



interactive analyzer

understanding data

Proaktives Datenqualitätsmanagement

Basierend auf fortschrittlicher Data-Mining-Technologie findet der Synesis Interactive Analyzer völlig autonom und ohne manuelle Definition von Regeln, Filtern oder Kriterien selbst hoch komplexe Fehlermuster in Ihren Daten. Die Software liefert Korrekturvorschläge und unterstützt Sie bei der raschen Behebung der Inkonsistenzen. Einfache Bedienung, vielfältige Integrationsmöglichkeiten und das flexible Lizenzmodell machen die Software für Unternehmen jeder Größe attraktiv.

The screenshot displays the Interactive Analyzer software interface. A 'Data Consistency Monitor' dialog box is open, showing the file path 'C:\Users\dorneich\Documents\automobilB.txt' and settings for 'Group field' (WERKSTATTFALL_ID), 'Max. number of findings' (50), and 'Min. deviation strength' (50). The main window shows a data table with columns for 'WERKSTATTFALL_ID', 'ITEM', and 'ITEM'. Below the table, a list of association rules is displayed, such as 'Replace 'KOSTEN=1EUR' by '=KOSTEN=1-100EUR'. This is 275 times more probable.' and 'Replace 'BEFUND4=8386' by '=BEFUND4=6578'. This is 81 times more probable.' The interface also includes a 'Selected associations' section with radio buttons for 'superset', 'only group IDs', 'intersection', and 'entire records'.

Business Case

Die permanente Sicherstellung von Datenqualität und Datenkonsistenz geschäftsrelevanter Datensammlungen ist eine grundlegende Herausforderung für Manager und Fachkräfte aus IT, Qualitätssicherung, Logistik, Kundenmanagement und vielen anderen Bereichen. Experten schätzen die von fehlerhaften Daten verursachten volkswirtschaftlichen Gesamtkosten alleine in den USA auf mehr als 600 Milliarden US \$ jährlich (Eckerson, 2002): Datenfehler können fehlerhafte Produkte und Dienstleistungen hervorbringen, die Effektivität von Marketing- und Vertriebskampagnen verringern oder die Ursache für mangelhaftes Kundenmanagement und unzufriedene Kunden sein.

Aus diesem Grund wurden vielerorts unternehmensweite oder bereichsspezifische 'Data Governance' Funktionen geschaffen, deren Hauptaufgabe die permanente Einhaltung festgelegter Datenqualitäts-Standards ist. Dennoch besteht weithin eine gewaltige Diskrepanz zwischen den hohen Geldsummen, die für Hardware, Software, Administration und Datenerfassung ausgegeben werden, und der Aufmerksamkeit, die der Qualität der erfassten Daten gewidmet wird. In Summe führt das zu ineffizientem Mitteleinsatz und Unzufriedenheit mit den Systemen.

Datenqualitätsprobleme entstehen durch unvollständige, unrichtige oder inkonsistente Daten. Eine Software-basierte Lösung zum Datenqualitäts-Monitoring sollte daher einige oder alle der folgenden Funktionen beinhalten:

- **„Data Profiling“** – Datenüberprüfung mit dem Ziel, Probleme und Herausforderungen zu identifizieren.
- **Daten-Standardisierung** – ein **Regelsystem** für Datenaufnahme und –verarbeitung, das sicherstellt, dass die Daten definierte Qualitätsregeln einhalten.
- **„Matching or Linking“** – systematische Datenabgleiche, so dass ähnliche aber leicht unterschiedliche Datensätze zusammengeführt werden können (**Dublettenbereinigung**).
- **Monitoring** – systematische und regelmäßig ablaufende Datenqualitätsüberwachung und Generierung von Berichten oder (semi-)automatisch ablaufende **Korrekturmaßnahmen**

[Dieser Text basiert auf dem Wikipedia-Artikel http://en.wikipedia.org/wiki/Data_quality]

Datenqualitäts-Monitoring mit dem Interactive Analyzer

Traditionelle Methoden zur Aufdeckung und Behebung von Datenqualitätsproblemen betrachten einzelne Datenmerkmale auf fehlende oder aus dem üblichen Wertebereich herausfallende Werte. Ferner werden manuell definierte Filter und ‘Business Rules’ eingesetzt. All diese Methoden unterstützt der Interactive Analyzer auch. Als Alleinstellungsmerkmal werden darüber hinaus Data-Mining-basierte Methoden angeboten. Diese finden völlig selbständig, ohne dass manuelle Filter und Business Rules definiert werden müssen, subtile Inkonsistenzen zwischen den Werten verschiedener Merkmale eines Datensatzes. Zu jeder gefundenen mutmaßlichen Inkonsistenz findet die Software schließlich die – statistisch begründet – besten Korrekturvorschläge. Die Vorteile des Interactive Analyzer Ansatzes sind:

