

Masterarbeit

Erkennung von Fahrspuren mittels Fahrzeugtrajektorien aus Luftaufnahmen

im Master-Studiengang Informatik
der Hochschule Furtwangen

Steffen Schmid

Zeitraum: Wintersemester 2018

Prüfer: Prof. Dr. Christoph Reich

Zweitprüfer: Stefan Kaufmann

Firma: IT-Designers GmbH

Betreuer: Stefan Kaufmann

Eidesstattliche Erklärung

Hiermit versichere ich, die vorliegende Arbeit selbstständig und unter ausschließlicher Verwendung der angegebenen Literatur und Hilfsmittel erstellt zu haben.

Die Arbeit wurde bisher in gleicher oder ähnlicher Form keiner anderen Prüfungsbehörde vorgelegt und auch nicht veröffentlicht.

Esslingen, den 30. Oktober 2018

Unterschrift

Zitat

„Some fancy quote“

FooBar Muman

Danksagung

Kurzfassung

Schlagworte:

Inhaltsverzeichnis

Kurzfassung	v
1 Einleitung	1
1.1 Rahmen der Arbeit	1
1.1.1 Das Projekt MEC-View	1
1.1.2 MEC-View Luftbeobachtung	1
1.2 Motivation und Ziele	1
1.3 Aufbau dieser Arbeit	1
2 Grundlagen	2
2.1 Verkehrsanalyse mittels Luftaufnahmen	2
2.2 Rekonstruktion von Fahrzeugtrajektorien aus Luftaufnahmen	2
2.3 Datenaufbereitung und Bereinigung	2
2.4 Clusteranalyse	2
2.4.1 Eigenschaften von Cluster-Sets und Clustern	4
2.4.2 Cluster-Algorithmen	6
2.5 Ähnlichkeitsmaße zum Vergleich von Fahrzeugtrajektorien	12
2.5.1 HU Distanz	13
2.5.2 Hausdorff Distanz	14
2.5.3 Longest-Common-Subsequence	15
2.6 Untersuchung möglicher Straßentopologien	16
3 Stand der Technik	17
3.1 Clustering von Trajektoriedaten	17
3.2 Erkennung von Fahrspuren	17
3.3 Klassifizierung von Fahrspuren	17
3.4 Defizite vorhandener Lösungen und benötigte Neuerungen	17
4 Clustering von Fahrzeugtrajektorien	18
4.1 Vorverarbeitung der Roh-Trajektorien	18
4.2 Gruppierung der Trajektorien	18
5 Fahrbahn-Bestimmung aus Trajektorie-Clustern	19
6 Fahrbahn Klassifizierung	20
7 Realisierung LaneDetection in MEC-View TrackerApplication Software	21

8	Ergebnisse und Auswertung	22
9	Zusammenfassung und Fazit	23
	Literaturverzeichnis	24

Abbildungsverzeichnis

2.1	Rohdaten (links) und erwünschtes Clustering-Ergebnis (rechts)	3
2.2	Ablauf einer Clusteranalyse	4
2.3	Agglomeratives Clustering dargestellt als Dendrogramm und geschachteltes Cluster-Diagramm	7
2.4	Funktionsweise von k-Means	8
2.5	Darstellung des EM-Cluster-Algorithmus über mehrere Iterationen	10
2.6	Schritte des DBSCAN Algorithmus	11
2.7	Erreichbarkeit in DBSCAN	12
2.8	Trajektorien im 2-dimensionalen Raum	13

Listings

1 Einleitung

Staus und zäh fließender Verkehr sind sowohl auf Schnell- und Autobahnen, als auch in Städten ein großes Problem und Ärgerniss für Autofahrer. Sie kosten diese nicht nur wertvolle Zeit, sondern auch viel Geld. Laut einer Studie von [Cookson u. a.] kostet Stau jeden deutschen Autofahrer pro Jahr durchschnittlich 1770 €. In Summe ergeben sich hieraus beinahe 80 Milliarden Euro an Kosten. Stau ist allerdings nicht nur finanziell für Privatpersonen oder auch Unternehmen ein großer Faktor, sondern er erhöht auch das Unfallrisiko und trägt maßgeblich zur schlechten Luftqualität in Innenstädten bei. Aufgrund längerer Fahrzeiten und der häufigen Be- und Entschleunigung, steigt der Kraftstoffverbrauch der Fahrzeuge und dadurch auch die Schadstoffbelastung in der Luft [Hemmerle, 2016].

Um Stau so gut wie möglich vermeiden zu können, muss man den Verkehr verstehen. Nötig ist ein Verständnis des Straßenverkehrs als Ganzes, sowie der Auswirkungen, welche einzelne Verkehrsteilnehmer und deren Verhalten, auf diesen haben. Hierzu ist das Erstellen von Simulationen sowie die Auswertung realer Verkehrsaufkommen unerlässlich. Die auf diese Weise gesammelten Erkenntnisse bilden die Grundlage, um Straßenabschnitte, insbesondere auch in Innenstädten, intelligent zu gestalten. Des Weiteren können sie eingesetzt werden, um beispielsweise Ampelschaltungen in Städten zu optimieren, wovon auch bestehende Infrastrukturen profitieren können.

Diese Arbeit beschäftigt sich mit der Realisierung einer automatischen Fahrspurerkennung aus Luftaufnahmen, welche bei der Analyse von Spurwechselvorgängen zum Einsatz kommt. Hierzu werden die Trajektoriedaten von Fahrzeugen ausgewertet.

1.1 Rahmen der Arbeit

1.1.1 Das Projekt MEC-View

1.1.2 MEC-View Luftbeobachtung

1.2 Motivation und Ziele

1.3 Aufbau dieser Arbeit

2 Grundlagen

In diesem Kapitel werden die für das Verständnis und die Durchführung der Arbeit benötigten Grundlagenthemen vorgestellt. Nach einer kurzen Erläuterung der Möglichkeiten der Verkehrsanalysen mittels Luftaufnahmen, wird daher darauf eingegangen, auf welche Weise die in dieser Arbeit verwendeten Fahrzeugtrajektorien ermittelt werden. Anschließend werden Methoden vorgestellt, welche zur Bereinigung der gewonnenen Daten verwendet werden können. Als wichtiges Mittel zur Identifizierung von Fahrspuren aus Trajektorien werden zudem verschiedene Cluster-Algorithmen und Distanzmaße vorgestellt.

2.1 Verkehrsanalyse mittels Luftaufnahmen

2.2 Rekonstruktion von Fahrzeugtrajektorien aus Luftaufnahmen

Die in dieser Arbeit verwendeten Fahrzeugtrajektorien stammen aus der Anwendung “Tracker-Application” des MEC-View Teilprojektes *Luftbeobachtung*. Nachfolgend wird beschrieben, wie diese aus den Videoaufnahmen rekonstruiert werden.

