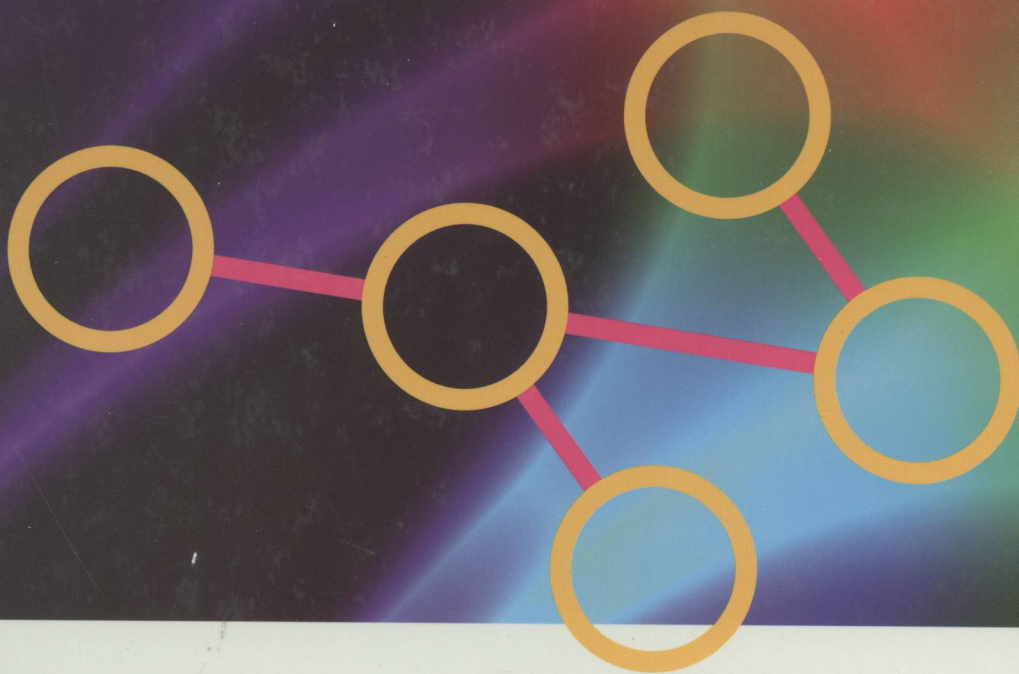


WILEY SERIES IN COMPUTATIONAL STATISTICS



Christian Borgelt,  
Matthias Steinbrecher and Rudolf Kruse

# GRAPHICAL MODELS

REPRESENTATIONS FOR LEARNING,  
REASONING AND DATA MINING,  
2ND EDITION

 WILEY

0212  
B732  
E-2

# Graphical Models

## Representations for Learning, Reasoning and Data Mining

### Second Edition

---

Christian Borgelt

*European Centre for Soft Computing, Spain*

Matthias Steinbrecher & Rudolf Kruse

*Otto-von-Guericke University Magdeburg, Germany*



E2009003794

 **WILEY**

A John Wiley and Sons, Ltd., Publication

This edition first published 2009  
© 2009, John Wiley & Sons Ltd

*Registered office*

John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, West Sussex, PO19 8SQ, United Kingdom

For details of our global editorial offices, for customer services and for information about how to apply for permission to reuse the copyright material in this book please see our website at [www.wiley.com](http://www.wiley.com).

The right of the author to be identified as the author of this work has been asserted in accordance with the Copyright, Designs and Patents Act 1988.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as permitted by the UK Copyright, Designs and Patents Act 1988, without the prior permission of the publisher.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

Designations used by companies to distinguish their products are often claimed as trademarks. All brand names and product names used in this book are trade names, service marks, trademarks or registered trademarks of their respective owners. The publisher is not associated with any product or vendor mentioned in this book. This publication is designed to provide accurate and authoritative information in regard to the subject matter covered. It is sold on the understanding that the publisher is not engaged in rendering professional services. If professional advice or other expert assistance is required, the services of a competent professional should be sought.



***Library of Congress Cataloging-in-Publication Data***  
Record on file

A catalogue record for this book is available from the British Library.

