The Data Profiler tests for valid and invalid values? Sentence Style OK for chapter title?

- The Data Profiler tests
- How testing results are stored in the EMES Not in the doc. Values that match more than one test Important assumptions about validity
- Invalidity tests
- Another definition of invalid
- How a value can be both valid and invalid

Nat in the doc.

could be "Both valid and invalid" If a chapter combines conceptual and reference information; typically, the reference information is placed at the end of the chapter. or is that a completely But, not a head!

Suggest a new organization: · About the Water Profiler (see suggested topic)
· Important assumptions about validity

- . Invalidity tests
- Validity tests (to be written)
- Values that match more than one test Deta Profiler test descriptions (changed

from "The Data Profiler tests")

The tests are about the Physical Element Value.

Does the Wata Profiler check only for the validity/involidity of the Physical, Element value? The stitle should reflect the content. As written, the title suggests that the Data Profiler checks more than the validity/invalidity at more than one type of value.

Does a value motch a test, or is the Data Profiler testing a Value against a test and then returns the test name The Data Profiler tests Consider: This list describes the Data Profiler validity and invalidity tests, in the order they are performed. The Deata Profiler -> Note that the invalidity tests precede the validity tests. tests for invalid and If a value matches one - and only one, not more of these tests, the test name (Null, Invalid Valid values. The value, Invalid character, etc.) appears in the profile results as the REASON in the Data Quality tests described are in Summary, and as the DESCRIPTION in the Common Values reports. bold and If a value matches multiple tests, it appears in the profile results with one of the descriptions the order that they are Tisted in "Values that match more than one test" Possible to move the performed. Tests for invalid values Test: NULL C suggest a readi period inside the quotation mark a suggest a reading What is the test? If the record format defines NULL, the Data Profiler compares the Physical Element value to the ·· NULL value. If the value is NULL, that's how it's treated. For more information, see "Interpreting a wat in the value as NULL" - contractions are being used necessary? here, but not elsewhere Test: Invalid value These are all Dela Profiler tests. Data profiler does a comparison between the Physical Element value and the Validation Speq invalid values If it matches and invalid value, it is treated as invalid. Physical Element Test: Invalid value The Physical Element value is compared against the character set specified in the Validation Spec. If the value contains any characters that are not listed in Value Must Only Use These Characters, it is treated as invalid. compared against the values in matches any of the values tooked up in the invalid looking file is this value Test: Invalid lookup The value is looked up in the invalid lookup file. If the value is found there, it is treated as invalid. Plysical Element! Test: Invalid expression compared I'm not sure The Physical Element value is tested against each of the invalid expressions in order. If the invalid What's being expression returns NULL or an error, the value is treated like it did not trigger the expression. If invalid expressions return a non-zero result, the value is treated as invalid. 1 is the index of the tested - the first valid expression in the list. Physical Element value How is the expression triggered? Is the test comparing the Physical Element value to a list of invalid or expressions? Not sure what 4 hus has to do bloes the Physical with the test. expressione? Element value trigger explessions?

Tests for valid values & suggest a heading

Test: Ignore this value when computing statistics

The Physical Element value is compared against the list of valid values in the Validation Spec. If Physical Element value is considered valid but is not included in the calculation of the Histogram, Total values, Valid Values, Invalid values. Mean value Distinct Values Union of the Histogram, Value, Maximum Value, Empty Values, Blank Nalues, Normal Values, Ascending Pairs, Descending Pairs, or Cross-Field Relationships. This setting is in the Add valid values dialog

<u>Value:</u>			Base type:
	lue as a DML litera when computing		
203CH DECOM	ALCOHOLOGICA (CERTIFICACIÓN MATERIAL PARAMETER)	Control of the Contro	
	OK	Cancel	Help

Test: is valid There's a test called is-invalid and a function with the same?

The Physical Element value is tested by calling the the DML function is_valid. If is_valid returns 0 t what calls the function bold OK? (False), the value is treated as invalid.

NOTE: If the Physical Element DML definition contains the is_valid_XXX field valdity function, the is_valid_XXX function is currently ignored.

"field validity

Not sure why this information is prelevant. The Physical Element DML function is called . What does the DML

Test: Valid value

definition have to do with the function? compared ogsi The Physical Element value is looked at in comparison to the valid values in the Validation Spec. If the value matches any of these valid values, the value is treated as valid.

Physical Element

Test: Valid range

against

The Physical Element value is compared to the valid ranges in the Validation Spec. If the value falls within any of the ranges (including if it matches either endpoint of any range), the value is treated as valid.

