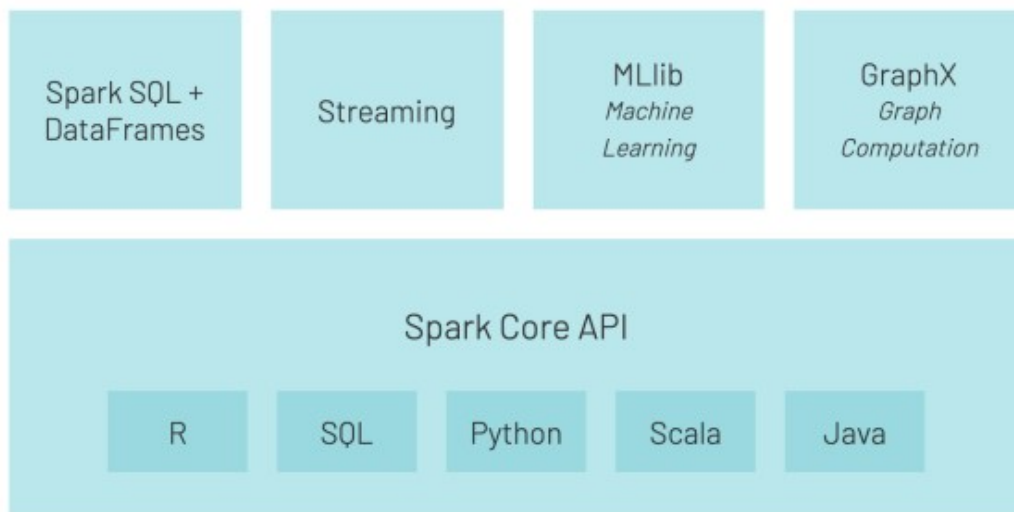Yesterday we took a closer look into Spark Scala with notebooks in Azure Databricks and how to handle data engineering. Today we will look into the Spark SQL and DataFrames that is using Spark Core API.



"Spark SQL is a spark module for structured data processing and data querying. It provides programming abstraction called DataFrames and can also serve as distributed SQL query engine. It enables unmodified Hadoop Hive queries to run up to 100x faster on existing deployments and data. It also provides powerful integration with the rest of the Spark ecosystem (e.g.: integrating SQL query processing with machine learning)." (Apache Spark Tutorial).

Start your Azure Databricks workspace and create new Notebook. I named mine as: *Day22_SparkSQL* and set the language: *SQL*. Now let's explore the functionalities of Spark SQL.

**1.Loading Data**

We will load data from */databricks-datasets* using Spark SQL, R and Python languages. The CSV dataset will be **data_geo.csv** in the following folder:

```
%scala
display(dbutils.fs.ls("/databricks-datasets/samples/population-vs-price"))
```

**1.1. Loading using Python**

```
%python
data = spark.read.csv("/databricks-datasets/samples/population-vs-
price/data_geo.csv", header="true", inferSchema="true")
```

And materialize the data using to create a view with name *data_geo_py*:

```
%python
data.createOrReplaceTempView("data_geo_py")
```

And run the following SQL Statement:

```
SELECT * FROM data_geo_py LIMIT 10
```

**1.2. Loading using SQL**

```
DROP TABLE IF EXISTS data_geo;

CREATE TABLE data_geo
USING com.databricks.spark.csv
OPTIONS (path "/databricks-datasets/samples/population-vs-price/data_geo.csv",
```

```
header "true", inferSchema "true")
```

And run the following SQL Statement:

```
SELECT * FROM data_geo LIMIT 10
```

### 1.3. Loading using R

```r
%r
library(SparkR)
data_geo_r <- read.df("/databricks-datasets/samples/population-vs-price/
data_geo.csv", source = "csv", header="true", inferSchema = "true")
registerTempTable(data_geo_r, "data_geo_r")
```

Cache the results:

```
CACHE TABLE data_geo_r
```

And run the following SQL Statement:

```
SELECT * FROM data_geo_r LIMIT 10
```



All three DataFrames are the same (unless additional modification are done; like: dropping rows with null values, etc).

