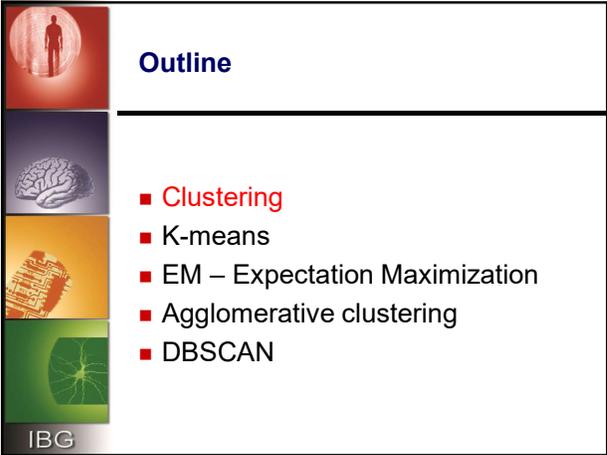


**Signal & Data Analysis in Neuroscience
2020
Clustering**

Izhar Bar-Gad
Room: 408 Phone: 7141 Email: izhar.bar-gad@biu.ac.il

The slide features a grid of four images: a red silhouette of a person, a green brain with neural connections, a blue brain, and a yellow circuit board.



Outline

- Clustering
- K-means
- EM – Expectation Maximization
- Agglomerative clustering
- DBSCAN

IBG

The slide includes a vertical sidebar with four icons: a person silhouette, a brain, a circuit board, and a neural network.



Acknowledgements & Links

- Many slides adapted from presentations by Dana Cohen (BIU) and David Sontag (NYU)
- Links
http://home.deib.polimi.it/matteucc/Clustering/tutorial_html/index.html

IBG

The slide includes a vertical sidebar with four icons: a person silhouette, a brain, a circuit board, and a neural network.






Clustering

Clustering of **data** is a method by which large sets of data is **grouped** into clusters of smaller sets of similar data.

Example: ●●●●●●●●●●

The balls of same color are clustered into a group as shown below : ●●●● ●●●● ●●●●

Thus, clustering means grouping of data or dividing a large data set into smaller data sets of some similarity.

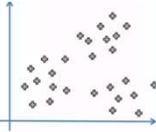
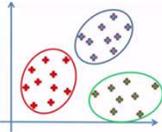
IBG

4






Clustering


➔ Clustering ➔


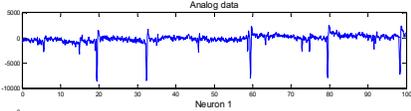
IBG

5

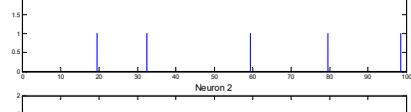




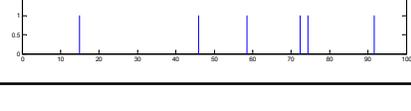

Example: Spike sorting I



Neuron 1



Neuron 2



IBG

6

Example: Spike sorting II

The slide displays a 2x2 grid of scatter plots showing the projection of neural data onto the 1st and 2nd principal components (PCs). The top-left plot is labeled 'Projection on 1st PC' and the bottom-left plot is labeled 'Projection on 2nd PC'. Two clusters are highlighted in the 2D projection: a green cluster and a pink cluster. To the right of each scatter plot is a corresponding amplitude waveform (Amplitude in μV vs. Offset in ms) for that cluster. The green waveform shows a sharp peak, while the pink waveform shows a broader peak. The bottom-left plot also shows a waveform for a single neuron. The IBG logo is in the bottom left corner.

Classification vs. Clustering

- **Classification**
 - Supervised learning
 - Partition examples into groups according to pre-defined categories/classes
 - Requires labeled data for training
- **Clustering**
 - Unsupervised learning
 - Partition examples into groups when no pre-defined categories/classes are available
 - Only instances required, but no labels

The IBG logo is in the bottom left corner.

A "good" clustering algorithm

- **Internal criterion:**
 - The intra-cluster similarity is high
 - The inter-cluster similarity is low

depends on representation & similarity measure
- **External criterion:**

The quality of a clustering is also measured by its ability to discover hidden patterns or latent classes or in comparison to "gold standard"

The IBG logo is in the bottom left corner.



How hard is clustering?



