

# Clustering: Models and Algorithms

Shikui Tu

2019-02-28

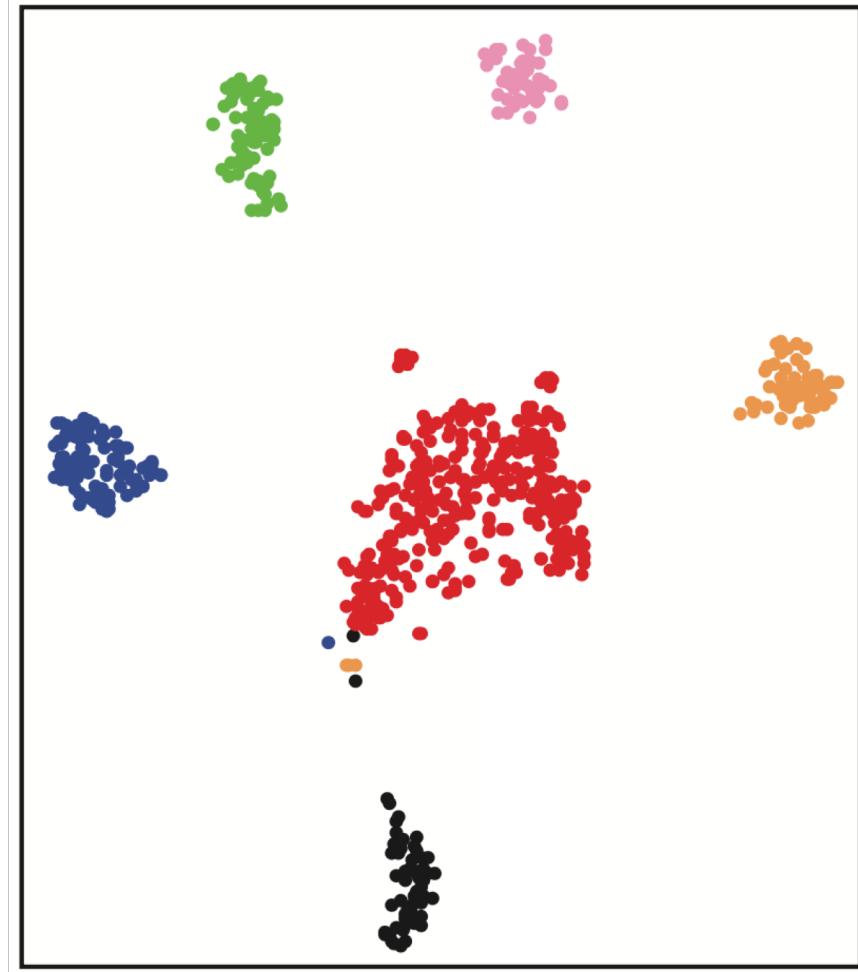
# Outline

- Clustering
  - K-mean clustering, hierarchical clustering
- Adaptive learning (online learning)
  - CL, FSCL, RPCL
- Gaussian Mixture Models (GMM)
- Expectation-Maximization (EM) for maximum likelihood

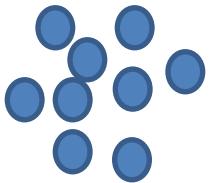
# What is clustering?

例子：不同类型的癌细胞会各自聚在一起

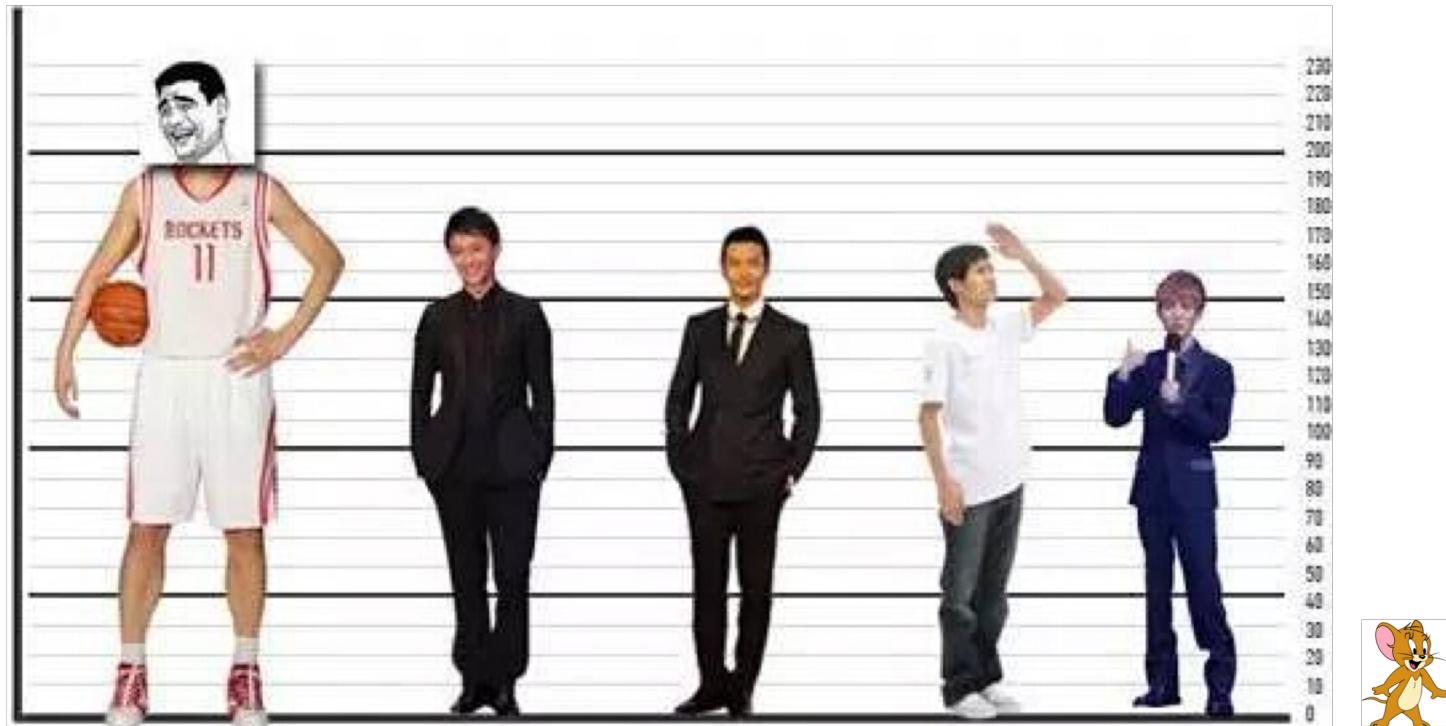
物以类聚



# How to represent a cluster



- 例如：将每个人的身高记下来



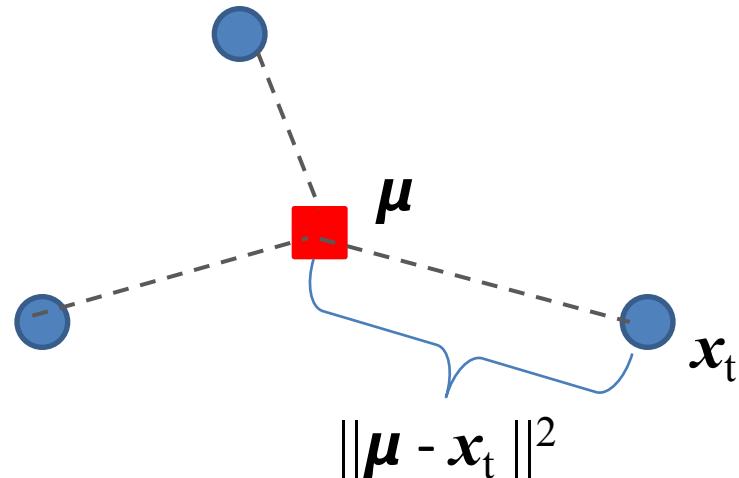
但是，如果只能记一个身高数值...

平均值

总误差最小

# How to define error?

Square distance:



$$\|\mu - x_1\|^2 + \|\mu - x_2\|^2 + \|\mu - x_3\|^2$$

可以证明：当  $\mu$  是所有数据点的均值时，平方距离和最小

# Matrix derivatives

$$\left[ \frac{\partial \mathbf{x}}{\partial y} \right]_i = \frac{\partial x_i}{\partial y} \quad \left[ \frac{\partial x}{\partial \mathbf{y}} \right]_i = \frac{\partial x}{\partial y_i} \quad \left[ \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right]_{ij} = \frac{\partial x_i}{\partial y_j}$$

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \quad (69)$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T \quad (70)$$

$$\frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{b}}{\partial \mathbf{X}} = \mathbf{b} \mathbf{a}^T \quad (71)$$

