**Chapter 24.** # December 2008

Welcome to the December 2008 edition of IBM InfoSphere Information Server Developer's Notebook. This month we answer the question;

How do I perform looping, and before/after record processing within QualityStage custom rule sets?

*Excellent question! This edition of (IBM) InfoSphere Information Server Developer's Notebook (IISDN) is the last in the current series on the topic of QualityStage custom rules sets, and covers advanced topics, including; looping, before and after row processing, unhandled data, null record output, and more.*

In the example put forth in this document, we process a large structured and custom record format, much as we would the PDF and HTTP file read examples from months prior. We use the looping and before/after record processing capabilities of QualityStage rule sets, to read data structures of these types.

## Software versions

All of these solutions were *developed and tested* on (IBM) InfoSphere Information Server (IIS) version 8.1, using the Microsoft Windows XP/SP2 platform to support IIS client programs, and a RedHat Linux Advanced Server 4 (RHEL 4) FixPak U6 32 bit SMP server (Linux kernel version 2.6.9-67.EL-smp) to support the IIS server side components.

IBM InfoSphere Information Server allows for a single, consistent, and accurate view of data across the full width of the corporate enterprise, be it relational or non-relational, staged or live data. As a reminder, the IBM InfoSphere Information Server product contains the following major components;

WebSphere Business Glossary Anywhere™, WebSphere Information Analyzer™, WebSphere Information Services Director™, WebSphere DataStage™, WebSphere QualityStage™, WebSphere Metadata Server and Metabridges™, WebSphere Metadata Workbench™, WebSphere Federation Server™, Classic Federation™, Event Publisher™, Replication Server™, Rational Data Architect™, DataMirror Transformation Server™, and others.

Obviously, IBM InfoSphere Information Server is a large and capable product, addressing many strategic needs across the enterprise, and supporting different roles and responsibilities.

## 24.1  Terms and core concepts

In the two prior editions of (IBM) InfoSphere Information Server Developer's Notebook (IISDN), October/2008 and November/2008, we covered QualityStage related topics; more specifically, the Investigate and then Standardize stages. Also in the November/2008 edition, we covered QualityStage custom rule sets, including;

– How to create a new QualityStage custom rule set, how to provision same, and how to use this new custom rule set inside an Information Server DataStage/QualityStage component Job.

– How to work with 3 of the components to a QualityStage rule set, specifically; Classification files, Dictionary files, and Pattern Action files.

– We created two new QualityStage custom rule sets-

   • A simple rule set to standardize a company name.

     This was a simple rule, removing common words like (The, A, and An) from the leading portion of a name; this example being used to demonstrate all of the surrounding techniques of creating a new rule set, testing it, and making use of it within a Job.

   • A simple rule set to standardize a U.S. based telephone number.

     This example got more into the Pattern Action language; more pattern action language commands, as well as simple and user-defined pattern classes.

     By example, we also detailed the general scheme behind standardizing data; How to structure a custom rule, how to approach the tokenization of data and related.

In this edition of IISDN, we build the third and final custom rule set (for a while anyway). This rule set demonstrates looping, and before/after row processing. (Before and after each larger input record that is sent to the custom rule set.)

**Note:** This is one of the few editions of IISDN where you really must read and master a prior edition; specifically in this case, you must read and understand the November/2008 edition of IISDN, QualityStage Custom Rule Sets I.

Figure 24-1 displays the sample data and sample Job used in this example. Comments related to this sample data include;

– Figure 24-1 displays a subset, and the hardest part, of a production rule set we encountered. In effect, known column names are presented in the format;

     colN=Some multi-word value

– The approved/accepted column names are from a list of known column names. The examples used in Figure 24-1 include colA, colB, colC, and colD. And these column name, column value pairs are separated by commas. (The last column name, column value pair ends with no comma.)

– Known column name, column value pairs may appear in any order. So data may appear colA, colB, colC, colD, or colD, colB, colA, colC, or whatever.

And previously unknown column names must be sent to an "UnhandledData" column. Examples displayed in Figure 24-1, include; colX, colY, and colZ, but any unknown column name should produce this behavior.

