

Data Visualization

Context

Data visualization is viewed by many disciplines as a modern equivalent of visual communication. It involves the creation and study of the visual representation of data, meaning «information that has been abstracted in some schematic form, including attributes or variables for the units of information». It is one of the steps in data analysis or data science.

The primary goal of data visualization is to communicate information clearly and efficiently via statistical graphics, plots and information graphics. Effective visualization helps users analyse and reason about data and evidence. It makes complex data more accessible, understandable and usable.

Data visualization is strategic. All its stake is to help to decision making because the future of company is on a line : data which would not be right could lead to bad decision.

In our projet, the data visualization is very important because it can help clients to visualize their electrical consumption over the year of their home.

DBeaver and Tableau

Why

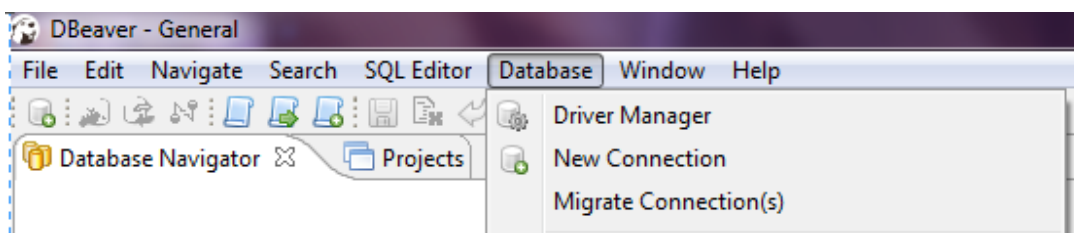
Tableau is a data visualization tool. Like we said, data visualization is very important nowadays so we decided to choose Tableau which is one of the more famous dataviz's tool.

With Tableau, The user is not dependent anymore. He can create himself his analyses and his displays to publish them then and spread them.

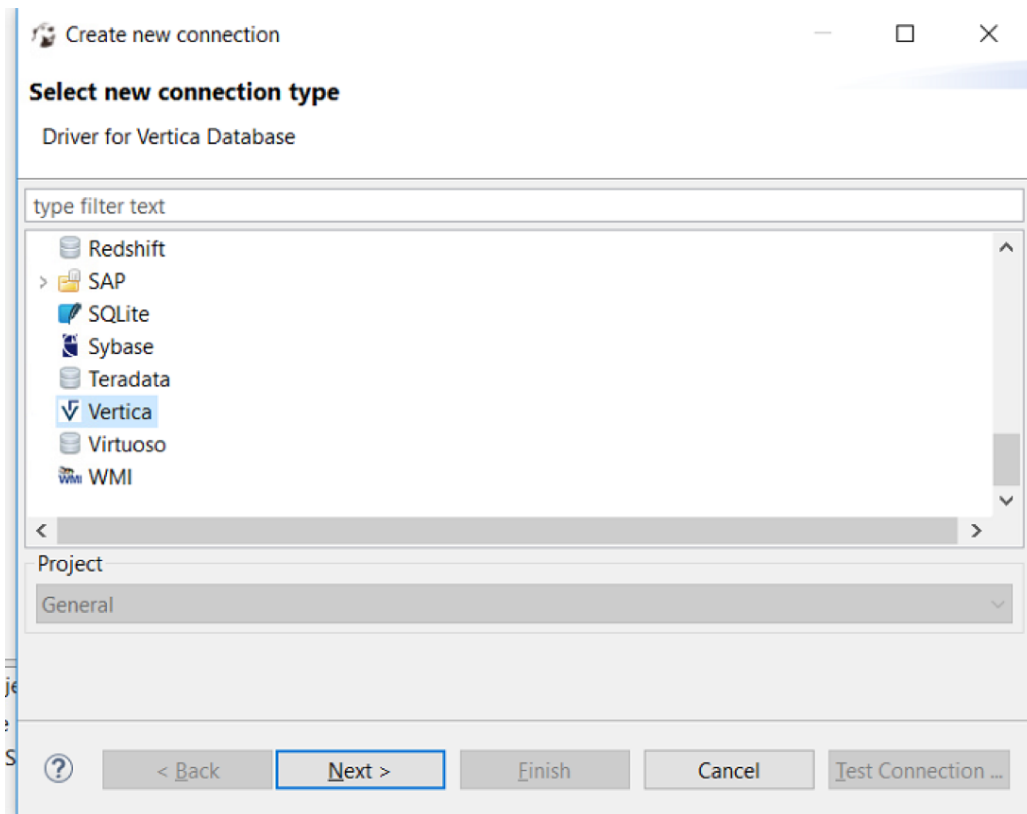
The DBeaver is an SQL client and a database administration tool that we used to connect to vertica.

