

Wege aus dem Datenlabyrinth

- Datenqualität auf dem Prüfstand -

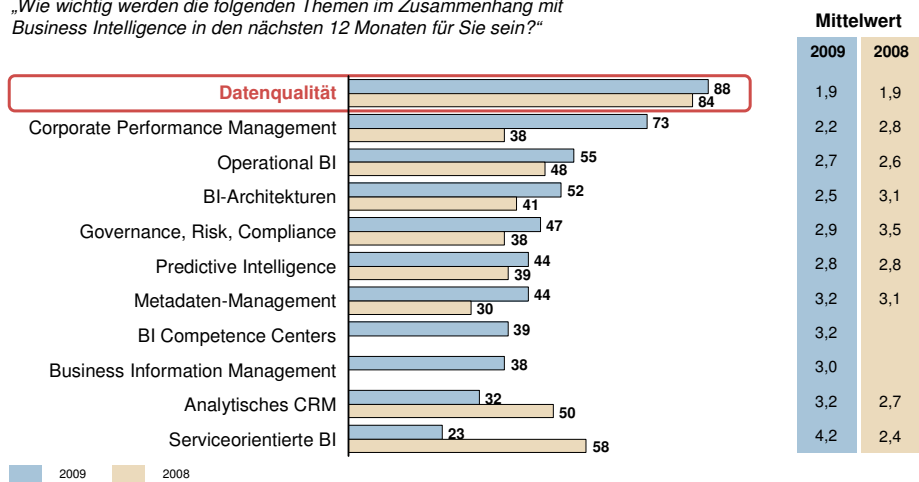
17. November 2009



Studie IT-Trends 2009 in Deutschland: Das BI-Top-Thema ist Datenqualität

Business Intelligence: Bedeutung einzelner Themen [%]

„Wie wichtig werden die folgenden Themen im Zusammenhang mit Business Intelligence in den nächsten 12 Monaten für Sie sein?“



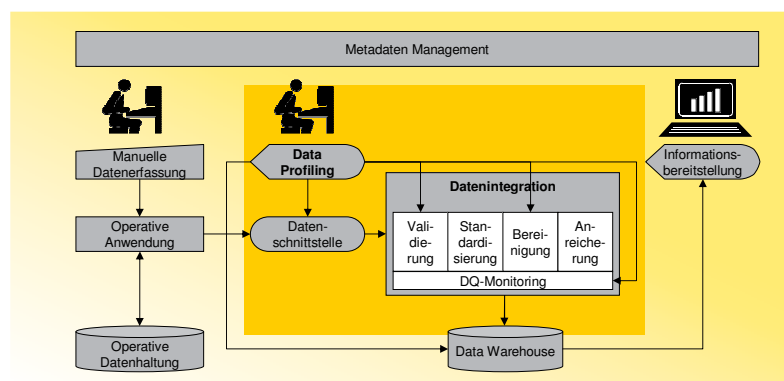
Quelle: Capgemini IT-Trends 2009; Basis: Befragte, die Business Intelligence für eines der 3 wichtigsten Themen halten (n = 35)



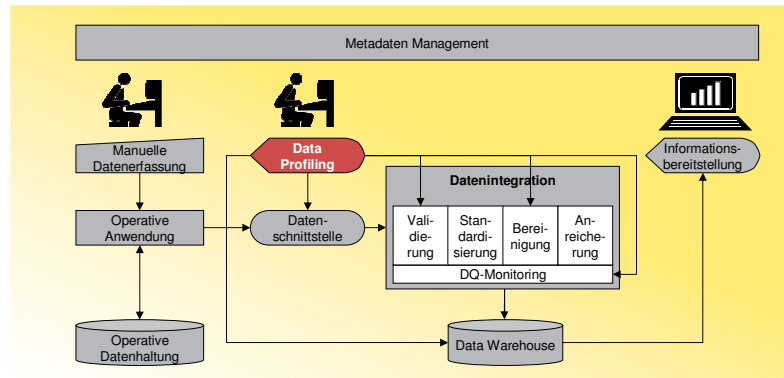
AGENDA

- Data Profiling
- Datenintegrationsprozess
 - Validierung
 - Standardisierung und Bereinigung von Adressen
 - Datenanreicherung
 - DQ-Monitoring
- Fazit

