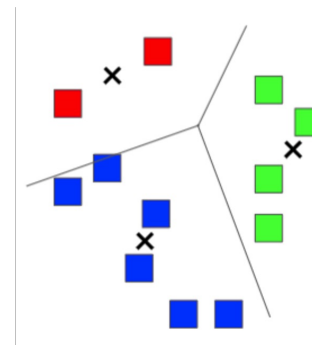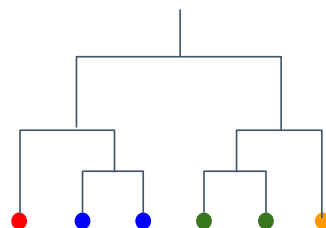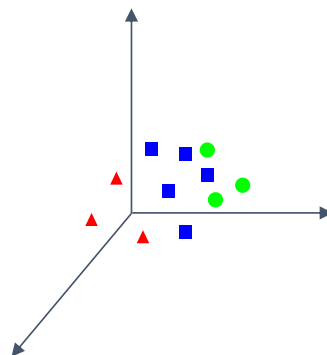# Clustering 2

## EE412: Foundation of Big Data Analytics

# Announcements

- Homeworks
  - HW1 (due: 10/05)
  - HW2 (will be posted at 10/10)

# Recap

- Clustering
  - Curse of dimensionality
  - Clustering strategies
- Hierarchical Clustering
  - Euclidean vs. non-Euclidean
  - Centroids vs. clustroids
- $k$-means Clustering
  - $k$-means++
  - Selection of $k$

# Outline

1. **<u>BFR Algorithm</u>**
2. BFR Algorithm: Process
3. CURE Algorithm
4. GRGPF Algorithm
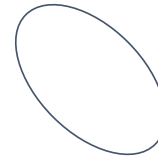
# BFR (Bradley, Fayyad, and Reina) Algorithm

- Variant of $k$-means for very large (disk-resident) data sets
- Assumes that clusters are normally distributed in a Euclidean space
  - Standard deviations in different dimensions may vary
  - Clusters are axis-aligned ellipses
- **Goal:** Find cluster centroids
  - Point assignment can be done in a second pass through the data



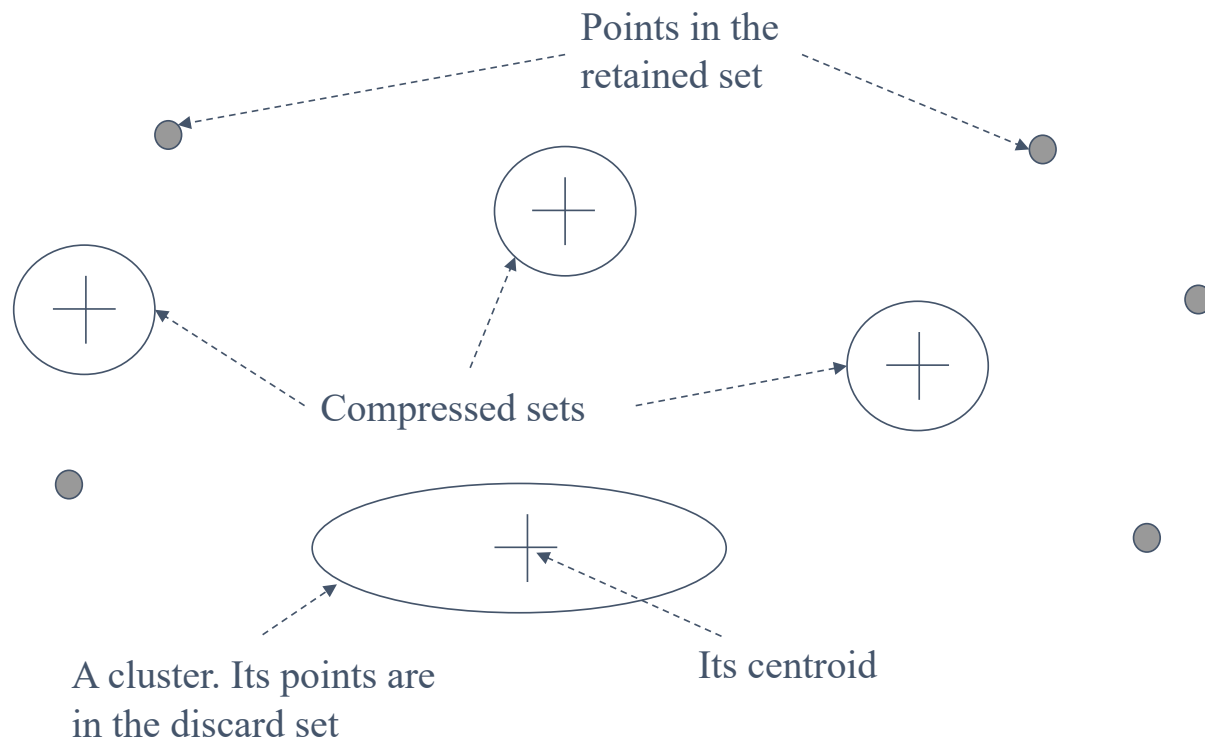OK                    OK                    Not OK

# Main Idea

- **Idea:** Keep summary statistics of groups of points
  - Points are read from disk one main-memory-full at a time
  - Most points from previous memory loads are summarized
  - Changes memory requirement from $O(\text{data})$ to $O(\text{clusters})$
- **3 sets:** Discard set, Compressed set, and Retained set

