# Introduction to Spoon

Spoon is the UI of the Pentaho ETL Tool or PDI.  When you first open and run it, the Default "Perspective" is "Data integration".  This is the perspective we will be working with.

The Data Integration perspective allows you to create transformations, jobs, and inspect your data allowing for iterative updates as you work.   Let us get ourselves familiar with Spoon.



1. Toolbar
2. Connect menu
3. Sub-toolbar
4. Design and View Tabs (Or sometimes called "Left pane")
5. Canvas

| Component | Name | Description |
|---|---|---|
| **1** | Toolbar | Single-click access to common actions such as create a new file, opening existing documents, save and save as. |
| **2** | Connect Menu | Create and connect to repositories for centrally storing your ETL jobs and transformations. We are not going to use this functionality. |
| **3** | Sub-toolbar | Provides buttons for quick access to common actions specific to the transformation or job such as Run, Preview, and Debug. |
| **4** | Design and View Tabs | The **Design tab** of the Explore pane provides an organized list of transformation steps or job entries used to build transformations and jobs. Transformations are created by simply dragging transformation steps from the Design tab onto the canvas and connecting them with hops to describe the flow of data.<br><br>The **View tab** of the Explore pane shows information for each job or transformation. This includes information such as available database connections and which steps and hops are used.<br><br>In the image, the Design tab is selected. |
| **5** | Canvas | Main design area for building transformations and jobs describing the ETL activities you want to perform. |

| Icon | Description |
|------|-------------|
| ⊡ | Create a new job or transformation |
| 📂 | Open transformation/job from file if you are not connected to a repository or from the repository if you are connected to one |
| ▤ | Explore the repository |
| 💾 | Save the transformation/job to a file or to the repository |
| 📝 | Save the transformation/job under a different name or file name (Save as) |
| ⊜ ▾ | Switch between the different perspectives.<br>• Data Integration — Create ETL transformations and jobs<br>• Schedule — Manage scheduled ETL activities on the Pentaho Server |
| ▷ ▾ | Run transformation/job and set run options; runs the current transformation from XML file or repository |
| ❚❚ | Pause transformation |
| ☐ | Stop transformation |

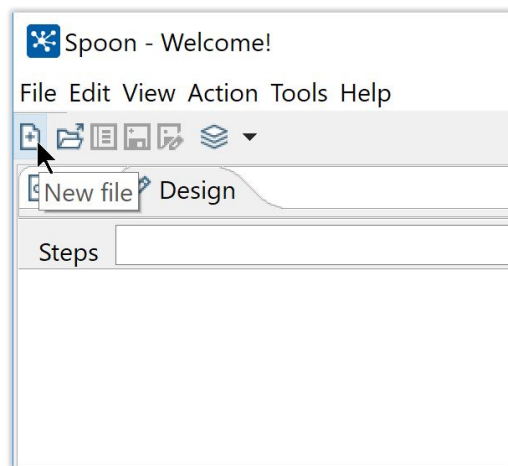| | |
|---|---|
| 👁 | Preview transformation: runs the current transformation from memory. You can preview the rows that are produced by selected steps. |
| 🐞 | Run the transformation in debug mode; allows you to troubleshoot execution errors |
| ▷ | Replay the processing of a transformation |
| 🗶 | Verify transformation |
| 🗐 | Run an impact analysis of a transformation on the database |
| 🗗 | Generate the SQL that is needed to run the transformation. |
| 🗗 | Launch the database explorer allowing you to preview data, run SQL queries etc. on the database. |
| 🗄 | Show execution results pane |
| 🔒 | Lock transformation |

# Exercise 1 - Copy Text file to Table

In this exercise, the objective is to familiarise you with basic operations of the Spoon tool in PDI. You should complete this exercise before attending the lecture.
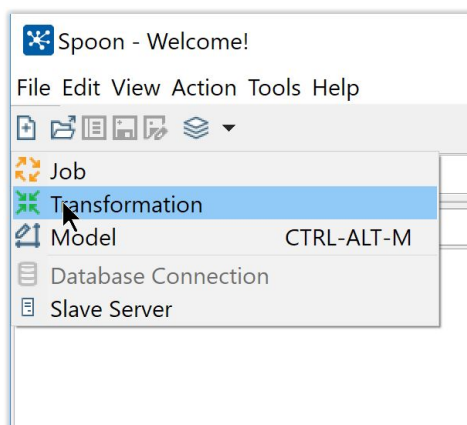
Prerequisite:
1. Start up your Postgres server from Bitnami Mapp/Wapp stack.
2. Note down the ipaddress, port number, user and password for connecting to Postgres server. Make sure you have a database named **postgres** in the server already. If not, create one.
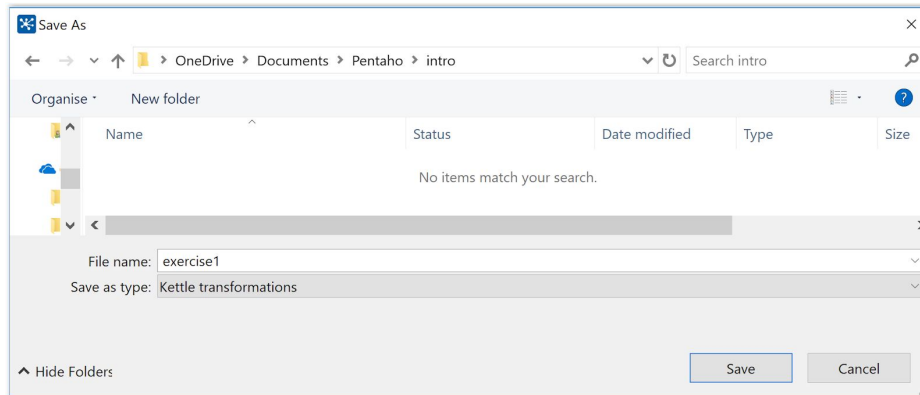
STEPS

1. Click on New file icon

2. Select Transformation from drop down list

3. Save transformation file with name "exercise1"

Transformation steps are saved with ".ktr": file extension.

