# Introduction to Machine Learning



## Marco Cristani

"...what we want is a machine that can learn from experience.

Alan Turing, 1947

DEEP LEARNING CONSPIRACY IN NATURE, 521, P 436-444
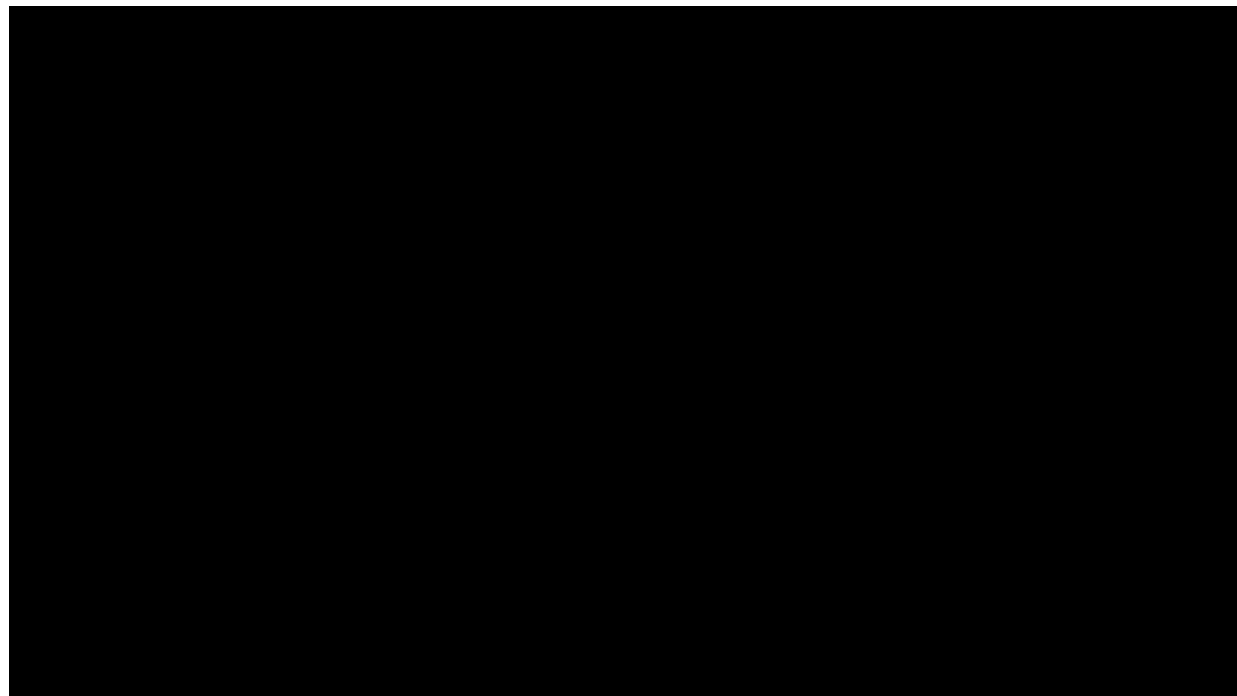
JÜRGEN SCHMIDHUBER'S HOME PAGE
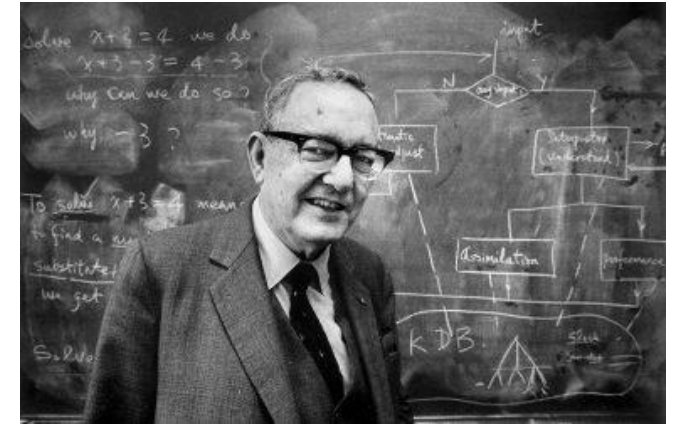
Yoshua Bengio

Geoffrey Hinton

Yann LeCun

# What is machine learning?

- "Learning is any process by which a system improves performance from experience"
  - Herbert Simon (*Turing Award* 1975, *Nobel Prize* in
    Economics 1978)

- Definition by Tom Mitchell (1998):
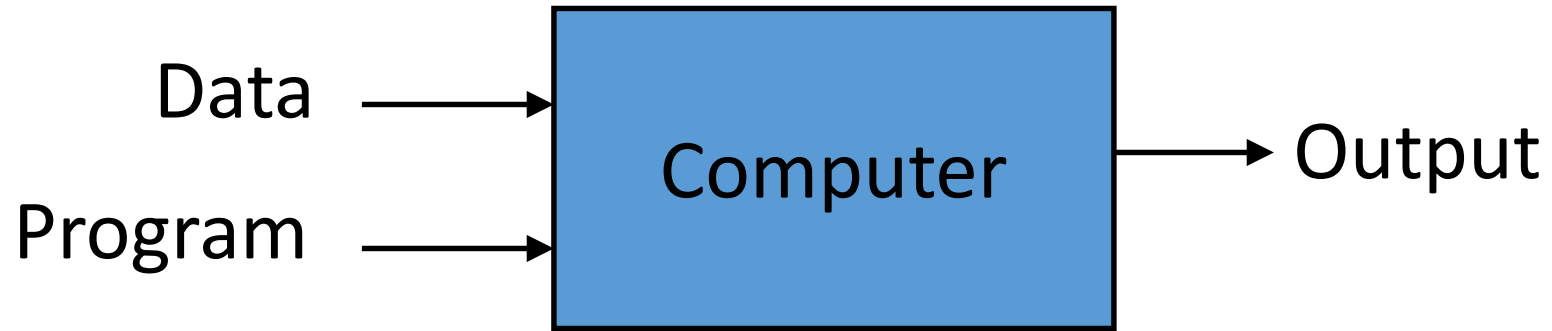
  Machine Learning is the study of algorithms that
  - improve their performance *P*
  - at some task *T*
  - with experience *E*

  A well-defined learning task is given by *<P, T, E>*

# Difference w.r.t. traditional programming

**Traditional Programming**

Data → Computer

Program → Computer

Computer → Output

**Machine Learning**

Data → Computer

Output → Computer

Computer → Program

# Magic?

**No, more like gardening**

- **Seeds** = Algorithms
- **Nutrients** = Data
- **Gardener** = You
- **Plants** = Programs

# When Do We Use Machine Learning?

- ML is used when:
  - Human expertise does not exist (navigating on Mars, forecasting)

  - Humans can't explain their expertise (recognition)

  - Models must be customized (personalized medicine)

  - Models are based on huge amounts of data (genomics)

*A case of handwritten recognition: what makes a «2»?*

# Some more examples of tasks that are best solved by using a learning algorithm

- Recognizing patterns:
  - Facial identities or facial expressions
  - Handwritten or spoken words
  - Medical images
- Generating patterns:
  - Generating images or motion sequences
- Recognizing anomalies:
  - Unusual credit card transactions
  - Unusual patterns of sensor readings in a nuclear power plant
- Forecasting:
  - Future stock prices or currency exchange rates

# Some *bachelor* projects in machine learning

- Expression recognition
  - (Dr. *Luca Brunelli*, now in Statwolf, data science, https://www.statwolf.com/)



Angry  NULL  Disgust Fear  Happy  Sadness Surprise

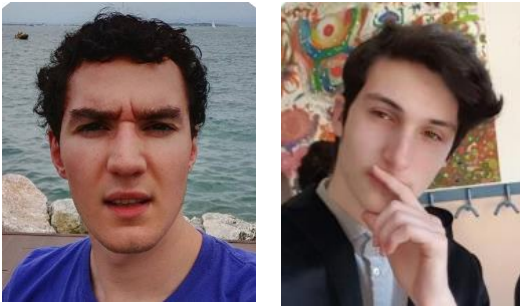# Some *bachelor* projects in machine learning

- Scene recognition

**Predictions:**

- **Type of environment:** indoor
- **Semantic categories:** office:0.61, home_office:0.13,
- **SUN scene attributes:** enclosedarea, nohorizon, cloth, man-made, electricindoorlighting, matte, research, sterile

# Some *bachelor* projects in machine learning
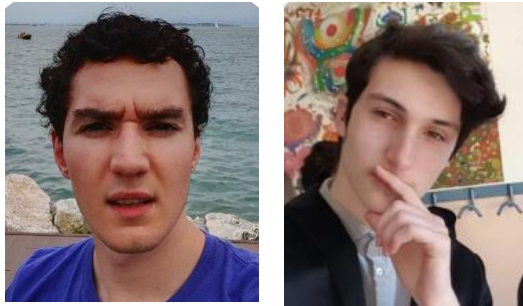
- creating 3D objects from 2D images



(Carlo Veronesi, Nicholas Merci, their ongoing bachelor thesis)

2D image

3D *generated* images

a particular of the 3D body

# Some *bachelor* projects in machine learning

- creating 3D objects from 2D images (2)



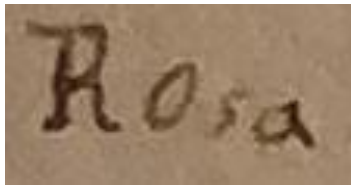(Carlo Veronesi, Nicholas Merci, their ongoing bachelor thesis)
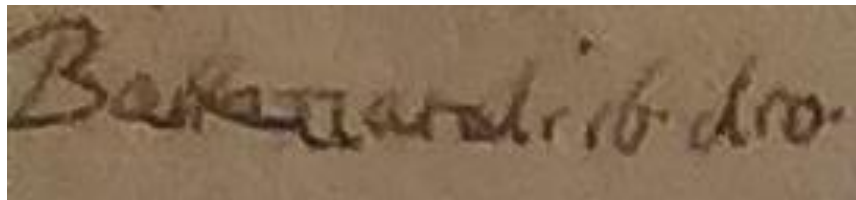
2D image

3D *generated* images

a particular of the 3D body

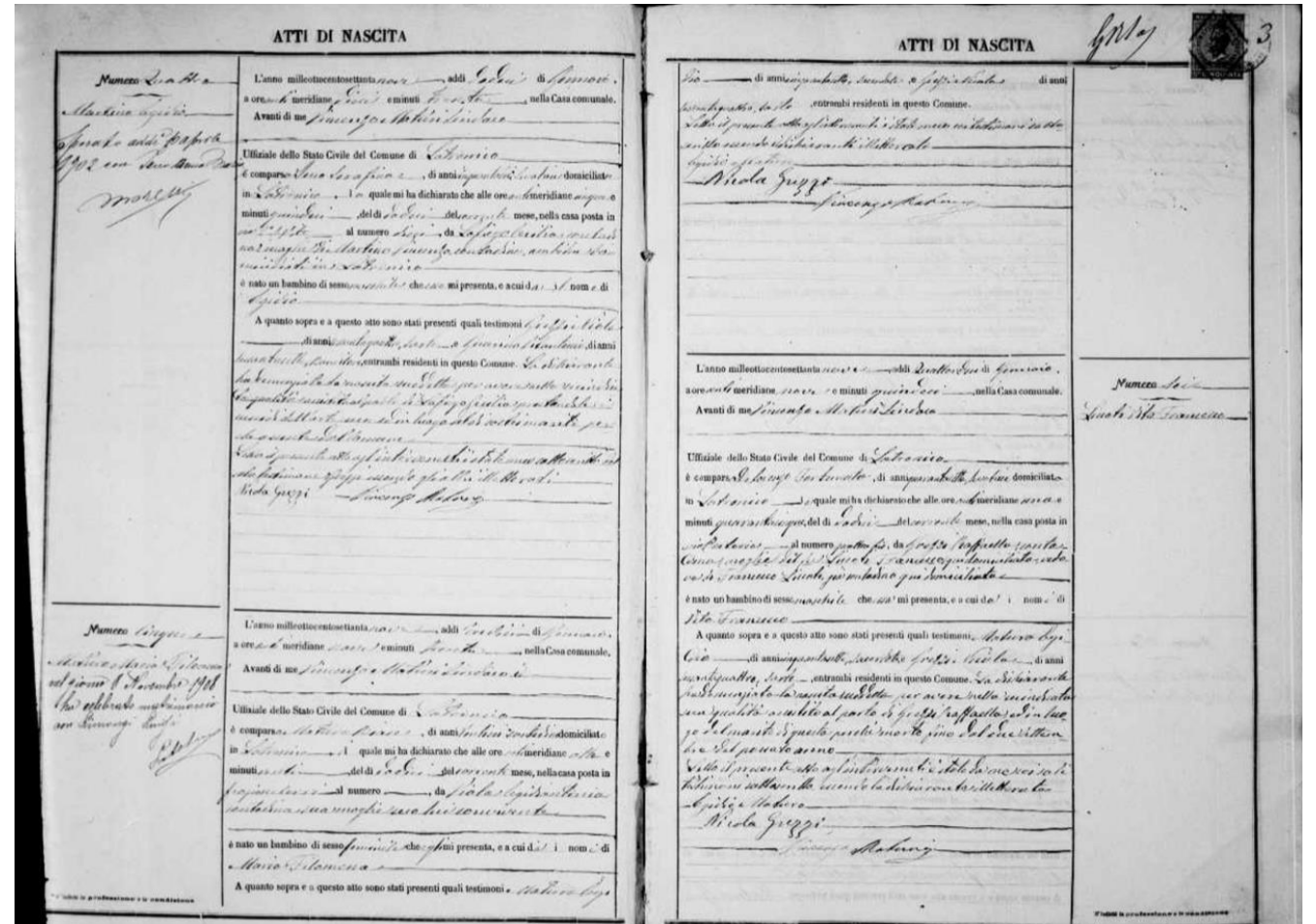# Some *bachelor* projects in machine learning

