

# Restaurant Recommendation Challenge

Making prediction by Hadoop and Spark

Bohua Jia

University of Stavanger, Norway  
261858@stud.uis.no

Dipanjana Banik

University of Stavanger, Norway  
261860@stud.uis.no

## ABSTRACT

Prediction is an indispensable part of people's lives, ranging from weather forecasts to national development prospects. A decision tree is an entry-level mechanics algorithm that, while simple and easy to understand, can make very complex and accurate forecasts. Google used decision trees to make market forecasts around 2014 and made money for doing so. This article will present some core functions for the decision tree we implemented using Python. Moreover, using Hadoop mrjob and Spark to handle large datasets. We aim to make a performance comparison of the different data frames.

## KEYWORDS

Decision Tree, MrJob, Spark

## 1 INTRODUCTION

*Overview.* Enjoying local cuisine could be one of the driving forces behind tourism, and providing travellers with information on local cuisine and restaurant recommendations could be beneficial to them when determining where to go. They are the sorting and searching functions that assist travellers in locating the best local cuisine and restaurants among the vast amount of information available. They can also assist tourists in deciding where to travel.

*Research Problem.* In today's modern world, we check the ratings of a restaurant before planning to dine. Today, people are more aware than ever of the link between health and nutrition and are more cognizant of what foods they eat when they dine at restaurants. Another modern trend is local foods, with a special emphasis on sustainability. So it is very important to improve the ratings or recommend a restaurant to a customer. Restaurant marketing is the process of getting people to visit a restaurant. Restaurant marketing creates loyalty, provides data for research and analytics, and allows restaurants to better understand their ideal customer profile. As a result, restaurant marketing is extremely important: it allows us to attract new consumers, increase our client base, and turn existing drop-in customers into regulars. Furthermore, the COVID-19 epidemic has emphasized the importance of this even more.

*Contribution summary.* Using the restaurant recommendation dataset as a starting point, we deployed Machine Learning techniques. We created our own Decision Tree algorithm from scratch and utilized it in the MapReduce and Spark framework to achieve our goals. In addition, we used the sklearn and Spark decision trees to forecast suggestions for a consumer or a vendor to employ. Specifically, Hadoop and Spark are data processing frameworks built to

deal with large amounts of data. After that, we looked at the algorithms' performance and demonstrated that our implementation is scalable, and we came to the same conclusion.

The source of our project is available here:  
[https://github.com/ProjBH/DAT500\\_project\\_mrjob\\_decision\\_tree](https://github.com/ProjBH/DAT500_project_mrjob_decision_tree)

## 2 BACKGROUND

### 2.1 Hadoop

Hadoop [1][2] is a distributed file system and framework that uses the MapReduce paradigm to analyze and convert very large data collections.[3] The most important feature of Hadoop is that it can distribute data to different hosts in the same cluster, thereby greatly improving work efficiency. A Hadoop cluster can grow processing ability, storage capacity, and IO throughput by simply adding commodity servers. In other words, the more servers in the Hadoop cluster, the more powerful it has.

### 2.2 HDFS

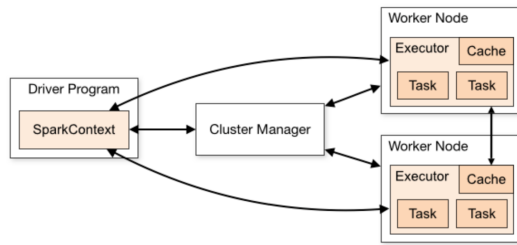
HDFS (Hadoop distributed file system) separates file system metadata from application data. The NameNode, a dedicated server in HDFS, holds system metadata. DataNodes are additional servers that store application data. All servers are ultimately linked and use TCP-based protocols to communicate with others.

#### 2.2.1 How does it work.

In HDFS, the file's content is divided into large chunks, and each chunk is replicated at several DataNodes separately. The NameNode is responsible for keeping track of the namespace tree and mapping file chunks to DataNodes (the physical location of file data). When an HDFS client wants to read a file, it first asks the NameNode for the locations of the data chunks that make up the file and then reads the contents of those chunks from the DataNode nearest to the client. The client asks the NameNode to choose three DataNodes to host the block copies when writing data. The client then uses a pipeline to write data to the DataNodes. Each cluster has a single NameNode in the present configuration. As each DataNode can perform many application activities parallelly. So, the cluster can include thousands of DataNodes and ten times more clients than Datanodes.

### 2.3 Spark

An open-source parallel processing framework that allows in-memory processing to speed up large-scale data analysis applications. Spark is implemented in Scala [4][5]. Unlike Hadoop, Spark is capable of storing intermediate data in memory.[6]



**Figure 1: Apache Spark architecture**

### 2.3.1 Driver.

The driver consists of a Spark session. The Spark session takes the task and divides it into smaller tasks that are handled by the worker nodes.

### 2.3.2 Worker Nodes.

Each worker node receives a task distributed by the driver and executes the task. Those worker nodes are known as a cluster.

### 2.3.3 Cluster Manager.

The cluster manager communicates with both the driver and the worker nodes to: manage resource allocation, manage program division and manage program execution [7].

## 2.4 Spark MLlib

Spark MLlib is a Spark Core package that offers machine learning algorithms and MLlib is built for ease of use, scalability, and interaction with other tools. It includes scalable and quick implementations of standard machine learning models. Spark Machine Learning library MLlib contains several different applications like - Regression for Predictions, Clustering, Classification Algorithms, Dimensionality Reduction, and Regression for Predictions. Spark's base computing foundation is a significant benefit. Working on a large-scale machine learning project becomes much easier due to this [9, 10].

## 2.5 Scikit Learn library (sklearn)

Scikit Learn library has a module function DecisionTreeClassifier() for easily implementing decision tree classifier. In our project, we used the decision tree from sklearn to compare the prediction result.