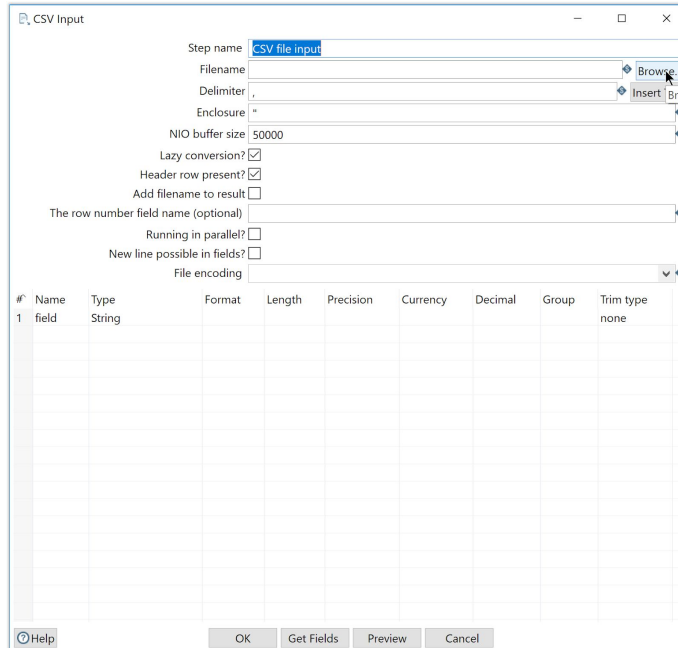4. Drag "CSV file input" from left pane to the canvas



At the left pane, click on "> Input" to reveal drop down of steps for input.
Click on "CSV file input" icon and drag.

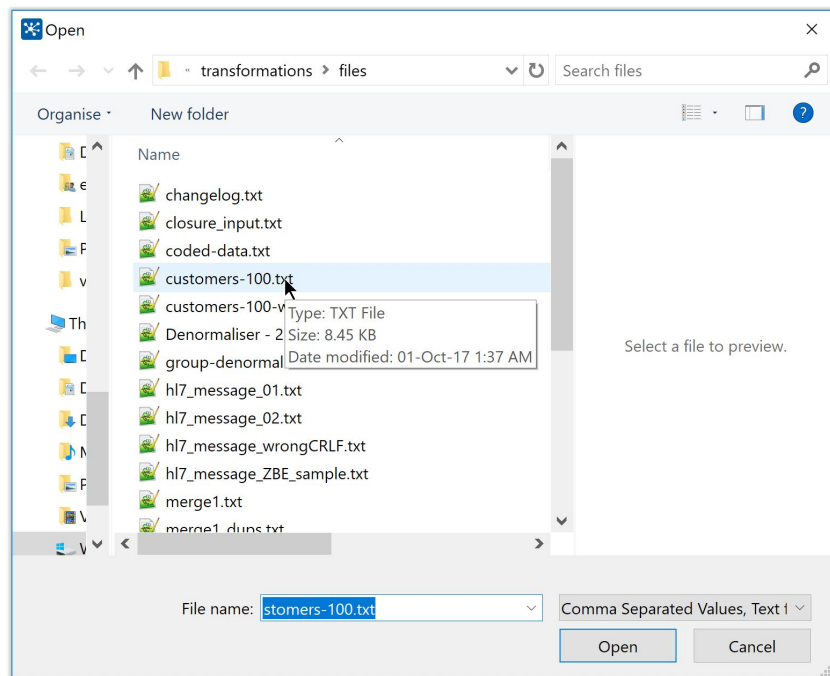You may also use the search box and enter "CSV" and the filtered list will show



Drag the icon to the canvas and release

5. Edit the "CSV file input" step by double clicking on it

Click on browse and navigate to folder data-integration/samples/transformation/files



Choose customers-100.txt and click on "Open"

Click on Get fields. You will get a new dialog box



Notice the field names are not separated into multiple lines
due to the wrong delimiter ", " vs "; "
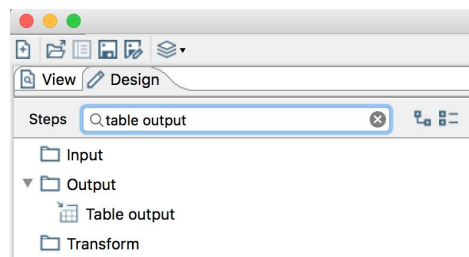
Click on "Cancel" on dialog "Sample size"



Go ahead and change the delimiter to " ; " and click on Get fields again.
The fields are now correctly detected and data type is set to String. The dialog box "Sample size" is asking if you would like to set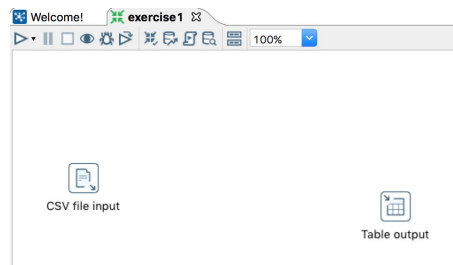 the data type automatically by sampling the first 100 rows. Click on "Cancel" for now as you want it to remain as STRING.
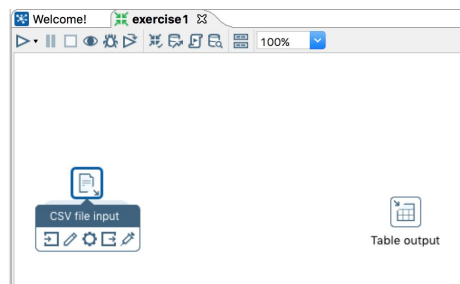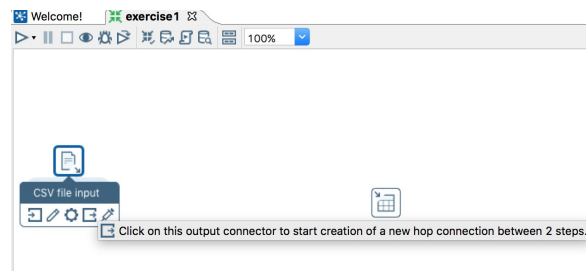
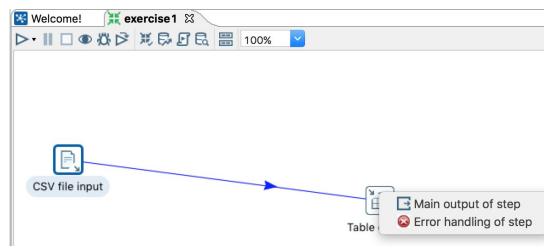6. Drag and drop the Table output step
   a.



