



## Wissensentdeckung...

...wird laut Gartner Group (1999) in den nächsten 10 Jahren sprunghaft ansteigen: bei n Projekten 1999 werden 2009 schon n<sup>4</sup> Projekte bestehen.

... dient den Firmen vor allem dazu, ihren Kunden besser „zuzuhören“ und die Beziehung zu ihnen zu verbessern.

Man kann aber auf der Grundlage vorhandener Daten noch vieles andere verbessern.



# Kdnuggets 2002 Poll

## Industries/fields where you currently apply data mining: [608 votes total]

Banking (77)		13%
Biology/Genetics/Proteomics (32)		5%
Direct Marketing/Fundraising (42)		7%
eCommerce/Web (53)		9%
Entertainment (10)		2%
Fraud Detection (51)		8%
Insurance (36)		6%
Investment/Stocks (17)		3%
Manufacturing (28)		5%
Pharmaceuticals (31)		5%
Retail (36)		6%
Scientific data (51)		8%
Security (14)		2%
Supply Chain Analysis (21)		3%
Telecommunications (56)		9%

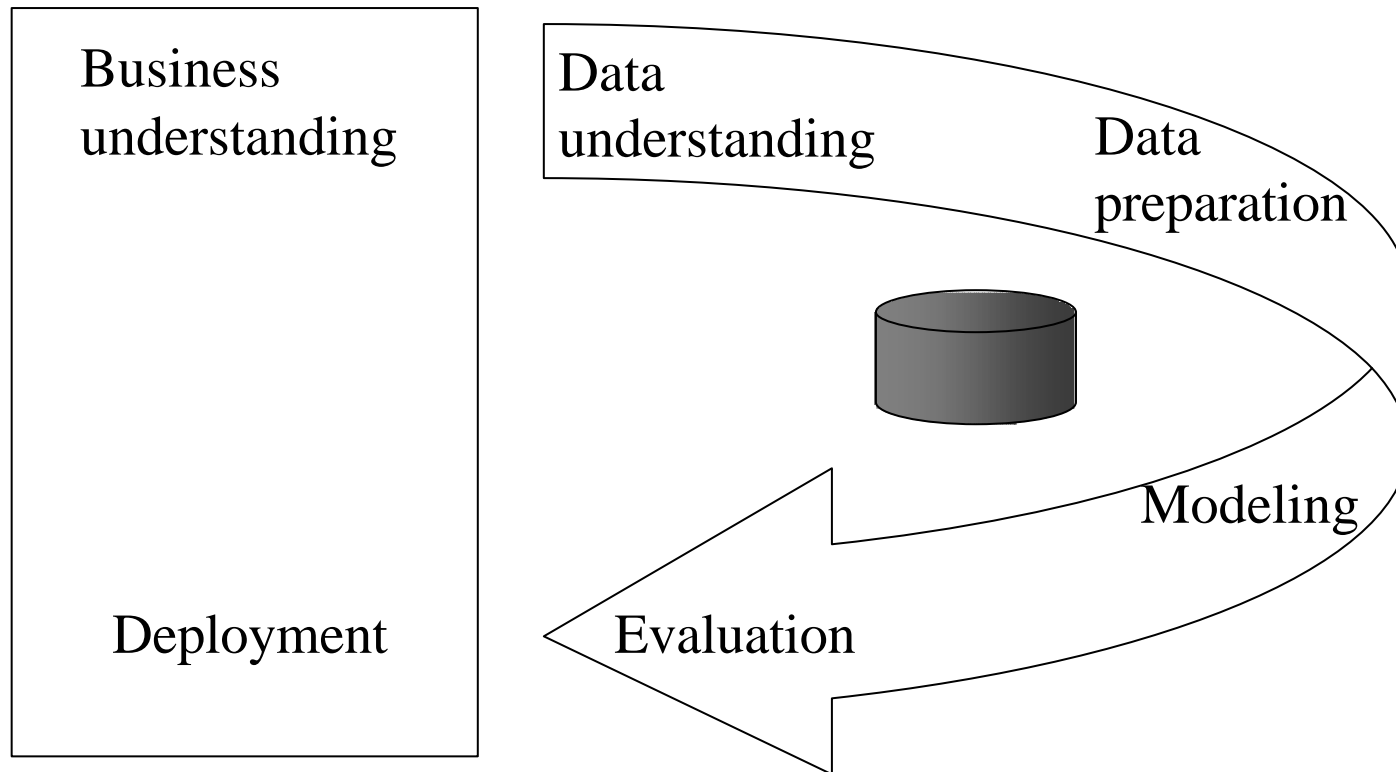


## Häufig benutzte Systeme/ Verfahren

- SPSS CLEMENTINE (128)
- Weka (101)
- SAS (100)
- CART MARS (89)
- SAS Enterprise Miner (67)
- IBM Intelligent Miner (35)
- Statistische Verfahren
- Verfahren des maschinellen Lernens
- Data cube und Frequent Itemsets



