



Master's Programme in Data Science

Spark project

Mikko Saukkoriipi 013877851

November 26, 2020

UNIVERSITY OF HELSINKI
FACULTY OF SCIENCE
P. O. Box 68 (Pietari Kalmin katu 5)
00014 University of Helsinki

Tiedekunta — Fakultet — Faculty		Koulutusohjelma — Utbildningsprogram — Degree programme	
Faculty of Science		Master's Programme in Data Science	
Tekijä — Författare — Author			
Mikko Saukkoriipi 013877851			
Työn nimi — Arbetets titel — Title			
Spark project			
Työn laji — Arbetets art — Level		Aika — Datum — Month and year	
Report		November 26, 2020	
		Sivumäärä — Sidantal — Number of pages	
		13	
Tiivistelmä — Referat — Abstract			
<p>This is my Spark project document. The given project questions were:</p> <ol style="list-style-type: none">1. For the first data set (data-1.txt), provide the value of the median of the data set. You should provide the exact median, not an approximation, and sorting the complete data set is not an acceptable solution.2. The second data set (data-2.txt) contains the matrix A. Your task is to calculate $A * A^T * A$ and provide the resulting matrix as your answer in the same format as the input matrix.			
Avainsanat — Nyckelord — Keywords			
Spark, Quickselect, Matrices			
Säilytyspaikka — Förvaringsställe — Where deposited			
Muita tietoja — Övriga uppgifter — Additional information			

Contents

1	Introduction	1
2	Finding the median	3
2.1	Searching median by counting	3
2.2	Quickselect	4
3	Matrix multiplication	7
3.1	Matrix multiplication with SQL	7
3.2	Matrix multiplication with MLib	7
4	Conclusions	11
	References	13

1. Introduction

The spark project was more demanding and time-consuming than I expected. In the last period, I had done the course Big Data Platforms where we also used PySpark, and I thought that I know the basics of Spark. Unfortunately, I soon found out that I knew a lot less about Spark programming and it's libraries than I had thought. At the same time, I was fortunate to know Amdahl's law and how the program must be splittable, to be efficient in distributed architecture. I tried to keep this in mind when designing the program, but it was surprisingly tricky, and it will need much more practice.

For this project, I spent approximately 40 hours. The first 8 hours I spent on familiarising myself with how to connect to Ukko2, how to send files with ssh, how to write code in terminal (used vim), and how to write scripts. Before this project, I had only written code in Jupyter Notebook, and this was the first time writing scripts. The rest of the hours I spent on reading Spark manuals and testing on what works and what doesn't.

There was lots of confusion about the allowed tools. On the second last day, it was told that Mllib is permitted, which was a great help. In the future, I would hope for more precise guidance regarding this kind of thing.

I would have hoped that we would have more time to spent on this project. Now we had to write essays simultaneously and do projects and homework for other classes. As a result, I had to reduce the time I spent on reading articles and writing the essays.

Unfortunately, I got sick three days before the deadline, but I did what I could, and I am pretty satisfied with the results.

