

Big Data / Data Mining

Einführung

Christoph Menzer
Fachgruppe Informatik

26. Oktober 2020

Motivation

- Warum Data Mining?
- Was ist Data Mining?
- Was können Sie erwarten?

Überblick

- Maschinelles Lernen
- Vorgehensmodelle
- Anwendungsklassen
- Lernstrategien
- Analysetechniken
- Datenvorverarbeitung
- Werkzeuge

Grundlagen

- Grundbegriffe
- Datentypen
- Abstands- und Ähnlichkeitsmaße

RapidMiner

- RapidMiner Studio
- Funktionsweise
- Attribute
- Datentypen

Literatur

- ▶ in jeder Sekunde werden Tausende oder Millionen von Daten gemessen, erhoben und gespeichert
- ▶ Warum? Weil wir es können und es kaum oder kein Geld kostet
- ▶ sowohl im beruflichen als auch im privaten Umfeld
- ▶ Daten und (vor allem) daraus abgeleitete Information sind Handelsware und zu Produktionsfaktor geworden

Data Mining ist der Prozess des Entdeckens bedeutsamer neuer Zusammenhänge, Muster und Trends durch die Analyse großer Datensätze mittels Mustererkennung sowie statistischer und mathematischer Verfahren. (Ittner 2019)

Data Mining [ist] eine Sammlung von Techniken, Methoden und Algorithmen für die Analyse von Daten, die somit auch Grundtechniken für neuere und komplexere Ansätze, wie das Business Intelligence oder auch Big Data darstellen. (Cleve und Lämmel 2014, S. 3)

- ▶ Techniken und Methoden zur Auswertung großer Datenmengen
- ▶ Einführung in Prozesse und Techniken der Mustererkennung in *strukturierten Daten* (Teilgebiet Data Mining, vgl. Text Mining, Web Mining, Künstliche Intelligenz (KI))
- ▶ Grundbegriffe und Vorgehensweisen des Data Mining
- ▶ Anwendungsklassen und Einsatzmöglichkeiten
- ▶ Verfahren zur Datenanalyse
- ▶ Datenselektion und -integration, Datensäuberung, Datenreduktion, Datentransformation (Datenvorverarbeitung, DVV)
- ▶ Bewertung und Visualisierung der Resultate

- ▶ computerbasierte Lernverfahren, die aus Eingabeinformationen Wissen (vgl. S. 15) generieren können
- ▶ einfachste Lernstrategie: Auswendiglernen (abspeichern, z.B. in Datenbank; hardcoded source code)
- ▶ angestrebt wird aber: Verständnis von Zusammenhängen und Hintergründen (Muster oder Abhängigkeiten)
- ▶ Induktives Lernen: aus Beispielen lernen und verallgemeinern

- ▶ statistische Verfahren zur Aufdeckung von Zusammenhängen

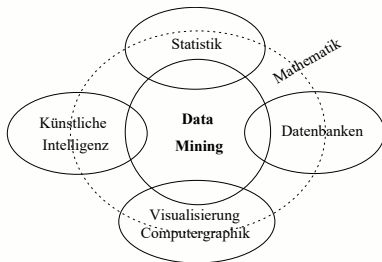
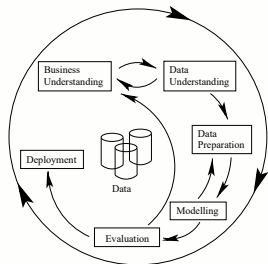
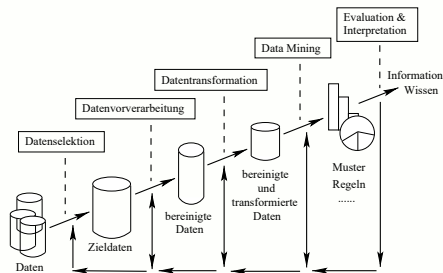


Abbildung: Interdisziplinarität (Quelle: Cleve und Lämmel 2014, S. 12)

„Data Mining ist ein typisch iterativer Prozess: Man probiert etwas aus, berechnet ein Resultat und prüft dieses“ (Cleve und Lämmel 2014, S. 197)



(a) CRISP-Modell



(b) Ablauf nach Fayyad

Abbildung: Quelle: Cleve und Lämmel 2014, S. 5 f.

