

dnc



ETAPAS CRISP-DM

MATERIAL COMPLEMENTAR

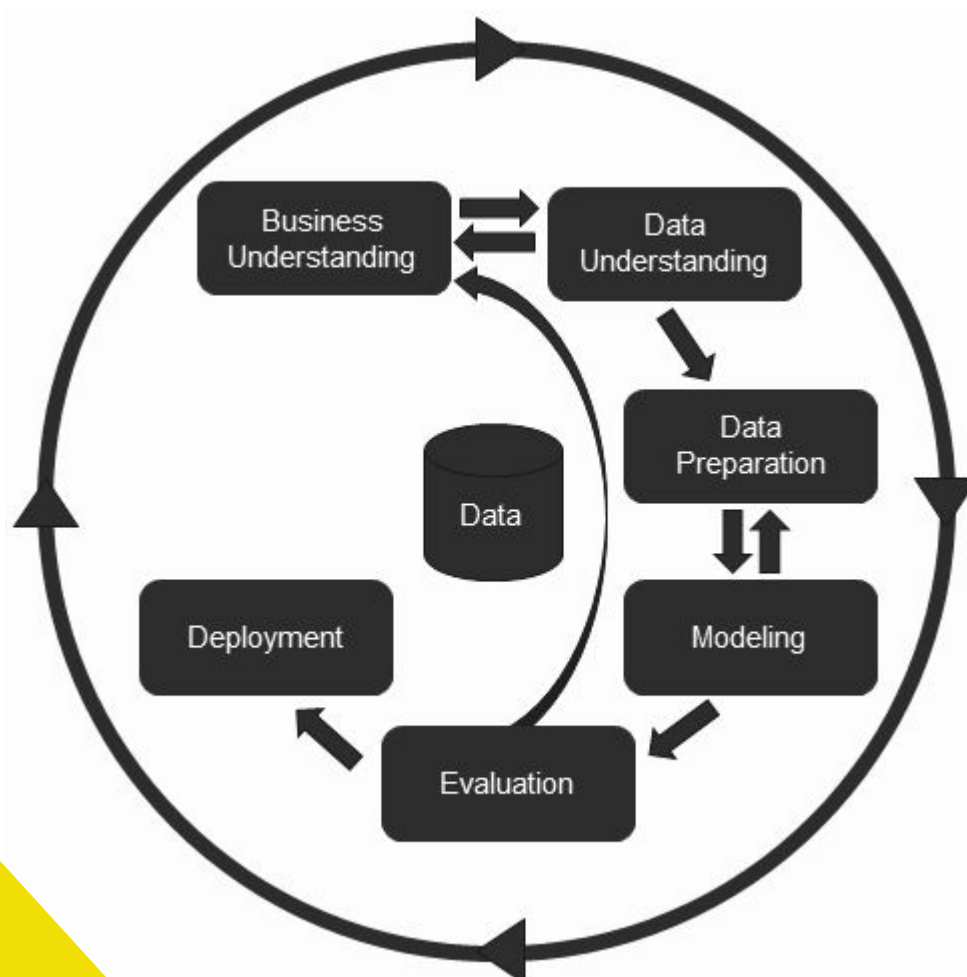


**FAÇA DOWNLOAD DESTE ARQUIVO E UTILIZE-O
ATÉ O FINAL DO PROJETO COMO GUIA OFICIAL!**

CRISP-DM

Cross Industry Standard Process for Data Mining

Processo criado em meados de 1996, o CRISP DM é hoje a metodologia mais utilizada nas empresas para guiar projetos de data science. Tem como objetivo auxiliar os profissionais de mineração de dados principalmente, mas também trabalhando desde sua etapa inicial com o entendimento do problema de negócio.





ETAPAS DO PROCESSO



Business Understanding

A primeira etapa trata de entender o contexto da empresa, quais os desafios a serem resolvidos. É fundamental nessa etapa traçar os objetivos a serem alcançados e expectativas do cliente em relação à entrega. Defina aqui também quais as métricas de sucesso a serem alcançados. Conheça as pessoas envolvidas: saiba quem são as pessoas do time de tecnologia, quem é o responsável pela área de negócios e os tomadores de decisão.

Mapeie os estão as coisas: nessa etapa explore onde estão armazenados e organizados os dados, quais ferramentas a empresa utiliza e acessos que precisa para o projeto.



Data Understanding

A compreensão dos dados é uma das etapas que se gasta mais tempo por ter que coletar e organizar, iniciando realmente o trabalho de mineração dos dados. É importante o profissional validar a qualidade das informações obtidas e a viabilidade de se desenvolver o projeto.

É a fase de se fazer o estudo e levantamento dos conteúdos, conseguindo mapear se existem anomalias, quais os dados faltantes, como será a extração e qual o custo da análise. Uma das competências que ressalta é o lado investigativo de buscar respostas através das informações garimpadas e ter um olhar analítico de entender o que os dados representam.



ETAPAS DO PROCESSO



Data Preparation

Depois do estudo e análise, chegou a hora de tratar esses dados brutos. Mas o que seria tratar os dados? É onde vai organizar quais os dados que serão modelados para o projeto, garantir que as informações estão alinhadas de acordo com o formato alinhado e o que será relevante para cruzar nas informações.

Aqui existem 4 etapas de preparação:

- **Data Selection:** a partir dos dados brutos, selecionar o que será modelado e documentar tudo;
- **Data Cleaning:** limpeza de fato, eliminando incoerências;
- **Construct Data:** com os dados finais obtidos na limpeza, vai identificar o que precisa construir e complementar para suas análises, criar campos ou tabelas, por exemplo;
- **Integrating Data:** fazer a integração de fontes distintas se for preciso.



ETAPAS DO PROCESSO



Modeling

Pronto, agora é a hora de modelar!

Depois que fez a mineração inicial, limpou e tratou seus dados, o próximo passo é desenvolver o seu modelo.

Para isso, precisa decidir quais as técnicas de data mining que serão aplicadas, fazer testes em amostras e verificar se o modelo está com um bom desempenho computacional para seguir. Estes dados minerados serão os insumos para os algoritmos na predição do modelo do negócio.



Evaluation

Após criar modelos e algoritmos, nessa etapa fará a avaliação dos resultados e quais as variações podem surgir. É importante observar se o modelo criado atende aos objetivos iniciais traçados na etapa de entendimento do negócio.

É uma etapa que pode apresentar uma necessidade de modificar o modelo por não atender aos critérios estabelecidos, não tendo assim uma boa performance.



ETAPAS DO PROCESSO



Deployment

Com o modelo desenvolvido, chegou a hora de colocá-lo em produção!

Na etapa de deployment, implementa-se o protótipo e o cliente faz o acompanhamento de sua performance após sua entrega.

Esse é o momento que a equipe entrega o modelo e organiza o conhecimento para poder ser repassado ao usuário.



CRISP-DM 1.0

Step-by-step data mining guide

Pete Chapman (NCR), Julian Clinton (SPSS), Randy Kerber (NCR),
Thomas Khabaza (SPSS), Thomas Reinartz (DaimlerChrysler),
Colin Shearer (SPSS) and Rüdiger Wirth (DaimlerChrysler)



This document describes the CRISP-DM process model and contains information about the CRISP-DM methodology, the CRISP-DM reference model, the CRISP-DM user guide, and the CRISP-DM reports, as well as an appendix with additional related information. This document and information herein are the exclusive property of the partners of the CRISP-DM consortium: NCR Systems Engineering Copenhagen (USA and Denmark), DaimlerChrysler AG (Germany), SPSS Inc. (USA), and OHRA Verzekeringen en Bank Groep B.V. (The Netherlands).

Copyright © 1999, 2000

All trademarks and service marks mentioned in this document are marks of their respective owners and are as such acknowledged by the members of the CRISP-DM consortium.

Foreword

CRISP-DM was conceived in late 1996 by three “veterans” of the young and immature data mining market. DaimlerChrysler (then Daimler-Benz) was already ahead of most industrial and commercial organizations in applying data mining in its business operations. SPSS (then ISL) had been providing services based on data mining since 1990 and had launched the first commercial data mining workbench—Clementine®—in 1994. NCR, as part of its aim to deliver added value to its Teradata® data warehouse customers, had established teams of data mining consultants and technology specialists to service its clients’ requirements.

At that time, early market interest in data mining was showing signs of exploding into widespread uptake. This was both exciting and terrifying. All of us had developed our approaches to data mining as we went along. Were we doing it right? Was every new adopter of data mining going to have to learn, as we had initially, by trial and error? And from a supplier’s perspective, how could we demonstrate to prospective customers that data mining was sufficiently mature to be adopted as a key part of their business processes?

A standard process model, we reasoned, non-proprietary and freely available, would address these issues for us and for all practitioners.

A year later, we had formed a consortium, invented an acronym (CRoss-Industry Standard Process for Data Mining), obtained funding from the European Commission, and begun to set out our initial ideas. As CRISP-DM was intended to be industry-, tool-, and application-neutral, we knew we had to get input from as wide a range as possible of practitioners and others (such as data warehouse vendors and management consultancies) with a vested interest in data mining. We did this by creating the CRISP-DM Special Interest Group (“The SIG,” as it became known). We launched the SIG by broadcasting an invitation to interested parties to join us in Amsterdam for a day-long workshop: We would share our ideas, invite them to present theirs, and openly discuss how to take CRISP-DM forward.

On the day of the workshop, there was a feeling of trepidation among the consortium members. Would no one be interested enough to show up? Or, if they did, would they tell us they really didn’t see a compelling need for a standard process? Or that our ideas were so far out of step with others’ that any idea of standardization was an impractical fantasy?

The workshop surpassed all our expectations. Three things stood out:

- Twice as many people turned up as we had initially expected
- There was an overwhelming consensus that the industry needed a standard process and needed it now
- As attendees presented their views on data mining from their project experience, it became clear that although there were superficial differences—mainly in demarcation of phases and in terminology—there was tremendous common ground in how they viewed the process of data mining



By the end of the workshop, we felt confident that we could deliver, with the SIG's input and critique, a standard process model to service the data mining community.

Over the next two and a half years, we worked to develop and refine CRISP-DM. We ran trials in live, large-scale data mining projects at Mercedes-Benz and at our insurance sector partner, OHRA. We worked on the integration of CRISP-DM with commercial data mining tools. The SIG proved invaluable, growing to over 200 members and holding workshops in London, New York, and Brussels.

By the end of the EC-funded part of the project—mid-1999—we had produced what we considered a good-quality draft of the process model. Those familiar with that draft will find that a year later, although now much more complete and better presented, CRISP-DM 1.0 is by no means radically different. We were acutely aware that, during the project, the process model was still very much a work-in-progress; CRISP-DM had only been validated on a narrow set of projects. Over the past year, DaimlerChrysler had the opportunity to apply CRISP-DM to a wider range of applications. SPSS' and NCR's Professional Services groups have adopted CRISP-DM and used it successfully on numerous customer engagements covering many industries and business problems. Throughout this time, we have seen service suppliers from outside the consortium adopt CRISP-DM, repeated references to it by analysts as the de facto standard for the industry, and a growing awareness of its importance among customers (CRISP-DM is now frequently referenced in invitations to tender and in RFP documents). We believe our initiative has been thoroughly vindicated, and while future extensions and improvements are both desirable and inevitable, we consider CRISP-DM Version 1.0 sufficiently validated to be published and distributed.

CRISP-DM has not been built in a theoretical, academic manner working from technical principles, nor did elite committees of gurus create it behind closed doors. Both these approaches to developing methodologies have been tried in the past, but have seldom led to practical, successful, and widely adopted standards. CRISP-DM succeeds because it is soundly based on the practical, real-world experience of how people conduct data mining projects. And in that respect, we are overwhelmingly indebted to the many practitioners who contributed their efforts and their ideas throughout the project.

The CRISP-DM consortium
August 2000

Table of contents

| | | |
|-----|--|----|
| I | Introduction | 6 |
| 1 | The CRISP-DM methodology | 6 |
| 1.1 | Hierarchical breakdown | 6 |
| 1.2 | Reference model and user guide | 7 |
| 2 | Mapping generic models to specialized models | 7 |
| 2.1 | Data mining context | 7 |
| 2.2 | Mappings with contexts | 8 |
| 2.3 | How to map | 8 |
| 3 | Description of parts | 9 |
| 3.1 | Contents | 9 |
| 3.2 | Purpose | 9 |
| II | The CRISP-DM reference model | 10 |
| 1 | Business understanding | 13 |
| 1.1 | Determine business objectives | 14 |
| 1.2 | Assess situation | 14 |
| 1.3 | Determine data mining goals | 16 |
| 1.4 | Produce project plan | 16 |
| 2 | Data understanding | 17 |
| 2.1 | Collect initial data | 18 |
| 2.2 | Describe data | 18 |
| 2.3 | Explore data | 18 |
| 2.4 | Verify data quality | 19 |
| 3 | Data preparation | 20 |
| 3.1 | Select data | 21 |
| 3.2 | Clean data | 21 |
| 3.3 | Construct data | 21 |
| 3.4 | Integrate data | 22 |
| 3.5 | Format data | 22 |

| | | |
|-----|---------------------------------------|----|
| 4 | Modeling | 23 |
| 4.1 | Select modeling technique | 24 |
| 4.2 | Generate test design | 24 |
| 4.3 | Build model | 24 |
| 4.4 | Assess model | 25 |
| 5 | Evaluation | 26 |
| 5.1 | Evaluate results | 26 |
| 5.2 | Review process | 27 |
| 5.3 | Determine next steps | 27 |
| 6 | Deployment | 28 |
| 6.1 | Plan deployment | 28 |
| 6.2 | Plan monitoring and maintenance | 29 |
| 6.3 | Produce final report | 29 |
| 6.4 | Review project | 29 |
| III | The CRISP-DM user guide | 30 |
| 1 | Business understanding | 30 |
| 1.1 | Determine business objectives | 30 |
| 1.2 | Assess situation | 32 |
| 1.3 | Determine data mining goals | 35 |
| 1.4 | Produce project plan | 36 |
| 2 | Data understanding | 37 |
| 2.1 | Collect initial data | 37 |
| 2.2 | Describe data | 39 |
| 2.3 | Explore data | 40 |
| 2.4 | Verify data quality | 41 |
| 3 | Data preparation | 42 |
| 3.1 | Select data | 42 |
| 3.2 | Clean data | 43 |
| 3.3 | Construct data | 44 |
| 3.4 | Integrate data | 46 |
| 3.5 | Format data | 46 |

| | | |
|-----|--|----|
| 4 | Modeling | 47 |
| 4.1 | Select modeling technique | 47 |
| 4.2 | Generate test design | 49 |
| 4.3 | Build model | 49 |
| 4.4 | Assess model | 50 |
| 5 | Evaluation | 51 |
| 5.1 | Evaluate results | 52 |
| 5.2 | Review process | 53 |
| 5.3 | Determine next steps | 53 |
| 6 | Deployment | 54 |
| 6.1 | Plan deployment | 54 |
| 6.2 | Plan monitoring and maintenance | 55 |
| 6.3 | Produce final report | 55 |
| 6.4 | Review project | 56 |
| IV | The CRISP-DM outputs | 57 |
| 1 | Business understanding | 57 |
| 2 | Data understanding | 58 |
| 3 | Data preparation | 60 |
| 4 | Modeling | 60 |
| 5 | Evaluation | 62 |
| 6 | Deployment | 62 |
| 7 | Summary of dependencies | 64 |
| V | Appendix | 65 |
| 1 | Glossary/terminology | 65 |
| 2 | Data mining problem types | 66 |
| 2.1 | Data description and summarization | 66 |
| 2.2 | Segmentation | 67 |
| 2.3 | Concept descriptions | 68 |
| 2.4 | Classification | 69 |
| 2.5 | Prediction | 70 |
| 2.6 | Dependency analysis | 70 |

Introduction

The CRISP-DM methodology

1.1 Hierarchical breakdown

The CRISP-DM methodology is described in terms of a hierarchical process model, consisting of sets of tasks described at four levels of abstraction (from general to specific): phase, generic task, specialized task, and process instance (see figure 1).

