27. SAS Club

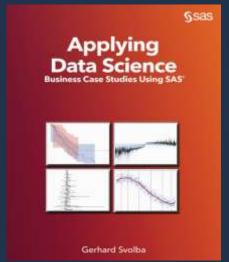
Buchpräsentation:

Applying Data Science

Business Case Studies Using SAS

Gerhard Svolba, Franz Helmreich, Gernot Engel, Matthias Svolba, Mihai Paunescu

Wien, 23. November 2017 – ARES Tower, Wien







Agenda

- 10 mal "Data Science in Action"
 - Supervised Machine Learning Methoden
 - Unsupervised Machine Learning Methoden
 - Simulationen

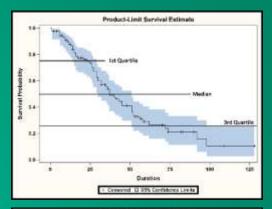
- SAS Viya Offenheit für unterschiedliche Benutzertypen
 - Gernot Engel, Franz Helmreich, Matthias Svolba, Gerhard Svolba
- SAS Tipps und Tricks Session
 - Mihai Paunescu, Gerhard Svolba

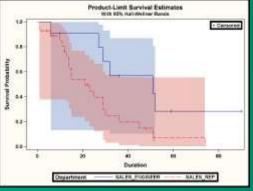


Data Science in Action: #1

Performing Headcount Survival Analysis for Employee Retention

Can assumptions about the average length of time intervals be made, even if most of the endpoints have not yet been observed?



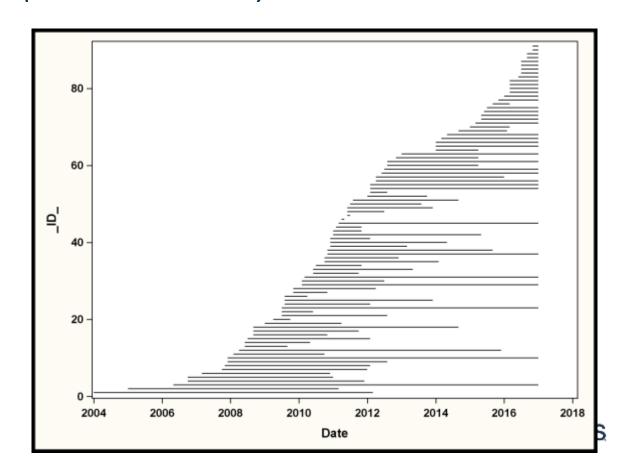


Survival analysis methods: Kaplan-Meier estimates
Cox Proportional Hazards regression
Survival Data Mining
Company Confidential - For Internal Use City

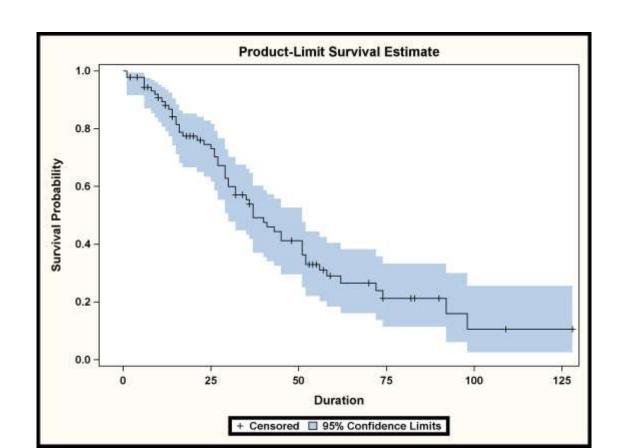


Nicht zu allen Mitarbeitern haben wir ein "Ereignis-Datum" (Glücklicherweise)

- Betrachten der Karrieren pro Mitarbeiter
 - Unterschiedliche Länge
 - Kündigung oder "zensiert"