   b.

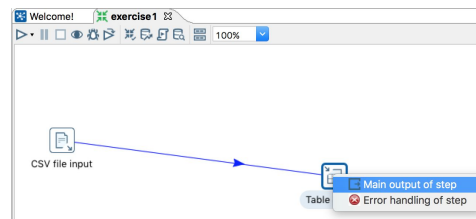7. Connect the CSV file input to Table output
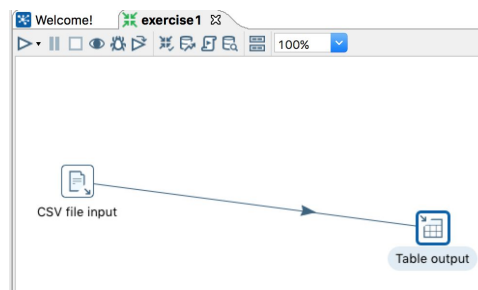


Click on the first step



Click and drag the output connector icon



Drag the arrow point over the second step



Release the mouse button and a tooltip will appear.
Click on "Main output of step"

**NOTE:**
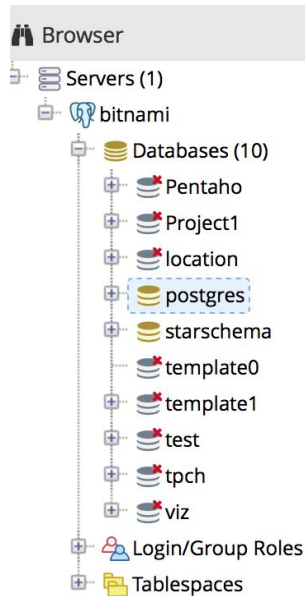Other ways to connect the 2 steps using the mouse with clickable scroll button:
Go to https://help.pentaho.com/Documentation/6.0/0L0/0Y0/030/010
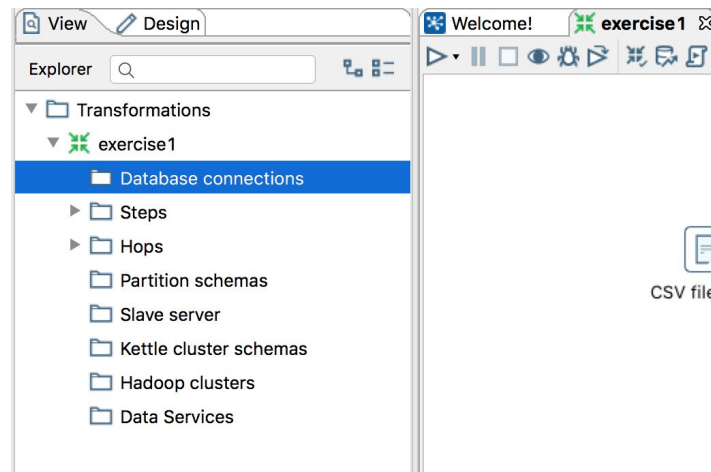Paragraph: More About Hops

1. Use <SHIFT + left-click> on first step and drag to second step

2. left-click on the source step, hold down the middle mouse button, and drag the

