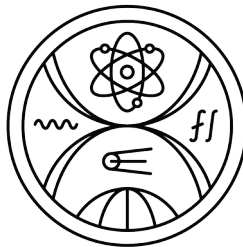# Introduction to Computational Intelligence

## Unsupervised feedforward neural networks

**Igor Farkaš**
Centre for Cognitive Science
Comenius University Bratislava

# Feature extraction

- starts from an initial set of measured data

- the number of input variables for a dataset = dimensionality

- builds derived values (features) intended to be informative and non-redundant

- linear or nonlinear FE methods

- facilitating the subsequent learning and generalization steps

- in some cases leading to better human interpretations

- feature extraction entails dimensionality reduction

- nonlinear methods offer more variability of solutions

- FE methods may be based on neural networks

- We show two examples: PCA (FE) and SOM (data visualisation)

# Hebbian learning – a single neuron

- Canadian psychologist Donald Hebb (1949) postulated:

*"When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased."*
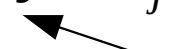
- Let us assume one linear neuron with $n$ inputs:

$$y = \sum_{j=1}^{n} w_j x_j = \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}$$

- According to Hebb's postulate:

$$w_j(t+1) = w_j(t) + \alpha \, y \, x_j \quad \text{for} \quad j = 1, 2, \ldots, n$$

where $\alpha$ is the learning rate.

- Oja's rule: $\quad w_j(t+1) = w_j(t) + \alpha \, y \, x_i - \alpha \, y^2 \, w_j(t)$
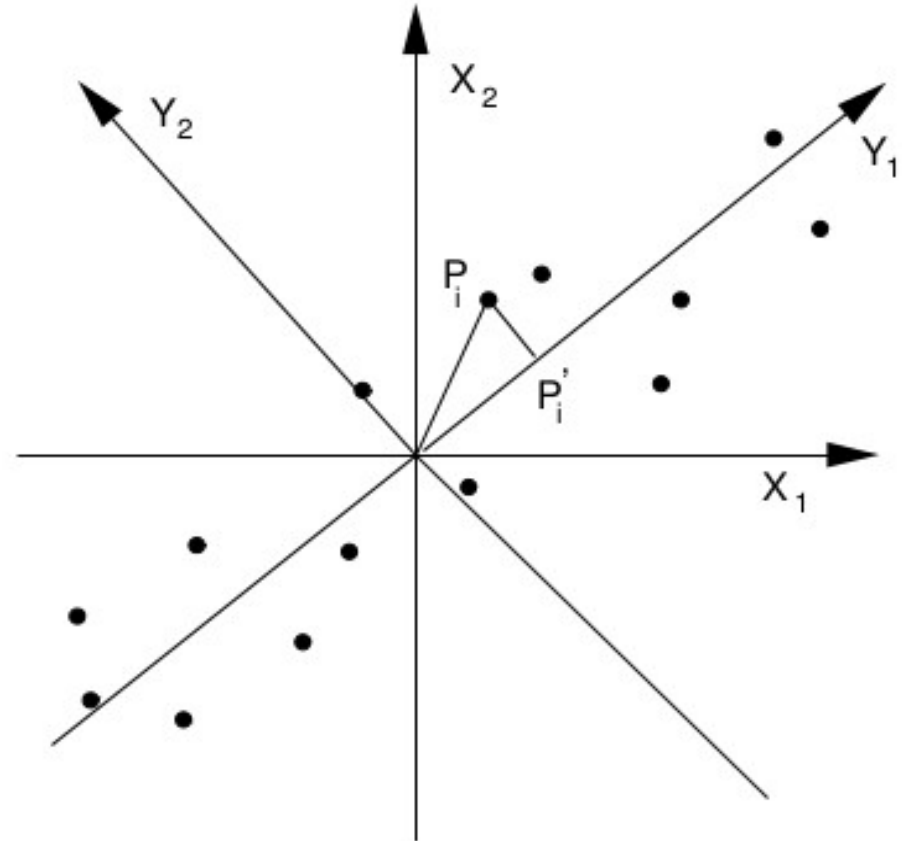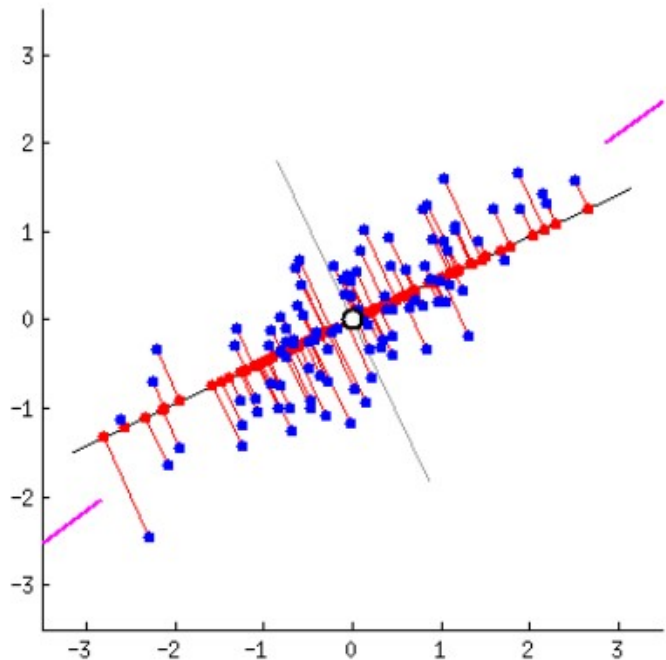
stabilizing term

# Single neuron behavior

Properties of linear neuron trained by Oja's rule on zero-meaned data vectors (i.e. $\langle \mathbf{x} \rangle = \mathbf{0}$):

- weight vector $w$ converges to match the 1st eigenvector of data covariance matrix $\mathbf{R} = \langle \mathbf{x}\mathbf{x}^T \rangle$

- after convergence, $w$ maximizes $\langle y^2 \rangle$

- after convergence, the vector norm $\|w\| = 1$.

