



Fachhochschule der Wirtschaft  
FHDW Bielefeld

Studienarbeit im Modul Classification and Clustering

## Kundensegmentierung mit RetailRocket und Online Retail II

Vorgestellt von:

**Christian Roth**

Senner Straße 68

33647 Bielefeld

`christian.roth@edu.fhdw.de`

Studiengang:

Wirtschaftsinformatik mit Schwerpunkt Data Science (*M.Sc.*)

Prüferin:

Prof. Dr. Yvonne Gorniak

Eingereicht am:

08. September 2025 in Bielefeld

## Gendererklärung

Aus Gründen der besseren Lesbarkeit wird auf die gleichzeitige Verwendung der Sprachformen männlich, weiblich und divers (m/w/d) verzichtet. Sämtliche Personenbezeichnungen gelten gleichermaßen für alle Geschlechter.

# Abstract

Diese Arbeit untersucht die Kundensegmentierung auf Basis zweier offener, datenschutzkonformer E-Commerce-Datensätze: RetailRocket mit Clickstream-Ereignissen und Online Retail II mit Transaktionen. Ziel ist es, robuste und geschäftlich interpretierbare Segmente zu gewinnen, die sowohl Nutzungsverhalten als auch Kaufmuster abbilden. Methodisch orientiert sich die Studie am CRISP-DM-Zyklus. Beide Quellen werden in ein einheitliches Ereignisschema überführt und durch ein gemeinsames Merkmalset beschrieben, das unter anderem RFM-Metriken, Konversionsraten, Zeitmuster, Diversität und Zwischenkaufintervalle umfasst. Ergänzend wird Text Mining auf Kategoriepfaden eingesetzt.

Zur Bildung der Cluster werden zwei kontrastierende Verfahren angewandt: K-Means als zentroidbasiertes und HDBSCAN als dichtebasiertes Verfahren. Die Qualität der Lösungen wird mit Silhouette- und Davies–Bouldin-Index sowie mit dem DBCV-Index bewertet, die Stabilität über Bootstrapping und den Adjusted Rand Index geprüft und die Clustertendenz mit der Hopkins-Statistik abgesichert. Darüber hinaus werden die Verfahren durch Ablationsanalysen und systematische Parameterabstimmungen optimiert.

Die Ergebnisse zeigen, dass sich konsistente Kundensegmente aus unterschiedlichen Datenwelten ableiten lassen und dass bestimmte Merkmalsgruppen die Clusterbildung besonders stark beeinflussen. Die Arbeit liefert eine reproduzierbare Referenzpipeline für Kundensegmentierung im CRM-Kontext und macht transparent, welche Zielkonflikte zwischen K-Means und HDBSCAN in Praxisdaten auftreten.

# Inhaltsverzeichnis

<b>Gendererklärung</b>	<b>i</b>
<b>Abstract</b>	<b>ii</b>
<b>Inhaltsverzeichnis</b>	<b>iii</b>
<b>Abbildungsverzeichnis</b>	<b>iv</b>
<b>Tabellenverzeichnis</b>	<b>v</b>
<b>1 Einleitung</b>	<b>1</b>
1.1 Problemstellung . . . . .	1
1.2 Zielsetzung . . . . .	1
1.3 Vorgehensweise . . . . .	1
<b>2 Grundlagen</b>	<b>2</b>
2.1 Kundensegmentierung im CRM Kontext . . . . .	2
2.2 Zwei unabhängige Datensichten und deren Harmonisierung im Lichte von CRISP-DM . . . . .	3
2.3 Datenquellen und Datenschutz RetailRocket und Online Retail II . . . . .	3
2.4 Verfahren und Qualitätsmaße im Überblick K Means, HDBSCAN und pas- sende Indizes . . . . .	4
2.5 Übergang zum Praxisteil Messäquivalenz und minimaler Merkmalsatz . . .	4
<b>3 Praxisteil</b>	<b>5</b>
3.1 Datenverständnis RetailRocket Ereignisschema und Eigenschaften . . . . .	6
3.2 Datenverständnis Online Retail II Transaktionsschema und Eigenschaften .	6
3.3 Vereinheitlichung beider Datensätze gemeinsames Ereignisschema . . . . .	6
3.4 Datenaufbereitung und Feature Engineering RFM Konversionsproxys Zeit- muster Diversität Interkaufintervalle . . . . .	6
3.5 Text Mining auf Kategoriepfaden TF IDF optional Item2Vec . . . . .	6