– Lastly, any output column values should be defaulted to some value when NULL.
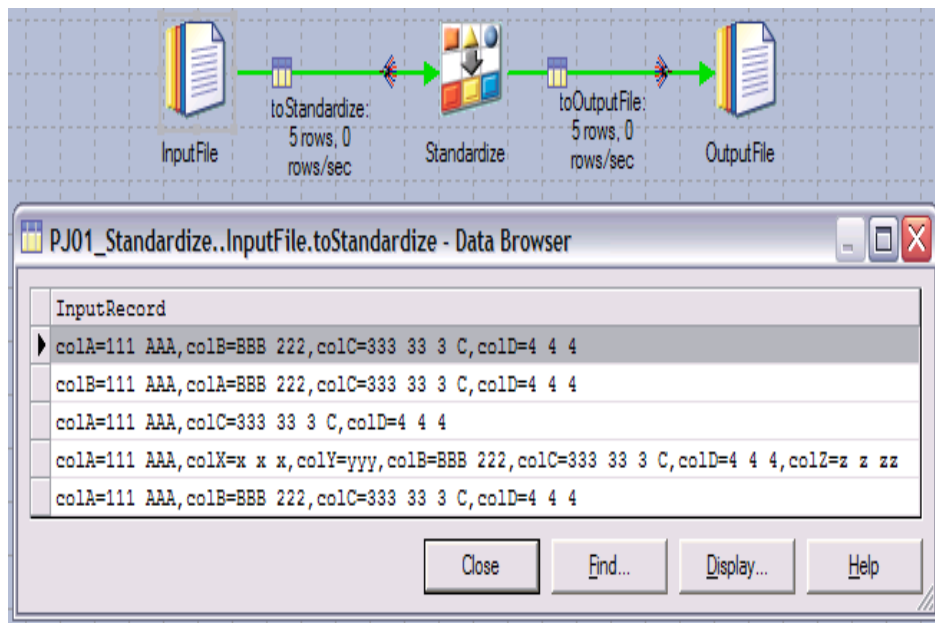


*Figure 24-1   Sample data set used in this example, sample Job.*

# 24.2  Create this example

In order to create the example QualityStage custom rule set detailed above, complete the following steps;

*Remember to use consult the prior month's edition of IISDN as a guide.*

1. Using your favorite editor, create a sample data set with the conditions as displayed in Figure 24-1, and accompanying text.

2. Log on to the DataStage/QualityStage Designer program.

   Create a new custom rule set, and then a new Parallel Job that uses this new custom rule set. The input data source should output one column, which is processed by the Standardize stage displayed in Figure 24-1. To accomplish this, the Format TAB of the Sequential File stage should include;

   – Field defaults -> Quote -> None.

   – Field defaults -> Delimiter -> None.

   – Record level -> Final delimiter -> End.

   – Record Level -> Record delimiter -> UNIX newline.

3. Alter the contents of the Classification file to equal that as displayed in Figure 24-2. A code review follows.



*Figure 24-2   Classification File for this example.*

A code review for Figure 24-2 includes;

   – Basically we enter each known column name value in this file, and give it the user-defined pattern value of "T".

Page 8.

– The single advantage of this step is that we may later test for any known
   column name by looking for "T". This as opposed to saying something like,

   {psuedo code}

   if ((col = "colA") or (col = "colB") or (col = "colC") ... etcetera

   We say single advantage, because the Classification file can provide other
   services, just not any we need for this example.

4. Alter the contents of the Dictionary file to equal that as displayed in
   Figure 24-3. A code review follows.



*Figure 24-3   Dictionary File for this example.*

A code review for Figure 24-3 includes;

– Basically this dictionary outputs 5 columns; colA, colB, colC, colD, and
   UnhandledData.

– UnhandledData is not a keyword, but is an output column name that is
   commonly used.

5. Alter the contents of the Pattern Action file to equal that as displayed in
   Figure 24-4, Figure 24-5, and Figure 24-6. A code review follows.