## 2.6 Pandas and Spark Datafram

Dataframes contain a table of data with rows and columns. However, Spark Dataframe and Pandas Dataframe are not the same things. Pandas is an open-source Python library based on the NumPy library. It is a Python package that lets us handle numerical data and time series using various data structures and operations. Its primary goal is to make data import and analysis easier. Pandas DataFrame is a two-dimensional, size-mutable tabular data format with labelled axes that might be different (rows and columns). So there are the three fundamental components of a Pandas DataFrame: data, rows, and columns [21].

Spark engine designed for large-scale data processing. It outperforms other cluster computing systems in terms of speed, such as Hadoop. In addition, Spark has several high-level APIs so that

it is simple to write parallel jobs in Spark. In Spark, DataFrames have distributed data sets that are organized into rows and columns. Each column in a DataFrame is given a name and a type [21].

When we need to manipulate large volumes of data, we utilize Spark. There are several nodes in a Spark DataFrame. Unlike Pandas, Spark DataFrame ensures fault tolerance. Spark DataFrame is also immutable. Furthermore, Spark DataFrame operations are quicker than pandas since they run parallel on all computers' cores.

## 2.7 Related works

### 2.7.1 Supervised learning.

The machine learning task of learning a function that translates an input to an output based on example input-output pairs is known as supervised learning. [11] It uses labeled training data and a collection of training examples to infer a function.[12] Each example in supervised learning is a pair of an input item and the desired output value. A supervised learning algorithm examines the training data and generates an inferred function that can be applied to new cases. The algorithm will be able to accurately determine the class labels for unseen examples in the best-case scenario. This necessitates the learning algorithm to "reasonably" generalize from the training data to unknown conditions.

### 2.7.2 Random forest.

Random Forest is a computationally efficient technique that can operate quickly over large datasets.[13] Random decision forests are also known as random forests. It is an ensemble learning approach for classification, regression, and other problems that work by training a large number of decision trees. For classification tasks, the random forest's output is the class chosen by the majority of trees. For regression tasks, the mean or average prediction of the individual trees is returned.[14][15] Random forests correct for decision trees' habit of overfitting to their training set.[16] Random forests generally produce better results than decision trees. However, data characteristics can affect their performance.[17][18]

## 3 RELEVANT ALGORITHMS

### 3.1 Decision Tree

A decision tree is a form of supervised learning algorithm that may be used to solve issues in both regression and classification. The algorithm creates rules using training data that may represent a tree structure. Like any other tree representation, it has a root node, internal nodes, and leaf nodes. For example, the condition on attributes is represented by the internal node, the branches represent the condition's outcomes, and the leaf node represents the class label. Following the if-else style principles, we can start at the top with the root node and work down to the leaf node to obtain the classification. The class label for our classification challenge is the leaf node where we stop. Both category and numerical data can be used in a decision tree. On the other hand, other machine learning methods cannot deal with categorical input and require encoding to numeric values[20]

One of the predictive modeling methods used in machine learning is decision trees. Utilizing a Decision Tree aims to build a training model that can predict the class or value of a target variable by learning basic decision rules from training data. In Decision Trees,

we begin at the tree's root to forecast a class label for a record. The values of the root attribute are compared to the values of the record's attribute. We follow the branch corresponding to that value and proceed to the next node based on the comparison[7, 8].

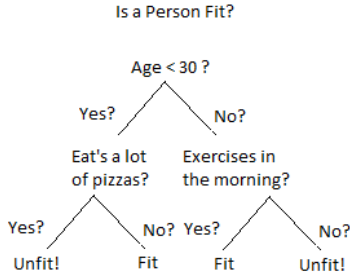


Figure 2: Decision Tree[8]

Classification trees and Regression trees are the two primary forms of Decision Trees that we can create. The fundamental problem in a Decision Tree is determining the attribute for each level's root node, which is known as attribute selection. The Gini index and Information gain are the two most well-known attributes.

The entropy varies as we use a node in a decision tree to split the training examples into smaller groups, and information gain measures this change. Entropy is a measurement of a random variable's ambiguity; it describes the impurity of a random collection of instances. The greater the entropy, the more information is contained. On the other hand, the Gini index is a statistic that measures how often a randomly picked element is mistakenly identified. It indicates that elements of a lower Gini index characteristic should be preferred.[19]

In this project, we have implemented a decision tree algorithm that can recommend a vendor to a customer or a customer to a vendor in both ways. We have implemented this algorithm on Hadoop architecture to see how the clustering works in a distributed environment. In mapper, we used our own decision tree, which we built from scratch and spark implementation; we used the spark MLlib to predict the result.

## 4 IMPLEMENTATION

This section will present the datasets utilized in the research and the implementation of the decision tree in this part. And also outline the difficulties encountered and the decisions that have been made to solve those challenges in the course of work.

### 4.1 Dataset

The dataset used in this article is Restaurant Recommendation Challenge from kaggle[22]. The dataset contains several tables. However, during the project, only two tables were involved that are train\_full and test\_full, as these two tables include every information offered by the other tables. The main difference between train\_full and test\_full is a column called target, which would be the prediction the decision tree will make. Besides the target column, both tables

contain 72 columns. The train\_full table has 3.1 GB in size, and the test\_full table has 902.52 MB.

### 4.2 Preprocessing

The goal of the preprocessing stage in our project is to fill the missing values in the testing and the training dataset. These two datasets contain various information about customers and vendors. As per our observation, we found that many rows have missing values and for our project, we did not drop any of the columns and rows with missing values. We applied some different techniques to fill the missing values. Also, in preprocessing stage, we had to deal with some datetime column types. At the final preprocessing stage, we changed all the categorical values to numerical so that the implementation of the decision tree could be done correctly and gives the prediction and accuracy more precisely.