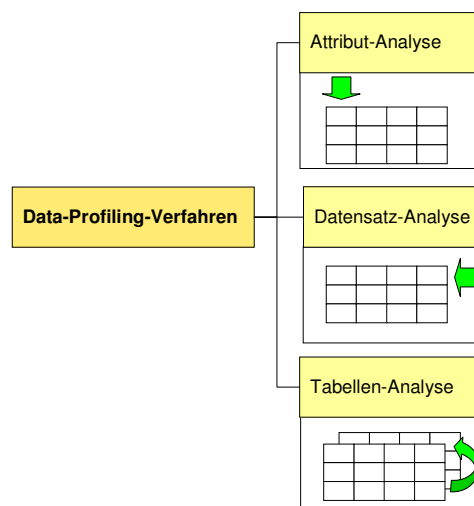
Bausteine eines erfolgreichen Datenqualitätsmanagements



Data Profiling automatisiert den Prozess der Analyse von Dateninhalten



Nur die richtige Kombination der Data-Profiling-Verfahren führt zum Erfolg



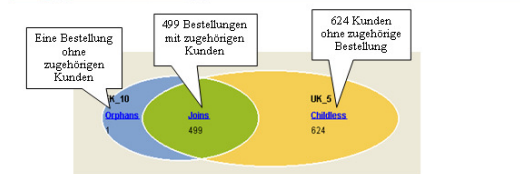
Die referenziellen Abhängigkeiten zwischen Tabellen werden zur Fehlersuche genutzt (Tabellen-Analyse)

Standardanalysen

- referenzielle Abhängigkeiten
- Kardinalitäten
- redundante Attribute
- Anzahl Datensätze

Here are the referential analysis results for BESTELLUNG, which has 8 columns and 500 rows.

Relationship	Type	Documented?	Discovered?	Local Attributes	Remote Attributes	Remote Relation	Cardinality Range	# Orphans	% Compliant
FK_10	Foreign Key	No	Yes	KUNDENCODE	KUNDEN_ID	KUNDEN_STAMM	1..-1	1	99.8%



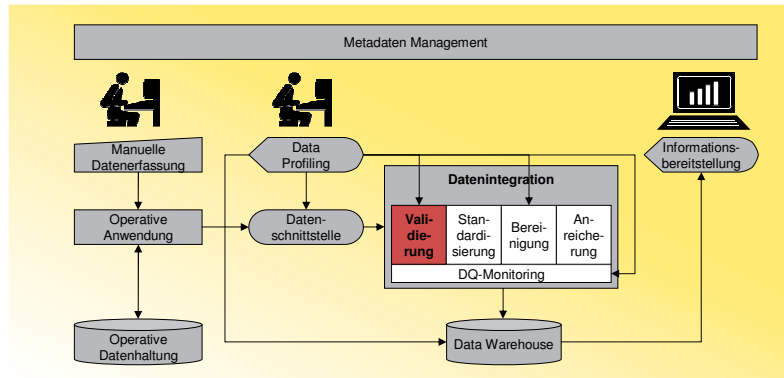
Gruppierung	Wert?	# Zeilen	Ø Zeilen
Monat	Januar 2008	214222	214222
Monat	Februar 2008	237123	225.673
Monat	März 2008	227528	226.291
Monat	April 2008	214789	223.416
Monat	Mai 2008	241927	227.118
Monat	Juni 2008	216456	225.341
Monat	Juli 2008	123546	210.799

Bei der Durchführung des Projekts sollten die vier wichtigsten Faktoren für ein effizientes und erfolgreiches Data Profiling beachtet werden

