# Clustering methods

Ali Madani
Farnoosh Khodakarami

# DataSets

Features(Attribute/variable)

| ID | Address | # Bed | #Bath | ... | School Score | Year Build | Crime Rate |
|----|---------|-------|-------|-----|--------------|------------|------------|
|    |         |       |       | ... |              |            |            |
|    |         |       |       | ... |              |            |            |

Data records (samples)

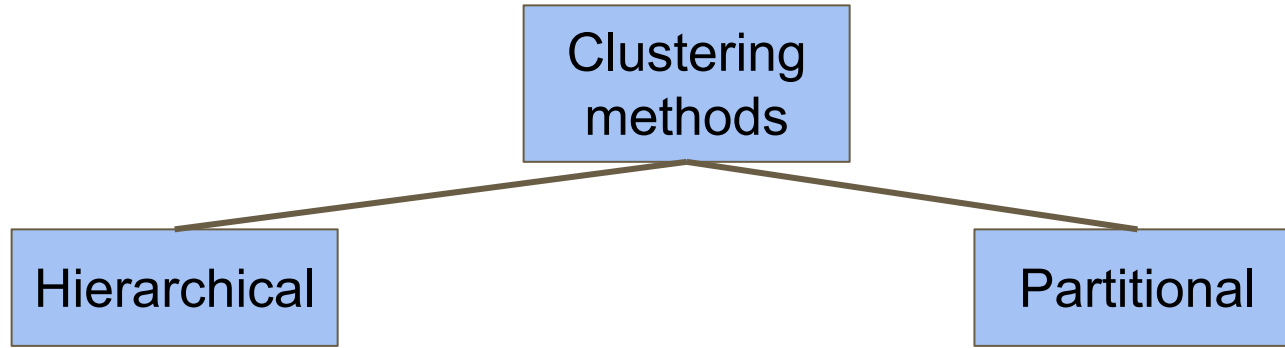**# Features = Dimension of dataset**

# Unsupervised versus supervised learning

- **Supervised learning**
  - Trying to predict label or value of data points
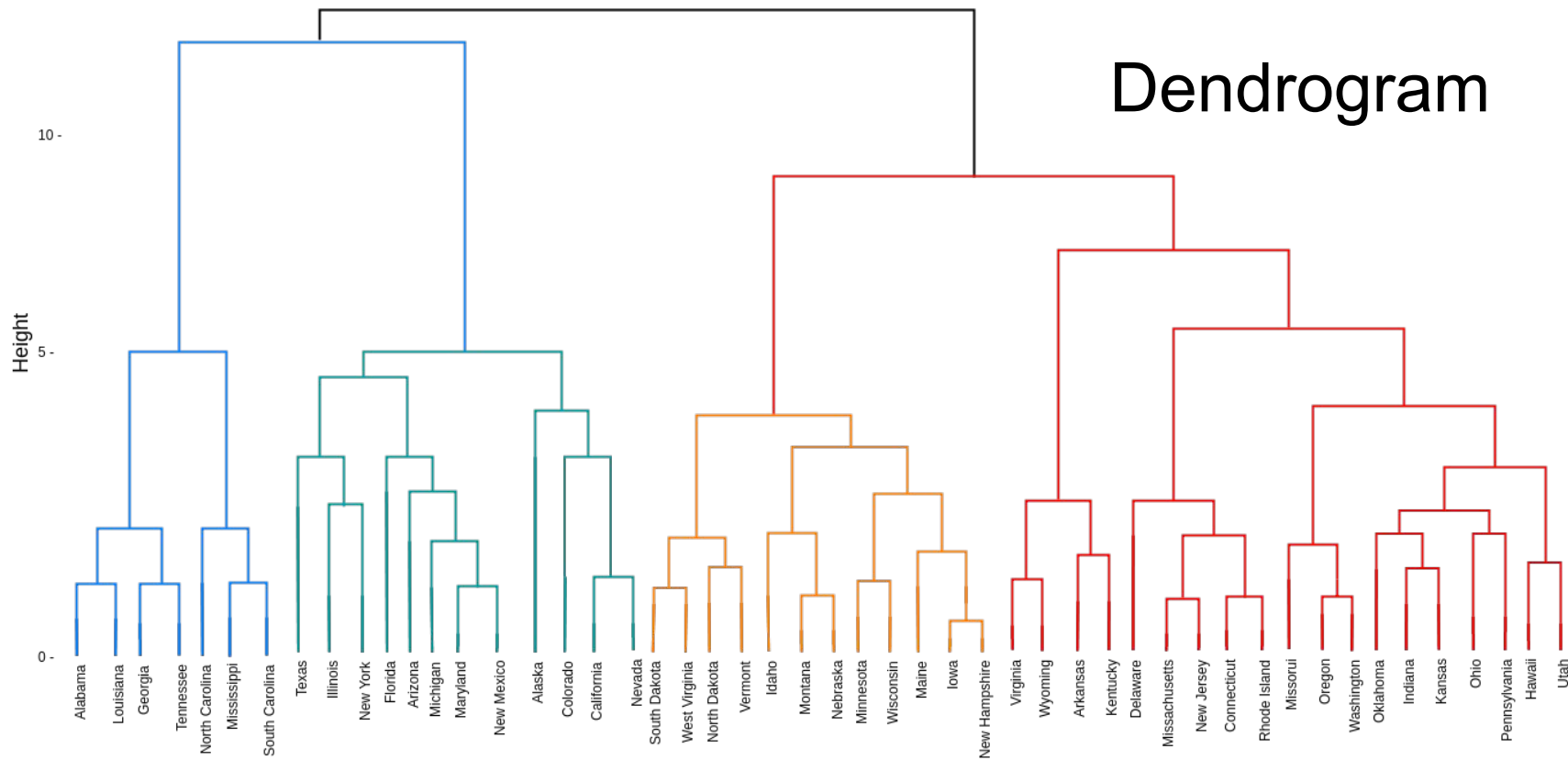
- **Unsupervised learning**
  - Unlabeled input data

