

**ΟΙΚΟΝΟΜΙΚΟ
ΠΑΝΕΠΙΣΤΗΜΙΟ
ΑΘΗΝΩΝ**



**ATHENS UNIVERSITY
OF ECONOMICS
AND BUSINESS**

Εξόρυξη γνώσης από Βάσεις Δεδομένων και τον Παγκόσμιο Ιστό

Ενότητα # 4: Unsupervised Learning (Clustering)

Διδάσκων: Μιχάλης Βαζιργιάννης

Τμήμα: Προπτυχιακό Πρόγραμμα Σπουδών
“Πληροφορικής”



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ & ΘΡΗΣΚΕΥΜΑΤΩΝ, ΠΟΛΙΤΙΣΜΟΥ & ΑΘΛΗΤΙΣΜΟΥ
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στα πλαίσια του εκπαιδευτικού έργου του διδάσκοντα.
- Το έργο «**Ανοικτά Ακαδημαϊκά Μαθήματα στο Οικονομικό Πανεπιστήμιο Αθηνών**» έχει χρηματοδοτήσει μόνο τη αναδιαμόρφωση του εκπαιδευτικού υλικού.
- Το έργο υλοποιείται στο πλαίσιο του Επιχειρησιακού Προγράμματος «Εκπαίδευση και Δια Βίου Μάθηση» και συγχρηματοδοτείται από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο) και από εθνικούς πόρους.



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ & ΘΡΗΣΚΕΥΜΑΤΩΝ, ΠΟΛΙΤΙΣΜΟΥ & ΑΘΛΗΤΙΣΜΟΥ
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



Άδειες Χρήσης

- Το παρόν εκπαιδευτικό υλικό υπόκειται σε άδειες χρήσης Creative Commons.
- Οι εικόνες προέρχονται



Σκοποί ενότητας

Εισαγωγή και εξοικείωση με τις μεθόδους,
Clustering, K-means, Expectation Maximization
(EM), Spectral Clustering.

Περιεχόμενα ενότητας

- Clustering
- K-means
- Expectation Maximization (EM)
- Spectral Clustering

**ΟΙΚΟΝΟΜΙΚΟ
ΠΑΝΕΠΙΣΤΗΜΙΟ
ΑΘΗΝΩΝ**



**ATHENS UNIVERSITY
OF ECONOMICS
AND BUSINESS**

Clustering

Μάθημα: Εξόρυξη γνώσης από Βάσεις Δεδομένων και τον Παγκόσμιο Ιστό

Ενότητα # 4: Unsupervised Learning (Clustering)

Διδάσκων: Μιχάλης Βαζιργιάννης

Τμήμα: Προπτυχιακό Πρόγραμμα Σπουδών “Πληροφορικής”

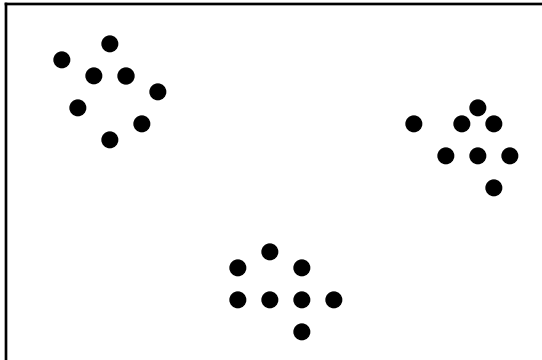
Supervised vs. Unsupervised Learning

- **Unsupervised learning (clustering)**
 - The class labels of training data are unknown
 - Given a set of measurements, observations, etc. establish the existence of clusters in the data
- **Supervised learning (classification)**
 - **Supervision:** The training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations
 - New data is classified based on the training set
- **Semi-supervised clustering**
 - Learning approaches that use **user input** (i.e. constraints or labeled data)
 - Clusters are defined so that user-constraints are satisfied

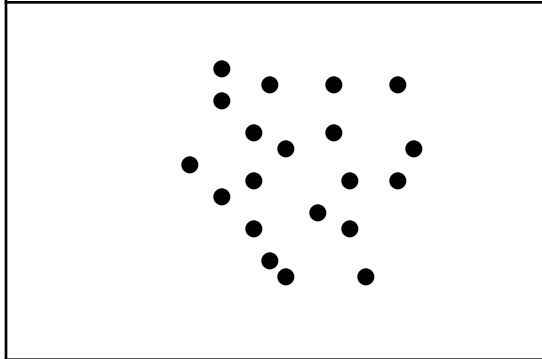
Clustering

- “automated detection of group structure in data”
 - Typically: partition N data points into K groups (clusters) such that the points in each group are more similar to each other than to points in other groups
 - descriptive technique (contrast with predictive)
 - for real-valued vectors, clusters can be thought of as clouds of points in p -dimensional space

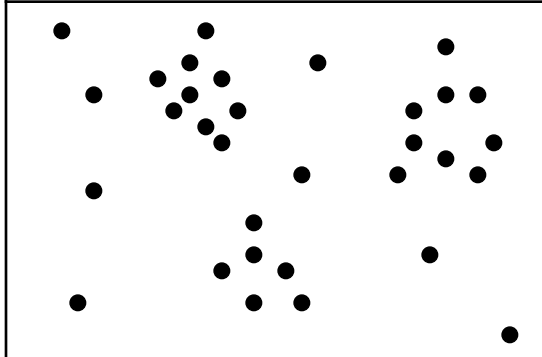
Clustering



Sometimes easy



Sometimes impossible



and sometimes in between

Why is Clustering useful?

- “Discovery” of new knowledge from data
 - Contrast with supervised classification (where labels are known)
 - Long history in the sciences of categories, taxonomies, etc
 - Can be very useful for summarizing large data sets
 - For large n and/or high dimensionality
- Applications of clustering
 - Discovery of new types of galaxies in astronomical data
 - Clustering of genes with similar expression profiles
 - Cluster pixels in an image into regions of similar intensity
 - Segmentation of customers for an e-commerce store
 - Clustering of documents produced by a search engine
 - many more

General Issues in Clustering

- Representation:
 - What types of clusters are we looking for?
- Score:
 - The criterion to compare one clustering to another
- Optimization
 - Generally, finding the optimal clustering is NP-hard
 - Greedy algorithms to optimize score are widely used
- Other issues
 - Distance function, $D(x(i),x(j))$ critical aspect of clustering, both
 - distance of pairs of objects
 - distance of objects from clusters
 - How is K selected?
 - Different types of data
 - Real-valued versus categorical
 - Attribute-valued vectors vs. n^2 distance matrix

Clustering Methods

- **Partitional algorithms**
 - K-Means, PAM, CLARA, CLARANS [Ng and Han, VLDB 1994]
- **Hierarchical algorithms**
 - CURE [Guha et al, SIGMOD'98], BIRCH [Zhang et al, SIGMOD'96], CHAMELEON [IEEE Computer, 1999]
- **Density based algorithms**
 - DENCLUE [Hinneburg, Keim, KDD'98], DBSCAN [Ester et al, KDD 96]
- **Subspace Clustering**
 - CLIQUE [Agrawal et al, SIGMOD'98], PROCLUS [Agrawal et al, SIGMOD'99], ORCLUS: [Aggarwal, and Yu, SIGMOD' 00], DOC: [Procopiuc, Jones, Agarwal, and Murali, SIGMOD'02]
- **Locally adaptive clustering techniques**
 - LAC
- **Spectral clustering**
 - [Ng, Jordan, Weiss], [Shi/Malik], [Scott/Longuet-Higgins], [Perona/ Freeman]

Partitional Algorithms: Basic Concept

- Partitional method:

- Partition the data set into a set of k disjoint partitions (clusters).

- Problem Definition:

- Given an integer k , find a partitioning of k clusters that optimizes the chosen partitioning criterion

**ΟΙΚΟΝΟΜΙΚΟ
ΠΑΝΕΠΙΣΤΗΜΙΟ
ΑΘΗΝΩΝ**



**ATHENS UNIVERSITY
OF ECONOMICS
AND BUSINESS**

K-means

Μάθημα: Εξόρυξη γνώσης από Βάσεις Δεδομένων και τον Παγκόσμιο Ιστό

Ενότητα # 4: Unsupervised Learning (Clustering)

Διδάσκων: Μιχάλης Βαζιργιάννης

Τμήμα: Προπτυχιακό Πρόγραμμα Σπουδών “Πληροφορικής”

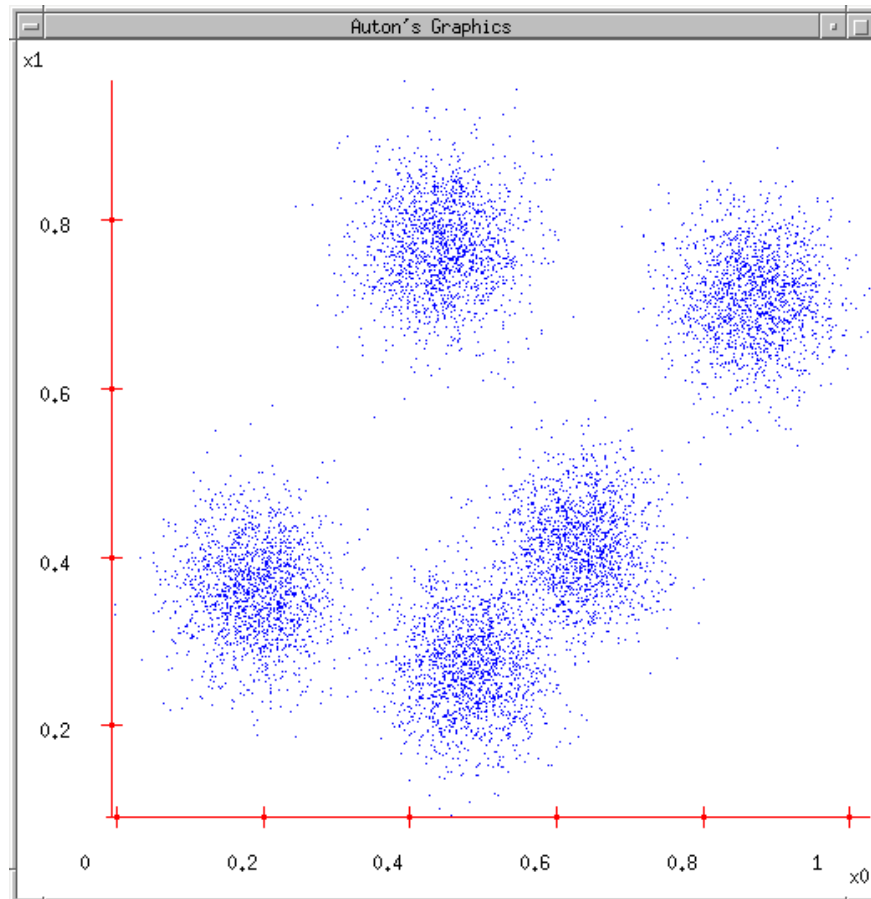
K-means Clustering

- basic idea:
 - Score = $wc(C)$ = sum-of-squares within cluster distance
 - start with randomly chosen cluster centers $c_1 \dots c_k$
 - repeat until no cluster memberships change:
 - assign each point x to cluster with nearest center
 - find smallest $d(\underline{x}, \underline{c}_i)$, over all $\underline{c}_1 \dots \underline{c}_k$
 - recompute cluster centers over data assigned to them
 - $\underline{c}_i = 1/(n_i) \sum_{x \in C_i} \underline{x}$
- algorithm terminates (finite number of steps)
 - decreases $\text{Score}(X, C)$ each iteration membership changes
- converges to local maxima of $\text{Score}(X, C)$
 - not necessarily the global maxima ...
 - different initial centers (seeds) can lead to diff local maxs

K-means Complexity

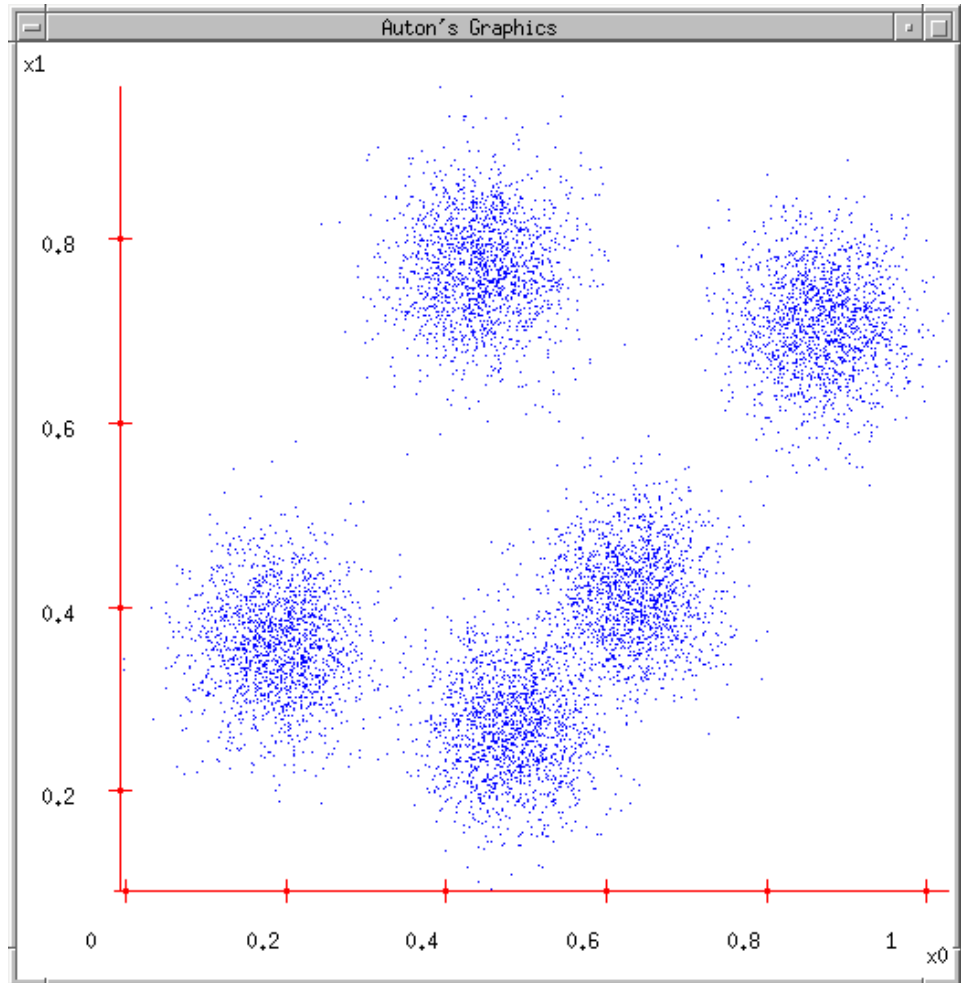
- **time complexity:** $O(I e \underline{n} k)$ \ll exhaustive's n^k
 - I = number of iterations (steps)
 - e = cost of distance computation ($e=p$ for Euclidian dist)
- **speed-up tricks** (especially useful in early iterations)
 - use nearest $x(i)$'s as cluster centers instead of mean
 - reuse of cached dists from size n^2 dist mat D (lowers effective “ e ”)
 - k-medoids: use one of $x(i)$'s as center because mean not defined
 - recompute centers as points reassigned
 - useful for large n (like online neural nets) & more cache efficient
 - PCA: reduce effective “ e ” and/or fit more of X in RAM
 - “condense”: reduce “ n ” by replace group with prototype
 - even more clever data structures (see work by Andrew Moore, CMU)

K-means example (courtesy of Andrew Moore, CMU)



