

# Data Mining

-  
Marketing-Schlagwort  
oder  
ernstzunehmende Innovation?

**Hans-Peter Höschel,  
SAS Institute, Heidelberg**



# Datamining als Marketing-Schlagwort

➤ Wunsch: **grosse Datenmengen** auswerten,  
Tera- und **Gigabytes**, mit “modernen”  
Methoden

➤ **Gigantomanie** von technischen  
Softwareanbietern und Fachjournalisten

➤ **Seriöse** Datamining Anbieter bieten

⇒

**Optimales Aufwand-Nutzen-Verhältnis bei  
verschiedenen Anwendungssituationen**

- » durch abgestuft leistungsfähige Algorithmen
- » seriöse Aufwand-Nutzen Kalkulation

# Typische Datamining Fragen

## ➤ **Datamining-Werbung**

- » "In 20% der Fälle, bei denen ein spezieller Markentoaster verkauft wurde, kauften die Kunden auch passende Küchenhandschuhe und Tischdecken."

## ➤ **Datamining-Wunsch: Datamining beantwortet uns auf Knopfdruck heute die Fragen, die wir morgen stellen wollten.**

## ➤ **Datamining-Realität**

- » seit vielen Jahren funktionsfähige Teil- und Speziallösungen und nunmehr
- » einige neuere Ansätze durch leistungsfähigere DV und spezialisierte Algorithmen

# Datamining - Realität praktisch einsetzbare Software

## ➤ **Rationaler Kern beim Datamining**

- » bekannte und neue Algorithmen möglichst einfach vom Endanwender bei grossen Datenmengen anwenden.

## ➤ **Tendenzen aktueller Datamining Lösungen**

### » **1. Abweichungsanalyse**

- > “intelligente” SQL-Algorithmen = “Einfaches” Datamining:

### » **2. Gruppieren ohne Zielvariable**

- > Visuell interaktiv & automatisch (letzteres bekannt als Clustern)

### » **3. Ursache-Wirkungsanalyse mit Zielvariable**

- > Automatische Response Analyse
- > Entscheidungsbäume & Segmentation
- > Neuronale Netze

# Datamining Anwendungsfelder in Marketing und Produktion

## ➤ **Anwendungsgebiete**

### » **Database(d)-Marketing**

- > Kundenklassifikation Basis: Verkaufs- u. sozio-ökonomische Daten
- > Kaufwahrscheinlichkeiten für bestimmte Produkte

### » **Produkt-Design, TQM - Total Quality Management**

- > bedarfsgerechte Entwicklung neuer Produkte
- > Qualitätskontrolle Abweichungsanalyse

### » **Controlling**

- > Abweichungsanalyse

## ➤ **Endanwender: Kenntnisse Datenanalyse notwendig. Aber Vorteil:**

## ➤ **Keine Annahmen über Zufallsverteilungen**

# Wie funktioniert Datamining ?

## ➤ **Vorarbeiten: Daten bereitstellen**

- » Extraktion, Prüfung, Korrektur, Selektion, Transformation  
= über **80% des Gesamtaufwandes** in grossen Projekten

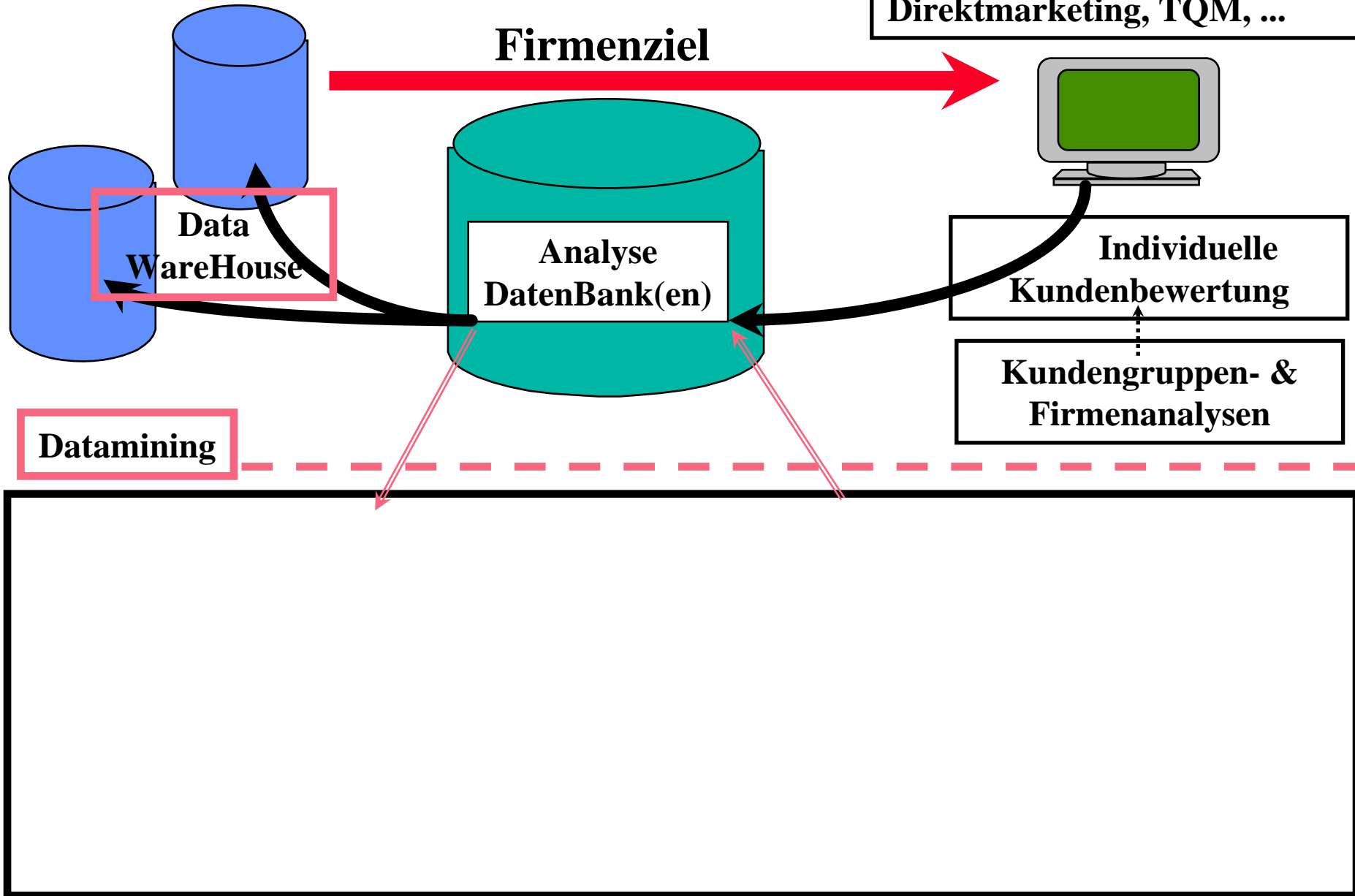
## ➤ **Datamining:**

- » Abweichungsanalysen
- » Klassifikation: Zuordnung von Daten zu Klassen
- » Clustering: Bildung von Klassen ähnlicher Daten
- » Entdecken von Abhängigkeiten und Trends

