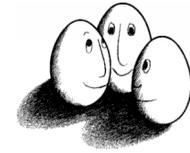


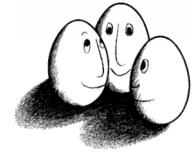
Der CRISP-DM Prozess für Data Mining



Wozu einen standardisierten Prozess?

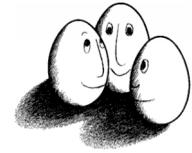
Der Prozess der Wissensentdeckung muss verlässlich und reproduzierbar sein – auch für Menschen mit geringem Data Mining Hintergrundwissen.

- Rahmen für “Speicher-” und “Wiedereinsetzpunkte”
- Hilft bei der Planung und der Verwaltung der Analyseaufgaben
- Leichter Einstieg



CRISP-DM Standard

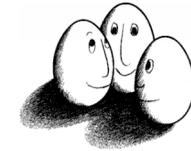
- Cross Industry Standard Process for Data Mining
- Initiative begann September 1996
- Gefördert durch die EU
- Mehr als 200 Mitglieder der CRISP-DM SIG weltweit
 - DM Anbieter - SPSS, NCR, IBM, SAS, SGI, Data Distilleries, Syllogic, Magnify...
 - Berater - Cap Gemini, ICL Retail, Deloitte & Touche...
 - Anwender - BT, ABB, Lloyds Bank, AirTouch, Experian...



