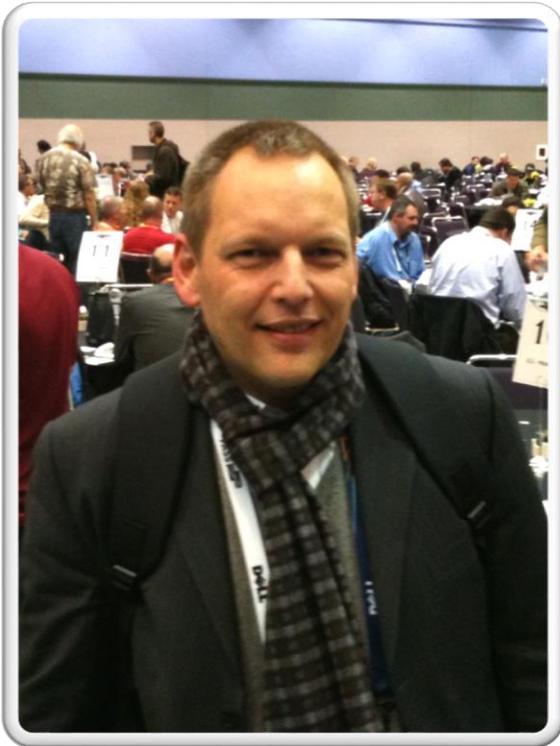# Datenqualität mit SSIS und DQS

## Alexander Karl

PASS

# Speaker

Alexander  Karl

**.net - CDE**          **SQL + BI Consultant**

**Microsoft** CERTIFIED *Trainer*

**Microsoft** CERTIFIED *IT Professional* | Database Administrator 2008
Server Administrator on Windows Server® 2008
Database Administrator on SQL Server® 2005

# Agenda

- Definitionen rund um „Datenqualität"

- bisherige Lösungen in T-SQL und CLR

- Lösung in SSIS

- Lösung mit DQS

# „Datenqualität"

- **Glaubwürdigkeit**
  sind die Daten in einem korrekten Intervall ?

- **Interpretierbarkeit**
  werden Daten fortlaufend gleich dargestellt ?

- **Schlüsselintegrität**
  Schlüsseleindeutigkeit / Ref. Integrität

- **Nützlichkeit**
  Zeitnähe, Vollständigkeit

# das liefert T-SQL  (1/3)

- **Keine NULL Werte als Tabelleneintrag**
  ( default-Values )

- **Constraints zw. mehreren Spalten**
  ( Eintrittsdatum  <=  Austrittsdatum )

- **maximale Textlänge**
  ??   *min. Textlänge* bei Nachnamen z.B. 2

# das liefert T-SQL  (2/3)

- Passende  T-SQL  Funktionen
  -- IsNumeric
  -- IsDate
  -- *IsText*  ??

