

Steven Helmis | Robert Hollmann

Webbasierte Datenintegration

VIEWEG+TEUBNER RESEARCH

Ausgezeichnete Arbeiten zur Informationsqualität

Herausgeber:

Dr. Marcus Gebauer

Bewertungskommission des
Information Quality Best Master Degree Award 2008:

Prof. Dr. Holger Hinrichs, FH Lübeck (Kommissionsvorsitz)

Dr. Marcus Gebauer, WestLB AG und Vorsitzender der DGIQ

Prof. Dr. Knut Hildebrand, HS Darmstadt

Bernhard Kurpicz, OrgaTech GmbH

Prof. Dr. Jens Lüssem, FH Kiel

Michael Mielke, Deutsche Bahn AG und Präsident der DGIQ

Prof. Dr. Felix Naumann, HPI, Uni Potsdam

Prof. Dr. Ines Rossak, FH Erfurt

Die Deutsche Gesellschaft für Informations- und Datenqualität e.V. (DGIQ) fördert und unterstützt alle Aktivitäten zur Verbesserung der Informationsqualität in Gesellschaft, Wirtschaft, Wissenschaft und Verwaltung. Zu diesem Zweck befasst sie sich mit den Voraussetzungen und Folgen der Daten- und Informationsqualität. Sie fördert zudem durch Innovation und Ausbildung die Wettbewerbsfähigkeit der Unternehmen sowie die des unternehmerischen und akademischen Nachwuchses in Deutschland.

Die vorliegende Schriftenreihe präsentiert ausgezeichnete studentische Abschlussarbeiten in der Daten- und Informationsqualität. Veröffentlicht werden hierin die Siegerarbeiten des jährlich stattfindenden „Information Quality Best Master Degree Award“.

Steven Helmis | Robert Hollmann

Webbasierte Datenintegration

Ansätze zur Messung und Sicherung
der Informationsqualität in heterogenen
Datenbeständen unter Verwendung
eines vollständig webbasierten Werkzeuges

Mit einem Geleitwort von Dr. Marcus Gebauer

VIEWEG+TEUBNER RESEARCH

Bibliografische Information der Deutschen Nationalbibliothek
Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der
Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über
<<http://dnb.d-nb.de>> abrufbar.



Gedruckt mit freundlicher Unterstützung
der Information Quality Management Group

1. Auflage 2009

Alle Rechte vorbehalten

© Vieweg+Teubner | GWV Fachverlage GmbH, Wiesbaden 2009

Lektorat: Christel A. Roß

Vieweg+Teubner ist Teil der Fachverlagsgruppe Springer Science+Business Media.
www.viewegteubner.de



Das Werk einschließlich aller seiner Teile ist urheberrechtlich geschützt.
Jede Verwertung außerhalb der engen Grenzen des Urheberrechtsgesetzes
ist ohne Zustimmung des Verlags unzulässig und strafbar. Das gilt
insbesondere für Vervielfältigungen, Übersetzungen, Mikroverfilmungen
und die Einspeicherung und Verarbeitung in elektronischen Systemen.

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in
diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme,
dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung
als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Umschlaggestaltung: KünkelLopka Medienentwicklung, Heidelberg

Druck und buchbinderische Verarbeitung: STRAUSS GMBH, Mörlenbach

Gedruckt auf säurefreiem und chlorfrei gebleichtem Papier.

Printed in Germany

ISBN 978-3-8348-0723-6

Geleitwort

Als Vorsitzender der Deutschen Gesellschaft für Informations- und Datenqualität (DGIQ e.V.) bin ich glücklich darüber, dass Sie dieses Buch in Ihren Händen halten. Das vorliegende Buch ist Ausdruck unseres Bestrebens, dem wissenschaftlichen Nachwuchs die Möglichkeit zu eröffnen, ihre Arbeiten einem breiten Publikum darstellen zu können.

Dass Sie gerade diese Arbeit vorfinden, ist Ergebnis eines strengen Auswahlprozesses, den die DGIQ mit dem zum ersten Mal ausgeschriebenen „Information Quality Best Master Degree Award 2008“ durchgeführt hat. Studenten waren aufgefordert, ihre Abschlussarbeiten zum Thema Informationsqualität in diesem Wettbewerb durch ihre begutachtenden Professoren einreichen zu lassen. Vertreter aus Wissenschaft, Forschung und Industrie haben diese akademischen Abschlussarbeiten begutachtet.