To connect DBeaver to Vertica, you need to follow these:

In DBeaver, select Database > New Connection



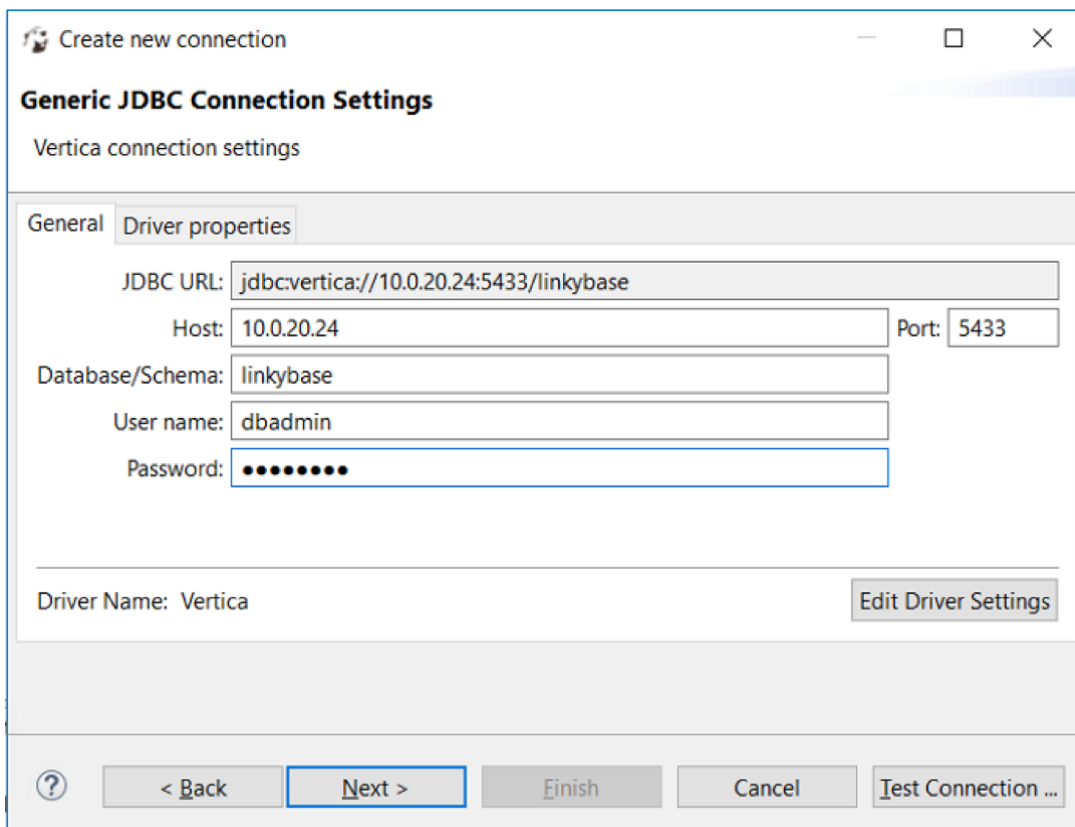
In the Create New Connection window, select Vertica and click Next.



In the general tab, enter you database credentials.

Then, in Edit Driver Settings, select add file and select the Vertica JDBC file you downloaded.

You can check if your connector is working by clicking Test Connection. If your connection to Vertica is successful, a message appears. Click on Next button next.



In Default schema, select linky_target because data are in linky_target.

Create new connection

Finish connection creation

General connection settings.

Connection name:

Connection type:

Connection folder:

Security

☒ Save password locally

Miscellaneous

☒ Show system objects

☐ Show utility objects

☐ Read-only connection

Filters

Connection

Auto-commit: ☒

Isolation level:

Default schema:

Keep-Alive:

Bootstrap queries:

From here, you can run queries and visually explore your Vertica database. When you are satisfy of your queries, you can create a view which can be used in Tableau.

DBeaver - General - [<Vertica - LinkyBase> Script]

File Edit Navigate Search SQL Editor Database Window Help

Auto Vertica - LinkyBase <None> 200

Database Navigator

- Vertica - LinkyBase
 - linky_config
 - Tables
 - linky_reject
 - linky_target
 - stream_clusters
 - stream_events
 - stream_load_specs
 - stream_lock
 - stream_microbatch_histo
 - stream_microbatch_sourc
 - stream_microbatches
 - stream_scheduler
 - stream_scheduler_histo
 - stream_sources
 - stream_targets
 - Views
 - linky_target_view
 - Indexes
 - Procedures

Project - General

Name	DataSource	Preview
SQL Scripts		
Scripts.sq Vertica - Li...	<empt>	

```
SELECT
LEFT(ID,4) AS Foyer,
TO_TIMESTAMP(CONCAT ('2017',SUBSTR(ID,5)),'YYYYMMDDHHMISS') AS DateTime,
CAST (SUBSTRING(ID, 16) AS INT) AS Conso
FROM linky_config.linky_target LIMIT 15;
```

Result

SELECT LEFT(ID,4) AS Foyer, Enter a SQL expression to filter results (use Ctrl+Space)

