



Master DSAI

Advanced Deep Learning

2022-2023

From Artificial Intelligence to a first neuron model then to MLP, CNN and RNN

Frederic Precioso
(INRIA – CNRS – UCA, Team MAASAI)

Disclaimer

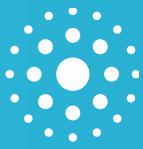
If any content in this presentation is yours but is not correctly referenced or if it should be removed, please contact me and I will correct it.

Overview

- **Context & Vocabulary**
- Math Basics
- Simple Models
- Spatial Invariance (CNN)
- Sequential Invariance (RNN)

Overview

- **Context & Vocabulary**
 - *What is Artificial Intelligence?*
 - *Machine Learning without Maths*
 - *Machine Learning VS Data Mining*
 - *Machine Learning VS Data Science*
 - *Machine Learning VS Statistics*
- Math Basics
- Simple Models
- Spatial Invariance (CNN)
- Sequential Invariance (RNN)



CONTEXT & VOCABULARY



What is Artificial intelligence?



How is Artificial intelligence defined?

univ-cotedazur.fr

- The term ***Artificial Intelligence***, as a research field, was coined at the conference on the campus of Dartmouth College in the summer of **1956**, even though the idea was around since Antiquity: Hephaestus built automatons of metal to work for him or protect others (as in the picture on the right *Talos protecting goddess Europa*), the Golem made of clay in Jewish folklore (on the left), Yan Shi has built humanoid automatons 10th century BC during Zhou dynasty, etc.



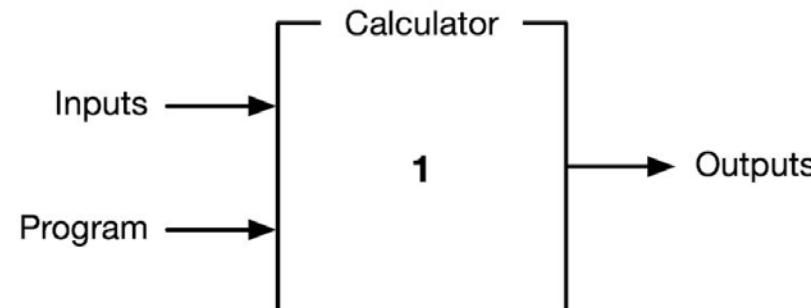
- Closer to the Dartmouth conference but still before, the first manifesto on Artificial Intelligence, an unpublished report "***Intelligent Machinery***", written by Alan Turing in **1948**. He already distinguished two different approaches to AI, which may be termed "***top-down***" and "***bottom-up***" (*now more commonly called knowledge-driven AI and data-driven AI respectively*).

(sources: Wikipedia, <https://www.greeklegendsandmyths.com/automatons.html> , Stanford Encyclopedia of Philosophy: <https://plato.stanford.edu/entries/artificial-intelligence/>, http://www.alanturing.net/turing_archive/pages/Reference%20Articles/what_is_AI/What%20is%20AI02.html)



How is Artificial intelligence defined?

- “*top-down*” or *knowledge-driven AI*
 - cognition = high-level phenomenon, independent of low-level details of implementation mechanism
 - Evolutionary Algorithms (1954, 1957, 1960), Knowledge Representation, Reasoning (1959, 1970), Expert Systems (1970), Logic, Automata, Intelligent Agent Systems (1990)...



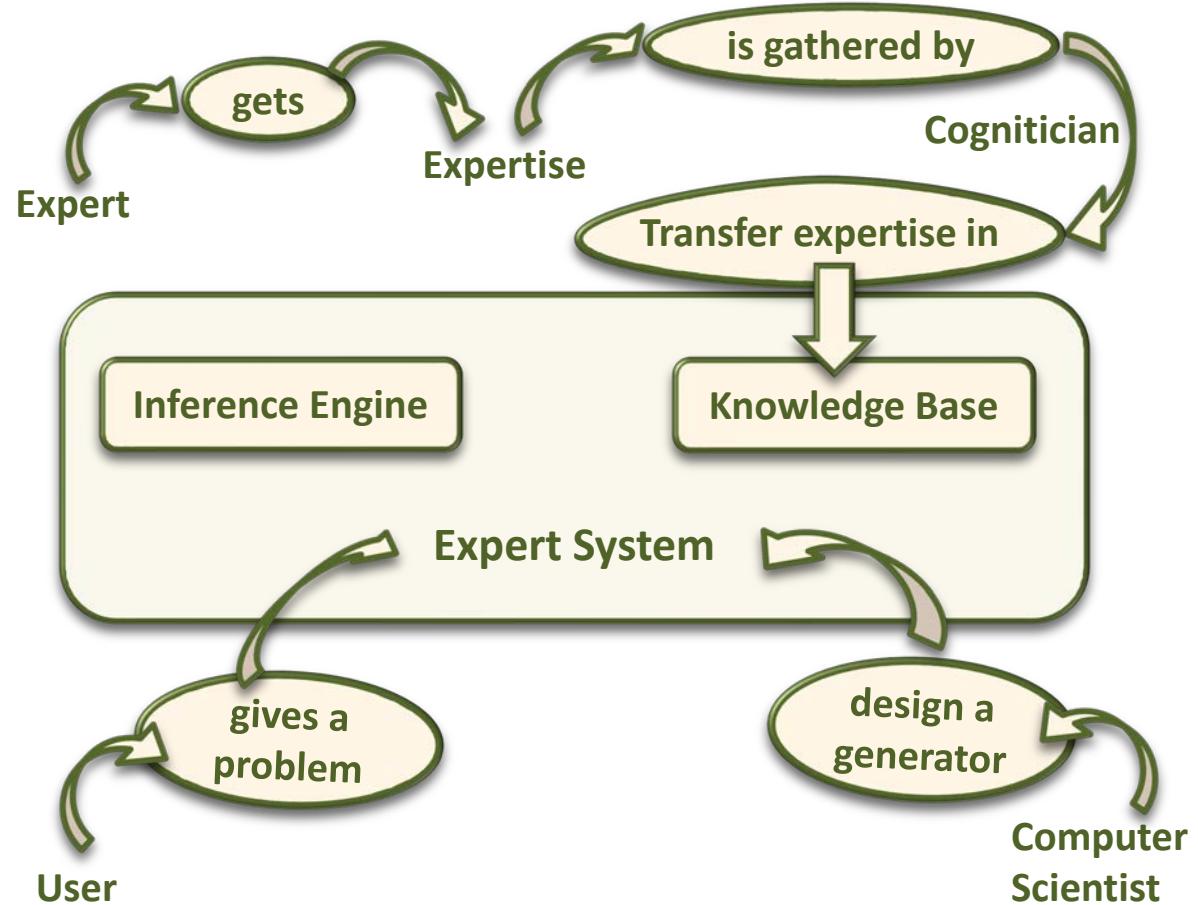
(1) Hypothetical-deductive machines

(Figure from: *Neurons spike back The invention of inductive machines and the artificial intelligence controversy*”, D. Cardon, J.-P. Cointet, A. Mazières, Translated by Elizabeth Libbrecht In Réseaux Volume 211, Issue 5, 2018, pages 173 to 220)



Artificial Intelligence, Top-Down

- Example of an expert system:



Artificial Intelligence, Top-Down

- Expert system:

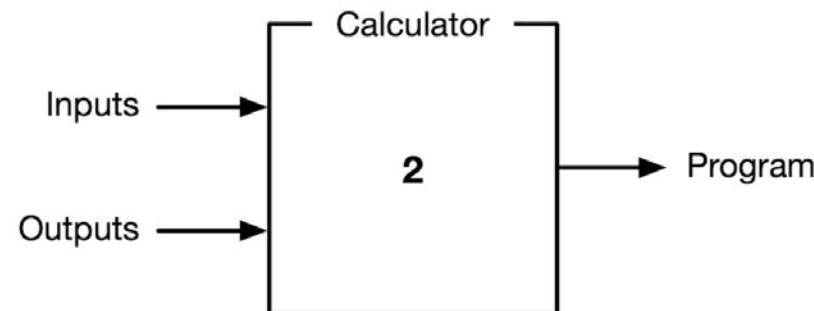
The screenshot shows a news article from RBR Insider. The header features the RBR logo (an orange 'rb' icon) and navigation links for AI, Manufacturing, Supply Chain, Robo Dev, Healthcare, CRO, Events, and All Topics. On the right, there are 'Login | Join RBR Insider' and a search icon. The main title of the article is 'A Cyclist's Encounter with an Indecisive Google Self-Driving Car'. Below the title is a subtitle: 'A bicyclist recently had a two-minute standoff with a Google self-driving car at a four-way stop in Austin, Texas. So what happened? We explain.' The byline indicates it was published on AUGUST 26, 2015, by STEVE CROWE.

The screenshot shows a news article from STAT. A red 'BRIEF' button is visible in the top left corner. The main title is 'STAT: IBM's Watson gave 'unsafe and incorrect' cancer treatment advice'. Below the title, the author is listed as Meg Bryant and the date as PUBLISHED July 26, 2018. The 'SHARE IT' section is partially visible at the bottom. The 'Dive Brief:' section contains two bullet points: 'A STAT review of internal IBM documents suggests the company's Watson supercomputer wrongly advised doctors on how to treat patients' cancers.' and 'The documents — slides presented by then-IBM Watson Health deputy chief health officer Andrew Norden in June and July of last year — include "multiple examples of unsafe and incorrect treatment".'



How is Artificial intelligence defined?

- **"bottom-up" or data-driven AI**
 - opposite approach, start from data to build incrementally and mathematically mechanisms taking decisions
 - First neuron (1943), first neural network machine (1950), neucognitron (1975), Decision Trees (1983), Backpropagation (1984-1986), Random Forest (1995), Support Vector Machine (1995), Boosting (1995), Deep Learning (1998/2006)...



(2) inductive machines

(Figure from: *Neurons spike back The invention of inductive machines and the artificial intelligence controversy*", D. Cardon, J.-P. Cointet, A. Mazières, Translated by Elizabeth Libbrecht In Réseaux Volume 211, Issue 5, 2018, pages 173 to 220)



How is Artificial intelligence defined?

- *AI is originally defined in 1956, by Marvin Lee Minsky:*
*“The construction of computer programs doing tasks, that are, **for the moment**, accomplished **more satisfactorily** by human beings because they require **high level mental processes** such as: learning. perceptual organization of memory and critical reasoning”.*
- There are so the “artificial” side with the usage of computers or sophisticated electronic processes and the side “intelligence” associated with its goal to imitate the (human) behavior.

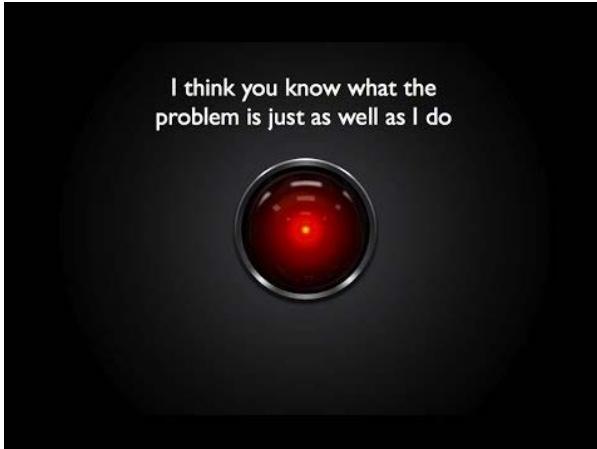
Why Artificial Intelligence is so difficult to grasp?

- Frequently, when a technique reaches **mainstream use**, it is **no longer considered as artificial intelligence**; this phenomenon is described as the *AI effect*: "AI is whatever hasn't been done yet." (*Larry Tesler's Theorem*)
-> e.g. Path Finding (GPS), Chess electronic game, Alpha Go...
- Consequently, AI domain is continuously evolving and so very difficult to grasp.



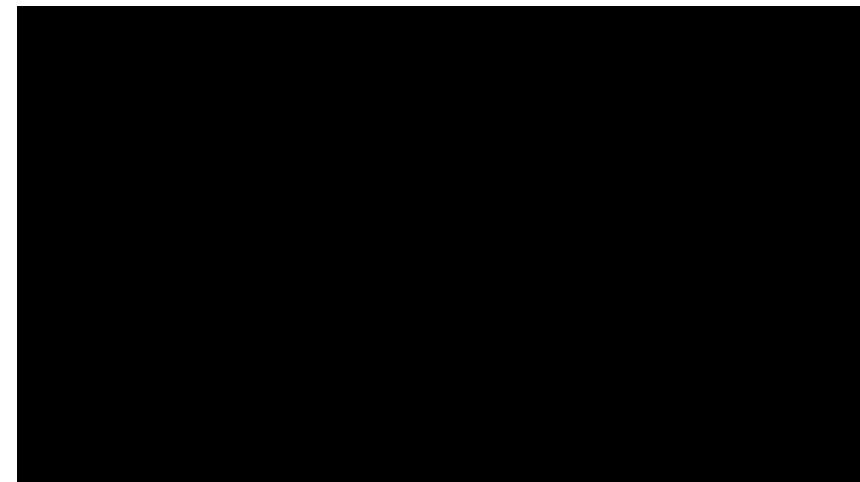
How is Artificial intelligence defined?

- The concept of ***strong artificial intelligence*** makes reference to a machine capable not only of producing intelligent behavior, but also to experience a feeling of a real sense of itself, “real feelings” (whatever may be put behind these words), and "an understanding of its own arguments".



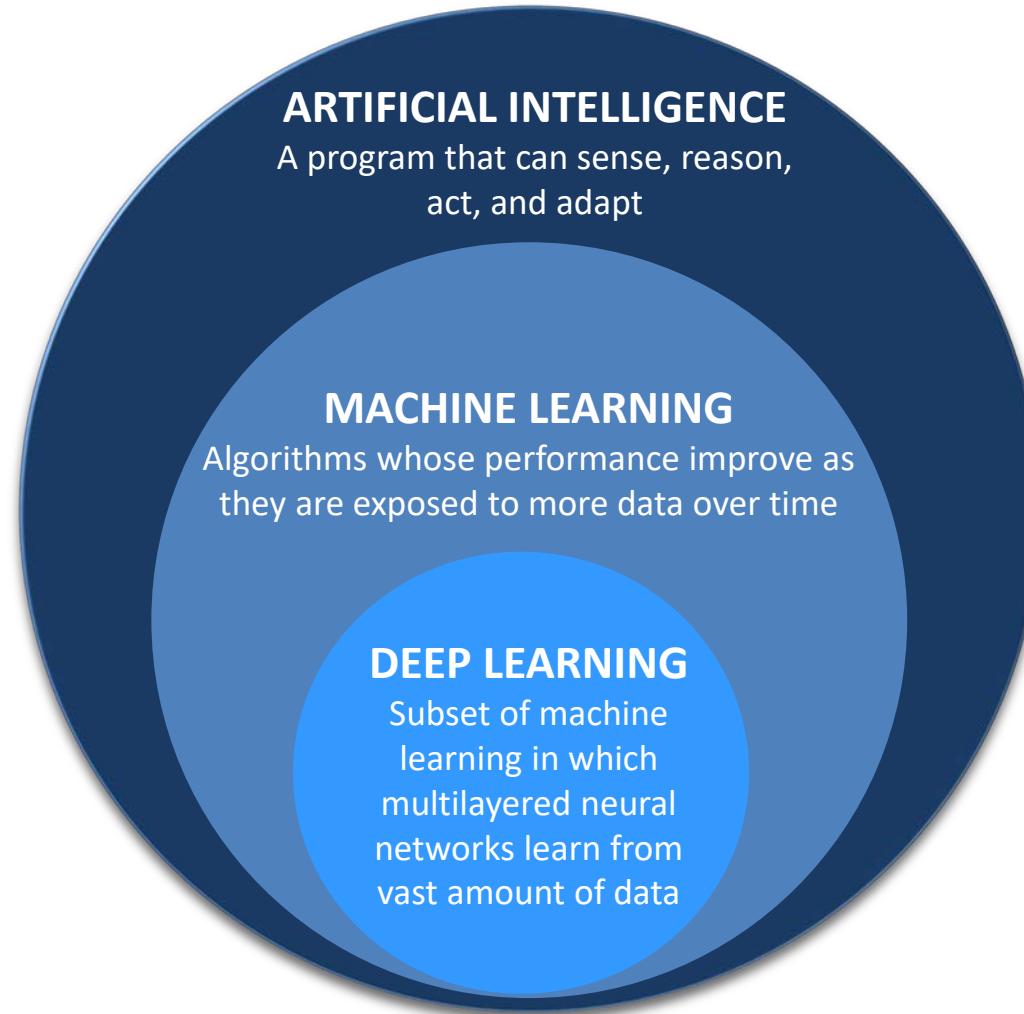
How is Artificial intelligence defined?

- The notion of ***weak artificial intelligence*** is a pragmatic approach of engineers: targeting to build more autonomous systems (to reduce the cost of their supervision), algorithms capable of solving problems of a certain class, etc. But this time, the machine *simulates* the intelligence, it seems to act as if it was smart.





AI vs Machine Learning vs Deep Learning





But what is Machine Learning?

Machine Learning

$$\begin{pmatrix} \mathbf{x} \end{pmatrix} \xrightarrow{f(\mathbf{X}, \alpha)} y$$

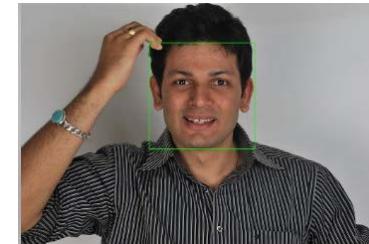
\mathbf{x}



Face detection



y



Scores prediction



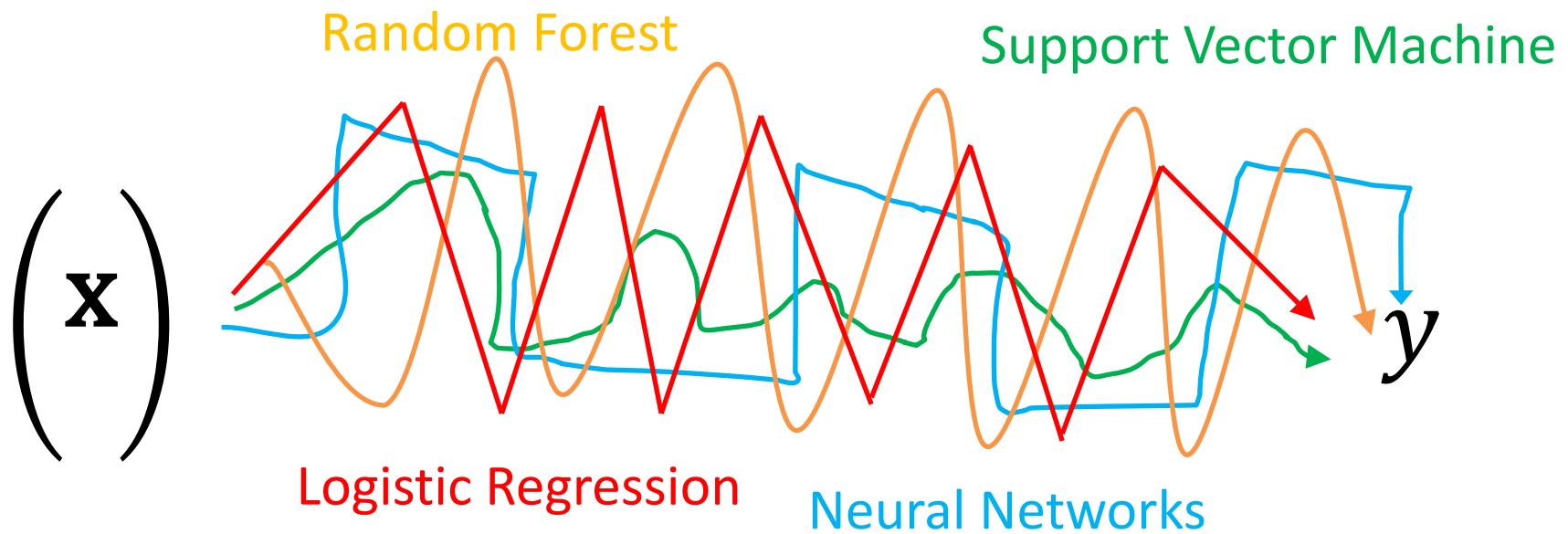
Sport bets

Voice recognition



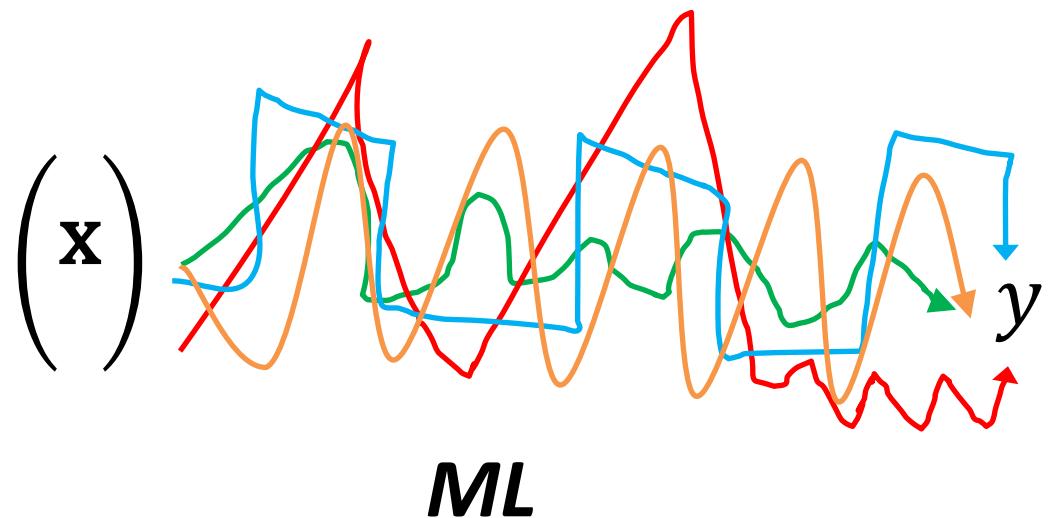
Machine Learning

$$\begin{pmatrix} \mathbf{x} \end{pmatrix} \xrightarrow{f(\mathbf{X}, \alpha) ?} y$$



Machine Learning

$$\begin{pmatrix} \mathbf{x} \end{pmatrix} \xrightarrow{f(\mathbf{X}, \alpha) ?} y$$

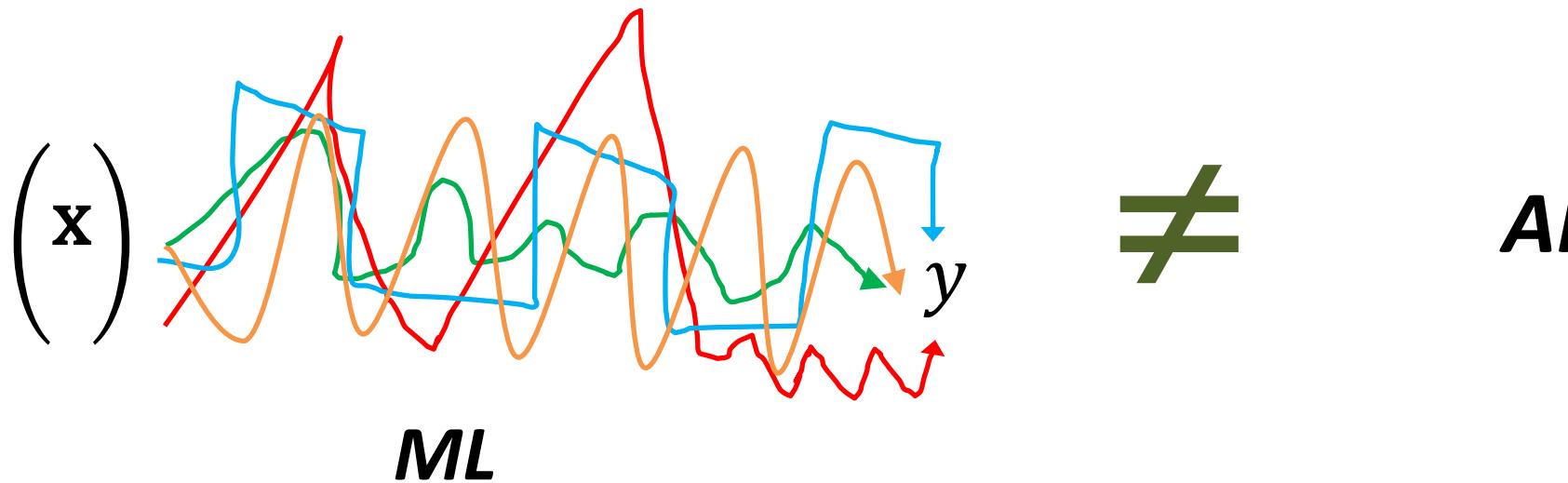


“Weather Forecasting”

Cf. Question trolley dilemma

Machine Learning

$$\begin{pmatrix} \mathbf{x} \end{pmatrix} \xrightarrow{f(\mathbf{X}, \alpha) ?} y$$



Francis Bach at *Frontier Research and Artificial Intelligence Conference*: “Machine Learning is not AI”

(https://erc.europa.eu/sites/default/files/events/docs/Francis_Bach-SEQUOIA-Robust-algorithms-for-learning-from-modern-data.pdf

<https://webcast.ec.europa.eu/erc-conference-frontier-research-and-artificial-intelligence-25#>)

Beware of the diversion!



Trolley dilemma

Back to weather forecast

- **The impact of summer weather on our economy!**
 - The weather influences the sales of products and impacts the revenues of companies in various sectors.
 - According to a study by Credoc, **1°C above seasonal norms leads to a 1.5% increase in tourist numbers**, which leaves professionals in this sector on the lookout for rain or shine.
 - More specifically, in the hotel sector, **a drop of only 1°C compared to seasonal norms in June, for example, will reduce sales by 8 to 10%**.

Back to weather forecast

- Should we penalize algorithms that predict bad weather?
- Are they unethical?
- If in doubt, should we favor positive predictions for tourism? But then what about shops selling sweaters and coats?

- What is the solution?

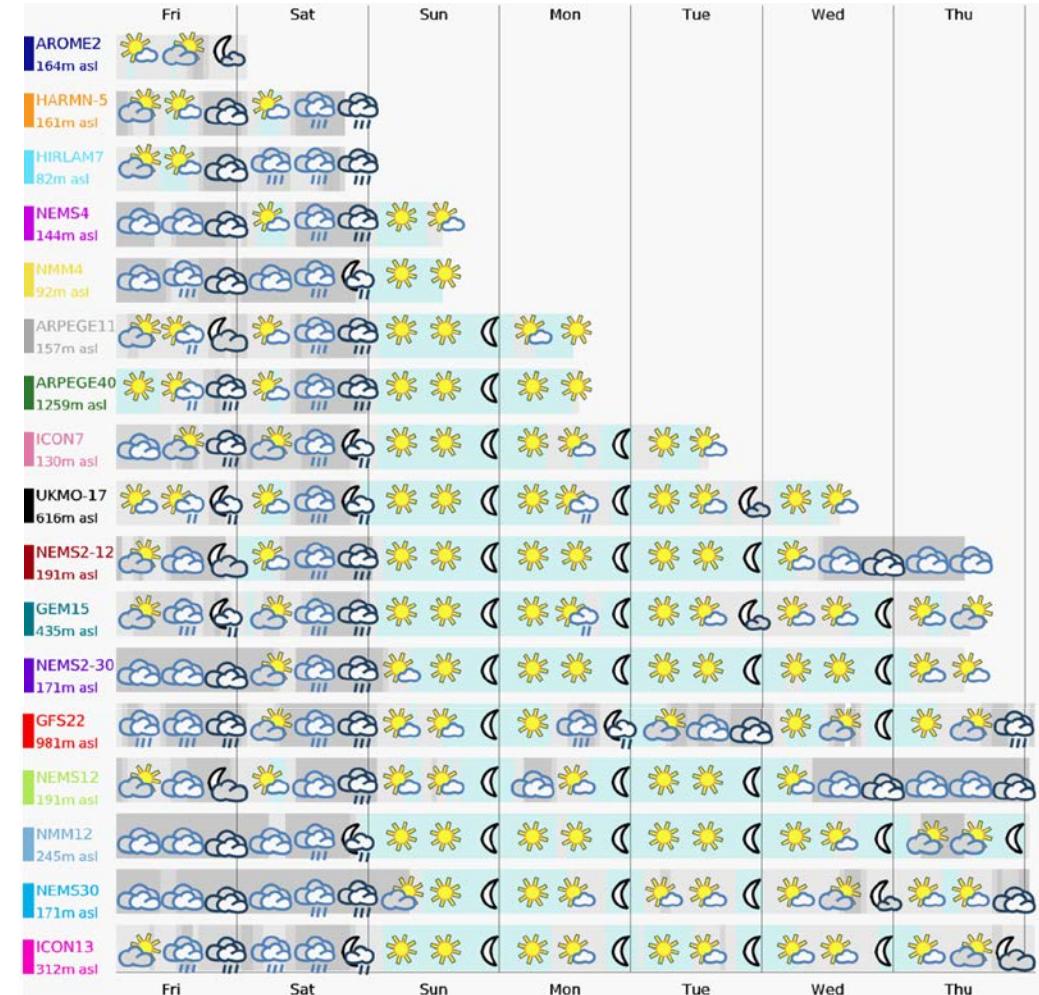
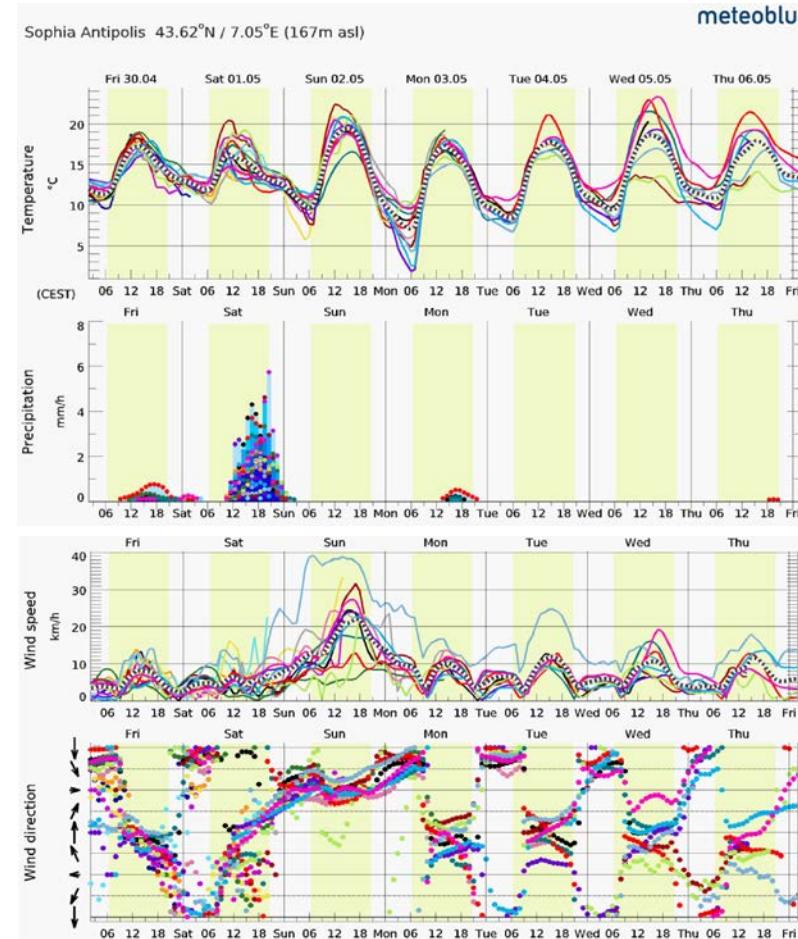
Back to weather forecast



A solution: Defining a confidence indice



Back to weather forecast



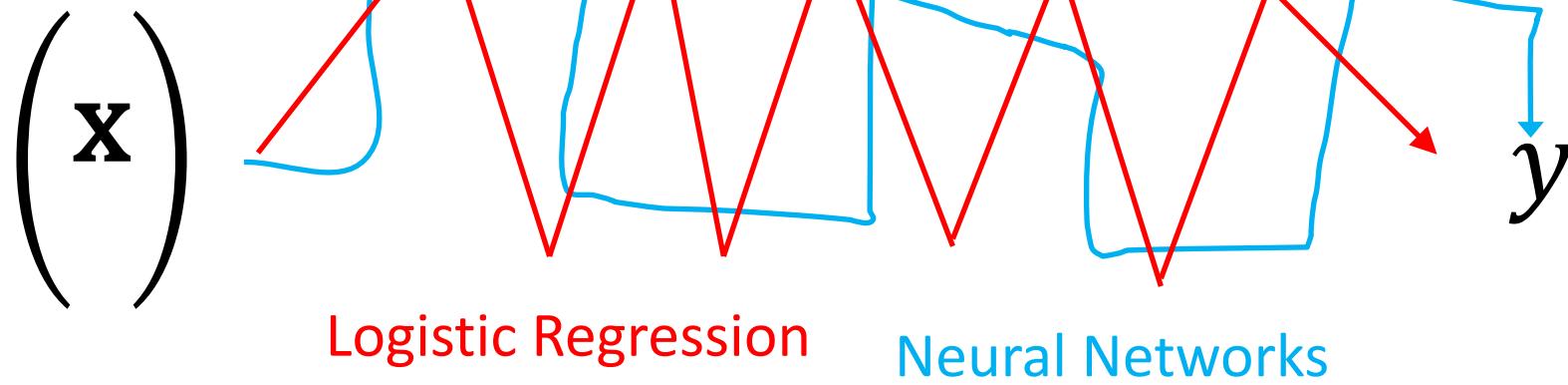
A solution: Combining different models

Partial conclusions

- A lot of what works today is not strictly speaking AI (compared to the original definition)
- It is the same as predicting the weather, predicting earthquakes or predicting the state of financial markets by very large mathematical models (hypothetico-deductive vs. inductive).
- Are we asking these models to be ethical?
- What does it mean that a mathematical function is ethical (or not)? Is \sqrt{x} more ethical than cosine? (*even if I have my own opinion on that point ;-)*)
- We can only measure the model on examples for which we know the result!

Machine Learning

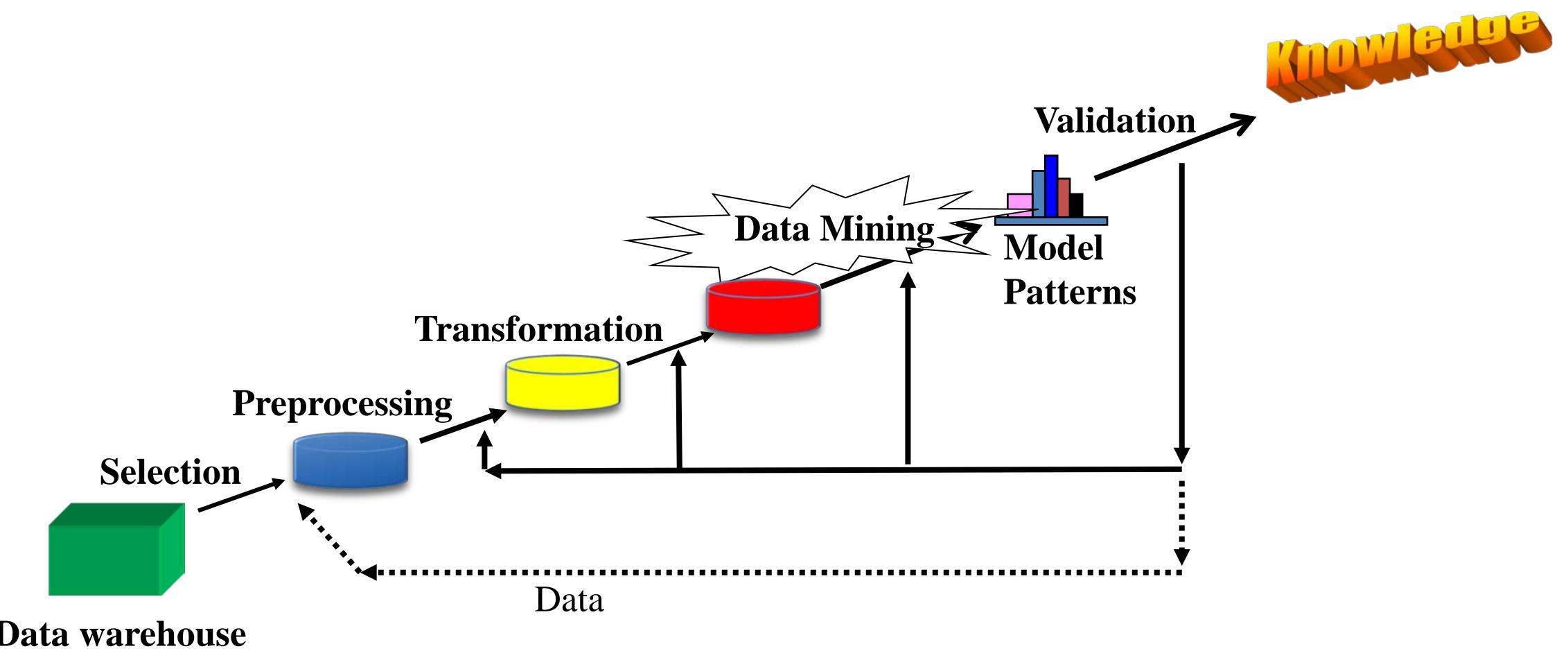
$$\begin{pmatrix} \mathbf{x} \end{pmatrix} \xrightarrow{f(\mathbf{X}, \alpha) ?} y$$



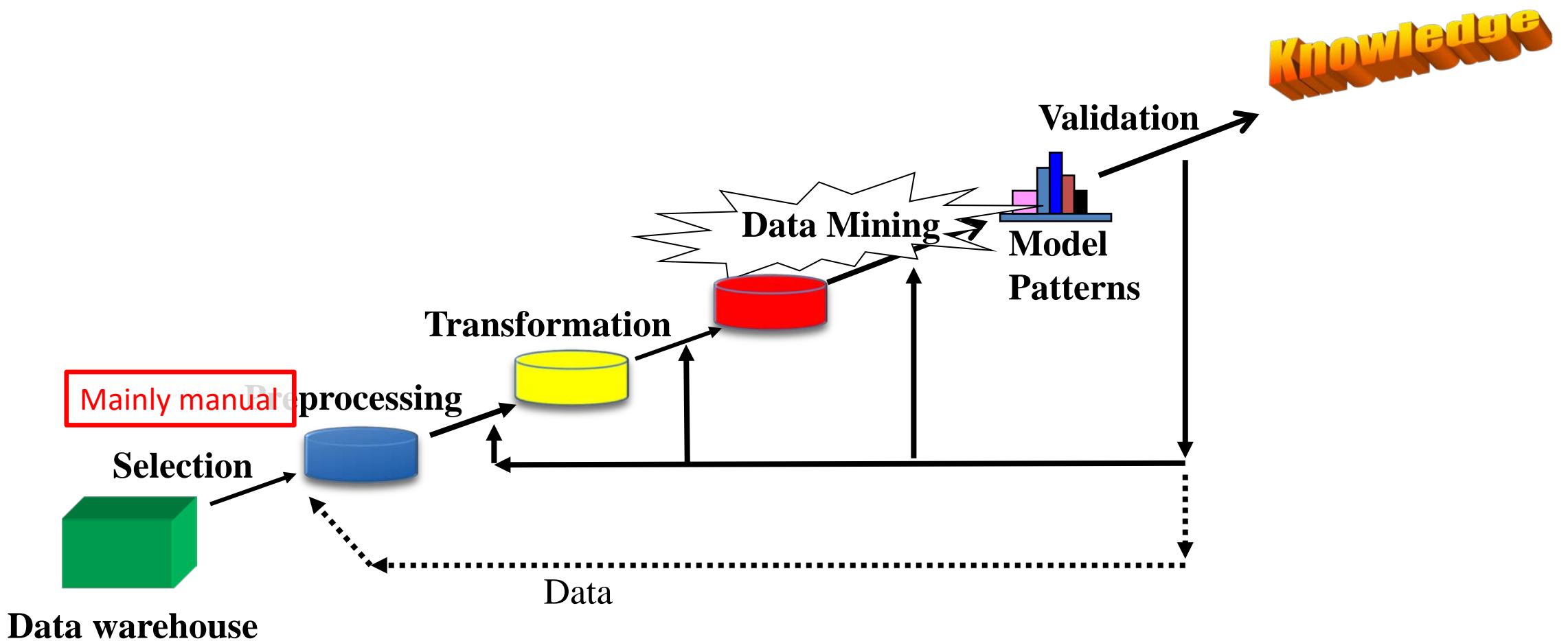


Machine learning VS Data Mining?

