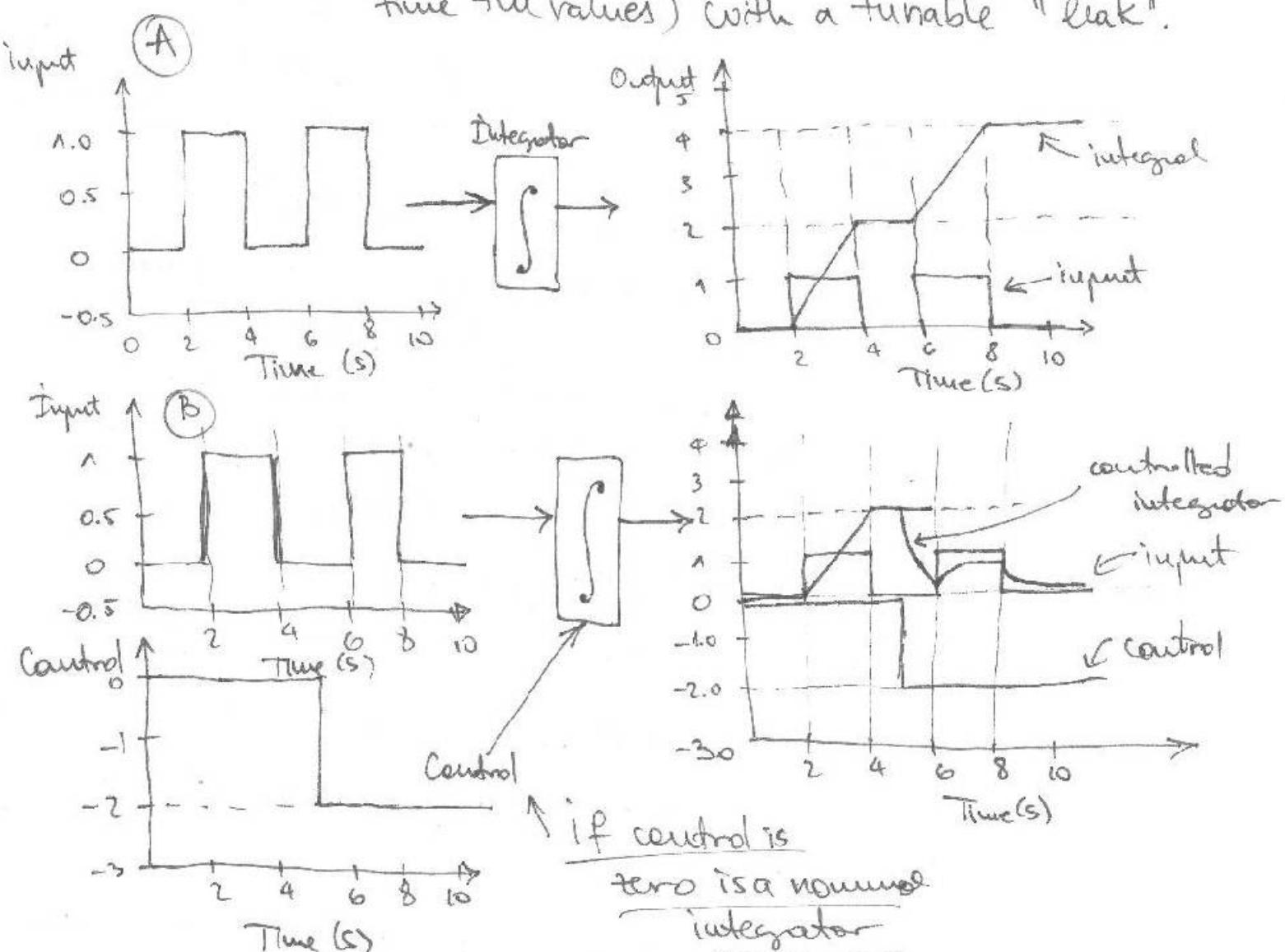


Spiking Neural Computation

- Neural Engineering Framework (NEF), 2003, 2013
Eliasmith, Anderson ; (Neuro-software open source)
- NEF as a general methodology that allows you to build large-scale, biologically plausible, neural models of cognition;
- does it make assumptions about what specific functions the brain performs ;
- it is a set of 3 principles that help determine how the brain performs some given function;
- NEF is a "neural computer"
 - * if you have a guess about the high-level function of a certain brain area or how neurons there respond NEF provides a way to connect groups of neurons ("populations" or "layers") to realize that function.
- The core of NEF is based on 3 principles :
 1. Representation(s): are defined by the combination of nonlinear encoding and weighted linear decoding.
 2. Transformations: are functions of the variables represented by neural populations.
 3. Dynamics : characterized by considering neural representations as states of a dynamical system \Rightarrow control theory analysis.

- The basic processing unit in NEF is typically a population of neurons, rather than single cells.
- Let's describe the 3 principles through an example:
a "controlled integrator".
 - * a simple neural circuit functioning as a simple memory
 - * it can be loaded with data, hold the information and then erase it
 - * this is in fact an integrator (i.e. accumulates in time the values) with a tunable "leak".



- (A). Standard integrator: continuously adds its input to its current state, the mathematical integral of the input.
- (B). controlled integrator: values less than 0 for the control makes the integrator forget (the speed of forgetting is proportional with the control variable).

NEF offers the possibility to solve problems with neurons:

- * specify the properties of the neurons;
- * specify the values to be represented and the functions to be computed;
↓
- * learns / solves for the connection weights between populations that will implement the desired function

NEF works with:

- feed-forward and recurrent connections.
- allows for complex dynamical systems such as integrators, oscillators, Kalman filters.
- incorporates local error-driven learning rules allowing for online adaptation and optimization of responses.

NEF has been used for:

- model visual attention,
- inductive reasoning,
- reinforcement learning,
- building the world's largest functional brain model using 2.5 million neurons to perform eight different tasks by interpreting visual input and producing hand-written output via a simulated 6-muscle arm.

Spaun ←



RBM Deep
Belief Network
(vision system)

+

working memory (cortex) including motor part)

+

action selection (basal ganglia)

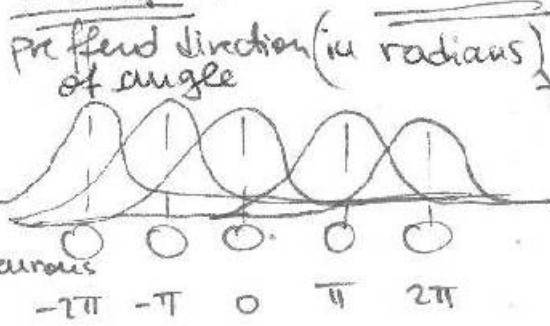
+

information routing (thalamus)

Representation in NEF

- distributed representations.
- differentiates between the activity of neurons in a population and the value being represented.
- the value is usually thought of as a vector x
- to map between x and activity a every

Example: neuron i has an encoding vector e_i



* the encoding is the preferred direction vector for that neuron : the vector for which the neuron will fire most strongly (tuning curves).

- given a value x , the activity of neuron i is:

$$a_i = G(\alpha_i \cdot e_i \cdot x + b_i)$$

where

α_i - gain parameter

e_i - encoding vector \rightarrow Kernel

b_i - bias

G - nonlinear activation function

(sigmoid, tanh, or LIF neurons)

- in order to also decode the value represented by a neural population activation we use a linear decoder d_i :

$$\hat{x} = \sum_i a_i \cdot d_i$$

d_i - is a set of weights that maps the activity back into a value, an estimate of x called \hat{x} .

- finding d_i is an optimization problem (least-squares minimization), we need to find the set of weights that minimizes the difference between x and its estimate \hat{x} :

input mapping to state interval (population) mapping

$$Y_j = \sum_{\mathbb{X}} a_j \cdot X$$

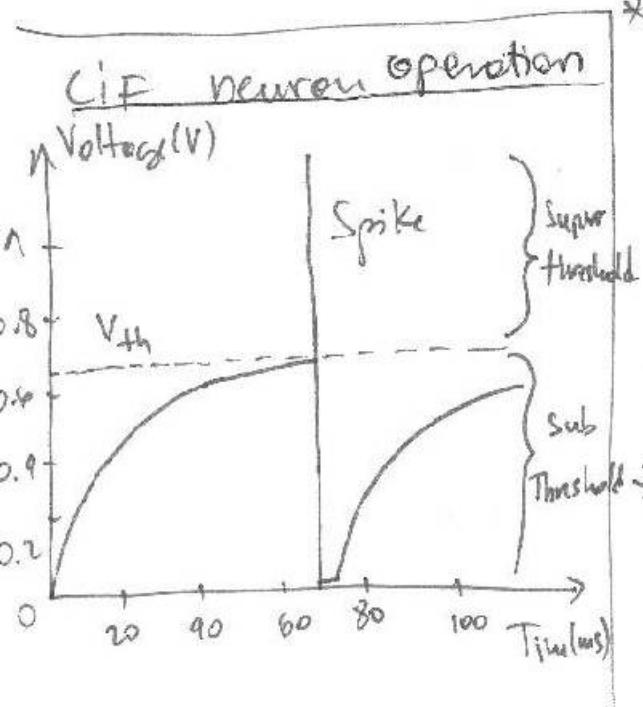
$$\Gamma_{ij} = \sum_{\mathbb{X}} a_{ij} \cdot a_j$$

$$d = \Gamma_{ij}^{-1} \cdot Y_j$$

- having this decoder allows you to determine how accurately a group of neurons is representing same value.

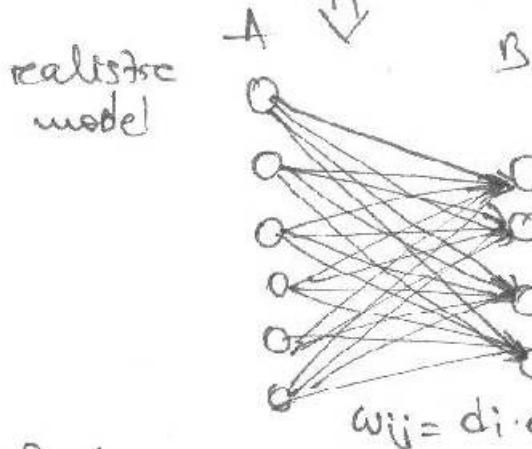
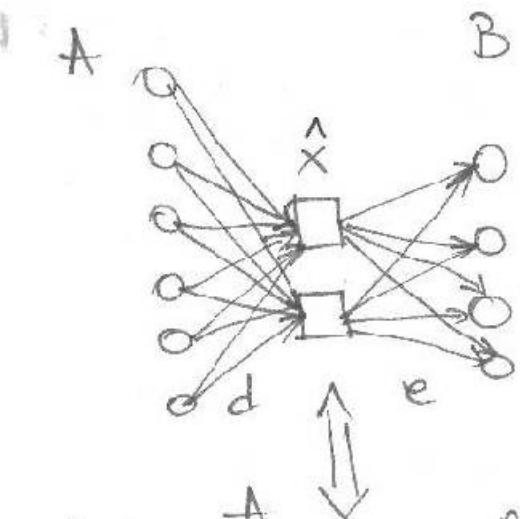
Transformation in NEF

- in order to do computations we need to connect neural populations together.
- consider 2 populations A and B and you want to pass information from A to B:
 - * if you set A to represent the value 0.2, then the synaptic connection weights between A and B should cause B to also represent value 0.2; Obviously this is computing the function $f(x) = x$.

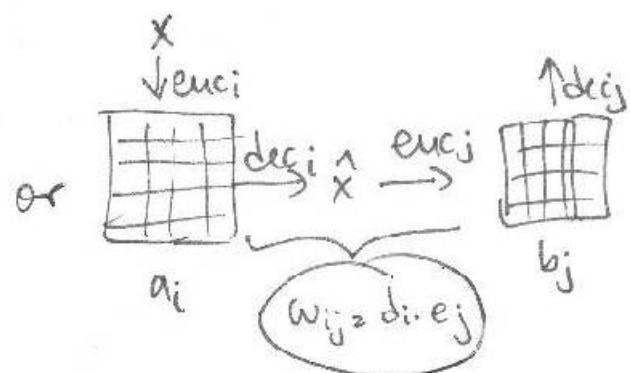


- * this cannot be realized by simply connecting neurons in A to neurons in B (might be different numbers) and the nonlinearity G makes this inaccurate.

- To create this in an accurate manner let's first assume we have an intermediate group of perfectly ideal linear neurons as shown next:



d is the set of weights to compute \hat{x}
 $\hat{x} = \sum_i a_i d_i$
 using the encoder e of B we can make B represent x
 $a_i = G(d_i e_i \hat{x} + b_i)$



- Such an approach is not limited to computing $f(x) = x$ in fact adjusting the weights $d^{f(x)}$ one can approximate any function $f(x)$.

- This approach can compute nonlinear functions with a single layer of connections \rightarrow no error back-propagation is required!

$$d^{f(x)} = \Gamma^{-1} \cdot \Upsilon^{f(x)}$$

$$w_{ij} = \alpha e_j \cdot d_i$$

where

$$\Gamma_{ij} = \sum_x a_i a_j$$

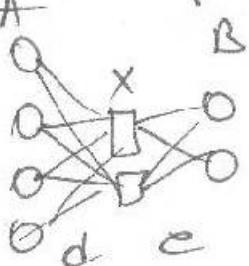
intra population
connections

$$\Upsilon = \sum_x a_j f(x)$$

activity to function
mapping

- this is possible because distributed representations allow for complex functions to be computed in a single set of connections \Rightarrow similar to SVM there is a projection from high-dimensional space to do computation
- random variations in $a_i \text{ and } b_i$ are required in
$$a_i = G(x_i \cdot e_i \cdot x + b_i)$$
because the function $f(x)$ will be approximated out of linear sums of the tuning curves

$$\rightarrow \text{more diverse tuning curves leads to a better approximation!}$$
- Such an approach is fast to simulate because instead of multiplying the activity of A by the full weight matrix you multiply by the decoders and then multiply by the encoders to produce the input to B.

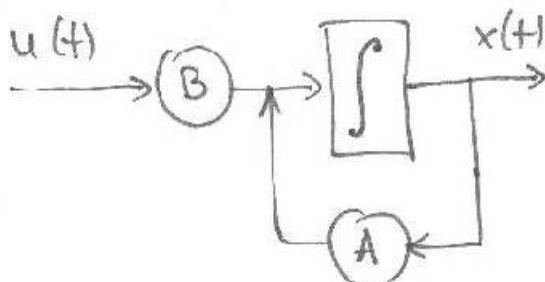


Savings in memory & simulation time!

Dynamics in NEF

- the representation and transformation are sufficient to produce models capable of computing functions of the form $y = f(x)$.
- NEF also offers the possibility to compute dynamic functions of the form $\boxed{\frac{dx}{dt} = A(x) + B(u)}$
when x is the value being represented, u is some input, A and B are arbitrary functions
- a special case, is returning to our initial example, the integrator, described by $\boxed{\frac{dx}{dt} = u \Rightarrow x(t) = \int u(t) dt}$
 - } if $u = 0$, the system maintains the current state, because $\frac{dx}{dt} = 0$
 - } if $u > 0$, the system will store an increasing value
 - } if $u < 0$, the system will store a decreasing value

Analogy to control systems:



LTI state equation

$$\dot{x}(t) = Ax(t) + Bu(t)$$

- A - determines what aspects of the current state affect the future state
- B - maps the input to the state space
- The integrator block - the transfer function / here is perfect integration

Deep Learning Vs. NEF

- connectionist model vs. biologically plausible dynamical model
- basic formulation:

$$x \rightarrow f \rightarrow y \quad \text{Error, } \mathcal{E} = \frac{1}{2}(y-x)^2$$

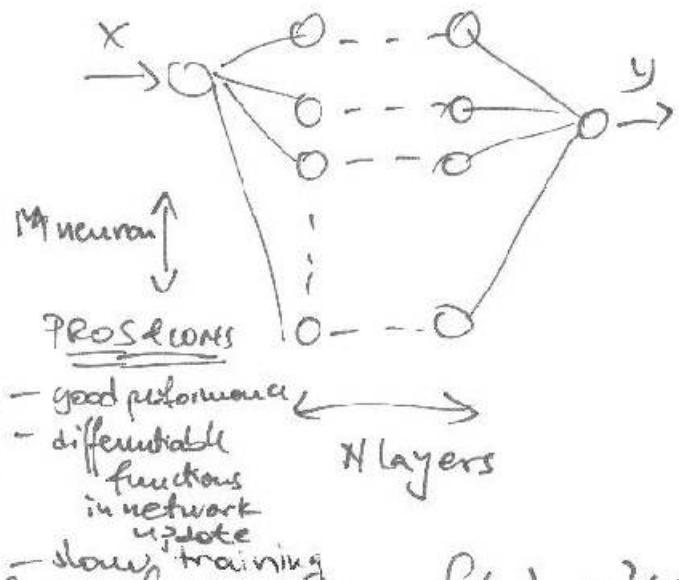
Gradient descent

$$\nabla f = \frac{\partial \mathcal{E}}{\partial w}$$

in deep nets

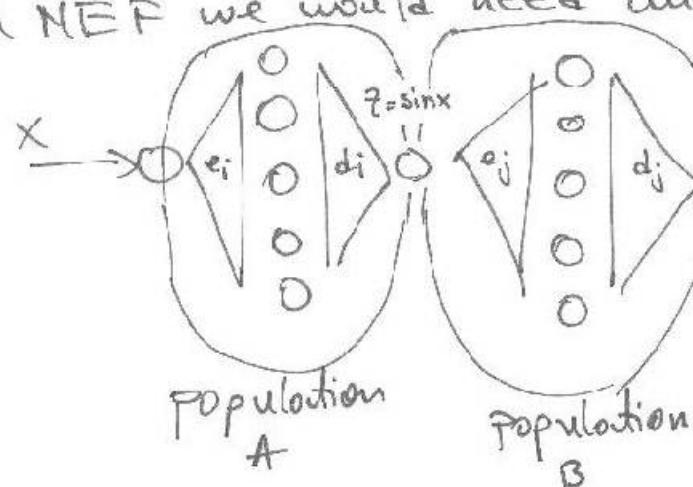
is easy to compute but also error prone, and problematic in large networks.

Another problem is that is computationally expensive.



Let's consider a simple problem of learning $f(x) = \sin^2(x)$

In NEF we would need an architecture as:



PROS & CONS

- this is fast
- interpretable
- requires more "design details" by the modeller

Math add-on : Least-squares optimization (LSO)

- consider n points (pairs) (x_i, y_i) $i=1 \dots n$.
- find the parameters of a model to fit the data.
- the model is of the form $f(x, \beta)$ where β contains the n parameters of the model.
- the fit of a model is measured by a residual defined as a difference between the actual value y_i and the prediction $f(x, \beta)$:

$$r_i = y_i - f(x, \beta)$$

- LSO finds the optimal β by minimizing the sum of squared residuals:

$$S = \sum_{i=1}^n r_i^2$$

- Example : linear least squares model

$$f(x, \beta) = \beta_0 + \beta_1 x$$

$$\frac{\partial S}{\partial \beta_j} = 2 \sum_i r_i \frac{\partial r_i}{\partial \beta_j} \text{ and because}$$

$r_i = y_i - f(x_i, \beta)$ the gradient equations

become

$$-2 \sum_i r_i \frac{\partial f(x_i, \beta)}{\partial \beta_j} = 0, j=1, \dots$$

This method is used to find the decoding weights d_i for $\hat{x} = \sum_i d_i x_i$ in NEF.

Immunological Computation I Artificially Immune Systems

Immunity: the reaction to foreign substances (pathogens) which includes

- * primary response
- * secondary response

Immune System: organs, cells, molecules

coordinated response
to a pathogen

↓
Immune response

Antigen (fungi, bacteria, viruses, other protozoa)

- * have the capability to trigger an immune response
- * depending on:
 - foreignness
 - molecular size
 - chemical composition

The Biological Immune System (BIS)

- detect and respond to foreign substances

- inherently distributed

- fault-tolerant

- complex behavior through interactions among its components

- self / non-self discrimination

Two types of immunity:

- (1) * innate immunity
- (2) * adaptive immunity

① Innate immunity

- initiates & regulates immune responses

② Adaptive immunity (Acquired immunity)

- * directed against specific invaders
- * consists mainly of lymphocytes (white blood cells):

} → B cells
} → T cells

- * lymphocytes recognize & destroy certain substances

- * substances that trigger a lymphocyte response are called antigens / immunogens

- * adaptive immune responses are antigen specific!

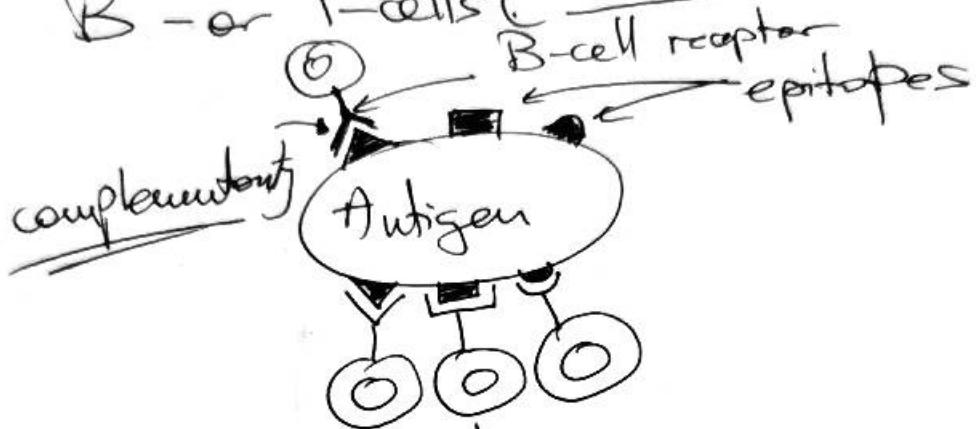
Biological defense mechanisms:

Nonspecific defense mechanism	Innate Immunity	Specific defense mechanism
Autonomic Barrier		Adaptive Immunity
<ul style="list-style-type: none"> • Skin • Mucous membranes • Secretions of skin & membranes 	<ul style="list-style-type: none"> • Phagocytic white blood cells • Antibacterial proteins • Inflammatory response 	<ul style="list-style-type: none"> • Lymphocytes • Antibodies

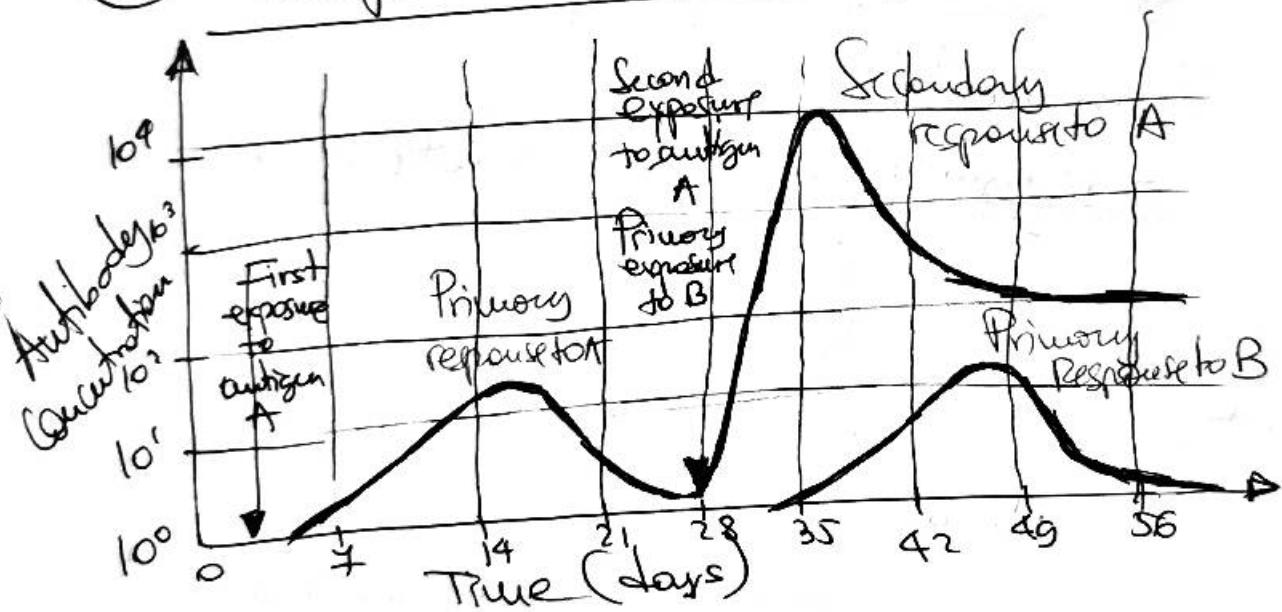
Immune System Dynamics

(A) Immune recognition (matching & binding)

- molecular bonds established between the surface of the antigen and receptors of B- or T-cells (complementary)



(B) Response to antigens



(C) T-cell maturation

- * produced initially by bone marrow
- * initially inert, they need to go through a maturity process in thymus

* during maturity they develop a T-cell receptor to recognize non-self peptides using selection (positive & negative)

- ④ B-cell Proliferation (affinity maturation)
- * when binding to an antigen B-cells undergo proliferation & differentiation
 - * the result are B-cells that are used as:
 - } (A) Memory cells : remember previous exposure to antigen to trigger faster response
 - (B) Plasma cells : secrete large quantities of antibodies

Immunological Computation

1. Pattern matching

- recognize antigens & generate specific responses
- recognition based on binding molecular shape
electrostatic charge

2. Feature extraction

- immune receptors (B, T cells) bind to sequences/ portions of peptides and not to complete antigens (i.e. epitopes)
- filter to extract matching segments in presence of molecular noise.

3. Learning & memory

- learns through interaction with environment
- proliferation of cells to determine memory & adaptive responses

④ Diversity

- selection (clonal selection) are testing for new detectors configurations
- highly combinatorial process to find optimal receptors
- exploration & exploitation for reproducing promising individuals

⑤ Distributed Processing

- not centrally controlled (i.e. nervous system) ^{opposite to}
- local detection & response
- distributed behavior of molecules & cells that make decisions in local environment

⑥ Self-Regulation

- adaptive response intensity depending on attack strength
- strong attacks need many resources
- after attack regulates to stop generating cells (new resources) and release the used resources.

⑦ Self-Protection

- self-defending (auto-immune diseases) anomalies as

Artificial Immune Systems

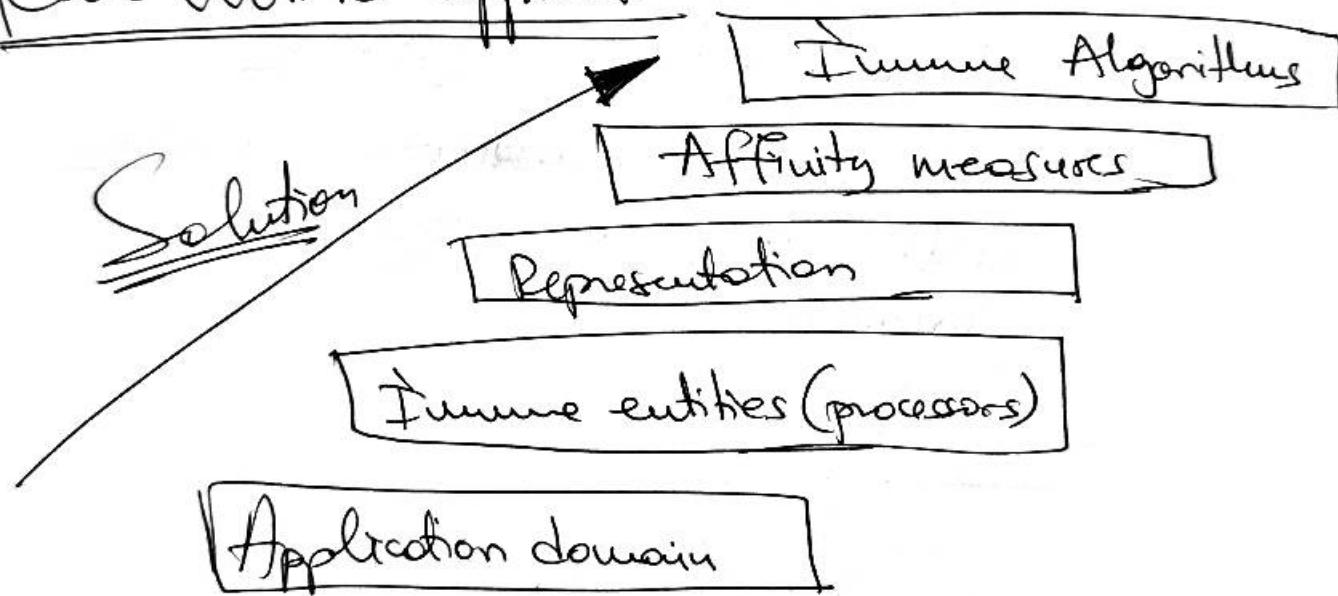
- mathematical & computational modelling of immunology
- algorithm design & implementation

1. Negative selection

2. Clonal Selection

3. Immune Networks

Real-World applications



Examples : Computer security, Robot control,
Defect/Fault diagnosis, Gaming, Software
Testing

Other types of Machine Learning

Models and Systems:

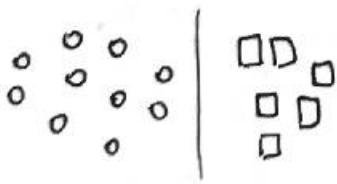
Statistical &
Stochastic ML

Support Vector Machines (SVM)

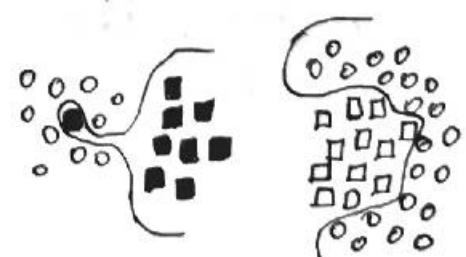
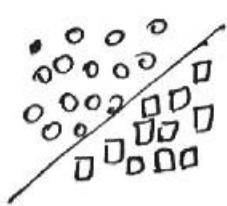
* Motivation:

- MLPs are simple but they lack optimality and they converge slowly (remember backprop).
- in a classification problem classes might not be linearly separable, and classes will overlap \Rightarrow we moved from single neuron to a MLP
- what if we can construct a nonlinear separation boundary by using a linear boundary in a large, transformed version of the feature space?

- SVMs: construct a hyperplane as the decision surface in such a way that the margin of separation between positive & negative examples is maximized for each training sample (binary classification) in feed forward networks, a binary learning machine basically.



Linearly separable classes



Non-linearly separable classes

A. Optimized hyperplane for linearly separable patterns

Consider the training sample $\{(x_i, d_i)\}_{i=1}^N$ where x_i is the input pattern for the i -th example and d_i is the corresponding target output (class label)

We assume that the pattern (class) represented by $d_i = +1$ and the pattern $d_i = -1$ are "linearly separable".

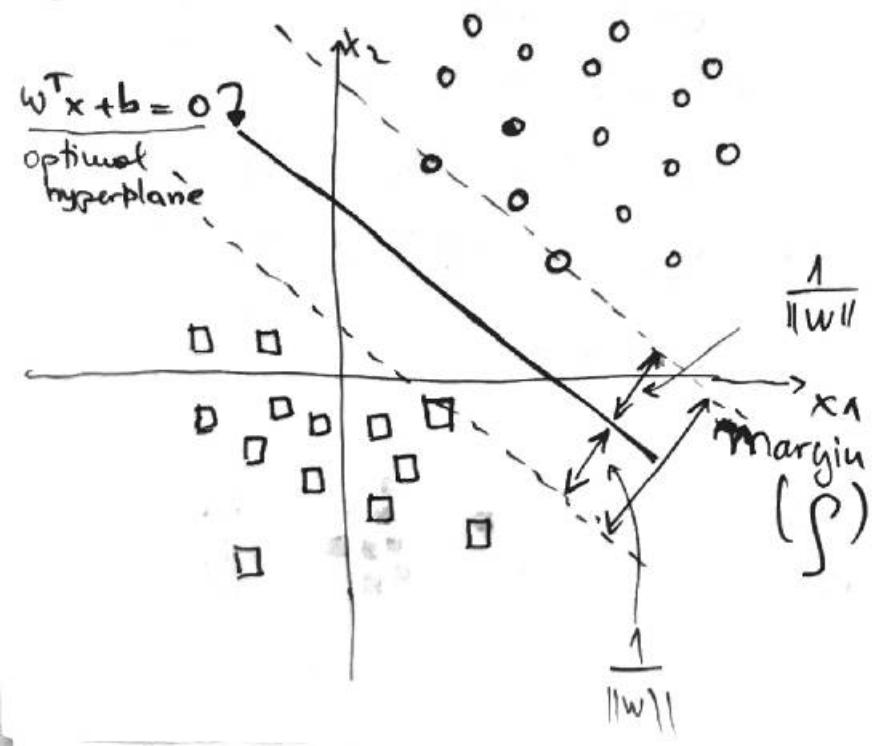
The equation of a decision surface (hyperplane) that does the separation is :

$$w^T \cdot x + b = 0$$

where x = input vector, w = adjustable weight vector, b = bias. So, we can write:

$$\begin{cases} w^T x_i + b \geq 0 & \text{for } d_i = +1 \\ w^T x_i + b < 0 & \text{for } d_i = -1 \end{cases}$$

For a 2D case



The goal of the SVM is to find a particular hyperplane for which the separation margin is maximized!

↓
find the optimal hyperplane

If we let w_0 and b_0 be the optimal values of the weight vector and bias, respectively, the optimal hyperplane representing a multidimensional linear decision surface is defined as:

$$w_0^T \cdot x + b_0 = 0$$

The discriminant function

$$g(x) = w_0^T \cdot x + b_0$$

gives an algebraic distance from x to the optimal hyperplane, where x can be represented as

$$x = x_p + r \cdot \frac{w_0}{\|w_0\|}$$

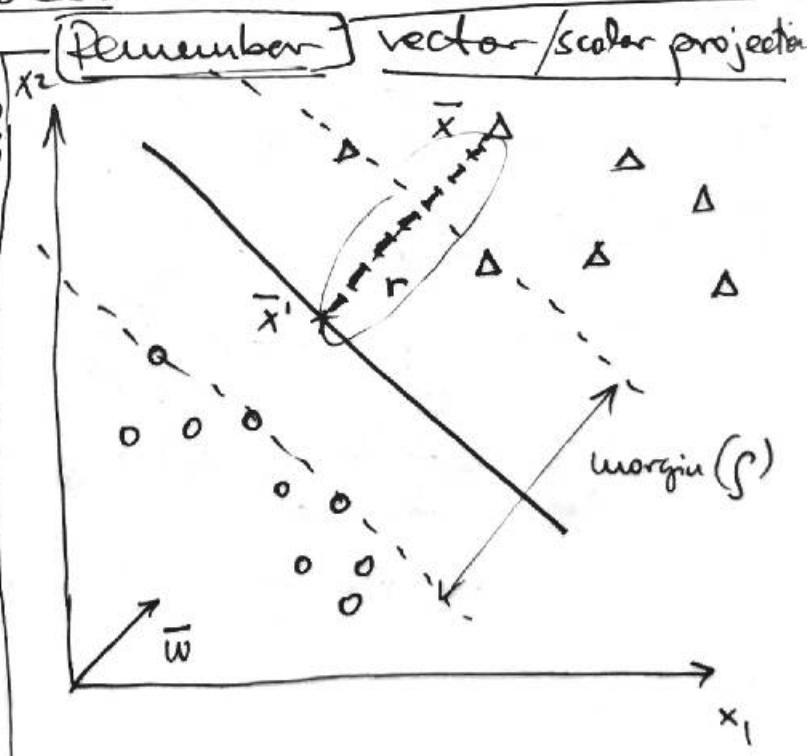
where x_p - the normal projection of x on the hyperplane and r - the desired algebraic distance; r is positive if x is on the positive side of the hyperplane and negative if x is on the negative side.

What is the Euclidean distance from \bar{x} to the decision boundary?

The shortest distance from a point to a hyperplane is perpendicular to the plane, so parallel to \bar{w} . A unit vector in this direction is $\frac{\bar{w}}{\|\bar{w}\|}$.

The worked line (---) is a translation of $r \cdot \frac{\bar{w}}{\|\bar{w}\|}$. Then the projection \bar{x}' is

$$\bar{x}' = \bar{x} + r \cdot \frac{\bar{w}}{\|\bar{w}\|}$$



In our case if we consider $g(x) = w_0^T x + b_0$ and
 $x = x_p + r \cdot \frac{w_0}{\|w_0\|}$ then

$$x_p = x - r \cdot \frac{w_0}{\|w_0\|} \text{ and because } x_p \text{ is on the hyperplane, } g(x_p) = 0 \text{ so,}$$

$$g(x_p) = w_0^T \left(x - r \frac{w_0}{\|w_0\|} \right) + b_0 = 0$$

$$w_0^T \cdot x - w_0^T \cdot r \cdot \frac{w_0}{\|w_0\|} + b_0 = 0$$

$$(w_0^T \cdot x + b_0) - w_0^T \cdot w_0 \cdot \frac{r}{\|w_0\|} = 0$$

$$w_0^T \cdot w_0 \cdot \frac{r}{\|w_0\|} = w_0^T \cdot x + b_0 \Rightarrow r = \frac{w_0^T \cdot x + b_0}{\|w_0\|}$$

but $\|x\| = \sqrt{x^T \cdot x}$ in Euclidean space dot product is inner product

$$r = \frac{g(x)}{\|w_0\|} \quad \underline{\text{q.e.d}}$$

The points closest to the optimal hyperplane are called "support vectors" and have an important role in SVM (& they are the honestest to classify).

As we shown that

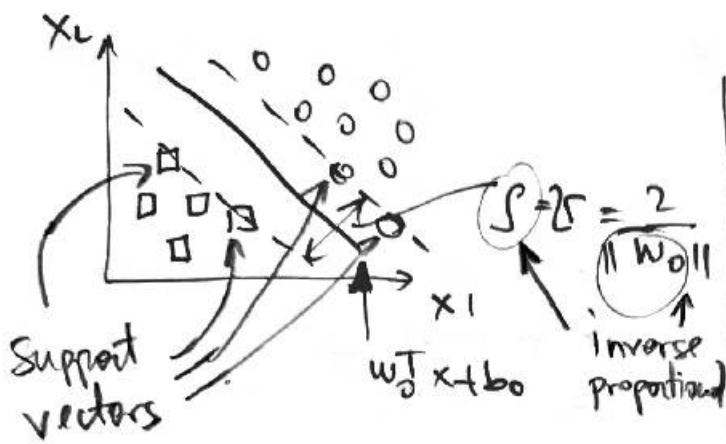
$$r = \frac{g(x)}{\|w_0\|}$$

the algebraic distance from a support vector $x^{(s)}$ to the optimal hyperplane is

$$r = \frac{g(x^{(s)})}{\|w_0\|} = \begin{cases} \frac{1}{\|w_0\|} & \text{if } d^{(s)} = +1 \\ -\frac{1}{\|w_0\|} & \text{if } d^{(s)} = -1 \end{cases}$$

margin of separation

$$\text{so } P = 2r = \frac{2}{\|w_0\|}$$



Maximizing the margin of separation p is equivalent to minimizing the Euclidean norm of the weight vector.

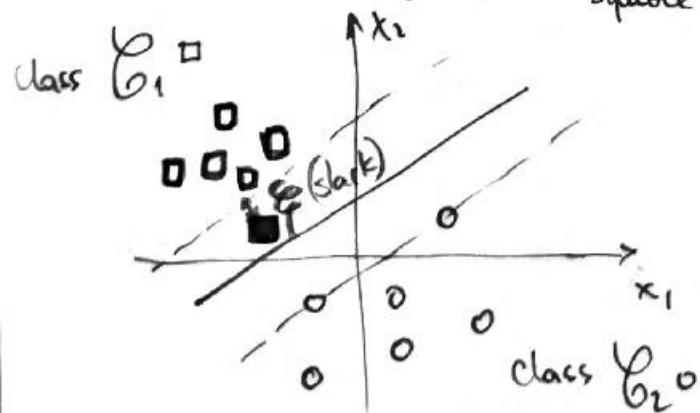
B. Optimal hyperplane for non-linearly separable patterns

We want to find a hyperplane which minimizes the probability of classification error, averaged over the training samples.

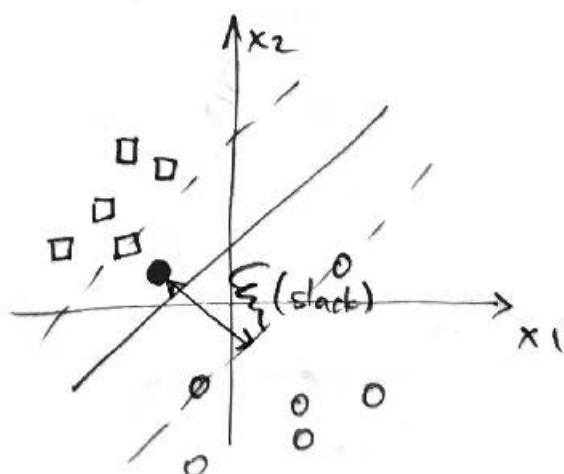
The margin of separation between classes is said to be soft if a data point (x_i, d_i) violates the following condition:

$$d_i \cdot (w^T x_i + b) \geq +1, i=1, 2, \dots, N$$

- ① The point falls inside the separation region but on the correct side of the decision surface (bold point square)



- ② The point falls on the wrong side of the decision surface (bold point circle)



In order to measure the deviation of the data points from the condition of separability we introduce the slack variable ξ_i such that:

$$d_i \cdot (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, i=1,2,\dots,N$$

falling on the correct side $0 < \xi_i \leq 1$

falling on the wrong side $\xi_i > 1$

Learning in SVM

The task is to find the separating hyperplane for which the miscalcification error, averaged over the training data, is minimized.

