

Introduction to *shallow* Neural Networks (Multi-Layer Perceptron, MLP)

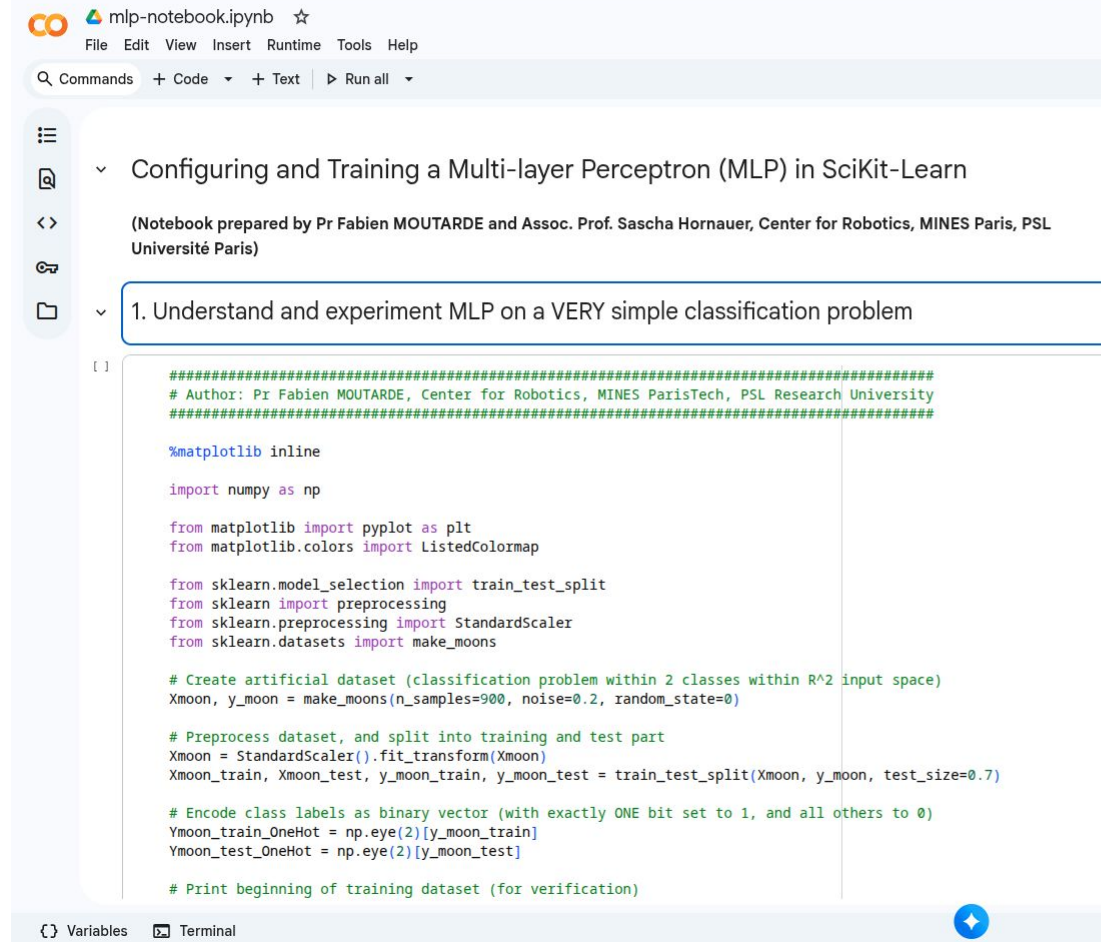
Assoc. Prof. Sascha Hornauer

Center for Robotics
Mines Paris
PSL Université

sascha.hornauer@minesparis.psl.eu

For the Practical

- Install Anaconda
- Install Jupyter Lab
- Install scikit-learn, matplotlib



The screenshot shows a JupyterLab interface with a notebook titled 'mlp-notebook.ipynb'. The notebook has a menu bar with 'File', 'Edit', 'View', 'Insert', 'Runtime', 'Tools', and 'Help'. Below the menu bar is a search bar and tabs for 'Commands', 'Code', 'Text', and 'Run all'. The notebook content is organized into sections: 'Configuring and Training a Multi-layer Perceptron (MLP) in SciKit-Learn' and '1. Understand and experiment MLP on a VERY simple classification problem'. The first cell of code is selected, showing the following Python code:

```
#####  
# Author: Pr Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL Research University  
#####  
  
%matplotlib inline  
  
import numpy as np  
  
from matplotlib import pyplot as plt  
from matplotlib.colors import ListedColormap  
  
from sklearn.model_selection import train_test_split  
from sklearn import preprocessing  
from sklearn.preprocessing import StandardScaler  
from sklearn.datasets import make_moons  
  
# Create artificial dataset (classification problem within 2 classes within R^2 input space)  
Xmoon, y_moon = make_moons(n_samples=900, noise=0.2, random_state=0)  
  
# Preprocess dataset, and split into training and test part  
Xmoon = StandardScaler().fit_transform(Xmoon)  
Xmoon_train, Xmoon_test, y_moon_train, y_moon_test = train_test_split(Xmoon, y_moon, test_size=0.7)  
  
# Encode class labels as binary vector (with exactly ONE bit set to 1, and all others to 0)  
Ymoon_train_OneHot = np.eye(2)[y_moon_train]  
Ymoon_test_OneHot = np.eye(2)[y_moon_test]  
  
# Print beginning of training dataset (for verification)
```

Neural Networks: from biology to engineering

- **Understanding and modelling of brain**



- **Imitation to reproduce high-level functions**



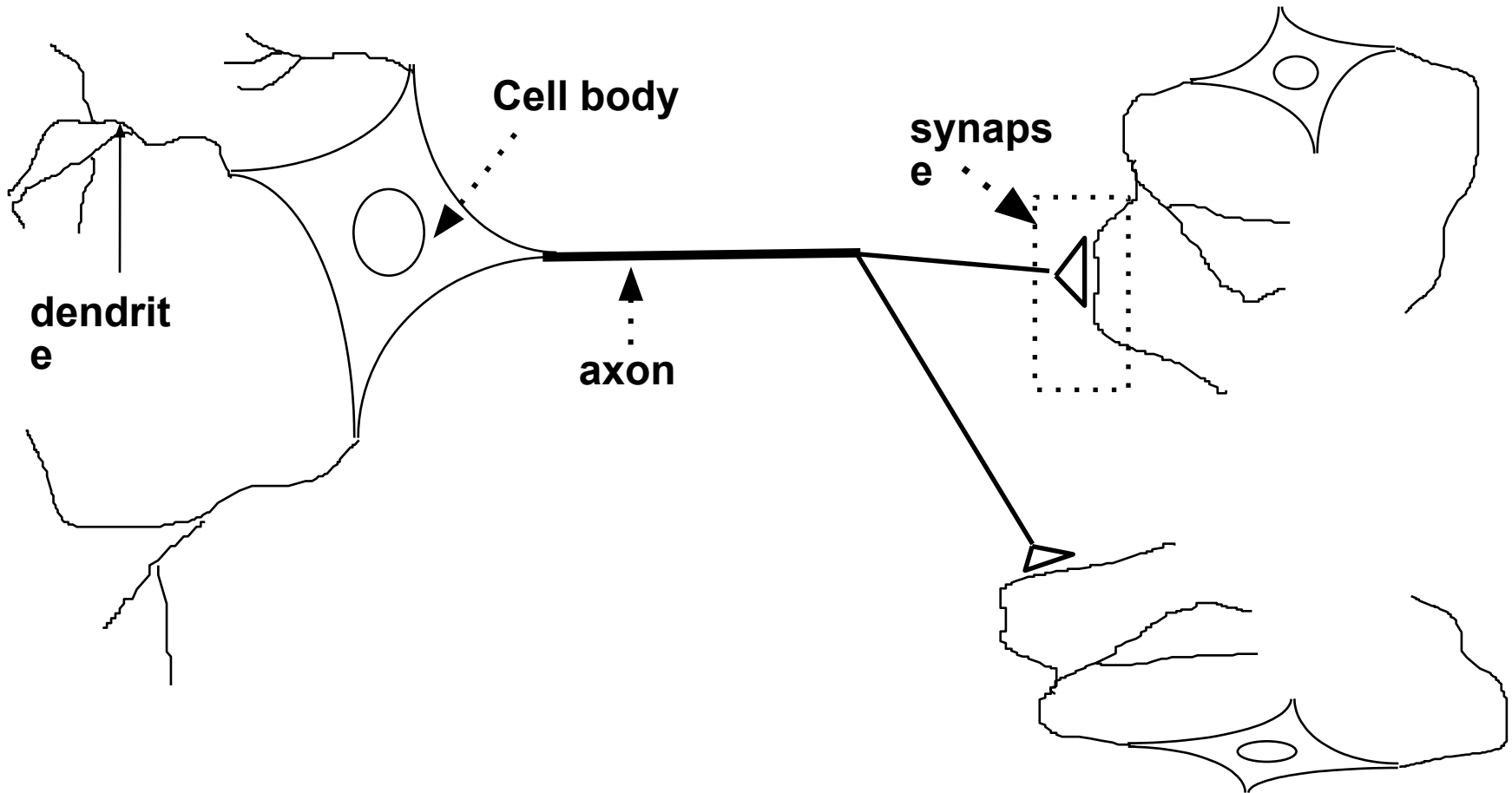
- **Mathematical tool for engineers**

Modelling any input-output function by “learning” from examples:

- Pattern recognition
- Voice recognition
- Classification, diagnosis

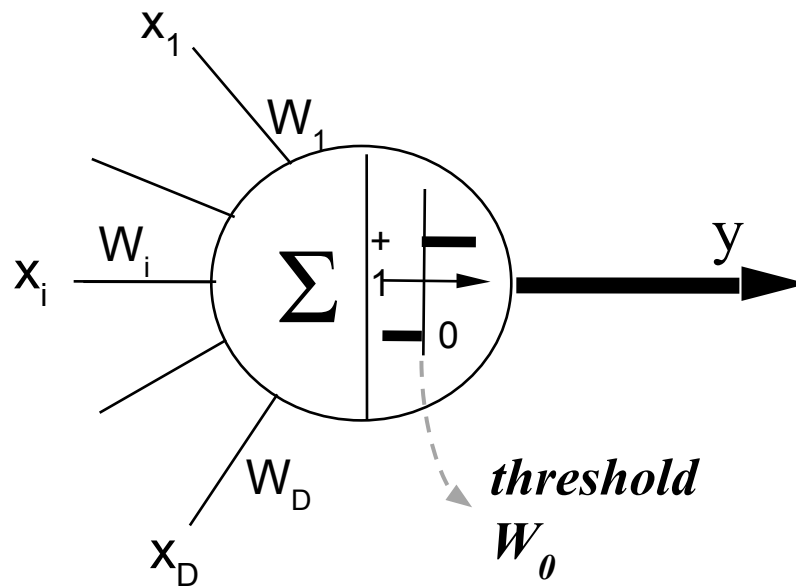
- Identification
- Forecasting
- Control, regulation