- **Automatische Fehlersuche ohne Voreinstellung von Fehlerkriterien:** Die Software findet eigenständig die Inkonsistenzkriterien, sie findet **alle statistisch** in der Datenbasis **signifikanten Kriterien** und sie findet **nur die statistisch signifikanten Kriterien**. Dieser ‘Hypothesen-freie’ Ansatz erspart nicht nur viel Arbeitszeit, er stellt auch sicher, dass kein wichtiges Kriterium bei der Definition der anzuwendenden Filter-Kriterien übersehen wird, dass aber auch nicht ‚nach Bauchgefühl‘ Inkonsistenzkriterien definiert werden, die sich auf der Datenbasis gar nicht rechtfertigen lassen.
- **Automatische Identifikation von ‚Top-Problemen‘:** Alle gefundenen mutmaßlichen Inkonsistenzen werden mit objektiven statistischen Signifikanzzahlen bewertet und quantifiziert, so dass sich geordnete Listen von ‘Top-Problemen’ bilden lassen. Dadurch wird die Begründung der Einstufung als Inkonsistenz und der vorgenommenen Korrektur reproduzierbar, auditierbar und jederzeit nachvollziehbar und kann automatisch protokolliert werden – zum Beispiel um regulatorische Vorgaben zu erfüllen.
- **Identifikation hochkomplexer Datenzusammenhänge:** Die eingesetzten Data-Mining-Verfahren können auch Muster und Zusammenhänge entdecken und verwenden, die wegen Ihrer Komplexität oder wegen der Größe oder Komplexität der Datenbasis mit traditionellen Methoden nicht entdeckt werden können.

Insgesamt ist der Data-Mining-Ansatz des Interactive Analyzer sehr gut komplementär zu bestehenden Methoden der Datenqualitätssicherung einsetzbar: er füllt die Lücken, die andere Methoden auf den Gebieten der Objektivität, Vollständigkeit, Analysegeschwindigkeit, Nachvollziehbarkeit, Dokumentation und Automatisierbarkeit haben.

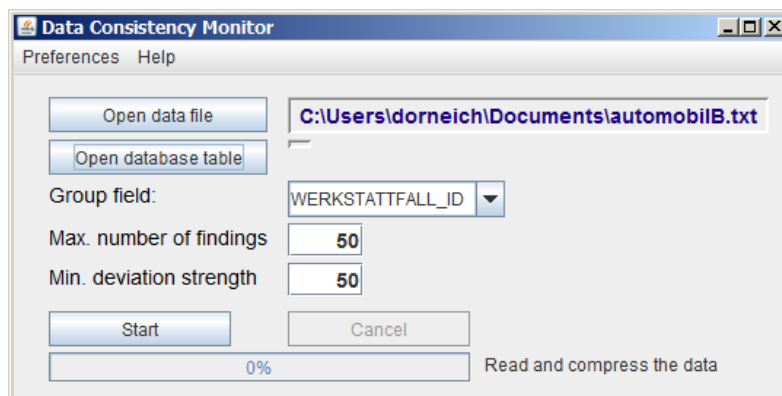
Alleinstellungsmerkmale des Interactive Analyzer

Im Vergleich zu anderen Methoden und Software-Werkzeugen aus dem Bereich ‘Data Quality Assurance’ bietet Interactive Analyzer die folgenden Differenzierungsmerkmale:

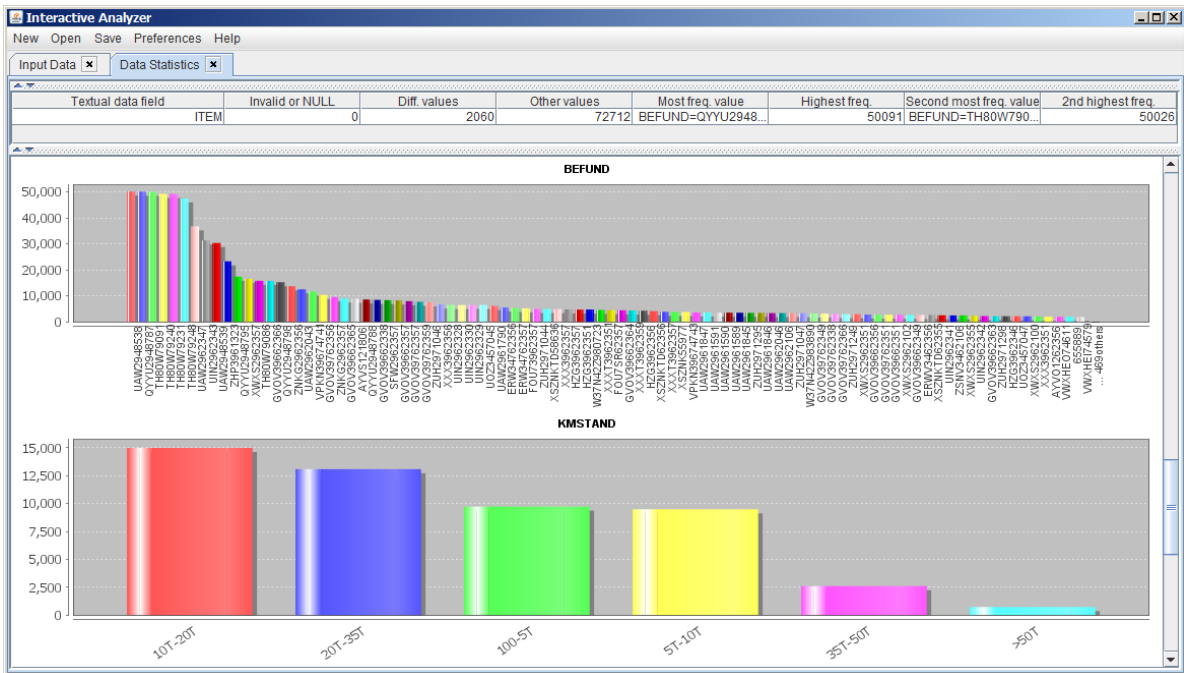
- Völlig **autonom**, ohne manuelle Definition von Regeln oder Kriterien, findet Interactive Analyzer auch **hoch komplexe Fehlerzusammenhänge**.
- **Leichte Bedienbarkeit**: Einfache Oberflächen, leichter Zugriff auf Datenbanktabellen, Datenbank-Sichten, Excel-Sheets oder Flachtextdateien. **Schneller Überblick über die Güte der Datenqualität**
- **Sehr günstige Total Cost of Ownership**: die geringen Hardware-, Administrations- und Einarbeitungsanforderungen und das flexible Preismodell machen den Einsatz des Interactive Analyzer für Unternehmen jeder Größe interessant.
- **Geschwindigkeit und Interaktivität**: die Datenkonsistenzanalyse dauert typischerweise nicht länger als einige Minuten.
- **Leistungsfähigkeit und Skalierbarkeit** durch seine komprimierte 'In-Memory'-Datenhaltung und seine leistungsfähigen, parallelisierten ablaufenden Analyse-Algorithmen.
- **Integration in bestehende IT-Strukturen und Prozesse**: Interactive Analyzer interagiert mit Datenbanken, Reporting-Systemen und Unternehmensapplikationen über standardisierte Schnittstellen wie JDBC, Web Services (SOA) und XML. Die Software ist nutzbar in Form von
 - automatisierten Prozessen und regelmäßig ablaufenden Batch-Job, Systemservices oder Datenbankprozessen.
 - selbständigen Anwendungsprogramme mit intuitiven grafischen Benutzeroberflächen
 - Java Applets, in Unternehmensportale und Reporting-Anwendungen integriert
 - Web Services im Rahmen von SOA-Architekturen
 - Funktionalen Erweiterungen ('stored procedures') in Oracle, IBM DB2 usw.

Anwendungsbeispiel: Garantie- und Kulanz-Daten eines Automobilherstellers

Im folgenden Anwendungsbeispiel wird eine Tabelle mit Hilfe des Interactive Analyzer™ Data Consistency Monitor auf Datenqualität, Datenkonsistenz und Auffälligkeiten untersucht, die Daten zu Werkstattaufenthalten von Fahrzeugen im Rahmen von Garantiefällen oder Kulanzreparaturen enthält. Zu jedem Werkstattfall enthält die Tabelle mehrere Datenzeilen mit Informationen über das Fahrzeug (Typschlüssel, Sonderausstattungen, km-Stand usw.), den internen Fehlerspeicher-Status zum Reparaturzeitpunkt sowie die abgerechneten Reparaturbefunde und –kosten. Die Daten sind OEM unabhängig und künstlich erzeugt, wurden aber realen Anwendungsprojekten in der Automobilindustrie nachempfunden.