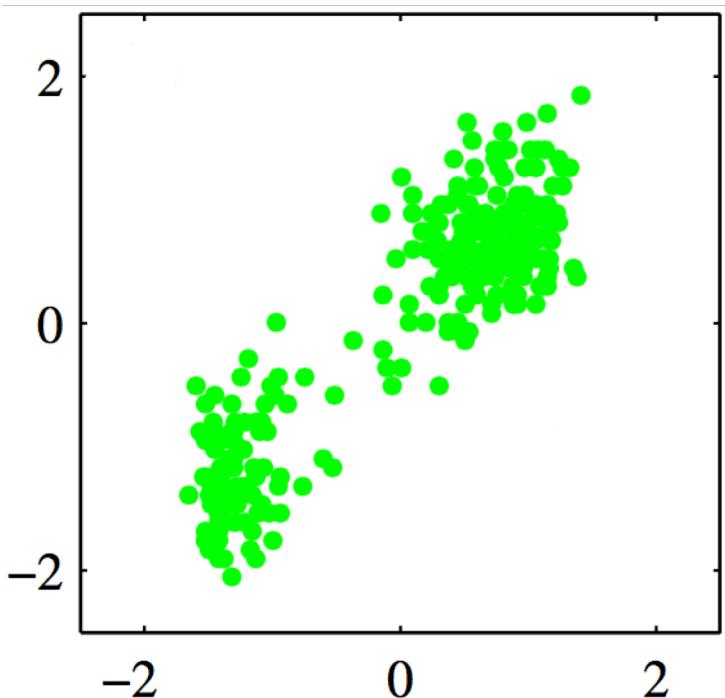
$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{X}} = \frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{a}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{a}^T \quad (72)$$

$$\frac{\partial \det(\mathbf{X})}{\partial \mathbf{X}} = \det(\mathbf{X}) (\mathbf{X}^{-1})^T \quad (49)$$

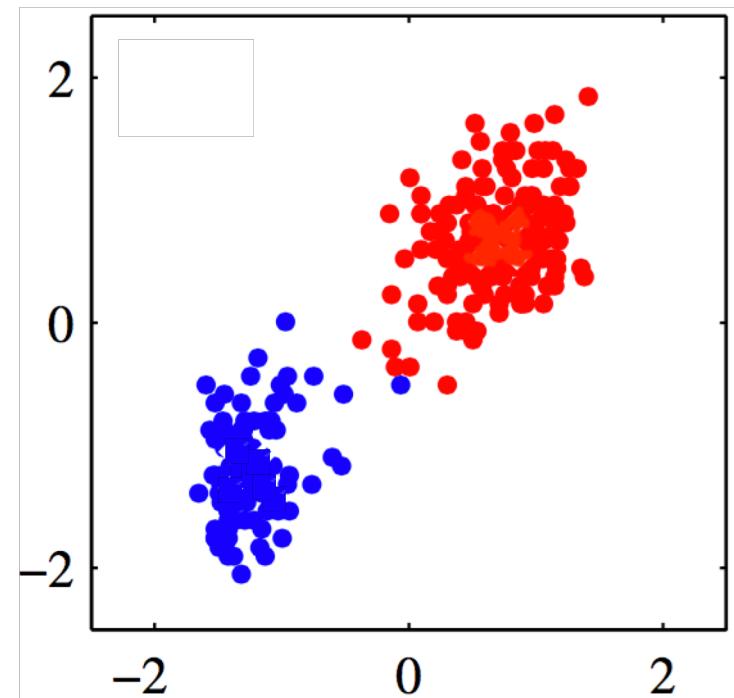
$$\frac{\partial \mathbf{Y}^{-1}}{\partial x} = -\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial x} \mathbf{Y}^{-1} \quad (59)$$

# Clustering the data

We have the following data:



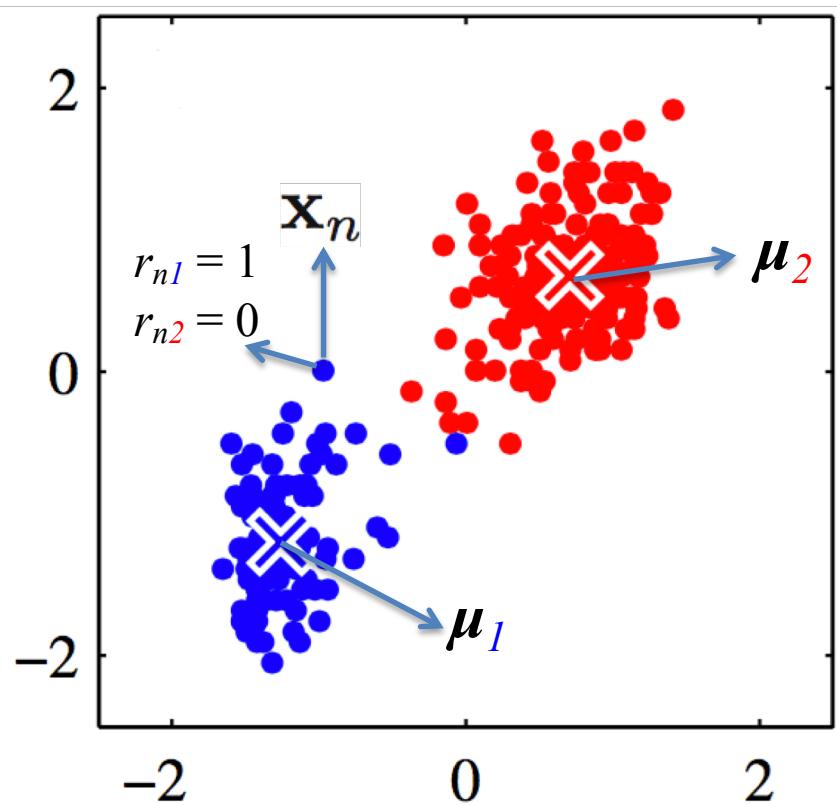
We want to cluster the data into two clusters (**red** and **blue**)



How?

# Minimize the sum of square distances $J$

minimize 
$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$



$r_{nk} = 1$  if and only if data point  $\mathbf{x}_n$  is assigned to cluster  $k$ ;  
otherwise  $r_{nk} = 0$ .

$k = 1, 2$ ;  $K = 2$  clusters

$n = 1, \dots, N$ ;  
 $N$ : the total number of points.

We need to calculate  $\{r_{nk}\}$  and  $\{\boldsymbol{\mu}_k\}$ .

If we know  $r_{n1}$ ,  $r_{n2}$  for all  $n=1,\dots,N$

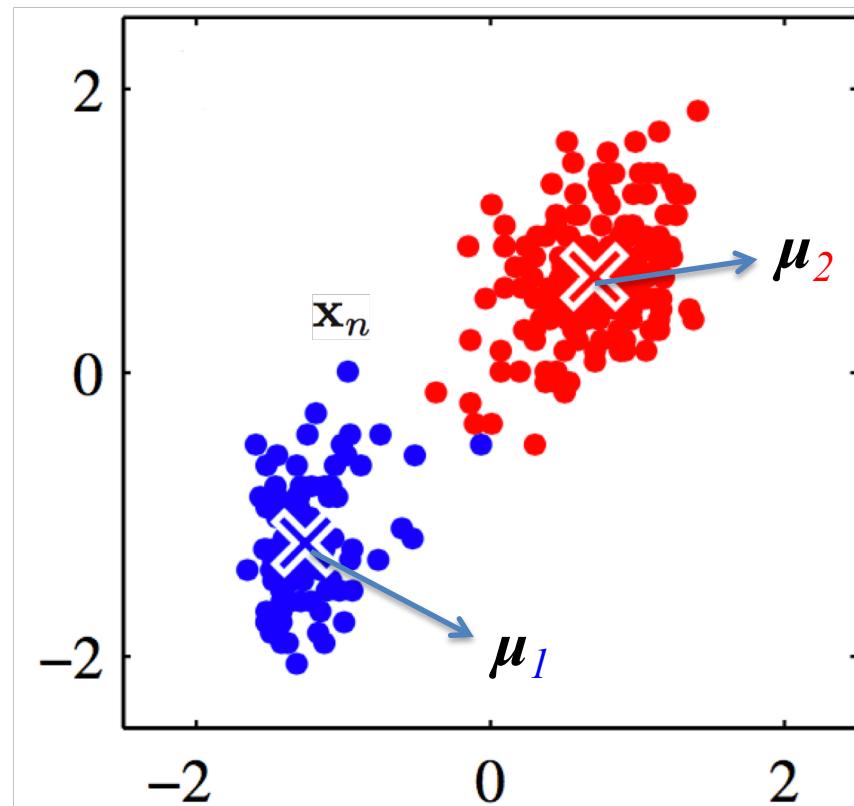
Since the points have been assigned to cluster 1 or cluster 2, we calculate

$\mu_1$  = mean of the points in cluster 1

$\mu_2$  = mean of the points in cluster 2

Or formally

$$\mu_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$



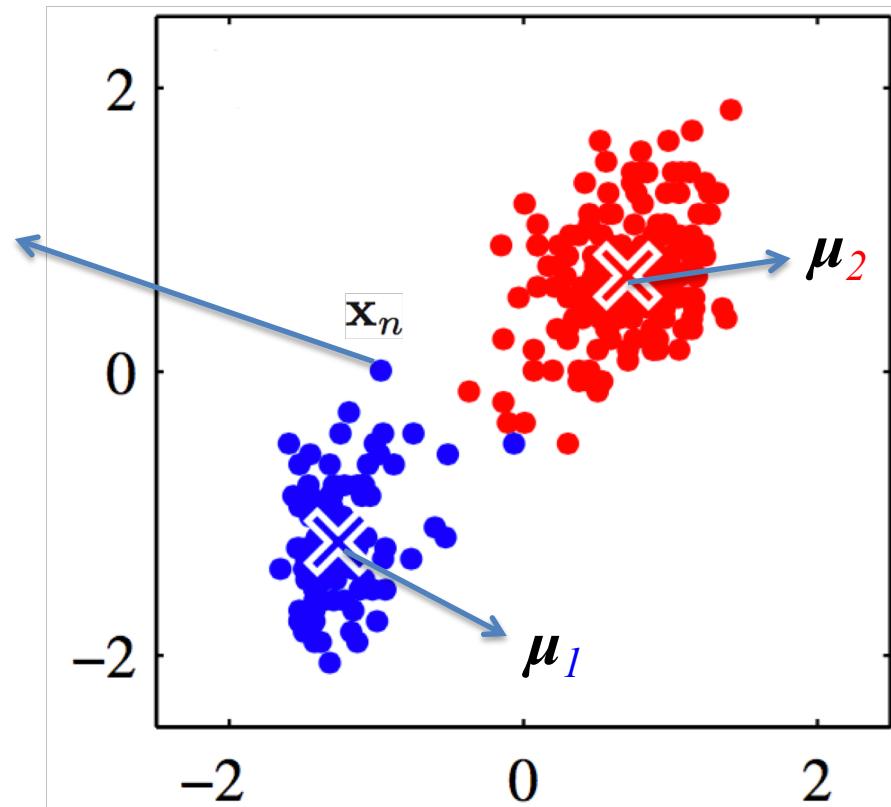
We call it the **M Step**.