# Why do we use clustering?

- Group data points based on their similarities using the given features
  - Identify similar data points (within the same group)
  - Identify dissimilar data points (within different groups)

- Some methods also identifies outliers and points that cannot be assigned to any cluster

- Investigating differences between the groups
  - Survival of breast cancer patients

# Categories of clustering methods that we focus on

# Agglomerative hierarchical clustering



Dendrogram

https://online.visual-paradigm.com/diagrams/templates/dendrogram/cluster-dendrogram/

# Steps of agglomerative hierarchical clustering

1.  Computing dissimilarity or similarity between every pair of data points
2.  Using linkage function to group objects into hierarchical cluster tree
    a.  linkage function determines the distance between sets of data points as a function of the pairwise distances between data points in the groups.
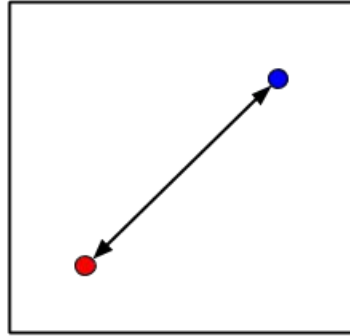3.  Deciding where to cut the hierarchical tree into clusters.

# Common distance measures used in clustering

- Euclidean distance
- Manhattan distance
- Minkowski distance
- Chebychev distance
- Cosine similarity
- Hamming distance
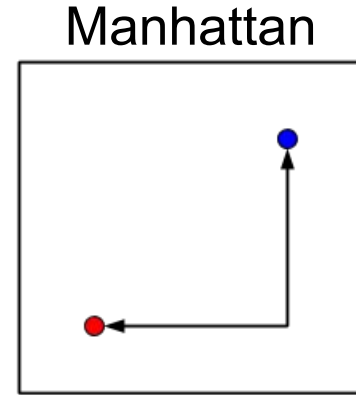- Binary distance
  - Jaccard index
  - Hamming distance

# Euclidean distance

Euclidean

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$
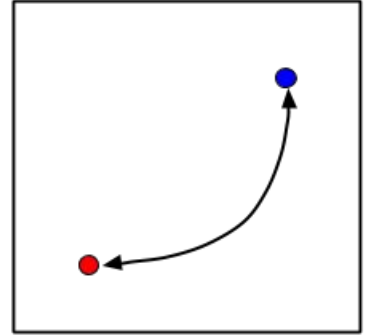
# Manhattan distance

Manhattan



$$d = |x_1 - x_2| + |y_1 - y_2|$$

# Minkowski distance

Minkowski

$$d = \left( |x_1 - x_2|^p + |y_1 - y_2|^p \right)^{1/p}$$
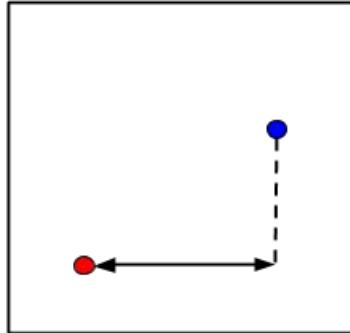
# Chebychev distance

$$d = max(|x_1 - x_2|, |y_1 - y_2|)$$

Chebychev

# Jaccard index

$$J(A, B) = \frac{A \cap B}{A \cup B}$$

Jaccard

"G o l a n g"