# Three Classes of Points

*no longer use*

- **Discard set**
  *summarize 하기 충분히 가까운 애들.*
  - Points close enough to a centroid to be **summarized**
- **Compressed set**
  *=> 서로 뭉칠만큼 가깝지만, existing centroid랑 가까지는 않은것. (mini-clustered)*
  - Groups of points that are close together but not to any existing centroid
  - These points are **summarized** but not assigned to a cluster
- **Retained set**
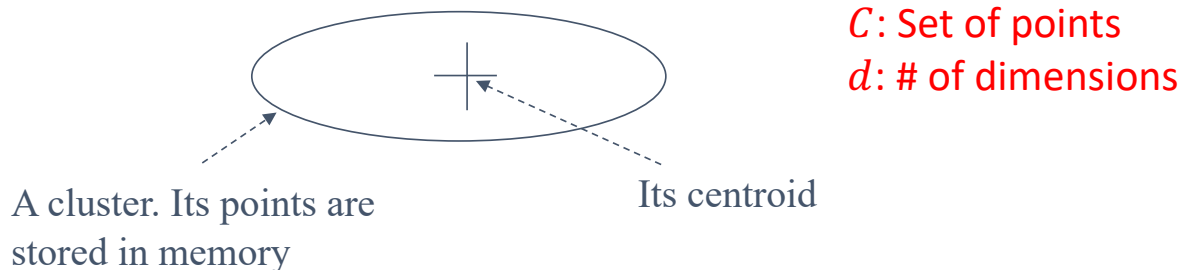  - Isolated points waiting to be assigned to a compressed set

# Cluster Visualization



Points in the retained set

Compressed sets

A cluster. Its points are in the discard set

Its centroid

# Summarizing Sets of Points

처음 custom에 어떻게 하기 ↗ num of dimensions

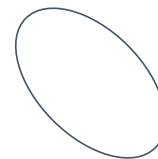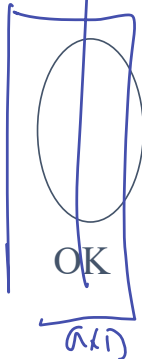- Discard or compressed set is **summarized** by $2d + 1$ values
  - 1개 The number of points, $N$
  - 1개 The vector SUM, where $\text{SUM}_i$ is the sum of the coordinates of all points
  - d개 The vector SUMSQ, where $\text{SUMSQ}_i$ is the sum of squared coordinates
    - That is, $\text{SUM}_i = \sum_{k \in C} x_{ki}$ and $\text{SUMSQ}_i = \sum_{k \in C} x_{ki}^2$

SUM₂ =) sum of all data of dimension 2



$C$: Set of points
$d$: # of dimensions

A cluster. Its points are stored in memory

Its centroid

# Summarizing Sets of Points

- We can compute the average and variance of a cluster
    - Average in dimension $i$ (i.e., the **centroid**) is $\text{SUM}_i/N$
    - Variance in dimension $i$ is $\text{SUMSQ}_i/N - (\text{SUM}_i/N)^2$
        - Because $\text{Var}(X) = E[X^2] - E[X]^2$
    - Standard deviation is the square root of variance

=) elipse를 나타내는 enodoh info

- This is based on the assumption of "axis-aligned" clusters
    - Without it, we need to store the $d \times d$ covariance matrix



OK

axis

OK

Not OK

axis

# Benefits of the Representation

- Easy to add a new point to a cluster
  - Increase $N$ by 1
  - Add the vector to SUM
  - Add the squares of components to SUMSQ
- Also easy to combine two sets
  - Add corresponding values of $N$, SUM, SUMSQ

# Pop Quiz

- Represent the cluster of points (5, 1), (6, −2), and (7, 0)
  - N = ?    $3$
  - SUM = ?    $(18, -1]$
  - SUMSQ = ?    $25+36+49 = 110$    $1+4 = 5$
- Compute the variance of the first dimension
  - Variance = ?

# Outline

Compressed set =) discard set로 정입시키는 없는다.
그냥 개별히 위슬뿔

# Overview of the Algorithm

1. Initialize $k$ clusters/centroids (as in $k$-means)
2. **for** each chunk **in** a data file
3.     **for** each point **in** the chunk
4.         Assign it to a cluster if it is sufficiently close to the cluster
5.     Cluster the remaining points, creating new clusters
6.     Try to merge new clusters with any of the existing clusters

여기로 CS와 합칠지 고려

가깝지 않으면 그냥 remaining으로 넘기고
remaining들 끼리 모아서 CS로 산출자. 이 CS를 아주 CS와 합칠지 고려.
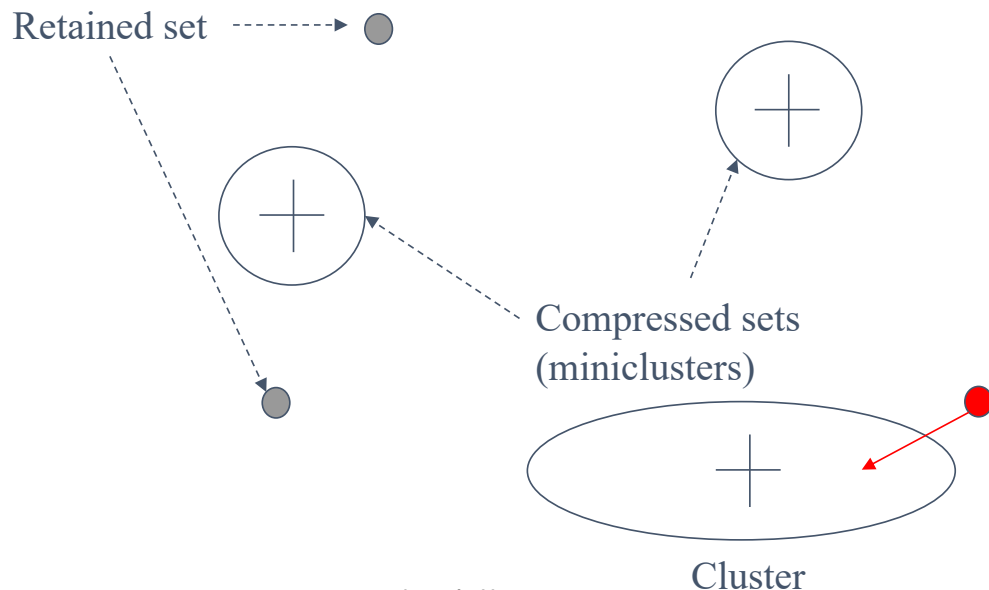여기에 포함되지 않은 애들은 RS임

# Selection of the Initial Centroids

→really important.

