

Clustering (Part 3)

OUTLINE

① k-means clustering

- Review of popular algorithms: LLoyd's, k-means++.
- Coreset-based algorithm: MR-kmeans.

② k-median clustering (sketch)

③ Clustering evaluation

k-means clustering

Problem definition

WE FOCUS ON: **k -means clustering for Euclidean spaces**

Given a pointset P of N points in \mathbb{R}^D , under the standard L_2 -distance function d , and an integer $k > 0$, determine a set $S \subset \mathbb{R}^D$ of k centers which minimizes

$$\Phi_{\text{kmeans}}(P, S) = \sum_{x \in P} (d(x, S))^2.$$

Observations:

- k-means clustering aims at **minimizing cluster variance**
- Unlike the formulation presented in the previous set of slides, now the centers are not **required to belong to the input P** .

Properties of Euclidean spaces

Let $X = (X_1, X_2, \dots, X_D)$ and $Y = (Y_1, Y_2, \dots, Y_D)$ be two points in \mathbb{R}^D . Recall that their L_2 -distance is

$$d(X, Y) = \sqrt{\sum_{i=1}^D (X_i - Y_i)^2} \triangleq \|X - Y\|.$$

Definition

The **centroid** of a set P of N points in \mathbb{R}^D is

$$c(P) = \frac{1}{N} \sum_{X \in P} X,$$

Sum component-wise

where the sum is component-wise. Note that $c(P)$ does not necessarily belong to P .

Fact

The **centroid** $c(P)$ of a set $P \subset \mathbb{R}^D$ is the point of \mathbb{R}^D which minimizes the sum of the square distances to all points of P .

Lloyd's algorithm

In 2006 it was included in the top-10 most influential
determining algorithms

Input: Set P of N points from \mathbb{R}^D , integer $k > 1$

Output: Set $S \subset \mathbb{R}^D$ of k centers with small $\Phi_{\text{kmeans}}(P, S)$.

$S = \{c_1, c_2, \dots, c_k\} \leftarrow$ arbitrary set of k points in \mathbb{R}^D ;

$\Phi \leftarrow \Phi_{\text{kmeans}}(P, S)$;

while (true) **do**

for $1 \leq i \leq k$ **do**

$C_i \leftarrow$ points of P closest to c_i ; // break ties arbitrarily

$c_i \leftarrow$ centroid of C_i ;

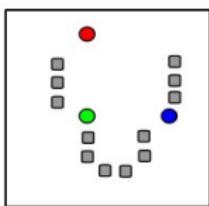
 Let $S = \{c_1, c_2, \dots, c_k\}$;

if ($\Phi_{\text{kmeans}}(P, S) < \Phi$) **then** $\Phi \leftarrow \Phi_{\text{kmeans}}(P, S)$;

else return S ;

$$c_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$$

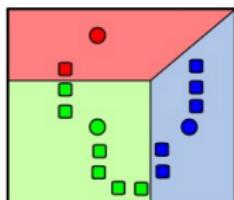
Example



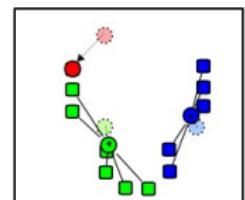
Initial centers

Iteration 1:

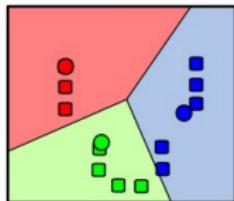
Partition



Calculation of centroids



Iteration 2:



Weaknesses of Lloyd's algorithm

- ① The number of iterations can be exponential in the input size (in the worst case, $\geq 2^{\Omega(\sqrt{N})}$).

COUNTERACTION: Fix number of iterations or stop iterations when the improvement of the objective function drops below a fixed threshold.

- ② The algorithm may be trapped into a local optimum with objective function arbitrarily far from optimal value.

COUNTERACTION: Make a careful selection of initial centers (which improves also convergence).

Rather than selecting the initial set of centers in an arbitrary or random fashion,
we can use k-means++

k-means++ (2007)

Let P be a set of N points from \mathbb{R}^D , and let $k > 1$ be an integer.

Algorithm k-means++ computes a (initial) set S of k centers for P using the following procedure (note that $S \subseteq P$):

$c_1 \leftarrow$ random point chosen from P with uniform probability;

$S \leftarrow \{c_1\}$;

for $2 \leq i \leq k$ **do**

foreach $x \in P - S$ **do** $\pi(x) \leftarrow (d(x, S))^2 / \sum_{y \in P-S} (d(y, S))^2$;

$c_i \leftarrow$ random point in $P - S$ according to distribution $\pi(\cdot)$;