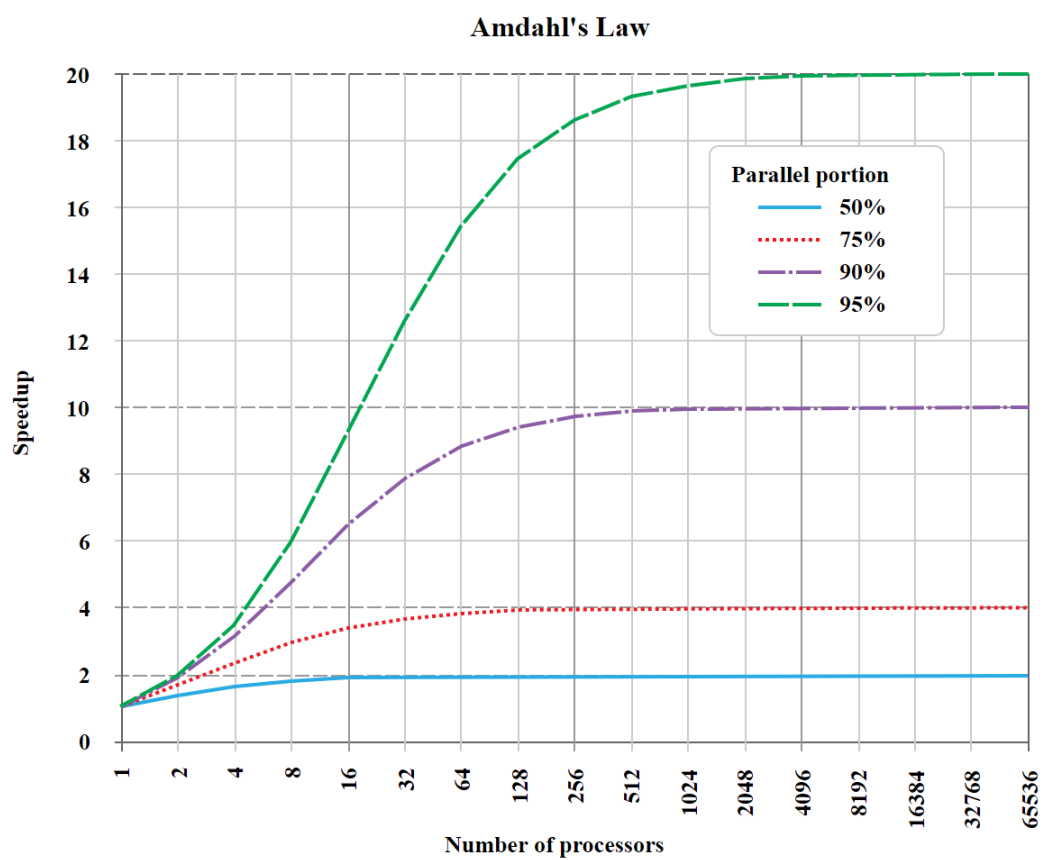


Figure 1.1: Amdahl's Law [2]

2. Finding the median

Finding the median without sorting was a trickier task than I expected. Next, I will introduce the models creating process and the results I got. I planned two different algorithms for this exercise, and the Quickselect style algorithm was implemented in Spark.

2.1 Searching median by counting

My first plan was to start counting the place in the set for the numbers. In other words, count how many numbers there are in the set that are less than the given number. Then continue this process until we found numbers in places $k/2$ and $k/2+1$, where k is number of numbers in list. Then finally calculate average of these middle numbers. I did implement this algorithm in Python3, and code is below.

```
import pandas as pd
from numpy import median

# Read the data
df = pd.read_csv("data-1-sample.txt", header=None)
df = df.rename(columns={"index": "value"})
df1 = df

# Count number of values less than value in row
def count_less(values, row):
    values = [value for value in values if value <= row[0]]
    cnt_values = len(values)
    return cnt_values

# Make list of all values
values = list(df1[0])

# Count how many values have smaller value
df1["cnt_less"] = df1.apply(lambda row: count_less(values, row), axis=1)

# Select two middle values (works if Len is even)
df1 = df1.loc[(df1["cnt_less"] <= df1.shape[0]/2+1) & (df1["cnt_less"] >= df1.shape[0]/2)]

# Show df
df1

]:
```

	0	cnt_less
214	40.480026	500
703	40.497400	501

```
print("True median:", df[0].median())
print("My median:", sum(df1[0])/2)

True median: 40.48871315055
My median: 40.48871315055
```

This model would have needed two improvements to be efficient. The first im-

provement would be to add a break when middle numbers are found and not to run it through all numbers. The second improvement would be pre-filtering, where to drop part of the numbers. This could be done, for example, by dropping 30% of the smallest and largest numbers. This would have reduced unnecessary calculations.

I did not implement this prototype and its improvements because I started to worry about this model's efficiency. It is easy to distribute, but it does lots of unnecessary calculations. For this reason, I did not write this algorithm to Spark but changed my plan to the Quickselect algorithm.

2.2 Quickselect

I had read about the Quickselect algorithm before, but this was the first time I implemented it. Unfortunately, it was not as efficient as I had hoped for this task. It would have become more efficient with prefiltering and indexing moving of the pivot.

According to Wikipedia [3], Quickselect can reach the performance of $O(n \log n)$ to $O(n)$, but in the worst case, performance is $O(n^2)$. With Quickselect, I was able to get a sample set median in less than 30 seconds, but with the full dataset, it did not finish in 2 hours. Therefore I don't have a median for the full set, but I am sure that the algorithm works. I guess that with a given full-size dataset, the performance is close to $O(n^2)$, which is too slow.