#### 4.2.1 Filling empty values.

There must be no empty or null value in the dataset to build the decision tree as a prerequisite. So in the first stage, we filled the missing values with two different strategies. For the string datatype column, we filled the values with the most frequent values; for the integer and double type column, we filled the values with median strategies.

---

```

1 for dataTypes in sparkDF.dtypes:
2     if dataTypes[1] == "string":
3         frequentValueinColumnName = \
4             sparkDF.freqItems([dataTypes[0]], 0.60)
5         frequentValue = frequentValueinColumnName.collect()[0][
6             dataTypes[0] + "_freqItems"
7         ]
8     if len(frequentValue) == 0:
9         # if there is no missing value
10        # then continue to next logic
11        continue
12    else:
13        # get the most frequent value and
14        # replace the null/missing filed with that value
15        sparkDF = \
16            sparkDF.fillna(str(frequentValue[0]),
17                           subset=[dataTypes[0]])
18    elif dataTypes[1] == "int":
19        # fill the null/missing value with median
20        imputer = Imputer(
21            inputCols=[dataTypes[0]],
22            outputCols=[dataTypes[0]],
23        ).setStrategy("median")
24        sparkDF = imputer.fit(sparkDF).transform(sparkDF)
25    elif dataTypes[1] == "double":
26        # fill the null/missing value with median
27        imputer = Imputer(
28            inputCols=[dataTypes[0]],
29            outputCols=[dataTypes[0]],
30        ).setStrategy("median")
31        sparkDF =
32        imputer.fit(sparkDF).transform(sparkDF)
  
```

---

The empty datetime values replaced with the 00:00:00 timestamp.

```

1 import calendar
2 for eachDay in calendar.day_name:
3     eachDay = eachDay.lower()
4     sparkDF = sparkDF.fillna(
5         value="00:00:00",
6         subset=[
7             eachDay + "_from_time2",
8             eachDay + "_from_time1",
9             eachDay + "_to_time2",
10            eachDay + "_to_time1",
11        ],
12    )

```

#### 4.2.2 Replacing values.

Also, we replaced some of the column values with different symbols to use for a later stage.

```

1 # replace vendor_tag column value from , to -
2 # replace vendor_tag column value from , to -
3 sparkDF = sparkDF.withColumn(
4     "vendor_tag",
5     regexp_replace("vendor_tag", ",", "-")
6 )
7 # replace vendor_tag_name column value from , to -
8 sparkDF = sparkDF.withColumn(
9     "vendor_tag_name",
10    regexp_replace("vendor_tag_name", ",", "-")

```

#### 4.2.3 Function for counting number of missing cell value.

We also wrote a function to count the number of empty cell values in a column. However, we did not run it during the program execution because counting empty values in a column of big data set costs much time and is not practical, so we used the above methods to fill the dataset.

```

1 def countIfColumnisNull(sparkDF, columnName):
2     checkNumberOfEmptyCells = sparkDF.select(
3         [count(
4             when(
5                 isnan(c) | col(c).isNull(), c)
6             ).alias(c) for c in sparkDF.columns]
7     )
8     totalEmptyCellCount = \
9     checkNumberOfEmptyCells.collect()[0][columnName]
10    return totalEmptyCellCount

```

#### 4.2.4 Save the output in HDFS.

After finishing the filling the missing values we saved the output in HDFS directory.

```

1 sparkDF.coalesce(1).write.format(
2     "com.databricks.spark.csv").option(
3     "header", "true").mode(
4     "overwrite").save(

```

```

"/dis_materials/testfull_filledValues_output"
)

```

#### 4.2.5 Cleaning.

After thoroughly filling all missing data, the next step is to massage the dataset. In this step, we only keep the reasonable proper columns, merge and do calculations between related columns to make the data set tidier, but the information is more concentrated. All information from each row value has a unique key, then passed this key-value pair to the reducer.

```

1 def mapper(self, _, line):
2     random.seed(datetime.now())
3     key = random.random()
4     rows = csv.reader([line])
5     for row in rows:
6         if row[1] != "gender":
7             ...
8             create_at_x = row[4].split(' ')[0]
9             ...
10            sunday_from_time1 = pd.to_datetime(row[30])
11            sunday_to_time1 = pd.to_datetime(row[31])
12            sunday_from_time2 = pd.to_datetime(row[32])
13            sunday_to_time2 = pd.to_datetime(row[33])
14            temSun = ((sunday_to_time1-sunday_from_time1)+
15            (sunday_to_time2-sunday_from_time2))
16            temSun = temSun.total_seconds()/3600
17            sunday = str(round(temSun,1)).split(' ')[-1]
18            ...
19            primary_tags = row[58]
20            primary_tags = re.sub(r"[\\"\\'\\'"]", '', primary_tags)
21            ...
22            tem = vendor_tag.split('-')
23            vendor_tag_sum = str(len(tem))
24            ...
25            target = row[72]
26            tu = (customer_id,...,target)
27            yield (key,tu)
28
29 def reducer(self, key, values):
30     yield (key, list(values))