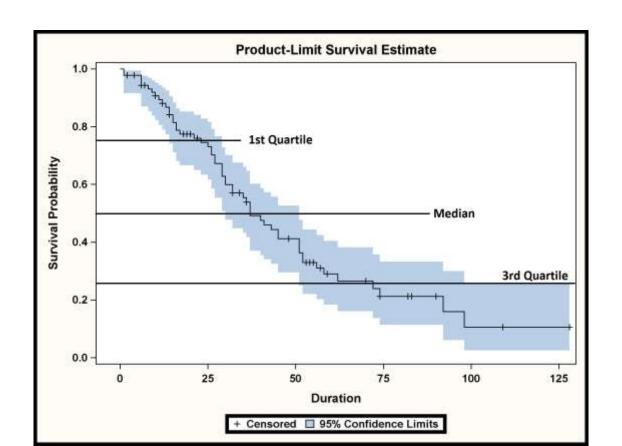


Survival Kurve (mit Konfidenzband) für alle Mitarbeiter



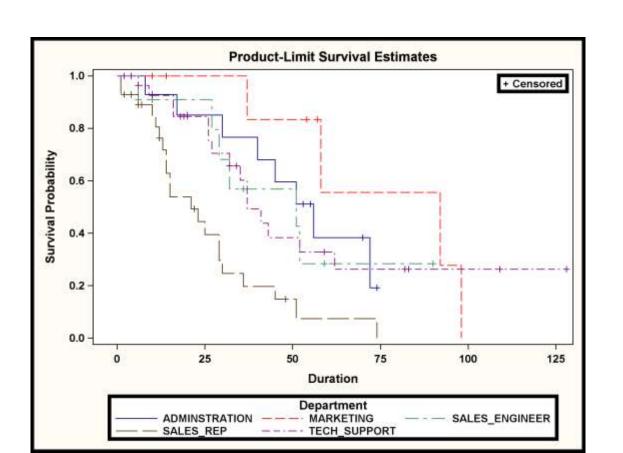


Interpretation der Survival Kurse anhand der Quartile





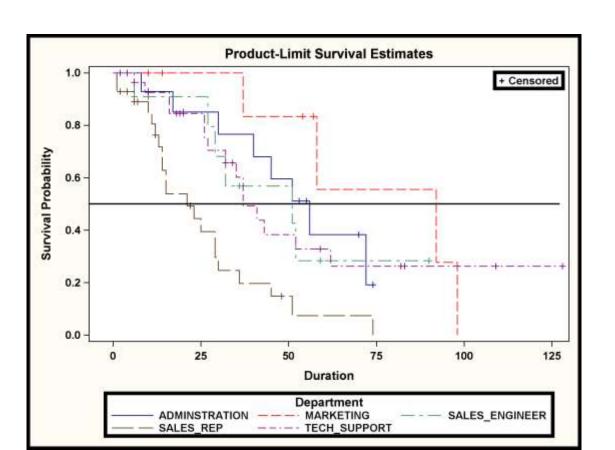
Survival-Kurve pro Abteilung





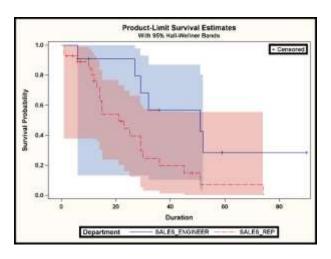
Survival-Kurve pro Abteilung

Referenz-Linie für den Median

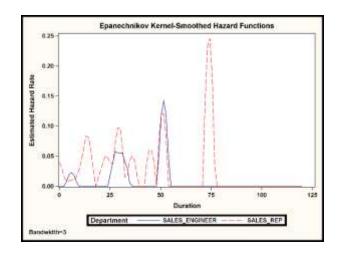




Die Kaplan Meier Methode und die Cox Proportional Hazards Regression verarbeitet zensierte Beobachtungen



Kaplan Meier Methods und Cox Proportional Hazards Regression: Sales engineers haben eine bessere "survival time" als sales representatives.



Betrachten der Hazard Kurven: Es gibt ein hohes Risiko die Sales Engineers nach 26 und 50 Monaten zu verlieren.



Time-to-Event Analyse mit SAS/STAT Procedures

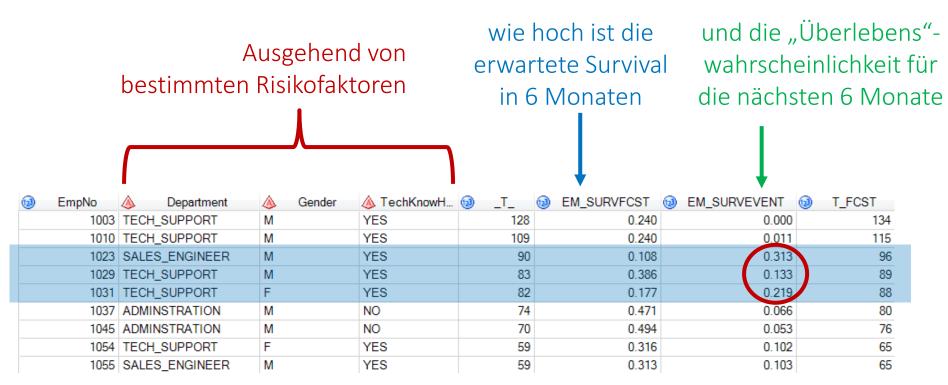
Proc LIFETEST und PROC PHREG

```
proc lifetest data=employees outsurv = survplot
              plots=(hazard(bandwidth=3 maxtime=120)
              survival(cb=hw));
 time duration*status(1);
 strata department;
 where department in ("sales rep", "sales engineer");
run;
PROC PHREG DATA=Employees outest = ParamEstimates;
 CLASS department gender TechKnowHow StartPeriod/ PARAM=effect REF=first;
 MODEL Duration*Status(1) = department gender / SELECTION=stepwise;
 OUTPUT OUT=surv pred survival=SurvPred
                      Atrisk = ObsAtRsik
                               =DisplacmLikelihood;
                       T<sub>1</sub>D
RUN;
```



"Wie lange wird Gerhard Svolba noch in unserem Unternehmen sein?"

Vorhersage der Verweildauer für indivudelle Mitarbeiter



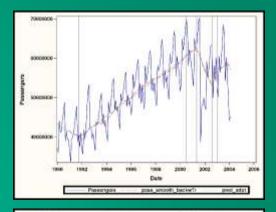


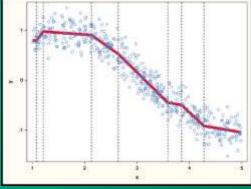
Data Science in Action: #2

Detecting Structural
Changes and Outliers in
Longitudinal Data

Can events and changes in the course over time be automatically detected?

