



Der Business Analytics Club für SAS User

## 24. SAS CLUB

WIEN, 18. NOVEMBER 2014  
25HOURS-HOTEL



## AGENDA

- Mehr als linear oder logistisch? – Quantils Regressionen und Adaptive Splines in SAS  
*Dr. Mihai Paunescu*

- Moderne SAS Analytik Architekturen

*Mag. Gernot Engel, Ing. Phillip Manschek,  
DI Rainer Sternecker, Dr. Gerhard Svolba*

- Eine Liebeserklärung an den SAS Enterprise Guide

*Mag. Bernadette Fabits*

- SAS Analytics 4 U – Neue Möglichkeiten für Studierende und Dozenten

*DI (FH) Bettina Brandl*



Der Business Analytics Club für SAS User

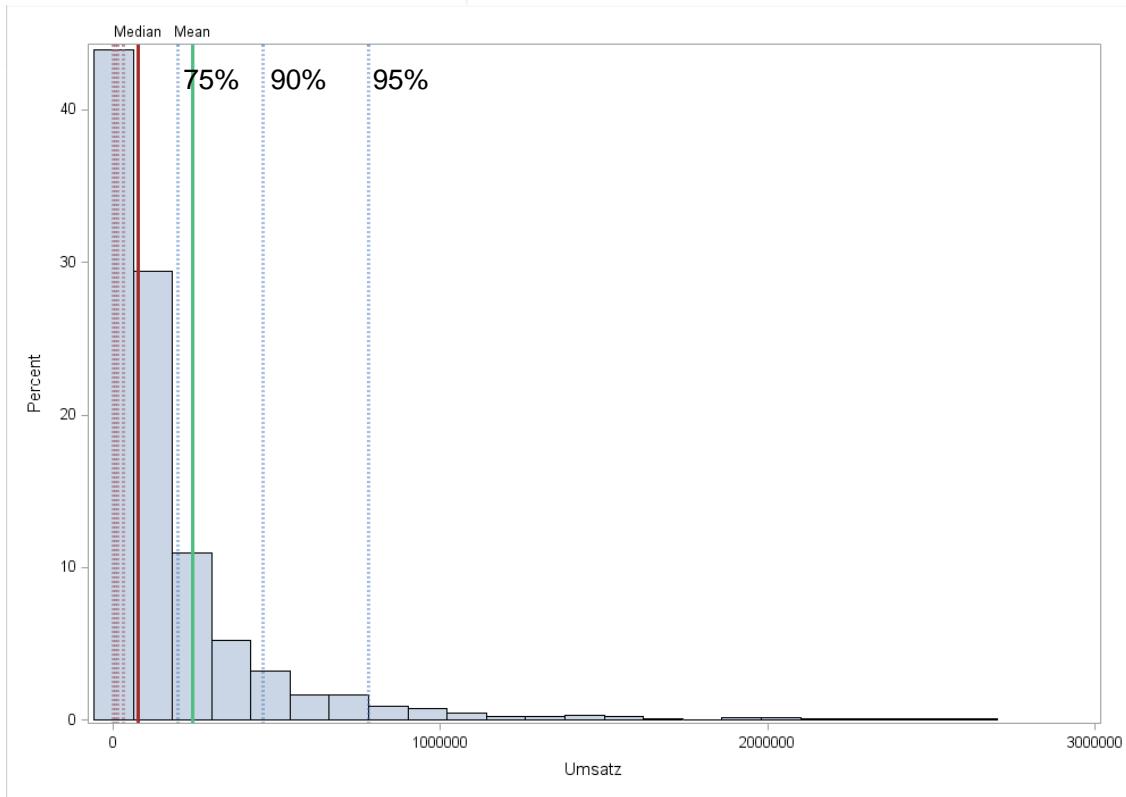
## MEHR ALS LINEAR ODER LOGISTISCH? QUANTILS REGRESSIONEN UND ADAPTIVE SPLINES IN SAS

MIHAI PAUNESCU

```
proc univariate data=dat;
ods select moments quantiles;
var sales;
run;
```

Basic Statistical Measures			
Location		Variability	
Mean	241825.8	Std Deviation	823581
Median	75875.1	Variance	6,78E+16
Mode	.	Range	18605198
		Interquartile Range	168416

Quantiles (Definition 5)	
Level	Quantile
100% Max	18605230.40
99%	2699730.30
95%	779877.70
90%	458327.90
75% Q3	198564.00
50% Median	75875.10
25% Q1	30147.60
10%	13607.10
5%	5598.60
1%	500.75
0% Min	32.55



Quantiles (Definition 5)	
Level	Quantile
100% Max	18605230.40
99%	2699730.30
95%	779877.70
90%	458327.90
75% Q3	198564.00
50% Median	75875.10
25% Q1	30147.60
10%	13607.10
5%	5598.60
1%	500.75
0% Min	32.55

```

ods graphics on / width=1000px height=700px;
proc sgplot data=dat;
  where sales < 2699730.30;
  histogram sales;
  refline 779877.70 458327.90 198564.00 / axis=x lineattrs=(color=bigb pattern=dot thickness=3);
  refline 30147.60 13607.10 5598.60 500.75 /axis=x lineattrs=(color=dapk pattern=dot thickness=3);
  refline 75875.10 /axis=x label='Median' lineattrs=(color=brown thickness=3);
  refline 241825.8 /axis=x label='Mean' lineattrs=(color=big thickness=3);
run;

```

```
proc quantreg data=dat;
  model Sales= / quantile=(0.5);
run;
```

Parameter Estimates						
Parameter	DF	Estimate	Standard Error	95% Confidence Limits	t Value	Pr >  t
Intercept	1	75875.10	2.654.220	70.670.691 81.079.509	28.59	<.0001

Quantiles (Definition 5)	
Level	Quantile
100% Max	18605230.40
99%	2699730.30
95%	779877.70
90%	458327.90
75% Q3	198564.00
<b>50% Median</b>	<b>75875.10</b>
25% Q1	30147.60
10%	13607.10
5%	5598.60
1%	500.75
0% Min	32.55

```
proc quantreg data=dat ci=resampling ;
  class cust_grp;
  model Sales=cust_grp / quantile=(0.5) seed=12345;
run;
```

Parameter Estimates							
Parameter	DF	Estimate	Standard Error	95% Confidence Limits	t Value	Pr >  t	
Intercept	1	47335.60	3.071.371	41.313.237 53.357.963	15.41	<.0001	
Cust_grp	I	88078.60	6.903.957	74.541.278 101615.92	12.76	<.0001	
Cust_grp	O	21870.35	5.324.410	11.430.214 32.310.486	4.11	<.0001	
Cust_grp	N	0.0000	0.0000	0.0000 0.0000	.	.	

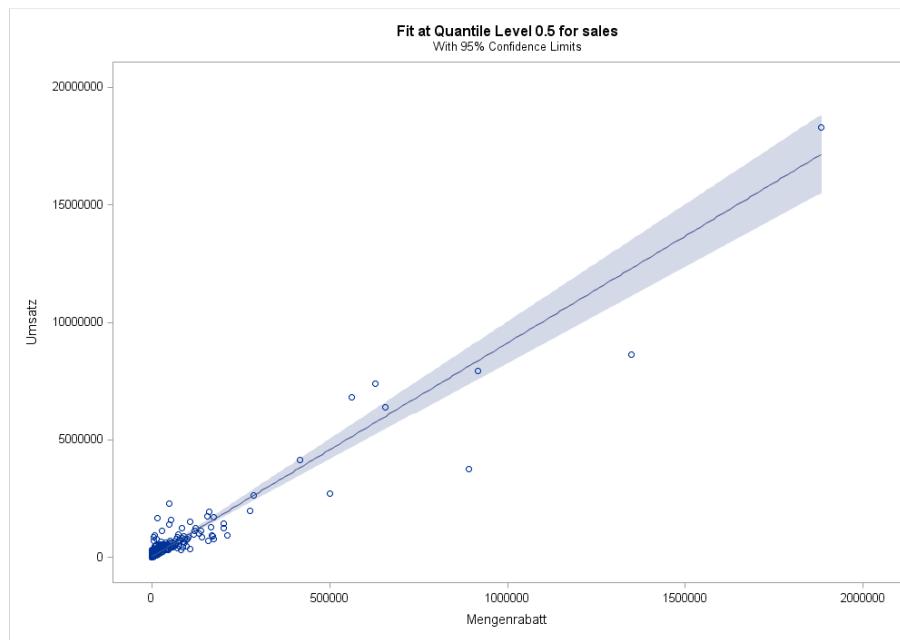
```
proc means data=dat P1 median P90;
  var sales;
  class Cust_grp;
run;
```

Analysis Variable : Sales		
Cust_grp	N Obs	Median
I	793	135414.2
O	1024	69062.08
N	1004	47323.3

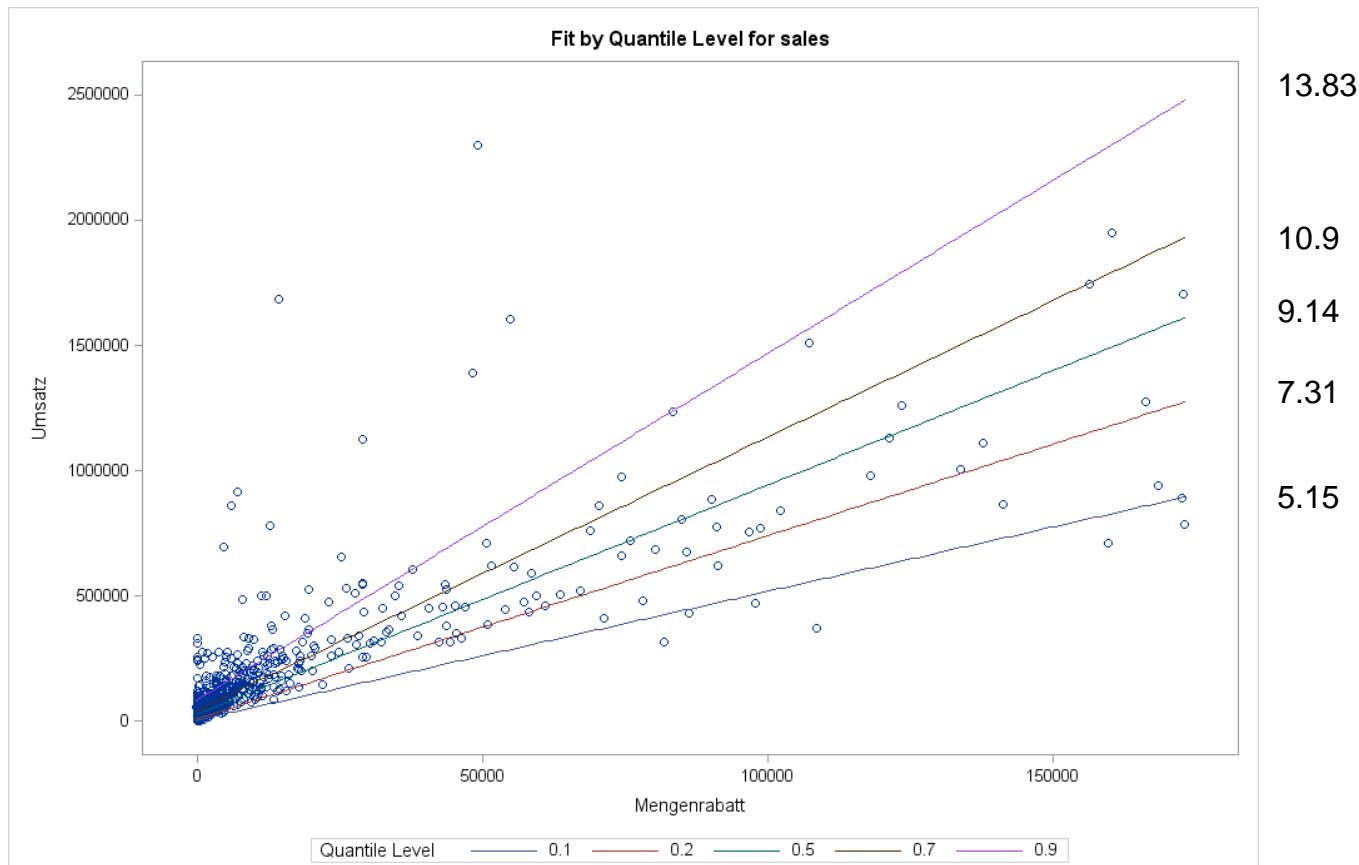
$$47335.6 + 88078.6 = 135414.2$$

```
proc quantreg data=dat ci=resampling plots=(fitplot);  
  where cust_grp='n';  
  model Sales=Rabatt_Menge / quantile=(0.5) seed=12345;  
run;
```

