

Masterarbeit  
in der Angewandten Informatik  
Nr. AI-2021-MA-014

# **Generisches Clustern hoch-komplexer Produktdaten**

Hannes Dröse

Abgabedatum: 29.04.2022

Prof. Dr. Ines-Kerstin Rossak  
Dipl.-Inform. Robert Queck

## **Zusammenfassung**

Die Kurzfassung wird zumeist als letzter Abschnitt geschrieben. Sie soll auf einer Seite den Hintergrund bzw. die Motivation sowie die zentralen Ergebnisse der Arbeit zusammenfassen. Der Zweck dieses Texts ist es, einem Recherchierenden oder Informationssuchenden zu einem bestimmten Thema zu signalisieren, ob er in diesem Werk für ihn relevante Informationen finden wird: Lohnt es sich, die Arbeit zu bestellen oder zu kaufen? Was ist neu?

Die Kurzfassung ist das sichtbare Aushängeschild der Arbeit und dient der Aufnahme in Referenzdatenbanken, in Online-Literatur-Shops u. Ä. Schauen Sie sich die Abstracts in der Fachliteratur (etwa in Fachbüchern, am Anfang von Zeitschriftenartikeln, in Konferenz- Proceedings, in ACM- oder IEEE-Literaturdatenbanken, im OPAC der Bibliothek usw.) an, um ein Gefühl dafür zu entwickeln, was hineingehört (siehe Abschnitt 2.4 Materialsammlung, S. 2).

Was definitiv nicht hineingehört, sind Gliederungsübersichten, chronologische Arbeitsberichte u. dgl.

### **Abstract**

The abstract has to be provided in english as well.

# Inhaltsverzeichnis

<b>Abbildungsverzeichnis</b>	<b>V</b>
<b>Tabellenverzeichnis</b>	<b>VI</b>
<b>1 Einleitung</b>	<b>1</b>
1.1 Hintergrund . . . . .	1
1.2 Themenfindung . . . . .	1
1.3 Kernfrage und Ablauf . . . . .	2
<b>2 Clusteranalyse</b>	<b>4</b>
2.1 Begriff und Einordnung . . . . .	4
2.2 Notation . . . . .	4
2.3 Distanz- und Ähnlichkeitsmaße . . . . .	5
2.3.1 Definition . . . . .	5
2.3.2 Numerische Attribute . . . . .	6
2.3.3 Kategorische Attribute . . . . .	8
2.3.4 String-Attribute . . . . .	10
2.3.5 Gemischte Attribute . . . . .	12
2.4 Clustering-Verfahren . . . . .	14
2.4.1 Partitionierendes Clustering . . . . .	14
2.4.2 Hierarchisches Clustering . . . . .	16
2.5 Cluster-Validität . . . . .	20
2.5.1 Überblick . . . . .	20
2.5.2 External Indices . . . . .	21
2.5.3 Internal Indices . . . . .	23
<b>3 Konzeption</b>	<b>26</b>
3.1 Überblick . . . . .	26
3.2 Datenquellen & -sets . . . . .	27
3.2.1 Akeneo-PIM . . . . .	27
3.2.2 Icecat . . . . .	31
3.2.3 Datenset . . . . .	32
3.3 Clustering . . . . .	33
3.3.1 Clustering-Verfahren . . . . .	33
3.3.2 Distanzfunktion . . . . .	34
3.4 Evaluation . . . . .	38
3.4.1 Kriterien . . . . .	38
3.4.2 Aspekte . . . . .	39
3.4.3 Vorgehen . . . . .	41

<b>4</b>	<b>Implementierung</b>	<b>42</b>
4.1	Überblick . . . . .	42
4.2	Akeneo-PIM . . . . .	42
4.2.1	Installation und Deployment . . . . .	42
4.2.2	Akeneo-Client und -Cache . . . . .	42
4.3	Datenset . . . . .	42
4.4	Clustering . . . . .	43
4.4.1	Überblick zu externen Libraries . . . . .	43
4.4.2	Clustering-Package . . . . .	43
4.4.3	Distanzfunktionen . . . . .	43
4.4.4	Datenvorbereitung . . . . .	43
4.5	Evaluation . . . . .	44
<b>5</b>	<b>Auswertung</b>	<b>45</b>
5.1	Datenset "Smartphone-Hüllen" . . . . .	45
5.1.1	Verarbeitung multi-kategorischer Attribute . . . . .	45
5.1.2	Attribut-Auswahl . . . . .	46
5.1.3	Attribut-Gewichtung . . . . .	46
5.2	Datenset "Smartphones" . . . . .	46
5.2.1	Verarbeitung multi-kategorischer Attribute . . . . .	46
5.2.2	Attribut-Auswahl . . . . .	47
5.2.3	Attribut-Gewichtung . . . . .	47
5.3	Kombiniertes Datenset . . . . .	47
5.3.1	Verwendung aller Attribute . . . . .	47
5.3.2	Verwendung gemeinsamer Attribute . . . . .	47
<b>6</b>	<b>Fazit und Ausblick</b>	<b>49</b>
	<b>Quellenverzeichnis</b>	<b>VII</b>
	<b>Anhang</b>	<b>X</b>
	<b>Selbstständigkeitserklärung</b>	<b>XI</b>

## Abbildungsverzeichnis

2.1	Bsp. eines Dendrograms . . . . .	17
3.1	Grober Architektur-Entwurf für den Praxisteil . . . . .	26
3.2	Schematische Darstellung der Funktionen von Akeneo-PIM [Aken22b] . . . . .	27
3.3	Akeneo-PIM Webfrontend bei der Bearbeitung eines Produktes . . . . .	28
5.1	Stabilität und Qualität der Hüllen mit multi-kategorischem Attribut . . . . .	45
5.2	Stabilität und Qualität der Hüllen mit String-Attributen . . . . .	45
5.3	Stabilität und Qualität der Smartphones mit multi-kategorischen Attributen . . . . .	46
5.4	Stabilität und Qualität der Smartphones mit String-Attributen . . . . .	48

## Tabellenverzeichnis

2.1	Eigenschaften der Abstandfunktion $d$ . . . . .	5
2.2	Übersicht über die gängigen Maße aus der Minkowski-Familie [Cha07] . . . . .	6
2.3	Kontingenz der Datenpunkt-Paare in $Y$ und $Y'$ [Ste04] . . . . .	22
3.1	Attribut-Typen in Akeneo-PIM [Aken22h] . . . . .	30
3.2	Einteilung der Akeneo-Typen in übergeordnete Datenklassen . . . . .	34
3.3	Vorverarbeitung der numerischen Attribute in Akeneo . . . . .	35
5.1	Clustering der Hüllen mit multi-kategorischem Attribut . . . . .	45
5.2	Clustering der Hüllen mit String-Attributen . . . . .	46
5.3	Clustering der Smartphones mit multi-kategorischen Attributen . . . . .	46
5.4	Clustering der Smartphones mit String-Attributen . . . . .	46

# 1 Einleitung

## 1.1 Hintergrund

Diese Masterarbeit ist in Kooperation mit der adesso SE entstanden. adesso ist ein deutsches IT-Beratungs- und Dienstleistungsunternehmen, welches 1997 gegründet worden ist mit Hauptsitz in Dortmund. Seit der Gründung ist die Firma sehr kontinuierlich gewachsen. Mittlerweile sind über 5.800 Mitarbeiter an 44 Standorten in 10 verschiedenen europäischen Ländern hier beschäftigt.

Über die Jahre sind nicht nur die Mitarbeiterzahlen gestiegen, sondern auch die Menge an abgedeckten Branchen, in denen das Unternehmen tätig ist. In den letzten Jahren sind zunehmend Projekte im E-Commerce-Sektor umgesetzt worden. Mit Beginn von 2022 kümmert sich eine eigene Abteilung des Unternehmens explizit um das Thema E-Commerce und Retail.

In diesen Bereich fällt die Programmierung und Betreuung von Online-Shops sowie das umfangreiche Thema des Product-Information-Managements (PIM). Dabei geht es um das Verwalten und Aufbereiten von Produktdaten für verschiedene Anwendungen wie Warenhaltung, Marketing oder Bestellabwicklung. adesso erwägt in Zukunft ein eigenes Product-Information-Management-System (PIM-System) zu entwickeln. In Vorbereitung dessen steht die Evaluierung der Machbarkeit verschiedener Ansätze und Ideen rund um das Product-Information-Management an. Dies ist auch der Aufhänger für diese Masterarbeit.

## 1.2 Themenfindung

Die initiale Idee bestand darin zu überprüfen, ob es mögliche Anwendungen des Verfahrens der Clusteranalyse im Zusammenhang mit dem Product-Information-Management gibt. Beim Clustering handelt es sich um die Einteilung von Objekten in Gruppen (sog. Cluster), sodass sich die Objekte in der gleichen Gruppe ähnlicher sind als Objekte in den anderen Gruppen. [King15, Kap. 1.1 What Is a Cluster?]

In einem ersten Brainstorming sind verschiedene Anwendungen der Clusteranalyse diskutiert worden. Es folgte eine intensive Literaturrecherche zu den besprochenen Anwendungen, die im Folgenden kurz diskutiert werden.

Die erste Idee war das automatische Kategorisieren von Produkten. Zu dieser Anwendung gibt es kaum Literatur – lediglich einige hoch-experimentelle Ansätze. Z.B. nutzten Chen et al. [ChZhYi19] ein neuronales Netz zum Vorschlagen von Produktkategorien zu Produkten bestehend aus Titel und Beschreibung.



Eine weitere Idee war die Nutzung der Clusteranalyse zur Entwicklung eines Produktempfehlungsalgorithmus. Zu diesem Thema gibt es ganze Reihe an Arbeiten:

- Oh und Kim [OhKi19] implementierten einen Empfehlungsalgorithmus für Cloud Computing-Angebote. Die verschiedenen Anbieter wurden nach den Anforderung der auszuführenden Applikation geclustert, um so das am besten passende Angebot zu finden.
- Cui [Cui21] verbesserte die Klick- und Kaufraten in einem Online-Shop durch die Kombination von Assoziationsregeln mit Fuzzy Clustering. Dadurch konnten passende Warenkörbe schneller identifiziert und vorgeschlagen werden.
- Zuletzt sei die Arbeit von Kumar et al. erwähnt [KRRT01]. Sie näherten sich mit einem möglichst allgemeinen mathematischen Modell an das Thema. Auf dieser Basis identifizierten sie Faktoren, die ein Empfehlungsalgorithmus erfüllen sollte. Werden diese befolgt, so können bereits mit wenigen Daten über das Nutzerverhalten effektive Empfehlungen gegeben werden. Das Clustern der Produkte beschleunigt diesen Prozess, da Erkenntnisse zu einem Produkt auf die gesamte Gruppe übertragbar sind.

Die letzte Idee beschäftigte sich mit der Suche von Produkten. Ein sehr spannendes Paper hierzu verfassten Kou und Lou [KoLo12]. Sie verbesserten die Klickraten einer Websuchmaschine durch die Nutzung von Clustering. Ihr Algorithmus nutzte die am meisten geklickten Webseiten zu jedem Suchbegriff als initiale Schwerpunkte für die Cluster. Anschließend sind die übrigen Webseiten rund um diese Schwerpunkte geclustert worden. Dadurch konnten neu hinzugekommene Seiten direkt in die Suche integriert werden. Außerdem erhöhte sich die Relevanz der gefundenen Suchergebnisse.

Davon abgesehen gibt es für diesen Bereich kaum Arbeiten, in denen Suchfunktionen mit Clustering kombiniert worden sind. Das liegt sicherlich daran, dass klassische probabilistische Suchalgorithmen (der bekannteste ist der BM25 [RoZa09]) sehr gute Ergebnisse liefern. Ansätze wie von Kou und Lou sind in dem Zusammenhang eher als experimentell anzusehen.

## 1.3 Kernfrage und Ablauf

Parallel zu der beschriebenen Recherche ist ebenfalls ein Blick in typische PIM-System wie z.B. Akeneo-PIM geworfen worden. Ziel war es zu verstehen, wie die Produktdaten in solchen System typischerweise vorliegen. Dies ist nötig, um einzuschätzen, ob und wie das Clustering der Produktdaten in solchen Systemen möglich ist.

Zum einen weisen die Produktdaten einen hohen Grad an Struktur auf. Die Attribute sind fest definiert. Umfangreiche Constraints sorgen für die Einhaltung des definierten Schemas. Zum anderen kommen sehr verschiedenartige Datentypen vor wie z.B. Textfelder, Einfach- und Mehrfachauswahl, numerische Daten mit verschiedenen Einheiten. Zudem sind viele Attribute lokalisierbar und weisen je nach Sprache andere Werte auf. Die klassische Clusteranalyse arbeitet allerdings nur mit numerischen Daten [King15, Kap. 1.2 Capturing the Clusters]. Somit war abzusehen, dass das Clustern selbst kein einfacher Prozess sein wird.

Dadurch verschob sich der Fokus der Arbeit. Bevor mögliche Anwendungen der Clusteranalyse evaluiert werden, stellt sich zuerst die Frage, ob so vielfältige Produktdaten überhaupt geclustert werden können. Die Kernfrage der Arbeit lautet also:

*Ist das effektive Clustern hoch-komplexer Produktdaten möglich – und wenn ja, wie?*

Diese Fragestellung war der Ausgangspunkt weiterer Recherchen, Experimente und Versuche. Sie werden in den folgenden Kapiteln erläutert.

*Kapitel 2* gibt einen Überblick über die wissenschaftlichen Grundlagen und Erkenntnisse des Themenfeldes.

Auf dieser Basis wird in *Kapitel 3* eine Konzeption dargelegt, wie das Clustern durchgeführt werden kann. Ebenso werden Versuche geplant, die zur Überprüfung des Konzeptes dienen.

*Kapitel 4* beschreibt die Umsetzung des dargelegten Konzeptes. Thema dabei sind die Implementierung nötiger Software-Komponenten, Tools und Techniken, die zur Überprüfung der Fragestellung angewendet worden sind.

Die Auswertung und Evaluation der Versuche erfolgt in *Kapitel 5*.

Im finalen *Kapitel 6* wird das Fazit gezogen, ob die Kernfrage beantwortet werden kann. Ebenso wird ein Ausblick für weitere Versuche und Fragestellungen in der Zukunft gegeben.

## 2 Clusteranalyse

### 2.1 Begriff und Einordnung

King definiert die *Clusteranalyse* als die “[...] Generierung eines Klassifizierungsschemas, welches Individuen in eine feste Anzahl an Gruppen einteilt, so dass sich die Individuen innerhalb einer Gruppe auf eine Art und Weise ähnlich sind und unähnlich denen in anderen Gruppen” [King15, Kap. 1.1 What Is a Cluster?]. Diese Gruppen werden auch als Cluster bezeichnet.

Dieser Prozess des Clusterings ist eine Methode des *unüberwachten Lernens* (*unsupervised learning*) – einem Teilgebiet des maschinellen Lernens. Papp et al. schreiben dazu: “Machine Learning beschäftigt sich mit der Entwicklung von Methoden, die Modelle von Zusammenhängen aus Daten erlernen, anstatt sie *per Hand* zu implementieren” [PWMO19, Kap. 5 Statistik-Grundlagen]. Ferner geben sie an, dass die Unterschiede zur Statistik fließend sind [PWMO19, Kap. 5 Statistik-Grundlagen]. Unüberwachtes Lernen bedeutet dabei, dass die verwendeten Daten nicht im Vorhinein gekennzeichnet sind (Unlabeled Data). Stattdessen werden Ähnlichkeiten und Muster in den Daten selbst gesucht ohne eine vorgegebene Zielgröße. Häufig dienen diese Analysen als erste Schritte der *Data Exploration*, aus denen im Anschluss neue Erkenntnisse und Anwendungen abgeleitet werden. [PWMO19, Kap. 5.2.3 Unüberwachtes Lernen]

Allgemein wird die Clusteranalyse als eine Form des Data Minings gesehen. Laut Bissantz und Hagedorn “beschreibt [Data Mining] die Extraktion implizit vorhandenen, nicht trivialen und nützlichen Wissens aus großen, dynamischen, relativ komplex strukturierten Datenbeständen” [BiHa09]. Mit “Wissen” meinen sie dabei in den Daten implizit enthaltene Muster, welche für den Anwender interessant sind und mit einer hinreichenden Wahrscheinlichkeit tatsächlich in den Daten existieren. [BiHa09]

### 2.2 Notation

Bevor konkrete Algorithmen und Verfahren der Clusteranalyse vorgestellt werden, müssen einige grundlegende Sachverhalte und ihre mathematische Definition geklärt werden:

Die Objekte, welche es zu clustern gilt, werden in dieser Arbeit auch als *Datenpunkte* oder *Produkte* (da dies die Anwendung im Praxisteil ist) bezeichnet. Mathematisch werden sie i.d.R. mittels  $x$  definiert. Diese Produkte bestehen aus einer festen Menge an Attributen (z.B. Farbe, Gewicht etc.). Mathematisch ist so ein Datenpunkt als ein Vektor definiert und jedes Element im Vektor steht für die Ausprägung eines spezifischen Attributes. Mittels Superskript werden spezifische Attribute eines Datenpunktes angesprochen.  $x^i$  bezeichnet also das  $i$ -te Attribut von  $x$ .

$$x = (a|a \text{ is attribute}) \quad (2.1)$$

Die Produkte oder Datenpunkte sind Teil eines *Datensets*  $X$ . Hierbei handelt es sich um eine simple Menge dieser Datenpunkte  $x \in X$ . Mittels Subskript werden einzelne Punkte des Sets angesprochen (z.B.  $x_1, x_2$  oder  $x_i$ ). Wenn nicht anders angegeben, gilt  $n = |X|$ .

Als *Cluster*  $C$  wird nun eine Teilgruppe aus dem Datenset bezeichnet:  $C \subset X$ . Allerdings können Cluster auch "hierarchisch" strukturiert sein. D.h. ein Cluster kann sowohl Datenpunkte als auch andere Cluster enthalten.

$$C = \{e|e = C_i \vee e \in X\} \quad (2.2)$$

Die *Menge aller Cluster* wird mit  $K = \{C_1, C_2, \dots\}$  bezeichnet. Dabei wird die Anzahl der Cluster mit  $k = |K|$  angegeben.