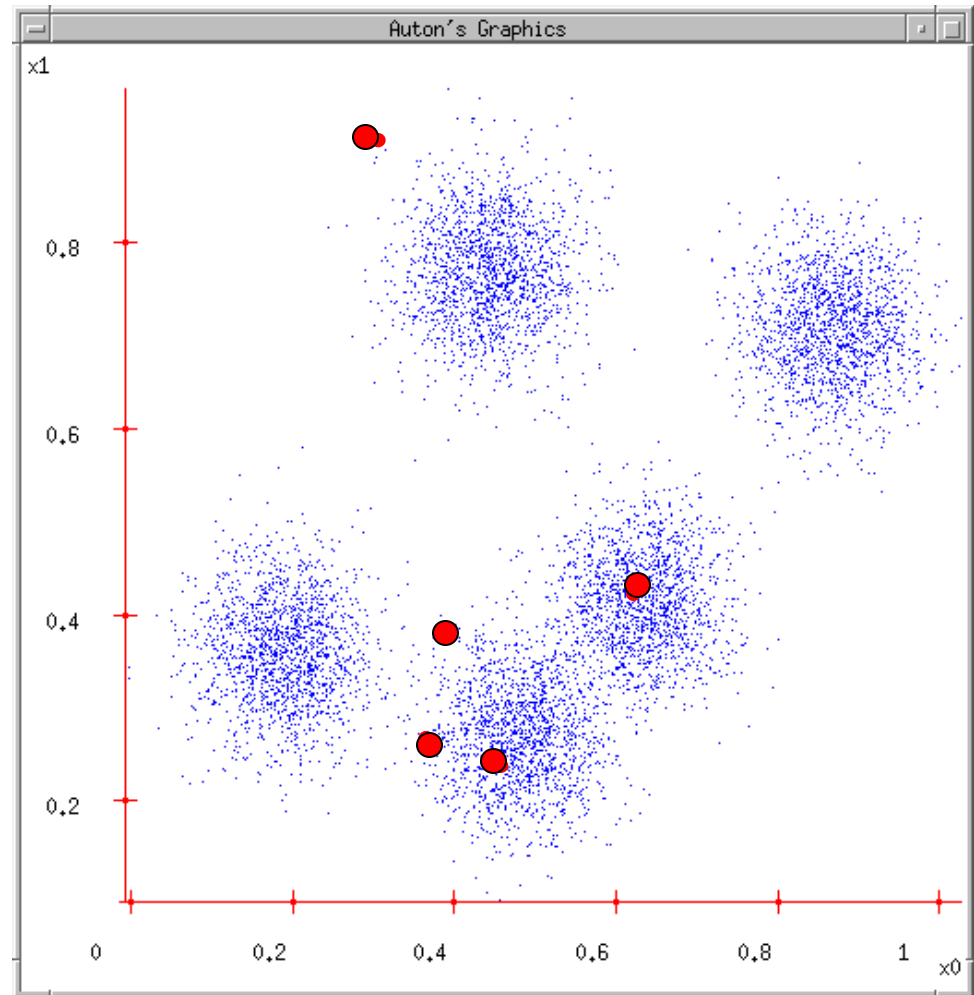
K-means

1. Ask user how many clusters they'd like.
(e.g. $K=5$)



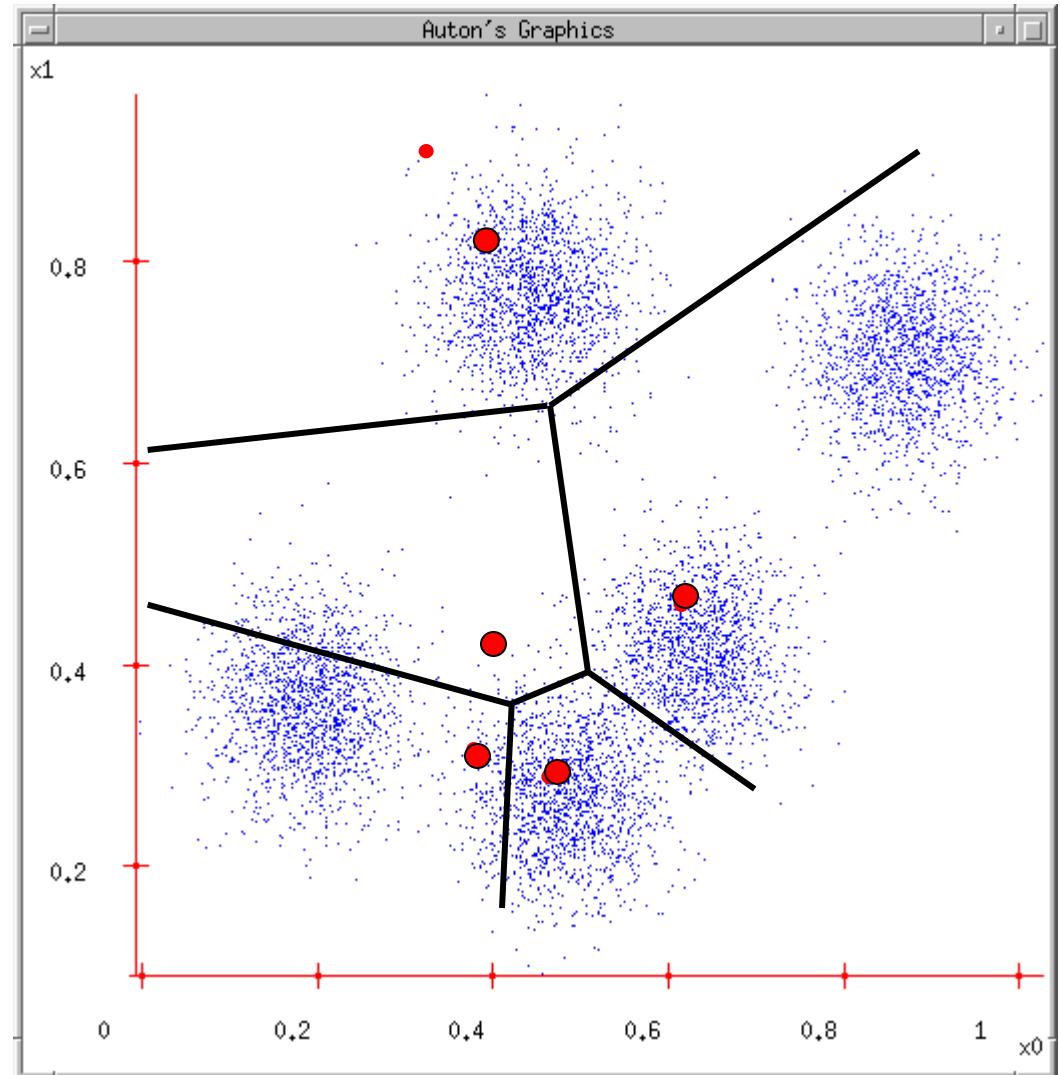
K-means

1. Ask user how many clusters they'd like. (e.g. $K=5$)
2. Randomly guess K cluster Center locations



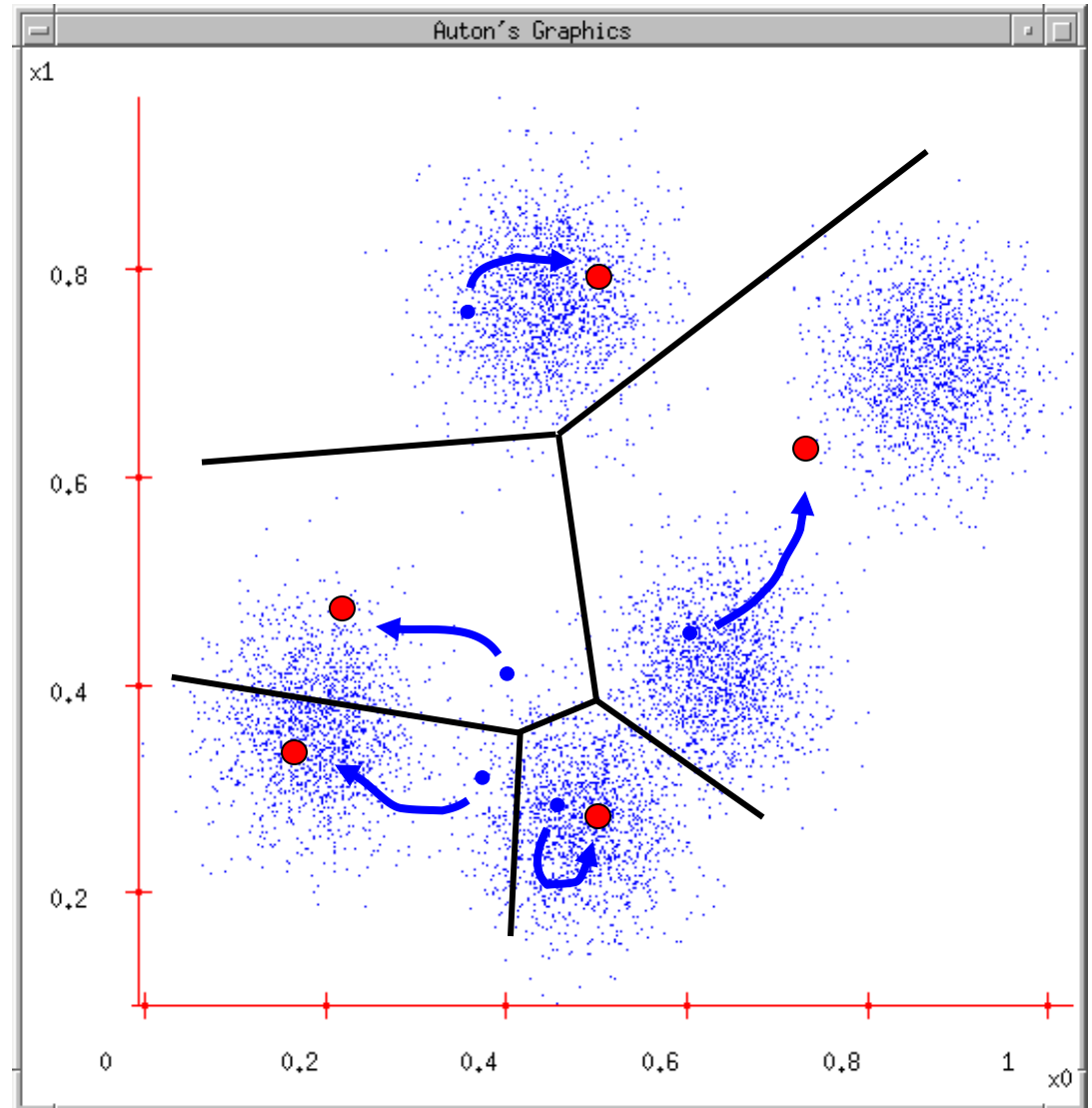
K-means

1. Ask user how many clusters they'd like. (e.g. $K=5$)
2. Randomly guess K cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)



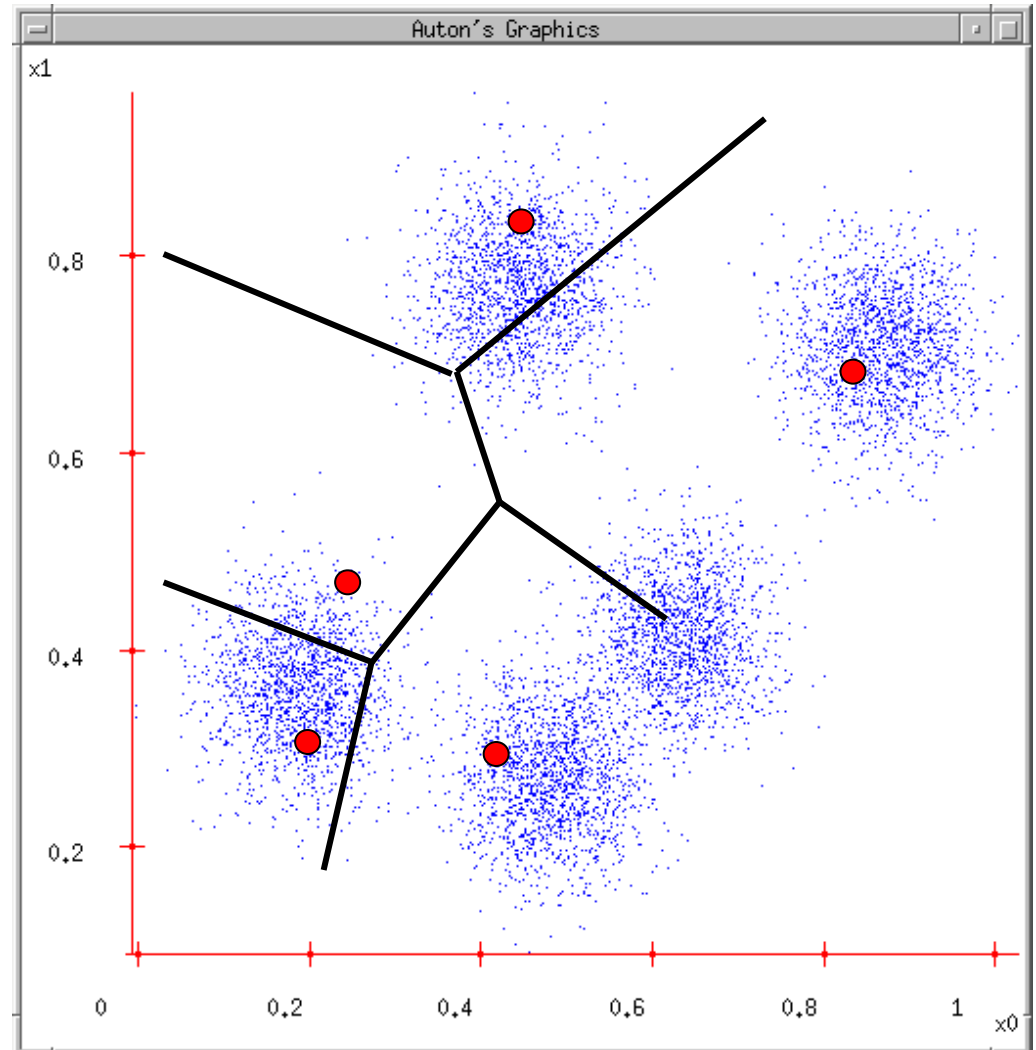
K-means

1. Ask user how many clusters they'd like. (e.g. $k=5$)
2. Randomly guess k cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns



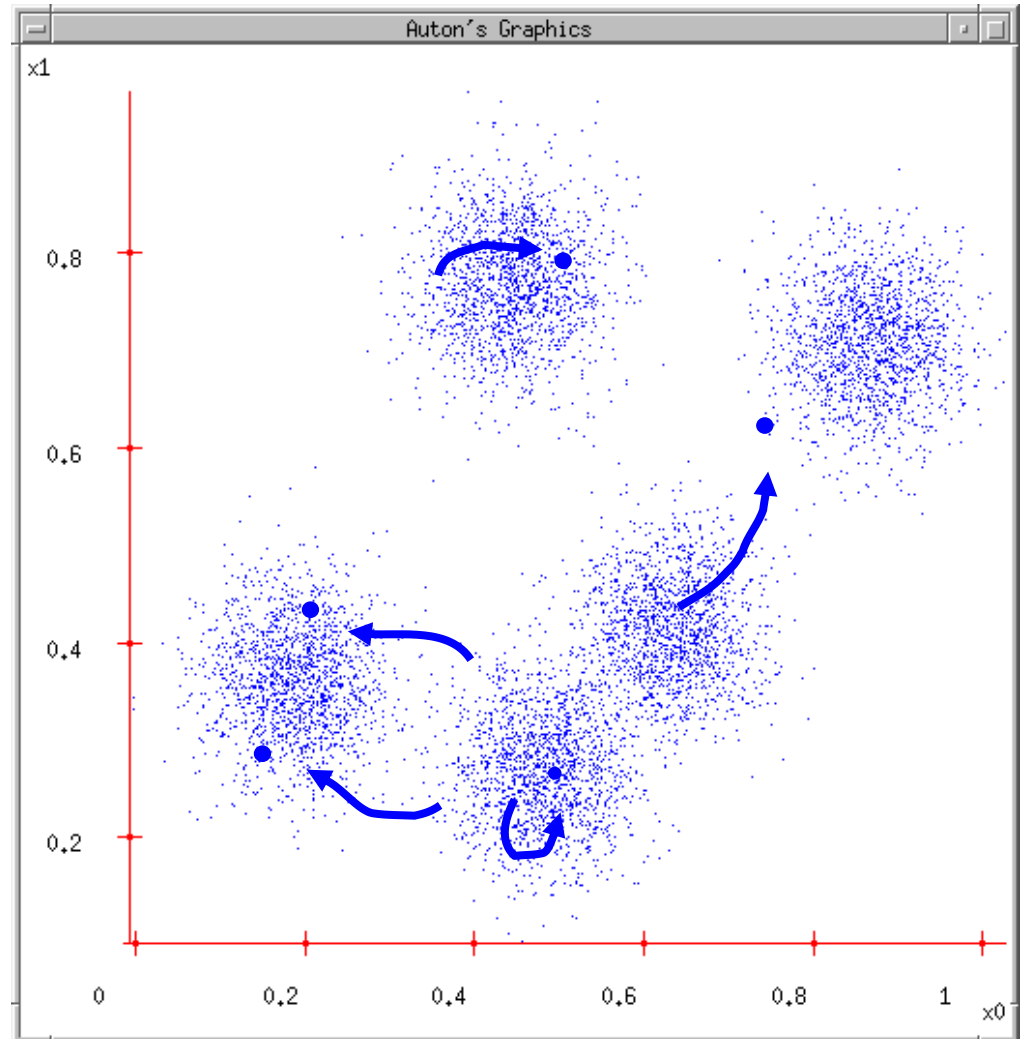
K-means

1. Ask user how many clusters they'd like. (e.g. $k=5$)
2. Randomly guess k cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns
5. New Centers => new boundaries
6. Repeat until no change!



K-means

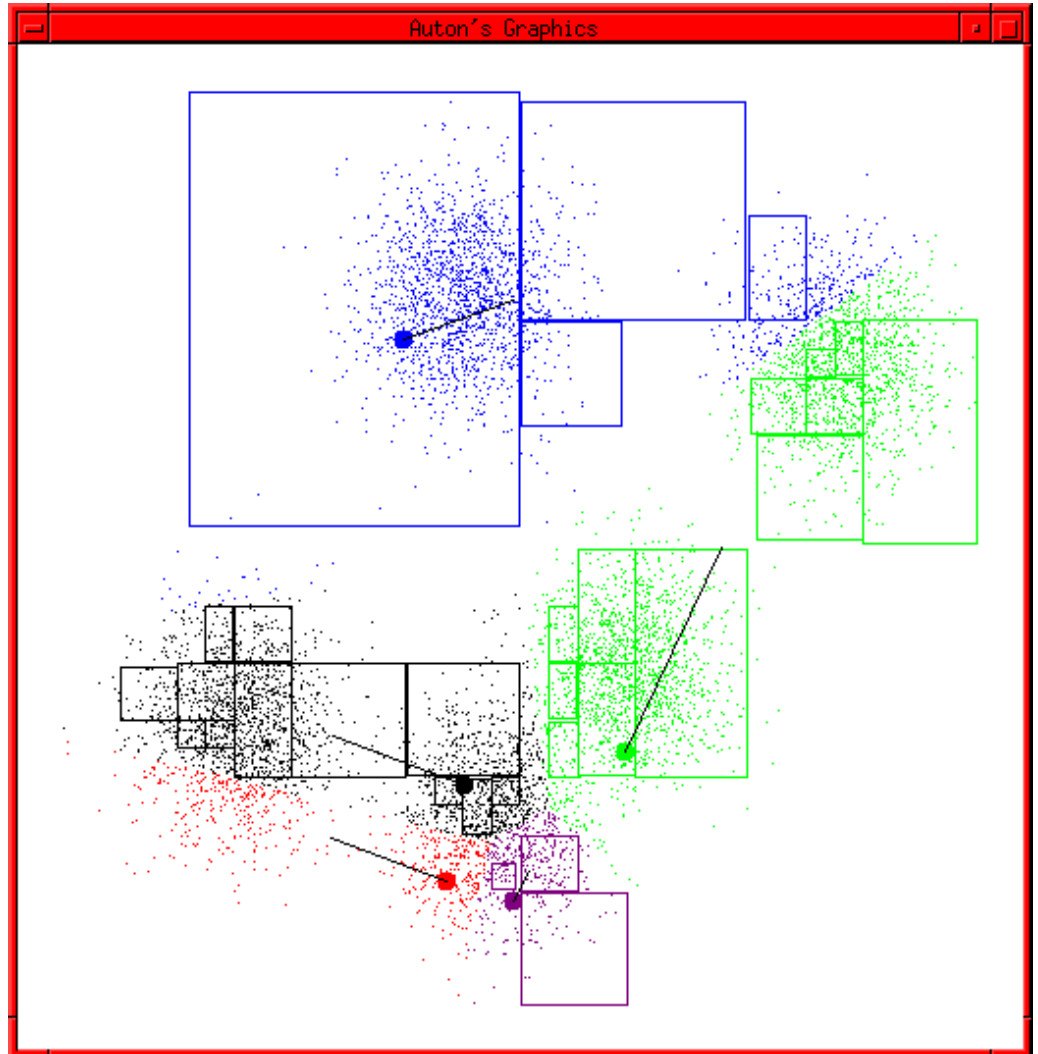
1. Ask user how many clusters they'd like. (e.g. $k=5$)
2. Randomly guess k cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns...
5. ...and jumps there
6. ...Repeat until terminated!



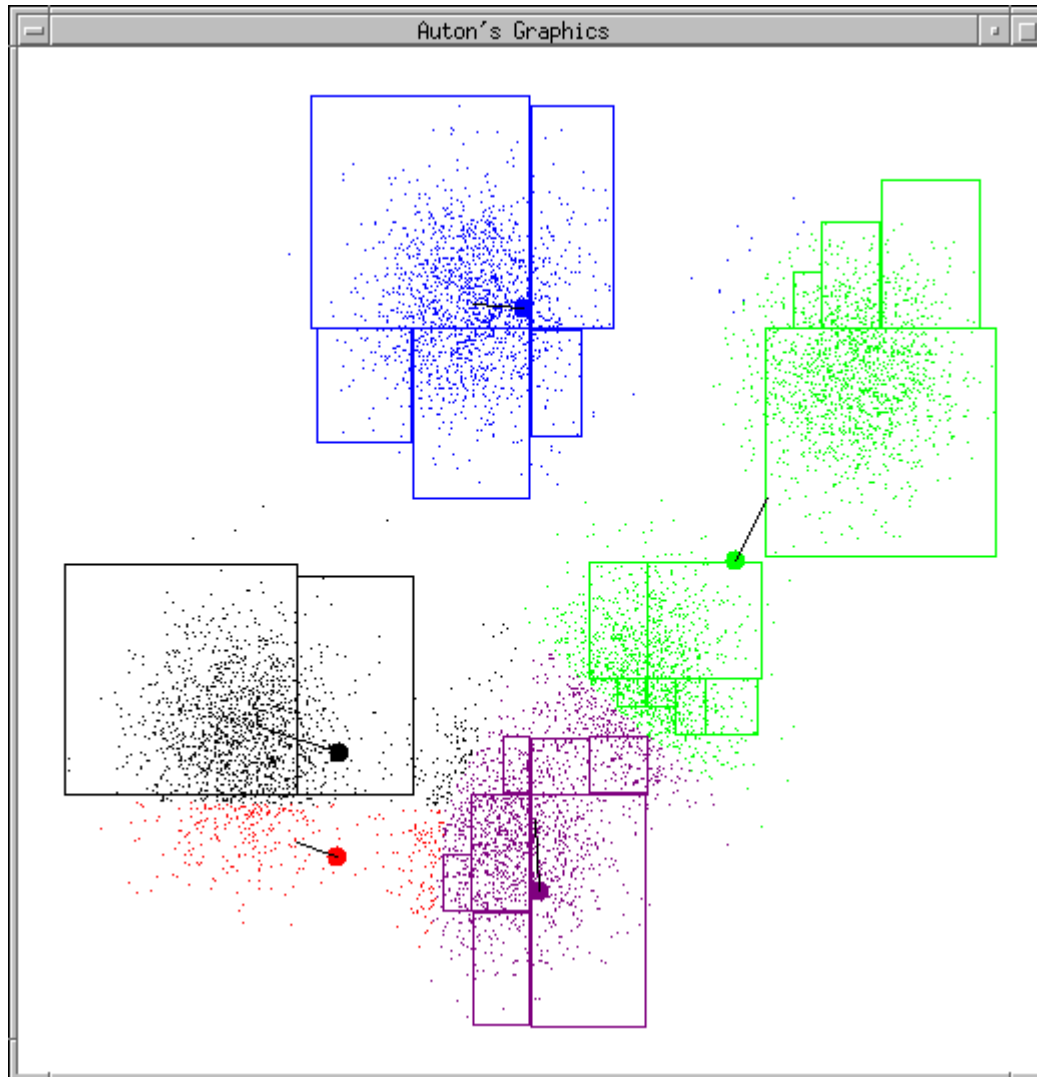
Accelerated Computations

Example generated by Pelleg and Moore's accelerated k-means

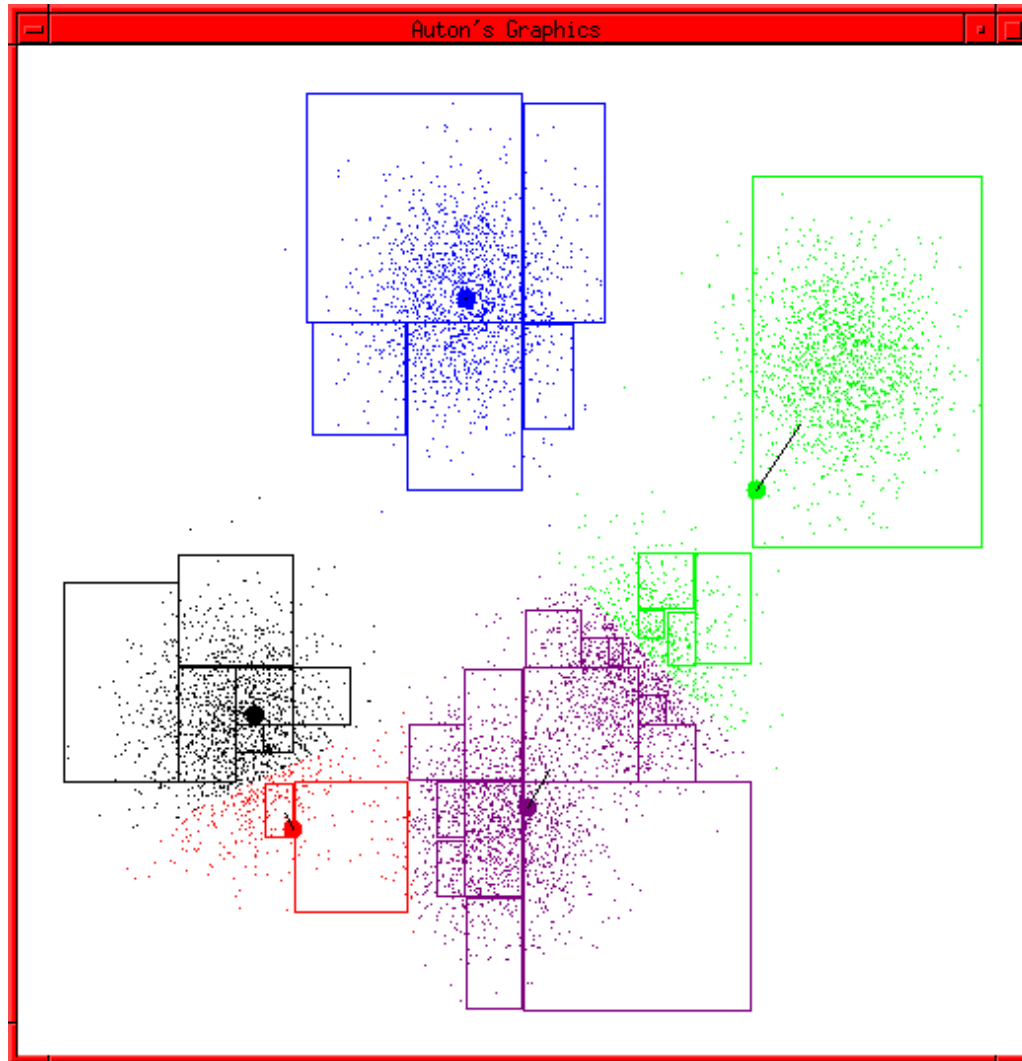
*Dan Pelleg and Andrew Moore.
Accelerating Exact k-means
Algorithms with Geometric
Reasoning. Proc. Conference
on Knowledge Discovery in
Databases 1999, (KDD99)
(available on
www.autonlab.org/pap.html)*



