**A new Stratebi white paper**

www.stratebi.com

Aug 2013

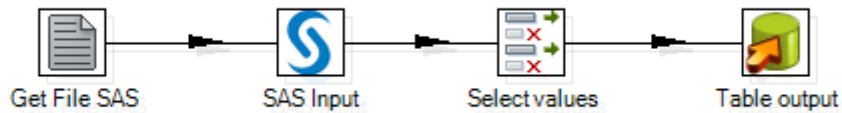# Pentaho & SAS: Getting data from SAS and exploit it into Pentaho

In this post we try to unveil the capabilities of the new Pentaho Data Integration SAS Input step. This new feature was included in the latest stable version of PDI (4.4) and is very useful for those corporations which use SAS as corporative Business Analytics tool and want to exploit the information into Pentaho BI Suite.

This new add-on is an evidence of Pentaho strategy focused on expanding their products and tools with new capabilities

Our main goal in this document is reading a SAS file using PDI. In a second stage we will make use of the information read using AgileBI plug-in included in PDI.
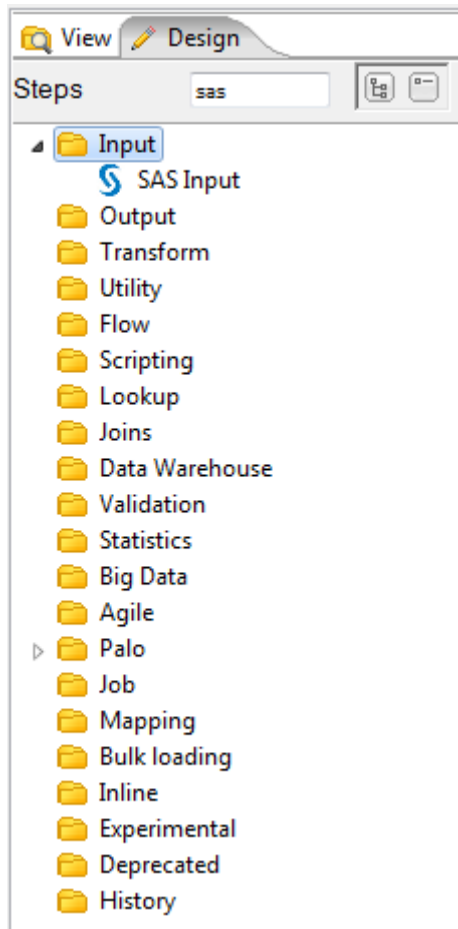
We are going to use as sample data the results of *"Estimating and modelling relative survival using SAS"* research carried out by Paul Dickman on 2004. Dickman's study includes Finnish patients diagnosed with colon carcinoma between 1975 and 1994. Here is the full document for detailed reference http://biostat3.net/download/sas/readme.pdf
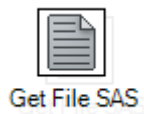
Sample data contains:

- Age at diagnosis
- Date of diagnosis
- Date of exit
- Month of diagnosis
- Sex
- Clinical stage at diagnosis
- Anatomical subsite of tumour
- Survival time in completed months
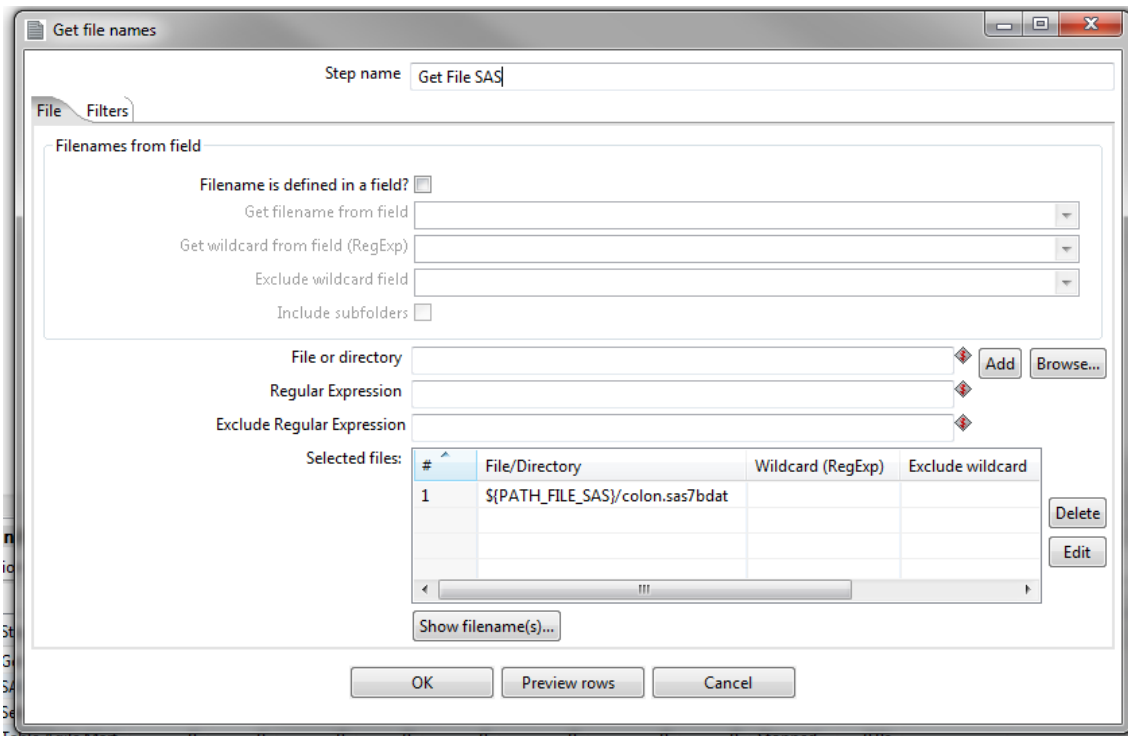- Survival time in completed years

