

Exercise Sheet 7

Mario Rodríguez Ruiz

June 1, 2017

Contents

1	Explain your system	4
2	Basic Resilient Distributed Dataset (RDD)	4
2.1	Create two RDD objects of a, b and do the following tasks. Words should be remained in the results of join operations.	4
2.1.1	Perform rightOuterJoin and fullOuterJoin operations between a and b. Briefly explain your solution.	4
2.1.2	Using map and reduce functions to count how many times the character "s" appears in all a and b	6
3	DataFrame API and Spark SQL	8
3.1	Replace the null value(s) in column points by the mean of all points	9
3.2	Replace the null value(s) in column dob and column last_name by "unknown" and "-" respectively	9
3.3	In the dob column, there exist several formats of dates, e.g. October 14, 1983 and 26 December 1989. Let's convert all the dates into DD-MM-YYYY format where DD, MM and YYYY are two digits for day, two digits for months and four digits for year respectively	10
3.4	Insert a new column age and calculate the current age of all students	12
3.5	Let's consider granting some points for good performed students in the class. For each student, if his point is larger than 1 standard deviation of all points, then we update his current point to 20, which is the maximum	14
3.6	Create a histogram on the new points created in the task 5	14
4	BONUS: Querying and Updating Documents	16
4.1	Use dataset students.txt for this exercise. Complete several tasks using basic querying, updating, and inserting mongoDB statements.	16
4.1.1	Display only students' last name and students' course gpas. For the course gpas, only display 5 elements starting at the 5th index in the course gpas array	16
4.1.2	Display only students' last name and students' SAT test scores . . .	17
4.1.3	Add new empty list called evaluations:[] to all students using one mongoDB command	18

List of Figures

2.1	Perform a right outer join of self and other	5
2.2	Perform a right outer join of self and other	6
2.3	How many times the character 's'	7
3.1	DataFrame initial	8
3.2	Replace the null value(s) in column 'points'	9
3.3	Replace the null value(s) in column 'dob' and column 'last_name'	10
3.4	Convert all the dates into DD-MM-YYYY format	12
3.5	Calculate the current age of all students	13
3.6	Calculate the standard deviation	14
3.7	Histogram on the new points	15
4.1	Students collections	16

4.2	Display only students' last name and students' course gpas	17
4.3	Display only students' last name and students' SAT test scores	17
4.4	MongoDB Shell	18
4.5	Add new empty list called evaluations:[]	18

List of Tables

1.1	My system	4
-----	---------------------	---

1 Explain your system

Home System	
Machine	Asus Notebook ROG G60Jx
Operating System	Windows 10 Pro 64-bit
CPU	Intel Core i7 720QM @1.60GHz
Number of cores	4
Number of threads	8
RAM	16GB @665MHz (9-9-9-24)
Programming language version Python	v3.6.1
Programming language version Java	v1.8
Hadoop	v2.6.0
Spark	v2.1.1

Table 1.1: My system

2 Basic Resilient Distributed Dataset (RDD)

2.1 Create two RDD objects of a, b and do the following tasks. Words should be remained in the results of join operations.

2.1.1 Perform rightOuterJoin and fullOuterJoin operations between a and b. Briefly explain your solution.

– rightOuterJoin() –

Perform a right outer join of self and other.

```
1 rightOuterJoin(self, other, numPartitions=None)
```

For each element (k, w) in *other*, the resulting RDD will either contain all pairs $(k, (v, w))$ for v in this, or the pair $(k, (None, w))$ if no elements in *self* have key k .

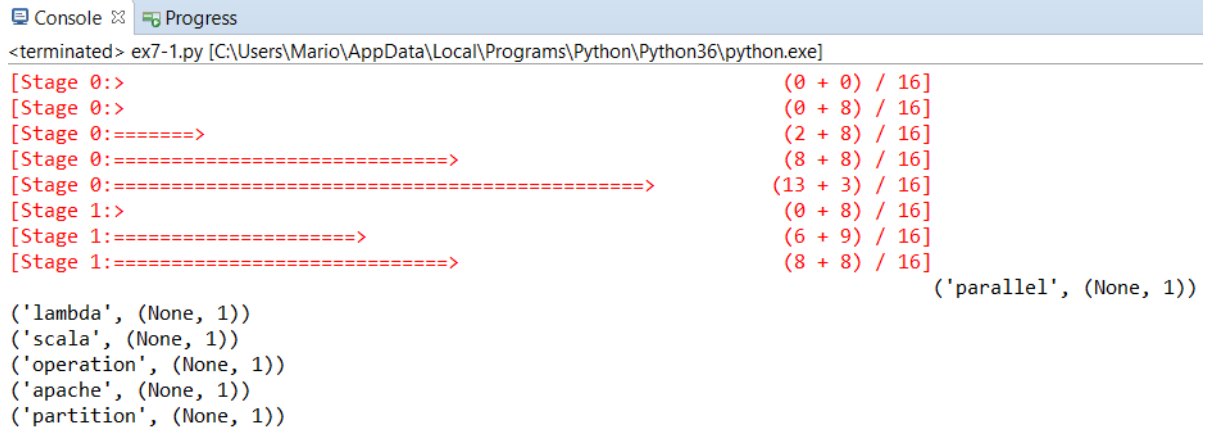
```
1 # Imports the PySpark libraries
2 from pyspark import SparkConf, SparkContext
3
4 # The 'os' library allows us to read the environment variable
5   SPARK_HOME defined in the IDE environment
6 import os
7
8 # Configure the Spark context to give a name to the
9   application
10 sparkConf = SparkConf().setAppName("ex71")
11 sc = SparkContext(conf = sparkConf)
12
13 a = ["spark", "rdd", "python", "context", "create", "class"]
14 b = ["operation", "apache", "scala", "lambda", "parallel", "partition"]
15
16 c1 = sc.parallelize(a).map(lambda a: (a, 1))
```

```

15 c2 = sc.parallelize(b).map(lambda a: (a, 1))
16
17 for i in c1.rightOuterJoin(c2).collect(): print (i)