```

The output:

```

0.00018626161646613829 [{"K5760F8", "Male", "1", "1", "2019-01-17", "2019-01-17", "0", "Home", "84", "-8.4023", "-78.56", "-1.005",
"0.0707", "Restaurants", "2.0", "13.0", "1.0", "15", "0.0", "Yes", "0.0", "1.0", "1", "11", "88", "4.3", "12.0", "12.0",
"12.0", "12.5", "7.0", "15.0", "1.0", "40", "2018-09-16", "2020-04-07", "3", "0", "0"]}
0.0010880878473140229 [{"UJCZTD", "Male", "1", "1", "2019-10-03", "2019-11-16", "1", "Home", "154", "0.05985", "0.086977", "-0.0
72", "0.0094", "Restaurants", "2.0", "15.0", "1.0", "150", "0.0", "Yes", "0.0", "1.0", "1", "11", "88", "4.5", "13.2", "13.2", "13.2",
"13.2", "13.2", "11.2", "13.0", "1.0", "45", "2019-01-17", "2020-04-02", "3", "1", "0"]}
0.001088083518518392 [{"MUVAC26", "Female", "1", "1", "2019-08-26", "2019-10-01", "0", "Other", "113", "-0.6164", "-0.12354", "0
.6384", "0.5209", "Restaurants", "2.0", "15.0", "1.0", "15", "0.0", "Yes", "0.0", "1.0", "1", "11", "88", "4.7", "12.0", "12.0", "12
.0", "12.0", "12.0", "10.0", "11.0", "1.0", "44", "2018-10-24", "2020-04-07", "3", "0", "0"]}

```

Figure 3: Cleaning the data

#### 4.2.6 Formatting.

In this step, we are going to format the previous output. All the keys will be removed, and the rest of the data will be written as a CSV type of file.

```

1 for line in sys.stdin:
2     line = line.strip()

```

```

3     value = line.split("\t")
4     if len(value)>=2:
5         value = ((value[1].lstrip("[")
6                 .rstrip("]")).replace(" ", ""))
7         value = value.replace("\\""", "")
8         print(value)

```

The output:

```

116F5, Male, 1, 1, 2018-03-15, 2018-03-15, 0, Week, 104, -0.688, -76.5, -0.83244, 0.656, Restaurants, 2, 0, 2, 0, 1, 0, 10, 0, 0, Yes, 0, 0, 1, 0, 1, 11, EN, 4,
5, 13, 5, 13, 5, 13, 5, 13, 5, 11, 5, 13, 5, [primary_tags:134], 1, 0, 1, 2018-10-19, 2020-04-03, 3, 0, 0
116F5, Male, 1, 1, 2018-03-15, 2018-03-15, 0, Week, 105, -0.688, -76.5, -0.968, 0.888, Restaurants, 2, 0, 15, 0, 1, 0, 12, 0, 0, Yes, 0, 0, 1, 0, 1, 11, EN, 4,
5, 13, 5, 13, 5, 13, 5, 13, 5, 11, 5, 13, 5, [primary_tags:134], 1, 0, 1, 2018-10-19, 2020-04-03, 2, 0, 0
116F5, Male, 1, 1, 2018-03-15, 2018-03-15, 0, Week, 106, -0.688, -76.5, 0.6187, 0.5273, Restaurants, 2, 0, 15, 0, 1, 0, 10, 0, 0, Yes, 0, 0, 1, 0, 1, 11, EN, 4,
5, 13, 2, 13, 2, 13, 2, 13, 2, 11, 2, 13, 2, [primary_tags:32], 1, 0, 1, 2018-10-19, 2020-04-06, 3, 0, 0

```

Figure 4: Formatting the data

### 4.3 Converting categorical data to numerical

Categorical characteristics are string-type data that humans can easily understand. Machines, on the other hand, cannot directly understand categorical data. As a result, categorical data must be converted into numerical data that machines can understand. Machine learning models are unable to understand categorical data. As a result, the conversion to numerical format is required. There are several methods for converting category data to numerical data. The two most famous method is Label Encoding and One Hot Encoding. Our project used the Label Encoding method to convert all the string type features into numerical features. Categorical to numerical conversion is mainly used for the Python scikit-learn and Spark MLlib package to make the decision tree. In Hadoop, we use the Decision Tree, which was built from scratch and it can process categorical features without any conversion. Some categorical features we have reduced to one feature. For example, some columns in both the test and training dataset have a four column value for each day(*sunday\_from\_time1*, *sunday\_to\_time1*, *sunday\_from\_time2*, *sunday\_to\_time2*). So we reduced it to only one column (*sunday\_opening\_time*) and made the value only as a numerical opening hour.

```

1 for eachDay in calendar.day_name:
2     eachDay = eachDay.lower()
3     subset = [
4         eachDay + "_from_time2",
5         eachDay + "_from_time1",
6         eachDay + "_to_time2",
7         eachDay + "_to_time1",
8     ]
9     for col_name in subset:
10        sparkDF = sparkDF.withColumn(
11            col_name, date_format(
12                col(
13                    col_name), "HH:mm:ss"
14                ).cast("timestamp")
15        )
16    sparkDF = sparkDF.withColumn(
17        eachDay + "_opening_time",
18        round(
19            (
20                unix_timestamp(col(subset[3]))

```

```

21        - unix_timestamp(col(subset[1]))
22        + unix_timestamp(col(subset[2]))
23        - unix_timestamp(col(subset[0]))
24    )
25    / 3600,
26    2,
27    ),
28    )
29    for col_name in subset:
30        sparkDF = sparkDF.drop(col_name)

```

We have calculated some of the vendors' features manually and converted them into a numerical column.

```

1 sparkDF = sparkDF.withColumn(
2     "vendor_tag", size(
3         split(
4             col("vendor_tag"), r"\-")
5     )
6     sparkDF = sparkDF.withColumn(
7         "vendor_tag_name", size(
8             split(
9                 col(
10                    "vendor_tag_name"
11                ), r"\-")
12            )
13        )

```

Finally, we have indexed all the columns with string datatype by the Label Encoding method. The translation of categorical data into numerical data that a computer can interpret is referred to as label encoding. In this system, each type of data is allocated a number beginning with 0. For this purpose, we used the Label Encoder from sklearn to translate these characteristics into numerical values.

```

1 for dataTypes in sparkDF.dtypes:
2     if dataTypes[1] == "string":
3         indexer = StringIndexer(
4             inputCol=dataTypes[0],
5             outputCol=dataTypes[0] + "_indexed"
6         )
7         sparkDF = indexer.fit(sparkDF).transform(sparkDF)
8         sparkDF = sparkDF.drop(col(dataTypes[0]))
9         sparkDF = sparkDF.withColumnRenamed(
10            dataTypes[0] + "_indexed", dataTypes[0])

```

### 4.4 Decision Tree

A decision tree is a proper supervised learning technique. It's essentially a categorization issue. It is a diagram in the shape of a tree that depicts a line of activity. It is made up of nodes and leaf nodes. The conclusion is drawn using these nodes and leaf nodes.