Die vorliegende „Doppel“-Arbeit von Robert Hollmann und Steven Helmig zeichnet sich insbesondere durch das geschlossene Bild ihrer Forschung aus. Neben der Betrachtung der Datenqualität in der heutigen typischen heterogenen Systemlandschaft, bei gleichzeitig komplizierter werdenden Systemarchitekturen, stehen auch die semantischen Herausforderungen für die Datenqualität im Fokus. Die Darstellung, wie Daten zu konsolidieren, zu bereinigen, und der Kunde dabei auch noch konstruktiv mit Hilfe eines Werkzeuges in diesen Prozess einzubinden ist, ist den Kandidaten herausragend gelungen.

Besonders freue ich mich, dass wir mit dem Verlag Vieweg+Teubner nun die Siegerarbeiten in einer Schriftenreihe jährlich veröffentlichen können. Für die Initiative des Verlages möchte ich mich recht herzlich bedanken.

Offenbach, 27. August 2008

Dr. Marcus Gebauer

Vorwort

Jeder Prozess ist mit der Erzeugung von Daten verbunden. So generieren unzählige heterogene IT-Systeme in immer globaler werdenden Wirtschaftsunternehmen Tag für Tag Millionen von Datensätze, die in richtigen Zusammenhang gebracht, wertvolle Informationen für das Unternehmen und die Entscheidungsfindung in Selbigem von existentieller Bedeutung sein können. Ganzheitliche Sichten auf die Gesamtinformationen sind von großem Wert für Entscheidungsprozesse. Hinzu kommt die überwältigende Informationsflut diverser Internetquellen, die eine Bereicherung für einen Datenbestand bedeuten kann, aber auch Probleme erzeugt und so die Daten(-qualität) in einem Bestand nachhaltig schädigen kann. Diese Heraus- und Anforderungen an Informationssysteme sind nur unter Beachtung von ausreichender Informationsqualität zu erreichen und zu bewältigen. Bei der Integration der heterogenen Datenquellen für eine ganzheitliche Sicht spielt die Qualität der Daten vor- und nach der Integration eine bedeutende Rolle. Datenqualitätsmängel, wie z.B. Duplikate, mangelnde Reputation oder Überalterung bzw. Inkonsistenz können den Informationsgewinn eines integrierten Systems empfindlich stören. Das führt zu falschen bzw. nicht zielgerichteten Entscheidungen deren Grundlage diese „schlechten“ Daten bildeten. Jedes Informationssystem kämpft mit solchen Problemen, die schon im jeweiligen System selbst entstehen, die ihre Manifestation oft jedoch bei der Integration zeigen. Nach und nach werden von Softwareherstellern und Unternehmen diese Probleme und die Vorteile von hoher Datenqualität erkannt. Jedoch besteht auf diesem Gebiet enormer Handlungsbedarf bei Benutzern wie auch Systemherstellern. Vorreiter hier sind einige Unternehmen und Experten, die in einem Zusammenschluss der Deutschen Gesellschaft für Informations- und Datenqualität (DGIQ) diese Probleme offensiv thematisieren, Lösungen vorstellen und für die Problematik sensibilisieren.

Auch die Autoren der vorliegenden Arbeiten konnten so, an einer immer lebendiger agierenden „IQ-Society“ partizipieren. Im Rahmen ihres Studiums und ihren Masterarbeiten befassten sich die Autoren eingehend mit dem Thema Datenqualität und der Analyse und Beseitigung von Datenqualitätsmängeln. Im Ergebnis konnten sie einen Überblick über IST-Stand, Vorgehensweisen und Chancen der Datenqualitässicherung erarbeiten. Gemeinsam entstanden Ansätze zur Umsetzung sowie die prototypische Implementierung eines Datenqualitätswerkzeu-

ges, das ausschließlich auf aktuellen Webtechnologien aufbaut. Mit diesen Ideen konnte der erste Platz des „DGIQ Best Master Degree Award“, der im Jahre 2007 ausgeschrieben wurde, erreicht werden. Mit dieser Veröffentlichung wollen die Autoren zur weiteren Sensibilisierung für gute bzw. schlechte Informationsqualität beitragen und konkrete Lösungen für die Sicherung einer solchen Qualität aufzeigen. Denn nur valide Information schaffen einen Vorteil, den jedes Unternehmen auf einem globalen Markt für sich in Anspruch nehmen möchte. Wir wollen Sie mit unseren Arbeiten dazu motivieren sich aktiv an der Diskussion zum Thema zu beteiligen. Vorteile für sich und Ihr Unternehmen zu erkennen und vielleicht auch Teil der „IQ-Society“ zu werden.

