

# K-means Clustering

Tuesday, May 2, 2017 9:53 AM

$$\{x_i\}_{i=1}^n, x_i \in \mathbb{R}^p$$

Task: learn  $z_i \in \{1, \dots, K\}$  (cluster assignments)

eg  $C_1 = \{\underline{x}_1, \underline{x}_3, \underline{x}_2\}$   $z_1=1, z_2=2, z_3=1$

$C_2 = \{\underline{x}_2, \underline{x}_5, \underline{x}_6\}$   $z_4=3$

$C_3 = \{\underline{x}_4, \underline{x}_8, \underline{x}_7\}$

also learn  $m_k \in \mathbb{R}^p$ , cluster centers,  $k=1, \dots, K$

Objective of K-means

$$\min_{z, m} J(z, m) = \sum_{i=1}^n \|x_i - m_{z_i}\|_2^2$$

↳ distortion

Write this:  $J(z, m) = \sum_{k=1}^K \underbrace{\sum_{i: z_i=k} \|x_i - m_k\|_2^2}_{\text{fix } z_i \text{'s}}$

min  $m_k$  = average w/in clusters

## Lloyd's Algorithm

(1) initialize  $m_k$  arbitrarily

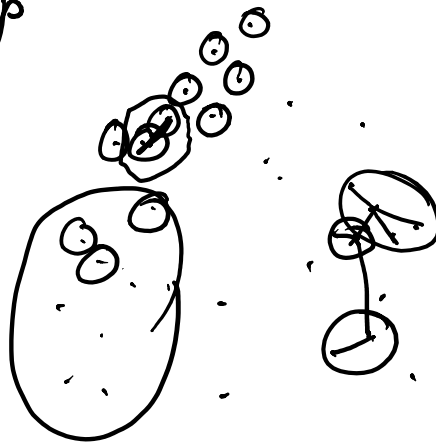
(2) Alternate

(a) Update  $z_i \leftarrow \arg \min_k \|x_i - m_k\|_2^2$

(b)  $m_k \leftarrow \frac{\sum_{i=1}^n 1\{z_i=k\} x_i}{\sum_{i=1}^n 1\{z_i=k\}}$

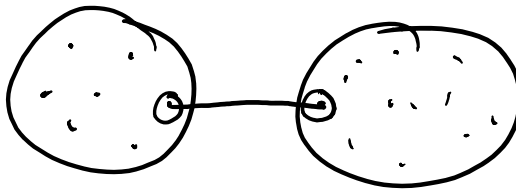
## Agglomerative Clustering: Bottom-up

- (1) Start w/ all data points in own cluster
- (2) Find clusters most similar:  $C_1$  &  $C_2$  (\*)
- (3) Merge  $C_1$  &  $C_2$  goto 2



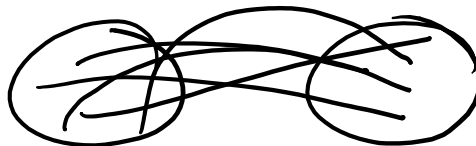
## Cluster similarity

Single linkage:



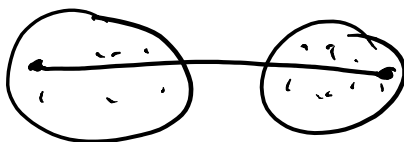
$$d_{sl}(C_1, C_2) = \min_{x \in C_1, y \in C_2} d(x, y)$$