- Birth records digitalization towards family tree estimation



Rosa



Battezzardi ib tro  (???)

# Sample applications

- Web search
- Computational biology
- Finance
- E-commerce
- Space exploration
- Robotics
- Information extraction
- Social networks
- Debugging
- *[Your favorite area]*

# Defining the learning task

*Improve on task T, with respect to performance metric P, based on experience E*

T: Playing checkers
P: Percentage of games won against an arbitrary opponent
E: Playing practice games against itself

T: Driving on four-lane highways using vision sensors
P: Average distance traveled before a human-judged error
E: A sequence of images and steering commands recorded while observing a human driver.

T: Recognizing hand-written words
P: Percentage of words correctly classified
E: Database of human-labeled images of handwritten words

T: Categorize email messages as spam or legitimate.
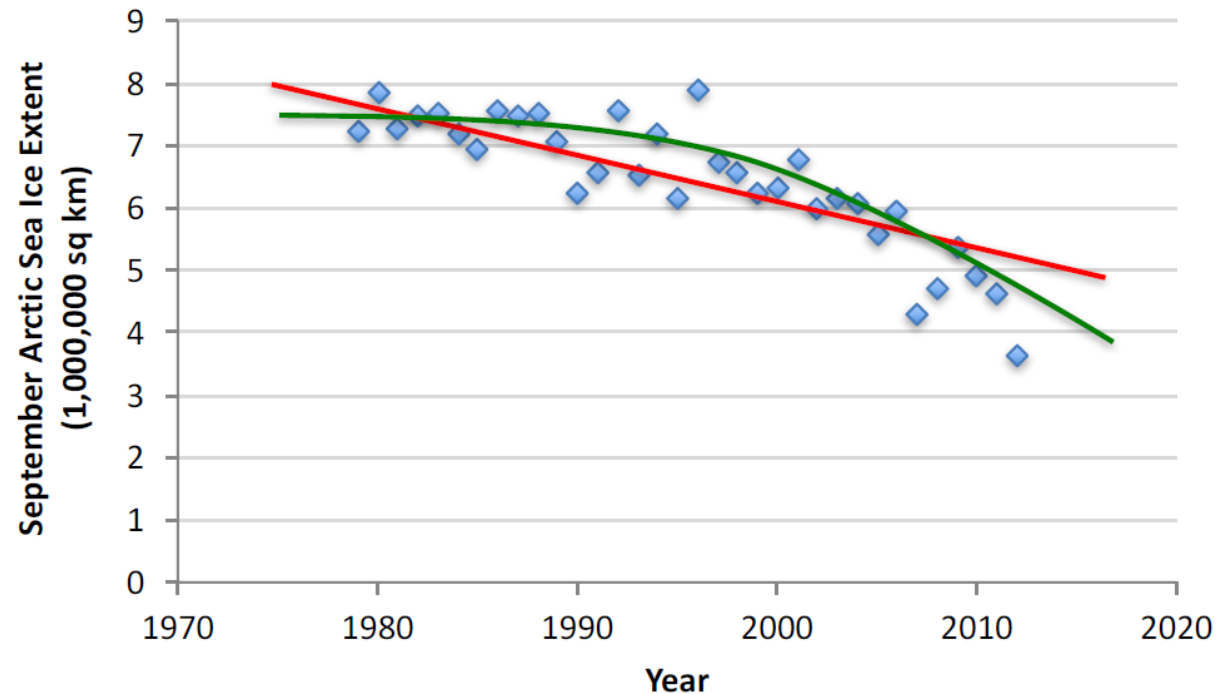P: Percentage of email messages correctly classified.
E: Database of emails, some with human-given labels

# Types of learning

- **Supervised (inductive) learning (= *regression|sup.ed classification*)**
  - Given: training data + desired outputs (labels)
- **Unsupervised learning**
  - Given: training data (without desired outputs)
- **Semi-supervised learning**
  - Given: training data + a few desired outputs
- **Reinforcement learning**
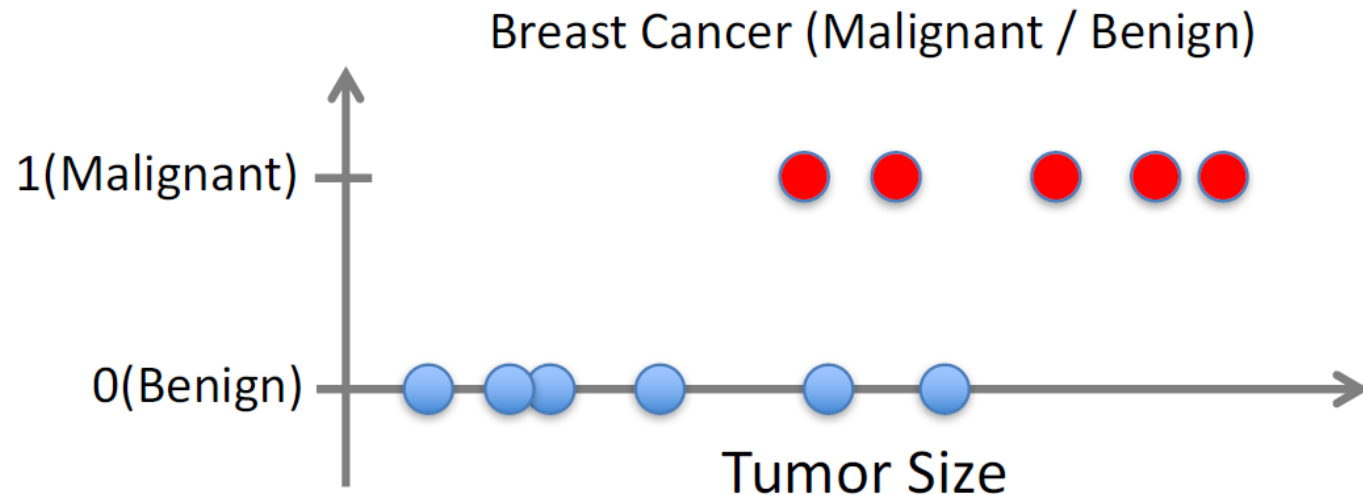  - Rewards from sequence of actions

# Supervised Learning: Regression

- Given (x1, y1), (x2, y2), ..., (xn, yn)

- Learn a function f(x) to predict y given x

- – y is real-valued == regression

# Supervised Learning: Classification

- Given (x1, y1), (x2, y2), ..., (xn, yn)
- Learn a function f(x) to predict y given x
- – y is categorical == (supervised) classification



Breast Cancer (Malignant / Benign)

1(Malignant)

0(Benign)

Tumor Size

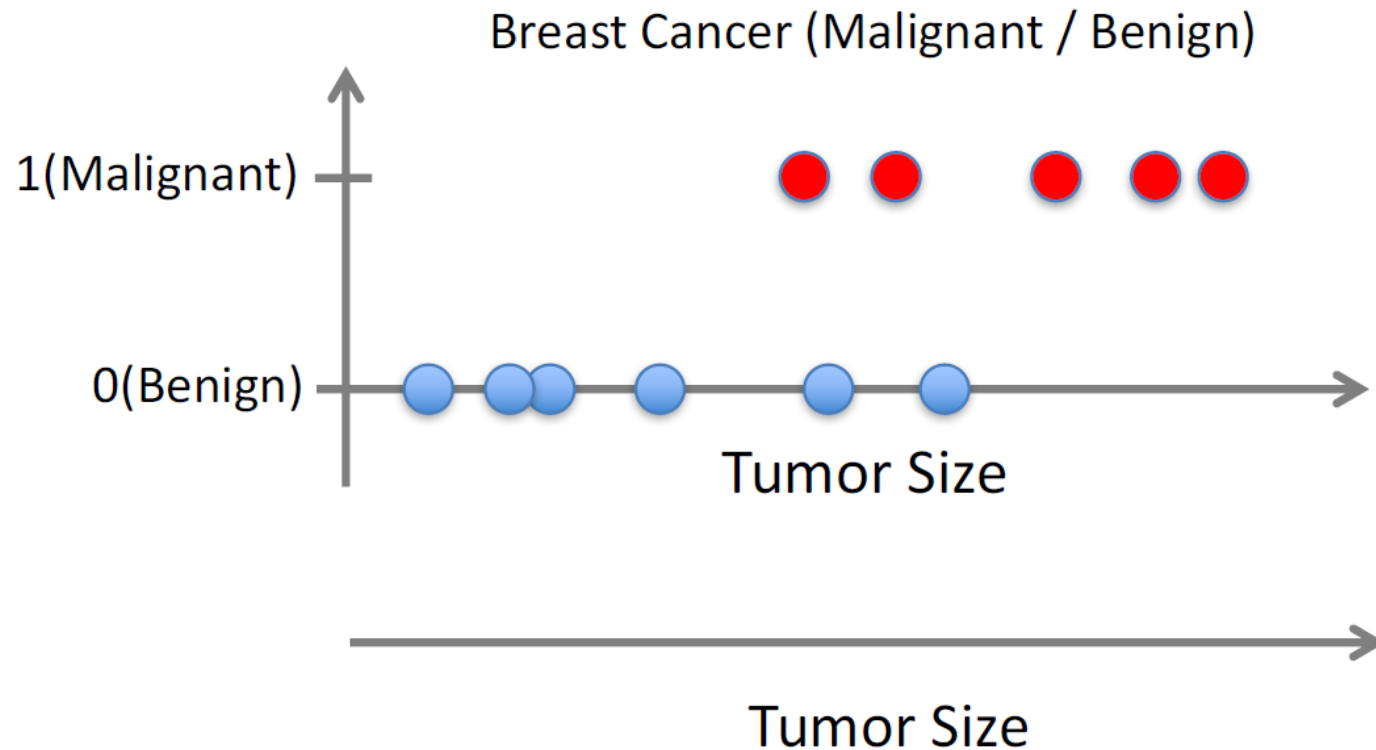# Supervised Learning: Classification

- Given (x1, y1), (x2, y2), ..., (xn, yn)

- Learn a function f(x) to predict y given x

- – y is categorical == (supervised) classification



Breast Cancer (Malignant / Benign)