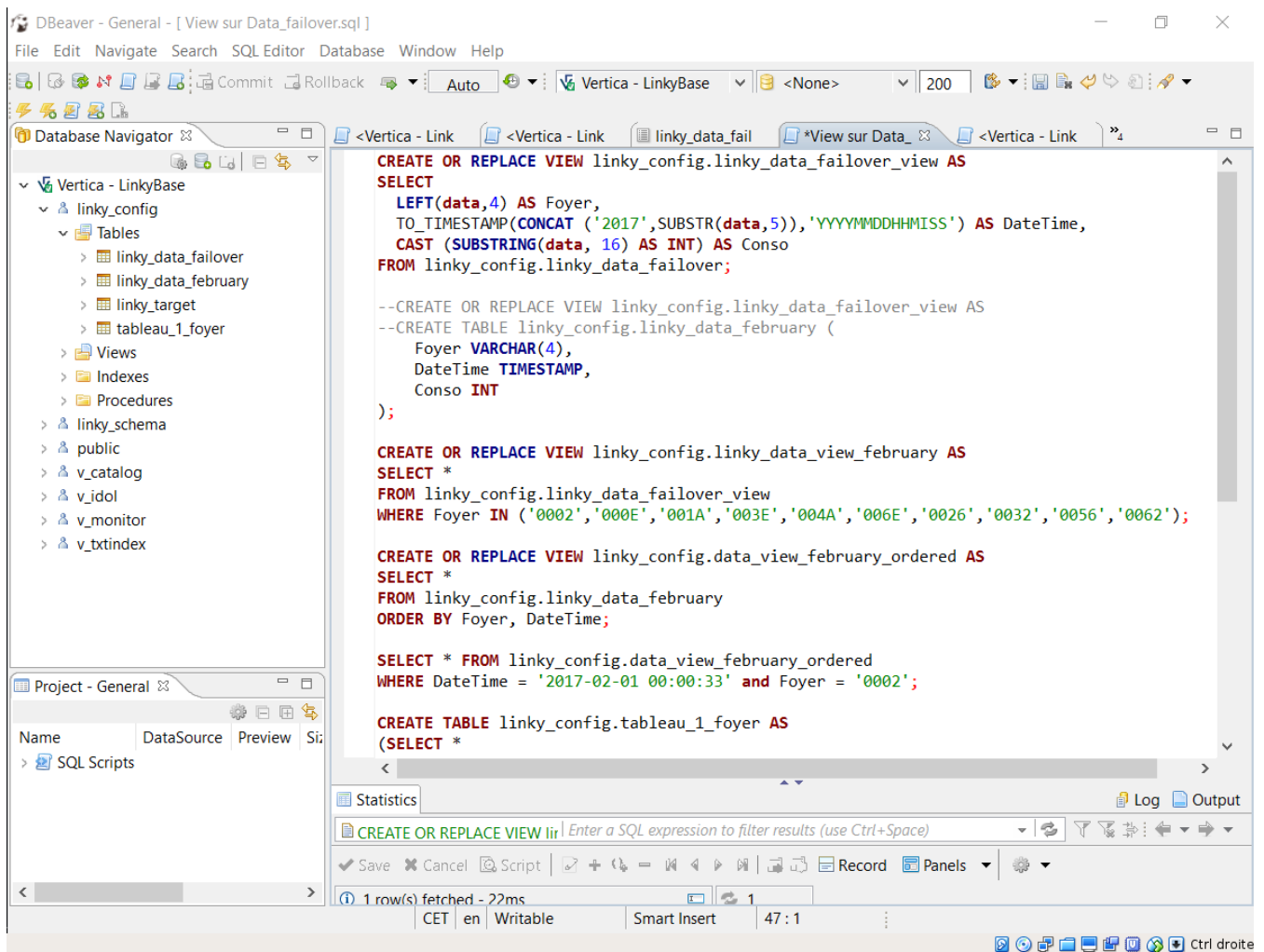
	Foyer	DateTime	Conso
1	001A	2017-02-12 04:44:34	110
2	001A	2017-02-12 04:44:35	110
3	001A	2017-02-12 04:44:38	110
4	001A	2017-02-12 04:44:43	110
5	001A	2017-02-12 04:44:47	110
6	001A	2017-02-12 04:44:50	110
7	001A	2017-02-12 04:44:56	110
8	001A	2017-02-12 04:45:04	110
9	001A	2017-02-12 04:45:08	110

15 row(s) fetched - 565ms

Save Cancel Script

Record Panels Grid Text

CET en Writable Smart Insert 1:1



We want to show you how to use Tableau now that we made our queries on DBeaver.

Tableau uses ODBC to connect to Vertica. The ODBC drivers for Vertica are part of a Vertica client package.

You need to download and install the ODBC driver on Windows.

This is the home page when you open Tableau.



On the left, click on HP Vertica.

