

Clustering Algorithms

Hierarchical Clustering

k -Means Algorithms

CURE Algorithm

Methods of Clustering

- ◆ **Hierarchical (Agglomerative):**
 - ◆ Initially, each point in cluster by itself.
 - ◆ Repeatedly combine the two "nearest" clusters into one.
- ◆ **Point Assignment:**
 - ◆ Maintain a set of clusters.
 - ◆ Place points into their "nearest" cluster.

Hierarchical Clustering

- ◆ Two important questions:
 1. How do you determine the “nearness” of clusters?
 2. How do you represent a cluster of more than one point?

Hierarchical Clustering --- (2)

- ◆ **Key problem**: as you build clusters, how do you represent the location of each cluster, to tell which pair of clusters is closest?
- ◆ **Euclidean case**: each cluster has a *centroid* = average of its points.
 - ◆ Measure intercluster distances by distances of centroids.

Example



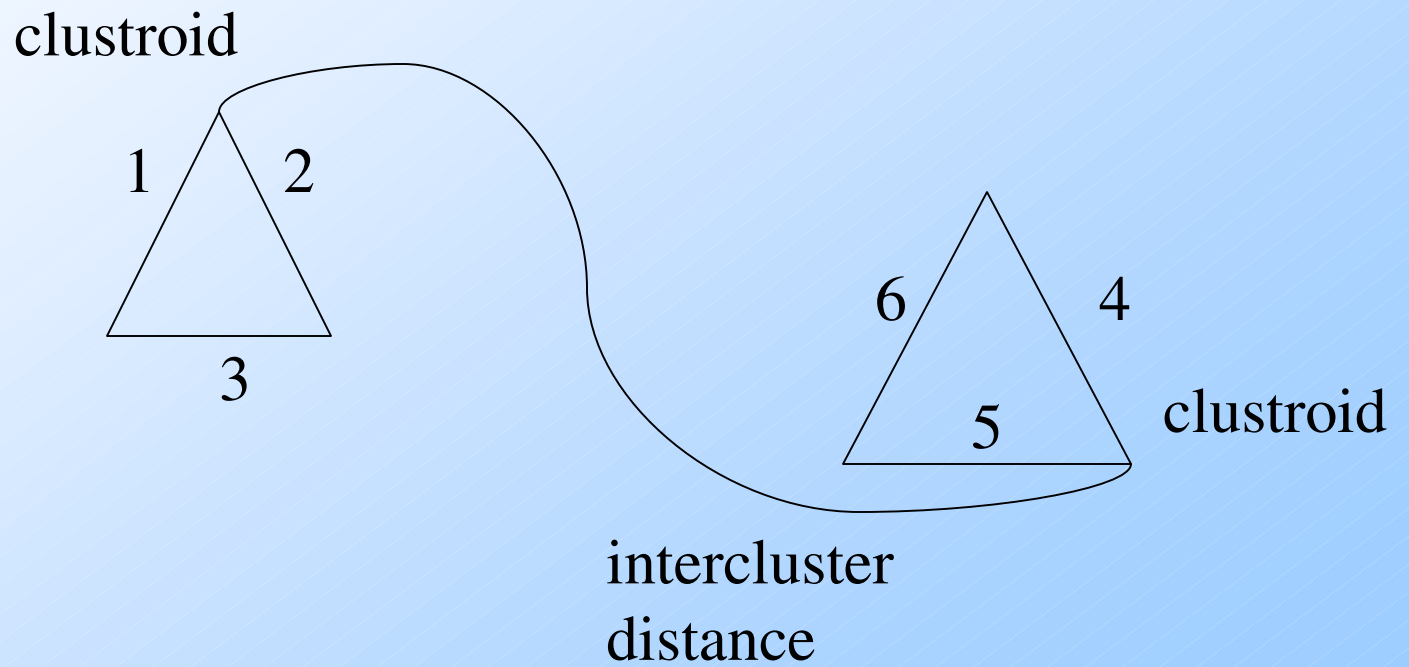
And in the Non-Euclidean Case?

- ◆ The only “locations” we can talk about are the points themselves.
 - ◆ I.e., there is no “average” of two points.
- ◆ Approach 1: *clustroid* = point “closest” to other points.
 - ◆ Treat clustroid as if it were centroid, when computing intercluster distances.

“Closest” Point?

- ◆ Possible meanings:
 1. Smallest maximum distance to the other points.
 2. Smallest average distance to other points.
 3. Smallest sum of squares of distances to other points.
 4. Etc., etc.

Example



Other Approaches to Defining “Nearness” of Clusters

- ◆ **Approach 2:** intercluster distance = minimum of the distances between any two points, one from each cluster.
- ◆ **Approach 3:** Pick a notion of “cohesion” of clusters, e.g., maximum distance from the clustroid.
 - ◆ Merge clusters whose *union* is most cohesive.

Return to Euclidean Case

- ◆ Approaches 2 and 3 are also used sometimes in Euclidean clustering.
- ◆ Many other approaches as well, for both Euclidean and non.

k -Means Algorithm(s)