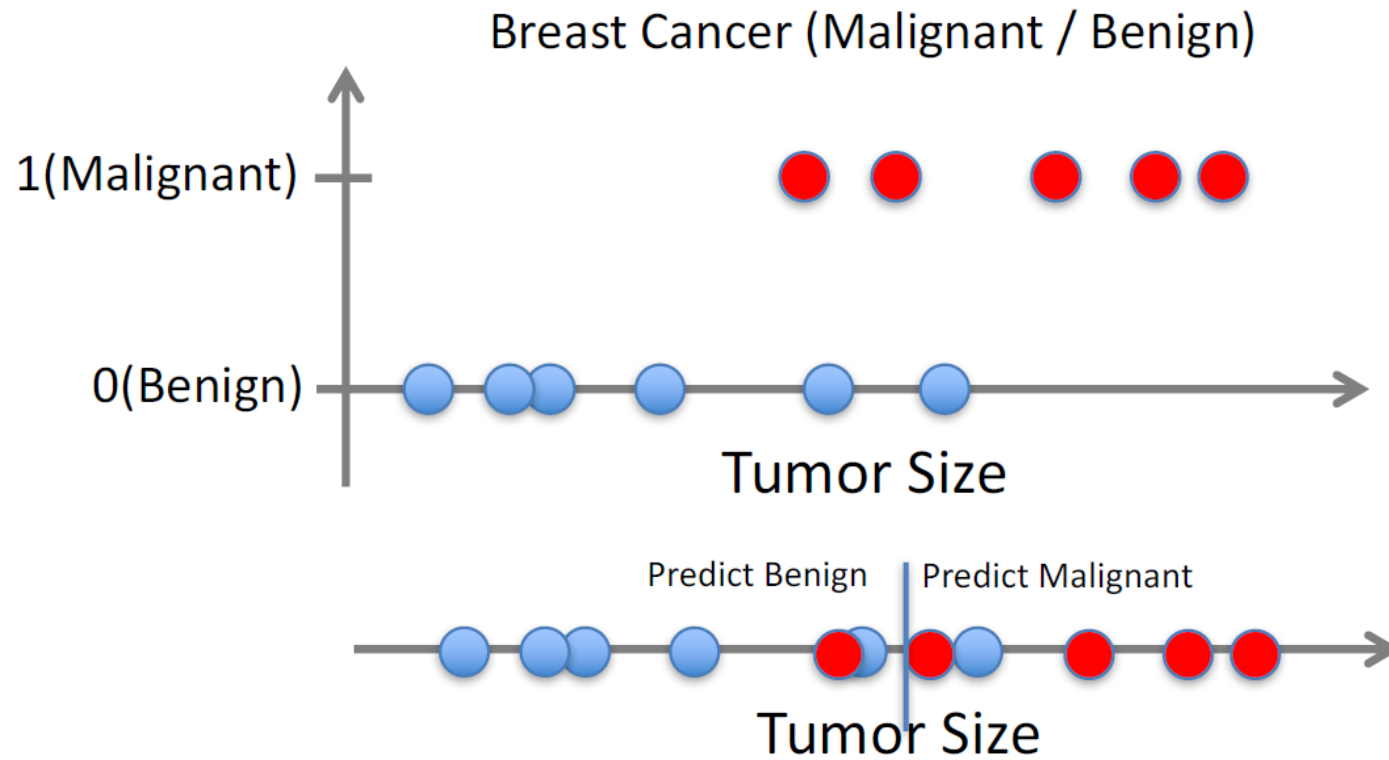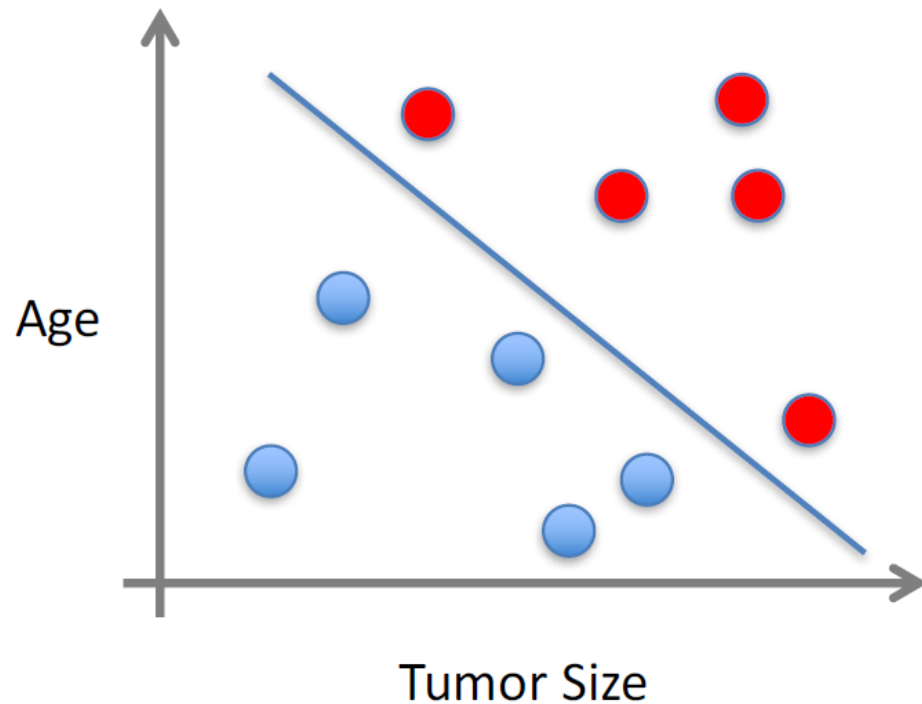# Supervised Learning: Classification

- Given (x1, y1), (x2, y2), ..., (xn, yn)
- Learn a function f(x) to predict y given x
- – y is categorical == (supervised) classification

Breast Cancer (Malignant / Benign)

1(Malignant)

0(Benign)

Tumor Size

Predict Benign | Predict Malignant
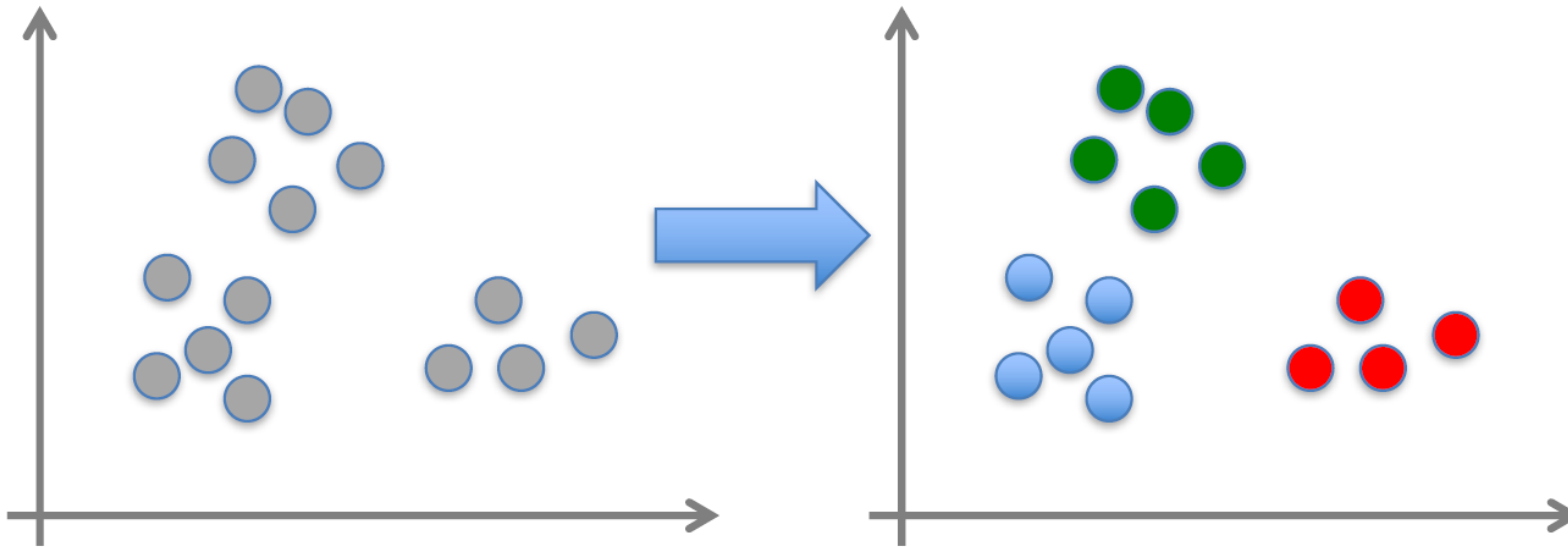
Tumor Size

# Supervised Learning

- x can be multi-dimensional

    – Each dimension corresponds to an attribute



- Clump Thickness
- Uniformity of Cell Size
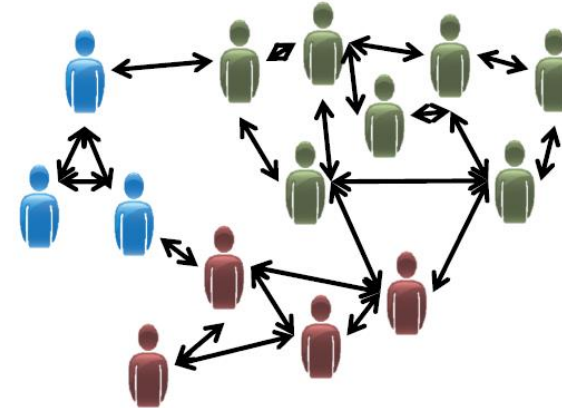- Uniformity of Cell Shape

...

# Unsupervised Learning

- Given x1, x2, ..., xn (without labels)
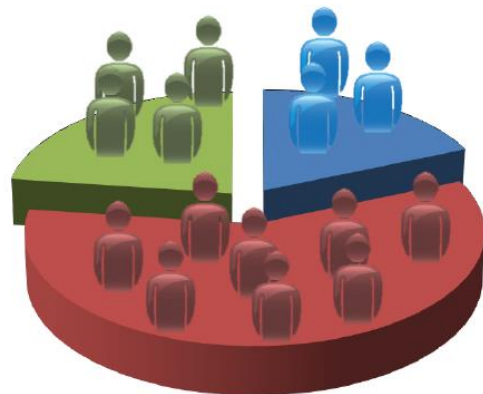- Output hidden structure behind the x's
– e.g., clustering