## ➤ **Umsetzen in Aktionen Marketing und Produktion**

## ➤ **Gesamt-Aufwand: erheblich, aber er rentiert sich.**

# Datamining im Data Warehouse

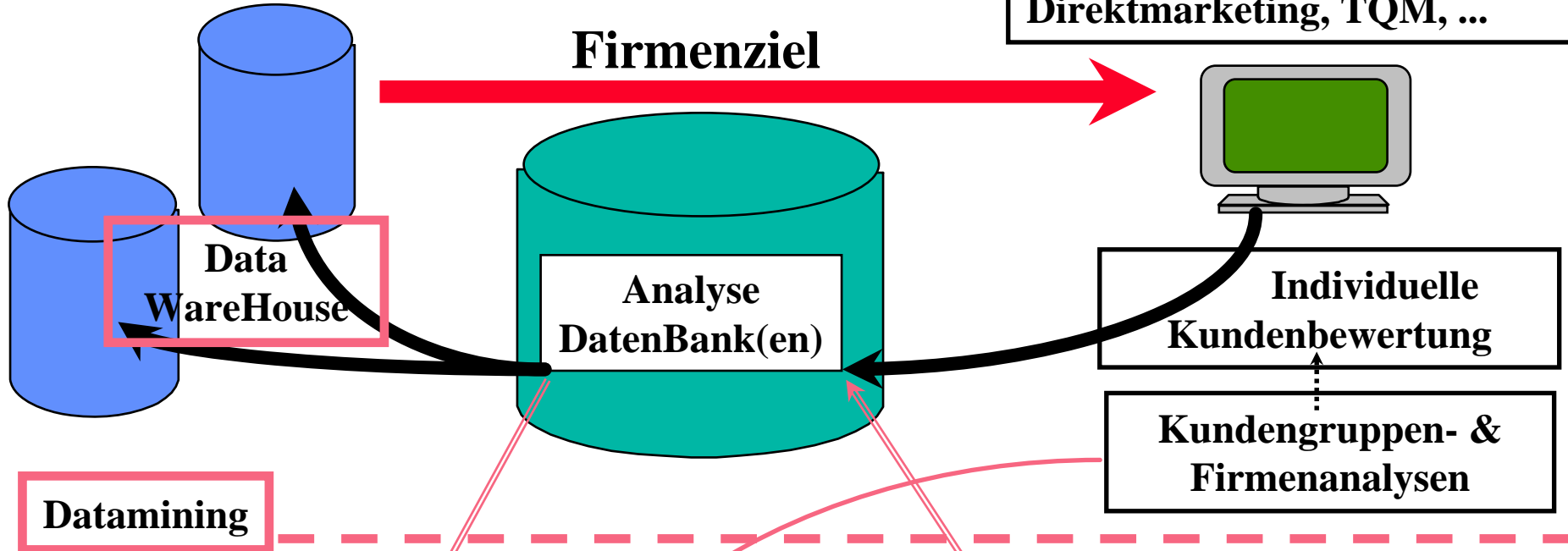


# Datamining im Data Warehouse

Kundenkontakt & Produktion

Direktmarketing, TQM, ...

Firmenziel



Datamining

Stichproben Trainings- & Prüfdaten  
Datenaufbereitung

Visuelle Analysen  
SAS/Insight  
SAS/Spectrview

1. Automatische lineare Regression  
(Automatische Response Analyse  
ARA Application SASD)

2. Kundensegmentierung  
(SAS: Tree Application)

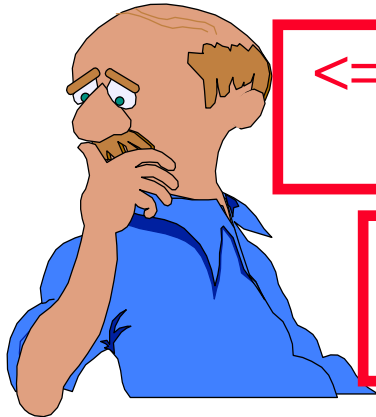
3. Neuronale Netze (SAS NNA)

4. Trendberechnungen (SAS/ETS)



# Datamining Schritt 1: Konzept und Ziele

## Kundenanalyse mit Punktebewertung - Beispiel Marketingaktion Testversand

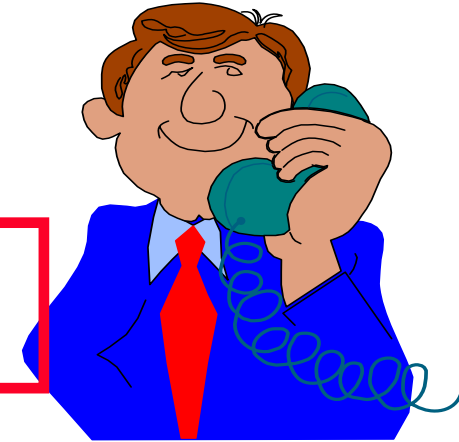


**$\leq 0$  Punkte = antwortet nicht**

z.B. 49500 Kunden

**Kunde antwortet: 1000 Punkte  $\Rightarrow$**

z.B. 500 Kunden



**Welche Faktoren ???** (beeinflussen das Kundenverhalten?)

oder besser: Aus welchen Daten kann man das Kundenverhalten  
**vorausberechnen?**

**Kundendaten** z.B.: Alter, Geschlecht, Beruf, Bildung,  
Kinderzahl, Wohnungsgrösse, Umsätze in verschiedenen  
Warengruppen, Umsätze zu bestimmten Zeiten, .....



# Nutzen durch Datenanalyse & Datamining

- **Katalogwerbung 1 Million Kunden**
- **Werbebrief 10 DM. Antworten erbringen ca. 1000 DM Deckungsbeitrag. Die Antwortquote steige von 1% auf 2% bei Selektion der 10% “besten” Kunden.**
- **Gewinn ohne Selektion: 1 Mio Werbebriefe: Kosten = 10 Mio DM. Antworten 1% = 10.000Kunden=> \*1000DM => Ertrag = 10 Mio DM. Gewinn = 0 DM**
- **mit Selektion durch Datenanalyse: 10% von 1 Mio = 100000 Werbebriefe: Kosten= 1 MioDM. Antwort 2%= 2000 => \* 1000DM=> Ertrag= 2 Mio. Gewinn = 1 Mio**