Schlussendlich gibt es noch das *Clustering-Ergebnis*  $Y$ . Dies ist ein Vektor der Länge  $n$ , welches jeden Datenpunkt im Set einem entsprechenden Cluster zuordnet.

$$Y = (l|l_i = \text{label for } x_i \in X) \quad (2.3)$$

## 2.3 Distanz- und Ähnlichkeitsmaße

### 2.3.1 Definition

Die Forderung, dass Objekte in gleichen Clustern sich "ähnlich" sein sollen und unähnlich zu den Objekten in anderen Clustern, muss in irgendeiner Form quantifiziert werden. Dies erfolgt über die Berechnung der "Nähe" (engl. proximity) der Objekte zueinander. Dazu werden sog. Abstands- bzw. Distanzmaße verwendet. [KaRo09, Kap. 1.2 Types of Data and How to Handle Them]

Mathematisch wird dieses Distanzmaß mittels der Funktion  $d(x_1, x_2)$  (engl. distance) definiert, welche die Distanz zwischen zwei Datenpunkten  $x_1$  und  $x_2$  als skalaren Wert zurückgibt. Zusätzlich geben Kaufmann und Rousseeuw folgende Eigenschaften für  $d$  an [KaRo09, Kap. 1.2.1 Interval-Scaled Variables]:

Tabelle 2.1: Eigenschaften der Abstandfunktion  $d$

1.	$d(x_1, x_2) \geq 0$	Distanzen sind stets positiv
2.	$d(x_1, x_1) = 0$	zwei gleiche Objekte haben immer einen Abstand von 0
3.	$d(x_1, x_2) = d(x_2, x_1)$	die Distanzfunktion ist kommutativ bzw. symmetrisch
4.	$d(x_1, x_3) \leq d(x_1, x_2) + d(x_2, x_3)$	Distanzen geben stets den kürzesten Weg an

Anstatt des Abstandes kann alternativ auch die Ähnlichkeit zweier Objekte berechnet werden. Solche Ähnlichkeitsmaße  $s(x_1, x_2)$  (engl. similarity) sind häufig im Intervall  $[0; 1]$  angegeben, wobei die 1 maximale Ähnlichkeit angibt. Ist nun eine entsprechende Distanzfunktion z.B. durch Normalisierung ebenfalls im Intervall  $[0; 1]$  definiert, so können Distanzen und Ähnlichkeiten einfach ineinander überführt werden [KaRo09, Kap. 1.2.3 Similarities]:

$$d(x_1, x_2) = 1 - s(x_1, x_2) \quad (2.4)$$

Distanzen und Ähnlichkeiten sind dadurch beliebig austauschbar. Daher ist dieses Vorgehen stets empfehlenswert.

Es existieren eine Vielzahl an Distanz- und Ähnlichkeitsmaßen. Welche zur Anwendung kommen, hängt maßgeblich von den Attributen und vor allem den Attribut-Typen ab, aus denen die Datenpunkte bestehen. [KaRo09, Kap. 1.2 Types of Data and How to Handle Them]

### 2.3.2 Numerische Attribute

Numerische Attribute sind im Allgemeinen alle Arten von (rationale) Zahlen [KaRo09, Kap. 1.2 Types of Data and How to Handle Them] mit stetigen (engl. continuous) Werten [Hua98]. Es kann sich dabei sowohl um Intervall- als auch um Verhältnisskalen handeln (engl. interval and ratio data). [Bosl12, Kap. 1 Basic Concepts of Measurement]

Datenpunkte, die nur aus numerischen Attributen bestehen, sind die klassischste Form von Daten, mit denen Clusteranalyse durchgeführt wird. Die meisten Algorithmen und Verfahren sind damit speziell auf diesen Attribut-Typ ausgelegt. [King15, Kap. 1.2 Capturing the Clusters]

#### Minkowski-Familie

Der Abstand zwischen reinen numerischen Vektoren wird mittels Maßen aus der sog. Minkowski-Familie berechnet. Das älteste dieser Metriken ist der "euklidische Abstand" (auch pythagoreische Metrik genannt), welcher den Abstand zwischen zwei Punkten über eine gerade Linie zwischen diesen beiden bestimmt. Im 19. Jahrhundert verallgemeinerte Hermann Minkowski dieses Maß zu einer ganzen Familie. [Cha07]

Tabelle 2.2: Übersicht über die gängigen Maße aus der Minkowski-Familie [Cha07]

Name	$p$ -Norm	Formel
allgemeine Form		$d(x_1, x_2) = \sqrt[p]{\sum_{i=1}^n  x_1^i - x_2^i ^p}$
Manhattan-Abstand	$p = 1$	$d(x_1, x_2) = \sum_{i=1}^n  x_1^i - x_2^i $
euklidischer Abstand	$p = 2$	$d(x_1, x_2) = \sqrt{\sum_{i=1}^n  x_1^i - x_2^i ^2}$
Chebyshev-Abstand	$p = \infty$	$d(x_1, x_2) = \max  x_1^i - x_2^i $

Die allgemeine Form ist über den Ausdruck  $p$  parametrisiert. Die konkreten Vertreter der Familie kommen nun durch ein spezifisch gewähltes  $p$  zustande. Daher wird bei diesen Maßen

auch von der sog.  $p$ -Norm gesprochen. Eine höhere  $p$ -Norm steht dabei im Allgemeinen für eine robustere bzw. eindeutigere Bestimmung des Abstandes. [King15, Kap. 12.3 Which Proximity Measure Should Be Used?]

Dies sind bei weitem nicht alle Arten von Maßen für numerische Attribute. Cha [Cha07] hat eine Vielzahl solcher Maße zusammengetragen, analysiert und miteinander verglichen. Dabei kam er zu folgenden Erkenntnissen:

Die meisten dieser Maße sind Abwandlungen eines Vertreters aus der Minkowski-Familie bzw. lassen sich auf einen solchen abbilden (z.B. basieren einige Maße auf der Pearson-Korrelation, welche ihrerseits wiederum auf dem euklidischen Abstand beruht). In seinen anschließenden Versuchen zeigte er, dass Maße, die auf einander basieren, auch stets ähnliche Abstände liefern. [Cha07]

Daraus lässt sich ableiten, dass die Vertreter der Minkowski-Familie ausreichend sind, um den Abstand numerischer Vektoren zu bestimmen.

## Vorverarbeitung

Häufig macht eine Vorverarbeitung der numerischen Attribute Sinn, um verschiedene Arten von Problemen zu vermeiden. [KaRo09, Kap. 1.2.1 Interval-Scaled Variables]

**Logarithmische Attribute** Die Abstandsmaße der Minkowski-Familie gehen i.d.R. von linearen Verhältnissen der Zahlen zueinander aus (Punkte mit einem Abstand von 2 sind doppelt so weit entfernt, wie Punkte mit einem Abstand von 1). Es kann aber vorkommen, dass die vorliegenden Daten eigentlich einen logarithmischen Zusammenhang aufweisen (der Abstand von 10 zu 20 ist exakt gleichbedeutend mit dem Abstand von 100 zu 200). Damit die Abstandsmaße richtig arbeiten können, sollten solche Attribute vorher in eine lineare Abbildung überführt werden [KaRo09, Kap. 1.2.5 Nominal, Ordinal, and Ratio Variables]. Dies wird folgendermaßen bewerkstelligt:

$$x'^i = \log x^i \quad (2.5)$$

**Standardisierung** Weitere Probleme können entstehen, wenn die absoluten Werte in den verschiedenen Dimensionen der Vektoren stark voneinander abweichen. Angenommen das Datenset besteht aus zwei-dimensionalen Vektoren. Die Werte in der ersten Dimension liegen im Bereich  $[0; 1]$  und in der zweiten im Bereich  $[100; 200]$ . Dadurch würden die Unterschiede zwischen den Werten in der zweiten Dimension viel stärker ins Gewicht fallen, da hier die Differenzen zwischen den Werten bis zu 100 Mal höher ausfallen. Um eine Gleichgewichtung zu erzeugen, sollten diese Attribute jeweils auf die gleiche Art standardisiert werden. [KaRo09, Kap. 1.2.1 Interval-Scaled Variables]

Milligan und Cooper [MiCo88] haben eine Vielzahl solcher Methoden der Standardisierung analysiert und experimentell miteinander verglichen.

$$x'^i = \frac{x^i - \text{avg } X^i}{\sigma(X^i)} \quad (2.6)$$

Die Formel 2.6 zeigt die Standardisierung mittels des klassischen z-Wertes. Die auftretenden Werte werden dadurch auf die 0 zentriert und mit einer Varianz von 1 drumherum verteilt. Dieses Maß ist besonders anfällig für Outliers (stark von der Hauptmenge abweichende Werte). In Milligan und Coopers Versuchen schnitt es daher eher mittelmäßig ab. [MiCo88]

$$x'^i = \frac{x^i - \text{avg } X^i}{\text{mad}(X^i)} \quad (2.7)$$

$$\text{mad}(X^i) = \frac{1}{n} \sum_j |x_j^i - \text{avg } X^i| \quad (2.8)$$

Kaufmann und Rousseeuw empfehlen eine modifizierte Form des z-Wertes nach Hartigan [Hart75]. Statt durch die Standardabweichung  $\sigma(X^i)$  zu teilen, wird durch die mittlere absolute Abweichung (engl. *mean absolute deviation*)  $\text{mad}(X^i)$  dividiert. Dieses Vorgehen soll den Einfluss von Outliern dämpfen. [KaRo09, Kap. 1.2.1 Interval-Scaled Variables]

$$x'^i = \frac{x^i - \min X^i}{\max X^i - \min X^i} \quad (2.9)$$

Eine weitere Methode ist die Skalierung aller Werte auf das Intervall  $[0; 1]$ . Obwohl Outliers bei dieser Methode ebenfalls die restlichen Werte in einen kleineren Bereich komprimieren, war dieser Ansatz in Milligan und Coopers Versuchen äußerst robust und lieferte über alle Versuche akzeptable Ergebnisse. [MiCo88]

Kaufmann und Rousseeuw geben zu bedenken, dass die unterschiedlichen Gewichtungen der Werte echte Zusammenhänge der realen Welt widerspiegeln könnten. Durch Standardisierung gehen diese entsprechend verloren. Weiterhin ist sogar denkbar, bewusst einige Attribute stärker zu gewichten als anderen (indem sie mit entsprechenden Gewichten multipliziert werden). Dies sind allerdings eher speziellere Fälle, in denen im Vorfeld bereits Domänenwissen zum Datenset vorliegen muss. [KaRo09, Kap. 1.2.1 Interval-Scaled Variables]

### 2.3.3 Kategorische Attribute

Kategorische Attribute bilden die zweite große Gruppe von Attribut-Typen. Sie zeichnen sich dadurch aus, dass als Wertausprägungen Labels verwendet werden, die verschiedene Kategorien repräsentieren. Diese Labels können in Form von Zeichenketten (z.B. “red”, “green”, “blue”), Bool-Werten (*true/false*, Ja/Nein, 0/1) oder auch diskreten Ziffern auftreten. Es ist eine weitere Unterscheidung in die zwei Teilgruppen ordinale und nominale Attribute möglich. [Bos12, Kap. 5 Categorical Data]

#### Ordinale Attribute

Lassen sich die vorkommenden Labels eines Attributes in eine sinnvolle Reihenfolge bringen (z.B. “XS”, “S”, “M”, “L”, “XL”), so werden sie als *ordinal* bezeichnet. Wichtig dabei ist, dass nur eine sinnvolle Reihenfolge besteht. Eine Aussage über die Intervalle zwischen den Labels (“XS” zu “S” ist genauso weit wie “M” zu “L”) oder über deren quantitativen

Verhältnisse ( $3 \cdot \text{“XS”} = 2 \cdot \text{“S”}$ ) macht aber keinen Sinn. [Bos12, Kap. 1 Basic Concepts of Measurement]

Für die Berechnung des Abstandes gibt es grundsätzlich zwei Möglichkeiten: Zum einen können ordinale Attribute genauso wie nominale verarbeitet werden (siehe nächster Abschnitt). Zum anderen kann die Reihenfolge der Labels verwendet werden, um die Labels in numerische Werte umzuwandeln. Anschließend kann der Abstand mittels der Maße für numerische Attribute berechnet werden. [KaRo09, Kap. 1.2.5 Nominal, Ordinal, and Ratio Variables]

Die Umwandlung erfolgt dabei so, dass die Labels entsprechend ihrer natürlichen Reihenfolge beginnend mit 1 nummeriert werden. Der finale numerische Wert berechnet sich dann wie folgt:

$$x'^i = \frac{x^i - 1}{|\{x \in X^i\}| - 1} \quad (2.10)$$

$|\{x \in X^i\}|$  ist dabei die Menge an auftretenden Labels. Dadurch liegen alle Werte erneut im Intervall  $[0; 1]$ , welches in  $|\{x \in X^i\}| - 1$  gleichmäßige Schritte unterteilt ist. [KaRo09, Kap. 1.2.5 Nominal, Ordinal, and Ratio Variables]

### Nominale Attribute

Die Verarbeitung von nominalen Attributen geschieht über die vorherige Umwandlung in sog. binäre Attribute. Diese Attribute sind dadurch gekennzeichnet, dass sie nur zwei verschiedene Werte annehmen können (i.d.R. 0/1 bzw. *false/true*) [KaRo09, Kap. 1.2.5 Nominal, Ordinal, and Ratio Variables]. Hier gilt es zwei Fälle zu unterscheiden:

Kommen in den Labels eines Attributs maximal zwei verschiedene Werte vor (z.B. Ja/Nein, schwarz/weiß etc.), so werden diese in ein einzelnes binäres Attribut umgewandelt. Solche Attribute werden auch als *symmetrische* binäre Attribute bezeichnet. [KaRo09, Kap. 1.2.4 Binary Variables]

Bei mehr als zwei Wertausprägungen gibt es hingegen wiederum zwei Möglichkeiten:

- Zum einen können die vorkommenden Labels auf genau zwei verschiedene Wertausprägungen komprimiert werden, sodass wieder ein einzelnes binäres Attribut verwendet werden kann. Dieser Vorgang geht aber mit einem Informationsverlust einher. [KaRo09, Kap. 1.2.5 Nominal, Ordinal, and Ratio Variables]
- Zum anderen kann für jedes vorkommende Label ein eigenes binäres Attribut angelegt werden. Das Attribut, welches das tatsächlich vorkommende Label eines Datenpunktes repräsentiert, wird entsprechend auf *true* gesetzt und alle anderen auf *false*. Diese Umwandlung resultiert in sog. *asymmetrische* binäre Attribute. [KaRo09, Kap. 1.2.4 Binary Variables]

Kennzeichnend für die symmetrischen binäre Attribute ist, dass beide Wertausprägung eine konkrete Eigenschaft des Objektes darstellen (z.B. 0 = schwarz und 1 = weiß). Asymmetrische zeigen hingegen nur das Vorhanden- oder Nicht-Vorhandensein einer bestimmten



Eigenschaft an und müssen stets im Kontext der gesamten Gruppe betrachtet werden. Für beide Arten gibt es jeweils eigene Distanzmaße. [KaRo09, Kap. 1.2.4 Binary Variables]

**Simple Matching** Das sog. Simple Matching ist für die Bewertung der Ähnlichkeit von Objekten mit ausschließlich symmetrischen binären Attributen geeignet [KaRo09, Kap. 1.2.4 Binary Variables]. Es berechnet sich wie folgt:

$$s(x_1, x_2) = \frac{|x_1 \cap x_2| + |\bar{x}_1 \cap \bar{x}_2|}{|x_1 \cup \bar{x}_1 \cup x_2 \cup \bar{x}_2|} \quad (2.11)$$

Verglichen werden also die Menge an Attributen mit der gleichen Wertausprägung in beiden Objekten (beide *true*:  $|x_1 \cap x_2|$  und beide *false*:  $|\bar{x}_1 \cap \bar{x}_2|$ ) im Verhältnis zur Anzahl aller Attribute. Es handelt sich damit um ein Ähnlichkeitsmaß im Intervall  $[0; 1]$ , wobei die 1 für perfekte Übereinstimmung steht. [KaRo09, Kap. 1.2.4 Binary Variables]

Es gibt von diesem Verfahren eine Reihe von Abwandlungen, welche die Ausprägungen bestimmter Wertekombinationen gesondert gewichten. Alle diese Varianten basieren im Kern auf dem Simple Matching und liefern stets ähnliche Ergebnisse. [KaRo09, Kap. 1.2.4 Binary Variables]

**Jaccard-Koeffizient** Für die Bewertung der Ähnlichkeit von Objekten mit asymmetrischen binären Attributen eignet sich hingegen der Jaccard-Koeffizient besser [KaRo09, Kap. 1.2.4 Binary Variables]. Er berechnet sich wie folgt:

$$s(x_1, x_2) = \frac{|x_1 \cap x_2|}{|x_1 \cup x_2|} \quad (2.12)$$

Im Gegensatz zum Simple Matching werden hier alle Attribute, in denen beide Objekte den Wert *false* aufweisen, ignoriert. Die Übereinstimmung wird nur mit den Attributen berechnet, wo wenigstens eines von beiden Objekten den Wert *true* aufweist. [KaRo09, Kap. 1.2.4 Binary Variables]

*Hintergrund:* Angenommen ein kategorisches Attribut weist im Datenset bis zu drei verschiedene Labels auf ("red", "green" und "blue").  $x_1$  hat das Label "red" und  $x_2$  hat "blue". Nach der Umwandlung in binäre Werte wäre  $x_1 = (1, 0, 0)$  und  $x_2 = (0, 0, 1)$ . Das Simple Matching würde nun eine Ähnlichkeit von  $\frac{1}{3}$  berechnen, da beiden Objekten gemeinsam ist, dass sie *nicht* "green" sind. Bei asymmetrischen binären Attributen sind aber zwangsläufig mit steigender Anzahl an vorkommenden Labels die meisten Werte 0 bzw. *false*. Dies treibt das Ergebnis vom Simple Matching künstlich in die Höhe, macht aber inhaltlich in diesem Kontext keinen Sinn. [KaRo09, Kap. 1.2.4 Binary Variables]

### 2.3.4 String-Attribute

Manche Autoren (z.B. [RaRa11]) definieren eine weitere Gruppe sog. String-Attribute. Streng genommen sind dies ebenfalls nominale bzw. kategorische Attribute. Sie zeichnen sich aber durch eine hohe Variation an Werten aus, sodass häufig keine einzige Wertausprägung mehr als einmal vorkommt. Ein klassisches Beispiel wäre der Produkttitel. Es wird i.d.R.

keine zwei Produkte mit dem gleichen Titel geben. Dies führt die Betrachtung dieser Attribute als “kategorisch” ad-absurdum, wenn jeder Datenpunkt seine eigene individuelle Kategorie definiert. Sollen diese Werte dennoch auf ihre Ähnlichkeit geprüft werden, so gibt es eine Reihe verschiedener Vorgehensweisen.

## String-Metrics

Eine Methode besteht aus dem Vergleich übereinstimmender Buchstaben zwischen zwei String-Werten. Dazu existieren eine Reihe sog. “String-Metriken” von verschiedenen Autoren. [CRF03]

Hamming [Hamm50] entwickelte 1950 ein Verfahren zur Fehlererkennung binärer Codes. Es bildet die Grundlage des sog. *Hamming-Abstands*, welches sich später in abgewandelter Form im Bereich der String-Metrics als Distanzmaß etablierte [RaRa11]. Der Hamming-Abstand ist die Anzahl an Buchstaben zwischen zwei Zeichenketten gleicher Länge, welche **nicht** identisch sind. Er kann nur zum Vergleich von Strings mit exakt identischer Anzahl an Buchstaben verwendet werden und liefert Abstände von 0 bis  $n$  (Anzahl an Buchstaben im jeweiligen String).