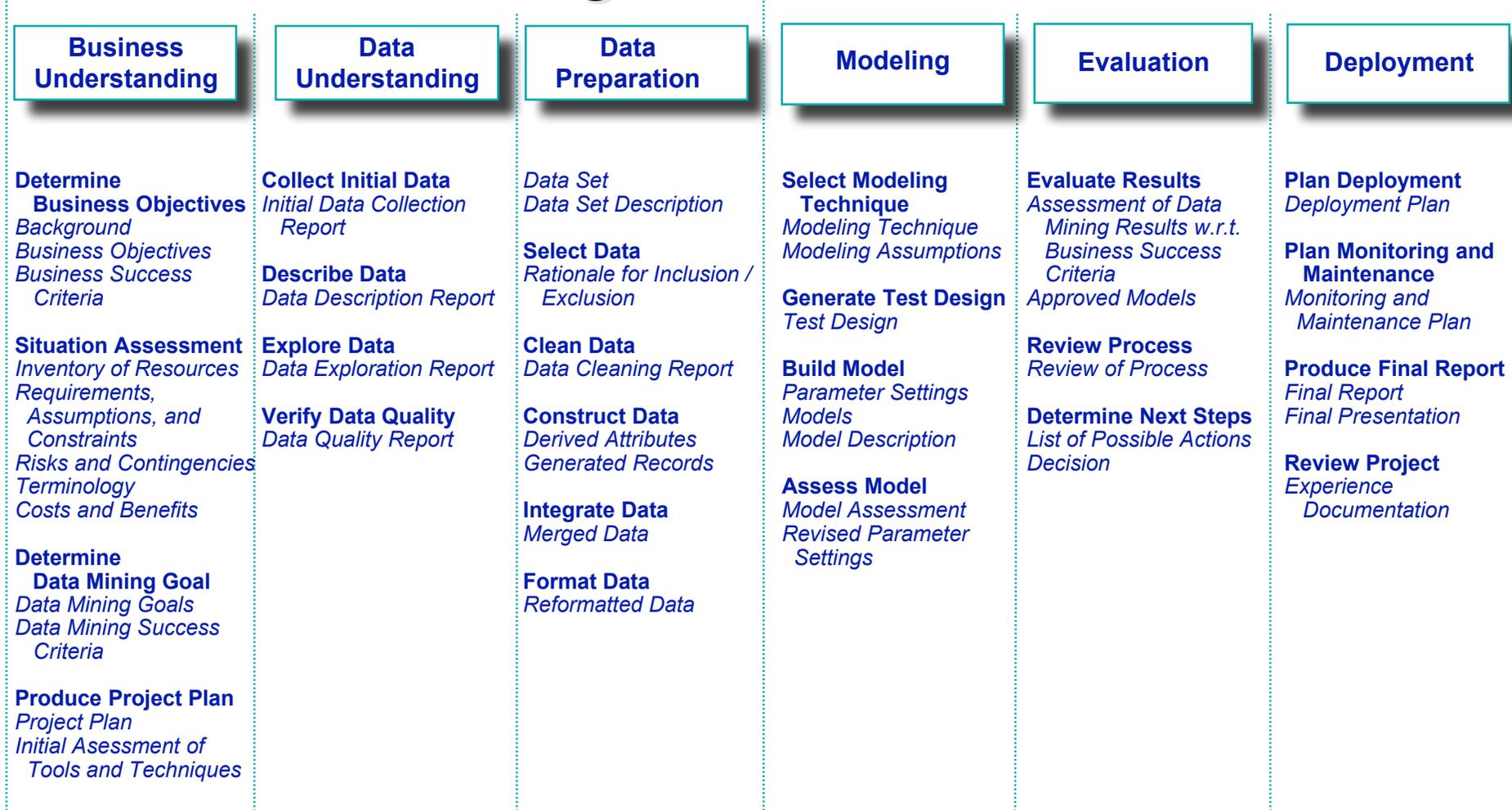
CRISP-DM

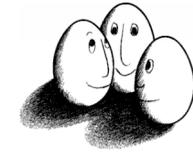
- Nicht-proprietär
- Neutral bezüglich
 - Anwendungsfeldern
 - Applikationen
- Als Leitlinie zu verstehen
- Erfahrungssammlung:
 - Analysetemplates



