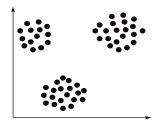
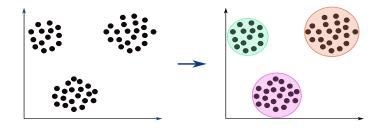
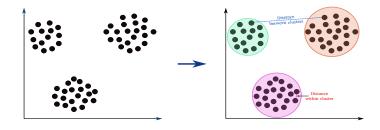
By Van Dinh Tran

November 5, 2023







Definition: A function f is called a clustering function if it takes a set of data points (D) and a distance metric (d) as input and returns a partition (P) on D.

$$P = f(d, D)$$

- $\bullet \ P = \{C_1, C_2, \ldots, C_k\} \ s.t \ C_1 \cup C_2 \cup \ldots \cup C_k = D; \ C_i \cap C_j = \Phi$
- C_i are called clusters

Definition: A function f is called a clustering function if it takes a set of data points (D) and a distance metric (d) as input and returns a partition (P) on D.

$$P = f(d, D)$$

- $\bullet \ P = \{C_1, C_2, \ldots, C_k\} \ s.t \ C_1 \cup C_2 \cup \ldots \cup C_k = D; \ C_i \cap C_j = \Phi$
- C_i are called clusters

Three desirable properties of a clustering algorithm.

- Scale invariance: $f(d, D) = f(\alpha d, D)$, α is a scalar
- \bullet Richness: any possible P on D should be a possible outcome of f
- Consistency: If we apply f on D using d' and it reduces within distances and increases between distances, f(d, D) = f(d', D) = P.

Definition: A function f is called a clustering function if it takes a set of data points (D) and a distance metric (d) as input and returns a partition (P) on D.

$$P = f(d, D)$$

- $\bullet \ P = \{C_1, C_2, \ldots, C_k\} \ s.t \ C_1 \cup C_2 \cup \ldots \cup C_k = D; \ C_i \cap C_j = \Phi$
- C_i are called clusters

Three desirable properties of a clustering algorithm.

- Scale invariance: $f(d, D) = f(\alpha d, D)$, α is a scalar
- Richness: any possible P on D should be a possible outcome of f
- Consistency: If we apply f on D using d' and it reduces within distances and increases between distances, f(d, D) = f(d', D) = P.

[5ex] No function that satisfies those three properties

⇒ Constraint relaxation

• Common distance metrics

- Euclidean distance: $d(x, y) = \sqrt{\sum_{i=1}^{n} (x_i y_i)^2}$
- Manhattan distance: $d(x,y) = \sum_{i=1}^{n} |x_i y_i|$
- $\blacktriangleright \ \ \text{Peason correlation distance:} \ \ d\big(x,y\big) = 1 \frac{\sum_{i=1}^n (x_i \hat{x})(y_i \hat{y})}{\sqrt{\sum_{i=1}^n (x_i \hat{x})^2 \sum_{i=1}^n (y_i \hat{y})^2}}$
- ▶ Kernel-based distance

Common distance metrics

- Euclidean distance: $d(x,y) = \sqrt{\sum_{i=1}^{n} (x_i y_i)^2}$
- \blacktriangleright Manhattan distance: $d(x,y) = \sum_{i=1}^n |x_i y_i|$
- $\blacktriangleright \ \ \text{Peason correlation distance:} \ \ d\big(x,y\big) = 1 \frac{\sum_{i=1}^n (x_i \hat{x})(y_i \hat{y})}{\sqrt{\sum_{i=1}^n (x_i \hat{x})^2 \sum_{i=1}^n (y_i \hat{y})^2}}$
- Kernel-based distance

• Metrics to measure clustering performance

- ▶ Internal metrics
 - Measures of cohesion and separation: Within-cluster sum of squares (WSS) and Between-cluster sum of squares (BSS)
 - Measures of validity: Silhouette score, Davies-Bouldin index, etc
- External metrics: Rand index, adjusted rand index

• Some types of clustering methods

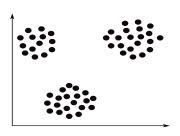
- Centroid models (e.g k-means): represents each cluster by a single mean vector, each point is assigned to the cluster whose nearest centroid.
- Distribution models: clusters are modeled using statistical distributions, such as multivariate normal distributions used by the expectation-maximization algorithm.
- Connectivity models (e.g hierarchical clustering): build models based on distance connectivity
- Density models (e.g DBSCAN) defines clusters as connected dense regions in the data space.

Dataset $D = \{x_i | x_i \in R^n\}$, a distance metric d

- 1. Specify number of clusters K
- 2. Initialize K data points as cluster centroids
- 3. Keep iterating until there is no change to the centroids
 - 3.1 Assign each data point to the closest cluster centroid.
 - 3.2 Compute the centroids for the clusters by taking the average of the all data points that belong to each cluster.

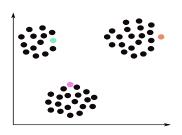
Dataset $D = \{x_i | x_i \in R^n\}$, a distance metric d

- 1. Specify number of clusters K
- 2. Initialize K data points as cluster centroids
- 3. Keep iterating until there is no change to the centroids
 - 3.1 Assign each data point to the closest cluster centroid.
 - 3.2 Compute the centroids for the clusters by taking the average of the all data points that belong to each cluster.