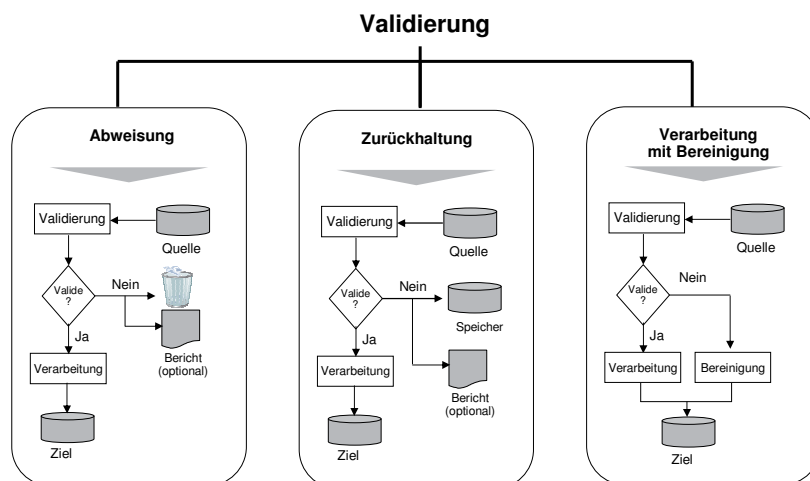
Erfolgsfaktoren

- methodisches Vorgehen
- fundierte Kenntnisse über die Fachlichkeit und die zugehörigen Prozesse
- optimale Zusammensetzung und Erfahrung des Teams
- richtige Kombination und Einsatz der verschiedenen Verfahren