# CRISP-DM Process Model



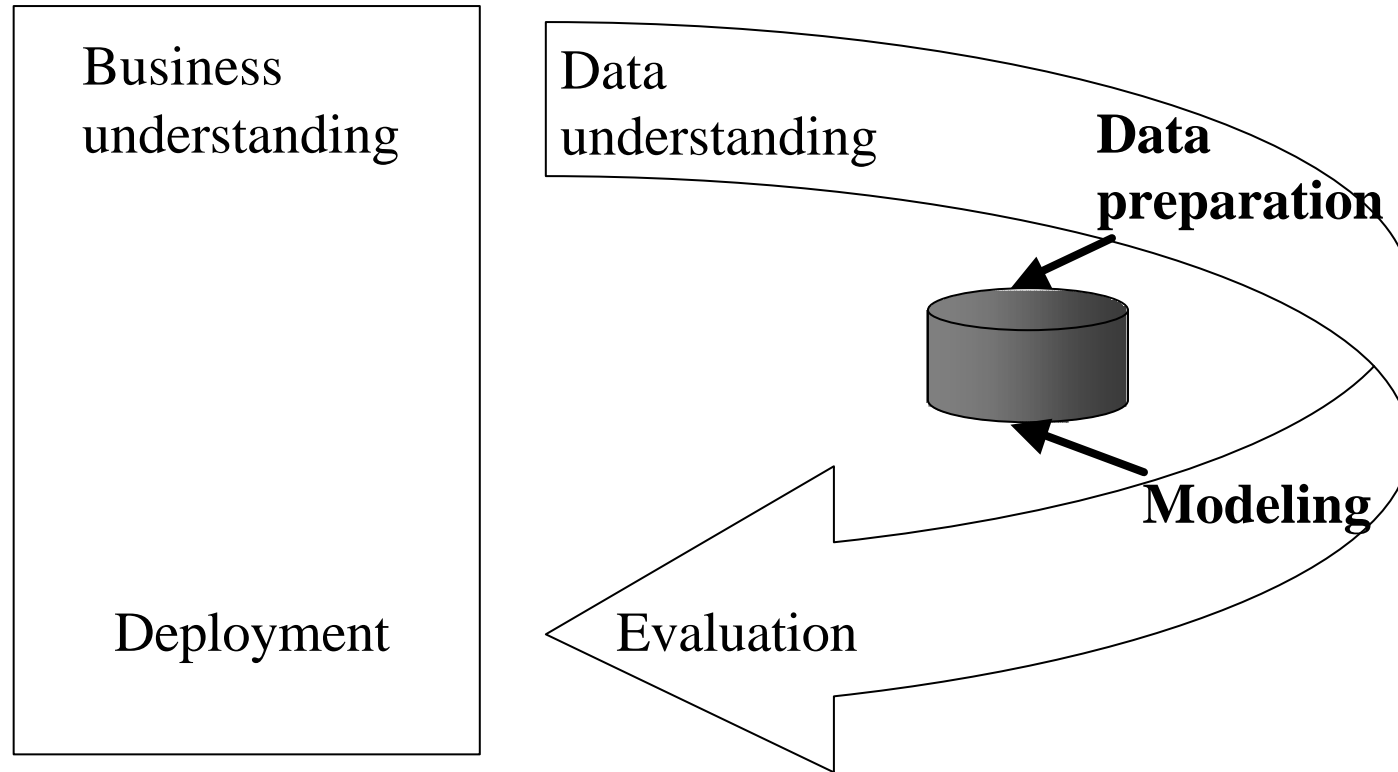


# Ausgangssituation

- Datenvorverarbeitung muss besser unterstützt werden
  - 80 % des Aufwands der Datenanalyse sind Vorverarbeitung
  - Bessere Daten bringen bessere Analyseergebnisse
- Dokumentation der Vorverarbeitung fehlt meist
  - Ähnliche Fälle müssen jedesmal neu erstellt werden
  - Erfahrungen der Mitarbeiter werden nicht weitergegeben
- Datenanalyse ist eine Anfragesprache an Datenbanken
  - Datenanalyse sollte in DBMS integriert sein
  - Anfragen für die Datenanalyse müssen hoch optimiert werden



# Mining Mart Focus





# Mining Mart Ziele

- Operatoren für die Vorverarbeitung
  - direkt auf der Datenbank
  - maschinelles Lernen für die Vorverarbeitung
- Dokumentation
  - der Daten
  - der Fälle
- Wiederverwendung von abstrahierten Fällen

# Mining Mart Ansatz

- Metadaten zur Beschreibung von
  - Daten,
  - Operatoren und
  - Fällen (Sequenzen von Operatoren)
- Compiler, der Metadaten in ausführbaren SQL-Code übersetzt oder externe Verfahren aufruft
- Sammlung von Fällen in Form von operationalen Metadaten





