



# Moderne Data Mining –Techniken und -Anwendungen

Dr. Reinhard Strüby

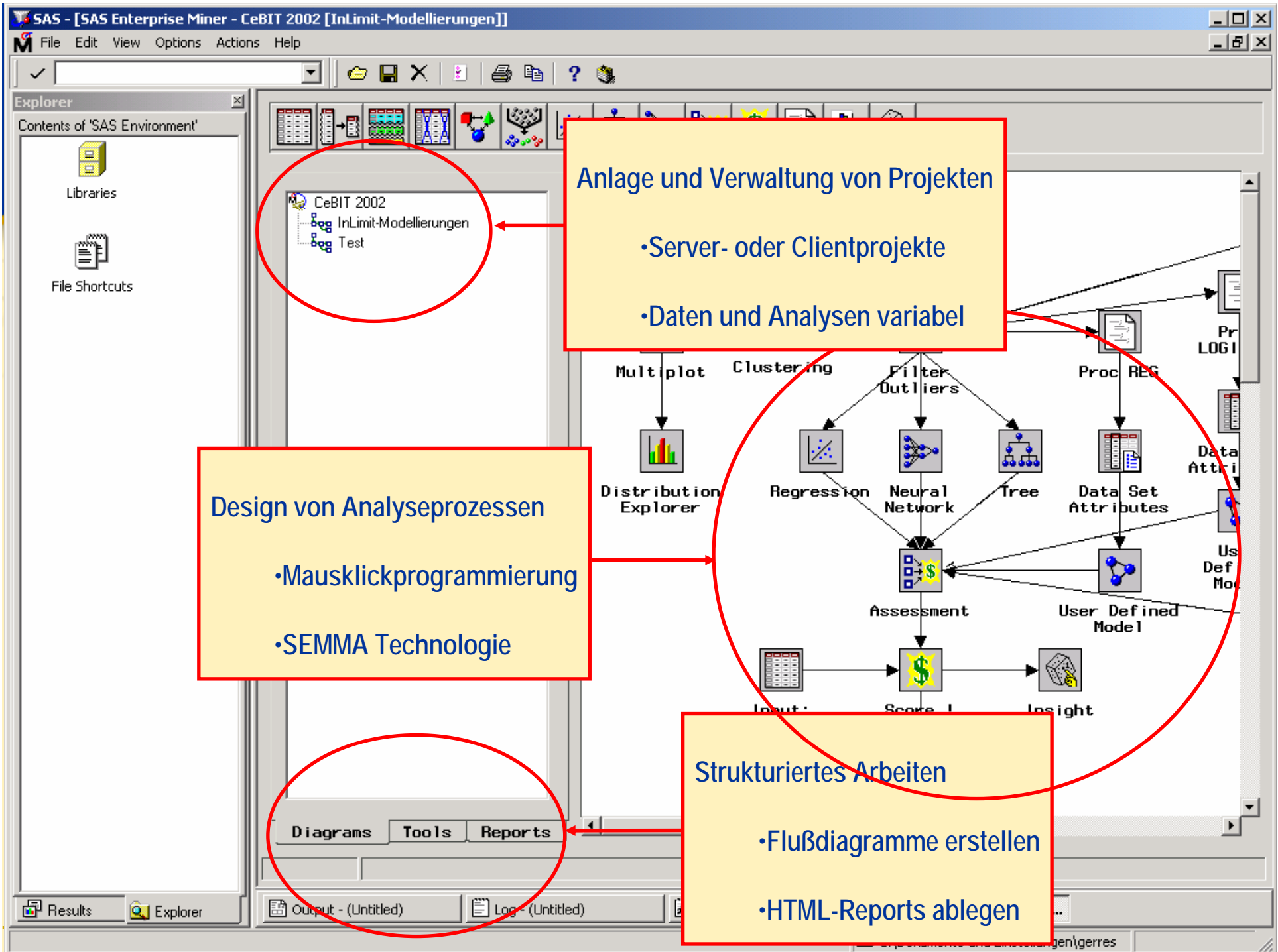
Business Competence Center

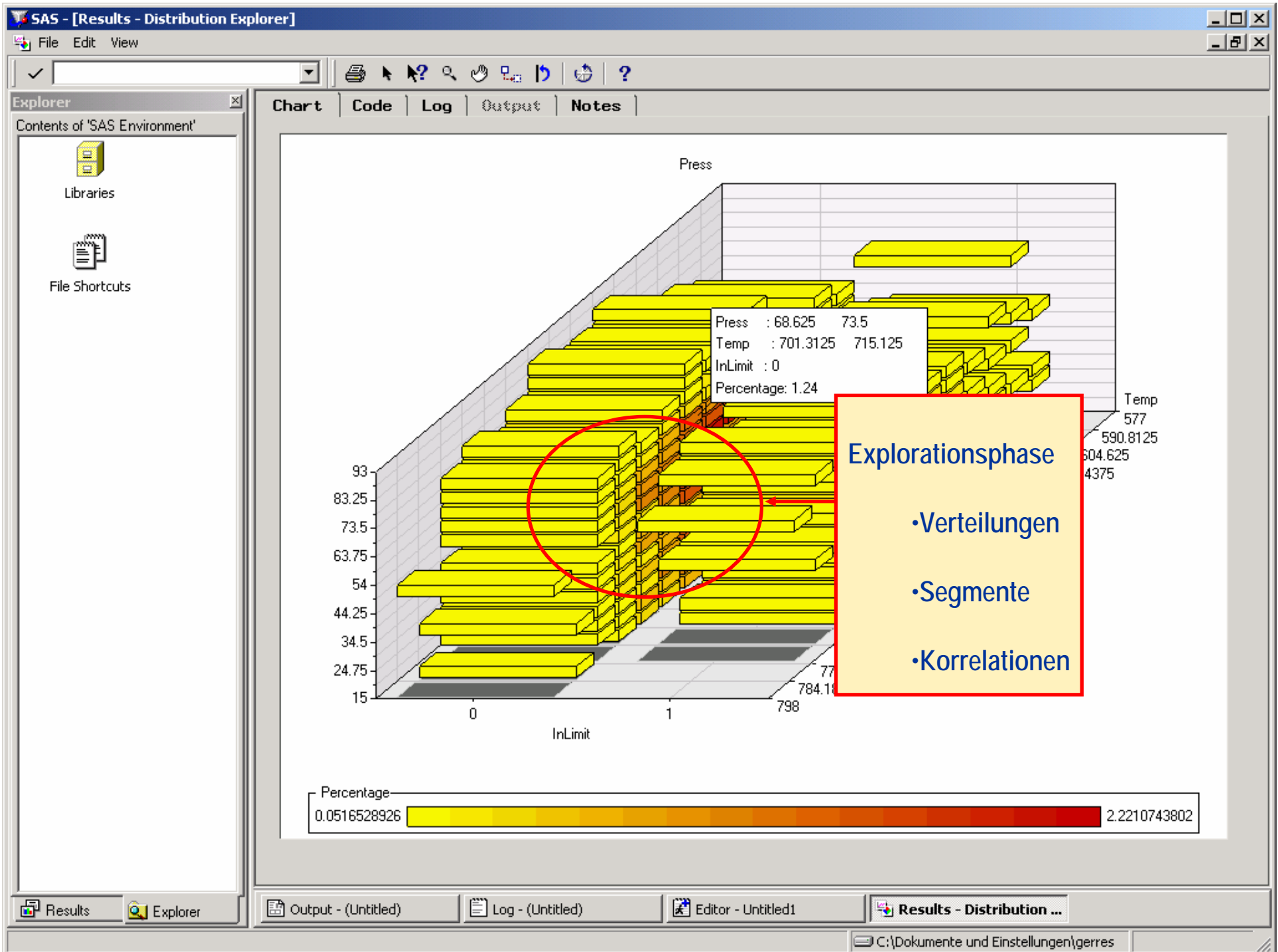
SAS Deutschland

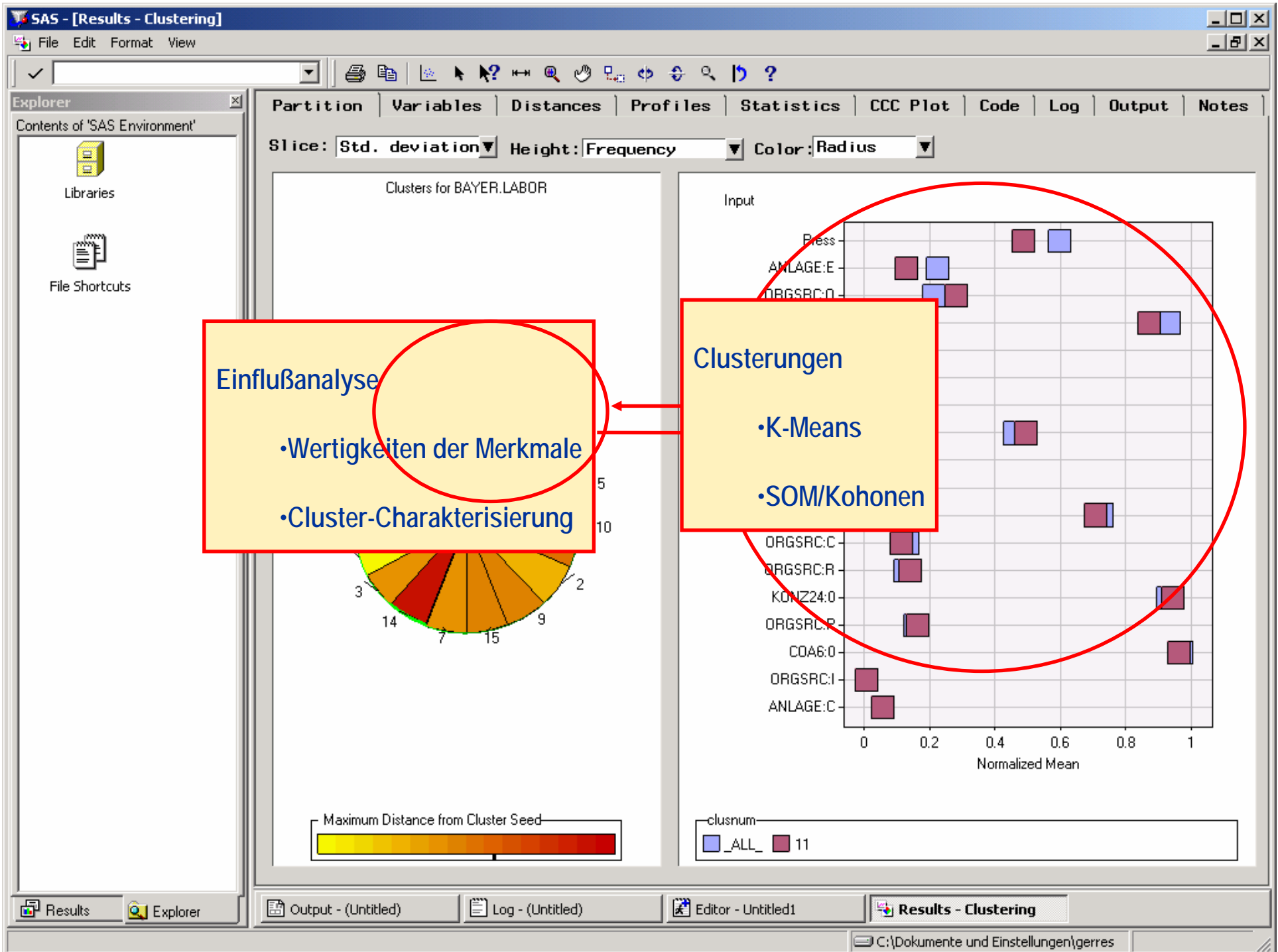
*The Power to Know.*

# SAS Data Mining: **SEMMA** Technologie

- **S**tichproben
  - Zufällig, geschichtet
- **E**xploration
  - Grafisch, interaktiv
- **M**odifikationen
  - Ausreißer, Transformationen, Imputationen
- **M**odellierungen
  - Regressionen, Entscheidungsbäume, Neuronale Netze
- **A**ssessment
  - Liftcharts, Profitcharts





