Clustering based on Loïc Cerf's slides (UFMG)



Marc Plantevit

March 11, 2021 UCBL – LIRIS – DM2L







LIRIS Clustering Approaches

- Partition-based algorithms: build several partitions then assess them w.r.t. some criteria.
- Hierarchy-based algorithms: create a hierarchical decomposition of the objects w.r.t. some criteria.
- Density-based algorithms: based on the notions of density and connectivity.

Characteristics

- extensibility
- ability to handle different data types
- prior for parameter settings
- ability to handle noisy data and outliers









Density-based Clustering: DBSCAN

6 Conclusion

Marc Plantevit







- **B** Hierarchical Clustering
- Openation Density-based Clustering: DBSCAN

6 Conclusion

Marc Plantevit

Inductive database vision

Querying a clustering:

$$\{X \in P \mid \mathcal{Q}(X, \mathcal{D})\}$$

where:

• $\mathcal D$ is a set of objects $\mathcal O$ associated with a similarity measure,

•
$$P$$
 is $\{(C_1, \ldots, C_k) \in (2^{\mathcal{O}})^k \mid \begin{cases} \forall i = 1..k, C_i \neq \emptyset \\ \forall j \neq i, C_i \cap C_j \neq \emptyset \end{cases} \}, \cup_{l=1}^k C_l = \mathcal{O} \end{cases}$

• Q is a function to optimize. It quantifies how similar are pairs of objects in a same cluster and how dissimilar are those in two different clusters

LIRIS Inductive database vision

Querying a clustering with k-means:

$$\{X \in P \mid \mathcal{Q}(X, \mathcal{D})\}$$

where:

• $\mathcal D$ is a set of objects $\mathcal O$ associated with a similarity measure,

•
$$P$$
 is $\{(C_1, \ldots, C_k) \in (2^{\mathcal{O}})^k \mid \begin{cases} \forall i = 1..k, C_i \neq \emptyset \\ \forall j \neq i, C_i \cap C_j \neq \emptyset \end{cases} \}$ where $\cup_{l=1}^k C_l = \mathcal{O}$
 $k \in \mathbb{N} \setminus \{0\}$ is fixed,

• Q is the maximization of the sum, over all objects, of the similarities to the centers of the assigned clusters:

$$(C_1,\ldots,C_k)\mapsto \sum_{i=1}^k\sum_{o\in C_i}s(o,\frac{\sum_{o\in C_i}}{|C_i|})$$

LIRIS Inductive database vision

Querying a clustering with k-means:

$$\{X \in P \mid \mathcal{Q}(X, \mathcal{D})\}$$

where:

• $\mathcal D$ is a set of objects $\mathcal O$ associated with a similarity measure,

•
$$P$$
 is $\{(C_1, \ldots, C_k) \in (2^{\mathcal{O}})^k \mid \begin{cases} \forall i = 1..k, C_i \neq \emptyset \\ \forall j \neq i, C_i \cap C_j \neq \emptyset \end{cases}$ where $\bigcup_{l=1}^k C_l = \mathcal{O}$
 $k \in \mathbb{N} \setminus \{0\}$ is fixed,

• Q is the maximization of the sum, over all objects, of the similarities to the centers of the assigned clusters:

$$(C_1,\ldots,C_k)\mapsto \sum_{i=1}^k\sum_{o\in C_i}s(o,\mu_i)$$

1

Exact algorithm

Input: $\mathcal{O}, \mathcal{D}, k \in \mathbb{N} \setminus \{0\}$ **Output:** the clustering of \mathcal{O} maximizing f: the sum, over all objects, of the similarities to the centers of the assigned clusters $\mathcal{C}_{\max} \leftarrow \emptyset$ $f_{max} \leftarrow -\infty$ for all k-clustering C of O do if $f(\mathcal{C}, \mathcal{D}) > f_{\max}$ then $f_{\max} \leftarrow f(\mathcal{C}, \mathcal{D})$ $\mathcal{C}_{\max} \leftarrow \mathcal{C}$ end if end for

 $\mathsf{output}(\mathcal{C}_{\mathsf{max}})$

Number of *k*-clusterings

Question

How many k-clusterings are enumerated?

Marc Plantevit

Number of *k*-clusterings

Question

How many *k*-clusterings are enumerated? The Stirling number of the second kind, i. e., $\frac{1}{k!} \sum_{t=0}^{k} (-1)^t \binom{k}{t} (k-t)^n = O(k^n)$.

LIRIS $\frac{k-\text{means}}{k-\text{means principles}}$

k-means is a greedy iterative approach that always converges to a local maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters.

LIRIS

k-means principles

k-means is a greedy iterative approach that always converges to a local maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters.