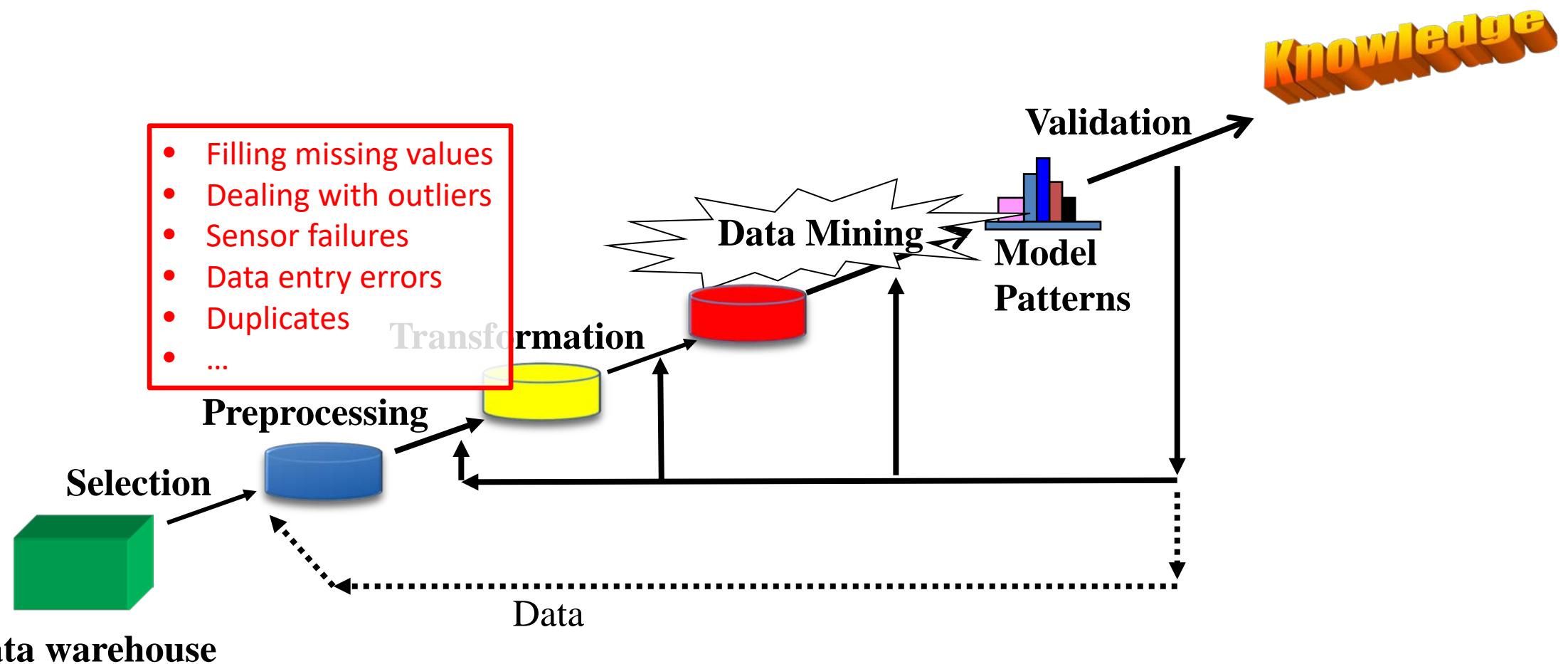
Data Mining Workflow



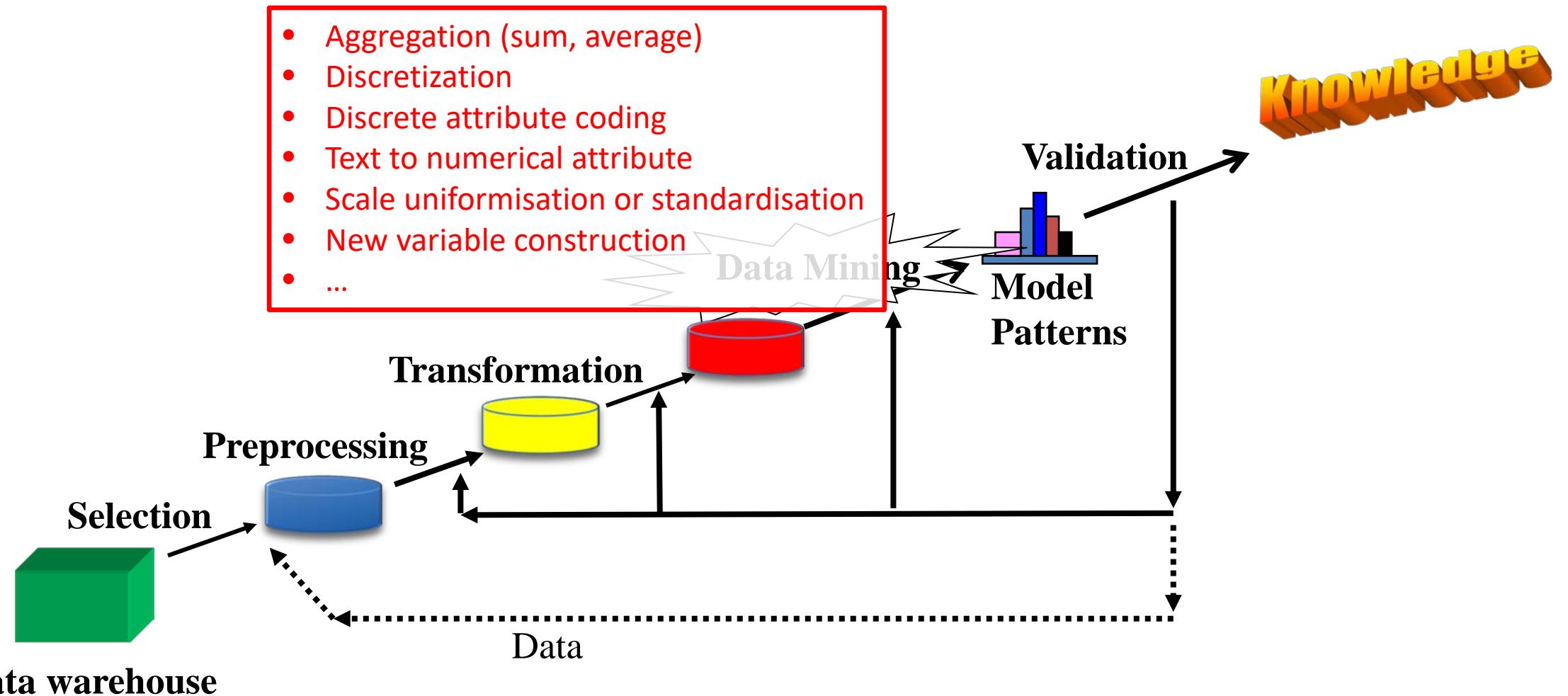
Data Mining Workflow



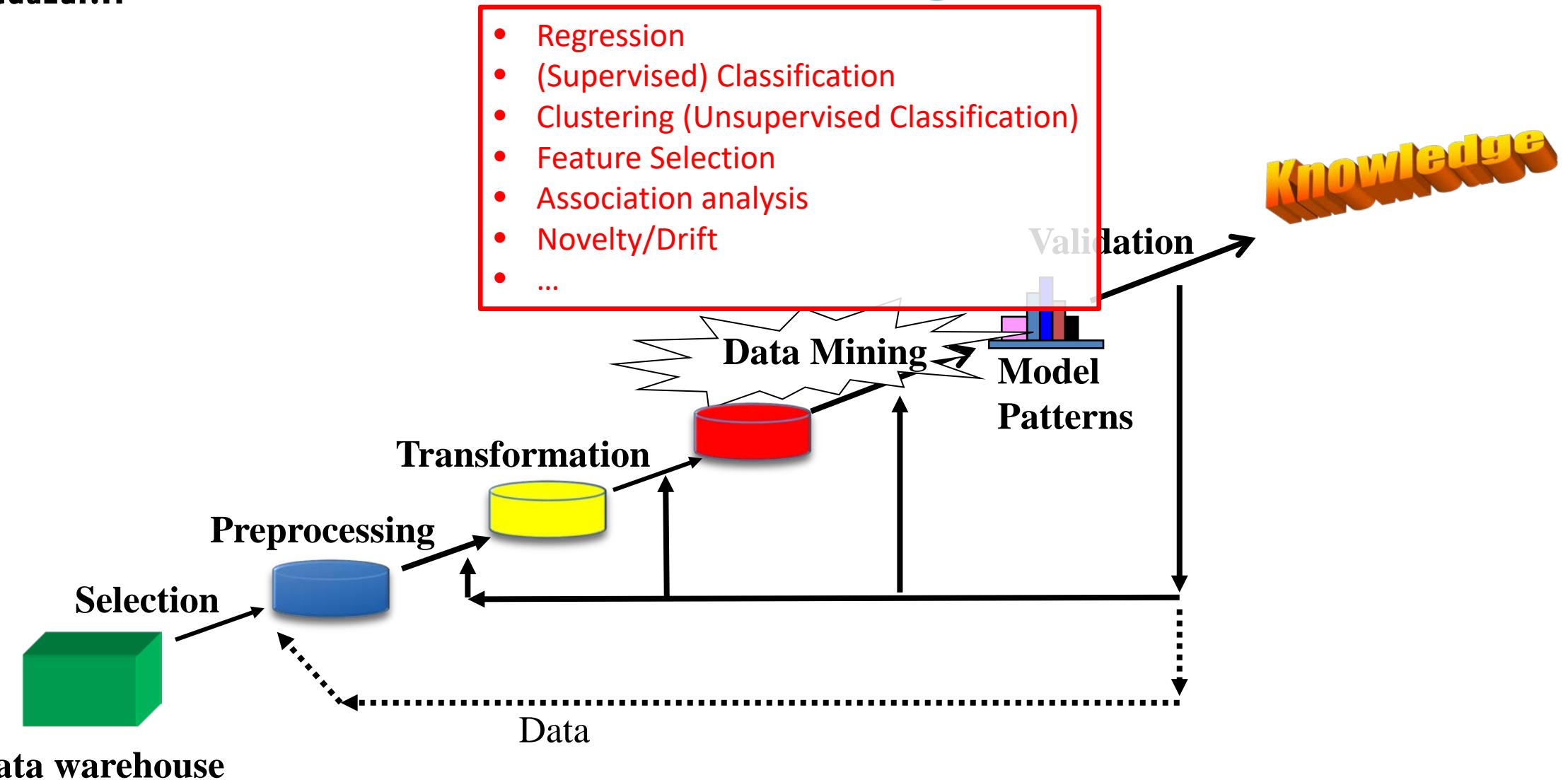
Data Mining Workflow



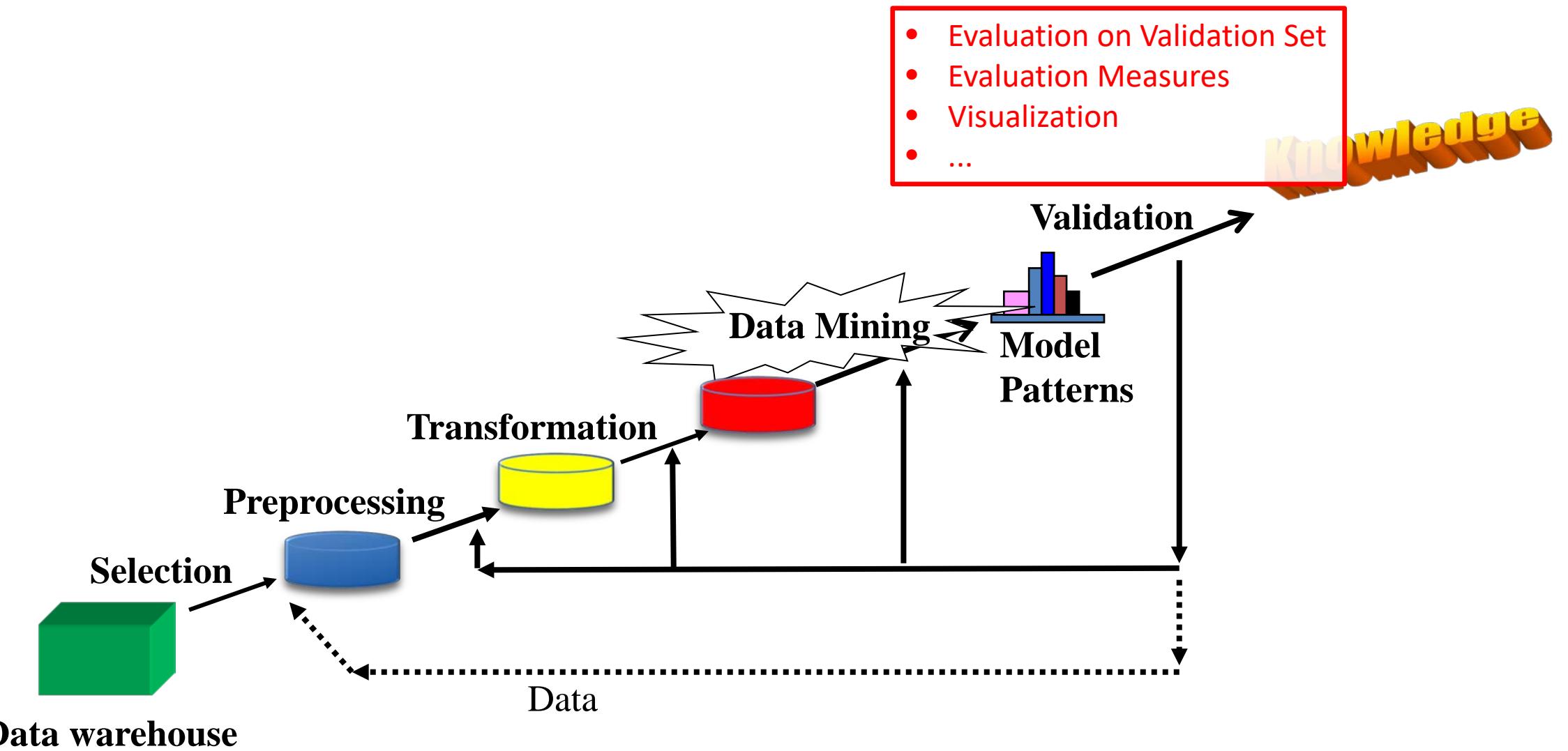
Data Mining Workflow

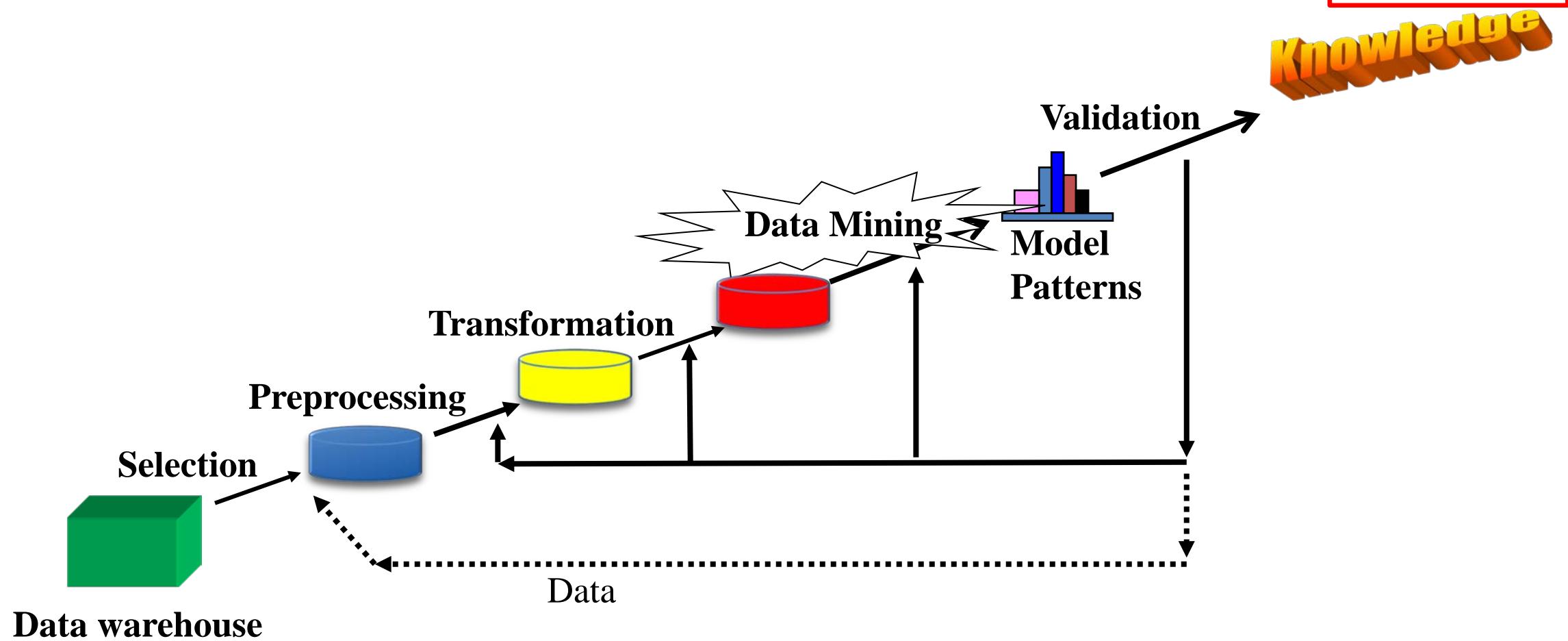


Data Mining Workflow



Data Mining Workflow





- Visualization
- Reporting
- Knowledge
- ...

Data Mining Workflow

Problems

- Regression
- (Supervised) Classification
- Density Estimation / Clustering
(Unsupervised Classification)
- Feature Selection
- Association analysis
- Anomaly/Novelty/Drift
- ...

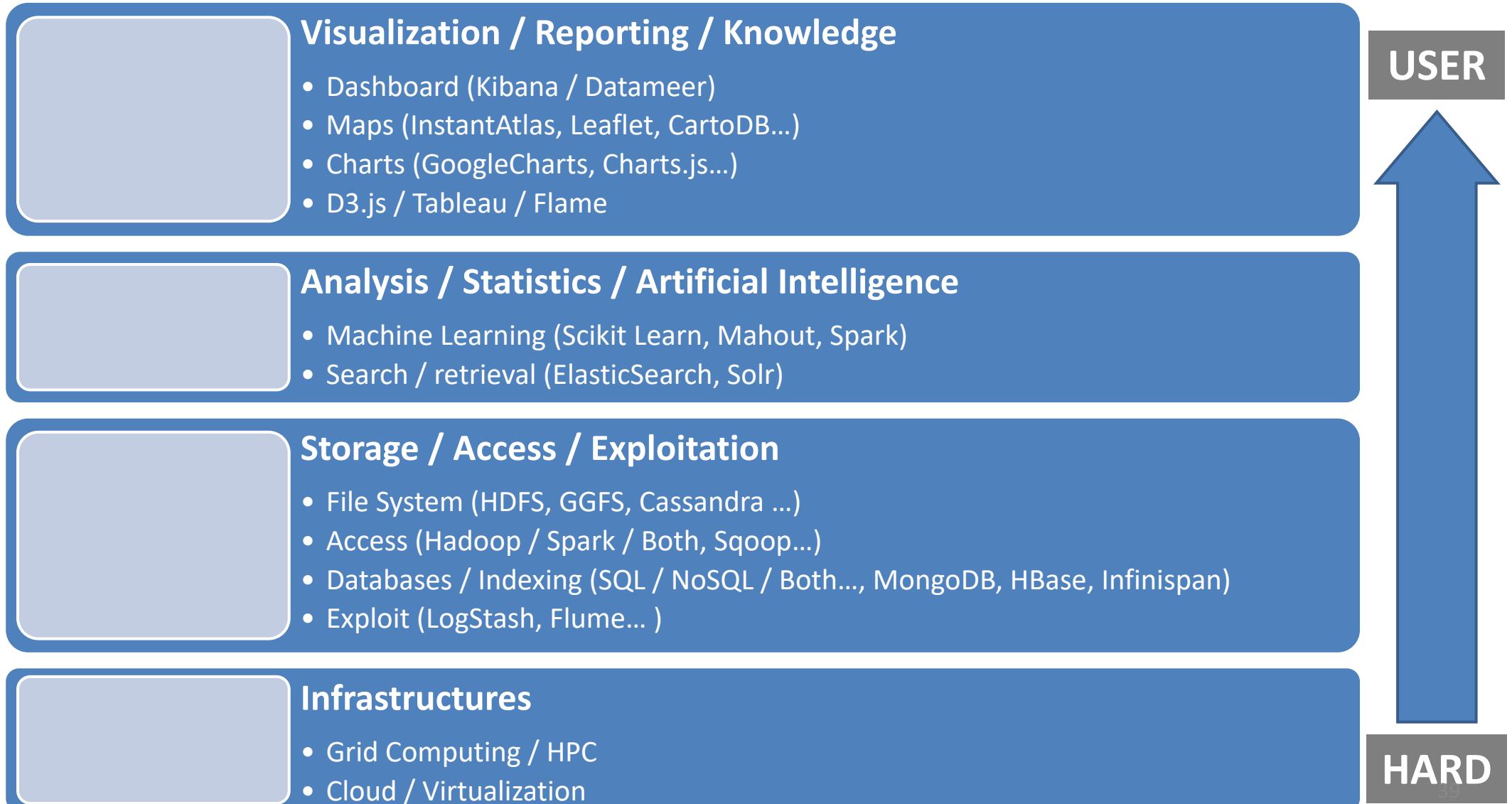
Possible Solutions

- Machine Learning
 - Support Vector Machine
 - Artificial Neural Network
 - Boosting
 - Decision Tree
 - Random Forest
 - ...
- Statistical Learning
 - Gaussian Models (GMM)
 - Naïve Bayes
 - Gaussian processes
 - ...
- Other techniques
 - Galois Lattice
 - ...



Machine learning VS Data science?

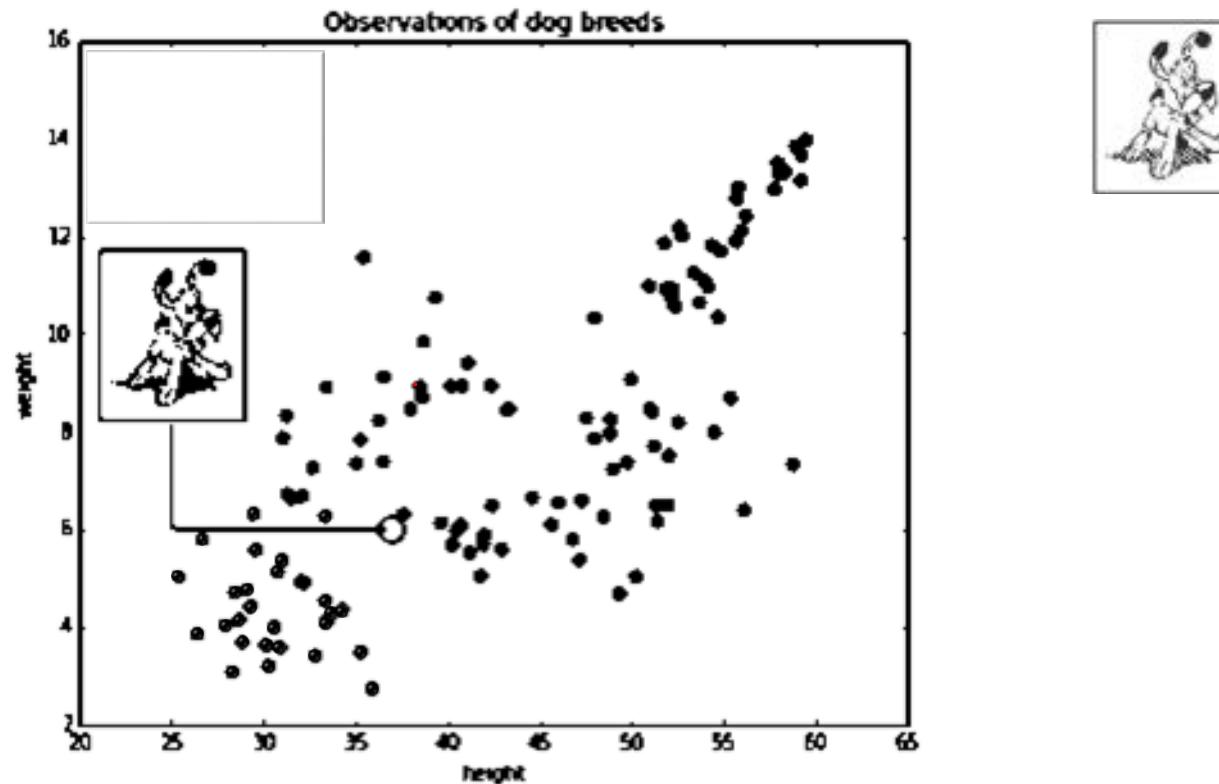
Data Science Stack





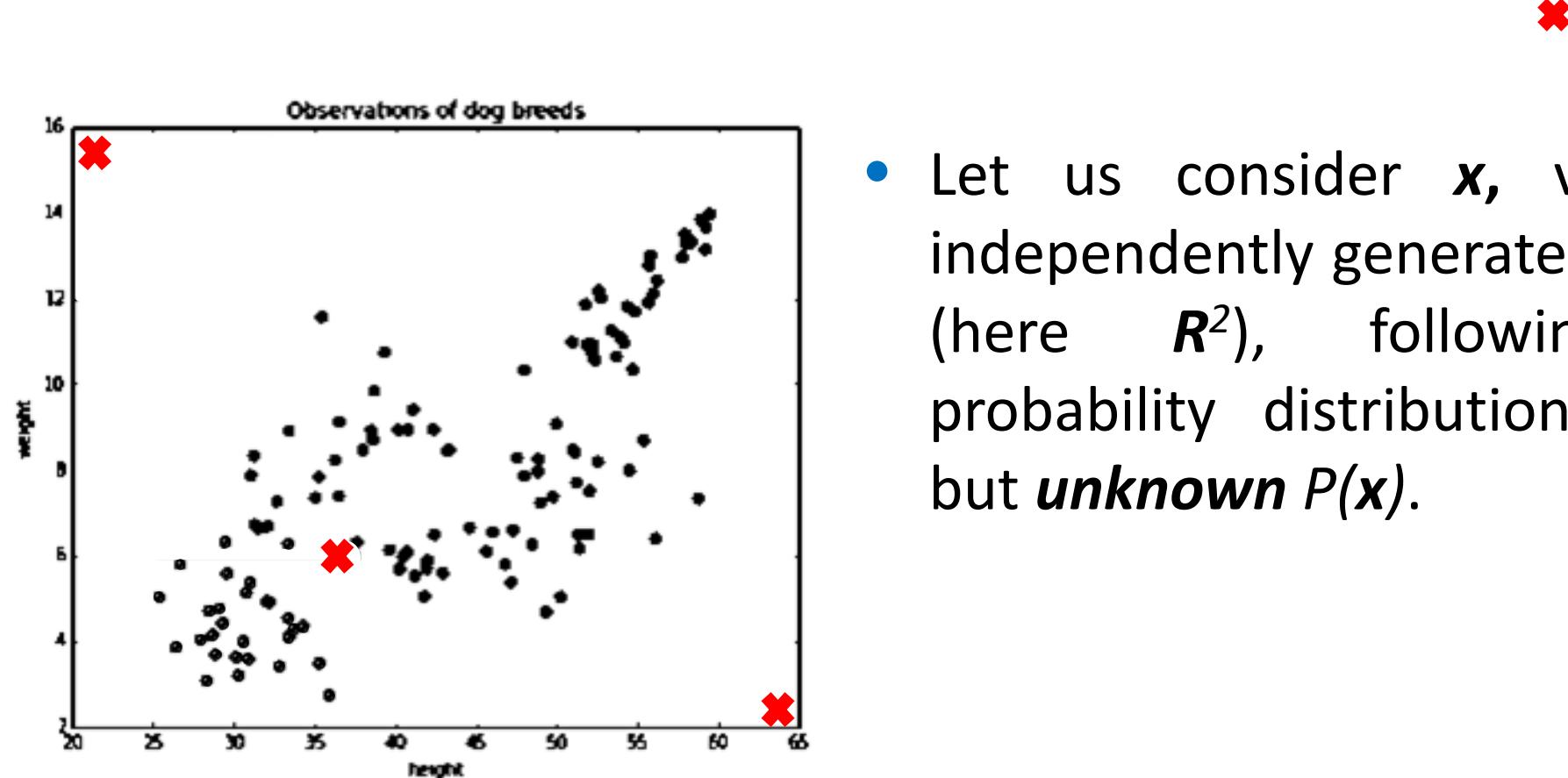
Machine Learning is Statistics?

What breed is that Dogmatix (Idéfix)?



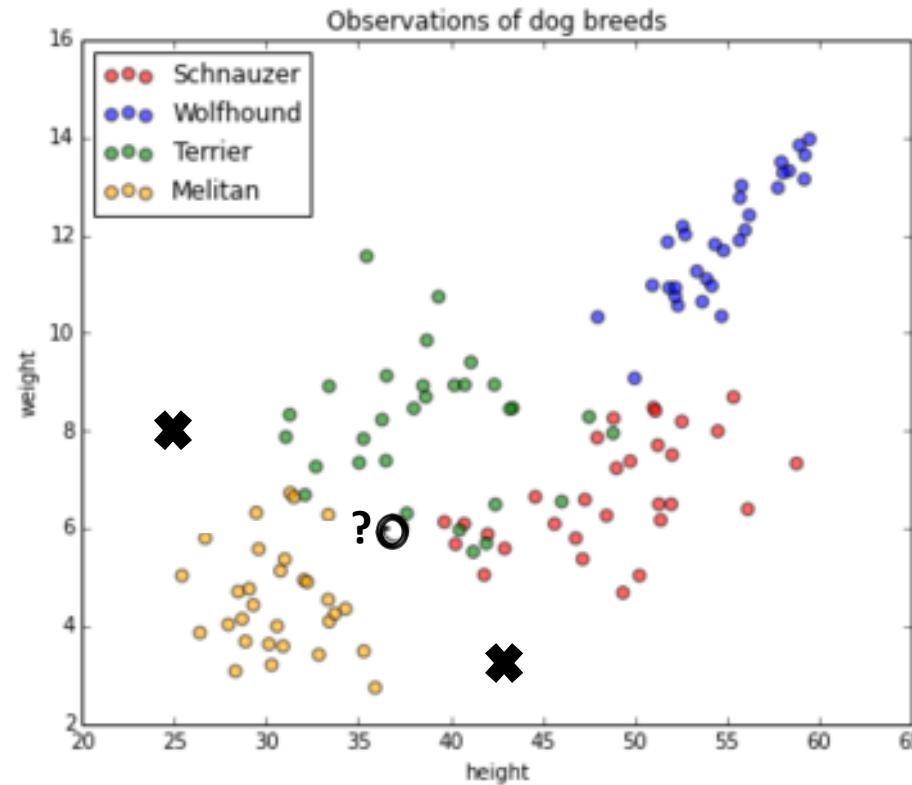
The illustrations of the slides in this section come from the blog "Bayesian Vitalstatistix: What Breed of Dog was Dogmatix?"

Does any real dog get this height and weight?



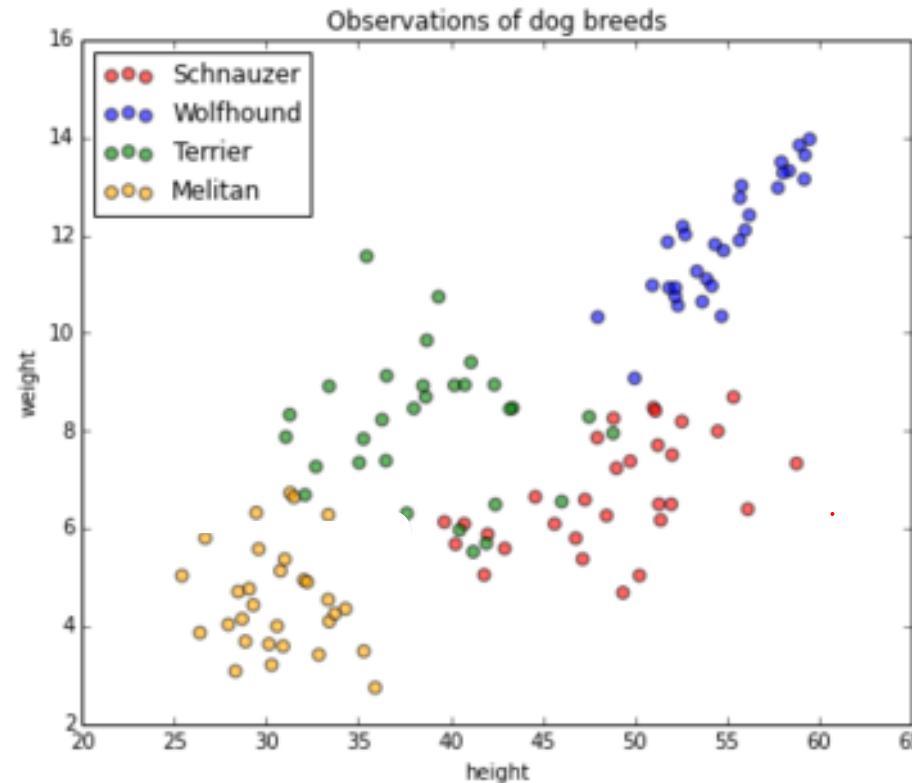
- Let us consider x , vectors independently generated in \mathbb{R}^d (here \mathbb{R}^2), following a probability distribution fixed but *unknown* $P(x)$.

What should be the breed of these dogs?



- An Oracle assigns a value y to each vector x following a probability distribution $P(y/x)$ also fixed but *unknown*.

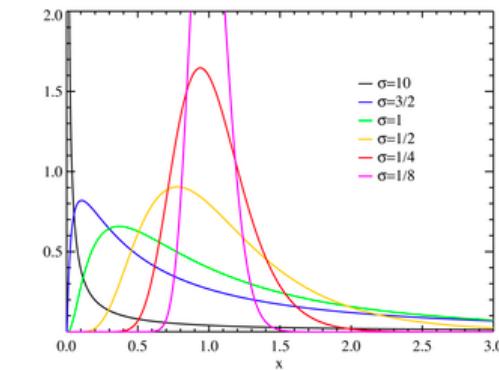
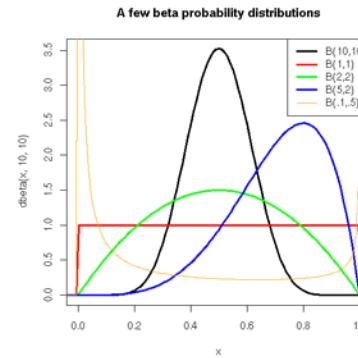
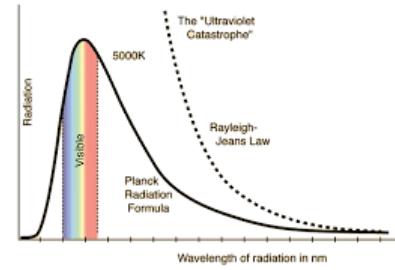
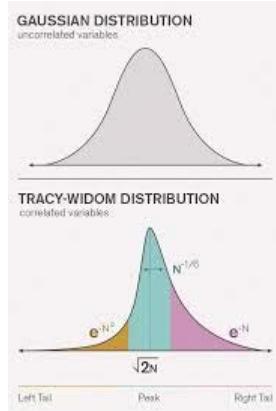
An oracle provides me with examples?



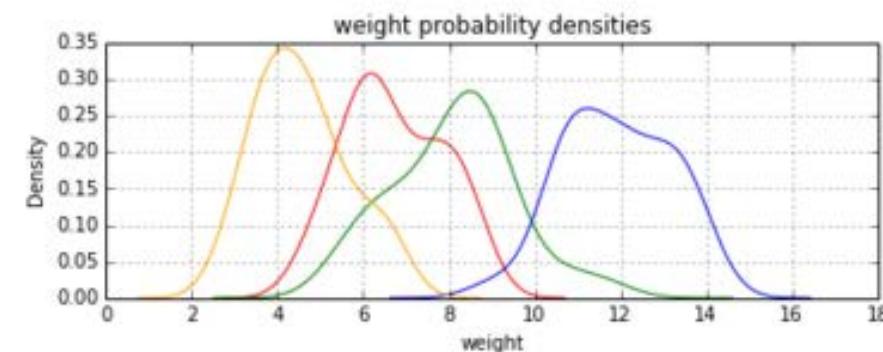
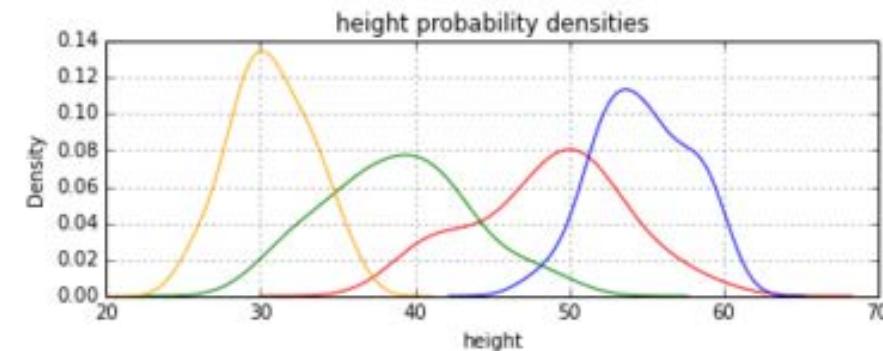
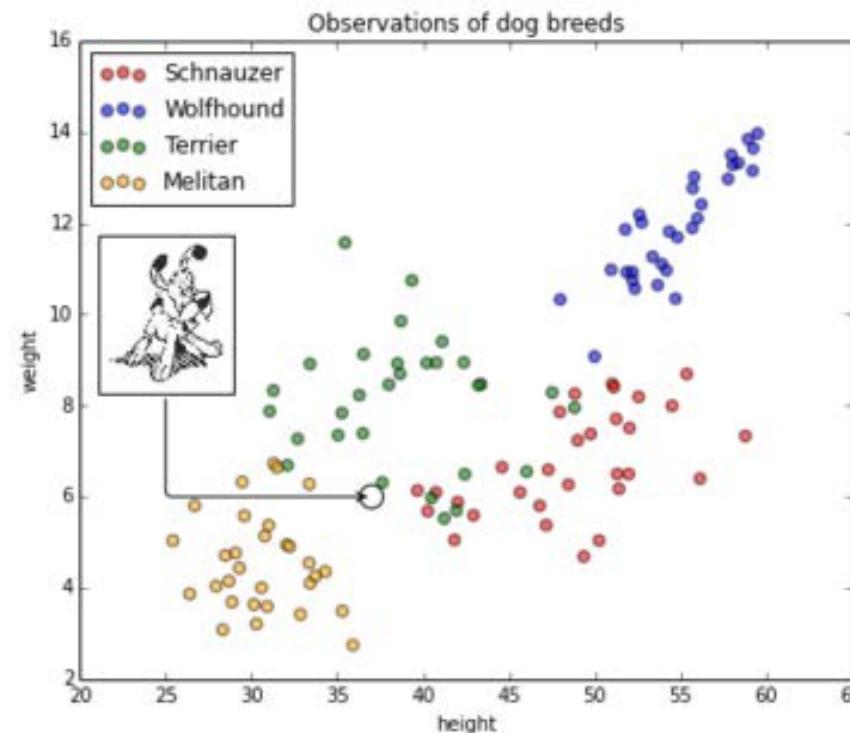
- Let S be a training set
$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\},$$
with m training samples i.i.d. which follow the joint probability
$$P(x, y) = P(x)P(y|x).$$



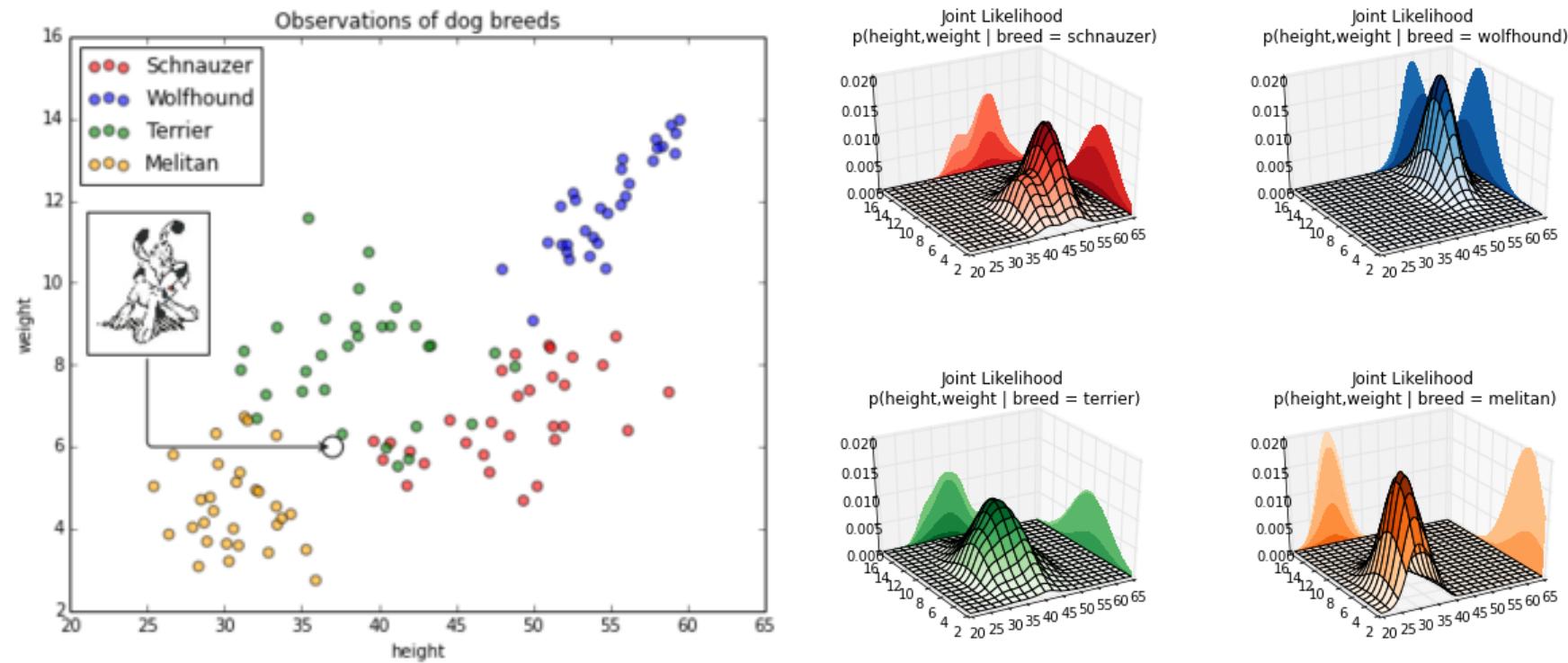
Statistical solution: Models, Hypotheses on data distribution...



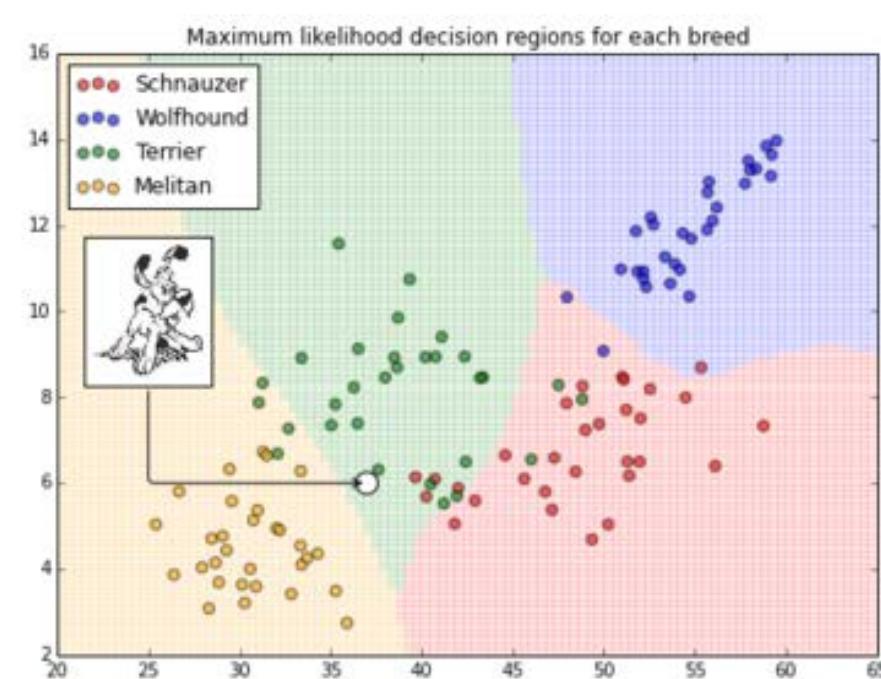
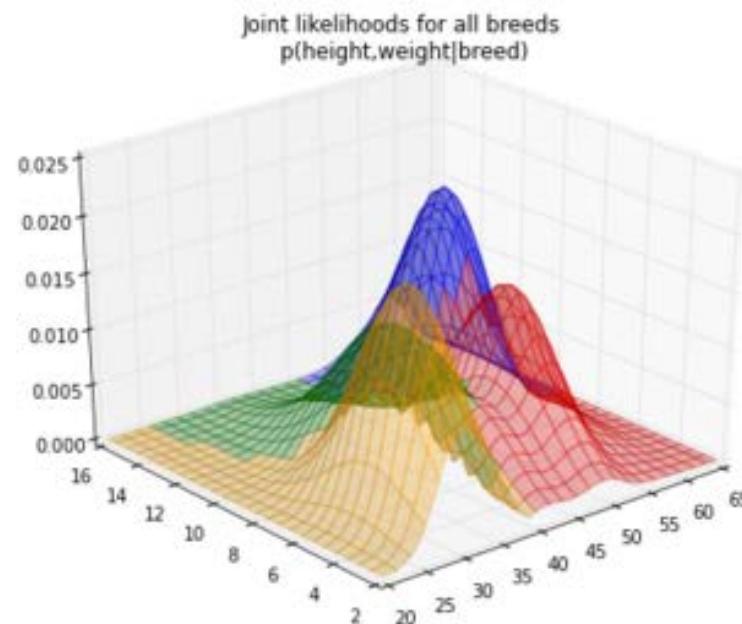
Statistical solution $P(\text{height}, \text{weight}|\text{breed})$...



Statistical solution $P(\text{height}, \text{weight} | \text{breed})...$

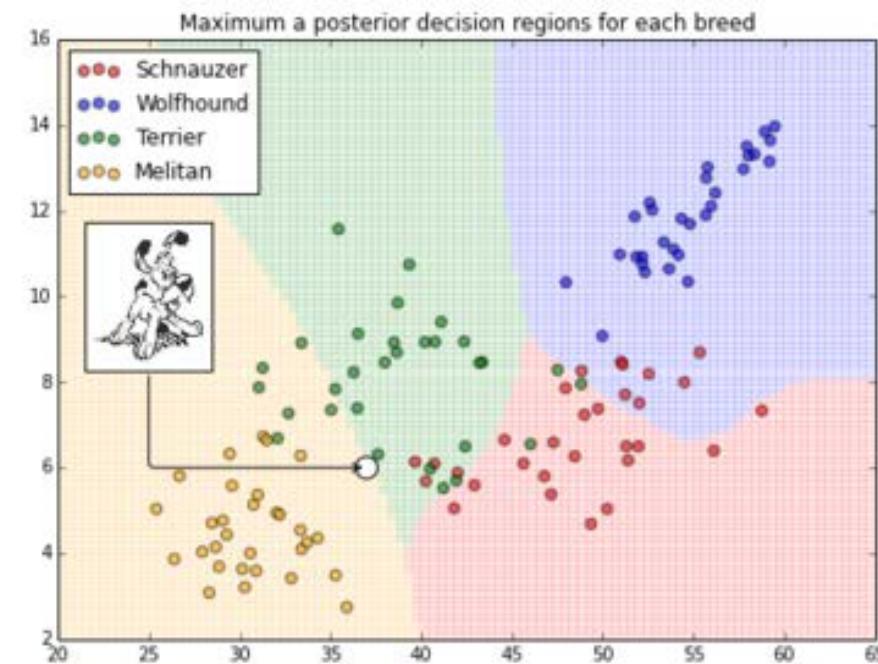
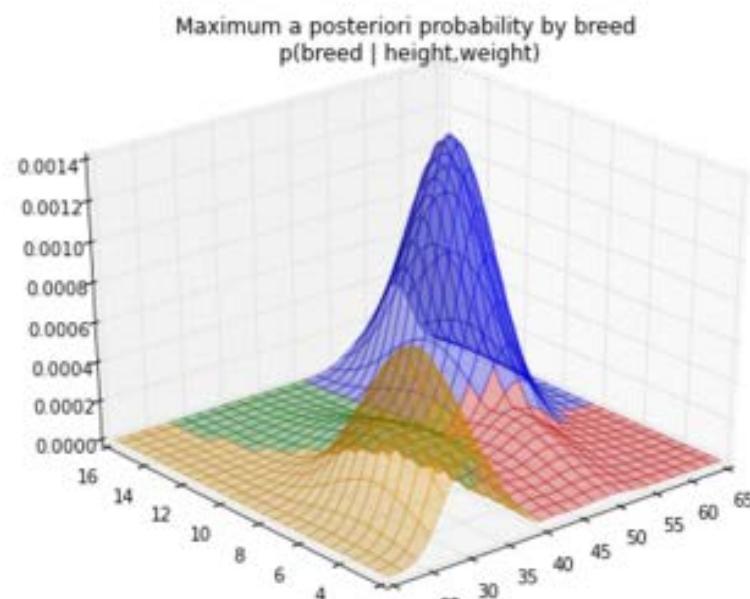


Statistical solution $P(\text{height}, \text{weight} | \text{breed})$...

