

⇒ Cluster Analysis

- partitions large dataset into groups acc. to similarity.
- used For outlier detection.
- main focus → distance based cluster-analysis
- unsupervised learning, learning by observation.
- Requirements of clustering:
 - ~ Scalability (should work on a large scale database)
 - ~ Ability to deal diff. data types: numeric, binary, nominal, categorical, ordinal or mixture of them.
 - ~ discovery of arbitrarily shaped clusters: should be able to detect arbitrarily shaped clusters & not only spherical ones.
 - ~ Domain knowledge & input params.
 - ~ Dealing with noisy data.
 - ~ Incremental clustering & independence from input order.
 - ~ clustering high-dimensional data
 - ~ constraint-based clustering.
 - ~ Interpretability & usability.
- Clustering approach
 - ~ Partitioning criteria: with or without hierarchy
 - ~ separⁿ of clusters: clusters may or may not be mutually exclusive.
 - ~ similarity measures: distance function, or connectivity/density/contiguity based.
 - ~ dist. based ⇒ adv. of optimization
 - ~ connectivity ⇒ " of arbitrary shaped clusters^{formal}
 - ~ clustering space: which parameters/attributes to consider while clustering.

⇒

K-means:-

- Form random k -partitions.
- repeat:
 - ~ compute centroid of each cluster
 - ~ Assign obj to nearest centroid / clusteruntil no change.
- $O(nkt)$ terminates at local optimal.
- Weakness:
 - sensitive to noise / outliers
 - Need to specify k in advance.
- ⇒ Variations
 - ~ selection of initial k means
 - ~ dissimilarity calculations
 - ~ strategies to calc. cluster means.

⇒

K-modes

- Variation of K-means, which can cluster nominal data.
- replace mean by mode.
- a different dissimilarity funcⁿ
- freq. based method to update mode.

⇒ K-medoids

→ basic idea: In K-means, rather than mean value, select an actual object most similar to mean value of each cluster.

→ NP-hard.

→ Partitioning Around Medoids (PAM)

→ randomly select initial K - objects o_1, \dots, o_K

repeat:

replace o_i by o_j if quality increased

~~random~~ \rightarrow all possible points.

until no change in quality

⇒ effective for small dataset, not scalable

⇒ Density-Based Clustering

→ to find arbitrarily shaped clusters.

→ handles noise

→ one scan

→ needs density parameters for termination

→ DBSCAN:

→ Two parameters: ϵ : radius of ϵ neighborhood.
 $MinPts$: density threshold.

→ core-objects: An object with contains at least $MinPts$ objects in its ϵ neighborhood

→ clustering: ^{use} core-objects & its ϵ neigh., to form dense region i.e clusters.

→ some terminologies:

↳ Directly-density reachable (ddr):

core-object q , object p ; p is ddr from q if p belongs to ϵ -neighborhood of q .

↳ Density-reachable (dr):

~~For two points p & q are dr~~

A point q is ~~dr~~ dr from p if

for n points p_1, p_2, \dots, p_n , (where

$p_1 = p$ & $p_n = q$) p_{i+1} is ddr from p_i .

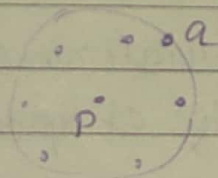
↳ Density connected:

Two points p & q are dc. if \exists a point

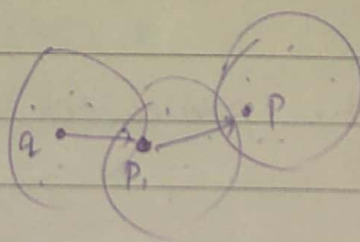
o such that p is dr from o &

q is also dr from o .

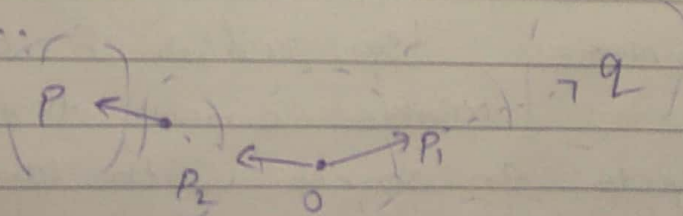
\Rightarrow ddr:



\Rightarrow dr:



\Rightarrow dc:



⇒ DBSCAN Algorithm:

* (w.r.t ϵ & minPts)

✓ points P in database:

↪ retrieve all density-reachable points of p *

if p is core-point:
it is a cluster

else

do nothing.