- Neu ab 2012
  -- Try_Cast
  -- Try_Convert
  -- Try_Parse

# das liefert T-SQL  (3/3)

- **exakte Übereinstimmung**
  -- WHEN  IN  (  … )
  -- JOIN
  -- MERGE

- **ähnliche Übereinstimmung**
  -- Soundex            (T-SQL)
  -- Difference         (T-SQL)
  -- Fuzzy Lookup    (SSIS  Enterprise)
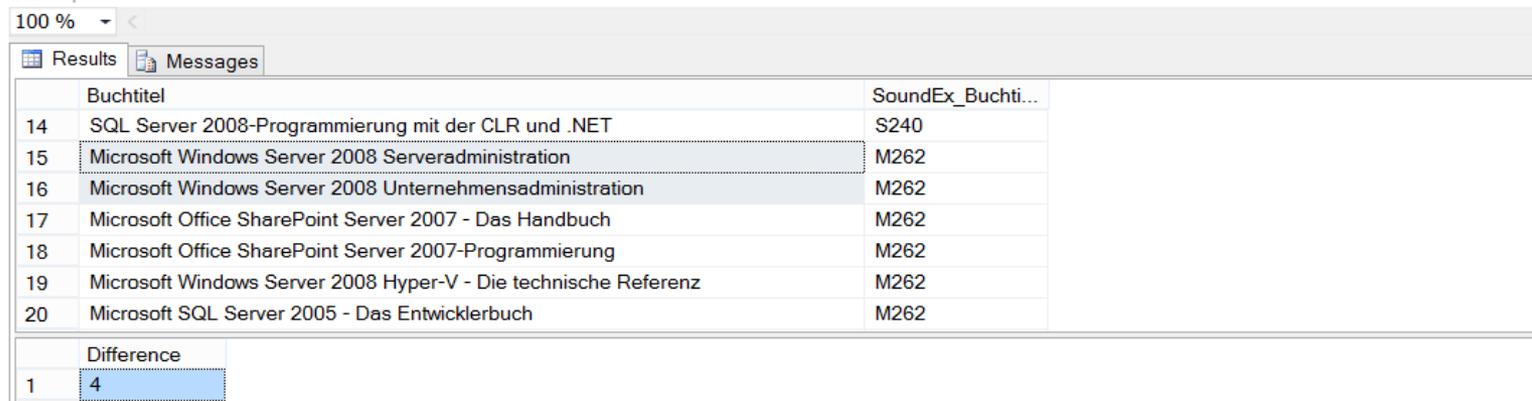  -- Fuzzy Grouping (SSIS  Enterprise)

# DEMO

# Demo result

```sql
1  -- Demo 01 TSQL
2  -->> passt der Import-String in die Zieltabelle ??
3  USE <Database_Name>
4  GO
5
6  SELECT *
7  FROM    INFORMATION_SCHEMA.Columns
8
9  ------ alternative/ früher
10 SELECT t.name     as 'Tab_Name'
11      , c.name     as 'Col_Name'
12      , c.*
13 FROM   sys.tables t
14 join   sys.columns c
15 ON     t.object_id = c.object_id
16
17 ------
18
19 SELECT b.Buchtitel, Len(b.Buchtitel) as 'Len_Buchtitel'
20 FROM   dbo.Buch b
21 ORDER  by Len(b.Buchtitel) DESC
22
```

# Demo result

```
24  -- Demo 01 TSQL
25  -->> gibt es phonetisch ähnliche Einträge
26
27  SELECT b.Buchtitel, SoundEx(b.Buchtitel) as 'SoundEx_Buchtitel'
28  FROM   dbo.Buch b
29  ORDER  by 2 DESC
30
31  SELECT Difference( 'Microsoft Windows Server 2008 Serveradministration'
32                   , 'Microsoft Windows Server 2008 Unternehmensadministration'
33                   ) as 'Difference'
34
35  --  http://support.microsoft.com/kb/100365
36
```

100 %

**Results** | **Messages**

|    | Buchtitel | SoundEx_Buchti... |
|----|-----------|-------------------|
| 14 | SQL Server 2008-Programmierung mit der CLR und .NET | S240 |
| 15 | Microsoft Windows Server 2008 Serveradministration | M262 |
| 16 | Microsoft Windows Server 2008 Unternehmensadministration | M262 |
| 17 | Microsoft Office SharePoint Server 2007 - Das Handbuch | M262 |
| 18 | Microsoft Office SharePoint Server 2007-Programmierung | M262 |
| 19 | Microsoft Windows Server 2008 Hyper-V - Die technische Referenz | M262 |
| 20 | Microsoft SQL Server 2005 - Das Entwicklerbuch | M262 |

|   | Difference |
|---|------------|
| 1 | 4 |

# workaround  mit  CLR

- ## CLR Functions

  mit  RegEx
  --  Textsuche  ( -- *IsText* )
  mit  Levenshtein
  --  ähnliche Übereinstimmung

- ## !!  Verwendung im Dataflow

# Regular Expressions

| Anwendung | Regex |
|---|---|
| Buchstaben & Ziffern | \w |
| nur Buchstaben | ^[a-zA-Z] |
| Zahlzeichen | \d |
| ISBN | ^978\-\d\-\d{5}\-\d{3}\-\d |
| email | \b[A-Z0-9._%+-]+@[A-Z0-9.-]+\.[A-Z]{2,4}\b |