An iteration consists in two steps:

E Each object is assigned to the cluster whose center is the most similar (thus defining a clustering);

LIRIS

k-means principles

k-means is a greedy iterative approach that always converges to a local maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters.

An iteration consists in two steps:

- E Each object is assigned to the cluster whose center is the most similar (thus defining a clustering);
- M The center of each cluster is updated to the mean of the objects assigned to it.

LIRIS k-means principles

k-means is a greedy iterative approach that always converges to a local maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters.

An iteration consists in two steps:

- E Each object is assigned to the cluster whose center is the most similar (thus defining a clustering);
- M The center of each cluster is updated to the mean of the objects assigned to it.

Initially, the centers of the clusters are randomly drawn. The procedure stops when, from an iteration to the next one, the centers of the clusters have not changed much (or at all).



Dataset:



Marc Plantevit

Clustering 10 / 52





Marc Plantevit

Clustering 10 / 52









Clustering 10

10 / 52





Marc Plantevit

Clustering 10 / 52





Marc Plantevit

Clustering 10 / 52

LIRIS

3-means with $|\mathcal{A}| = 2$: illustration

3-means clustering of the objects in a two-dimensional space using the Euclidean distance.



² 11 / 52

.

k-means algorithm

Input: $\mathcal{O}, \mathcal{D}, k \in \mathbb{N} \setminus \{0\}$

Output: a clustering of \mathcal{O} *locally* maximizing the sum, over all objects, of the similarities to the centers of the assigned clusters $(\mu_i)_{i=1..k} \leftarrow \operatorname{random}(\mathcal{D})$

repeat

LIRIS

$$(C_i)_{i=1..k} \leftarrow \operatorname{assign_cluster}(\mathcal{O}, \mathcal{D}, (\mu_i)_{i=1..k}) \\ (c, (\mu_i)_{i=1..k}) \leftarrow \operatorname{update_centers}(\mathcal{D}, (C_i)_{i=1..k}, (\mu_i)_{i=1..k})) \\ \operatorname{until} c \\ \operatorname{output}((C_i)_{i=1..k}) \end{cases}$$

Marc Plantevit

LIRIS assign_cluster

Input: $\mathcal{O}, \mathcal{D}, (\mu_i)_{i=1..k} \in (\mathbb{R}^{|\mathcal{A}|})^k$ Output: $(C_i)_{i=1..k}$ the clustering of \mathcal{O} such that $\forall i = 1..k, \forall j \neq i, \forall o \in C_i, s(o, \mu_i) \geq s(o, \mu_j)$ for all $o \in \mathcal{O}$ do $a \leftarrow \arg \max_{i=1..k} s(o, \mu_i)$ $C_a \leftarrow C_a \cup \{o\}$ end for return $((C_i)_{i=1..k})$

Marc Plantevit

Complexity of assign_cluster

Question

Assuming the computation of a similarity is linear in the number of attributes |A|, what is the complexity of **assign_cluster**?

Marc Plantevit

Clustering 14 / 52

Complexity of assign_cluster

Question

Assuming the computation of a similarity is linear in the number of attributes |A|, what is the complexity of **assign_cluster**? $O(k|O \times A|)$.

update_centers

Input: $\mathcal{D}, (C_i)_{i=1..k}$ a clustering of $\mathcal{O}, (\mu_i)_{i=1..k} \in (\mathbb{R}^{|\mathcal{A}|})^k$ **Output:** $c \in \{$ false, true $\}$ indicating whether the convergence is reached, $(\mu'_i)_{i=1..k} \in (\mathbb{R}^{|\mathcal{A}|})^k$ such that $\forall i = 1..k, \mu'_i = \frac{\sum_{o \in C_i} o}{|C_i|}$ $c \leftarrow true$ for $i = 1 \rightarrow k \operatorname{do}_{\mu'_i} \leftarrow \frac{\sum_{o \in C_i} o}{|C_i|}$ if $\mu'_i \neq \mu_i$ then $c \leftarrow false$ end if end for return($c, (\mu'_i)_{i=1, k}$)

Complexity of *k*-means

Question

LIRIS

Assuming the computation of a similarity is linear in the number of attributes |A|, what is the complexity of **assign_cluster**? $O(k|O \times A|)$.

Question

What is the complexity of update_centers?

Complexity of *k*-means

Question

LIRIS

Assuming the computation of a similarity is linear in the number of attributes $|\mathcal{A}|$, what is the complexity of **assign_cluster**? $O(k|\mathcal{O} \times \mathcal{A}|)$.

Question

What is the complexity of **update_centers**? $O(|\mathcal{O} \times \mathcal{A}|)$.

Complexity of *k*-means

Question

LIRIS

Assuming the computation of a similarity is linear in the number of attributes $|\mathcal{A}|$, what is the complexity of **assign_cluster**? $O(k|\mathcal{O} \times \mathcal{A}|)$.

