

**Christian Gottermeier**

## Data Mining

Modellierung und Durchführung ausgewählter Fallstudien  
mit dem SAS Enterprise Miner

**Diplomarbeit**

## **Bibliografische Information der Deutschen Nationalbibliothek:**

Bibliografische Information der Deutschen Nationalbibliothek: Die Deutsche Bibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de/> abrufbar.

Dieses Werk sowie alle darin enthaltenen einzelnen Beiträge und Abbildungen sind urheberrechtlich geschützt. Jede Verwertung, die nicht ausdrücklich vom Urheberrechtsschutz zugelassen ist, bedarf der vorherigen Zustimmung des Verlanges. Das gilt insbesondere für Vervielfältigungen, Bearbeitungen, Übersetzungen, Mikroverfilmungen, Auswertungen durch Datenbanken und für die Einspeicherung und Verarbeitung in elektronische Systeme. Alle Rechte, auch die des auszugsweisen Nachdrucks, der fotomechanischen Wiedergabe (einschließlich Mikrokopie) sowie der Auswertung durch Datenbanken oder ähnliche Einrichtungen, vorbehalten.

Copyright © 2003 Diplomica Verlag GmbH  
ISBN: 9783832472177

**Christian Gottermeier**

## **Data Mining**

**Modellierung und Durchführung ausgewählter Fallstudien mit dem SAS Enterprise Miner**



---

Christian Gottermeier

# Data Mining

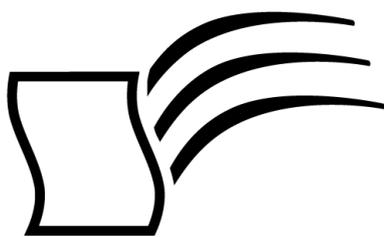
*Modellierung und Durchführung ausgewählter  
Fallstudien mit dem SAS Enterprise Miner*

**Diplomarbeit**

**Ruprecht-Karls-Universität Heidelberg**

**Fachbereich Wirtschafts- und Sozialwissenschaften**

**Abgabe Januar 2003**



***Diplom.de***

Diplomica GmbH ———  
Hermannstal 119k ———  
22119 Hamburg ———

Fon: 040 / 655 99 20 ———  
Fax: 040 / 655 99 222 ———

agentur@diplom.de ———  
www.diplom.de ———

ID 7217

Gottermeier, Christian: Data Mining - Modellierung und Durchführung ausgewählter Fallstudien mit dem SAS Enterprise Miner

Hamburg: Diplomica GmbH, 2003

Zugl.: Ruprecht-Karls-Universität Heidelberg, Universität, Diplomarbeit, 2003

---

Dieses Werk ist urheberrechtlich geschützt. Die dadurch begründeten Rechte, insbesondere die der Übersetzung, des Nachdrucks, des Vortrags, der Entnahme von Abbildungen und Tabellen, der Funksendung, der Mikroverfilmung oder der Vervielfältigung auf anderen Wegen und der Speicherung in Datenverarbeitungsanlagen, bleiben, auch bei nur auszugsweiser Verwertung, vorbehalten. Eine Vervielfältigung dieses Werkes oder von Teilen dieses Werkes ist auch im Einzelfall nur in den Grenzen der gesetzlichen Bestimmungen des Urheberrechtsgesetzes der Bundesrepublik Deutschland in der jeweils geltenden Fassung zulässig. Sie ist grundsätzlich vergütungspflichtig. Zuwiderhandlungen unterliegen den Strafbestimmungen des Urheberrechtes.

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Die Informationen in diesem Werk wurden mit Sorgfalt erarbeitet. Dennoch können Fehler nicht vollständig ausgeschlossen werden, und die Diplomarbeiten Agentur, die Autoren oder Übersetzer übernehmen keine juristische Verantwortung oder irgendeine Haftung für evtl. verbliebene fehlerhafte Angaben und deren Folgen.