Due to the fact that we are working at data level this new feature is available a transformation step. Below is a screenshot with the new SAS Input step.
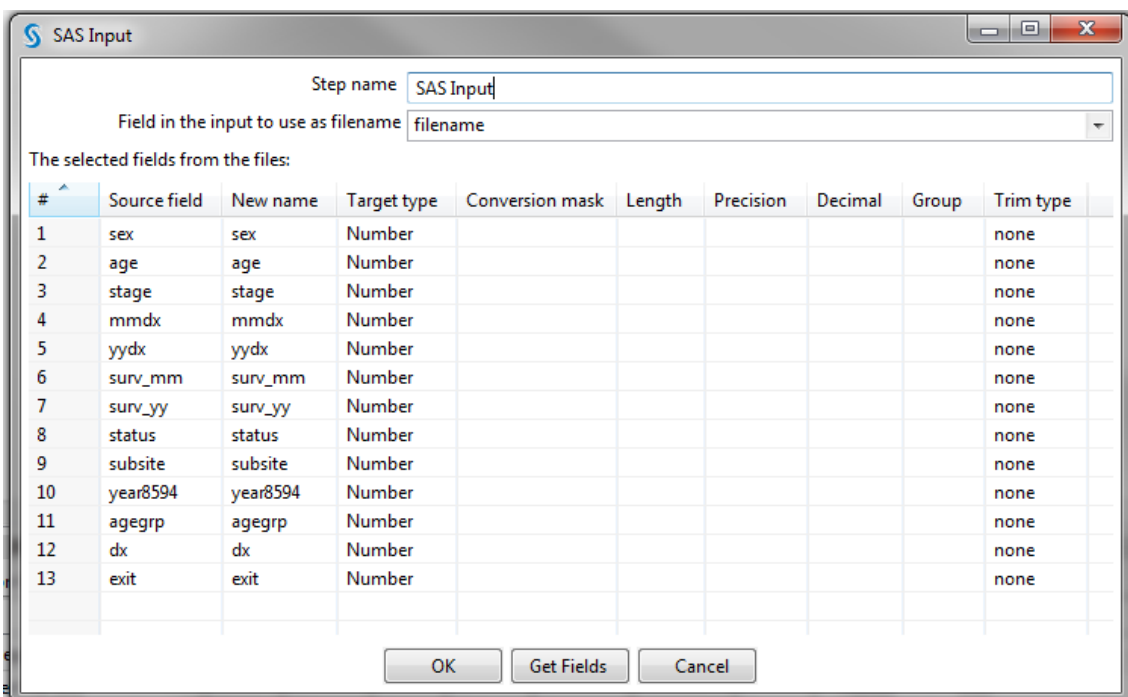
Below are listed the Pentaho Data Integration components used in this document:

- **Get File Names:** This step is used to indicate the path and name of SAS source file.
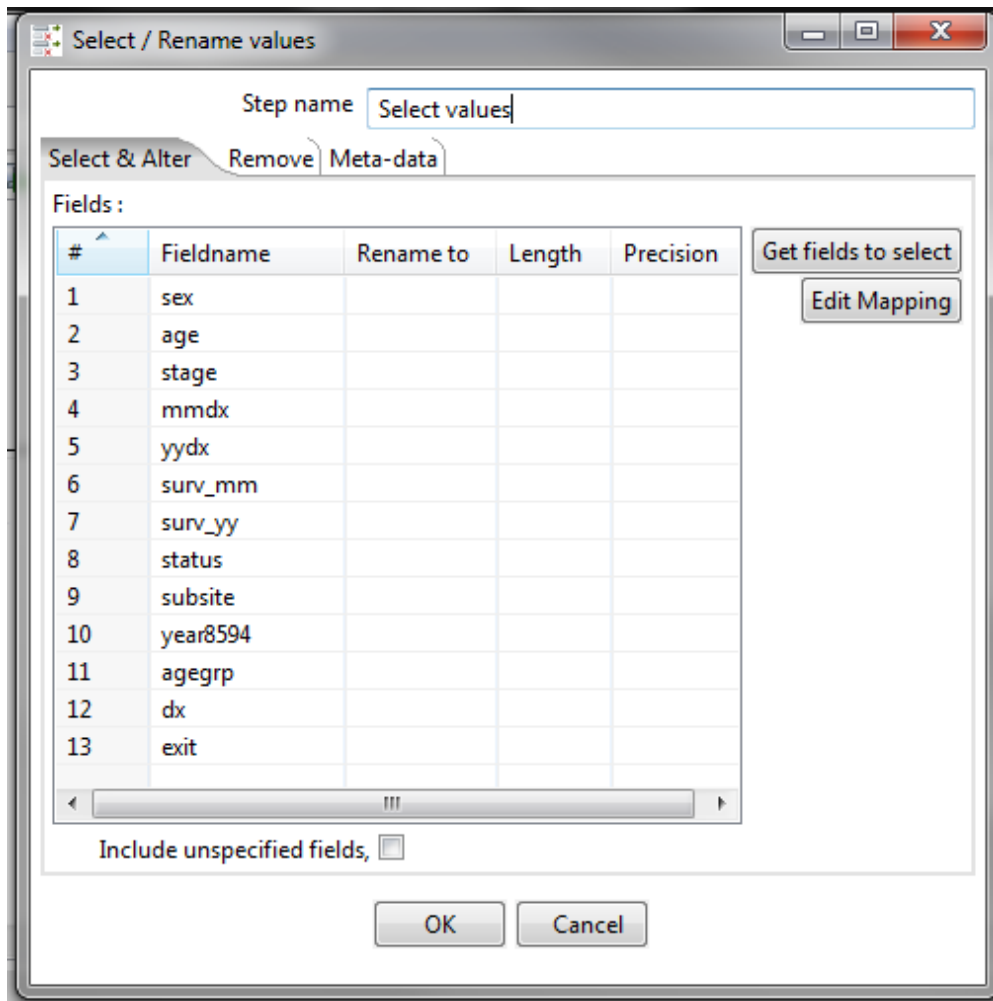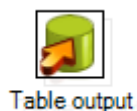


