



Novel Approach to Seat Matrix Prediction using Hadoop

Ravikiran M D^{1*}, Poonam Gouli²

^{1,2}Department of Computer Science & Engineering, RVCE, Bengaluru

Corresponding Author: ravikirandevadiga@gmail.com

DOI: <https://doi.org/10.26438/ijcse/v7si13.2532> | Available online at: www.ijcseonline.org

Abstract— Big data is a term that is used to describe huge data sets having large, varied and complex structure with the hardness of analysing, storing and visualizing for further analysis or results. The analysis into large amounts of data to expose secret correlations and hidden patterns can be termed as big data analytics. [1] Big Data has proved to be beneficial as it helps to gain richer and deeper insights into the underlying mass of data. Common Entrance Test is a flat form for the students to opt for colleges to pursue Under graduation. Every year approximately 150000 students take up CET. Thus, these students will compete for seat among 220 engineering colleges across Karnataka, that are enrolled to the CET cell. Student can opt for a College based on availability of seat or branch for the rank they obtain in CET. Predictive analytics is the basic enabler for big data. On a day to day basis, Businesses collect large quantity of customer data which is used by predictive analytics along with historical data, coupled with customer insight, to forecast future events. For predicting a college several tools are taken into consideration they are: HBase for database, MapReduce for data processing mahout's distributed naive Bayes classification for classifying and training data.

Keywords—CET, HADOOP, PIG, MAHOUT, MAPREDUCE

I. INTRODUCTION

Allocation of seats to a college as per KEA norms is an elaborate process. The patterns and the schemes that are followed for the allotment each year are subjected to timely reforms that are imposed by the governing organization. Two of the recently followed strategies are Offline seat matrix display and Online seat allotment in terms of mock and actual placement.

This paper comprises of an overview of big data's scope, methods, samples, challenges associated with it. The project also aims at putting an end to the eagerness of the parents along with the forthcoming candidates who are waiting for their turn to get into a reputed college to pursue their graduation / post-graduation. This tool will help the students/parents to understand the probability of getting a seat in a college. Section I contains the Introduction of “*Novel approach to seat matrix prediction using Hadoop*”, Section II contains the Related work done for the paper, Section III contains the Experiment and Data, Section IV contains the Methodology, Section V contains Results, Section VI contains conclusion and future Scope, Section VII contains References.

II. RELATED WORK

Literature survey describes the existing and established theory and research thus simplifying the process of achieving the end result of the project. It also helps in identifying the deviations in the analysis that are done over the same concept. Here, a list of all the papers that were referred in the due course of the project as been listed.

In the distributed environment [2] mahout works well. Mahout effectively scales in the cloud using the Apache Hadoop library and provides the developer with a ready-to-use framework for data mining tasks on large volumes of data. The two basic deciding factors considered from the candidates end in the admission process are: AIEEE rank and family pressure. This makes it difficult for them to decide of another branch in case they don't get a seat in the desired branch.

Regular evaluation of the teaching patterns of universities [3] plays an important role in improving management levels. The methodology used in college teaching has the characteristics of the Markov chain. Using the Markov chain, it selects four ranks of comprehensive evaluation results in different specialties in order to carry out an analysis in a state shift; it then carries out a gradual assessment of specialized teaching through the progressive matrix of the transfer matrix and its degrees of efficiency;

finally, according to the vector of a steady state, it forecasts the long-term effect of specialized teaching.

Seat matrix prediction is a non-linear classification problem. This paper [4] proposes a new approach for Prediction of college entrance examination (CEE) based on the learning algorithm of vector machines. Data to be collected include ranking in all subjects, CEE score, the number of college admission plan and relevant data of the past two years. With all the data gathered a training set is formed.

III. EXPERIMENTS AND DATA

Relevant details should be given including experimental design and the technique (s) used along with appropriate statistical methods used clearly along with the year of experimentation (field and laboratory).

3.1 Data Description

Data set was downloaded from the kea website. Dataset used is in PDF format. Inside the PDF format the data was there in xml format. The dataset contains Rank, CET NO, Name of the student, Category, College/Course etc. from KEA Website.

3.2 Research Challenge

One of the main challenges encountered while beginning the project was that of Data Collection, as it wasn't available in the desired format. It was placed in such a manner that the main PDF document housed has an XML document inside it, which can be referred to as semi structured format. This cannot be converted directly into csv format. Data wrangling method from Visual Basic script was used to convert the original document into csv format. Next challenge that came in was that that algorithm selection. Since a large collection of libraries were available in Mahout such as distributed naive based classification method was chosen as the base algorithm.

3.3 Motivation

Big Data is the most important asset of any organization which is further processed to produce useful information [5]. Big Data techniques are widely used in industrial sectors to generate patterns that are helpful for earning more profits and expansion of business. One such tool can be to make appropriate predictions in the education sector where in decisions pertaining to college allotment can be made. It can serve to be an effective means of information to both institutions and aspiring students to represent and understand the statistics.