Test: Valid pattern

against

The Physical Element value is compared to the valid patterns in the Validation Spec. If there's any match, it's treated as valid.

Consider: If the Physical Element value matches a valid pattern in the Validation Spee, the Physical Element Value is valid.
The Data Profiler tests for valid and invalid values

Test: Valid expression
The Physical ry
the fire
The physical ry
the physi The Physical Element value is tested against each of the valid expressions in order 1 is the index of the first expression in the lat. Values for expressions that return NULL or an error are treated like they did not trigger the expression. If expressions return a non-zero result, the value is treated as

Test: Valid lookup Eleme matches a value The value is looked up in the valid lookup file. If the value found in the valid lookup file, it is treated as valid.

Values that match more than one test This is conceptual

When the Data Profiler tests have run on all the values in the dataset for a given Physical Element, the Data Profiler analyzes the results and decides how to group the statistics for values that matched more than one test. For matching values, the DESCRIPTION value in the Data Quality Summary of the REASON value in the Common Value reports may be one of the following (in addition to those sited in the Data Profiler tests)!

Valid <

Both valid and invalid V

Multiple reasons for invalidity

Multiple reasons for validity

Fails all definitions of validity

A description of valid mean that the Physical Element value did not match any invalidity test, and there are no valid values, ranges, patterns, expressions, or lookups defined. invalid?

Both valid and invalid

verified? One value of the Physical Element being profiled is valid in one record while the same value is invalid in another record. This can happen when one or more validity or invalidity expressions refer to a Physical Element in the record other than the Physical Element being tested. while

The result can be one of these:

- One value in one record causes a validity expression to result in True while the same value in 0another record causes a validity expression to result in False.
- One value in one record causes an invalidity expression to result in True while the same value in another record causes an invalidity expression to result in False.

and Common John Valid

- One value in one record causes a validity expression to result in True while the same value in another record causes an invalidity expression to result in True()
- One value in one record causes an invalidity expression to result in True while the same value in another record causes an invalidity expression to result in True

Not in the Could it be valid and invalid For more information, see "How a value can be both valid and invalid"

Multiple reasons for validity

Multiple instances of a Physical Element value triggers more than one test for validiity. This can happen when multiple validity expressions refer to fields in teh record other than the field being tested.

The result can be that a value is valid in one record because it triggers the first validity expression in one record but triggers the second validity expression in another record. that?

-Multiple reasons for invalidity

The meaning of multiple reason for invalidity depends on the setting of Reasons for Invalidity when the profile. For more information, see "The setting of Reasons for invalidity"

Lif Report first reason only was selected when when running a job, multiple reasons for invalidity means the value is invalid for one reason in one or more records, and invalid for a different reason in one or more other records. This can occur when an invalidity expression uses the value of another Physical Element in the record. For example, a value might trigger the third invalidity expression in one record but trigger the fourth invalidity expression in another record.

If Report all reasons was selected when running a job, multiple reasons for invalidity can means the same thing as was described in the preceding paragraph (Or, it can mean that the testing of a Physical Element value in one record had at least two of the following results:

The value matches an invalid value

- The value contained an invalid character
- The value is in a lookup file of invalid values
- The value triggered an expression that indicates invalidity.
- The value is invalid for the data type defined in the dataset's record format
- The value fails all definitions of validity()

Fails all definitions of validity Invalid?

Valid values, valid ranges, valid patterns, valid expressions, and/or valid lookups are defined, however) the value does not match any of them (And, the value passes the is_valid test and is not flagged Ignore this Value when Computing Results.

should be capped, but
should be capped, but
may be this way in the UI.
The Data Profiler tests

The Data Profiler tests

The Data Profiler tests for valid and invalid values

AB INITIO CONFIDENTIAL AND PROPRIETARY - DO NOT COPY

Ignore this value when computing statistics?

Important assumptions about validity

I would avoid meluding

would avoid man. There are four important assumptions the Data Profiler makes about validity:

The number of items.

If a value matches both a valid and an invalid criteria.

Validation of Validatio If a value matches both a valid and an invalid criteria due to contradictory specification in the Validation Spec, the value is considered as to be invalid

If all four sections of the Valids tab of the Validation Spec are blank, and no valid lookup file is specified, every value is assumed to be valid, unless it is evaligible in a list. Considered? of the following:

On the Invalids tab of the Validation Spec, the value is in the list of Invalid Values, the value triggers at least one of the Expressions That Trigger Invalidity, or contains a character not listed in Values Must Use Only These Characters.