→ spatial indexing → $O(n \log n)$ else $O(n^2)$

⇒ OPTICS:

↪ users have to select ^{/provide} parameters for DBSCAN

↪ OPTICS: outputs a ^{ordering} of clusters.

↪ To construct clusters, object are processed in a ^{select} specific order.

↪ Order: Object that is ^{core-}dist with lowest ϵ
(i.e. high density clusters are first)

↪ core-distance: For an object p , ^{core-}distance ϵ' is the minimum value for which \hat{p} is a core-object

↪ reachability-dist: minimum radius value that makes p desc. from q

$$= \max \{ \text{core-distance}(q), \text{dist}(p, q) \}.$$

⇒ DENCLUE:-

→ uses kernel density estimation (non-parametric)

→ general idea:

- ~ treat an observed object as an indicator of high-probability density in its surrounding
- ~ probability density \propto dist. From this point to observed objects

For a random variable f ,
kernel density approxⁿ of P.D.F is

$$f_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x_i - x}{h}\right)$$

h → smoothing parameter

K → kernel funcⁿ: $\int_{-\infty}^{\infty} K(u) \cdot du = 1$

$$K(u) = K(u)^2$$

freq used : $K\left(\frac{x_i - x}{h}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - x}{h}\right)^2}$

kernel funcⁿ :

(Gaussian with $\mu=0$ & $\sigma^2=1$)

~ a point x^* is density attractor

if → its a local maximum of estimated density func.

↳ $f(x^*) \geq$ some threshold.

these non-trivial $x^*(s)$ are centers of clusters
Some funcⁿ(s):

Influence of y on $x = f_{\text{gaussian}}(x, y) = \frac{-(d(x, y))^2}{e^{2\sigma^2}}$

Total " on $x = f_{\text{gaussian}}^D(x) = \sum_{i=1}^N f(x, x_i)$

gradient of x in direction of x_i = $V f_{\text{gaussian}}^D(x, x_i) = \sum_{i=1}^N (x - x_i) \frac{-(d(x, x_i))^2}{e^{2\sigma^2}}$

$$= \sum_{i=1}^N (x - x_i) \cdot f_{\text{gaussian}}(x, x_i)$$

⇒ STING:

- divide the spatial area into ^{rectangular} cells
 - Arrange them in an hierarchy
 - ↳ highest level cell \Rightarrow partitioned to smaller cells
 - ↳ Use top-down approach to answer queries
 - remove irrelevant cells from further consideration.
- Advantages: Query-independent, incremental update, $O(k)$, $k \rightarrow$ number of grids at lowest level.

⇒ CLIQUE:-

- both density based & grid based.
- ↳ partitions m -dimensional data into non-overlapping rectangular units & each dimensional is partitioned into eq. number of intervals (of equal length)
- ↳ frac of data points \geq input parameter contained in unit \Rightarrow for a unit to be dense
- ↳ cluster \approx maximal set of connect dense units in a subspace.
- ↳ Adv.:- automatically finds subspace of highest dimensionality
 - :- insensitive to order of records in input
 - :- linearly scalable.
- Disadv.:- trade off between accuracy & simplicity of method.

⇒ Evaluation of clustering:-

⇒ Assesing Clustering Tendency:-

we check for existence of non-random because a clustering algo. may return random & non-meaningful clusters.

we statistical tests for spatial randomness:-

⇒ Hopkins Statistic:-

dataset D

step 1: sample n random points p_1, \dots, p_n

$\forall i \in 1 \text{ to } n$

$$x_i = \min_{V \in D} d(p_i, V)$$

step 2: n considers n points q_1, \dots, q_n

$\forall i \in 1 \text{ to } n$

$$y_i = \min_{V \in D - \{q_1, q_2, \dots, q_n\}} d(q_i, V)$$

step 3:

$$H = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i + \sum_{i=1}^n y_i}$$

$H > 0.5$, (reject clustering)

⇒ Determining # clusters:

↪ determining how many clusters to form, is a really difficult task.

⇒ most simple method ⇒ # clusters = $\sqrt{n/2}$
cluster size = $\sqrt{2n}$

→ Elbow Method

clusters ↑ ⇒ sum of intra-cluster variance ↓

$\text{var}(k)$ = total variance when k clusters are formed.

plot curve of var against k . The First/most significant turning point of curve is suggest optimal value of k .