As you can see, we connect on the server 4 to access Vertica (cf Architecture part)

HP Vertica

Serveur : Port :

Base de données :

Entrez les informations de connexion à la base de données :

Nom d'utilisateur :

Mot de passe :

[SQL initial...](#) [Connexion](#)

This is the page you see once you are connected. You just have to drag and drop one table at the indicated position

Tableau - Classeur1

Fichier Données Serveur Fenêtre Aide

Connexions: 10.0.20.24 HP Vertica

Base de données: LinkyBase

Schéma: linky_config

Table:

- data_view_febru...bruary_ordered)
- linky_data_failov...ky_data_failover)
- linky_data_failov...ta_failover_view)
- linky_data_febru...y_data_february)
- linky_data_july (l...g.linky_data_july)
- linky_data_view...a_view_february)
- linky_target (link...nfig.linky_target)
- tableau_1_foyer (...tableau_1_foyer)
- tableau_2_foyers...ableau_2_foyers)
- Nouvelle requête SQL personnalisée

linky_config

Faites glisser des tables ici

Trier les champs Ordre de la sou Afficher les alias Afficher les ch... 10000 lignes

Source de données Feuille 1

Once you drag and drop your table, you can explore the data as you can see on the next picture. In our case, we wanted to visualize the home's id, the date and the home's consumption

Tableau - Classeur1

Fichier Données Serveur Fenêtre Aide

Connexions: 10.0.20.24 HP Vertica

Base de données: LinkyBase

Schéma: linky_config

Table:

- data_view_february...w_february_ordered)
- linky_data_failover (l...ig.linky_data_failover)
- linky_data_failover_vi...ky_data_failover_view)
- linky_data_february (l...g.linky_data_february)
- linky_data_july (linky_config.linky_data_july)
- linky_data_view_febr...data_view_february)
- linky_target (linky_config.linky_target)
- tableau_1_foyer (linky...onfig.tableau_1_foyer)
- tableau_2_foyers (lin...fig.tableau_2_foyers)
- Nouvelle requête SQL personnalisée

tableau_2_foyers...

Connexion: Direct

Filtres: 0 Ajouter

tableau_2_foyers

Trier les champs Ordre de la s Afficher les Afficher les ... 1000 lignes

Abc	#	
tableau_2_fov... tableau_2_foyers	tableau_2_f...	
Foyer	Conso	
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0
0056	01/02/2017 00:0...	0

Source de données Feuille 1

For visualizing your data, click on **Feuille 1** at the bottom of Tableau. You can now see the following page.