1. In einem ersten Schritt verschaffen wir uns einen Überblick über die in den Daten enthaltenen Informationen und die statistischen Verteilungen der einzelnen Merkmalswerte. Die Daten erscheinen gut gepflegt: es sind weder größere Datenlücken noch fehlerhafte Merkmalswerte erkennbar.



2. Nun wollen wir mit Hilfe des Data-Mining-Ansatzes herausfinden, ob dennoch Inkonsistenzen in den Daten enthalten sind. Tatsächlich findet Interactive Analyzer Dutzende von Auffälligkeiten (siehe folgende Abbildung). Die ersten ca. 25 Auffälligkeiten wurden allesamt als hoch-signifikant eingestuft (der χ^2 -Konfidenztest liefert eine Verlässlichkeit von >99.9%). Eine Teilmenge dieser Auffälligkeiten betrifft Kombinationen von Fahrzeug-Baureihen (Typschlüsseln), Befund-Schlüsseln und Sonderausstattungen, die es in dieser Kombination eigentlich nicht geben dürfte. Diese Muster wurden durch Anklicken mit der Maus blau markiert.

3. Mit Hilfe der drei Knöpfe ‚Show Support‘, ‚Explore Support‘ und ‚Export Support‘ lassen sich die betroffenen Fahrzeuge und Reparaturvorgänge in den Originaldaten anschauen, in Charts visualisieren sowie in eine Textdatei, ein Excel-Sheet oder eine Datenbank exportieren, um die Daten gegebenenfalls zu korrigieren.

freq	item frequencies	lift	c	c	chi ² conf	item1	item2	item3
2	{771*;26177*}	0.002509009428827506	1.000	BEFUND=838694	KOSTEN=<1EUR	
2	{789*;26177*}	0.0024517696699949396	1.000	BEFUND=998748	KOSTEN=<1EUR	
2	{4949*;15370*}	6.657102182009259E-4	1.000	BEFUND=XSZNKTD58636	TYPSCHEUESSEL=UU58	
2	{2645*;15370*}	0.002491190827883843	1.000	BEFUND=XSZNKTD62355	TYPSCHEUESSEL=UU58	
2	{4105*;15370*}	0.0024077465553298777	1.000	BEFUND=XSZNKTD62356	TYPSCHEUESSEL=UU58	
2	{1201*;15370*}	0.002743213879996988	1.000	BEFUND=XSZNKTD62357	TYPSCHEUESSEL=UU58	
2	{2440*;10473*}	0.0019815982706506818	1.000	BEFUND=ZSNV3462106	TYPSCHEUESSEL=UV56	
2	{784*;26177*}	0.002467405956155622	1.000	BEFUND4=8386	KOSTEN=<1EUR	
2	{1079*;15370*}	0.0030533826412200023	1.000	FEHLERSPEICHER=GVOV397...	TYPSCHEUESSEL=UU58	
2	{2028*;15370*}	0.0016245561488542318	1.000	FEHLERSPEICHER=XSZNKTD...	TYPSCHEUESSEL=UU58	
2	{6522*;15370*}	5.051517739767528E-4	1.000	FEHLERSPEICHER=XZH5873	TYPSCHEUESSEL=UU58	
2	{1412*;15370*}	0.002333286026824633	1.000	SONDERAUSSTATT=112	TYPSCHEUESSEL=UU58	
2	{1388*;15370*}	0.0023736310301703044	1.000	SONDERAUSSTATT=177	TYPSCHEUESSEL=UU58	
2	{2998*;15370*}	0.0010989325783443572	1.000	SONDERAUSSTATT=1HC	TYPSCHEUESSEL=UU58	
2	{7344*;10473*}	0.0013167483062057907	1.000	SONDERAUSSTATT=52W	TYPSCHEUESSEL=UV56	
2	{7344*;8359*}	0.0016497553548143612	1.000	SONDERAUSSTATT=52W	TYPSCHEUESSEL=UV16	
2	{7383*;8359*}	0.001641040678011197	1.000	SONDERAUSSTATT=679	TYPSCHEUESSEL=UV16	
2	{5305*;8359*}	0.0011419230278752752	1.000	SONDERAUSSTATT=70B	TYPSCHEUESSEL=UV16	
2	{10203*;10473*}	0.002843340064914827	1.000	SONDERAUSSTATT=70V	TYPSCHEUESSEL=UV56	
2	{1708*;23341*}	0.001270191550645069	1.000	TYPSCHEUESSEL=UU98	SONDERAUSSTATT=176	
2	{8359*;8416*}	7.198077070910569E-4	1.000	TYPSCHEUESSEL=UV16	BEFUND=GVOV39662338	
2	{8359*;23747*}	0.0010204070683249816	1.000	TYPSCHEUESSEL=UV16	SONDERAUSSTATT=301	
2	{8359*;10327*}	0.0017598242460186893	1.000	TYPSCHEUESSEL=UV16	SONDERAUSSTATT=7IC	
2	{10473*;23747*}	2.036088676627643E-4	1.000	TYPSCHEUESSEL=UV56	SONDERAUSSTATT=301	
2	{10473*;11687*}	4.1371607601503065E-4	1.000	TYPSCHEUESSEL=UV56	BEFUND=VPKN39674741	
2	{864*;19936*}	0.002939847326258843	1.000	TYPSCHEUESSEL=UV57	SONDERAUSSTATT=680	
2	{3015*;6466*}	0.0025974878673956743	1.000	TYPSCHEUESSEL=UV58	BEFUND=UIN2962330	
2	{3015*;6466*}	0.0025974878673956743	1.000	TYPSCHEUESSEL=UV58	BEFUND=UIN2962329	
3	{1618*;23341*;26177*}	0.0025937925587822953	0.986	BEFUND=996535	SONDERAUSSTATT=176	KOSTEN=<1EUR
3	{1618*;20595*;26177*}	0.0029396315666199344	0.986	BEFUND=996535	SONDERAUSSTATT=1HO	KOSTEN=<1EUR
3	{1618*;19945*;26177*}	0.0030354330466050407	0.986	BEFUND=996535	SONDERAUSSTATT=5FY	KOSTEN=<1EUR
3	{1618*;19936*;26177*}	0.003036803376531779	0.986	BEFUND=996535	SONDERAUSSTATT=680	KOSTEN=<1EUR
3	{1618*;19719*;26177*}	0.003070222228030709	0.986	BEFUND=996535	SONDERAUSSTATT=535	KOSTEN=<1EUR