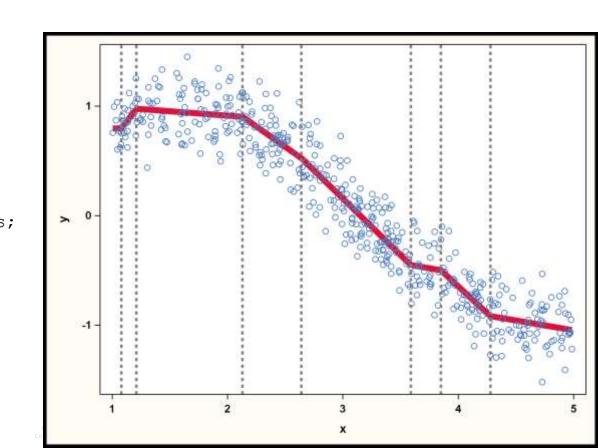
Smoothing Of Longitudinal Data
Multivariate Adaptive Regression Splines
Automatic Breakpoint Detection
Automatic Detection of Outliers with ARIMA Models





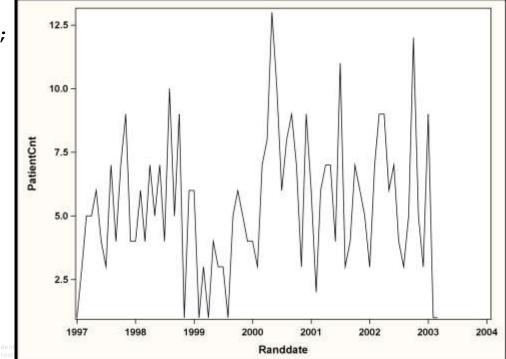


Multivariate Adaptive Regression Splines mit der ADAPTIVEREG Procedure

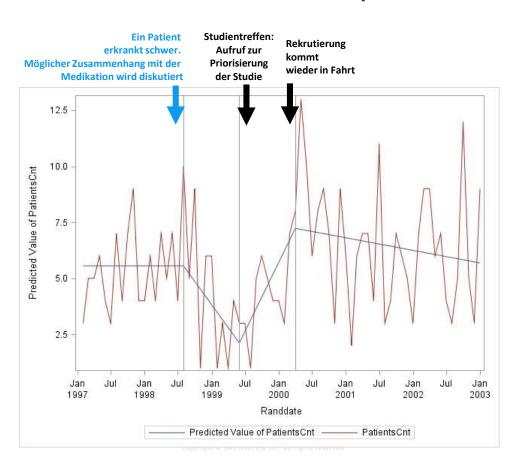


Anzahl der rekrutierten Patienten in einer klinischen Studie im Zeitverlauf

proc adaptivereg data=patients recruitment plots=all;

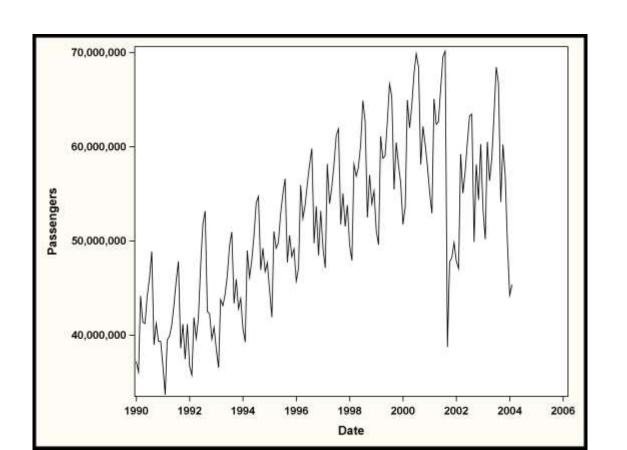


Was ist zu bestimmten Zeitpunkten in meiner klinischen Studie passiert?





Anzahl der Flug-Passagiere in den Jahren 1990 bis 2004

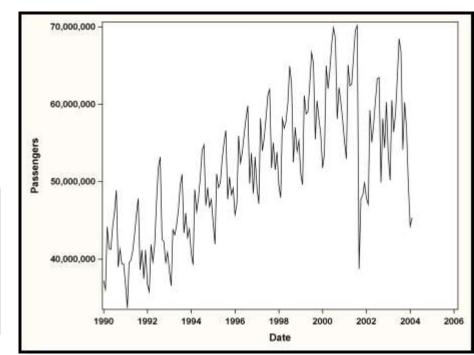




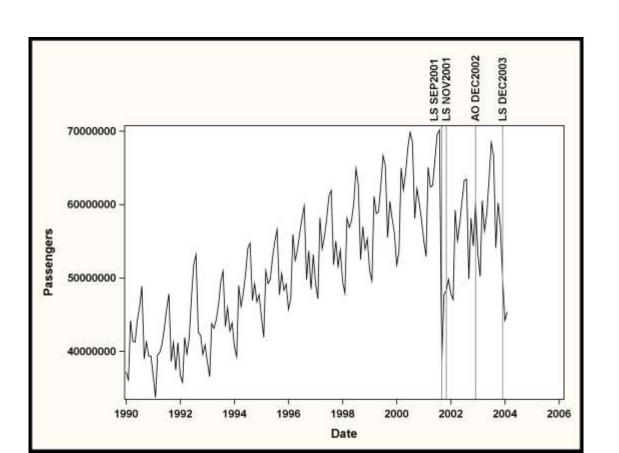
Automatische Ausreißer-Erkennung mit der X13-Procedure

```
proc x13 data=flights_911 date=date;
  var passengers;
  arima model=( (0,1,1)(0,1,1) );
  outlier;
run;
```

Regression Model Parameter Estimates						
For Variable Passengers						
Туре	Parameter	NoEst	Estimate	Standard Error	t Value	Pr > t
Automatically Identified	LS SEP2001	Est	-17993818	1113414.76	-16.16	<.0001
	LS NOV2001	Est	5939640.53	1123179.20	5.29	<.0001
	AO DEC2002	Est	5039786.79	1111284.48	4.54	<.0001
	LS DEC2003	Est	-8531934.1	1249512.29	-6.83	<.0001