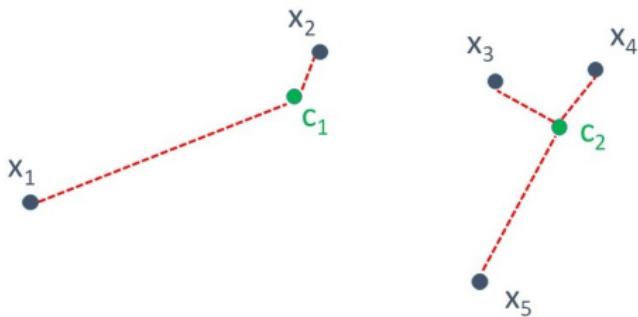
$S \leftarrow S \cup \{c_i\}$

return S

$$\sum_{x \in P-S} \pi(x) = 1$$

Obs. The set S returned by k-means++ is a subset of P

k-means++



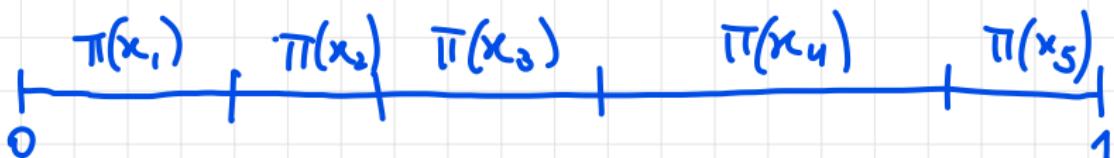
Selection of c_3

$$\pi(x_1) = \frac{(d(x_1, c_1))^2}{(d(x_1, c_1))^2 + (d(x_2, c_1))^2 + (d(x_3, c_2))^2 + (d(x_4, c_2))^2 + (d(x_5, c_2))^2}$$

k-means++

Suppose that we know $\pi(x) \neq x \in P-S$

How do we determine the next center



Draw a random number $r \in [0, 1]$ with uniform probability



Then select x_i as next center

k-means++

Theorem (Arthur-Vassilvitskii'07)

Let S be the set of centers returned by k-means++ for a pointset P . We have

$$E[\Phi_{\text{kmeans}}(P, S)] \leq 8(\ln k + 2) \cdot \Phi_{\text{kmeans}}^{\text{opt}}(P, k).$$

IMPORTANT REMARK: S and, consequently, $\Phi_{\text{kmeans}}(P, S)$, are random variables which depend on the random center selection!

Observations

What do we know so far about k-means clustering

* IN PRACTICE

k-means++ followed by some iterations of Lloyd's algorithm yields a good-quality solution. This approach can be implemented efficiently, even in MapReduce, but for very large points and if many iterations of Lloyd's are executed, then performance degraded

* IN THEORY : if we want constant approximation =

{ - 50 approximation (solve what efficiently)
- 6.36 approximation (inefficiently)
) they cannot handle large datasets

* For general metric spaces, getting up
good solutions efficiently is
even harder

k-means clustering for big data

How can we compute a “good” k-means clustering of a pointset P that is too large for a single machine?

Main alternatives:

- Distributed (e.g., MapReduce) implementation of Lloyd’s algorithm and k-means++.
(The next two exercises explore this alternative.)
- Coreset-based approach, similar in spirit to what was done for k-center. This approach lends itself to efficient distributed (e.g., MapReduce) as well as streaming implementations.

Lloyd's and k-means++ in MapReduce

Do the following two exercises assuming that the algorithms are applied to a set P of N points from \mathbb{R}^D , with D constant, and that the target number of clusters is $k \in (1, \sqrt{N}]$.

Exercise

Show how to implement Lloyd's algorithm in MapReduce using $O(1)$ rounds per iteration, local space $M_L = O(\sqrt{N})$ and aggregate space $M_A = O(N)$. (Hint: Make the set of centers and Φ global variables.)

Exercise

Show how to implement algorithm k-means++ in MapReduce using $O(k)$ rounds, local space $M_L = O(\sqrt{N})$ and aggregate space $M_A = O(N)$.

Observations

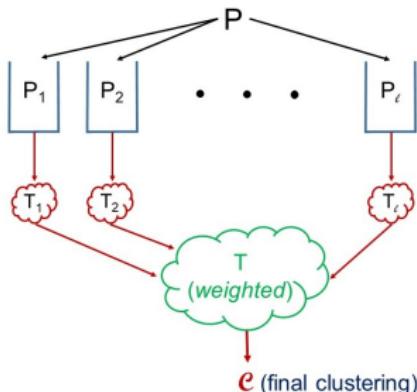
- * Spark (mllib library) provides implementations of Lloyd's and k-means++ algorithm for k-means++. Spark uses a parallel variant called k-means||
- * k-means++ yields also good BICRITERIA solutions

$k\text{-means} \Rightarrow S$ of $k' > k$ centers, say $k' = \beta k$
 $\beta \in O(1)$ is sufficient to ensure that $\beta > 1$

$$\Phi_{k\text{-means}}(P, S) = O(1) \Phi_{k\text{-means}}^{\text{opt}}(P, k)$$

Coreset-based approach for k-means

We use a (composable) coresset-based approach similar as for k-center.



- The local coresets T_i 's and the final clustering C are computed using a sequential algorithm for k-means.
- Each point of T is given a weight equal to the number of points in P it represents, so to account for their cumulative contribution to the objective function

Need a weighted variant of k-means clustering.

Weighted k-means clustering

Weighted variant of the k-means clustering problem

Input:

- Set P of N points from \mathbb{R}^D .
- Weight $w(x) > 0$ for every $x \in P$.
- Target number k of clusters.

Output: Set S of k centers in \mathbb{R}^D minimizing the following objective function.

$$\Phi_{\text{kmeans}}^w(P, S) = \sum_{x \in P} w(x) \cdot (d(x, S))^2.$$

↑
weighted obj. function

Weighted k-means clustering: observations

- The unweighted k-means clustering problem is a special case of the weighted variant where all weights are equal to 1.
- For integer weights, most known algorithms can be adapted to solve the weighted variant. It is sufficient to regard each point x as representing $w(x)$ identical copies of itself. The same approximation guarantees are maintained.