   *Do not enter the line numbers; line numbers appear for the benefit of the code
   review that follows.*

```
A' RS02_DATAFORMAT01.PAT - Vi for Windows               _ □ ✗

File  Edit  Macro  Options  Help

 1   ;;QualityStage v8.0
 2   ;---------------------------------------------------
 3   ; Parser Rules
 4   ;---------------------------------------------------
 5
 6   ; SEPLIST    " ~`!@#$%^&*()_-+=(){}[]|\\:;\"<>,.?/'"
 7   ; STRIPLIST  " ~`!@$^*_+=(){}[]|\\:;<>?"
 8
 9   \PRAGMA_START
10   SEPLIST        " ~`!@#$%^&*()_-+=(){}[]|\\:;\"<>,.?/'"
11   STRIPLIST      " ~`!@$^*_+(){}[]|\\:;<>?"
12   \PRAGMA_END
13
14
15   \POST_START
16   \POST_END
17
18
19   ;---------------------------------------------------
20
21
22   ; &                                              ;
23   **                                               ;
24   COPY        "null"    {colA}                      ;
25   COPY        "null"    {colB}                      ;
26   COPY        "null"    {colC}                      ;
27   COPY        "null"    {colD}                      ;
28   COPY        "N"       vAnyUnhandled               ;
29
30
31   & | \= | ** | \,                                 ;
32   CALL MYSUB_MID                                    ;
33   REPEAT                                            ;
34
35
36   & | \= | **                                      ;
37   CALL MYSUB_LAST                                   ;
38
39
40   ** | [ vAnyUnhandled = "N" ]                      ;
41   COPY        "null"    {UnhandledData}             ;
42
43
44   ;---------------------------------------------------
45
```

*Figure 24-4   Pattern Action file, page 1 of 3.*

```
45    ;
46
47    \SUB MYSUB_MID
48
49
50    T = "colA" | \= | ** | \,                              ;
51    COPY_S       [3]         {colA}                         ;
52    RETYPE       [1]         0                              ;
53    RETYPE       [2]         0                              ;
54    RETYPE       [3]         0                              ;
55    RETYPE       [4]         0                              ;
56    RETURN                                                  ;
57
58    T = "colB" | \= | ** | \,                              ;
59    COPY_S       [3]         {colB}                         ;
60    RETYPE       [1]         0                              ;
61    RETYPE       [2]         0                              ;
62    RETYPE       [3]         0                              ;
63    RETYPE       [4]         0                              ;
64    RETURN                                                  ;
65
66    T = "colC" | \= | ** | \,                              ;
67    COPY_S       [3]         {colC}                         ;
68    RETYPE       [1]         0                              ;
69    RETYPE       [2]         0                              ;
70    RETYPE       [3]         0                              ;
71    RETYPE       [4]         0                              ;
72    RETURN                                                  ;
73
74    T = "colD" | \= | ** | \,                              ;
75    COPY_S       [3]         {colD}                         ;
76    RETYPE       [1]         0                              ;
77    RETYPE       [2]         0                              ;
78    RETYPE       [3]         0                              ;
79    RETYPE       [4]         0                              ;
80    RETURN                                                  ;
81
82
83    ! T | \= | ** | \, | [ vAnyUnhandled = "N" ] ;
84    COPY         "Y"         vAnyUnhandled                  ;
85    COPY         [1]         vTemp1                         ;
86    CONCAT       "="         vTemp1                         ;
87    COPY_S       [3]         vTemp2                         ;
88    CONCAT       vTemp1      {UnhandledData}                ;
89    CONCAT       vTemp2      {UnhandledData}                ;
90    RETYPE       [1]         0                              ;
91    RETYPE       [2]         0                              ;
92    RETYPE       [3]         0                              ;
93    RETYPE       [4]         0                              ;
94    RETURN                                                  ;
95
96    ! T | \= | ** | \, | [ vAnyUnhandled = "Y" ] ;
97    COPY         [1]         vTemp1                         ;
98    CONCAT       "="         vTemp1                         ;
99    COPY_S       [3]         vTemp2                         ;
100   CONCAT       ","         {UnhandledData}                ;
101   CONCAT       vTemp1      {UnhandledData}                ;
102   CONCAT       vTemp2      {UnhandledData}                ;
103   RETYPE       [1]         0                              ;
104   RETYPE       [2]         0                              ;
105   RETYPE       [3]         0                              ;
106   RETYPE       [4]         0                              ;
107   RETURN                                                  ;
108
109
110   \END_SUB
111
```