### 2.Viewing DataFrame

Viewing DataFrame is done by simple SELECT statement, the ANSI SQL Standard. E.g.:

```
SELECT City
,`2014 Population estimate`
,`2015 median sales price`
,`State Code` AS State_Code
FROM data_geo
WHERE `State Code` = 'AZ';
```

```
1  SELECT City
2  ,`2014 Population estimate`
3  ,`2015 median sales price`
4  ,`State Code` AS State_Code
5  FROM data_geo
6  WHERE `State Code` = 'AZ';
```

▸ (1) Spark Jobs

| | City | 2014 Population estimate | 2015 median sales price | State_Code |
|---|---|---|---|---|
| 1 | Chandler | 254276 | null | AZ |
| 2 | Gilbert[20] | 239277 | null | AZ |
| 3 | Glendale | 237517 | null | AZ |
| 4 | Mesa | 464704 | null | AZ |
| 5 | Peoria | 166934 | null | AZ |
| 6 | Phoenix | 1537058 | 206.1 | AZ |
| 7 | Scottsdale | 230512 | null | AZ |

Showing all 10 rows.

Command took 1.12 seconds -- by tomaz.kastrun@gmail.com at 21/12/2020, 20:51:82 on databbricks_cl1_Stand

You can also combine all three DataFrames that were imported using three different languages (SQL, R, Python).

```
SELECT *, 'data_geo_SQL' AS dataset  FROM data_geo
UNION
SELECT *, 'data_geo_Python' AS dataset  FROM data_geo_py
UNION
SELECT *, 'data_geo_R' AS dataset FROM data_geo_r
ORDER BY `2014 rank`, dataset
LIMIT 12
```



```
1  SELECT *, 'data_geo_SQL' AS dataset  FROM data_geo
2  UNION
3  SELECT *, 'data_geo_Python' AS dataset  FROM data_geo_py
4  UNION
5  SELECT *, 'data_geo_R' AS dataset FROM data_geo_r
6  ORDER BY `2014 rank`, dataset
7  LIMIT 12
```

▸ (2) Spark Jobs

| | 2014 rank | City | State | State Code | 2014 Population estimate | 2015 median sales price | dataset |
|---|---|---|---|---|---|---|---|
| 1 | 1 | New York[6] | New York | NY | 8491079 | 388.6 | data_geo_Python |
| 2 | 1 | New York[6] | New York | NY | 8491079 | 388.6 | data_geo_R |
| 3 | 1 | New York[6] | New York | NY | 8491079 | 388.6 | data_geo_SQL |
| 4 | 2 | Los Angeles | California | CA | 3928864 | 434.7 | data_geo_Python |
| 5 | 2 | Los Angeles | California | CA | 3928864 | 434.7 | data_geo_R |
| 6 | 2 | Los Angeles | California | CA | 3928864 | 434.7 | data_geo_SQL |
| 7 | 3 | Chicago | Illinois | IL | 2722389 | 192.5 | data_geo_Python |

Showing all 12 rows.

Command took 1.33 seconds -- by tomaz.kastrun@email.com at 21/12/2020, 20:57:11 on databbricks_cl1_Standard

### 3.Running SQL

### 3.1. Date and Time functions

```
SELECT
 CURRENT_TIMESTAMP() AS now
,date_format(CURRENT_TIMESTAMP(), "L") AS Month_
,date_format(CURRENT_TIMESTAMP(), "LL") AS Month_LeadingZero
,date_format(CURRENT_TIMESTAMP(), "y") AS Year_
,date_format(CURRENT_TIMESTAMP(), "d") AS Day_
,date_format(CURRENT_TIMESTAMP(), "E") AS DayOFTheWeek
,date_format(CURRENT_TIMESTAMP(), "H") AS Hour
,date_format(CURRENT_TIMESTAMP(), "m") AS Minute
,date_format(CURRENT_TIMESTAMP(), "s") AS Second
```