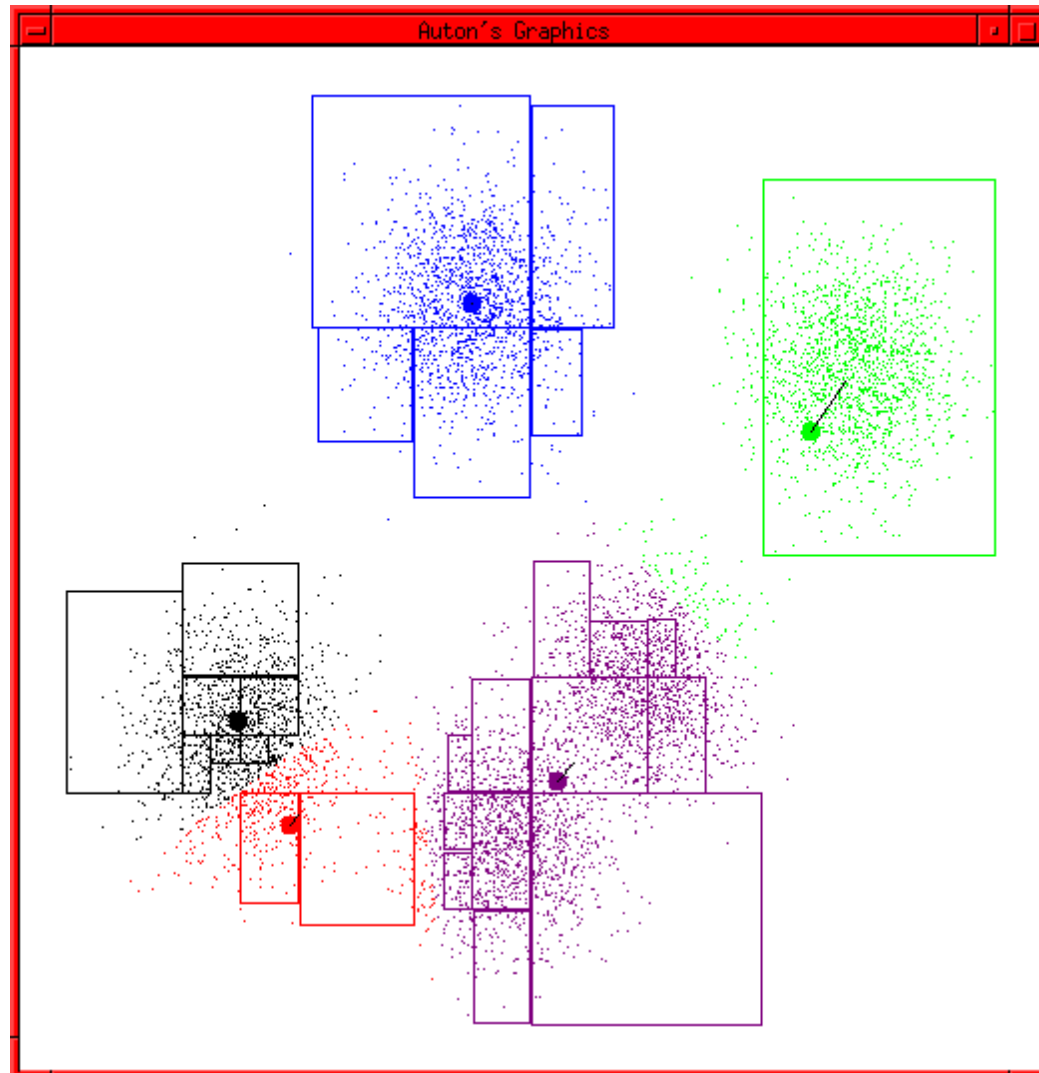
K-means continues...



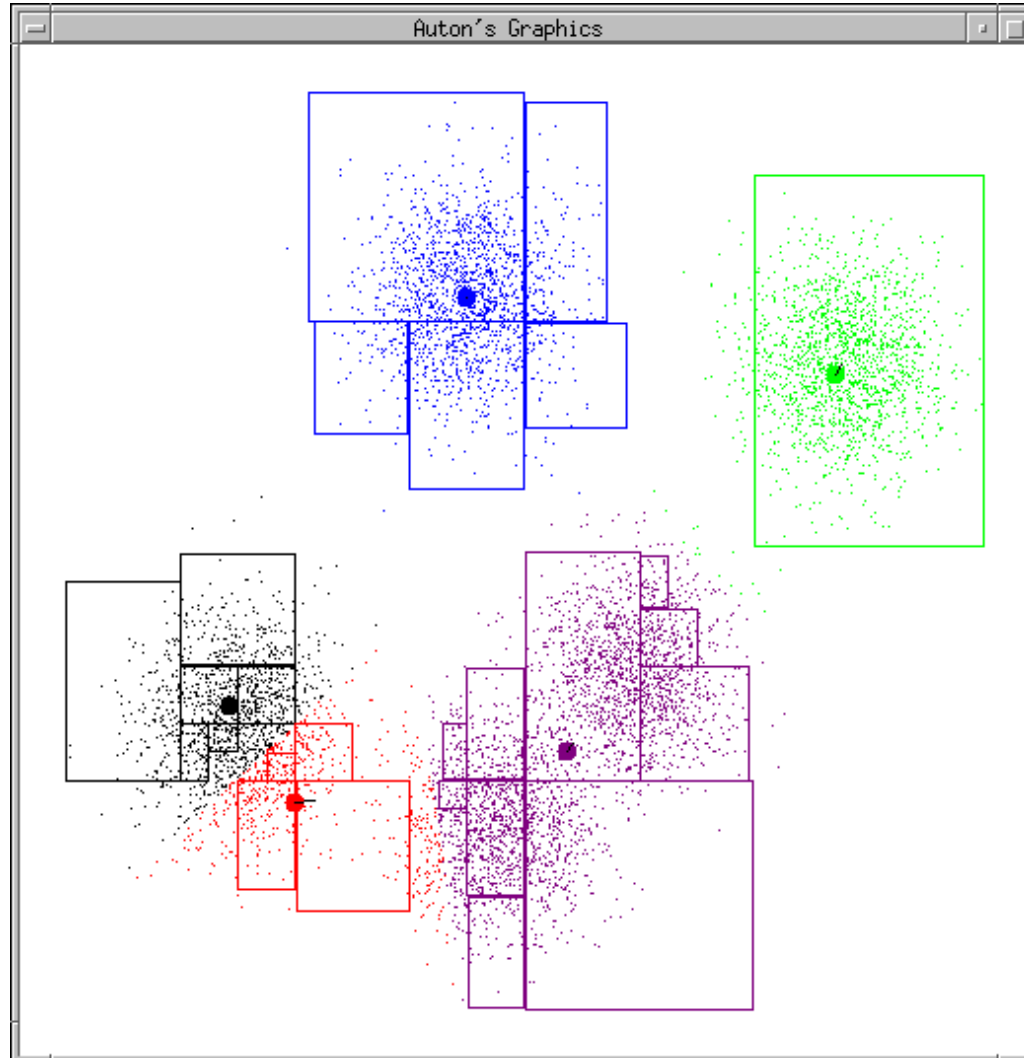
K-means continues...



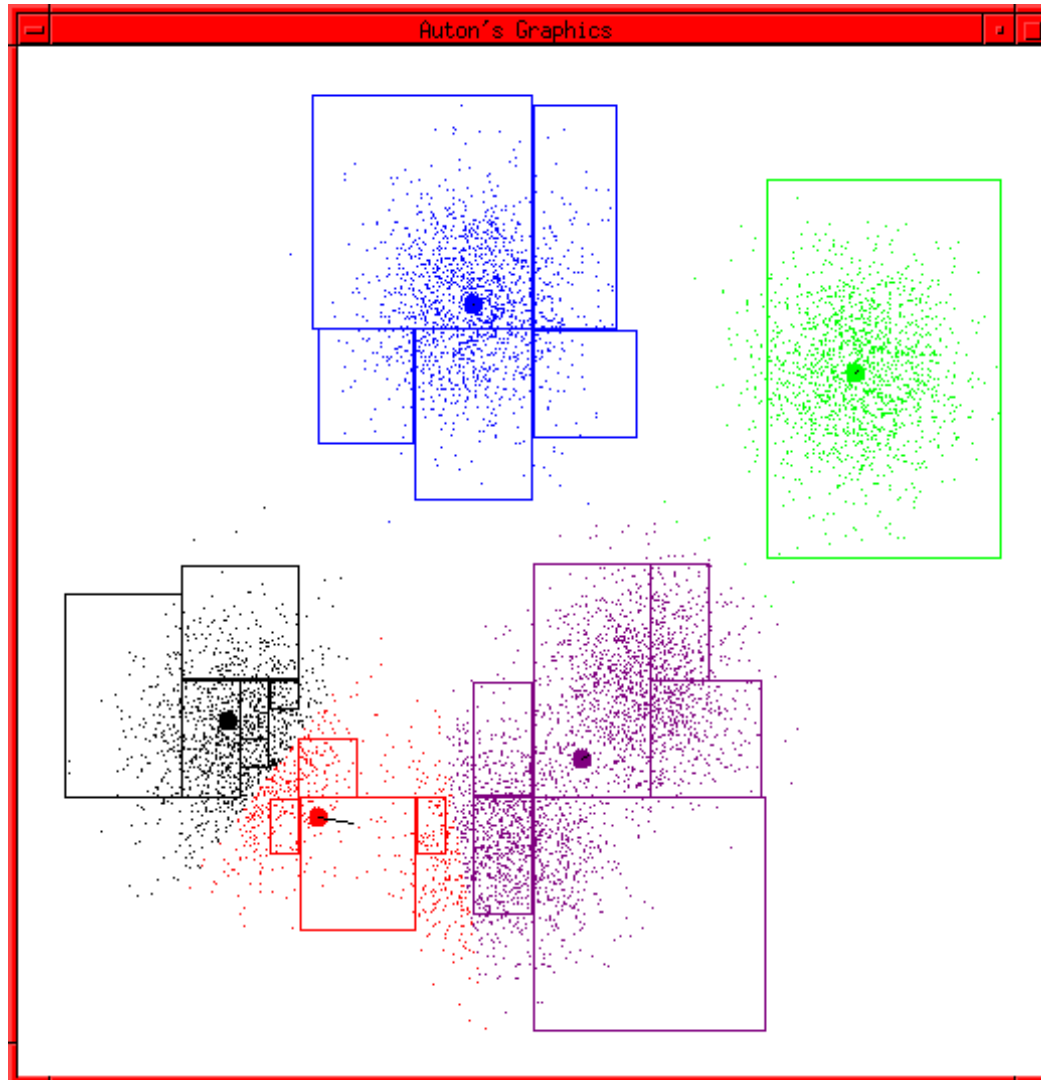
K-means continues...



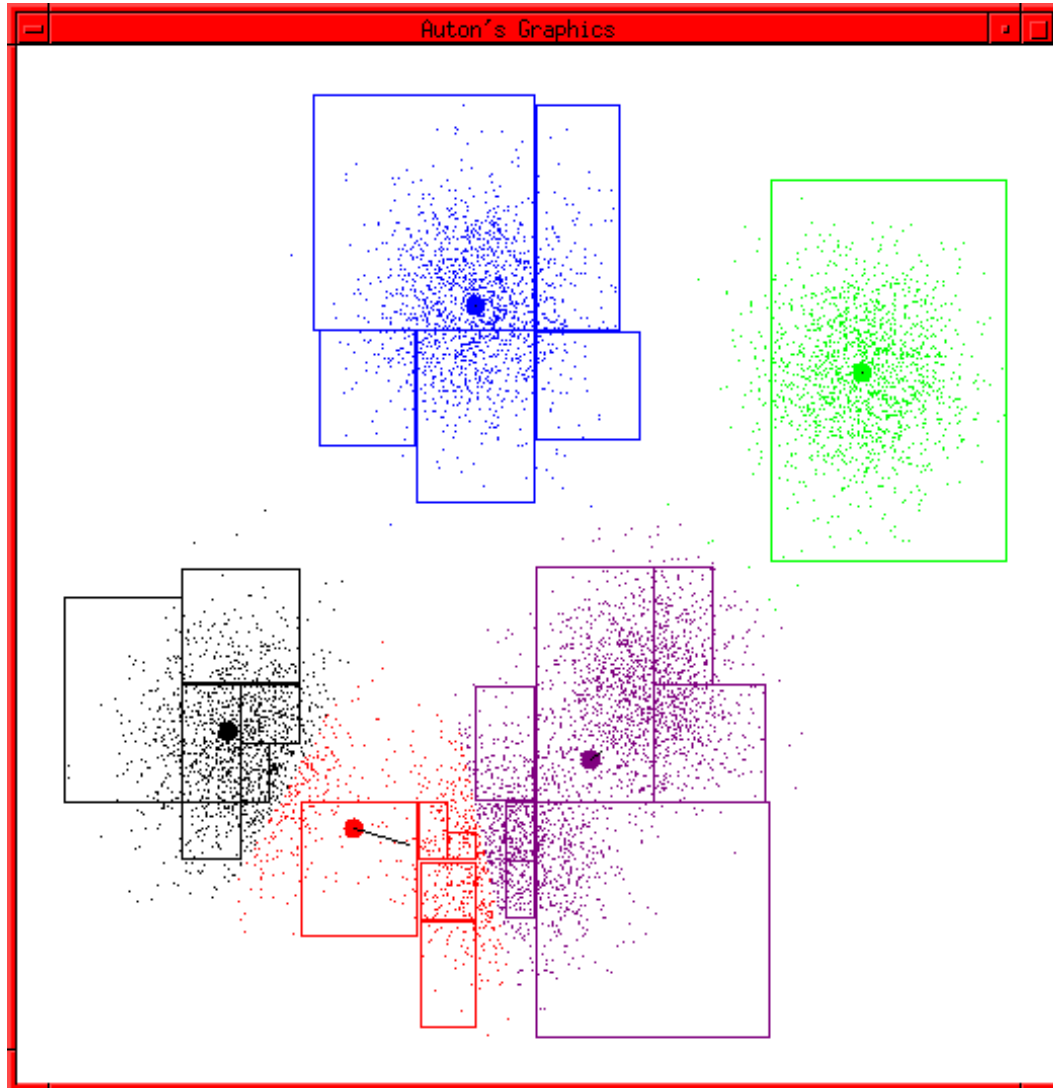
K-means continues...



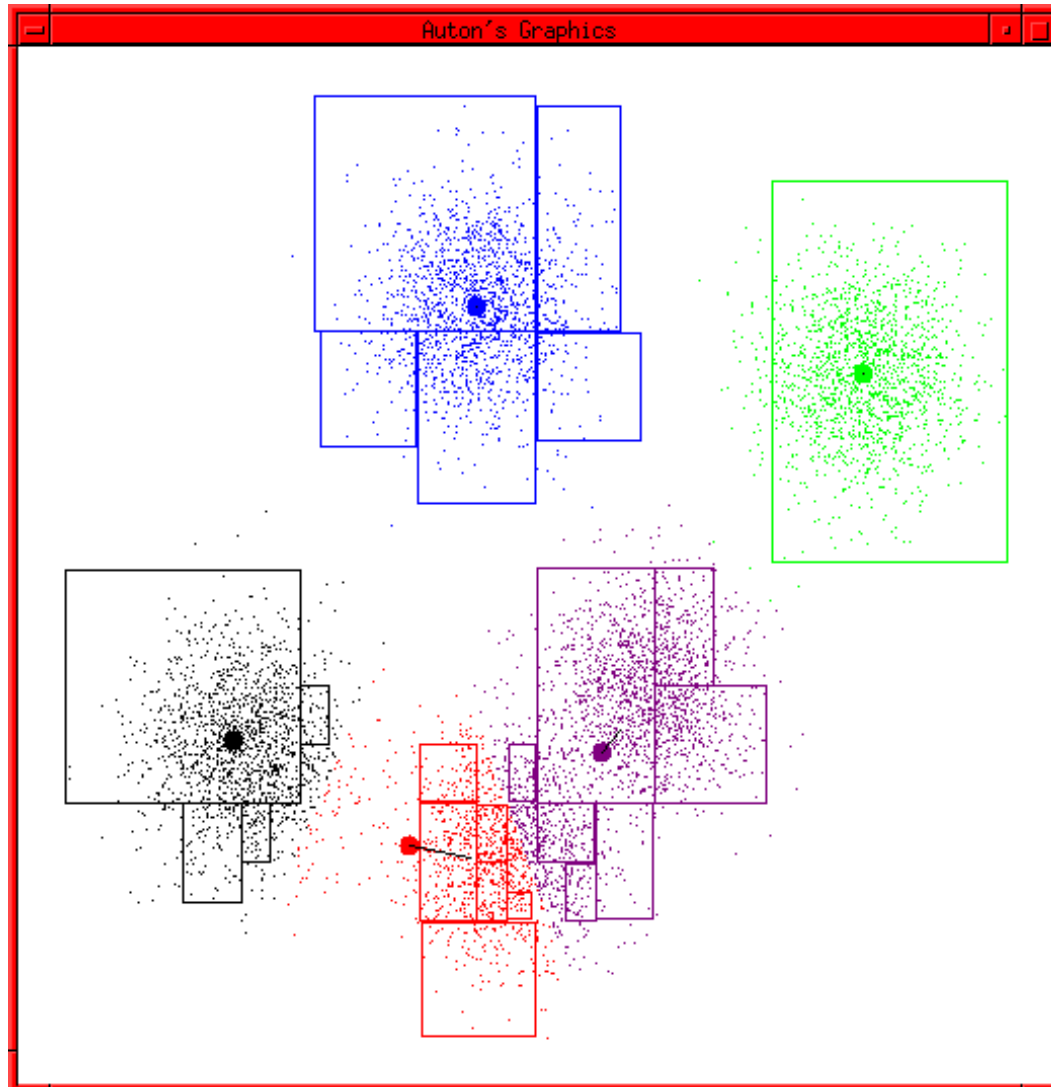
K-means continues...



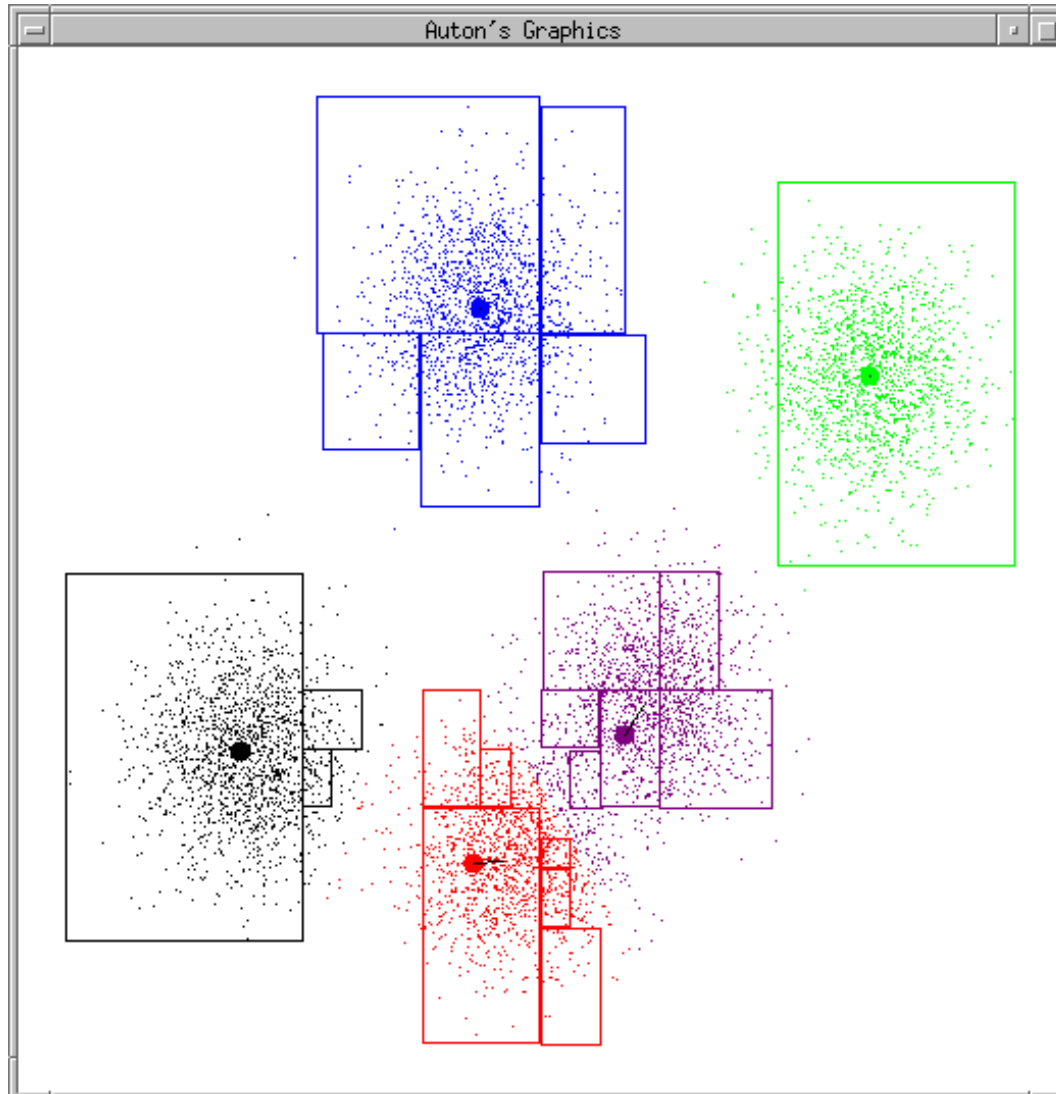
K-means continues...