- The $k$ initial centroids can be selected as in $k$-means
  - Take $k$ random points (not a good way, try it)
  - Take a small random sample, cluster it, and use the centroids
  - Take a sample; Pick a random point, and then $k-1$ more points
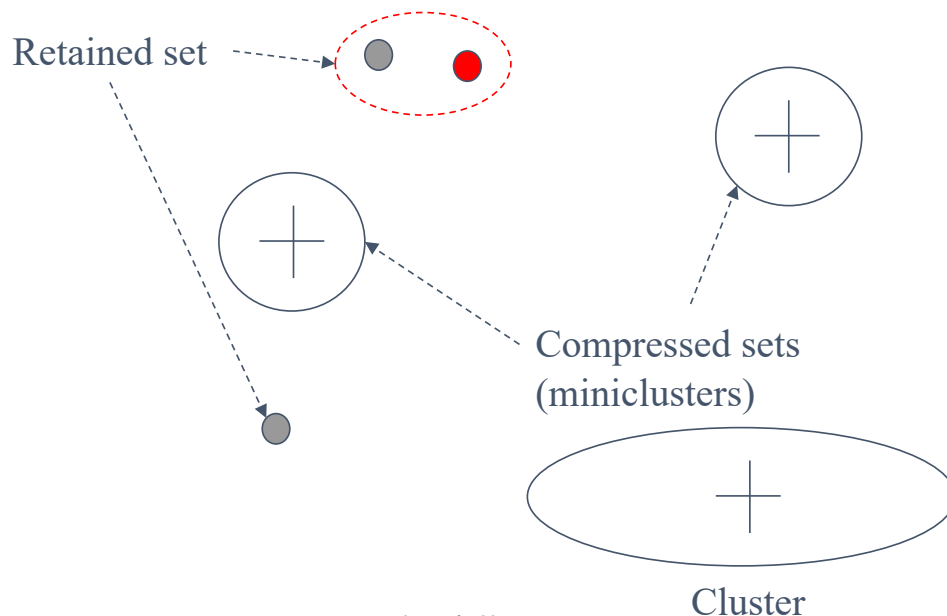    - Each as far from the previously selected points as possible

# Processing a Chunk of Points (1/5)

- All points that are "sufficiently close" to the centroid of a cluster:
    - Added to that cluster



Retained set

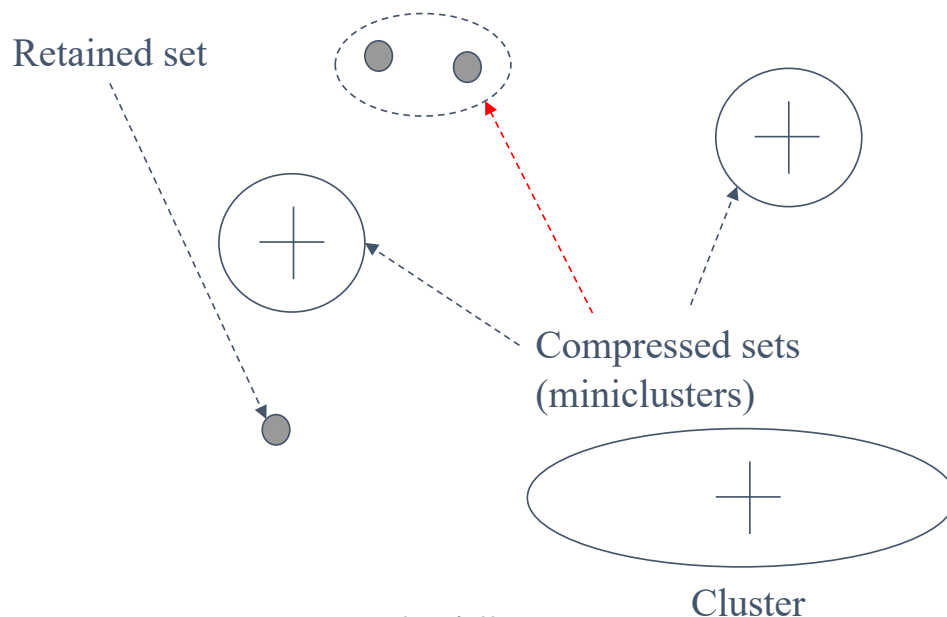Compressed sets
(miniclusters)