CRISP: Lebenszyklus in 6 Etappen

1. Verstehen der Aufgabe (Business understanding)
2. Verständnis der Daten (Data understanding)
3. Datenvorbereitung (Data preparation)
4. Modellbildung (Modeling)
5. Evaluation (Evaluation)
6. Einsatz im und Konsequenzen für das Unternehmen (Deployment)

Prognosemethoden

- ▶ Vorhersage numerischer Werte
- ▶ Klassifikation (Vorhersage)

Entdeckungsmethoden

- ▶ Cluster-Analyse (Zerlegung)
- ▶ Assoziationsanalyse (Zusammenhänge)

Unüberwachtes Lernen

- ▶ zu entdeckende Muster sind unbekannt
- ▶ Gruppierung oder Klassifikation ist nicht gegeben
- ▶ nicht vergleichbar
- ▶ z.B. Clustering

Überwachtes Lernen

- ▶ Resultat ist für Beispiele gegeben
- ▶ vergleichbar
- ▶ z.B. Vorhersage numerischer Werte, Klassifikation
- ▶ Annahme: Beispieldaten sind repräsentativ

- ▶ Lineare Regression
- ▶ k-Nearest Neighbour
- ▶ Bayes-Klassifikator
- ▶ Entscheidungsbaumlernen
- ▶ Support Vector Machine
- ▶ Künstliches Neuronales Netz
- ▶ k-Means
- ▶ A-Priori-Algorithmus

- ▶ schwierigste Aufgabe
- ▶ Ausreißer (fehlende, falsche oder widersprüchliche Werte)
- ▶ Datentransformation, Skalierung
- ▶ Qualität der Daten entscheidend für Ergebnis
- ▶ etwa 80% des Gesamtaufwandes
- ▶ Phase „Data understanding“ von großer Bedeutung
- ▶ Im ersten Schritt einfache statistische Tests durchführen
- ▶ Arten der Datenvorbereitung
 - Datenselektion und -integration
 - Datensäuberung
 - Datenreduktion
 - Datentransformation

- ▶ Texteditoren
- ▶ Tabellenkalkulationsprogramme
- ▶ ...
- ▶ KNIME (Konstanz Information Miner)
- ▶ WEKA (Waikato Environment for Knowledge Analysis)
- ▶ JavaNNS für SNNS (Stuttgarter Neuronale Netze Simulator)
- ▶ RapidMiner
- ▶ IBM SPSS Modeler
- ▶ R
- ▶ ...

In vielen Data-Warehouse- und Datenbank-Systemen sind Data-Mining-Komponenten integriert.

- ▶ Texteditoren
- ▶ Tabellenkalkulationsprogramme
- ▶ ...
- ▶ KNIME (Konstanz Information Miner)
- ▶ WEKA (Waikato Environment for Knowledge Analysis)
- ▶ JavaNNS für SNNS (Stuttgarter Neuronale Netze Simulator)
- ▶ **RapidMiner**
- ▶ IBM SPSS Modeler
- ▶ R
- ▶ ...

In vielen Data-Warehouse- und Datenbank-Systemen sind Data-Mining-Komponenten integriert.

Grundlagen

Definition

Mit *Daten* bezeichnet man eine Ansammlung von Zeichen mit der dazugehörigen Syntax.

- ▶ Plural des Wortes Datum
- ▶ Interpretation als Informationseinheit
- ▶ Man unterscheidet
 - Unstrukturierte Daten
 - Semistrukturierte Daten
 - Strukturierte Daten

Aus Daten entstehen Informationen dadurch, dass diese Daten eine Bedeutung bekommen.

Definition

Eine *Information* ist ein Datum, welches mit einer Bedeutung gekoppelt ist.

- ▶ zweckbestimmte Interpretation von Daten
- ▶ Daten werden erst dann zur Information, wenn sie im Kontext betrachtet werden und eine Bedeutung erhalten

Definition

Eine Information in Verbindung mit der Fähigkeit, diese zu benutzen, wird als *Wissen* bezeichnet

- ▶ eine Information wird zu Wissen, wenn man mit ihr etwas anzufangen weiß

Definition

Data Mining (Datenschürfen) ist die Extraktion von Wissen aus Daten.