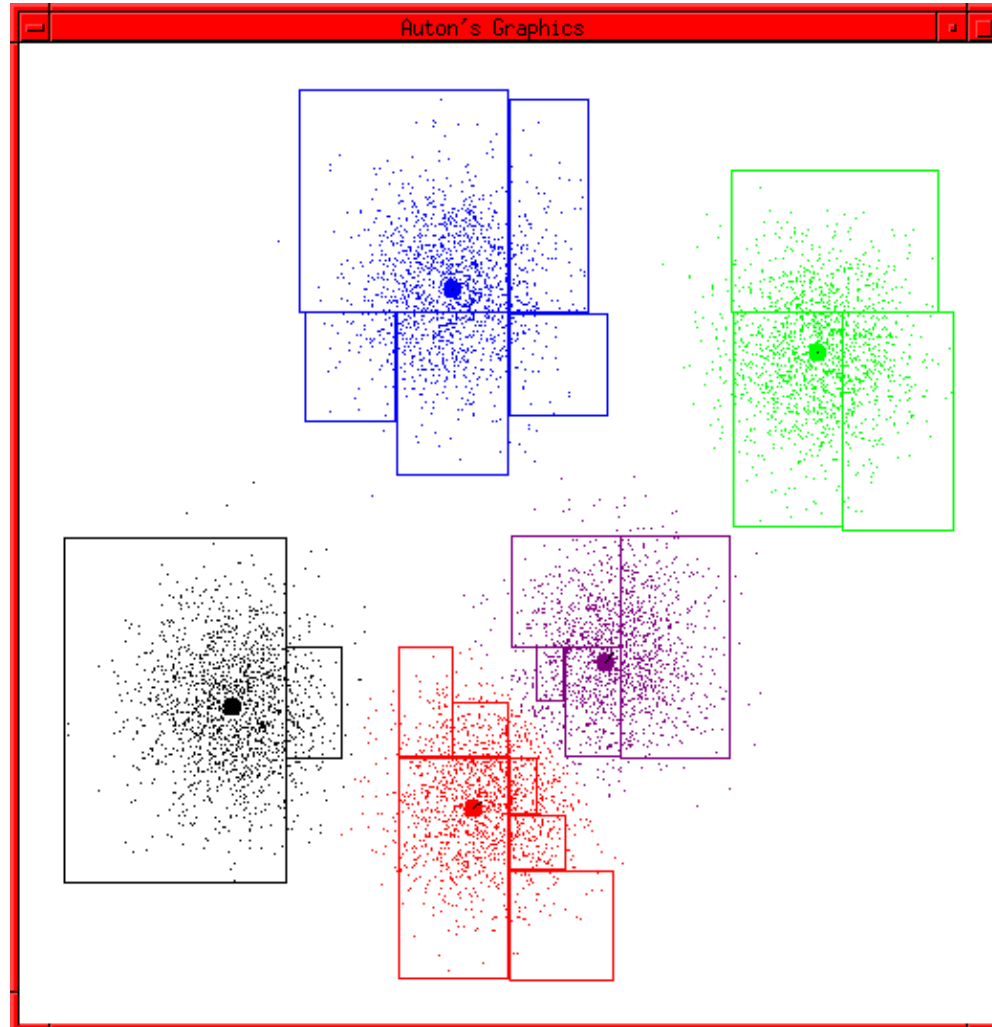
K-means continues...



K-means continues...



K-means terminates





Image



Clusters on color

K-means clustering of RGB (3 value) pixel color intensities, $K = 11$ segments
(courtesy of David Forsyth, UC Berkeley)

Issues in K-means clustering

- Simple, but useful
 - tends to select compact “isotropic” cluster shapes
 - can be useful for initializing more complex methods
 - many algorithmic variations on the basic theme
- Choice of distance measure
 - Euclidean distance
 - Weighted Euclidean distance
 - Many others possible
- Selection of K
 - “screen diagram” - plot SSE versus K, look for knee
 - Limitation: may not be any clear K value

**ΟΙΚΟΝΟΜΙΚΟ
ΠΑΝΕΠΙΣΤΗΜΙΟ
ΑΘΗΝΩΝ**



**ATHENS UNIVERSITY
OF ECONOMICS
AND BUSINESS**

Expectation Maximization (EM)

Μάθημα: Εξόρυξη γνώσης από Βάσεις Δεδομένων και τον Παγκόσμιο Ιστό

Ενότητα # 4: Unsupervised Learning (Clustering)

Διδάσκων: Μιχάλης Βαζιργιάννης

Τμήμα: Προπτυχιακό Πρόγραμμα Σπουδών “Πληροφορικής”

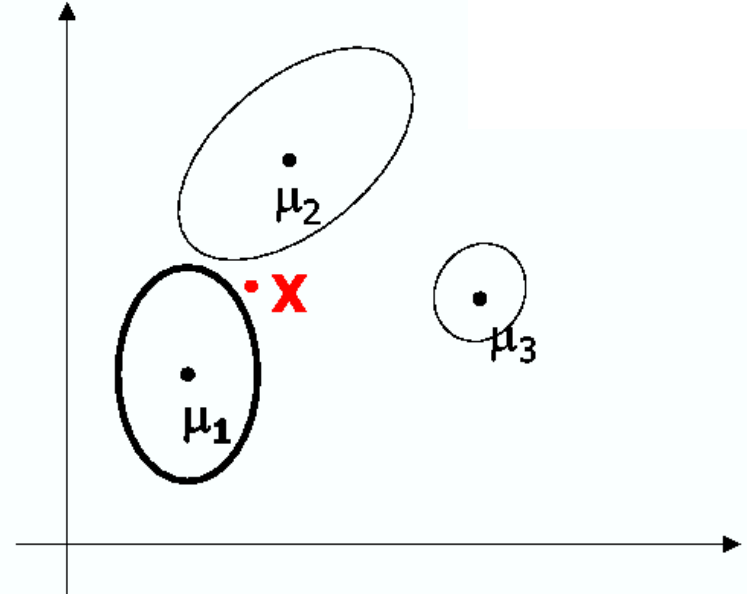
Expectation Maximization

- (EM) algorithm is an iterative method for finding maximum likelihood estimates of parameters in statistical models that depend on unobserved latent variables.
- Assume X the data observed - We assume the data are produced by K different classes/processes represented by respective weights there fore w_k

$$f(x) = \sum_k w_k f_k(x|q_k)$$

Gaussian Mixture Models (GMM)

- Assume the the components are normal distributions $N(\mu_k, \Sigma_k)$
 - often assume diagonal covariance: $\Sigma_{jj} = \sigma_j^2, \Sigma_{i \neq j} = 0$
 - or sometimes even simpler: $\Sigma_{jj} = \sigma^2, \Sigma_{i \neq j} = 0$
- $f(x) = \sum_{k=1 \dots K} w_k f_k(x; \theta_k)$ with $\theta_k = \langle \mu_k, \Sigma_k \rangle$ or $\langle \mu_k, \sigma_k \rangle$
- generative model:
 - - randomly choose a component
 - selected with probability w_k
 - - generate $x \sim N(\mu_k, \sigma_k)$
 - - note: μ_k & σ_k both d-dim vectors



Learning Mixture Models from Data

- **Score function** Log-likelihood $L(\theta)$
 - $L(\theta) = \log p(X|\theta) = \log \sum_H p(X,H|\theta)$
 - H = hidden variables (cluster memberships of each x)
 - $L(\theta)$ cannot be optimized directly
- **EM Procedure**
 - General technique for maximizing log-likelihood with missing data
 - For mixtures
 - *E-step*: compute “memberships” $p(k | x) = w_k f_k(x; \theta_k) / f(x)$
 - *M-step*: pick a new θ to max expected data log-likelihood
 - Iterate: guaranteed to climb to (local) maximum of $L(\theta)$

Expectation maximization (EM)

The [Expectation-maximization algorithm](#) computes missing memberships of data points in a chosen distribution model.

- **Expectation step**

- initial guesses for the parameters in our mixture model,
- compute "partial membership" of each data point in each constituent distribution.
- By calculating expectation for the membership variables of each data point.

- **Example.**

- Data set resulting from a sum of two Gaussian distributions.

$$P(\mathbf{x}_i) = (1 - f)N(\mathbf{x}_i | m_1, S) + fN(\mathbf{x}_i | m_2, S)$$

- f is the mixing coefficient in $(0,1]$, assume σ is known and constant.
- For each data point i , compute a membership value for each of the two Gaussians

$$y_{1,i}(\mathbf{x}_i) = \frac{(1 - f)N(\mathbf{x}_i | m_1, S)}{(1 - f)N(\mathbf{x}_i | m_1, S) + fN(\mathbf{x}_i | m_2, S)}$$

- and similarly for $y_{2,i}$

Expectation maximization (EM)

- **The maximization step**
- With expectation values for group membership
- - Re-compute estimates of distribution parameters.

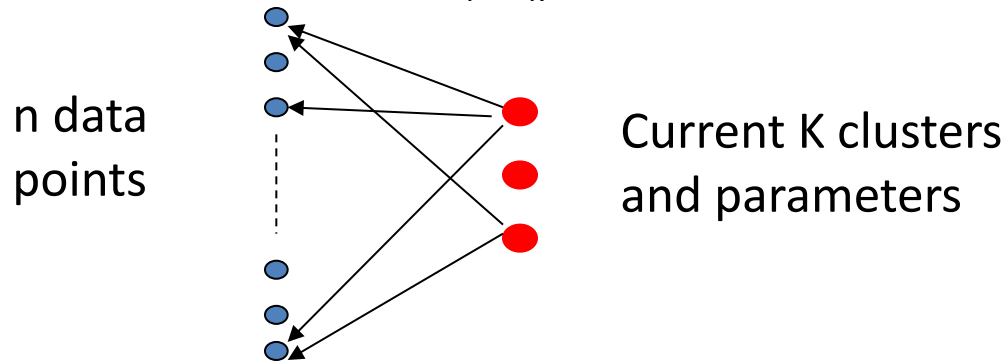
$$\left\{ \begin{array}{l} f = \frac{\sum_i y_{i,2}}{N} \\ \mu_1 = \frac{\sum_i y_{i,1} x_i}{\sum_i y_{i,1}} \end{array} \right.$$

– N is the total number of data points.

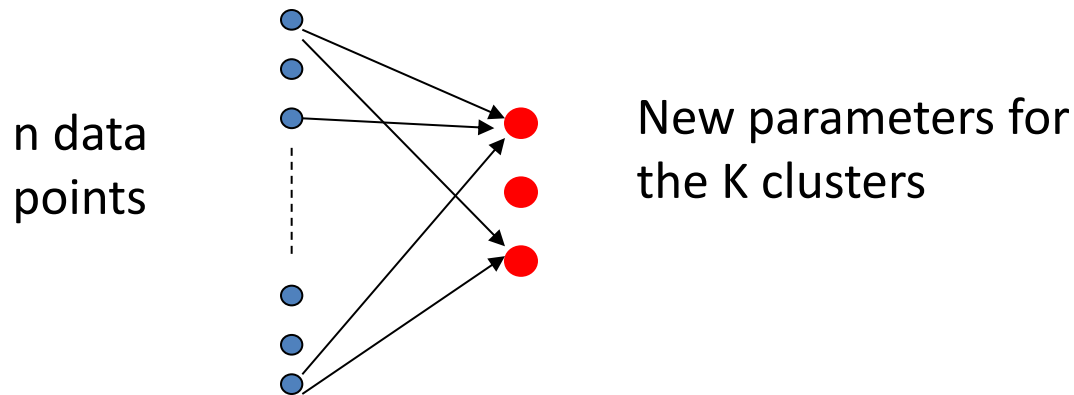
- - back to the Expectation step: Re-compute new membership values.
- - repeated until change in the mixture model parameters below threshold

The EM Step

E step: Compute memberships $p(x_i/\vartheta_k)$



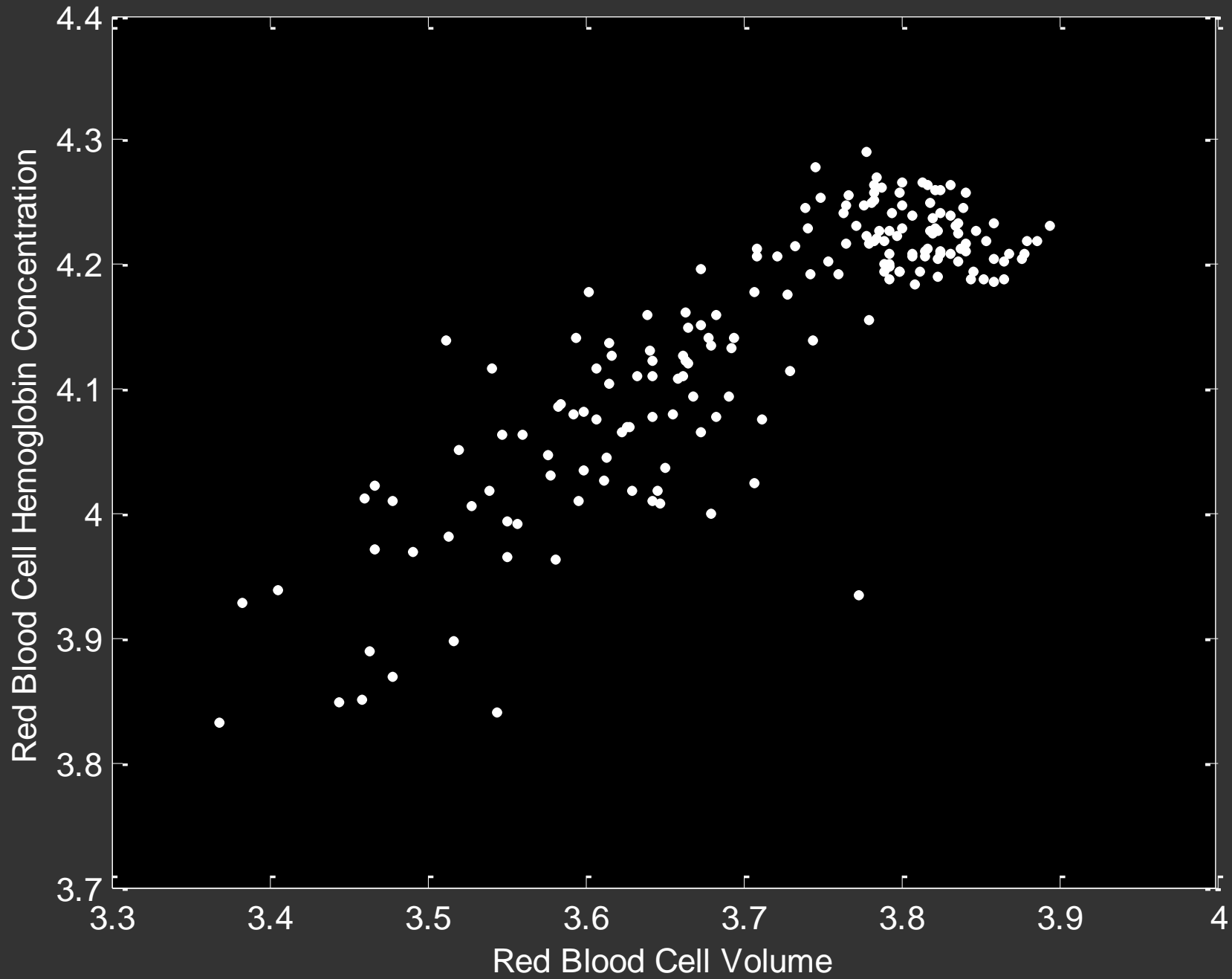
M step: Compute q_i , given n data points and memberships



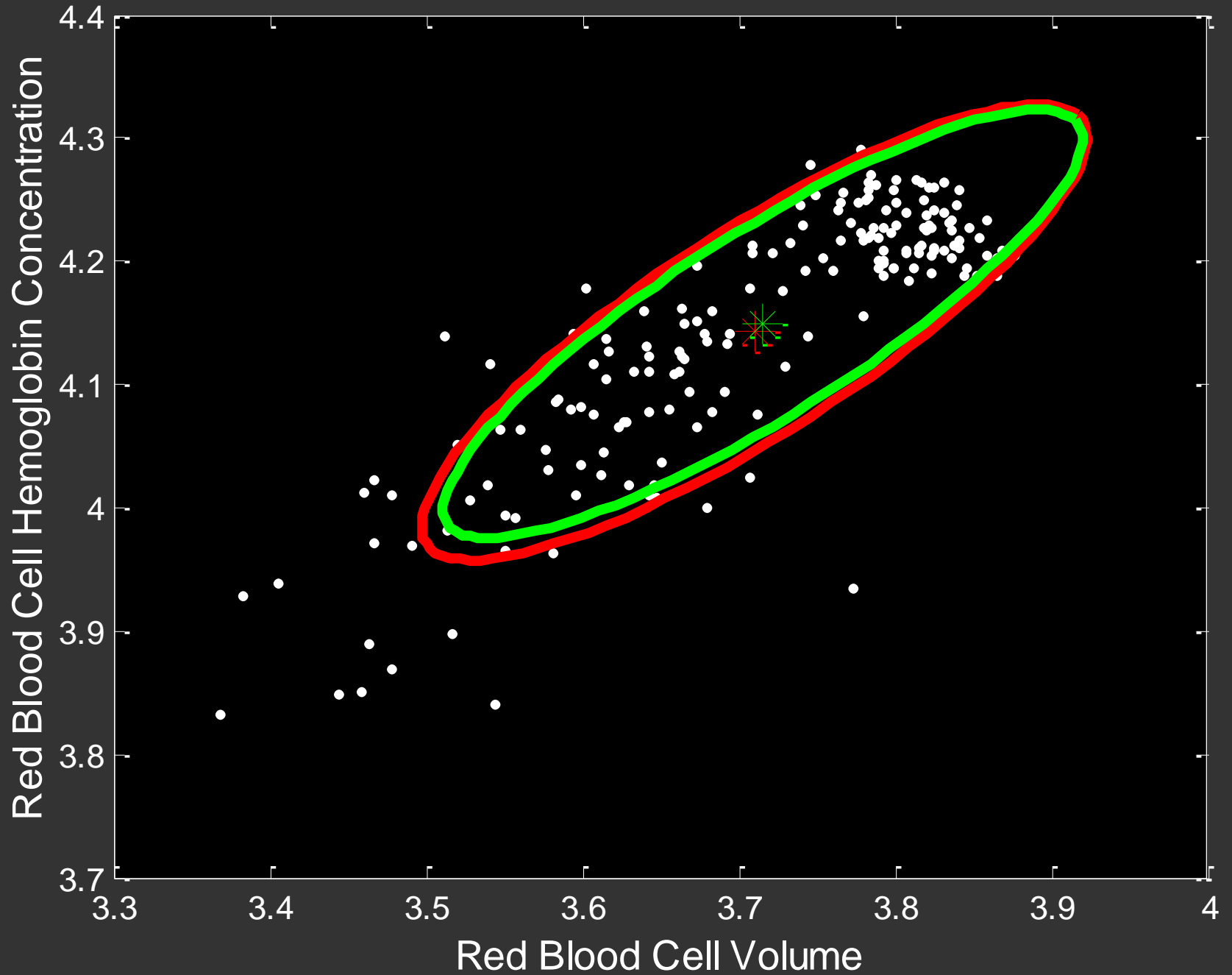
Comments on Mixtures and EM Learning

- - Complexity of each EM iteration
 - Depends on the probabilistic model being used
 - e.g., for Gaussians, E-step: $O(nK)$, M-step: $O(Knp^2)$
 - Sometimes E or M-step is not closed form
 - => can requires numerical methods at each iteration
- - K-means interpretation
 - Gaussian mixtures with isotropic (diagonal, equi-variance) Σ_k 's
 - Approximate the E-step by choosing most likely cluster (instead of using membership probabilities)