# Unsupervised Learning


Organize computing clusters
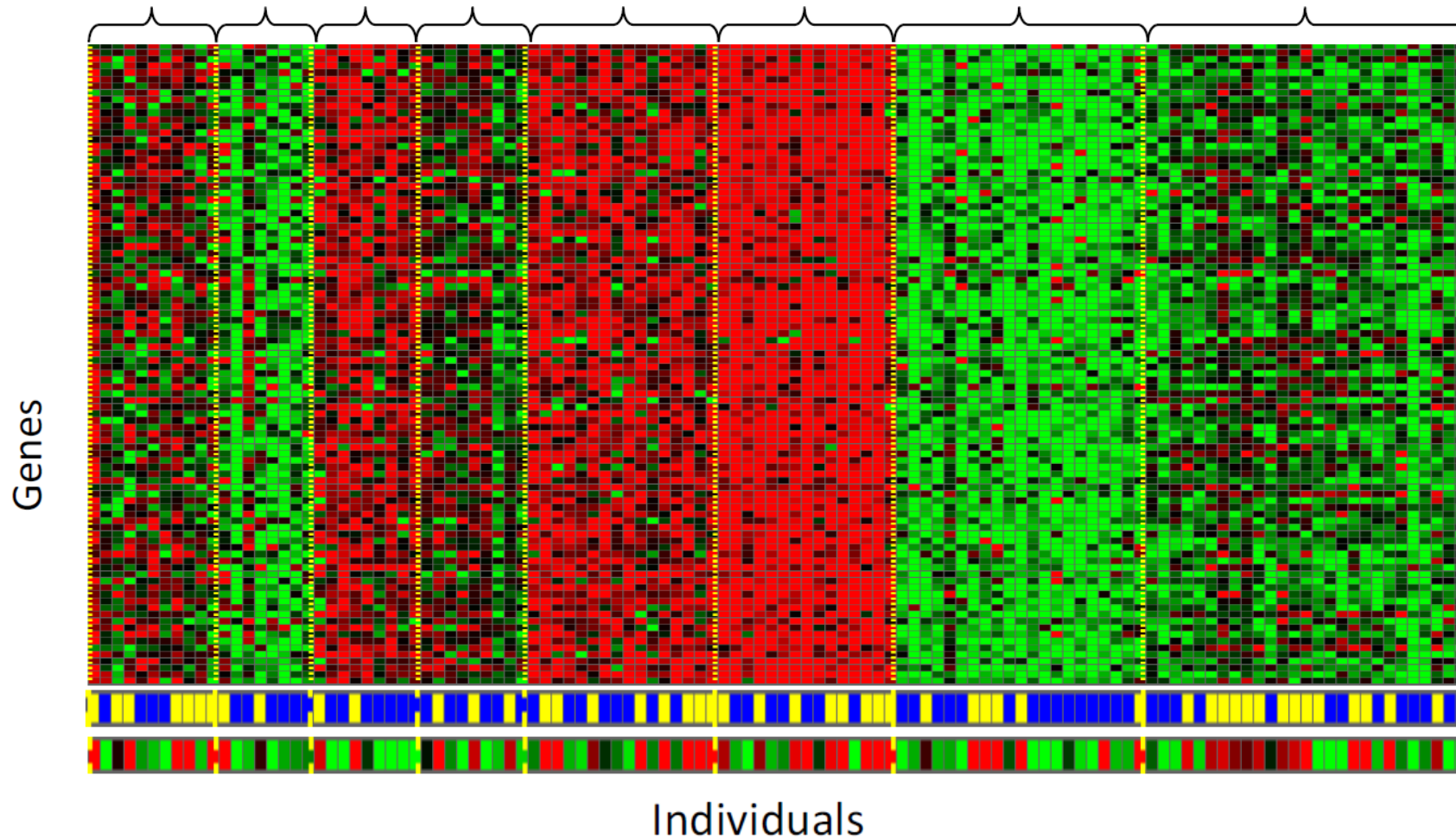

Social network analysis


Market segmentation


Astronomical data analysis

Image credit: NASA/JPL-Caltech/E. Churchwell (Univ. of Wisconsin, Ma
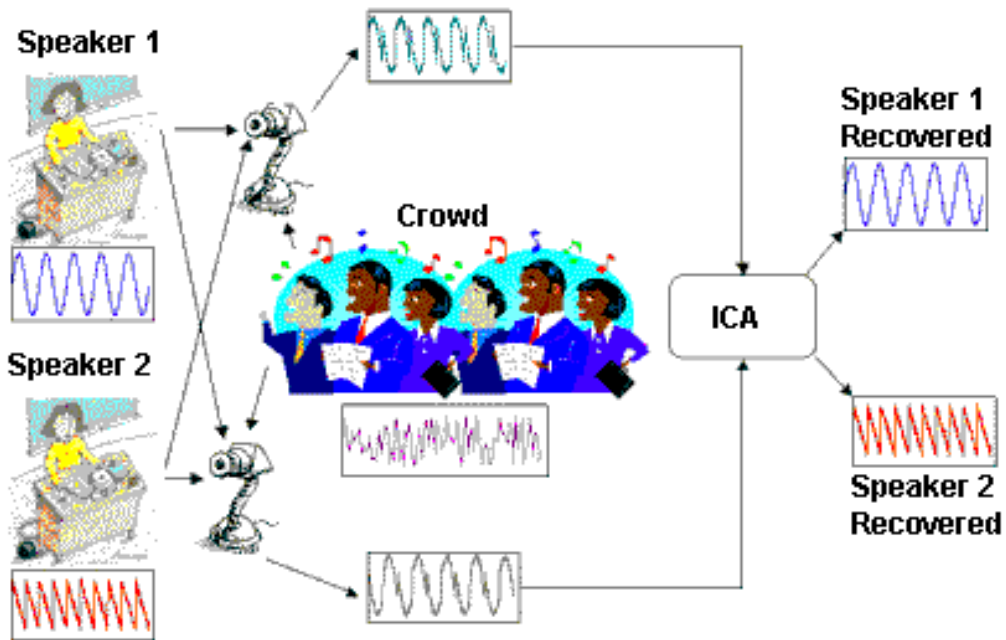
# Unsupervised Learning

- Genomics application: group individuals by genetic similarity

# Unsupervised Learning

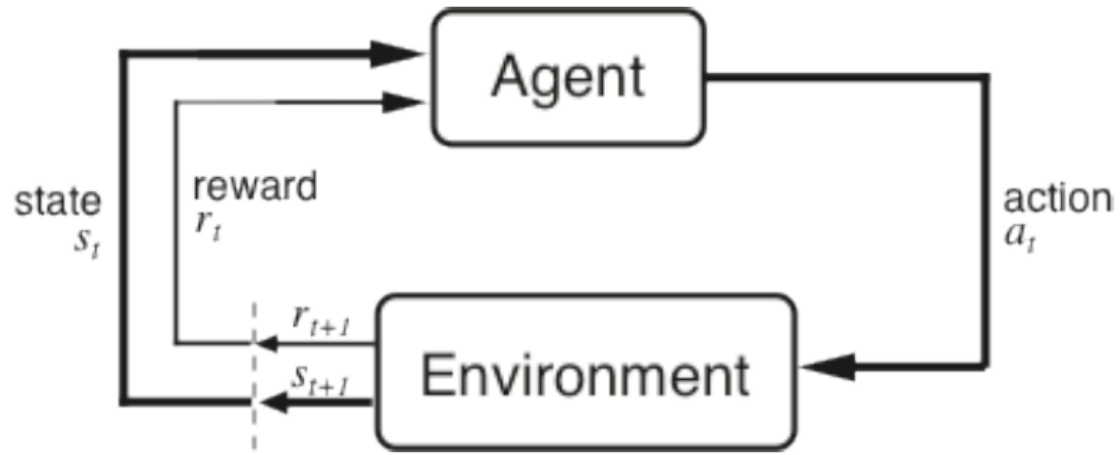- Independent component analysis – separate a combined signal into its original sources





Input video (two people speaking together)

JOHN    ALL    RORY

Video source: Team Coco, https://www.youtube.com/watch?v=UT7h4nRcWjU

# Reinforcement learning

- Given a sequence of states and actions with (delayed) rewards, output a policy
  - Policy is a mapping from states → actions that tells you what to do in a given state
- Examples:
  - Credit assignment problem
  - Game playing
  - Robot in a maze
  - Balance a pole on your hand

# The Agent-Environment Interface



Agent and environment interact at discrete time steps $\,:\; t = 0, 1, 2, K$

Agent observes state at step $t:\quad s_t \in S$

produces action at step $t:\quad a_t \in A(s_t)$

gets resulting reward $:\quad r_{t+1} \in \Re$

and resulting next state $:\quad s_{t+1}$

# Reinforcement learning

# ML in a Nutshell

- Tens of thousands of machine learning algorithms

- Hundreds new every year

- Every machine learning algorithm has three components:
  - **Representation**
  - **Evaluation**
  - **Optimization**

# Representation

- Decision trees
- Sets of rules / Logic programs
- Instances
- Graphical models (Bayes/Markov nets)
- Neural networks
- Support vector machines
- Model ensembles
- Etc.

# Evaluation

- Accuracy
- Precision and recall
- Squared error
- Likelihood
- Posterior probability
- Cost / Utility
- Margin
- Entropy
- K-L divergence
- Etc.

# Optimization

- Combinatorial optimization
  - E.g.: Greedy search
- Convex optimization
  - E.g.: Gradient descent
- Constrained optimization
  - E.g.: Linear programming

# Types of Learning

- **Supervised (inductive) learning**
  - Training data includes desired outputs

- **Unsupervised learning**
  - Training data does not include desired outputs

- **Semi-supervised learning**
  - Training data includes a few desired outputs

- **Reinforcement learning**
  - Rewards from sequence of actions

# Inductive Learning

- **Given** examples of a function *(X, F(X))*

- **Predict** function *F(X)* for new examples *X*
  - Discrete *F(X)*: Classification
  - Continuous *F(X)*: Regression
  - *F(X)* = Probability(*X*): Probability estimation

# ML in Practice

- Understanding domain, prior knowledge, and goals
- Data integration, selection, cleaning, pre-processing, etc.
- Learning models
- Interpreting results
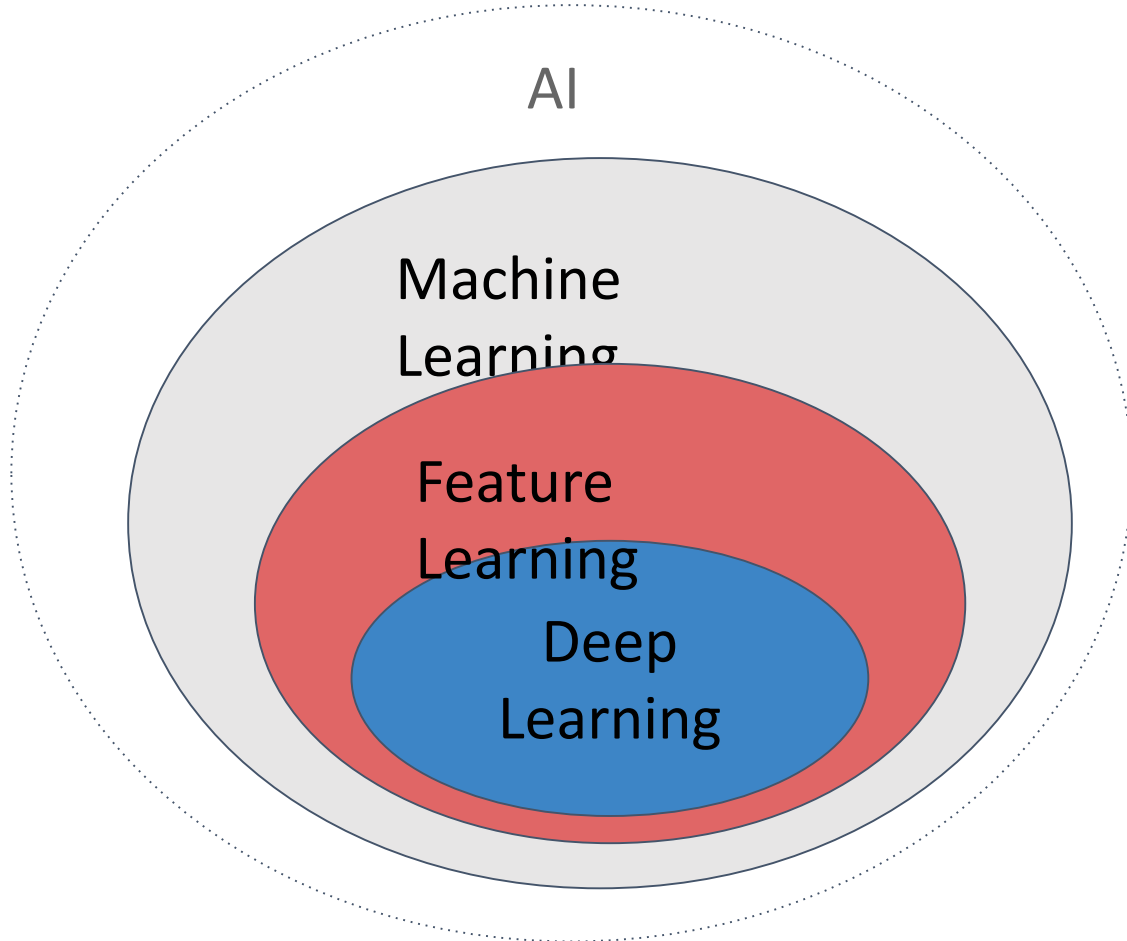- Consolidating and deploying discovered knowledge
- Loop

# Deep Neural Networks

Thanks to: [Deep Learning by Google - Take machine learning to the next level](#)

# What is Deep Learning



**Deep Learning** (DL) has emerged around the '10 as a general <u>tool to solve recognition problems</u> in:

- computer vision
- speech recognition
- robotics
- *discovering* new medicines
- *understanding* natural language
- *understanding* documents
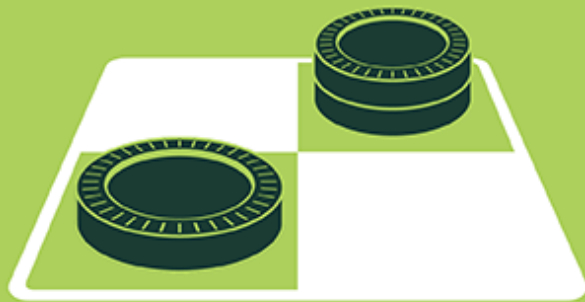- ranking
- … and many other applications!