   hop to the second step

3. Select both steps, then right-click and choose New Hop

4. Use <CTRL + left-click> to select two steps then right-click on the step and

   choose New Hop

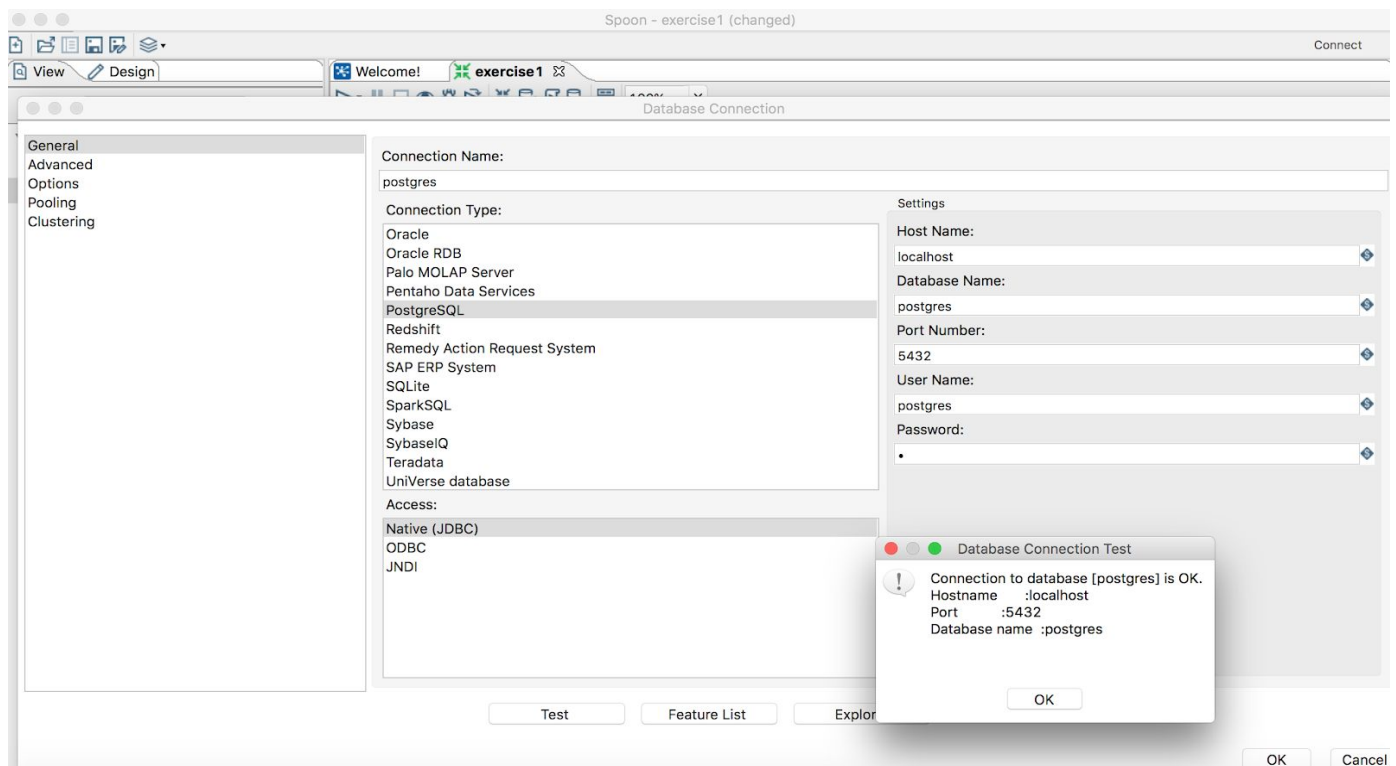8. **Create a database and connect to it**
Open pgadmin and create a database named **postgres** , if it doesn't exist already.
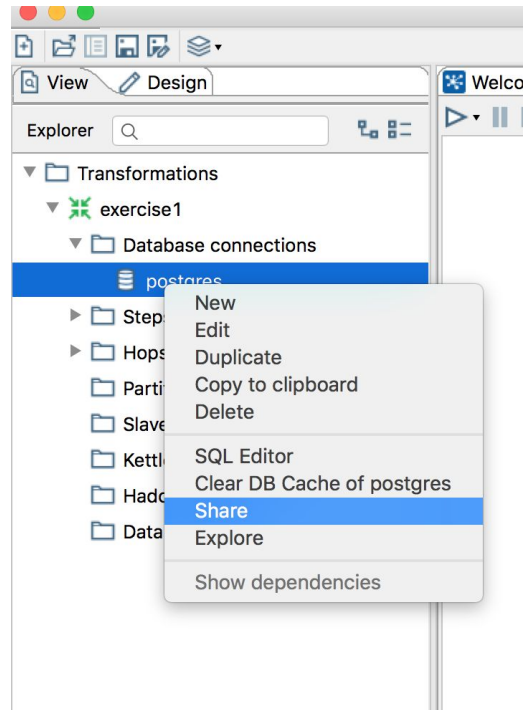


Go back to Spoon.

Click on the View tab and expand "Database connections". It should be
empty, since you have not associated any database to Spoon yet.
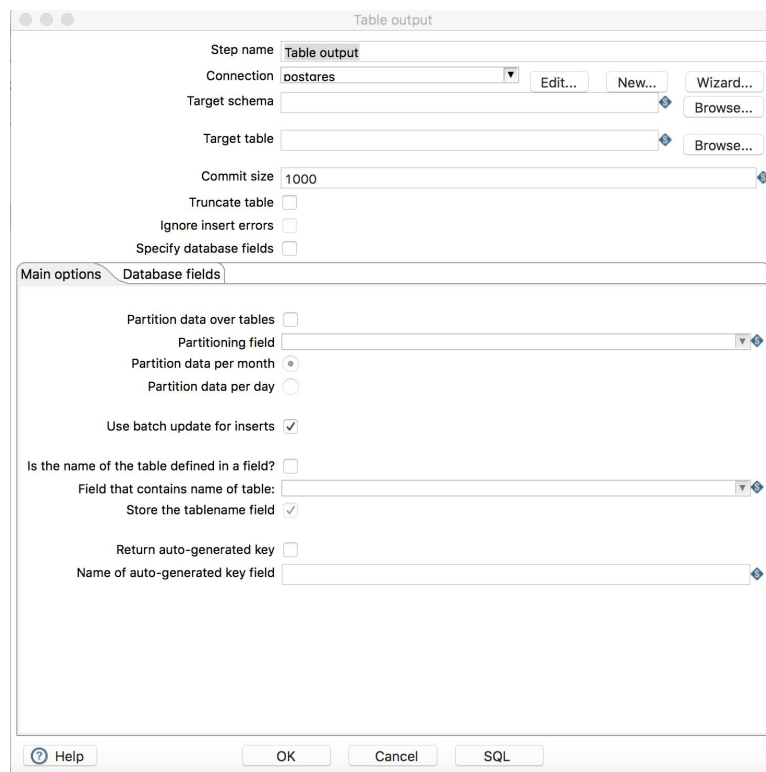


