



**ataccama**

# ONE Desktop Workshop

## Scoring & Explanations

Prepared for: v15.4.x

Prepared by: Ataccama

Dated: October 2024



## Contents of this Document

1. Introduction	3
2. Tasks	3
2.1. Create a plan and add a data source	3
2.2. Add Alter Format and Simple Scoring	3
2.3. Configure the Simple Scoring step	4
3. Conclusion	4
Correct answers, hints, and useful tips	5

# 1. Introduction

Assigning a score value for a record can be one of the ways to distinguish good records from bad ones. The scoring element can be imagined as 'penalty points' assigned for failing certain rules, validity checks, or given conditions.

This workshop will explain data quality by using the **Simple scoring** step, where business rules and corresponding scores are defined.

## 2. Tasks

Let's revisit our business rules that were defined in our Profiling workshop (**ONE Desktop Workshop – Data Profiling**). We now have some slightly updated rules:

- **src\_sin** should be exactly 9 digits, should be numbers only, and cannot start with a zero.
- **meta\_last\_update** must not be in the future.

We will expand on this and define these scoring rules:

- For all attributes, if they pass the business rule, give a score of 0 and an explanation of "SIN\_VALID" or "META\_LAST\_UPDATE\_VALID" as appropriate.
- if **src\_sin** is null, give a score of 10,000,000 and the explanation "SIN\_NULL."
- If **src\_sin** is longer than 9 characters, but after removing non-numbers is 9 characters, give a score of 200 and the explanation "SIN\_CHARS."
- if **src\_sin** after converting to a number is less than 9 digits, give a score of 100,000 and the explanation "SIN\_TOO\_SHORT". (This includes cases with beginning zeros)
- If **meta\_last\_update** is in the future, give a score of 9,000,000 and the explanation "META\_LAST\_UPDATE\_IN\_FUTURE."

### 2.1. Create a plan and add a data source.

- › In the **plans** folder, create a new plan called **04\_party\_scoring.plan**
- › Add **party\_full\_1.csv** into the plan.

## 2.2. Add Alter Format and Simple Scoring

- › Find the Alter **format** step and add it to the plan.
- › Connect the '**out**' endpoint of the **Text File Reader** to the '**in**' endpoint of the **Alter Format** step.
- › Open the **Alter Format** step and add 2 new columns: **sco\_default** (INTEGER) and **exp\_default** (STRING). Leave the expressions blank.
- › Add a **Simple Scoring** step and connect it to the data flow

## 2.3. Configure the Simple Scoring step.

After having gone through the previous workshops, you should be able to work out how to fill in the **Simple Scoring** step.



*Have a try and work out how to fill in the **Simple Scoring** step based on the requirements at the beginning of the chapter.*

*Consult the previous workshops for tips on expressions.*

*Don't forget the **CTRL+Space** content assist and the F1 help on expressions.*

*You can also use the expression debug function and templates. Good luck!*

When finished, end the flow with a text file writer, name the file **party\_full\_1\_score.csv**, and store it on **data \ out**. Open the file and check the results.

## 3. Conclusion

We have come to the end of this workshop, where we have performed some simple scoring on our data. The lower the score, the better the quality of the data.



*Looking at the data, can you think of more data quality checks to be put into place?*

## Correct answers, hints, and useful tips

Here is how the **Simple Scoring** step should be configured:

Simple Scoring

General

Advanced

Other

Default Score Column:

sco\_default

Default Explain Column:

exp\_default

Scoring Cases:

	Condition	Description	Explanation	Score
1	src_sin is null	sin is null	"SIN_NULL"	100000000
2	length(src_sin) > 9 and length(tostring(tointeger(trashNonDigits(src_sin)))) = 9		"SIN_CHARS"	200
3	length(tostring(tointeger(trashNonDigits(src_sin)))) < 9		"SIN_TOO_SHORT"	100000
4	length(tostring(tointeger(trashNonDigits(src_sin)))) = 9		"SIN_VALID"	0
5	today(meta_last_update, 'yyyy-MM-dd') > today()		"META_UPDATE_IN_FUTURE"	9000000
6	today(meta_last_update, 'yyyy-MM-dd') <= today()		"META_LAST_UPDATE_VALID"	0
*				

Add

Fill Columns...