

Topic 1: CLUSTERING

CMSC 35400/STAT 37710 Machine Learning
Risi Kondor, The University of Chicago

Clustering

In modern ML, more often than not, the inputs are high dimensional real vectors:

$$\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d.$$

Each x_i is called a **feature** (**covariate** in Stats).

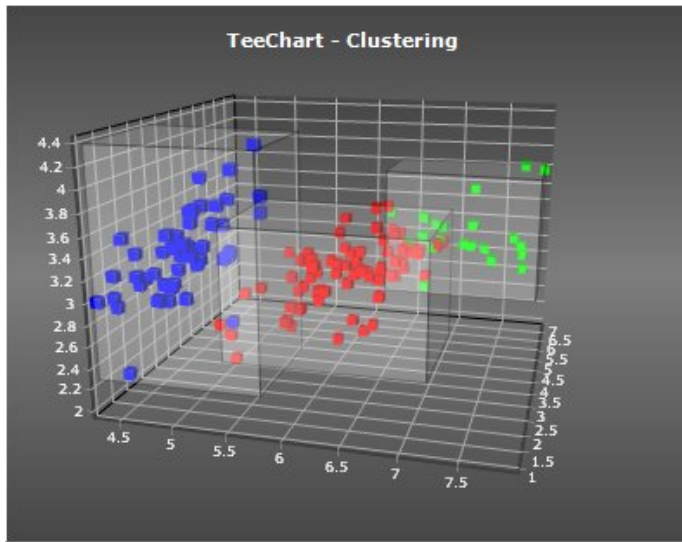
Example: $x_1 = \text{age}$, $x_2 = \text{weight}$, $x_3 = \text{blood pressure}$, ...

Example: $x_i = \text{intensity of a pixel } i \text{ in an image}$

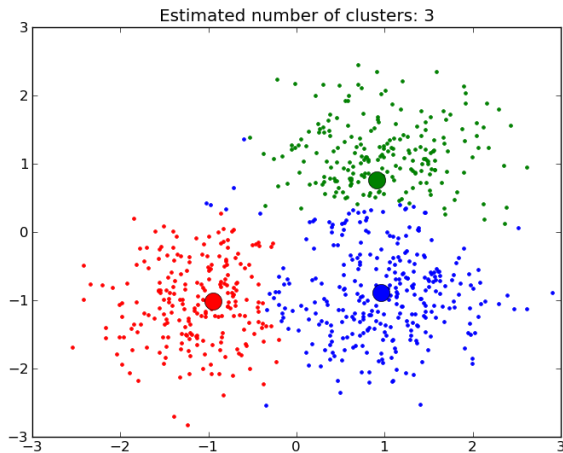
It often makes sense to ask whether a dataset $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ can be partitioned into a small number of **clusters** of similar datapoints.

→ Clustering is a typical unsupervised learning problem.

Clustering

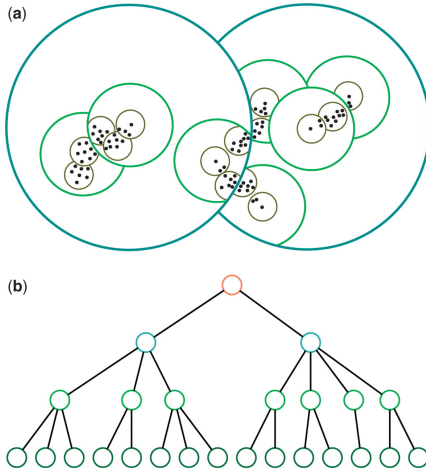


Clustering



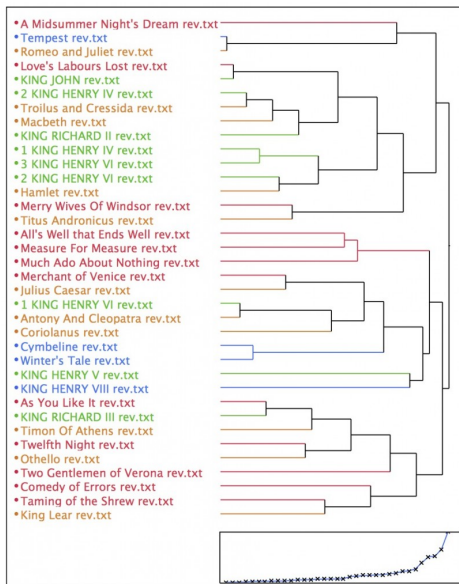
Cluster representatives indicated.