Lloyd's algorithm with weights

Modifications to the algorithm

- * Computation of a centroid c_i for C_i becomes

$$c_i = \frac{1}{\sum_{x \in C_i} w(x)} \sum_{x \in C_i} w(x) \cdot x$$

vector
↓ scalar

- * Objective function

$$\phi_{k\text{means}}(P, S) \rightarrow \phi_{k\text{means}}^w(P, S)$$

Example of weighted centroid



centroid (red)

$$\left(\frac{4+6+4+5}{9}, \frac{16+10+2+2}{9} \right) =$$
$$\left(\frac{19}{9}, \frac{30}{9} \right) \approx (2.11, 3.33)$$

k-means++ algorithm with weights

Modifications to the algorithm

S = current set of centers

For each $x \in P - S$ the probability of being selected as next center is

$$\pi(x) = \frac{w(x) \cdot (d(x, S))^2}{\sum_{y \in P - S} w(y) \cdot (d(y, S))^2}$$

Coreset-based MapReduce algorithm for k-means

$\text{MR-kmeans}(\mathcal{A})$: MapReduce algorithm for k-means (described in the next slide) which uses a sequential algorithm \mathcal{A} for k-means, as an argument (functional approach!).

In principle, one can instantiate \mathcal{A} with any sequential algorithm, as long as the following characteristics are met:

- \mathcal{A} solves the more general weighted variant.
- \mathcal{A} requires space proportional to the input size (if larger space is needed, the analysis of M_L for the MapReduce algorithm must be amended).

MR-kmeans(\mathcal{A})

Input Set P of N points in \Re^D , integer $k > 1$, sequential k-means algorithm \mathcal{A} .

Output Set S of k centers in \Re^D which is a good solution to the k-means problem on P .

Round 1:

- Map Phase: Partition P arbitrarily in ℓ subsets of equal size P_1, P_2, \dots, P_ℓ .
- Reduce Phase: for every $i \in [1, \ell]$ separately, run \mathcal{A} on P_i (with unit weights) to determine a set $T_i \subseteq P_i$ of k centers, and define
 - For each $x \in P_i$, the proxy $\tau(x)$ as x 's closest center in T_i .
 - For each $y \in T_i$, the weight $w(y)$ as the number of points of P_i whose proxy is y .

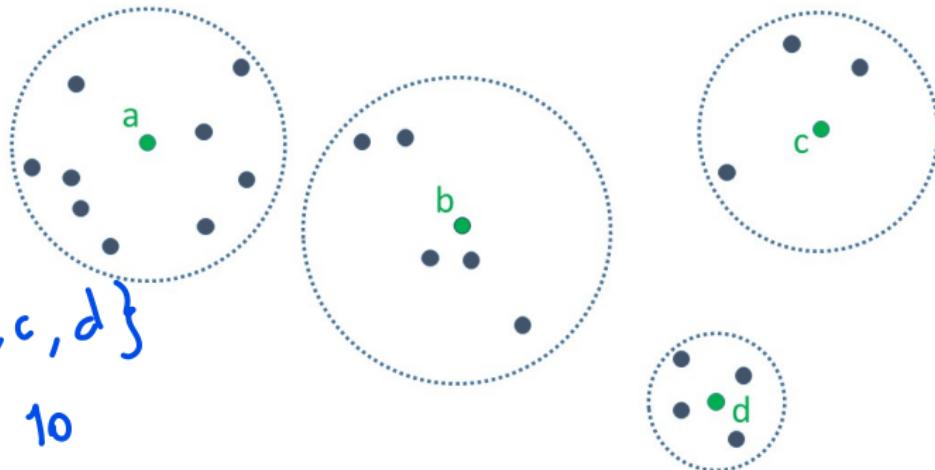
MR-kmeans(\mathcal{A})

Round 2:

- Map Phase: empty.
- Reduce Phase: gather the coresset $T = \cup_{i=1}^{\ell} T_i$ of $\ell \cdot k$ points, together with their weights, and run, using a single reducer, \mathcal{A} on T (with the given weights) to determine a set $S = \{c_1, c_2, \dots, c_k\}$ of k centers, which is then returned as output..

Example for Round 1

Set T_i (green points) computed in a partition P_i



$$T_i = \{a, b, c, d\}$$

$$w(a) = 10$$

$$w(b) = 6$$

$$w(c) = 4$$

$$w(d) = 5$$

if any of a, b, c, d is not in P_i
we subtract 1 from its weight

Analysis of MR-kmeans(\mathcal{A})

Assume $k \leq \sqrt{N}$. By setting $\ell = \sqrt{N/k}$, it is easy to see that MR-kmeans(\mathcal{A}) requires Same analyses as for MR-FFT

- Local space $M_L = O(\max\{N/\ell, \ell \cdot k\}) = O(\sqrt{N \cdot k}) = o(N)$
- Aggregate space $M_A = O(N)$

The quality of the solution returned by the algorithm is stated in the following theorem, which is proved in the next few slides.

