

ANALYSIS OF WEATHER USING HADOOP

**Akshaya Sri J^{*1}, Mohindran S. R.^{*2}, Rohinth. S^{*3}, Sanchez Innocencia D^{*4},
Dr. M. Marimuthu^{*5}, Dr. P. Velvadivu^{*6}**

^{*1,2,3,4}III M.Sc (Integrated), Decision and Computing Sciences,
Coimbatore Institute of Technology, Coimbatore, India

^{*5,6}Assistant Professor, Department of Computing,
Coimbatore Institute of Technology, Coimbatore, India.

ABSTRACT

The Temperature is the one which is mainly considered for the changing weather conditions. This paper deals with the analysis of Temperature, where we use Big Data features like MapReduce to calculate the Average Temperature with respect to its year. We use HIVE, which is an open-source data warehousing software used for reading large datasets. Using HIVE, we write queries to fetch the data and we can also update it. Further, we use Random Forest Regression Technique, a Machine Learning Algorithm to classify and predict the temperature based on the classification. The dataset for this analysis is taken from the Data.gov website which contained data for about 85 years and we have taken the data for about 30 years from 1990 – 2020. The above analysis is done using Hadoop and Pyspark.

Keywords: MapReduce, HIVE, Random Forest Regression, Hadoop, PySpark

I. INTRODUCTION

Data has been growing since its existence in various domains like share market, social media etc. The data may be of any kind like structured, unstructured and semi-structured. BigData is one which provides a lot of techniques and tools for processing efficiently any kind of data and storing it in tools which helps to store large data. Big data can handle data sets with sizes beyond the ability of commonly used software tools to capture and process the data.

The change of Temperature plays a vital role to predict the Weather, where temperature of a day determines the Weather conditions. Forecasting of Weather is very much useful in many sectors through which they get affected like the Agricultural Sector, Water Resources, Air Traffic and also major causes in the field of Tourism. Forecasting of Temperature will be much useful for the people of these sectors so that they may make a calculation of the Weather Conditions and prepare accordingly. For our analysis, we first use HIVE queries for table creation, retrieval of data and then use the MapReduce programming framework to make the analysis on temperature and finally we build a Machine Learning Model using PYSPARK for prediction.

II. LITERATURE REVIEW

In the paper “Novel Weather Data Analysis Using Hadoop and MapReduce”, the author has considered two important Big Data techniques which are Hadoop and Mapreduce to analyze the large dataset. They even developed a system that adopts the historical weather data by using the Big Data Techniques mentioned.

Chouksey P in his paper “A Review of Weather Data Analytics using Big Data” has explained the weather parameters which are to be used for analysis. It suggested that weather data is mostly analyzed using MapReduce. So he stated that comparison study of weather data using MapReduce and Spark is very much needed.

Shiva Naga Jyothi CH along with his colleagues in his paper “Weather Data Analysis Using Hadoop MapReduce and Spark” presents the analysis of weather data by calculating maximum, minimum, average temperatures and top coldest and hottest stations in a particular location or year using various weather parameters.

III. PRELIMINARIES

A. MapReduce

MapReduce is a programming framework that allows us to perform distributed and parallel processing on large data sets in a distributed environment. Then, the reducer aggregates those intermediate data tuples into a smaller set of tuples or key-value pairs which is the final output. Here using Mapreduce, we find the average temperature

across different years. The Mapper takes in the Weather input file as an Input and emits the Year and Temperature. The Reducer emits the data for Average Temperature.

B. Apache Hive

Now-a-days, the use of Hive Metastore is one of the big constants in big data implementation. Hive ships with the metastore service. This service allows you to manage your metastore as any other database. You can host this service on any of the popular RDBMS (e.g. MySQL, PostgreSQL etc.). This allows you to query your metastore with simple SQL queries, along with provisions of backup and disaster recovery. We have installed hive with mysql as a metastore and loaded the data. Here, we have used HIVE queries to Load and Create tables, to find the Maximum and Minimum values with respect to its years and also to find the second maximum temperature with respect to its year.

IV. MACHINE LEARNING ALGORITHM

A. Data Preprocessing

Data preprocessing is a data mining technique that involves transforming raw data into an understandable format. Since, there are high chances of having missing values in the dataset, unwanted attributes etc. We have cleaned it up using preprocessing techniques. After taking the dataset, we have dropped two variables of our dataset which won't affect the further analysis. We have also checked if there are any missing values in the dataset.

B. PySpark

Today, PySpark is being used by the majority of Data Scientists because of the library set it contains. PySpark is one which links the Python API to the spark core. We are implementing the machine learning algorithm using pyspark. For our dataset, we have used Random forest Regression to classify the Data and then predict the temperature.