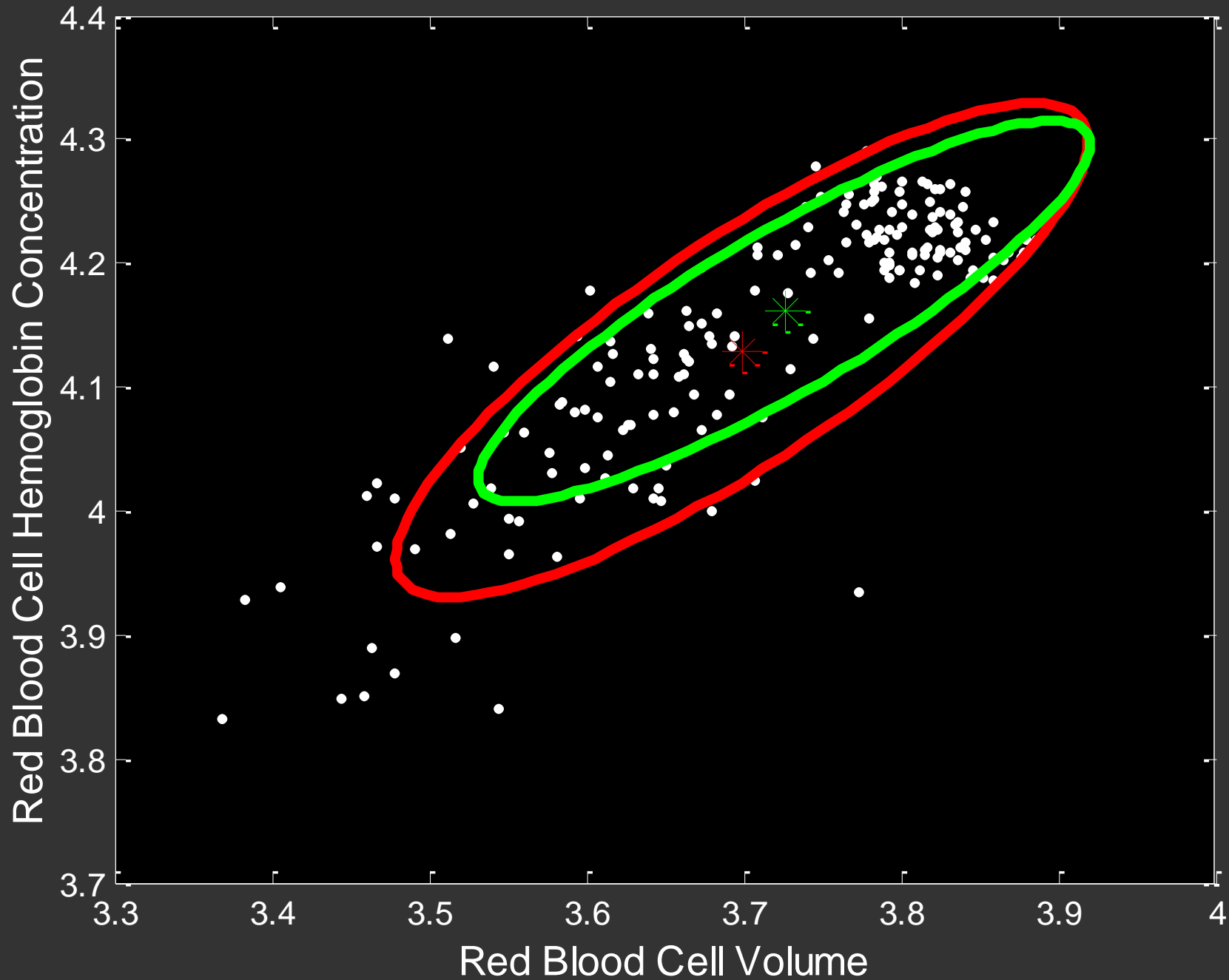
ANEMIA PATIENTS AND CONTROLS



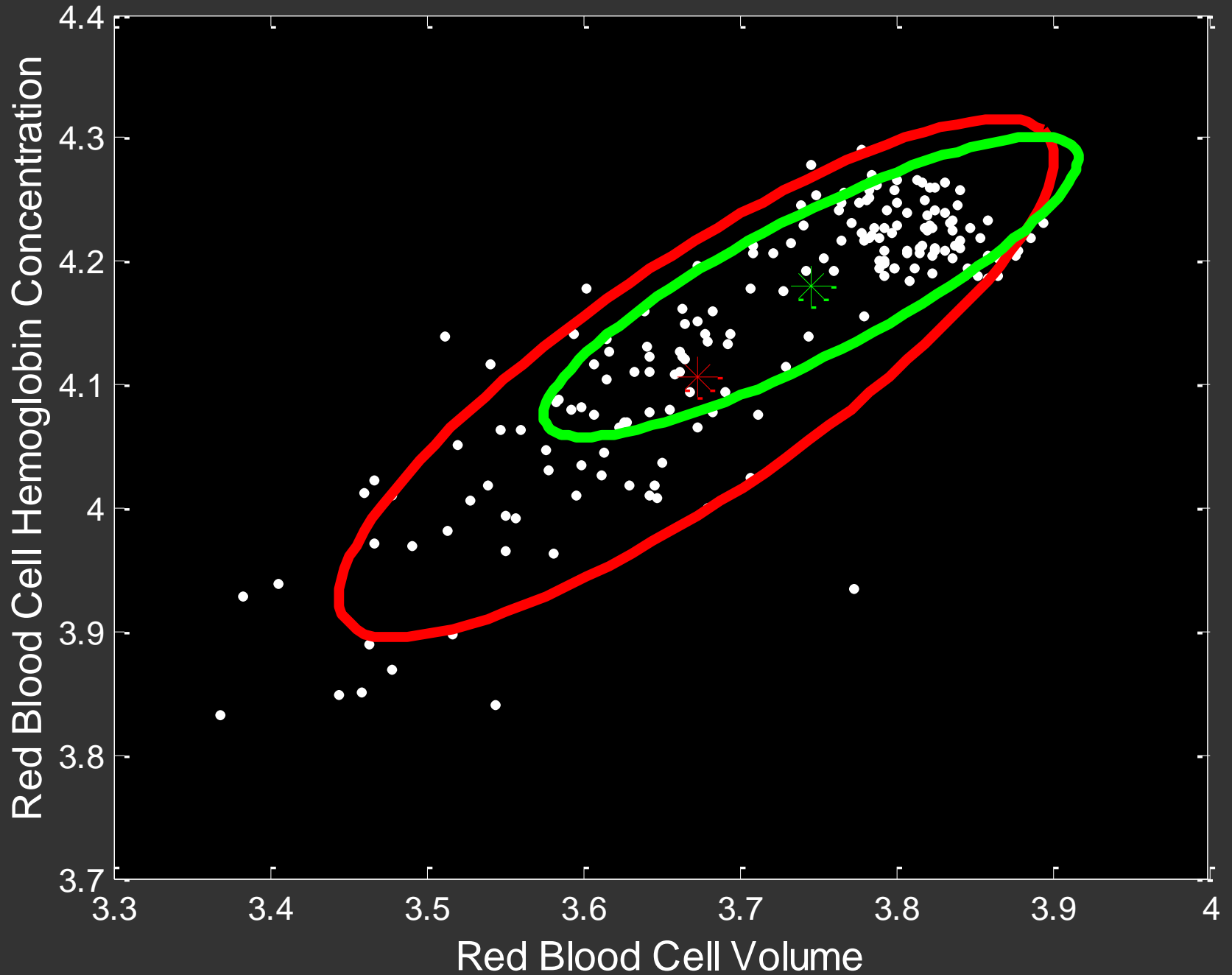
EM ITERATION 1



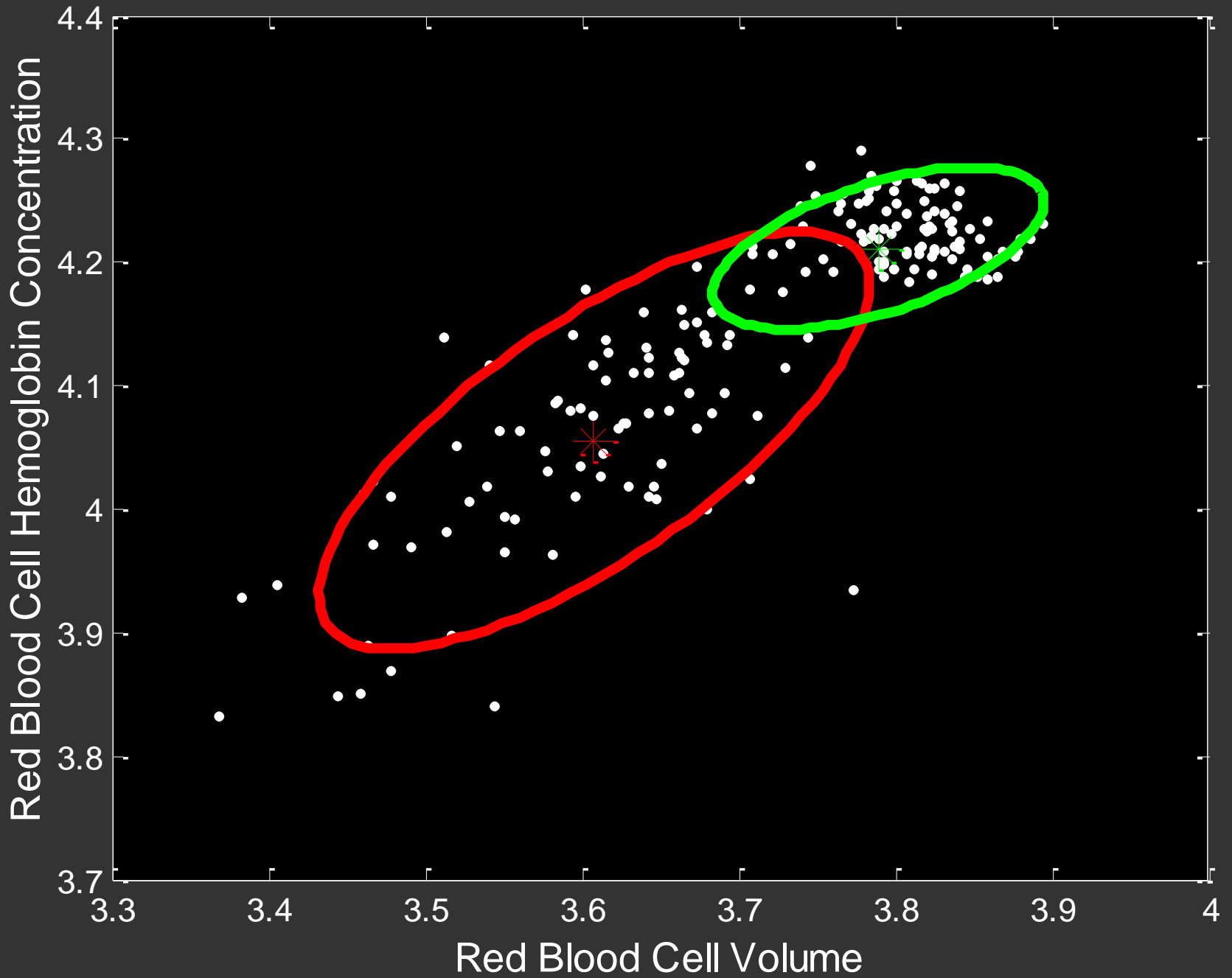
EM ITERATION 3



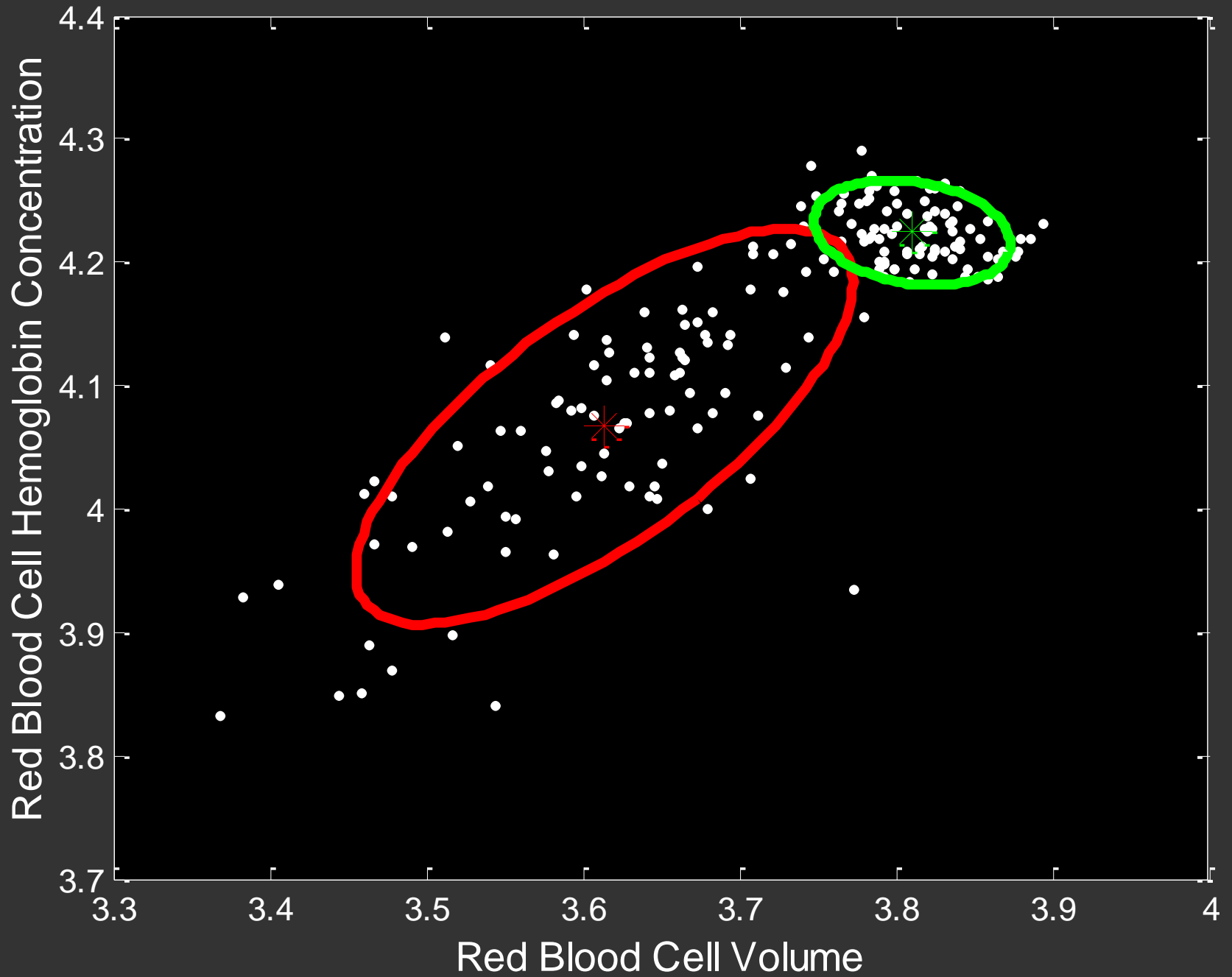
EM ITERATION 5



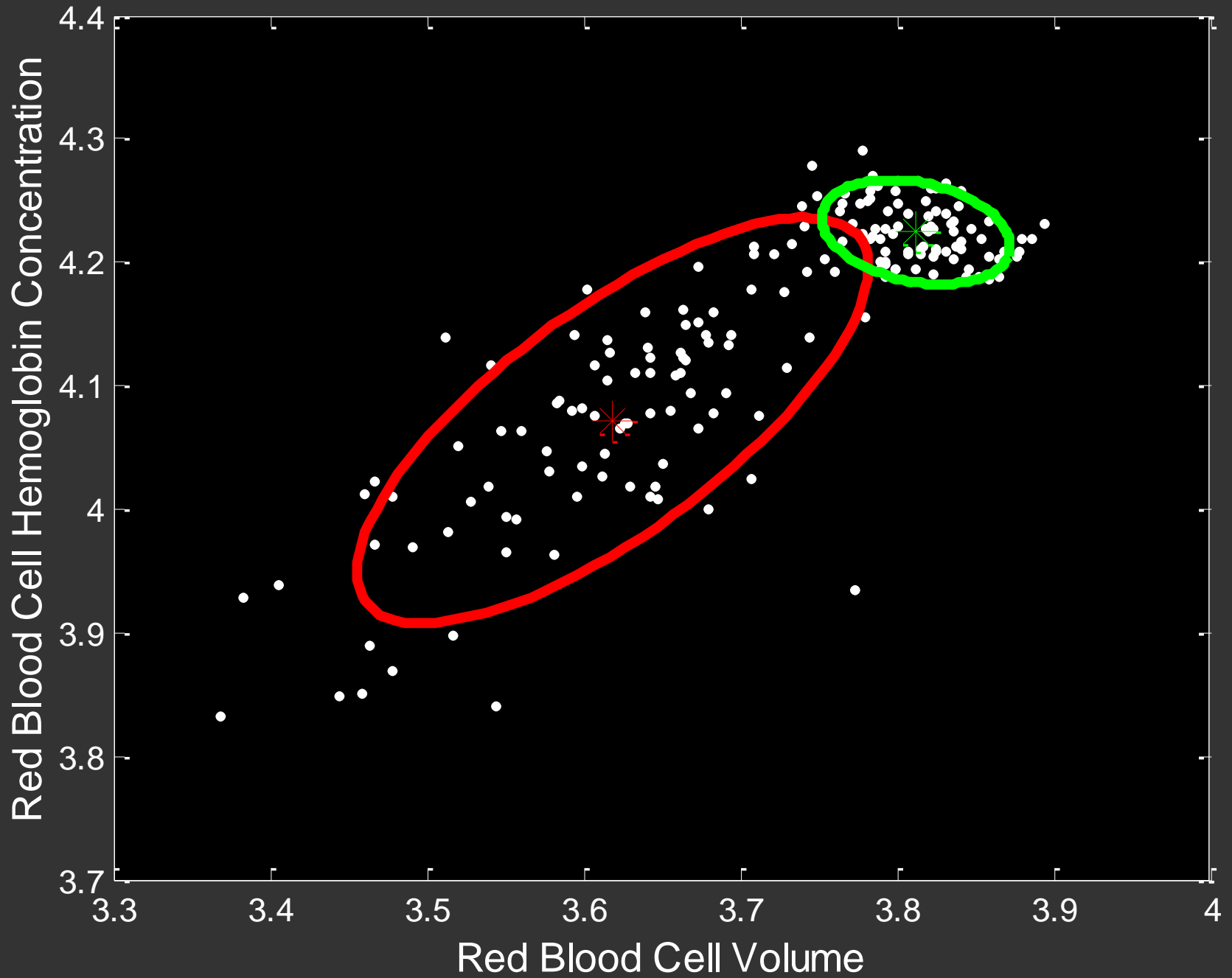
EM ITERATION 10



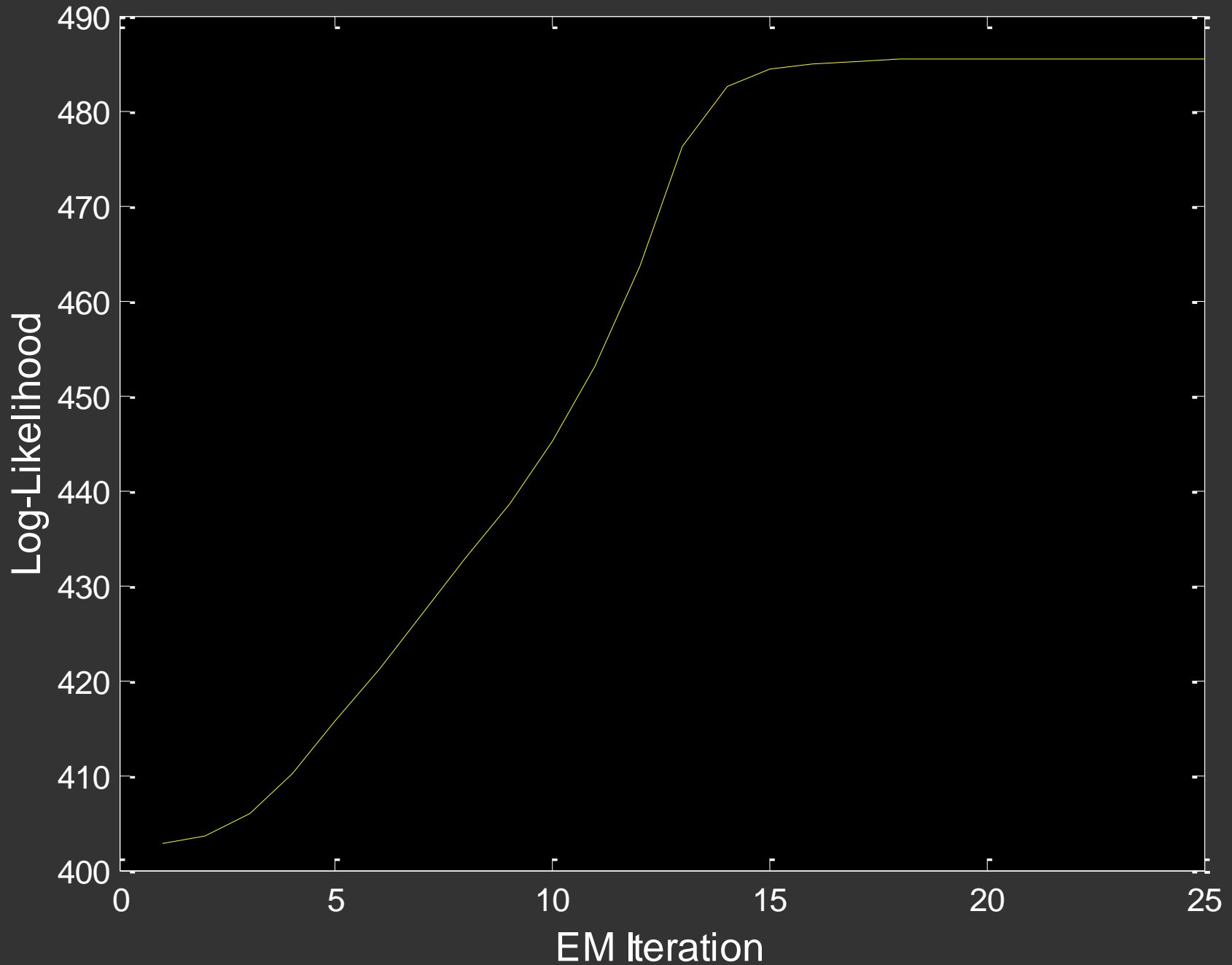
EM ITERATION 15



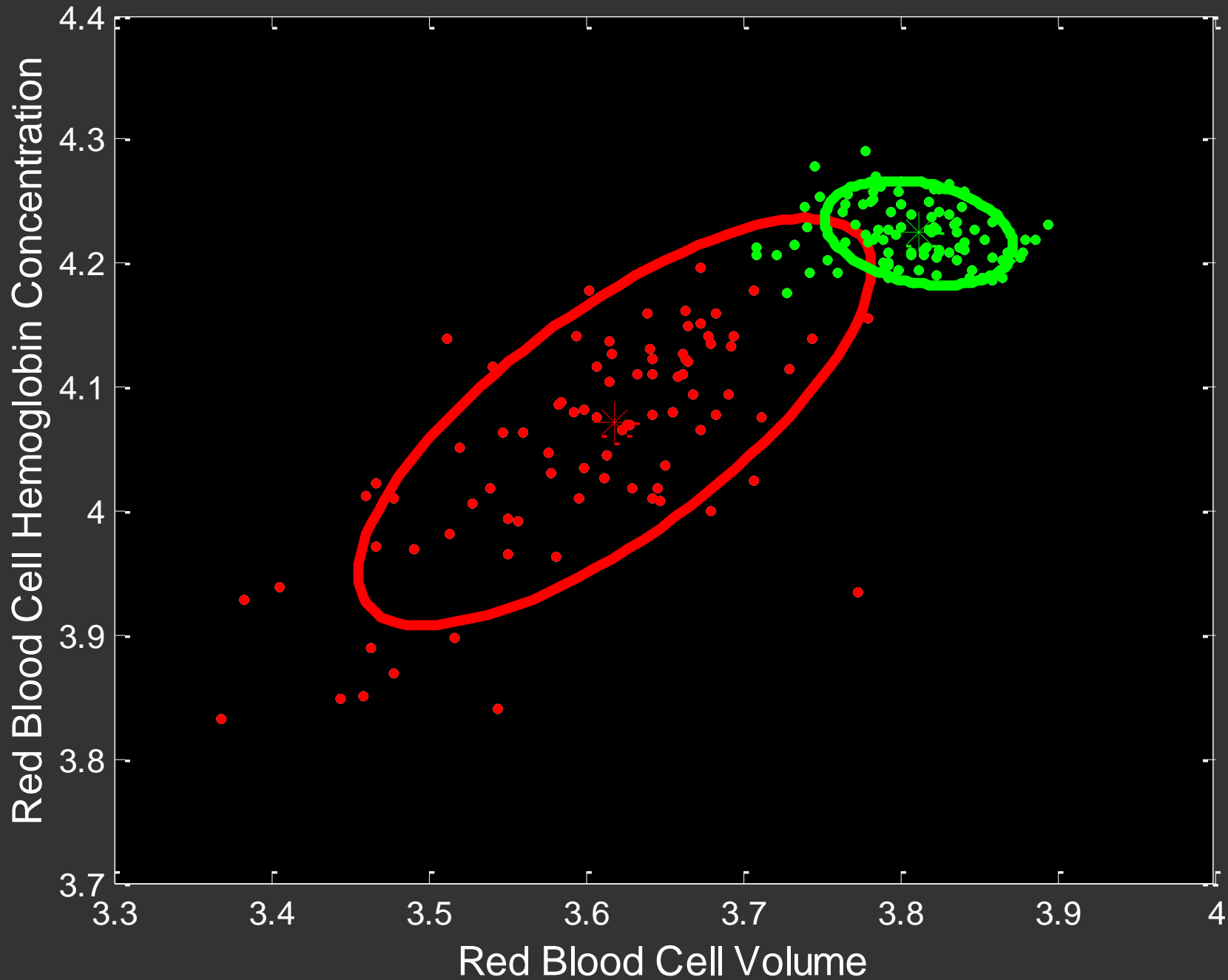
EM ITERATION 25



LOG-LIKELIHOOD AS A FUNCTION OF EM ITERATIONS

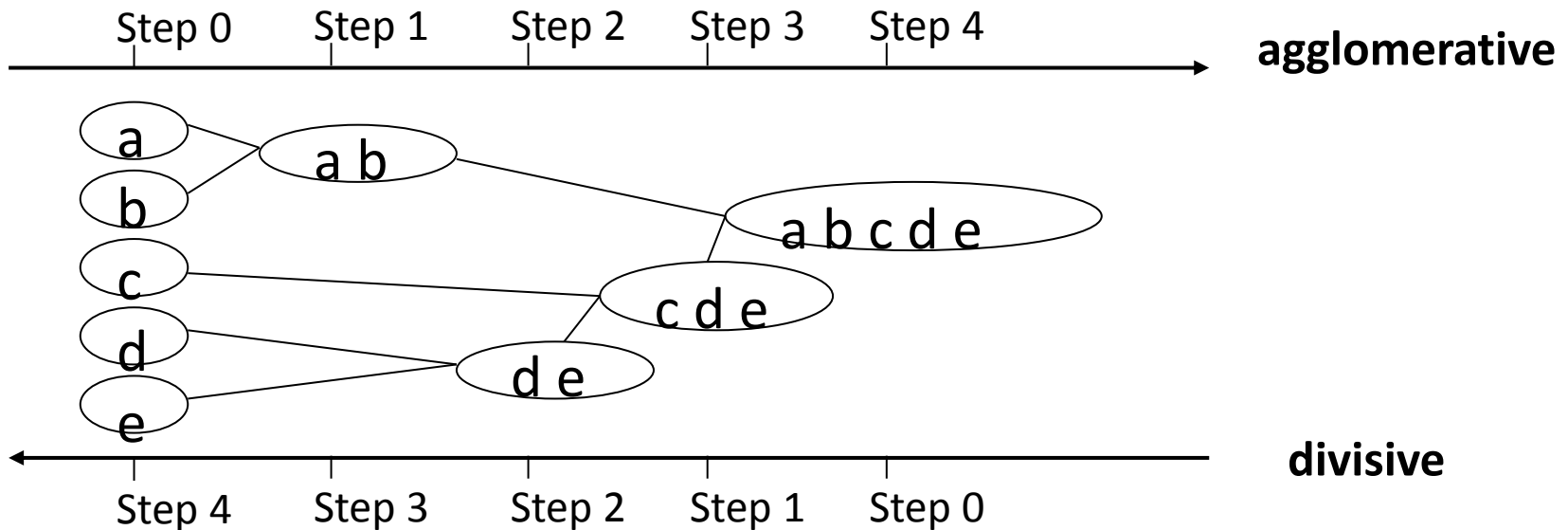


ANEMIA DATA WITH LABELS



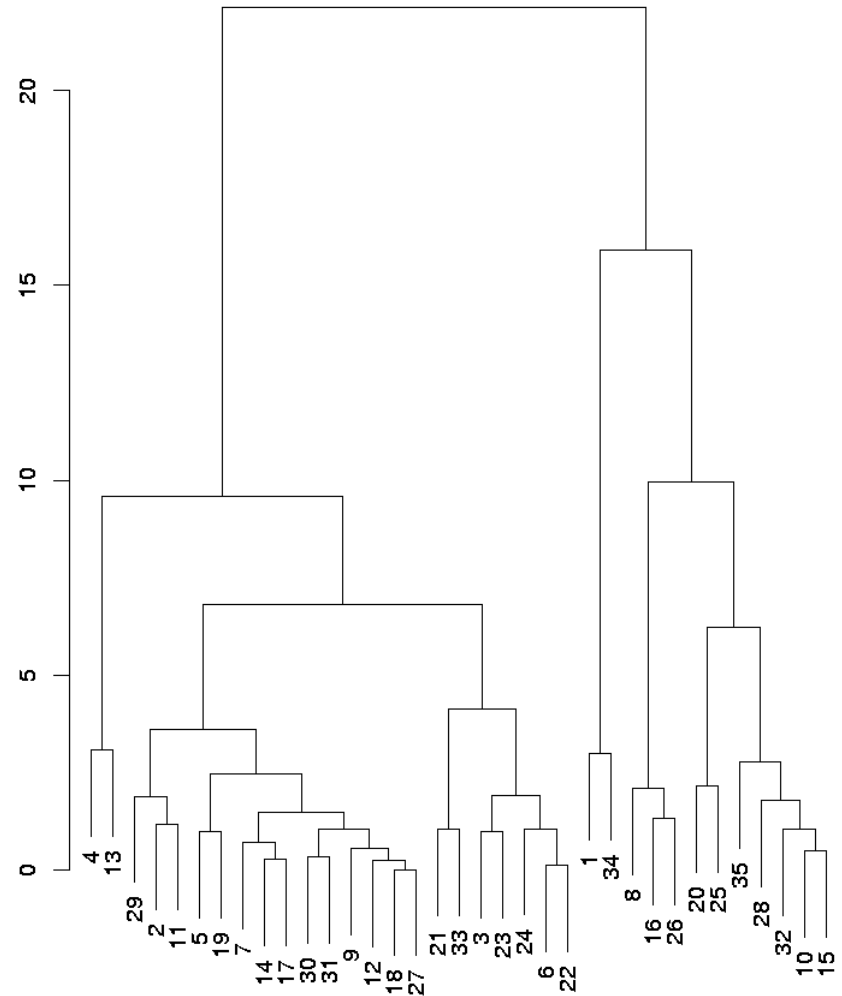
Hierarchical Clustering

- **Two basic approaches:**
 - merging smaller clusters into larger ones (agglomerative),
 - splitting larger clusters (divisive)
- **visualize both via “dendograms”**
 - shows nesting structure
 - merges or splits = tree nodes