Ein Verfahren, welches auch für Strings unterschiedlicher Länge geeignet ist, ist der *Levensthein-Abstand*. Zur Berechnung wird gezählt, wie viele Buchstaben eingefügt, entfernt oder umgewandelt werden müssen, um einen String in den anderen umzuwandeln. Dabei wird stets die kleinstmögliche Zahl an Ersetzungen bevorzugt. Die Ergebnisse liegen damit ebenfalls zwischen 0 (geringster Abstand, da kein einziger Schritt zur Umwandlung nötig ist) und  $n$ , was die Anzahl an Buchstaben des längeren Strings ist. Der Nachteil dieses Verfahrens ist die recht aufwendige Berechnung, da stets mehrere Arten der Buchstabentransformationen durchgerechnet werden müssen, um die kürzeste Kette an Transformationen zu finden. [RaRa11]

Mit deutlich weniger Rechenaufwand kommt die *Jaro-Metrik* [Jaro95] aus. Sie berechnet sich aus der Menge an Buchstaben, welche in beiden Wörtern in der gleichen Reihenfolge vorkommen, zu der Menge aller Buchstaben je Wort (Dies ist eine verkürzte Darstellung. Siehe [Jaro95] oder [CRF03] für die genaue Berechnung). Das Ergebnis gibt die Ähnlichkeit beider Strings zwischen 0 (komplett unähnlich) und 1 (maximale Übereinstimmung). [CRF03]

Eine Variante des eben erklärten Maßes ist die *Jaro-Winkler-Metrik*, welcher eine etwas höhere Ähnlichkeit feststellt, wenn die ersten Buchstaben (Prefix) identisch sind. [CRF03]

In den Versuchen von Cohen [CRF03] lieferte der Jaro-Winkler-Abstand die besten Ergebnisse bei der Verwendung für das Clustering.

## Tokenization

Strings-Metrics sind eher für den Vergleich kürzerer Zeichenketten geeignet. Für ganze Textzeilen und -blöcke werden eher Verfahren aus dem Document Retrieval angewandt. [StKaKu00]

Wesentlich im Document Retrieval ist die Zerlegung eines Textes in die einzelnen Wörter, aus denen sich der Text zusammensetzt. Da in vielen Sprachen das gleiche Wort je nach Fall

verschiedene Endungen (engl. stemming) aufweisen kann, werden diese Endungen häufig entfernt (z.B. mit Verfahren wie Porters Stemming Algorithmus). Ebenso weisen alle Sprachen kleine Füllwörter (z.B. Artikel, Bindewörter etc.) auf, die sehr häufig in eigentlich jedem Text vorkommen. Auch diese werden i.d.R. entfernt (engl. stop-word removal), da sie keine Relevanz für die Unterscheidung der Texte haben. Das Ergebnis ist eine Menge aus den Grundwörtern eines Textes, welche auch “Tokens” genannt werden. [StKaKu00; und CRF03]

Nach dieser “Tokenization” können zwei Texte z.B. auf ihre Ähnlichkeit überprüft werden, indem man die auftretenden Token beider Texte miteinander vergleicht. Für die Berechnung können bereits erklärte Metriken wie der Jaccard-Koeffizient eingesetzt werden. Die Tokens repräsentieren dann “Kategorien” oder Themenfelder der Texte. [CRF03]

Eine andere Möglichkeit besteht darin, die Tokens in numerische Vektoren umzuwandeln. Jedes Element im Vektor repräsentiert eines der Token und der Wert des Elementes entspricht der Anzahl, wie häufig dieses Token im ursprünglichen Text vorkommt. Dieser Vorgang wird auch als Transformation der Texte in das “Vector-Space-Model” bezeichnet. [StKaKu00]

Diese numerischen Vektoren können nun wie numerische Attribute verarbeitet werden. Allerdings existieren auch gesonderte Metriken speziell für den Bereich des Document Retrieval. Bspw. wird statt des euklidischen Abstandes häufig der “cosine”-Abstand zur Bestimmung der Distanz verwendet. [StKaKu00]

$$\text{cosine}(x_1, x_2) = \frac{|x_1 \cdot x_2|}{|x_1| \cdot |x_2|} \quad (2.13)$$

### 2.3.5 Gemischte Attribute

In einem Datenset können nicht nur Attribute des gleichen Types (z.B. ausschließlich numerisch) vorkommen. Gibt es Attribute mit unterschiedlichen Typen, so spricht man von gemischten Daten oder gemischten Attributen. Der Umgang mit diesen ist sehr vielfältig. [KaRo09, Kap. 1.2.6 Mixed Variables]

#### Umwandlung aller Attribute in den selben Typ

Eine der einfachsten Varianten besteht darin, alle Attribute in den gleichen Typ umzuwandeln. [KaRo09, Kap. 1.2.6 Mixed Variables]

**Umwandlung in numerische Werte** Für ordinale Attribute ist diese Umwandlung bereits beschrieben worden. Nominale Attribute können in eine künstliche Reihenfolge gebracht werden, wodurch sie ihren ordinalen Brüdern gleichen und entsprechend umgewandelt werden können. Diese Transformation verzerrt u.U. die Daten, da nun z.B. “rote” und “grüne” Produkte als ähnlicher betrachtet werden als “rote” und “blaue” (aufgrund der künstlichen Reihenfolge). Nichtsdestotrotz kann dieses Verfahren gute Ergebnisse produzieren. [KaRo09, Kap. 1.2.6 Mixed Variables]

Die Transformation von String-Attributen in numerische Vektoren ist ebenfalls bereits dargestellt worden.

**Umwandlung in kategorische Werte** Auch der umgekehrte Weg ist denkbar. Ordinale Attribute können als nominal betrachtet und verarbeitet werden. Ebenso können numerische Werte durch Diskretisierung in feste Gruppen oder Bänder eingeteilt werden (z.B. in Altersspannen “unter 18”, “18-35” usw.). [KaRo09, Kap. 1.2.6 Mixed Variables]

Diese Umwandlung geht aber mit einem massiven Verlust an Informationen einher. Bspw. brachte die Diskretisierung numerischer Attribute in den Versuchen von Milligan und Cooper die schlechtesten Clustering-Ergebnisse. [MiCo88]

### Subspace-Clustering

Ein ganz anderer Ansatz besteht darin, das Clustering mehrmals durchzuführen, jeweils mit jedem Attribut-Typen separat. Anschließend werden die Ergebnisse verglichen und entsprechende Erkenntnisse abgeleitet. Allerdings ist hier der Interpretationsspielraum sehr groß, sodass i.d.R. andere Arten des Vorgehens empfehlenswerter sind. [KaRo09, Kap. 1.2.6 Mixed Variables]

Dieser Ansatz wird auch als *Subspace-Clustering* bezeichnet und kann sehr weit ausgeführt werden. So werden teilweise mehrmalige Clusterings mit verschiedenen Kombinationen an Attributen (beliebigen Types) durchgeführt und bewertet. Solche Ansätze sind sehr aufwendig sowohl in Sachen Laufzeit als auch in der Interpretation der eigentlichen Ergebnisse. [JiCh17]

### Kombinierte Distanzfunktion

Schließlich besteht die Möglichkeit, verschiedene Distanzmaße zu kombinieren. Die verschiedenen Attribute also mit einer Metrik ihres Types zu verrechnen und diese Ergebnisse anschließend zu addieren. [KaRo09, Kap. 1.2.6 Mixed Variables]

Ein geläufiges Beispiel ist das kombinierte Distanzmaß des K-Prototypes-Algorithmus [Huan98], welches wie folgt definiert ist:

$$d(x_1, x_2) = d_{num}(x_1^{num}, x_2^{num}) + w \cdot d_{cat}(x_1^{cat}, x_2^{cat}) \quad (2.14)$$

Dieser Algorithmus kann numerische und kategorische Attribute gemeinsam verarbeiten. Die spezifischen Distanzmaße für jeden Attribut-Typ sind dabei frei wählbar (im Original: euklidisch für die numerischen und Simple Matching für die kategorischen). Der Faktor  $w$  soll Unterschiede in der Gewichtung der Distanzmaße ausgleichen. So hat der euklidische Abstand umso höhere Werte, je mehr numerische Attribute vorkommen. Das Simple Matching hingegen liegt immer im Intervall  $[0; 1]$ . [Huan98]

## 2.4 Clustering-Verfahren

### 2.4.1 Partitionierendes Clustering

#### Überblick

Partitionierende Verfahren betrachten die Clusterzuteilung als ein Minimierungsproblem. Auf verschiedene Arten starten die jeweiligen Verfahren mit einer initialen Gruppierung der Datenpunkte. Anschließend werden die Clusterzuteilungen wiederholt verändert bis ein lokales Optimum erreicht ist [King15, Kap. 4.1 Partition Clustering – Introduction]. Wichtig zu erwähnen ist dabei, dass jeder Datenpunkt nur genau einem Cluster zugeordnet werden kann [KaRo09, Kap. 1.3.1 Partitioning Methods]. Die konkreten Verfahren unterscheiden sich dabei in der Art der Initialisierung der Cluster, dem genauen Ablauf der Minimierung sowie in der Definition, wann ein Clustering “besseres” ist [King15, Kap. 4.1 Partition Clustering – Introduction]. Ebenso gibt es spezielle Verfahren für das Clustern von Daten, welche nicht ausschließlich numerisch sind. [Huan98]

Kennzeichnend für diese Verfahren ist außerdem, dass die Anzahl der erwarteten Cluster im Vorfeld bekannt sein muss. Dies ist zum Teil problematisch, wenn kaum Informationen zum Datenset vorliegen. Eine mögliche Lösung ist das mehrmalige Clustern mit verschiedenen Werten für die Clusteranzahl. Anschließend kann mittels verschiedener Metriken (siehe Abschnitt Cluster-Validität) eine adäquate Clusteranzahl aus den Ergebnissen ausgewählt werden. Ebenso könnte zu erst ein hierarchisches Clustering zur Bestimmung der Clusterzahl durchgeführt werden. [King15, Kap. 4.1 Partition Clustering – Introduction]

Diese Verfahren neigen dazu in lokalen Minima während des Clusterings hängen zu bleiben. Deswegen ist es empfehlenswert, das jeweilige Verfahren mehrmals mit unterschiedlichen Start-Partitionierungen durchzuführen und anschließend das “beste” Ergebnis aus diesen Durchläufen auszuwählen. Trotz der mehrfachen Durchläufe sind partitionierende Verfahren allgemein recht effizient in ihrer Laufzeit. Die Komplexität beträgt  $\mathcal{O}(n \cdot k \cdot l)$ .  $n$  ist die Anzahl an Datenpunkten,  $k$  die Menge an gesuchten Clustern und  $l$  die mehrmalige Wiederholung des Verfahrens.  $k$  und  $l$  sind typischerweise recht kleine Werte, die auch mit steigender Anzahl von Datenpunkten im Datenset nicht zunehmen. Daher handelt es sich unterm Strich um eine lineare Laufzeit von  $\mathcal{O}(n)$ . [Huan98]

#### K-Means

K-Means ist das bekannteste partitionierende Clustering-Verfahren und gilt häufig als Inbegriff der Clusteranalyse selbst. [Huan98]

Jedes Cluster wird in diesem Verfahren durch einen Mittelpunkt  $c$  repräsentiert. Dieser Mittelpunkt oder Centroid ist kein tatsächlicher Punkt des Datensets [StKaKu00], sondern wird aus den Punkten des Clusters berechnet. Konkret wird für jedes einzelne Attribut der Durchschnitt (engl. mean) aus den Werten der Cluster-Mitglieder für das jeweilige Attribut gebildet [King15, Kap. 4.5 K-Means Algorithm]:

$$c(C) = \frac{1}{|C|} \sum_{x_i \in C} x_i \quad (2.15)$$

Durch das Verändern der Clusterzuordnungen ändern sich entsprechend auch die Mittelpunkte der Cluster. Ziel dieser Verschiebung ist, dass die Summe der Abstände der Datenpunkte zu ihrem jeweiligen Mittelpunkt minimal wird. Der Abstand wird mit einem Distanzmaß für numerische Attribute (meist der euklidische Abstand) berechnet. [King15, Kap. 4.5 K-Means Algorithm]

Dieses Verfahren funktioniert also ausschließlich mit numerischen Attributen, da nur hier aus den Werten ein Durchschnitt berechnet werden kann. [Huan98]

Steinbach et al. [StKaKu00] beschreiben den typischen Ablauf des Verfahrens folgendermaßen:

1. Wahl von  $k$  Punkten aus dem Datenset als initiale Mittelpunkte
2. Zuordnung aller Datenpunkte des Datensets zum jeweils nächstgelegenen Mittelpunkt (mit dem geringsten Abstand)
3. Neuberechnung aller Mittelpunkte aus den Durchschnitten der zugeordneten Datenpunkte (wodurch sich diese verschieben)
4. Wiederholung der Schritte 2 und 3; solange bis sich die Clusterzuordnungen nicht mehr verändern oder ein Höchstlimit an Iterationen überschritten ist

Die Wahl der initialen Punkte kann z.B. durch die Wahl der  $k$  ersten oder letzten Punkte im Datenset geschehen. Ebenso möglich ist die Wahl jedes  $\frac{n}{k}$ -ten Punktes. Am geläufigsten ist wohl die komplett zufällige Wahl von Punkten aus dem Set. [King15, Kap. 4.3 The Initial Partition]

Zu diesem Basis-Verfahren gibt es eine große Vielfalt an Varianten und Abwandlungen. Bspw. muss die Neuberechnung der Mittelpunkte nicht erst am Ende eines Laufes durch das gesamte Datenset geschehen. Stattdessen können die Mittelpunkte auch unmittelbar, wenn ein neuer Punkt zum Cluster zugeordnet oder entfernt wird, aktualisiert werden. King [King15, Kap. 4.5 K-Means Algorithm] bezeichnet dieses alternative Vorgehen als "MacQueens-Methode". Viele weitere Varianten sind in der Literatur beschrieben worden. Im folgenden wird nur eine kleine relevante Auswahl dessen präsentiert.

## K-Medoids

Der K-Medoids ist eine Variante des K-Means-Algorithmus, bei dem die Centroids der Cluster tatsächliche Datenpunkte des Datensets sein müssen. Anstatt den Durchschnitt der jeweiligen Attribute aller Punkte im Cluster zu berechnen, werden die Abstände aller Punkte eines Clusters zueinander betrachtet. Der Punkt, welcher in der Summe den kleinsten Abstand zu allen anderen ausweist, muss zwangsläufig "mittig" im Cluster liegen. Fortan gilt dieser Punkt als Centroid des Clusters. Dieses Verfahren ist rechenaufwendiger als der normale K-Means. Allerdings konvergieren die Cluster häufig in weniger Durchläufen zu einem stabilen Optimum. Deshalb kann dieses Verfahren, auf den gesamten Ablauf betrachtet, dennoch Performance-Vorteile zum klassischen K-Means bringen. [ArVa16]

## K-Modes

Der K-Modes verarbeitet ausschließlich kategoriale Attribute – statt numerischen. Das Grundprinzip des K-Means bleibt bestehen. Als Distanzmaß muss eine entsprechend geeignete Metrik verwendet werden (im Original das Simple Matching). Außerdem erfolgt die Berechnung der Cluster-Mittelpunkte nicht über den Durchschnitt, da zwischen kategorischen Attributen kein Durchschnitt berechnet werden kann. Daher bilden sich die Centroids aus dem Modus (engl. mode), also dem am häufigsten auftretenden Label im jeweiligen Attribut. [Huan98]

Dieser Algorithmus verwendet außerdem den “MacQueens-Ansatz”. D.h. dass sich der Modus für jedes Attribut direkt aktualisiert, wenn ein Objekt dem Cluster neu zugeordnet oder entfernt wird. Auch hier wird dieser Prozess mehrmals durchlaufen, solange bis sich die Clusterzuteilungen nach einem kompletten Durchlauf durch das Datenset nicht mehr verändert haben. [Huan98]

## K-Prototypes

Der K-Prototypes kombiniert den K-Means- und K-Modes-Algorithmus zum Clustern gemischter Datensets. Die Centroids berechnen sich entsprechend aus dem Durchschnitt der Werte numerischer und aus dem Modus der Werte kategorischer Attribute. Die Distanz wird mittels einer kombinierten Distanzfunktion (siehe der entsprechende Abschnitt) berechnet. [Huan98]

## 2.4.2 Hierarchisches Clustering

### Überblick

Hierarchische Verfahren produzieren eine ineinander verschachtelte Partitionierung der Datenpunkte in Cluster. Dabei werden zwei grundlegende Arten von Verfahren unterschieden [StKaKu00; und King15, Kap. 3.1 Hierarchical Clustering – Introduction]:

- *Top-down- oder divisive Verfahren* beginnen mit einem Cluster, welches alle Datenpunkte enthält. Anschließend wird dieses Cluster in zwei kleinere Cluster aufgesplittet, danach werden die beiden resultierenden Cluster wiederum jeweils in zwei kleinere aufgeteilt usw. Dieses Verfahren wird solange fortgesetzt, bis nur noch Cluster übrig bleiben, die genau einen Datenpunkt enthalten.
- *Bottom-up- oder agglomerative Verfahren* starten exakt andersrum mit jeweils einem Cluster für jeden Datenpunkt. Anschließend werden immer wieder zwei Cluster zu einem größeren kombiniert, solange bis im letzten Schritt ein Cluster entsteht, welches wiederum alle Datenpunkte enthält.