C. Random Forest Regression

The Random forest is also sometimes called a random decision forest. It is an ensemble learning technique used for solving supervised learning tasks such as classification and regression. An advantageous feature of RF is that it can overcome the overfitting problem across its training dataset. Also an additional advantage is that it does not need any preprocessing initially and accepts even the missing values and predicts the output. Here, using Random Forest, we have first classified and calculated the accuracy score.

V. RESULTS AND DISCUSSION

A. MapReduce Program

Using MapReduce, we have taken the full dataset and then summed up the dataset year wise for the Temperature, Dew Point and the wind Speed taking into consideration the Month. And further we have calculated the Average of Temperature, Dew Point and the Wind Speed with respect to its years from 1990 - 2020 along with its station ID.

The Output is shown below.

```

WARNING: Use of this script to execute hdfs command is deprecated.
          Please use the hdfs command for it.

WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication.util.KerberosUtil (file:/usr/local/hadoop/share/hadoop/common/10/hadoop-common-2.8.0.jar) to method sun.security.auth.config.getScheme()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
11/01/21 12:12:38 WARN util.NativeCodeLoader: Unable to load native hadoop library for your platform... using builtin-java classes where applicable
cat000000000_1990 12.155083 12.438499 26.722337
cat000000000_1991 12.155083 12.438499 26.722337
cat000000000_1992 12.155083 12.438499 26.722337
cat000000000_1993 12.155083 12.438499 26.722337
cat000000000_1994 12.155083 12.438499 26.722337
cat000000000_1995 12.155083 12.438499 26.722337
cat000000000_1996 12.155083 12.438499 26.722337
cat000000000_1997 12.155083 12.438499 26.722337
cat000000000_1998 12.155083 12.438499 26.722337
cat000000000_1999 12.155083 12.438499 26.722337
cat000000000_2000 12.155083 12.438499 26.722337
cat000000000_2001 12.155083 12.438499 26.722337
cat000000000_2002 12.155083 12.438499 26.722337
cat000000000_2003 12.155083 12.438499 26.722337
cat000000000_2004 12.155083 12.438499 26.722337
cat000000000_2005 12.155083 12.438499 26.722337
cat000000000_2006 12.155083 12.438499 26.722337
cat000000000_2007 12.155083 12.438499 26.722337
cat000000000_2008 12.155083 12.438499 26.722337
cat000000000_2009 12.155083 12.438499 26.722337
cat000000000_2010 12.155083 12.438499 26.722337
cat000000000_2011 12.155083 12.438499 26.722337
cat000000000_2012 12.155083 12.438499 26.722337
cat000000000_2013 12.155083 12.438499 26.722337
cat000000000_2014 12.155083 12.438499 26.722337
cat000000000_2015 12.155083 12.438499 26.722337
cat000000000_2016 12.155083 12.438499 26.722337
cat000000000_2017 12.155083 12.438499 26.722337
cat000000000_2018 12.155083 12.438499 26.722337
cat000000000_2019 12.155083 12.438499 26.722337
cat000000000_2020 12.155083 12.438499 26.722337

```

B. Apache Hive

Using HIVE, we have loaded and stored the data in the database and then created the table using the HIVE query and also created 3 more queries in HIVE for retrieving the data we want. Through which we find the Year, where we have the Maximum and Minimum Temperature.

This is the output of the Maximum Temperature with respect to its years using HIVE SQL.

```
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes-per-reducer=number
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=number
In order to set a constant number of reducers:
  set mapreduce.job.reducers=number
Starting Job = job_161907804554_0002, Tracking URL = http://ml18088.proxy/application_161907804554_0002/
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_161907804554_0002
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2021-04-22 13:49:19.074 Stage-1 map = 0%, reduce = 0%
2021-04-22 13:49:20.293 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 3.89 sec
2021-04-22 13:49:43.001 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 6.46 sec
MapReduce Total cumulative CPU time: 6 seconds 460 msec
Ended Job = job_161907804554_0002
MapReduce Jobs Launched:
  Stage: Stage-1: Map: 1; Reduce: 1; Cumulative CPU: 6.46 sec; HDFS Read: 1221923 HDFS Write: 1110 SUCCESS
Total MapReduce CPU Time Spent: 6 seconds 460 msec
OK
1001099999 1990 48.5
1001099999 1991 47.2
1001099999 1992 44.8
1001099999 1993 46.2
1001099999 1994 46.6
1001099999 1995 45.5
1001099999 1996 46.6
1001099999 1997 46.2
1001099999 1998 47.7
1001099999 1999 46.8
1001099999 2000 47.2
1001099999 2001 47.2
1001099999 2002 56.6
1001099999 2003 49.0
1001099999 2004 46.9
1001099999 2005 51.4
1001099999 2006 47.5
1001099999 2007 46.5
1001099999 2008 46.5
1001099999 2009 52.0
1001099999 2010 47.3
1001099999 2011 48.4
1001099999 2012 47.8
1001099999 2013 47.4
1001099999 2014 49.1
1001099999 2015 47.6
1001099999 2016 51.0
1001099999 2017 51.0
1001099999 2018 51.7
1001099999 2019 53.7
1001099999 2020 46.7
Time taken: 39.97 seconds, fetched: 31 row(s)
```