ISBN 978-0-470-72210-7

Typeset in 10/12 cmr10 by Laserwords Private Limited, Chennai, India  
Printed in Great Britain by TJ International Ltd, Padstow, Cornwall

# Graphical Models

---

# Wiley Series in Computational Statistics

## Consulting Editors:

Paolo Giudici  
*University of Pavia, Italy*

Geof H. Givens  
*Colorado State University, USA*

Bani K. Mallick  
*Texas A & M University, USA*

---

*Wiley Series in Computational Statistics* is comprised of practical guides and cutting edge research books on new developments in computational statistics. It features quality authors with a strong applications focus. The texts in the series provide detailed coverage of statistical concepts, methods and case studies in areas at the interface of statistics, computing, and numerics.

With sound motivation and a wealth of practical examples, the books show in concrete terms how to select and to use appropriate ranges of statistical computing techniques in particular fields of study. Readers are assumed to have a basic understanding of introductory terminology.

The series concentrates on applications of computational methods in statistics to fields of bioinformatics, genomics, epidemiology, business, engineering, finance and applied statistics.



# Preface

Although the origins of graphical models can be traced back to the beginning of the 20th century, they have become truly popular only since the mid-eighties, when several researchers started to use Bayesian networks in expert systems. But as soon as this start was made, the interest in graphical models grew rapidly and is still growing to this day. The reason is that graphical models, due to their explicit and sound treatment of (conditional) dependences and independences, proved to be clearly superior to naive approaches like certainty factors attached to if-then-rules, which had been tried earlier.

Data Mining, also called Knowledge Discovery in Databases, is another relatively young area of research, which has emerged in response to the flood of data we are faced with nowadays. It has taken up the challenge to develop techniques that can help humans discover useful patterns in their data. In industrial applications patterns found with these methods can often be exploited to improve products and processes and to increase turnover.

This book is positioned at the boundary between these two highly important research areas, because it focuses on learning graphical models from data, thus exploiting the recognized advantages of graphical models for learning and data analysis. Its special feature is that it is not restricted to probabilistic models like Bayesian and Markov networks. It also explores relational graphical models, which provide excellent didactical means to explain the ideas underlying graphical models. In addition, possibilistic graphical models are studied, which are worth considering if the data to analyze contains imprecise information in the form of sets of alternatives instead of unique values.

Looking back, this book has become longer than originally intended. However, although it is true that, as C.F. von Weizsäcker remarked in a lecture, anything ultimately understood can be said briefly, it is also evident that anything said too briefly is likely to be incomprehensible to anyone who has not yet understood completely. Since our main aim was comprehensibility, we hope that a reader is remunerated for the length of this book by an exposition that is clear and self-contained and thus easy to read.

Christian Borgelt, Matthias Steinbrecher, Rudolf Kruse  
Oviedo and Magdeburg, March 2009

# Contents

<b>Preface</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Data and Knowledge . . . . .	2
1.2 Knowledge Discovery and Data Mining . . . . .	5
1.2.1 The KDD Process . . . . .	6
1.2.2 Data Mining Tasks . . . . .	7
1.2.3 Data Mining Methods . . . . .	8
1.3 Graphical Models . . . . .	10
1.4 Outline of this Book . . . . .	12
<b>2 Imprecision and Uncertainty</b>	<b>15</b>
2.1 Modeling Inferences . . . . .	15
2.2 Imprecision and Relational Algebra . . . . .	17
2.3 Uncertainty and Probability Theory . . . . .	19
2.4 Possibility Theory and the Context Model . . . . .	21
2.4.1 Experiments with Dice . . . . .	22
2.4.2 The Context Model . . . . .	27
2.4.3 The Insufficient Reason Principle . . . . .	30
2.4.4 Overlapping Contexts . . . . .	31
2.4.5 Mathematical Formalization . . . . .	35
2.4.6 Normalization and Consistency . . . . .	37
2.4.7 Possibility Measures . . . . .	39
2.4.8 Mass Assignment Theory . . . . .	43
2.4.9 Degrees of Possibility for Decision Making . . . . .	45
2.4.10 Conditional Degrees of Possibility . . . . .	47
2.4.11 Imprecision and Uncertainty . . . . .	48
2.4.12 Open Problems . . . . .	48
<b>3 Decomposition</b>	<b>53</b>
3.1 Decomposition and Reasoning . . . . .	54
3.2 Relational Decomposition . . . . .	55