"G o p h e r"

# Hamming distance
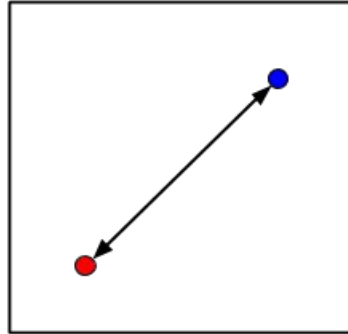
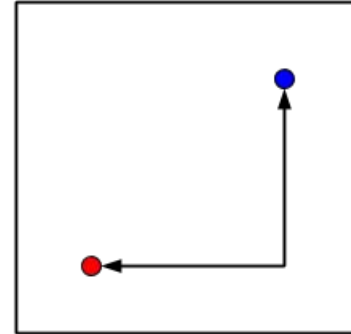d=*number of bits that they are different*
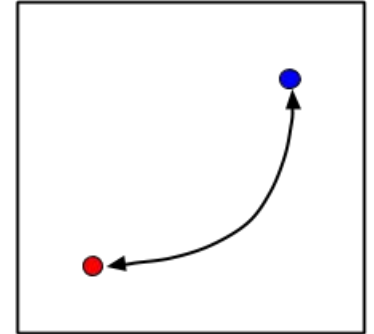
Hamming

"Golang"

"Gopher"

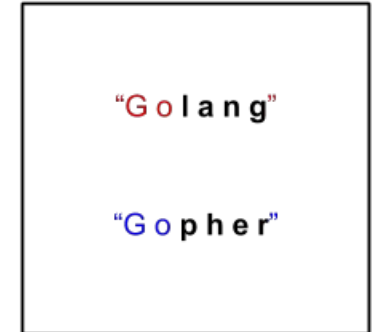# Different distance measures

Euclidean     Manhattan     Minkowski
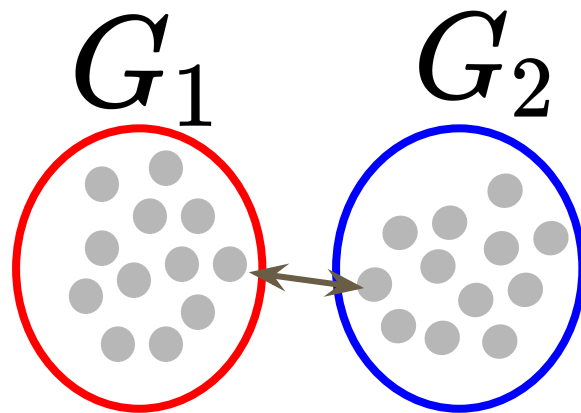
Chebychev     Jaccard     Hamming

"Golang"

"Gopher"

"Golang"

"Gopher"

# Some of widely-used linkage functions

- Single
- Complete
- Average
- Median
- Centroid

# Some of widely-used linkage functions

- **Single**
- Complete
- Average
- Median
- Centroid

$$G_1 \qquad G_2$$

$$D(G_1, G_2) = min(d(x,y)), x \in G_1, y \in G_2$$

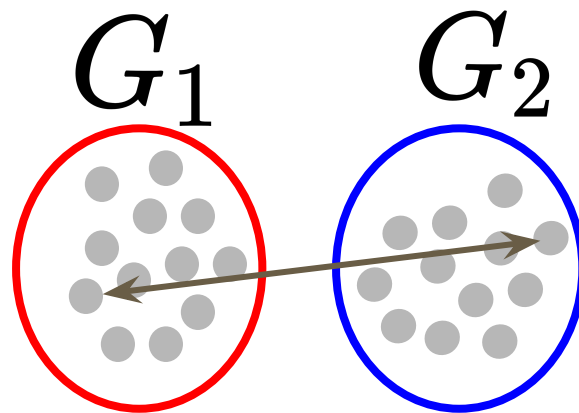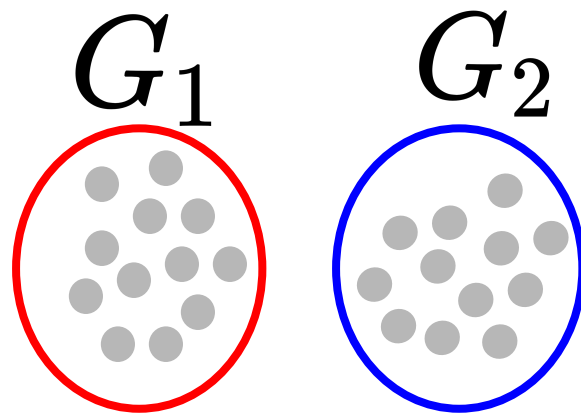# Some of widely-used linkage functions

