

## CENTER-BASED CLUSTERING

- SETTING:  $N$  OBSERVATIONS  $\vec{x}_1, \dots, \vec{x}_N \in \mathbb{R}^d$
- GOAL: FIND CLUSTERS (SUBSETS/GROUPS) OF SIMILAR POINTS
- OUTPUT:
  - NO DEMOGRAPHY
  - UNLIKE HIERARCHICAL CLUSTERING
    - # OF CLUSTERS MUST BE SPECIFIED IN ADVANCE
    - (k+1) CLUSTERING OF DATA (LIKELY) DOES NOT RESULT FROM SPLITTING A CLUSTER IN (k) - CLUSTERING RESULT
    - CANNOT BE CONSTRUCTED AS AN ITERATIVE ALGORITHM
      - SPECIFY AN OBJECTIVE AND (TRY TO) OPTIMIZE
  - CLUSTER CENTERS  $\{c_1, \dots, c_k\} \in \mathbb{R}^d$ ,  $k = \#$  DESIRED CLUSTERS  
AND PARTITION  $\mathcal{C} = \{C_1, \dots, C_k\}$ ,  $C_i \subset \{\vec{x}_1, \dots, \vec{x}_N\}$  GIVES  
SAMPLES IN CLUSTER  $i$ .

### METHODS:

#### k-MEANS:

CENTERS: 
$$c_i = \frac{1}{N_i} \sum_{x_j \in C_i} \vec{x}_j, \quad N_i = |C_i| \quad (\text{ASSUMING EUCLIDEAN DISTANCE})$$

#### OBJECTIVE:

$$\Phi(\mathcal{C}) = \sum_{j=1}^k \sum_{x_i \in C_j} d^2(x_i, \vec{c}_j)$$

← TYPICALLY EUCLIDEAN DISTANCE  
← MINIMIZE SUM OF SQUARES

SQUARED DISTANCE PLACES

GREATER EMPHASIS ON LARGE

DISTANCES!

## k-CENTER

OBJECTIVE:

$$\Phi_{\text{center}}(C) = \max_{j=1, \dots, k} \max_{x_i \in G_j} d(x_i, c_j)$$

- \* MINIMIZE THE MAXIMUM DISTANCE BETWEEN ANY DATA POINT AND ITS CENTER
- \* CENTERS  $c_1, \dots, c_k$  DO NOT HAVE EXPLICIT REPRESENTATIONS
  - FOR FIXED  $k$  THIS CAN BE EXPRESSED AS A GRAPH THEORY PROBLEM, NP-COMPLEX ALGORITHMS
  - DETERMINING  $k$  IS NP-HARD
- \* CENTERS  $c_1, \dots, c_k$  GENERATE VORONOI PARTITIONS
- \* REQUIRE APPROXIMATE ALGORITHMS
  - NO GLOBAL OPTIMAL SOLUTIONS w/o A LOT OF LUCK!

## k-MEDIAN / k-MEDOIDS

• CENTERS:

$$c_i = \underset{\tilde{x} \in G_i}{\text{argmin}} \sum_{x_j \in G_i, \{x_j\}} d(\tilde{x}, x_j)$$


• OBJECTIVE:

$$\Phi_{\text{k-med}}(C) = \sum_{j=1}^k \sum_{x_i \in G_j} d(\tilde{x}_i, c_j)$$

↑ USED DISTANCE (NOT SQUARED)  
AS OPPOSED TO k-MEANS

- \* CENTERS ARE GIVEN BY ACTUAL DATA POINTS (REPRESENTATIVES!)
- \* A VARIANT CALLED PARTITIONING AROUND MEDOIDS (PAM)  
ALSO EXISTS.

