

ARTIFICIAL INTELLIGENCE

LEI/3, LMA/3, MBE/1

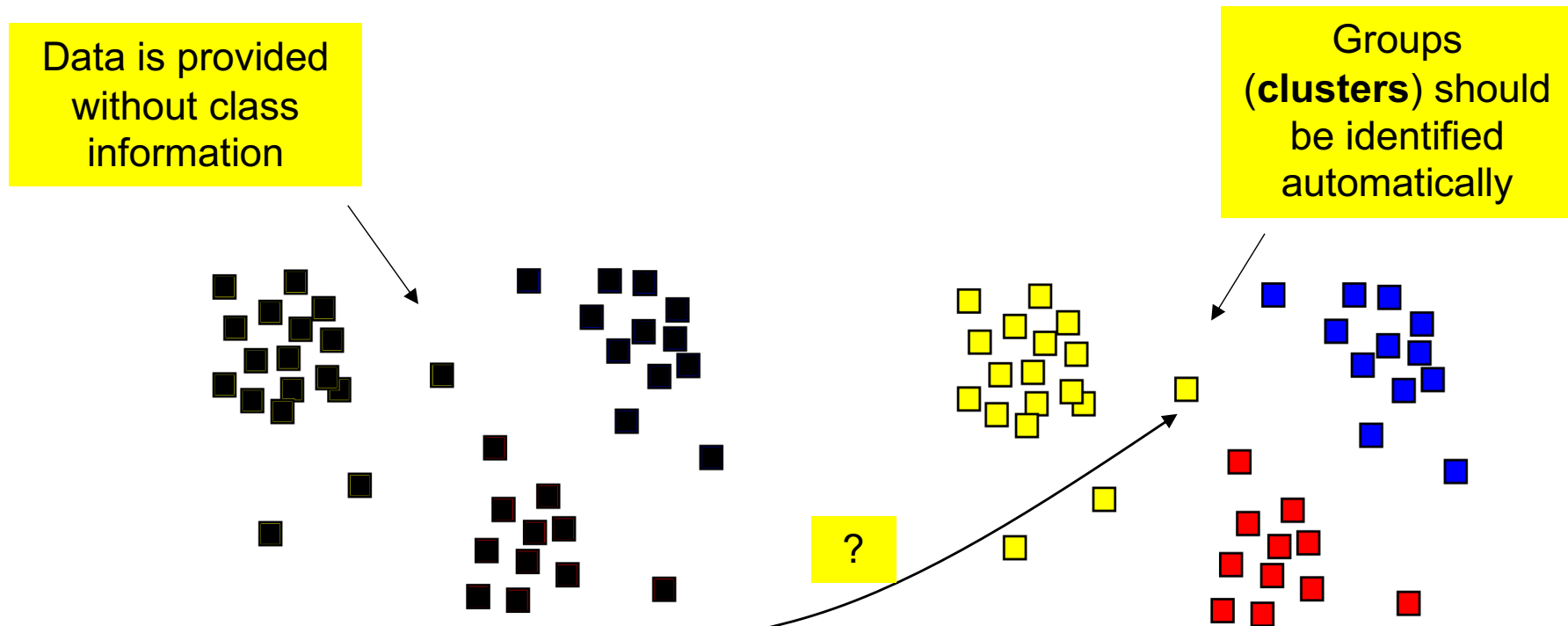
University of Beira Interior, Department of Informatics

Hugo Pedro Proença

hugomcp@di.ubi.pt, 2022/23

Unsupervised Learning

- This concept is associated to learning without a “supervisor”
 - It also known as self-organization, or cluster analysis
- The basic idea is that, instead of attempting to mimic the behavior of the supervisor, to identify commonalities in the data



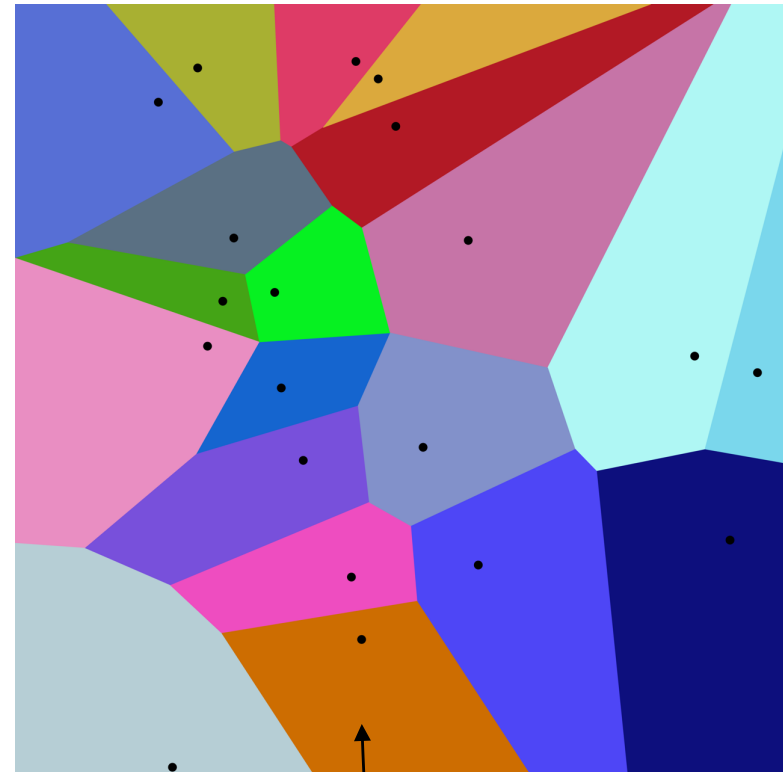
- The notion of "cluster" cannot be objectively defined, which justifies different clustering algorithms.