Theorem

Suppose that \mathcal{A} is an α -approximation algorithm for weighted k -means clustering, for some $\alpha > 1$, and let S be the set of k centers returned by MR-kmeans(\mathcal{A}) on input P . Then:

$$\Phi_{\text{kmeans}}(P, S) = O(\alpha^2) \Phi_{\text{kmeans}}^{\text{opt}}(P, k).$$

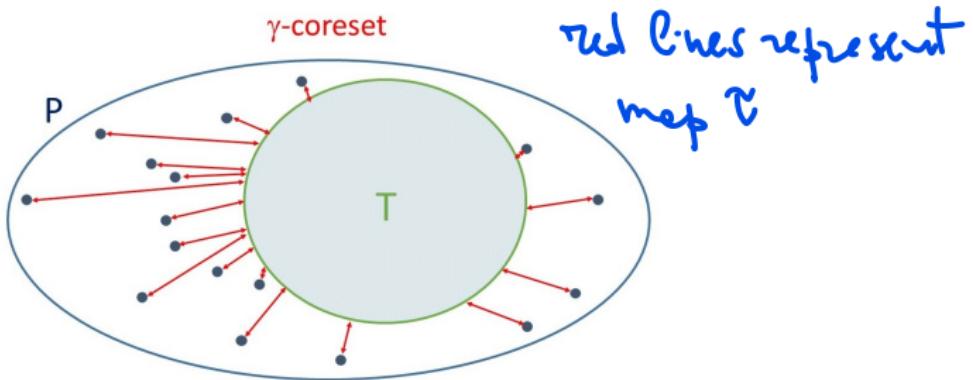
That is, MR-kmeans(\mathcal{A}) is an $O(\alpha^2)$ -approximation algorithm.

Analysis of MR-kmeans(\mathcal{A})

Definition

Given a pointset P , a coresset $T \subseteq P$ and a proxy function $\tau : P \rightarrow T$, we say that T is a γ -coreset for P, k and the k-means objective if

$$\sum_{p \in P} (d(p, \tau(p)))^2 \leq \gamma \cdot \Phi_{\text{kmeans}}^{\text{opt}}(P, k).$$



The smaller γ and the better T represents P

Proof of Theorem

Proof idea



Claim 1 The coresets T computed by $\text{MR-kmeans}(A)$

together with the map τ , is an α -core set

Recall $T = \bigcup_{i=1}^k T_i$; T_i = d-approximation for kmeans on P_i

$$\forall x \in P_i : \tau(x) \in T_i \text{ closest point of } T_i \text{ to } x$$

$$\Rightarrow d(x, \tau(x)) = d(x, T_i)$$

Proof of Theorem

- * S_i^{OPT} = optimal solution of k-means on P_i
- * S^{OPT} = optimal solution of k-means on P

$$\begin{aligned}
 \sum_{x \in P} (d(x, \tau(x)))^2 &= \sum_{i=1}^l \sum_{x \in P_i} (d(x, \tau_i(x)))^2 \\
 &= \sum_{i=1}^l \sum_{x \in P_i} (d(x, T_i))^2 \\
 &\leq \sum_{i=1}^l \alpha \sum_{x \in P_i} (d(x, S_i^{\text{OPT}}))^2 \\
 &\quad (\text{since } T_i \text{ is an } \alpha\text{-approx. sol.}) \\
 &= \alpha \sum_{i=1}^l \sum_{x \in P_i} (d(x, S_i^{\text{OPT}}))^2
 \end{aligned}$$

Proof of Theorem

$$\leq \alpha \sum_{i=1}^l \sum_{x \in P_i} (d(x, S^{OPT}))^2$$

$$= \alpha \sum_{x \in P} (d(x, S^{OPT}))^2$$

$$= \alpha \bigcirc_{k\text{-means}}^{OPT} (P, k)$$

$\Rightarrow T$ with proxy function $\gamma(\cdot)$ is an α -consist for P , k , and k -means \square

Proof of Theorem

Claim 2 If T with proxy function $\ell(\cdot)$ is an α coreset for P , k and k -means, then the solution S (computed by running A on T (weighted according to γ)) is such that

$$\phi_{k\text{mean}}(P, S) = O(\alpha^2) \overline{\phi}_{k\text{means}}^{\text{opt}}(P, k)$$

WE SKIP THIS PROOF

Proof of Theorem

Remark : It can be shown, in general, that if T is a γ -cover set, and we run an α -approximate algorithm A on T we get an $O((1+\gamma)\alpha)$ approximation.

The theorem follows immediately from the 2 claims

Proof of Theorem

Proof of Theorem

Observations

- Two different k-means sequential algorithms can be used in R1 and R2. For example, one could use k-means++ in R1, and k-means++ plus LLoyd's in R2.
- In practice, the algorithm is fast and accurate.
- If $k' > k$ centers are selected from each P_i in R1 (e.g., through k-means++), the quality of T , hence the quality of the final clustering, improves.

especially when P
comes from a space of
low dimensionality

k-median clustering (sketch)

k-median clustering

Given a pointset P from a metric space (M, d) and an integer $k > 0$, the **k-median clustering** problem requires to determine a set $S \subset P$ of k centers which minimizes:

$$\Phi_{\text{kmedian}}(P, S) = \sum_{x \in P} d(x, S).$$

- $O(1)$ sequential approximation algorithms for k-median are known but the most accurate ones are computationally intensive.
- The most popular sequential algorithm, **PAM: Partitioning Around Medoids** (a.k.a. k-medoids algorithm) is based on a local search optimization strategy.
- A coresnet-based approach similar to MR-kmeans can be employed where the use of accurate sequential strategies is confined to a coresnet much smaller than the input.

