

# The Dataset

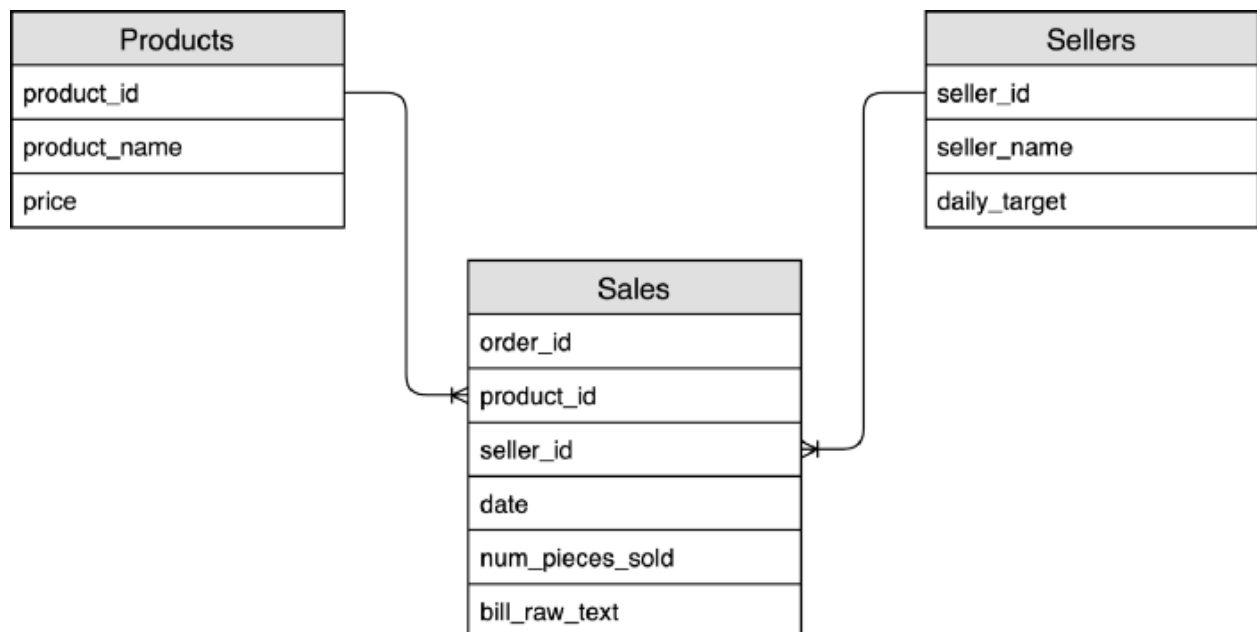
Let's describe the dataset briefly that we will use: it consists of three tables coming from the database of a shop, with products, sales, and sellers. The complete materials needed for this exercise can be downloaded from our public cloud here:

<https://cloudshellsatipytics.blob.core.windows.net/public-share/spark-exercise-materials.zip>

You will find 3 CSV files describing the data and the format of this data so that it is easier for you to see it; these CSV's are of no use for the exercise. They are simply an example of the data. The actual data you will need to parse is in PARQUET, stored in 3 folders ( Product\_parquet, sales\_parquet, sellers\_parquet ).

## The Data Model:

While you explore the below Data schema, feel free to open the CSV files and study the format and the data in general to get an understanding of the model.



**Sales Table:**

Each row in this table is an order, and every order can contain only one product. Each row stores the following fields:

- order\_id: The order ID
- product\_id: The single product sold in the order. All orders have exactly one product)
- seller\_id: The selling employee ID that sold the product
- num\_pieces\_sold: The number of units sold for the specific product in the order
- bill\_raw\_text: A string that represents the raw text of the bill associated with the order
- date: The date of the order.

**Products Table:**

Each row represents a distinct product. The fields are:

- product\_id: The product ID
- product\_name: The product name
- price: The product price

**Sellers Table:**

This table contains the list of all the sellers:

- seller\_id: The seller ID
- seller\_name: The seller name
- daily\_target: The number of items (regardless of the product type) that the seller needs to hit his/her quota. For example, if the daily target is 100,000, the employee needs to sell 100,000 products; He can hit the quota by selling 100,000 units of product\_0 or by selling 30,000 units of product\_1 and 70,000 of product\_2.

## The Exercises:

**Exercises will be shared at the beginning of the pair programming session.**

Thank you for taking the time and spending the effort  
with us!