# Introduction to Data Mining

## Lecture #15: Clustering-2

**U Kang**
**Seoul National University**

# In This Lecture

- Learn the motivation and advantage of BFR, an extension of K-means to very large data

- Learn the motivation and advantage of CURE, an extension of K-means to clusters of arbitrary shapes
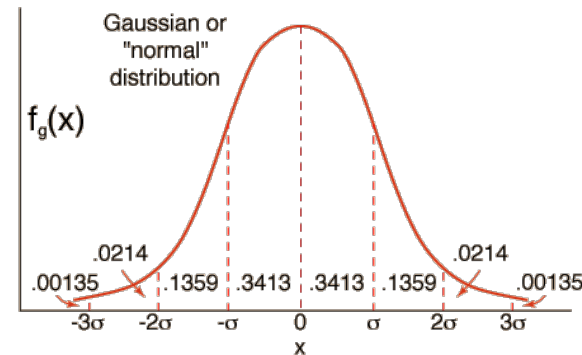
# Outline

➡ ☐ **BFR Algorithm**
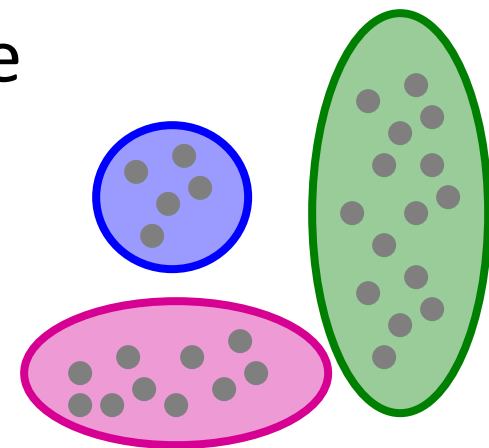
 ☐ CURE Algorithm

*BFR:*
*Extension of k-means to large data*

# BFR Algorithm

- **BFR** [Bradley-Fayyad-Reina] is a variant of *k*-means designed to handle **very large** (disk-resident) data sets

- **Assumes** that clusters are normally distributed around a centroid in a Euclidean space
  - Standard deviations in different dimensions may vary
    - Clusters are axis-aligned ellipses

- **Efficient way to summarize clusters** (want memory required O(clusters) and not O(data))

# BFR Algorithm

- Points are read from disk 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:

  - Take $k$ random points

  - Take a small random sample and cluster optimally

  - 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

**3 sets of points which we keep track of:**

- **Discard set (DS):**

    - Points close enough to a centroid to be summarized
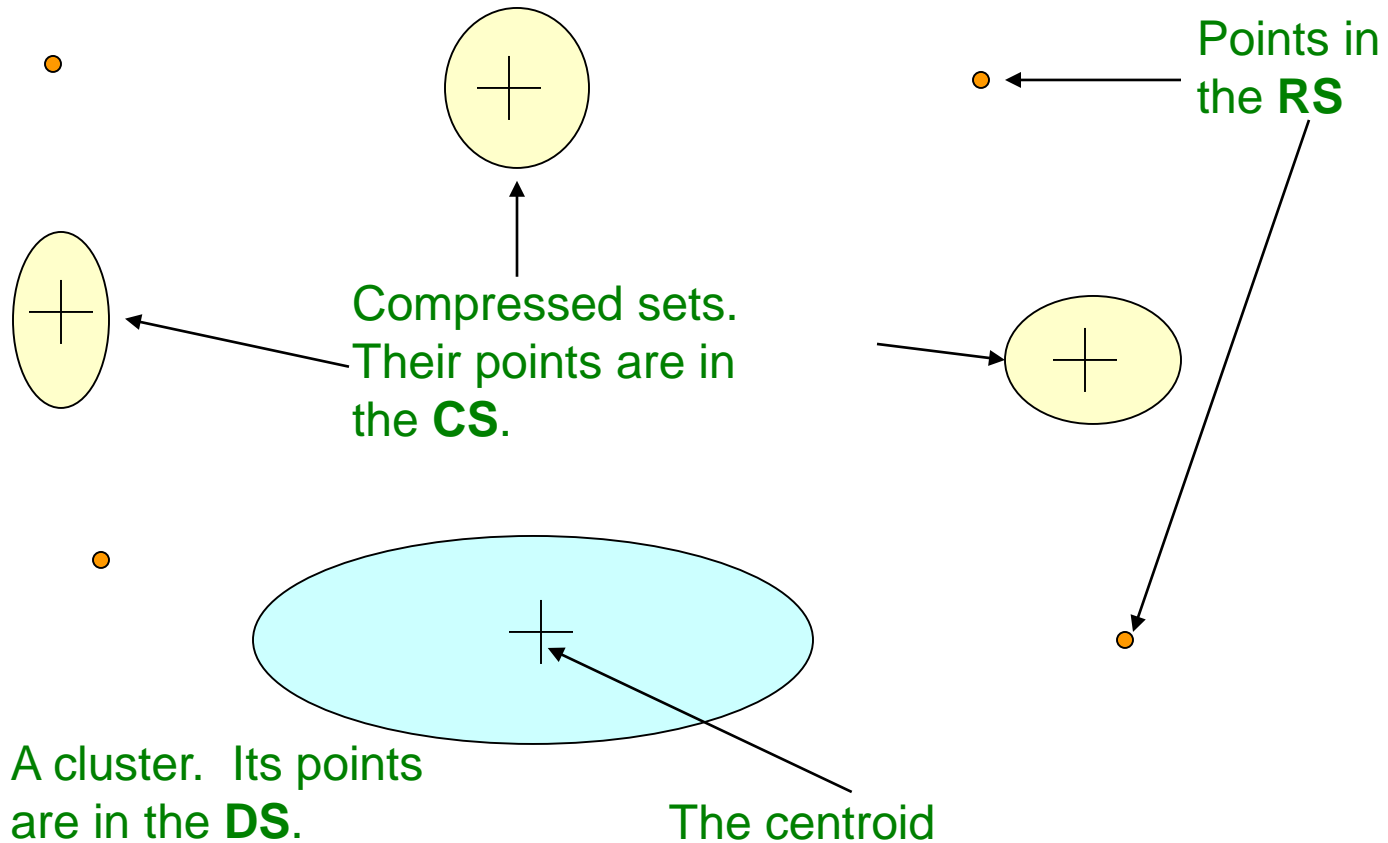
- **Compression set (CS):**

    - Groups of points that are close together but not close to any existing centroid

    - These points are summarized, but not assigned to a cluster

- **Retained set (RS):**

    - Isolated points waiting to be assigned to a compression set
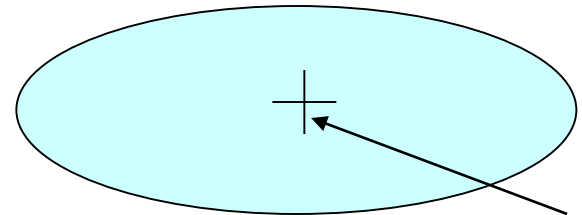
# BFR: "Galaxies" Picture

Points in the **RS**

Compressed sets. Their points are in the **CS**.

A cluster. Its points are in the **DS**.

The centroid

**Discard set (DS):** Close enough to a centroid to be summarized
**Compression set (CS):** Summarized, but not assigned to a cluster
**Retained set (RS):** Isolated points

# Summarizing Sets of Points

**For each cluster, the discard set (DS) is <u>summarized</u> by:**

- The number of points, *N*

- The vector *SUM*, whose $i^{th}$ component is the sum of the coordinates of the points in the $i^{th}$ dimension

- The vector *SUMSQ*: $i^{th}$ component = sum of squares of coordinates in $i^{th}$ dimension

A cluster.
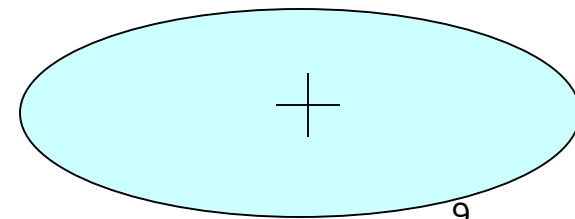All its points are in the **DS**.