Das bedeutet, dass jeder Datenpunkt mehreren verschiedenen Clustern zugeteilt ist, je nachdem welche Stufe in der Hierarchie betrachtet wird. Alternativ können die Ergebnisse auch als Baumstruktur in einem sog. Dendrogram dargestellt werden. [StKaKu00]

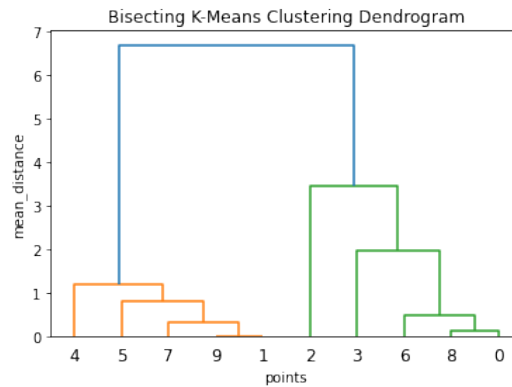


Abbildung 2.1: Bsp. eines Dendrograms

Ein großer Vorteil von diesem Verfahren ist, dass die Anzahl an gesuchten Clustern nicht vorher bekannt sein muss. Stattdessen wird eine Clusterzuteilung für alle möglichen Clusterzahlen ( $1 \leq k \leq n$ ) im Zuge der Berechnung ermittelt [KaRo09, Kap. 1.3.2 Hierarchical Methods]. Dadurch können die verschiedenen Hierarchie-Ebenen mit geeigneten Metriken (siehe Abschnitt Cluster-Validität) bewertet und anschließend z.B. eine geeignete Clusteranzahl für die Durchführung eines partitionierenden Verfahrens ausgewählt werden. [King15, Kap. 4.1 Partition Clustering – Introduction]

Weiterhin bringt die hierarchische Struktur vielfältige Möglichkeiten der Analyse. So kann nicht nur festgestellt werden, welche Objekte in verschiedene Cluster sortiert werden, sondern auch welche der Cluster sich ähnlicher sind (aufgrund ihrer Position in der Hierarchie). Das Ermitteln von Outliern wird ebenfalls möglich, da diese meistens recht früh in der Hierarchie in ihr eigenes Cluster platziert werden. [DoBi22]

Ein großer Nachteil dieser Verfahren ist, dass die Clusterzuteilungen im Laufe der Berechnung nicht mehr nachträglich korrigiert werden. Nachdem also zwei Datenpunkte in einem Top-Down-Verfahren in verschiedene Cluster gesplittet worden sind, werden sie in den nachfolgenden Hierarchie-Ebenen niemals wieder im selben Cluster landen (für Bottom-Up-Verfahren gilt das gleiche entsprechend umgekehrt). Dies führt u.U. zu einer ungeeigneten Clusterzuteilung je nach betrachteter Hierarchie-Ebene. [KaRo09, Kap. 1.3.2 Hierarchical Methods]

### Agglomeratives Clustering (Bottom-Up)

Wie bereits erklärt beginnt das agglomerative Clustering mit jedem Datenpunkt in einem eigenen Cluster. Anschließend werden stets die beiden “am nächsten gelegen” bzw. “ähnlichsten” Cluster zu einem größeren zusammengefasst. Dieser Prozess wird solange wiederholt, bis sich zum Schluss alle Datenpunkte in einem großen Cluster befinden. [StKaKu00]

Von zentraler Bedeutung dabei ist die Frage, wie zwei Cluster miteinander auf ihre Ähnlichkeit untersucht werden. Für die vorkommenden Datenpunkte wird dazu ein geeignetes Distanz- oder Ähnlichkeitsmaß gewählt. Die Berechnung der Ähnlichkeit zwischen zwei Clustern (welche ja potenziell aus mehreren Punkten bestehen können) geschieht dann mit Hilfe sog. Fusionsverfahren (engl. linkage functions). Ein bekannter Vertreter ist z.B. der “Single-Link”, welcher den Abstand zwischen zwei Clustern über den Abstand der



zwei am nächsten gelegenen Punkte der Cluster definiert [DoBi22]. Nachfolgend sind einige typische Vertreter aufgelistet:

**Single-Link**  $d(C_1, C_2) = \min_{x_i \in C_1, x_j \in C_2} d(x_i, x_j)$

Abstand wird durch die beiden dichtesten Punkte definiert; neigt zur Bildung von langen Ketten (chaining effect); anfällig gegen Outlier

**Complete-Link**  $d(C_1, C_2) = \max_{x_i \in C_1, x_j \in C_2} d(x_i, x_j)$

Abstand wird durch die beiden am entferntesten liegenden Punkte definiert; neigt zur Bildung von vielen sehr kleinen Clustern; anfällig gegen Outlier

**Average-Link**  $d(C_1, C_2) = \frac{1}{|C_1|} \frac{1}{|C_2|} \sum_{x_i \in C_1} \sum_{x_j \in C_2} d(x_i, x_j)$

durchschnittlicher Abstand aller Punkte der beiden Cluster zueinander; ähnlich dem Mittelpunkt von k-Means, aber viel aufwendiger, da immerzu jeder Punkt eines Clusters mit jedem anderen Punkt der anderen Clustern verglichen werden muss

**Centroid-Link**  $d(C_1, C_2) = d(\lfloor \frac{1}{|C_1|} \sum_{x_i \in C_1} x_i \rfloor, \lfloor \frac{1}{|C_2|} \sum_{x_j \in C_2} x_j \rfloor)$

Abstand der Mittelpunkte der Cluster zueinander (Berechnung des Mittelpunktes wie bei K-Means); ähnliches Verhalten wie Average-Link, aber effizienter zu berechnen, da der Mittelpunkt nur einmal pro Cluster berechnet werden muss

**Ward-Methode**  $d(C_1, C_2) = \frac{|C_1||C_2|}{|C_1|+|C_2|} d(\lfloor \frac{1}{|C_1|} \sum_{x_i \in C_1} x_i \rfloor, \lfloor \frac{1}{|C_2|} \sum_{x_j \in C_2} x_j \rfloor)$

berechnet die Zunahme der Varianz durch den jeweiligen Merge; neigt zur Bildung von Kreis- oder Kugelförmigen Clustern; recht effizient; anfällig gegen Outlier

Da für den Centroid-Link und die Ward-Methode die Mittelpunkte der Cluster gebildet werden müssen, sind diese Verfahren primär nur für numerische Datensets geeignet [DoBi22]. Die Verwendung dieser Verfahren für gemischte Attribute (z.B. durch Verwendung des Modus für kategoriale) ist in der Literatur nicht beschrieben worden.

Ein großer Nachteil dieser Verfahren ist die Laufzeit. Da im ersten Schritt jeder Datenpunkt mit jedem anderen verglichen werden muss (es sind ja zu Beginn  $n$  Cluster), beträgt die Laufzeit immer mindestens  $\mathcal{O}(n^2)$ . Durch geschickte Implementierung, kann diese Laufzeit auch für einige Verfahren (z.B. für Single- [Sibs73] oder Complete-Link [Defa77]) tatsächlich erreicht werden. Häufig liegen die Laufzeiten aber bei  $\mathcal{O}(n^2 \log n)$  oder sogar  $\mathcal{O}(n^3)$ . [DoBi22]

## Divisives Clustering (Top-Down)

**Überblick** Divisives Clustering startet mit allen Datenpunkten in einem großen Cluster. Anschließend werden immerwieder Zweier-Splits der Cluster durchgeführt, solange bis final jeder Datenpunkt in genau einem eigenen Cluster platziert ist. [KaRo09, Kap. 6 Divisive Analysis (Program DIANA)]

Lange Zeit wurde dieser Ansatz des hierarchischen Clusterings ignoriert, da das Verfahren auf den ersten Blick sehr rechenintensiv erscheint. Im aller ersten Zustand, wenn alle Punkte in einem Cluster sind, gibt es theoretisch  $2^n - 1$  verschiedene Möglichkeiten dieses große Cluster in zwei kleinere zu teilen. Eine solche exponentielle Laufzeit könnte niemals für Datensets auch nur mittlerer Größe verwendet werden. Allerdings ist es nicht nötig alle dieses Splits zu betrachten und durchzurechnen. Mittlerweile gibt es einige Ansätze für diese Art von Clustering-Verfahren, welche dieses Problem geschickt umgehen. [KaRo09, Kap. 6 Divisive Analysis (Program DIANA)]

Außerdem zeigt die Praxis, dass Top-down-Clustering i.d.R. bessere und intuitivere Cluster generiert. Der Grund ist vermutlich, dass wir Menschen beim Clustering per Hand automatisch eine Art Top-down-Denkweise an den Tag legen. [King15, Kap. 3.3 Agglomerative versus Divisive Clustering]

**DIANA** DIANA steht für “divisive analysis” und wurde von Kaufman und Rousseeuw [KaRo09, Kap. 6 Divisive Analysis (Program DIANA)] nach einer Idee von Macnaughton-Smith et al. [MWDM64] entwickelt.

Sie beschreiben den Ablauf über die Analogie einer politischen Partei, welche sich durch Meinungsverschiedenheiten aufsplittet. Zuerst tritt das Mitglied mit den größten Differenzen zum Rest der Partei aus derselben aus und gründet eine eigene. Anschließend entscheiden sich nach und nach die restlichen Mitglieder, ob sie ebenfalls “die Fronten wechseln” oder in der ursprünglichen Partei verbleiben wollen [KaRo09, Kap. 6 Divisive Analysis (Program DIANA)]. Aufgrund dieser Analogie wird das Verfahren auch als “Splitterpartei-Methode” bezeichnet [RaRa11].

Konkret lässt sich das Verfahren von Kaufman und Rousseeuw, wie folgt beschreiben:

1. Start mit allen Punkten in einem großen Cluster
2. Berechnung des durchschnittlichen Abstandes jedes Datenpunktes zu den restlichen Punkten im Cluster
3. der Punkt mit der größten mittleren Abweichung verlässt das große Cluster und bildet ein neues
4. alle restlichen Punkte aktualisieren ihren durchschnittlichen Abstand
5. die Punkte, deren durchschnittlicher Abstand zum Cluster größer ist als ihr Abstand zum neu herausgelösten (also der Splittergruppe), verlassen das ursprüngliche Cluster ebenfalls und werden dem neuen Cluster zugeordnet
6. Wiederholung der Schritte ab 4., solange bis keine Punkte mehr das Cluster wechseln
7. aus den entstandenen Clustern wird das größte (mit dem größten Durchmesser) ausgewählt und die Schritte ab 2. wiederholt, solange bis jeder Datenpunkt in einem eigenen Cluster gelandet ist

Der *durchschnittliche Abstand*  $md$  eines Datenpunktes  $x_i$  zu den anderen Punkten  $x_j$  im Cluster  $C$  berechnet sich aus dem Durchschnitt der einzelnen Abstände zu den jeweils anderen Punkten. Diese Abstände werden mit einem geeigneten Distanzmaß berechnet. Damit ist dieses Verfahren für jede Art von Datentypen geeignet, zu denen auch ein Distanzmaß konzipiert werden kann. Mathematisch berechnet sich  $md$  für den Datenpunkt  $x_i \in C$ , wie folgt:

$$md(x_i, C) = \frac{1}{|C| - 1} \sum_{x_j \in C, x_i \neq x_j} d(x_i, x_j) \quad (2.16)$$

Im Ablauf des Verfahrens ist von einem *Clusterdurchschnitt* die Rede. Es dient zur Bestimmung, welches Cluster als nächstes geteilt werden soll. Der Clusterdurchschnitt  $cd$  ist der Abstand der beiden Datenpunkte im Cluster, die am weitesten voneinander entfernt sind:

$$cd(C) = \max_{x_i \in C, x_j \in C} d(x_i, x_j) \quad (2.17)$$

Die genaue Laufzeit des Verfahrens ist nirgendwo beschrieben worden. Am maßgeblichsten müsste der Split des ersten großen Clusters für die Laufzeit sein. Hier werden initial alle Datenpunkte miteinander verrechnet, sodass mindestens eine Laufzeit von  $\mathcal{O}(n^2)$  zustande kommt. Die nachfolgenden Splits sind zunehmend weniger aufwendig in der Berechnung, da die Cluster mit jedem Durchlauf kleiner werden. Dadurch sinkt die Menge an Datenpunkten, die miteinander verrechnet werden müssen, immer weiter ab.

Rajalingam und Ranjini [RaRa11] verglichen die Performance einiger Bottom-up-Verfahren mit einem Top-down-Verfahren, welches dem beschriebenen DIANA entweder identisch oder sehr ähnlich ist. In ihren Versuchen war das Top-down-Verfahren fast doppelt so schnell in der Berechnung der Cluster.

**Bisecting K-Means** Der Bisecting K-Means wurde von Steinbach et al. [StKaKu00] zum Clustern von Dokumenten entwickelt. Ihr ursprünglicher Plan war der Vergleich von agglomerativen Verfahren mit dem klassischen K-Means-Algorithmus. In diesem Prozess entwickelten sie die Idee für das divisive Verfahren – Bisecting K-Means.

Initial befinden sich wieder alle Datenpunkte in einem großen Cluster. Für den Split in zwei kleinere Cluster wird nun der klassische K-Means-Algorithmus mit  $k = 2$  ausgeführt. Aus den resultierenden kleineren Clustern gilt es nun einen Kandidaten für den nächsten Split auszuwählen. Das kann das Cluster mit dem größeren Durchschnitt oder mit der höheren Anzahl an Datenpunkten sein (in ihren Versuchen ergab sich zwischen beiden Vorgehen kein Unterschied). Der Split wird wiederum durch K-Means mit  $k = 2$  ausgeführt usw. Das Verfahren endet, wie gewohnt, wenn alle Punkte in ihrem eigenen Cluster gelandet sind. [StKaKu00]

Das Verfahren besticht durch seine Laufzeit  $\mathcal{O}(n)$ , welche durch den K-Means für den allerersten Split bestimmt wird. Alle folgenden Splits sind durch die kleinere Menge an Datenpunkten je Cluster zu vernachlässigen. [StKaKu00]

## 2.5 Cluster-Validität

### 2.5.1 Überblick

Die Clusteranalyse selbst ist ein Verfahren des unüberwachten Lernens und findet interne Muster in Daten ohne Referenz zu externen Zuweisungen (Labels). Dennoch finden verschiedene Clusteringverfahren unterschiedliche Gruppenzuteilungen. Auch die verschiedenen Parameter, die für ein jeweiliges Clusteringverfahren gesetzt werden, beeinflussen das Ergebnis. Daher bedarf es Methoden zur Evaluation und Vergleich der Clusterings miteinander. [RAAQ11]

**Externe Indizes** (engl. External Indices) sind Metriken, welche die berechnete Gruppenzuteilung mit einer extern vorgegeben Zuteilung vergleicht. Das heißt, es gibt eine erwartete Art der Gruppierung, welche nicht Teil des Datensets selbst ist. Diese Indizes messen nun den

Grad der Übereinstimmung zwischen berechnetem und gewünschtem Clustering-Ergebnis. [RAAQ11]

**Interne Indizes** (engl. Internal Indices) messen die Qualität des Clusterings ohne externe Informationen. Die Cluster sollten möglichst “kompakt” gruppiert und die verschiedenen Gruppen “gut von einander getrennt” sein. Diese Indizes versuchen diese Anforderungen zu quantifizieren. [RAAQ11; und King15, Kap. 8.1. Cluster Validity – Introduction]

Die **Stabilität** ist ebenfalls ein wichtiger Faktor. Hierbei wird das Datenset mehrmals mit verschiedenen Modifikationen geclustert. Solche Modifikationen können die Änderung einzelner Werte oder das Weglassen ganzer Spalten sein. Ein “stabiles” Clustering erzeugt auch mit diesen Veränderungen ähnliche Ergebnisse. Für die konkrete Bewertung der Stabilität werden externe oder interne Indizes verwendet und ihre Ergebnisse über die verschiedenen Clusterings miteinander verglichen. [King15, Kap. 8.1. Cluster Validity – Introduction]

Manche Autoren (z.B. [RAAQ11] und auch [King15]) führen formal noch sog. **Relative Indizes** an. Hiermit wird die Performance verschiedener Clustering-Verfahren und verschiedener Meta-Parameter der Clusterings über das gleiche Datenset verglichen. In der Praxis erfolgt dies aber stets über den Vergleich der externen oder internen Indizes der Clustering-Ergebnisse miteinander.

## 2.5.2 External Indices

### Überblick

Zur Bewertung der Übereinstimmung zweier Clustering-Ergebnisse können grundsätzlich Metriken aus dem Bereich der Klassifikation verwendet werden [HuAr85]. Beispiele dafür wären Maße wie die Entropie oder ihr Gegensatz die Reinheit (engl. purity), welche analysieren, wie viele falsche zu richtigen Zuordnungen innerhalb einer Klasse aufgetreten sind [RAAQ11]. Ebenso das F-Maß (engl. F-Measure), welches aus dem Bereich des Document Retrievals bekannt ist. Es simuliert durchgeführte “Suchen” und vergleicht die gefundenen mit den erwarteten Suchergebnissen [RAAQ11; und StKaKu00].

Das Problem mit diesen Maßen ist, dass die Benennung der Labels im Ergebnis von entscheidender Bedeutung in der Bewertung ist. Nehmen wir an, wir haben ein Datenset  $X$  mit vier Datenpunkten. Die berechnete Clusterzuordnung ist  $Y = (0, 0, 1, 1)$  und die erwartete Zuordnung ist  $Y' = (1, 1, 0, 0)$ . Metriken der Klassifikation würden eine Übereinstimmung von null Prozent feststellen, da keines der Labels in  $Y$  und  $Y'$  den gleichen Klassen zugeordnet worden ist. Im Kontext der Clusteranalyse weisen  $Y$  und  $Y'$  aber eine perfekte Übereinstimmung auf, denn es geht alleine um die Zuordnung der Datenpunkte in die gleichen Gruppen. Wie diese Gruppen benannt werden (in diesem Beispiel 0 und 1), spielt dabei keine Rolle. Daher sind für die Clusteranalyse eigene Metriken entwickelt worden, welche unabhängig von der Benennung der Labels die Übereinstimmung zwischen den Zuordnungen berechnen können. [Rand71; und HuAr85]

## Rand-Index

Der Rand-Index war einer der ersten Metriken, die speziell für die Clusteranalyse entwickelt worden sind [King15, Kap. 8.4 Indices of Cluster Validity]. William M. Rand veröffentlichte dieses Maß 1971 [Rand71] und es ist immer noch eines der populärsten und weitverbreitetsten (bzw. vor allem die Verbesserungen dieses Indexes siehe nächster Abschnitt). [HuAr85; und Stei04]

Um unabhängig von der Benennung der Labels zu werden, vergleicht man nicht wie in der Klassifikation die gefundenen Labels eins zu eins miteinander. Stattdessen werden alle möglichen Paarungen der Datenpunkte in  $X$  betrachtet und die Zuordnung der Paare in  $Y$  und  $Y'$  miteinander verglichen. [Rand71]

Tabelle 2.3: Kontingenz der Datenpunkt-Paare in  $Y$  und  $Y'$  [Stei04]

$Y \setminus Y'$	Paar im gleichen Cluster	Paar in anderem Cluster
Paar im gleichen Cluster	$a$	$b$
Paar in anderem Cluster	$c$	$d$

Die Tabelle zeigt die möglichen Fälle: Ein Paar aus Datenpunkten kann entweder dem gleichen Cluster oder zwei unterschiedlichen Clustern zugeordnet worden sein. Haben sowohl  $Y$  als auch  $Y'$  das Paar jeweils in das gleiche Cluster (mit dem gleichen Label wie auch immer dieses benannt ist) zugeordnet, so wird dieses Paar zu  $a$  gezählt. Haben beide jeweils ein unterschiedliches Cluster zugeordnet, so wird das Paar in  $d$  gezählt usw. [Stei04]

$$rand = \frac{a + d}{a + b + c + d} = \frac{a + d}{\binom{n}{2}} \quad (2.18)$$

Der Rand-Index teilt nun die Menge an Paaren, welche die gleiche Zuordnung erhalten haben ( $a$  und  $d$ ) durch die Anzahl aller möglichen Paarungen. Das Ergebnis ist ein Wert zwischen 0 (keine Übereinstimmung) und 1 (perfekte Übereinstimmung). [Rand71]

Der Rand-Index weist eine Reihe von Problemen auf, welche teilweise von Rand selbst erkannt oder später ermittelt worden sind:

- Je höher die Anzahl an Clustern ist, desto höher (näher an der 1) liegt der Rand-Index standardmäßig. Das kommt daher, dass bei einer hohen Anzahl an Clustern die meisten Paare unterschiedlichen Clustern zugewiesen sein müssen. Das treibt den Index künstlich in die Höhe. [Stei04; teilweise zitiert nach Rand71]
- Zum Vergleich verschiedener Clustering-Verfahren werden häufig Monte-Carlo-Simulationen mit großen Mengen an zufällig generierten Datensätzen und Clusterings verwendet. Werden zwei Zufalls-Zuordnungen miteinander verglichen, so gibt der Rand-Index keine Werte nahe der 0, was aber wünschenswert wäre. Das liegt daran, dass die Anzahl an zufälligen Übereinstimmungen nicht adäquat herausgerechnet wird. [HuAr85; und King15, Kap. 8.4 Indices of Cluster Validity]

## Adjusted Rand-Index

Aus den genannten Problemen haben eine Vielzahl von Autoren versucht, eine bessere Variante des Rand-Indexes zu finden. Vor allem das “Herausrechnen” von angeblichen hohen Übereinstimmungen bei zufälligen Zuteilungen (engl. correction for chance) ist eines der Hauptanliegen gewesen [HuAr85]. Von allen vorgestellten Lösungen scheint die Variante von Hubert und Arabie [HuAr85] diejenige mit den wünschenswertesten Eigenschaften zu sein. (siehe die Versuche z.B. von Steinley [Ste04])

$$arand = \frac{rand - rand_{expected}}{rand_{max} - rand_{expected}} \quad (2.19)$$

Im Allgemeinen geht es darum, den Rand-Index zu berechnen und anschließend den “erwarteten Rand-Index” für zufällige Zuordnungen davon abzuziehen. Verschiedene Autoren haben nun unterschiedliche Methoden hergeleitet, um diesen “erwarteten Rand-Index” zu berechnen. [HuAr85]

$$arand = \frac{\binom{n}{2}(a+d) - [(a+b)(a+c) + (c+d)(b+d)]}{\binom{n}{2}^2 - [(a+b)(a+c) + (c+d)(b+d)]} \quad (2.20)$$

Hubert und Arabie berechnen den “erwarteten Rand-Index” aus der Menge aller möglichen Permutationen der Paarungen und ihrer Übereinstimmung  $\frac{(a+b)(a+c) + (c+d)(b+d)}{\binom{n}{2}}$ . Die gegebene Formel zeigt eine breits vereinfachte Umstellung von Steinley. [HuAr85; und vereinfachte Formel aus Ste04]

Dieser “Adjusted Rand-Index” liefert Werte zwischen  $-1$  und  $1$ . Negative Ergebnisse oder Werte um die  $0$  zeigen, dass die Übereinstimmung der beiden Clusterings eher auf den Zufall zurückzuführen und damit nicht signifikant ist. Höhere Zahlen nahe der  $1$  stehen nach wie vor für eine perfekte Übereinstimmung, welche zusätzlich statistische Signifikanz aufweist. Auch der Bias zu hohen Werten bei einer hohen Anzahl an Clustern wird dadurch ausgeglichen. [HuAr85; und Ste04]

## 2.5.3 Internal Indices

In der Literatur werden eine Vielzahl von Indizes beschrieben, welche die Qualität der gefundenen Cluster messen können. Einen Überblick dazu geben Rendón et al. [RAAQ11]. An dieser Stellen sollen nur ein paar Vertreter vorgestellt werden, welche für diese Arbeit von besonderer Relevanz sind.

### Silhouettenkoeffizient

Der Silhouettenkoeffizient wurde von Peter J. Rousseeuw entwickelt [Rous87] und ist ein häufig verwendetes Maß für die Qualität von Clusterings.

$$s(x) = \frac{b(x) - a(x)}{\max a(x), b(x)} \quad (2.21)$$

$$a(x) = d(x, C_x) \quad (2.22)$$

$$b(x) = \min_{C_x \neq C_i} d(x, C_i) \quad (2.23)$$

Für jeden Datenpunkt  $x$  im Datenset wird die Silhouetten-Weite  $s(x)$  berechnet. Diese ergibt sich aus der Differenz zwischen dem durchschnittlichen Abstand zu allen Datenpunkten im selben Cluster ( $a(x)$ ) und dem durchschnittlichen Abstand zu allen Datenpunkten im direkt benachbarten Cluster ( $b(x)$ ). Anschließend wird der Wert zwischen  $-1$  und  $1$  normiert. [Rous87]

- Eine negative Silhouetten-Weite gibt an, dass der Datenpunkt näher am benachbarten als an seinem eigenen Cluster liegt und somit falsch zugeordnet worden ist.
- Werte um die 0 zeigen, dass der Datenpunkt fast mittig zwischen beiden Clustern platziert ist.
- Befindet sich die Weite nahe der 1, so liegt der Datenpunkt mittig in seinem eigenen Cluster und das benachbarte Cluster ist ordentlich weit entfernt.

$$sil = \frac{1}{n} \sum_{i=1}^n s(x_i) \quad (2.24)$$

Der Koeffizient ergibt sich schließlich aus dem Durchschnitt aller Silhouetten-Weiten aller Datenpunkte. Werte nahe der 1 deuten auf kompakte und wohl-separierte Cluster hin. [Rous87]

In den Versuchen von Rendón et al [RAAQ11] erwies sich dieser Index als einer der besten Indikatoren für ein gutes Clustering-Ergebnis. Er ist die direkte mathematische Definition für die Anforderung, dass Datenpunkte im gleichen Cluster möglichst ähnlich und zu den Punkten der anderen Cluster möglichst unähnlich sein sollen. Ein Nachteil dieses Indizes ist die quadratische Laufzeit, da jeder Punkt mit jedem anderen verglichen werden muss.

### Davies-Bouldin Index

Der Davies-Bouldin Index wurde von seinen Namesgebern David L. Davies und Donald W. Bouldin [DaBo79] entwickelt. Ziel war es, eine Metrik zu konstruieren, welche die durchschnittliche Ähnlichkeit zwischen benachbarten Clustern berechnet.

$$dbi = \frac{1}{k} \sum_{i=1}^k \max_{i \neq j} \frac{d(c_i, C_i) + d(c_j, C_j)}{d(c_i, c_j)} \quad (2.25)$$

Die Ähnlichkeit zwischen den Clustern  $K$  ( $K$  steht hier für die Menge an gefundenen Clustern  $C_i$ ) berechnet sich aus dem durchschnittlichen Abstand der Punkte eines Clusters zu ihrem Cluster-Mittelpunkt ( $d(c_i, C_i)$  und  $d(c_j, C_j)$ ) geteilt durch den Abstand der beiden Cluster-Mittelpunkte zueinander ( $d(c_i, c_j)$ ). [DaBo79]

Je kleiner der Wert des Indexes, desto enger liegen die Datenpunkte um ihren Cluster-Mittelpunkt im Verhältnis zum benachbarten Cluster. Möglichst niedrige Werte nahe der 0 sind also als optimal anzusehen. [DaBo79]

Der große Vorteil von diesem Verfahren ist die geringere Laufzeit, da die Punkte der Cluster nur mit ihrem Mittelpunkt verrechnet werden. Nachteilig ist, dass die Be- und Verrechnung der Mittelpunkte primär nur für numerische Vektoren definiert ist. In den Versuchen von Rendón et al. [RAAQ11] schnitt dieser Index genauso gut ab wie der Silhouettenkoeffizient.



## 3 Konzeption

### 3.1 Überblick

Zur Beantwortung der Kernfrage wird ein Konzept für das Clustern komplexerer Datenstrukturen erarbeitet und anschließend evaluiert. Diese Herleitung soll aber nicht nur theoretischer Natur sein. Daher wird die Überprüfung des Konzeptes an einem möglichst realen Beispiel mit realitätsnahen Daten erfolgen. Die folgende Grafik zeigt einen Überblick zum Praxisteil dieser Arbeit:

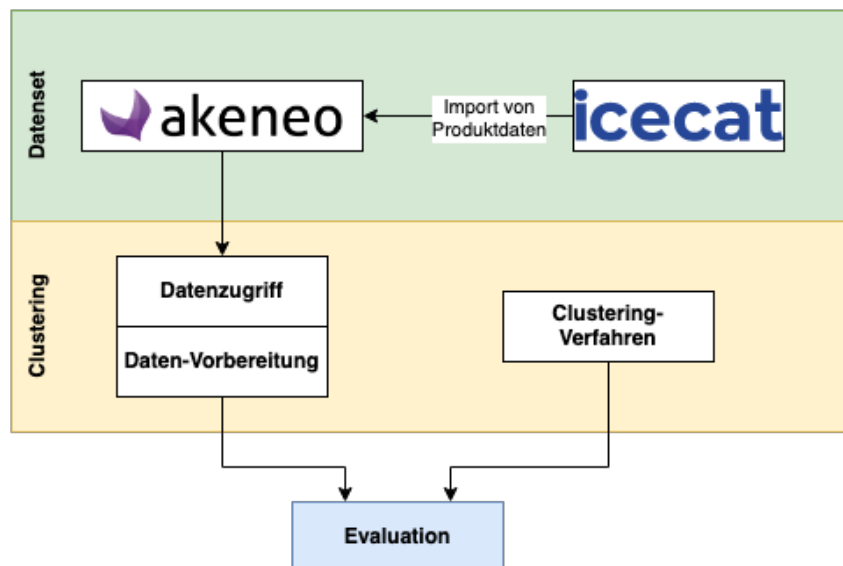


Abbildung 3.1: Grober Architektur-Entwurf für den Praxisteil

Um ein authentisches *Dataset* in den typischen Datenstrukturen zu erhalten, wird eine Instanz des Product-Information-Management-Systems Akeneo-PIM aufgesetzt und mit Produktdaten aus dem umfangreichen Online-Katalog Icecat gefüllt. Im Abschnitt Datenquellen & -sets werden beide Komponenten genauer erläutert. Außerdem werden genaue Anforderungen für die Auswahl der Produkte, welche importiert werden, definiert.

Parallel wird ein Konzept für ein *Clustering-Verfahren* erarbeitet, welches die vielfältigen Datenstrukturen verarbeiten kann. Dies setzt i.d.R. eine entsprechende Vorverarbeitung der Daten voraus. Das entsprechende Konzept wird im Abschnitt Clustering-Verfahren hergeleitet.

In der *Evaluation* kommt schließlich das erstellte Dataset mit dem erarbeiteten Clustering-Verfahren zusammen. Es wird überprüft, inwiefern das Verfahren "sinnvolle" Cluster in den Produktdaten findet. Eine genaue Definition von "sinnvollem" Clustering sowieso die konkrete Berechnung entsprechender Metriken usw. wird im Abschnitt Evaluation dargelegt.

## 3.2 Datenquellen & -sets

### 3.2.1 Akeneo-PIM

Als realistische Datenquelle ist das Product-Information-Management-System Akeneo-PIM <https://www.akeneo.com/de/akeneo-pim-community-edition/> ausgewählt worden.

#### Überblick

Der Ursprung von Product-Information-Management-Systemen liegt in gedruckten Produktkatalogen. Speziell im letzten Jahrhundert waren diese die zentrale Form, wie Produktdaten verteilt und verwaltet worden sind. Mit beginnender Digitalisierung und dem Aufkommen des E-Commerce in den 1990er-Jahren wurden zunehmend computer-basierte Systeme zur Warenhaltung und Bestellabwicklung, sog. ERP-Systeme, verwendet. Anfangs wurden in diesen Systemen auch die Produktdaten gespeichert. Dieses Vorgehen stieß aber zunehmend an Grenzen, da die schiere Vielfalt an Daten (Produkte, Bestellungen, Transaktionen, Warenbestand usw.) schwer zu verwalten war. Mitte der 2000er kamen spezielle Systeme zum Einsatz, die lediglich die Produktdaten enthielten. Somit konnte der “Produktkatalog” zentral verwaltet, angereichert, übersetzt und aufbereitet werden. Weitere Systeme (z.B. besagte ERP-Systeme) konnten nun auf diesen zentralen Katalog zugreifen und müssen selbst nur minimale Informationen zu den Produkten vorhalten. PIM-Systeme sitzen daher heute häufig im Zentrum der Unternehmens-Architektur und dienen als zentrale Datenquelle für die angeknüpften Systeme. [Pimc22]

“Akeneo-PIM” der gleichnamigen Firma ist ein solches PIM-System, welches heute breite Anwendung findet. Es ist in einer betreuten Enterprise Edition erhältlich sowie als Open-Source-Variante “Akeneo-PIM Community Edition” (verfügbar unter <https://github.com/akeneo/pim-community-dev>), welche selbst gehostet werden muss. Die Community-Variante ist frei erhältlich und daher sehr populär. [Aken22a]

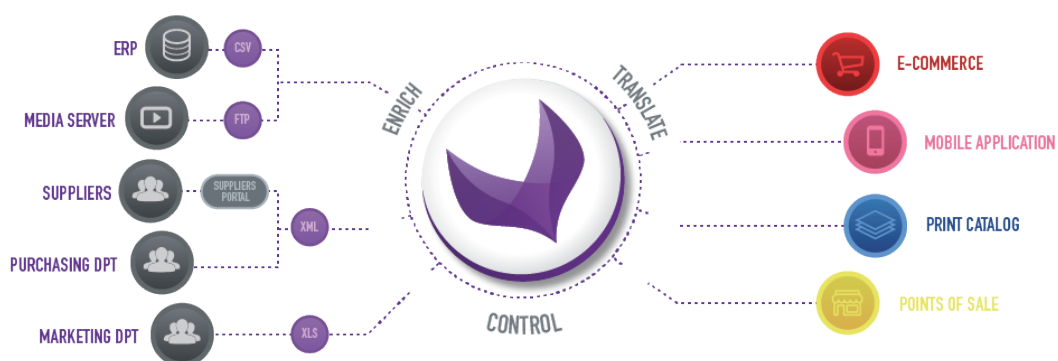


Abbildung 3.2: Schematische Darstellung der Funktionen von Akeneo-PIM [Aken22b]

Ein weiterer Grund für die weite Verbreitung ist, dass Akeneo-PIM über diverse Schnittstellen an externe Systeme angebunden werden kann. Ebenso unterstützt es eine Vielzahl an Import- und Export-Möglichkeiten in unterschiedliche Datenformaten. Es gibt außerdem

den “Akeneo App Store” (<https://marketplace.akeneo.com/>), wo viele Anbieter Plugins veröffentlichen, um die Funktionen und Anbindungen von Akeneo-PIM zu erweitern. [Aken22b]

Diese Gründe führten zur Verwendung von Akeneo-PIM für diese Arbeit.

## Systemübersicht

Akeneo-PIM ist eine PHP-Applikation, welche auf dem Symfony-Framework basiert. Für die Datenspeicherung wird die Verwendung von einer MySQL-Datenbank empfohlen. Zum schnelleren Suchen und Abrufen kommt Elasticsearch zum Einsatz. [Aken22c] Hinzu kommen weitere Technologien zur Generierung des Webfrontends. Um diesen verschiedenen Tools Herr zu werden, stellt Akeneo Konfigurationen für Containerization-Tools wie Docker-Compose zur Verfügung. [Aken22d]

Das System kann über ein Webfrontend genutzt werden. Hier lassen sich verschiedene Nutzer und Applikationen hinterlegen, welche verschiedene Berechtigungen zum Abrufen und Bearbeiten der Daten bekommen können.

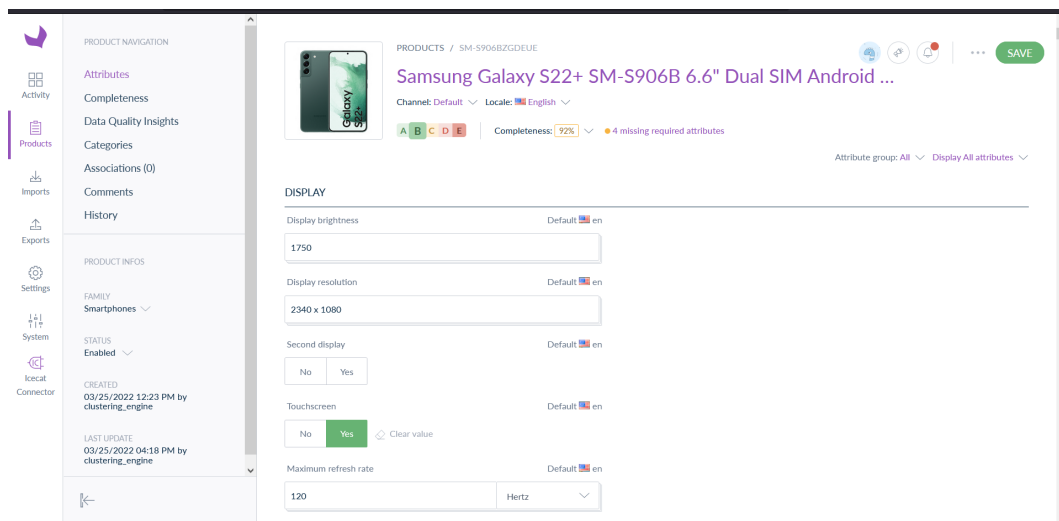


Abbildung 3.3: Akeneo-PIM Webfrontend bei der Bearbeitung eines Produktes

Eine Restful-API steht ebenfalls zur Verfügung, um Daten mittels HTTP-Requests abzurufen und zu verändern. In der Regel werden die Daten im JSON-Format transferiert und zurückgegeben. [Aken22e]

## Datenstrukturen

Innerhalb von Akeneo-PIM kommen eine Vielzahl von verschiedenen Datenstrukturen zum Einsatz. Im Folgenden werden nur die für diese Arbeit wichtigsten Strukturen kurz erklärt.

**Product & Product Values** Die Produkte sind das zentrale Element in Akeneo-PIM. Aus ihnen leiten sich eine Reihe weiterer Strukturen ab. [Aken22f]

Jedes Produkt wird genau einer “Family” zugeordnet. Sie bestimmt, welche Attribute dieses Produkt aufweisen sollte.

Ein Produkt kann einer oder mehreren “Categories” zugeordnet sein. Der Sinn dabei ist der Aufbau z.B. von Navigationsstrukturen. Hierbei ist die Zuordnung in mehrere verschiedene Kategorien möglich.

Die “ProductValues” sind die tatsächlichen Wertausprägungen eines Produktes. Sie werden in Form einer Hash-Map angegeben. Der Key der Einträge ist stets der “code” eines bestimmten “Attributes” (eine weitere Datenstruktur). Dieses Attribut wird separat definiert und legt verschiedene Constraints für den “ProductValue” fest. Dabei bietet Akeneo volle Flexibilität. Jedes Produkt kann jedes beliebige Attribut mit einem entsprechenden Wert aufweisen.

```

1 {
2   "identifier": "1111111195",
3   "family": "clothing",
4   "categories": ["tshirts"],
5   "values": {
6     "ean": [
7       {"locale": null, "scope": null, "data": "1234567890207"}
8     ],
9     "size": [
10      {"locale": null, "scope": null, "data": "s"}
11    ],
12    // ...
13  },
14  "enabled": true,
15  "created": "2017-10-05T11:25:48+02:00",
16  "updated": "2017-10-05T11:25:48+02:00",
17  // ...
18 }

```

**Category** Mit einer “Category” können Produkte in verschiedene Kategorien sortiert werden. Wichtig ist, dass “Categories” beliebig verschachtelt und hierarchisch sortiert werden können. Es gibt kein Limit für die Tiefe der Verschachtelung. [Aken22g]

```

1 {
2   "code": "tvs_projectors",
3   "parent": "master",
4   "labels": {
5     "en_US": "TVs and projectors",
6     // ...
7   }
8 }

```

*Hinweis:* Mit Ausnahme der “Products” haben die meisten Entitäten einen “code”, welcher als eindeutiger Identifier fungiert.