- Hence, the linear Hebbian neuron projects input vectors to the direction that maximizes the data discrimination capability (maximum information preservation).

- If more units are used, they can be trained to perform PCA (e.g. using GHA algorithm)

# Linear feature extraction

X = original space (2D)
Y = projection space (2D)

Direction Y1 captures maximum
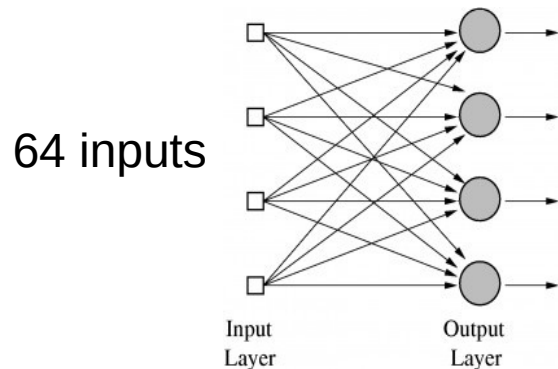variance in original 2D data.

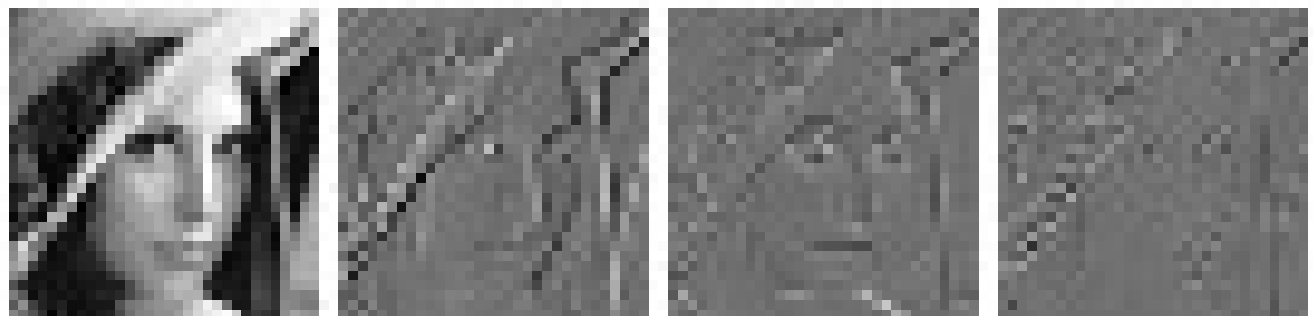# Example of PCA application

- Image of 256x256 pixels, converted to gray image (0...255)

- NN trained on sub-images 8x8 (input size)

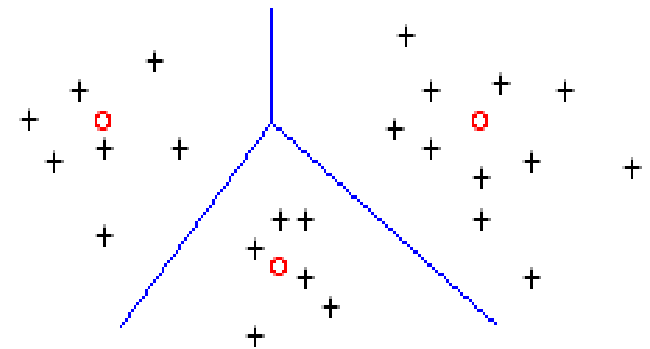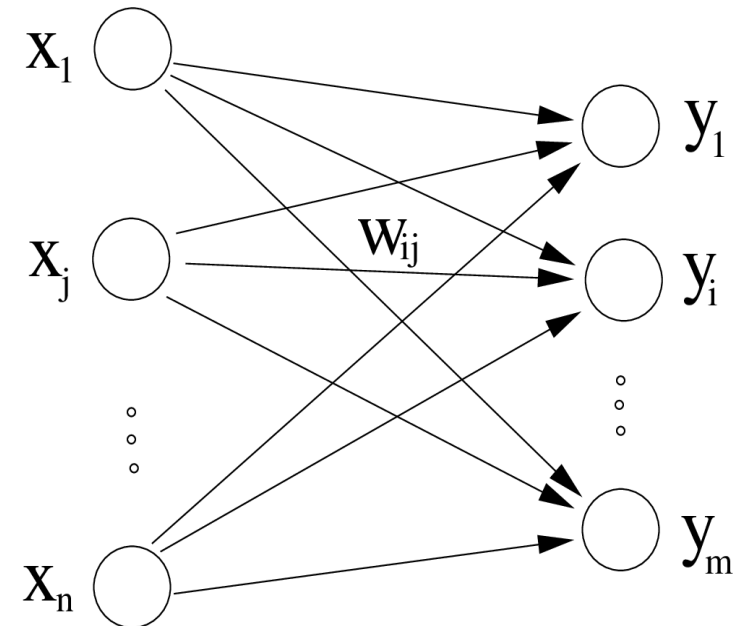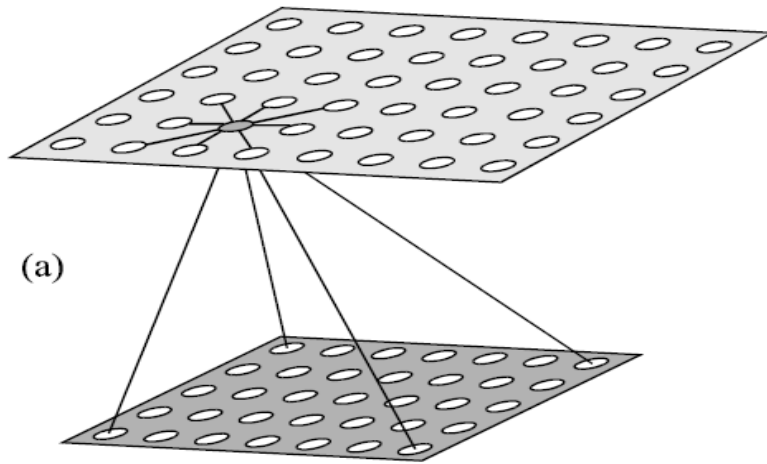- 8 largest eigenvectors (weight vectors): 



