# WEEK 8 ASSIGNMENT 1

**Large-Scale Data Storage Systems – DATA-5400 | Spring 2020**

**Christina Morgenstern**

---

For this assignment, I am using Apache Spark on my Bitnami Hadoop Virtual Machine. For running the queries, I will be using Python, as I am most familiar with this programming language.

Start Bitnami and check for availability of the package `pyspark` which comes preinstalled with Bitnami. However, by default, `pyspark` is not found.



I need to create a soft-link to Python using the following command:
```
sudo ln /usr/bin/python3 /usr/bin/python
```



Now `pyspark` is available for use with Spark. Entering the command `pyspark` starts Spark using Python version 3.5.3.

## Part 1) Word count.

Exit Spark with command `exit()` and list files and directories on local VM in Bitnami.



The text file `testfile.txt` previously created contains three sentences. This file will be used for the word count in Spark.

Start `pyspark` again and run word count command:



The first transformation retrieves the testfile.txt from the local directory.

The flatMap function splits the sentences into individual tokens based on spaces.
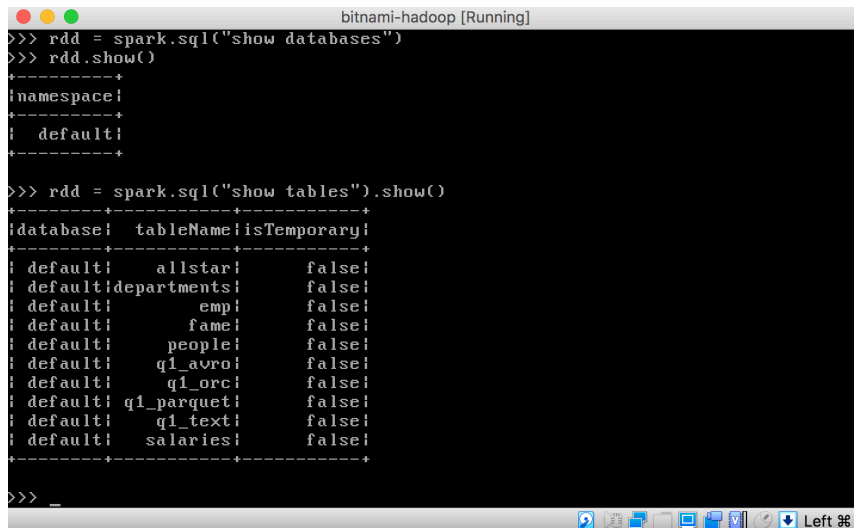
The Map function creates a count list for each word.

The reduceByKey function adds up similar words.

The results are as expected, the strings 'This', 'sentence' and 'is' occur each three times, the numbers '1', '2' and '3' each once. I guess (' ',3) means the period at the end of the sentences.

**Part 2) Using Spark to run SQL queries on the Baseball dataset.**
To run SQL commands in Spark, a variable e.g. `rdd` is created and assigned to `spark.sql` followed by the SQL command in brackets and double quotes. These SQL commands are the same as previously encountered in Hive.
`rdd = spark.sql("show tables").show()` lists all stored tables. This includes tables, that were created in earlier assignments.
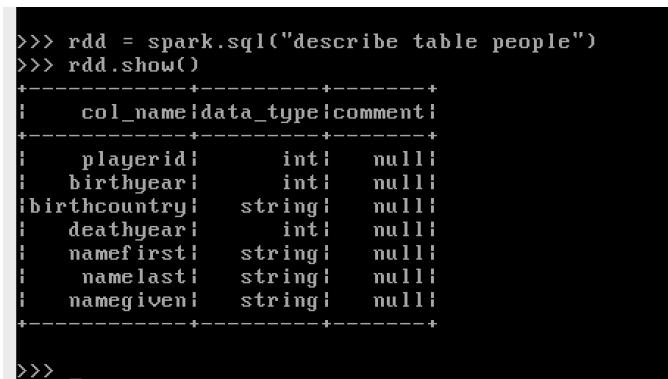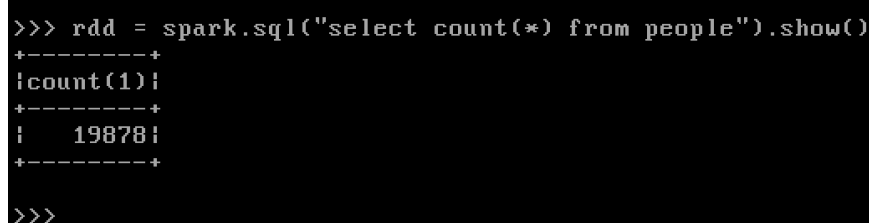


For the following queries, the people table is used. Use the `describe table` command to remind myself of the column names.



From the baseball players table, called people, the number of entries is retrieved as follows:



The total number of baseball players is 19.878. This result confirms the previous result. Although this query took slightly longer than a Hive query.

To get the number of players born before 1960, the following query was run:

```
>>> rdd = spark.sql("select count(*) from people where birthyear < 1960 ").show(
)
+--------+
|count(1)|
+--------+
|   12536|
+--------+
```

12.536 players were born before 1960.

To get the number of players born in or after 1960, the following query was run:

```
>>> rdd = spark.sql(" select count(*) from people where birthyear >= 1960").show
()
+--------+
|count(1)|
+--------+
|    7227|
+--------+

>>>
```

7.227 players were born in or after 1960.

To get the number of players born in the USA use the following command. Note that for string comparison single quotes are used.

```
>>> rdd = spark.sql("select count(*) from people where birthcountry = 'USA'").sh
ow()
+--------+
|count(1)|
+--------+
|   17254|
+--------+

>>>
```

17.254 players were born in the USA.

To get the number of players born outside the USA use the following command.

```
>>> rdd = spark.sql("select count(*) from people where birthcountry != 'USA'").s
how()
+--------+
|count(1)|
+--------+
|    2624|
+--------+

>>> _
```

2.624 players were born outside the US.