3.2.1	A Simple Example . . . . .	55
3.2.2	Reasoning in the Simple Example . . . . .	57
3.2.3	Decomposability of Relations . . . . .	61
3.2.4	Tuple-Based Formalization . . . . .	63
3.2.5	Possibility-Based Formalization . . . . .	66
3.2.6	Conditional Possibility and Independence . . . . .	70
3.3	Probabilistic Decomposition . . . . .	74
3.3.1	A Simple Example . . . . .	74
3.3.2	Reasoning in the Simple Example . . . . .	76
3.3.3	Factorization of Probability Distributions . . . . .	77
3.3.4	Conditional Probability and Independence . . . . .	78
3.4	Possibilistic Decomposition . . . . .	82
3.4.1	Transfer from Relational Decomposition . . . . .	83
3.4.2	A Simple Example . . . . .	83
3.4.3	Reasoning in the Simple Example . . . . .	84
3.4.4	Conditional Degrees of Possibility and Independence . . . . .	85
3.5	Possibility versus Probability . . . . .	87
<b>4</b>	<b>Graphical Representation</b>	<b>93</b>
4.1	Conditional Independence Graphs . . . . .	94
4.1.1	Axioms of Conditional Independence . . . . .	94
4.1.2	Graph Terminology . . . . .	97
4.1.3	Separation in Graphs . . . . .	100
4.1.4	Dependence and Independence Maps . . . . .	102
4.1.5	Markov Properties of Graphs . . . . .	106
4.1.6	Markov Equivalence of Graphs . . . . .	111
4.1.7	Graphs and Decompositions . . . . .	114
4.1.8	Markov Networks and Bayesian Networks . . . . .	120
4.2	Evidence Propagation in Graphs . . . . .	121
4.2.1	Propagation in Undirected Trees . . . . .	122
4.2.2	Join Tree Propagation . . . . .	128
4.2.3	Other Evidence Propagation Methods . . . . .	136
<b>5</b>	<b>Computing Projections</b>	<b>139</b>
5.1	Databases of Sample Cases . . . . .	140
5.2	Relational and Sum Projections . . . . .	141
5.3	Expectation Maximization . . . . .	143
5.4	Maximum Projections . . . . .	148
5.4.1	A Simple Example . . . . .	149
5.4.2	Computation via the Support . . . . .	151
5.4.3	Computation via the Closure . . . . .	152
5.4.4	Experimental Evaluation . . . . .	155
5.4.5	Limitations . . . . .	156



<b>6</b>	<b>Naive Classifiers</b>	<b>157</b>
6.1	Naive Bayes Classifiers . . . . .	157
6.1.1	The Basic Formula . . . . .	157
6.1.2	Relation to Bayesian Networks . . . . .	160
6.1.3	A Simple Example . . . . .	161
6.2	A Naive Possibilistic Classifier . . . . .	162
6.3	Classifier Simplification . . . . .	164
6.4	Experimental Evaluation . . . . .	164
<b>7</b>	<b>Learning Global Structure</b>	<b>167</b>
7.1	Principles of Learning Global Structure . . . . .	168
7.1.1	Learning Relational Networks . . . . .	168
7.1.2	Learning Probabilistic Networks . . . . .	177
7.1.3	Learning Possibilistic Networks . . . . .	183
7.1.4	Components of a Learning Algorithm . . . . .	192
7.2	Evaluation Measures . . . . .	193
7.2.1	General Considerations . . . . .	193
7.2.2	Notation and Presuppositions . . . . .	197
7.2.3	Relational Evaluation Measures . . . . .	199
7.2.4	Probabilistic Evaluation Measures . . . . .	201
7.2.5	Possibilistic Evaluation Measures . . . . .	228
7.3	Search Methods . . . . .	230
7.3.1	Exhaustive Graph Search . . . . .	230
7.3.2	Greedy Search . . . . .	232
7.3.3	Guided Random Graph Search . . . . .	239
7.3.4	Conditional Independence Search . . . . .	247
7.4	Experimental Evaluation . . . . .	259
7.4.1	Learning Probabilistic Networks . . . . .	259
7.4.2	Learning Possibilistic Networks . . . . .	261
<b>8</b>	<b>Learning Local Structure</b>	<b>265</b>
8.1	Local Network Structure . . . . .	265
8.2	Learning Local Structure . . . . .	267
8.3	Experimental Evaluation . . . . .	271
<b>9</b>	<b>Inductive Causation</b>	<b>273</b>
9.1	Correlation and Causation . . . . .	273
9.2	Causal and Probabilistic Structure . . . . .	274
9.3	Faithfulness and Latent Variables . . . . .	276
9.4	The Inductive Causation Algorithm . . . . .	278
9.5	Critique of the Underlying Assumptions . . . . .	279
9.6	Evaluation . . . . .	284