```

The Figure 2.1 shows the execution result of the previous code.



```

Console Progress
<terminated> ex7-1.py [C:\Users\Mario\AppData\Local\Programs\Python\Python36\python.exe]
[Stage 0:> (0 + 0) / 16]
[Stage 0:> (0 + 8) / 16]
[Stage 0:====> (2 + 8) / 16]
[Stage 0:=====> (8 + 8) / 16]
[Stage 0:=====> (13 + 3) / 16]
[Stage 1:> (0 + 8) / 16]
[Stage 1:====> (6 + 9) / 16]
[Stage 1:=====> (8 + 8) / 16]
('parallel', (None, 1))
('lambda', (None, 1))
('scala', (None, 1))
('operation', (None, 1))
('apache', (None, 1))
('partition', (None, 1))

```

Figure 2.1: Perform a right outer join of self and other

– fullOuterJoin() –

Perform a right outer join of self and other.

```

1 fullOuterJoin(other, numPartitions=None)

```

For each element (k, v) in *self*, the resulting RDD will either contain all pairs $(k, (v, w))$ for w in *other*, or the pair $(k, (v, None))$ if no elements in *other* have key k .

```

1 # Imports the PySpark libraries
2 from pyspark import SparkConf, SparkContext
3
4 # The 'os' library allows us to read the environment variable
5   SPARK_HOME defined in the IDE environment
6 import os
7
8 # Configure the Spark context to give a name to the
9   application
10 sparkConf = SparkConf().setAppName("ex7-1-2")
11 sc = SparkContext(conf = sparkConf)
12
13 a = ["spark", "rdd", "python", "context", "create", "class"]
14 b = ["operation", "apache", "scala", "lambda", "parallel", "partition"]
15
16 c1 = sc.parallelize(a).map(lambda a: (a, 1))
17 c2 = sc.parallelize(b).map(lambda a: (a, 1))
18
19 for i in c1.fullOuterJoin(c2).collect(): print (i)

```

Similarly, for each element (k, w) in *other*, the resulting RDD will either contain all *pairs* $(k, (v, w))$ for v in *self*, or the *pair* $(k, (None, w))$ if no elements in *self* have key k .

The Figure 2.2 shows the execution result of the previous code.

```

<terminated> ex7-1-2.py [C:\Users\Mario\AppData\Local\Programs\Python\Python36\python.exe]
[Stage 1:=====>                                     (8 + 8) / 16]
[Stage 1:=====>                                     (9 + 7) / 16]
[Stage 1:=====>                                     (16 + 0) / 16]
('python', (1, None))
('spark', (1, None))
('context', (1, None))
('create', (1, None))
('parallel', (None, 1))
('lambda', (None, 1))
('class', (1, None))
('rdd', (1, None))
('scala', (None, 1))
('operation', (None, 1))
('apache', (None, 1))
('partition', (None, 1))

```

Figure 2.2: Perform a right outer join of self and other

2.1.2 Using map and reduce functions to count how many times the character "s" appears in all a and b

```

1 # Imports the PySpark libraries
2 from pyspark import SparkConf, SparkContext
3
4 # The 'os' library allows us to read the environment variable
   SPARK_HOME defined in the IDE environment
5 import os
6
7 # Configure the Spark context to give a name to the
   application
8 sparkConf = SparkConf().setAppName("ex7-1-2")
9 sc = SparkContext(conf = sparkConf)
10
11 a = ["spark", "rdd", "python", "context", "create", "class"]
12 b = ["operation", "apache", "scala", "lambda", "parallel", "
   partition"]
13
14 c1 = sc.parallelize(a)
15 c2 = sc.parallelize(b)

```

– How many times the character 's' appears in a –

```

1 char = 's'
2 # how many times the character "s" appears in a

```

```

3 var = c1.flatMap(lambda x: x).map(lambda x: (x==char, 1)).
  reduceByKey(lambda x, y: x + y).collect()
4
5 print("\nHow many times the character 's' appears in 'a':\t",
  var[1][1])

```

– How many times the character 's' appears in b –

```

1 char = 's'
2 # how many times the character "s" appears in b
3 var = c2.flatMap(lambda x: x).map(lambda x: (x==char, 1)).
  reduceByKey(lambda x, y: x + y).collect()
4
5 print("\nHow many times the character 's' appears in 'b':\t",
  var[1][1])

```

– How many times the character 's' appears in all a and b –

```

1 char = 's'
2 # how many times the character "s" appears in all a and b
3 c3 = c1.union(c2)
4 var = c3.flatMap(lambda x: x).map(lambda x: (x==char, 1)).
  reduceByKey(lambda x, y: x + y).collect()

```

Console Progress

<terminated> ex7-1-2.py [C:\Users\Mario\AppData\Local\Programs\Python\Python36\python.exe]

How many times the character 's' appears in 'a': 3

[Stage 2:> (0 + 8) / 8]
[Stage 3:> (0 + 8) / 8]