#### 4.4.1 Entropy.

$$E(s) = \sum_{i=1}^c -p_i \log_2 p_i$$

The entropy of a decision tree measures the purity of the splits. It determines the level of impurity and disorganization. Therefore, it is beneficial to make decisions in a decision tree. For example, it



assists in predicting which node will divide first based on entropy levels.

---

```

1 def calculate_entropy(data):
2     label_column = data[:, -1]
3     _, counts = numpy.unique(label_column,
4     return_counts=True)
5     pro = counts / counts.sum()
6     entropy = sum(pro * -numpy.log2(pro))
7     return entropy

```

---

#### 4.4.2 Making prediction.

The prediction is made by applying the row of data with the decision tree.

---

```

1 def make_prediction(row, tree):
2     question = list(tree.keys())[0]
3     feature_name, comparison_operator, value =
4     question.split(" ")
5     # ask question
6     try:
7         # feature is continuous
8         if comparison_operator == "<=":
9             if row[feature_name] <= float(value):
10                 answer = tree[question][0]
11             else:
12                 answer = tree[question][1]
13         # feature is categorical
14         else:
15             if str(row[feature_name]) == value:
16                 answer = tree[question][0]
17             else:
18                 answer = tree[question][1]
19         # base case
20         if not isinstance(answer, dict):
21             return answer
22         # recursive part
23         else:
24             residual_tree = answer
25             return make_prediction(row, residual_tree)

```

---

#### 4.4.3 Calculate accuracy.

After using the training data to train the decision tree, it is necessary to have a function to test the accuracy of the result.

---

```

1 def calculate_accuracy(df, tree):
2     df["classification"] = df.apply(make_prediction,
3     axis=1, args=(tree,))
4     df["classification_correct"] = df["classification"]
5     == df["target"]
6     accuracy = df["classification_correct"].mean()
7     return accuracy

```

---

#### 4.4.4 Splitting training dataset.

Since the test dataset has no target column, the only way to determine whether the predictions are correct is to split the training set into two subsets—one for training the decision tree and another

for testing the predictions. In that case, we decided to split the dataset randomly. It can avoid constantly training the tree with a similar dataset if the accuracy is not satisfactory.

---

```

1 def train_test_split(df, sub_size):
2     if isinstance(sub_size, float):
3         test_size = round(sub_size * len(df))
4     indices = df.index.tolist()
5     test_indices = random.sample(population=indices,
6     k=test_size)
7     train_df = df.loc[test_indices]
8     test_df = df.drop(test_indices)
9     return train_df, test_df

```

---

## 4.5 Applying Decision Tree algorithm with MapReduce framework

### 4.5.1 Splitting.

Since the data size is too large for a standard python file, splitting the training dataset has to be done by mrjob. We decided to use a random function with a current timestamp as the random seed to produce the random number to indicate whether the current row belongs to the sub-training set or sub-testing set.

---

```

1 def mapper(self, _, line):
2     random.seed(datetime.now())
3     tag = random.random()
4     line = line.strip()
5     row = line.split(',')
6     if row[0] != "customer_id":
7         num = random.randint(0,10)
8         if num <=1:
9             key = 'Test'+str(tag)
10            yield (key,row)
11        else:
12            key = 'Train'+str(tag)
13            yield (key,row)
14
15 def reducer(self, key, values):
16     yield key, list(values)

```

---

After sub-datasets have been produced, the output again needs to format.

### 4.5.2 Formatting.

The main difference between formatting for sub-datasets and original datasets is when the reading line from the input file, the keys are different. And in this step, we will read the key separately. Since the mrjob cannot produce more than one file once, we must use two similar mappers to handle it and create datasets for sub-training and sub-testing.

---

```

1 for line in sys.stdin:
2     line = re.sub("[\\\"'", "", line)
3     #find('Train') for sub_train dataset
4     if line.find('Test') == 0:
5         line = line.strip()
6         value = line.split("\t")

```

---

```

7         if len(value)!=0:value =
8             ((value[1].lstrip("[")).rstrip("]"))
9             .replace(" ", "")
10            print(value)

```

#### 4.5.3 Building Decision Tree.

This step will build the decision tree with sub-training dataset and using the function implemented from scratch. Since the decision tree cannot be made for every line from the input file, we saved the line's information to a local variable after reading all lines and then used it to train and build the decision tree.

```

1 import ...
2 import zipimport
3
4 importer = zipimport.zipimporter('DT.mod')
5 dt = importer.load_module('DecisionTree')
6 column_name = ['customer_id',..., 'target']
7 # An empty tuple
8 csv_context = ()
9 for line in sys.stdin:
10     line = line.strip()
11     row = line.split(',')
12     if 'customer_id' not in str(row):
13         added_value_tuple = (row,)
14         csv_context = csv_context + added_value_tuple
15
16 df = pd.DataFrame(csv_context, columns=column_name)
17 # change data types
18 df = df.astype({'sunday': np.float64 ...})
19 tree = dt.decision_tree_algorithm(df)
20 print('%s\t%s' % ('tree', tree))