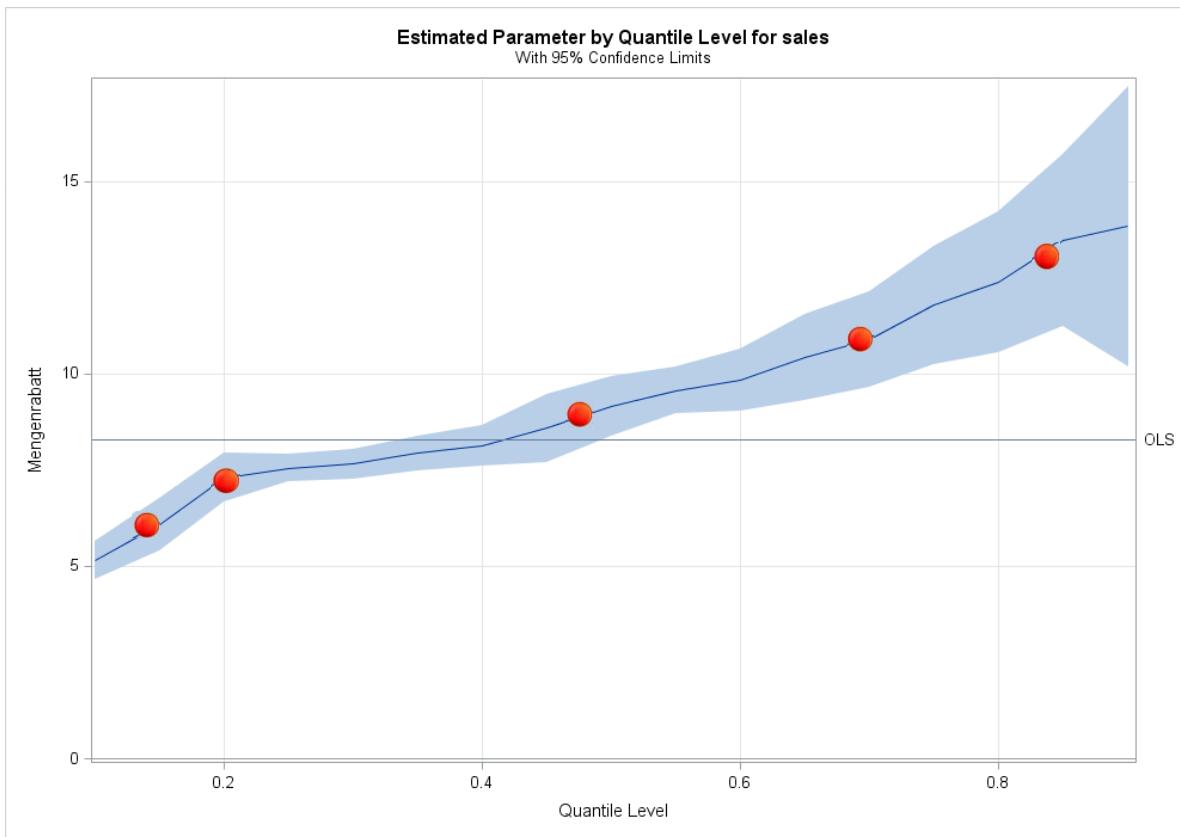
Parameter Estimates							
Parameter	DF	Estimate	Standard Error	95% Confidence Limits		t Value	Pr >  t
Intercept	1	28552.1	5873.371	17026.581	40077.615	4.86	<.0001
Rabatt_Menge	1	9.1	0.4522	8.2082	9.9829	20.12	<.0001



```
proc quantreg data=dat ci=resampling plots=(fitplot);  
ods select fitplot;  
where cust_grp='n' and Rabatt_Menge < 200000;  
model Sales=Rabatt_Menge / quantile=(0.1 0.2 0.5 0.7 0.9) seed=12345;  
run;
```



```
proc quantreg data=dat ci=resampling;
ods select quantplot;
where cust_grp='n' and Rabatt_Menge < 200000;
model Sales=Rabatt_Menge
/quantile=(0.1 to 0.9 by 0.05) plot=(quantplot /unpack ols) seed=1268 ;
run;
```



13.83

10.9

9.14

7.31

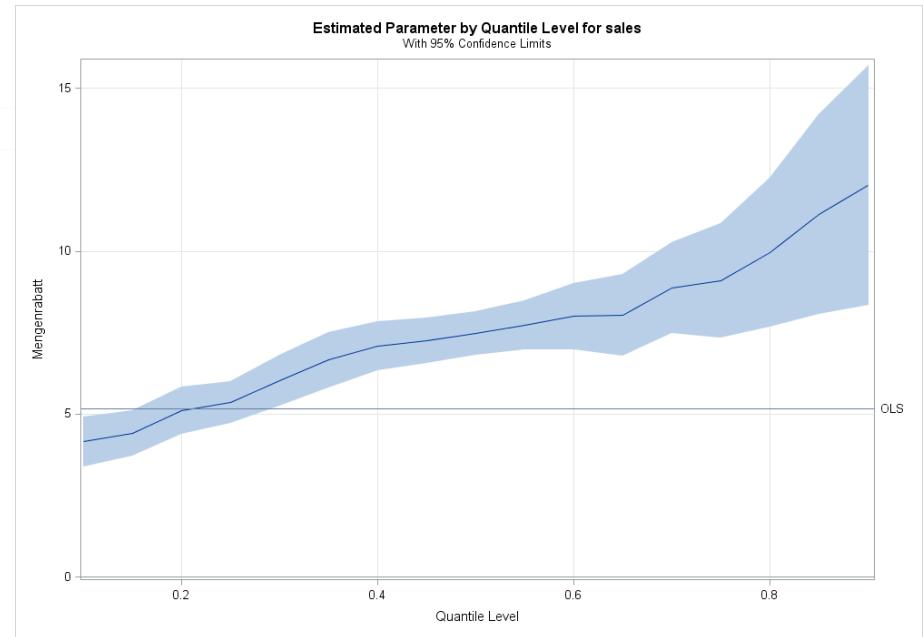
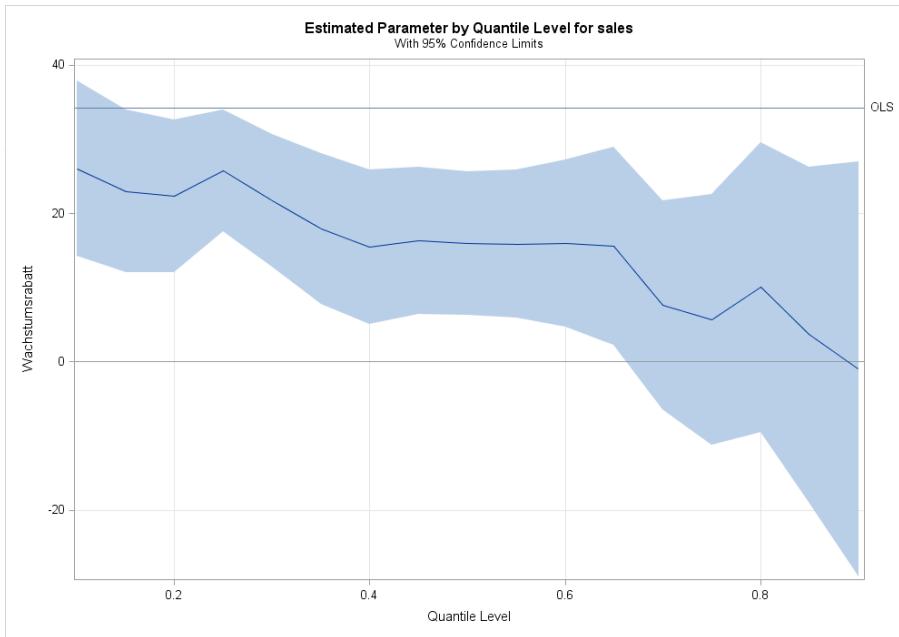
5.15