Diplomica GmbH

<http://www.diplom.de>, Hamburg 2003

Printed in Germany

# INHALTSVERZEICHNIS

<b>1. Einführung</b> .....	<b>1</b>
<b>2. Data Mining</b> .....	<b>3</b>
2.1 Definitionen und Erklärungen.....	3
2.2 Einführung in die wichtigsten Verfahren.....	5
2.2.1 Data Mining als interdisziplinäre Wissenschaft.....	5
2.2.1.1 Multivariate Analysemethoden .....	6
2.2.1.1.1 Regressionsanalyse.....	6
2.2.1.1.2 Clusteranalyse.....	6
2.2.1.2 Künstliche Intelligenz (KI) und maschinelles Lernen.....	6
2.2.1.2.1 Entscheidungsbaumverfahren.....	7
2.2.1.2.2 Künstliche neuronale Netze (KNN).....	7
2.2.1.2.3 Selbstorganisierende Karten (SOM) / Kohonen-Netze .....	7
2.2.1.3 Assoziations- und Sequenzanalyse.....	8
2.2.2 Alternative Einordnungsmöglichkeiten.....	8
2.2.2.1 Überwachtes vs. unüberwachtes Lernen .....	8
2.2.2.2 Parametrische vs. nichtparametrische Verfahren .....	9
2.3 Architekturüberlegungen .....	10
2.3.1 Data Warehouse (DWH) und Data Marts .....	10
2.3.2 Integration mit Data Mining.....	11
2.3.3 OLAP .....	11
2.3.4 OLAP und Data Mining .....	12
2.4 Einsatzgebiete .....	13
2.4.1 Customer Relationship Management (CRM).....	14
2.4.2 Text Mining.....	14
2.4.2 Web Mining.....	15
<b>3. Pre-Processing</b> .....	<b>17</b>
3.1 Partitionierung der Daten .....	17
3.1.1 Trainings-, Validierungs- und Testdaten.....	17
3.1.2 Seltene Zielereignisse.....	17
3.1.3 Massiv große oder beschränkt kleine Datensätze .....	18
3.1.3.1 Cross Validation .....	18

3.1.3.2 Sampling.....	18
3.2 Variablenselektion oder das Problem hoher Dimensionalität.....	19
3.3 Fehlende Werte .....	19
3.4 Transformationsprozesse .....	20
<b>4. Die Methoden.....</b>	<b>21</b>
Vorbemerkungen: Grundproblematik Generalisierbarkeit .....	21
4.1 Regressionsanalyse .....	22
4.1.1 Einführung in die lineare Regression .....	22
4.1.1.1 Lineare Einfachregression .....	22
4.1.1.1.1 Schätzung der Koeffizienten.....	22
4.1.1.2 Lineare Mehrfachregression.....	23
4.1.1.3 Annahmen des linearen Regressionsmodells .....	24
4.1.2 Logistische Regression.....	24
4.1.2.1 Einführung in die logistische Regression.....	24
4.1.2.2 Der Rechenansatz der logistischen Regression .....	25
4.1.2.3 Schätzung der Koeffizienten .....	25
4.1.3 Variablenauswahlverfahren.....	26
4.2 Clusteranalyse .....	26
4.2.1 Einführung in die Clusteranalyse .....	26
4.2.2 K-Means-Verfahren .....	27
4.2.3 Der K-Means-Algorithmus .....	27
4.3 Entscheidungsbaumverfahren .....	28
4.3.1 Aufbau eines Entscheidungsbaums.....	28
4.3.1.1 Algorithmen.....	29
4.3.1.2 Auswahlmaße .....	29
4.3.1.2.1 Informationsgewinn und Entropie.....	30
4.3.1.2.2 Gini-Index.....	30
4.3.1.2.3 $\chi^2$ -Maß.....	31
4.3.1.3 Stoppkriterien .....	32
4.3.2 Pruning .....	32
4.3.3 Surrogat-Splits für das Einfügen fehlender Werte .....	33
4.3.4 Wälder: Bagging und Boosting.....	33
4.4 Künstliche neuronale Netze .....	35
4.4.1 Einführung in die künstlichen neuronalen Netze .....	35

