

Projet Data Integration

16.01.2021

Réalisé par: BELLARABI Ibtihal

Vue d'ensemble

Parmi les étapes les plus importantes du business Intelligence est la collecte et l'alimentation des données en utilisant les outils ETL (**EXTRACTION- TRANSFORMATION- LOAD**).

Il consiste à collecter les infos nécessaires depuis les différentes sources **extraites, transformées et chargées** dans un entrepôt de données.

Objectifs

1. Se familiariser avec les tâches d'intégration des données
2. nettoyage des données et centralisation dans une datawarehouse

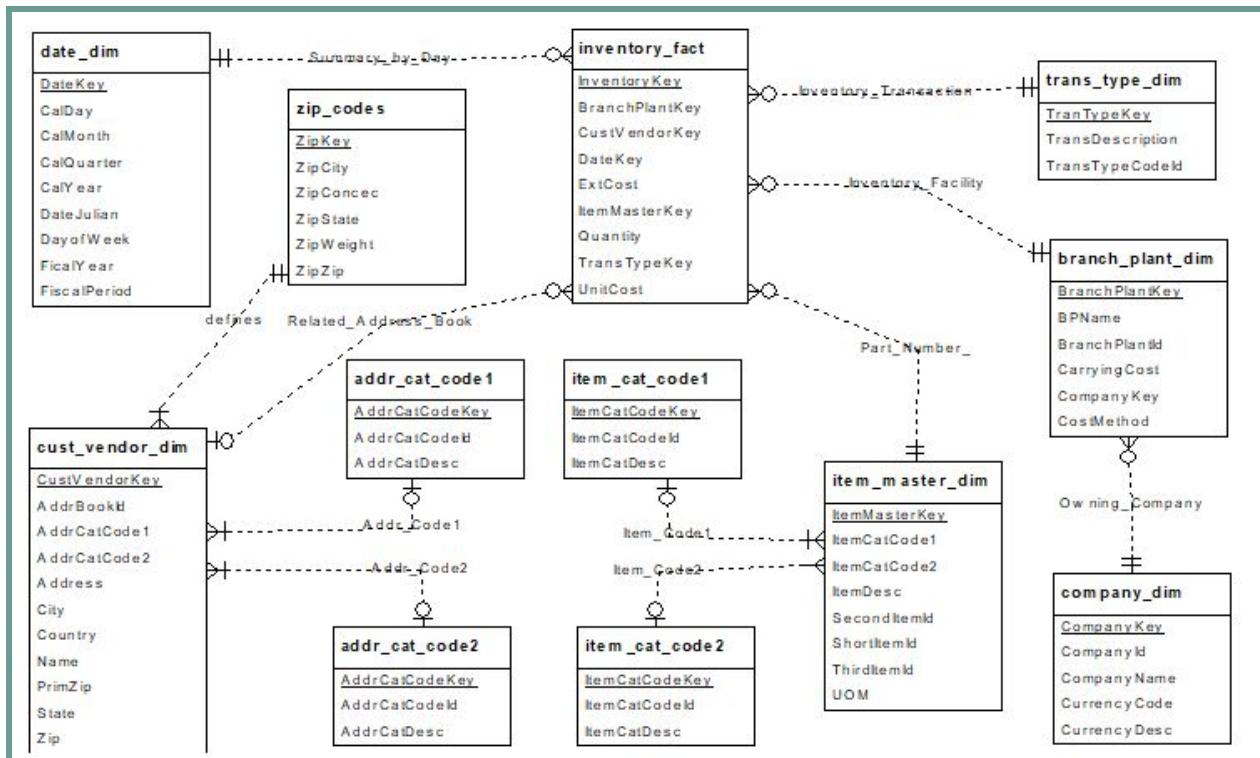
Outils

Pentaho Data Integration (interface graphique Spoon)

Grandes étapes

I. Cration de la data warehouse (de la table Inventory_Fact)

Snowflake Schema



Query 1 Inventory_Fact

```
1
2 •  create database if not exists inventory_fact;
3 •  use inventory_fact;
4
5 •  drop table if exists inventory_fact ;
6 •  drop table if exists branch_plant_dim ;
7 •  drop table if exists cust_vendor_dim ;
8 •  drop table if exists item_master_dim ;
9 •  drop table if exists addr_cat_code1 ;
10 •  drop table if exists addr_cat_code2 ;
11 •  drop table if exists item_cat_code1 ;
12 •  drop table if exists item_cat_code2 ;
13 •  drop table if exists company_dim ;
14 •  drop table if exists zip_codes ;
15 •  drop table if exists date_dim ;
16 •  drop table if exists trans_type_dim ;
17
18 •  drop table if exists Currency_Dim ;
19
20 •  CREATE TABLE Currency_Dim (Currency_ID VarChar(3),Exchange_Rate Decimal(8,2));
```

Output :

#	Time	Action	Message
5708	05:17:43	Insert into INVENTORY_FACT (INVENTORYKEY,BRANCHPLANTKEY,DATEKEY,ITEMMASTERKEY,TRANS...	1 row(s) affected
5710	05:17:43	Insert into INVENTORY_FACT (INVENTORYKEY,BRANCHPLANTKEY,DATEKEY,ITEMMASTERKEY,TRANS...	1 row(s) affected
5712	05:17:43	Insert into INVENTORY_FACT (INVENTORYKEY,BRANCHPLANTKEY,DATEKEY,ITEMMASTERKEY,TRANS...	1 row(s) affected
5715	05:17:43	Insert into INVENTORY_FACT (INVENTORYKEY,BRANCHPLANTKEY,DATEKEY,ITEMMASTERKEY,TRANS...	1 row(s) affected
5717	05:17:43	Insert into INVENTORY_FACT (INVENTORYKEY,BRANCHPLANTKEY,DATEKEY,ITEMMASTERKEY,TRANS...	1 row(s) affected
5719	05:17:43	Insert into INVENTORY_FACT (INVENTORYKEY,BRANCHPLANTKEY,DATEKEY,ITEMMASTERKEY,TRANS...	Running...

II. ETL Opérations

1. Data Source

we need to perform cleaning operations on two data sources. Both data sources provide facts for the *Inventory_Fact* table of the inventory data warehouse. We need to create a transform for each data source.

we create the **Microsoft Excel input**

The dialog box has tabs at the top: Files, Sheets, Content, Error Handling, Fields, and Additional output fields. The 'Files' tab is selected. In the 'Spread sheet type (engine)' dropdown, 'Excel 97-2003 XLS (JXL)' is chosen. The 'File or directory' field contains 'C:\Users\hp\Desktop\Projet\ExcellInventory.xls'. There are buttons for 'Add' and 'Browse...'. Below this are fields for 'Regular Expression', 'Exclude Regular Expression', and 'Password'. A table titled 'Selected files:' lists one item: '# 1 File/Directory C:\Users\hp\Desktop\Projet\ExcellInventory.xls'. Columns include 'File/Directory', 'Wildcard (RegExp)', 'Exclude wildcard', and 'Requires password'. Buttons for 'Delete' and 'Edit' are on the right. At the bottom, there are sections for 'Accept filenames from previous steps' (checkboxes for 'Accept filenames from previous step' and 'Step to read filenames from'), a 'Field in the input to use as filename' dropdown, and a 'Show filename(s)...' button. At the very bottom are 'OK', 'Preview rows', and 'Cancel' buttons.


 Microsoft Excel input
Step name **Inventory_EXCEL**

Files Sheets Content Error Handling Fields Additional output fields

#	Name	Type	Length	Precision	Trim type	Repeat	Format	Currency
1	BranchPlantKey	Integer	-1	-1	none	N		
2	Date	Date	-1	-1	none	N		
3	ItemMasterKey	Integer	-1	-1	none	N		
4	TransTypeKey	Integer	-1	-1	none	N		
5	CustVendorKey	Integer	-1	-1	none	N		
6	UnitCost	Number	-1	-1	none	N		
7	Currency	String	-1	-1	none	N		
8	Quantity	Integer	-1	-1	none	N		
9								
1..								
1..								

[Get fields from header row...](#)

[Help](#) OK Preview rows Cancel


 Examine preview data

Rows of step: Inventory_EXCEL (15 rows)

#	BranchPlantKey	Date	ItemMasterKey	TransTypeKey	CustVendorKey	UnitCost	Currency	Quantity
1	1	2012/08/06 00:00:00.000	19	1	5	11,22	USD	10
2	3	2012/08/13 00:00:00.000	17	3	7	22,33	JPY	20
3	5	2012/08/20 00:00:00.000	9	5	1	33,11	USD	30
4	7	2012/08/27 00:00:00.000	7	1	3	33,22	JPY	40
5	9	2012/09/03 00:00:00.000	15	1	11	22,01	GBP	50
6	11	2012/09/10 00:00:00.000	5	3	15	11,03	USD	20
7	13	2013/09/18 00:00:00.000	9	5	13	15,04	GBP	22
8	15	2012/09/24 00:00:00.000	<null>	3	7	17,09	USD	18
9	17	2012/10/01 00:00:00.000	5	3	5	33,02	JPY	15
1..	19	2013/10/08 00:00:00.000	7	1	1	22,11	USD	16
1..	5	2012/10/15 00:00:00.000	1	3	15	17,11	JPY	17
1..	7	2012/10/22 00:00:00.000	3	1	13	21,22	USD	15
1..	9	2012/10/29 00:00:00.000	5	5	11	33,21	USD	18
1..	13	2012/11/05 00:00:00.000	15	<null>	9	44,11	JPY	13
1..	15	2012/11/12 00:00:00.000	13	1	99	12,11	USD	14

we create the

Microsoft Access input

Microsoft Access input

Step name

File Content Fields Additional output fields

Filenames from field

Filename is defined in a field?

Get filename from field

File or directory Add Browse

Regular Expression

Exclude Regular Expression

Selected files:

#	File/Directory	Wildcard (RegEx)
1	C:\Users\hp\Desktop\Projet\AccessInventory.mdb	

Show filename(s)... Delete Edit

Help OK Preview rows Cancel

Microsoft Access input

Step name

File Content Fields Additional output fields

#	Name	Column	Type	Format	Length	Precision	Currency	Decimal	Group	Trim type	Repeat
1	PurchaseDay	PurchaseDay	Integer							none	N
2	PurchaseYear	PurchaseYear	Integer							none	N
3	UnitCost	UnitCost	Number							none	N
4	Quantity	Quantity	Integer							none	N
5	Currency	Currency	String							none	N
6	BranchPlantKey	BranchPlantKey	Integer							none	N
7	ItemMasterKey	ItemMasterKey	Integer							none	N
8	CustVendorKey	CustVendorKey	Integer							none	N
9	TransTypeKey	TransTypeKey	Integer							none	N
1..	PurchaseMonth	PurchaseMonth	Integer							none	N

Get fields

Help OK Preview rows Cancel



#	PurchaseDay	PurchaseYear	UnitCost	Quantity	Currency	BranchPlantKey	ItemMasterKey	CustVendorKey	TransTypeKey	PurchaseMonth
1	10	2011	45,2	10	USD	2	2	4	2	1
2	28	2013	16,1	20	JPY	4	8	18	4	2
3	14	2011	15,7	30	GBP	6	6	14	4	2
4	21	2011	25,6	40	EUR	8	<null>	16	4	3
5	4	2011	54,5	50	EUR	10	14	18	4	4
6	3	2011	24,3	60	USD	12	1	12	2	1
7	8	2011	78,1	70	USD	14	8	14	4	2
8	31	2011	14,45	80	USD	16	4	6	2	1
9	18	2011	22,0	90	USD	18	20	8	2	4
1..	11	2011	12,56	12	JPY	20	18	<null>	2	4
1..	28	2011	18,26	25	JPY	18	12	2	2	3
1..	7	2011	19,11	35	USD	16	10	99	2	3
1..	14	2011	20,77	45	EUR	16	10	12	2	3
1..	17	2011	26,89	55	JPY	14	2	10	4	1
1..	21	2011	25,36	65	USD	8	6	10	4	2
1..	7	2011	14,12	75	GBP	2	4	8	4	2

Spoon - Project_Transformation (changed)

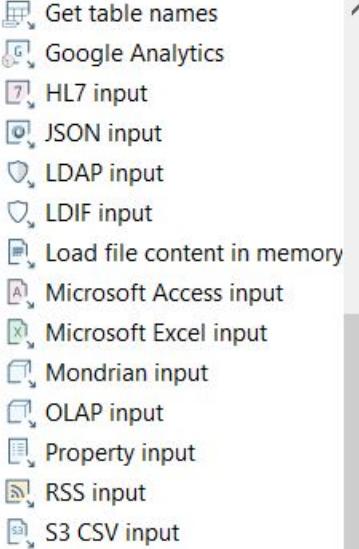
File Edit View Action Tools Help



Welcome! Project_Transformation   Inventory_EXCEL

 Project_Access



- Get table names
- Google Analytics
- HL7 input
- JSON input
- LDAP input
- LDIF input
- Load file content in memory
- Microsoft Access input
- Microsoft Excel input
- Mondrian input
- OLAP input
- Property input
- RSS input
- S3 CSV input

∞



Reject a record if any field is null.