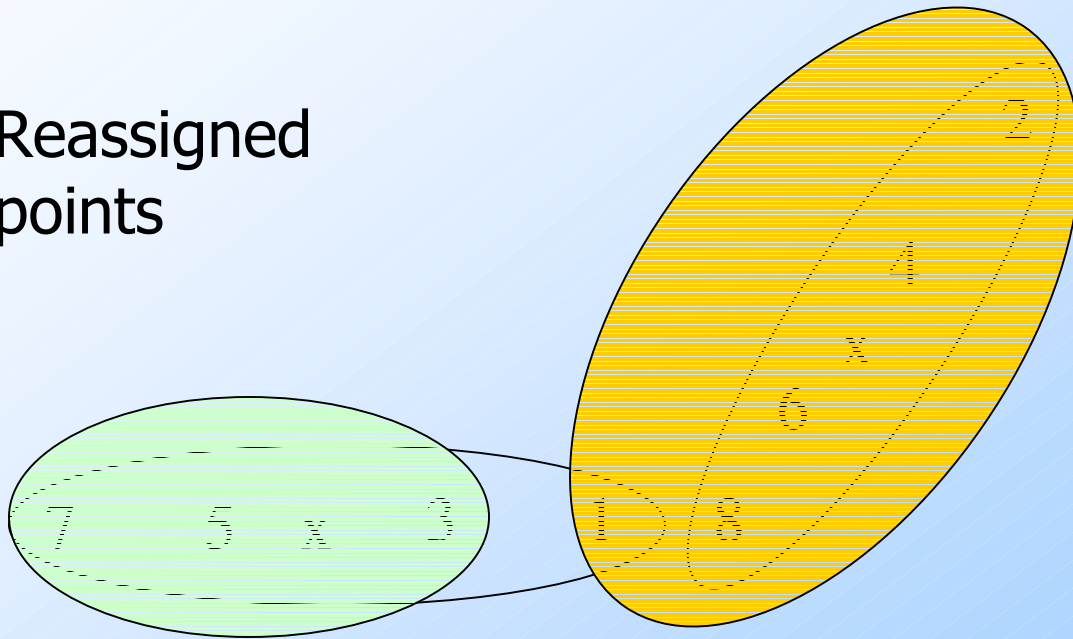
- ◆ Assumes Euclidean space.
- ◆ Start by picking k , the number of clusters.
- ◆ Initialize clusters by picking one point per cluster.
 - ◆ For instance, pick one point at random, then $k-1$ other points, each as far away as possible from the previous points.

Populating Clusters

1. For each point, place it in the cluster whose current centroid it is nearest.
2. After all points are assigned, fix the centroids of the k clusters.
3. **Optional**: reassign all points to their closest centroid.
 - ◆ Sometimes moves points between clusters.