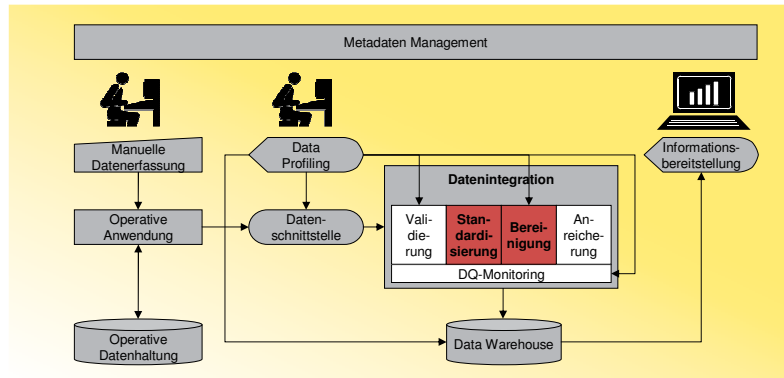
Validierung



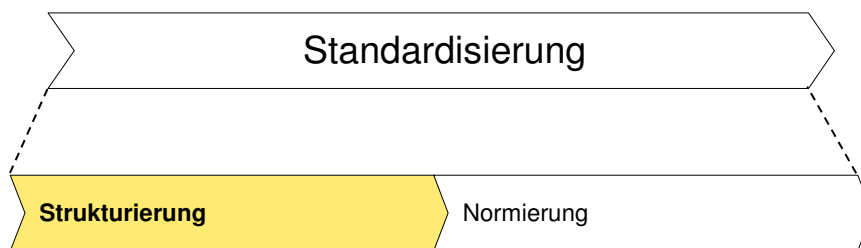
Verfahren zur Validierung



Standardisierung und Bereinigung von Adressen

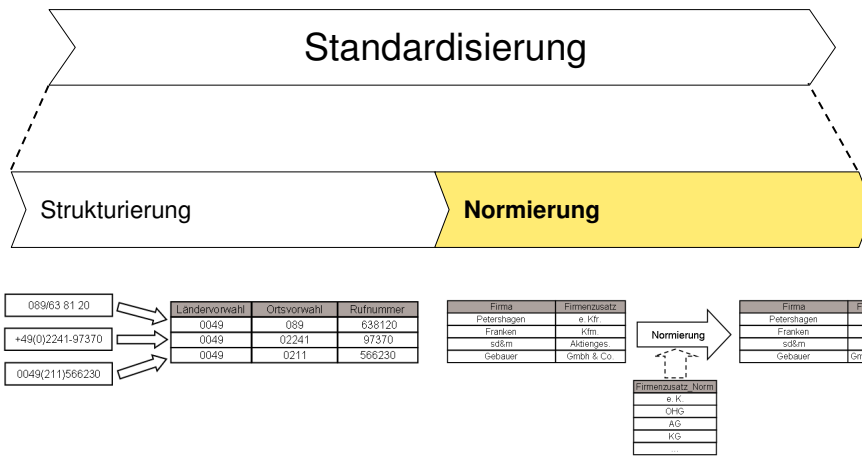


Am besten lassen sich standardisierte Adressen bereinigen



	Ländervorwahl	Ortsvorwahl	Rufnummer
099/63 81 20	0049	099	638120
+49(0)2241-97370	0049	02241	97370
0049(211)566230	0049	0211	566230

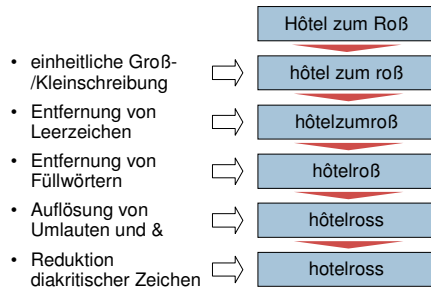
Am besten lassen sich standardisierte Adressen bereinigen



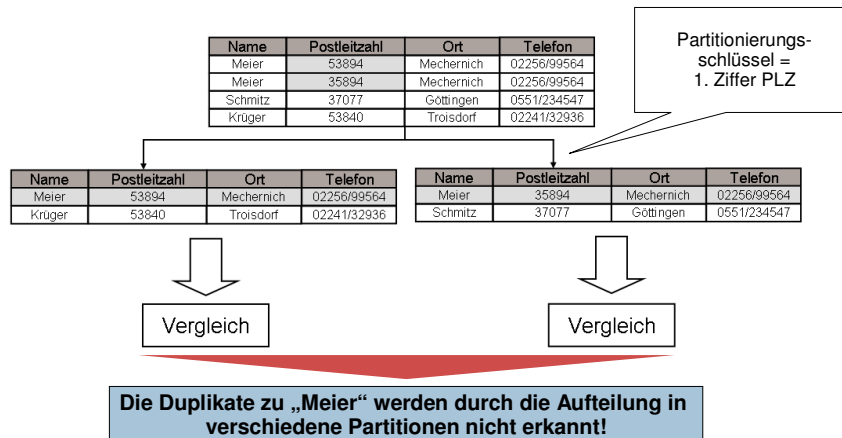
Der Prozess zur Entfernung von Duplikaten gliedert sich in vier Schritte



Die Vorverarbeitung legt den Grundstein für den Erfolg



Die Partionierung steigert die Performance, birgt aber Risiken



Zur Erkennung von Duplikaten werden die Adressen miteinander verglichen



Vorname	Name	Strasse	PLZ
Ulrich	Meier	Mülheimer Straße	53840
Ulrich	Mayer	Mülheimer Straße	33480

Vergleichsfunktion

- ✓ „Der Vorname ist an weniger als drei Stellen unterschiedlich..“
- ✓ „Der Nachname ist phonetisch ähnlich“
- ✗ „Die ersten zwei Ziffern der PLZ sind gleich“

Entscheidungsfunktion

Duplikate

Nicht-Duplikate

Vorname	Name	Strasse	PLZ
Ulrich	Meier	Mülheimer Straße	53840
Ulrich	Mayer	Mülheimer Straße	33480