# If we know $\mu_1, \mu_2$

We should assign point  $\mathbf{x}_n$  to cluster 1, because

$$\|\mathbf{x}_n - \mu_1\|^2 < \|\mathbf{x}_n - \mu_2\|^2$$

Then,  $r_{n1} = 1$   
 $r_{n2} = 0$

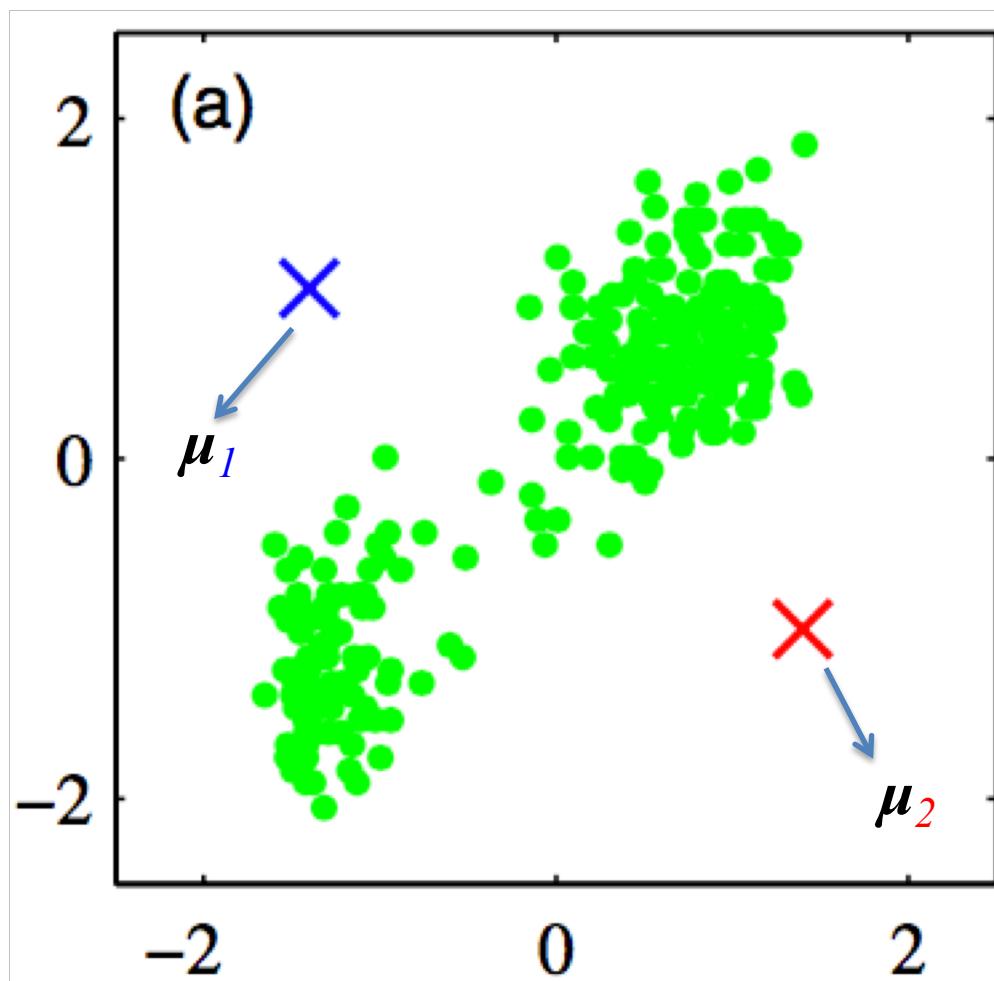


Or formally

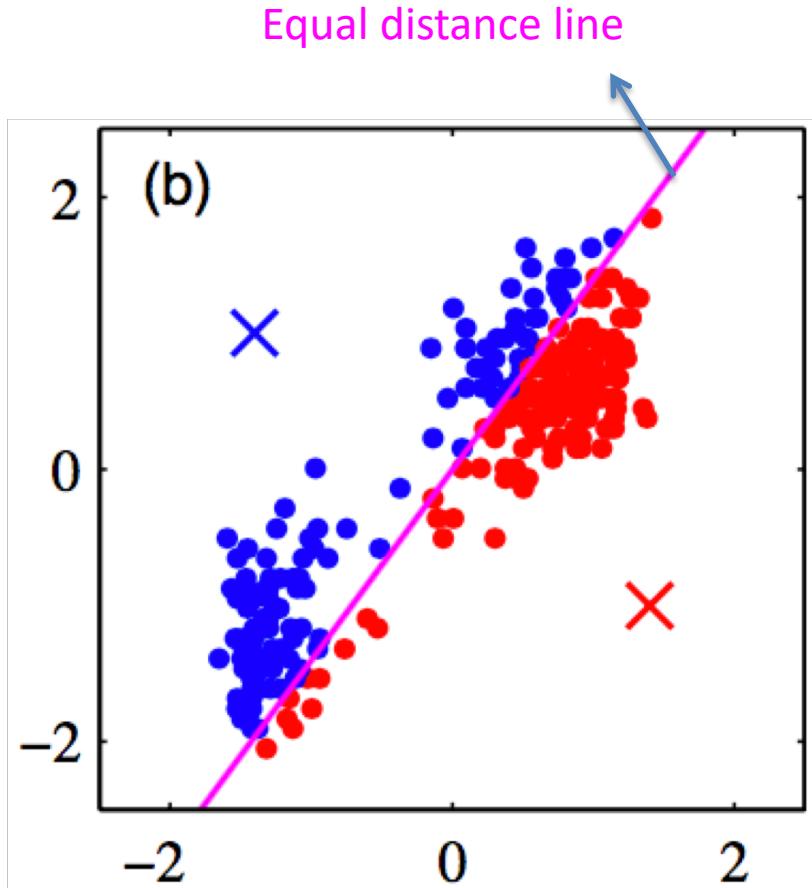
$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$

We call it the **E Step**

# Initialization



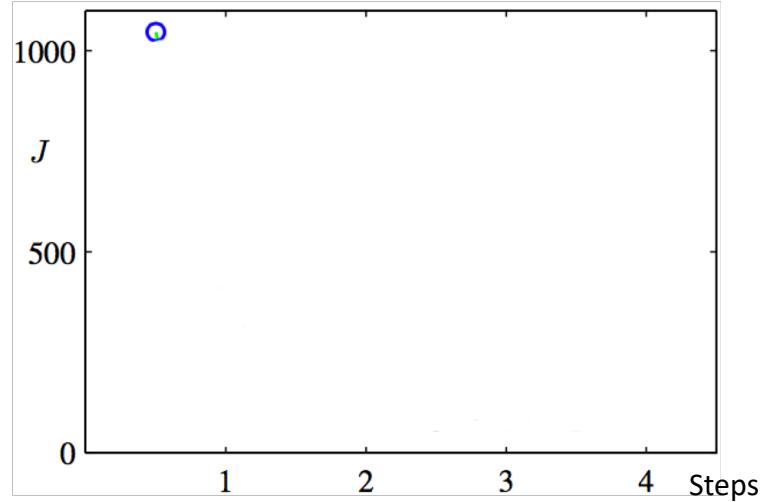
Given  $\mu_1, \mu_2$ , calculate  $r_{n1}, r_{n2}$  for all  
 $n=1, \dots, N$



