

## 1. Case study: “the e-commerce”

Consider the problem described in the “Final Project B”.

In this practical class we will develop some tools to assist the “market basket analysis” and to automatically generate the proper “.basket” format for the Orange algorithms’ processing.

The tools will be implemented using exclusively PostgreSQL capabilities and Python.

- a) Take a look at the “z\_dataset\_sample.txt” and notice that it contains events, taken from the “Final Project B” dataset, recorded during the first hour of the first day of January the 2012.

## 2. Create a database and populate it with the e-commerce dataset

- a) Execute the “00\_script\_CREATE\_DB.txt”, to create “db\_e\_commerce\_sample” database.
- b) Complete the “01\_script\_CREATE\_SCHEMA.txt”, to create the “track” table that will be populated with the dataset (cf., “z\_dataset\_sample.txt”).
- c) Adjust the “02\_script\_IMPORT\_POPULATE\_SCHEMA.txt” to import the dataset into the “track” table. Use the “pgAdmin” tool to confirm that the dataset was properly imported.

## 3. Build views to assist dataset analyzes (for the “market basket”)

Consider the file: “a03\_CREATE\_VIEW.txt”.

- a) Create the “v\_cookie\_session\_number\_of\_events” view that aggregates (groups) cookies and sessions and computes the total number of events.
- b) Create the “v\_cookie\_number\_of\_sessions” view that aggregates cookies and computes the total number of sessions (for each cookie) and also the total sum of events for each session.  
*Suggestion:* build this view from the previous (“v\_cookie\_session\_number\_of\_events”).
- c) Create the “v\_number\_of\_cookies\_number\_of\_sessions” view that aggregates the number of sessions and computes the total number of cookies (visitors) at each session. *Suggestion:* build this view from the previous.
- d) Create the “v\_number\_of\_events\_per\_session\_number\_of\_cookies” view that aggregates the number of events per session and computes the number of cookies (visitors).  
*Suggestion:* build this view from the “v\_cookie\_session\_number\_of\_events”.
- e) Use the file “z\_VIEW\_expected.txt” to test your views against the expected result.

#### 4. Export a basket for the extraction of “association rules”

- Consider the “v\_export” view that is implemented at the end of the “a03\_CREATE\_VIEW.txt” file. This is an example of a possible basket that we may now build from the previous views. Remove the comments and implement your own view with the information that you consider the most appropriate to build the “basket”. *Suggestion:* for now you may just use the provided view in order to validate the whole process (cf., next items).
- Consider the “04\_script\_EXPORT\_DATA.txt” and guarantee that the information to be exported adheres to the <TID, PID> “market basket” problem formulation.
- If you adopted the provided “v\_export” view you can validate your result against the expected one by comparing the result with the data in file “z\_dataset\_sample\_OUT\_expected.txt”.

#### 5. ... info about the “.basket” format and “string normalization”

- Before proceeding take a look at the “.basket” format as describe in “Basket Format” section in: <http://orange.readthedocs.io/en/latest/reference/rst/Orange.data.formats.html>
- Also take a look at the “product\_gui” data in file “z\_dataset\_sample\_OUT.txt”. Notice that there is some “string normalization” work to do! For example a string (in “product\_gui”) that contains the “=” character may get the “.basket” parser to fail. Why?  
  
... also all strings should be converted to the same (lower or upper) case; accents and white spaces eliminated to reduce the possibility of different strings representing the same product. Identify potential problems; recall that this dataset is a small sample from a much larger one.

#### 6. Generate the “.basket” format for sparse matrix representation

Consider the file: “\_goPy\_transform.py”.

- Complete the provided file in order to automatically normalize strings and generate “.basket” file.
- You can validate your implementation checking your result against the expected one as provided in “z\_dataset\_sample\_OUT\_expected.basket”. *Note:* be aware that your files may have the same data but following a different order from the one provided in the “\*\_expected.\*” files.

#### 7. Build your “market basket analysis” (e.g., with Orange workflows)

- Open “Orange”, select “Options \ Add-ons [x] Orange3-Associate” \ <Ok>. Build the workflows: Data\“File” >> Data\“Data Table” >> Associate\“Frequent Itemset”, and Data\“File (1)” >> Data\“Data Table” >> Associate\“Association Rules”. Explore the “\*test\*.basket” datasets. Explore “support”, “confidence” and “number of rules”.
- For **programmatically implementing** the workflows you may consider the **baseline example** provided in the folder “AssociationRulesExample”.