The centroid

# Summarizing Points: Comments

- **$2d + 1$** values represent any size cluster
  - ❑ **$d$** = number of dimensions

- Average in **each dimension** (<span style="color:magenta">the centroid</span>) can be calculated as **$SUM_i / N$**
  - ❑ **$SUM_i$** = $i^{th}$ component of SUM

- Variance of a cluster's discard set in dimension $i$ is: **$(SUMSQ_i / N) - (SUM_i / N)^2$**
  - ❑ And standard deviation is the square root of that

- **Next step: Actual clustering**

**Note:** Removing the "axis-aligned" clusters assumption would require storing full covariance matrix to summarize the cluster. So, instead of **SUMSQ** being a **$d$**-dim vector, it would be a **$d \times d$** matrix, which is too big!

U Kang

9

# The "Memory-Load" of Points

**Processing the "Memory-Load" of points (1):**

- **1)** Find those points that are "**sufficiently close**" to a cluster centroid and add those points to that cluster and the **DS**

    - ❑ These points are so close to the centroid that they can be summarized and then discarded

- **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**

**Discard set (DS):** Close enough to a centroid to be summarized.
**Compression set (CS):** Summarized, but not assigned to a cluster
**Retained set (RS):** Isolated points

# The "Memory-Load" of Points

**Processing the "Memory-Load" of points (2):**

- **3) DS set:** Adjust statistics of the clusters to account for the new points
  - Update *N*s, *SUM*s, *SUMSQ*s

- **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

**Discard set (DS):**  Close enough to a centroid to be summarized.
**Compression set (CS):**  Summarized, but not assigned to a cluster
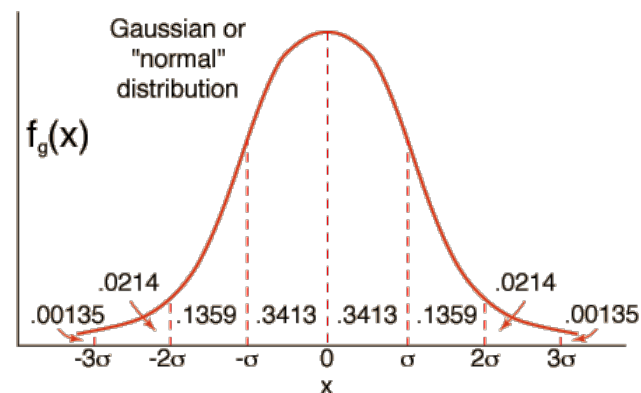**Retained set (RS):** Isolated points

# A Few Details…

- **Q1) How do we decide if a point is "close enough" to a cluster that we will add the point to that cluster?**

- **Q2) How do we decide whether two compressed sets (CS) deserve to be combined into one?**

# How Close is Close Enough?

- **Q1) We need a way to decide whether to put a new point into a cluster (and discard)**

- **BFR suggests two ways:**
  - High likelihood of the point belonging to currently nearest centroid (and, the point far from all other centroids)
  - The **Mahalanobis distance** is small (< t)



U Kang

# Mahalanobis Distance

- **Normalized Euclidean distance from centroid**

- For point $(x_1, ..., x_d)$ and centroid $(c_1, ..., c_d)$
    1. Normalize in each dimension: $y_i = (x_i - c_i) / \sigma_i$
    2. Take sum of the squares of the $y_i$
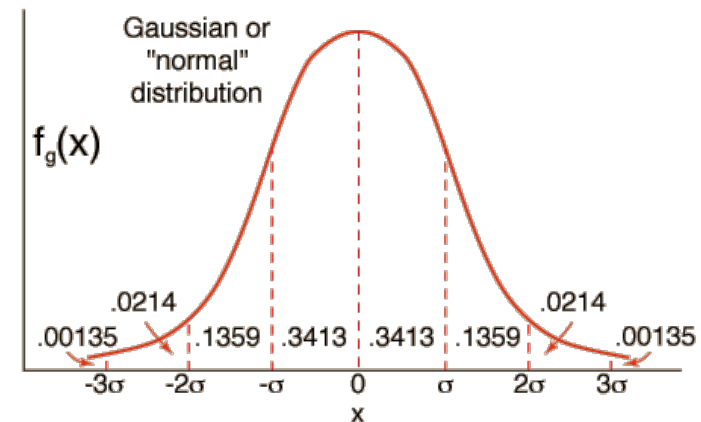    3. Take the square root

$$d(x, c) = \sqrt{\sum_{i=1}^{d} \left( \frac{x_i - c_i}{\sigma_i} \right)^2}$$

$\sigma_i$ ... standard deviation of points in the cluster in the $i$th dimension

