

INF553 Foundations and Applications of Data Mining

Summer 2019

Assignment 2

Deadline: June 10th 11:59 PM PST

1. Overview of the Assignment

In this assignment, you will implement the **SON** algorithm using the Apache Spark Framework. You will develop a program to find frequent itemsets in two datasets, one simulated dataset and one real-world dataset generated from Yelp dataset. The goal of this assignment is to apply the algorithms you have learned in class on large datasets more efficiently in a distributed environment.

2. Requirements

2.1 Programming Requirements

- a. **You must use Python to implement all tasks.** There will be **10% bonus** for each task if you also submit a Scala implementation and both your Python and Scala implementations are correct.
- b. **You are required to only use Spark RDD** in order to understand Spark operations more deeply. You will not get any point if you use Spark DataFrame or DataSet.

2.2 Programming Environment

Python 3.6, Scala 2.11 and Spark 2.3.3

We will use these library versions to compile and test your code. There will be a 20% penalty if we cannot run your code due to the library version inconsistency.

2.3 Write your own code

Do not share code with other students!!

For this assignment to be an effective learning experience, you must write your own code! We emphasize this point because you will be able to find Python implementations of some of the required functions on the web. Please do not look for or at any such code!

TA will combine all the code we can find from the web (e.g., Github) as well as other students' code from this and other (previous) sections for plagiarism detection. We will report all detected plagiarism.

2.4 What you need to turn in

Your submission must be a zip file with name: **firstname_lastname_hw2.zip** (all lowercase). You need to pack the following files in the zip file (see Figure 1):

- a. two Python scripts, named: (all lowercase)

firstname_lastname_task1.py **firstname_lastname_task2.py**

b1. [OPTIONAL] two Scala scripts, named: (all lowercase)

firstname_lastname_task1.scala

firstname_lastname_task2.scala

b2. [OPTIONAL] one jar package, named: **firstname_lastname_hw2.jar** (all lowercase)

c. You don't need to include your results. We will grade on your code with our testing data (data will be in the same format).

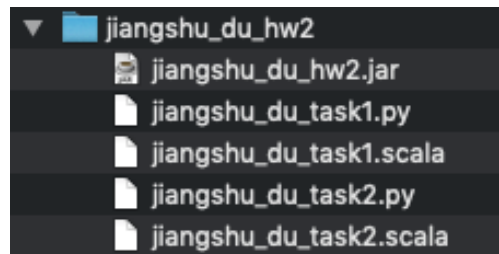


Figure 1: Submission Structure

3. Datasets

In this assignment, you will use one simulated dataset and one real-world. In task 1, you will build and test your program with a small simulated CSV file that has been provided to you. In task 2, you will build and test your program with a large real-world data which generated by a subset using business.json and review.json from the Yelp dataset (<https://www.yelp.com/dataset>) with the same structure as the simulated data. For saving the time, the dataset has been provided to you.

Figure 2: Input Data Format

1	user_id,business_id
2	0x1,0x6b
3	0x1,0x6a
4	0x1,0x67
5	0x1,0x61
6	0x1,0x69
7	0x1,0x64

4. Tasks

In this assignment, you will implement the **SON algorithm** to solve all tasks (Task 1 and 2) on top of Apache Spark Framework. You need to find **all the possible combinations of the frequent itemsets** in any given input file within the required time. You can refer to the Chapter 6 from the Mining of Massive Datasets

book and concentrate on section 6.4 – Limited-Pass Algorithms. You need to use A-Priori algorithm to process each chunk of the data.

4.1 Task 1: Simulated data (6 pts)

There is a CSV file(test_data.csv) posted on the Blackboard.You can use this test file to debug your code.

In this task, you need to build two kinds of market-basket model.