```

#### 4.5.4 Calculating accuracy.

This step will test the accuracy of the predictions made by the decision tree. However, before testing the predictions, the output of the building decision tree must merge with the sub-testing dataset. The mrjob does not support taking two input files at once. Once merged, the decision tree can be identified by its key when reading the input file. And then can use the function implemented in the scratch to get accuracy on the tree.

```

1 importer = zipimport.zipimporter('DT.mod')
2 dt = importer.load_module('DecisionTree')
3 column_name = ['customer_id',..., 'target']
4 trees = []
5 csv_context = ()
6 for line in sys.stdin:
7     if line.find('tree') == 0:
8         line = line.strip()
9         line = line.split('\t')
10        tree = line[1]
11        tree = ast.literal_eval(tree)
12        trees.append(tree)
13    else:
14        line = line.strip()
15        row = line.split(',')

```

```

16         if 'customer_id' not in str(row):
17             added_value_tuple = (row,)
18             csv_context = csv_context + added_value_tuple
19
20 df = pd.DataFrame(csv_context, columns=column_name)
21 df = df.astype({'sunday': np.float64, ...})
22 for i in range(len(trees)):
23     accuracy = dt.calculate_accuracy(df, trees[i])
24     print('%s\t%s' % ('accuracy', accuracy))

```

#### 4.5.5 Making predictions.

If the accuracy of the decision tree is satisfied, then we also need to combine the decision tree output with the testing dataset like in the last step. Finally, call the function implemented in the scratch to make predictions.

```

1 importer = zipimport.zipimporter('DT.mod')
2 dt = importer.load_module('DecisionTree')
3 column_name = ['customer_id',..., 'location_number_obj']
4 trees = []
5 csv_context = ()
6 for line in sys.stdin:
7     if line.find('tree') == 0:
8         line = line.strip()
9         line = line.split('\t')
10        tree = line[1]
11        tree = ast.literal_eval(tree)
12        trees.append(tree)
13    elif len(trees)!=0:
14        line = line.strip()
15        row = line.split(',')
16        if 'customer_id' not in str(row):
17            csv_context = (row,)
18            df = pd.DataFrame(csv_context,
19                             columns=column_name)
20            df = df.astype({'sunday': np.float64, ...})
21            x = dt.make_prediction(df.iloc[0], trees[0])
22            row.append(x)
23            print(row)

```

## 4.6 Applying Decision Tree algorithm using MapReduce framework and Sklearn package

In this session, we build the decision tree using the Sklearn package with mrjob, and it uses the numerical data as input. Unlike the self-implemented decision tree, the decision tree from Sklearn does not require that many preprocessing steps.

#### 4.6.1 Data sets preparation.

In the data preparation, we read every input line and give it a unique key using a random number with the current timestamp as a random seed. And then add an identified tag with the key. Then pass every row to the reducer and print the result as the new dataset.

```

1 class sklearn_train_preparation_map(MRJob):
2     def mapper(self, _, line):

```

```

3      random.seed(datetime.now())
4      key = random.random()
5      rows = csv.reader([line])
6      for row in rows:
7          # + 'test' for testing side
8          yield (key+' train',row)
9
10     def reducer(self, key, values):
11         yield (key, list(values))

```

After producing both new training and new testing datasets, merge them into one dataset.

#### 4.6.2 Making prediction.

When reading data from the input file, the data is split into two tuples based on the different key values—one for training and another for testing. Finally, applying the function from the Sklearn package to build and train the decision tree and use the tree to predict.

```

1 training_set = ()
2 testing_set = ()
3 column_name = ['customer_id',..., 'target']
4
5 for line in sys.stdin:
6     line = line.strip()
7     values = line.split('\t')
8     key = values[0].split(' ')
9     if key == 'train':
10         value = values[1].split(',')
11         training_set = training_set + (row,)
12     elif key == 'test':
13         value = values[1].split(',')
14         testing_set = testing_set + (row,)
15
16 df_train = pd.DataFrame(training_set, columns=column_name)
17 df_test = pd.DataFrame(testing_set, columns=column_name)
18 X = df_train.drop("target", axis=1)
19 y = df_train.target
20 X_train, X_test, y_train, y_test =
21 train_test_split(X, y, test_size=0.3)
22 clf = DecisionTreeClassifier()
23 clf = clf.fit(X_train, y_train)
24
25 for i in range(len(df_test)):
26     test_pred = clf.predict(df_test[i])
27     df_test[i].append(test_pred)
28     print(df_test[i])

```

We were limited on time and opted to focus on implementation in Spark, hence the mrjob for Sklearn decision tree has not yet been properly debugged.

## 4.7 Applying Decision Tree algorithm using PySpark and Spark MLlib

PySpark is the Python API for Apache Spark and MLlib is a wrapper over the PySpark and Spark's machine learning (ML) library. In our

project we used the DecisionTreeClassifier from the PySpark ml classification package.

#### 4.7.1 Loading the data.

First, we initialize the spark session and load the preprocessed data from HDFS. We are now using the preprocessed numerical only data to make the Decision Tree.

```

1 flieList = {
2     0: "hdfs://namenode:9000/test_categoricalToNumerical.csv",
3     1: "hdfs://namenode:9000/train_categoricalToNumerical.csv",
4 }
5
6 spark = \
7 SparkSession.builder.appName(
8     "Spark Decision Tree").getOrCreate()
9 sparkDFTrain = spark.read.csv(
10     flieList[1],
11     header=True,
12     inferSchema=True,
13 )
14 sparkDFTest = \
15 spark.read.csv(
16     flieList[0],
17     header=True,
18     inferSchema=True,
19 )

```

#### 4.7.2 Making Decision Tree.

Before making the Decision Tree, we need to select some of the features and a label column. From our dataset, the label column is the target and except for the target column, we use all the other columns as our features. We used the MLlib VectorAssembler function and VectorAssembler is a feature transformer that merges multiple columns into a vector column. It combines all the feature columns into a vector. After that we selected the features and labels as a model data frame.

```

1 featureColumns = (
2     "customer_id",
3     "gender",
4     "status_x",
5     ...
6 )
7 # converting the feature list into a vector
8 assembler = VectorAssembler(
9     inputCols=featureColumns, outputCol="features")
10 outputTrain = assembler.transform(sparkDFTrain)
11 # only features and label for making decision tree
12 modelDFTrain = outputTrain.select("features", "target")