2.3 Datenaufbereitung und Bereinigung

2.4 Clusteranalyse

Die Clusteranalyse (kurz Clustering) ist ein wichtiges Werkzeug zur Auswertung von Daten unterschiedlichster Art. Sie stellt dabei kein konkretes Vorgehen oder einen Algorithmus dar, sondern beschreibt ein allgemeines Problem, welches auf unterschiedlichste Weise gelöst werden kann. Grundsätzlich ist das Ziel der Clusteranalyse, Datenobjekte aufgrund ihrer Eigenschaften und Beziehungen untereinander so zu gruppieren, dass sich die Objekte einer Gruppe möglichst stark ähneln und sich von den Objekten anderer Gruppen möglichst stark unterscheiden. Je höher die *Homogenität* in einem Cluster und die *Differenz* zwischen den Clustern, desto besser ist die gewählte Clustering Methode. Der Einsatz von Clustering ist in vielen Anwendungsgebieten und in den unterschiedlichsten wissenschaftlichen Disziplinen sehr beliebt, um ein Verständnis für Daten zu erhalten beziehungsweise

diese anschließend weiter verarbeiten zu können. So kommt die Clusteranalyse unter anderem in den Feldern des maschinellen Lernens, der Mustererkennung, Bildanalyse, der Biologie (Taxonomie) oder im Bereich Data Mining zum Einsatz. [Tan u. a., 2007]

Die Clusteranalyse hat viel mit dem Problem der Klassifizierung von Daten gemein, insofern sie Datenobjekten Label zuordnet. Im Gegensatz zu *überwachten* Klassifizierungsansätzen, wie dem heute populären überwachten Lernen, leiten Cluster-Algorithmen die Label allerdings alleine aus den vorhandenen Daten ab. Es kommen keine Vergleichsobjekte mit bekannten, händisch vergebenen Labels zum Einsatz. Aus diesem Grund wird die Clusteranalyse auch häufig als *unüberwachte Klassifizierung* bezeichnet. [Tan u. a., 2007]

Das Konzept eines *Clusters* ist nicht genau definiert, was in einer Vielzahl an unterschiedlichen Ansichten und Algorithmen resultiert, welche sich jeweils für andere Anwendungsfälle eignen und verschiedene Eigenschaften besitzen. Hieraus ergibt sich auch die Tatsache, dass Clustering keine selbsttätiger Prozess ist, welcher sich auf einheitliche Weise auf unterschiedliche Probleme anwenden lässt. Jedes Problem erfordert die individuelle und sorgfältige Auswahl eines passenden Algorithmus, eines Distanzmaßes und der richtigen Parameter. Die Bestimmung dieser geschieht iterativ und nicht selten nach dem Prinzip des *Trial and Error*. In Abbildung 2.1 ist beispielhaft ein Datensatz (links) mit – für den Menschen intuitiv ersichtlich – 7 unterschiedlichen Clustern (rechts) dargestellt. Nach [Jain, 2010] kann allerdings kein verfügbarer Clustering Algorithmus diese alle erkennen. [Jain u. a., 1999; Tan u. a., 2007]

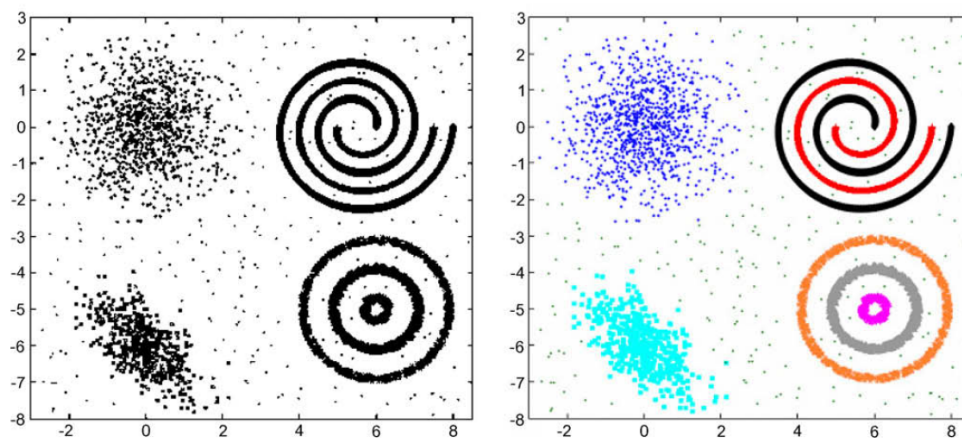


Abb. 2.1: Rohdaten (links) und erwünschtes Clustering-Ergebnis (rechts) [Jain, 2010]

Aufgrund der Limitationen, welche alle Cluster-Algorithmen besitzen, muss der Analyst sich vor deren Anwendung intensiv mit den zu verarbeitenden Daten beschäftigen. Er muss ein Verständnis dafür besitzen, welche Struktur die Daten besitzen, beziehungsweise annehmen können, und nach welchen Mustern zu suchen ist. Besonders wichtiger ist zudem auch die Auswahl der richtigen, das heißt relevanten, Datenmerkmale (*“Feature Selektion”*) und die Wahl deren Repräsentation (*“Feature Transformation”*). Die Selektion und

gegebenenfalls Transformation der Daten muss in einem Vorverarbeitungsschritt geschehen, dessen Qualität einen maßgeblichen Einfluss auf das finale Clustering Ergebnis hat. Basierend auf vorangegangener Beschreibung und [Jain u. a., 1999], lässt sich der Ablauf einer Clusteranalyse wie folgt darstellen:

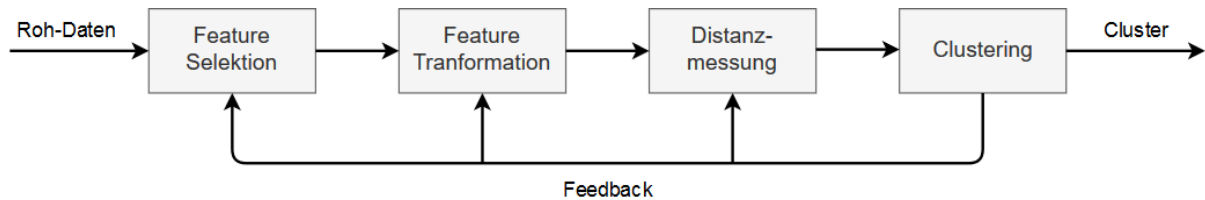


Abb. 2.2: Ablauf einer Clusteranalyse

[Jain, 2010] nennt einige weitere Herausforderungen, welchen man sich bei der Clusteranalyse bewusst sein muss:

- Daten können Ausreißer enthalten. Wie sollen diese behandelt werden?
- Die Anzahl der Zielcluster ist üblicherweise nicht bekannt.
- Validierung der gefundenen Cluster

2.4.1 Eigenschaften von Cluster-Sets und Clustern

Das aus einer Analyse resultierende Cluster-Set und die einzelnen Cluster selbst, können in verschiedene Kategorien unterteilt werden beziehungsweise unterschiedliche Eigenschaften besitzen. Nachfolgend sind die wichtigsten basierend auf [Tan u. a., 2007] und [Jain u. a., 1999; Jain, 2010] aufgeführt.

Cluster-Sets

Bei Cluster-Sets kann grundsätzlich zwischen nachfolgenden Eigenschaften unterschieden werden.