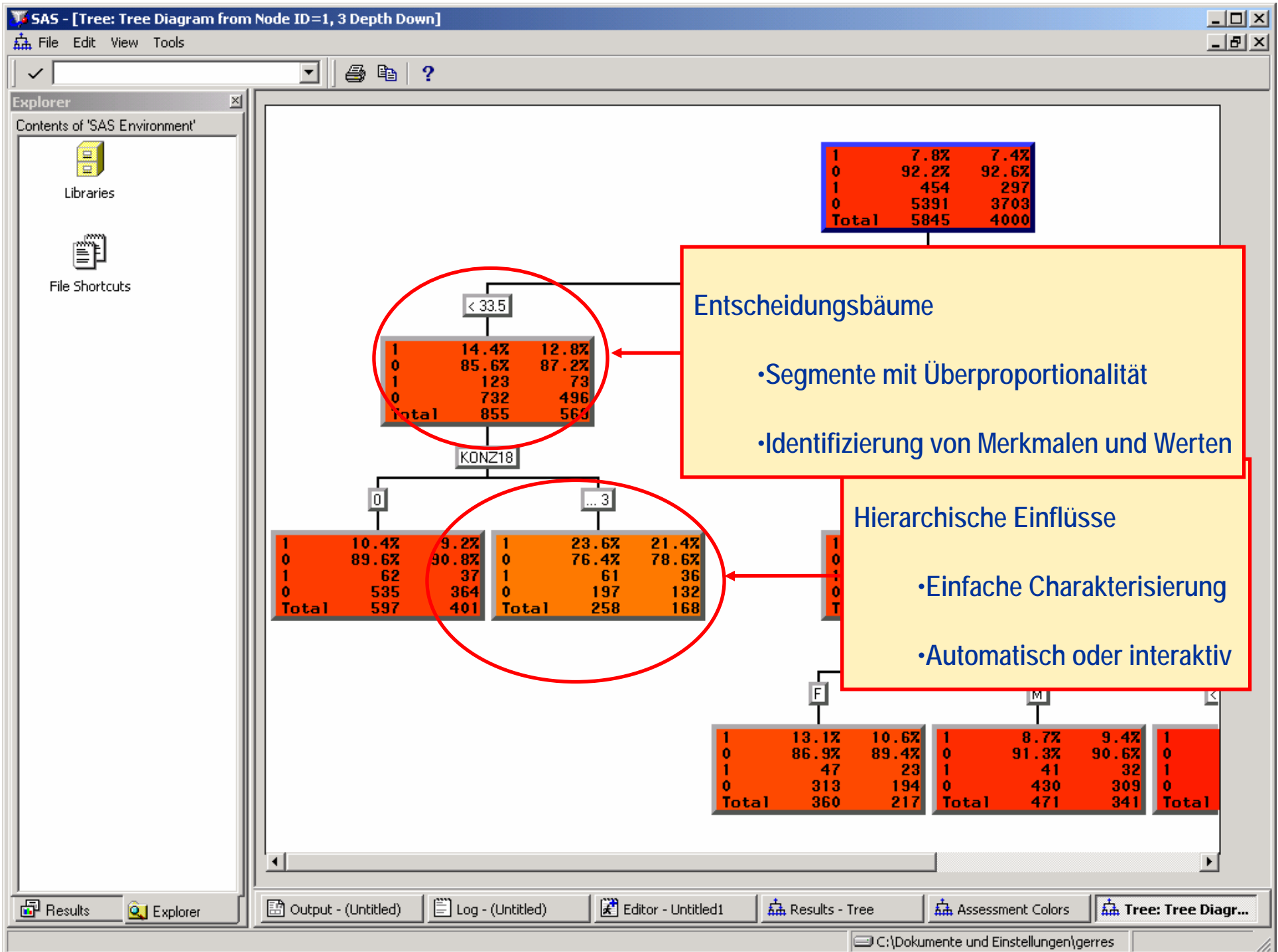


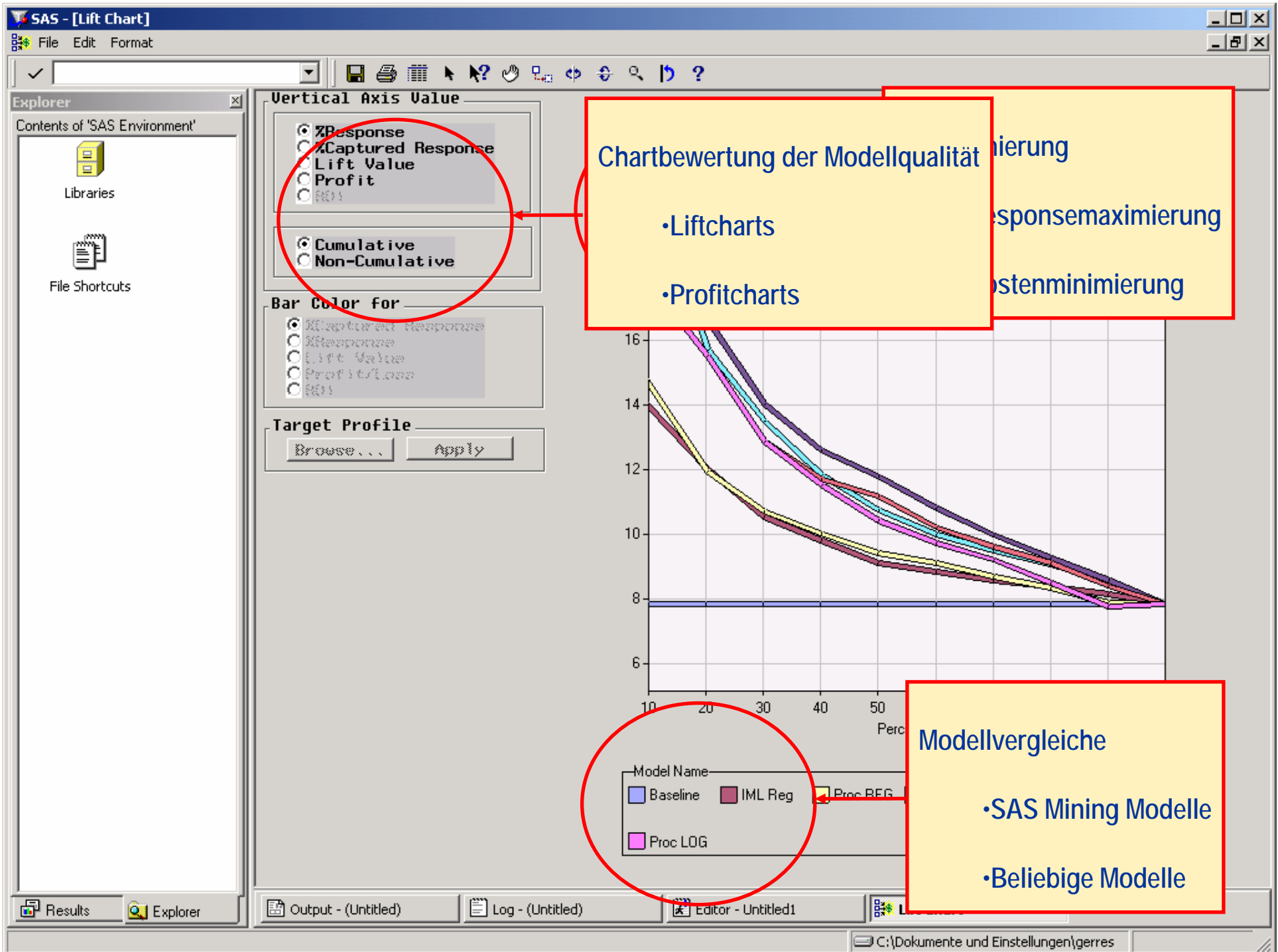
### Einflußanalyse

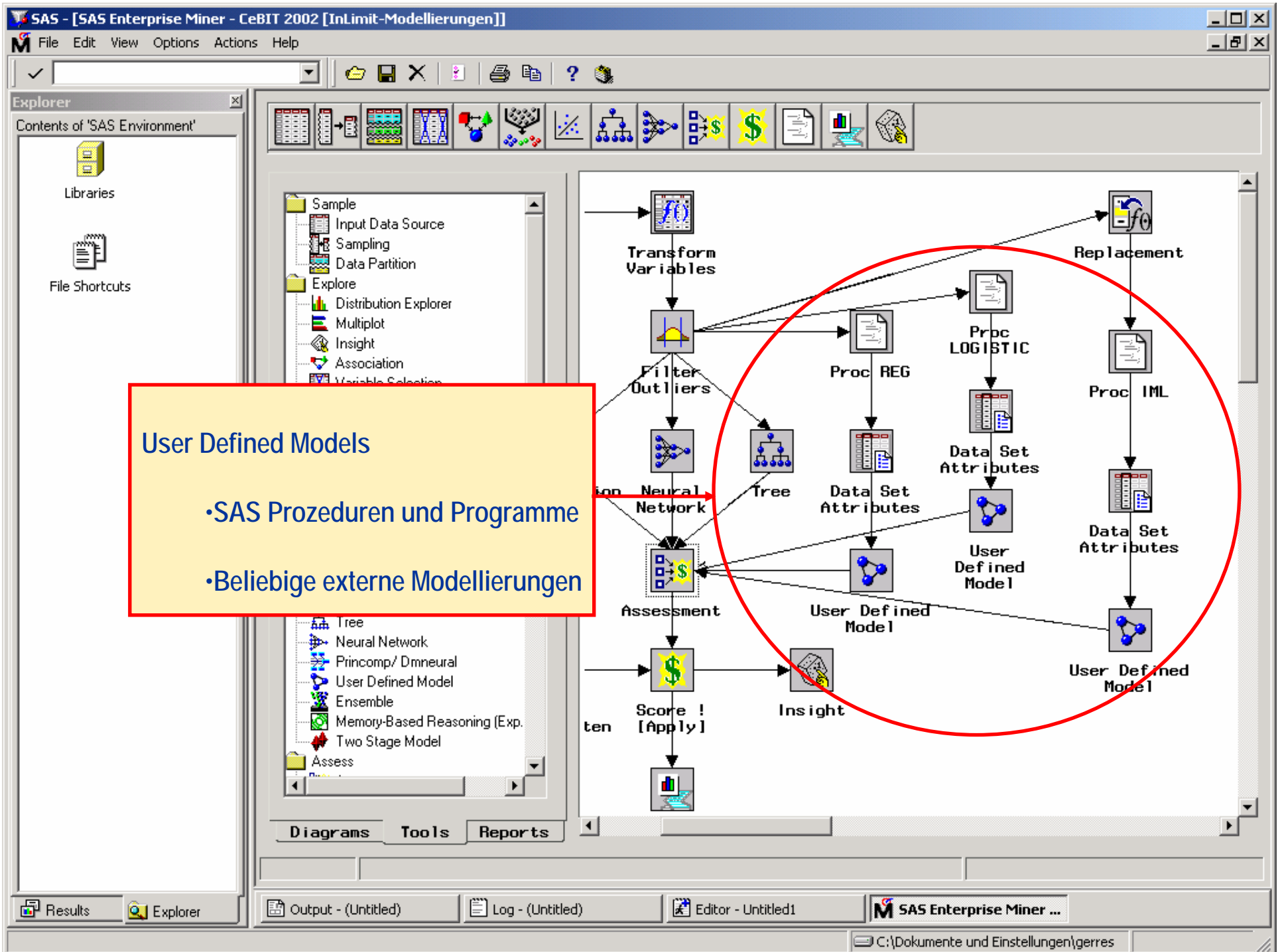
- Wertigkeiten der Merkmale
- Cluster-Charakterisierung