At the top level, the data mining process is organized into a number of phases; each phase consists of several second-level generic tasks. This second level is called generic because it is intended to be general enough to cover all possible data mining situations. The generic tasks are intended to be as complete and stable as possible. Complete means covering both the whole process of data mining and all possible data mining applications. Stable means that the model should be valid for yet unforeseen developments like new modeling techniques.

The third level, the specialized task level, is the place to describe how actions in the generic tasks should be carried out in certain specific situations. For example, at the second level there might be a generic task called clean data. The third level describes how this task differs in different situations, such as cleaning numeric values versus cleaning categorical values, or whether the problem type is clustering or predictive modeling.

The description of phases and tasks as discrete steps performed in a specific order represents an idealized sequence of events. In practice, many of the tasks can be performed in a different order, and it will often be necessary to repeatedly backtrack to previous tasks and repeat certain actions. Our process model does not attempt to capture all of these possible routes through the data mining process because this would require an overly complex process model.

The fourth level, the process instance, is a record of the actions, decisions, and results of an actual data mining engagement. A process instance is organized according to the tasks defined at the higher levels, but represents what actually happened in a particular engagement, rather than what happens in general.

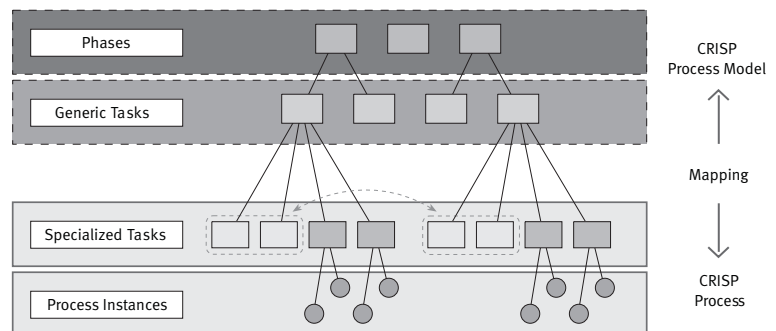


Figure 1: Four level breakdown of the CRISP-DM methodology

1.2 Reference model and user guide

Horizontally, the CRISP-DM methodology distinguishes between the reference model and the user guide. The reference model presents a quick overview of phases, tasks, and their outputs, and describes what to do in a data mining project. The user guide gives more detailed tips and hints for each phase and each task within a phase, and depicts how to carry out a data mining project.

This document covers both the reference model and the user guide at the generic level.

Mapping generic models to specialized models

2.1 Data mining context


The data mining context drives mapping between the generic and the specialized level in CRISP-DM. Currently, we distinguish between four different dimensions of data mining contexts:

- The **application domain** is the specific area in which the data mining project takes place
- The **data mining problem type** describes the specific class(es) of objective(s) that the data mining project deals with (see also Appendix 2)
- The **technical aspect** covers specific issues in data mining that describe different (technical) challenges that usually occur during data mining
- The **tool and technique** dimension specifies which data mining tool(s) and/or techniques are applied during the data mining project

Table 1 below summarizes these dimensions of data mining contexts and shows specific examples for each dimension.

| Dimension | Data Mining Context | | | |
|-----------|---------------------|-------------------------------|------------------|--------------------|
| | Application Domain | Data Mining Problem Type | Technical Aspect | Tool and Technique |
| Examples | Response Modeling | Description and Summarization | Missing Values | Clementine |
| | Churn Prediction | Segmentation | Outliers | MineSet |
| | ... | Concept Description | ... | Decision Tree |
| | | Classification | | ... |
| | | Prediction | | |
| | | Dependency Analysis | | |

Table 1: Dimensions of data mining contexts and examples



A specific data mining context is a concrete value for one or more of these dimensions. For example, a data mining project dealing with a classification problem in churn prediction constitutes one specific context. The more values for different context dimensions are fixed, the more concrete is the data mining context.

2.2 Mappings with contexts

We distinguish between two different types of mapping between generic and specialized level in CRISP-DM.

Mapping for the present: If we only apply the generic process model to perform a single data mining project, and attempt to map generic tasks and their descriptions to the specific project as required, we talk about a single mapping for (probably) only one usage.

Mapping for the future: If we systematically specialize the generic process model according to a pre-defined context (or similarly systematically analyze and consolidate experiences of a single project toward a specialized process model for future usage in comparable contexts), we talk about explicitly writing up a specialized process model in terms of CRISP-DM.

Which type of mapping is appropriate for your own purposes depends on your specific data mining context and the needs of your organization.

2.3 How to map

The basic strategy for mapping the generic process model to the specialized level is the same for both types of mappings:

- Analyze your specific context
- Remove any details not applicable to your context
- Add any details specific to your context
- Specialize (or instantiate) generic contents according to concrete characteristics of your context
- Possibly rename generic contents to provide more explicit meanings in your context for the sake of clarity

Description of parts

3.1 Contents

The CRISP-DM process model (this document) is organized into five different parts:

- Part I is this introduction to the CRISP-DM methodology, which provides some general guidelines for mapping the generic process model to specialized process models
- Part II describes the CRISP-DM reference model, its phases, generic tasks, and outputs
- Part III presents the CRISP-DM user guide, which goes beyond the pure description of phases, generic tasks, and outputs, and contains more detailed advice on how to perform data mining projects
- Part IV focuses on the reports to be produced during and after a project, and suggests outlines for these reports. It also shows cross references among outputs and tasks.
- Part V is the appendix, which includes a glossary of important terminology and a characterization of data mining problem types

3.2 Purpose

Users and readers of this document should be aware of the following instructions:

- If you are reading the CRISP-DM process model for the first time, begin with part I, the introduction, in order to understand the CRISP-DM methodology, all of its concepts, and how different concepts relate to each other. In further readings, you might skip the introduction and only return to it if necessary for clarification.
- If you need fast access to an overview of the CRISP-DM process model, refer to part II, the CRISP-DM reference model, either to begin a data mining project quickly or to get an introduction to the CRISP-DM user guide
- If you need detailed advice in performing your data mining project, part III, the CRISP-DM user guide, is the most valuable part of this document. Note: if you have not read the introduction or the reference model first, go back and read these first two parts.
- If you are at the stage of data mining when you write up your reports, go to part IV. If you prefer to generate your deliverable descriptions during the project, move back and forth between parts III and IV as desired.
- Finally, the appendix is useful as additional background information on CRISP-DM and data mining. Use the appendix to look up various terms if you are not yet an expert in the field.

II The CRISP-DM reference model

The current process model for data mining provides an overview of the life cycle of a data mining project. It contains the phases of a project, their respective tasks, and the relationships between these tasks. At this description level, it is not possible to identify all relationships. Relationships could exist between any data mining tasks depending on the goals, the background, and the interest of the user—and most importantly—on the data.

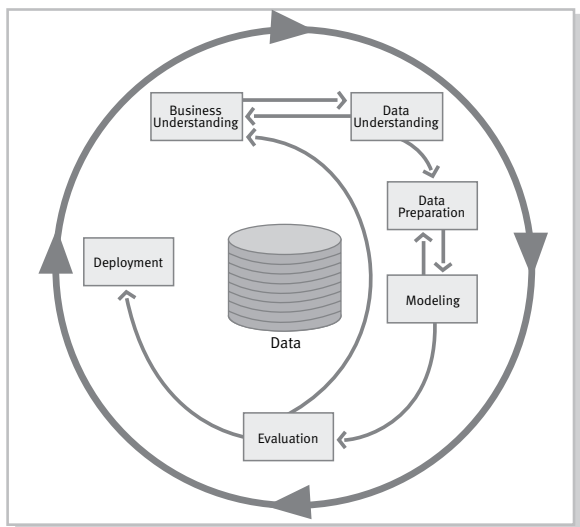


Figure 2: Phases of the CRISP-DM reference model

The life cycle of a data mining project consists of six phases, shown in **Figure 2**. The sequence of the phases is not rigid. Moving back and forth between different phases is always required. The outcome of each phase determines which phase, or particular task of a phase, has to be performed next. The arrows indicate the most important and frequent dependencies between phases.

The outer circle in Figure 2 symbolizes the cyclical nature of data mining itself. Data mining does not end once a solution is deployed. The lessons learned during the process and from the deployed solution can trigger new, often more-focused business questions. Subsequent data mining processes will benefit from the experiences of previous ones. In the following, we briefly outline each phase:

Business understanding

This initial phase focuses on understanding the project objectives and requirements from a business perspective, then converting this knowledge into a data mining problem definition and a preliminary plan designed to achieve the objectives.

Data understanding

The data understanding phase starts with initial data collection and proceeds with activities that enable you to become familiar with the data, identify data quality problems, discover first insights into the data, and/or detect interesting subsets to form hypotheses regarding hidden information.

Data preparation

The data preparation phase covers all activities needed to construct the final dataset [data that will be fed into the modeling tool(s)] from the initial raw data. Data preparation tasks are likely to be performed multiple times and not in any prescribed order. Tasks include table, record, and attribute selection, as well as transformation and cleaning of data for modeling tools.

Modeling

In this phase, various modeling techniques are selected and applied, and their parameters are calibrated to optimal values. Typically, there are several techniques for the same data mining problem type. Some techniques have specific requirements on the form of data. Therefore, going back to the data preparation phase is often necessary.

Evaluation

At this stage in the project, you have built a model (or models) that appears to have high quality from a data analysis perspective. Before proceeding to final deployment of the model, it is important to thoroughly evaluate it and review the steps executed to create it, to be certain the model properly achieves the business objectives. A key objective is to determine if there is some important business issue that has not been sufficiently considered. At the end of this phase, a decision on the use of the data mining results should be reached.

Deployment

Creation of the model is generally not the end of the project. Even if the purpose of the model is to increase knowledge of the data, the knowledge gained will need to be organized and presented in a way that the customer can use it. It often involves applying “live” models within an organization’s decision making processes—for example, real-time personalization of Web pages or repeated scoring of marketing databases. Depending on the requirements, the deployment phase can be as simple as generating a report or as complex as implementing a repeatable data mining process across the enterprise. In many cases, it is the customer, not the data analyst, who carries out the deployment steps. However, even if the analyst will carry out the deployment effort, it is important for the customer to understand up front what actions need to be carried out in order to actually make use of the created models.