# **Mahalanobis Distance**

- If clusters are normally distributed in **$d$** dimensions, then after transformation, one standard deviation **=** $\sqrt{d}$

- Accept a point for a cluster if its M.D. is **<** t  (a parameter), e.g. **2** standard deviations
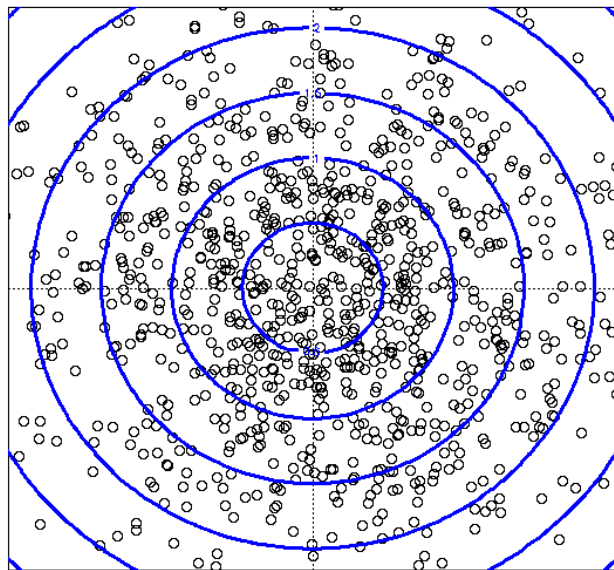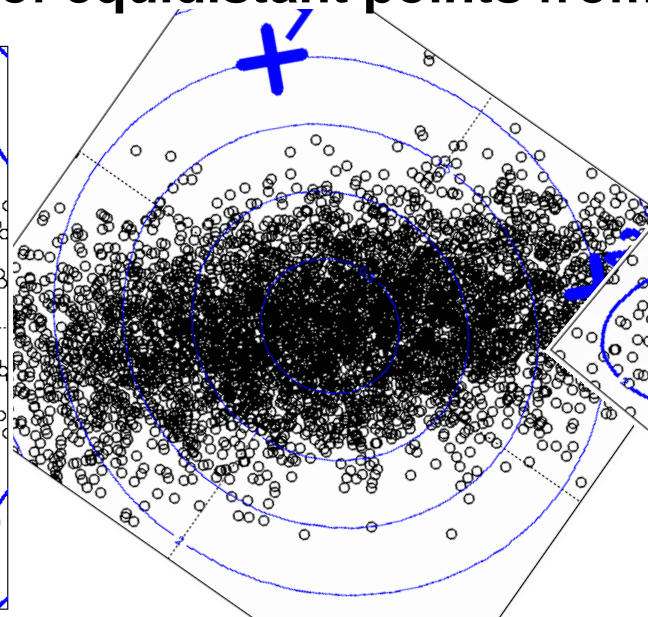