Example

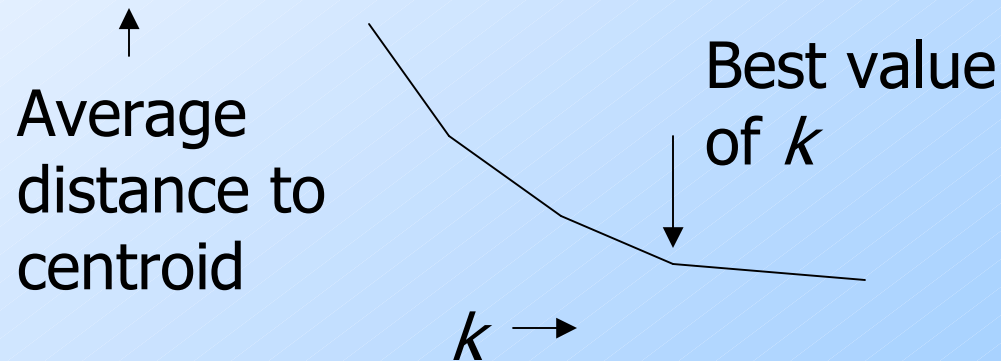
Reassigned
points



Clusters after first round

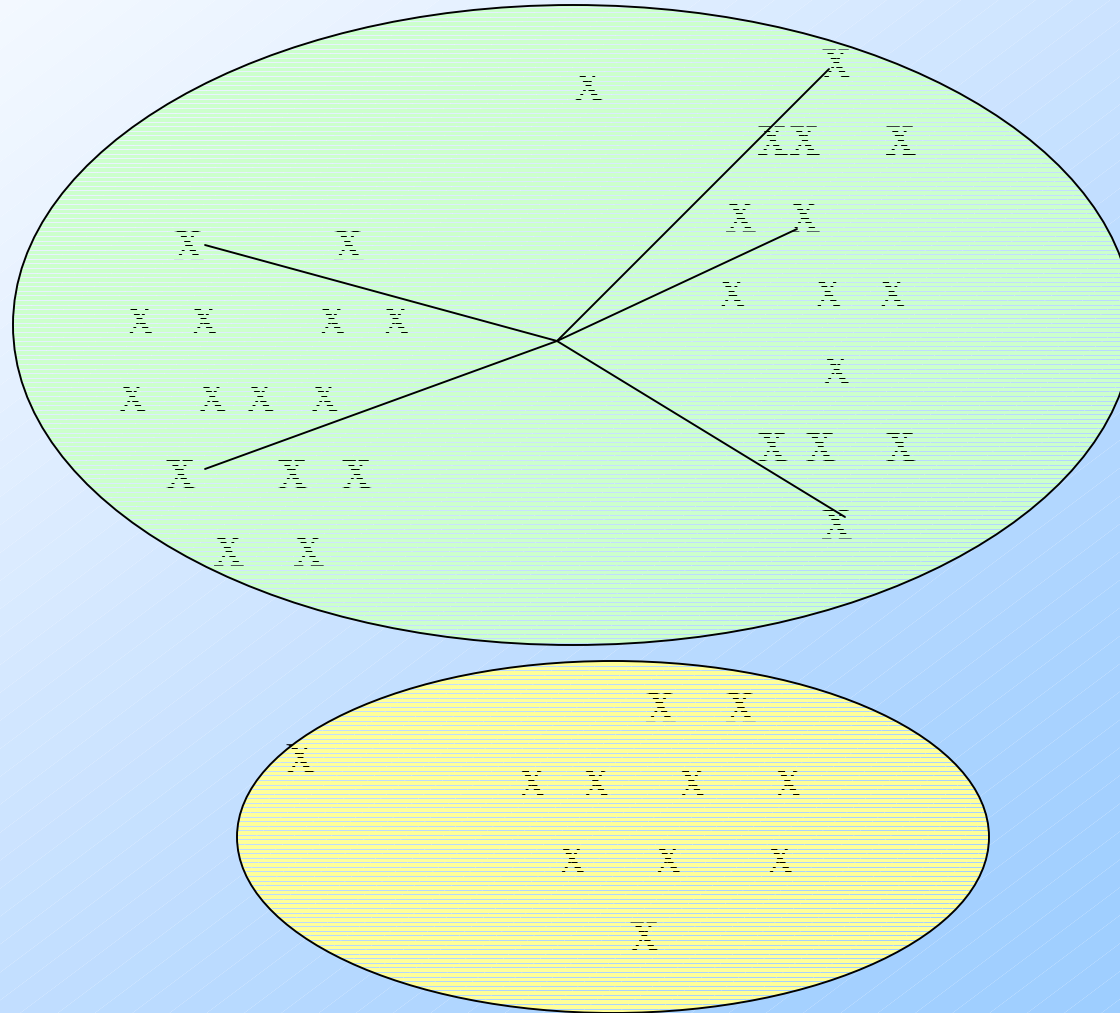
Getting k Right

- ◆ Try different k , looking at the change in the average distance to centroid, as k increases.
- ◆ Average falls rapidly until right k , then changes little.



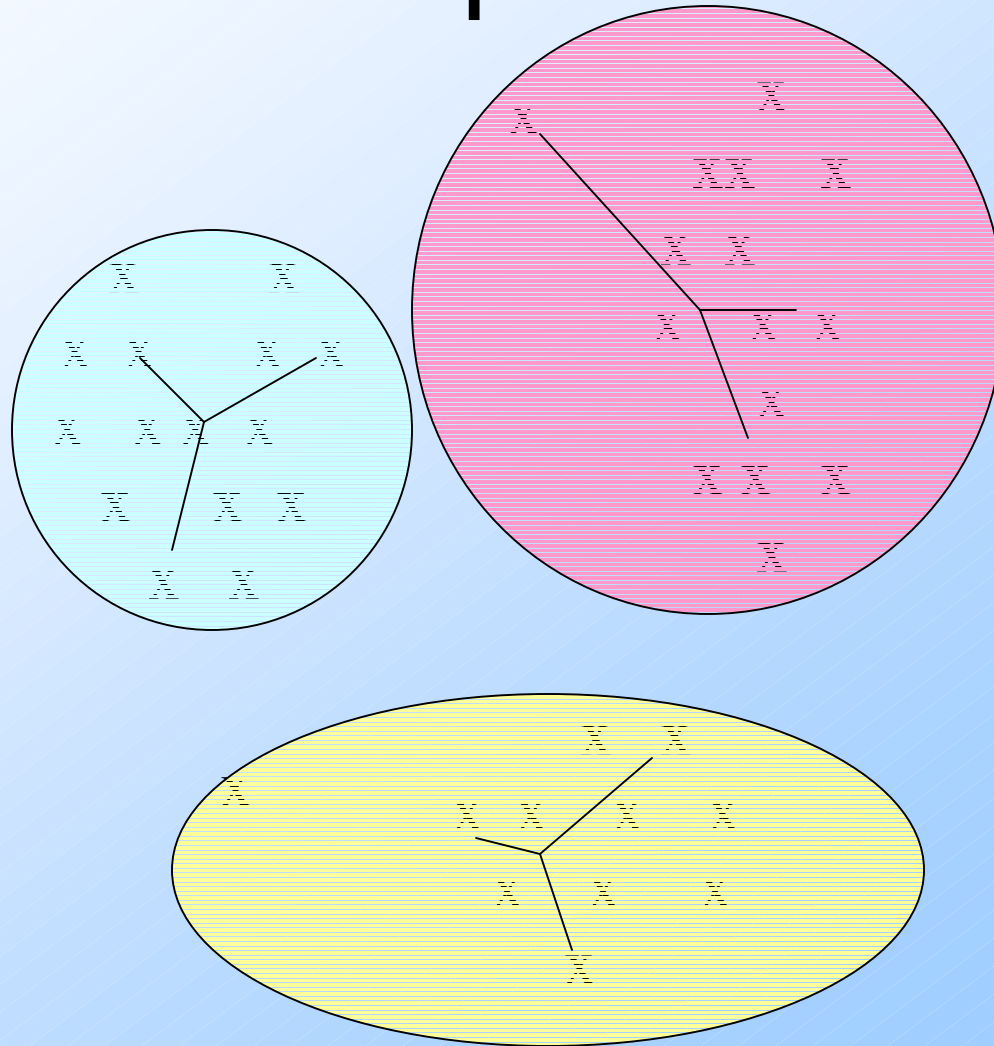
Example

Too few;
many long
distances
to centroid.



