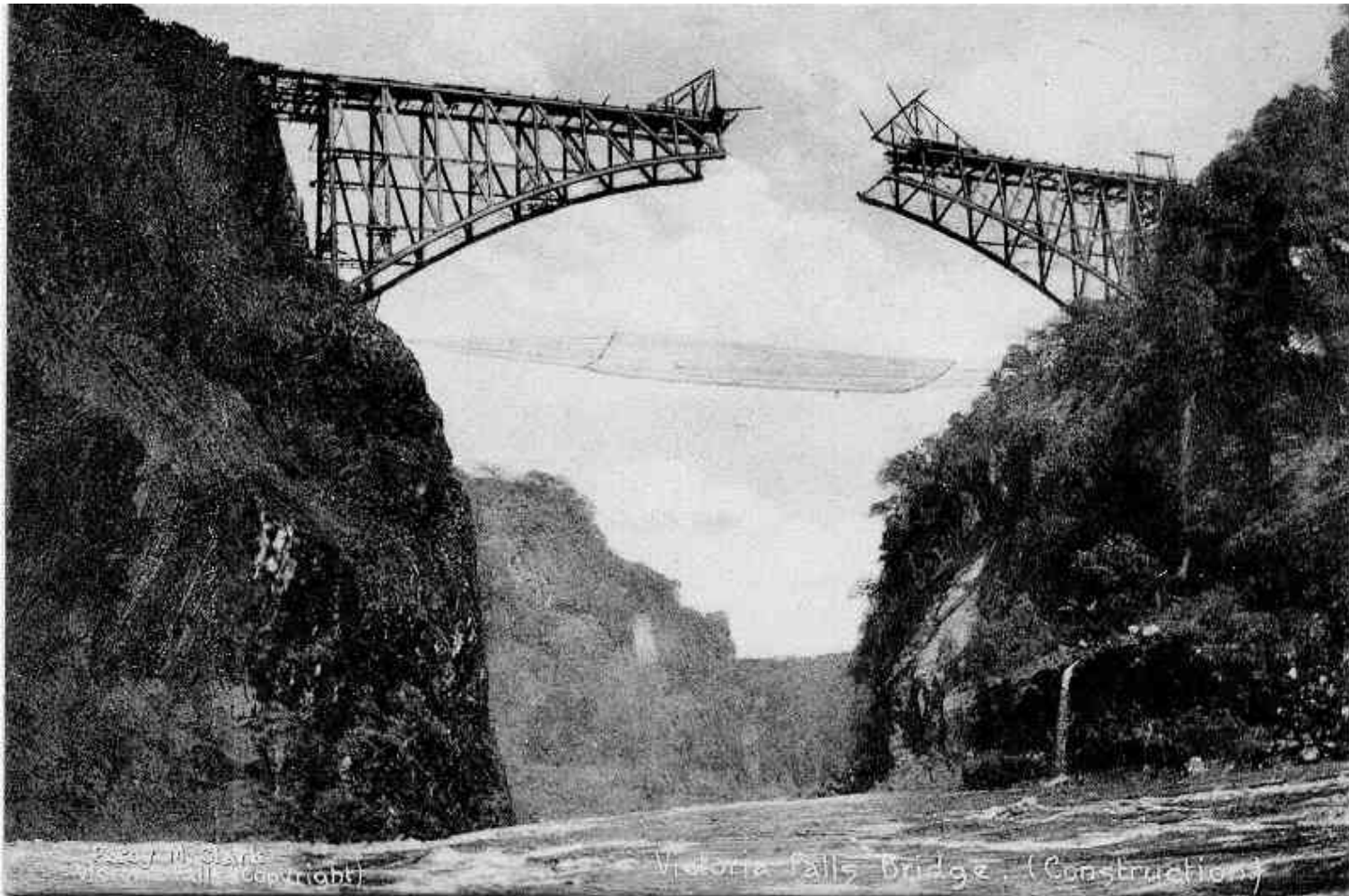


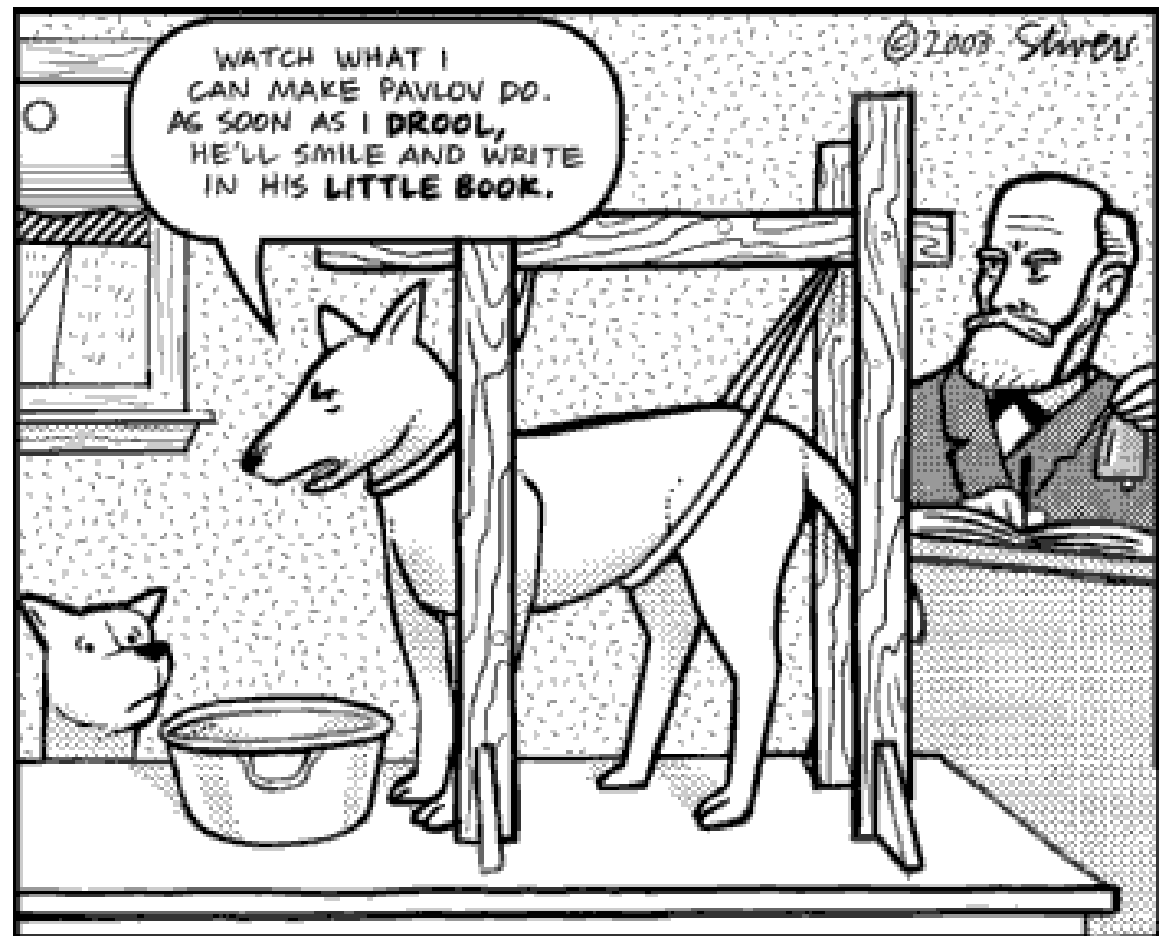
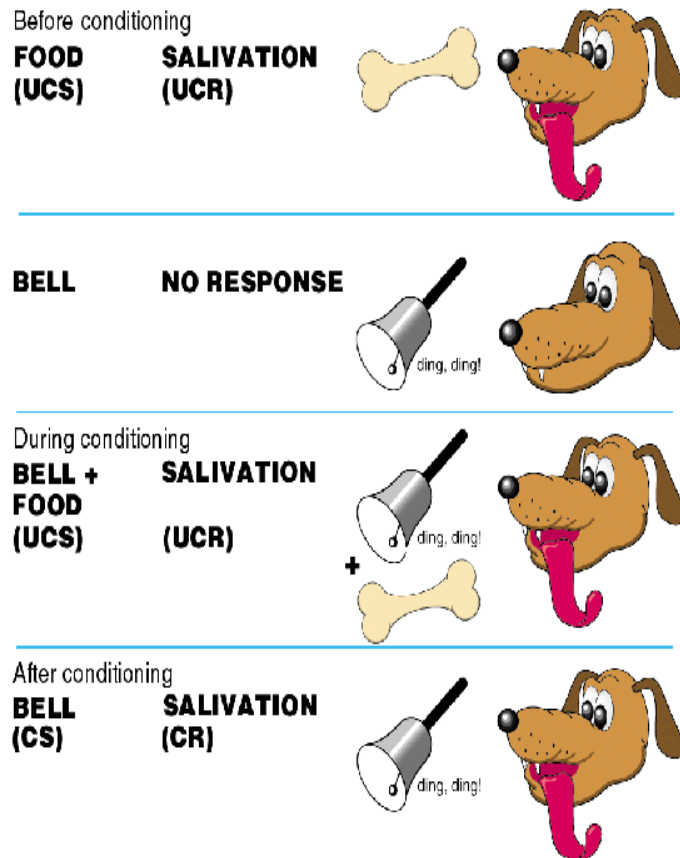
Computational Neuroscience



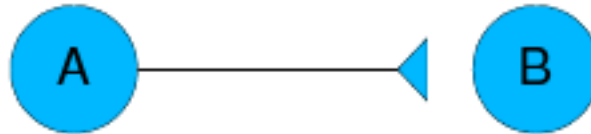
Structure – Dynamics – Implementation – Algorithm – Computation - Function

Learning at psychological level

- Classical conditioning



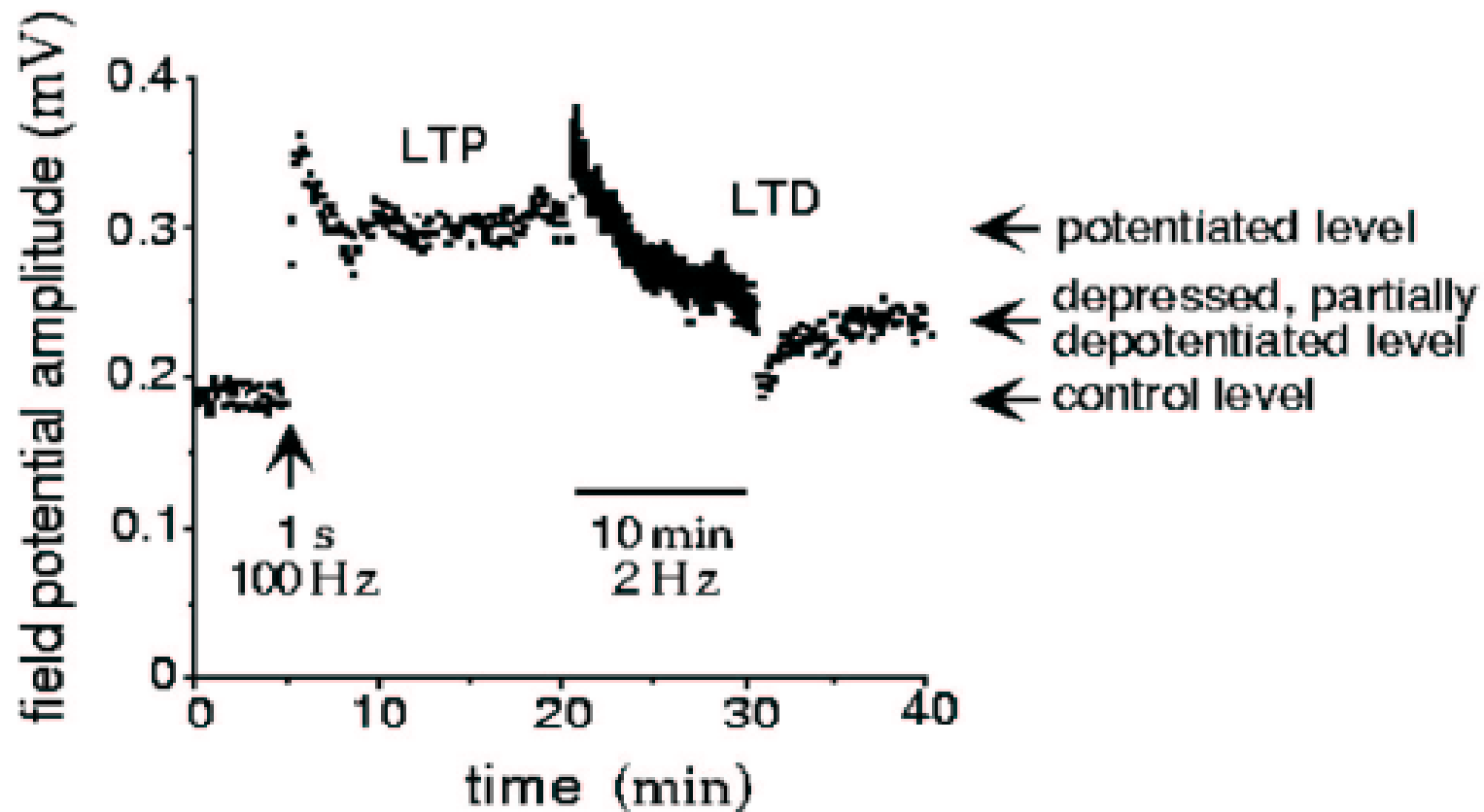
Hebb's rule



When an axon of cell A is near enough to excite cell B, and repeatedly or consistently takes part in firing it, some growth process or metabolic change takes part in one or both cells such that A's efficiency, as one of cell firing B, is increased"
(Hebb, The Organization of Behavior, 1949)