# Mining Mart Ansatz

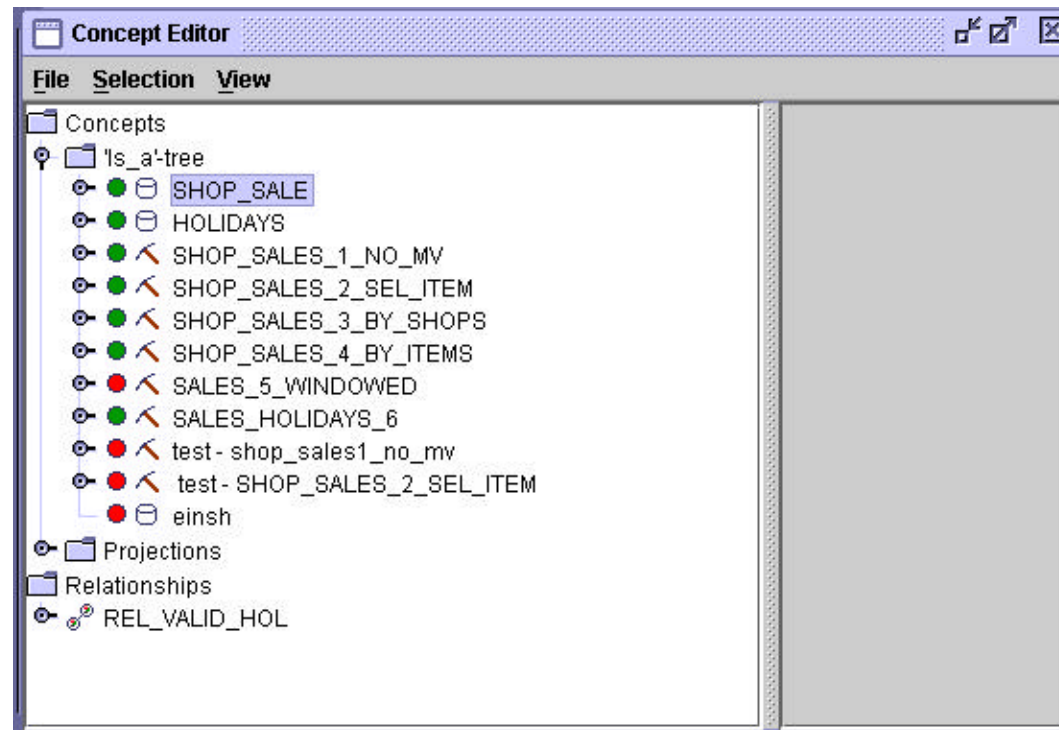
- Unterstützung der Benutzer bei der Vorverarbeitung
  - Erfolgreiche Fälle als Vorlage in Form von Metadaten
  - Anwendbarkeitsbedingungen bei Operatoren
- Offener Ansatz
  - Fälle mit Beschreibung im Internet
  - Erweiterbarkeit durch Einbinden neuer Operatoren
- Publikation der Metadaten - Geheimhaltung der relationalen Geschäftsdaten (incl. DB-Schema)



# Meta Modell für Metadaten

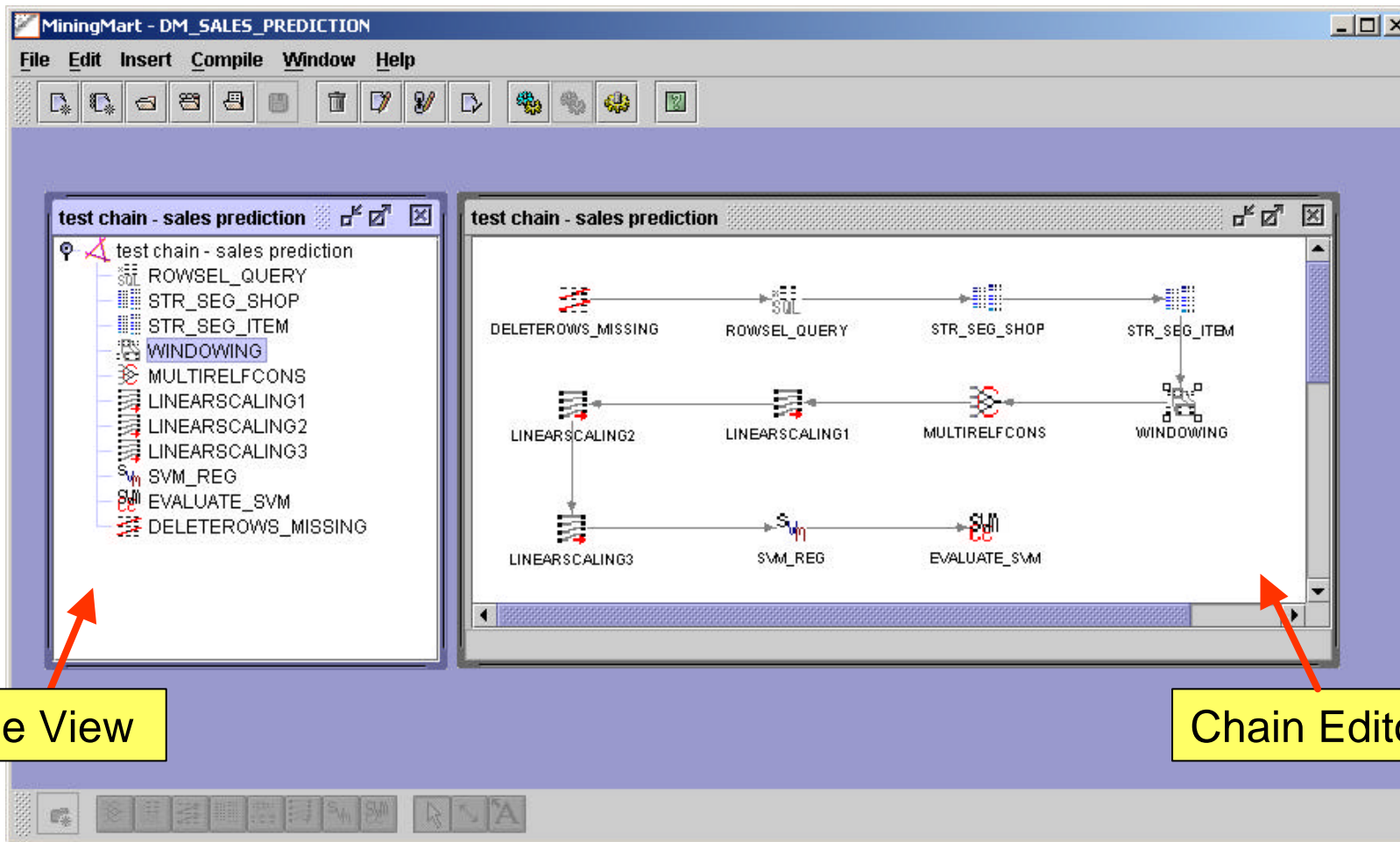
<p>Das begriffliche Modell beschreibt die Objekte und Klassen der Anwendung</p>	<p>Das Fallmodell beschreibt Operatorketten</p>
<p>Das relationale Modell bescheibt die Datenbank</p>	<p>Das Ausführungsmodell generiert SQL statements oder Aufrufe externer Verfahren</p>