Unsupervised Learning

- There are different families of methods to perform clustering:
 - **Connectivity** models, in which models are built based on distance connectivity
 - Hierarchical clustering
 - **Centroid** models, that represent clusters by mean vectors (i.e., centroids)
 - K-means
 - **Distribution** models, where clusters are modelled according to statistical distributions
 - DBSCAN
 - **Neural** models, where networks implement a form of PCA that finds appropriate feature subspaces
 - Self-Organizing Map (SOM)
- **Clusters Evaluation**
 - **Internal** Evaluation, when the model is evaluated based on the data that was clustered itself
 - Davies-Bouldin index: $DB = \frac{1}{N} \sum_{i=1}^N \max_{j \neq i} \frac{\sigma_i + \sigma_j}{d(c_i, c_j)}$
where “c” represents one centroid, “σ” is the average distance of the elements in one cluster to its centroid and “N” is the number of clusters
 - **External** Evaluation, when the model is evaluated based on new data, typically with class labels
 - Purity: $P = \frac{1}{N} \sum_{i=1}^M \max_{d \in D} |m \cap d|$
where “M” represents the set of clusters, and “D” is the labeled data

K-Means

- It is the most used clustering algorithm, due to its effectiveness and easiness of implementation.
 - Aims to partition “**n**” observations into “**k**” clusters
 - Each observation belongs to the nearest cluster centroid, which is the prototype of the cluster.
 - This results in a partitioning of the data space into **Voronoi cells**.
 - A Voronoi diagram is a partitioning of a plane into regions based on distance to points in a subset of the plane.
 - These points (a.k.a. prototypes) determine the shape of the corresponding Voronoi cell.
 - For each prototype there is a corresponding region consisting of all points closer to that seed than to any other. These regions are called Voronoi cells.



Positions in each cell (color) are closest to the corresponding centroid than to any other

K-Means

- For K-Means, the value of “K” must be given beforehand
 - There are different heuristics to automatically find the optimal value of “K”, but depend on the specific problem considered
- Having a data set \mathbf{X} : $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$

1. Initialize (randomly) “K” centroids $\boldsymbol{\mu}$: $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k\}$
2. While (\neg stopping_criterion($\boldsymbol{\mu}, \mathbf{X}$))

1. For every \mathbf{x}_i :

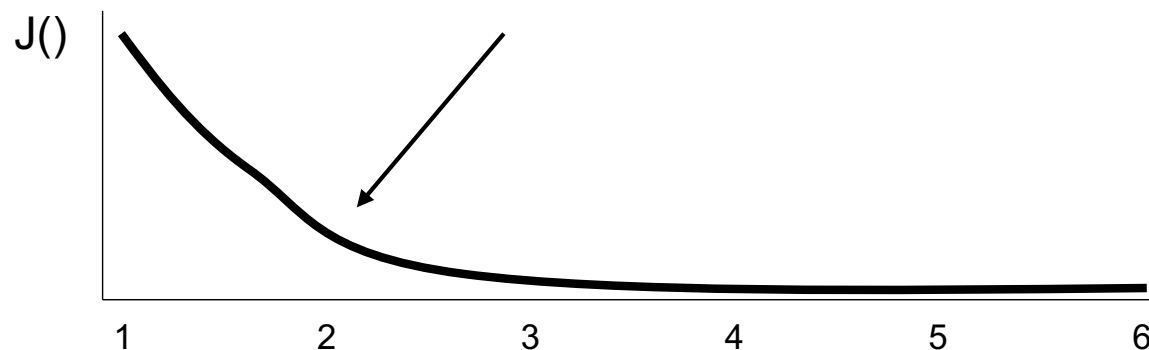
$$\mathbf{c}_i = \arg \min_j d(\mathbf{x}_i, \boldsymbol{\mu}_j) \quad // \text{cluster assignment}$$

2. For every $\boldsymbol{\mu}_i$:

$$\boldsymbol{\mu}_i = \frac{1}{n_i} \sum_j \mathbf{x}_j \mid \mathbf{x}_j \text{ assigned to } \mathbf{c}_i \quad // \text{centroid update}$$