- Single
- **Complete**
- Average
- Median
- Centroid

$G_1$    $G_2$

$$D(G_1, G_2) = max(d(x, y)), x \in G_1, y \in G_2$$

# Some of widely-used linkage functions

- Single
- Complete
- **Average**
- Median
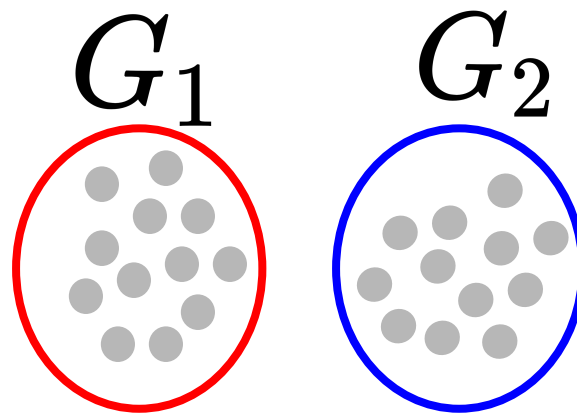- Centroid

$$G_1 \qquad G_2$$

$$D(G_1, G_2) = \frac{1}{|N_{G_1}||N_{G_2}|} \Sigma_x \Sigma_y d(x, y)$$

$$x \in G_1, y \in G_2$$

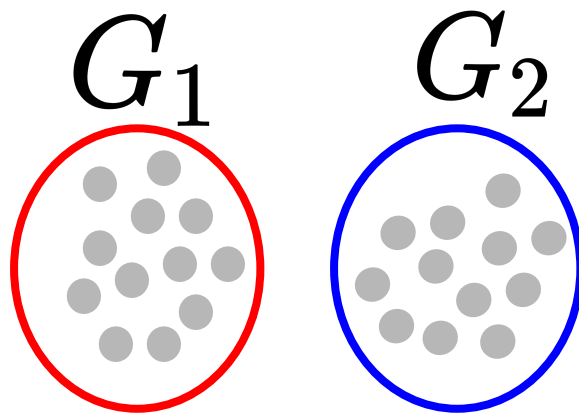# Some of widely-used linkage functions

- Single
- Complete
- Average
- **Median**
- Centroid

$$G_1 \qquad G_2$$

Weighted center of mass distance (WPGMC)
*Note*. appropriate for Euclidean distances only

# Some of widely-used linkage functions

- Single
- Complete
- Average
- Median
- **Centroid**

$$G_1 \qquad G_2$$

$$D(G_1, G_2) = \|c_{G_1} - c_{G_2}\|$$

*Note*. appropriate for Euclidean distances only

# Let's see how agglomerative hierarchical clustering works

Dendrogram



Height

Data points should be rearranged for building the dendrogram

# Let's see how agglomerative hierarchical clustering works

Dendrogram

Height

a   e   c   b   d

Data points should be rearranged for building the dendrogram

a
e
c
b
d

# Let's see how agglomerative hierarchical clustering works

Dendrogram

Height

# Let's see how agglomerative hierarchical clustering works

# Hierarchical versus partitional clustering

- Partitional clustering: division of the set of data points into clusters
  - each data object belongs to one subset

- Hierarchical clustering is a set of nested clusters
  - organized as a tree

# K-means clustering

Steps of k-means clustering

1) Choosing k

# K-means clustering

Steps of k-means clustering

1) Choosing k
2) Randomly selecting k data points (as initial centers)
   a) Final result depend on these points
      i) Because k-means converges to one of many possible local minima

# K-means clustering

Steps of k-means clustering

1) Choosing k
2) Randomly selecting k data points (as initial centers)
3) Calculating distance of every data point to the chosen centers in step (2)

# K-means clustering

Steps of k-means clustering

1) Choosing k
2) Randomly selecting k data points (as initial centers)
3) Calculating distance of every data point to the chosen centers in step (2)
4) Assign each data point to the nearest center

