# * clustering

cluster → a no. of similar things that occur together.

→ technique in which data points are arranged in similar groups dynamically without any pre-assignⁿ of groups.

Appⁿ →
- Image Processing
- Recommedation Engine
- Insurance

# * properties of clustering algo.

- scalability
- ability to deal with diff data types
- Minimal to requirements for domain knowledge for determine input parameters
- Interpretability & usability

# * Typical Requirements → clustering in data mining

All these ⎨ + {
- Scalability.
- Discovery of clusters with arbitrary slope.
}

# * Problems with clustering
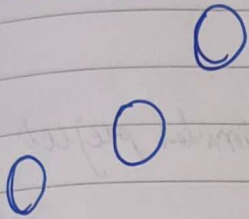- not address all requirⁿ
- dealing with large no. of dimensions & data items → problematic
- effectiveness depends on distance
- distance measure should it define it.
- result of clustering → interpreted n diff ways

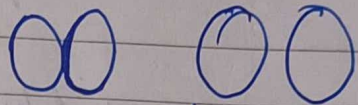# Types of Clustering

## a) well seperated clusters

a cluster of
- set of points such than
  any point in a cluster in closes
  to every other point is cluster



- threashold → used to
  specify all objects in
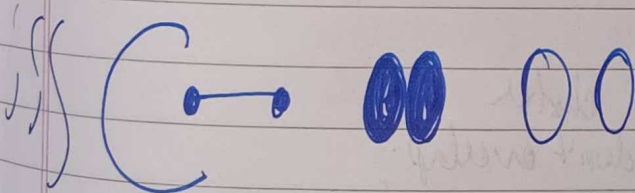  clust or sufficiently close

## b) prototype based clusters.

- set of object → is closer
  to prototype or cental
  of cluster.



- data → numerical → centroid.
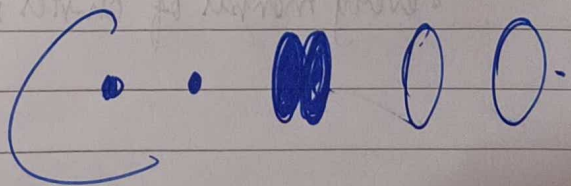- data → categorical → prototype
- kmeans & kmedoids are example
- .

## c) contiguity based clusters.

point in cluster is close to one or
more point in cluster than any
other



## d) Density based.

- cluster → dense region of parts
  seperated by low-density region
  from other regions of hight density

- used when cluster are irregular
  or interhoind & when herself
  outlier are preset

# ✱ K-Means -

- heuristic method.
- supervised learning algo.
- solve clustering problems.
- group are unlabelled dataset in diff clusters.
- $K \rightarrow$ defines no. of predefined clusters needs to be created.

  $K=2 \rightarrow 2$ clusters

- each dataset belongs only one group has similar project
- minimizes sum b/w data points & clusters
- distance calculation $\rightarrow$ Euclidean distance

Adv $\rightarrow$ efficient in compution.
$\rightarrow$ easily to implement.

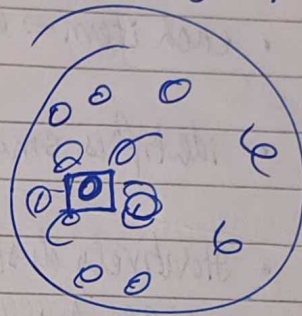Disadv $\rightarrow$ only when mean is defined
$\rightarrow$ need to specify $K$, no. of clusters, in advance
$\rightarrow$ Trouble with noisy data & outliers
$\rightarrow$ not suitable to discove clustel with non complex
shapes

## Properties —

- always $K$ clusters
- always atleast one item in cluster
cluster
- are non - hierarchial & don't overlap.
- every member of cluster is closes to cluster.

**\# k - medoids -**

- each cluster represented by one of objectes in cluster.

- Data points are choosen by medoids.

- classic partition clustering techniqus that groups data set of n objects into k groups

- known as priori

- less delicate to outlier

- convex shape not required.

- More robust to noises.

**\# Hierarcheal Clustery**

- method of cluster analysis in which data points are arranged in hierachy of clusters.

Adv → • simple to implement
- easy & results in hierachy (adl info)
- doesn't need prespecify no. of clusts.