# Overview

- History
- Preliminaries: logistic classification
- Training
- Deep networks
- Regularization
- Architectures
  - Convolutional networks
  - Embeddings
  - Recurrent models

# History

# ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.

# MACHINE LEARNING

Machine learning begins to flourish.

# DEEP LEARNING

Deep learning breakthroughs drive AI boom.

1950's    1960's    1970's    1980's    1990's    2000's    2010's

Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

# Everything can be optimized in Computer Science

- Given a problem to solve P, it can be formalized as {P,C,F}
    - P := the problem formulation

    - C = {$c_1$,$c_1$,...,$c_n$} := set of configurations, each one of them representing a possible solution to P

    - f:C→R := function which provides a goodness measure of the configuration w.r.t. the problem to be solved

- Casting the problem via minimization means to maximize or minimize the function f in the C space, *independently on the implied meaning of P*
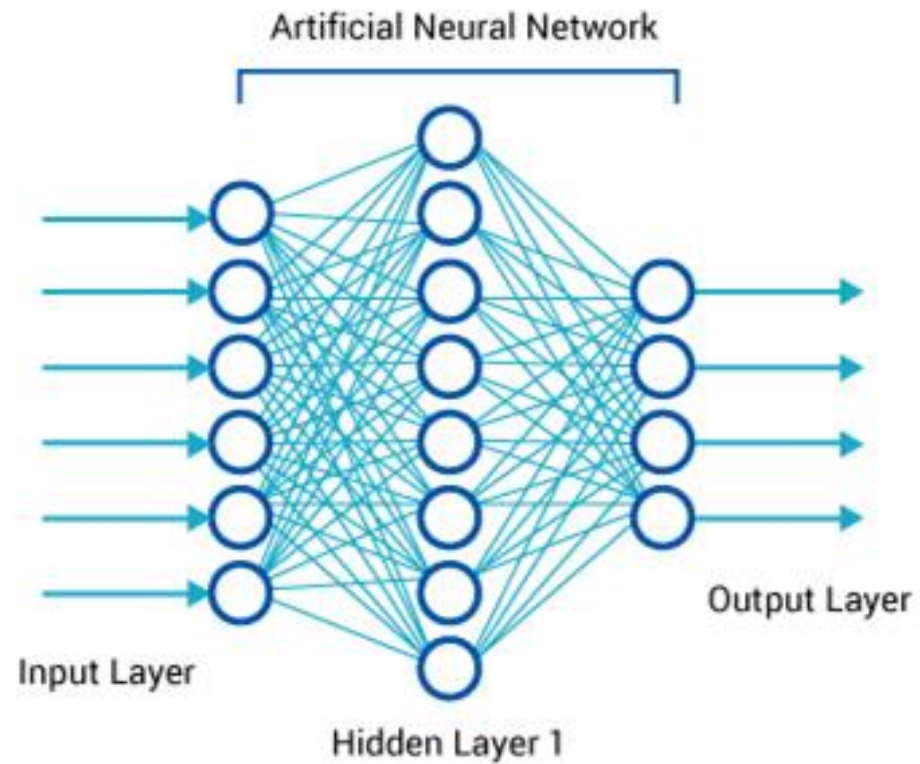
# Minimization: to be used *always*?

- **Problem P$_1$**: *sort numbers* x$_1$,x$_2$,...,x$_N$ *in increasing order*
  - In this case, minimization could be left apart
  - In facts, there is at least one algorithm (e.g., quicksort) which brings directly to the best (in sense of the function f) configuration
- **Problem P$_2$**: foresee the stocks' trend
  - Much more difficult to formalize into an algorithm
  - Minimization comes to help [Yong et al. 2015]

[Yong et al. 2015] Hu, Yong, et al. "Application of evolutionary computation for rule discovery in stock algorithmic trading: A literature review." *Applied Soft Computing* 36 (2015): 534-551.

# Inside the minimization approach

- The main goal of an optimization approach is that of exploring the configuration space C looking for the best configuration given the function f (obviously avoiding the brute force way!)
- The set of configurations C give a space to explore (very often, a manifold)
- Optimizing means to explore the manifold by iterative approaches (e.g., the gradient descent family of strategies)
- The more the manifold is complex (non convex, multimodal), the more often local minima are met

# Neural Networks [1943 - McCulloch & Pitts]

- Optimization approaches which scale *very well* with data
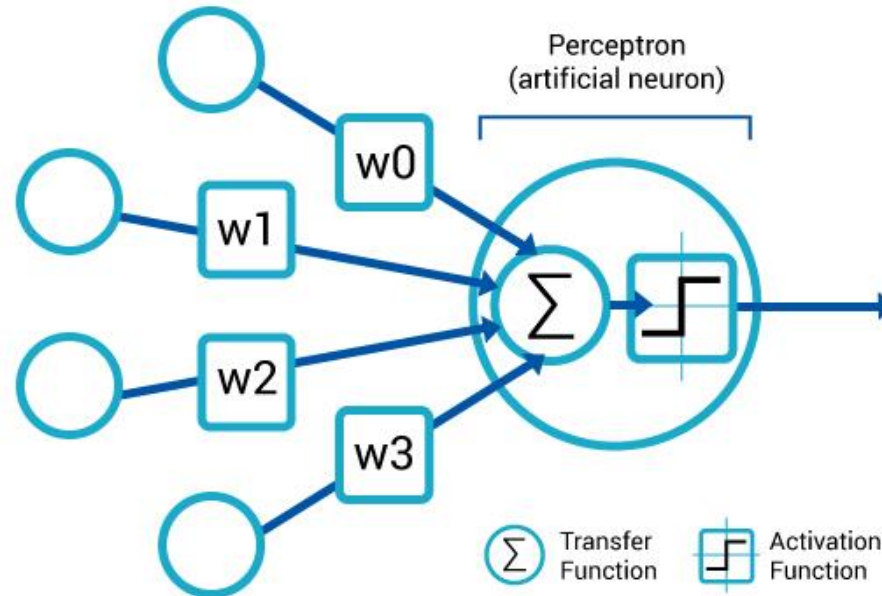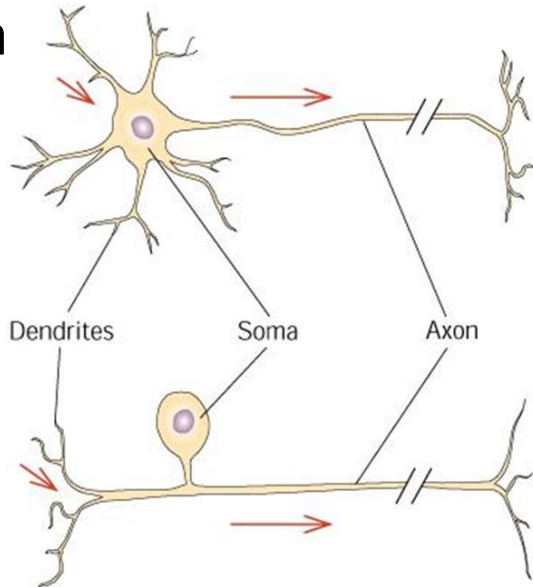- We are talking about *artificial neurons* and *layered computation*

Artificial Neural Network

Input Layer

Hidden Layer 1

Output Layer

# Neural Networks - Neurons

NN are composed by artificial neurons (1943 - McCulloch & Pitts).
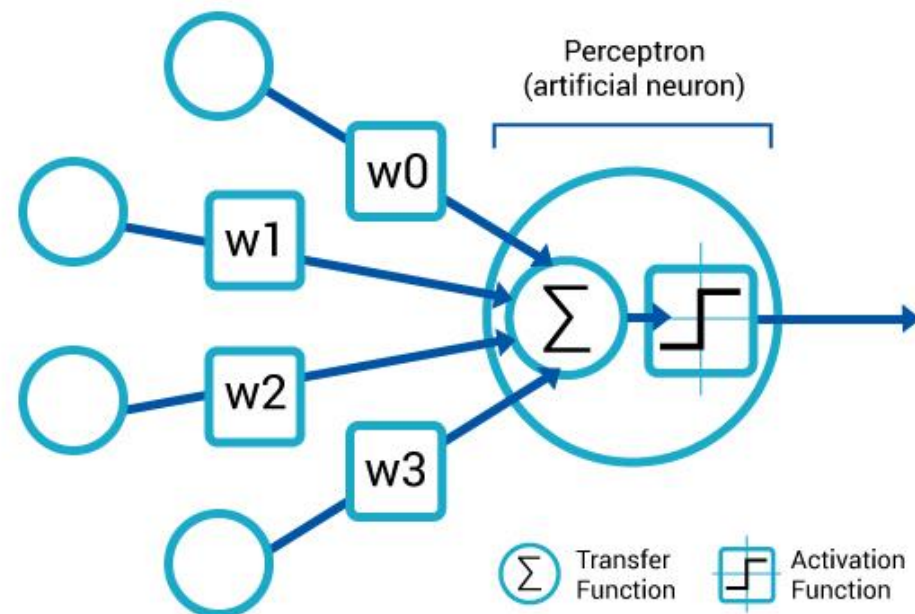
Each neuron has:

- *dendrites* (**inputs**)
- a *nucleus/soma/perceptron* (**transfer fuction + activation function**)
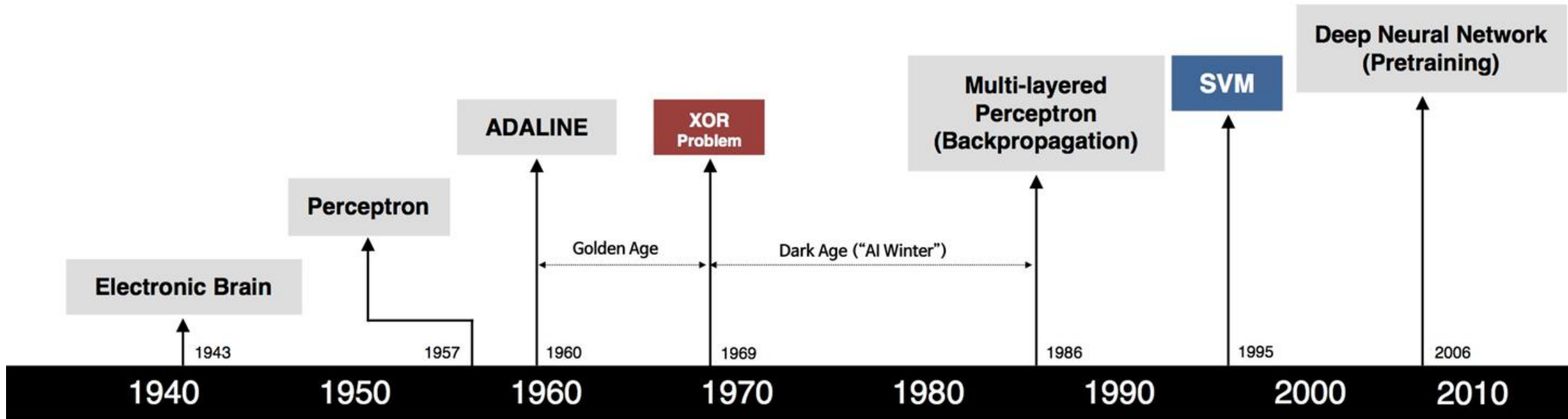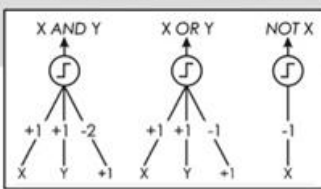- a

# Neural Networks - Neurons (2)

The information flow is **unidirectional**:

- The neuron get **inputs** (electric potentials) from the *dendrites*, that weight them ($w_i$'s)
- In the *nucleus*, the weighted inputs are summed together (the **transfer** $\Sigma$ of the whole information coming from the dendrites)
- In the *nucleus*, the summation flows into an **activation function**, which may inhibit, diminish or amplify it
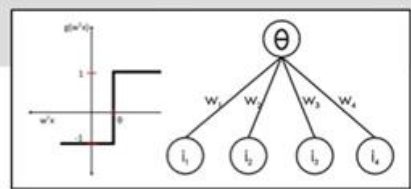- The **output** of the activation function

A History of Deep Learning — timeline from 1940 to 2010.