Biological neurons

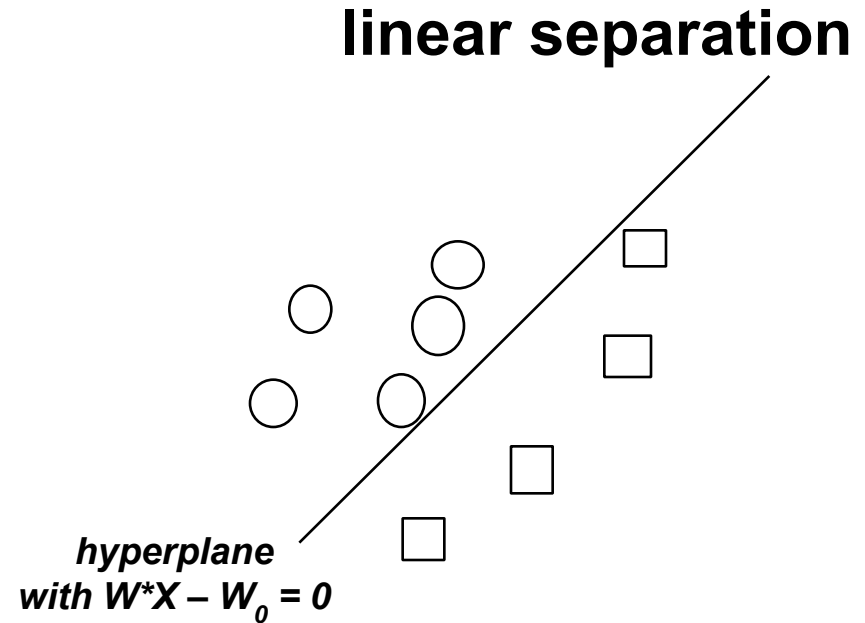
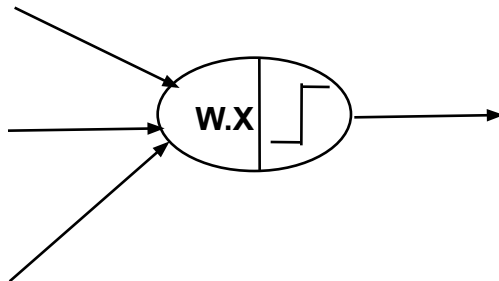


- **Electric signal: dendrites** ☐ **cell body** ☐ **> axon** ☐ **synapses**

- Mc Culloch & Pitts (1943)
 - Simple model of neuron
 - goal: model the brain

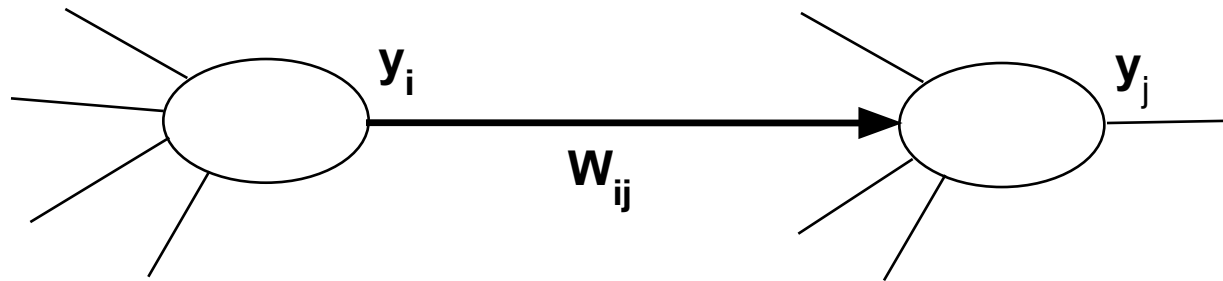


Linear separation by a single neuron



- Hebb rule (1949)

”*Cells that fire together wire together*”, ie synaptic weight increases between neurons that activate simultaneously



$$W_{ij}(t + dt) = W_{ij}(t) + \lambda y_i(t) y_j(t)$$

First formal Neural Networks *(en français : Réseaux de Neurones)*

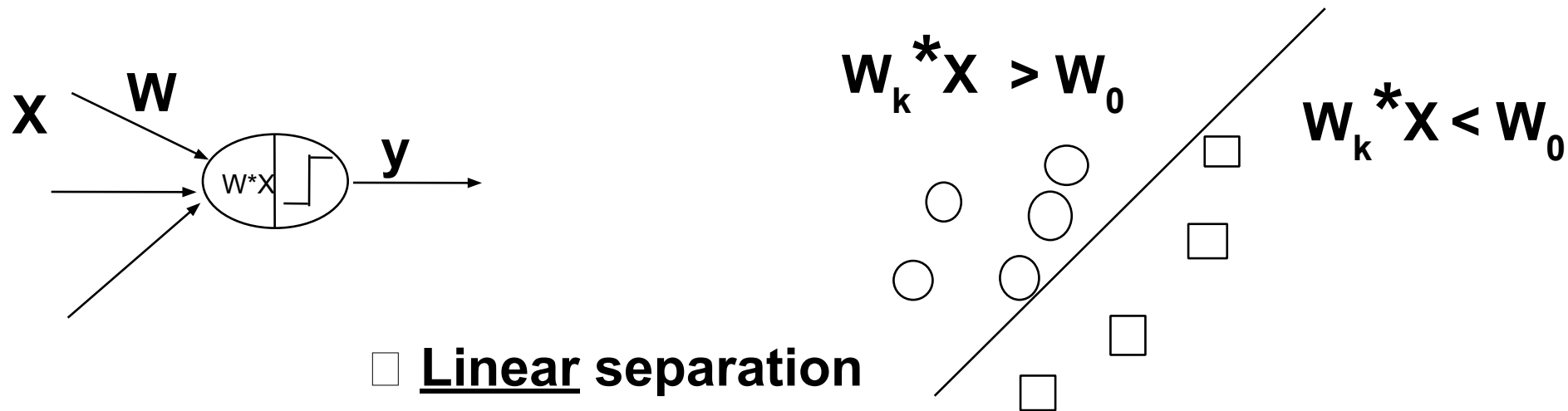
- **PERCEPTRON (Rosenblatt, 1957)**
- **ADALINE (Widrow, 1962)**

**Formal neuron of Mac Culloch & Pitts
+
Hebb rule for learning**



**Possible to “learn” Boolean functions
by training from examples**

Training of Perceptron



Training algorithm:

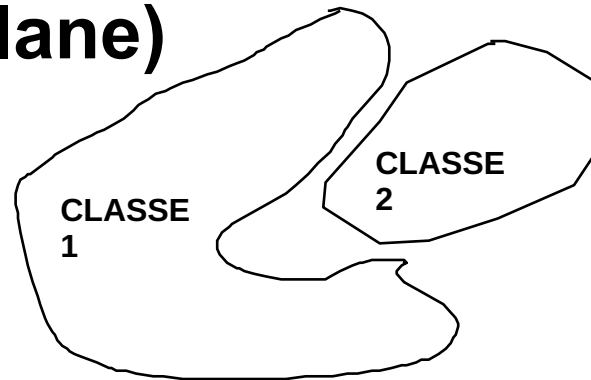
$$\begin{aligned} W_{k+1} &= W_k + vX && \text{if } X \text{ incorrectly classified (v: target value)} \\ W_{k+1} &= W_k && \text{if } X \text{ correctly classified} \end{aligned}$$

- **Convergence if linearly-separable problem**
- **?? if NOT linearly-separable**

Limits of first models, necessity of “hidden” neurons

- **PERCEPTRONS**, book by Minsky & Papert (1969)
Detailed study on Perceptrons and their intrinsic limits:
 - can NOT learn some types of Boolean functions (even simple one like XOR)
 - can do ONLY LINEAR separations

**But many classes cannot be linearly-separated
(by a single hyper-plane)**



- **Necessity of several layers in the Neural Network**
- **Requires new training algorithm**

USE OF DERIVABLE NEURONS

+

APPLY GRADIENT DESCENT METHOD



- **GRADIENT BACK-PROPAGATION** (Rumelhart 1986, Le Cun 85)
(en français : *Rétro-propagation du gradient*)
 - Can train hidden neurons
- Empirical solutions for MANY real-world applications
- Some strong theoretical results:
Multi-Layer Perceptrons are **UNIVERSAL**
(and parsimonious) approximators
- around years 2000's: still used, but much less popular than SVMs and boosting

2nd recent « revival »: Deep-Learning

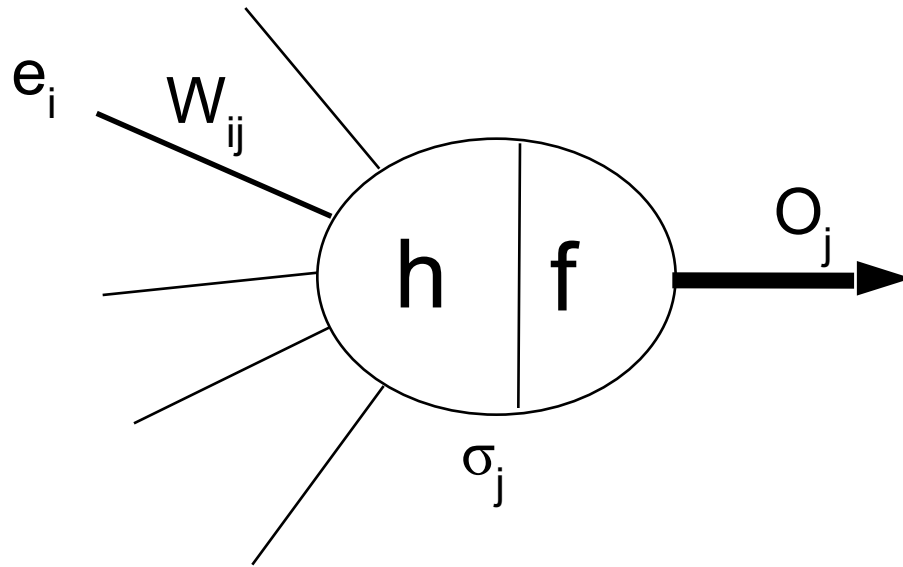
- Since ~2006, rising interest for, and excellent results with "deep" neural networks, consisting in MANY layers:
 - Unsupervised "intelligent" initialization of weights
 - Standard gradient descent, and/or fine-tuning from initial values of weights
 - Hidden layers □ learnt hierarchy of features
- In particular, since ~2013 dramatic progresses in visual recognition (and voice recognition), with deep Convolutional Neural Networks

DEFINITIONS OF FORMAL NEURONS

In general: a processing “unit” applying a simple operation to its inputs, and which can be “connected” to others to build a networks able to realize any input-output function

“Usual” definition: a “unit” computing a weighted sum of its inputs, and then applying some non-linearity (sigmoid, ReLU, ...)