Disadv → • largeclusters breaks
- diff to handle diff sizedclusters & shapes
- sensetive to noise & outliees
- can't be chaged or deleted once done

# Hierachical Clustering
## Types

| Agglomerative Clustering | Divisive Clustering |
|---|---|
| • bottom up approach | • town down approach |
| • each item → own cluster | • all item → one cluster |
| • identifies small cluster | • identifies large cluster |
| • Iteratively clustered all merged together | • Large clusters are successively divided |
| • also know AGNES (Agglomerative Nesting) | • also known DIANA (divise analysis) |

**✱ Dendogram**

- Diagram representing tree of hierachy.
- The similar two objects are less is the height of link that joins them

- major info lost here

# * DBSCAN

- Density Based spatial clustering of Application with Noise.
- groups together closely packed points.
- unsupervised learning.

$\epsilon \rightarrow$ radius of neighborhood around point x.
minpts $\rightarrow$ required to form a density cluster

Adv $\rightarrow$
- don't need to specify no. of clusters
- flexible in shape & size of clusters
- Able to deal with noise & outliers
- Ability to identify uneven shape.
- easy $\rightarrow$ someone who knows dataset $\rightarrow$ to set parameters.

Disadv $\rightarrow$
- Input parameters $\rightarrow$ difficult to determine
- some situation very sensitive to input parameter.
- confused $\rightarrow$ border point belong to too cluster
- result $\rightarrow$ distance metric
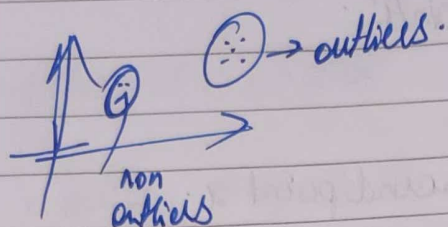- hard ro guess correct parameters

# * Spectral clustering $\rightarrow$
- solves by creating clusters with arbitary & non-linear shapes, no assumption $\rightarrow$ shape of clusters.

Adv $\rightarrow$
- no assumption $\rightarrow$ satistices of cluster
- Easy to implement
- Good clustering result
- fast to sparse data set of several thousand elements

Disadv $\rightarrow$
- maybe sensative to choice of parameters
- computationally expensive for large subset

# ✳ Outlier Analysis —

- satistical observation i.e. marked differently in value from others in sample.



→ outliers.

non outliers

## → Types —

### 1) Global
↓
- a data obj is called global outlier.
- if it deviates from rest of dataset

### 2) Contextual (Conditional)

- deviates significantly on context of object
- only with context → outlier
- generalization of local outliers

### 3) Collective Outliers —
- Obj as a whole deviates significantly from entire dataset. it is collective.
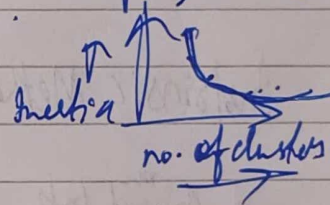
- Challenges in Outlier detection

→ modeling normal objects & outliers
→ appln specific outlier detection
→ handling noise in " "
→ Understandability

## Local Outlier Factor (LOF)

* concept of local density
* locality → k nearest neighbors
* distance b/w k's used to determine density
* points that lower local density ⇒ their neighbors are outliers
* unsupervised anomaly detection - method

## Elbow method -

* used to determine no. of clusters
* Inertia is sum of all distance of data point from centroid of clusters.

algo → * start with any value of k & perform k-mean algo.
* Determine total inertia
* Increase k by 1 & carry out step 1 & 2 until inertia is not significant.

Inertia ⌐⌐⌐ no. of clusters

- Measuring Clustering Quality

  - Extrinsic ⟶ Ground truth available
                ⟶ reward behaviour

    - cluster homogenity ⟹ purer cluster ⟶ better clustering
    - cluster completeness ⟶ If 2 two obj ⟶ same category
                                              ↓
                                          same cluster
    - Rag Bag ⟹ obj cant be merged with other obj.
    - Small cluster preservation ⟹ splitting small category into
                                    Small is more
                                    harmful than
                                    large ⟶ small.

# Intrinsic Method –

- Ground truth not available.
- evaluate how well clusters separated ⟹ completeness
- 
- Silhouette coefficient ⟹ defines goodness of clustering
                                              technique.

  values ⟶ $-1$ to $1$

  $1$ ⟶ clusters well part & distinguished
  $0$ ⟶ cluster indifferent
  $-1$ ⟶ clusters are assigned in wrong way

  Silhouette Score ⟶ $\dfrac{b-a}{max(a,b)}$