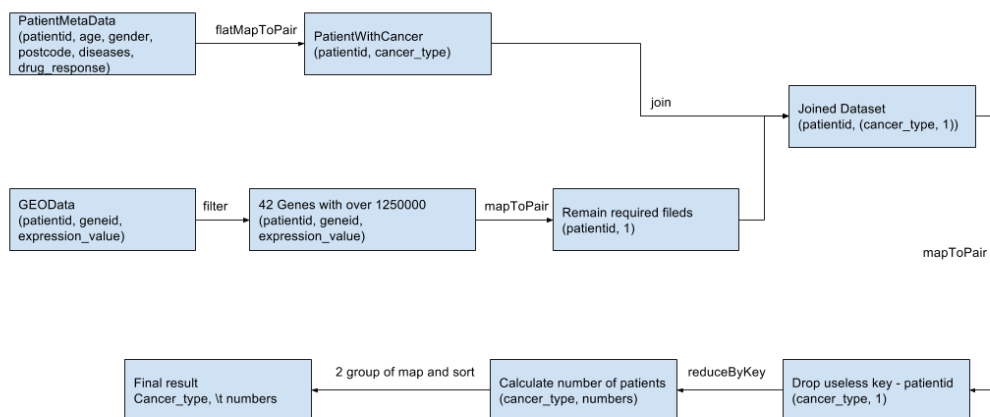


# COMP5349 Assign3

## Task1

In this task, the PatientMetaData file is read in Spark as PairRDD, followed by a flatMapToPair operation to filter the patients who get cancer. After that, the GEO file is read, followed by a filter operation to filter 42 gene with valid express value. In addition, a map function is used to transfer the result to (patientid, 1).

A join operation is executed upon two PairRDDs processed above based on the same key - patientid followed by a mapToPair transformation to adjust the joined result. Then reduceByKey operation is executed based on patientid and two sortByKey function is followed to get the required format.



Performance:

Dataset	Runtime	Run Command On Large Dataset
Small	12s	spark-submit \ --class assign3.Task1 \ --master yarn \ --num-executors 3 \ comp5349assign3-0.0.1.jar \ hdfs://soit-hdp-pro-1.ucc.usyd.edu.au:8020/share/genedata/large/ \ #input path hdfs://soit-hdp-pro-1.ucc.usyd.edu.au:8020/user/cshe6391/assign3/task1large #output path
Large	14s	

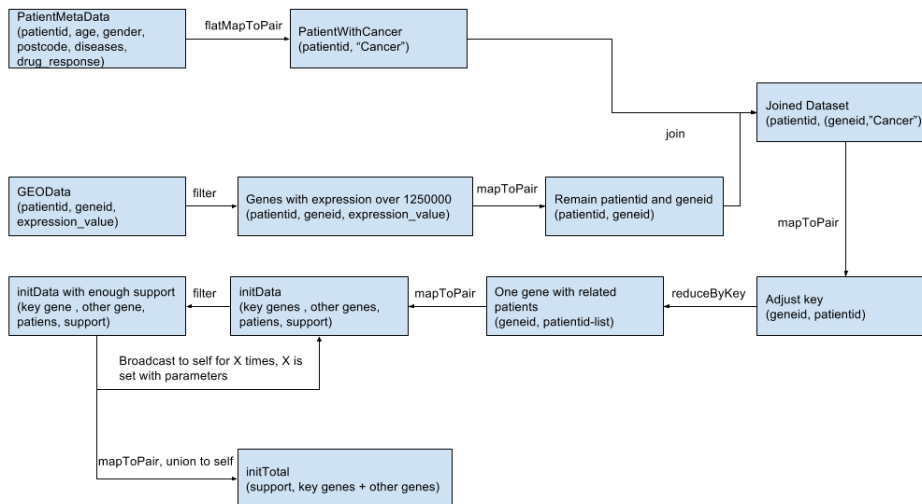
N.B. Remove comment after hash tag before execute command.

## Task2

In this task, the PatientMetaData.txt is read and filtered based on the cancer list. On the other hand, the GEO.txt is read and filtered with valid expression value followed by a map to pair function to drop unrelated fields. A join operation is executed based on patientid.

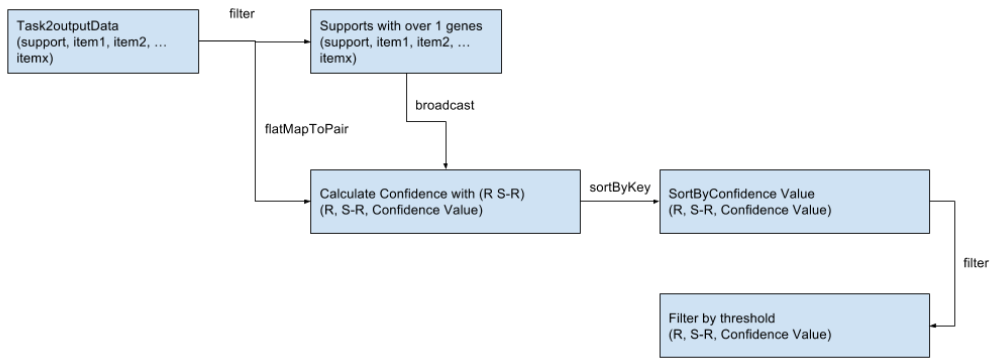
After that, the PairRDD is grouped and reduced by geneid in that patients are collected together. InitData is the initial dataset of loop followed by a filter to drop some entry without enough support. In each loop, the result is broadcast in that they can be used in the next loop.

FrequentItemMapper is a customised mapToPair function to help mine frequent items.



Dataset	Runtime	Run Command On Large Dataset
Small	23s	spark-submit \ --class assign3.Task2 \ --master yarn \ --num-executors 10 \ #in this experiment 10 executors are used
Large	147s	comp5349assign3-0.0.1.jar \ hdfs://soit-hdp-pro-1.ucc.usyd.edu.au:8020/share/genedata/large/ \ #input path hdfs://soit-hdp-pro-1.ucc.usyd.edu.au:8020/user/cshe6391/assign3/task2large \ #output path 0.3 \ #set support percentage 4 \ #set max iteration times

## Task3



Performance:

Dataset	Runtime	Run Command On Large Dataset
Small	13s	spark-submit \ --class assign3.Task3 \ --master yarn \ --num-executors 3 \ comp5349assign3-0.0.1.jar \ hdfs://soit-hdp-pro-1.ucc.usyd.edu.au:8020/user/cshe6391/assign3/task2 \ #input path hdfs://soit-hdp-pro-1.ucc.usyd.edu.au:8020/user/cshe6391/assign3/task3large \ #output path 0.6 \ #set confidence threshold
Large	14s	

N.B. Remove comment after hash tag before execute command.

## Appendix

HDFS Location:

Item	Path
Task 1 small	/user/cshe6391/assign3/task1small
Task 1 large	/user/cshe6391/assign3/task1large
Task 2 small	/user/cshe6391/assign3/task2small
Task 2 large	/user/cshe6391/assign3/task2large
Task 3 small	/user/cshe6391/assign3/task3small
Task 3 large	/user/cshe6391/assign3/task3large