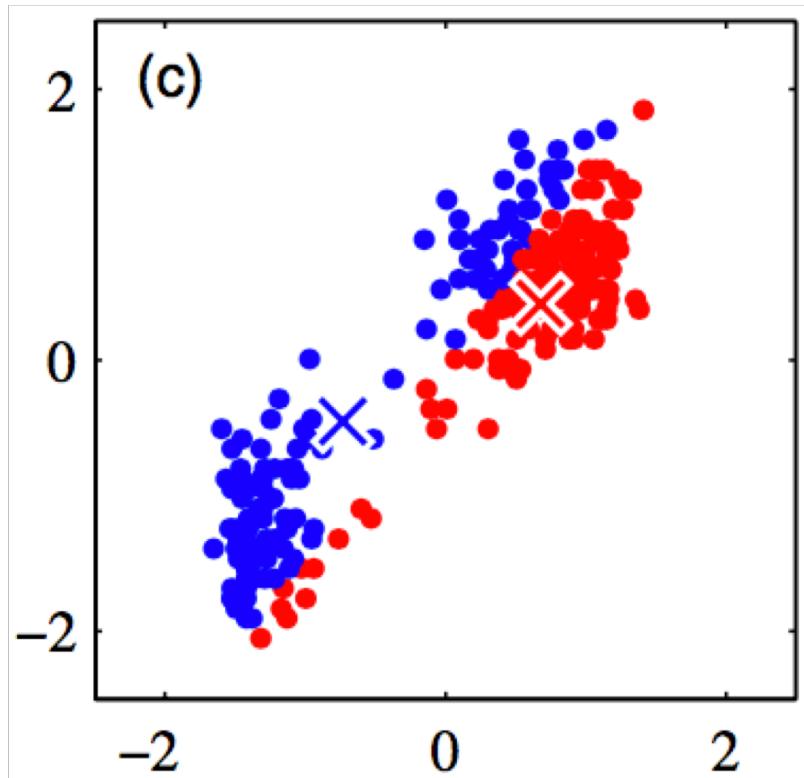
E Step

Assign the points to the nearest cluster:

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$



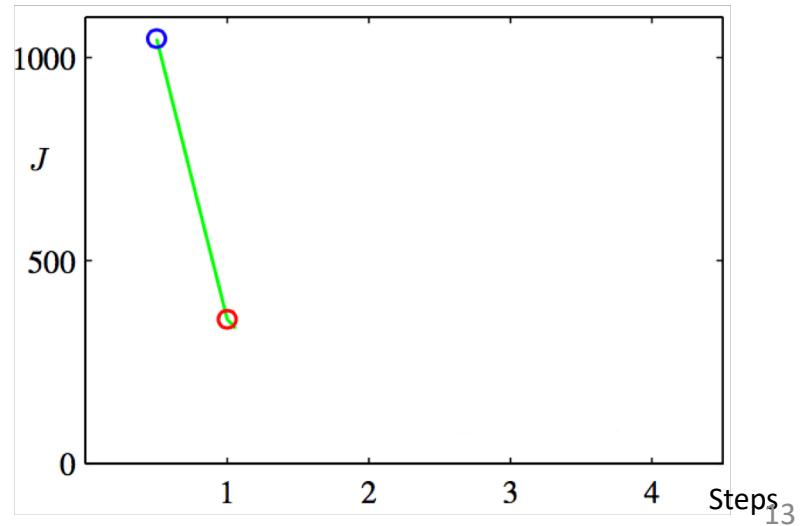
Given  $r_{n1}$ ,  $r_{n2}$ , calculate  $\mu_1, \mu_2$



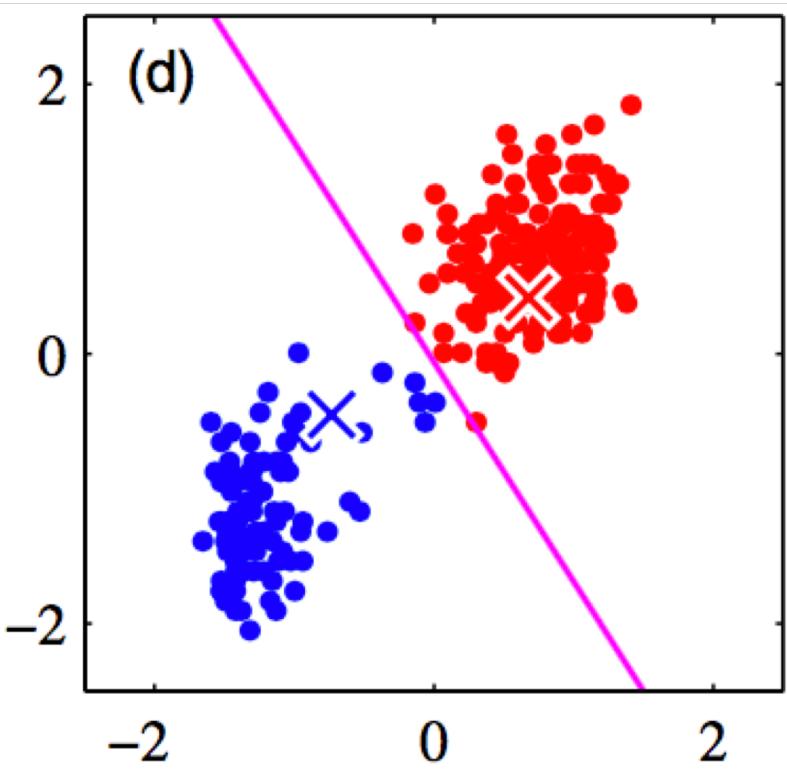
M Step

Calculate the means of the points in each cluster:

$$\mu_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$



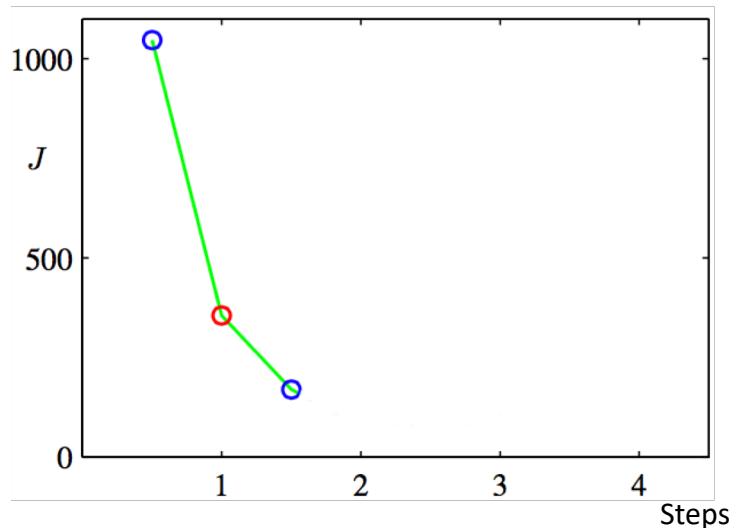
Given  $\mu_1, \mu_2$ , calculate  $r_{n1}, r_{n2}$  for all  
 $n=1, \dots, N$



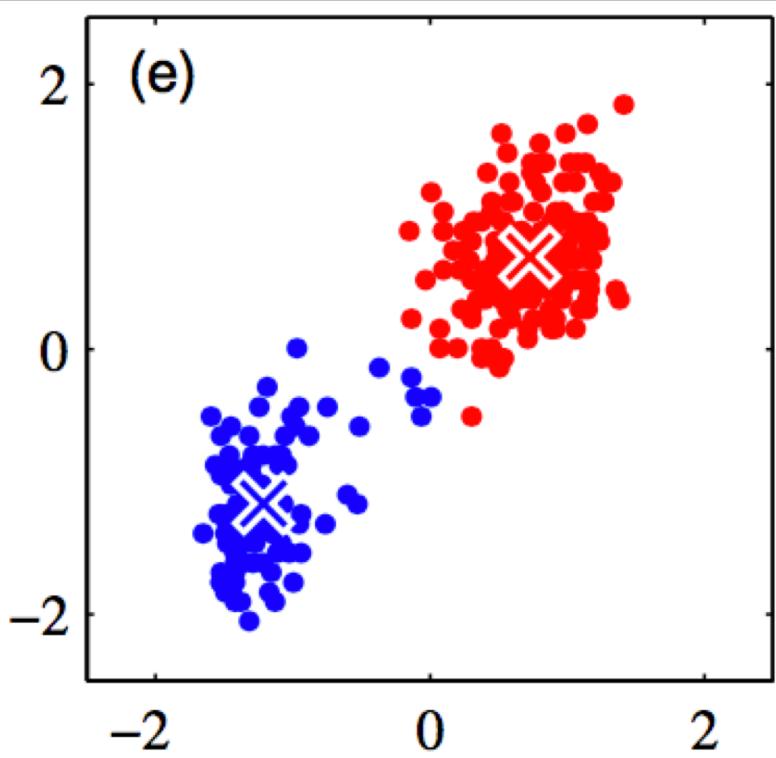
E Step

Assign the points to the nearest cluster:

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \| \mathbf{x}_n - \boldsymbol{\mu}_j \|^2 \\ 0 & \text{otherwise.} \end{cases}$$



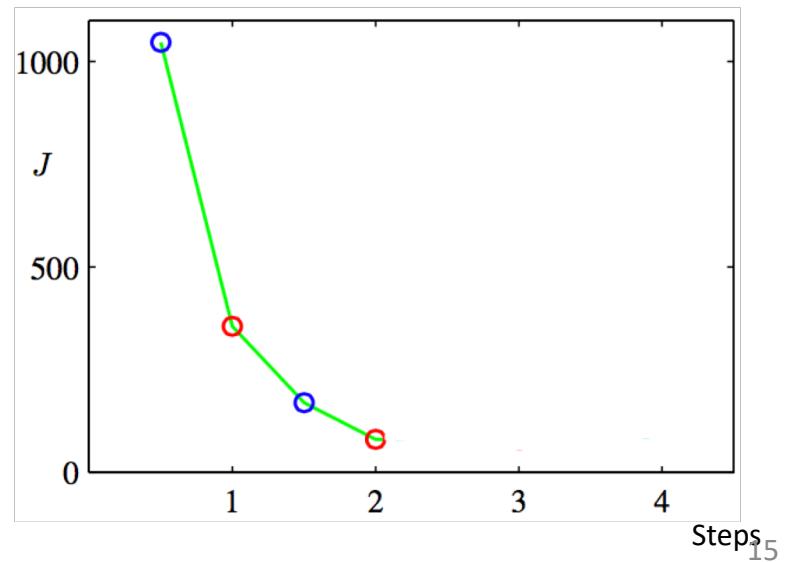
Given  $r_{n1}$ ,  $r_{n2}$ , calculate  $\mu_1, \mu_2$