Hierarchisch vs. Partitioniert Von *hierarchischen* Cluster-Sets wird gesprochen, wenn die einzelnen Cluster verschachtelt sind und dabei eine Baum-Struktur bilden. Cluster sind hingegen *partitioniert*, wenn keine Überlagerungen zwischen ihnen existiert.

Exklusiv vs. Überlappend vs. Fuzzy *Exklusive* Cluster-Sets liegen vor, wenn jedem Datenwert ein oder kein Zielcluster zugeordnet wird. Im Gegensatz hierzu können bei *überlappenden* Cluster-Sets Objekte einer oder mehrerer Gruppen angehören. Bei dem sogenannten *Fuzzy* oder *Soft* Cluster-Sets, gehört ein Datenobjekt einem Cluster mit einer bestimmten Wahrscheinlichkeit oder Gewicht an. Algorithmen, welche Daten eine Wahrscheinlichkeit für die Zugehörigkeit zu einem Cluster zuweisen, werden *probabilistische* Cluster-Algorithmen genannt.

Komplett vs. Partielle Von *kompletten* Cluster-Sets wird gesprochen, wenn jedes Element der Eingangsdaten einem Cluster zugeordnet wird. Bei *partiellen* Sets ist dies nicht der Fall. Hier kann ein bestimmter Anteil an Datenwerten als Ausreißer markiert werden, welche keine Gruppe besitzen.

Cluster

Da, wie oben erwähnt, nicht klar definiert ist, was ein Cluster ausmacht, können auch diese unterschiedliche Eigenschaften besitzen. Die wichtigsten Cluster-Arten sind nachfolgend erläutert.

Klar separierte Cluster Unter *klar separierten* Clustern versteht man solche, in welchen jedes Datenelement einen geringeren Abstand zu allen anderen Elementen des Clusters hat, als zu Elementen außerhalb des Clusters. Diese idealistische Definition eines Clusters ist nur dann erfüllt, wenn die in den Daten enthaltenen Cluster einen großen Abstand voneinander haben. Dies ist in der Realität allerdings selten der Fall.

Prototyp basierte Cluster Von einem *Prototyp basierten* Cluster wird gesprochen, wenn alle Elemente einer Gruppe einen geringeren Abstand zu einem Prototyp oder Referenzwert des Clusters besitzen, als zu denen anderer Gruppierungen. Ein solcher Prototyp ist üblicherweise der Mittelwert der Datenelemente eines Clusters (*Centroid*).

Graphen basierte Cluster Die Definition eines *Graphen basierten* Clusters kann immer dann verwendet werden, wenn Daten als vernetzter Graph dargestellt werden. In einem solchen sind die Elemente Knoten und die Kanten repräsentieren Beziehungen zwischen ihnen. Ein Cluster in einem solchen Graphen ist definiert als Menge von Knoten, welche untereinander verbunden sind, jedoch keine Verbindungen zu Elementen außerhalb des Clusters haben.

Dichte basierte Cluster *Dichte basierte* Cluster sind definiert als Regionen mit einer hohen Dichte an Objekten, welche von Regionen umgeben sind, welche eine geringe Objektdichte besitzen. Elemente, welche in einer solchen Region mit geringen Dichte liegen, welche aus Ausreißer interpretiert. Dichte Bereiche werden üblicherweise gefunden, indem die Nachbarschaften von Elementen untersucht werden.

Konzeptionelle Cluster Eine sehr allgemeine Definition eines Clusters ist die der *konzeptionellen* Gruppen. Hiermit ist gemeint, dass die Elemente eines Clusters einige gemeinsame Eigenschaften besitzen. Dies schließt die oben genannten Cluster-Arten mit ein, lässt sich allerdings beliebig erweitern. So sind beispielsweise in Abbildung XXX x) konzeptionelle Cluster dargestellt, die die Form zweier Kreise und eines Dreiecks haben. Um solche Muster erkennen zu können, würde ein Algorithmus eine besondere Definition eines Clusters benötigen.

2.4.2 Cluster-Algorithmen

Um mit den oben beschriebenen unterschiedlichen Cluster-Set und Cluster Definitionen umgehen zu können, existieren verschiedene Clustering-Modelle. Einige wichtige Clustering-Ansätze sind die Vernetzungs-Modelle, Centroid-basierte-Modelle, Verteilungs-Modelle oder Dichte-Modelle. Für jedes dieser Modelle existieren unterschiedliche Algorithmen. Im Folgenden werden diese Modelle und jeweils exemplarisch ein Algorithmus der diese vertritt vorgestellt.

Vernetzungs-Modelle

Vernetzungs-Modelle werden auch häufig *hierarchische Cluster-Modelle* genannt. Sie beruhen auf der Annahme, dass Elemente, welche nahe beieinander liegen, eine höhere Gemeinsamkeit besitzen als solche, welche weiter voneinander entfernt sind. Zur Bestimmung der Nähe zwischen Elementen benötigen Vernetzungs-Modelle, wie auch andere Cluster-Modelle, eine Definition von Distanz. Diese legt ein sogenanntes *Distanzmaß* fest. Zusätzlich ist ein *Link-Kriterium* notwendig, welches bestimmt, wie genau die Entfernung zwischen zwei Clustern ermittelt wird. Übliche Link-Kriterien sind *Minimum-Linkage*, welches die minimale Distanz zwischen den Objekten der Cluster als Distanz verwendet, oder *Maximum-Linkage* beziehungsweise *Average-Linkage*. [Jain u. a., 1999; George Seif, 2018]

Grundsätzlich teilen sich hierarchische Cluster-Algorithmen in zwei Gruppen auf: *Agglomerative* (Bottom-Up) und *Divisive* (Top-Down) Algorithmen. Agglomerative Ansätze weisen zu Beginn des Cluster-Vorgangs jedem Datenelement eine eigene Gruppe zu und vereinen diese anschließend. Bei divisiven Ansätzen werden hingegen zu Beginn alle Elemente in einem Cluster zusammengefasst und diese in den nachfolgenden Schritten geteilt.

Als Beispiel wird anschließend der **agglomerative-hierarchische Cluster-Algorithmus** genauer vorgestellt. Sein Vorgehen lässt sich sehr gut anhand sogenannter Dendrogramme oder geschachtelter Cluster-Diagramme darstellen (siehe Abbildung 2.3)

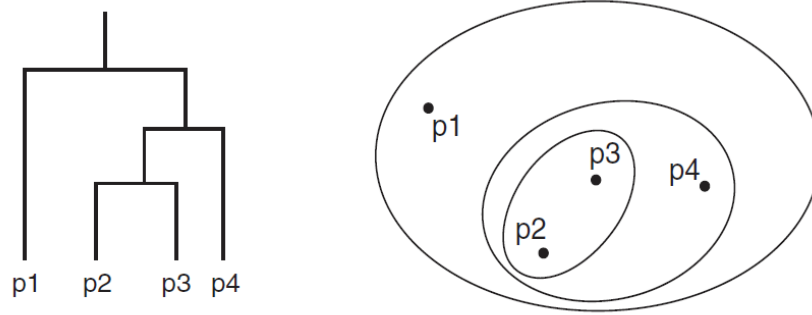


Abb. 2.3: Agglomeratives Clustering dargestellt als Dendrogramm und geschachteltes Cluster-Diagramm