Question

What is the complexity of **update_centers**? $O(|\mathcal{O} \times \mathcal{A}|)$.

Question

What is the complexity of k-means if $t \in \mathbb{N}$ iterations are necessary to converge?

Complexity of *k*-means

Question

LIRIS

Assuming the computation of a similarity is linear in the number of attributes $|\mathcal{A}|$, what is the complexity of **assign_cluster**? $O(k|\mathcal{O} \times \mathcal{A}|)$.

Question

What is the complexity of **update_centers**? $O(|\mathcal{O} \times \mathcal{A}|)$.

Question

What is the complexity of k-means if $t \in \mathbb{N}$ iterations are necessary to converge? $O(tk|\mathcal{O} \times \mathcal{A}|)$.



Worst-case scenarios require $2^{\Omega(|\mathcal{O}|)}$ iterations to converge but a smoothed analysis gives a polynomial complexity.



Worst-case scenarios require $2^{\Omega(|\mathcal{O}|)}$ iterations to converge but a smoothed analysis gives a polynomial complexity.

The low complexity of *k*-means is its greatest advantage.

Marc Plantevit

Clustering 17 / 52



• Convergence towards a *local* maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters;

LIRIS $\frac{k-\text{means}}{\text{Limitations of }k-\text{means}}$

- Convergence towards a *local* maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters;
- Sensitivity to outliers (k-medoids replaces the means by medians);

LIRIS Limitations of k-means

- Convergence towards a *local* maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters;
- Sensitivity to outliers (k-medoids replaces the means by medians);
- Tendency to produce equi-sized clusters;
k-means LIRIS Limitations of *k*-means

- Convergence towards a *local* maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters;
- Sensitivity to outliers (k-medoids replaces the means by medians);
- Tendency to produce equi-sized clusters; 0

.

The number of clusters must be known beforehand.

LIRIS Limitations of k-means

- Convergence towards a *local* maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters;
- Sensitivity to outliers (k-medoids replaces the means by medians);
- Tendency to produce equi-sized clusters;

.

• The number of clusters must be known beforehand.

LIRIS *k-means* The elbow method

Plot a measure of the quality of the k clusters (e.g., the sum, over all objects, of the similarities to the centers of the assigned clusters) when k increases. Choose k after a large drop of the growth.

k-means

The elbow method

Plot a measure of the quality of the k clusters (e.g., the sum, over all objects, of the similarities to the centers of the assigned clusters) when k increases. Choose k after a large drop of the growth.

More principled method exist and can be seen as variants (finding the best trade-off between quality and compression).

k-means

LIRIS The elbow method

Δ

Plot a measure of the quality of the k clusters (e.g., the sum, over all objects, of the similarities to the centers of the assigned clusters) when k increases. Choose k after a large drop of the growth.

More principled method exist and can be seen as variants (finding the best trade-off between quality and compression).

If the quadratic time complexity of a hierarchical agglomeration is not prohibitive, the number of clusters can be determined from the dendrogram.

Marc Plantevit

Clustering 19 / 52

LIRIS Limitations of k-means

- Convergence towards a *local* maximum of the sum, over all objects, of the similarities to the centers of the assigned clusters;
- Sensitivity to outliers (k-medoids replaces the means by medians);
- Tendency to produce equi-sized clusters;
- The number of clusters must be known beforehand.

k-means

LIRIS

Tendency to produce equi-sized clusters



Marc Plantevit

Clustering 21 / 52







- Hierarchical Clustering
- Openation Description (Intersection) Descript

5 Conclusion

Marc Plantevit

Clustering 22 / 52

LIRIS EM EM EM

The dataset \mathcal{D} is seen as a random sample from a $|\mathcal{A}|$ -dimensional random variable O.

LIRIS EM assumptions

The dataset \mathcal{D} is seen as a random sample from a $|\mathcal{A}|$ -dimensional random variable O.

This probability density function is given as a mixture model of the $k \in \mathbb{N} \setminus \{0\}$ clusters $(C_i)_{i=1..k}$:

$$f(o) = \sum_{i=1}^{k} f_i(o) P(C_i)$$

, where $P(C_i)$ is the probability to belong to the cluster C_i and f_i is the probability density function of this cluster whose type of distribution is chosen beforehand.

LIRIS <u>
Maximum likelihood estimation</u>

EM searches a parametrization θ of f (i.e., $(P(C_i)_{i=1..k})$ and the parametrization of the $(f_i)_{i=1..k}$) so that the likelihood that \mathcal{D} is indeed a random sample of O is maximized:

 $\arg\max_{\theta} P(\mathcal{D}|\theta) \ .$

LIRIS Maximum likelihood estimation

EM searches a parametrization θ of f (i.e., $(P(C_i)_{i=1..k})$ and the parametrization of the $(f_i)_{i=1..k}$) so that the likelihood that \mathcal{D} is indeed a random sample of O is maximized:

$$rgmax_{ heta} P(\mathcal{D}| heta)$$
 .