General formal neuron



e_i : inputs of neuron
 σ_j : potential of neuron
 O_j : output of neuron

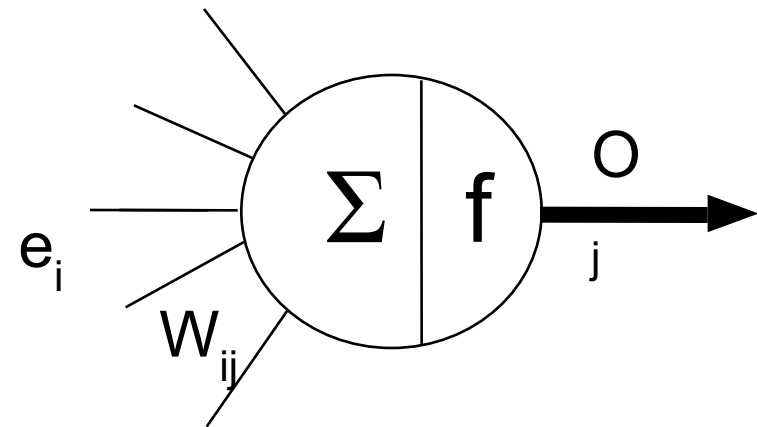
W_{ij} : (synaptic) weights
 h : input function (computation of potential = Σ , dist, kernel, ...)
 f : activation (or transfer) function

$$\sigma_j = h(e_i, \{W_{ij}, i=0 \text{ à } k_j\})$$

$$O_j = f(\sigma_j)$$

The combination of particular h and f functions defines the *type* of formal neuron

PRINCIPLE



$$O_j = f\left(W_{0j} + \sum_{i=1}^{n_j} W_{ij}e_i\right)$$

W_{0j} = "bias"

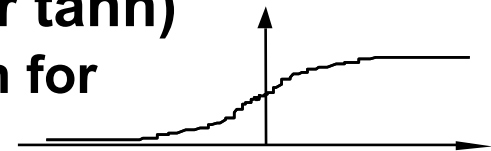
ACTIVATION FUNCTIONS

- **Threshold (Heaviside or sign)**

□ binary neurons

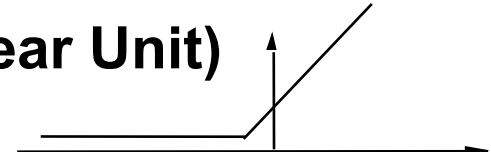
- **Sigmoid (logistic or tanh)**

□ most common for MLPs

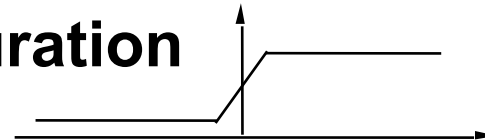


- **Identity** □ linear neurons

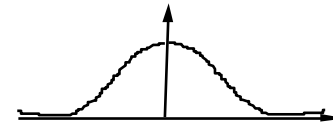
- **ReLU (Rectified Linear Unit)**



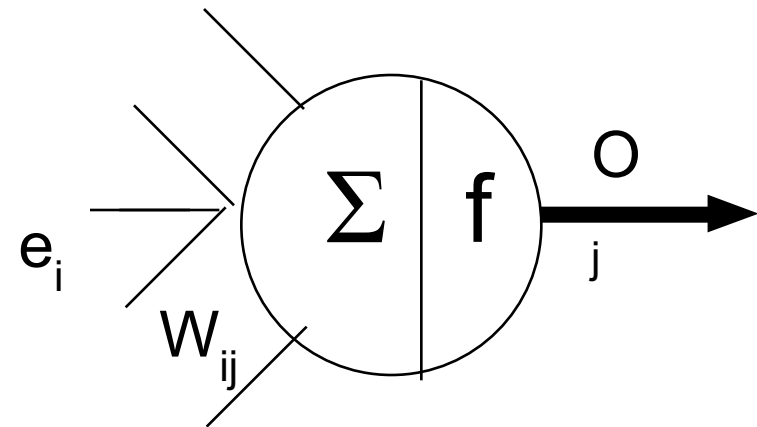
- **Saturation**



- **Gaussian**



PRINCIPLE



$$O_j = f\left(W_{0j} + \sum_{i=1}^{n_j} W_{ij}e_i\right)$$

W_{0j} = "bias"

ACTIVATION FUNCTIONS

- Threshold (Heaviside or sign)

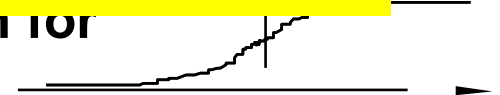
□ binary neurons

Named after
Olivier Heaviside!

- Sigmoid

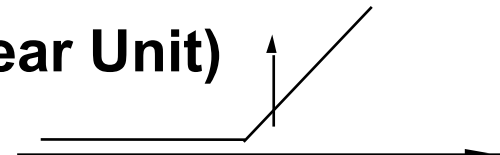
□ most common for

MLPs

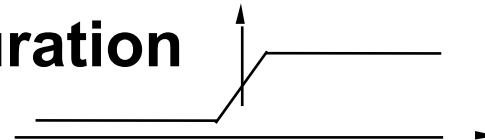


- Identity □ linear neurons

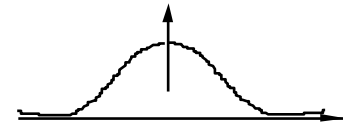
- ReLU (Rectified Linear Unit)



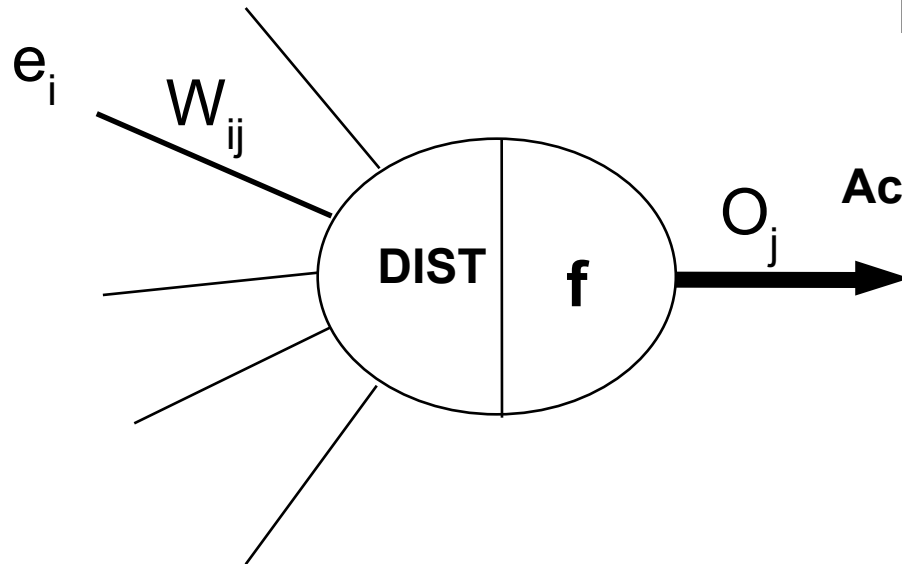
- Saturation



- Gaussian



“Distance” formal neurons



Input function:
$$h\left(\begin{pmatrix} e_1 \\ \dots \\ e_k \end{pmatrix}, \begin{pmatrix} W_{1j} \\ \dots \\ W_{kj} \end{pmatrix}\right) = \sum_i (e_i - W_{ij})^2$$

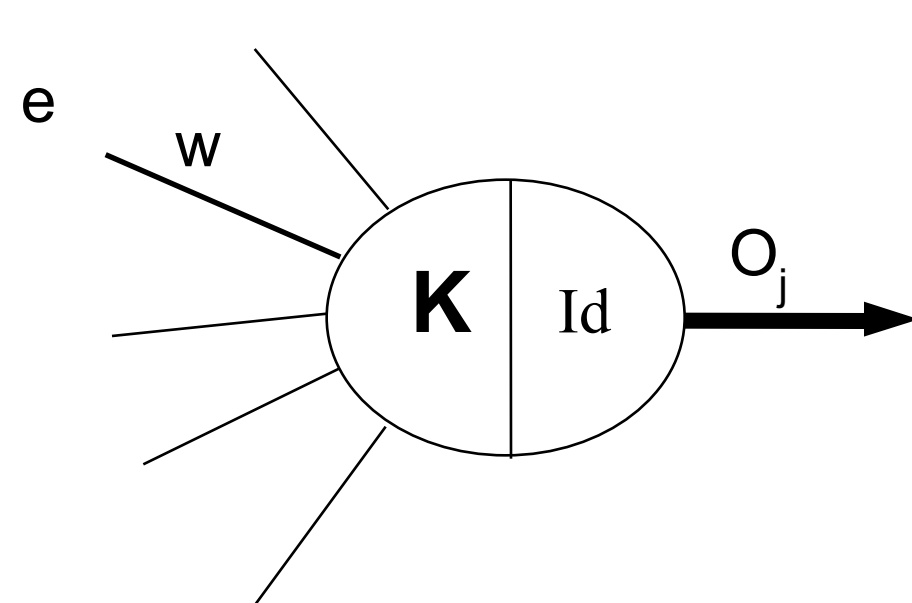
Activation function = Identity or Gaussian

$$f(x) = \|x - w\|$$

$$f(x) = \exp\left(-\frac{\|x - w\|^2}{2\sigma^2}\right)$$

The potential of these neurons is the Euclidian DISTANCE between input vector $(e_i)_i$ and weight vector $(W_{ij})_i$

Kernel-type formal neurons



Input function: $h(e, w) = K(e, w)$

- Constraints on K:
 - K symmetric and "positive" in Mercer sense

Activation function = Identity

Examples of possible kernels:

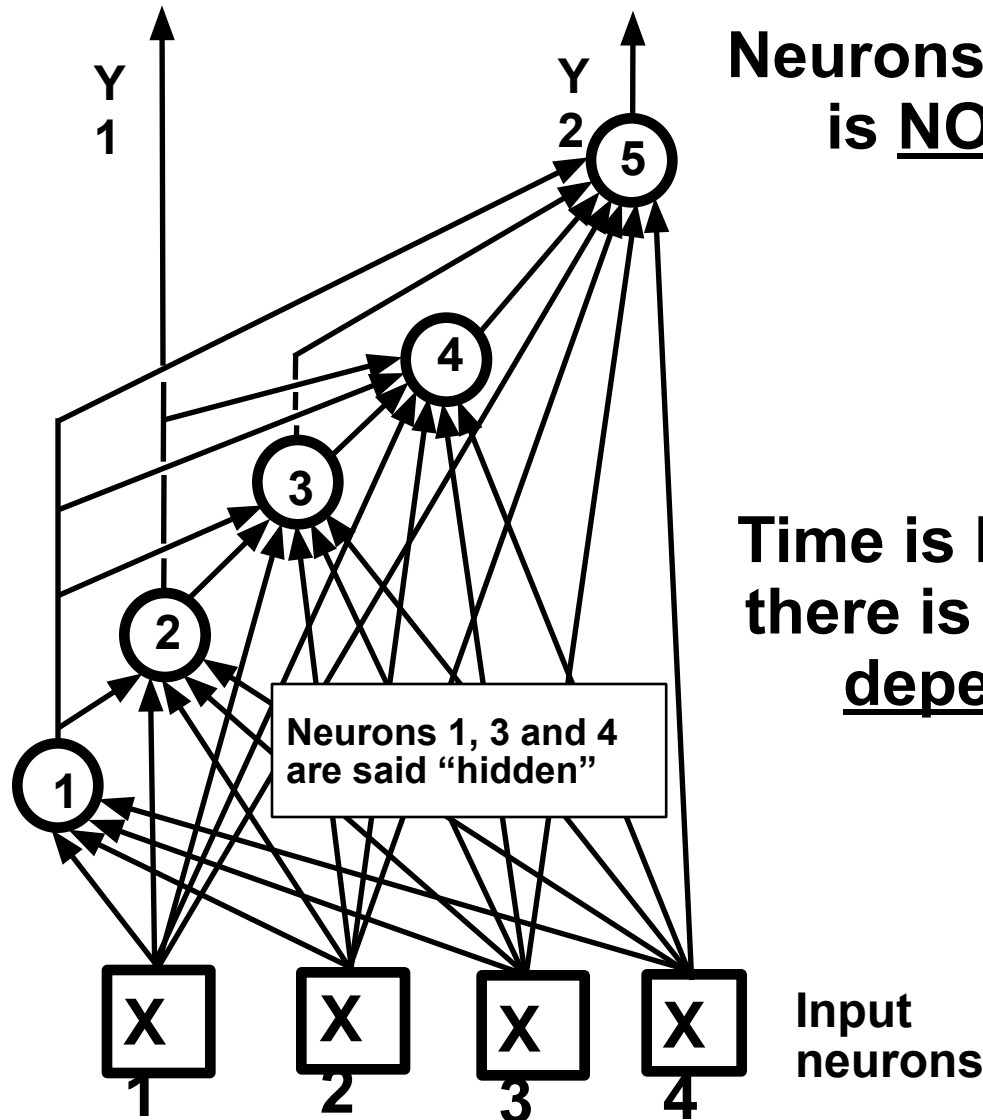
- Polynomial: $K(u, v) = [u * v + 1]^p$
- Radial Basis Function: $K(u, v) = \exp(-||u - v||^2 / 2\sigma^2)$
 - *equivalent to distance-neuron+gaussian-activation*
- Sigmoid: $K(u, v) = \tanh(u * v + \theta)$
 - *equivalent to summing-neurons+sigmoid-activation*

TWO FAMILIES OF NETWORKS

- **FEED-FORWARD NETWORKS**
(en français, “*réseaux non bouclés*”):
NO feedback connection,
The output depends only on current input (NO memory)
- **FEEDBACK OR RECURRENT NETWORKS**
(en français, “*réseaux bouclés*”):
Some internal feedback/backwards connection
☐ output depends on current input
AND ON ALL PREVIOUS INPUTS (some memory inside!)