How many times the character 's' appears in 'b': 1

[Stage 4:> (0 + 8) / 16]
[Stage 4:==> (1 + 8) / 16]
[Stage 4:=====> (8 + 8) / 16]
[Stage 4:=====> (13 + 3) / 16]
[Stage 5:> (0 + 8) / 16]
[Stage 5:=====> (7 + 8) / 16]
[Stage 5:=====> (8 + 8) / 16]
[Stage 5:=====> (14 + 2) / 16]

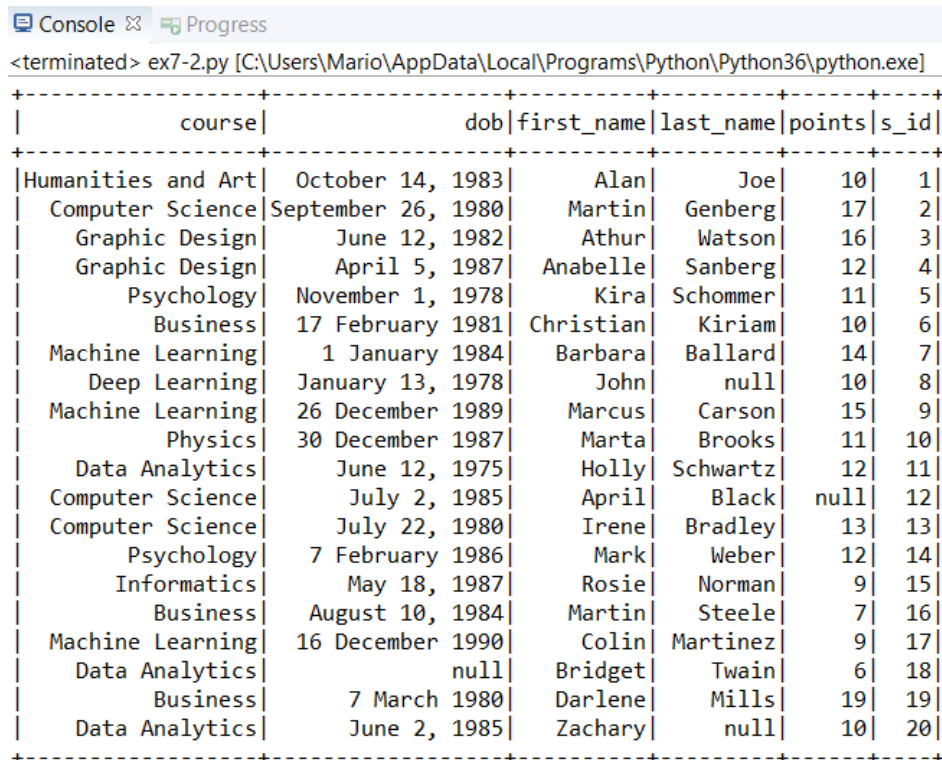
How many times the character 's' appears in all 'a' and 'b': 4

Figure 2.3: How many times the character 's'

3 DataFrame API and Spark SQL

Use dataset *students.json* for this exercise. First creating DataFrames from the dataset and do several tasks as follows:

```
1 import json
2 # Imports the PySpark libraries
3 from pyspark import SparkConf, SparkContext
4 from pyspark.sql import SQLContext
5
6 # The 'os' library allows us to read the environment variable
7   SPARK_HOME defined in the IDE environment
8 import os
9
10 # Configure the Spark context to give a name to the
11   application
12 sparkConf = SparkConf().setAppName("ex7-1-2")
13 sc = SparkContext(conf = sparkConf)
14 sqlContext = SQLContext(sc)
15
16 data = sqlContext.read.json("../datasets/students.json")
17 data.show()
```



Console Progress

<terminated> ex7-2.py [C:\Users\Mario\AppData\Local\Programs\Python\Python36\python.exe]

course	dob	first_name	last_name	points	s_id
Humanities and Art	October 14, 1983	Alan	Joe	10	1
Computer Science	September 26, 1980	Martin	Genberg	17	2
Graphic Design	June 12, 1982	Athur	Watson	16	3
Graphic Design	April 5, 1987	Anabelle	Sanberg	12	4
Psychology	November 1, 1978	Kira	Schommer	11	5
Business	17 February 1981	Christian	Kiriam	10	6
Machine Learning	1 January 1984	Barbara	Ballard	14	7
Deep Learning	January 13, 1978	John	null	10	8
Machine Learning	26 December 1989	Marcus	Carson	15	9
Physics	30 December 1987	Marta	Brooks	11	10
Data Analytics	June 12, 1975	Holly	Schwartz	12	11
Computer Science	July 2, 1985	April	Black	null	12
Computer Science	July 22, 1980	Irene	Bradley	13	13
Psychology	7 February 1986	Mark	Weber	12	14
Informatics	May 18, 1987	Rosie	Norman	9	15
Business	August 10, 1984	Martin	Steele	7	16
Machine Learning	16 December 1990	Colin	Martinez	9	17
Data Analytics	null	Bridget	Twain	6	18
Business	7 March 1980	Darlene	Mills	19	19
Data Analytics	June 2, 1985	Zachary	null	10	20

Figure 3.1: DataFrame initial

3.1 Replace the null value(s) in column points by the mean of all points

```
1 from pyspark.sql.functions import col, when, mean
2 [...]
3 mean_points = round(data.select([mean('points')]).head()[0],
4     3)
5 data.show()
6
7 ch_null_points = data.withColumn("points",
8     when(col("points").isNull(), mean_points).
9     otherwise(col("points")))
10 ch_null_points.show()
```