$$\tau_w \frac{dw_i}{dt} = v \cdot u_i$$

Hebb's rule (?) in an experiment at population level

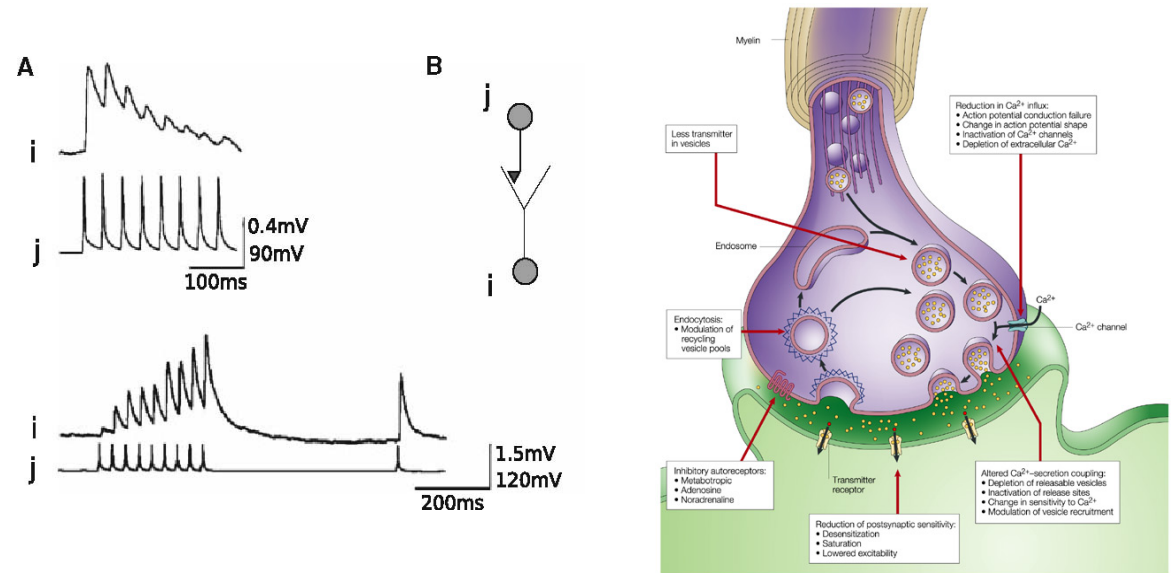


LTP – long term potentiation

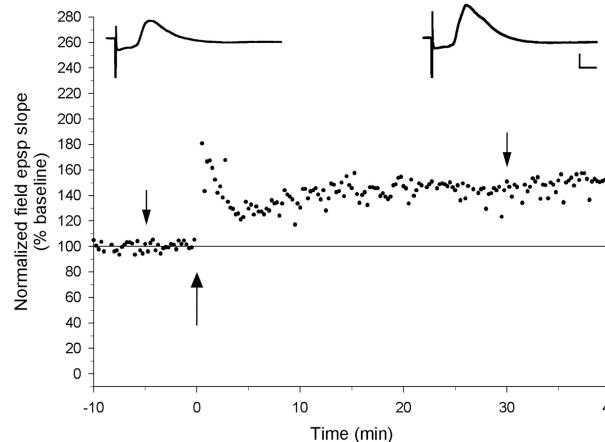
LTD – long term depression

Plasticity in the neural system

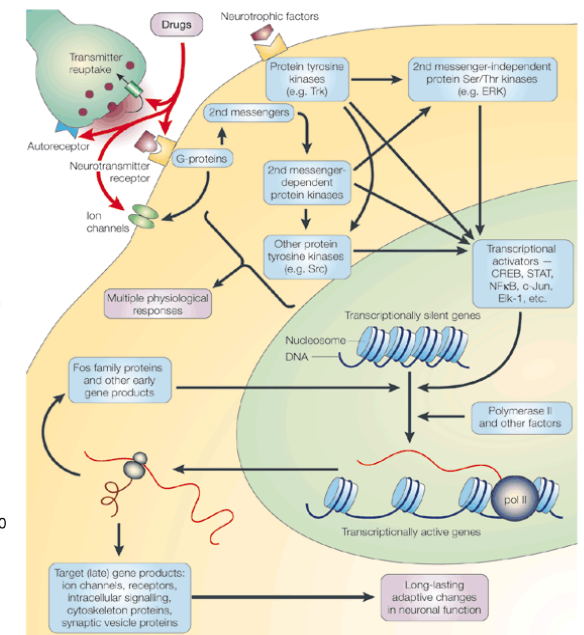
- Position of the plasticity: synapses, firing thresholds of the postsynaptic neurons (excitability)
- Potentiation, depression
- STP: Calcium dynamics, transmitter depletion
range < 1 minute



- LTP: geneexpression (induction, expression, maintenance), NMDA magnesium-block
range > 1 minute

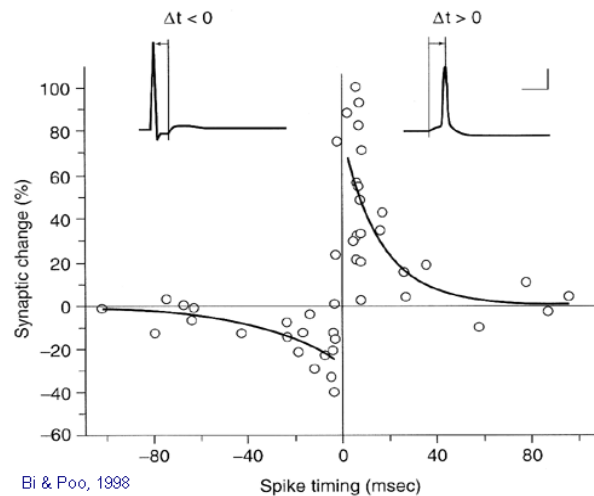


- Correlation between the molecular and the psychological level



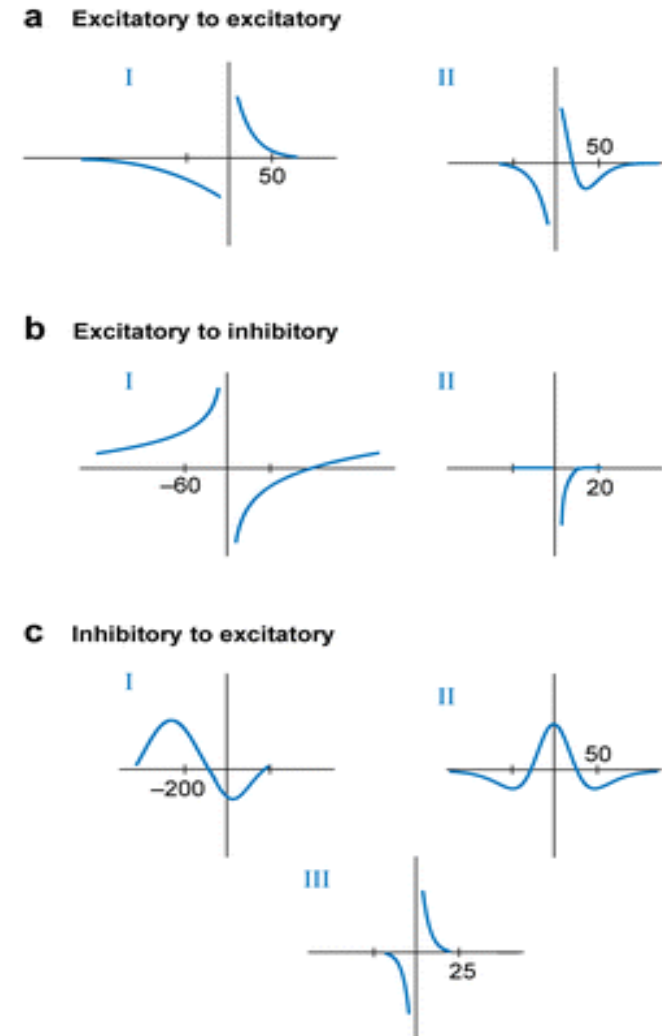
Spike-time dependent plasticity (STDP)

- Timing-dependent plasticity:
 - If the postsynaptic spike follows the presynaptic in short time window (causal order) the connection strength increases
 - Spiking in the anticausal order decreases the connection strength
 - Many more variables



- A Hebb-szabály formalizációja:
lineáris ráta-modellben

$$\tau_w \frac{d\mathbf{w}}{dt} = \nu \mathbf{u}$$



Caporale N, Dan Y. 2008.
Annu. Rev. Neurosci. 31:25–46.

Stabilized Hebb rules

- Problems with the Hebb rule:
 - Weights can only increase
 - No competition between the synapses – inputselectivity can not be implemented
- Simple solution: ceiling to the weights
- BCM: stabilizing with postsynaptic excitability

$$\tau_w \frac{d \mathbf{w}}{dt} = v \mathbf{u} (v - \theta_u) \qquad \tau_\theta \frac{d \theta_u}{dt} = v^2 - \theta_u$$

- Synaptic normalization

- Subtractive normalization

$$\tau_w \frac{d \mathbf{w}}{dt} = v \mathbf{u} - \frac{v (\mathbf{1} \cdot \mathbf{u}) \mathbf{1}}{N_u}$$

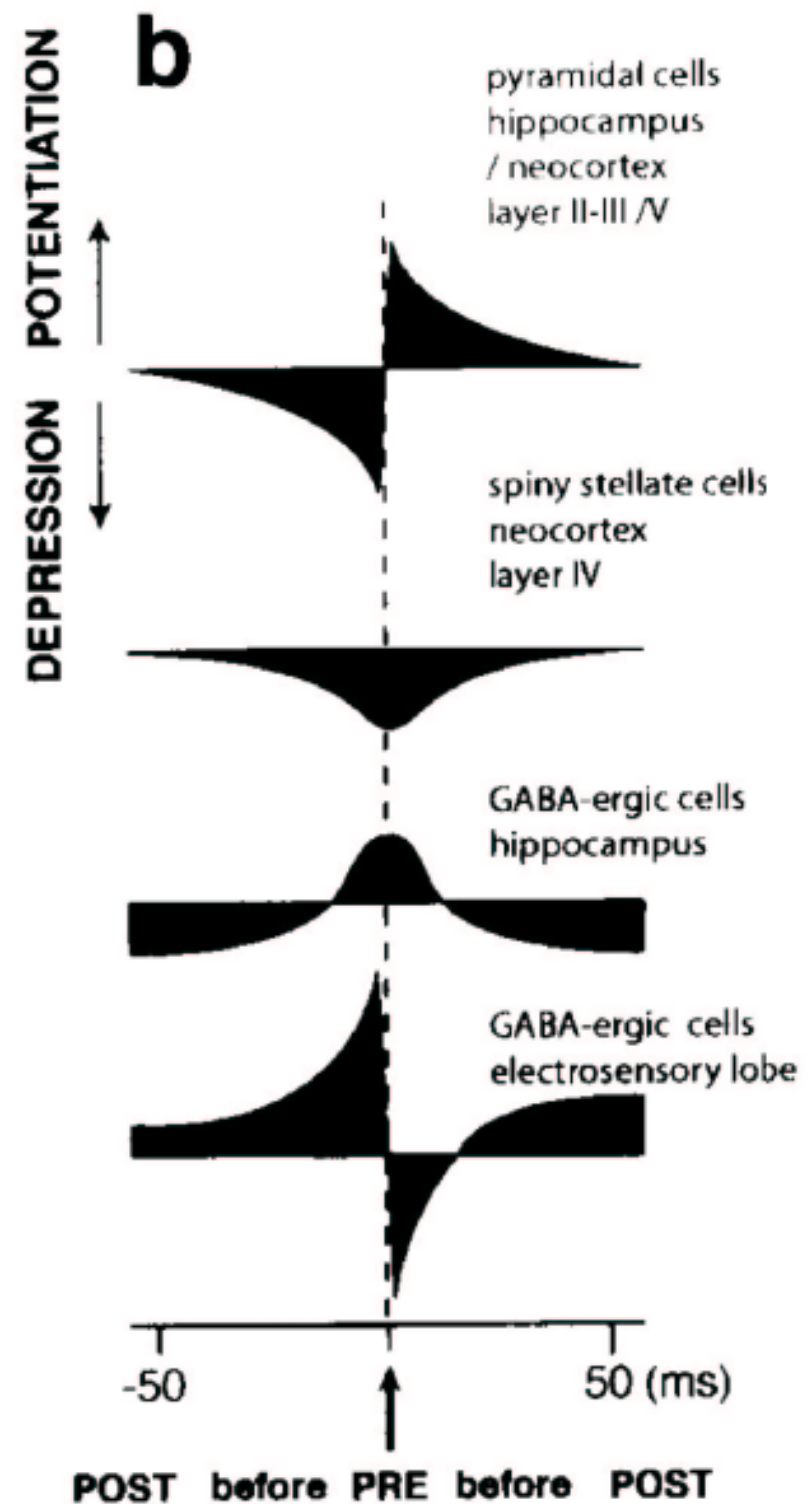
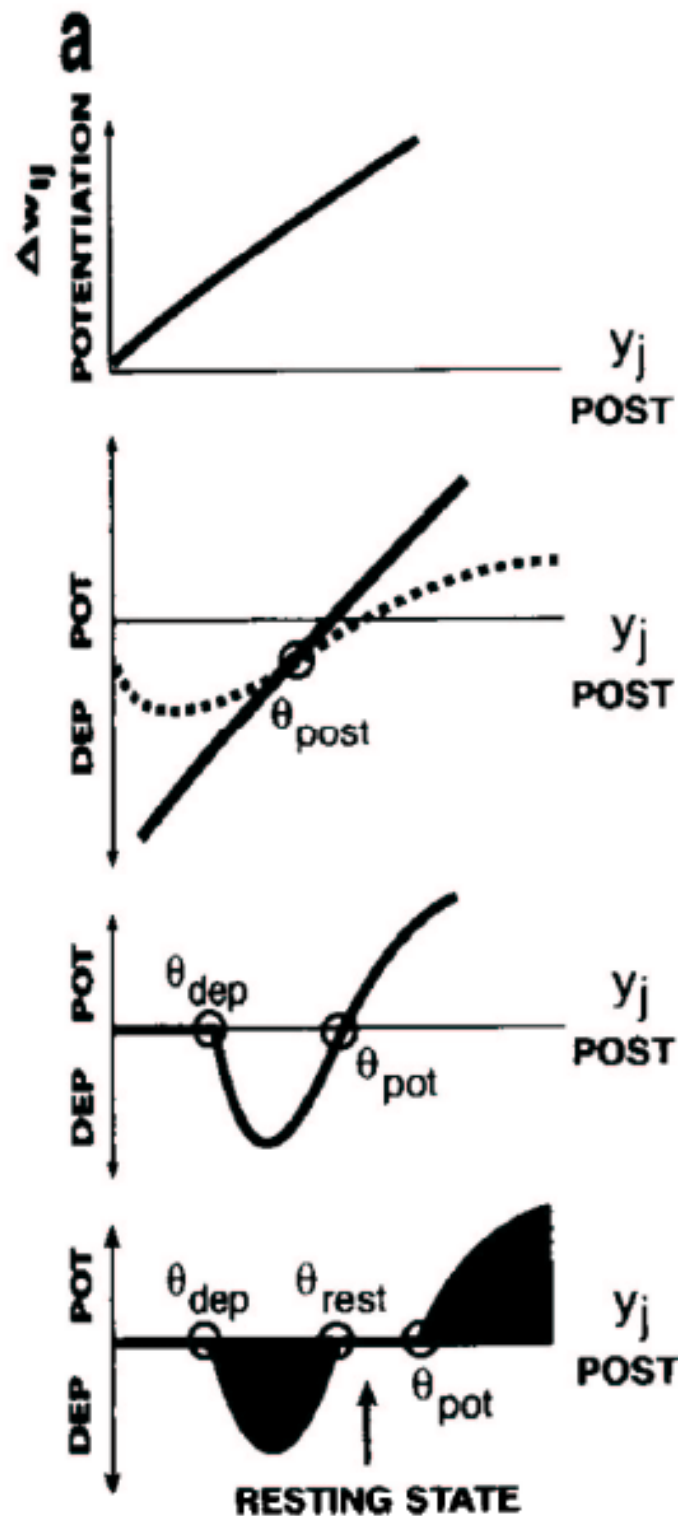
Global rule, but results in observed connection patterns (Ocular dominance)