Cluster

# Processing a Chunk of Points (2/5)

- The points that are not sufficiently close to any centroid:
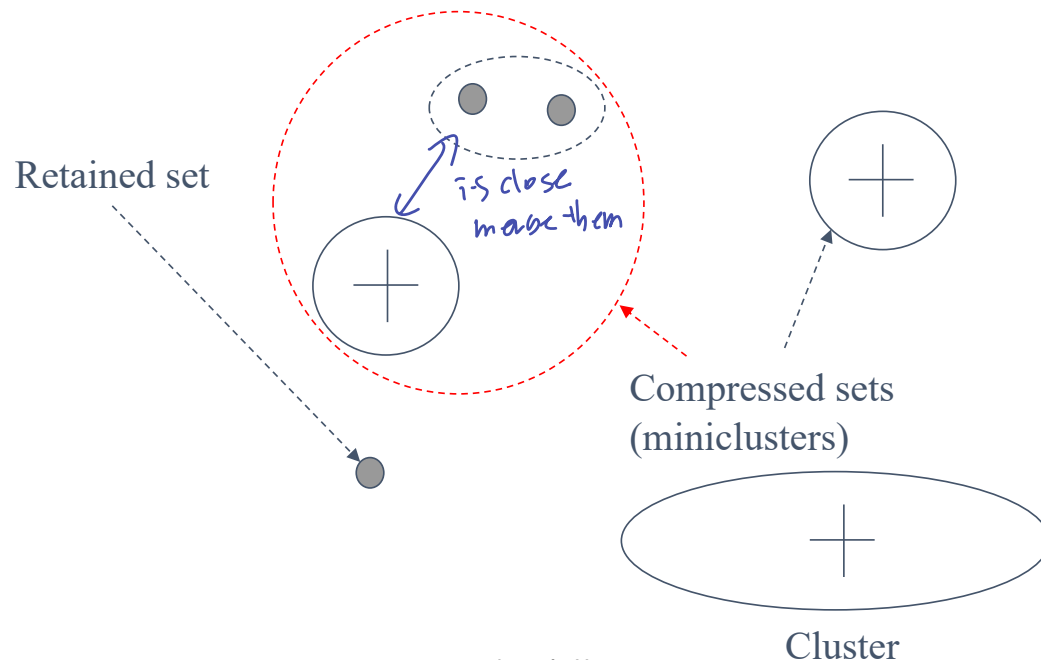  - Clustered with the points in the retained set using any clustering algorithm

Retained set

Compressed sets
(miniclusters)

Cluster

# Processing a Chunk of Points (3/5)

- Clusters of ≥ 2 points are summarized and become miniclusters
- Singleton clusters remain in the retained set
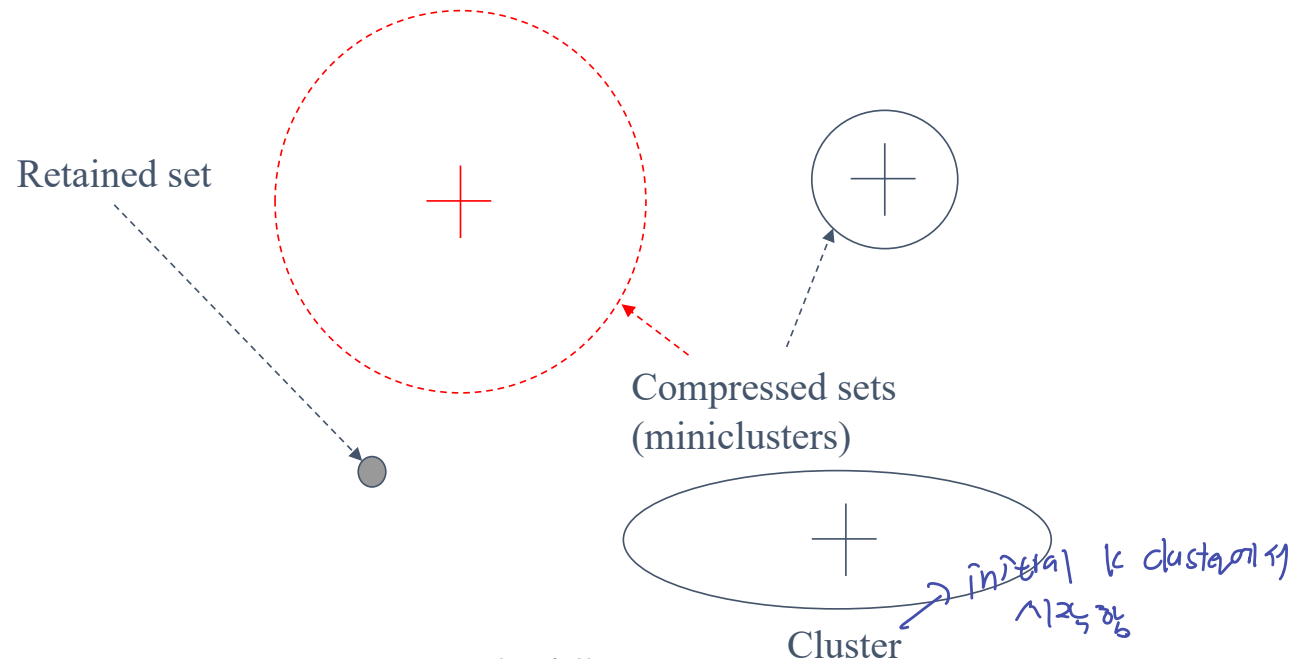
Retained set

Compressed sets
(miniclusters)

Cluster

# Processing a Chunk of Points (4/5)

- Cluster the new miniclusters with the old miniclusters



Retained set

is close
marge them

Compressed sets
(miniclusters)

Cluster

# Processing a Chunk of Points (5/5)

- Points assigned to a cluster or a minicluster are written to disk



Retained set

Compressed sets
(miniclusters)

Cluster

초 initial k cluster에서
시작함

# After Processing All Chunks

- At the last round, what to do with compressed and retained sets?
- **Option 1:** Treat them as outliers and never cluster them  *strong way*
- **Option 2:** Assign each of them to the nearest cluster  그냥 버려두기. (remove all outliers)
  - For the compressed set, combine each minicluster with the nearest cluster