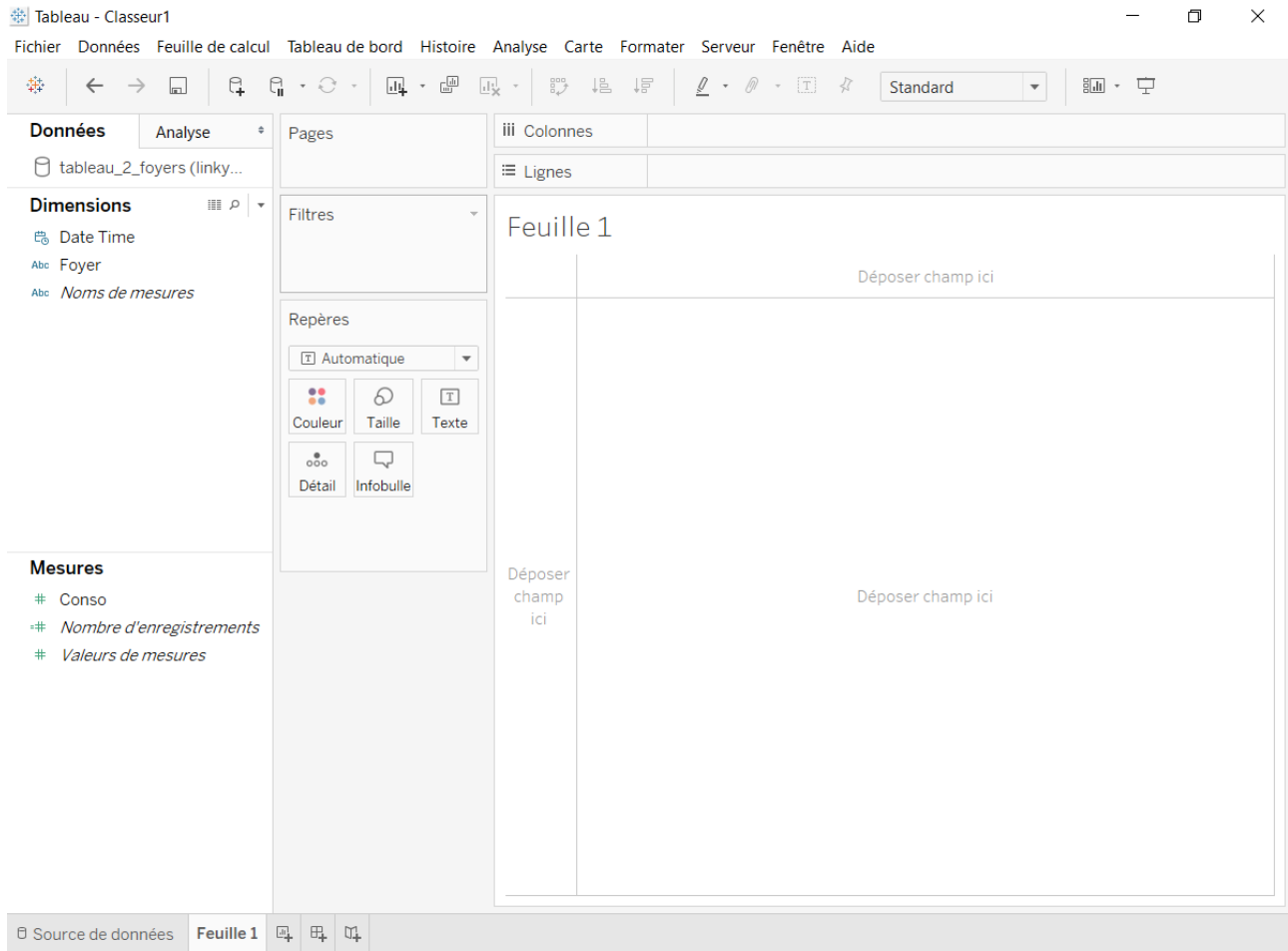


Tableau is really simple and easy to use. You can drag and drop what you need to see into « colonnes » and « Lignes ».

You can also apply filters also by dragging and dropping into the « Filtre » section.

For example, we apply a filter on the date. We only want to select only the two first weeks of February.

Filtrer [Date Time] ✕

Dates relatives

Plage de dates

Date de début

Date de fin

Spéciale

Plage de dates ☐ Afficher les heures

01/02/2017
03/03/2017

Charger le domaine

☐ Inclure des valeurs null

Réinitialiser

OK

Annuler

Appliquer

We can also apply a filter on the home to select.

Filtrer [Foyer] ✕

Général Caractère générique Condition Premiers

☒ Sélectionner dans la liste
☐ Liste de valeurs personnalisées
☐ Utiliser tout

Saisir le texte de recherche

☐ 0056
☐ 0062

Tout

Aucun

☐ Exclure

Résumé

Champ : [Foyer]
Sélection : 0 valeur(s) sur 2 sélectionnée(s)
Caractère générique : Tout
Condition : Aucun
Limite : Aucun