*Figure 24-5   Pattern Action file, page 2 of 3.*

```
111
112
113
114
115     \SUB MYSUB_LAST
116
117
118     T = "colA" ¦ \= ¦ **                           ;
119     COPY_S        [3]        {colA}                ;
120     RETYPE        [1]        0                     ;
121     RETYPE        [2]        0                     ;
122     RETYPE        [3]        0                     ;
123     RETURN                                         ;
124
125     T = "colB" ¦ \= ¦ **                           ;

126     COPY_S        [3]        {colB}                ;
127     RETYPE        [1]        0                     ;
128     RETYPE        [2]        0                     ;
129     RETYPE        [3]        0                     ;
130     RETURN                                         ;
131
132     T = "colC" ¦ \= ¦ **                           ;
133     COPY_S        [3]        {colC}                ;
134     RETYPE        [1]        0                     ;
135     RETYPE        [2]        0                     ;

136     RETYPE        [3]        0                     ;
137     RETURN                                         ;
138
139     T = "colD" ¦ \= ¦ **                           ;
140     COPY_S        [3]        {colD}                ;
141     RETYPE        [1]        0                     ;
142     RETYPE        [2]        0                     ;
143     RETYPE        [3]        0                     ;
144     RETURN                                         ;
145
146
147     ! T ¦ \= ¦ ** ¦ [ vAnyUnhandled = "N" ]        ;
148     COPY          "Y"        vAnyUnhandled         ;
149     COPY          [1]        vTemp1                ;
150     CONCAT        "="        vTemp1                ;
151     COPY_S        [3]        vTemp2                ;
152     CONCAT        vTemp1     {UnhandledData}       ;
153     CONCAT        vTemp2     {UnhandledData}       ;
154     RETYPE        [1]        0                     ;
155     RETYPE        [2]        0                     ;
156     RETYPE        [3]        0                     ;
157     RETURN                                         ;
158
159     ! T ¦ \= ¦ ** ¦ [ vAnyUnhandled = "Y" ]        ;
160     COPY          [1]        vTemp1                ;
161     CONCAT        "="        vTemp1                ;
162     COPY_S        [3]        vTemp2                ;
163     CONCAT        ","        {UnhandledData}       ;
164     CONCAT        vTemp1     {UnhandledData}       ;
165     CONCAT        vTemp2     {UnhandledData}       ;
166     RETYPE        [1]        0                     ;
167     RETYPE        [2]        0                     ;
168     RETYPE        [3]        0                     ;
169     RETYPE        [4]        0                     ;
170     RETURN                                         ;
171
172
173     \END_SUB
174
175
176
```

:set nu

176:176:1  74x60        Command

*Figure 24-6   Pattern Action File, page 3 of 3.*

[placeholder]

A code review for Figure 24-4, Figure 24-5, and Figure 24-6 includes;

– Lines 10-11, specify that we field separate (token separate) on commas and the equal symbol, but do not strip them from the patterns that we evaluate later.

This is key, because our basic pattern is,

ColumnName=Column value of 1 or more words, [comma]

We will use the hard presence of the equal symbol, and then comma to delimit our unknown and multi-word value for (column value).

We also separate and strip spaces, to avoid receiving extra spaces in our outputted values. Stripping spaces, but then also copying with spaces (COPY_S) allows for this to occur.

**Note:** Not certain what we are discussing here? As a test after you complete this example, call to strip spaces from the SEPLIST and/or STRIPLIST, and see what your outputted data looks like.

– Lines 15-16 are included for discussion, but are non-functional in this example.