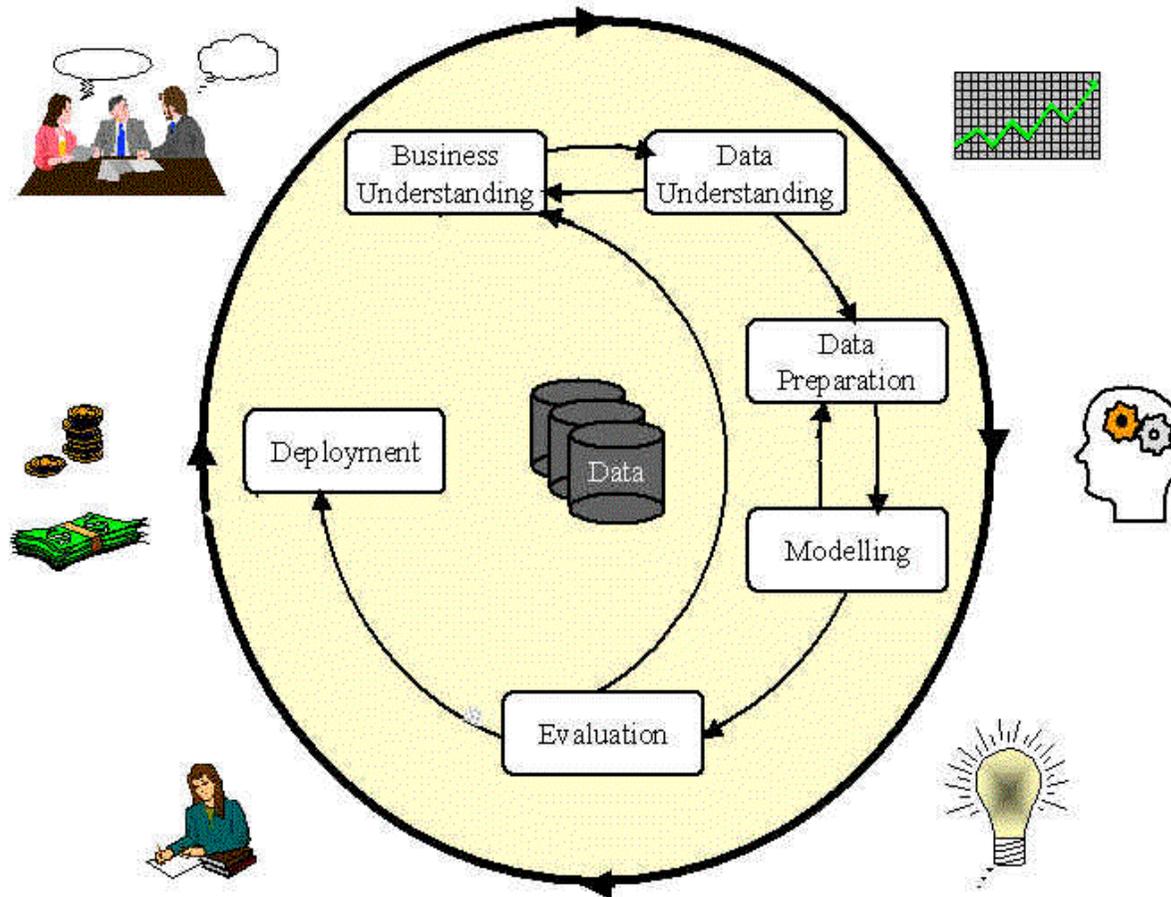


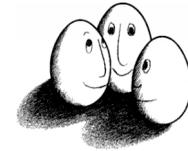
Phasen und Aufgaben





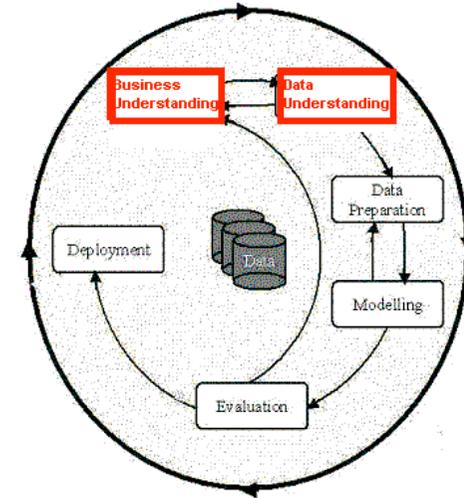
Phasen im DM – Prozess



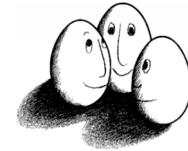


Phasen im DM – Prozess (1 & 2)

- Business Understanding:
 - Aufstellung des (Geschäfts-)ziels
 - Aufstellung des Data Mining Ziels
 - Aufstellung von Erfolgskriterien

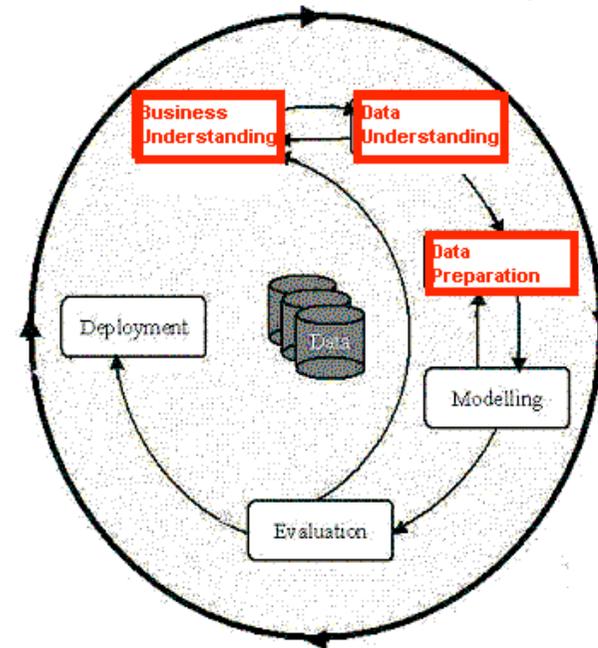


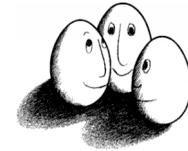
- Data Understanding
 - Untersuchung der Daten und Prüfung der Qualität
 - Outlier entdecken



Phasen im DM – Prozess (3)

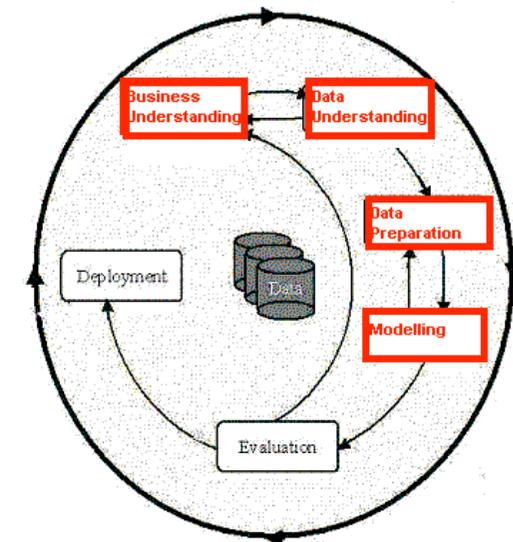
- Data Preparation:
 - Benötigt häufig bis zu 90% der Gesamtzeit
 - Datensammlung
 - Beurteilung
 - Konsolidierung und Säuberung
 - Tabellenlinks, Aggregationstiefe, fehlende Werte, ...
 - Datenauswahl
 - Ausreißer?
 - Sampling?
 - Welche Variablen?
 - Transformationen – Erzeugung neuer Variablen