```

With a max depth of 4, we used two types of impurity methods to make the comparison.

```

1 dtc = DecisionTreeClassifier(
2     labelCol="target",
3     featuresCol="features",

```



```

4      maxBins=5000,
5      maxDepth=4,
6      # impurity="gini"
7      impurity="entropy",)

```

#### 4.7.3 Make prediction and calculate accuracy.

After making the decision tree object, we did the same task for our testing set and separated the features and label. First, we fit the training dataset into the Decision Tree Classifier (dtc) model and then predict by transforming the model's test set. MulticlassClassificationEvaluator function gives us the final accuracy result which we will discuss in the next section.

```

1  dtcModel = dtc.fit(modelDFTrain)
2  modelDFTest = outputTest.select("features", "target")
3  dtcPred = dtcModel.transform(modelDFTest)
4  evaluator = MulticlassClassificationEvaluator(
5      labelCol="target",
6      predictionCol="prediction",
7      metricName="accuracy")
8  dtcAcc = evaluator.evaluate(dtcPred)
9  print("accuracy", dtcAcc)

```

## 5 EXPERIMENTAL EVALUATION

In this chapter, we will analyze the outcomes of the algorithms constructed and applied to our dataset. We will talk about the probable weaknesses of our implementations and some approaches to make them better moving forward. The algorithms running time will also be described, and their performance will be evaluated as a function of changes in the hardware configurations and algorithm parameters. We will also analyze the outcomes of the decision tree algorithms with MrJob and Spark on the dataset. We will also talk about the potential flaws in our implementations and how we can fix them.

### 5.1 Cluster Setup

Our cluster configuration consists of a network of four virtual machines (VMs), each of which has 4GB of RAM and two core CPUs. Four virtual machines (VMs) have been setup as a single cluster consisting of one master node and three subordinate nodes. The Hadoop 3.2.1 distribution has been deployed on the cluster, and Yarn is being used as the resource manager. The version of Spark that has been installed is 3.1.3. Yarn is also used as a cluster manager for Spark tasks, which we find to be quite useful. All of the testing and debugging were made in namenode and.

### 5.2 Performance evaluation of Spark

We observed that the Spark job runs faster than Hadoop during the experiment. We also measured the execution time and accuracy between two methods of Spark Decision Tree and Hadoop MapReduce. Our experiment shows that when we use the node impurity as Gini full data set gives lower accuracy than node impurity as Entropy.

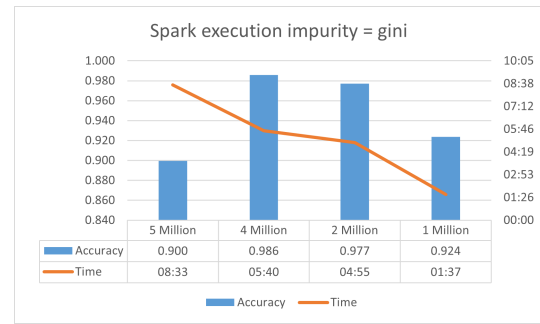


Figure 5: Spark execution comparison 1

Also we observed that time consumption is also slightly high in entropy. As we reduce the dataset sample size, accuracy also increases and execution time decreases. However, when we used the data chunk size of 2million, the result is pretty much the same in both Gini and Entropy impurity.

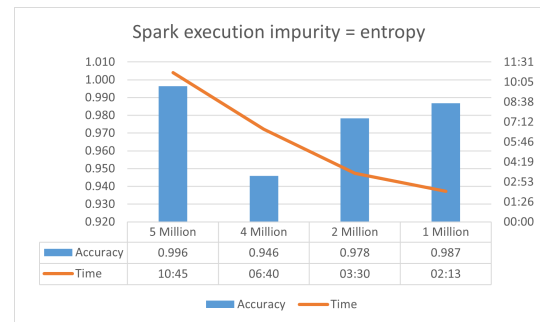


Figure 6: Spark execution comparison 2

### 5.3 Performance evaluation of Hadoop

Unlike Spark Hadoop execution time takes much time, our observation implies that the accuracy result remains the same even if we test with smaller chunks of data. On a full dataset, Hadoop execution time takes about 50min to complete and gives the output, and it is five times slower than Spark maximum execution time ( 11min) compared with Spark. Nevertheless, as we can see that accuracy remains constant, we can imply that Hadoop may be used where the data prediction needs to be more precise.

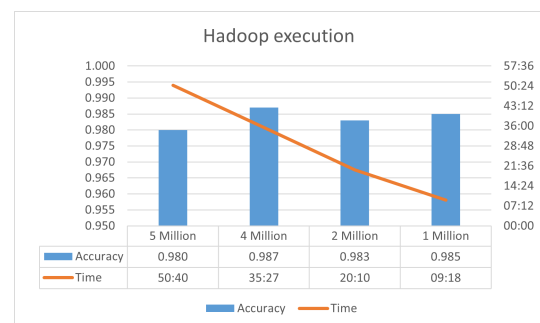


Figure 7: Hadoop execution analysis

## 5.4 Challenges

### 5.4.1 Converting the mrjob output.

The output of mrjob is plain text, so the first thing we need to solve is converting the output of mrjob into the correct format or the correct data architecture. We encounter this kind of problem at the very beginning. Then we develop formatting mrjob to overcome it. Even when we need to pass the plain text as a decision tree which builds from scratch, we find a way to read the text as the dictionary type, but when we try to pass the text as a decision tree object for the Sklearn package, we do not find any possible way to that. So, instead of passing the tree object from different mrjob, we decided to build the tree when all data have been read. Furthermore, we think that may cause a problem when using mrjob since the data may overflow the RAM, but that will not be a problem for Spark. So, we think the improvement can be made by shifting the mrjob to Spark in the future.