Mit der passenden Metrik lässt sich die Ähnlichkeit messen



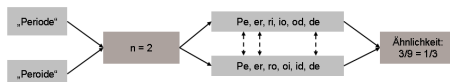
Textuelle Werte

- Phonetik
 - Soundex
 - Kölner Phonetik
- Editierabstand
 - Levenshtein-Distanz
- n-Gramme
- Rekursiver Feldabgleich
- WHIRL
- Ontologien

Numerische Werte

- Differenzbildung
- Hamming-Abstand
- Ontologien
- Konvertierung in textuelle Werte

Bi-Gramme (n = 2)



Hamming-Abstand



Aus den Duplikaten entsteht ein konsolidierter „Goldener Datensatz“

Vorverarbeitung

Partitionierung

Duplikaterkennung

Konsolidierung

Konsolidierung

... identischer Duplikate

... ergänzender Duplikate

... widersprechender Duplikate

Maßnahmen

- Löschung der überzähligen Datensätze

Artikel ID	Bezeichnung	Lieferant ID	Artikel ID	Bezeichnung	Lieferant ID
100	Schlüssel, 13 mm	234	202	Nuss, VK, 13 mm	649
102	Schlüssel, 15 mm	243	102	Schlüssel, 15 mm	243

Artikel ID	Bezeichnung	Lieferant ID
100	Schlüssel, 13 mm	234
102	Schlüssel, 15 mm	243
202	Nuss, VK, 13 mm	649

- Zusammenführung der Datensätze

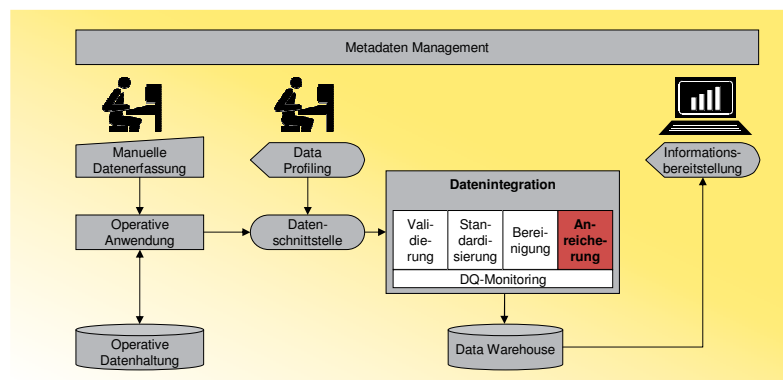
Artikel ID	Grp ID	Bezeichnung	Preis Netto	Artikel ID	Bezeichnung	Lieferant ID
100	1	Schlüssel, 13 mm	4,39	100	Schlüssel, 13 mm	234
102	1	Schlüssel, 15 mm	4,73	102	Schlüssel, 15 mm	243
202	1	Nuss, VK, 13 mm	3,67	202	Nuss, VK, 13 mm	649

Artikel ID	Grp ID	Bezeichnung	Preis Netto	Lieferant ID
100	1	Schlüssel, 13 mm	4,39	234
102	1	Schlüssel, 15 mm	4,73	243
202	1	Nuss, VK, 13 mm	3,67	649

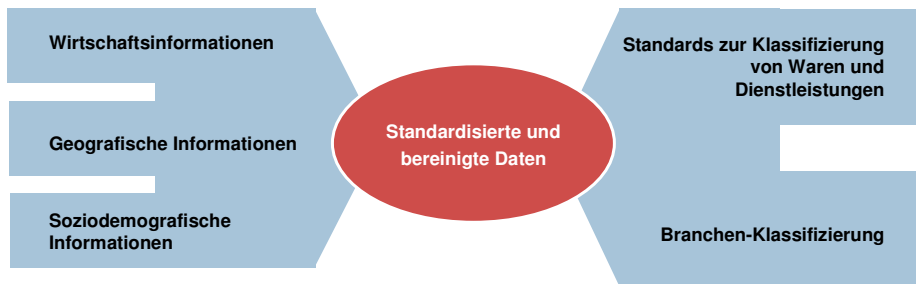
- Konfliktvermeidende Verfahren
- Konfliktlösende Verfahren