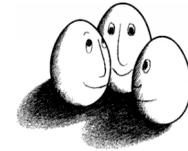




Phasen im DM – Prozess (4)

- Model Building
(Data Mining Schritt)
 - Auswahl der Modellierungstechnik(en) wird auf Basis des erstellten Data Mining Ziels und der Daten getroffen
 - Modellbildung ist ein iterativer Prozess – welcher sich für überwachtes und unüberwachtes Lernen unterscheidet
 - Modelle können beschreibend oder vorhersagend sein
- Vielzahl von Verfahren während der Vorlesung

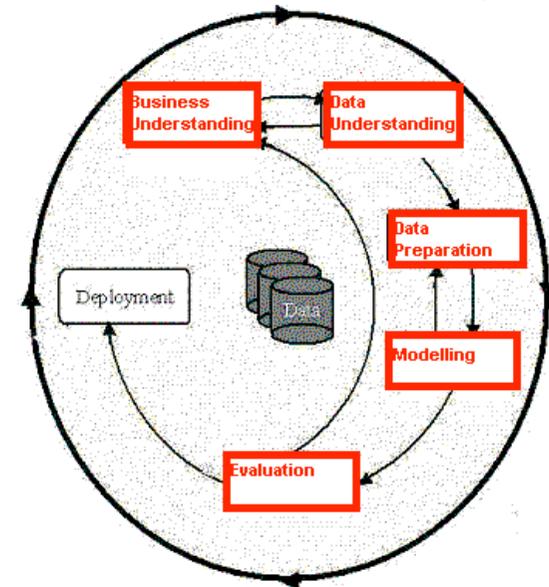


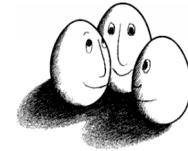


Phasen im DM – Prozess (5)

■ Modell Evaluation

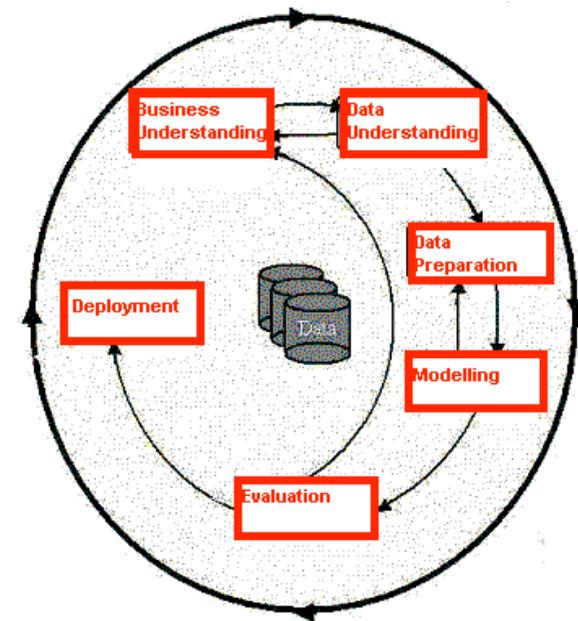
- Evaluation des Modells: wie gut arbeitet es auf Testdaten?
- Evaluationsmethoden und –kriterien hängen vom Modelltypen ab:
 - beispielsweise eine Fehlermatrix für Klassifikationsmodelle oder der mittlere quadratische Fehler für Regressionsmodelle
- Interpretation des Modells: ob wichtig oder nicht, ob einfach oder schwer hängt erneut von den Zielen und den Modelltypen ab

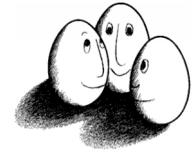




Phasen im DM – Prozess (6)

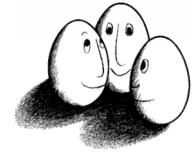
- Deployment
 - Bestimmung, wie die Ergebnisse genutzt werden können
 - Wer wird sie nutzen?
 - Wie oft werden sie genutzt?
- Anwendung der Ergebnisse durch
 - Vorhersagen auf Datenbanken
 - Anwendung als (Geschäfts-)regeln
 - interaktive on-line Vorhersagen