Figure 3 presents an outline of phases accompanied by generic tasks (bold) and outputs (italic). In the following sections, we describe each generic task and its outputs in more detail. We focus our attention on task overviews and summaries of outputs.

| Business Understanding | Data Understanding | Data Preparation | Modeling | Evaluation | Deployment |
|--|--|---|---|--|--|
| Determine Business Objectives <i>Background Business Objectives Business Success Criteria</i> | Collect Initial Data <i>Initial Data Collection Report</i> | Select Data <i>Rationale for Inclusion/Exclusion</i> | Select Modeling Techniques <i>Modeling Technique Modeling Assumptions</i> | Evaluate Results <i>Assessment of Data Mining Results w.r.t. Business Success Criteria Approved Models</i> | Plan Deployment <i>Deployment Plan</i> |
| Assess Situation <i>Inventory of Resources Requirements, Assumptions, and Constraints Risks and Contingencies Terminology Costs and Benefits</i> | Describe Data <i>Data Description Report</i> | Clean Data <i>Data Cleaning Report</i> | Generate Test Design <i>Test Design</i> | Review Process <i>Review of Process</i> | Plan Monitoring and Maintenance <i>Monitoring and Maintenance Plan</i> |
| Determine Data Mining Goals <i>Data Mining Goals Data Mining Success Criteria</i> | Explore Data <i>Data Exploration Report</i> | Construct Data <i>Derived Attributes Generated Records</i> | Build Model <i>Parameter Settings Models Model Descriptions</i> | Determine Next Steps <i>List of Possible Actions Decision</i> | Produce Final Report <i>Final Report Final Presentation</i> |
| Produce Project Plan <i>Project Plan Initial Assessment of Tools and Techniques</i> | Verify Data Quality <i>Data Quality Report</i> | Integrate Data <i>Merged Data</i> | Assess Model <i>Model Assessment Revised Parameter Settings</i> | | Review Project <i>Experience Documentation</i> |
| | | Format Data <i>Reformatted Data Dataset Dataset Description</i> | | | |

Figure 3: Generic tasks (bold) and outputs (italic) of the CRISP-DM reference model

1 Business understanding

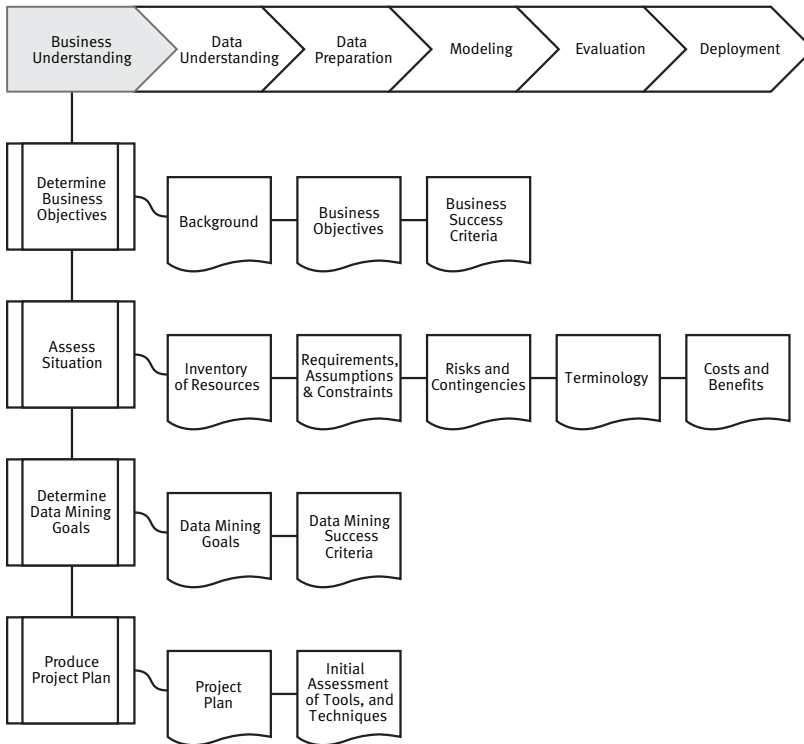



Figure 4: Business Understanding



1.1 Determine business objectives

Task

Determine business objectives

The first objective of the data analyst is to thoroughly understand, from a business perspective, what the customer really wants to accomplish. Often the customer has many competing objectives and constraints that must be properly balanced. The analyst's goal is to uncover important factors, at the beginning, that can influence the outcome of the project. A possible consequence of neglecting this step is to expend a great deal of effort producing the right answers to the wrong questions.

Outputs

Background

Record the information that is known about the organization's business situation at the beginning of the project.

Business objectives

Describe the customer's primary objective, from a business perspective. In addition to the primary business objective, there are typically other related business questions that the customer would like to address. For example, the primary business goal might be to keep current customers by predicting when they are prone to move to a competitor. Examples of related business questions are "How does the primary channel used (e.g., ATM, branch visit, Internet) affect whether customers stay or go?" or "Will lower ATM fees significantly reduce the number of high-value customers who leave?"

Business success criteria

Describe the criteria for a successful or useful outcome to the project from the business point of view. This might be quite specific and able to be measured objectively, for example, reduction of customer churn to a certain level, or it might be general and subjective, such as "give useful insights into the relationships." In the latter case, it should be indicated who makes the subjective judgment.

1.2 Assess situation

Task

Assess situation

This task involves more detailed fact-finding about all of the resources, constraints, assumptions, and other factors that should be considered in determining the data analysis goal and project plan. In the previous task, your objective is to quickly get to the crux of the situation. Here, you want to expand upon the details.

Outputs

Inventory of resources

List the resources available to the project, including personnel (business experts, data experts, technical support, data mining experts), data (fixed extracts, access to live, warehoused, or operational data), computing resources (hardware platforms), and software (data mining tools, other relevant software).

Requirements, assumptions, and constraints

List all requirements of the project, including schedule of completion, comprehensibility and quality of results, and security, as well as legal issues. As part of this output, make sure that you are allowed to use the data.

List the assumptions made by the project. These may be assumptions about the data that can be verified during data mining, but may also include non-verifiable assumptions about the business related to the project. It is particularly important to list the latter if it will affect the validity of the results.

List the constraints on the project. These may be constraints on the availability of resources, but may also include technological constraints such as the size of dataset that it is practical to use for modeling.

Risks and contingencies

List the risks or events that might delay the project or cause it to fail. List the corresponding contingency plans, what action will be taken if these risks or events take place.


Terminology

Compile a glossary of terminology relevant to the project. This may include two components:

- (1) A glossary of relevant business terminology, which forms part of the business understanding available to the project. Constructing this glossary is a useful “knowledge elicitation” and education exercise.
- (2) A glossary of data mining terminology, illustrated with examples relevant to the business problem in question

Costs and benefits

Construct a cost-benefit analysis for the project, which compares the costs of the project with the potential benefits to the business if it is successful. The comparison should be as specific as possible. For example, use monetary measures in a commercial situation.



1.3 Determine data mining goals

Task

Determine data mining goals

A business goal states objectives in business terminology. A data mining goal states project objectives in technical terms. For example, the business goal might be “Increase catalog sales to existing customers.” A data mining goal might be “Predict how many widgets a customer will buy, given their purchases over the past three years, demographic information (age, salary, city, etc.), and the price of the item.”

Outputs

Data mining goals

Describe the intended outputs of the project that enable the achievement of the business objectives.

Data mining success criteria

Define the criteria for a successful outcome to the project in technical terms—for example, a certain level of predictive accuracy or a propensity-to-purchase profile with a given degree of “lift.” As with business success criteria, it may be necessary to describe these in subjective terms, in which case the person or persons making the subjective judgment should be identified.

1.4 Produce project plan

Task

Produce project plan

Describe the intended plan for achieving the data mining goals and thereby achieving the business goals. The plan should specify the steps to be performed during the rest of the project, including the initial selection of tools and techniques.

Outputs

Project plan

List the stages to be executed in the project, together with their duration, resources required, inputs, outputs, and dependencies. Where possible, make explicit the large-scale iterations in the data mining process—for example, repetitions of the modeling and evaluation phases.

As part of the project plan, it is also important to analyze dependencies between time schedule and risks. Mark results of these analyses explicitly in the project plan, ideally with actions and recommendations if the risks are manifested.

Note: the project plan contains detailed plans for each phase. Decide at this point which evaluation strategy will be used in the evaluation phase.

The project plan is a dynamic document in the sense that at the end of each phase, a review of progress and achievements is necessary and a corresponding update of the project plan is recommended. Specific review points for these updates are part of the project plan.

Initial assessment of tools and techniques

At the end of the first phase, an initial assessment of tools and techniques should be performed. Here, for example, you select a data mining tool that supports various methods for different stages of the process. It is important to assess tools and techniques early in the process since the selection of tools and techniques may influence the entire project.

2 Data understanding

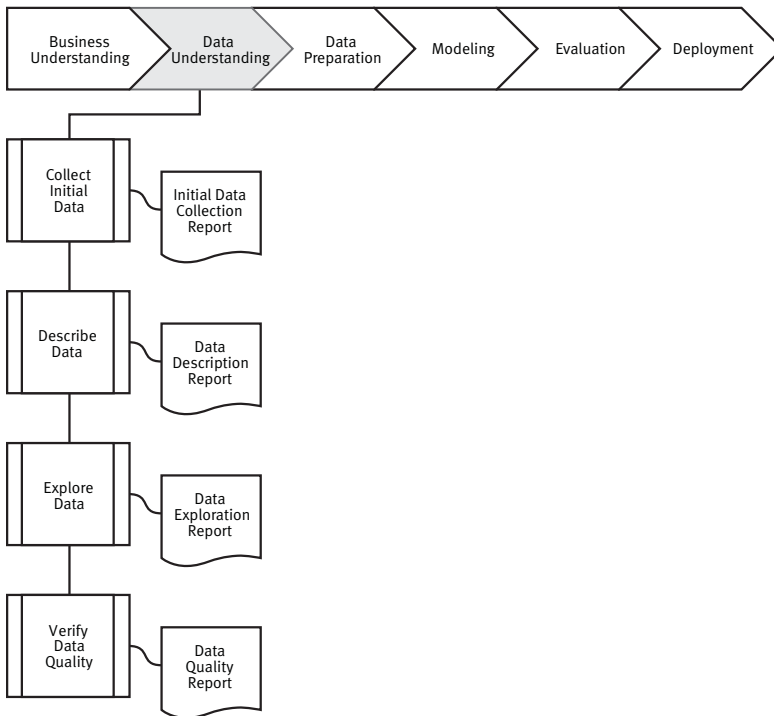


Figure 5: Data understanding



2.1 Collect initial data

Task

Collect initial data

Acquire the data (or access to the data) listed in the project resources. This initial collection includes data loading, if necessary for data understanding. For example, if you use a specific tool for data understanding, it makes perfect sense to load your data into this tool. This effort possibly leads to initial data preparation steps.

Note: if you acquire multiple data sources, integration is an additional issue, either here or in the later data preparation phase.

Output

Initial data collection report

List the dataset(s) acquired, together with their locations, the methods used to acquire them, and any problems encountered. Record problems encountered and any resolutions achieved. This will aid with future replication of this project or with the execution of similar future projects.

2.2 Describe data

Task

Describe data

Examine the “gross” or “surface” properties of the acquired data and report on the results.

Output

Data description report

Describe the data that has been acquired, including the format of the data, the quantity of data (for example, the number of records and fields in each table), the identities of the fields, and any other surface features which have been discovered. Evaluate whether the data acquired satisfies the relevant requirements.

2.3 Explore data

Task

Explore data

This task addresses data mining questions using querying, visualization, and reporting techniques. These include distribution of key attributes (for example, the target attribute of a prediction task) relationships between pairs or small numbers of attributes, results of simple aggregations, properties of significant sub-populations, and simple statistical analyses. These analyses may directly address the data mining goals; they may also contribute to or refine the data description and quality reports, and feed into the transformation and other data preparation steps needed for further analysis.

Output**Data exploration report**

Describe results of this task, including first findings or initial hypothesis and their impact on the remainder of the project. If appropriate, include graphs and plots to indicate data characteristics that suggest further examination of interesting data subsets.

2.4 Verify data quality**Task****Verify data quality**

Examine the quality of the data, addressing questions such as: Is the data complete (does it cover all the cases required)? Is it correct, or does it contain errors and, if there are errors, how common are they? Are there missing values in the data? If so, how are they represented, where do they occur, and how common are they?

Output**Data quality report**

List the results of the data quality verification; if quality problems exist, list possible solutions. Solutions to data quality problems generally depend heavily on both data and business knowledge.

3 Data preparation

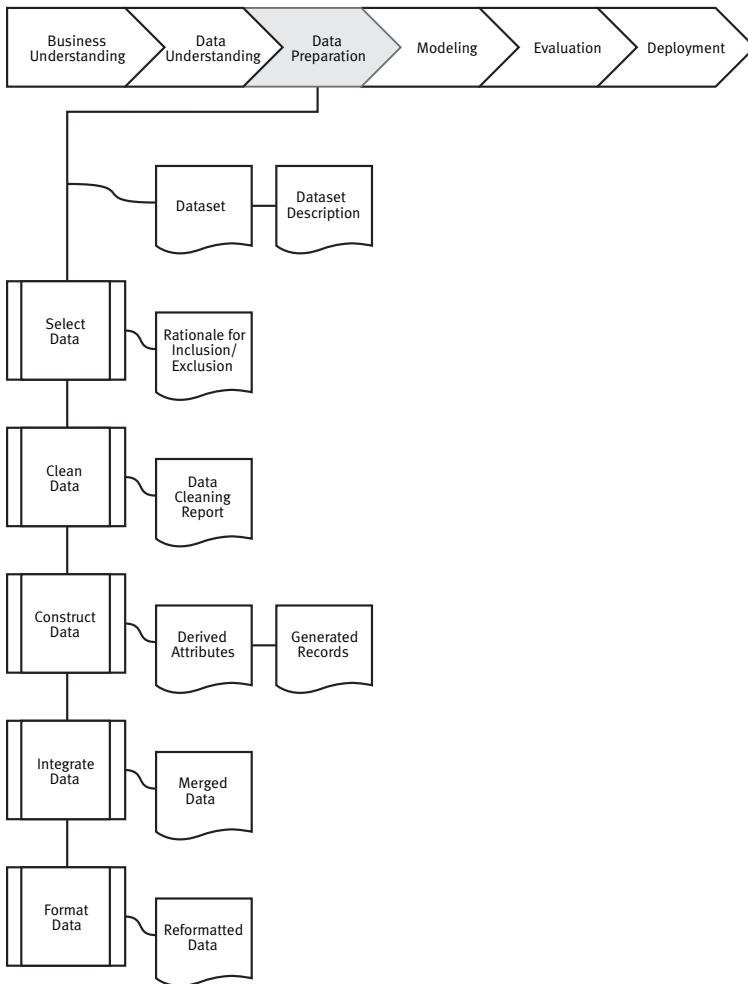


Figure 6: Data preparation

Outputs

Dataset

These are the dataset(s) produced by the data preparation phase, which will be used for modeling or the major analysis work of the project.

Dataset description

Describe the dataset(s) that will be used for the modeling and the major analysis work of the project.

3.1 Select data

Task

Select data

Decide on the data to be used for analysis. Criteria include relevance to the data mining goals, quality, and technical constraints such as limits on data volume or data types. Note that data selection covers selection of attributes (columns) as well as selection of records (rows) in a table.

Output

Rationale for inclusion/exclusion

List the data to be included/excluded and the reasons for these decisions.

3.2 Clean data

Task

Clean data

Raise the data quality to the level required by the selected analysis techniques. This may involve selection of clean subsets of the data, the insertion of suitable defaults, or more ambitious techniques such as the estimation of missing data by modeling.

Output

Data cleaning report

Describe what decisions and actions were taken to address the data quality problems reported during the Verify Data Quality task of the Data Understanding phase. Transformations of the data for cleaning purposes and the possible impact on the analysis results should be considered.

3.3 Construct data

Task

Construct data

This task includes constructive data preparation operations such as the production of derived attributes or entire new records, or transformed values for existing attributes.