## QUANTILS- REGRESSION

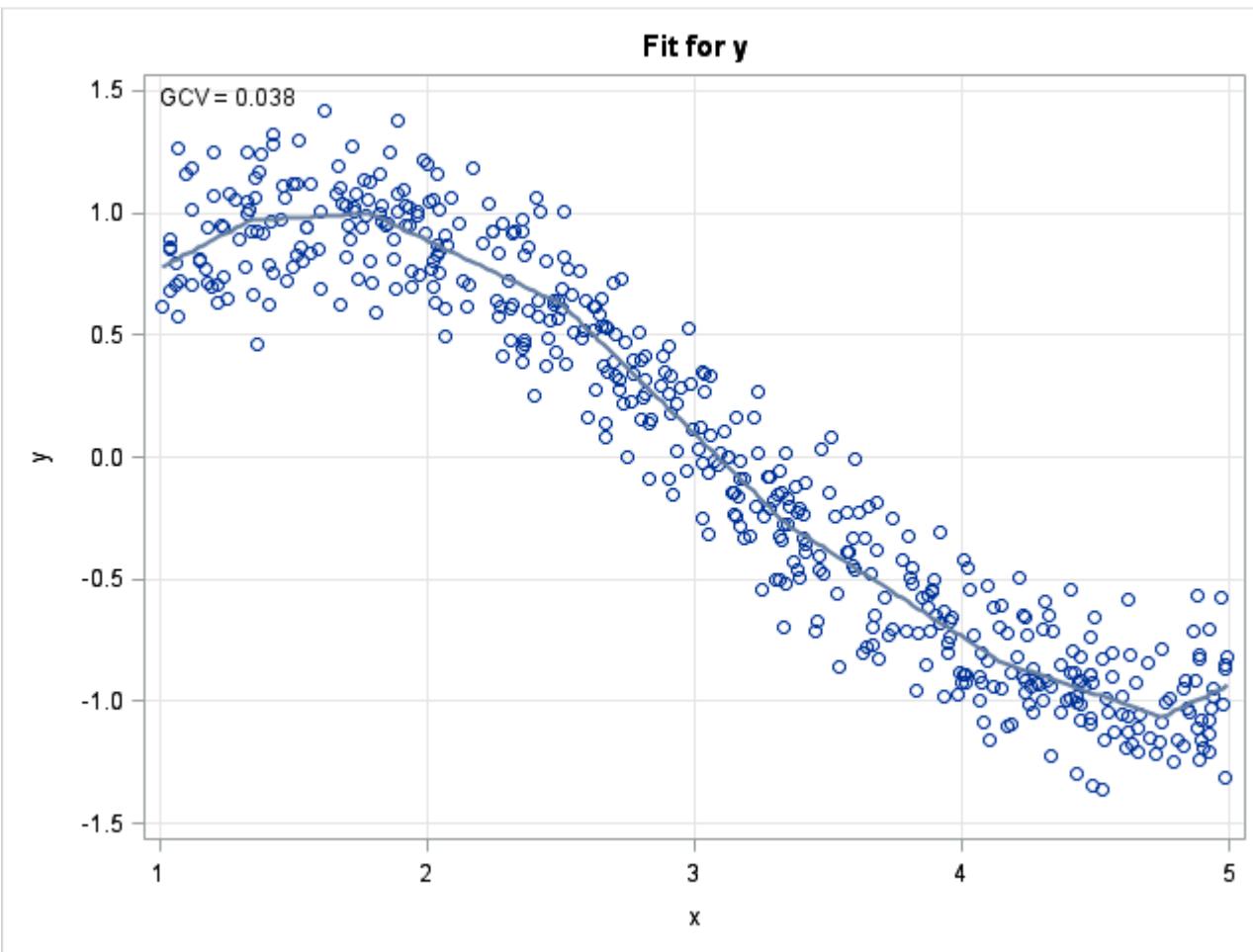
## MULTIVARIATE QUANTILSREGRESSION

sas®club

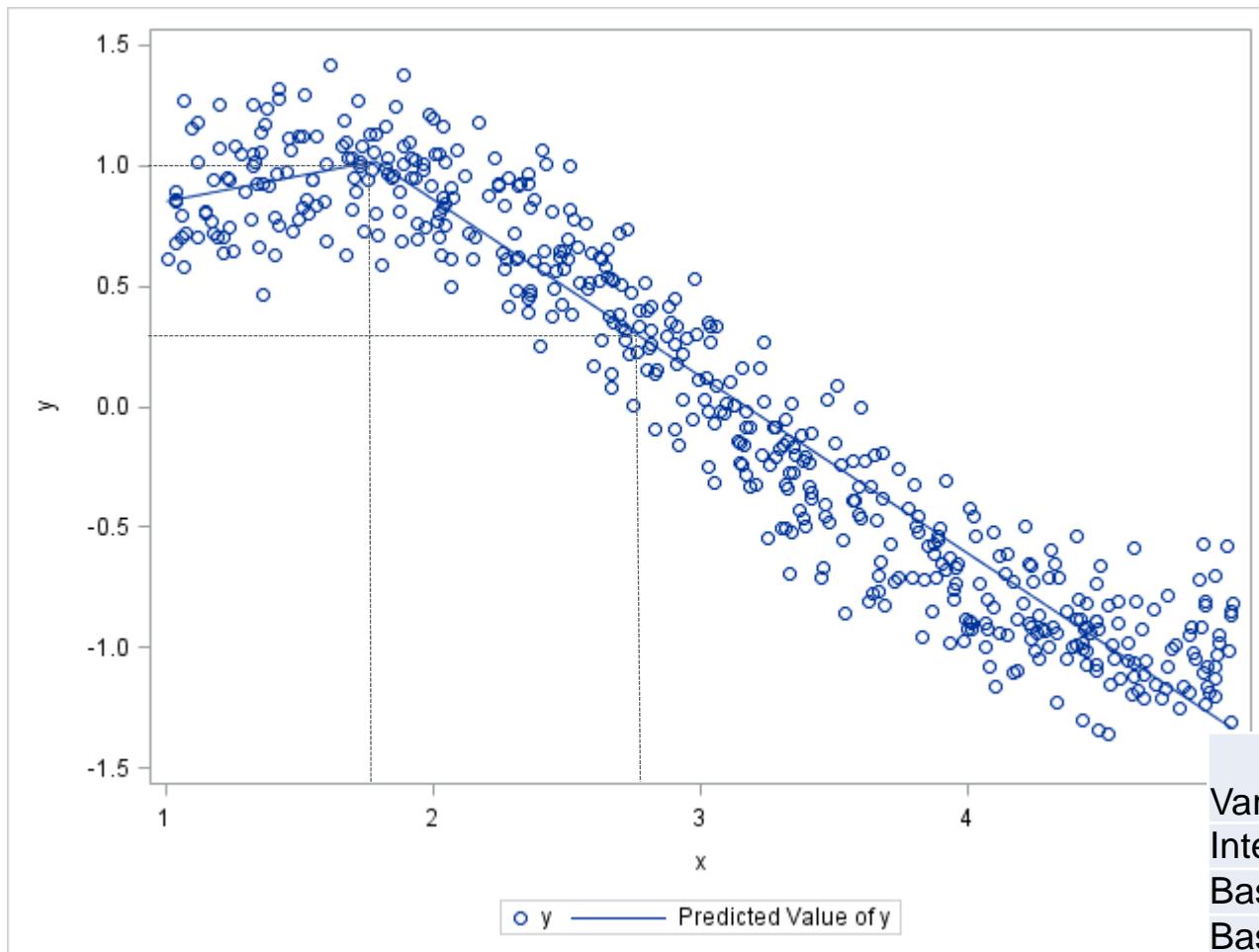
```
proc quantreg data=dat ci=resampling;  
where cust_grp='n';  
model Sales=Rabatt_Wachstum Rabatt_Menge /  
quantile=(0.1 to 0.9 by 0.05) plot=(quantplot /unpack ols) seed=1268 ;  
run;
```



- Verwendung:
  - Zusammenhänge entdecken für extreme Bereiche der Zielvariable
  - Robuste Medianschätzung gegenüber Ausreißer, ohne Verteilungsannahmen
- Hinweise:
  - Ist nicht äquivalent zu linearen Regressionen für Segmente von Beobachtungen
  - Schneller: HPQUANTSELECT ab SAS/STAT 13.2



```
proc adaptivereg  
plots=all  
details=bases;  
  
model y = x;  
  
run;
```

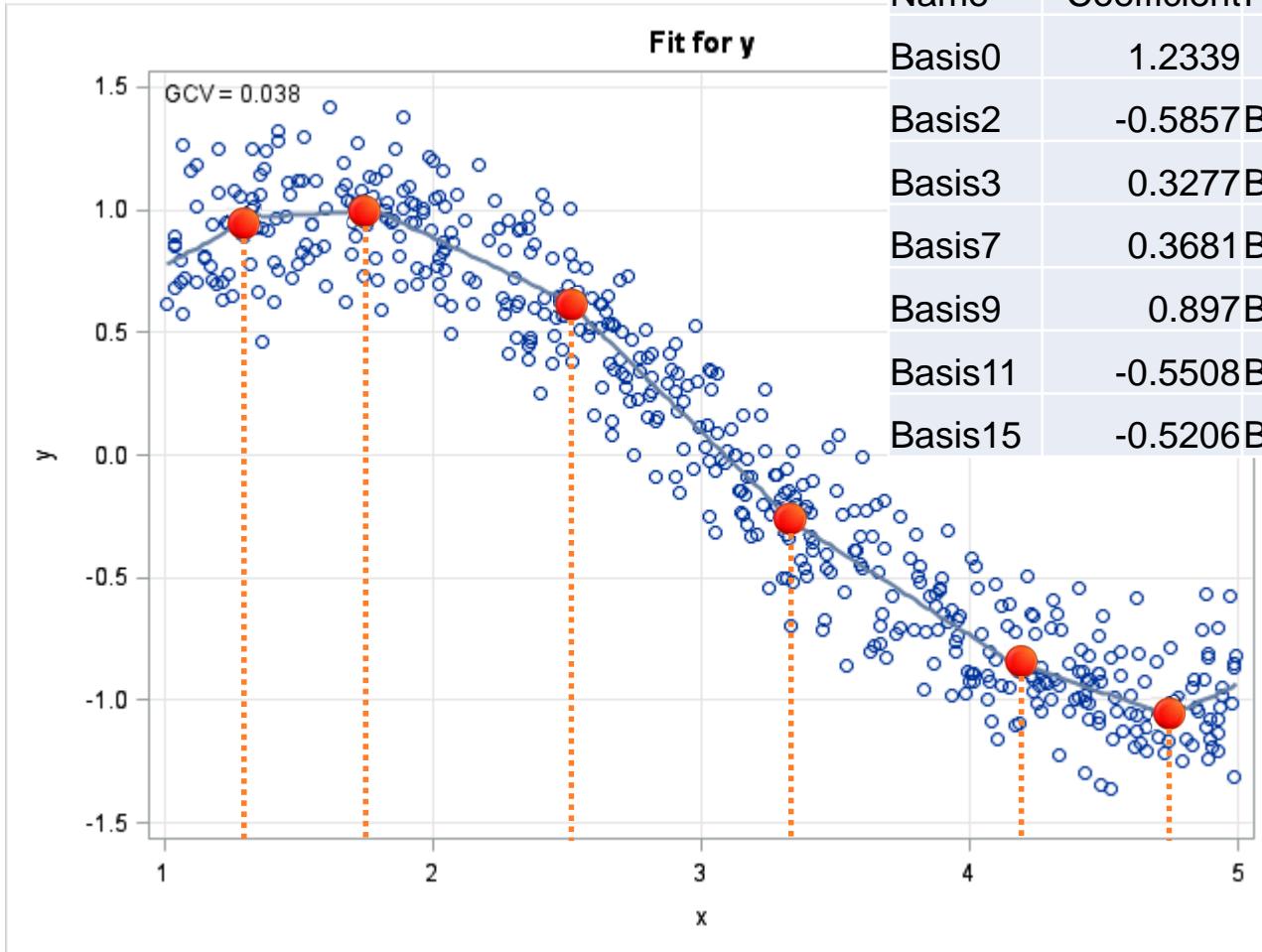


Basis1	$\text{MAX}(x - 1.8, 0)$
Basis2	$\text{MAX}(1.8 - x, 0)$

```
data ds2;
set ds;
Basis1 = max(x-1.8, 0);
Basis2 = max(1.8 - x, 0);
run;

proc reg data=ds2;
model y = basis1 basis2;
run;
```

Variable	Parameter			
	Estimate	t Value	Pr >  t	
Intercept	1.02	52.04	<.0001	
Basis1	-0.73	-67.82	<.0001	
Basis2	-0.21	-3.31	0.001	

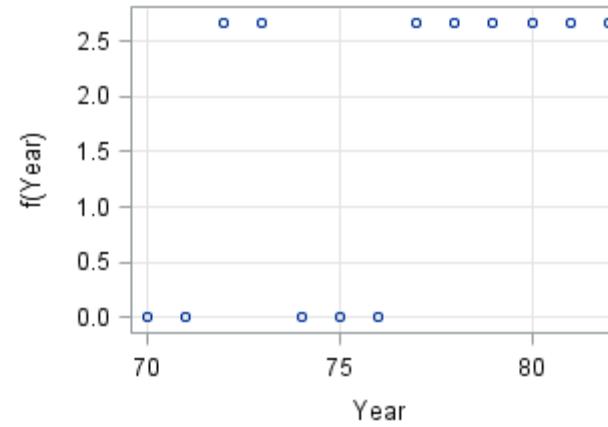
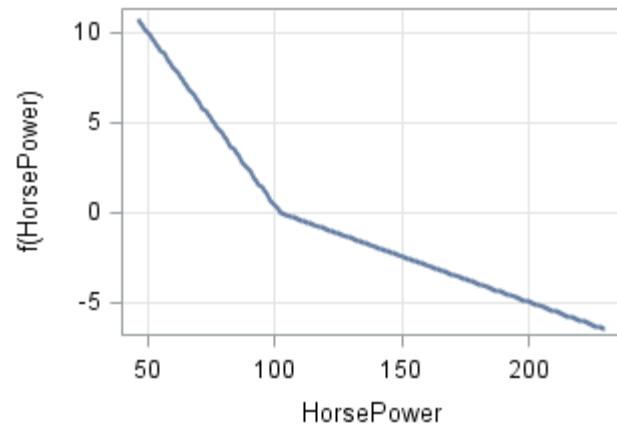


