Aggregations, JOINs and Nested Queries

Apache Spark™ and Databricks® allow you to create on-the-fly data lakes.

In this lesson you:

- · Use basic aggregations.
- · Correlate two data sets with a join
- · Use subselects

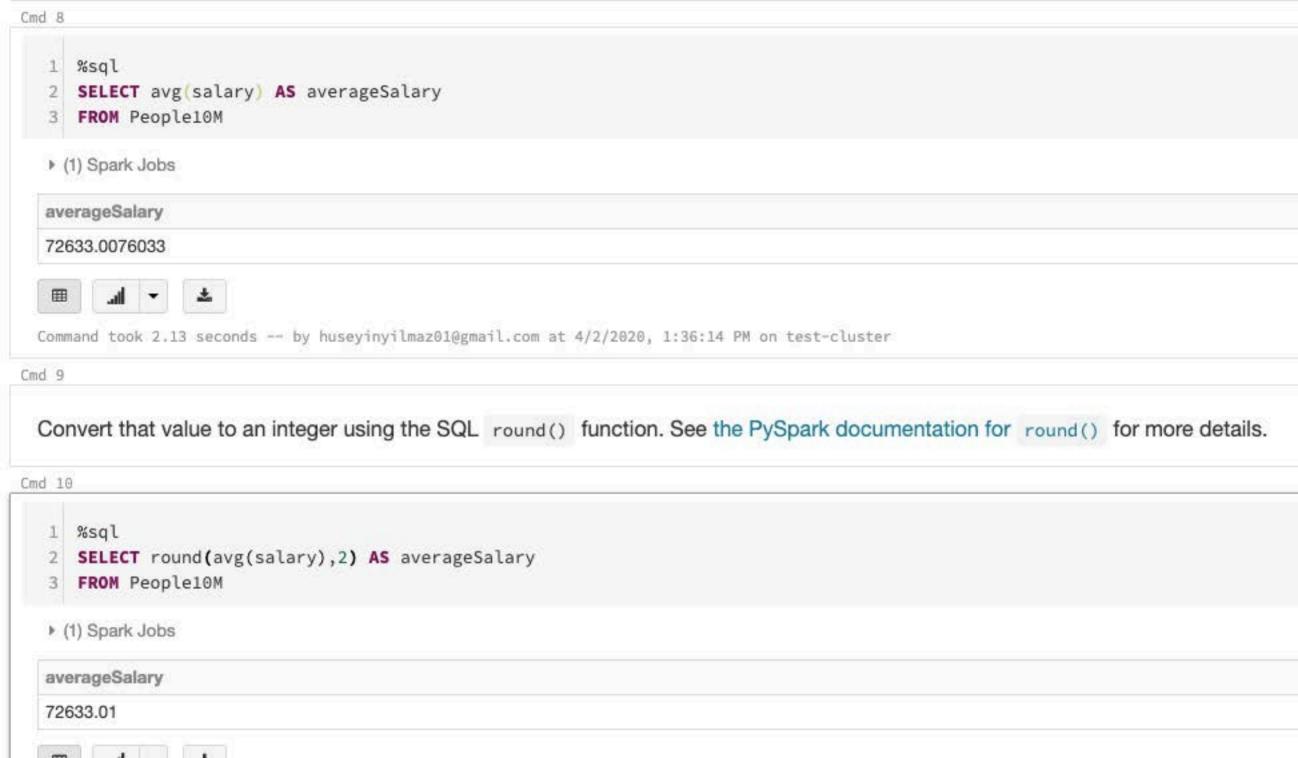
Basic aggregations

Using built-in Spark functions, you can aggregate data in various ways.

Run the cell below to compute the average of all salaries in the People10M table.



By default, you get a floating point value.



Cmd 12		
1 %sql 2 SELECT max(salary) AS max, min(salary) AS m	n round(avg(salary)) AS average	
FROM People10M (1) Spark Jobs	in, Found avg Sacury// No average	
3 FROM People10M	min win	average

CHILD TT

Joining two tables

Correlate the data in two data sets using a SQL join.

The People10M table has 10 million names in it.

How many of the first names appear in Social Security data files?

To find out, use the SSANames table with first name popularity data from the United States Social Security Administration.

For every year from 1880 to 2014, SSANames lists the first names of people born in that year, their gender, and the total number of people given that name.

By joining the People10M table with SSANames, weed out the names that aren't represented in the Social Security data.

(In a real application, you might use a join like this to filter out bad data.)

Start by taking a quick peek at what SSANames looks like.

Cmd 16

1 %sql

SELECT * FROM SSANames

▶ (1) Spark Jobs

firstName	gender	total	year
Jennifer	F	54336	1983
Jessica	F	45278	1983
Amanda	F	33752	1983
Ashley	F	33292	1983
Sarah	F	27228	1983
Melissa	F	23472	1983
Nicole	F	22392	1983
Stephanie	F	22323	1983
Heather	F	20749	1983

Showing the first 1000 rows.