**Family** Eine “Family” definiert die Attribute, welches ein zugehöriges “Product” aufweisen sollte. Dies ist allerdings keine strikte Zuteilung. Produkte können auch Attribute aufweisen, welche gar nicht in ihrer “Family” vorgesehen sind. Ebenso können in der Family erforderliche

(engl. required) Attribute hinterlegt werden. Produkte müssen diese aber nicht erfüllen. Das Produkt weist dann lediglich einen niedrigeren “Completeness Score” auf. Dadurch kann in Akeneo auch mit unvollständigen Produkten gearbeitet werden, bevor man diese weiter anreichert. [Aken22g]

```

1 {
2   "code": "camcorders",
3   "attributes": ["description", "name", "optical_zoom"],
4   "attribute_requirements": {
5     "ecommerce": ["description", "name", "optical_zoom"],
6     "mobile":    ["name"],
7     // ...
8   },
9   "labels": {
10    "en_US": "Camcorders",
11    // ...
12  }
13 }
```

Die Erforderlichkeit von Attributen kann für jeden Channel einzeln definiert werden. Für eine Erklärung zu Channels siehe Abschnitt Channel, Currency, Locale.

**Attribute & Attribute Group** Die verschiedenen Wertausprägungen eines Produktes müssen einem entsprechenden “Attribute” zugeordnet sein. Ein solches Attribut muss zuerst definiert werden, bevor Produkte einen Wert darin aufweisen können. [Aken22g]

```

1 {
2   "code": "auto_exposure",
3   "type": "pim_catalog_boolean",
4   "group": "technical",
5   "localizable": false, // erlaubt andere Werte in anderen Sprachen
6   "scopable": false,    // erlaubt andere Werte auf anderen Channels
7   "unique": false,      // verhindert die mehrfache Verwendung eines Wertes
8   // ... diverse weitere Constraints
9 }
```

Jedes Attribut kann einer “AttributeGroup” zugeordnet werden. Dies erlaubt das bessere Strukturieren der Attribute. [Aken22g]

Jedes Attribut ist einem `type` zugeordnet, welches den Datentyp anzeigt. Z.B. der Typ `pim_catalog_boolean` zeigt einen booleschen Wert an, `pim_catalog_text` eine Textzeile usw. Zusätzlich können in einem Attribut auch weitere Constraints für die jeweiligen Typen definiert werden. Bspw. kann mit dem Constraint `decimals_allowed` festgelegt werden, ob ein numerisches Attribut Nachkommastellen erlaubt oder nicht [Aken22g]. Auf diese Constraints wird allerdings nicht weiter eingegangen, da sie für diese Arbeit keine Rolle spielen. Die folgende Tabelle zeigt eine Auflistung aller möglicher Typen.

Tabelle 3.1: Attribut-Typen in Akeneo-PIM [Aken22h]

Akeneo-Typ	Erklärung
<code>pim_catalog_identifier</code>	eindeutige Id
<code>pim_catalog_text</code>	Textzeile
<code>pim_catalog_textarea</code>	mehrzeiliger Textblock
<code>pim_catalog_simpleselect</code>	einfache Auswahl aus einer Liste an Optionen

Akeneo-Typ	Erklärung
<code>pim_catalog_multiselect</code>	mehrfache Auswahl aus einer Liste an Optionen
<code>pim_catalog_boolean</code>	boolescher Wert
<code>pim_catalog_date</code>	Datum
<code>pim_catalog_number</code>	einfache rationale Zahl
<code>pim_catalog_metric</code>	rationale Zahl mit Recheneinheit
<code>pim_catalog_price_collection</code>	Preise je Währung
<code>pim_catalog_image</code>	Bild
<code>pim_catalog_file</code>	sonstige Datei (z.B. Datenblatt als PDF)
<code>pim_reference_data_simpleselect</code>	einfache Auswahl einer "ReferenceEntity"
<code>pim_reference_data_multiselect</code>	mehrfache Auswahl von "ReferenceEntities"

*Hinweis:* "ReferenceEntities" sind in dieser Arbeit ebenfalls nicht verwendet worden und können daher ignoriert werden.

**Channel, Currency, Locale** Um maximale Flexibilität bei der Verwaltung der Daten zu gewährleisten, können verschiedene "Channels" oder "Scopes" definiert werden. Für jeden "Channel" können die erforderlichen Währungen ("Currencies") und Lokalisierungen ("Locales") hinterlegt werden. [Aken22i]

```

1  {
2    "code": "ecommerce",
3    "currencies": ["USD", "EUR"],
4    "locales": ["de_DE", "en_US", "fr_FR"],
5    "category_tree": "master",
6    "labels": {
7      "en_US": "E-Commerce",
8      // ...
9    }
10 }
```

"Attributes" können als `localizable` definiert werden. Also treten verschiedene Werte je nach gewählter Sprache und Region auf (z.B. der Produktname). Außerdem können sie als `scopable` festgelegt werden. In dem Fall können zu jedem "Channel" zusätzlich nochmal unterschiedliche Werte hinterlegt werden. Bspw. gibt es die Produktbeschreibung in den verschiedenen Sprachen und dazu für jedes Produkt eine kürzere Beschreibung im Channel "mobile" und eine längere im Channel "print". [Aken22g]

**Measurement Family** Für den Attribut-Typ "`pim_catalog_metric`" werden eine Zahl und eine Einheit hinterlegt. Jede Einheit wird einer bestimmten "Measurement Family" zugeordnet. Dort sind alle verfügbaren Einheiten sowie deren Umrechnung ineinander hinterlegt. [Aken22g]

### 3.2.2 Icecat

Akeneo-PIM ist nur eine Software, in der Produkte strukturiert gespeichert und aufbereitet werden können. Tatsächliche Produktdaten müssen aber erst importiert werden. Eine Möglichkeit wäre die Kreation eines komplett künstlichen Sets an Produkten. Das Ziel

dieser Arbeit ist es aber, eine möglichst realitätsnahe Evaluation durchzuführen. Daher wird für den Import von realistischen Produktdaten auf Icecat (<https://icecat.biz/de>) zurückgegriffen.

Icecat ist ein sehr umfangreicher offener Online-Produktkatalog. Über 33.000 Hersteller veröffentlichen hier die Daten ihrer Produkte. Zum heutigen Stand (April 2022) sind über 11 Mio. Produkte bei Icecat gelistet. Viele dieser Daten sind nur für Premium-Nutzer zugänglich. Allerdings gibt es eine ganze Reihe von Unternehmen (u.a. große Player wie Samsung), die als “Icecat-Sposoren” ihre Produktdaten frei zugänglich zur Verfügung stellen. [Icec22a]

Icecat bietet Herstellern damit die Möglichkeit, sehr schnell ihre neuen Produkte zu verbreiten, sodass Shops und Händler direkt die Daten automatisiert in ihr Sortiment aufnehmen können. Zusätzlich werden die Produktdaten nach der UNSPSC-Taxonomy (eine international geläufige Taxonomie zur Definition von Produkteigenschaften) von Icecat standardisiert. [Icec22b]

Aus diesem umfangreichen Fundus wird ein Datenset zur Beantwortung der Kernfrage erstellt. Icecat ist zudem besonders attraktiv, da es einen Importer für Akeneo gibt (<https://marketplace.akeneo.com/extension/akeneo-icecat-connector>), welcher den Import von Produktdaten größtenteils automatisieren soll.

### 3.2.3 Datenset

#### Anforderungen

Für die Auswahl von Produkten werden zunächst einige Anforderungen definiert:

Die Menge an verschiedenen Attributen könnte sich sehr unterschiedlich auf das Clustering auswirken. Daher sollte es sowohl Produkte mit sehr **wenigen** und mit sehr **vielen** Attributen geben.

Daraus folgt, dass Produkte aus **verschiedenen Kategorien** gewählt werden sollten. Als zusätzliche Herausforderung könnten diese Kategorien am besten **miteinander verwandt** sein, sodass die Fähigkeit einer eindeutigen Trennung der Kategorien überprüft werden kann.

Innerhalb einer Produktkategorien sollten **verschiedene Variationen** des gleichen Produktes sowie **verwandte Varianten** (etwa mehrere verschiedene Generationen des Produktes) auftreten.

Zuletzt sollten **Duplikate** enthalten sein. In realen Datensets kommt häufig dasselbe Produkt mehrmals vor, etwa durch geringfügige Unterschiede der Daten beim Import aus verschiedenen Quellen. Ein Clustering-Verfahren sollte in der Lage sein, solche Duplikate zu erkennen.

## Umsetzung

Aus den Anforderungen ist folgendes Vorgehen abgeleitet worden:

Es werden Smartphones der Firma Samsung aus der Galaxy S-Reihe importiert. Hier werden Exemplare der letzten drei Generationen ausgewählt. In jeder Generation gibt es unterschiedliche Varianten (die Standard-Ausführung, “Plus”, “Ultra” und “Fan-Edition”), welche verschiedene Bildschirmgrößen und etwas unterschiedliche Komponenten verbaut haben. [DE22]

Smartphones weisen sehr viele verschiedene Attribute auf und eine einfache Suche auf Icecat zeigte bereits mehrere Duplikate der Samsung-Geräte auf der Plattform.

Als nächstes werden noch Produkte aus der Kategorie “Smartphone-Hüllen” ergänzt. Es werden Hüllen von verschiedenen Herstellern für die besagten Samsung Galaxy Smartphones ausgewählt. Solche Hüllen haben eine viel geringere Menge an Attributen, sind aber trotzdem eng mit den Smartphones im Set verwandt.

Es ist damit sogar denkbar, einige Clustering mit dem Ziel durchzuführen, die Smartphones mit ihren zugehörigen Hüllen in die gleiche Gruppe einteilen zu können.

## 3.3 Clustering

### 3.3.1 Clustering-Verfahren

#### Anforderungen

Um aus den vielfältigen Ansätzen des Clusterings ein geeignetes Verfahren auswählen zu können, werden zunächst einige genauere Anforderungen für die Fragestellung dieser Arbeit definiert:

Ein **hierarchisches Verfahren** ist für diese Arbeit wünschenswert. Es ist bezeichnend, dass Akeneo beliebige Verschachtelungen der “Categories” erlaubt. In der Praxis ist die Einteilung von Produkten in Kategorien stark vom gewünschten Detail-Grad abhängig. Ein hierarchisches Verfahren bildet dieses Verhalten exakter ab. Außerdem müssen partitionierende Verfahren die gesuchte Anzahl an Clustern vorab übergeben bekommen. Damit sind zuerst andere Analysen notwendig, bevor der richtige Wert ermittelt ist. Eine hierarchische Einteilung gibt dem Anwender die Wahl, wie weit aufgefächert die Ergebnisse sein sollen.

Ebenso ist die **Laufzeit** ein wichtiger Faktor, da in der Praxis tendenziell sehr große Mengen an Produkten auftreten können.

Von Vorteil ist weiterhin, wenn das Verfahren mit relativ vielen **verschiedenen Datentypen** umgehen kann. Die meisten vorherigen Transformationen verringern tendenziell den Informationsgehalt der Daten (bspw. die Diskretisierung numerischer Attribute in eine geringere Anzahl an Bändern). [KaRo09, Kap. 1.2.6 Mixed Variables]



## Ansatz

Agglomerative Verfahren weisen immer mindestens eine quadratische Laufzeit auf [Sibs73]. Diverse Ansätze sind dagegen meistens effizienter (siehe Versuche von [RaRa11]) und liefern tendenziell für uns Menschen besser nachvollziehbare Einteilungen. [King15, Kap. 3.3 Agglomerative versus Divisive Clustering]

Das Verfahren Bisecting K-Means wird in dem Zusammenhang als sehr potent angesehen [StKaKu00]. Allerdings muss dieser Ansatz modifiziert werden, um auch mit gemischten Attributen arbeiten zu können. Da zum K-Means eine Variante für gemischte Attribute verfügbar ist (K-Prototypes [Huan98]), wird dieser Ansatz mit dem “Bisecting”-Prinzip kombiniert. Das eingesetzte Verfahren ist also ein **Bisecting K-Prototypes**. Aus der Literatur ist diese Kombination noch nicht bekannt und damit neuartig. Mit diesem Ansatz können sehr viele Arten von Attributen verarbeitet werden. Mehr dazu im folgenden Abschnitt.

### 3.3.2 Distanzfunktion

#### Problemstellung

Grundsätzlich kann K-Prototypes mit gemischten Attributen (numerisch und kategorisch) arbeiten. Die Produktdaten in Akeneo-PIM können allerdings noch weitere Ausprägungen annehmen. Die folgende Tabelle zeigt die grobe Einteilung der Attribut-Typen in Typ-Klassen:

Tabelle 3.2: Einteilung der Akeneo-Typen in übergeordnete Datenklassen

Klasse	Akeneo-Typ
numerisch	pim_catalog_date, pim_catalog_number, pim_catalog_metric, pim_catalog_price_collection
kategorisch	pim_catalog_boolean, pim_catalog_simpleselect, pim_reference_data_simpleselect
multi-kategorisch	pim_catalog_multiselect, pim_reference_data_multiselect
string	pim_catalog_text, pim_catalog_textarea
Datei	pim_catalog_image, pim_catalog_file
sonstige	pim_catalog_identifier

Für die numerischen und kategorischen Attribute muss überprüft werden, ob eine Art “Vorverarbeitung” vor dem Clustering sinnvoll ist. Ansonsten können diese direkt verwendet werden.

Die Klasse “multi-kategorisch” beschreibt die Möglichkeit aus einer gegebenen Liste an Optionen mehrere auswählen zu können. Bspw. könnte eine Liste von Materialien gegeben sein und das jeweilige Produkt markiert alle vorkommenden (z.B. Silikon, PET & Glas). Eine solche Datenklasse ist in der Literatur nicht beschrieben. Der Umgang mit diesen Attributen muss also gesondert erarbeitet werden.

Die “string”-Klasse ist streng genommen klassisch “kategorisch”. Allerdings handelt es sich hier um freie Textfelder (z.B. der Produkttitel), welche bei jedem Produkt völlig frei gefüllt

werden kann. Somit ist die Betrachtung als “Kategorien” nicht zielführend. Wenn bspw. jedes Produkt einen individuellen Titel hat, dann gibt es genau so viele “Kategorien” wie Produkte. Dies ist keine sinnvolle Form der Datenverarbeitung. Auch für diese Klasse muss eine gesonderte Art des Umgangs gefunden werden.

Die Klasse “Datei” wird im Rahmen dieser Arbeit ignoriert. Die Analyse von Bildern oder Textdokumenten ist ein weites Feld mit vielen verschiedenen Ansätzen. Die tatsächliche Relevanz bspw. eines Produktbildes für das Clustering ist aber fraglich. In der Theorie müssten alle Eigenschaften eines Produktes in den anderen Attributen ebenfalls abgebildet sein. Im Bild liegen diese Information aber äußerst unstrukturiert vor, wenn überhaupt.

Der Akeneo-Typ “pim\_catalog\_idenfier” ist die Id eines Produktes. Sie hat keine Relevanz für das Clustering. Ebenso können in Akeneo Attribute als **unique** gekennzeichnet werden. Dies ist sinnvoll für das Speichern mehrerer verschiedener Identifier (also in Akeneo wird das besagte Identifier-Attribut genutzt, aber die Id im ERP-System ist eine andere, die hier ebenfalls hinterlegt ist). Diese Identifier geben aber keine Ähnlichkeitsinformationen, sondern sind rein logistische Werte. Sie werden im Clustering stets ignoriert werden.

Ein weiteres Problem ist, dass die Definition der Attribut-Anforderungen in den Akeneo-Families keinen bindenden Charakter haben. Das heißt, jedes Produkt kann Werte für jedes beliebige Attribut aufweisen oder auch nicht. Der Umgang mit **null**-Values wird also ebenfalls eine zentrale Rolle einnehmen.

## Ansatz

**numerische Attribute** Viele der numerischen Attribute enthalten ihre Daten auf eine implizite Art und Weise (z.B. ein Datum). Die folgende Tabelle zeigt für die verschiedenen Typen in Akeneo, ob und wie diese Daten extrahiert werden.

Tabelle 3.3: Vorverarbeitung der numerischen Attribute in Akeneo

Akeneo-Typ	Vorverarbeitung
pim_catalog_date	Umwandlung in Unix-Timestamp, Normalisierung
pim_catalog_number	Normalisierung
pim_catalog_metric	Umrechnung in Standard-Unit des Attributs, Normalisierung
pim_catalog_price_collection	Filter nach einer Währung z.B. USD, Normalisierung

Die Verarbeitungsschritte sind recht selbsterklärend. Alle Attribute werden auf das Intervall zwischen  $[0, 1]$  normalisiert, um eine stärkere Gewichtung von Attributen mit tendenziell höheren Zahlen (z.B. Unix-Timestamps) zu vermeiden. Für die Normalisierung werden nur die tatsächlich vorkommenden Werte genutzt. Werte, welche in den Constraints der Akeneo-Attribute definiert sind (z.B. **date\_min** und **date\_max**), werden nicht betrachtet.

$$x^i = \frac{x^i - \min X^i}{\max X^i - \min X^i} \quad (3.1)$$

**kategorische Attribute** Es könnte bei ordinalen Attributen sinnvoll sein, diese u.U. als numerische Attribute anzusehen. Allerdings geben die Daten in Akeneo keinen direkten Rückschluss her, ob es sich z.B. bei einem Single-Select eigentlich um ein ordinales Attribut handelt. Somit müssten hier alle Attribute (knapp hundert) händisch analysiert werden, was den Rahmen dieser Arbeit weit gesprengt hätte.

Der K-Prototypes verlangt keine spezielle Vorverarbeitung kategorischer Attribute abseits von der Umwandlung der Labels in symmetrische bzw. asymmetrische binäre Werte. Der Akeneo-Typ “Bool” ist theoretisch symmetrisch. Praktisch kann aber in Akeneo jedes Attribut auch `null` sein. Somit können tatsächlich drei Wertausprägungen (*true*, *false* oder `null`) vorkommen. Daher werden diese Attribute genauso wie alle anderen (z.B. Single-Selects) in asymmetrische binäre Attribute umgewandelt.

