# Survey Using Big Data Tools for Doing Rainfall Prediction

**A. Saranya, R. Anandan**

*Abstract--- Big data collect large volume of data and gives the analysis of the data obtained using the Hadoop Framework in a fast and more efficient manner which is very accurate and makes our work very easy. Rainfall data is collected from the data set obtained and using Sqoop tool we get the data from MySql to HDFS architecture. The HDFS is mainly the database of Hadoop architecture which stores the data and distributes it to the various tools of Hadoop. It performs the accumulation of data in a way which makes decision analysis for the final output easy to obtain. Hence in our project we are focusing on getting the data from the data set and storing it in HDFS to get the analysis by Hadoop framework using Hive, Map Reduce and Pig to get the output of result which consists of all the analysis of rainfall in a city for all the years and give a clear perspective of the state of rainfall in the city at any moment of the year and also the analysis is shown through bar graph and pie chart which makes the understanding of analysis a little easier and it brings a little more significance to our work.*

*Keywords--- Hadoop, MySql, HDFS, Hive, MapReduce, Pig, Sqoop.*

## I. INTRODUCTION

Big Data is a concept which defines large amount of data which can be structured or unstructured. It is widely used as the amounts of records in a database are very high which can only be handled by a machine and not by human. It provides the platform which makes the handling of such large data much easier, the processing done using big data is much faster and easier. Plus there is no limit in the size of data which can be there for processing. Big Data works on the concept of three V's Velocity, Volume, Variety which clearly highlight the importance of Big Data.

Hadoop is a framework which provides a platform to work on Big Data. The Big Data tools Hive, Pig and Map-Reduce all are used to get the best analysis of the data no matter how large the size of the data is. The framework allows us to get of the data analysis by all of the three tools. The sqqop interacts with the database MySql and it fetches the data to the HDFS from where on all the data gets distributed to all the tools.

The tools get the data from the HDFS and the data is broken into packets or splitted so a separate analysis can be done of the Dataset. Then all the data is integrated and a combined analysis of the data is done. All the three tools use there analysis mechanisms to get the best possible output and the mechanism used to operate the data are either by query or programming methodology. The final output is very helpful with all the analysis. It is efficient and reliable

**A. Saranya,** Research Scholar, Department of Computer Science & Engineering, VELS Institute of Science and Technology.

Assistant Professor (O.G), Department of Software Engineering, SRM Institute of Science & Technology, (e-mail: aksaranya@gmail.com)

**Dr.R. Anandan,** Professor, Department of Computer Science & Engineering, VELS Institute of Science and Technology. (e-mail: anandan.se@velsuniv.ac.in)

which handles the data even if there is power loss and data consistency is maintained.

## II. SURVEY DETAILS

In paper[1] Rainfall information is gathered to anticipate the tempest notices from the hydrological information. This is considered as research thought as it expends immense number of records from the circulated framework. The paper depicts a answer for dealing with the data in light of spatial transient qualities utilizing Map Reduce Framework. The workload is characterized utilizing Support Vector Machine (SVM). It utilizes include determination and decrease calculation related with the dataset.

In paper[2] This manuscript was motivated by Congress where amid a talk of tempest water BMP plan, two inquiries were postured: 1. What precipitation force ought to be utilized for such plan? 2. To what extent (amid a tempest) should a BMP play out its water quality capacity? A more thorough inquiry is "All things considered, how hard would it be able to rain?" This paper tends to this last inquiry and sets up a system to address the other two inquiries. In this composition the experimental hyetograph that relates profundity and length, and along these lines whether a tempest is front stacked, back stacked.

In paper[3] Rainstorm Prediction using Support Vector Machine in Hadoop Cluste Precipitation information is gathered to foresee the tempest notices from the hydrological information. This is considered as an examination thought as it devours colossal number of records from the appropriated framework. This paper depicts a novel answer for deal with the information in view of spatial transient attributes utilizing a Map Reduce Framework. The workload is grouped utilizing bolster vector machine (SVM).