Réinitialiser

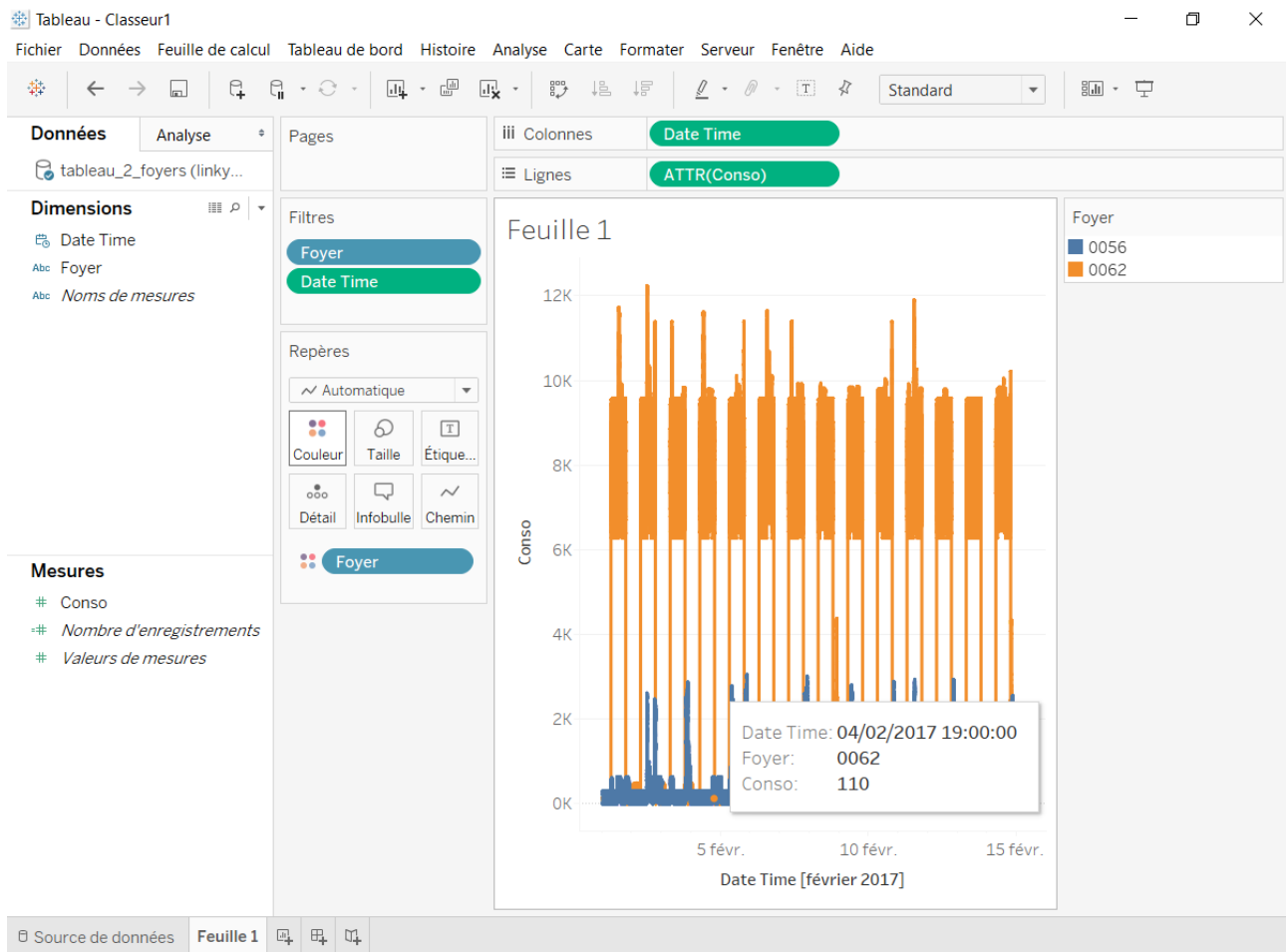
OK

Annuler

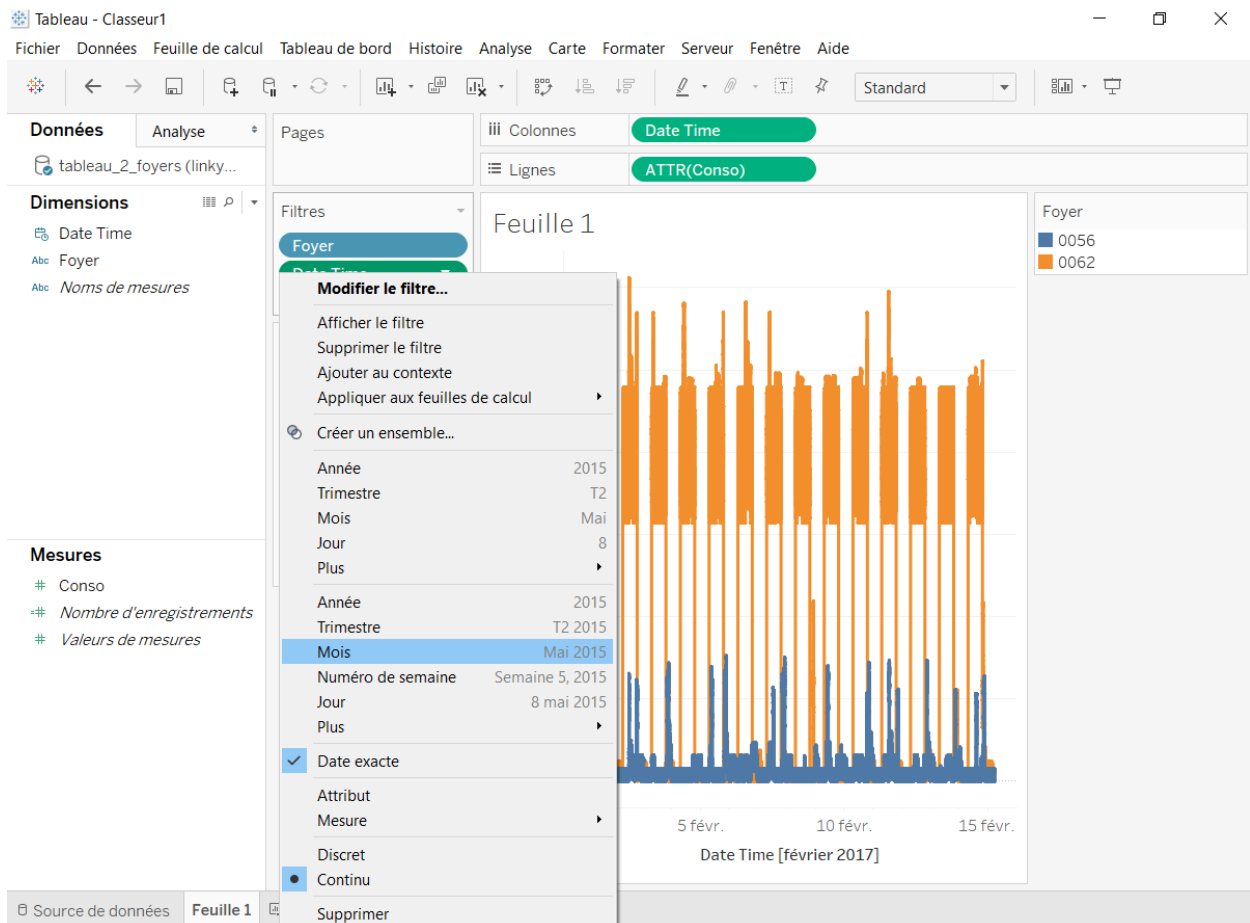
Appliquer

Here is the result of all the filter we apply on our data. you can see the electrical consumption of the two home we selected in our filter at the date we chose.