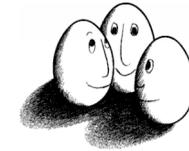
Beispiel: Analyse von Kundendaten

- Gegeben: Kundendaten eines Telefonanbieters
- Aufgabe: Bestimmung der Faktoren, welche Kunden zu “schlechten” Kunden machen, d.h. solchen Kunden mit schlechter Zahlungsmoral
- Screenshots etc. mit Hilfe von RapidMiner – Experimenten erstellt



Business & Data Understanding (1 & 2)

- Geschäftsziel: Reduktion der Zahl der Kunden, welche niemals zahlen werden und damit Erhöhung des Gewinns
- Erfolgskriterium: Reduktion dieser Anzahl um mind. 30%
- Sammlung der Daten bereits erfolgt:
 - Exceltabelle: jede Zeile ein Kunde, jede Spalte eine Eigenschaft (Merkmal, Attribut oder Variable)
 - Merkmale beschreiben Kunden und deren Verhalten
 - Die Spalte “Neverpay” soll vorhergesagt werden



YALE@annex8

File Edit View Experiment Tools Help

ExampleSet

meta data view
 data view
 plot view

ExampleSet (5519 examples, 0 special attributes, 50 regular attributes) View Filter (5519 / 5519): all

Gender	Agentid	SalePEDay	CreditLevel	phonedepo...	vericheck	paymentSta...	Connection...	Connection...	PrePay_Atte...	Neverpay	CreditScore	ModelType	TIC
Female	VIM_3841	0	3	No	U	PAYMENT C	vic_3470	3	0	0	Unknown	P	N
Male	VIM_3320	0	3	No	U	PAYMENT C	vic_1748	12	0	0	Unknown	P	N
Female	vim_3761	0	3	No	U	PAYMENT C	vic_3614	21	0	0	Unknown	P	N
Female	vimd3861	0	3	No	U	PAYMENT C	vic_3470	3	0	0	Unknown	P	N
Female	vimd3778	0	3	No	U	PAYMENT C	vic_3892	22	0	0	Unknown	P	N
Female	vim_2549	0	3	No	U	PAYMENT C	vic_3892	11	0	0	Unknown	P	N
Female	vimd3753	0	3	No	U	PAYMENT C	vic_3892	4	0	0	Unknown	P	N
Female	vim_3623	0	3	No	U	PAYMENT C	vicd1748	4	0	0	Unknown	P	N
Male	vim_2549	0	3	No	U	PAYMENT C	vicd3614	5	0	0	Unknown	P	N
Female	vim_3879	0	3	No	U	PAYMENT C	vicd3614	5	0	0	Unknown	P	Y
Male	vim_3623	0	3	No	U	PAYMENT C	vicd3614	5	0	0	Unknown	P	N
Male	vim_3623	0	3	No	U	PAYMENT C	vicd3470	10	0	0	Unknown	P	N
Female	vimd3281	0	3	No	U	PAYMENT C	vicd3614	5	0	0	Unknown	P	N
Male	VIM_1915	0	3	No	U	PAYMENT C	vicd3884	6	0	0	Unknown	P	N
Female	vim_3843	0	3	No	U	PAYMENT C	vicd3470	9	0	0	Unknown	P	N
Male	vim_3884	0	3	No	U	PAYMENT C	vicd3893	11	0	0	Unknown	P	N
Female	vimd3859	0	3	No	U	PAYMENT C	vicd3892	10	0	0	Unknown	P	N