### Clusterungen

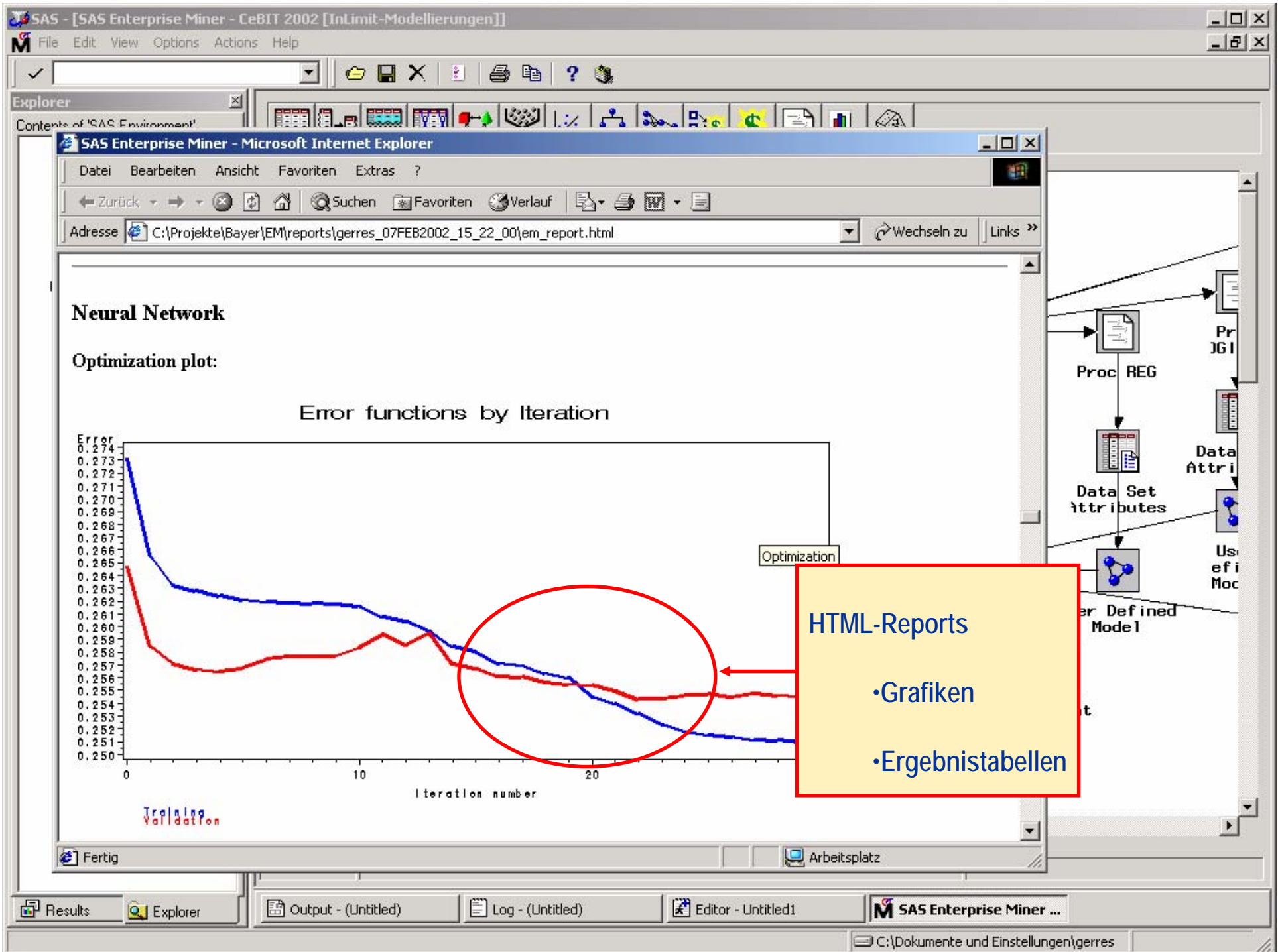
- K-Means
- SOM/Kohonen

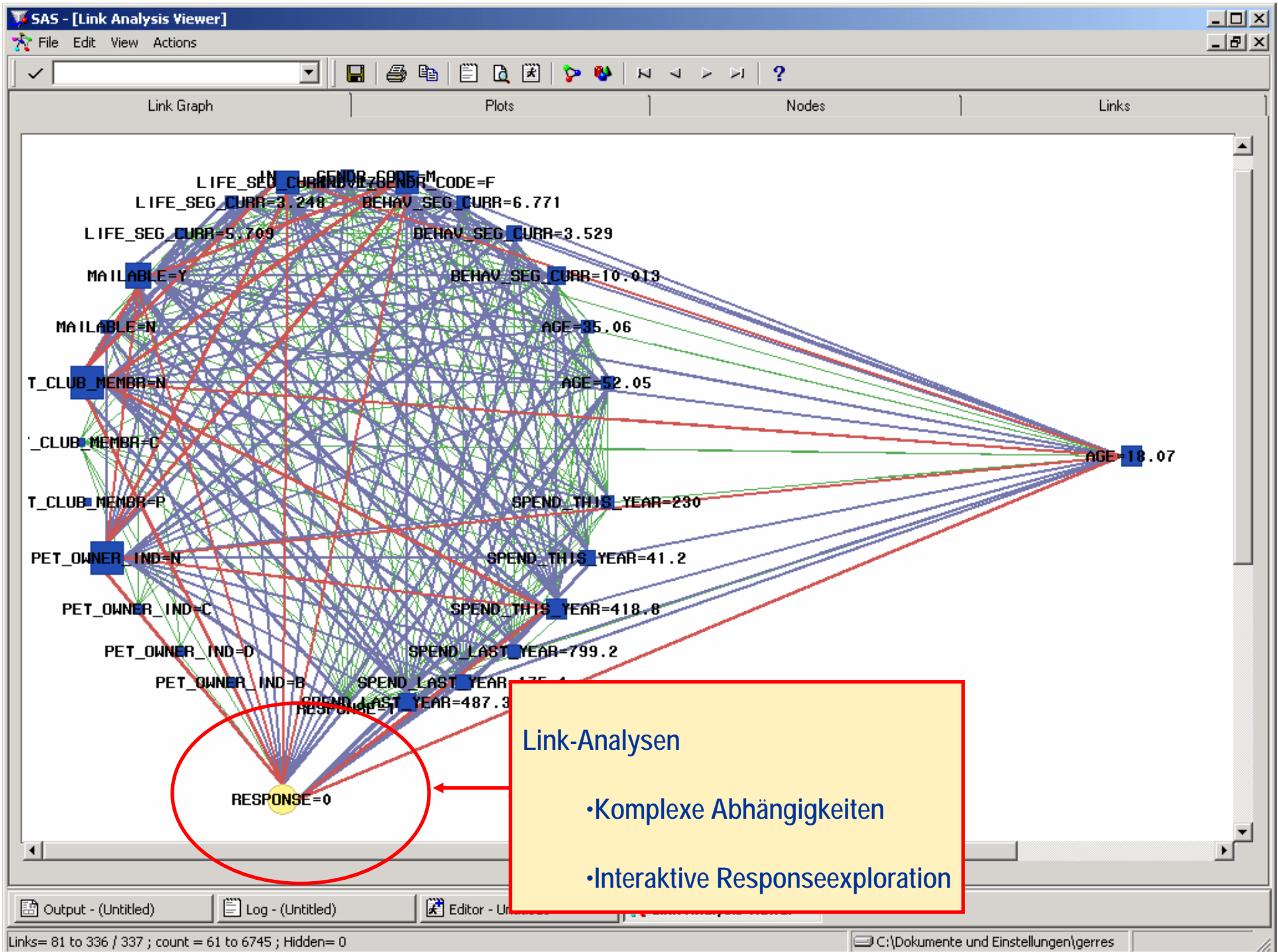










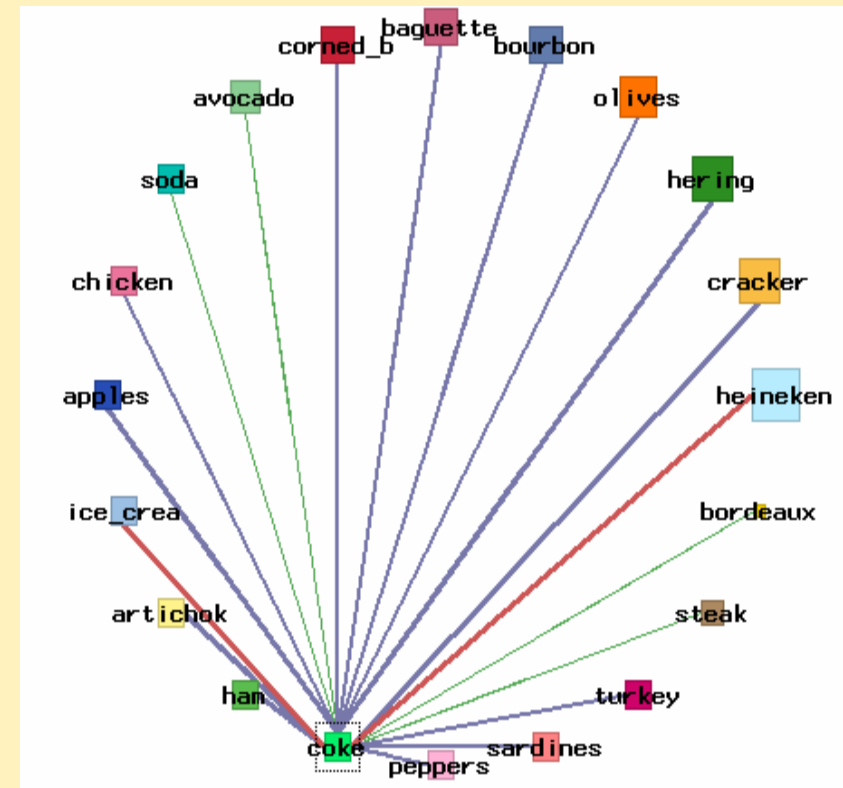


# Assoziationen: Warenkorb

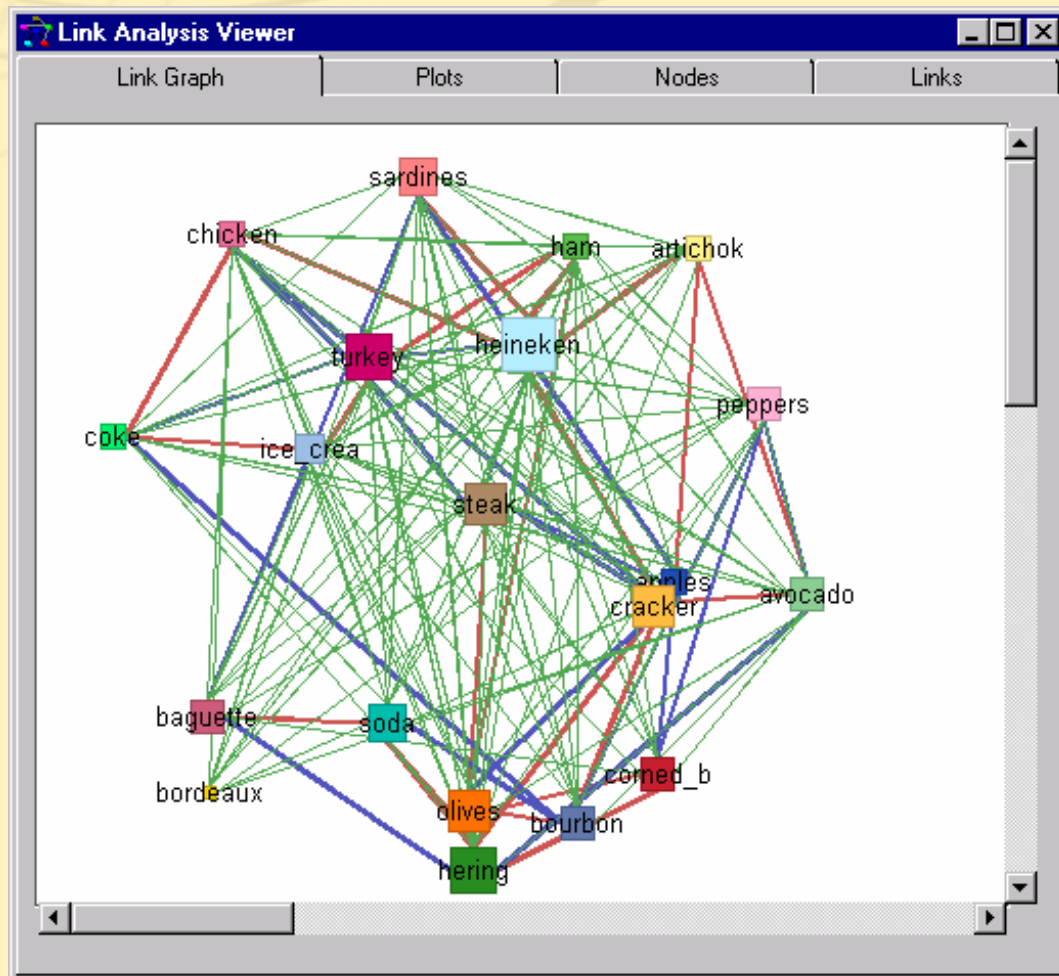
## Assoziationen

	Relations	Lift	Support(%)	Confidence(%)	Transaction Count	Rule
1	2	2.38	21.98	74.32	220.00	coke ==> ice_crea
2	2	2.38	21.98	70.29	220.00	ice_crea ==> coke
3	2	1.91	21.08	58.13	211.00	avocado ==> artichok
4	2	1.91	21.08	69.18	211.00	artichok ==> avocado
5	2	1.68	14.69	49.66	147.00	sardines ==> coke
6	2	1.68	14.69	49.66	147.00	coke ==> sardines
7	2	1.66	11.79	51.98	118.00	steak ==> apples
8	2	1.66	11.79	37.58	118.00	apples ==> steak
9	2	1.65	22.08	78.09	221.00	turkey ==> olives
10	2	1.65	22.08	46.72	221.00	olives ==> turkey
11	2	1.63	15.08	51.01	151.00	sardines ==> ice_crea
12	2	1.63	15.08	48.24	151.00	ice_crea ==> sardines
13	2	1.62	25.07	78.93	251.00	soda ==> cracker
14	2	1.62	25.07	51.43	251.00	cracker ==> soda
15	2	1.55	13.39	47.35	134.00	turkey ==> ham

## Link Analysis

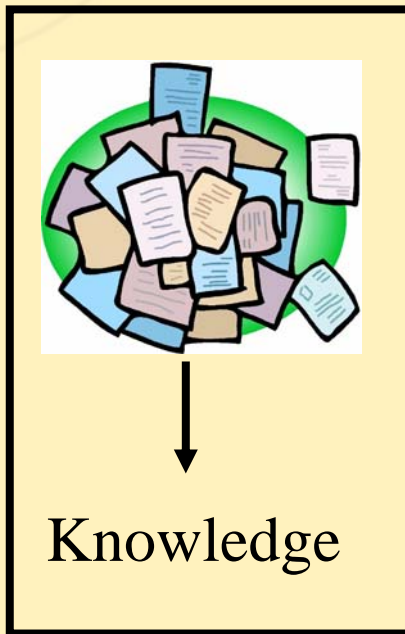


# Warenkorb-Assoziationen (MDS-Plot)



- Visualisierung der Beziehungen in einem Warenkorb
- Farben und Linienstärken entsprechend Konfidenz und Lift