# Levenshtein

W http://de.wikipedia.org/wiki/Levenshtein-Distanz

W Levenshtein-Distanz... ×

Benutzerkonto anlegen   Anmelden

Artikel   Diskussion        Lesen   Bearbeiten   Versionsgeschichte    Suchen

**WIKIPEDIA**
Die freie Enzyklopädie

## Levenshtein-Distanz

Die **Levenshtein-Distanz** (auch **Editierdistanz**) zwischen zwei Zeichenketten ist die minimale Anzahl von Einfüge-, Lösch- und Ersetz-Operationen, um die erste Zeichenkette in die zweite umzuwandeln. Benannt ist die Distanz nach dem russischen Wissenschaftler Wladimir Lewenstein, der sie 1965 einführte. Mathematisch ist die Levenshtein-Distanz eine Metrik auf dem Raum der Symbolsequenzen.

Beispielsweise ist die Levenshtein-Distanz zwischen „Tier" zu „Tor" 2. Eine mögliche Folge von zwei Operationen ist:

1. Tier
2. Toer (Ersetze i durch o)
3. Tor (Lösche e)

## Beispiel [Bearbeiten]

Das Verfahren lässt sich nun leicht in einer Tabelle durchführen. Hier ein Beispiel mit den beiden Zeichenketten „Tier" und „Tor".

```
   ε T o r
ε  0 1 2 3
T  1 0 1 2
i  2 1 1 2
e  3 2 2 2
r  4 3 3 2
```

Hauptseite
Themenportale
Von A bis Z
Zufälliger Artikel

▼ Mitmachen
   Artikel verbessern

13 | Dez. 2013   **Datenqualität mit SSIS und DQS**

# codeplex

# DEMO

**Datenqualität mit SSIS und DQS**

# Demo result

```
 1    -- Demo 02 CLR
 2   use test_CLR
 3
 4    -- aktivieren der CLR-Integration
 5    EXEC sp_configure 'clr enabled', 1
 6    go
 7    reconfigure
 8    go
 9   ---------
10
11    --Listing C# Assembly
12   using System;
13   using System.Data;
14    ...
15
16   --> mit obigem Listing (xmlSaveProc.cs) muss mittels Compiler eine .dll erstellt werden.
17    --  csc /t:library filename.cs
18    --  evtl vorher path configurieren auf C:\WINDOWS\Microsoft.NET\Framework\v2.0.50727
19    go
20
21    -- Registrieren der Assembly
22   CREATE ASSEMBLY CLR_Levenshtein
23    FROM 'C:\ <folder> \CLR_Levenshtein.dll'
24    WITH Permission_Set = Safe
25
26    -- überprüfung
27    SELECT * FROM sys.assemblies
28    SELECT * FROM sys.assembly_files
```

# Demo result

```
CLR_Levenshtein.cs    ✕
StoredFunctions                          ▼    Levenshtein(SqlString S1, SqlString S2)
1   using System;
2   using System.Data;
3   using System.Data.SqlClient;
4   using System.Data.SqlTypes;
5   using Microsoft.SqlServer.Server;
6
7   public partial class StoredFunctions
8   {
9       [Microsoft.SqlServer.Server.SqlFunction(IsDeterministic = true, IsPrecise = false)]
10      public static SqlDouble Levenshtein(SqlString S1, SqlString S2)
11      {
12          if(S1.IsNull)
13              S1 = new SqlString("");
14          if(S2.IsNull)
15              S2 = new SqlString("");
16          String SC1 = S1.Value.ToUpper();
17          String SC2 = S2.Value.ToUpper();
18          int n = SC1.Length;
19          int m = SC2.Length;
20
21          int[,] d = new int[n + 1, m + 1];
22          int cost = 0;
23
24          if (n + m == 0) {
25              return 100;
26          } else if (n == 0) {
27              return 0;
28          } else if (m == 0) {
29              return 0;
30          }
31          for (int i = 0; i <= n; i++)
32              d[i, 0] = i;
33          for (int j = 0; j <= m; j++)
34              d[0, j] = j;
35          for (int i = 1; i <= n; i++)
36          {
37              for (int j = 1; j <= m; j++)
38              {
39                  if (SC1[i - 1] == SC2[j - 1])
40                      cost = 0;
41                  else
42                      cost = 1;
43                  d[i, j] = System.Math.Min(System.Math.Min(d[i - 1, j] + 1, d[i, j - 1] + 1), d[i - 1, j - 1] + cost);
44              }
45          }
46
47          double percentage = System.Math.Round((1.0 - ((double)d[n, m]/(double)System.Math.Max(n,m))) * 100.0,2);
48          return percentage;
49      }
50  };
```

