

## ANALYZING IMPACTS OF WEATHER DATA USING APACHE HADOOP FRAMEWORK

**Dr S. Anitha** Assistant professor, PG Department of Computer Science, Dwaraka Doss Goverdhan Doss Vaishnav College, Arumbakkam, India. mail: [anitasenthil@gmail.com](mailto:anitasenthil@gmail.com)

**M. Kalaivani** Assistant professor, PG Department of Computer Science, Dwaraka Doss Goverdhan Doss Vaishnav College, Arumbakkam, India. mail: [kalaivani@dgvaishnavcollege.edu.in](mailto:kalaivani@dgvaishnavcollege.edu.in)

### **Abstract:**

In data science era, the scale of weather data is enormous and rising rapidly. Apache Hadoop is a fast and efficient framework which has been used in many applications in big data field. However, for the large-scale weather dataset, the traditional algorithms are not capable enough to satisfy the genuine application requirements efficiently. Hadoop is a framework which deals with Big and Huge variety of datasets which supports processing components that collectively called Hadoop Ecosystem. This paper proposes efficient weather data analyses are carried out by Apache MapReduce and Apache Pig in Hadoop framework. Weather datasets are taken from NCDC Database for this proposed research. The impacts of Weather analysis are obtained from both Mapreduce and Apache pig and they were compared.

### **Keywords:**

Big data Analytics, Hadoop Ecosystem, Apache MapReduce, NCDC Datasets

## **1. Introduction**

In the atmospheric sciences, weather data is really rich and valued, which requires a mass of scientific computing, and provides services to the communities. Climate data are dramatically increasing in volume and complexity, since users of these data in the scientific community and the public are rapidly increasing [1]. The paper deals with the general architecture of HADOOP including the details of its numerous components. For analyzing the large number of data, traditional data analysis techniques have failed to carry out analysis on larger data sets effectively. Newly, the dominant platform that has proved in processing hefty sets of data is Hadoop, which is considered to be operative for distributed file processing and distributed storage of wider range of data. The main component of Hadoop is HDFS and Mapreduce. MapReduce is a programming model for computing bigger data sets and HDFS is a Hadoop Distributed File System that stores data in the type of memory blocks and distributes the data across cluster of nodes. In this paper, Apache Mapreduce is used to analyze the NCDC weather dataset.

This paper is organized as follows: Section 2 includes the review of literature and an outline of the paper. Section 3 undertakes the research methodology and Section 4 discusses the Results and Discussion of this study. Section 5 concludes the proposed method and scope for future work.

## **2. Review of Literature**

MapReduce is a key technology of using cloud computing to process a big amount of data. It is a parallel programming model and an associated implementation for processing and generating huge datasets in a broad variety of real world tasks proposed by Google. It is not only a programming model, but also a task scheduling model. It is composed of two essential functions: Map and Reduce, defined by users. A Map function is used to handle every Input and convert it as an intermediate key/value pair. A Reduce function is specified to combine all of the intermediate value with the same middle key [2]. Map Reduce is an open-source framework for processing of a huge quantity data across the collection of processors using the high-level languages. These components afford easy way of using languages, graphical interfaces and also for handling data on thousands of workstations. Hadoop cluster is a set of machines together networked in a place. Data processing and all occur within the cloud of machines. User can process jobs to Hadoop from the desktop processor in remote location from its cluster [Google MapReduce is typically utilized to perform distributed computing on clusters

of computers. Thereby, the effect originally achieved only by expensive high-performance computer can be achieved by low-cost computing services.