Console Progress

<terminated> ex7-2.py [C:\Users\Mario\AppData\Local\Programs\Python\Python36\python.exe]

Humanities and Art	October 14, 1983	Alan	Joe	10.0	1
Computer Science	September 26, 1980	Martin	Genberg	17.0	2
Graphic Design	June 12, 1982	Athur	Watson	16.0	3
Graphic Design	April 5, 1987	Anabelle	Sanberg	12.0	4
Psychology	November 1, 1978	Kira	Schommer	11.0	5
Business	17 February 1981	Christian	Kiriam	10.0	6
Machine Learning	1 January 1984	Barbara	Ballard	14.0	7
Deep Learning	January 13, 1978	John	null	10.0	8
Machine Learning	26 December 1989	Marcus	Carson	15.0	9
Physics	30 December 1987	Marta	Brooks	11.0	10
Data Analytics	June 12, 1975	Holly	Schwartz	12.0	11
Computer Science	July 2, 1985	April	Black	11.737	12
Computer Science	July 22, 1980	Irene	Bradley	13.0	13
Psychology	7 February 1986	Mark	Weber	12.0	14
Informatics	May 18, 1987	Rosie	Norman	9.0	15
Business	August 10, 1984	Martin	Steele	7.0	16
Machine Learning	16 December 1990	Colin	Martinez	9.0	17
Data Analytics	null	Bridget	Twain	6.0	18
Business	7 March 1980	Darlene	Mills	19.0	19
Data Analytics	June 2, 1985	Zachary	null	10.0	20

Figure 3.2: Replace the null value(s) in column 'points'

3.2 Replace the null value(s) in column dob and column last_name by "unknown" and "- -" respectively

```
1 from pyspark.sql.functions import col, when, mean
2 [...]
3 ch_null_dob = ch_null_points.withColumn("dob",
4     when(col("dob").isNull(), "unknown").
```

```

5         otherwise(col("dob"))
6
7 ch_null_lastName = ch_null_dob.withColumn("last_name",
8         when(col("last_name").isNull(), "--").
9         otherwise(col("last_name")))
10 ch_null_lastName.show()

```

Console Progress

<terminated> ex7-2.py [C:\Users\Mario\AppData\Local\Programs\Python\Python36\python.exe]

course	dob	first_name	last_name	points	s_id
Humanities and Art	October 14, 1983	Alan	Joe	10.0	1
Computer Science	September 26, 1980	Martin	Genberg	17.0	2
Graphic Design	June 12, 1982	Athur	Watson	16.0	3
Graphic Design	April 5, 1987	Anabelle	Sanberg	12.0	4
Psychology	November 1, 1978	Kira	Schommer	11.0	5
Business	17 February 1981	Christian	Kiriam	10.0	6
Machine Learning	1 January 1984	Barbara	Ballard	14.0	7
Deep Learning	January 13, 1978	John	--	10.0	8
Machine Learning	26 December 1989	Marcus	Carson	15.0	9
Physics	30 December 1987	Marta	Brooks	11.0	10
Data Analytics	June 12, 1975	Holly	Schwartz	12.0	11
Computer Science	July 2, 1985	April	Black	11.737	12
Computer Science	July 22, 1980	Irene	Bradley	13.0	13
Psychology	7 February 1986	Mark	Weber	12.0	14
Informatics	May 18, 1987	Rosie	Norman	9.0	15
Business	August 10, 1984	Martin	Steele	7.0	16
Machine Learning	16 December 1990	Colin	Martinez	9.0	17
Data Analytics	unknown	Bridget	Twain	6.0	18
Business	7 March 1980	Darlene	Mills	19.0	19
Data Analytics	June 2, 1985	Zachary	--	10.0	20

Figure 3.3: Replace the null value(s) in column 'dob' and column 'last_name'

3.3 In the dob column, there exist several formats of dates, e.g. October 14, 1983 and 26 December 1989. Let's convert all the dates into DD-MM-YYYY format where DD, MM and YYYY are two digits for day, two digits for months and four digits for year respectively

```

1 import json, re
2 # Imports the PySpark libraries
3 from pyspark import SparkConf, SparkContext
4 from pyspark.sql.functions import col, when, mean, udf
5 from pyspark.sql.types import StringType
6 [...]
7 # Returns the number of the month that corresponds to a
  character string

```

```

8 def monthToNumber(string):
9     m = {
10         'jan': 1,
11         'feb': 2,
12         'mar': 3,
13         'apr': 4,
14         'may': 5,
15         'jun': 6,
16         'jul': 7,
17         'aug': 8,
18         'sep': 9,
19         'oct': 10,
20         'nov': 11,
21         'dec': 12
22     }
23     s = string.strip()[:3].lower()
24     out = m[s]
25     return out
26
27 # Changes the way a date is represented (DD-MM-YYYY)
28 def dateFormatChange(dat):
29     # Return all non-overlapping matches of pattern in string,
30     # as a list of strings
31     line = re.findall(r"[\w']+", dat)
32     if "unknown" not in line:
33         for i in line:
34             # If the string is a digit
35             # otherwise it is the month
36             if i.isdigit():
37                 # If the number of elements of the digit is greater
38                 # than two is that it is a year, otherwise it is the
39                 # day
40                 if len(i)>2:
41                     y = i
42                 else:
43                     d = i
44                 else:
45                     m = str(monthToNumber(i))
46                     line = d + "-" + m + "-" + y
47     else:
48         line = ''.join(line)
49     return line
50
51 def newDateValues(value):
52     return dateFormatChange(value)
53
54 udfValueToCategory = udf(newDateValues, StringType())