# SEMMA - die SAS<sup>®</sup> Data Mining Methodik

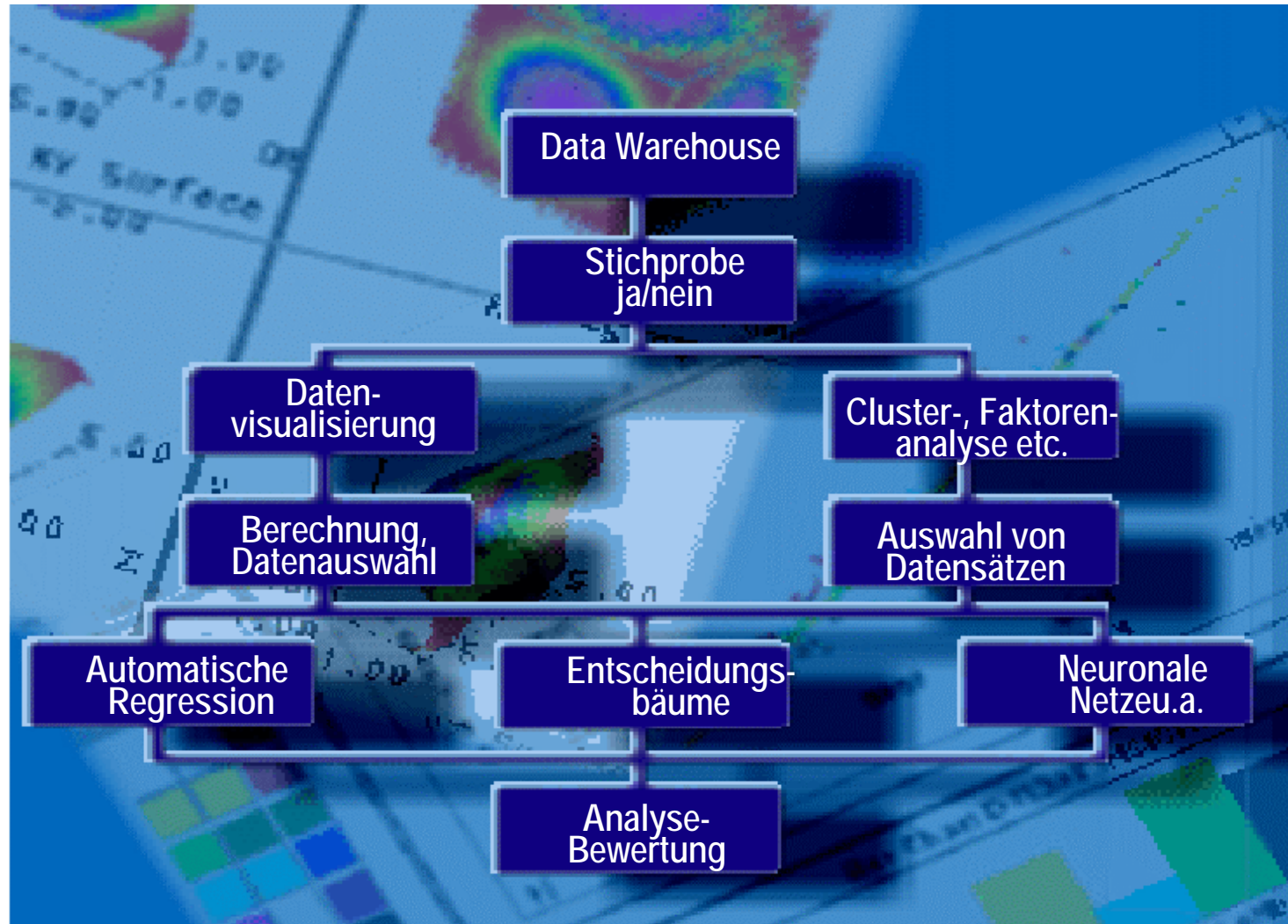
Stichproben

Exploration

Modifikation

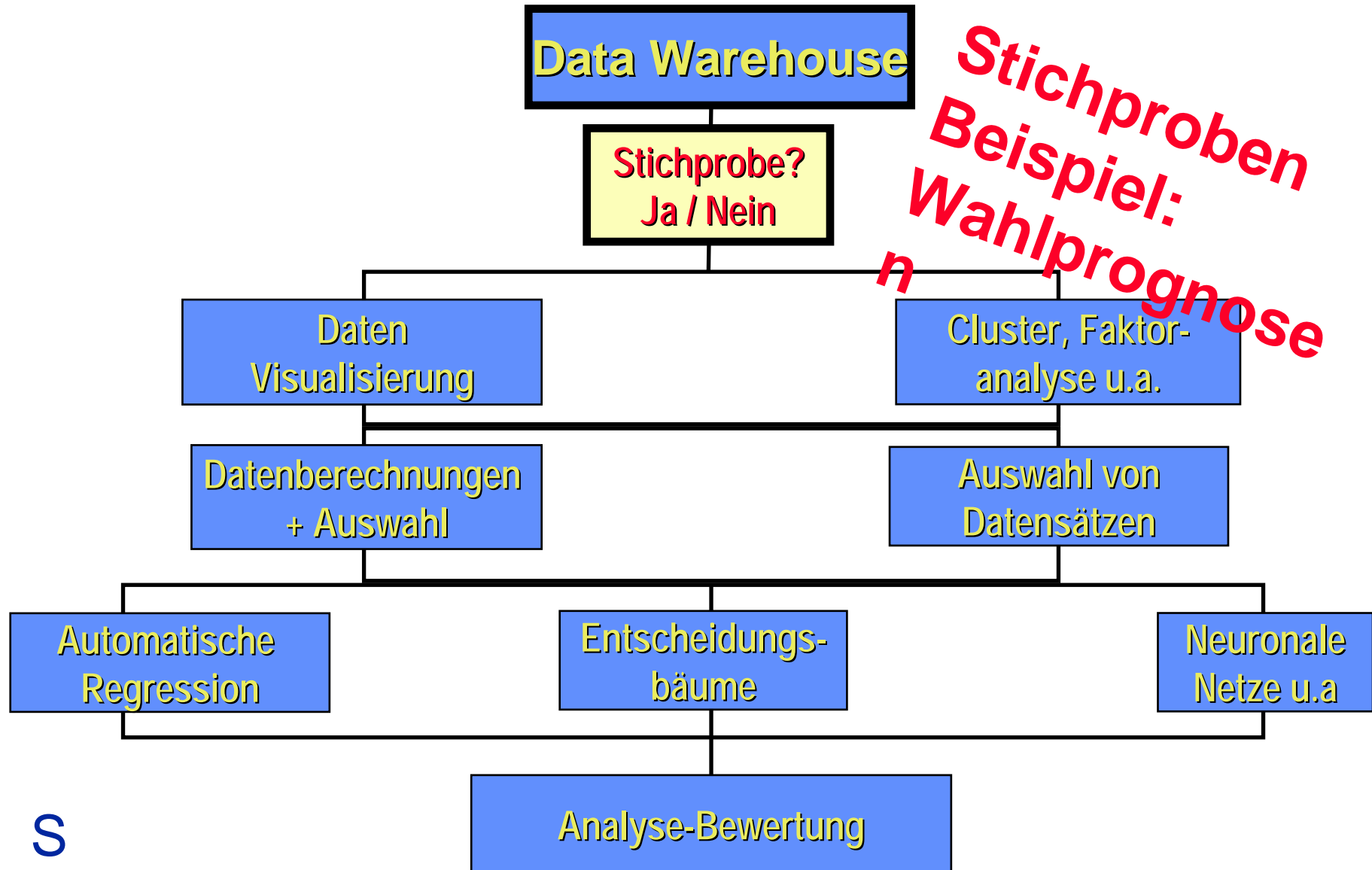
Modellierung

Auswertung



# SEMMA - die Datamining Technologie 1

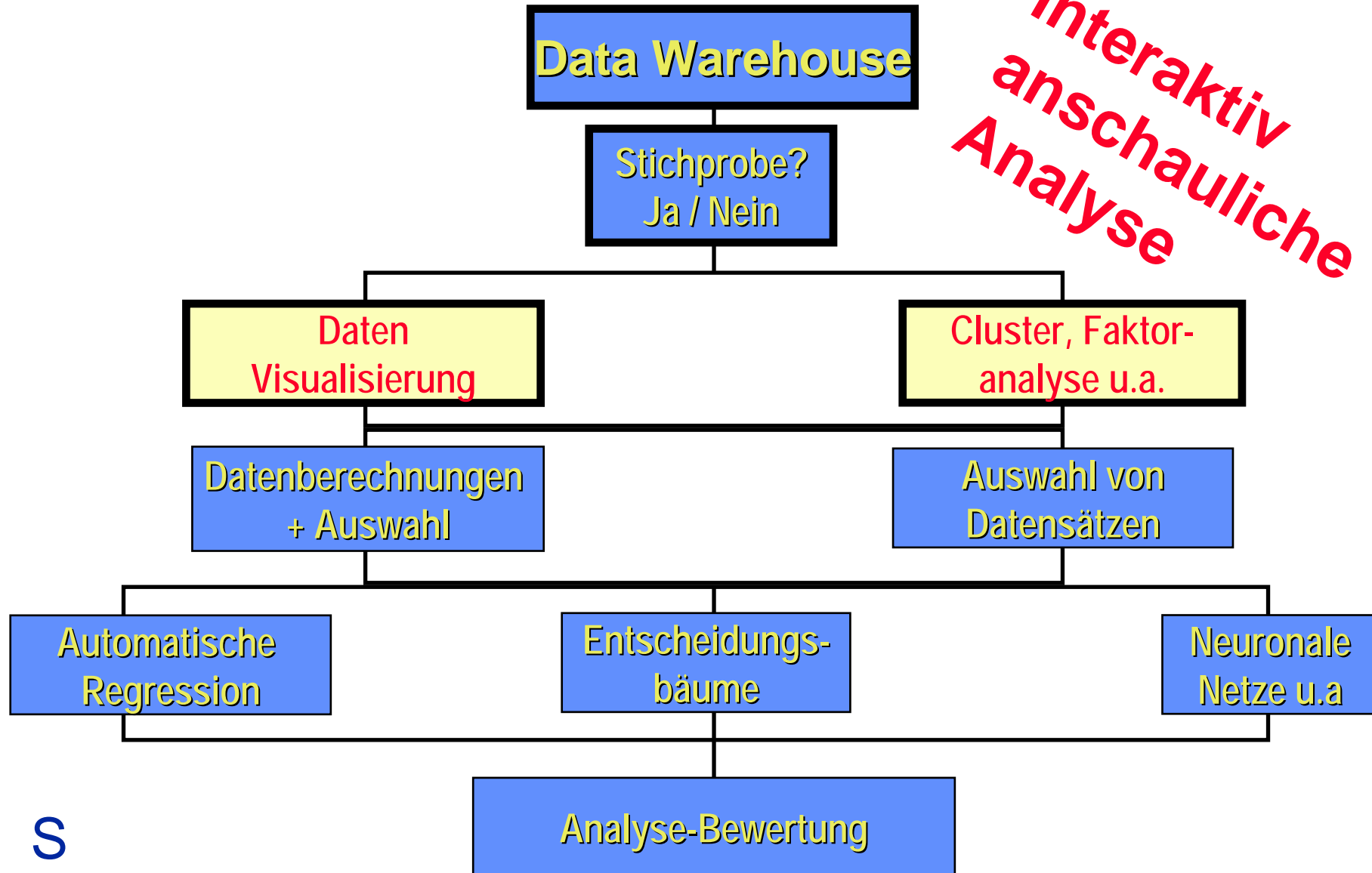
## STICHPROBEN



# SEMMA- die Datamining Technologie 2

## EXPLORATION

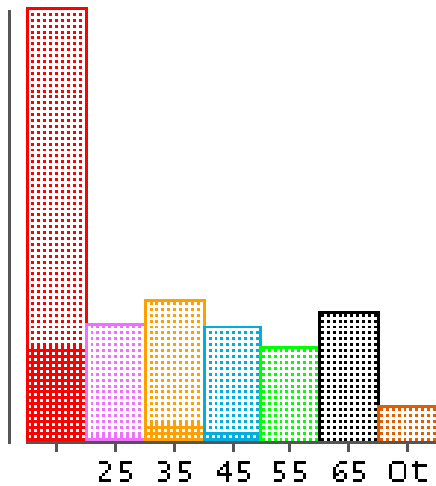
*interaktiv  
anschauliche  
Analyse*



# Datamining Exploration: Visuelle Analyse

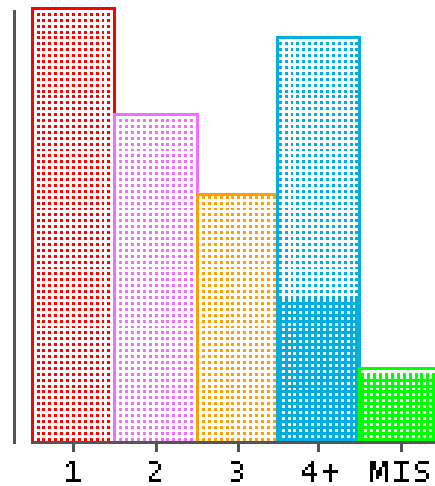
## Databased Marketing Testaktion

F  
r  
e  
q  
u  
e  
n  
c



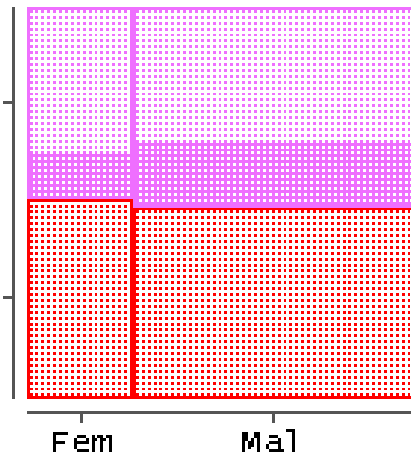
ALTER

F  
r  
e  
q  
u  
e  
n  
c