A. For the linear case the learning problem boils down to a constrained optimization problem:

Given the training sample $\{(\mathbf{x}_i, d_i)\}_{i=1}^N$, find the optimal weights w and b so that

$$d_i \cdot (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i=1,2,\dots,N$$

and the weight vector w minimizes the cost function

$$\phi(w) = \frac{1}{2} \cdot w^T \cdot w$$

B. In the nonlinear case given the training sample $\{(x_i, d_i)\}_{i=1}^N$ find the optimum values of the weight vector w and bias b such that they satisfy the constraint:

$$d_i \cdot (w^T \cdot x_i + b) \geq 1 - \xi_i \quad i=1,2,\dots,N$$

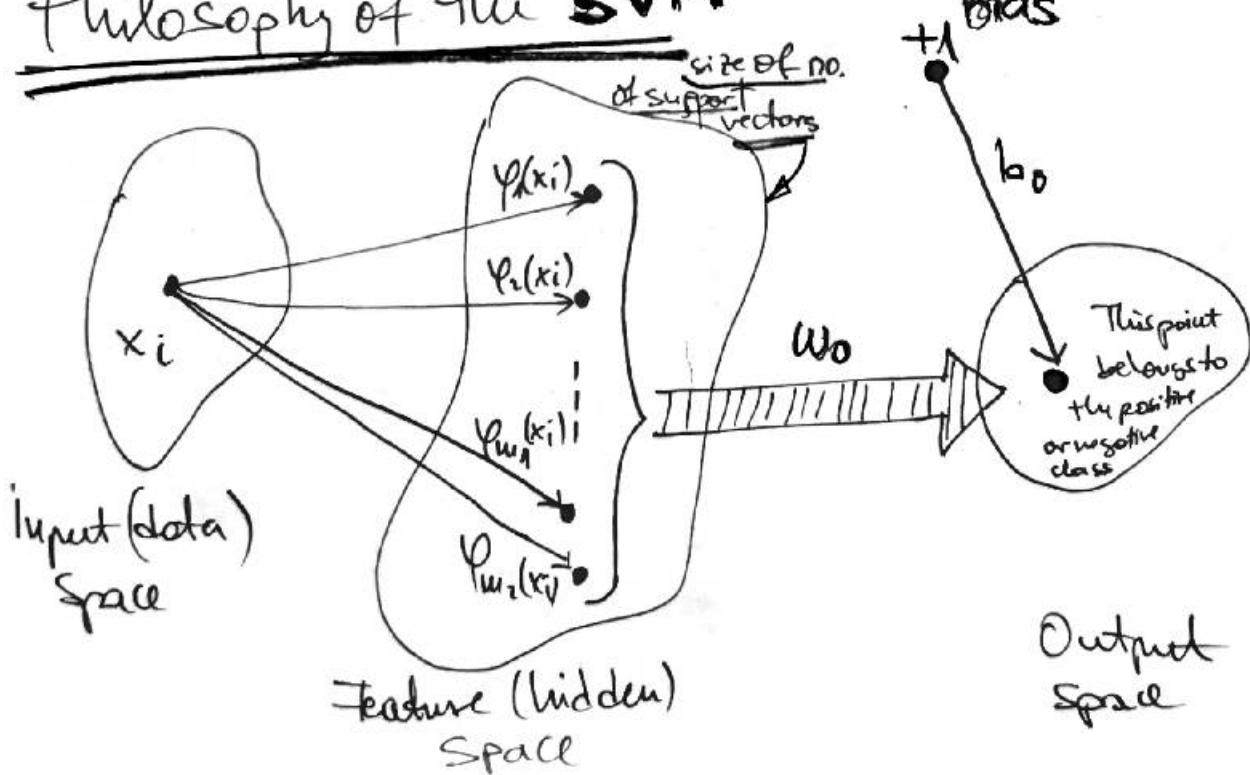
$$\xi_i \geq 0 \text{ for all } i$$

and such that the weight vector w and the slack variables ξ_i minimize the cost function:

$$\phi(w, \xi_i) = \frac{1}{2} \cdot w^T \cdot w + C \sum_{i=1}^N \xi_i$$

C - user-defined positive parameter

Philosophy of the SVM



Unsupervised Learning - AIML Course

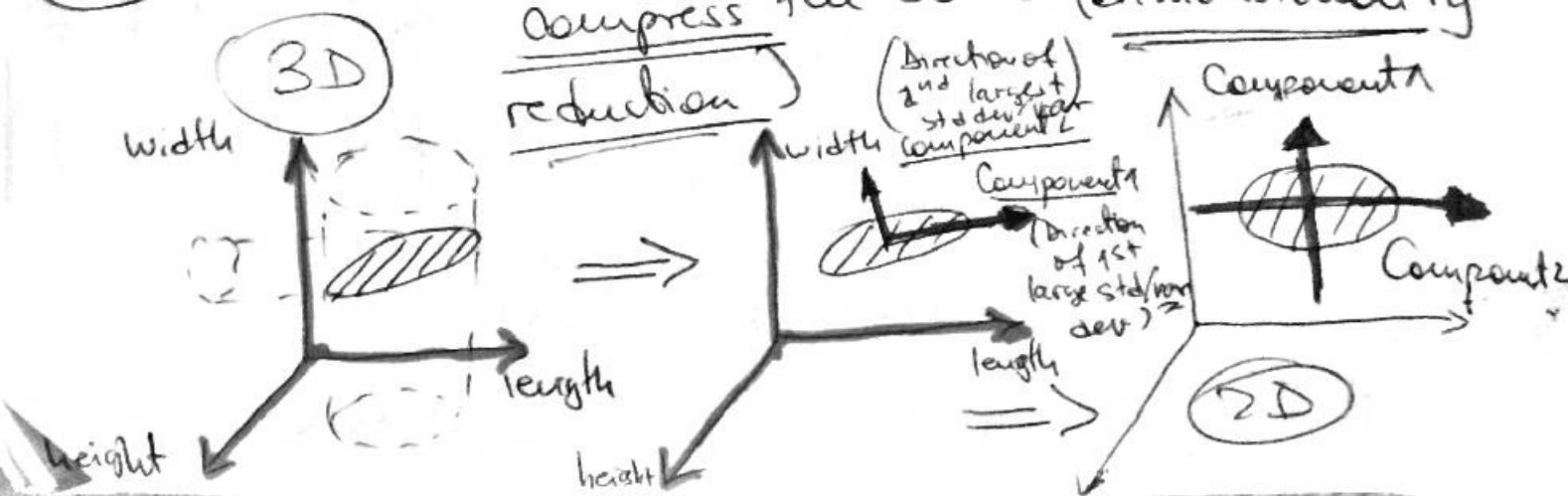
Recall we can have 3 types of learning algorithms:

- Supervised learning (learning with a teacher
⇒ we have an error signal).
- unsupervised learning (learning without a teacher ⇒ we do not have an error signal).
clustering, learning structure
- reinforcement learning (learn through the interaction with the environment ⇒ reward signal to adjust policy).

Unsupervised Learning assumes:

- self-organization to produce output vectors by updating weights;
- does not require labels, discovers classes or clusters
- learns the statistical regularities of the data;
- developing a representation of the internal structure of data or compress the data (dimensionality reduction)

TCA



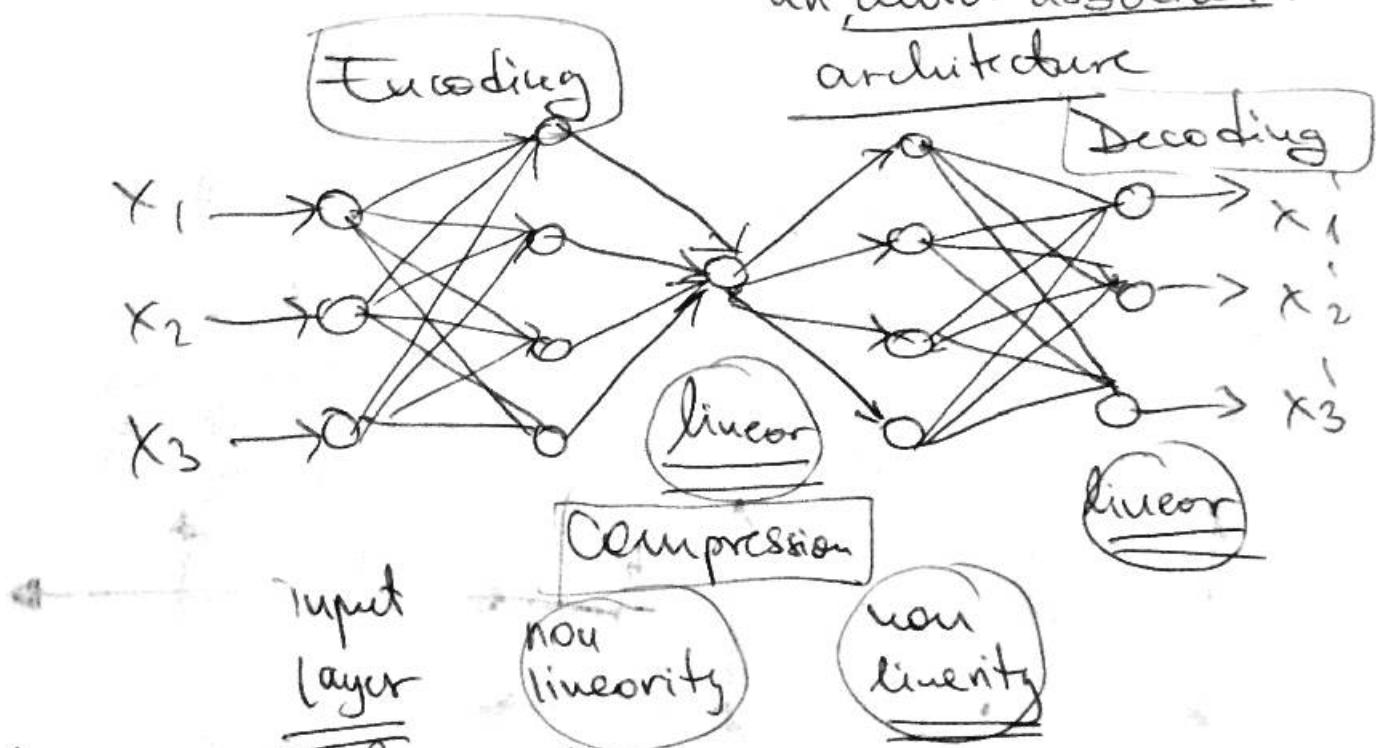
PCA - finds a suitable low-dimensional representation of a high-dimensional space:

- provides a sequence of best linear approximations to a given high-dim observation;
- uses a low-dimensional manifold first 2 components (i.e. largest variance) to represent the most variability in the input space;
- it applies a linear transformation so its application is limited

If a solution

multilinear PCA

- achieved with a neural network - an MLP with an auto-associative architecture



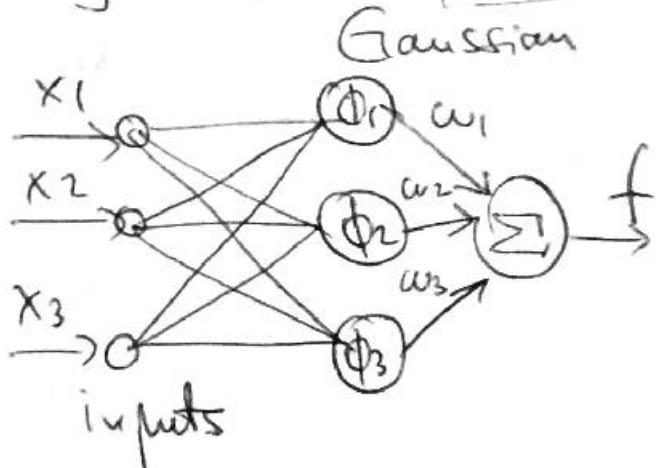
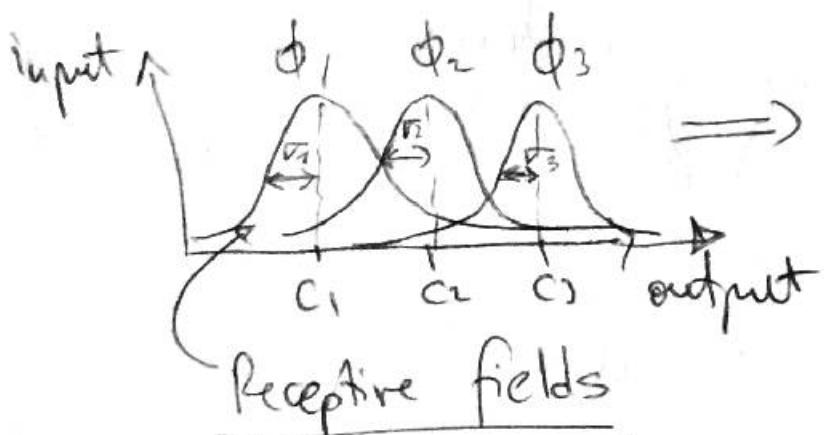
Architecture for: pattern recognition, data compression

Radial Basis Functions (RBF)

- when solving nonlinearly separable pattern-classifier there is a benefit to map the input space to a higher dimensional space;
- a nonlinear mapping can help to transform a nonlinear separable classification into a linearly separable one with high probability
- RBFs offer an approach to approximate any arbitrarily function if a sufficient number of radial-basis functions is given with appropriate parameters:

$$f(x) = \sum_{i=1}^N w_i \cdot \phi\left(\|x - c_i\|\right) \quad -\frac{x^2}{\beta}$$

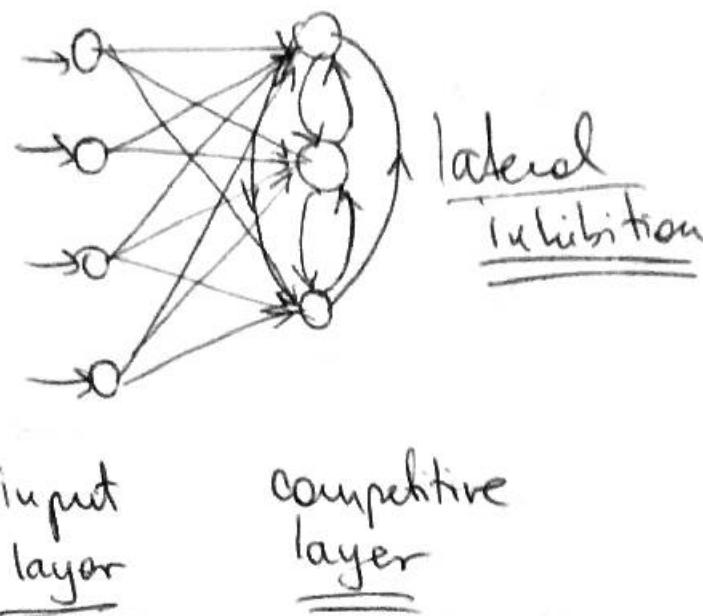
where ϕ can take many shapes, e.g. $\phi(x) = e^{-\frac{x^2}{\beta}}$



- area of the input space where the neuron will elicit a response

Vector Quantization (VQ)

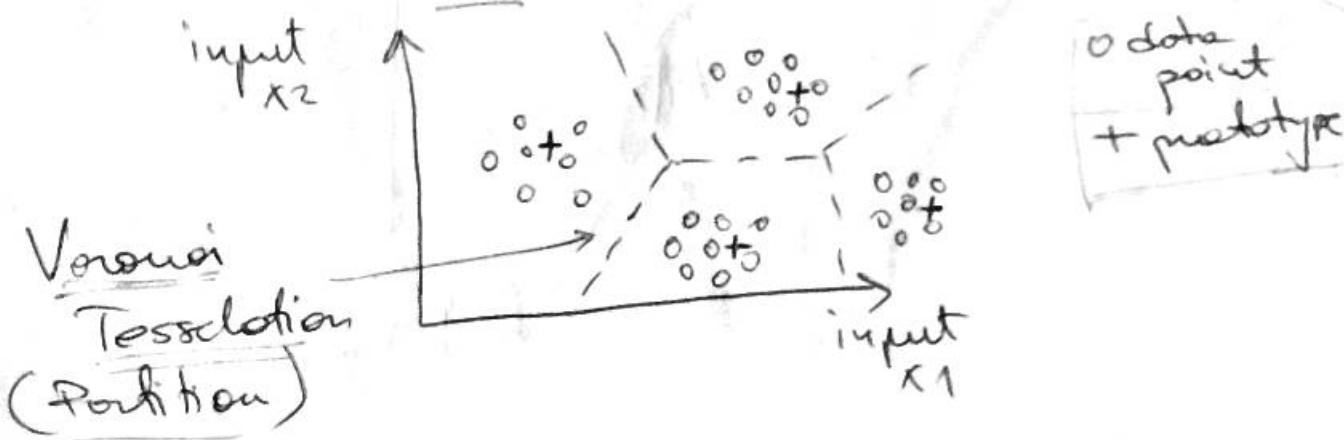
- in neurobiology synapses are strengthened or weakened in a competitive process;
- in the learning process there is competition among neurons which to fire first;
- the neuron with the greatest total input will "win" the competition and will fire, whereas the others will switch off in a process named "Winner-Take-All" (WTA)
- Typical network to implement such a behavior :



- The basic mechanism is to find a winning unit and update its weight to make it more likely to win in the future if similar input is given to the network.

VQ is a form of competitive learning

- finds discovered structure in the input data
- is a form of lossy compression
(some information in the input is lost in the process.)



Algorithm:

- #1 choose the no of clusters, M
- #2. Randomly initialize prototypes w_1, w_2, \dots, w_M
- #3. Repeat until "good enough"
 - #4. Randomly pick an input x
 - #5. Find the winning prototype $\underline{w_K}$

Competition process

$$\|w_K - x\| \leq \|w_i - x\| \text{ for all nodes}$$

- #6. Update weights for the winner

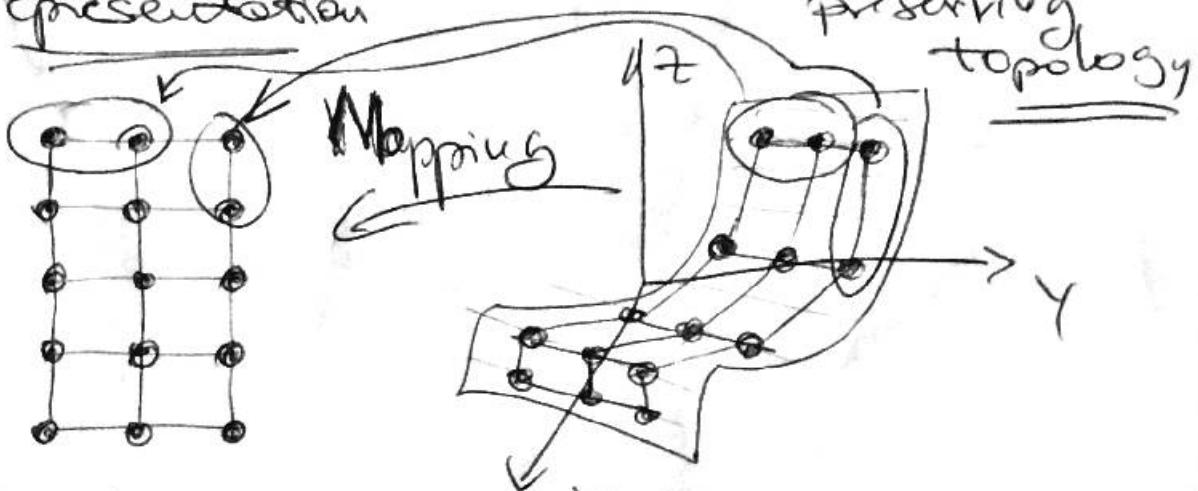
$$w_K(t+1) = w_K(t) + \eta \cdot (x - w_K(t))$$

$\|\cdot\|$ or $\|\cdot\|_2$ distance, norm or metric (e.g. Euclidean)

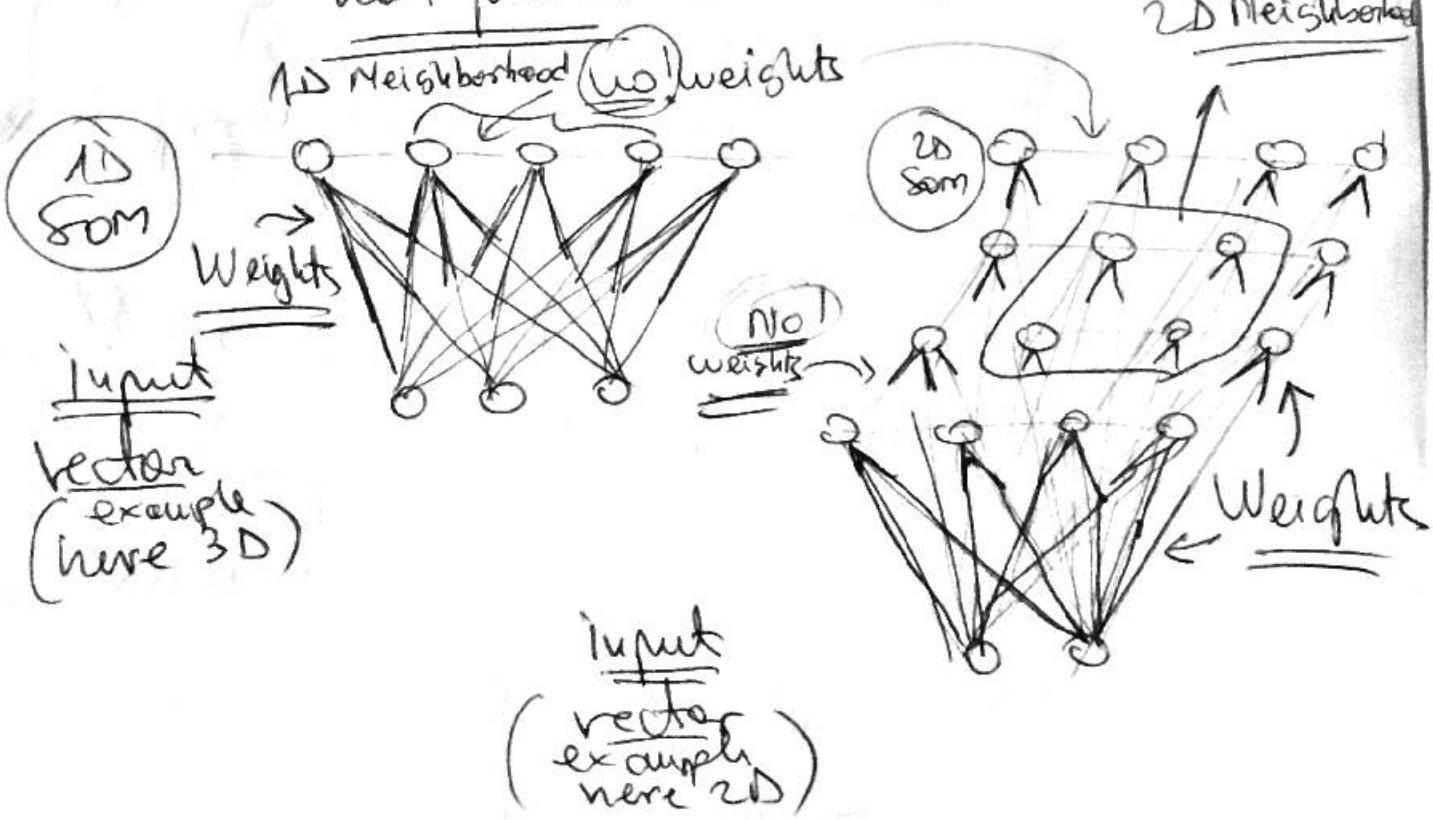
η - learning rate

Self-Organizing Maps (SOM)

- motivated by the retina-cortex mapping
- topologically preserving mapping: close points in the input space will be close also in the representation



- Think about a squashed flower $3D \rightarrow 2D$
- we can have 1D or 2D SOM
 - 1D - forced to converge mathematically
 - 2D - not forced to converge



- SOM can develop an internal representation of the input and create new classes automatically

Algorithm:

#1. Initialize weights w_{ij} , define neighborhood function $\phi(i, j)$ which changes with $\text{Edu} \rightarrow \text{dist}$ & $T(t)$

→ #2. Select input x and find the winner (Best Matching Unit)

$$\|x - w_k\| \leq \|x - w_j\| \text{ for all nodes } j \neq i$$

#3. Update weights for all units j given the winner i

$$w_j(t+1) = w_j(t) + \eta \phi(i, j)(x - w_j(t))$$

η - learning rate

#4. Repeat from #2 until convergence and update η and ϕ

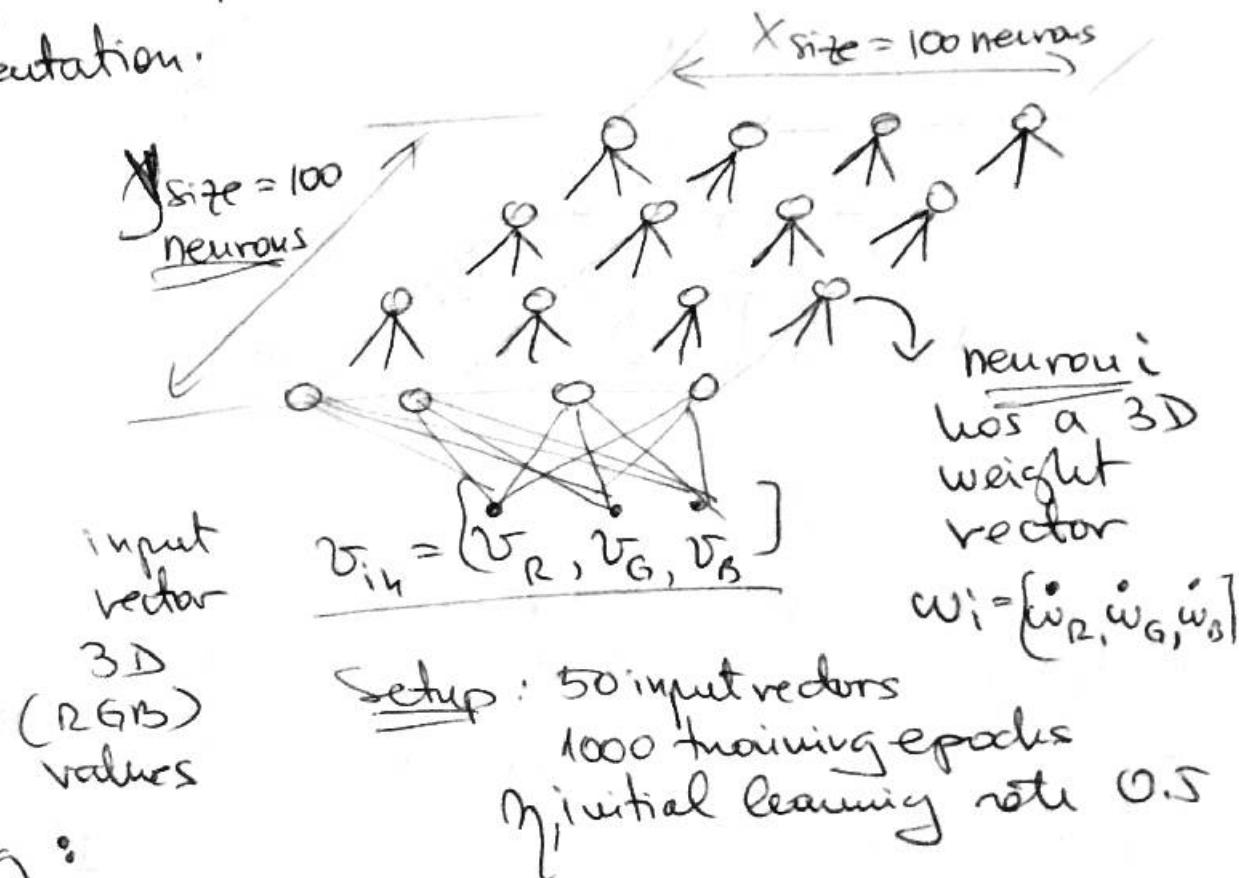
Typical problems:

* Mapping to 2D - faulty convergence
"fuzz" appears in the weights

* Mapping to a too low dimension
3D cube to 2D - losing a dimension

SOM example (dme)

Ex: Mapping colors from 3D colorspace (RGB) to a 2D representation.



- #1. Initialize each neuron's weight randomly in the range of the input 0 - 255 (range of colors)
- #2. Pick a random input vector v_{ih} from dataset and present it to the network
- #3. Find the winner / Best Matching Unit (BMU) by comparing each neuron's weight vector to the input
- #4. Compute the radius of the neighborhood function and adjust it.
- #5. Update the weights of the neurons in the neighborhood of the BMU closer to the input.
- #6. Repeat #2 for N epochs = 1000

Implementation details:

(A) Computing the winning neuron (BMU)

- use Euclidean distance to find the "closest" neuron to the input

$$2\text{Norm} \quad \|v_{in} - w_{BMU}\| = \min_i \{ \|v_{in} - w_i\| \}$$

$$\text{where } \|\cdot\| = D = \sqrt{\sum_{i=0}^D (v_{in}[i] - w[i])^2}$$

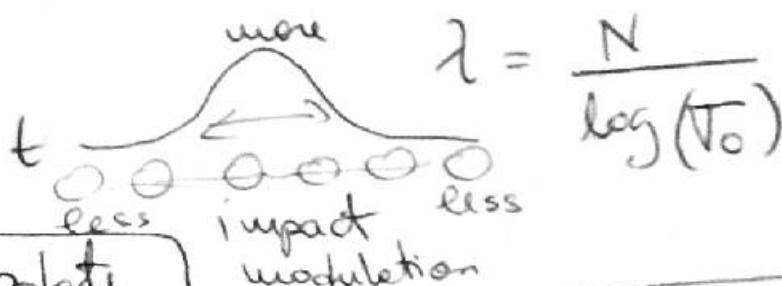
(B) The neighborhood function shrinks in time

$$\sigma(t) = \sigma_0 \cdot e^{-\left(\frac{t}{\lambda}\right)} \quad t = 1, 2, 3, \dots$$

σ_0 - initial value

λ - time constant.

$$\sigma_0 = \frac{\max(x_{size}, y_{size})}{2}$$



$$\lambda = \frac{N}{\log(\sigma_0)}$$

(C) Weight update

- each neuron in the winner's vicinity gets updated
- the closer a neuron is to the winner, the stronger the update

$$w(t+1) = w(t) + \eta(t) \cdot \phi(t) (v_{in}(t) - w(t))$$

$$\phi(t) = e^{-\frac{-dist^2(t)}{2\sigma^2(t)}}$$

(x_{BMU}, y_{BMU})

$$\eta(t) = \eta_0 \cdot e^{-\frac{(t-t_0)^2}{2\lambda^2}}$$

$\lambda = \frac{N}{\log(\sigma_0)}; \eta_0 = 0.5$

dist - topological distance between a neuron and the winner (BMU)

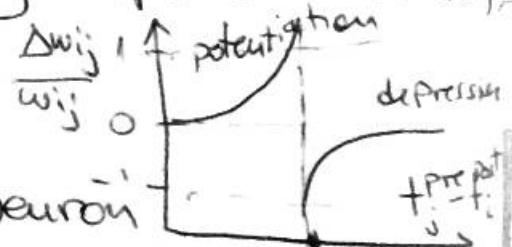
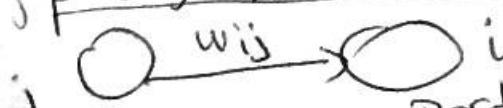
$$dist = \sqrt{(x_{BMU} - x_i)^2 + (y_{BMU} - y_i)^2}$$

Hopfield Networks (HN)

- * based on Hebbian learning principle;
- * capable of store and recover associations
- * pre-/and post-synaptic neurons where the activation patterns (time) is the core mechanism;
- * can be used as associative or content addressable memory.
- * uses time in the dynamics using feedback loops (recurrent network);
- * recognize learned patterns & can reconstruct corrupted patterns.

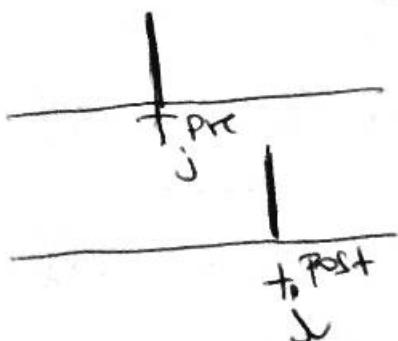
* STDP (Spike Timing Dependent Plasticity)

presynaptic neuron



post synaptic neuron

if pre
and post
spikes



if pre and post synaptic action potentials are correlated then there is an enhancement in synaptic strength.

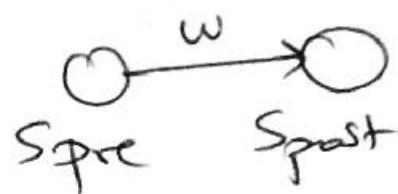
* Abstracting from Biology \rightarrow HN neuron

$$w = S_{post} \cdot S_{pre}$$

choose
in
weight

state
Post-synaptic
neuron

state
Pre-synaptic
neuron



for binary neurons
states {0, 1}

	Pre			Post		
Post	0	1		0	1	
0	0	0		0	x	x
1	0	1		1	x	1

Δw

not OK!

for bipolar neurons

states {-1, 1}

	Pre			Post		
Post	-1	1		-1	1	
-1	1	-1		-1	1	
1	-1	1		1	-1	1

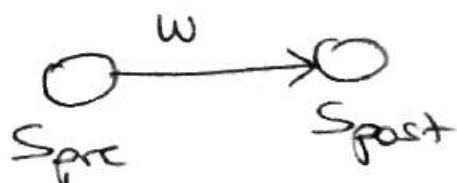
$\Rightarrow -1$

OK

This representation allows for both correlation & anticorrelation of the states.

But how can such a system learn patterns?

- consider only 2 neurons



Hebbian learning produces a weight w such that a learned input S_{input} maps onto the corresponding output S_{target} .

Algorithm

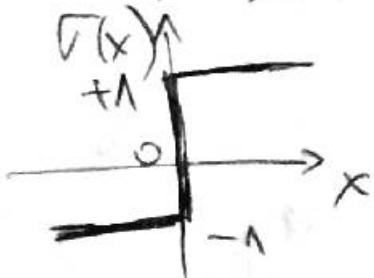


#1. Connect input S_{input} to S_{pre}
 $\circlearrowleft S_{\text{pre}} = S_{\text{input}}$

#2. Squashing/activation function Γ

Binary threshold/
Step function

$$\Gamma(x) = \begin{cases} +1, & x \geq 0 \\ -1, & x < 0 \end{cases}$$



#3. Weight calculation with Hebbian Learning

$$w = s_{\text{target}} \cdot S_{\text{input}}$$

Here we can learn only 1 pattern

correlated	$+1 \rightarrow +1$ & $-1 \rightarrow -1$
anticorrelated	$+1 \rightarrow -1$ & $-1 \rightarrow +1$

#4. substitute #3 in #2 and we get

$$S_{\text{post}} = \Gamma(s_{\text{target}} \cdot S_{\text{input}} \cdot S_{\text{pre}})$$

#5. But $S_{\text{pre}} = S_{\text{input}}$ (check #1)

$$S_{\text{post}} = \Gamma(s_{\text{target}} \cdot S_{\text{input}} \cdot S_{\text{input}})$$

$= 1$ because

$$(-1)(-1) = 1$$

$$(1)(1) = 1$$

(q.e.d)

$$S_{\text{post}} = \Gamma(s_{\text{target}})$$

Assume we want to learn multiple patterns and correct errors in the input?

* The learning rule can be extended to a network of neurons:

$$W = \overline{s_{\text{post}}} \times \overline{s_{\text{pre}}} = \begin{pmatrix} \text{matrix} \\ \vdots & \ddots \end{pmatrix}^T$$

\overline{s} - vector s

$$\overline{s_{\text{target}}} = \nabla (W \cdot \overline{s_{\text{inp}}})$$

* Such a network can remember a vector as input.

Let's assume we learn a vector and then we flip some bits and feed it to the network.

Let's learn $\overline{s_{\text{in}}} = (-1 -1 -1 1 -1 +1) \in \mathbb{R}_+^6$

$$\overline{s_{\text{target}}} = (-1 1 -1 -1 1) \in \mathbb{R}_+^5$$

A Hopfield network will converge to a learned pattern given the noisy (corrupted) pattern

Consider $\overline{s_x} = (-1 \boxed{1} -1 1 -1 \boxed{-1})$

2 differences from $\overline{s_{\text{in}}}$ (learned)

$$\overline{s_{\text{out}}} = \nabla (W \cdot \overline{s_x}) \rightarrow \nabla \left(\begin{pmatrix} 1 & 1 & 1 & -1 & 1 & -1 \\ -1 & -1 & -1 & 1 & -1 & 1 \\ 1 & 1 & 1 & -1 & 1 & -1 \\ 1 & 1 & 1 & -1 & 1 & -1 \\ -1 & -1 & -1 & 1 & -1 & 1 \\ -1 & -1 & -1 & 1 & -1 & 1 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \\ -1 \\ 1 \\ -1 \\ -1 \end{pmatrix} \right)$$

$$W = \overline{s_{\text{post}}} \times \overline{s_{\text{pre}}} = \overline{s_{\text{target}}} \otimes \overline{s_{\text{in}}}$$

where

$$\begin{cases} \overline{s_{\text{post}}} = \overline{s_{\text{target}}} \\ \overline{s_{\text{pre}}} = \overline{s_{\text{in}}} \end{cases}$$

Computed as $W = \overline{s_{\text{target}}}^T \otimes \overline{s_{\text{in}}}$

$$\Rightarrow \nabla(W \cdot \vec{s}_X) = \nabla \begin{pmatrix} -2 \\ +2 \\ -2 \\ -2 \\ +2 \end{pmatrix} = \begin{pmatrix} -1 \\ 1 \\ -1 \\ -1 \\ 1 \end{pmatrix} = \vec{s}_{\text{out}} = \vec{s}_{\text{target}}$$

If 4 out of 6 bits are correct the network finds the correct pattern!

How to learn more than a pattern?

- simply add weight matrices of all patterns \Rightarrow bit patterns will interfere
- the one pattern which is closest to the input is called attractor.
- use orthogonal vectors (dot prod. is zero) to avoid confusion (interference).
- the number of patterns that can be learned is 0.15 · N (N-no. of neurons) for reliable performance.

Hopfield network example pseudocode

Pattern recovery task

Steps:

- #1. Acquire training patterns and testing patterns (distorted)
- train-pat [n-pat][idx]
- test-pat [n-pat][idx]
- $n\text{-pat}$ - no of patterns
 idx - val of current pixel in patterns

- #2. Use bipolar representation $\{-1, 1\}$

Training phase:

- #3. for $pt = 1 : n\text{-pat}$
- #4. for $i = 1 : n\text{-neurons}$
- #5. for $j = 1 : n\text{-neurons}$
- #6. $W[i][j] += \text{train-pat}[pt][i] \cdot \text{train-pat}[pt][j]$
- \uparrow equivalent to
 $w_{ij}(+1) = w_{ij}(1) + \Delta w_{ij}$ where $\Delta w_{ij} = s_i \cdot s_j$
- $w_{jj} = w_{ii} = 0$
- \uparrow pre post
- #7. end for j
- #8. end for i
- #9. end for pt

Recovery phase:

- #10. for iter = 1 : max_iters
- #11. for ptd = 1 : n_pat
- #12. for $n_1 = 1 : n\text{-neurons}$
- #13. $out[n_1] = 0$
- #14. for $n_2 = 1 : n\text{-neurons}$
- #15. $out[n_1] += W[n_1][n_2] \cdot \text{test-pat}[ptd][n_2]$
- \uparrow were you can use a sign() function as activation
- #16. end for n_2
- #17. end for n_1
- #18. end for ptd

#19. if (out == one of the testing patterns) ← converging condition

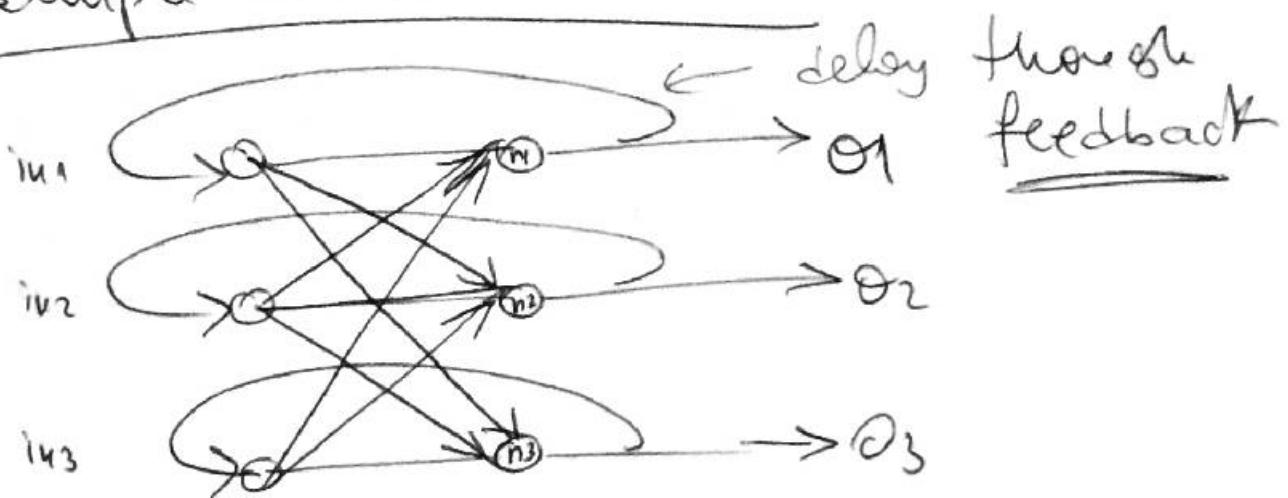
#20. break

#21. else testing_pat = out ; ← feed back the output

#22. end for iter

#23. print recovered pattern

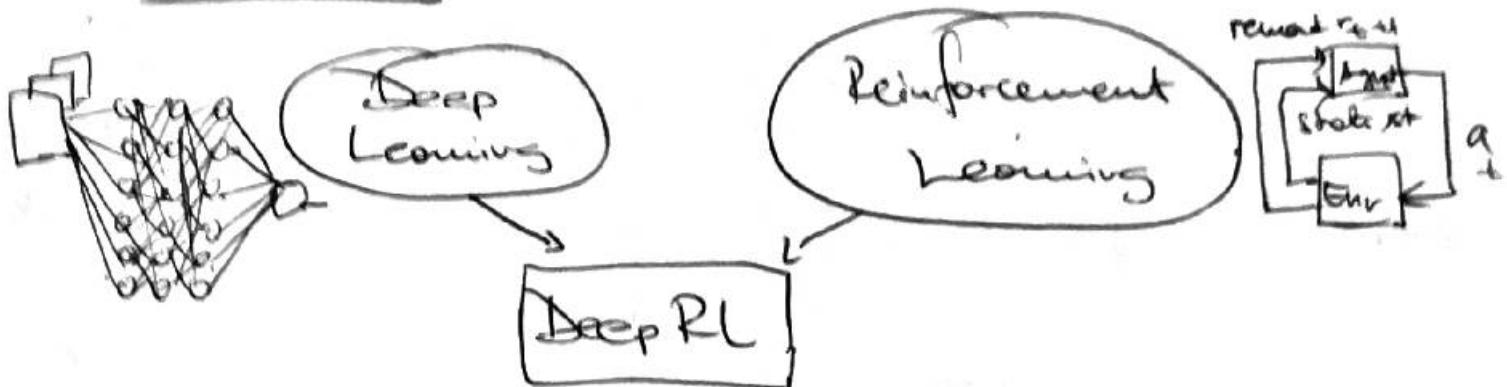
Sample network architecture



HTML Class Deep Reinforcement Learning

A practical overview

The Basics:



* use deep nets to represent value function ($V(s)$) or policy ($\pi(s)$) or the model.

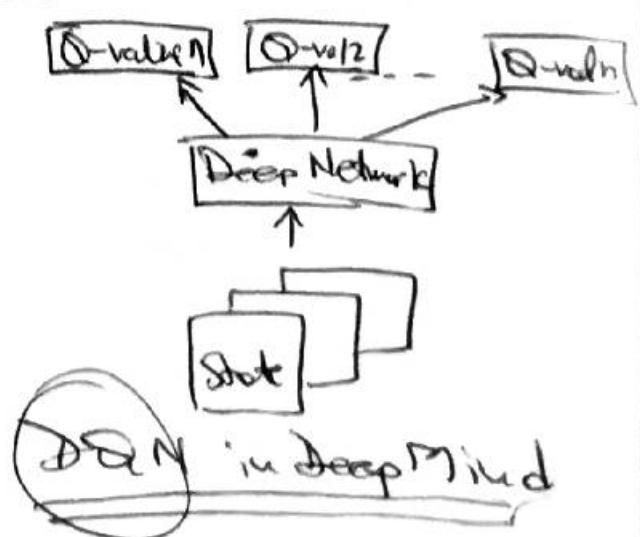
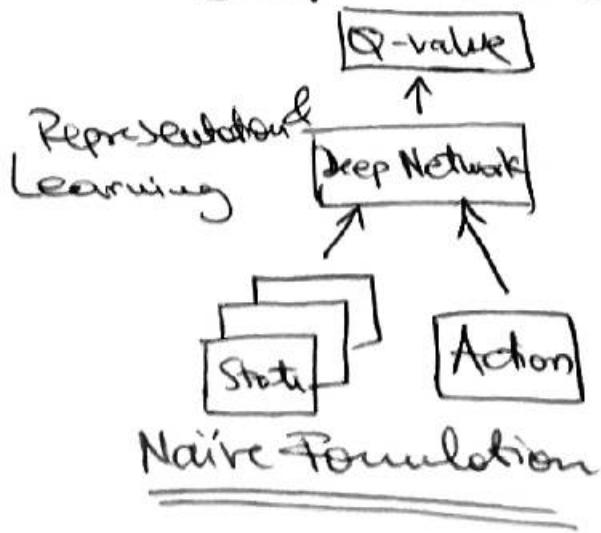
* use SGD to optimize $V(s)$, $\pi(s)$ or model end-to-end.