Regression Spline Model after Backward Selection			
Name	Coefficient Parent	Variable	Knot
Basis0	1.2339	Intercept	
Basis2	-0.5857Basis0	x	1.7865
Basis3	0.3277Basis0	x	4.1447
Basis7	0.3681Basis0	x	3.3424
Basis9	0.897Basis0	x	4.7489
Basis11	-0.5508Basis0	x	2.5061
Basis15	-0.5206Basis0	x	1.3345

```
proc adaptivereg data=autompg plots=all details=bases;
class cylinders year origin;
model mpg = cylinders displacement horsepower
weight acceleration year origin / additive;
run;
```

Basis Information	
Name	Transformation
Basis0	1
Basis1	Basis0*MAX(Weight - 3139,0)
Basis2	Basis0*MAX(-3139 - Weight,0)
<b>Basis3</b>	<b>Basis0*NOT(MISSING(HorsePower))</b>
<b>Basis4</b>	<b>Basis0*MISSING(HorsePower)</b>
<b>Basis5</b>	<b>Basis3*MAX(HorsePower - 102,0)</b>
<b>Basis6</b>	<b>Basis3*MAX(-102 - HorsePower,0)</b>
Basis7	Basis0*(Year = 80 OR Year = 82 OR Year = 81 OR Year = 79 OR Year = 78 OR Year = 77 OR Year = 73 OR Year = 72)
Basis8	Basis0*NOT(Year = 80 OR Year = 82 OR Year = 81 OR Year = 79 OR Year = 78 OR Year = 77 OR Year = 73 OR Year = 72)
Basis9	Basis0*MAX(Displacement - 85,0)
Basis10	Basis0*MAX(-85 - Displacement,0)
Basis11	Basis0*MAX(Displacement - 97,0)
Basis12	Basis0*MAX(-97 - Displacement,0)
Basis13	Basis0*MAX(Acceleration - 21,0)
Basis14	Basis0*MAX(-21 - Acceleration,0)
Basis15	Basis3*MAX(Displacement - 105,0)
Basis16	Basis3*MAX(-105 - Displacement,0)

Basis Information		
Name	Coefficient	Transformation
Basis3	-4.03	Basis0*NOT(MISSING(HorsePower))
Basis4	0	Basis0*MISSING(HorsePower)
Basis5	-0.05	Basis3*MAX(HorsePower - 102,0)
Basis6	0.19	Basis3*MAX(102 - HorsePower,0)
Basis7	2.67	Basis0*(Year = 80 OR Year = 82 OR Year = 81 OR Year = 79 OR Year = 78 OR Year = 77 OR Year = 73 OR Year = 72)
Basis8	0	Basis0*NOT(Year = 80 OR Year = 82 OR Year = 81 OR Year = 79 OR Year = 78 OR Year = 77 OR Year = 73 OR Year = 72)



PROC Adaptivereg

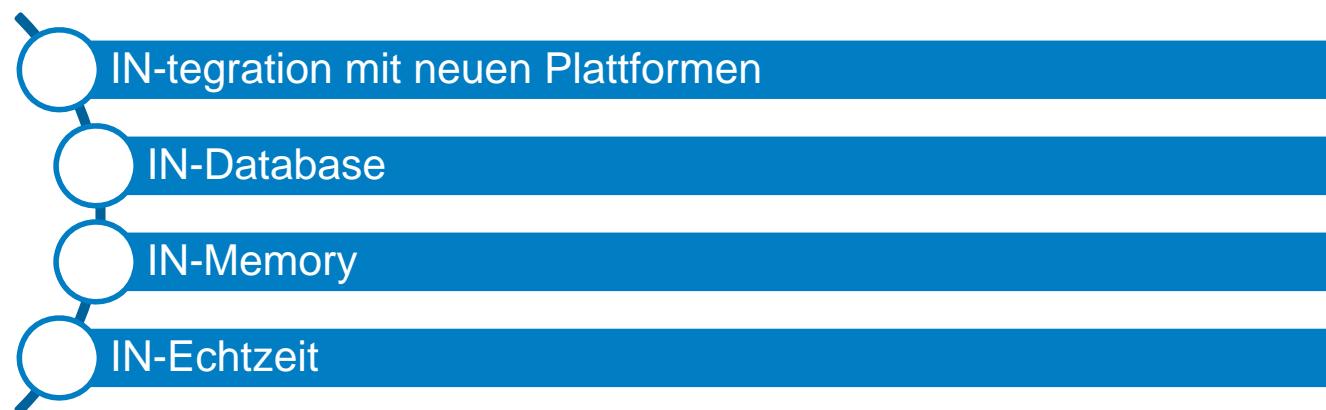
<http://support.sas.com/resources/papers/proceedings13/457-2013.pdf>

Weitere Nicht-Parametrische Regressionen: PROC GAM, PROC LOESS



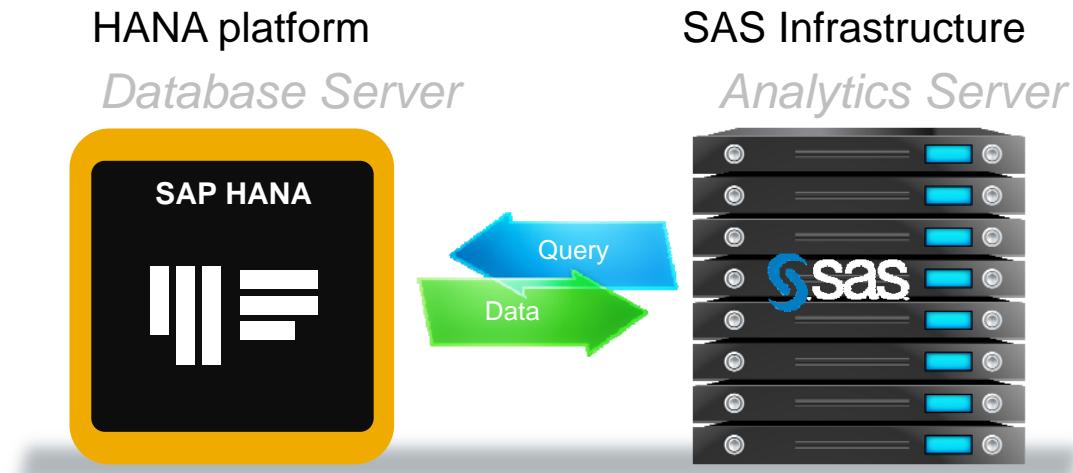
Der Business Analytics Club für SAS User

## MODERNE SAS ANALYTIK ARCHITEKTUREN

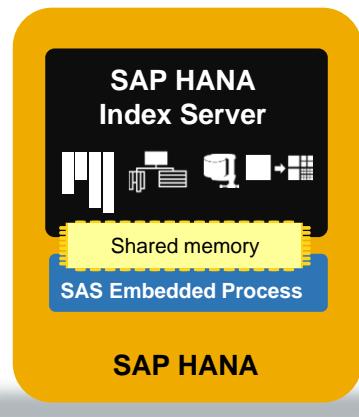


- SAS/ACCESS Interface to
  - PostGreSQL
  - Vertica
  - PI System
  - Impala
  - SAP Hana
  - Hadoop

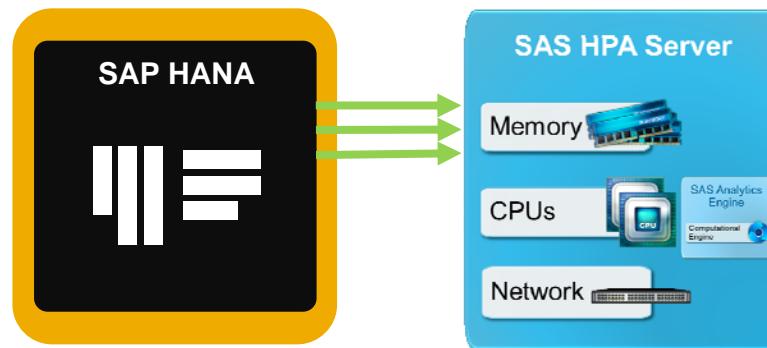
- SAS/ACCESS Interface to SAP Hana



- SAS Scoring Accelerator for SAP HANA



- SAS High-Performance Predictive Modeling Workbench for SAP HANA





Der Business Analytics Club für SAS User



## HADOOP UND SAS – STATUS UND AUSBLICK

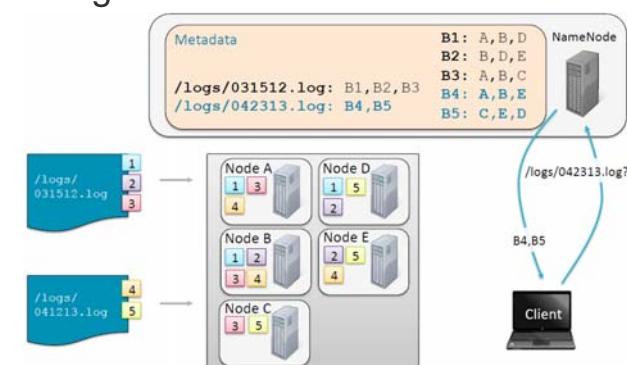
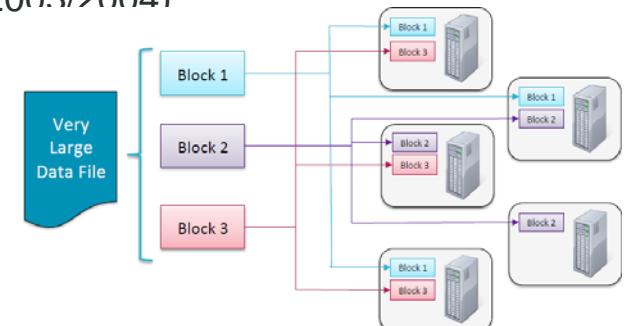
WIEN, OKT 2014  
GERNOT ENGEL, SAS AUSTRIA

## Historie:

- \*basiert auf **Google** Forschungspapieren für ein verteiltes Dateisystem(2003/2004)
- \*von **Yahoo** für die Entwicklung einer WWW-Suchmaschine aufgegriffen
- \*an die **Apache** Foundation übergeben, seitdem open-source (2009)

## Technologie:

- Prinzip der **horizontalen Skalierung** auf kostengünstiger Hardware („Scale out“)
- Prinzip der **Datenlokalität**: der Programmcode wird auf die Cluster-Nodes mit den zugehörigen Daten verteilt, dort verteilt ausgeführt und (Teil-)Ergebnisse wieder zusammengeführt
- Structure on read / late binding
- HDFS**: Hadoop Distributed File System
- Map/Reduce**: ein Verarbeitungsverfahren nach dem „Teile und herrsche“-Ansatz