Gaussian or "normal" distribution

$f_g(x)$

.0214

.00135     .1359   .3413   .3413   .1359     .00135

-3σ     -2σ     -σ     0     σ     2σ     3σ

x

# Picture: Equal M.D. Regions
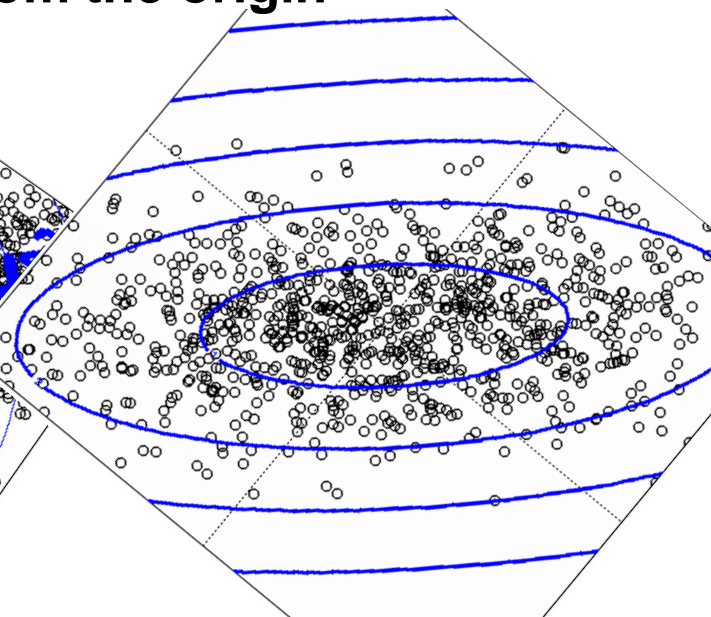
- **Euclidean vs. Mahalanobis distance**

**Contours of equidistant points from the origin**



**Uniformly distributed points, Euclidean distance**

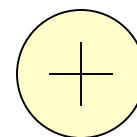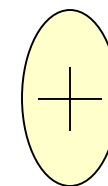**Normally distributed points, Euclidean distance**

**Normally distributed points, Mahalanobis distance**

# Should 2 CS clusters be combined?

**Q2) Should 2 CS subclusters be combined?**

- Compute the variance of the combined subcluster
    - *N*, *SUM*, and *SUMSQ* allow us to make that calculation quickly
- Combine if the combined variance is small (< s)

# Outline

☑ BFR Algorithm

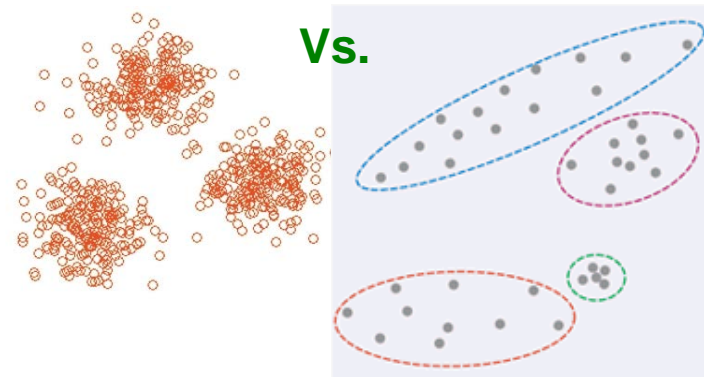➡ ☐ **CURE Algorithm**

*CURE:*
*Extension of k-means to clusters*
*of arbitrary shapes*

# The CURE Algorithm

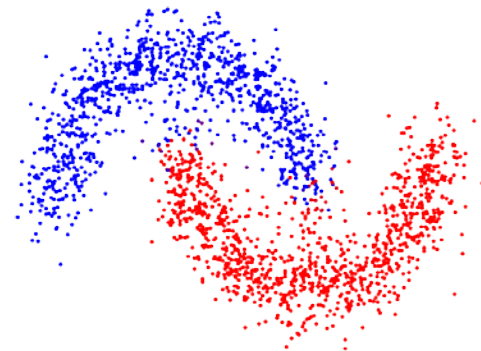- **Problem with BFR/*k*-means:**

  - Assumes clusters are normally distributed in each dimension

  - And axes are fixed – ellipses at an angle are *not OK*

  Vs.

- **CURE (Clustering Using REpresentatives):**

  - Assumes a Euclidean distance

  - Allows clusters to assume any shape

  - **Uses a collection of representative points to represent clusters**

# Starting CURE

**2 Pass algorithm. Pass 1:**

- **1) Pick a random sample of points that fit in main memory**

- **2) Initial clusters:**
  - ❑ Cluster these points hierarchically – group nearest points/clusters