# K-means clustering

Steps of k-means clustering

1) Choosing k
2) Randomly selecting k data points (as initial centers)
3) Calculating distance of every data point to the chosen centers in step (2)
4) Assign each data point to the nearest center
5) Calculate new center of each cluster

# K-means clustering

Steps of k-means clustering

1) Choosing k
2) Randomly selecting k data points (as initial centers)
3) Calculating distance of every data point to the chosen centers in step (2)
4) Assign each data point to the nearest center
5) Calculate new center of each cluster
6) Repeat steps (3) to (5) until convergence
   a) No point in changing cluster membership

# Implementing k-means manually (simple example)

1) Choosing k (good guess k=3)

# Implementing k-means manually (simple example)

1) Choosing k (good guess k=3)
2) Randomly selecting 3 data points
   (as initial centers)

# Implementing k-means manually (simple example)

1) Choosing k (good guess k=3)
2) Randomly selecting 3 data points (as initial centers)
3) Calculating distance of every data point to the chosen centers in step (2)

# Implementing k-means manually (simple example)

1) Choosing k (good guess k=3)
2) Randomly selecting 3 data points (as initial centers)
3) Calculating distance of every data point to the chosen centers in step (2)
4) Assign each data point to the nearest center

# Implementing k-means manually (simple example)

1) Choosing k (good guess k=3)
2) Randomly selecting 3 data points (as initial centers)
3) Calculating distance of every data point to the chosen centers in step (2)
4) Assign each data point to the nearest center
5) Calculate new center of each cluster

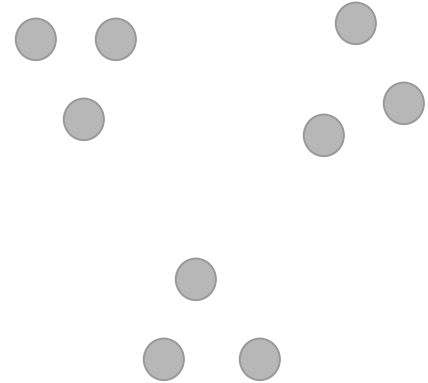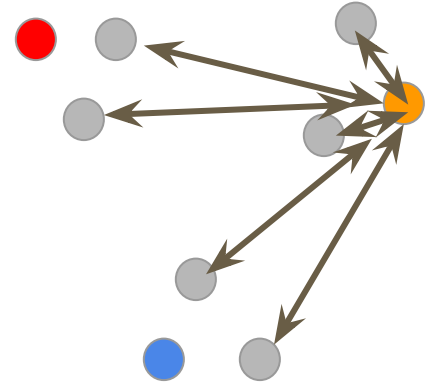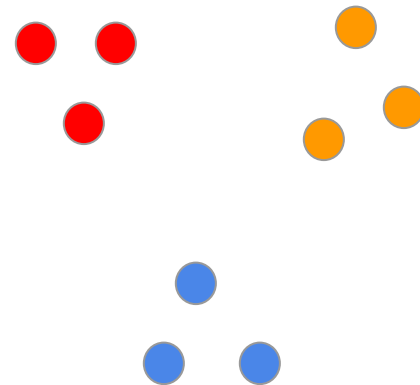# Implementing k-means manually (simple example)

1) Choosing k (good guess k=3)
2) Randomly selecting 3 data points (as initial centers)
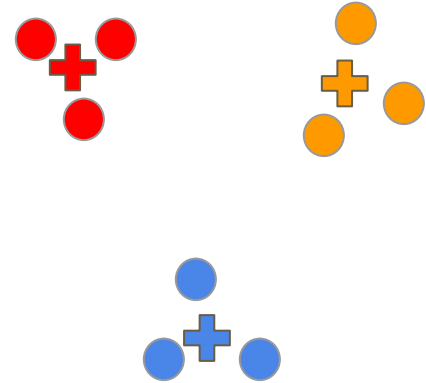3) Calculating distance of every data point to the chosen centers in step (2)
4) Assign each data point to the nearest center
5) Calculate new center of each cluster
6) No need for repetition
   a) Good choice of initial centers
      i) Fast convergence

# Elbow method for selecting optimal K

Distortion is the sum of squared distances from each point to its assigned center

Distortion Score Elbow for KMeans Clustering

$elbow at k = 7, score = 26333.181$

Distortion

fit time (seconds)

k

# Hierarchical k-means (divisive hierarchical clustering)

- Start with all the data points

- Implement k-means (k=2) in an iterative manner

- Until reaching singletons (data points)

# Hierarchical k-means (divisive hierarchical clustering)

Issue with this process:
- Some real clusters could be dismissed because of iterative k-means