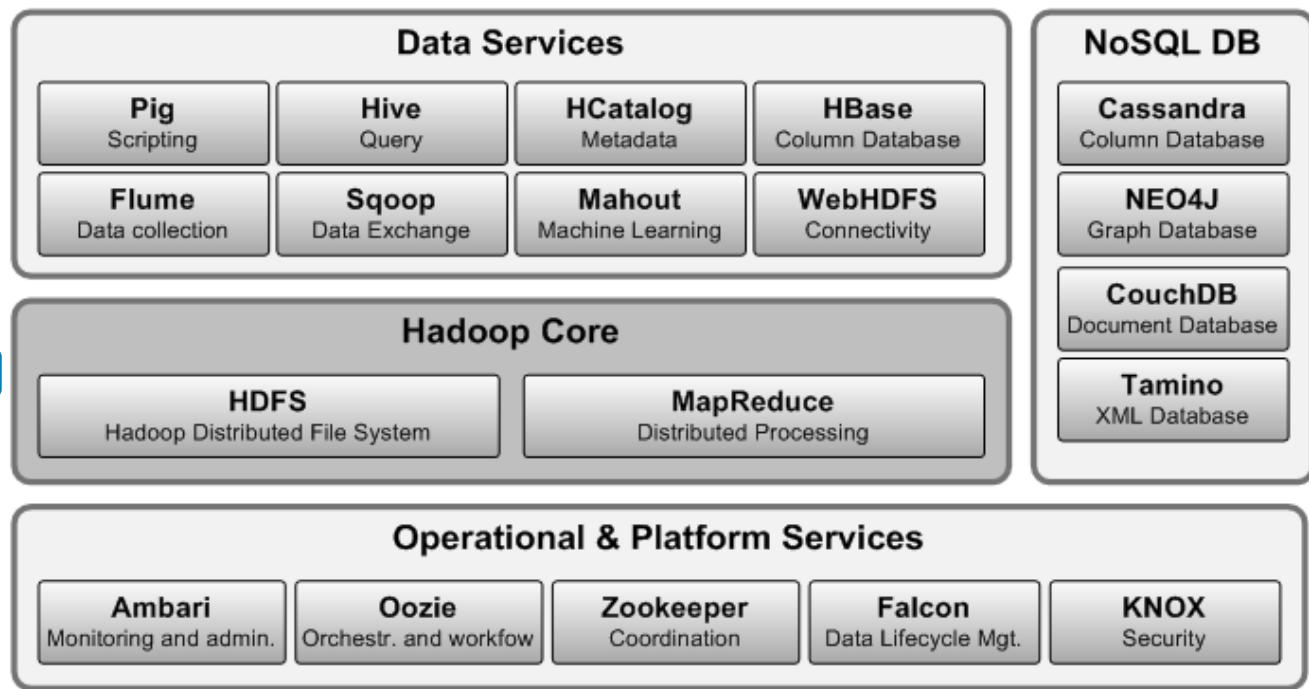
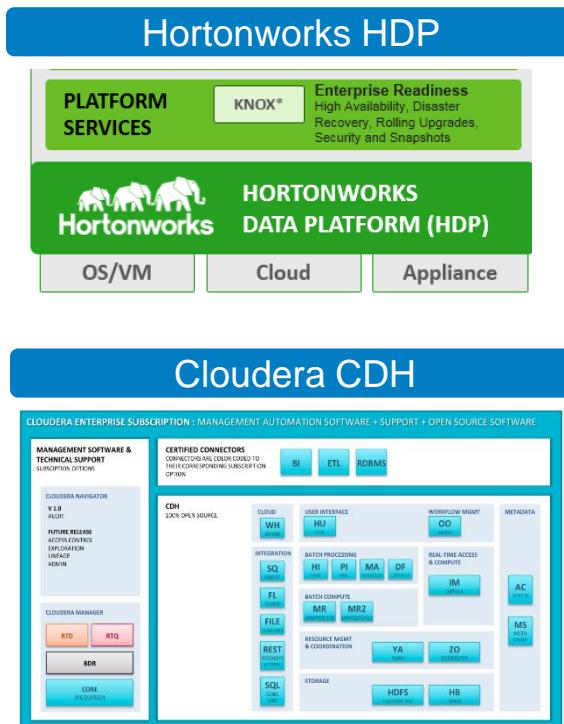
Mehr Daten

Andere Daten

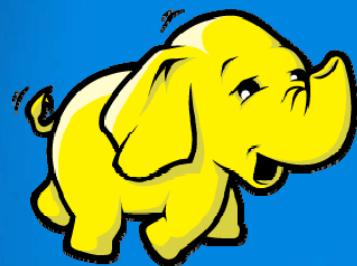
Kostengünstig

## DAS HADOOP ÖKOSYSTEM

## HADOOP IST EINE PROJEKTPLATTFORM MIT HDFS UND MAP/REDUCE ALS KERNMODULEN

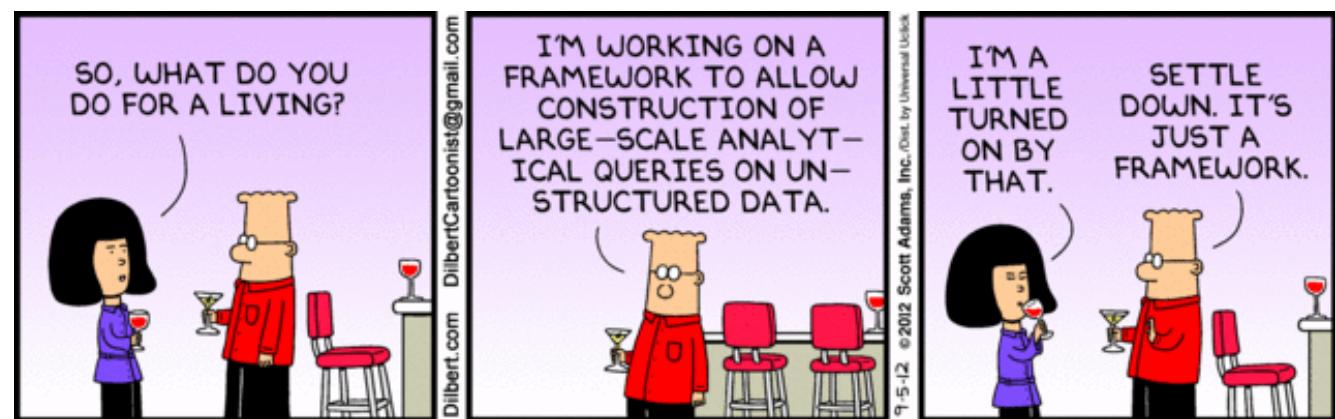


SAS Support  
 Tier1 : Cloudera, Hortonworks  
 Tier2 : MapR, Pivotal, IBM ..

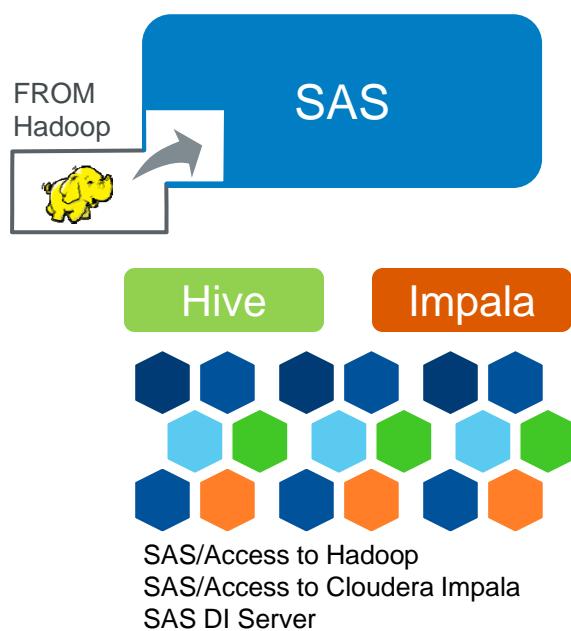


## VISION - WARUM UNTERSTÜTZT SAS HADOOP?

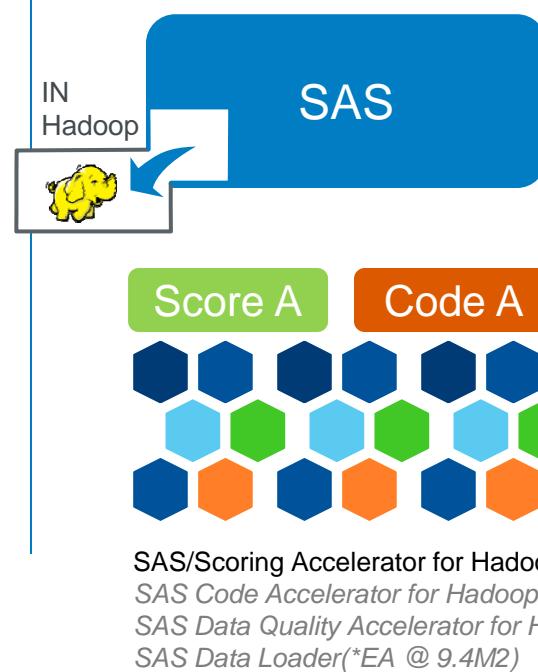
To be the Analytic and Data Management solution of choice for Hadoop.



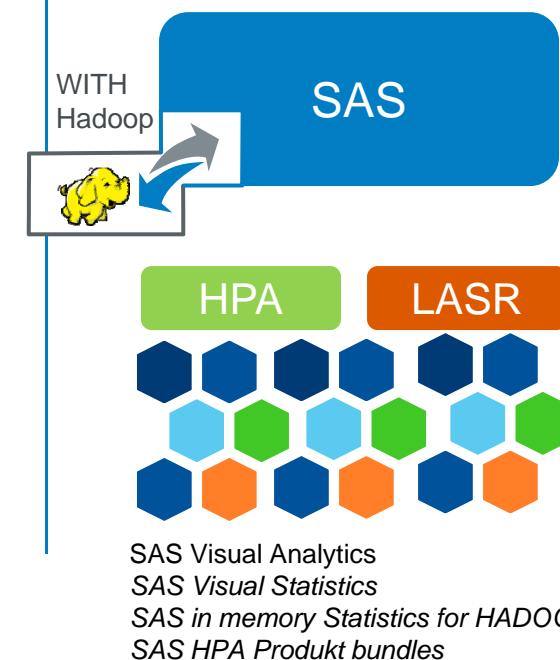
**SAS/Access to Hadoop Push some SAS processing from Hadoop into SAS**



Embedded Process - Push SAS data processing to Hadoop with Map Reduce



In-Memory Analytics - Use Hadoop for Storage persistence and commodity computing.



## FUNKTIONEN IM DATENMANAGEMENT FÜR HADOOP

### Base SAS

- Map Reduce + Pig Scripting + HDFS Kommandos

### SAS Access to Hadoop

- Hive, Hive2 + eigene Metadaten („Information Maps für HDFS Dateien“)
- Pushdown von Procs: FREQ, RANK, REPORT, SORT, SUMMARY/MEANS & TABULATE

### SAS Access to Impala (Cloudera)

### SAS Data Integration Studio Transformationen



- Read/Write HDFS files
- Submit HiveQL code
- Execute Map/Reduce code
- Submit Pig Latin
- Transfer data to/from Hadoop using Hadoop utilities
- SQL transforms pushed down with Access to Hadoop engine
- Submit DQ and DS2 code with Code Accelerator

