# DATA MINING 2
# Advanced Clustering

Riccardo Guidotti

a.a. 2020/2021

K-Means Extensions
# Bisecting K-Means

# Bisecting K-means

- Variant of K-Means that can produce a hierarchical clustering
- The number of clusters K must be specified.
- Start with a unique cluster containing all the points.

1: Initialize the list of clusters to contain the cluster containing all points.
2: **repeat**
3:    Select the cluster with the highest SSE to the list of clusters
4:    **for** $i = 1$ to *number_of_iterations* **do**
5:       Bisect the selected cluster using basic 2-Means
6:    **end for**
7:    Add the two clusters from the bisection to the list of clusters.
8: **until** Until the list of clusters contains $K$ clusters
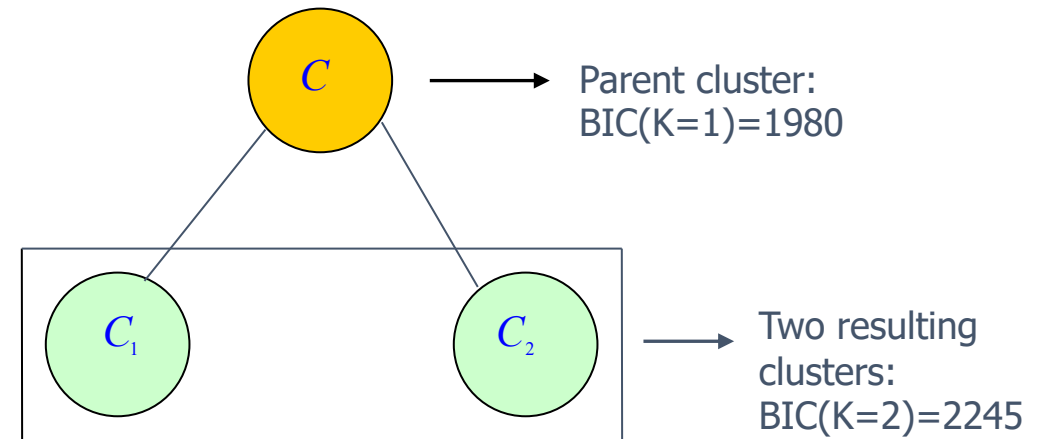
# Bisecting K-means Limitations

- The algorithm can be also exhaustive and terminating at a singleton clusters if K is not specified.

- Terminating at singleton clusters
  - Is time consuming
  - Singleton clusters are meaningless (i.e., over-splitting)
  - Intermediate clusters are more likely to correspond to real classes

- Bisecting K-Means do not uses any criterion for stopping bisections before singleton clusters are reached.

K-Means Extensions

# X-Means

# Bayesian Information Criterion (BIC)

- A strategy to stop the Bisecting algorithm when meaningful clusters are reached to avoid over-splitting.

- The **BIC** can be adopted as **splitting criterion** of a cluster in order to decide whether a cluster should split or no.

- BIC **measures the improvement** of the cluster structure between a cluster and its two children clusters.

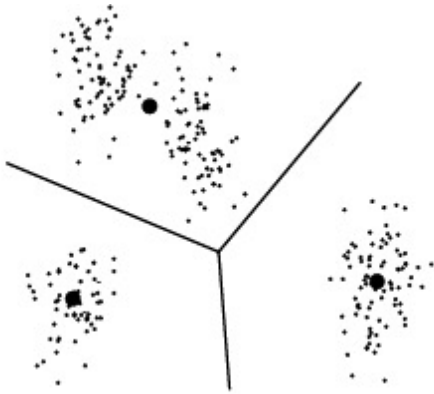- If the BIC of the parent is less than BIC of the children than we accept the bisection.

$C$

Parent cluster:
BIC(K=1)=1980

$C_1$      $C_2$

Two resulting clusters:
BIC(K=2)=2245

# X-Means

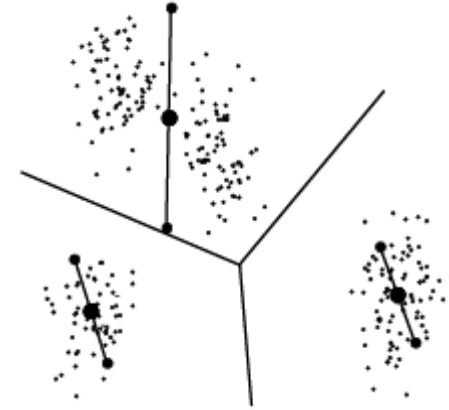For $k$ in a given range $[r_1, r_{max}]$:

1. Improve Params: run K-Means with with the current $k$.

2. Improve Structure: recursively split each cluster in two (Bisecting 2-Means) and use *local BIC* to decide to keep the split. Stop if the current structure does not respect *local BIC* or the number of clusters is higher than $r_{max}$.