# The Concept Editor



- Definieren und editieren von Begriffen und Relationen
- Abbildung von Begriffen und Relationen auf die Datenbank

# The Case Editor



Tree View

Chain Editor

# Setting up an SVM Step

**SVM\_REG - SupportVectorMachineForRegression**

InputConcept: SALES\_HOLIDAYS\_6 [Change]

Target Attribute: SCALED\_WINDOW2 [Change]

**Kernel Type:** Anova [Change]

Sample Size: 200

LossFunction Pos: 1

LossFunction Neg: 20

C: .01

Epsilon: .5

Output Attribute: PREI

**Predicting Attributes:**

- SCALED\_WEEK
- SCALED\_WINDOW1
- ADVENT\_48\_51
- OSTERN

[Add] [Insert] [Delete]

[Save] [Cancel] [Close]

**select new Kernel Type**

- dot
- polynomial
- neutral
- radial
- anova

[select] [cancel]

# The Internet Case Base

The screenshot shows a Netscape browser window displaying a web page titled 'DM\_SALES\_PREDICTION'. The page is divided into several sections:

- InfoLayer**: A blue header bar.
- Overview**: A link to the overview page.
- Concepts**: A list of concepts including:
  - Object
    - Step
    - Case
    - ParameterObject
      - Concept
      - MultiColumnFeature
      - BaseAttribute
      - Value
    - BA\_CONCEPT\_T
    - ColumnSet
    - Column
    - Operator
    - Parameter
    - STEPSEQUENCE\_T
    - ColumnDatatype
    - ConceptualDatatype
    - BA\_COLUMN\_T
    - CONCEPT\_CASE\_T
    - User
- Administration**: A link to the login page.
- DM\_SALES\_PREDICTION**: A table showing case details:
 

CA_ID	1000000467
CA_NAME	DM_SALES_PREDICTION
CA_MODE	FINAL
CA_POPULATION	0
CA_OUTPUT	0
- Step**: A list of steps including:
  - DELETEROWS\_MISSING
  - EVALUATE\_SVM
  - LINEARSCALING
  - MULTIRELFCONS
  - ROWSEL\_QUERY
  - STR\_SEG\_ITEM
  - STR\_SEG\_SHOP
  - SVM\_REG
  - WINDOWING
- Concept**: A list of concepts including:
  - DMTIME
  - DM\_HOLIDAY
  - DELETED\_MISSING\_VALUES
  - ROWSEL\_NEW
  - SEG\_SHOPS
  - SEG\_ITEMS
  - WINDOWED\_NEW
  - MULTIRELFEATURECONS

## Benutzer des Metamodells

- Der Datenbankadministrator liefert das relationale Modell.
- Der Anwender liefert das begriffliche Modell.
- Die Datenanalyseexpertin liefert das Fallmodell oder passt es an.  
Die ersten Fälle werden vom Mining Mart Projekt erstellt.
- Das MiningMart Projekt entwickelt Operatoren bzw. zu vorhandenen deren Metadaten.



# Abstraktionsebenen

- Formalismus zur Repräsentation von Sachbereichsbegriffen, Relationen und KDD-Prozess
  - Was ist ein Begriff, eine Relation, ein KDD-Prozess, ein Operator
- Sprache zur Repräsentation eines bestimmten Sachbereichs, bestimmter Relationen und eines bestimmten Falles
  - Welche Begriffe, Relationen habe ich in einem bestimmten Fall?
- Implementierung des Formalismus und der Sprache als Datenbanktabellen, über denen der Compiler arbeitet und auf die die GUI zugreifen.



# Wir sind keine Benutzer!

- Wie sieht das Metamodell aus?  
Welche Repräsentation braucht man, um
  - den allgemeinen KDD-Prozess,
  - die Operatoren mit ihren Parametern, Bedingungen und Zusicherungen,
  - Sachbereiche (Begriffe)
  - einen bestimmten Fall mit seinen
  - bestimmten Parametersetzungen
  - und einem bestimmten Sachbereich zu beschreiben?
- Wie realisiert man das Metamodell?

M4

Instanz



# Relational Data Model

## Formalismus zur Beschreibung der Datenbank

- Column
  - Supertype: Attribute
  - Subtypes: None
  - Attributes: name (of the column), dataType
  - Associations: belongsToColumnSet, keys, correspondsToBaseAttribute
- ColumnSet (meist eine Tabelle oder Sicht)
- ColumnStatistics
- ColumnSetStatistics
- Key, PrimaryKey, ForeignKey



# Conceptual Data Model

## Formalismus zur Beschreibung des Sachbereichs

- Concept
  - Supertype: Class
  - Attributes: name, subConceptRestriction
  - Associations: isA, correspondsToColumnSet, FromConcept, ToConcept Constraints
- Relationship
- FeatureAttribute
- Value
- UserInput
- RoleRestriction
- DomainDataType