Applications:

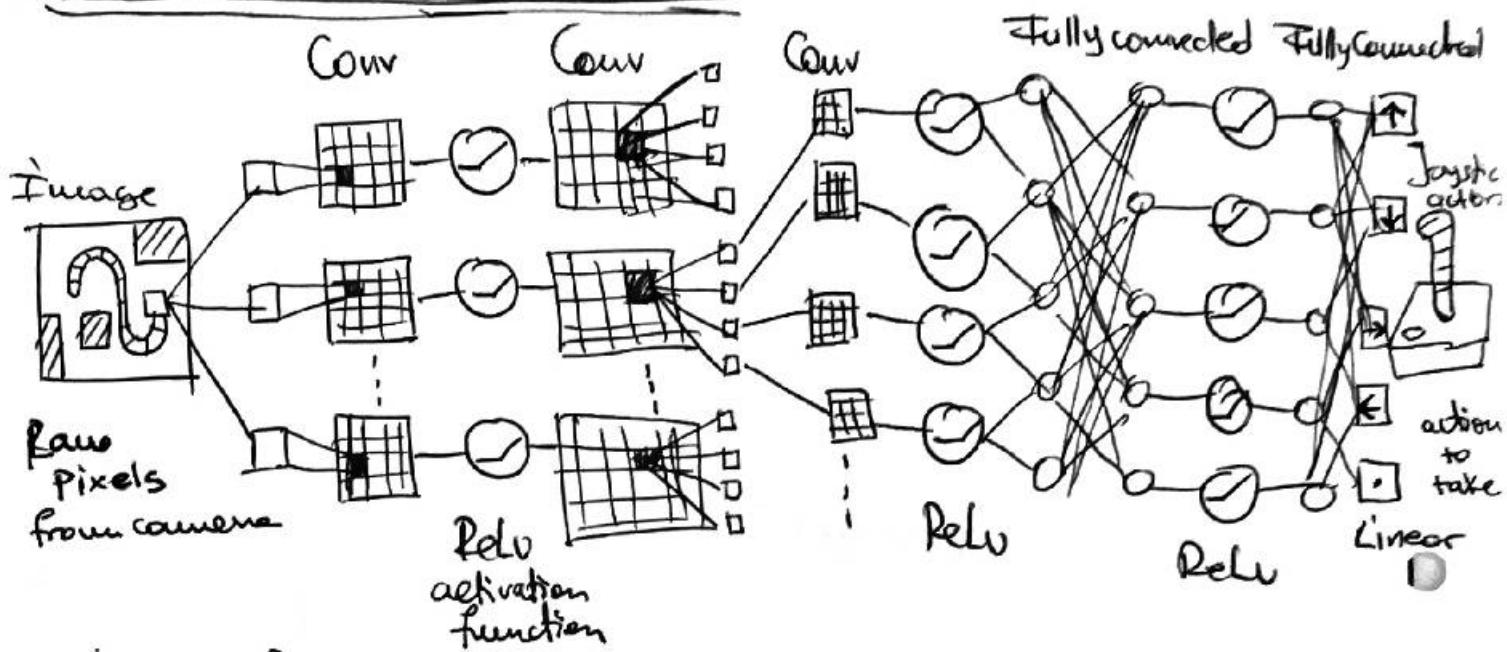
I Athen Game player (December 2013)

II "Human Level Control through deep reinforcement learning" \Rightarrow Deep Q-Network (DQN)

DeepMind \Rightarrow AlphaGo (2016)



DQN - architecture



Loss function definition

$$L = \frac{1}{2} [r(s) + \underbrace{\max_{a'} (Q(s', a'))}_{\text{target}} - \underbrace{Q(s, a)}_{\text{prediction}}]^2$$

Q table update algorithm :

- #1 Do a feedforward pass on current s to predict Q for all a.
- #2. Do a feedforward pass for next state s' and calculate $\max_{a'} Q(s', a')$.
- #3. Set $\hat{Q} = r + \gamma \max_{a'} Q(s', a')$.
- #4. Update weights using backprop.

Weights update for Value Iteration

- represent value function depending on weights w

$$Q(s, a, w) \approx Q''(s, a)$$
- define objective function by MSE (Mean Squared Error) in Q -values:

$$L(w) = E \left[\left(r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w) \right)^2 \right]$$

expectation operator (mean value) ↓ target
Learning rule:

$$\frac{\partial L(w)}{\partial w} = E \left[\left(r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w) \right) \cdot \frac{\partial Q(s, a, w)}{\partial w} \right]$$

Optimize using SGD on $\frac{\partial L(w)}{\partial w}$

Weights update using Policy Iteration

* same representation of value, depending on weights as in value iteration:

$$Q(s, a, w) \approx Q^{\pi}(s, a)$$

- Define objective function by MSE in Q-values:

$$L(w) = E \left[(r + \gamma \cdot Q(s', a', w) - Q(s, a, w))^2 \right]$$

expectation operator (mean value) \rightarrow target
 gradient update rule using SGD on

$$\frac{\partial L(w)}{\partial w} = E \left(r + \gamma Q(s', a', w) - Q(s, a, w) \right) \frac{\partial Q(s, a, w)}{\partial w}$$

Stability in Deep RL

- Oscillations & divergence in naive networks for Deep RL:

- policy changes with slight changes in Q-values - given the data distribution.

- Rewards and Q-values scales are not known \Rightarrow large unstable backpropagation updates.

DQN handles stability:

(A) use experience replay:

- break correlations in data.
- learn from past policies.

(B) freeze target Q-network:

- avoid oscillations.
- break connection between Q-network and target.

(C) provide stable gradients:

- clip rewards, typically $[-1, +1]$
- normalize network adaptively.

- (R) Build a dataset from agent's own experience
- take actions according to greedy policy π
 - store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory D .
 - sample minibatch of transitions (s, a, r, s') from D .
 - optimize MSE between Q-network and Q-learning targets:

$$L(w) = \mathbb{E}_{s, a, r, s' \sim D} \left[(r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w))^2 \right]$$

B.

- Fix parameters used in Q-learning target
- compute Q-learning targets w.r.t old fixed parameter w^-

$$r + \gamma \max_{a'} Q(s', a', w^-)$$

- optimize MSE between Q-network and Q-learning targets:

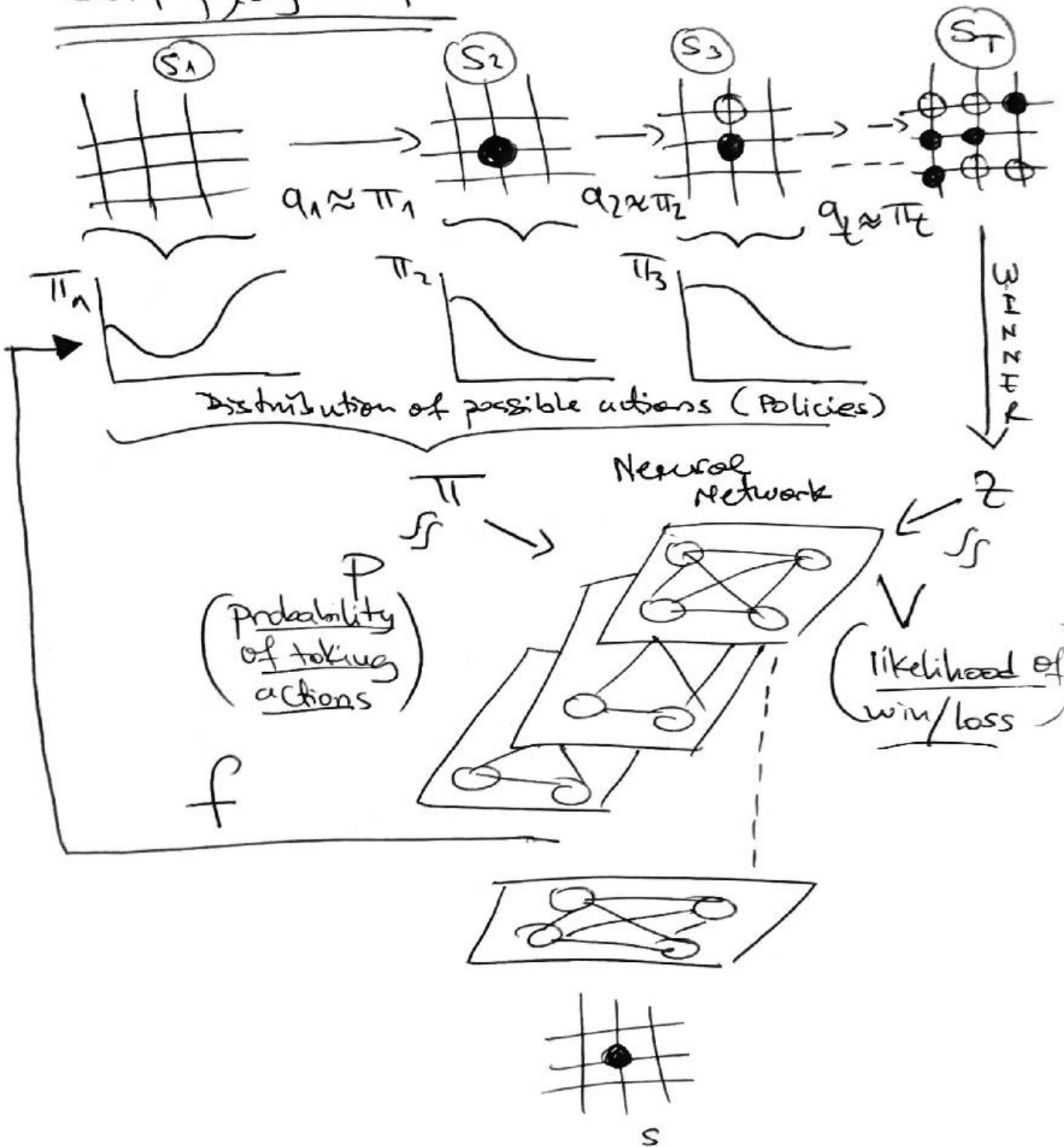
$$L(w) = \mathbb{E}_{s, a, r, s' \sim D} [(r + \gamma \max_{a'} Q(s', a', w)) - Q(s, a, w)]^2$$

- periodically update $w^- \leftarrow w$
-

Deep RL from "blank slate"
("tabula rasa")

- Alpha Go zero!
- Improve planning and evaluation with feedback from self-play (No, zero human game total!)

Self play Loop



Worked example

MLP-BP

x

$$\text{net}_0 = x_i \\ \text{Out}_0 = \text{net}_0$$

w_{00}

$$\text{net}_2 = w_{00} \cdot x + w_{01} \cdot y + w_{02} \cdot \text{bias}$$

$$\text{Out}_2 = \tanh(\text{net}_2)$$

w_{01}

w_{02}

w_{10}

w_{11}

w_{12}

$$\text{Net}_1 = y_i \\ \text{Out}_1 = \text{Net}_1$$

$$\text{net}_3 = w_{10} \cdot x + w_{11} \cdot y + w_{12} \cdot \text{bias};$$

$$\text{Out}_3 = \tanh(\text{net}_3)$$

w_{20}

w_{21}

w_{22}

$$\text{Net}_4 = w_{20} \cdot x + w_{21} \cdot y + w_{22} \cdot \text{bias};$$

$$\text{Out}_4 = \tanh(\text{net}_4)$$

w_{30}

w_{31}

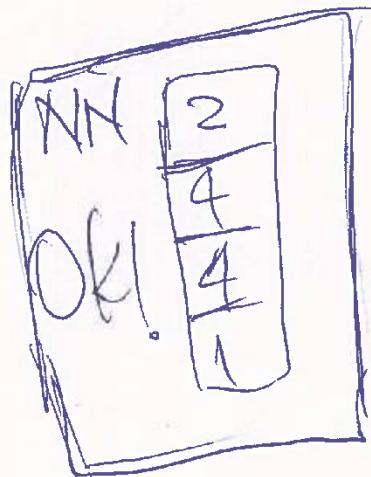
w_{32}

$$\text{Net}_5 = w_{30} \cdot x + w_{31} \cdot y + w_{32} \cdot \text{bias}$$

$$\text{Out}_5 = \tanh(\text{net}_5)$$

bias

h



NN

2

4

4

1

OK!

bias

i

h

bias

j

bias

i

bias

j

i

j

bias

i

bias

j

bias

$$\delta_{10} = \delta_{\text{out}} = 2 \cdot (\text{target} - \text{out}) \quad (6)$$

$$\delta_9 = \delta_{10} \cdot w_{003} \cdot (1 - \text{out}_9^2) \quad \downarrow$$

$$\delta_8 = \delta_{10} \cdot w_{002} \cdot (1 - \text{out}_8^2) \quad (1)$$

$$\delta_7 = \delta_{10} \cdot w_{001} \cdot (1 - \text{out}_7^2)$$

$$\delta_6 = \delta_{10} \cdot w_{000} \cdot (1 - \text{out}_6^2)$$

Backpropagation

$$\delta_5 = w_{h33} \cdot \delta_9 \cdot (1 - \text{out}_5^2) + w_{h23} \cdot \delta_8 \cdot (1 - \text{out}_5^2) + w_{h13} \cdot \delta_7 \cdot (1 - \text{out}_5^2) + w_{h03} \cdot \delta_6 \cdot (1 - \text{out}_5^2)$$

$$\delta_4 = w_{h32} \cdot \delta_9 \cdot (1 - \text{out}_4^2) + w_{h22} \cdot \delta_8 \cdot (1 - \text{out}_4^2) + w_{h12} \cdot \delta_7 \cdot (1 - \text{out}_4^2) + w_{h02} \cdot \delta_6 \cdot (1 - \text{out}_4^2)$$

$$\delta_3 = w_{h31} \cdot \delta_9 \cdot (1 - \text{out}_3^2) + w_{h21} \cdot \delta_8 \cdot (1 - \text{out}_3^2) + w_{h11} \cdot \delta_7 \cdot (1 - \text{out}_3^2) + w_{h01} \cdot \delta_6 \cdot (1 - \text{out}_3^2)$$

$$\delta_2 = w_{h30} \cdot \delta_9 \cdot (1 - \text{out}_2^2) + w_{h20} \cdot \delta_8 \cdot (1 - \text{out}_2^2) + w_{h10} \cdot \delta_7 \cdot (1 - \text{out}_2^2) + w_{h00} \cdot \delta_6 \cdot (1 - \text{out}_2^2) \quad (II)$$

$$\delta_1 = w_{34} \cdot \delta_5 + w_{21} \cdot \delta_4 + w_{11} \cdot \delta_3 + w_{01} \cdot \delta_2 \quad (III)$$

$$\delta_0 = w_{30} \cdot \delta_5 + w_{20} \cdot \delta_4 + w_{10} \cdot \delta_3 + w_{00} \cdot \delta_2$$

Update weights

Output:

$$\Delta w_{g04} = \eta \cdot \delta_{10} \cdot 1;$$

$$\Delta w_{003} = \eta \cdot \delta_{10} \cdot \text{out}_9;$$

$$\Delta w_{002} = \eta \cdot \delta_{10} \cdot \text{out}_8;$$

$$\Delta w_{001} = \eta \cdot \delta_{10} \cdot \text{out}_7;$$

$$\Delta w_{000} = \eta \cdot \delta_{10} \cdot \text{out}_6;$$

Hidden

$$\Delta w_{h33} = \eta \cdot \sqrt{g} \cdot \text{Out}_5$$

$$\Delta w_{h23} = \eta \cdot \sqrt{g} \cdot \text{Out}_5$$

$$\Delta w_{h13} = \eta \cdot \sqrt{g} \cdot \text{Out}_5$$

$$\Delta w_{h03} = \eta \cdot \sqrt{g} \cdot \text{Out}_5$$

$$\Delta w_{h32} = \eta \cdot \sqrt{g} \cdot \text{Out}_4$$

$$\Delta w_{h22} = \eta \cdot \sqrt{g} \cdot \text{Out}_4$$

$$\Delta w_{h12} = \eta \cdot \sqrt{g} \cdot \text{Out}_4$$

$$\Delta w_{h02} = \eta \cdot \sqrt{g} \cdot \text{Out}_4$$

$$\Delta w_{h31} = \eta \cdot \sqrt{g} \cdot \text{Out}_3$$

$$\Delta w_{h21} = \eta \cdot \sqrt{g} \cdot \text{Out}_3$$

$$\Delta w_{h11} = \eta \cdot \sqrt{g} \cdot \text{Out}_3$$

$$\Delta w_{h01} = \eta \cdot \sqrt{g} \cdot \text{Out}_3$$

$$\Delta w_{h30} = \eta \cdot \sqrt{g} \cdot \text{Out}_2$$

$$\Delta w_{h20} = \eta \cdot \sqrt{g} \cdot \text{Out}_2$$

$$\Delta w_{h10} = \eta \cdot \sqrt{g} \cdot \text{Out}_2$$

$$\Delta w_{h00} = \eta \cdot \sqrt{g} \cdot \text{Out}_2$$

Input

$$\Delta w_{i31} = \eta \cdot \sqrt{g} \cdot y$$

$$\Delta w_{i21} = \eta \cdot \sqrt{g} \cdot y$$

$$\Delta w_{i11} = \eta \cdot \sqrt{g} \cdot y$$

$$\Delta w_{i01} = \eta \cdot \sqrt{g} \cdot y$$

$$\Delta w_{i30} = \eta \cdot \sqrt{g} \cdot x$$

$$\Delta w_{i20} = \eta \cdot \sqrt{g} \cdot x$$

$$\Delta w_{i10} = \eta \cdot \sqrt{g} \cdot x$$

$$\Delta w_{i00} = \eta \cdot \sqrt{g} \cdot x$$

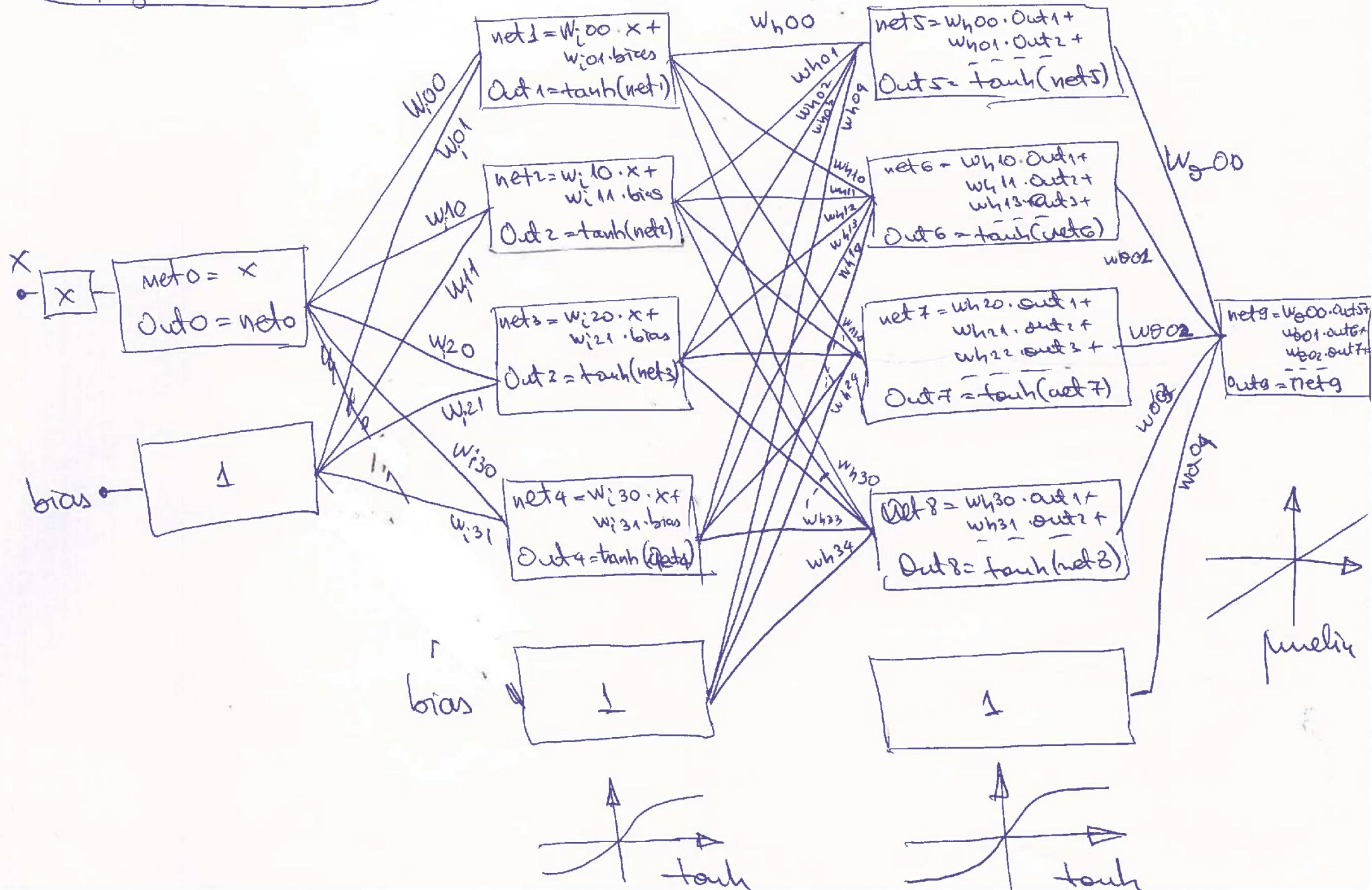
$$W_{ij} + = \Delta w_{ij}, i$$

- update

weights

+ bias

Worked example Regression MLP



Backpropagation

$$\delta_g = \text{Jout} = 2 \cdot (\text{target} - \text{out})$$

$$\delta_8 = \delta_g \cdot w_{g03} \cdot (1 - \text{Out}_8^2)$$

$$\delta_7 = \delta_g \cdot w_{g02} \cdot (1 - \text{Out}_7^2)$$

$$\delta_6 = \delta_g \cdot w_{g01} \cdot (1 - \text{Out}_6^2)$$

$$\delta_5 = \delta_g \cdot w_{g00} \cdot (1 - \text{Out}_5^2)$$

①

$$\delta_4 = w_{h03} \cdot \delta_5 \cdot (1 - \text{Out}_4^2) + w_{h13} \cdot \delta_6 \cdot (1 - \text{Out}_4^2) + w_{h23} \cdot \delta_7 \cdot (1 - \text{Out}_4^2) + w_{h33} \cdot \delta_8 \cdot (1 - \text{Out}_4^2)$$

$$\delta_3 = w_{h02} \cdot \delta_5 \cdot (1 - \text{Out}_3^2) + w_{h12} \cdot \delta_6 \cdot (1 - \text{Out}_3^2) + w_{h22} \cdot \delta_7 \cdot (1 - \text{Out}_3^2) + w_{h32} \cdot \delta_8 \cdot (1 - \text{Out}_3^2)$$

$$\delta_2 = w_{h01} \cdot \delta_5 \cdot (1 - \text{Out}_2^2) + w_{h11} \cdot \delta_6 \cdot (1 - \text{Out}_2^2) + w_{h21} \cdot \delta_7 \cdot (1 - \text{Out}_2^2) + w_{h31} \cdot \delta_8 \cdot (1 - \text{Out}_2^2)$$

$$\delta_1 = w_{h00} \cdot \delta_5 \cdot (1 - \text{Out}_1^2) + w_{h10} \cdot \delta_6 \cdot (1 - \text{Out}_1^2) + w_{h20} \cdot \delta_7 \cdot (1 - \text{Out}_1^2) + w_{h30} \cdot \delta_8 \cdot (1 - \text{Out}_1^2)$$

②

$$\delta_0 = w_{i00} \cdot \delta_1 + w_{i10} \cdot \delta_2 + w_{i20} \cdot \delta_3 + w_{i30} \cdot \delta_4$$

Weight update

Output:

$$\Delta w_{g04} = \eta \cdot \delta_g \cdot \text{bias};$$

$$\Delta w_{g03} = \eta \cdot \delta_g \cdot \text{Out}_8;$$

$$\Delta w_{g02} = \eta \cdot \delta_g \cdot \text{Out}_7;$$

$$\Delta w_{g01} = \eta \cdot \delta_g \cdot \text{Out}_6;$$

$$\Delta w_{g00} = \eta \cdot \delta_g \cdot \text{Out}_5;$$

Hidden:

$$\Delta W_{h34} = \eta \cdot (\sqrt{5}) \cdot \cancel{\text{Out}_1 \times \cancel{\text{Out}_2 \times \cancel{\text{Out}_3 \times \cancel{\text{Out}_4}}} \cdot 1}$$

$$\Delta W_{h24} = \eta \cdot (\sqrt{7}) \cdot \text{bias}$$

$$\Delta W_{h14} = \eta \cdot (\sqrt{6}) \cdot \text{bias}$$

$$\Delta W_{h04} = \eta \cdot (\sqrt{5}) \cdot \text{bias}$$

$$\Delta W_{h33} = \eta \cdot \sqrt{8} \cdot \text{Out}_4$$

$$\Delta W_{h23} = \eta \cdot \sqrt{7} \cdot \text{Out}_4$$

$$\Delta W_{h13} = \eta \cdot \sqrt{6} \cdot \text{Out}_4$$

$$\Delta W_{h03} = \eta \cdot \sqrt{5} \cdot \text{Out}_4$$

$$\Delta W_{h32} = \eta \cdot \sqrt{8} \cdot \text{Out}_3$$

$$\Delta W_{h22} = \eta \cdot \sqrt{7} \cdot \text{Out}_3$$

$$\Delta W_{h12} = \eta \cdot \sqrt{6} \cdot \text{Out}_3$$

$$\Delta W_{h02} = \eta \cdot \sqrt{5} \cdot \text{Out}_3$$

$$\Delta W_{h31} = \eta \cdot \sqrt{8} \cdot \text{Out}_2$$

$$\Delta W_{h21} = \eta \cdot \sqrt{7} \cdot \text{Out}_2$$

$$\Delta W_{h11} = \eta \cdot \sqrt{6} \cdot \text{Out}_2$$

$$\Delta W_{h01} = \eta \cdot \sqrt{5} \cdot \text{Out}_2$$

$$\Delta W_{h30} = \eta \cdot \sqrt{8} \cdot \text{Out}_1$$

$$\Delta W_{h20} = \eta \cdot \sqrt{7} \cdot \text{Out}_1$$

$$\Delta W_{h10} = \eta \cdot \sqrt{6} \cdot \text{Out}_1$$

$$\Delta W_{h00} = \eta \cdot \sqrt{5} \cdot \text{Out}_1$$

Input: $\Delta W_{i30} = \eta \cdot \sqrt{4} \cdot x$

$$\Delta W_{i20} = \eta \cdot \sqrt{3} \cdot x$$

$$\Delta W_{i10} = \eta \cdot \sqrt{2} \cdot x$$

$$\Delta W_{i00} = \eta \cdot \sqrt{1} \cdot x$$

Update

$$w_{ij} += \Delta w_{ij}$$

* ARMA* models for Time Series prediction or

Time Series Forecasting using Stochastic Models

A time series is a series of data points indexed in time order.

||→ Time series analysis: methods & tools to analyze data in order to extract statistics & other characteristics of data;

||→ Time series forecasting: using a model to predict future values based on previous values.

- We previously learnt about MLPs which were able to learn functions out of (input, output) pairs using supervised learning. The functions could be arbitrarily non-linear.
- The simpler approach is to use linear models, which for forecasting offer good results at low computational cost, simplicity and with ease of understanding.

We will now look at linear and non-linear Stochastic time series models.

• Auto regressive Moving Average (ARMA) Models

- is a combination of AR(p) and MA(q) models
- in an AR(p) model the future value of variable is assumed to be a linear combination of p past observations and a random error:

$$Y_t = C + \sum_{i=1}^p \varphi_i Y_{t-i} + \epsilon_t$$

P - order of the model

$$Y_t = C + \varphi_1 Y_{t-1} + \varphi_2 Y_{t-2} + \dots + \varphi_p Y_{t-p} + \epsilon_t$$

where Y_t - actual value, ϵ_t - random error at time t

φ_i - model parameters

A MA(q) model uses past q errors as variables to predict y_t :

$$Y_t = \mu + \sum_{j=1}^q \theta_j e_{t-j}$$

$$Y_t = \mu + \theta_1 e_{t-1} + \theta_2 e_{t-2} + \dots + \theta_q e_{t-q} + \epsilon_t$$

Where μ - mean of the series, θ_j - model parameters
q - order of the model

$\epsilon_t \sim N(0, \sigma^2)$ random, normally distributed noise

The ARMA(p, q) model can be mathematically described as:

$$Y_t = C + \epsilon_t + \sum_{i=1}^p \varphi_i Y_{t-i} + \sum_{j=1}^q \theta_j e_{t-j}$$

Usually the models φ_i and θ_j are polynomials

$$\varphi(L) = 1 - \sum_{i=1}^p \varphi_i L^i \text{ and } \theta(L) = 1 + \sum_{j=1}^q \theta_j L^j$$

$L Y_t = Y_{t-1}$

where L is the lag operator.

- Nonlinear models : NAR (Nonlinear Auto Regressive) and NARMA (Nonlinear Autoregressive Moving Average)

- A simple extension to nonlinear systems can be done as a NAR model, generically defined as:

$$y_t = h(y_{t-1}, y_{t-2}, y_{t-3}, \dots, y_{t-p}) + \epsilon_t$$

where $h(\cdot)$ is a smooth unknown function with the assumption that the best (i.e. min mean square error) prediction of y_t given previous $y_{t-1}, y_{t-2}, \dots, y_{t-p}$ is its conditional mean:

$\hat{y}_t = \mathbb{E}(y_t | y_{t-1}, \dots, y_{t-p})$ ← expectation operator

$$\hat{y}_t = \mathbb{E}(y_t | y_{t-1}, \dots, y_{t-p}) = h(y_{t-1}, \dots, y_{t-p})$$

- The feed forward network able to implement a NAR model is a nonlinear approximation of h such as: activation function

$$\hat{y}_t = h(y_{t-1}, \dots, y_{t-p}) = \sum_{i=1}^I w_i \cdot f_i \left(\sum_{j=1}^P w_{ij} y_{t-j} \right)$$

- NARMA generalize ARMA, with h unknown

$$y_t = h(y_{t-1}, y_{t-2}, \dots, y_{t-p}, \epsilon_{t-1}, \dots, \epsilon_{t-q})$$

- Then we have a recurrent neural network to approximate it

$$\hat{y}_t = h(y_{t-1}, \dots, y_{t-p}) = \sum_{i=1}^I w_i \cdot f_i \left(\underbrace{\sum_{j=1}^P w_{ij} y_{t-j}}_{\text{AR}} + \underbrace{\sum_{j=1}^q w_{ij} (y_{t-j} - \hat{y}_{t-j})}_{\text{MA}} \right) + \epsilon_t$$

• Auto regressive Integrated Moving Average (ARIMA)

- ARMA models are limited to stationary time series.
- in applications such as socio-economic systems and business the phenomena are non-stationary.
- ARIMA is a generalization of ARMA.

ARIMA (p, d, q) is described mathematically as:

$$\varphi(L) \cdot (1-L)^d y_t = \Theta(L) \varepsilon_t \quad \rightarrow \text{finite differencing level}$$

p - order of auto regression

q - order of moving average

d - order of the integrated model

d - controls differencing $\rightarrow d=1$ typical

$d=0 \Rightarrow$ ARMA model

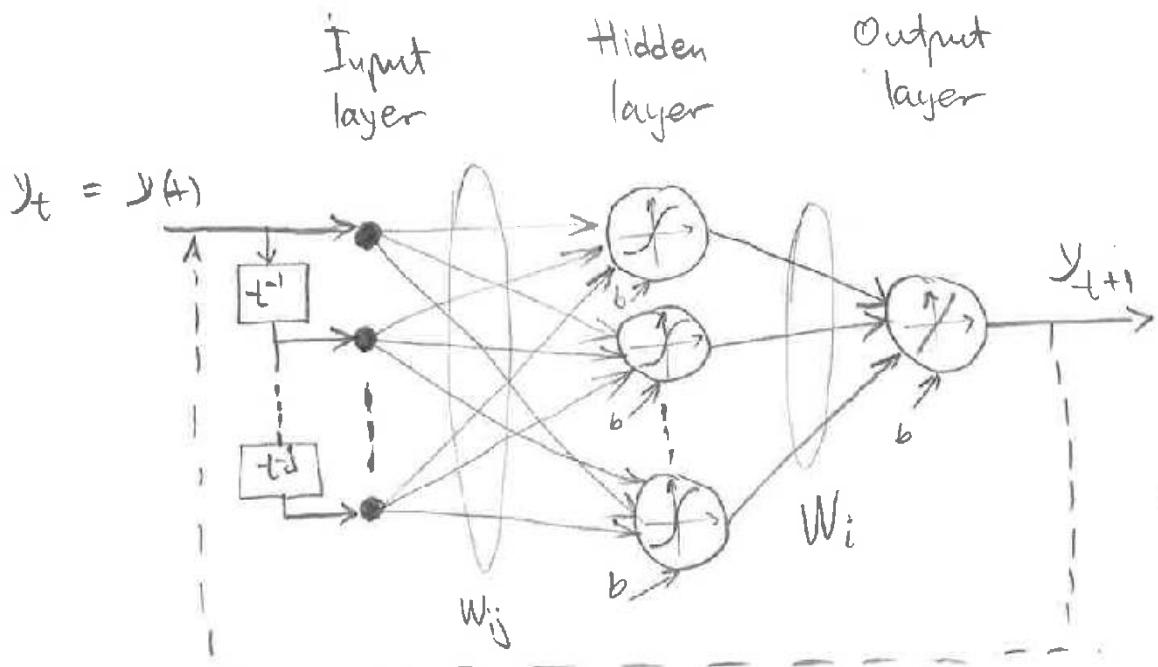
$$y_t = \frac{\Theta(L) \cdot \varepsilon_t}{\varphi(L) \cdot (1-L)^d}$$

where

$$\left\{ \begin{array}{l} \varphi(L) = 1 - \sum_{i=1}^p \varphi_i L^i \\ \Theta(L) = 1 + \sum_{j=1}^q \Theta_j L^j \end{array} \right. \quad \text{where } (L) \text{ is the lag operator (delay)}$$

which has $y_t = L y_{t+1}$
the property of

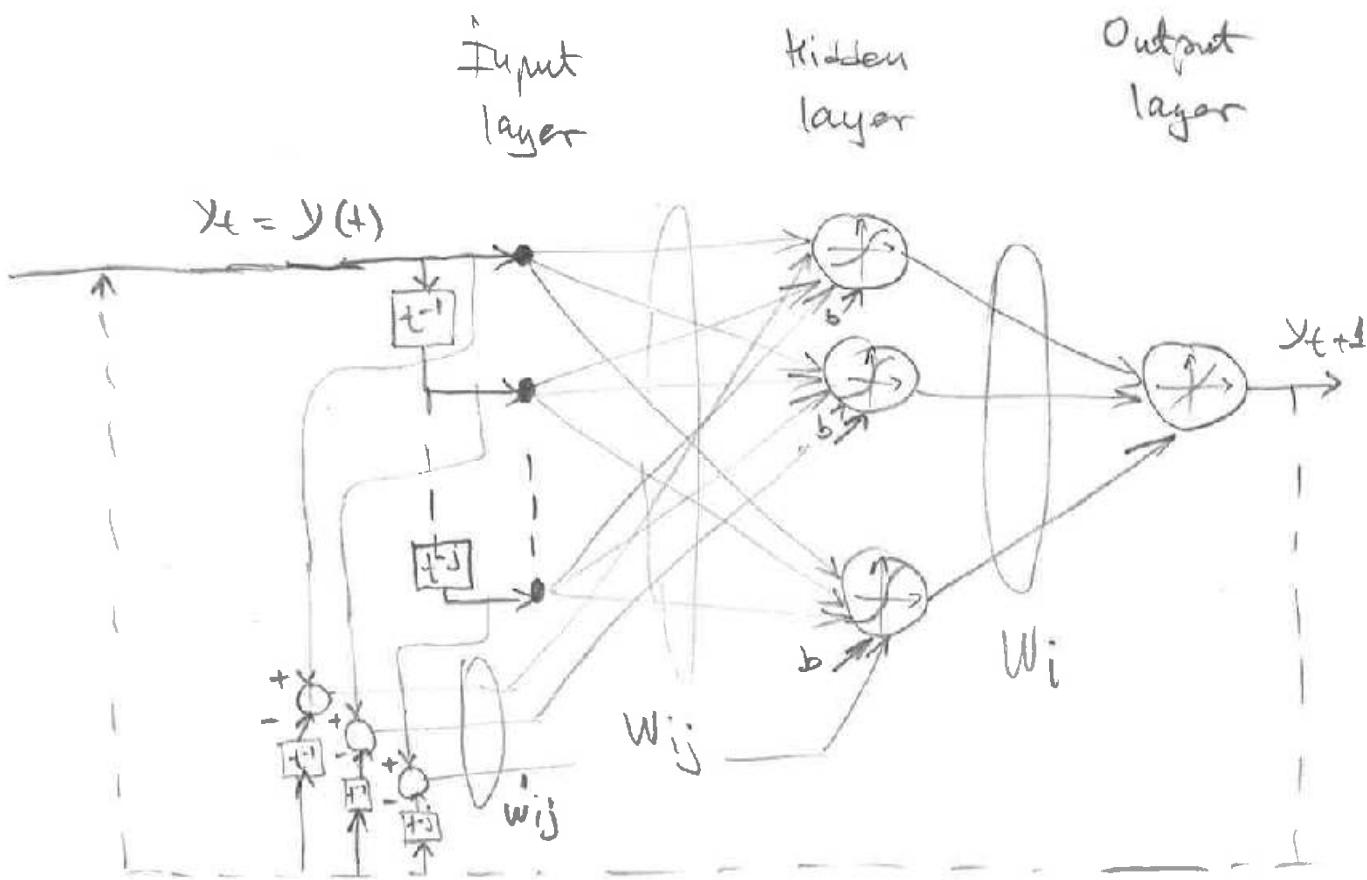
NAR network architecture example



Feedback is open for training & closed for prediction.

$$\hat{y}_t = h(y_{t-1}, \dots, y_{t-p}) = \sum_{i=1}^i w_i \cdot f\left(\sum_{j=1}^p w_{ij} y_{t-j}\right)$$

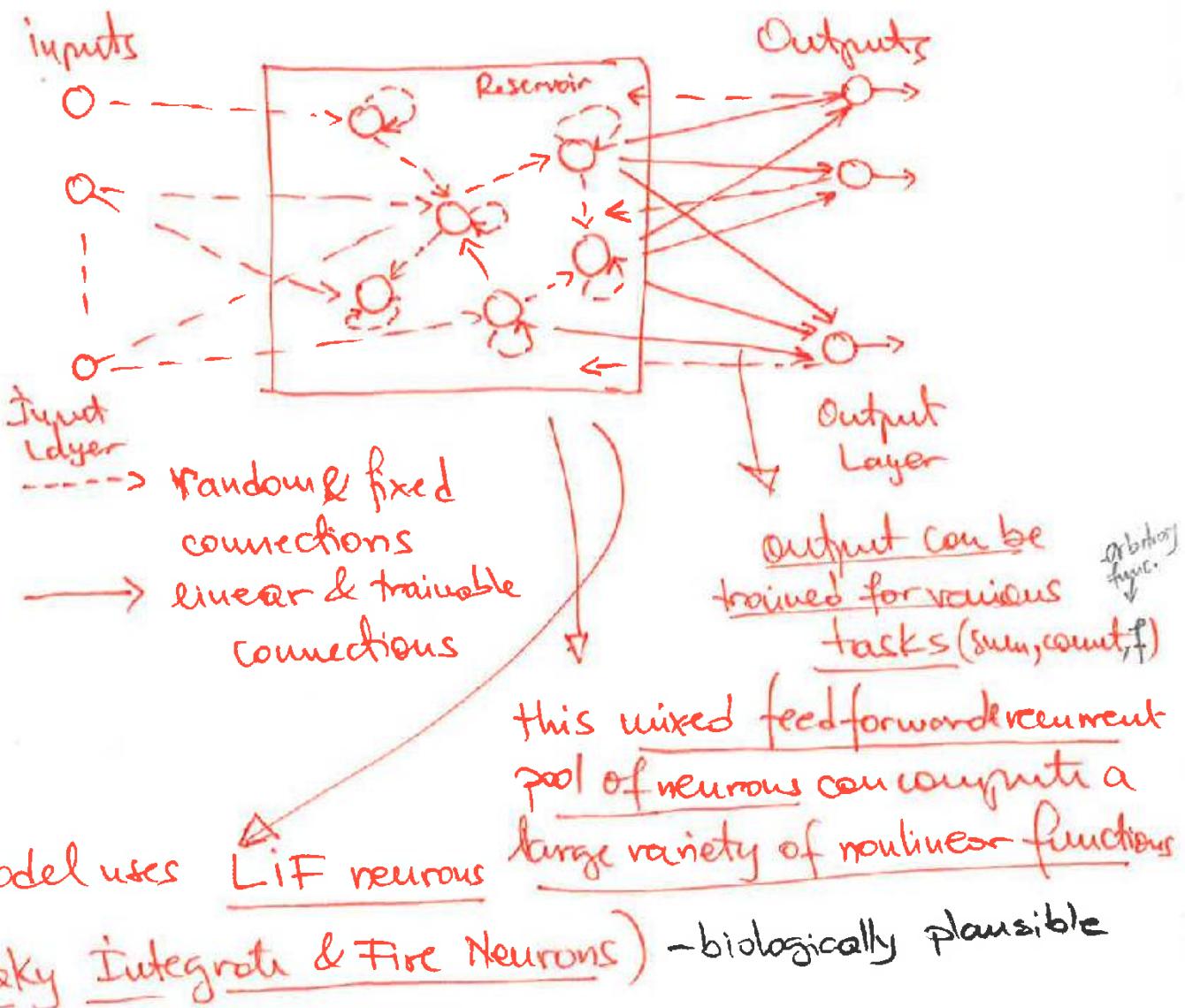
NARMA network architecture



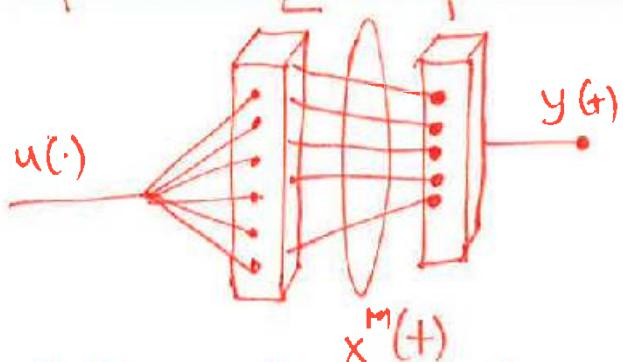
Liquid State Machines (LSM_s)

- models computations in cortical microcircuits of the brain; Wolfgang Maass developed it.
- model of real-time computations on continuous streams of data; (Reservoir Computing)
- use a high-dimensional dynamical system and have the inputs continuously perturbing it!
 - project inputs to a higher dimensional space.
 - internal dynamics of a recurrent spiking neural network.
 - the liquid itself doesn't generate any output it is just a "reservoir" for the inputs.
 - there is a readout function to compute the output of the LSM.
 - Given a time series input the LSM can produce a time series, and to obtain the desired values one needs to adjust the weights on the connections between reservoir and output.

The generic model is depicted in the following figure:



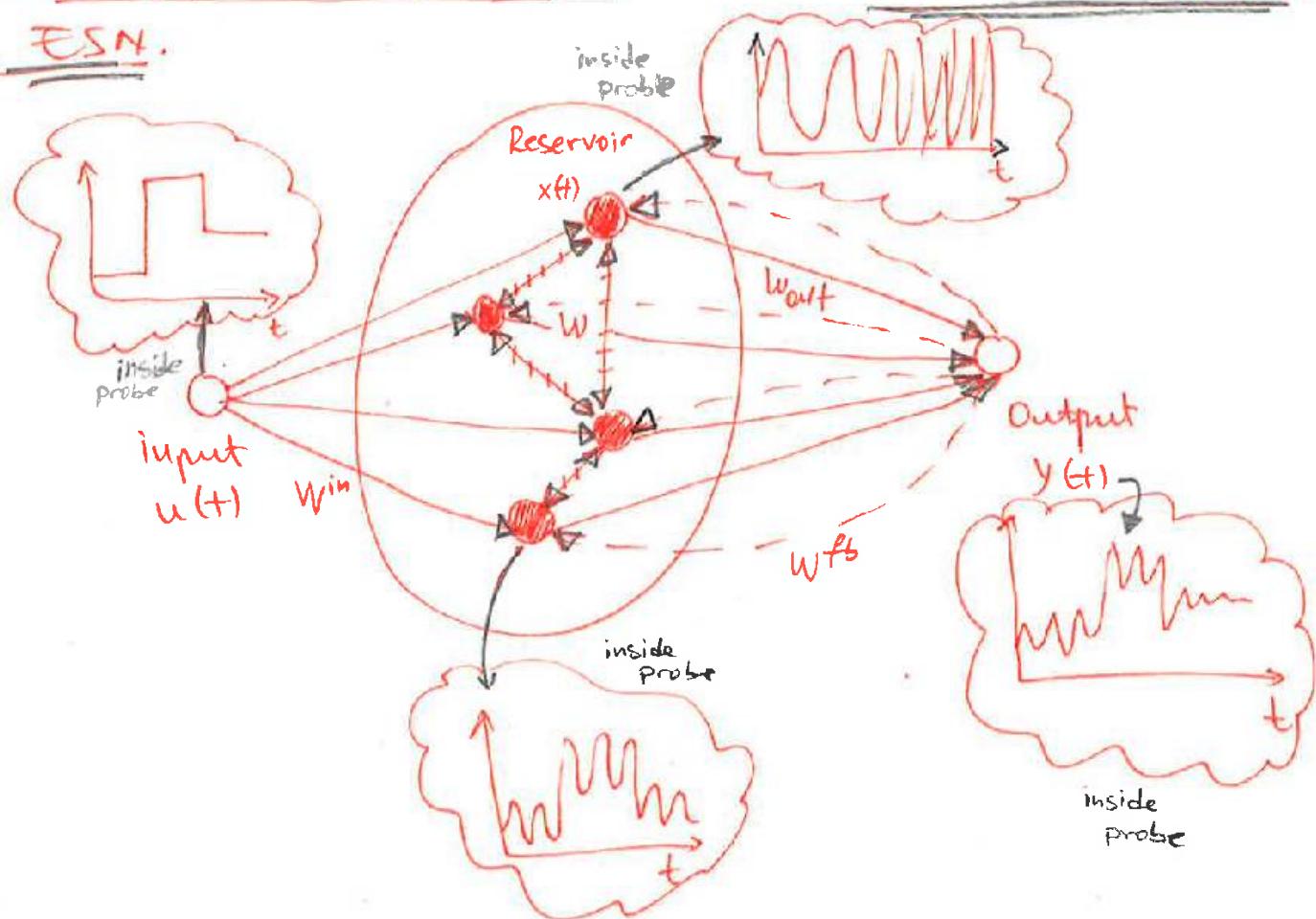
Formal definition



- A filter maps between two functions of time $u(\cdot) \rightarrow y(\cdot)$
- A LSM, $M = \langle L^M, f^M \rangle$, L^M = filter, f^M = readout
- State at time t , $x^M(t) = (L^M \cdot u)(t)$
- Output at time t , $y(t) = (Mu)(t) = (f^M)(x^M(t))$

In practical applications LIF (Leaky Integration and Fire neuron) dynamics are too complex and not necessary.