3. Store the actual configuration with a global BIC calculated on the whole configuration

4. If $k > r_{max}$ stop and return the best model w.r.t. the *global BIC*.
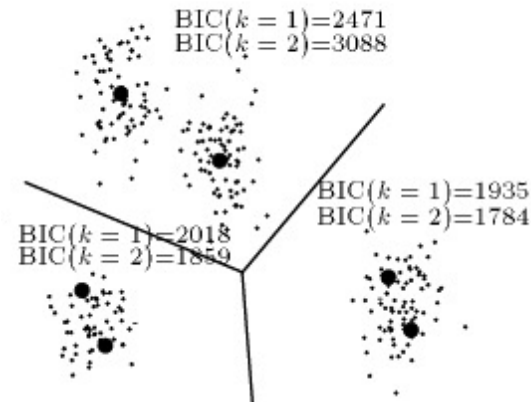
# X-Means

1. K-means with k=3

2. Split each centroid in 2 children moved a distance proportional to the region size in opposite direction (random)
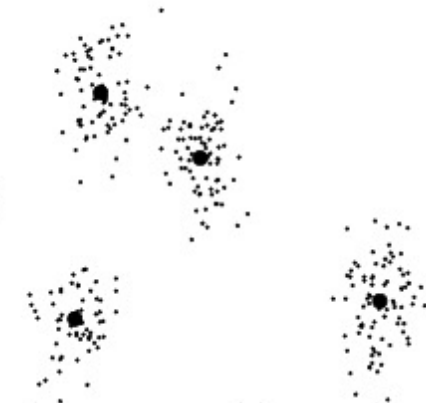
3. Run 2-means in each region locally

4. Compare BIC of parent and children

4. Only centroids with higher BIC survives

$\text{BIC}(k=1)=2471$
$\text{BIC}(k=2)=3088$

$\text{BIC}(k=1)=1935$
$\text{BIC}(k=2)=1784$

$\text{BIC}(k=1)=2018$
$\text{BIC}(k=2)=1859$

# BIC Formula in X-Means

- The BIC score of a data collection is defined as (Kass and Wasserman, 1995):

$$BIC(M_j) = \hat{l}_j(D) - \frac{p_j}{2} \log R$$

- $\hat{l}_j(D)$ is the log-likelihood of the dataset D

- $p_j$ is a function of the number of independent parameters: centroids coordinates, variance estimation.

- $R$ is the number of points of a cluster, M is the number of dimensions

- Approximate the probability that the clustering in $M_j$ is describing the real clusters in the data

# BIC Formula in X-Means

- Adjusted Log-likelihood of the model.

- **The likelihood that the data is "explained by" the clusters** according to the spherical-Gaussian assumption of K-Means

$$BIC(M_j) = \hat{l}_j(D) - \frac{p_j}{2} \log R$$

- Focusing on the set $D_n$ of points which belong to centroid $n$

$$\hat{l}(D_n) = -\frac{R_n}{2} \log(2\pi) - \frac{R_n \cdot M}{2} \log(\hat{\sigma}^2) - \frac{R_n - K}{2}$$
$$+ R_n \log R_n - R_n \log R$$

- It estimates how closely to the centroid are the points of the cluster.

K-Means Origins
# Expectation Maximization

# Model-based Clustering (probabilistic)

- In order to understand our data, we will assume that there is a generative process (a model) that creates/describes the data, and we will try to find the model that **best fits** the data.
  - Models of different complexity can be defined, but we will assume that our model is a distribution from which data points are sampled
  - **Example**: the data is the height of all people in **Greece**

- In most cases, a single distribution is not good enough to describe all data points: **different parts of the data follow a different distribution**
  - **Example**: the data is the height of all people in Greece and China
  - We need a mixture model
  - Different distributions correspond to different clusters in the data.