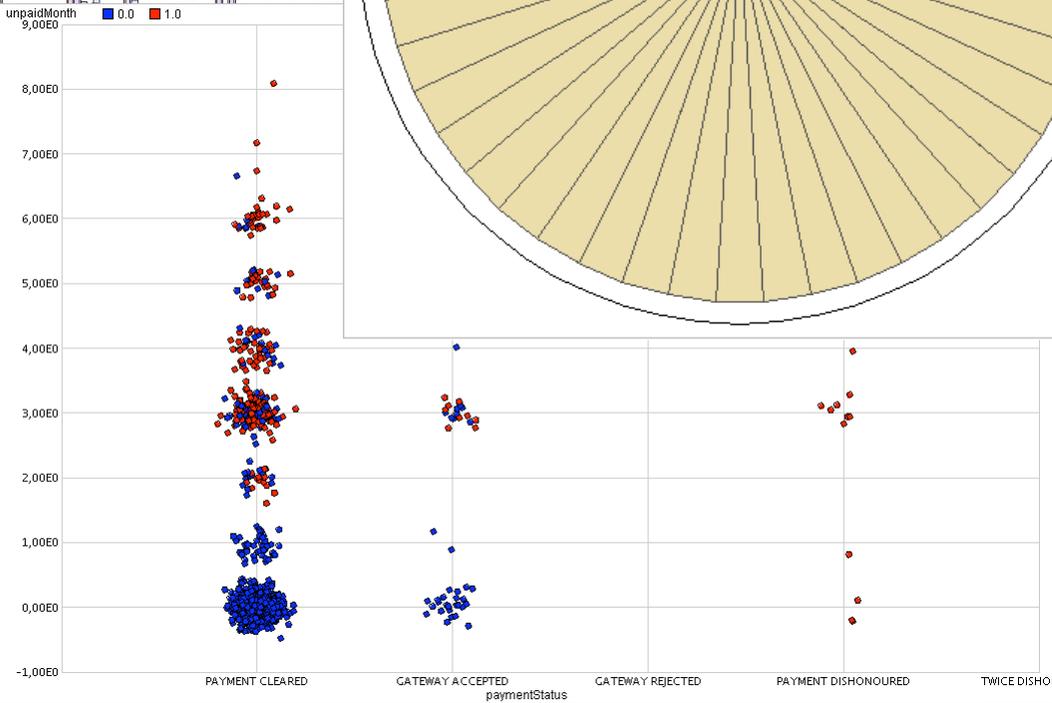
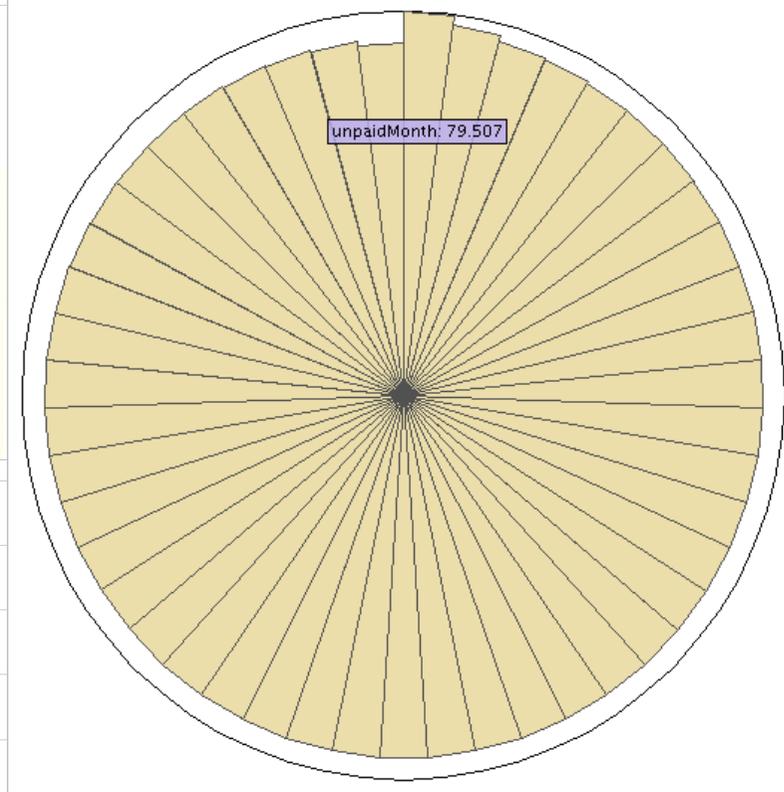
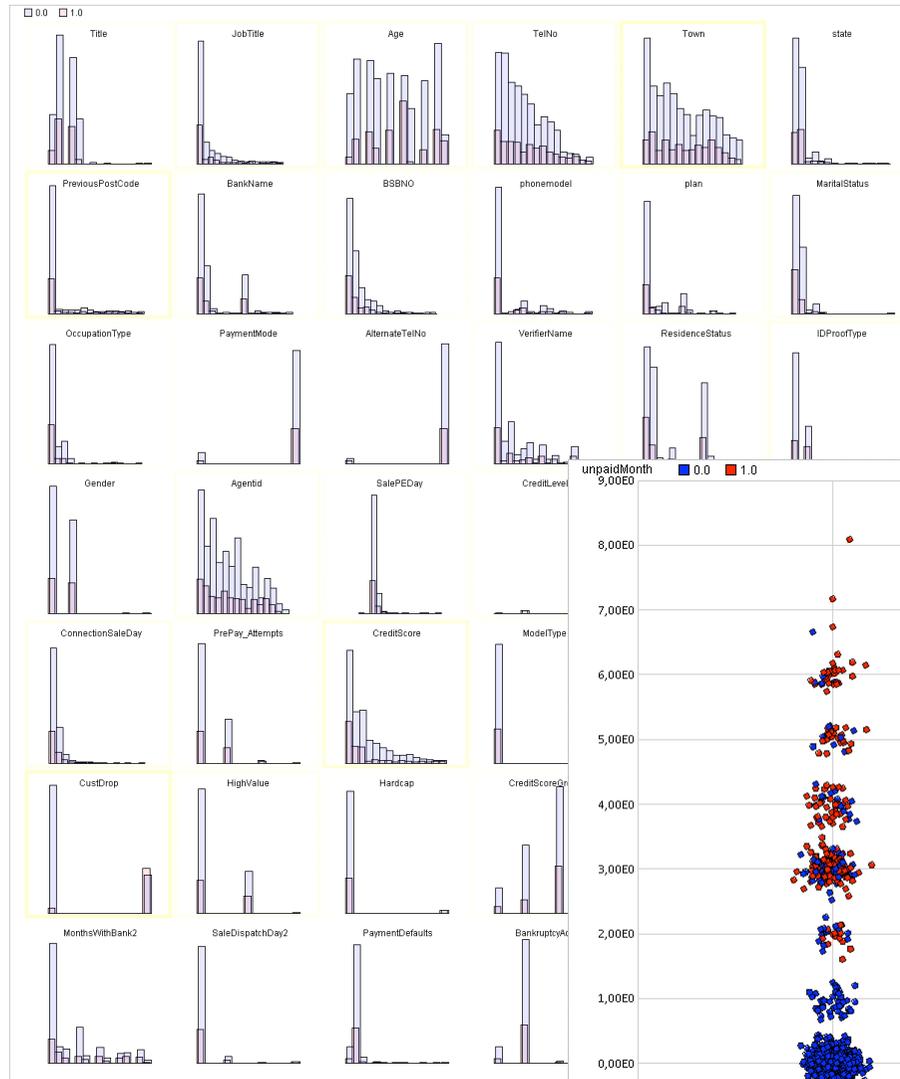
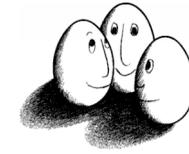
The data for the attribute 'Neverpay'.

