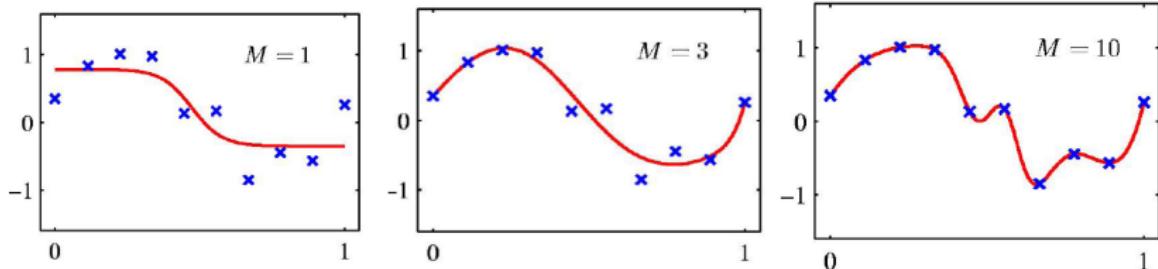
Momentum

Add a fraction (α =momentum) of the last change to the current change:

$$\Delta w_{ij} = \eta \cdot \Delta_i \cdot x_j \cdot \Phi'(a_i) + \alpha \cdot \Delta w_{ij}$$

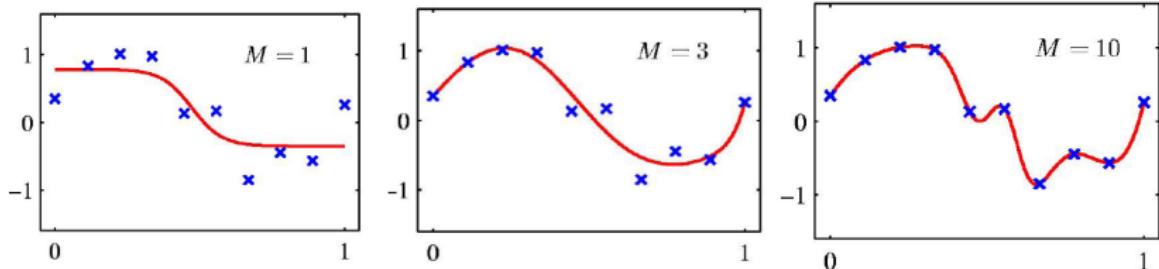


Overfitting



- **Generalization:** To establish a balance between correct responses for the training patterns and unseen new patterns.
- **Memorization:** When the model memorizes training samples instead of learning the descriptive common patterns.
- **Overfitting:** Weak generalization. It happens when the network complexity is more than the problem complexity.

Overfitting



- **Generalization:** To establish a balance between correct responses for the training patterns and unseen new patterns.
- **Memorization:** When the model memorizes training samples instead of learning the descriptive common patterns.
- **Overfitting:** Weak generalization. It happens when the network complexity is more than the problem complexity.
- In neural networks, the model complexity is specified by the number of neurons and weights.

Overfitting

- Overfitting: fitting a model to a particular training data set does not guarantee good prediction performance on unseen test data.
- If you train the net for too long, then you run the risk of overfitting.

Overfitting

- Overfitting: fitting a model to a particular training data set does not guarantee good prediction performance on unseen test data.
- If you train the net for too long, then you run the risk of overfitting.
- Lack of sufficient training data increases the risk of overfitting.

Overfitting

- Overfitting: fitting a model to a particular training data set does not guarantee good prediction performance on unseen test data.
- If you train the net for too long, then you run the risk of overfitting.
- Lack of sufficient training data increases the risk of overfitting.
- Example: consider a single neuron with 5 inputs and the following training set:

x_1	x_2	x_3	x_4	x_5	y
1	1	0	0	0	2
2	0	1	0	0	4
3	0	0	1	0	6
4	0	0	0	1	8

$$\hat{y} = \sum_{i=1}^5 w_i \cdot x_i$$

$$\overline{W} = [2, 0, 0, 0, 0] \quad \text{Correct}$$

$$[0, 2, 4, 6, 8] \quad \text{Overfitted}$$

Avoid Overfitting

- One possible approach is to reduce the size of the network.
 - However, large networks have the potential to be more powerful than small networks.

Avoid Overfitting

- One possible approach is to reduce the size of the network.
 - However, large networks have the potential to be more powerful than small networks.
- Provide more training samples (not always possible).

Avoid Overfitting

- One possible approach is to reduce the size of the network.
 - However, large networks have the potential to be more powerful than small networks.
- Provide more training samples (not always possible).
- Stop learning before overfitting happens.

Avoid Overfitting

- One possible approach is to reduce the size of the network.
 - However, large networks have the potential to be more powerful than small networks.
- Provide more training samples (not always possible).
- Stop learning before overfitting happens.
- Use regularization terms to dynamically adjust network complexity.