Since the dataset is assumed to be a random sample from O (i. e., independent and identically distributed as O), the objective becomes the computation of:

$$\arg\max_{\theta} \prod_{o \in \mathcal{O}} f(o)$$
.

LIRIS Maximum likelihood estimation

EM searches a parametrization θ of f (i.e., $(P(C_i)_{i=1..k})$ and the parametrization of the $(f_i)_{i=1..k}$) so that the likelihood that \mathcal{D} is indeed a random sample of O is maximized:

$$rgmax_{ heta} P(\mathcal{D}| heta)$$
 .

Since the dataset is assumed to be a random sample from O (i. e., independent and identically distributed as O), the objective becomes the computation of:

$$rg\max_{ heta}\prod_{o\in\mathcal{O}}f(o)$$
 .

It usually is hard to analytically compute $\arg \max_{\theta} \prod_{o \in \mathcal{O}} f(o)$.

Marc Plantevit

EM is a greedy iterative approach that always converges to a local maximum of $P(\mathcal{D}|\theta)$.

EM is a greedy iterative approach that always converges to a local maximum of $P(D|\theta)$.

An iteration consists in two steps:

E Given θ , the posterior probabilities of each object to belong to each cluster is computed;

EM is a greedy iterative approach that always converges to a local maximum of $P(D|\theta)$.

An iteration consists in two steps:

- E Given θ , the posterior probabilities of each object to belong to each cluster is computed;
- M θ is updated to reflect these probabilities.

EM is a greedy iterative approach that always converges to a local maximum of $P(D|\theta)$.

An iteration consists in two steps:

E Given θ , the posterior probabilities of each object to belong to each cluster is computed;

M θ is updated to reflect these probabilities.

Initially, the parametrization of θ is randomly drawn and $\forall i = 1..k, P(C_i) = \frac{1}{k}$. The procedure stops when, from an iteration to the next one, the parametrization has not changed much (or at all).

Marc Plantevit

LIRIS Expectation step

Р(

Given θ , the posterior probability of an object $o \in O$ to belong to a cluster C_i is:

$$C_{i}|o) = \frac{P(C_{i} \land o)}{P(o)}$$

$$= \frac{P(o|C_{i})P(C_{i})}{\sum_{a=1..k} P(o \land C_{a})}$$

$$= \frac{P(o|C_{i})P(C_{i})}{\sum_{a=1..k} P(o|C_{a})P(C_{a})}$$

$$= \frac{f_{i}(o)P(C_{i})}{\sum_{a=1..k} f_{a}(o)P(C_{a})} \cdot$$

Marc Plantevit

Clustering 26 / 52

$\mathbf{LIRIS} \stackrel{\text{EM}}{\text{Maximization step (1/2)}}$

The distribution of a cluster usually is assumed multivariate normal, thus parametrized with a *location* (the center of the cluster) and a *covariance matrix*.

LIRIS $\frac{M}{Maximization step (1/2)}$

The distribution of a cluster usually is assumed multivariate normal, thus parametrized with a *location* (the center of the cluster) and a *covariance matrix*.

Given $(P(C_i|o))_{i=1..k,o\in\mathcal{O}}$, the location of the cluster C_i is updated to the weighted sample mean μ_i :

$$\frac{\sum_{o \in \mathcal{O}} P(C_i|o)o}{\sum_{o \in \mathcal{O}} P(C_i|o)}$$

$\frac{1}{Maximization step (2/2)}$

Given $(P(C_i|o))_{i=1..k,o\in\mathcal{O}}$, the covariance of the cluster C_i between the random variables O_a and O_b is updated to the weighted sample covariance:

$$\frac{\sum_{o \in \mathcal{O}} P(C_i|o)(o_a - \mu_{i,a})(o_b - \mu_{i,b})}{\sum_{o \in \mathcal{O}} P(C_i|o)} \ .$$

$\frac{1}{Maximization step (2/2)}$

Given $(P(C_i|o))_{i=1..k,o\in\mathcal{O}}$, the covariance of the cluster C_i between the random variables O_a and O_b is updated to the weighted sample covariance:

$$\frac{\sum_{o \in \mathcal{O}} P(C_i|o)(o_a - \mu_{i,a})(o_b - \mu_{i,b})}{\sum_{o \in \mathcal{O}} P(C_i|o)}$$

Given $(P(C_i|o))_{i=1..k,o\in\mathcal{O}}$, the prior probability of belonging to the cluster C_i is updated to:

$$\frac{\sum_{o\in\mathcal{O}}P(C_i|o)}{|\mathcal{O}|}$$

LIRIS $\frac{\mathsf{EM}}{\mathsf{EM}}$ with $|\mathcal{A}| = 1$ and k = 2: illustration

EM clustering of the objects in a one-dimensional space using the Euclidean distance.