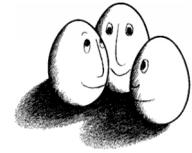
Save...

02.04.2007 19:09:48: Checking experimental setup...
 02.04.2007 19:09:48: Inner operators are ok.
 02.04.2007 19:09:48: Checking i/o classes...
 02.04.2007 19:09:48: i/o classes are ok. Experiment output: ExampleSet.
 02.04.2007 19:09:48: Experiment ok.

19:11:05

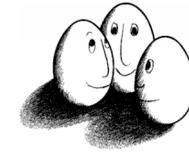
Start | D:\localhome\DMV\yale | YALE@annex8 | Microsoft PowerPoint - [D... | 19:11





Data Preprocessing (3)

- Glück: Die Daten liegen bereits in 1 (!) Tabelle vor. Hauptteil der Arbeit also schon erledigt!
- Ansonsten: ETL (z.B. MiningMart)
- Noch nötig:
 - Umwandlung des numerischen Zielattributs Neverpay in ein nominales Klassifikationsattribut
 - Definition des Zielattributs Neverpay
 - Ersetzung fehlender Werte durch “Unknown”
 - Löschen sinnloser Merkmale
 - Ziehen eines stratifizierten Samples



YALE@annex8 (xval.xml*)

File Edit View Experiment Tools Help

Operator Tree

- Root
 - Experiment
 - ExcelExampleSource
 - ExcelExampleSource
 - AttributeSubsetPreprocessing
 - AttributeSubsetPreprocessing
 - Numeric2Polynomial
 - Numeric2Polynomial
 - ChangeAttributeType
 - ChangeAttributeType
 - MissingValueReplenishment
 - MissingValueReplenishment
 - RemoveUselessAttributes
 - RemoveUselessAttributes
 - StratifiedSampling
 - StratifiedSampling
 - XValidation
 - XValidation
 - Y-NaiveBayes
 - Y-NaiveBayes
 - OperatorChain
 - OperatorChain