```

```

53         new_dates = ch_null_lastName.withColumn("dob",
54           udfValueToCategory("dob"))
new_dates.show()

```

Console Progress

<terminated> ex7-2.py [C:\Users\Mario\AppData\Local\Programs\Python\Python36\python.exe]

course	dob	first_name	last_name	points	s_id
Humanities and Art	14-10-1983	Alan	Joe	10.0	1
Computer Science	26-9-1980	Martin	Genberg	17.0	2
Graphic Design	12-6-1982	Athur	Watson	16.0	3
Graphic Design	5-4-1987	Anabelle	Sanberg	12.0	4
Psychology	1-11-1978	Kira	Schommer	11.0	5
Business	17-2-1981	Christian	Kiriam	10.0	6
Machine Learning	1-1-1984	Barbara	Ballard	14.0	7
Deep Learning	13-1-1978	John	--	10.0	8
Machine Learning	26-12-1989	Marcus	Carson	15.0	9
Physics	30-12-1987	Marta	Brooks	11.0	10
Data Analytics	12-6-1975	Holly	Schwartz	12.0	11
Computer Science	2-7-1985	April	Black	11.737	12
Computer Science	22-7-1980	Irene	Bradley	13.0	13
Psychology	7-2-1986	Mark	Weber	12.0	14
Informatics	18-5-1987	Rosie	Norman	9.0	15
Business	10-8-1984	Martin	Steele	7.0	16
Machine Learning	16-12-1990	Colin	Martinez	9.0	17
Data Analytics	unknown	Bridget	Twain	6.0	18
Business	7-3-1980	Darlene	Mills	19.0	19
Data Analytics	2-6-1985	Zachary	--	10.0	20

Figure 3.4: Convert all the dates into DD-MM-YYYY format

3.4 Insert a new column age and calculate the current age of all students

```

1 from pyspark.sql.functions import udf, datediff, to_date, lit,
   unix_timestamp
2 [...]
3 date_fmt = '%Y-%m-%d'
4 # Current date
5 today = datetime.today()
6 today = today.strftime(date_fmt)
7
8 # Converts the days into years and returns them as integers
9 def ageValues(value):
10     if(value == None):
11         return value
12     else:
13         return int(value/365)

```

```

14
15 udfValue = udf(ageValues, IntegerType())
16
17 # Difference of days between the current date and the student's
   # date
18 age = new_dates.withColumn("age",
19                             datediff(to_date(lit(today)),
20                                     to_date(unix_timestamp('dob', "dd-MM-yyyy").cast("
                                         timestamp"))))
21
22 # Update age to years
23 age = age.withColumn("age", udfValue("age"))

```

Console Progress

<terminated> ex7-2.py [C:\Users\Mario\AppData\Local\Programs\Python\Python36\python.exe]

course	dob	first_name	last_name	points	s_id	age
Humanities and Art	14-10-1983	Alan	Joe	10.0	1	33
Computer Science	26-9-1980	Martin	Genberg	17.0	2	36
Graphic Design	12-6-1982	Athur	Watson	16.0	3	34
Graphic Design	5-4-1987	Anabelle	Sanberg	12.0	4	30
Psychology	1-11-1978	Kira	Schommer	11.0	5	38
Business	17-2-1981	Christian	Kiriam	10.0	6	36
Machine Learning	1-1-1984	Barbara	Ballard	14.0	7	33
Deep Learning	13-1-1978	John	--	10.0	8	39
Machine Learning	26-12-1989	Marcus	Carson	15.0	9	27
Physics	30-12-1987	Marta	Brooks	11.0	10	29
Data Analytics	12-6-1975	Holly	Schwartz	12.0	11	41
Computer Science	2-7-1985	April	Black	11.737	12	31
Computer Science	22-7-1980	Irene	Bradley	13.0	13	36
Psychology	7-2-1986	Mark	Weber	12.0	14	31
Informatics	18-5-1987	Rosie	Norman	9.0	15	30
Business	10-8-1984	Martin	Steele	7.0	16	32
Machine Learning	16-12-1990	Colin	Martinez	9.0	17	26
Data Analytics	unknown	Bridget	Twain	6.0	18	null
Business	7-3-1980	Darlene	Mills	19.0	19	37
Data Analytics	2-6-1985	Zachary	--	10.0	20	32

Figure 3.5: Calculate the current age of all students

3.5 Let's consider granting some points for good performed students in the class. For each student, if his point is larger than 1 standard deviation of all points, then we update his current point to 20, which is the maximum

```

1 # Calculation of the standard deviation
2 sdv = round(age.select([stddev('points')]).head()[0], 3)
3 # Calculation of the standard average or mean
4 mean_points = round(age.select([mean('points')]).head()[0], 3)
5 # if his point is larger than 1 standard deviation of all
   points,
6 # then we update his current point to20,
7 dev = age.withColumn("points",
8     when(col("points") > (mean_points+sdv+1), 20).
9     otherwise(col("points")))