Outputs

Derived attributes

Derived attributes are new attributes that are constructed from one or more existing attributes in the same record. Example: $\text{area} = \text{length} * \text{width}$.

Generated records

Describe the creation of completely new records. Example: Create records for customers who made no purchase during the past year. There was no reason to have such records in the raw data, but for modeling purposes it might make sense to explicitly represent the fact that certain customers made zero purchases.

3.4 Integrate data

Task

Integrate data

These are methods whereby information is combined from multiple tables or records to create new records or values.

Output

Merged data

Merging tables refers to joining together two or more tables that have different information about the same objects. Example: a retail chain has one table with information about each store's general characteristics (e.g., floor space, type of mall), another table with summarized sales data (e.g., profit, percent change in sales from previous year), and another with information about the demographics of the surrounding area. Each of these tables contains one record for each store. These tables can be merged together into a new table with one record for each store, combining fields from the source tables.

Merged data also covers aggregations. Aggregation refers to operations in which new values are computed by summarizing information from multiple records and/or tables. For example, converting a table of customer purchases where there is one record for each purchase into a new table where there is one record for each customer, with fields such as number of purchases, average purchase amount, percent of orders charged to credit card, percent of items under promotion, etc.

3.5 Format data

Task

Format data

Formatting transformations refer to primarily *syntactic* modifications made to the data that do not change its meaning, but might be required by the modeling tool.

Output

Reformatted data

Some tools have requirements on the order of the attributes, such as the first field being a unique identifier for each record or the last field being the outcome field the model is to predict.

It might be important to change the order of the records in the dataset. Perhaps the modeling tool requires that the records be sorted according to the value of the outcome attribute. Commonly, records of the dataset are initially ordered in some way, but the modeling algorithm needs them to be in a fairly random order. For example, when using neural networks, it is generally best for the records to be presented in a random order, although some tools handle this automatically without explicit user intervention.

Additionally, there are purely syntactic changes made to satisfy the requirements of the specific modeling tool. Examples: removing commas from within text fields in comma-delimited data files, trimming all values to a maximum of 32 characters.

4 Modeling

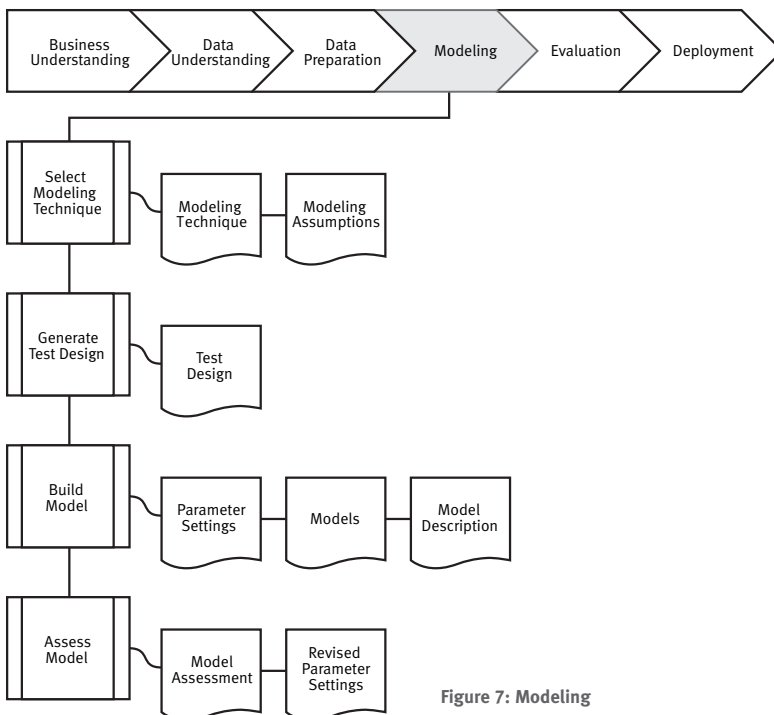


Figure 7: Modeling



4.1 *Select modeling technique*

Task **Select modeling technique**

As the first step in modeling, select the actual modeling technique that is to be used. Although you may have already selected a tool during the Business Understanding phase, this task refers to the specific modeling technique, e.g., decision-tree building with 5.0, or neural network generation with back propagation. If multiple techniques are applied, perform this task separately for each technique.

Outputs **Modeling technique**

Document the actual modeling technique that is to be used.

Modeling assumptions

Many modeling techniques make specific assumptions about the data—for example, that all attributes have uniform distributions, no missing values allowed, class attribute must be symbolic, etc. Record any such assumptions made.

4.2 *Generate test design*

Task **Generate test design**

Before we actually build a model, we need to generate a procedure or mechanism to test the model's quality and validity. For example, in supervised data mining tasks such as classification, it is common to use error rates as quality measures for data mining models. Therefore, we typically separate the dataset into train and test sets, build the model on the train set, and estimate its quality on the separate test set.

Output **Test design**

Describe the intended plan for training, testing, and evaluating the models. A primary component of the plan is determining how to divide the available dataset into training, test, and validation datasets.

4.3 *Build model*

Task **Build model**

Run the modeling tool on the prepared dataset to create one or more models.

Outputs **Parameter settings**

With any modeling tool, there are often a large number of parameters that can be adjusted. List the parameters and their chosen values, along with the rationale for the choice of parameter settings.

Models

These are the actual models produced by the modeling tool, not a report.

Model descriptions

Describe the resulting models. Report on the interpretation of the models and document any difficulties encountered with their meanings.

4.4 Assess model

Task

Assess model

The data mining engineer interprets the models according to his domain knowledge, the data mining success criteria, and the desired test design. The data mining engineer judges the success of the application of modeling and discovery techniques technically; he contacts business analysts and domain experts later in order to discuss the data mining results in the business context. Please note that this task only considers models, whereas the evaluation phase also takes into account all other results that were produced in the course of the project.

The data mining engineer tries to rank the models. He assesses the models according to the evaluation criteria. As much as possible, he also takes into account business objectives and business success criteria. In most data mining projects, the data mining engineer applies a single technique more than once, or generates data mining results with several different techniques. In this task, he also compares all results according to the evaluation criteria.

Outputs

Model assessment

Summarize results of this task, list qualities of generated models (e.g., in terms of accuracy), and rank their quality in relation to each other.

Revised parameter settings

According to the model assessment, revise parameter settings and tune them for the next run in the Build Model task. Iterate model building and assessment until you strongly believe that you have found the *best* model(s). Document all such revisions and assessments.

5 Evaluation

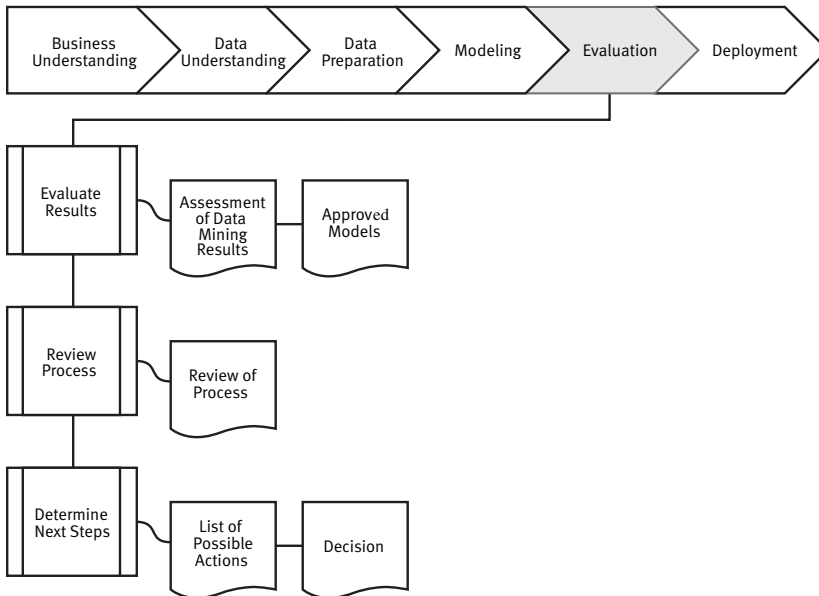


Figure 8: Evaluation

5.1 Evaluate results

Task

Evaluate results

Previous evaluation steps dealt with factors such as the accuracy and generality of the model. This step assesses the degree to which the model meets the business objectives and seeks to determine if there is some business reason why this model is deficient. Another option is to test the model(s) on test applications in the real application, if time and budget constraints permit.

Moreover, evaluation also assesses other data mining results generated. Data mining results involve models that are necessarily related to the original business objectives and all other findings that are not necessarily related to the original business objectives, but might also unveil additional challenges, information, or hints for future directions.

Outputs

Assessment of data mining results with respect to business success criteria

Summarize assessment results in terms of business success criteria, including a final statement regarding whether the project already meets the initial business objectives.

Approved models

After assessing models with respect to business success criteria, the generated models that meet the selected criteria become the approved models.

5.2 Review process

Task

Review process

At this point, the resulting models appear to be satisfactory and to satisfy business needs. It is now appropriate to do a more thorough review of the data mining engagement in order to determine if there is any important factor or task that has somehow been overlooked. This review also covers quality assurance issues—for example: Did we correctly build the model? Did we use only the attributes that we are allowed to use and that are available for future analyses?

Output

Review of process

Summarize the process review and highlight activities that have been missed and those that should be repeated.

5.3 Determine next steps

Task

Determine next steps

Depending on the results of the assessment and the process review, the project team decides how to proceed. The team decides whether to finish this project and move on to deployment, initiate further iterations, or set up new data mining projects. This task includes analyses of remaining resources and budget, which may influence the decisions.

Outputs

List of possible actions

List the potential further actions, along with the reasons for and against each option.

Decision

Describe the decision as to how to proceed, along with the rationale.

6 Deployment

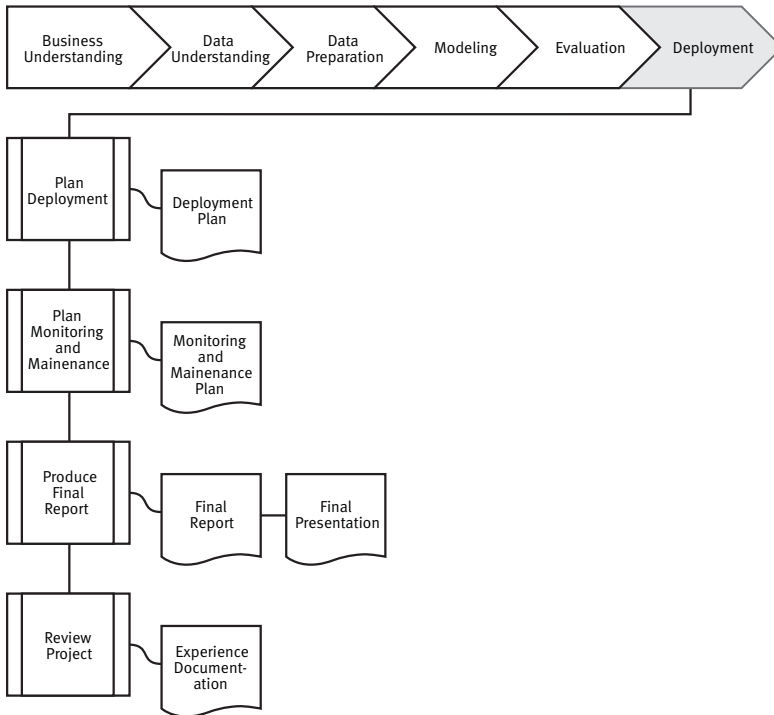


Figure 9: Deployment

6.1 Plan deployment

Task

Plan deployment

This task takes the evaluation results and determines a strategy for deployment. If a general procedure has been identified to create the relevant model(s), this procedure is documented here for later deployment.

Output

Deployment plan

Summarize the deployment strategy, including the necessary steps and how to perform them.

6.2 Plan monitoring and maintenance

Task **Plan monitoring and maintenance**

Monitoring and maintenance are important issues if the data mining result becomes part of the day-to-day business and its environment. The careful preparation of a maintenance strategy helps to avoid unnecessarily long periods of incorrect usage of data mining results. In order to monitor the deployment of the data mining result(s), the project needs a detailed monitoring process plan. This plan takes into account the specific type of deployment.

Output **Monitoring and maintenance plan**

Summarize the monitoring and maintenance strategy, including the necessary steps and how to perform them.

6.3 Produce final report

Task **Produce final report**

At the end of the project, the project team writes up a final report. Depending on the deployment plan, this report may be only a summary of the project and its experiences (if they have not already been documented as an ongoing activity) or it may be a final and comprehensive presentation of the data mining result(s).

Outputs **Final report**

This is the final written report of the data mining engagement. It includes all of the previous deliverables, summarizing and organizing the results.

Final presentation

There will also often be a meeting at the conclusion of the project at which the results are presented to the customer.

6.4 Review project

Task **Review project**

Assess what went right and what went wrong, what was done well and what needs to be improved.

Output **Experience documentation**

Summarize important experience gained during the project. For example, pitfalls, misleading approaches, or hints for selecting the best suited data mining techniques in similar situations could be part of this documentation. In ideal projects, experience documentation also covers any reports that have been written by individual project members during previous phases of the project.



III The CRISP-DM user guide

1 Business understanding

1.1 Determine business objectives

Task Determine business objectives

The first objective of the analyst is to thoroughly understand, from a business perspective, what the customer really wants to accomplish. Often the customer has many competing objectives and constraints that must be properly balanced. The analyst's goal is to uncover important factors at the beginning of the project that can influence the final outcome. A likely consequence of neglecting this step would be to expend a great deal of effort producing the correct answers to the wrong questions.

Output Background

Collate the information that is known about the organization's business situation at the start of the project. These details not only serve to more closely identify the business goals to be achieved but also serve to identify resources, both human and material, that may be used or needed during the course of the project.