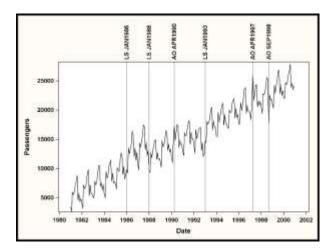




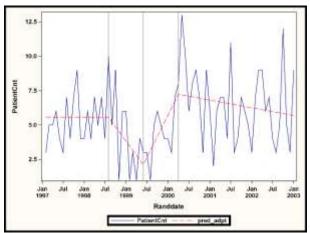


Automatisches Erkennen von Breakpoints und Ausreißern

Anwenden von analytischen Methoden zum Erkennen von Zeitpunkten, wo der Verlauf der Daten vom "normalen" Muster abweicht.



Erkennen von Shifts und Pulse Events mit ARIMA Modellen



Verwenden von Multivariate Adaptive Regression Splines zum Auffinden von Bruchpunkten

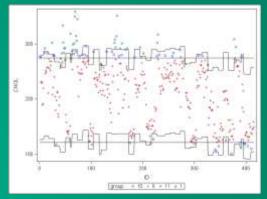


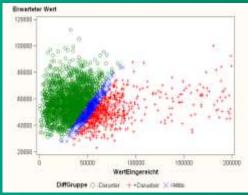
Data Science in Action: #3

Proving a reference value that considers all available co-information

Can analytics help me to reduce the "Yes, but ... " sentences in my business dicussions?

Linear Regression
Decision Trees
Time Series Analysis





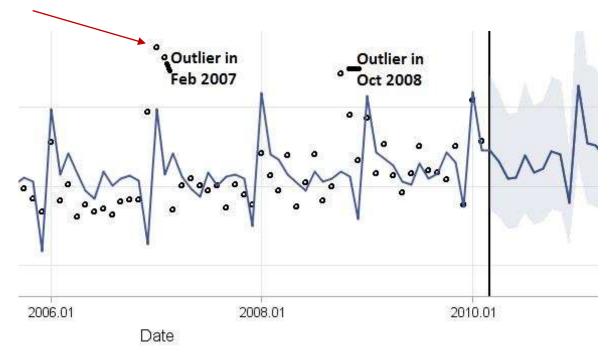


"Ja, aber …. im Jänner haben wir immer deutlich mehr Ereignisse"

Plausibler Wert für Jänner 2007

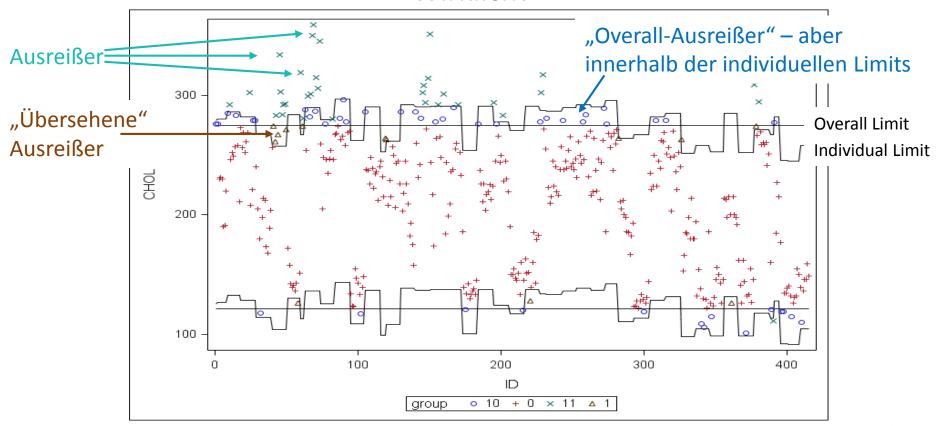
weil Jänner-Werte immer

höher sind





"Alle deren Wert größer x ist, sind Ausreißer! - Wirklich?"



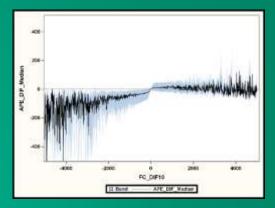


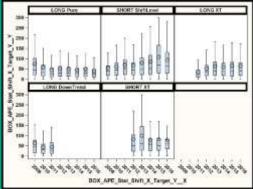
Data Science in Action: #4

Explaining Forecast Errors and Deviations

Do the demand planners really improve forecast accuracy with their manual overwrites?

Linear Regression Quantile Regression Descriptive Statistics





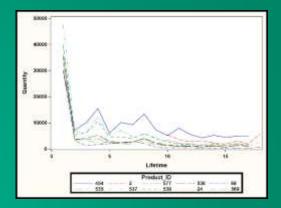


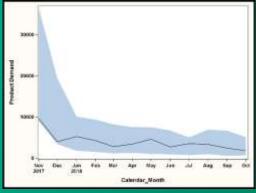
Data Science in Action: #5

Forecasting the Demand for New Products

Can the expected demand of products that are introduced only right now be estimated for forecast planning?