By default, the colors of the two home is identical but we can also apply a filter on the colors.

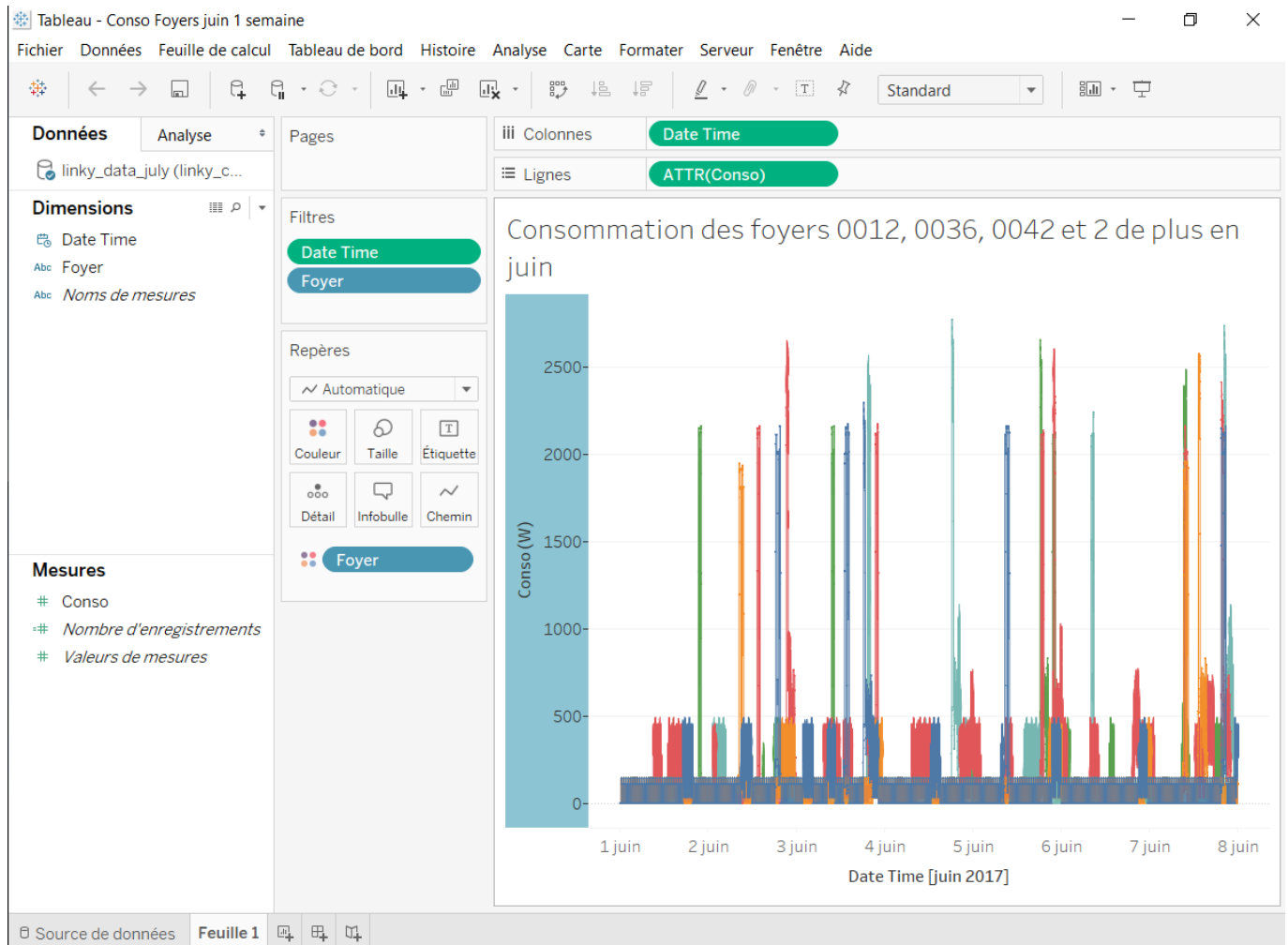


You can modify your filters whenever you want by clicking on the filter you want to change.



This is an other example of visualization of our data.

You can see the electrical consumption of several home during one week in June.



HPE asked us to realize a proof of project for EDF which is installing new intelligent electric meter in all the country in order to recover the data without sending a person retrieve them. The data are stored on servers. What EDF wants to do is to analyse all the data to be able to do prediction for the future and also proposing a consumption's visualization to their client. EDF aim to be able to suggest their clients a solution to reduce their invoice. As you can see, we have a pretty good visualization of the data of several home for a week and we think that these results are satisfying and can be showed to client asking about their home's consumption.

Amelioration

In order to improve our project in the futur, we think that use Spark will allow to improve our performance.