# Expectation Maximization Algorithm

- Initialize the values of the parameters in $\Theta$ to some random values

- Repeat until convergence
    - **E-Step**: Given the parameters $\Theta$ **estimate** the membership probabilities $\mathrm{P}\left(G_j\middle|x_i\right)$
    - **M-Step:** Given the probabilities $\mathrm{P}\left(G_j\middle|x_i\right)$, calculate the parameter values $\Theta$ that (in expectation) **maximize** the data likelihood

- **Examples**
    - **E-Step**: Assignment of points to clusters
        - K-Means: hard assignment, EM: soft assignment
    - **M-Step**: Parameters estimation
        - K-Means: Computation of centroids, EM: Computation of the new model parameters

# EM in K-Means

centroids

- Initialize the values of the parameters in Θ to some random values (randomly select the centroids)

- Repeat until convergence
  - **E-Step**: Given the parameters Θ (given the centroids) **estimate** the membership probabilities $\mathrm{P}(G_j|x_i)$ (assign points to clusters based on distances with the centroids)
  - **M-Step:** Given the probabilities $\mathrm{P}(G_j|x_i)$ (given the membership of points to clusters, i.e., 100% probability of belonging to a cluster) calculate the parameter values Θ that (in expectation) **maximize** the data likelihood (calculate the new centroids as mean values, i.e., those that minimize the distances with the other points in the cluster)

# Expectation Maximization Algorithm

---

**Algorithm 9.2** EM algorithm.

---

1: Select an initial set of model parameters.
   (As with K-means, this can be done randomly or in a variety of ways.)
2: **repeat**
3:   **Expectation Step** For each object, calculate the probability that each object belongs to each distribution, i.e., calculate $prob(distribution\ j|\mathbf{x}_i, \Theta)$.
4:   **Maximization Step** Given the probabilities from the expectation step, find the new estimates of the parameters that maximize the expected likelihood.
5: **until** The parameters do not change.
   (Alternatively, stop if the change in the parameters is below a specified threshold.)

---

K-Means Brother

# Mixture Gaussian Model

# Gaussian Distribution

- Example: the data is the height of all people in Greece
- Experience has shown that this data follows a Gaussian (Normal) distribution

$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- $\mu$ = mean, $\sigma$ = standard deviation

# Mixture Gaussian Model

- What is a model?
  - A Gaussian distribution is defined by the mean $\mu$ and the standard deviation $\sigma$
  - We define our model as the pair of parameters $\theta = (\mu, \sigma)$

- More generally, a model is defined as a vector of parameters $\theta$

- We want to find the normal distribution $N(\mu, \sigma)$ that best fits our data
  - Find the best values for $\mu$ and $\sigma$
  - But what does "best fit" mean?

# Maximum Likelihood Estimation (MLE)

- Suppose that we have a vector $X = \{x_1, \dots, x_n\}$ of values
- We want to fit a Gaussian model $N(\mu, \sigma)$ to the data
- Probability of observing a point $x_i$

$$P(x_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

- Probability of observing all points (we assume independence)

$$P(X) = \prod_{i=1}^{n} P(x_i) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

- We want to find the parameters $\theta = (\mu, \sigma)$ that maximizes the probability $P(X|\theta)$

# Maximum Likelihood Estimation (MLE)

- The probability $P(X|\theta)$ as a function of $\theta$ is the **Likelihood** function

$$L(\theta) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

- It is usually easier to work with the **Log-Likelihood** function

$$LL(\theta) = -\sum_{i=1}^{n} \frac{(x_i-\mu)^2}{2\sigma^2} - \frac{1}{2}n\log 2\pi - n\log\sigma$$

- Thus, the Maximum Likelihood Estimation for the Gaussian Model consists in finding the parameters $\mu, \sigma$ that maximize $LL(\theta)$

$$\mu = \frac{1}{n}\sum_{i=1}^{n} x_i = \mu_X$$

Sample Mean

$$\sigma^2 = \frac{1}{n}\sum_{i=1}^{n}(x_i-\mu)^2 = \sigma_X^2$$

Sample Variance

# Maximum Likelihood Estimation (MLE)

- Note: these are also the most likely parameters given the data.

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$$

- If we have no prior information about $\theta$, or $X$, then maximizing $P(\theta|X)$ is the same as maximizing $P(X|\theta)$.

# Mixture of Gaussians

• Suppose that you have the heights of people from Greece and China and the distribution looks like the figure below (dramatization)



(a) Probability density function for the mixture model.