In paper[4] Precipitation profundities for different spans and frequencies, alluded to as profundity length recurrence (DDF) in this report, have numerous employments. A typical utilization of DDF is for the plan of structures that control and course restricted spillover, for example, parking areas, storm channels, and ducts. Another utilization of DDF is to drive waterway stream models that join precipitation attributes. Exact DDF gauges are vital for practical basic plans at stream intersections and for creating solid surge forecast models.

In paper[5] In the course of recent years, the creators and numerous others at Google have executed several unique reason calculations that procedure a lot of crude information, for example, crept records, web ask for logs,

and so on., to figure different sorts of determined information, for example, upset lists, different portrayals of the chart structure of web reports, synopses of the quantity of pages slithered per have, the arrangement of most incessant questions in a given day, and so on. Most such calculations are reasonably clear.

In paper[6] Most precipitation information is put away in positions that are difficult to break down and mine. In these configurations, the measure of information is colossal. In this paper, they propose procedures to compress the crude precipitation information into a model that encourages storm investigation and mining, and diminishes the information estimate. The outcome is to change over crude precipitation information into significant tempest driven information, which is then put away in a social database for simple investigation and mining.

### III. PROPOSED METHODOLOGY

In this project we are reading rainfall facts by way of using Hadoop device at the side of a few Hadoop ecosystems like HDFS, map reduce, sqoop, hive and pig. Maximum rainfall data is saved in formats which is not easy to understand and refine. The data volume is very high in the codes provided, by using the usage of these tools we are able to technique no hassle of information, no information misplaced problem, we can get excessive throughput, maintenance value additionally very less and it's miles an open source software, it is compatible on all of the systems given that it is java based. the output obtained by using processing huge chunks of information thru Hadoop structure is very useful for analyzing and predicting purposes.

*Data Source*

Dataset is obtained from the net which acts as the primary source of data. The data input is the input given by the client to the administrator which refines the data from the dataset using the Hadoop architecture.

### IV. SYSTEM ARCHITECTURE

When the rainfall data is provided initially stored in data base by using sqoop the stored in database is fetched and it is given to hadoop distributed file system. The data will be provided to different bigdata tools like hive, pig and mapreduce.
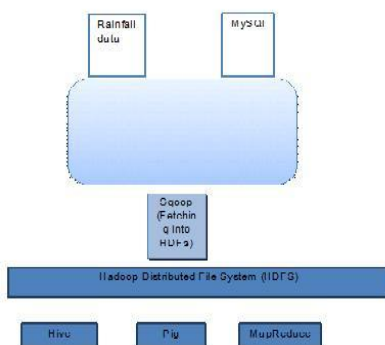


**Fig. 1: Architecture diagram**

### V. EXPERIMENTAL SETUP

Initially the data source which I took for prediction is sent to sqoop for importing into HDFS. Then the data is being accessed and by using appropriate data analytic tools analysis is carried out and the output is generated.
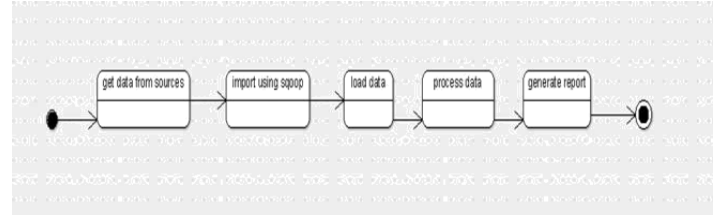