- ▶ Extraktion von implizitem Wissen
- ▶ dieses sollte nicht trivial, sondern nützlich sein
- ▶ weitgehend automatischer Ablauf (Vorbereitung der Daten aber oft manuell)
- ▶ Teilgebiete
 - Text Mining (unstrukturierte Daten)
 - Web Mining (semistrukturierte Daten)
 - Data Mining im engeren Sinn (strukturierte Daten)

Datensatz, Instanz, Objekt oder Muster

- ▶ im Data Mining weitgehend synonym
- ▶ Damit ist stets ein Datensatz in einer Datenbanktabelle gemeint, der ein Objekt oder eine Instanz anhand einer Menge von Merkmalswerten charakterisiert

Wetter-Beispiel

(sunny, hot, high, false, no)

- = math. Quintupel
- ▶ Auffassung als Objekt
- ▶ dieses Objekt ist eine Instanz der Klasse aller Wetter-Situationen

- ▶ Metrische
 - berechenbar (z.B. Mittelwerte)
 - miteinander vergleichbar
 - sortierbar (Ordnungsrelation, z.B. „<“ oder „>“)
- ▶ Ordinale
 - nicht berechenbar, aber sortierbar
- ▶ Nominale
 - nicht berechenbar
 - unterliegen keinerlei Rangfolge, d.h. sie lassen sich nicht sortieren

Metrische Datentypen können weiter unterteilt werden:

- ▶ Diskrete
 - endliche Anzahl von Werten, z.B. Schulnoten
- ▶ Kontinuierliche
 - beliebiger Zahlenwert innerhalb des Definitionsbereiches
 - im Computer auch diskret, beliebige Genauigkeit angenommen

Weiter wird nach der zugrundeliegenden Skala unterschieden:

- ▶ Intervallbasierte
 - Nullpunkt willkürlich
 - eingeschränkt berechenbar (20°C nicht viermal wärmer als 5°C)
- ▶ Verhältnisbasierte
 - natürlicher Nullpunkt
 - berechenbar, vergleichbar (60km doppelt so weit wie 30km)
 - abhängig von Maßeinheit (60km ungleich 60cm)
- ▶ Absolutskalenbasierte
 - wie Verhältnisbasierte, aber unabhängig von Maßeinheit

Achtung

Numerische Werte sind nicht dringend metrisch (z.B. Kundennummer, PLZ).

Beispiele

- ▶ metrisch, Intervallskala
 - Jahreszahlen
 - Temperatur in Celsius, Fahrenheit
- ▶ metrisch, Verhältnisskala
 - Entfernung
 - Körpergröße
- ▶ metrisch, Absolutskala
 - Lebensjahre
 - Zahl der Kinder

Beispiele

► ordinal

- Schulnoten {1,2,3,4,5,6}
- Meinung {stimme zu, stimme teilweise zu, stimme nicht zu}

► nominal

- Geschlecht {weiblich, männlich}
- Haarfarbe {blond, dunkelblond, braun, rot, schwarz, grau}

- ▶ im Data Mining werden Datensätze oft miteinander verglichen
- ▶ Ähnlichkeit zweier Datensätze muss quantifizierbar sein
- ▶ meist über sogenannte Abstandsmaße realisiert, die die „Unähnlichkeit“ der Datensätze quantifizieren

Definition

Der *Abstand* zweier Datensätze v und w ist durch das *Abstandsmaß*

$$\text{dist}(v, w)$$

definiert.

Definition

Die *Ähnlichkeit* zweier Datensätze v und w ist in Abhängigkeit von einem Abstandsmaß dist definiert.

$$\text{simil}(v, w) = f(\text{dist}(v, w))$$

- ▶ Je größer die Distanz zwischen den Datensätzen, desto geringer ist die Ähnlichkeit

Eine Abstandsfunction (auch *Distanzfunktion* genannt) muss folgende Eigenschaften erfüllen:

1. $\text{dist}(x, y) \geq 0$
2. $\text{dist}(x, x) = 0$ bzw. $\text{dist}(x, y) = 0$ gdw. $x = y$
3. $\text{dist}(x, y) = \text{dist}(y, x)$ (kommutativ)
4. $\text{dist}(x, y) \leq \text{dist}(x, z) + \text{dist}(z, y)$ (transitiv)

Typische Distanzfunktionen:

Hamming-Distanz $\text{dist}_H(v, w) = \text{count}_i(v_i \neq w_i)$

Euklidische Distanz $\text{dist}_E(v, w) = \sqrt{\sum_i (v_i - w_i)^2}$

Manhattan-Distanz $\text{dist}_{Man}(v, w) = \sum_i |v_i - w_i|$

Maximum-Distanz $\text{dist}_{Max}(v, w) = \max(|v_i - w_i|)$

RapidMiner

- ▶ Umgebung für maschinelles Lernen und Data-Mining
- ▶ Programmiersprache: Java, damit weitgehend plattformunabhängig
- ▶ grafische Benutzeroberfläche
- ▶ mehr als 1500 Algorithmen des Maschinellen Lernens, über 30 interaktive Visualisierungen
- ▶ Austauschbarkeit von Datenvorverarbeitungsschritten, Lernverfahren und Evaluierungsverfahren
 - Durchführen von Vergleichen unterschiedlicher Verfahren
 - Kombination/Verschachteln von Verfahren
- ▶ RapidMiner Studio *Open Core* unter AGPL-Lizenz
- ▶ Free- (limitiert), Professional-, Enterprise-, Educational-Lizenzen