- Brute force: consider all possible clusters, and pick the one that has best inter and intra cluster distance properties.
- Given n points tested in in k clusters, the number of different configurations is:
$$\frac{k^n}{k!}$$
- Enumerating using brute force is too hard leading to iterative optimization algorithms

IBG



Clustering methods



- **Hierarchical**
 - Agglomerative (bottom-up)
 - Divisive (top-down)
- **Partitioning**
 - K-means
 - Mixture of Gaussians

IBG



Outline



- Clustering
- **K-means**
- EM – Expectation Maximization
- Agglomerative clustering
- DBSCAN

IBG



K-means

An iterative clustering algorithm:

- **Initialize:** Pick K random cluster centers
- **Alternate:**
 1. Assign data points to closest cluster center
 2. Change the cluster center to the average of its assigned points
- **Stop:** No assignments change occurs



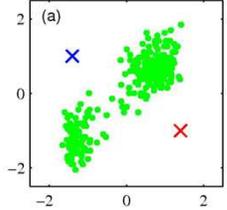
IBG

13



K-means – example 1

- **Initialize:** Pick 2 random cluster centers



(a)



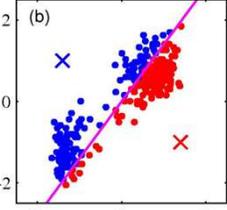
IBG

14



K-means – example 2

- **Alternate 1:** Assign all data points to closest center



(b)



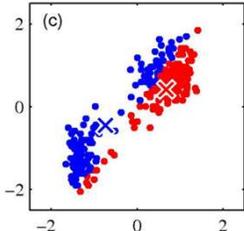
IBG

15



K-means – example 3

- Alternate 2: Change cluster center to the average of the assigned points

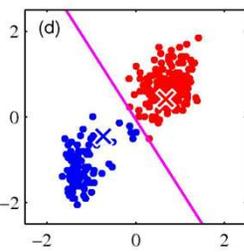


IBG

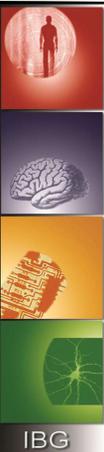


K-means – example 4

- Repeat alternate stages.



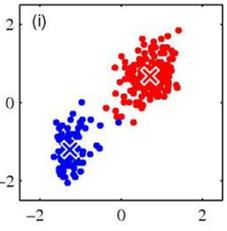
IBG



K-means – example 5

- Finish: Convergence.

K means is guaranteed to converge !



IBG



K-mean – formal description

- Start with random guess of where the K cluster centers $\mathbf{m}_1 \dots \mathbf{m}_K$ are
- Repeat the following until cluster centers stop changing:
 - assign each data point to the nearest cluster:

$$p(n, k) = 1 \text{ if data point } \mathbf{x}^{(n)} \text{ is closer to } \mathbf{m}_k \text{ than to any other } \mathbf{m}_{j \neq k}$$
 - move each cluster center to the **mean** of all data points assigned to it:

$$\mathbf{m}_k = \frac{\sum_n p(n, k) \mathbf{x}^{(n)}}{\sum_j p(j, k)}$$

← Vector sum of all data points assigned to cluster k
← Count of all data points assigned to cluster k

$$\mathbf{m}_k = \sum_n w(n, k) \mathbf{x}^{(n)} \text{ where } w(n, k) \triangleq \frac{p(n, k)}{\sum_j p(j, k)}$$




IBG



K-means converges

Objective

$$\min_{\mu} \min_C \sum_{i=1}^k \sum_{x \in C_i} |x - \mu_i|^2$$

1. Fix μ , optimize C :

Step 1 of kmeans

$$\min_C \sum_{i=1}^k \sum_{x \in C_i} |x - \mu_i|^2 = \min_C \sum_i^n |x_i - \mu_{x_i}|^2$$
2. Fix C , optimize μ :

$$\min_{\mu} \sum_{i=1}^k \sum_{x \in C_i} |x - \mu_i|^2$$
 - Take partial derivative of μ_i and set to zero, we have
$$\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$$

Step 2 of kmeans

Kmeans takes an alternating optimization approach, each step is guaranteed to decrease the objective – thus guaranteed to converge

[Slide from Alan Fern]



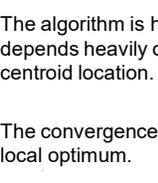
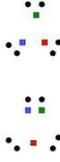


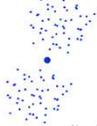
IBG



K-mean – heuristic algorithm

- The algorithm is heuristic and depends heavily on initial centroid location.
- The convergence may occur to a local optimum.