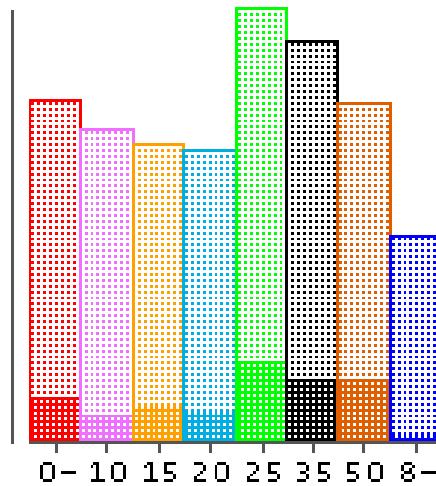
HHGROESS

R  
E  
S  
P  
O  
N  
S



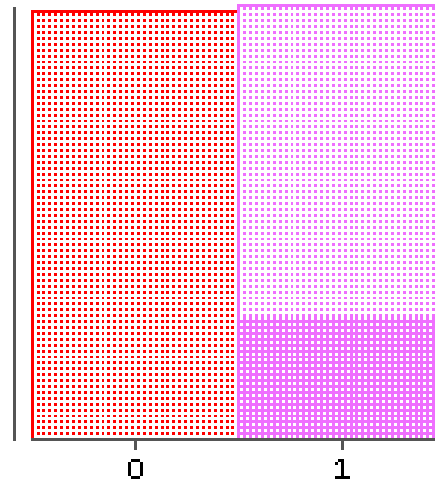
GESCHL

F  
r  
e  
q  
u  
e  
n  
c



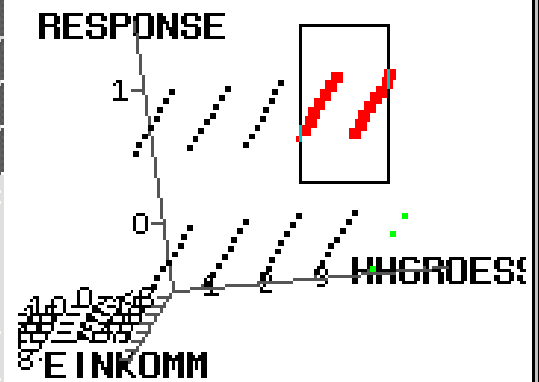
EINKOMM

F  
r  
e  
q  
u  
e  
n  
c



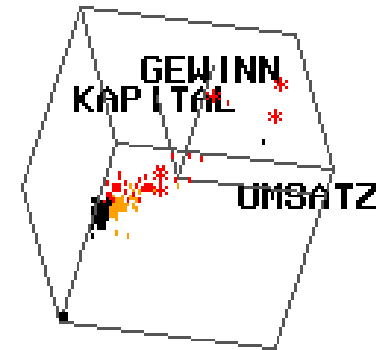
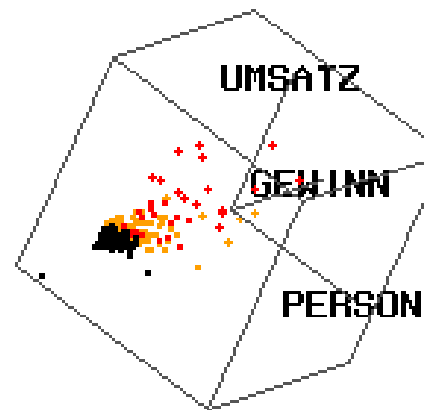
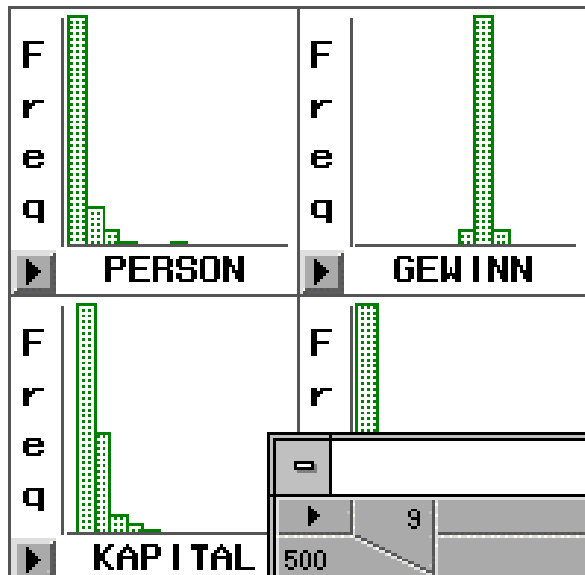
RESPONSE

R  
E  
S  
P  
O  
N  
S

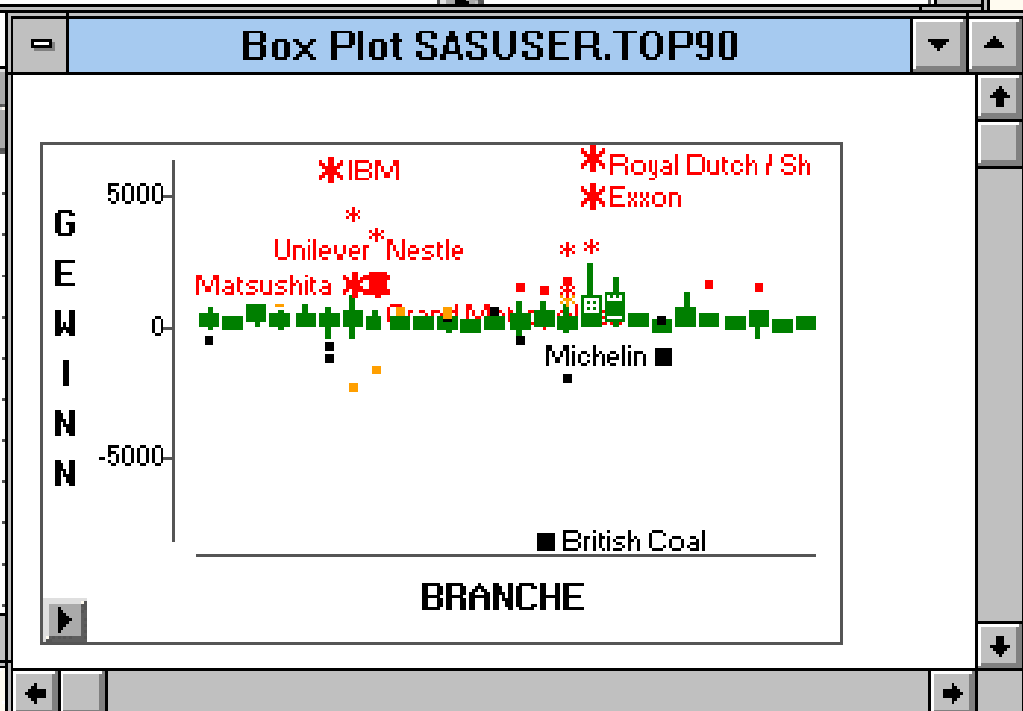


# Datamining Exploration: Visuelle Analyse

## Unternehmensanalyse



9		LAND
* 322	U.S.	
■ 323	U.S.	
■ 324	U.S.	
■ 325	Italy	
* 326	Italy	
■ 327	New Zealand	
* 328	U.S.	
■ 329	Japan	
■ 330	Britain	
■ 331	U.S.	
* 332	U.S.	