- Oja-rule

$$\tau_w \frac{d \mathbf{w}}{dt} = v \mathbf{u} - \alpha v^2 \mathbf{u}$$

Local rule, but can generate the observed patterns

Rate-based
and
spike-time
dependent
learning rules

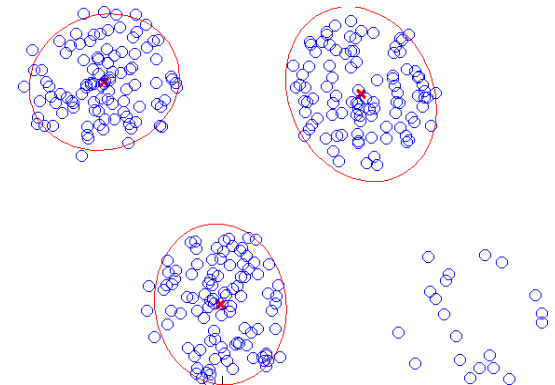
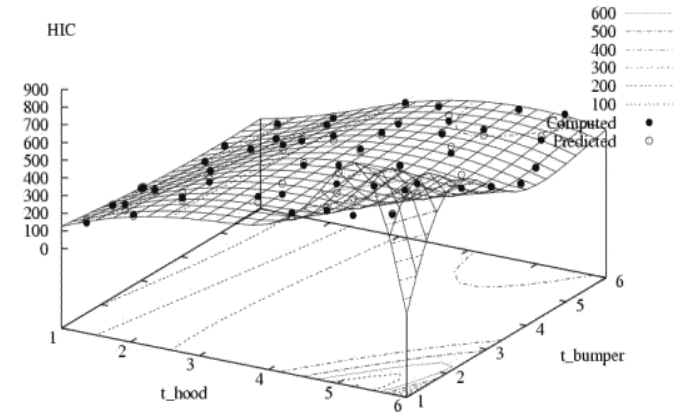


The mathematical formalization

- Tuning of the model parameters based on data
- Two level dynamics
 - Variables (input-output transformation) - fast
 - Parameters - slow
- Memory vs learning
 - Memory is a simple recall, without change of the representation
 - Learning, continuous refinement of the representation accompanied by output generation
- Main task: prediction of the future, based on the past

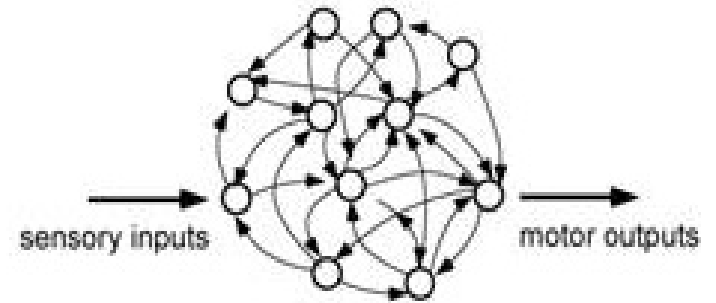
Three basic types of learning

- Supervised
 - data: input output pairs
 - Aim, function approximation, classification
- Reinforcement
 - data: observed states, rewards
 - Aim: optimal strategy for maximization of reward
- Unsupervised, representational
 - Data: set of input
 - Aim: optimal representation / model finding
- Combination of them

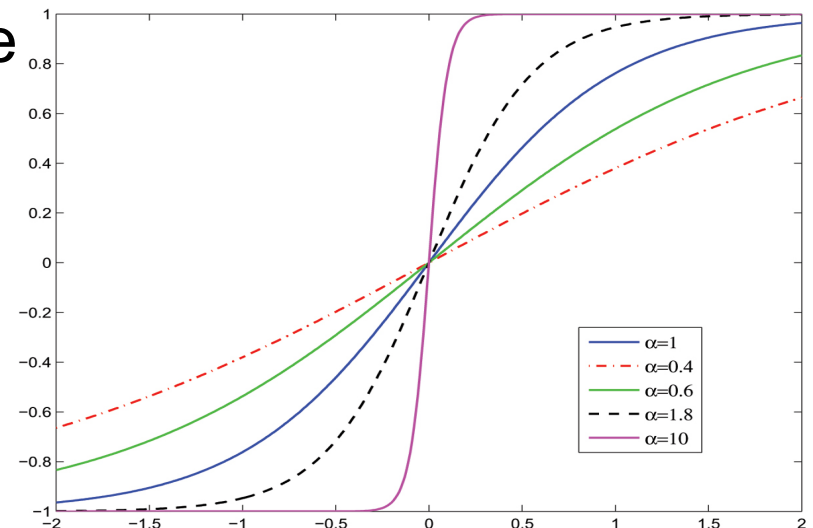


Learning in neural systems

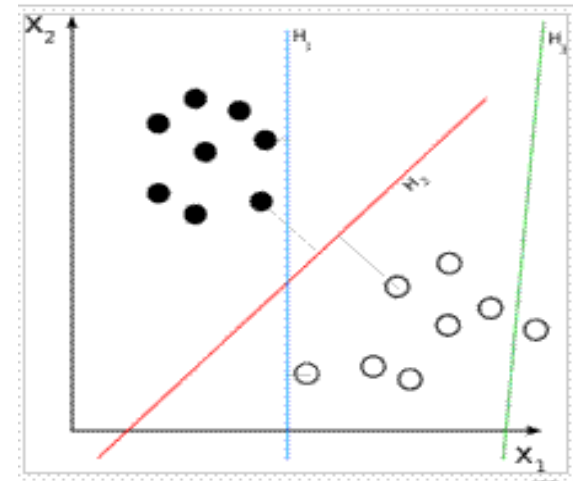
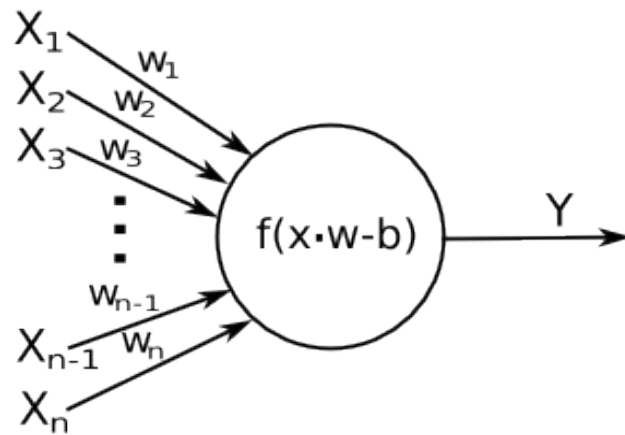
- Single neuron
- Feedforward network
- Recurrent network
- today: rate model
 - Parameters: weights, thresholds
 - Transfer functions
 - Step function: H (Heavyside)
 - Sigmoid
 - Linear neuron



$$y = f(\mathbf{x}\mathbf{w} - \theta)$$



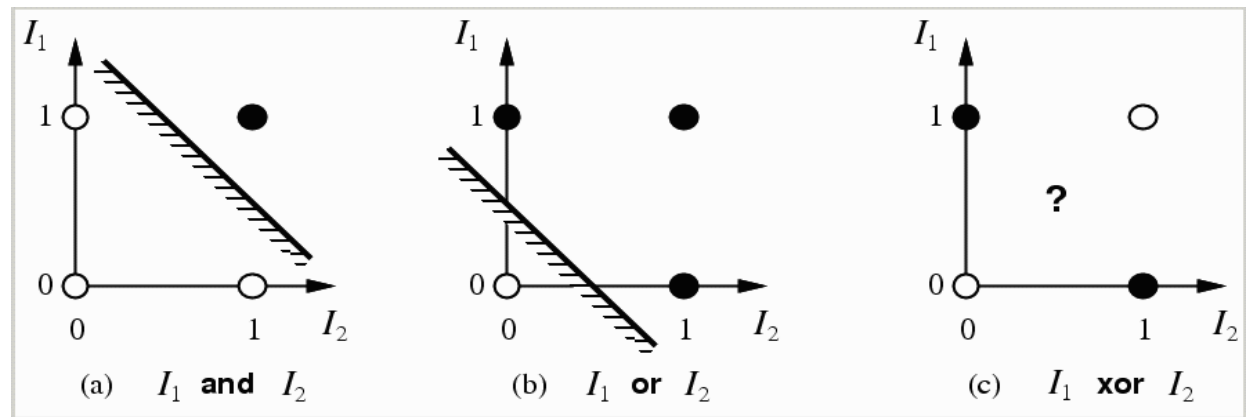
Perceptron



- Binary neurons: linear separation
 - In two dimension, the separator line:

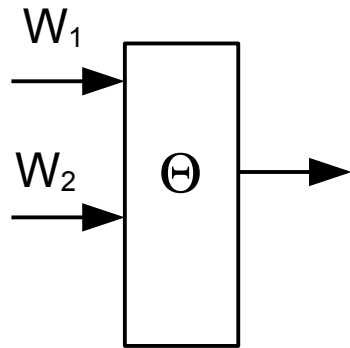
$$\theta = x_1 w_1 + x_2 w_2 \longrightarrow x_2 = \frac{-w_1}{w_2} x_1 + \frac{\theta}{w_2}$$

- Logical functions

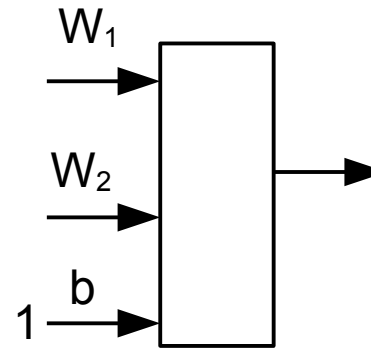


Error-correcting learning rules

Rosenblatt-algorithm – a binary neuron



$$y = H(w_1 + w_2 - \Theta)$$



$$y = H(w_1 + w_2 + b)$$

The threshold is transformed out by introducing a constant 1 input $b = -\Theta$ (**bias**)
Learning the **bias** is equivalent to the learning of a weight.

Error-correcting learning rules

- Using the actual and the correct answer, and the distance between them
- Rosenblatt-algorithm – binary neurons

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \epsilon (t^m - y(\mathbf{x}^m)) \mathbf{x}^m$$

- Delta-rule

- Continuous activity – gradient-method

$$w_b(t+1) = w_b(t) - \epsilon \frac{\partial E}{\partial w_b} \quad E = \frac{1}{2} \sum_m^{N_s} (t^m - y(\mathbf{x}^m))^2 \quad \frac{\partial E}{\partial w_b} = - \sum_m^{N_s} (t^m - y(\mathbf{x}^m)) \mathbf{x}^m$$

approximation for linear neurons:

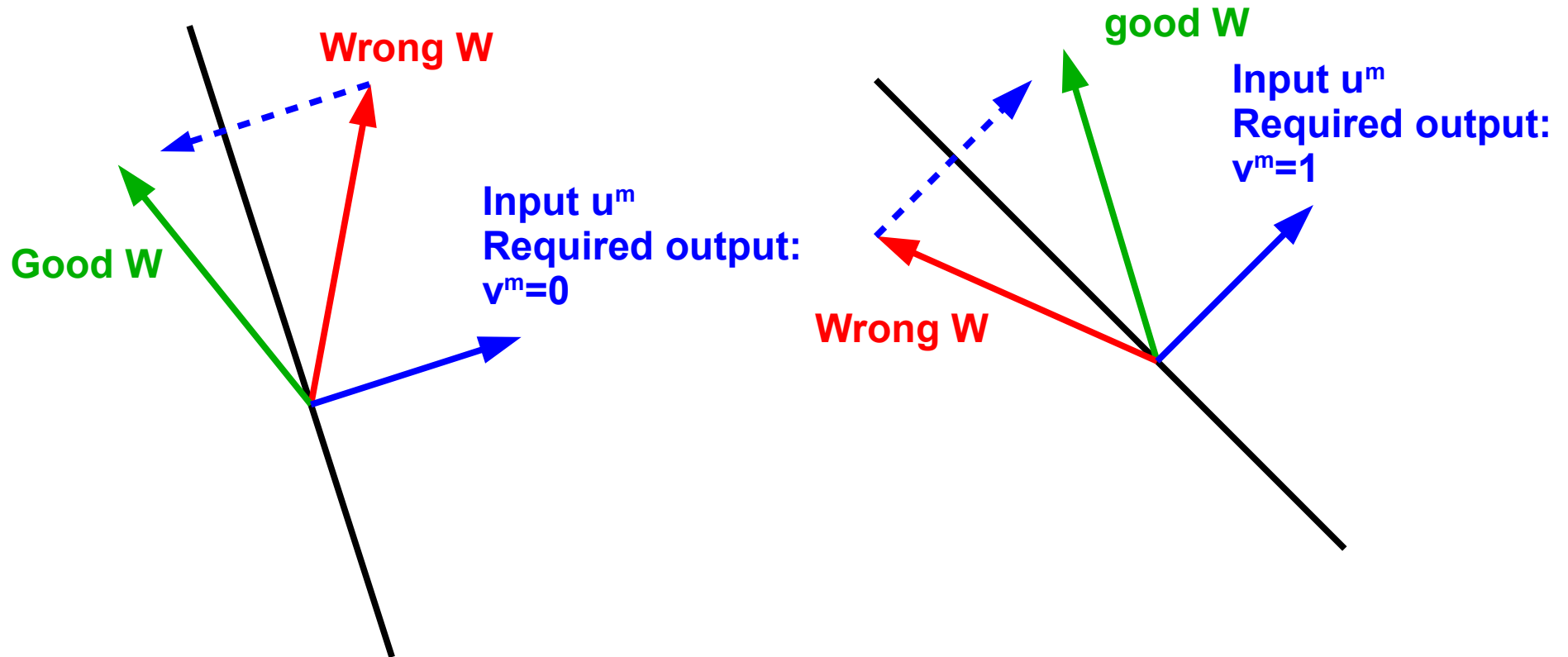
$$\mathbf{w}(t+1) = \mathbf{w} + \epsilon (t^m - y(\mathbf{x}^m)) \mathbf{x}^m$$