### SAS Federation Server

- Datenvirtualisierung und Zugriffsschutz für Hadoop und andere Datenquellen

### SAS Event Stream Processing Engine

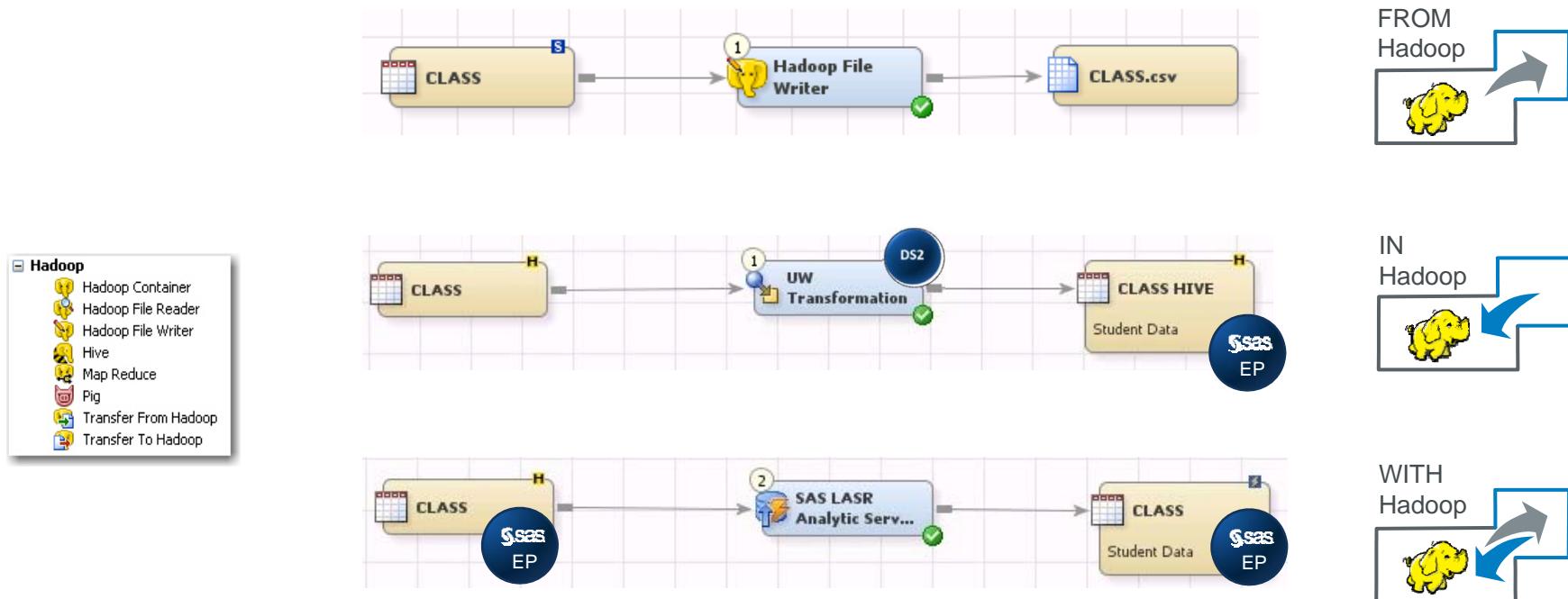
- Hadoop Adapter für SAS ESP, um Data-Streams im HDFS zu speichern

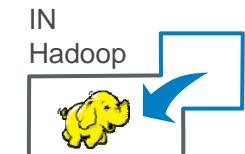
```
/* SUBMIT SQL QUERY TO HIVE */
LIBNAME MYHDP HADOOP PORT=10000 SERVER=HDPSRV02;
PROC SQL;
  INSERT INTO MYHDP.CARS_HIVE
    SELECT MAKE, MODEL, MSRP FROM SASHHELP.CARS;
QUIT;

/* SUBMIT PIG SCRIPT, HDFS COMMAND, MR JOB */
FILENAME CFG "C:\SAMPLE_DATA\HADOOP_CONFIG.XML";
FILENAME PIGCODE1 "C:\SAMPLE_DATA\PIG_CD.TXT";
PROC HADOOP OPTIONS=CFG;
  PIG CODE=PIGCODE1;
  HDFS DELETE="/USER/HADOOP/OUTPUT_MR1";
  MAPREDUCE INPUT="/..." OUTPUT="/..."
    JAR="..." MAP="..." REDUCE="...";
RUN;

/* COPY FILE FROM HDFS TO LOCAL SAS */
FILENAME CFG "C:\SAMPLE_DATA\HADOOP_CONFIG.XML";
PROC HADOOP OPTIONS=CFG USERNAME="HADOOP"
  PASSWORD="XXXX";
HDFS COPYTOLOCAL="/USER/HADOOP/TESTFOLDER"
OUT="C:\SAMPLE_DATA\" ;
RUN;
```

## SAS PRODUKTE FÜR ALLE DREI EBENEN





### SAS Scoring Accelerator for Hadoop

Ausführen von Scoring Modellen aus EM und STAT Projekten

### SAS Data Quality Accelerator for Hadoop (EA in 9.4M2)

Ausführen von DQ Routinen (Parse, Standardize, Gender Analysis, Identification, Match Code...)

### SAS Code Accelerator for Hadoop

Ausführen von DataStep2 Code

### SAS Data Loader for Hadoop (EA in 9.4M2)

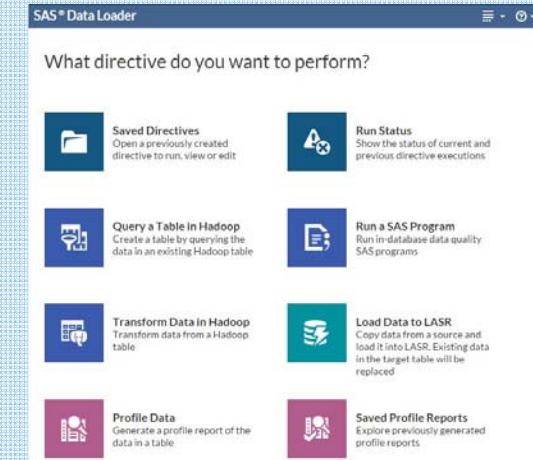
Web-basiertes In-Database Tool für Endanwender (z.B. aus den Fachbereichen)

#### Steckbrief SAS Data Loader 2.1

- Point & Click Datenmanagement-Routinen, die in Hadoop ausgeführt werden (Abfragen, Tabellen anlegen, SAS Programme)
- In-Memory Beladung: Transfer von Hadoop-Daten in den LASR Server
- “Profile Data” für Hadoop-Daten
- HTML5 basiertes Interface
- Wird als vApp ausgeliefert (läuft in VMWare Player)

#### Roadmap

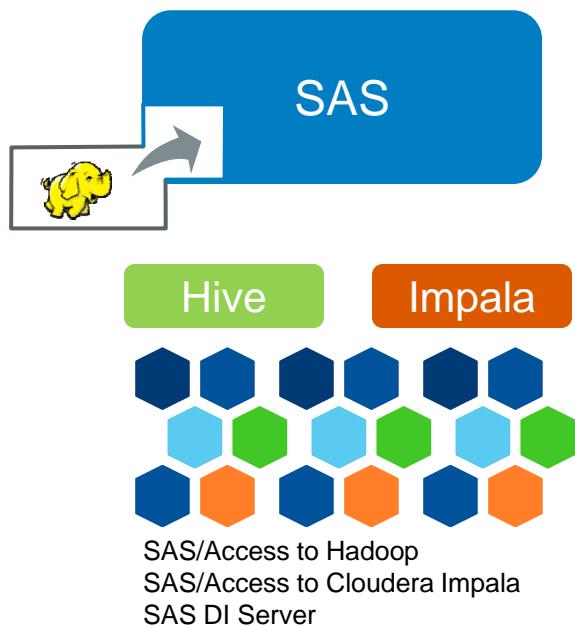
- Joins in Hadoop
- Sqoop Support
- Data Quality Verfahren
- Metadaten / Lineage
- Security
- Monitoring



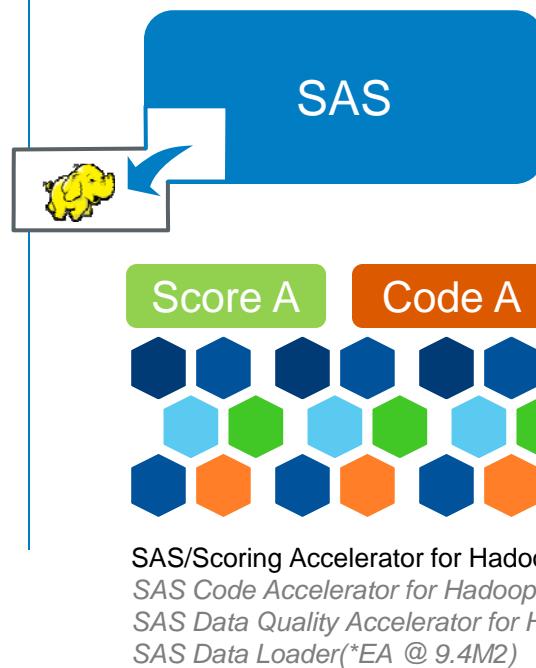
## SAS & HADOOP

## BASIS TECHNOLOGIEN

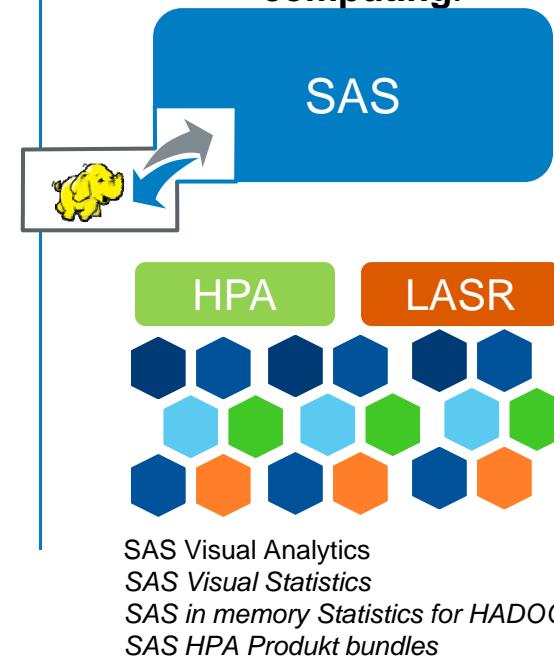
SAS/Access to Hadoop Push some SAS processing from Hadoop into SAS



Embedded Process - Push SAS data processing to Hadoop with Map Reduce



**In-Memory Analytics - Use Hadoop for Storage persistence and commodity computing.**

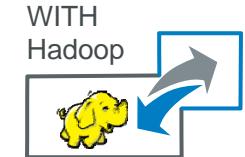


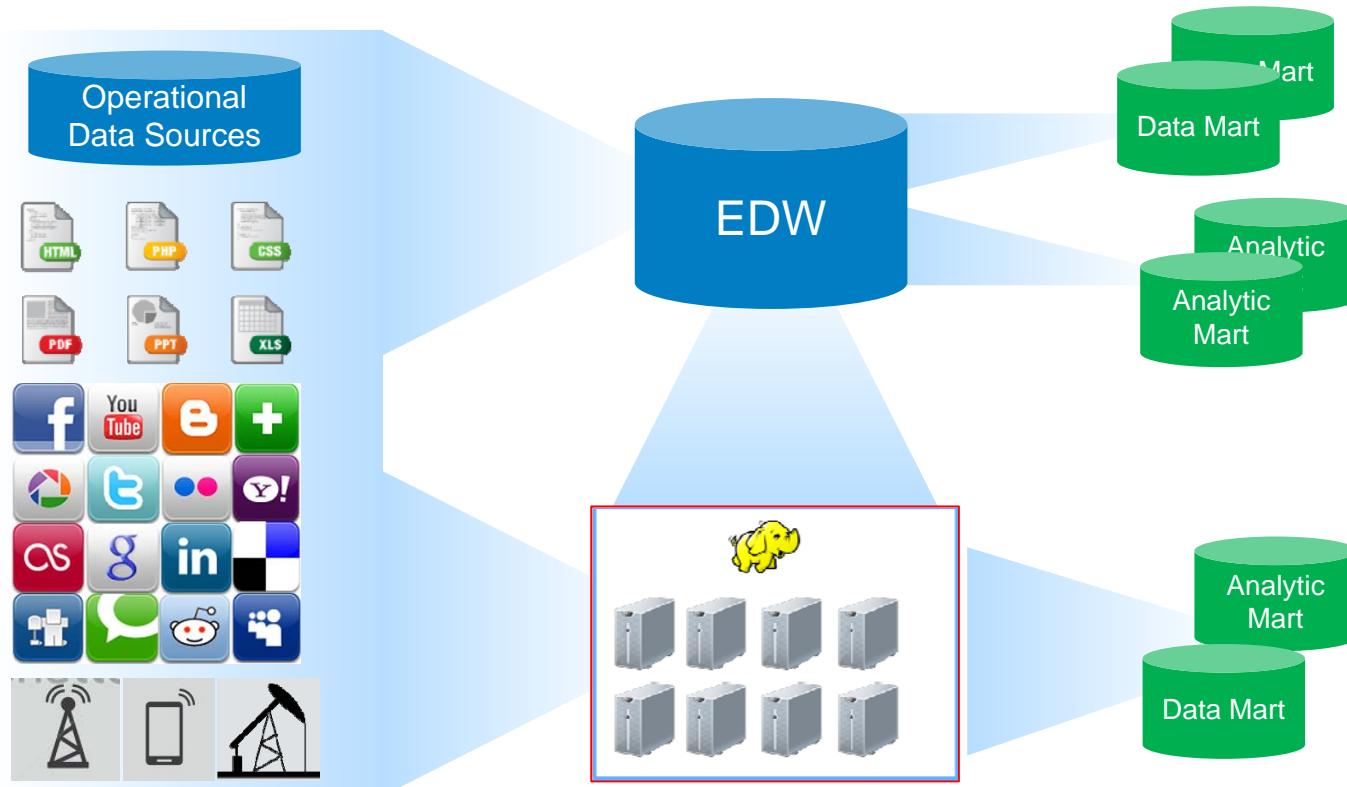
The screenshot shows the SAS Studio interface. On the left, the file 'demo\_short\_with\_ODS\_10032013\_v3.smpUSE.sas' is open in the code editor, displaying SAS code for data processing and statistical analysis. The results viewer on the right displays several tables and a bar chart. One table is titled 'Pairwise Correlations for Table WORK.CARDATAP' and another is 'Summary Statistics for Table WORK.CARDATAP'. The bar chart, titled 'Formatted', shows the frequency of various categories across different columns.