Poisson Regression Cluster Analysis Similarity Search



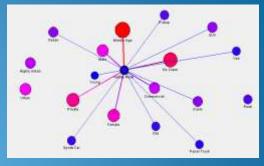


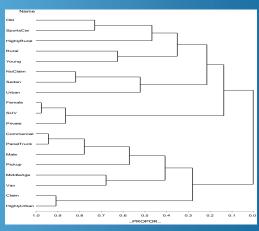


Data Science in Action: #6

Listening to Your Data – Discover Relationships with Unsupervised **Analysis Methods**

Can your data tell you stories about your analysis subjects, even if you don't ask explicitly?





Unsupervised machine learning methods: association analysis variable clustering



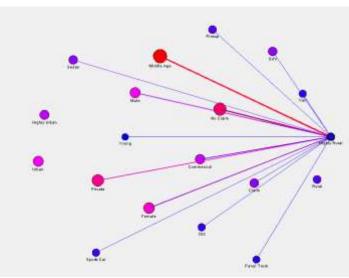
Lassen Sie ihre Daten sprechen!

Auffinden von Zusammenhängen in Ihren Analysedaten

• Daten aus der KFZ-Versicherung mit 6 Eigenschaften pro Versicherungsnehmer

Variable	Feature
AGE	YOUNG, MIDLIFE, OLD
GENDER	MALE, FEMALE
DENSITY	HIGHLY URBAN, URBAN, HIGHLY RURAL, RURAL
CAR_TYPE	VAN, SPORTS CAR, SUV, SEDAN, PICK UP
CAR_USAGE	PRIVATE, COMMERICIAL
CLM_FLAG	CLAIM, NO CLAIM

 Anwenden von unsupervised machine learning (Assoziationsanalyse) um Zusammenhänge zwischen den Eigenschaften aufzudecken.

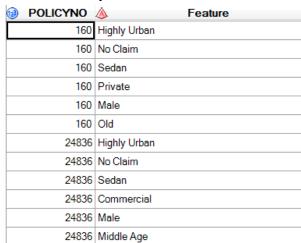


Trauen Sie sich! Transponieren Sie die Daten, so wie Sie es sonst typischerweise nicht tun.

One-Row-Per-Subject



Multple-Row-Per-Subject Key-Value Tabelle





Lassen Sie ihre Daten sprechen!

Männer fahren kaum Sportwägen?

Regel 278 besagt, dass Sportwägen nur in 2,54 % der Fälle von Männern gefahren werden (erwartet wären 46 %)

index	A RULE	≜ _LHAND	_RHAND	@ COUNT (SUPPORT @	EXP_CONF (CONF	LIFT	19
267	Commercial ==> Sports Car	Commercial	Sports Car	200.00	1.94	11.44	5.28	0.48	S
268	Rural ==> Claim	Rural	Claim	102.00	0.99	26.66	6.52	0.2	\$
269	Claim ==> Rural	Claim	Rural	102.00	0.99	15.18	3.71	0.2	ţ
270	Young ==> Highly Urban	Young	Highly Urban	10.00	0.10	34.93	8.33	0.2	1
271	Highly Rural ==> Claim	Highly Rural	Claim	32 00	0.31	26.66	6.30	0.2	1
272	Claim ==> Highly Rural	Claim	Highly Rural	32.00	0.31	4.93	1.17	0.24	\$
273	Van> Female	Van	Female	117.00	1.14	53.82	12.70	0.24	1
274	Female> Van	Female	Van	117.00	1,14	8.94	2.11	0.2	1
275	Panel Truck> Female	Panel Truck	Female	40.00	0.39	53.82	4.69	0.0	9
276	Male> SUV	Male	SUV	99,00	0.96	27.98	2.08	0.0	7
277	SUV> Male	SUV	Male	99.00	0.96	46.18	3.43	0.0	7
278	Sports Car> Male	Sports Car	Male	30.00	0.29	46.18	2.54	0.00	5

- Kann anzeigen, dass in unserer Datenbasis tatsächlich Sportwägen in erster Linie von Frauen gefahren warden.
- Möglicherweise bietet ein Mitbewerber eine Polizze für Männer zu einem deutlich besseren Preis an.
- Ein fachliche Erklärung kann sein, dass der Sportwagen das 2. oder 3. Auto in der Familie ist, und dieser aus steuerlichen Gründe auf die Ehefrau registriert ist.
- Kann auch ein Trigger für eine detailliertere Analyse der Datenqualität sein.

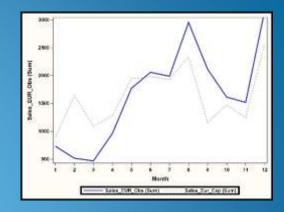


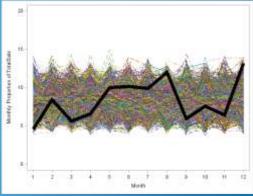
Data Science in Action: #7

Checking the Alignment with Predefined Pattern

Which customers show a behavior that is far from what you expected?

Chi2 independency test Benford's law Time Series Similarity



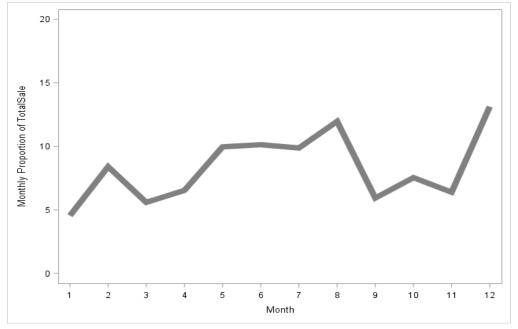




"Welche meiner Verkäufer halten sich kaum an unsere Vorgaben?"

Der Bedarf an "Sub-Contracts" für ein Cateringunternehmen variert im Verlauf eines Kalenderjahres

Verkäufer sind angehalten, entsprechend dieses Musters Verträge zu akquirieren.