 Filter rows

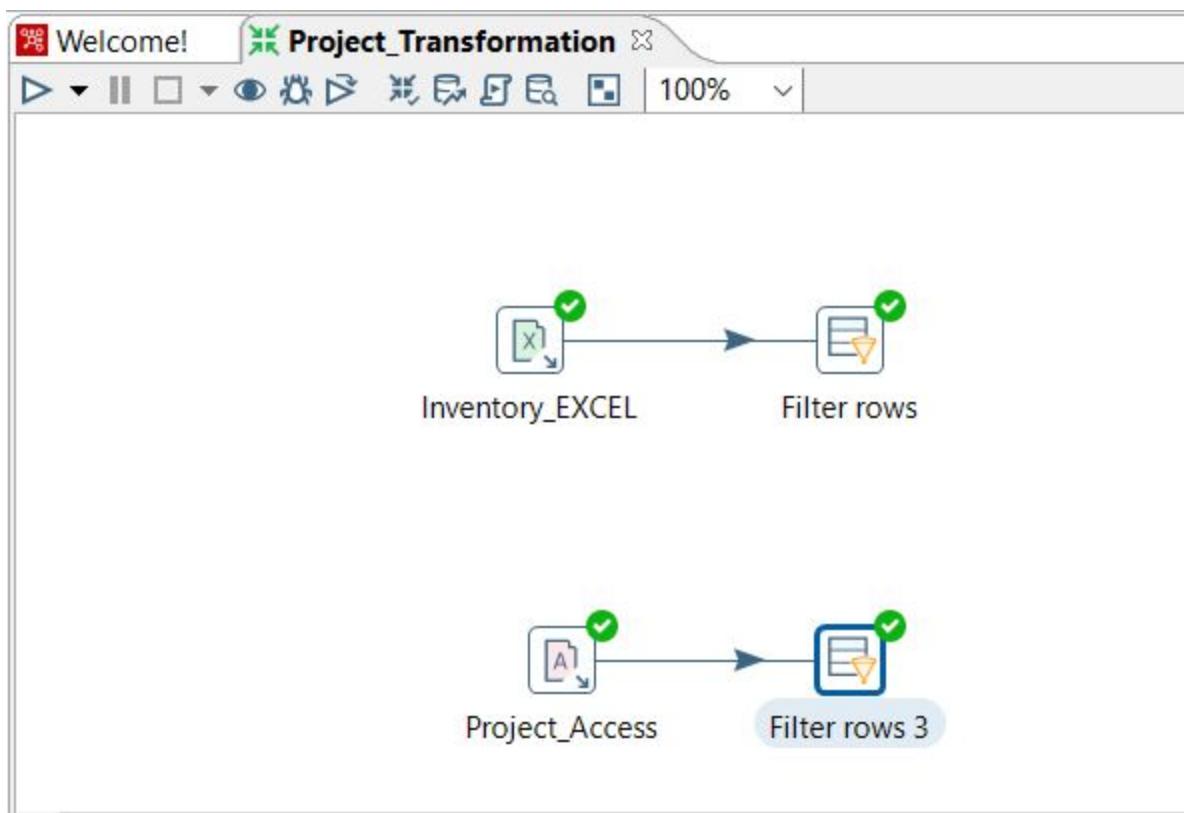
Step name:

Send 'true' data to step:

Send 'false' data to step:

The condition:

PurchaseDay IS NOT NULL
AND
PurchaseYear IS NOT NULL
AND
UnitCost IS NOT NULL
AND
Quantity IS NOT NULL
AND
Currency IS NOT NULL
AND
BranchPlantKey IS NOT NULL
AND
ItemMasterKey IS NOT NULL
AND
CustVendorKey IS NOT NULL
AND
TransTypeKey IS NOT NULL
AND
PurchaseMonth IS NOT NULL



After running the transformation, this is the executing results

The screenshot shows the Talend Studio interface with a transformation project named "Project_Transformation". The flow consists of two main steps: "Inventory_EXCEL" (an Excel source icon) connected to "Filter rows" (a filter icon), and "Project_Access" (an Access source icon) connected to "Filter rows 3" (another filter icon). Below the flow, the "Execution Results" tab is selected, displaying a table of data with the following columns: #, BranchPlantKey, Date, ItemMasterKey, TransTypeKey, CustVendorKey, UnitCost, Currency, and Quantity. The data rows are as follows:

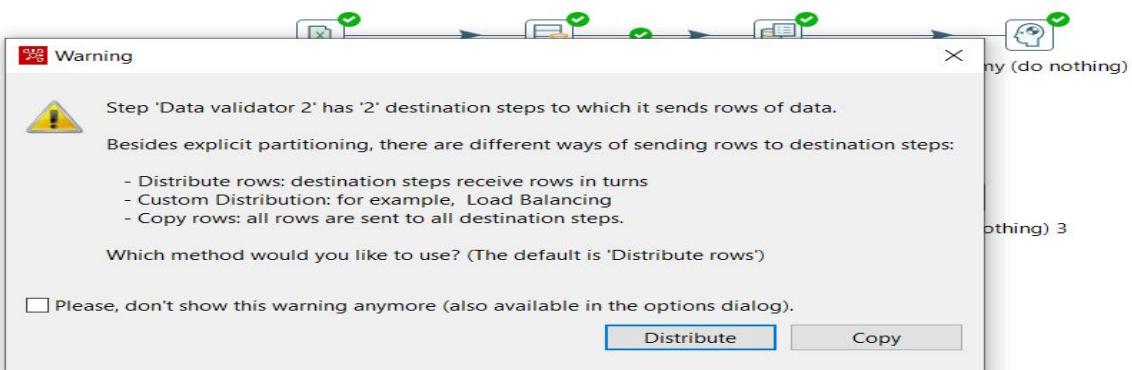
#	BranchPlantKey	Date	ItemMasterKey	TransTypeKey	CustVendorKey	UnitCost	Currency	Quantity
1	1	2012/08/06...	19	1	5	11,22	USD	10
2	3	2012/08/13...	17	3	7	22,33	JPY	20
3	5	2012/08/20...	9	5	1	33,11	USD	30
4	7	2012/08/27...	7	1	3	33,22	JPY	40
5	9	2012/09/03...	15	1	11	22,01	GBP	50
6	11	2012/09/10...	5	3	15	11,03	USD	20
8	17	2012/10/01...	5	3	5	33,02	JPY	15
1..	5	2012/10/15...	1	3	15	17,11	JPY	17
1..	7	2012/10/22...	3	1	13	21,22	USD	15
1..	9	2012/10/29...	5	5	11	33,21	USD	18
1..	15	2012/11/12...	13	1	99	12,11	USD	14
7	13	2013/09/18...	9	5	13	15,04	GBP	22
9	19	2013/10/08...	7	1	1	22,11	USD	16

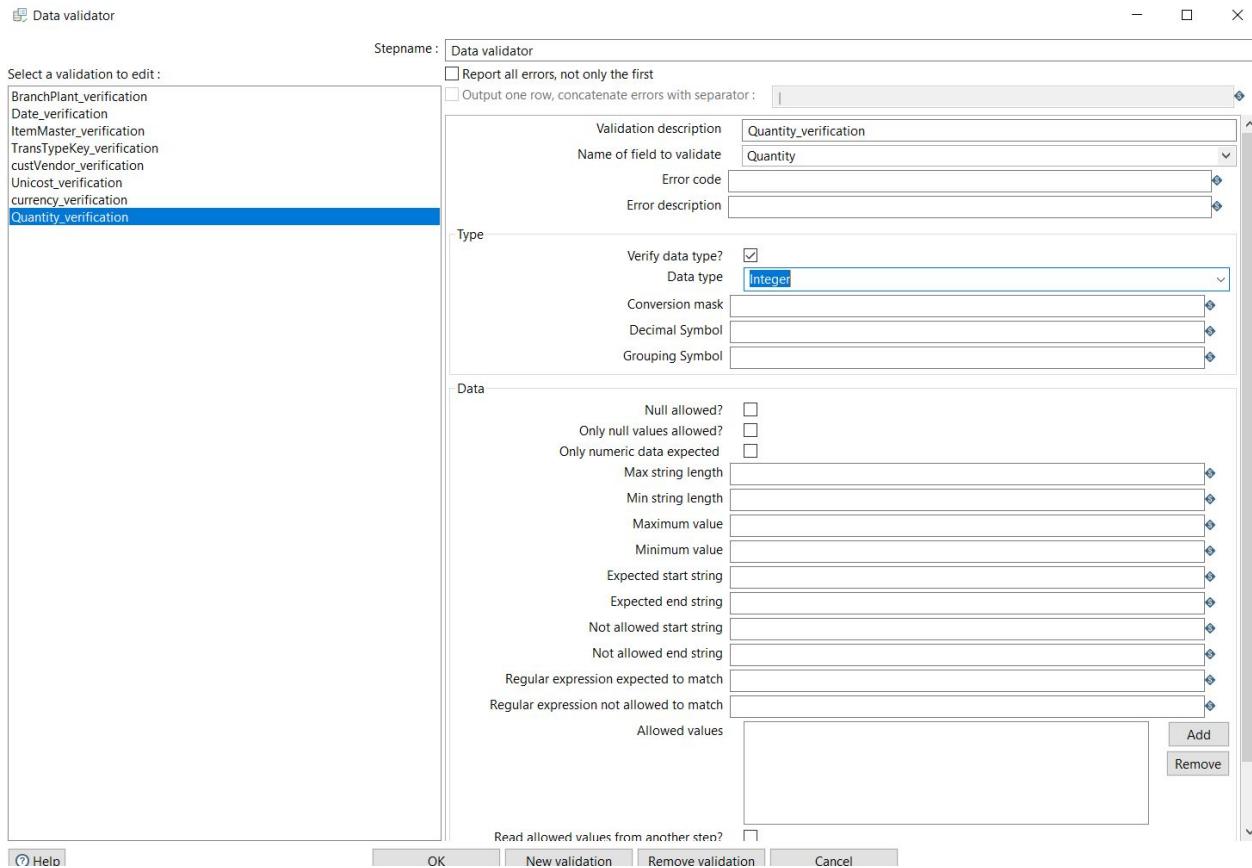
Execution Results

#	PurchaseDay	PurchaseYear	UnitCost	Quantity	Currency	BranchPlantKey	ItemMasterKey	CustVendorKey	TransTypeKey	PurchaseMonth
1	10	2011	45,2	10	USD	2	2	4	2	1
2	28	2013	16,1	20	JPY	4	8	18	4	2
3	14	2011	15,7	30	GBP	6	6	14	4	2
4	4	2011	54,5	50	EUR	10	14	18	4	4
5	3	2011	24,3	60	USD	12	1	12	2	1
6	8	2011	78,1	70	USD	14	8	14	4	2
7	31	2011	14,45	80	USD	16	4	6	2	1
8	18	2011	22,0	90	USD	18	20	8	2	4
9	28	2011	18,26	25	JPY	18	12	2	2	3
1..	7	2011	19,11	35	USD	16	10	99	2	3
1..	14	2011	20,77	45	EUR	16	10	12	2	3
1..	17	2011	26,89	55	JPY	14	2	10	4	1
1..	21	2011	25,36	65	USD	8	6	10	4	2
1..	7	2011	14,12	75	GBP	2	4	8	4	2

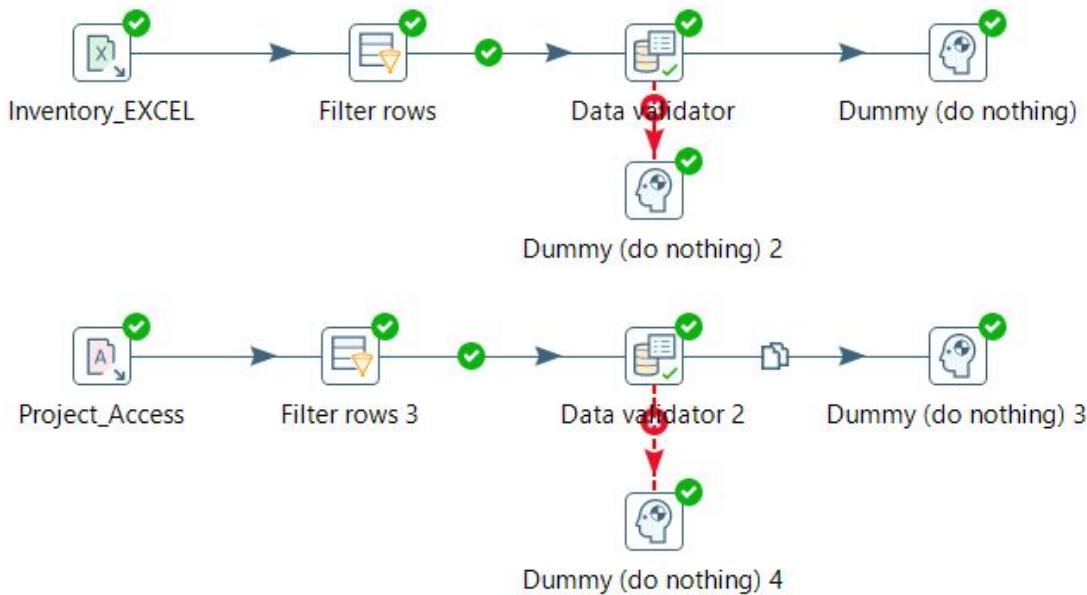
Reject a record if any field value does not match its data type.