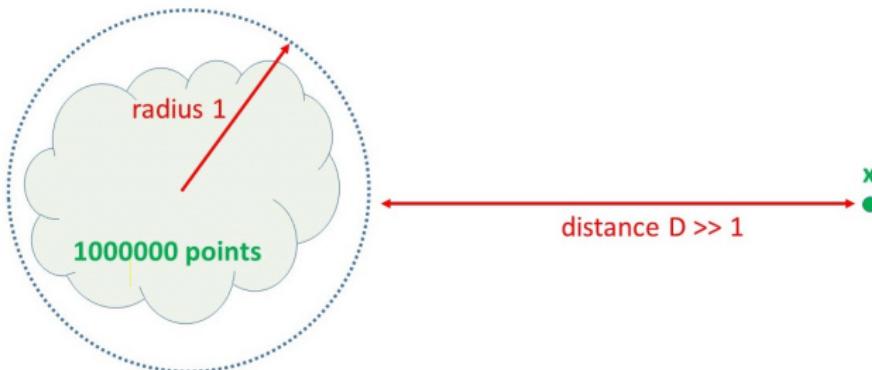
k-center vs k-means vs k-median

Some final considerations on the 3 center-based clustering problems.

- All 3 problems tend to return **spherically-shaped clusters**.
- They **differ in the center selection strategy** since, given the centers, the best clustering around them is the same for the 3 problems.
- Suitable **coreset-based strategies** can be employed for the 3 problems, featuring:
 - Ability to handle huge inputs.
 - Tunable performance-accuracy tradeoffs.
 - Very good behaviour in practice.
- The selection of the best k is often difficult, and this is a drawback common to the 3 problems.

k-center vs k-means vs k-median

The following example shows that different impact that outliers can have for the 3 problems.



For $k \geq 2$, the relevance of the outlier x in the 3 problems is different

- **k-center:** x must be always selected as center!
- **k-means:** x must be selected as center if $D \geq 1000$!
- **k-median:** x must be selected as center if $D \geq 1000000$!

Exercises

Exercise

Let S be the set of k centers returned by k-means++. Using Markov's inequality determine a value $\alpha > 1$ such that

$$\Pr(\Phi_{\text{kmeans}}(P, S) \leq \alpha \cdot \Phi_{\text{kmeans}}^{\text{opt}}(P, k)) \geq 1/2.$$

Recall that for a real-valued nonnegative random variable X with expectation $E[X]$ and a value $a > 0$, the Markov's inequality states that

$$\Pr(X \geq a) \leq \frac{E[X]}{a}.$$

Exercises

The next exercise shows how to make the approximation bound of algorithm k-means++ to hold with high probability rather than in expectation.

Exercise

From before, we know that if S is the set of centers computed by k-means++ for P , then there is a value $\alpha > 1$ such that

$$\Pr(\Phi_{\text{kmeans}}(P, S) \leq \alpha \cdot \Phi_{\text{kmeans}}^{\text{opt}}(P, k)) \geq 1/2.$$

Show that the above probability can be made $\geq 1 - 1/N$, if we run several *independent instances* of k-means++, and return the set of centers which yields the minimum value of the objective function, among the ones computed in the various runs.

Clustering evaluation

Unsupervised Evaluation

- Let P be a set of N points from a metric space (M, d)
- Let $\mathcal{C} = (C_1, C_2, \dots, C_k)$ be a partition of P into k clusters (i.e., a k -clustering ignoring cluster centers, if available).

Goal of unsupervised evaluation: assess the *quality* of \mathcal{C} without reference to external information or ground truth.

Quality measures:

- Specific objective function, if any, which was used to select \mathcal{C} among all feasible clusterings. For example k-center/k-means/k-median objective functions which, however, capture only intra-cluster similarity.
- Silhouette coefficient, which captures both intra-cluster similarity and inter-cluster dissimilarity.

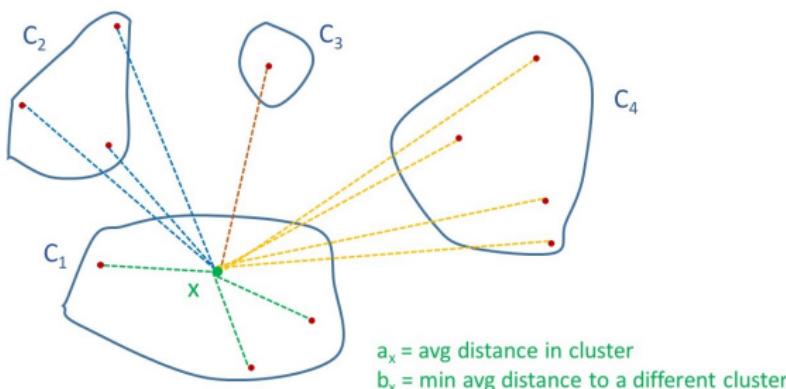
Observation: the silhouette coefficient is often used to select the most suitable number of clusters k , for a given clustering algorithm.