Activities Organization

- Develop organizational charts identifying divisions, departments, and project groups. The chart should also identify managers' names and responsibilities
- Identify key persons in the business and their roles
- Identify an internal sponsor (financial sponsor and primary user/domain expert)
- Indicate if there is a steering committee and list members
- Identify the business units which are affected by the data mining project (e.g., Marketing, Sales, Finance)

Problem area

- Identify the problem area (e.g., marketing, customer care, business development, etc.)
- Describe the problem in general terms
- Check the current status of the project (e.g., Check if it is already clear within the business unit that a data mining project is to be performed, or whether data mining needs to be promoted as a key technology in the business)
- Clarify prerequisites of the project (e.g., What is the motivation of the project? Does the business already use data mining?)
- If necessary, prepare presentations and present data mining to the business

- Identify target groups for the project result (e.g., Are we expected to deliver a report for top management or an operational system to be used by naive end users?)
- Identify the users' needs and expectations

Current solution

- Describe any solution currently used to address the problem
- Describe the advantages and disadvantages of the current solution and the level to which it is accepted by the users

Output

Business objectives

Describe the customer's primary objective, from a business perspective. In addition to the primary business objective, there are typically a large number of related business questions that the customer would like to address. For example, the primary business goal might be to keep current customers by predicting when they are prone to move to a competitor, while a secondary business objective might be to determine whether lower fees affect only one particular segment of customers.

Activities

- Informally describe the problem to be solved
- Specify all business questions as precisely as possible
- Specify any other business requirements (e.g., the business does not want to lose any customers)
- Specify expected benefits in business terms

Beware!

- Beware of setting unattainable goals—make them as realistic as possible.

Output

Business success criteria

Describe the criteria for a successful or useful outcome to the project from the business point of view. This might be quite specific and readily measurable, such as reduction of customer churn to a certain level, or general and subjective, such as "give useful insights into the relationships." In the latter case, be sure to indicate who would make the subjective judgment.

Activities

- Specify business success criteria (e.g., Improve response rate in a mailing campaign by 10 percent and sign-up rate by 20 percent)
- Identify who assesses the success criteria

Remember! Each of the success criteria should relate to at least one of the specified business objectives.

Good idea! Before starting the situation assessment, you might analyze previous experiences of this problem—either internally, using CRISP-DM, or externally, using pre-packaged solutions.

1.2 Assess situation

Task **Assess situation**

This task involves more detailed fact-finding about all of the resources, constraints, assumptions, and other factors that should be considered in determining the data analysis goal and in developing the project plan.

Output **Inventory of resources**

List the resources available to the project, including personnel (business and data experts, technical support, data mining experts), data (fixed extracts, access to live warehoused or operational data), computing resources (hardware platforms), and software (data mining tools, other relevant software).

Activities **Hardware resources**

- Identify the base hardware
- Establish the availability of the base hardware for the data mining project
- Check if the hardware maintenance schedule conflicts with the availability of the hardware for the data mining project
- Identify the hardware available for the data mining tool to be used (if the tool is known at this stage)

Sources of data and knowledge

- Identify data sources
- Identify type of data sources (online sources, experts, written documentation, etc.)
- Identify knowledge sources
- Identify type of knowledge sources (online sources, experts, written documentation, etc.)
- Check available tools and techniques
- Describe the relevant background knowledge (informally or formally)

Personnel sources

- Identify project sponsor (if different from internal sponsor as in Section 1.1.1)
- Identify system administrator, database administrator, and technical support staff for further questions
- Identify market analysts, data mining experts, and statisticians, and check their availability
- Check availability of domain experts for later phases

Remember!

Remember that the project may need technical staff at odd times throughout the project, for example during data transformation.

Output**Requirements, assumptions, and constraints**

List all requirements of the project, including schedule of completion, comprehensibility, and quality of results and security, as well as legal issues. As part of this output, make sure that you are allowed to use the data.

List the assumptions made by the project. These may be assumptions about the data, which can be verified during data mining, but may also include non-verifiable assumptions related to the project. It is particularly important to list the latter if they will affect the validity of the results.

List the constraints made on the project. These constraints might involve lack of resources to carry out some of the tasks in the project in the time required, or there may be legal or ethical constraints on the use of the data or the solution needed to carry out the data mining task.

Activities**Requirements**

- Specify target group profile
- Capture all requirements on scheduling
- Capture requirements on comprehensibility, accuracy, deploy ability, maintainability, and repeatability of the data mining project and the resulting model(s)
- Capture requirements on security, legal restrictions, privacy, reporting, and project schedule

Assumptions

- Clarify all assumptions (including implicit ones) and make them explicit (e.g., to address the business question, a minimum number of customers with age above 50 is necessary)
- List assumptions on data quality (e.g., accuracy, availability)
- List assumptions on external factors (e.g., economic issues, competitive products, technical advances)
- Clarify assumptions that lead to any of the estimates (e.g., the price of a specific tool is assumed to be lower than \$1,000)
- List all assumptions regarding whether it is necessary to understand and describe or explain the model (e.g., how should the model and results be presented to senior management/sponsor)



Constraints

- Check general constraints (e.g., legal issues, budget, timescales, and resources)
- Check access rights to data sources (e.g., access restrictions, password required)
- Check technical accessibility of data (operating systems, data management system, file or database format)
- Check whether relevant knowledge is accessible
- Check budget constraints (fixed costs, implementation costs, etc.)

Remember!

The list of assumptions also includes assumptions at the beginning of the project, i.e., what the starting point of the project has been.

Output

Risks and contingencies

List the risks, that is, the events that might occur, impacting schedule, cost, or result. List the corresponding contingency plans: what action will be taken to avoid or minimize the impact or recover from the occurrence of the foreseen risks.

Activities

Identify risks

- Identify business risks (e.g., competitor comes up with better results first)
- Identify organizational risks (e.g., department requesting project doesn't have funding for the project)
- Identify financial risks (e.g., further funding depends on initial data mining results)
- Identify technical risks
- Identify risks that depend on data and data sources (e.g., poor quality and coverage)

Develop contingency plans

- Determine conditions under which each risk may occur
- Develop contingency plans

Output

Terminology

Compile a glossary of terminology relevant to the project. This should include at least two components:

- (1) A glossary of relevant business terminology, which forms part of the business understanding available to the project
- (2) A glossary of data mining terminology, illustrated with examples relevant to the business problem in question

Activities

- Check prior availability of glossaries; otherwise begin to draft glossaries
- Talk to domain experts to understand their terminology
- Become familiar with the business terminology

Output**Costs and benefits**

Prepare a cost-benefit analysis for the project, comparing the costs of the project with the potential benefits to the business if it is successful

Activities

- Estimate costs for data collection
- Estimate costs of developing and implementing a solution
- Identify benefits (e.g., improved customer satisfaction, ROI, and increase in revenue)
- Estimate operating costs

Good idea!

The comparison should be as specific as possible, as this enables a better business case to be made.

Beware!

Remember to identify hidden costs, such as repeated data extraction and preparation, changes in workflows, and time required for training.

1.3 Determine data mining goals

Task**Determine data mining goals**

A business goal states objectives in business terminology; a data mining goal states project objectives in technical terms. For example, the business goal might be, “Increase catalog sales to existing customers,” while a data mining goal might be, “Predict how many widgets a customer will buy, given their purchases over the past three years, relevant demographic information, and the price of the item.”

Output**Data mining goals**

Describe the intended outputs of the project that enable the achievement of the business objectives. Note that these are normally technical outputs.

Activities

- Translate the business questions to data mining goals (e.g., a marketing campaign requires segmentation of customers in order to decide whom to approach in this campaign; the level/size of the segments should be specified).
- Specify data mining problem type (e.g., classification, description, prediction, and clustering). For more details about data mining problem types, see Appendix 2.



| | |
|-------------------|---|
| Good idea! | It may be wise to re-define the problem. For example, modeling product retention rather than customer retention when targeting customer retention delivers results too late to affect the outcome. |
| Output | Data mining success criteria Define the criteria for a successful outcome to the project in technical terms, for example a certain level of predictive accuracy or a propensity-to-purchase profile with a given degree of “lift.” As with business success criteria, it may be necessary to describe these in subjective terms, in which case the person or persons making the subjective judgment should be identified. |
| Activities | <ul style="list-style-type: none">■ Specify criteria for model assessment (e.g., model accuracy, performance and complexity)■ Define benchmarks for evaluation criteria■ Specify criteria which address subjective assessment criteria (e.g., model explain ability and data and marketing insight provided by the model) |
| Beware! | Remember that the data mining success criteria are different than the business success criteria defined earlier. Remember it is wise to plan for deployment from the start of the project. |

1.4 Produce project plan

| | |
|---------------|---|
| Task | Produce project plan Describe the intended plan for achieving the data mining goals and thereby achieving the business goals. |
| Output | Project plan List the stages to be executed in the project, together with their duration, resources required, inputs, outputs, and dependencies. Wherever possible, make explicit the large-scale iterations in the data mining process—for example, repetitions of the modeling and evaluation phases. As part of the project plan, it is also important to analyze dependencies between time schedule and risks. Mark results of these analyses explicitly in the project plan, ideally with actions and recommendations for actions if the risks are manifested. Although this is the only task in which the project plan is directly named, it nevertheless should be consulted continually and reviewed throughout the project. The project plan should be consulted at minimum whenever a new task is started or a further iteration of a task or activity is begun. |

Activities

- Define the initial process plan and discuss the feasibility with all involved personnel
- Combine all identified goals and selected techniques in a coherent procedure that solves the business questions and meets the business success criteria
- Estimate the effort and resources needed to achieve and deploy the solution. (It is useful to consider other people's experience when estimating timescales for data mining projects. For example, it is often postulated that 50-70 percent of the time and effort in a data mining project is used in the Data Preparation Phase and 20-30 percent in the Data Understanding Phase, while only 10-20 percent is spent in each of the Modeling, Evaluation, and Business Understanding Phases and 5-10 percent in the Deployment Phase.)
- Identify critical steps
- Mark decision points
- Mark review points
- Identify major iterations

Output

Initial assessment of tools and techniques

At the end of the first phase, the project team performs an initial assessment of tools and techniques. Here, it is important to select a data mining tool that supports various methods for different stages of the process, since the selection of tools and techniques may influence the entire project.

Activities

- Create a list of selection criteria for tools and techniques (or use an existing one if available)
- Choose potential tools and techniques
- Evaluate appropriateness of techniques
- Review and prioritize applicable techniques according to the evaluation of alternative solutions

2 Data understanding

2.1 Collect initial data

Task

Collect initial data

Acquire the data (or access to the data) listed in the project resources. This initial collection includes data loading, if necessary for data understanding. For example, if you intend to use a specific tool for data understanding, it is logical to load your data into this tool.

Output

Initial data collection report

Describe all the various data used for the project, and include any selection requirements for more detailed data. The data collection report should also define whether some attributes are relatively more important than others.

Remember that any assessment of data quality should be made not just of the individual data sources but also of any data that results from merging data sources. Because of inconsistencies between the sources, merged data may present problems that do not exist in the individual data sources.

Activities

Data requirements planning

- Plan which information is needed (e.g., only for given attributes, or specific additional information)
- Check if all the information needed (to solve the data mining goals) is actually available

Selection criteria

- Specify selection criteria (e.g., Which attributes are necessary for the specified data mining goals? Which attributes have been identified as being irrelevant? How many attributes can we handle with the chosen techniques?)
- Select tables/files of interest
- Select data within a table/file
- Think about how long a history one should use (e.g., even if 18 months of data are available, only 12 months may be needed for the exercise)

Beware!

Be aware that data collected from different sources may give rise to quality problems when merged (e.g., address files merged with a customer database may show inconsistencies of format, invalidity of data, etc.).

Insertion of data

- If the data contain free text entries, do we need to encode them for modeling or do we want to group specific entries?
- How can missing attributes be acquired?
- How can we best extract the data?

Good idea!

Remember that some knowledge about the data may be available from non-electronic sources (e.g., from people, printed text, etc.).

Remember that it may be necessary to preprocess the data (time-series data, weighted averages, etc.).

2.2 Describe data

Task**Describe data**

Examine the “gross” properties of the acquired data and report on the results.

Output**Data description report**

Describe the data that has been acquired, including the format of the data, the quantity of the data (e.g., the number of records and fields within each table), the identities of the fields, and any other surface features that have been discovered.

Activities**Volumetric analysis of data**

- Identify data and method of capture
- Access data sources
- Use statistical analyses if appropriate
- Report tables and their relations
- Check data volume, number of multiples, complexity
- Note if the data contain free text entries

Attribute types and values

- Check accessibility and availability of attributes
- Check attribute types (numeric, symbolic, taxonomy, etc.)
- Check attribute value ranges
- Analyze attribute correlations
- Understand the meaning of each attribute and attribute value in business terms
- For each attribute, compute basic statistics (e.g., compute distribution, average, max, min, standard deviation, variance, mode, skewness, etc.)
- Analyze basic statistics and relate the results to their meaning in business terms
- Decide if the attribute is relevant for the specific data mining goal



- Determine if the attribute meaning is used consistently
- Interview domain experts to obtain their opinion of attribute relevance
- Decide if it is necessary to balance the data (based on the modeling techniques to be used)

Keys

- Analyze key relationships
- Check amount of overlaps of key attribute values across tables

Review assumptions/goals

- Update list of assumptions, if necessary

2.3 Explore data

Task

Explore data

This task tackles the data mining questions that can be addressed using querying, visualization, and reporting techniques. These analyses may directly address the data mining goals. However, they may also contribute to or refine the data description and quality reports, and feed into the transformation and other data preparation steps needed before further analysis can occur.

Output

Data exploration report

Describe the results of this task, including first findings or initial hypotheses and their impact on the remainder of the project. The report may also include graphs and plots that indicate data characteristics or point to interesting data subsets worthy of further examination.

Activities

Data exploration

- Analyze properties of interesting attributes in detail (e.g., basic statistics, interesting sub-populations)
- Identify characteristics of sub-populations

Form suppositions for future analysis

- Consider and evaluate information and findings in the data descriptions report
- Form a hypothesis and identify actions
- Transform the hypothesis into a data mining goal, if possible
- Clarify data mining goals or make them more precise. A “blind” search is not necessarily useless, but a more directed search toward business objectives is preferable.
- Perform basic analysis to verify the hypothesis

2.4 Verify data quality

Task

Verify data quality

Examine the quality of the data, addressing questions such as: Is the data complete (does it cover all the cases required)? Is it correct or does it contain errors? If there are errors, how common are they? Are there missing values in the data? If so, how are they represented, where do they occur, and how common are they?