Feed-forward networks

(en français : réseau “NON-bouclé”)

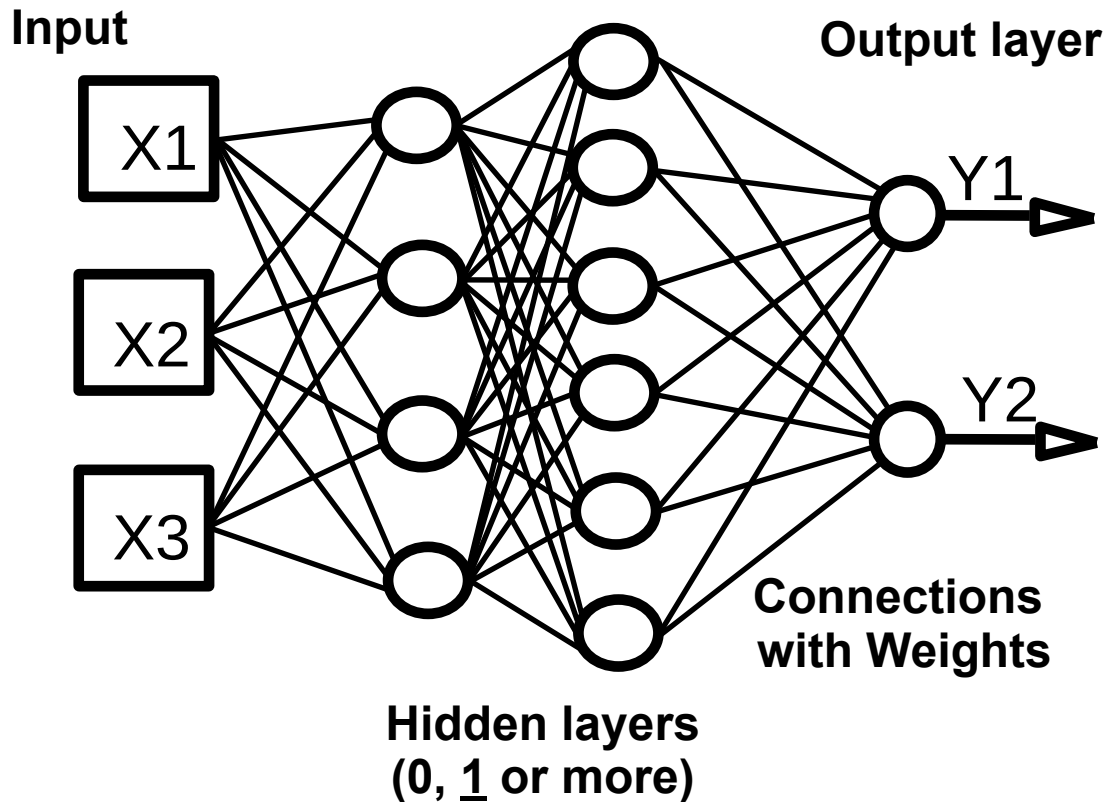


Neurons can be ordered so that there is NO “backwards” connection



Time is NOT a functional variable, i.e. there is NO MEMORY, and the output depends only on current input

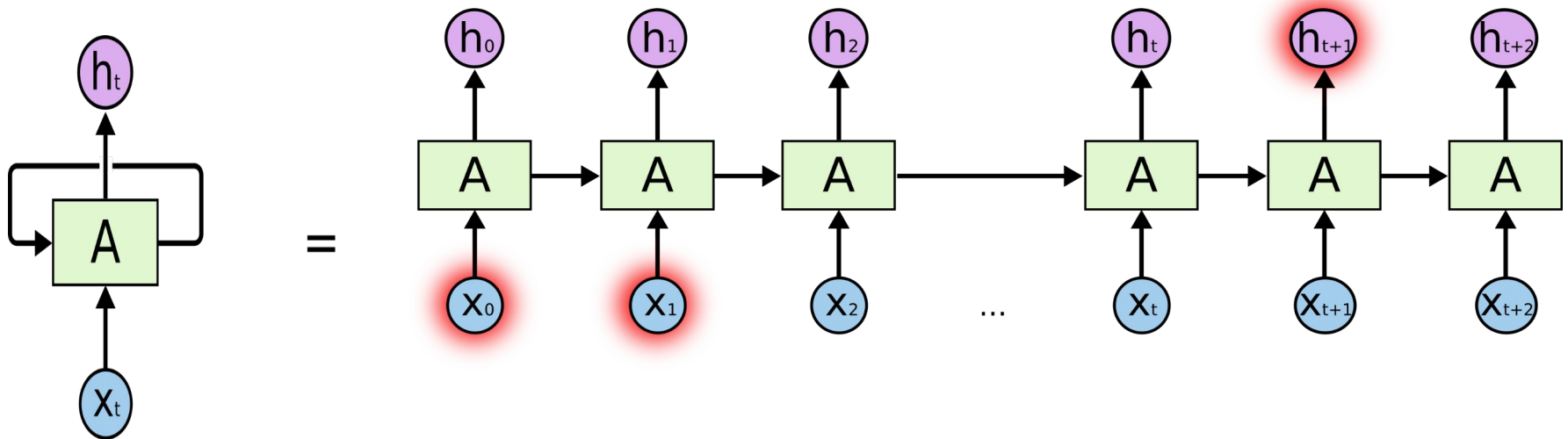
Feed-forward Multi-layer Neural Networks



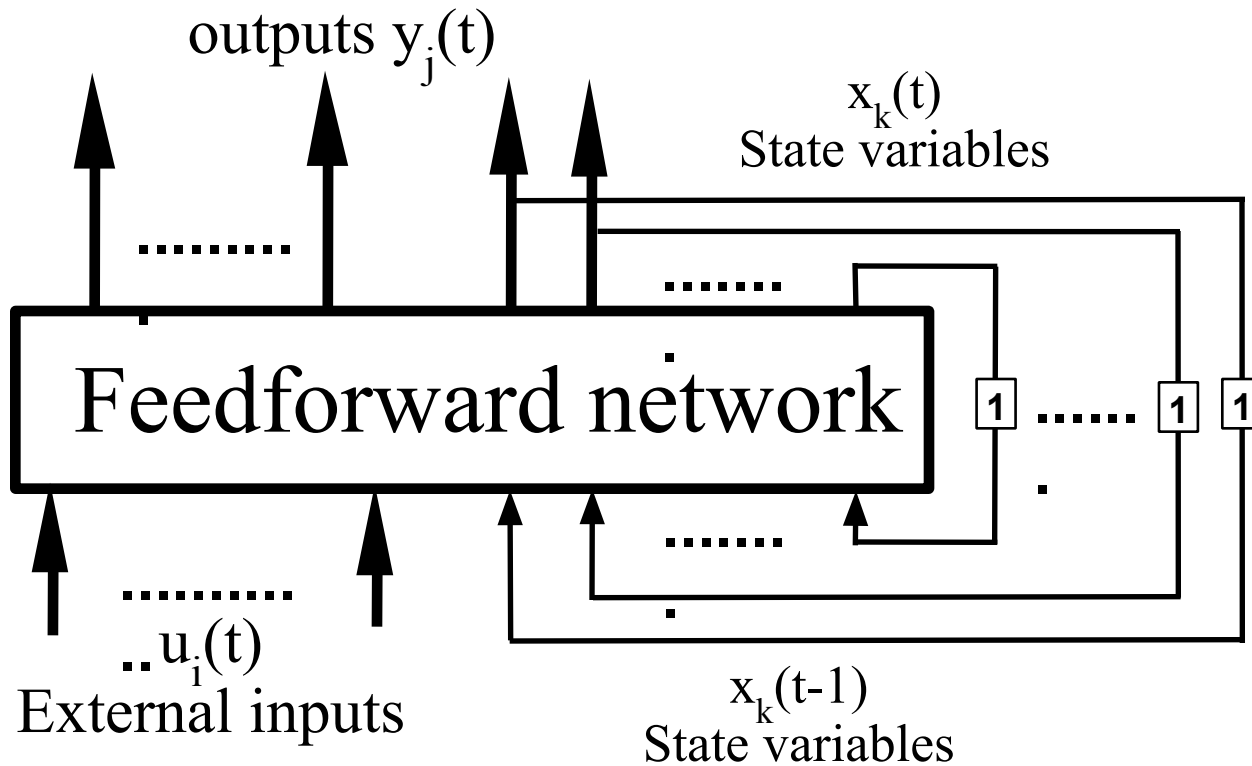
For “Multi-Layer Perceptron” (MLP),
neurons type generally “summing with sigmoid activation”
[terme français pour MLP : “Réseau Neuronal à couches”]

Recurrent Neural Networks

- Input is new data **and** the output from a previous calculation
- Unrolled, these *Recurrent Neural Networks* can have problem with long-term dependencies



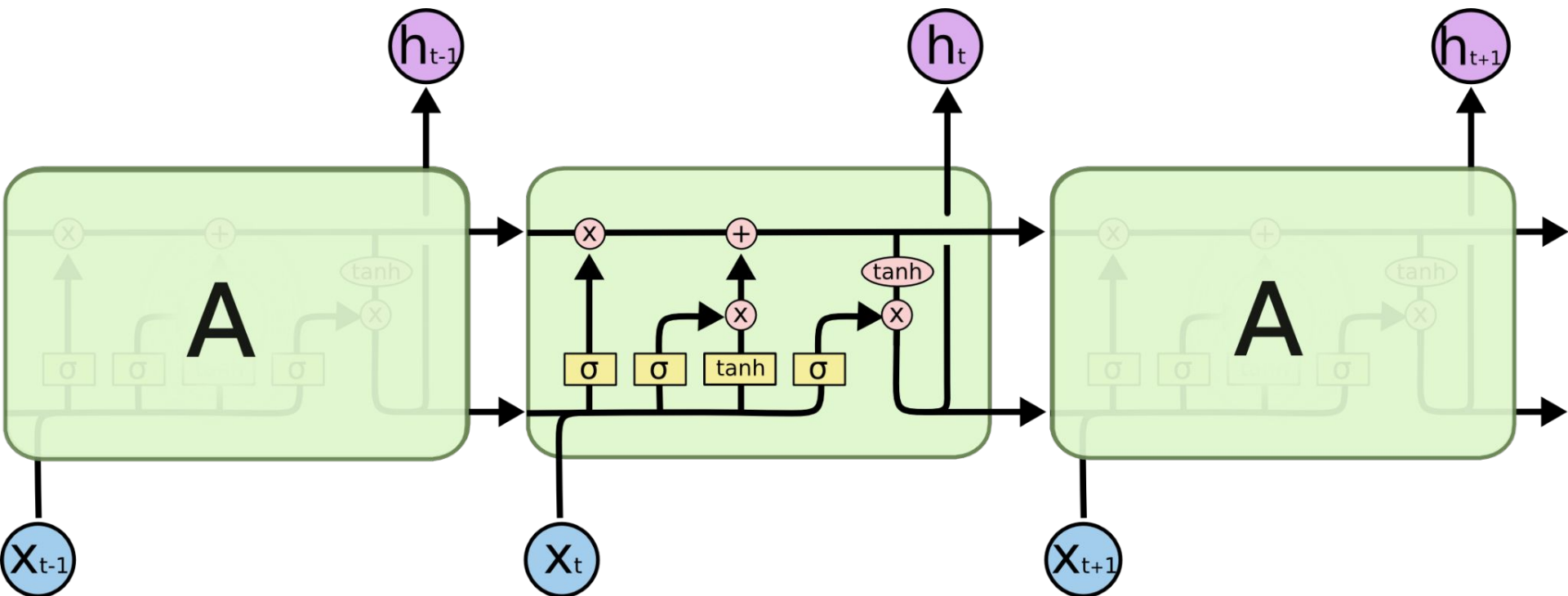
Canonical form of Recurrent Neural Networks