Hierarchical clustering

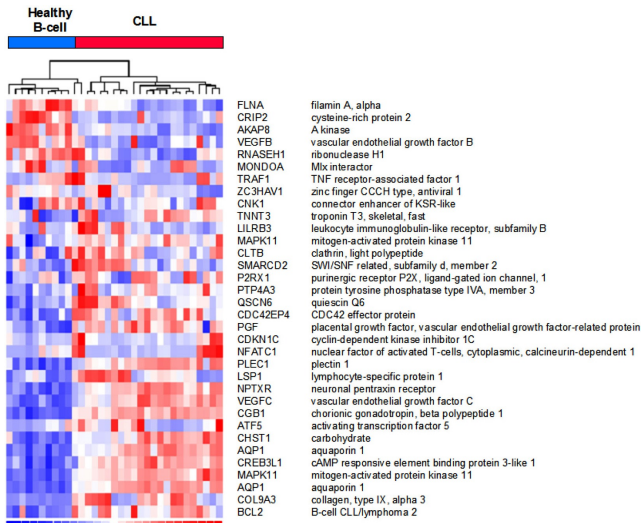


Cutting the tree at any level gives a flat clustering. Thanks to this freedom, don't have to decide the number of clusters in advance.

Hierarchical clustering

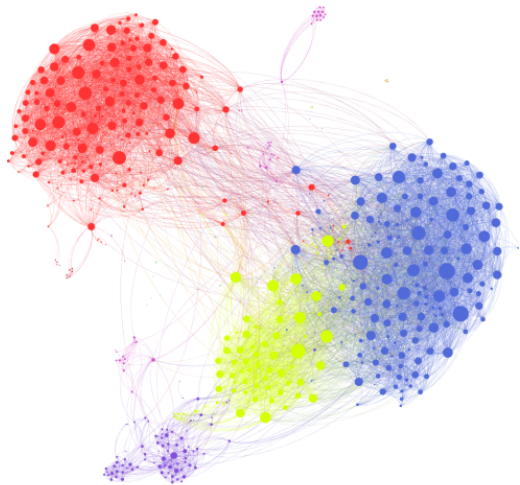


Hierarchical clustering



[Pallasch et al., Blood, 2009]

Clustering of nodes in a graph



Also known as **graph partitioning** (these are somebody's Facebook friends).

Clustering: the good

Clustering is important because

- It is a natural thing to want to do with large data.
- Can reveal a lot about the structure of data → exploratory data analysis.
e.g., finding new types of stars, patients with similar disease profiles, ...
- Allows us to compress data by replacing points by their cluster representatives (called **vector quantization**).
- Key part of finding structure in large graphs & networks.

Clustering: the bad

- Unsupervised problem → always harder to formalize.
- Ill-defined: different objective functions possible, no clear winner. Even after we've clustered the data it's hard to say whether the clustering is good or bad → subjective.
- What is the “correct” number of clusters? Also subjective. Often data is very ambiguous in this regard.
- End users may attribute too much significance to the clusters with unforeseeable consequences.
- Compared to supervised ML, the theory is in its infancy.

Outline

1. Flat clustering: k -means
2. Hierarchical clustering: agglomerative clustering
3. Model based clustering: mixture of Gaussians

Flat clustering

Flat clustering

Input: the datapoints $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$;
the desired number of clusters $k \in \mathbb{N}$.

Output: k disjoint sets C_1, C_2, \dots, C_k whose union is $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$.

Clustering is driven by a distance metric, d . In the simplest case it is just the Euclidean distance

$$d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\| = \left(\sum_{i=1}^d (x_i - x'_i)^2 \right)^{1/2}.$$

Let's assign each cluster a representative point \mathbf{m}_i . Depending on context, we might or might not require \mathbf{m}_i to be one of the $\mathbf{x}_1, \dots, \mathbf{x}_n$ datapoints.

Cost functions

Start with a **cost function** (in this context also called **distortion**) that our algorithm tries minimize:

- Max distance to cluster center:

$$J_{\max} = \max_{i \in \{1, \dots, k\}} \max_{\mathbf{x} \in C_i} d(\mathbf{x}, \mathbf{m}_i).$$

- Average distance to cluster center:

$$J_{\text{avg}} = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} d(\mathbf{x}, \mathbf{m}_i).$$

- Average squared distance to cluster center:

$$J_{\text{avg}^2} = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} d(\mathbf{x}, \mathbf{m}_i)^2.$$

- Sum of squared intra-cluster distances:

$$J_{\text{IC}} = \sum_{i=1}^k \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \sum_{\mathbf{x}' \in C_i} d(\mathbf{x}, \mathbf{x}')^2.$$

(Prove that $J_{\text{IC}} \sim J_{\text{avg}^2}$)

The k -means algorithm

Problem: find C_1, C_2, \dots, C_k and centroids $\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_k \in \mathbb{R}^d$ that minimize

$$J_{\text{avg}^2} = \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} d(\mathbf{x}, \mathbf{m}_j)^2,$$

where $d(\mathbf{x}, \mathbf{m}_j) = \|\mathbf{x} - \mathbf{m}_j\|$.

This is an **optimization problem**.

- Is it continuous? **No**. Is it combinatorial? **No**. \rightarrow **Mixed**.
- Is it convex? **No**.
- How do we solve it? **Alternating minimization strategy**.