Get File SAS

- **SAS Input:** At the initial stage we should click on Get Fields button to retrieve the name and properties of the fields included in the SAS file. Then if we consider it necessary we could change the properties identified by PDI.
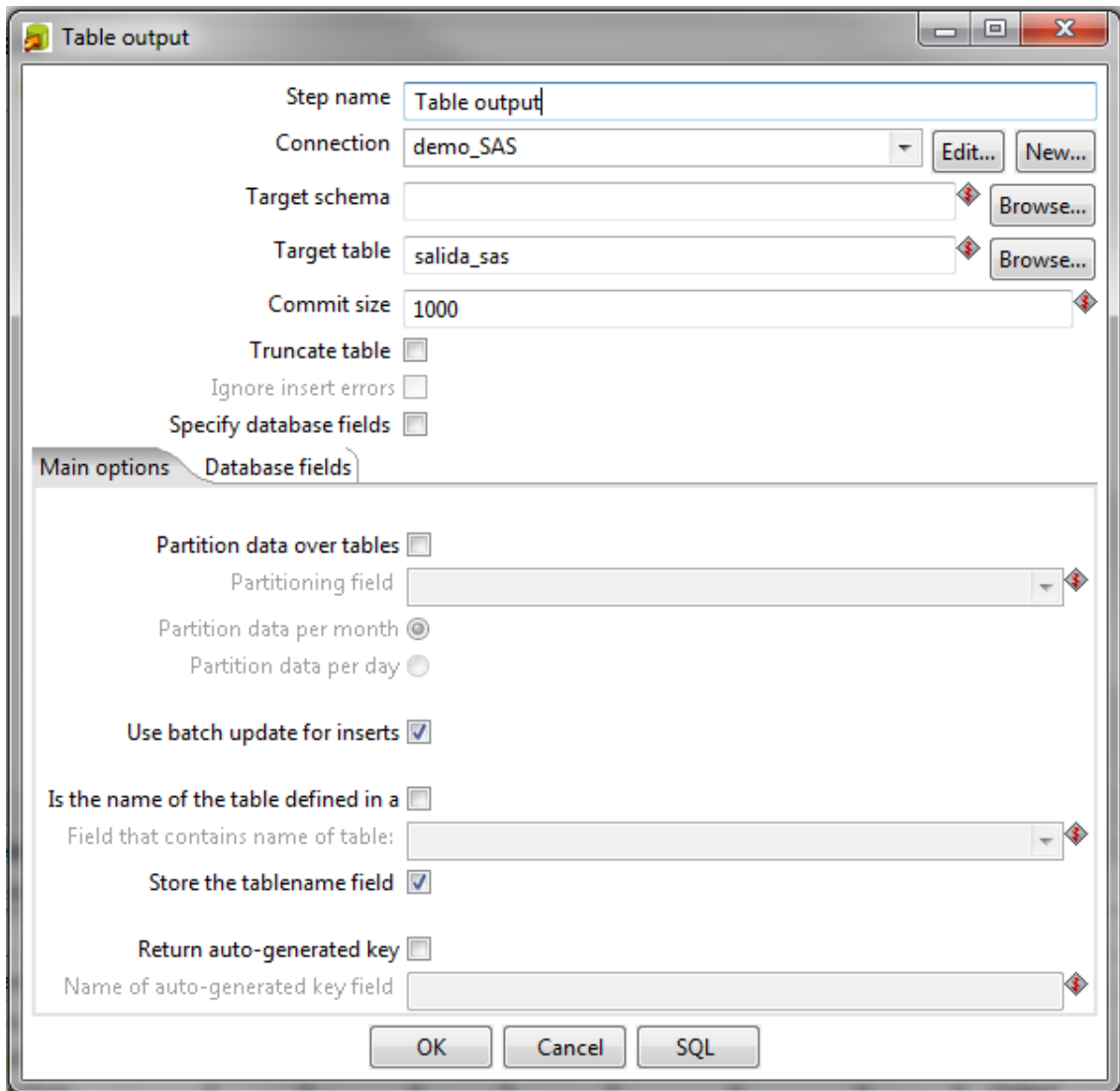


SAS Input



| # | Source field | New name | Target type | Conversion mask | Length | Precision | Decimal | Group | Trim type |
|---|---|---|---|---|---|---|---|---|---|
| 1 | sex | sex | Number | | | | | | none |
| 2 | age | age | Number | | | | | | none |
| 3 | stage | stage | Number | | | | | | none |
| 4 | mmdx | mmdx | Number | | | | | | none |
| 5 | yydx | yydx | Number | | | | | | none |
| 6 | surv_mm | surv_mm | Number | | | | | | none |
| 7 | surv_yy | surv_yy | Number | | | | | | none |
| 8 | status | status | Number | | | | | | none |
| 9 | subsite | subsite | Number | | | | | | none |
| 10 | year8594 | year8594 | Number | | | | | | none |
| 11 | agegrp | agegrp | Number | | | | | | none |
| 12 | dx | dx | Number | | | | | | none |
| 13 | exit | exit | Number | | | | | | none |

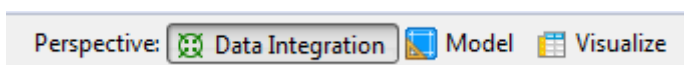- **Select values:** In this step we filter the fields from the stream.