Reconstructions:



64 inputs



Input Layer    Output Layer

Neuron outputs:



6

# Competitive learning

- a kind of unsupervised learning
- Features:
  - linear neurons
  - winner: $y_c = \max_i \{ w_i^T . x \}$
    - i.e. best matching unit $c$
  - winner-take-all adaptation:
    $$\Delta w_c = \alpha(x - w_c) \quad \alpha \in (0,1)$$
    $$\| w_c \| \leftarrow 1$$
  - risk of "dead" neurons
- algorithm: in each iteration
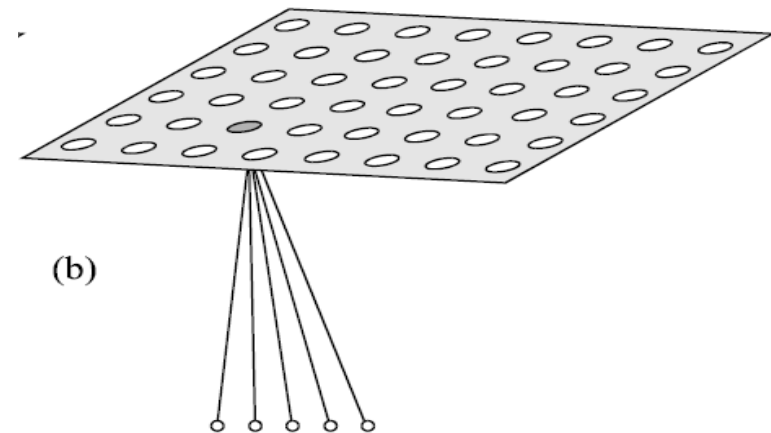  - find winner, adapt its weights
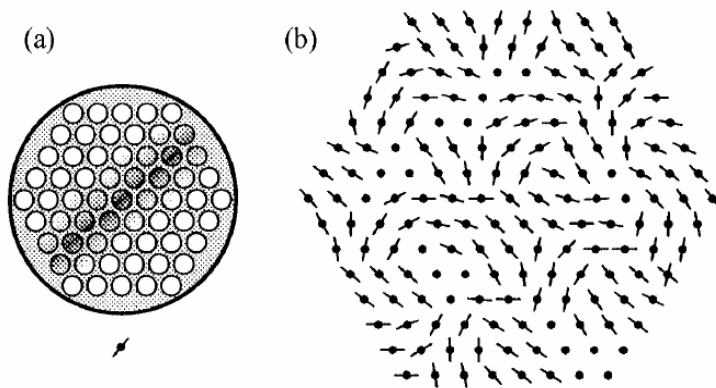- useful for clustering
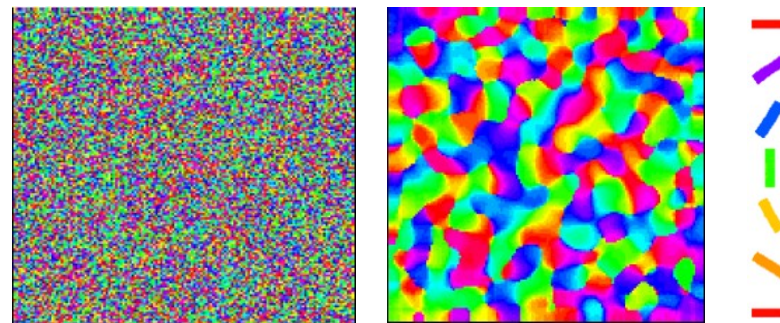
# Feature mapping