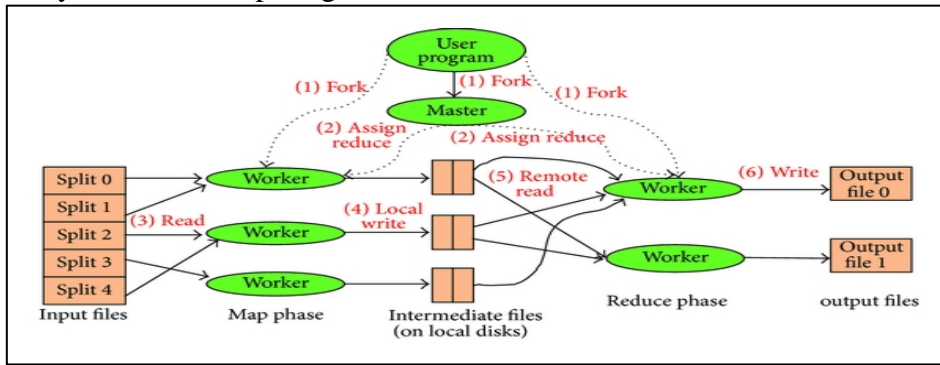


Fig. 1: Mapreduce paradigm

## 2. Research Methodology

In this research, Analyzing weather data of Fairbanks, Alaska is utilized to find cold and hot days using MapReduce and Hadoop. The National Climatic Data Center (NCDC) is the world's largest active archive of weather data. In this paper, this research work consists of the following five phases,

1. Pre-processing the NCDC datasets
2. Feature selection from the datasets
3. Loading the preprocessed dataset to HDFS
4. Analysis of NCDC weather data set with Mapreduce and Pig
5. Comparing and Evaluating the results.



Fig-2: Mapreduce Function

**Step1:** for Compiled the Java File: `javac -classpath /home/student3/hadoop-common-2.6.1.jar:/home/student3/hadoop-mapreduce-client-core-2.6.1.jar:/home/cloudera/commons-cli-2.0.jar -d . MaxTemperature.java MaxTemperatureMapper.java MaxTemperatureReducer.java`

**Step 2:** Created the JAR file: `jar -cvf hadoop-project.jar *.class`

**Step 3:** Executed the jar file: `hadoop jar hadoop-project.jar MaxTemperature /home/student3/Project/ /home/student3/Project_output111`

**Step 4:** Copy the output file to local hdfs `dfs -copyToLocal /home/student3/Project_output111/part-r-00000`

## 3. Description of Dataset

In this paper, Apache Map-Reduce weather analysis algorithm is applied in NCDC dataset for analyzing weather datasets to predict the weather condition for a particular year. The NCDC weather dataset is downloaded for year 1930 and loaded it in HDFS system. MapReduce and Pig algorithm is implemented in dataset to find the Min, Max, avg temperature for different stations. Maximum and Minimum temperature are retrieved and are used to find the cold and hot data respectively. NCDC (National Climatic Data Center) is collected across more than 116 weather stations and more than 1000 observations centers. The data is unstructured, which becomes a challenging task to analyze it. Weather sensors are gathering weather statistics throughout the world in a huge volume of log data. NCDC weather data is unstructured and record-oriented. In this dataset, each row has lots of fields like longitude, latitude, daily max-min temperature, daily average temperature, etc. temperature is taken as the main element.

## 4. Results and Discussion

The data were cleaned, preprocessed, and then fed into Mapreduce algorithm and Apache Pig. Hadoop is installed in pseudo distributed mode. The performance evaluation between the pig and Mapreduce is depicted in Fig-3 and below are the commands for the performance of Hadoop. And

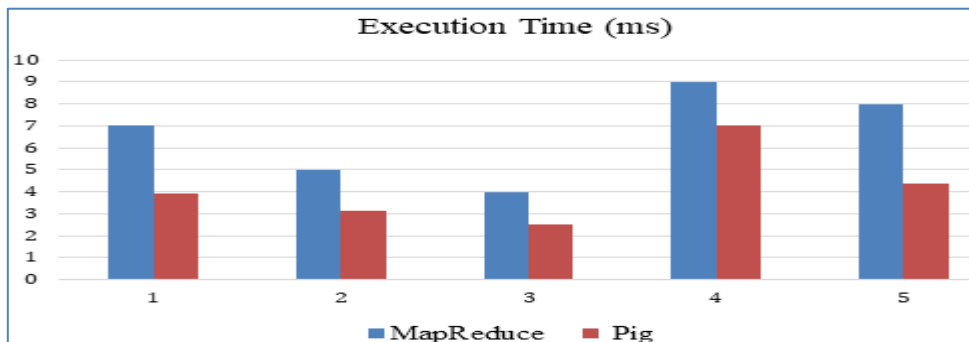
according to the data analyzing speed and efficiency Apache Pig proves to be better than Mapreduce. Comparison of state of art Literature survey on various research works are listed in Table-3.