# Hierarchical k-means (divisive hierarchical clustering)

Issue with this process:
- Some real clusters could be dismissed because of iterative k-means

# Hierarchical k-means (divisive hierarchical clustering)

Issue with this process:
- Some real clusters could be dismissed because of iterative k-means

We like these to
be one cluster

# Affinity propagation

- Number of clusters (k) does not need to be determined

# Affinity propagation

- Number of clusters (k) does not need to be determined
- Using similarity matrix
    - Need pairwise similarity between data points
        - Similarity between data points $i$ and $j$: negative euclidean distance between data points $i$ and $j$

# Affinity propagation

- Number of clusters (k) does not need to be determined
- Using similarity matrix
  - Need pairwise similarity between data points
    - Similarity between data points $i$ and $j$: negative euclidean distance between data points $i$ and $j$
- Determines clusters and representative data points for the clusters (exemplars)

# Affinity propagation

- Number of clusters (k) does not need to be determined
- Using similarity matrix
  - Need pairwise similarity between data points
    - Similarity between data points $i$ and $j$: negative euclidean distance between data points $i$ and $j$
- Determines clusters and representative data points for the clusters (exemplars)
- Iterative massage passing
  - Data points compete to become exemplars
    - Exemplars are real data points not centers as in k-means

# Affinity propagation

- Number of clusters (k) does not need to be determined
- Using similarity matrix
    - Need pairwise similarity between data points
        - Similarity between data points $i$ and $j$: negative euclidean distance between data points $i$ and $j$
- Determines clusters and representative data points for the clusters (exemplars)
- Iterative massage passing
    - Data points compete to become exemplars
        - Exemplars are real data points not centers as in k-means
- Initialization independent

# Message matrices

- **Responsibility matrix**: $R(i,k)$
  - Responsibility massage from $i$ to $k$
    - Accumulated evidence for how well-suited $k$ is to serve s the exemplar for $i$

# Message matrices

- **Responsibility matrix**: $R(i,k)$
  - Responsibility massage from $i$ to $k$
    - Accumulated evidence for how well-suited $k$ is to serve as the exemplar for $i$
- **Availability matrix**: $A(i,k)$
  - Availability massage from $k$ to $i$
    - Accumulated evidence for how appropriate it is for $i$ to choose $k$

# Sending responsibility and availability



Frey, Brendan J., and Delbert Dueck. "Clustering by passing messages between data points." *science* 315.5814 (2007): 972-976.

# Iterations of affinity propagation for 2 dimensional data points



INITIALIZATION · ITERATION #1 · ITERATION #2 · ITERATION #3 · ITERATION #4 · ITERATION #5 · ITERATION #6 · CONVERGENCE

non-exemplar ▬▬▬ exemplar

Frey, Brendan J., and Delbert Dueck. "Clustering by passing messages between data points." *science* 315.5814 (2007): 972-976.

# DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

Let's look at clusters from another point of view:
- ● Clusters is a maximal set of density-connected points

Ester, Martin, et al. "A density-based algorithm for discovering clusters in large spatial databases with noise." *Kdd*. Vol. 96. No. 34. 1996.

# DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

Let's look at clusters from another point of view:
- Clusters is a maximal set of density-connected points

Two main parameters:
- **Epsilon (Eps)**: Maximum radius of neighbourhood
- **Minimum number of points (MinPts):** Minimum number of points in the Eps-neighbourhood of a data point

Ester, Martin, et al. "A density-based algorithm for discovering clusters in large spatial databases with noise." *Kdd*. Vol. 96. No. 34. 1996.

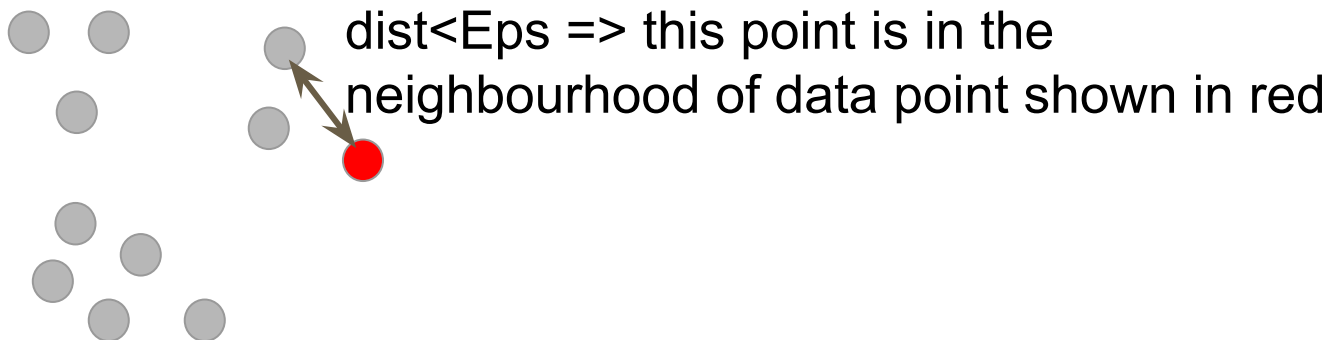# DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