Die Autoren

Steven Helmis, Robert Hollmann

Inhaltsverzeichnis

Abbildungsverzeichnis	XIII
Tabellenverzeichnis	XVII
Abkürzungsverzeichnis	XIX

I Datenbereinigung und Konsolidierung von heterogenen Datenbeständen

– Steven Helmis – 1

1 Einleitung	3
1.1 Motivation	4
1.2 Zielsetzung der Arbeit	5
1.3 Aufbau der Arbeit	5
2 Datenqualität	7
2.1 Datenqualität definieren	7
2.2 Datenfehler	8
2.3 Qualitätskriterien	11
2.4 Methoden zur Einstufung der Qualität	14
3 Dimensionen und Architektur der Informationsintegration	25
3.1 Verteilung	25
3.2 Heterogenität	26
3.3 Autonomie	28
3.4 Integrationsarchitektur	29
4 Data Cleaning	35
4.1 Datenanalyse	36
4.2 Normalisierung und Validierung	39
4.3 Record Matching	40

4.4	Record Merging	42
5	Konzeption des Data Cleaning Toolkits	49
5.1	Bewertung und Analyse existierender Systeme	49
5.2	Anforderungsanalyse	52
5.3	Architektur Data Cleaning Toolkit	54
5.4	Funktionsumfang	55
6	Implementierung	63
6.1	Datenbankentwicklung	63
6.2	Webentwicklung	71
6.3	Probleme während der Implementierungsphase	77
7	Zusammenfassung und Ausblick	79
	Literaturverzeichnis	81
II	Auffinden und Bereinigen von Duplikaten in heterogenen Datenbeständen	
	– Robert Hollmann –	89
8	Einleitung	91
8.1	Motivation	92
8.2	Zielstellungen dieser Arbeit	93
8.3	Gliederung dieser Arbeit	94
9	Informationen, Daten und Wissen- ein Definitionsversuch	95
9.1	Begriffsdefinitionen	96
9.2	Herkunft von Daten und Informationen	98
9.3	Beschaffenheit von Daten und Zugriff auf Informationen	98
10	Informationsintegration im Fokus der Datenqualität	103
10.1	Ist-Stand in Unternehmen- Notwendigkeit der Integration	103
10.2	Informations- und Datenqualität	105
10.3	Sicherung der Datenqualität	114
10.4	Kosten der Datenqualität	115
11	Duplikate in Datenbeständen	117
11.1	Dubletten und deren Identifikation	117

11.2 Ein Framework zur Objektidentifikation	118
11.3 Das Dilemma der Dublettensuche	120
12 Konkrete Verfahren zur Dublettenuauffindung und Klassifikation	125
12.1 Ähnlichkeitsmessungen und Klassifikation	125
12.2 Ähnlichkeitsbestimmung bei Tupeln in einem Datenbestand . . .	126
12.3 Vorselektion für die Dublettensuche	142
13 Konzept der Datenqualitätsanwendung „DCT“	147
13.1 Zielstellung der Applikation	147
13.2 Anforderungsanalyse	148
13.3 Technologiemoell	157
13.4 Datenbankmodell	160
13.5 Applikationsarchitektur	164
13.6 Applikationsstruktur	166
13.7 Entwicklung einer Benutzeroberfläche	169
14 Implementierung, ausgewählte Algorithmen- und Datenstrukturen	173
14.1 „DCT“- Der Verbindungsmanager	173
14.2 „DCT“- Der Workspace-Table Manager	176
14.3 „DCT- Data Profiling“	177
14.4 „DCT“-Plausibilitätskontrolle	180
14.5 „DCT“- Auffinden von Duplikaten	182
15 Fazit und Ausblick	187
Literaturverzeichnis	189
16 Anhang	195