Avoid Overfitting

- One possible approach is to reduce the size of the network.
 - However, large networks have the potential to be more powerful than small networks.
- Provide more training samples (not always possible).
- Stop learning before overfitting happens.
- Use regularization terms to dynamically adjust network complexity.
- Use ensemble methods.

Avoid Overfitting

- One possible approach is to reduce the size of the network.
 - However, large networks have the potential to be more powerful than small networks.
- Provide more training samples (not always possible).
- Stop learning before overfitting happens.
- Use regularization terms to dynamically adjust network complexity.
- Use ensemble methods.
- use random dropout technique for hidden neurons.

Regularization

- Since a larger number of parameters causes overfitting, a natural approach is to constrain the model to use fewer non-zero parameters.
- The most applied regularization is adding the penalty $\lambda||W||$ to the loss function .

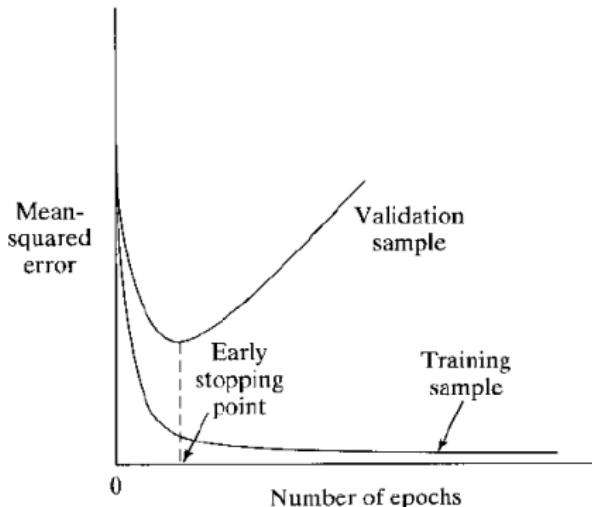
$$L = \frac{1}{2}(Y - y)^2 + \lambda||W||$$

- Therefore the learning rule is re-written as:

$$\Delta w_{ij} = \eta \cdot \Delta_i \cdot x_j \cdot \Phi'(a_i) - \eta \cdot \lambda \cdot w_{ij}$$

Early stopping method of training

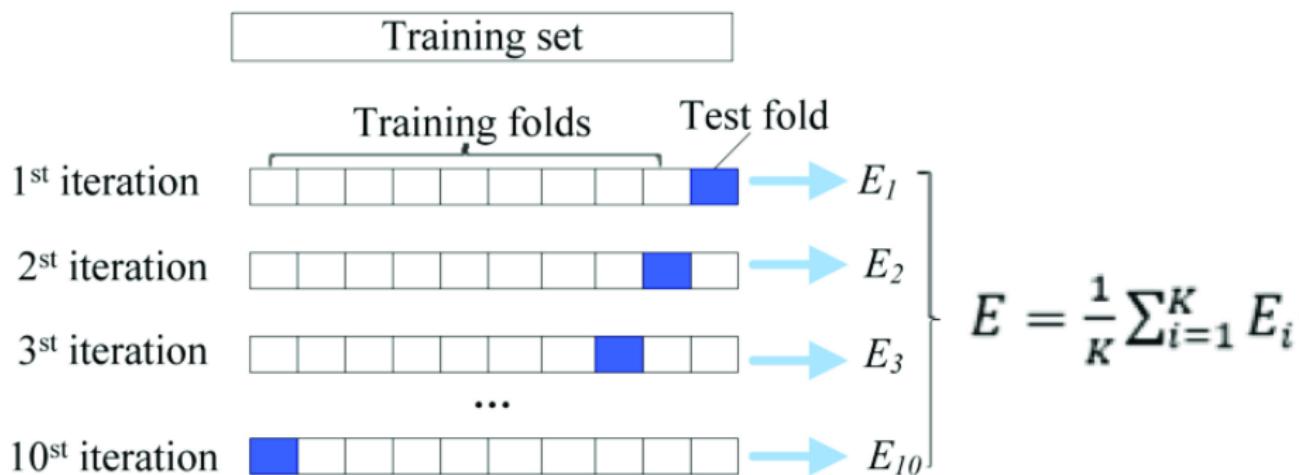
- Split training samples into a training set (80%) and a validation set (20%).



- Stop learning when the loss decreases on train set but increases on validation set.