# Demo result

```sql
1  -- Demo 02 CLR
2  -- Ausführung
3
4  CREATE Function fn_Levenshtein( @S1 nvarchar(4000) , @S2 nvarchar(4000))
5  RETURNS float
6  AS
7    EXTERNAL NAME CLR_Levenshtein.StoredFunctions.Levenshtein
8    --            DLL            .Class           .Function name
9
10
11 SELECT dbo.fn_Levenshtein( 'SQL Server Internals'
12                          , 'SQL Server Integration'
13                          )  as 'percentDiff'
14
```
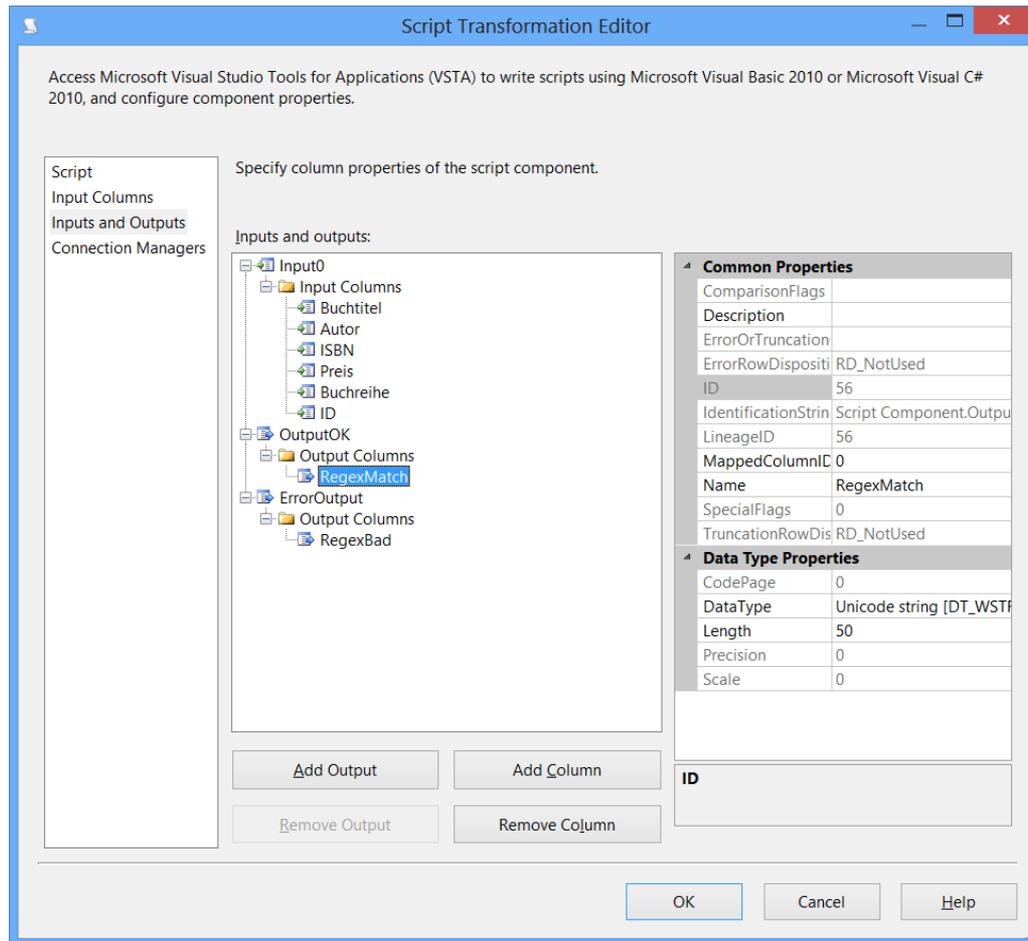
100 %  ▾