Abbildungsverzeichnis

2.1	Klassifikation von Daten-Qualitäts-Problemen	8
2.2	Konzeptionelles Gerüst der Datenqualität	12
2.3	Qualitäts-Dimensionen	13
2.4	Allgemeine Hierarchie	22
3.1	Orthogonale Dimensionen der Informationsintegration	30
3.2	Mediator-Wrapper-Architekturen	33
5.1	Drei-Schichten-Architektur des DCT	54
5.2	Modulabschnitte des DCT	56
5.3	Funktionsübersicht im Detail	56
5.4	Informationen zum Laden der Daten	58
5.5	Qualitätsmerkmale der WST	59
5.6	Spaltenzuordnung für den Vergleich mit Referenz	60
5.7	Standardisierung von Attributen	61
5.8	Ergebnis eines Vergleichs mit Referenzdaten	61
6.1	ER-Modell Metadaten	64
6.2	Übersicht der implementierten Prozeduren und Funktionen	65
8.1	Allgemeine Architektur einer Integrationslösung	92
9.1	Daten, Information und Wissen	95
9.2	Semiotisches Dreieck	98
9.3	Übersicht über mögliche Datenquellen	99
9.4	Einteilung der Datenbeschaffenheit	99
9.5	Schlüsselbeziehung zweier relationaler Tabellen	101
10.1	Heterogene IT-Landschaft als weitverbreiteter IST-Stand in Unternehmen	105
10.2	Datenqualität in Analogie zur industriellen Fertigung	106
10.3	Qualitätsdimensionen	108
10.4	Bewertung der Qualität von Daten aus verschiedene Sichten	110

10.5	Datenqualitätsprobleme im Kontext der Integration	111
10.6	Zyklus des TDQM	115
11.1	Generisches Modell zur Identifizierung von Objekten	119
11.2	Konflikt zwischen den Zielen der Dublettensuche	121
11.3	Zusammenhang zwischen relevanten und gefundenen Datensätzen	122
12.1	Übersicht über die Duplikaterkennung abgeändert	126
12.2	Vektorraummodell für die Ähnlichkeitsbestimmung	133
12.3	Dublettenidentifizierung mit Hilfe externer Daten	135
12.4	Aufbau einer Hashspeicherstruktur	137
12.5	Mögliche Klassifikation von Clusterverfahren	138
12.6	Gegenüberstellung von hierarchischen und partitionierenden Clusterverfahren	139
12.7	Hierarchische Clustering Verfahren in der Übersicht	140
12.8	DBSCAN Algorithmus zur dichte-basierten Erzeugung von Clustern	142
12.9	Ablauf der Sorted Neighborhood Methode	144
13.1	DCT-UseCases in UML-Notation	151
13.2	Funktionsübersicht „DCT“	152
13.3	Übersicht über die anfallenden Daten des „DCT“	156
13.4	Client-Server-Architektur von Webanwendungen	158
13.5	Technologiemodell des „DCT“	160
13.6	Entity-Relationship-Modell des „DCT“	163
13.7	Architektur der Anwendung „DCT“	166
13.8	Struktur einer „MVC“ Webanwendung	167
13.9	MVC-Struktur des „DCT“	167
13.10	Klassendiagramm der Anwendung „DCT“	169
13.11	Screendesign des „DCT“	170
13.12	Frei positionierbare Fenster des „DCT“	171
13.13	Tooltip zur visuellen Unterstützung der Bedienung	171
14.1	Verbindungsmanager Übersicht des „DCT“	174
14.2	„DCT- Verbindung bearbeiten“	174
14.3	Auswahl zu importierender Tabellenattribute	175
14.4	Benennung der Zielattribute und Datentypendefinition	176
14.5	Zusammenfassung vor dem Laden in den Workspace	176
14.6	Zugriff auf die DQ-Funktionen des „DCT“ via WST-Manager	177
14.7	Datenqualitätsübersicht des „DCT“	178

14.8	Musterermittlung im „DCT“	180
14.9	Anzeigen aller Datensätze mit nicht belegten Werten in einer definierten Spalte	180
14.10	Mapping zwischen Referenz- und Workspacetabelle	181
14.11	Gefundene Inkonsistenzen beim Referenzdatenvergleich	182
14.12	Auswahl der Attribute für den Test auf Dubletten	183
14.13	Anzeige aller potentiellen Duplikate im „DCT“	184
16.1	DCT Klassendiagramm	197