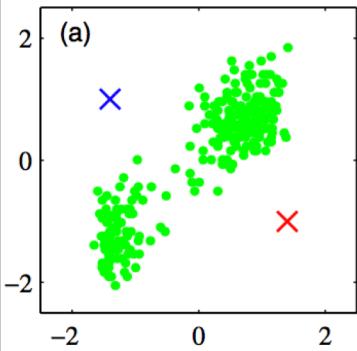
M Step

Calculate the means of the points in each cluster:

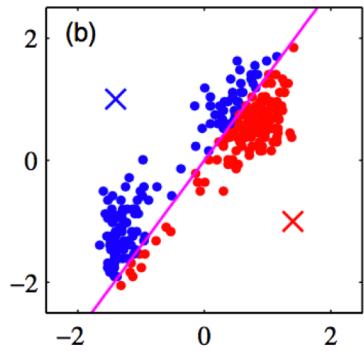
$$\mu_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$



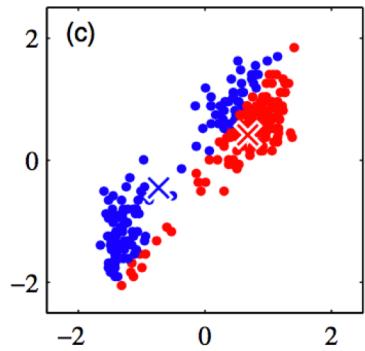
Initialization



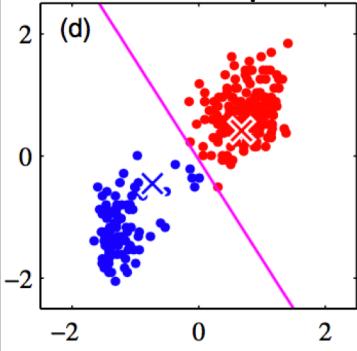
E-Step



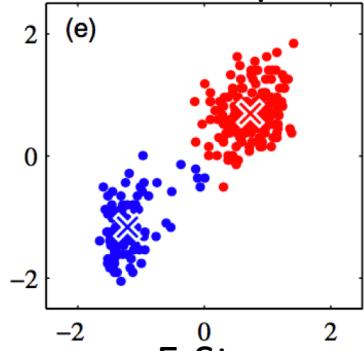
M-Step



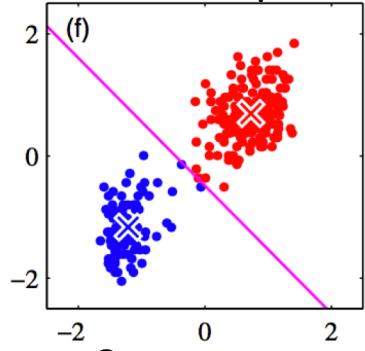
E-Step



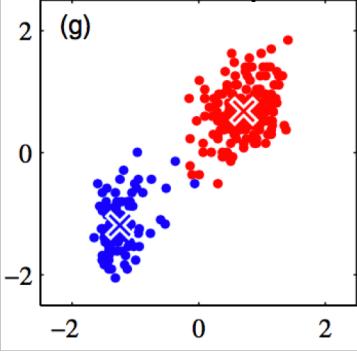
M-Step



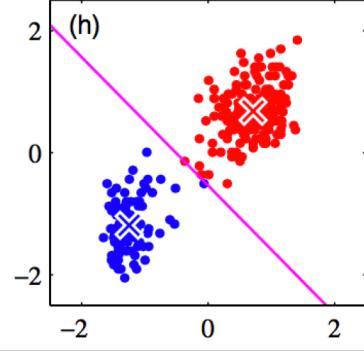
E-Step



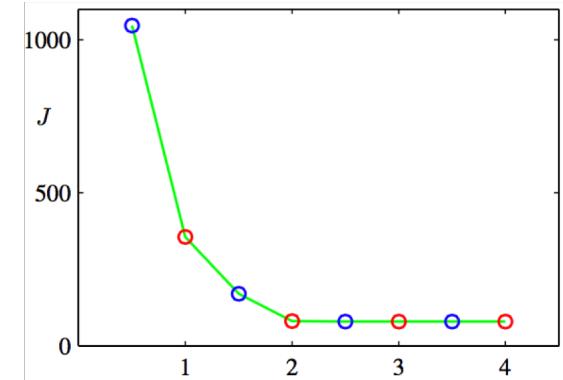
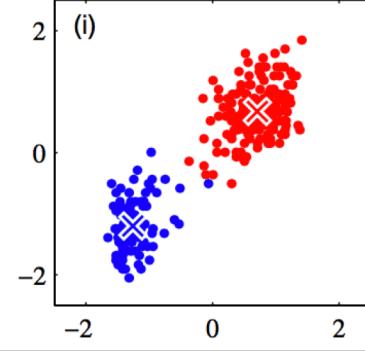
M-Step



E-Step



Convergence



If  $J$  does not change, or  $\{\mu_1, \mu_2\}$  do not change, then the algorithm converges.

# K均值法小结

- 初始化均值点  $\mu_1, \dots, \mu_k$
- 迭代如下
  - 把每个数据点按照就近原则分配给相应的  $\mu_i$
  - 把  $\mu_i$  更新为所分配的数据点的均值
- 迭代停止，如果聚类分配不变

Initialize  $\mathbf{m}_i, i = 1, \dots, k$ , for example, to  $k$  random  $\mathbf{x}^t$

Repeat

For all  $\mathbf{x}^t \in \mathcal{X}$

$$b_i^t \leftarrow \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$

For all  $\mathbf{m}_i, i = 1, \dots, k$

$$\mathbf{m}_i \leftarrow \sum_t b_i^t \mathbf{x}^t / \sum_t b_i^t$$

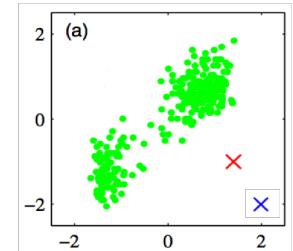
Until  $\mathbf{m}_i$  converge

# Basic ingredients

- Model or structure
- Objective function
- Algorithm
- Convergence

# Questions for K-mean algorithm

- Does it find the global optimum of  $J$ ?
  - No, the nearest local optimum, depending on initialization
- If Euclidean distance is not good for some data, do we have other choices?
- Can we assign each data point to the clusters probabilistically?
- If  $K$  (the total number of clusters) is unknown, can we estimate it from the data?



# Outline

- Clustering
  - K-mean clustering, **hierarchical clustering**
- Adaptive learning (online learning)
  - CL, FSCL, RPCL
- Gaussian Mixture Models (GMM)
- Expectation-Maximization (EM) for maximum likelihood

# Hierarchical Clustering

- $k$ -means clustering requires
  - $k$
  - Positions of initial centers
  - A distance measure between points (e.g. Euclidean distance)
- Hierarchical clustering requires a measure of distance between *groups* of data points

# Hierarchical Clustering

- Agglomerative clustering
- A very simple procedure:
  - Assign each data point into its own group
  - Repeat: look for the two closest groups and merge them into one group
  - Stop when all the data points are merged into a single cluster

# Distance Measure

- Distance between data points  $a$  and  $b$ :
  - $d(a, b)$
- Group A and B
  - Single-linkage