Select values



- **Table output:** Finally we save the sample data into a MySQL table. If the table doesn't exist PDI provides a SQL button which creates an auto-generated SQL code ready to create the table.


Table output

Once we have data stored in a table, we will easily visualize them using PDI. At the top right there is a Perspective selection toolbar. Up to now during the ETL design process the perspective selected was Data Integration option.



By selecting Model option we could observe at a glance the data of our model in an OLAP view (Analysis tab), besides it is also available the option of building a report with a wizard (Reporting tab).
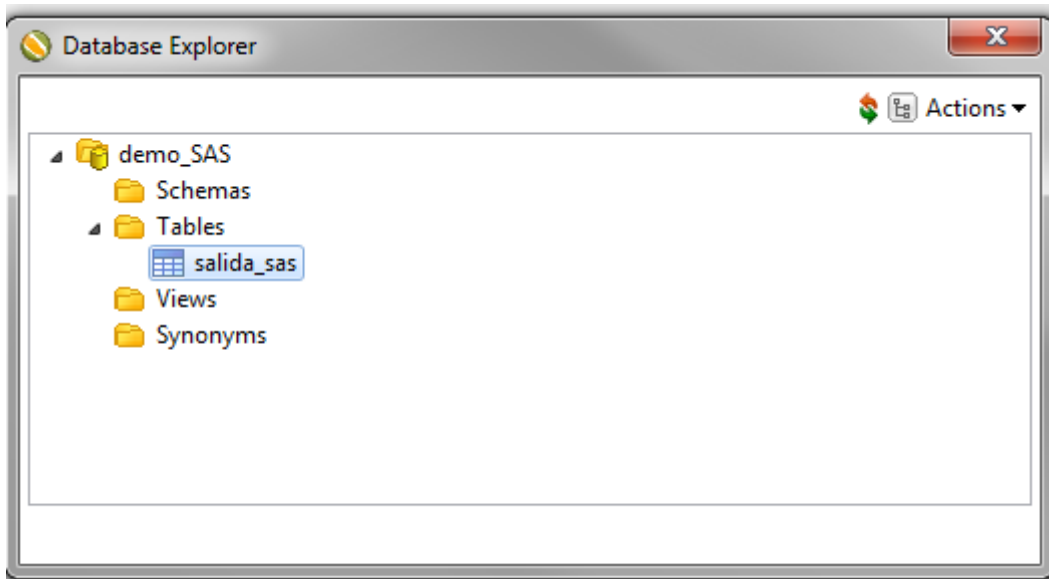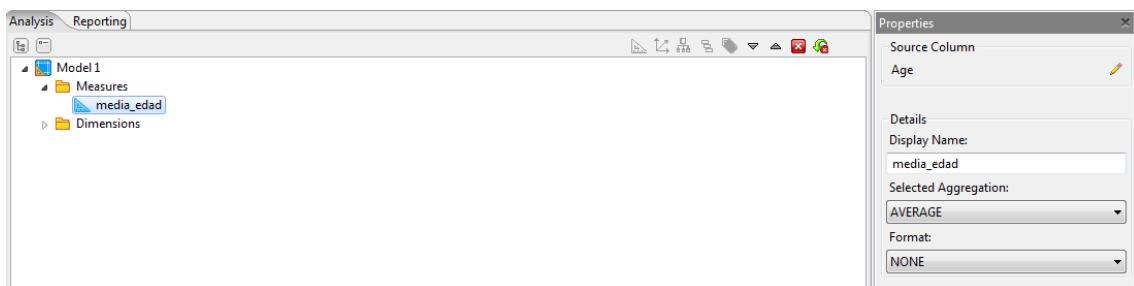
Analysis perspective



Reporting perspective

First we will assign a Model Name to our model and define a Datasource. Our data origin will be the MySQL table previously defined.
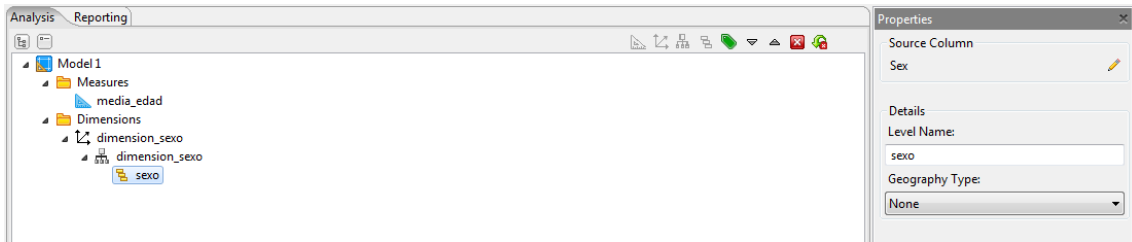


Next, having selected the Analysis scene we start designing a multidimensional structure. It is necessary to define at least one measure and one dimension to make use of the information in an OLAP view.
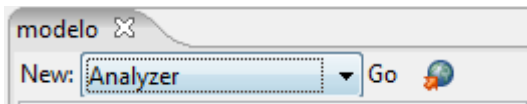
We define average age measure (media_edad) with an aggregation function of average. We select the field saved into the database table named age, as is evident with this indicator we could analyze the average age of the population.