**Set the name of the connection**
**Fill in the configuration for the postgres database and click on "Test" and OK.**

***Common Errors:** You must test the connection before you use it. Since, we often put wrong
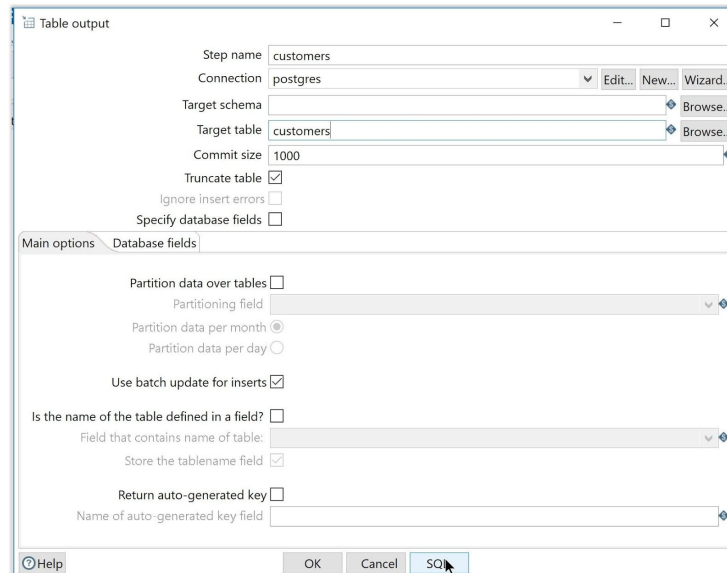password, username, wrong port etc.

Right Click on the postgres connection, and click Share. It is a good practice to Share the connection. As we create new transformations, the shared connection will automatically be taken as default connection.

9. Edit the Table output step

Change step name to "customers"



Set the table name.  Leave the other options as it is. We need to create a new table by executing a SQL Query.

Click on SQL button and a suggested SQL query dialog will appear. If the query dialog does not appear, then it might be that customers table already exists in the database. You need to delete it manually from pgadmin, if that is the case.
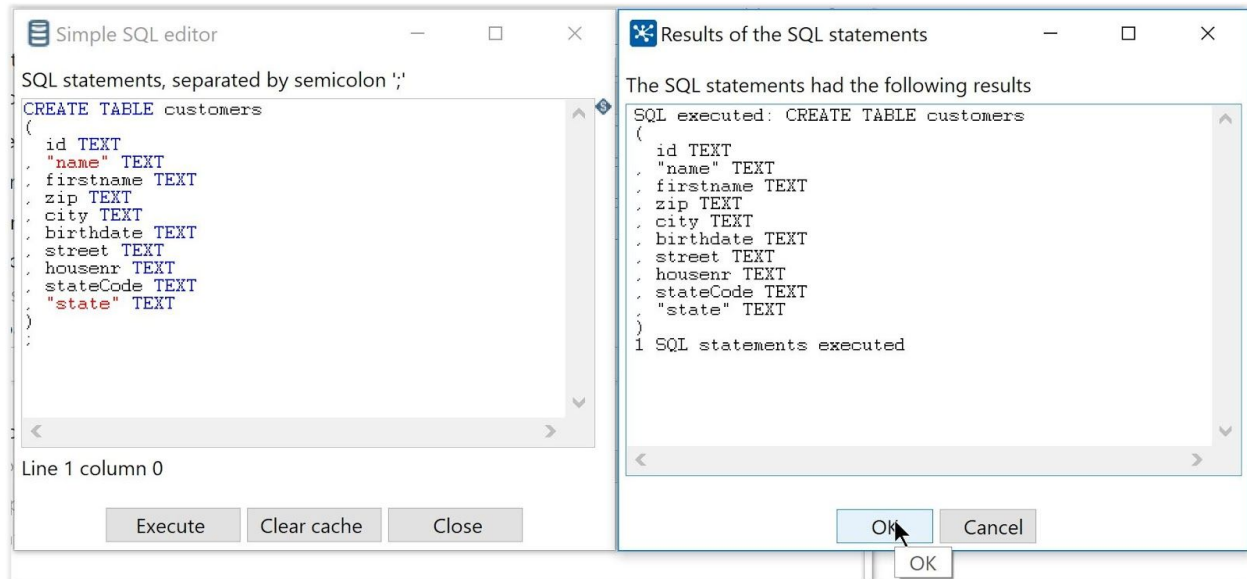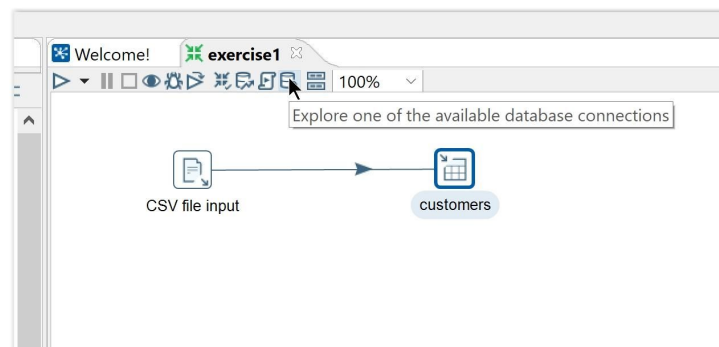


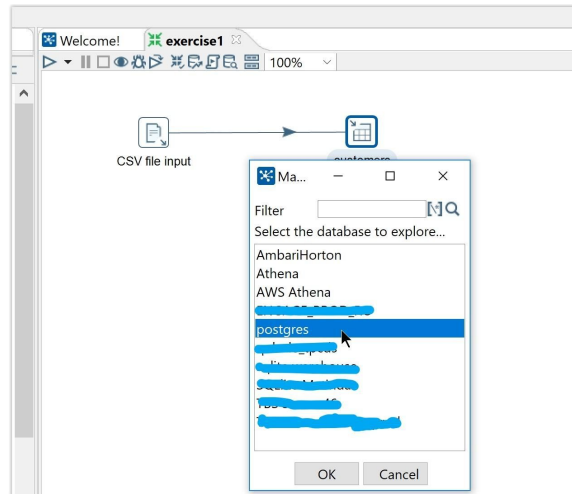Check the SQL query and edit it if necessary.  Click on "Execute"

Check to see there is no error on execution. Click OK.