<b>10 Visualization</b>	<b>287</b>
10.1 Potentials . . . . .	288
10.2 Association Rules . . . . .	289
<b>11 Applications</b>	<b>295</b>
11.1 Diagnosis of Electrical Circuits . . . . .	295
11.1.1 Iterative Proportional Fitting . . . . .	296
11.1.2 Modeling Electrical Circuits . . . . .	297
11.1.3 Constructing a Graphical Model . . . . .	299
11.1.4 A Simple Diagnosis Example . . . . .	301
11.2 Application in Telecommunications . . . . .	304
11.3 Application at Volkswagen . . . . .	307
11.4 Application at DaimlerChrysler . . . . .	310
<b>A Proofs of Theorems</b>	<b>317</b>
A.1 Proof of Theorem 4.1.2 . . . . .	317
A.2 Proof of Theorem 4.1.18 . . . . .	321
A.3 Proof of Theorem 4.1.20 . . . . .	322
A.4 Proof of Theorem 4.1.26 . . . . .	327
A.5 Proof of Theorem 4.1.28 . . . . .	332
A.6 Proof of Theorem 4.1.30 . . . . .	335
A.7 Proof of Theorem 4.1.31 . . . . .	337
A.8 Proof of Theorem 5.4.8 . . . . .	338
A.9 Proof of Lemma 7.2.2 . . . . .	340
A.10 Proof of Lemma 7.2.4 . . . . .	342
A.11 Proof of Lemma 7.2.6 . . . . .	344
A.12 Proof of Theorem 7.3.1 . . . . .	345
A.13 Proof of Theorem 7.3.2 . . . . .	346
A.14 Proof of Theorem 7.3.3 . . . . .	347
A.15 Proof of Theorem 7.3.5 . . . . .	350
A.16 Proof of Theorem 7.3.7 . . . . .	351
<b>B Software Tools</b>	<b>353</b>
<b>Bibliography</b>	<b>359</b>
<b>Index</b>	<b>383</b>

# Chapter 1

## Introduction

Due to modern information technology, which produces ever more powerful computers and faster networks every year, it is possible today to collect, transfer, combine, and store huge amounts of data at very low costs. Thus an ever-increasing number of companies and scientific and governmental institutions can afford to compile huge archives of tables, documents, images, and sounds in electronic form. The thought is compelling that if you only have enough data, you can solve any problem—at least in principle.

A closer examination reveals though, that data alone, however voluminous, are not sufficient. We may say that in large databases we cannot see the wood for the trees. Although any single bit of information can be retrieved and simple aggregations can be computed (for example, the average monthly sales in the Frankfurt area), general patterns, structures, and regularities usually go undetected. However, often these patterns are especially valuable, for example, because they can easily be exploited to increase turnover. For instance, if a supermarket discovers that certain products are frequently bought together, the number of items sold can sometimes be increased by appropriately arranging these products on the shelves of the market (they may, for example, be placed adjacent to each other in order to invite even more customers to buy them together, or they may be offered as a bundle).