Two main parameters:
- **Epsilon (Eps)**: Maximum radius of neighbourhood
- **Minimum number of points (MinPts):** Minimum number of points in the Eps-neighbourhood of a data point

dist<Eps => this point is in the neighbourhood of data point shown in red

# DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

Two main parameters:
- **Epsilon (Eps)**: Maximum radius of neighbourhood
- **Minimum number of points (MinPts):** Minimum number of points in the Eps-neighbourhood of a data point
  - Then that data point will be called a **core** point

Different types of data points:
- **Core**
- **Border**: within epsilon distance does not meet MinPts criteria
  - But there is at least 1 core point within the epsilon distance
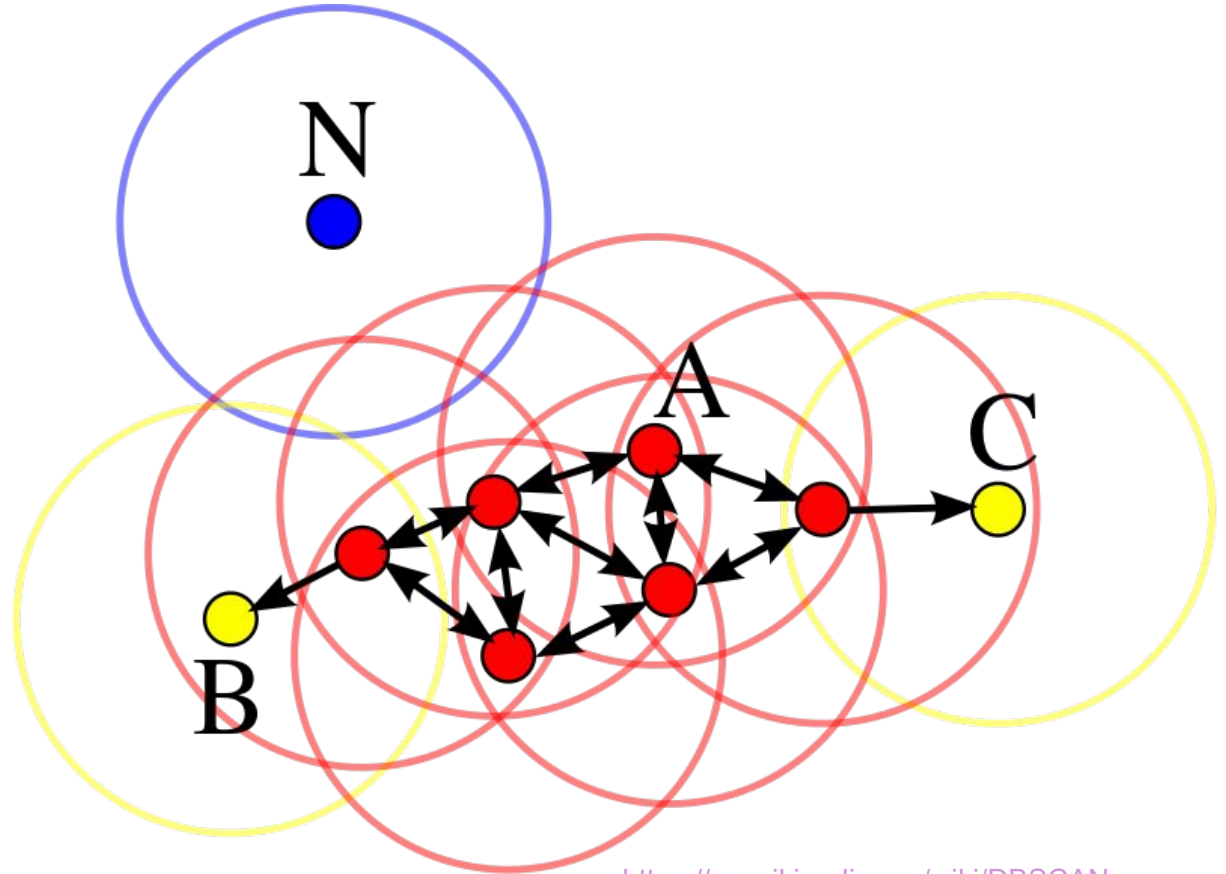- **Noise (Outlier)**: Not assigned to any clusters