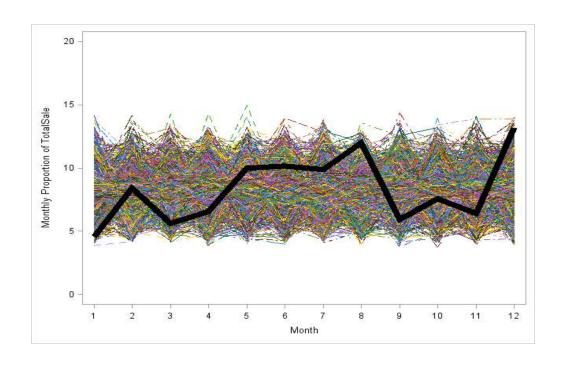




Anzeige der Jahresverläufe pro Verkäufer hilft nicht wirklich

Kein klares Bild.

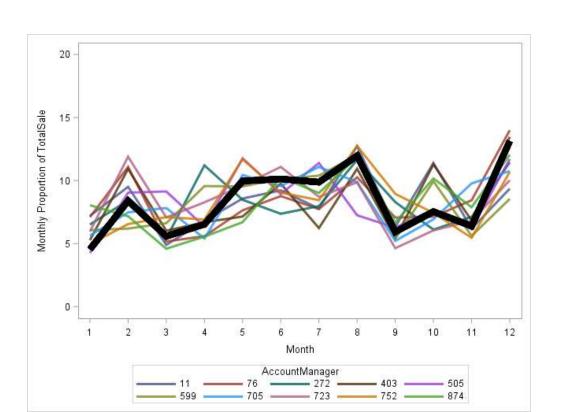
Unmöglich, alle Linien einzeln durchzusehen.





Ranking der Verkäufer mit analytischen Methoden (1)

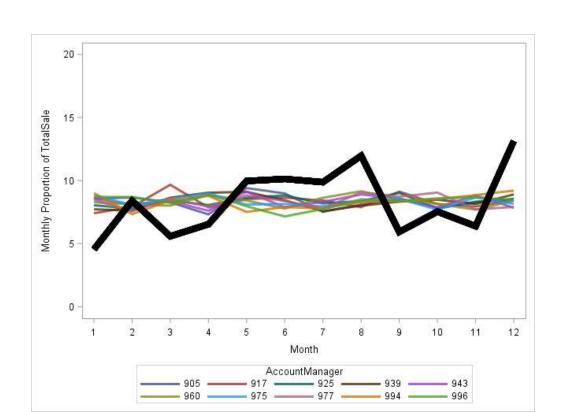
Top 10 Verkäufer bzgl. "Alignment" mit der Vorgabe





Ranking der Verkäufer mit analytischen Methoden (2)

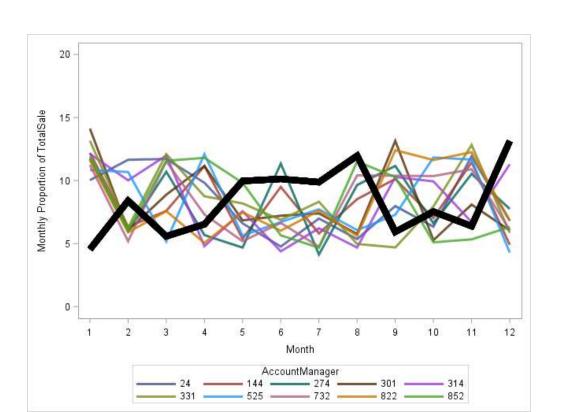
Top 10 Verkäufer, für die es keine saisonale Variation gibt.





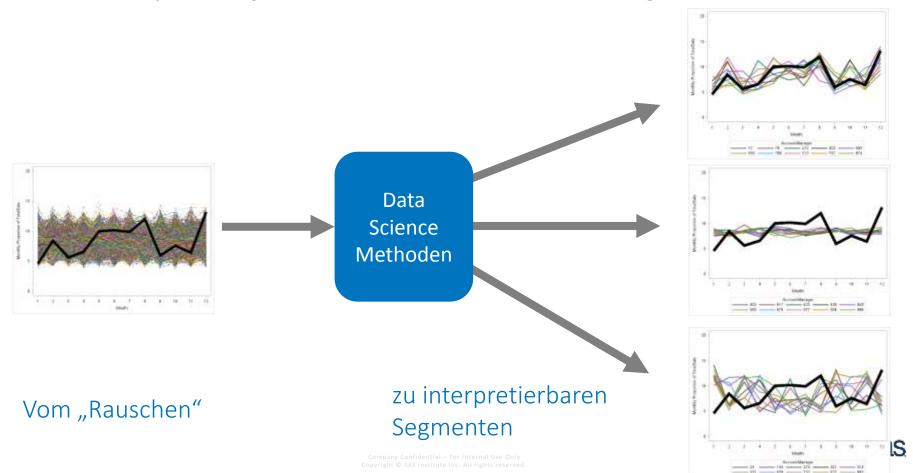
Ranking der Verkäufer mit analytischen Methoden (3)

Top 10 Verkäufer die "gegen" das Muster arbeiten





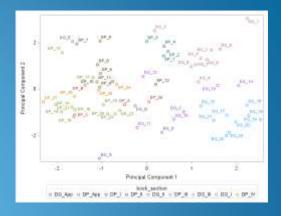
Analytik hilft mir, ein klareres Bild zu gewinnen!