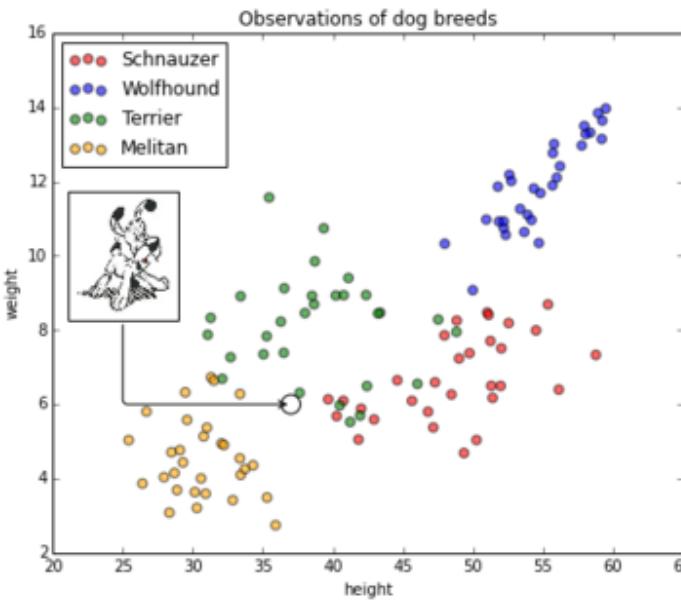




Statistical solution: Bayes, $P(\text{breed}|\text{height}, \text{weight})\dots$



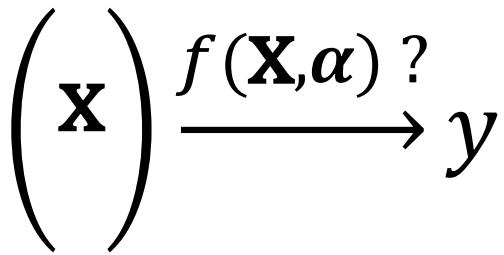
Machine Learning



- we have a learning machine (i.e. an algorithm) which can provide a family of functions $\{f(\mathbf{x}; \boldsymbol{\alpha})\}$, where $\boldsymbol{\alpha}$ is a set of parameters.

$$\begin{pmatrix} \mathbf{x} \end{pmatrix} \xrightarrow{f(\mathbf{X}, \boldsymbol{\alpha}) ?} y$$

The problem of (Machine) Learning



- The problem of learning consists in finding the function (among the $\{f(\mathbf{x}; \alpha)\}$) which provides the best approximation \hat{y} of the true label y given by the Oracle.
- “**best**” is defined in terms of minimizing a specific *error measure/cost/loss related to your problem/objectives*
 $L((\mathbf{x}, y), \alpha) \in [a; b]$.

The problem of (Machine) Learning

- Thus, the objective is to minimize the (*real*) **Risk**, i.e. the expectation of the error cost:

$$R(\alpha) = \int L((x, y), \alpha) dP(x, y)$$

where $P(x, y)$ is unknown.

- The training set $S = \{(x_i, y_i)\}_{i=1,\dots,m}$ is built through an i.i.d. sampling according to $P(x, y)$. Since we cannot compute $R(\alpha)$, we look for minimizing the **Empirical Risk** instead:

$$R_{emp}(\alpha) = \frac{1}{m} \sum_{k=1}^m L((x_i, y_i), \alpha)$$



Machine Learning fundamental Hypothesis

$S = \{(x_i, y_i)\}_{i=1, \dots, m}$ is built through an *i.i.d.* sampling according to $P(x, y)$.

Machine Learning  *Statistics*

Train through Cross-Validation

Machine Learning  *Statistics*

Training set & Test set have to be distributed according to the same law

Vapnik learning theory (1995)

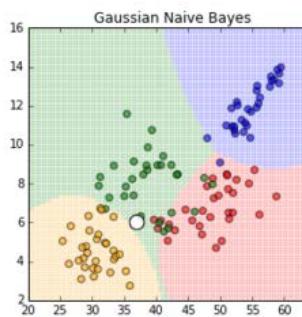
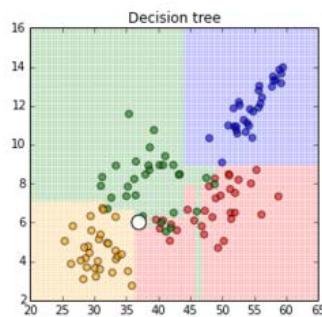
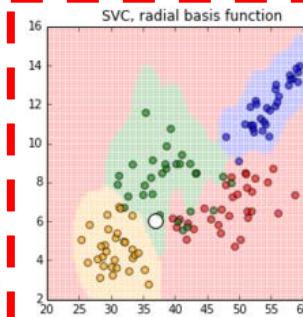
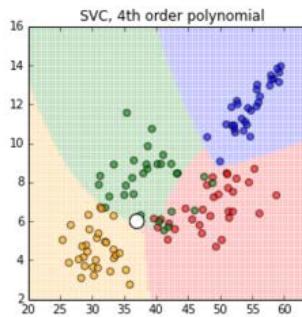
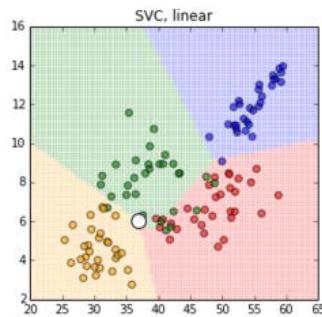
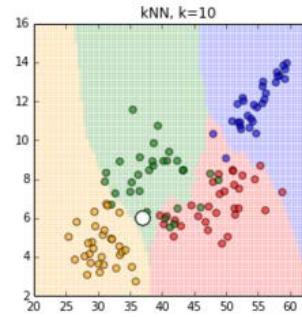
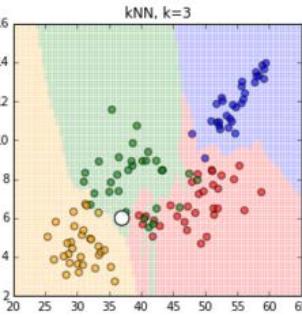
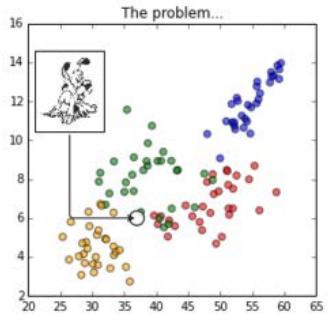
Vapnik had proven the following equation $\forall m$ with a probability at least equal to $1 - \eta$:

$$R(\alpha_m) \leq R_{emp}(\alpha_m) + (b - a) \sqrt{\frac{d_{VC}(\ln(2m/d_{VC}) + 1) - \ln(\eta/4)}{m}}$$

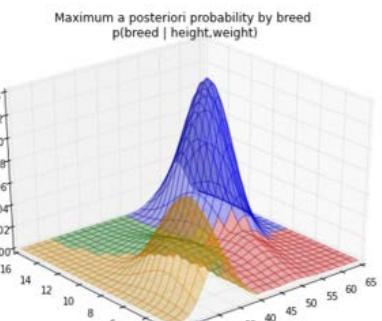
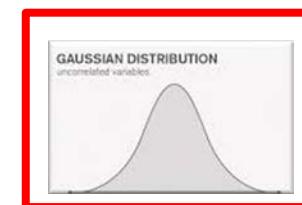
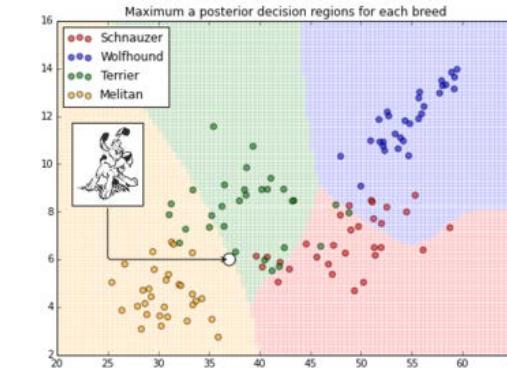
Training Error *Generalization Error*

Thus minimizing the **Risk** depends on minimizing the **Empirical Risk** and the **Generalization Error** of the model which depends on m (the number of training sample), and d_{VC} (the complexity of the model family chosen, also called *Vapnik-Chervonenkis Dimension*).

Machine Learning vs Statistics

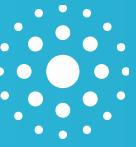


VS

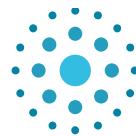


Machine Learning vs Statistics

- Machine Learning and Statistics can be considered as two sides of the same coin: Statistical Learning
- In Statistics, you make explicit assumptions on the model driving the data
- In Machine Learning, you do implicit Statistics; you make very few very general assumptions on the data but no assumptions on the model itself



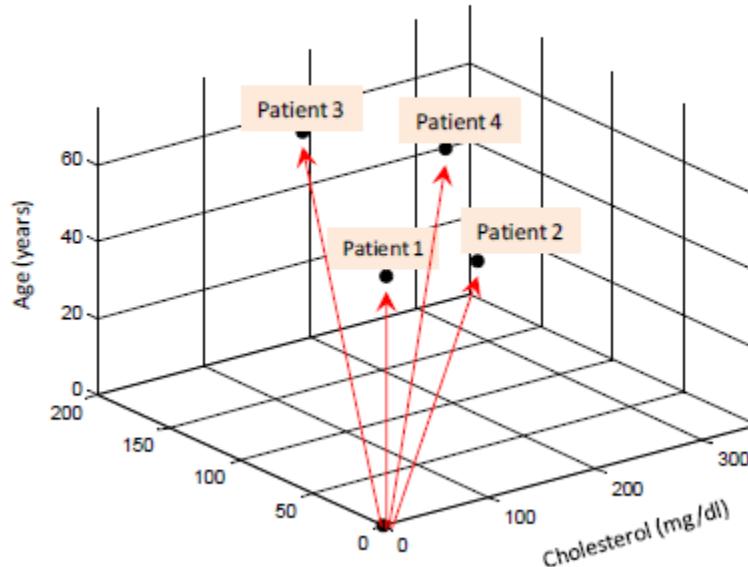
MATHEMATICAL BASICS



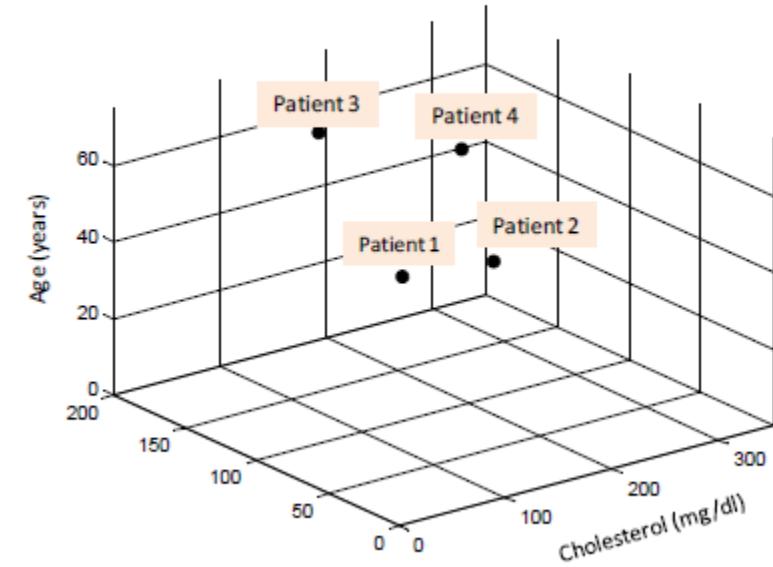
How to represent samples geometrically? Vectors & points in n-dimensional space (\mathbb{R}^n)

Patient id	Cholesterol (mg/dl)	Systolic BP (mmHg)	Age (years)	Tail of the vector	Arrow-head of the vector
1	150	110	35	(0,0,0)	(150, 110, 35)
2	250	120	30	(0,0,0)	(250, 120, 30)
3	140	160	65	(0,0,0)	(140, 160, 65)
4	300	180	45	(0,0,0)	(300, 180, 45)

Vector representation



Point representation





Basic operation on vectors in \mathbb{R}^n

1. Multiplication by a scalar

Consider a vector $\mathbf{a} = (a_1, a_2, \dots, a_n)$ and a scalar c

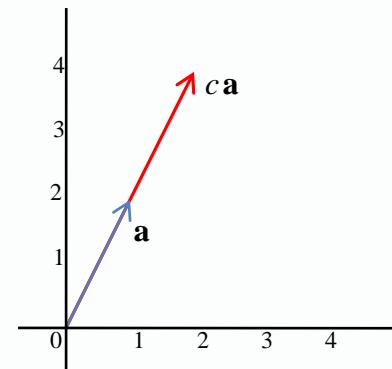
Define: $c\mathbf{a} = (ca_1, ca_2, \dots, ca_n)$

When you multiply a vector by a scalar, you “stretch” it in the same or opposite direction depending on whether the scalar is positive or negative.

$$\mathbf{a} = (1, 2)$$

$$c = 2$$

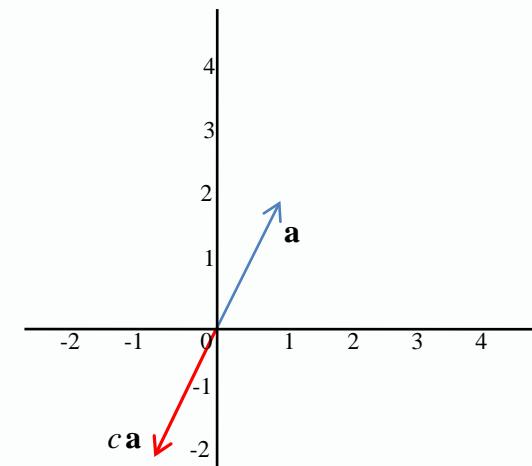
$$c\mathbf{a} = (2, 4)$$



$$\mathbf{a} = (1, 2)$$

$$c = -1$$

$$c\mathbf{a} = (-1, -2)$$



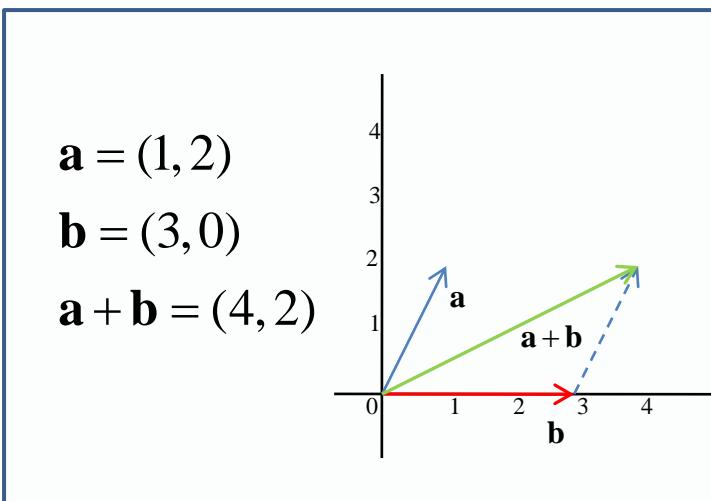


Basic operation on vectors in \mathbb{R}^n

2. Addition

Consider vectors $\mathbf{a} = (a_1, a_2, \dots, a_n)$ and $\mathbf{b} = (b_1, b_2, \dots, b_n)$

Define: $\mathbf{a} + \mathbf{b} = (a_1 + b_1, a_2 + b_2, \dots, a_n + b_n)$



Recall addition of forces in classical mechanics.



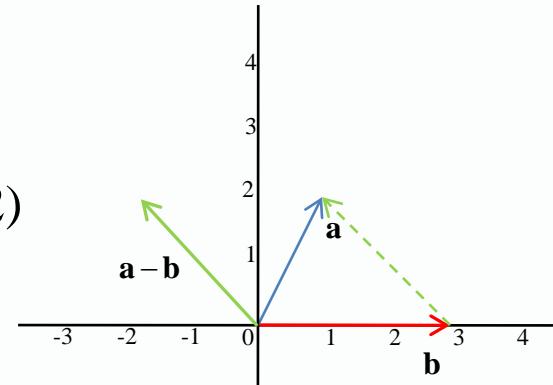
Basic operation on vectors in \mathbb{R}^n

3. Subtraction

Consider vectors $\mathbf{a} = (a_1, a_2, \dots, a_n)$ and $\mathbf{b} = (b_1, b_2, \dots, b_n)$

Define: $\mathbf{a} - \mathbf{b} = (a_1 - b_1, a_2 - b_2, \dots, a_n - b_n)$

$$\begin{aligned}\mathbf{a} &= (1, 2) \\ \mathbf{b} &= (3, 0) \\ \mathbf{a} - \mathbf{b} &= (-2, 2)\end{aligned}$$



What vector do we need to add to \vec{b} to get \vec{a} ? I.e., similar to subtraction of real numbers.



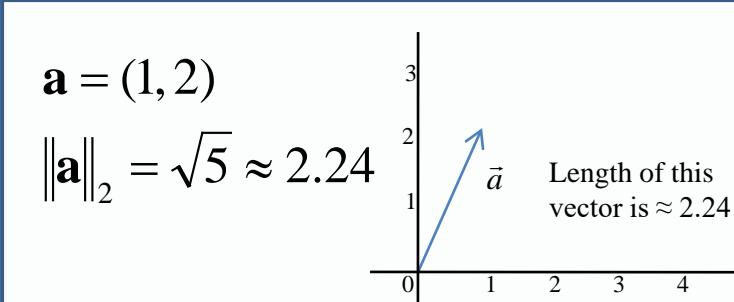
Basic operation on vectors in \mathbb{R}^n

4. Euclidian length or L2-norm

Consider a vector $\mathbf{a} = (a_1, a_2, \dots, a_n)$

Define the L2-norm: $\|\mathbf{a}\|_2 = \sqrt{a_1^2 + a_2^2 + \dots + a_n^2}$

We often denote the L2-norm without subscript, i.e. $\|\mathbf{a}\|$



L2-norm is a typical way to measure length of a vector; other methods to measure length also exist.



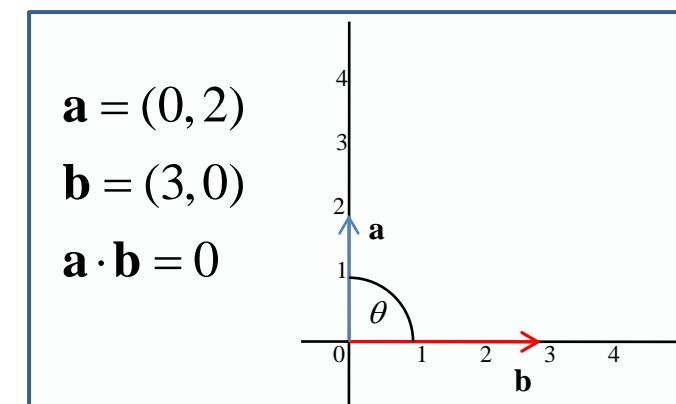
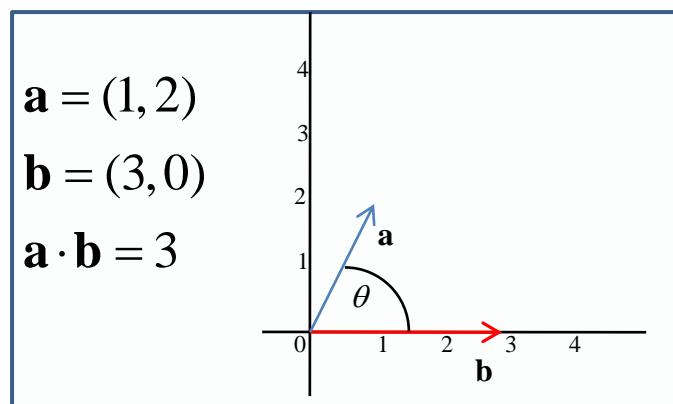
Basic operation on vectors in \mathbb{R}^n

5. Dot product

Consider vectors $\mathbf{a} = (a_1, a_2, \dots, a_n)$ and $\mathbf{b} = (b_1, b_2, \dots, b_n)$

Define dot product: $\mathbf{a} \cdot \mathbf{b} = a_1b_1 + a_2b_2 + \dots + a_nb_n = \sum_{i=1}^n a_i b_i$

The law of cosines says that $\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\|_2 \|\mathbf{b}\|_2 \cos \theta$ where θ is the angle between \mathbf{a} and \mathbf{b} . Therefore, when the vectors are perpendicular $\mathbf{a} \cdot \mathbf{b} = 0$.





Basic operation on vectors in \mathbb{R}^n

5. Dot product (continued)

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2 + \dots + a_n b_n = \sum_{i=1}^n a_i b_i$$

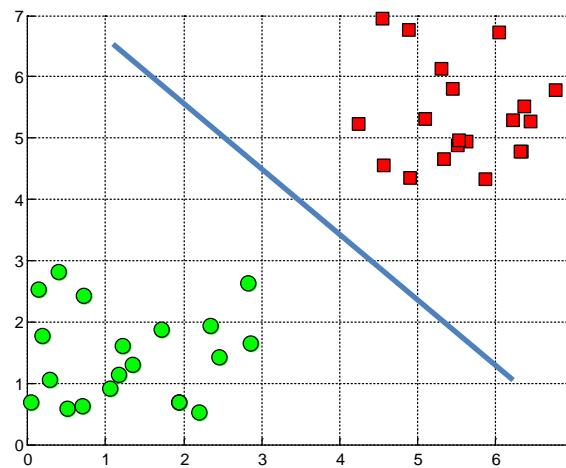
- Property: $\mathbf{a} \cdot \mathbf{a} = a_1 a_1 + a_2 a_2 + \dots + a_n a_n = \|\mathbf{a}\|_2^2$
- In the classical regression equation $y = \mathbf{w} \cdot \mathbf{x} + b$
the response variable y is just a dot product of the vector representing patient characteristics (\mathbf{X}) and the regression weights vector (\mathbf{w}) which is common across all patients plus an offset b .



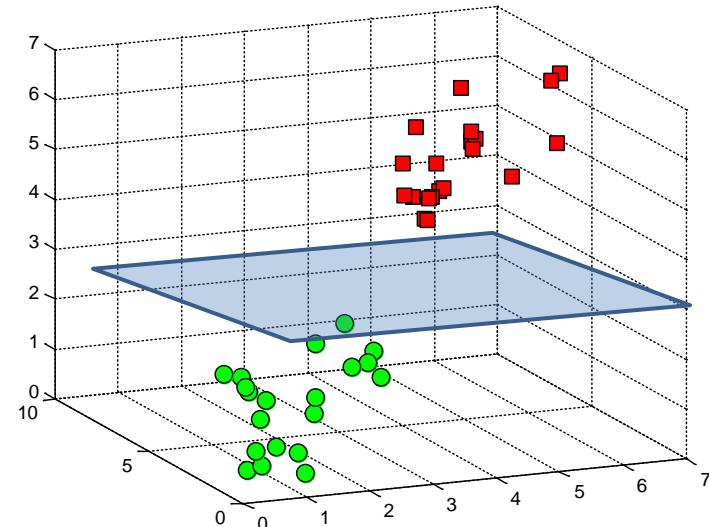
Hyperplanes as decision surfaces

- A hyperplane is a linear decision surface that splits the space into two parts;
- It is obvious that a hyperplane is a binary classifier.

A hyperplane in \mathbb{R}^2 is a line

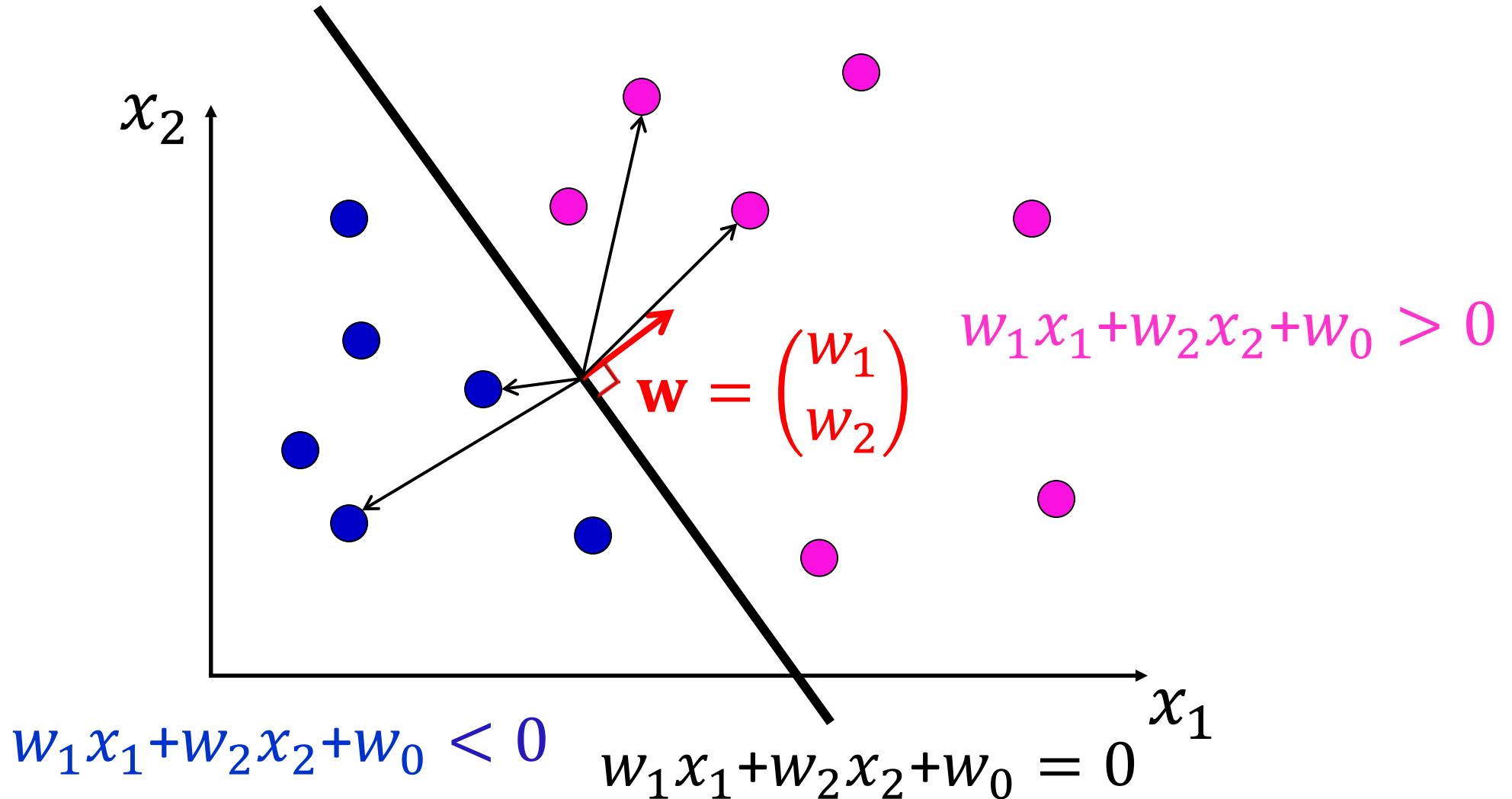


A hyperplane in \mathbb{R}^3 is a plane

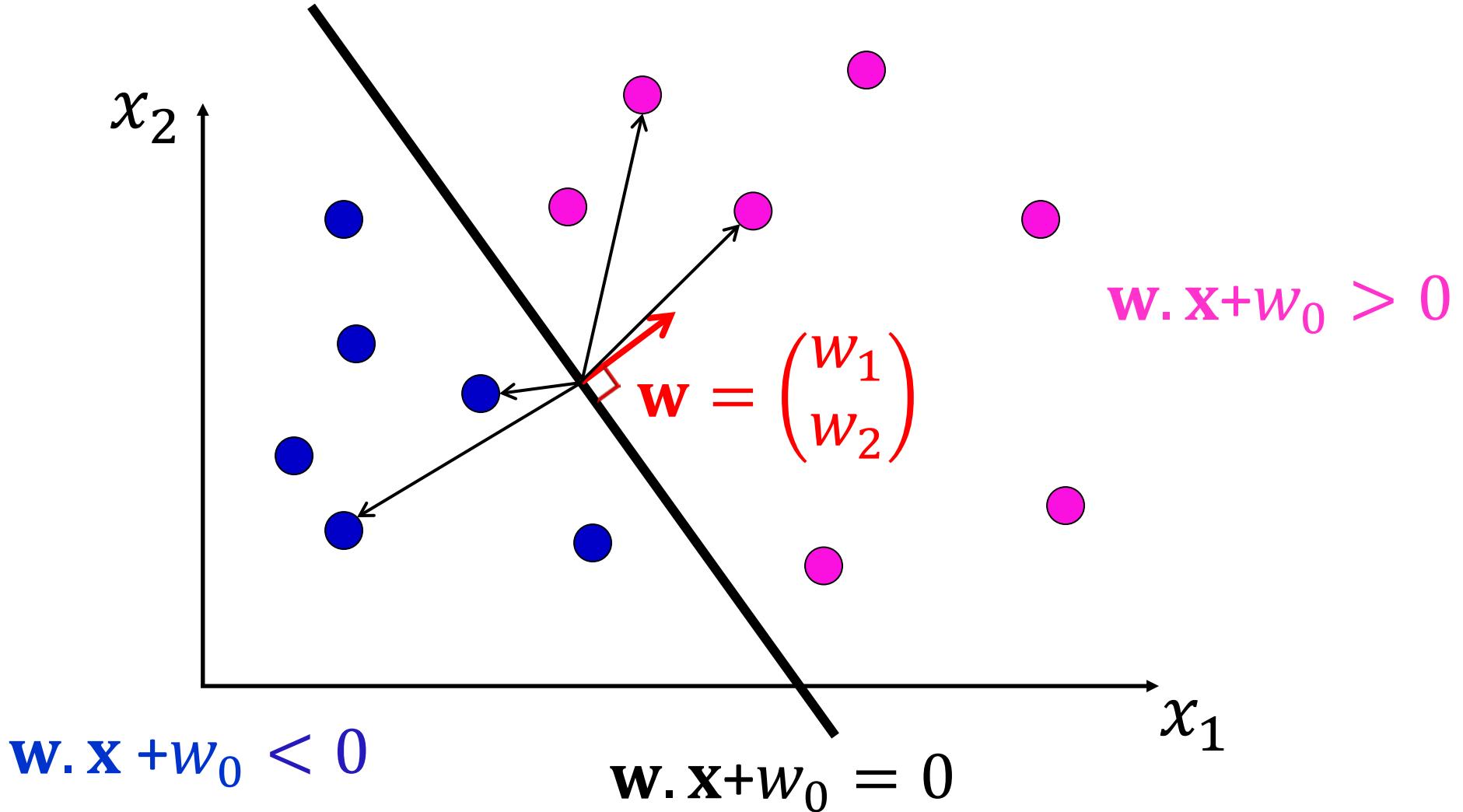


A hyperplane in \mathbb{R}^n is an $n-1$ dimensional subspace

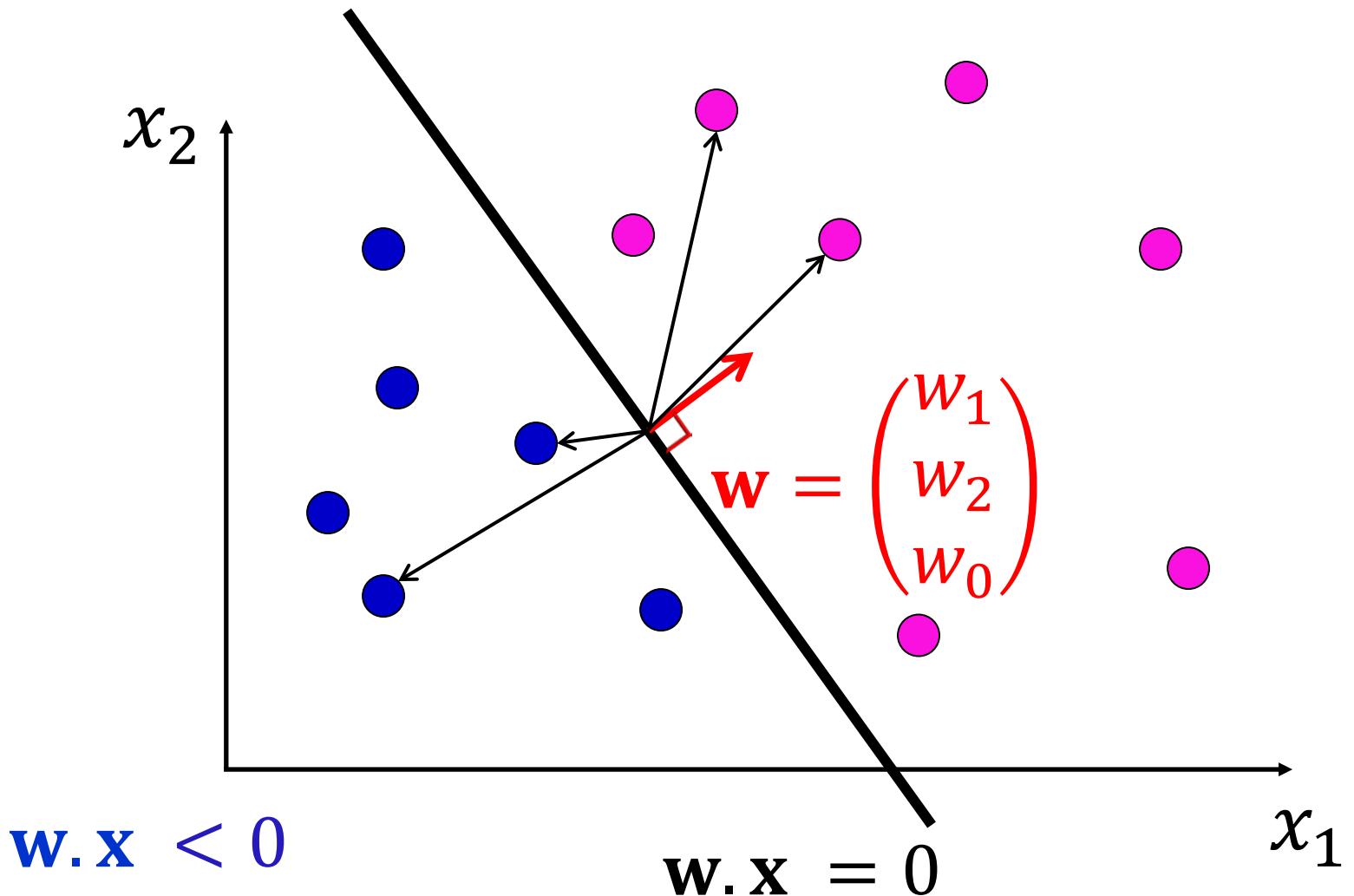
Geometry and Algebra



Geometry and Algebra

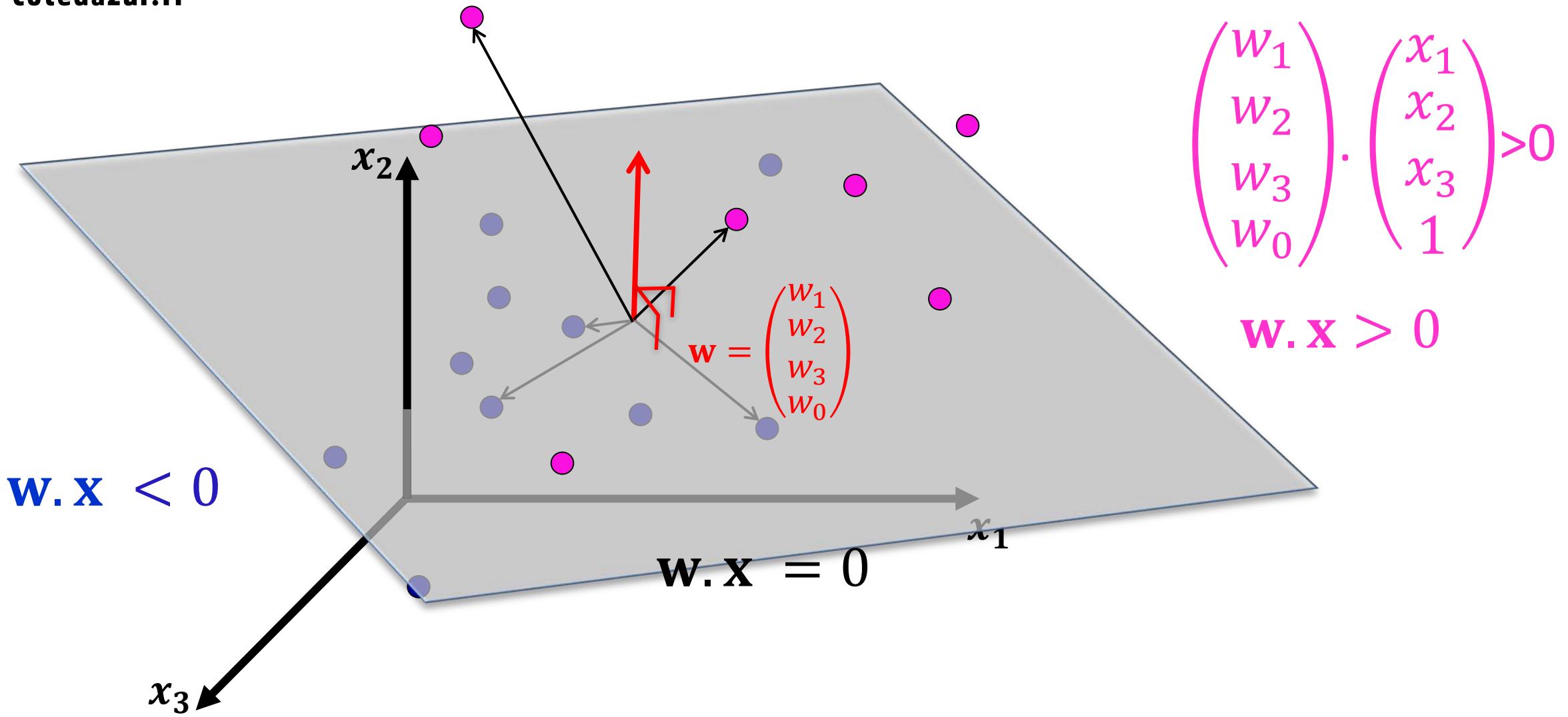


Simplification



$$\begin{pmatrix} w_1 \\ w_2 \\ w_0 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix} > 0$$
$$w \cdot x > 0$$

Geometry and Algebra



Derivative

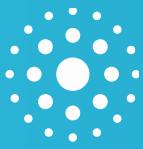
- Derivative of a scalar function (of one variable)
 - $(ax)' = a$
 - $(ax + b)' = a$
 - $(g(f(x)))' = g'(f(x))f'(x)$ **Chain rule**

Gradient operator

- Gradient of a multivariate function (\mathbf{x} is a vector)
 - $\nabla_{\mathbf{x}}(\mathbf{a}\mathbf{x}) = \mathbf{a}$
 - $\nabla_{\mathbf{x}}(\mathbf{a}\mathbf{x} + b) = \mathbf{a}$
 - $\nabla_{\mathbf{x}}(g(f(\mathbf{x}))) = \nabla_{\mathbf{x}}g(f(\mathbf{x})) \nabla_{\mathbf{x}}f(\mathbf{x}) \quad (\text{Chain rule})$
 - $\nabla_{\mathbf{x}_i}(\mathbf{v} \cdot \mathbf{x}) = \nabla_{\mathbf{x}_i}(v_1 \cdot x_1 + v_2 \cdot x_2 + \cdots + v_n \cdot x_n) = v_i$

Overview

- Machine Learning vs Statistics
- Math Basics
- **Simple Model**
- From Simple to Complex



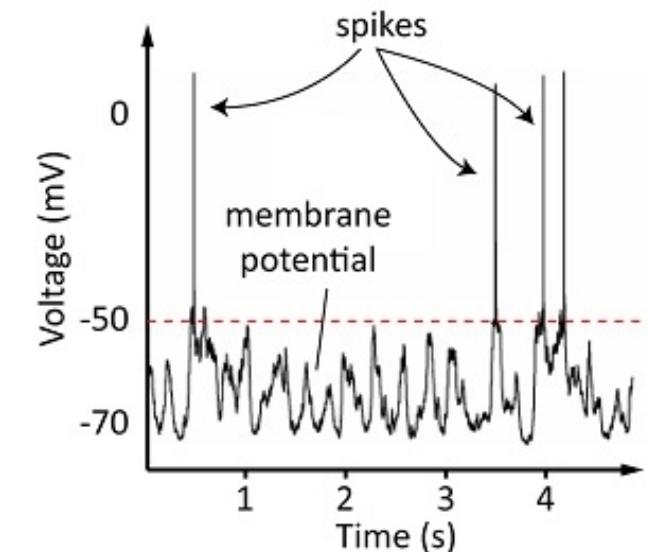
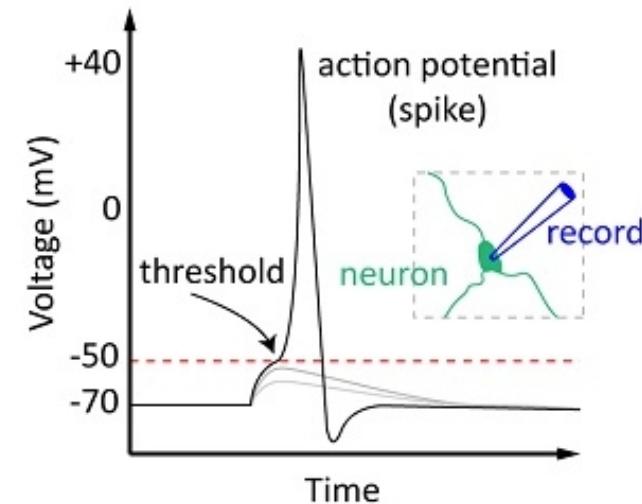
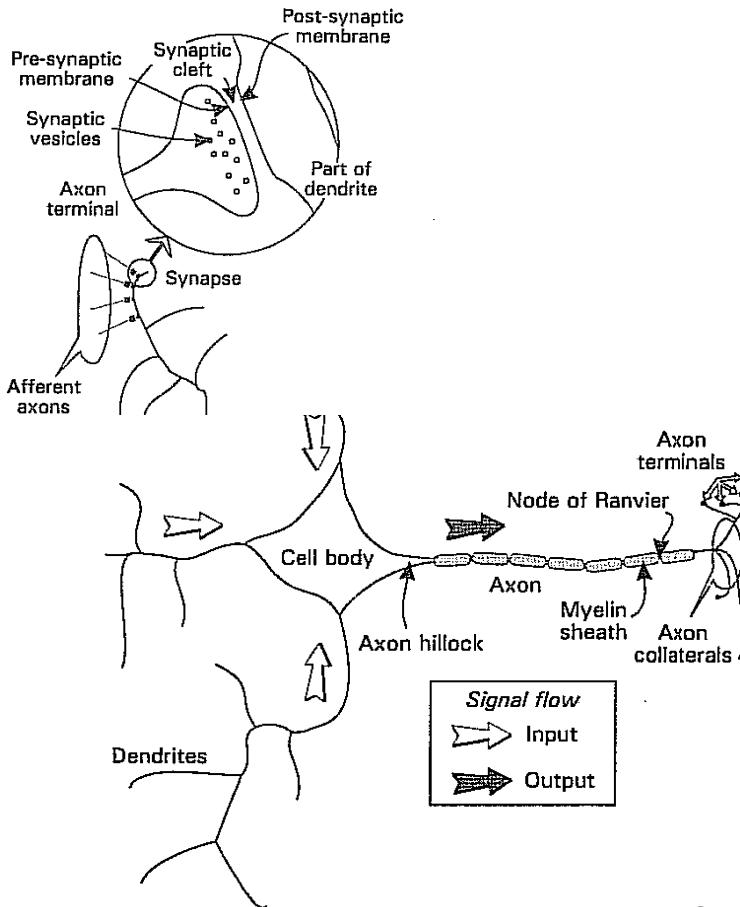
SIMPLE MODEL



Initial Model: Perceptron

Biological neuron

- Before we study artificial neurons, let's look at a biological neuron



First, biological neurons

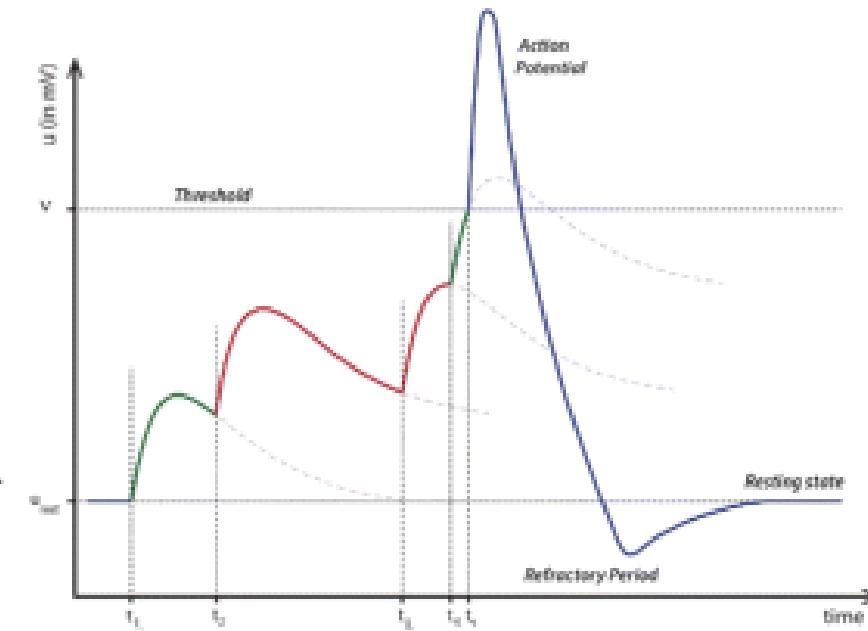
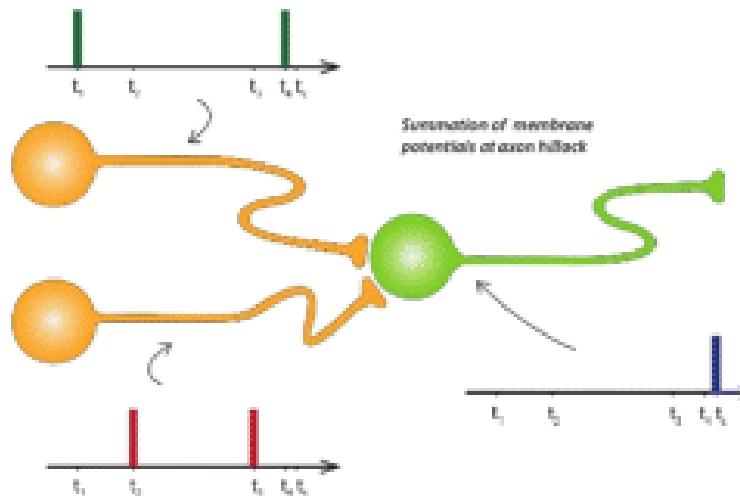
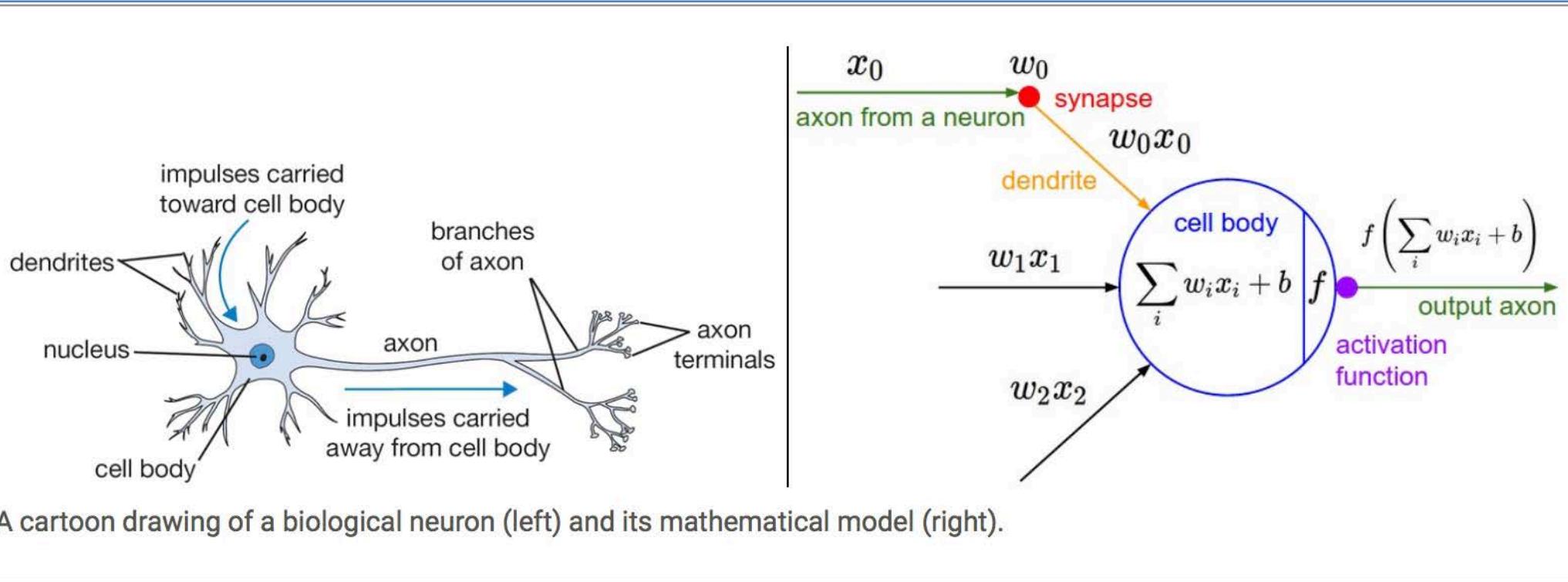


Figure from: Iakymchuk, T., et al. "Simplified spiking neural network architecture and STDP learning algorithm applied to image classification". In Journal of Image Video Proc. 2015, 4 (2015).