```

Console Progress

<terminated> ex7-2.py [C:\Users\Mario\AppData\Local\Programs\Python\Python36\python.exe]

course	dob	first_name	last_name	points	s_id	age
Humanities and Art	14-10-1983	Alan	Joe	10.0	1	33
Computer Science	26-9-1980	Martin	Genberg	20.0	2	36
Graphic Design	12-6-1982	Athur	Watson	20.0	3	34
Graphic Design	5-4-1987	Anabelle	Sanberg	12.0	4	30
Psychology	1-11-1978	Kira	Schommer	11.0	5	38
Business	17-2-1981	Christian	Kiriam	10.0	6	36
Machine Learning	1-1-1984	Barbara	Ballard	14.0	7	33
Deep Learning	13-1-1978	John	--	10.0	8	39
Machine Learning	26-12-1989	Marcus	Carson	15.0	9	27
Physics	30-12-1987	Marta	Brooks	11.0	10	29
Data Analytics	12-6-1975	Holly	Schwartz	12.0	11	41
Computer Science	2-7-1985	April	Black	11.737	12	31
Computer Science	22-7-1980	Irene	Bradley	13.0	13	36
Psychology	7-2-1986	Mark	Weber	12.0	14	31
Informatics	18-5-1987	Rosie	Norman	9.0	15	30
Business	10-8-1984	Martin	Steele	7.0	16	32
Machine Learning	16-12-1990	Colin	Martinez	9.0	17	26
Data Analytics	unknown	Bridget	Twain	6.0	18	null
Business	7-3-1980	Darlene	Mills	20.0	19	37
Data Analytics	2-6-1985	Zachary	--	10.0	20	32

Figure 3.6: Calculate the standard deviation

3.6 Create a histogram on the new points created in the task 5

```

1 import matplotlib.pyplot as plt
2 [...]

```

```

3 sqlContext.registerDataFrameAsTable(dev, "bb")
4 columns_num = [2, 4]
5 # Extract the columns that interest (points and names)
6 df2 = dev.select(*(dev.columns[i] for i in columns_num))
7
8 # Compute histogram
9 def histogram(par):
10     points = list(map(lambda x: x[0], par))
11     name = map(lambda x: x[1], par)
12     plt.xlabel('Points')
13     plt.ylabel('Names')
14     plt.title('Task-6 Histogram')
15     plt.barh(range(len(points)), points, color = 'blue')
16     plt.yticks(range(len(points)), name)
17     plt.show()
18
19 # Put in RDD format to work with the histogram
20 par = df2.rdd.map(lambda x: (x[1], x[0]))
21 histogram(par.take(20))

```

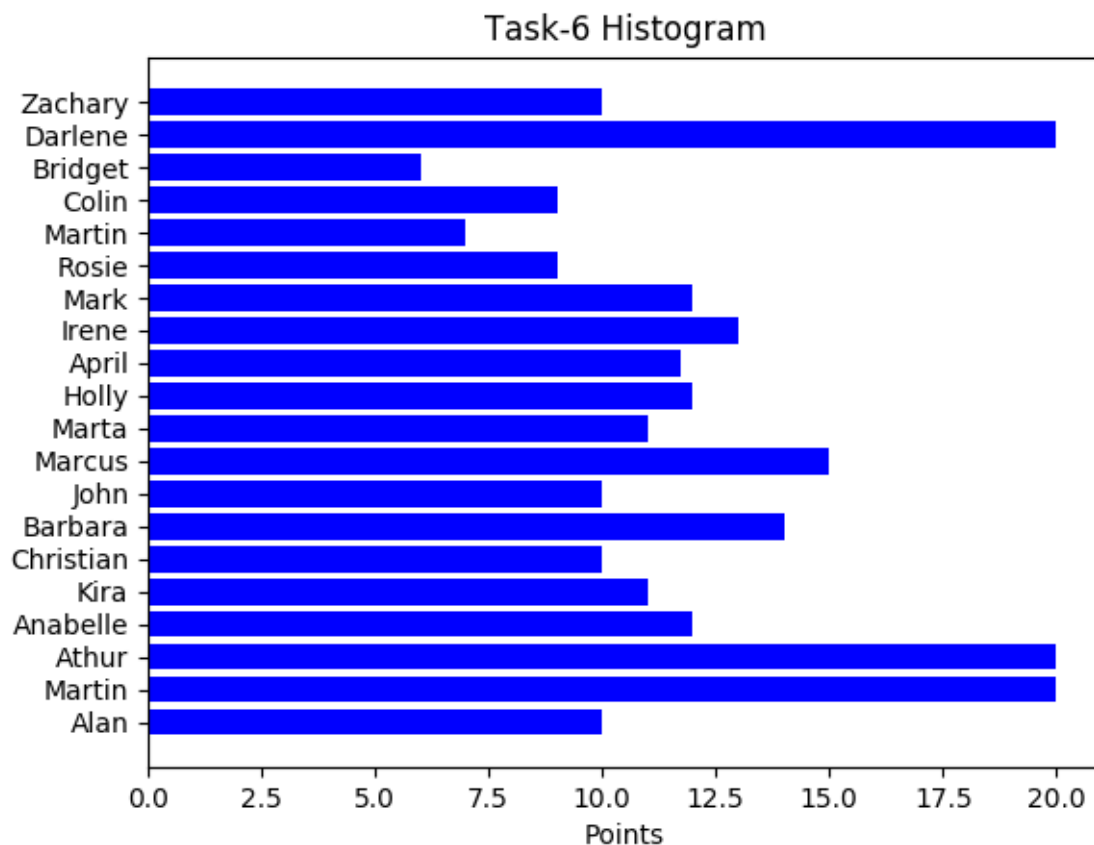


Figure 3.7: Histogram on the new points

4 BONUS: Querying and Updating Documents

4.1 Use dataset students.txt for this exercise. Complete several tasks using basic querying, updating, and inserting mongoDB statements.