To do the verification, we create a data validator for all the fields, it checks the type of data and filter the admitted rows

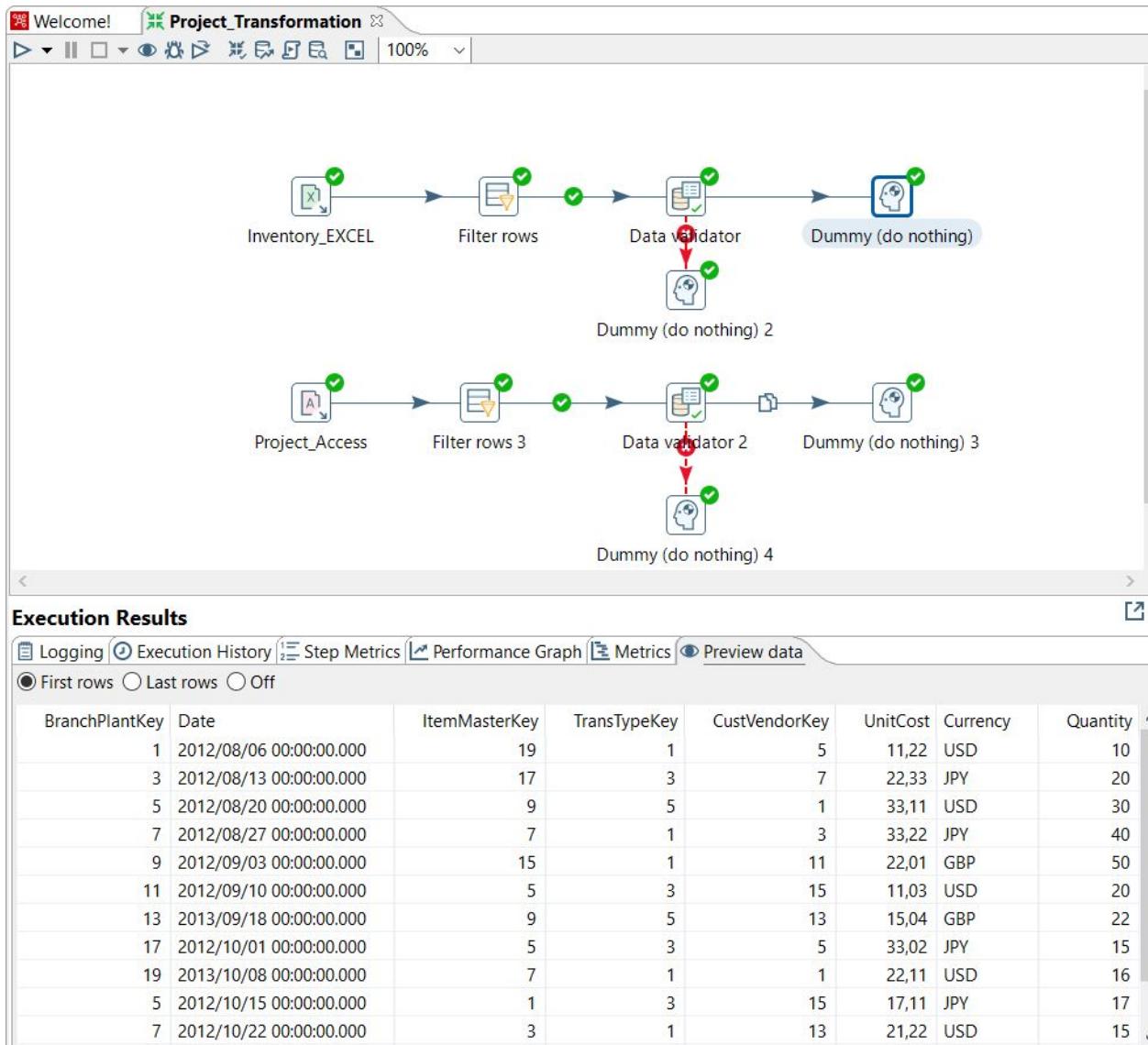


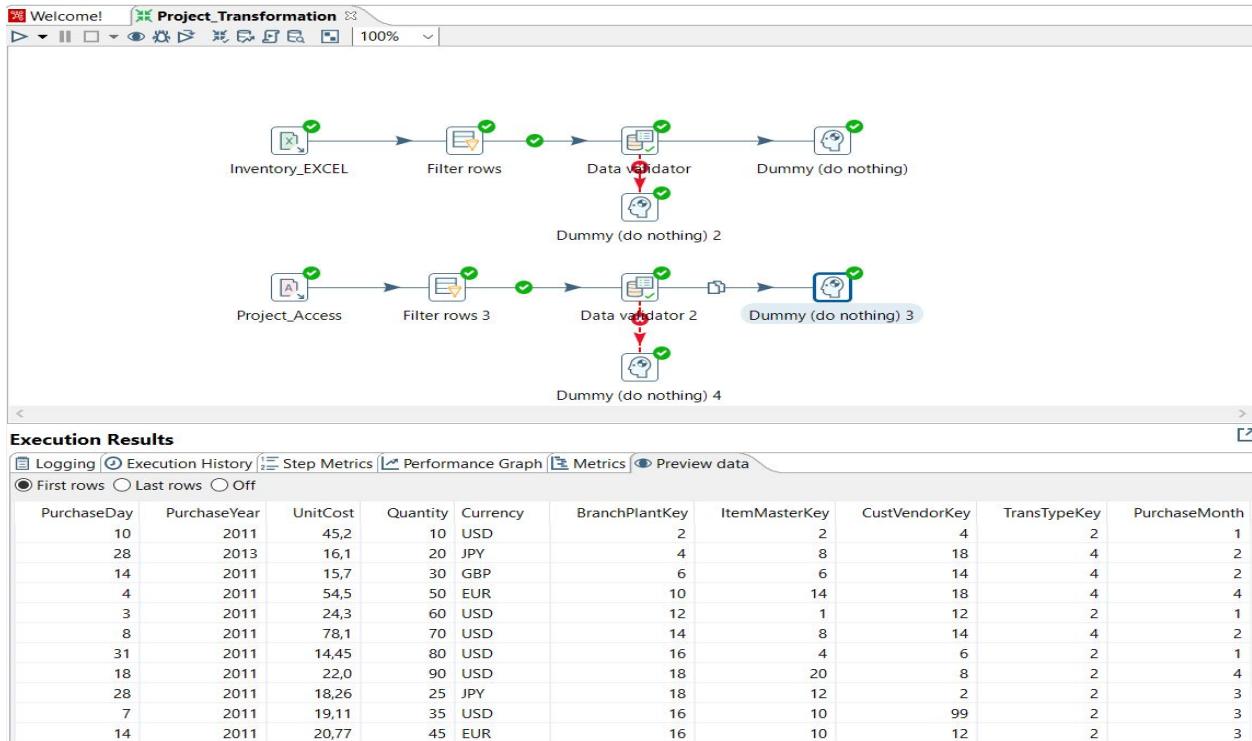


After running the transformation



All data types are suiting the type of their field





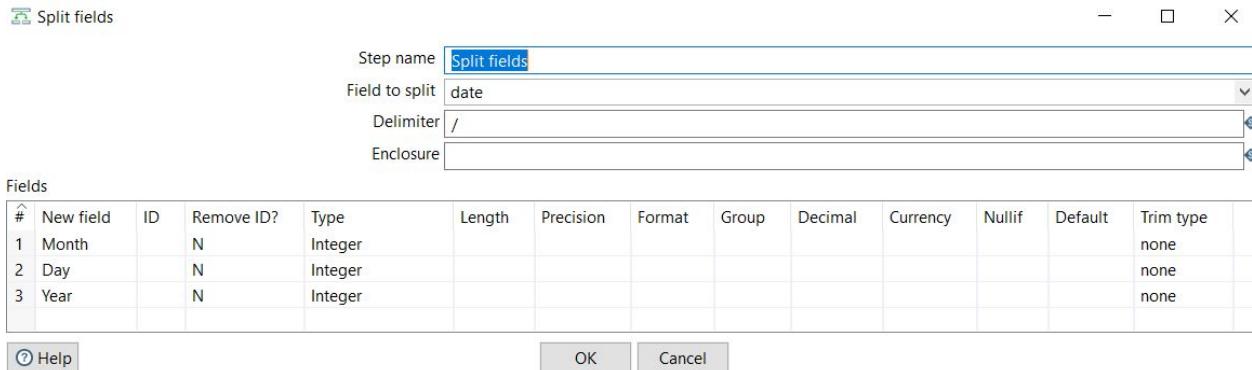
Reject invalid dates: the combination of month, day, and year should be a valid date (including leap year processing) that exists in the *Date_Dim* table.

Parsing Dates with Excel Data Sources

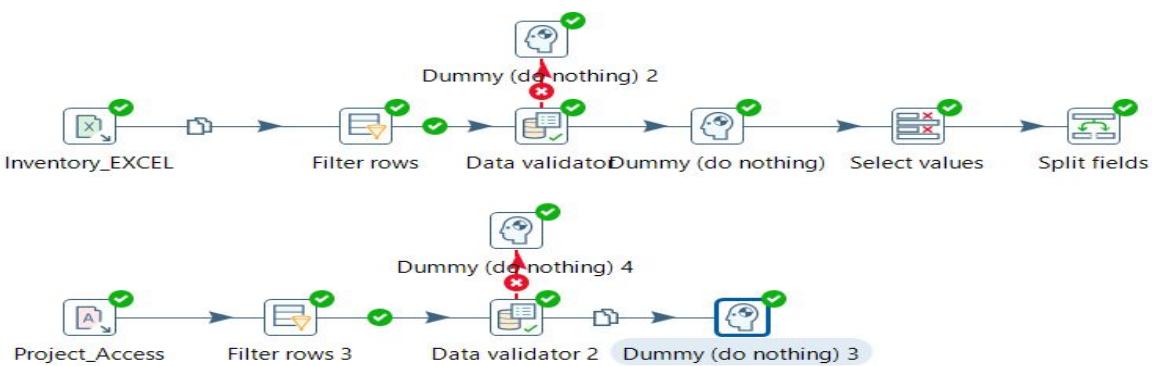
To do this task we need to parse the date field and split it into year, day, month fields , the check if they are valid or not.

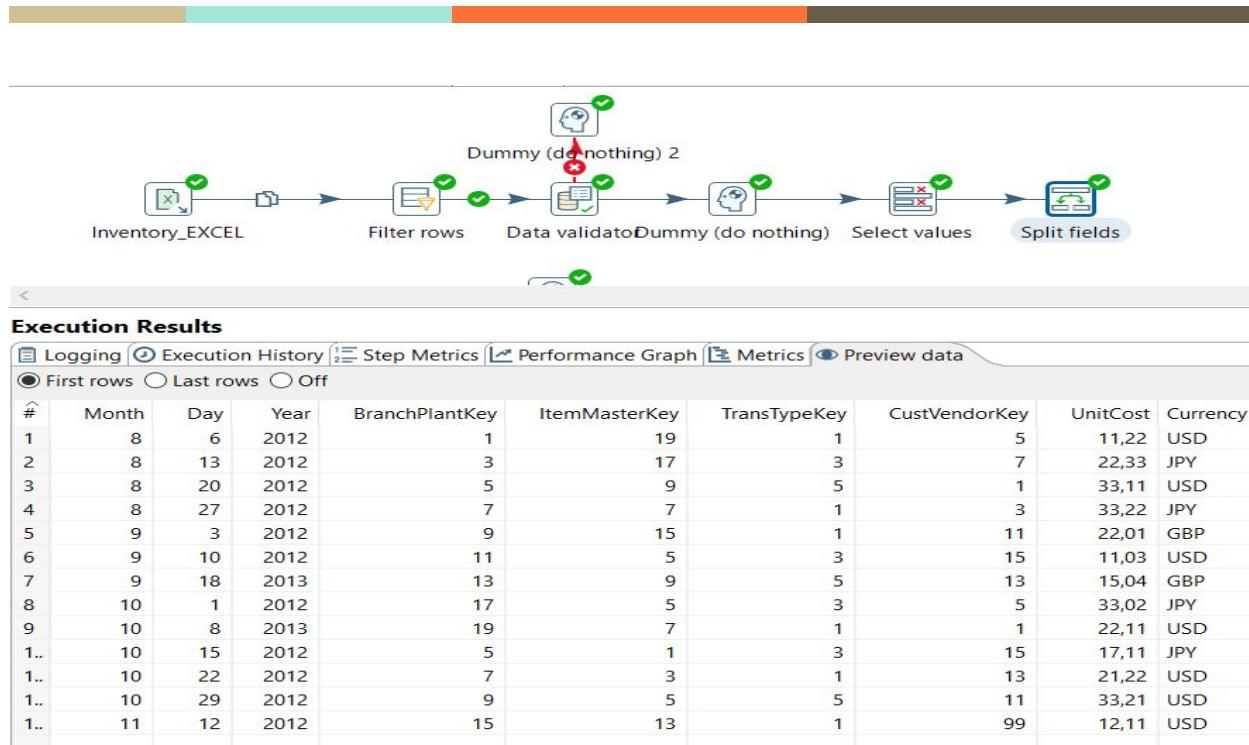
we create selected values step and fill in the changes of the field

Then, we create the split step



After running the transformation





We create a data validator to put validations conditions of the components of date

The Data Validator configuration dialog is shown with the following settings:

- Type** section:
 - Validation description: Month_verification
 - Name of field to validate: Month
 - Error code: (empty)
 - Error description: (empty)
- Data** section:
 - Verify data type?
 - Data type: Integer
 - Conversion mask: (empty)
 - Decimal Symbol: (empty)
 - Grouping Symbol: (empty)
 - Null allowed?
 - Only null values allowed?
 - Only numeric data expected?
 - Max string length: (empty)
 - Min string length: (empty)
 - Maximum value: 12
 - Minimum value: 1