**Fig. 2: Information moved from first stage to report stage**

### VI. Experimental Results & Analysis

The data set consist of monthly, weekly, hourly based rain information of particular state. In that state the data is divided in to region wise and for all the region the information is provided in above informed manner. First fig show the unstructured data then second fig shows converted unstructured to structured data
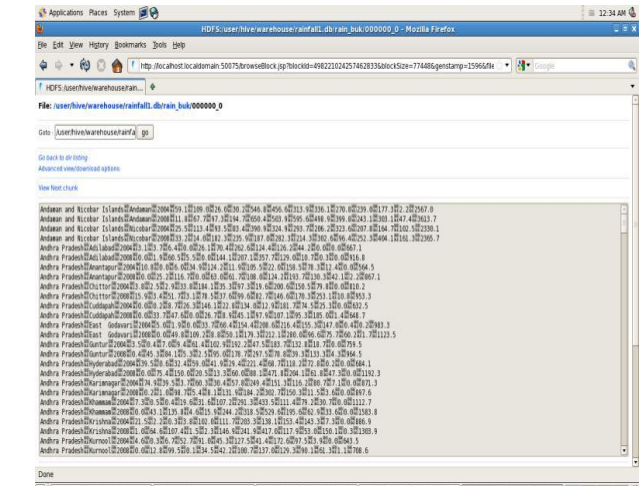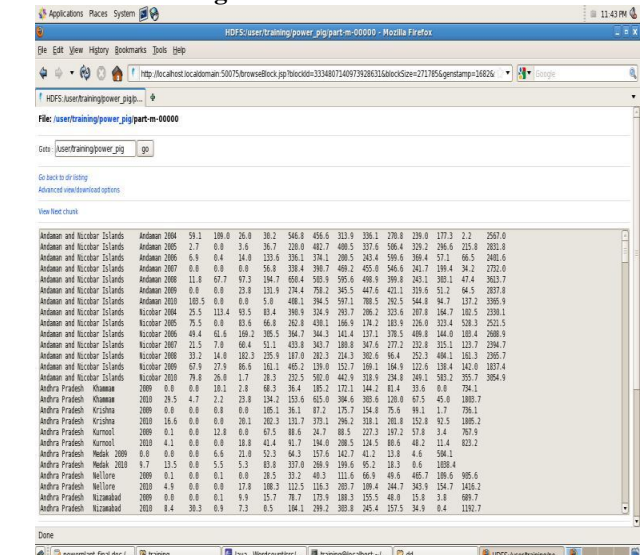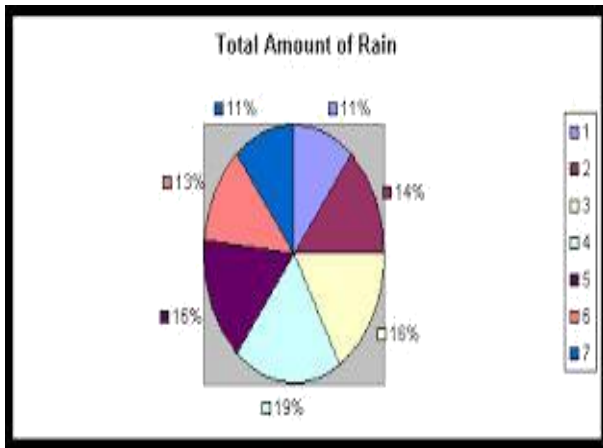


**Fig. 3: Unstructured Data**



**Fig. 4: Structured data**
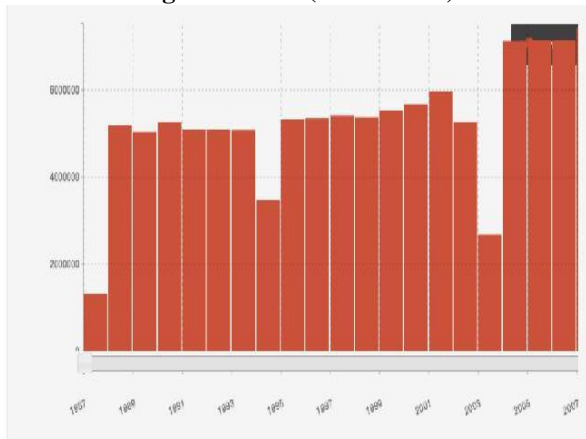
**Fig. 5: Results (Theoretical)**



**Fig. 6: Results (Graphs & Pie-Chart)**

The results are obtained after processing the dataset through sqoop and getting it into the HDFS which allows us to analyze the data through giving query in hive, commands in pig, and programming in Map Reduce which allows us to get the analysis of the rainfall data which is very easily understandable and analysis makes the prediction part very much easier.

The results obtained are very much accurate as the analysis are obtained using the Hadoop tools which are very efficient even for a large amount of data.

The result is the analysis of the rainfall in different periods in a particular state or a particular part of a district. The result can be obtained in a theoretical as well as a graphical manner which can be used to compare the data of a particular district with itself and as well as with a different state or a district.