Average linkage:



$$d_{al}(C_1, C_2) = \frac{1}{|C_1||C_2|} \sum_{x \in C_1, y \in C_2} d(x, y)$$

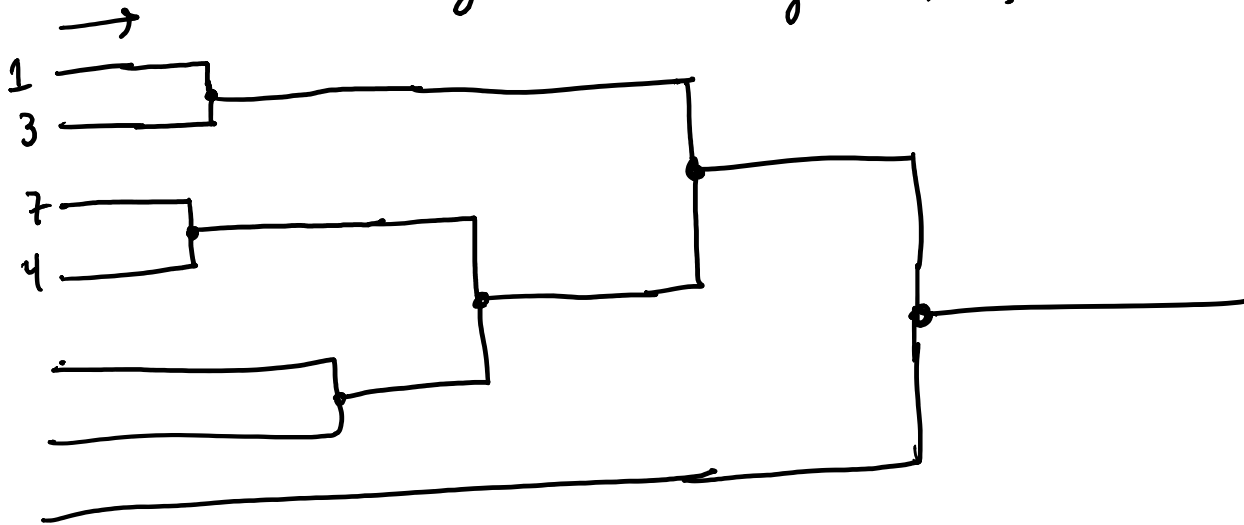
Complete linkage:  $d_{cl}(C_1, C_2) = \max_{x \in C_1, y \in C_2} d(x, y)$



General observations: sl tends to yield "unbalanced"

"clusters" and  $cl$  yields "balanced clusters"

Dendrogram Visualization tool where we start w/  
 $n$  lines and merge when merge clusters



Divisive Clustering Top-down

Option 1: recursively apply  $k$ -means w/  $k=2$

Con 1: initialization of  $k$ -means really matters

Con 2: algorithm can violate monotonicity  
of distortion

Option 2: greedy approach

(1) Start w/ 1 cluster

(2) Repeat

- Choose a cluster  $G$
  - Remove point most dissimilar from  $G$
- Starts cluster  $H$

$$\left. \begin{array}{l} \text{repeat} \\ \left\{ \begin{array}{l} - \text{remove } x^* = \operatorname{argmax}_{x \in G} \frac{1}{|G|-1} \sum_{g \in G \setminus \{x\}} d(x, g) - \frac{1}{|H|} \sum_{h \in H} d(x, h) \\ \text{add } x^* \rightarrow H \end{array} \right. \end{array} \right\}$$

# Document Vectorization

Tuesday, May 2, 2017 9:59 AM

## Bag-of-words models

- ▷ Construct a dictionary of words : enumerate all words in our corpus (all documents)
- ▷ vectorize a document,  $D$ ,

$X_{Di}$  = measure of frequency of word  $i$  in  $D$ .

eg  $X_{Di} = 1\{i \in D\}$  (binary TF)

$$\begin{aligned}\|X_D - X_{D'}\|_2^2 &= \sum_i (1\{i \in D\} - 1\{i \in D'\})^2 \\ &= \sum_i 1\{i \in D, i \notin D' \text{ or } i \in D', i \notin D\}\end{aligned}$$

other frequencies: count of word,  $\frac{\text{count}}{\text{total \# of words in } D}$

## Inverse document freq

$$\text{idf}(i) = \log\left(\frac{\# \text{ of doc's in corpus}}{\{D : i \in D\}}\right)$$

## Term-frequency Inverse-document-frequency

$$\text{tfidf}(i, D) = \text{tf}(i, D) \cdot \text{idf}(i)$$