These lines allow for an after record processing capability, each individually received record after another. In other words, if you wanted to check what you produced for each row, before actually outputting data, you would put code between Line 15 and Line 16. If you have 99 input records, code contained between Line 15 and Line 16 would execute 99 times.

However, this is a special code construct within QualityStage, and may only be used for features like SOUNDEX, NYSIIS, and related.

You may have your own custom after row processing block, but you do not accomplish this need via Line 15 and Line 16; you accomplish it via a technique similar to Line 40, discussed below.

– Lines 22-28 form our before record processing block.

In effect, Lines 22-28 are executed for/before every new input record we receive.

Line 22 is actually commented out, but represents the common use case. The "&" simple pattern class matches any single token, and would thus cause the lines that follow to execute for every new input record we receive.

We present an alternate solution in Line 23, with the simple pattern class of "**". The simple pattern class "**", matches zero or more tokens, in the

event we wished to also process entirely NULL input records, a topic which is not pursued further in this document. ("&" does not match entirely NULL input records, "**" does.)

While Line 23 is our actual pattern line, Lines 24-28 are the commands that execute when this pattern is matched.

**Note:** Recall that this pattern is matched on every new input record. Because this pattern is also listed first, it serves as our *before record pattern block*.

Huh? If you want to perform any kind of processing before each input row that is received, this is the technique that you use. Here we use this technique to set a default value for our output columns. Only if these output columns are overwritten later as part of normal processing will their default values be changed.

Lines 24-27 copy the hard coded value of "null" to our output columns; colA, colB, colC, and colD.

Line 28 copies the value of "null" to a user defined variable with the name, vAnyUnhandledData. vAnyUnhandledData is not a keyword; we could of just as easily called this variable, Bob.

– Lines 31-33, and then Lines 36-37 are where the main processing is performed.

Recall that the input data is in the general form,

colA=Word word,colB=Word 22,colC=More words still,colD=word

Line 31 will match the first 3 column name, column values pairs in the input data stream. A code review of Line 31 follows;

& is the built in simple pattern class for any single token, and will match colA, colB, or any unknown column name (value).

\= is the simple pattern class to match the equal symbol.

** is the simple pattern class to match zero or more (words).

\, is the simple pattern class to match the comma.

Line 36 differs from Line 31, in that Line 36 does not require the trailing comma, as is the case with the last column name, column value pair on any input record.

Line 32 calls to invoke a user defined pattern action language sub-routine with the name, MYSUB_MID.

Line 33 (REPEAT) will cause the sub-routine invoked on Line 32, to repeat until it finds no more pattern matches.

**Note:** In effect, this sub-routine (MYSUB_MID) will nibble off (look for, find, and then remove matching tokens off) of the leading portion of the input record, until there are no more pattern matches.

This is because this sub-routine invocation is followed by the REPEAT command. If the REPEAT command were not present on Line 33, then this sub-routine (MYSUB_MID) would only invoke one time, as is the case with the sub-routine invocation on Line 36 (MYSUB_LAST).

What these two user defined sub-routines do, is determined by their contents, listed below.

We say l*eading portion* above, because the default pattern matching process reads the input string from left to right. If or when necessary, you may choose to read input string left to right, right to left, or a mixture of both. You may even choose to just look for specific tokens found anywhere inside the larger input string.

– Lines 40-41 form our (user defined and configurable) after record processing block.

In effect, Lines 40-41 are executed for/after every new input record we receive.

Line 40 also offers a new and specific use case. While the "**" simple pattern matches every row, the next portion of this line displays how we can evaluate user defined variables.

[ vAnyUnhandled = "N" ]

Will be true of the user defined variable, vAnyUnhandled, is equal to "N". This clause is treated the same as any simple or user defined pattern class; if true, run the pattern action commands that follow, otherwise do not.

In effect, we use the block to set the output column UnHandledData to equal a default value, just as we did for output columns, colA, colB, etcetera. Logically we could have set and maintained the value in

UnHandledData in the same manner as colA, colB, etcetera, but this provided us an example to demonstrate after row processing.

– Lines 47-110 provide the contents of the user defined sub-routine with the name, MYSUB_MID.