Would be better to have one cluster here

... and two clusters here





IBG



Distances

The results of the clustering depend on the distance

- Euclidean distance (L_2 norm):

$$L_2(\vec{x}, \vec{x}') = \sum_{i=1}^m (x_i - x'_i)^2$$

- L1 norm distance:

$$L_1(\vec{x}, \vec{x}') = \sum_{i=1}^m |x_i - x'_i|$$




IBG



K-means

Strength of the K-means:

- *Efficient: $O(tkn)$, n objects, k clusters, and t iterations. $k, t \ll n$.*

Weakness of the k-means:

- *Requires mean, what about categorical data?*
- *Need to specify k , in advance.*
- *Unable to handle noisy data and outliers*
- *Not suitable for clusters of non-convex shapes.*

In an Euclidian world a convex shape is a set of points containing all line segments between each pair of its points





IBG



Soft (fuzzy) K-mean

- Clustering typically assumes that each instance is given a "hard" assignment to exactly one cluster.
- Does not allow uncertainty in class membership or for an instance to belong to more than one cluster.
- *Soft clustering* gives probabilities that an instance belongs to each of a set of clusters.
- Each instance is assigned a probability distribution across a set of discovered categories (probabilities of all categories must sum to 1).





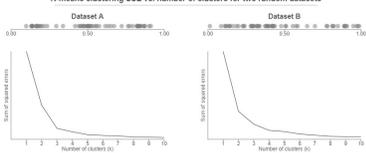
IBG



Number of clusters

- The number of clusters (K) is unknown and typically cannot be set unambiguously.
- A simple method is the “elbow” method which locates the point of change in the additional reduction of “within cluster sum of squares”.

K-means clustering SSE vs. number of clusters for two random datasets






Number of clusters

The Silhouette method - how similar a point is to its own cluster (cohesion) compared to other clusters (separation)

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))} \text{ if } |C_i| > 1$$

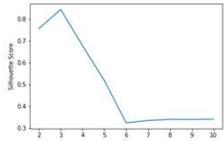
and

$$s(i) = 0 \text{ if } |C_i| = 1$$

For each data point $i \in C_i$ (data point i in the cluster C_i), let

$$a(i) = \frac{1}{|C_i| - 1} \sum_{j \in C_i, j \neq i} d(i, j)$$

For each data point $i \in C_i$, we now define

$$b(i) = \min_{j \neq i} \frac{1}{|C_j|} \sum_{j \in C_j} d(i, j)$$



(wikipedia.org)



Outline

- Clustering
- K-means
- EM – Expectation Maximization
- Agglomerative clustering
- DBSCAN








A better algorithm: Mixture-of-Gaussians clustering

When the data vectors are clustered, it is more appropriate to fit a distribution with multiple peaks. Consider the mixture-of-Gaussians distribution:

$$p(\mathbf{x}; \mathbf{m}_1, \Sigma_1, \dots, \mathbf{m}_K, \Sigma_K) = \frac{1}{K} \sum_{k=1}^K g_k(\mathbf{x}; \mathbf{m}_k, \Sigma_k)$$

↑ mixture distribution ↑ Gaussian distributions, with means and covariances \mathbf{m}_k, Σ_k

How do we fit such a distribution to a set of data vectors $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$? If we knew which Gaussian is "responsible" for each data vector, we could compute the mean and covariance separately for each Gaussian – from the vectors it is responsible for. This suggests the following iterative algorithm (the EM algorithm):

Iterate the following two steps until convergence:

- Expectation (E-step):** Compute $P(x_i | E(g))$ for each example given the current model, and probabilistically re-label the examples based on these posterior probability estimates.
- Maximization (M-step):** Re-estimate the model parameters from the probabilistically re-labeled data.

IBG






Expectation Maximization

1. Compute the probability $p(n,k)$ that data point n came from Gaussian k , and the normalized weights $w(n,k)$ which sum to 1 for each Gaussian:

$$p(n,k) = \frac{g_k(\mathbf{x}^{(n)})}{\sum_j g_j(\mathbf{x}^{(n)})} \quad w(n,k) = \frac{p(n,k)}{\sum_j p(j,k)}$$
2. Re-compute the mean and covariance of all data points that Gaussian k is responsible for, using $w(n,k)$ as weights:

$$\mathbf{m}_k = \sum_n w(n,k) \mathbf{x}^{(n)} \quad \Sigma_k = \sum_n w(n,k) (\mathbf{x}^{(n)} - \mathbf{m}_k)(\mathbf{x}^{(n)} - \mathbf{m}_k)^T$$
3. Repeat until \mathbf{m}, Σ no longer change.