# How Close is Close Enough?

*얼마나 가까워야 하는가.*

- Need a way to decide whether to put a new point into a cluster

  *→ basically euclidian.*

- BFR compares the **Mahalanobis distance** with a threshold

  - Exploit the assumption that points are normally distributed
  - Euclidean distance from the centroid $c$ normalized by standard dev. $\sigma_i$
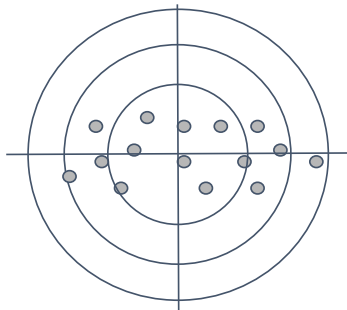
  *assumption (ellips shape)*

  - **Definition:**

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

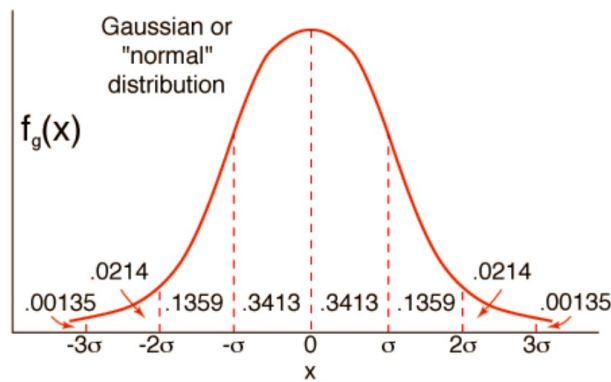*std를 고려해서 진짜 분포까지의 거리를 더 알맞게.*

Euclidean distance:

Mahalanobis distance:

# Assigning a Point to a Cluster

- Choose a cluster whose centroid has the least Mahalanobis distance
  - *likelihood*
- Add a point $p$ if the distance is less than a threshold
  - E.g., if threshold = 4
  - Then, Pr(value being 4 standard deviations from mean) $< 10^{-6}$



Source: Stanford CS246 (2022)

# When to Merge Two Clusters?

- Compute the variance of the combined subcluster
- Combine if the combined variance is below some threshold
- **Many alternatives:** E.g., considering density

# Outline

1. BFR Algorithm
2. BFR Algorithm: Process
3. **CURE Algorithm**
4. GRGPF Algorithm

Assumption 이나 연결계 Strong 하려면
scale을 위해 sacrifice 하는거임.
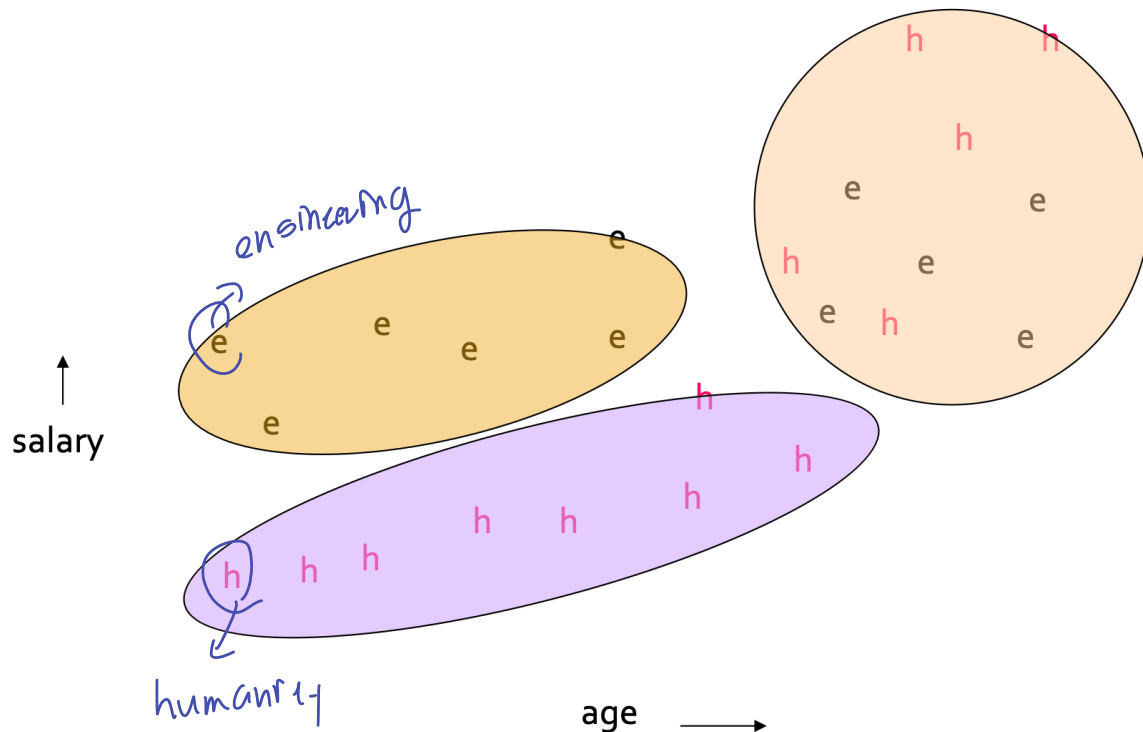(ex, k-initial point 잡는거 super important함)

# CURE (Clustering Using REpresentatives)

- CURE is a 2-pass algorithm for large disk-resident data
- No assumption about the shape of clusters
  - No need to be normally distributed in each dimension   *How is it possible?*
- Uses a collection of representative points to represent clusters   *=) not statistics*
  - No centroids   *more naive way, but effective*
- Assumes a Euclidean distance, with $k$ (# of clusters) given