K-Means

- **Stopping criteria.** There are a number of different possibilities
 - **Simplistic:** Predefine a **number of iterations**
 - Might be “*too many*”, or “*too few*”, depending of the complexity of the feature space
 - **Elaborate 1:** Evaluate clusters **stationarity** and stop when the changes in clusters positions between consecutive iterations is less than a small threshold.
 - **Elaborate 2:** Evaluate samples **assignments** and stop when no samples (or a very small number) of samples changes its centroid between consecutive iterations.
- **Choose the value of “K”**
 - **Elbow method.**
 - Define a cost function $J()$ and repeat the clustering procedure for a growing number of clusters. Define “K” as the value where the curvature of $J()$ is maximal



K-Means

• Distance Functions

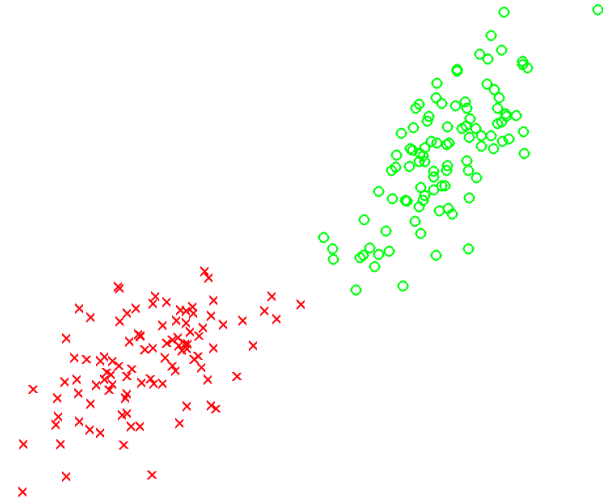
- Different functions can be used, as long as they met the properties of being a “metric”
- A metric on a set X is a function $d : X \times X \rightarrow [0, \infty)$, where for all $x, y, z \in X$, the following conditions are satisfied:
 - $d(\mathbf{x}, \mathbf{y}) \geq 0$ // non-negativity or separation axiom
 - $d(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} = \mathbf{y}$ // identity of indiscernibles
 - $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ // symmetry
 - $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ // triangle inequality

• Examples:

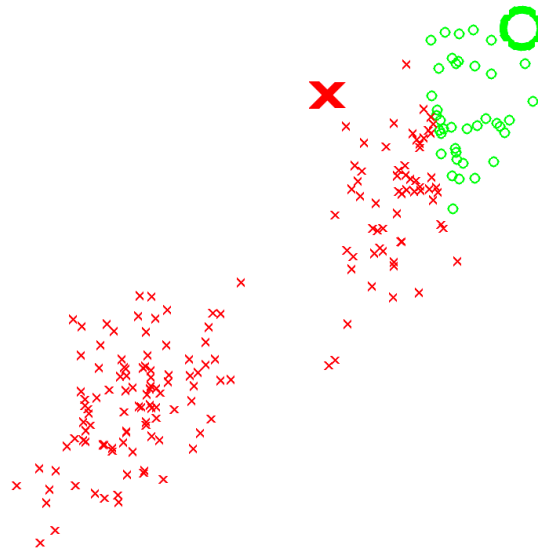
- Euclidean distance: $d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum (x_i - y_i)^2}$
- Manhattan distance: $d(\mathbf{x}, \mathbf{y}) = \sum (|x_i - y_i|)$
- Chebyshev distance: $d(\mathbf{x}, \mathbf{y}) = \max |x_i - y_i|$

K-Means: Example

- Consider the following synthetic dataset:

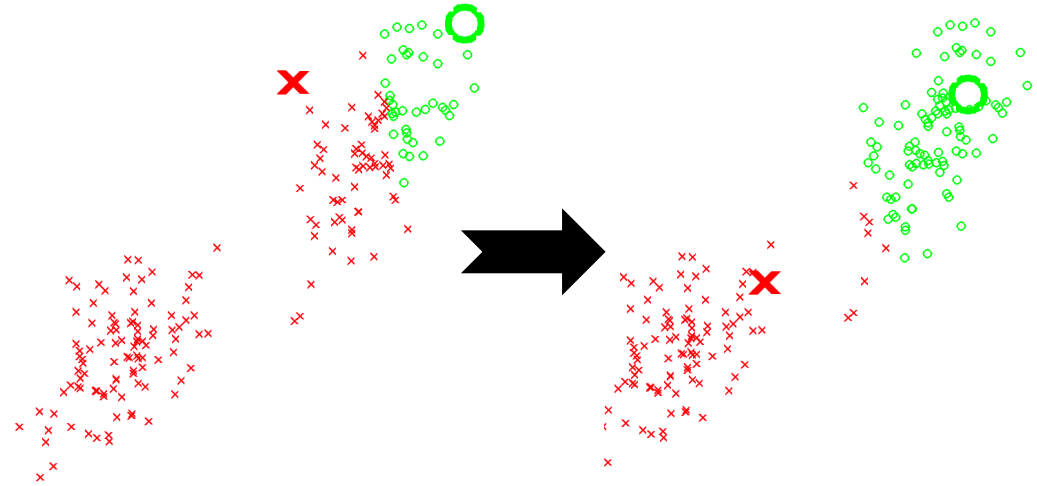


- Random initialization of 2 clusters:

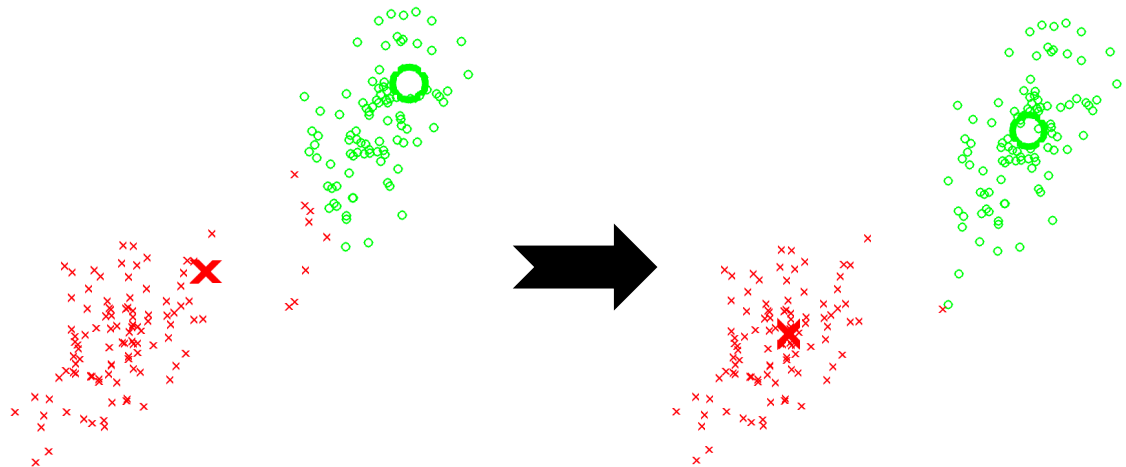


K-Means: Example

- K-Means: Iteration 1

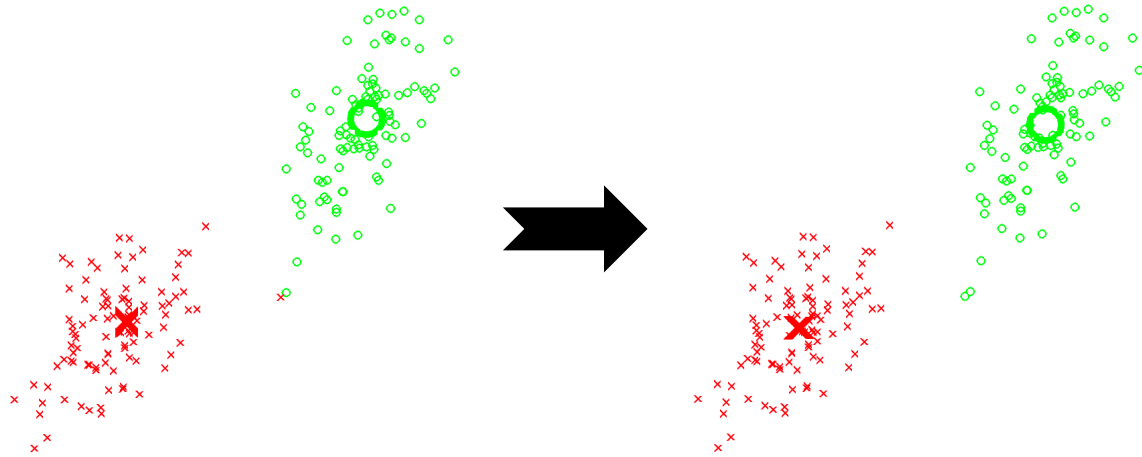


- K-Means: Iteration 2

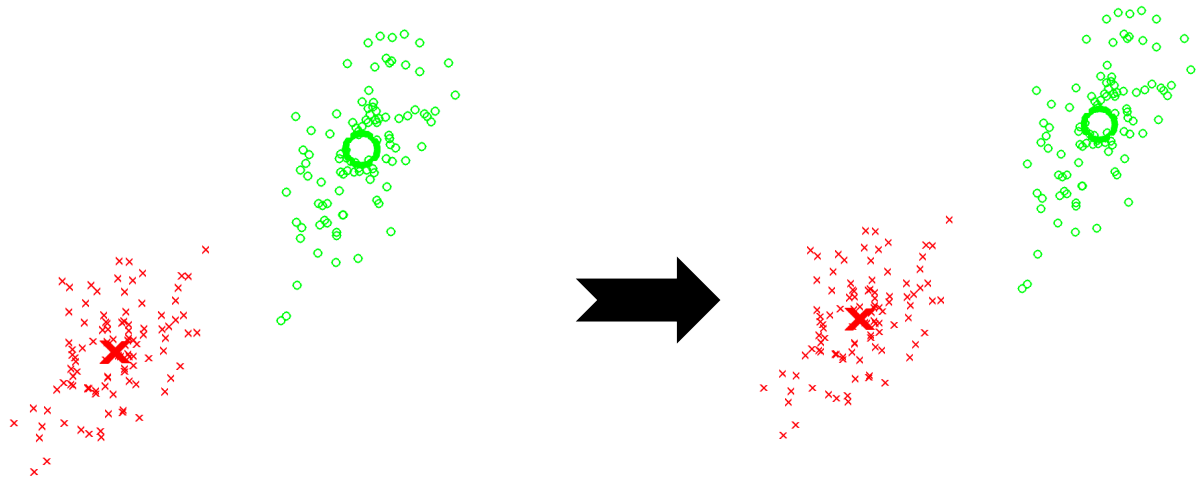


K-Means: Example

- K-Means: Iteration 3

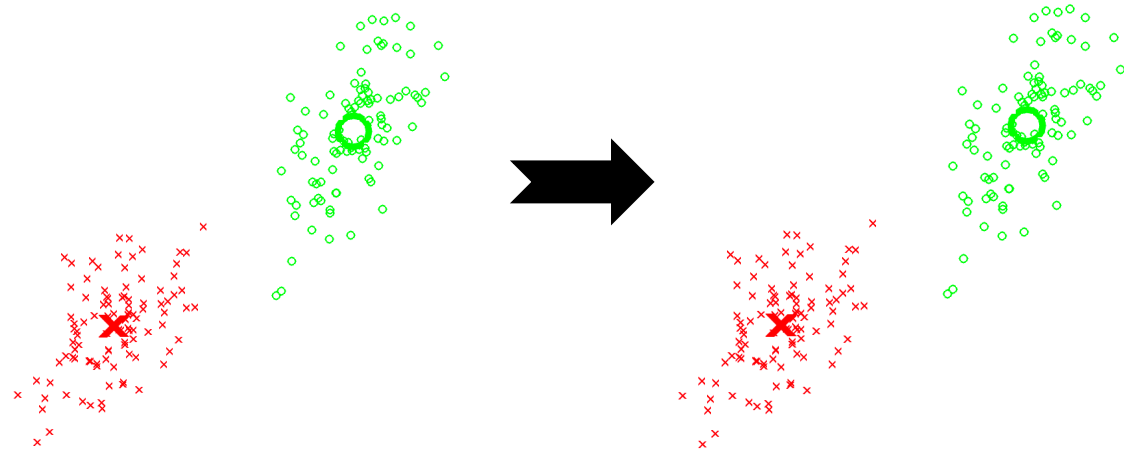


- K-Means: Iteration 4

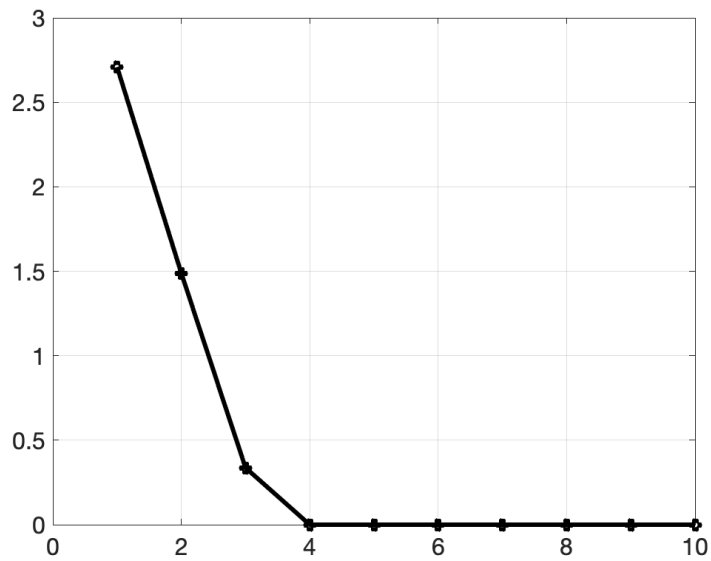


K-Means: Example

- K-Means: Iteration 5

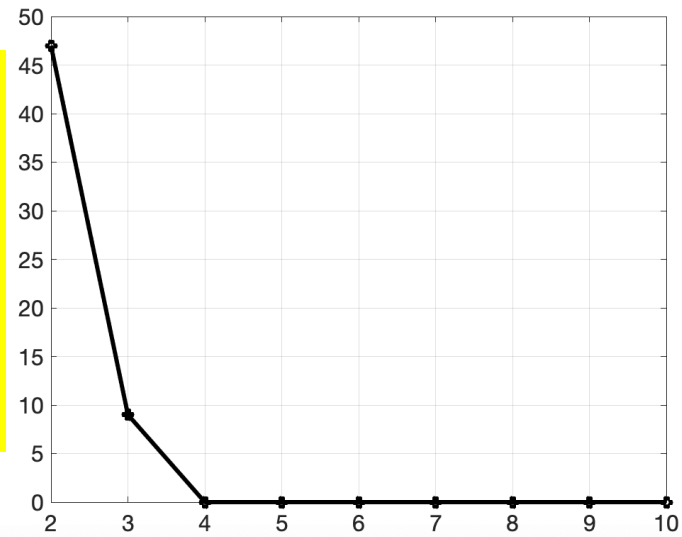


Sum of changes in centroid positions



Iteration

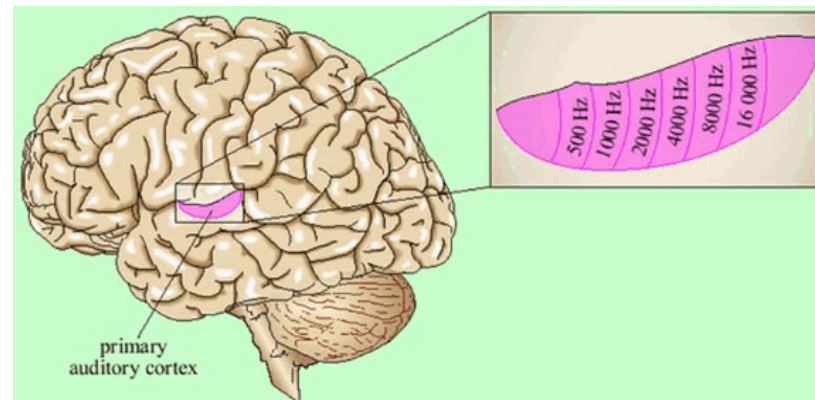