Tools

# Datamining Exploration: Visuelle Analyse Erkennung von Konzentrationsmustern

Planes Volume Render

LEVEL Plane => Auto

LATITUDE Plane => Auto

LNGITUDE Plane => Auto

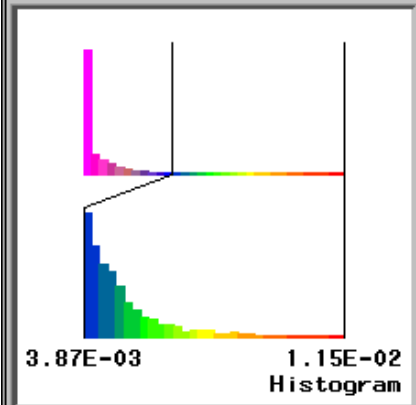
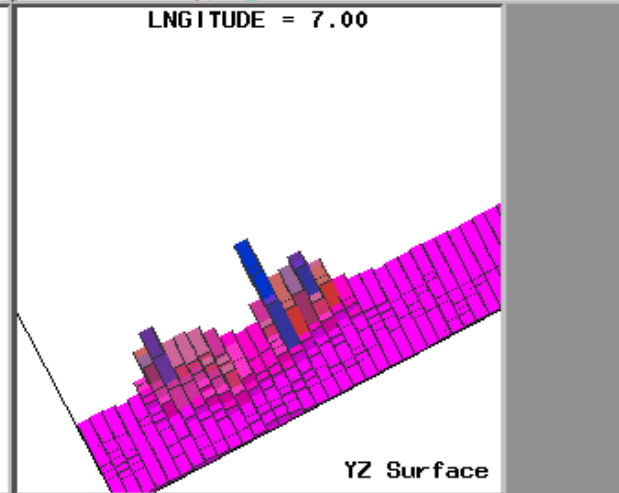
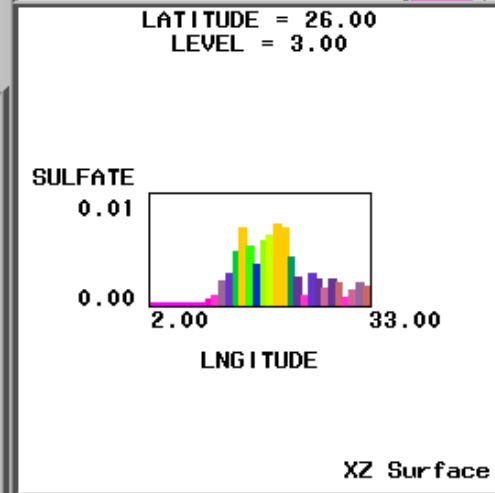
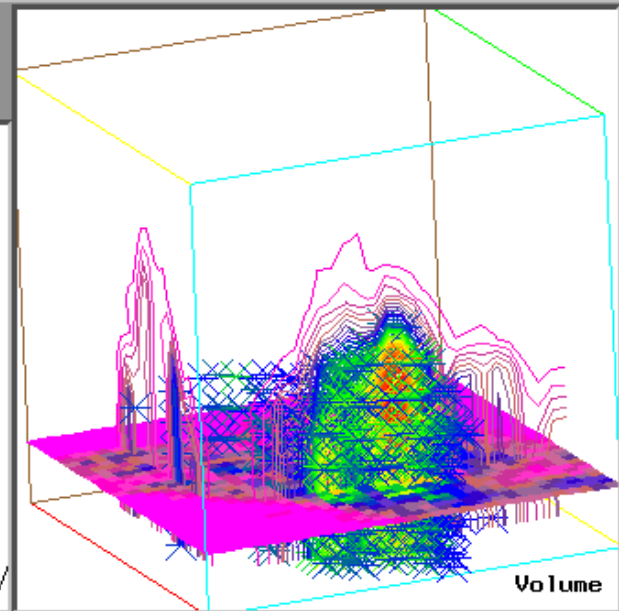
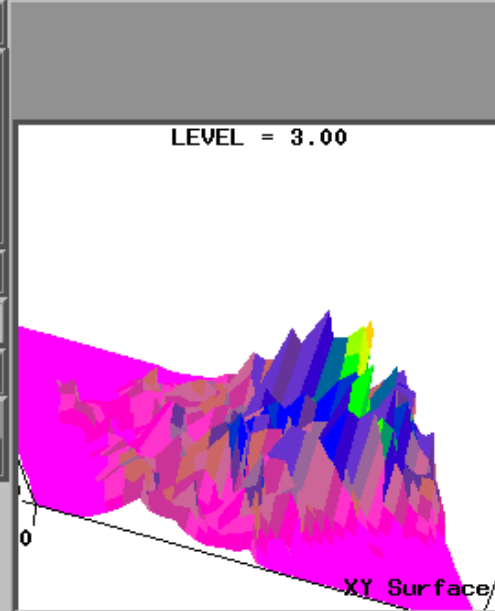
NOR plane => Auto