Validation description	Day_verificationn
Name of field to validate	Day
Error code	
Error description	
Type	
Verify data type?	<input checked="" type="checkbox"/>
Data type	Integer
Conversion mask	
Decimal Symbol	
Grouping Symbol	
Data	
Null allowed?	<input type="checkbox"/>
Only null values allowed?	<input type="checkbox"/>
Only numeric data expected	<input type="checkbox"/>
Max string length	
Min string length	
Maximum value	31
Minimum value	1

Data validator

Stepname :

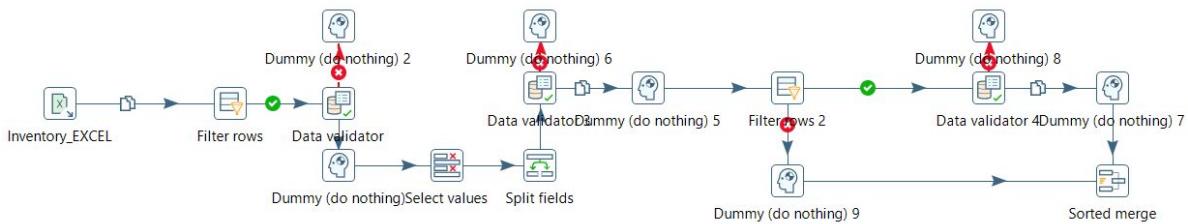
Select a validation to edit :

- Report all errors, not only the first
- Output one row, concatenate errors with separator :

Month_verification	Validation description	Year_verification
Day_verificationn	Name of field to validate	Year
Year_verification	Error code	
	Error description	
Type		
Verify data type?	<input checked="" type="checkbox"/>	
Data type	Integer	
Conversion mask		
Decimal Symbol		
Grouping Symbol		
Data		
Null allowed?	<input checked="" type="checkbox"/>	
Only null values allowed?	<input type="checkbox"/>	
Only numeric data expected	<input type="checkbox"/>	
Max string length		
Min string length		
Maximum value	2100	
Minimum value	1900	
Expected start string		
Expected end string		
Not allowed start string		
Not allowed end string		
Regular expression expected to match	[0-9][0-9][0-9][0-9]	
Regular expression not allowed to match		

To verify the leap year

we use a filter row to select the rows with the month 2 and we fix the marge of the values of the day using a data validator, and then we merge all the records in a sorted merge.



Filter rows

Step name **Filter rows 2**

Send 'true' data to step: **Data validator 4**

Send 'false' data to step: **Dummy (do nothing) 9**

The condition:

<input type="text" value="Month"/>	=	<input type="text"/>
		<input type="text" value="2"/> (Integer)

Data validator

Stepname : **Data validator 4**

Select a validation to edit : **Day**

Stepname : **Data validator 4**

Report all errors, not only the first
 Output one row, concatenate errors with separator :

Type

Validation description : Day
Name of field to validate : Day
Error code :
Error description :
Verify data type?
Data type : Integer
Conversion mask :
Decimal Symbol :
Grouping Symbol :
Data

Null allowed?
Only null values allowed?
Only numeric data expected?
Max string length :
Min string length :
Maximum value : 29
Minimum value : 1
Expected start string :
Expected end string :
Not allowed start string :
Not allowed end string :
Regular expression expected to match :
Regular expression not allowed to match :
Allowed values :

Read allowed values from another step?

OK New validation Remove validation Cancel

Reject invalid foreign key references: the Customer vendor key, branch plant key, transtype key, and Item master key must be valid references to rows of the respective tables in the inventory data warehouse.

BranchPlantKey	BranchPlantId	CompanyKey	CarryingCost	CostMethod	BPName
10	10	3	0.11	07	Branch Plant 10
11	11	3	0.17	07	Branch Plant 11
12	12	3	0.16	07	Branch Plant 12
13	13	4	0.05	07	Branch Plant 13
14	14	4	0.18	07	Branch Plant 14
15	15	4	0.10	07	Branch Plant 15
16	16	4	0.08	07	Branch Plant 16
17	17	5	0.21	07	Branch Plant 17
18	18	5	0.14	07	Branch Plant 18
19	19	5	0.14	07	Branch Plant 19
20	20	5	0.23	07	Branch Plant 20
NULL	NULL	NULL	NULL	NULL	NULL

Type	Validation description	BranchPlant_verification
	Name of field to validate	BranchPlantKey
	Error code	
	Error description	
Data	Verify data type?	<input checked="" type="checkbox"/>
	Data type	Integer
	Conversion mask	
	Decimal Symbol	
	Grouping Symbol	
	Null allowed?	<input type="checkbox"/>
	Only null values allowed?	<input type="checkbox"/>
	Only numeric data expected	<input type="checkbox"/>
	Max string length	
	Min string length	
	Maximum value	20
	Minimum value	1

ItemMasterKey	ShortItemId	SecondItemId	ThirdItemId	ItemCatCode1	ItemCatCode2	ItemDesc	UOM
10	10	Second Part 10	Thrid Part 10	3	2	Part Description 10	EA
11	11	Second Part 11	Thrid Part 11	6	7	Part Description 11	EA
12	12	Second Part 12	Thrid Part 12	6	6	Part Description 12	EA
13	13	Second Part 13	Thrid Part 13	4	3	Part Description 13	EA
14	14	Second Part 14	Thrid Part 14	3	2	Part Description 14	EA
15	15	Second Part 15	Thrid Part 15	5	8	Part Description 15	EA
16	16	Second Part 16	Thrid Part 16	6	1	Part Description 16	EA
17	17	Second Part 17	Thrid Part 17	4	4	Part Description 17	EA
18	18	Second Part 18	Thrid Part 18	7	4	Part Description 18	EA
19	19	Second Part 19	Thrid Part 19	4	5	Part Description 19	EA
20	20	Second Part 20	Thrid Part 20	3	3	Part Description 20	EA

Validation description	ItemMaster_verification
Name of field to validate	ItemMasterKey
Error code	
Error description	
Type	<p>Verify data type? <input checked="" type="checkbox"/></p> <p>Data type Integer</p> <p>Conversion mask</p> <p>Decimal Symbol</p> <p>Grouping Symbol</p>
Data	<p>Null allowed? <input type="checkbox"/></p> <p>Only null values allowed? <input type="checkbox"/></p> <p>Only numeric data expected <input type="checkbox"/></p> <p>Max string length</p> <p>Min string length</p> <p>Maximum value 20</p> <p>Minimum value 1</p>

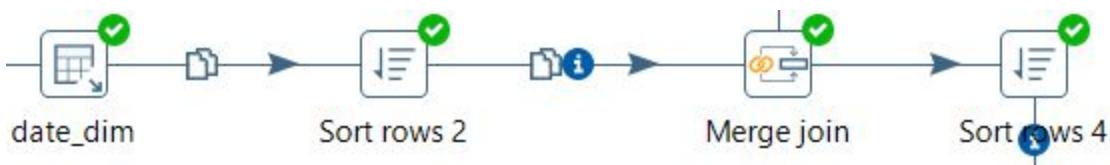
TransTypeKey	TransTypeCodeId	TransDescription
1	IA	Inventory Adjustment
2	IT	Inventory Transfer
3	IS	Inventory Simple Issue
4	OV	Purchase Order Receipt
5	AR	Sales Order Shipment
NUL	NUL	NUL

Validation description	TransTypeKey_verification
Name of field to validate	CustVendorKey
Error code	
Error description	
Type	<p>Verify data type? <input checked="" type="checkbox"/></p> <p>Data type Integer</p> <p>Conversion mask</p> <p>Decimal Symbol</p> <p>Grouping Symbol</p>
Data	<p>Null allowed? <input type="checkbox"/></p> <p>Only null values allowed? <input type="checkbox"/></p> <p>Only numeric data expected <input type="checkbox"/></p> <p>Max string length</p> <p>Min string length</p> <p>Maximum value 5</p> <p>Minimum value 1</p>

Validation description	<input type="text" value="custVendor_verification"/>
Name of field to validate	<input type="text" value="CustVendorKey"/>
Error code	<input type="text"/>
Error description	<input type="text"/>
Type	<p>Verify data type? <input checked="" type="checkbox"/></p> <p>Data type <input type="text" value="Integer"/></p> <p>Conversion mask <input type="text"/></p> <p>Decimal Symbol <input type="text"/></p> <p>Grouping Symbol <input type="text"/></p>
Data	<p>Null allowed? <input type="checkbox"/></p> <p>Only null values allowed? <input type="checkbox"/></p> <p>Only numeric data expected <input type="checkbox"/></p> <p>Max string length <input type="text"/></p> <p>Min string length <input type="text"/></p> <p>Maximum value <input type="text" value="20"/></p> <p>Minimum value <input type="text" value="1"/></p>

After validation, you should perform the following processing steps. These steps will enable the data to be loaded into the *Inventory_Fact* table of the inventory data warehouse.

For data source 1, the month, day, and year fields should be used to find a matching row in the *Date_Dim* table in the inventory data warehouse. After finding the matching row, the *Date_Key* value in the *Date_Dim* row should be used for the *Date_Key* value in the *Inventory_Fact* table.



Merge join

Step name: **Merge join**

First Step: Sort rows

Second Step: Sort rows 2

Join Type: INNER

Keys for 1st step:

#	Key field
1	Day
2	Month
3	Year

Keys for 2nd step:

#	Key field
1	CalDay
2	CalMonth
3	CalYear

For data source 2, the Purchase Date field should be parsed into its day, month, and year components. These components should be used to find a matching row in the *Date_Dim* table. Then, the *Date_Key* value in the matching *Date_Dim* row should be used for the *Date_Key* value in the *Inventory_Fact* table. See the explanation in the following section about parsing dates in Excel data sources.



Merge join

Step name: **Merge join 7**

First Step: Sort rows 14

Second Step: Sort rows 11

Join Type: INNER

Keys for 1st step:

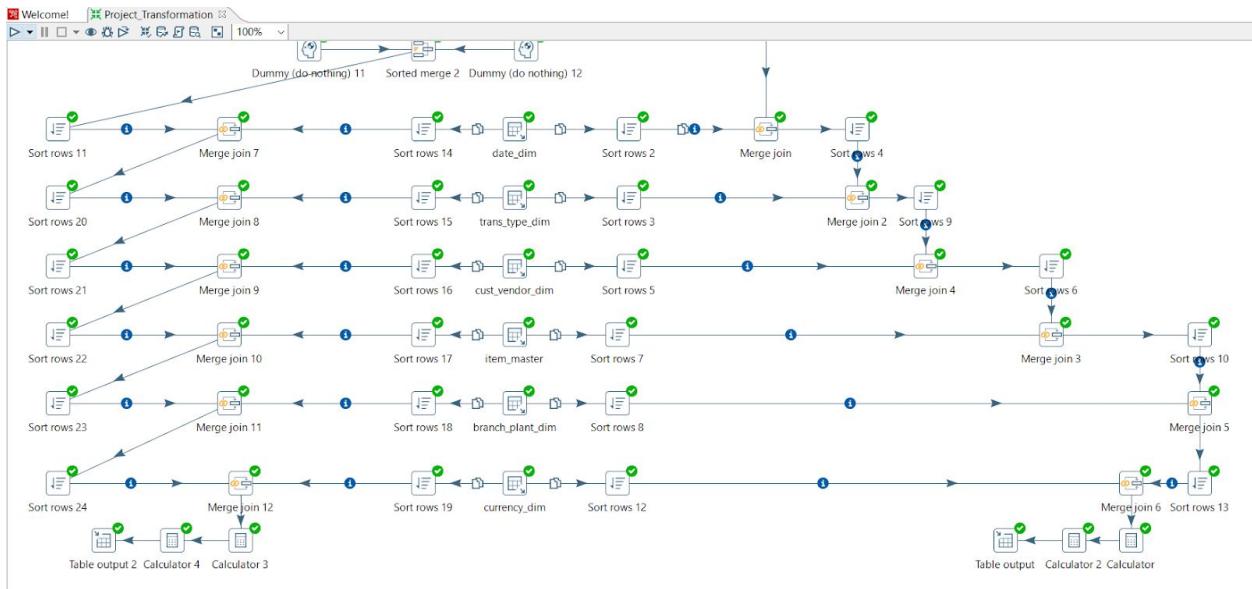
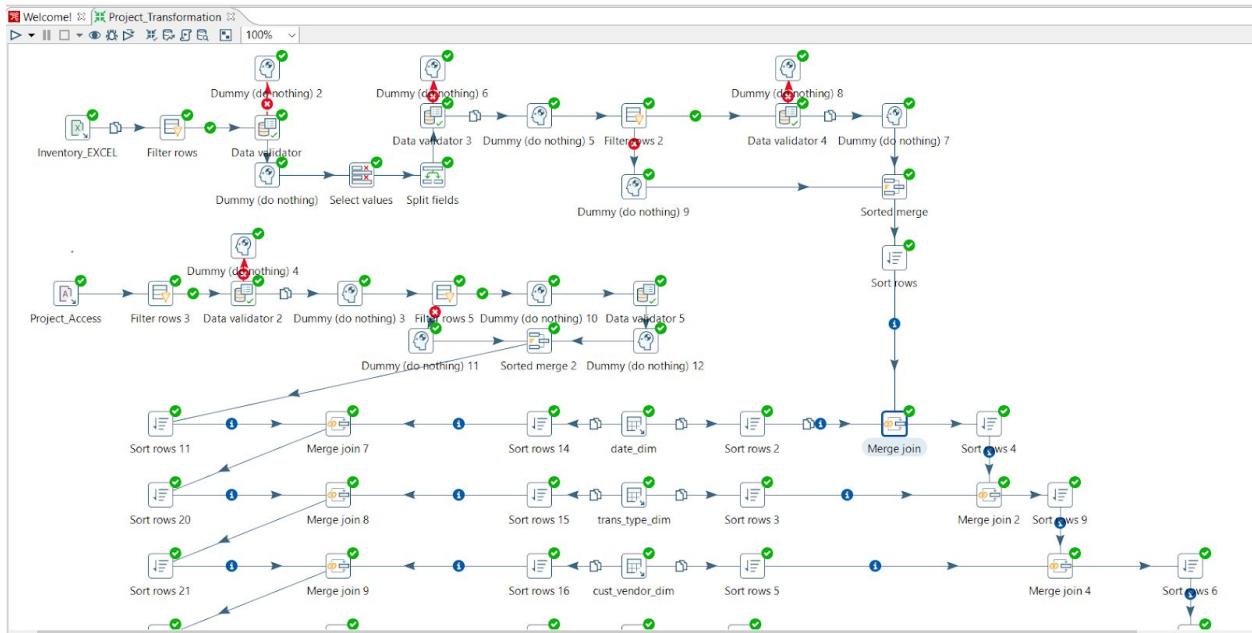
#	Key field
1	CalDay
2	CalMonth
3	CalYear