The pie chart shows the total amount of rainfall in Tamil Nadu (example) in the previous years. The pie chart results directly show us the region with most rainfall and the year in which the rainfall was very high.

The theoretical as well as graphical and pie chart representation makes the analysis understanding very easier, which allows the prediction for future resources very easy. The results show us the complete analysis of the rainfall in a particular state in Tamil Nadu.

## VII.    RESULTS & ANALYSIS

The analysis of the results obtained can be done very easily because the output clearly depicts the area with high or low rainfall. The state in Tamil Nadu which experiences most of the rainfall goes up to 19% and the least rainfall in a

state goes up to 11% which makes us easy to predict the future data as well as get the clear image of rainfall in a state in Tamil Nadu, whether the rainfall makes it good for development of housing facilities or not. The analysis can be done that there is too much fluctuation of rainfall in Tamil Nadu over the coming years, which is clearly depicted in the diagram. So, it's not a good time in the year to get some projects going in Tamil Nadu.

The theoretical representation shows us the complete theoretical analysis of rainfall in a State in different years over a period. Precisely it gives us the data of how the rainfall has been in major parts of a state notably Andhra Pradesh, Andaman & Nicobar Islands in a period of time. The analysis or any comparison after this can be easily done on the basis of the output. The analysis shows us that Andaman has a rainfall percentage than Andhra in the coming years and we can clearly predict that this pattern will be followed in the coming years.

## VIII.    CONCLUSION& FUTURE WORK

In this work, the storm class records are designed and implemented within the map reduce characteristic. To evaluate the rainfall information in Hadoop environment. Hadoop environment is hive, pig, map reduce, if you want evaluation to find some deep evaluation on rainfall, which clearly validate the efficacy of our fashions in phrases of different programs and furthermore reveal a few exhilarating rainfall water portal locating. In destiny the sparks one hundred times quicker than Hadoop.

"We are using spark as we are able to get result hundred instances faster than Hadoop. The secret is that it runs in-memory at the cluster and that it is not tied to map. This makes repeated access to the same entry a lot quicker. Spark can run as a standalone or on top of Hadoop YARN, in which it is able to examine information at once from HDFS."

## REFERENCE PAPERS

[1]  Rainfall analysis and rainstorm prediction using MapReduce Framework C.P Shabariram ; K.E. Kannammal ; T. Manoj praphakar Computer Communication and Informatics (ICCCI), 2016 International Conference on 30 May 2016.

[2]  Rainstorm Prediction using Support Vector Machine in Hadoop Cluster C.P Shabariram data mining and knowledge engineering, vol7 2015.

[3]  Depth-Duration Frequency of Precipitation for Texas, U.S. Department of the Interior, U.S. Geological Survey, Page Contact Information: Contact      USGS, December 07 2016.

[4]  Map Reduce: Simplified Data Processing on Large Clusters, Jeffrey Dean and Sanjay Ghemawat, OSDI 2004.

[5]  Extracting storm-centric characteristics from raw rainfall data for storm analysis and mining, Kulsawasd Jitkajornwanich, Ramez Elmasri, John McEnery, Chengkai Li, Univ. of Texas at Arlington, Arlington, TX, ACM Digitstal library, nov 2012.

[6]  Shabariram, C.P. (2015).Rainstorm Prediction using Support Vector Machine in Hadoop Cluster. Data Mining and Knowledge Engineering, 7(9), 316-320.

[7] Jitkajornwanich, K., Elmasri, R., McEnery, J., & Li, C. (2012, November). Extracting storm-centric characteristics from raw rainfall data for storm analysis and mining. In Proceedings of the 1st ACM SIGSPATIAL International Workshop on Analytics for Big Geospatial Data (pp. 91-99). ACM.

[8] Dean, J., & Ghemawat, S. (2008). MapReduce: simplified data processing on large clusters. Communications of the ACM, 51(1), 107-113.

[9] Cleveland, T. G., & Thompson, D. B. (2008). Rainfall Intensity in Design. In World Environmental and Water Resources Congress 2008: Ahupua'A (pp. 1-10).

[10] Asquith, W. H. (1998). Depth-duration frequency of precipitation for Texas. US Department of the Interior, US Geological Survey.

506