# CURE: Pass 1

1.  Pick a random sample of data   *random pick*

2.  Cluster them in main memory using hierarchical clustering
    - Merge two clusters when they have close pairs of points   *↳ remove shape assumption ✗ concept of centroids*

3.  Pick representative points from each cluster
    - For each cluster, pick a sample of points, as dispersed as possible   *far as possible*
    - Move them a fraction of distance, e.g. 20%, toward the centroid

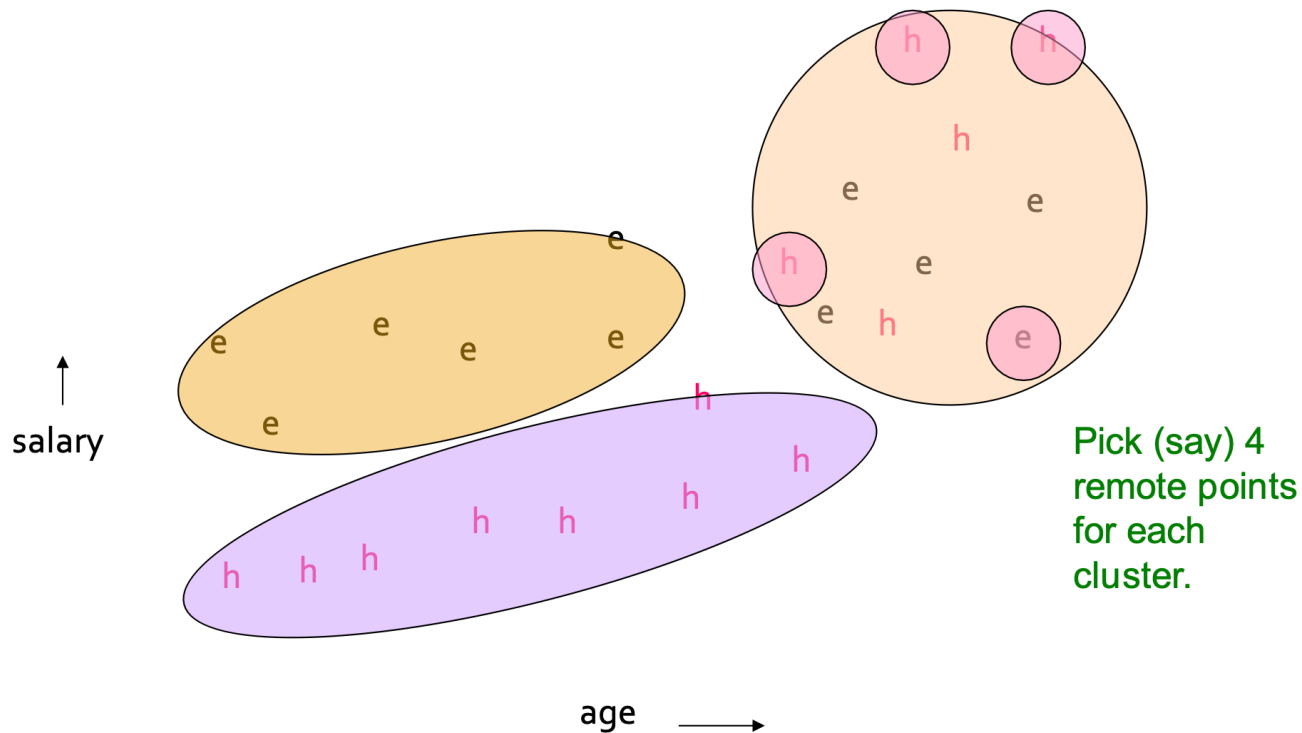4.  Merge clusters with the closest pair of representatives