Then, artificial neurons

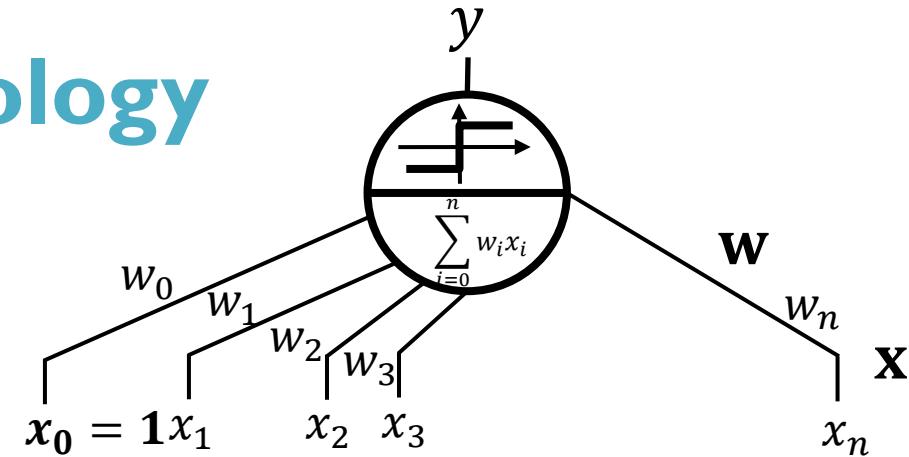


Pitts & McCulloch (1943), binary inputs & activation function f thresholding

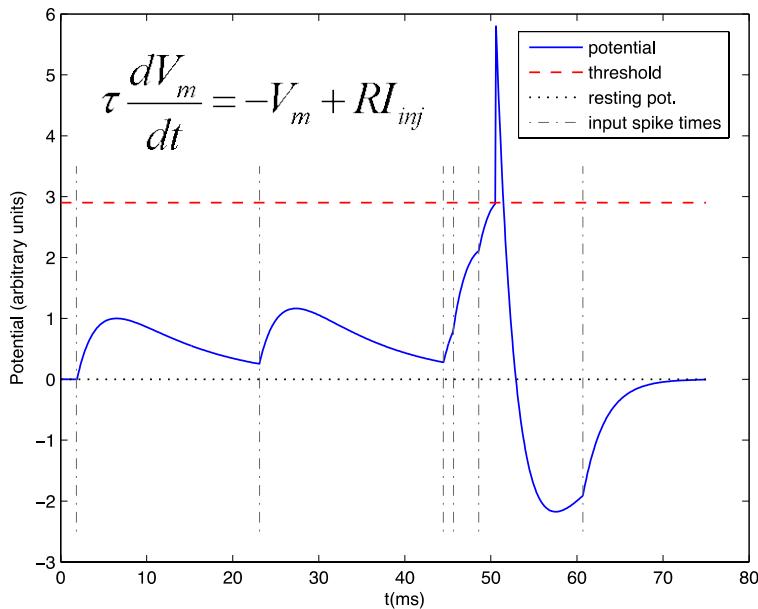
Rosenblatt (1956), real inputs & activation function f thresholding



Artificial vs biology

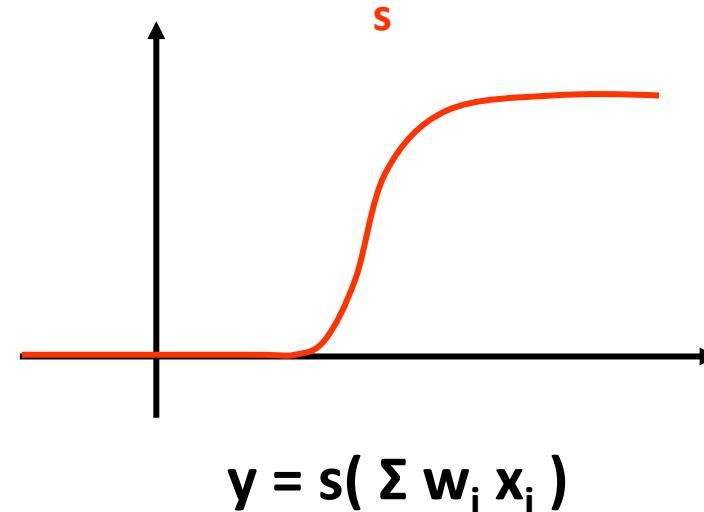


Spike-based description



Gradient descent: KO

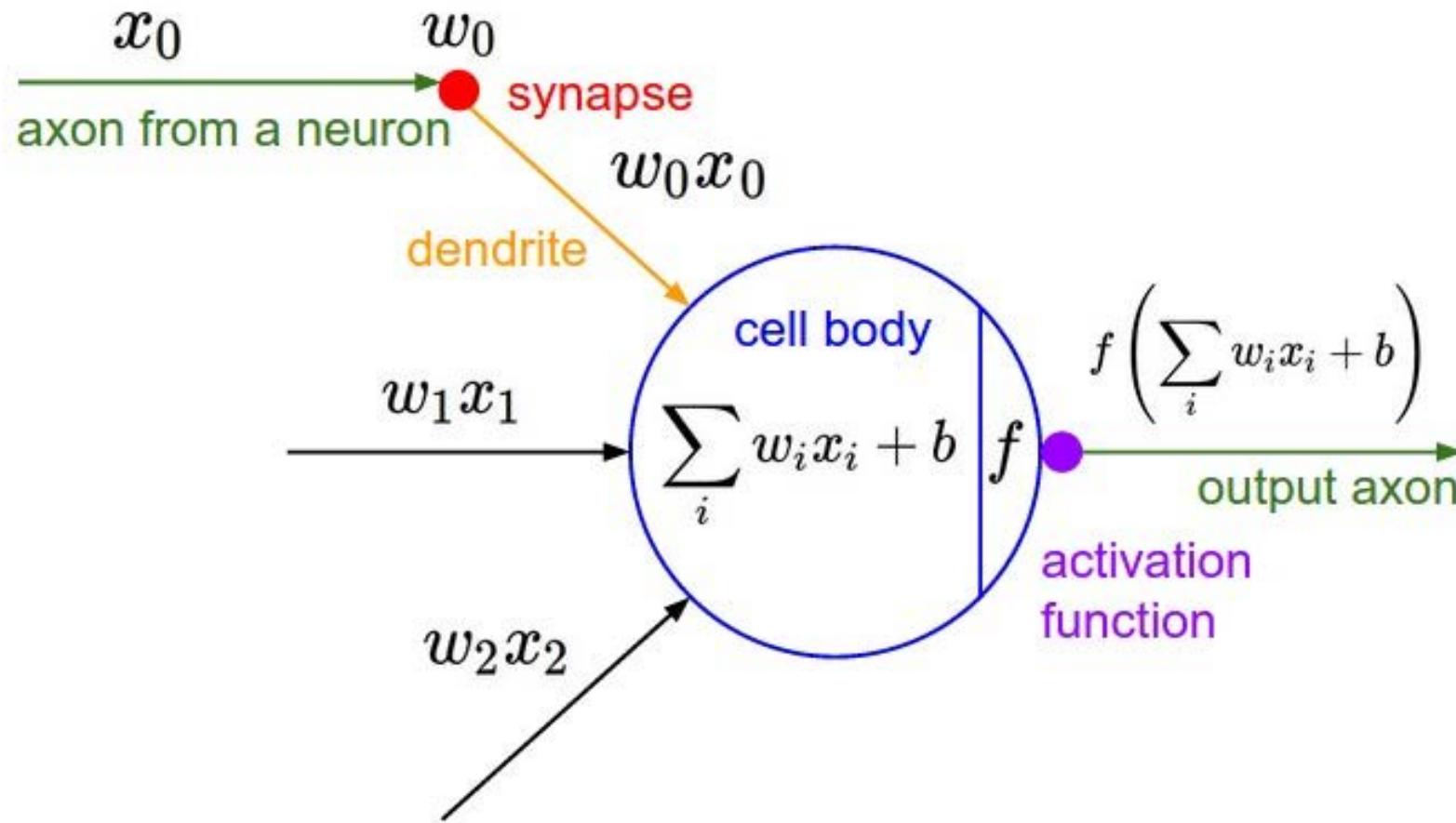
Rate-based description *Steady regime*



Gradient descent: OK

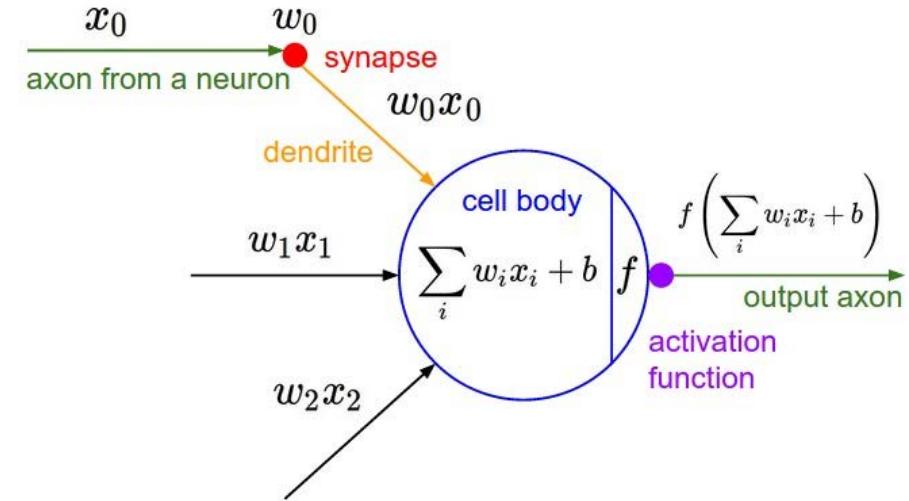
A single artificial neuron

- It is a **very simple abstract** of a biological neuron (McCulloch & Pitts, 1943)



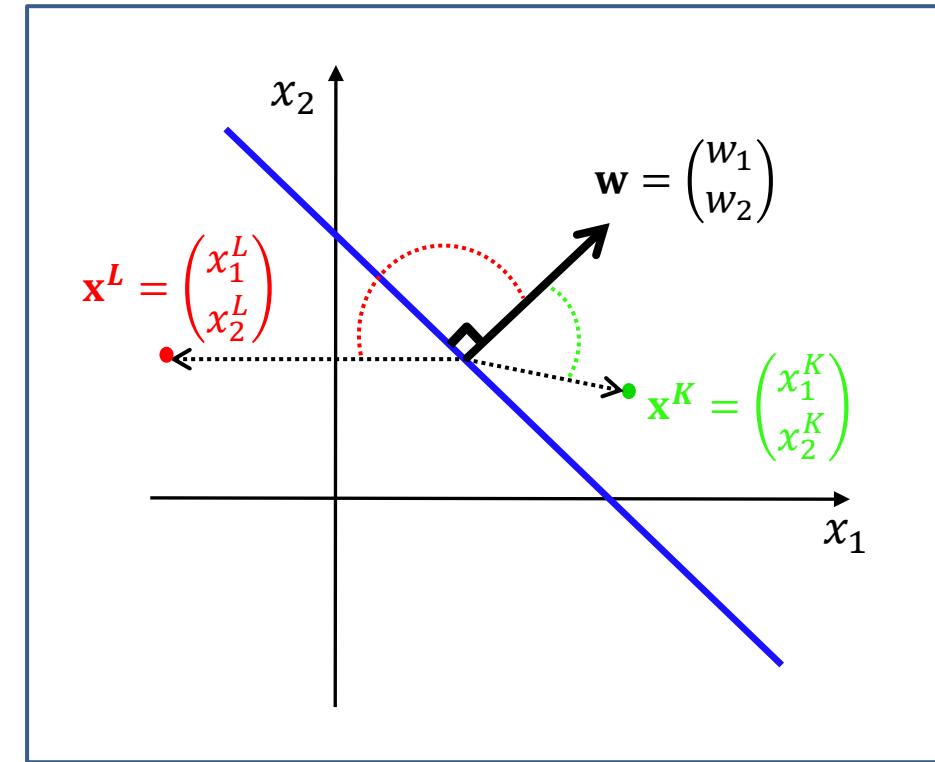
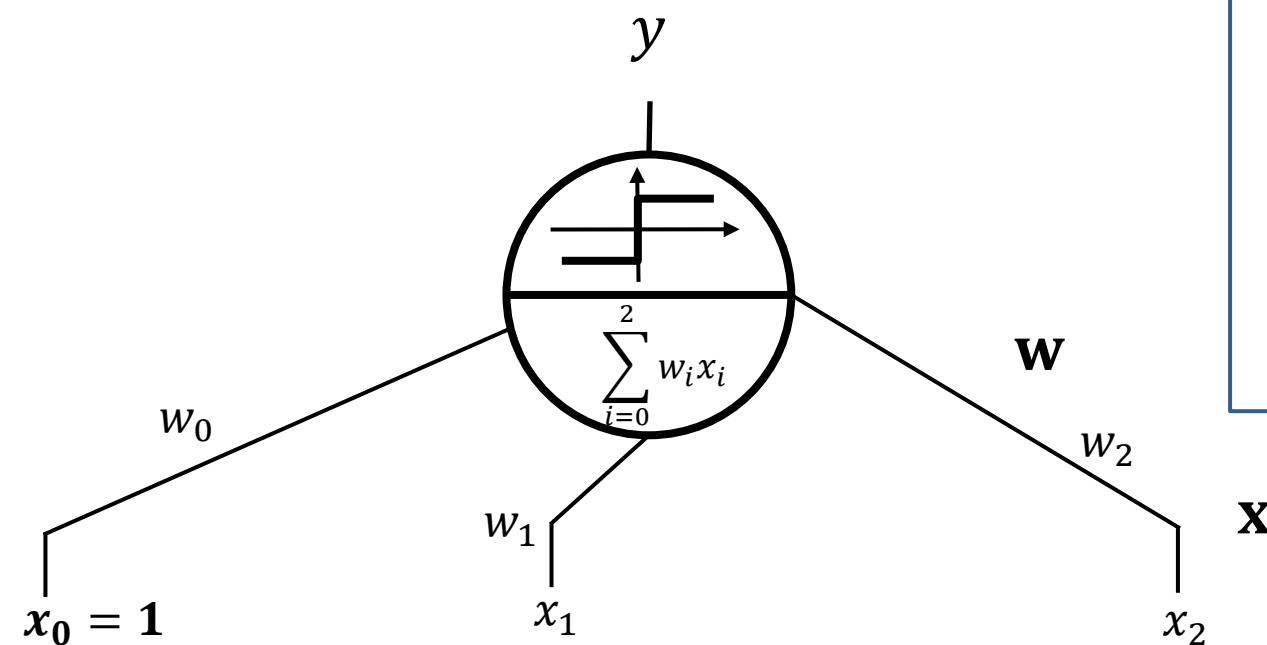
A single artificial neuron

- Each input x has an associated **weight w** which can be modified
- Inputs x corresponds to signals from other neuron axons
 - x_0 - Bias are ‘special’ inputs, with weight w_0
- Weights **W** corresponds to synaptic modulation (i.e. something like strength/amount of neurotransmitters)
- The summation corresponds to ‘cell body’
- The activation function corresponds to axon hillock - computes some function f of the weighted sum of its inputs
- So, output $y=f(z)$, corresponds to axon signal



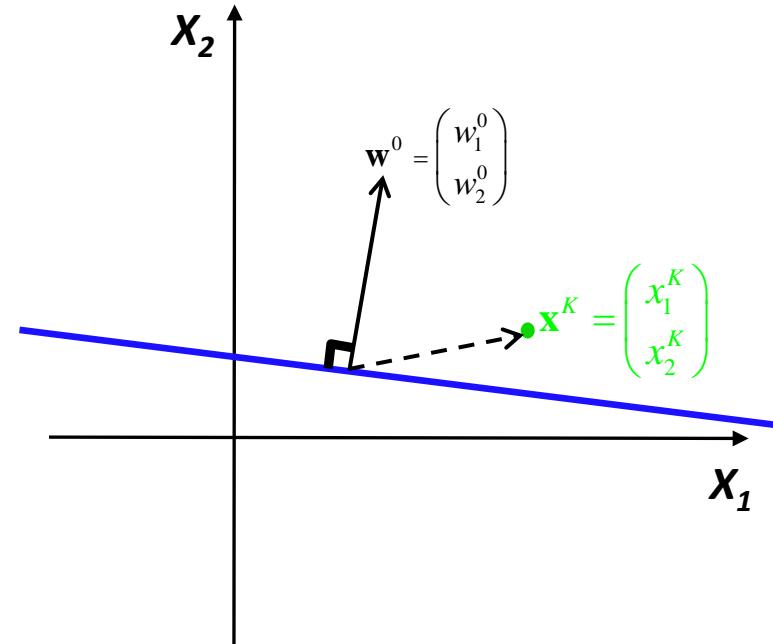
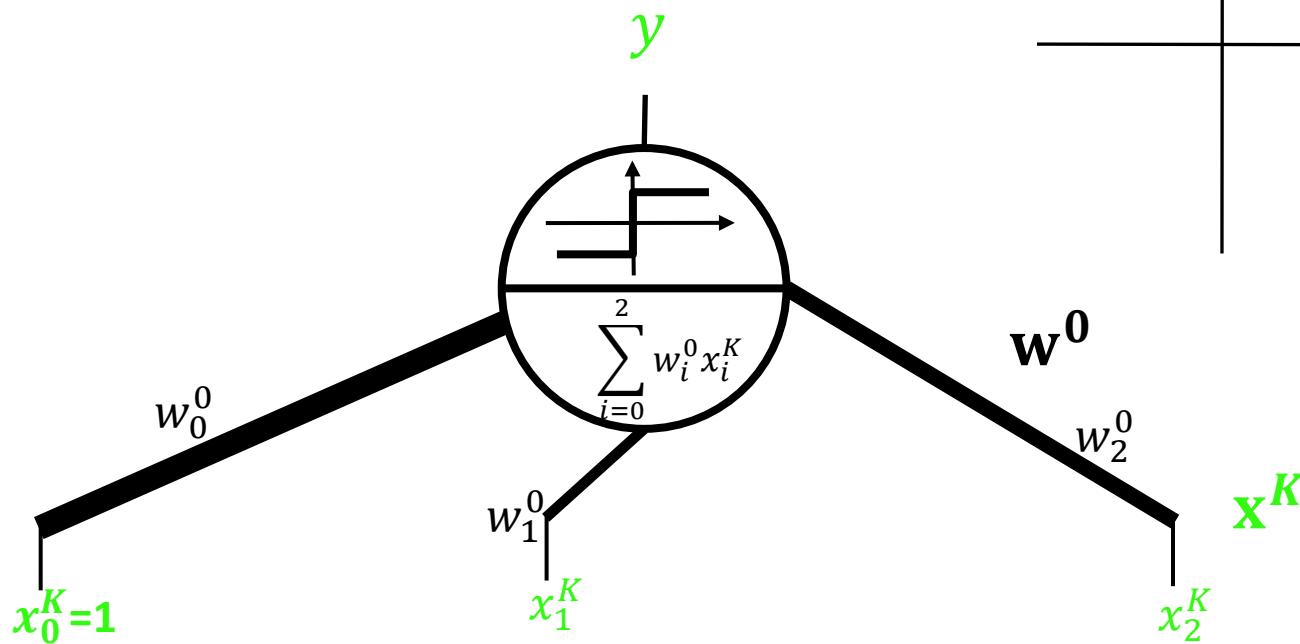


Artificial neuron = a linear classifier



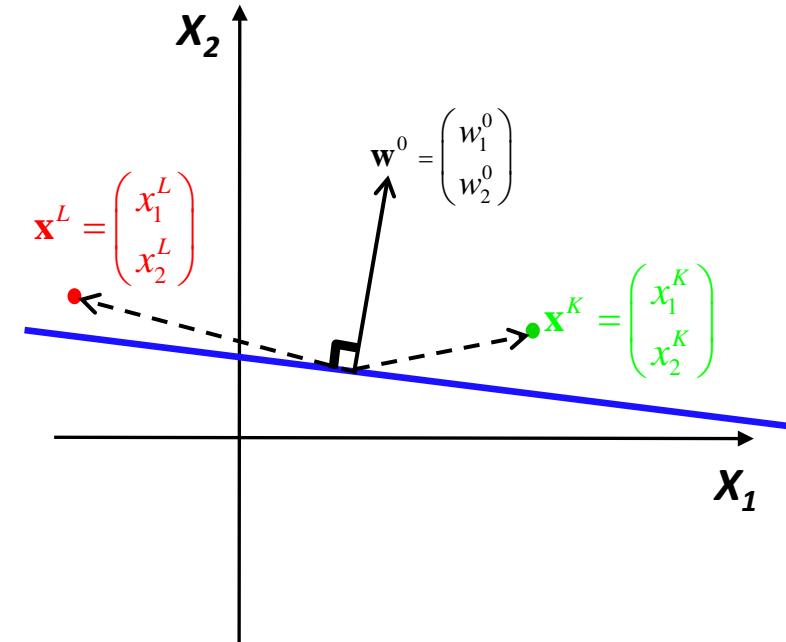
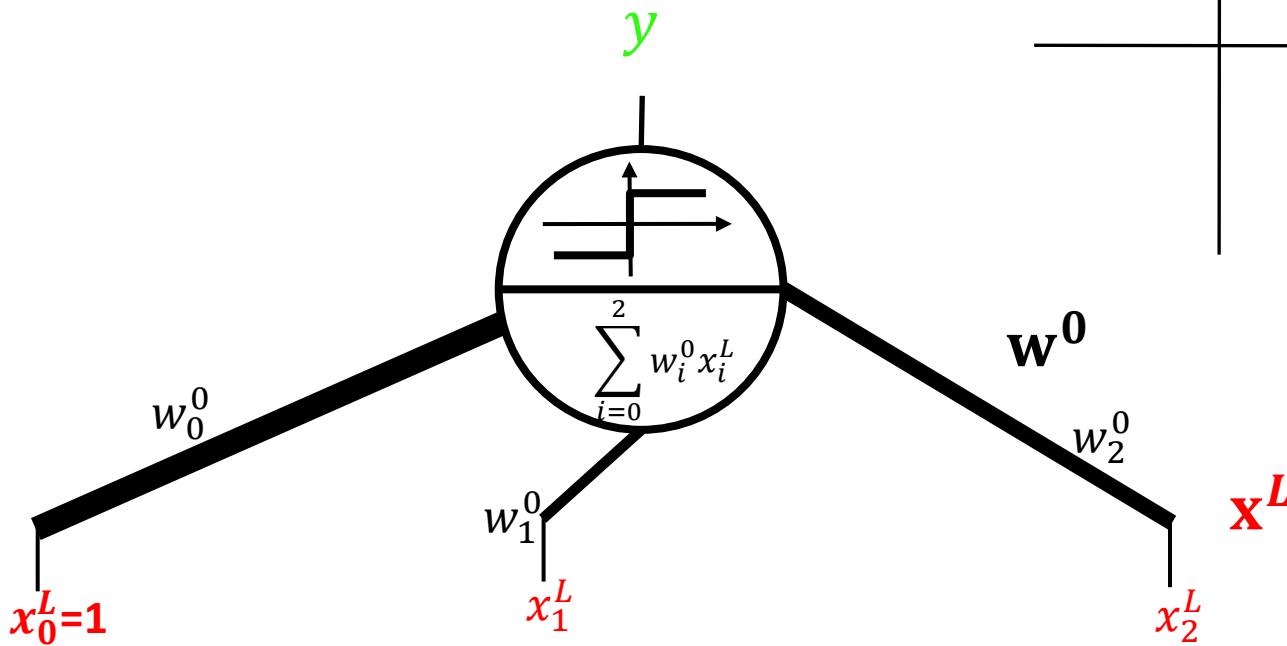
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:



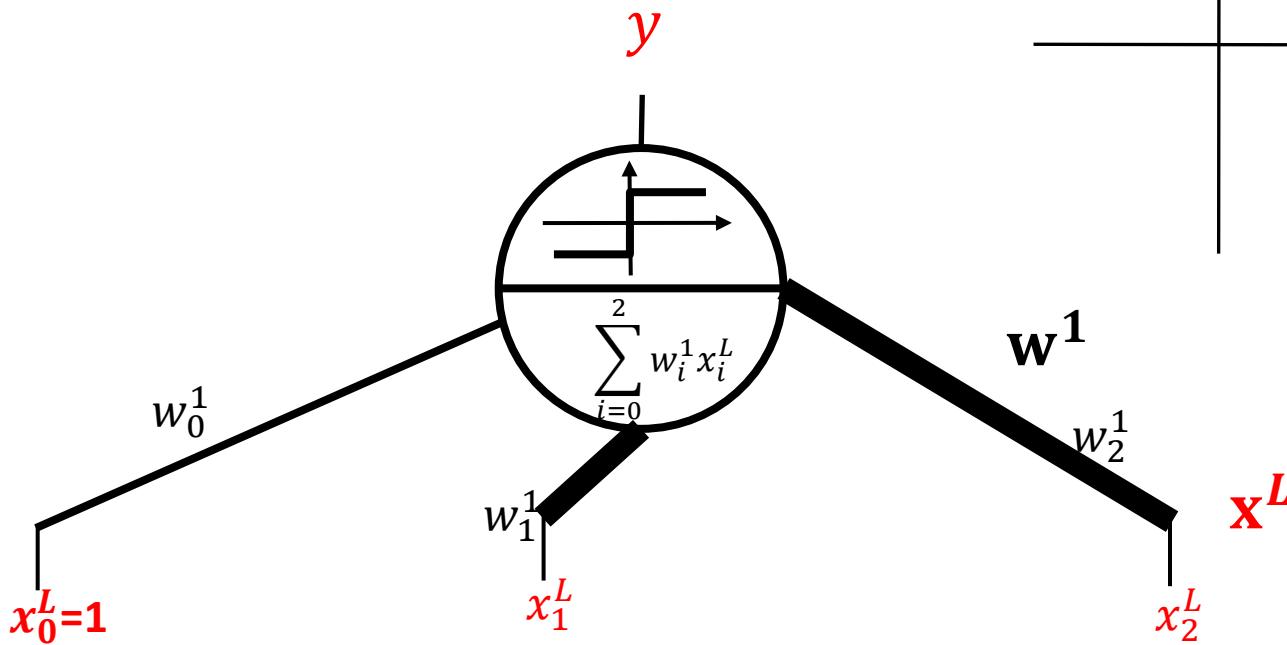
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:



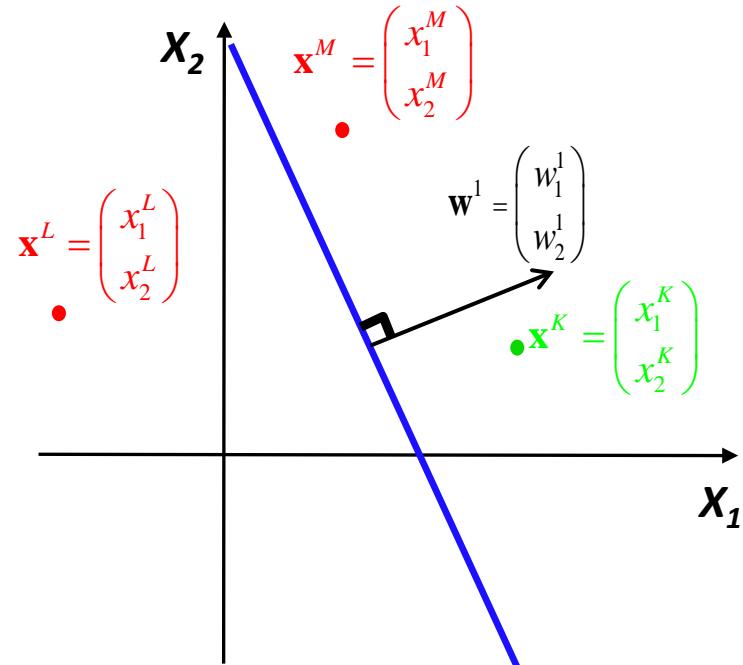
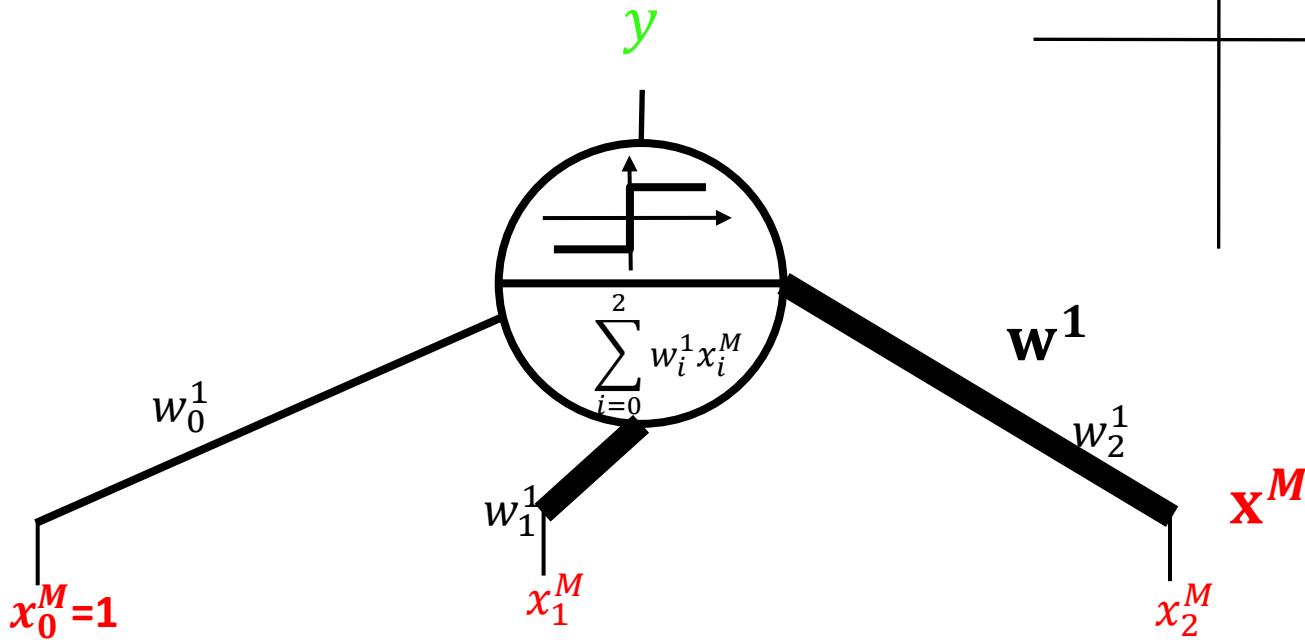
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:



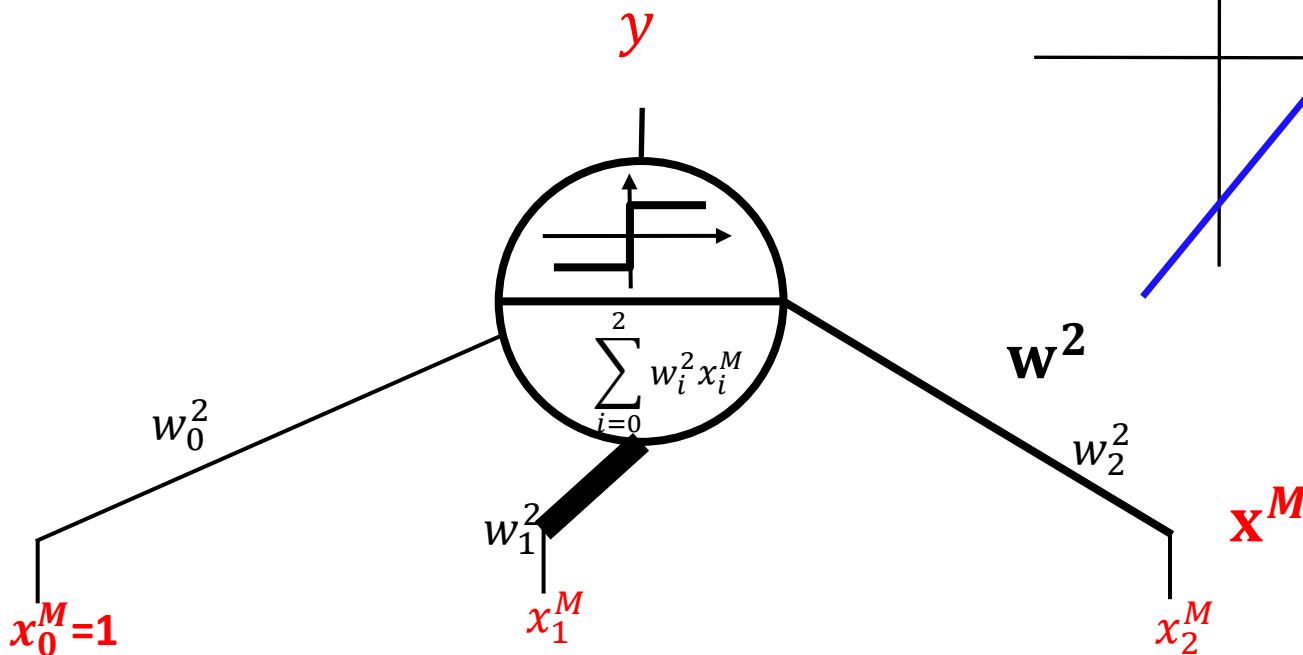
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:



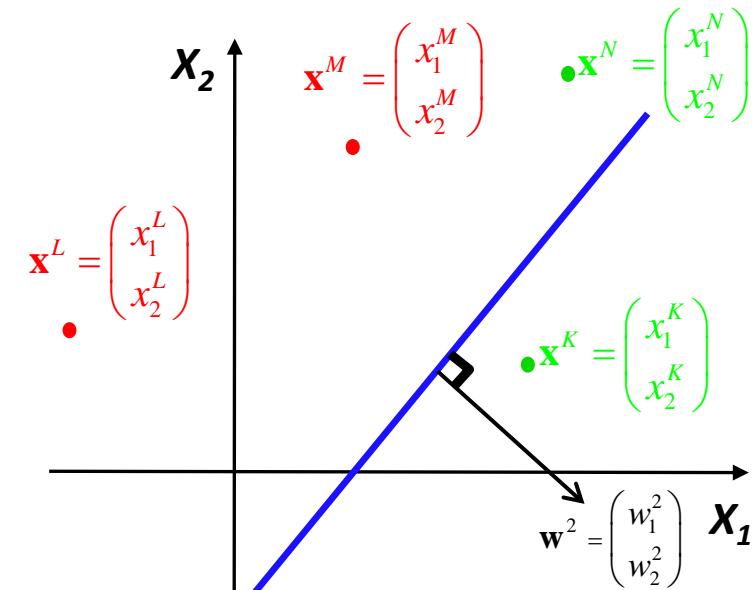
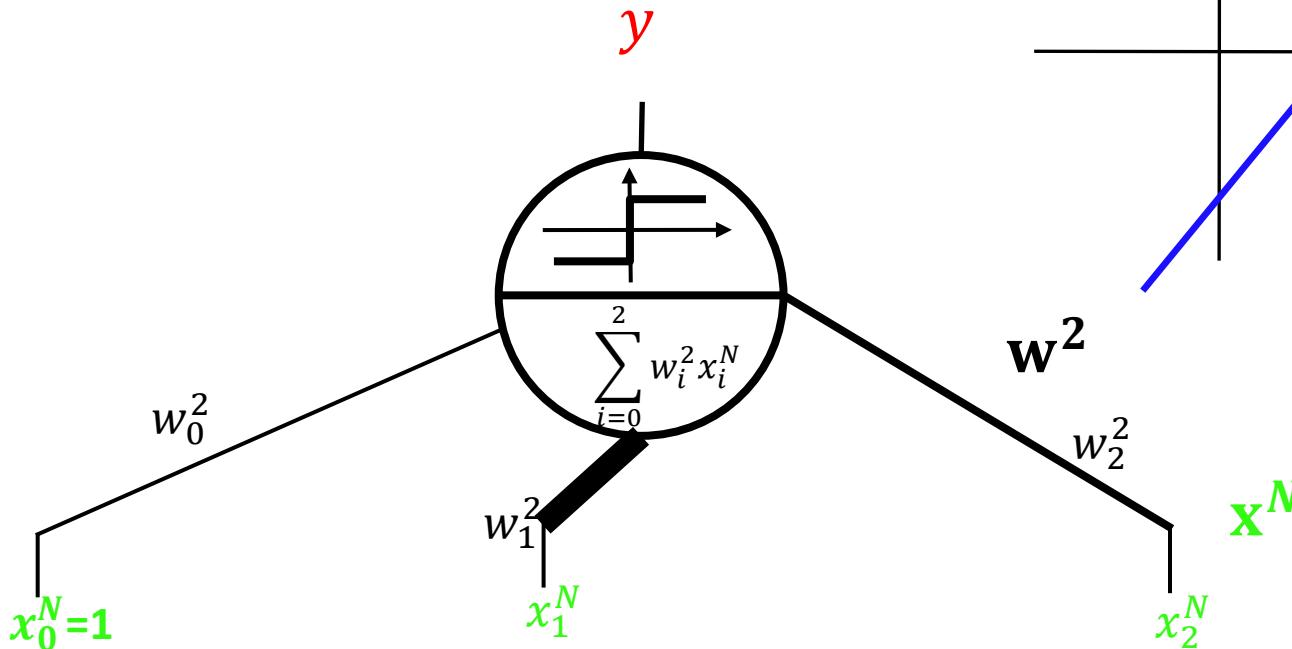
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:



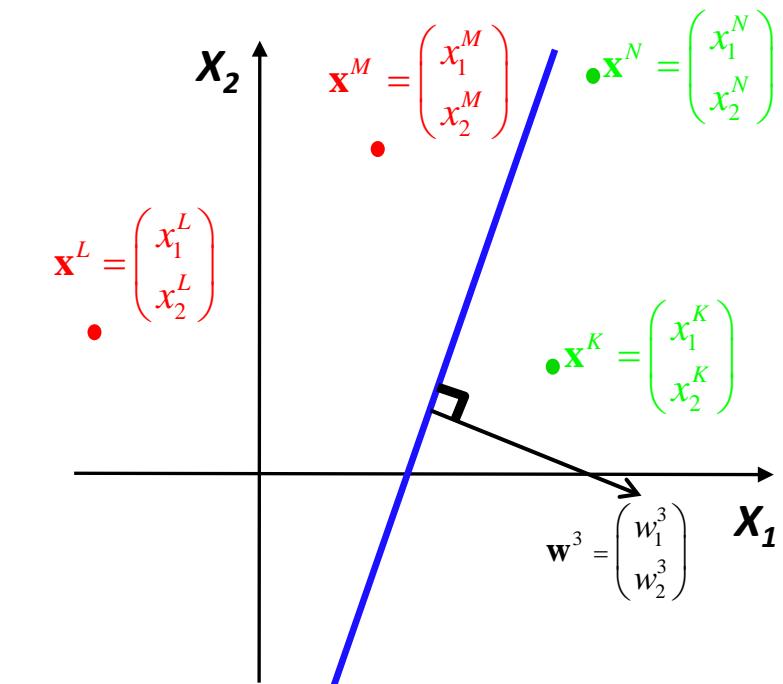
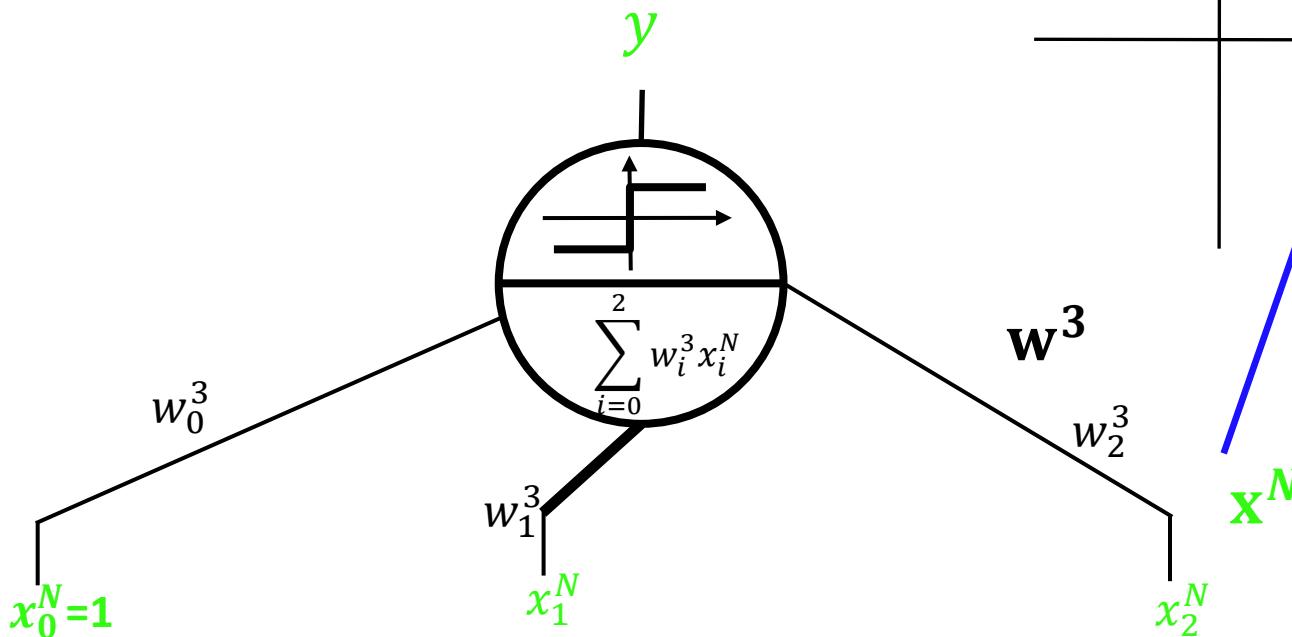
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:



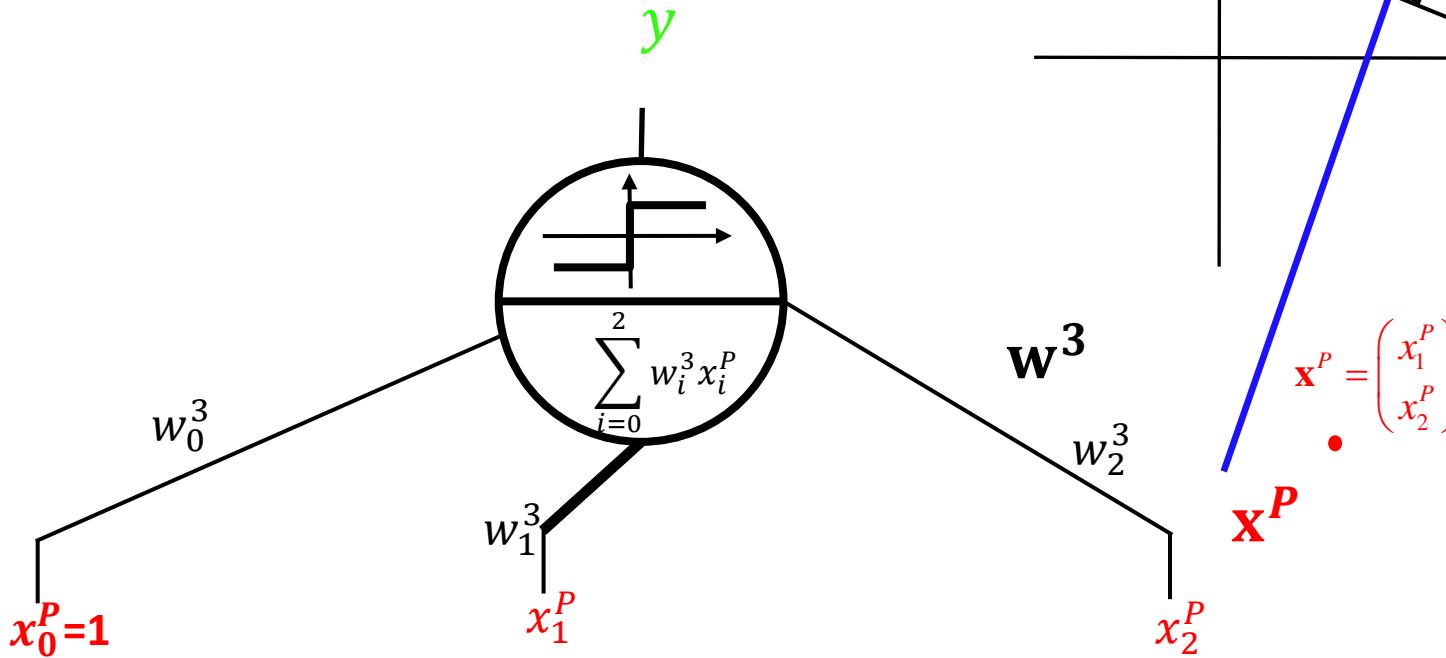
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:



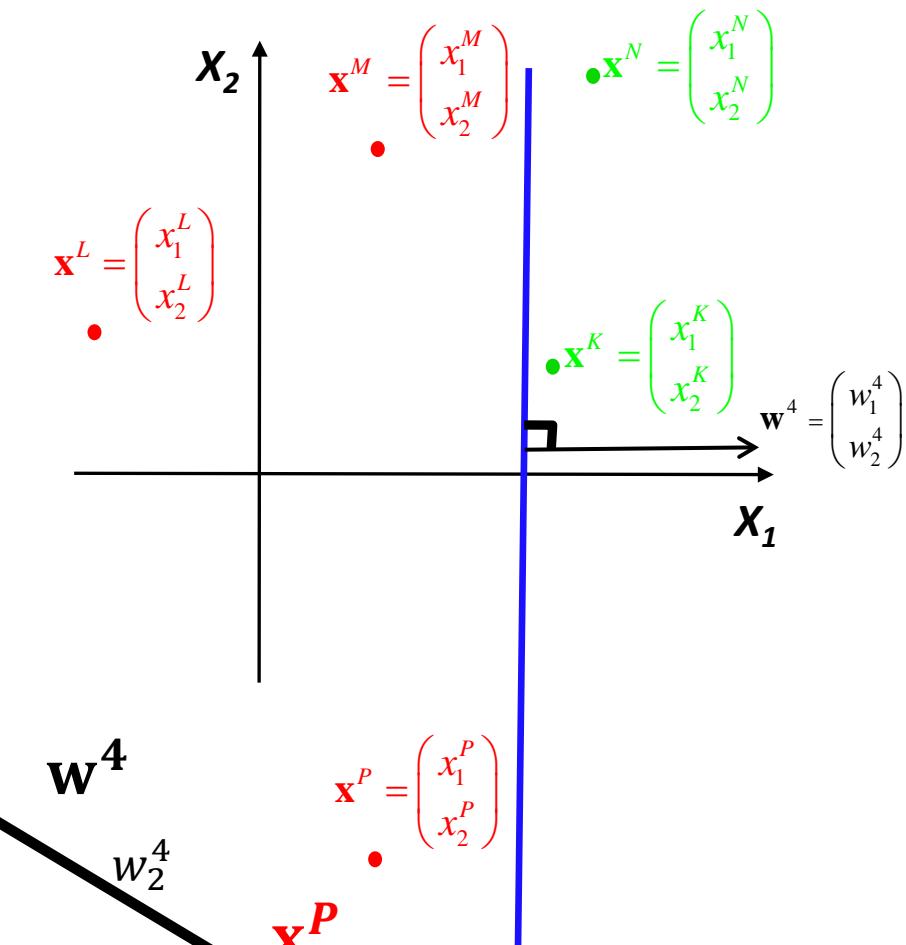
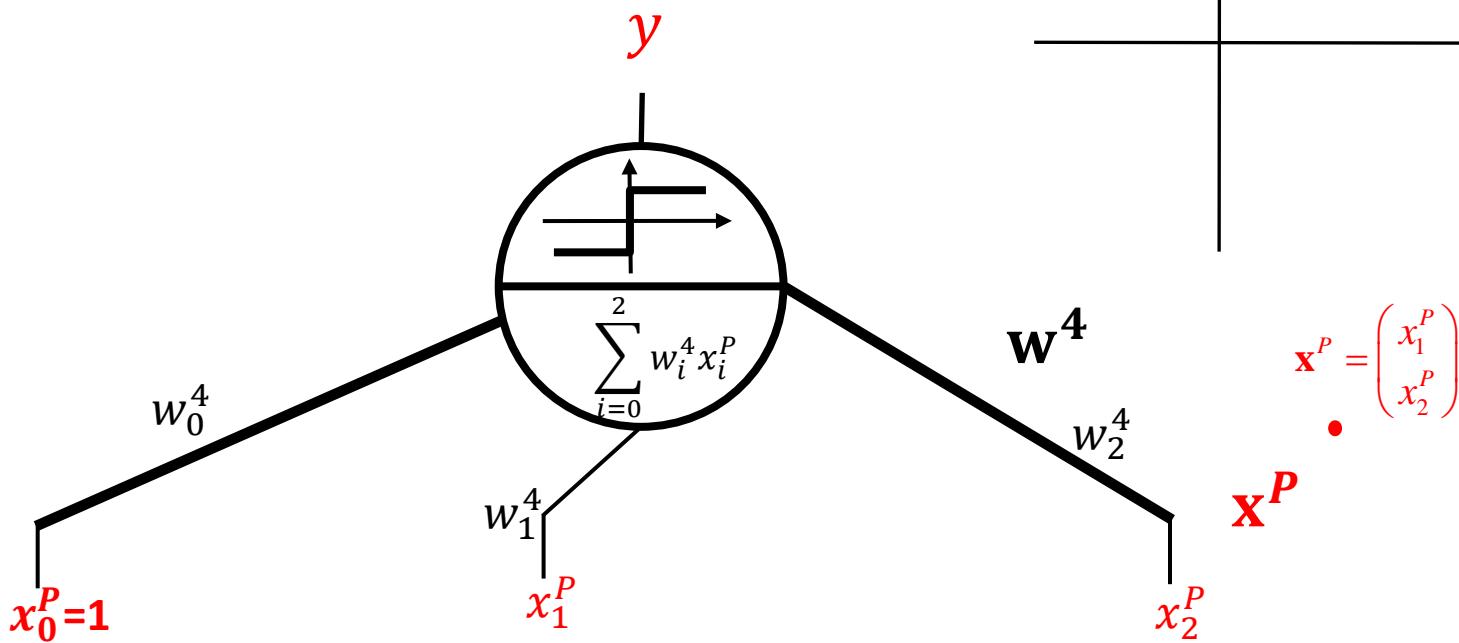
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:



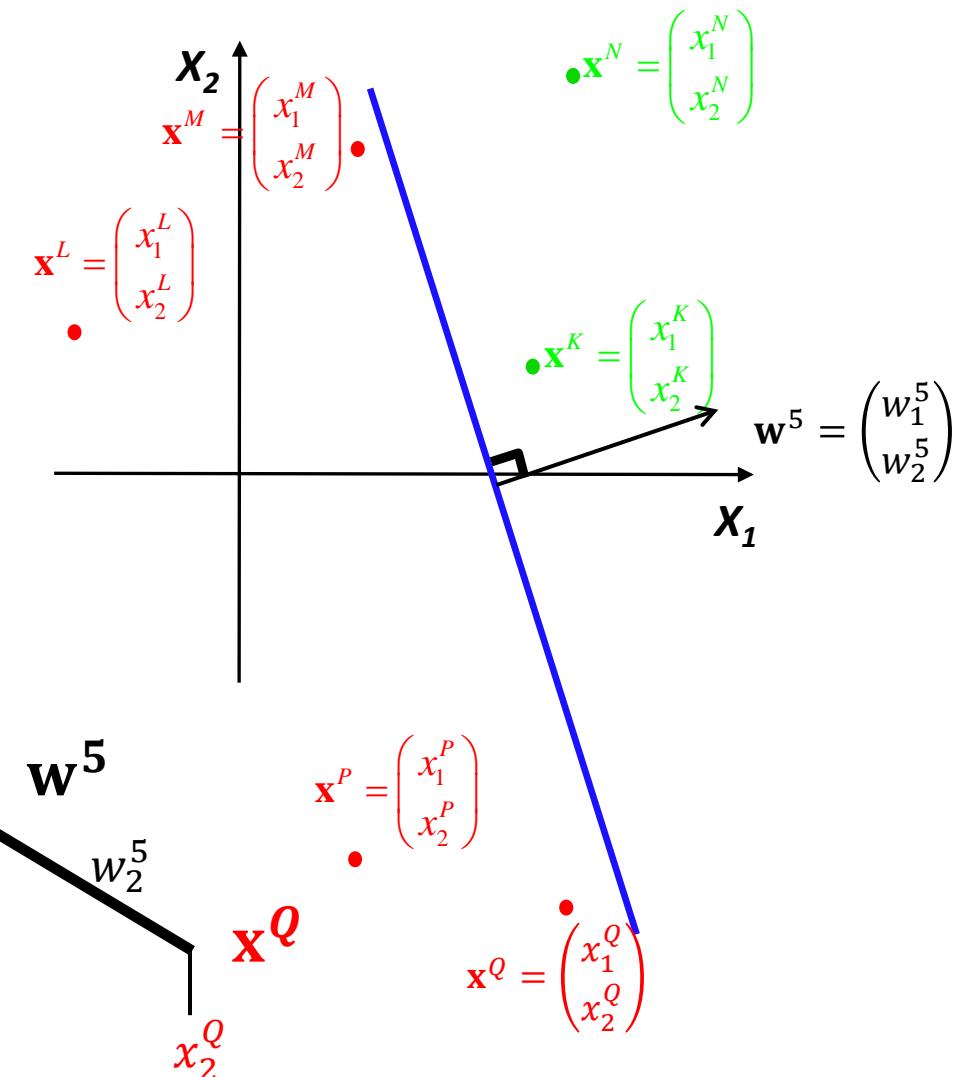
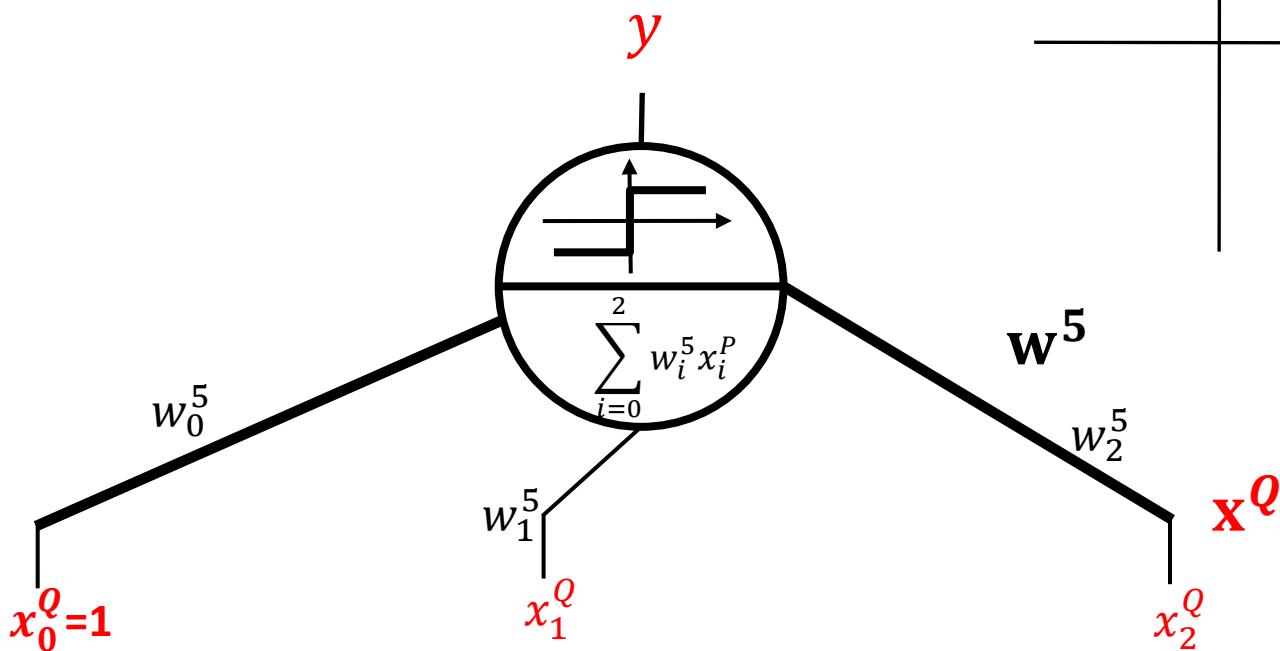
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:

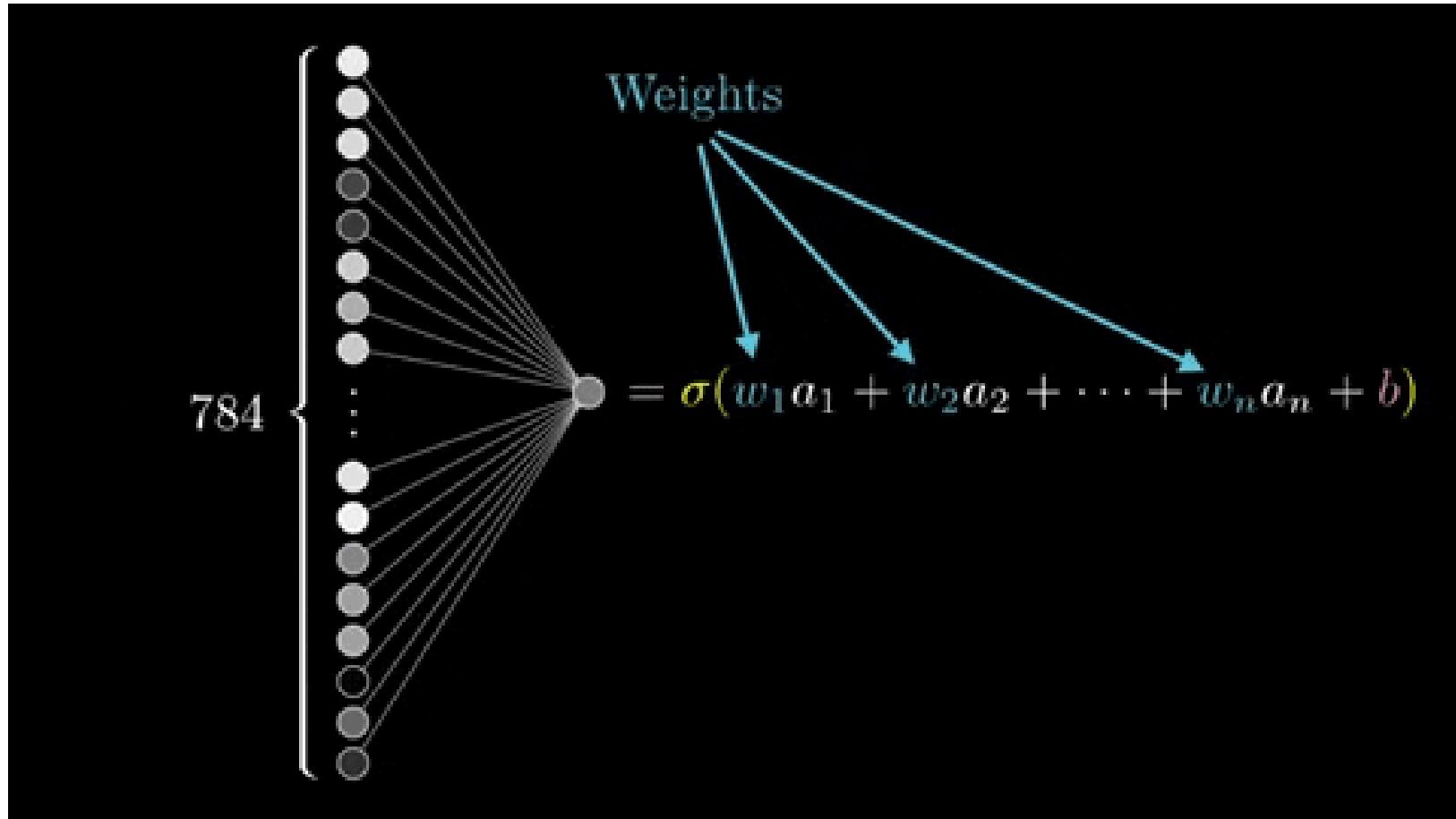


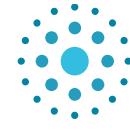
One artificial neuron (perceptron)

At training you want to set the weights,
so that your training samples are correctly classified:



One neuron = simple linear decision





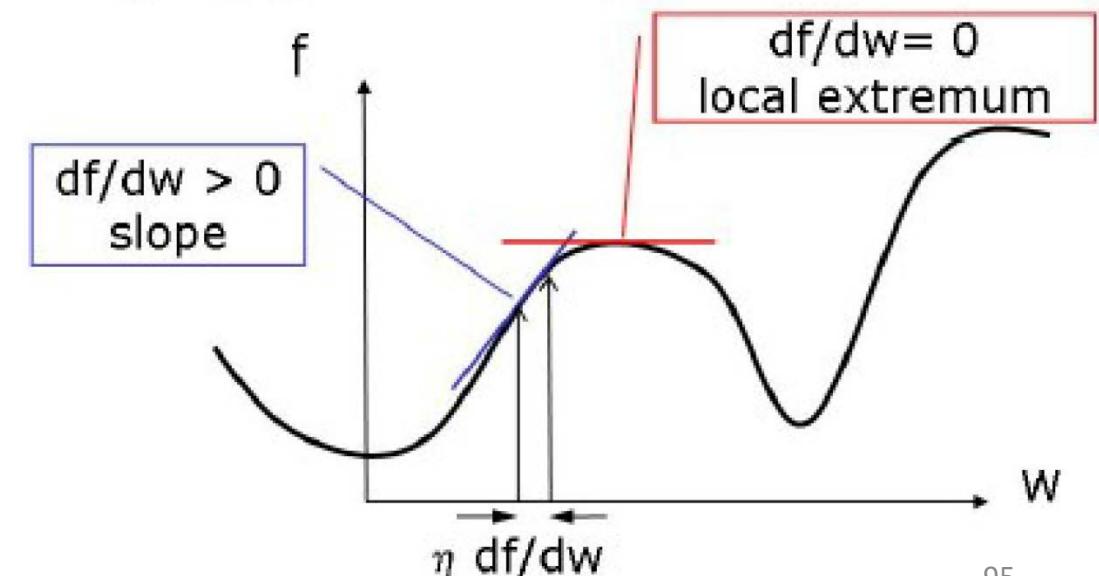
Perceptron: Rosenblatt's Algorithm (1956-1958)

Perceptron Algorithm

- Pick initial weight vector (including w_0), e.g. (0, 0,...,0)
- Repeat until all points are correctly classified
 - *Repeat for each point*
 - Calculate $y^i \mathbf{w} \mathbf{x}^i$ for point i
 - If $y^i \mathbf{w} \mathbf{x}^i > 0$, the point is correctly classified
 - Else change the weights to increase the value of $y^i \mathbf{w} \mathbf{x}^i$; change in weight proportional to $y^i \mathbf{x}^i$

Gradient Ascent

- Why pick $y^i \mathbf{x}^i$ as increment to weights?
- To maximize scalar function of one variable $f(\mathbf{w})$
 - Pick initial \mathbf{w}
 - Change \mathbf{w} to $\mathbf{w} + \eta \frac{df}{d\mathbf{w}}$ ($\eta > 0$, small)
 - Until f stops changing ($\frac{df}{d\mathbf{w}} \approx 0$)



Gradient Ascent

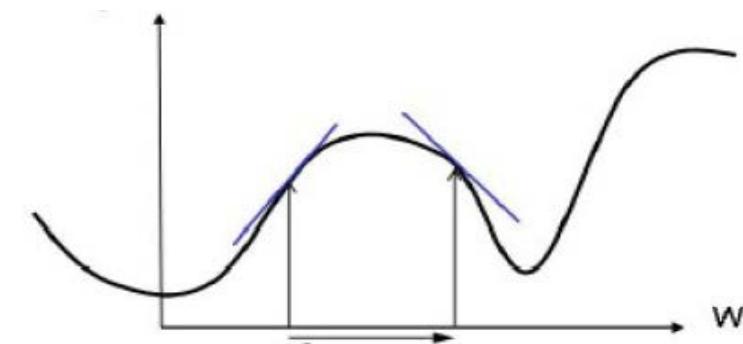
- To maximize a multivariate function $f(\mathbf{w})$
 - Pick initial \mathbf{w}
 - Change \mathbf{w} to $\mathbf{w} + \eta \nabla f_{\mathbf{w}}$ ($\eta > 0$, small)
 - Until f stops changing ($\nabla f_{\mathbf{w}} \approx 0$)
- Find local maximum, unless function is globally convex

$$\nabla f_{\mathbf{w}} = \left[\frac{\partial f}{\partial \mathbf{w}_1}, \quad \dots \quad , \frac{\partial f}{\partial \mathbf{w}_n} \right]$$

Gradient Ascent

- To maximize a multivariate function $f(\mathbf{w})$
 - Pick initial \mathbf{w}
 - Change \mathbf{w} to $\mathbf{w} + \eta \nabla f_{\mathbf{w}}$ ($\eta > 0$, small)
 - Until f stops changing ($\nabla f_{\mathbf{w}} \approx 0$)
- Find local maximum, unless function is globally convex
- If f is non-linear, the learning rate η has to be chosen very carefully
 - Too small \Rightarrow slow convergence
 - Too big \Rightarrow oscillations

$$\nabla f_{\mathbf{w}} = \left[\frac{\partial f}{\partial \mathbf{w}_1}, \dots, \frac{\partial f}{\partial \mathbf{w}_n} \right]$$



Gradient Ascent

- Maximize margin of misclassified points

$$f(\mathbf{w}) = \sum_{\text{on } \mathbf{i} \text{ misclassified points}} y^i \mathbf{w} \mathbf{x}^i$$

$$\nabla_{\mathbf{w}} f(\mathbf{w}) = \sum_{\text{on } \mathbf{i} \text{ misclassified points}} y^i \mathbf{x}^i$$

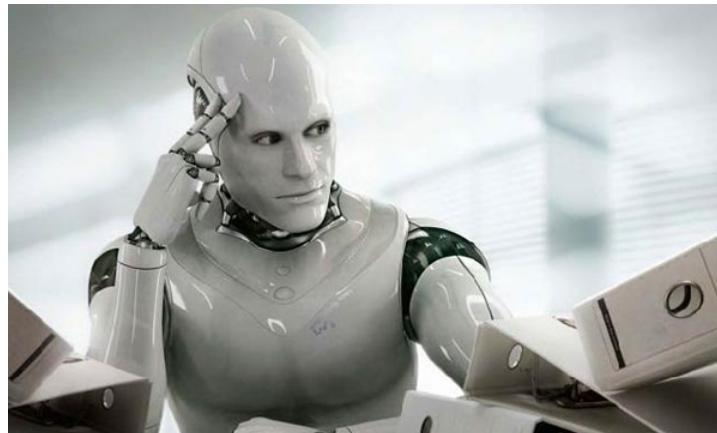
- Off-line training: Compute, at each iteration, the gradient as sum over all training points
- On-line training: Approximate gradient by one of the terms in the sum: $y^i \mathbf{x}^i$
(principle of the *Stochastic Gradient Descent*, SGD)

Gradient Descent

This is the “learning” of machines in deep learning

→ Even alpha go using this approach.

What People imagine



In reality

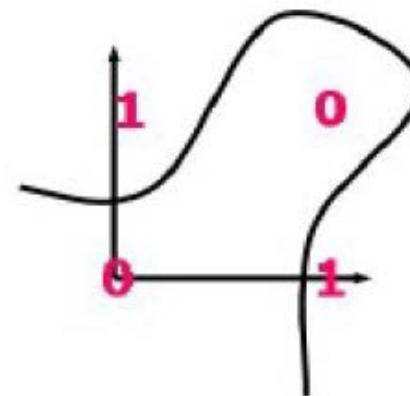
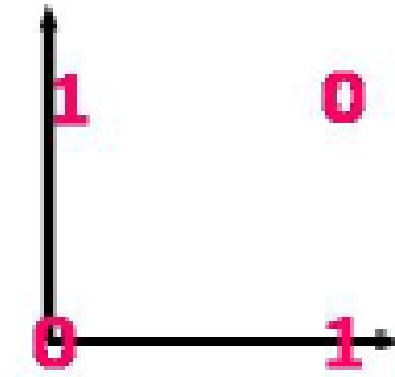


Perceptron Algorithm

- Each change of w decreases the error on a specific point. However, changes for several points are correlated, that is different points could change the weights in opposite directions. Thus, this iterative algorithm requires several loops to converge.
- Guarantee to find a separating hyperplane if one exists – if data is linearly separable.
- If data are not linearly separable, then this algorithm loops indefinitely.

Beyond Linear Separability

- Values of the XOR boolean function cannot be separated by a single perceptron unit [Minsky and Papert, 1969].

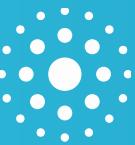


Minsky, M. and Papert, S. (1969). *Perceptrons: An Introduction to Computational Geometry*. MIT Press.



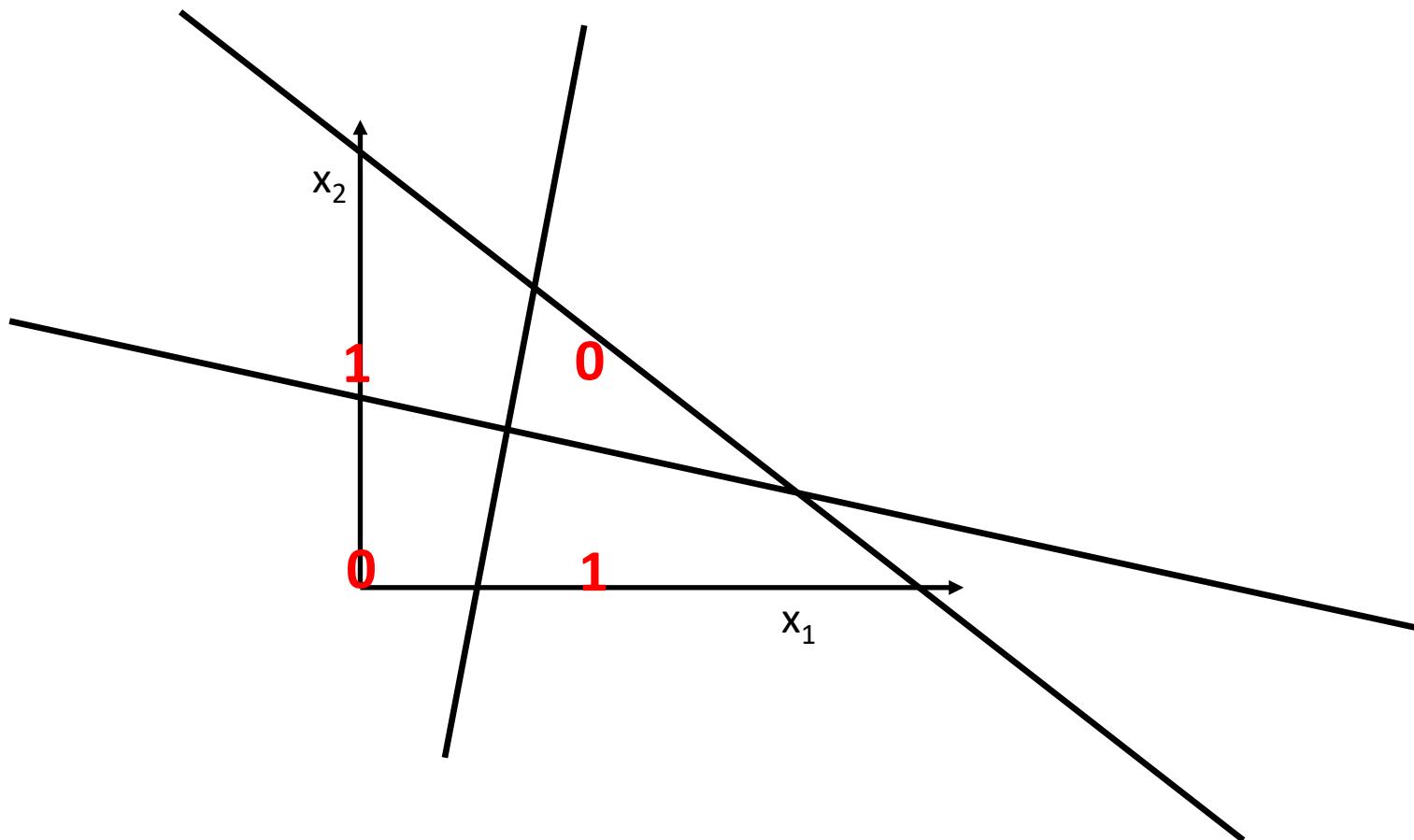
Overview

- Machine Learning vs Statistics
- Math Basics
- Simple Model
- **From Simple to Complex**

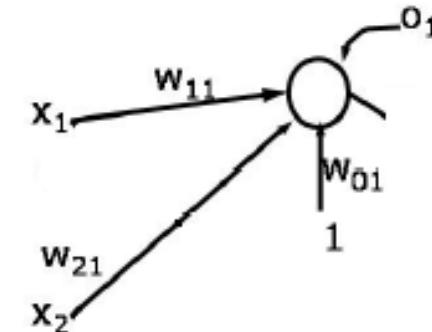
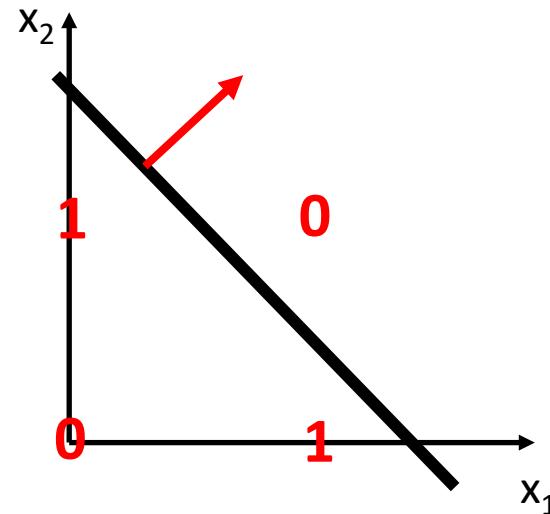


FROM SIMPLE TO COMPLEX

Problem which cannot be solved with a unique straight line

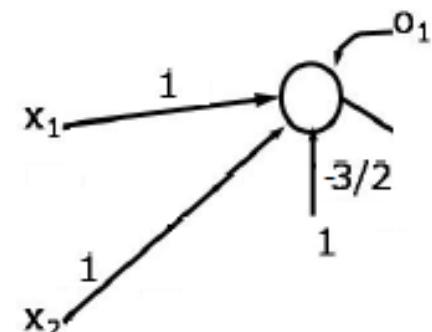


Problem which cannot be solved with a unique straight line

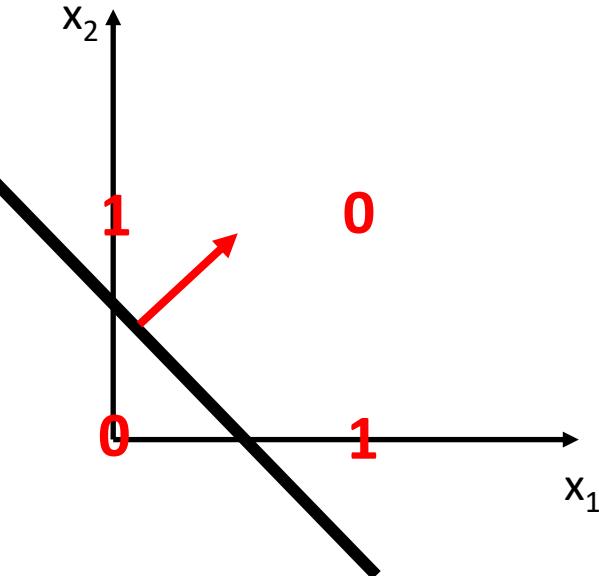


$$w_{01} = -\frac{3}{2} \quad w_{11} = w_{21} = 1$$

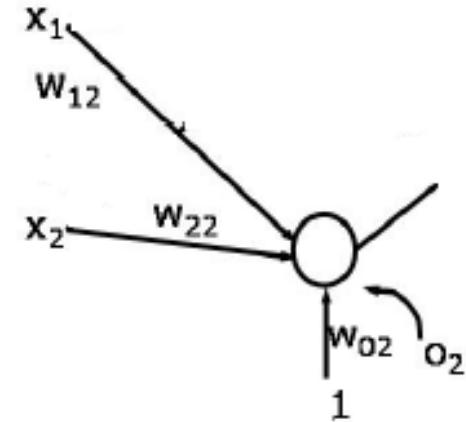
x_1	x_2	o_1
0	0	0
0	1	0
1	0	0
1	1	1



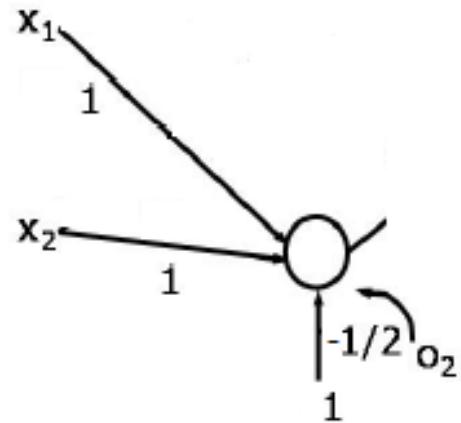
Problem which cannot be solved with a unique straight line



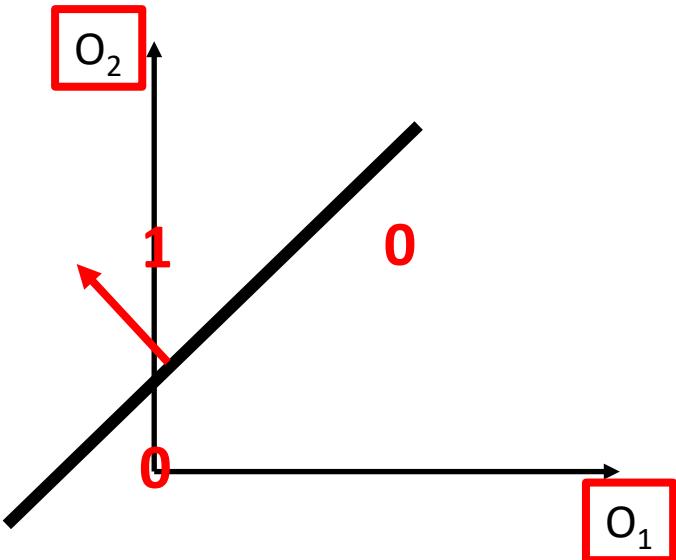
x_1	x_2	o_1	o_2
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	1



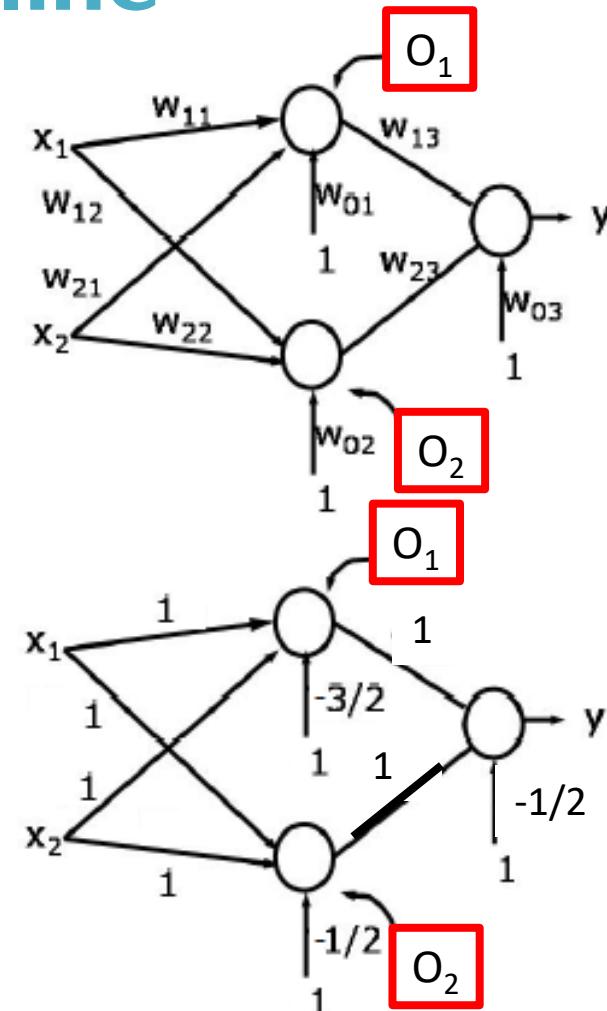
$$w_{02} = -1/2 \quad w_{12} = w_{22} = 1$$



Problem which cannot be solved with a unique straight line



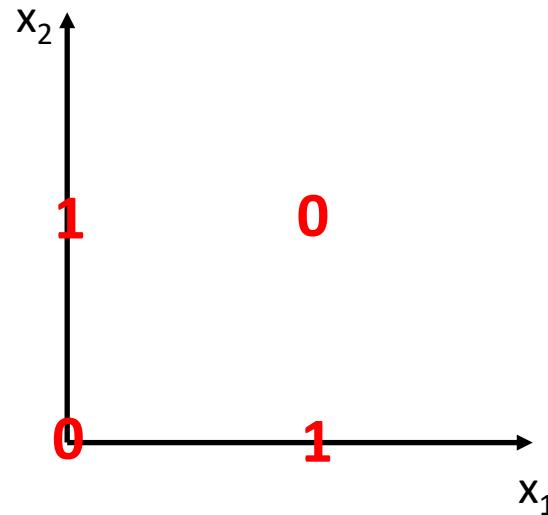
o_1	o_2	y
0	0	0
0	1	1
0	1	1
1	1	0



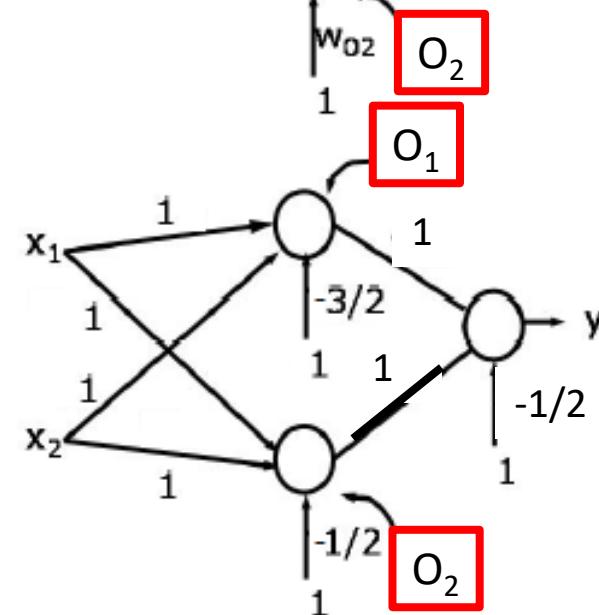
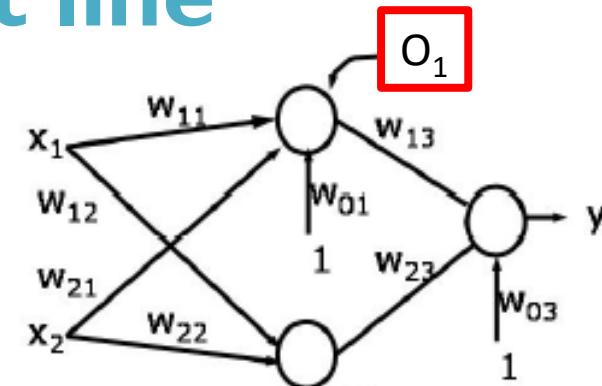
$$w_{23}O_2 + w_{13}O_1 + w_{03} = 0$$

$$w_{03} = -1/2, w_{13} = -1, w_{23} = 1$$

Problem which cannot be solved with a unique straight line



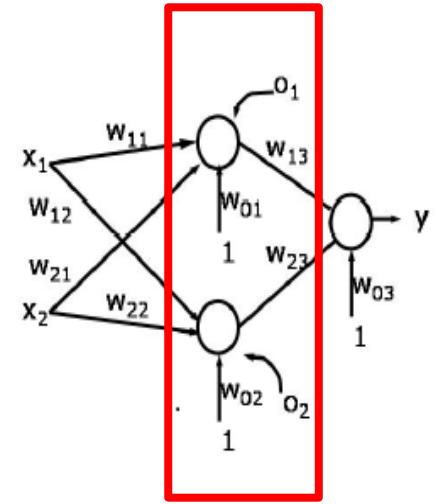
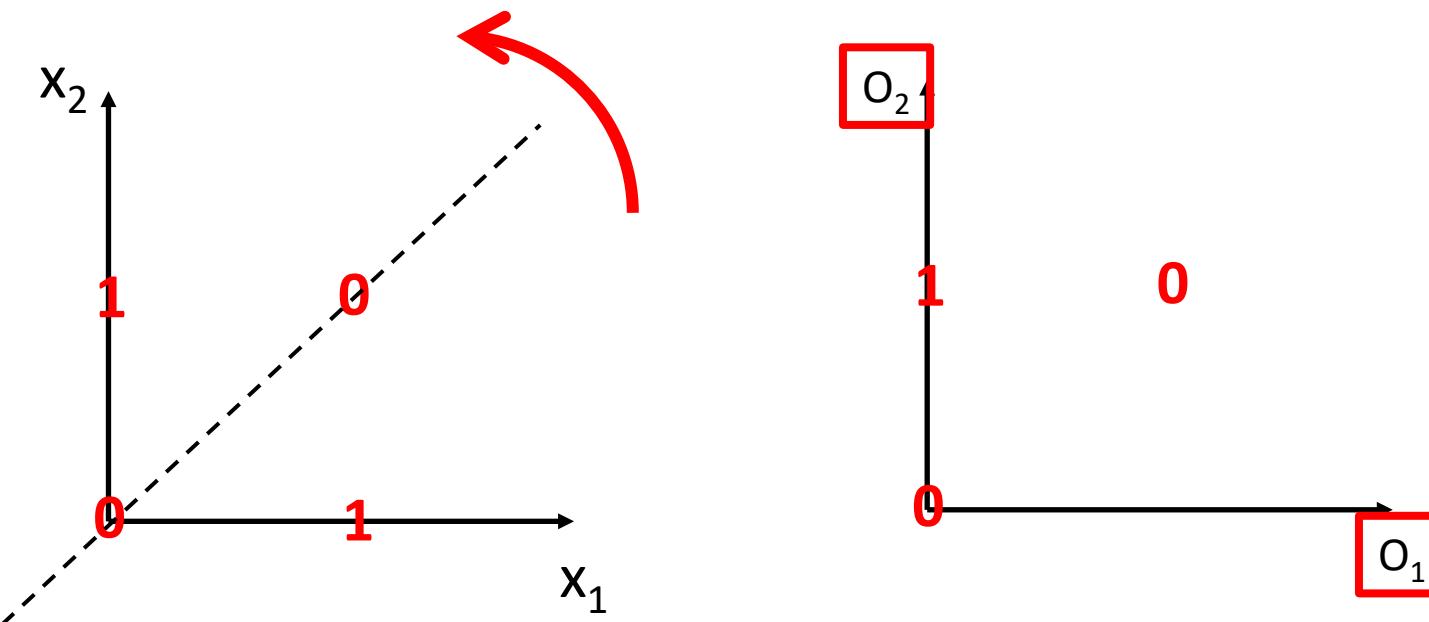
x_1	x_2	o_1	o_2	y
0	0	0	0	0
0	1	0	1	1
1	0	0	1	1
1	1	1	1	0



$$w_{23}O_2 + w_{13}O_1 + w_{03} = 0$$

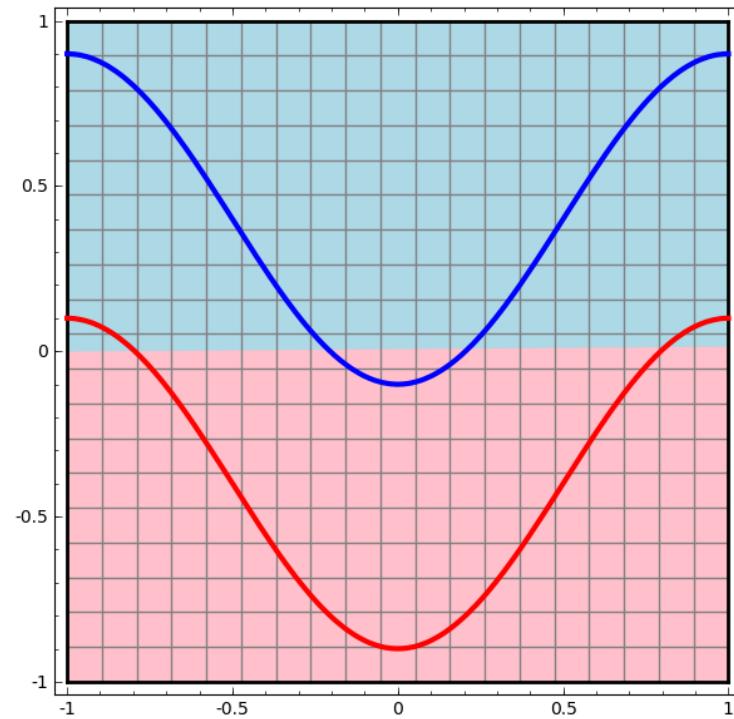
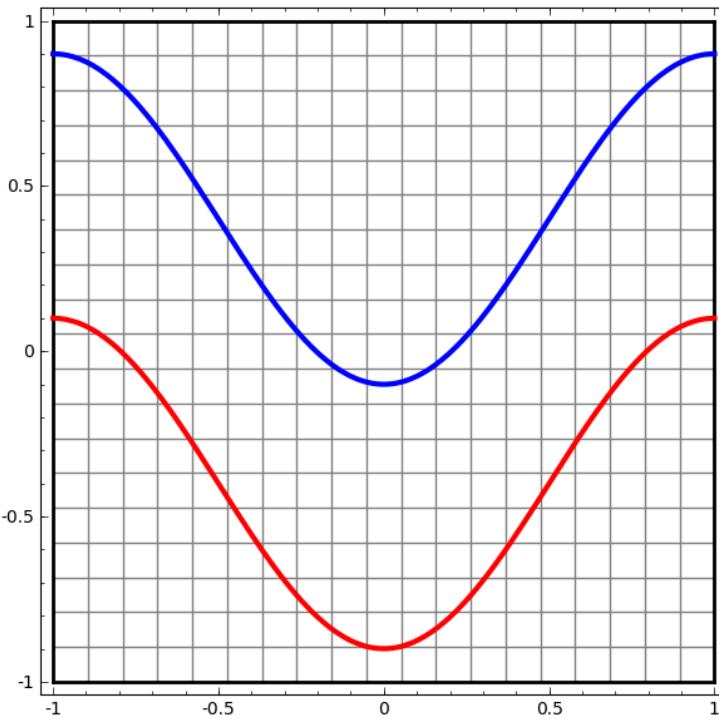
$$w_{03} = -1/2, w_{13} = -1, w_{23} = 1$$