Output

Data quality report

List the results of the data quality verification; if there are quality problems, list possible solutions.

Activities

- Identify special values and catalog their meaning

Review keys, attributes

- Check coverage (e.g., whether all possible values are represented)
- Check keys
- Verify that the meanings of attributes and contained values fit together
- Identify missing attributes and blank fields
- Establish the meaning of missing data
- Check for attributes with different values that have similar meanings (e.g., low fat, diet)
- Check spelling and format of values (e.g., same value but sometimes beginning with a lower-case letter, sometimes with an upper-case letter)
- Check for deviations, and decide whether a deviation is “noise” or may indicate an interesting phenomenon
- Check for plausibility of values, (e.g., all fields having the same or nearly the same values)

Good idea!

Review any attributes that give answers that conflict with common sense (e.g., teenagers with high income levels).

Use visualization plots, histograms, etc. to reveal inconsistencies in the data.

Data quality in flat files

- If data are stored in flat files, check which delimiter is used and whether it is used consistently within all attributes
- If data are stored in flat files, check the number of fields in each record to see if they coincide



Noise and inconsistencies between sources

- Check consistencies and redundancies between different sources
- Plan for dealing with noise
- Detect the type of noise and which attributes are affected

Good idea!

Remember that it may be necessary to exclude some data since they do not exhibit either positive or negative behavior (e.g., to check on customers' loan behavior, exclude all those who have never borrowed, do not finance a home mortgage, those whose mortgage is nearing maturity, etc.).

Review whether assumptions are valid or not, given the current information on data and business knowledge.

3 Data preparation

Output

Dataset

These are the dataset(s) produced by the data preparation phase, used for modeling or for the major analysis work of the project.

Output

Dataset description

This is the description of the dataset(s) used for the modeling or for the major analysis work of the project.

3.1 Select data

Task

Select data

Decide on the data to be used for analysis. Criteria include relevance to the data mining goals, quality, and technical constraints such as limits on data volume or data types.

Output

Rationale for inclusion/exclusion

List the data to be used/excluded and the reasons for these decisions.

Activities

- Collect appropriate additional data (from different sources—in-house as well as externally)
- Perform significance and correlation tests to decide if fields should be included
- Reconsider Data Selection Criteria (See Task 2.1) in light of experiences of data quality and data exploration (i.e., may wish include/exclude other sets of data)
- Reconsider Data Selection Criteria (See Task 2.1) in light of experience of modeling (i.e., model assessment may show that other datasets are needed)
- Select different data subsets (e.g., different attributes, only data which meet certain conditions)

- Consider the use of sampling techniques (e.g., A quick solution may involve splitting test and training datasets or reducing the size of the test dataset, if the tool cannot handle the full dataset. It may also be useful to have weighted samples to give different importance to different attributes or different values of the same attribute.)
- Document the rationale for inclusion/exclusion
- Check available techniques for sampling data

Good idea!

Based on Data Selection Criteria, decide if one or more attributes are more important than others and weight the attributes accordingly. Decide, based on the context (i.e., application, tool, etc.), how to handle the weighting.

3.2 Clean data

Task

Clean data

Raise the data quality to the level required by the selected analysis techniques. This may involve the selection of clean subsets of the data, the insertion of suitable defaults, or more ambitious techniques such as the estimation of missing data by modeling.


Output

Data cleaning report

Describe the decisions and actions that were taken to address the data quality problems reported during the Verify Data Quality Task. If the data are to be used in the data mining exercise, the report should address outstanding data quality issues and what possible effect this could have on the results.

Activities

- Reconsider how to deal with any observed type of noise
- Correct, remove, or ignore noise
- Decide how to deal with special values and their meaning. The area of special values can give rise to many strange results and should be carefully examined. Examples of special values could arise through taking results of a survey where some questions were not asked or not answered. This might result in a value of 99 for unknown data. For example, 99 for marital status or political affiliation. Special values could also arise when data is truncated—e.g., 00 for 100-year-old people or all cars with 100,000 km on the odometer.
- Reconsider Data Selection Criteria (See Task 2.1) in light of experiences of data cleaning (i.e., you may wish to include/exclude other sets of data).



Good idea! Remember that some fields may be irrelevant to the data mining goals and, therefore, noise in those fields has no significance. However, if noise is ignored for these reasons, it should be fully documented as the circumstances may change later.

3.3 Construct data

Task

Construct data

This task includes constructive data preparation operations such as the production of derived attributes, complete new records, or transformed values for existing attributes.

Activities

- Check available construction mechanisms with the list of tools suggested for the project
- Decide whether it is best to perform the construction inside the tool or outside (i.e., which is more efficient, exact, repeatable)
- Reconsider Data Selection Criteria (See Task 2.1) in light of experiences of data construction (i.e., you may wish include/exclude other sets of data)

Output

Derived attributes

Derived attributes are new attributes that are constructed from one or more existing attributes in the same record. An example might be: $\text{area} = \text{length} * \text{width}$.

Why should we need to construct derived attributes during the course of a data mining investigation? It should not be thought that only data from databases or other sources should be used in constructing a model. Derived attributes might be constructed because:

- Background knowledge convinces us that some fact is important and ought to be represented although we have no attribute currently to represent it
- The modeling algorithm in use handles only certain types of data—for example we are using linear regression and we suspect that there are certain non-linearities that will be not be included in the model
- The outcome of the modeling phase suggests that certain facts are not being covered

Activities

Derived attributes

- Decide if any attribute should be normalized (e.g., when using a clustering algorithm with age and income, in certain currencies, the income will dominate)
- Consider adding new information on the relevant importance of attributes by adding new attributes (for example, attribute weights, weighted normalization)

- How can missing attributes be constructed or imputed? [Decide type of construction (e.g., aggregate, average, induction).]
- Add new attributes to the accessed data

Good idea!

Before adding Derived Attributes, try to determine if and how they ease the model process or facilitate the modeling algorithm. Perhaps “income per person” is a better/easier attribute to use than “income per household.” Do not derive attributes simply to reduce the number of input attributes.

Another type of derived attribute is the single-attribute transformation, usually performed to fit the needs of the modeling tools.

Activities

Single-attribute transformations

- Specify necessary transformation steps in terms of available transformation facilities (for example, change a binning of a numeric attribute)
- Perform transformation steps

Good idea!

Transformations may be necessary to change ranges to symbolic fields (e.g., ages to age ranges) or symbolic fields (“definitely yes,” “yes,” “don’t know,” “no”) to numeric values. Modeling tools or algorithms often require them.

Output

Generated records

Generated records are completely new records, which add new knowledge or represent new data that is not otherwise represented (e.g., having segmented the data, it may be useful to generate a record to represent the prototypical member of each segment for further processing).

Activities

Check for available techniques if needed (e.g., mechanisms to construct prototypes for each segment of segmented data).



3.4 *Integrate data*

Task

Integrate data

These are methods for combining information from multiple tables or other information sources to create new records or values.

Output

Merged data

Merging tables refers to joining together two or more tables that have different information about the same objects. At this stage, it may also be advisable to generate new records. It may also be recommended to generate aggregate values.

Aggregation refers to operations where new values are computed by summarizing information from multiple records and/or tables.

Activities

- Check if integration facilities are able to integrate the input sources as required
- Integrate sources and store results
- Reconsider Data Selection Criteria (See Task 2.1) in light of experiences of data integration (i.e., you may wish to include/exclude other sets of data)

Good idea!

Remember that some knowledge may be contained in non-electronic format.

3.5 *Format data*

Task

Format data

Formatting transformations refers primarily to syntactic modifications made to the data that do not change its meaning, but might be required by the modeling tool.

Output

Reformatted data

Some tools have requirements on the order of the attributes, such as the first field being a unique identifier for each record or the last field being the outcome field the model is to predict.

Activities

Rearranging attributes

Some tools have requirements on the order of the attributes, such as the first field being a unique identifier for each record or the last field being the outcome field the model is to predict.

Reordering records

It might be important to change the order of the records in the dataset. Perhaps the modeling tool requires that the records be sorted according to the value of the outcome attribute.

Reformatted within-value

- These are purely syntactic changes made to satisfy the requirements of the specific modeling tool
- Reconsider Data Selection Criteria (See Task 2.1) in light of experiences of data cleaning (i.e., you may wish to include/exclude other sets of data)

4 Modeling

4.1 Select modeling technique

Task

Select modeling technique

As the first step in modeling, select the actual initial modeling technique. If multiple techniques are to be applied, perform this task separately for each technique.

Remember that not all tools and techniques are applicable to each and every task. For certain problems, only some techniques are appropriate (See Appendix 2, where techniques appropriate for certain data mining problem types are discussed in more detail). “Political requirements” and other constraints further limit the choices available to the data mining engineer. It may be that only one tool or technique is available to solve the problem at hand—and that the tool may not be absolutely the best, from a technical standpoint.

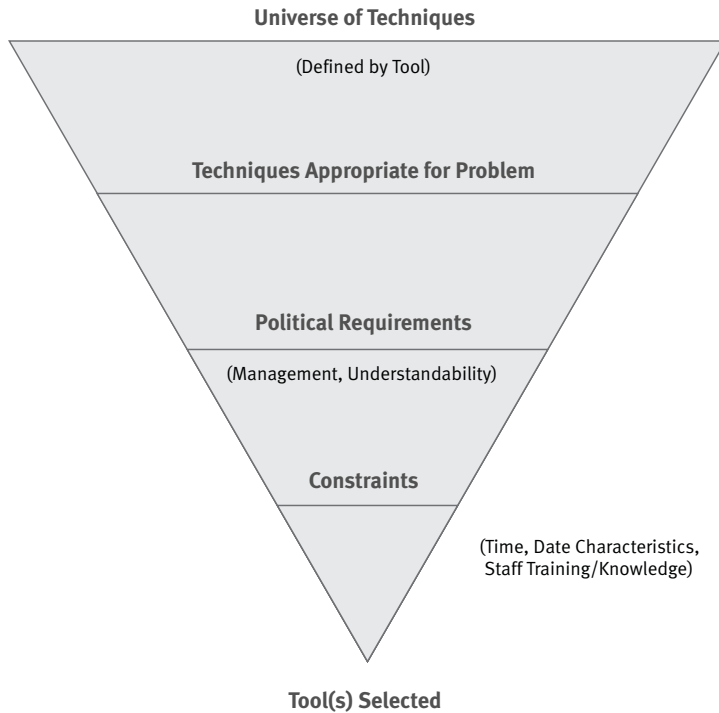


Figure 10:
Universe of Techniques

| | |
|-------------------|--|
| Output | <p>Modeling technique</p> <p>Record the actual modeling technique that is used.</p> |
| Activities | Decide on appropriate technique for exercise, bearing in mind the tool selected. |
| Output | <p>Modeling assumptions</p> <p>Many modeling techniques make specific assumptions about the data.</p> |
| Activities | <ul style="list-style-type: none"> ■ Define any built-in assumptions made by the technique about the data (e.g., quality, format, distribution) ■ Compare these assumptions with those in the Data Description Report ■ Make sure that these assumptions hold and go back to the Data Preparation Phase, if necessary |

4.2 *Generate test design*

Task **Generate test design**

Prior to building a model, it is necessary to define a procedure to test the model's quality and validity. For example, in supervised data mining tasks such as classification, it is common to use error rates as quality measures for data mining models. Therefore, the test design specifies that the dataset should be separated into training and test sets. The model is built on the training set and its quality estimated on the test set.

Output **Test design**

Describe the intended plan for training, testing, and evaluating the models. A primary component of the plan is to decide how to divide the available dataset into training data, test data, and validation test sets.

- Activities**
- Check existing test designs for each data mining goal separately
 - Decide on necessary steps (number of iterations, number of folds, etc.)
 - Prepare data required for test

4.3 *Build model*

Task **Build model**

Run the modeling tool on the prepared dataset to create one or more models.

Output **Parameter settings**

With any modeling tool, there are often a large number of parameters that can be adjusted. List the parameters and their chosen values, along with the rationale for the choice.

- Activities**
- Set initial parameters
 - Document reasons for choosing those values

Output **Models**

Run the modeling tool on the prepared dataset to create one or more models.

- Activities**
- Run the selected technique on the input dataset to produce the model
 - Post-process data mining results (e.g., edit rules, display trees)

Output

Model description

Describe the resulting model and assess its expected accuracy, robustness, and possible shortcomings. Report on the interpretation of the models and any difficulties encountered.

Activities

- Describe any characteristics of the current model that may be useful for the future
- Record parameter settings used to produce the model
- Give a detailed description of the model and any special features
- For rule-based models, list the rules produced, plus any assessment of per-rule or overall model accuracy and coverage
- For opaque models, list any technical information about the model (such as neural network topology) and any behavioral descriptions produced by the modeling process (such as accuracy or sensitivity)
- Describe the model's behavior and interpretation
- State conclusions regarding patterns in the data (if any); sometimes the model reveals important facts about the data without a separate assessment process (e.g., that the output or conclusion is duplicated in one of the inputs)

4.4 Assess model

Task

Assess model

The model should now be assessed to ensure that it meets the data mining success criteria and passes the desired test criteria. This is a purely technical assessment based on the outcome of the modeling tasks.

Output

Model assessment

Summarize results of this task, list qualities of generated models (e.g., in terms of accuracy), and rank their quality in relation to each other.

Activities

- Evaluate results with respect to evaluation criteria
- Test result according to a test strategy (e.g.: Train and Test, Cross-validation, bootstrapping, etc.)
- Compare evaluation results and interpretation
- Create ranking of results with respect to success and evaluation criteria
- Select best models
- Interpret results in business terms (as far as possible at this stage)
- Get comments on models by domain or data experts
- Check plausibility of model

- Check effect on data mining goal
- Check model against given knowledge base to see if the discovered information is novel and useful
- Check reliability of result
- Analyze potential for deployment of each result
- If there is a verbal description of the generated model (e.g., via rules), assess the rules: Are they logical, are they feasible, are there too many or too few, do they offend common sense?
- Assess results
- Get insights into why a certain modeling technique and certain parameter settings lead to good/bad results

Good idea! “Lift Tables” and “Gain Tables” can be constructed to determine how well the model is predicting.