(*Visual Workflow for Predictive Analytics* | RapidMiner© Studio 2019; RapidMiner Pricing 2019)

File Edit Process View Connections Settings Extensions Help

Views: Design Results Turbo Prep Auto Model

Find data, operators, etc. All Studio

Repository

- Import Data
- URI_Meta
- 07_Clustering
- 08_Other
- Templates
 - Churn Modeling
 - Credit Risk Modeling
 - Counterparty Risk Data (v1)
 - Credit Risk Modeling (v1)
 - Direct Marketing
 - Geographic Distances
 - LIR Chart
 - Market Basket Analysis

Operators

Search for Operators

- Data Access (53)
- Blending (79)
- Cleansing (26)
- Modeling (156)
 - Predictive (61)
 - Lazy (2)
 - Bayesian (2)
 - Trees (9)
 - Rules (5)
 - Neural Nets (4)
 - Ensembles (8)

Get more operators from the Marketplace

Process

Process

100%

Retrieve

Set Role

Filter Examples

Select Attributes

Log to Data

Optimize Parameter...

Apply Model

Parameters

Retrieve

repository entry: repository Risk Data

Help

Retrieve

RapidMiner Studio Core

Tags: Load, Import, Read, Datasets, Examples, Example Set, Table, Repository, Data Access

Synopsis

This Operator can access stored information in the Repository and load them into the Process.

[Jump to Tutorial Process](#)

Leverage the Wisdom of Crowds to get operator recommendations based on your process design!

Activate Wisdom of Crowds

Grafische Benutzeroberfläche

- ▶ modularer Aufbau mit Knoten und einheitlichen Schnittstellen zwischen den Knoten die eine Erweiterung der Software ermöglichen
- ▶ Repository zur einheitlichen Speicherung von Datensets und Modellen
- ▶ Nachteil: geringere Performance
- ▶ Dashboard mit Operatoren, Repository, Prozess, Parametern, Daten-Editor und Hilfe-Seiten

Process

- ▶ ein Rapid-Miner-Prozess besteht aus einer Operatoren-Kette
- ▶ jeder Operator hat Parameter, die eingestellt werden können
- ▶ zwischen den Operatoren können sogenannte *IO Objects* (Input-Output-Objekte) weitergereicht werden
- ▶ die IO Objects können im internen *Repository* gespeichert werden

Process



ExampleSet

- ▶ *ExampleSet* ist das Objekt, in dem Daten beschrieben und gespeichert werden
- ▶ ein *ExampleSet* ist eine Menge von *Examples* (Datensatz, math. Tupel)
- ▶ ein *Example* besteht aus einer Menge von Attributen
- ▶ jedes *Example* eines *ExampleSets* hat die gleichen Attribute
- ▶ ein *ExampleSet* ist im Prinzip eine Tabelle, die Zeilen sind die *Examples* und die Spalten die Attribute

File Edit Process View Connections Settings Extensions Help

Views: Design Results Turbo Prep Auto Model

Find data, operators, etc. All Studio

PerformanceVector (Performance) ExampleSet (Log to Data) Log Optimize Parameters (Grid)

Result History ExampleSet (Apply Model)