Key	Value	Type
> (1) ObjectId("592f48df18d1b5f8d6a1fda6")	{ 6 fields }	Object
> (2) ObjectId("592f48df18d1b5f8d6a1fda7")	{ 6 fields }	Object
> (3) ObjectId("592f48df18d1b5f8d6a1fda8")	{ 6 fields }	Object
> (4) ObjectId("592f48df18d1b5f8d6a1fda9")	{ 6 fields }	Object
> (5) ObjectId("592f48df18d1b5f8d6a1fdaa")	{ 6 fields }	Object
> (6) ObjectId("592f48df18d1b5f8d6a1fdab")	{ 6 fields }	Object
> (7) ObjectId("592f48df18d1b5f8d6a1fdac")	{ 6 fields }	Object
> (8) ObjectId("592f48df18d1b5f8d6a1fdad")	{ 6 fields }	Object
> (9) ObjectId("592f48df18d1b5f8d6a1fdae")	{ 6 fields }	Object
> (10) ObjectId("592f48df18d1b5f8d6a1fdaf")	{ 6 fields }	Object
> (11) ObjectId("592f48df18d1b5f8d6a1fdb0")	{ 6 fields }	Object
> (12) ObjectId("592f48df18d1b5f8d6a1fdb1")	{ 6 fields }	Object
> (13) ObjectId("592f48df18d1b5f8d6a1fdb2")	{ 6 fields }	Object
> (14) ObjectId("592f48df18d1b5f8d6a1fdb3")	{ 6 fields }	Object
> (15) ObjectId("592f48df18d1b5f8d6a1fdb4")	{ 6 fields }	Object
> (16) ObjectId("592f48df18d1b5f8d6a1fdb5")	{ 6 fields }	Object
> (17) ObjectId("592f48df18d1b5f8d6a1fdb6")	{ 6 fields }	Object
> (18) ObjectId("592f48df18d1b5f8d6a1fdb7")	{ 6 fields }	Object
> (19) ObjectId("592f48df18d1b5f8d6a1fdb8")	{ 6 fields }	Object
> (20) ObjectId("592f48df18d1b5f8d6a1fdb9")	{ 6 fields }	Object
> (21) ObjectId("592f48df18d1b5f8d6a1fdbb")	{ 6 fields }	Object
> (22) ObjectId("592f48df18d1b5f8d6a1fdbb")	{ 6 fields }	Object
> (23) ObjectId("592f48df18d1b5f8d6a1fdbb")	{ 6 fields }	Object
> (24) ObjectId("592f48df18d1b5f8d6a1fdbb")	{ 6 fields }	Object
> (25) ObjectId("592f48df18d1b5f8d6a1fdbb")	{ 6 fields }	Object
> (26) ObjectId("592f48df18d1b5f8d6a1fdbb")	{ 6 fields }	Object
> (27) ObjectId("592f48df18d1b5f8d6a1fdbb")	{ 6 fields }	Object
> (28) ObjectId("592f48df18d1b5f8d6a1fdbb")	{ 6 fields }	Object
> (29) ObjectId("592f48df18d1b5f8d6a1fdbb")	{ 6 fields }	Object