Results | Messages

| | percentDiff |
|---|---|
| 1 | 72,73 |

# SSIS  custom components

- C#  und  VB .NET

- „frei programmierbar"

- synchron sehr performant

- optimale Verwendung im Dataflow

- konfigurierbar via „toolDB"

# SSIS custom components

# SSIS  custom components

```
105  ⊟       Public Overrides Sub Input0_ProcessInputRow(ByVal Row As Input0Buffer)
106               '
107               Dim RegexPattern As String
108               RegexPattern = "^978\-\d\-\d{5}\-\d{3}\-\d$"
109               ''--  ^      >> Beginn der Übereinstimmung
110               ''--  \-     >> "-" (= minus) ist Sonderzeichen, desh. \ als Escape
111               ''--  \w     >> beliebiges Wortzeichen
112               ''--  \d     >> beliebige Ziffer
113               ''--  $      >> Ende der Übereinstimmung
114
115               If Not Regex.IsMatch(Row.ISBN, RegexPattern) Then
116
117                   Me.ErrorOutputBuffer.AddRow()
118                   Me.ErrorOutputBuffer.RegexBad = Row.ISBN
119
120                   Row.RegexMatch = "!!"
121               Else
122
123                   Row.RegexMatch = "ok"
124               End If
125
126               '
127           End Sub
128
129   End Class
```

# DEMO

# Demo result



**Datenqualität mit SSIS und DQS**

# Demo  result

OutputOK Data Viewer at DFT_Demo_03

Detach    Copy Data

| Buchtitel | Autor | ISBN | Preis | . | ID | RegexMatch |
|---|---|---|---|---|---|---|
| Konfigurieren von Windows Server 2008 Active Directory | 4 | 978-3-86645-940-3 | 79,00 | | 1 | ok |
| Konfigurieren einer Windows Server 2008-Netzwerkinfrastruktur | 11 | 978-3-86645-942-7 | 79,00 | | 2 | ok |
| Microsoft Windows Server 2008 Serveradministration | 12 | 978-3-86645-946-5 | 79,00 | | 3 | ok |
| Konfigurieren von Windows 7 MCTS | 12 | 978-3-86645-980-9 | 79,00 | | 4 | ok |
| Microsoft Windows Server 2008 Unternehmensadministration | 12 | 978-3-86645-947-2 | 79,00 | | 5 | ok |
| Microsoft Office SharePoint Server 2007 - Das Handbuch | 3 | 978-3-86645-117-9 | 59,00 | | 6 | ok |
| Microsoft Office SharePoint Server 2007-Programmierung | 18 | 978-3-86645-642-6 | 39,90 | | 7 | ok |
| Internetinformationsdienste (IIS) 7.0 - Die technische Referenz | 10 | 978-3-86645-924-3 | 79,00 | | 8 | ok |
| Cloud Computing mit der Windows Azure Platform | 17 | 978-3-86645-533-7 | 39,90 | | 9 | ok |
| Konfigurieren der Windows Server-Virtualisierung | 14 | 978-3-86645-952-6 | 79,00 | | 10 | ok |
| Microsoft Windows Server 2008 Hyper-V - Die technische Referenz | 9 | 978-3-86645-926-7 | 79,00 | | 11 | ok |
| Implementieren und Warten von SQL Server 2008 MCTS | 5 | 978-3-86645-932-8 | 79,00 | | 12 | ok |
| SQL Server 2008 - Database Development | 5 | 978-0-7356-2639-3 | 70,00 | | 13 | !! |
| Business Intelligence und Reporting mit SQL Server 2008 | 1 | 978-3-86645-657-0 | 59,00 | | 14 | ok |
| Business Intelligence mit Office 2007 und SQL Server | 1 | 978-3-86645-637-2 | 49,90 | | 15 | ok |
| SQL Server 2008 | 13 | 978-3-86645-519-1 | 39,90 | | 16 | ok |