# Visualizing different types of data points in DBSCAN
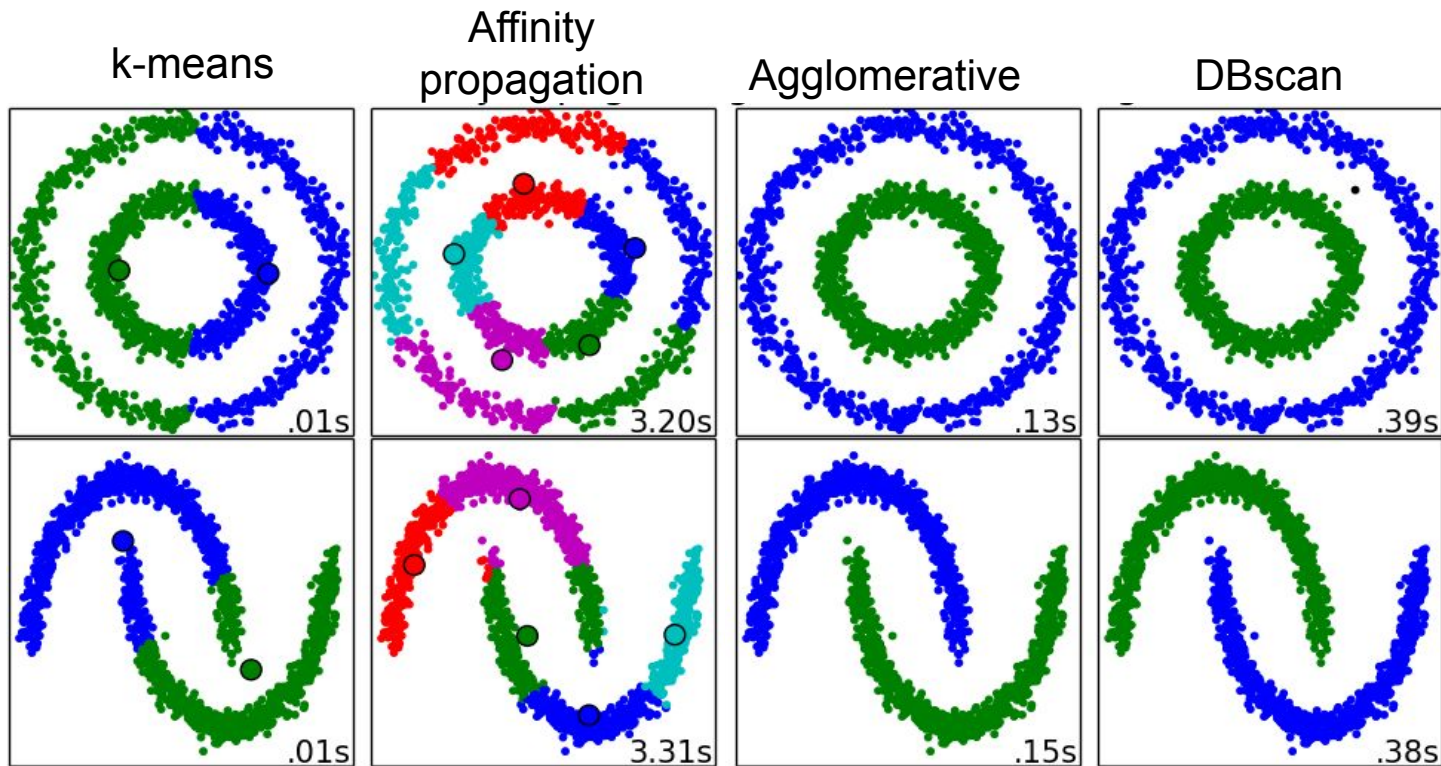
**A**: Core
**B** and **C**: Border
**N**: outlier

# Some disadvantages of DBSCAN

- May not work well with clusters of similar densities
    - But great for separating clusters of low and high densities

- May not be great with very high dimensional data

# Comparison of clustering methods



k-means     Affinity propagation     Agglomerative     DBscan

# Comparison of clustering methods



k-means     Affinity propagation     Agglomerative     DBscan