Im ersten Schritt des Algorithmus werden alle Datenpunkte als separate Cluster markiert. Diesen Schritt repräsentieren die Blätter des Dendrogramms. Anschließend muss ein Distanzmaß und ein Link-Kriterium gewählt werden. Das am häufigsten verwendete Distanzmaß ist sicherlich der euklidische Abstand, welcher die Distanz zwischen zwei Punkten oder Vektoren im n -dimensionalen Raum bestimmt. Er ist definiert durch die Formel 2.1.

$$dist(p, q) = \|q - p\|_2 = \sqrt{\sum_{i=1}^n (q_i - p_i)^2} \quad (2.1)$$

Wird als Link-Kriterium beispielsweise *Minimum-Linkage* gewählt, ist dieses definiert als:

$$link(P, Q) = \min\{dist(p, q) : p \in P, q \in Q\} \quad (2.2)$$

Hierbei entsprechen P und Q zwei Clustern, welche die Elemente $p \in P$ und $q \in Q$ enthalten. Auf Basis des gewählten Link-Kriteriums kann nun eine Distanz-Matrix für die einzelnen Cluster erstellt werden. Die zwei Cluster mit minimalem Abstand voneinander werden anschließend zusammengeführt und die vorherigen Schritte werden wiederholt, bis nur noch ein Cluster (Wurzel des Dendrogramms) beziehungsweise die gewünschte Clusteranzahl übrig ist. [George Seif, 2018; Tan u. a., 2007]

Bei den meisten Varianten des agglomerativen Clusterings muss der Nutzer die Anzahl der Zielcluster im voraus festlegen, was problematisch ist, da diese meist nicht bekannt ist. Umgangen werden kann dies nur, indem ein Link-Kriterium gewählt wird, das ab

einer bestimmten Distanz zwischen den Clustern diese nichtmehr fusioniert [George Seif, 2018].

Die Zeitkomplexität des agglomerativen Clusterings beträgt bestenfalls $O(m^2 \log m)$, weshalb die Menge der Daten, welche mit ihm verarbeitet werden können erheblich begrenzt ist [Tan u. a., 2007].

Centroid-Modelle

Centroid basierte Cluster-Modelle betrachten im Gegensatz zu hierarchischen Modellen nicht die Distanz zwischen Clustern, sondern die Entfernung von Objekten zu Referenzpunkten, sogenannten *Centroids*.

Ein Beispiel für einen Centroid-Cluster-Algorithmus ist **k-Means**. Dieser ist aufgrund seines Alters, seiner Einfachheit und der vielen Weiterentwicklungen wohl der bekannteste Cluster-Algorithmus überhaupt.

Das Ziel von k-Mean ist es, für eine n-dimensionale Punktmenge $X = \{x_1 \dots x_n\}$ ein Cluster-Set $C = \{c_1 \dots c_k\}$ zu finden, welches die Summe der quadratischen Abweichung (Gleichung 2.3) zwischen allen Punkten in einem Cluster und deren Mittelwert μ_k (Centroids) minimiert.

$$J(c_k) = \sum_{k=1}^K \sum_{x_i \in c_k} \|x_i - \mu_k\|^2 \quad (2.3)$$

Eine Lösung für dieses Problem zu finden, ist NP-Schwer. Aus diesem Grund ist k-Means ein approximativer Ansatz, welcher nicht garantieren kann, ein globales Minimum zu finden. Die Funktionsweise des Algorithmus ist in Abbildung 2.4 dargestellt. Die Kreuze entsprechen hierbei den Centroids, welche sich über die Iterationen hinweg verschieben.

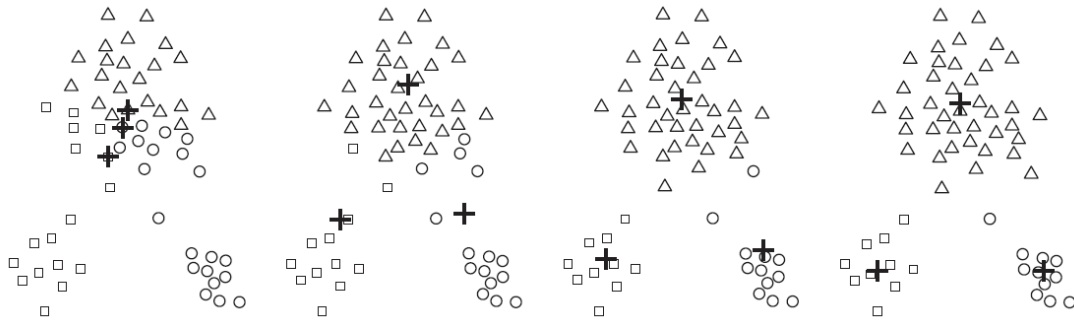


Abb. 2.4: Funktionsweise von k-Means [Tan u. a., 2007]

Ausgehend von der Punktmenge X und der gesuchten Cluster-Anzahl k , werden im ersten Schritt k zufällig positionierte Centroids μ_k definiert. Anschließend wird für alle Punkte x_i der nächstgelegene Centroid μ_j gesucht.

$$j = \arg \min(\text{dist}(x_i, \mu_j)) \quad (2.4)$$

x_i wird daraufhin Mitglied in Cluster C_j . Als Distanzmaß (dist) kann hier wieder der euklidische Abstand (Gleichung 2.1) verwendet werden oder aber auch beliebige andere sinnvolle Metriken. Nachdem alle Punkte x_i einem Cluster zugewiesen wurden, werden die Centroid Positionen neu bestimmt. Hierzu wird der Durchschnitt aller Punkte eines Clusters berechnet:

$$c_j = \frac{1}{n} \sum_{x_j \in C_j} x_j \quad (2.5)$$

Diese zwei Schritte werden mehrfach wiederholt, bis das Ergebnis konvergiert, das heißt die Zuweisungen sich nur noch geringfügig ändern. [Jain, 2010]

Der primäre Nachteil des k-Means Algorithmus ist, das auch bei ihm die Anzahl der Zielcluster spezifiziert werden muss. Desweiteren ist sein Ergebnis aufgrund der zufälligen Initialisierung der Centroids nicht deterministisch. Vorteil von k-Means ist hingegen, dass seine Zeitkomplexität bei $O(n)$ liegt.

Um die genannten Nachteile, zumindest in Teilen, umgehen zu können, existieren diverse Weiterentwicklungen des k-Mean Algorithmus. So stammen beispielsweise von [Hamerly u. a.] und [Pelleg u. a.] die Algorithmen *g-Means* beziehungsweise *x-Means*, welche die Clusteranzahl k auf Basis mehrerer k-Means Durchläufe und statistischer Kennzahlen bestimmen.