- Größe der Knoten proportional zum Verkaufsertrag

# Text Mining



- Aufdecken und Verwenden von Wissen in Dokumenten
- Mustererkennung
- Beziehungen zwischen Dokumenten und Begriffen
- Transformation von Text in quantitative Informationen

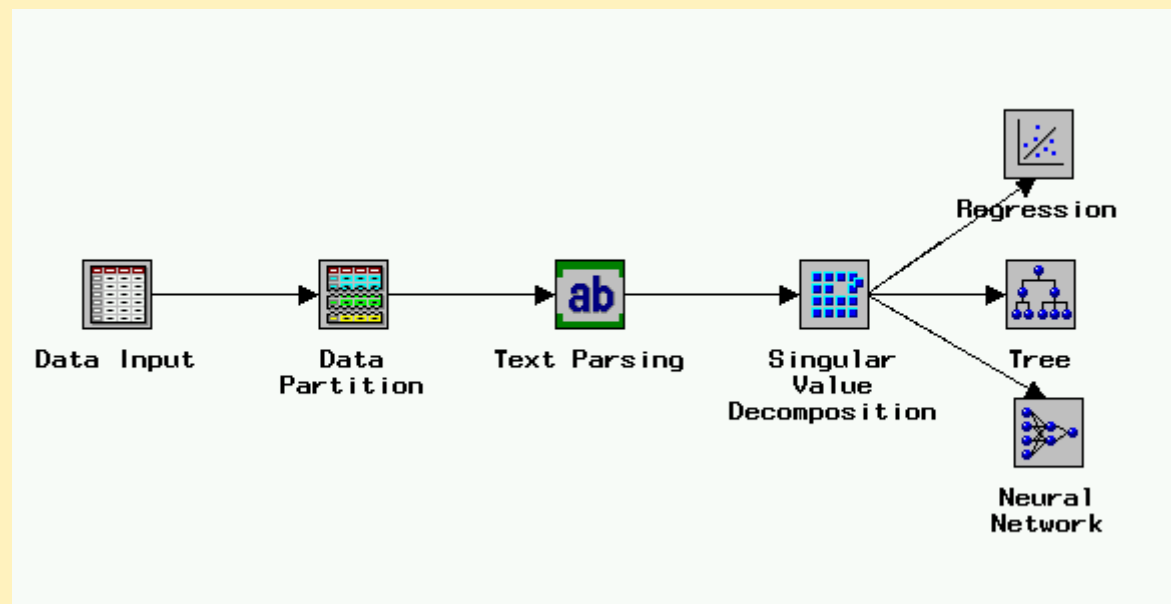
## Text Mining

- Ähnlichkeiten von Dokumenten
  - Gemeinsamkeiten in Verwendung von Wortkombinationen
  - Ähnlichkeiten von Wörtern
    - basierend auf ihrem Vorkommen in ähnlichen Dokumenten
  
- Verwendung externer Quellen
  - Navigations-Sequenzen bei Webnutzung
  - Akzeptanz von Themen für spezielle Kundengruppen
  - Kundencluster aus analogem Navigationsverhalten

# Text Mining Anwendungen

- Klassifikation und Clusterung
  - Email-Filterung
  - Mail & Help Desk Routing
  - Kategorisierung in Wissensdatenbanken
  
- Prognose
  - Aktienkurs aus gegenwärtigen Geschäftsberichten
  - Kundenzufriedenheit
  - Servicekosten anhand von Service Log Daten in Textform

# Text Mining Prozeß





## Singular Value Decomposition (SVD)

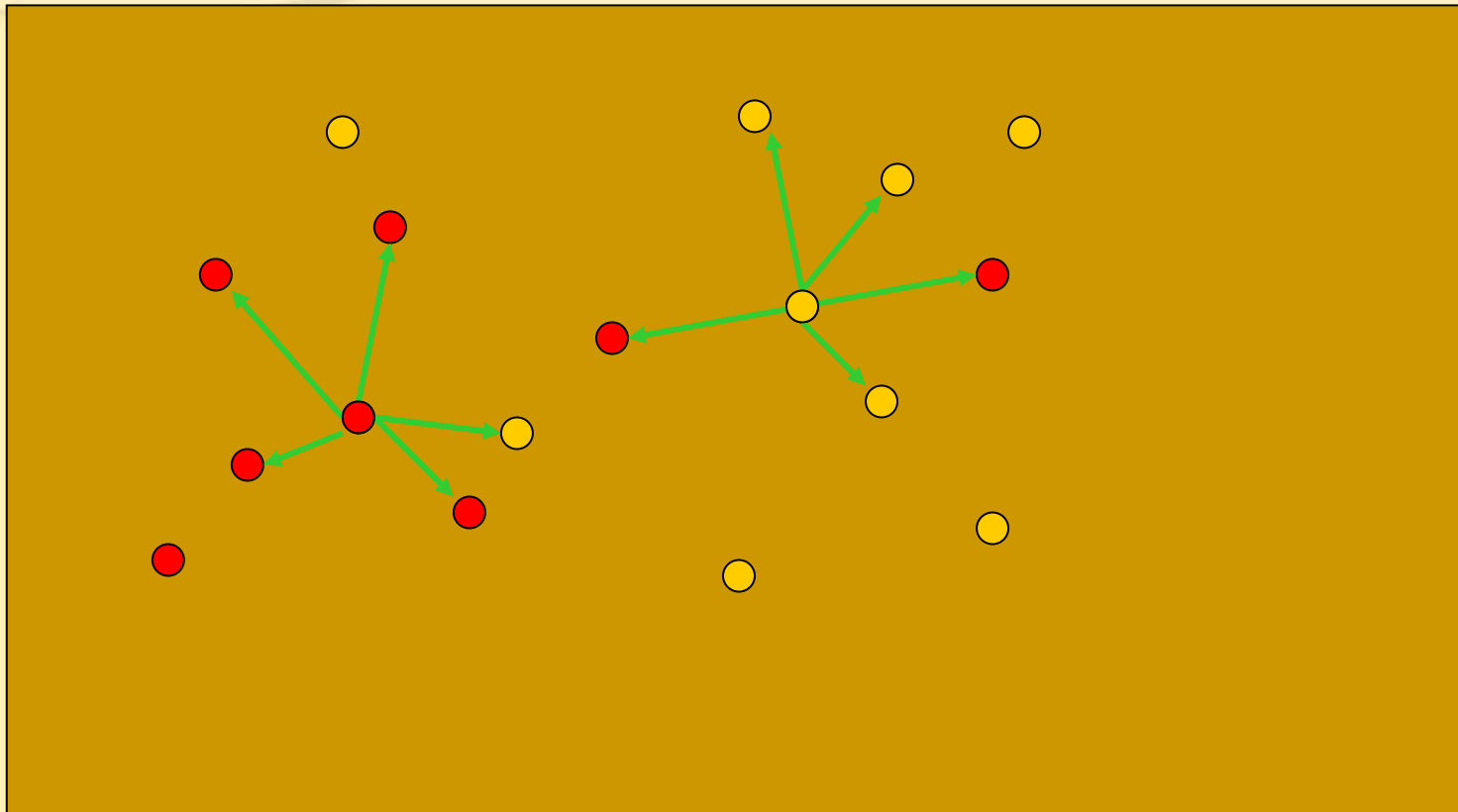
- Dimensionsreduzierung
  - Informationsmatrix
- Latente Semantik Analyse
  - Zusammenhang zwischen Begriffen und Dokumenten
  - Dokumente mit unterschiedlichen Wortmustern aber ähnlicher Bedeutung