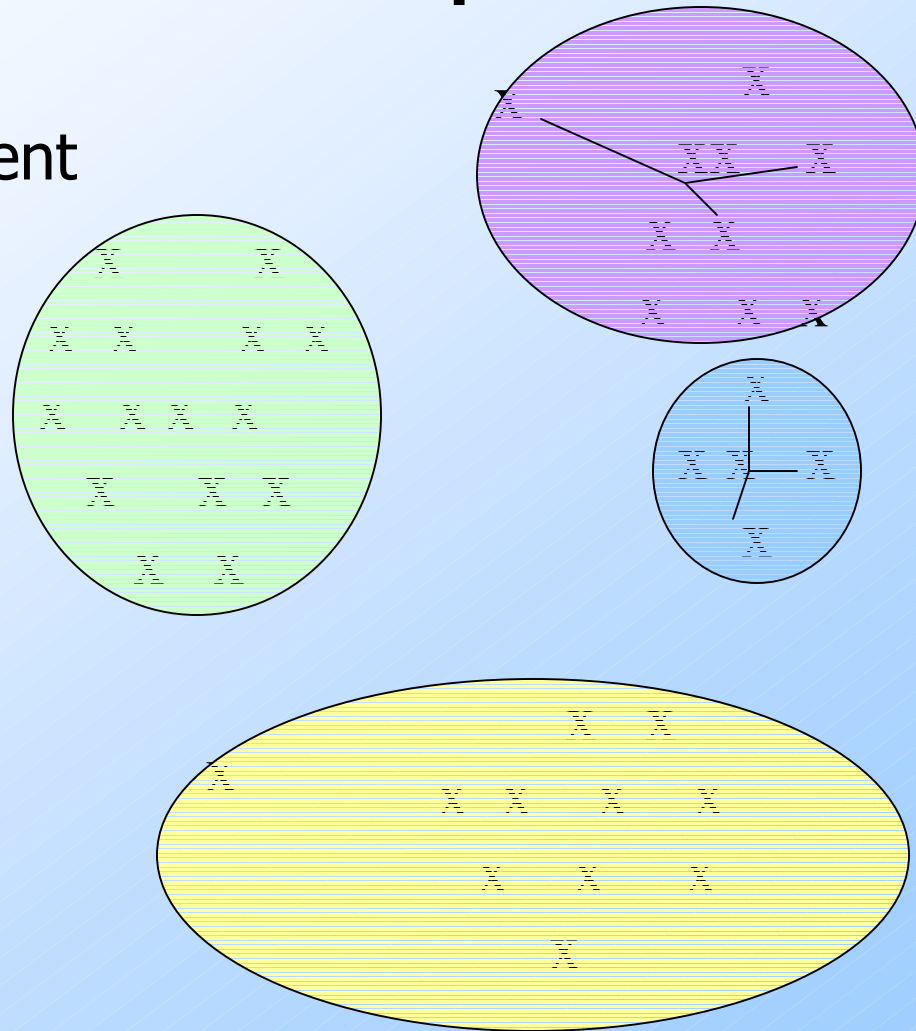
Example

Just right;
distances
rather short.



Example

Too many;
little improvement
in average
distance.



BFR Algorithm

- ◆ BFR (**Bradley-Fayyad-Reina**) is a variant of k -means designed to handle very large (disk-resident) data sets.
- ◆ It assumes that clusters are normally distributed around a centroid in a Euclidean space.
 - ◆ Standard deviations in different dimensions may vary.

BFR --- (2)

- ◆ Points are read one main-memory-full at a time.
- ◆ Most points from previous memory loads are summarized by simple statistics.
- ◆ To begin, from the initial load we select the initial k centroids by some sensible approach.

Initialization: k -Means

- ◆ Possibilities include:
 1. Take a small random sample and cluster optimally.
 2. Take a sample; pick a random point, and then $k - 1$ more points, each as far from the previously selected points as possible.

Three Classes of Points

1. The *discard set* : points close enough to a centroid to be represented statistically.
2. The *compression set* : groups of points that are close together but not close to any centroid. They are represented statistically, but not assigned to a cluster.
3. The *retained set* : isolated points.

Representing Sets of Points

- ◆ For each cluster, the discard set is represented by:
 1. The number of points, N .
 2. The vector SUM, whose i^{th} component is the sum of the coordinates of the points in the i^{th} dimension.
 3. The vector SUMSQ: i^{th} component = sum of squares of coordinates in i^{th} dimension.