4.4.2 Netzwerkarchitektur .....	35
4.4.2.1 Multilayer Perceptron (MLP).....	35
4.4.2.2 Radiale-Basisfunktionen-Netze (RBF-Netze).....	39
4.4.3 Lernregel .....	42
4.4.3.1 Gradientenabstiegsverfahren .....	43
4.4.3.1.1 Probleme bei Gradientenverfahren .....	43
4.4.3.2 Backpropagation.....	44
4.4.3.3 Konjugierter Gradientenabstieg.....	46
4.4.3.4 Newton-Verfahren.....	46
4.4.3.5 Levenberg-Marquard.....	46
4.4.4 Regulierbarkeit .....	46
4.4.4.1 Early Stopping .....	46
4.4.4.2 Weight Decay .....	47
4.4.5 Selbstorganisierende Karten (SOM) / Kohonen-Netze.....	47
4.4.5.1 Prinzipien der selbstorganisierenden Karten .....	47
4.4.5.2 Lernverfahren der selbstorganisierenden Karten.....	48
4.5 Assoziations- und Sequenzanalyse .....	49
4.5.1 Einführung in die Assoziationsregeln .....	49
4.5.1.1 Support .....	50
4.5.1.2 Konfidenz .....	50
4.5.1.3 Lift .....	50
4.5.2 Sequenzmuster .....	51
<b>5. Modellbewertung .....</b>	<b>52</b>
5.1 Bewertung der Klassifizierungsleistung .....	52
5.2 Draw Lift Charts .....	53
<b>6. Fallstudien.....</b>	<b>55</b>
6.1 Fallstudie A: Optimierung einer Mailing-Aktion .....	55
6.2 Fallstudie B: Funktionsweise der KNN .....	58
6.2.1 Auswahl der Netzwerkarchitektur bei NRBF-Netzen.....	58
6.2.2 Auswahl des Lernverfahrens bei MLP-Netzwerkarchitekturen.....	60
6.2.3 Early Stopping.....	63
6.3 Fallstudie C: Entscheidungsbaumverfahren.....	64
6.3.1 Bestimmung des Auswahlmaßes.....	64
6.3.2 Bagging .....	67

<b>7. Zusammenfassung.....</b>	<b>68</b>
<b>Anhang .....</b>	<b>V</b>
<b>Abbildungsverzeichnis.....</b>	<b>XLIX</b>
<b>Literaturverzeichnis.....</b>	<b>LII</b>
<b>Abkürzungsverzeichnis .....</b>	<b>LIV</b>

# 1. Einführung

Entscheidungen sind ein Akt des menschlichen Verhaltens, bei denen eine Festlegung für eine unter mehreren Möglichkeiten stattfindet. Da bei diesen Handlungen die Berufung auf Traditionen oder Autoritäten oftmals nicht möglich ist, wurde schon früh auf verschiedenste Hilfsmittel zurückgegriffen. So ließ sich Julius Cäsar von einem Würfelergewinn leiten, General Wallenstein von einem Astrologen beraten oder es wurden Prognosen mit Hilfe von Glaskugeln, Spielkarten oder dem Stand der Sterne getroffen.