Problem which cannot be solved with a unique straight line



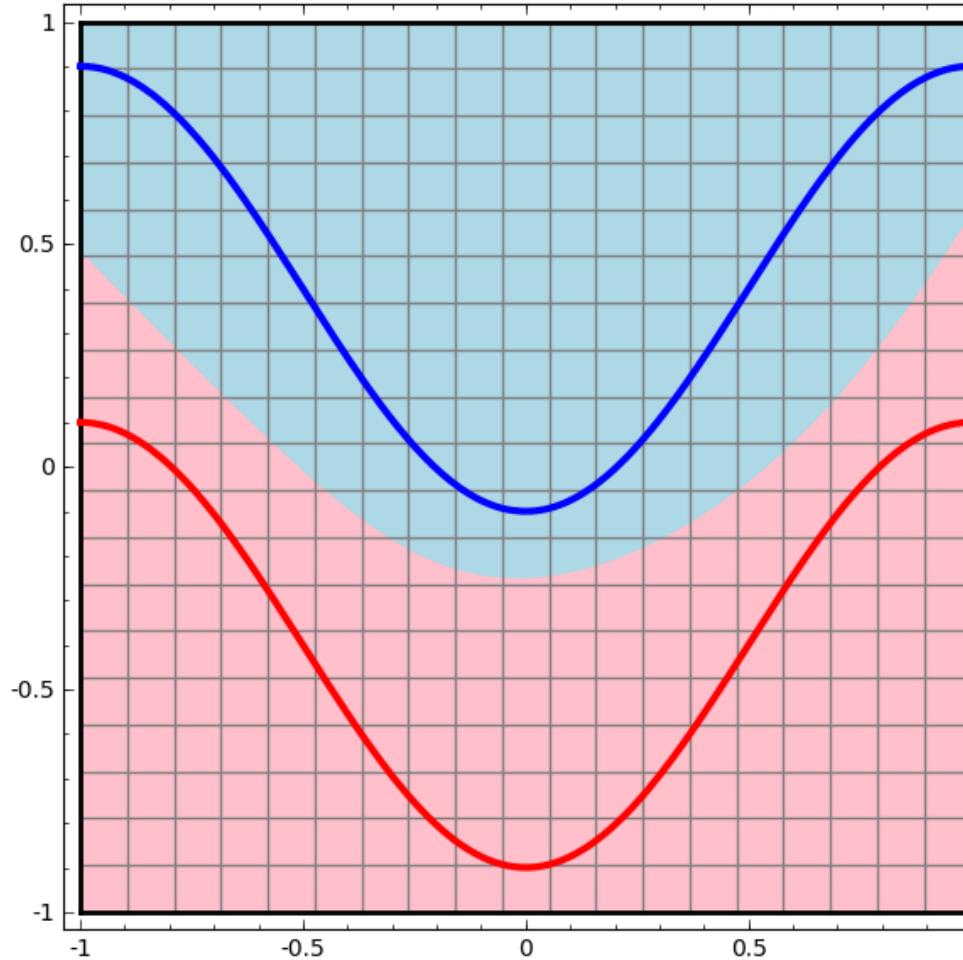
Adding a hidden layer of neurons has fold the input space

One perceptron

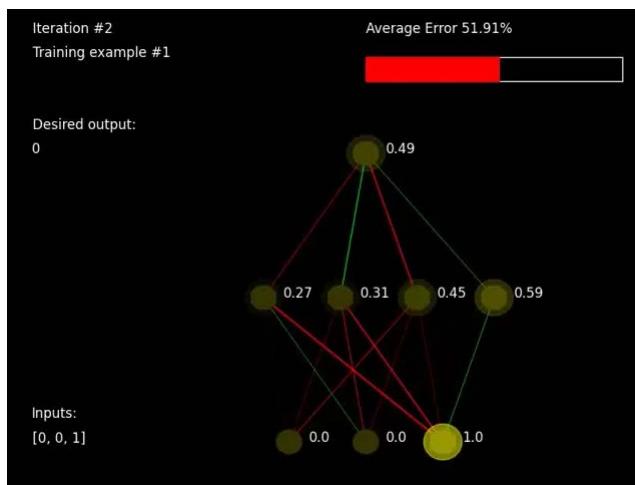
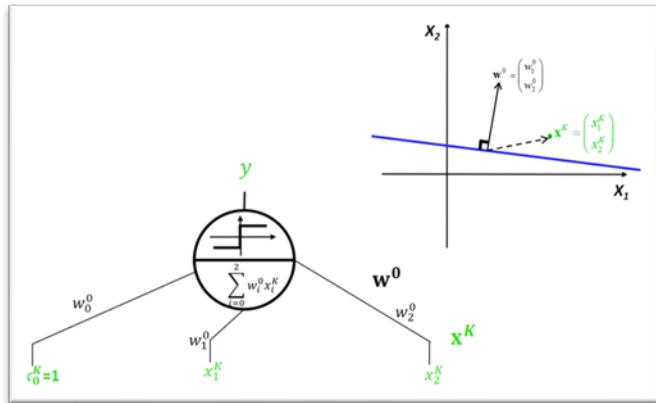


Illustrations from: <http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/>

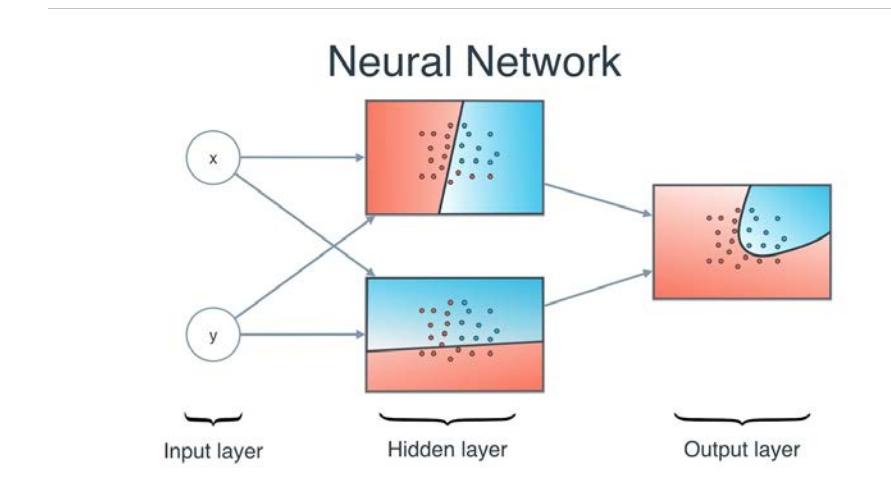
Multi-Layer Perceptron, manifold disentanglement



Illustrations from: <http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/>

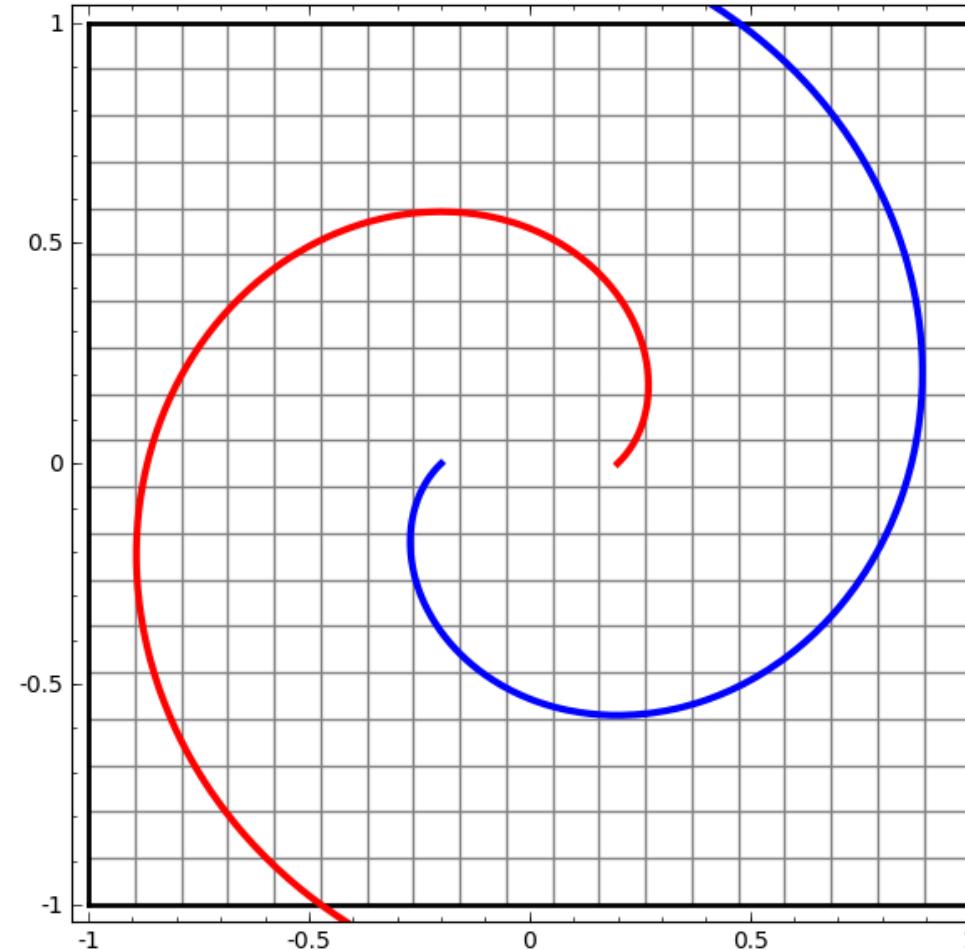


From perceptron to network



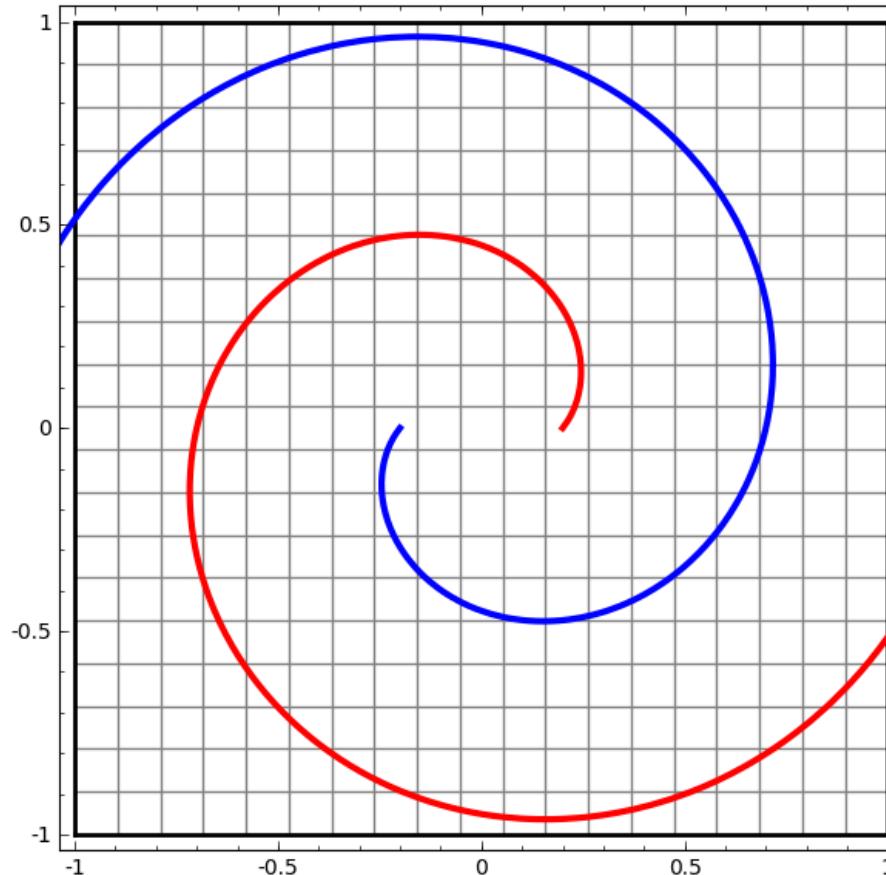
@tachyeonz: A friendly introduction to neural networks and deep learning.

Multi-Layer Perceptron, manifold disentanglement



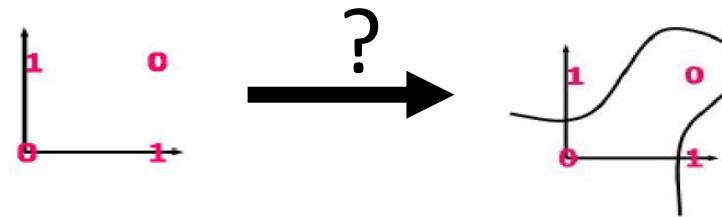


Multi-Layer Perceptron, manifold disentanglement

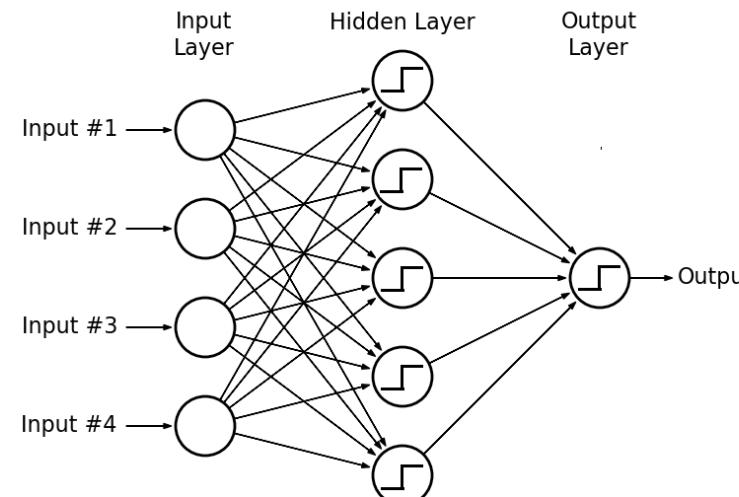


Single Perceptron Unit

- **Perceptron** only learns linear function [Minsky and Papert, 1969]



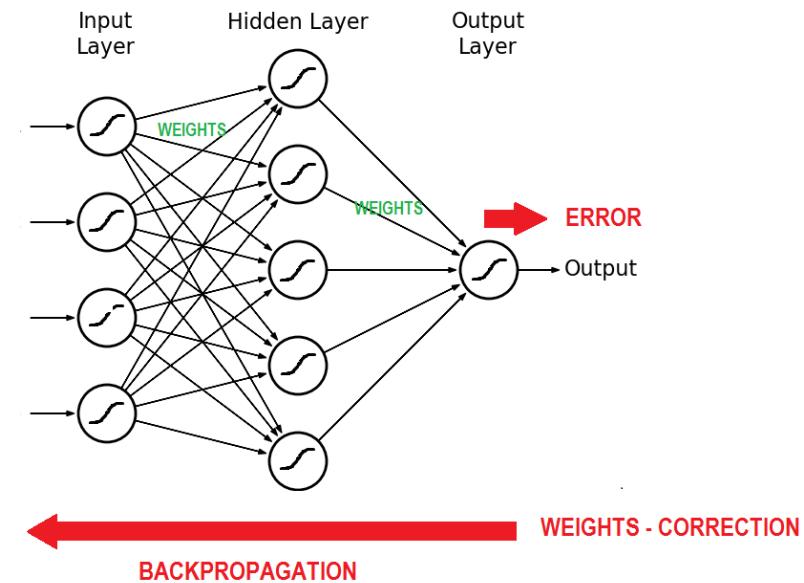
- Non-linear function needs layer(s) of neurons → Neural Network
- Neural Network = input layer + hidden layer(s) + output layer





Multi-Layer Perceptron

- Training a neural network [Rumelhart et al. / Yann Le Cun et al. 1985]
- **Unknown parameters: weights on the synapses**
- Minimizing a cost function: some metric between the predicted output and the given output

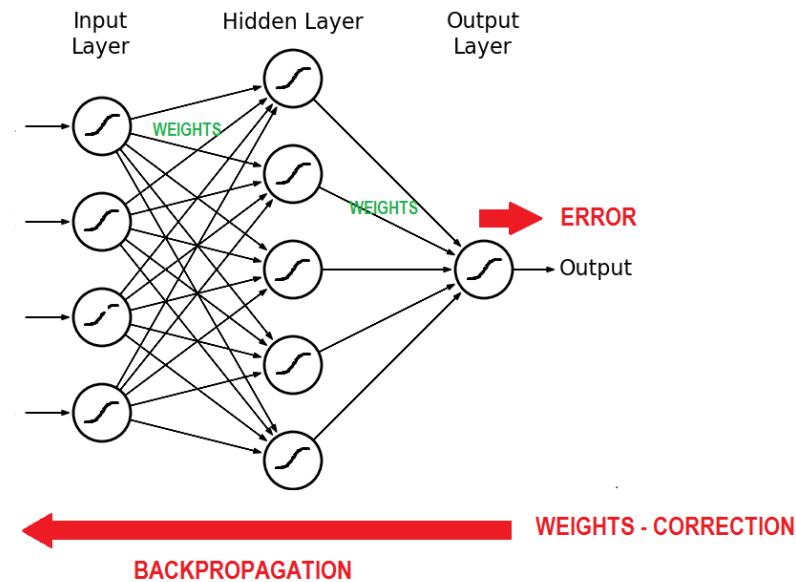


- Step function: non-continuous functions are replaced by a continuous non-linear ones



Multi-Layer Perceptron

- Minimizing a cost function: some metric between the predicted output and the given output



- Equation for a network of 3 neurons (i.e. 3 perceptrons):

$$y = s(w_{13}s(w_{11}x_1 + w_{21}x_2 + w_{01}) + w_{23}s(w_{12}x_1 + w_{22}x_2 + w_{02}) + w_{03})$$

Multi-Layer Perceptron

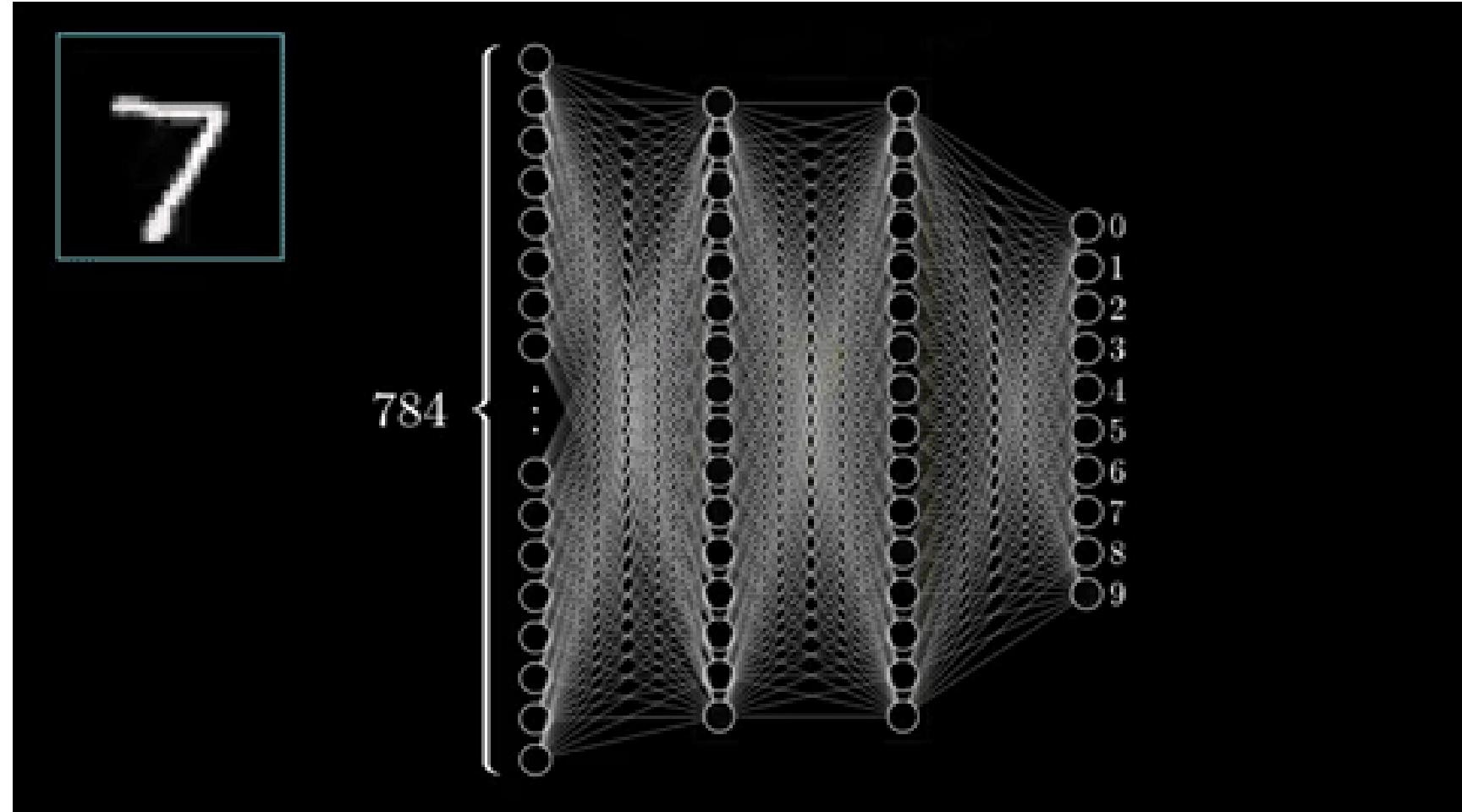
Theorem (Cybenko (1989), Hornik & Stinchcombe & White (1989))
A neural network with one single hidden layer is a **universal approximator**: it can represent any continuous function on compact subsets of \mathbb{R}^n

- 2 layers is enough ... theoretically:

“...networks with one internal layer and an arbitrary continuous sigmoidal function can approximate continuous functions with arbitrary precision providing that no constraints are placed on the number of nodes or the size of the weights”

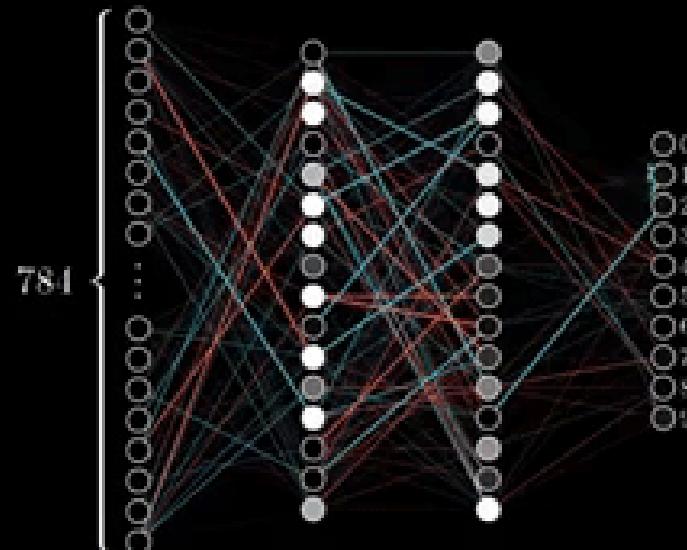
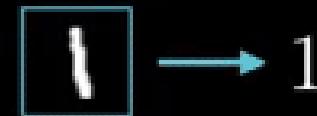
- **But *no efficient learning rule* is known and the size of the hidden layer is *exponential* with the complexity of the problem which is unknown beforehand.**

Solution: Neural Network

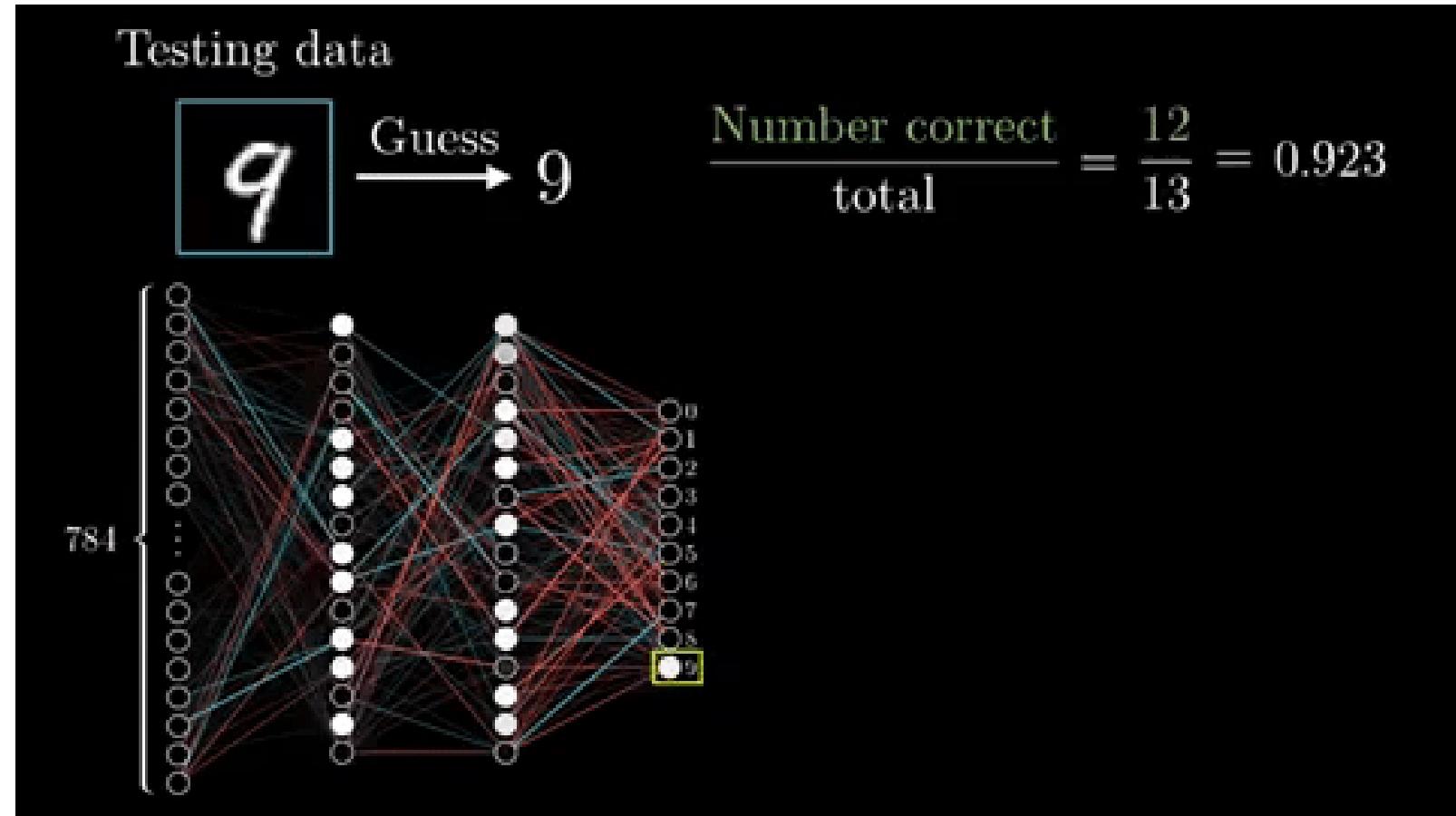


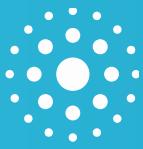
Training

Training in
progress...



Inference/Test





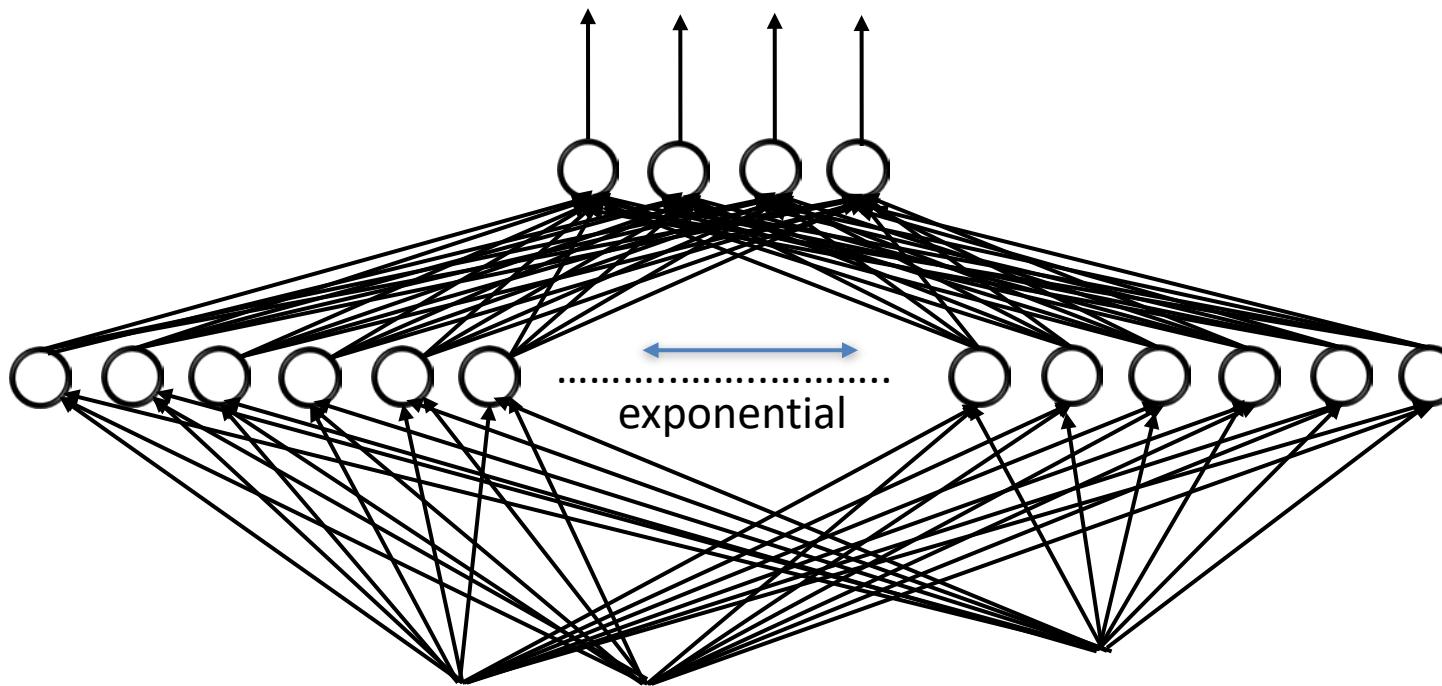
DEEP LEARNING PRINCIPLES



Deep representation origins

univ-cotedazur.fr

- **Theorems** (Cybenko (1989), Hornik & Stinchcombe & White (1989), Allan Pinkus (1999))
A neural network with one single hidden layer is a “universal approximator”, it can represent any continuous function on compact subsets of $\mathbb{R}^n \Rightarrow 1$ hidden layer is enough...but the hidden layer size may be exponential for error ε (or even infinite for error 0), and there is no efficient learning rule known.



Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4), 303-314.

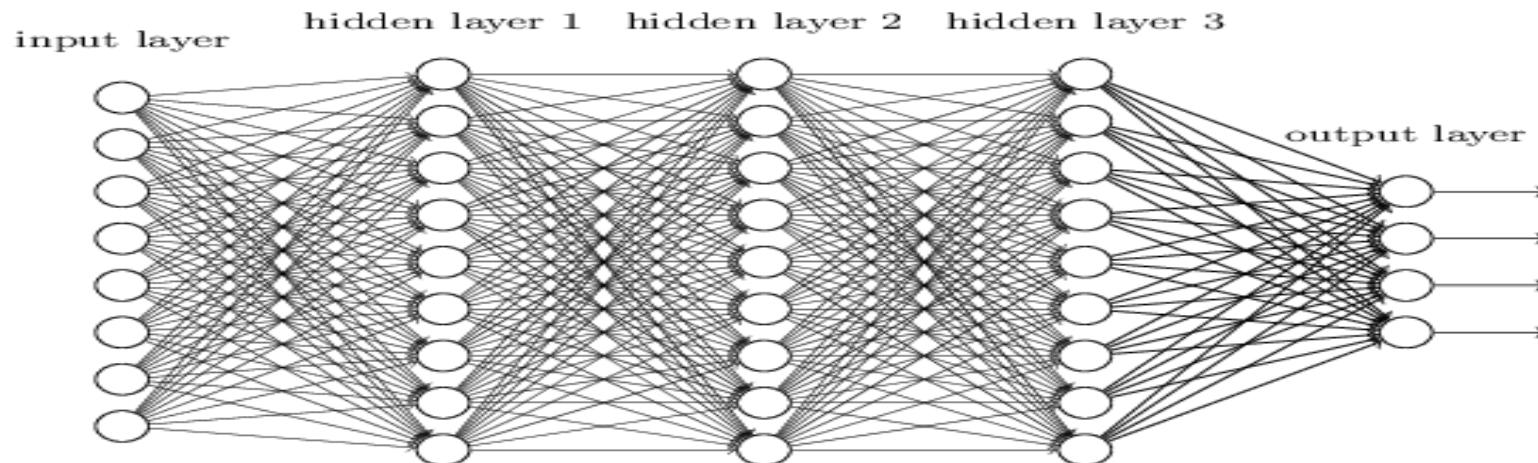
Hornik, K., Stinchcombe, M., & White, H. (1989). Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5), 359-366.

Pinkus, Allan (1999). "Approximation theory of the MLP model in neural networks". *Acta Numerica*. 8: 143–195.



Deep representation origins

- **Theorem Hastad (1986), Bengio et al. (2007)** Functions representable compactly with k layers may require exponentially size with $k-1$ layers



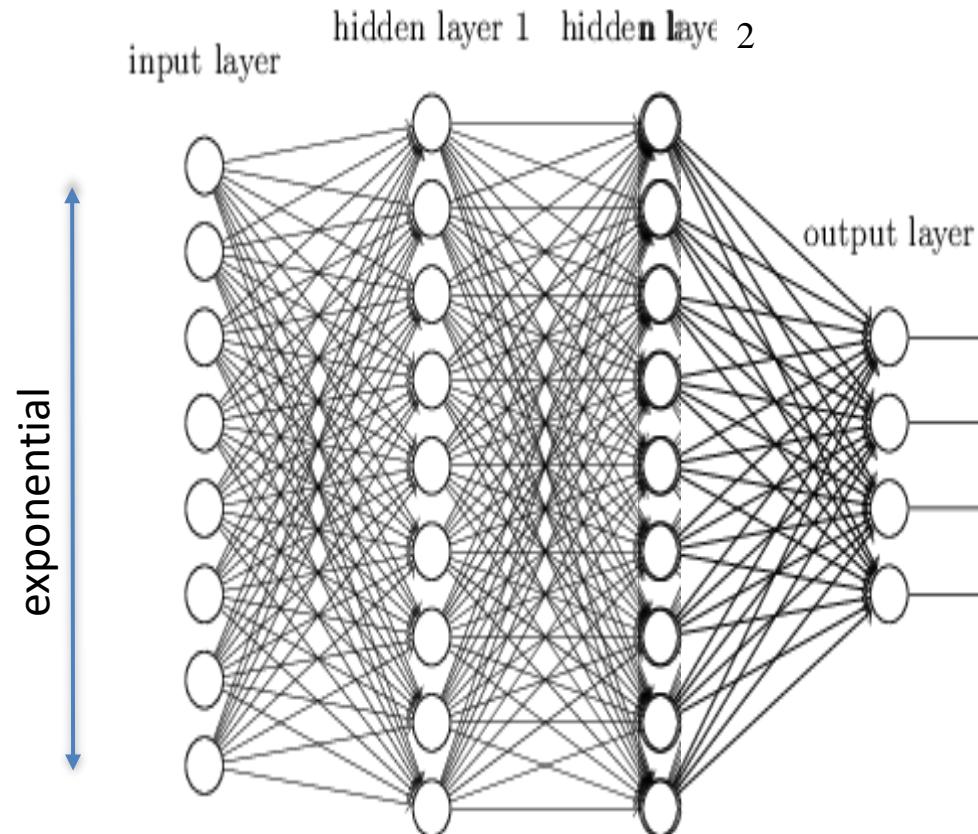
Johan T. Håstad. Computational Limitations for Small Depth Circuits. MIT Press, Cambridge, MA, 1987.

Bengio, Y., & LeCun, Y. (2007). Scaling learning algorithms towards AI. *Large-scale kernel machines*, 34(5), 1-41.



Deep representation origins

- **Theorem Hastad (1986), Bengio et al. (2007)** Functions representable compactly with k layers may require exponentially size with $k-1$ layers





Deep Neural Network Approximation Theory (2021)

“...deep networks provide exponential approximation accuracy—i.e., the approximation error decays exponentially in the number of nonzero weights in the network—of the multiplication operation, polynomials, sinusoidal functions, and certain smooth functions.

Moreover, this holds true even for one-dimensional oscillatory textures and the Weierstrass function—a fractal function, neither of which has previously known methods achieving exponential approximation accuracy.

We also show that in the approximation of sufficiently smooth functions finite-width deep networks require strictly smaller connectivity than finite-depth wide networks.”

Elbrächter, D., Perekrestenko, D., Grohs, P., & Bölcskei, H. (2021). Deep neural network approximation theory. IEEE Transactions on Information Theory, 67(5), 2581-2623.

The Blessing of dimensionality: Thomas Cover's Theorem (1965)

Cover's theorem states: A complex pattern-classification problem cast in a high-dimensional space nonlinearly is more likely to be linearly separable than in a low-dimensional space (repeated sequence of Bernoulli trials).

The Blessing of dimensionality: Thomas Cover's Theorem (1965)

Cover's theorem states: A complex pattern-classification problem cast in a high-dimensional space nonlinearly is more likely to be linearly separable than in a low-dimensional space.
(repeated sequence of Bernoulli trials)

The number of groupings that can be formed by $(l-1)$ -dimensional hyperplanes to separate N points in two classes is:

$$O(N, l) = 2 \sum_{i=0}^l \frac{(N-1)!}{(N-1-i)!i!}$$

Notice: The total number of possible groupings is 2^N

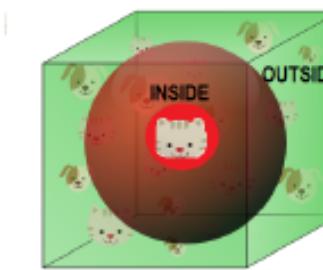


The curse of dimensionality [Bellman, 1956]

univ-cotedazur.fr

- Euclidian distance is not relevant in high dimension: $d \geq 10$
 - ➊ look at the examples at distance at most r
 - ➋ the hypersphere volume is too small: practically empty of examples

$$\frac{\text{volume of the sphere of radial } r}{\text{hypersphere of } 2r \text{ width}} \xrightarrow{d \rightarrow \infty} 0$$



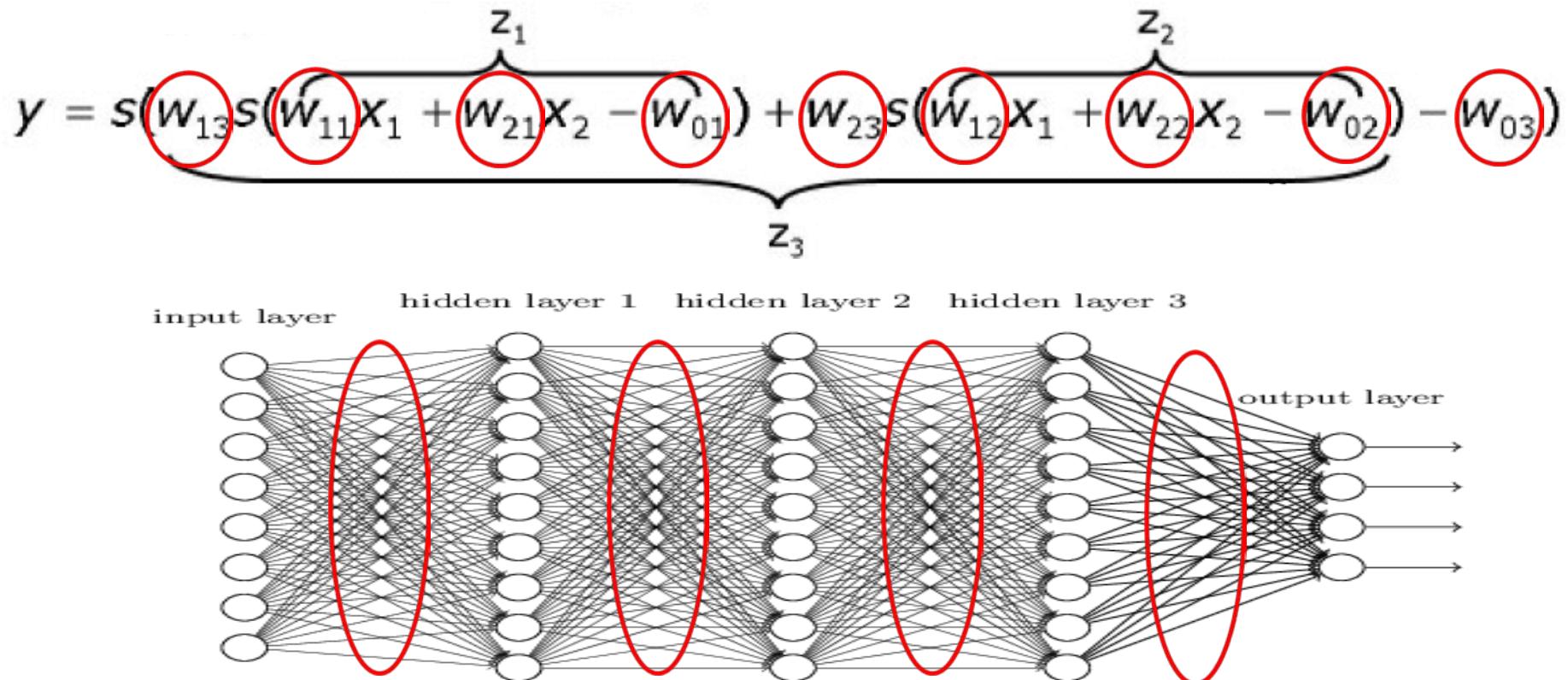
- ➌ need a number of examples exponential in d

Remark

Specific care for data representation

Structure the network?