The output at time t depend not only on external inputs $U(t)$, but also (via internal “state variables”) on the whole sequence of previous inputs (and on initialization of state variables)

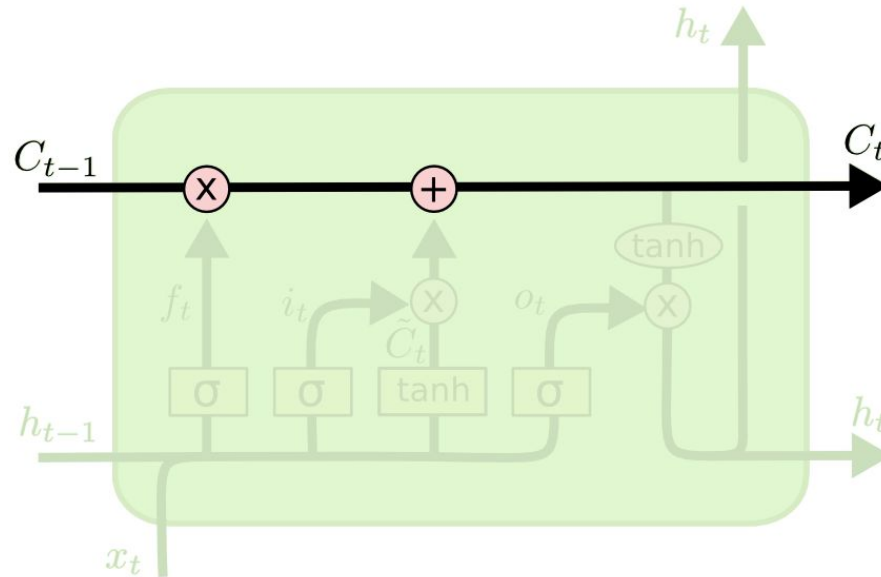
LSTMs learn how much data to keep

Four little networks keep track of a state and add information as necessary



X = Element-Wise Multiplication
 + = Element-Wise Addition

LSTM Steps

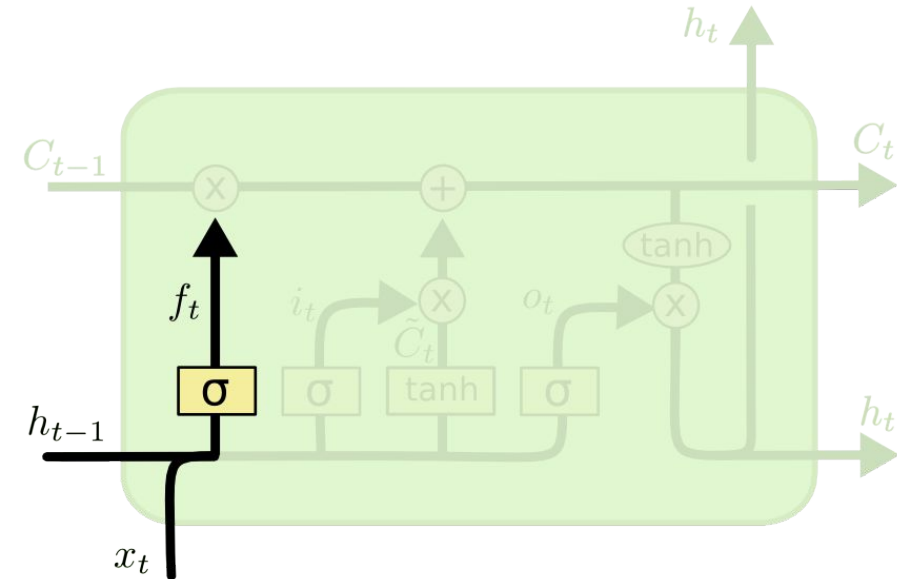


C = Cell State

Works like long term memory

From: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

LSTM Steps - Forget Gate



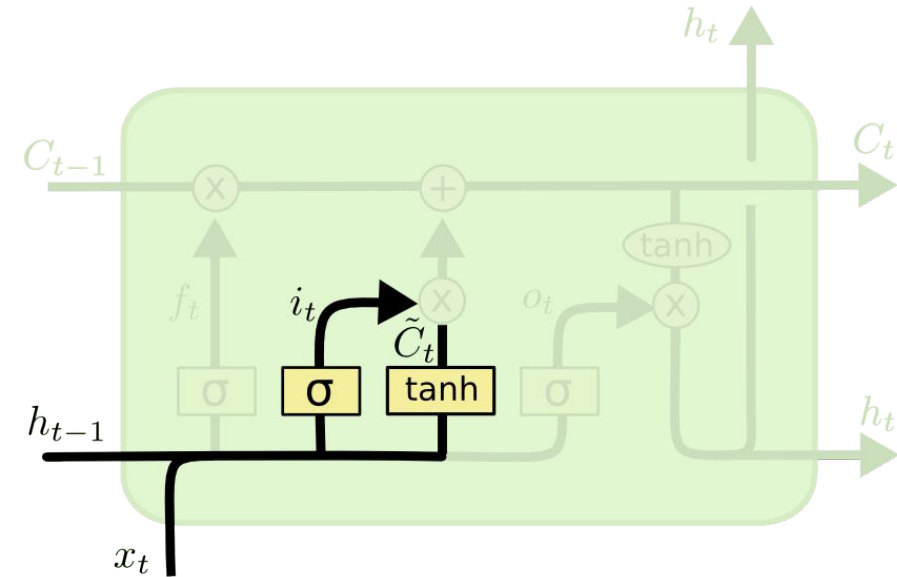
σ = Sigmoid Layer

$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

h = Hidden State

From: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

LSTM Steps - Input Gate



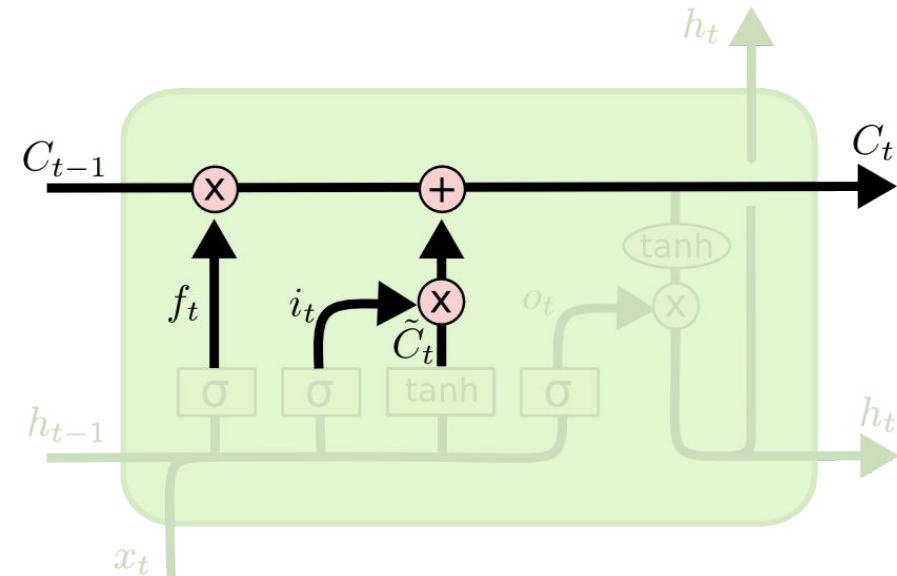
$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

\tilde{C}_t = New Candidate Cell Value

From: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

LSTM Steps

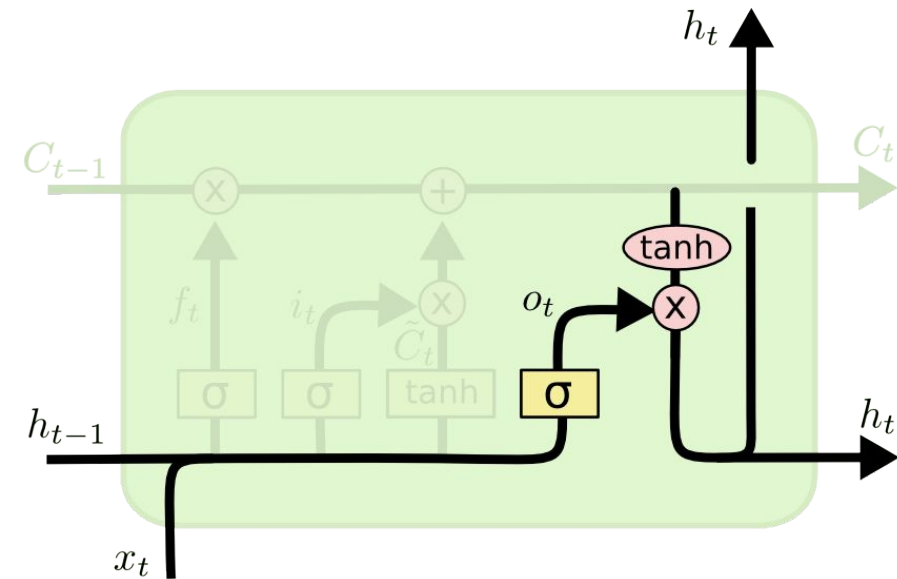


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

C_t = New Cell State

From: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

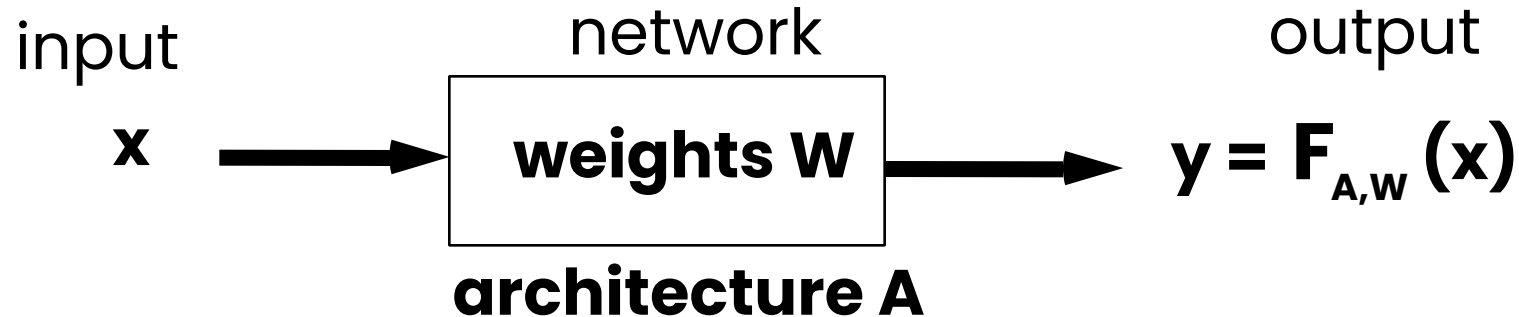
LSTM Steps



$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

From: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>



- Two modes:
 - Training: based on examples of (input,output) couples, the network modifies
 - Its parameters W (synaptic weights of connections)
 - And also potentially its architecture A
(by creating/eliminating neurons or connections)
 - Inference:
computation of output associated to a given input
(architecture and weights remaining frozen)