**On the SQL dialog select all and copy the SQL Query and paste it in a notepad++ for your reference in future.  This is good practice.**
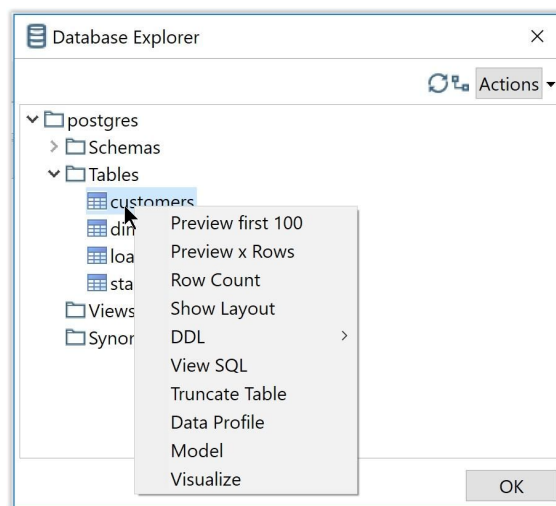
10. Check the database to see that your table is created.



Click on Database connection icon to browse the connections.

Choose the postgres connection that you have made before.

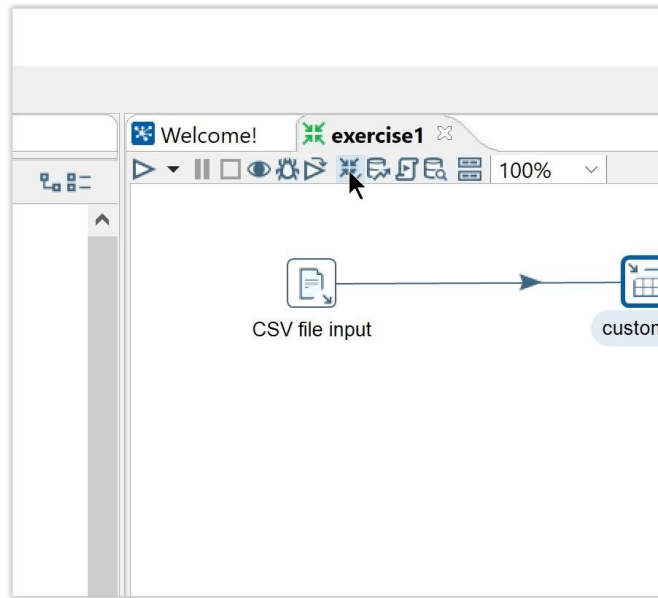Navigate to the table and right-click to view the table.
There should be **no rows** in your table as you have just created the table.
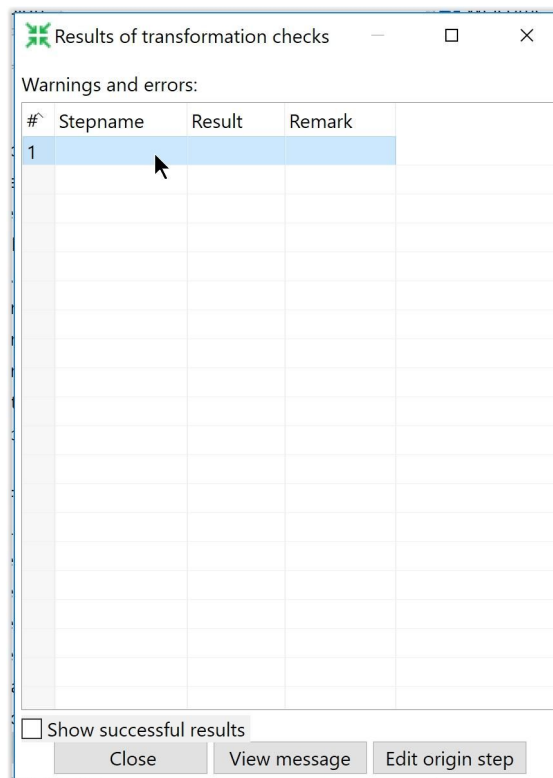
10. Verify the transformation
"Verify transformation" helps you to check that all steps are properly configured and all input fields and output fields connect from step to step.
***Common Errors:** You may encounter:
1. **Errors:** Fields from the output of a step does not exist in the input of the next step.
   Output field name of a step does not match with the input field name of the next step.
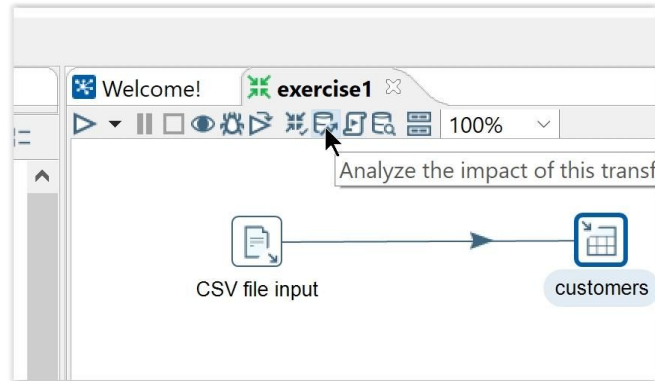2. **Warnings:** You may ignore them most of the time.

Click on the "Verify" icon



You should get a blank table, showing that no errors were found.