Artikel ID	Grp ID	Bezeichnung	Preis Netto	Artikel ID	Bezeichnung	Preis Netto
100	1	Schlüssel, 13 mm	4,39	100	Schlüssel, 13 mm	3,89
102	1	Schlüssel, 15 mm	4,73	102	Schlüssel, 15 mm	4,02
202	1	Nuss, VK, 13 mm	3,67	202	Nuss, VK, 13 mm	3,56

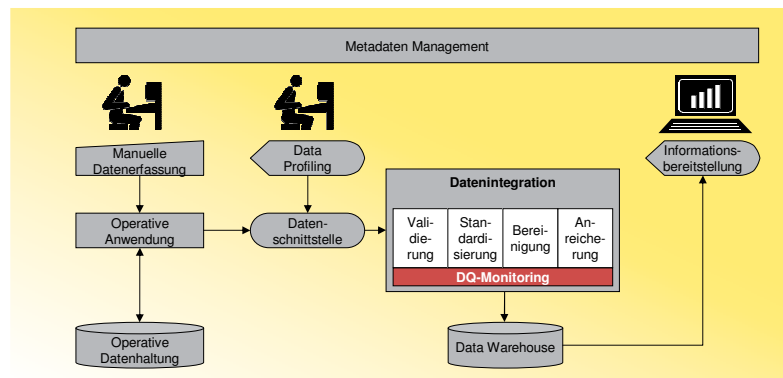
Datenanreicherung



Daten werden angereichert um ihren Informationswert zu steigern

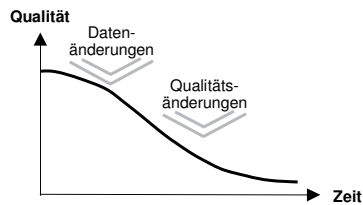


Fortlaufendes Monitoring der Datenqualität sichert den nachhaltigen Erfolg des Datenqualitätsmanagements

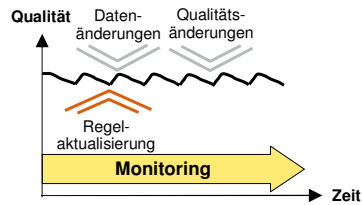


Voraussetzung für die Erhaltung der erreichten Datenqualität ist fortlaufendes Monitoring

Datenqualitäts-Lifecycle (ohne Monitoring)



Datenqualitäts-Lifecycle (mit Monitoring)



Das Monitoring im Betrieb muss automatisiert erfolgen.

DQ-Kennzahlen aus dem Monitoring werden den DQ-Kriterien zugeordnet

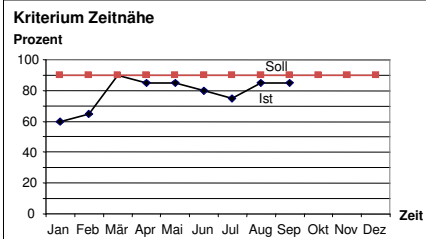
Darstellungsformen der Ergebnisse aus dem Monitoring

Ampeldarstellung

	SOLL	IST	AMPEL
Korrektheit	95%	96%	○ ○ ●
Vollständigkeit	90%	88%	○ ● ○
Zuverlässigkeit	95%	75%	● ○ ○