## Memory-Based Reasoning

- **“When I see a bird that walks like a duck, swims like a duck and quacks like a duck – I call that bird a duck.” Richard C. Cushing**
- $k$ -Nearest Neighbor Algorithmus
- Startprobe: Basisprognose durch Ähnlichkeiten
- Weitere Variable bekannt
  - Nachführung der Probe
  - Präzisierung der Prognose
- Ähnlichkeiten zu iterativen Segmentierungen

# Memory-Based Reasoning



# Memory-Based Reasoning

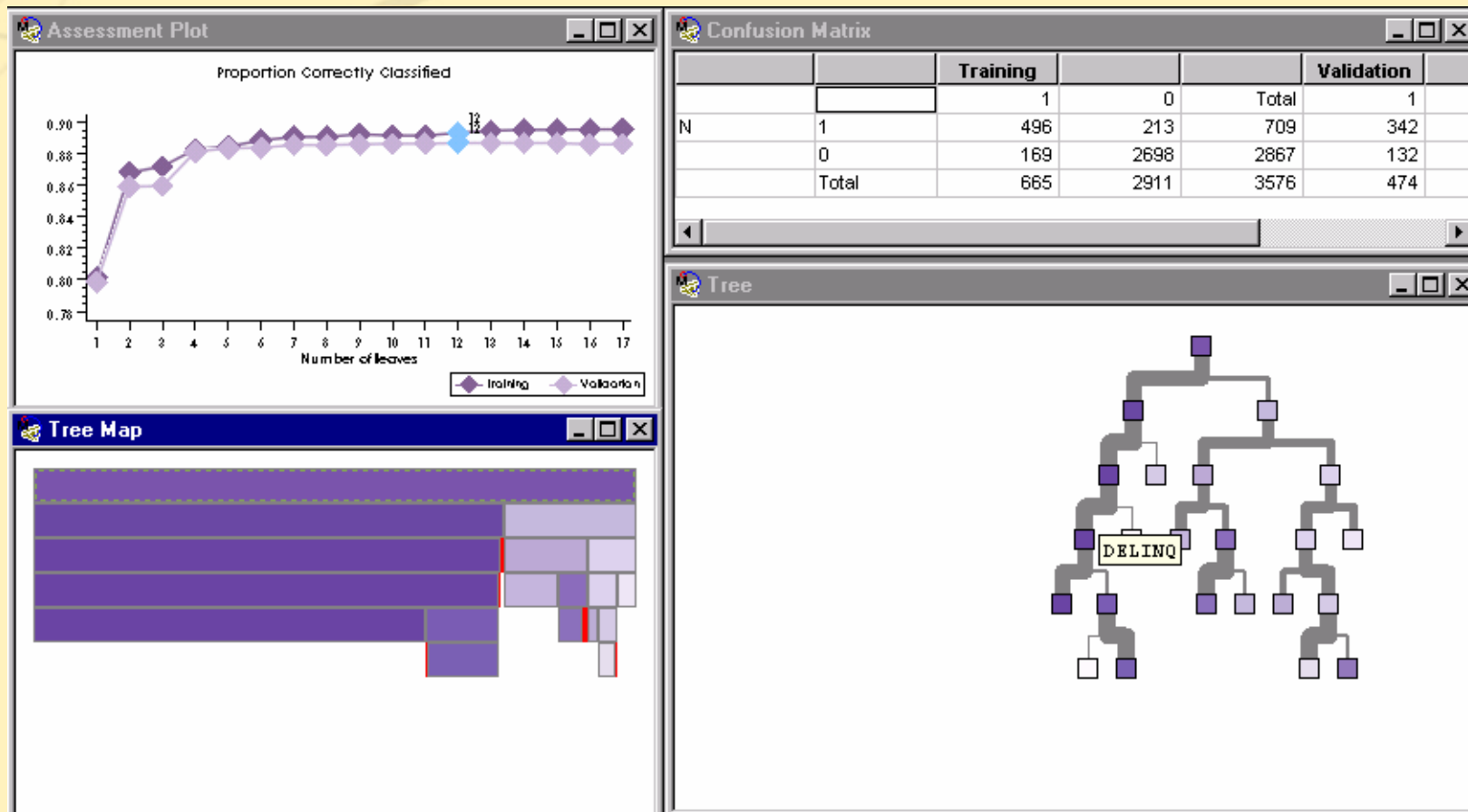
Observation ID	Target : Purchase (Y/N)	Observation Ranking Based on the Distance to the Probe (1: closest 5: farthest)
7	Y	3
12	N	2
35	Y	5
108	Y	1
334	N	4

k	Observation ID of Nearest Neighbors	Target Value of Nearest Neighbors	Posterior Probabilities of the Probe
1	108	Y	prob(Y) = 100% prob(N) = 0%
2	108, 12	Y, N	prob(Y) = 50% prob(N) = 50%
3	108, 12, 7	Y, N, Y	prob(Y) = 67% prob(N) = 33%
4	108, 12, 7, 334	Y, N, Y, N	prob(Y) = 50%; prob(N) = 50%
5	108, 12, 7, 334, 35	Y, N, Y, N, Y	prob(Y) = 60% prob(N) = 40%

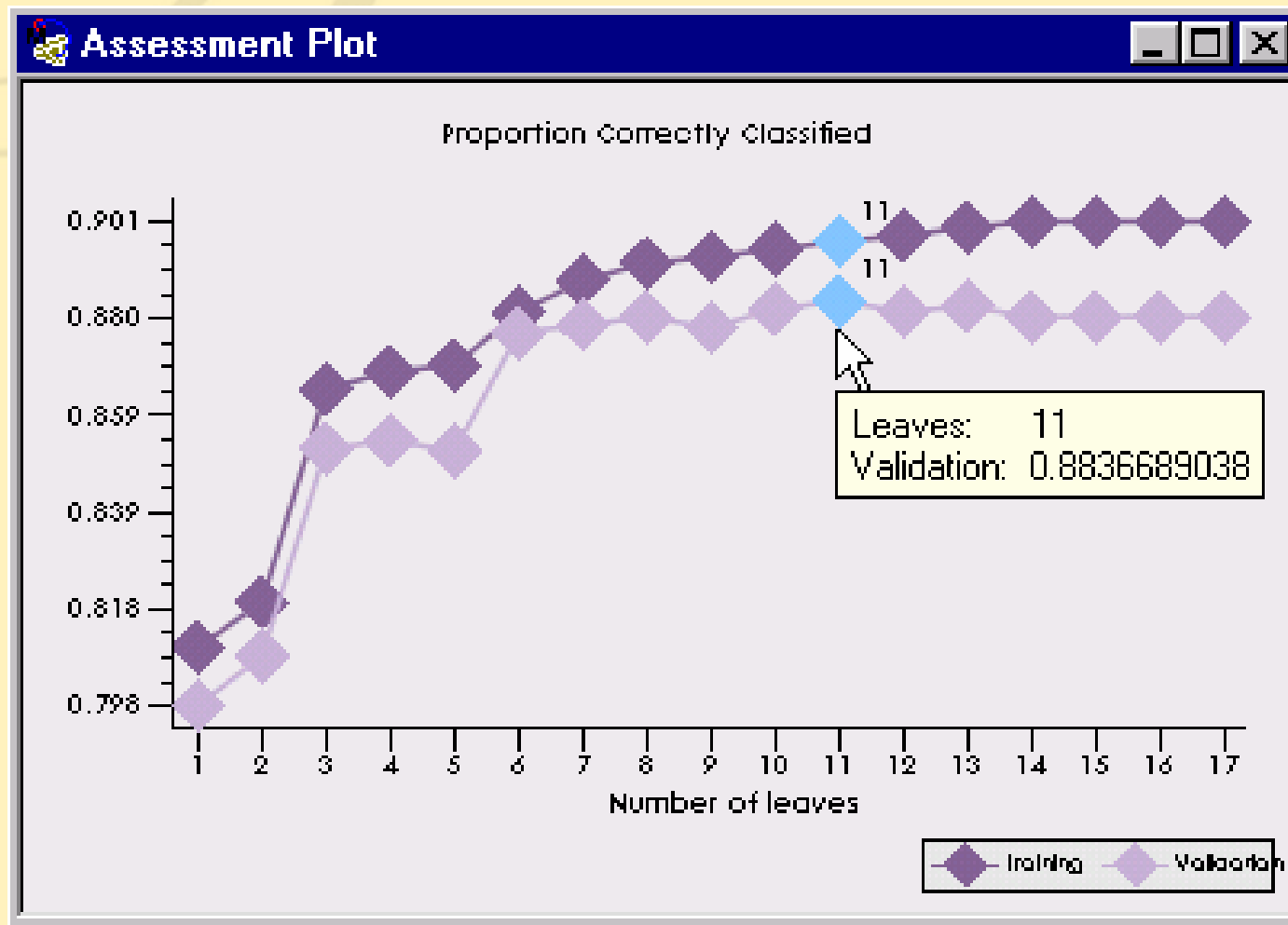
## Memory-Based Reasoning im Vergleich

- Alternative zu anderen prediktiven Modellierungstechniken
  - Vorteilhaft bei erkennbaren Strukturen und Mustern
- Modellanpassung mit zeitlicher Dimension
- Einfache Interpretation
- Hoher Speicherbedarf
- Aufwendige Modellnutzung
  - Keine Formel, kein SAS Base Code
  - Rechenintensive Abstandskalkulationen

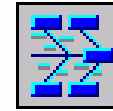
# MFC Tree Browser



# MFC Tree Assessment Plot



# Princomp / Dmneural



Princomp/  
Dmneural

- Additive nichtlineare Modelle
- Wenn Kollinearitäten Modellgüte einschränken
- Hauptkomponenten mit höchstem Bestimmtheitsmaß auf Zielgrößen
- Bucket-Hauptkomponenten als Input-Variablen einer Neuronalen-Netz-Modellierung
- Binäre oder intervallskalierte Zielvariable



## NEURAL oder DMNEURAL ?

- **NEURAL** für kleinere Datenmengen
  - Feingliedrige Modellierung komplexer nichtlinearer Zusammenhänge
  - Zeitintensiv
  
- **DMNEURAL** für große Datenmengen
  - Größere Modellierung
  - Hohe Performance und numerische Stabilität