11. Run an impact analysis of the transformation
You want to know before running the transform what impact it has on the resulting tables
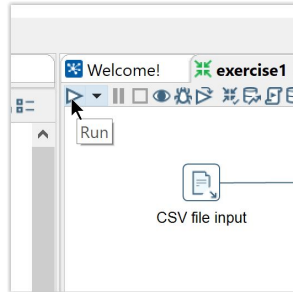
Click on "Impact analysis" Icon



**\*\*\*Important/Common Errors**:
1. The last step "customers" which is an Table output step is set to TRUNCATE the table on every run. Make sure that this is the correct setting for the transformation that you want to do. This is okay if you are updating intermediate tables (if any) such as staging table (staging tables will be explained in the lecture), but **not okay** for Dimension tables and Fact table. Loading data into Dimension and fact table MUST NOT TRUNCATE. Since we are not creating dimension or fact table here, we are choosing truncate option.
2. Note the fields that will be written into the table and the source step of the value of each field. For this simple transform is only 2 steps. In a complex transform this will help you to debug.

12. Running the transformation
There are many options on running the transformation. Most important is that Pentaho will take the ktr file ("exercise1.ktr") to run. If there are changes you have made you need to save it. The name of the transformation, "exercise1" on the tab will be bold if changes are not saved. In addition the run transformation step will ask you again to save the changes.

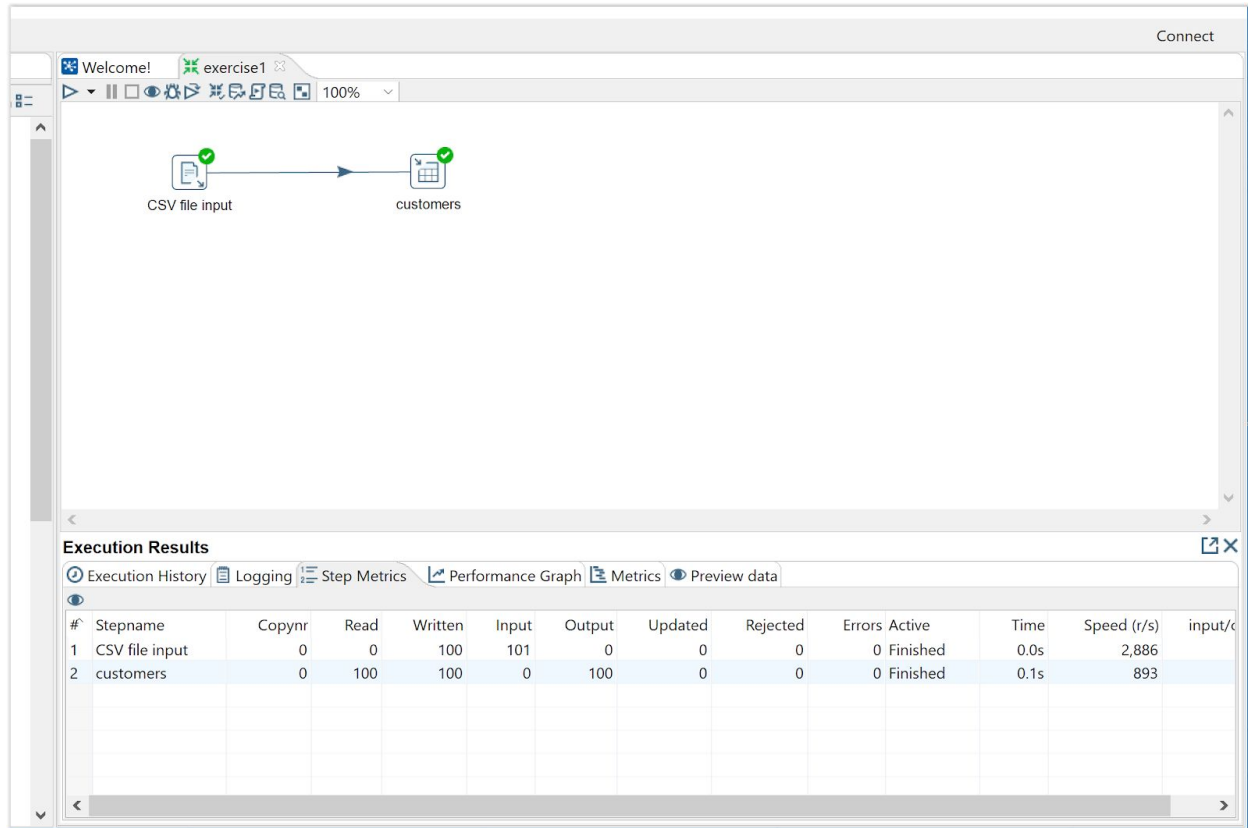Click on "Run" Icon



Click on "Run"

Note "exercise1" is BOLD , which means you have made changes but have not saved the file.

Click Yes



If you click on cancel, it will remind you to save the file first.

A new pane will appear at the bottom half of the canvas.