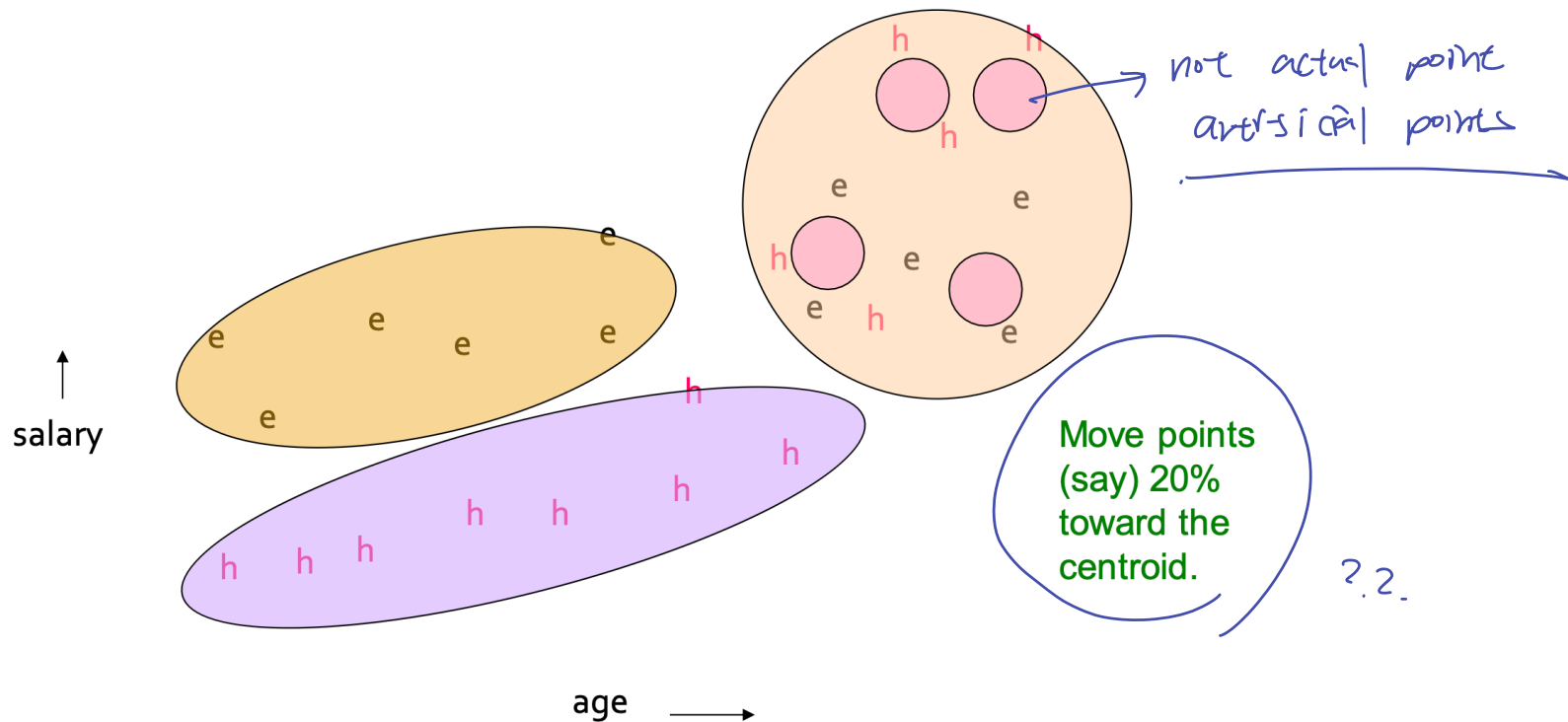# Example: Picking Dispersed Points



Source: Stanford CS246 (2022)

# Example: Picking Dispersed Points



Pick (say) 4 remote points for each cluster.

salary

age

Source: Stanford CS246 (2022)

# Example: Picking Dispersed Points



not actual point
artificial points

salary

age

Move points (say) 20% toward the centroid.

?.?.

Source: Stanford CS246 (2022)

# CURE: Pass 2

1. Rescan the whole dataset and visit each point $p$ in the data set

2. Place it in the "closest cluster"
   - Find the closest representative point to $p$
   - Assign $p$ to the representative's cluster

# Why to 20% Move Inward?

- Suppose that initial sample is large enough
- Some of the representatives will be on the boundary of clusters
  - Moving them towards the centroid
- Large, dispersed clusters will shrink more than small, dense ones
- As a result, the algorithm favors a small, dense cluster



기본적으로 가장 진행되면
dispense 의 단적이 더 점들을
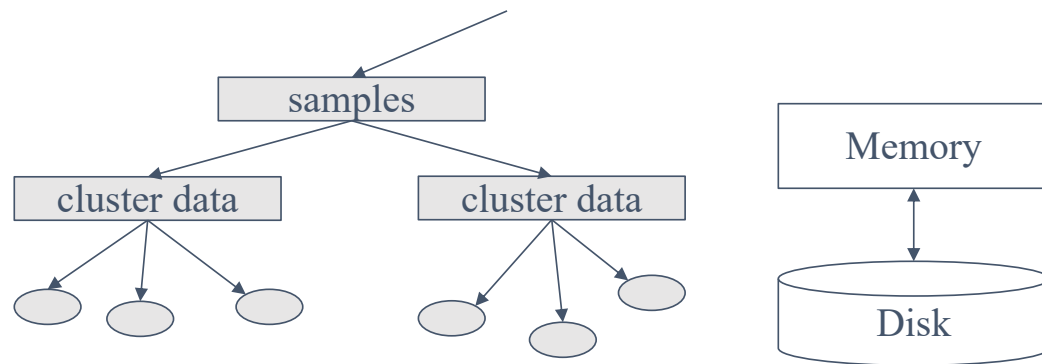질 많을 한다.

→ 결과적에
small 인 객들이
불거책

Source: Stanford CS246 (2022)

# Outline

1. BFR Algorithm
2. BFR Algorithm: Process
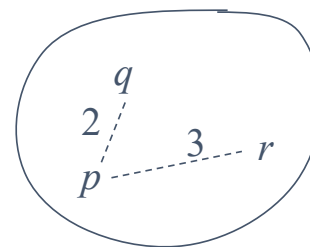3. CURE Algorithm
4. **GRGPF Algorithm**

# GRGPF Algorithm

- Does not require a Euclidean space
- Represents clusters by **well-chosen sample points** in memory
- Organizes clusters hierarchically, as a tree (not covered today)
  - New point is assigned to a cluster by passing it down the tree

# Cluster Representation

- How to represent (or summarize) a cluster in GRGPF
  - The number of points, $N$
  - The clustroid $c$   *x assume   euclidian*
  - The **rowsum** of the clustroid
    - Sum of the squares of the distances from $p$ to each point in the cluster
  - The $k$ points that are **nearest** to the clustroids, and their rowsums
  - The $k$ points that are **furthest** from the clustroids, and their rowsums

*rowsum* of point $p$ in cluster $C = \sum_{c \in C} d(p,c)^2$

*x proportional to num of data*

# Justification of the Representation

rowsum =) smallest rowsum (clustroid)

- Clustroid is the point in the cluster with the smallest rowsum

- Why the $k$ points nearest to the clustroids?

  - If the clustroid changes, the new clustroid would be one of them
  - $p$ becomes the new clustroid if $\text{rowsum}(p) < \text{rowsum}(c)$

- Why the $k$ points farthest to the clustroids?