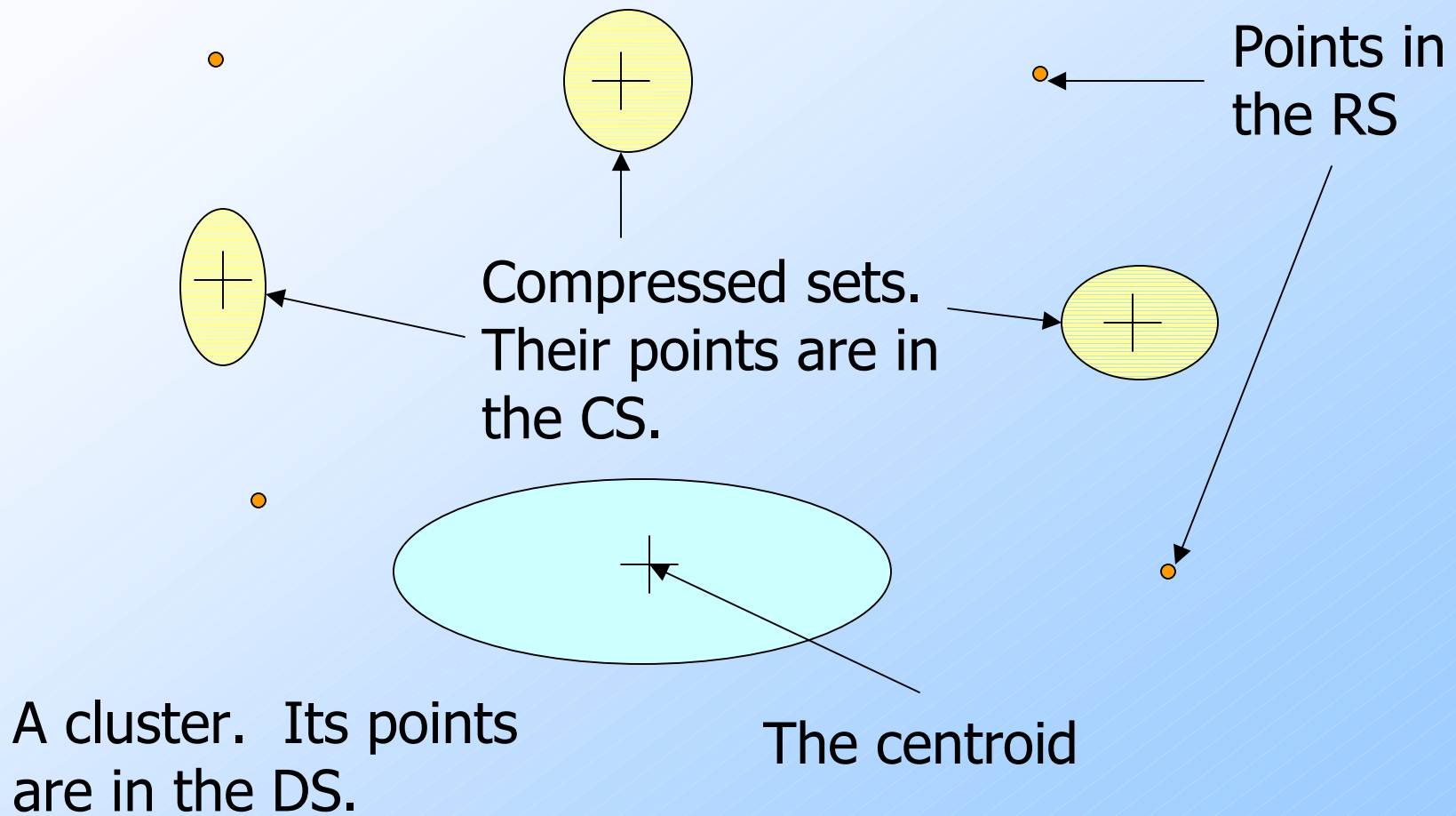
Comments

- ◆ $2d + 1$ values represent any number of points.
 - ◆ d = number of dimensions.
- ◆ Averages in each dimension (centroid coordinates) can be calculated easily as SUM_i / N .
 - ◆ $SUM_i = i^{\text{th}}$ component of SUM.

Comments --- (2)

- ◆ Variance of a cluster's discard set in dimension i can be computed by:
$$(\text{SUMSQ}_i / N) - (\text{SUM}_i / N)^2$$
- ◆ And the standard deviation is the square root of that.
- ◆ The same statistics can represent any compression set.

"Galaxies" Picture



Processing a “Memory-Load” of Points

1. Find those points that are “sufficiently close” to a cluster centroid; add those points to that cluster and the DS.
2. Use any main-memory clustering algorithm to cluster the remaining points and the old RS.
 - ◆ Clusters go to the CS; outlying points to the RS.

Processing --- (2)

3. Adjust statistics of the clusters to account for the new points.
4. Consider merging compressed sets in the CS.
5. If this is the last round, merge all compressed sets in the CS and all RS points into their nearest cluster.

A Few Details . . .

- ◆ How do we decide if a point is “close enough” to a cluster that we will add the point to that cluster?
- ◆ How do we decide whether two compressed sets deserve to be combined into one?

How Close is Close Enough?

- ◆ We need a way to decide whether to put a new point into a cluster.
- ◆ BFR suggest two ways:
 1. The *Mahalanobis distance* is less than a threshold.
 2. Low likelihood of the currently nearest centroid changing.

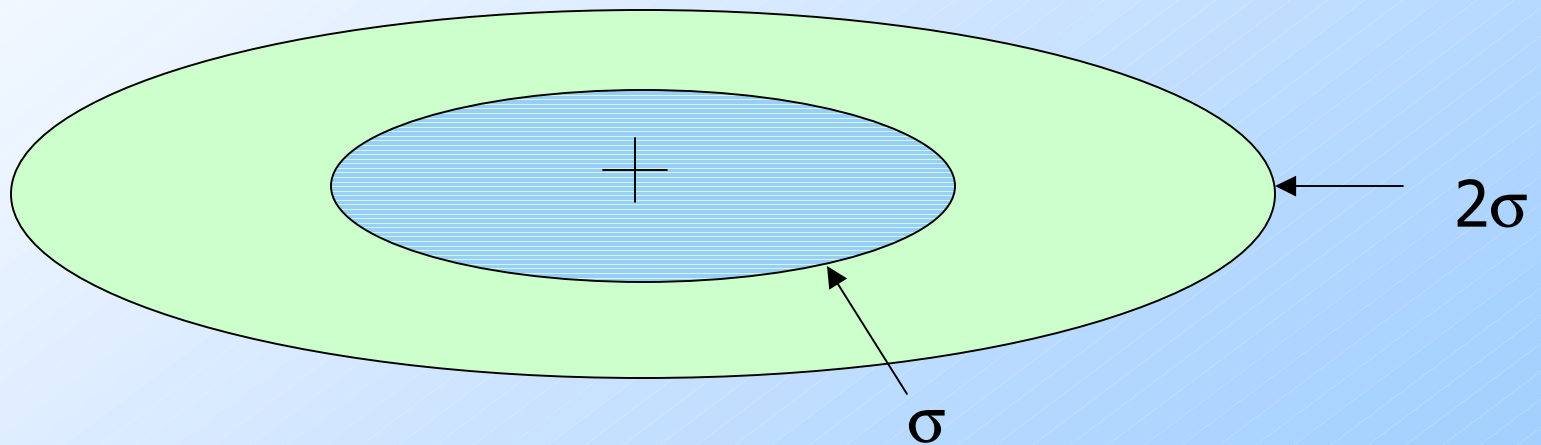
Mahalanobis Distance

- ◆ Normalized Euclidean distance.
- ◆ For point (x_1, \dots, x_k) and centroid (c_1, \dots, c_k) :
 1. Normalize in each dimension: $y_i = (x_i - c_i) / \sigma_i$
 2. Take sum of the squares of the y_i 's.
 3. Take the square root.