Distributions-Modelle

Distributions-Cluster-Modelle basieren auf der Verwendung von statistischen Wahrscheinlichkeitsverteilungen wie beispielsweise der Gauß-Verteilung. Cluster werden darüber definiert, wie wahrscheinlich es ist, dass Objekte der selben Verteilung angehören. Problematisch ist die Verwendung dieser Cluster-Methodik, da sie anfällig für das Problem des "*Overfitting*" ist, wenn die Komplexität der verwendeten Modelle nicht beschränkt wird. Zudem ist die Annahme, dass vielen realen Datensätzen ein statistisches Verteilungsmodell zugrundeliegt, gefährlich. Ist diese These jedoch berechtigt, haben die Modelle den Vorteil, dass sie neben der Zuweisung von Objekten zu Clustern auch Korrelationen zwischen einzelnen Attributen aufzeigen können. [Anders Drachen, 2014]

Nachfolgend wird der bekannteste Vertreter der Distributions-Cluster-Algorithmen vorgestellt: das *Expectation-maximization* (EM) Verfahren unter Verwendung sogenannter

Gaussian-Mixture-Models (GMM). Die Funktionsweise des EM-Algorithmus hat grundsätzlich viel gemein mit der des k-Mean Ansatzes. Es wird ebenfalls mit einer festen Anzahl zufällig initialisierter Modelle gestartet, welche anschließend über mehrere Iterationen an die Daten angepasst werden. Im Gegensatz zu k-Means, sind die gewählten Modelle hingegen Gauß-Verteilungen, welche zwei Parameter besitzen: ihren Mittelwert und die Standardabweichung. Das Vorgehen des EM-Algorithmus ist nachfolgend, basierend auf [George Seif, 2018], beschrieben und in Abbildung 2.5 grafisch dargestellt.

- 1) Wahl der Clusteranzahl k und Initialisierung der Gauß-Modelle für die entsprechenden Cluster.
- 2) Berechnung der Wahrscheinlichkeit, dass ein Datenpunkt zu einem Cluster gehört. Je näher ein Datenpunkt dem Zentrum einer Gauß-Verteilung ist, desto höher die Wahrscheinlichkeit für dessen Zugehörigkeit.
- 3) Basierend auf den Wahrscheinlichkeiten werden die Parameter der Verteilungen neu berechnet. Hierzu wird die gewichtete Summe der Datenpunkt-Positionen errechnet. Die Gewichte entsprechen dabei den Wahrscheinlichkeiten, dass ein Element zu einem Cluster gehört. Hierdurch werden die Gauß-Modelle automatisch den in den Daten enthaltenen Clustern angepasst.
- 4) Wiederholung der Schritte 2) und 3), bis das Clustering-Ergebnis konvergiert.

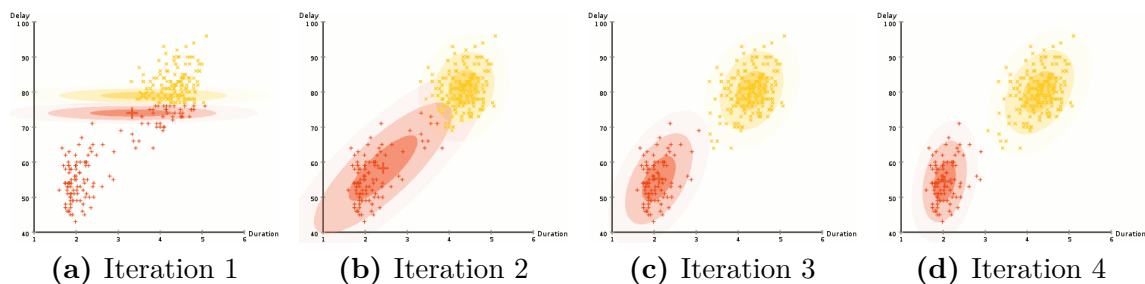


Abb. 2.5: Darstellung des EM-Cluster-Algorithmus über mehrere Iterationen [George Seif, 2018]

Ziel des EM-Algorithmus ist es, die Parameter der Gauß-Modelle so zu optimieren, dass diese die Verteilung der Daten bestmöglich beschreiben. Am Ende des Clusterings besitzt jeder Datenwert die Zugehörigkeit-Wahrscheinlichkeiten für die einzelnen Cluster. Ein Element wird jenem Cluster zugeordnet, für welches es die höchste Wahrscheinlichkeit besitzt.

Dichte-Modelle

Dichte basierte Cluster sind, wie oben beschrieben, definiert als Regionen hoher Objektdichte, welche von Bereichen geringer Dichte umgeben sind. Dichte-Clustering-Modelle

suchen nach eben solchen Regionen. Großer Vorteil der Algorithmen dieser Klasse ist, dass sie Cluster beliebiger Formen finden können, nicht auf die Vorgabe einer Clusteranzahl angewiesen sind und mit Ausreißern umgehen können.

Als Vertreter der Dichte-basierten Ansätze wird nachfolgend der **DBSCAN** Algorithmus (*Density-Based Spatial Clustering of Applications with Noise*), wie in [Gao, 2012] beschrieben, vorgestellt. Er verwendet als Maß für die Dichte einer Region die sogenannte ϵ -Nachbarschaft (*Eps*). Diese selektiert für ein Objekt p alle Objekte, welche innerhalb des Radius ϵ um dieses liegen:

$$N_\epsilon(p) = \{q | \text{dist}(p,q) \leq \epsilon\} \quad (2.6)$$

Eine ϵ -Nachbarschaft besitzt eine hohe Dichte, wenn in ihr mindestens *MinPts* Objekte liegen.

Basierend auf der Definition von *Eps*, werden die in einem Datensatz vorhandenen Elemente in drei Klassen unterteilt. Sie sind entweder *Kern*-, *Rand*- oder *Ausreißer*- Objekte. Ein Kernobjekt hat mindestens *MinPts* andere Punkte in *Eps*. Randobjekte besitzen weniger als *MinPts* in *Eps*, liegen aber in der Nachbarschaft eines Kernobjektes. Ausreißerobjekte sind weder Kern- noch Randobjekte.

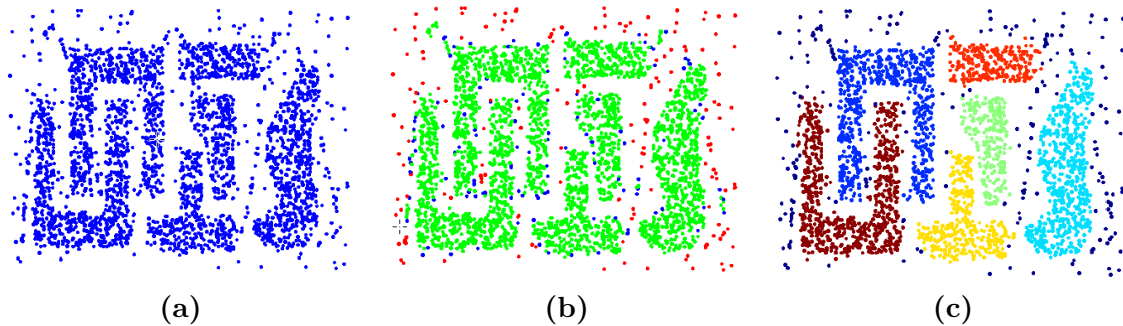


Abb. 2.6: Schritte des DBSCAN Algorithmus, a) Rohdaten, b) Klassifizierung in Kern- (grün), Rand- (blau) und Ausreißer- (rot) Punkte, c) Cluster Ergebnis [Gao, 2012]

Auf Basis der drei Objektklassen, lässt sich das Prinzip der dichte-basierten *Erreichbarkeit* definieren. Ein Objekt q ist von p *direkt* erreichbar, wenn p ein Kernobjekt ist und q in dessen *Eps* liegt. In Abbildung 2.7 gilt dies beispielsweise für p und p_2 . Zwei Elemente sind *indirekt* erreichbar, wenn sie über eine Reihe von Zwischenschritten (direkte Relationen) verbunden sind (transitiv). Dies ist in Abbildung 2.7 für q und p der Fall.