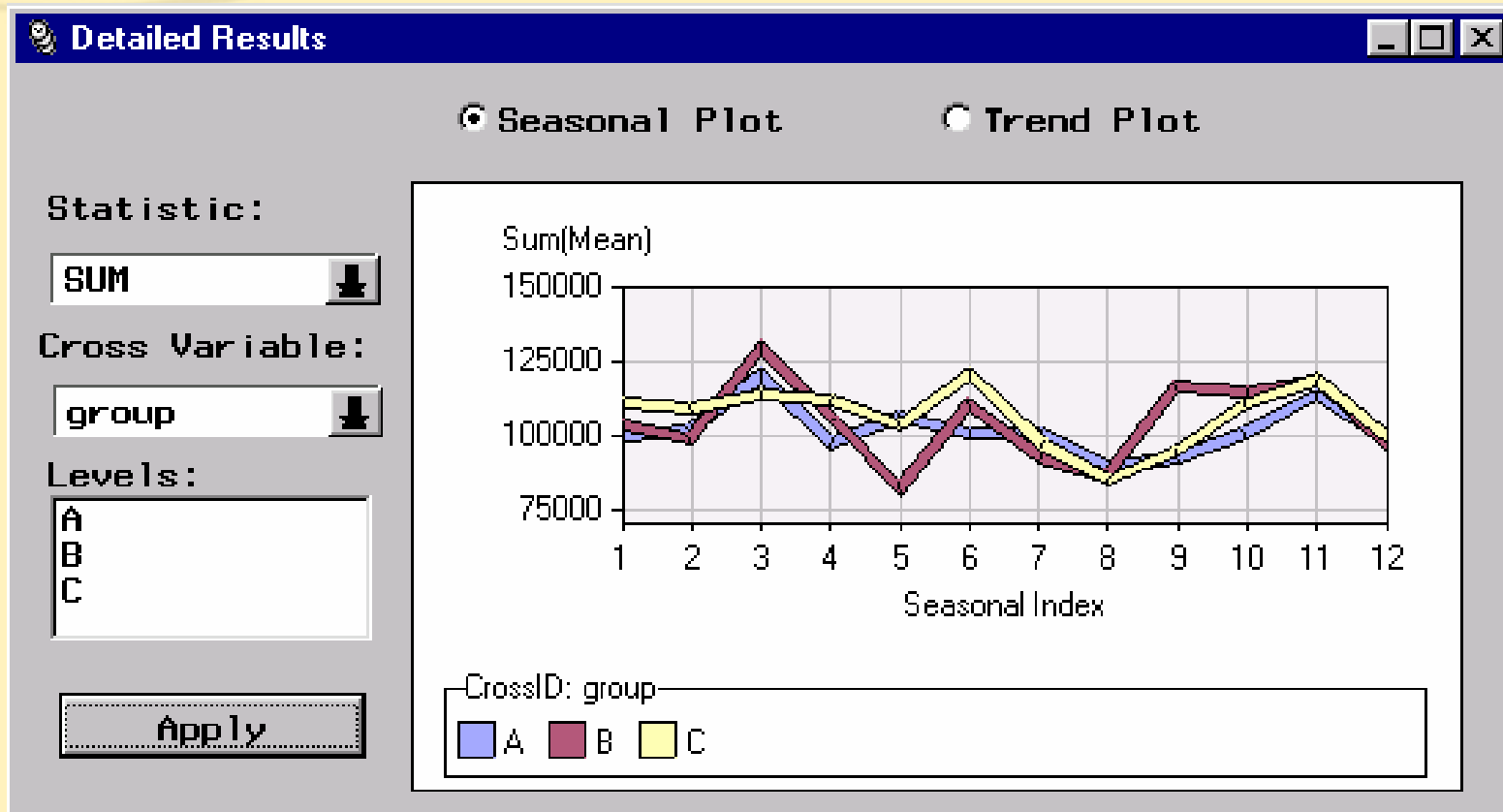
Time Series

## Time Series Node

- Operative Transaktionsdaten in Zeitreihen überführen
- Trendberechnungen
- Saisonale Effekte aufdecken

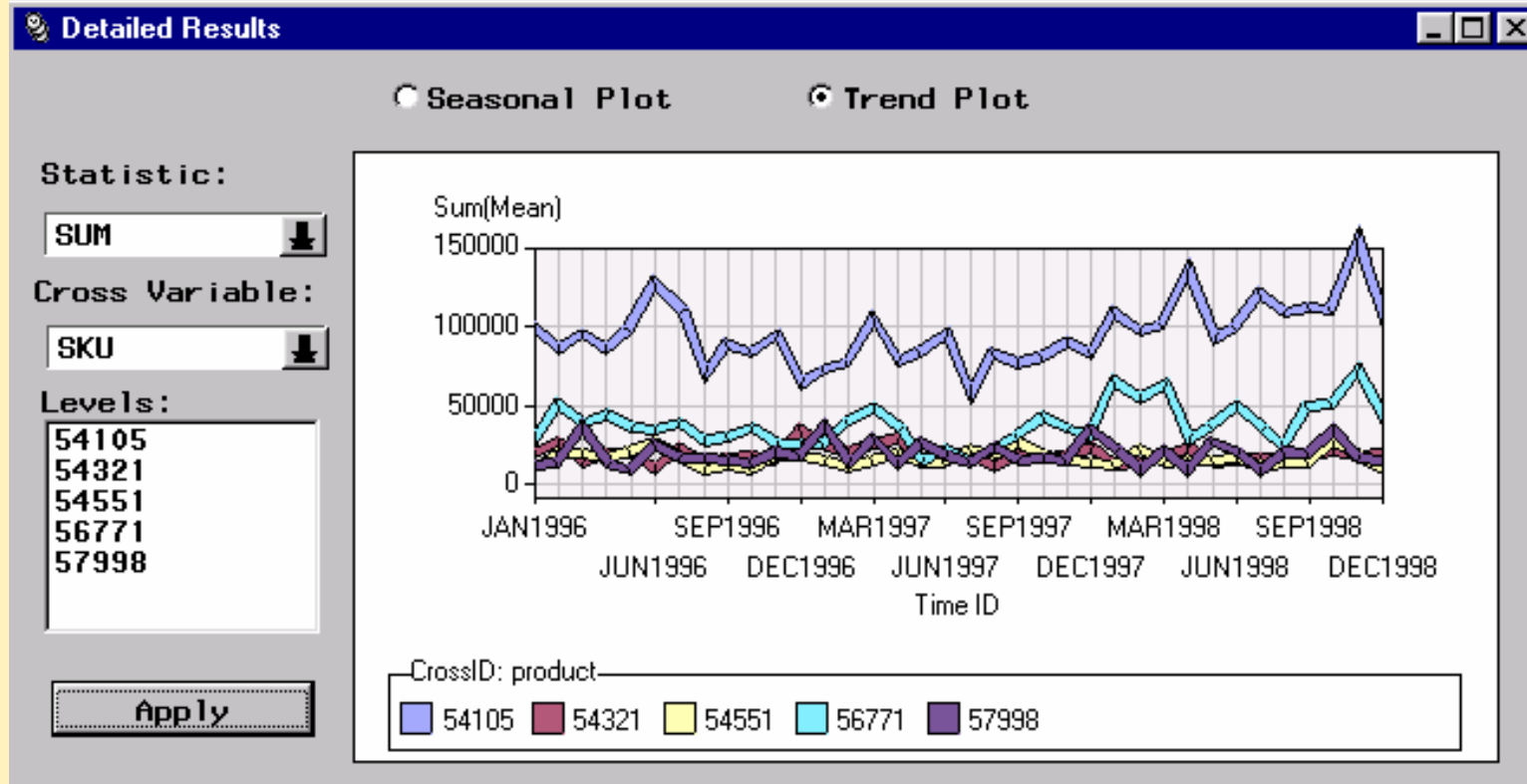
# Seasonal Plot

Mittlere Zielwerte für jede Saison



# Trend Plot

Mittlere Zielwerte für jedes Intervall des Zyklus



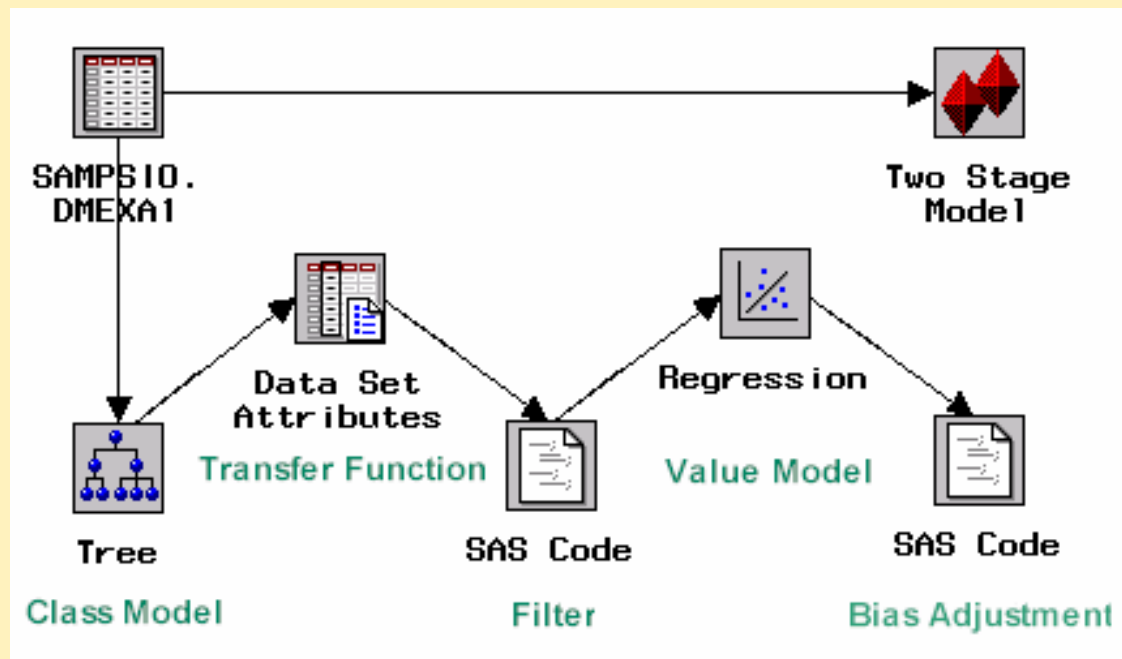
## Zweistufenmodelle

- Kombinierte Vorhersage von nominaler und intervallskalierter Zielgröße
- Beispiel: Kundenbewertung
  - Kunde kauft: ja/nein
  - Wenn Kunde kauft: Umsatzprognose

## Zweistufenmodelle

- Automatische Modellwahl für beide Teilmodelle
- Modellverknüpfung kann gesteuert werden
  - Transfer-Funktionen
  - Filter-Optionen
- Intervall-Modell wird korrigiert durch Posteriori-Wahrscheinlichkeiten des Class-Modells
- Erzeugter Scorecode kombiniert beide Modelle
- Assessment-Plots für das zusammengesetzte Modell

# Erhebliche Vereinfachung



# Typische Anwendungen für Mehrstufenmodelle

- Marketing Automation / Kampagnenmanagement
  - Basisselektionen
  - Auswahl der Kommunikationskanäle
  - Cell Splits
  - Testkampagnen
  - Responsebewertungen
  - Optimierte mehrstufige Kampagnen



# Score Code Deployment

## C\*Score und J\*Score

- Automatische Übersetzung
  - Input: SAS Dastep Score Code generiert durch Enterprise Miner (kein manueller Code)
  - Ouput: Score Code in C und Java Syntax
  