Mahalanobis Distance --- (2)

- ◆ If clusters are normally distributed in d dimensions, then after transformation, one standard deviation = \sqrt{d} .
 - ◆ I.e., 70% of the points of the cluster will have a Mahalanobis distance $< \sqrt{d}$.
- ◆ Accept a point for a cluster if its M.D. is $<$ some threshold, e.g. 4 standard deviations.

Picture: Equal M.D. Regions



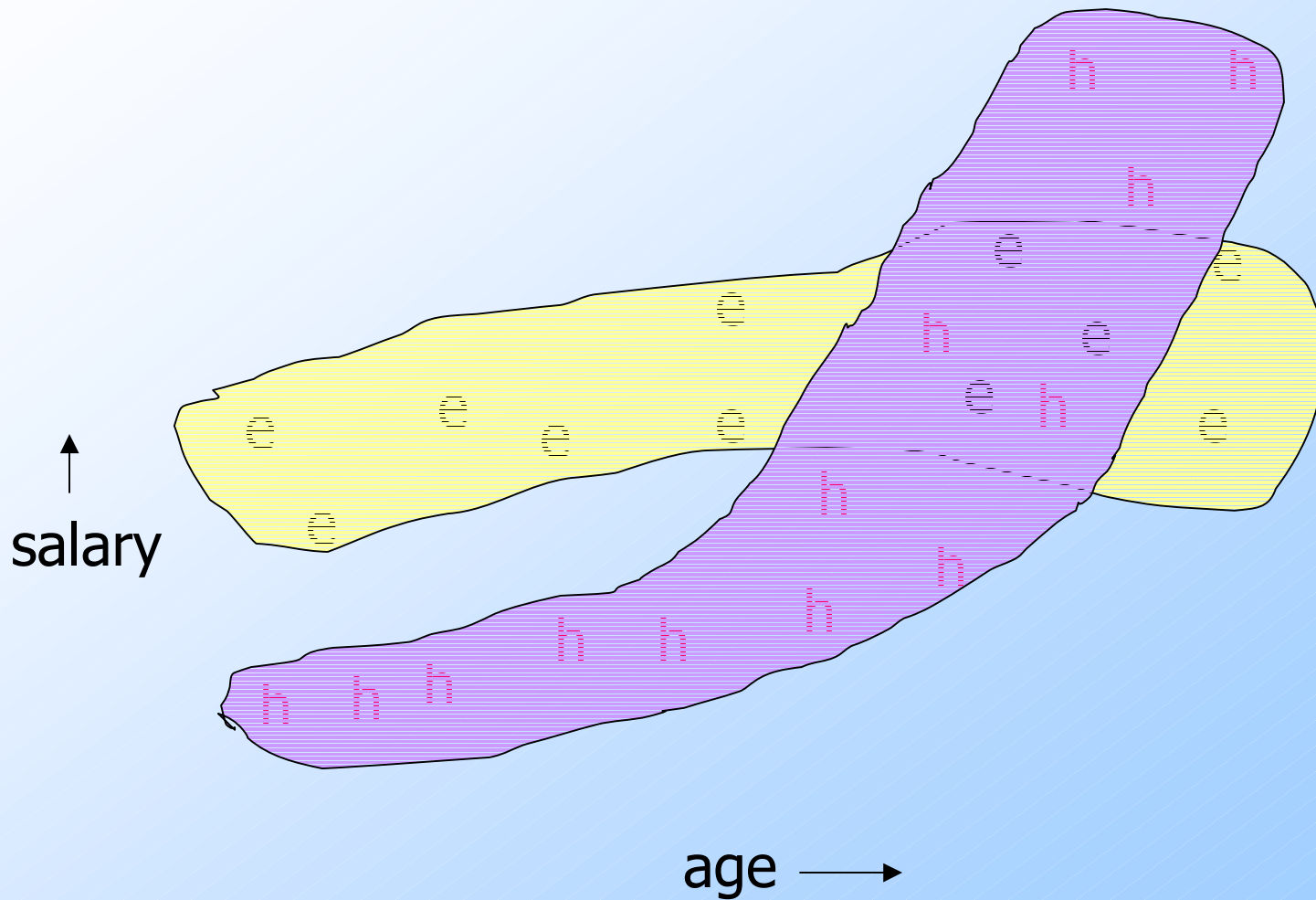
Should Two CS Subclusters Be Combined?

- ◆ Compute the variance of the combined subcluster.
 - ◆ N , SUM, and SUMSQ allow us to make that calculation.
- ◆ Combine if the variance is below some threshold.

The CURE Algorithm

- ◆ Problem with BFR/ k -means:
 - ◆ Assumes clusters are normally distributed in each dimension.
 - ◆ And axes are fixed --- ellipses at an angle are *not* OK.
- ◆ CURE:
 - ◆ Assumes a Euclidean distance.
 - ◆ Allows clusters to assume any shape.