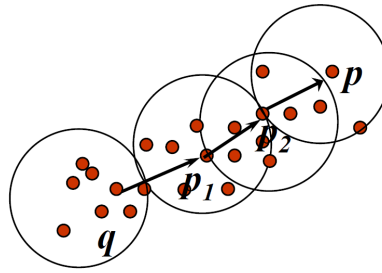


Abb. 2.7: Erreichbarkeit in DBSCAN [Gao, 2012]

Der DBSCAN Algorithmus lässt sich, basierend auf den obigen Definitionen, informell wie folgt beschreiben:

- 1) Unterteilung der Objekte in die drei Objektklassen. (Abb. 2.6 b))
- 2) Aussortierung der Ausreißer-Objekte.
- 3) Wahl eines nicht zugewiesenen Kernobjektes.
- 4) Erstellung eines neuen Clusters für das Kernobjekt und alle von ihm ausgehend direkt oder indirekt erreichbaren Objekte
- 5) Wiederholung der Schritte 3) und 4), bis alle Kern- und Randobjekte einem Cluster zugewiesen sind. (Abb. 2.6 c))

DBSCAN besitzt die oben beschriebenen Vorteile Dichte-basierter Cluster-Algorithmen. Dank einer Zeitkomplexität von $O(n \log n)$ kann er außerdem auch auf große Datensätze angewendet werden. Nachteil des Ansatzes ist hingegen, dass er schlecht mit Clustern umgehen kann, welche unterschiedliche Dichten besitzen.

2.5 Ähnlichkeitsmaße zum Vergleich von Fahrzeugtrajektorien

Bei der Clusteranalyse ist neben der Wahl des passenden Cluster-Algorithmus insbesondere die Entscheidung, welches Ähnlichkeits- beziehungsweise Distanzmaß verwendet wird, ausschlaggebend. Im obigen Abschnitt wurden bereits die euklidische Distanz (Gleichung 2.1) als ein mögliches Distanzmaß definiert. Dieses kann jedoch nur zur Bestimmung der Distanz zwischen n -dimensionalen Punkten im euklidischen Raum verwendet werden. Dies gilt ebenso für andere einfache Maße wie die Manhattan-Distanz oder die Pearson-Distanz.

Um Fahrzeugtrajektorien korrekt gruppieren zu können, ist ein Ähnlichkeitsmaß notwendig, welches je nach Anforderungen die unterschiedlichen Aspekte der Trajektorien vergleicht. Häufig werden die Eigenschaften Lage, Form und Länge hierzu herangezogen. In der Literatur werden diverse Maße zum Vergleich von Trajektorien vorgestellt. Diese besitzen alle unterschiedliche Eigenschaften, Vor- und Nachteile.

Nachfolgend werden exemplarisch drei Ähnlichkeitsmaße vorgestellt, anhand welcher ersichtlich ist, welche Abwägungen bei der Wahl des Maßes gemacht werden müssen. In allen drei Fällen werden die Trajektorien als Reihen 2-dimensionaler Punkte mit Länge n interpretiert: $t_i = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$. Der n -te Punkt einer Trajektorie ist gegeben über $t_i(n)$ und deren Punkt-Länge über $len(t_i)$. Die Menge der zu vergleichenden Trajektorien ist $T = \{t_1, t_2, \dots, t_m\}$. Abbildung 2.8 zeigt eine Auswahl möglicher Trajektorien.

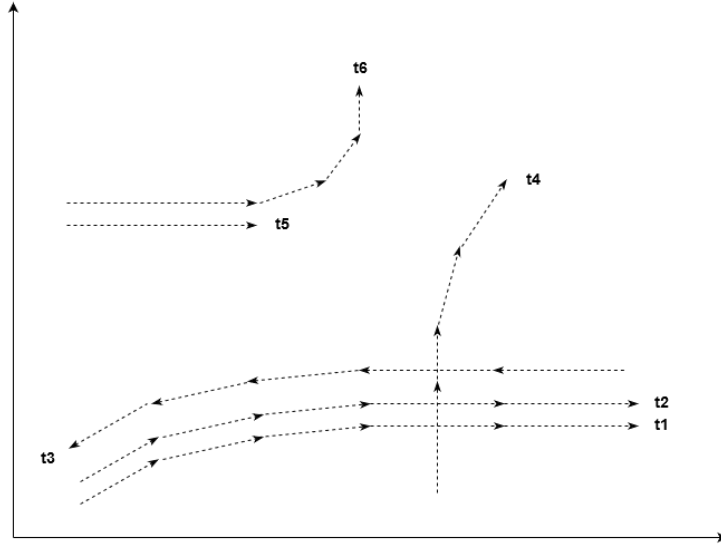


Abb. 2.8: Trajektorien im 2-dimensionalen Raum

2.5.1 HU Distanz

Die HU Distanz wurde erstmals in der Arbeit *“Similarity based vehicle trajectory clustering and anomaly detection”* von [Hu u. a., 2005] verwendet. Es ist ein sehr einfaches Distanzmaß, welches auf der mittleren euklidischen Distanz zwischen zwei Trajektorien basiert. Berechnet wird die HU Distanz für zwei Trajektorien t_1 und t_2 wie folgt:

$$D_{HU}(t_1, t_2) = \frac{1}{N} \sum_{n=1}^N dist(t_1(n), t_2(n)) \quad (2.7)$$

$$wobei N = \min(len(t_1), len(t_2)) \quad (2.8)$$

Aus dieser Formel lassen sich die Vor- und Nachteile der HU Distanz ableiten. Der klare Vorteile der HU Distanz ist deren Einfachheit und die Effizienz von $O(n)$. Nachteil ist hingegen, dass das Distanzmaß nur gut funktioniert, wenn die Trajektorien bestimmte Kriterien erfüllen. So sollten Trajektorien, welche einem Cluster angehören, auch immer möglichst auf selber Höhe beginnen, damit deren mittlerer Abstand nicht, aufgrund einer Verschiebung, erhöht wird. Außerdem ist es notwendig, die Abstände zwischen den Punkten der Trajektorien auf die selbe Länge zu bringen, damit beim paarweisen Vergleich

immer Elemente verglichen werden, welche gleichweit vom Start der Spuren entfernt sind. Diese Eigenschaften der Trajektorien müssen über einen Vorverarbeitungsschritt geschaffen werden. Problematisch bei der Verwendung der HU Distanz ist außerdem, dass beim Vergleich zweier Trajektorien immer nur die ersten N Punkte (s. Gleichung 2.8) betrachtet werden. Die kann dazu führen, dass zwei Trajektorien, welche zu Beginn fast identisch sind und später auseinanderlaufen, trotzdem einen hohen Ähnlichkeitswert besitzen (siehe t_5 und t_6 in Abbildung 2.8).