Training principle for Neural Networks

- **Supervised training = adaptation of synaptic weights of the network so that its output is close to target value for each example**
- **Given n examples $(X_p; D_p)$, and the network outputs $Y_p = \text{NN}(X_p)$, the average quadratic error is**

$$E(W) = \sum_p (Y_p - D_p)^2$$

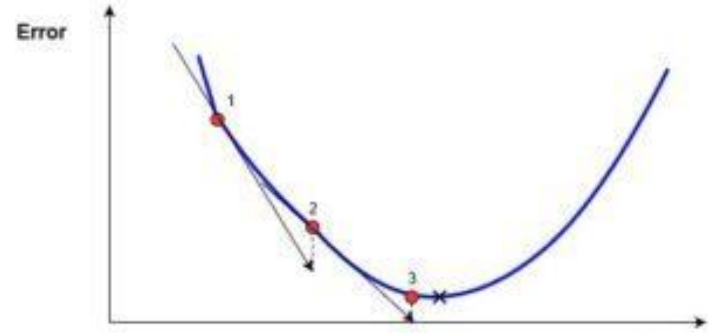
Training ~ finding $W^* = \text{ArgMin}(E)$, ie minimize the cost function $E(W)$

- **Generally this is done by using gradient descent (total, partial or stochastic):**

$$W(t+1) = W(t) - \lambda \cdot \text{grad}_W(E) \quad [+ \mu(t)(W(t) - W(t-1))]$$

Usual training algo for Multi Layer Perceptrons (MLP)

**Training a Neural Network
= optimizing values of its weights&biases**



- **Random initialization**
- Training by **Stochastic Gradient Descent** (SGD), using *back-propagation*:
 - Input 1 (or a few) random training sample(s)
 - Propagate
 - Calculate error (loss)
 - Back-propagate through all layers from end to input, to compute gradient and update weights

Smart method for efficient computing of gradient (w.r.t. weights) of a Neural Network cost function, based on chain rule for derivation.

Cost function is $Q(t) = \sum_m \text{loss}(Y_m, D_m)$, where m runs over training set examples

Usually, $\text{loss}(Y_m, D_m) = ||Y_m - D_m||^2$ [quadratic error]

Total gradient:

$$W(t+1) = W(t) - \lambda(t) \text{grad}_W(Q(t)) + \mu(t)(W(t) - W(t-1))$$

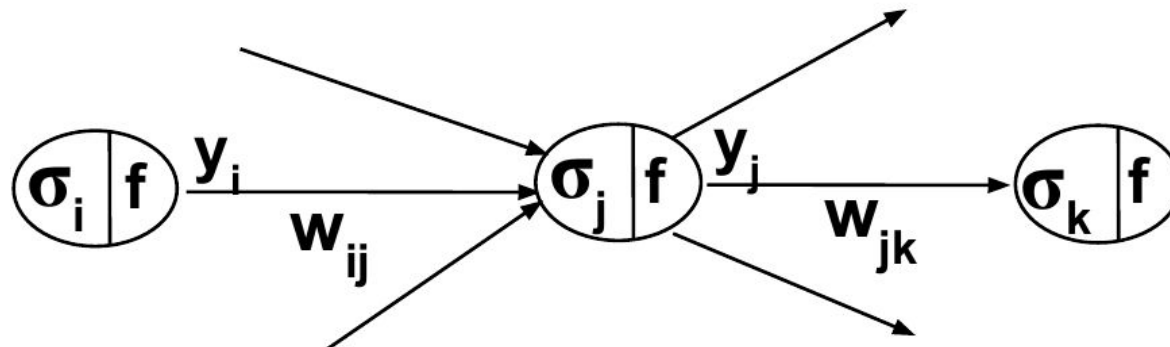
Stochastic gradient:

$$W(t+1) = W(t) - \lambda(t) \text{grad}_W(Q_m(t)) + \mu(t)(W(t) - W(t-1))$$

where $Q_m = \text{loss}(Y_m, D_m)$, is error computed on only ONE example randomly drawn from training set at every iteration and
 $\lambda(t)$ = learning rate (fixed, decreasing or adaptive), $\mu(t)$ = momentum

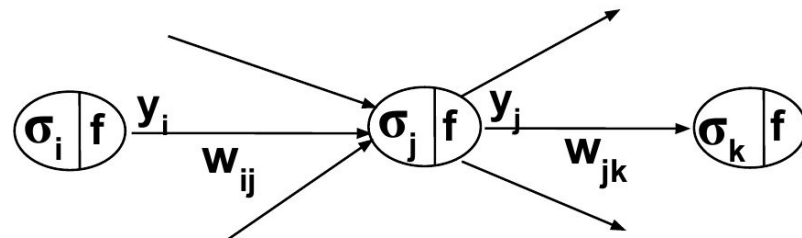
Now, how to compute dQ_m/dW_{ij} ?

Find update for weight



$$W_{ij}(t + 1) = W_{ij}(t) - \lambda(t) \frac{\partial E_m}{\partial W_{ij}}$$

$$E_m = ||Y_m - D_m||_2$$



$$dE_m/dW_{ij} = (dE_m/d\sigma_j)(d\sigma_j/dW_{ij}) = (dE_m/d\sigma_j) y_i$$

Let $\delta_j = (dE_m/d\sigma_j)$. Then

$$\mathbf{W}_{ij}(t+1) = \mathbf{W}_{ij}(t) - \lambda(t) y_i \delta_j$$

(and $W_{oj}(t+1) = W_{oj}(t) - \lambda(t)\delta_j$)

If neuron j is output, $\delta_j = (dE_m/d\sigma_j) = (dE_m/dy_j)(dy_j/d\sigma_j)$ with $E_m = ||Y_m - D_m||^2$

$$\text{so } \delta_j = 2(y_j - D_j)f'(\sigma_j) \text{ if neuron j is an output}$$

Otherwise, $\delta_j = (dE_m/d\sigma_j) = \sum_k (dE_m/d\sigma_k)(d\sigma_k/d\sigma_j) = \sum_k \delta_k (d\sigma_k/d\sigma_j) = \sum_k \delta_k W_{jk} (dy_j/d\sigma_j)$

$$\text{so } \delta_j = (\sum_k W_{jk} \delta_k) f'(\sigma_j) \text{ if neuron j is "hidden"}$$

□ all the δ_j can be computed successively from last layer to upstream layers by “error backpropagation” from output

Intermediate steps and tips

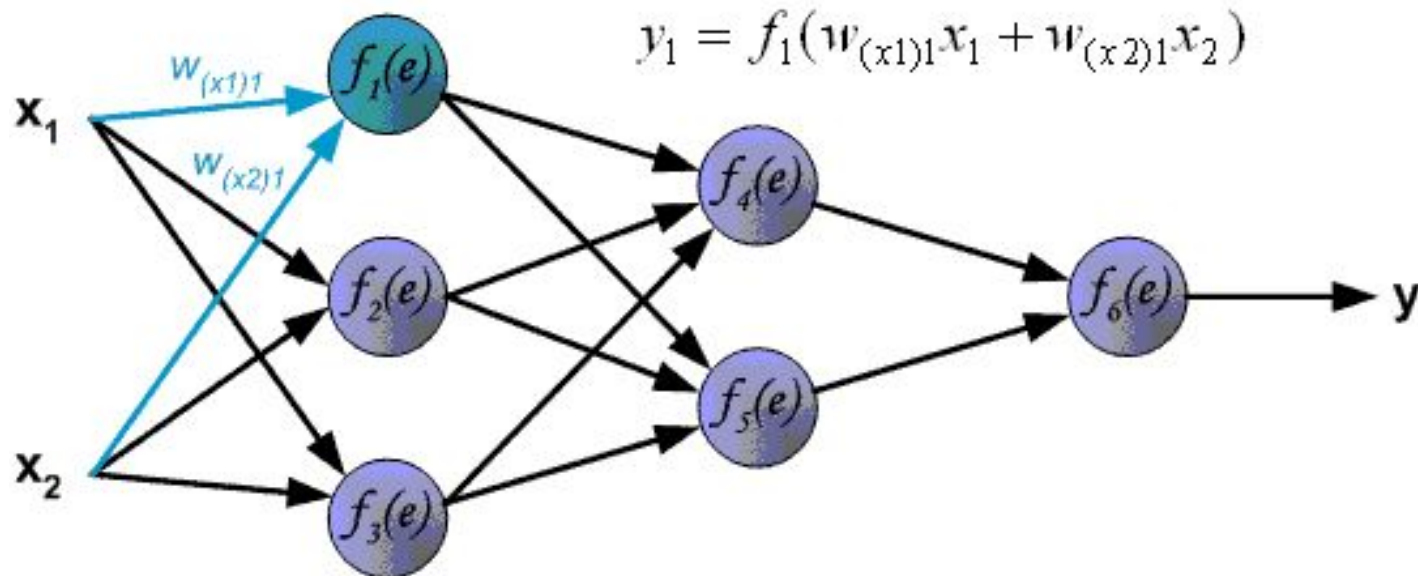
$$\sigma_j = \sum_i W_{ij} y_i$$

$$\frac{dy_j}{d\sigma_j} = f'(\sigma_j)$$

$$\frac{d\sigma_k}{d\sigma_j} = \frac{d\sigma_k}{dy_j} \frac{dy_j}{d\sigma_j} = W_{jk} f'(\sigma_j)$$

Animated illustration of Back-Propagation

FP



Animated GIF from the good tutorial

<https://medium.com/datadriveninvestor/what-is-gradient-descent-intuitively-42f10dfb293f>

Universal approximation theorem

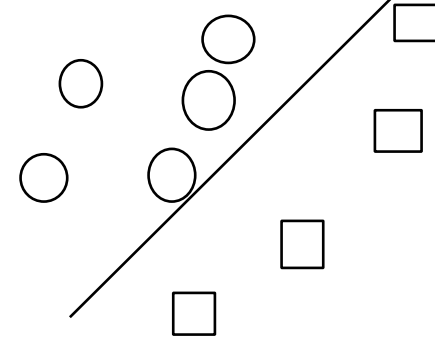
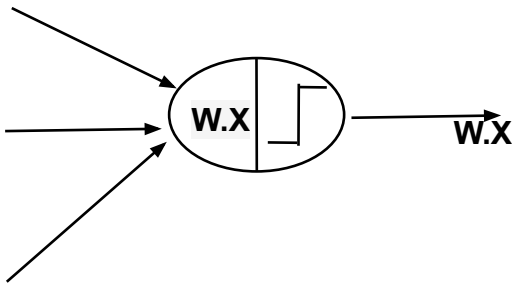
Cybenko 1989

- For any continuous function F defined and bounded on a bounded set, and for any ε , there exists a layered Neural Network with **ONLY ONE HIDDEN LAYER** (of *sigmoid* neurons) which approximates F with error $< \varepsilon$

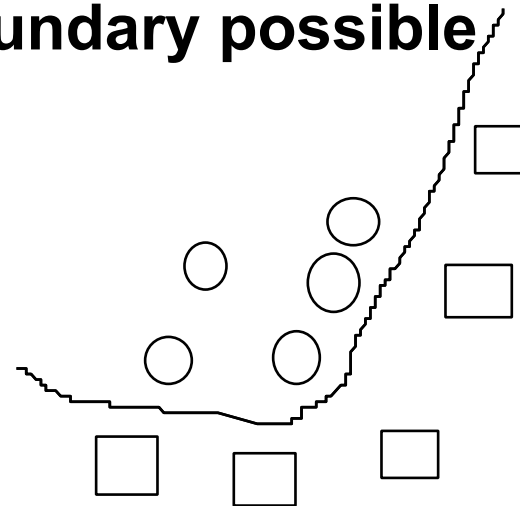
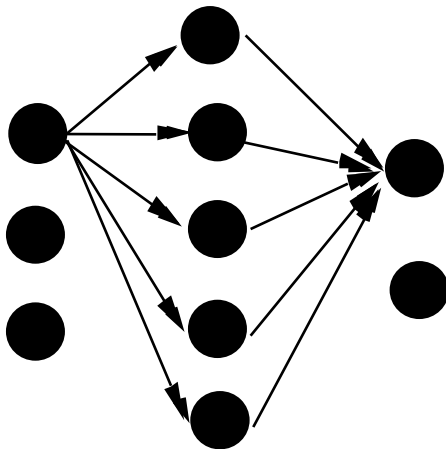
...But the theorem does not provide any clue about how to find this one_hidden-layer NN, nor about its size!
And the size of hidden layer might be huge...