(b) 20,000 points generated from the mixture model.

**Figure 9.2.** Mixture model consisting of two normal distributions with means of -4 and 4, respectively. Both distributions have a standard deviation of 2.

# Mixture of Gaussians

- In this case the data is the result of the **mixture** of two Gaussians
    - One for Greek people, and one for Chinese people
    - Identifying for each value which Gaussian is most likely to have generated it will give us a clustering.



(a) Probability density function for the mixture model.

(b) 20,000 points generated from the mixture model.

**Figure 9.2.** Mixture model consisting of two normal distributions with means of -4 and 4, respectively. Both distributions have a standard deviation of 2.

# Mixture Model

- A value $x_i$ is generated according to the following process:
  - First select the nationality
    - With probability $\pi_G$ select Greek, with probability $\pi_C$ select China ($\pi_G + \pi_C = 1$)

  - Given the nationality, generate the point from the corresponding Gaussian
    - $P(x_i|\theta_G) \sim N(\mu_G, \sigma_G)$ if Greece
    - $P(x_i|\theta_C) \sim N(\mu_C, \sigma_C)$ if China

# Mixture Model

Assign a point to a cluster. In K-Means they are the membership: hard assignment.

Describe a cluster. In K-Means they are the centroids.

- Our model has the following parameters

$$\Theta = (\pi_G, \pi_C, \mu_G, \mu_C, \sigma_G, \sigma_C)$$

Mixture probabilities    Distribution Parameters

- For value $x_i$, we have:

$$P(x_i|\Theta) = \pi_G P(x_i|\theta_G) + \pi_C P(x_i|\theta_C)$$

- For all values $X = \{x_1, \dots, x_n\}$

$$P(X|\Theta) = \prod_{i=1}^{n} P(x_i|\Theta)$$

- We want to estimate the parameters that **maximize** the Likelihood

# Mixture Model

- Once we have the parameters $\theta = (\pi_G, \pi_C, \mu_G, \sigma_G, \mu_C, \sigma_C)$, we can **estimate** the **membership probabilities** $P(G|x_i)$ and $P(C|x_i)$ for each point $x_i$:

- This is the probability that point $x_i$ belongs to the Greek or the Chinese population (cluster)

$$
\begin{aligned}
P(G|x_i) &= \frac{P(x_i|G)P(G)}{P(x_i|G)P(G) + P(x_i|C)P(C)} \\
&= \frac{P(x_i|G)\pi_G}{P(x_i|G)\pi_G + P(x_i|C)\pi_C}
\end{aligned}
$$

# Mixture of Gaussians as EM

- Initialize the values of the parameters in $\theta$ to some random values

- Repeat until convergence
  - **E-Step**: Given the parameters $\Theta$ **estimate** the membership probabilities $P(G|x_i)$ and $P(C|x_i)$.
  - **M-Step:** Calculate the parameter values $\Theta$ that (in expectation) **maximize** the data likelihood.

$$\pi_G = \frac{1}{n} \sum_{i=1}^{n} P(G|x_i) \qquad \pi_C = \frac{1}{n} \sum_{i=1}^{n} P(C|x_i)$$

Fraction of population in G,C

$$\mu_C = \sum_{i=1}^{n} \frac{P(C|x_i)}{n * \pi_C} x_i \qquad \mu_G = \sum_{i=1}^{n} \frac{P(G|x_i)}{n * \pi_G} x_i$$

MLE Estimates if $\pi$'s were fixed

$$\sigma_C^2 = \sum_{i=1}^{n} \frac{P(C|x_i)}{n * \pi_C} (x_i - \mu_C)^2 \qquad \sigma_G^2 = \sum_{i=1}^{n} \frac{P(G|x_i)}{n * \pi_G} (x_i - \mu_G)^2$$

DBSCAN Evolution
# OPTICS

# When DBSCAN Works Well



Original Points

Clusters

- **Resistant to Noise**

- **Can handle clusters of different shapes and sizes**

# When DBSCAN Does NOT Work Well



**Original Points**



(MinPts=4, Eps=9.75).



(MinPts=4, Eps=9.92)

- **Varying densities**

- **High-dimensional data**

# OPTICS

- OPTICS: Ordering Points To Identify the Clustering Structure
  - Produces a special order of the dataset wrt its density-based clustering structure.
  - This cluster-ordering contains info equivalent to the density-based clusterings corresponding to a broad range of parameter settings.
  - Good for both automatic and interactive cluster analysis, including finding intrinsic clustering structure.
  - Can be represented graphically or using visualization techniques.