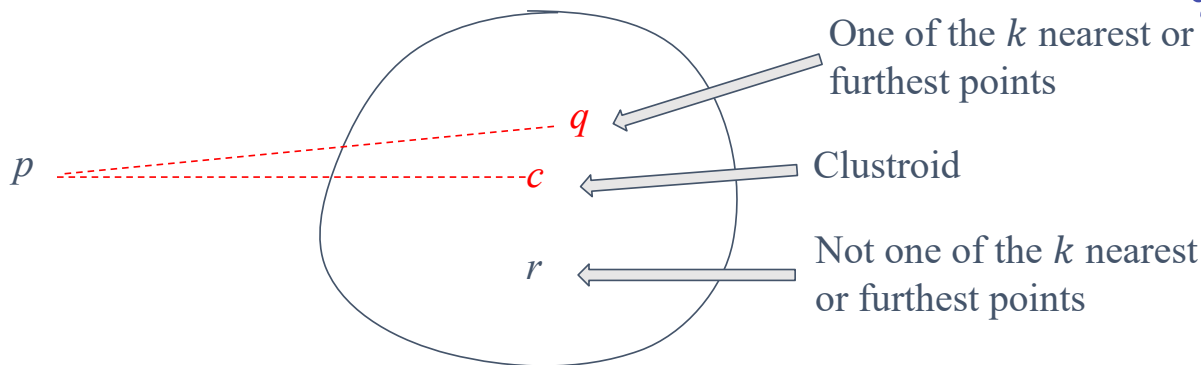  - Used to determine whether two clusters are close enough to merge

# Adding a Point to a Cluster

- How can we add a point $p$ to a cluster?
  - Add 1 to $N$
  - For each $q \in \{\text{clustroid}\} \cup k$ nearest points $\cup k$ furthest points
  - Update rowsum$(q)$ as $\underline{\text{rowsum}(q) + d^2(p, q)}$ //
- What if $p$ needs to be included in the representation?
  - We cannot compute this exactly without going to disk → 만들 전들의 점들가 없는데 $p$의 rowsum을 알 수가 없다.



One of the $k$ nearest or furthest points

Clustroid

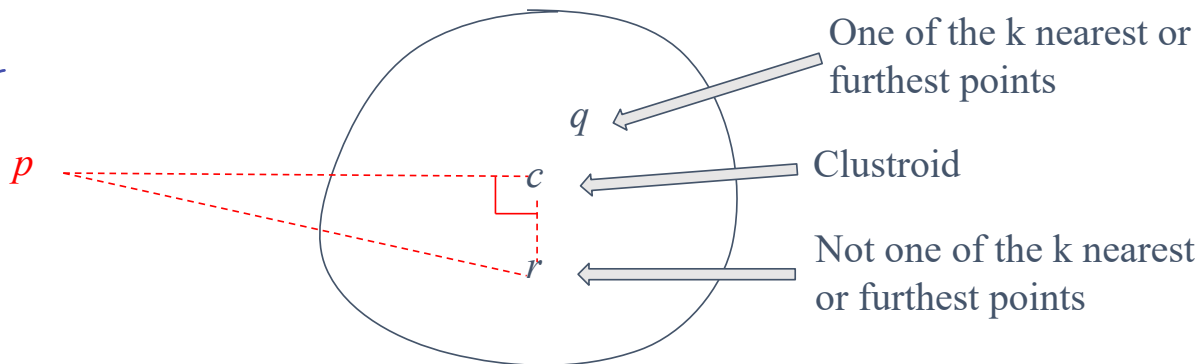Not one of the $k$ nearest or furthest points

# Estimating the Rowsum

- **Estimation:** $\text{rowsum}(c) + N \times d^2(p, c)$   는 점들까지의 거리가 없으므로 가능

- With the curse of dimensionality, almost all angles are right angles

- Thus, $d^2(p, r) \approx d^2(p, c) + d^2(c, r)$ by the Pythagorean theorem

증명. ⎡ 거리?
    ⎣ angle?



$p$

$q$ → One of the k nearest or furthest points

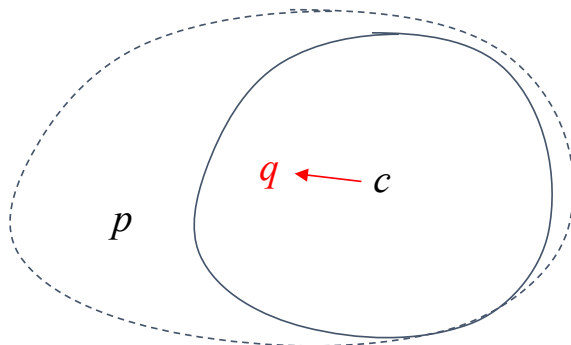$c$ → Clustroid

$r$ → Not one of the k nearest or furthest points

# Possibly Updating the Clustroid

- If $\mathrm{rowsum}(p) < \mathrm{rowsum}(c)$, make $p$ the new clustroid
- Eventually, the true clustroid may not be one of the $k$ closest points
  - Cluster representation needs to be recomputed periodically from disk

# Other Details of GRGPF

- See the textbook for other details:
  - How to initialize the cluster tree
  - How to use the tree for each new point
  - How to split a cluster
  - How to merge clusters (with the $k$ furthest points)

# Summary

1. **BFR Algorithm**
   - Cluster representation
   - Three classes of sets

2. **BFR Algorithm: Process**

3. **CURE Algorithm**
   - 2-pass algorithm

4. **GRGPF Algorithm**
   - Rowsum
   - Estimation of a rowsum