- **3) Pick representative points:**
  - ❑ For each cluster, pick a sample of points, as dispersed as possible
  - ❑ From the sample, pick representatives by moving them (say) 20% toward the centroid of the cluster
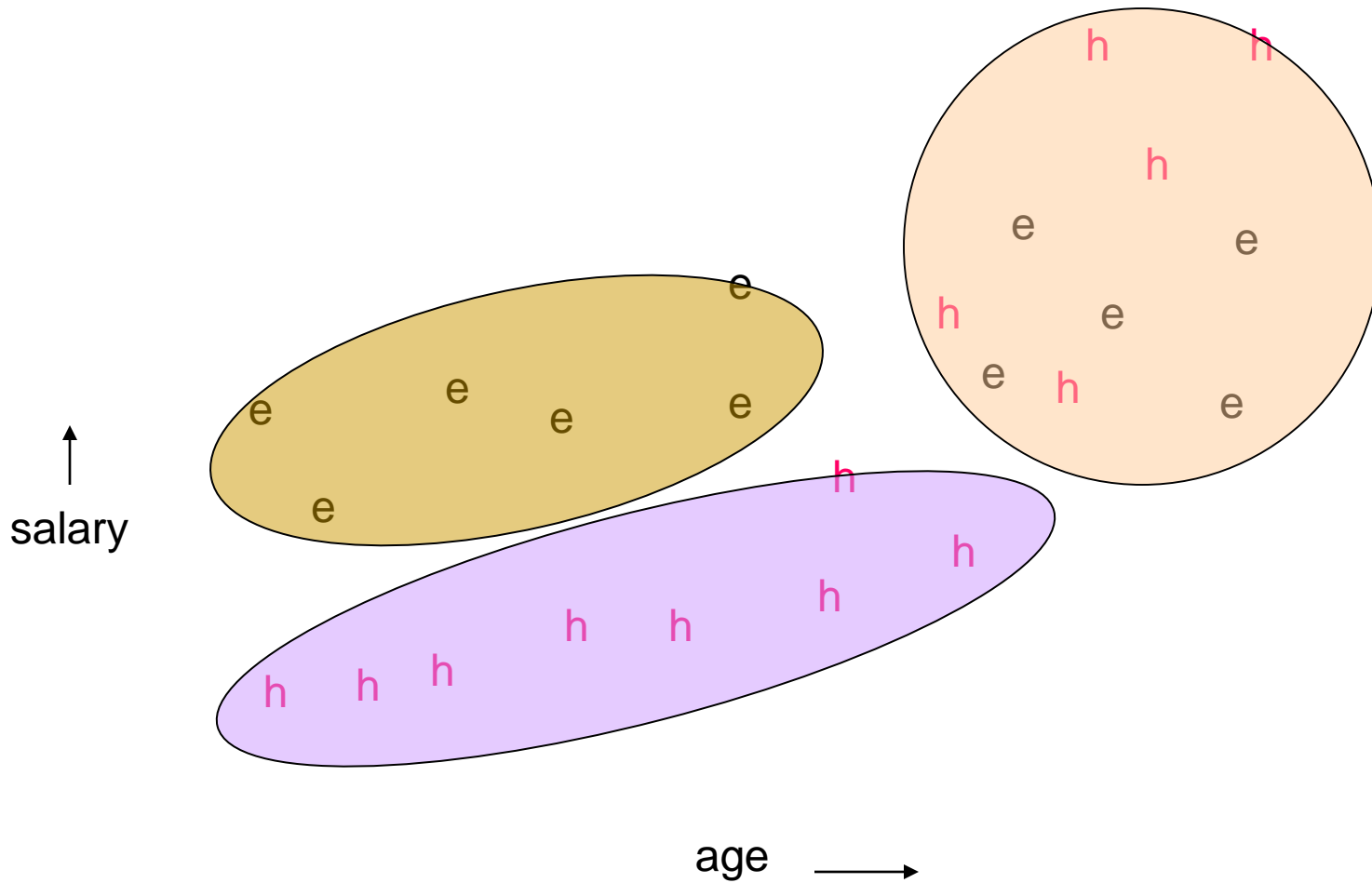
# Starting CURE