data: automobiliB_20090703.adf Selected associations: 25

superset only group IDs Show support Explore support Export support
 intersection entire records

Diese Korrektur kann manuell, das heißt durch einen menschlichen Fachexperten, erfolgen. Interactive Analyzer kann aber auch eigenständig, aufbauend auf einer statistischen Assoziationsanalyse der Datenbasis, Korrektorempfehlungen vorschlagen. In der folgenden Abbildung werden Korrektorempfehlungen für die blau markierte Auffälligkeit angezeigt. Offensichtlich wurde in einem Reparaturfall ein Reparaturbefund („8386“), der normalerweise mit Kosten von knapp unter 100 € (275 mal) oder knapp über 100 € (481 mal) verbucht wird, fälschlicherweise als kostenfrei verbucht (Kosten<1€).

The screenshot shows the 'Interactive Analyzer' software interface. The main window displays a table with columns for 'freq', 'item frequencies', 'lift', 'c.c.', 'chi²', 'conf', 'item1', 'item2', and 'item3'. A 'Correction Hints' dialog box is open, displaying the following information:

- Replace 'KOSTEN=1EUR' by '=KOSTEN=101-500EUR'. This is 481 times more probable.
- Replace 'KOSTEN=1EUR' by '=KOSTEN=1-100EUR'. This is 275 times more probable.
- Replace 'BEFUND4=8386' by '=BEFUND4=6588'. This is 123 times more probable.
- Replace 'BEFUND4=8386' by '=BEFUND4=6578'. This is 81 times more probable.
- Replace 'BEFUND4=8386' by '=BEFUND4=3878'. This is 64 times more probable.

The dialog box also includes an 'OK' button and a close button (X). The main window shows the 'Selected associations' section with '1' selected, and options for 'superset', 'intersection', 'only group IDs', and 'entire records'.

Zusammenfassung

Interactive Analyzer findet Auffälligkeiten, Datenfehler und Inkonsistenzen, die sich erst in der Zusammenschau mehrerer Datenmerkmale als Auffälligkeit herausstellen und deshalb von den gängigen Datenqualitäts-Überwachungsmethoden nicht gefunden werden. Darüber hinaus liefert Interactive Analyzer durch quantitative Analysen begründete Korrekturhinweise und erlaubt so eine effektive, automatisierte oder zumindest Software-unterstützte Behebung von Datenfehlern. Und die ersten aussagefähigen Ergebnisse erhalten sie meist schon am ersten Einsatztag, ohne komplexe Hardwareanforderungen und überbordende Lizenzkosten.

Kontakt

Dr. Ansgar Dorneich
 Fichtenstraße 7
 71088 Holzgerlingen
 Email: info@i-analyzer.com
 Web: www.i-analyzer.com