- Minsky-paper 1969: the neural networks can solve only linear problems

Error-correcting learning rules

Rosenblatt-algorithm – binary neurons

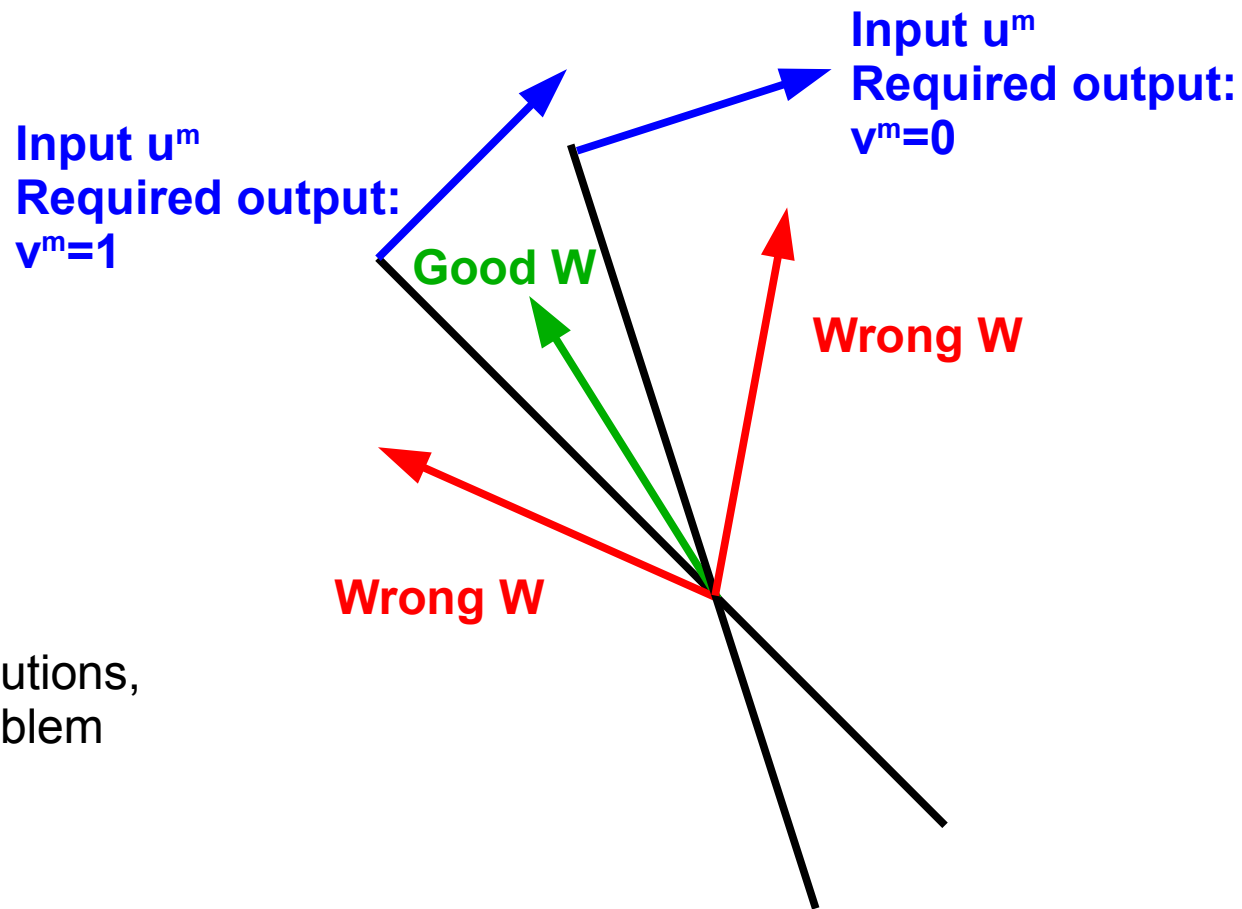
$$\mathbf{w} \rightarrow \mathbf{w} + \epsilon (\nu^m - \nu(\mathbf{u}^m)) \mathbf{u}^m$$



Error-correcting learning rules

Rosenblatt-algorithm – a binary neuron

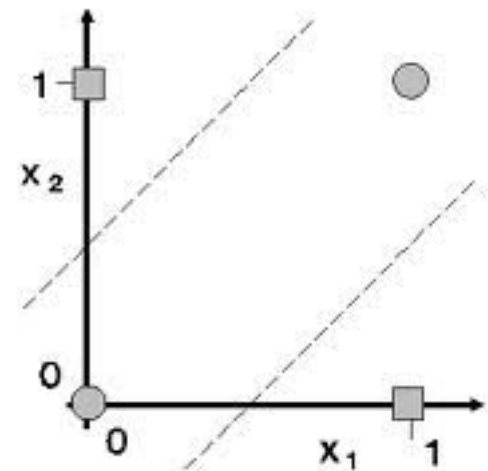
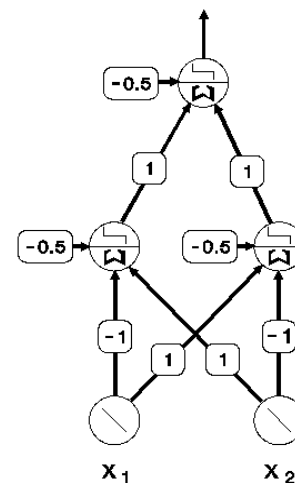
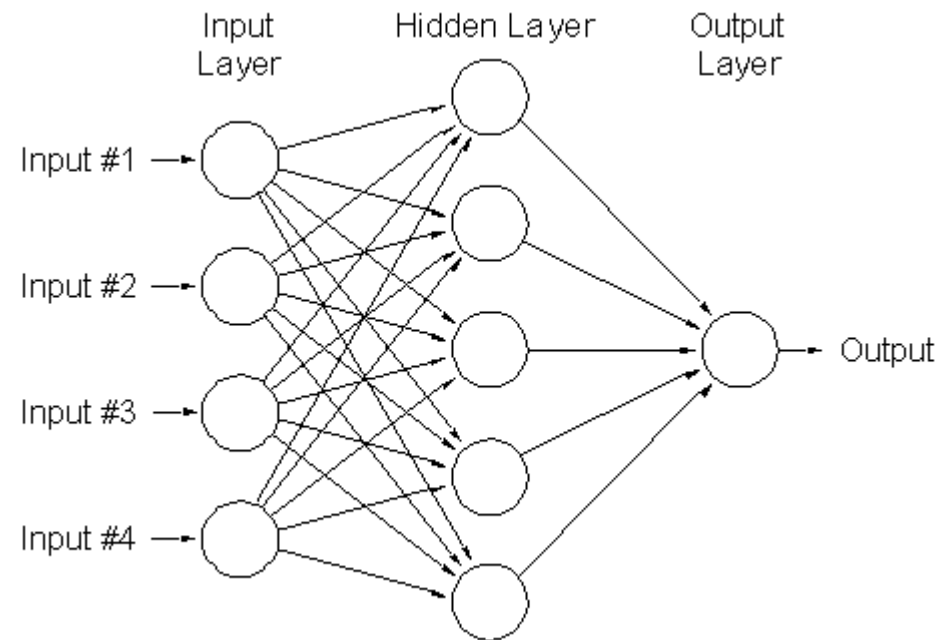
$$\mathbf{w} \rightarrow \mathbf{w} + \epsilon (\nu^m - \nu(\mathbf{u}^m)) \mathbf{u}^m$$



Cone of solutions,
Konvex problem

Multi-Layer Perceptron

- Nonlinear separation
- regression
- Dense in \mathbb{R}^2 with only 1 hidden layer
- Its representational capabilities is increasing by the number of hidden layers
- Feedforward structures in the neural system, exaple: visual system



What could be represented by a simple, one layered, feed forward network called perceptron?

It is able to learn many functions, but there are some exceptions such as XOR.

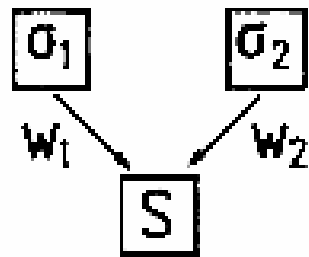


Fig. 5.2. Architecture of the XOR gate.

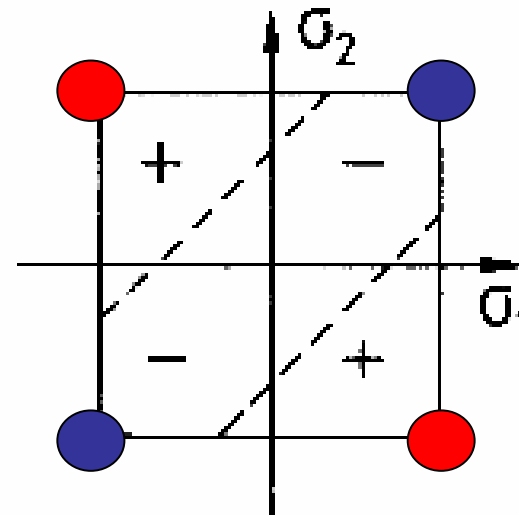


Fig. 5.3. Disconnected areas of equal output value in the space of input variables.

Problem:

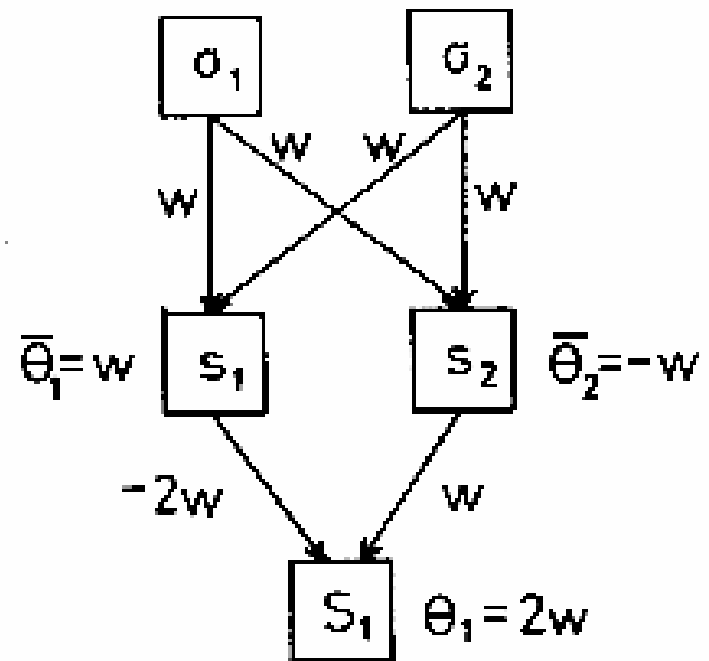
The linearly inseparable functions are more numerous as the dimension of the problem increases.

In two dimensions the problem can be transformed: this requires a two layered network

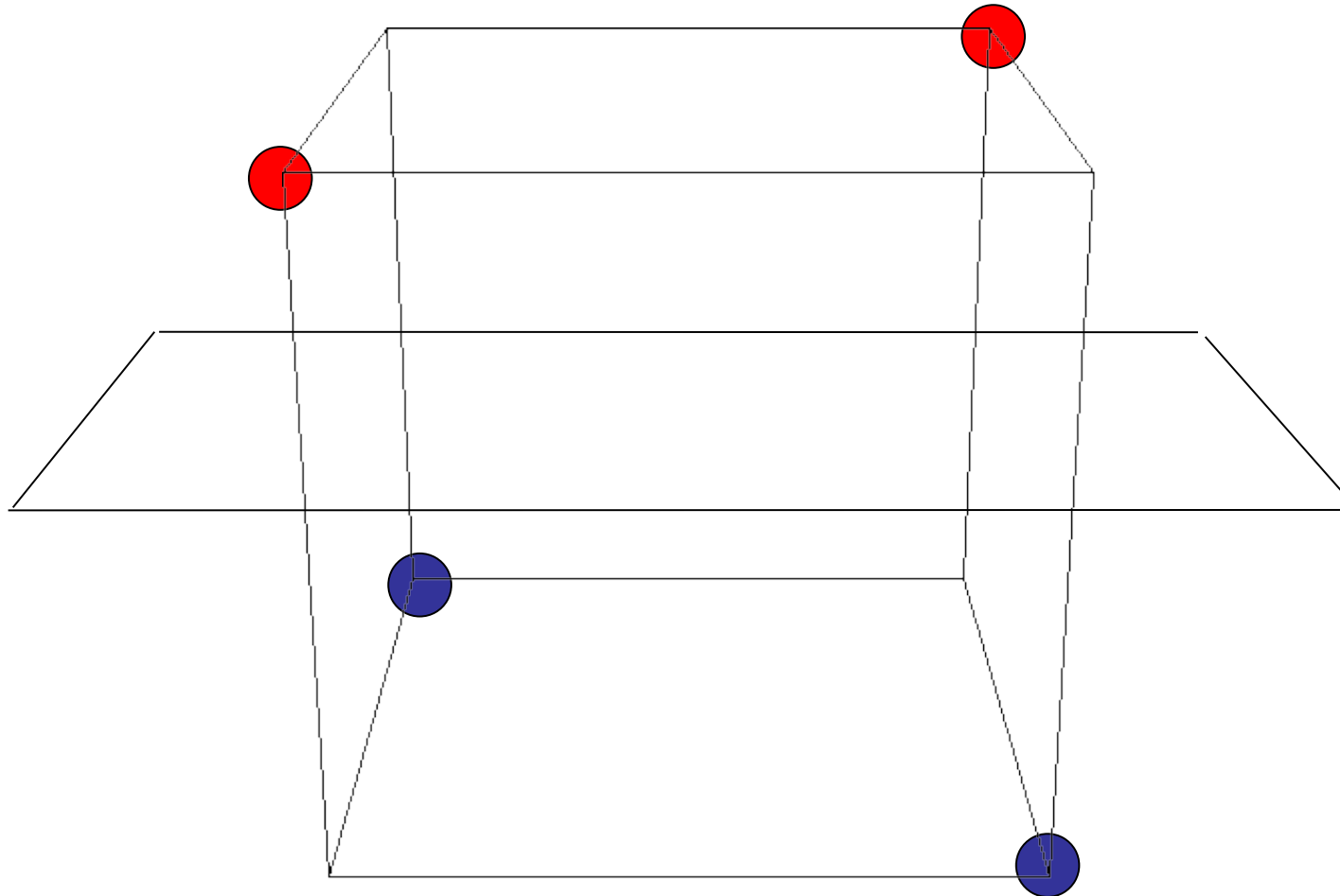
With this two layered network,
all the two dimensional
Boolean-functions can be learned.

But in higher dimensions?

The weights and the thresholds
appropriate to the XOR solution:



A possible solution: increasing the embedding dimension



In three dimension the XOR problem is linearly separable. As the embedding dimension increases, the fraction of linearly inseparable logical functions vanishes

Soft-max output layer, to represent probabilities

The logistic function maps the activation onto the $[0, 1]$ interval:

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

A softmax layer is a generalization of the logistic function for multiple variables:

$$p_i = y_i = \frac{e^{\frac{z_i}{T}}}{\sum e^{\frac{z_i}{T}}}$$

$$\frac{\partial y_i}{\partial z_i} = y_i(1 - y_i)$$

Non local, the activity of the neurons depends on the input of all neurons, not only its own input

T: temperature parameter

T \rightarrow inf: uniform distribution

T \rightarrow 0: converge to Max function

Its derivative is smooth and local,
Similar to the derivative of the logistic function

Error backpropagation

- Input: z
- Required output (target): t_n
- Actual output: y
- Partial derivatives of the error function:
- Gradient method, which converges to a local minimum of the error function.

$$z = \mathbf{x}\mathbf{w} + b$$

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

$$\frac{\partial E}{\partial w_{bi}} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w_{bi}}$$

$$\frac{\partial E}{\partial y} = t^n - y(\mathbf{x})$$

$$\frac{dy}{dz} = y(1 - y)$$

$$\frac{\partial z}{\partial w_i} = x_i$$

$$w_i(t+1) = w_i(t) + (y - t^n) y(1 - y) x_i$$

$$\delta_j = y_j(1 - y_j) \sum w_{jq} \delta_q$$

Slow convergence along highly correlated variables

Problem:

In case of strongly correlated variables, the gradient is many times almost Perpendicular to the direction of the minimum.