Keys for 2nd step:

#	Key field
1	PurchaseDay
2	PurchaseMonth
3	PurchaseYear

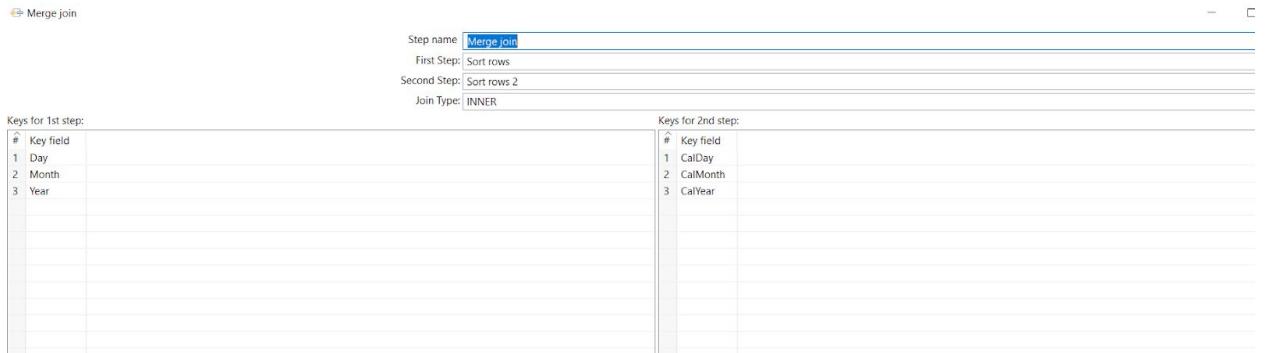
Merge Join Order

we join the tables with the two data sources in the following order .

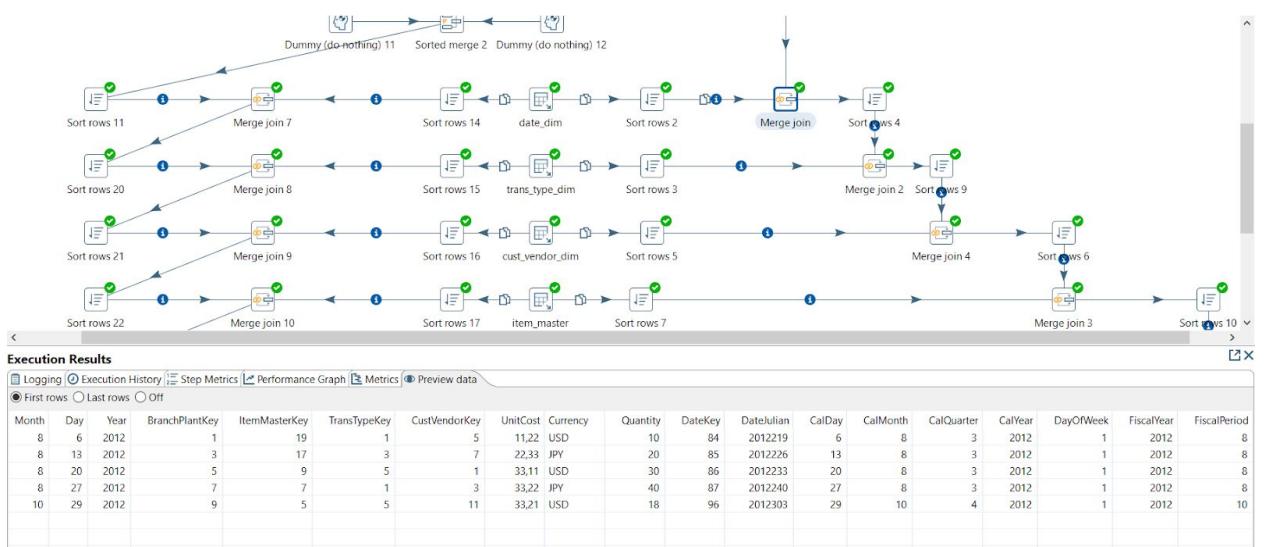


- Date_Dim

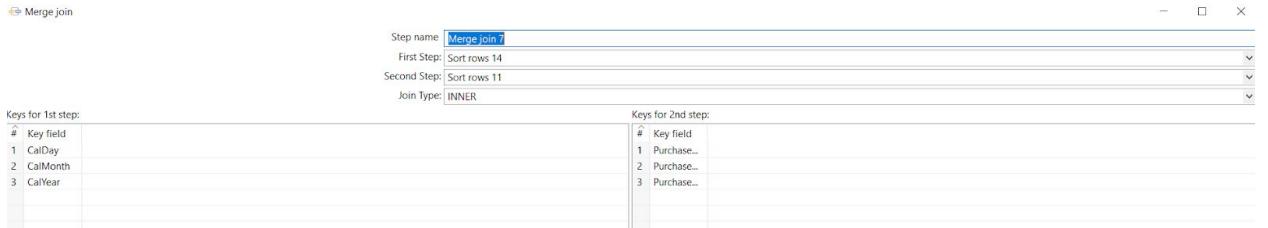
Excel source 1



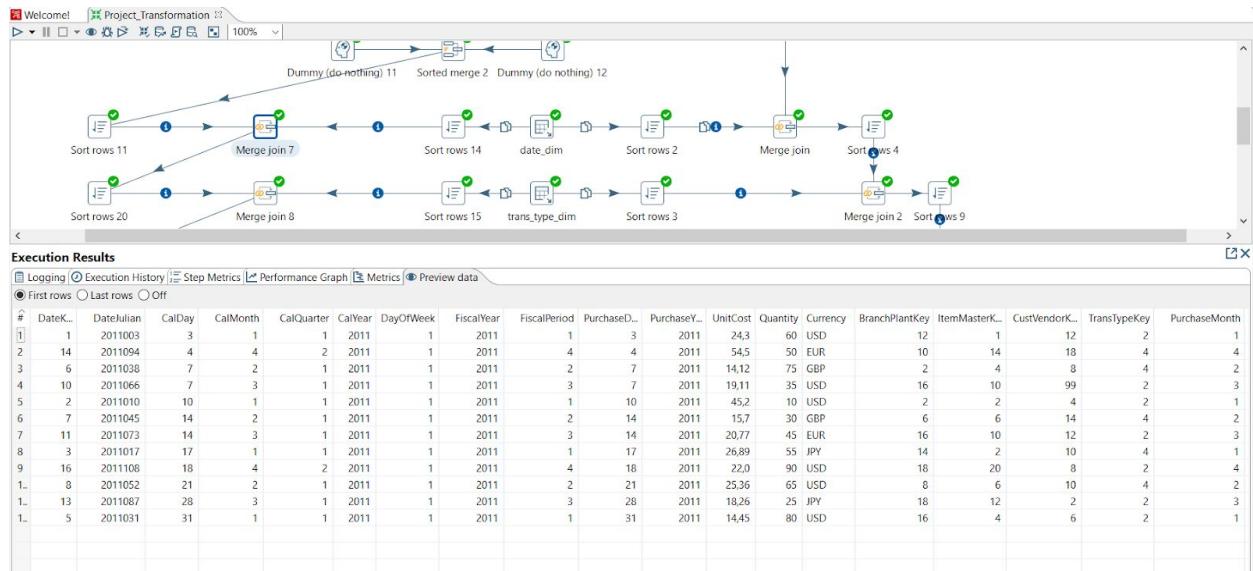
After running the transformation



Access source 2

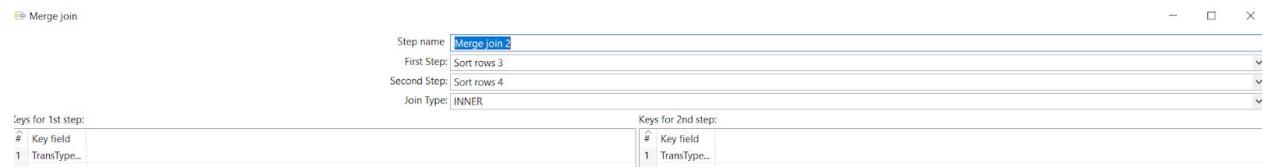


After running the transformation

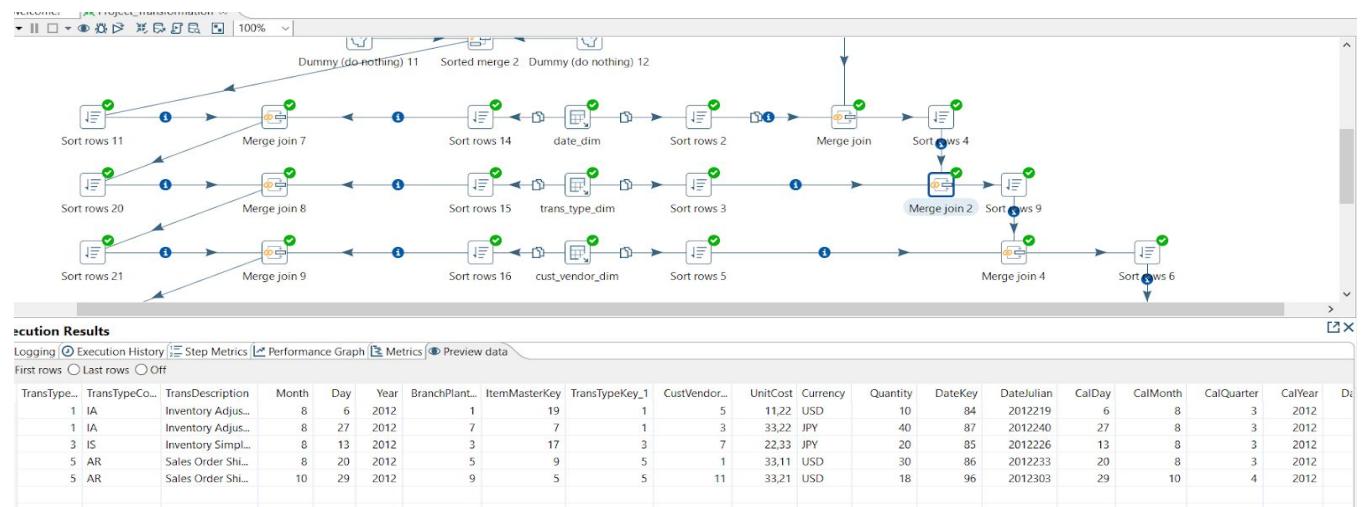


• Trans_Type_Dim

Excel source 1



After running the transformation



Access source 2

Merge join

Step name: Merge join 3

First Step: Sort rows 15

Second Step: Sort rows 20

Join Type: INNER

Keys for 1st step:

Key field
1 TransType...

Keys for 2nd step:

Key field
1 TransType...