The k -means algorithm

Let γ_i be the cluster that \mathbf{x}_i is assigned to, i.e., $C_j = \{ \mathbf{x}_i \mid \gamma_i = j \}$.

- Given the $\gamma_1, \gamma_2, \dots, \gamma_n$ cluster assignments, J_{avg^2} is minimized by setting

$$\mathbf{m}_j \leftarrow \frac{1}{|C_j|} \sum_{i: \gamma_i = j} \mathbf{x}_i \quad j = 1, 2, \dots, k.$$

- Given the $\mathbf{m}_1, \dots, \mathbf{m}_k$ cluster centroids, J_{avg^2} is minimized by setting

$$\gamma_i = \operatorname{argmin}_{j \in \{1, 2, \dots, k\}} d(\mathbf{x}, \mathbf{m}_j) \quad i = 1, 2, \dots, n.$$

The k -means algorithm

```
{ $\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_k$ }  $\leftarrow k$  random points in  $\Omega$  ;  
while(convergence){  
     $C_1, C_2, \dots, C_k \leftarrow \emptyset$  ;  
    for  $i=1$  to  $n$  {                                     // Assign each point to the closest center  
         $\hat{j} \leftarrow \arg \min_{j \in \{1, \dots, k\}} d(\mathbf{x}_i, \mathbf{m}_j)$  ;  
         $C_{\hat{j}} \leftarrow C_{\hat{j}} \cup \{\mathbf{x}_i\}$  ;  
    }  
    for  $j=1$  to  $k$                                        // Recompute cluster centers  
         $\mathbf{m}_j \leftarrow \frac{1}{|C_j|} \sum_{\mathbf{x} \in C_j} \mathbf{x}$  ;  
}
```

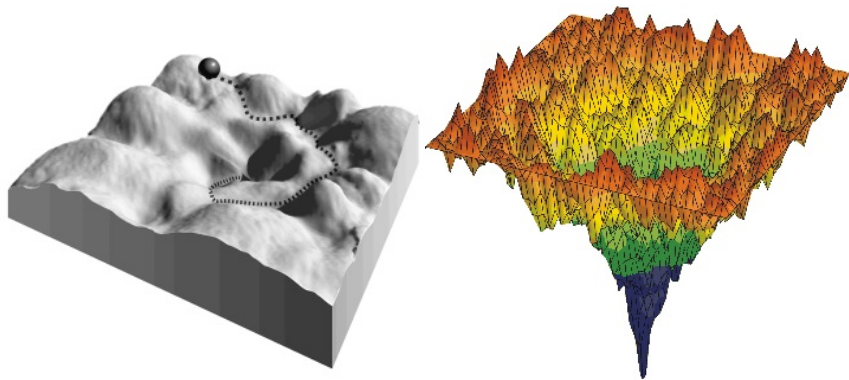
The k -means algorithm

- Probably the most popular clustering algorithm.
- Effectively does alternating minimization on

$$J_{\text{avg}^2} = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} d(\mathbf{x}, \mathbf{m}_i)^2.$$

- Converges in a finite number of steps (Why?) but best upper bound is n^{kd} [Inaba et al., 1989].
- Finding the optimal clustering is NP-hard for general d (even for $k = 2$) or general k (even $d = 2$) [Dasgupta et al., 2009]
- There is no guarantee that the algorithm converges to the globally optimal solution (in most cases it won't). This is a serious problem. Often end up with some clusters only having a single datapoint. Solutions:
 - Random restarts
 - Merge clusters that are too small
 - Split clusters that are too large
 - Annealing and other methods for dealing with complicated energy surfaces
 - etc.

Local vs. global minima



Complicated energy landscapes with lots of local minima are the bane of modern science (ML, optimization, protein folding, etc.).

k -means++

Arthur and Vassilvitskii (2007)

k -means++

choose \mathbf{m}_1 uniformly at random from $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$

for($i = 2$ to k) {

 choose \mathbf{m}_i from $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ with probability

$$p(\mathbf{m}_i = \mathbf{x}_j) = \frac{(D_{i-1}(\mathbf{x}_j))^2}{\sum_{\ell} (D_{i-1}(\mathbf{x}_{\ell}))^2}$$

 where

$$D_{i-1}(\mathbf{x}) = \min_{p \in \{1, 2, \dots, i-1\}} \|\mathbf{x} - \mathbf{m}_p\|.$$

}

Run k -means initialized with $(\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_k)$ as usual

k -means++

Theorem [Arthur and Vassilvitskii (2007)] Let m_1, m_2, \dots, m_k be the initial cluster centers returned by the k -means++ initialization procedure. Then

$$\mathbb{E}[J_{\text{avg}^2}(m_1, m_2, \dots, m_k)] \leq 8(\ln k + 2)J_{\text{avg}^2}^*,$$

where $J_{\text{avg}^2}^*$ is the minimum of J_{avg^2} over all possible clusterings.