Die HU Distanz kann aufgrund der genannten Einschränkungen nur in speziellen Fällen oder unter Verwendung eines Vorverarbeitungsschrittes angewandt werden. Sie liefert ansonsten suboptimale Clustering Ergebnisse.

2.5.2 Hausdorff Distanz

Die Hausdorff Distanz ist ein komplexeres Maß zur Bestimmung der Ähnlichkeit zwischen zwei Trajektorien. Sie misst grundsätzlich den Abstand zwischen zwei nicht-leeren, ungeordneten Teilmengen A und B und ist für Trajektorien definiert über die Gleichungen [Atev u. a., 2010]:

$$D_{HD}(t_1, t_2) = \max(h(t_1, t_2), h(t_2, t_1)) \quad (2.9)$$

$$h(t_1, t_2) = \max_{i \in t_1} \min_{j \in t_2} \text{dist}(i, j) \quad (2.10)$$

$h(t_1, t_2)$ wird als gerichtete Hausdorff Distanz *von t_1 nach t_2* bezeichnet. Sie findet die maximale Distanz einer Trajektorie zum nächsten Punkt der anderen Trajektorie [Huttenlocher u. a., 1993]. Da h gerichtet ist, gilt $h(t_1, t_2) \neq h(t_2, t_1)$. Aus diesem Grund wird die Hausdorff Distanz *zwischen* zwei Trajektorien mittels D_{HD} bestimmt. dist kann ein beliebiges Maß für die Distanz zweier Punkte sein, wie beispielsweise die euklidische Distanz. Grundsätzlich lässt sich über die Hausdorff Distanz die Form zweier Trajektorien vergleichen. Diese sind ähnlich, wenn jeder Punkt einer Trajektorie einen nahegelegenen Punkt in der Vergleichsbahn besitzt.

Vorteil der Hausdorff Distanz im Vergleich zur HU Distanz ist, dass diese immer vollständige Trajektorien vergleicht und nicht nur Teile. Außerdem ist bei ihrer Verwendung keine Vorverarbeitung in Form von Resampling et cetera notwendig. Problematisch ist das Distanzmaß hingegen, da es mit ungeordneten Sets arbeitet und somit im Fall von Trajektorien deren Orientierung nicht beachtet. Zwei parallel aber in entgegengesetzte Richtungen laufende Trajektorien würden nach Hausdorff daher eine hohe Ähnlichkeit besitzen. Zudem kann das Distanzmaß schlecht mit Ausreißern umgehen, da bereits ein einzelner dieser Punkte, bei ansonsten identischen Trajektorien, zu einer beliebig kleinen Ähnlichkeit führen kann. Von Nachteil ist auch, dass die Zeitkomplexität der Hausdorff-Distanz bei $O(n \cdot m)$ liegt.

2.5.3 Longest-Common-Subsequence

Das *Longest-Common-Subsequence* (LCSS) Ähnlichkeitsmaß basiert auf dem allgemeinen Problem der Findung einer längsten gemeinsamen Subsequenz zwischen zwei Sequenzen. Da Trajektorien, nach obiger Definition, lediglich Punktfolgen sind, lässt sich das Verfahren sehr gut auf diese anwenden. Aufgrund einiger kleiner Erweiterungen des Basis-Algorithmus, besitzt das LCSS Ähnlichkeitsmaß einige besondere Eigenschaften. Der LCSS Algorithmus für Trajektorien ist grundsätzlich wie folgt definiert [Vlachos u. a., 2002]:

$$LCSS_{\epsilon, \delta}(t_1, t_2) = \begin{cases} 0 & \text{if } t_1 \text{ or } t_2 \in \emptyset \\ 1 + LCSS_{\epsilon, \delta}(t'_1, t'_2) & \text{if } dist(t_1(n), t_2(m)) < \epsilon \\ & \wedge |n - m| \leq \delta \\ \max(LCSS_{\epsilon, \delta}(t'_1, t_2), LCSS_{\epsilon, \delta}(t_1, t'_2)) & \text{otherwise} \end{cases} \quad (2.11)$$

Hierbei gilt $t'_1 = \{t_1(0), \dots, t_1(n-1)\}$. Die Parameter ϵ und δ bestimmen das Vergleichsverhalten des Algorithmus. Über ϵ wird definiert, wieweit zwei Punkte maximal voneinander entfernt liegen können, um immer noch als “übereinstimmend” zu gelten. δ bestimmt hingegen, wieweit man in der Zeit gehen darf, um einen übereinstimmenden Punkt zu finden. Abbildung XXX veranschaulicht die Bedeutung von ϵ und δ . Da die obige LCSS Funktion nur ein diskretes Zählmaß definiert, ist das eigentliche LCSS Ähnlichkeitsmaß üblicherweise wie folgt definiert [Vlachos u. a., 2002]:

$$D_{LCSS}(t_1, t_2) = 1 - \frac{LCSS(t_1, t_2)}{\min(len(t_1), len(t_2))} \quad (2.12)$$

Vorteile der LCSS Ähnlichkeitdefinition sind, dass sie mit kompletten Trajektorien arbeitet und robust gegenüber Ausreißern ist, da nicht für alle Punkte Übereinstimmungen in den Trajektorien gefunden werden müssen. Über ϵ und δ kann die “Strenge” des Algorithmus geregelt werden. Zudem berücksichtigt das LCSS Maß die Orientierung der Trajektorien, solange δ nicht zu hoch gewählt wird. Die rekursive Definition des LCSS Algorithmus aus Gleichung 2.11 lässt sich mittels dynamischer Programmierung mit Zeitkomplexität $O(n \cdot m)$ berechnen.

Übersicht Ähnlichkeitsmaße

Anhand der drei ausgewählten und oben exemplarisch beschriebenen Ähnlichkeitsmaße, ist bereits ersichtlich, dass die Wahl eines passenden Maßes nicht trivial ist. Es muss die Qualität und Form der Daten berücksichtigt werden und abgewogen werden, in wie weit es möglich beziehungsweise gewünscht ist, die Daten vorzuverarbeiten. Das primäre

Auswahlkriterium ist allerdings natürlich die situationsabhängige Definition von “Ähnlichkeit”: Sind sich Trajektorien ähnlich, wenn sie lediglich die selbe Form haben und ansonsten irgendwo im Raum liegen? Sind sie sich ähnlich, wenn sie die selbe Form haben und im Raum nahe beieinander liegen? Ist ihre Orientierung relevant? Dies sind wichtige Fragen, welche vor der Wahl eines Ähnlichkeitsmaßes geklärt werden müssen. Da die Maße als Distanzfunktionen bei der Clusteranalyse verwendet werden, ist ihr Verhalten ausschlaggebend für den Erfolg der Gruppierung.