Another Reservoir Computing technique sharing the basic idea of LSTM is the Echo State Network or ESN.



ESN implementation

$u(t)$ - input signal

W_{in} - inputs weight vector (random, fixed)

W - reservoir recurrent weight vector (random, fixed)

W_{out} - output weight vector (trainable)

W^{fb} - feedback (optional) weight vector (fixed)

$x(t)$ - state vector of the network

$f(\cdot)$ - tanh activation function (for example)

$$x(t+1) = f(W_{in}^{in}u(t) + W \cdot x(t) + W^{fb} \cdot y(t))$$

Artificial Intelligence & Machine Learning

Class 1

Introduction to AIML

What is AI?

According to wikipedia: the ability of a computer to learn from specific data or experimental observation.

In general : AI describes a set of nature-inspired computational methodologies & approaches to address complex problems which are hard to solve or model mathematically or through traditional methods.

But what is Intelligence?

60's

- Chess playing as ultimate challenge now even a phone CPU using searching & sorting can do the task

2015

- A PC using a programmed ANN beats GO world master 4/6 games

But why is it artificial?

Transferring principles of computation & data processing from nature to computer systems. But also putting math at work in statistical learning.

Intelligence can be viewed as a ^{superposition} sequence of steps applied to an input, by a system, to compute an output. → See robotics for example!

In computational theory: A Turing Machine = a mathematical model of a hypothetical computing machine that uses a set of rules to determine a result from a set of input variables. \Rightarrow The Turing Test

Focus of the class:

\Rightarrow Sorting & searching in graphs (Traditional computation)

- * Strategies for sorting
- * pros & cons at scale
-
- * graph traversals / graph processing
- * dynamic programming

\Rightarrow Neural networks (Supervised learning)

- * from biology to technical systems
- * learning in artificial neurons
- * from single neurons to neural networks
- * learning in neural networks
backpropagation
- * supervised learning tips & tricks

\Rightarrow Unsupervised learning

- * introduction to unsupervised learning
- * learning structure from data:
SOM, VD_i, RBF
- * extract underlying clusters:
VD_i
- * reconstruct distorted patterns
after learning them: Hopfield

Tasks:

- * regression
- * classification

||=> Deep Neural Learning

- * Fundamentals of deep learning
- * Common architectures for deep learning
- * Building blocks of deep networks
- * Major deep network architectures

||=> Technical implementations

- * Recurrent neural networks LSTM, ESN, LST_m
- * Time series prediction (ARMA, NARMA)
- * SVM

||=> Reinforcement Learning

- * Introduction to RL, how is it different from supervised & unsupervised learning
- * Q - learning as a practical way to instantiate RL

||=> Evolutionary Programming

- * Introduction to evolutionary programming
- * Genetic algorithms

||=> Fuzzy Inference Systems

- * Introduction to fuzzy logic
- * Fuzzy control systems

⇒ Online distributed streaming machine learning

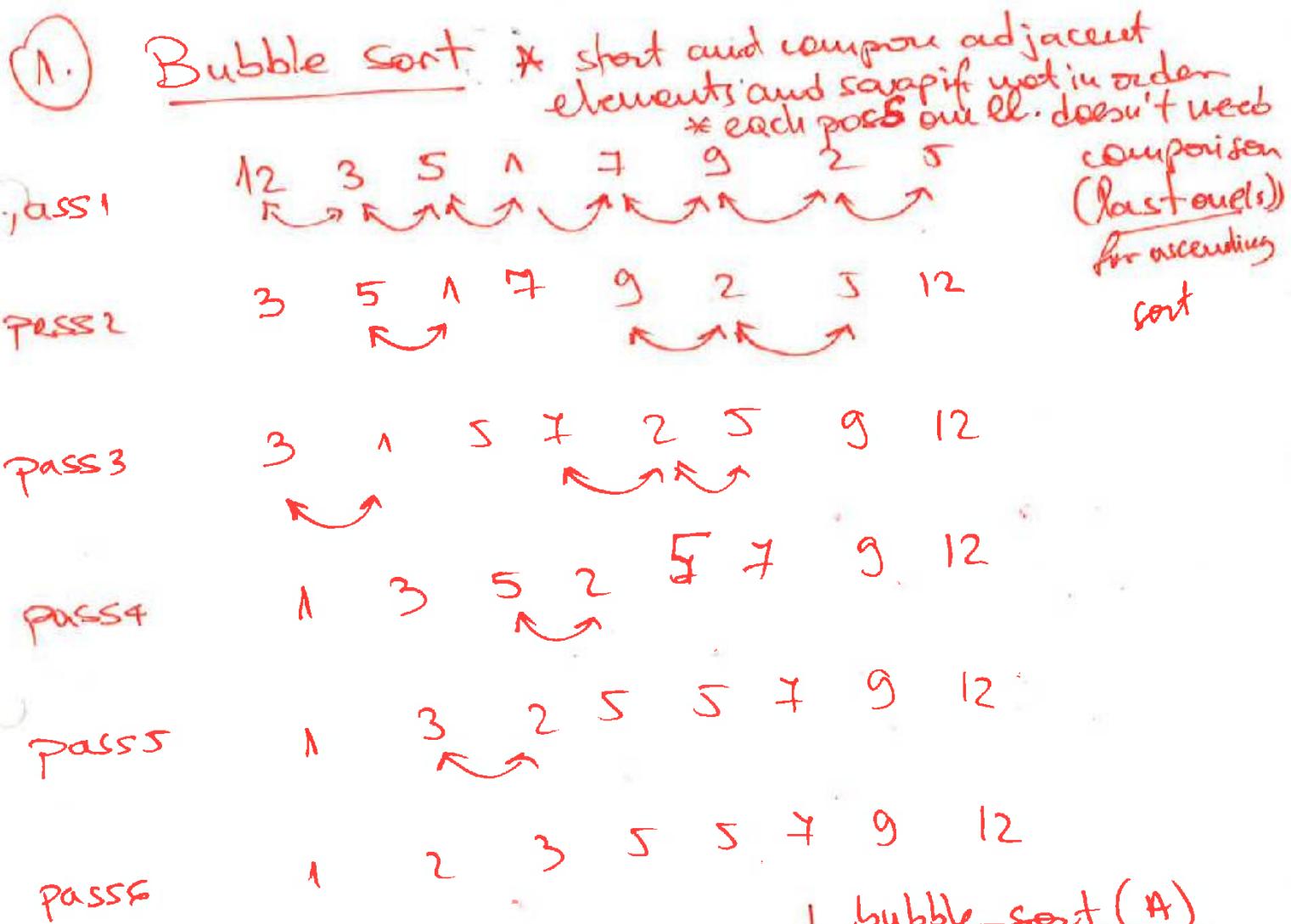
- * Big Data Analytics
- * Stream processors & building online ML on top
- * ML for Big Data on streams
 - | → VHTT Classifiers
 - | → AMR Regressors
 - | → Perceptrons

⇒ Spiking Neural Computation

- computing with biologically plausible neurons

Traditional Computation

Sorting : organizing elements in a sequence in a certain order \Rightarrow metric is complexity, the Big O notation for time execution \Rightarrow efficient execution time & resource-wise



bubble-sort(A)

- \rightarrow 15 moves
- \rightarrow 42 compares

- * Worst case runtime $\Theta(n^2)$
- * overhead for large, nearly-sorted data

```

do
    swapped = false
    for i = 1 to len(A)
        if (A(i-1) > A(i)) then
            swap(A(i-1), A(i))
            swapped = true
        endif
    endfor
    while (swapped)
  
```

② Insertion sort

`insertion_sort(A)`

`for i = 1 to len(A)`

`keyval = A[i]`

`j = i - 1`

`while (A[j] > keyval && j ≥ 0)`

`A[j + 1] = A[j]`

`j = j - 1`

`end while`

`A[j + 1] = keyval`

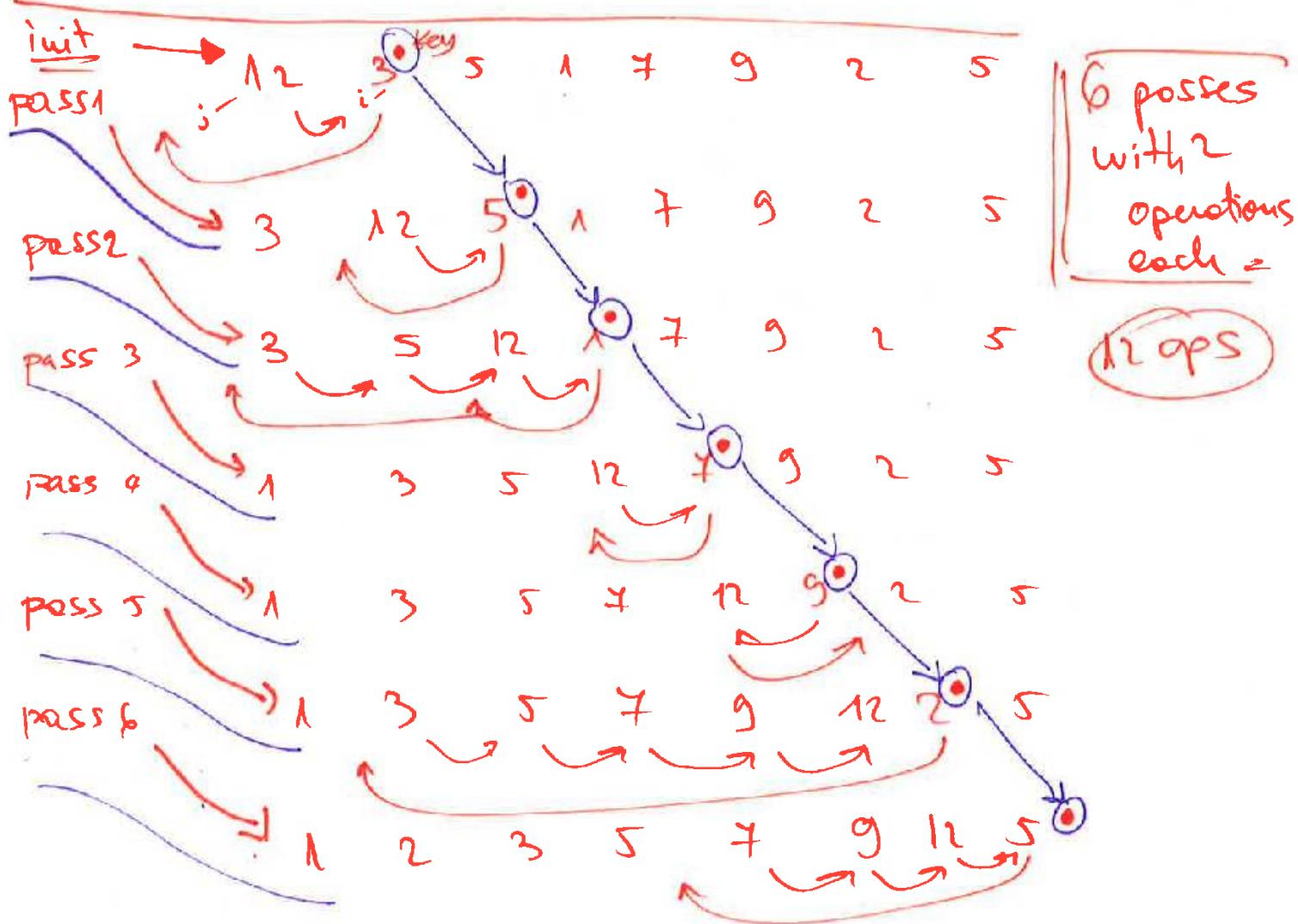
`end for`

Reversed order

Worst case time: $O(n^2)$

Best runtime: $O(n)$

already sorted



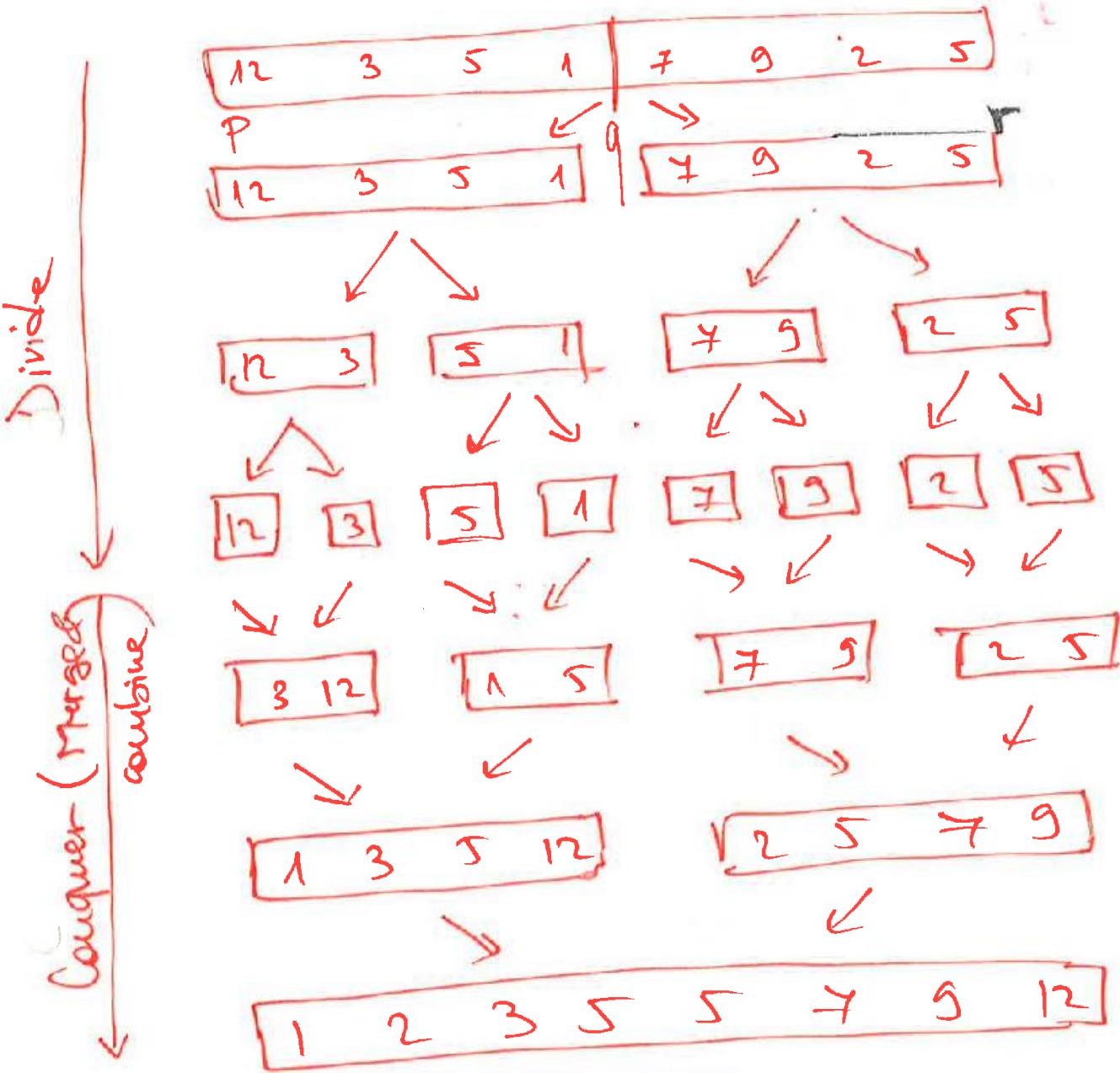
③

Merge Sort

- * Divide & Conquer
- * uses recursion

Analysis:

- Time complexity: $\Theta(n \log n)$
- Space complexity: $\Theta(n)$
- Number of moves: $\Theta(n^2)$

Code $A[]$ = input array P = left limit R = right limit q = middle L = left position R = right positionmerge-sort (A, P, R)if ($P < R$) $q = \frac{P+R}{2}$ merge-sort (A, P, q)merge-sort ($A, q+1, R$)merge (A, P, q, R)

end if

• merge(A, P, Q, R)

$n_1 = Q - P + 1$; left partition size

$n_2 = R - Q$; right partition size

for $i = 1$ to n_1

| $L[i] = A[P+i-1]$;

end for

for $j = 1$ to n_2

| $R[j] = A[Q+j]$;

end for

$i = j = 1$

for $k = P$ to R

| if $L[i] \leq R[j]$

| $A[k] = L[i]$

| $i++$

| else

| $A[k] = R[j]$

| $j++$

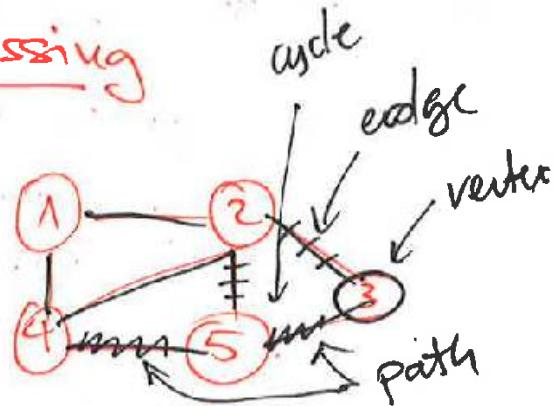
| endif

end for

Graphs & graph processing

- a powerful representation & processing tool;

- elements:



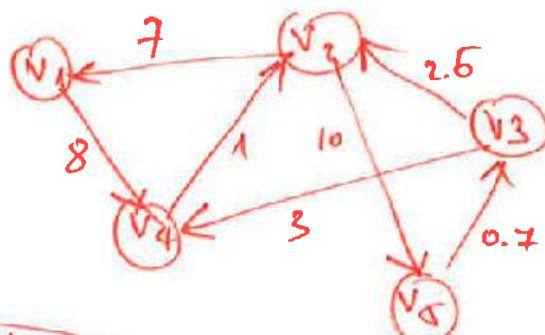
* Vertex = node, represented by a key and also some additional information (payload)

* edge = arc, connecting vertices & depicts a relationship between them
can be one/two way

* weight = arcs/edges can be weighted so that there is a cost associated with an edge \rightarrow can have different meanings

* path - a sequence of vertices connected by edges in a graph

* Cycle = a path that starts & ends in the same vertex



Representing connectivity:

Adjacency matrix

	v1	v2	v3	v4	v5
v1		7		8	
v2	7		2.5	1	10
v3		2.5		3	0.7
v4	8	1	10		
v5		10	0.7		

Adjacency list

id = "v1"	adj = {v2:7, v4:8}
v2	
v3	
v4	
v5	id = "v5" adj = {v2:10, v3:0.7}

Algorithms on graphs

- graphs can model many types & processes in information systems, both biological & technical
- graph traversal is an interesting problem and refers to visiting all nodes in a graph
BFS, DFS, Dijkstra, A*

BFS (Breadth First Search)

BFS(G, s)

s.parent = null

s.color = gray

s.distance = 0

enqueue(\emptyset, s)

while $Q \neq \emptyset$

n = dequeue(\emptyset)

for $v = n.\text{all Neighbors}$

if $v.\text{color} == \text{white}$

$v.\text{color} = \text{gray}$

$v.\text{distance} = n.\text{distance} + d(n, v)$

$v.\text{parent} = n$

$n.\text{neighbors} = n.\text{neighbors} + v$

enqueue(\emptyset, v)

endif

~~v.color = black~~

end for

~~$n.\text{color} = \text{black}$~~

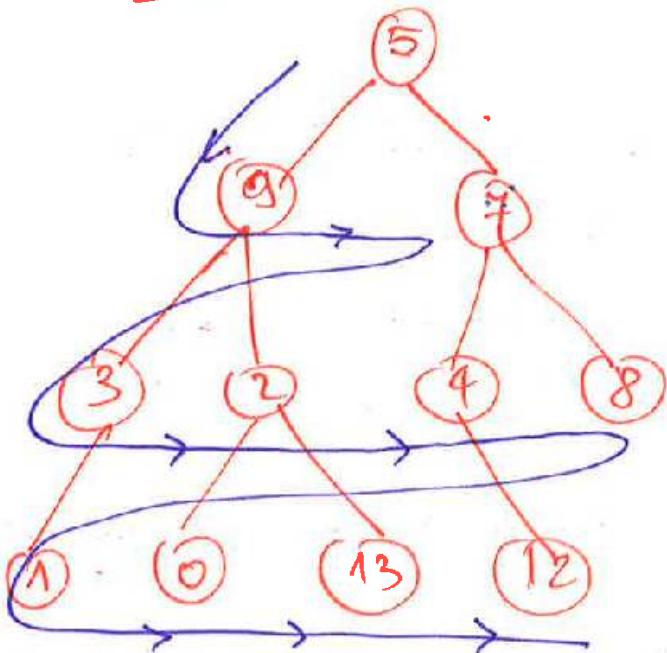
end while

- * BFS colors the vertices (graph coloring):
- * white - undiscovered
- * gray - visited unprocessed
- * black - visited & processed

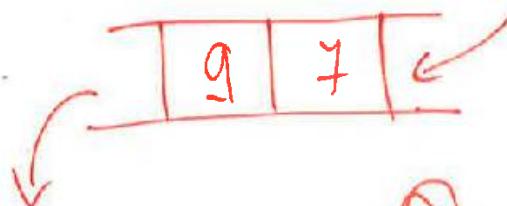
// Each node represented by triplets:

{	• parent_id
• color	
• distance	}

Example



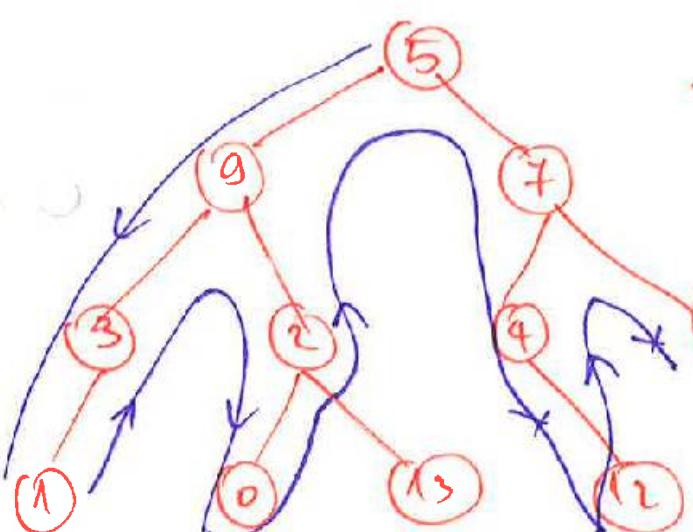
5 9 7 3 2 4 8 1 0 13 12



Queue

Depth First Search (DFS)

- similar to BFS but add vertices at the beginning instead of the end of the queue
- uses recursion



5 9 3 1 2 0 13 7 4 12 8

put



Stack

get

Dijkstra's algorithm

- algorithm to compute shortest path among vertices in a graph
- extends BFS with associating a cost on each edge \Rightarrow Cheapest First Search (CFS)

Elements : each vertex has 3 embedded properties

color = {white, grey, black}
cost = {positive numeric value}
parent = id of the parent node

Algorithm:

$n(\text{all}).\text{color} = \text{white}$; all nodes unvisited & unprocessed
 $n(\text{all}).\text{cost} = \infty$; cost unknown
 $n(s).\text{parent} = 0$; start node has no parent
 $Q = \{\text{all nodes}\}$; all nodes are in the queue

while $Q \neq \{\}$.

$n = \text{dequeue}(Q)$; extract a node
 $n.\text{color} = \text{grey}$; visited & not processed

for $c = n.\text{all_children}$

if $c.\text{cost} > n.\text{cost} + \text{edge}(c, n)$
 $c.\text{cost} = n.\text{cost} + \text{edge}(c, n)$
 $c.\text{parent} = n$

endif

end for

$n.\text{color} = \text{black}$

end while

AIML Class

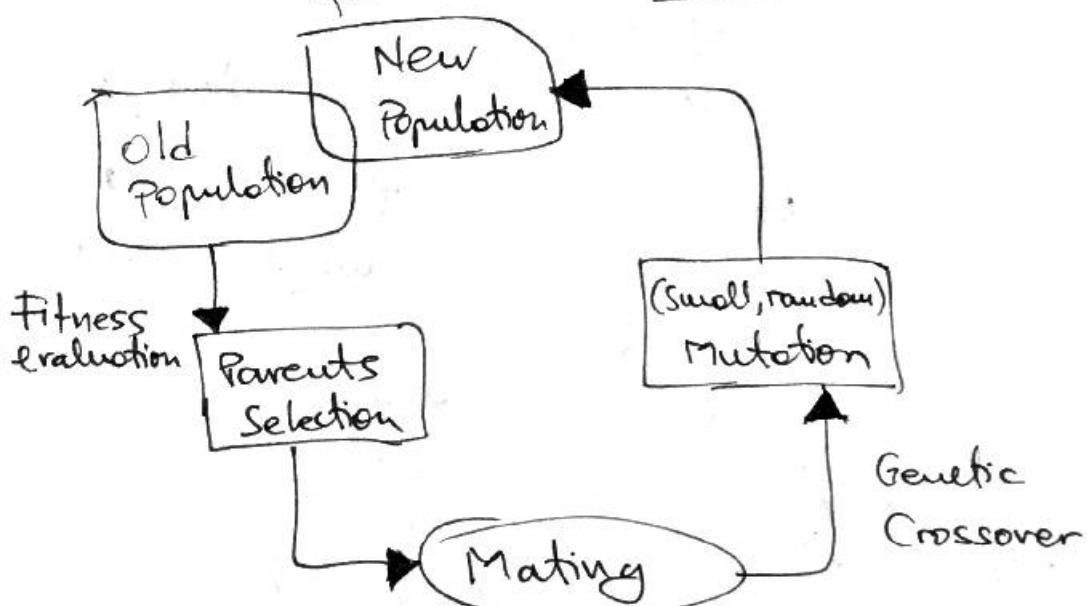
GA & EP

Genetic Algorithms

Evolutionary Programming

★ Evolutionary Programming (EP) is a method for simulating evolution; rooted in Darwin's (1859) theory of evolution, natural selection and survival of the best/fittest.

- There are 3 main lines of research:
 1. Genetic Algorithms (GA) *chromosome
 2. Evolution Strategies (ES) *individual
 3. Evolutionary Programming (EP) *specie
- The core idea of all of them:
 - population of competing solutions subject to random alterations.
 - competition to be retained as parents for successive reproduction epochs.



- EP is an abstraction from biological theory of computation/evolution used for optimization problems.

—EP emerged as an alternative to AI:

- doesn't emulate neural computation or human behaviors.
 - models processes that generate organisms or increase intellect over time:
 - intelligence defined here as ability to achieve goals and evolution is optimization.

* Genetic Algorithms (GA)

- introduced formally in 1970 by John Holland.
 - attractive for optimization problems.
 - work well on mixed (continuous and discrete) combinatorial problems.

Basic Foundation:

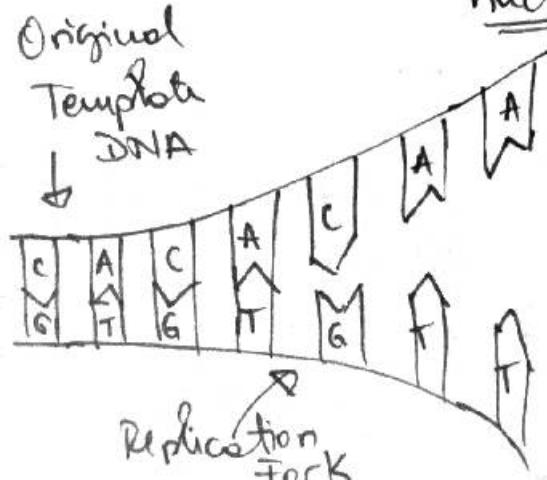
- the gene is a part of the DNA that encodes information; chromosomes contain multiple genes.
 - the human genome has \sim 20500 genes.
 - the DNA sequence is composed of complementary nucleotides:
These bases form

These bases form
pairs between the

2 strands:

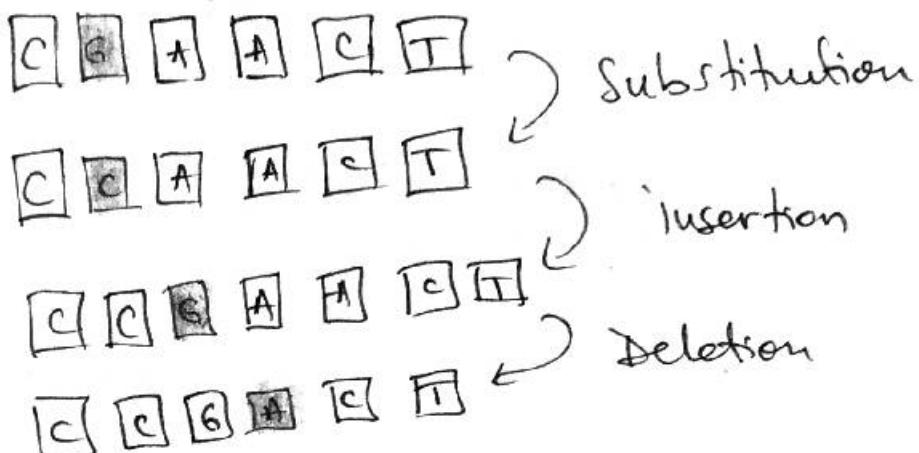
A with T

C with G

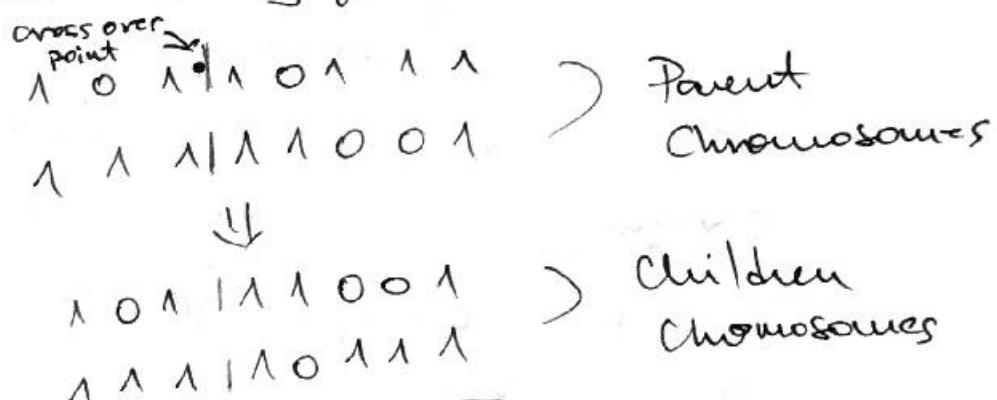


- during reproduction subprocesses can occur :
 - mutations (substitution, insertion, deletion)
 - cross-over

Mutation



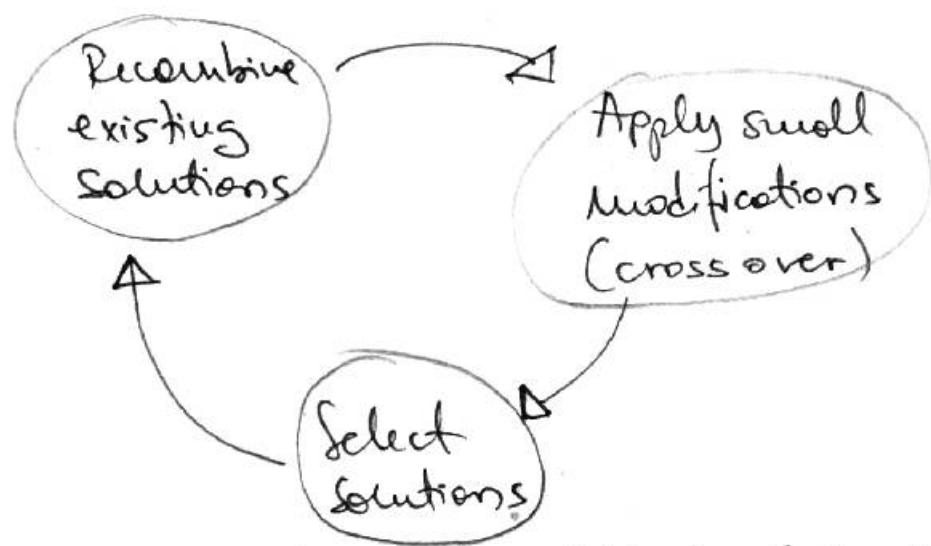
Cross-over (assume binary genes)



But how can we apply this in technology?

- consider the problem of searching a large state-space or n-dimensional surface.
- GAs offer advantage over traditional optimization techniques (BFS, DFS, Dijkstra etc.) \Rightarrow but they also tend to be computationally expensive.
- GAs are adaptive search algorithms based on genetics which offer an intelligent random search.

- although randomized GAs are not random
- they exploit historical information to direct the search into a region of better performance.
- it is seen as a 3-step iterative process of stochastic exploration:



- GAs are less susceptible to fall in local minima than gradient methods.
- GAs can optimize nonlinear, discontinuous functions without an analytical formulation.

Elements of GA:

4. Encoding

- a chromosome is composed of genes and should represent the solution.
- usually a binary string.
- encoding is problem dependent!

B. Population (of chromosomes)

- the size of the population is problem dependent
- initial population is generated randomly to allow spanning entire search space.
- the size determines the problem complexity and exploration ability.

c. Fitness function

- evaluates how "good" a potential solution is relative to others.
- usually returns a value reflecting how optimal a solution is.
- this is the most important design aspect of the GA !

D. Selection Criteria

- Selecting parents for recombination and crossovers:
- usually select better parents (fitness wise) to produce better children.
- uses elitism or other similar strategies
 - * proportional selection;
 - * tournament selection;
 - * rank based selection.

Save best solution for next generation

E. Reproduction / mutation operators

- crossover operates on parent chromosomes generating offspring chromosomes

- typically choose a random crossover point and copy from first parent with the remaining genes from the second parent.
- after crossover mutation takes place.
- mutation ensures that the search doesn't fall in a local minima by randomly changing offspring resulting from crossover.

Using GAs in practice: Optimizing Nonlinear Functions

Problem: Given a nonlinear function f

$$f(x) = x \cdot \sin(10\pi x) + 1 \text{ find the } x \text{ value such that } f(x) \text{ is maximum.}$$

Considering $x \in [-1, 2]$, the problem assumes finding x_0 such that $f(x_0) \geq f(x)$ $\forall x \in [-1, 2]$ and $x \in [-1, 2]$.

Let's look first at the analytical approach.

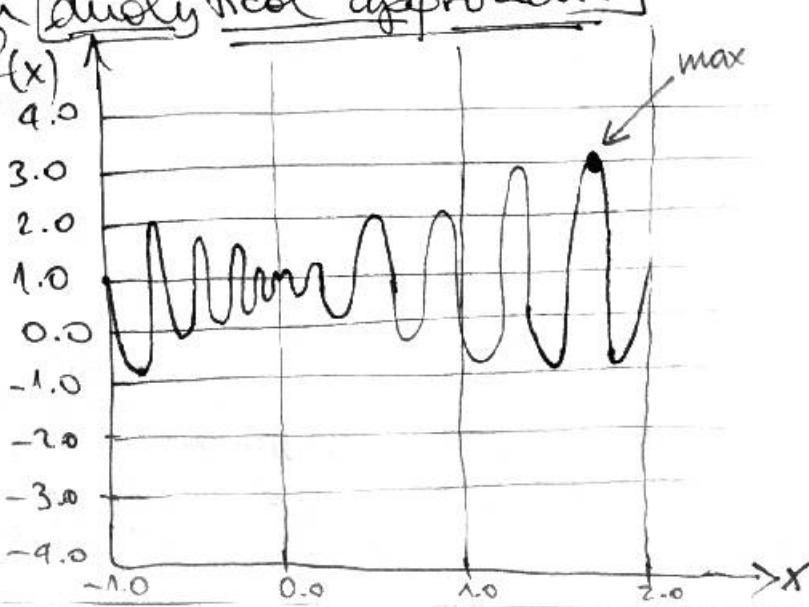
- compute the derivative $f'(x)$

$$f'(x) = \frac{d}{dx} f(x) = x' \sin(10\pi x) + x (\sin(10\pi x))'$$

$$f'(x) = 1 \cdot \sin(10\pi x) + 10\pi x \cos(10\pi x)$$

we make $f'(x) = 0$

$$\sin(10\pi x) + 10\pi x \cos(10\pi x) = 0$$



we divide with $\cos(10\pi x) > 0$ so we have

$$\frac{\sin(10\pi x)}{\cos(10\pi x)} + 10\pi x = 0$$

$$\tan(10\pi x) = -10\pi x$$

This eq. has an infinite nr. of solutions given by

$$\left\{ \begin{array}{l} x_i = \frac{2i-1}{20} + ei, i=1,2,\dots \\ x_0 = 0 \\ x_i = \frac{2i+1}{20} + ei, i=-1,-2,\dots \end{array} \right.$$

Because we look at $x \in [-1, 1]$ then $f(x)$ reaches the max for: $x = 1.85 \Rightarrow f(1.85) = 1.85 \cdot \sin(18.5\pi) + 1$

$$f(1.85) = 1.85 + \sin\left(18\pi + \frac{\pi}{2}\right) + 1$$

$f(1.85) = 2.85$ which is the max value

The GA approach

Step 1 Encoding

A. We binary encoding to represent x .

B. choose the dimension (nr. of bits)

for our case we have 3 types of values (see analytic formulation)

- ① we would like 6 places after the decimal point (precision)

— so, we should split the $[-1, 2]$ interval
in 3×10^6 equal sized ranges.

$$2^{21} = 2097152 < \underbrace{3000000}_{3 \times 10^6} < 2^{22} = 4194304$$



so, we use a 22 bit genes
representation for chromosomes

$$\langle b_{21}, b_{20}, b_{19}, \dots, b_0 \rangle$$

c. Perform the mapping from binary representation
of chromosomes to real values $[-1, 2]$

step 1 : convert binary to decimal

$$\left(\langle b_{21}, b_{20}, \dots, b_0 \rangle \right)_2 = \left(\sum_{i=0}^{21} b_i \cdot 2^i \right)_{10} = x^*$$

step 2 : encode value in interval $[x_{\min}, x_{\max}]$

$$x^* = (x - x_{\min}) \cdot \left(\frac{2^n - 1}{x_{\max} - x_{\min}} \right)$$

where $\begin{cases} x_{\min} = -1 \\ x_{\max} = 2 \end{cases}$ in our case



$$x = x_{\min} + x^* \cdot \frac{x_{\max} - x_{\min}}{2^n - 1}$$

and for our case

$$x = -1 + x^* \cdot \frac{3}{2^{22} - 1}$$

Ex: to which x the next chromosome corresponds?

$$c = (1|0|0|0101110110101000|111)_{b_1 b_2 b_3 b_4 b_5}$$

$$x^* = 1 \cdot 2^{21} + 0 \cdot 2^{20} + 0 \cdot 2^{19} + \dots + 1 \cdot 2^4 + \dots + 1 \cdot 2^1 + 1 \cdot 2^0 \Rightarrow$$

$$x^* = 2288964$$

But we know that

$$x = -1.0 + x^* \frac{3}{2^{22}-1}$$

$$\Downarrow$$

$$x = -1.0 + \frac{3 \cdot 2288964}{2^{22}-1}$$

$$\underline{x = 0.637197.}$$

We can easily see

$$c = (\underbrace{0 \quad \dots \quad 0}_{22 \text{ bits/gene}}) = -1$$

$$c = (\overbrace{1 \quad \dots \quad 1}) = 2.$$

Step 2 Initialize the population

- randomly initialize the 22 bits in each chromosome of the population.
- define the number of chromosomes in the population.

Step 3 Define fitness function, e

- e is applied on binary values c and is equivalent to function f : $e(c) = f(x)$

— e can be seen as the environment
rating potential solutions in terms of
their fitness.

e.g.: Consider 3 chromosomes:

$$\left\{ \begin{array}{l} C_1 = (1 \overset{21}{0} 0 0 1 0 1 1 1 0 1 1 0 1 0 1 0 0 0 1 1 1) \\ C_2 = (0 0 0 0 0 0 0 1 1 1 0 0 0 0 0 0 0 1 0 0 0 0 0) \\ C_3 = (1 1 1 0 0 0 0 0 0 1 1 1 1 1 1 1 0 0 0 1 0 1) \end{array} \right.$$

corresponding to

$$\left\{ \begin{array}{l} x_1 = 0.637157 \\ x_2 = -0.953843 \\ x_3 = 1.627888 \end{array} \right.$$

for which the evaluation / fitness function gives the following results:

$$e(c_1) = f(x_1) = 1.586345$$

$$e(c_2) = f(x_2) = 0.048878$$

$$e(c_3) = f(x_3) = 2.250650$$

↗ chromosome 3 is the fittest!

Step 4] Crossover

e.g. — consider c_2 and c_3 for crossover.
— the crossover point is after gene / bit index 17.

cross over point bit 17

$$c_2 = (0\ 0000 | 01110000000100000)$$

$$c_3 = (11100 | 0000011111000101)$$

\downarrow crossover 2 resulting children

$$c'_2 = (00000 | 00000111111000101)$$

$$c'_3 = (11100 | 01110000000100000)$$

if we evaluate the fitness of the 2 children

$$e(c'_2) = f(x'_{c_2}) = f(-0.008113) = 0.940$$

$$e(c'_3) = f(x'_{c_3}) = f(1.66) = 2.489$$

$$\left\{ \begin{array}{l} e(c_2) = 0.078878 \\ e(c_3) = 2.250650 \end{array} \right.$$

Second offspring has a better fitness than both parents!

Step 5 [Mutation]

- alters one / more genes with a probability equal to the mutation rate

e.g.: assume gene 17th of c_3 was selected for mutation:

$$c_3 = (1110\boxed{0}00000111111000101)$$

\downarrow mutation

$$c'_3 = (11101000\overline{0}11111000101)$$

$$x_3' = -1.0 + x_3^{1*} \frac{3}{2^{22}-1} = 1.721638$$

$$e(c_3') = f(x_3') = -0.082257$$

but $e(c_3) = 2.250650$



in this case: mutation determined a decrease of the fitness value!

- If we mutate gene 11 in chromosome 3:

$$c_3'' = (\overset{\circ}{1110000001} 111111000101)$$

with

$$x_3'' = 1.630818 \Rightarrow e(c_3'') = f(x_3'') = 2.34355$$



in this case: improvement of fitness over initial value!