Unter wirtschaftlichen Gesichtspunkten sind Entscheidungen eine rationale Wahl zwischen mehreren Möglichkeiten, wobei der Entscheidungsprozess als tragendes Element der ökonomischen Tätigkeit herausgestellt wird. Gerade in diesem Umfeld wird die Entscheidungsfindung – nun allerdings wissenschaftlich fundiert und mit weitreichenden Konsequenzen – durch folgende Verfahren unterstützt: Analysemethoden wie Benchmarking, Lebenszyklus- oder Erfahrungskurvenkonzept und Prognoseverfahren wie die Delphi-Methode oder die Szenario-Technik. Allerdings sind die meisten dieser Verfahren i.d.R. auf spezielle Problemstellungen ausgerichtet. Ganzheitliche Lösungsansätze werden seit den 60er Jahren zur Unterstützung des Managements bereitgestellt. Mit Hilfe von Informationssystemen soll die Entscheidungsfindung verbessert werden. Häufig wechselnde Schlagworte wie z.B. Management Information System (MIS) oder Decision Support System (DSS) konnten allerdings noch keine durchschlagenden Erfolge erzielen. Seit Mitte der 90er Jahre wurden mit neuen konzeptionellen Ansätzen, die meist unter dem Oberbegriff „Business Intelligence“ zusammengefasst werden, erfolgsversprechende Lösungen zum Aufbau entscheidungsorientierter Informationssysteme (EIS) etabliert. EIS setzen sich dabei aus Werkzeugen zur Selektion und Speicherung entscheidungsrelevanter Informationen (Data Warehouse) sowie zur entscheidungsunterstützenden Modellierung (OLAP-Tools) zusammen. Eine konsequente Umsetzung des Data Warehouse Gedanken führt zu immensen Datensammlungen, die, um die Archivierung nicht zum Selbstzweck werden zu lassen, dann auch ausgewertet werden sollen. An dieser Stelle setzt Data Mining an.

In Kapitel 2 werden die Grundzüge des Data Mining dargestellt, eine Verbindung zu Data Warehouse und OLAP gezogen und die Einsatzgebiete skizziert, in denen sich Data Mining durchgesetzt hat. In Kapitel 3 wird der erste wichtige Schritt, der vor der eigentlichen Modellierung stattfinden sollte, das Pre-Processing, erläutert. Die Modelle und die damit verbundenen Methodiken der Data Mining-Verfahren werden in Kapitel 4 vorgestellt. Stets wird eine Verbindung zum SAS<sup>®</sup> Enterprise Miner<sup>™</sup> gesucht und so eine

Anpassung der dort verankerten Möglichkeiten an die Theorie vorgenommen. Die Vorgehensweise der Modellbewertung und die dafür existierenden Kriterien werden in Kapitel 5 dargestellt. Die praktische Umsetzung der Data Mining-Modelle wird anhand verschiedener Fallstudien im sechsten Kapitel gezeigt. Dafür werden die von der SAS<sup>®</sup> Institute Inc. erstellten Fälle bearbeitet. Diese Daten sind stark idealisiert, d.h. sofort analysierbar und deshalb sehr gut geeignet, um die einzelnen Schritte Pre-Processing, Modellierung der einzelnen Verfahren und Modellbewertung durchzuführen.

## 2. Data Mining

### 2.1 Definitionen und Erklärungen

Durch die Entwicklung der Informationstechnologie ist die Datenerhebung, -speicherung und -verwaltung stark ausgeprägt. Datenbanken der Größenordnung Gigabyte oder Terabyte sind weit verbreitet. Schätzungen zufolge verdoppeln sich die weltweit vorhandenen Informationen alle 20 Monate, in Datenbanken ist die Rate noch größer. Allerdings konnten die Möglichkeiten der manuellen Analysetechnik mit dieser Entwicklung nicht Schritt halten. Die Folge ist, dass die Datenmengen nicht zu aussagefähigen Informationen verdichtet werden. Der verschärfte Wettbewerbsdruck zwingt aber zur Nutzung aller Informationsquellen.

Der Begriff Data Mining bezeichnet eine relativ neue Forschungs- und Anwendungsrichtung, obwohl die Bestandteile schon lange existieren. Verfahren der klassischen statistischen Datenanalyse, Anwendungen aus der künstlichen Intelligenz (KI), der Mustererkennung und des maschinellen Lernens wurden in das sog. Knowledge Discovery in Databases (KDD) integriert. Wie der Name schon assoziiert, besteht auch eine starke Verbindung zur Datenbanktechnologie. KDD und Data Mining werden in dieser Arbeit und auch meistens in der Literatur simultan verwendet.

Die Aufgabe des Data Minings ist die Entwicklung, Implementierung und Ausführung von Datenanalysemethoden. Ein Fokus wird dabei auf sehr große Datensätze mit komplexen Strukturen gelegt. Die vier Hauptaufgaben sind:

- a.) Vorhersage- und Klassifikationsmodelle,
- b.) Segmentierungen oder Clusterung,
- c.) Dimensionsreduktion und
- d.) Assoziationsanalysen.