Most commonly used tabs are:
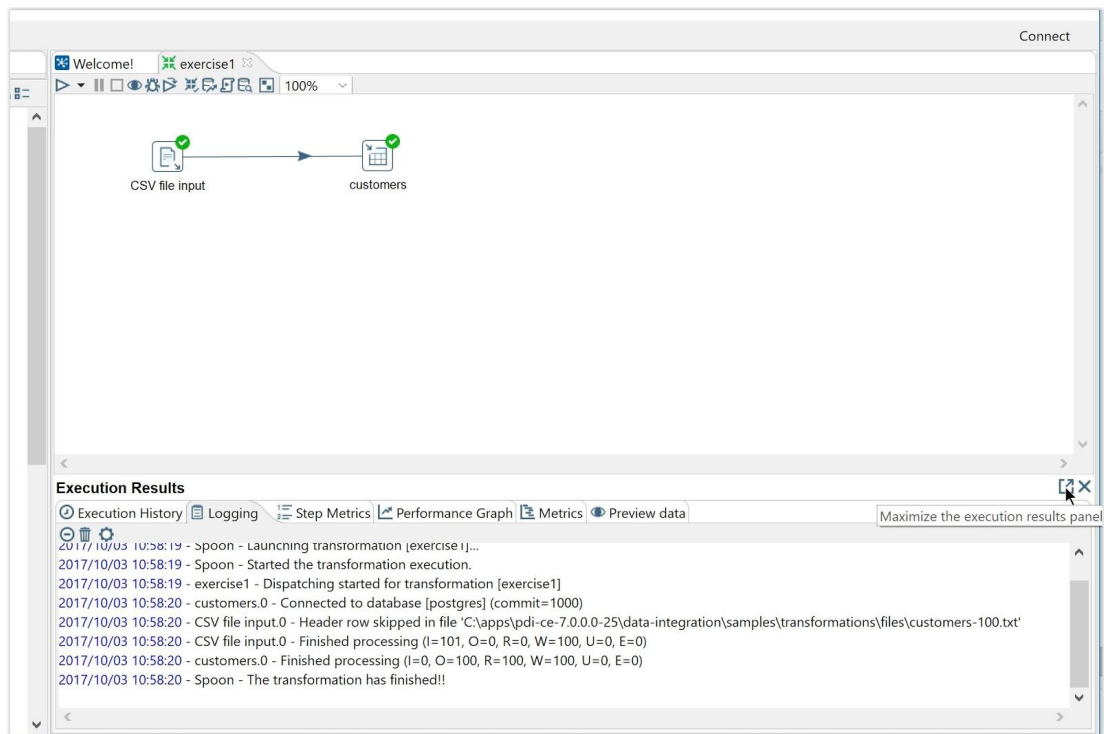1. Step metrics
2. Logging
3. Metrics

- Step metrics Tab



Every step is an object with an input and output. When the input is from internal steps it is counted as a **Read**, if the input is from outside of the transform (from a text file in this case), it is counted as an **Input** row count. When the output of a step is to another step that the rows counted are in the **Write** counter and if the output is to an external output, in this case it is a table in postgres, it is counted as an **Output**.
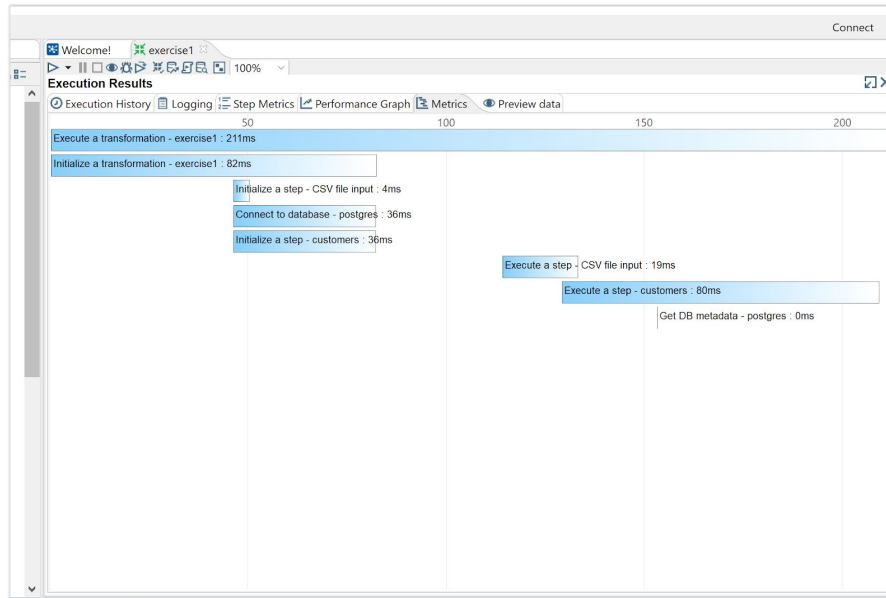
- Logging Tab



Logging tab contains the log of the run.  You change the logging level at the Run dialog.  The most detailed setting can be very verbose.  Note you can click on the full pane Icon to see the full log.
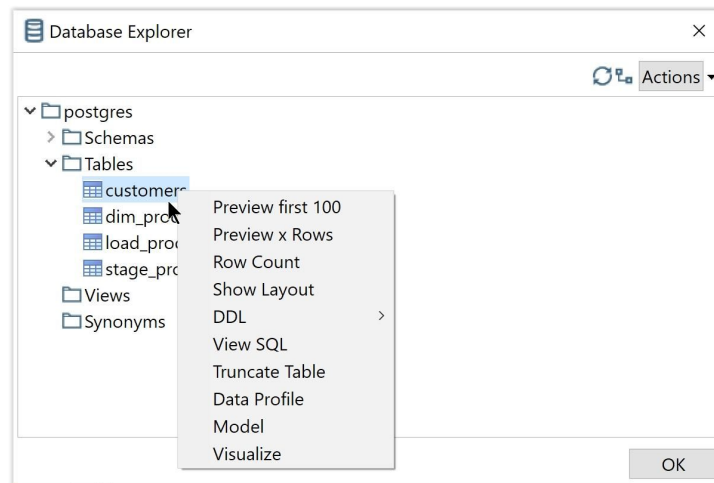


Full pane to view the logs.

● Metrics Tab



The Metrics tab shows a Gantt chart after your transformation executes. It contains information such as how long it takes to connect to a database, how much time is spent executing a SQL query, or how long it takes to load a transformation.

## 13. Check the table by browsing the table



Click on Browse Database connection Icon and navigate to customers table and right-click. Chose preview first 100 rows.

Check to see that all rows are there.

## Appendix:

Watch this 5-minute introductory youtube video. It illustrates how to join two tables in Spoon.
https://www.youtube.com/watch?v=RGmm_xXUwrM

User guide for Spoon -
http://wiki.pentaho.com/display/EAI/Spoon+User+Guide

A comprehensive documentation of various steps you can use in Spoon-
-https://wiki.pentaho.com/display/EAI/Pentaho+Data+Integration+Steps

A comprehensive documentation of various jobs you may want to use -
https://wiki.pentaho.com/display/EAI/Pentaho+Data+Integration+Job+Entries