Generic GA algorithm pseudocode

- #1. Initialize chromosomes in the population P
- #2. Evaluate $e(c) = f(x)$
- #3. Select "good" individual
- #4. Repeat
 - #5. for $i = 0$ to λ (# of children to produce)
 - #6. select P_1, P_2 parents ^{out} of P -
 - #7. generate x_i by crossover P_1, P_2
 - #8. mutate x_i
 - #9. evaluate $f(x_i)$
 - #10. add x_i to P'
 - #11. end
 - #12. select new P out of P' (different selection strategies)
 - #13. forget P'
- #14. Until fitness is "good" enough

In our example to maximize $f(x) = x \sin(10\pi x) + 1$

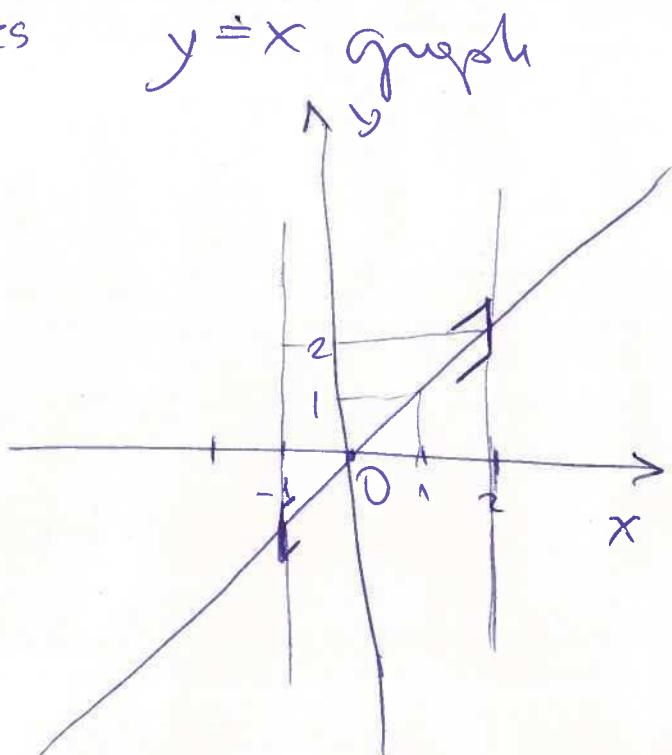
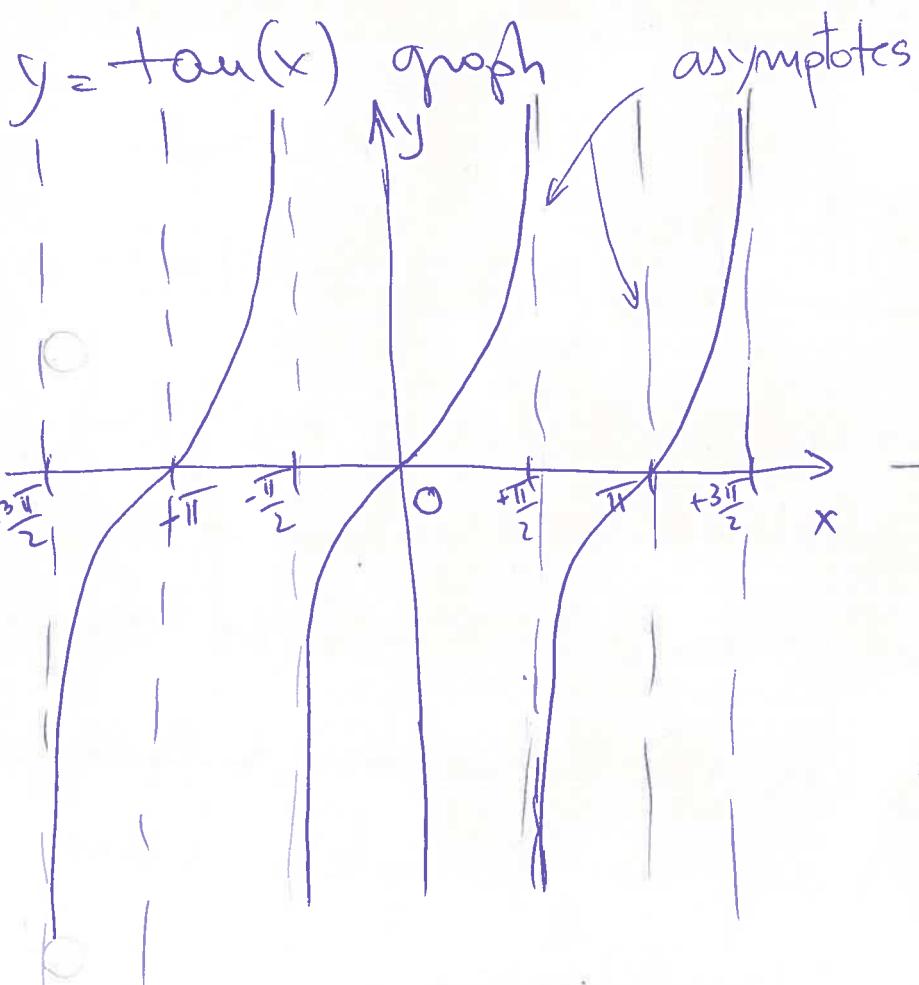
after 150 generations the fittest chromosome is

$c_{max} = (111100110100010000101)$ which corresponds
to $x_{max} = 1.850000$ and $f(x_{max}) = 2.850.000$.

Trigonometric equation proof

Appendix 1

$\tan(10\pi x) = -10\pi x$ brought to
the simplest form
 $\tan(x) = x$



the eq. has a solution
in every interval of
the following form

$$\left[\frac{(2k-1)\pi}{2}, \frac{(2k+1)\pi}{2} \right] = I(k)$$

$k \in \mathbb{N}$

\Leftrightarrow

$$\frac{(2k-1)\pi}{2} < x < \frac{(2k+1)\pi}{2}$$

When x is in the $I(k)$
and since $x = y$ is
clear that y is in the
 $I(k)$ and so

$$|y| < \frac{(1+2 \cdot |k|) \cdot \pi}{2}$$

So if we put
things together

- for values of x close to the left end of $I(k)$

$$\tan(x) < -\frac{(2|k|+1)\pi}{2} < x \Rightarrow \tan(x)-x < 0$$

- for values of x close to the right end of $I(k)$

$$\tan(x) > \frac{(2|k|+1)\pi}{2} > x \Rightarrow \tan(x)-x > 0$$

Since $\tan(x)$ is continuous on $I(k)$ there must be at least one value of x between where $\tan(x)-x=0$ so $\tan(x)=x$ so if we consider we have one solution for each k there are infinitely many solutions to $\tan(x)=x$.

Appendix 2 General mapping function u-bits of represent.

$$x^* = (x - x_{\min}) \cdot \frac{2^n - 1}{x_{\max} - x_{\min}}$$

$$x^* = \frac{x(2^u - 1)}{x_{\max} - x_{\min}} - \frac{x_{\min}(2^u - 1)}{x_{\max} - x_{\min}}$$

$$x^*(x_{\max} - x_{\min}) = (2^{u-1})(x - x_{\min})$$

$$x - x_{\min} = \frac{x^*(x_{\max} - x_{\min})}{2^u - 1}$$

$$x = x_{\min} + \frac{x^*(x_{\max} - x_{\min})}{2^u - 1}$$

Fuzzy Logic

Ex 1 Worked example on Mandani FIS.

Evaluate the risk of a software engineering project by referring only 2 criteria, the project funding and the project staffing.

Step 1 : I/O variables identification

2 inputs : project funding (%) (PF)
project staffing (%) (PS)

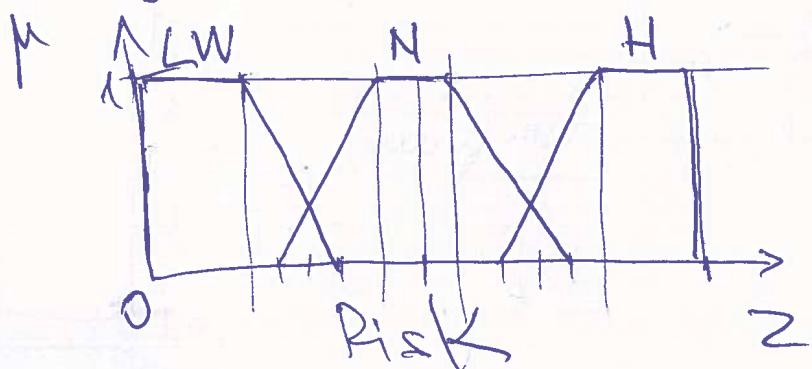
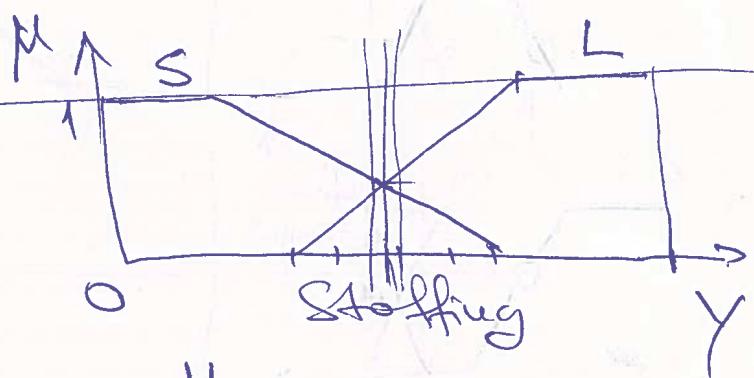
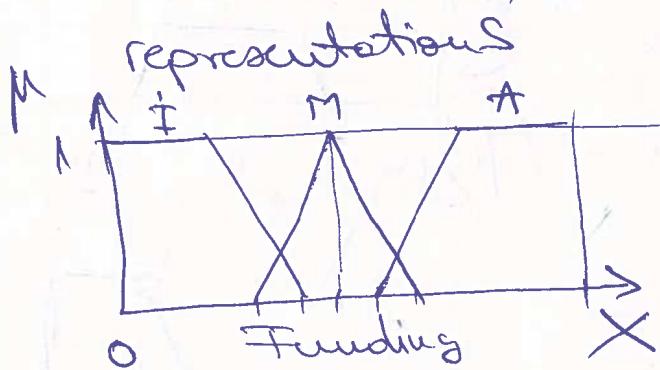
1 output : project risk (%) (PR)

universe of discourse \rightarrow fuzzy sets

$X = \{ \text{inadequate (I)}, \text{marginal (M)}, \text{adequate (A)} \}$ for PF

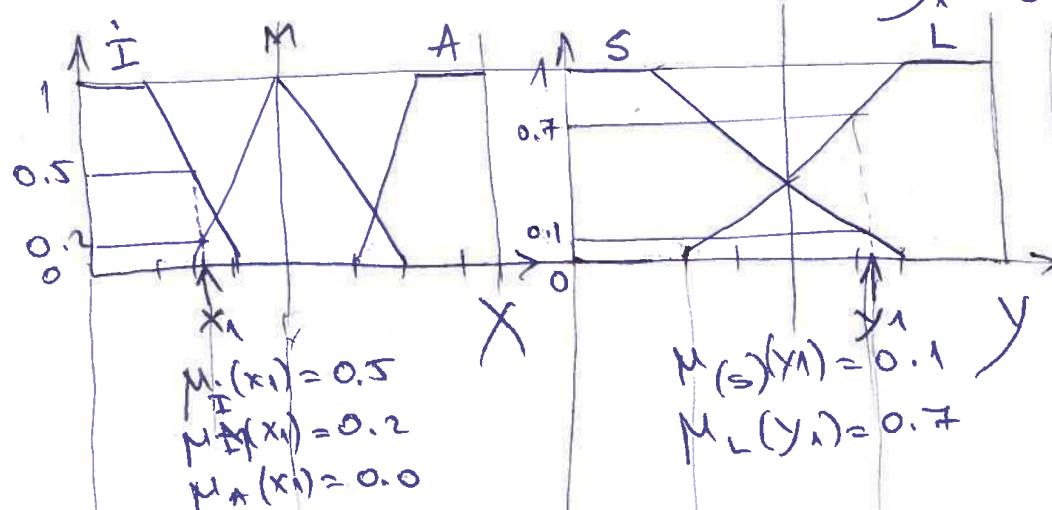
$Y = \{ \text{small (S)}, \text{large (L)} \}$ for PS

$Z = \{ \text{low (LW)}, \text{normal (N)}, \text{high (H)} \}$ for PR



Step 2 Fuzzification

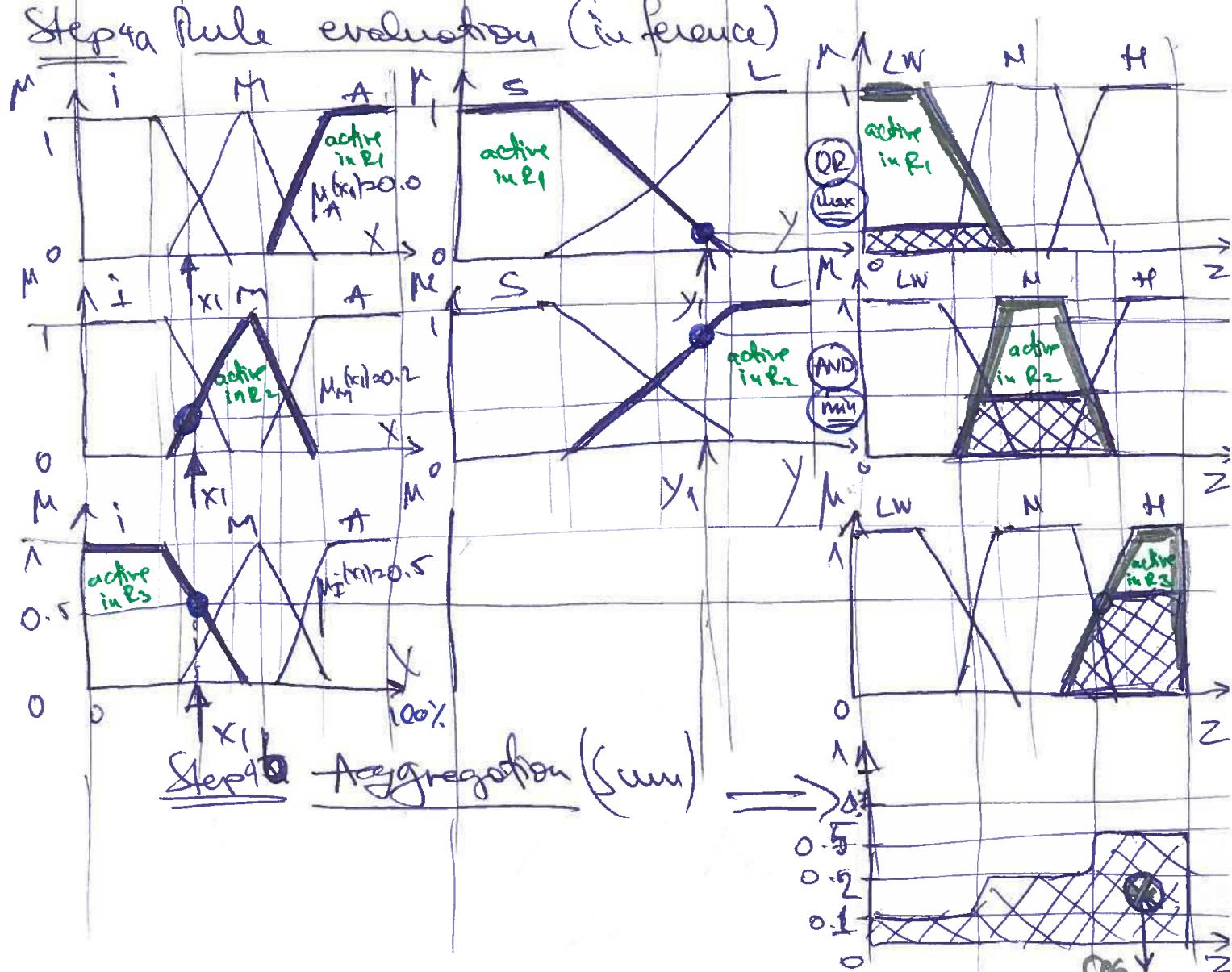
ex $x_1 = 35\% \text{ PF}$
 $y_1 = 60\% \text{ PS}$



Step 3 Rule design (example rules)

- R₁: IF PF is I OR PS is S THEN PR is LW.
 R₂: IF PF is M AND PS is S THEN PR is N.
 R₃: IF PF is I THEN PR is H.

Step 4a Rule evaluation (Inference)



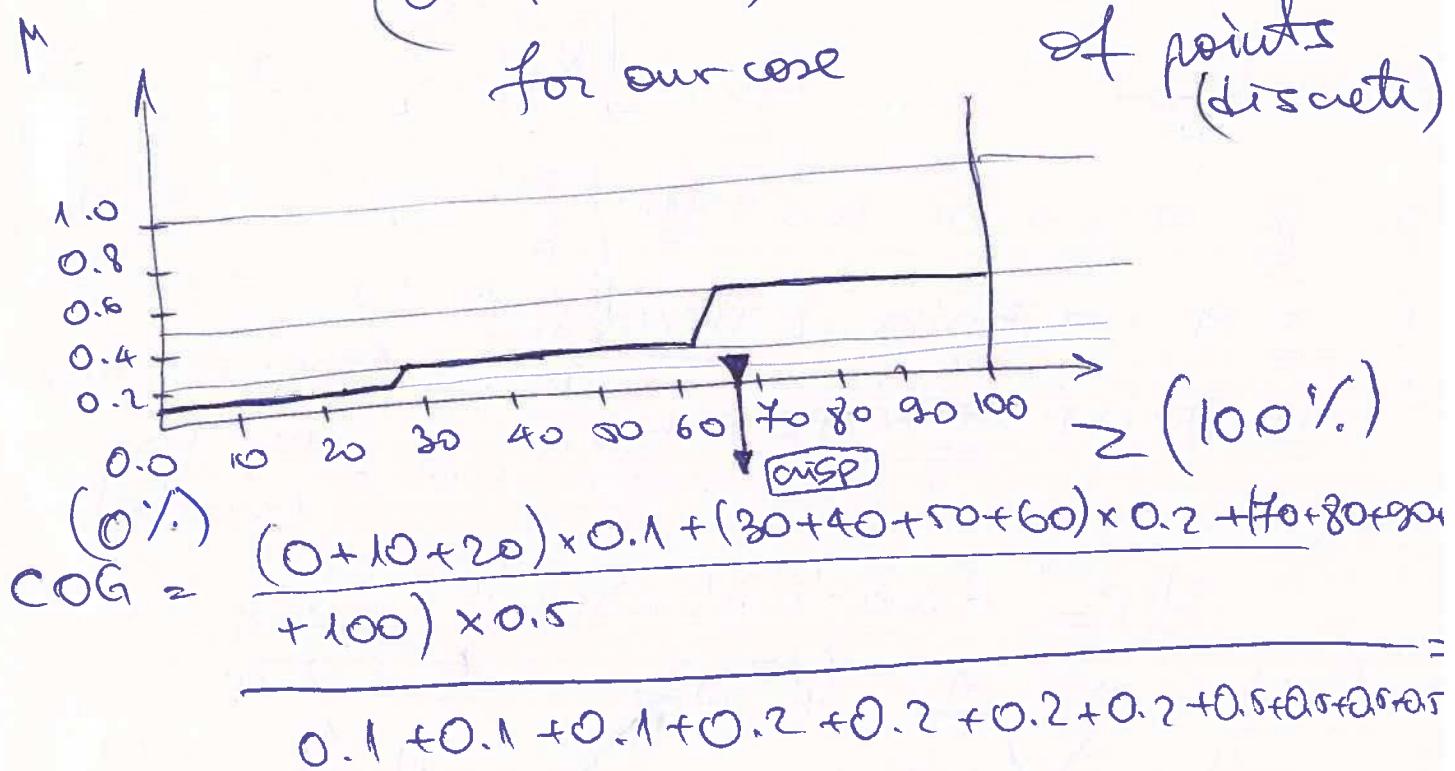
Step 5 Defuzzification (using COG)

$$COG = \frac{\int_a^b \mu_A(x) \cdot x \, dx}{\int_a^b \mu_A(x) \, dx} \rightarrow COG = \frac{\sum_{x=a}^b \mu_A(x) \cdot x}{\sum_{x=a}^b \mu_A(x)}$$

(continuous)

for our case

Over a sample
of points
(discrete)



$$\underline{COG = 67.4 \%}$$

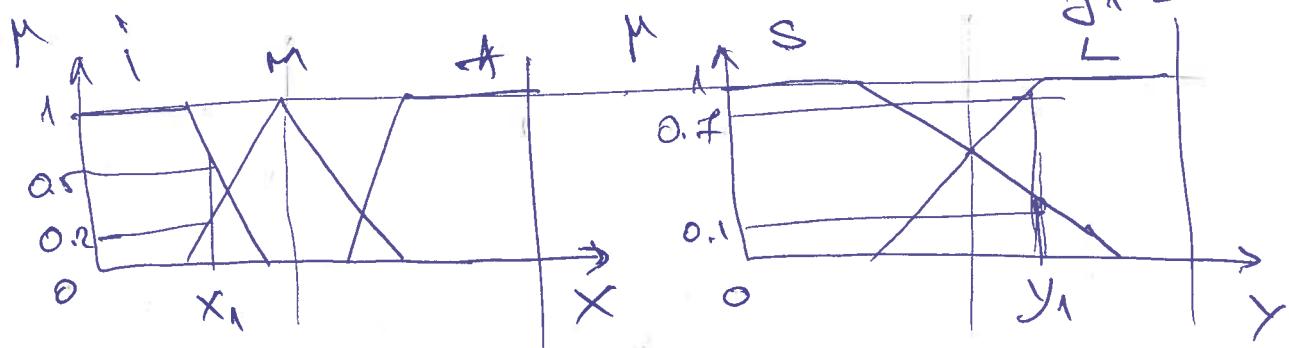
Ex ② Analogous problem solved with TSK FIS

- the inference time is shortened by eliminating the defuzzification and by using negligibles as outputs

- Singleton = fuzzy set with $\mu = 1$ at a single point in the universe of discourse and 0 elsewhere.

The output of the TSK is a function of the input variables!

Step 2 Fuzzification

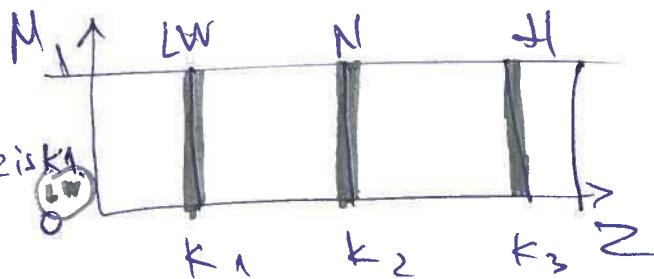


$$x_1 = 35\% \text{ PF}$$

$$y_1 = 60\% \text{ PS}$$

Step 3 Rule design

R1: IF PF is A OR PS is S THEN PR is k1



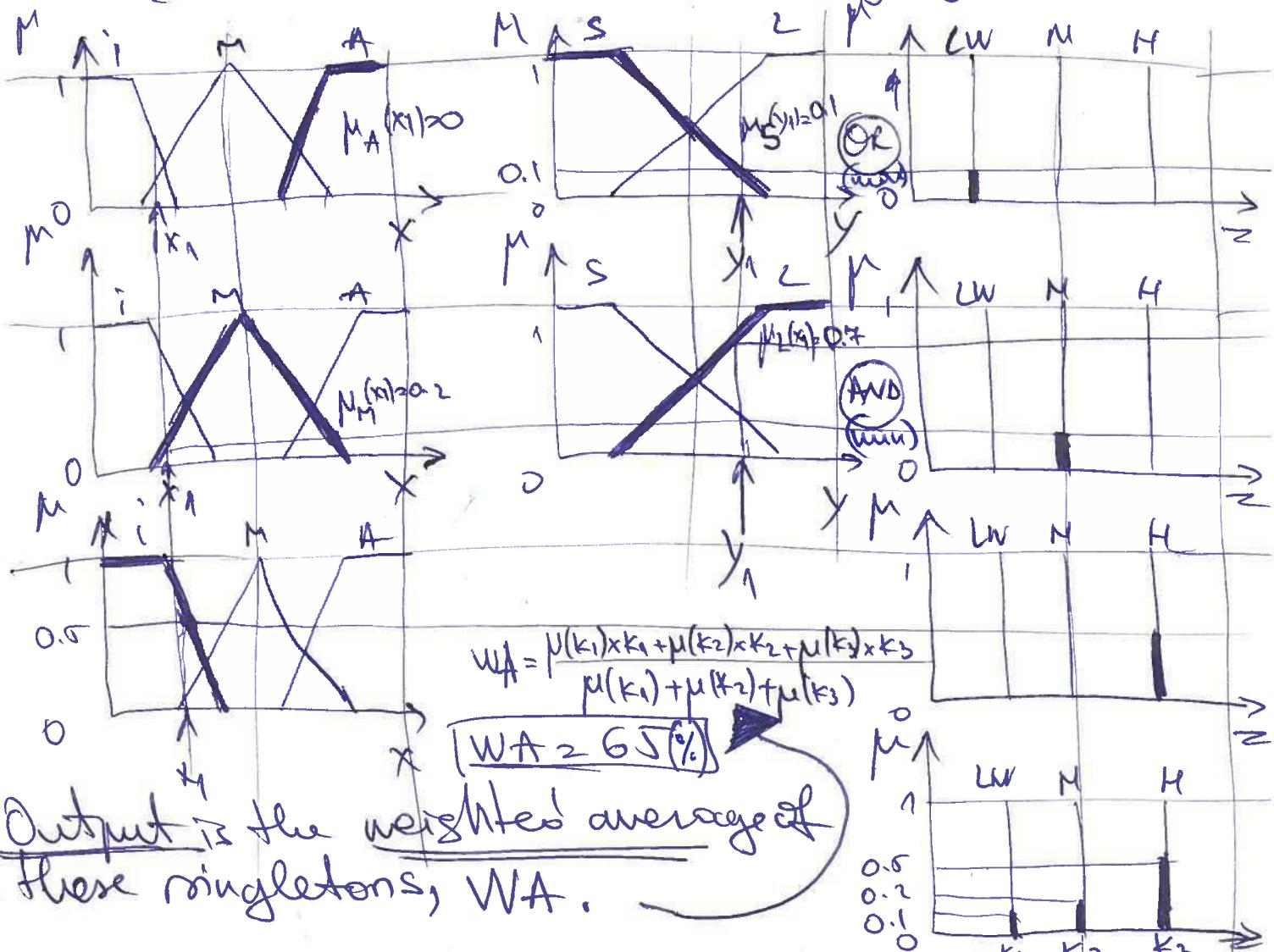
R2: IF PF is M AND PS is L THEN PR is k2

R3: IF PF is H THEN PR is k3

Output representation
in TSK
(singletons)

$$k_1 = 20, k_2 = 50, k_3 = 80$$

Step 4 Rule evaluation & aggregation



AIML Course 9 | Fuzzy Logic & Fuzzy Inference Systems

Intro:

- Data representation and computation inspired by the way humans describe phenomena through language
- Uses imprecision and involves modifiers of linguistic terms (e.g. "quite", "fairly", "very") and uncertainty:
 - Q: "How comfortable is this room?"
 - A: "It's quite cold here.
It's fairly warm.
It's too hot for me."
- Trades off between significance and precision in computing with words.

Formalism and terminology:

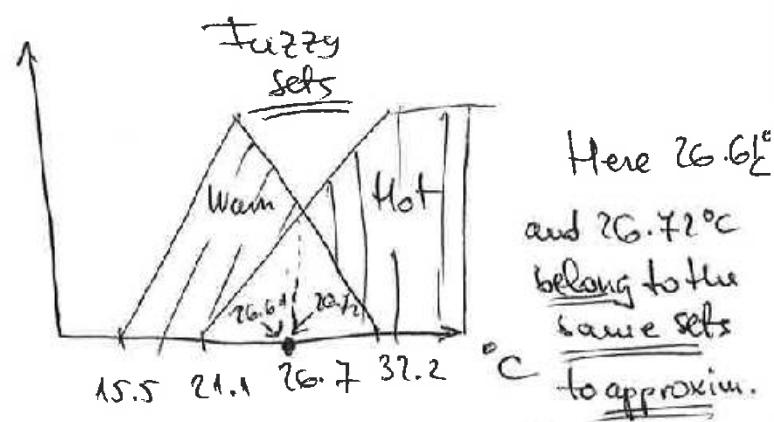
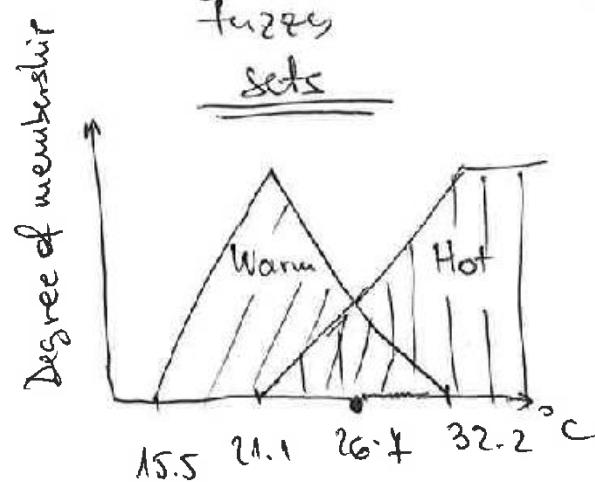
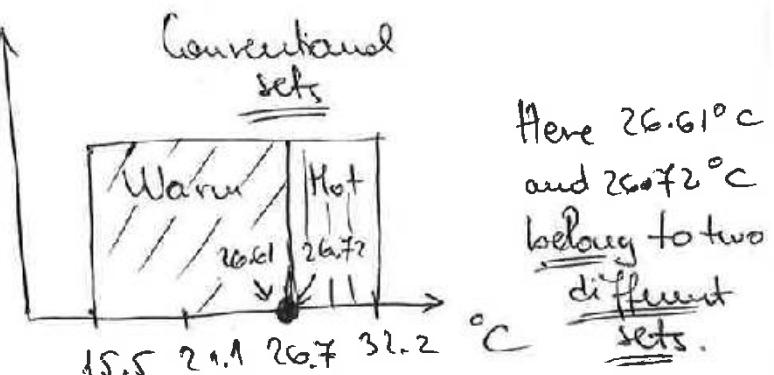
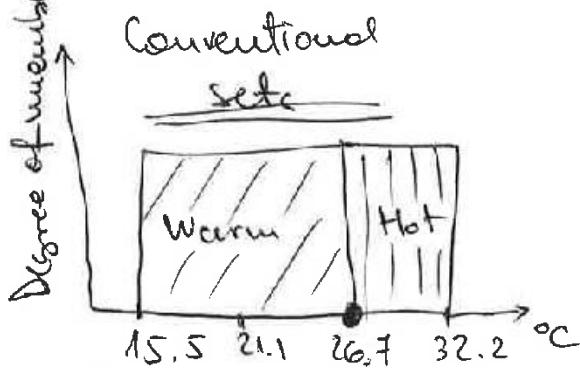
- Introduced by Lotfi Zadeh, 1965, "Fuzzy Sets"
- fuzzy logic is a polyvalent (multi-valued) logic
- Boolean logic is a special case of fuzzy logic
- fuzzy logic originates from set theory:

* crisp sets = conventional set, the degree of membership of any element is 0 or 1.

* fuzzy sets = sets whose elements have degrees of membership typically between 0 and 1.

Intuition

For example is 26.7°C warmer or hot?



Crisp values
(temperature)

Using fuzzy logic

Fuzzification:

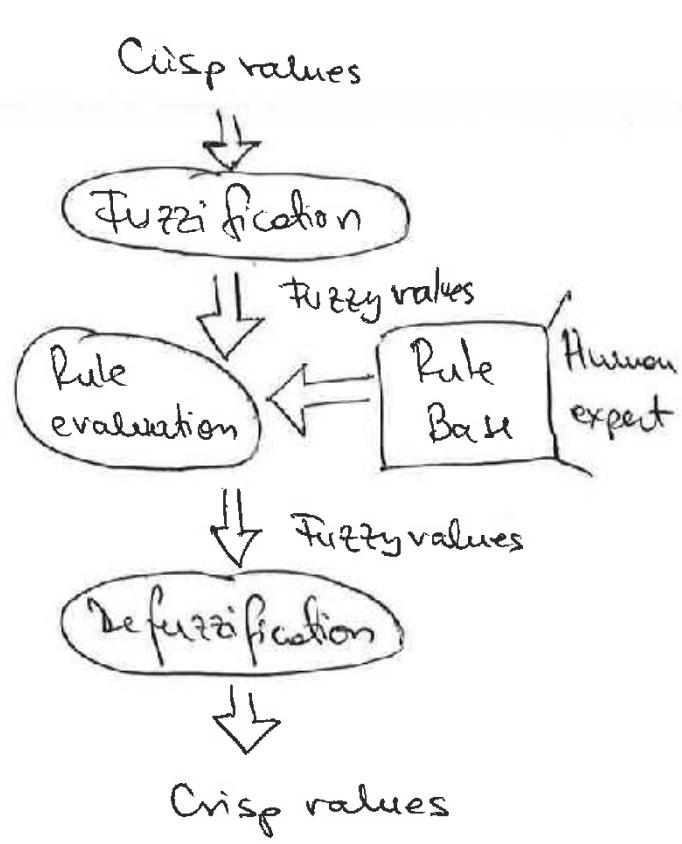
- fuzzy sets & logic
- linguistic variables

Rule evaluation:

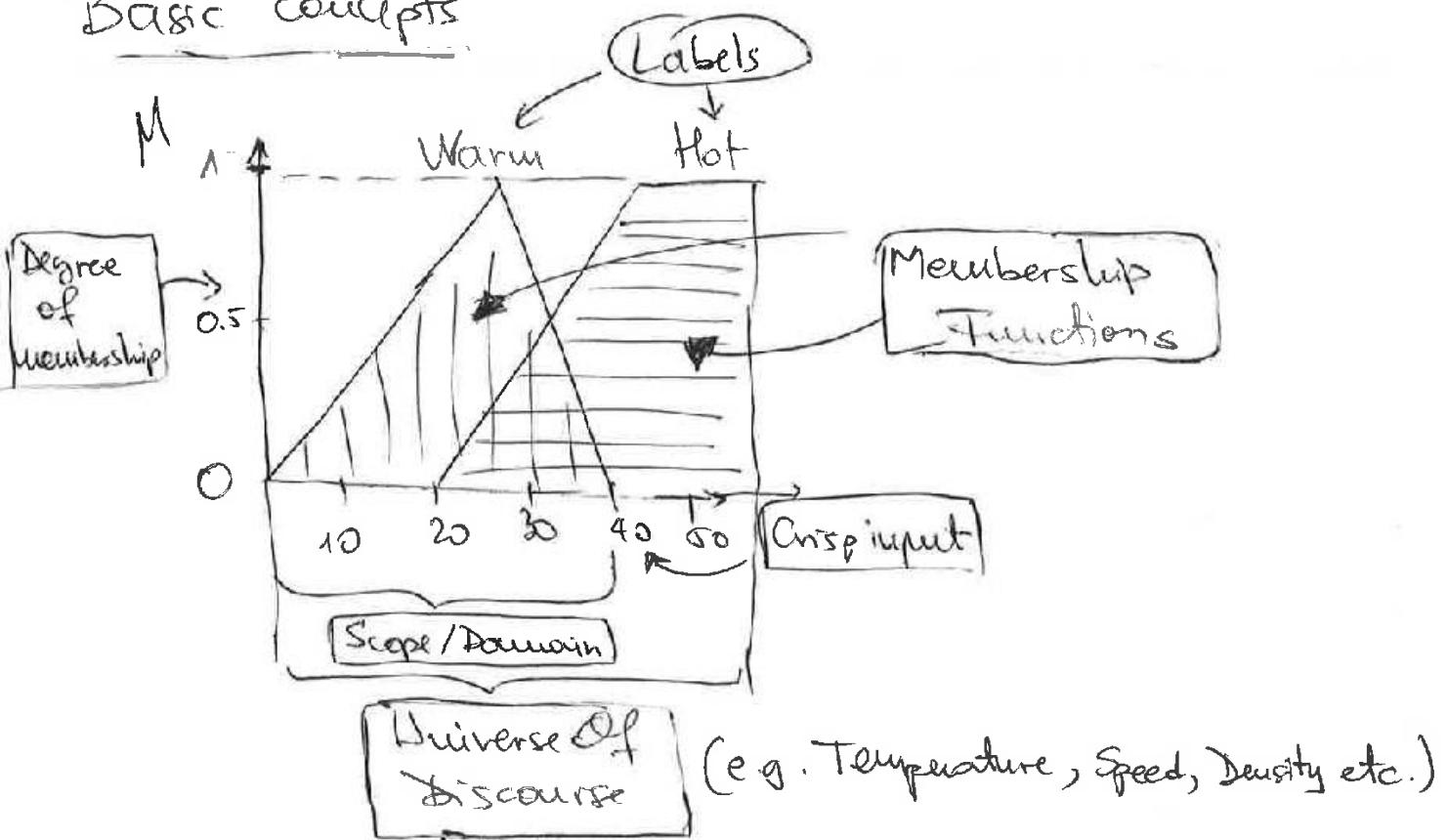
- fuzzy set operations
- fuzzy relations & inference
- fuzzy composition & reasoning

Defuzzification:

- extract crisp information from fuzzy sets



Basic concepts



Definitions:

A fuzzy set A corresponding to a universe of discourse X can be represented by an ordered set of pairs:

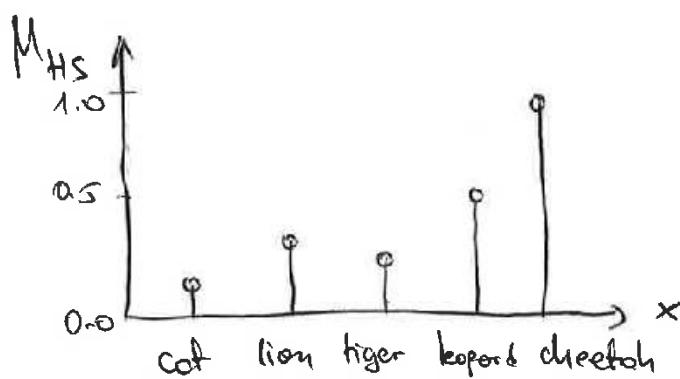
$$A = \{ (x, \mu_A(x)) \mid x \in X \}$$

E.g. Universe of discourse: cat families species

$$X = \{ \text{"cat"}, \text{"lion"}, \text{"tiger"}, \text{"leopard"}, \text{"cheetah"} \}$$

Fuzzy set HS for animals with high speed

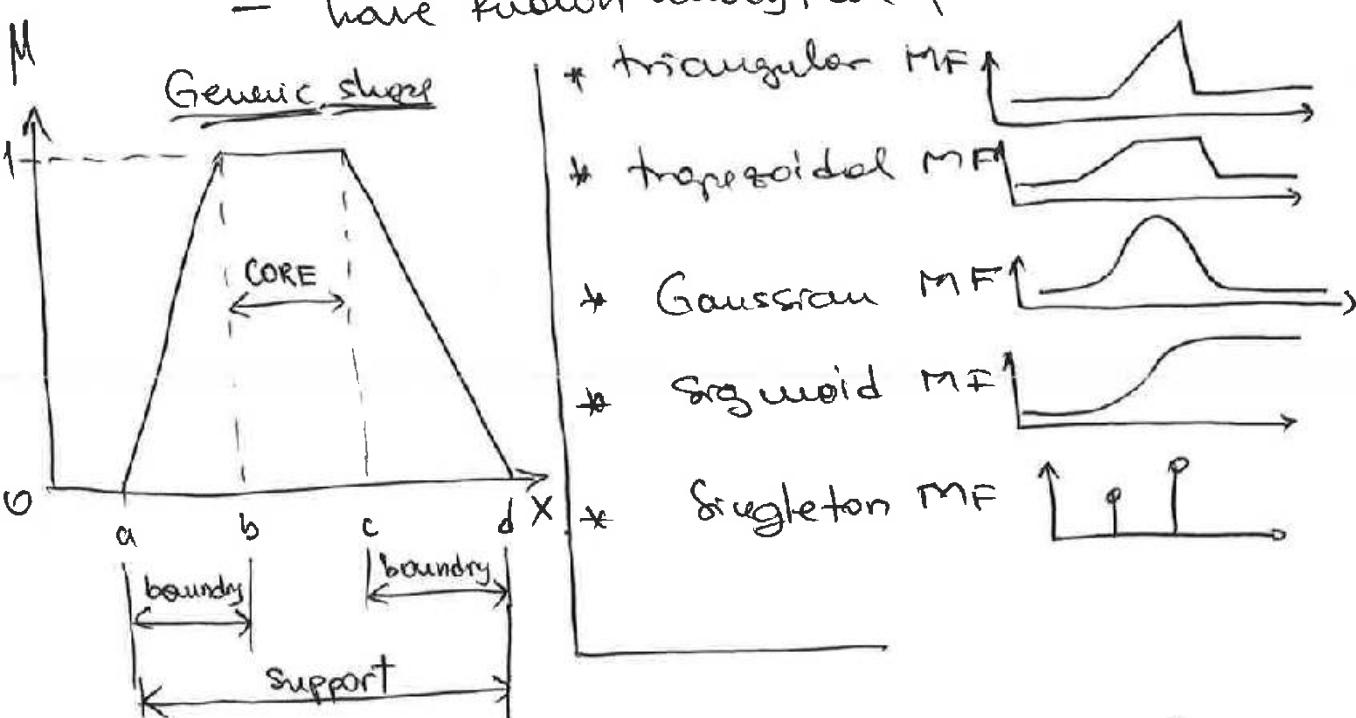
$$\text{HS} = \{ (x, \mu_{\text{HS}}(x)) \mid x \in X \} = \{ (\text{cat}, 0.1), (\text{lion}, 0.3), (\text{tiger}, 0.2), (\text{leopard}, 0.5), (\text{cheetah}, 0.9) \}$$



Intervals of fuzzy logic and fuzzy inference systems.

Membership functions

- mathematical functions that give numerical meaning to a fuzzy set.
- map crisp inputs from a domain to membership degrees in [0,1].
- have known analytical formulations :



support : $\text{supp}(A) = \{x \in X \mid \mu_A(x) > 0\}$

boundary : bnd : $\text{bnd}(A) = \{x \in X \mid 0 < \mu_A(x) < 1\}$

core : $\text{core}(A) = \{x \in X \mid \mu_A(x) = 1\}$

Singleton : fuzzy set with 1 element $\text{supp}(A) = \text{core}(A) = \{x_0\}$

Operations on membership functions

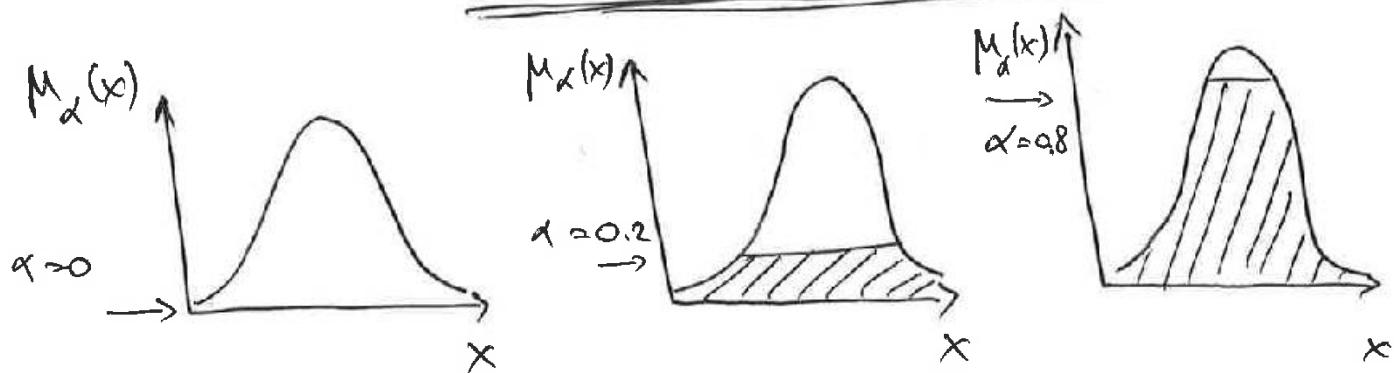
Alpha cut set A_α is a cusp set of a fuzzy set

where

$$A_\alpha = \{x \in X \mid \mu_A(x) \geq \alpha\}$$

and if $A_\alpha = \{x \in X \mid \mu_A(x) > \alpha\}$ we

have a strong alpha cut set.