Dataset:



LIRIS $\frac{\mathsf{EM}}{\mathsf{EM}}$ with $|\mathcal{A}| = 1$ and k = 2: illustration

EM clustering of the objects in a one-dimensional space using the Euclidean distance.



Marc Plantevit

LIRIS $\frac{\mathsf{EM}}{\mathsf{EM}}$ with $|\mathcal{A}| = 1$ and k = 2: illustration

EM clustering of the objects in a one-dimensional space using the Euclidean distance.

Iteration 5:



LIRIS $\frac{\mathsf{EM}}{\mathsf{EM} \text{ with } |\mathcal{A}| = 2 \text{ and } k = 3}$

EM clustering of the objects in a two-dimensional space using the Euclidean distance.

Dataset:



LIRIS $\frac{\mathsf{EM}}{\mathsf{EM} \text{ with } |\mathcal{A}| = 2 \text{ and } k = 3}$

EM clustering of the objects in a two-dimensional space using the Euclidean distance.

Iteration 1:



LIRIS $\frac{\mathsf{EM}}{\mathsf{EM} \text{ with } |\mathcal{A}| = 2 \text{ and } k = 3}$

EM clustering of the objects in a two-dimensional space using the Euclidean distance.

Iteration 36:



EM

EM algorithm with mixture of Gaussians

Input: $\mathcal{O}, \mathcal{D}, k \in \mathbb{N} \setminus \{0\}$

Output: a fuzzy clustering of \mathcal{O} corresponding to posterior probabilities of a *locally* maximized likelihood of a mixture of Gaussians

$$\begin{array}{l} (\mu_i)_{i=1..k} \leftarrow \operatorname{random}(\mathcal{D}) \\ (\Sigma_i)_{i=1..k} \leftarrow (I, \dots, I) \\ (P(C_i))_{i=1..k} \leftarrow (\frac{1}{k}, \dots, \frac{1}{k}) \\ \textbf{repeat} \\ (P(C_i|o))_{i=1..k,o\in\mathcal{O}} \leftarrow \\ \text{expectation}(\mathcal{O}, \mathcal{D}, (\mu_i)_{i=1..k}, (\Sigma_i)_{i=1..k}, (P(C_i))_{i=1..k}) \\ (c, (\mu_i)_{i=1..k}, (\Sigma_i)_{i=1..k}, (P(C_i))_{i=1..k}) \leftarrow \\ \text{maximization}(\mathcal{D}, (P(C_i|o))_{i=1..k,o\in\mathcal{O}}, (\mu_i)_{i=1..k}) \\ \textbf{until } c \\ \textbf{output}((P(C_i|o))_{i=1..k,o\in\mathcal{O}}) \end{array}$$

Marc Plantevit

Clustering 31 / 52

LIRIS expectation

Input: $\mathcal{O}, \mathcal{D}, (\mu_i)_{i=1..k} \in (\mathbb{R}^{|\mathcal{A}|})^k, (\Sigma_i)_{i=1..k} \in$ $(\mathbb{R}^{|\mathcal{A}| \times |\mathcal{A}|})^k, (P(C_i))_{i=1..k} \in [0,1]^k$ **Output:** $(P(C_i|o))_{i=1..k,o\in\mathcal{O}}$ the fuzzy assignment of the objects in \mathcal{O} to the clusters given by the mixture of Gaussians parametrized with $(\mu_i)_{i=1..k}, (\Sigma_i)_{i=1..k}, (P(C_i))_{i=1..k}$ for all $o \in \mathcal{O}$ do for $i = 1 \rightarrow k$ do $P(C_i|o) \leftarrow \frac{f_i(o)P(C_i)}{\sum_{i=1}^{k} f_i(o)P(C_i)}$ end for end for $\operatorname{return}((P(C_i|o))_{i=1, k, o \in \mathcal{O}})$