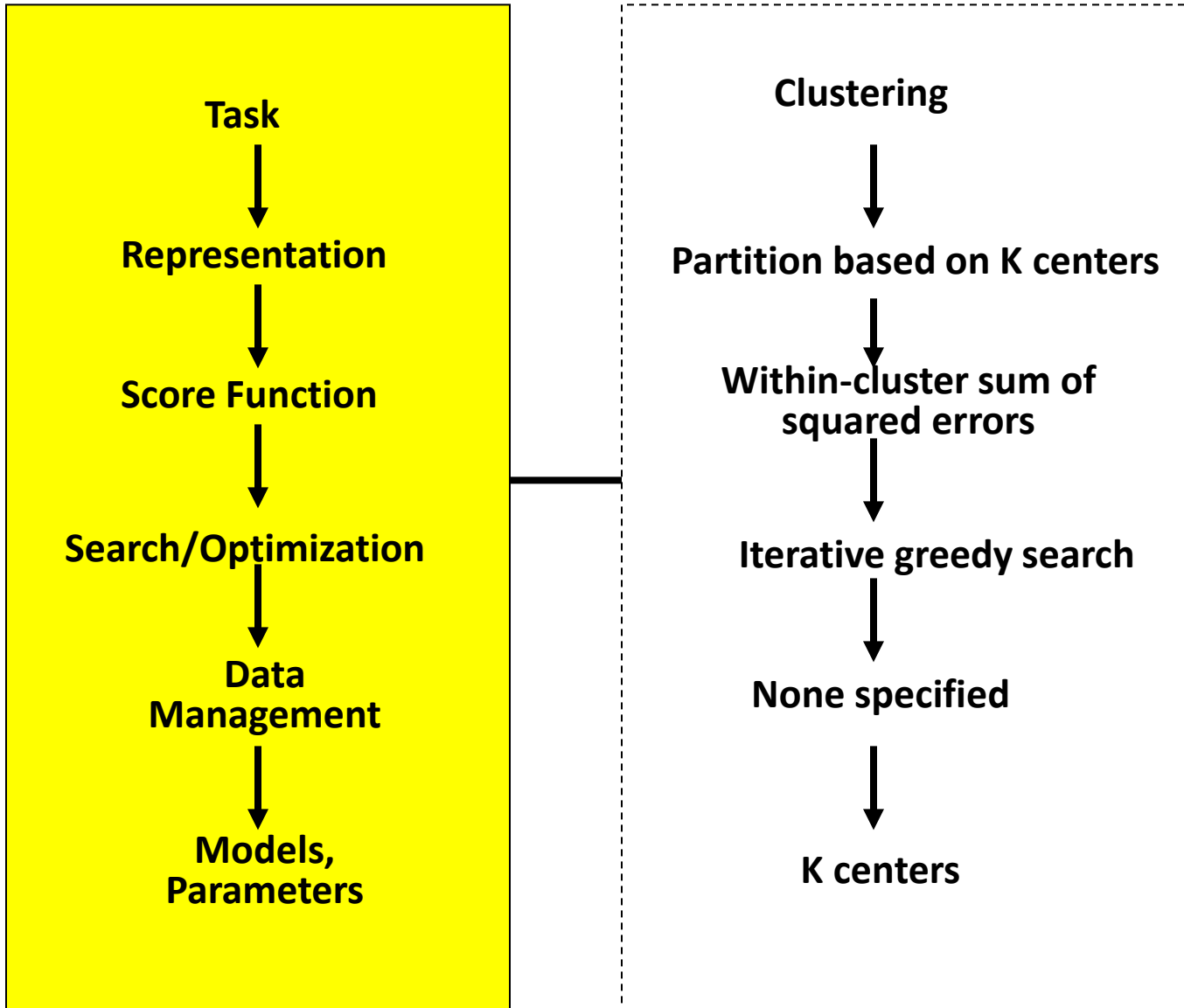
Hierarchical Clustering: Complexity

- **Quadratic algorithms**
- **Running time** can be improved using
 - sampling
 - [Guha et al, SIGMOD 1998]
 - [Kollios et al, ICDE 2001]
 - or using the triangle inequality (when it holds)

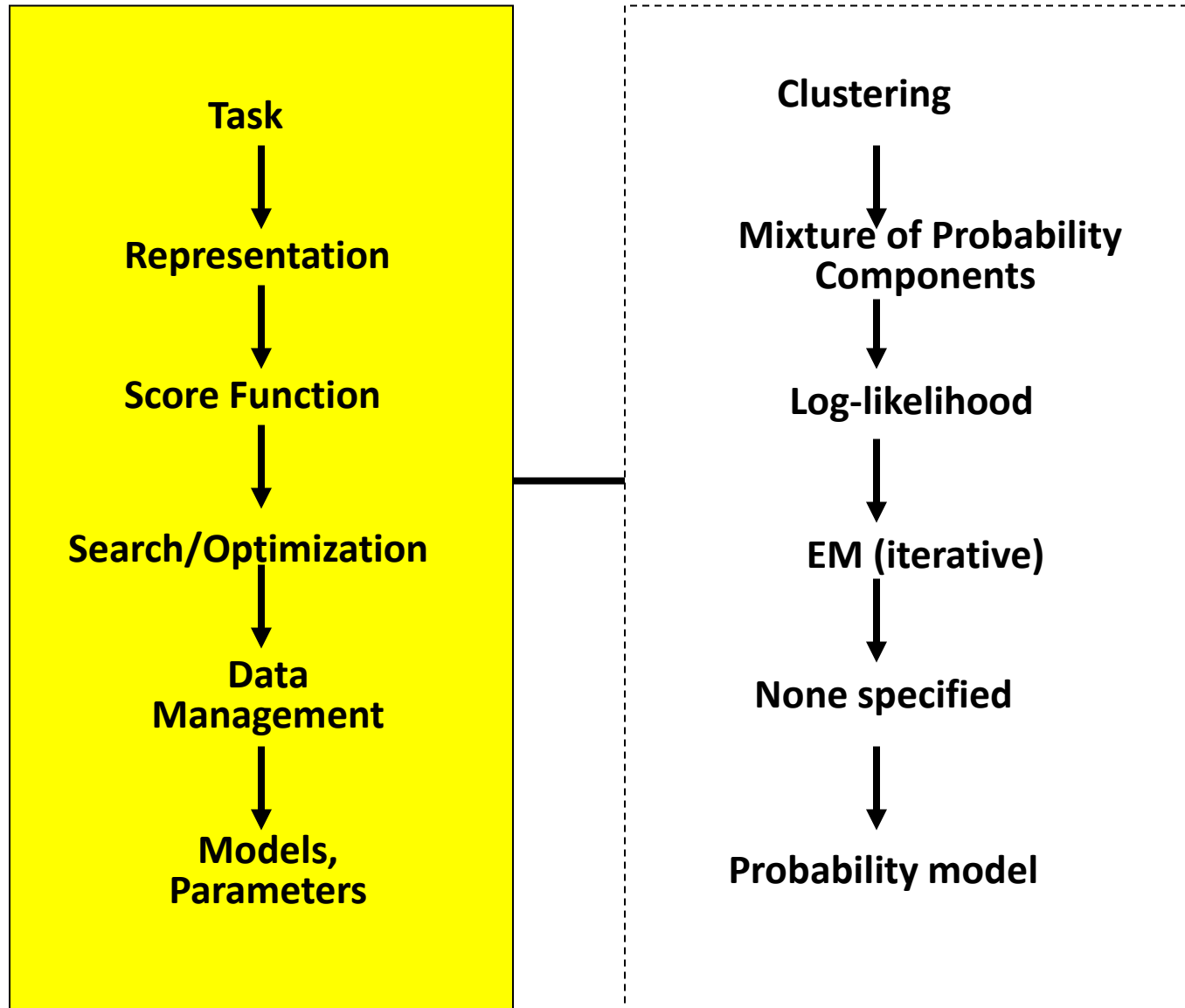


*based on slides by Padhraic Smyth UC, Irvine

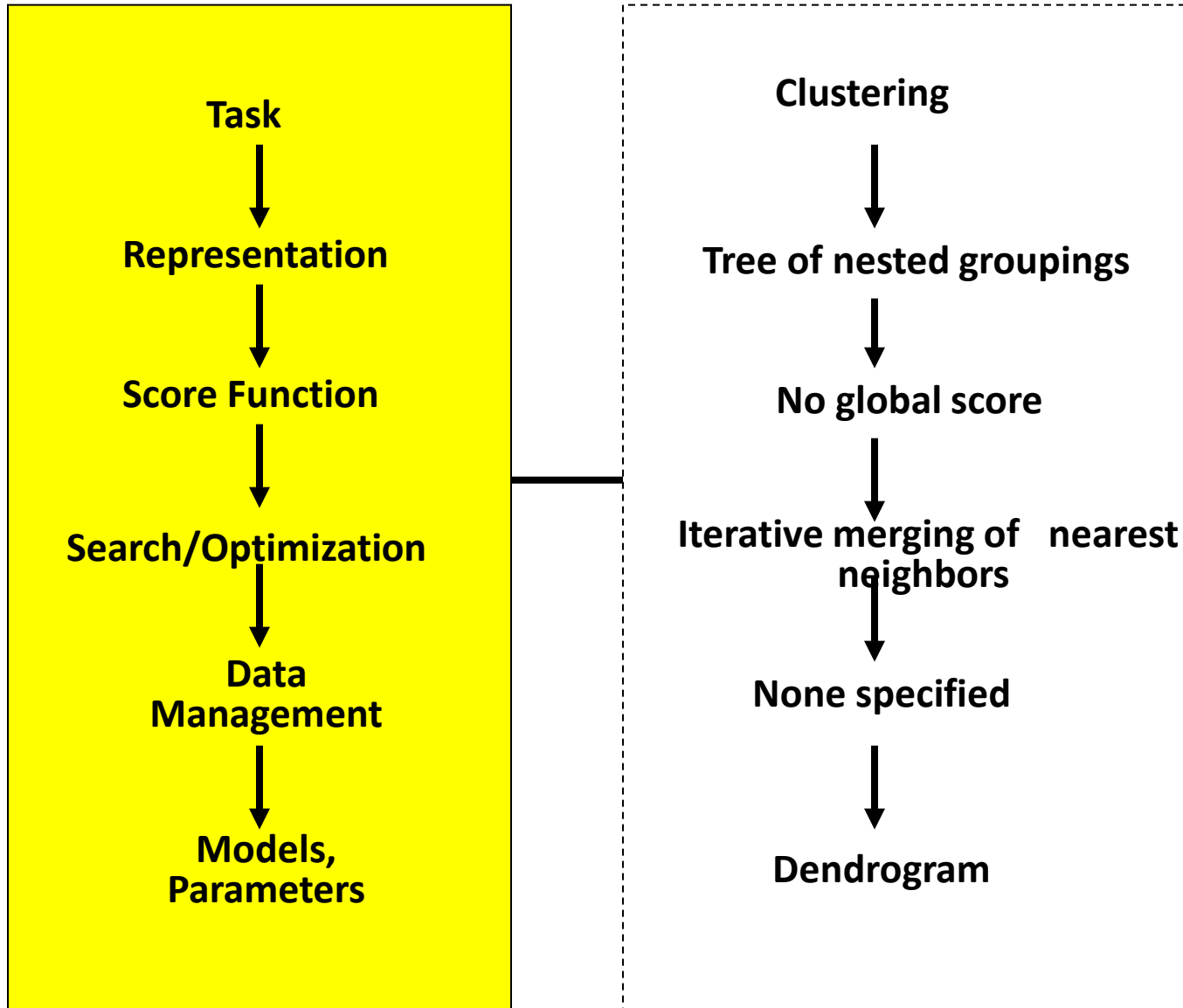
K-Means Clustering



Probabilistic Model-Based Clustering



Single-Link Hierarchical Clustering



**ΟΙΚΟΝΟΜΙΚΟ
ΠΑΝΕΠΙΣΤΗΜΙΟ
ΑΘΗΝΩΝ**



**ATHENS UNIVERSITY
OF ECONOMICS
AND BUSINESS**

Association Rules

Μάθημα: Εξόρυξη γνώσης από Βάσεις Δεδομένων και τον Παγκόσμιο Ιστό

Ενότητα # 5: Unsupervised Learning (Clustering)

Διδάσκων: Μιχάλης Βαζιργιάννης

Τμήμα: Προπτυχιακό Πρόγραμμα Σπουδών “Πληροφορικής”

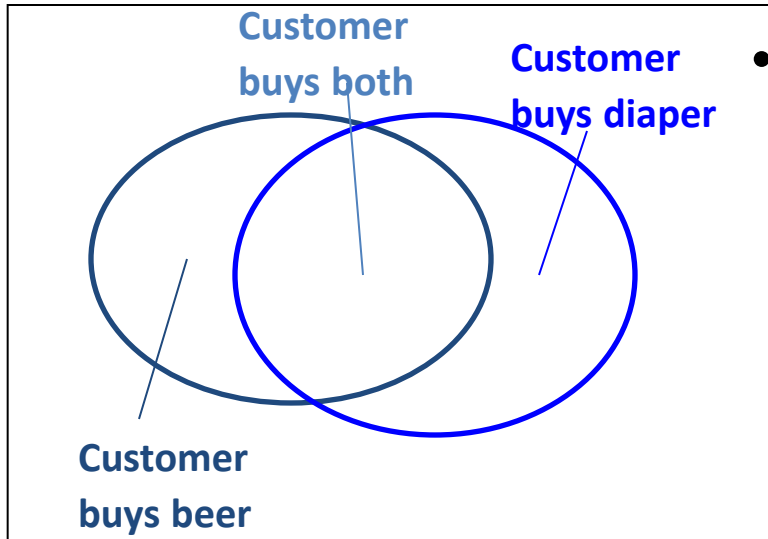
What Is Association Mining?