## COMPARISON

| k-CENTER   | k-MEDOID   | k-MEANS   |
|--|--|---|
| <p>NEED TO COMPUTE CENTERS FROM DATA</p> <p>- NEED <math>\vec{x}_1, \dots, \vec{x}_n</math></p>  | <p>ONLY<br/>✓<br/>USES ORIGINAL DISTANCE AT EACH STEP</p> <p>- ONLY REQUIRES D</p> <p>- DON'T NEED QUANT. DATA</p> | <p>NEED TO COMPUTE CENTERS FROM DATA</p> <p>- NEED <math>\vec{x}_1, \dots, \vec{x}_n</math></p> <p>- NEED METRIC: QUANTITATIVE DATA</p> |
| <p>USES <math>d(x_i, c_j)</math></p> <p>- LESS SENSITIVE TO OUTLIERS</p>   | <p>USES <math>d(x_i, c_j)</math></p> <p>- LESS SENSITIVE TO OUTLIERS</p>   | <p>USES <math>d^2(x_i, c_j)</math></p> <p>- GREATER EMPHASIS ON LARGE DISTANCES</p> <p>- SENSITIVE TO OUTLIERS</p>                      |
| <p>CENTERS ESTIMATED FROM DATA</p> <p>- CENTERS NEED NOT BE REPRESENTATIVE</p>  | <p>CENTERS CHOSEN FROM DATA</p> <p>"REPRESENTATIVE OBJECTS"</p>  | <p>CENTERS ESTIMATED FROM DATA</p> <p>- CENTERS NEED NOT BE REPRESENTATIVE</p>  |

## MINIMIZING OBJECTIVES

### SIMPLE (COMBINATORIAL) ALGORITHM

FOR FIXED  $K$ , WE CAN TRY ALL POSSIBLE CLUSTERING ASSIGNMENTS

- # OF DISTINCT ASSIGNMENTS

$$S(N, k) = \frac{1}{k!} \sum_{j=1}^k (-1)^{k-j} \binom{k}{j} j^N$$

$\nwarrow$   
 LABEL SWITCHING

$$S(10, 4) = 34,105$$

$$S(19, 4) \approx 10^{13}$$

EXPONENTIAL IN  $N!$

\* MUST USE ITERATIVE/APPROX. / GREEDY ALGORITHMS

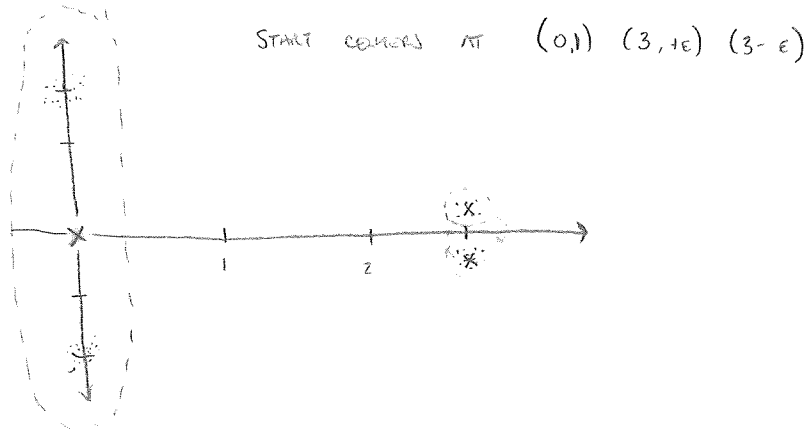
# Ex: k-means (Lloyd's Algorithm)

- 1) START w/ k CENTERS
- 2) CLUSTER EACH POINT TO ITS NEAREST CENTER
- 3) FIND THE NEW CLUSTER CENTERS

$$* C_i = \frac{1}{N_i} \sum_{x_j \in C_i} \vec{x}_j \quad \text{IN CASE OF EUCLIDEAN DISTANCE}$$

- 4) REPEAT 2, 3 TO CONVERGENCE!

a)



b)



CENTERS A 0, 5, 10



| CENTER | MEMBERS |
|--------|---------|
| 2      | 2, 3    |
| 5      | $\phi$  |
| 8      | 7, 8    |



?

- \* COLLAPSE OF ALGORITHM IF TOO MANY CLUSTERS CHOSEN (TRY MANY k)  
IF UNKNOWN
- \* POOR FIT IF BAD INITIAL CONDITIONS CHOSEN (TRY MANY ICs!)