Output **Revised parameter settings**
According to the model assessment, revise parameter settings and tune them for the next run in the Build Model task. Iterate model building and assessment until you find the best model.

Activities Adjust parameters to produce better models.


5 Evaluation

Previous evaluation steps dealt with factors such as the accuracy and generality of the model. This step assesses the degree to which the model meets the business objectives, and seeks to determine if there is some business reason why this model is deficient. It compares results with the evaluation criteria defined at the start of the project.

A good way of defining the total outputs of a data mining project is to use the equation:

$$\text{RESULTS} = \text{MODELS} + \text{FINDINGS}$$

In this equation, we are defining that the total output of the data mining project is not just the models (although they are, of course, important) but also the findings, which we define as anything (apart from the model) that is important in meeting the objectives of the business or important in leading to new questions, lines of approach, or side effects (e.g., data quality problems uncovered by the data mining exercise). Note: Although the model is directly connected to the business questions, the findings need not be related to any questions or objectives, as long as they are important to the initiator of the project.



5.1 Evaluate results

Task

Evaluate results

This step assesses the degree to which the model meets the business objectives, and seeks to determine if there is some business reason why this model is deficient. Another option is to test the model(s) on test applications in the real application, if time and budget constraints permit.

Moreover, evaluation also assesses other generated data mining results. Data mining results cover models that are related to the original business objectives and all other findings. Some are related to the original business objectives while others might unveil additional challenges, information, or hints for future directions.

Output

Assessment of data mining results with respect to business success criteria

Summarize assessment results in terms of business success criteria, including a final statement related to whether the project already meets the initial business objectives.

Activities

- Understand the data mining results
- Interpret the results in terms of the application
- Check effect on for data mining goal
- Check the data mining result against the given knowledge base to see if the discovered information is novel and useful
- Evaluate and assess results with respect to business success criteria (i.e., has the project achieved the original Business Objectives)
- Compare evaluation results and interpretation
- Rank results with respect to business success criteria
- Check effect of result on initial application goal
- Determine if there are new business objectives to be addressed later in the project, or in new projects
- State recommendations for future data mining projects

Output

Approved models

After accessing models with respect to business success criteria, select and approve the generated models that meet the selected criteria.

5.2 Review process

Task

Review process

At this point, the resulting model appears to be satisfactory and appears to satisfy business needs. It is now appropriate to make a more thorough review of the data mining engagement in order to determine if there is any important factor or task that has somehow been overlooked. At this stage of the data mining exercise, the Process Review takes the form of a Quality Assurance Review.

Output

Review of process

Summarize the process review and list activities that have been missed and/or should be repeated.

Activities

- Provide an overview of the data mining process used
- Analyze the data mining process. For each stage of the process ask:
 - Was it necessary?
 - Was it executed optimally?
 - In what ways could it be improved?
- Identify failures
- Identify misleading steps
- Identify possible alternative actions and/or unexpected paths in the process
- Review data mining results with respect to business success criteria

5.3 Determine next steps

Task

Determine next steps

Based on the assessment results and the process review, the project team decides how to proceed. Decisions to be made include whether to finish this project and move on to deployment, to initiate further iterations, or to set up new data mining projects.

Output

List of possible actions

List possible further actions along with the reasons for and against each option.

- Activities**
- Analyze the potential for deployment of each result
 - Estimate potential for improvement of current process
 - Check remaining resources to determine if they allow additional process iterations (or whether additional resources can be made available)
 - Recommend alternative continuations
 - Refine process plan

Output **Decision**
Describe the decisions made, along with the rationale for them.

- Activities**
- Rank the possible actions
 - Select one of the possible actions
 - Document reasons for the choice

6 Deployment

6.1 Plan deployment

Task **Plan deployment**
This task starts with the evaluation results and concludes with a strategy for deployment of the data mining result(s) into the business.

Output **Deployment plan**
Summarize the deployment strategy, including necessary steps and how to perform them.

- Activities**
- Summarize deployable results
 - Develop and evaluate alternative plans for deployment
 - Decide for each distinct knowledge or information result
 - Determine how knowledge or information will be propagated to users
 - Decide how the use of the result will be monitored and its benefits measured (where applicable)
 - Decide for each deployable model or software result
 - Establish how the model or software result will be deployed within the organization's systems
 - Determine how its use will be monitored and its benefits measured (where applicable)
 - Identify possible problems during deployment (pitfalls to be avoided)

6.2 Plan monitoring and maintenance

Task **Plan monitoring and maintenance**

Monitoring and maintenance are important issues if the data mining results become part of the day-to-day business and its environment. A careful preparation of a maintenance strategy helps to avoid unnecessarily long periods of incorrect usage of data mining results. In order to monitor the deployment of the data mining result(s), the project needs a detailed plan for monitoring and maintenance. This plan takes into account the specific type of deployment.

Output **Monitoring and maintenance plan**

Summarize monitoring and maintenance strategy, including necessary steps and how to perform them.

- ### **Activities**
- Check for dynamic aspects (i.e., what things could change in the environment?)
 - Decide how accuracy will be monitored
 - Determine when the data mining result or model should not be used any more. Identify criteria (validity, threshold of accuracy, new data, change in the application domain, etc.), and what should happen if the model or result could no longer be used. (update model, set up new data mining project, etc.).
 - Will the business objectives of the use of the model change over time? Fully document the initial problem the model was attempting to solve.
 - Develop monitoring and maintenance plan.

6.3 Produce final report

Task **Produce final report**

At the end of the project, the project team writes up a final report. Depending on the deployment plan, this report may be only a summary of the project and its experience, or a final presentation of the data mining result(s).

Output **Final report**

At the end of the project, there will be at least one final report in which all the threads are brought together. As well as identifying the results obtained, the report should also describe the process, show which costs have been incurred, define any deviations from the original plan, describe implementation plans, and make any recommendations for future work. The actual detailed content of the report depends very much on the intended audience.

Activities

- Identify what reports are needed (slide presentation, management summary, detailed findings, explanation of models, etc.)
- Analyze how well initial data mining goals have been met
- Identify target groups for report
- Outline structure and contents of report(s)
- Select findings to be included in the reports
- Write a report

Output

Final presentation

As well as a final report, it may be necessary to make a final presentation to summarize the project—maybe to the management sponsor, for example. The presentation normally contains a subset of the information contained in the final report, structured in a different way.

Activities

- Decide on target group for the final presentation and determine if they will already have received the final report
- Select which items from the final report should be included in final presentation

6.4 Review project

Task

Review project

Assess what went right and what went wrong, what was done well, and what needs to be improved.

Output

Experience documentation

Summarize important experience gained during the project. For example, pitfalls, misleading approaches, or tips for selecting the best-suited data mining techniques in similar situations could be part of this documentation. In ideal projects, experience documentation also covers any reports that have been written by individual project members during the project.

Activities

- Interview all significant people involved in the project and ask them about their experience during the project
- If end users in the business work with the data mining result(s), interview them: Are they satisfied? What could have been done better? Do they need additional support?
- Summarize feedback and write the experience documentation
- Analyze the process (things that worked well, mistakes made, lessons learned, etc.)
- Document the specific data mining process (How can the results and the experience of applying the model be fed back into the process?)
- Generalize from the details to make the experience useful for future projects

IV The CRISP-DM outputs

This section contains brief descriptions of the purpose and the contents of the most important reports. Here, we focus on reports that are meant to communicate the results of a phase to people not involved in this phase (and possibly not involved in this project). These are not necessarily identical to the outputs as described in the reference model and the user guide. The purpose of these outputs is mostly to document results while performing the project.

1 Business understanding

The results of the Business Understanding phase can be summarized in one report. We suggest the following sections:

Background

The Background section provides a basic overview of the project context. This lists what area the project is working in, what problems have been identified, and why data mining appears to provide a solution.

Business objectives and success criteria

The Business Objectives section describes the goals of the project in business terms. For each objective, Business Success Criteria, i.e., explicit measures for determining whether or not the project succeeded in its objectives, should be provided. This section should also list objectives that were considered but rejected. The rationale of the selection of objectives should be given.

Inventory of resources

The Inventory of Resources section aims to identify personnel, data sources, technical facilities, and other resources that may be useful in carrying out the project.

Requirements, assumptions, and constraints

This section lists general requirements for the project's execution: type of project results, assumptions made about the nature of the problem and the data being used, and constraints imposed on the project.

Risks and contingencies

This section identifies problems that may occur in the project, describes the consequences, and states what actions can be taken to minimize such risks.

Terminology

The Terminology section allows people unfamiliar with the problems being addressed by the project to become more familiar with them.



Costs and benefits

This section describes the costs of the project and predicted business benefits if the project is successful (e.g., return on investment). Other less tangible benefits (e.g., customer satisfaction) should also be highlighted.

Data mining goals and success criteria

The Data Mining Goals section states the results of the project that enable the achievement of the business objectives. As well as listing the probable data mining approaches, the success criteria for the results in data mining terms, should also be listed.

Project plan

This section lists the stages to be executed in the project, together with their duration, resources required, inputs, outputs, and dependencies. Where possible, it should make explicit the large-scale iterations in the data mining process—for example, repetitions of the modeling and evaluation phases.

Initial assessment of tools and techniques

This section gives an initial view of what tools and techniques are likely to be used and how. It describes the requirements for tools and techniques, lists available tools and techniques, and matches them to requirements.

2 Data understanding

The results of the Data Understanding phase are usually documented in several reports. Ideally, these reports should be written while performing the respective tasks. The reports describe the datasets that are explored during data understanding. For the final report, a summary of the most relevant parts is sufficient.

Initial data collection report

This report describes how the different data sources identified in the inventory were captured and extracted.

Topics to be covered:

- Background of data
- List of data sources with broad area of required data covered by each
- For each data source, method of acquisition or extraction
- Problems encountered in data acquisition or extraction

Data description report

Each dataset acquired is described in this report.

Topics to be covered:

- Each data source described in detail
- List of tables (may be only one) or other database objects
- Description of each field, including units, codes used, etc.

Data exploration report

This report describes the data exploration and its results.

Topics to be covered:

- Background, including the broad goals of data exploration. For each area of exploration undertaken:
 - Expected regularities or patterns
 - Method of detection
 - Regularities or patterns found, expected and unexpected
 - Any other surprises
 - Conclusions for data transformation, data cleaning, and any other pre-processing
 - Conclusions related to data mining goals or business objectives
 - Summary of conclusions

Data quality report

This report describes the completeness and accuracy of the data.

Topics to be covered:

- Background, including broad expectations about data quality. For each dataset:
 - Approach taken to assess data quality
 - Results of data quality assessment
 - Summary of data quality conclusions



3 Data preparation

The reports in the data preparation phase focus on the pre-processing steps that produce the data to be mined.

Dataset description report

This report provides a description of the dataset (after pre-processing) and the process by which it was produced.

Topics to be covered:

- Background, including broad goals and plan for pre-processing
- Rationale for inclusion/exclusion of datasets. For each included dataset:
 - Description of the pre-processing, including the actions that were necessary to address any data quality issues
 - Detailed description of the resulting dataset, table by table and field by field
 - Rationale for inclusion/exclusion of attributes
 - Discoveries made during pre-processing, and any implications for further work
 - Summary and conclusions

4 Modeling

The outputs produced during the Modeling phase can be combined into one report. We suggest the following sections:

Modeling assumptions

This section defines any explicit assumptions made about the data and any assumptions that are implicit in the modeling technique to be used.

Test design

This section describes how the models are built, tested, and evaluated.

Topics to be covered:

- Background—outlines the modeling undertaken and its relationship to the data mining goals. For each modeling task:
 - Broad description of the type of model and the training data to be used
 - Explanation of how the model will be tested or assessed
 - Description of any data required for testing
 - Plan for production of test data if any
 - Description of any planned examination of models by domain or data experts
 - Summary of test plan

Model description

This report describes the delivered models and overviews the process by which they were produced.

Topics to be covered:

- Overview of models produced. For each model:
 - Type of model and relationship to data mining goals
 - Parameter settings used to produce the model
 - Detailed description of the model and any special features. For example:
 - For rule-based models, list the rules produced plus any assessment of per-rule or overall model accuracy and coverage
 - For opaque models, list any technical information about the model (such as neural network topology) and any behavioral descriptions produced by the modeling process (such as accuracy or sensitivity)
 - Description of the model's behavior and interpretation
 - Conclusions regarding patterns in the data (if any). Sometimes the model will reveal important facts about the data without a separate assessment process (e.g., that the output or conclusion is duplicated in one of the inputs).
- Summary of conclusions

Model assessment

This section describes the results of testing the models according to the test design.

Topics to be covered:

- Overview of assessment process and results, including any deviations from the test plan. For each model:
 - Detailed assessment, including measurements such as accuracy and interpretation of behavior
 - Any comments on models by domain or data experts
 - Summary assessment of models
 - Insights into why a certain modeling technique and certain parameter settings led to good/bad results
 - Summary assessment of complete model set



5 Evaluation

Assessment of data mining results with respect to business success criteria

This report compares the data mining results with the business objectives and the business success criteria.

Topics to be covered:

- Review of business objectives and business success criteria (which may have changed during and/or as a result of data mining). For each business success criterion:
 - Detailed comparison between success criterion and data mining results
 - Conclusions about achievability of success criterion and suitability of data mining process
- Review of project success:
 - Has the project achieved the original business objectives?
 - Are there new business objectives to be addressed later in the project or in new projects?
 - Conclusions for future data mining projects

Review of process

This section assesses the effectiveness of the project and identifies any factors that may have been overlooked that should be taken into consideration if the project is repeated.

List of possible actions

This section makes recommendations regarding the next steps in the project.

6 Deployment

Deployment plan

This report specifies the deployment of the data mining results.

Topics to be covered:

- Summary of deployable results (derived from Next Steps report)
- Description of deployment plan

Monitoring and maintenance plan

The monitoring and maintenance plan specifies how the deployed results are to be maintained. Topics to be covered:

- Overview of results deployment and indication of which results may require updating (and why). For each deployed result:
 - Description of how updating will be triggered (regular updates, trigger event, performance monitoring)
 - Description of how updating will be performed
- Summary of the results updating process

Final report

The final report is used to summarize the project and its results.