$$d(A, B) = \min_{a \in A, b \in B} d(a, b)$$

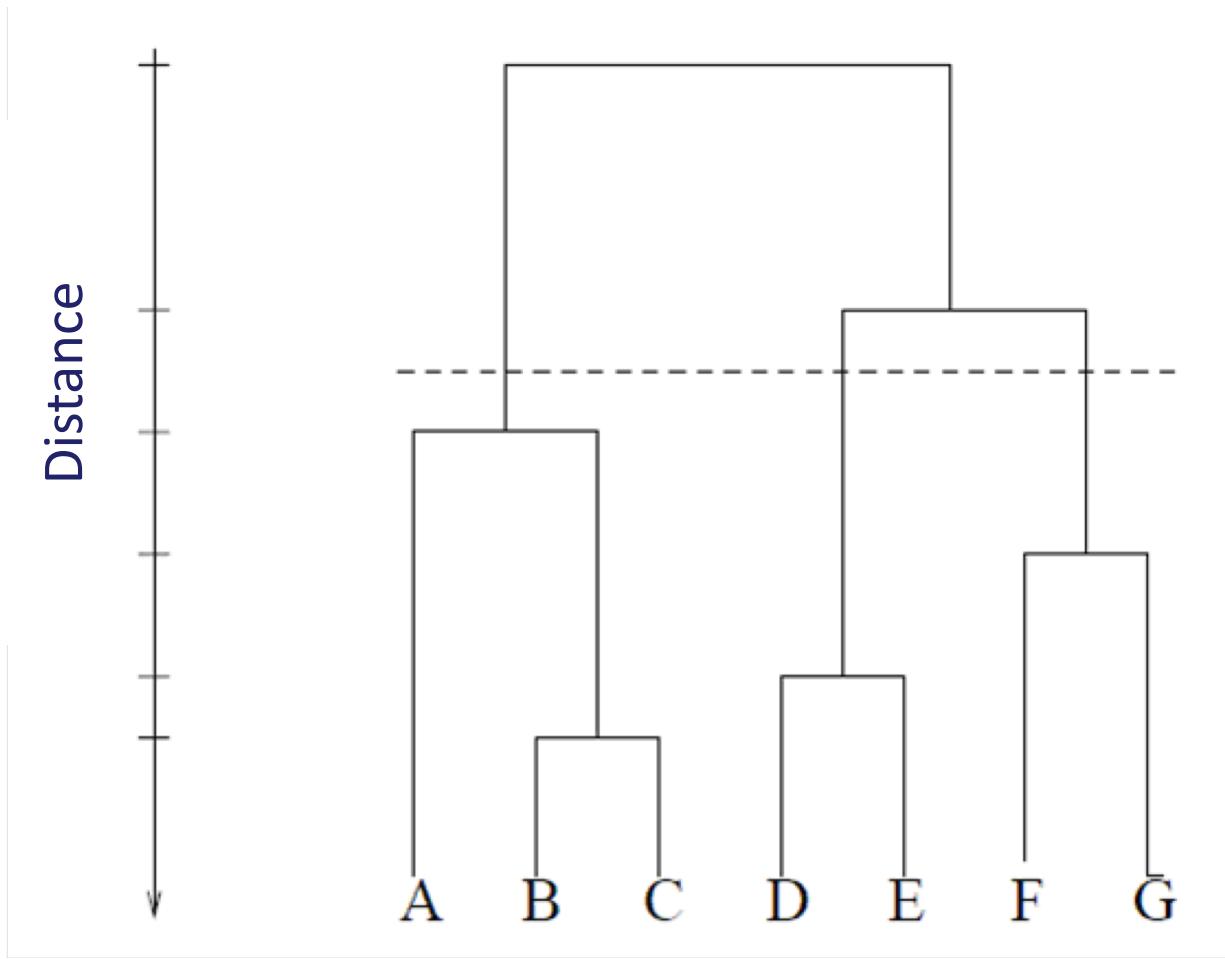
– Complete-linkage

$$d(A, B) = \max_{a \in A, b \in B} d(a, b)$$

– Average-linkage

$$d(A, B) = \frac{\sum_{a \in A, b \in B} d(a, b)}{|A| \cdot |B|}$$

# Dendrogram



# Outline

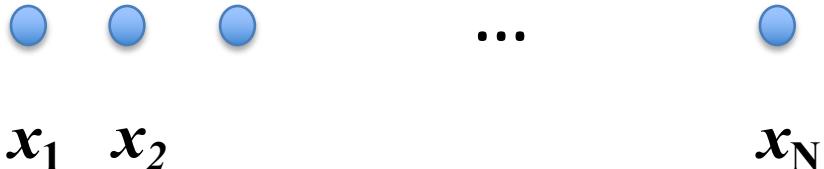
- Clustering
  - K-mean clustering, hierarchical clustering
- Adaptive learning (online learning)
  - CL, FSCL, RPCL
- Gaussian Mixture Models (GMM)
- Expectation-Maximization (EM) for maximum likelihood

# From batch to adaptive

- Given a batch of data points



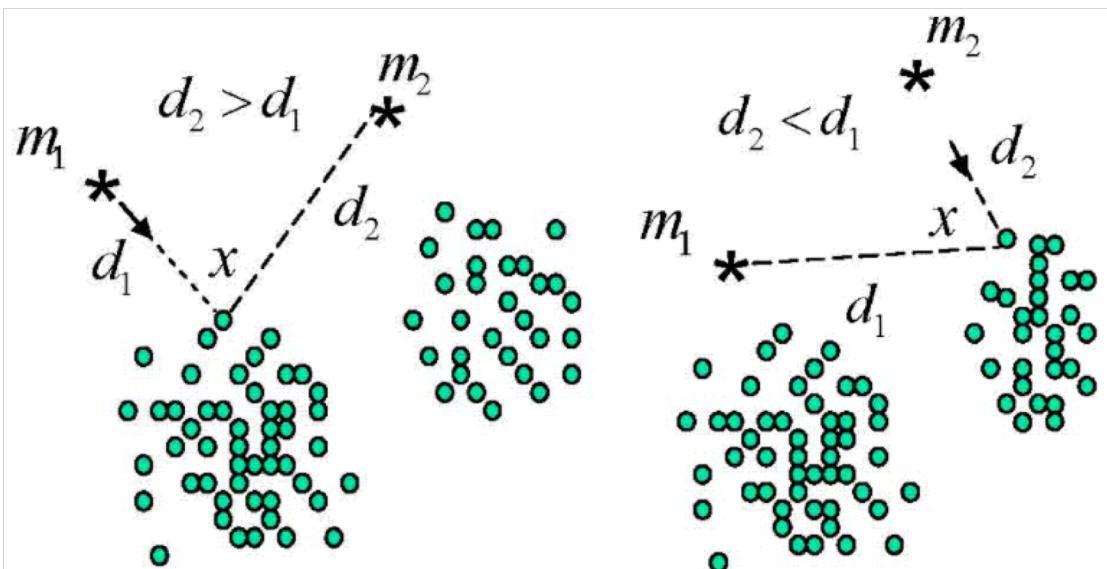
- Data points come one by one:



# Competitive learning

- Data points come one by one:

$x_1 \quad x_2 \quad \dots \quad x_N$



(a)  $m_1$  is the winner

(b)  $m_2$  is the winner

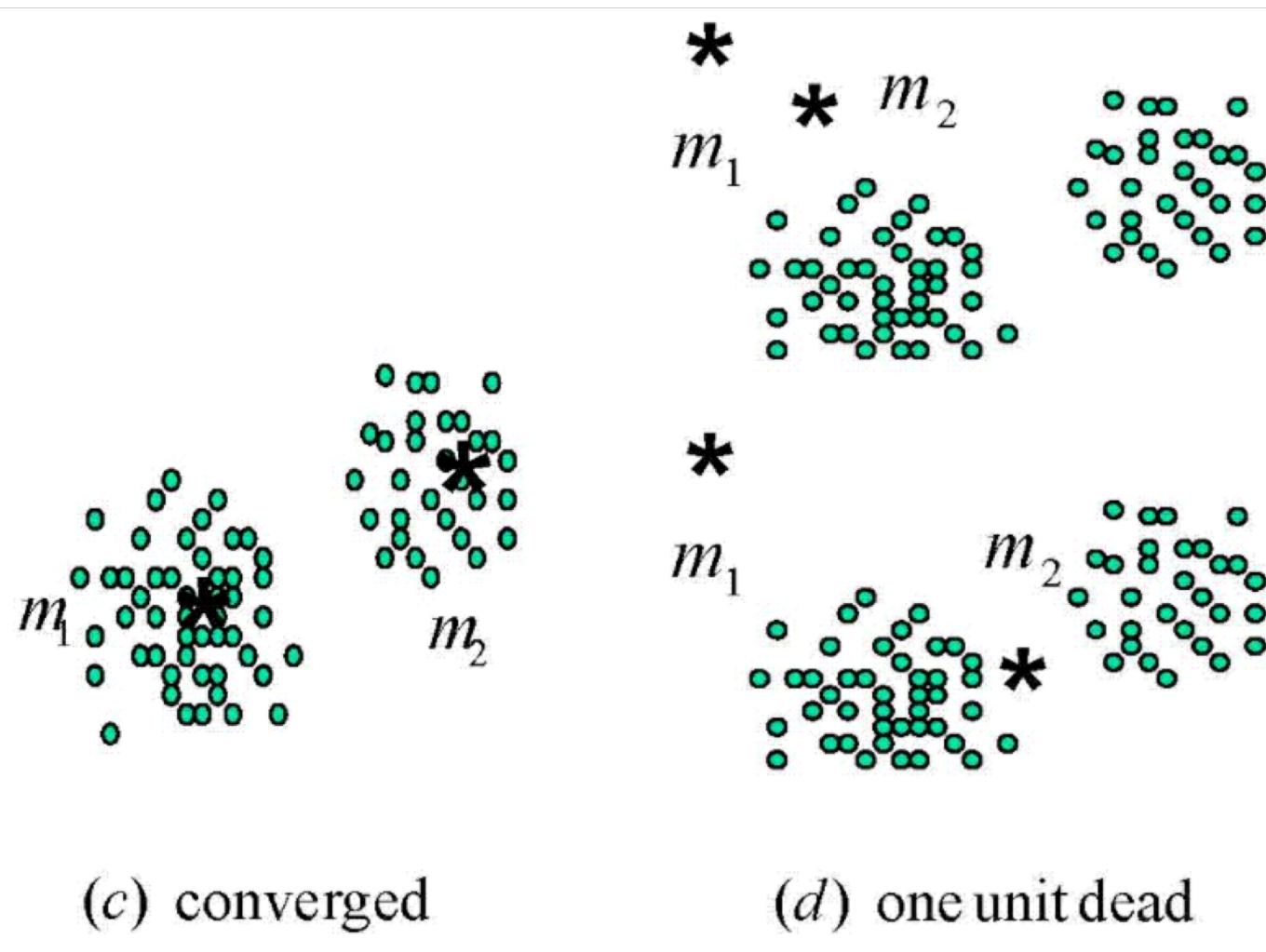
$$\varepsilon_t(\theta_j) = \|x_t - m_j\|^2$$

$$p_{j,t} = \begin{cases} 1, & \text{if } j = c, \\ 0, & \text{otherwise;} \end{cases}$$