Possible solutions:

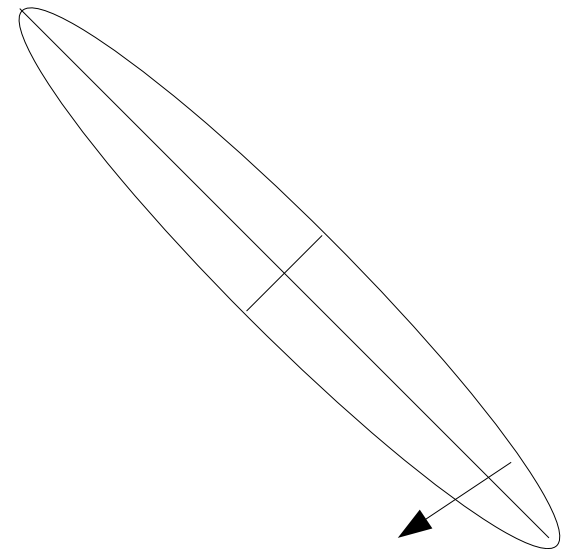
Momentum method,

Hessian matrix

Hessian free optimalizator

Conjugate gradient method

Adaptive steplength



Recurrent Networks

Reservoir computing:

Context reverberation

Echo state network

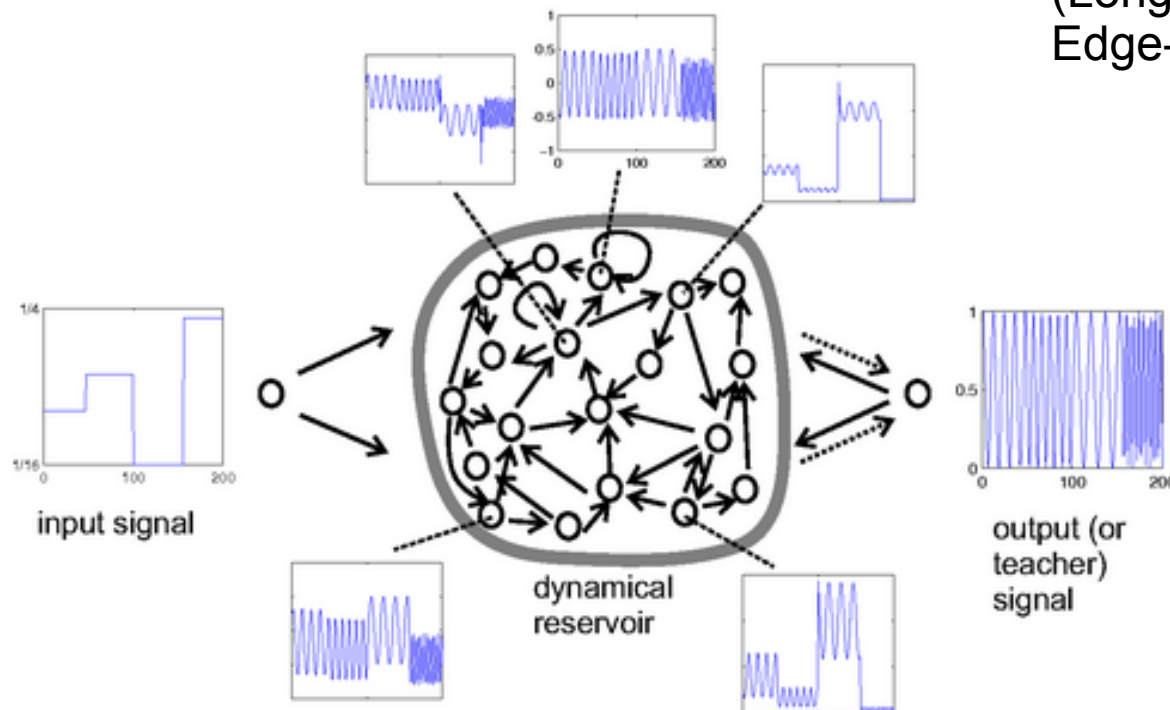
Fluid state networks

Possible learning techniques:

Error-back propagation in time

Fixed reservoir,
trainable perceptron

(Long short term memory:
Edge-of-chaos)



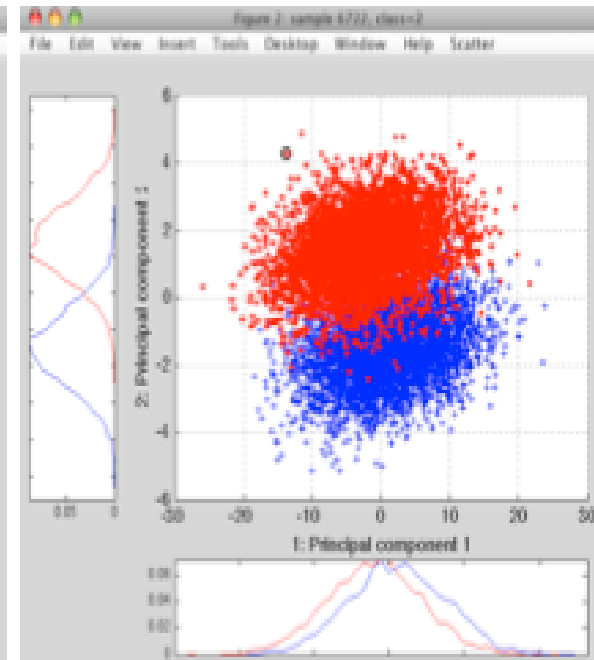
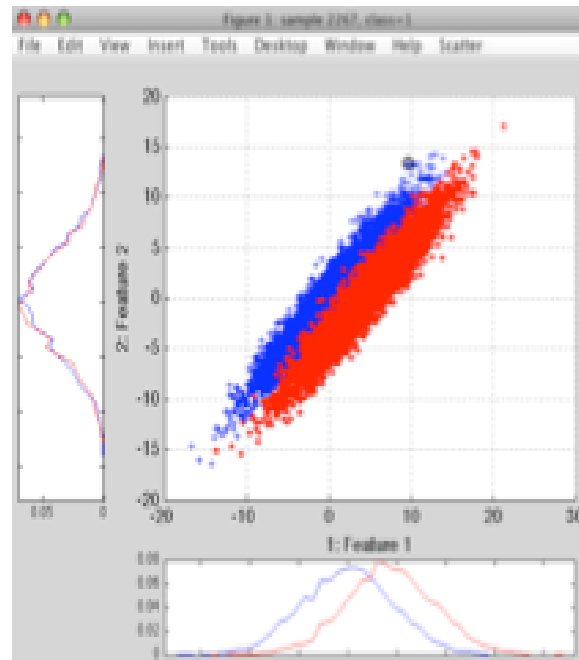
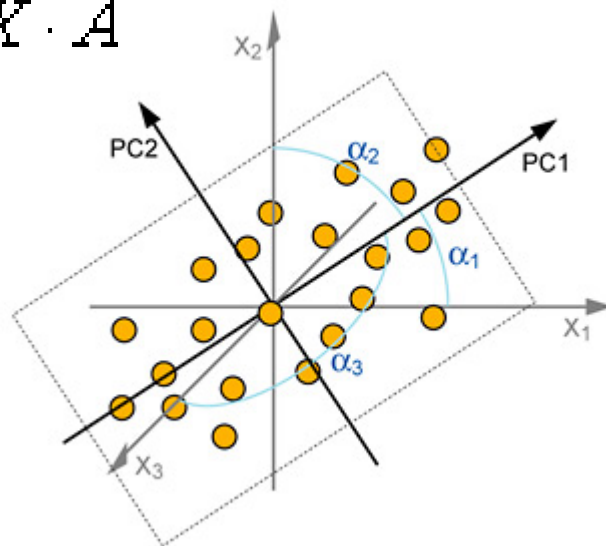
Principal component analysis

$$\bar{x}_k = \frac{1}{K} \sum_{k=1}^K x_{k\alpha} \quad \sigma_k = \sqrt{\frac{1}{K-1} \sum_{k=1}^K (x_{k\alpha} - \bar{x}_k)^2}$$

$$S = \left(\text{covar}[y_i, y_j]_{i,j=1}^N \right) \approx \left(\frac{1}{K-1} \sum_{k=1}^K y_k y_k^T \right)_{i,j=1}^N = \frac{1}{K-1} Y^T \cdot Y$$

$$S \cdot X' = \lambda \cdot X'$$

$$Z = X \cdot A$$



Principal component network, derivation of Oja's rule:

$$w_i(n+1) = w_i(n) + \eta y(\mathbf{x}(n))x_i(n)$$

$$w_i(n+1) = \frac{w_i(n) + \eta y(\mathbf{x}(n))x_i(n)}{(\sum_{j=1}^m [w_j(n) + \eta y(\mathbf{x}(n))x_j(n)]^p)^{1/p}}$$

$$w_i(n+1) = \frac{w_i(n)}{(\sum_j w_j^p)^{1/p}} + \eta \left(\frac{y(n)x_i(n)}{(\sum_j w_j^p)^{1/p}} - \frac{w_i(n) \sum_j y(n)x_j(n)w_j(n)}{(\sum_j w_j^p)^{1+1/p}} \right) + O(\eta^2)$$

$$y(\mathbf{x}(n)) = \sum_{j=1}^m x_j(n)w_j(n) \qquad \|\mathbf{w}\| = (\sum_{j=1}^m w_j^p)^{1/p} = 1$$

$$w_i(n+1) = w_i(n) + \eta y(n)(x_i(n) - w_i(n)y(n))$$

The Oja-rule and the principal component analysis

y output
x input
w weight matrix

$$y = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$$

Oja's rule

$$\Delta \mathbf{w} = \alpha (\mathbf{x}y - y^2 \mathbf{w})$$

Substituting y:

$$\Delta \mathbf{w} = \alpha (\mathbf{x} \mathbf{x}^T \mathbf{w} - \mathbf{w}^T \mathbf{x} \mathbf{x}^T \mathbf{w} \mathbf{w})$$

Assuming $\text{mean}(\mathbf{x}) = 0$ & averaging over the input $\Rightarrow \mathbf{x} \mathbf{x}^T = \mathbf{C}$

The convergence point: $\Delta \mathbf{w} = 0$

$$\mathbf{C} \mathbf{w} - \mathbf{w}^T \mathbf{C} \mathbf{w} \mathbf{w} = 0$$

This is an eigenvector equation!

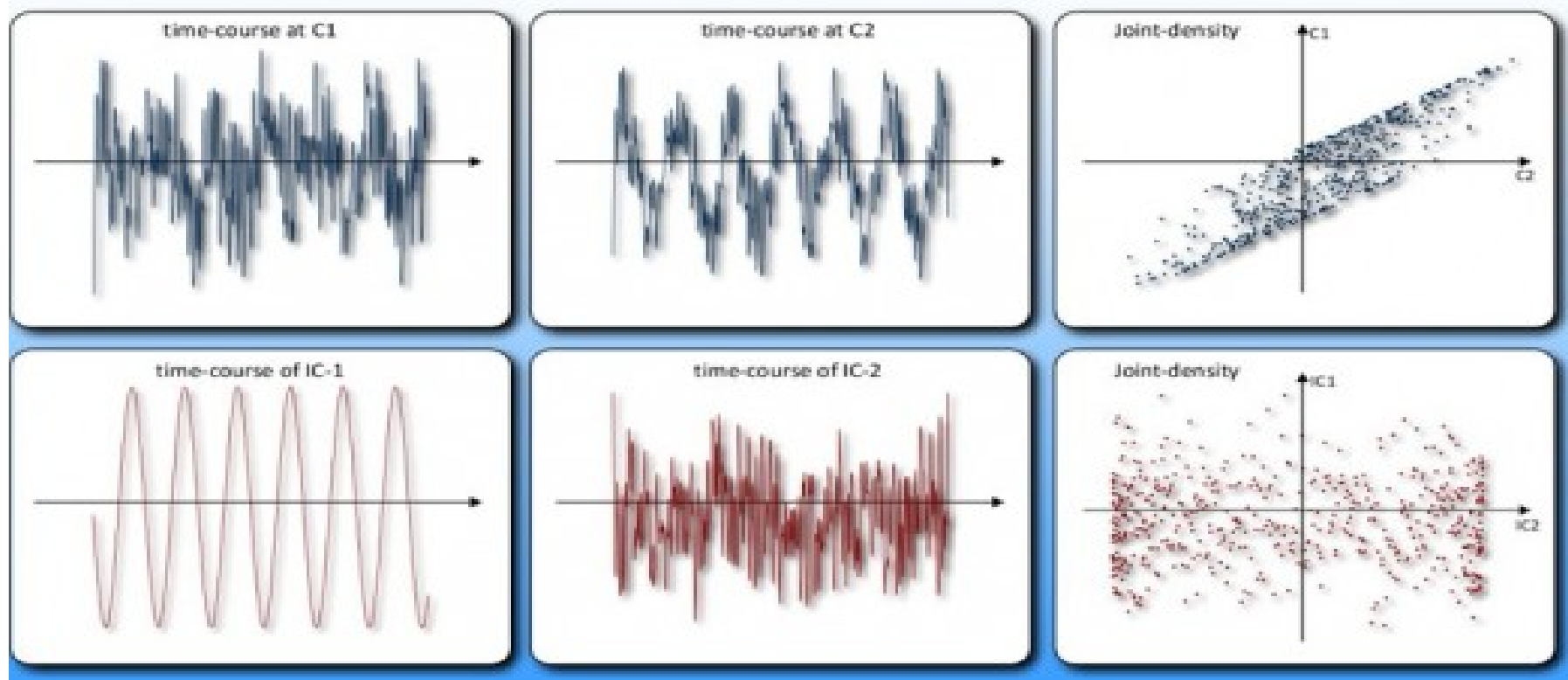
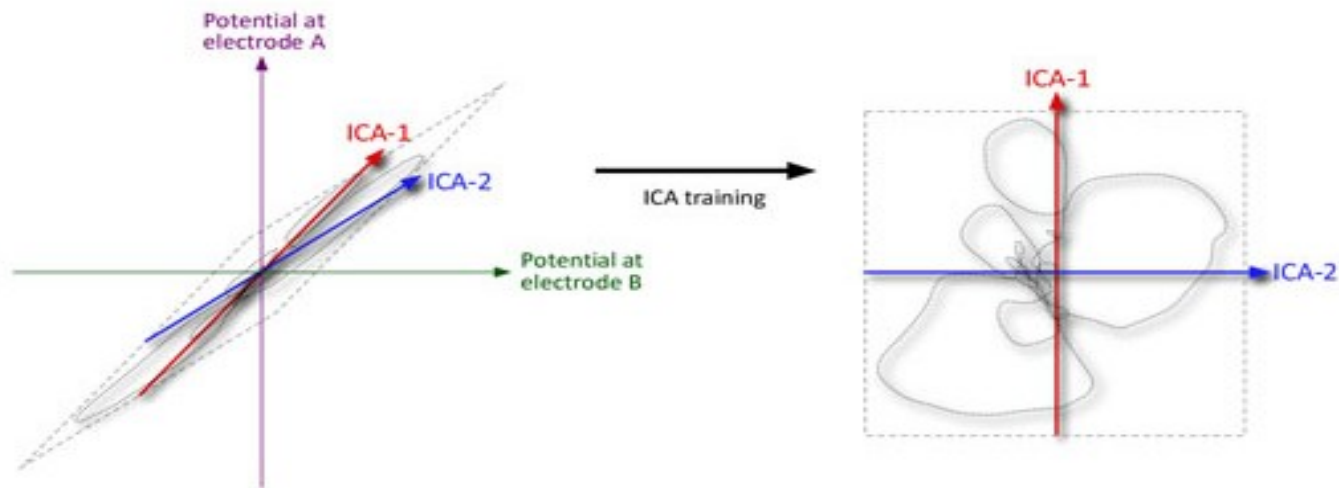
$\mathbf{w}^T \mathbf{C} \mathbf{w}$ is scalar!