## 3.2. Built-in functions

```sql
SELECT
 COUNT(*) AS Nof_rows
,SUM(`2014 rank`) AS Sum_Rank
,AVG(`2014 rank`) AS Avg_Rank
,SUM(CASE WHEN `2014 rank` > 150 THEN 1 ELSE -1 END) AS Sum_case
,STD(`2014 rank`) as stdev
,MAX(`2014 rank`) AS Max_Val
,MIN(`2014 rank`) AS Min_Val
,KURTOSIS (`2014 rank`) as Kurt
,SKEWNESS(`2014 rank`) AS Skew
,CAST(SKEWNESS(`2014 rank`) AS INT) AS Skew_cast
 FROM data_geo
```



## 3.3. SELECT INTO

You can also store results using SELECT INTO statement, with table being predifined:

```sql
DROP TABLE IF EXISTS tmp_data_geo;
CREATE TABLE tmp_data_geo (`2014 rank` INT, State VARCHAR(64), `State Code`
VARCHAR(2));



INSERT INTO tmp_data_geo
     FROM data_geo SELECT
                `2014 rank`
                ,State
                ,`State Code`
WHERE   `2014 rank` >= 50 AND `2014 rank` < 60 AND `State Code` = "C";



SELECT * FROM tmp_data_geo;
```

```
1   DROP TABLE IF EXISTS tmp_data_geo;
2   CREATE TABLE tmp_data_geo (`2014 rank` INT, State VARCHAR(64), `State Code` VARCHAR(2));
3
4
5   INSERT INTO tmp_data_geo
6       FROM data_geo SELECT
7                   `2014 rank`
8                   ,State
9                   ,`State Code`
10  WHERE   `2014 rank` >= 50 AND `2014 rank` < 60 AND `State Code` = "C";
11
12
13  SELECT * FROM tmp_data_geo;
```

▶ (3) Spark Jobs

| | 2014 rank ▲ | State ▲ | State Code ▲ |
|---|---|---|---|
| 1 | 56 | California | CA |
| 2 | 52 | California | CA |
| 3 | 59 | California | CA |
| 4 | 57 | California | CA |

Showing all 4 rows.

▦   📊 ▾   ⬇

Command took 3.77 seconds -- by tonaz.kastrun@gmail.com at 21/12/2020, 21:48:45 on databbricks_cl1_Standard

### 3.4. JOIN

```
SELECT
  dg1.`State Code`
 ,dg1.`2014 rank`
 ,dg2.`State Code`
 ,dg2.`2014 rank`
FROM data_geo AS dg1
JOIN data_geo AS dg2
ON dg1.`2014 rank` = dg2.`2014 rank`+1
AND dg1.`State Code` = dg2.`State Code`
WHERE dg1.`State Code` = "CA"
```

```
1   SELECT
2     dg1.`State Code`
3    ,dg1.`2014 rank`
4    ,dg2.`State Code`
5    ,dg2.`2014 rank`
6   FROM data_geo AS dg1
7   JOIN data_geo AS dg2
8   ON dg1.`2014 rank` = dg2.`2014 rank`+1
9   AND dg1.`State Code` = dg2.`State Code`
10  WHERE dg1.`State Code` = "CA"
```

▶ (1) Spark Jobs

| | State Code ▲ | 2014 rank ▲ | State Code ▲ | 2014 rank ▲ |
|---|---|---|---|---|
| 1 | CA | 285 | CA | 284 |
| 2 | CA | 186 | CA | 185 |
| 3 | CA | 154 | CA | 153 |
| 4 | CA | 36 | CA | 35 |
| 5 | CA | 137 | CA | 136 |
| 6 | CA | 185 | CA | 184 |
| 7 | CA | 138 | CA | 137 |

Showing all 16 rows.