The contents of any user defined sub-routine follow the same basic rules as the normal content areas of the Pattern Action file, with pattern match lines, and pattern action commands. User defined sub-routines may even call other user defined sub-routines.

The MYSUB_MID user defined sub-routine contains 6 (count) pattern match lines, located on; Lines 50, 58, 66, 74, 83 and 96. The pattern match blocks beginning on Lines 50, 58, 66 and 74, process a *known* input column names ("T"). The pattern match blocks beginning on Lines 83 and 96, process previously unknown column names("! T", not T).

The blocks starting on Lines 50, 58, 66, and 74 are similar to one another. For example; Line 51 copies the column value to a given/correct output column, and Lines 52-55 remove the received/identified tokens from further consideration. (These tokens are removed because we have fully processed them. We're done with them.)

**Note:** The RETURN command on Line 56, along with the REPEAT of the sub-routine on Line 33, is what forms our looping construct, allowing us to walk the larger input string.

The blocks starting on Lines 83 and 96 are similar to one another.

In effect, the block starting on Line 83 is executed upon receiving the first unknown column name. The block starting on Line 96 executes for any second and subsequent unknown column name. The only difference between these two blocks is how we format the output column, UnHandledData, which we append to consecutively in a very precise and well formatted manner.

Line 84 changes the values of the user defined variable named, vAnyUnhandled to "Y".

Lines 85-89 are just so that we may precisely control spacing and formatting of the value(s) in our output column.

– Lines 115-173 provide the contents of the user defined sub-routine with the name, MYSUB_LAST.

The user defined sub-routines MYSUB_MID and MYSUB_LAST are *very similar*. (The only difference is the presence of the check for the trailing and field delimiting comma.) We offer the current solution as is, preferring to leave a support legacy of a simpler and cleaner example, versus a more complex, exotic, and possibly harder to understand example for those who follow us.

Because this sub-routine nearly mirrors the preceding one, no further code review is provided.