E modification
for Independent
Component analysis

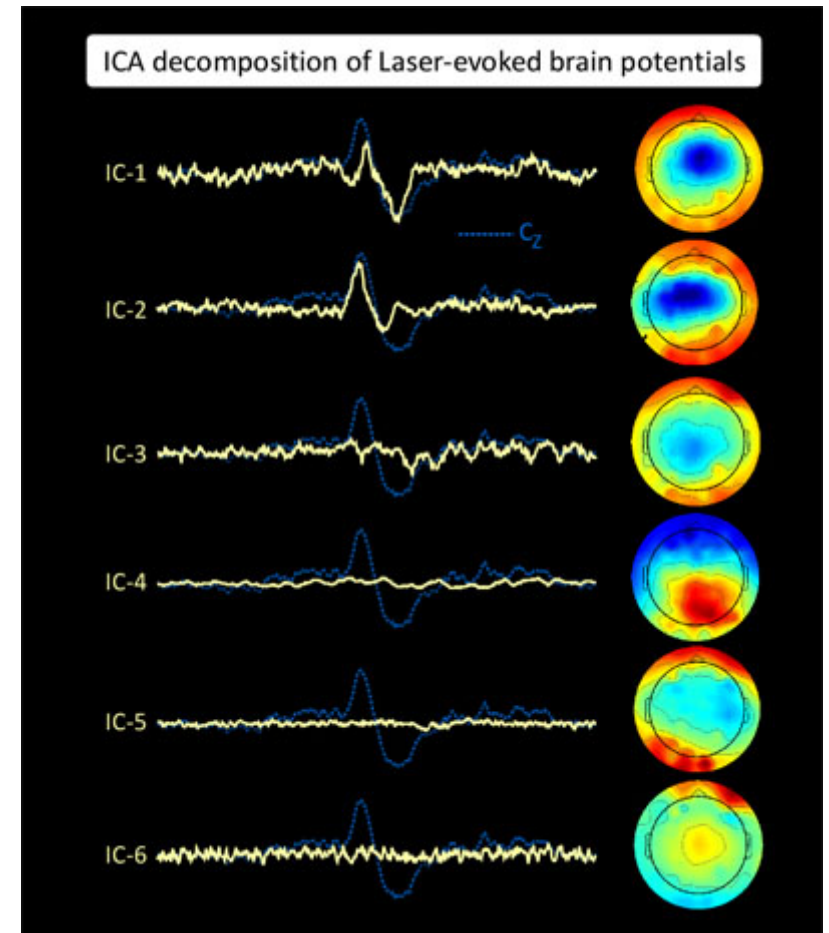
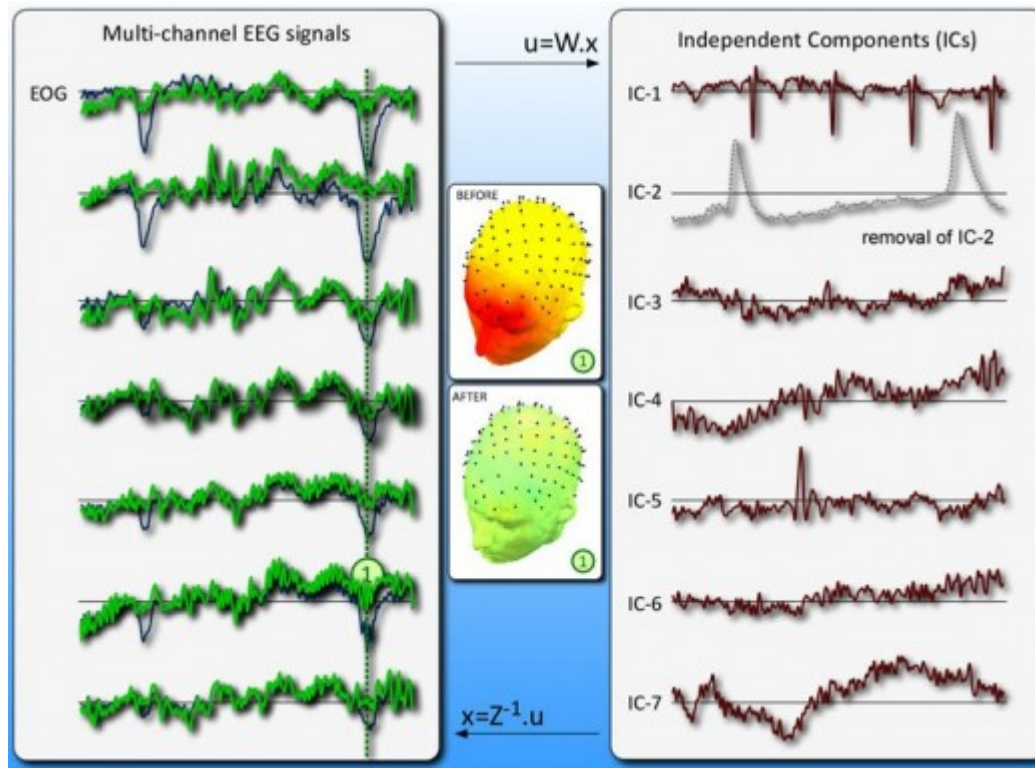
$$\Delta \mathbf{w} = \alpha (\mathbf{x}y^3 - \mathbf{w})$$

If $\mathbf{C} = \mathbf{I}$

Independent component analysis (ICA)

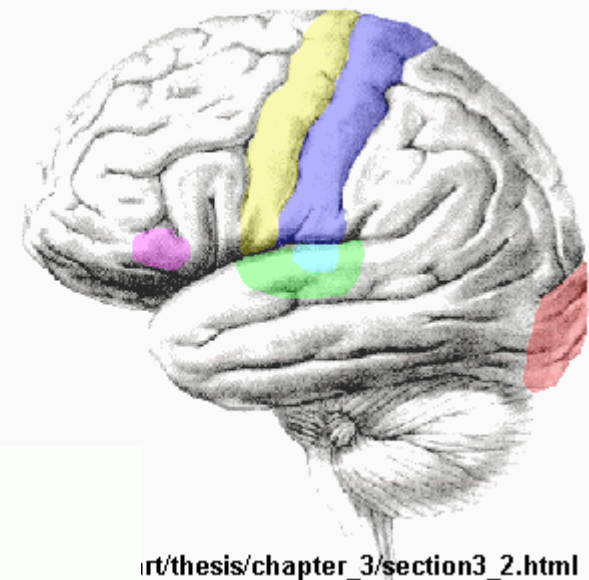


Independent component analysis (ICA)

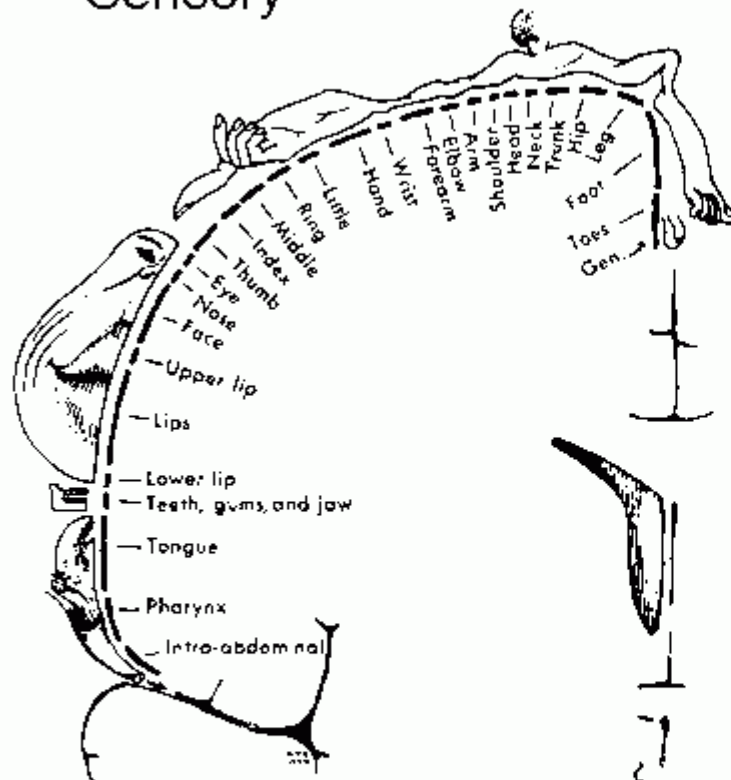


The somatosensory map

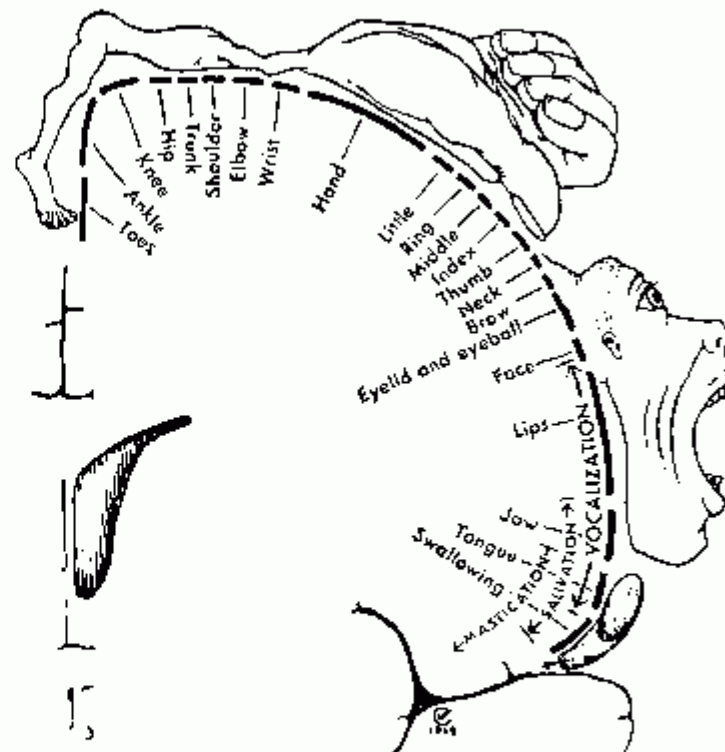
- Somatosensory
- Visual
- Auditory
- Gustatory
- Olfactory
- Motor



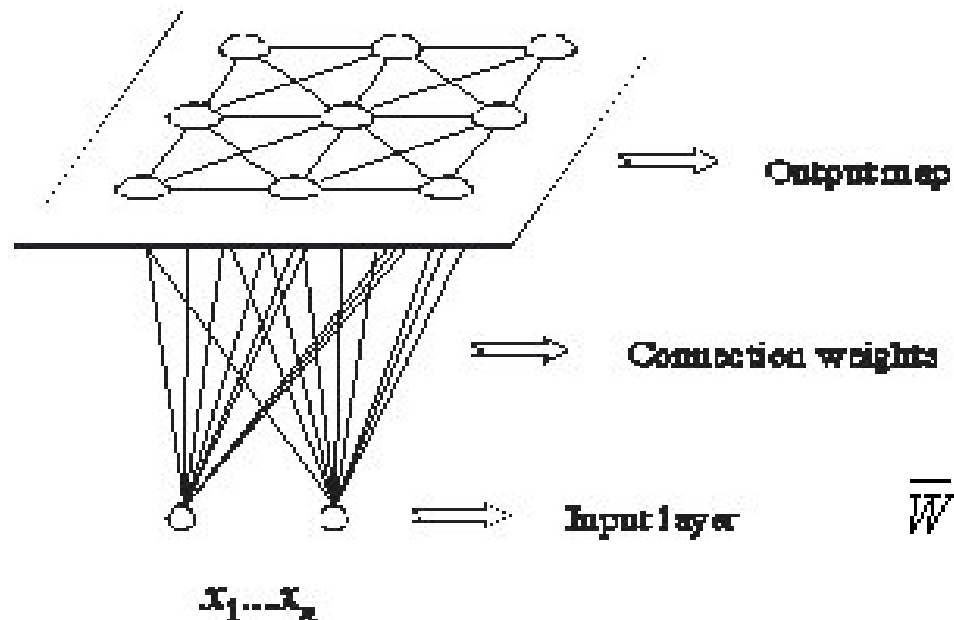
Sensory



Motor



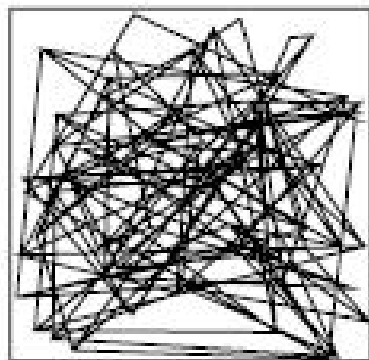
Kohonen's self-organizing map



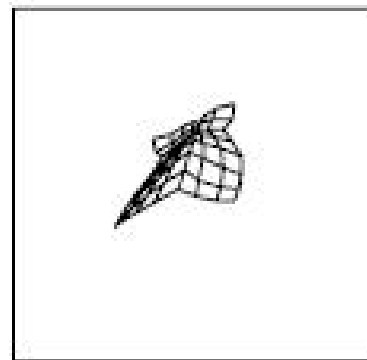
$$D_{ik} = \sqrt{\sum_{n=1}^N (x_{kn} - w_n^i)^2}$$

Winner take all

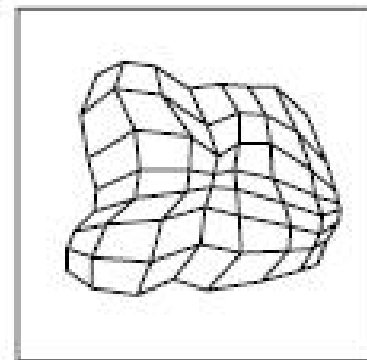
$$\bar{W}^i(t+1) = \bar{W}^i(t) + \eta(t) N(c, r) [\bar{X}_k - \bar{W}^i(t)]$$