Tabellenverzeichnis

1.1	Statistische Kennzahlen	4
2.1	Bsp. single-source Problem auf Daten-Ebene	9
2.2	Bsp. single-source Problem auf Schema-Ebene	10
2.3	Bsp. multi-source Probleme auf Schema- und Datenebene	11
2.4	Werte zweier Suchmaschinen	16
2.5	Bsp. Simple Additive Weighting Methode	19
2.6	Bsp. Simple Additive Weighting Methode mit Idealen	20
2.7	Bsp. Rangfolge mit TOPSIS	21
2.8	AHP Vergleichsskala	23
4.1	Ergebnis einer Dublettensuche bei Adressen	43
4.2	Beispieltabelle 1 (t1)	44
4.3	Beispieltabelle 2 (t2)	44
4.4	Ergebnis INNER JOIN	45
4.5	Ergebnis aus Kombination OUTER JOIN mit UNION	46
4.6	Resultset - INTERSECT	47
4.7	Resultset - EXCEPT	47
6.1	Transformation von Straßennamen	67
6.2	Transformation von Telefonnummern	69
9.1	Beispiel für relationale Datenspeicherung- Tabelle Kunden	100
9.2	Beispiel für relationale Datenspeicherung- Tabelle Kontakt	100
10.1	Kostenarten bei der DQ- Sicherung	116
11.1	Mögliche Ergebnisse einer Dublettenerkennung	121
12.1	Mögliche Operationen für die Berechnung der Levenshtein-Distanz	127
12.2	Levenshteindistanz zwischen „Maik“ und „Mike“	128
12.3	Übersicht über die Soundex kodierten Laute	130
12.4	Laute nach der „Kölner Phonetik“ kodiert	131

12.5	Beispiel für mögliche Tokenbewertung	132
13.1	Datenbanktabelle „meta_connection“	161
13.2	Datenbanktabelle „meta_ws_table“	162
13.3	Datenbanktabelle „meta_mapping“	162
13.4	Datenbanktabelle „Entwickelte Prozeduren und Funktionen“	165
16.1	Kölner Phonetik Ersetzungsregeln	195

Abkürzungsverzeichnis

AHP	Analytic Hierarchy Process
Ajax	Asynchronous JavaScript and XML
ASCII	American Standard Code for Information Interchange
ASP	Application Service Provider
BI	Business Intelligence
BNF	Backus-Naur-Form
CRM	Customer Relationship Management
CSS	Cascading Style Sheets
CSV	Comma Separated Values
CWM	Common Warehouse Metamodel
DB	Datenbank
DBMS	Datenbank Management System
DBS	Datenbank System
DCT	Data Cleaning Toolkit
DD	Data Dictionary
DDL	Data Definition Language
DEA	Data Envelopment Analysis
DGIQ	Deutsche Gesellschaft für Informations- und Datenqualität
DIN	Deutsches Institut für Normung
DML	Data Manipulation Language
DQ	Datenqualität
ER	Entity-Relationship
ERP	Enterprise Resource Planning
ETL	Extraktion, Transformation, Laden

FDH	Free Disposable Hull
GD	gefundene Datensätze die als Duplikat markiert wurden
HTML	Hypertext Markup Language
HTTP	Hypertext Transfer Protocol
HTTPS	Hypertext Transfer Protocol Secure
IDs	Identifikationen
IQC	Information Quality Criterion
IQS	Information Quality Scores
IR	Information Retrieval
IT	Informationstechnologie
JMS	Java Message Service
KMU	kleine und mittelständige Unternehmen
LDAP	Lightweight Directory Access Protocol
MADM	Multi-Attribute Decision-Making
MDBS	Multidatenbanksystem
MVC	Model-View-Controller
OCRA	Operational Competitiveness Rating
ODBC	Open Database Connectivity
OLEDB	Object Linking and Embedding Database
PDMS	Peer Daten Management System
PK	Primary Key
RAID	Redundant Array of Inexpensive Disks
RD	relevante Datensätze, die in der Realität Duplikate darstellen
ROI	Return on Investment

SAW	Simple Additive Weighting
SFA	Stochastic Frontier Analysis
SNMP	Simple Network Management Protocol
SQL	Structured Query Language
SSIS	SQL Server Integration Services
TOPSIS	Technique for Order Preference by Similarity to Ideal Solution
Web 2.0	populärer Sammelbegriff für Server- und Client-basierte Internettechnologie
WHIRL	Word-based Heterogenous Information Retrieval Logic
WS	Work-Space
WST	Work-Space-Table
XML	Extensible Markup Language
XQuery	XML Query Language

Teil I

Datenbereinigung und Konsolidierung von heterogenen Datenbeständen

– Steven Helmis –

1 Einleitung

Technischer Fortschritt und Globalisierung führten in den vergangenen Jahren zu einem expandierenden Datenaufkommen. Durch den Einsatz von modernen Datenbanktechnologien wird die effektive Verwaltung und Speicherung der Daten mit verschiedenartigen Strukturen im Giga- und Terabyte-Bereich jedoch beherrschbar. Durch das Kommunikationsmedium Internet beispielsweise erfolgt ein weltweiter Zugriff auf verteilte Informationen. Die Folge ist eine Informationsflut, die neben relevanten auch redundante und inkonsistente Fakten beinhaltet. Möchte man solche Daten in unternehmensinternen Informationssystemen einsetzen, stellen Angaben zur Aktualität oder der Vertrauenswürdigkeit einen relevanten Faktor zur Entscheidungsfindung dar.