3.6	Clustering Ansatz 1 K-Means Auswahl der Clusterzahl über Silhouette und Davies Bouldin . . . . .	6
3.7	Clustering Ansatz 2 HDBSCAN Parameterwahl und DBCV . . . . .	6
3.8	Stabilität und Robustheit Bootstrapping und Adjusted Rand Index . . . . .	6
3.9	Segmentprofiling KPI Profile und Top Kategorien je Cluster . . . . .	6
3.10	Reproduzierbarkeit Notebook Struktur Hyperparameter Versionierung . . . . .	6
<b>4</b>	<b>Evaluation und Ergebnisse</b>	<b>7</b>
4.1	Clusterqualität je Datensatz RetailRocket und Online Retail II . . . . .	7
4.2	Vergleich der Lösungen K-Means gegenüber HDBSCAN . . . . .	7
4.3	Sensitivitätsanalyse Feature Varianten und Hyperparameter . . . . .	7
4.4	Grenzen Datenqualität Sparsity Cold Start . . . . .	7
<b>5</b>	<b>Fazit und Ausblick</b>	<b>8</b>
5.1	Beantwortung der Forschungsfrage . . . . .	8
5.2	Nutzen für CRM, Kampagnensteuerung und nächste Schritte . . . . .	8
5.3	Ausblick weitere Text Mining Merkmale und Uplift Experimente . . . . .	8
	<b>Literaturverzeichnis</b>	<b>9</b>
	<b>Anhang</b>	<b>10</b>
	Anhangsverzeichnis . . . . .	10
	<b>A: Zero-Trust-Modell</b>	<b>11</b>
	A1: Zero-Trust-Säulen . . . . .	11
	<b>B: Erste Analyse der Daten</b>	<b>12</b>
	B1: Lorem . . . . .	12
	<b>Glossar</b>	<b>13</b>

# Abbildungsverzeichnis

# Tabellenverzeichnis

# 1 Einleitung

## 1.1 Problemstellung

Unternehmen segmentieren Kunden oft je Datenquelle, etwa nur im Clickstream oder nur in Transaktionen. So entstehen Cluster mit begrenzter Übertragbarkeit. CRM benötigt jedoch stabile, interpretierbare Segmente, die Nutzungs- und Kaufverhalten zusammenführen und ohne personenbezogene Daten auskommen.

## 1.2 Zielsetzung

Die Arbeit vereinheitlicht RetailRocket und Online Retail II zu einem Ereignisschema und beschreibt Kunden mit einem gemeinsamen Merkmalset aus RFM, Konversionsraten, Zeitmustern, Diversität, Zwischenkaufintervallen sowie Text-Mining auf Kategoriepfaden. Darauf aufbauend werden K-Means und HDBSCAN angewandt und verglichen. Qualität, Stabilität und Clustertendenz werden mit Silhouette-Index, Davies-Bouldin-Index, DBCV, Bootstrapping, Adjusted Rand Index und der Hopkins-Statistik bewertet.

## 1.3 Vorgehensweise

Die Untersuchung folgt CRISP-DM mit Datenverständnis und EDA, Vereinheitlichung, Merkmalerstellung, Modellierung, Evaluation und Profiling zu CRM-Segmenten. Leitend ist die Frage, wie robust und übertragbar Kundencluster aus RetailRocket und Online Retail II sind, die mit K-Means und HDBSCAN gebildet werden, gemessen mit Silhouette-Index, Davies-Bouldin-Index, DBCV und Adjusted Rand Index, und wie sie zu geschäftlich interpretierbaren Segmenten profiliert werden können.



## 2 Grundlagen

CRISP-DM bildet den Prozessrahmen dieser Arbeit. Wirth und Hipp schreiben: „The life cycle of a data mining project is broken down in six phases“ (Wirth und Hipp (2000), S. 4). Die Phasen reichen von Business Understanding über Data Understanding, Data Preparation, Modeling und Evaluation bis zu Deployment; sie werden iterativ durchlaufen (Wirth und Hipp (2000), S. 4–7). Wirth und Hipp sprechen von „this highly iterative, creative process with many parallel activities“ und halten fest: „it is never the case that a phase is completely done before the subsequent phase starts“ (Wirth und Hipp (2000), S. 9). Dieser Rahmen stellt sicher, dass Ziele, Daten, Modelle und Bewertung systematisch aufeinander aufbauen.