Vorhersage und Klassifikation unterscheiden sich hauptsächlich in den Skalierungsmaßen. Ist die abhängige Variable intervallskaliert, spricht man von Vorhersagemodellen. Klassifikationsmodelle dagegen besitzen einen binären, ordinalen oder nominalen Regressand. Man spricht von Klassifikation, weil die Klassenwahrscheinlichkeiten angegeben werden sollen. Vertreter dieser Modelle sind u.a. die Regressionsanalyse, das Entscheidungsbaumverfahren oder die künstlichen neuronalen Netze.

Die Besonderheit der Clusterung ist das Fehlen einer abhängigen Variablen. Als Modelle, die eine Klassenunterscheidung aufgrund der Homogenität bzw. Heterogenität der unabhängigen Variablen vornehmen, können beispielsweise das K-Means-Verfahren oder sog. selbstorganisierende Karten, die Kohonen-Netze, angeführt werden.

Die Dimensionsreduktion ist in komplexen Systemen unerlässlich, um eine Visualisierung der Problemstellung zu ermöglichen. Auch wenn eine Reduktion auf zwei bis drei Dimensionen, die eine grafische Anschauung ermöglicht, nicht gelingt, ist eine Variablenselektion das Instrument, um Modelle zu vereinfachen und damit die Fähigkeit zur Modellentwicklung zu erhöhen (vgl. auch Abschnitt 3.2).

Eine Assoziationsanalyse soll Beziehungen zwischen Variablen aufdecken. Als Verfahren werden Assoziationsregeln, Sequenzmuster oder Link-Analysen verwendet.

In den Vorhersage- und Klassifikationsmodellen liegt der Schwerpunkt dieser Arbeit, sowohl bei der Analyse der Modellierung als auch bei den Fallstudien, bei denen ausschließlich diese Modelle untersucht werden<sup>1</sup>.

An dieser Stelle sollte auch auf die Schwierigkeiten eingegangen werden, eine allgemeingültige Definition oder eine einheitliche Terminologie zu finden. Entgegen der Gleichstellung der Begriffe KDD und Data Mining (siehe oben) existiert auch die Sichtweise, KDD als den Gesamtprozess der Analyse anzusehen und Data Mining als Synonym für die einzelnen eingesetzten Methoden zu bezeichnen. Fayyad<sup>2</sup>, der das Thema Data Mining seit den Anfängen zu Beginn der 90er Jahre prägte, stellte folgenden Zusammenhang zwischen beiden Begriffen her:

*“KDD is the non-trivial process of identifying valid, novel, potential useful and ultimately understandable pattern in data.”*

*“Data Mining is a simple step in the KDD process that under acceptable computational efficiency limitations enumerates structures (patterns or models) over the data.”*

Data Mining mit seinen Methoden aus den verschiedensten Gebieten ist eine interdisziplinäre Wissenschaft. Schon deshalb wird es je nach Schwerpunkt und wissenschaftlicher Herkunft verschiedene Sichtweisen, Ansätze und vor allem Notationen geben. Darüber hinaus finden sich nicht nur in der Literatur, sondern auch in der praktischen Umsetzung der verschiedenen Anbieter von Data Mining-Software unterschiedliche Ansätze, die Teilgebiete ausdrücklich ausgrenzen oder neue Methoden einführen. Der Ansatz dieser Arbeit besteht nicht darin, alle Methoden und Sichtweisen darzustellen, sondern eine Einführung in die Problematik zu bieten, aufgrund derer die nachher folgenden Fallstudien durchgeführt werden.

Die deutsche Data Mining-Übersetzung Datenmustererkennung nennt die Identifizierung von Mustern als Ziel des Data Mining-Prozesses, während der Begriff KDD eine

---

<sup>1</sup> Der Ablauf bei der Erstellung von Vorhersage- und Klassifikationsmodellen kann im Anhang in Kapitel A.1 nachvollzogen werden.

<sup>2</sup> Küppers (1999), S. 23.