Number of instances that changed assignment



Iteration

Unsupervised Learning: SOMs

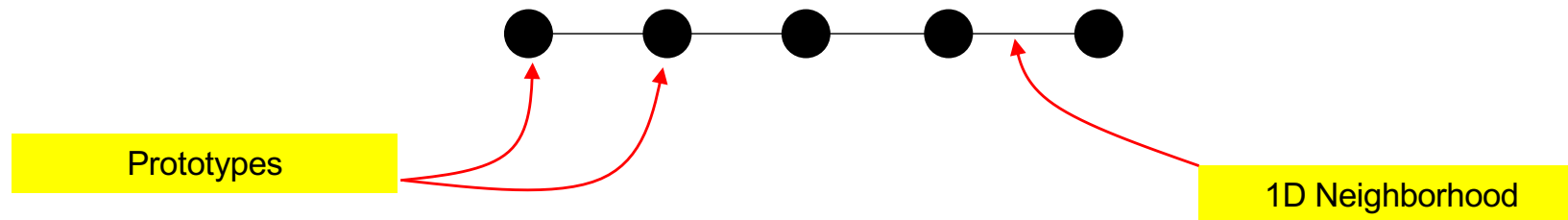
- The **Self-Organizing Map (SOM)** algorithm was one of the earliest neural network models (proposed by Kohonen in 1984), as a way to explain the spatial organization of the brain's functions, as observed especially in the cerebral cortex.
 - As an example, sound signals of different frequencies are mapped to the primary auditory cortex in which **neighboring neurons respond to similar frequencies.**