**2 Pass algorithm. Pass 1:**

- **4) Merge clusters**
    - Merge two clusters that are sufficiently close (<t)
        - Cluster distance: minimum distance of representative points
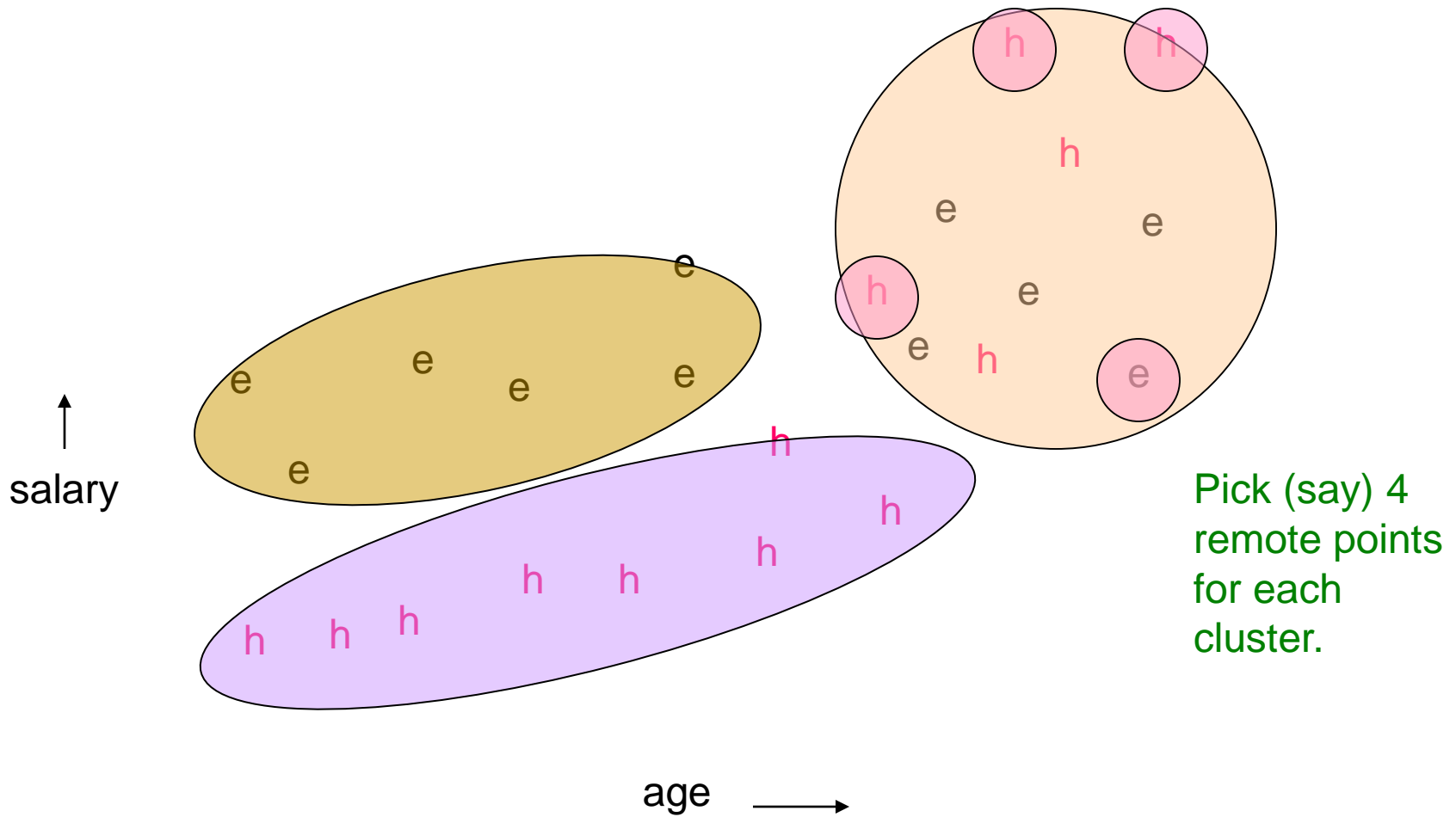    - Repeat, until there are no more sufficiently close clusters
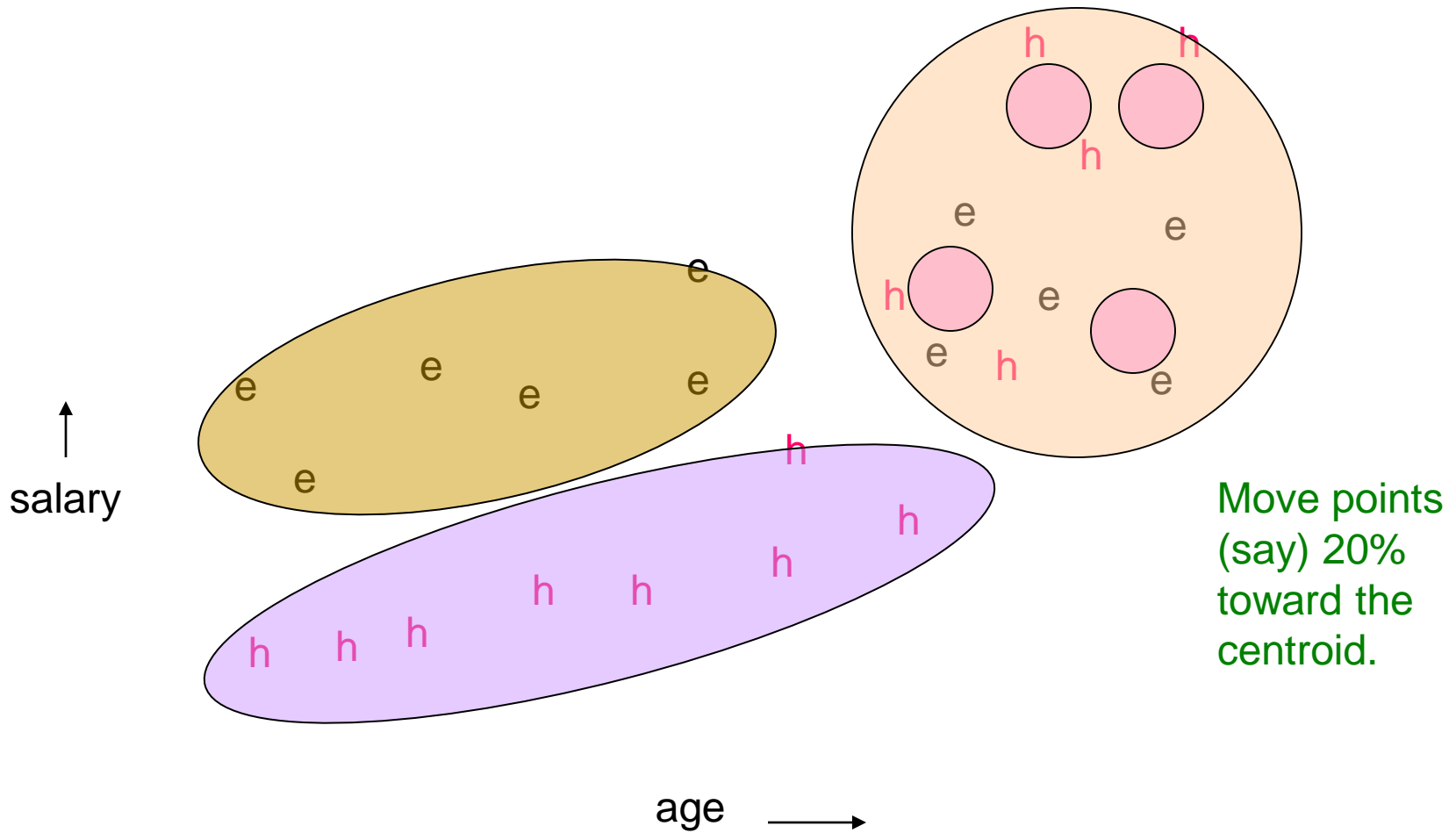
# Example: Initial Clusters