biologically motivated models

model with extracted features
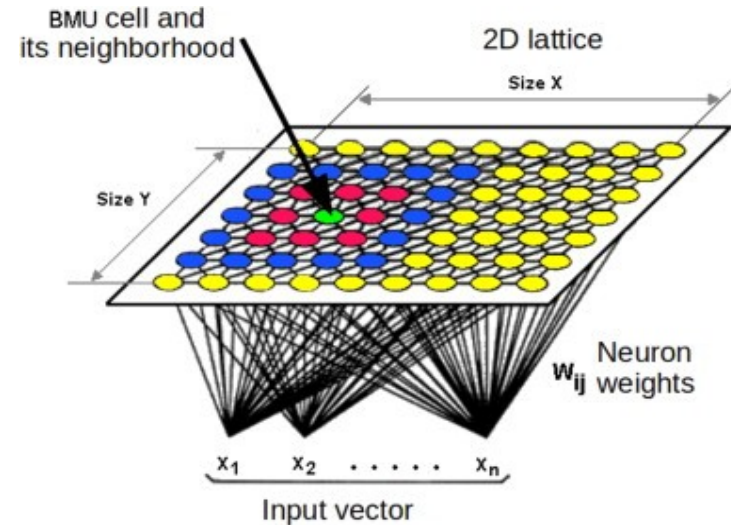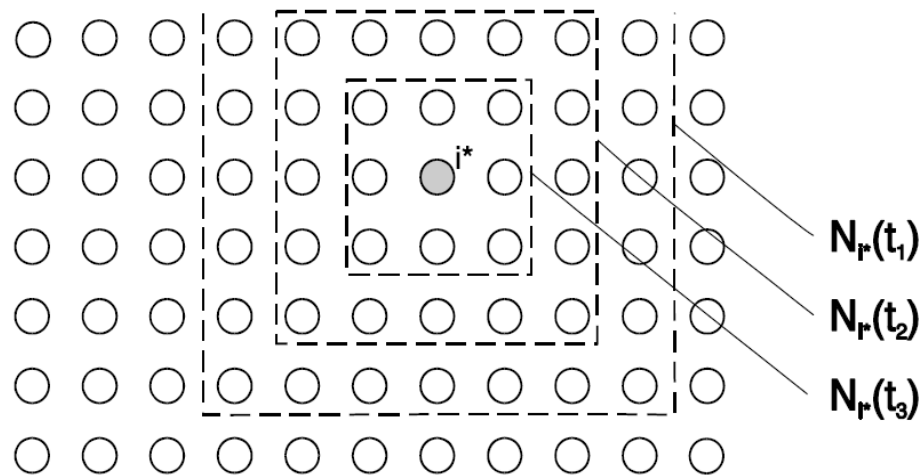


e.g. mapping from retina to cortex ->
orientation map

- introduced topology of neurons in the map
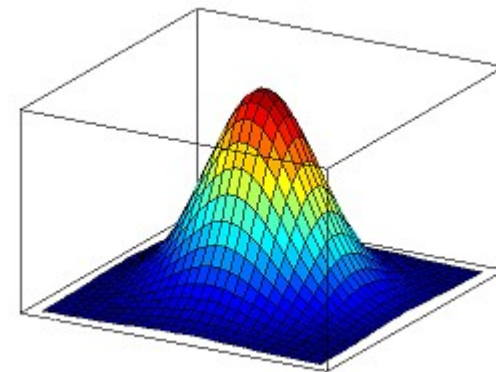- winner-take-most due to neuron cooperation

# SOM: Neighborhood function

- computationally efficient substitute for lateral interactions
- neurons adapt only within the winner neighborhood
- neighborhood radius decreases in time

- rectangular neighborhood:



$N_r(t_1)$

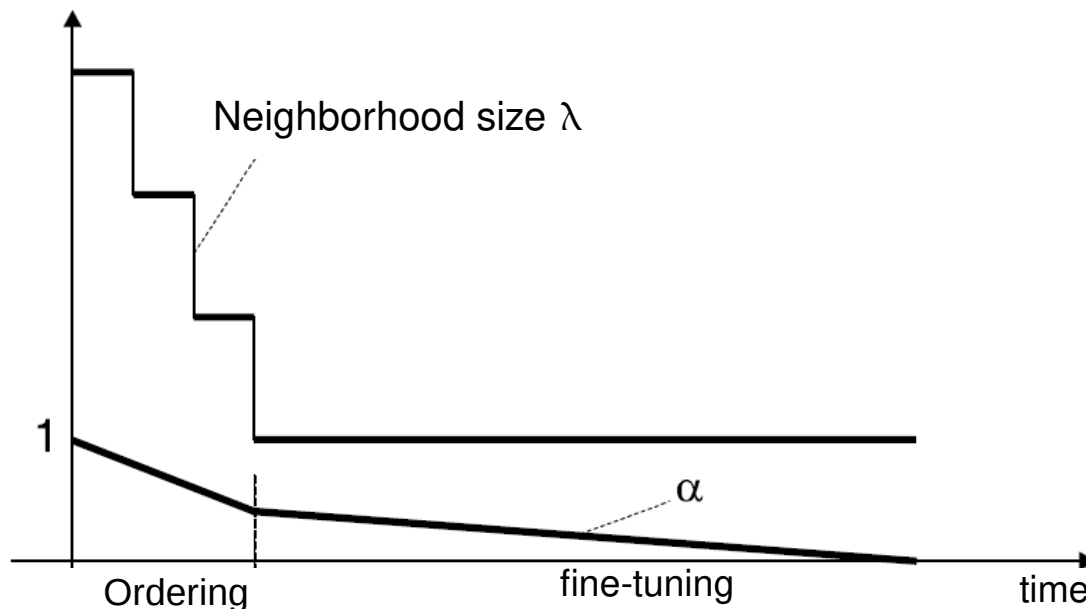$N_r(t_2)$

$N_r(t_3)$



- alternative: Gaussian neighborhood

# SOM algorithm

- randomly choose an input $x$
- find winner $i*$ for $x$
- adapt weights within the neighborhood

$$i^* = \arg\min_i \|x - w_i\|$$
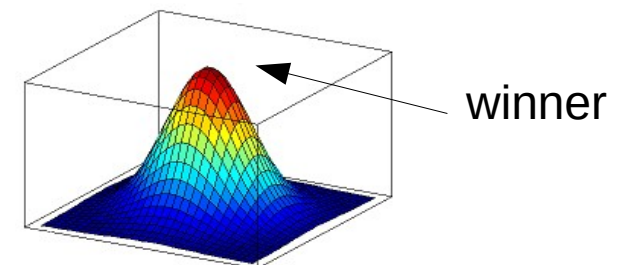
$$w_i(t+1) = w_i(t) + \alpha(t) h(i^*, i)[x(t) - w_i(t)]$$