▦   📊 ▾   ⬇

Command took 2.12 seconds -- by tomaz.kastrun@gmail.com at 21/12/2020, 21:19:37 on databbricks_cl1_Standard

### 3.5. Common Table Expressions

```
WITH cte AS (
    SELECT * FROM data_geo
    WHERE `2014 rank` >= 50 AND `2014 rank` < 60
)
SELECT * FROM cte;
```

### 3.6. Inline tables

```
SELECT * FROM  VALUES
("WA", "Seattle"),
("WA", "Tacoma"),
("WA", "Spokane") AS data(StateName, CityName)
```



### 3.7. EXISTS

```
WITH cte AS (
    SELECT * FROM data_geo
    WHERE `2014 rank` >= 50 AND `2014 rank` < 60
)
SELECT *
FROM data_geo as dg
WHERE
  EXISTS (SELECT * FROM cte WHERE cte.city = dg.city)
AND NOT EXISTS (SELECT * FROM cte WHERE cte.city = dg.city AND `2015 median
sales price` IS NULL )
```



### 3.8.Window functions

```
SELECT
```

```
 City
,State
,RANK() OVER (PARTITION BY State ORDER BY `2015 median sales price`) AS rank
,`2015 median sales price` AS MedianPrice
FROM data_geo
WHERE
  `2015 median sales price` IS NOT NULL;
```

```
1  SELECT
2    City
3    ,State
4    ,RANK() OVER (PARTITION BY State ORDER BY `2015 median sales price`) AS rank
5    ,`2015 median sales price` AS MedianPrice
6  FROM data_geo
7  WHERE
8    `2015 median sales price` IS NOT NULL;
```

▶ (3) Spark Jobs

|   | City | State | rank | MedianPrice |
|---|------|-------|------|-------------|
| 1 | Mobile | Alabama | 1 | 122.5 |
| 2 | Montgomery | Alabama | 2 | 129 |
| 3 | Huntsville | Alabama | 3 | 157.7 |
| 4 | Birmingham | Alabama | 4 | 162.9 |
| 5 | Tucson | Arizona | 1 | 178.1 |
| 6 | Phoenix | Arizona | 2 | 206.1 |
| 7 | Little Rock | Arkansas | 1 | 131.8 |

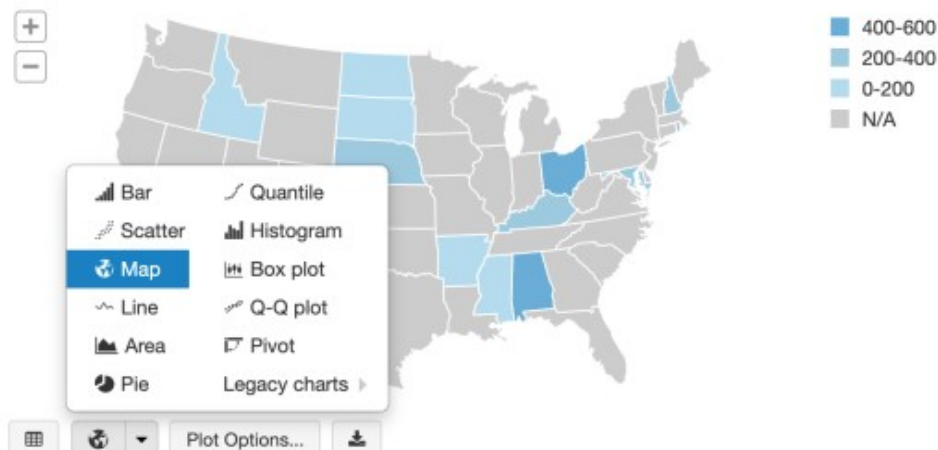Showing all 109 rows.

### 4. Exploring the visuals

Results of a SQL SELECT statements that are returned as a table, can also be visualised. Given the following SQL Statement:

```
SELECT `State Code`, `2015 median sales price` FROM data_geo
```

in the result cell you can select the plot icon and pick Map.



```
1  SELECT `State Code`, `2015 median sales price` FROM data_geo
```