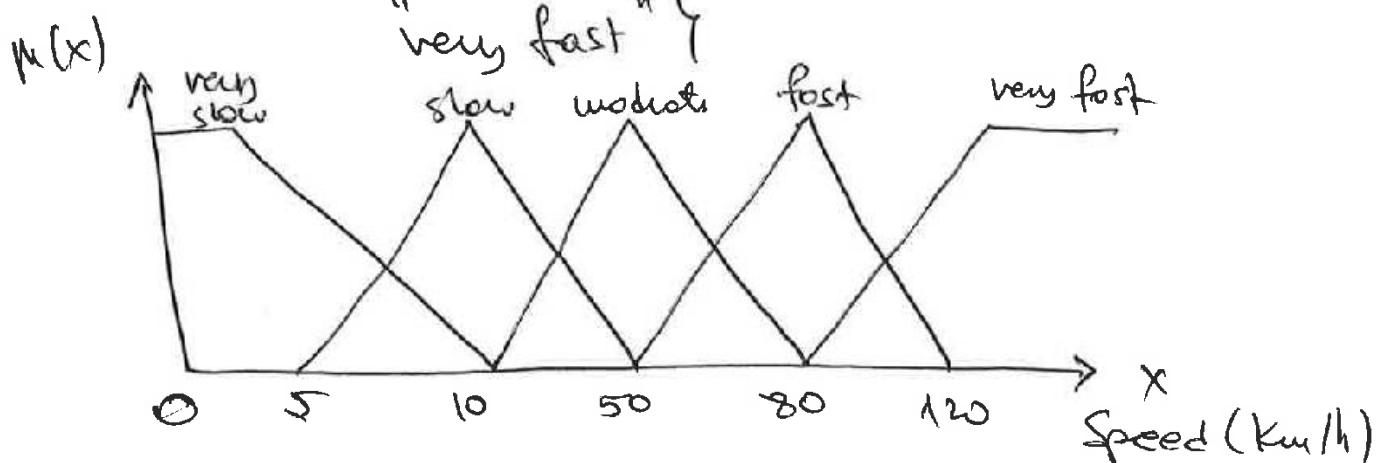


Linguistic variables

- * identifiers for membership functions.
- * define partitions of the universe of discourse.
- * depend on the application.

E.g. a variable "speed" for a control system

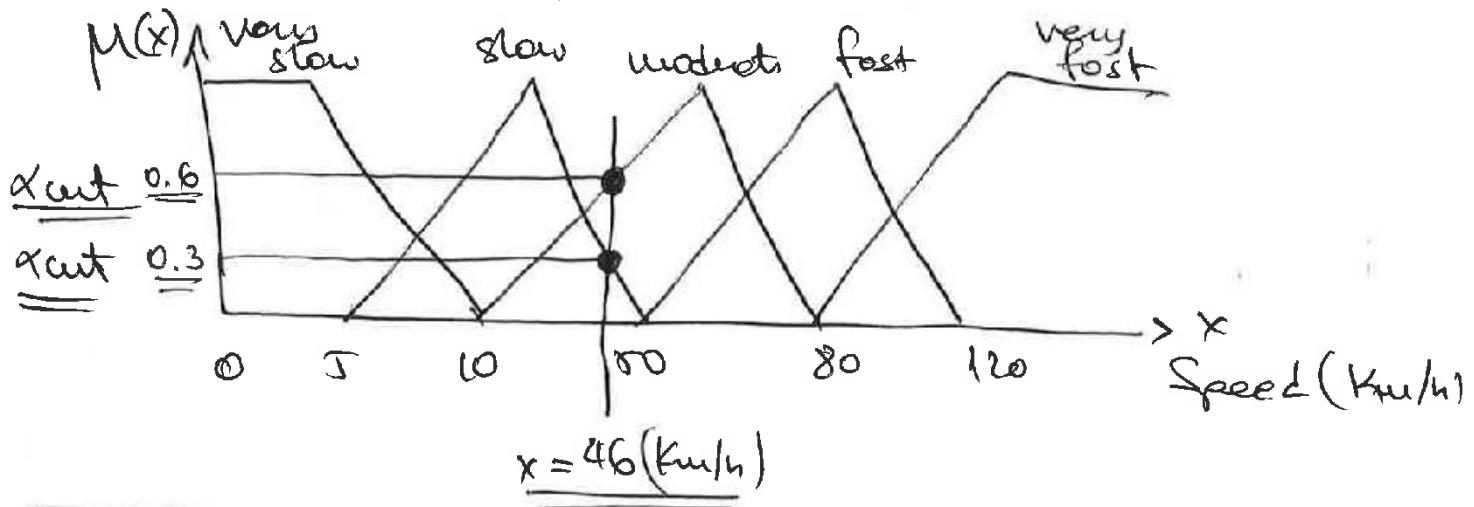
$$S = \{ \text{"very slow", "slow", "moderate", "fast", "very fast"} \}$$



- Linguistic variables \Rightarrow computing with words

- assign a membership degree for a crisp value for each label of a linguistic variable.

- E.g. Speed partitioning for an autonomous car.



Fuzzification: finds degrees of membership for crisp values:

$$\mu_{\text{very slow}}(x=46) = 0.0 \quad (\text{no intersection with } \underline{\text{very slow MF}})$$

$$\mu_{\text{slow}}(x=46) = 0.3$$

$$\mu_{\text{modest}}(x=46) = 0.6$$

$$\mu_{\text{fast}}(x=46) = 0.0 \quad (\text{no intersection with } \underline{\text{fast MF}})$$

$$\mu_{\text{very fast}}(x=46) = 0.0 \quad (\text{no intersection with } \underline{\text{very fast MF}})$$

Rule evaluation

- rules in fuzzy logic are using first order logic:
- IF-THEN clauses.

Antecedent (Premise)

Consequent (Conclusion)

IF x is A AND y is B THEN z is C

Linguistic variables
(e.g. Speed, temperature, acceleration)

Fuzzy Operator
(AND, OR, NEG)

Labels

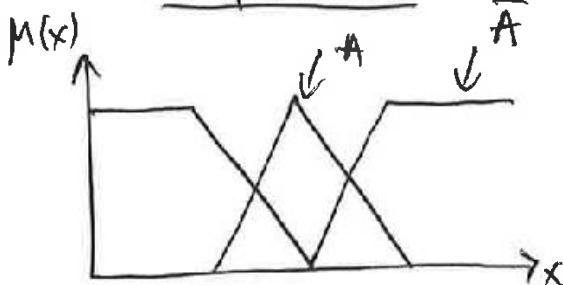
- if there are several input variables x_1, \dots, x_n with several connections (AND, OR) the IF-THEN clauses can be written as:

IF x_1 is A_{r_1} AND x_2 is A_{r_2} OR x_3 is A_{r_3} ... AND x_n is A_{r_n}
THEN y is B_r

where $r = 1 \dots n$ is the rule number

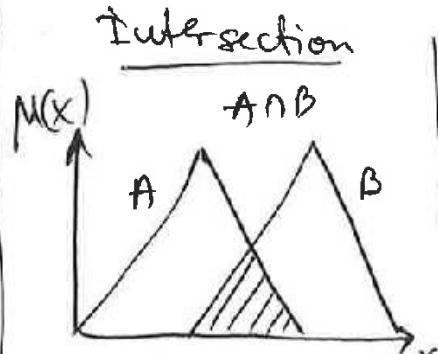
Basic operators

Complement



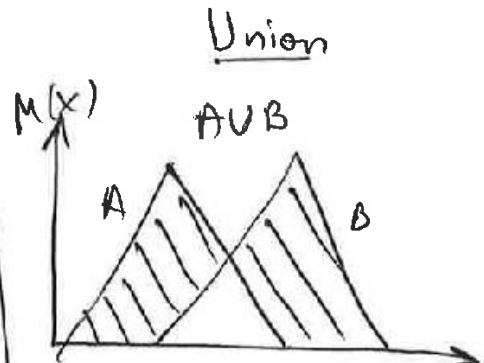
$$\mu_A(x) = 1 - \mu_{\bar{A}}(x)$$

Intersection



$$\begin{aligned} \text{MIN: } & \min\{\mu_A(x), \mu_B(x)\} \\ \text{PROD: } & \mu_A(x)\mu_B(x) \end{aligned}$$

Union

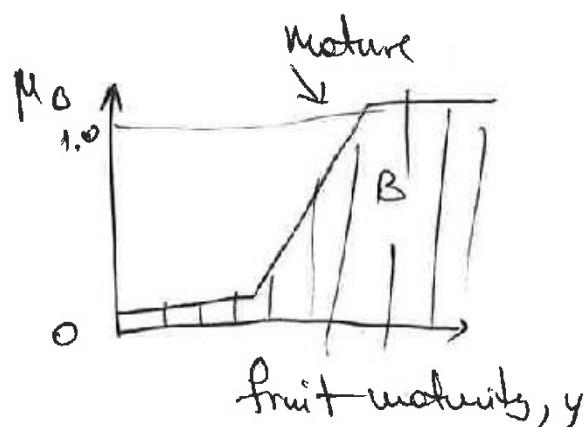
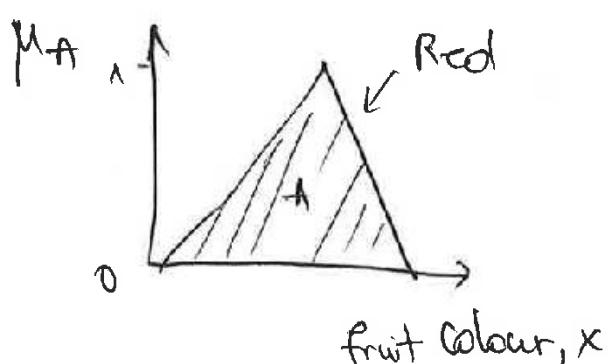


$$\begin{aligned} \text{MAX: } & \max\{\mu_A(x), \mu_B(x)\} \\ \text{ASUM: } & \mu_A(x) + \mu_B(x) - \mu_A(x)\mu_B(x) \end{aligned}$$

Fuzzy Implication

- * used to model dependencies, correlations or connections among variables.
- * generalization of a fuzzy set from 2D to 3D.
- * describes the degree of association among variables.

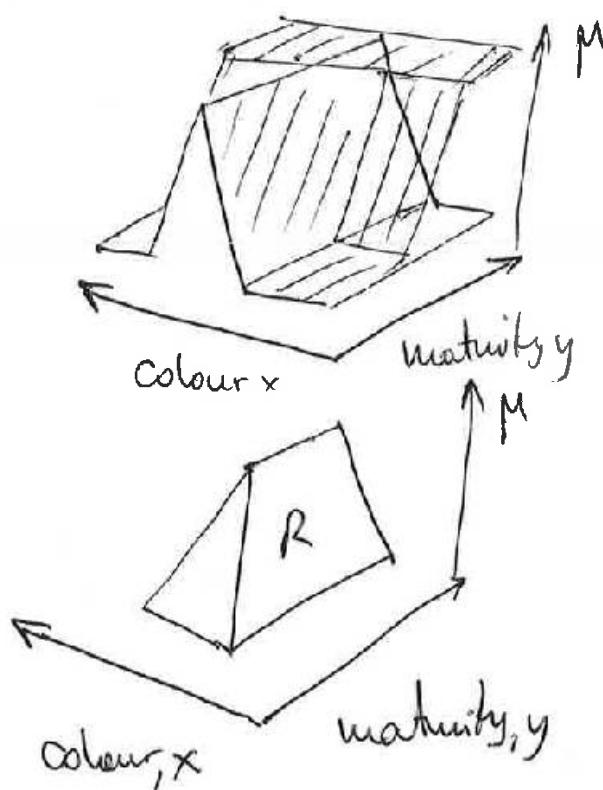
E.g. relation between color of a fruit x and the grade of maturity y



$$R: (x=A) \rightarrow (y=B)$$

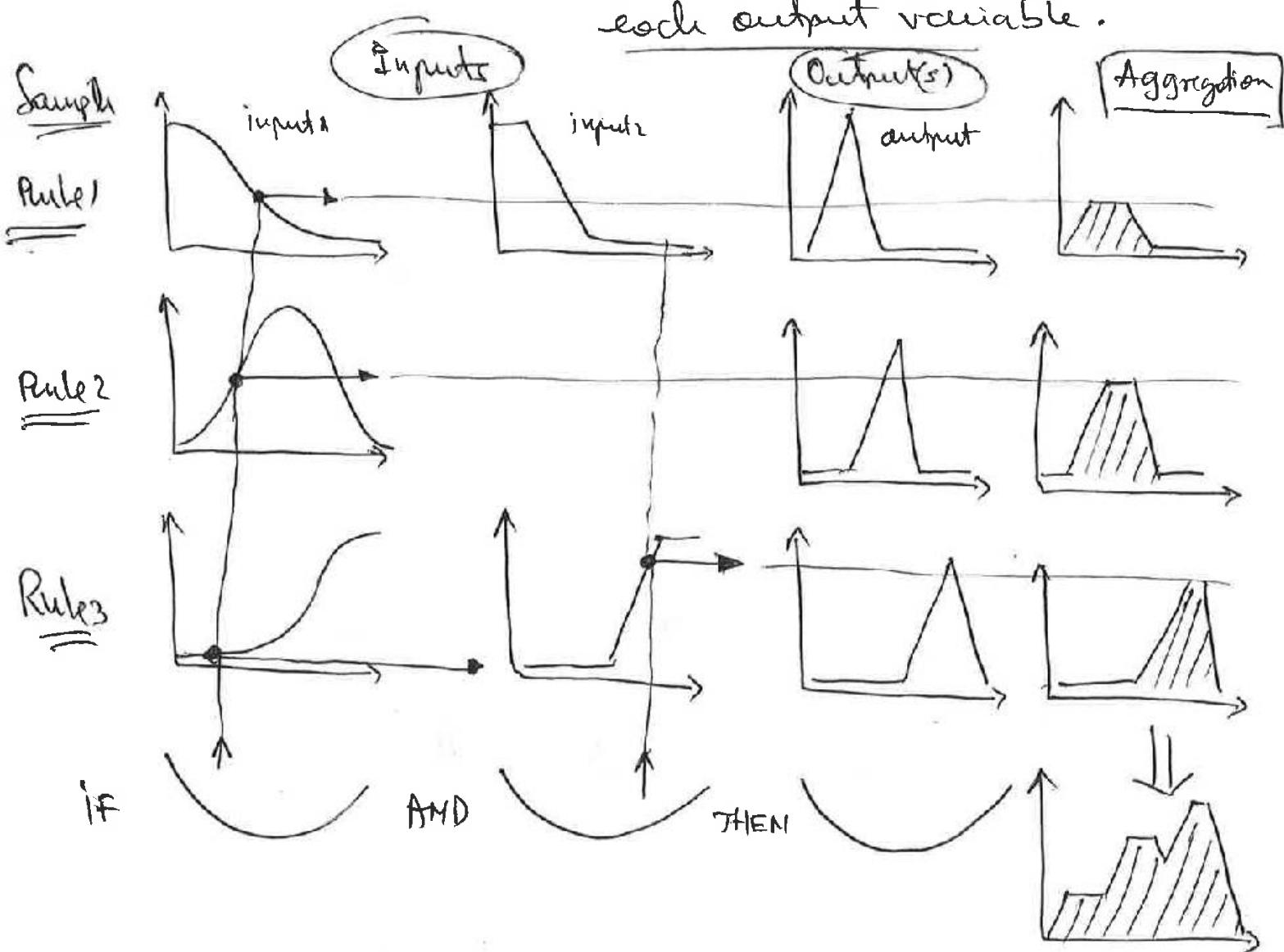
'IF x is A THEN y is B '

"if the colour is red then,
the fruit is mature."

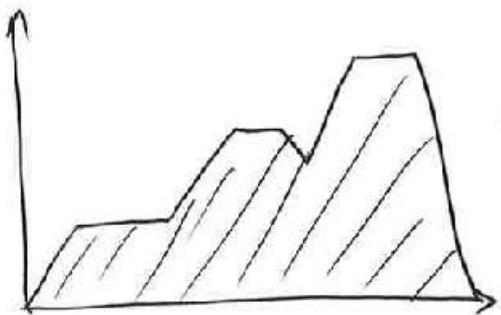


Fuzzy aggregation

- decisions are taken based on testing all the rules.
- rules must be combined for inference.
 - outputs are combined in a single fuzzy set.
 - once for each output variable.
 - the output is a fuzzy set for each output variable.



Defuzzification

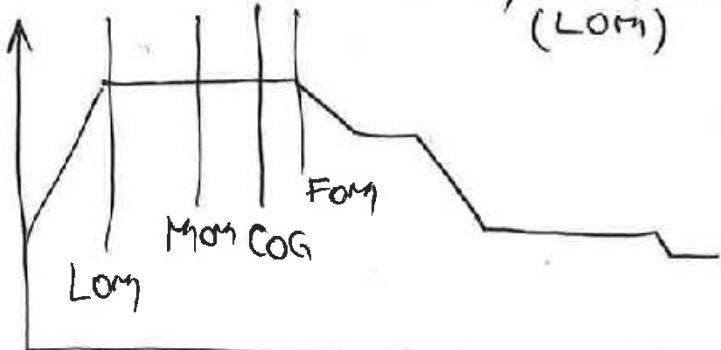


$$COG = \frac{\int_{x_{min}}^{x_{max}} x \cdot \mu_A(x) dx}{\int_{x_{min}}^{x_{max}} \mu_A(x) dx}$$

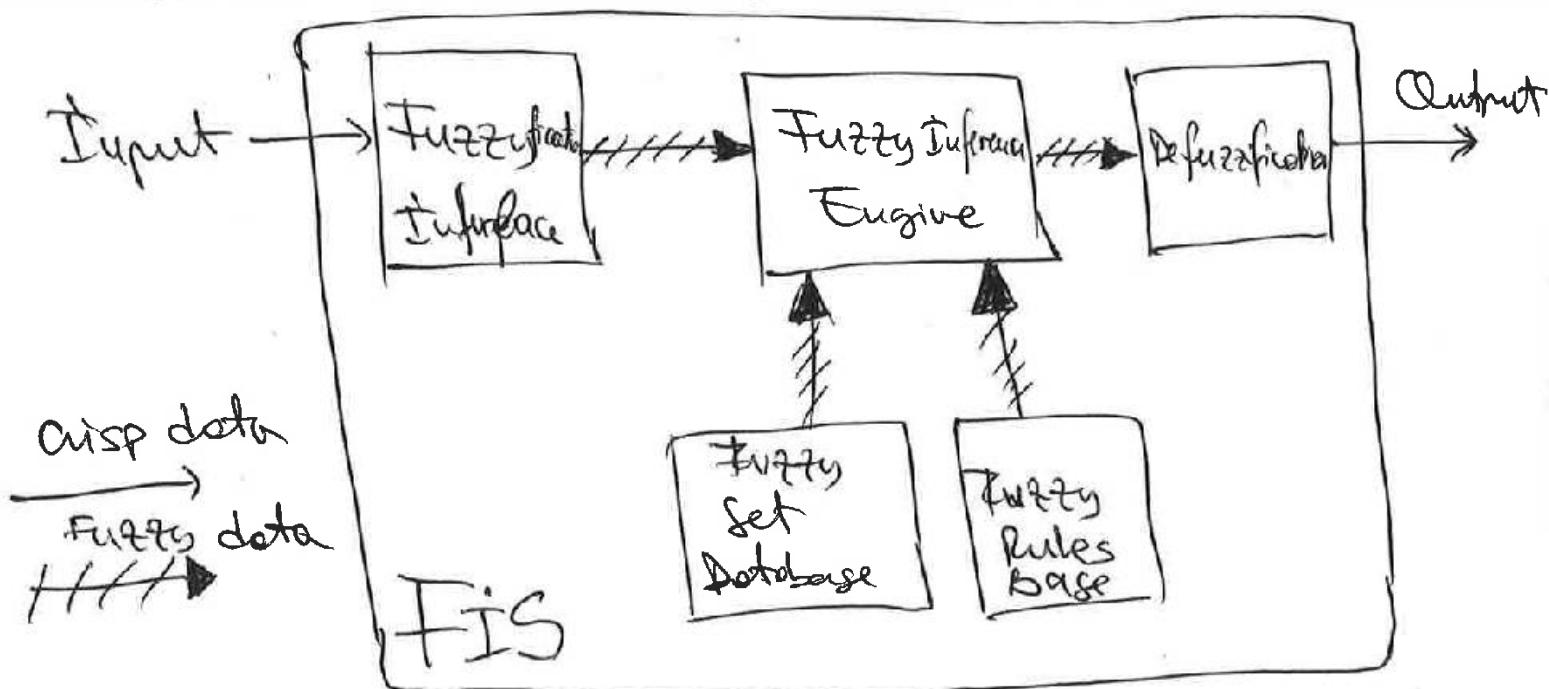
Example :
methods

- * after inference the result is a fuzzy value ;
- * the result needs to be mapped to cisp ;
- * done using the MF of the output variable,
- * various ~~mechanisms~~ :

center of gravity, center of maximum, mean of maxima, first of maxima (FOM), last of maxima (LOM)



From fuzzy logic to fuzzy inference systems



Worked example of FIS

Problem: Evaluate the risk of a software engineering project by referring only 2 criteria : the project funding (PF) and the project staffing (PS).

Step 1: Identify input/output variables

2 inputs: project funding (%) (PF)
project staffing (%) (PS)

1 output: project risk (%) (PR)

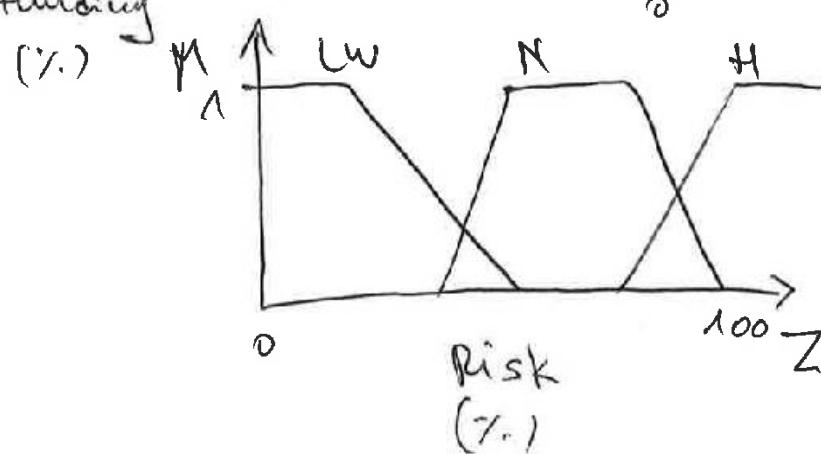
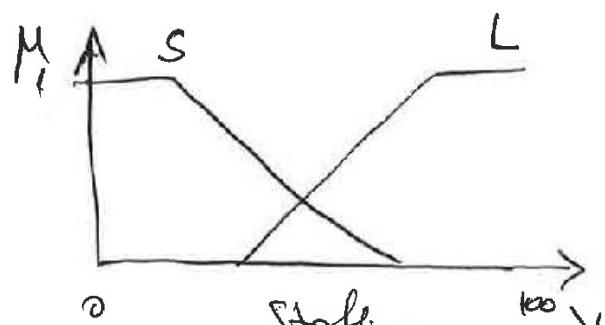
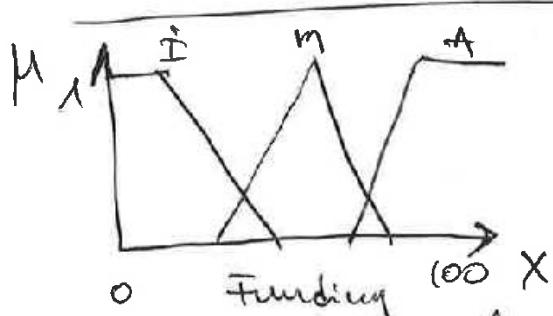
Universe of discourse \rightarrow fuzzy sets

$X = \{ \text{inadequate (I)}, \text{marginal (M)}, \text{adequate (A)} \}$ for PF

$Y = \{ \text{small (S)}, \text{large (L)} \}$ for PS

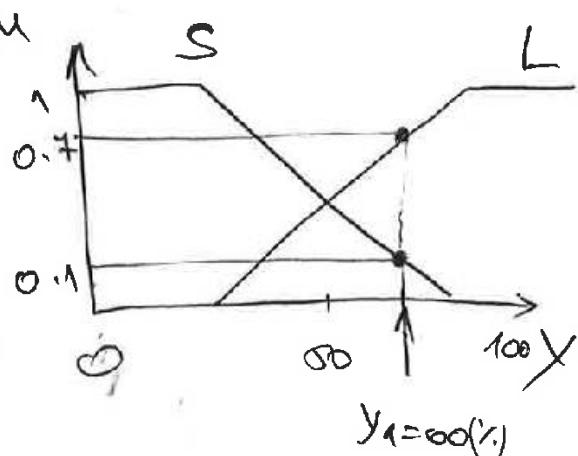
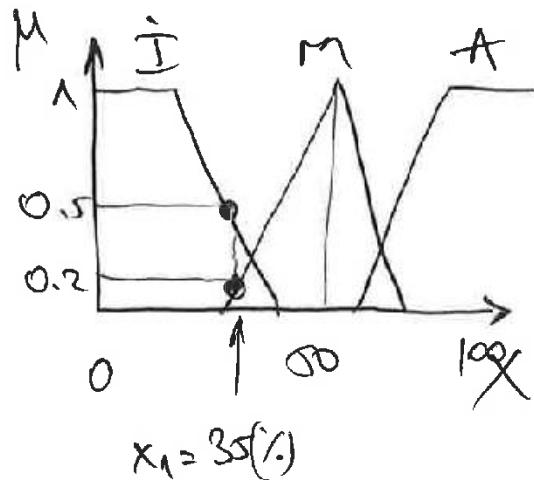
$Z = \{ \text{low (LW)}, \text{normal (N)}, \text{high (H)} \}$ for PR

Membership functions



Step 2: Fuzzification

Consider $x_1 = 35\% \text{ PF}$ and $y_1 = 60\% \text{ PS}$



$$\mu_I(x_1=35) = 0.5$$

$$\mu_S(y_1=60) = 0.7$$

$$\mu_M(x_1=35) = 0.2$$

$$\mu_L(y_1=60) = 0.1$$

$$\mu_A(x_1=35) = 0.0$$

Step 3: Rule design (here for example only 3)

Rule 1: IF PF is A OR PS is S THEN PR is L.

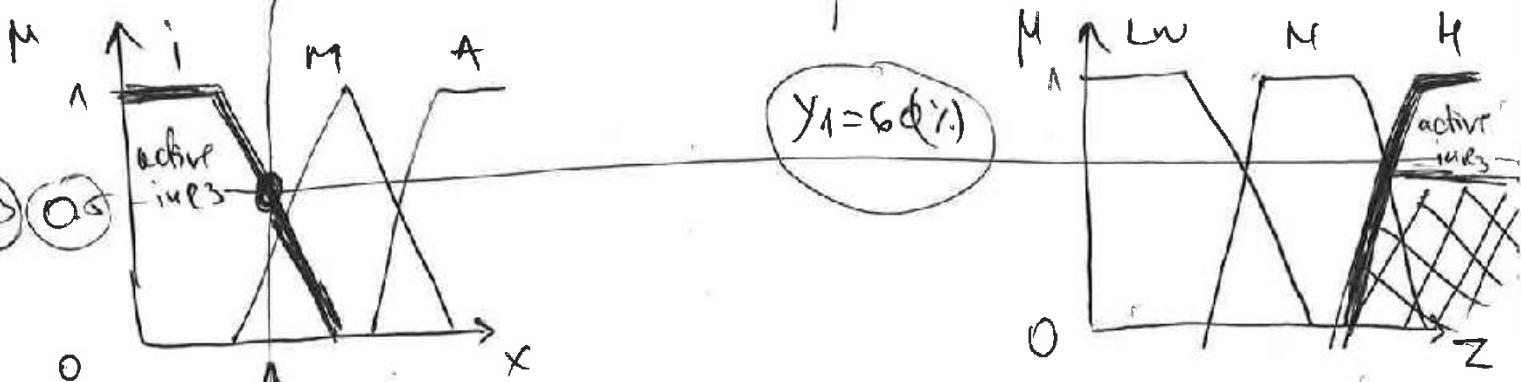
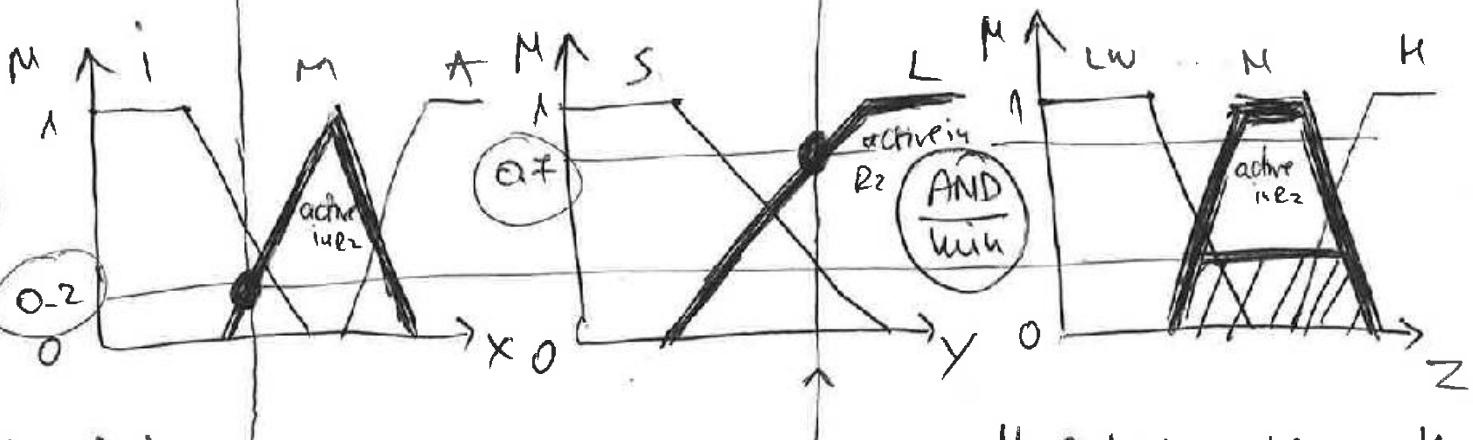
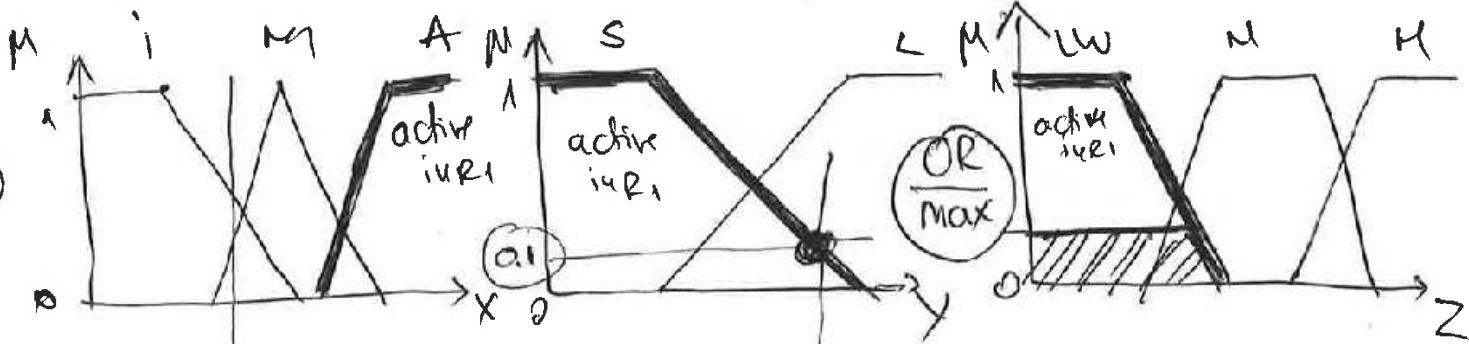
Rule 2: IF PF is M AND PS is L THEN PR is H.

Rule 3: IF PF is I THEN PR is H.

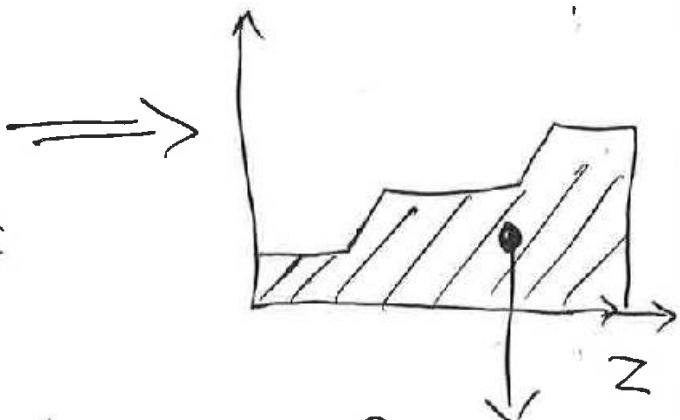
Observations: One can generate many rules to make sure one covers the entire space of solutions.

* Some rules are redundant, through logical conditions they mutually include or exclude.

Step 4 : Rule evaluation (inference)



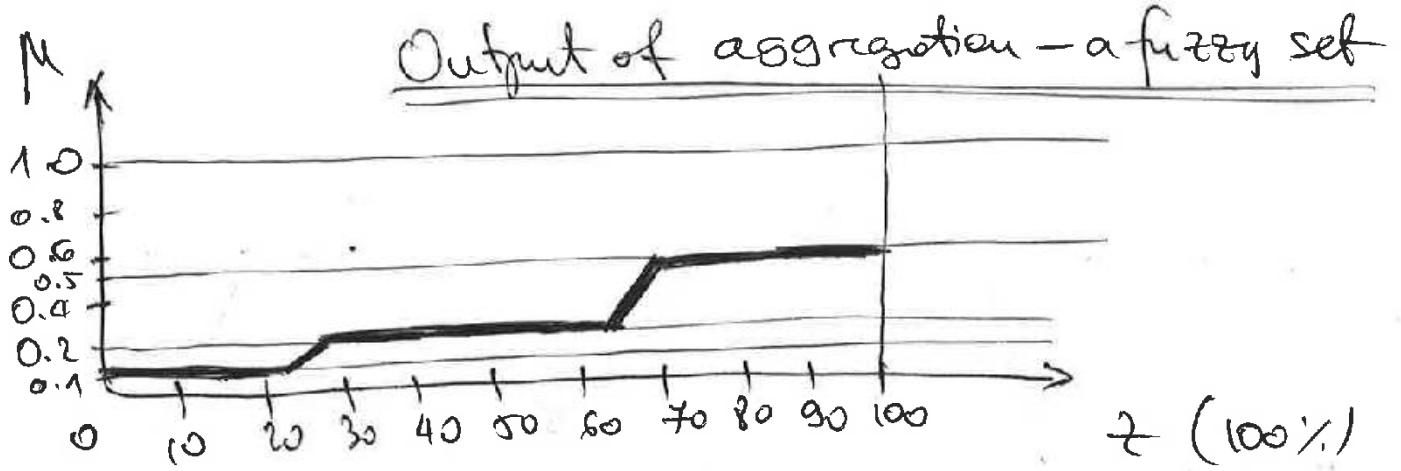
Step 5 : Aggregation (sum)



Step 6 : Defuzzification (using COG)

$$COG = \frac{\int_a^b x \mu_A(x) dx}{\int_a^b \mu_A(x) dx} \Rightarrow COG_{\text{discrete}} = \frac{\sum_{x=a}^b x \mu_A(x)}{\sum_{x=a}^b \mu_A(x)}$$

Continuous



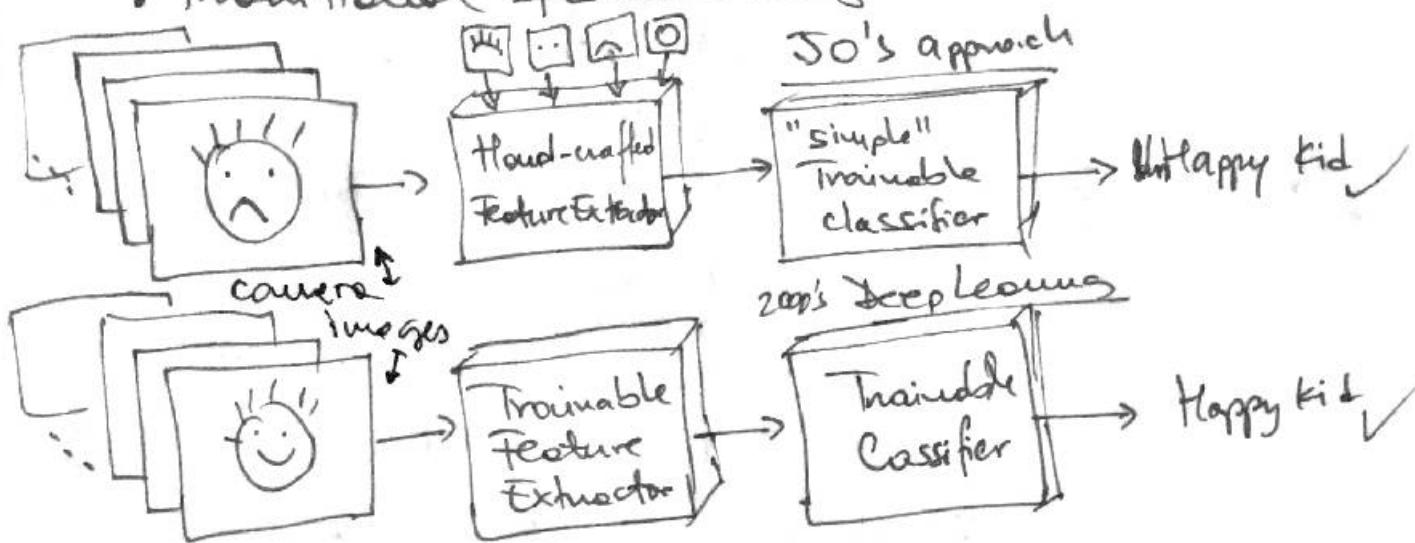
$$\text{COG} = \frac{(0+10+20) \times 0.1 + (30+40+50+60) \times 0.2 + (70+80+50+100) \times 0.5}{(0.1+0.1+0.1) + (0.2+0.2+0.2+0.2) + (0.5+0.5+0.5+0.5)}$$

$$\text{COG} = 67.9 \text{ (\%.)}$$

AIML Class: Deep (Neural) Learning

Deep Learning = Learning Representations/features

- Traditional pattern recognition



- Deep learning learns hierarchical representations
 - hierarchy with increasing level of abstraction
 - each stage is a trainable feature transformation
 - inspired by the hierarchical processing in the visual cortex with multiple pathways and stages (HMAX model)
 - low-level features are shared among categories
 - high-level features are more global and more invariant
- 3 Types of training:

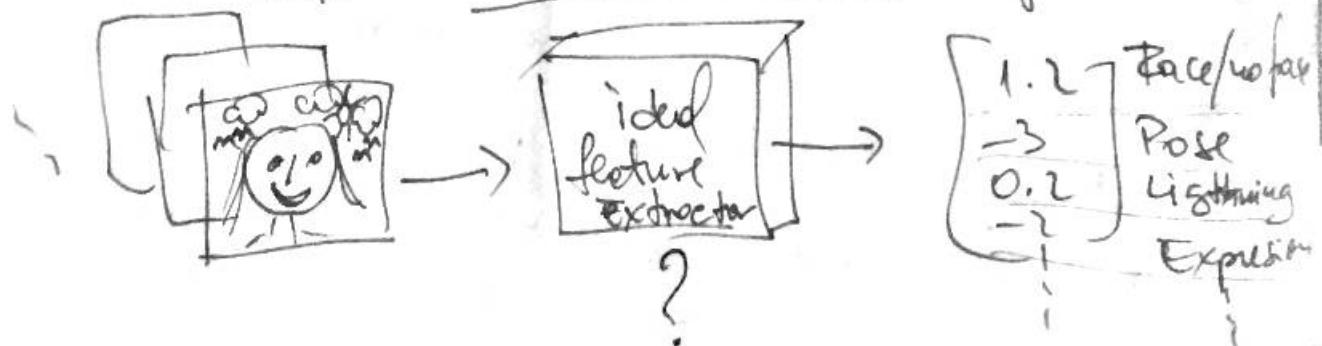
- (A). purely supervised:
 - init parameters randomly
 - train in supervised mode using backpropagation
 - e.g. Speech & image recognition

- (b) • Unsupervised, layerwise and a supervised classifier
- train each layer unsupervisedly
 - train a classifier on top, keeping all other layers fixed
- good when training set with labels are limited

- (c) • Unsupervised, layerwise and a global supervised fine-tuning
- train each layer unsupervisedly
 - add a classifier layer and retrain the whole network supervisedly
- good when labelset is poor

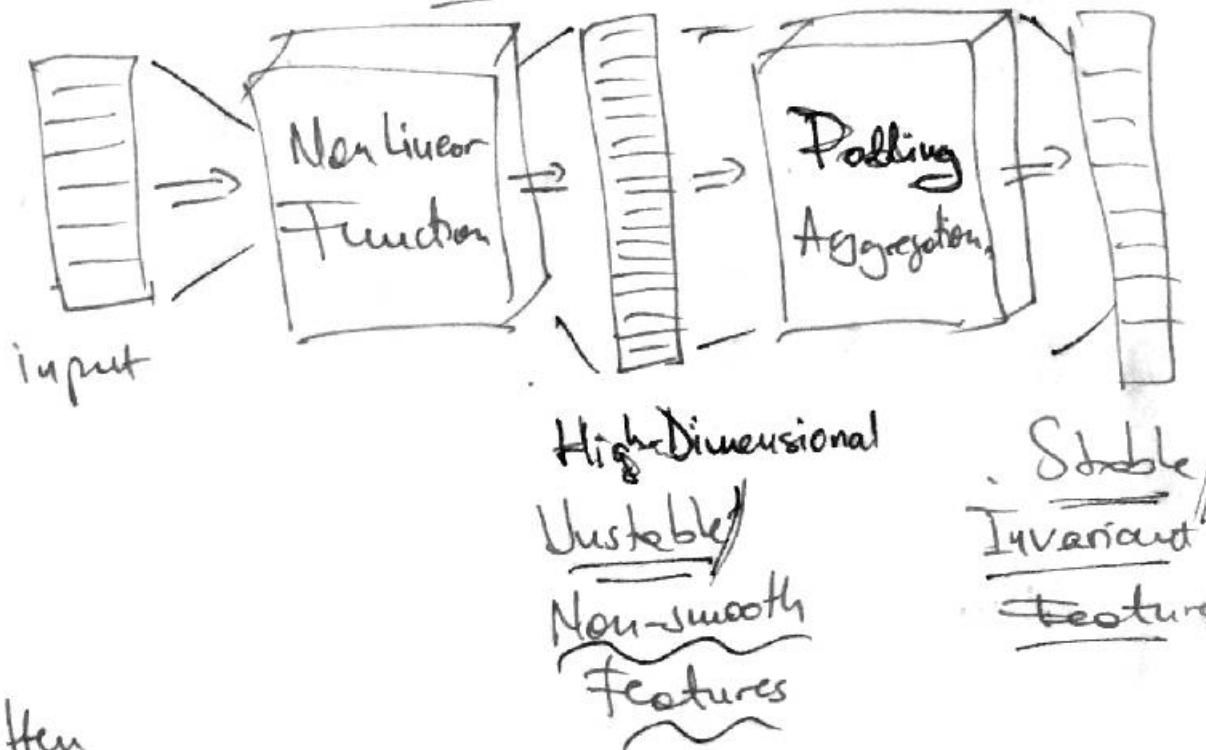
The manifold hypothesis

- natural data lives in a low-dimensional (nonlinear) manifold
- there is no general method to learn functions that turn images (for example) into a low-dimensional representation

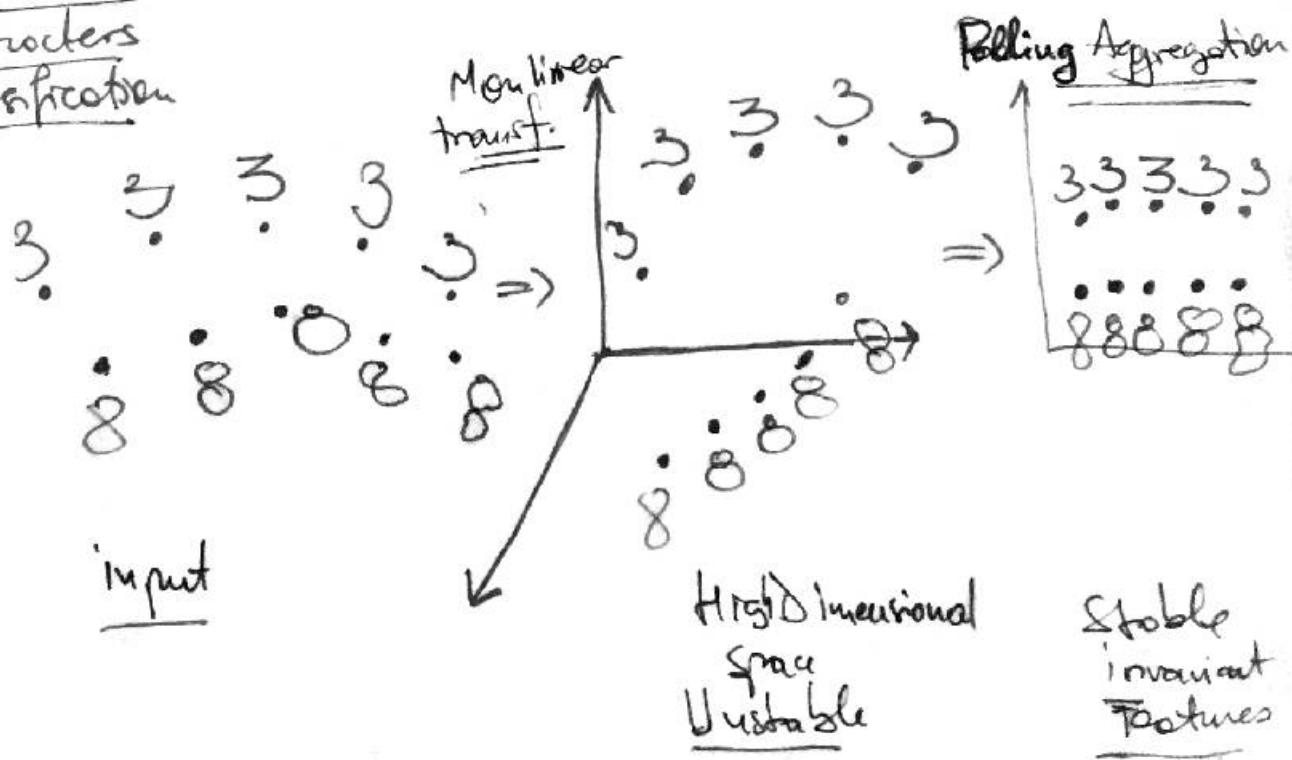


The inner workings of a deep learning network

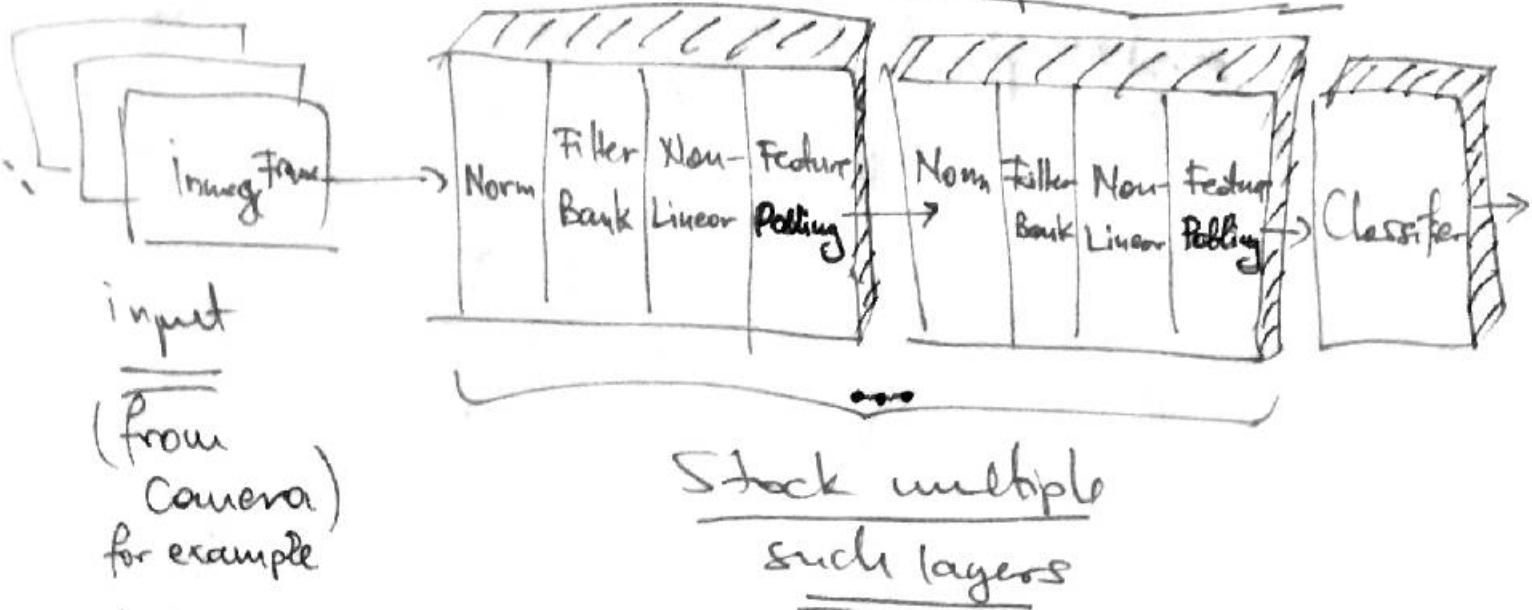
- embeds the input into a high(er) dimensional space
 - in this new space features that were not separable might become separable
 - bring together things that are semantically similar



E.g.: hand-written
characters
classification



Generic architecture of a Deep Network



Norm: - normalization: filtering, average removal;

Filter: - dimension expansion;

Bank:

Non: - activation function;

Linearity: - WTA (Winner-Take-All).