Zu beachten ist aber, dass kategorische Attribute auf ihre Ähnlichkeit und numerische auf ihre Distanz überprüft werden. Da die numerischen im Intervall  $[0; 1]$  liegen, können die Ergebnisse der Ähnlichkeitsmaße (z.B. Jaccard-Koeffizient) einfach invertiert werden.

$$d(x_1^{cat}, x_2^{cat}) = 1 - \frac{|x_1^{cat} \cap x_2^{cat}|}{|x_1^{cat} \cup x_2^{cat}|} \quad (3.2)$$

**multi-kategorische Attribute** Der naheliegendste Ansatz besteht darin, die verschiedenen auftretenden Kombinationen an gewählten Optionen in eigene Kategorien zu fassen und sie wie normale kategorische Attribute zu behandeln. Dies geht allerdings mit einem Verlust an Informationen einher: Angenommen ein Produkt besteht aus den Materialien {silicone, pet} und ein anderes aus {silicone, glass}. Beide würden nun unterschiedlichen Kategorien zugeordnet und ihre Ähnlichkeit ist 0, obwohl sie zumindest eines der Materialien gemeinsam haben. Besser wäre, wenn dieses Attribut z.B. mittels Jaccard-Koeffizient analysiert werden würde. Hier würde eine Ähnlichkeit von  $\frac{1}{3}$  herauskommen.

Dieser Ansatz ist neuartig und so noch nicht beschrieben worden. Daher wird er im Rahmen dieser Arbeit mit dem naheliegenderem Ansatz verglichen werden.

**String-Attribute** Für diese Attribute gibt sehr vielfältige Ansätze:

Die einfache Verwendung als kategorische Attribute fällt, wie bereits in der Problemstellung dargelegt, weg.

Eine weitere Möglichkeit wäre der Vergleich mittels String-Metrics (wie Jaro-Winkler oder Levensthein-Distance). Dieser Ansatz birgt ein Hauptproblem: Der K-Prototypes-Algorithmus berechnet für jedes Cluster einen Mittelpunkt. Die String-Metrics sind aber nur für einen paarweisen Vergleich von Strings geeignet. Ein “Mittelpunkt” kann hieraus nicht abgeleitet werden. Die Verwendung des Modus (wie bei kategorischen Attributen) scheidet ebenfalls aus. Wenn jeder Produkttitel individuell ist, dann ist der Modus immer 1. Dieser Ansatz funktioniert für das gewählte Clustering-Verfahren also nicht ohne weiteres.

Schließlich besteht der Ansatz, mit Tokenization zu arbeiten. In der Literatur wird dieser Ansatz aber lediglich für Datensets beschrieben, deren Datenpunkte aus ausschließlich einem String-Wert bestehen. Datensets mit mehreren String-Attributen oder in gemischter Form sind hingegen noch nicht zusammen verarbeitet worden.

Zur Lösung kam die Idee folgende auf: Zuvor wurde ein Ansatz zur Evaluation multi-kategorischer Attribute hergeleitet. Dieser Ansatz könnte für Strings ebenfalls verwendet werden. Zerlegt man einen solchen String in Tokens, so erhält man eine Art multi-kategorischen Wert. Bsp.: aus “Samsung Galaxy S20 128GB” wird {samsung, galaxi, s20, 128gb}. Solche Tokens lassen sich wieder mittels Jaccard-Koeffizienten vergleichen.

Auch dieser Ansatz wird mit der Verwendung der Strings als einfaches kategorisches Attribut verglichen werden.

### Mathematische Formulierung

Aus den genannten Ansätzen lässt sich insgesamt folgende Formell für die finale Distanz-Funktion ableiten:

$$d(x_1, x_2) = \frac{d_{num}(x_1, x_2) + n_{cat} \cdot d_{cat}(x_1, x_2) + d_{mul}(x_1, x_2) + |x_1^{num_{null}}| + |x_2^{num_{null}}|}{|attributes \text{ in } x_1 \cup x_2|} \quad (3.3)$$

$$d_{num}(x_1, x_2) = \sum_{i \in num} |x_1^i - x_2^i| \quad (3.4)$$

$$d_{cat}(x_1, x_2) = 1 - \frac{|x_1^{cat} \cap x_2^{cat}|}{|x_1^{cat} \cup x_2^{cat}|} \quad (3.5)$$

$$d_{mul}(x_1, x_2) = \sum_{i \in mul} \frac{|x_1^i \cap x_2^i|}{|x_1^i \cup x_2^i|} \quad (3.6)$$

Die numerischen Attribute werden mittels Manhattan-Distanz verrechnet. Da die Attribute vorher normalisiert worden sind, kann so maximal eine Distanz von 1 je numerischem Attribut entstehen. Die kategorischen Attribute werden mit dem beschriebenen inversen Jaccard-Koeffizienten berechnet (da asynchron). Der Koeffizient wird außerdem mit der Anzahl an kategorischen Attributen  $n_{cat}$  multipliziert, um ihn mit den anderen Metriken gleich zu gewichten. Die multi-kategorischen Attribute werden jeweils einzeln mittels inversem Jaccard-Koeffizienten verglichen.

Wie bereits ausführlich dargelegt, können in den Produktdaten von Akeneo immerzu null-Values vorkommen. Die kategorischen und multi-kategorischen Distanzmaße können dank des Jaccard-Koeffizienten problemlos damit umgehen. Die Minkowski-Metriken sind allerdings nur für tatsächlich vorhandene numerische Werte in beiden Produkten definiert. Daher wird  $d_{num}$  ausschließlich in diesem Fall aufgerufen. Für den Fall, dass null-Values in den numerischen Attributen von einem der beiden Produkte vorkommen, wird für jedes dieser Attribute 1 addiert. Die Attribute sind sich ja maximal unähnlich, da das eine Produkt es definiert und das andere nicht. In der Formel ist das mit  $|x_1^{num_{null}}|$  und  $|x_2^{num_{null}}|$  hinterlegt.

Der Divisor  $|attributes \text{ in } x_1 \cup x_2|$  sorgt für eine Normalisierung der Distanz in das Intervall  $[0; 1]$ . Dadurch werden Unterschiede zwischen Produkten, welche alleine aus der verschiedenen Anzahl an definierten Attributen entstehen, ausgeglichen.

Eine vereinfachte (serialisierte) Form dieser Funktion sieht folgendermaßen aus:

$$d(x_1, x_2) = \frac{\sum d'(x_1^i, x_2^i)}{|\text{attributes in } x_1 \cup x_2|} \quad (3.7)$$

$$d'(x_1^i, x_2^i) = \begin{cases} 1 & , i \text{ is numerical} \wedge (x_1^i \text{ is null} \vee x_2^i \text{ is null}) \\ |x_1^i - x_2^i| & , i \text{ is numerical} \\ 0 & , i \text{ is categorical} \wedge x_1^i = x_2^i \\ 1 & , i \text{ is categorical} \wedge x_1^i \neq x_2^i \\ 1 - \frac{|x_1^i \cap x_2^i|}{|x_1^i \cup x_2^i|} & , i \text{ is multi - categorical} \end{cases} \quad (3.8)$$

## 3.4 Evaluation

### 3.4.1 Kriterien

Das hergeleitete Clustering-Verfahren in seinen Varianten sollte auf seine Validität überprüft werden. Speziell wird getestet, ob mit diesem Verfahren “sinnvolle” Cluster gefunden werden. Die verschiedenen Aspekte dieser Sinnhaftigkeit werden nun anhand von drei Hauptkriterien definiert: Stabilität, Qualität und Erkennungsfähigkeit.

#### Stabilität

Der K-Prototypes-Algorithmus nutzt wie alle Verfahren dieser Klasse ein Initialisierungsverfahren, welches auf dem Zufall beruht. Da es sich um ein klassisches Minimierungsverfahren mit eventuellen lokalen Minima handelt, können mehrmalige Durchläufe über das gleiche Datenset verschiedene Clusterzuteilungen finden. Ein Verfahren, welches keinerlei Determinismus aufweist und je nach Uhrzeit komplett andere Ergebnisse liefert, ist allerdings in der Praxis nicht zu gebrauchen. Zudem sollte eine wohldefinierte Distanzfunktion in der Lage sein, die Datenpunkte eindeutig genug voneinander zu trennen, sodass trotz verschiedener Startpunkte die Zuteilung in Cluster stets ähnlich abläuft.

Zur Prüfung der Stabilität wird also wie folgt vorgegangen: Das Clustering wird stets mehrmals hintereinander ausgeführt. Anschließend wird die Ähnlichkeit der gefundenen Cluster mittels Adjusted-Rand-Index berechnet. Da ein hierarchisches Verfahren verwendet wird, wird jede Hierarchie-Ebene der verschiedenen Durchläufe betrachtet und anschließend der Durchschnitt aus allen Ähnlichkeitsmessungen über alle Hierarchie-Stufen und alle Durchläufe berechnet. Die Ähnlichkeit der Cluster sollte dabei nahe 100% sein.

#### Qualität

Ziel des Clusterings ist es, ordentlich voneinander getrennte Gruppen zu finden. Dies ist zum einen wichtig, um mit den Ergebnissen überhaupt weiterarbeiten und belastbare Aussagen zu den Gruppierungen finden zu können. Zum anderen sind wohl-separierte Cluster aber auch ein Garant für eine gewisse Resilienz vor neuen Datenpunkten, welche in Zukunft zum Datenset hinzugefügt werden. Schließlich ist es nicht sinnvoll, wenn sich jedes mal bei einem

neu hinzugefügten Produkt die Cluster stark verschieben. Sind sie ordentlich voneinander getrennt, so ist dies unwahrscheinlicher.

Daher wird ein gefundenes Clustering-Ergebnis mittels eines internen Index auf seine Qualität überprüft. Zum Einsatz können entweder der Silhouetten-Koeffizient oder Davies-Bouldin-Index kommen. Um bei der Evaluation keine zusätzlichen Fehlerquellen zu erzeugen, wird für die Berechnung der Metriken auf externe Bibliotheken zurückgegriffen werden. Es wird dabei der Index zum Einsatz kommen, der besser nutzbar ist.

Da das Clustering ein hierarchisches Verfahren verwendet, wird der interne Index für jede sinnvolle Hierarchie-Stufe ( $2 \leq k \leq n$ ) berechnet und anschließend der Durchschnitt gebildet.

### Erkennungsfähigkeit

Stabile und wohl-separierte Cluster zu finden, sind sehr sinnvolle Kriterien. Allerdings bewerten sie eher das Clustering an sich ohne einen wirklichen Bezug zur Realität. Dies soll mit diesem dritten Kriterium umgesetzt werden.

**Generationen und Modelle** In dem erstellten Datenset kommen objektiv erkennbare natürliche “Gruppen” vor. Die verschiedenen Generationen der Smartphones sowie die verschiedenen Modelle innerhalb einer Serie sind für uns Menschen offensichtliche Unterscheidungsmerkmale dieser Produkte. Daher sollte das Clustering-Verfahren ebenfalls in der Lage sein, diese Unterteilung zu finden.

Konkret werden die Generationen und Modelle in den “Akeneo-Categories” hinterlegt werden. Nach einem Cluster-Durchlauf werden anschließend die Hierarchie-Ebenen für die jeweiligen Kategorien abgefragt (z.B. wenn es drei Generationen im Datenset gibt, dann sollte das Clustering bei  $k = 3$  jedes Produkt der gleichen Generation in die gleiche Gruppe positioniert haben). Der Vergleich der Übereinstimmung mit “Akeneo-Categories” und Clustering-Ergebnis erfolgt über den Adjusted-Rand-Index.

**Duplikate** Außerdem befinden sich unter den Smartphones einige Duplikate. Diese werden ebenfalls mit einer entsprechenden “Category” in Akeneo markiert. Anschließend wird geprüft, ob die Duplikate über alle Hierarchie-Ebenen stets dem gleichen Cluster zugeordnet worden sind (was man von einem sinnvollen Clustering-Verfahren erwarten würde). Je nach Implementierung könnte hier der Jaccard-Koeffizient zum Einsatz kommen. Der Vergleich erfolgt nun über die Menge an zugeteilten Clustern beider Produkte. Diese sollte bei beiden sehr ähnlich sein und der Koeffizient entsprechend eine Ähnlichkeit nahe der 1 feststellen.

### 3.4.2 Aspekte

Für die Evaluation spielen eine Reihe von Teilaspekten bei der Auswertung eine besondere Rolle:

## Verarbeitung multi-kategorischer Attribute

Im vorherigen Abschnitt ist ein alternativer Ansatz zur Verarbeitung multi-kategorischer Attribute vorgestellt worden, welcher mit weniger Informationsverlust einhergeht. Daher sollte diese Art der Verarbeitung bessere Cluster produzieren, als wenn die multi-kategorischen Attribute in normale kategorische Attribute umgewandelt werden.

Zusätzlich kann diese Art der Verarbeitung auch auf String-Attribute angewendet werden, welches ebenfalls zu besseren Ergebnissen führen sollte.

## Auswahl der Attribute

Produkte in einem PIM-System weisen gerne mehrere hundert Attribute auf. Dadurch stellt sich die Frage, wie sich die Menge an Attributen auf das Clustering auswirkt. Dieser Einfluss wird zum einen mit den verschiedenen Produktfamilien (Smartphones und Hüllen) überprüft. Zum anderen könnte aber auch das gezielte Weglassen von vorhandenen Attributen ebenfalls die Ergebnisse verändern. Dabei sind verschiedene Ebenen denkbar:

Manche Typen von Attributen könnten wichtiger sein als andere. Z.B. stecken in String-Attribute potenziell mehr Informationen als in einem einfachen kategorischen Attribut. Andererseits könnte die "Unstrukturiertheit" der Strings diese Informationen verwaschen. Daher werden die **verschiedenen Typen** in unterschiedlichen Kombinationen geclustert und verglichen.

Über die Families in Akeneo können Attribute als **optional oder erforderlich** gekennzeichnet werden. Attribute mit vielen null-Werten könnten für das Clustering womöglich keine Rolle spielen. Daher werden die Cluster mit und ohne Verwendung optionaler Attribute überprüft.

Als letztes wird überprüft, ob eine **menschliche "Vorauswahl"** der Attribute sinnvollere Ergebnisse liefern kann. Um den Rahmen dieser Arbeit nicht komplett zu sprengen, wird der Autor die Auswahl nach persönlichem Ermessen treffen. Es ist aber denkbar, in Zukunft z.B. externe Statistiken heranzuziehen, um vor allem Attribute zu verwenden, welche für Kunden von besonderer Relevanz sind.

## Gewichtung der Attribute

In der gesamten bisherigen Herleitung ist stets die "Gleichgewichtung" der Attribute angestrebt worden (z.B. durch Normalisierung der numerischen Attribute etc.). Es ist aber denkbar, dass der komplette Ausschluss von Attributen mit zu viel Verlust an Informationen einhergeht. Daher wird gesondert überprüft, ob sich bessere Cluster bilden, wenn bestimmte Attribute (z.B. die menschliche "Vorauswahl" im vorherigen Absatz) im Verhältnis zu den anderen übergewichtet werden.

Die Übergewichtung lässt sich bewerkstelligen, indem die entsprechenden Attribute im Datenset während der Vorverarbeitung dupliziert werden.

### 3.4.3 Vorgehen

Die genannten Aspekte werden mit Hilfe der definierten Kriterien und Metriken evaluiert. Das Vorgehen ist dabei wie folgt:

Als erstes werden nur die “Smartphone-Hüllen” geclustert und in den genannten Aspekten überprüft. Da hier die Anzahl an Attributen überschaubar ist, können Fehler in der Implementierung leichter identifiziert und verbessert werden. Auch lässt sich der Einfluss einzelner Attribute auf das Endergebnis leichter bewerten.

Diese initialen Ergebnisse werden nun mit dem Datenset, welches nur aus den “Samsung Smartphones” besteht, gegengeprüft. Hier kommen deutlich mehr Attribute vor, sodass interessant zu sehen sein wird, ob sich die Ergebnisse aus “einfacheren” Produkten auch auf komplexere übertragen lassen.

Schließlich werden alle Produkte (Smartphones und Hüllen) gleichzeitig geclustert mithilfe der zuvor gewonnen Erkenntnisse (bspw. welche Attribut-Typen oder Gewichtungen die besten Ergebnisse geliefert haben etc.). Ziel hierbei ist es, dass die Ergebnisse des Clusterings der Teilsets und des Gesamtsets möglichst identisch sind. Ist das der Fall, so ist das gefundene Clustering-Verfahren tatsächlich “generisch”, da nicht einmal eine vorherige Trennung der Produktfamilien nötig ist.

Eine letzte Teilfrage wird in diesem Zusammenhang ebenfalls evaluiert: Wird das Clustering nur mit Attributen durchgeführt, welche in beiden Familien vorkommen, so könnten u.U. die “passenden” Smartphones und Hüllen in gemeinsame Gruppen sortiert werden. Dieser Versuch stellt eher eine Spekulation dar. Die Überprüfung ist in diesem Rahmen aber leicht durchführbar und bildet damit den Abschluss der Evaluation.