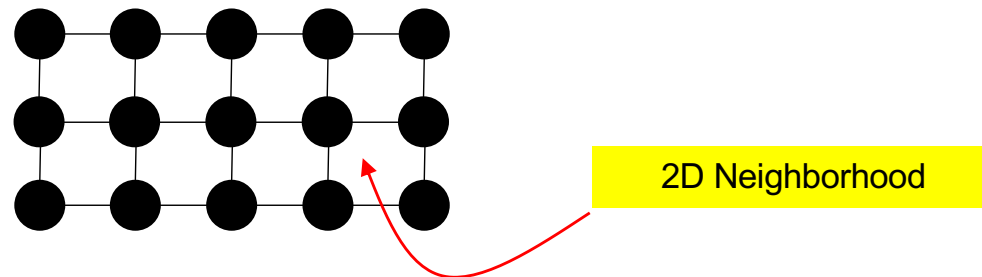
- Following the ideas of spatially ordered line detectors (von der Malsburg, 1973) and the neural fields (Amari, 1980), the key points in SOMs were:
 - introduce a model composed of two interacting subsystems of different natures, where **the key is a competitive neural network** that implements the winner-take-all function,
 - Design a synaptic **plasticity model** for the neurons in learning. In practice terms, the **learning is restricted spatially** to the local neighborhood of the most active neurons.