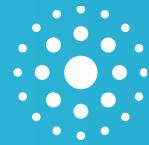
- Can we put any structure reducing the space of exploration and providing useful properties (invariance, robustness, sequentiality...)?





Enabling factors

- Why do it now ? Before 2006, training deep networks was unsuccessful because of practical aspects
 - faster CPU's
 - parallel CPU architectures
 - advent of GPU computing
 - Advances in ML/Optim (1995 -> 2005)
- Hinton, Osindero & Teh « A Fast Learning Algorithm for Deep Belief Nets », Neural Computation, 2006
- Bengio, Lamblin, Popovici, Larochelle « Greedy Layer-Wise Training of Deep Networks », NIPS'2006
- Ranzato, Poultney, Chopra, LeCun « Efficient Learning of Sparse Representations with an Energy-Based Model », NIPS'2006
- Results...
 - 2009, sound, interspeech +~24%
 - 2011, text, +~15% without linguistic at all
 - 2012, images, ImageNet +~20%
 - 2020, molecules/graphs, AlphaFold, +~24%
(sorry in French, <https://www.youtube.com/watch?v=OGewxRMME8o>)

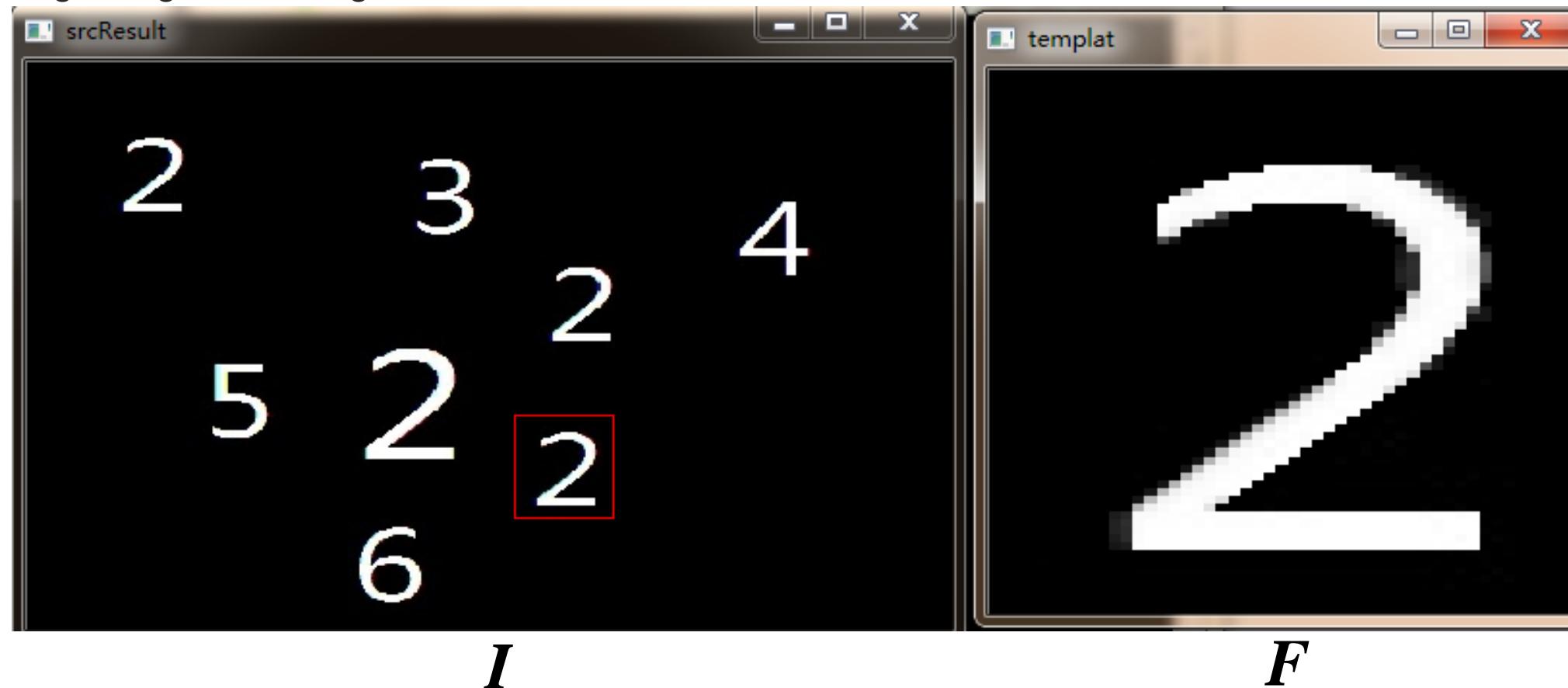


Invariance to spatial changes

**CONVOLUTIONAL NEURAL NETWORKS
(AKA CNN, ConvNET)**

Do you know template matching?

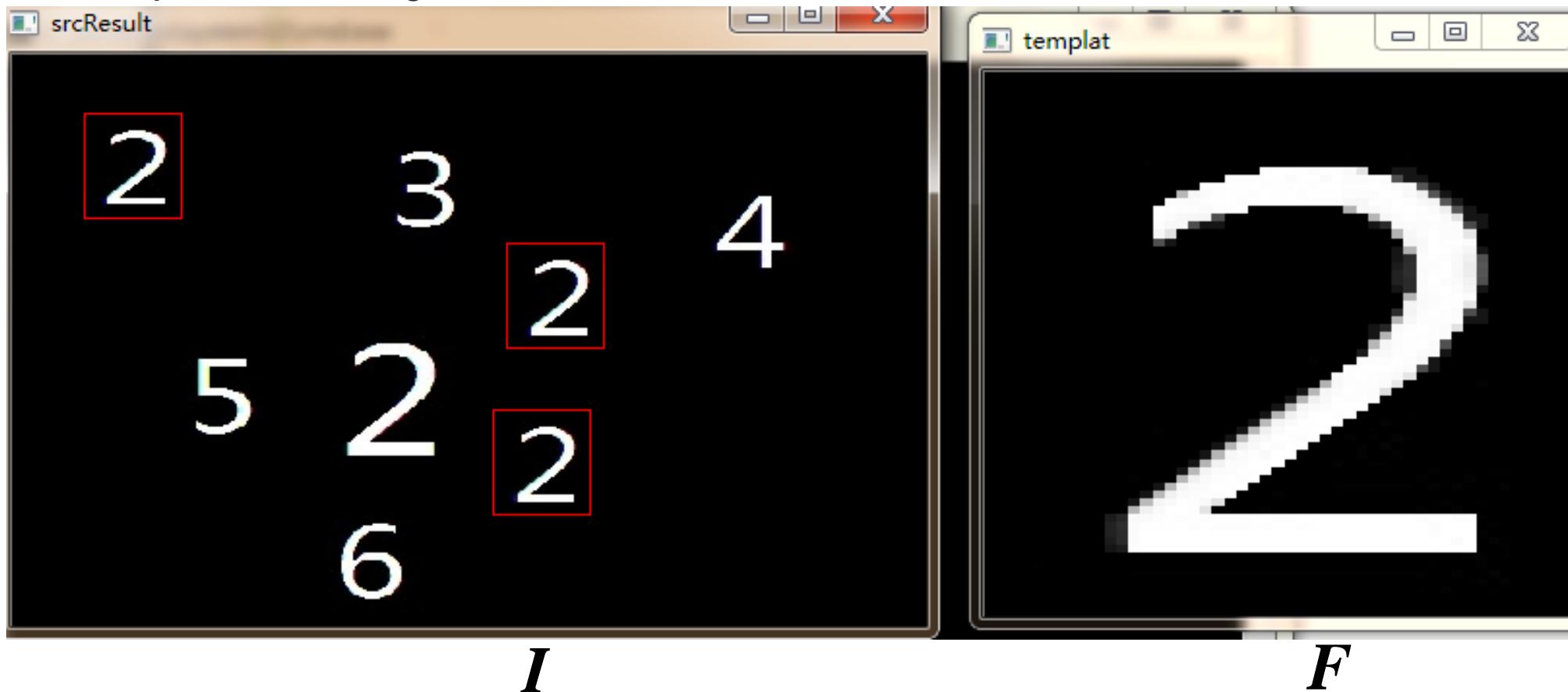
Single target matching results:



(picture from <https://programmer.group/5ca5086fed10b.html>)

Do you know template matching?

Multi-objective matching result:



(picture from <https://programmer.group/5ca5086fed10b.html>)

Do you know template matching?

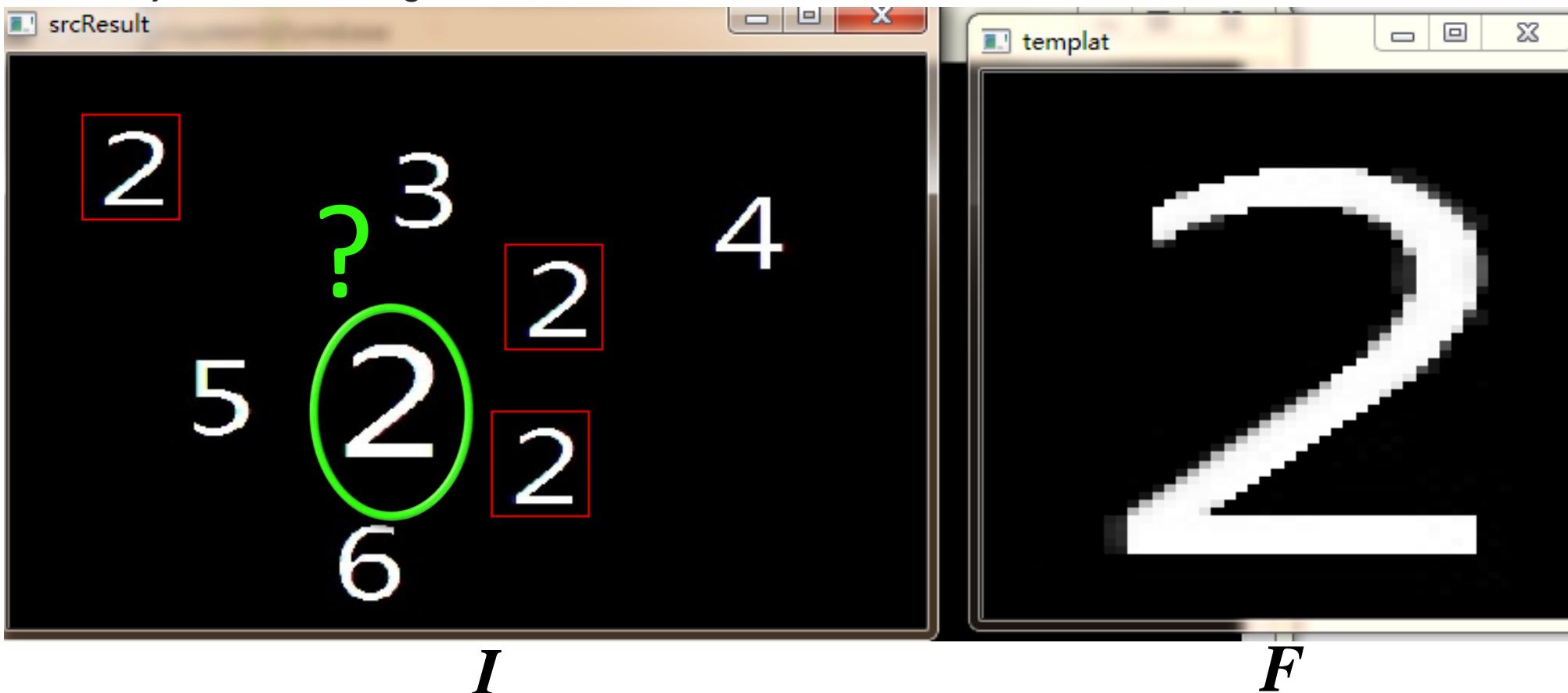
How to compute Template matching? Cross-correlation!

$$F \circ I(x, y) = \sum_{j=-N}^N \sum_{i=-N}^N F(i, j)I(x+i, y+j)$$

When this is maximal, the Filter is perfectly positioned with respect to the pattern to detect.

Do you know template matching?

Multi-objective matching result:



Are there more universal templates?

Multi-objective matching result:



(picture from <https://programmer.group/5ca5086fed10b.html>)

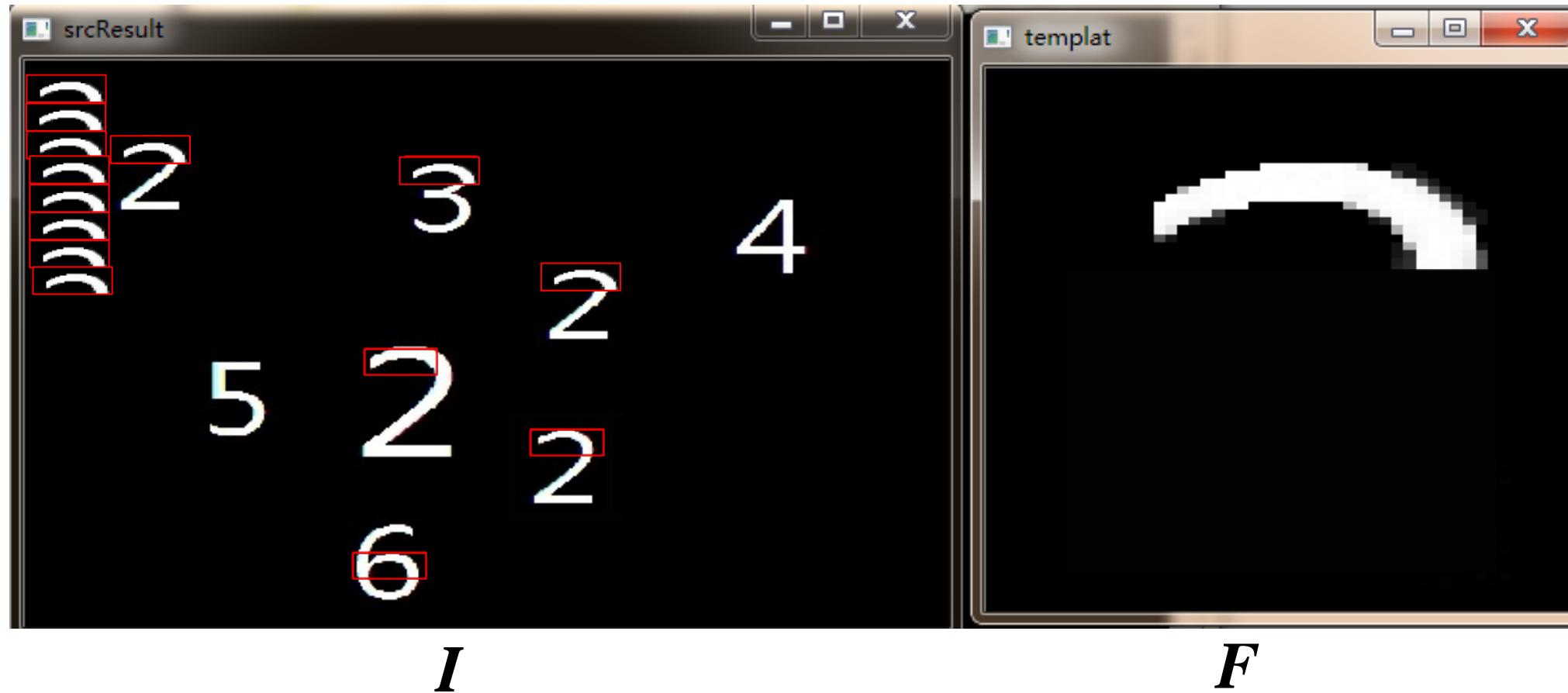
Are there more universal templates?

Multi-objective matching result:



How to locate these templates?

By scanning I , we acquire **translation invariance**



How to define these templates?

By learning what filter maximizes local cross-correlations



Cross-correlation or Convolution?

- How to define these templates?

By learning what filter maximizes **local cross-correlations**

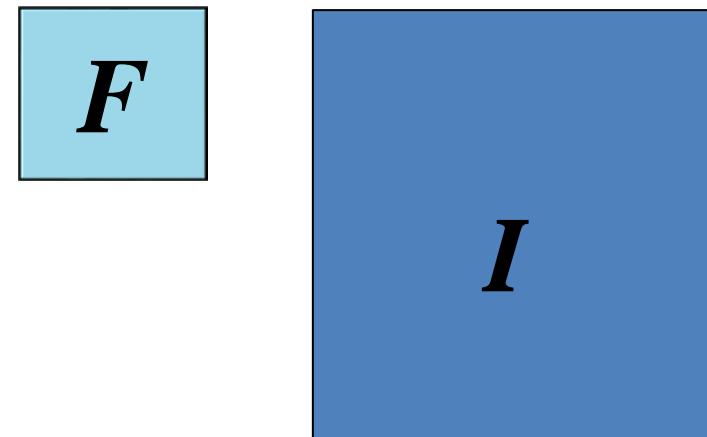
$$F \circ I(x, y) = \sum_{j=-N}^N \sum_{i=-N}^N F(i, j)I(x+i, y+j)$$

Or **local convolutions**

$$F * I(x, y) = \sum_{j=-N}^N \sum_{i=-N}^N F(i, j)I(x-i, y-j)$$

Cross-correlation or Convolution?

- Convolution:
 - Flip the filter in both dimensions (bottom to top, right to left)
 - Then apply cross-correlation

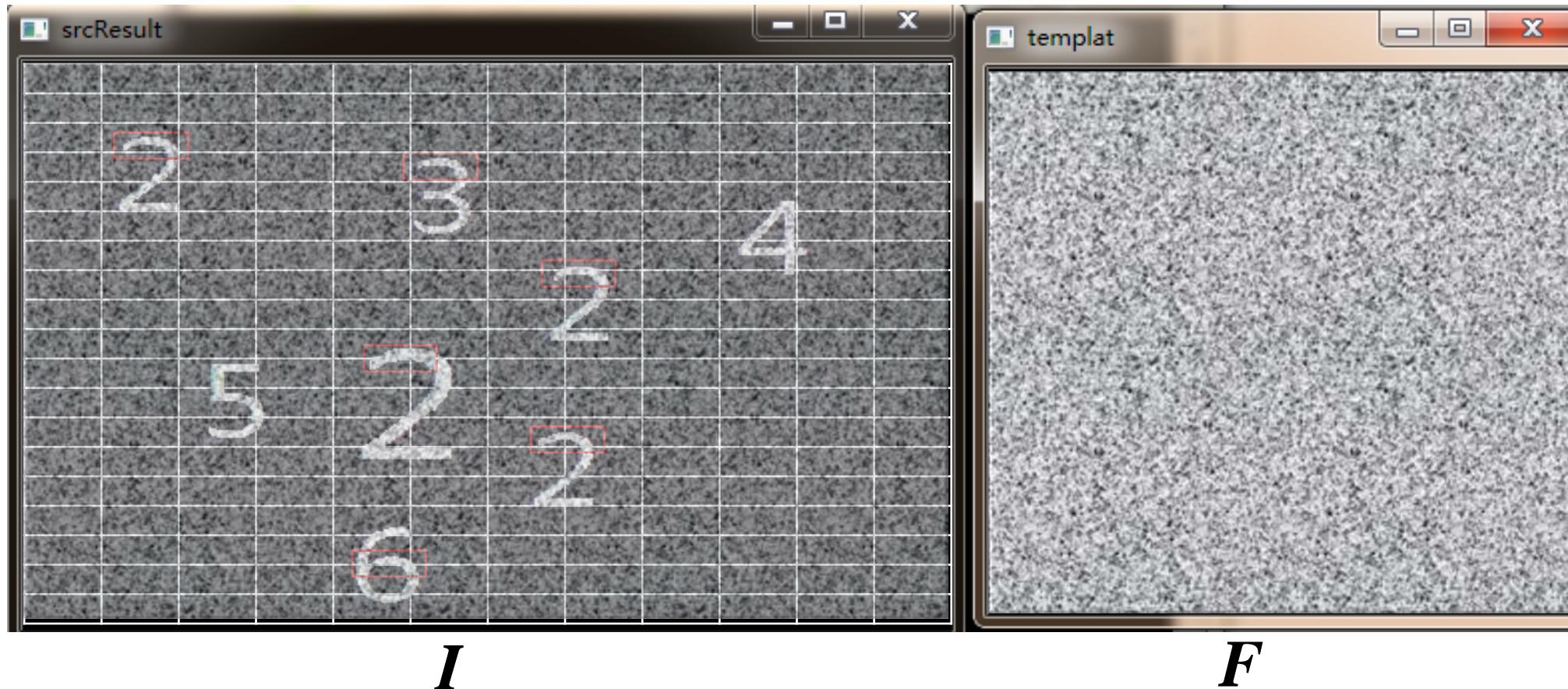


Cross-correlation or Convolution?

- Convolution:
 - Flip the filter in both dimensions (bottom to top, right to left)
 - Then apply cross-correlation
- Actually, The Convolutional operation widely used in Convolution Neural Network (CNN) is commonly and largely considered as a **misnomer**
 - The operation that is used is strictly speaking a correlation instead of convolution.
 - Both the operators have **translational invariance and locality**.

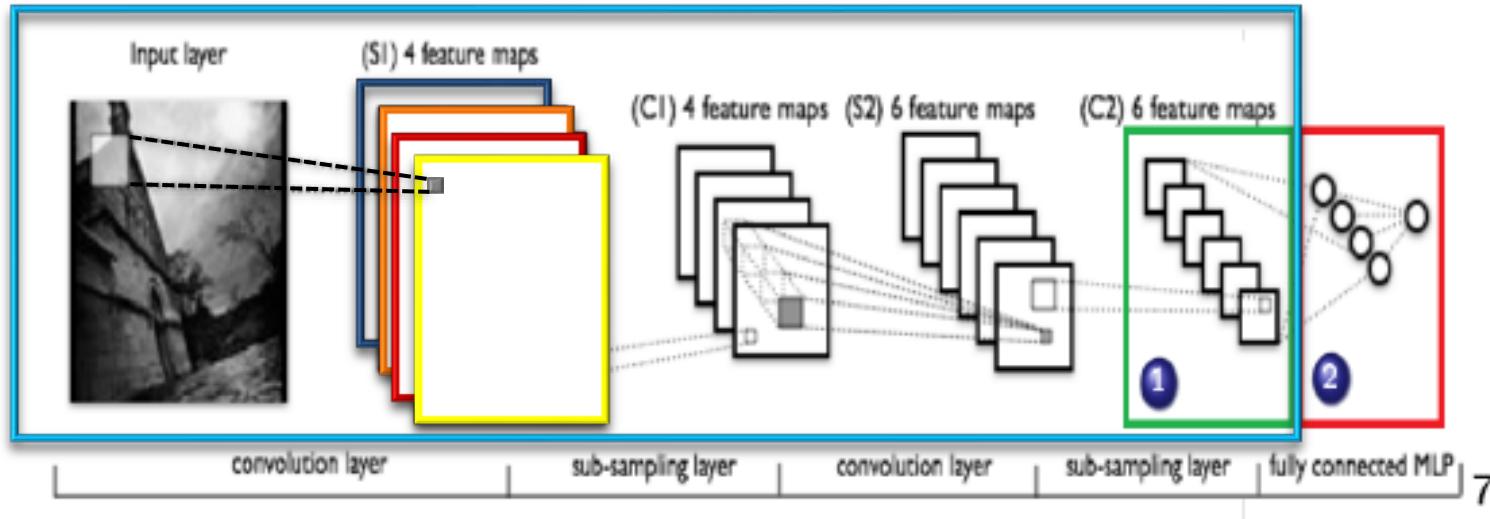
How to define these templates?

And instead of scanning I , let us just **duplicate** the template **into a grid**.



Deep representation by CNN

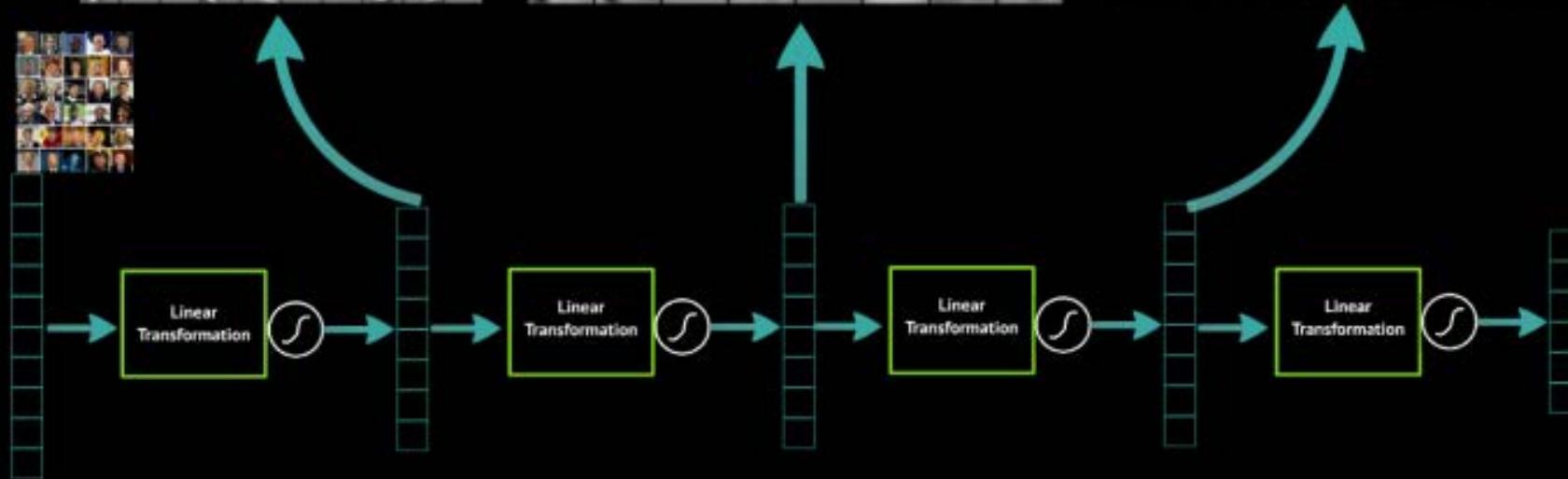
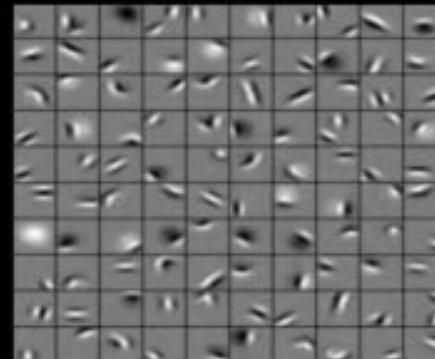
- feature map = result of the convolution
- convolution with a filter extract characteristics (*edge detectors*)
- extract parallelised characteristics at each layer



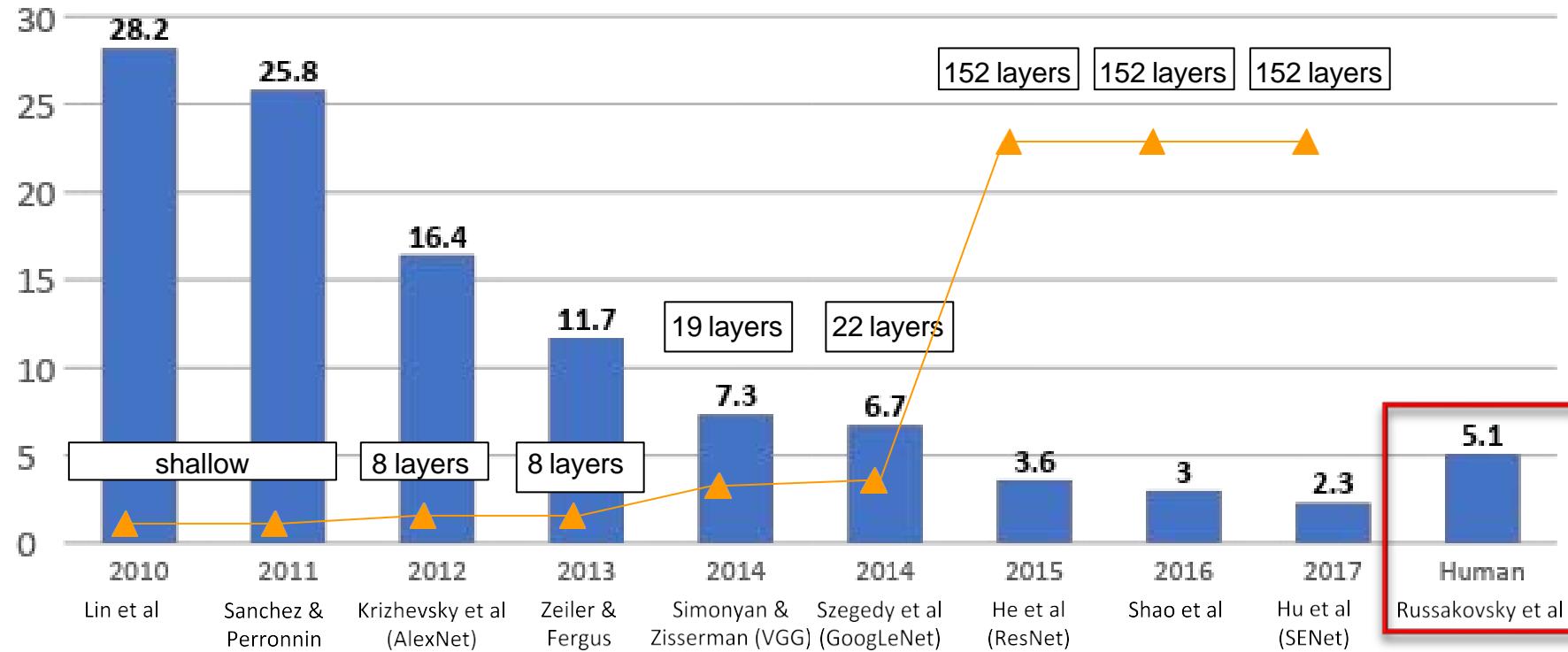
- ➊ final representation of our data
- ➋ classifier (MLP)

Deep representation by CNN

Deep Learning learns layers of features



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Invariance to sequential changes

RECURRENT NEURAL NETWORKS (AKA RNN, LSTM / GRU)

Borrowed from COLAH's Blog, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

From "Introduction to Deep Learning", by Professor Qiang Yang, from The Hong Kong University of Science and Technology

And from lecture on RNN of Prof Adriana Kovashka

Some pre-RNN captioning results



This is a picture of one sky, one road and one sheep. The gray sky is over the gray road. The gray sheep is by the gray road.



Here we see one road, one sky and one bicycle. The road is near the blue sky, and near the colorful bicycle. The colorful bicycle is within the blue sky.



This is a picture of two dogs. The first dog is near the second furry dog.

Results with Recurrent Neural Networks



"man in black shirt is playing guitar."



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

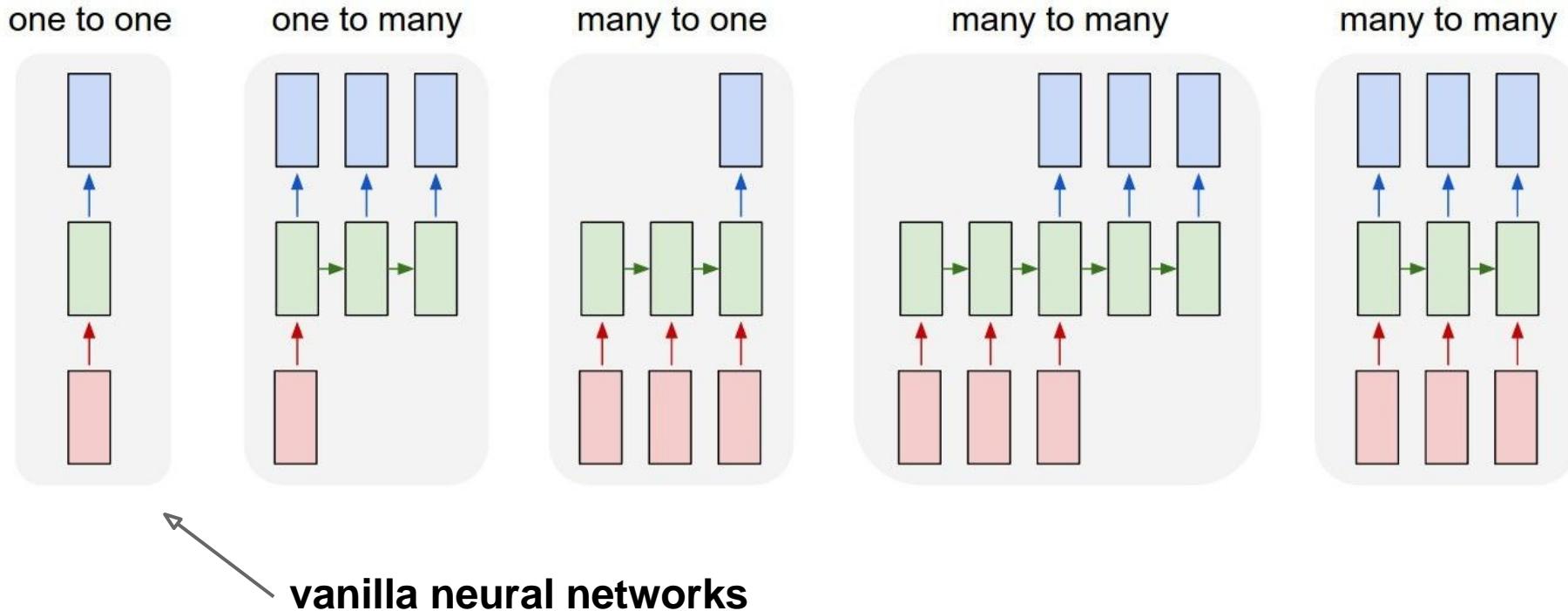


"two young girls are playing with lego toy."

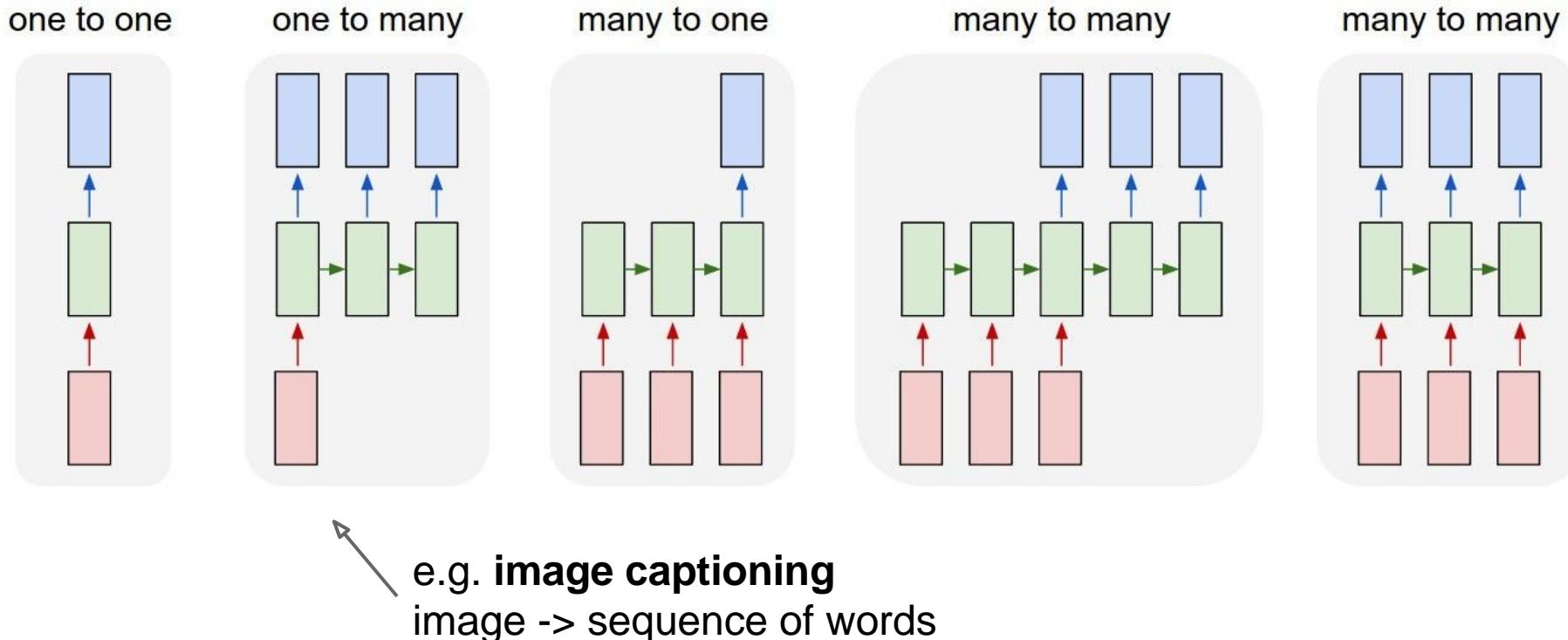


"boy is doing backflip on wakeboard."

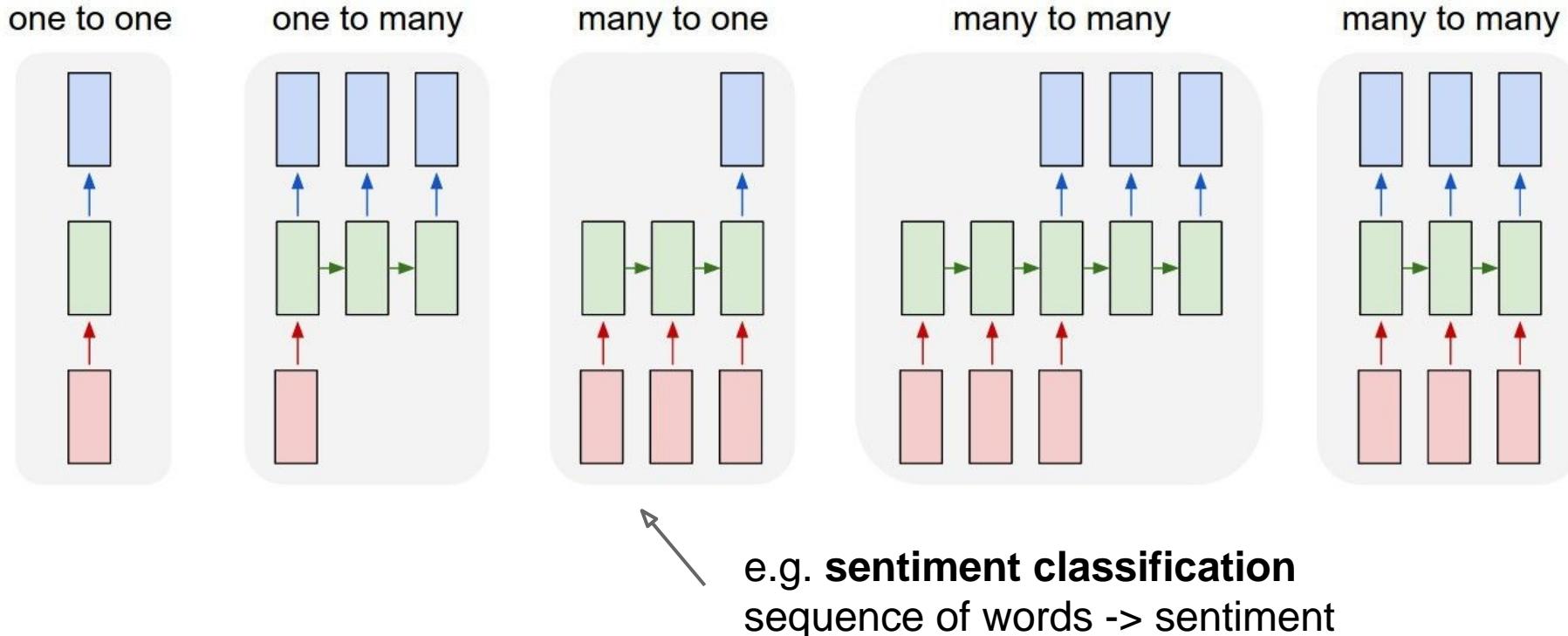
Recurrent Networks offer a lot of flexibility:



Recurrent Networks offer a lot of flexibility:

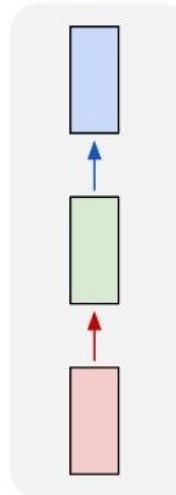


Recurrent Networks offer a lot of flexibility:

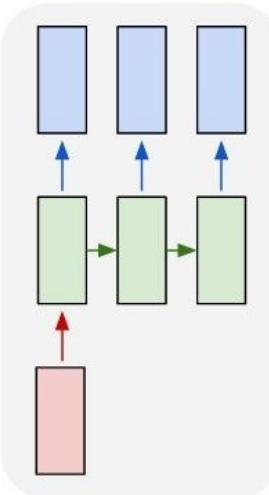


Recurrent Networks offer a lot of flexibility:

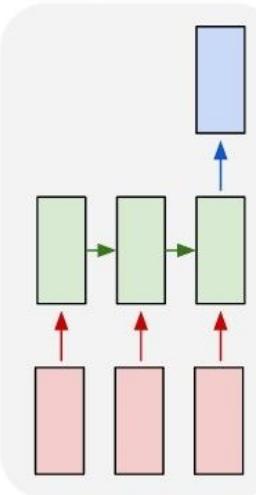
one to one



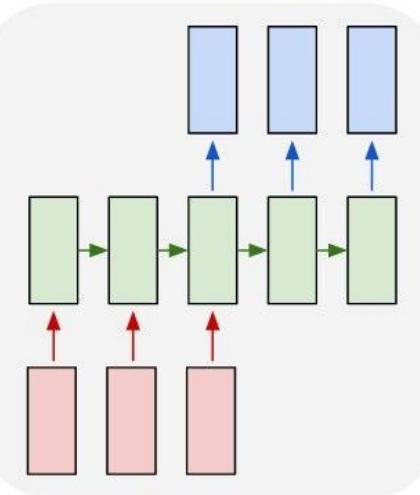
one to many



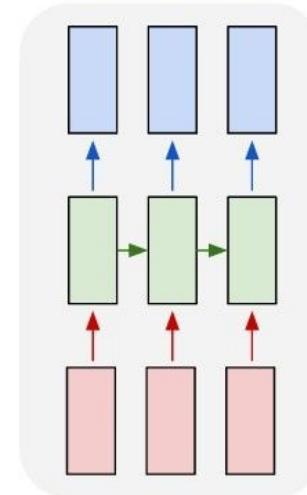
many to one



many to many

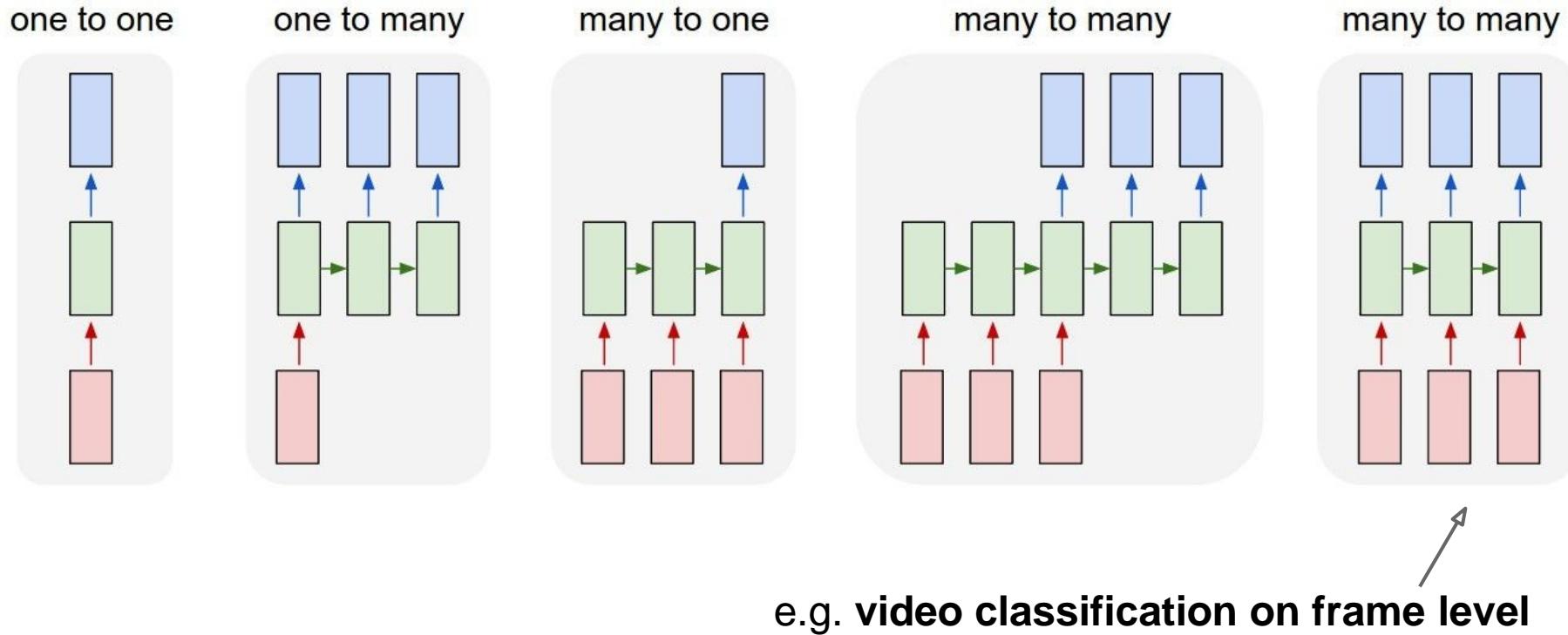


many to many



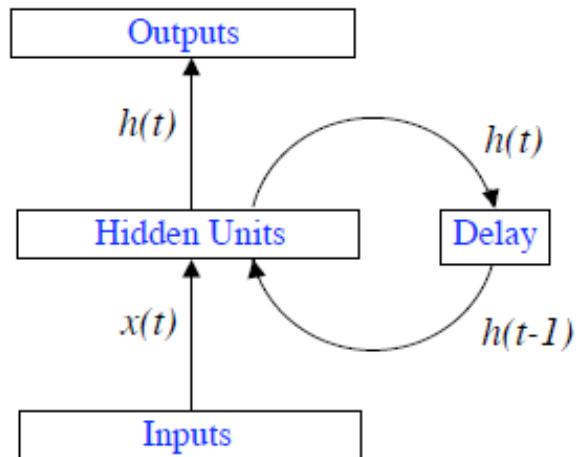
e.g. **machine translation**
seq of words -> seq of words

Recurrent Networks offer a lot of flexibility:



What are RNNs?