Dabei ist die mangelnde Datenqualität der immer wiederkehrende Auslöser, der Customer Relationship Management (CRM) Projekte oder Enterprise Resource Planning (ERP) Systeme versagen lässt oder nicht den zu erwartenden Vorteil erbringt. Falsche, fehlende oder veraltete postalische Informationen sind in vielen Fällen die Ursache von Beeinträchtigungen, die schwerwiegende Konsequenzen nach sich ziehen. Inkorrekte Adressen sind verantwortlich für das Scheitern von internen und externen Kommunikationsprozessen in Unternehmen [UNI03]. Aber nicht nur personengebundene Daten können Mängel aufweisen. Auch bei der Präsentation naturwissenschaftlicher Daten sind Konflikte, die auf Grund von differenzierenden Erfassungsmethoden und diversen syntaktischen und semantischen Heterogenitäten zu einer verminderten Qualität des Datenbestandes führen nicht auszuschließen [MaJBL05].

Betrachtet man die zunehmende Globalisierung in mittelständischen Unternehmen, wächst die Nachfrage an Business Intelligence (BI) Produkten zum Sammeln und Aufbereiten von Daten. Das zeitnahe Darstellen von geschäftsrelevanten Informationen, wie z.B. den Ablauf oder die Ergebnisse laufender bzw. abgeschlossener Geschäftsprozesse fördert strategische Unternehmensentscheidungen. Liegen hier fehlerhafte Informationen vor, sind die Konsequenzen oft mit erhöhten Kosten verbunden. Dieser Aspekt erfordert eine stringente Datenqualitätssicherung während des *ETL-Prozesses*.

1.1 Motivation

Die durch ein BI-Werkzeug erzeugten Informationen bilden den Schlüssel für geschäftliche Innovationen. Kei Shen behauptet in [She06]: „*Unternehmen, die Informationsintegration mit maximaler Effizienz vorantreiben, erzielen mit fünf Mal höherer Wahrscheinlichkeit mehr Wertschöpfung.*“ Diese Aussage wird bekräftigt durch die Studie „Business Intelligence im Mittelstand“ von Dirk Fridrich und Dr. Carsten Bange vom Business Application Research Center [FB07]. In dieser Untersuchung wurden mittelständische deutsche Unternehmen mit einem Jahresumsatz zwischen 50 Millionen und einer Milliarde Euro und einer Mitarbeiteranzahl zwischen 100 und 10000 befragt. 279 ausgefüllte Fragebögen konnten zur Auswertung herangezogen werden. Es zeigte sich, dass bereits 48 % der Befragten Firmen Software zur Unternehmenssteuerung einsetzen und 40 % eine Anschaffung planen. Weiterhin stellte sich heraus, dass Datenqualität die wichtigste Eigenschaft als auch der bedeutendste Kritikpunkt an der Business-Intelligence-Software ist. Durch diese Studie wird ersichtlich, dass Lösungen in den Bereichen Berichtswesen, Planung, Datenanalyse, Budgetierung oder Konsolidierung in Unternehmen betreut bzw. noch benötigt werden. Für eine erfolgreiche Realisierung ist jedoch ein aktueller und qualitativ hochwertiger Datenbestand Voraussetzung.

2006	gesamt	pro Stunde
Geborene:	672.724	78
Gestorbene:	821.627	95
Eheschließung:	373.681	43
Ehescheidungen (2005):	201.693	23
Zuzüge über die Grenzen Deutschlands:	707.352	82
Fortzüge über die Grenzen Deutschlands:	628.399	73
Gewerbeanmeldungen:	895.144	104

Tabelle 1.1: Statistische Kennzahlen [SÄD07, Bund06]

Anhand der Tabelle 1.1 mit statistischen Kennzahlen der Bundesrepublik Deutschland (vgl. [SÄD07], [Bun06]) wird die Änderungsgeschwindigkeit von adressbezogenen Daten dargestellt. Betrachtet man die rapiden Veränderungen in den verschiedenen Sektoren, ist davon auszugehen, dass Adressdatenbanken keine 100-prozentigen aktuellen Anschriften enthalten. Der momentane Zustand der Daten kann nur über diverse Software, wie sie in Abschnitt 5.1 vorgestellt wird ermittelt und gegebenenfalls modifiziert werden.