Multi-layer (MLP) v.s. single-layer (perceptron)

Single-layer □ one linear separation by neuron



Multi-layer: any shape of boundary possible



ADVANTAGES

- Universal approximators (& classifiers)
- Fast to compute
- Robustness to data noise
- Rather easy to train and program

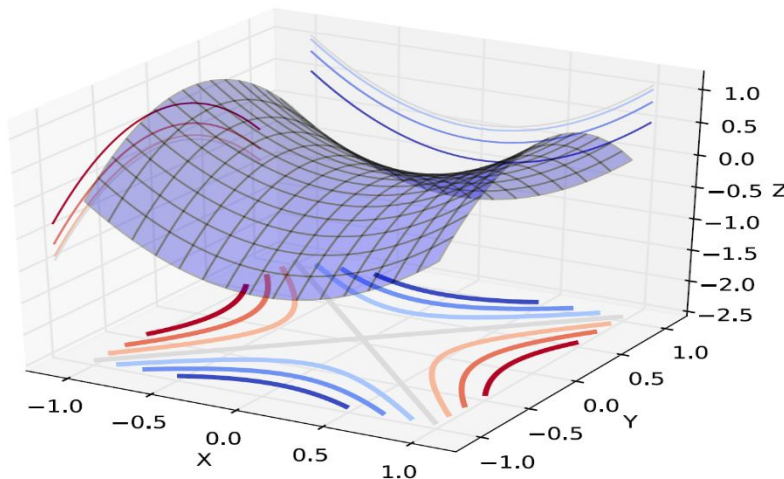
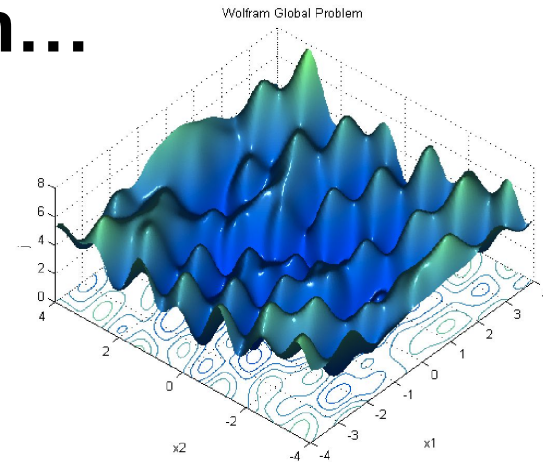
DRAWBACKS

- *Choice of ARCHITECTURE (# of neurons in hidden layer) is CRITICAL, and empiric!*
- Many other critical hyper-parameters
(learning rate, # of iterations, initialization of weights, etc...)
- *Many local minima in cost function*
- Blackbox: difficult to interpret the model

Why gradient descent works *despite non-convexity?*

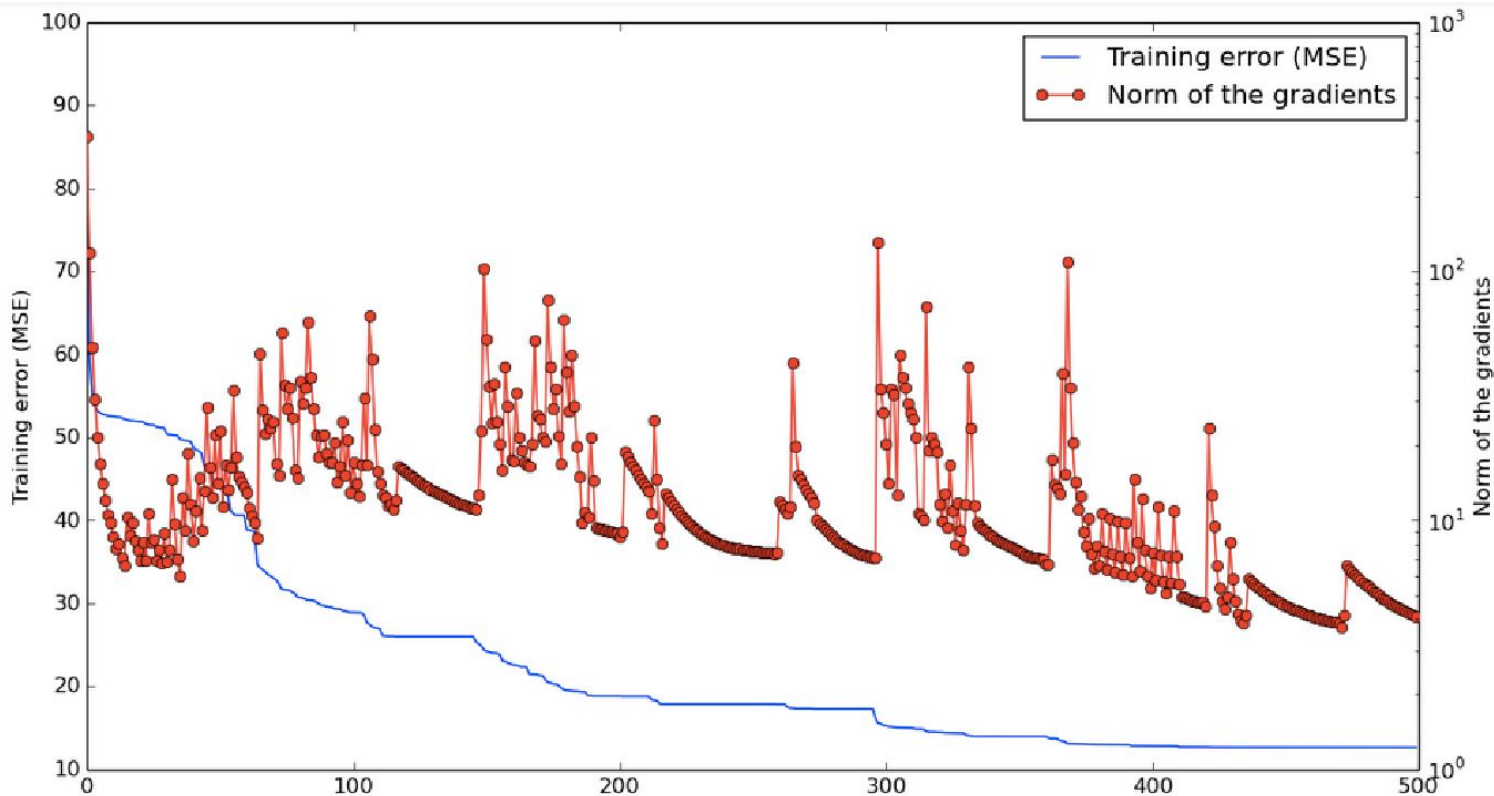
- Local minima dominate in low-Dim...

- ...but recent work has shown saddle points dominate in high-Dim



- Furthermore, most local minima are close to the global minimum

Saddle points in training curves



- **Oscillating between two behaviors:**
 - Slowly approaching a saddle point
 - Escaping it

METHODOLOGY FOR SUPERVISED TRAINING OF MULTI-LAYER NEURAL NETWORKS

Training set vs. TEST set

- Space of possible input values usually infinite, and training set is only a FINITE subset
- Zero error on all training examples \neq good results on whole space of possible inputs (cf generalization error \neq empirical error...)

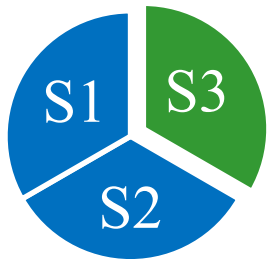


- Need to collect enough and representative examples
- Essential to keep aside a subset of examples that shall be used only as TEST SET for estimating final generalization (when training finished)
- Need also to use some “validation set” independant from training set, in order to tune all hyper-parameters (layer sizes, number of iterations, etc...)

Optimize hyper-parameters by "VALIDATION"

To **avoid over-fitting** and **maximize generalization**, absolutely essential to use some VALIDATION estimation, for optimizing training hyper-parameters (and stopping criterion):

- either use a *separate validation dataset* (random split of data into Training-set + Validation-set)
- or use CROSS-VALIDATION:
 - Repeat k times: train on $(k-1)/k$ proportion of data + estimate error on remaining $1/k$ portion
 - Average the k error estimations



3-fold cross-validation:

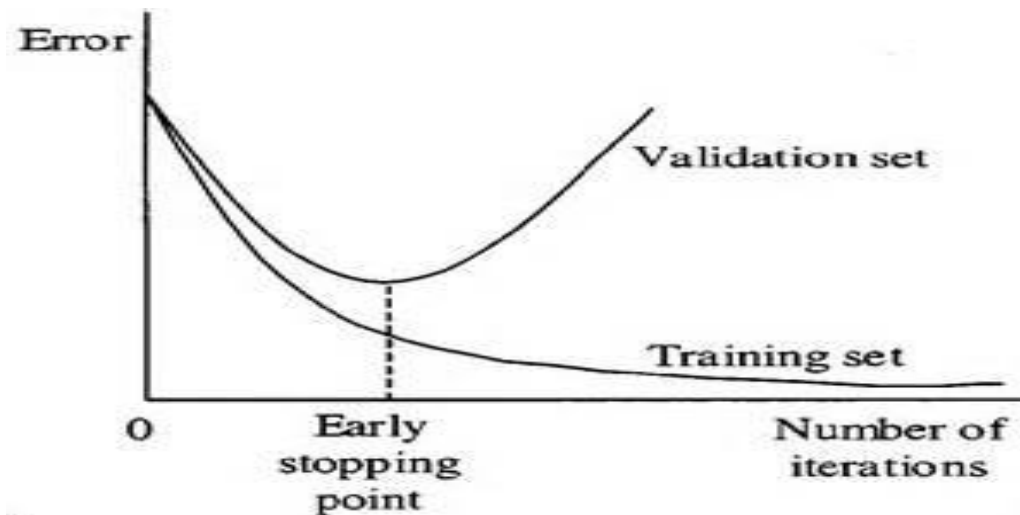
- Train on $S1 \cup S2$ then estimate err_{S3} error on $S3$
- Train on $S1 \cup S3$ then estimate err_{S2} error on $S2$
- Train on $S2 \cup S3$ then estimate err_{S1} error on $S1$
- Average validation error: $(\text{err}_{S1} + \text{err}_{S2} + \text{err}_{S3})/3$