- Association rule mining:
 - Finding frequent patterns, associations, correlations, or causal structures among sets of items or objects in transaction databases, relational databases, and other information repositories.
- Applications:
 - Basket data analysis, cross-marketing, catalog design, loss-leader analysis, clustering, classification, etc.
- Examples.
 - Rule form: “Body \rightarrow Head [support, confidence]”.
 - $\text{buys}(x, \text{“diapers”}) \rightarrow \text{buys}(x, \text{“beers”})$ [0.5%, 60%]
 - $\text{major}(x, \text{“CS”}) \wedge \text{takes}(x, \text{“DB”}) \rightarrow \text{grade}(x, \text{“A”})$ [1%, 75%]

Association Rule: Basic Concepts

- Given: (1) database of transactions, (2) each transaction is a list of items (purchased by a customer in a visit)
- Find: all rules that correlate the presence of one set of items with that of another set of items
 - E.g., *98% of people who purchase tires and auto accessories also get automotive services done*
- Applications
 - * \Rightarrow *Maintenance Agreement* (What the store should do to boost Maintenance Agreement sales)
 - *Home Electronics* \Rightarrow * (What other products should the store stocks up?)
 - Attached mailing in direct marketing
 - Detecting “ping-pong”ing of patients, faulty “collisions”

Rule Measures: Support and Confidence



- Find all the rules $X \& Y \Rightarrow Z$ with minimum confidence and support
 - support, s , probability that a transaction contains $\{X \& Y \& Z\}$
 - confidence, c , conditional probability that a transaction having $\{X \& Y\}$ also contains Z

| Transaction ID | Items Bought |
|----------------|--------------|
| 2000 | A,B,C |
| 1000 | A,C |
| 4000 | A,D |
| 5000 | B,E,F |

Let minimum support 50%, and minimum confidence 50%, we have

- $A \Rightarrow C$ (50%, 66.6%)
- $C \Rightarrow A$ (50%, 100%)

Mining Association Rules—An Example

| Transaction ID | Items Bought |
|----------------|--------------|
| 2000 | A,B,C |
| 1000 | A,C |
| 4000 | A,D |
| 5000 | B,E,F |

Min. support 50%
Min. confidence 50%

| Frequent Itemset | Support |
|------------------|---------|
| {A} | 75% |
| {B} | 50% |
| {C} | 50% |
| {A,C} | 50% |

For rule $A \Rightarrow C$:

support = support($\{A \Rightarrow C\}$) = 50%

confidence = support($\{A \Rightarrow C\}$)/support($\{A\}$) = 66.6%

The Apriori principle:

Any subset of a frequent itemset must be frequent

Mining Frequent Itemsets: the Key Step

- Find the *frequent itemsets*: the sets of items that have minimum support
 - A subset of a frequent itemset must also be a frequent itemset
 - i.e., if $\{AB\}$ is a frequent itemset, both $\{A\}$ and $\{B\}$ should be a frequent itemset
 - Iteratively find frequent itemsets with cardinality from 1 to k (k -itemset)
- Use the frequent itemsets to generate association rules.

The Apriori Algorithm

- Join Step: C_k is generated by joining L_{k-1} with itself
- Prune Step: Any $(k-1)$ -itemset that is not frequent cannot be a subset of a frequent k -itemset
- Pseudo-code:
 - C_k : Candidate itemset of size k
 - L_k : frequent itemset of size k

$L_1 = \{\text{frequent items}\};$

for $(k = 1; L_k \neq \emptyset; k++)$ **do begin**

C_{k+1} = candidates generated from L_k ;

for each transaction t in database **do**

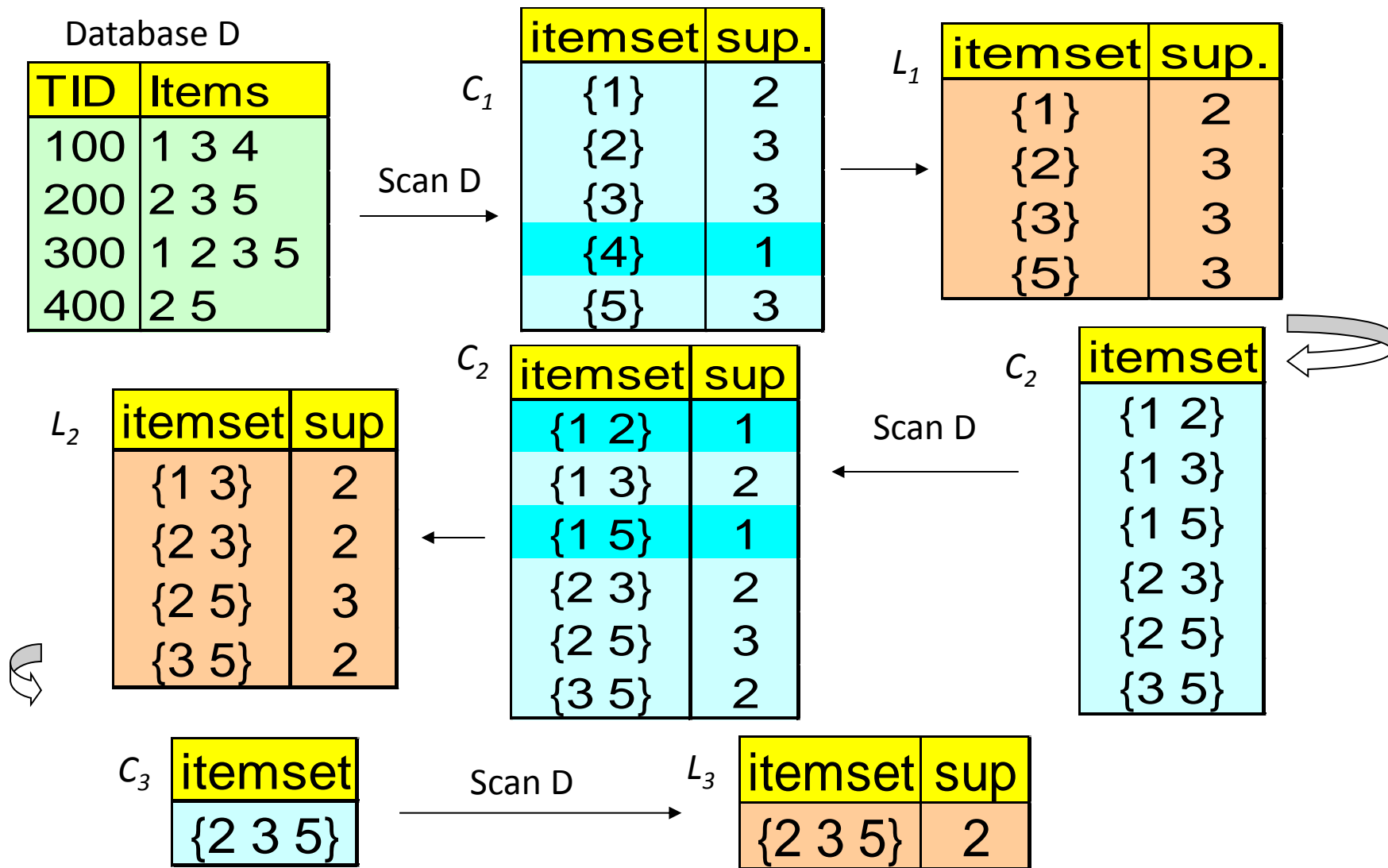
 increment the count of all candidates in C_{k+1}
 that are contained in t

L_{k+1} = candidates in C_{k+1} with min_support

end

return $\cup_k L_k$;

The Apriori Algorithm — Example



Example of Generating Candidates

- $L_3 = \{abc, abd, acd, ace, bcd\}$
- Self-joining: $L_3 * L_3$
 - $abcd$ from abc and abd
 - $acde$ from acd and ace
- Pruning:
 - $acde$ is removed because ade is not in L_3
- $C_4 = \{abcd\}$

**ΟΙΚΟΝΟΜΙΚΟ
ΠΑΝΕΠΙΣΤΗΜΙΟ
ΑΘΗΝΩΝ**



**ATHENS UNIVERSITY
OF ECONOMICS
AND BUSINESS**

Spectral Clustering

Μάθημα: Εξόρυξη γνώσης από Βάσεις Δεδομένων και τον Παγκόσμιο Ιστό

Ενότητα # 4: Unsupervised Learning (Clustering)

Διδάσκων: Μιχάλης Βαζιργιάννης

Τμήμα: Προπτυχιακό Πρόγραμμα Σπουδών “Πληροφορικής”

Spectral Clustering

Given Graph $G=(V,E)$ undirected:

Vertex Set $V=\{v_1, \dots, v_n\}$, Edge e_{ij} between v_i and v_j
we assume weight $w_{ij}>0$ for e_{ij}

$|V|$: number of vertices

d_i degree of v_i : $d_i = \sum_{v_j \in V} w_{ij}$

$$v(V) = \sum_{v_i \in V} d_i$$

for $A \subset V$ $\bar{A} = V - A$

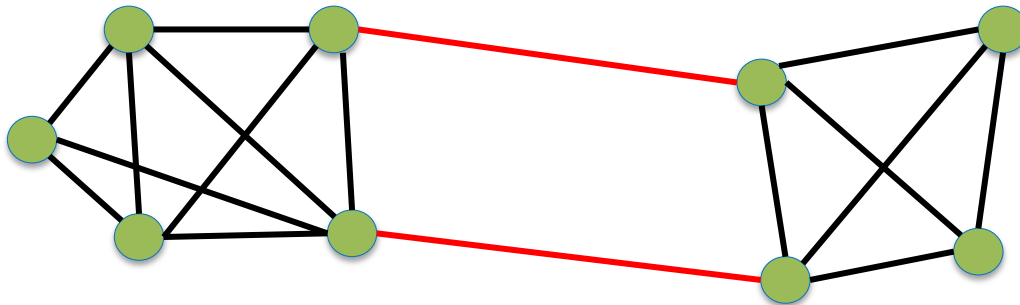
Given $A, B \subset V$ & $A \cap B = \emptyset$ $w(A, B) = \sum_{v_i \in A, v_j \in B} w_{ij}$

D: Diagonal matrix where $D(i,i)=d_i$

W: Adjacency matrix $W(i,j)=w_{ij}$

Graph-Cut

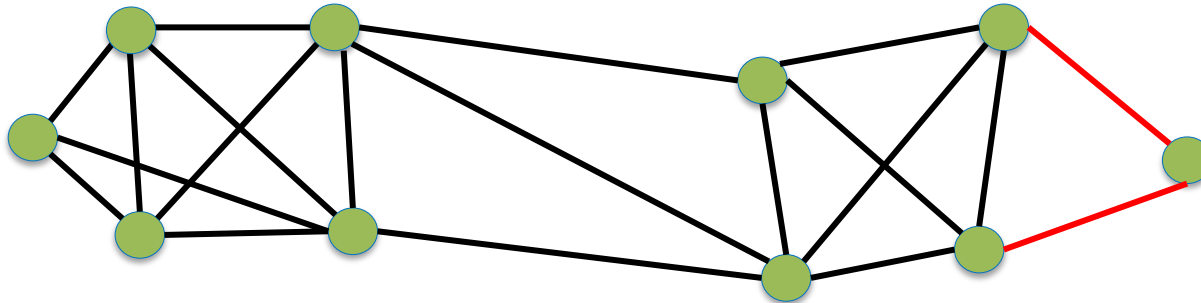
- For k clusters:
 - undirected graph: $1/2$ we count twice each edge



- Min-cut: Minimize the edges' weight a cluster shares with the rest of the graph

Min-Cut

- Easy for $k=2$: $\text{Mincut}(A_1, A_2)$
 - Stoer and Wagner: “A Simple Min-Cut Algorithm”
- In practice one vertex is separated from the rest
 - The algorithm is drawn to outliers



Normalized Graph Cuts

We can normalize by the size of the cluster (size of sub-graph) :

number of Vertices (Hagen and Kahng, 1992):

$$Ratiocut(A_1, \dots, A_k) = \sum_{i=1}^k \frac{cut(A_i, \overline{A_i})}{|A_i|}$$

sum of weights (Shi and Malik, 2000) :

$$Ncut(A_1, \dots, A_k) = \sum_{i=1}^k \frac{cut(A_i, \overline{A_i})}{v(A_i)}$$

Optimizing these functions is NP-hard

Spectral Clustering provides solution to a relaxed version of the above

From Graph Cuts to Spectral Clustering

For simplicity assume $k=2$:

Define $f: V \rightarrow \mathbb{R}$ for Graph G :

$$f_i = \begin{cases} 1 & v_i \in A \\ -1 & v_i \in \bar{A} \end{cases}$$

Optimizing the original cut is equivalent to an optimization of:

$$\begin{aligned} & \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2 \\ &= \sum_{v_i \in A, v_j \in \bar{A}} w_{ij} (1 - (-1))^2 + \sum_{v_i \in \bar{A}, v_j \in A} w_{ij} (-1 - 1)^2 \\ &= \mathbf{8 * cut}(A, \bar{A}) \end{aligned}$$

Graph Laplacian

How is the previous useful in Spectral clustering?

$$\begin{aligned} & \sum_{i,j=1}^n w_{ij}(f_i - f_j)^2 \\ &= \sum_{i,j=1}^n w_{ij}f_i^2 - 2 \sum_{i,j=1}^n w_{ij}f_i f_j + \sum_{i,j=1}^n w_{ij}f_j^2 \\ &= \sum_{i,j=1}^n d_i f_i^2 - 2 \sum_{i,j=1}^n w_{ij}f_i f_j + \sum_{i,j=1}^n d_j f_j^2 \\ &= 2 \left(\sum_{i,j=1}^n d_i f_i^2 - \sum_{i,j=1}^n w_{ij}f_i f_j \right) \\ &= 2(\mathbf{f}^T \mathbf{D} \mathbf{f} - \mathbf{f}^T \mathbf{W} \mathbf{f}) = 2\mathbf{f}^T (\mathbf{D} - \mathbf{W}) \mathbf{f} = 2\mathbf{f}^T \mathbf{L} \mathbf{f} \end{aligned}$$

\mathbf{f} : a single vector with the cluster assignments of the vertices

$\mathbf{L} = \mathbf{D} - \mathbf{W}$ the Laplacian of a graph

Properties of L

L is

Symmetric

Positive

Semi-definite ($x^T Lx$, x is a non negative vector)

The smallest eigenvalue of L is 0

The corresponding eigenvector is $\mathbb{1}$

L has n non-negative, real valued eigenvalues

$$0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

Two Way Cut from the Laplacian

We could solve $\min_f f^T L f$ where $f \in \{-1, 1\}^n$

NP-Hard for discrete cluster assignments

Relax the constraint to $f \in R^n$:

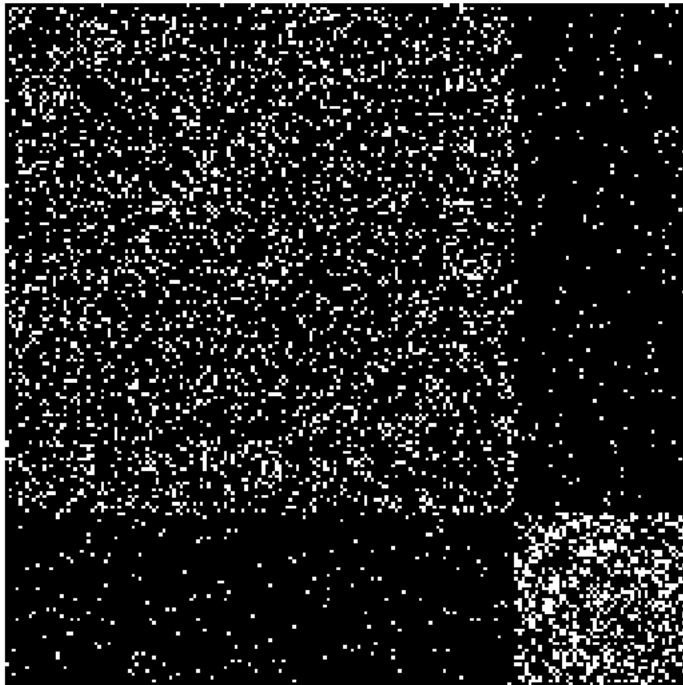
$$\min_f f^T L f \text{ subject to } f^T \mathbf{1} = 0$$

Solution: (**Rayleigh-Ritz Theorem**)

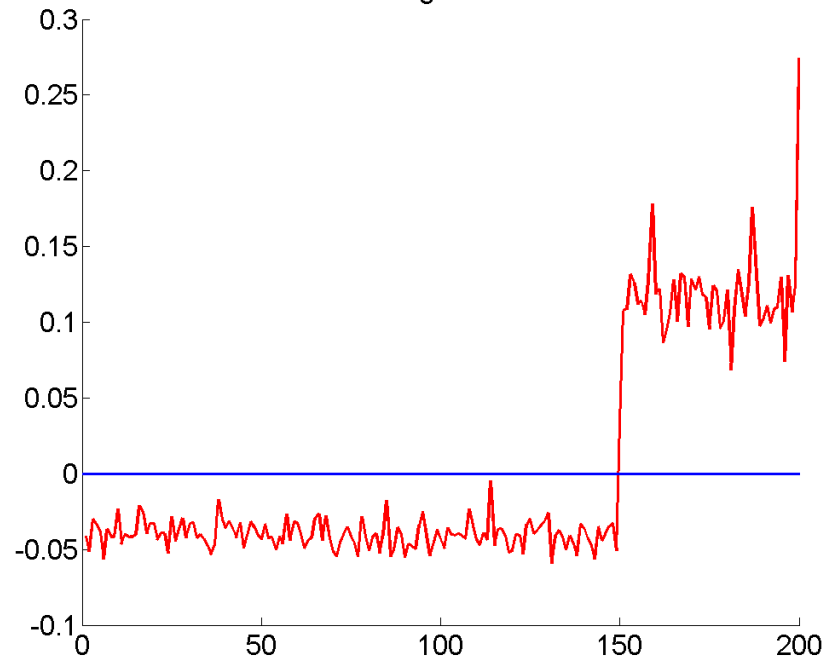
- the eigenvector corresponding to smallest eigenvalue: 0 TRIVIA as it offers no information
- we use the second eigenvector as an approximation ($f_i > 0$ vertex belongs to cluster A, $f_i < 0$ to the other cluster V-A)

Example

Adjacency Matrix



2nd Eigenvector



Ratio Cut

$$\text{Ratiocut}(A_1, \dots, A_k) = \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{|A_i|}$$

Define $f: V \rightarrow \mathbb{R}$ for Graph G :

$$f_i = \begin{cases} \sqrt{\frac{|\bar{A}|}{|A|}} & v_i \in A \\ -\sqrt{\frac{|A|}{|\bar{A}|}} & v_i \in \bar{A} \end{cases}$$

$$\sum_{i,j=1}^n w_{ij}(f_i - f_j)^2 = 2\text{cut}(A, \bar{A}) \left(\sqrt{\frac{|\bar{A}|}{|A|}} + \sqrt{\frac{|A|}{|\bar{A}|}} + 2 \right)$$

$$= 2|V|\text{Ratiocut}(A, \bar{A})$$

Ratio Cut

We have $\min_f f^T L f$ subject to

$$f^T \mathbf{1} = 0, f^T f = n$$

$$f^T \mathbf{1} = \sum_i^n f_i = \sum_{v_i \in A} \sqrt{\frac{|\bar{A}|}{|A|}} + \sum_{v_i \in \bar{A}} -\sqrt{\frac{|A|}{|\bar{A}|}} = |A| \sqrt{\frac{|\bar{A}|}{|A|}} - |\bar{A}| \sqrt{\frac{|A|}{|\bar{A}|}} = 0$$

$$f^T f = \sum_i^n f_i^2 = |\bar{A}| + |A| = n$$

The second smallest eigenvalue of $L f = \lambda f$ approximates the solution

Normalized Cut

- $Ncut(A_1, \dots, A_k) = \sum_{i=1}^k \frac{cut(A_i, \bar{A}_i)}{v(A_i)}$

- Define $f: V \rightarrow \mathbb{R}$ for Graph G :

-

$$f_i = \begin{cases} \sqrt{\frac{v(\bar{A})}{v(A)}} & v_i \in A \\ -\sqrt{\frac{v(A)}{v(\bar{A})}} & v_i \in \bar{A} \end{cases}$$

- $\sum_{i,j=1}^n w_{ij} (f_i - f_j)^2 =$
 $2cut(A, \bar{A}) \left(\sqrt{\frac{v(\bar{A})}{v(A)}} + \sqrt{\frac{v(A)}{v(\bar{A})}} + 2 \right)$
 $= 2v(V)Ncut(A, \bar{A})$

Normalized Cut

Similarly we come to : $\min_f f^T L f$
subject to $f^T D \mathbf{1} = 0$, $f^T D f = v(V)$

Assume $h = D^{1/2} f$

$\min_h h^T D^{-1/2} L D^{-1/2} h$ subject to
 $h^T D^{1/2} \mathbf{1} = 0$, $h^T h = v(V)$

The answer is in the eigenvector of the second smallest eigenvalue
of $L_{sym} = D^{-1/2} L D^{-1/2}$
Shi and Malik (2000)

L_{sym} is the normalized Laplacian

has n non-negative, real valued eigenvalues

$$0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

Multi-Way Graph Partition

Define $f_{ij} = \begin{cases} \frac{1}{\sqrt{|A_j|}} & v_i \in A_j \\ 0 & \text{otherwise} \end{cases}$

we have a vector indicating the cluster a vertex belongs to

Similarly to the other equations we can deduce:

$$f_i^T L f_i = cu(A_i, \bar{A}_i) / |A_i|$$

$$\sum_{i=1}^k f_i^T L f_i = \sum_{i=1}^k (F^T L F)_{ii} = Tr(F^T L F)$$

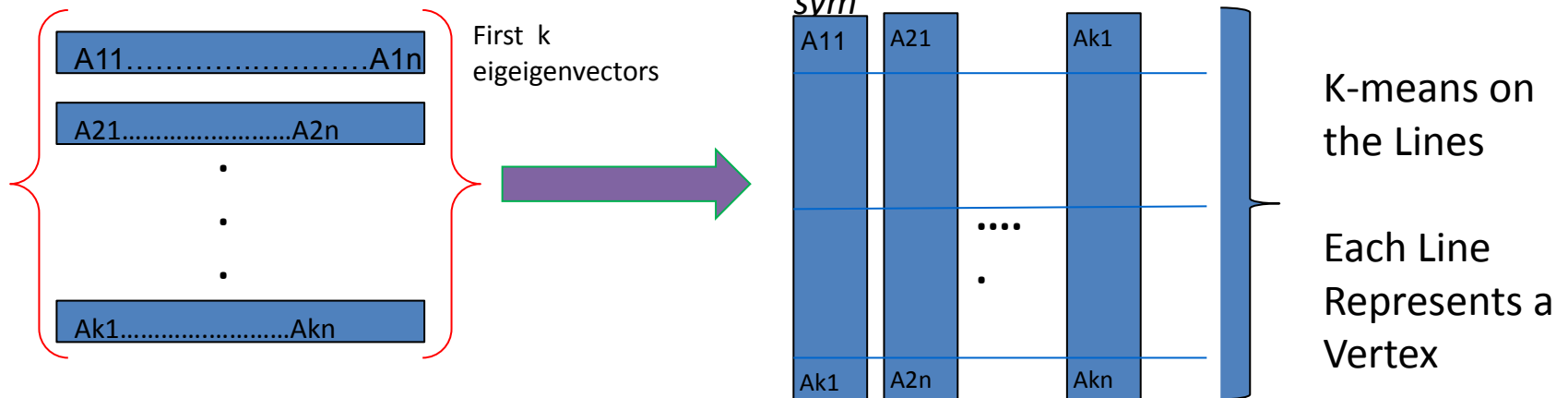
Where Tr is the Trace of a Matrix

So now the RatioCut becomes:

$$min \sum_{i=1}^k (F^T L F)_{ii} \text{ subject to } F^T F = I$$

Multi-Way Graph Partition

- The solution can now be given by the first k eigenvectors of L as columns
- The real values need to be converted to cluster assignments
 - We use k -means to cluster the rows
 - We can substitute L with L_{sym}



Algorithm for $k > 2$

Compute Laplacian (L, L_{sym}).

