

Cookies erleichtern die Bereitstellung unserer Dienste. Mit der Nutzung unserer Dienste erklären Sie sich damit einverstanden, dass wir Cookies verwenden. [Weitere Informationen](#) [OK](#)

[INFRASTRUKTUR](#) [DATA SOURCING](#) [ANALYTICS](#) [BEST PRACTICES](#) [INDUSTRIE 4.0](#) [RECHT & SICHERHEIT](#)

Sie befinden sich hier: [Definitionen](#)



Definition

Was ist der k-Means-Algorithmus?

18.07.18 | Autor / Redakteur: [Stefan Lubert](#) / [Nico Litzel](#)

(Bild: © aga7ta -
stock.adobe.com)

Der k-Means-Algorithmus ist ein Rechenverfahren, das sich für die Gruppierung von Objekten, die sogenannte Clusteranalyse, einsetzen lässt. Dank der effizienten Berechnung der Clusterzentren und dem geringen Speicherbedarf eignet sich der Algorithmus für die Analyse großer Datenmengen, wie sie im Big-Data-Umfeld üblich sind.

Der k-Means-Algorithmus ist eines der am häufigsten verwendeten mathematischen Verfahren zur Gruppierung von Objekten (Clusteranalyse). Der Algorithmus ist in der Lage, eine große Menge ähnlicher Objekte mit einer vorher bekannten Anzahl von Gruppen in Cluster zu unterteilen. Da es sich um ein sehr effizientes Verfahren handelt, das mit verschiedenen Datentypen zurecht kommt, und der Speicherbedarf gering ist, eignet sich der k-Means-Algorithmus für die Datenanalyse im Big-Data-Umfeld.

Die Laufzeit des Algorithmus ist linear zur Anzahl der vorhandenen Datenpunkte. Die Anzahl der Schleifendurchläufe zur Ermittlung der Clusterzentren ist klein. Die Idee für den k-Means-Algorithmus stammt ursprünglich von Hugo Steinhaus aus dem Jahr 1957. Erst 1983 wurde der Algorithmus in einer Informatik-Zeitschrift veröffentlicht. Es besteht eine große Ähnlichkeit zwischen dem k-Means-Algorithmus und dem sogenannten Expectation-Maximization-Algorithmus. Für den k-Means-Algorithmus ist die Anzahl der Clusterzentren vorab festzulegen.

existieren verschiedene Erweiterungen wie k-Means++ oder der k-Median-Algorithmus.

Ziele der Clusteranalyse mit dem k-Means-Algorithmus

Die Clusteranalyse ist ein Verfahren, mit dem sich eine bestimmte Anzahl von Objekten in homogene Gruppen einteilen lässt. Ziel ist es, dass die verschiedenen Objekte in einer Gruppe sich nach der erfolgten Einteilung möglichst ähnlich sind. Die Eigenschaft, dass Objekte in Dimensionen eingeteilt sind, nennen sich Cluster. Der k-Means-Algorithmus lässt sich für mehrdimensionale Objekte anwenden und nähert sich durch ständig wiederholende Neuberechnungen den jeweiligen Clusterzentren an, bis keine signifikante Veränderung mehr stattfindet.

Typischer Ablauf des k-Means-Algorithmus

Der k-Means-Algorithmus findet Anwendung auf Daten oder Objekte in einem n-dimensionalen Raum. Das Verfahren verläuft in folgenden Schritten:

1. Wahl von k-Punkten als Anfangszentren der Berechnung
2. Zuordnung der Datenpunkte zu den verschiedenen Clustern auf Basis des Abstandes zu den Zentren
3. Neuberechnung der Clusterzentren
4. Wiederholung ab Schritt 2 – bis sich die Lage der Zentren nicht mehr ändert

In den ersten Durchläufen des Algorithmus treten noch große Änderungen der Zentren auf. In zunehmenden Schleifendurchläufen werden die Veränderungen immer kleiner. Wenn ein effizienter Ablauf des Algorithmus ist die Wahl der Anfangszentren.

Probleme bei der Anwendung des k-Means-Verfahrens

Ein Problem bei der Anwendung des k-Means-Algorithmus stellt die Vorgabe der Anzahl der Cluster und die Wahl der Anfangszentren dar. Die gefundene Lösung ist stark von den gewählten Startpunkten abhängig. Die Wahl der Startpunkte ist die Kenntnis einer Clusterstruktur notwendig, die aus Voranalysen der Daten stammen könnte. In der Praxis erfolgt der Durchlauf des Algorithmus mit unterschiedlichen Startwerten. Die verschiedenen Ergebnisse liefern Hinweise auf eine möglichst plausible Struktur der Cluster. Wenn eine ungeeignete Anzahl von Clusterzentren als Startwerte, können sich unter Umständen komplett andere Lösungen oder ungeeignete Clustereinteilungen ergeben.