In a similar way we proceed to create a dimension containing the sex of each patient. This dimension will only include one hierarchy and one level called sex and the table field is also named sex in the database.
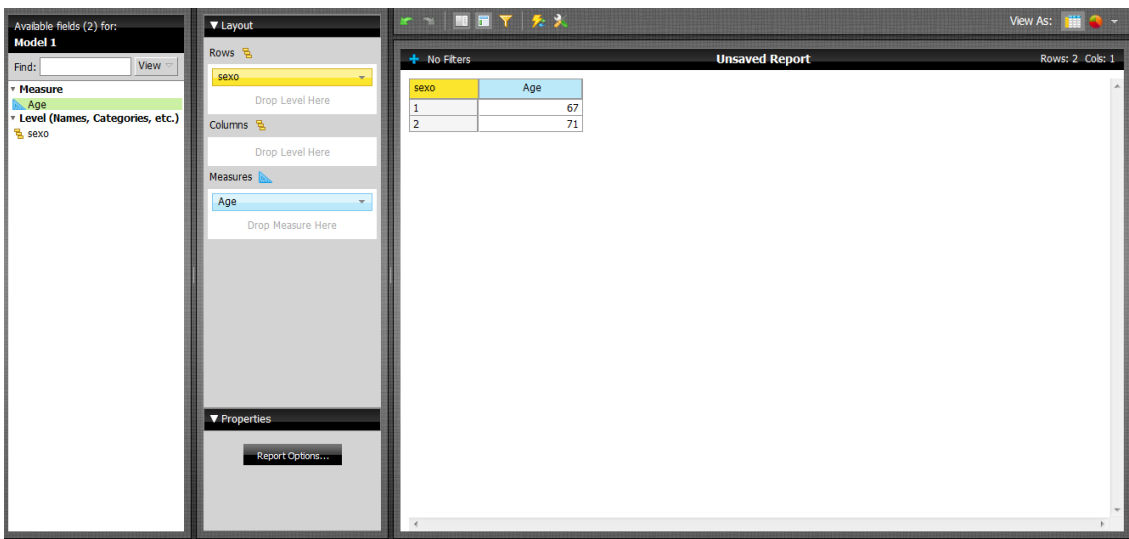
In order to visualize the OLAP structure we have developed we should only click on Go button to launch Visualize perspective.
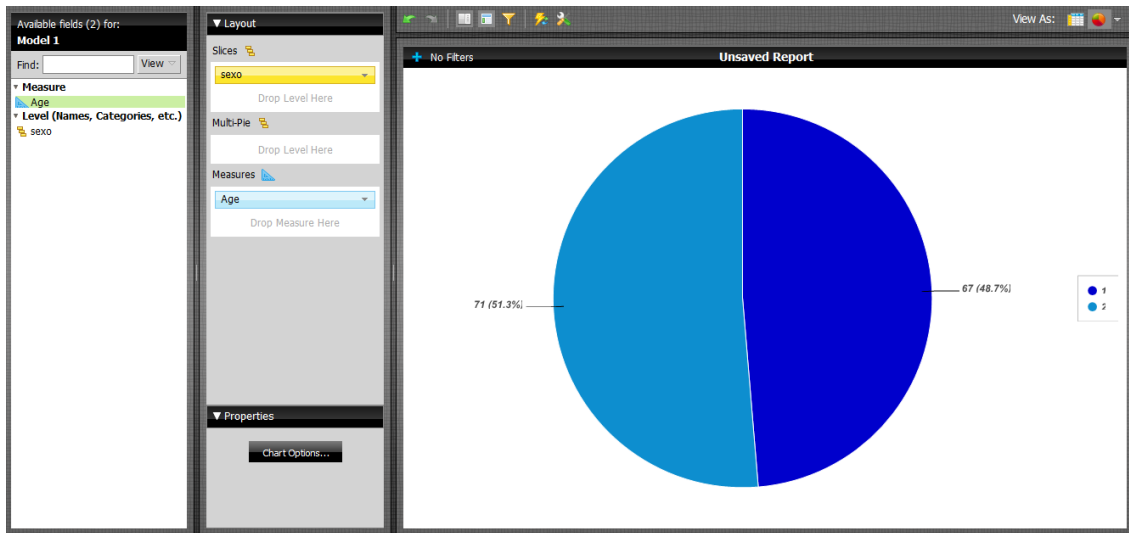


Situated on this perspective we could drag and drop into rows and columns the measures and dimensions we have created in the previous stage. At the top right of the tool exists a *View as* toolbar to switch between table and chart format.