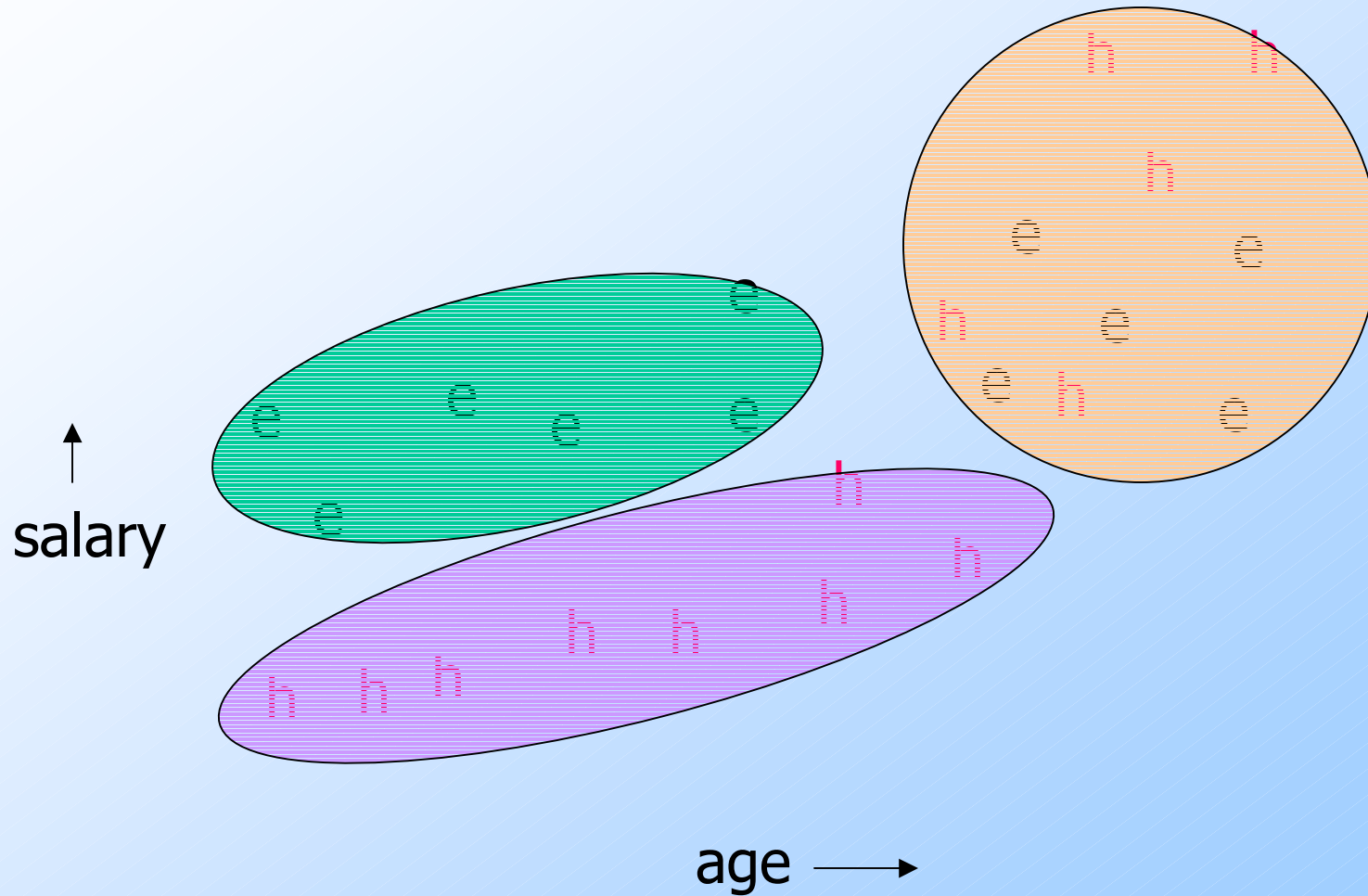
Example: Stanford Faculty Salaries



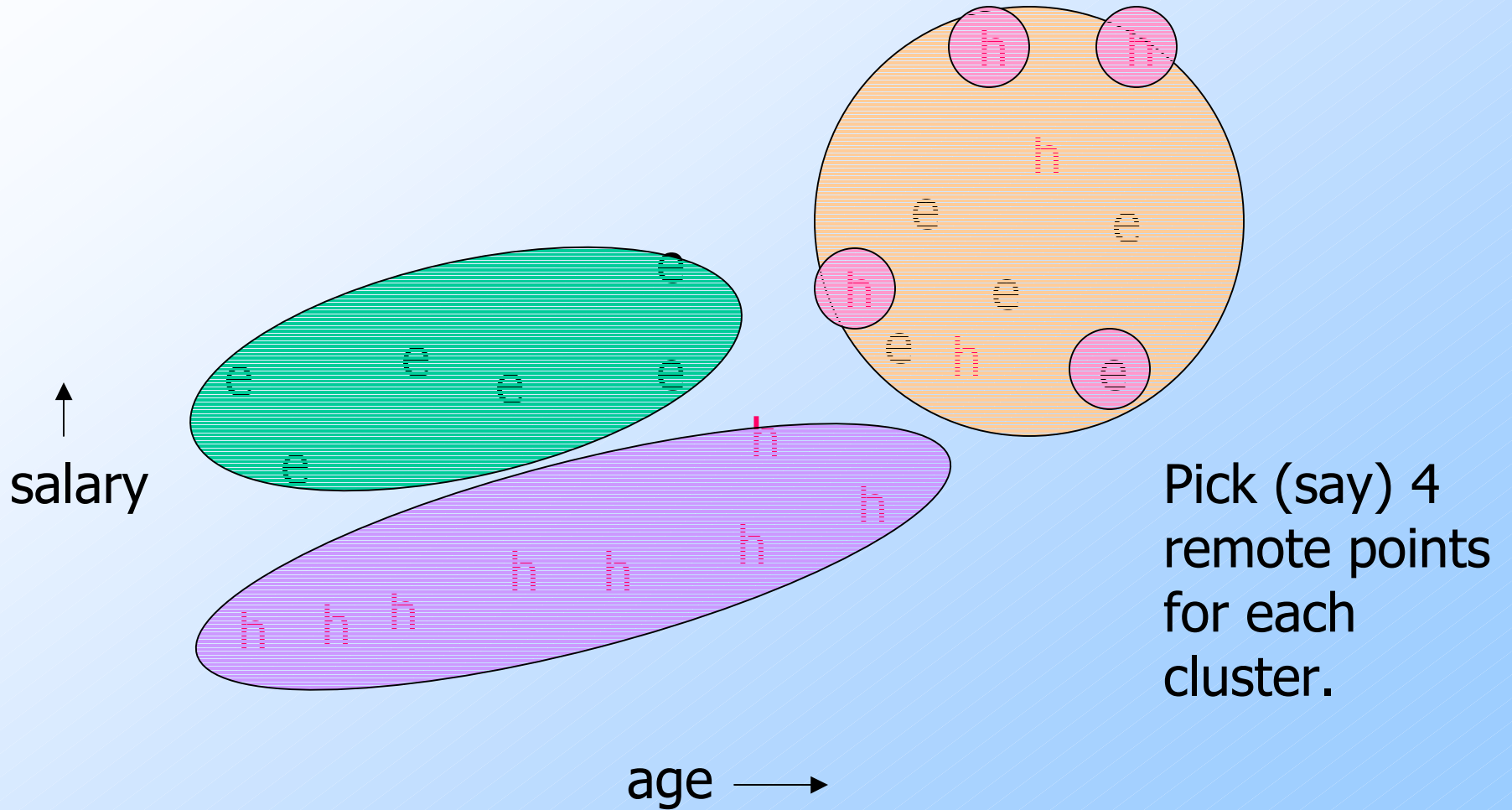
Starting CURE

1. Pick a random sample of points that fit in main memory.
2. Cluster these points hierarchically --- group nearest points/clusters.
3. For each cluster, pick a sample of points, as dispersed as possible.
4. From the sample, pick representatives by moving them (say) 20% toward the centroid of the cluster.