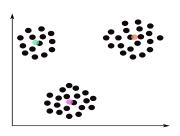
Dataset $D = \{x_i | x_i \in R^n\}$, a distance metric d

- 1. Specify number of clusters K
- 2. Initialize K data points as cluster centroids
- 3. Keep iterating until there is no change to the centroids
 - 3.1 Assign each data point to the closest cluster centroid.
 - 3.2 Compute the centroids for the clusters by taking the average of the all data points that belong to each cluster.



Dataset $D = \{x_i | x_i \in R^d\}$, a distance metric d

- 1. Specify number of clusters K
- 2. Initialize K data points as cluster centroids
- 3. Keep iterating until there is no change to the centroids
 - 3.1 Assign each data point to the closest cluster centroid.
 - 3.2 Compute the centroids for the clusters by taking the average of the all data points that belong to each cluster.



- Total distances:

$$\begin{split} T &= \sum_{k=1}^K \sum_{i \in C_k} (\sum_{j \in C_k} d_{ij} + \sum_{j \neq C_k} d_{ij}) \\ &= \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \in C_k} d_{ij} + \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \neq C_k} d_{ij} \\ &= T_w + T_b \end{split}$$

- Total distances:

$$\begin{split} T &= \sum_{k=1}^K \sum_{i \in C_k} (\sum_{j \in C_k} d_{ij} + \sum_{j \neq C_k} d_{ij}) \\ &= \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \in C_k} d_{ij} + \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \neq C_k} d_{ij} \\ &= T_w + T_b \end{split}$$

$$\mathrm{Normalized} \ T_w = \sum_{k=1}^K \frac{1}{n_k} \sum_{i,j \in C_k} d_{ij} = \sum_{k=1}^K \frac{1}{n_k} \sum_{i,j \in C_k} \|x_i - x_j\|^2 = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2$$

- Total distances:

$$\begin{split} T &= \sum_{k=1}^{K} \sum_{i \in C_k} (\sum_{j \in C_k} d_{ij} + \sum_{j \neq C_k} d_{ij}) \\ &= \sum_{k=1}^{K} \sum_{i \in C_k} \sum_{j \in C_k} d_{ij} + \sum_{k=1}^{K} \sum_{i \in C_k} \sum_{j \neq C_k} d_{ij} \\ &= T_w + T_b \end{split}$$

$$\text{Normalized } T_w = \sum_{k=1}^K \frac{1}{n_k} \sum_{i,j \in C_k} d_{ij} = \sum_{k=1}^K \frac{1}{n_k} \sum_{i,j \in C_k} \|x_i - x_j\|^2 = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2$$

- Idea:

$$(T_w \to min; T_b \to max) \Longleftrightarrow T_w \to min$$

- Total distances:

$$\begin{split} T &= \sum_{k=1}^K \sum_{i \in C_k} (\sum_{j \in C_k} d_{ij} + \sum_{j \neq C_k} d_{ij}) \\ &= \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \in C_k} d_{ij} + \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \neq C_k} d_{ij} \\ &= T_w + T_b \end{split}$$

$$\text{Normalized } T_w = \sum_{k=1}^K \frac{1}{n_k} \sum_{i,j \in C_k} d_{ij} = \sum_{k=1}^K \frac{1}{n_k} \sum_{i,j \in C_k} \|x_i - x_j\|^2 = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2$$

- Idea:

$$(T_w \to min; T_b \to max) \Longleftrightarrow T_w \to min$$

- Objective function:

$$\underset{\mu,C}{\operatorname{argmin}} \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2$$

- Total distances:

$$\begin{split} T &= \sum_{k=1}^K \sum_{i \in C_k} (\sum_{j \in C_k} d_{ij} + \sum_{j \neq C_k} d_{ij}) \\ &= \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \in C_k} d_{ij} + \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \neq C_k} d_{ij} \\ &= T_w + T_b \end{split}$$

$$\text{Normalized } T_w = \sum_{k=1}^K \frac{1}{n_k} \sum_{i,j \in C_k} d_{ij} = \sum_{k=1}^K \frac{1}{n_k} \sum_{i,j \in C_k} \|x_i - x_j\|^2 = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2$$

- Idea:

$$(T_w \to min; T_b \to max) \Longleftrightarrow T_w \to min$$

- Objective function:

$$\underset{\mu,C}{\operatorname{argmin}} \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2$$

Note: K-means is guaranteed to converge, but not same result and speed due to different initialization.

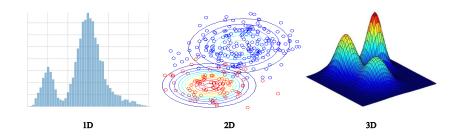
Gaussain mixture models (GMMs)

https://www.naftaliharris.com/blog/visualizing-k-means-clustering/

Initialization

- Classical clustering methods:
 - Work well with low-dimensional data
 - Poorly perform with high-dimensional data due to the sparsity (curse of dimensionality)
- Solution with high dimensional data:
 - ▶ Dimension reduction + clustering
 - ► Autoencoder + clustering (sequential or simultaneous manner)
 - ▶ Autoencoder is key for most deep clustering methods

Gaussain mixture models (GMMs)



Gaussain mixture models (GMMs)