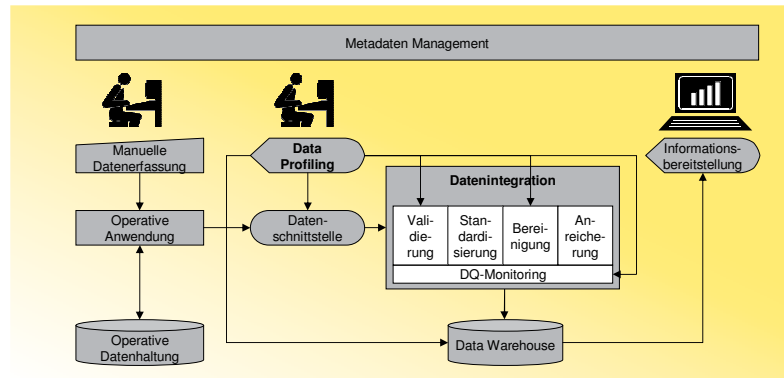
DQ-Kennzahlen/-Kriterien werden in Bezug auf Schwellwerte betrachtet analog zu Dashboards im Umfeld von Business Intelligence.

Zeitreihenbetrachtung



Die zeitliche Entwicklung von DQ-Kennzahlen/-Kriterien wird betrachtet, häufig mit Angabe einer Ziellinie oder eines Toleranzbereichs.

Bausteine eines erfolgreichen Datenqualitätsmanagements



Produktauswahl

<i>Kategorie</i>	Data Profiling	Datenvalidierung	Datenbereinigung	Name & Adresse	Dubletten	DQ-Monitoring
Anbieter						
Oracle	X	X	(OEM Trillium)	X	X	
IBM	X	X	X	X	X	X
Informatica	X	X	X	X	X	X
SAP	X	X	X	X	X	X
DataFlux	X	X	X	X	X	X
Human Inference	X	X	X	X	X	X
Trillium	X	X	X	X	X	X
Uniserv	X	X	X	X	X	X
Innovative Systems	X	X	X	X	X	
Pitney Bowes (Group 1)	X	X	X	X	X	X

Fazit

Die Hoffnung vieler Unternehmen auf Lösung dieses Problems durch die bloße Einführung neuer Systeme oder standhaftem Ignorieren der Problematik schwindet und macht endlich Platz für **wirksame Maßnahmen**.

- Der Schlüssel zum Erfolg ist ein proaktives **Datenqualitätsmanagement**.
- Insbesondere das **Data Profiling** erkennt viele Schwachstellen in den Quelldaten und sollte zu Beginn eines Projektes eingesetzt werden.
- Durch **Standardisierung, Bereinigung** und **Anreicherung** der Daten kann deren Wert für ein Unternehmen sehr stark erhöht werden.
- Das **Data Quality Monitoring** bietet im produktiven Betrieb die Sicherheit, dass Veränderungen und neue Qualitätsprobleme in den Daten rechtzeitig erkannt werden.

Es gibt keine Ausreden mehr für schlechte Daten – fangen Sie an!

Die Expertise zum Datenqualitätsmanagement in BI-Projekten wurde in einem Fachbuch gebündelt

Detlef Apel, Wolfgang Behme, Rüdiger Eberlein, Christian Merighi
“Datenqualität erfolgreich steuern”

Buchteil:	Theorie	Technische Umsetzung	Projektpraxis
Inhalte:	<ul style="list-style-type: none">• Datenqualität• Ursachen und Ausprägungen schlechter Datenqualität• Auswirkungen schlechter Datenqualität• Organisation• Referenzarchitektur für Business-Intelligence-Anwendungen• Kennzahlen zur Messung der Datenqualität	<ul style="list-style-type: none">• Verbesserung der Datenqualität im Quellsystem• Data Profiling• Erfolgreiche Datenvalidierung und -filterung• Standardisierung und Bereinigung• Datenanreicherung• Verbesserung der Datenqualität in der Bereitstellung und Präsentation• Metadaten Management• Data Quality Monitoring• Produktauswahl/-integration	<ul style="list-style-type: none">• Datenqualitätsmanagement in einer Studie• Datenqualitätsmanagement in der Spezifikation• Datenqualitätsmaßnahmen in der Konstruktionsphase• Steuerung der Datenqualität in der Realisierung• Steuerung der Datenqualität im Betrieb

