# Brute Force,Forward & Backward Feature Selection Algorithm

| Version | 1.0 |
|---|---|
| Document Title | Final Report of Big Data Group Project |
| Group Members | Tanuj Gyan, Sayanti Ray, Raviprakash Shreeraksha, Babak Shabani |
| Date of Submission | 1-May-2018 |
| Instructor Name | Dr. Johnson P Thomas |
| TA Name | Vasanthi Mudunuri |

**Brute Force Feature Selection Algorithm:**

This is a trial and error approach where we consider all possible subset of features and find the accuracy of these subsets with the help of Machine Learning Algorithms. All Machine Learning Algorithms used are supervised machine learning algorithms and so we need to have both Test and Train data to calculate accuracy.

So we have used Map Reduce of Hadoop to implement Brute Force Algorithm.First we will consider the number of features given and then calculate all possible subsets of these features,avoiding the redundancy of the subsets.In the Mapper, we have taken each subset of features and then in the reducer we output the subset.

Below are the snapshots of Map Reduce program implementing Brute Force Algorithm.



Fig 1: Map Reduce program in execution



Fig 2: Creation of Subsets

Above screenshots display the execution of MapReduce program and creation of subsets for 5 features.

Fig 3: Output

## Challenges Faced :

➢ There are total 760 features and we had to create all possible subsets ,but the data is large and it was not feasible to run Brute Force algorithm as the system would crash.

➢ When implementing Decision tree ,after Stage 8 or 12 the data contained for that task was very large and exceeded the task size.Please refer the below snapshot.

```
>>> from pyspark.mllib.regression import LabeledPoint
>>> from pyspark.mllib.tree import DecisionTree
>>> from pyspark.mllib.classification import LogisticRegressionWithSGD
>>> from pyspark.mllib.classification import SVMWithSGD
>>>
>>> train_data = sc.textFile('hdfs://hadoop1:9000/CS5433/Group_Project/train_data.csv').map(lambda x: x.split(',')).map(lambda line: [float(v
>>> train_data =train_data.collect()
>>> train_labels = sc.textFile('hdfs://hadoop1:9000/CS5433/Group_Project/train_labels.csv')
>>> train_labels =train_labels.collect()
>>> train_labels=[int(x) for x in train_labels]
>>> train_data_list=zip(train_labels,train_data)
>>> train_data_list=sc.parallelize(train_data_list)
>>>
>>> a=train_data_list.map(lambda l: LabeledPoint(l[0], l[1:]))
>>>
>>> a=a.collect()
18/05/01 19:08:53 WARN TaskSetManager: Stage 12 contains a task of very large size (27430 KB). The maximum recommended task size is 100 KB.
>>>
```

## Forward Algorithm :

Initially the feature set is subjected to be empty and then in FFS we keep adding the features as we progress.

Here we used a threshold of 0.03 to maintain the accuracy within a limit, without that value obtained from subset of features will give the maximum accuracy which is equivalent to the original accuracy of the machine learning algorithm.

**Backward Algorithm:**

Initially the feature set is subjected to be full and then in BFS we keep subtracting the features as we progress.

Both Forward and Backward Algorithm use feature importance to select the important features in the feature set.In the Forward algorithm we keep on adding the important feature based on feature importance and in the Backward algorithm we keep on removing the unimportant feature.