salary

age

# Example: Pick Dispersed Points



salary

age

Pick (say) 4 remote points for each cluster.

# Example: Pick Dispersed Points

salary ↑

age →

Move points (say) 20% toward the centroid.
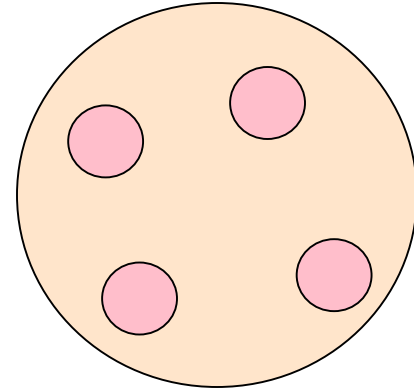
# Finishing CURE

## Pass 2:

- Now, rescan the whole dataset and visit each point **p** in the data set

**p**

- **Place it in the "closest cluster"**

  - Normal definition of "closest":
    Find the closest representative to **p** and assign it to representative's cluster

# Summary: Clustering

- **Clustering:** Given a **set of points**, with a notion of **distance** between points, **group the points** into some number of *clusters*

- **Algorithms:**

  - Agglomerative **hierarchical clustering**:
    - Centroid and clustroid

  - *k*-means:
    - Initialization, picking *k*

  - **BFR**

  - **CURE**

# Questions?