Auch problematisch für den k-Means-Algorithmus sind Datenmengen, die sich überlappend in Teilen nahtlos ineinander übergehen. In diesen Fällen ist das k-Means-Verfahren in der Lage, die verschiedenen Gruppen zuverlässig voneinander zu trennen. Daten mit überlappenden Clusterstrukturen werden ebenfalls nicht unterstützt. Sind Ausreißer in den Daten

vorhanden, können diese das Ergebnis stark verfälschen, da k-Means keine Ausreißer und jedes Objekt einem Cluster zuordnet. In diesen Fällen ist vor der Auswertung des k-Means-Algorithmus eine Bereinigung der Daten (Noisereduktion) durchzuführen.

Die Erweiterung k-Means++

Für den k-Means-Algorithmus existiert die Erweiterung k-Means++. Sie beschreibt ein Verfahren, mit dem die Clusterzentren als Startwerte nicht mehr zufällig, sondern nach einer Vorschrift gewählt werden. Sind die Clusterzentren als Startwerte bestimmt, konvergiert der anschließend ausgeführte k-Means-Algorithmus sehr schnell und in wenigen Schleifendurchläufen. Typischerweise lässt sich eine Verdopplung der Geschwindigkeit erreichen.

Anwendungen des k-Means-Algorithmus

Der k-Means-Algorithmus findet in vielen Bereichen Anwendung. Aufgrund seiner Einfachheit und des geringen Speicher- und Rechenbedarfs eignet er sich für die Datenanalyse großer Datenmengen im Big-Data-Umfeld. In der Bildverarbeitung wird k-Means häufig zur Segmentierung der Bilddaten eingesetzt. Mit den Ergebnissen des Algorithmus ist beispielsweise die Trennung von Vorder- und Hintergrund oder das Erkennen von Objekten möglich.

Ein weiterer wichtiger Anwendungsbereich sind das Marketing und die Analyse des Kundenverhaltens. k-Means findet in vorliegenden Kundendaten verschiedene Gruppen von Kunden mit jeweils ähnlichem Kundenverhalten. Die von k-Means ermittelten Homogenen Gruppen lassen sich klar voneinander trennen. Mithilfe spezifischer Marketingmaßnahmen können die Gruppen effizienter und mit größerer Erfolgsaussicht ansprechbar.



Über den Autor

Stefan Luber

Definition

Was ist Customer Experience?

KOMMENTAR ZU DIESEM ARTIKEL ABGEBEN

ANONYM MITDISKUTIEREN ODER EINLOGGEN **ANMELDEN**

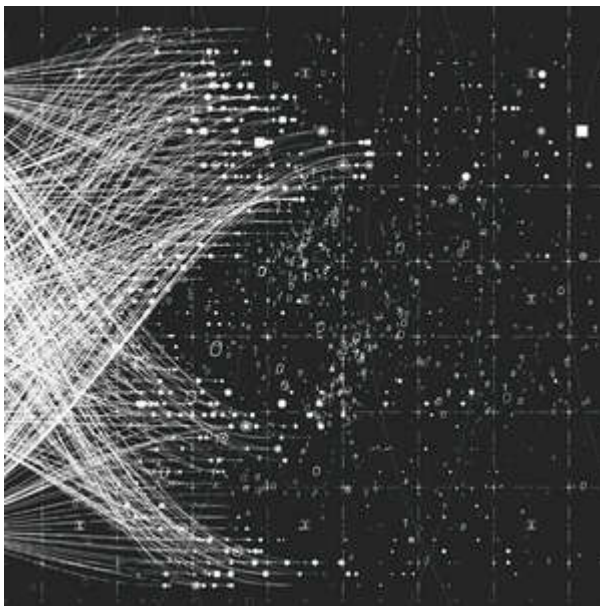


Name eingeben...

Zur Wahrung unserer Interessen speichern wir zusätzlich zu den o.g. Informationen die IP-Adresse. Dies dient ausschließlich dem Zweck, dass Sie als Urheber des Kommentars identifiziert werden können. Rechtliche Grundlage ist die Wahrung gemäß Art. 6 Abs. 1 lit. f) DSGVO.

Kommentieren

AKTUELLE BEITRÄGE ZU DIESEM THEMA

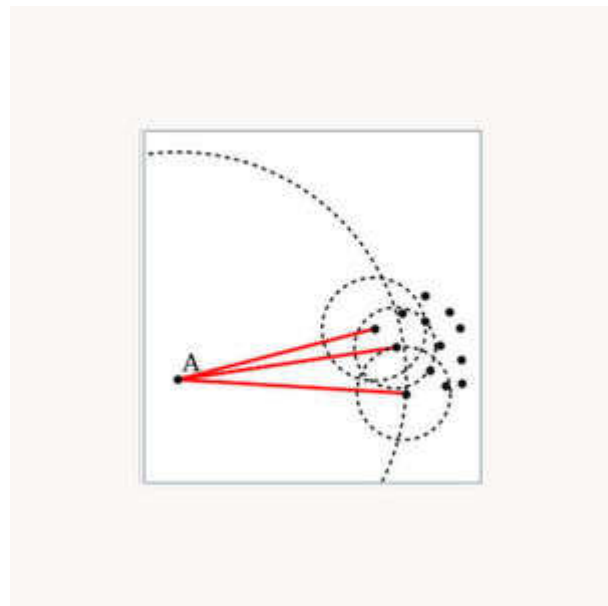


[Algorithmen](#)