Unsupervised evaluation: silhouette

For a point $x \in P$ belonging to some cluster C_i let

$$\begin{aligned} a_x &= d_{\text{sum}}(x, C_i) / |C_i| \\ b_x &= \min_{j \neq i} d_{\text{sum}}(x, C_j) / |C_j|, \end{aligned}$$

where $d_{\text{sum}}(x, C)$ denotes the sum of the distances between x and the points of a cluster C (i.e., $d_{\text{sum}}(x, C) = \sum_{y \in C} d(x, y)$).



Unsupervised evaluation: silhouette

Definition: silhouette coefficient

Let $\mathcal{C} = (C_1, C_2, \dots, C_k)$ be a partition of P into k clusters. For any $x \in P$ the silhouette coefficient for x is

$$s_x = \frac{b_x - a_x}{\max\{a_x, b_x\}},$$

To measure the quality of \mathcal{C} we define the average silhouette coefficient as

$$s_{\mathcal{C}} = \frac{1}{|P|} \sum_{x \in P} s_x.$$

Unsupervised evaluation: silhouette

Note that

$$-\max\{a_x, b_x\} \leq b_x - a_x \leq \max\{a_x, b_x\}$$

$$\Rightarrow s_x \in [-1, 1] \Rightarrow s_c \in [-1, 1]$$

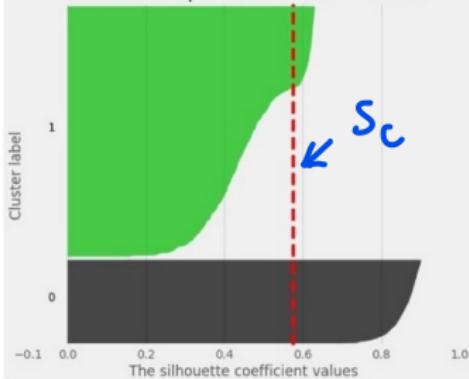
s_x close to $-1 \Rightarrow a_x \gg b_x \Rightarrow x$ is not clustered well

s_x close to $1 \Rightarrow a_x \ll b_x \Rightarrow x$ is clustered well

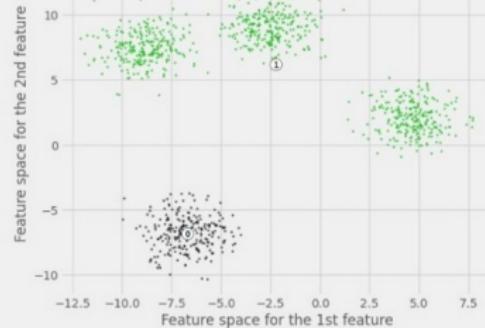
If s_c is close to 1 then C is a good clustering (\therefore e.g. most points are clustered well)

Silhouette analysis for KMeans clustering on sample data with n_clusters = 2

The silhouette plot for the various clusters.

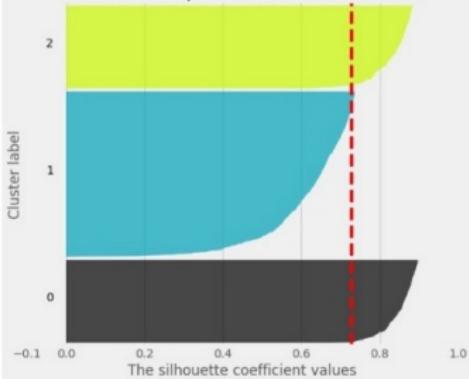


The visualization of the clustered data.

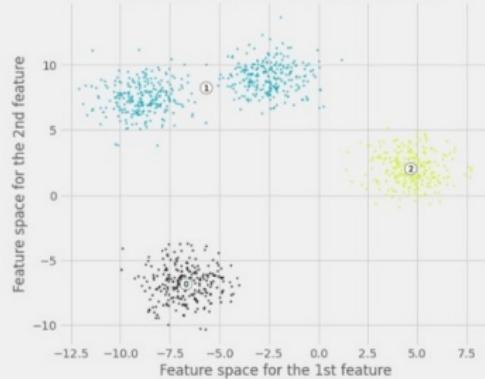


Silhouette analysis for KMeans clustering on sample data with n_clusters = 3

The silhouette plot for the various clusters.

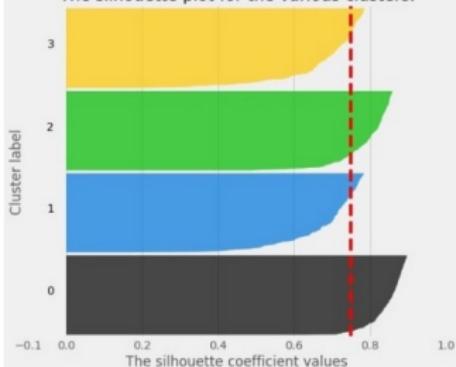


The visualization of the clustered data.

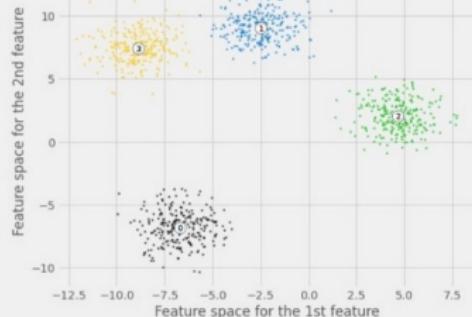


Silhouette analysis for KMeans clustering on sample data with n_clusters = 4

The silhouette plot for the various clusters.