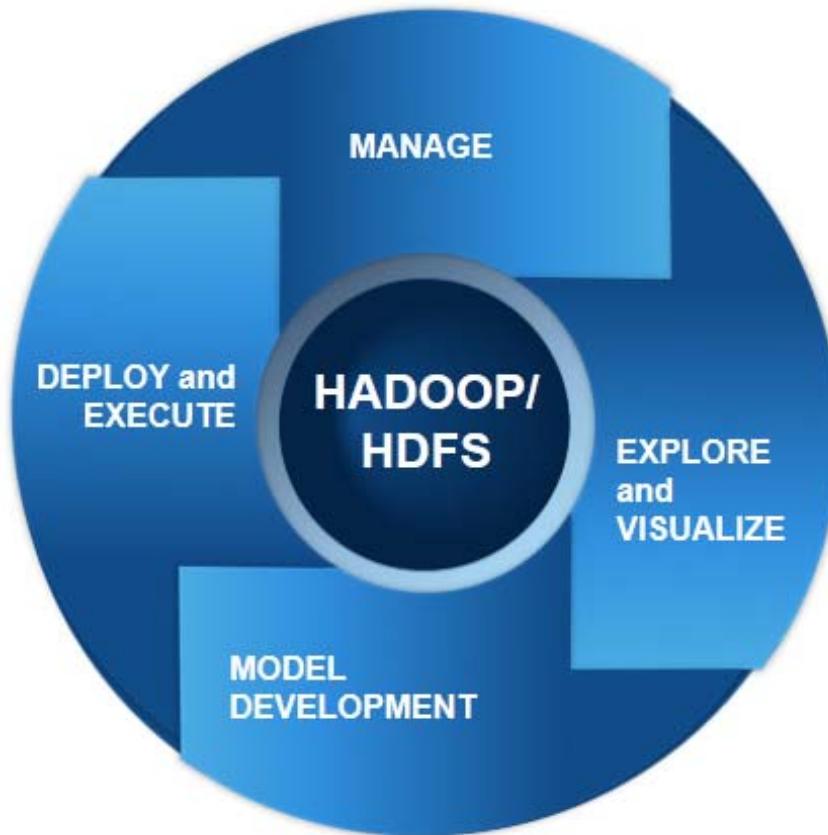
Column	Row	mraquisitionretailaverageprice	mraquisitionretailaverageprice	mraquisitionretailaverageprice	mraquisitionretailaverageprice
1	1	0.0003	0	0	0
2	2	0.0003	1.0000	0.0027	0.0162
3	3	0.0104	0.0027	1.0000	0.0002
4	4	0.0102	0.0102	0.0002	1.0000
5	5	0.0278	0.0236	0.0016	0.0450
6	6	0.0258	0.0251	0.0046	0.0057
7	7	0.0118	0.0007	0.0131	0.0032
8	8	0.0112	0.0112	0.0081	0.0088

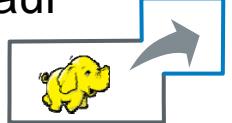
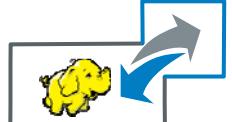
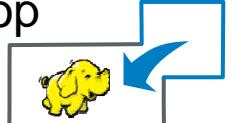
position	badbuy	Column	Min	Max	N	Sum	Mean	Std Dev.	Std Err.	Coefficients of Variation	Number Missing
ADESA	N	avgOdo	2514.14	97157	12441	152129149	14073	8032.88	4547499	31.9450	0
ADESA	Y	avgOdo	2661.14	33340	2163	25657473.8	13916	4266.17	217334	31.9561	0
MANNICH	N	avgOdo	705.07	30724	30239	347241300.4	13004	3443.37	23.2943	30.1050	0
MANNICH	Y	avgOdo	1870.13	47410	4718	25577270.7	13226	4320.08	62.8054	32.4528	0
OTHER	N	avgOdo	2118.89	48223	15433	25779023.4	14390	4058.38	48.7480	34.9059	0
OTHER	Y	avgOdo	1208.28	42628	2088	29779023.4	14390	4058.38	106.99	33.7003	0

- Linear Regression
- Multiple Regression
- Logistic Regression
- Analysis of Variance
- K-means Clustering
- Decision Trees
- Random Woods
- Generalized Linear Models
- Density-based Spatial Clustering (DBSCAN)
- Text pre-processing
- Text parsing
- Singular Value Decomposition
- Topic analysis
- Associations
- Implicit and explicit recommendations
- Group By processing







1. **Data Management:** SAS optimiert und erleichtert den Zugriff auf Daten in Hadoop 
2. **In-Memory Analytics:** SAS erweitert und beschleunigt Analytik auf Hadoop-Daten. 
3. **In-Database Processing:** SAS verlagert (analytische) SAS Funktionalität in das Hadoop Cluster. 

- SAS und HADOOP Informationen:

[http://www.sas.com/de\\_de/software/sas-hadoop.html](http://www.sas.com/de_de/software/sas-hadoop.html) - [http://www.sas.com/en\\_us/software/sas-hadoop.html](http://www.sas.com/en_us/software/sas-hadoop.html)

Code Beispiele:

<http://support.sas.com/resources/papers/proceedings14/SAS033-2014.pdf>

- Interessante White papers:

[http://www.sas.com/en\\_us/whitepapers/big-data-analytics-hadoop-107049.html](http://www.sas.com/en_us/whitepapers/big-data-analytics-hadoop-107049.html)

[http://www.sas.com/en\\_us/whitepapers/bringing-power-of-sas-to-hadoop-105776.html](http://www.sas.com/en_us/whitepapers/bringing-power-of-sas-to-hadoop-105776.html)

BARC: Big data analytics in der DACH region:

[http://www.sas.com/de\\_de/whitepapers/ba-wp-barc-big-data-analytics-2014-2298353.html](http://www.sas.com/de_de/whitepapers/ba-wp-barc-big-data-analytics-2014-2298353.html)

- Tom White: Hadoop : The Definitive Guide (O'Reilly)

<http://shop.oreilly.com/product/0636920021773.do>

- Edward Capriolo: Programming Hive (O'Reilly)

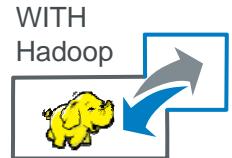
<http://shop.oreilly.com/product/0636920023555.do>

- Next Steps, wie Starten?

Test SAS&HADOOP - SAS Testlizenzen, Testumgebungen, Erfahrungsaustausch...

Kontakt : [gernot.engel@sas.com](mailto:gernot.engel@sas.com), [rainer.sternecker@sas.com](mailto:rainer.sternecker@sas.com)

## SAS VISUAL ANALYTICS



### Central Entry Point



#### DATA BUILDER

- Join data from multiple sources
- Create calculated and derived columns
- Load data



#### ADMINISTRATOR

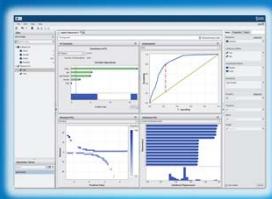
- Monitor SAS® LASR™ Analytic server
- Load/unload data
- Manage security



#### EXPLORER

- Perform ad-hoc analysis and data discovery
- Apply advanced analytics

### Integration



#### VISUAL STATISTICS

- Take VA analytics one step further
- Perform statistical modeling and classification



#### DESIGNER

- Create dashboard style reports for web or mobile



#### MOBILE BI

- Native iOS and Android applications that deliver interactive reports

**SAS® LASR™ ANALYTIC SERVER**

- Demo Visual Analytics / Visual Statistics



## SAS IN-DATABASE CAPABILITIES

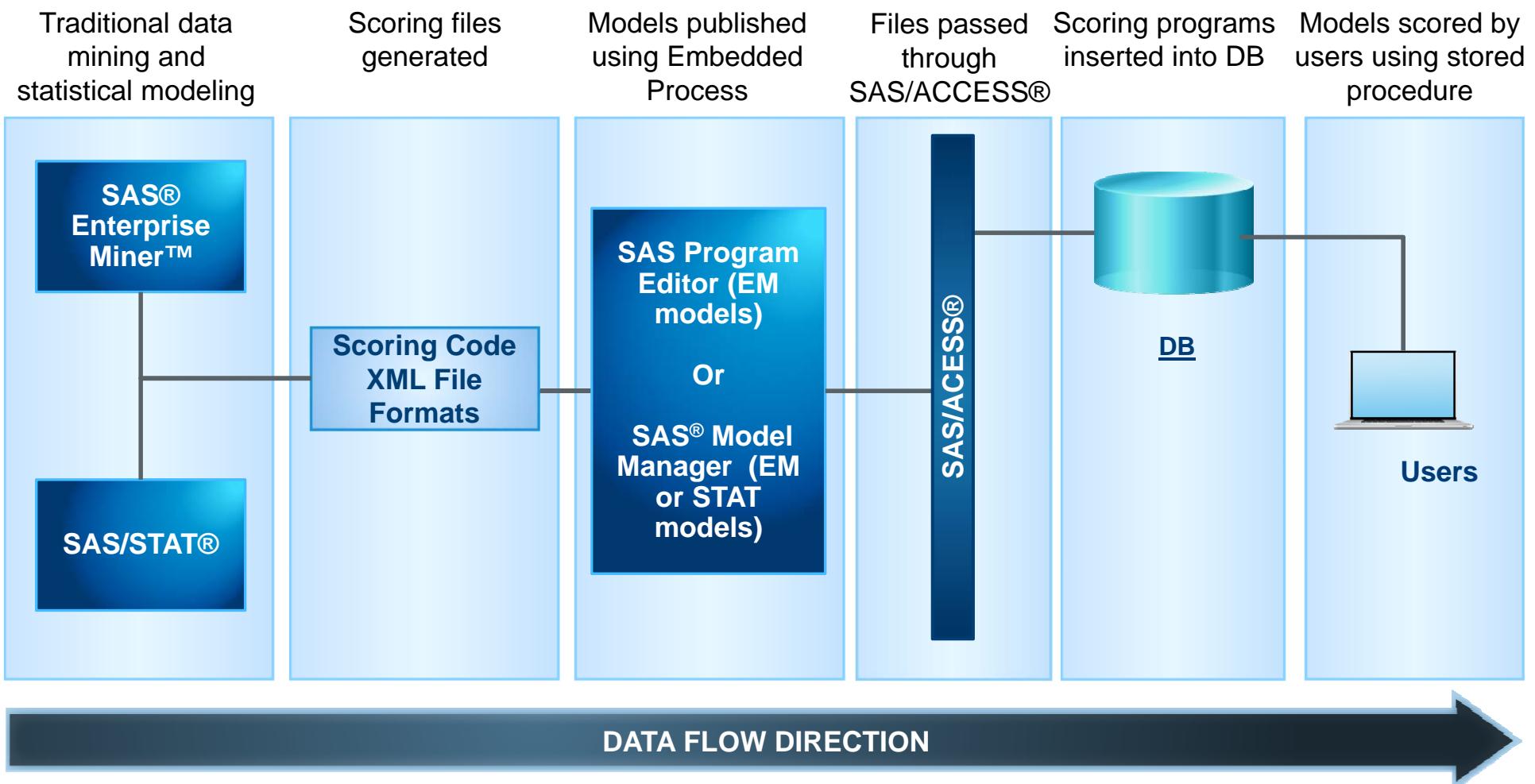
SQL Pushdown  
Reporting / OLAP  
ETL Integration  
Analytic Discovery  
Model Deployment & Management  
Data Manipulation/ Data Quality

