# Airline case study: data transformation

José Luis Puente Bodoque

January, 2023

**Abstract**

Pentaho Data Integration is for applying a ETL process on data input in order to populate a data warehouse. In this report we describe in depth the logic behind each of those processes on the airline case study.

# Contents

# 1 Building dimension tables

The table 1 shows a reasonable correspondence between the data source and dimension tables.

| Input file name | Dimension table |
|---:|:---|
| airport.csv | airport |
| airport_city_state.csv | |
| fare.csv | fare |
| channel.csv | payment_channel |
| customer.csv | passenger |
| flight.csv | airplane |
| hour.csv | hour |

**Table 1:** Correspondence between input and dimension

Below we step into the explanation of each ETL process.

## 1.1 Airport

The figure 1 shows the transformation performed to build the airport dimension table.
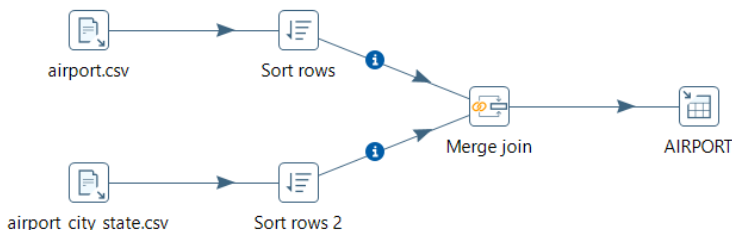


**Figure 1:** ETL structure of airport dimension

Taking a look `airport.csv` we notice state field is missing although we does find in `airport_city_state.csv`. In order to fetch state field to the main stream we use **Merge join** step with data coming from the two CSV input files. We set join type as INNER so that only rows having the same *city* keys in both sources be included in the result (figure 2). If we are careful to read the web documentation on Merge Join step we will run into a note that says: "Input rows are
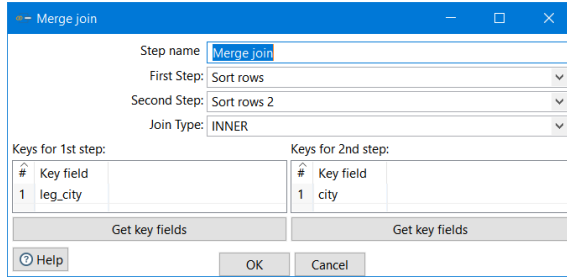
**Figure 2:** Merge Join step

expected to be sorted on the specified key fields". Thus we have to include one **Sort Row** step before each Merge Join step.

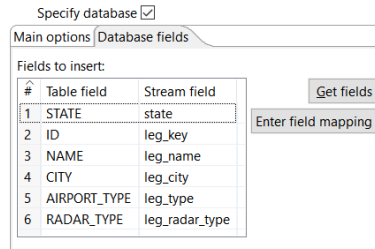The figure 3 shows the field mapping between stream fields and table fields.



**Figure 3:** Field mapping of airport dimension

## 1.2 Fare

The figure 4 shows the ETL process performed to build the fare dimension table. It consists simply of the extraction and data load of input data into the output table.



**Figure 4:** ETL strcture of fare dimension

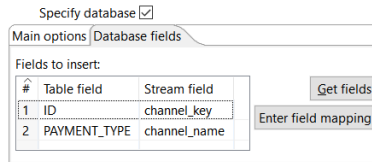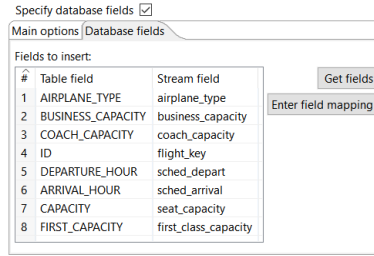The figure 5 shows the field mapping between stream fields and table fields.

**Figure 5:** Field mapping of fare dimension

## 1.3 Hour

The figure 6 shows the ETL process performed to build the hour dimension table. It consists simply of the extraction and data load of input data into the output table..



**Figure 6:** ETL structure of hour dimension

`hour.csv` is a handmade file to save all the possible combinations of hour and moment of day. It is not provided but it's essential for finding arrival and departure time identifiers later. The figure 7 shows the header and the first 15 rows as sample.



**Figure 7:** Sample of `hour.csv`

The figure 8 shows the field mapping between stream fields and table fields.

**Figure 8:** Field mapping of hour dimension

## 1.4  Payment Channel

The figure 9 shows the ETL process performed to build the payment channel dimension table. It consists simply of the extraction and data load of input data into the output table.



channel.csv                    PAYMENT_CHANNEL

**Figure 9:** ETL structure for payment channel dimension

The figure 10 shows the field mapping between stream fields and table fields.



**Figure 10:** Field mapping of payment channel dimension

## 1.5  Passenger

The figure 11 shows the ETL process performed to build the passenger dimension table. It consists simply of the extraction and data load of input data into the output table.

**Figure 11:** ETL structure of passenger dimension

The figure 12 shows the field mapping between stream fields and table fields.



**Figure 12:** Field mapping of passenger dimension

## 1.6 Airplane

The figure 13 shows the ETL process performed to build the airplane dimension table. It consists simply of the extraction and data load of input data into the output table.



**Figure 13:** ETL structure of airplane dimension

The figure 14 shows the field mapping between stream fields and table fields.

**Figure 14:** Field mapping of airplane dimension

# 2 Building the fact table

`frequentflyer.csv` file aims to be the input of the fact table. However there's an issue with the file content that in fact we talked about in the first report.

## 2.1 Issue with the input

`frequentflyer.csv` lacks the departure and arrival hour keys. The keys to look for are either in `hour.csv` or in the hour dimension table. This drives us to evaluate two possible solutions for loading the data.

## 2.2 ETL load process

The first option to load the data consists of using the **Database Lookup** step (figure 15) and the second one, the **Stream Lookup** one (figure 16). Both works fine but are quite slow. The structure presents a neck bottle in the last output table step that slows down the pace of implementation and delays the process up to 5 minutes. The figure 17 shows the metrics of the transformation. May the first option is more convenient in business environments because the keys we need are always in a certain dimension table and not necessarily in the data source.



**Figure 15:** First option

**Figure 16:** Second option



**Figure 17:** Step metrics for fact table

The figure 18 shows how we have looked up departure and arrival hours on flight keys.



**Figure 18:** First Stream lookup step

The figure 19 shows how we have looked up the id_arrival key on hour key and sched_arrival key.

**Figure 19:** Second Stream lookup 2

The figure 20 shows how we have looked up the id_depart key on both hour key and sched_depart key.



**Figure 20:** Third Stream lookup 3

**Add Sequence** step generates an incremental sequence of integer values for the surrogate key of the fact table..

Finally the figure 21 shows the mapping between the fields coming from the main stream and fact table fields.

**Note:** the input table has 7257 rows so we must change the commit size to 7257 or greater to avoid an error.

**Figure 21:** Field mapping of fact table

# 3 Ailine data warehouse

This is how our data warehouse for the airline case study looks like.



**Figure 22:** Fact table

**Figure 23:** Airport dimension table



**Figure 24:** Airplane dimension table



**Figure 25:** Passenger dimension table

**Figure 26:** Fare dimension table



**Figure 27:** Payment Channel dimension table



**Figure 28:** Hour dimension table