**Electronic Brain** — 1943

**Perceptron** — 1957

**ADALINE** — 1960

**XOR Problem** — 1969

**Multi-layered Perceptron (Backpropagation)** — 1986

**SVM** — 1995

**Deep Neural Network (Pretraining)** — 2006

Golden Age | Dark Age ("AI Winter")

1940 · 1950 · 1960 · 1970 · 1980 · 1990 · 2000 · 2010

S. McCulloch – W. Pitts
- Adjustable Weights
- Weights are not Learned

F. Rosenblatt
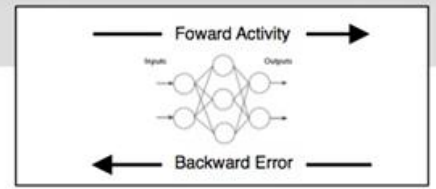- Learnable Weights and Threshold
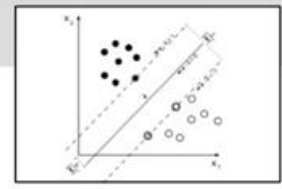
B. Widrow – M. Hoff

M. Minsky – S. Papert
- XOR Problem

D. Rumelhart – G. Hinton – R. Wiliams
- Solution to nonlinearly separable problems
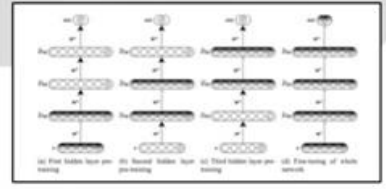- Big computation, local optima and overfitting

V. Vapnik – C. Cortes
- Limitations of learning prior knowledge
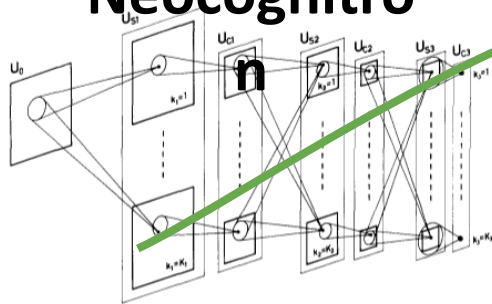- Kernel function: Human Intervention

G. Hinton – S. Ruslan
- Hierarchical feature Learning
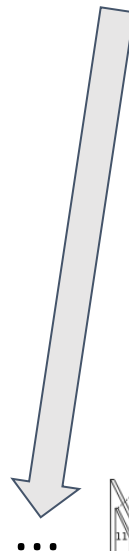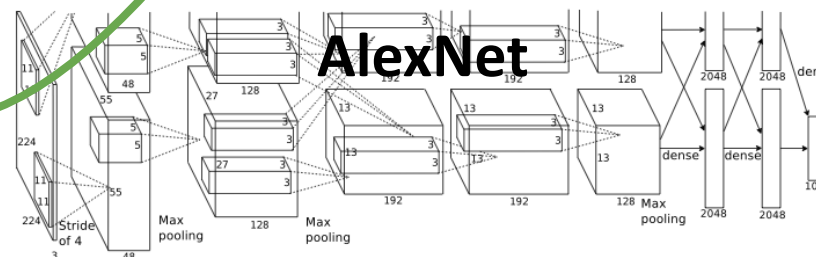
# Neural Networks - the renaissance



Le Cun's
**LeNet-5**

**What happened?**

Fukushima's
**Neocognitron**

Krizhevsky's
**AlexNet**

...

1980        1990        2000        2010        2020

http://people.csail.mit.edu/torralba/research/drawCNN/drawNet.html
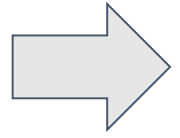
# Neural Networks and Deep Learning
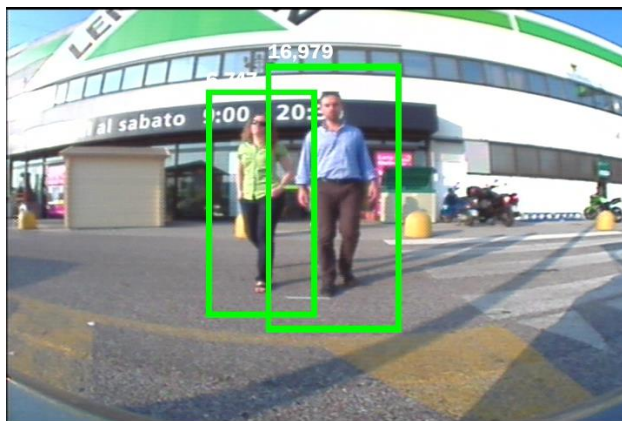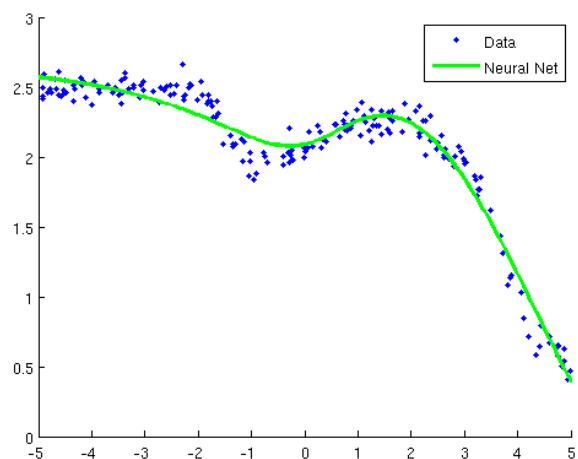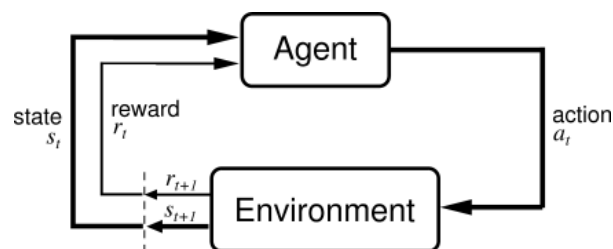


Neural Networks → GPUs + Data → **Deep** Neural Networks

# Supervised Classification

- Traditional kind of problem the NN do solve
  - Regression
  - Ranking
  - Reinforcement Lea

Labels {'a', 'b', 'c', 'd', 'e'}

# Preliminaries: logistic classification

# Logistic Classifier

- It assigns a *score y* to the input **x** through a linear model (W,**b**)
- The score helps to identify the class label that wins

$$\mathbf{x}W + \mathbf{b} = \mathbf{y}$$

1xF  FxC        1xC        1xC

To be trained via a training procedure

*x = input or feature vector; **F** = number of features; **W** = weights matrix;*
*C = number of classes **b** = biases; **y** = output or logits/scores vector*

# Logistic Classifier - the score is not enough

$$x W + b = y \implies s(y)$$



$$\begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix}$$

scores

$$\begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$$

probabilities

$s(y)$ is the **SOFTMAX** function

# Softmax function

- Converts scores into probability
  distributions
    - $\mathbb{R} \rightarrow (0,1)$
    - *Open* codomain!
- The softmax function highlights the
  largest values and suppress values which
  are significantly below the maximum
  value



$$\begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix} \implies s(y_i) = \frac{e^{y_i}}{\sum_{j=1}^{C} e^{y_j}} \implies \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$$
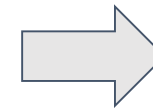
scores

probabiliti
es

# One-Hot Encoding (OHE)

- → *There is only one correct label for each input sample*
- → *There is the need to evaluate the classification result*
- OHE encodes labels for a C-class problem in $R^C$, indicating the c-th class label with 1, the rest 0's

  needs representation ... icient especially in the case C is very big (thousands or millions of classes...!)

**Which distance measure?**

MSE

Cross-Entropy

$S(\mathbf{y})$

$$\begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$$

our results

?

$$\begin{bmatrix} 1.0 \\ 0.0 \\ 0.0 \end{bmatrix}$$

g.t. labels

# Cross-Entropy

GOAL: computes the distance between two probability vectors

- Non symmetric function $D(\mathbf{s}, \mathbf{L}) \neq D(\mathbf{L}, \mathbf{s})$

$D(\mathbf{s},\mathbf{L})$
$$\begin{cases} \to 0 & \text{high similarity} \\ \\ \to \infty & \text{low similarity} \end{cases}$$

- The *log(·)* makes the training faster to converge than the other alternative MSE function $(\sum(s(y)-l)^2)$

$s(\mathbf{y})$

$\begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$

$\mathbf{L}$

$\begin{bmatrix} 1.0 \\ 0.0 \\ 0.0 \end{bmatrix}$

$$D(\mathbf{s}, \mathbf{L}) = - \sum_{j=1}^{C} L_i \log s_i \quad = 0.36$$

# Résumé



**Multinomial Logistic Classification**

$$D(s(\mathbf{x}W+\mathbf{b}),\ \mathbf{L})$$

# Training

# Gradient Descent

GOAL: search for the nearest local minimum of a function F

IDEA: iterate on the parameter set proportionally to the negative of the function gradient

$$\theta_{t+1} = \theta_t - \alpha \nabla F(\theta_t), t \geq 0$$

such that

$$F(\theta_0) \geq F(\theta_1) \geq F(\theta_2) \geq \ldots$$

*$\vartheta$ = set of parameters*

# Gradient Descent

- <u>GOAL</u>: minimize a loss function

- Needs to compute the entire training set performance of our linear model, that consists in N inputs (which is, in general, very big)

- Needs to minimize a loss function, which depends on W (big matrix) and **b**

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^{N} D(s(\mathbf{x}_i W + \mathbf{b}), \mathbf{L}_i)$$

Loss = average cross-entropy

$$D(\mathbf{s},\mathbf{L}) \begin{cases} \to 0 & \text{good} \\ \\ \to \infty & \text{bad} \end{cases}$$

$$\mathcal{L}(w_i, w_j)$$

$w_j$

learning step

$$-\alpha \nabla \mathcal{L}(w_i, w_j)$$

W

# Stochastic Gradient Descent (SGD)

- <u>IDEA</u>: use a random subset (*batch*) of the data (of a given *size*) to compute an approximation of the gradient of the loss function
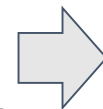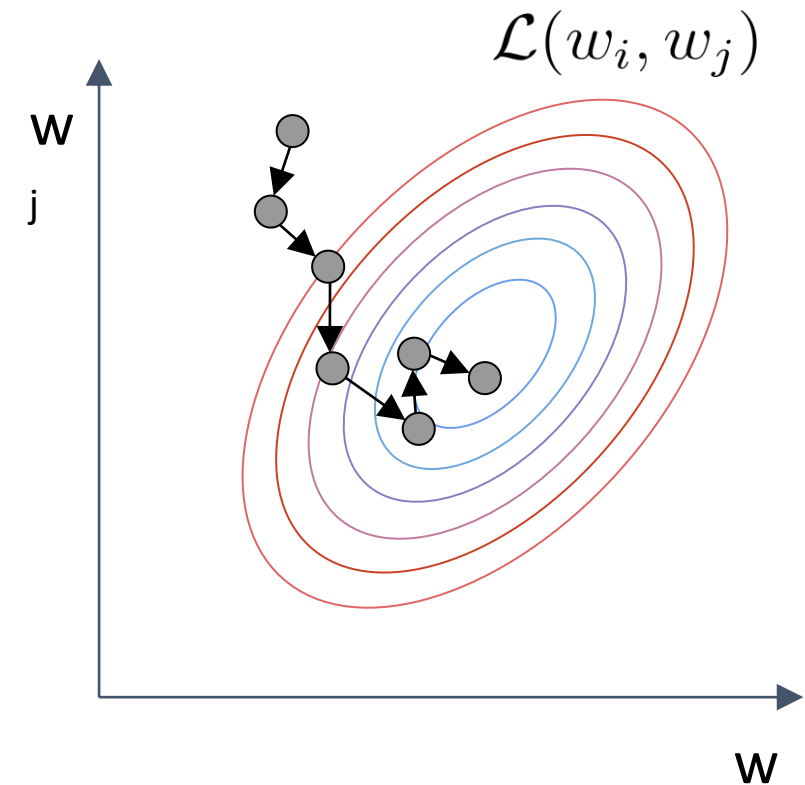
Iterative implementation of the GD algorithm
At each step, a new batch is extracted

- Pros:
  - simple but sufficiently effective
  - fast (depending on batch size)
  - scale the problem with data and model size
- Cons:
  - needs more iterations to converge

$\mathcal{L}(w_i, w_j)$

$w_j$

$w_i$

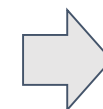... but tricks to ameliorate SGD are present!