## Conceptual Data Model cont'd

- BaseAttribute
  - SuperType: FeatureAttribute
  - Associations:
    - domainDataType: is it from the raw data or has it been created?
    - isPartOfMultiColumnFeature: pointer to a set of columns which together form a Feature



# Case Model

## Formalismus zur Beschreibung von KDD-Prozessen

- Case
  - Attributes: name,
    - case mode -- {test, final},
    - caseInput -- list of entities from the conceptual model,
    - caseOutput -- concept, normally the input to the data mining step
    - Documentation - free text
  - Associations:
    - listOfSteps - aggregation of steps,
    - population - concept from caseInput, the one the analysis deals with,
    - targetAttributes - FeatureAttribute to which the data analysis is applied

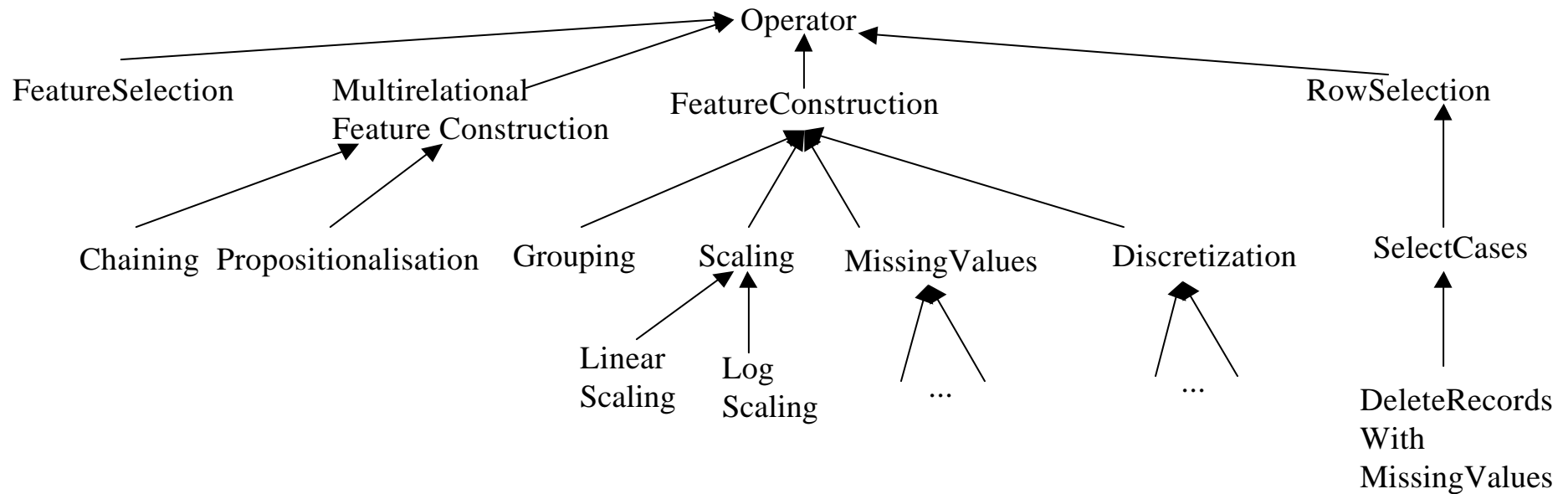


## Case Model cont'd

- Step
  - Attributes: name
  - Associations: belongsToCase, embedsOperator, predecessor, successor
- Operator
  - Attributes with values {yes, no}:
    - loopable -- apply operator several times with changed parameters,
    - multi-stepable - operator delivers several results which will be processed separately in parallel,
    - manual - using no external algorithm
  - Associations:
    - parameters forming the input of the operator
    - conditions -- to be checked given the data,
    - constraints -- to be checked without access to data,
    - assertions - will be true after operator execution