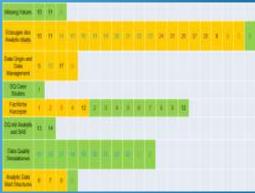


Data Science in Action: #8

Topic Search Documents and Clustering

Can I automatically find clusters of documents with similar content?





Text Mining
Text Parsing (Synonyme, Stemming, Stop-Listen)
Term by Document Weights



Kann ich ähnliche Kapitel erkennen, ohne die Bücher (von Gerhard ©) erst lesen zu müssen?

Topic > +access.+file.+text.+relational.+relational database



PAGE 104 Data Preparation for Analytics Using SAS Chapter 13; Accessing Data PAGE 103 Part 3 Data Mart Coding and Content Chapter 13 Access Transposing One- and Multiple-Rows-per-Subject Data Structures 115 Chapter 15 Transposing Longitudinal Data 131 Chapter 16 Transformations of Chapter 17 Transformations of Categorical Variables 161 Chapter 18 Multiple Interval-Scaled Observations per Subject 179 Chapter 19 Multiple Categorical Variables 161 Chapter 18 Multiple Categorical Observations of Categorical Variables 161 Chapter 18 Multiple Categorical Observations of Categorical Variables 161 Chapter 18 Multiple Categorical Observations of Categorical Variables 161 Chapter 18 Multiple Interval-Scaled Observations per Subject 179 Chapter 19 Multiple Categorical Variables 161 Chapter 18 Multiple Interval-Scaled Observations per Subject 179 Chapter 19 Multiple Categorical Variables 161 Chapter 18 Multiple Interval-Scaled Observations per Subject 179 Chapter 19 Multiple Categorical Variables 161 Chapter 18 Multiple Interval-Scaled Observations per Subject 179 Chapter 19 Multiple Categorical Variables 161 Chapter 18 Multiple Interval-Scaled Observations per Subject 179 Chapter 19 Multiple Interval-Scaled Observations per Subject 179 Chapter 19 Multiple Interval-Scaled Observations per Subject 179 Chapter 19 Multiple Interval-Scaled Observations per Subject 19



PAGE 38 Data Preparation for Analytics Using SAS Chapter 5: The Origin of Data PAGE 43 Part 2 Data Structures and Data Modeling Chapter 5 The Models 45 Chapter 7 Analysis Subjects and Multiple Observations 51 Chapter 8 The One-Row-per-Subject Data Mart 61 Chapter 9 The Multiple-Rows-p Data Structures for Longitudinal Analysis 77 Chapter 11 Considerations for Data Marts 89 Chapter 12 Considerations for Predictive Modeling 95 Introdu



Multiple-Rows-per-Subject Data | Sets 372 B.6 Selected Features of the SAS Language for Data Management 375 B.7 Benefits of the SAS Macro Language for Data Management 375 B.7 Benefits of the SAS Management 375 B.7 Benefits of



PAGE 382 Data Preparation for Analytics Using SAS Appendix B: The Power of SAS for Analytic Data Preparation PAGE 381 Appendix B The Power of 369B.1 Motivation B.2 Overview 370 B.3 Extracting Data from Source Systems 371 B.4 Changing the Data Mart Structure: Transposing 371 B.5 Data Mar

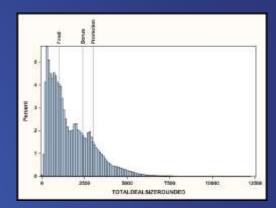
PAGE 178 Data Preparation for Analytics Using SAS Chapter 17: Transformations of Categorical Variables PAGE 177 Chapter 17 Transformations Introduction 17.2 General Considerations for Categorical Variables 162 17.3 Derived Variables 164 17.4 Combining Categories 166 17.5 Dummy Codin Multidimensional Categorical Variables 172 17.7 Lookup Tables and External Data 176 17.1 Introduction In this chapter we will deal with transformation 40 Data Quality for Analytics Using SAS Chapter 3: Data Availability 41 Chapter 3: Data Availability 3.1 Introduction 32 3.2 General Considerations 32 Re: data availability 32 Availability and usability 32 Effort to make data available 33 Dependence on the operational process 33 Availability and alignment in t of Historic Data 34 Categorization and examples of historic data 34 The length of the history 35 Customer event histories 35 Operational systems and a

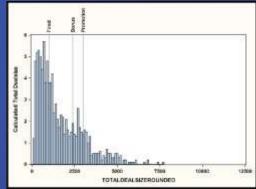
Data Science in Action: #9

Using Monte Carlo Simulations to Understand the Outcome Distribution

When the sales manager looks at the project pipeline, does the sum of weighted averages give him or her a full picture?

Monte Carlo simulations
Mathematical programming



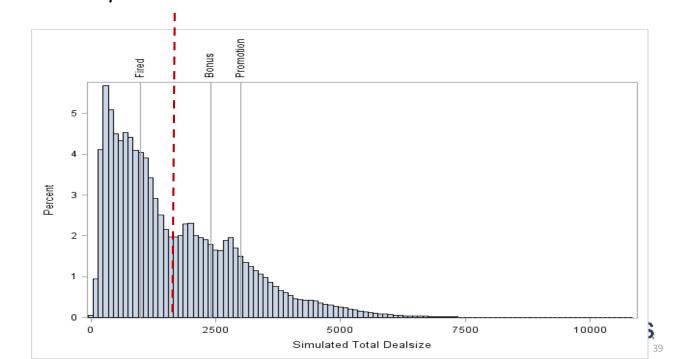




Wird der Sales Manager seinen Job behalten?