MANAGEMENT

ADMINISTRATION

SECURITY

## IN-DATABASE SCORING ACCELERATOR



High-Performance Statistics	High-Performance Data Mining	High-Performance Text Mining	High-Performance Optimization	High-Performance Econometrics	High-Performance Forecasting
<ul style="list-style-type: none"><li>HPLOGISTIC</li><li>HPREG</li><li>HPLMIXED</li><li>HPNLMOD</li><li>HPSPLIT</li><li>HPGENSELECT</li></ul>	<ul style="list-style-type: none"><li>HPREDUCE</li><li>HPNEURAL</li><li>HPFOREST</li><li>HP4SCORE</li><li>HPDECIDE</li></ul>	<ul style="list-style-type: none"><li>HPTMINE</li><li>HPTMSCORE</li></ul>	<ul style="list-style-type: none"><li>OPTLSO</li><li>Select features in<ul style="list-style-type: none"><li>OPTMILP</li><li>OPTLP</li><li>OPTMODEL</li></ul></li></ul>	<ul style="list-style-type: none"><li>HPCOUNTREG</li><li>HPSEVERITY</li><li>HPQLIM</li></ul>	<ul style="list-style-type: none"><li>HPFORECAST</li></ul>

Common procedures (HPDS2, HPDMDB, HPSAMPLE, HPSUMMARY,  
HPIMPUTE, HPBIN, HPCORR)



### SAS<sup>®</sup> ESP

- ermöglicht und erleichtert die Verarbeitung von Massendaten (Sensoren, Telematik, Netzwerkverkehr, ...)
- kann Daten in Bewegung verarbeiten, ohne sie speichern zu müssen
- kann als Frontend für beliebige Systeme eingesetzt werden
- skaliert sehr gut (mehrere Mio. Events/Sek), bei extrem niedrigen Antwortzeiten (Millisekunden)
- arbeitet mit Datenflüssen (gerichteten Graphen) deren Knoten u.a. SQL-Operationen (wie aggregate, join, union, compute, filter, copy) und pattern matching repräsentieren

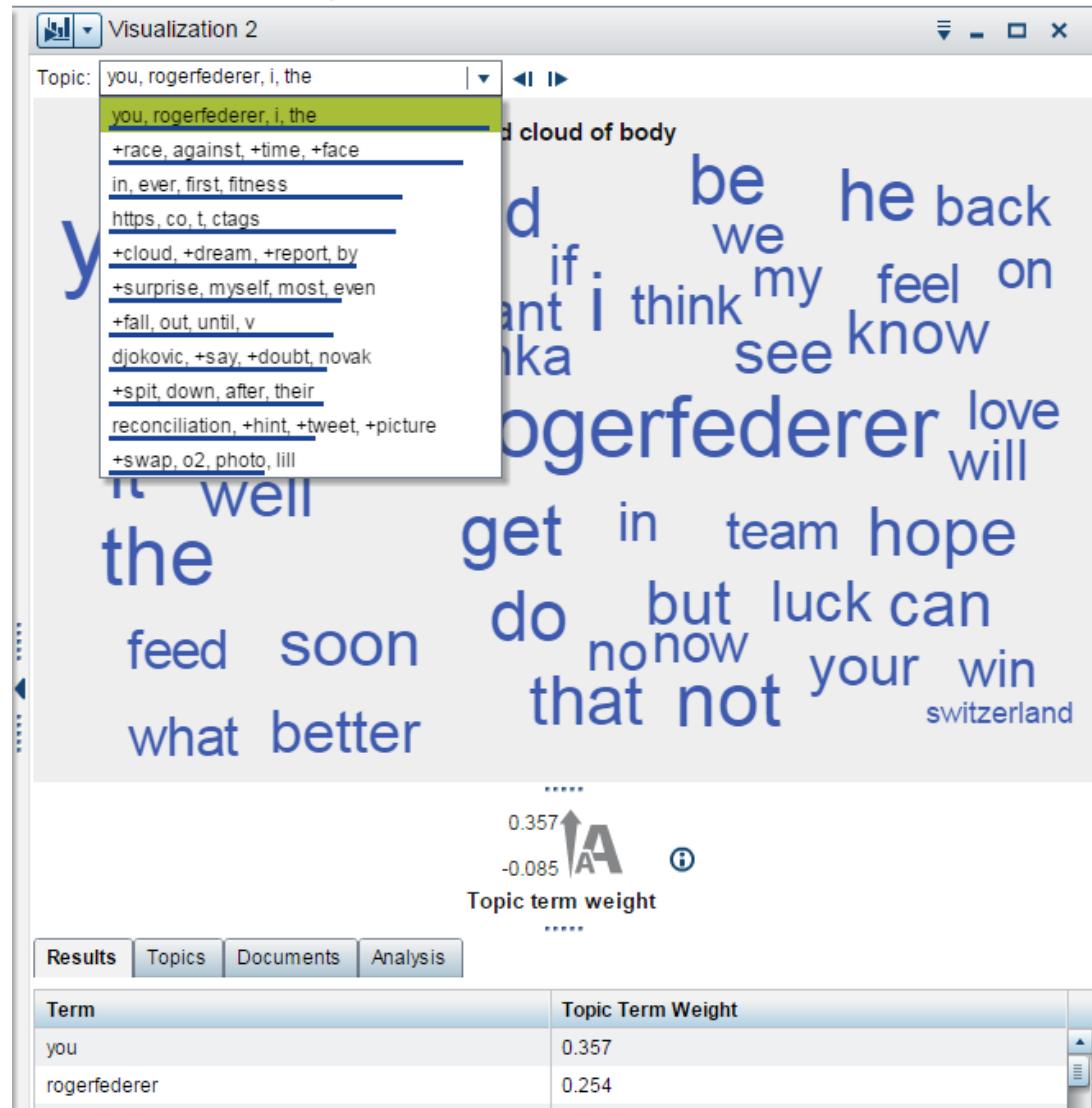


IN-ECHTZEIT

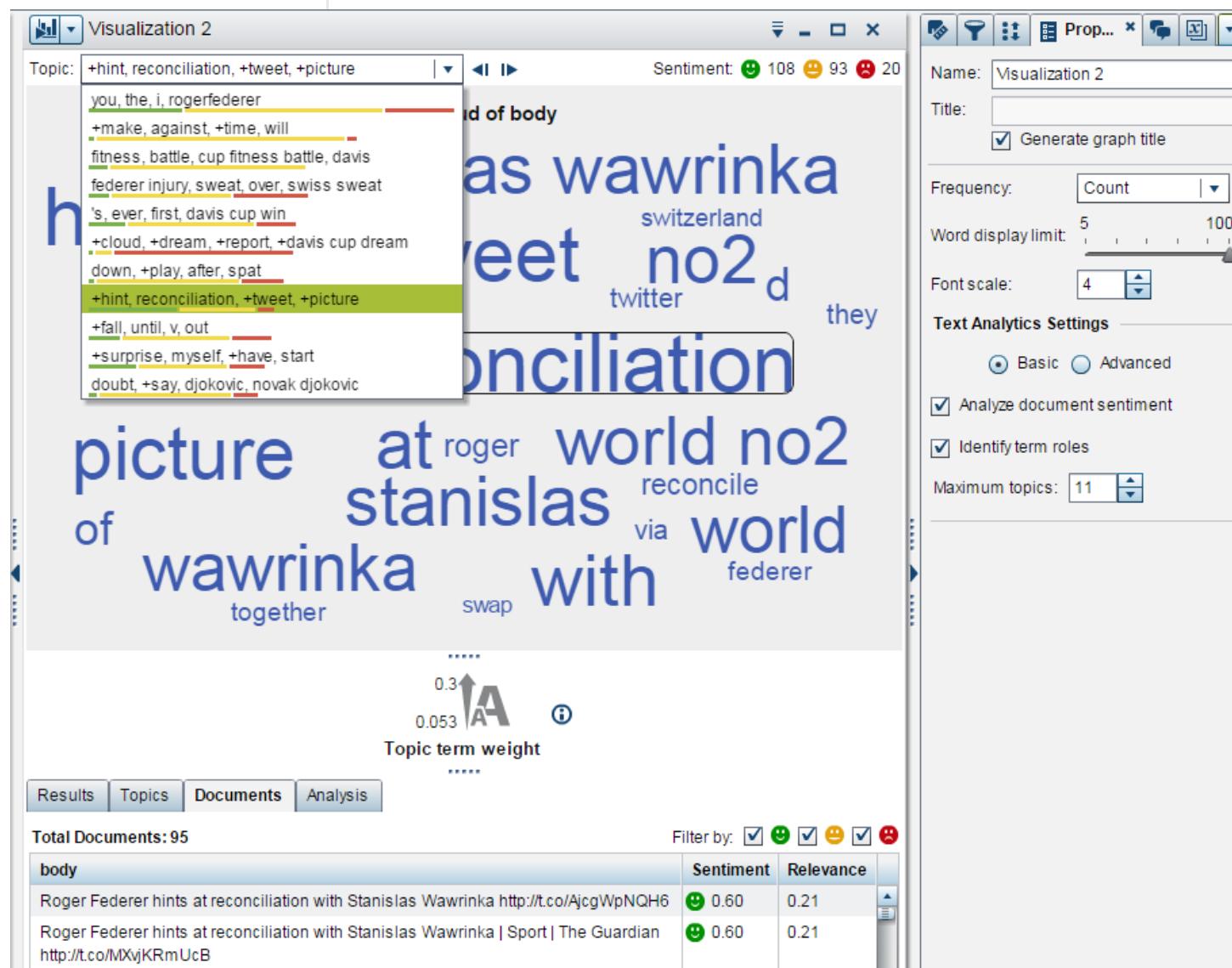
AUCH IN SAS VISUAL ANALYTICS

sas®club  
SAS INSTITUTE INC.

- Demo VA in Echtzeit



## WORD CLOUDS – SENTIMENT



Results	Topics	Documents	Analysis
Term	Topic Term Weight	Role	
twitter	0.083	Proper noun with ambiguous classi...	▲
via	0.081	Preposition	
they	0.078	Pronoun	
together	0.077	Adverb	
federer	0.076	Proper noun with ambiguous classi...	
switzerland	0.072	Geographical place	☰
swap	0.065	Verb	▼