Unsupervised Learning: SOMs

- Computationally, a **Self-Organizing Map (SOM)** is an artificial neural network, trained according to the **unsupervised learning paradigm**
 - The topology of the network defines typically in a **low dimensional hyperspace**:
 - 1D Lines: e.g., 5 cells



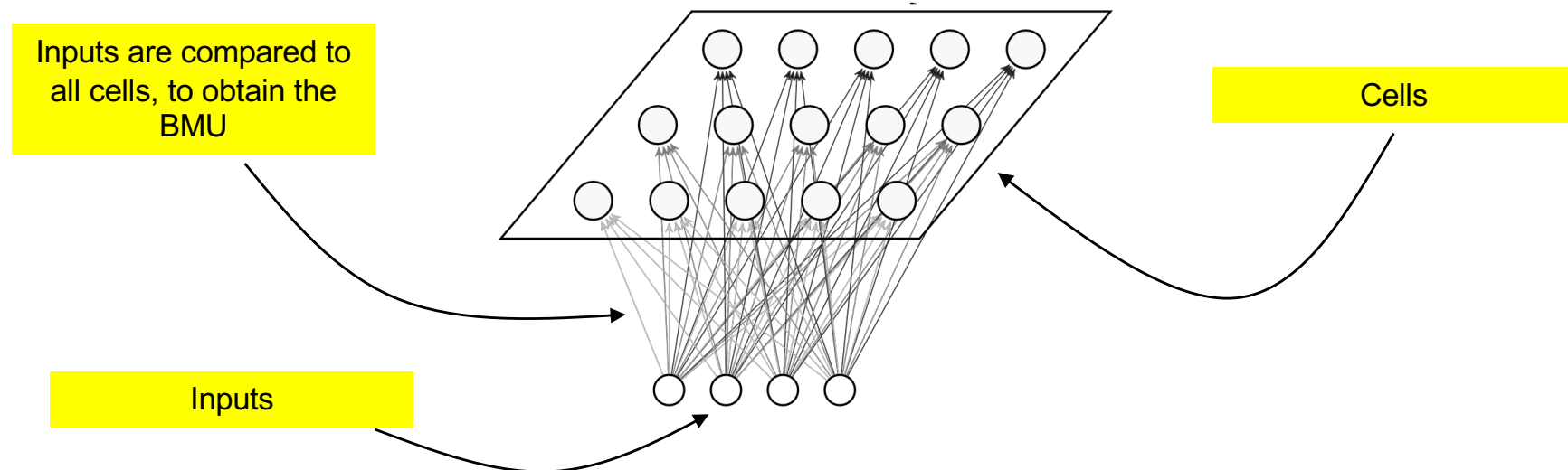
- 2D Grids: e.g., 3 x 5 cells



- 3D Spaces: ...

Unsupervised Learning: SOMs

- The cells of SOMs form the set of prototypes (clusters), and are autonomously inferred, according to a competitive learning paradigm
 - Having a set of input instances (learning set), each one lying in a (typically) high hyperspace, the idea is:
 - find the SOM prototype that best represents each instance and “*assign the prototype to it*”.
 - Called the **BMU** (Best Matching Unit)
 - Move (organize) the cells in the SOM, such that the BMU and its neighbors adjust their configurations (weights) towards the input instance
- Conceptually, there are only two entities in SOMs:



Unsupervised Learning: SOMs

• SOM Learning Algorithm

1. Initialize the weights of the cells (c_i) randomly
2. While ! Stop_criterion()
 1. For each input instance $\mathbf{x}^j \in \mathcal{R}^n$
 1. For each cell c_i , get the distance between \mathbf{x}^j and c_i

$$d(\mathbf{x}^j, c_i)$$

e.g., Euclidean
Distance

2. Find the BMU

$$b = \arg \min_i d(\mathbf{x}^j, c_i)$$

3. Update the weights of the BMU cell:

$$\boldsymbol{\theta} = \boldsymbol{\theta} + \gamma (\boldsymbol{\theta} - \mathbf{x})$$

Learning Rate

5. Update the weights of the neighbour cells (c^k):

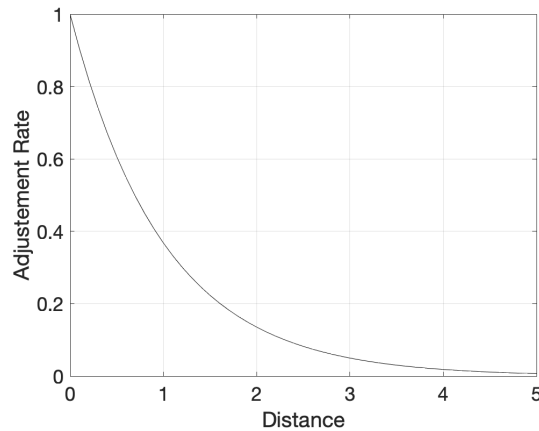
$$\boldsymbol{\theta} = \boldsymbol{\theta} + \gamma \exp\left(-\frac{d(c_k, c_b)}{\sigma}\right) (\boldsymbol{\theta} - \mathbf{x})$$

Neighborhood Constant

Unsupervised Learning: SOMs

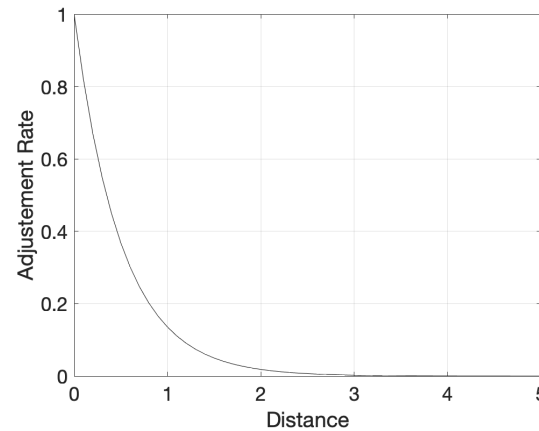
- The shape of the neighborhood function determines how much (and how many) cells adjacent to the BMU are updated.