ProjectI	DealSize	Proba-
D	(1000 \$)	bility
1	1500	10%
2	10	65%
3	500	20%
4	50	50%
5	100	40%
6	30	90%
7	10	60%
8	150	20%
9	200	25%
10	180	10%
11	900	10%
12	750	20%
13	600	10%
14	320	20%
15	100	40%
16	50	80%
17	2000	5%
18	400	20%
19	2500	10%
20	1700	15%

Gewichtetes Mittel: \$ 1.661.500

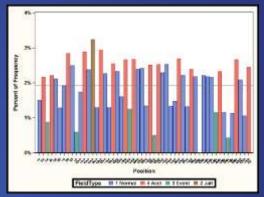


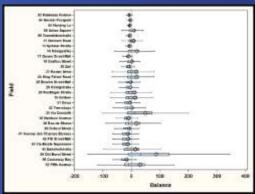
Data Science

in Action: #10

Studying Complex Systems – Simulating the Monopoly Board Game

How can you simulate complex environments to get insight in the most frequent processes?





Monte Carlo Simulations



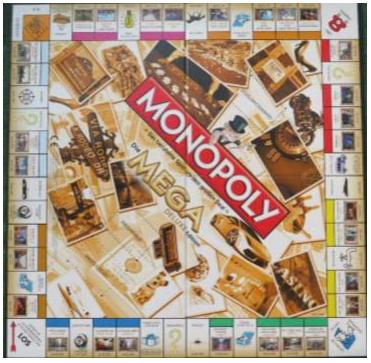
Das Monopoly Spiel ist vielen Frameworks im Geschäftsleben gar nicht so unähnlich



Komplexe Regeln



Zusätzliche Anweisungen



Rahmenwerk von Möglichkeiten und Ereignissen Alleignissessen



Monetäre Dimension



Dynamische Komponenten

ZufälligeKomponenten

Simulation komplexer Prozesse erlaubt mir Einblick in Zusammenhänge (die ich sonst nicht gesehen hätte)



Würfel-Summe





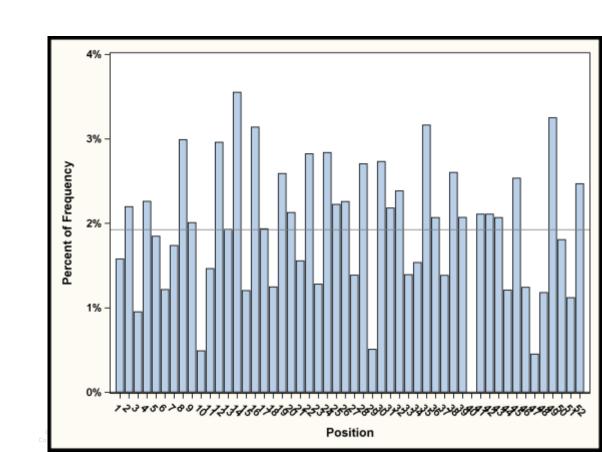
Gehe ins Gefängnis!



Ereignis Felder



Accelerator Würfel



Further Links and Downloads



- Cases #1-2, 4-7, 9-10:
 - http://www.sascommunity.org/wiki/Applying Data Science Business Case Studies Using SAS
 - http://www.sascommunity.org/wiki/DOWNLOAD SECTION: Applying Data Science Business Case Studies Using SAS
- #1 Survival
 - SAS/STAT® 14.2 User's Guide. The LIFETEST Procedure. http://support.sas.com/documentation/onlinedoc/stat/142/lifetest.pdf (accessed 1 March 2017).
 - Allison, P. 1995. Survival Analysis Using SAS®: A Practical Guide, Second Edition. Cary, NC: SAS Institute Inc.
- #2 Detecting Breakpoints and Outliers
 - Kuhfeld, W., and W. Cai. 2013. "Introducing the New ADAPTIVEREG Procedure for Adaptive Regression." SAS Global Forum Proceedings. http://support.sas.com/resources/papers/proceedings13/457-2013.pdf (Paper 457-2013).
- #3 Individual Reference Values: http://www.sascommunity.org/wiki/Data Quality for Analytics
- #4 Forecast Error Analysis
 - SGF2018 Paper 1673 Getting More Insight into Your Forecast Errors with the GLMSELECT and QUANTSELECT Procedures
 - KSFE 2015: Gerhard Svolba: Mehr als linear oder logistisch ausgewählte Möglichkeiten neuer Regressionsmethoden in SAS Download the <u>presentation</u> and the <u>paper</u>



Further Links and Downloads (Forts.)

- #6 Feature Data Mining: http://www.sascommunity.org/wiki/Data Preparation for Analytics
- #8 Text Mining
 - KSFE 2017: Beitrag "SAS Text Analytics findet Zusammenhänge in Texten Ergebnisse eines Selbstversuchs"
 - <u>SAS Club</u> 2015: SAS Contextual Analysis in Action Erfahrungen aus einem Selbstversuch
- #9 Sales Manager Simulation
 - <u>SAS Club</u>: 2016, Mihai Paunescu: Simulationen und Mathematische Programmierung mit SAS
 - KSFE 2018 (to be prepared)
- #10 Monopoly Simulation
 - KSFE 2017: "Gewinnen beim Monopoly® Spiel Alles nur Zufall? Oder gibt es doch ein paar Muster, die man kennen sollte?"
 - SAS Club 2007: Simulationen und Monte-Carlo Analysen mit SAS