Figure 4.1: Students collections

4.1.1 Display only students' last name and students' course gpas. For the course gpas, only display 5 elements starting at the 5th index in the course gpas array

```
1 from pymongo import MongoClient
2 client = MongoClient('mongodb://localhost:27017/')
3 [...]
4 db = client['Bonus7']
5 collection = db['students']
6 cur = collection.find({}, {"name.last":1, "course_gpas":{"
    $slice": [5, 5]}})
7
8 for doc in cur:
9     print(doc)
```



```

Console Progress
<terminated> ex1.py [C:\Users\Mario\Anaconda3\python.exe]
{'_id': ObjectId('592f48df18d1b5f8d6a1fda6'), 'name': {'last': 'Holland'}, 'course_gpas': [3.49, 3.96, 2.22, 2.24, 2.9]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fda7'), 'name': {'last': 'Mills'}, 'course_gpas': [3.49, 3.02, 2.91, 3.74, 2.95]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fda8'), 'name': {'last': 'Cotter'}, 'course_gpas': [2.51, 3.31, 2.5, 3.99, 3.13]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fda9'), 'name': {'last': 'Aldwell'}, 'course_gpas': [2.56, 2.82, 3.9, 3.4, 2.55]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdaa'), 'name': {'last': 'Walker'}, 'course_gpas': [2.74, 3.29, 3.79, 2.79, 3.79]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdab'), 'name': {'last': 'Barrett'}, 'course_gpas': [3.23, 3.94, 3.86, 2.48, 3.91]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdac'), 'name': {'last': 'Edwards'}, 'course_gpas': [2.88, 3.41]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdad'), 'name': {'last': 'Caldwell'}, 'course_gpas': [3.29, 2.94, 3.43, 2.97, 3.6]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdae'), 'name': {'last': 'Roberson'}, 'course_gpas': [3.98, 2.5, 2.95, 3.59, 3.05]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdaf'), 'name': {'last': 'Bass'}, 'course_gpas': [2.37, 3.85, 2.14, 2.89, 2.94]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb0'), 'name': {'last': 'Reyes'}, 'course_gpas': [2.69, 3.48, 3.79, 2.35, 2.67]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb1'), 'name': {'last': 'Lynch'}, 'course_gpas': [3.73, 2.74, 3.2, 3.66, 3.89]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb2'), 'name': {'last': 'Hampton'}, 'course_gpas': [2.24, 2.03, 3.56, 2.36, 2.48]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb3'), 'name': {'last': 'Maxwell'}, 'course_gpas': [3.91, 2.62, 3.34, 3.84, 3.17]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb4'), 'name': {'last': 'Arnold'}, 'course_gpas': [2.26, 2.43, 2.64, 3.69, 3.87]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb5'), 'name': {'last': 'Mills'}, 'course_gpas': [3.03, 2.66, 3.77, 3.38, 3.97]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb6'), 'name': {'last': 'Martinez'}, 'course_gpas': [3.55, 2.81, 3.99, 2.71, 3.13]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb7'), 'name': {'last': 'Ford'}, 'course_gpas': [2.22, 3.23, 3.6, 3.85, 2.09]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb8'), 'name': {'last': 'Washington'}, 'course_gpas': [2.34, 3.05, 3.57]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb9'), 'name': {'last': 'Steele'}, 'course_gpas': [2.63, 3.61, 2.88, 3.18, 2.77]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb9'), 'name': {'last': 'Norman'}, 'course_gpas': [3.52, 2.55, 3.84, 3.33, 3.85]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdbb'), 'name': {'last': 'Bradley'}, 'course_gpas': [2.81, 2.42, 3.12, 3.55, 3.44]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdbb'), 'name': {'last': 'Weber'}, 'course_gpas': [3.61, 2.41, 2.26, 2.98, 3.16]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdbd'), 'name': {'last': 'Black'}, 'course_gpas': [3.78, 3.05, 2.33, 3.05, 3.69]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdbd'), 'name': {'last': 'Schwartz'}, 'course_gpas': [3.02, 3.53, 2.14, 3.36, 2.5]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdbf'), 'name': {'last': 'Brooks'}, 'course_gpas': [2.79, 3.87, 2.34, 3.3, 2.42]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdc0'), 'name': {'last': 'Carson'}, 'course_gpas': [3.66, 2.42, 3.74, 2.66, 2.11]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdc1'), 'name': {'last': 'Henderson'}, 'course_gpas': [3.69, 3.58, 2.35, 3.11, 3.8]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdc2'), 'name': {'last': 'Ballard'}, 'course_gpas': [3.08, 3.95, 2.11, 3.82, 3.9]}

```

Figure 4.2: Display only students' last name and students' course gpas

4.1.2 Display only students' last name and students' SAT test scores

```

1 from pymongo import MongoClient
2 client = MongoClient('mongodb://localhost:27017/')
3 [...]
4 db = client['Bonus7']
5 collection = db['students']
6 cur = collection.find({"test_scores":
7     {'$elemMatch': {"test": "SAT", "score": {"$gt": 700} }},
8     {"name.last": 1, "test_scores":
9     {"$slice": -1}}).sort("test_scores.score", -1)
10
11 for doc in cur:
12     print(doc)

```

```

Console Progress
<terminated> ex1.py [C:\Users\Mario\Anaconda3\python.exe]
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb9'), 'name': {'last': 'Norman'}, 'test_scores': [{'test': 'SAT', 'score': 786}]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb9'), 'name': {'last': 'Walker'}, 'test_scores': [{'test': 'SAT', 'score': 783}]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb9'), 'name': {'last': 'Weber'}, 'test_scores': [{'test': 'SAT', 'score': 766}]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb9'), 'name': {'last': 'Black'}, 'test_scores': [{'test': 'SAT', 'score': 763}]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb9'), 'name': {'last': 'Mills'}, 'test_scores': [{'test': 'SAT', 'score': 759}]}
{'_id': ObjectId('592f48df18d1b5f8d6a1fdb9'), 'name': {'last': 'Steele'}, 'test_scores': [{'test': 'SAT', 'score': 718}]}

```

Figure 4.3: Display only students' last name and students' SAT test scores

4.1.3 Add new empty list called evaluations:[] to all students using one mongoDB command

```
1 # Add new empty list called evaluations:[]
2 collection.updateMany({}, { "$set" : {"evaluation":[]} })
3 # Add "eval comment": "This student is very clever" in "
  evaluation"
4 collection.updateMany({}, { "$push" : {"evaluation.0.
  eval_comment":"This student is very clever"}}
5 )
6 # Add "eval comment": "This student always submits exercises
  ontime" in "evaluation"
7 collection.updateMany({}, { "$push": {"evaluation.0.
  eval_comment":"This student always submits exercises ontime
  "}}
8 )
```

The version of pymongo has given me quite a few problems with the updateMany function, so I had to perform those three functions using the MongoDB console as Figure 4.4 shows.

```
Símbolo del sistema - mongo
> db.students.updateMany({}, { "$set" : {"evaluation":[]} })
{ "acknowledged" : true, "matchedCount" : 29, "modifiedCount" : 29 }
> db.students.updateMany({}, { "$push" : {"evaluation.0.eval_comment":"This student is very clever"}} )
{ "acknowledged" : true, "matchedCount" : 29, "modifiedCount" : 29 }
> db.students.updateMany({}, { "$push" : {"evaluation.0.eval_comment":"This student always submits exercises ontime"}} )
{ "acknowledged" : true, "matchedCount" : 29, "modifiedCount" : 29 }
>
```

Figure 4.4: MongoDB Shell

[illegible]

Figure 4.5: Add new empty list called evaluations:[]