- update SOM parameters (neighborhood, learning rate)
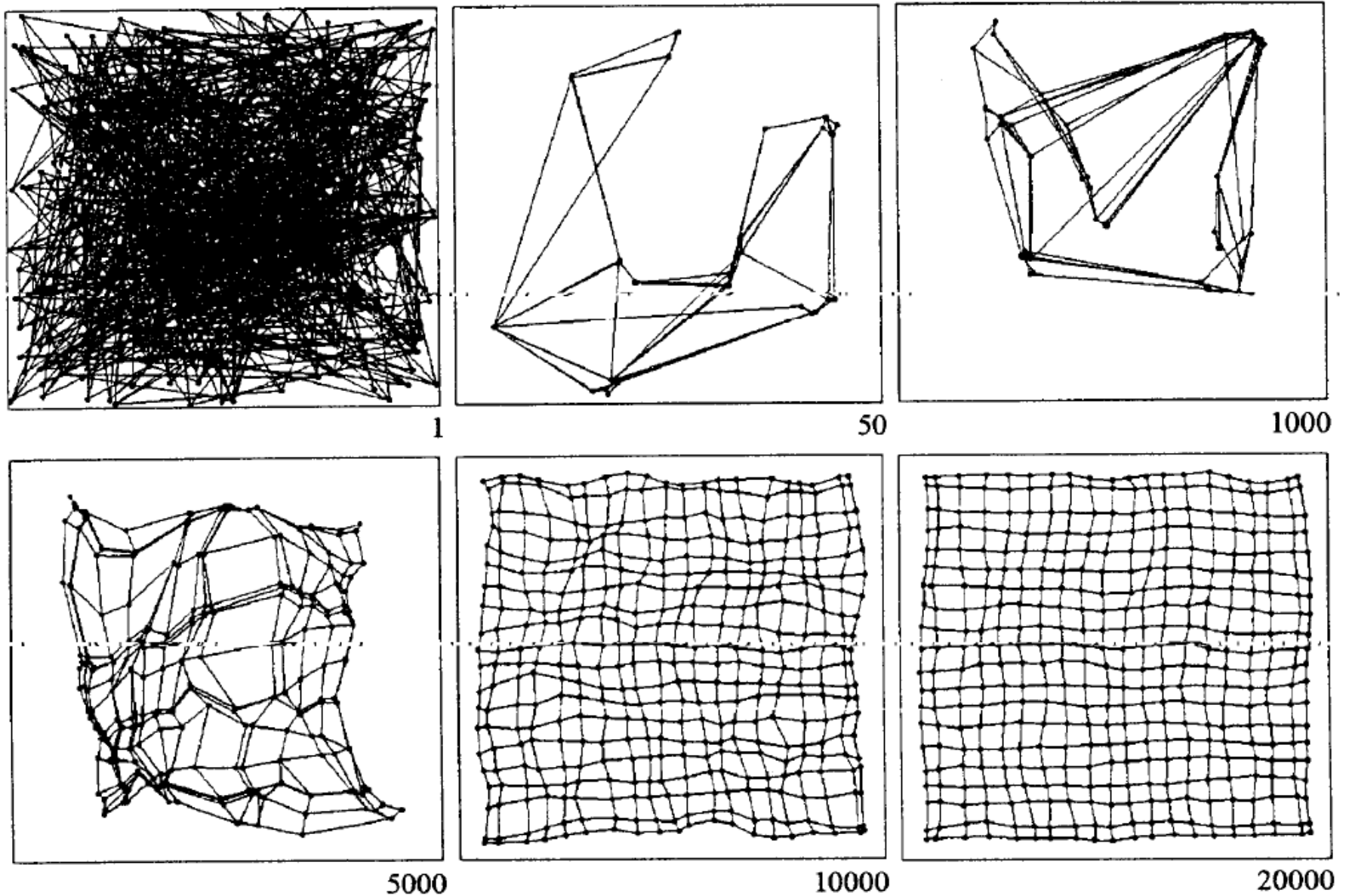- repeat until stopping criterion is met

Neighborhood size $\lambda$

Input-output mapping:

$X \rightarrow \{1,2,...,m\}$    or

$X \rightarrow Y$    $y = [y_1, y_2,...,y_m]$

where e.g. $y_i = \exp(-\|x - w_i\|)$

winner

1

Ordering     fine-tuning     time

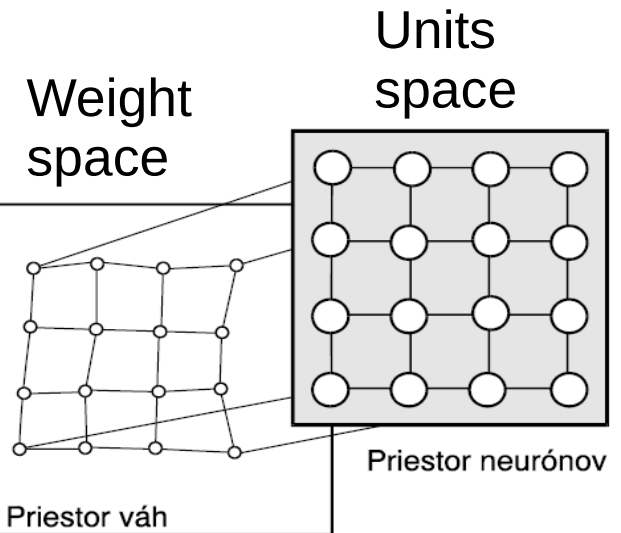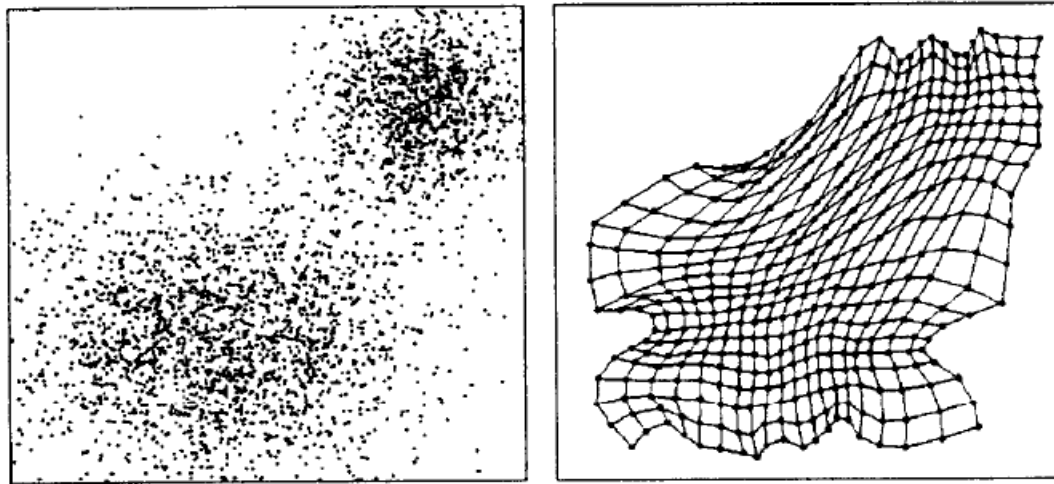$\alpha$

# Example: 2D inputs, 20x20 neurons



Weight vectors (knots) become topographically ordered at the end of training

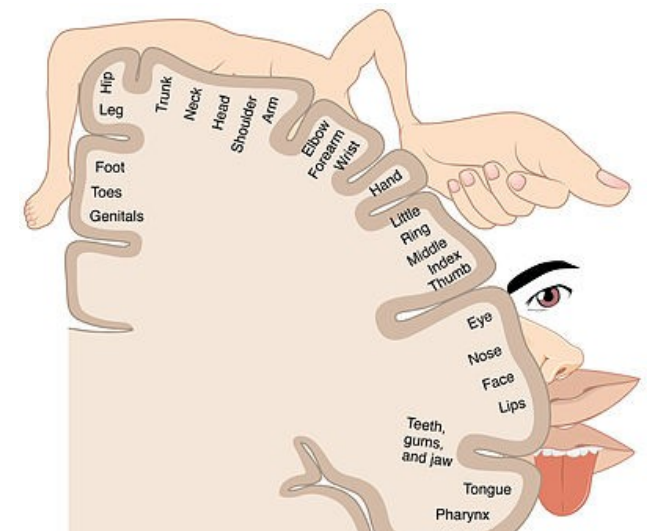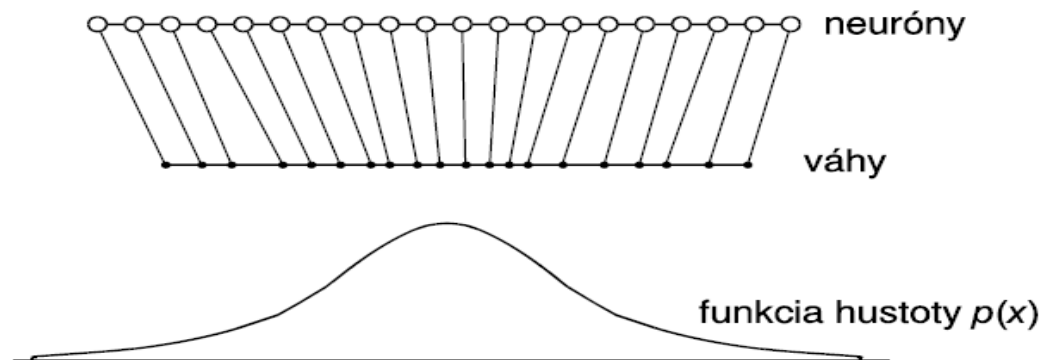# SOM approximates input data distribution

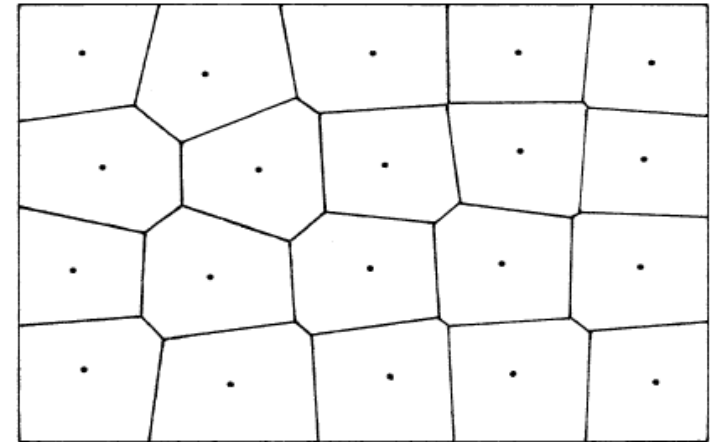- by allocating more neurons to denser areas