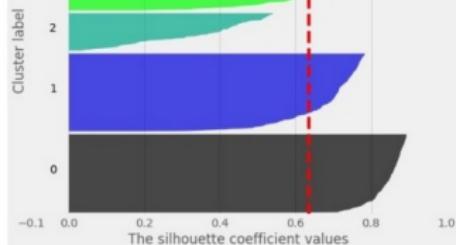


The visualization of the clustered data.

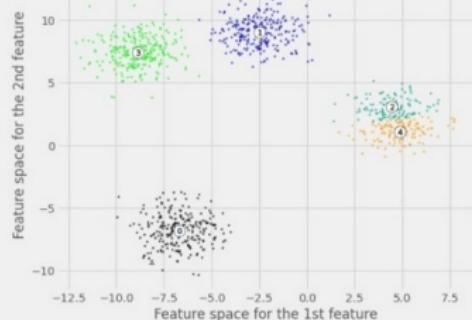


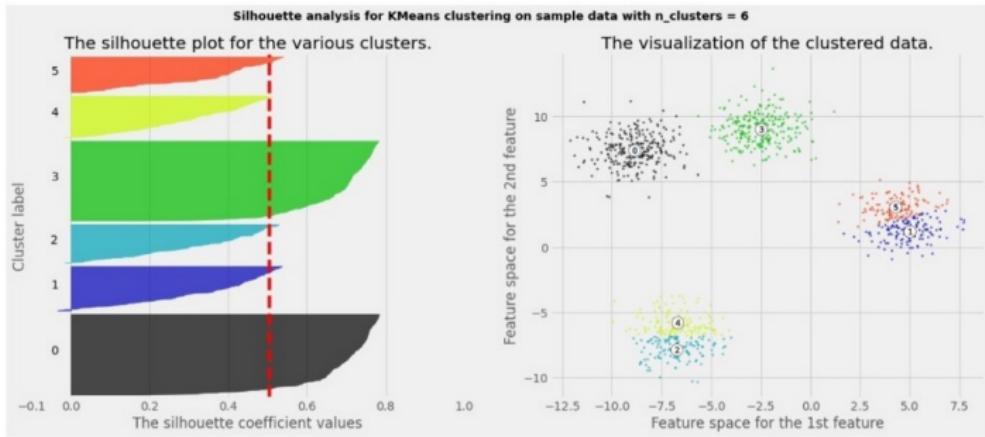
Silhouette analysis for KMeans clustering on sample data with n_clusters = 5

The silhouette plot for the various clusters.



The visualization of the clustered data.





Computational considerations

- The straightforward computation of s_C requires $\Theta(|P|^2)$ distance computations, unfeasible for very large $|P|$.
- Simplified computation for cosine and squared Euclidean distances ($\Theta(|P|k)$ distance computations). Methods are provided in the Spark library.
- There exist heuristics to approximate s_C , some with provable guarantees (from Padova!).

Approximating the sum of distances

Consider a point x and a set C . The following method can be used to approximate $d_{\text{sum}}(x, C)$. Let $n = |C|$.

- ① Fix a target sample size $t \in [1, n]$. *(Typically, $t \ll n$)*
- ② Select a sample $S_t \subset C$, where each $y \in C$ is independently included in S_t with probability t/n (Poisson sampling). $E[S_t] = t$
- ③ Compute the approximation

$$\tilde{d}_{\text{sum}}(x, C, t) = \frac{n}{t} \sum_{y \in S_t} d(x, y)$$

Remark. $\tilde{d}_{\text{sum}}(x, C, t)$ is a random variable, whose value depends on the random sample S_t .

Proposition

$\tilde{d}_{\text{sum}}(x, C, t)$ is an unbiased estimator of $d_{\text{sum}}(x, C)$, i.e.,

$$E[\tilde{d}_{\text{sum}}(x, C, t)] = d_{\text{sum}}(x, C).$$

Proof of Proposition

For each $y \in C$ define a random variable $V_x(y)$

$$V_x(y) = \begin{cases} d(x, y) & \text{if } y \in S_t \text{ with prob. } t/n \\ 0 & \text{if } y \notin S_t \text{ with prob } 1 - t/n \end{cases}$$

$$\begin{aligned}\tilde{d}_{sum}(x, C, t) &= \frac{n}{t} \sum_{y \in S_t} d(x, y) \\ &= \frac{n}{t} \sum_{y \in C} V_x(y)\end{aligned}$$

Proof of Proposition

$$\begin{aligned}
 E[\tilde{d}_{\text{sum}}(s, c, t)] &= E\left[\frac{n}{t} \sum_{y \in C} V_x(y)\right] \\
 &= \frac{n}{t} \sum_{y \in C} E[V_x(y)]
 \end{aligned}$$

$$E[V_x(y)] = \frac{t}{n} d(x, y)$$

\Rightarrow

$$\begin{aligned}
 E[\tilde{d}_{\text{sum}}(x, c, t)] &= \frac{n}{t} \cdot \sum_{y \in C} \frac{t}{n} d(x, y) \\
 &= \sum_{y \in C} d(x, y) = d_{\text{sum}}(x, c)
 \end{aligned}$$
□

Approximating the sum of distances

Observations:

- The absolute error of the estimate $|\tilde{d}_{\text{sum}}(x, C, t) - d_{\text{sum}}(x, C)|$ can be upper bounded analytically as a function of t (see [EW04]).
- While the theoretical bound is somewhat weak, in practice the error is rather low even with small t .