After running the transformation

Execution Results

Logging Execution History Step Metrics Performance Graph Metrics Preview data

First rows Last rows Off

#	TransTypeKey	TransTypeCodeId	TransDescription	DateKey	DateJulian	CalDay	CalMonth	CalQuarter	CalYear	DayOfWeek	FiscalYear	FiscalPeriod	Purchas
1	2	IT	Inventory Transfer	1	2011003	3	1	1	2011	1	2011	1	1
2	2	IT	Inventory Transfer	10	2011066	7	3	1	2011	1	2011	3	
3	2	IT	Inventory Transfer	2	2011010	10	1	1	2011	1	2011	1	
4	2	IT	Inventory Transfer	11	2011073	14	3	1	2011	1	2011	3	
5	2	IT	Inventory Transfer	16	2011108	18	4	2	2011	1	2011	4	
6	2	IT	Inventory Transfer	13	2011087	28	3	1	2011	1	2011	3	
7	2	IT	Inventory Transfer	5	2011031	31	1	1	2011	1	2011	1	
8	4	OV	Purchase Order Receipt	14	2011094	4	4	2	2011	1	2011	4	
9	4	OV	Purchase Order Receipt	6	2011038	7	2	1	2011	1	2011	2	
1..	4	OV	Purchase Order Receipt	7	2011045	14	2	1	2011	1	2011	2	
1..	4	OV	Purchase Order Receipt	3	2011017	17	1	1	2011	1	2011	1	

Cust_Vendor_Dim

Excel source 1

Merge join

Step name: Merge join 4

First Step: Sort rows 5

Second Step: Sort rows 9

Join Type: INNER

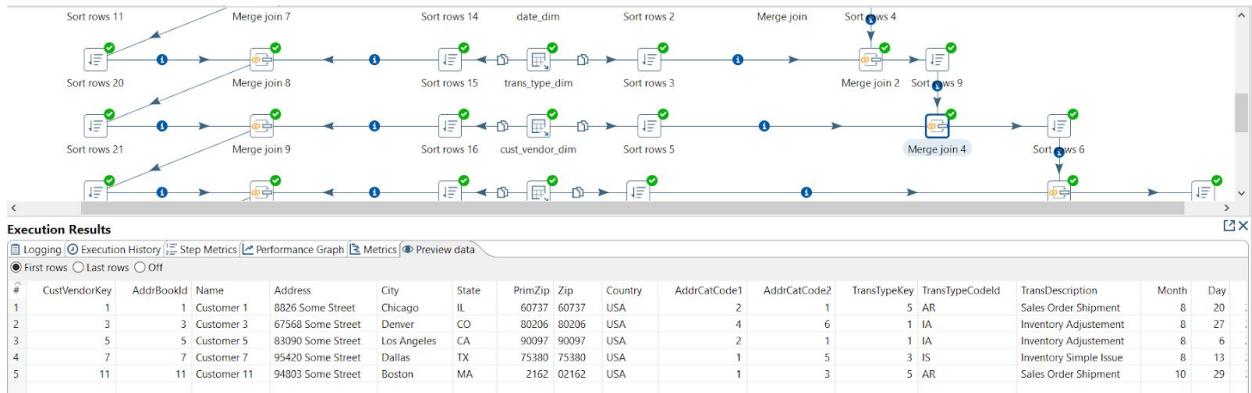
Keys for 1st step:

Key field
1 CustVendor...

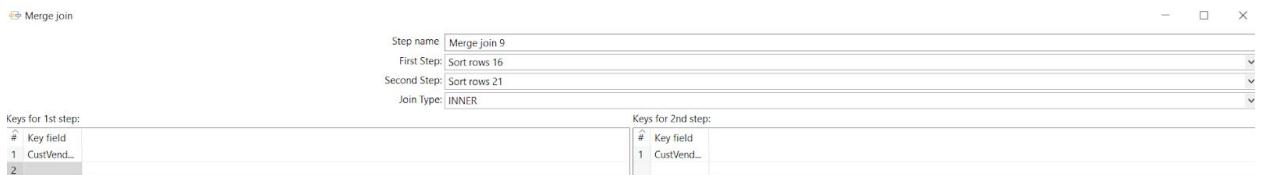
Keys for 2nd step:

Key field
1 CustVendor...

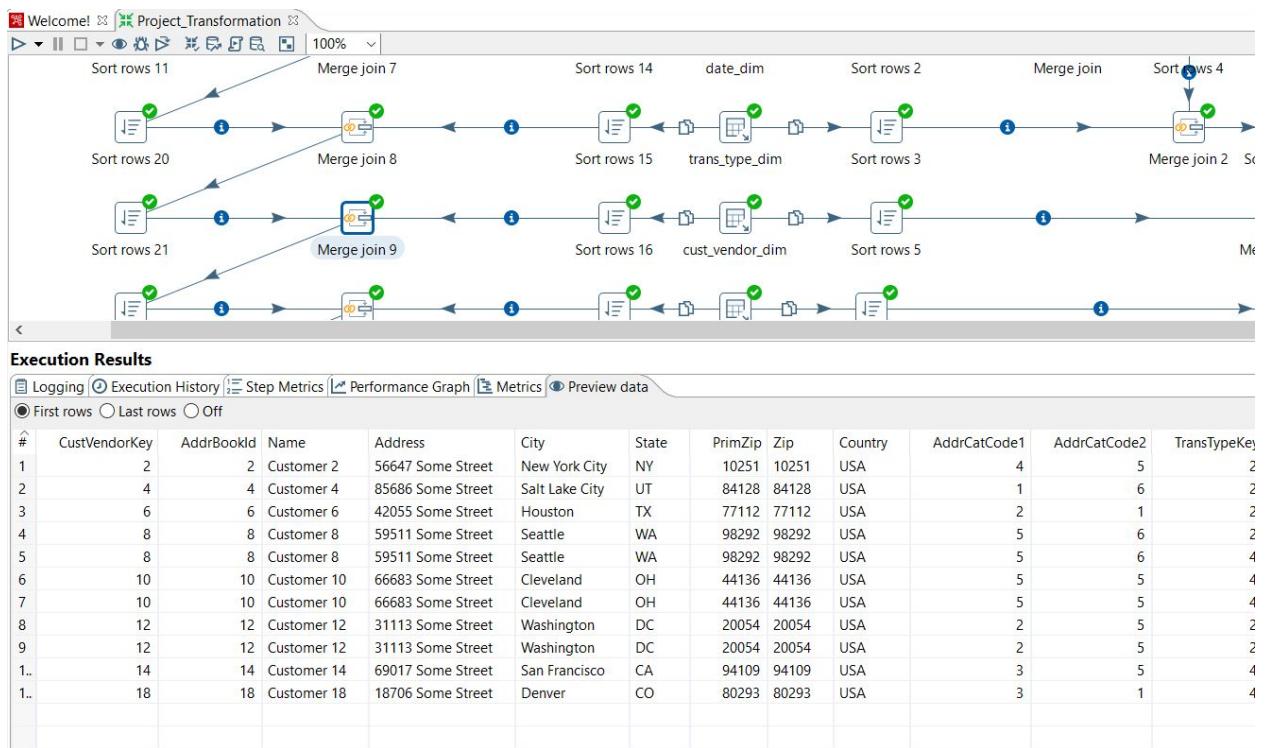
After running the transformation



Access source 2

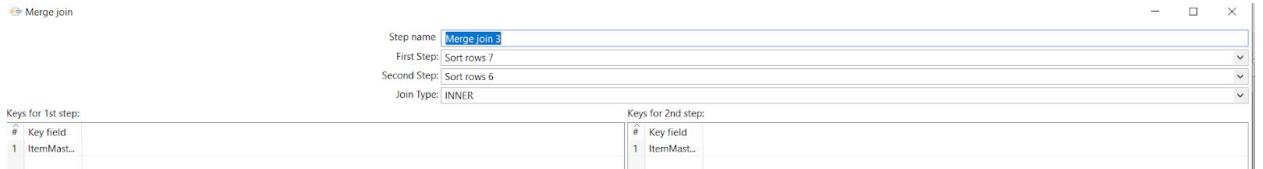


After running the transformation

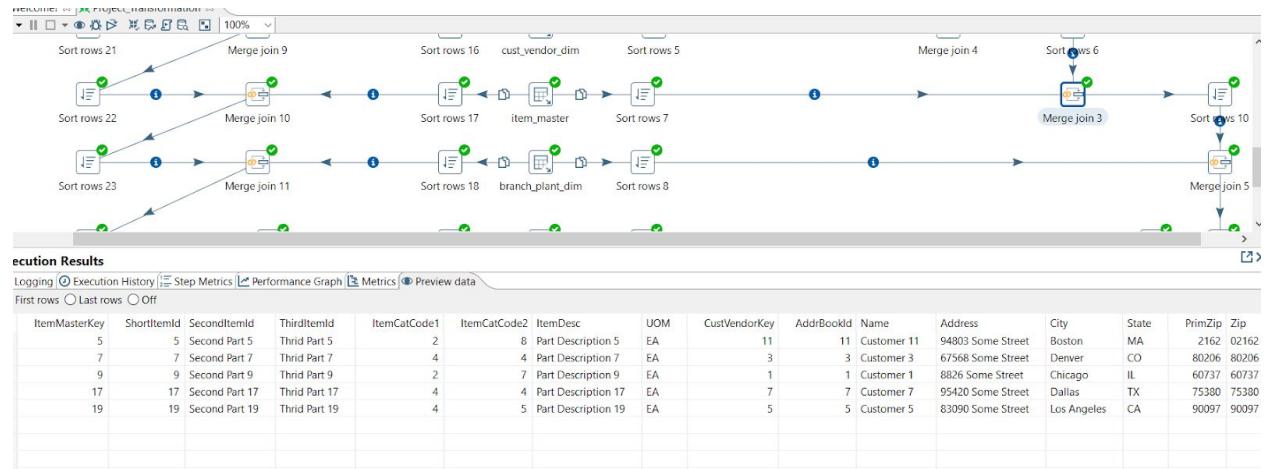


· Item_Master_Dim

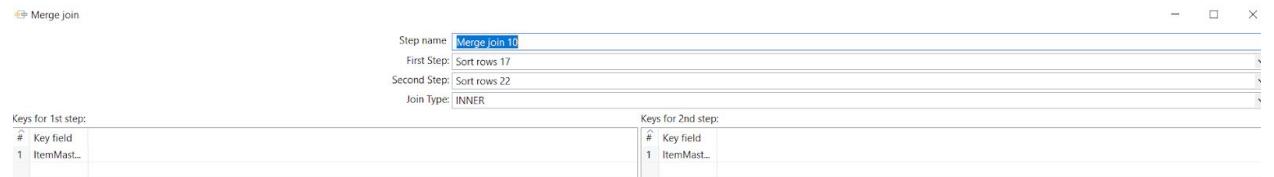
Excel source 1



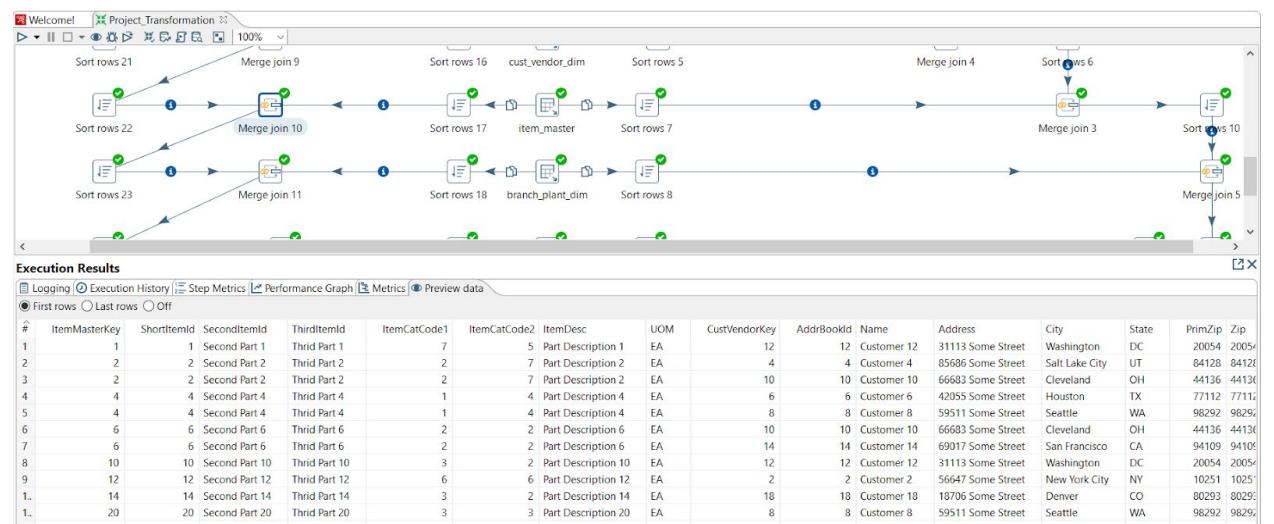
After running the transformation



Access source 2



After running the transformation



Branch_Plant_Dim

Excel source 1

Merge join

Step name: Merge join 5

First Step: Sort rows 10
Second Step: Sort rows 8
Join Type: INNER

Keys for 1st step:
Key field
1 BranchPla...

Keys for 2nd step:
Key field
1 BranchPla...

After running the transformation

Welcome! Project_Transformation 100%

Execution Results

#	ItemMasterKey	ShortItemId	SecondItemId	ThirdItemId	ItemCatCode1	ItemCatCode2	ItemDesc	UOM	CustVendorKey	AddrBookId	Name	Address	City	State	PrimZip	Zip
1	19	19	Second Part 19	Third Part 19	4	5	Part Description 19	EA	5	Customer 5	83090 Some Street	Los Angeles	CA	90097	90097	
2	17	17	Second Part 17	Third Part 17	4	4	Part Description 17	EA	7	Customer 7	95420 Some Street	Dallas	TX	75380	75380	
3	9	9	Second Part 9	Third Part 9	2	7	Part Description 9	EA	1	Customer 1	8826 Some Street	Chicago	IL	60737	60737	
4	7	7	Second Part 7	Third Part 7	4	4	Part Description 7	EA	3	Customer 3	67588 Some Street	Denver	CO	80206	80206	
5	5	5	Second Part 5	Third Part 5	2	8	Part Description 5	EA	11	Customer 11	94803 Some Street	Boston	MA	2162	02162	

Access source 2

Merge join

Step name: Merge join 11

First Step: Sort rows 18
Second Step: Sort rows 23
Join Type: INNER

Keys for 1st step:
Key field
1 BranchPla...

Keys for 2nd step:
Key field
1 BranchPla...