Next, get an idea of how many distinct names there are in each of our tables, with a quick count of distinct names. Cmd 18 1 %sql SELECT count (DISTINCT firstName) FROM People10M (1) Spark Jobs count(DISTINCT firstName) 5113 .d → ± Command took 5.59 seconds -- by huseyinyilmaz01@gmail.com at 4/2/2020, 2:01:36 PM on test-cluster Cmd 19 1 %sql 2 SELECT count(DISTINCT firstName) FROM SSANames (1) Spark Jobs count(DISTINCT firstName) 93889 Command took 3.98 seconds -- by huseyinyilmaz01@gmail.com at 4/2/2020, 2:01:39 PM on test-cluster

By introducing two more temporary views, each one consisting of distinct names, the join will be easier to read/write.

```
Cmd 21
     %sql
     CREATE OR REPLACE TEMPORARY VIEW SSADistinctNames AS
       SELECT DISTINCT firstName AS ssaFirstName
       FROM SSANames;
  5
     CREATE OR REPLACE TEMPORARY VIEW PeopleDistinctNames AS
        SELECT DISTINCT firstName
       FROM People10M
 OK
 Command took 0.42 seconds -- by huseyinyilmaz01@gmail.com at 4/2/2020, 2:02:25 PM on test-cluster
Cmd 22
 Next, join the two tables together to get the answer.
Cmd 23
  1 %sql
  2 SELECT firstName
  3 FROM PeopleDistinctNames
     INNER JOIN SSADistinctNames ON firstName = ssaFirstName
  ▶ (4) Spark Jobs
  firstName
  Alayna
```

Melaine

Но	w many are there?
Cmd 2	5
1	%sql
2	SELECT count(*)
3	FROM PeopleDistinctNames
4	
•	(1) Spark Jobs
co	ount(1)
50	96
-	

Nested Queries

Joins are not the only way to solve the problem.

A sub-select works as well.

1	%sql
2	SELECT count(firstName)
3	FROM PeopleDistinctNames
4	WHERE firstName IN (
5	SELECT ssaFirstName FROM SSADistinctNames
6)
b (1) Spark Jobs
	unt(firstName)

Exercise 1

In the tables above, some of the salaries in the People10M table are negative

These salaries represent bad data.

Your job is to convert all the negative salaries to positive ones, and then sort



Cmd 30

Hint: See the Apache Spark documentation, built-in functions.

Step 1

Create a temporary view called PeopleWithFixedSalaries, where all the neg

```
Cmd 31
```

%sql

-- TODO 3

create or replace temporary view PeopleWithFixedSalaries as select firstName, lastName, abs(salary) as salary from People10M

Step 2

Starting with the table PeopleWithFixedSalaries, create another view called PeopleWithFixedSalariesSorted where:

- 1. The data set has been reduced to the first 20 records
- 2. The records are sorted by the column salary in ascending order

```
1 %sql
2 -- TODO
3
4 create or replace temporary view PeopleWithFixedSalariesSorted as
5 select * from PeopleWithFixedSalaries
6 order by salary asc
7 limit(20)
```

Command took 0.02 seconds -- by huseyinyilmaz01@gmail.com at 4/2/2020, 2:29:26 PM on test-cluster

Cmd 35

OK

Cmd 34

Exercise 2

As a refinement, assume that all salaries under \$20,000 represent bad rows and filter them out.

Additionally, categorize each person's salary into \$10K groups.

Cmd 37

Step 1

Create a temporary view called PeopleWithFixedSalaries20K where:

- Start with the table PeopleWithFixedSalaries
- 2. The data set excludes all records where salaries are below \$20K
- 3. The data set includes a new column called salary10k, that should be the salary in groups of 10,000. For example:
 - A salary of 23,000 should report a value of "2"
 - A salary of 57,400 should report a value of "6"
 - A salary of 1,231,375 should report a value of "123"

```
Cmd 38
```

```
%sql
2 -- TODO
3 create or replace temporary view PeopleWithFixedSalaries20K as
4 select *, round(salary/10000) as salary10k from PeopleWithFixedSalaries
5 where salary >= 20000
6 order by salary10k
```

OK

Exercise 3

Using the People10M table, count the number of females named Caren who were born before March 1980.

Cmd 41

Step 1

Starting with the table People10M, create a temporary view called Carens where:

- The result set has a single record
- 2. The data set has a single column named total
- 3. The result counts only
 - Females (gender)
 - First Name is "Caren" (firstName)
 - Born before March 1980 (birthDate)

```
1 %sql
2 -- TODO
3
4 create or replace temporary view Carens as
5 select count(*) as total from People10M
```

OK

WHERE firstName = 'Caren' AND birthDate < to_timestamp('1980-03-01')

Challenge Exercise 4

Use the SSANames table to find the most popular first name for girls in 1885, 1915, 1945, 1975, and 2005.

Cmd 6

Step 1

Create a temporary view called HistoricNames where:

- The table HistoricNames is created using a single SQL query.
- 2. The result has three columns:
 - firstName
 - o year
 - o total
- Hint: Explore the data before crafting your solution.

```
+
Cmd 7
  1 %sql
     create or replace temporary view HistoricNames as
     SELECT year, firstName, total
     FROM (SELECT firstName, year, total,
  5
                  rank() OVER (PARTITION BY year ORDER BY total DESC) as rank
           FROM (select firstName, year, max(total) as total from SSANames
  6
                 where gender = "F" and year in (1885, 1915, 1945, 1975, 2005)
  7
                 group by year, firstName
                 order by total) ) tmp
  9
     WHERE rank = 1 ORDER BY total asc
 10
```