Some Neural Networks training "tricks"

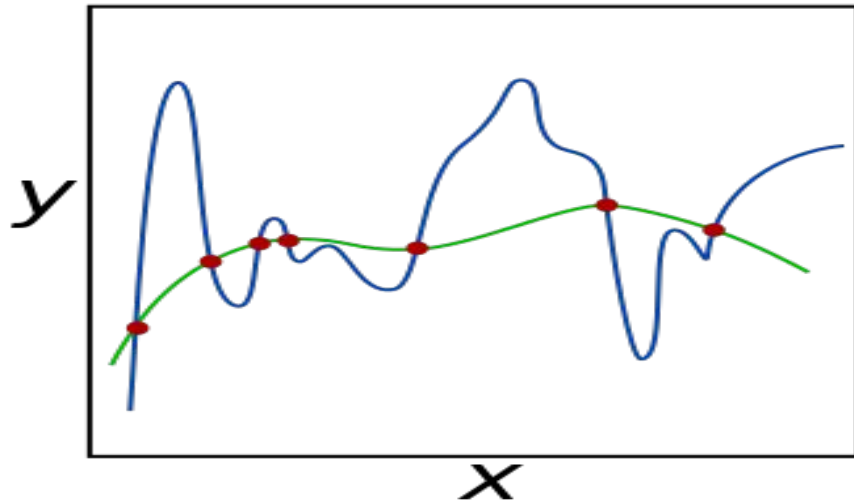
- Importance of input normalization
(zero mean, unit variance)
- Importance of weights initialization
random but SMALL and prop. to $1/\sqrt{\text{nbInputs}}$
- Decreasing (or adaptive) learning rate
- Importance of training set size
If a Neural Net has a LARGE number of free parameters,
□ train it with a sufficiently large training-set!
- Avoid overfitting by Early Stopping of training iterations
- Avoid overfitting by use of L1 or L2 regularization

Avoid overfitting by **EARLY STOPPING**

- For Neural Networks, a first method to avoid overfitting is to **STOP LEARNING** iterations as soon as the *validation_error* stops decreasing
- Generally, not a good idea to decide the number of iterations beforehand. Better to **ALWAYS USE EARLY STOPPING**



Avoid overfitting using regularization penalty (weight decay)



Trying to fit too many free parameters with not enough information can lead to overfitting

Regularization = penalizing too complex models

Often done by adding a special term to cost function

For neural network, the regularization term is just norm L2 or L1 of vector of all weights:

$$K = \sum_m (\text{loss}(Y_m, D_m)) + \beta \sum_{ij} |W_{ij}|^p \quad \text{with } p=2 \text{ (L2) or } p=1 \text{ (L1)}$$

□ name “**Weight decay**”

MLP hyper-parameters

- **Number and sizes of hidden layers!!**
- **Activation functions**
- **Learning rate (& momentum) [optimizer]**
- **Number of gradient iterations!! (& early_stopping)**
- **Regularization factor**
- **Weight initialization**

Tuning hyper-parameters of MLPs in practice

- Use 'adam' optimizer
- Test/compare ***WIDELY VARIED HIDDEN LAYER SIZES***
(typically 30;100;300;1000;30-30;100-100)
- Test/compare ***SEVERAL INITIAL LEARNING RATES***
(typically 0.1;0.03;0.01;0.003;0.001)
- Make sure **ENOUGH ITERATIONS** for convergence
(typically >200 epochs), but **EARLY STOPPING**
on validation_error to avoid overfitting
(☐ check by plotting learning curves!!)

Some Bonus slides....

Network Losses

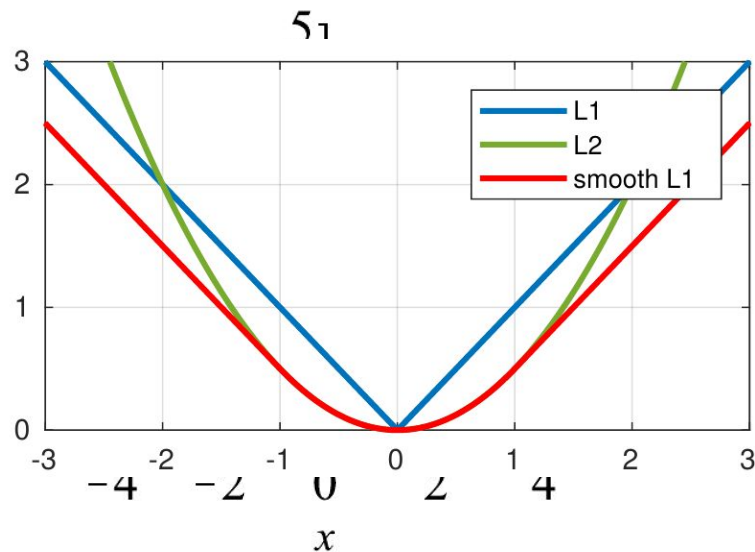
- For Networks: Best to find **differentiable** function
- Error must have **gradient** to point **towards** 'more correct' solution

Loss Application

- For classification: What is the probability distribution of the instance belonging to a certain class?
- For image segmentation: Calculate loss per pixel.
- **Main losses:** Solve the desired task
- **Auxiliary losses:** Achieve desired network behavior, potentially unrelated to final goal
- For tracking, Bounding Box position and size regression

Smooth L1 Loss (Huber loss with Delta=1)

- Less sensitive to outliers than L2, avoids *exploding* gradients better
- Still smooth around 0 opposite to L1



$$loss(x, y) = \begin{cases} 0.5(x - y)^2, & \text{if } |x - y| < 1 \\ |x - y| - 0.5, & \text{otherwise} \end{cases}$$

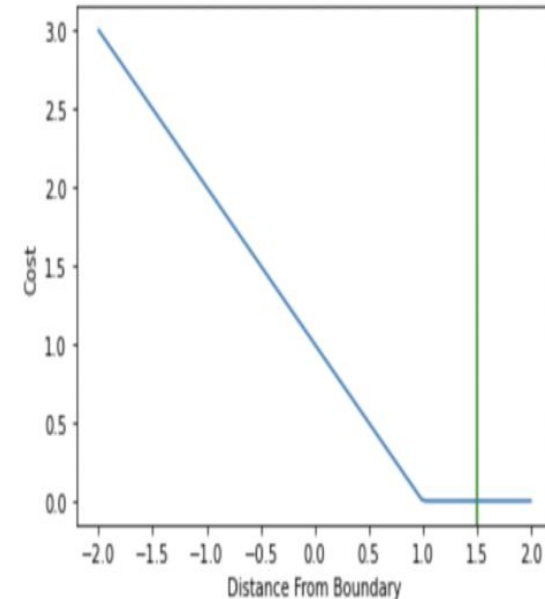
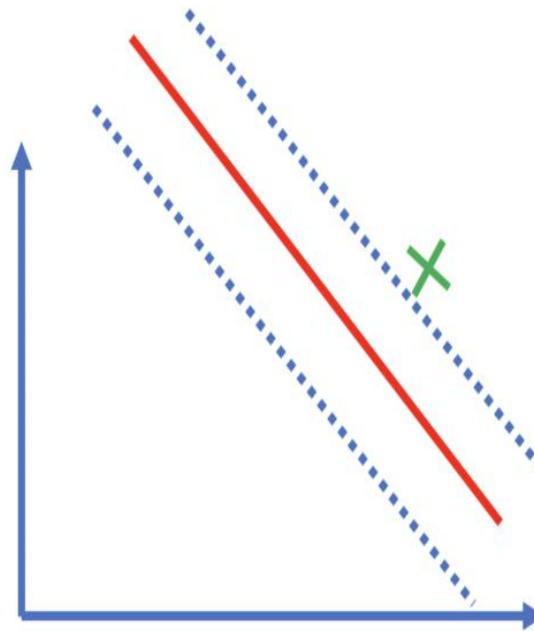
Prediction should be further away from decision boundary?

- Important for Support Vector Machines with binary prediction

$$l(y) = \max(0, 1 - t \cdot y)$$

Hinge-Loss:

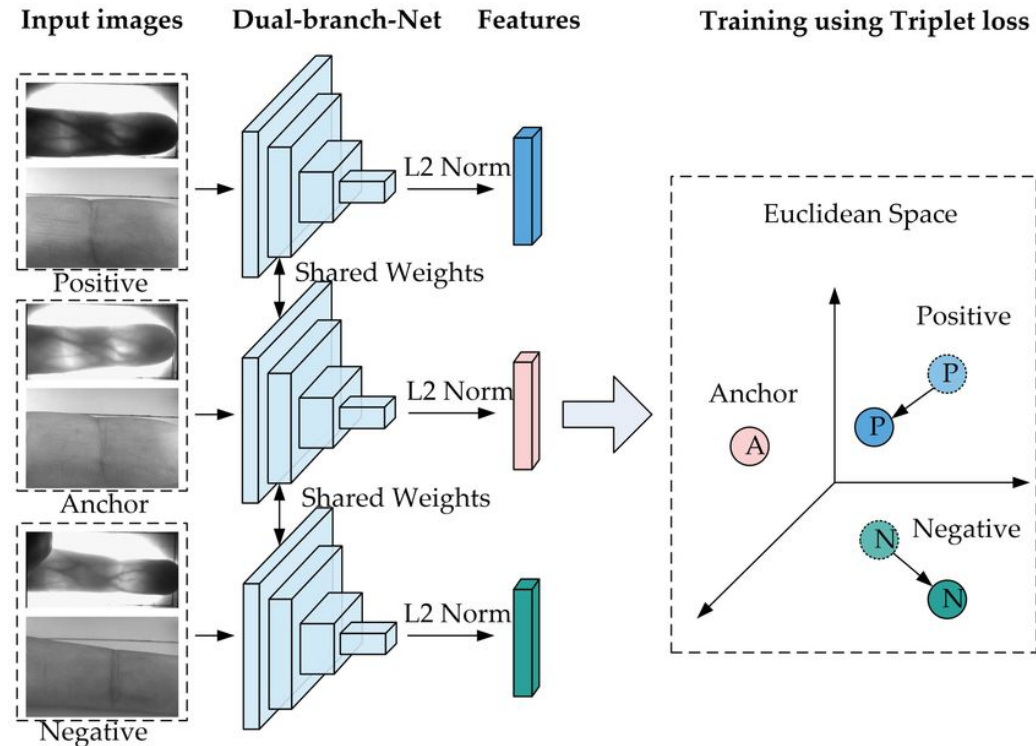
- t = Ground Truth ± 1
- y = Prediction



Compare: <https://programmatically.com/understanding-hinge-loss-and-the-svm-cost-function/>

More Specialized Losses Exist

- Negative Log-Likelihood Loss
- Binary Cross Entropy Loss
- Margin Ranking Loss
- Triplet Loss
- Cosine Embedding Loss
- Kullback-Leibler Divergence Loss



Some (old) references on (shallow, i.e. non deep) Neural Networks

- ***Réseaux de neurones : méthodologie et applications***, G. Dreyfus et al., Eyrolles, 2002.
- ***Réseaux de neurones formels pour la modélisation, la commande, et la classification***, L. Personnaz et I. Rivals, CNRS éditions, collection Sciences et Techniques de l'Ingénieur, 2003.
- ***Réseaux de neurones : de la physique à la psychologie***, J.-P. Nadal, Armand Colin, 1993.