

Landmark FN-DBSCAN: An Efficient Density-based Clustering Algorithm with Fuzzy Neighborhood*

Hao Liu, Satoshi Oyama, Masahito Kurihara, Haruhiko Sato
(Graduate School of Information Science and Technology Hokkaido University)

1. Introduction and Summary

Among various clustering techniques, the density-based clustering algorithms have several advantages as follows:

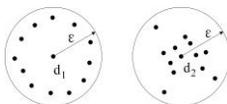
1. The number of clusters is not required before carrying out clustering.
2. The detected clusters can be represented in an arbitrary shape.
3. Outliers can be detected and removed.

However, in general, the parameters of density-based clustering algorithms are usually difficult to select. So, in order to make the density-based clustering algorithms more robust, the extension with fuzzy set theory has attracted a lot of attentions recently. The fuzzy neighborhood DBSCAN (FN-DBSCAN) is a typical one with this idea. But FN-DBSCAN usually requires a time complexity of $O(n^2)$ where n is the number of data in the data set. This implies that the algorithm is not suitable for the work with large scale data sets.

In this study, we proposed a novel clustering algorithm called landmark fuzzy neighborhood DBSCAN (landmark FN-DBSCAN), which can provide a similar result to FN-DBSCAN, but only requires a linear time and space complexity to the size of input data set.

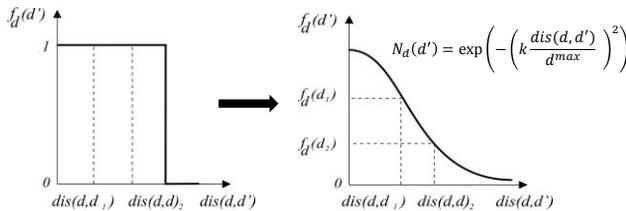
2. Related Work

- 1) DBSCAN (Using a crisp neighborhood function):

$$N_d(d') = \begin{cases} 1, & \text{if } \text{dis}(d, d') \leq \varepsilon \\ 0, & \text{otherwise} \end{cases}$$


In the right figure, data d_1 and d_2 have the same cardinality value which is 12, but the density values are not the same.

- 2) FN-DBSCAN (using a fuzzy neighborhood function)



By applying the fuzzy neighborhood function, FN-DBSCAN can provide different cardinality values for d_1 and d_2 in the above mentioned figure.

3. The Proposed Method: Landmark FN-DBSCAN

- 1) Landmark Concept

Given a data set $D(D = [d_{ij}]_{n \times m})$, a landmark, l , is a triplet:

$$l = \langle V, N_f^l(l), \mu \rangle$$

- V is a m -dimensional vector referred as a data vector in D .
- $N_f^l(l) = \{d \in D | f_l(d) \geq \varepsilon_1\}$, called fuzzy neighborhood.
- μ is a positive real number, called the membership level.

- 2) Measuring membership degrees

• Landmark & data:

$$f_l(d) = \exp\left(-\left(r \cdot k \cdot \frac{\text{dis}(V, d)}{\Delta d^{\text{max}}}\right)^2\right)$$

• Landmark & landmark:

$$f_{l_1}(l_2) = \exp\left(-\left(k \cdot \frac{\text{dis}(V_1, V_2)}{\Delta d^{\text{max}}}\right)^2\right)$$

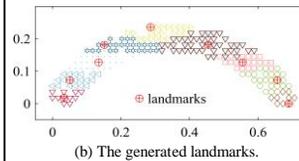
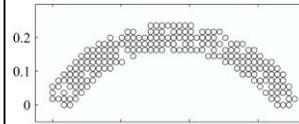
• Membership level:

$$\mu = \sum_{d \in N_f^l(l)} f_l(d)$$

3) The landmark FN-DBSCAN Algorithm

The algorithm consists of three steps:

- I. Divide a data set into several subsets represented by generated "landmarks" (Algorithm 1).
- II. Execute a modified version of FN-DBSCAN on the generated landmark set and output the index of landmarks.
- III. Label data by the index of landmarks.