However, to find these patterns and thus to exploit more of the information contained in the available data turns out to be fairly difficult. In contrast to the abundance of data there is a lack of tools to transform these data into useful knowledge. As John Naisbett remarked [Fayyad *et al.* 1996]:

We are drowning in information, but starving for knowledge.

As a consequence a new area of research has emerged, which has been named *Knowledge Discovery in Databases (KDD)* or *Data Mining (DM)* and which has taken up the challenge to develop techniques that can help humans to discover useful patterns and regularities in their data.

In this introductory chapter we provide a brief overview on knowledge discovery in databases and data mining, which is intended to show the context of this book. In a first step, we try to capture the difference between “data” and “knowledge” in order to attain precise notions by which it can be made clear why it does not suffice just to gather data and why we must strive to turn them into knowledge. As an illustration we will discuss and interpret a well-known example from the history of science. Secondly, we explain the process of discovering knowledge in databases (the *KDD process*), of which data mining is just one, though very important, step. We characterize the standard data mining tasks and position the work of this book by pointing out for which tasks the discussed methods are well suited.

## 1.1 Data and Knowledge

In this book we distinguish between *data* and *knowledge*. Statements like “Columbus discovered America in 1492” or “Mrs Jones owns a VW Golf” are *data*. For these statements to qualify as data, we consider it to be irrelevant whether we already know them, whether we need these specific pieces of information at this moment, etc. For our discussion, the essential property of these statements is that they refer to single events, cases, objects, persons, etc., in general, to single instances. Therefore, even if they are true, their range of validity is very restricted and thus is their usefulness.

In contrast to the above, *knowledge* consists of statements like “All masses attract each other.” or “Every day at 17:00 hours there runs an InterCity (a specific type of train of German Rail) from Magdeburg to Braunschweig.” Again we neglect the relevance of the statement for our current situation and whether we already know it. The essential property is that these statements do not refer to single instances, but are general laws or rules. Therefore, provided they are true, they have a wide range of validity, and, above all else, they allow us to make predictions and thus they are very useful.

It has to be admitted, though, that in daily life statements like “Columbus discovered America in 1492.” are also called knowledge. However, we disregard this way of using the term “knowledge”, regretting that full consistency of terminology with daily life language cannot be achieved. Collections of statements about single instances do not qualify as knowledge.

Summarizing, data and knowledge can be characterized as follows:

### Data

- refer to single instances  
(single objects, persons, events, points in time, etc.)
- describe individual properties
- are often available in huge amounts  
(databases, archives)

- are usually easy to collect or to obtain  
(for example cash registers with scanners in supermarkets, Internet)
- do not allow us to make predictions

### Knowledge

- refers to *classes* of instances  
(*sets* of objects, persons, events, points in time, etc.)
- describes general patterns, structures, laws, principles, etc.
- consists of as few statements as possible  
(this is an objective, see below)
- is usually hard to find or to obtain  
(for example natural laws, education)
- allows us to make predictions

From these characterizations we can clearly see that usually knowledge is much more valuable than (raw) data. It is mainly the generality of the statements and the possibility to make predictions about the behavior and the properties of new cases that constitute its superiority.

However, not just any kind of knowledge is as valuable as any other. Not all general statements are equally important, equally substantial, equally useful. Therefore knowledge must be evaluated and assessed. The following list, which we do not claim to be complete, names some important criteria:

### Criteria to Assess Knowledge

- correctness (probability, success in tests)
- generality (range of validity, conditions for validity)
- usefulness (relevance, predictive power)
- comprehensibility (simplicity, clarity, parsimony)
- novelty (previously unknown, unexpected)

In science correctness, generality, and simplicity (parsimony) are at the focus of attention: One way to characterize science is to say that it is the search for a minimal correct description of the world. In business and industry greater emphasis is placed on usefulness, comprehensibility, and novelty: the main goal is to get a competitive edge and thus to achieve higher profit. Nevertheless, none of the two areas can afford to neglect the other criteria.