Marc Plantevit

LIRIS expectation

Input: $\mathcal{O}, \mathcal{D}, (\mu_i)_{i=1..k} \in (\mathbb{R}^{|\mathcal{A}|})^k, (\Sigma_i)_{i=1..k} \in$ $(\mathbb{R}^{|\mathcal{A}| \times |\mathcal{A}|}_{\perp})^k, (P(C_i))_{i=1..k} \in [0,1]^k$ **Output:** $(P(C_i|o))_{i=1..k,o\in\mathcal{O}}$ the fuzzy assignment of the objects in \mathcal{O} to the clusters given by the mixture of Gaussians parametrized with $(\mu_i)_{i=1,k}, (\Sigma_i)_{i=1,k}, (P(C_i))_{i=1,k}$ for all $o \in \mathcal{O}$ do for $i = 1 \rightarrow k$ do $P(C_i|o) \leftarrow \frac{(\det(\Sigma_i)e^{(o-\mu_i)^T \Sigma_i^{-1}(o-\mu_i)})^{\frac{-1}{2}}P(C_i)}{\sum_{i=1}^{k} (\det(\Sigma_i)e^{(o-\mu_i)^T \Sigma_a^{-1}(o-\mu_i)})^{\frac{-1}{2}}P(C_i)}$ end for end for return $((P(C_i|o))_{i=1} k_{o \in \mathcal{O}})$

LIRIS <u> EM</u> <u> Complexity of expectation</u>

Question

What is the complexity of computing $(det(\Sigma_i))_{i=1..k}$?

LIRIS <u> Complexity of expectation</u>

Question

What is the complexity of computing $(\det(\Sigma_i))_{i=1..k}$? $kO(|\mathcal{A}|^3)$.

LIRIS Complexity of expectation

Question

What is the complexity of computing $(\det(\Sigma_i))_{i=1..k}$? $kO(|\mathcal{A}|^3)$.

Question

What is the complexity of computing $(\Sigma_i^{-1})_{i=1..k}$ once $(\det(\Sigma_i))_{i=1..k}$ is known?

LIRIS Complexity of expectation

Question

What is the complexity of computing $(\det(\Sigma_i))_{i=1..k}$? $kO(|\mathcal{A}|^3)$.

Question

What is the complexity of computing $(\Sigma_i^{-1})_{i=1..k}$ once $(\det(\Sigma_i))_{i=1..k}$ is known? $O(k|\mathcal{A}|^2)$.

LIRIS Complexity of expectation

Question

What is the complexity of computing $(\det(\Sigma_i))_{i=1..k}$? $kO(|\mathcal{A}|^3)$.

Question

What is the complexity of computing $(\Sigma_i^{-1})_{i=1..k}$ once $(\det(\Sigma_i))_{i=1..k}$ is known? $O(k|\mathcal{A}|^2)$.

Question

What is the complexity of computing one Mahalanobis distance, $(o, o') \mapsto (o - o')^T \Sigma_i^{-1} (o - o')$, once Σ_i^{-1} is known?
LIRIS Complexity of expectation

Question

What is the complexity of computing $(\det(\Sigma_i))_{i=1..k}$? $kO(|\mathcal{A}|^3)$.

Question

What is the complexity of computing $(\Sigma_i^{-1})_{i=1..k}$ once $(\det(\Sigma_i))_{i=1..k}$ is known? $O(k|\mathcal{A}|^2)$.

Question

What is the complexity of computing one Mahalanobis distance, $(o, o') \mapsto (o - o')^T \Sigma_i^{-1} (o - o')$, once Σ_i^{-1} is known? $O(|\mathcal{A}|^2)$.

LIRIS Complexity of expectation

Question

What is the complexity of computing $(\det(\Sigma_i))_{i=1..k}$? $kO(|\mathcal{A}|^3)$.

Question

What is the complexity of computing $(\Sigma_i^{-1})_{i=1..k}$ once $(\det(\Sigma_i))_{i=1..k}$ is known? $O(k|\mathcal{A}|^2)$.

Question

What is the complexity of computing one Mahalanobis distance, $(o, o') \mapsto (o - o')^T \Sigma_i^{-1} (o - o')$, once Σ_i^{-1} is known? $O(|\mathcal{A}|^2)$.

Question

What is the complexity of expectation?

Marc Plantevit

FM LIRIS **Complexity of expectation**

Question

What is the complexity of computing $(\det(\Sigma_i))_{i=1,k}$? $kO(|\mathcal{A}|^3)$.

Question

What is the complexity of computing $(\sum_{i=1}^{-1})_{i=1,k}$ once $(\det(\Sigma_i))_{i=1..k}$ is known? $O(k|\mathcal{A}|^2)$.

Question

What is the complexity of computing one Mahalanobis distance, $(o, o') \mapsto (o - o')^T \Sigma_i^{-1} (o - o')$, once Σ_i^{-1} is known? $O(|\mathcal{A}|^2)$.

Question

What is the complexity of **expectation**? $O(k|\mathcal{A}|^2(|\mathcal{O}| + |\mathcal{A}|))$.

