# Introduction to Data Science
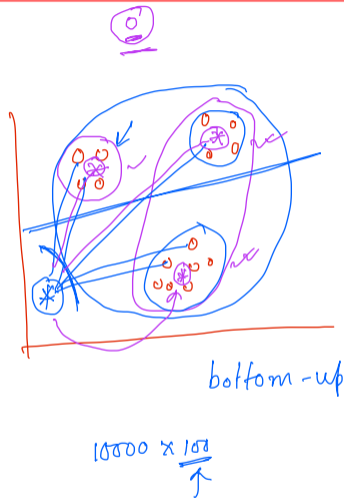
# Clustering

**Arijit Mondal**

**Dept. of Computer Science & Engineering**

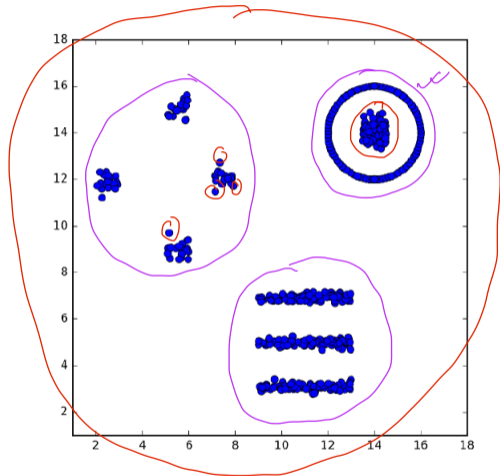**Indian Institute of Technology Patna**

`arijit@iitp.ac.in`

# Introduction

- Clustering is the problem of grouping of points by similarity.
- Splitting of points based on similarity
- Applications
  - Hypothesis development
  - Modeling over smaller subset of data
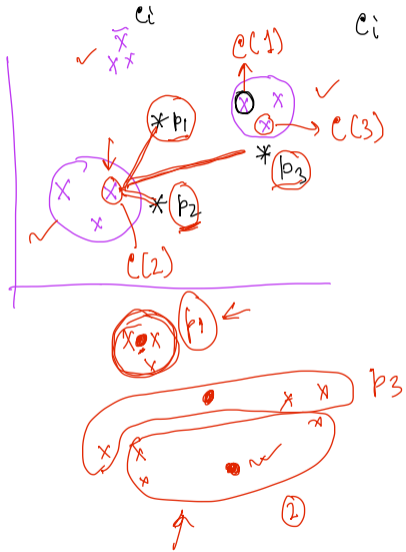  - Data reduction
  - Outliers detection

bottom-up

10000 × 100

# Example

image source: Data Science Design Manual

$k=1$
$k=2$
$\boxed{k=3}$

$$\min\left\{\frac{d(c_i, p_1)}{p_1}, \; d(c_i, p_2), \; d(c_i, p_3)\right\}$$

label $c_i$ with the class of $p_j$ which is closest to $c_i$

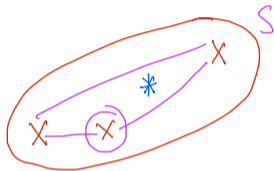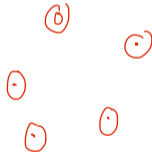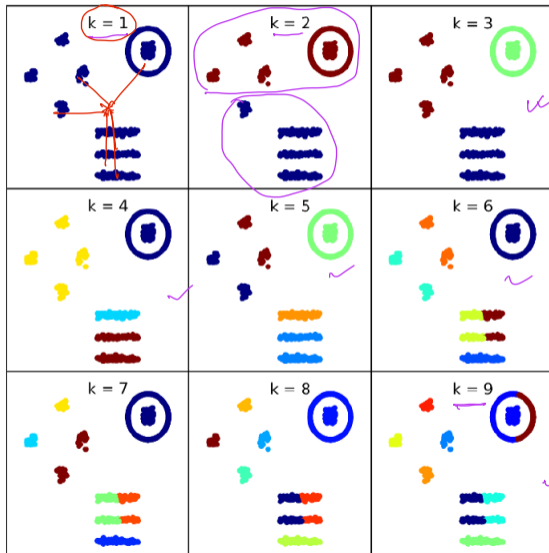$p_i \leftarrow$ Average of the points in the current cluster

# Example

image source: Data Science Design Manual

# Centers vs Centroids

- Centroids: $C_d = \dfrac{1}{|S|} \sum\limits_{p \in S} p[d]$

- Centers: $\underset{c \in S}{\arg\min} \sum\limits_{i=1}^{n} d(c, p_i)$

image source: Data Science Design Manual

Black → E - O

Blue → E > 0

image source: Data Science Design Manual

# Agglomerative clustering

- A bottom-up approach
- Combining similar items
- Distance measures - $C_1, C_2$ are some clusters



$k = 3, 2$

# Agglomerative clustering

- A bottom-up approach
- Combining similar items
- Distance measures - $C_1, C_2$ are some clusters
  - Nearest neighbor - $d(C_1, C_2) = \min\limits_{x \in C_1, y \in C_2} \|x - y\|$

# Agglomerative clustering
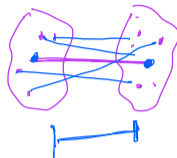
- A bottom-up approach
- Combining similar items
- Distance measures - $C_1, C_2$ are some clusters
  - Nearest neighbor - $d(C_1, C_2) = \min\limits_{x \in C_1, y \in C_2} \|x - y\|$
  - Average link - $d(C_1, C_2) = \dfrac{1}{|C_1| \cdot |C_2|} \sum\limits_{x \in C_1} \sum\limits_{y \in C_2} \|x - y\|$