## Tycho Brahe and Johannes Kepler

Tycho Brahe (1546–1601) was a Danish nobleman and astronomer, who in 1576 and in 1584, with the financial support of Frederic II, King of Denmark and Norway, built two observatories on the island of Sen, about 32 km to

the north-east of Copenhagen. Using the best equipment of his time (telescopes were unavailable then—they were used only later by Galileo Galilei (1564–1642) and Johannes Kepler (see below) for celestial observations) he determined the positions of the sun, the moon, and the planets with a precision of less than one minute of arc, thus surpassing by far the exactitude of all measurements carried out earlier. He achieved in practice the theoretical limit for observations with the unaided eye. Carefully he recorded the motions of the celestial bodies over several years [Greiner 1989, Zey 1997].

Tycho Brahe gathered data about our planetary system. Huge amounts of data—at least from a 16th century point of view. However, he could not discern the underlying structure. He could not combine his data into a consistent scheme—to some extent, because he adhered to the geocentric system. He could tell exactly in what position Mars had been on a specific day in 1585, but he could not relate the positions on different days in such a way as to fit his highly accurate observational data. All his hypotheses were fruitless. He developed the so-called Tychonic planetary model, according to which the sun and the moon revolve around the earth, but all other planets revolve around the sun, but this model, though popular in the 17th century, did not stand the test of time. Today we may say that Tycho Brahe had a “data mining” or “knowledge discovery” problem. He had the necessary data, but he could not extract the knowledge contained in it.

Johannes Kepler (1571–1630) was a German astronomer and mathematician and assistant to Tycho Brahe. He advocated the Copernican planetary model, and during his whole life he endeavored to find the laws that govern the motions of the celestial bodies. He strove to find a mathematical description, which, in his time, was a virtually radical approach. His starting point were the catalogs of data Tycho Brahe had compiled and which he continued in later years. After several unsuccessful trials and long and tedious calculations, Johannes Kepler finally managed to condense Tycho Brahe’s data into three simple laws, which have been named after him. Having discovered in 1604 that the course of Mars is an ellipse, he published the first two laws in “*Astronomia Nova*” in 1609, the third ten years later in his principal work “*Harmonica Mundi*” [Feynman *et al.* 1963, Greiner 1989, Zey 1997].

1. Each planet moves around the sun on an elliptical course, with the sun at one focus of the ellipse.
2. The radius vector from the sun to the planet sweeps out equal areas in equal intervals of time.
3. The squares of the periods of any two planets are proportional to the cubes of the semi-major axes of their respective orbits:  $T \sim a^{\frac{3}{2}}$ .

Tycho Brahe had collected a large amount of celestial data, Johannes Kepler found the laws by which they can be explained. He discovered the hidden knowledge and thus became one of the most famous “data miners” in history.

Today the works of Tycho Brahe are almost forgotten. His catalogs are merely of historical value. No textbook on astronomy contains extracts from his measurements. His observations and minute recordings are raw data and thus suffer from a decisive disadvantage: They do not provide us with any insight into the underlying mechanisms and therefore they do not allow us to make predictions. Kepler's laws, however, are treated in all textbooks on astronomy and physics, because they state the principles that govern the motions of planets as well as comets. They combine all of Brahe's measurements into three fairly simple statements. In addition, they allow us to make predictions: If we know the position and the velocity of a planet at a given moment, we can compute, using Kepler's laws, its future course.

## 1.2 Knowledge Discovery and Data Mining

How did Johannes Kepler discover his laws? How did he manage to extract from Tycho Brahe's long tables and voluminous catalogs those simple laws that revolutionized astronomy? We know only fairly little about this. He must have tested a large number of hypotheses, most of them failing. He must have carried out long and complicated computations. Presumably, outstanding mathematical talent, tenacious work, and a considerable amount of good luck finally led to success. We may safely guess that he did not know any universal method to discover physical or astronomical laws.

Today we still do not know such a method. It is still much simpler to gather data, by which we are virtually swamped in today's "information society" (whatever that means), than to obtain knowledge. We even need not work diligently and perseveringly any more, as Tycho Brahe did, in order to collect data. Automatic measurement devices, scanners, digital cameras, and computers have taken this load from us. Modern database technology enables us to store an ever-increasing amount of data. It is indeed as John Naisbett remarked: We are drowning in information, but starving for knowledge.