[Analytics-Methoden – von deskriptiven Analysen bis Machine-Learning-Algorithmen](#)

Dem Datenanalysten stehen zahlreiche Methoden zur Verfügung. Der folgende Artikel erläutert einige Methoden – von statistischen, deskriptiven Methoden bis zu Supervised und Unsupervised Machine Learning.

[lesen](#)



[Grundlagen Statistik & Algorithmen, Teil 7](#)

[So deckt der Local Outlier Factor Anomalien auf](#)

Um Trends zu erkennen, wird oft die Clusteranalyse herangezogen. Der k-Means-Algorithmus etwa zeigt an, wo sich Analyseergebnisse in einer Normalverteilung ballen. Für manche Zwecke ist es aber aufschlussreicher, Ausreißer zu untersuchen, denn sie bilden die Antithese zum „Normalen“, etwa im Betrugswesen. Der Local-Outlier-Factor-Algorithmus (LOF) ist in der Lage, den Abstand von Ausreißern zu ihren Nachbarn zu berechnen und deckt so Anomalien auf.

[lesen](#)



[Grundlagen Statistik & Algorithmen, Teil 7](#)

Der sich Clu: Ber: Spe die Um: in d

[lesen](#)



Machine Learning

So kommen Sie mit den passenden Algorithmen zum Ziel

Einfache Algorithmen sind das tägliche Brot vieler Programmierer. Wer etwa eine Software erstellt, die den Preis eines Hauses basierend auf der Größe ausgibt, schreibt normalerweise einen Algorithmus, der abhängig vom Input (Hausgröße) einen bestimmten Output (Preis) berechnet. Wenn es um Machine Learning geht, steigen die Anforderungen an die Programmierung und die Mathematik.

lesen



Kommentar von Lars Milde, Tableau

Top 10 der Business Intelligence Trends für das Jahr 2017

2016 lagen Self-Service-Analysen im Trend. Viele Unternehmen führten den modernen Business-Analytics-Ansatz ein, bei dem IT und Geschäftsbetrieb zusammenarbeiten, um die eigenen Daten optimal zu nutzen. Die IT begann, skalierbare und ausbaufähige Technologien zu nutzen, Geschäftsanwender teilten ihre Daten und arbeiteten gemeinsam daran.

lesen

Lei

Ta

Be

Vor

10.

als

unt

und

les

MEHR ZUM THEMA



Algorithmen

Analytics-Methoden – von deskriptiven Analysen bis Machine-Learning-Algorithmen

mehr...



Data Science

IoT-Basics – die technische Basis von Big Data

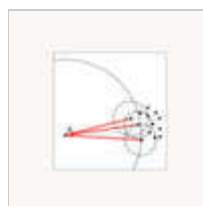
mehr...



Definition

Was ist BigTable?

mehr...



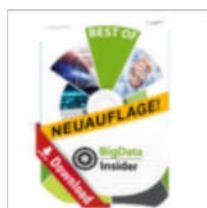
Grundlagen Statistik & Algorithmen, Teil 7

So deckt der Local Outlier Factor Anomalien auf

mehr...



PASSENDE WHITEPAPER & WEBCASTS



[Alle Whitepaper](#)

Diese aktuellen Themen bewegen die Big-Data-Welt

Das BEST OF BigData-Insider

[mehr...](#)

[Alle Webcasts](#)



Technology-Update für IT-Manager

CIOBRIEFING 12/2018

[mehr...](#)



Dieser Beitrag ist urheberrechtlich geschützt. Sie wollen ihn für Ihre Zwecke verwenden? Kontaktieren Sie uns über: [support.vogel](mailto:support.vogel@vogel.com)



Oder kontaktieren Sie uns [direkt](#)

BigData-Insider ist eine Marke der Vogel Communications Group. Unser gesamtes Angebot finden Sie [hier](#)

[AGB](#) | [EWG](#) | [Hilfe](#) | [Kundencenter](#) | [Media](#) | [Datenschutz](#) | [Impressum](#)
Copyright © 2020 Vogel Communications Group

©garrykillian – stock.adobe.com; gemeinfrei; Kernel Machine.svg / Alisneaky, svg version by User:Zirguezi / CC BY-SA 4.0; © Dmitry Nikolaev - stock.adobe.com; Tableau;