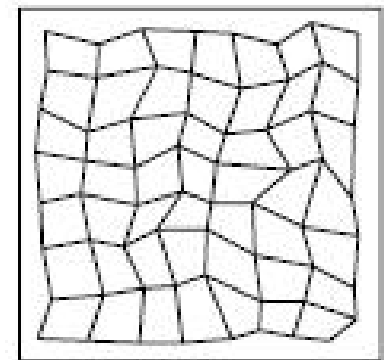
Iteration 0



Iteration 200



Iteration 600



Iteration 1900

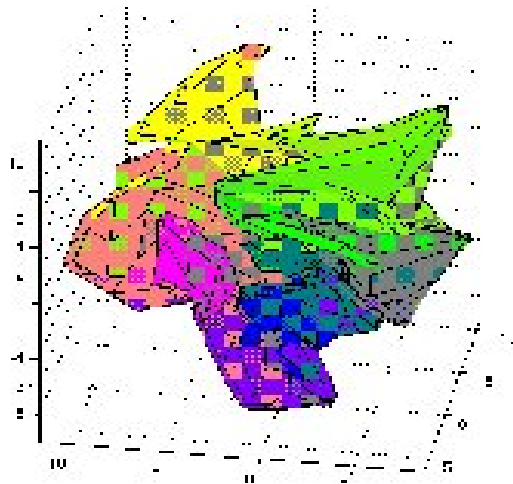
Kohonen's self-organizing map

Generates feature maps

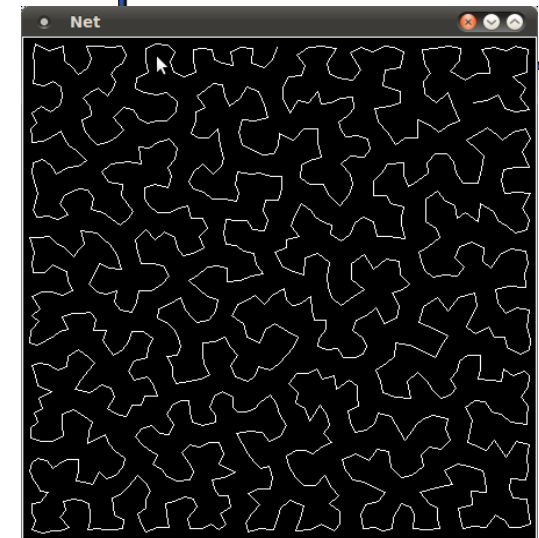
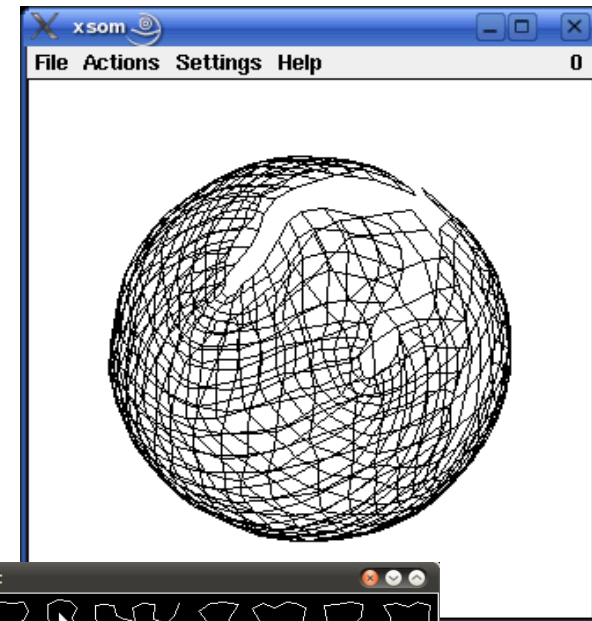
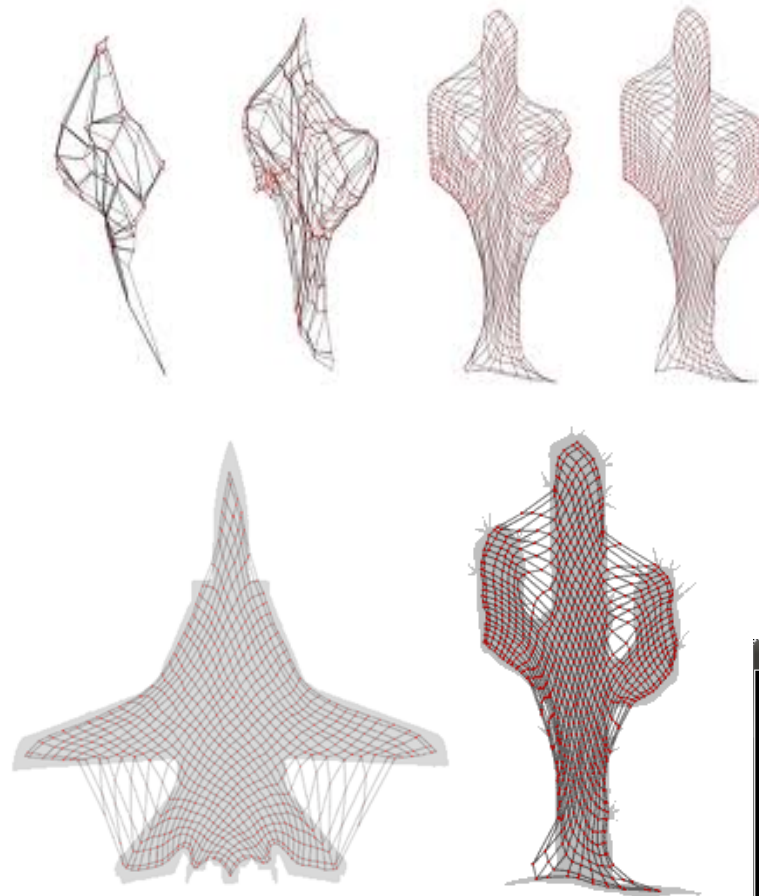
Captures the structure in the inputs



a) 2D projection

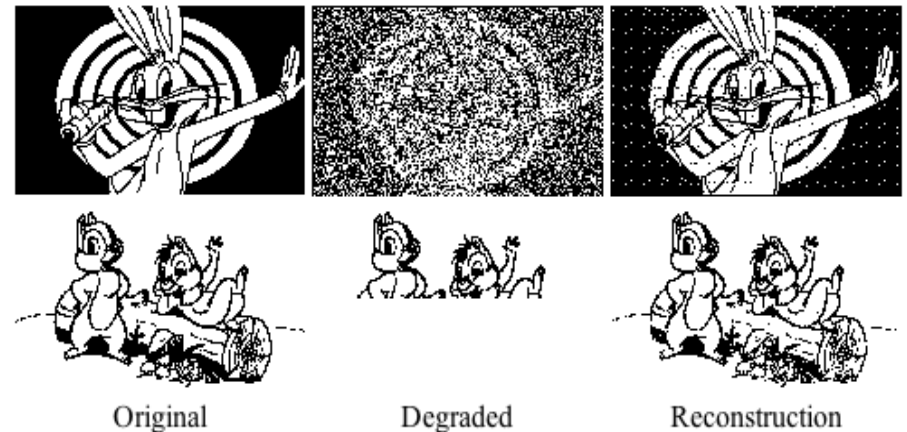
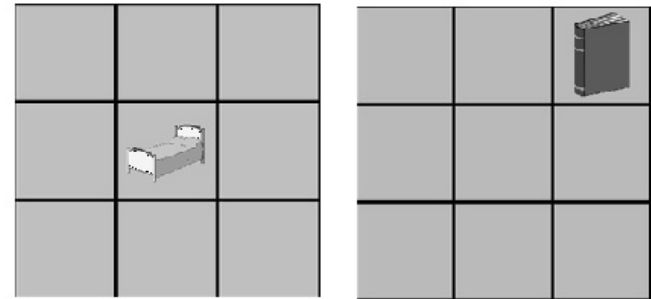


b) 3D projection



Associative memory

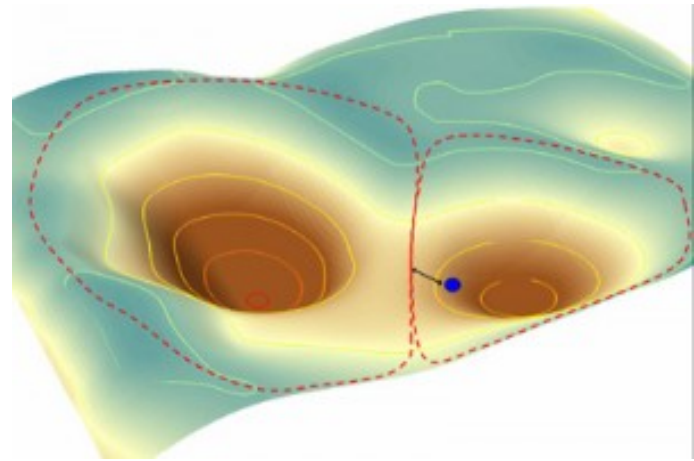
- Heteroassociation
 - Exemple: place-object
- Autoassociation
 - From partial pattern the original
- Difference between the memory of a computer and an AM is the addressing



- Capacity: how many patterns can be stored and retrived (non-unique definition)
- Stability: for each patterns, we want the most simmlar (non-unique definition)

Attractor networks

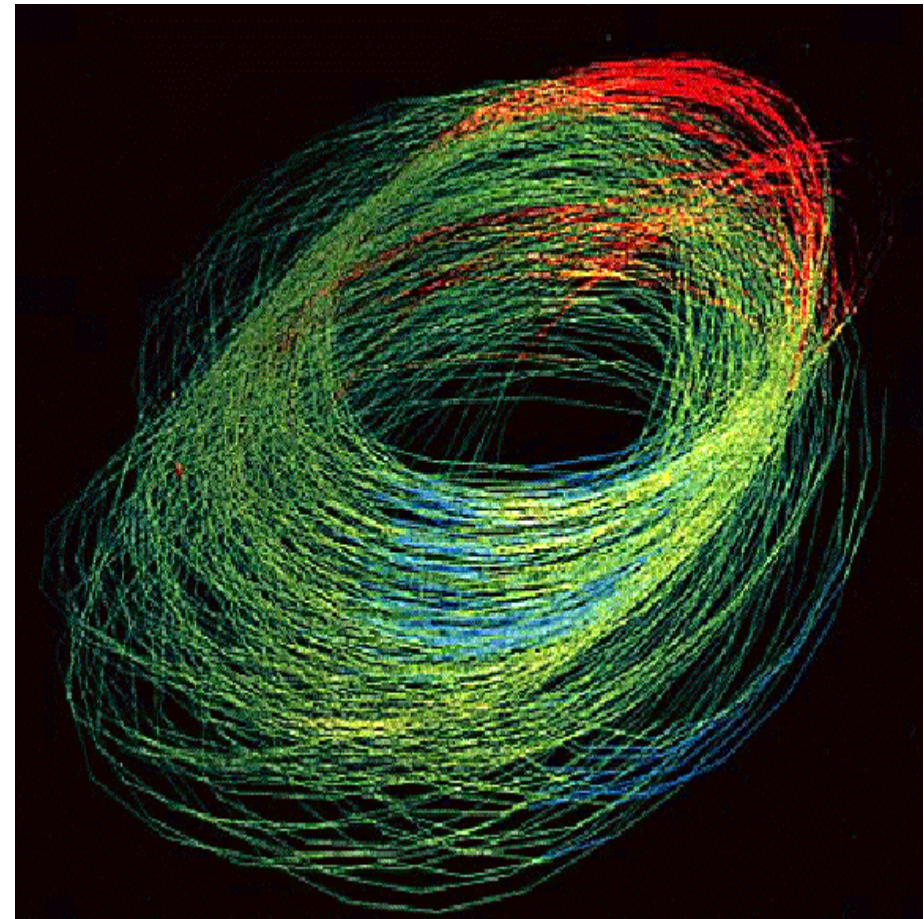
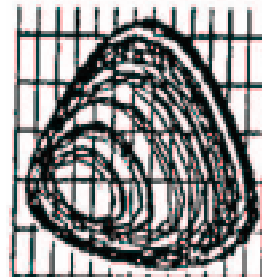
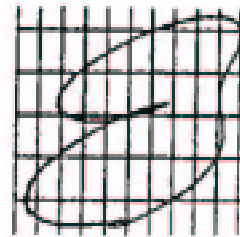
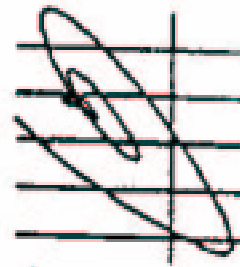
- Types of attractors
 - Fixed points
 - Periodic
 - Chaotic
- Basin of attractions
- Implementation: recurrent neural networks
- Attractor formation: synaptic weights
 - Offline learning
 - Online learning
 - One-shot learning
- Retrieval: convergence from an initial condition



MODIFICATIONS OF THE CLASSICAL SCENARIO

Fixed points vs. strange attractors

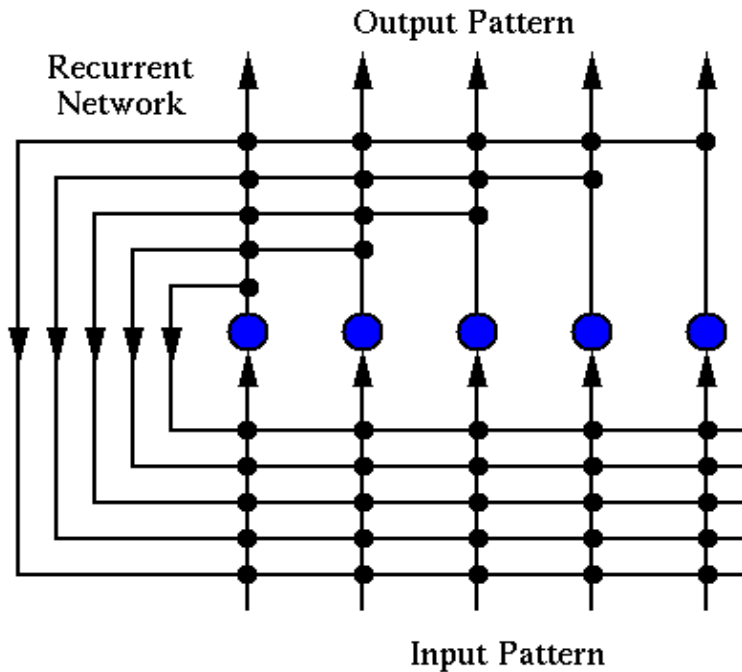
*Not only fixed points
but also limit cycle
or strange attractors
might be involved*



- neurophysiological experiments, theoretical studies
- structural conditions for the possible mechanism for the generation of rhythms and chaos can be given based on the notions of qualitative stability and instability

W. J. Freeman

Hopfield-networks



- Associative memory
- Binary MCP-neurons
- Patterns: binary vectors
- Symmetric weight matrix
- Dale's law is violated: a cell can be excitatory and inhibitory parallel
- Recurrent networks in the brain (dominantly):
hippocampus CA3 region, ...
Hebb's rule

- Offline learning patterns to learn:
 $\{s^1 \dots s^N\}$

$$W_{ij} = \frac{1}{N} \sum_n s_i^s s_j^s$$

$$\mathbf{x}^{t+1} = \text{sgn}(\mathbf{W} \mathbf{x}^t - \boldsymbol{\theta}) \quad x_k^{t+1} = \text{sgn}\left(\sum_i^K W_{ik} x_i^t - \theta_k\right)$$

- Iterations: synchronous or sequential

The dynamics of the Hopfield network

- Stability-analysis of a nonlinear system: The definition of the „energy” of a state by introducing a Lyapunov-function:
 - Bounded
 - The iterative dynamics always decreases (or increases) it.