I had plans on how to improve my QuickSelect algorithm. The first improvement would be to do pre-filtering, similar way as with my counting algorithm. The second improvement would be to make pivot smarter with indexing. Also, I started to think if QuickSelect even is efficient in distributed infrastructure with multiple cores. Unfortunately, I did not have more time to dig deeper into this problem.

```

median.py • matrix.py •
C: > Users > saukk > Desktop > Distributed data infrastructure > Spark project > median.py > ...
1 from pyspark import SparkContext, SparkConf
2 import os
3 import sys
4 import math
5
6 # Datasets path on shared group directory on Ukko2.
7 dataset = "/wrk/group/grp-ddi-2020/datasets/data-1-sample.txt"
8 #dataset = "/wrk/group/grp-ddi-2020/datasets/data-1.txt"
9
10 conf = (SparkConf())
11     .setAppName("mpsaukko")    ##change app name to your username
12     .setMaster("spark://10.251.52.13:7077")
13     .set("spark.cores.max", "10")    ##dont be too greedy ;)
14     .set("spark.rdd.compress", "true")
15     .set("spark.executor.memory", "32G")
16     .set("spark.broadcast.compress", "true")
17 sc = SparkContext(conf=conf)
18
19 # Read the data
20 data = sc.textFile(dataset)
21 data = data.map(lambda s: float(s))
22 data.cache()
23
24 # Define functio find_median
25 def find_median(data):
26     count = data.count()
27     middle = round(count / 2) + 1
28
29     # If number of numbers in list is odd
30     if (count % 2 != 0):
31         return quickselect(data, middle)
32
33     # If number of numbers in file is even, then calculate average of two middle numbers
34     else:
35         middle1 = quickselect(data, middle)
36         middle2 = quickselect(data, middle - 1)
37         return (middle1 + middle2) / 2
38
39
40 def quickselect(array, k):
41
42     # Select pivot number
43     pivot = array.takeSample(False, 1, seed=0)
44     pivot = pivot[0]
45
46     # Create two arrays. One with values more than pivot and one with values greater than pivot.
47     numbers_less = array.filter(lambda x: x < pivot)
48     numbers_greater = array.filter(lambda x: x > pivot)
49
50     # If middle value found, then return pivot
51     if k == (numbers_less.count() + 1):
52         return pivot
53
54     # If k<= number_less, then run quickselect again.
55     if k <= (numbers_less.count()):
56         return quickselect(numbers_less, k)
57
58     # If k > number_less, then run quickselect again.
59     if k > (numbers_less.count() + 1):
60         return quickselect(numbers_greater, k - (numbers_less.count() + 1))
61
62 # Find and print median
63 median = find_median(data)
64 print(" ")
65 print("----- Median is: " + str(median) + " -----")
66 print(" ")

```

Figure 2.1: My Spark median finding code

```
[25]: import pandas as pd
      from numpy import median

[26]: # Read the data
      A = pd.read_csv("data-2-sample.txt", sep=" ", header=None)

      # Drop last column
      A = A.iloc[:, :-1]

      # Find A transpose
      AT = A.T

      # Calculate AAT
      AAT = A.dot(AT)

      # Calculate AATA
      AATA = AAT.dot(A)

      # Show first row
      AATA.head(1)
```

Out[26]:

9	...	990	991	992	993	994	995	996	997	998	999
0552e+11	...	1.273313e+11	1.309680e+11	1.288276e+11	1.265437e+11	1.259945e+11	1.266185e+11	1.279121e+11	1.272910e+11	1.232555e+11	1.258357e+11

Figure 2.2: Python code to confirm correct result for matrix computation sample set

3. Matrix multiplication

The matrix multiplication task was more straightforward than the median calculation. The median calculation task restriction to avoid sorting in, by all means, made the designing of the algorithm challenging, but with matrix multiplication, there was no such issue. Now we did not have such limitations, but the limits came from the tools that Spark was able to offer. My first plan was to calculate multiplications with spark.sql queries, but later it was revealed that Mllib is also allowed to use. After this, I changed the plan and implemented the whole algorithm based on Mllib.

3.1 Matrix multiplication with SQL

My first plan was to do this task bases on Spark SQL queries. Matrix multiplication can be calculated with Spark SQL with code below (code is not tested).

```
results = spark.sql("SELECT A.row row, B.col col, SUM(A.val * B.val) val  
FROM A JOIN B ON A.col = B.row GROUP BY A.row, B.col")
```

As can be seen, this multiplication made this way is pretty straight forward. But there is a problem with efficiency because it can not be distributed to Spark workers efficiently. Fortunately, it was revealed that we are allowed to use Mllib [1], and so I abandoned this idea.

3.2 Matrix multiplication with MLib

I had no previous experience with matrix computations on Spark nor Mllib, which revealed a more challenging task than expected. Mllib offers all the needed tools for this task, but it requires a sophisticated combination of the Mllib matrices and operations.