Attached    Total rows: 0, buffers: 0                    Rows displayed = 31

# SSIS Data Profiling Task

- <xml>  output

- Data Profiling Viewer

- Xquery für „handmade Analyse"

# DEMO

# Demo result

# Demo result

```xml
Dataprofiling_Result.xml   ×
 1    <?xml version="1.0"?>
 2    <DataProfile xmlns:xsd="http://www.w3.org/2001/XMLSchema"
 3                 xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
 4                 xmlns="http://schemas.microsoft.com/sqlserver/2008/DataDebugger/">
 5      <ProfileVersion>1.0</ProfileVersion>
 6      <DataSources>...</DataSources>
14      <DataProfileInput>...</DataProfileInput>
55      <DataProfileOutput>
56        <Profiles>
57          <ColumnValueDistributionProfile IsExact="true" ProfileRequestID="ValueDistReq">
58            <DataSourceID>{33E8D61B-6415-4970-8B54-FAFD156C27DD}</DataSourceID>
59            <Table DataSource="localhost" Database="PerformanceDB" Schema="dbo" Table="View_1" RowCount="50" />
60            <Column Name="promotion_id" SqlDbType="Int" MaxLength="0" Precision="10" Scale="0" LCID="-1"
61                    CodePage="0" IsNullable="true" StringCompareOptions="32768" />
62            <NumberOfDistinctValues>2</NumberOfDistinctValues>
63            <ValueDistribution>
64              <ValueDistributionItem>
65                <Value>0</Value>
66                <Count>10</Count>
67              </ValueDistributionItem>
68              <ValueDistributionItem>
69                <Value>1160</Value>
70                <Count>40</Count>
71              </ValueDistributionItem>
72            </ValueDistribution>
73          </ColumnValueDistributionProfile>
74          <ColumnLengthDistributionProfile ProfileRequestID="LengthDistReq" IsExact="true">
```

# Demo result

```sql
1  -- Dataprofiling
2  -- ValueDistributionProfile
3
4  Declare @file     varChar(255)
5  Set     @file     = 'C:\ <folder> \Dataprofiling_Result.xml'
6
7  Declare @charVar varChar(max)
8        , @nameSp  varChar(400)
9        , @sqlCmd  varChar(400)
10       , @xmlVar  xml
11
12 Declare @tmpTable Table (col1 varchar(max))
13
14 Set @nameSp  = ' xmlns:xsd="http://www.w3.org/2001/XMLSchema"
15                  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
16                  xmlns="http://schemas.microsoft.com/sqlserver/2008/DataDebugger/"'
17
18 Set @sqlCmd  = ' Select * From OPENROWSET ( BULK ''' + @file + ''', SINGLE_BLOB ) AS x ';
19
20 Insert Into @tmpTable
21 exec( @sqlCmd )
22
23 Set  @xmlVar = ( select Top(1) CAST( Replace( col1, @nameSp, '' )   as xml ) from @tmpTable );
24 ---- select @xmlVar
```

# Demo result

```sql
26  Declare @i int;  -- idoc
27  Execute sp_xml_preparedocument @i OutPut
28                                 , @xmlVar
29
30  SELECT *
31  FROM   OpenXML ( @i, '/DataProfile/DataProfileOutput/Profiles/ColumnValueDistributionProfile
32                                                 /ValueDistribution/ValueDistributionItem' )
33  WITH  ( ProfileRequest  nvarChar(100)  '../../@ProfileRequestID'
34        , SchemaName        nvarchar(100)  '../../Table/@Schema'
35        , TableName         nvarChar(100)  '../../Table/@Table'
36        , RowCnt            nvarChar(100)  '../../Table/@RowCount'
37        , ColumnName        nvarChar(100)  '../../Column/@Name'
38        , ColumnType        nvarChar(100)  '../../Column/@SqlDbType'
39        , DistinctValues    nvarChar(100)  '../../NumberOfDistinctValues'
40        , Value_item        nvarChar(100)  'Value'
41        , Count_item        nvarChar(100)  'Count'
42        )
43
44  Execute sp_xml_removedocument @i
```

