











- A clustering is a set of clusters
- Important distinction between hierarchical and partitional sets of clusters
 - Partitional (unnested) Clustering
 - A division of data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset
 - Hierarchical (nested) clustering
 - A set of nested clusters organized as a hierarchical tree
 - Each node (cluster) in the tree (except for the leaf
 - nodes) is the union of its children (subclusters), and the root of the tree is the cluster containing all the objects



8



9









<section-header><list-item><list-item><list-item><table-container>

 Types of Clusters

 • Graph-Based (Contiguity-Based) Clusters

 • A cluster is a set of points such that a point in a cluster is closer (or more similar) to one or more other points in the cluster than to any point not in the cluster.

 • Other Contiguity-Date

 • Other Contiguity-Date

14



15









Clustering Algorithms

- Techniques to introduce many of the concepts involved in cluster analysis:
 - K-means and its variants
 - a prototype-based, partitional clustering technique that attempts to find a user-specified number of clusters (K), which are represented by their centroids
 - Agglomerative Hierarchical clustering
 - a collection of closely related clustering techniques that produce a hierarchical clustering by starting with each point as a singleton cluster and then repeatedly merging the two closest clusters until a single, all encompassing cluster remains
 - Density-based clustering

19

K-means Clustering

- Prototype-based clustering techniques create a onelevel partitioning of the data objects.
 - K-means
 - defines a prototype in terms of a centroid, which is usually the mean of a group of points, and is typically applied to objects in a continuous n-dimensional space.
 - K-medoid
 - defines a prototype in terms of a medoid, which is the most representative point for a group of points, and can be applied to a wide range of data since it requires only a proximity measure for a pair of objects.
 - While a centroid almost never corresponds to an actual data point, a medoid, by its definition, must be an actual data point

20

K-means Clustering

- · Partitional clustering approach
- Number of clusters, K, must be specified
- Each cluster is associated with a centroid (center point)
- Each point is assigned to the cluster with the closest centroid
- The basic algorithm is very simple:
- 1: Select K points as the initial centroids
- 2: repeat
- 3: Form K clusters by assigning all points to the closest centroid.
- 4: Recompute the centroid of each cluster.
- 5: **until** The centroids don't change

21



22





K-means Objective Function

- A common objective function (used with Euclidean distance measure) is <u>Sum of Squared</u> Error (<u>SSE</u>), which is also known as scatter
 - For each point, the error is the distance to the nearest cluster center

To get SSE, we square these errors and sum them. $SSE = \sum_{k} \sum_{i=1}^{k} dist(c_{i}, x)^{2} \qquad c_{i} = \frac{1}{k} \sum_{i=1}^{k} x$

$$\sum_{i=1}^{n} \sum_{x \in C_i} dist(c_i, x)^2 \qquad c_i = \frac{1}{m_i} \sum_{x \in C_i}$$

- *x* is a data point in cluster C_i , c_i is the centroid for cluster C_i , m_i is the number of objects in the *i*th cluster
- SSE improves in each iteration of K-means until it reaches a local or global minima.

25

K-means Objective Function

- To illustrate, the centroid of a cluster containing the three two-dimensional points, (1,1), (2,3), and (6,2),
- Centroid is

1

$$c_i = \frac{1}{m_i} \sum_{x \in C_i} x$$

$$((1+2+6)/3, (1+3+2)/3) = (3, 2)$$

26

K-means Objective Function K-means is not restricted to data in Euclidean space Consider document data and the cosine similarity measure. Assume that the document data is represented as a document-term matrix Objective is to maximize the similarity of the documents in a cluster to the cluster centroid; this quantity is known as the cohesion of the cluster. The analogous quantity to the total SSE is the total cohesion, which is given by SSE = ∑ cosine(c_i, x) / (x is a data point in cluster C_i, c_i is the centroid for cluster C_i

27





28











33







35

have only one.

Copyright 2000 N. AYDIN. All rights reserved.





Limitations of K-means · K-means has problems when clusters are of differing Sizes - Densities - Non-globular shapes · K-means has problems when the data contains outliers. - One possible solution is to remove outliers before clustering

41



38



Limitations of K-means: Differing Sizes K-means (3 Clusters) **Original Points**



Limitations of K-means: Non-globular Shapes

44



45









Hierarchical Clustering



Hierarchical Clustering

- Two main types of hierarchical clustering - Agglomerative:
 - Start with the points as individual clusters
 - At each step, merge the closest pair of clusters until only one cluster (or k clusters) left
 - Divisive:
 - Start with one, all-inclusive cluster
 - At each step, split a cluster until each cluster contains an individual point (or there are k clusters)
- Traditional hierarchical algorithms use a similarity or distance matrix

- Merge or split one cluster at a time

50



51

49





Steps 1 and 2





54







































68



69



70

Hierarchical Clustering: Group Average

- Compromise between Single and Complete Link ٠
- Strengths
 - Less susceptible to noise
- Limitations Biased towards globular clusters

Cluster Similarity: Ward's Method

- Similarity of two clusters is based on the increase in squared error when two clusters are merged
 - Similar to group average if distance between points is distance squared
- · Less susceptible to noise
- Biased towards globular clusters
- Hierarchical analogue of K-means - Can be used to initialize K-means



Hierarchical Clustering: Time and Space requirements

- O(N²) space since it uses the proximity matrix. - N is the number of points.
- O(N³) time in many cases
 - There are N steps and at each step the size, N², proximity matrix must be updated and searched
 - Complexity can be reduced to $O(N^2 \mbox{log}(N)$) time with some cleverness

74

Hierarchical Clustering: Problems and Limitations

- Once a decision is made to combine two clusters, it cannot be undone
- No global objective function is directly minimized
- Different schemes have problems with one or more of the following:
 - Sensitivity to noise
 - Difficulty handling clusters of different sizes and non-globular shapes

DBSCAN

DBSCAN is a density-based algorithm.

Density = number of points within a specified radius (Eps)

A point is a core point if it has at least a specified number of points

· These are points that are at the interior of a cluster

A border point is not a core point, but is in the neighborhood of a

A noise point is any point that is not a core point or a border point

- Breaking large clusters

(MinPts) within Eps

core point

· Counts the point itself

75

77



76



78















or clusterings

87



86

Measures of Cluster Validity Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following two types. Supervised: Used to measure the extent to which cluster labels match externally supplied class labels. Entropy Often called external indices because they use information external to the data Unsupervised: Used to measure the goodness of a clustering structure without respect to external information. Sum of Squared Error (SSE) Often called internal indices because they only use information in the data







































105

Final Comment on Cluster Validity

- "The validation of clustering structures is the most difficult and frustrating part of cluster analysis.
- Without a strong effort in this direction, cluster analysis will remain a black art accessible only to those true believers who have experience and great courage."

Algorithms for Clustering Data, Jain and Dubes

 H. Xiong and Z. Li. Clustering Validation Measures. In C. C. Aggarwal and C. K. Reddy, editors, Data Clustering: Algorithms and Applications, pages 571–605. Chapman & Hall/CRC, 2013.