6. Save, compile, and test your work. A successful result appears similar to that as displayed in Figure 24-7.

> **Note:** When creating and working with QualityStage custom rules sets, remember to save your work often.
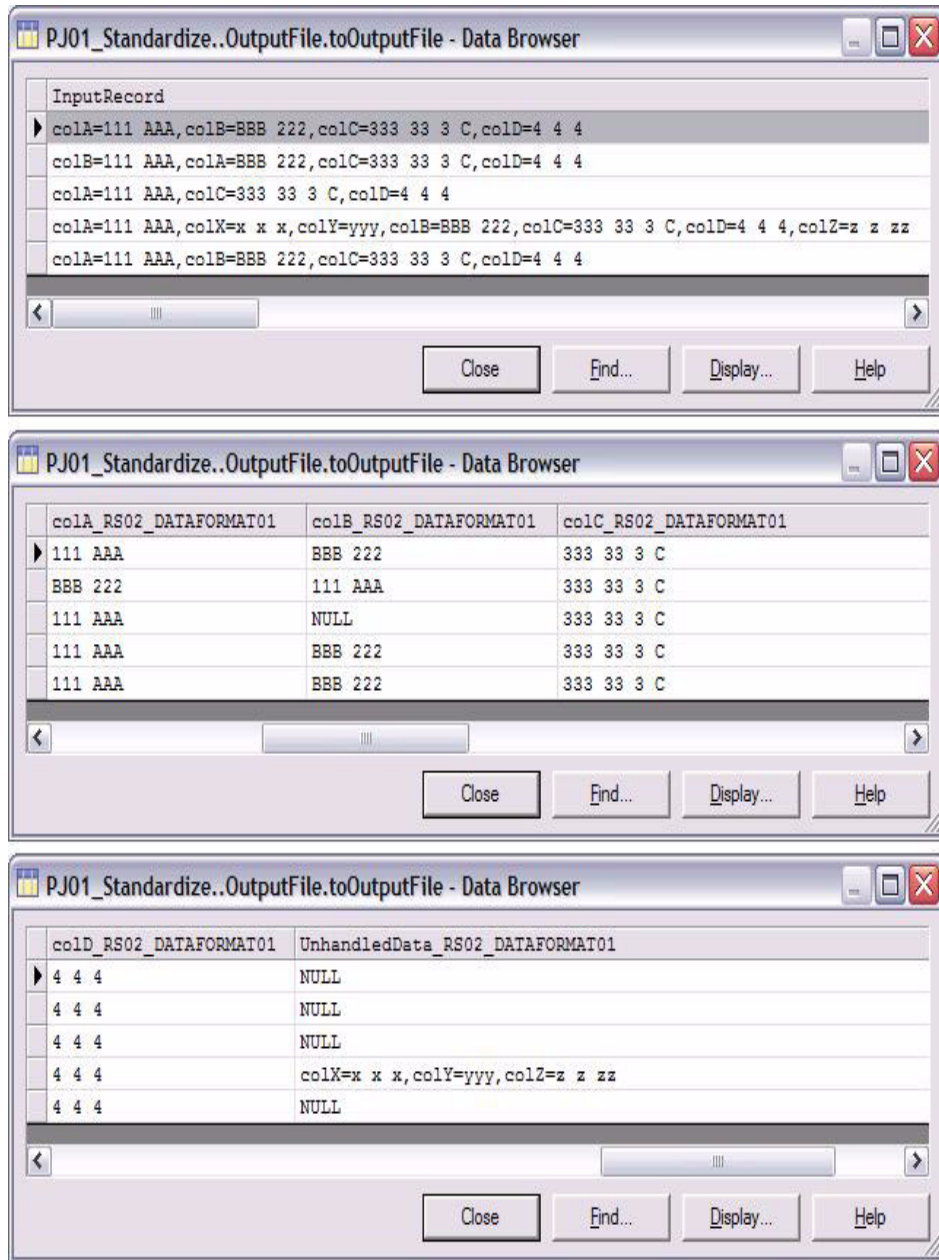
*Figure 24-7   Successful result from test.*

# 24.3  In this document, we reviewed or created:

We examined the (IBM) InfoSphere Information Server (IIS), QualityStage component, Standardize stage. Specifically, we created and tested one more advanced QualityStage custom rule set; a very powerful technology with nearly endless applications to standardize and cleanse data of any source or purpose.

More specifically, we demonstrated looping, before and after row processing, and handling NULL and (variably user defined) UnHandled Data.

**Persons who help this month.**

Andy Wilson, and Robert Dickson.

### Additional resources:

The IBM InfoSphere Information Server, QualityStage component, product tutorial.

### Legal statements:

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both. These and other IBM trademarked terms are marked on their first occurrence in this information with the appropriate symbol (® or $^{TM}$ ), indicating trademarks that were owned by IBM at the time this information was published. A complete and current list of IBM trademarks is available on the Web at http://www.ibm.com/legal/copytrade.shtml

Other company, product or service names may be trademarks or service marks of others.

### Special attributions:

The listed trademarks of the following companies require marking and attribution:

Adobe, the Adobe logo, PostScript, and the PostScript logo are either registered trademarks or trademarks of Adobe Systems Incorporated in the United States, and/or other countries.

Cell Broadband Engine is a trademark of Sony Computer Entertainment, Inc. in the United States, other countries, or both and is used under license therefrom.

Java and all Java-based trademarks are trademarks of Sun Microsystems, Inc. in the United States, other countries, or both.

Microsoft, Windows, Windows NT, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

Microsoft trademark guidelines

Intel, Intel logo, Intel Inside, Intel Inside logo, Intel Centrino, Intel Centrino logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium, and Pentium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

Intel trademark information

UNIX is a registered trademark of The Open Group in the United States and other countries.

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both.

ITIL is a registered trademark, and a registered community trademark of the Office of Government Commerce, and is registered in the U.S. Patent and Trademark Office

IT Infrastructure Library is a registered trademark of the Central Computer and Telecommunications Agency, which is now part of the Office of Government Commerce.

Other company, product, or service names may be trademarks or service marks of others.