**Table-2 output of weather data analysis.**

```

In hdfs environment: Create the temporary content file in the input directory:
[cloudera@quickstart ~]$ hdfs dfs -mkdir weather_dir
Put the file.txt into hdfs:
[cloudera@quickstart ~]$ hdfs dfs -put /home/cloudera/Desktop/wd.txt weather_dir/
[cloudera@quickstart ~]$ hdfs dfs -ls weather_dir/
Found 1 items
-rw-r--r-- 1 cloudera cloudera 41881 2019-10-09 22:16 weather_dir/wd.txt
To see the content of the file:
[cloudera@quickstart ~]$ hdfs dfs -cat weather_dir/wd.txt

23907 20150101 2.423 -98.08 30.62 2.2 -0.6 0.8 0.9 6.2 1.47 C 3.7 1.1 2.5
99.9 85.4 97.2 0.369 0.308 -99.000 -99.000 -99.000 7.0 8.1 -9999.0 -
9999.0 - 9999.0 23907 20150102 2.423 -98.08 30.62 3.5 1.3 2.4 2.2 9.0
1.43 C 4.9 2.3 3.1 100.0 98.8 99.8 0.391 0.327 -99.000 -99.000 -99.000
7.1 7.9 -9999.0 -9999.0 - 9999.0 23907 20150103 2.423 -98.08 30.62 15.9
2.3 9.1 7.5 2.9 11.00 C 16.4 2.9 7.3 100.0 34.8 73.7 0.450 0.397 -99.000 -
99.000 -99.000 7.6 7.9 -9999.0 -9999.0 - 9999.0 23907 20150104 2.423 -
98.08 30.62 9.2 -1.3 3.9 4.2 0.0 13.24 C 12.4 -0.5 4.9 82.0 40.6 61.7 0.414
0.352 -99.000 -99.000 -99.000 7.3
To see the output:
[cloudera@quickstart ~]$ hdfs dfs -ls out/
Found 2 items
-rw-r--r-- 1 cloudera cloudera 0 2019-10-09 22:20 out/_SUCCESS
-rw-r--r-- 1 cloudera cloudera 4632 2019-10-09 22:20 out/part-r-00000
[cloudera@quickstart ~]$ hdfs dfs -cat out/part-r-00000
1 The Day is Cold Day :20200101 -21.8
2 The Day is Cold Day :20200102 -23.4
3 The Day is Cold Day :20200103 -25.4
4 The Day is Cold Day :20200104 -26.8
5 The Day is Cold Day :20200105 -28.8
6 The Day is Cold Day :20200106 -30.0
7 The Day is Cold Day :20200107 -31.4
8 The Day is Cold Day :20200108 -33.6
9 The Day is Cold Day :20200109 -26.6
10 The Day is Cold Day :20200110 -24.3
    
```



**Fig. 3: Comparison of runtime –Mapreduce and Apache Pig**

**Table 3: Comparison of State of Art Literature Survey**

Sl.no.	Literature survey	Tools used on Weather Dataset
1	Dhyani et al,2014	Hadoop, Spark
2	Riyaz et al,2015	Hadoop, MapReduce
3	Dagade et al ,2015	Hadoop, Apache Spark, Hive
4	Navadia et al,2017	Hadoop

5	Chouksey et al,2017	Hadoop, Hive and Pyspark.
6	Suryanarayana, et al,2019	Hadoop, MapReduce
7	Priyanka Dinesh et al,2020	Hadoop, MapReduce, Spark.
8	Gupta et al,2021	Hadoop, MapReduce, Flask
9	Michael Raj et al,2022	R Package
10	Parameshachari, B. D et al,2022	Hadoop, MapReduce,