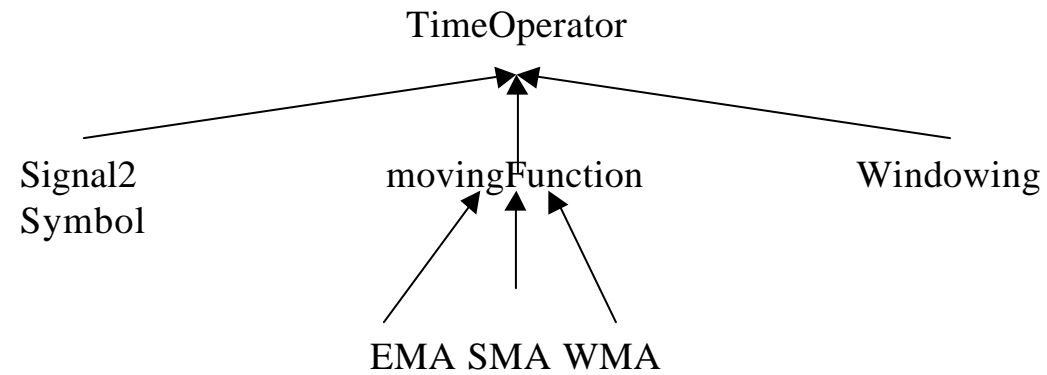


# Operatoren des Metamodells





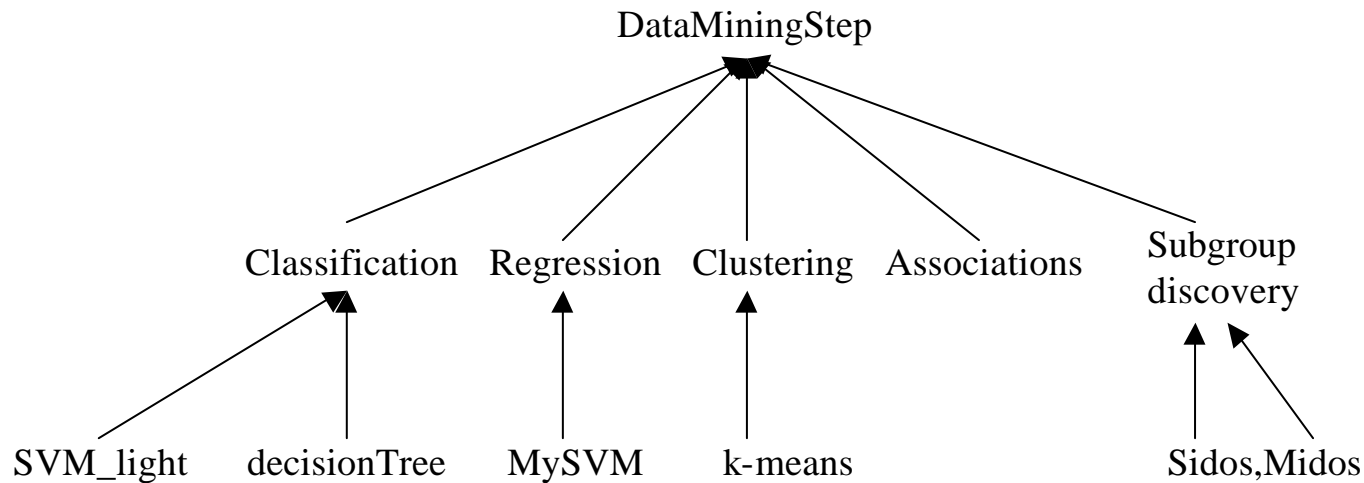
# Time Operators in M4







# Lernoperatoren des Metamodells



*Neu*

Lernoperatoren sind auch Vorverarbeitungoperatoren!  
 Beispiel: C4.5 zur Discretisierung oder Ersetzung fehlender Werte.



## M4 Modell

- Sie haben Ausschnitte aus dem Formalismus gesehen, der für die Definition eines bestimmten Modells benötigt werden.
- Es ist das Metamodell für die Metadaten .  
In Auszeichnungssprachen (SGML, XML) entspricht dies dem Formalismus der DTD, nicht einer DTD, sondern den formalen Mitteln, eine DTD zu definieren.
- Der Formalismus beinhaltet Klassen mit Vererbung, Attribute und Assoziationen.
- Die Definitionen des Metametamodells sind für alle Sprachen in diesem Formalismus gültig.

# Implementierung von M4

- Sowohl M4 als auch ein bestimmter Fall wird in Datenbanktabellen gespeichert.
- Die M4-Tabellen sind Teil des Systems. Sie steuern den Compiler und an sie schließen die GUI an.
- Die Tabellen des konkreten Falles werden jeweils per GUI oder per Datenbankzugriff gefüllt.

## M4 -- Base Attribute

- Basisattribute entsprechen auf der begrifflichen Ebene einer Spalte der relationalen Ebene.
- Tabelle BASATTRIB\_T hat die Attribute:
  - BA\_ID: M4 ID vom Typ Integer
  - BA\_NAME: Name vom Typ String
  - BA\_CONDTID: Typ des Basisattributs, Fremdschlüssel für die Tabelle CONDATYPE\_T, die die Begriffe des M4-Modells enthält
  - BA\_ATTRIBTYPE: ‚DB‘ (Rohdaten) oder ‚MINING‘ (durch Preprocessing erzeugte Spalte)
  - BA\_MCFID: Fremdschlüssel für die Tabelle MCFEATURE\_T, falls dieses Basisattribut gemeinsam mit anderen ein multi-column feature bildet
  - BA\_VALID: ‚YES‘, ‚NO‘ zeigt an, ob es die entsprechende Spalte in den Rohdaten oder in einer erzeugten Sicht gibt.

## Verbindende Tabellen

- BA\_COLUMN\_T verbindet Basisattribute mit éntsprechenden Spalten des relationalen Modells. Attribute:
  - BAC\_ID: ID in dieser Tabelle
  - BAC\_BAID: Fremdschlüssel auf BASEATTRIB\_T
  - BAC\_COLID: Fremdschlüssel auf COLUMN\_T des relationalen Modells
- BA\_CONCEPT\_T verbindet Basisattribute mit ihren Begriffen.
  - BC\_ID: ID in dieser Tabelle
  - BC\_BAID: Fremdschlüssel auf BASEATTRIB\_T
  - BC\_CONID: Fremdschlüssel auf CONCEPT\_T, die Tabelle der Begriffe



# Operatortabellen – für alle Fälle gleich

- OPERATOR\_T: allgemeine Definition der Operatoren, OP\_PARAMS\_T: Eingabe der Operatoren,
- OP\_CONSTR\_T: Bedingungen für die Anwendung eines Operators, die ohne Daten geprüft werden kann
- OP\_COND\_T: Bedingungen für die Anwendung eines Operators, die anhand konkreter Daten geprüft werden muss
- OP\_ASSERT\_T: Aussagen über einen Operator (Nachbedingung)

# Realisierung der Operatordefinition Missing Values

OP\_PARAMS\_T

Input: TheConcept, TargetAttribute,  
PredictingAttributes

Output: FilledAttribute

Condition: TargetAttribute is a BA  
with missing values

Constraint: TargetAttribute and  
PredictingAttributes  
belong to TheConcept

Assertion: FilledAttribute belongs to  
TheConcept,  
FilledAttribute is a BA wo. MV