Recurrent neural networks (RNNs) are connectionist models with the ability to selectively pass information across sequence steps, while processing sequential data one element at a time.



The simplest form of **fully recurrent neural network** is an MLP with the previous set of hidden unit activations feeding back into the network along with the inputs

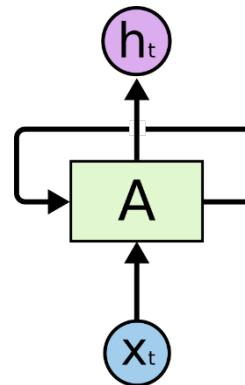
Allow a 'memory' of previous inputs to persist in the network's internal state, and thereby influence the network output

$$h(t) = f_H(W_{IH}x(t) + W_{HH}h(t - 1))$$

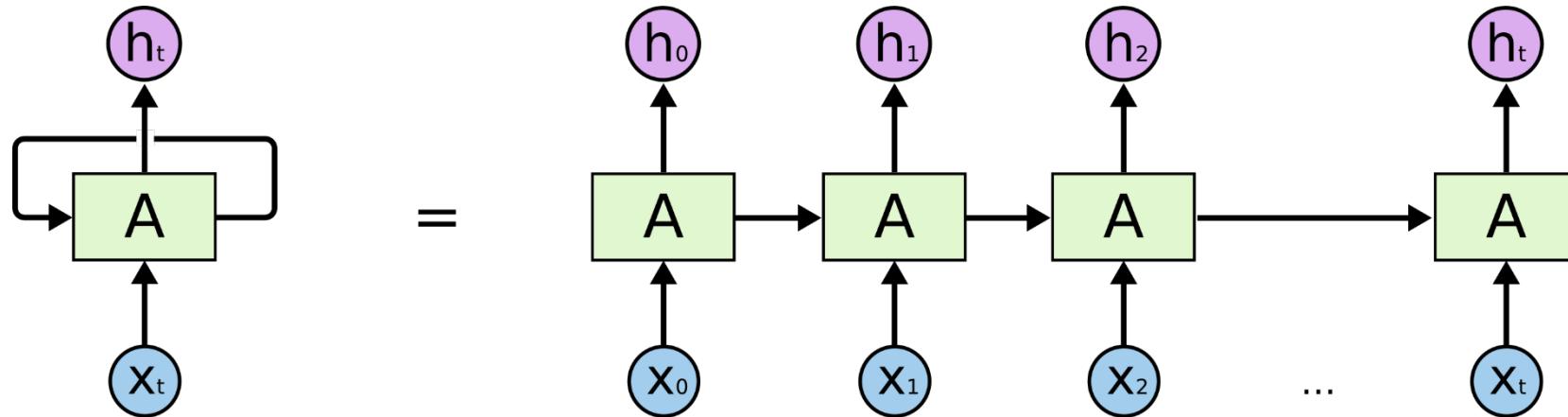
$$y(t) = f_O(W_{HO}h(t))$$

f_H and f_O are the activation function for hidden and output unit; W_{IH} , W_{HH} , and W_{HO} are connection weight matrices which are learnt by training

Recurrent Neural Networks have loops.



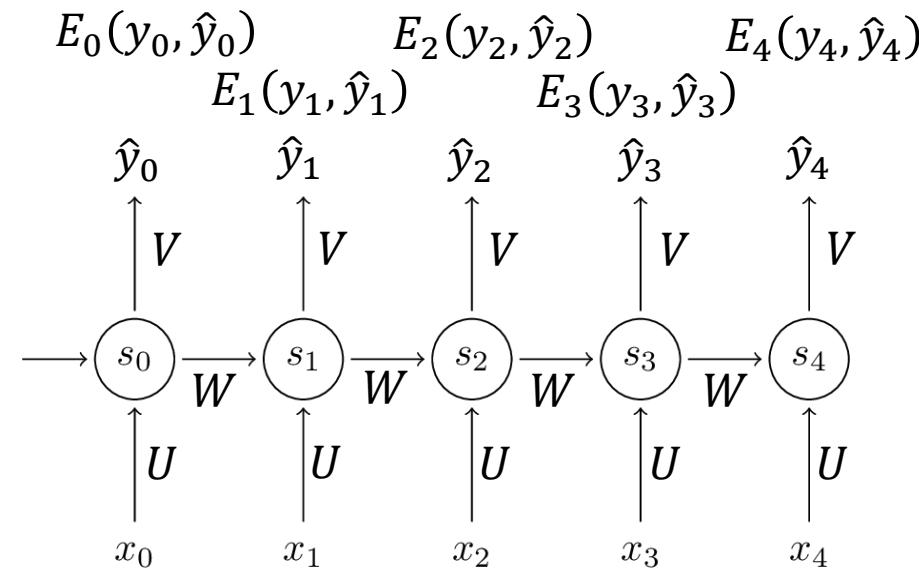
An unrolled recurrent neural network.



In the last few years, there have been incredible success applying RNNs to a variety of problems: speech recognition, language modeling, translation, image captioning...

What are RNNs?

- The recurrent network can be converted into a feed-forward network by ***unfolding over time***



An unfolded recurrent network. Each node represents a layer of network units at a single time step. The weighted connections from the input layer to hidden layer are labelled ' W_{IH} ', those from the hidden layer to itself (i.e. the recurrent weights) are labelled ' W_{HH} ' and the hidden to output weights are labelled ' W_{HO} '. **Note that the same weights are reused at every time step. Bias weights are omitted for clarity.**

What are RNNs?

- Training RNNs (determine the parameters)

Back Propagation Through Time (BPTT) is often used to learn the RNN
BPTT is an extension of the back-propagation (BP)

- The output of this RNN is \hat{y}_t

$$s_t = \tanh(Ux_t + Ws_{t-1})$$
$$\hat{y}_t = \text{softmax}(Vs_t)$$

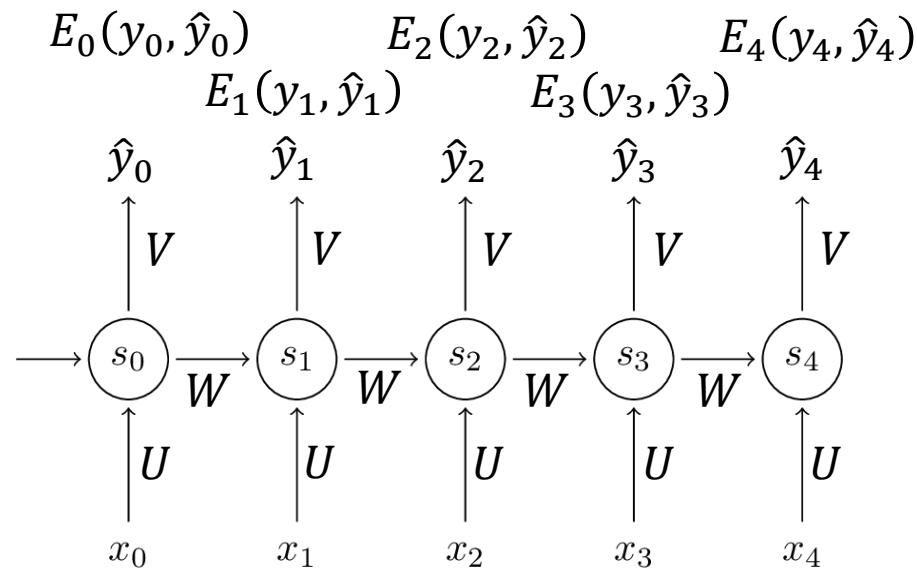
- The loss/error function of this network is

$$E_t(y_t, \hat{y}_t) = -y_t \log \hat{y}_t$$

The error at each time step

$$E(y, \hat{y}) = \sum_t E_t(y_t, \hat{y}_t)$$

the total loss is the sum of the errors at each time step



What are RNNs?

- Training RNNs (determine the parameters)
 - ✓ The gradients of the error with respect to our parameters
Just like we sum up the errors, we also sum up the gradients at each time step for one training example. For parameter W , the gradient is

$$\frac{\partial E}{\partial W} = \sum_t \frac{\partial E_t}{\partial W}$$

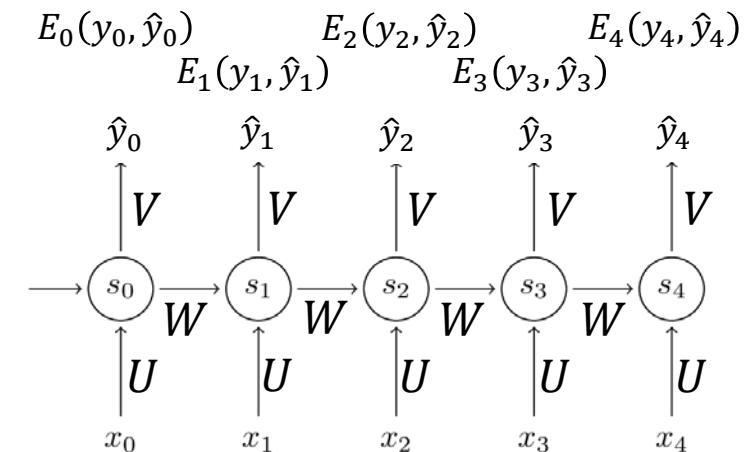
- ✓ The gradient at each time step

we use time 3 as an example

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial W} \quad \xrightarrow{\text{Chain Rule}}$$

$s_3 = \tanh(Ux_3 + Ws_2)$ s₃ depends on W and s_2 , we cannot simply consider s_2 a constant

$$\frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial W} \quad \xrightarrow{\text{Apply Chain Rule again on } s_k}$$

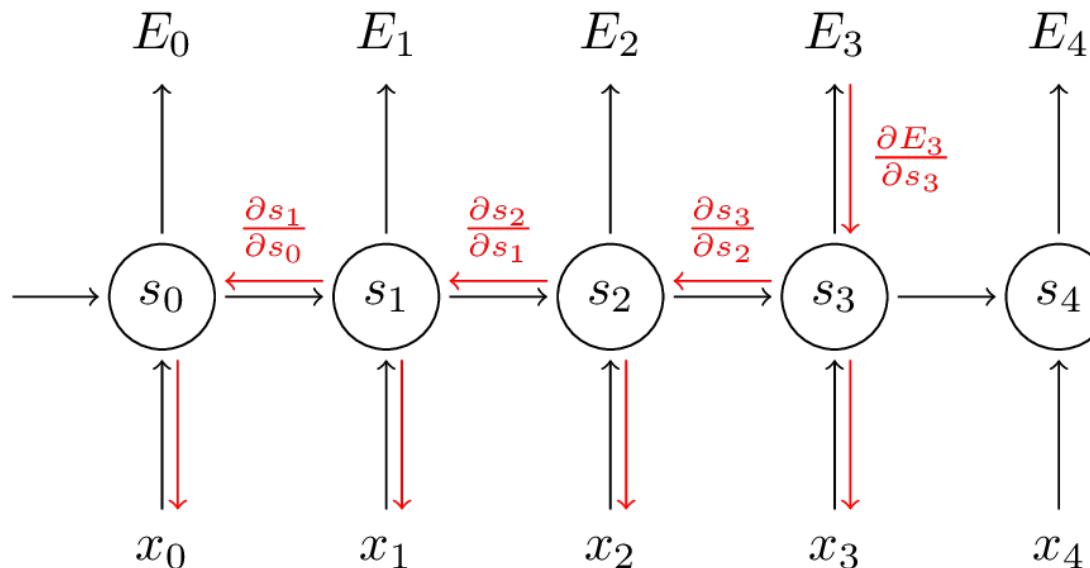


What are RNNs?

- Training RNNs (determine the parameters)

$$\frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial W} \quad \longrightarrow \quad \frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \prod_{j=k+1}^3 \frac{\partial s_j}{\partial s_{j-1}} \frac{\partial s_k}{\partial W}$$

Because W is used in every step up to the output we care about, we need to backpropagate gradients from $t = 3$ through the network all the way to $t = 0$



What are RNNs?

- The vanishing gradient problem

To understand why, let's take a closer look at the gradient we calculated above:

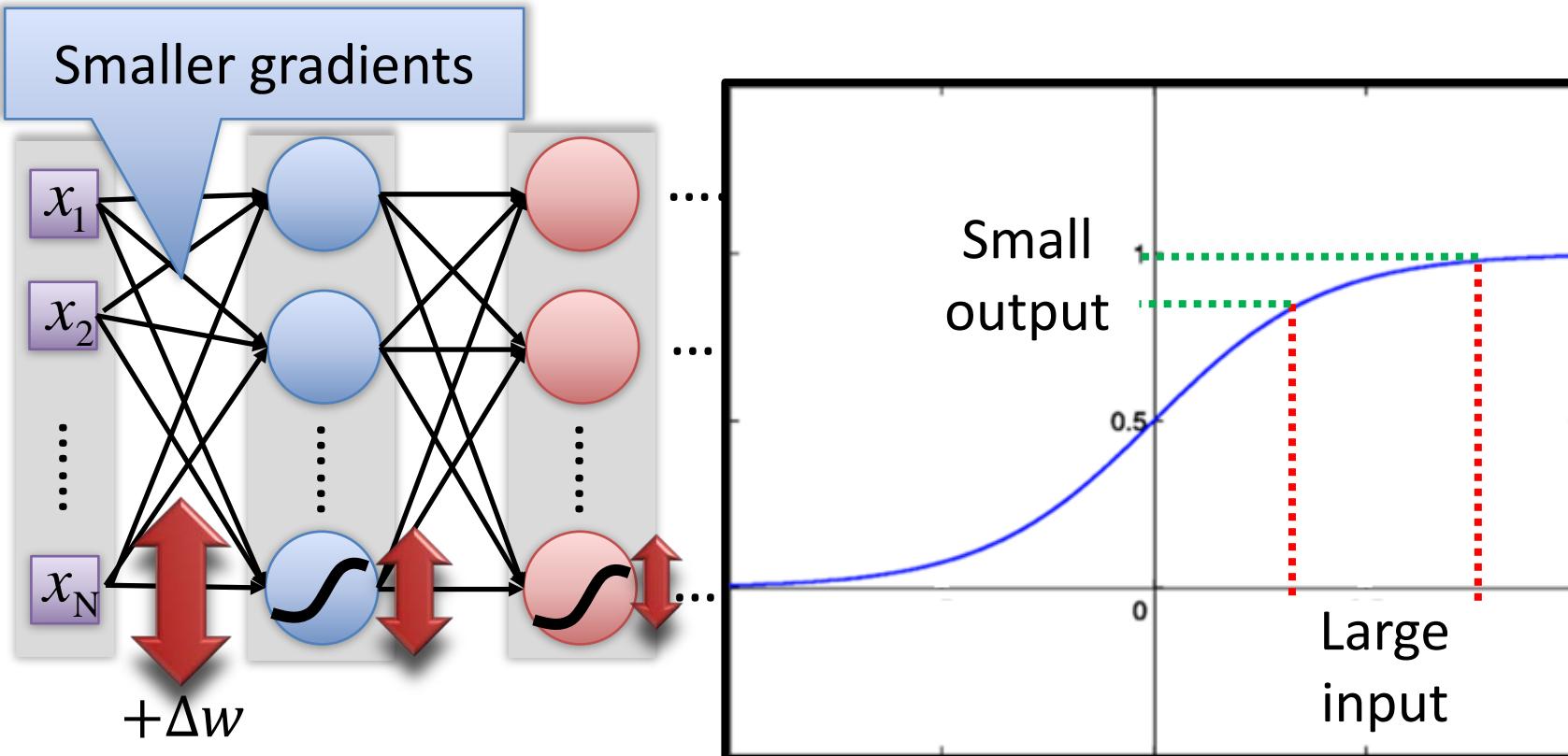
$$\frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \underbrace{\prod_{j=k+1}^3 \frac{\partial s_j}{\partial s_{j-1}}}_{\text{ }} \frac{\partial s_k}{\partial W}$$

Because the layers and time steps of deep neural networks relate to each other through multiplication, derivatives are susceptible to vanishing

Gradient contributions from “far away” steps become zero, and the state at those steps **does not contribute** to what you are learning: You end up not learning long-range dependencies.

Recipe of Deep Learning

- *Vanishing Gradient Problem*



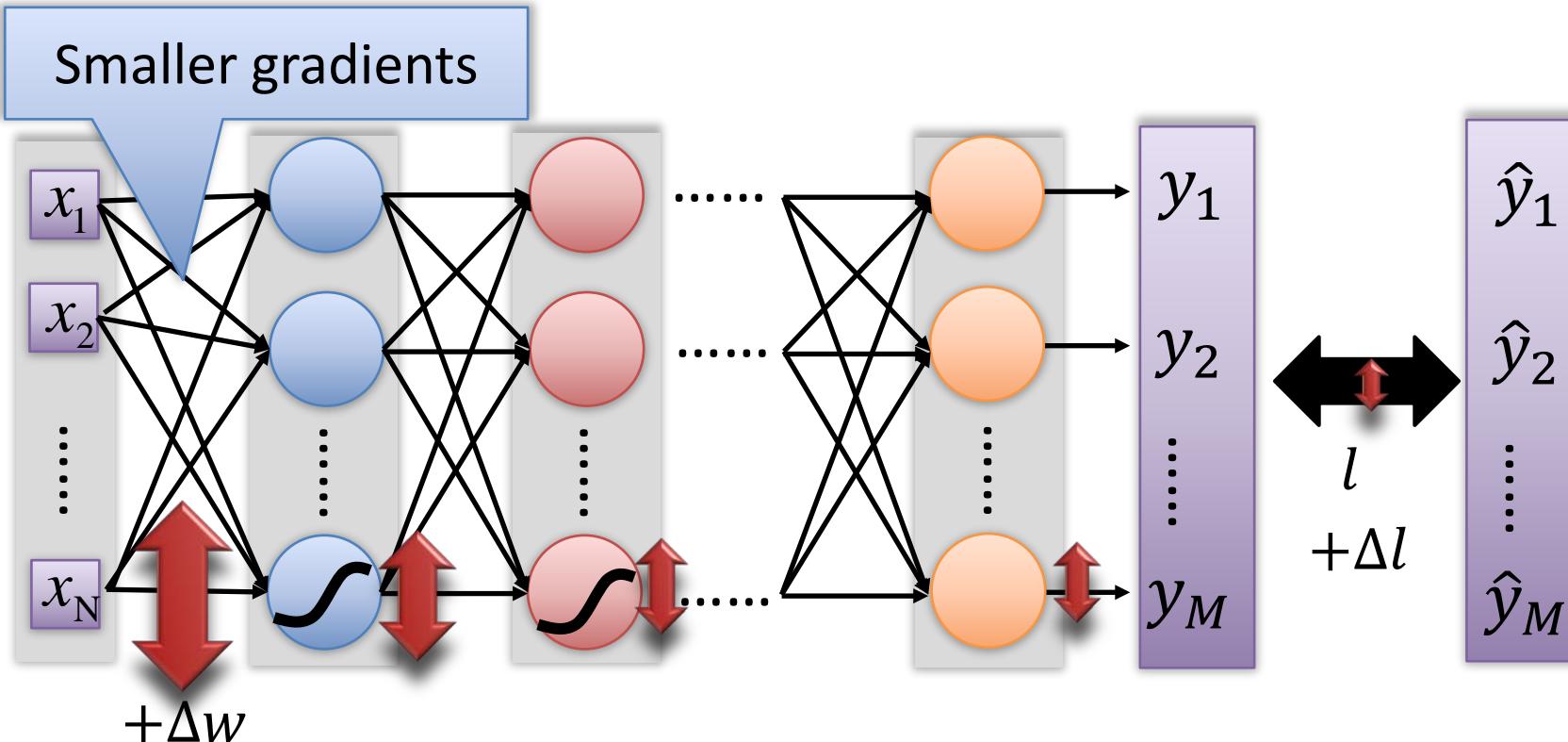
Intuitive way to compute the derivatives ...

$$\frac{\partial l}{\partial w} = ? \quad \frac{\Delta l}{\Delta w}$$



Recipe of Deep Learning

- *Vanishing Gradient Problem*



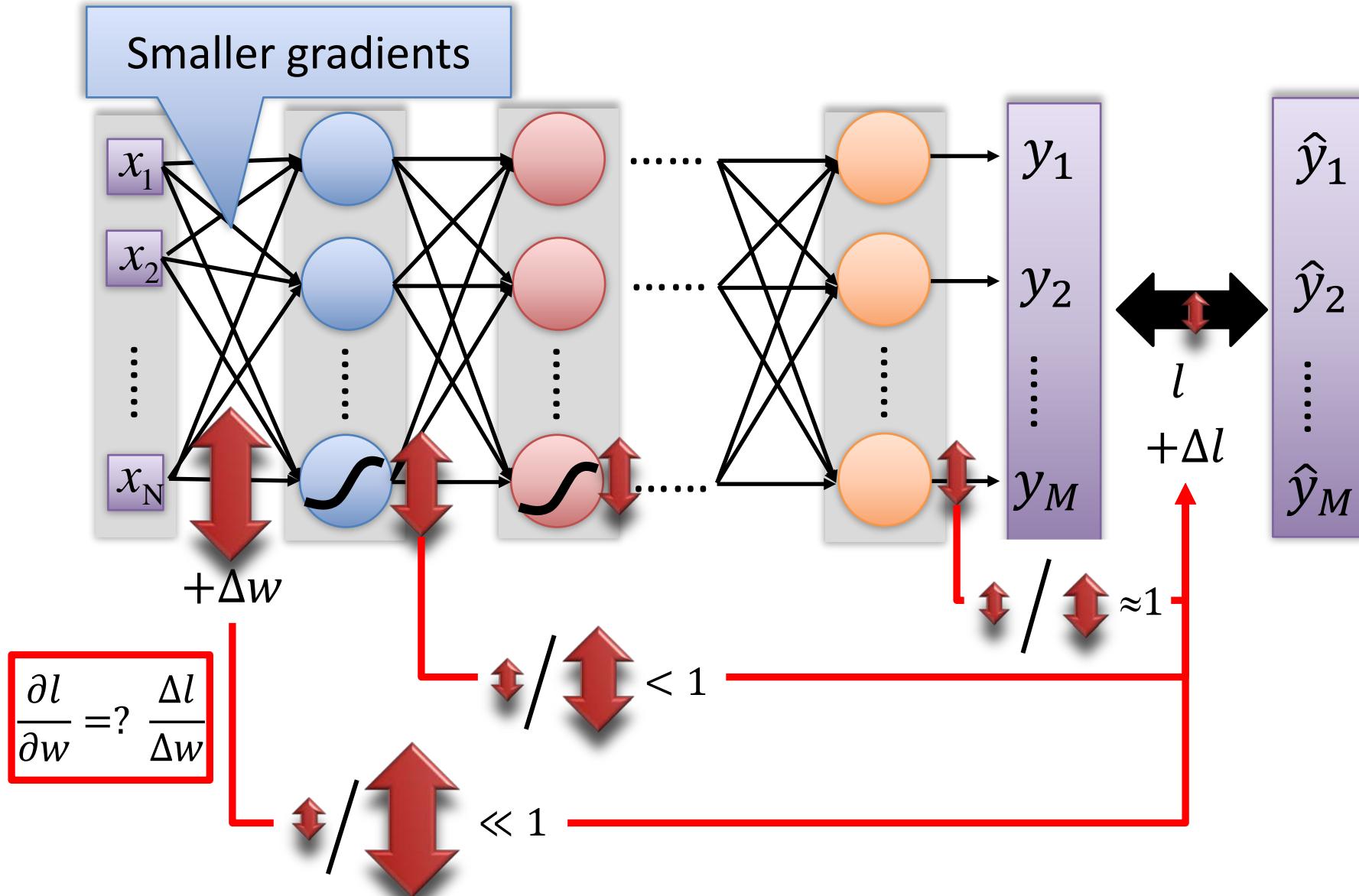
Intuitive way to compute the derivatives ...

$$\frac{\partial l}{\partial w} = ? \quad \frac{\Delta l}{\Delta w}$$



Recipe of Deep Learning

- *Vanishing Gradient Problem*



What are RNNs?

- The vanishing gradient problem

To understand why, let's take a closer look at the gradient we calculated above:

$$\frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \underbrace{\prod_{j=k+1}^3 \frac{\partial s_j}{\partial s_{j-1}}}_{\text{ }} \frac{\partial s_k}{\partial W}$$

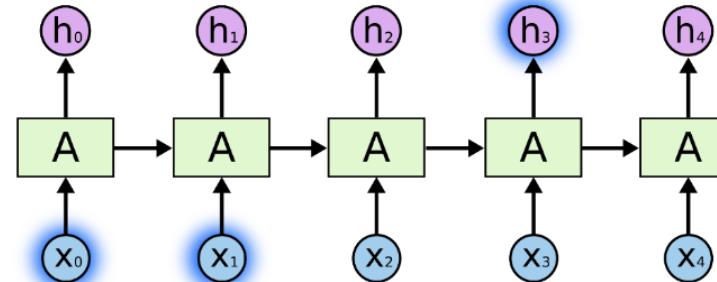
Because the layers and time steps of deep neural networks relate to each other through multiplication, derivatives are susceptible to vanishing

Gradient contributions from “far away” steps become zero, and the state at those steps **does not contribute** to what you are learning: You end up not learning long-range dependencies.

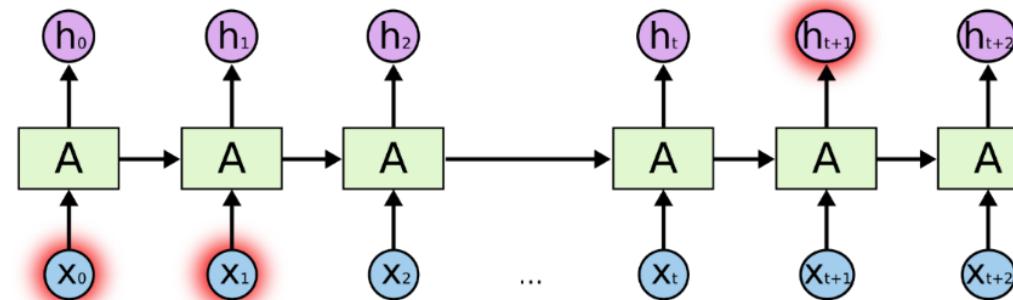


The Problem of Long-Term Dependencies

univ-cotedazur.fr

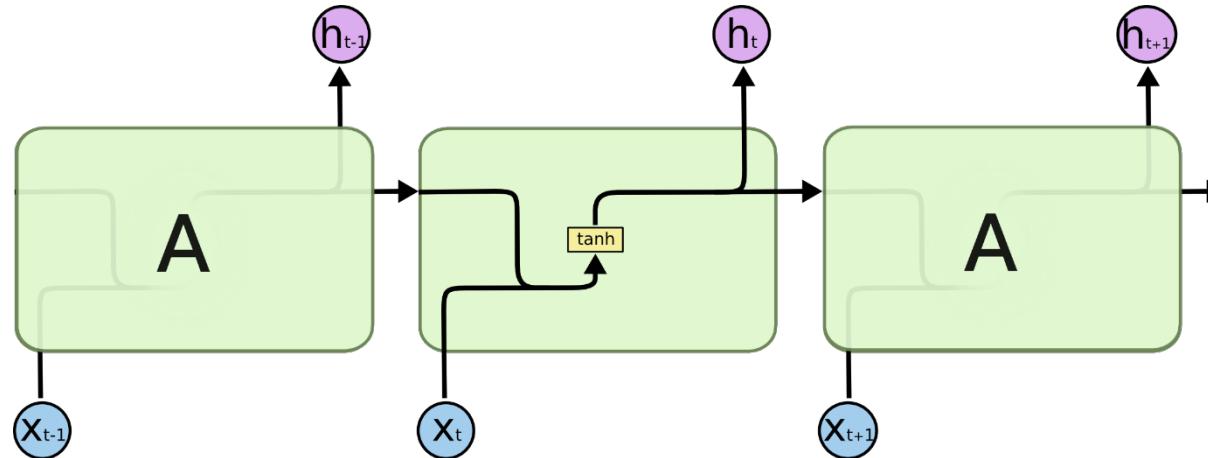


If we are trying to predict the last word in “the clouds are in the (sky),” we don’t need any further context – it’s pretty obvious the next word is going to be sky.



*Consider trying to predict the last word in the text “I grew up in France... I speak fluent (**French**).” Recent information suggests that the next word is probably the name of a language, but if we want to narrow down which language, we need the context of France, from further back.*

The repeating module in a standard RNN contains a single layer.



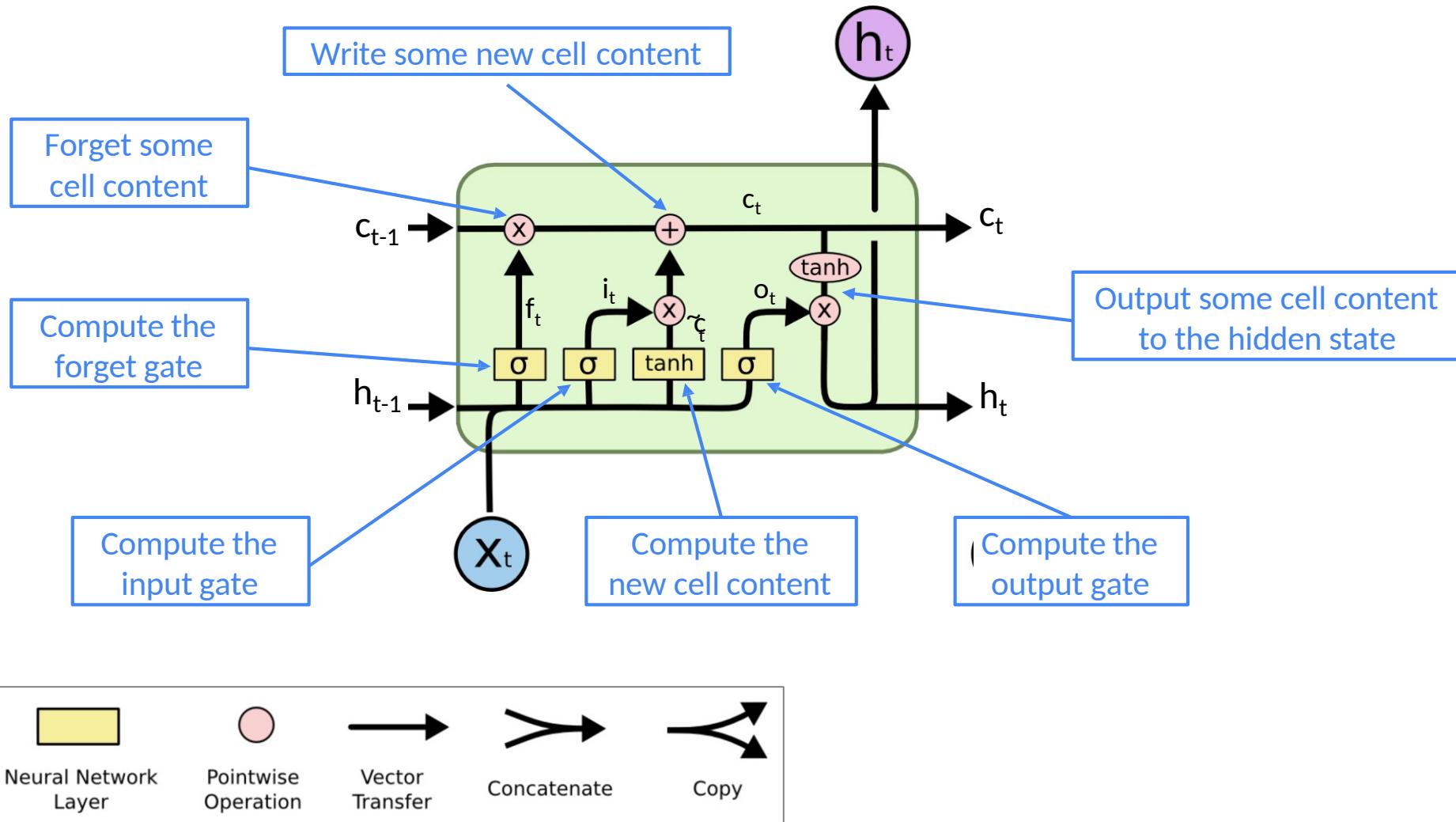
Unfortunately, as that gap grows, RNNs become unable to learn to connect the information (cf. vanishing gradients)

The problem was explored in depth by Hochreiter (1991) and Bengio, et al. (1994), who found some pretty fundamental reasons why it might be difficult.

Thankfully, LSTMs do not have this problem! (Hochreiter & Schmidhuber, 1997)

Review on your own: Long Short-Term Memory (LSTM)

univ-cotedazur.fr You can think of the LSTM equations visually like this:



Source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Review on your own: Long Short-Term Memory (LSTM)

univ-cotedazur.fr We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep t :

Forget gate: controls what is kept vs forgotten, from previous cell state

Input gate: controls what parts of the new cell content are written to cell

Output gate: controls what parts of cell are output to hidden state

Sigmoid function: all gate values are between 0 and 1

$$f^{(t)} = \sigma(W_f h^{(t-1)} + U_f x^{(t)} + b_f)$$

$$i^{(t)} = \sigma(W_i h^{(t-1)} + U_i x^{(t)} + b_i)$$

$$o^{(t)} = \sigma(W_o h^{(t-1)} + U_o x^{(t)} + b_o)$$

New cell content: this is the new content to be written to the cell

Cell state: erase (“forget”) some content from last cell state, and write (“input”) some new cell content

Hidden state: read (“output”) some content from the cell

$$\tilde{c}^{(t)} = \tanh(W_c h^{(t-1)} + U_c x^{(t)} + b_c)$$

$$c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

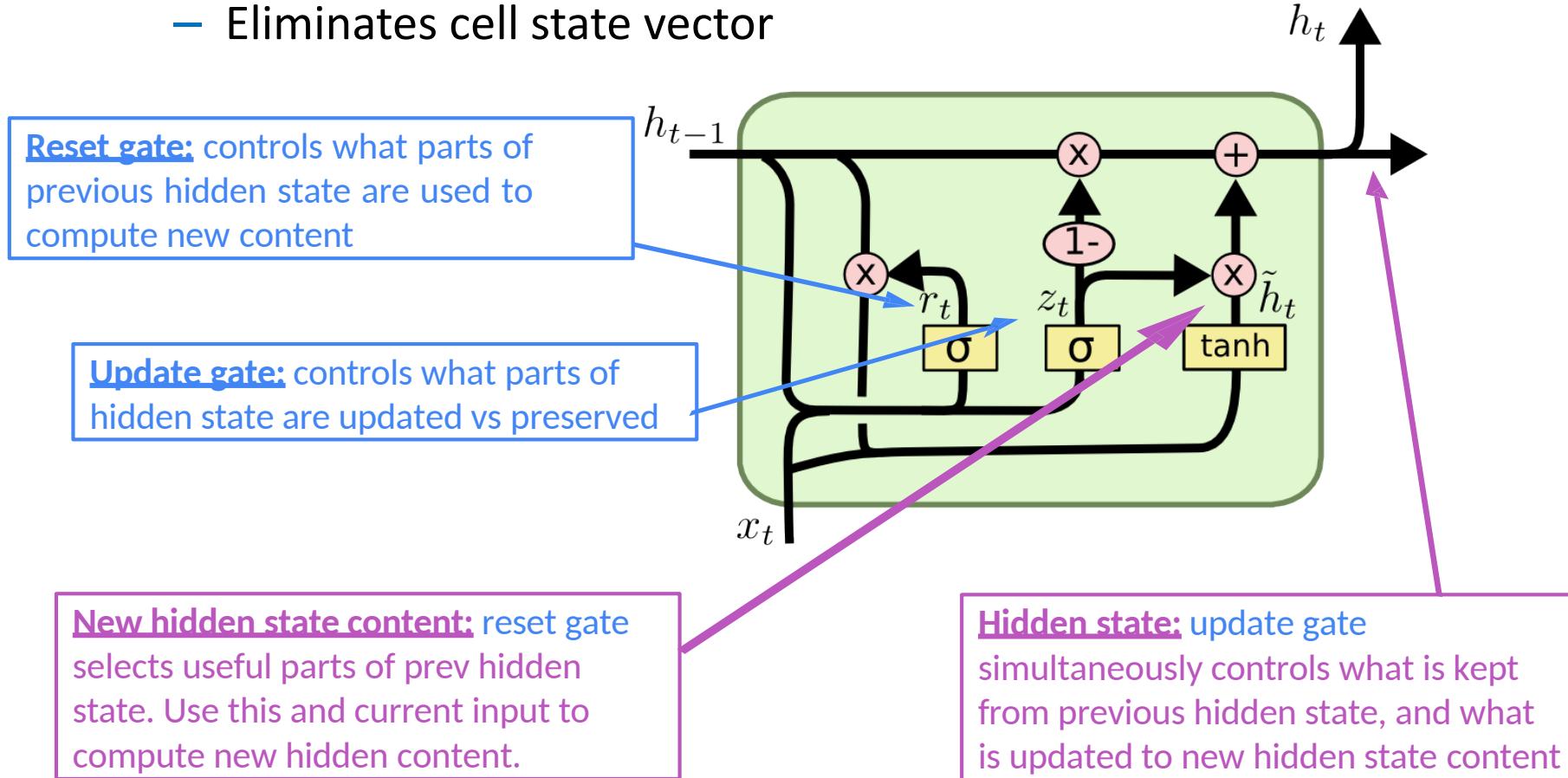
$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

Gates are applied using element-wise product

All these are vectors of same length n

Gated Recurrent Unit (GRU)

- Alternative RNN to LSTM that uses fewer gates ([Cho, et al., 2014](#))
 - Combines forget and input gates into “update” gate.
 - Eliminates cell state vector



Review on your own: Gated Recurrent Units (GRU)

univ-cotedazur.fr

- Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- On each timestep t we have input $\mathbf{x}^{(t)}$ and hidden state $\mathbf{h}^{(t)}$ (no cell state).

Update gate: controls what parts of hidden state are updated vs preserved

$$\mathbf{u}^{(t)} = \sigma(\mathbf{W}_u \mathbf{h}^{(t-1)} + \mathbf{U}_u \mathbf{x}^{(t)} + \mathbf{b}_u)$$

Reset gate: controls what parts of previous hidden state are used to compute new content

$$\mathbf{r}^{(t)} = \sigma(\mathbf{W}_r \mathbf{h}^{(t-1)} + \mathbf{U}_r \mathbf{x}^{(t)} + \mathbf{b}_r)$$

New hidden state content: reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

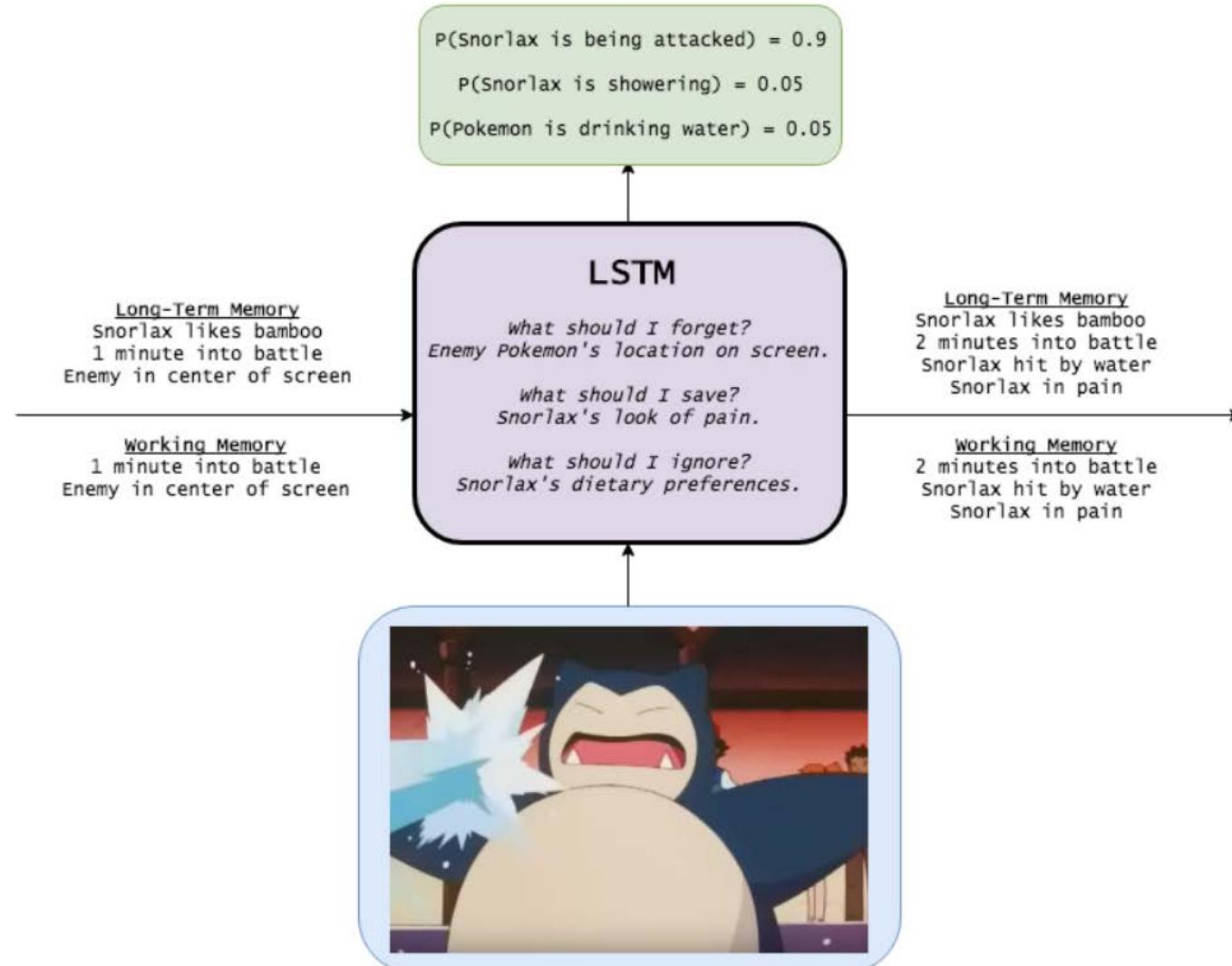
$$\tilde{\mathbf{h}}^{(t)} = \tanh(\mathbf{W}_h (\mathbf{r}^{(t)} \circ \mathbf{h}^{(t-1)}) + \mathbf{U}_h \mathbf{x}^{(t)} + \mathbf{b}_h)$$

$$\mathbf{h}^{(t)} = (1 - \mathbf{u}^{(t)}) \circ \mathbf{h}^{(t-1)} + \mathbf{u}^{(t)} \circ \tilde{\mathbf{h}}^{(t)}$$

Hidden state: update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

How does this solve vanishing gradient?
GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)

Cell State vs Hidden State





Activity

<http://blog.echen.me/2017/05/30/exploring-lstms/>



LSTM vs GRU

- Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used
- The biggest difference is that GRU is quicker to compute and has fewer parameters
- There is no conclusive evidence that one consistently performs better than the other
- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)
- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient



LSTMs: real-world success

- In 2013-2015, LSTMs started achieving state-of-the-art results for sequence modeling
 - Successful tasks include: handwriting recognition, speech recognition, machine translation, parsing, image captioning
 - LSTM became the dominant approach
- Starting in 2019, other approaches (e.g. Transformers) became more dominant for certain NLP tasks (will discuss next lecture)
 - For example in WMT (machine translation competition):
 - In WMT 2016, the summary report contains "RNN" 44 times
 - In WMT 2018, the report contains "RNN" 9 times and "Transformer" 63 times

Source: "Findings of the 2016 Conference on Machine Translation (WMT16)", Bojar et al. 2016, <http://www.statmt.org/wmt16/pdf/W16-2301.pdf>

Source: "Findings of the 2018 Conference on Machine Translation (WMT18)", Bojar et al. 2018, <http://www.statmt.org/wmt18/pdf/WMT028.pdf>