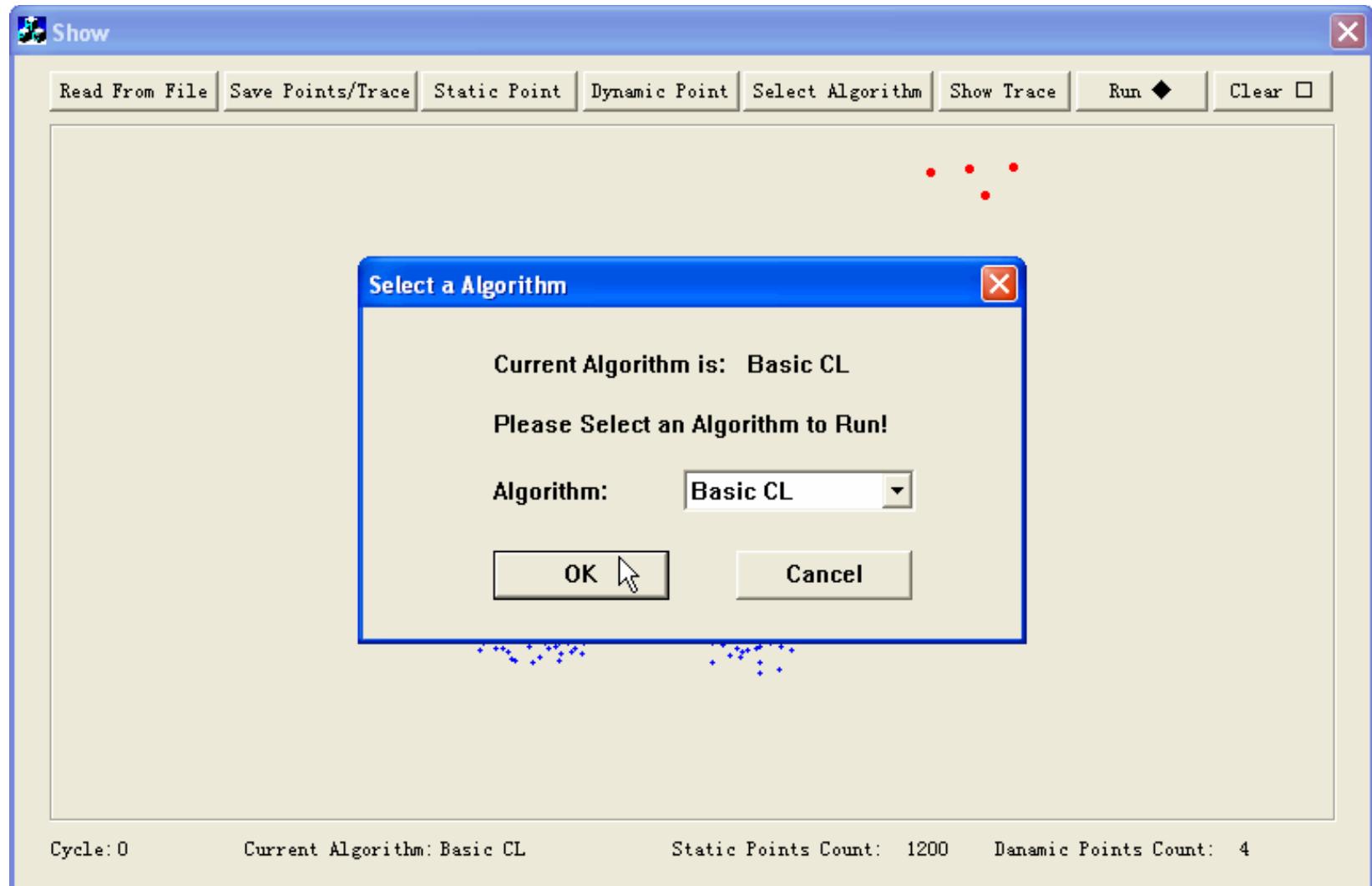
$$c = \arg \min_j \varepsilon_t(\theta_j).$$