## 4 Implementierung

### 4.1 Überblick

- Akeneo-PIM Instanz auf adesso VM
- AKeneo Client und Cache
- Import von Daten aus Icecat
- Implementierung Clustering-Bibliothek + weitere (z.B. Datenaufbereitung)
- Durchführung des Clusterings selbst

### 4.2 Akeneo-PIM

#### 4.2.1 Installation und Deployment

- Überblick zur Projekt-Initialisierung
- Installation auf adesso VM
- Icecat Importer
- Anlegen verschiedener Basis-Attribute für Icecat

#### 4.2.2 Akeneo-Client und -Cache

- Klasse zur Interaktion mit Akeneo-API, behandelt Authentifizierung, Auflösung von Listen etc.
- Requests dauern mitunter einige Sekunden
- Cache fragt alle wichtigen Endpunkte ab und speichert JSON-payload in JSON-Dateien
- Abfrage von Daten aus dem Cache (lädt und parsed JSON-Dateien in Python-Datenstrukturen "dacite" Package)

### 4.3 Datenset

- Icecat Taxonomy => Import der Attribute
- händische Auswahl aller zu importierenden Attribute
- Suche nach Produkten auf icecat.biz
- Beschränkung auf Sponsoren-Vendors
- Import der Produkte mittels SKU und EAN; dann Icecat Importer "Enrich Products"
- Analyse der Daten und Korrektur verschiedener Fehler:

- Multi-Select werden als Single-Select importiert
  - einige Single-Selects als Text
  - Import diverser “number” Attribute schlägt fehl => Import als Text und anschließend Umwandlung
  - Fix: Attribut dupliziert mit suffix “\_fixed” und korrektem Typ
  - Außerdem: fehlerhafter Import diverser metrischer Attribute (wahrscheinlich Unit nicht richtig erkannt), aber ist bei allen gleich fehlerhaft importiert => wird eh normalisiert, daher ignoriert
- Statistiken zu Datenset: Attribute, Categories, Families, Produkte

## 4.4 Clustering

### 4.4.1 Überblick zu externen Libraries

- sklearn, nicovk
- erlauben keine null-Values !
- Umsetzung von Multi-Selects nicht möglich
- Orientierung an diesen Bibliotheken, aber keine Verwendung

### 4.4.2 Clustering-Package

- KMeans & Bisecting KMeans mit generischer Centroid-Klasse
- Centroid-Klasse erklären
- KMeans => Überblick
- Bisecting KMeans => Überblick

### 4.4.3 Distanzfunktionen

- Implementierung der Centroid-Klasse
- Map aus Attribut-Code zu Wert in Datentypen: float, string, set of strings
- Verarbeitung des Typs entsprechend
- Tracking des Means und der Modes je nach Attribut => Generierung eines künstlichen Datenpunktes für Distanzfunktion
- Gewichtung der Attribute möglich

### 4.4.4 Datenvorbereitung

- Produkte parsen:
- mit Filtern
- Normalisierung von float
- Tokenization von strings => Umwandlung in set of strings, genauso wie Multi-Selects

## 4.5 Evaluation

- Verwendung von metrics Funktionen aus SKLearn
- Durchführung mit Python-Packages in Jupyter-Notebooks
- Generierung aller möglichen Versuche
- Stabilität => 10 mal wiederholen da KMeans mit random-init, Vergleich mit `adj_rand_index`
- Qualität => Silhouettenkoeffizient, da in SKLearn am flexibelsten einsetzbar
- Korrektheit => Vergleich mit Categories (Smartphone Serien und Modellen) mit `adj_rand_index`

## 5 Auswertung

### 5.1 Datenset “Smartphone-Hüllen”

#### 5.1.1 Verarbeitung multi-kategorischer Attribute

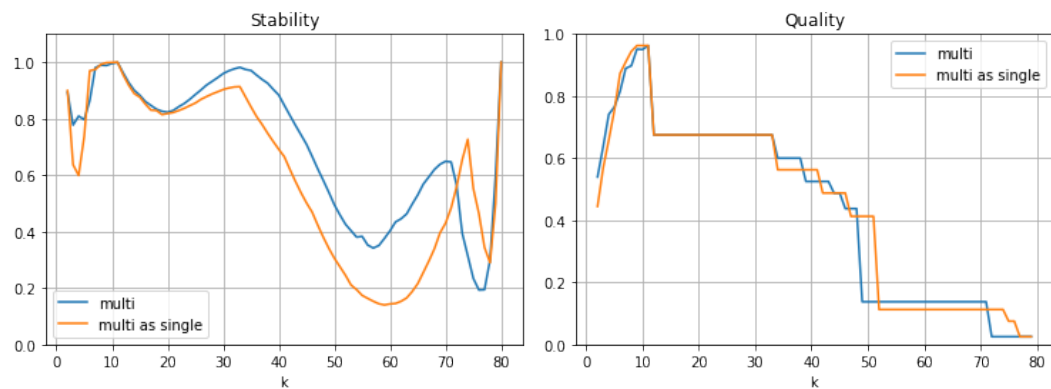


Abbildung 5.1: Stabilität und Qualität der Hüllen mit multi-kategorischem Attribut

Name	Stabilität	Qualität	Erkennung: Generation	Modell	$\emptyset$
multi	70.8%	43.9%	2.1%	4.3%	30.3%
multi als single	61.0%	44.5%	0.4%	4.3%	27.5%

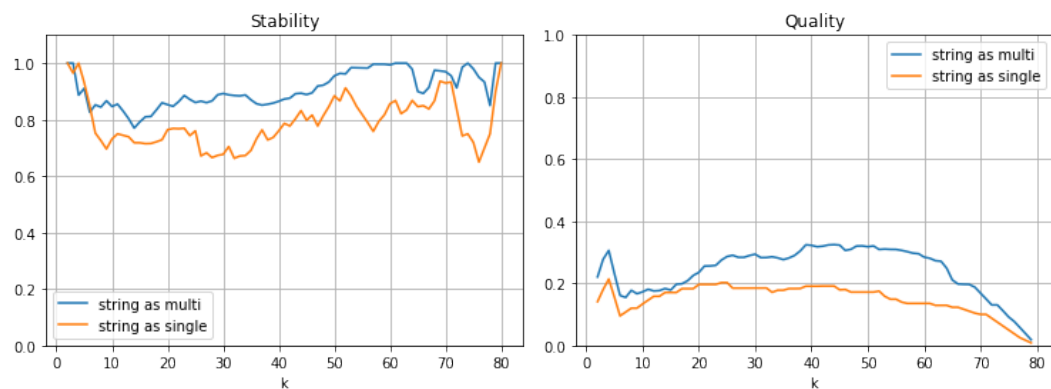


Abbildung 5.2: Stabilität und Qualität der Hüllen mit String-Attributen

Tabelle 5.2: Clustering der Hüllen mit String-Attributen

Name	Stabilität	Qualität	Erkennung: Generation	Modell	Ø
strings als multi	90.8%	24.1%	1.5%	17.5%	33.5%
strings als single	79.0%	15.0%	1.2%	7.5%	25.7%

### 5.1.2 Attribut-Auswahl

Vergleich nach Datentypen

Vergleich nach Erforderlichkeit

Vergleich nach menschlicher Auswahl

### 5.1.3 Attribut-Gewichtung

## 5.2 Datenset “Smartphones”

### 5.2.1 Verarbeitung multi-kategorischer Attribute

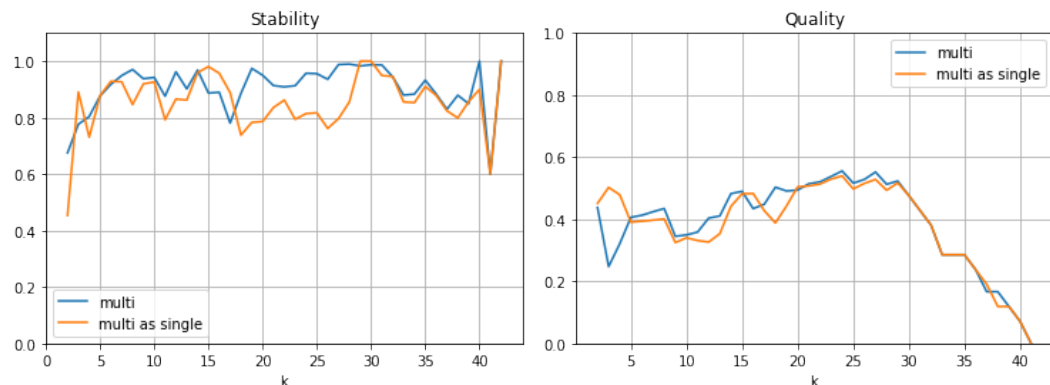


Abbildung 5.3: Stabilität und Qualität der Smartphones mit multi-kategorischen Attributen

Tabelle 5.3: Clustering der Smartphones mit multi-kategorischen Attributen

Name	Stabilität	Qualität	Erkennung: Generation	Modell	Duplikate	Ø
multi	90.5%	38.9%	39.8%	45.2%	89.2%	60.7%
multi als single	85.4%	38.5%	11.7%	30.0%	87.8%	50.7%

Tabelle 5.4: Clustering der Smartphones mit String-Attributen

Name	Stabilität	Qualität	Erkennung: Generation	Modell	Duplikate	Ø
strings als multi	88.2%	34.0%	64.3%	87.2%	100.0%	74.7%

Name	Stabilität	Qualität	Erkennung: Generation	Modell	Duplikate	Ø
strings als single	80.7%	23.7%	13.8%	58.6%	96.7%	54.7%

### 5.2.2 Attribut-Auswahl

Vergleich nach Datentypen

Vergleich nach Erforderlichkeit

Vergleich nach menschlicher Auswahl

### 5.2.3 Attribut-Gewichtung

## 5.3 Kombiniertes Datenset

### 5.3.1 Verwendung aller Attribute

### 5.3.2 Verwendung gemeinsamer Attribute

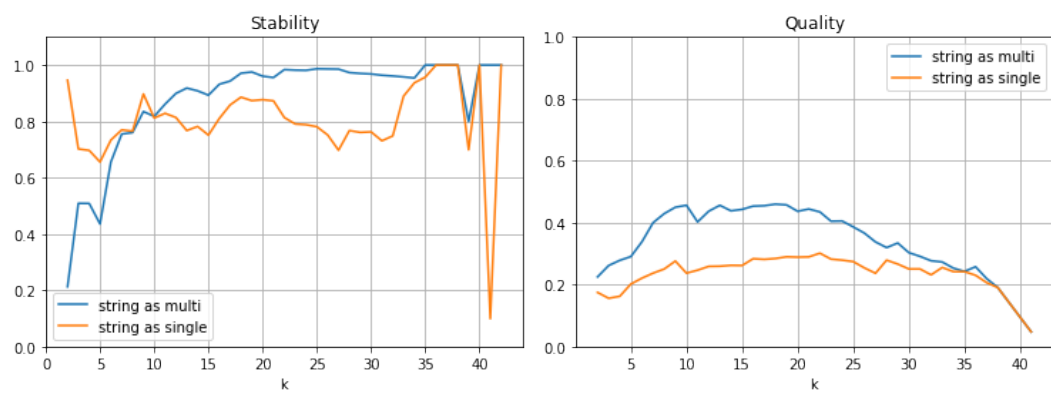


Abbildung 5.4: Stabilität und Qualität der Smartphones mit String-Attributen

## **6 Fazit und Ausblick**



## Quellenverzeichnis

- [Aken22a] AKENEO: *About Akeneo*. URL <https://www.akeneo.com/de/about-us/>. - abgerufen am 2022-04-04
- [Aken22e] AKENEO: *The REST API basics – Overview*. URL <https://api.akeneo.com/documentation/overview.html>. - abgerufen am 2022-04-04
- [Aken22g] AKENEO: *Concepts & resources – Catalog structure*. URL <https://api.akeneo.com/concepts/catalog-structure.html>. - abgerufen am 2022-04-04
- [Aken22h] AKENEO: *REST API endpoints – Create a new attribute*. URL [https://api.akeneo.com/api-reference-17.html#post\\_attributes](https://api.akeneo.com/api-reference-17.html#post_attributes). - abgerufen am 2022-04-04
- [Aken22c] AKENEO: *Recommended configuration*. URL [https://docs.akeneo.com/5.0/technical\\_architecture/technical\\_information/recommended\\_configuration.html](https://docs.akeneo.com/5.0/technical_architecture/technical_information/recommended_configuration.html). - abgerufen am 2022-04-04
- [Aken22d] AKENEO: *Install Akeneo PIM for development with Docker*. URL [https://docs.akeneo.com/5.0/install\\_pim/docker/installation\\_docker.html](https://docs.akeneo.com/5.0/install_pim/docker/installation_docker.html). - abgerufen am 2022-04-04
- [Aken22b] AKENEO: *Akeneo PIM in a nutshell*. URL <https://docs.akeneo.com/5.0/index.html>. - abgerufen am 2022-04-04
- [Aken22f] AKENEO: *Concepts & resources – Products*. URL <https://api.akeneo.com/concepts/products.html>. - abgerufen am 2022-04-04
- [Aken22i] AKENEO: *Concepts & resources – Target market settings*. URL <https://api.akeneo.com/concepts/target-market-settings.html>. - abgerufen am 2022-04-05
- [ArVa16] ARORA, PREETI ; VARSHNEY, SHIPRA ; u. a.: Analysis of k-means and k-medoids algorithm for big data. In: *Procedia Computer Science* Bd. 78, Elsevier (2016), S. 507–512
- [BiHa09] BISSANTZ, NICOLAS ; HAGEDORN, JÜRGEN: Data Mining (Datenmustererkennung). In: *Wirtschaftsinformatik* Bd. 51, Springer (2009), Nr. 1, S. 139–144
- [Bosl12] BOSLAUGH, SARAH: *Statistics in a nutshell: A desktop quick reference* : O'Reilly Media, 2012
- [Cha07] CHA, SUNG-HYUK: Comprehensive survey on distance/similarity measures between probability density functions. In: *City* Bd. 1 (2007), Nr. 2, S. 1
- [ChZhYi19] CHEN, HONGSHEN ; ZHAO, JIASHU ; YIN, DAWEI: Fine-grained product categorization in e-commerce. In: *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 2019, S. 2349–2352
- [CRF03] COHEN, WILLIAM W ; RAVIKUMAR, PRADEEP ; FIENBERG, STEPHEN E ; u. a.: A Comparison of String Distance Metrics for Name-Matching Tasks. In: *IJWeb*. Bd. 3 : Citeseer, 2003, S. 73–78
- [Cui21] CUI, YIMIN: Intelligent recommendation system based on mathematical modeling in personalized data mining. In: *Mathematical Problems in Engineering* Bd. 2021, Hindawi (2021)

- [DaBo79] DAVIES, DAVID L ; BOULDIN, DONALD W: A cluster separation measure. In: *IEEE transactions on pattern analysis and machine intelligence*, IEEE (1979), Nr. 2, S. 224–227
- [DE22] DE, SAMSUNG: *Galaxy S Serie*. URL <https://www.samsung.com/de/smartphones/galaxy-s/>. - abgerufen am 2022-04-05
- [Defa77] DEFAYS, DANIEL: An efficient algorithm for a complete link method. In: *The Computer Journal* Bd. 20, Oxford University Press (1977), Nr. 4, S. 364–366
- [DoBi22] DOGAN, ALICAN ; BIRANT, DERYA: K-centroid link: a novel hierarchical clustering linkage method. In: *Applied Intelligence* Bd. 52, Springer (2022), Nr. 5, S. 5537–5560
- [Hamm50] HAMMING, RICHARD W: Error detecting and error correcting codes. In: *The Bell system technical journal* Bd. 29, Nokia Bell Labs (1950), Nr. 2, S. 147–160
- [Hart75] HARTIGAN, JOHN A: *Clustering algorithms* : John Wiley & Sons, Inc., 1975
- [Huan98] HUANG, ZHEXUE: Extensions to the k-means algorithm for clustering large data sets with categorical values. In: *Data mining and knowledge discovery* Bd. 2, Springer (1998), Nr. 3, S. 283–304
- [HuAr85] HUBERT, LAWRENCE ; ARABIE, PHIPPS: Comparing partitions. In: *Journal of classification* Bd. 2, Springer (1985), Nr. 1, S. 193–218
- [Icec22a] ICECAT: *About Icecat*. URL <https://icecat.com/about-us/>. - abgerufen am 2022-04-05
- [Icec22b] ICECAT: *Icecat: open feed with product information, data-sheets for ecommerce*. URL <https://icecat.biz/en/menu/manufacturers/index.html>. - abgerufen am 2022-04-13
- [Jaro95] JARO, MATTHEW A: Probabilistic linkage of large public health data files. In: *Statistics in medicine* Bd. 14, Wiley Online Library (1995), Nr. 5-7, S. 491–498
- [JiCh17] JIA, HONG ; CHEUNG, YIU-MING: Subspace clustering of categorical and numerical data with an unknown number of clusters. In: *IEEE transactions on neural networks and learning systems* Bd. 29, IEEE (2017), Nr. 8, S. 3308–3325
- [KaRo09] KAUFMAN, LEONARD ; ROUSSEEUW, PETER J: *Finding groups in data: an introduction to cluster analysis*. Bd. 344 : John Wiley & Sons, 2009
- [King15] KING, RONALD S: *Cluster analysis and data mining: An introduction* : Stylus Publishing, LLC, 2015
- [KoLo12] KOU, GANG ; LOU, CHUNWEI: Multiple factor hierarchical clustering algorithm for large scale web page and search engine clickstream data. In: *Annals of Operations Research* Bd. 197, Springer (2012), Nr. 1, S. 123–134
- [KRRT01] KUMAR, RAVI ; RAGHAVAN, PRABHAKAR ; RAJAGOPALAN, SRIDHAR ; TOMKINS, ANDREW: Recommendation systems: A probabilistic analysis. In: *Journal of Computer and System Sciences* Bd. 63, Elsevier (2001), Nr. 1, S. 42–61
- [MiCo88] MILLIGAN, GLENN W ; COOPER, MARTHA C: A study of standardization of variables in cluster analysis. In: *Journal of classification* Bd. 5, Springer (1988), Nr. 2, S. 181–204
- [MWDM64] MACNAUGHTON-SMITH, P ; WILLIAMS, WT ; DALE, MB ; MOCKETT, LG: Dissimilarity analysis: a new technique of hierarchical sub-division. In: *Nature* Bd. 202, Nature Publishing Group (1964), Nr. 4936, S. 1034–1035

- [OhKi19] OH, YOORI ; KIM, YOONHEE: A resource recommendation method based on dynamic cluster analysis of application characteristics. In: *Cluster Computing* Bd. 22, Springer (2019), Nr. 1, S. 175–184
- [Pimc22] PIMCORE: *What, why and how of product information management*. URL <https://pimcore.com/en/what-is-pim>. - abgerufen am 2022-04-04
- [PWMO19] PAPP, STEFAN ; WEIDINGER, WOLFGANG ; MEIR-HUBER, MARIO ; ORTNER, BERNHARD ; LANGS, GEORG ; WAZIR, RANIA: *Handbuch Data Science: Mit Datenanalyse und Machine Learning Wert aus Daten generieren* : Carl Hanser Verlag GmbH Co KG, 2019
- [RAAQ11] RENDÓN, ERÉNDIRA ; ABUNDEZ, ITZEL ; ARIZMENDI, ALEJANDRA ; QUIROZ, ELVIA M: Internal versus External cluster validation indexes. In: *International Journal of computers and communications* Bd. 5 (2011), Nr. 1, S. 27–34
- [Rand71] RAND, WILLIAM M: Objective criteria for the evaluation of clustering methods. In: *Journal of the American Statistical association* Bd. 66, Taylor & Francis (1971), Nr. 336, S. 846–850
- [RaRa11] RAJALINGAM, N ; RANJINI, K: Hierarchical Clustering Algorithm - A Comparative Study. In: *International Journal of Computer Applications* Bd. 19, Citeseer (2011), Nr. 3, S. 42–46
- [Rous87] ROUSSEEUW, PETER J: Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. In: *Journal of computational and applied mathematics* Bd. 20, Elsevier (1987), S. 53–65
- [RoZa09] ROBERTSON, STEPHEN ; ZARAGOZA, HUGO: *The probabilistic relevance framework: BM25 and beyond* : Now Publishers Inc, 2009
- [Sibs73] SIBSON, ROBIN: SLINK: an optimally efficient algorithm for the single-link cluster method. In: *The computer journal* Bd. 16, Oxford University Press (1973), Nr. 1, S. 30–34
- [Ste104] STEINLEY, DOUGLAS: Properties of the Hubert-Arable Adjusted Rand Index. In: *Psychological methods* Bd. 9, American Psychological Association (2004), Nr. 3, S. 386
- [StKaKu00] STEINBACH, MICHAEL ; KARYPIS, GEORGE ; KUMAR, VIPIN: A comparison of document clustering techniques (2000)

## Anhang

- [Link zum Repo](#)
- [Überblick zum Repo](#)

## Selbstständigkeitserklärung

Ich, Hannes Dröse, versichere hiermit, dass ich die vorliegende Masterarbeit mit dem Thema

*Generisches Clustern hoch-komplexer Produktdaten*

selbstständig und nur unter Verwendung der angegebenen Quellen und Hilfsmittel angefertigt habe.

Erfurt, den 29.04.2022

.....

Hannes Dröse