Apache Spark is an open source cluster-computing framework. It provides an interface for programming entire clusters with implicit data parallelism and fault-tolerance.

Spark realizes a reading of the data at the level of the cluster, makes all the necessary operations for analysis, then writes the results at the same level. In spite of the fact that he speaks with the languages Scala, Java and Python, it makes best use its capacities with his native language, Scala.

Therefore, where MapReduce of Hadoop works by stage, Spark can work on all of the data at the same time. Spark is thus until ten times faster for the treatment(processing) and until hundred times faster to make the analysis in memory.

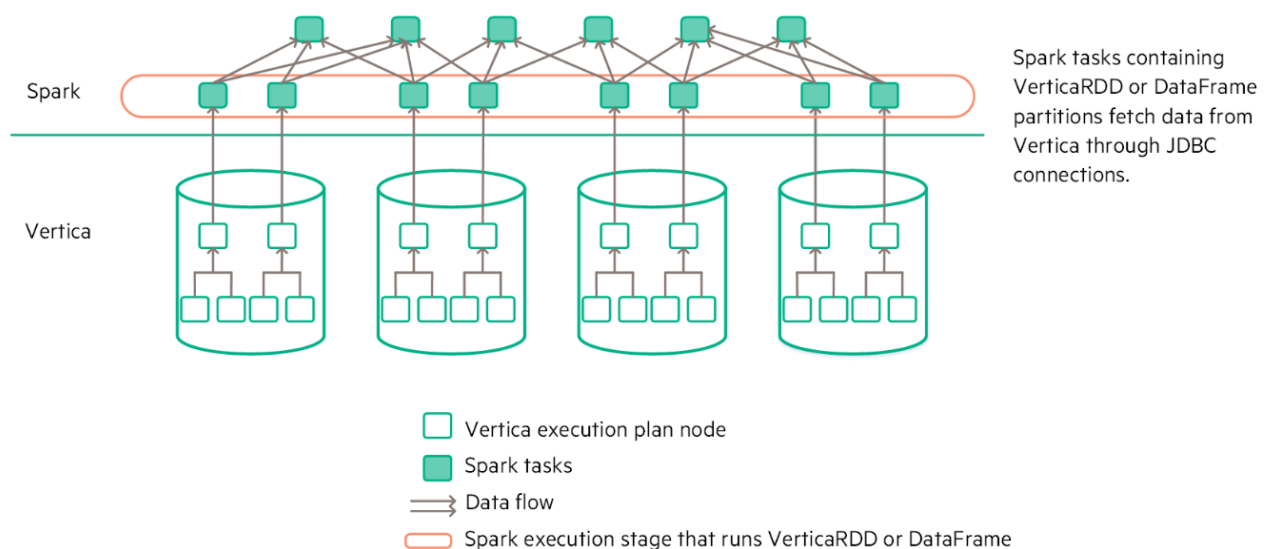
Spark executes all the operations of data analysis in memory and in real time. He does lean on records only when his memory is not sufficient anymore. On the contrary, with Hadoop the data are written on the record after each of the operations. This work in memory allows to reduce the latent periods between treatments, which explains such a speed.

However, Spark does not arrange a management system of file. It is necessary to supply with one (examples: Hadoop Distributed File System/Informix/Cassandra/OpenStack Swift/Amazon S313). It is advised to use him with Hadoop who remains at present the best global solution of storage thanks to its tools of administration, safety and monitoring more advanced.

In case of breakdown or of failure of the system: the objects of data are stored in what we call resilient distributed datasets (RDD) distribute on the cluster of data allowing the complete recovery of data.

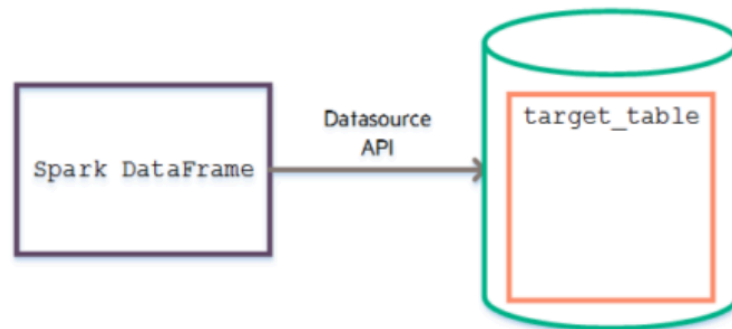
A RDD is a collection of data calculated from a source and preserved in memory lively (as long as the capacity allows it). One of the advantages brought by RDDs is in his capacity to keep enough information on the way a partition RDD was produced. In case of loss of a partition he is thus capable of recalculating it.

How Vertica and Spark work together:



Using the HPE Vertica Connector for Apache Spark, you can:

- Move large volumes of data from Spark DataFrames to Vertica tables using parallel read and write from HDFS.



- Save Spark data to Vertica with the DefaultSource API.
- Move data from Vertica to a Spark RDD or DataFrame.