Algorithm 1 LandmarkGeneration

Input: D, r, k, ε_1
Output: L

- 1: $L \leftarrow \emptyset$;
- 2: for all x in D do
- 3: find a landmark $l = (V, N, u) \in L$, such that $l, V = \min\{\text{dis}(l, V, x)\}$;
- 4: $u \leftarrow \exp\left(-\left(r \cdot k \cdot \frac{\text{dis}(l, V, x)}{\Delta d^{\text{max}}}\right)^2\right)$;
- 5: if $L = \emptyset$ or $u < \varepsilon_1$ then
- 6: $V \leftarrow x$; $N \leftarrow \{x\}$; $u \leftarrow 1$;
- 7: $l \leftarrow (V, N, u)$;
- 8: $L \leftarrow L \cup \{l\}$;
- 9: else
- 10: $l, N \leftarrow l, N \cup \{x\}$;
- 11: $l, u \leftarrow l, u + u$;
- 12: end if
- 13: end for

An Example of generated landmarks applying Algorithm 1.

- 4) Complexity Analysis

a) Time complexity:

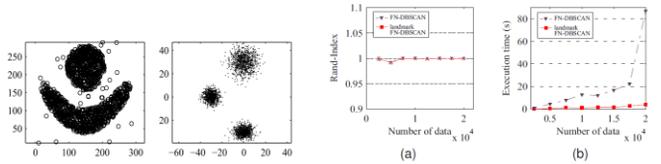
$$O(kn + k^2 + n) = O(kn + k^2) \xrightarrow{k \ll n} O(n)$$

b) Space complexity:

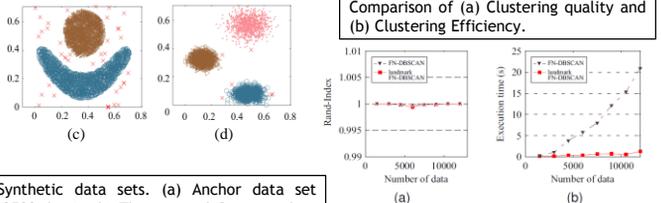
$$O(n + k + n + k) = O(n + k) \xrightarrow{k \ll n} O(n)$$

where n is the number of data and k is the number of the generated landmarks.

4. Experimental Results:



Results of Anchor data set ($r = 3$). Comparison of (a) Clustering quality and (b) Clustering Efficiency.



Synthetic data sets. (a) Anchor data set (2500 data). (b) Three-mixed Gaussian data set (3000 data). (c) The result for anchor data set. (d) The result for three-mixed Gaussian data set.

Results of Banana data set ($r = 3$). Comparison of (a) Clustering quality and (b) Clustering Efficiency.

Gaussian Data Set (3 Clusters).

Size	Landmark FN-DBSCAN Rand-Index					FN-DBSCAN Rand-Index	Landmark FN-DBSCAN execution (time/s)					FN-DBSCAN execution (time/s)
	r value	1.5	1.8	2.0	2.5		3.0	r value	1.5	1.8	2.0	
3000	0.9953	0.9976	0.9980	0.9982	0.9987	0.9991	0.13	0.13	0.16	0.22	0.27	2.00
4000	0.9869	0.9982	0.9985	0.9990	0.9994	0.9996	0.20	0.25	0.30	0.36	0.45	6.08
9000	0.9921	0.9981	0.9814	0.9999	0.9999	0.9999	0.20	0.20	0.23	0.28	0.55	12.34
12000	0.9926	0.9928	0.9969	0.9962	0.9995	0.9996	0.25	0.28	0.31	0.36	0.44	20.55
15000	0.9983	0.9984	0.9995	0.9995	0.9995	0.9996	0.34	0.36	0.41	0.45	0.58	24.05
18000	0.9924	0.9990	0.9994	0.9996	0.9996	0.9997	0.39	0.44	0.45	0.56	0.67	38.56
21000	0.9976	0.9995	0.9996	0.9995	0.9998	0.9998	0.45	0.52	0.56	0.67	0.80	47.28
24000	0.9955	0.9998	0.9998	0.9999	0.9999	0.9999	0.53	0.59	0.66	0.75	0.92	127.11
27000	0.9943	0.9973	1.0000	1.0000	1.0000	1.0000	0.58	0.66	0.70	0.88	1.06	137.14
30000	0.9992	0.9998	0.9998	0.9998	0.9999	0.9999	0.64	0.72	0.81	0.97	1.14	218.30

References: Hao Liu, Satoshi Oyama, Masahito Kurihara, Haruhiko Sato. "Landmark FN-DBSCAN: An Efficient Density-based Clustering Algorithm with Fuzzy Neighborhood," Journal of Advanced Computational Intelligence and Intelligent Informatics, Vol.17 No.1, pp: 60-73, 2013.