

CS 5433 - Big Data Management
Spring 2018
Big Data Group Project – Phase III

Group: 3

Baath, Harpinder
Polavarapu, Tejaswanth
Shankara, Akshay
Dandamudi, Madhu Kiran

Problem Understanding:

Feature selection is an important aspect of predictive analytics. It not only tends to reduce the complexity of an algorithm but improves performance and cost of data collection there by selecting sensitive variables which can provide significant predictive accuracy. The optimum feature subset selection is a NP hard problem which cannot be solved in polynomial time, genetic algorithms provide near optimal solution to this problem. Genetic algorithms are motivated from Darwin's theory of survival of the fittest and mimic biological reproduction process. It starts with random population of individuals (represented as chromosomes) and then evaluate the fitness of each individual. Two individuals are selected on the basis of their fitness value to serve as parents for crossover intending to produce better offsprings. Individuals with higher fitness value are more likely to be selected. Genetic algorithms need a lot of computational power and resources, fortunately, hadoop MapReduce allow parallelism which can be used to run genetic algorithms in parallel and obtain the desired result with improved efficiency.

In feature selection problem, each feature represents a gene and the collection of genes is a chromosome (subset of features). Each chromosome is represented by a string of 0s and 1s. 0; if feature is absent, otherwise 1.

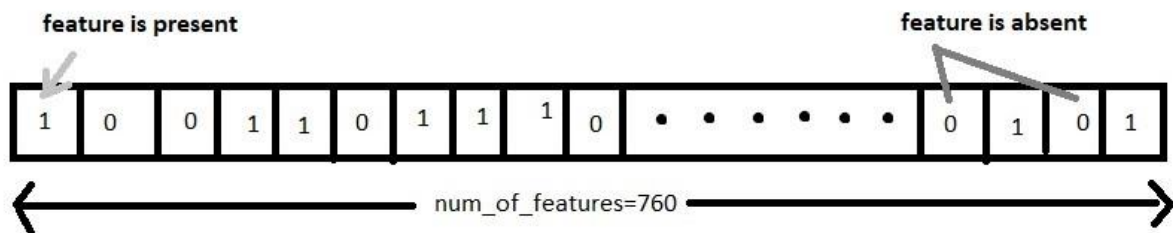


Figure 1: Representation of an individual

The population is a set of randomly selected chromosomes. With the help of the training and the test dataset obtained from the HDFS, the 'accuracy' (fitness) of each individual is calculated using machine learning algorithms such as decision tree, support vector machine etc. For each reproduction, the parents with the best fitness are selected from the pool of population using 'roulette wheel' phenomenon and a certain crossover criteria is used to produce new subsets which include features from both the parent subsets and some error (mutation). With each generation of population, we aim to achieve a better fitness and for this, we select a better set of individuals.

Implementation of GA in MapReduce

Genetic algorithm for feature selection problem would work in following steps:

1. Initialize the population with random individuals (subset of features).
2. Evaluate the fitness value of the individuals using MLib classification algorithms.
3. Select good solutions by using Tournament selection criteria (selection without replacement).

4. Create new individuals by recombining the selected population using uniform crossover.
5. Evaluate the fitness value of all offspring.
6. Repeat steps 3–5 until the convergence criteria is met.

Map()

Evaluation of the fitness function for the population (Steps 2 and 5) matches the MAP function, which has to be computed independent of other instances. As shown in the algorithm in Algorithm 1, the MAP evaluates the fitness of the given individual. It also keeps track of the the best individual and finally, writes it to a global file in the Distributed File System (HDFS). The client, which has initiated the job, reads these values from all the mappers at the end of the MapReduce and checks if the convergence criteria has been satisfied.

Algorithm 1: Map phase of each iteration of GA *map(key, value):*

```

individual ← ENCODED_INDIVIDUAL(key)
fitness ← CALCULATE_FITNESS(individual)
EMIT (individual, fitness)
//Keep track of the current best

if fitness > max_fitness then
    max_fitness ← fitness
    max_individual ← individual
end if

if all individuals have been processed then
    Write best individual to global file in HDFS
end if

```

Reduce()

We will be using Tournament selection without replacement. A tournament is conducted among S randomly chosen individuals and the winner is selected. The reduce function goes through the individuals sequentially, the individuals from the last round are buffered first. When the tournament window is full, selection and crossover is carried out as shown in the Algorithm 3. When the crossover window is full, we would use the Uniform Crossover operator. The value of S would be finalized by trying different values and evaluating the result.

Algorithm 3: Reduce phase of each iteration of the GA

```

Initialize processed_individuals ← 0,
tournArray [2· tSize], crossArray [cSize]
reduce(key, values):

while values.hasNext() do
    individual ← ENCODED_INDIVIDUAL(key)

```

```

fitness ← values.getValue()
if processed_individuals < tSize then

    //Wait for individuals to join in the tournament and put them for the last rounds
    tournArray [tSize + processed_individuals % tSize] ← individual
else

    //Conduct tournament over past window
    Selection_Crossover()
end if
processed_individuals ← processed_individuals + 1
if all individuals have been processed then

    //Cleanup for the last tournament windows
    for k=1 to tSize do

        Selection_Crossover()
        processed_individuals ← processed_individuals+1
    end for
end if
end while

Selection_Crossover:
crossArray[processed_individuals % cSize] ← TOURN(tournArray)
if (processed_individuals - tSize) % cSize = cSize - 1 then
    newIndividuals ← CROSSOVER(crossArray)
    for individual in newIndividuals do
        EMIT (individual, dummyFitness)
    end for
end if

```

Results

Results obtained from Decision Tree and Logistic Regression are different. In addition, the number of iterations each algorithm took to get best individual are also different, that can be observed in figures below.