Input: $\mathcal{D}, (P(C_i|o))_{i=1..k,o\in\mathcal{O}} \in [0,1]^{k|\mathcal{O}|}, (\mu_i)_{i=1..k} \in (\mathbb{R}^{|\mathcal{A}|})^k$ **Output:** $c \in \{$ false, true $\}$ indicating whether the convergence is reached, the new parametrization of the mixture of Gaussians $c \leftarrow true$ for $i = 1 \rightarrow k$ do $\mu_i' \leftarrow \frac{\sum_{o \in \mathcal{O}} P(C_i|i)o}{\sum_{o \in \mathcal{O}} P(C_i|i)}$ if $\mu'_i \neq \mu_i$ then $c \leftarrow false$ end if $\Sigma_{i}^{\prime} \leftarrow \frac{\sum_{o \in \mathcal{O}} P(C_{i}|o)(o-\mu_{i}^{\prime})(o-\mu_{i}^{\prime})^{T}}{\sum_{o \in \mathcal{O}} P(C_{i}|o)} \\ P(C_{i})^{\prime} \leftarrow \frac{\sum_{o \in \mathcal{O}} P(C_{i}|o)}{|\mathcal{O}|}$ end for return $(c, (\mu'_i)_{i=1,k}, (\Sigma'_i)_{i=1,k}, (P(C_i)')_{i=1,k})$

Complexity of EM

Question

What is the complexity of **expectation**? $O(k|\mathcal{A}|^2(|\mathcal{O}| + |\mathcal{A}|))$.

Question

What is the complexity of maximization?

EM LIRIS **Complexity of EM**

Question

What is the complexity of **expectation**? $O(k|\mathcal{A}|^2(|\mathcal{O}| + |\mathcal{A}|))$.

Question

What is the complexity of **maximization**? $O(k|\mathcal{O}||\mathcal{A}|^2)$.

Complexity of EM

Question

What is the complexity of **expectation**? $O(k|\mathcal{A}|^2(|\mathcal{O}| + |\mathcal{A}|))$.

Question

What is the complexity of **maximization**? $O(k|\mathcal{O}||\mathcal{A}|^2)$.

Question

What is the complexity of EM if $t \in \mathbb{N}$ iterations are necessary to converge?

Complexity of EM

Question

What is the complexity of **expectation**? $O(k|\mathcal{A}|^2(|\mathcal{O}| + |\mathcal{A}|))$.

Question

What is the complexity of **maximization**? $O(k|\mathcal{O}||\mathcal{A}|^2)$.

Question

What is the complexity of EM if $t \in \mathbb{N}$ iterations are necessary to converge? $O(tk|\mathcal{A}|^2(|\mathcal{O}| + |\mathcal{A}|))$.

LIRIS Diagonal covariance matrix

A lower complexity is obtained by assuming all attributes independent, i. e., all covariance matrices diagonal. The operations involving such a matrix become linear in $|\mathcal{A}|$ and the total time complexity of EM becomes $O(tk|\mathcal{O} \times \mathcal{A}|)$.

LIRIS Diagonal covariance matrix

A lower complexity is obtained by assuming all attributes independent, i. e., all covariance matrices diagonal. The operations involving such a matrix become linear in $|\mathcal{A}|$ and the total time complexity of EM becomes $O(tk|\mathcal{O} \times \mathcal{A}|)$.

However, if the attributes are not really independent, the obtained fuzzy clustering become much worse:



Marc Plantevit

/ 52

 $\frac{1}{k-\text{means as specialization of EM}}$

k-means is EM with f_i chosen as follows:

EM

$$\begin{cases} 1 \text{ if } C_i = \arg \max_{a=1..k} s(o, \mu_a) \\ 0 \text{ otherwise} \end{cases}$$

Marc Plantevit

.







3 Hierarchical Clustering

Density-based Clustering: DBSCAN

5 Conclusion

Marc Plantevit

Hierarchical Clustering

- Build a hierarchy of clusters (not an unique partition);
- The number of clusters k is not required as input;
- Use a distance matrix as clustering criteria
- An early-termination condition can be used (ex. nb clusters).

LIRIS Algorithm

Input: a sample of m objets x_1, \ldots, x_m .

- **()** The algorithm begins with m clusters (1 cluster = 1 object);
- O Merge the 2 clusters that are the closest.
- Ind If it remains only one cluster.
- Go to step 2.

Output: a dendogram



A hierarchy that can be split at a given level to form a partition.

- the hierarchy: a tree called dendogram
- the leaves = the objects

Marc Plantevit

LIRIS

Clustering

LIRIS Distance between clusters

- Distance between the centers (centroid method)
- Minimal distance among the pairs composed of objects from the two clusters (Single Link Method):

$$d(i,j) = \min_{x \in C_i; y \in C_j} d(x,y)$$

• Maximal distance among the pairs composed of objects from the two clusters (Complete Link Method):

$$d(i,j) = max_{x \in C_i; y \in C_j} d(x,y)$$

• Average distance among the pairs composed of objects from the two clusters (Average Linkage Method):

$$d(i,j) = avg_{x \in C_i; y \in C_j}d(x,y)$$

Marc Plantevit

Clustering

Pros:

- Conceptually simple.
- Theoretical properties well-known.