### 5.4.2 Importing third party packages.

We encounter a problem when importing the Numpy and Pandas package. Then we figure it out by using pip install relative packages. However, we also implement a decision tree from scratch, but we cannot solve the problem with previous experience. So finally, we find a way to zip the original class, rename the extension from .zip to .mod first, and use zipimporter to load the module.

### 5.4.3 Randomly split the dataset.

That was not a problem until we realized that once the decision tree is built with a lower accuracy rate, it will never be satisfied. Then we came up with randomly splitting the dataset and using different data to train the decision tree. So finally, this problem has been solved.

## 6 CONCLUSION

In this project, we used mrjob to create a Machine Learning algorithm, namely Decision Tree. We used an Apache Hadoop cluster as the working environment to apply the method to a collection of data with varying data sizes.

We compared the results provided on the same dataset using well-designed and efficient libraries such as Spark's mllib, Python's sklearn, and the self-built decision tree module to the results produced by the implemented method. Finally, we show how the prediction results compare using a bar chart. Our decision tree's predictions are not exceptionally efficient but comparable with other libraries.

For future work, we have opportunities to improve the performance of building the decision tree more efficiently. Furthermore, have the opportunities to find a better way to manage the stored data to avoid memory overflow.

## REFERENCES

- [1] Apache hadoop, <https://hadoop.apache.org/>
- [2] J. Venner, Pro Hadoop, Apress, June 2009.
- [3] Shvachko K, Kuang H, Radia S, et al. The hadoop distributed file system[C]//2010 IEEE 26th symposium on mass storage systems and technologies (MSST). Ieee, 2010: 1-10.
- [4] Scala programming language. <http://www.scala-lang.org>.
- [5] Zaharia M, Chowdhury M, Franklin M J, et al. Spark: Cluster computing with working sets[C]//2nd USENIX Workshop on Hot Topics in Cloud Computing (HotCloud 10). 2010.
- [6] Tomasz Wiktorski. 2019. Data-intensive Systems: Principles and Fundamentals using Hadoop and Spark. <https://doi.org/10.1007/978-3-030-04603-32>
- [7] [n. d.].Decision Trees in Machine Learning.<https://towardsdatascience.com/decision-trees-in-machine-learning-641b9c4e8052>
- [8] Decision Trees for Classification, <https://www.xoriant.com/blog/product-engineering/decision-trees-machine-learning-algorithm.html>
- [9] What is Machine Learning Library (MLlib) | Databricks,<https://databricks.com/glossary/what-is-machine-learning-library>
- [10] [n. d.].Spark MLlib for Scalable Machine Learning with Spark.[https://www.projectpro.io/article/spark-mlib-for-scalable-machine-learning-with-spark/339mccetoc\\_1fb3o9fdfs8](https://www.projectpro.io/article/spark-mlib-for-scalable-machine-learning-with-spark/339mccetoc_1fb3o9fdfs8)
- [11] Stuart J. Russell, Peter Norvig (2010) Artificial Intelligence: A Modern Approach, Third Edition, Prentice Hall ISBN 9780136042594.
- [12] Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar (2012) Foundations of Machine Learning, The MIT Press ISBN 9780262018258.
- [13] Oshiro T M, Perez P S, Baranauskas J A. How many trees in a random forest?[C]//International workshop on machine learning and data mining in pattern recognition. Springer, Berlin, Heidelberg, 2012: 154-168.
- [14] Ho, Tin Kam (1995). Random Decision Forests (PDF). Proceedings of the 3rd International Conference on Document Analysis and Recognition, Montreal, QC, 14-16 August 1995. pp. 278-282. Archived from the original (PDF) on 17 April 2016. Retrieved 5 June 2016.
- [15] Ho TK (1998). "The Random Subspace Method for Constructing Decision Forests" (PDF). IEEE Transactions on Pattern Analysis and Machine Intelligence. 20 (8): 832-844. doi:10.1109/34.709601.
- [16] Hastie, Trevor; Tibshirani, Robert; Friedman, Jerome (2008). The Elements of Statistical Learning (2nd ed.). Springer. ISBN 0-387-95284-5.
- [17] Piryonesi S. Madeh; El-Diraby Tamer E. (2020-06-01). "Role of Data Analytics in Infrastructure Asset Management: Overcoming Data Size and Quality Problems". Journal of Transportation Engineering, Part B: Pavements. 146 (2): 04020022. doi:10.1061/JPEODX.0000175. S2CID 216485629.
- [18] Piryonesi, S. Madeh; El-Diraby, Tamer E. (2021-02-01). "Using Machine Learning to Examine Impact of Type of Performance Indicator on Flexible Pavement Deterioration Modeling". Journal of Infrastructure Systems. 27 (2): 04021005. doi:10.1061/(ASCE)IS.1943-555X.0000602. ISSN 1076-0342. S2CID 233550030.
- [19] [n. d.].Decision tree introduction with example.<https://www.geeksforgeeks.org/decision-tree-introduction-example/>
- [20] [n. d.].Decision Tree Classifier in Python Sklearn with Example.<https://machinelearningknowledge.ai/decision-tree-classifier-in-python-sklearn-with-example/>
- [21] [n. d.].Difference Between Spark DataFrame and Pandas DataFrame.<https://www.geeksforgeeks.org/difference-between-spark-dataframe-and-pandas-dataframe/>
- [22] ANDRII.2022.Restaurant Recommendation Challenge.<https://www.kaggle.com/datasets/mrmorj/restaurant-recommendation-challenge>