PARAM-ID	OP_ID	...	IO	TYPE
P001	O001		IN	Concept
P002	O001		IN	BaseAttribute
P003	O001		IN	BaseAttributes

CONST_ID	CONST_OP_ID	TY PE	OBJ_1	OBJ_2
CS01	O001	IN	P002	P001

OPERATOR\_T

OP_ID	OP_NAME	...	OP_MANUAL
O001	Missing Values		YES

COND_ID	COND_OP_ID	TYPE	OBJ_1	OBJ_2
CD01	O001	HAS_NULLS	P002	

# Operator als Schritt in einem KDD-Prozess

- CASE\_T: ID, Name, Validity
- STEP\_T: verbindet CASE\_T und OPERATOR\_T  
zeigt auf eine Tabelle mit konkreten Parametern für diesen Schritt (PARAMETER\_T)

PARAMETER\_T

PAR_ID	PAR_OBJ_ID	PAR_OBJ_TYPE	PAR_OP_ID	PAR_STEP_ID	...
	Foreign link to *	{BA, MCF, CON,..}	Foreign link to OPERATOR_T	Foreign link to STEP_T	

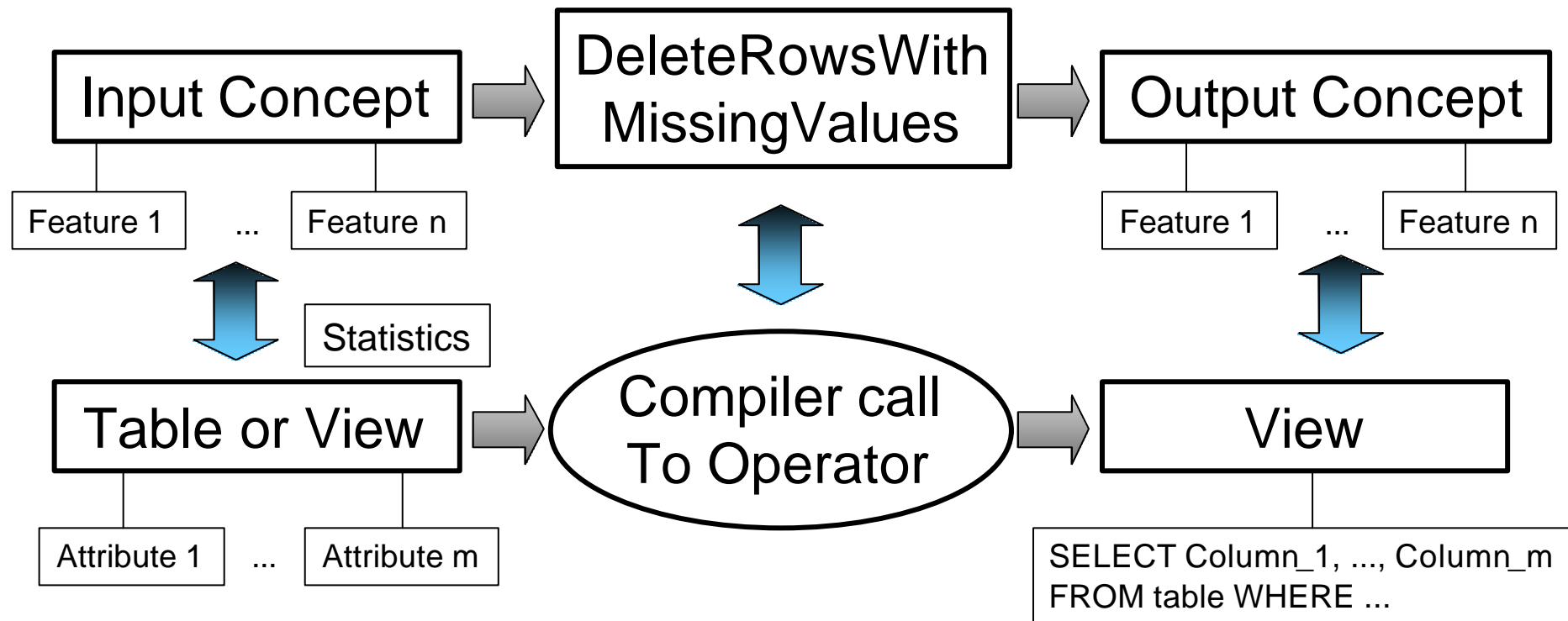
\* ∈ BASEATTRIB\_T  
MCFEATURE\_T  
VALUE\_T  
CONCEPT\_T  
RELATION\_T



## Was wissen Sie jetzt?

- Das Metamodell für KDD-Prozesse und Sachbereiche ist ein Formalismus, mit dem man bestimmte Prozesse und Begriffe definieren kann.
- Das Metamodell ist in Form von Datenbanktabellen gespeichert, die durch Fremdschlüsselbeziehungen verknüpft sind.
- KDD-Prozesse werden als Folge von Schritten beschrieben, wobei ein Schritt auf einen Operator zeigt.
- Ein Operator ist allgemein durch Tabellen beschrieben, die per Fremdschlüssel auf ihn zeigen. Beispiel: `OP_PARAMS_T`
- Ein Schritt zeigt auf einen Operator und wird durch `PARAMETER_T` beschrieben. `PARAMETER_T` zeigt auf Zeilen in Tabellen für Begriffe und Attribute.