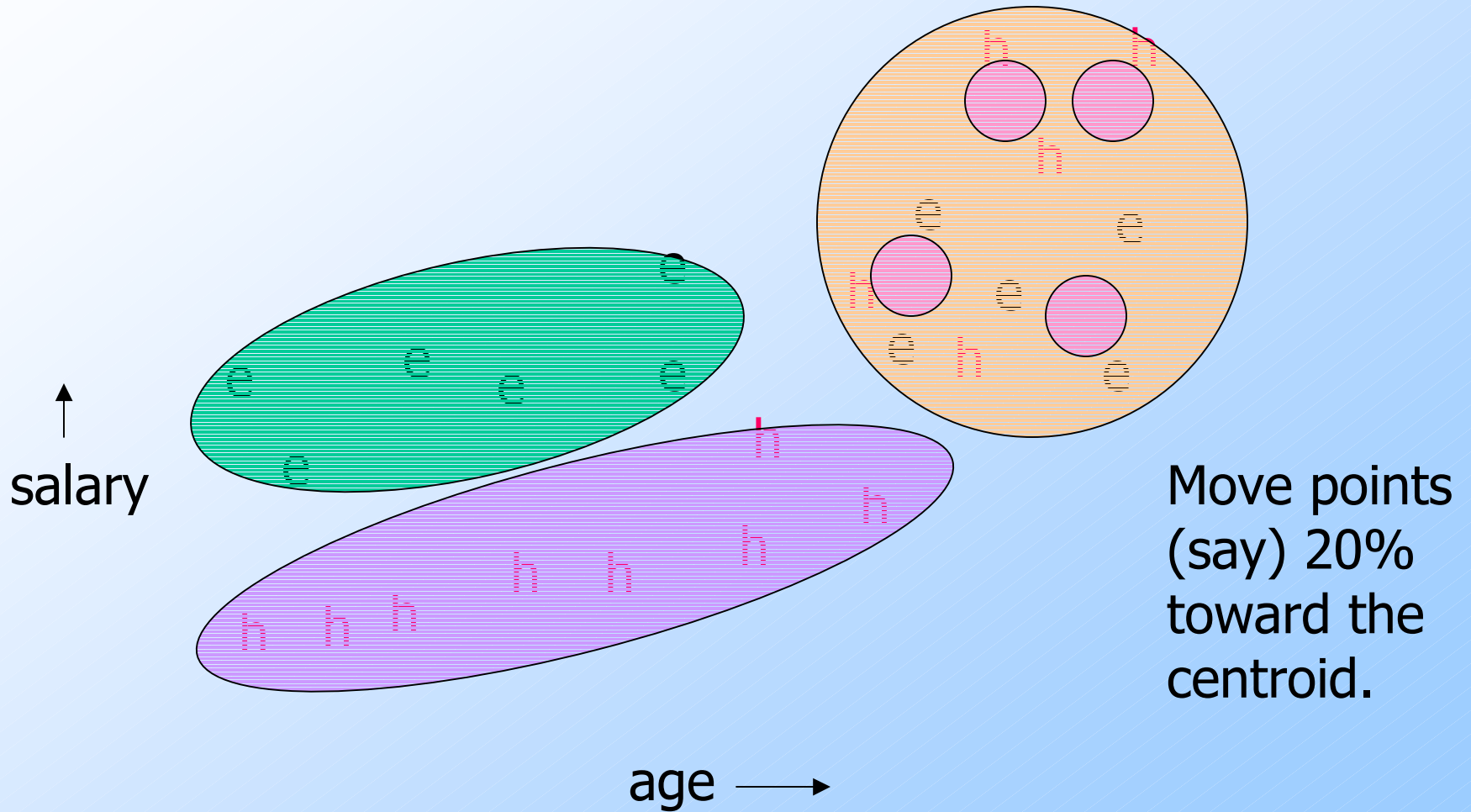
Example: Initial Clusters



Example: Pick Dispersed Points



Example: Pick Dispersed Points



Finishing CURE

- ◆ Now, visit each point p in the data set.
- ◆ Place it in the “closest cluster.”
 - ◆ Normal definition of “closest”: that cluster with the closest (to p) among all the sample points of all the clusters.