Compute the first k eigenvectors u_1, \dots, u_k of L .

*Let $U \in \mathbb{R}^{n \times k}$ the matrix containing the vectors
 u_1, \dots, u_k as columns.*

For $i = 1, \dots, n$,

let $y_i \in \mathbb{R}^k$ the vector corresponding to the i -th row of U .

*Cluster the points $y_i = 1, \dots, n \in \mathbb{R}^k$ with the k -means algorithm into
clusters C_1, \dots, C_k .*

Output: Clusters A_1, \dots, A_k with $A_i = \{j | v_j \in C_i\}$

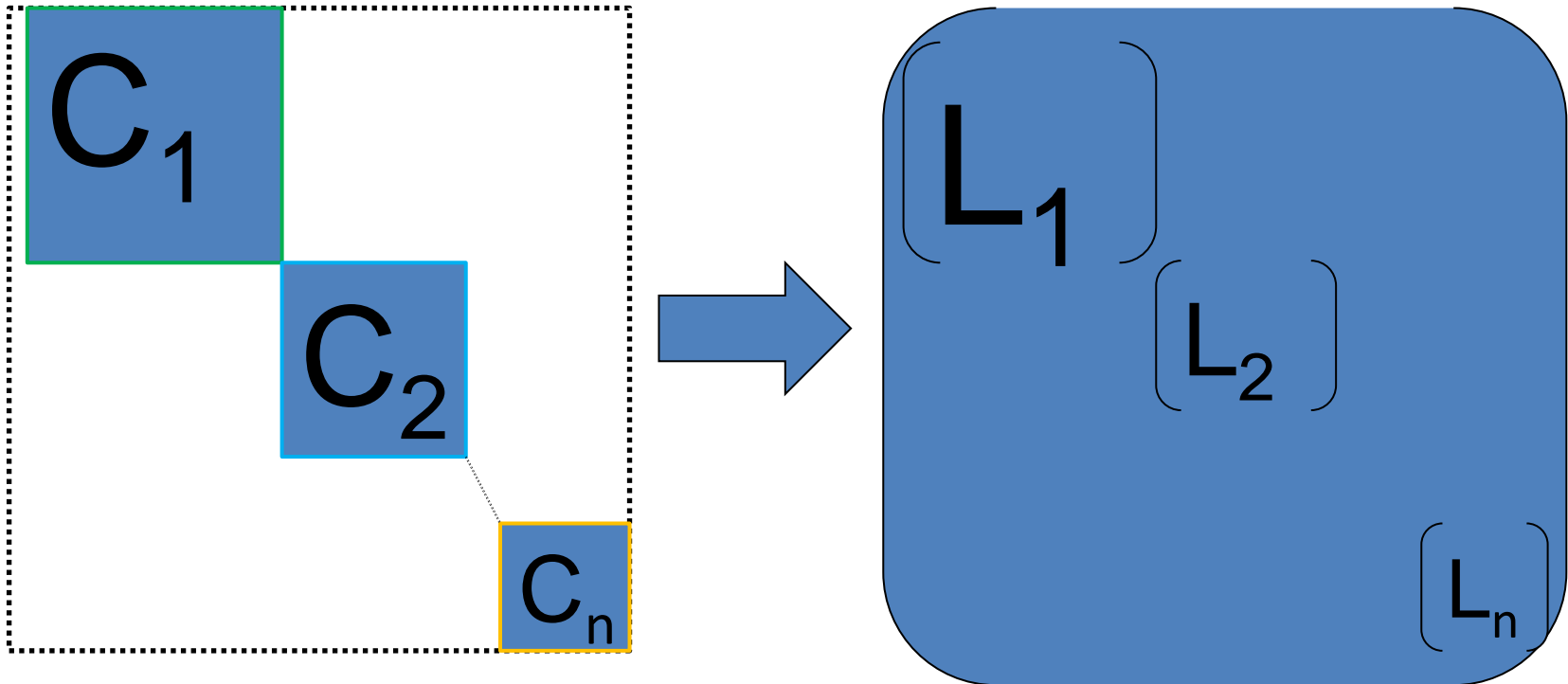
HOW DO WE CHOOSE k ?

We choose the k that maximizes the eigengap:

$$\Delta_k = |\lambda_k - \lambda_{k-1}| \text{ (Davis-Kahan Theorem)}$$

Ideally: for k connected components the Laplacian has k 0-eigenvalues

Laplacian-Eigenvectors-EigenValues



Everything sorted according to cluster : block diagonal form Matrix

L follows the same form composed on $L_1 \dots L_n$

Each L_i has the same properties as L : $n \times 0$ min eigenvalues etc..

Each "Second" eigenvector is a cut of C_i from the rest of the graph and holds a mapping (distance) of a vertex to the cluster i

Simple example

$$\begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix}$$

2 Eigenvectors

(1100) and
(0011)

Mapping vertices
in their clusters

**Permutation does not change
the result**

$$\begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix}$$

2 Eigenvectors

(1010) and
(0101)

Mapping vertices
to the same
clusters

**The cut remains the same
regardless of the ordering**

References

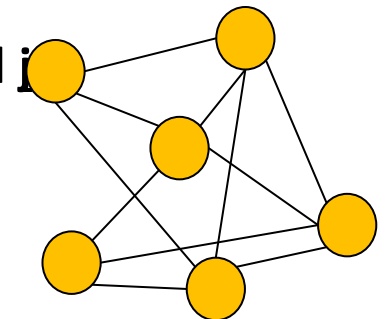
1. Ulrike von Luxburg, A Tutorial on Spectral Clustering, *Statistics and Computing*, 2007
2. Davis, C., W. M. Kahan (March 1970). The rotation of eigenvectors by a perturbation. III. *SIAM J. Numerical Analysis* 7
3. Shi, Jianbo, and Jitendra Malik. "Normalized cuts and image segmentation." *Pattern Analysis and Machine Intelligence, IEEE Transactions on* (2000).
4. Mechthild Stoer and Frank Wagner. 1997. A simple min-cut algorithm. *J. ACM*
5. Ng, Jordan & Weiss, K-means algorithm on the embeded eigen-space, *NIPS* 2001
6. Hagen, L. Kahng, , "New spectral methods for ratio cut partitioning and clustering," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on* , 1992

Modularity Based Methods

- Most of the community evaluation measures (e.g., conductance, cut-based measures), quantify the quality of a community based on
 - **Internal connectivity** (intra-community edges)
 - **External connectivity** (inter-community edges)
- **Question:** Is there any other way to distinguish groups of nodes with good community structure?
- **Random graphs** are not expected to present inherent community structure
- **Idea:** Compare the number of edges that lie **within a cluster** with the expected one in case of **random graphs** with the same degree distribution – **modularity measure**

Main idea

- **Modularity** function [Newman and Girvan '04], [Newman '06]
- Initially introduced as a measure for assessing the strength of communities
 - **Q = (fraction of edges within communities) –**
(expected number of edges within communities)
- What is the **expected** number of edges?
- Consider a configuration model
 - **Random graph** model with the same degree distribution
 - Let P_{ij} = probability of an edge between nodes i and j with degrees k_i and k_j respectively
 - Then $P_{ij} = k_i k_j / 2m$, where $m = |E| = \frac{1}{2} \sum_i k_i$



Formal definition of modularity

$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

where

- **A** is the adjacency matrix
- **k_i**, **k_j** the degrees of nodes **i** and **j** respectively
- **m** is the number of edges
- **C_i** is the community of node **i**
- **δ(.)** is the Kronecker function: 1 if both nodes **i** and **j** belong on the same community (**C_i = C_j**), 0 otherwise

[Newman and Girvan '04], [Newman '06]

Properties of modularity

$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

- **Larger modularity Q indicates better communities** (more than random intra-cluster density)
 - The community structure would be better if the number of internal edges exceed the expected number
- Modularity value is always **smaller than 1**
- It can also take **negative values**
 - E.g., if each node is a community itself
 - No partitions with positive modularity → No community structure
 - Partitions with large negative modularity → Existence of subgraphs with small internal number of edges and large number of inter-community edges

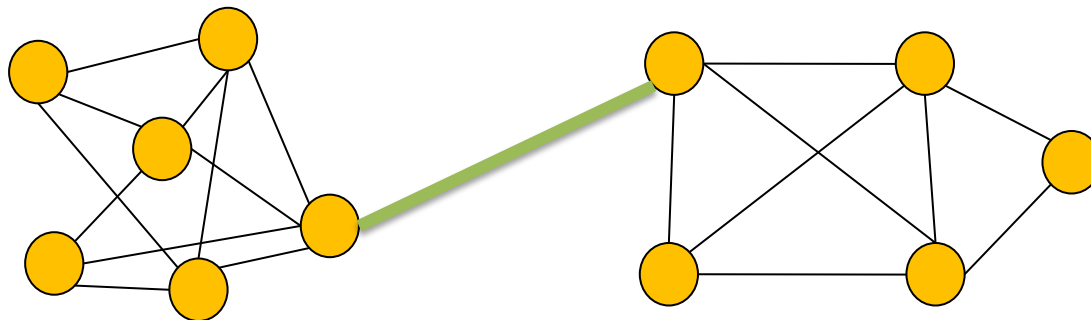
Applications of modularity

- Modularity can be applied:
 - As **quality function** in clustering algorithms
 - As **evaluation measure** for comparison of different partitions or algorithms
 - As a community detection tool itself
 - **Modularity optimization**
 - As criterion for reducing the size of a graph
 - **Size reduction preserving modularity [Arenas et al. '07]**

[Newman and Girvan '04], [Newman '06], [Fortunato '10]

Modularity-based clustering

- Modularity was first applied as a **stopping criterion** in the Newman-Girvan algorithm
- Newman-Girvan algorithm [**Newman and Girvan '04**]
 - A **divisive** algorithm (detect and remove edges that connect vertices of different communities)
 - **Idea:** try to identify the edges of the graph that are most between other vertices → responsible for connecting many node pairs
 - Select and remove edges based to the value of **betweenness centrality**
 - **Betweenness centrality:** number of **shortest paths** between every pair of nodes, that pass through an edge



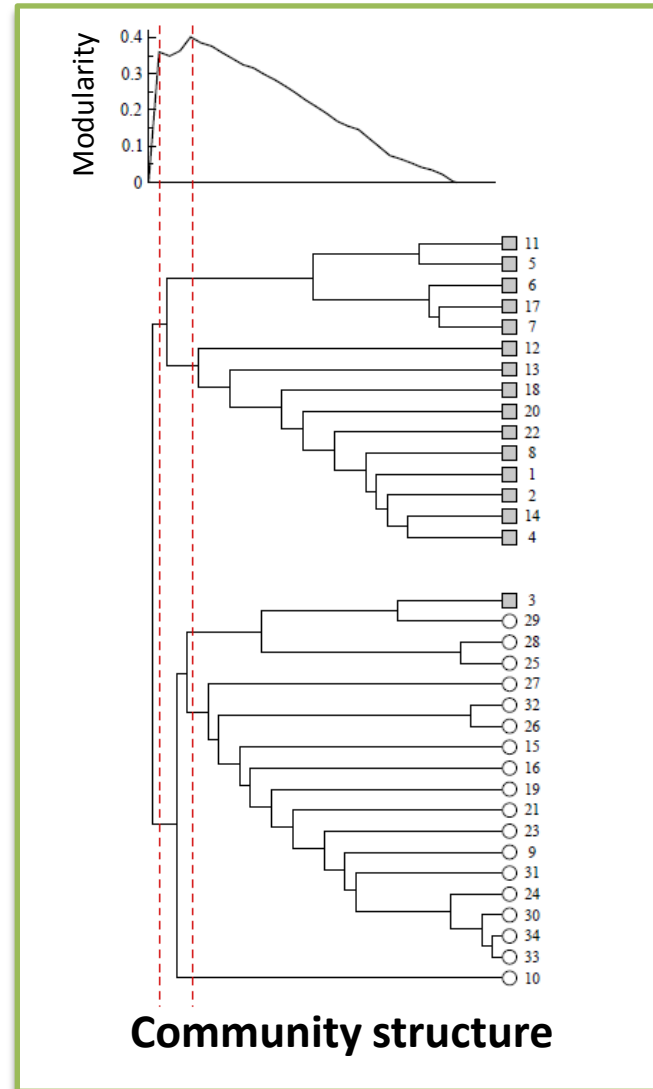
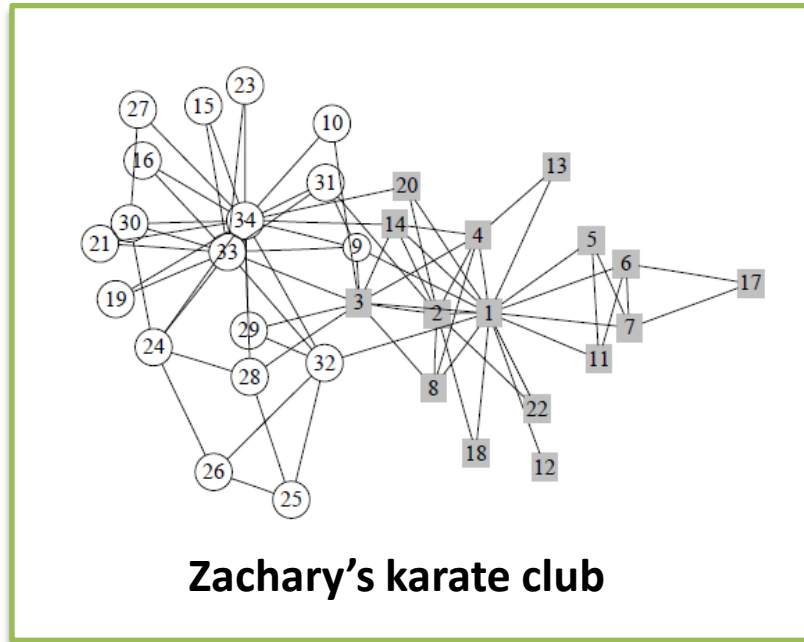
Edge betweenness is higher for edges that connect different communities

Newman-Girvan algorithm (1)

- **Basic steps:**
 1. Compute betweenness centrality for all edges in the graph
 2. Find and remove the edge with the highest score
 3. Recalculate betweenness centrality score for the remaining edges
 4. Go to step 2
- How do we know if the produced communities are **good ones** and stop the algorithm?
 - The output of the algorithm is in the form of a **dendrogram**
 - Use **modularity** as a criterion to cut the dendrogram and terminate the algorithm ($Q \approx 0.3-0.7$ indicates good partitions)
- Complexity: **$O(m^2n)$** (or **$O(n^3)$** on a sparse graph)

[Newman and Girvan '04], [Girvan and Newman '02]

Newman-Girvan algorithm (2)



[Newman and Girvan '04]

Modularity optimization

- High values of modularity indicate good quality of partitions
- **Goal:** find the partition that corresponds to the maximum value of modularity
- **Modularity maximization** problem
 - Computational difficult problem [**Brandes et al. '06**]
 - Approximation techniques and heuristics
- Four main categories of techniques
 1. Greedy techniques
 2. Spectral optimization
 3. Simulated annealing
 4. Extremal optimization

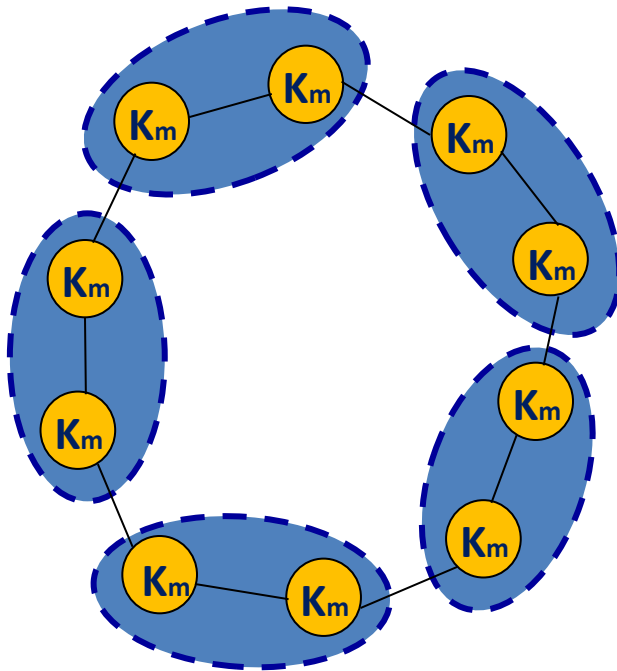
[Fortunato '10]

Greedy techniques

- Newman's algorithm [**Newman '04b**]
 - Agglomerative (bottom-up) hierarchical clustering algorithm
 - **Idea:** Repeatedly join pairs of communities that achieve the greatest increase of modularity (dendrogram representation)
 1. Initially, each node of the graph belongs on its own cluster (**n**)
 2. Repeatedly, join communities in pairs by adding edges
 - a. At each step, choose the pairs that achieve the **greatest increase** (or minimum decrease) of modularity
 - b. Consider only pairs of communities **between which there exist edges** (merging communities that do not share edges, it can never improve modularity)
 - Complexity: **$O((m+n) n)$** (or **$O(n^2)$** on a sparse graph)

Resolution limit of modularity

- **Resolution Limit** of modularity [Fortunato and Barthelemy '07]
- The method of modularity optimization may not detect communities with relatively small size, which depends on the total number of edges in the graph



- K_m are cliques with m edges ($m \leq \sqrt{|E|}$)
- K_m represent well-defined clusters
- However, the maximum modularity corresponds to clusters formed by two or more cliques
- It is difficult to know if the community returned by modularity optimization corresponds to a **single community** or a **union of smaller communities**

References (modularity)

- M.E.J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical Review E* 69(02), 2004.
- M.E.J. Newman. Modularity and community structure in networks. *PNAS*, 103(23), 2006.
- S.E. Schaeffer. Graph clustering. *Computer Science Review* 1(1), 2007.
- S. Fortunato. Community detection in graphs. *Physics Reports* 486 (3-5), 2010.
- M. Coscia, F. Giannotti, and D. Pedreschi. A classification for community discovery methods in complex networks. *Statistical Analysis and Data Mining* 4 (5), 2011.
- A. Arenas, J. Duch, A. Fernandez, and S. Gomez. Size reduction of complex networks preserving modularity. *New J. Phys.*, 9(176), 2007.
- M. Girvan and M.E.J. Newman. Community structure in social and biological networks. *PNAS* 99(12), 2002.
- U. Brandes, D. Delling, M. Gaertler, R. Gorke, M. Hoefer, Z. Nikoloski, and D. Wagner. On Modularity Clustering. *IEEE TKDE* 20(2), 2008.
- M.E.J. Newman. Fast algorithm for detecting community structure in networks. *Phys. Rev. E* 69, 2004.
- A. Clauset, M.E.J. Newman, and C. Moore. Finding community structure in very large networks. *Phys. Rev. E* 70, 2004.

References (modularity)

- M.E.J. Newman. Finding community structure in networks using the eigenvectors of matrices. *Phys. Rev. E* 74, 2006.
- R. Guimera, M. Sales-Pardo, L.A.N. Amaral. Modularity from Fluctuations in Random Graphs and Complex Networks. *Phys. Rev. E* 70, 2004.
- J. Duch and A. Arenas. Community detection in complex networks using Extremal Optimization. *Phys. Rev. E* 72, 2005.
- A. Arenas, J. Duch, A. Fernandez, and S. Gomez. Size reduction of complex networks preserving modularity. *New Journal of Physics* 9(6), 2007.
- E.A. Leicht and M.E.J. Newman. Community structure in directed networks. *Phys. Rev. Lett.* 100, 2008.
- V. Nicosia, G. Mangioni, V. Carchiolo, and M. Malgeri. Extending the definition of modularity to directed graphs with overlapping communities. *J. Stat. Mech.* 03, 2009.
- S. Muff, F. Rao, A. Caflich. Local modularity measure for network clusterizations. *Phys. Rev. E*, 72, 2005.
- S. Fortunato and M. Barthelemy. Resolution limit in community detection. *PNAS* 104(1), 2007.

**ΟΙΚΟΝΟΜΙΚΟ
ΠΑΝΕΠΙΣΤΗΜΙΟ
ΑΘΗΝΩΝ**

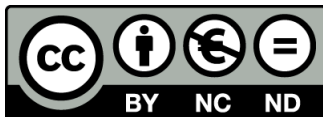


**ATHENS UNIVERSITY
OF ECONOMICS
AND BUSINESS**

Τέλος Ενότητας # 4

Μάθημα: Εξόρυξη γνώσης από Βάσεις Δεδομένων και τον Παγκόσμιο Ιστό, **Ενότητα # 4:** Unsupervised Learning (Clustering)

Διδάσκων: Μιχάλης Βαζιργιάννης, **Τμήμα:** Προπτυχιακό Πρόγραμμα Σπουδών “Πληροφορικής”



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ & ΘΡΗΣΚΕΥΜΑΤΩΝ, ΠΟΛΙΤΙΣΜΟΥ & ΑΘΛΗΤΙΣΜΟΥ
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