Contents:

- Summary of business understanding: background, objectives, and success criteria
- Summary of data mining process
- Summary of data mining results
- Summary of results evaluation
- Summary of deployment and maintenance plans
- Cost/benefit analysis
- Conclusions for the business
- Conclusions for future data mining

7 Summary of dependencies

The following table summarizes the main inputs to the deliverables. This does not mean that only the inputs listed should be considered—for example, the business objectives should be pervasive to all deliverables. However, the deliverables should address specific issues raised by their inputs.

| Phase | Deliverable | Refers To | Closely Related To |
|------------------------|---|---|-------------------------|
| Business Understanding | Background | | |
| | Business Objectives | Background | Terminology |
| | Business Success Criteria | Business Objectives | |
| | Inventory of Resources | | |
| | Requirements, Assumptions & Constraints | Business Objectives | |
| | Risks & Contingencies | Business Objectives; Business Success Criteria | |
| | Terminology | Background | Business Objectives |
| | Costs & Benefits | Business Objectives | Project Plan |
| | Data Mining Goals | Business Objectives; Requirements, Assumptions & Constraints | |
| | Data Mining Success Criteria | Business Success Criteria; Requirements; Assumptions & Constraints; Data Mining Goals | |
| | Project Plan | Business Objectives; Inventory of Resources; Requirements; Assumptions & Constraints; Risks & Contingencies | Costs & Benefits |
| Data Understanding | Initial Data Collection Report | Business Goals; Inventory of Resources; Data Mining Goals | |
| | Data Description Report | Business Goals; Initial Data Collection Report | Data Quality Report |
| | Data Quality Report | Business Goals; Initial Data Collection Report | Data Description Report |
| | Exploratory Analysis Report | Business Goals; Initial Data Collection Report | |
| Data Preparation | Dataset & Dataset Description | Business Goals; Data Mining Goals & Data Description Report; Data Quality Report; Exploratory Analysis Report | |
| Modeling | Test Design | Data Mining Goals; Data Mining Success Criteria | |
| | Models | Data Mining Goals | Parameter Settings |
| | Parameter Settings | Data Mining Goals | Models |
| | Model Description | Models; Parameter Settings; Test Design | |
| | Assessment | Data Mining Success Criteria; Test Design; Models | |
| Evaluation | Assessment w.r.t. Business Success Criteria | Business Success Criteria; Terminology | |
| | Review of Process | Business Goals; Assessment w.r.t. Business Success Criteria | |
| | Next Steps | Project Plan; Assessment w.r.t. Business Success Criteria | |
| Deployment | Deployment Plan | Business Goals; Requirements, Assumptions & Constraints | Maintenance Plan |
| | Maintenance Plan | Business Goals; Requirements, Assumptions & Constraints | Deployment Plan |
| | Final Report & Presentation | Business Goals; Terminology; Assessment w.r.t. Business Success Criteria | |
| | Experience Documentation | Project Plan; Review of Process | |

V Appendix

1 Glossary/terminology

Activity – Part of a task in the User Guide; describes actions to perform a task

CRISP-DM methodology – The general term for all concepts developed and defined in CRISP-DM

Data mining context – A set of constraints and assumptions, such as problem type, techniques or tools, application domain

Data mining problem type – A class of typical data mining problems, such as data description and summarization, segmentation, concept descriptions, classification, prediction, dependency analysis

Generic – A task that holds across all possible data mining projects

Model – The ability to apply algorithms to a dataset to predict target attributes; executable

Output – The tangible result of performing a task

Phase – A term for the high-level part of the CRISP-DM process model; consists of related tasks

Process instance – A specific project described in terms of the process model

Process model – Defines the structure of data mining projects and provides guidance for their execution; consists of reference model and user guide

Reference model – Decomposition of data mining projects into phases, tasks, and outputs

Specialized – A task that makes specific assumptions in specific data mining contexts

Task – A series of activities to produce one or more outputs; part of a phase

User guide – Specific advice on how to perform data mining projects



2 Data mining problem types

Usually, the data mining project involves a combination of different problem types, which together solve the business problem.

2.1 Data description and summarization

Data description and summarization aims at the concise description of characteristics of the data, typically in elementary and aggregated form. This gives the user an overview of the structure of the data. Sometimes, data description and summarization alone can be an objective of a data mining project. For instance, a retailer might be interested in the turnover of all outlets broken down by categories. Changes and differences in a previous period could be summarized and highlighted. This kind of problem would be at the lower end of the scale of data mining problems.

In almost all data mining projects, however, data description and summarization is a subordinate goal in the process, typically in its early stages. At the beginning of a data mining process, the user often knows neither the precise goal of the analysis nor the precise nature of the data. Initial exploratory data analysis can help users to understand the nature of the data and to form potential hypotheses for hidden information. Simple descriptive statistical and visualization techniques provide first insights into the data. For example, the distribution of customers by age and geographic regions suggests which parts of a customer group need to be addressed by further marketing strategies.

Data description and summarization typically occurs in combination with other data mining problem types. For instance, data description may lead to the postulation of interesting segments in the data. Once segments are identified and defined, a description and summarization of these segments is useful. It is advisable to carry out data description and summarization before any other data mining problem type is addressed. In this document, this is reflected in the fact that data description and summarization is a task in the data understanding phase.

Summarization also plays an important role in the presentation of final results. The outcomes of the other data mining problem types (e.g., concept descriptions or prediction models) may also be considered summarizations of data, but on a higher conceptual level.

Many reporting systems, statistical packages, OLAP, and EIS systems can cover data description and summarization but do usually not provide any methods to perform more advanced modeling. If data description and summarization is considered a stand-alone problem type and no further modeling is required, then these tools may be appropriate to carry out data mining engagements.

2.2 Segmentation

Segmentation aims at the separation of the data into interesting and meaningful subgroups or classes. All members of a subgroup share common characteristics. For instance, in shopping basket analysis, one could define segments of baskets depending on the items they contain.

Segmentation can be performed manually or semi-automatically. The analyst can hypothesize certain subgroups as relevant for the business question, based either on prior knowledge or on the outcome of data description and summarization. In addition, there are also automatic clustering techniques that can detect previously unsuspected and hidden structures in data that allow segmentation.

Segmentation can sometimes be a data mining objective. Then the detection of segments would be the main purpose of a data mining project. For example, all addresses in ZIP code areas with higher than average age and income might be selected for mailing advertisements for nursing home insurance.

Very often, however, segmentation is a step toward solving other problem types. Then, the purpose is to keep the size of the data manageable or to find homogeneous data subsets that are easier to analyze. Typically, in large datasets various influences overlay each other and obscure the interesting patterns. Then, appropriate segmentation makes the task easier. For instance, analyzing dependencies between items in millions of shopping baskets is very hard. It is much easier (and more meaningful, typically) to identify dependencies in interesting segments of shopping baskets—for instance, high-value baskets, baskets containing convenience goods, or baskets from a particular day or time.

Note: In the literature, there is some ambiguity in the meaning of certain terms. Segmentation is sometimes called clustering or classification. The latter term is confusing because some people use it to refer to the creation of classes, while others mean the creation of models to predict known classes for previously unseen cases. In this document, we restrict the term classification to the latter meaning (see below) and use the term segmentation for the former meaning, though classification techniques can be used to elicit descriptions of the segments discovered.

Appropriate techniques:

- Clustering techniques
- Neural networks
- Visualization



Example:

A car company regularly collects information about its customers concerning their socio-economic characteristics like income, age, sex, profession, etc. Using cluster analysis, the company can divide its customers into more understandable subgroups and analyze the structure of each subgroup. Specific marketing strategies are deployed for each separate group.

2.3 Concept descriptions

Concept description aims at an understandable description of concepts or classes. The purpose is not to develop complete models with high prediction accuracy, but to gain insights. For instance, a company may be interested in learning more about its loyal and disloyal customers. From a concept description of these concepts (loyal and disloyal customers) the company might infer what could be done to keep customers loyal or to transform disloyal customers to loyal customers.

Concept description has a close connection to both segmentation and classification. Segmentation may lead to an enumeration of objects belonging to a concept or class without providing any understandable description. Typically, segmentation is carried out before concept description is performed. Some techniques—conceptual clustering techniques, for example—perform segmentation and concept description at the same time.

Concept descriptions can also be used for classification purposes. On the other hand, some classification techniques produce understandable classification models, which can then be considered concept descriptions. The important distinction is that classification aims to be complete in some sense. The classification model needs to apply to all cases in the selected population. On the other hand, concept descriptions need not be complete. It is sufficient if they describe important parts of the concepts or classes. In the example above, it may be sufficient to get concept descriptions of those customers who are clearly loyal.

Appropriate techniques:

- Rule induction methods
- Conceptual clustering

Example:

Using data about the buyers of new cars and a rule induction technique, a car company could generate rules that describe its loyal and disloyal customers. Below are examples of the generated rules:

If *SEX = male and AGE > 51* then *CUSTOMER = loyal*
If *SEX = female and AGE > 21* then *CUSTOMER = loyal*
If *PROFESSION = manager and AGE < 51* then *CUSTOMER = disloyal*
If *FAMILY STATUS = bachelor and AGE < 51* then *CUSTOMER = disloyal*

2.4 Classification

Classification assumes that there is a set of objects characterized by some attributes or features that belong to different classes. The class label is a discrete (symbolic) value and is known for each object. The objective is to build classification models (sometimes called classifiers), that assign the correct class label to previously unseen and unlabeled objects. Classification models are mostly used for predictive modeling.

The class labels can be given in advance—defined by the user, for instance, or derived from segmentation. Classification is one of the most important data mining problem types that occurs in a wide range of applications. Many data mining problems can be transformed to classification problems. For example, credit scoring tries to assess the credit risk of a new customer. This can be transformed to a classification problem by creating two classes, good and bad customers. A classification model can be generated from existing customer data about their credit behavior. This classification model can then be used to assign new customers to one of the two classes and accept or reject them.

Classification has connections to almost all other problem types. Prediction problems can be transformed to classification problems by binning continuous class labels, since binning techniques allow transforming continuous ranges into discrete intervals. These discrete intervals, rather than the exact numerical values, are then used as class labels, and hence lead to a classification problem. Some classification techniques produce understandable class or concept descriptions. There is also a connection to dependency analysis because classification models typically exploit and elucidate dependencies between attributes.

Segmentation can either provide the class labels or restrict the dataset so that good classification models can be built. It is useful to analyze deviations before a classification model is built. Deviations and outliers can obscure the patterns that would allow a good classification model. On the other hand, a classification model can also be used to identify deviations and other problems with the data.

Appropriate techniques:

- Discriminant analysis
- Rule induction methods
- Decision tree learning
- Neural networks
- K nearest neighbor
- Case-based reasoning
- Genetic algorithms



Example:

Banks generally have information on the payment behavior of their credit applicants. Combining this financial information with other information about the customers, like sex, age, income, etc., it is possible to develop a system to classify new customers as good or bad customers (i.e., the credit risk in acceptance of a customer is either high or low).

2.5 Prediction

Another important problem type that occurs in a wide range of applications is prediction. Prediction is very similar to classification. The only difference is that in prediction the target attribute (class) is not a discrete qualitative attribute but a continuous one. The aim of prediction is to find the numerical value of the target attribute for unseen objects. In the literature, this problem type is sometimes called regression. If prediction deals with time-series data, then it is often called forecasting.

Appropriate techniques:

- Regression analysis
- Regression trees
- Neural networks
- K nearest neighbor
- Box-Jenkins methods
- Genetic algorithms

Example:

The annual revenue of an international company is correlated with other attributes like advertisement, exchange rate, inflation rate, etc. Having these values (or reliable estimates), the company can predict its expected revenue for the next year.

2.6 Dependency analysis

Dependency analysis consists of finding a model that describes significant dependencies (or associations) between data items or events. Dependencies can be used to predict the value of a data item given information on other data items. Although dependencies can be used for predictive modeling, they are mostly used for understanding. Dependencies can be strict or probabilistic.

Associations are a special case of dependencies, which have recently become very popular. Associations describe affinities of data items (i.e., data items or events which frequently occur together). A typical application scenario for associations is the analysis of shopping baskets. There, a rule like “in 30 percent of all purchases, beer and peanuts have been bought together” is a typical example for an association.

Algorithms for detecting associations are very fast and produce many associations. Selecting the most interesting ones is a challenge.

Dependency analysis has close connections to prediction and classification, for dependencies are implicitly used for the formulation of predictive models. There is also a connection to concept descriptions, which often highlight dependencies.

In applications, dependency analysis often co-occurs with segmentation. In large datasets, dependencies are seldom significant because many influences overlay each other. In such cases, it is advisable to perform a dependency analysis on more homogeneous segments of the data.

Sequential patterns are a special kind of dependencies in which the order of events is considered. In a shopping basket analysis, associations describe dependencies between items at a given time. Sequential patterns describe shopping patterns of one particular customer or a group of customers over time.

Appropriate Techniques:

- Correlation analysis
- Regression analysis
- Association rules
- Bayesian networks
- Inductive logic programming
- Visualization techniques

Example 1:

Using regression analysis, a business analyst has found that there are significant dependencies between the total sales of a product and both its price and the amount spent on advertising. This knowledge enables the business to reach the desired level of the sales by changing the product's price and/or the advertisement expenditure.

Example 2:

Applying association rule algorithms to data about car accessories, a car company has found that in 95 percent of cases, if a CD player is ordered, an automatic transmission is ordered as well. Based on this dependency, the car company decides to offer these accessories as a package, which leads to cost reduction.







To learn more, please visit www.spss.com. For SPSS office locations and telephone numbers, go to www.spss.com/worldwide.

SPSS is a registered trademark and the other SPSS products named are trademarks of SPSS Inc. All other names are trademarks of their respective owners.
© 2000 SPSS Inc. All rights reserved. CRISPMWP-1104

dnc



Data Science & Machine Learning

#HARDWORK

Esperamos que o material
tenha sido útil para você!
Foco nos estudos!