▶ (1) Spark Jobs

400-600
200-400
0-200
N/A

Bar    Quantile
Scatter    Histogram
Map    Box plot
Line    Q-Q plot
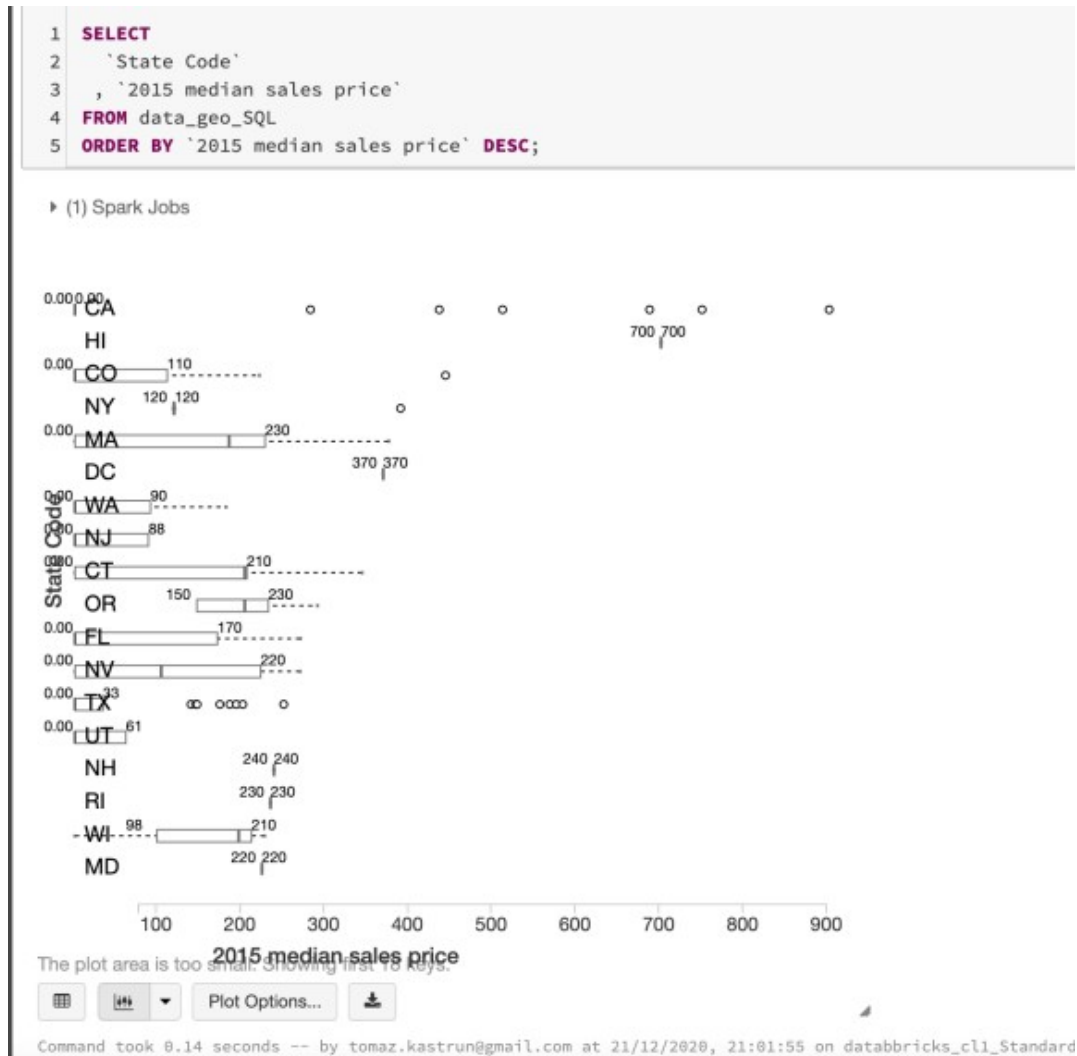Area    Pivot
Pie    Legacy charts ▸

Plot Options...

Furthermore, using "Plot Options…" you can change the settings of the variables on the graph,

aggregations and data series.

With additional query:

```
SELECT
  `State Code`
, `2015 median sales price`
FROM data_geo_SQL
ORDER BY `2015 median sales price` DESC;
```

you can also create a box-plot; again selecting the desired plot type.



There are also many other visuals available and much more SQL statements to explore and feel free to go a step further and beyond this blogpost.