## 5. Conclusion

Weather data analysis algorithm is applied on the NCDC datasets. The analysis shows Cold and hot days along with the temperature. Mapper and reducer classes of MapReduce are used for analyzing the dataset and the results were compared in terms of execution time. Pig is outperformed and gained less elapsed time than MapReduce. Numerous research has been carried out on weather analysis especially for its temperature. There is to examine all significant weather parameters like temperature, pressure and humidity. Also, the technology highlighting evaluation of Hadoop, MapReduce and Apache pig are more important to study which is better suited for weather data analysis. From this study, Apache Pig is the best tools for analyzing large number of instances in short execution time. Considering the merits and demerits of the proposed system for predicting impacts of weather dataset by utilizing the tools of Big Data Analytics.

## Reference

- [1] Dhyani, Bijesh, and Anurag Barthwal. "Big data analytics using Hadoop." *International Journal of Computer Applications* 108, no. 12 (2014): 0975-8887. "Novel Weather Data Analysis Using Hadoop and MapReduce" – A Case Study, 2019
- [2] Anitha, S., and Mary Metilda. "Apache Hadoop based effective sentiment analysis on demonetization and covid-19 tweets." *Global transitions proceedings* 3, no. 1 (2022): 338-342.
- [3] Suryanarayana, V., B. S. Sathish, A. Ranganayakulu, and P. Ganesan. "Novel weather data analysis using Hadoop and MapReduce—a case study." In *2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS)*, pp. 204-207. IEEE, 2019.
- [4] Chouksey, Priyanka, and Abhishek Singh Chauhan. "A review of weather data analytics using big data." *International Journal of Advanced Research in Computer and Communication Engineering* 6, no. 1 (2017): 365-368.
- [5] More, Priyanka Dinesh, Sunita Nandgave, and Megha Kadam. "Weather data analytics using hadoop with map-reduce." In *ICCCE 2019: Proceedings of the 2nd International Conference on Communications and Cyber Physical Engineering*, pp. 189-196. Springer Singapore, 2020.
- [6] Riyaz, P. A., and Surekha Mariam Varghese. "Leveraging map reduce with hadoop for weather data analytics." *OSR Journal of Computer Engineering (IOSR)* (2015).
- [7] Tf, Michael Raj, and Yaduvir Singh. "An Exploration on Big Data Analysis and Data Mining Methods." In *2022 International Conference on Futuristic Technologies (INCOFT)*, pp. 1-6. IEEE, 2022.
- [8] Kaul, Sameer. "Review of existing data mining techniques used for weather prediction." *Int J Res Appl Sci Eng Technol*, ISSN 5 (2017).
- [9] Dagade, Veershetty, Mahesh Lagali, Supriya Avadhani, and Priya Kalekar. "Big data weather analytics using hadoop." *International Journal of Emerging Technology in Computer Science & Electronics* 14, no. 2 (2015): 847-851.
- [10] Gupta, Rishi, Akhilesh Kumar Sharma, Oorja Garg, Krishna Modi, Shahreen Kasim, Zirawani Baharum, Hairulnizam Mahdin, and Salama A. Mostafa. "WB-CPI: Weather based crop prediction in India using big data analytics." *IEEE access* 9 (2021): 137869-137885.
- [11] Jain, Himanshi, and Raksha Jain. "Big data in weather forecasting: Applications and challenges." In *2017 International conference on big data analytics and computational intelligence (ICBDAC)*, pp. 138-142. IEEE, 2017.
- [12] Navadia, Sunil, Pintukumar Yadav, Jobin Thomas, and Shakila Shaikh. "Weather prediction: a novel approach for measuring and analyzing weather data." In *2017 International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud)(I-SMAC)*, pp. 414-417. IEEE, 2017.
- [13] Parameshachari, B. D. "Big data analytics on weather data: predictive analysis using multi node cluster architecture." *International Journal of Computer Applications* (2022): 0975-8887.