

SUBSCALE Algorithmus

William Mendat,¹ Max Ernst,² Steven Schall,³ Matthias Reichenbach⁴

Abstract: Das könnte das abstract sein

Keywords: C++; Subscale; Cluster

1 Überschrift/Heading

Hallo [?]

2 Überschrift/Heading

Hallo [?]

3 Berechnung der Dense Units

Clusters von benachbarten Punkten müssen zu einem minPoint Großen subsets kombiniert werden. Diese subsets werden Dense Units genannt. Dense Units werden in verschiedene Funktionen mit unterschiedlichen Zwecken berechnet. Das Hauptziel dieser Berechnung, ist die Bestimmung einer Teilmenge aus allen möglichen Kombinationen in einem subset. Dabei dürfen keine Wiederholungen der Dense Units auftreten. Die Kombinationen aus einem subset können mittels dem Binomialkoeffizient berechnen, $\binom{n}{k}$ dabei ist n die Anzahl der Elemente in dem subset und k die minimale Anzahl an punkten in einem subset (e.g. minPoints). Daraus entstehen k-Elemente großes subsets von einem n-Elemente set ohne Wiederholungen der Kombinationen. Wenn zum Beispiel ein Core Set aus folgenden Punkten besteht: [1, 5, 7, 9, 22] und die minimale Anzahl an punkten in einem Subset Drei ist, sind die ersten Drei Dense Units folgende: [1, 5, 7] und [1, 5, 9] und [1, 5, 22]. Die Formel zur Berechnung der Dense Units kann also folgendermaßen betrachtet werden: $\binom{|CS|}{minPoints}$ [KD14] Abbildung 1 zeigt ein weiteres Beispiel für n-Punkte in einem Core Set mit minPoints von 4.

¹ Hochschule Offenburg, Offenburg, Deutschland w mendat@stud-hs.offenburg.de

² Hochschule Offenburg, Offenburg, Deutschland w mendat@stud-hs.offenburg.de

³ Hochschule Offenburg, Offenburg, Deutschland s schall@stud-hs.offenburg.de

⁴ Hochschule Offenburg, Offenburg, Deutschland m reichen@stud-hs.offenburg.de

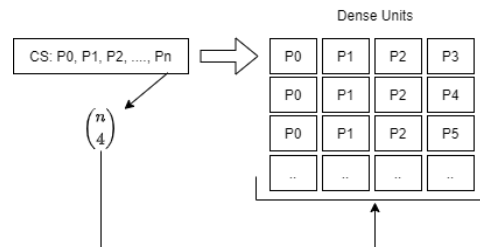


Abb. 1: Berechnung der Dense Units

4 Kollision von Dense Units

Da die Dense Units auch in unterschiedlichen Dimensionen existieren können, müssen diese als solche gekennzeichnet werden. In diesem Schritt ist das subspace clustering erkennbar. Direkt nach Berechnung der Dense Units, werden diese mit anderen Dense Units verglichen um Höhere subspaces zu finden. Um Kollisionen in den Dense Units festzustellen, muss als erstens jedem Punkt in dem Datenset eine hohe Zufallszahl zugewiesen werden. Dieser schritt wird vor dem ausführen des SUBSCALE Algorithmus durchgeführt. Für alle Dense Units, werden die Signaturen gebildet, diese Signatur ist die Summe aller Punkte in der Dense Unit. Anhand dieser Signatur werden die Dense Units in eine Tabelle eingetragen, mit den dazugehörigen Punkten und der Dimension. Bei einer Kollision der Signaturen, wird die Dimension zu dem bereits vorhanden Eintrag hinzugefügt (siehe Abbildung 2). [Ra]

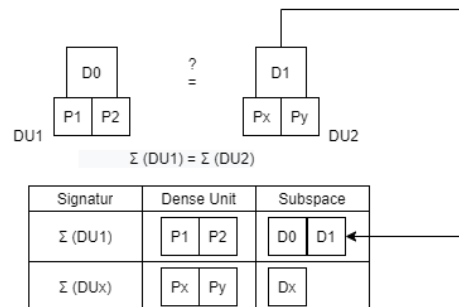


Abb. 2: Kollisionsauflösung von Dense Units

5 Abbildung Dense-Units auf Subspaces

Für das endgültige Clustering muss die Ausgabe des SUBSCALE-Algorithmus in die Struktur von den Clusterkandidaten abgebildet werden. Jeder Kandidat besteht aus Dimensionen und eindeutigen Punkt IDs. Dabei werden Punkte ausgesucht, welche in mehreren Dense-Units vertreten sind. Diese werden dann zusammen innerhalb eines Subspaces definiert, wobei

mehrfach vorkommende Punkte nur einmal eingetragen werden. Die hier zusammengeführten Punkte haben die Eigenschaft, dass sie ziemlich wahrscheinlich auch in dem jeweiligen Unterraum geclustert sind. Diese Eigenschaft macht diese Punkte zu günstigen Kandidaten für das endgültige Clustering, welches im folgenden Unterkapitel 6 näher beschrieben wird.

6 Abschließendes Clustering mit DBSCAN

Nach dem alle maximalen Subspaces identifiziert worden sind, wird zuletzt das Clustering durchgeführt, welches die maximalen Supspace Cluster finden soll. Zur Bewältigung dieser Aufgabe wird ein volldimensionaler Clustering-Algorithmus verwendet. Der Algorithmus, welcher in der von uns verwendeten Implementierung verwendet wird, nennt sich DBSCAN (Density-Based Spatial Clustering of Applications with Noise). DBSCAN ist ein auf die Dichteverbundenheit basierender Clustering-Algorithmus, der Cluster mit beliebiger Form findet und Rauschpunkte separat zurückliefert.

7 Verteilungsmöglichkeiten des SUBSCALE Algorithmus

[PLK21]

Literaturverzeichnis

- [KD14] Kaur, Amardeep; Datta, Amitava: SUBSCALE: Fast and Scalable Subspace Clustering for High Dimensional Data. In: 2014 IEEE International Conference on Data Mining Workshop. S. 621–628, 2014.
- [PLK21] Prinzbach, Jürgen; Lauer, Tobias; Kiefer, Nicolas: Accelerating Density-Based Subspace Clustering in High-Dimensional Data. In: 2021 International Conference on Data Mining Workshops (ICDMW). S. 474–481, 2021.
- [Ra] Ramin, Stanislav: Python Implementation of the SUBSCALE Algorithm. Bachelor thesis, Offenburg.