Machine Learning CS 4641-7641



Density-Based Clustering

Mahdi Roozbahani Georgia Tech

These slides are inspired based on slides from Jing Gao, Chao Zhang and Jiawei Han.

when you're a constant and you see d/dx



Outline

- Overview
- Basic Concepts
- The DBSCAN Algorithm
- Analysis of DBSCAN

Density-Based Clustering

Basic Idea

- Clusters are dense regions in the data space, separated by regions of lower density
- A cluster is defined as a maximal set of density-connected points
- 。 Detect arbitrarily shaped clusters
- Method
 - DBSCAN (<u>Density-Based Spatial</u> <u>Clustering of Applications with Noise</u>)



Outline

- Overview
- Basic Concepts
- The DBSCAN Algorithm
- Analysis of DBSCAN

High Density v.s. Low Density

- Two parameters
 - Eps (ε): Maximum radius of the neighborhood



- MinPts: Minimum number of points in the Eps-neighborhood of a point
- High density: ε-Neighborhood of an object contains at least MinPts of objects



Density of *p* is low Density of *q* is high

Core Points, Border Points, and Outliers



 $\varepsilon = 1$ unit, MinPts = 5

Given *ɛ* and *MinPts*, categorize the objects into three exclusive groups.

A point is a core point if it has more than a specified number of points (MinPts) within Eps—These are points that are at the interior of a cluster.

A border point has fewer than MinPts within Eps, but is in the neighborhood of a core point.

A noise point is any point that is not a core point nor a border point.

Practice:



Examples





Original Points

Point types: core, border and outliers

 ε = 10, MinPts = 4

Density-based related points

- Direct density reachability:
 - An object p is directly density-reachable from object q if (1) q is a core object; and (2) p is in q's ε-neighborhood



Density-based related points

- Density reachability:
 - A point **p** is density-reachable from a point **q** if there is a chain of points $p_1, \ldots, p_n, p_1 = q, p_n = p$ such that p_{i+1} is directly density-reachable from p_i

$$p_1 = q \rightarrow p_2 \rightarrow \dots \rightarrow p_n = q$$



Density-based related points

- Density connectivity:
 - A point *p* is density-connected to a point *q* if there is a point *o* such that both *p* and *q* are density-reachable from *o*







Outline

- Overview
- Basic Concepts
- The DBSCAN Algorithm
- Analysis of DBSCAN

The DBSCAN Algorithm



https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/

Outline

- Overview
- Basic Concepts
- The DBSCAN Algorithm
- Analysis of DBSCAN

DBSCAN is Sensitive to Parameters

Figure 8. DBScan results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.

Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.





Image Credit: George Karypis.



High value (what will happen?)

Clusters will merge and the majority of data points will be in the same cluster Low value (what will happen?)

A large part of data won't be clustered and considered as outliers. Because, they won't satisfy the number of pints to create a dense region

Do we need to define the number of clusters in DBSCAN?

Nope

Minimum number of Points (MinPts)



How about Eps? (Elbow effect)

- Idea is that for points in a cluster, their kth nearest neighbors are at roughly the same distance
- Noise points have the kth nearest neighbor at farther distance
- So, plot sorted distance of every point to its kth nearest neighbor



Here we have 3000 points and x-axis shows just a point index. Point indices are sorted in ascending order based on their 4th nearest neighbor distance

Elbow effect another example



minPts often does not have a significant impact on the clustering results

Erich Schuber et al

When DBSCAN Works Well

- Robust to noise
- Can detect arbitrarily-shaped clusters





Original Points

Clusters

When DBSCAN Does NOT Work Well

- Cannot handle varying densities
- Sensitive to parameters—hard to determine the best setting of parameters



Original Points





(MinPts=4, Eps=9.75)

Take-Home Messages

- The basic idea of density-based clustering
- The two important parameters and the definitions of neighborhood and density in DBSCAN
- Core, border and outlier points
- DBSCAN algorithm
- DBSCAN's pros and cons

Clustering Evaluation

- Internal measures for clustering evaluation
 - 。 Elbow method
 - 。Silhouette Coefficient
 - 。Graph-based measures (Beta-CV and Normalized cut)
 - Davies-Bouldin Index

We want intra-cluster datapoints to be as close as possible to each other and inter-clusters to be as far as possible from each other

The Davies-Bouldin Index



The Davies–Bouldin measure for a pair of clusters C_i and C_j is defined as the ratio ratio $D_{1,} D_{2,} D_{3} \iff D_{1} = \max \{ DB_{12}, DB_{13} \} \Rightarrow DB_{12} = \bigcup_{i \neq j} DB_{12} = \bigcup_{i \neq j} DB_{12} = \bigcup_{i \neq j} DB_{ij} = \bigcup_{i \neq j} DB_{ij} = \bigcup_{i \neq j} DB_{ij} = \bigcup_{i \neq j} DB_{ij}$

 DB_{ij} measures how compact the clusters are compared to the distance between the cluster means. The Davies–Bouldin index is then defined as

$$DB = \frac{1}{k} \sum_{i=1}^{k} D_i$$

a lower value means that the clustering is better