Case 1 (3 pts):

You will calculate the combinations of frequent businesses (as singletons, pairs, triples, etc.) that are qualified as frequent given a support threshold. You need to create a basket for each user containing the business ids reviewed by this user. If a business was reviewed more than once by a reviewer, we consider this product was rated only once. More specifically, the business ids within each basket are unique. The generated baskets are similar to:

user1: [business11, business12, business13, ...]

user2: [business21, business22, business23, ...]

user3: [business31, business32, business33, ...]

Case 2 (3 pts):

You will calculate the combinations of frequent users (as singletons, pairs, triples, etc.) that are qualified as frequent given a support threshold. You need to create a basket for each business containing the user ids that commented on this business. Similar to case 1, the user ids within each basket are unique. The generated baskets are similar to:

business1:[user11,user12,user13,...]

business2:[user21,user22,user23,...]

business3:[user31, user32, user33, ...]

Input format:

1. Case number: **Integer** that specifies the case. **1 for Case 1 and 2 for Case 2.**
2. Support: **Integer** that defines the minimum count to qualify as a frequent itemset.
3. Input file path: This is the path to the input file including path, file name and extension.
4. Output file path: This is the path to the output file including path, file name and extension.

Output format:

1. Runtime: **the total execution time from loading the file till finishing writing the output file**

You need to **print the runtime in the console** with the “Duration” tag, e.g., “Duration: 100”.

3. Apply the SON algorithm code to the filtered market-basket model;

Input format:

1. Filter threshold: **Integer that is** used to filter out qualified users
2. Support: **Integer** that defines the minimum count to qualify as a frequent itemset.
3. Input file path: This is the path to the input file including path, file name and extension.
4. Output file path: This is the path to the output file including path, file name and extension.

Output format:

1. Runtime: **the total execution time from loading the file till finishing writing the output file**

You need to **print the runtime in the console** with the "Duration" tag, e.g., "Duration: 100".

2. Output file

The output file format is the same with task 1. Both the intermediate results and final results should be saved in ONE output result file "firstname_lastname_task2.txt".

Execution example:

Python:

```
spark-submit firstname_lastname_task2.py <filter threshold> <support> <input_file_path>  
<output_file_path>
```

Scala:

```
spark-submit --class firstname_lastname_task2 firstname_lastname_hw2.jar <filter threshold>  
<support> <input_file_path> <output_file_path>
```

5. Evaluation Metric

Task 1:

Input File	Case	Support	Runtime (sec)
test_data.csv	1	4	<u><u><=100</u></u>
test_data.csv	2	7	<u><u><=200</u></u>

Task 2:

Input File	Filter Threshold	Support	Runtime (sec)
AZ_Yelp.csv	70	50	<u><u><=1000</u></u>

6. Grading Criteria

(% penalty = % penalty of possible points you get)

1. You can use your free 5-day extension separately or together.
2. There will be 10% bonus if you use both Scala and Python.
3. If you do not apply the SON algorithm, there will be no point for this assignment.
4. If we cannot run your programs with the command we specified, there will be 80% penalty.
5. If your program cannot run with the required Scala/Python/Spark versions, there will be 20% penalty.
6. If our grading program cannot find a specified tag, there will be no point for this question.

7. If the outputs of your program are unsorted or partially sorted, there will be 50% penalty.
8. If the header of the output file is missing, there will be 10% penalty.
9. We can regrade on your assignments within seven days once the scores are released. No argue after one week. There will be 20% penalty if our grading is correct.
10. There will be 20% penalty for late submission within a week and no point after a week.
11. There will be no point if the execution time in each task exceeds the runtime in Section 6 Evaluation Metric.
12. Only when your results from Python are correct, the bonus of using Scala will be calculated. There is no partially point for Scala. See the example below:

Example situations

Task	Score for Python	Score for Scala (10% of previous column if correct)	Total
Task1	Correct: 6 points	Correct: $6 * 10\%$	6.6
Task1	Wrong: 0 point	Correct: $0 * 10\%$	0.0
Task1	Partially correct: 3 points	Correct: $3 * 10\%$	3.3
Task1	Partially correct: 3 points	Wrong: 0	3.0