100 %

Results | Messages

|   | ProfileRequest | SchemaName | TableName | RowCnt | ColumnName | ColumnType | DistinctValues | Value_item | Count_item |
|---|---|---|---|---|---|---|---|---|---|
| 1 | ValueDistReq | dbo | View_1 | 50 | promotion_id | Int | 2 | 0 | 10 |
| 2 | ValueDistReq | dbo | View_1 | 50 | promotion_id | Int | 2 | 1160 | 40 |

# XML Task

```xml
<?xml version="1.0"?>
- <Buecher>
  - <Buch>
      <Buchtitel>Konfigurieren von Windows Server 2008 Active Directory</Buchtitel>
      <ISBN>978-3-86645-940-3</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Konfigurieren einer Windows Server 2008-Netzwerkinfrastruktur</Buchtitel>
      <ISBN>978-3-86645-942-7</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Microsoft Windows Server 2008 Serveradministration</Buchtitel>
      <ISBN>978-3-86645-946-5</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Konfigurieren von Windows 7 MCTS</Buchtitel>
      <ISBN>978-3-86645-980-9</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Microsoft Windows Server 2008 Unternehmensadministration</Buchtitel>
      <ISBN>978-3-86645-947-2</ISBN>
      <Preis>79.00</Preis>
    </Buch>
```

```xml
<?xml version="1.0"?>
- <Buecher>
  - <Buch>
      <Buchtitel>Konfigurieren von Windows Server 2008 Active Directory</Buchtitel>
      <ISBN>978-3-86645-940-3</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Konfigurieren einer Windows Server 2008-Netzwerkinfrastruktur</Buchtitel>
      <ISBN>978-3-86645-942-7</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Microsoft Windows Server 2008 Serveradministration</Buchtitel>
      <ISBN>978-3-86645-946-5</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Konfigurieren von Windows 7 MCTS</Buchtitel>
      <ISBN>978-3-86645-980-9</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Microsoft Windows Server 2008 Unternehmensadministration</Buchtitel>
      <ISBN>978-3-86645-947-2</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Microsoft Office SharePoint Server 2007 - Das Handbuch</Buchtitel>
      <ISBN>978-3-86645-117-9</ISBN>
      <Preis>59.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Microsoft Office SharePoint Server 2007-Programmierung</Buchtitel>
      <ISBN>978-3-86645-642-6</ISBN>
      <Preis>39.90</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Internetinformationsdienste (IIS) 7.0 - Die technische Referenz</Buchtitel>
      <ISBN>978-3-86645-924-3</ISBN>
      <Preis>79.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Cloud Computing mit der Windows Azure Platform</Buchtitel>
      <ISBN>978-3-86645-533-7</ISBN>
      <Preis>39.90</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>SQL Server 2008-Performance-Optimierung</Buchtitel>
      <ISBN>978-3-82732-778-9</ISBN>
      <Preis>39.90</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>SQL Server 2008 Business Intelligence Development </Buchtitel>
      <ISBN>978-0-73562-636-2</ISBN>
      <Preis>60.00</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>SQL Server 2008 Integration Services (Wrox Programmer to Programmer)</Buchtitel>
      <ISBN>978-0-47024-795-2</ISBN>
      <Preis>39.80</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>SQL Server Analysis Services 2008 with MDX (Wrox Programmer to Programmer)</Buchtitel>
      <ISBN>978-0-47024-798-3</ISBN>
      <Preis>39.80</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>SQL Server 2008 Reporting Services (Wrox Programmer to Programmer)</Buchtitel>
      <ISBN>978-0-47024-201-8</ISBN>
      <Preis>39.80</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>SQL Server 2008-Programmierung mit der CLR und .NET </Buchtitel>
      <ISBN>978-3-86645-436-1</ISBN>
      <Preis>39.90</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>Expert Cube Development with Microsoft SQL Server 2008 Analysis Services </Buchtitel>
      <ISBN>978-1-84719-722-1</ISBN>
      <Preis>29.90</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>The Data Warehouse Toolkit: The Complete Guide to Dimensional Modeling</Buchtitel>
      <ISBN>978-0-47120-024-6</ISBN>
      <Preis>49.90</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>The Microsoft Data Warehouse Toolkit</Buchtitel>
      <ISBN>978-0-47126-715-7</ISBN>
      <Preis>39.90</Preis>
    </Buch>
  - <Buch>
      <Buchtitel>The Data Warehouse ETL Toolkit: Practical Techniques </Buchtitel>
      <ISBN>978-0-76456-757-5</ISBN>
      <Preis>39.90</Preis>
    </Buch>
  </Buecher>
```