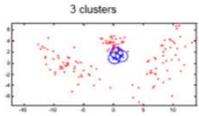
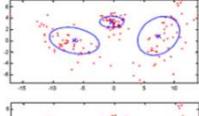
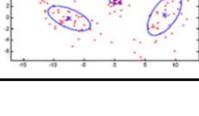
IBG



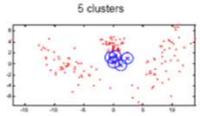
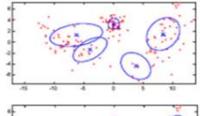
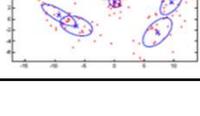



Mixture of gaussian - example

3 clusters

5 clusters

IBG

Comparison of the two algorithms

In both cases, we compute a quantity $p(n,k)$ that tells us how well data point n fits in cluster k . Then we compute the normalized weights

$$w(n,k) = p(n,k) / \sum_j p(j,k)$$

and re-compute the cluster centers according to weighted center-of-mass calculation

$$m_k = \sum_n w(n,k) \mathbf{x}^{(n)}$$

There are two differences:

1. In K-means the "fit" p is either 1 or 0, depending on which is the nearest cluster; In MOG, the values of p vary continuously between 0 and 1, and correspond to probabilities
2. In MOG we also re-compute the covariance matrix, which in turn affects how we determine the fit of data points to clusters; In K-means, the fit is always computed in the same way, corresponding to the assumption of circular clusters

Mixture-of-Gaussians vs. K-means clustering

The results are similar when the clusters are well-separated and roughly circular

IBG

Mixture-of-Gaussians vs. K-means clustering

The results are similar when the clusters are well-separated and roughly circular

But for more complex problems K-means can be fooled more easily

IBG



Outline

- Clustering
- K-means
- EM – Expectation Maximization
- **Agglomerative clustering**
- DBSCAN





IBG



Agglomerative clustering

- Bottom-up approach to hierarchical clustering
- Neighboring points generate a common cluster iteratively building larger clusters.
- Initialization: Each instance in its own cluster.
- Repeat:
 - Pick the two closest clusters
 - Merge them into a new cluster
- Stop: there is only one cluster
- Produces a family of clusterings represented by a dendrogram

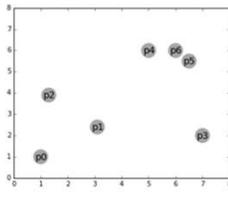
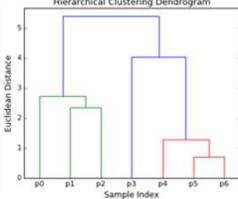




IBG

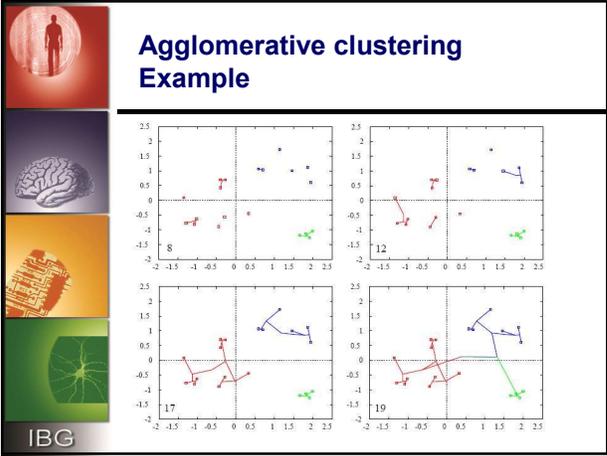


Agglomerative clustering Example

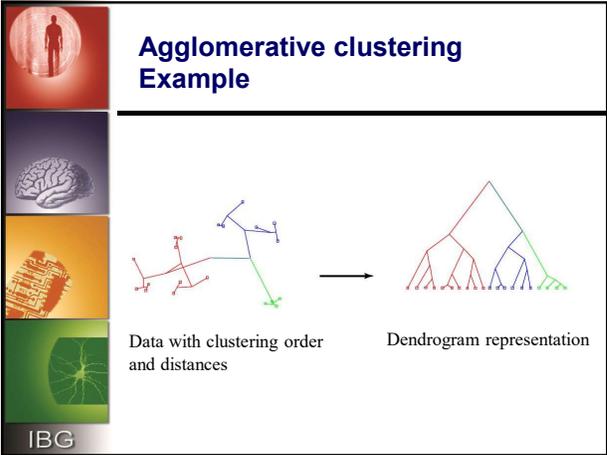



IBG

(kdnuggets.com)