This is the output of the Minimum Temperature with respect to its years using HIVE SQL

```
Starting Job = job_161907804554_0003, Tracking URL = http://ml18088.proxy/application_161907804554_0003/
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_161907804554_0003
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2021-04-22 13:50:14.167 Stage-1 map = 0%, reduce = 0%
2021-04-22 13:50:23.378 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 3.19 sec
2021-04-22 13:50:29.582 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 5.25 sec
MapReduce Total cumulative CPU time: 5 seconds 258 msec
Ended Job = job_161907804554_0003
MapReduce Jobs Launched:
  Stage: Stage-1: Map: 1; Reduce: 1; Cumulative CPU: 5.25 sec; HDFS Read: 1221923 HDFS Write: 1096 SUCCESS
Total MapReduce CPU Time Spent: 5 seconds 258 msec
OK
1001099999 1990 9.4
1001099999 1991 7.5
1001099999 1992 3.4
1001099999 1993 4.3
1001099999 1994 11.2
1001099999 1995 18.5
1001099999 1996 2.9
1001099999 1997 -6.2
1001099999 1998 -2.8
1001099999 1999 12.6
1001099999 2000 4.8
1001099999 2001 -8.8
1001099999 2002 7.2
1001099999 2003 1.6
1001099999 2004 6.6
1001099999 2005 12.1
1001099999 2006 12.1
1001099999 2007 9.2
1001099999 2008 12.8
1001099999 2009 9.9
1001099999 2010 14.2
1001099999 2011 11.4
1001099999 2012 8.1
1001099999 2013 16.7
1001099999 2014 12.3
1001099999 2015 8.8
1001099999 2016 16.9
1001099999 2017 11.9
1001099999 2018 11.4
1001099999 2019 11.9
1001099999 2020 8.9
Time taken: 27.282 seconds, fetched: 31 row(s)
```

This is the output of the 2nd Maximum Temperature with respect to its years using HIVE SQL

```
MapredLocal task succeeded
Launching Job 2 out of 2
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes-per-reducer=number
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=number
In order to set a constant number of reducers:
  set mapreduce.job.reducers=number
Starting Job = job_161911714887_0049, Tracking URL = http://ml18088.proxy/application_161911714887_0049/
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_161911714887_0049
Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
2021-04-23 00:42:39.024 Stage-2 map = 0%, reduce = 0%
2021-04-23 00:42:44.337 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 2.26 sec
2021-04-23 00:42:58.452 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 4.0 sec
MapReduce Total cumulative CPU time: 4 seconds 0 msec
Ended Job = job_161911714887_0049
MapReduce Jobs Launched:
  Stage: Stage-2: Map: 1; Reduce: 1; Cumulative CPU: 4.02 sec; HDFS Read: 1221885 HDFS Write: 991 SUCCESS
  Stage: Stage-2: Map: 1; Reduce: 1; Cumulative CPU: 4.0 sec; HDFS Read: 11949 HDFS Write: 880 SUCCESS
Total MapReduce CPU Time Spent: 4 seconds 28 msec
OK
1990 47.78
1991 47.18
1992 46.58
1993 45.38
1994 45.98
1995 43.48
1996 45.58
1997 45.58
1998 47.18
1999 45.38
2000 47.08
2001 46.88
2002 49.28
2003 48.28
2004 48.98
2005 49.08
2006 47.08
2007 46.28
2008 46.48
2009 48.18
2010 46.78
2011 47.98
2012 47.48
2013 46.48
2014 46.48
2015 47.58
2016 49.38
2017 46.48
2018 49.38
2019 51.98
2020 46.08
Time taken: 46.653 seconds, fetched: 31 row(s)
```