Table format
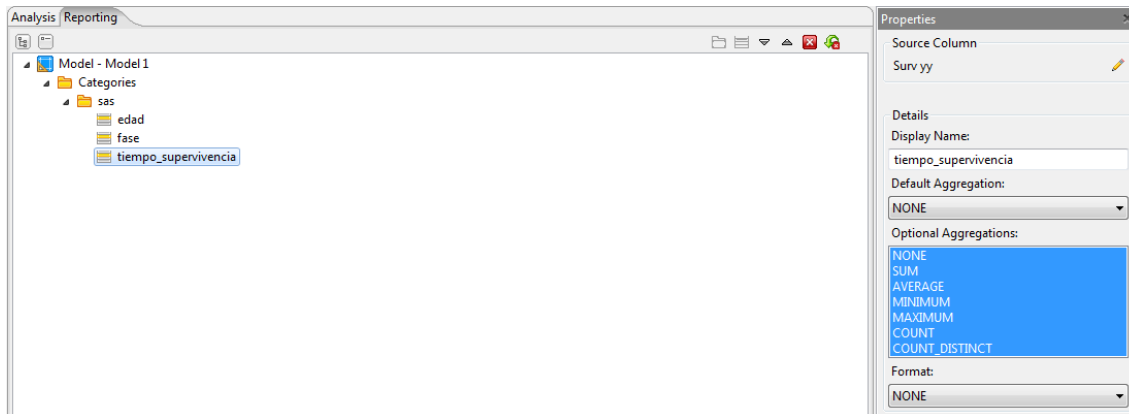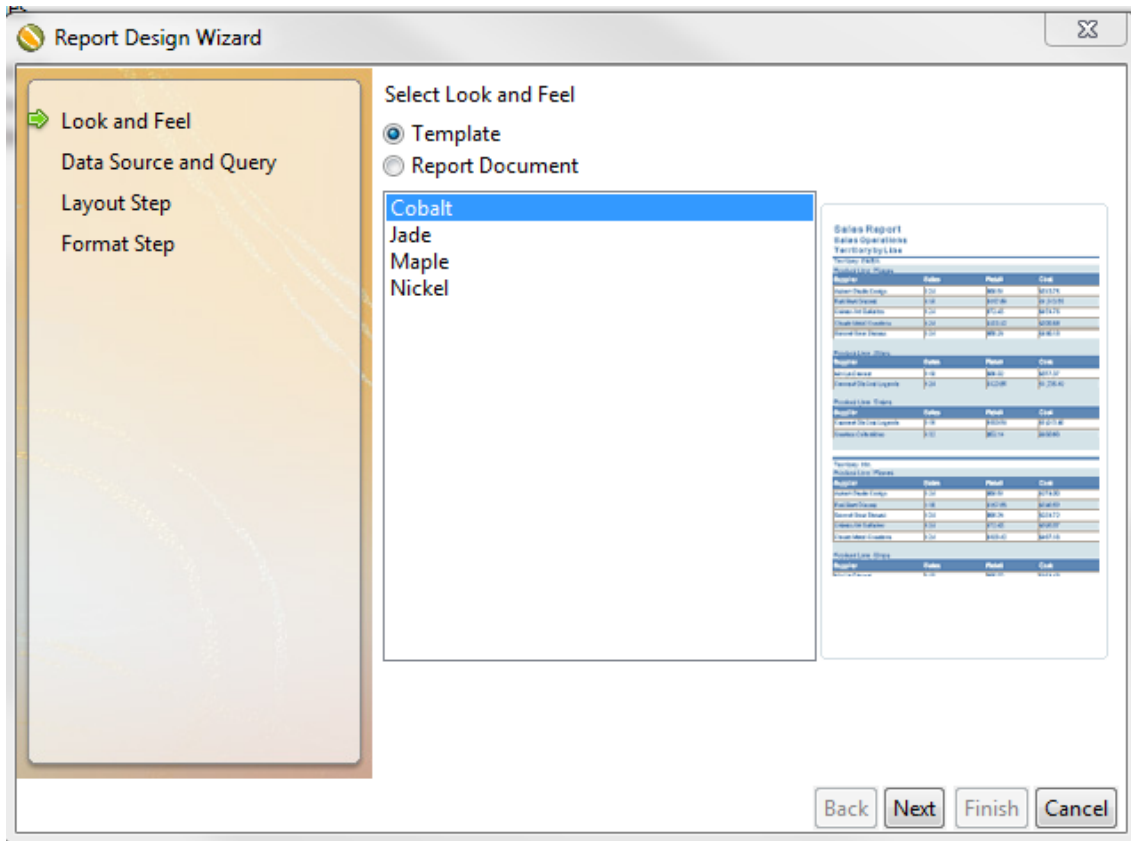