Point cloud Isosurface

Off On

X \* 0 +

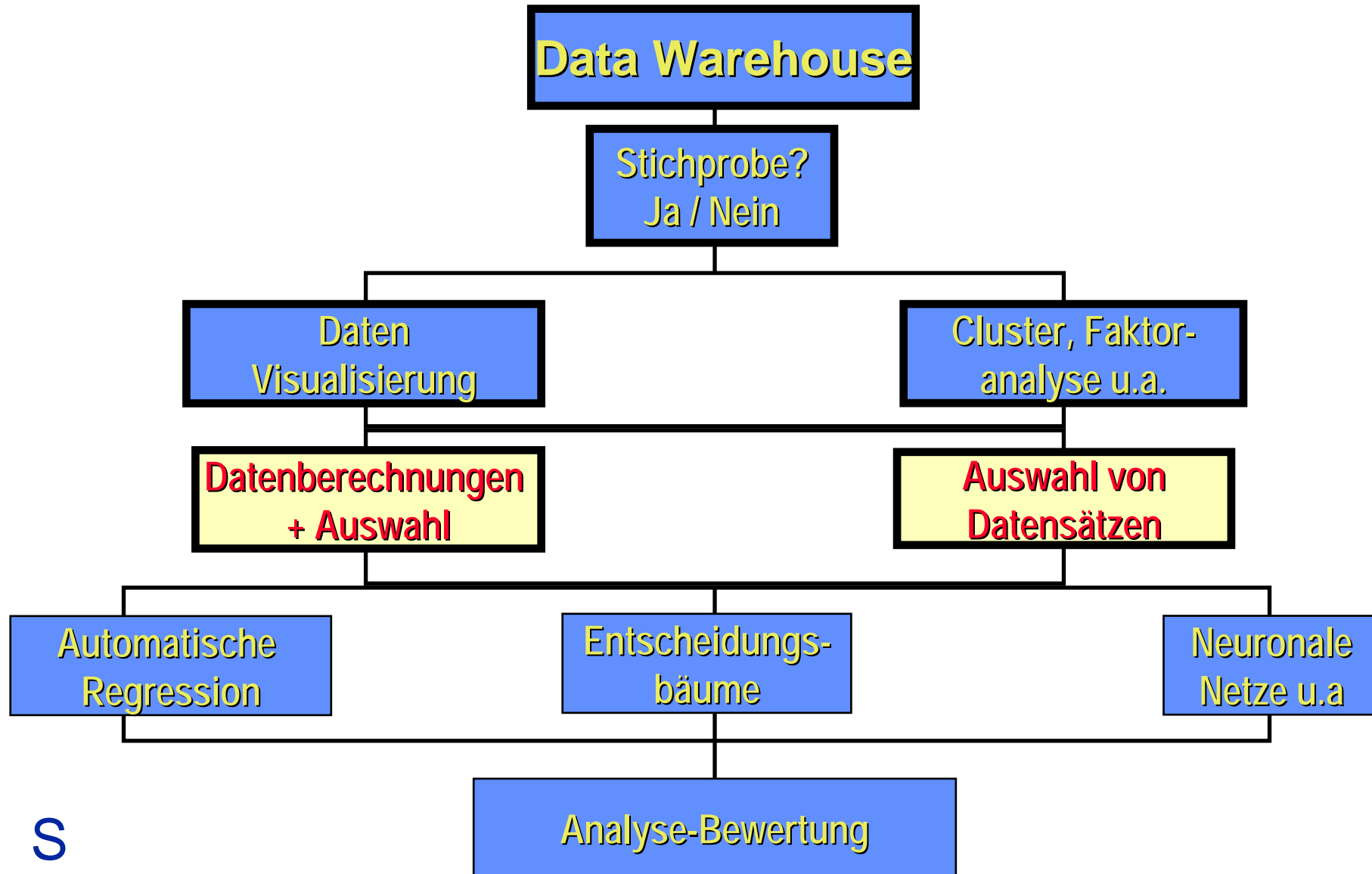
MARKER SIZE





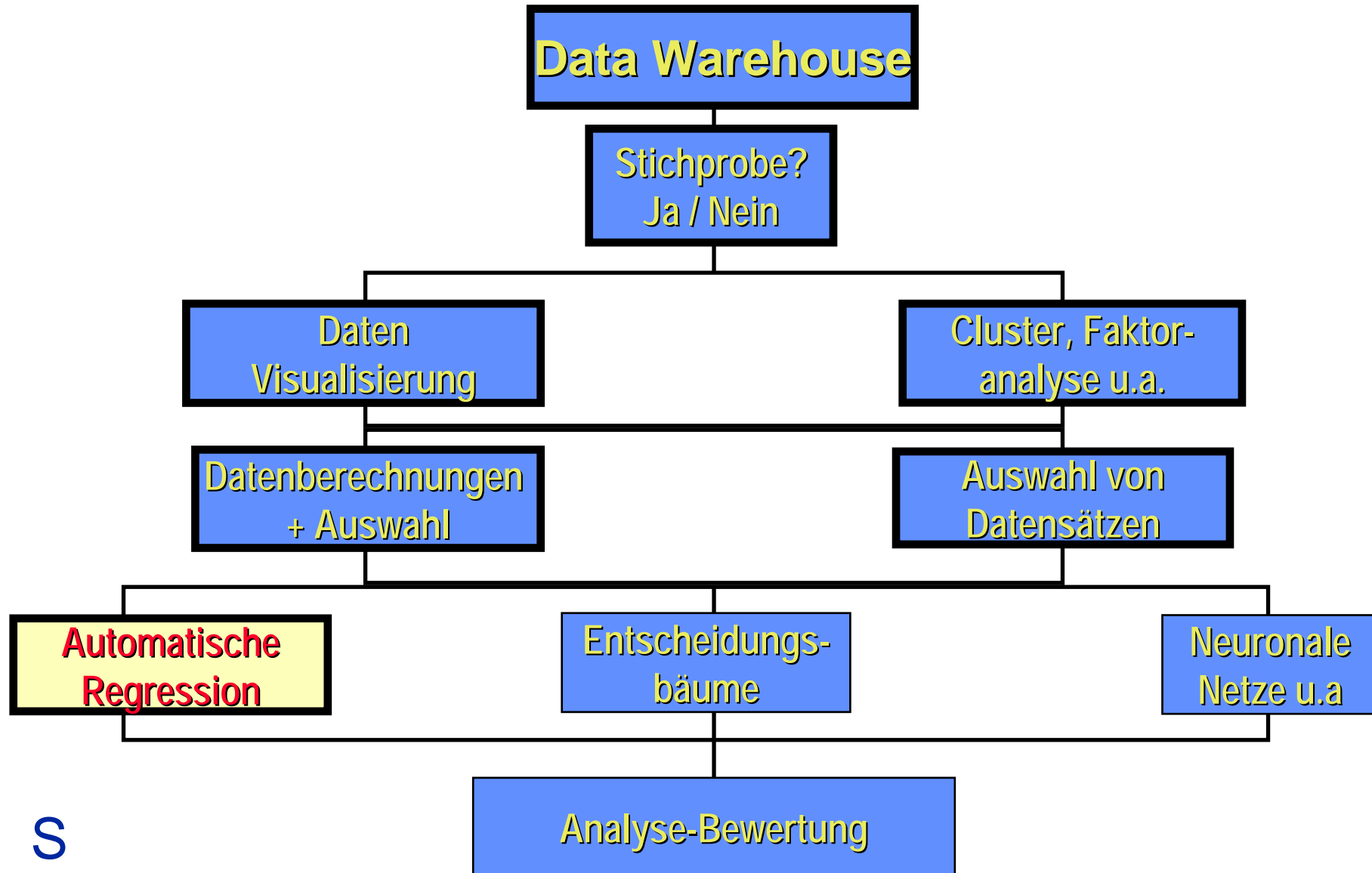
# SEMMA- die Datamining Technologie 3

## MANIPULATION



# SEMMA- die Datamining Technologie 4

## MODELLWAHL 1. Automatische Regression



# Datmining Modellwahl: mit linearer und logistischer Regression für Scoring

Response Scores

Minimum score: 0  
Maximum score: 1000

Base score: 440

New cutpoint : 547

CREDIT	HHSIZE	OCCUP
NO	1	.
YES	2	Blue Collar
	3	Other
	4+	White Collar

Automatische Auswahl wichtiger Merkmale  
und individuelle Punktebewertung (Scoring)

Splitnumber : 6

Goback

# Gütematrix aller Datamining-Verfahren mit Zielvariable: Regression, Entscheidungsbäume, Neuronale Netze

Training Data

Test Data

Actual Values

Calculated

NonResponse 0

Response 1

**70**

**30**

**35**

**65**

**78**

**22**

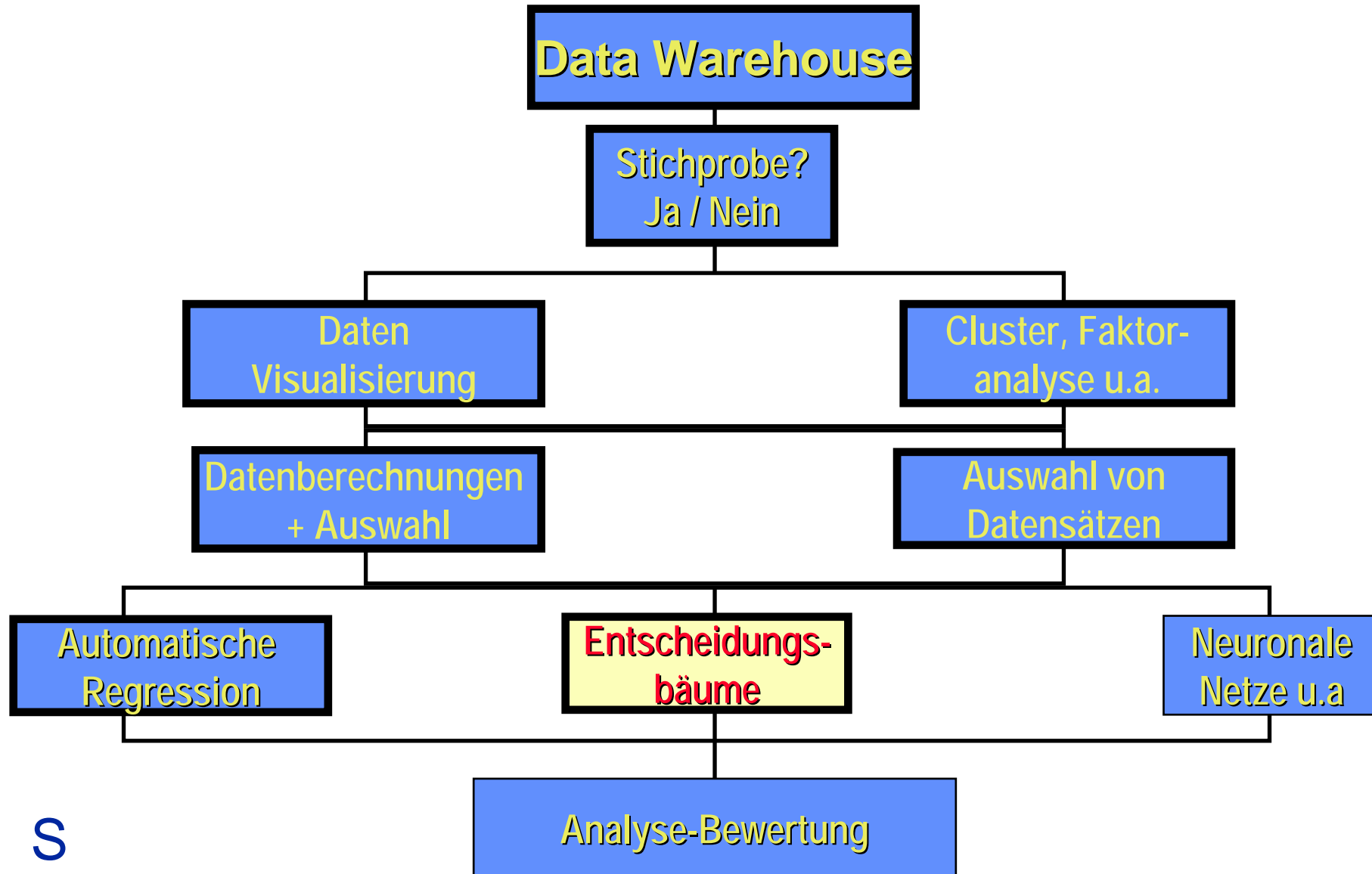
**40**

**60**

Goback

# SEMMA- die Datamining Technologie 4

## MODELLING 2. Entscheidungsbäume



# Datamining Modellwahl

## Entscheidungsbäume Treeanalysis

Split : RESPONSE  
 Values : NonResp Response  
 %of Node : 78% 22%  
 %of All : 78% 22%

Node : 3  
 Split : HHSIZE  
 Values : 1 2  
 %of Node : 84% 16%  
 %of All : 53% 10%

Node : 7  
 Split : INCOME  
 Values : 0-<8 10-<15 15-<20 20-<25  
 %of Node : 86% 14%  
 %of All : 33% 5.5%

REGION  
 N  
 % 19%  
 % 2.3%

Node : 11  
 Split : REGION  
 Values : W E  
 %of Node : 88% 12%  
 %of All : 24% 3.2%

Node : 13  
 Split : AGE  
 Values : MISSING 40-55 55+  
 %of Node : 75% 25%  
 %of All : 5.5% 1.8%