Cons:

- The clustering is definitive: erroneous decisions are impossible to modify later.
- Non-extensible method for large collections of objects $(\theta(n^2))$







- Hierarchical Clustering
- Density-based Clustering: DBSCAN

5 Conclusion



• For this kind of problem, the use of similarity (or distance) measures is less efficient than the use of neighborhood density

Density-based clustering

- Clusters are seen as dense regions separated by regions that are much less denser (noise)
- Two parameters:

LIRIS

- Eps: The maximum radius of the neighborhood
- MinPts: Minimum number of points within the Eps-neighborhood of a point.
- Neighborhood: $V_{Eps}(p) = \{q \in D \mid dist(p,q) \le Eps\}$
- A point p is directly density-accessible from q w.r.t. Eps, MinPts if

p MinPts = 5 Q Eps = 1 cm

 $P \in V_{\textit{Eps}}(q)$ and $|V_{\textit{Eps}}(q)| \ge \textit{MinPts}$

Clustering 46 / 52

- Accessibility: p is accessible from q w.r.t. Eps, MinPts if there exists p₁,..., p_n such that p₁ = q, p_n = p and p_{i+1} is directly accessible from p_i.
- Connexity: *p* is connected to *q* w.r.t. Eps and MinPts if there exists a point *o* such that *p* and *q* are accessible from *o*.





LIRIS DBSCAN: Density Based Spatial Clustering of Applications with Noise

- A cluster is the maximal set of connected points
- Cluster shapes are not necessary convex



Marc Plantevit

Clustering 48 / 52

DBSCAN Algorithm

- Choose p
- Retrieve all poinst that accessible from p (w.r.t. Eps and MinPts)
- If p is a center, then a cluster is created.
- If *p* is a limit, then there is not accessible point from *p*, **Skip to another point**

49 / 52

• Repeat until it remains no point.







- Hierarchical Clustering
- Openation Density-based Clustering: DBSCAN

Marc Plantevit

Clustering 50 / 52

LIRIS Summary

• *k*-means iteratively assigns each object to the cluster whose center is the most similar and recompute these centers as the mean of the objects they were assigned;

Clustering 51 / 52

LIRIS Summary

- *k*-means iteratively assigns each object to the cluster whose center is the most similar and recompute these centers as the mean of the objects they were assigned;
- Its strongest advantage is its low complexity;

Clustering 51 / 52

LIRIS Summary

- *k*-means iteratively assigns each object to the cluster whose center is the most similar and recompute these centers as the mean of the objects they were assigned;
- Its strongest advantage is its low complexity;
- Its worst drawback is its tendency to discover equi-sized clusters;

LIRIS Summary

- *k*-means iteratively assigns each object to the cluster whose center is the most similar and recompute these centers as the mean of the objects they were assigned;
- Its strongest advantage is its low complexity;
- Its worst drawback is its tendency to discover equi-sized clusters;
- *k*-means actually is a specialization of a whole class of algorithms called EM;

LIRIS Summary

- *k*-means iteratively assigns each object to the cluster whose center is the most similar and recompute these centers as the mean of the objects they were assigned;
- Its strongest advantage is its low complexity;
- Its worst drawback is its tendency to discover equi-sized clusters;
- *k*-means actually is a specialization of a whole class of algorithms called EM;
- They treat the dataset as a random sample of a multivariate random variable whose pdf is given as a mixture model;

LIRIS Summary

- *k*-means iteratively assigns each object to the cluster whose center is the most similar and recompute these centers as the mean of the objects they were assigned;
- Its strongest advantage is its low complexity;
- Its worst drawback is its tendency to discover equi-sized clusters;
- *k*-means actually is a specialization of a whole class of algorithms called EM;
- They treat the dataset as a random sample of a multivariate random variable whose pdf is given as a mixture model;
- They *locally* maximize the likelihood, i.e., the probability of observing the dataset given the parametrization of the mixture model;

Clustering 51 / 52

LIRIS Summary

- k-means iteratively assigns each object to the cluster whose center is the most similar and recompute these centers as the mean of the objects they were assigned;
- Its strongest advantage is its low complexity;
- Its worst drawback is its tendency to discover equi-sized clusters;
- k-means actually is a specialization of a whole class of algorithms called EM;
- They treat the dataset as a random sample of a multivariate random variable whose pdf is given as a mixture model;
- They *locally* maximize the likelihood, i.e., the probability of observing the dataset given the parametrization of the mixture model:
- They iteratively compute the expectation of the likelihood and update the parametrization so that this expectation is maximized.

The end.

Marc Plantevit

Clustering 52 / 52