Exercise

Let P be a pointset of N points from a metric space (M, d) and let $C \subseteq P$ be a subset of size n . Show how to compute $\tilde{d}_{\text{sum}}(x, C, t)$ efficiently in MapReduce, for all $x \in P$. Assume that t and n are known and

- $t \in O(\log N)$ but n can be large (up to N);
- each point $x \in P$ is represented by a pair $(ID_x, (x, f_x))$, where ID_x is a distinct integer in $[0, N - 1]$ and f_x is a binary flag with value 1 if $x \in C$, and 0 otherwise.

Approximating the Silhouette

Let $\mathcal{C} = (C_1, C_2, \dots, C_k)$ be a partition of P into k clusters.

Fix a sample size t and define $t_i = \min\{t, |C_i|\}$, for $1 \leq i \leq k$.

For $1 \leq i \leq k$ and $x \in C_i$, define

$$\tilde{a}_x = \tilde{d}_{\text{sum}}(x, C_i, t_i) / |C_i| \quad \text{and} \quad \tilde{b}_x = \min_{j \neq i} \tilde{d}_{\text{sum}}(x, C_j, t_j) / |C_j|,$$

The approximate silhouette coefficient of x is

$$\tilde{s}_x = \frac{\tilde{b}_x - \tilde{a}_x}{\max\{\tilde{a}_x, \tilde{b}_x\}},$$

The approximate average silhouette coefficient of \mathcal{C} is

$$\tilde{s}_{\mathcal{C}} = (1/|P|) \sum_{x \in P} \tilde{s}_x.$$

IMPORTANT: For cluster C_i compute only the sample s_i using parameter t_i

Exercise

Determine the expected number of distance computations needed to obtain \tilde{s}_c .

$\forall x \in P$ and $\forall i \quad 1 \leq i \leq k$

We must compute all distances between x and the sample s_i of C_i

\Rightarrow The total number of distance computations is

$$N \cdot \sum_{i=1}^k |S_i|$$

$$\begin{aligned}
 E\left[N \cdot \sum_{i=1}^K |S_i|\right] &= N \cdot \sum_{i=1}^K E[|S_i|] \\
 &= N \cdot \sum_{i=1}^K \frac{t_i}{|C_i|} \cdot |C_i| \\
 &= N \sum_{i=1}^K \min\{t_i, |C_i|\} \\
 &\leq N \cdot K \cdot t = \sigma(N^2)
 \end{aligned}$$

↓
 when K, t are
 sufficiently small

Other evaluation metrics

- **Supervised evaluation:** assessment of the quality of a clustering with reference to external information. For example, Entropy can be used to measure the coherence of a clustering with respect to a ground truth (e.g., class labels).
- **Clustering tendency:** assessment whether the data contain meaningful clusters, namely clusters that are unlikely to occur in random data. The Hopkins statistic can be used to this purpose.

Clustering evaluation is well described in [TSK06]

Summary of Parts 1, 2 and 3

- Metric space and prominent distance functions.
- Combinatorial optimization problem and approximation algorithm.
- k-center clustering:
 - Definition of the problem.
 - Farthest-First Traversal algorithm.
 - (Composable) coreset technique.
 - MR-Farthest-First Traversal.
- k-means clustering:
 - Definition of the problem.
 - Review of Lloyd's and k-means++ algorithms.
 - Weighted variant of the problem.
 - MR-kmeans
- k-median clustering (sketch)

Summary of Parts 1, 2 and 3

- Unsupervised clustering evaluation:
 - Silhouette coefficient.
 - Approximating the sum of distances.
 - Approximating the Silhouette.
 - Other evaluation metrics (sketch).

References

- LRU14 J. Leskovec, A. Rajaraman and J. Ullman. Mining Massive Datasets. Cambridge University Press, 2014. Sections 3.1.1 and 3.5, and Chapter 7
- BHK18 A. Blum, J. Hopcroft, and R. Kannan. Foundations of Data Science. Manuscript, June 2018. Chapter 7
- AV07 D. Arthur, S. Vassilvitskii. k-means++: the advantages of careful seeding. Proc. of ACM-SIAM SODA 2007: 1027-1035.
- B+12 B. Bahmani et al. Scalable K-Means++. PVLDB 5(7):622-633, 2012
- C+19 V. Cohen-Addad et al. Tight FPT Approximations for k-Median and k-Means. Proc. of ICALP 2019: 42:1-42:14.
- Wei16 D. Wei A Constant-Factor Bi-Criteria Approximation Guarantee for k-means++. Proc. of NIPS 2016: 604-612.
- TSK06 P.N.Tan, M.Steinbach, V.Kumar. Introduction to Data Mining. Addison Wesley, 2006. Chapter 8.
- EW04 D. Eppstein, J. Wang: Fast Approximation of Centrality. J. Graph Algorithms Appl. 8: 39-45 (2004).