After running the transformation

Welcome! Project_Transformation 100%

Execution Results

#	BranchPlantKey	BranchPlantId	CompanyKey	CarryingCost	CostMethod	BPName	ItemMasterKey	ShortItemId	SecondItemId	ThirdItemId	ItemCatCode1	ItemCatCode2	ItemDesc	UOM	CustVendor
1	2	2	1	0.14	07	Branch Plant 2	2	2	Second Part 2	Third Part 2	2	7	Part Description 2	EA	
2	2	2	1	0.14	07	Branch Plant 2	4	4	Second Part 4	Third Part 4	1	4	Part Description 4	EA	
3	6	6	2	0.01	07	Branch Plant 6	6	6	Second Part 6	Third Part 6	2	2	Part Description 6	EA	
4	8	8	2	0.0	07	Branch Plant 8	6	6	Second Part 6	Third Part 6	2	2	Part Description 6	EA	
5	10	10	3	0.14	07	Branch Plant 10	14	14	Second Part 14	Third Part 14	3	3	Part Description 14	EA	
6	12	12	3	0.16	07	Branch Plant 12	1	1	Second Part 1	Third Part 1	7	5	Part Description 1	EA	
7	14	14	4	0.18	07	Branch Plant 14	2	2	Second Part 2	Third Part 2	2	7	Part Description 2	EA	
8	16	16	4	0.08	07	Branch Plant 16	4	4	Second Part 4	Third Part 4	1	4	Part Description 4	EA	
9	16	16	4	0.08	07	Branch Plant 16	10	10	Second Part 10	Third Part 10	3	2	Part Description 10	EA	
10	18	18	5	0.14	07	Branch Plant 18	12	12	Second Part 12	Third Part 12	6	6	Part Description 12	EA	
11	18	18	5	0.14	07	Branch Plant 18	20	20	Second Part 20	Third Part 20	3	3	Part Description 20	EA	

Currency_Dim

Excel source 1

Merge join

Step name: Merge join 6

First Step: Sort rows 12
Second Step: Sort rows 13
Join type: INNER

Keys for 1st step:		Keys for 2nd step:	
# Key field	1 Currency...	# Key field	1 Currency

After running the transformation

Execution Results

#	Currency_ID	Exchange_Rate	ItemMasterKey	ShortItemId	SecondItemId	ThirdItemId	ItemCatCode1	ItemCatCode2	ItemDesc	UOM	CustVendorKey	AddrBookId	Name	Address	City
1	JPY	0.1	17	17	Second Part 17	Third Part 17	4	4	Part Description 17	EA	7	7	Customer 7	95420 Some Street	Dallas
2	JPY	0.1	7	7	Second Part 7	Third Part 7	4	4	Part Description 7	EA	3	3	Customer 3	67568 Some Street	Denver
3	USD	1.0	19	19	Second Part 19	Third Part 19	4	5	Part Description 19	EA	5	5	Customer 5	83090 Some Street	Los Angel
4	USD	1.0	9	9	Second Part 9	Third Part 9	2	7	Part Description 9	EA	1	1	Customer 1	8826 Some Street	Chicago
5	USD	1.0	5	5	Second Part 5	Third Part 5	2	8	Part Description 5	EA	11	11	Customer 11	94803 Some Street	Boston

Access source 2

Merge join

Step name: Merge join 12

First Step: Sort rows 19
Second Step: Sort rows 24
Join type: INNER

Keys for 1st step:		Keys for 2nd step:	
# Key field	1 Currency...	# Key field	1 Currency

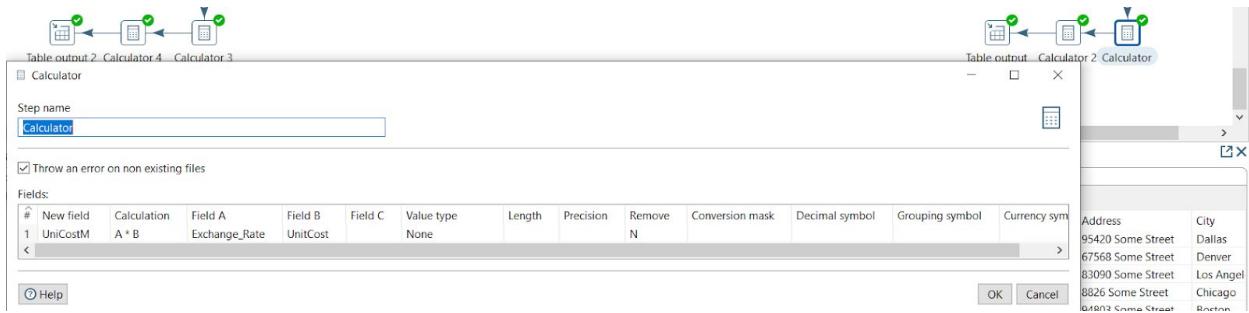
After running the transformation

Execution Results

#	Currency_ID	Exchange_Rate	BranchPlantKey	BranchPlantId	CompanyKey	CarryingCost	CostMethod	BPName	ItemMasterKey	Shc
1	EUR	1,4	10	10	3	0,11	07	Branch Plant 10	14	
2	EUR	1,4	16	16	4	0,08	07	Branch Plant 16	10	
3	GBP	1,54	2	2	1	0,14	07	Branch Plant 2	4	
4	GBP	1,54	6	6	2	0,01	07	Branch Plant 6	6	
5	JPY	0,1	14	14	4	0,18	07	Branch Plant 14	2	
6	JPY	0,1	18	18	5	0,14	07	Branch Plant 18	12	
7	USD	1,0	2	2	1	0,14	07	Branch Plant 2	2	
8	USD	1,0	8	8	2	0,0	07	Branch Plant 8	6	
9	USD	1,0	12	12	3	0,16	07	Branch Plant 12	1	
1..	USD	1,0	16	16	4	0,08	07	Branch Plant 16	4	
1..	USD	1,0	18	18	5	0,14	07	Branch Plant 18	20	

The *UnitCost* column in the *Inventory_Fact* table is computed as the currency conversion factor times the Unit Cost field in the data source.

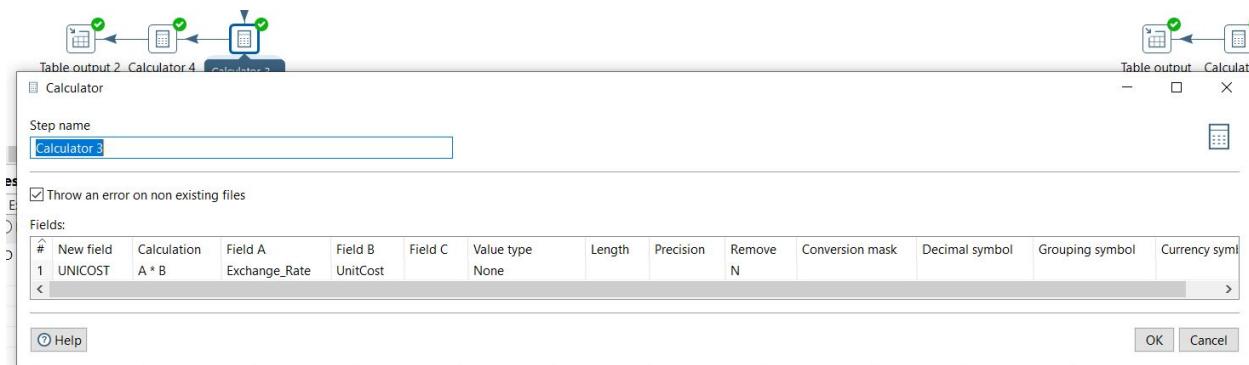
Excel source 1



After running the transformation

Currency_ID	Exchange_Rate	UnitCost	UniCostM
JPY	0,1	22,33	2,233
JPY	0,1	33,22	3,322
USD	1,0	11,22	11,22
USD	1,0	33,11	33,11
USD	1,0	33,21	33,21

Access source 2

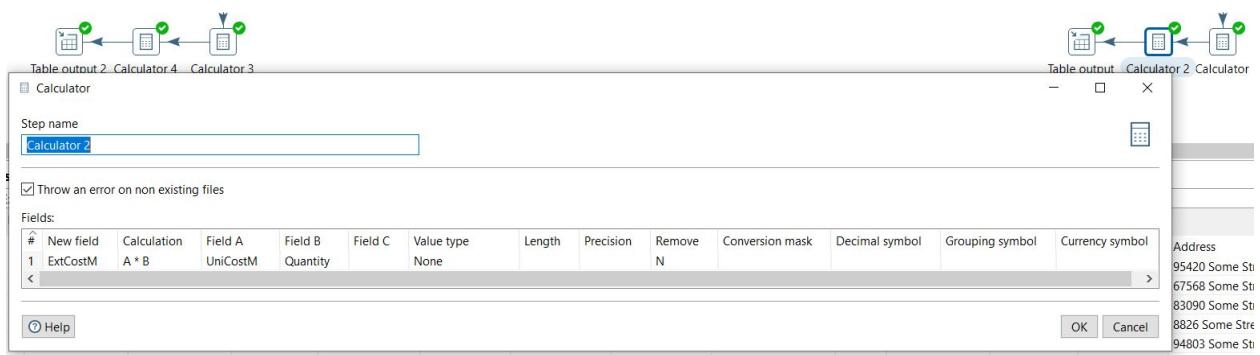


After running the transformation

Currency_ID	Exchange_Rate	UnitCost	UNICOST
EUR	1,4	54,5	76,3
EUR	1,4	20,77	29,078
GBP	1,54	14,12	21,7448
GBP	1,54	15,7	24,178
JPY	0,1	26,89	2,689
JPY	0,1	18,26	1,826
USD	1,0	45,2	45,2
USD	1,0	25,36	25,36
USD	1,0	24,3	24,3
USD	1,0	14,45	14,45
USD	1,0	22,0	22,0

The *ExtCost* column in the *Inventory_Fact* table is computed as the Unit Cost (after currency conversion) times the Quantity. We need to use a Calculator step.

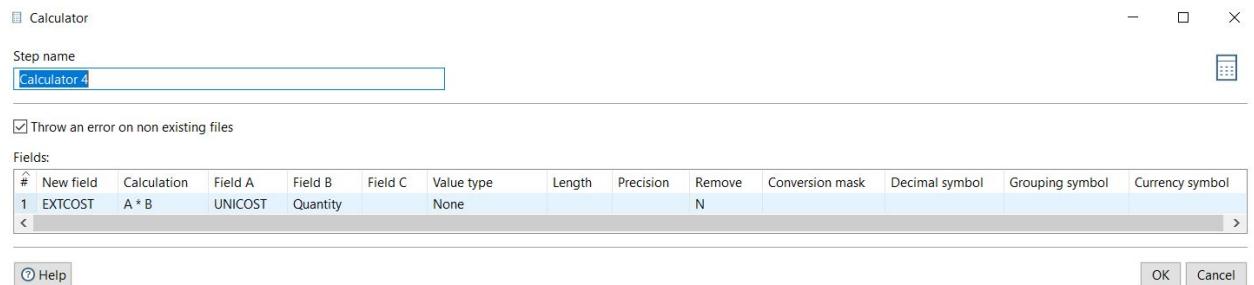
Excel source 1



After running the transformation

Currency	Quantity	UniCostM	ExtCostM
JPY	20	2,233	44,66
JPY	40	3,322	132,88
USD	10	11,22	112,2
USD	30	33,11	993,3
USD	18	33,21	597,78

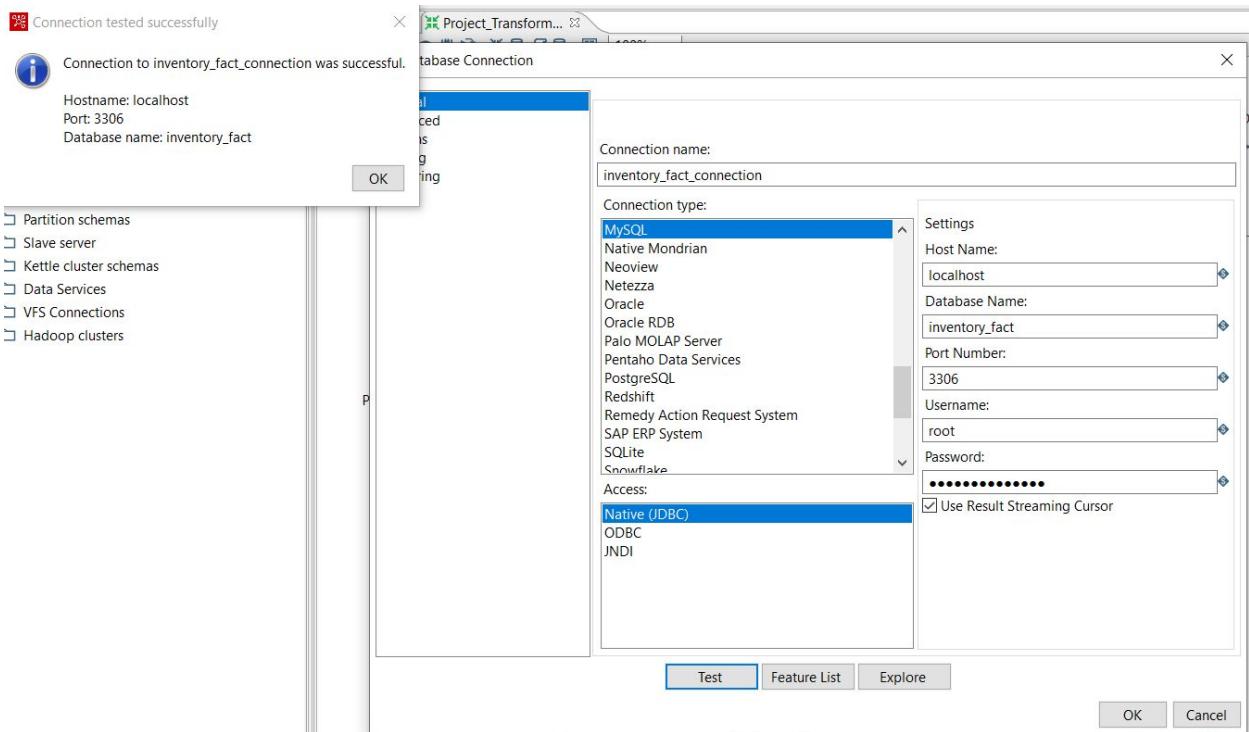
[Access source 2](#)



After running the transformation

Quantity	Currency	BranchPlantKey_1	ItemMasterKey_1	CustVendorKey_1	TransTypeKey_1	PurchaseMonth	UNICOST	EXTCOST
50	EUR	10	14	18	4	4	76,3	3815,0
45	EUR	16	10	12	2	3	29,078	1308,51
75	GBP	2	4	8	4	2	21,7448	1630,86
30	GBP	6	6	14	4	2	24,178	725,34
55	JPY	14	2	10	4	1	2,689	147,895
25	JPY	18	12	2	2	3	1,826	45,65
10	USD	2	2	4	2	1	45,2	452,0
65	USD	8	6	10	4	2	25,36	1648,4
60	USD	12	1	12	2	1	24,3	1458,0
80	USD	16	4	6	2	1	14,45	1156,0
90	USD	18	20	8	2	4	22,0	1980,0

Creation of database connection



the limit line of the original records in the data warehouse created previously.