Node : 14  
 Split : HOMEOWN  
 Values : N  
 %of Node : 85% 15%  
 %of All : 15% 2.6%

Node : 16  
 Split : GENDER  
 Values : Male  
 %of Node : 80% 20%  
 %of All : 7.0% 1.7%

Node : 17  
 Split : GENDER  
 Values : Female  
 %of Node : 90% 10%  
 %of All : 7.8% 0.87%

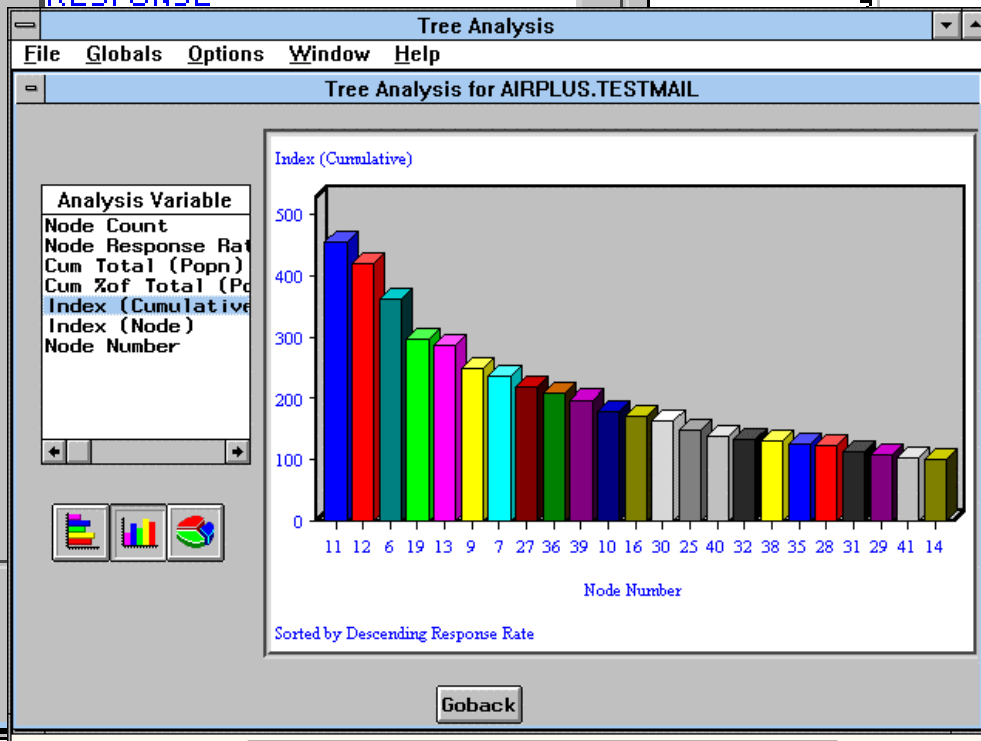
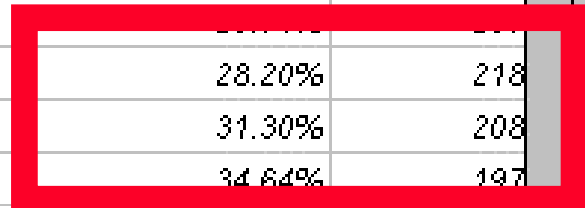
# Datamining Gewinnoptimierung durch Kundensegmentierung

Datasets	
AIRPLUS.TESTMAIL	↑
AIRPLUS.TESTSAMP	
MARKET.MAILSAMP	
MARKET.MAILSHOT	↓

Dependent Variable

RESPONSE

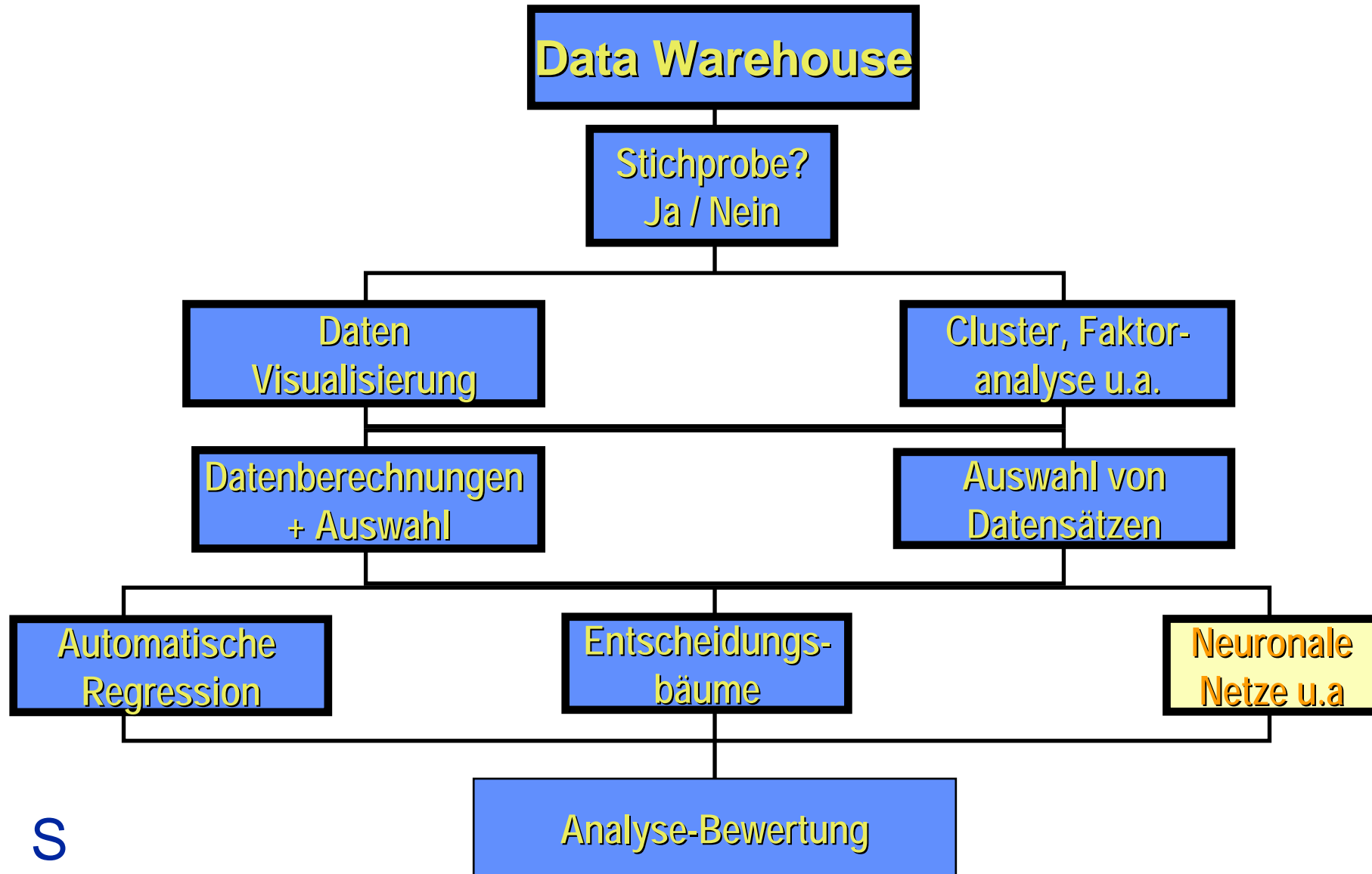
Node	Cum Total (Popn)	Cum %of Total (Popn)	Index (Cum)
11	85	1.21%	454
12	157	2.23%	419
6	490	6.95%	361
19	898	12.73%	295
13	1005	14.25%	287
9	1511	21.43%	249
7	1672	23.71%	237
1989	1989	28.20%	218
2207	2207	31.30%	208
2443	2443	34.64%	197
2983	2983	42.30%	178
3197	3197	45.33%	171
3442	3442	48.81%	164
4000	4000	56.99%	149



Produce Graph

# SEMMA- die Datamining Technologie 4

## MODELLING: 3. Neuronale Netze





# Datamining Modellwahl

## Neuronale Netze

### Training Data Set

SASUSER.SALES



### Inputs

REGION  
AMOUNT  
VISITS  
TIME

### Outputs

SALES

*Fine Tune**Start training*

### Compute Server

Local host



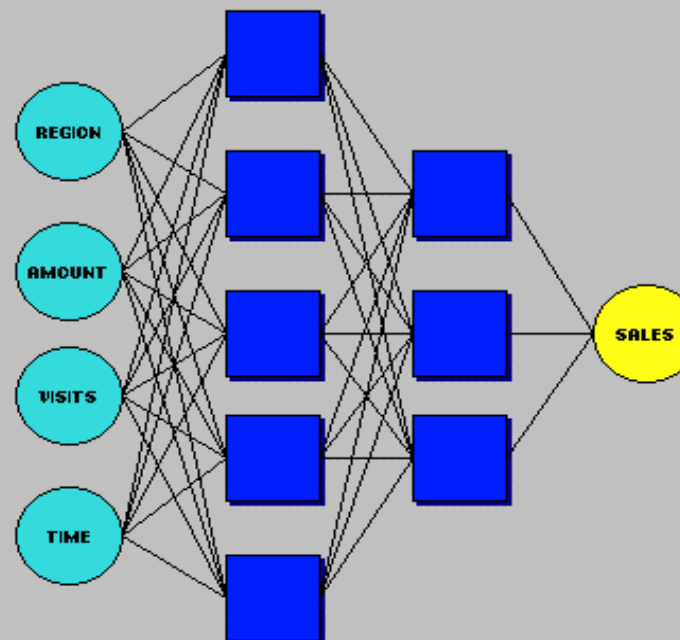
```
Random seed: 250
Preliminary runs: 5
RMSE: 0.04530847702865
RASE: 0.04173510495986
RFPE: 0.04861992590801
SBC : -377.389923537408
Finish: 06MAY96:12:30:12
```