If it took such a distinguished mind like Johannes Kepler several years to evaluate the data gathered by Tycho Brahe, which today seem to be negligibly few and from which he even selected only the data on the course of Mars, how can we hope to cope with the huge amounts of data available today? "Manual" analysis has long ceased to be feasible. Simple aids like, for example, representations of data in charts and diagrams soon reach their limits. If we refuse to simply surrender to the flood of data, we are forced to look for intelligent computerized methods by which data analysis can be automated at least partially. These are the methods that are sought for in the research areas called *Knowledge Discovery in Databases (KDD)* and *Data Mining (DM)*. It is true, these methods are still very far from replacing people like Johannes Kepler, but it is not entirely implausible that he, if supported by these methods, would have reached his goal a little sooner.

Often the terms *Knowledge Discovery* and *Data Mining* are used interchangeably. However, we distinguish them here. By *Knowledge Discovery in Databases (KDD)* we mean a process consisting of several steps, which is usually characterized as follows [Fayyad *et al.* 1996]:

Knowledge discovery in databases is the nontrivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data.

One step of this process, though definitely one of the most important, is *Data Mining*. In this step modeling and discovery techniques are applied.

### 1.2.1 The KDD Process

In this section we structure the KDD process into two preliminary and five main steps or phases. However, the structure we discuss here is by no means binding: it has proven difficult to find a single scheme that everyone in the scientific community can agree on. However, an influential suggestion and detailed exposition of the KDD process, which is close to the scheme presented here and which has had considerable impact, because it is backed by several large companies like NCR and DaimlerChrysler, is the CRISP-DM model (CRoss Industry Standard Process for Data Mining) [Chapman *et al.* 1999].

#### Preliminary Steps

- estimation of potential benefit
- definition of goals, feasibility study

#### Main Steps

- check data availability, data selection, if necessary, data collection
- preprocessing (usually 60–90% of total overhead)
  - unification and transformation of data formats
  - data cleaning  
(error correction, outlier detection, imputation of missing values)
  - reduction / focusing  
(sample drawing, feature selection, prototype generation)
- **Data Mining** (using a variety of methods)
- visualization  
(also in parallel to preprocessing, data mining, and interpretation)
- interpretation, evaluation, and test of results
- deployment and documentation



The preliminary steps mainly serve the purpose to decide whether the main steps should be carried out. Only if the potential benefit is high enough and the demands can be met by data mining methods, can it be expected that some profit results from the usually expensive main steps.

In the main steps the data to be analyzed for hidden knowledge are first collected (if necessary), appropriate subsets are selected, and they are transformed into a unique format that is suitable for applying data mining techniques. Then they are cleaned and reduced to improve the performance of the algorithms to be applied later. These preprocessing steps usually consume the greater part of the total costs. Depending on the data mining task that was identified in the goal definition step (see below for a list), data mining methods are applied (see farther below for a list), the results of which, in order to interpret and evaluate them, can be visualized. Since the desired goal is rarely achieved in the first go, usually several steps of the preprocessing phase (for example feature selection) and the application of data mining methods have to be reiterated in order to improve the result. If it has not been obvious before, it is clear now that KDD is an interactive process, rather than completely automated. A user has to evaluate the results, check them for plausibility, and test them against hold-out data. If necessary, he/she modifies the course of the process to make it meet his/her requirements.

### 1.2.2 Data Mining Tasks

In the course of time typical tasks have been identified, which data mining methods should be able to solve (although, of course, not every single method is required to be able to solve all of them—it is the combination of methods that makes them powerful). Among these are especially those named in the—surely incomplete—list below. We tried to characterize them not only by their name, but also by a typical question [Nakhaeizadeh 1998b].

- classification  
*Is this customer credit-worthy?*
- segmentation, clustering  
*What groups of customers do I have?*
- concept description  
*Which properties characterize fault-prone vehicles?*
- prediction, trend analysis  
*What will the exchange rate of the dollar be tomorrow?*
- dependence/association analysis  
*Which products are frequently bought together?*
- deviation analysis  
*Are there seasonal or regional variations in turnover?*