It it exists, that we showed, that the system will converge to a fixed point for every input patterns

- The Lyapunov-function of the Hopfield network:

$$E = -\frac{1}{2} \mathbf{x}^T \mathbf{W} \mathbf{x} - \boldsymbol{\theta} \mathbf{x}$$

- The learning builds attractors at the stored patterns, but not only there (spurious attractors)
- The HN can be used to optimize problems with quadratic forms

The capacity of the HN

- Information theoretical capacity

- The patterns can be considered as a set of random variables from Bernoulli-distribution

$$P(s_i^n = 1) = P(s_i^n = 0) = 0.5$$

- The convergence is requires for all patterns

$$\lim_{n \rightarrow \infty} P(s^a = \text{sgn}(\mathbf{W}\mathbf{s}^a)) = 1 \quad \forall a = 1 \dots M$$

- Than in can be shown (with many approximation) that

$$M \approx \frac{N}{2 \log_2 N}$$

- In case of the CA3 hippocampal subregion

- Cc. 200000 neurons, cc. 6000 patterns to store

- Different estimations

- Let us consider the sparsity of the patterns

$$P(s_i^n = 1) = \alpha$$

$$M \approx N \frac{1}{\alpha \log_2 \frac{1}{\alpha}}$$

Boltzmann machine

- To represent probability distributions – statistical interdependences between variables
- Stochastic state transition

$$\mathbf{I} = \mathbf{W}\mathbf{u} + \mathbf{M}\mathbf{v}$$

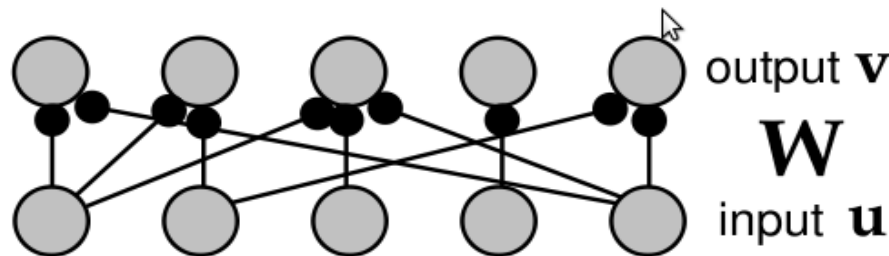
$$P(v_a^{t+1} = 1) = \frac{1}{1 + e^{-I_a}}$$

- The limit distribution

Energy:

$$E(\mathbf{v}) = -\mathbf{v}^T \mathbf{W}\mathbf{u} + \frac{1}{2} \mathbf{v}^T \mathbf{M}\mathbf{v}$$

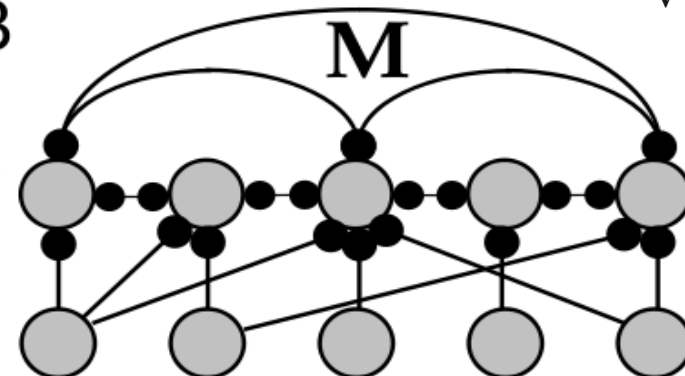
A



Boltzmann-distribution: $e^{-E(\mathbf{v})}$

$$P(\mathbf{v}) = \frac{e^{-E(\mathbf{v})}}{\sum_{\mathbf{v}} e^{-E(\mathbf{v})}}$$

B



Learning with Boltzmann machine

- Supervised learning, only for \mathbf{W} , similarly \mathbf{M}
- Error: Kullback-Leibler-divergency between the targeted and actually generated distribution

$$D_{KL}[P(\mathbf{v}|\mathbf{u}), P(\mathbf{v}|\mathbf{u}, \mathbf{W})] = \sum_{\mathbf{v}} P(\mathbf{v}|\mathbf{u}) \ln \frac{P(\mathbf{v}|\mathbf{u})}{P(\mathbf{v}|\mathbf{u}, \mathbf{W})}$$

independent of \mathbf{W}

average over the inputs: (instead of the $P(\mathbf{v}|\mathbf{u})$ weighted outputs)

$$\langle D_{KL} \rangle = -\frac{1}{N_s} \sum \ln P(\mathbf{v}^m | \mathbf{u}^m, \mathbf{W}) - K$$

- Gradient descent – for only one input

$$\frac{\partial \ln P(\mathbf{v}^m | \mathbf{u}^m, \mathbf{W})}{\partial W_{ij}} = v_i^m u_j^m - \sum_{\mathbf{v}} P(\mathbf{v} | \mathbf{u}^m, \mathbf{W}) v_i u_j^m$$

from the Boltzmann-distribution

- Delta-rule – the average for the all possible outputs is approximated with the actual one

$$W_{ij} \rightarrow W_{ij} + \epsilon_w (v_i^m u_j^m - v_i(\mathbf{u}^m) u_j^m)$$

two phases:

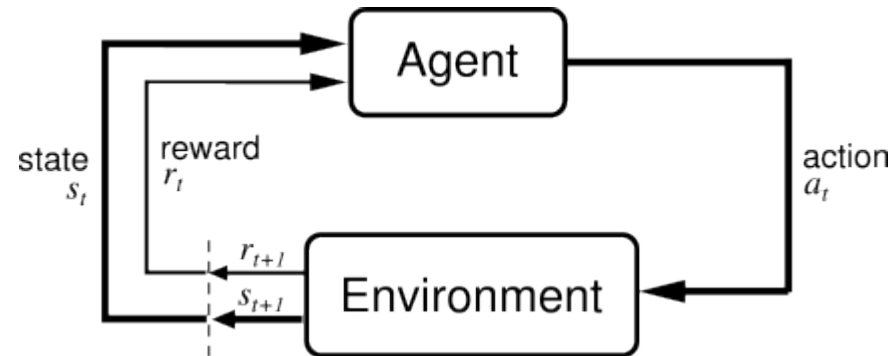
Hebbian

anti-Hebbian

- Unsupervised

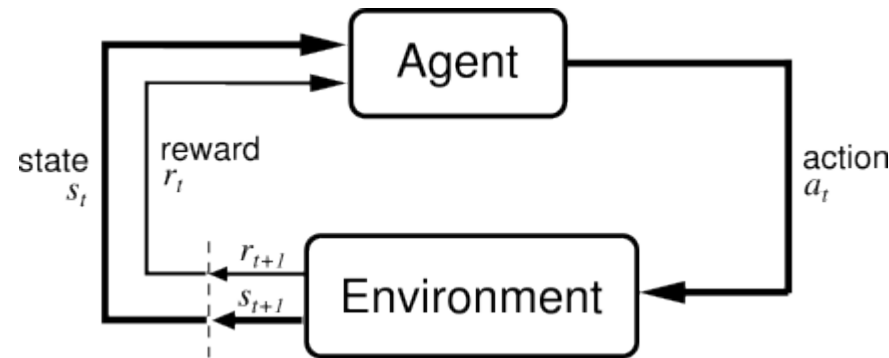
$$D_{KL}[P(\mathbf{u}), P(\mathbf{u}, \mathbf{W})]$$

Reinforcement Learning



- State space: the possible values of the sensory (or other input) variables
- Reward: in certain states we get information from the success
- Actions: the agent realizes a state transition (at least tries)
- aim: to maximize the reward in long run
- Value function: the utility of the states
- Representation of the value function:
 - Table (machine learning)
 - General function approximator, for example a feed-forward neural network
 - (embedden supervised learning)

Reinforcement learning



Solutions

- Model-based: state and state-transition representation value function and action policy

Types: Direct utility estimation (DUE)

Adaptive dynamic programming (ADP)

Temporal difference (TD)

- Model-free:

Q-learning: state-action pair value association

Strategy searching

Temporal difference learning

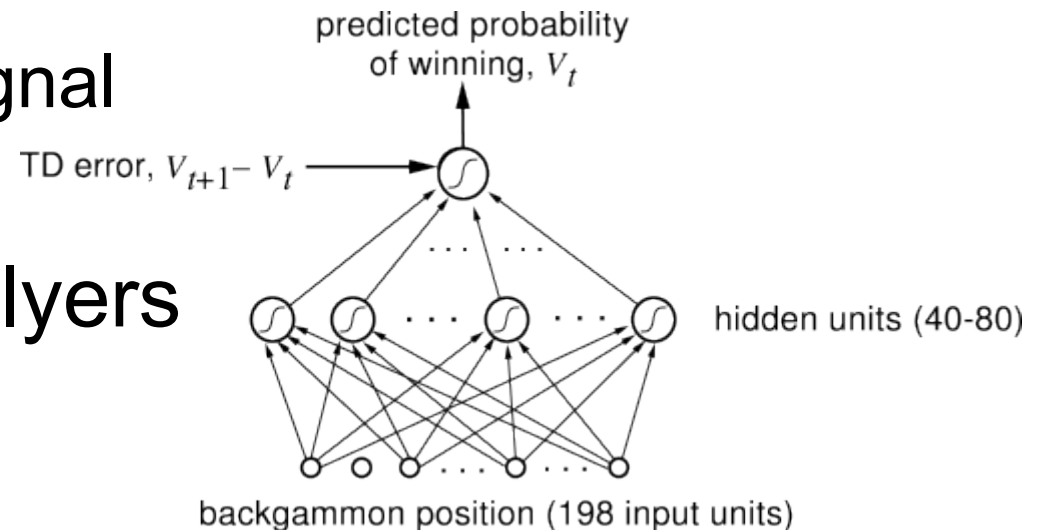
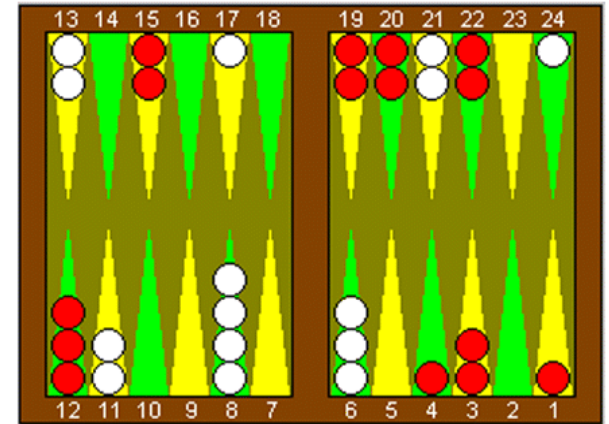
- Using the prediction error for the learning
- Updateing the value function in a neural representation:

$$w(\tau) \rightarrow w(\tau) + \epsilon \delta(t) u(t-\tau) \qquad \delta(t) = \sum_{\tau} r(t+\tau) - v(t)$$

- Calculation of the prediction error:
 - (IN principle we need to know the full reward in the future)
 - On step local approximation
$$\sum_{\tau} r(t+\tau) - v(t) \approx r(t) + v(t+1)$$
- If the environment is suitable, it converges to an optimal strategy
- The error can be back prpagated to the prvious states as well (eligibility trace, simmilar to the error back-propagation)
- Acion selection: exploration vs. exploitation

TD with feed-forward neural network

- Gerald Tesauro: TD-Gammon
 - Feedforward network
 - Input: the states which can be achieved by the possible actions
 - Output: values (winning probabilities)
- In each step, the error is calculated
 - Based on the reward signal
- Result: comparable with the best human players

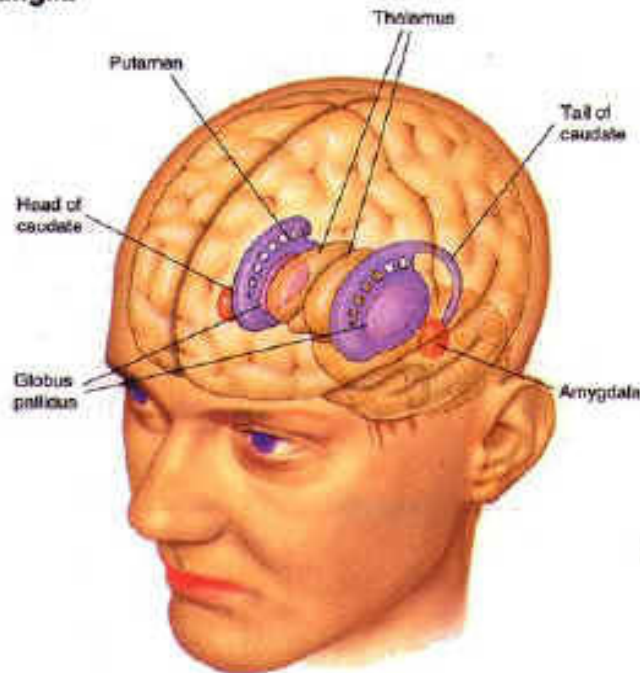


The effect of reward in dopaminergic cell of basal ganglia

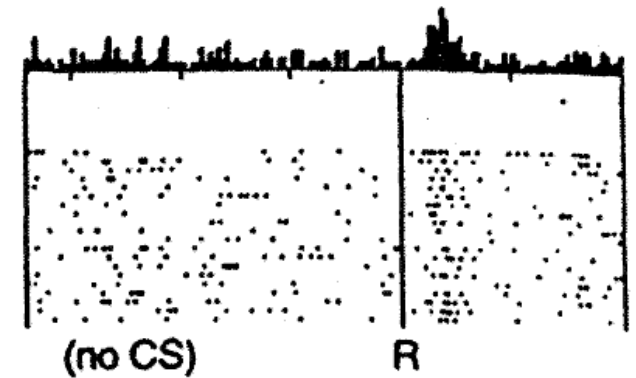
An interpretation:

Dopamine cells signal the difference between the expected and received reward.

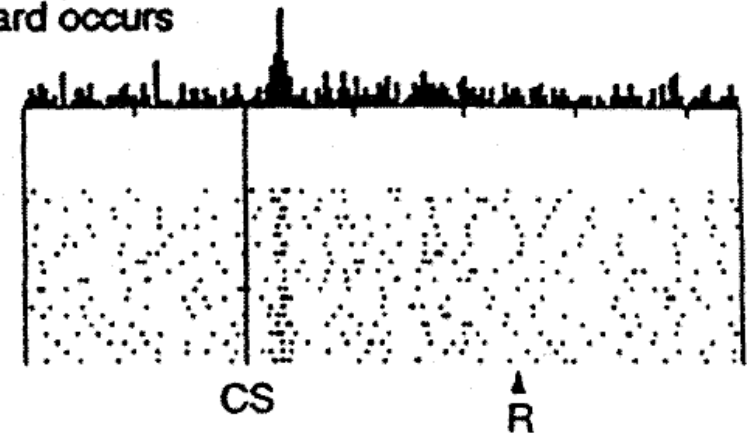
► The Basal Ganglia



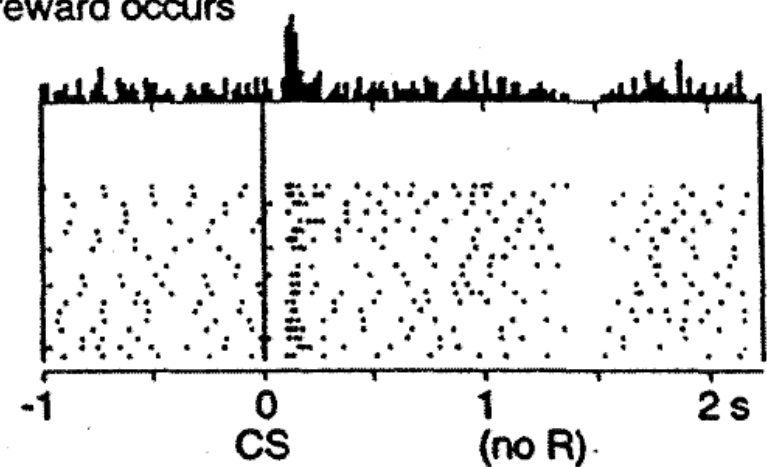
No prediction
Reward occurs



Reward predicted
Reward occurs



Reward predicted
No reward occurs



Problems to solve in learning systems

- Bias-variance dilemma
 - Structural error: the model (even with optimal parameters) can differ from the function to approximate (eg. Linear model fitting for cubic data)
 - Approximation error: infinite datapoints are needed for precise parameter tuning
- Accuracy vs. Predictive power
 - The models with too much parameters fits well, but generalize poorly
 - May have lower explanatory ability