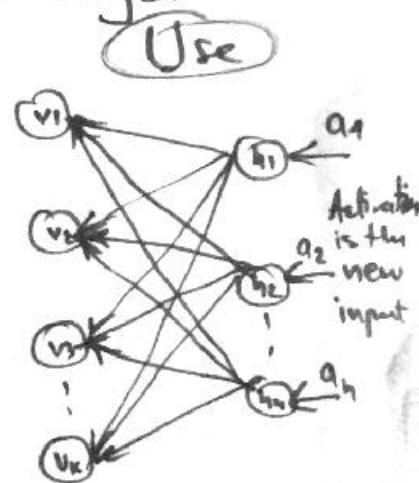
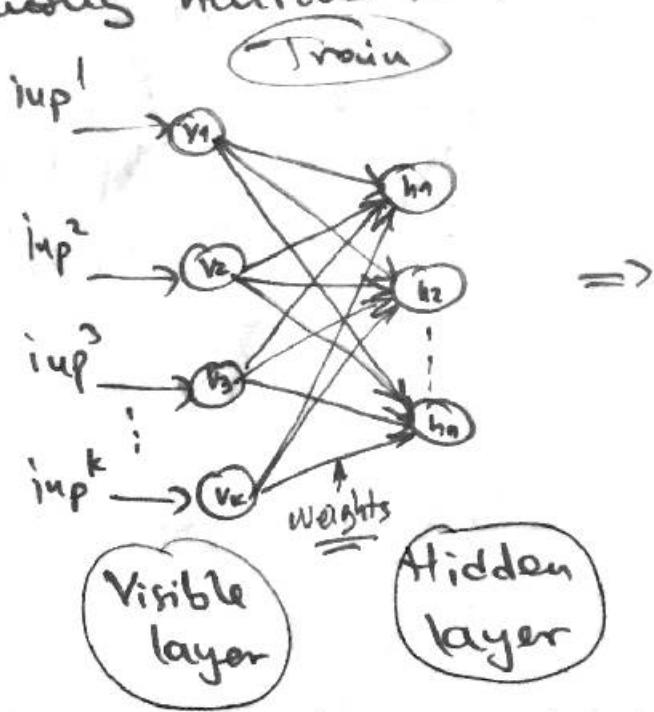
Pooling: aggregation over space or feature type.
(selection)

Deep Networks Building Blocks

Typically: Feedforward multi-layer neural nets → we discussed
 Now we look at: Restricted Boltzmann Machines (RBMs) in previous lectures
 ↳ Autoencoders

• RBMs (Restricted Boltzmann Machines)

- model probability distributions in data
- used for feature extraction
- used for dimensionality reduction
- the "restricted" comes from the fact that the network doesn't have connections among neurons in the same layer



Reconstruction
of the
input

- Use gradients in optimizing algorithm called contrastive divergence.
- Unsupervised learning.
- Usually size of input layer biggest size of hidden layer to create a bottleneck ⇒ a compact representation of the input space by extracting its structure.

(e.g. Handwritten
Recognition
MNIST
dataset)

- RBMs are energy based models, where the joint probability of visible (v) and hidden units (h) is proportional with the energy between them $E(v, h)$:

$$P(v, h) \propto e^{-E(v, h)}$$

- Manipulating the energy $E(v, h)$ we can generate the probability $P(v, h)$.

→ ~~we minimize the energy which maximizes the probability~~

→ this is done by learning the weight matrix W between input (visible) units and the hidden units (h):

$$E(v, h) = -v \cdot W \cdot h$$

→ Contrastive divergence algorithm is similar to gradient descent and assumes weight update:

$$W \leftarrow \lambda \cdot W \cdot cd(v_t)$$

where $cd(v_t)$ is the contrastive divergence of instance v_t and λ = learning rate

Unsupervised learning allows to learn distributions
(Remember SOM?!)

Observations:

- Iterative algorithm
- $cd(v_t)$ minimizes the difference between the real distribution of the data and the guess of the net.

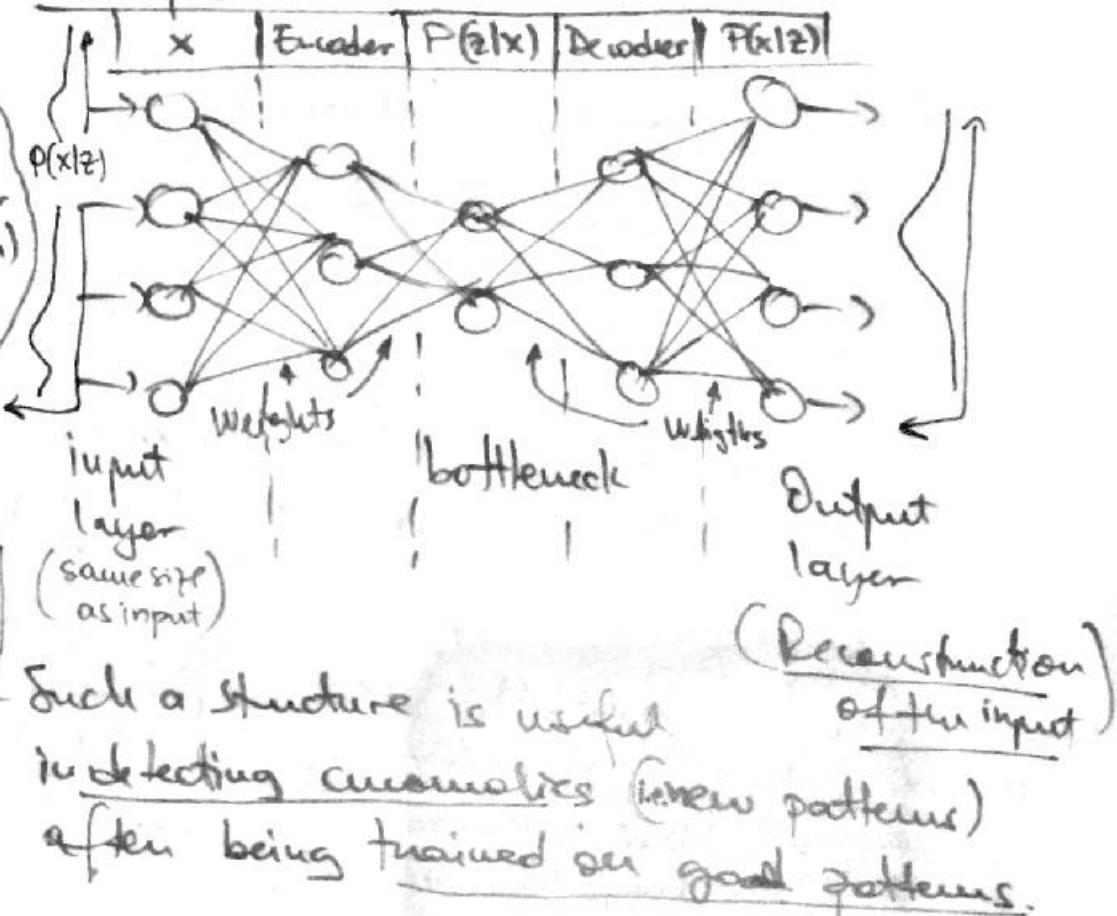
① Autoencoders

- a variant of feed forward neural nets.
- have an extra bias for calculating error of reconstructing the input.
- used to learn compressed representations of the data
 ⇒ their output is an efficient construction of the input
- difference from MLP is that autoencoders have an output layer with same size as input layer

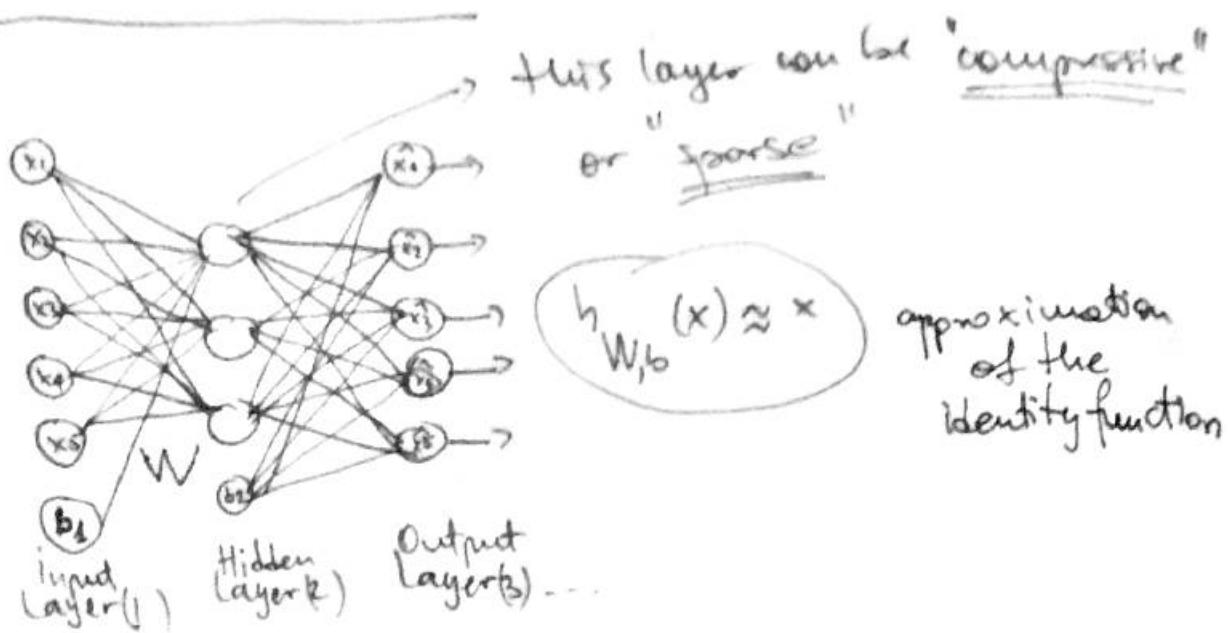
looking in depth:
 x = data instance
 $z \sim p(z)$ prior
 $x \sim p(x|z)$, generator (data distribution)
 The net learns an estimate of $p(z|x)$ to decode $P(x|z)$.

Learn the distribution of the data

- Sample network:



Autoencoders - add-ons



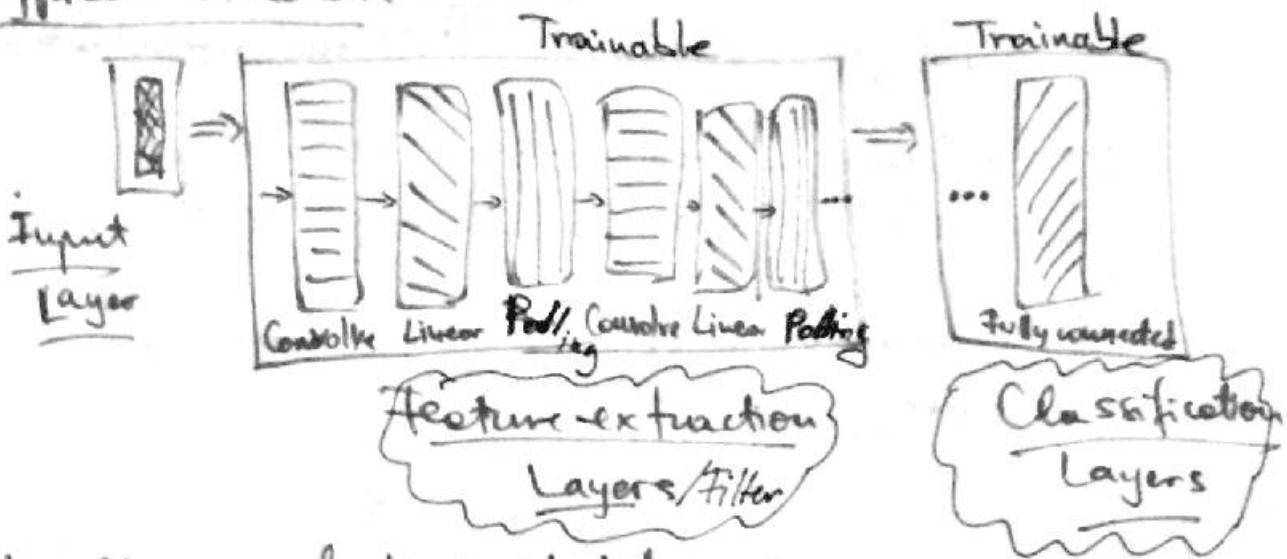
- $\{x_1, x_2, x_3, x_4, x_5\}$ - unlabeled set of training examples, there is no y_i
- the Autoencoder applies back propagation setting the targets to be equal to the inputs $y_i = x_i$
- it tries to learn unsupervisedly a function $h_{W,b}(x) \approx x$ which is an approximation to the identity function so that \hat{x} is similar to $x \Rightarrow$ but isn't this trivial??

Yes, but \Rightarrow due to the network structure ("bottleneck" in hidden layer) the network learns a "compressed" representation of the input: e.g. 10×10 px image as input, so $n=100$ input neurons, consider $k=50$ hidden neurons and $l=100$ output neurons \Rightarrow if there are correlations in the input the autoencoder(s) extracts that (similar to PCA)

* Convolutional Neural Networks (CNNs)

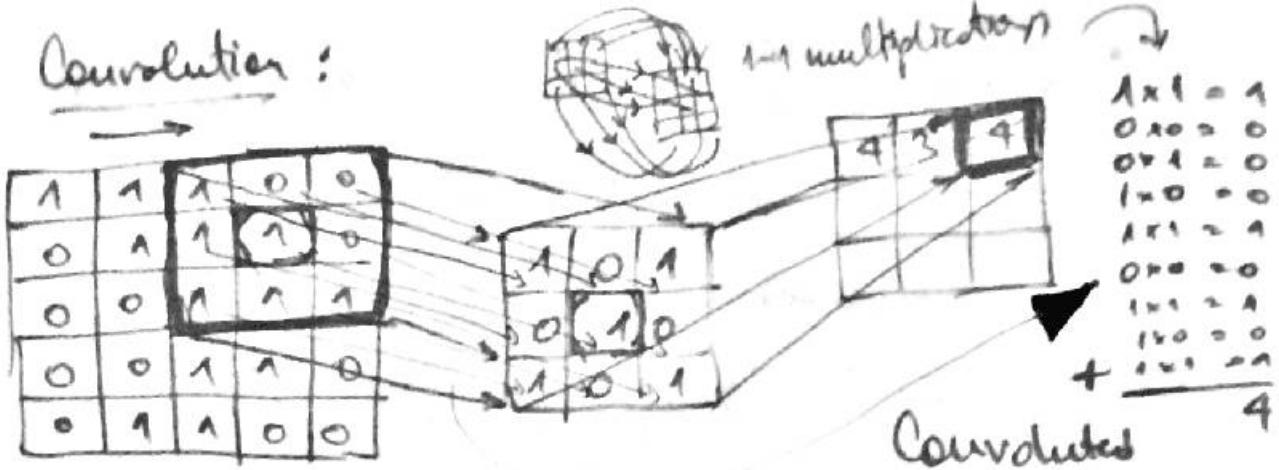
- learn high-level features using convolution
- used for object recognition & image classification
- build position & rotation invariant features
learn
out of the training data (it learns the structure).
 - ↳ So, ideal input are images and sounds with repeating patterns.
- supervised learning

Typical Structure:



- input layers : feed input data
- Convolution layers :- core of the CNN
 - dot product between the region of the neurons in the input layer and the weights to which they are locally connected in the output layer
 - feature detector layer
 - interpreted as a filter

Convolution:



- Input data kernel

(e.g. image B&W pixels)

Principle of convolution

Auto-element of the kernel (y_1) is placed over the source pixel (x_1). The source pixel is then replaced with a weighted sum of itself & nearby pixels. \rightarrow (Kernel \times pixel) = dot product

- Pooling layers: reduce data representation progressively and control overfitting, can implement different operations (mean, max...).
- fully connected layers: compute class scores.

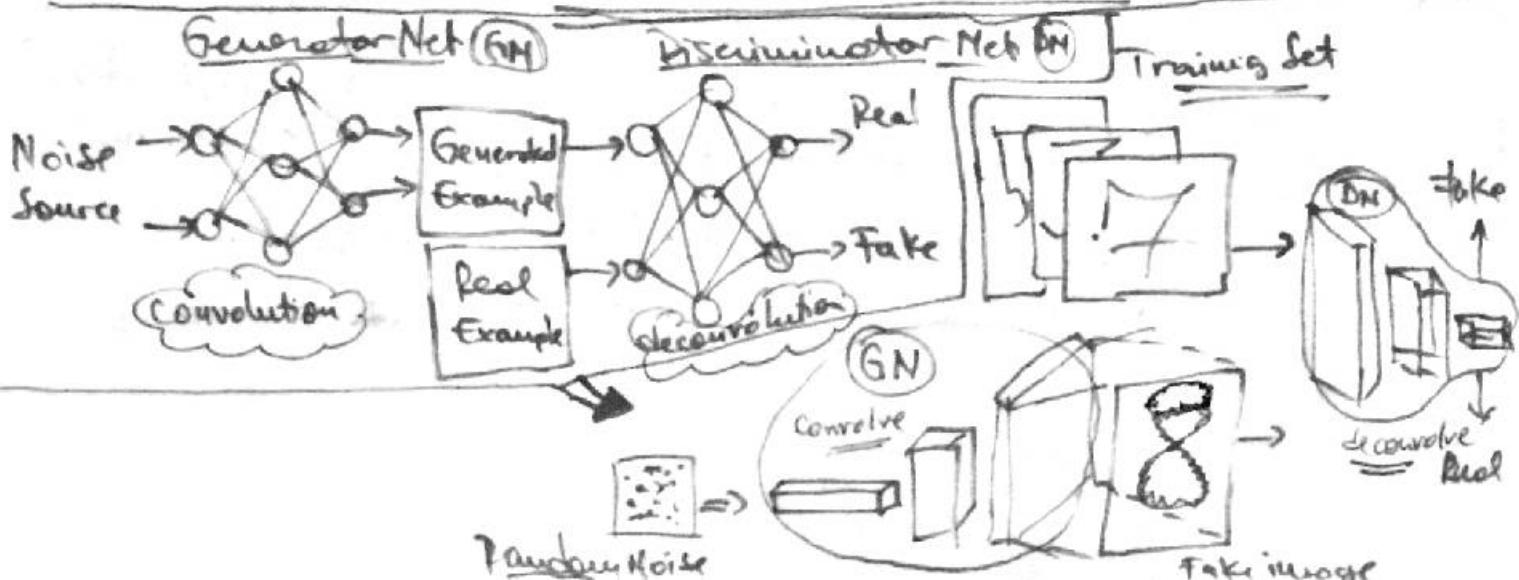
Major Deep Networks Architectures

• Deep Belief Networks (DBN)

- composed of layers of RBMs for the pretrain phase and feed-forward net for fine tuning phase
 - can extract high-level features of raw data
 - the RBMs learn the features unsupervisedly & reconstructs the input → feed into the feed-forward part is trained with backprop to make predictions on the labels (softmax).

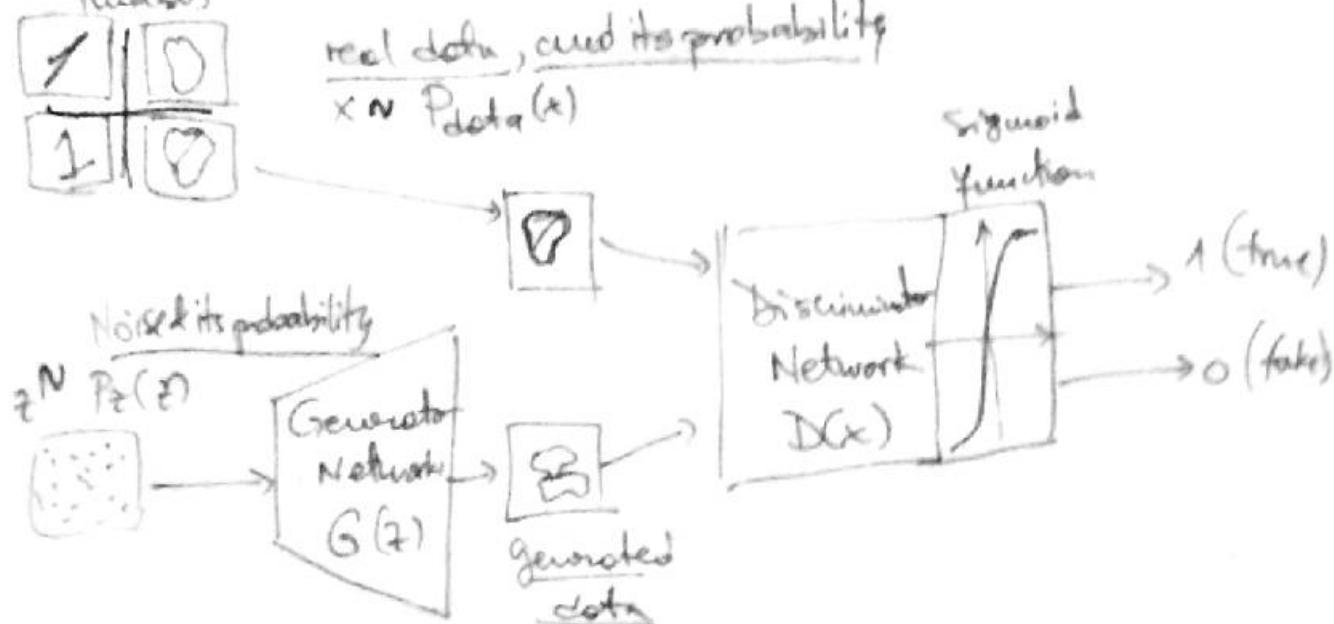
- Generative Adversarial Networks (GAN)

- used to synthesize (generate) new images out of other training examples
 - unsupervised learning to train two models in parallel
 - efficient representation of the input data
 - composed of a generative network and a discriminative network:



Generative Adversarial Networks - add-on

e.g.: handwritten characters
images



Training a GAN is a "fight" between generator & discriminator, or in math:

can be also seen as cross entropy

$$\min_{G} \max_{D} V(D, G)$$

wherever game theory? player minimax criteria

$$V(D, G) = E_{x \sim P_{\text{data}}(x)} [\log D(x)] + E_{z \sim P_z(z)} [\log (1 - D(G(z)))]$$

entropy that the data from real distribution passes through discriminator

(best case)
 D tries to maximize to 1

entropy that the data from random source passes through discriminator worst case (D) tries to minimize to 0

(D tries to minimize to 0)

- The discriminator tries to maximize V whereas the generator (G) tries to minimize it (i.e. minimum difference between real)

Training phases for GAN

Pass 1: Train discriminator & freeze generator
just a forward pass and no backpropagation
of information:

$D(x) \rightarrow$ update

$G(x) \rightarrow$ just compute output (no update)

Pass 2: Train generator & freeze discriminator

$D(x) \rightarrow$ just get output (no update)

$G(x) \rightarrow$ update

More explicit steps:

1. Train D on real data

2. Generate false inputs for generator G
and train D on fake data

3. Train G with the output of D

4. Repeat #1 to #3 for a number of epochs

If you would have a fully functional generator
you could duplicate almost anything: generate
news, books and novels with unimaginable stories
etc. :)) It would become real!

Additive rule

The entropy H of a discrete random variable X with possible values x_1, \dots, x_n and probability mass function $P(x)$ is:

$$H(X) = E[\bar{I}(X)] = E[-\ln(P(X))]$$

$$H(X|Y) = -\sum_{i,j} P(x_i, y_j) \log \frac{P(x_i, y_j)}{P(y_j)}$$

The amount of randomness or uncertainty in the variable X given the event Y .

Conditional entropy
operator

Expectation information content of X

$$H(X) = \sum_{i=1}^n P(x_i) I(x_i) = -\sum_{i=1}^n P(x_i) \log P(x_i)$$

b - base of logarithm

e.g. $\begin{cases} b=2 & \text{bits (Fibonacci number)} \\ b=e & \text{nats} \\ b=10 & \text{bans} \end{cases}$

Cross entropy

$$H(p, q) = E[-\log q] = H(p) + D_K(p||q)$$

$$H(p, q) = -\sum_x P(x) \log q(x)$$

$$D_K(p||q) = \sum_i P(i) \log \frac{P(i)}{q(i)}$$

AUROC & AUPRC

Precision

Reality			
		F	T
F	TN	FN	
T	FP	TP	

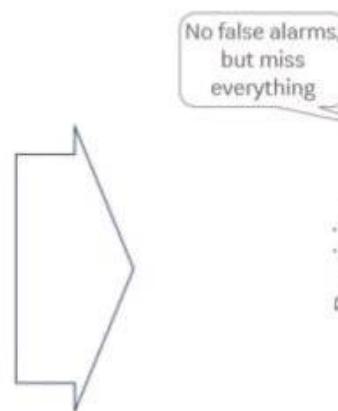
What fraction of the predicted true are actually true?

Recall

Reality			
		F	T
F	TN	FN	
T	FP	TP	

What fraction of the actual true are predicted to be true?

AUPRC



Precision

Recall

Note that TN is not considered in AUPRC, so it may be unhelpful when the TN rate is low... when there is low prevalence

Does my predictor catch most of the true positives without causing lots of false positives?

Catch all, but lots of false alarms

True Positive Rate

Reality			
		F	T
F	TN	FN	
T	FP	TP	

What fraction of the actual true are predicted to be true?

True Negative Rate

Reality			
		F	T
F	TN	FN	
T	FP	TP	

What fraction of the actually negative are predicted to be negative.

AUROC



True Positive Rate

False Positive Rate

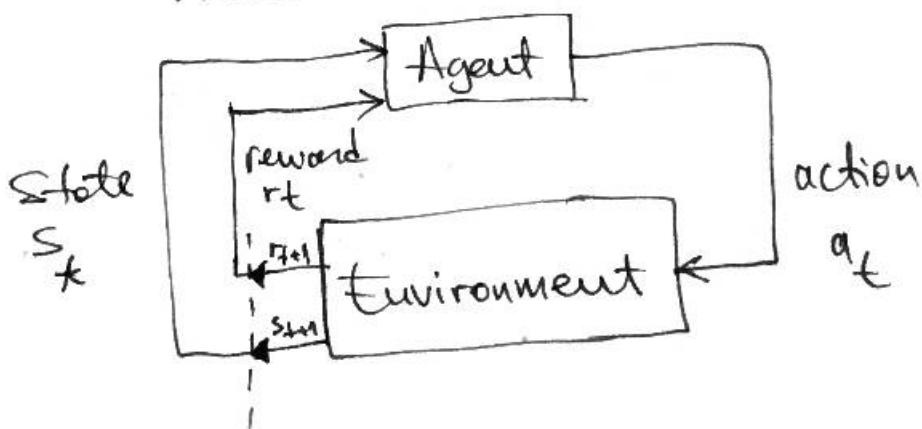
Catch all, but lots of false alarms

Does my predictor provide good results for positive and negative predictive value?

No false alarms, but miss everything

AIML Course : Reinforcement Learning

- * "If an action is followed by a satisfactory state of affairs, the tendency of the system to produce that action is strengthened or reinforced".
- * Richard Sutton, 1991 - inventor of Reinforcement Learning.
- * initially framed in adaptive optimal control.
- * solution for solving non-linear control problems for which there is no excessive resource (memory) or CPU-time consumption but rather the designer doesn't know how to program "the correct things/actions"; e.g. pole balancing, airplane stabilization.
- * RL has been seen as a derivative of supervised learning based on trial-and-error (and reward):
 - but there is no teacher;
 - output is quantified as positive or negative reward if system is closer or further from the goal respectively.
- * Basic formulation & elements:



1. Environment
2. Reward
3. Policy
4. Value function

- Net description:
- The policy: shows how to choose a good action for a given state;
 - The value function: shows how good / valuable a state is;
 - The reward: shows how much reward does being in a certain state bring;
- (RL) searches a mapping from state to action by trial-and-error.

The environment:

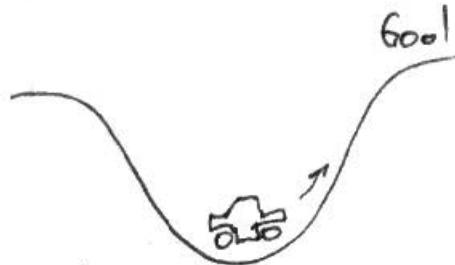
- must be observable through sensory readings & actions;
- must be deterministic, such that same actions lead to same results in repeated trials;
- these constraints are relaxed in more complex algorithms.

The reward:

- the reward $r(s)$, or the reinforcement signal depicts a mapping from state to action.
- the environment gives reward to the learning system (agent).
- the agent tries to maximize all expected future reward!
- There are at least 3 types of reward:
 - a. pure delayed reward: all reward is zero until the end state is reached and the sign of the final reward indicates success or failure: e.g. backgammon, inverted pendulum

b. minimum time to goal reward:

e.g. get a car uphill, the reward would be -1 at every step and 0 once the car reached the goal.



c. multi-player games:

→ two or more agents work simultaneously to achieve potentially opposing goals: e.g. predator-prey scenario.

The policy:

- the policy $\pi(s)$ is responsible to map from state to actions to be taken, this needs to be learnt → how to choose a good action for a given state.

The value function:

- the value function $V(s)$ shows how good / valuable a state is.
- the value is defined as the sum of the expected future rewards when starting from a state following a policy $\pi(s)$.
- the value can be represented as: a look-up table, equation or even a neural network (e.g. Deep Mind). System playing GO).

There are 2 models for value functions:

1. Finite horizon model:

$$V(s) = \sum_{t=1}^n r(s(t)), \text{ where } s(t+1) = \pi(s(t))$$

2. Infinite horizon model:

$$V(s) = \sum_{t=0}^{\infty} \gamma^t \cdot r(s(t)), \text{ where } s(t+1) = \pi(s(t))$$

discount factor

$$\gamma^t \in [0, 1]$$

Such problems can also be solved with Markov Decision Processes (MDP).

- these models assume the complete state of the world is available to the agent \rightarrow this is unrealistic.

An extension are the Partially Observable MDP or PoMDP.

- model the information available to the agent by using a function from the hidden states to the observables,
- the goal is to find a mapping from observations (not states!) to actions

In real problems:

- environments contain non-determinism
- agent's actions can fail
- sensors can be inaccurate
- MDPs can help to take decisions under uncertainty

↓
formulate the total reward!

We need to formulate the total reward from a policy as the sum of the discounted expected utility of each state visited by that policy.



The optimal policy is the one maximizing this equation!

There are several methods to solve this problem.

1. Value iteration

- if we know the true value of each state choose the action that maximizes the utility.
- but we don't know a state's initial value, just its immediate reward
- let's assume we have a look-up table of values $r(s)$ and we sweep through the table to update v'_t , the value at next step:

$$v'_{t+1} = \max_{\pi(s)} (r(s(t+1)) + \gamma V(s(t+1)))$$

or incrementally

$$\Delta v'_{t+1} = \max_{\pi(s)} (r(s(t+1)) + \gamma V(s(t+1)))$$

- the value iteration algorithm finds the optimal policy by taking into account long term reward $r(s)$
- but we need to know the results of all actions ($\max_{\pi(s)}$) in a given state and needs long time to converge

The algorithm for value iteration:

- #1. initialize $v(s)$ randomly for all states except target (goal).
- #2. initialize V with target value for the goal.
- #3. repeat
 - #4. initialize $\Delta = 0$
 - #5. for all $s \in S$ (states) in random order
 - calculate the difference $v - V(s)$
 - #6. calculate value $V(s) = \max_{\pi(s)} (r(s(t+1)) + \gamma^t V(s(t+1)))$
 - #7. calculate difference for update $\Delta = \max(\Delta, |v - V(s)|)$
 - #8. end for
 - #9. until $\Delta = 0$ or $\Delta \leq 0$

2. Continuous states, Residual Gradient Algorithm

- RL converges (guaranteed) for look-up tables yet is unstable for function approximation systems which can contain even small amounts of generalization.
- we can see the value function as $V(s) = V^*(s) + e(s)$ where $V^*(s)$ is the "perfect" value of the state and $e(s)$ is the error.
- a solution is to use Neural Networks to approximate $V(s)$! \implies

\Rightarrow such that

$$V^*(s(t)) = V(s(t), w(t)) = \text{NeuralNet}(s(t), w(t))$$

- the change in weights of the NN will be:

$$\Delta w(t) = -\eta \left(\max_{\pi(s)} (r(s(t+1)) + \gamma^t V(s(t+1))) - V(s(t)) \right) \frac{\partial V(s(t), w(t))}{\partial w(t)}$$

where

$$-\eta \left(\max_{\pi(s)} (r(s(t+1)) + \gamma^t V(s(t+1))) - V(s(t)) \right)$$

produces an error for current V

and

$$\frac{\partial V(s(t), w(t))}{\partial w(t)}$$
 is the change in V with respect to the weights (backprop)

- this method allows to "interpolate" between states and use RL even when not all states are explored

Final remarks:

— value iteration has weaknesses:

— \star long convergence time

— \star it computes the value of each state which we don't need, that's just useful to find the optimal policy.

— value iteration can be adopted to policy iteration.

Q-Learning

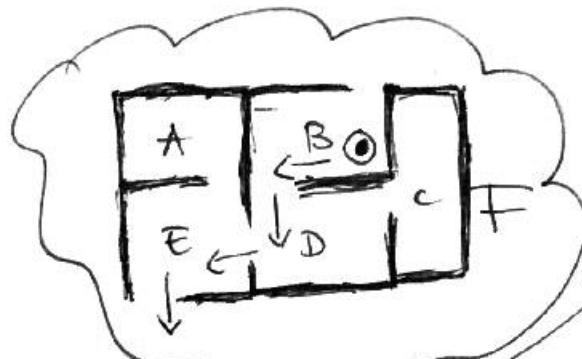
- value/policy iteration work well in finding an optimal policy but they assume that the agent has a great deal of domain knowledge.
- trade-off between learning time and a priori knowledge: Q-learning.
- model-free RL algorithm.
- agent doesn't have any model of the environment.
- agent has only knowledge of what states exist and what actions are possible in each state.
- The core idea is to map the tuple (state,action) to reward.
- We define a function $Q(s,a)$ as a reward for a given action in a state plus all future rewards along optimal actions:

$$Q(s,a) = r(s,a) + \gamma \max Q(s',a')$$

- this formulation allows the agent to "perform" only states $s(t)$ and $s(t+1)$ and not a large number of s_{t+1} .

Let's introduce the algorithm through an example

- a search and rescue robot is supposed to start in any of the rooms and find the best way to reach the outside world (#) from that room.



if γ^* close to 1
introduce delay
and wait for
greater reward

If γ^* close to 0
consider only
immediate
reward.

Algorithm :

- # 1. set parameter γ and reward matrix R
- # 2. initialize Q_1 with \emptyset
- # 3. for each period of exploration (episode)
 - # 4. select random initial state
 - # 5. do while not reached the goal (#)
 - select one among all possible actions for the current state
 - using this action, consider going to next state
 - get max Q value of this next state for all possible actions
 - calculate $Q(s, a) = r(s, a) + \gamma \max Q(s^*, a)$
 - set next state as current state
 - # 6. end

- again for continuous states we have a MLP to estimate the change in $Q(s, a)$:

$$\Delta Q(s(t), a(t)) = (r(s(t), a(t)) + \gamma \max Q(s(t+1), a(t+1)) - Q(s(t), a(t)))$$

and the change in weights is

$$\Delta w(t) = -\eta \Delta Q(s(t), a(t)) \frac{\partial Q(s(t), a(t))}{\partial w(t)}$$

A more general method to learn the value function and the function is TD-Learning or temporal difference learning.

- assume we have a RL system in which a sequence of actions is deterministic:



- we init states randomly.
- "true" information flows from right to left, 1000 steps until "true" information arrives at state 1.
- The formulation will be:

$$Q^*(s(t), a(t)) = r(s(t), a(t)) + \gamma \max Q(s(t+1), a(t+1))$$

$$Q^2(s(t), a(t)) = r(s(t), a(t)) + \gamma^1 r(s(t+1), a(t+1)) + \gamma^2 \max Q(s(t+2), a(t+2))$$

$$Q^3(s(t), a(t)) = r(s(t), a(t)) + \gamma^1 r(s(t+1), a(t+1)) + \gamma^2 r(s(t+2), a(t+2)) + \gamma^3 \max Q(s(t+3), a(t+3))$$

$$Q^n(s(t), a(t)) = r(s(t), a(t)) + \gamma^1 r(s(t+1), a(t+1)) +$$

$$+ \gamma^2 r(s(t+2), a(t+2)) +$$

$$+ \gamma^{n-1} r(s(t+n-1), a(t+n-1)) +$$

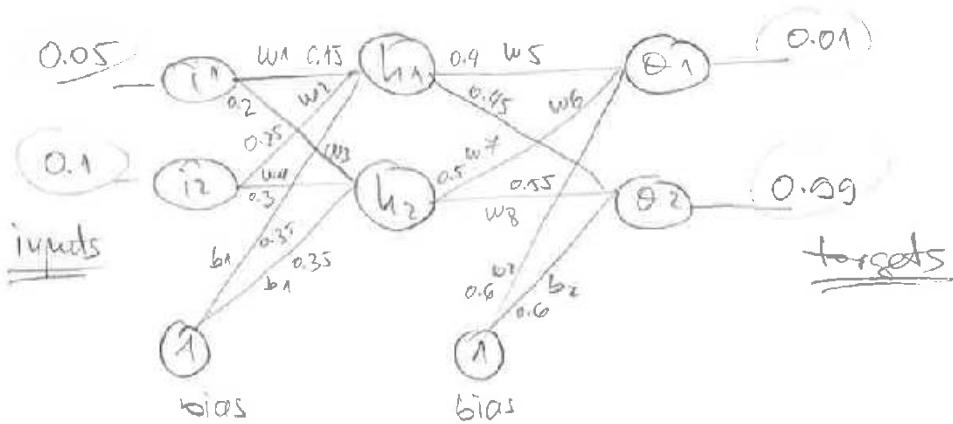
$$+ \gamma^n \max Q(s(t+n), a(t+n))$$

→ TD Learning learns a guess from a guess, if bootstraps.

→ TD Learning are naturally implemented online (e.g. important in robotic systems) and fully incrementally.