=

# XML Task

# DEMO

**Datenqualität mit SSIS und DQS**

# Fuzzy Grouping  + Lookup

# bisherige Lösungen

- Nachteil  „handmade"

- ggf. Portierungsaufwand

- umfangreichere Doku

- Benutzer-Einweisung in workflow

  >>  DQS  ( neu ab 2012 )

# DQS   Data Quality Service

# SQL Server 2012  Editionen

http://msdn.microsoft.com/de-de/library/cc645993.aspx

Von den SQL Server ... ×

▷ Server- und Enterprise-Entwicklung

▷ SQL Server

▷ SQL Server 2012

▷ Produktdokumentation

**Von den SQL Server 2012-Editionen unterstützte Funktionen**

| Funktionsname | Enterprise | Business Intelligence | Standard |
|---|---|---|---|
| Data Quality Services | Ja | Ja | |

# DQS Installer

```
C:\Program Files\Microsoft SQL Server\MSSQL11.SECOND\MSSQL\Binn>DQSInstaller.exe /?
Microsoft (R) DQS Installer Command Line Tool
Copyright (c) 2012 Microsoft. All rights reserved.

[12/12/2013 4:38:34 PM] DQS Installer started. Installation log will be written
to C:\Program Files\Microsoft SQL Server\MSSQL11.SECOND\MSSQL\Log\DQS_install.log

[12/12/2013 4:38:35 PM] Parsing DqsInstaller command line arguments.
usage DqsInstaller.exe [-install | -uninstall | -upgrade | -upgradedlls | -expor
tkbs | -importkbs] [<file name>] [-collation] | [-instance] 'instance name'


-install        - Install Data Quality Services in the provided in stance. (Default)
-uninstall      - Uninstall Data Quality Services from the provided instance.
-upgrade        - Upgrade Data Quality Services for the provided instance, to current version.
-upgradedlls    - Install Data Quality Services while skipping recreating the DQS databases and
                  only upgrade DQS DLLs.
-exportkbs      - Export all server knowledgebases.
-importkbs      - Import knowledgebases file to server.
-collation      - The collation of DB catalogs to install. The collation should be case insensitive.
<file name>     - The .dqsb file name used to import/export server backup data.
-instance       - Specify the SQL Server instance name that this installer will run against.
-?              - Show this usage message.
```

# DQS  Übersicht

# DQS Übersicht

# DQS   Übersicht

# DQS Übersicht

# DQS  zum Einstieg einige Begriffe

- Knowledgebase

- Domain Management

- Knowledge Discovery

- Matching Policy

1ˢᵗ „basic" domain rules

Domain Management

2ⁿᵈ bisherige Daten

Knowledge Discovery

## Aufbau der knowledge base

## Verwendung der knowledge base

Matching Policy

DQS Cleansing

# DEMO

# Demo Results

# Demo  Results

# Demo Results

1st „basic" domain rules

Domain Management

http://datamarket.azure.com/

Windows Azure Marketplace

2nd bisherige Daten

Knowledge Discovery

**Aufbau der knowledge base**

**Verwendung der knowledge base**

Matching Policy

DQS Cleansing

SSIS DQS
Matching Transformation

# DQS  Lösung

- ~~Nachteil  „handmade"~~        standardisiert

- ~~ggf. Portierungsaufwand~~  gut portierbar

- ~~umfangreichere Doku~~      Doku vorhanden

- ~~Benutzer-Einweisung~~        nachvollziehbarer
  ~~in workflow~~                         Workflow für Fachabt.

# DQS      Literaturempfehlung

download
## DQS step-by-step
german  /  english

# DQS      Literaturempfehlung

# vielen Dank