### 2.1 Kundensegmentierung im CRM Kontext

Clustering wird im CRM als Verfahren des unüberwachten Lernens eingesetzt, das natürliche Gruppierungen erkennt und Kunden in intern kohäsive Gruppen einordnet. „Clustering techniques identify meaningful natural groupings of records and group customers into distinct segments with internal cohesion“ (Tsipitsis und Chorianopoulos (2010), S. 39). Ziel ist die Bildung unterscheidbarer Kundentypologien, „so that they can be marketed more effectively“ (Tsipitsis und Chorianopoulos (2010), S. 40). Der Nutzen einer Segmentierung misst sich daran, ob die resultierenden Typologien transparent, aussagekräftig und handlungsleitend sind. „The value of each solution depends on its ability to represent transparent, meaningful, and actionable customer typologies“ (Tsipitsis und Chorianopoulos (2010), S. 129).

## 2.2 Zwei unabhängige Datensichten und deren Harmonisierung im Lichte von CRISP-DM

Business Understanding konkretisiert das Ziel auf robuste und übertragbare Segmente. Data Understanding führt zwei unabhängige, aber komplementäre Datensichten ein: eine verhaltensnahe Ereignissicht und eine wertbasierte Rechnungssicht. Erst die Kombination erlaubt zu prüfen, ob Clusterstrukturen auf gleich definierten Merkmalen in unterschiedlichen Datenwelten wiederkehren und damit eher domänenweit gültige Muster darstellen als datensatzspezifische Effekte (Wirth und Hipp (2000), S. 5; Tsiptsis und Chorianopoulos (2010), S. 39–40).

In Data Preparation werden beide Quellen in ein gemeinsames Ereignisschema überführt, um Messäquivalenz herzustellen. Praktisch umfasst dies konsistente Felder wie `customer_id`, `item_id`, `timestamp` und `event_type` sowie bei Bedarf `quantity`, `price`, `revenue` und eine Kategoriendarstellung. Damit lassen sich zentrale Konstrukte wie Recency, Frequency, Interkaufintervalle, Diversität und Zeitmuster in beiden Datensichten identisch operationalisieren (Wirth und Hipp (2000), S. 5–6).

In Modeling werden komplementäre Verfahren auf denselben Kernmerkmalen je Datensicht kalibriert, etwa ein zentroidbasiertes und ein dichte-basiertes Verfahren, um unterschiedliche Strukturannahmen abzudecken (Wirth und Hipp (2000), S. 6). Die Evaluation kombiniert interne Qualitätsmaße, die zur jeweiligen Methode passen, mit Stabilitäts- und Übertragbarkeitstests zwischen den Datensichten, bevor Ergebnisse reproduzierbar dokumentiert werden (Evaluation und Deployment, Wirth und Hipp (2000), S. 6–7, S. 9).

## 2.3 Datenquellen und Datenschutz RetailRocket und Online Retail II

Beide Datensätze sind offen zugänglich und anonymisiert, sodass in dieser Arbeit keine personenbezogenen Daten verarbeitet werden. RetailRocket bildet verhaltensnahe Ereignisfolgen ab, Online Retail II die wertbasierte Rechnungssicht; die Detailbeschreibung folgt erst im Praxisteil. Dieser kurze Abschnitt begründet lediglich die Verwendbarkeit und verankert

die Quellen im Prozessrahmen.

## **2.4 Verfahren und Qualitätsmaße im Überblick K Means, HDBSCAN und passende Indizes**

K Means steht als zentroidbasierte Baseline für kompakte Strukturen, HDBSCAN als dichtebasierter Ansatz für ungleich dichte Cluster und Rauschen. Bewertet wird mit zu den Verfahren passenden internen Indizes und mit Stabilitätstests; entscheidend ist die spätere Übertragbarkeit zwischen den Datensichten. Formelhafte Details und Parametertabellen entfallen hier und erscheinen erst mit den Ergebnissen.