Chart format

Then, we move to define a reporting metadata structure. At first, it is required the creation of a category and select the fields that we want to have available in further stages. These are the fields chosen:
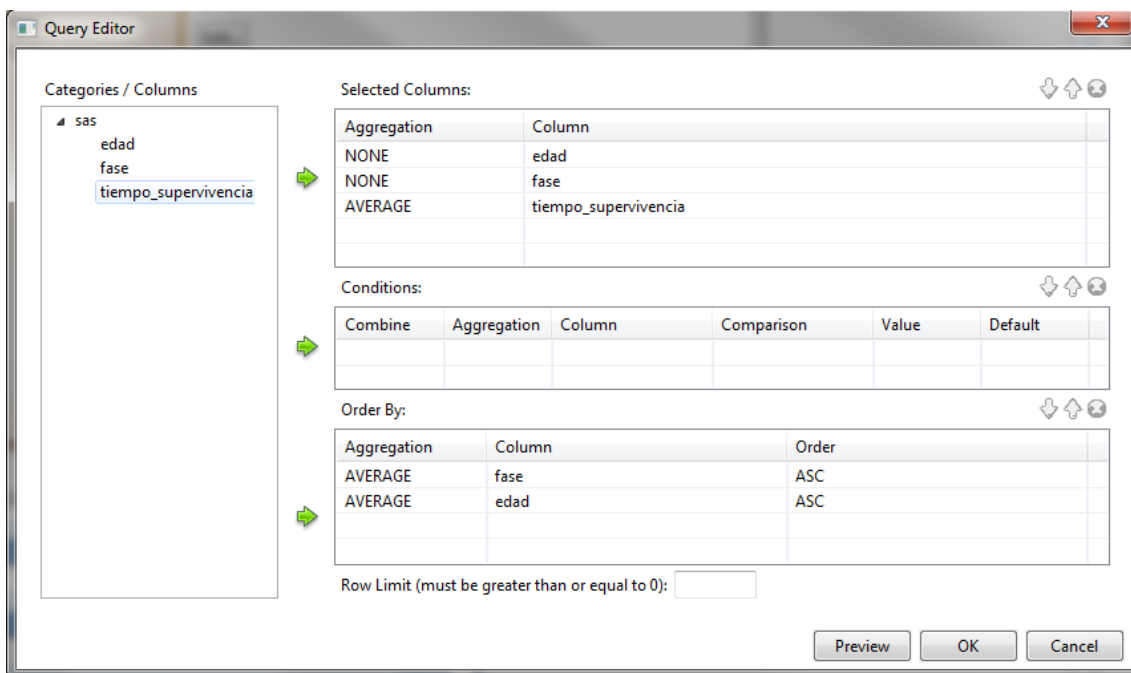
- Age (edad) → age field in database
- Clinical stage at diagnosis (fase) → stage field in database
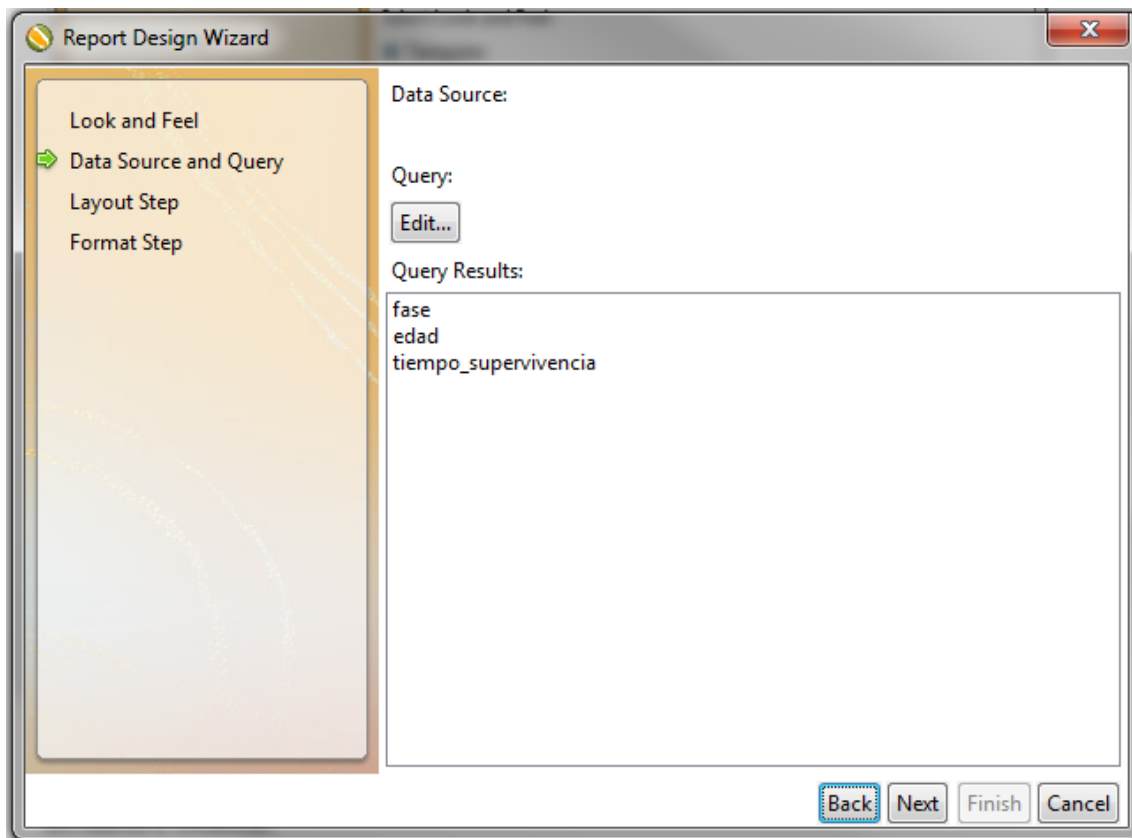- Survival time in completed years (tiempo_superviviencia) → surv_yy field in database