Open in Turbo Prep Auto Model Filter (34 / 34 examples): all

Row No.	prediction(D...	confidence...	confidence...	Long Term F...	Working Ca...	Debt Cash F...	Liability to E...	Net Debt to ...	Debt to Capl...	Long Term ...	Long Term
1	No	0.975	0.025	0.836	51.987	1.787	0.018	0.010	0.033	5.137	0.304
2	No	0.982	0.018	2.187	111.569	1.654	0.124	0.386	0.266	1.113	0.416
3	No	0.875	0.125	4.828	40.670	1.700	0.094	0.193	0.005	0.063	0.503
4	No	0.859	0.141	0.853	51.729	1.123	0.258	0.045	0.272	3.509	0.920
5	No	0.947	0.053	3.210	77.686	0.147	0.206	0.098	0.076	1.296	0.390
6	No	0.792	0.208	8.440	20.271	0.406	0.100	0.387	0.067	2.464	0.843
7	No	0.826	0.174	1.110	20.350	1.015	0.210	0.138	0.145	3.836	0.302
8	No	0.949	0.051	4.867	88.386	0.105	0.026	0.089	0.009	4.980	0.379
9	No	0.981	0.019	0.178	72.130	0.629	0.044	0.408	0.160	5.110	0.737
10	No	0.868	0.132	3.296	17.272	0.219	0.119	0.185	0.199	4.438	0.155
11	No	0.705	0.295	1.610	62.897	0.786	0.328	0.035	0.145	1.079	0.515
12	No	0.985	0.015	0.759	83.688	2.593	0.031	0.352	0.023	0.556	0.444
13	No	0.932	0.068	6.729	14.730	3.637	0.152	0.246	0.027	4.388	0.484
14	No	0.988	0.014	1.153	77.878	0.559	0.094	0.451	0.027	5.259	1.446
15	No	0.750	0.250	2.638	32.805	0.167	0.067	0.018	0.017	0.098	1.183
16	No	0.945	0.055	7.321	56.092	0.712	0.097	0.105	0.291	0.620	0.776
17	No	0.854	0.146	8.275	51.919	0.786	0.056	0.060	0.198	2.080	0.946
18	Yes	0.465	0.535	1.319	6.124	1.238	0.179	0.096	0.035	0.785	0.215

ExampleSet (34 examples, 3 special attributes, 19 regular attributes)

Repository

- Import Data
- ▼ V_Learner
 - 02_Preprocessing
 - 03_Validation
 - 04_Attributes
 - 05_Visualisation
 - 06_Meta
 - 07_Clustering
 - 08_Other
- ▼ Templates
 - Churn Modeling
 - Credit Risk Modeling
 - Counterparty Risk Data (x1)
 - Credit Risk Modeling (x1)
 - Direct Marketing
 - Geographic Distances
 - Lift Chart
 - Market Basket Analysis
 - Medical Fraud Detection
 - Operationalization
 - Outlier Detection
 - Predictive Maintenance
 - Price Risk Clustering
 - Web Analytics
- ▼ Time Series
- ▼ Tutorials
- ▼ Community Samples (connected)
 - DB (Legacy)
 - Local Repository (chrisoph)

1. Reguläre Attribute

- aus diesen Daten wird das Modell erlernt
- die Speziellen Attribute werden nicht für Lernoperationen verwendet

2. Spezielle Attribute

- **ID** eindeutiger Identifier zum Identifizieren eines Examples; wird vor allem in der Datenvorverarbeitung zum Verbinden von Datensätzen (Join) benötigt
- **Label** das Zielattribut welches vorhergesagt werden soll; für dieses Attribut wird das Modell trainiert
- **Prediction** die aus dem Modell erstellte Vorhersage des Zielattributs
- **Cluster** die Zugehörigkeit eines Examples zu einem Cluster
- **Weight** das Gewicht eines Examples
- **Batch** die Zugehörigkeit zu einem *Example Batch*

Value type	Name	Use
Nominal	nominal	Categorical non-numerical values, usually used for finite quantities of different characteristics
Numerical values	numeric	For numerical values in general
Integers	integer	Whole numbers, positive and negative Real numbers real Real numbers, positive and negative
Text	text	Random free text without structure
2-value nominal	binominal	Special case of nominal, where only two different values are permitted
multi-value nominal	polynominal	Special case of nominal, where more than two different values are permitted
Date Time	date_time	Date as well as time
Date	date	Only date
Time	time	Only time

(RapidMiner Studio Manual 2014)

Cleve, Jürgen und Uwe Lämmel (2014). *Data mining*. Studium. München: De Gruyter Oldenbourg. 306 S. ISBN: 978-3-486-72034-1 978-3-486-71391-6.

Ittner, Prof. Dr.-Ing. Andreas (2019). „Data Mining“. Vorlesung, Wintersemester 2019/20 (Hochschule Mittweida).

RapidMiner Pricing (2019). URL: <https://rapidminer.com/pricing/> (besucht am 12.09.2019).

RapidMiner Studio Manual (2014). URL: <https://docs.rapidminer.com/downloads/RapidMiner-v6-user-manual.pdf> (besucht am 12.09.2019).

Visual Workflow for Predictive Analytics | RapidMiner® Studio (2019). URL: <https://rapidminer.com/products/studio/> (besucht am 12.09.2019).