37



38

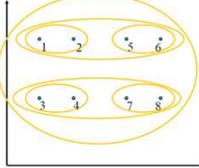
-
- How should we define “closest” for clusters with multiple elements?
 - Closest pair (single-link clustering)
 - Furthest pair (complete-link clustering)
 - Average of all pairs
 - Different choices create different clustering behaviors

39

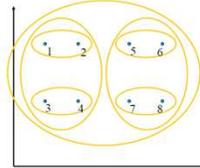


Agglomerative clustering Distances

Closest pair
(single-link clustering)



Farthest pair
(complete-link clustering)



[Pictures from Thorsten Joachims]



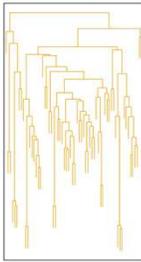


IBG

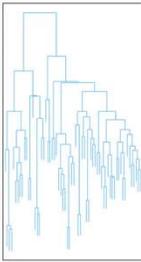


Agglomerative clustering Distances

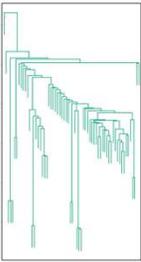
Average



Farthest



Nearest



Mouse tumor data from [Hastie et al.]





IBG



Agglomerative clustering Questions

- Will agglomerative clustering converge?
 - To a global optimum?
- Will it always find the true patterns in the data?
- How many clusters to pick?





IBG

Correlation matrix

Hierarchical

Raw

K-mean
K=?

IBG

Spellman, et al. 1998.
Gene expression data

Outline

- Clustering
- K-means
- EM – Expectation Maximization
- Agglomerative clustering
- **DBSCAN**

IBG

DBSCAN

- Density-based spatial clustering of applications with noise (DBSCAN) – Ester et al. 1996
- Discovers clusters of arbitrary shape
 - Group together points in high-density
 - Mark as outliers the points that lie alone in low-density regions
- Two main parameters:
 - ϵ (**epsilon**) - radius of the *neighborhood region*
 - **minPoints** - the minimum number of points that should be contained within that neighborhood.

IBG



- Epsilon neighborhood (N_ϵ) : set of all points within a distance ' ϵ '.
- Core point : A point that has at least 'minPoint' (including itself) points within its N_ϵ .
- Direct Density Reachable (DDR) : A point q is directly density reachable from a point p if p is core point and $q \in N_\epsilon$.
- Density Reachable (DR) : Two points are DR if there is a chain of DDR points that link these two points.
- Border Point: Point that are DDR but not a core point.
- Noise : Points that do not belong to any point's N_ϵ .

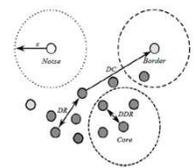
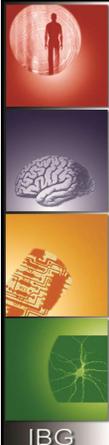
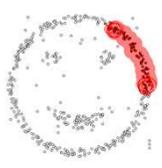


Figure 1: DBSCAN clustering with $\text{minPoints} = 4$



DBSCAN - algorithm

- Arbitrarily pick a point in the dataset which is not marked.
 - If there are at least **minPoint** points within a radius of ϵ to the point then we consider all these points to be part of the same cluster.
 - Expand cluster by breadth first repeating the neighborhood calculation for each neighboring point

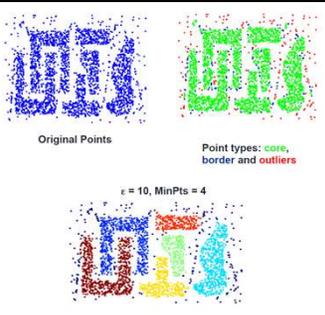


epsilon = 100
minPoints = 4

(kdnuggets.com)



DBSCAN



Original Points

Point types: core, border and outliers

epsilon = 10, MinPts = 4

Clusters



DBSCAN vs. Gaussian mixtures

- Advantages
 - Arbitrary distributions
 - Deals well with noise and outliers
- Disadvantages
 - Varying densities
 - Parameter dependence

DBSCAN	GaussianMixture
	
	
	
	
	
	

IBG (towardsdatascience.com)