As I had no previous experience in this kind of work, I spent most of the time reading the Mllib and Spark manuals. Different Mllib matrix classes are meant for different tasks and have different properties, and it was important to find the most efficient combination. With my first Mllib based code, I got the correct result for the

sample set in 1 minute. After changing the used matrices, reducing the number of steps, adding caching and tuning block matrices block sizes, I squeezed this time to 20 seconds.

Unfortunately, I was not able to get results for the full dataset. With 10 cores and 16gb of memory, it will end up with a no memory left error. This is a surprising result because all the steps are made with Spark's native tools. Spark offers so-called sparse and dense matrices, but our only option was to use a dense matrix. After further research, I found there are better libraries available for large-scale matrix multiplications. Unfortunately, we did not have a chance to use those tools in this task.

```

C:\Users> saukk > Desktop > Distributed data infrastructure > Spark project > matrix.py > ...
1  from pyspark import SparkContext, SparkConf, mllib
2  import os
3  import sys
4  from pyspark.sql import SparkSession
5  from pyspark.mllib.linalg.distributed import IndexedRowMatrix, CoordinateMatrix, MatrixEntry, RowMatrix
6
7  # Select the dataset
8  #dataset = "/wrk/group/grp-ddi-2020/datasets/data-2-sample.txt"
9  dataset = "/wrk/group/grp-ddi-2020/datasets/data-2.txt"
10
11  conf = (SparkConf()
12         .setAppName("mpsaukko")
13         .setMaster("spark://10.251.52.13:7077")
14         .set("spark.cores.max", "8") ##dont be too greedy ;)
15         .set("spark.rdd.compress", "true")
16         .set("spark.executor.memory", "64G")
17         .set("spark.broadcast.compress", "true"))
18
19  sc = SparkContext(conf=conf)
20  spark = SparkSession.builder.config(conf=conf).getOrCreate()
21
22  # Read the data
23  A = sc.textFile(dataset)
24  A = A.map(lambda s : [float(x) for x in s.split()])
25
26  # Zip index values with cell values
27  A = A.zipWithIndex().map(lambda x: (x[1], x[0]))
28
29  # Print step 1 ready. With full set 1min.
30  print(" ")
31  print("Step 1 ready")
32  print(" ")
33
34  # Conver A to IndexedRowMatrix
35  A = IndexedRowMatrix(A)
36
37  # Convert A to blockmatrices and set block size.
38  #A = A.toBlockMatrix(1000, 1000) # Works with sample set. Data to 1 block.
39  A = A.toBlockMatrix(100, 1000)
40
41  # Cache A, because it is used multiple times
42  A.cache()
43
44  # Print step 2 ready. With full set 3mins.
45  print(" ")
46  print("Step 2 ready")
47  print(" ")
48
49  # Next multiplications. We need to calculate A*AT*A.
50  # Size of the A is 1000000 x 1000 matrix so size of the A*AT would be 100000*100000.
51  # For better performance use matrix multiplication rule (A*AT)*A = A*(AT*A)
52
53  # Calculate A transpose
54  AT = A.transpose()
55
56  # Make first multiplication AT*A
57  ATA = AT.multiply(A)
58
59  # Print step 3 ready. With full set 9mins.
60  print(" ")
61  print("Step 3 ready")
62  print("ATA rows and cols:")
63  print("Rows:", ATA.numRows(), "Cols:", ATA.numCols())
64  print(" ")
65
66  # Make second multiplication A*ATA
67  AATA = A.multiply(ATA)
68
69  # Print step 4 ready
70  print(" ")
71  print("Step 4 ready")
72  print(" ")
73
74  # Convert AATA to indexRowMatrix
75  AATA = AATA.toIndexedRowMatrix()
76
77  # Get first row from AATA
78  first_row = AATA.rows.filter(lambda x: x.index == 0).first().vector.toArray()
79
80  # Print first row
81  print(first_row)

```

Figure 3.1: My Spark matrix multiplication code

4. Conclusions

It was an exciting project, and I learned lots of new and practical skills while writing these codes. Unfortunately, the number of exercises, essays, and projects during the last two has been total horror. Unfortunately, I could not get answers in either of the exercises with a full-size dataset. Nevertheless, I am happy that I could get both codes to work with sample data, especially when considering the limited amount of time I was able to spend on this project.

References

- [1] Apache spark mllib, 2020. [Online; accessed 26-November-2020].
- [2] Wikipedia contributors. Amdahl's law — Wikipedia, the free encyclopedia, 2020. [Online; accessed 26-November-2020].
- [3] Wikipedia contributors. Quickselect — Wikipedia, the free encyclopedia, 2020. [Online; accessed 24-November-2020].