C. Machine Learning Algorithm

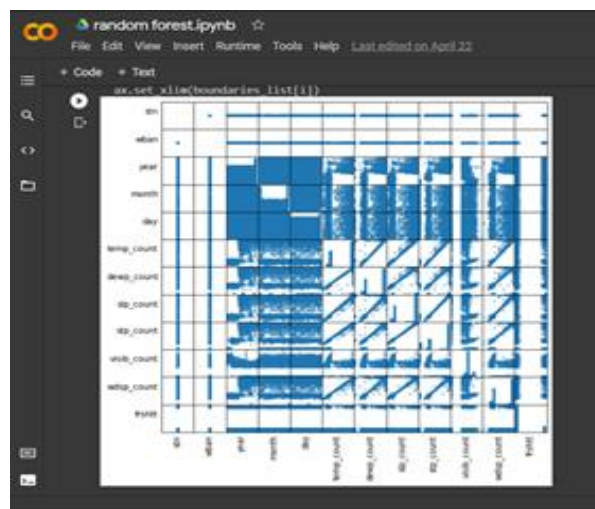
We have used the Random Forest Machine Learning Algorithm using the PySpark Environment by storing the data in PySpark. Here, we have found that the Accuracy of the Temperature is 0.77 and hence, we say that this is the best fit Model for the Analysis. The data has been read through PySpark and used for further analysis of Model fitting and prediction.

```
random forest.ipynb
File Edit View Insert Runtime Tools Help Last edited on April 22

+ Code + Test

from pyspark.sql import SparkSession
spark = SparkSession.builder.appName('random_forest').getOrCreate()
df = spark.read.csv('weather-mod.csv', header = True, inferSchema = True)
df.printSchema()

root
 |-- station: integer (nullable = true)
 |-- season: integer (nullable = true)
 |-- year: integer (nullable = true)
 |-- month: integer (nullable = true)
 |-- day: integer (nullable = true)
 |-- temp: double (nullable = true)
 |-- temp_count: integer (nullable = true)
 |-- dewp: double (nullable = true)
 |-- dewp_count: integer (nullable = true)
 |-- slp: double (nullable = true)
 |-- slp_count: integer (nullable = true)
 |-- stp: double (nullable = true)
 |-- stp_count: integer (nullable = true)
 |-- visib: double (nullable = true)
 |-- visib_count: integer (nullable = true)
 |-- whsp: double (nullable = true)
 |-- whsp_count: integer (nullable = true)
 |-- mxspd: double (nullable = true)
 |-- gust: double (nullable = true)
 |-- maxtemp: double (nullable = true)
 |-- mintemp: double (nullable = true)
 |-- prcp: double (nullable = true)
 |-- cndp: double (nullable = true)
 |-- frsttt: integer (nullable = true)
```



```
random forest.ipynb
File Edit View Insert Runtime Tools Help Last edited on April 22

+ Code + Test

av.set_will(boundsaries_list[1])

||
24.0
15.0
18.0
8.0
20.0
11.0
only showing top 50 rows

|| from pyspark.ml.feature import StringIndexer, IndexToString
labelInverse = IndexToString().setInputCol("prediction")
labelInverse.transform(rf_predictions).show()

from pyspark.ml.evaluation import MulticlassClassificationEvaluator

multi_evaluator = MulticlassClassificationEvaluator(labelCol = 'temp_index', metricName = 'accuracy')
print('Random Forest Classifier Accuracy:', multi_evaluator.evaluate(rf_predictions))

Random Forest Classifier Accuracy: 0.76719396667362
```

VI. CONCLUSION

Large amount of Weather Data has been collected every day which is being generated by the satellites, storing and processing of these large dataset seems to be very challenging. Hence, for storing these large data, we used an HIVE database and it's easy for retrieving the Maximum and Minimum Temperature and further more processes can also be done. For the Analysis, we have taken the main attributes as the Temperature, Dew Point and the Wind Speed, which is an essential attribute for Weather Forecasting. We used MapReduce and found the Average Temperature with respect to its Year. From the Analysis we have made, we conclude that for the Analysis of Weather, Temperature is one which is very important and also the Wind Speed and Dew Point are an important factor for the prediction of weather. And also we found that Random Forest Regression is one which is an appropriate Machine Learning Model to predict the Temperature.

VII. REFERENCES

- [1] <https://www.smartdatacollective.com/four-considerations-when-choosing-hadoop-distribution/>
- [2] <https://stackabuse.com/random-forest-algorithm-with-python-and-scikit-learn/>
- [3] <https://www.xajzkidx.cn/gallery/311-april2020.pdf>
- [4] <https://ieeexplore.ieee.org/document/8728444>
- [5] <http://www.bookmetrix.com/detail/chapter/e5f8dfce-be6e-4d9b-a642-a2dc16d9bf5b>
- [6] ["Novel Weather Data Analysis Using Hadoop and MapReduce" – A Case Study, 2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS).
- [7] "Big Data Weather Analytics Using Hadoop" an International Journal of Emerging Technology in Computer Science & Electronics (IJETCSE) ISSN: 0976-1353 Volume14 Issue 2.
- [8] Chouksey P., Chauhan A., "A Review of Weather Data Analytics using Big Data", IJARCCCE, ISSN: 2278- 1021 Volume-06, Issue01, Page No (365-368), January, 2017.