2D:



Weight space

Units space

Priestor neurónov

Priestor váh

1D:



neuróny

váhy

funkcia hustoty $p(x)$

# SOM performs two tasks simultaneously

**Input clustering**
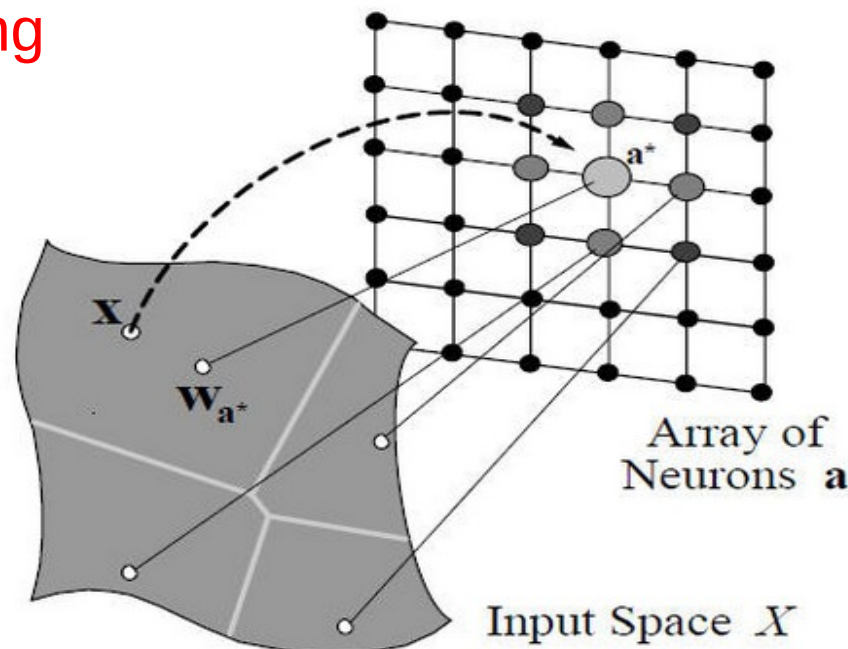
Voronoi compartments:

$$V_i = \{x \mid \|x - w_i\| < \|x - w_j\|, \ \forall j \neq i\}$$



Voronoi tessellation
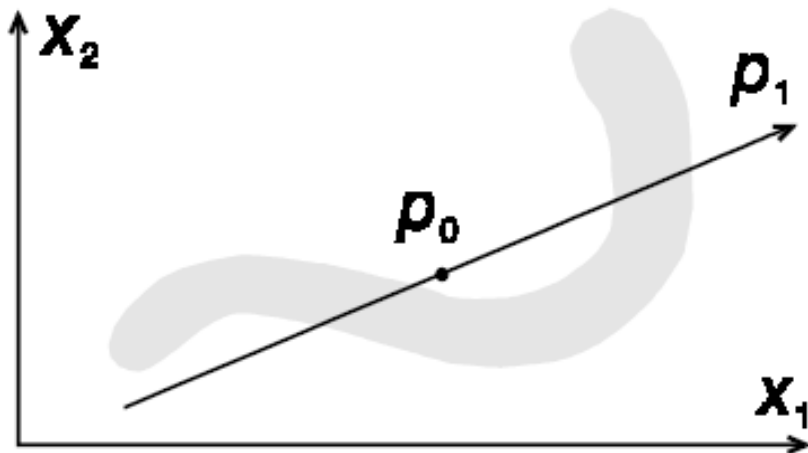of input space

**Topographic mapping**

i.e. similar stimuli
in input space are
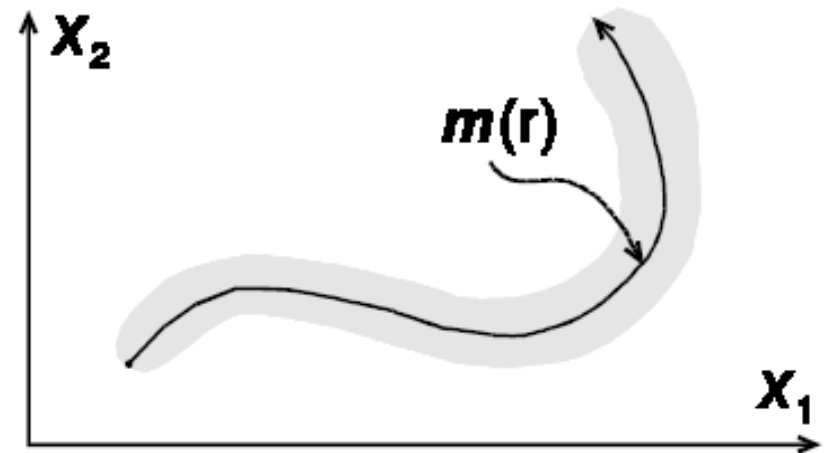mapped close
to each other in
the SOM