# SGD trick 1: momentum

- <u>GOAL</u>: improve the convergence of the optimizer exploiting the accumulated knowledge from previous steps

- <u>IDEA</u>: add a fraction of the previous update vector to the current update vector

$$M_0 = \nabla \mathcal{L}_0(w_i, w_j) = 0$$

$$M_t = \alpha \nabla \mathcal{L}_t(w_i, w_j) + \beta M_{t-1}$$

$$(w_i, w_j)_t = (w_i, w_j)_t - M_t$$



$\mathcal{L}(w_i, w_j)$

w$_j$

w$_i$

faster convergence and oscillation reduction

*$M$ = momentum; $\beta$: friction value (usually 0.9)*

# SGD trick 2: learning-rate decay

- GOAL: make the optimization more robust and accurate over time

- IDEA: apply a decay function to the learning rate or reduce it when the loss function reaches a plateau

$$\mathcal{L}(w_i, w_j)$$

w
j

w
i

α

tim

# SGD trick 3: z-normalization

- <u>GOAL</u>: avoid numerical instability

The values involved in the calculation of the gradient descent never get too big or too small



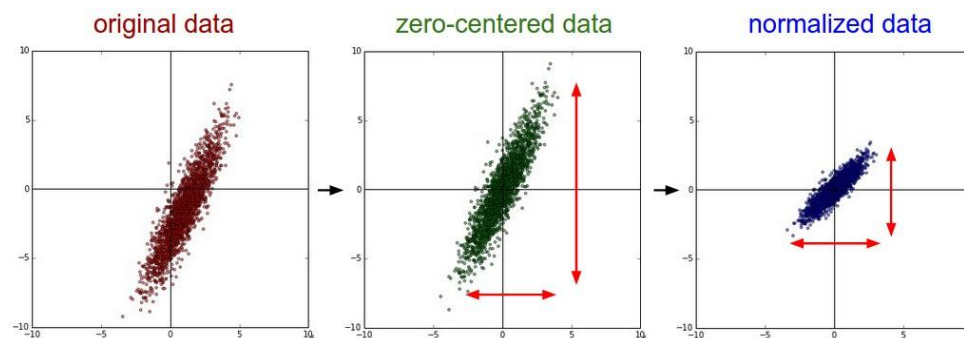original data     zero-centered data     normalized data

- <u>IDEA</u>: remove the mean and normalize over the variance of the *i*-th feature of the input vector $\mathbf{x}$

$$\hat{\mathbf{x}} = \frac{\mathbf{x} - E[\mathbf{x}]}{Var[\mathbf{x}]}$$

**z-normalization**

$$E[\mathbf{x}] = 0$$

$$\forall (i,j), Var[\mathbf{x}_i] = Var[\mathbf{x}_j]$$

# SGD trick 4: initialization

A random initialization of the weights and the zero-init of the biases is critical to get a good starting point for the training phase and the convergence of the SGD algorithm.

$$w_i = \mathcal{N}(\mu = 0, \sigma \to 0)$$

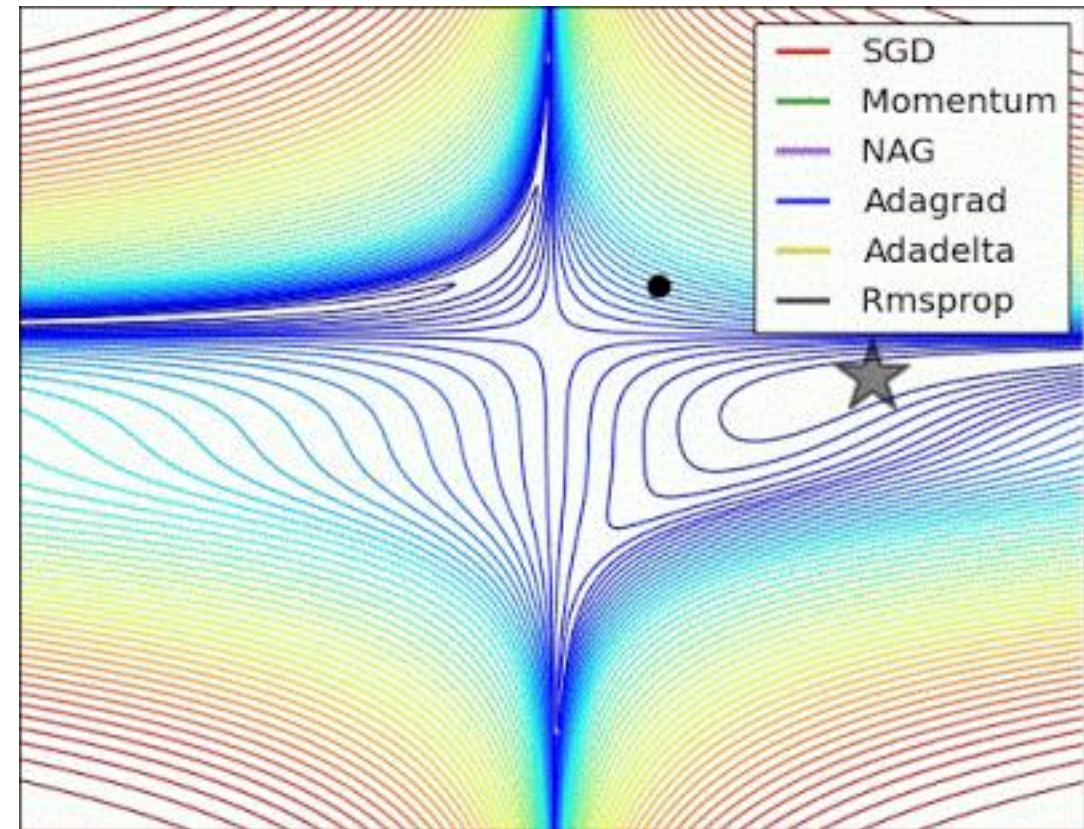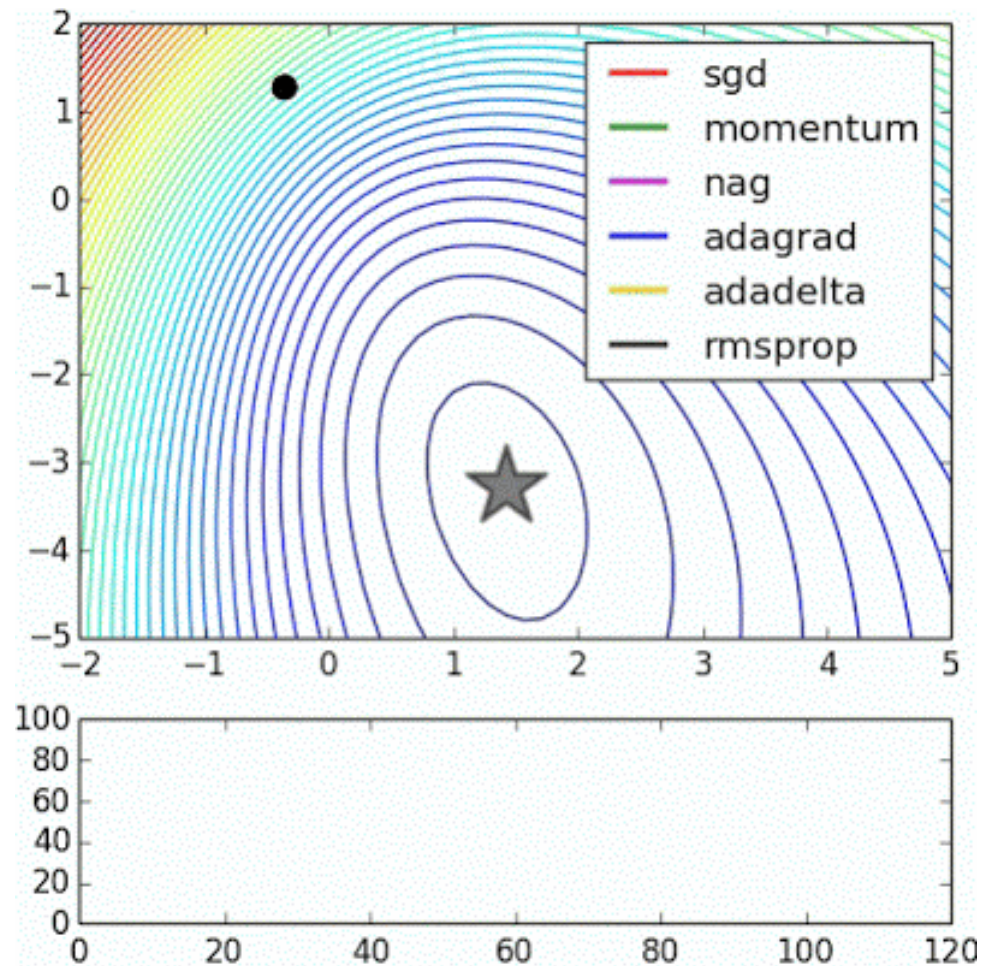$$\mathbf{b} = \begin{bmatrix} 0_1 & 0_2 & \dots & 0_C \end{bmatrix}$$