# Agglomerative clustering

- A bottom-up approach
- Combining similar items
- Distance measures - $C_1, C_2$ are some clusters
  - Nearest neighbor - $d(C_1, C_2) = \min\limits_{x \in C_1, y \in C_2} \|x - y\|$
  - Average link - $d(C_1, C_2) = \dfrac{1}{|C_1| \cdot |C_2|} \sum\limits_{x \in C_1} \sum\limits_{y \in C_2} \|x - y\|$
  - Nearest centroid - closest centroids, $|C_1| \cdot |C_2|$ point-pairs

# Agglomerative clustering
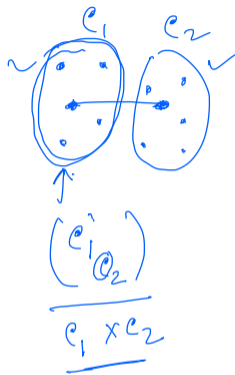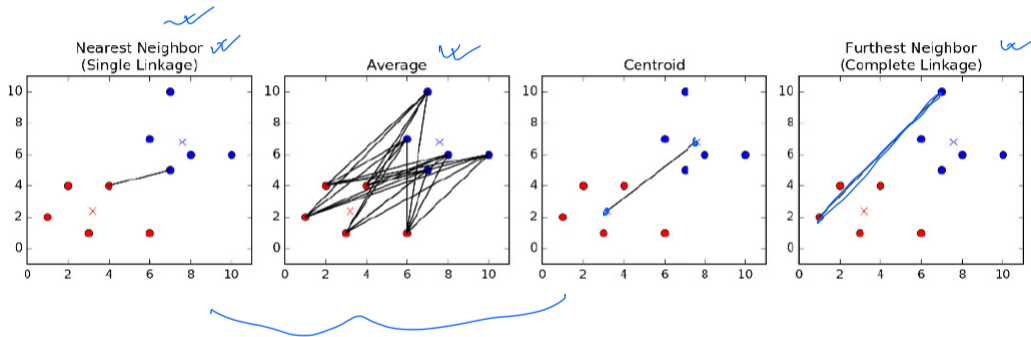
- A bottom-up approach
- Combining similar items
- Distance measures - $C_1, C_2$ are some clusters
  - Nearest neighbor - $d(C_1, C_2) = \min\limits_{x \in C_1, y \in C_2} \|x - y\|$
  - Average link - $d(C_1, C_2) = \dfrac{1}{|C_1| \cdot |C_2|} \sum\limits_{x \in C_1} \sum\limits_{y \in C_2} \|x - y\|$
  - Nearest centroid - closest centroids, $|C_1| \cdot |C_2|$ point-pairs
  - Furthest link - $d(C_1, C_2) = \max\limits_{x \in C_1, y \in C_2} \|x - y\|$
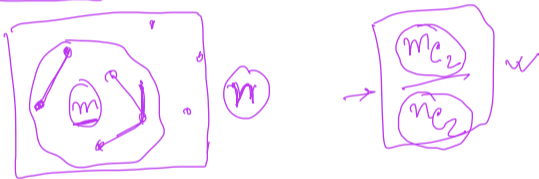
# Distance measures



image source: Data Science Design Manual

# Comparing clustering

- Jaccard similarity: Similarity between two sets, $J(s_1, s_2) = \dfrac{|s_1 \cap s_2|}{|s_1 \cup s_2|} \rightarrow \dfrac{common}{Union}$

- Jaccard distance: $1 - J(s_1, s_2)$. It is distance metric

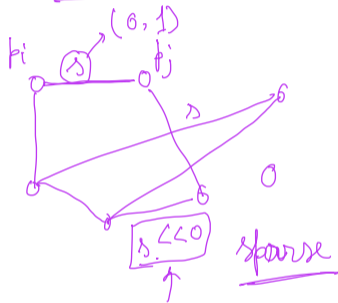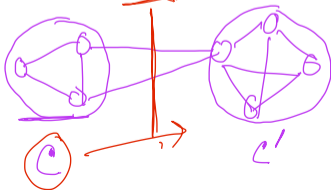- Rand index: ratio of compatible pairs to all possible pairs

# Similarity & Cuts

- Similarity measure - $S[i,j] = e^{-\beta \|p_i - p_j\|}$

- Similarity graph - It is based on similarity measure. Can be made sparse by applying thresholding on similarity values

- Cluster weight - $W(C) = \sum_{x \in C} \sum_{y \in C} S[i,j]$

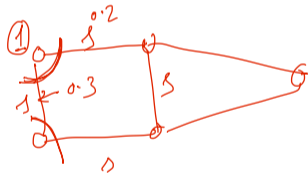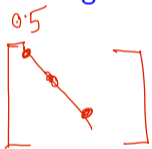- Cut weight - $W'(C) = \sum_{x \in C} \sum_{y \in V-C} S[i,j]$

- Conductance of cluster $C$ is $\dfrac{W'(C)}{W(C)}$

$e^{-x} \to 1 \quad x = 0$
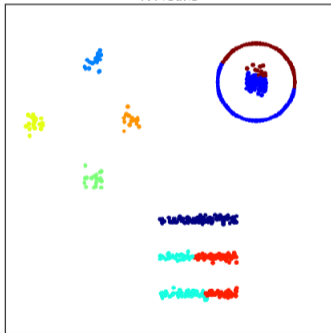
$\to 0 \quad x \to \infty$

# Spectral clustering

- Construct the Laplacian matrix $L = D - S$
  - $S$ - similarity matrix
  - $D$ - degree-weighted identity matrix, $D[i,i] = \sum_j S[i,j]$
- The most valuable eigenvectors for clustering here turn out to have the smallest non-zero eigenvalues
- Applying $k$-means clustering in this feature space produces connected clusters

# Spectral clustering