In order to gain knowledge of the metadata we have just generated, we select Report Wizard and click on Go button. Automatically we will be directed to Visualize perspective, here we have to choose a template.
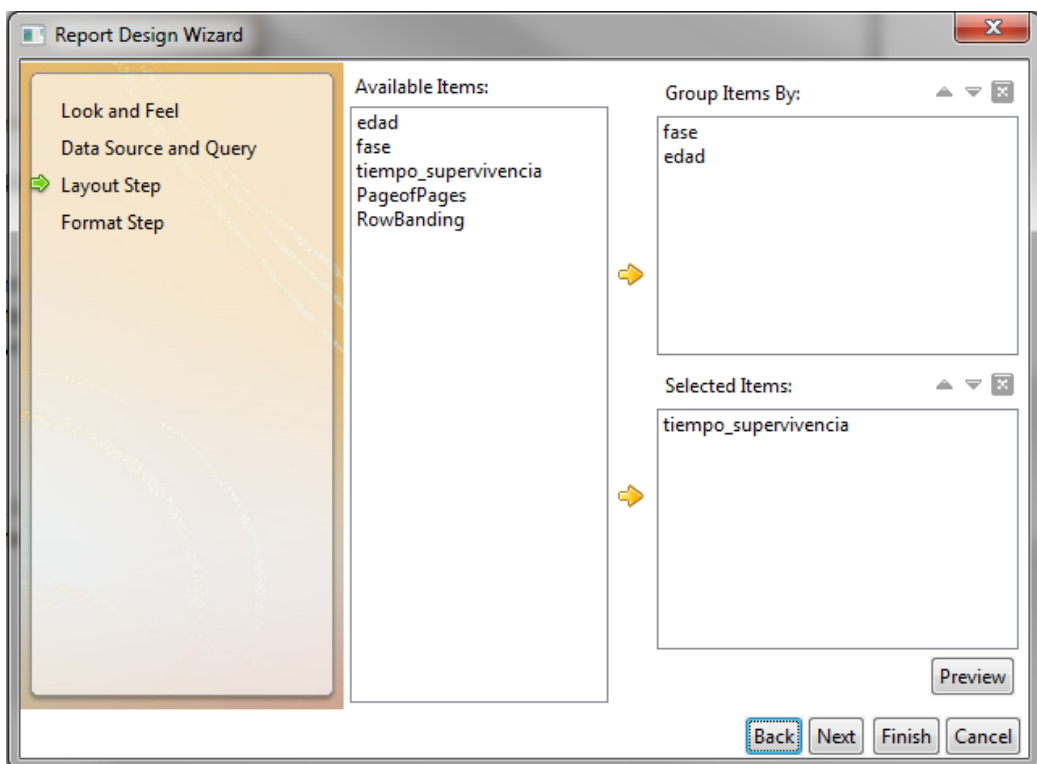
Next, we select the fields we want to show in the report. We choose age (edad), clinical stage at diagnosis (fase) and the average of survival time in completed years (tiempo_supervivencia) sorted in ascending order by clinical stage and patient age.
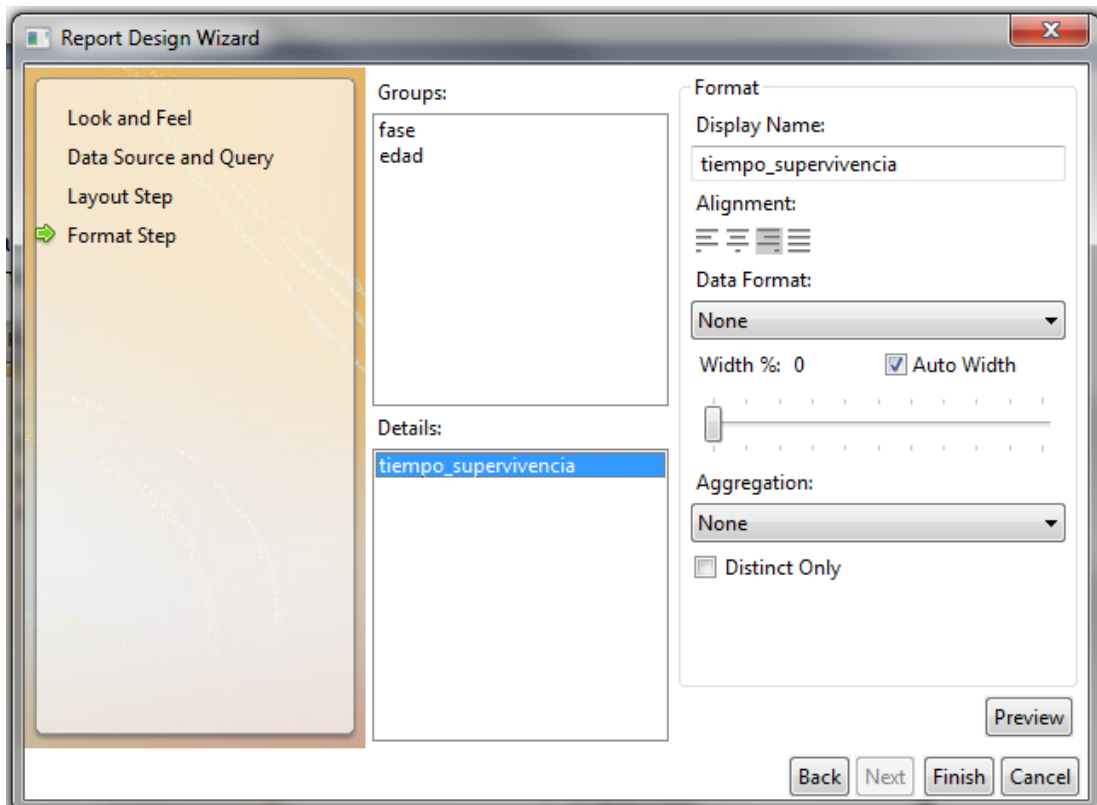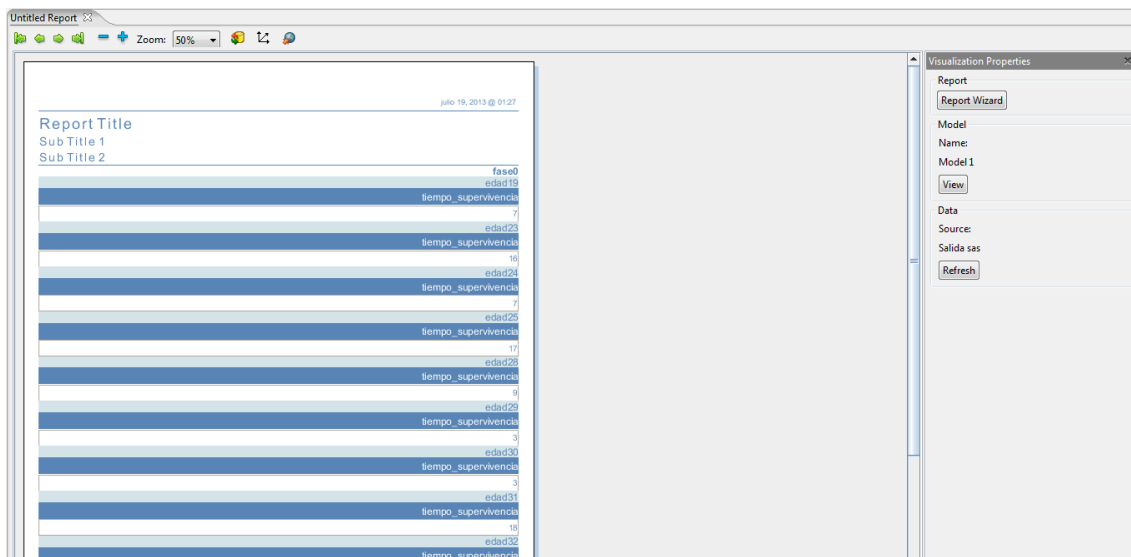
Then, we design the layout of the report. Our main objective is to show the results grouped by clinical stage and patient age.



After that we could change format properties such as alignment, width, format strings ...

Finally, we have our report finished in less than 5 minutes. It is possible to edit the report by pushing Report Wizard button located in the Visualization Properties palette located at the right side of the tool.



This powerful tool allows us to save the metadata structure previously generated with .xmi extension available to reuse it in future.