Back-prop. worked example on MLP

1



Parameters

Activation function in hidden and output neurons is logistic

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

$$\text{Learning rate } \eta = 0.5$$

Forward propagation

$$\text{net}_{h_1} = w_1 \cdot i_1 + w_2 \cdot i_2 = 0.05 \cdot 0.15 + 0.1 \cdot 0.25 = 0.3825$$

$$\text{out}_{h_1} = f(\text{net}_{h_1}) = \frac{1}{1+e^{-0.3825}} = \frac{1}{1.6821} = 0.6142$$

$$\text{net}_{h_2} = w_3 \cdot i_1 + w_4 \cdot i_2 = 0.05 \cdot 0.15 + 0.1 \cdot 0.25 = 0.3825$$

$$\text{out}_{h_2} = f(\text{net}_{h_2}) = \frac{1}{1+e^{-0.3825}} = \frac{1}{1.6770} = 0.5963$$

$$\text{out}_{\theta_1} = \frac{1}{1+e^{-\text{net}_{\theta_1}}}$$

$$\text{net}_{\theta_1} = w_5 \cdot \text{out}_{h_1} + w_6 \cdot \text{out}_{h_2} + b_1 \cdot 1 = 1.1438$$

$$\text{out}_{\theta_1} = \frac{1}{1+e^{-1.1438}} = 0.7583$$

$$\text{net}_{\theta_2} = w_7 \cdot \text{out}_{h_1} + w_8 \cdot \text{out}_{h_2} + b_2 \cdot 1 = 1.2093$$

$$\text{out}_{\theta_2} = \frac{1}{1+e^{-1.2093}} = \frac{1}{1.2099} = 0.7692$$

$$E_{\text{total}} = \sum \frac{1}{2} (\text{target} - \text{out})^2 \quad (\text{we use MSE})$$

$$E_{\text{total}} = E_{\theta_1} + E_{\theta_2}$$

$$E_{\theta_1} = \frac{1}{2} (\text{target}_{\theta_1} - \text{out}_{\theta_1})^2 = \frac{1}{2} \cdot (0.01 - 0.7583)^2 = 0.2799$$

Hidden layer

Output layer

Error out

$$E_{O2} = \frac{1}{2} (target_{O2} - out_{O2})^2 = \frac{1}{2} (0.09 - 0.7692)^2 = 0.0293$$

$$E_{\text{total}} = E_{O1} + E_{O2} = 0.2795 + 0.0293 = 0.3092$$

(2)

Backward propagation

We only look at w5 weight update in output layer

"chain" rule $\frac{\partial E_{\text{total}}}{\partial w5} = \boxed{\frac{\partial E_{\text{total}}}{\partial out_{O1}}} \cdot \boxed{\frac{\partial out_{O1}}{\partial net_{O1}}} \cdot \boxed{\frac{\partial net_{O1}}{\partial w5}}$

1.) $\frac{\partial E_{\text{total}}}{\partial out_{O1}} = \frac{\partial}{\partial out_{O1}} \left[\frac{1}{2} (target_{O1} - out_{O1})^2 + \frac{1}{2} (target_{O2} - out_{O2})^2 \right]$

$$\frac{\partial out_{O1}}{\partial net_{O1}}$$

no dependence on
out_{O2}

$\frac{\partial E_{\text{total}}}{\partial out_{O1}} = \frac{1}{2} \cdot 2 \cdot (target_{O1} - out_{O1}) \cdot (-1) + 0$

$$\frac{\partial out_{O1}}{\partial net_{O1}}$$

$\frac{\partial E_{\text{total}}}{\partial net_{O1}} = -(target_{O1} - out_{O1}) = -(0.01 - 0.7583) = 0.7483$

2.)

$$\frac{\partial out_{O1}}{\partial net_{O1}} = \frac{\partial f(\text{net}_{O1}(t))}{\partial \text{net}_{O1}} = f'(\text{net}_{O1}(t))$$

The derivative of the activation function $f(x) = \frac{1}{1+e^{-x}}$

is $\frac{d}{dx} f(x) = f(x)(1-f(x))$

$$\frac{\partial out_{O1}}{\partial net_{O1}} = out_{O1}(1-out_{O1}) = 0.7583(1-0.7583) = 0.1832$$

$$3.) \text{D}_{\text{net}h_1} = w_5 \cdot \text{out}_{h_1} + w_6 \cdot \text{out}_{h_2} + b_2 \cdot 1 \quad (3)$$

$$\frac{\partial \text{D}_{\text{net}h_1}}{\partial w_5} = 1 \cdot \text{out}_{h_1} + 0 + 0 = \text{out}_{h_1} = 0.6142$$

so $\frac{\Delta E_{\text{total}}}{\partial w_5} = 0.7483 \cdot 0.1832 \cdot 0.6142 = 0.0841$

And so the full weighted update for w_5

$$w_5(t+1) = w_5(t) - \eta \cdot \frac{\Delta E_{\text{total}}(t)}{\partial w_5(t)} = 0.4 - 0.5 \cdot 0.0841$$

$$\underline{w_5(t+1) = 0.3579}$$

similar to w_6, w_7, w_8

hidden layer weights update

$$\frac{\Delta E_{\text{total}}}{\partial w_i} = \frac{\Delta E_{\text{total}}}{\partial \text{out}_{h_i}} \cdot \frac{\partial \text{out}_{h_i}}{\partial w_i}$$

$$1.) \frac{\Delta E_{\text{total}}}{\partial \text{out}_{h_1}} = \frac{\Delta E_{O_1}}{\partial \text{out}_{h_1}} + \frac{\Delta E_{O_2}}{\partial \text{out}_{h_1}}$$

$$\frac{\Delta E_{O_1}}{\partial \text{out}_{h_1}} = \frac{\Delta E_{O_1}}{\partial w_{\text{net}h_1}} \cdot \frac{\partial \text{net}_{O_1}}{\partial \text{out}_{h_1}} = \frac{\Delta E_{O_1}}{\partial w_{\text{net}h_1}} \cdot \frac{\Delta \text{out}_{O_1}}{\partial \text{net}_{O_1}} \cdot \frac{\Delta \text{net}_{O_1}}{\partial \text{out}_{h_1}}$$

$$0.7483 \qquad 0.1832$$

$$\frac{\partial \text{out}_{01}}{\partial \text{out}_{h1}} = \frac{\partial [w_5 \cdot \text{auth}_1 + w_6 \cdot \text{auth}_2 + b_2 \cdot 1]}{\partial \text{out}_{h1}} = w_5 = 0.9 \quad (\text{we use time } +)$$

$$\frac{\partial E_{01}}{\partial \text{out}_{h1}} = 0.7 \times 33 \cdot 0.1732 \cdot 0.9 = 0.0548 \quad (4)$$

Similarly $\frac{\partial E_{02}}{\partial \text{out}_{h1}} = -0.0160$

$$\frac{\partial E_{\text{total}}}{\partial \text{out}_{h1}} = \frac{\partial E_{01}}{\partial \text{out}_{h1}} + \frac{\partial E_{02}}{\partial \text{out}_{h1}} = 0.0548 + (-0.0160) = 0.0388$$

(2.) $\frac{\partial \text{out}_{h1}}{\partial \text{out}_{h1}} = \text{out}_{h1} (1 - \text{out}_{h1}) = 0.6142 (1 - 0.6142) = 0.2369$

3.) $\frac{\partial \text{out}_{h1}}{\partial w_1} = \frac{\partial [w_1 \cdot i_1 + w_2 \cdot i_2 + b_2 \cdot 1]}{\partial w_1} = i_1 = 0.05$

Combining

$$\frac{\partial E_{\text{total}}}{\partial w_1} = 0.0388 \cdot 0.2369 \cdot 0.05 = 0.0004$$

weight update

$$w_1(t+1) = w_1(t) - \eta \cdot \frac{\partial E_{\text{total}}(t)}{\partial w_1(t)} = 0.15 - 0.5 \cdot 0.0004 = 0.1498$$

Similar w_2, w_3, w_4

(5)

For inputs 0.05 and 0.1 the net error was 0.3042. After running backprop of 1 step the error is 0.2940. If we repeat 10000 times error becomes 0.000035.

After 10000 epochs of training if we feed 0.05 and 0.1 the net will output 0.015 (vs. 0.01 target) and 0.084 (vs. 0.05 target).

Supervised neural computation

AI ML
Class

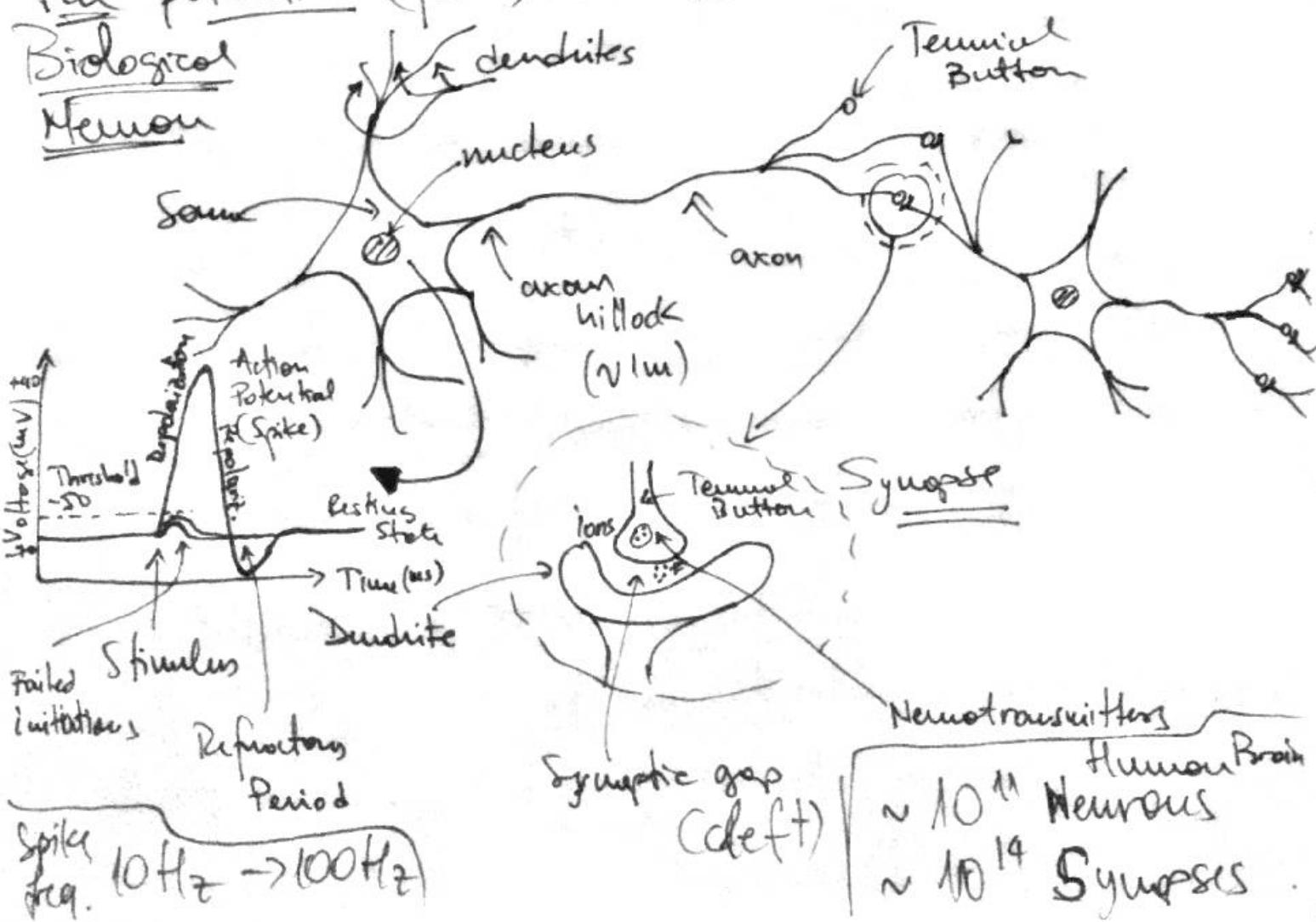
Introduction to neural networks

Neural systems have the ability to:

- extract meaning from complex, imprecise & often noisy data; ($N \sim 10^{14}$ of sensory input bands)
- learn patterns & trends;
- generalize & respond to unexpected inputs;

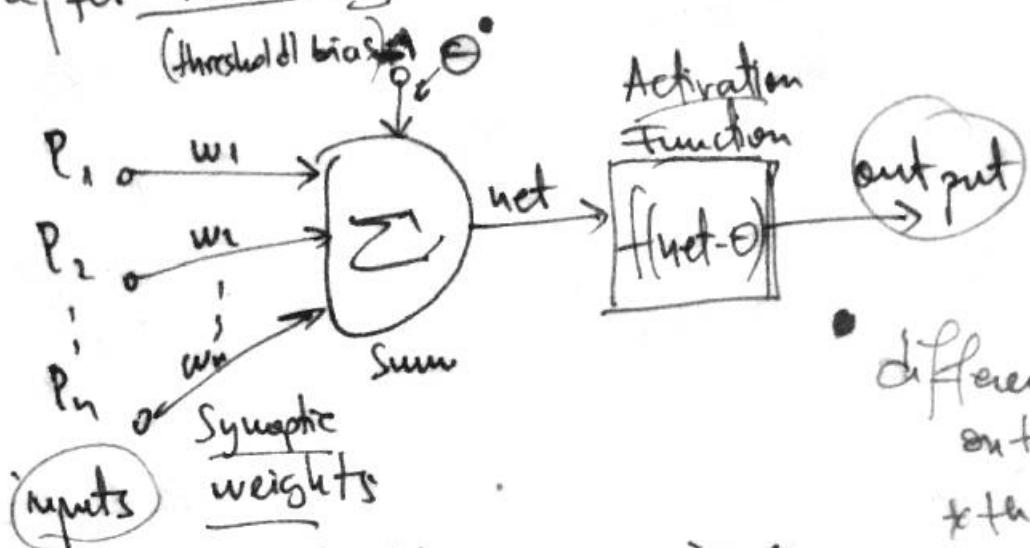
On short timescales, neurons act as devices which map inputs at synapses to a sequence of action potentials (spikes) at the output.

The Biological Neuron



The Artificial Neuron

- emulates the biological one & collects input signals (pre-synaptic spikes) and provides an output after reaching a threshold (post-synaptic spikes)



The net input to the neuron is :

$$\text{net} = \sum_{i=1}^n P_i \cdot w_i = P_1 \cdot w_1 + P_2 \cdot w_2 + \dots + P_n \cdot w_n$$

(θ, bias is an intrinsic property of the neuron when no input)

The output is computed as :

$$\text{out} = f(\text{net} - \theta)$$

↑ viewed as threshold to fire

- The neuron performs a nonlinear mapping from input to output typically from a high-dimensional input space to a low-dimensional output space ($1, n, m$)

- Mathematically is a weighted sum of the inputs and the application of an activation (squashing) function that determines the output of the neuron.

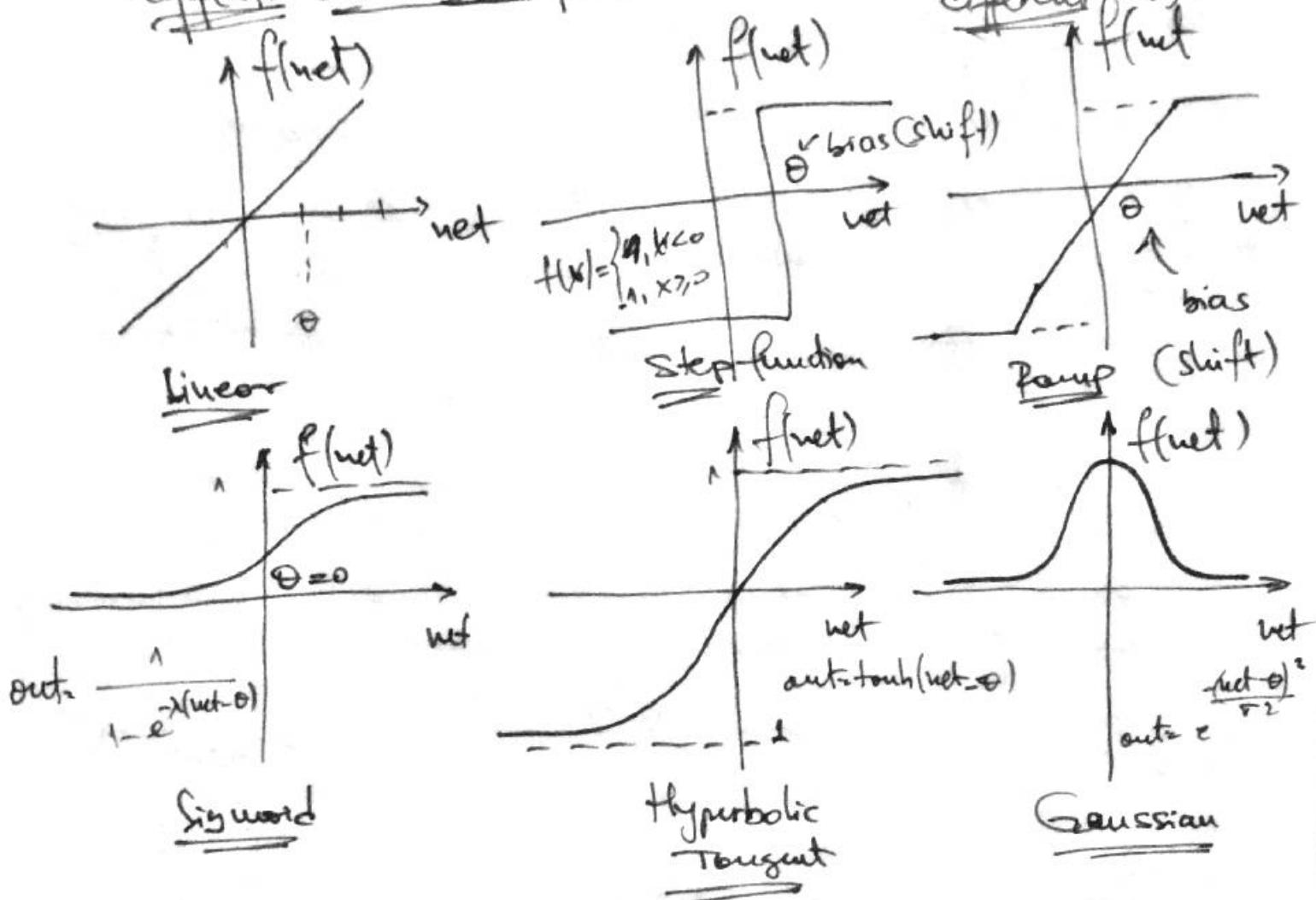
different views

on the bias

* threshold

* persistent activity of the neuron when no input

Typical nonlinearities (activation functions) are monotonically increasing functions with different effects on the output of the neuron. (They have to be differentiable).



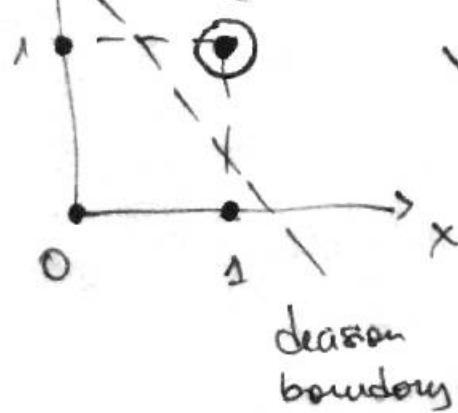
But, what can a single neuron compute?

→ Rosenblatt's definition
A 2 input perceptron with a single neuron with a step activation function is capable of separate / classify input patterns (simple e.g. $w_1 < 0, w_2 > 0$).

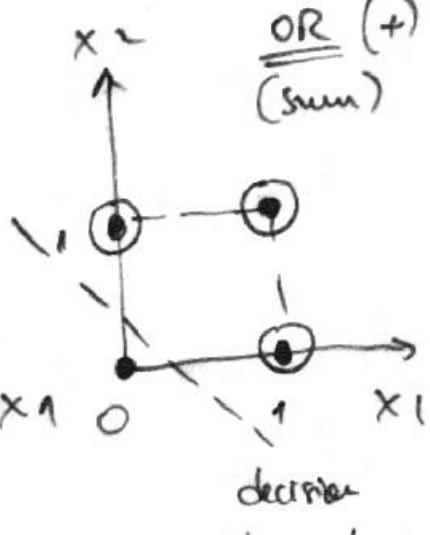
Such a classification implies a decision / separation boundary which is determined by the input vectors for which the net input is zero.

Considering linearly separable classes a single neuron perceptron can implement problems involving logic operations, such as AND and OR functions.

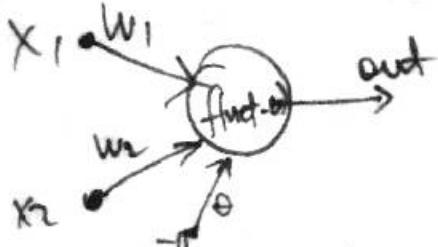
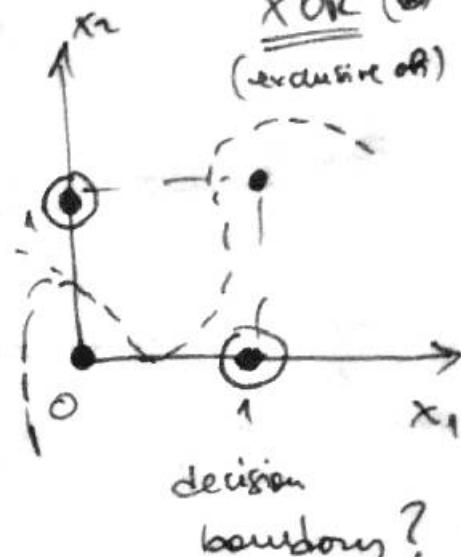
x_2 AND (\cdot)
(product)



OR (+)
(sum)



XOR (⊕)
(exclusive OR)



Sample solution:

$$w_1 = w_2 = 0.8$$

$$\theta = 1$$

Why?

$$0.8x_1 + 0.8x_2 - 1 < 0$$

for $x_1 = x_2 = 0$

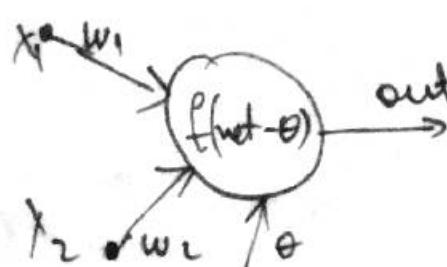
$$x_1 = 1 \quad x_2 = 0$$

$$x_1 = 0 \quad x_2 = 1$$

otherwise

$$0.8x_1 + 0.8x_2 - 1 > 0$$

if $x_1 \neq x_2 \neq 1$



Sample solution:

$$w_1 = w_2 = 0.8$$

$$\theta = 0.5$$

Why?

$$0.8x_1 + 0.8x_2 - 0.5 < 0$$

for $x_1 = x_2 = 0$

otherwise

$$0.8x_1 + 0.8x_2 - 0.5 > 0$$

for $x_1 = x_2 = 1$

$$x_1 = 1 \quad x_2 = 0$$

$$x_1 = 0 \quad x_2 = 1$$

We need a multi layer perceptron to cope with nonlinearity

Next lecture we will build a network for tasks

In order to build a linear decision boundary we need to train the system by updating the weights (for both inputs & bias).

Learning in artificial neurons

* training is based on Gradient Descent Learning Rule;

* minimize an error function of the mismatch between the target (expected) and the actual output of the neuron.

Error is given by:

$$E(t) = \sum_{p=1}^{P_T} (t_p (+) - \text{out}_p (+))^2$$

Supervision

t_p - target for pattern p used to train

out_p - output of the neuron for pattern p used to train

Squared error penalizes large errors more than small errors.

For each training pattern p the weight update follows the negative gradient in weight space to

minimize the mismatch.

The equation for gradient descent weight update

$$w_i(t) = w_i(t-1) + \Delta w_i(t)$$

$$\Delta w_i(t) = \eta \cdot \left(-\frac{\partial E_p(t)}{\partial w_i(t)} \right)$$

If classes are linearly separable the algorithm converges in a finite number of steps.

Learning rule in a single neuron

E = error signal for the entire training dataset

E_p = error for pattern p used in training

w_i = current set of weights

$$E = \sum_{\text{points in dataset}} (\text{true}_p - \text{out}_p)^2 \Rightarrow E_p = (\text{true}_p - \text{out}_p)^2$$

Given the error we compute the gradient (derivative) of E_p w.r.t input weights of the linear neuron (we consider a linear activation function):

$$\Delta w_i(t) = \eta G_i(t) \quad \text{with } G_i(t) = \frac{\partial E_p(t)}{\partial w_i(t)}$$

which we can rewrite using the "chain rule":

$$G_i(t) = \frac{\partial E_p(t)}{\partial w_i(t)} = \frac{\partial E_p(t)}{\partial \text{out}_p(t)} \cdot \frac{\partial \text{out}_p(t)}{\partial w_i(t)} \quad \text{with } n = \text{input size}$$

$$\bullet \text{out}_p(t) = \sum_{i=1}^n p_i(t) \cdot w_i(t)$$

And if we analyze both factors:

$$(1) \frac{\partial E_p(t)}{\partial \text{out}_p(t)} = \frac{\partial (\text{true}_p(t) - \text{out}_p(t))^2}{\partial \text{out}_p(t)} = \frac{2(\text{true}_p(t) - \text{out}_p(t))}{1}$$
$$= -2(\text{true}_p(t) - \text{out}_p(t))$$

and

$$(2) \frac{\partial \text{out}_p(t)}{\partial w_i(t)} = \frac{\partial \sum_{j=1}^n p_j(t) \cdot w_j(t)}{\partial w_i(t)} = p_i$$

only for $j=i$ the derivative exists

Combining (1) and (2) the gradient in the direction of the weight w_i is given by:

$$G_i(t) = \frac{\partial E_p(t)}{\partial w_i(t)} = \frac{\partial E_p(t)}{\partial \text{out}_p(t)} \cdot \frac{\partial \text{out}_p(t)}{\partial w_i(t)} = -2(\text{true}_p(t) - \text{out}_p(t)) \cdot p_i$$

$$\zeta_i(t) = -2(\text{true}_p(t) - \text{out}_p(t)) \cdot p_i$$

$$\Delta w_i(t) = (\eta) \cdot (-2) \cdot (\text{true}_p(t) - \text{out}_p(t)) \cdot p_i$$

used as weight adaption to minimize the error & learning rate tells us how fast One can use also non-linear activation functions as long as they are differentiable!

From single neurons to neural networks

(Single neuron perceptrons to multi-layer perceptrons)

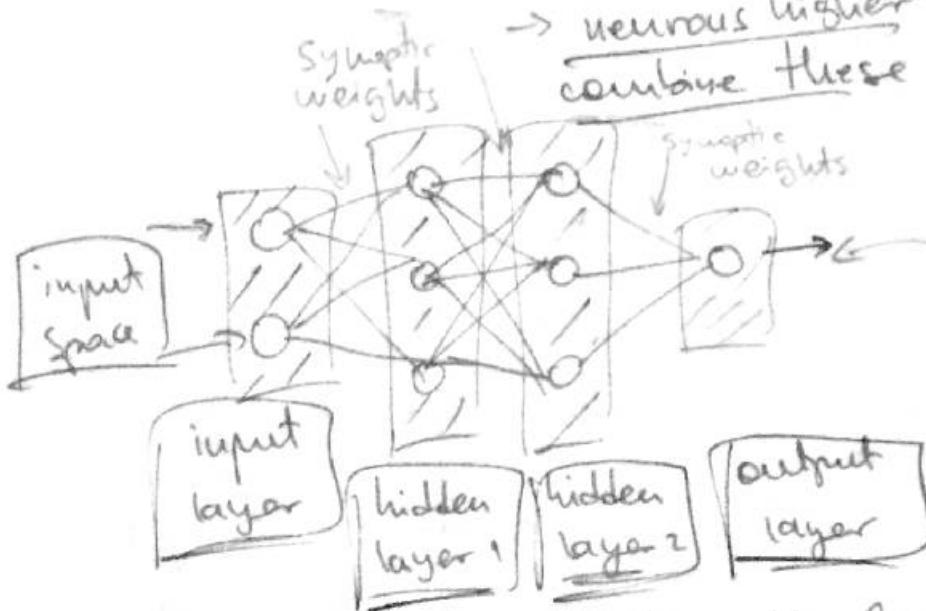
AIML
class

What if the task is to learn "arbitrary data" which doesn't fit in a single neuron transfer function?

→ combine neurons in networks

→ each neuron represent an aspect of the data

→ neurons higher up in the hierarchy combine these



Typical structure for a feed-forward network.

- The structure dictates the learning algorithm capabilities.
- We can have different processing models and learning capabilities depending on the network architecture: feed-forward, recurrent, deep etc.

Q

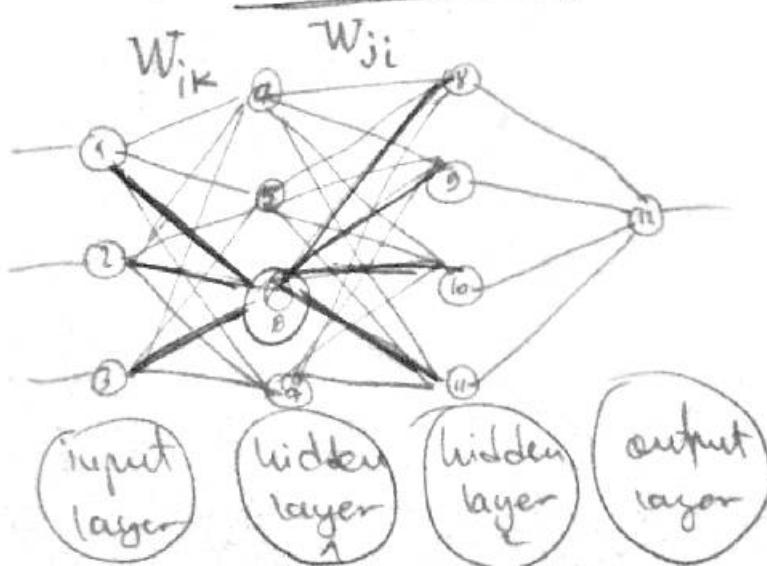
What if we try to apply the gradient descent rule we used for a single neuron? Will it work for a network?

As here we have multiple layers of neurons, each one has a different contribution. We know indeed the output error of the network but no individual neuron error.

The core idea is to backpropagate the output (network) error through the layers in order to update the weights.

Backpropagation / Error backpropagation

- searches the minimum of the error function in weight space using gradient descent.
- the particular combination of weights which minimizes the error function is considered to be the solution for learning a representation of the data.
- requires computing the gradient at each training iteration so we must guarantee that the error function is continuous and differentiable.
- computes the partial derivatives of the cost/error function w.r.t. to any weights in the network



$$k = \{1, 2, 3\}$$

$$i = 6$$

$$j = \{8, 9, 10, 11\}$$

Definitions:

- ① Error for unit i at time t

$$\delta_i(t) = -\frac{\partial E_i(t)}{\partial \text{net}_i(t)}$$

where the net input for neuron i is

$$\text{net}_i(t) = \sum_{k \in A} w_{ik}(t) \cdot \text{out}_k(t)$$

output
at
previous
layer

- ② The weight change for $w_{ik}(t)$ is

$$\Delta w_{ik}(t) = \frac{-\partial E_i(t)}{\partial w_{ik}(t)}$$

- We need to rewrite the weight change in terms of the error at the current neuron ($\delta_i(t)$) and the output of the previous layer ($\text{net}_i(t)$) this is the only information available at neuron i!

$$\Delta w_{ik}(t) = \frac{-\partial E_i(t)}{\partial w_{ik}(t)} = \frac{-\partial E_i(t)}{\partial \text{net}_i(t)} \cdot \frac{\partial \text{net}_i(t)}{\partial w_{ik}(t)}$$

$$\Delta w_{ik}(t) = \delta_i(t) \cdot \text{out}_k(t)$$

$$\text{where } \delta_i(t) = -\frac{\partial E_i(t)}{\partial \text{net}_i(t)}$$

$$\frac{\partial \text{net}_i(t)}{\partial w_{ik}(t)} = \frac{\partial \sum_{l \in A} w_{il}(t) \cdot \text{out}_l(t)}{\partial w_{ik}(t)} = \frac{\partial \text{out}_k(t)}{\partial w_{ik}(t)}$$

which is
non-zero
for $l = k$

- ③ Computing the forward activation of the network

- when input is applied to the network, we compute the output of all neurons through layers to the output

$$\text{out}_i(t) = f_i(\text{net}_i(t)) = f_i \left(\sum_{k \in A} w_{ik}(t) \cdot \text{out}_k(t) \right)$$

squashing function

④ Calculating the error at the output neuron

- Assuming that the expected / target value for the dataset point is target_o, then the error signal at the output, out_o, is

$$\delta_o(t) = -\frac{\partial E_o(t)}{\partial \text{net}_o(t)} = 2 \cdot (\text{target}_o - \text{out}_o)$$

(here we assumed a linear activation function)

⑤ Propagating the error back through the network

- after computing the error at the output of the net we propagate the signal back through the network
- the error at unit/neuron i is

$$\delta_i(t) = \frac{-\partial E_i(t)}{\partial \text{net}_i(t)}$$

and applying the "chain rule"

$$\delta_i(t) = \frac{-\partial E_i(t)}{\partial \text{net}_i(t)} = \sum_{j \in P} \frac{-\partial E_i(t)}{\partial \text{net}_j(t)} \cdot \frac{\partial \text{net}_i(t)}{\partial \text{out}_i(t)} \cdot \frac{\partial \text{out}_i(t)}{\partial \text{net}_i(t)}$$

$$1.) \frac{-\partial E_i(t)}{\partial \text{net}_j(t)} = \delta_j(t)$$

$$2.) \frac{\partial \text{net}_j(t)}{\partial \text{out}_i(t)} \cdot \frac{\partial \sum_{m \in P} w_{jm}(t) \cdot \text{out}_m(t)}{\partial \text{out}_i(t)} = w_{ji}$$

nonzero for m = i

$$3) \frac{\partial \text{out}_i(+)}{\partial \text{net}_i(+)} = \frac{\partial f(\text{net}_i(+))}{\partial \text{net}_i(+)} = f'(\text{net}_i(+))$$

- We combine 1), 2), 3) to calculate the error signal of neuron i in the network:

$$\delta_i(+) = \sum_{j \in P} [\delta_j(+) \cdot w_{ji}(+) \cdot f'(\text{net}_i(+))]$$

f is independent of j^{th} neuron so we can rewrite as:

$$\delta_i(+) = f'(\text{net}_i(+)) \cdot \sum_{j \in P} [\delta_j(+) \cdot w_{ji}(+)]$$

the error for neuron i : $(\delta_i(+))$ only depends on the known error in higher layers (posterior)

- Computing the weight update for weight anterior to posterior ($k \leftarrow i$)

$$\Delta w_{ik}(+) = \delta_i(+) \cdot \text{out}_k(+) \quad \text{and} \quad \delta_i(+) = f'(\text{net}_i(+)) \sum_{j \in P} \delta_j(+) \cdot w_{ji}(+)$$

The final synaptic weight update rule is:

$$\Delta w_{ik}(+) = f'(\text{net}_i(+)) \cdot \sum_{j \in P} \delta_j(+) \cdot w_{ji}(+) \cdot f(\text{net}_k(+))$$

or

$$\Delta w_{ik}(+) = f'(\text{net}_i(+)) \cdot \sum_{j \in P} \delta_j(+) \cdot w_{ji}(+) \cdot \text{out}_k(+)$$

We use this rule to train multi-layer neural nets for tasks as classification & regression.

Supervised neural learning \rightarrow tips & tricks

- the big problem with neural networks is parameter tuning \Rightarrow there is no optimal way to select network parameters.
- there are though some design choices for
 - (**) data pre-processing
 - (**) weight initialization
 - (**) error functions

① Weights initialization

- zero initialization \Rightarrow each neuron computes same output they will also learn based on similar gradients and there will be no symmetry break, and network will act as a single, highly redundant neuron
- to break symmetry initialize all neurons with small random numbers. (± 0.001)

② I/O normalization

- normalize each data dimension to be on approximately the same scale
- 2 ways:
 - * divide each dimension by its std. dev., once is zero centred
 - * normalize to $(-1, 1)$ on each dimension

③ Small weight updates

- modulated by the learning rate $\eta(t)$ typically $\eta > 0.001$

4. Avoiding local minima in weight space

- due to the gradient based learning rule the training can get stuck in local minima.
- a solution is to use Simulated Annealing to slightly perturb the weights
 - large perturbation at the beginning to get the network out of minimum and then slow/decay with the training process until settles in a global solution (usually marked by small Δw)

5. Network size

- the no. of inputs & outputs given by the problem, cannot be changed.
- no. of neurons, no. of layers, type of squashing functions?

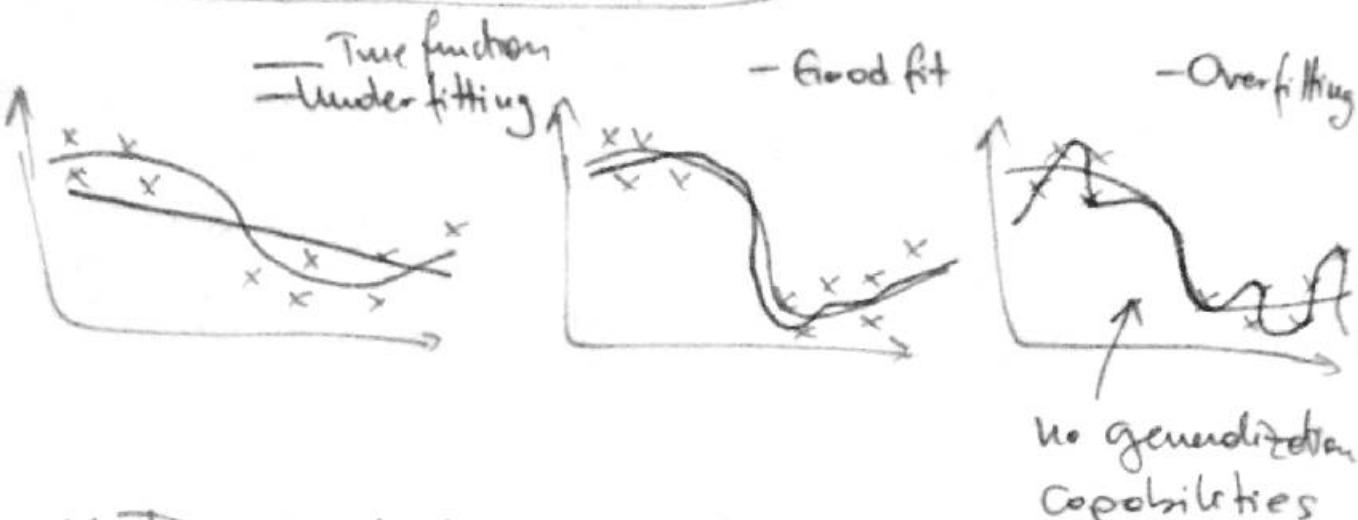
↓
a larger neuron count (capacity of the net)
will allow the net to represent more
complex functions



but there is no universal recipe!

- Some "rules of thumb":
 - size of the hidden layer should be between the input layer size and the output size
 - in a single hidden layer MLP you need $2 \times$ more neurons than inputs

⑥ Overtraining / overfitting

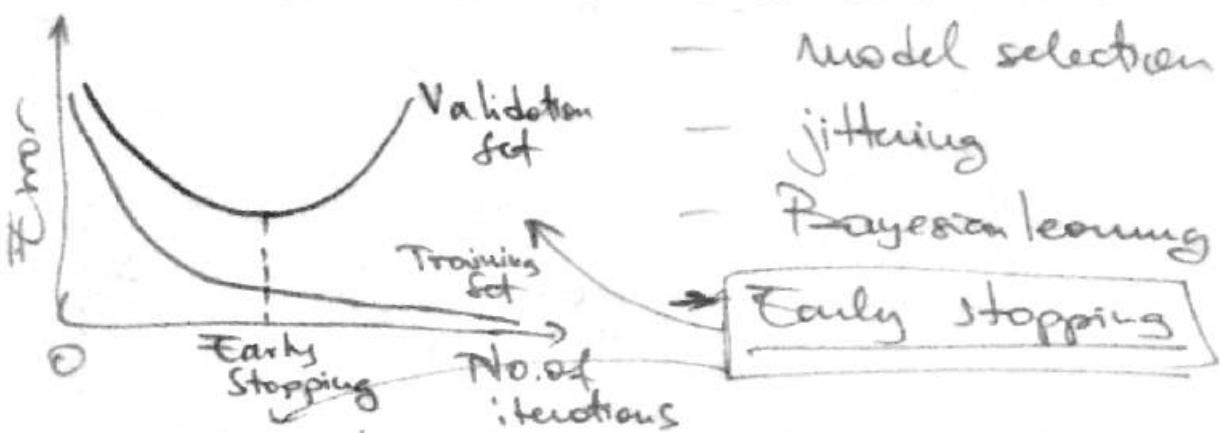


* To avoid this, use large amounts of data as a rule!

* other methods are also used:

- Model selection
- jittering
- Bayesian learning

Early stopping



- while training on the dataset the network becomes better and better but at some point in time it starts to represent each point! \rightarrow this is problematic \rightarrow overfitting!
- We need to find the time where the network finishes generalizing and starts learning individual points.
- Split dataset in **70% training + 30% testing**

(Training set)

(Validation Set)

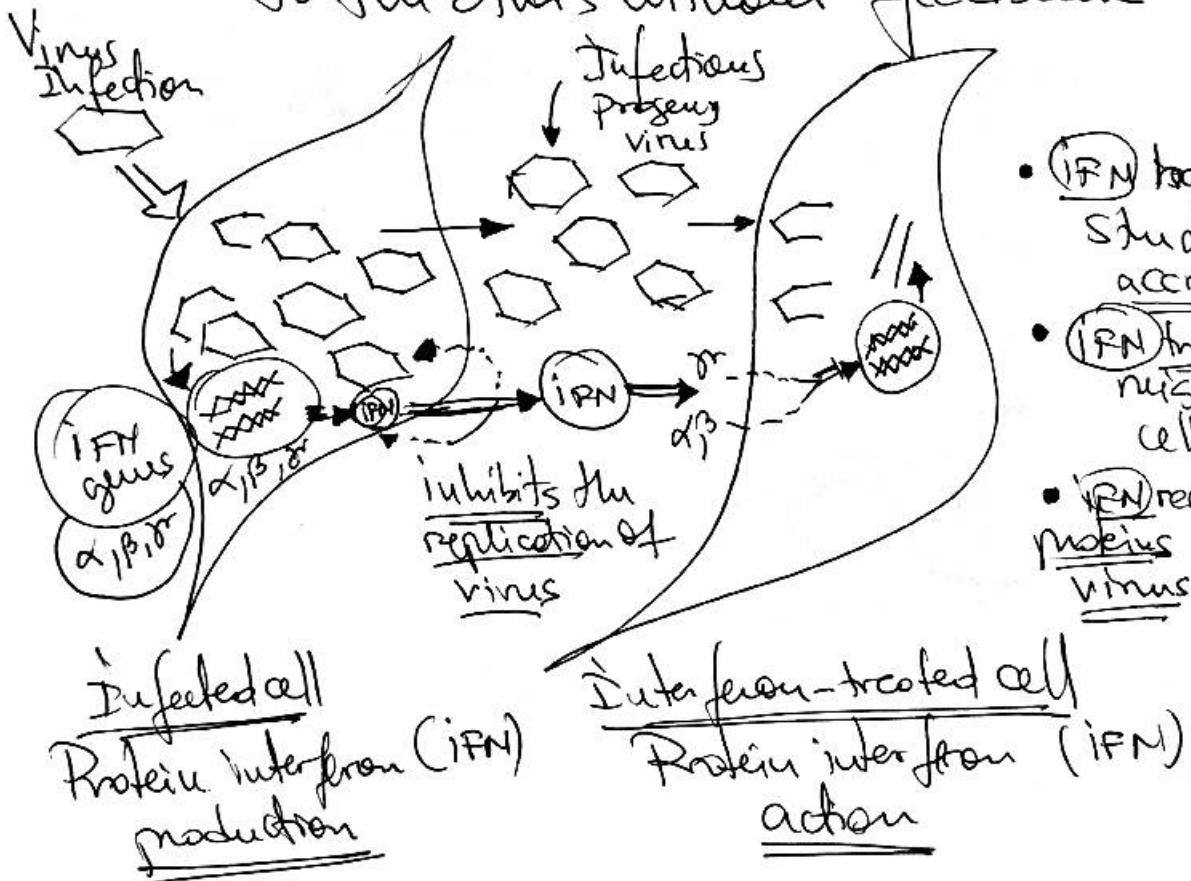
Signaling & Message Passing Mechanisms

in the Immune System

① Signal diffusion || Communication
② Dialogue Schemes

① Immune diffusion

- information passes from one immune component to the others without feedback



- IFN has some structure across cells
- IFN travels to neighboring cells
- IPN regulates protein inhibitor virus replication

② Immune Dialogue

- immune components exchange continuously molecular signals.
- sensitivity is context dependent.