# OPTICS: Extension from DBSCAN

- OPTICS requires two **parameters**:
  - $\varepsilon$, which describes the maximum distance (radius) to consider,
  - MinPts, describing the number of points required to form a cluster

- **Core point**. A point $p$ is a core point if at least MinPts points are found within its $\varepsilon$-neighborhood.

- **Core Distance**. It is the **minimum** value of radius required to classify a given point as a core point. If the given point is not a Core point, then it's Core Distance is undefined.

Eps = 6mm

MinPts = 5

Core_Distance(p) = 3mm

# OPTICS: Extension from DBSCAN

- **Reachability Distance**. The reachability distance between a point $p$ and $q$ is the **maximum** of the Core Distance of $p$ and the Distance between p and q.

- The Reachability Distance is not defined if $q$ is not a Core point. Below is the example of the Reachability Distance.

- In other words, if $q$ is within the core distance of $p$ then use the core distance, otherwise the real distance.



Eps = 6mm

MinPts = 5

Core_Distance(p) = 3mm

Reachability_Distance(q,p) = 7mm

Reachability_Distance(r,p) = 3mm

# OPTICS Pseudo-Code

- For each point *p* in the dataset
  - Initialize the reachability distance of *p* as undefined
- For each unprocessed point *p* in the dataset
  - Get the neighbors *N* of p
  - Mark *p* as processed and output to the *ordered list*
  - If *p* is a core point
    - Initialize a priority queue *Q* to get the closest point to *p* in terms of reachability
    - Call the function *update(N, p, Q)*
    - For each point *q* in *Q*
      - Get the neighbors *N'* of *q*
      - Mark q as processed and output to the *ordered list*
        - If *q* is a core point Call the function update(N', q, Q)

# OPTICS Pseudo-Code

- Function *update(N, p Q)*
  - Calculate the core distance for *p*
  - For each neighbor *q in N* (update the reachability)
    - If *q* is not processed
      - *new_rd* = reachability distance between *p* and *q*
    - If *q is not in Q*
      - *Q.insert(q, new_rd)*
    - Else
      - *If new_rd < q.rd*
        - *Q.move_up(q, new_rd)*

# OPTICS Output

- OPTICS outputs the points in a particular ordering, annotated with their smallest reachability distance.

- A reachability-plot (a special kind of dendrogram), the hierarchical structure of the clusters can be obtained easily.

- x-axis: the ordering of the points as processed by OPTICS

- y-axis: the reachability distance

- Points belonging to a cluster have a low reachability distance to their nearest neighbor, the clusters show up as valleys in the reachability plot. The deeper the valley, the denser the cluster.

# OPTICS Output

- Clusters are extracted
  1. by selecting a range on the x-axis after visual inspection,
  2. by selecting a threshold on the y-axis
  3. by different algorithms that try to detect the valleys by steepness, knee detection, or local maxima. Clustering obtained this way usually are hierarchical, and cannot be achieved by a single DBSCAN run.



https://scikit-learn.org/stable/auto_examples/cluster/plot_optics.html#sphx-glr-auto-examples-cluster-plot-optics-py

# OPTICS: The Radius Parameter

- Both core-distance and reachability-distance are undefined if no sufficiently dense cluster (w.r.t. ε) is available.

- Given a sufficiently large ε, this never happens, but then every ε-neighborhood query returns the entire database.

- Hence, the ε parameter is required to cut off the density of clusters that are no longer interesting, and to speed up the algorithm.

- The parameter ε is, strictly speaking, not necessary.

- It can simply be set to the maximum possible value.

- When a spatial index is available, however, it does play a practical role with regards to complexity.

- OPTICS abstracts from DBSCAN by removing this parameter, at least to the extent of only having to give the maximum value.

# References

- Advanced Clustering. Chapter 9.2.2. Introduction to Data Mining.

- Pelleg, Dan, and Andrew W. Moore. "X-means: Extending k-means with efficient estimation of the number of clusters." 2000.

- Mihael Ankerst; Markus M. Breunig; Hans-Peter Kriegel; Jörg Sander (1999). OPTICS: Ordering Points To Identify the Clustering Structure.