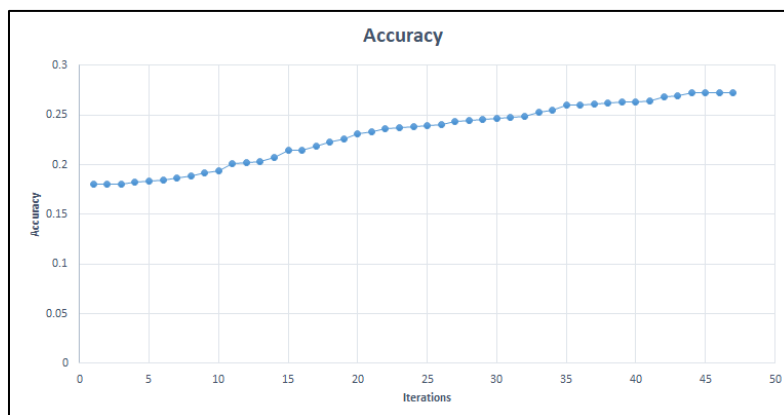


Figure 3: Effect of iterations on the Accuracy for logistic regression

Decision tree with 376 features got the accuracy of 0.272 and logistic regression with 412 features got the accuracy of 0.226.

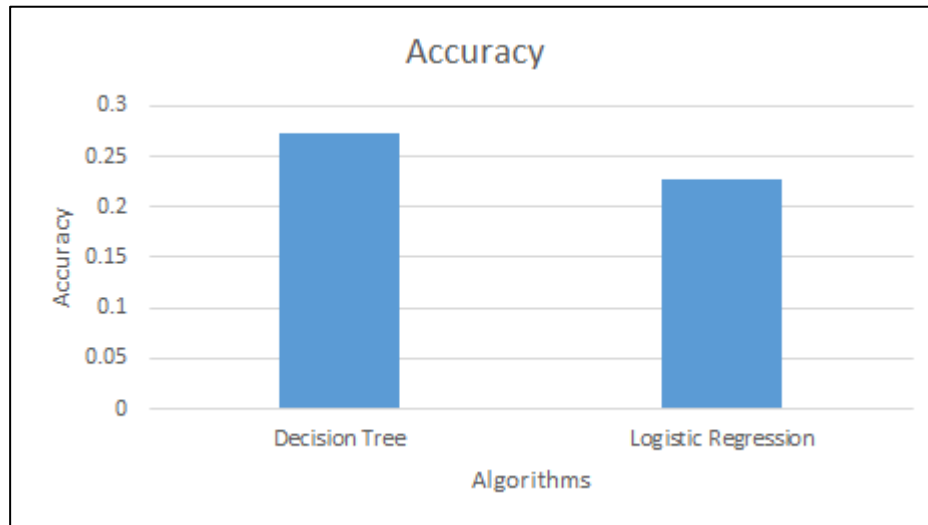


Figure 4: Prediction accuracy for optimal subset by two prediction algorithms

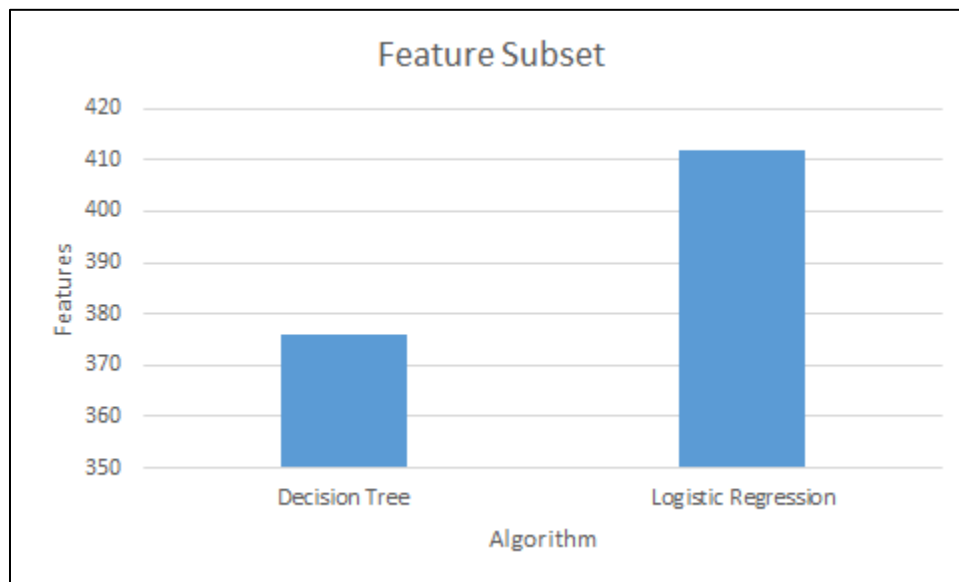


Figure 5: Size of optimal feature subset for two prediction algorithms

The optimal feature subset that we obtained from GA using decision tree algorithm is of size 376. Figure 6 shows the features in the subset. The subset obtained using logistic regression is of size 412, features included are shown in figure 7.