- Erleichterte Verwendung eines EM Model Score Code außerhalb der SAS Umgebung
  - Erfordert kein SAS System zur Laufzeit
  - Erlaubt Einbettung in externe Anwendungen

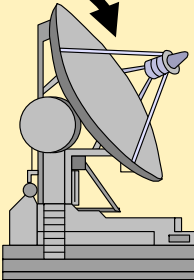
## Engine Diagnostic Tools bei der LH-Technik

- Fehlerfrüherkennung
- Verbrauchsoptimierter Motorbetrieb
- Vermeidung von Triebwerksstandläufen
- Optimierung des Arbeitsumfangs bei der Triebwerksüberholung

**Corrective Maintenance Action**



**ACMS**



**Data Link**

**Action Order**

**MS Eng. T/C Cust. OEM**

**On Request Output:**

- Input Data
- Trends
- Alert Messages

**Test Cell**

**IBM Mainframe Frankfurt**

ACM Ground System

Performance Analysis

- Aircraft
- Engine
- APU

Trend Recognition

**automatic Alert Messages**

# Engine Condition Monitoring - Funktionen

SLOATL  
EGT MARGIN

OIL SYSTEM  
MONITORING

EARLY  
FAILURE  
DETECTION

SOAP

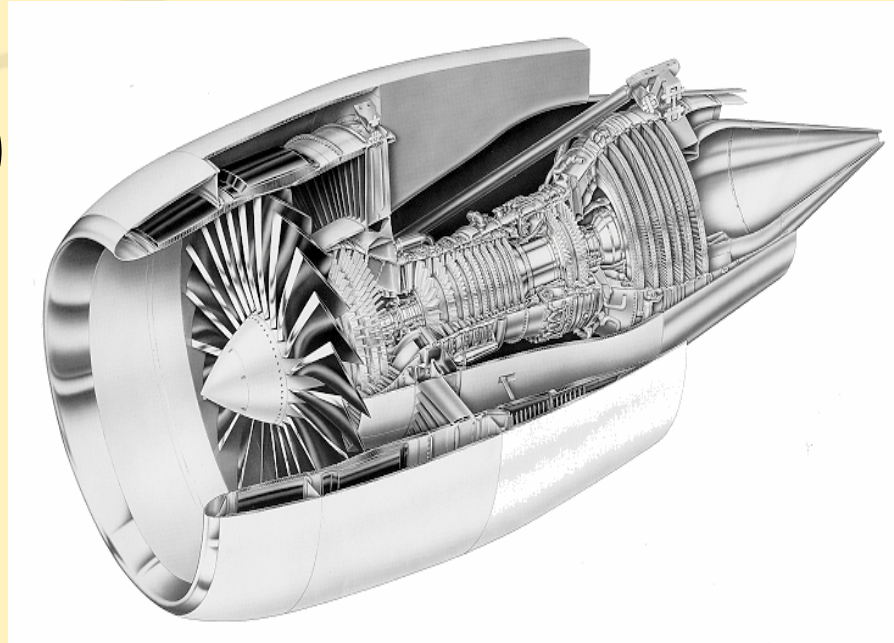
COLD  
FAN TRIM  
BALANCE

CONTROLS  
MONITORING

TEST CELL

MPA

DERATE  
STATISTICS



## Weiterentwicklung des Instandhaltungssystems

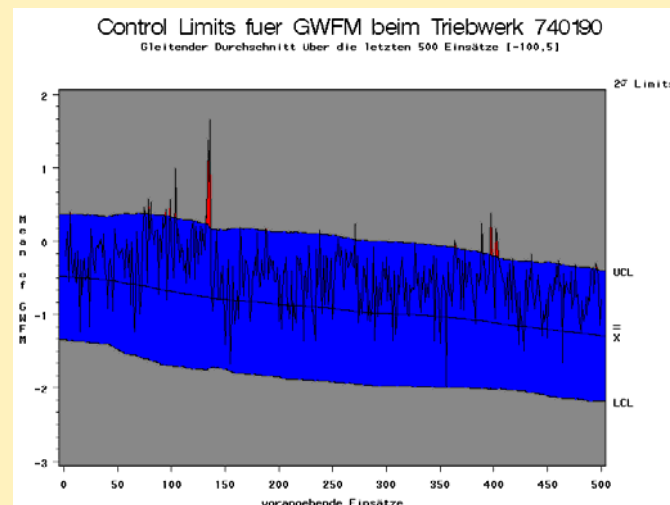
- Erhöhung der Sicherheitsreserven
- Verringerung des Kraftstoffverbrauchs
- Minimierung der Umweltbelastung
- Verringerung der Wartungskosten

## Analyseziele

- Erweiterung bisheriger Analysen durch komplexere nichtlineare Modelle mit Wechselwirkungen
- Selektion wesentlicher Variablen aus wachsender Zahl von Meßgrößen
- Erkennen von Zusammenhängen innerhalb von Schadensgruppen

## Ergebnisse

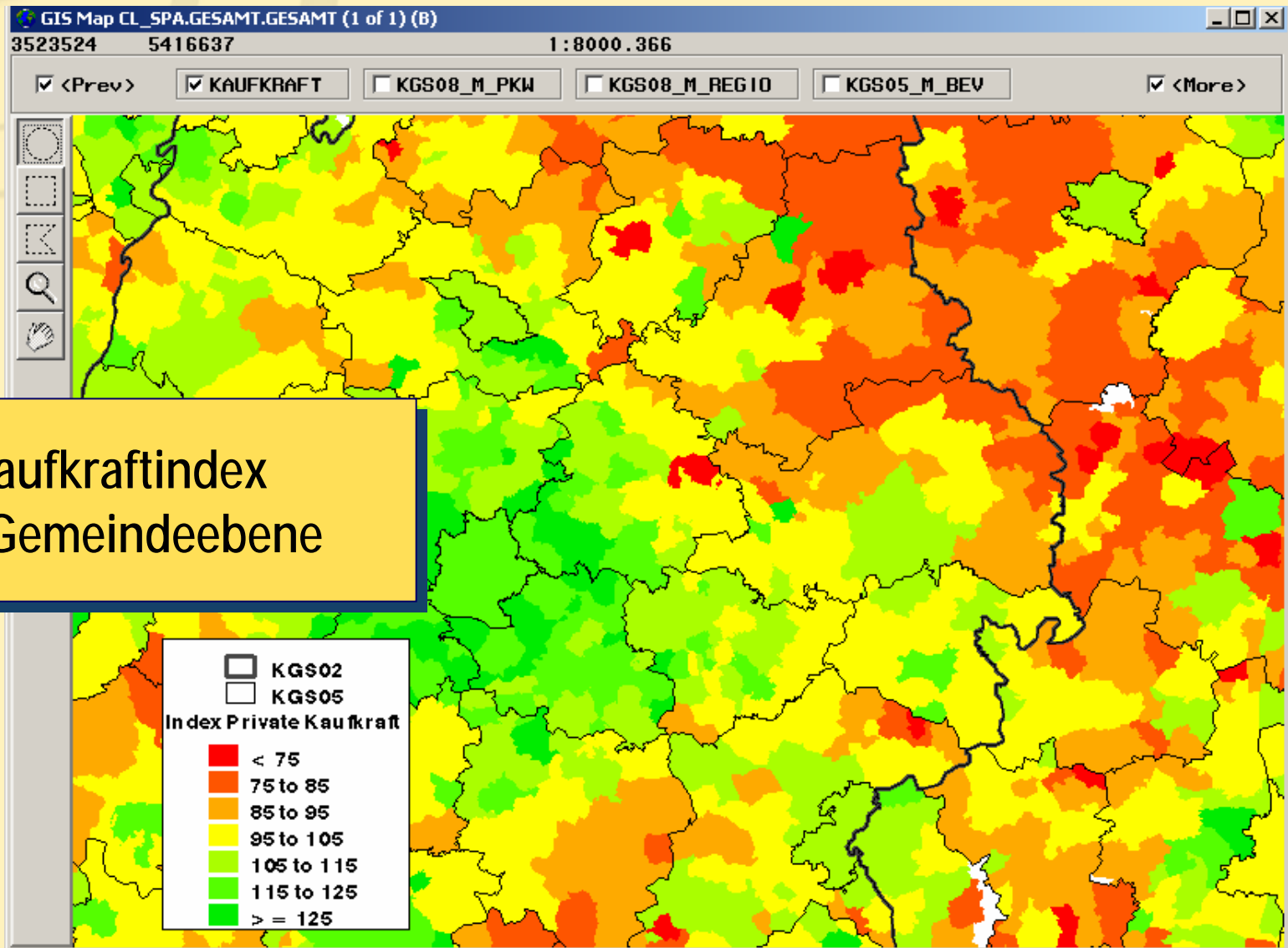
- Präzisere Analyseergebnisse durch Kombination mehrerer signifikanter Triebwerksparameter
- Reduzierung redundanter Warnmeldungen
- Erhöhte Flexibilität des Systems