Parameters XML Comment New Operator

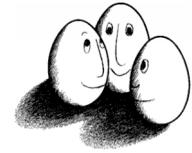
Key	Value
default	value
columns	Edit List (0)...
replenishment_value	Unknown

```

0.0: 2108 299
1.0: 378 527
]
(created by PerformanceEvaluator)
02.04.2007 19:20:15: Experiment finished successfully
    
```

19:20:38

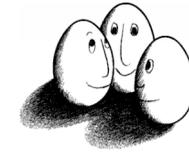
Start | D:\localhome\DMV\yale | YALE@annex8 (xval.x...) | Microsoft PowerPoint - [0... | 19:20



Model Building (4)

- Lernen unterschiedlicher Modelle, hier
 - Naïve Bayes
 - Entscheidungsbaumlerner
 - SVM
 - ...

- Hierzu ist nur der Austausch eines einzigen Operators nötig.



The screenshot displays the YALE software interface with the following components:

- Operator Tree:** A hierarchical list of operators including Root, Experiment, ExcelExampleSource, AttributeSubsetPreprocessing, Numeric2Polynomial, ChangeAttributeType, MissingValueReplenishment (highlighted), RemoveUselessAttributes, StratifiedSampling, XValidation, Y-naiveBayes (highlighted), and OperatorChain.
- Parameters Table:**

Key	Value
default	value
columns	Edit List (0)...
replenishment_value	Unknown
- Model View:** A decision tree diagram showing a root node 'CustDrop' branching into '= No' and '= Yes'. The '= Yes' branch leads to a leaf node '1.0 (13.0)'. The '= No' branch leads to a node 'paymentStatus', which further branches into five leaf nodes: '= PAYMENT CLEARANCE' (0.0 (406.0/92.0)), '= PAYMENT RECEIVED' (0.0 (17.0)), '= PAYMENT RELATED' (0.0 (0.0)), '= PAYMENT DISHONOURABLE' (1.0 (5.0)), and '= PAYMENT DISHONOURED' (0.0 (0.0)).
- Console Window:**

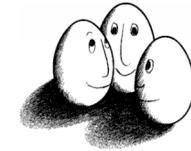
```

0.0: 2108 299
1.0: 378 527
]
(created by PerformanceEvaluator)
02.04.2007 19:20:15: Experiment finished successfully

+-. J48[0] (J48)
+-. OperatorChain[0] (OperatorChain)
+-. ModelApplier[0] (ModelApplier)
+-. PerformanceEvaluator[0] (PerformanceEvaluator)
02.04.2007 19:40:20: J48: Breakpoint reached

[2] J48: breakpoint reached after operator, press resume...

```

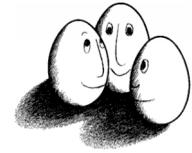


Model Evaluation (5)

- Evaluation der verschiedenen Modelle (Vorverarbeitungen) soll die Zahl der falschen Vorhersagen abschätzen...
- ... in einem “realen” Anwendungsfeld
- Trick: nur Teil der Daten zum Lernen verwenden, den Rest für die Fehlermessung

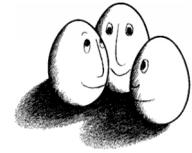
Fehlerreduktion
um mehr als 50%:
reicht!

classification_error: 10.76% +/- 1.50% (mikro: 10.76%)			
	true 0.0	true 1.0	class precision
pred. 0.0	3756	207	94.78%
pred. 1.0	387	1169	75.13%
class recall	90.66%	84.96%	



Deployment (6)

- Gelernte Modelle und Vorverarbeitungsketten können in Systeme / Anwendungen integriert werden
- Analyse der Modelle gibt einen Einblick in das Verhalten von Kunden
- Modelle können nachtrainiert werden bei signifikanter Erweiterung der Datenlage



Was wissen Sie jetzt?

- Der CRISP-Prozess ist ein standardisierter Prozess zur Durchführung von Data Mining
- Zunächst sollten Erfolgskriterien definiert werden
- Data Preprocessing verursacht üblicherweise den größten Aufwand
- Evaluation liefert Abschätzung für den Realfall
- Schritte sollten mehrfach besucht werden bei neuen Erkenntnissen