$$m_j^{new} = m_j^{old} + \eta p_{j,t}(x_t - m_j^{old}).$$

# When starting with “bad initializations”



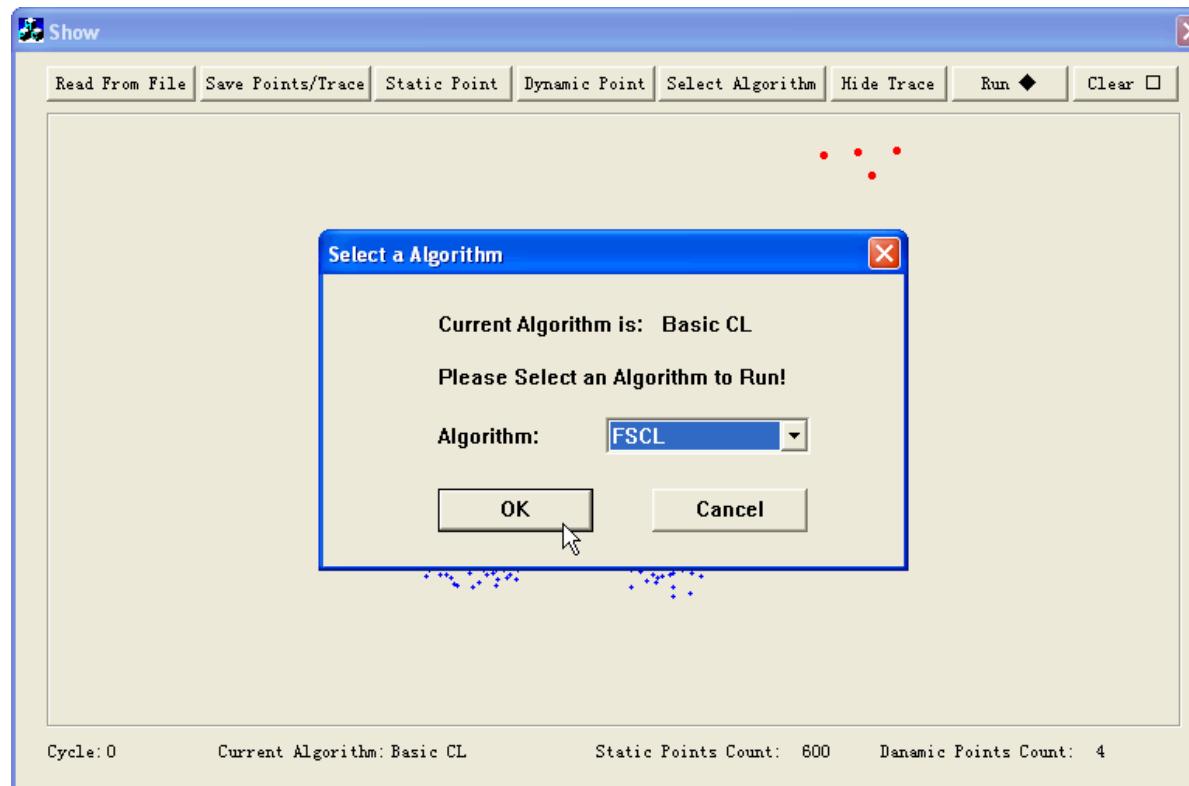
# A four-cluster case



# frequency sensitive competitive learning (FSCL) [Ahalt et al., 1990]

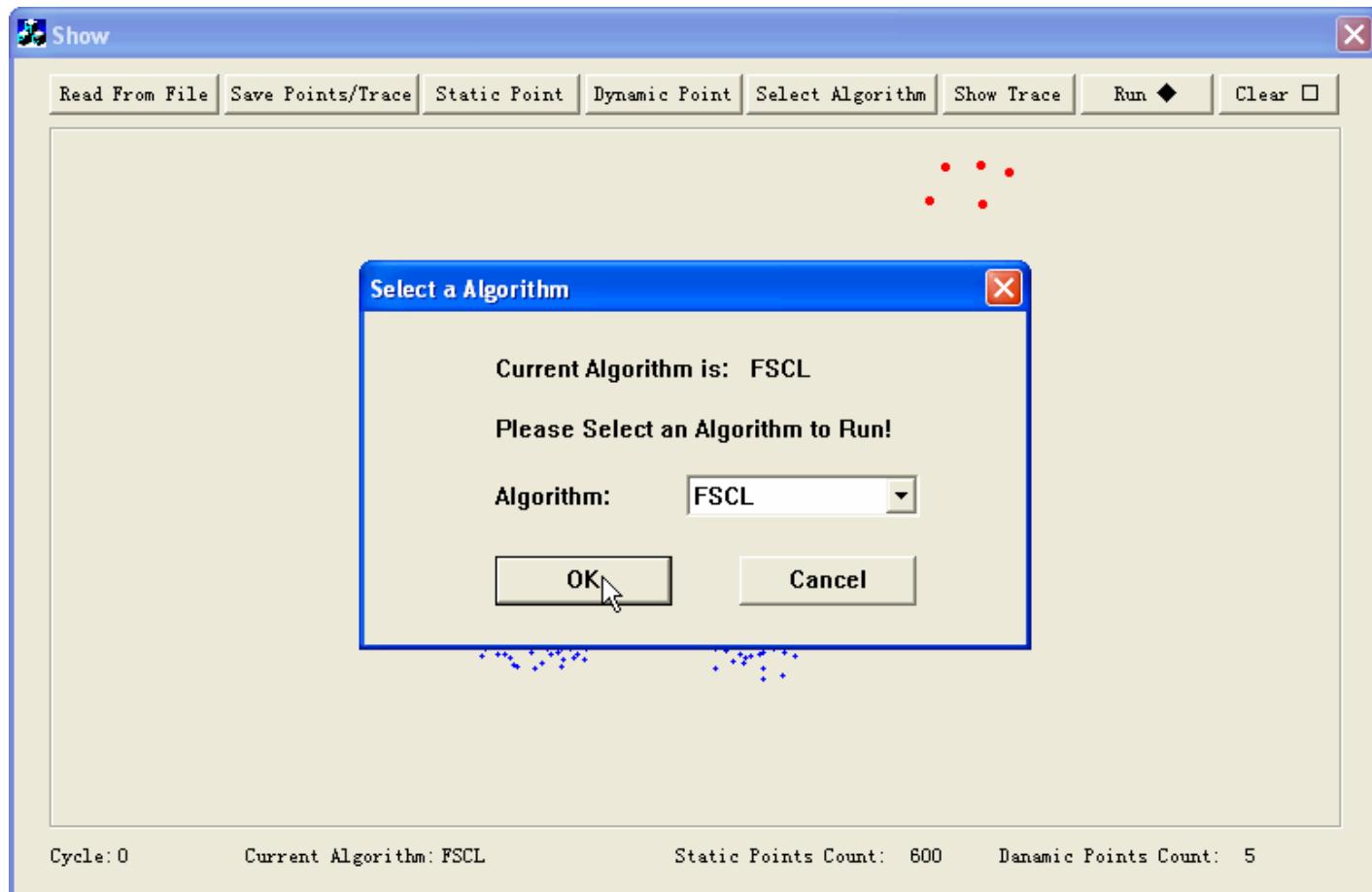
The idea is to penalize the frequent winners:

$$\varepsilon_t(\theta_j) = \alpha_j \|x_t - m_j\|^2$$



# FSCL is not good when there are extra centers

When k is pre-assigned to 5. the frequency sensitive mechanism also brings the extra one into data to disturb the correct locations of others



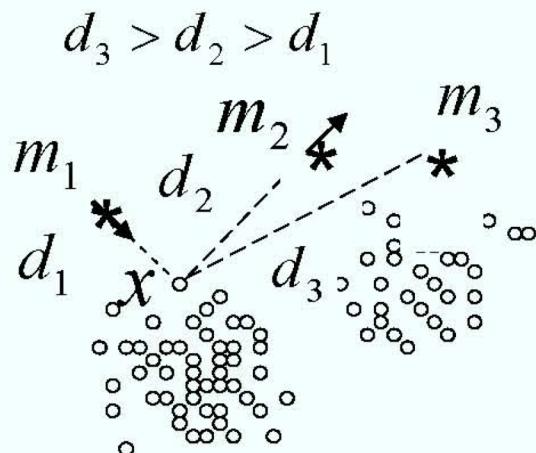
# Rival penalized competitive learning (RPCL)

(Xu, Krzyzak, & Oja, 1992 , 1993)

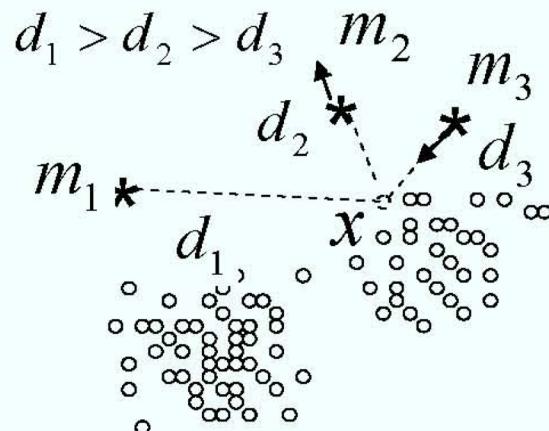
The RPCL differs from FSCL by implementing  $p_{j,t}$  as follows:

$$p_{j,t} = \begin{cases} 1, & \text{if } j = c, \\ -\gamma, & \text{if } j = r, \\ 0, & \text{otherwise,} \end{cases} \quad \left\{ \begin{array}{l} c = \arg \min_j \varepsilon_t(\theta_j), \\ r = \arg \min_{j \neq c} \varepsilon_t(\theta_j), \end{array} \right.$$

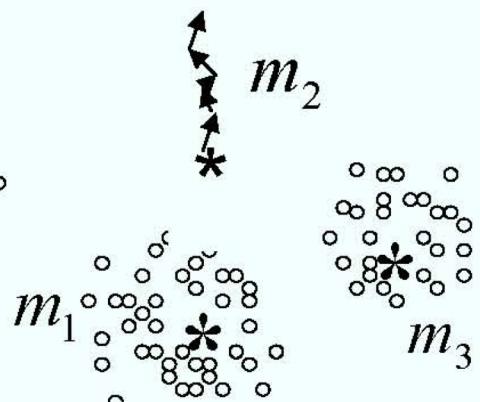
where  $\gamma$  approximately takes a number between 0.05 and 0.1 for controlling the penalizing strength.



(a)  $m_1$  is the winner  
 $m_2$  is the rival

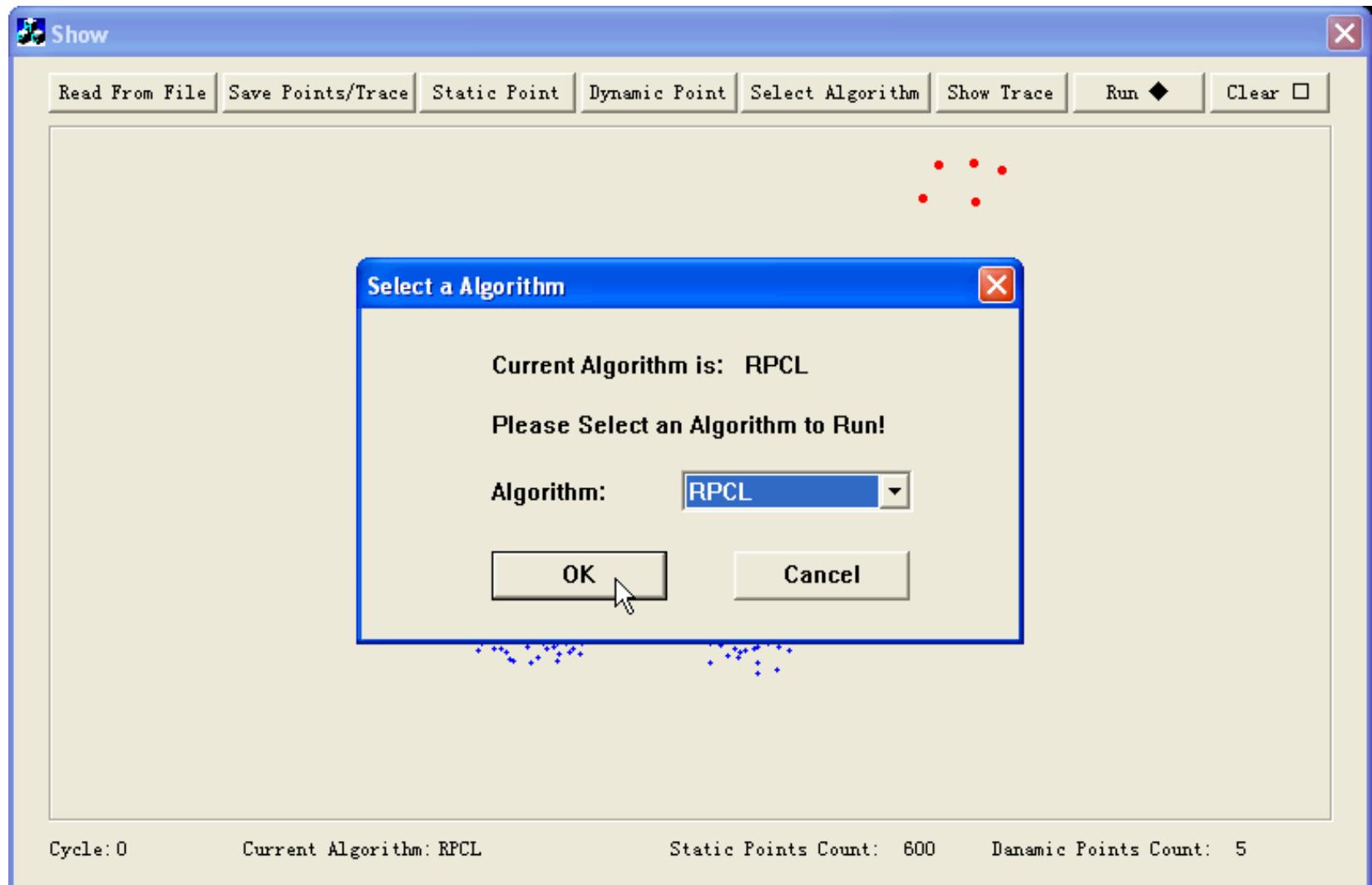


(b)  $m_3$  is the winner  
 $m_2$  is the rival



(c)  $m_1$  and  $m_3$  are converged  
 $m_2$  is driven far away

Rival penalized mechanism makes extra agents driven far away.



Thank you!