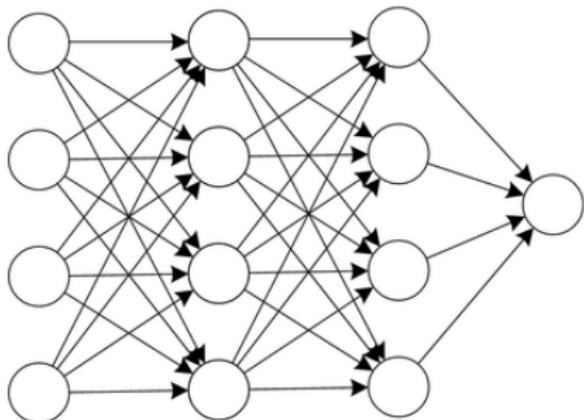
K-fold Cross-validation

- Split data into K disjoint folds.
- Iteratively, use nine folds for training and the remaining fold for validation.
- K-fold cross-validation is mainly used for tuning the meta-params such as learning rate, network size, ...
- K-fold cross-validation is also used for reporting performance on unsplited (train/test) data

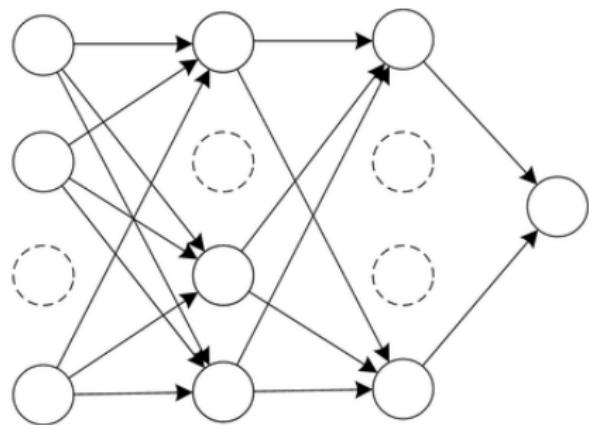


Dropout or Dropconnect

- Temporarily inactivate some of the hidden layer neurons (Dropout) or some of their weights(Dropconnect).
- These methods act like training multiple network inside one network (to decrease overfitting and increase generalization)



(a) Standard Neural Network

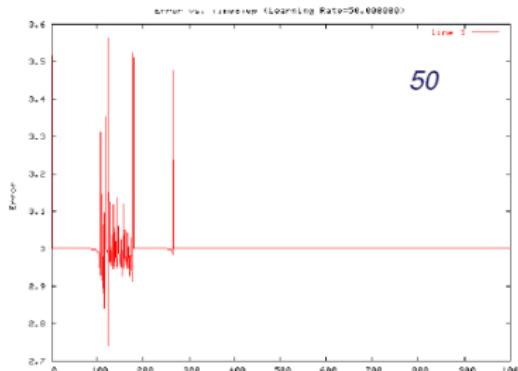
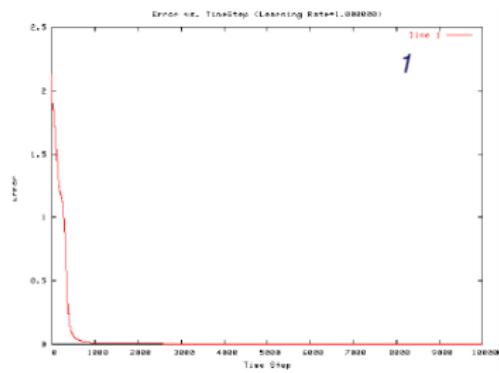
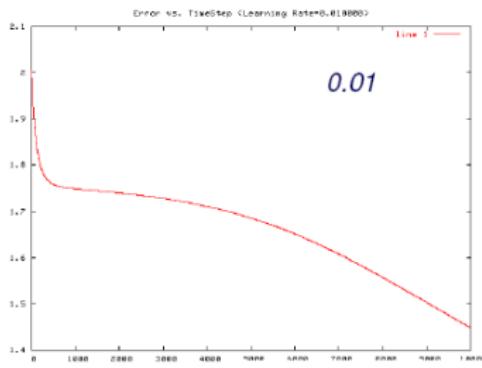


(b) Network after Dropout

The Vanishing and Exploding Gradient

- While increasing depth often reduces the number of parameters, the chain rule can cause in vanishing or exploding gradients.
- A sigmoid activation often encourages the vanishing gradient problem, because its derivative is less than 0.25.
- A ReLU activation unit is known to be less likely to create a vanishing gradient problem because its derivative is always 1 for positive values of the argument.
- **Adaptive learning rates** and **conjugate gradient methods** are other solutions.

Importance of Learning Rate



Learning mode

There are two basic modes of updating weights:

- The **pattern mode** in which weights are updated after the presentation of a single training pattern,
 - It is easier to fit into memory.
 - For larger datasets it can converge faster.
 - Due to frequent updates the steps taken towards the minima of the loss function have oscillations which can help getting out of local minimums.
- The **batch mode** in which weights are updated after a batch of patterns.
 - Less oscillations and noisy steps taken towards the global minima.
 - It can benefit from vectorization which increases the speed of processing.
 - It produces a more stable gradient descent convergence.