Wichtige Eigenschaften einiger in der Literatur häufig verwendeten Vergleichsmaße, inklusive der drei oben beschriebenen, sind nachfolgend nochmals in tabellarischer Form festgehalten.

Tab. 2.1: Parallelen Bienenkolonie und Cloud Load-Balancing

Ähnlichkeitsmaß	gerichtet	PreProc. nötig	Ausreißer-resistent
HU [Hu u. a., 2005]	✓	✓	✓
PCA [Bashir u. a., 2003]	✓	✓	✓
DTW [Keogh und Pazzani, 2000]	✓	✗	✗
HD [Chen u. a., 2011]	✗	✗	✗
mod. HD [Atev u. a., 2010]	✓	✗	✓
PF [Piciarelli und Foresti, 2006]	✓	✗	✗
LCSS [Vlachos u. a., 2002]	✓	✗	✓

2.6 Untersuchung möglicher Straßentopologien

3 Stand der Technik

3.1 Clustering von Trajektoriedaten

3.2 Erkennung von Fahrspuren

3.3 Klassifizierung von Fahrspuren

3.4 Defizite vorhandener Lösungen und benötigte Neuerungen

4 Clustering von Fahrzeugtrajektorien

4.1 Vorverarbeitung der Roh-Trajektorien

4.2 Gruppierung der Trajektorien

5 Fahrbahn-Bestimmung aus Trajektorie-Clustern

6 Fahrbahn Klassifizierung

7 Realisierung LaneDetection in MEC-View TrackerApplication Software

8 Ergebnisse und Auswertung

9 Zusammenfassung und Fazit

Literaturverzeichnis

- [Anders Drachen 2014] ANDERS DRACHEN: *Introducing Clustering II: Clustering Algorithms - GameAnalytics*. 2014. – URL <https://gameanalytics.com/blog/introducing-clustering-ii-clustering-algorithms.html>. – Zugriffsdatum: 2018-10-22
- [Atev u. a. 2010] ATEV, Stefan ; MILLER, Grant ; PAPANIKOLOPOULOS, Nikolaos P.: Clustering of vehicle trajectories. In: *IEEE Transactions on Intelligent Transportation Systems* (2010). – ISBN 15249050
- [Bashir u. a. 2003] BASHIR, FI ; KHOKHAR, AA ; SCHONFELD, D.: Segmented trajectory based indexing and retrieval of video data. In: *Proceedings 2003 International Conference on Image Processing (Cat. No.03CH37429)* Bd. 3, URL <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.84.7768{&}rep=rep1{&}type=pdf><http://ieeexplore.ieee.org/document/1246757/>, 2003, S. II-623-6. – ISBN 0-7803-7750-8
- [Chen u. a. 2011] CHEN, J ; WANG, R ; LIU, L ; , J S. ; ; CONTROL, Communications ; 2011 undefined: Clustering of trajectories based on Hausdorff distance. In: *ieeexplore.ieee.org* (2011). – URL <https://ieeexplore.ieee.org/abstract/document/6066483/>
- [Cookson u. a.] COOKSON, G ; RESEARCH, B Pishue I. ; FEBRUARY undefined ; 2017 undefined: Inrix global traffic scorecard. In: *jschultheis.de*. – URL http://jschultheis.de/wp-content/uploads/2018/02/INRIX_2017_Traffic_Scorecard_Report___German.pdf
- [Gao 2012] GAO, Jing: Clustering Lecture 4 : Density-based Methods. (2012), S. 1-16. – URL https://cse.buffalo.edu/~jing/cse601/fa12/materials/clustering_density.pdf
- [George Seif 2018] GEORGE SEIF: *The 5 Clustering Algorithms Data Scientists Need to Know*. 2018. – URL <https://towardsdatascience.com/the-5-clustering-algorithms-data-scientists-need-to-know-a36d136ef68>. – Zugriffsdatum: 2018-10-19
- [Hamerly u. a.] HAMERLY, G ; PROCESSING, C Elkan A. in neural information ; 2004 undefined: Learning the k in k-means. In: *papers.nips.cc*. – URL <http://papers.nips.cc/paper/2526-learning-the-k-in-k-means.pdf>

- [Hemmerle 2016] HEMMERLE, P.: Empirische physikalische Eigenschaften des übersättigten innerstädtischen Verkehrs und Energieeffizienz von Fahrzeugen. (2016). – URL <https://d-nb.info/1122559259/34>
- [Hu u. a. 2005] HU, Weiming ; FU, Zhouyu ; TAN, Tieniu: Similarity based vehicle trajectory clustering and anomaly detection. In: *Proceedings - International Conference on Image Processing, ICIP*, 2005. – ISBN 0780391349
- [Huttenlocher u. a. 1993] HUTTENLOCHER, DP ; ... , GA Klanderman IEEE T. on ; 1993, Undefined: Comparing images using the Hausdorff distance. In: *ieeexplore.ieee.org* (1993). – URL <https://ieeexplore.ieee.org/abstract/document/232073/>
- [Jain u. a. 1999] JAIN, A. K. ; MURTY, M. N. ; FLYNN, P. J.: Data clustering: a review. In: *ACM Computing Surveys* 31 (1999), sep, Nr. 3, S. 264–323. – URL <http://portal.acm.org/citation.cfm?doid=331499.331504>. – ISSN 03600300
- [Jain 2010] JAIN, Anil K.: Data clustering: 50 years beyond K-means. In: *Pattern Recognition Letters* 31 (2010), jun, Nr. 8, S. 651–666. – URL <https://www.sciencedirect.com/science/article/pii/S0167865509002323>. – ISSN 0167-8655
- [Keogh und Pazzani 2000] KEOGH, Eamonn J. ; PAZZANI, Michael J.: Scaling up dynamic time warping for datamining applications. In: *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '00*, URL <https://dl.acm.org/citation.cfm?id=347153><http://portal.acm.org/citation.cfm?doid=347090.347153>, 2000, S. 285–289. – ISBN 1581132336
- [Pelleg u. a.] PELLEG, D ; ICML, AW M. ; 2000 undefined: X-means: Extending k-means with efficient estimation of the number of clusters. In: *cs.uef.fi*. – URL <http://cs.uef.fi/~zhao/Courses/Clustering2012/Xmeans.pdf>
- [Piciarelli und Foresti 2006] PICIARELLI, C ; FORESTI, G. L.: On-line trajectory clustering for anomalous events detection. In: *Pattern Recognition Letters* 27 (2006), Nr. 15, S. 1835–1842. – URL <https://www.sciencedirect.com/science/article/pii/S0167865506000432>. – ISSN 01678655
- [Tan u. a. 2007] TAN, Pang-Ning u. a.: *Introduction to data mining*. Pearson Education India, 2007
- [Vlachos u. a. 2002] VLACHOS, M ; ENGINEERING, G Kollios D. ; 2002 undefined ; 2002 undefined: Discovering similar multidimensional trajectories. In: *ieeexplore.ieee.org* (2002). – URL <https://ieeexplore.ieee.org/abstract/document/994784/>