## **2.5 Übergang zum Praxisteil Messäquivalenz und minimaler Merkmalssatz**

Für einen fairen Vergleich werden beide Quellen in ein gemeinsames Ereignisschema überführt und ein identischer Kern von Merkmalen definiert. Dieser minimale, quellenübergreifende Merkmalssatz trägt die Modellierung und die Stabilitätsanalysen, während zusätzliche merkmals- oder datensatzspezifische Details erst in Kapitel 3 genutzt und begründet werden. So bleibt der Grundlagenteil schlank und führt direkt in die praktische Umsetzung.



### 3 Praxisteil

3.1 Datenverständnis RetailRocket Ereignisschema und Eigenschaften

3.2 Datenverständnis Online Retail II Transaktionsschema und Eigenschaften

3.3 Vereinheitlichung beider Datensätze gemeinsames Ereignisschema

3.4 Datenaufbereitung und Feature Engineering RFM Konversionsproxys Zeitmuster Diversität Interkaufintervalle

3.5 Text Mining auf Kategoriepfaden TF IDF optional Item2Vec

3.6 Clustering Ansatz 1 K-Means Auswahl der Clusterzahl über Silhouette und Davies Bouldin

3.7 Clustering Ansatz 2 HDBSCAN Parameterwahl und DBCV

3.8 Stabilität und Robustheit Bootstrapping und Adjusted Rand Index

3.9 Segmentprofiling KPI Profile und Top Kategorien je Cluster

## 4 Evaluation und Ergebnisse

4.1 Clusterqualität je Datensatz RetailRocket und Online Retail II

4.2 Vergleich der Lösungen K-Means gegenüber HDBSCAN

4.3 Sensitivitätsanalyse Feature Varianten und Hyperparameter

4.4 Grenzen Datenqualität Sparsity Cold Start

## 5 Fazit und Ausblick

Was eine [Multi-Faktor-Authentifizierung \(MFA\)](#) ist, wird im Glossar beschrieben. Auch `glspl` und `glslink` sind möglich.

### 5.1 Beantwortung der Forschungsfrage

### 5.2 Nutzen für CRM, Kampagnensteuerung und nächste Schritte

### 5.3 Ausblick weitere Text Mining Merkmale und Uplift Experimente

## Literaturverzeichnis

Tsiptsis, K., und Chorianopoulos, A. (2010). *Data Mining Techniques in CRM: Inside Customer Segmentation* (1. Aufl.). Wiley. <https://doi.org/10.1002/9780470685815>

Wirth, R., und Hipp, J. (2000). *CRISP-DM: Towards a Standard Process Model for Data Mining*. 11. <http://cs.unibo.it/~danilo.montesi/CBD/Beatriz/10.1.1.198.5133.pdf>

The CRISP-DM (CRoss Industry Standard Process for Data Mining) project proposed a comprehensive process model for carrying out data mining projects. The process model is independent of both the industry sector and the technology used. In this paper we argue in favor of a standard process model for data mining and report some experiences with the CRISP-DM process model in practice.



# Anhang

## Anhangsverzeichnis

<b>A</b> Zero-Trust-Modell .....	11
<b>A.1</b> Zero-Trust-Säulen .....	11
<b>B</b> Erste Analyse der Daten .....	12
<b>B.1</b> Ergebnisse der ersten Analyse .....	12

## **A: Zero-Trust-Modell**

### **A1: Zero-Trust-Säulen**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **B: Erste Analyse der Daten**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

### **B1: Lorem**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

# Glossar

**Multi-Faktor-Authentifizierung (MFA)** eine Sicherheitsmethode, die eine zusätzliche Verifikationsebene erfordert, beispielsweise durch eine Kombination aus Passwort und biometrischer Authentifizierung.

## Ehrenwörtliche Erklärung

Hiermit erkläre ich, dass ich die vorliegende Studienarbeitselbstständig angefertigt habe.  
Es wurden nur die in der Arbeit ausdrücklich benannten Quellen und Hilfsmittel benutzt.  
Wörtlich oder sinngemäß übernommenes Gedankengut habe ich als solches kenntlich gemacht.  
Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Bielefeld, den 08. September 2025



---

Christian Roth