# Comparison of SOM to PCA

- feature extraction and mapping, difference in feature representation

PCA                                                    SOM



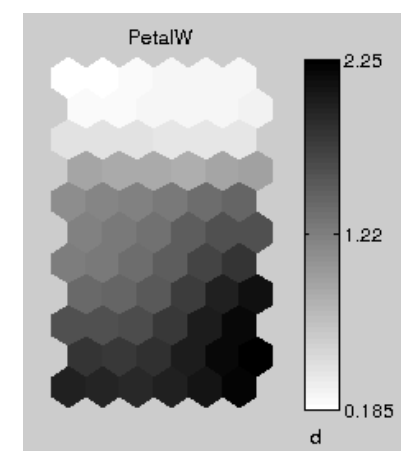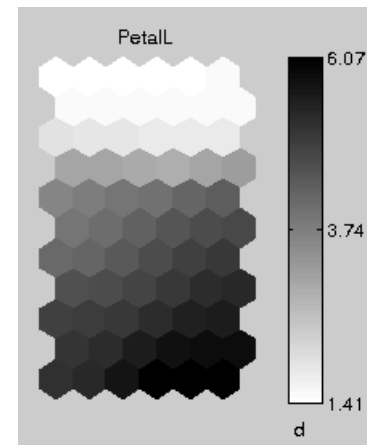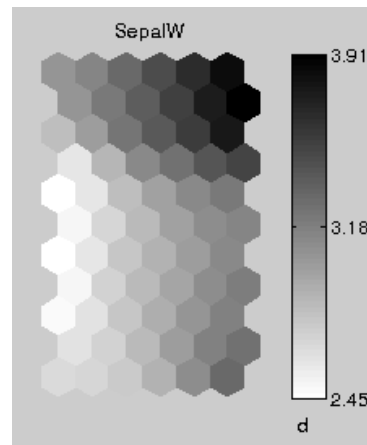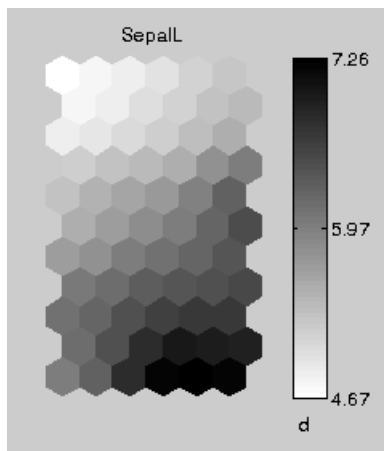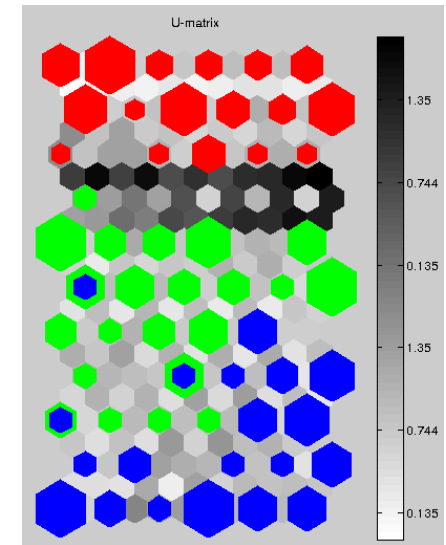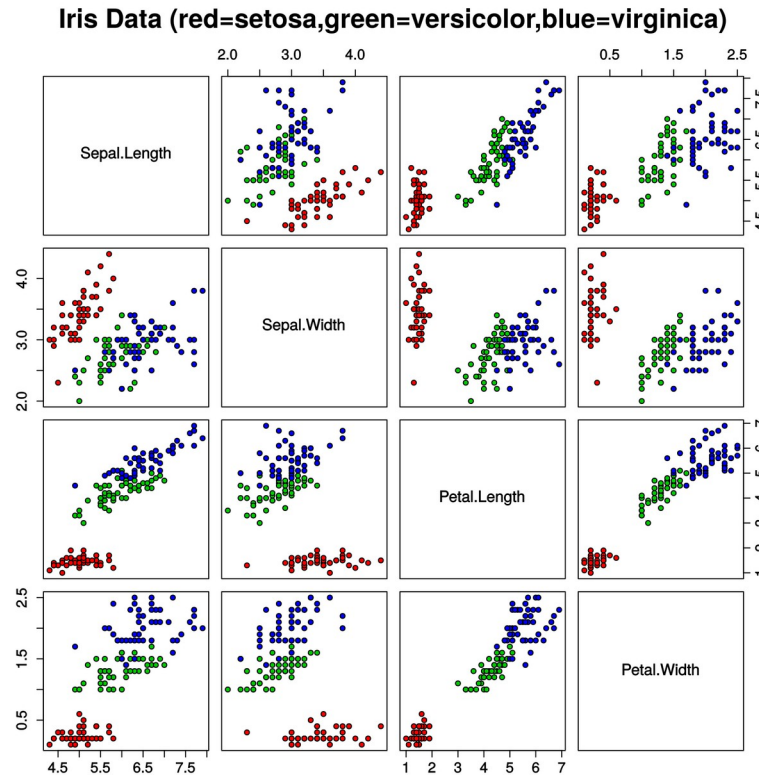(linear) principal components          (nonlinear) principal manifold
One unit represents 1 dimension       More units represent 1 dimension

SOM as a nonlinear generalization of PCA

# Application: Data visualization

- 7x10 SOM with hexa grid
- trained on 4D Iris data
- 3 classes
- Neuron labels generated
- according to votes
- Plots of component weights reveal an order

**Iris Data (red=setosa,green=versicolor,blue=virginica)**

# Summary

- Unsupervised methods extract features from data (in various ways)

- PCA – standard method for linear reduction of data dimensionality (preserves maximal variance)

- A single-layer neural net trained by Hebbian-like learning rule can implement PCA

- No need to calculate covariance matrix of input data

- SOM – a very popular algorithm for input clustering and topographic mapping

- useful for data visualization