# Compiler verbindet begriffliche und Datenbankebene

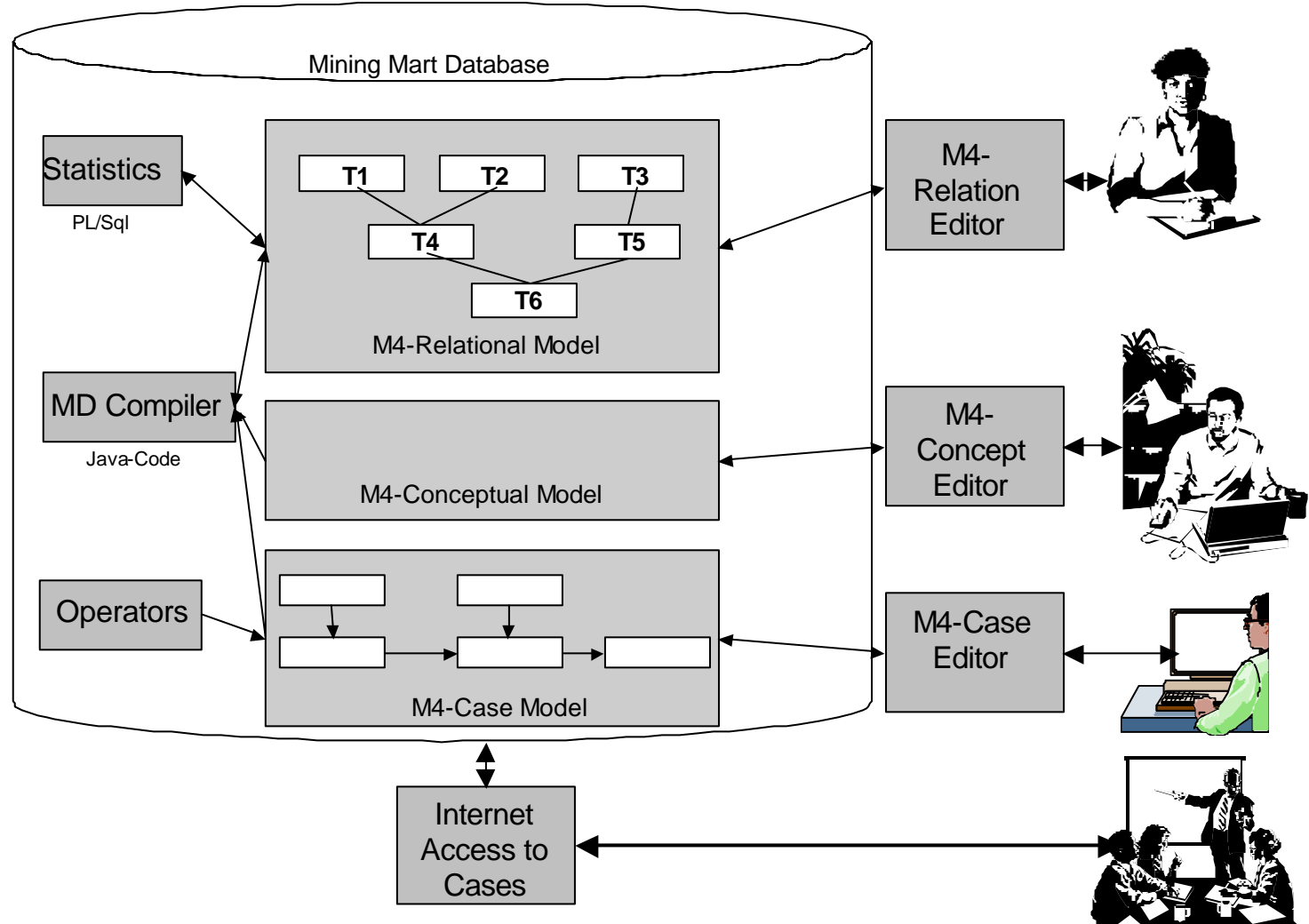




# Operator Applicability Information

- Case Editor
  - preparing output for operator application
  - property window for entering parameter values
  - ensures applicability of operator chains
- M4 Compiler
  - checks if inputs/parameters are available and valid
  - required data properties for applicability at runtime
  - assertions: avoid unnecessary database scans

# Systemarchitektur





# Anwendungen

- Versicherung (SwissLife)
  - Direktes Marketing
  - Analyse der Rückkäufe
- Telekommunikation (TILab, NIT)
  - Unbezahlte Rechnungen, Betrugsvorhersage
  - Unterstützung des Call Centers für Telekommunikationsdienste
- Handel (DM, holländische Zeitungen)
  - Abverkaufsprognose



# Abgeschlossene europäische Forschungsprojekte

- CRITIKAL Client-Server Rule Induction Technology for Industrial Knowledge Acquisition from Large Databases (1996 - 1998)
- KESO Knowledge Extraction
- CRISP-DM Cross Industry Standard Process for Data Mining (1997 - 2000)



# Aktuelle Forschung

- Andere Daten: Zeit, Raum, Genomsequenzen, Sätze
- Engere Anbindung an die Datenbank, Anfrageoptimierung
- Stärkere Unterstützung der frühen Phasen: Dateninspektion, Datenvorverarbeitung
- Metadaten zur Wiederverwendung von DM-Prozessen
- EU-Networks: Sol-EU-net, KDnet
- EU-Projekte: MetaL, Mining Mart, SPIN!, Cinq