### Architecture

Multilayer perceptron

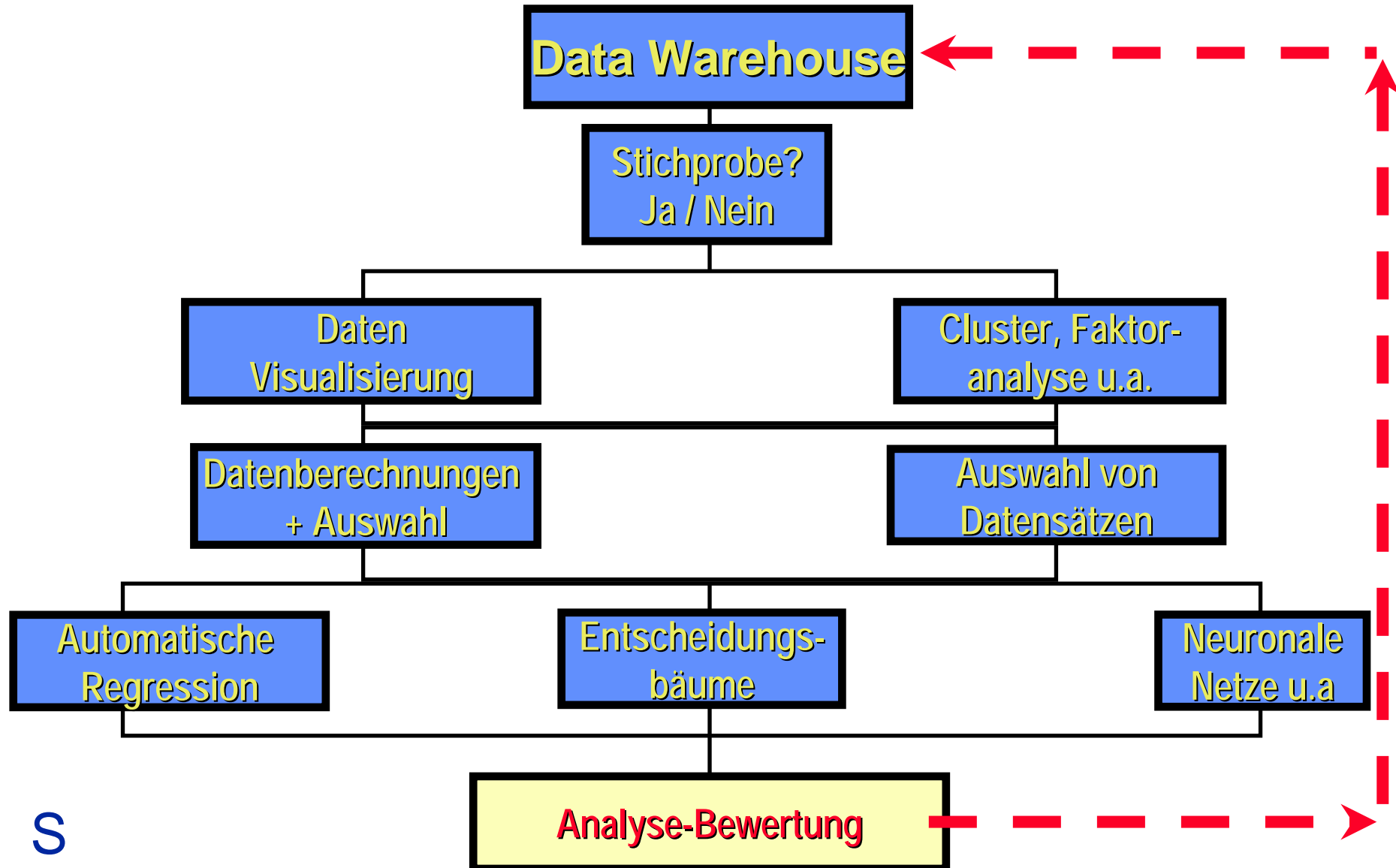


Diagram

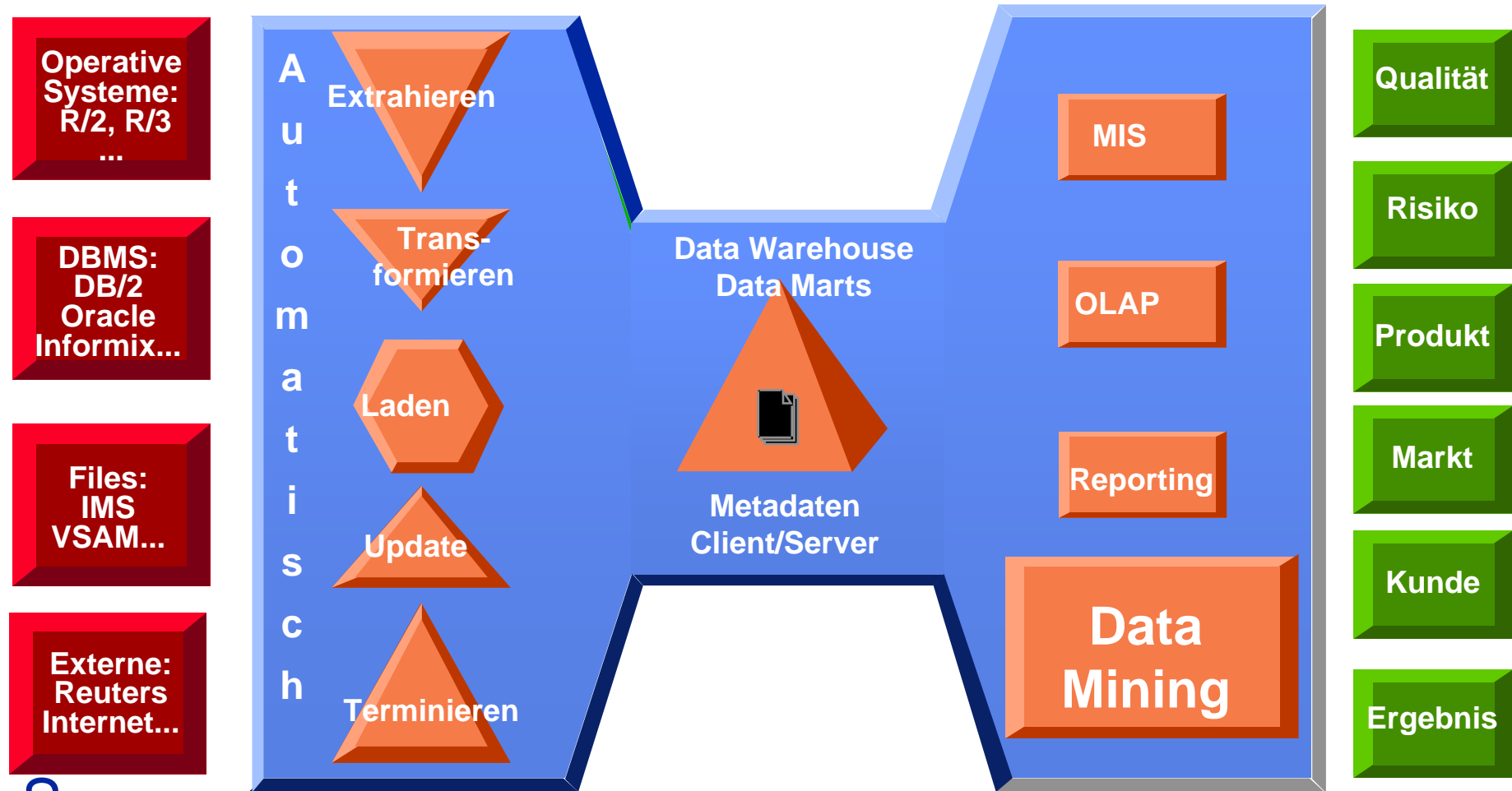


# SEMMA- die Datamining Technologie 5

## ANALYSE-BEWERTUNG



# Datamining im offenen Data Warehouse



S