## CHOOSING THE NUMBER OF CLUSTERS

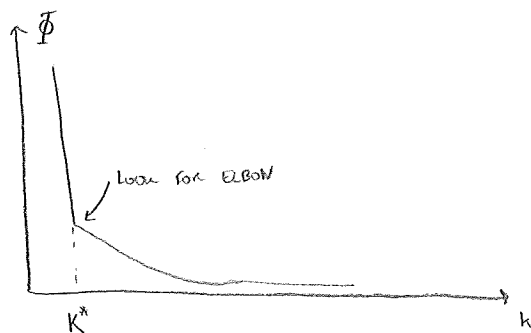
- A LOT OF DIFFERENT METHODS

- INFORMAL (SCREE PLOT) OPTION

ASSUMING OPTIMAL SOLUTIONS ARE FOUND

$$\Phi_{k\text{-means}}, \Phi_{k\text{-median}}, \Phi_{k\text{-med}}$$

ARE DECREASING IN THE # OF CLUSTERS



- COVARIANCE DECOMPOSITION FOR k-MEANS w/ EUCLIDEAN DISTANCE

- GIVEN CLUSTERING ASSIGNMENTS  $C(1), \dots, C(N) \in \{1, \dots, k\}$

$$\begin{aligned} N \hat{\Sigma} &= \sum_{j=1}^N (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})^T = \sum_{j=1}^k \sum_{\mathbf{x}_i \in C_j} (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \\ &= \sum_{j=1}^k \sum_{\mathbf{x}_i \in C(j)} \left( \begin{aligned} &(\mathbf{x}_i - \mathbf{c}_j)(\mathbf{x}_i - \mathbf{c}_j)^T + (\mathbf{c}_j - \bar{\mathbf{x}})(\mathbf{c}_j - \bar{\mathbf{x}})^T \\ &+ (\mathbf{x}_i - \mathbf{c}_j)(\mathbf{c}_j - \bar{\mathbf{x}})^T + (\mathbf{c}_j - \bar{\mathbf{x}})(\mathbf{x}_i - \mathbf{c}_j)^T \end{aligned} \right) \\ &\quad \sum_{\mathbf{x}_i \in C(j)} (\mathbf{x}_i - \mathbf{c}_j) = \mathbf{0} \\ &= \underbrace{\sum_{j=1}^k \sum_{\mathbf{x}_i \in C(j)} (\mathbf{x}_i - \mathbf{c}_j)(\mathbf{x}_i - \mathbf{c}_j)^T}_{\text{TOTAL WITHIN GROUP COVARIANCE}} + \underbrace{\sum_{j=1}^k N_j (\mathbf{c}_j - \bar{\mathbf{x}})(\mathbf{c}_j - \bar{\mathbf{x}})^T}_{\text{BETWEEN GROUP COVARIANCE}} \\ &= \Phi_{k\text{-means}}(C) + B(C) \end{aligned}$$

CALINSKI ; HARABAS (1974)

CHOOSE  $k$  MAXIMIZING

$$C(k) = \frac{\text{tr}(B(\zeta)) / (k-1)}{\text{tr}(\Phi(\zeta)) / (N-k)}$$

MEAN BETWEEN GROUP VARIANCE

MEAN WITHIN GROUP VARIANCE

\* NOTE CONNECTION / SIMILARITY TO ANOVA

GAP STATISTIC (TIBSHIRANI ET AL 2001)

$$G_{AP_N}(k) = \underbrace{E_N^* \left( \log \Phi_k(\zeta) \right) - \log \Phi_k(\zeta)}$$

GENERATE MULTIPLE

SAMPLES FROM A

SUITABLY CHOSEN NULL

DISTRIBUTION, MINIMIZE

$\Phi_k(\zeta)$  w/  $k$  CLUSTERS

TO ESTIMATE

CHOOSE  $k$  MAXIMIZING  $G_{AP_N}(k)$

\* ALSO RELATED TO HYPOTHESIS TESTING

NOTE: THE REFERENCE DISTRIBUTION IS GENERATED USING PCA