The value exists in a lookup table of invalid values.

The value is illegal for the record format.

If there is any entry in any section of the Valids tab, a value is treated as valid only if it matches one or more of the entries on the Valids tab (and is not explicitly invalid).

the value

Clamard values

Just mersoned

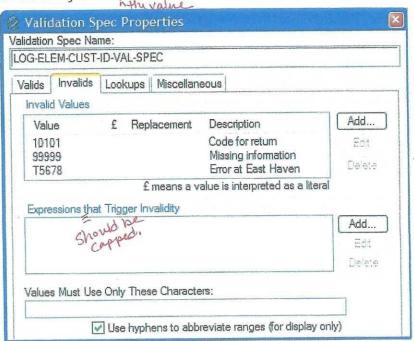
Just his topics

Invalidity tests

This topic presents a simplified look at the tests that cause a value to be treated as invalid. They are the same tests as those listed in "The Data Profiler tests", but without the details of the validity tests. Use this topic as a quick reference. In more information about the tests, see "The Wata Profiler tests.

are These are the Data Profiler invalidity tests listed in the order they are performed:

1. The value is compared against the invalid Values on the Invalid tab of the Validation Spec. If it matches any of these, it is treated as invalid.



make more plant to me before the

2. If the value contains any characters not listed in Values Must Only Use These Characters on the Invalids tab, the value is treated as invalid.

3. The value is looked up in the lookup file specified under Invalid Values Lookup on the Lookups

4. The value is tested in order against each of the expressions in the Expressions That Trigger Invalidity section of the Invalid tab. The value is treated as invalid if at least one expression returns a non-zero result. If the expression returns NULL or an error (it is treated as if the value did not trigger the expression.

5. The value is tested by calling the DML function is valid. If is valid returns 0 (False), the value is treated as invalid. What is calling the function?

If any valid value, valid range, valid pattern, valid expression, or valid lookup is defined but the value did not trigger any of those tests, the value is treated as invalid.

which tests? For valid value etc?

Alertion of invalid

Here is a definition of invalid that different stand

"Invalidity tests" and "The Tunderstand"

"The Tunderstand"

Here is a definition of invalid that differs in presentation (not in content) from those presented in "Invalidity tests" and "The Data Profiler tests". Use whichever definition is the easiest to

The Wata Profeler tests values for validity-x The conditions are listed in the same order that the Data Profile considers them.

A value is invalid if at least one of the following are true:

The value is in the list of invalid values.

Contains any characters not listed in the set of valid characters (but only if a set of valid characters is specified).

Watches a Key of a lookup table of invalid values.

Any of the invalidity expressions evaluates to True for that value.

The value is not valid according to the record format (specifically if the is_valid DML function returns False), and the value is not flagged Ignore in statistics in the object's Validation Spec.

The value is not listed as valid (where valid means that at least one of the following is True):

The value is in the list of valid values.

Falls within one of the valid ranges.

The value's It's pattern is in the list of valid patterns. Ignore this value when computing statistics?

Any of the validity expressions evaluate to True for that value.

The value matches a Key of a lookup table of valid values.

There is no explicit definition of valid (no valid values, no valid ranges, no valid patterns, no expression to indicate validity, and no lookup table specified for identifying valid values).

NOTE: The following condition makes a value neither valid nor invalid: the record format specifies NULL for a particular value, and the value found during profiling is equal to that value. See "interpreting a value the testing?

This list confuses as NUI

This list confuses as NUI

The Why list ma when when which the point is to discuss the point is the invalid?

The point is the invalid?

Conditions to be invalid?

Any value or Element

the value is 1 not invalid or not invalid. Not the valid. Im x?

About the Data Profiler

The Data Profiler performs tests to determine the invalidity and validity of a Physical Element value.

If a value matches one of the lests, the test name appears in the profile results. If a value metches multiple tests, the test name appears in the multiple tests, the test name appears in the profile results with one of the descriptions in profile results with one of the descriptions in "Values that match more than one test."

This is my suggestion for a new topic. I'm not sure if it's technically accurate. I took the information from the introduction in "The Clata Profiler tests." If the new topic is included, I would delete much of the information in the introduction for "The Clata Prefiler information in the introduction for "The Clata Prefiler tests" topic and replace it with a simple sentence that introduces the tests. I would also move the specifics about the information in the profile specifics about the information in the profile and to "Validity tests" (if written and the profile and to "Validity tests" (if written and the profile and to "Validity tests" (if written and the profile