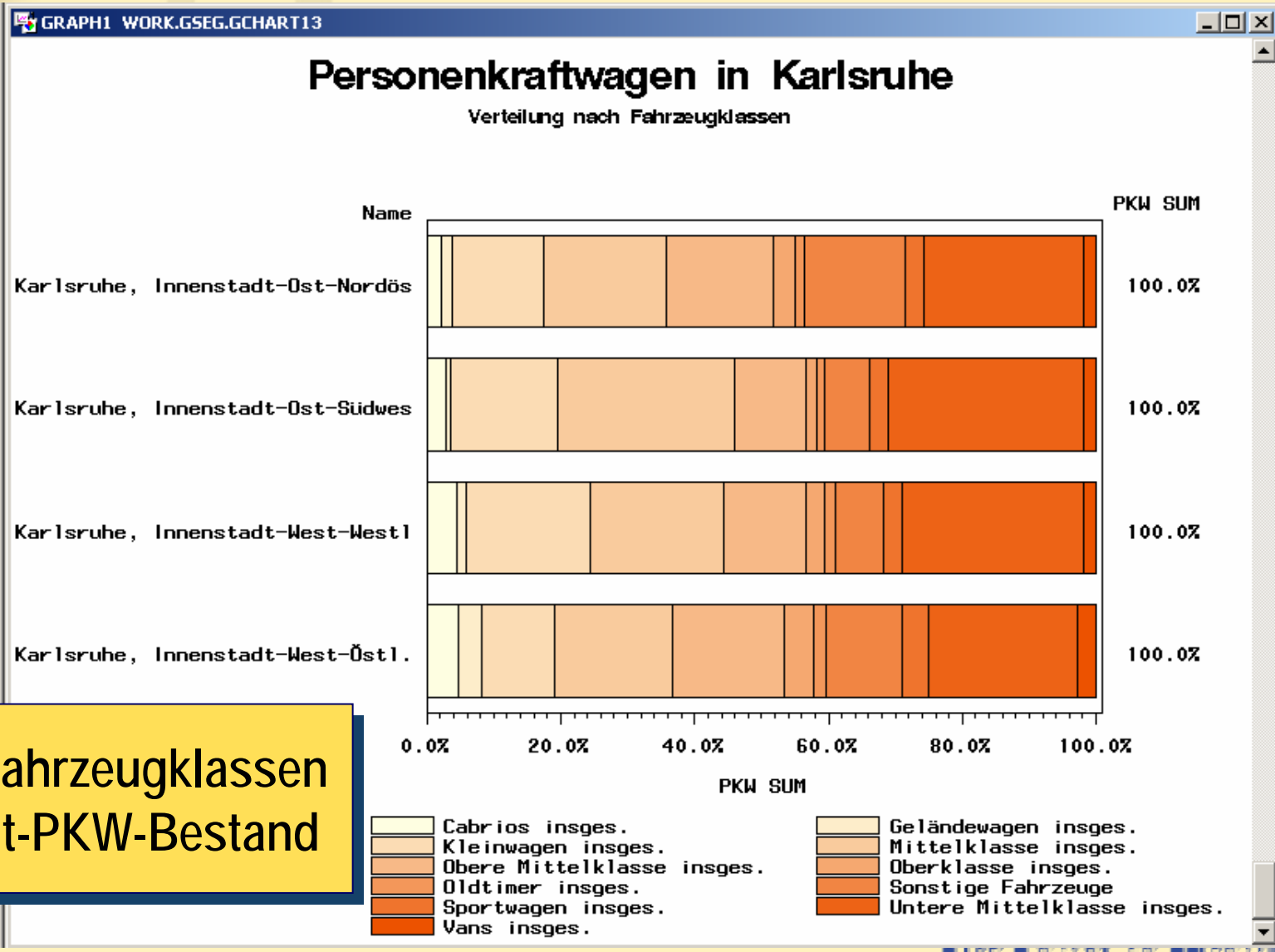


# Data Mining in der Genforschung

- Proc ALLELE
  - Analyse genetischer Markierungsdaten
- Proc CASECONTROL
  - Tests auf Relationen zwischen Markierung und Erkrankung
- Proc FAMILY
  - Transmission/Disequilibrium-Tests (TDT)
- Proc HAPLOTYPE
  - MLE für Haplotype-Häufigkeiten
- Proc PSMOOTH
  - Glättungsmethoden für multiples Testen

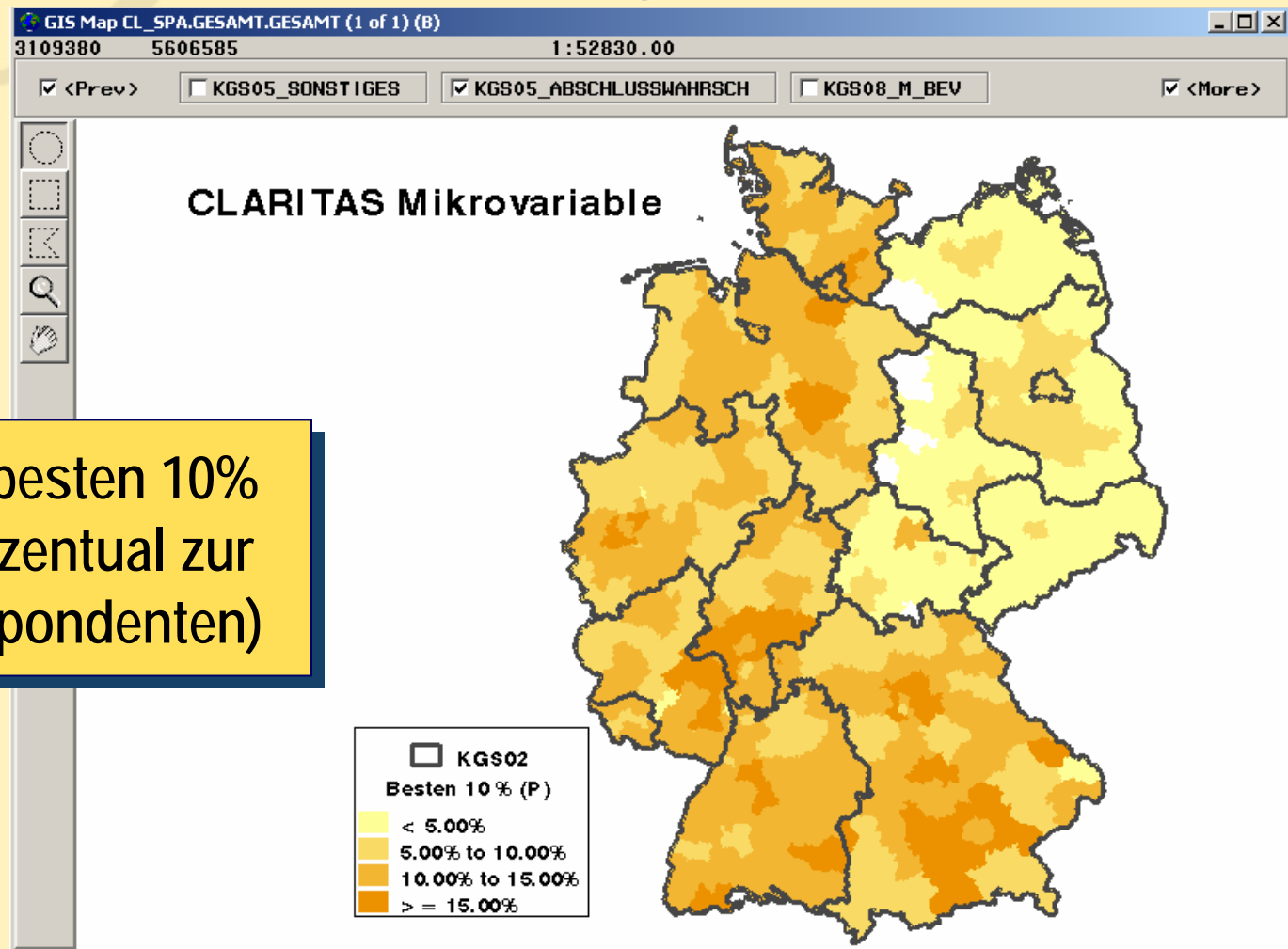






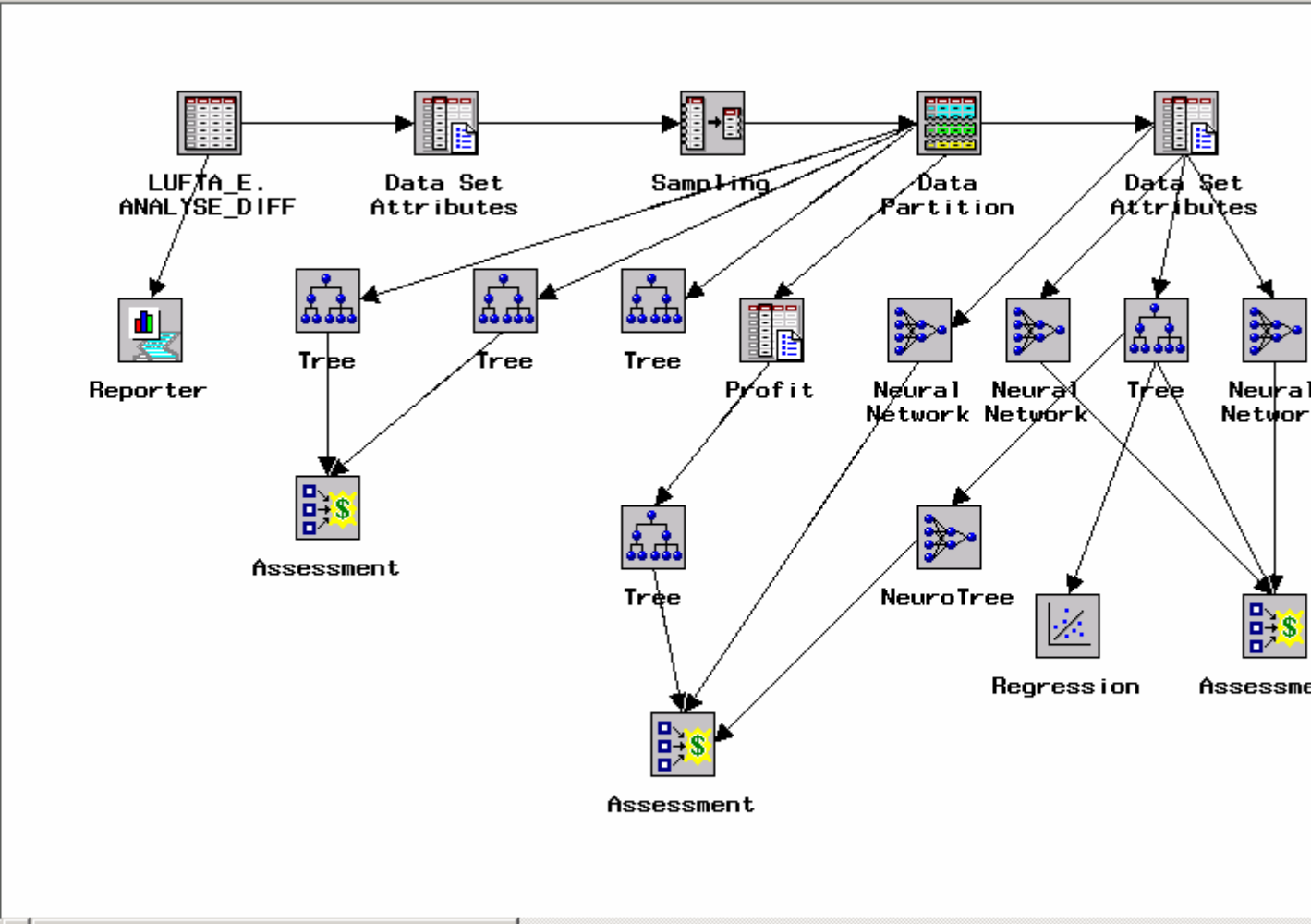
Anteil der Fahrzeugklassen  
am Gesamt-PKW-Bestand

# Abschlusswahrscheinlichkeit Lebensversicherung über € 100.000



Verteilung der besten 10%  
(Verteilung prozentual zur  
Anzahl der Respondenten)

- Sample
  - Input Data Source
  - Sampling
  - Data Partition
- Explore
  - Distribution Explorer
  - Multiplot
  - Insight
  - Association
  - Variable Selection
- Modify
  - Data Set Attributes
  - Transform Variables
  - Filter Outliers
  - Replacement
  - Clustering
  - SOM/Kohonen
- Model
  - Regression
  - Tree
  - Neural Network
  - User Defined Model
  - Ensemble
- Assess
  - Assessment
  - Score!
  - Reporter
- Utility



Diagrams Tools Reports

Connected to alpha

# SAS Enterprise Miner Version 5.0 (SAS Version 9)

The screenshot displays the SAS Enterprise Miner interface. On the left, a 'test' workflow is visible with nodes: test \*\*\* OPEN, test2, IDS, PART, CLUS, VARSEL, NEURAL, TREE, PRINCOMP, REG, PATH, and TIME. The main workspace shows a detailed flowchart starting with 'IDS' leading to 'PART'. From 'PART', the flow splits into three paths: one through 'VARSEL' to 'NEURAL', one through 'TREE', and one through 'PRINCOMP' to 'REG'. All three paths converge into a single 'PLOT' node. A secondary path from 'IDS' leads to 'PATH' and 'TIME'. A 'Property' table is visible on the left side of the interface.

Property	Value
NODEID	NEURAL
Data	
Network Arc...	NLP
Training Tec...	Default
Direct Conn...	No
Hidden Units	20

View: Basic

**Network Architecture**

Specify a network architecture used in training from the followings;  
 generalized linear model (GLIM), multilayer perceptron (MLP--Default), ordinary radial basis function with equal widths (ORBFEQ), ordinary radial



*The Power to Know.*