$$\mu = 0 \wedge \sigma \to 0 \implies \forall (i, j), w_i = w_j \pm \epsilon$$
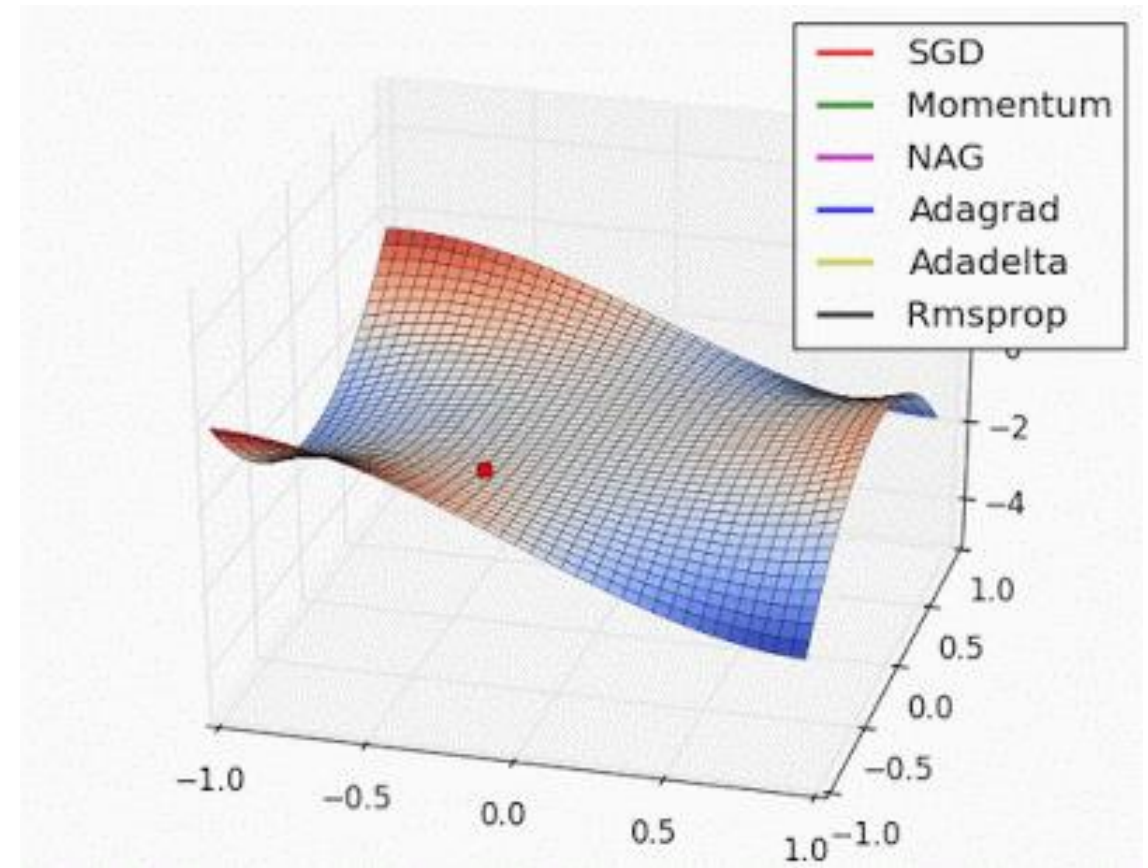
⬇

equal probability
of the weights
(no prior)

# Gradient Descent: graphical representation (2D)

# Gradient Descent: graphical representation (3D)

# SGD: tuning

**SGD**

Many hyperparameters:

- initial learning rate
- learning rate decay
- momentum
- batch size
- weight initialization



**AdaGrad**

SGD modification which implicitly applies

momentum and learning rate decay. It uses

**SGD**: H. Robinds and S. Monro, "A stochastic approximation method," Annals of Mathematical Statistics, vol. 22, pp. 400–407, 1951.

**AdaGrad**: J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online leaning and stochastic optimization," in COLT, 2010.

fewer parameters:

# From MLC to NN to Deep-NN

- **Multinomial Logistic Classification** (MLC)
  - *fast*: efficient computation due to the linear model
  - *stable*: small input variations generates small output variations and the derivative of the model is constant

but...
  - trains only a "small" parameter set
  - isn't responsive to non-linear inputs

- What if I concatenate multiple MLCs?

still a linear model...!

$$D(s(\mathbf{x}\mathbf{W}+\mathbf{b}), \mathbf{L})$$

$x_1 + x_2$

$x_1 * x_2$

| $\mathbf{x}$ | $W_1$ | + | $\mathbf{b}_1$ | = | $\mathbf{y}_1$ | $W_2$ | + | $\mathbf{b}_2$ | = | $\mathbf{y}_2$ |

| $\mathbf{x}$ | $W_3$ | + | $\mathbf{b}_3$ | = | $\mathbf{y}_3$ |

# From MLC to NN to Deep-NN

- GOAL: build a bigger and non linear model

- IDEA: concatenate many linear systems (MLC) and insert a **non linear function** between two consecutive MLCs, that is, the activation function
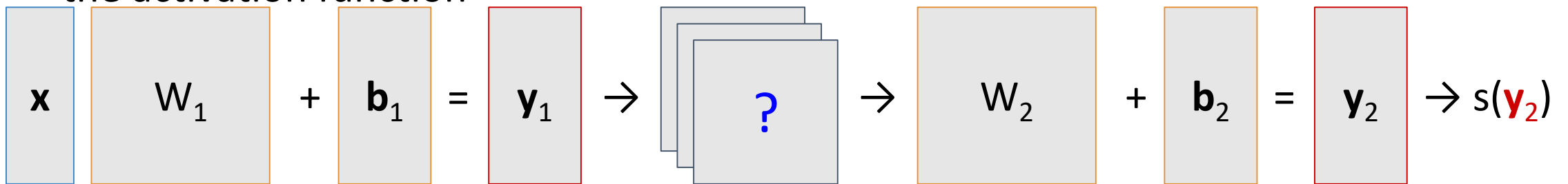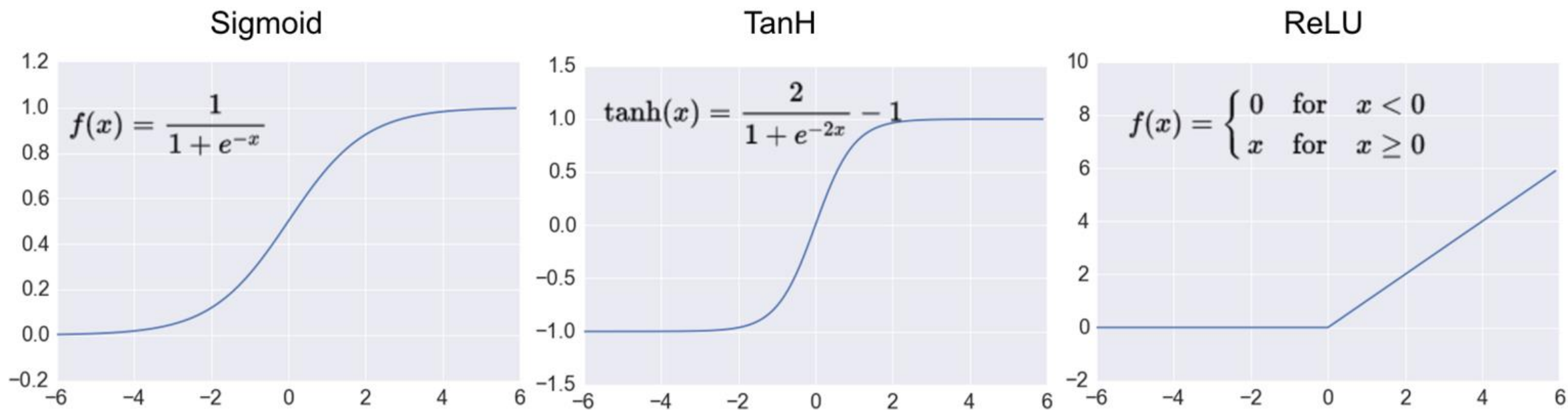
**Which function?**

ReLU    sigmoid

tanh



$\mathbf{x}$ | $W_1$ | + | $\mathbf{b}_1$ | = | $\mathbf{y}_1$ | → | ? | → | $W_2$ | + | $\mathbf{b}_2$ | = | $\mathbf{y}_2$ | → s($\mathbf{y}_2$)

# Choose the (non linear) activation function

### Sigmoid

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

### TanH

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

### ReLU

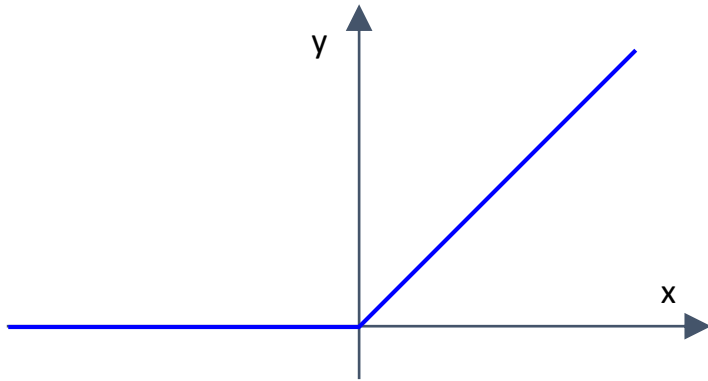$$f(x) = \begin{cases} 0 & \text{for} \quad x < 0 \\ x & \text{for} \quad x \geq 0 \end{cases}$$
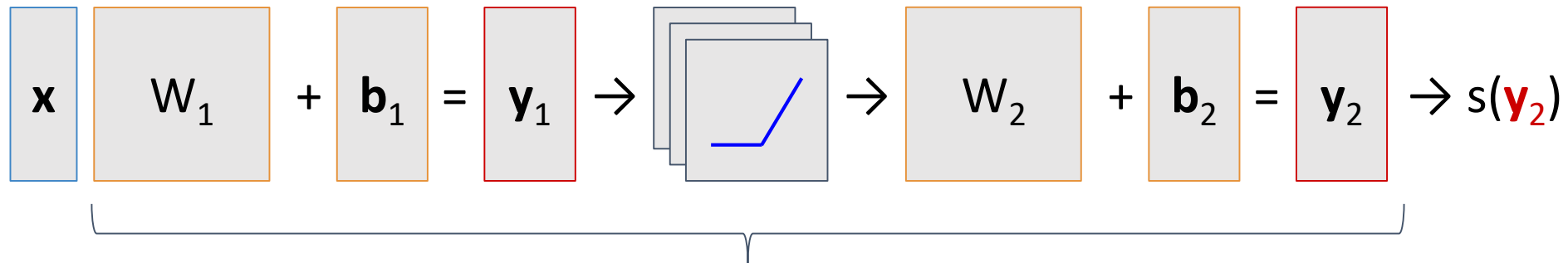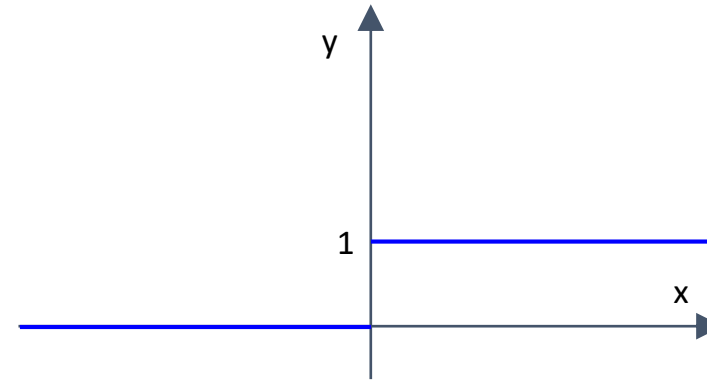
- A perceptron classic

- similar to sigmoid

- simplest function
- constant derivative

# ReLU: Rectified Linear Unit

$$ReLU(x) = max(0, x)$$

derivative



$$\boxed{\mathbf{x}} \; \boxed{W_1} \; + \; \boxed{\mathbf{b}_1} \; = \; \boxed{\mathbf{y}_1} \; \rightarrow \; \boxed{\diagup} \; \rightarrow \; \boxed{W_2} \; + \; \boxed{\mathbf{b}_2} \; = \; \boxed{\mathbf{y}_2} \; \rightarrow \; s(\mathbf{y}_2)$$

no more linear…!

# Résumé



non linear model

linear model

softmax

$\mathbf{x}$ | $W_1$ | + | $\mathbf{b}_1$ | = | $\mathbf{y}_1$ | $\rightarrow$ | H | $\rightarrow$ | $W_2$ | + | $\mathbf{b}_2$ | = | $\mathbf{y}_2$ | $\rightarrow$ s($\mathbf{y}_2$)

hidden layer

not very deep yet...

*H = number of "neurons" in the hidden layer*

# Practical example (1/2)



Image 28x28 pixels

#classes (C)

784

3

#features (F)

$x \quad * W1 \quad + b1$

$x \quad W_1 \quad + \quad b_1 \quad = \quad y_1 \quad \xrightarrow{} \quad s(y_1)$

1xF      FxC      1xC      1xC

# Practical example (2/2)



#hidden_neurons
(H)

$$\mathbf{x} * W_1 + \mathbf{b}_1 = \mathbf{y}_1 \rightarrow \text{ReLU}(\mathbf{y}_1) * W_2 + \mathbf{b}_2 = \mathbf{y}_2$$

# Examples of deep networks - conclusions