To insert the data extracted ,verified, transformed in the database, we use the table output step

Excel source 1

The screenshot shows the Table output configuration dialog on the left and the Data Flow Editor interface on the right.

Table output Configuration Dialog:

- Step name:** Table output
- Connection:** inventory_fact_connection
- Target schema:** inventory_fact
- Target table:** inventory_fact
- Commit size:** 1000
- Truncate table:**
- Ignore insert errors:**
- Specify database fields:**
- Main options:**
 - Partition data over tables:**
 - Partitioning field:** dropdown
 - Partition data per month:**
 - Partition data per day:**
 - Use batch update for inserts:**
 - Is the name of the table defined in a field?**
 - Field that contains name of table:** dropdown
 - Store the tablename field:**
 - Return auto-generated key:**
 - Name of auto-generated key field:** dropdown

Data Flow Editor:

The Data Flow Editor shows a data pipeline starting with an Excel source (Table output) connected to a Merge Join 6 (Sort rows 12). This is followed by two Calculator components (Calculator 2, Calculator), another Merge Join 6 (Sort rows 13), and finally a Table output component.

ItemCatCode2	ItemDesc	UOM	CustVendorKey	AddrBookId	Name	Address	City
4	Part Description 17	EA	7	7	Customer 7	95420 Some Street	Dallas
4	Part Description 7	EA	3	3	Customer 3	67568 Some Street	Denver
5	Part Description 19	EA	5	5	Customer 5	83090 Some Street	Los Angel
7	Part Description 9	EA	1	1	Customer 1	8826 Some Street	Chicago
8	Part Description 5	EA	11	11	Customer 11	94803 Some Street	Boston

The screenshot shows the Table output configuration dialog with the "Database fields" tab selected.

Table output Configuration Dialog:

- Step name:** Table output
- Connection:** inventory_fact_connection
- Target schema:** inventory_fact
- Target table:** inventory_fact
- Commit size:** 1000
- Truncate table:**
- Ignore insert errors:**
- Specify database fields:**

Main options: Database fields

Fields to insert:

#	Table field	Stream field
1	ItemMasterK...	ItemMasterK...
2	TransTypeK...	TransTypeKey
3	DateKey	DateKey
4	BranchPlan...	BranchPlantK...
5	CustVendor...	CustVendorK...
6	UnitCost	UniCostM
7	ExtCost	ExtCostM
8	Quantity	Quantity

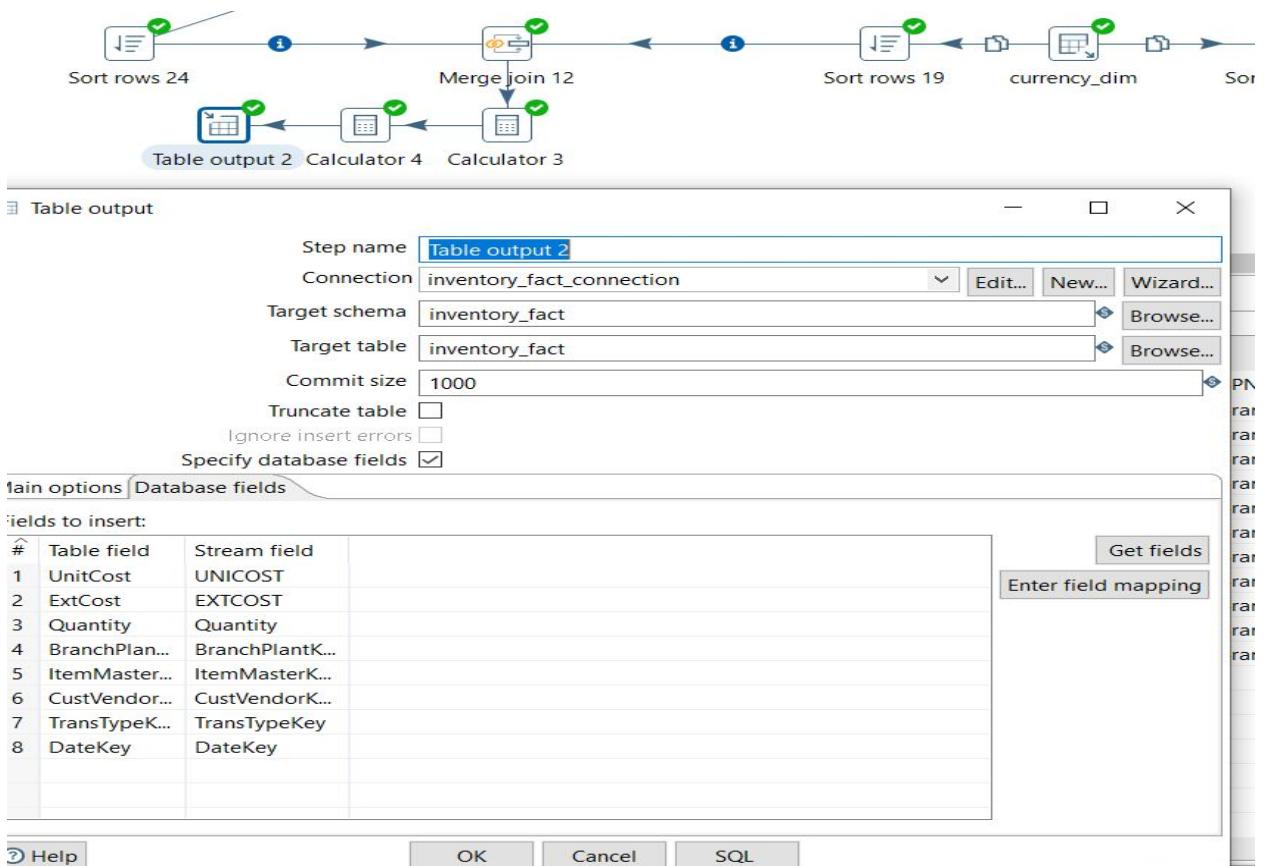
Buttons: Get fields, Enter field mapping, Help, OK, Cancel, SQL

After running the transformation

5 rows added

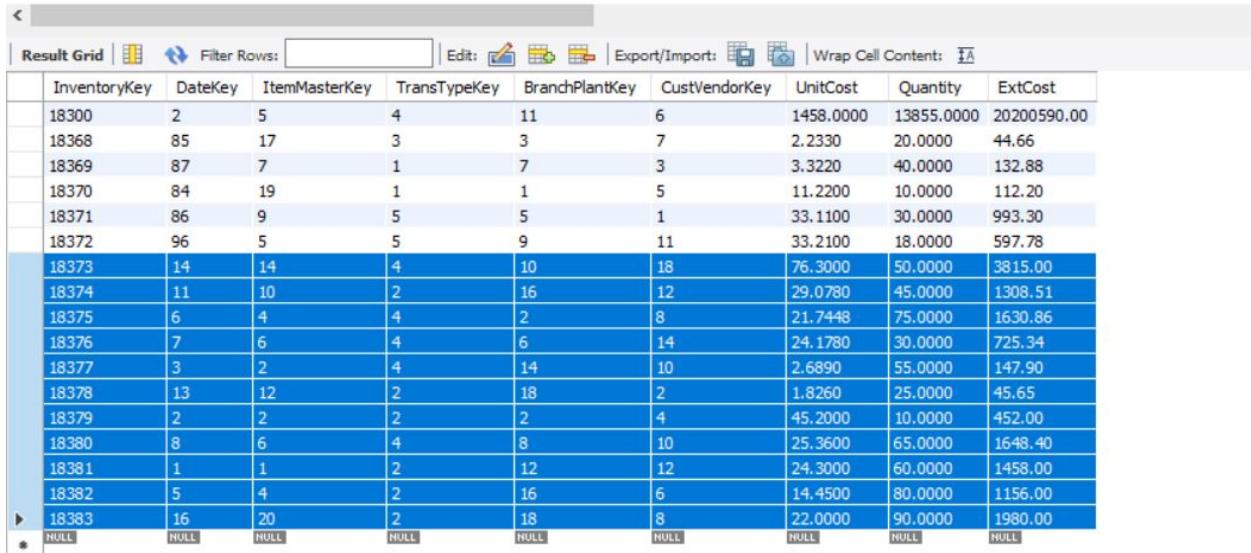
	InventoryKey	DateKey	ItemMasterKey	TransTypeKey	BranchPlantKey	CustVendorKey	UnitCost	Quantity	ExtCost
	18289	73	14	4	7	1	772.0000	9626.0000	7431272.00
	18290	70	6	4	18	18	628.0000	10734.0000	6740952.00
	18291	64	20	4	17	5	599.0000	14737.0000	8827463.00
	18292	9	11	4	3	13	1201.0000	6874.0000	8255674.00
	18293	4	1	4	20	12	604.0000	14153.0000	8548412.00
	18294	66	15	4	20	11	941.0000	9224.0000	8679784.00
	18295	82	5	4	4	16	1108.0000	11692.0000	12954736.00
	18296	52	15	4	11	10	620.0000	13557.0000	8405340.00
	18297	88	10	4	14	15	1036.0000	5396.0000	5590256.00
	18298	50	11	4	13	15	961.0000	12025.0000	11556025.00
	18299	66	6	4	8	10	919.0000	11717.0000	10767923.00
	18300	2	5	4	11	6	1458.0000	13855.0000	20200590.00
	18368	85	17	3	3	7	2.2330	20.0000	44.66
	18369	87	7	1	7	3	3.3220	40.0000	132.88
	18370	84	19	1	1	5	11.2200	10.0000	112.20
	18371	86	9	5	5	1	33.1100	30.0000	993.30
	18372	96	5	5	9	11	33.2100	18.0000	597.78
*	HULL	HULL	HULL	HULL	HULL	HULL	HULL	HULL	HULL

[Access source 1](#)



After running the transformation

11 rows added



The screenshot shows a 'Result Grid' window from the Pentaho Data Integration (Kettle) interface. The grid displays 11 rows of data with the following columns:

	InventoryKey	DateKey	ItemMasterKey	TransTypeKey	BranchPlantKey	CustVendorKey	UnitCost	Quantity	ExtCost
	18300	2	5	4	11	6	1458.0000	13855.0000	20200590.00
	18368	85	17	3	3	7	2.2330	20.0000	44.66
	18369	87	7	1	7	3	3.3220	40.0000	132.88
	18370	84	19	1	1	5	11.2200	10.0000	112.20
	18371	86	9	5	5	1	33.1100	30.0000	993.30
	18372	96	5	5	9	11	33.2100	18.0000	597.78
	18373	14	14	4	10	18	76.3000	50.0000	3815.00
	18374	11	10	2	16	12	29.0780	45.0000	1308.51
	18375	6	4	4	2	8	21.7448	75.0000	1630.86
	18376	7	6	4	6	14	24.1780	30.0000	725.34
	18377	3	2	4	14	10	2.6890	55.0000	147.90
	18378	13	12	2	18	2	1.8260	25.0000	45.65
	18379	2	2	2	2	4	45.2000	10.0000	452.00
	18380	8	6	4	8	10	25.3600	65.0000	1648.40
	18381	1	1	2	12	12	24.3000	60.0000	1458.00
	18382	5	4	2	16	6	14.4500	80.0000	1156.00
▶	18383	16	20	2	18	8	22.0000	90.0000	1980.00
*	HULL	HULL	HULL	HULL	HULL	HULL	HULL	HULL	HULL

Conclusion

During this project, I learnt how to extract the data from different types of sources (Excel and access sources), clean and filter the data based on conditions, using different step transformations available at the pentaho environment (excel / access input; filter rows; data validators; selected values; split fields) and then I was able to merge the data transformed into the data warehouse after matching the records with the values existing in the dimensions's tables and also calculating the measures using steps (merge join, sort table, calculator, table input).

And now the data is ready for the final step "the data mining".