Many cells updated

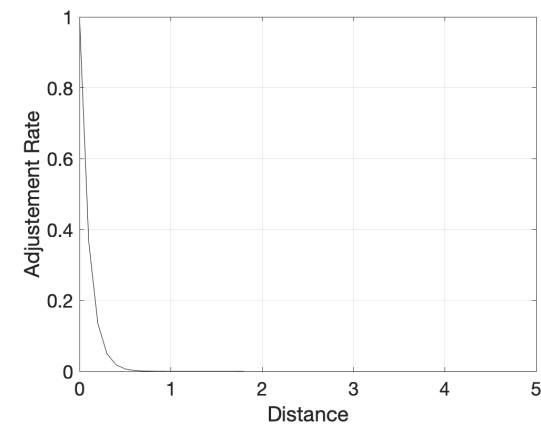


$\sigma = 1$

Few cells updated



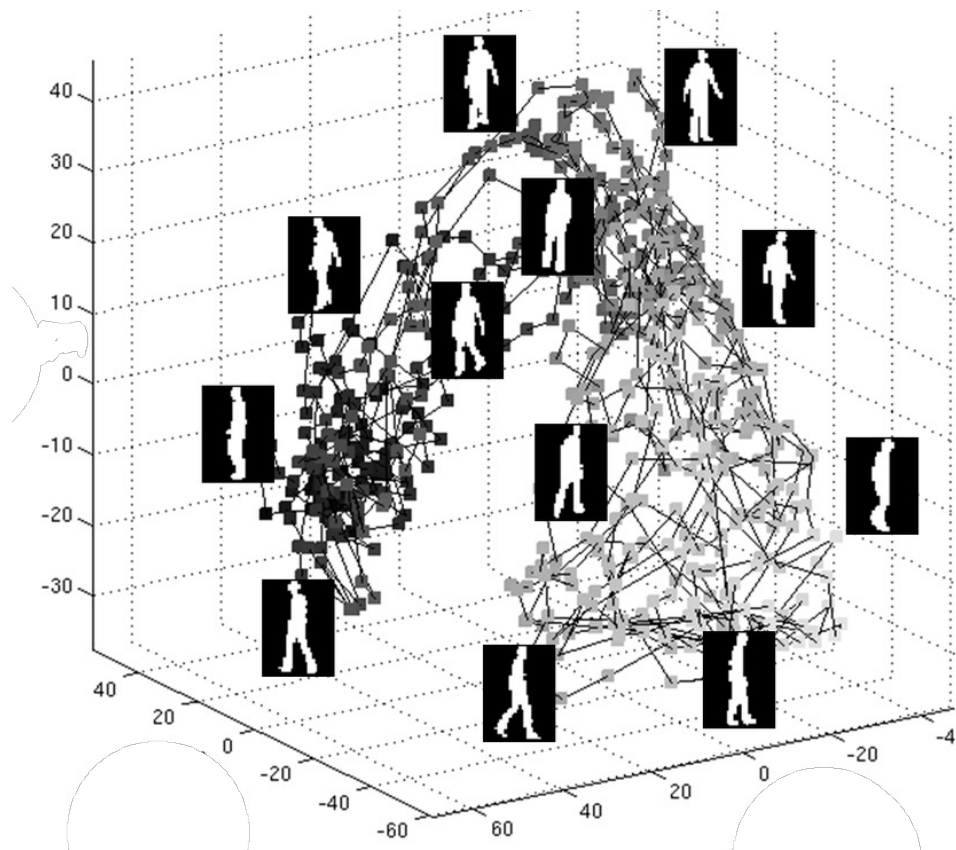
$\sigma = 0.5$



$\sigma = 0.1$

Unsupervised Learning: SOMs

- **More than just clustering, SOMs** build a low dimensional manifold, where high dimensional input samples are projected.
 - The relevant advantage of manifolds is the topological space property, i.e., neighbor elements “look-alike”, in opposition to elements that are farther in the destiny-space.



Machine Learning: SOM Exercise

- Consider the “[AR.tar](#)” dataset, available at the course web page.
 - It contains 3.315 [48 x 64] face images
 - We will use it to distinguish between “Male” and “Female” genders
- Implement a “Python” that:
 - Loads the set of images
 - Divide the set into two disjoint parts: “learning” and “test”
 - 90% for learning, 10% for test, randmly chosen
 - Builds a **SOM** manifold
- Check which SOM topology is more appropriate for distinguishing between...
 - Identities
 - Gender