3.4 Methodology

3.4.1 Extraction transformation load:

HBase is another example of an environment for non-relational data management that distributes massive data sets across the Hadoop framework. It has a column-based data layout that can be used to store and manipulate large data tables. It is a reasonable alternative as a persistent storage platform when running map reduce application. Well-timed updating, analysis and storage of the unstructured data will be conducted on a regular basis by utilizing HBase.

3.4.2 Training Phase:

Direct querying of unstructured data is not feasible in terms of big data. Hence a combined approach involving map reduce application along with mahout library for machine learning algorithms, text mining and binary search tree algorithms is made use of for achieving the task of querying into the big data.

3.4.3 Execution phase:

The two earlier stages are collaborated with each other.

A good walk-through of the previously stored data values is made at appropriate locations to determine what the future trend might be, for an y given scenario

3.4.4 Results:

For example, if a candidate has secured 150th rank and belongs to the GM category. Then this software will make a thorough analysis of the college allotment data pertaining to the earlier years. i.e. the allotment strategy followed for the 150th and its neighbouring rank holders to arrive at a predicted conclusion that the current candidate may or may not be allotted to the same college as that of the earlier pattern that was followed.

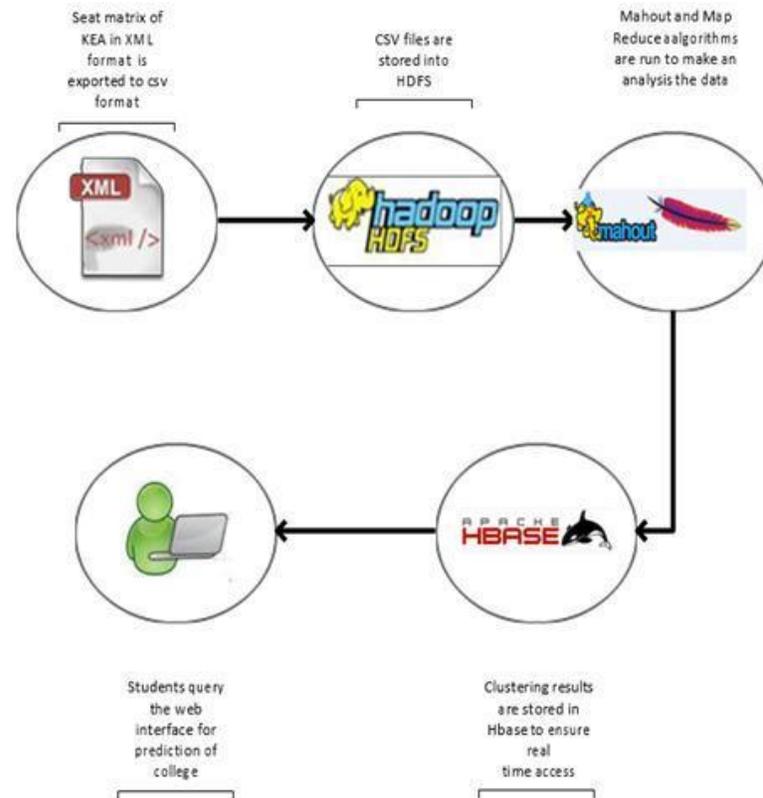


Fig 1 Methodology of seat matrix prediction

IV. ANALYSIS AND DESIGN

Steps followed to Prediction of Rank

1. Data Preprocessing
2. Store the KEA DATA to HDFS.
3. Mapper Phase.
4. Reducer Phase
5. Prediction phase with mahout algorithm,
6. Final Predicted Rank Stored in HBase.

- 1) Data Preprocessing: The Data was in PDF Format with unformatted Data Cells. PDF data is Converted to CSV with removing of unnecessary column values using data wrangling method.
- 2) Store the KEA DATA to HDFS: Converted CSV files from PDF will be saved to the HDFS file systems for Data post processing. The files will be divided into chunks of data in HDFS.
- 3) Mapper Phase: In the Mapper phase all CSV files from the HDFS will be collected together and key value pairs will be produced as a intermediate output
- 4) Reducer Phase: In the reducer phase intermediate output from mapper phase will be taken as input. And intermediate results will be sorted and results of reducer phase will be stored to HDFS file system.
- 5) Prediction phase with mahout algorithm: In the prediction phase previous results from the Reducer phase will be considered as input. we will user mahout machine learning module from apache to predict the rank.
- 6) Final Predicted Rank Stored in HBase: At last we will find out the predicted rank from the prediction algorithm we will store this results our own database in Hbase to make sure that remaining candidates will not get the same college as previous candidates predicted result.

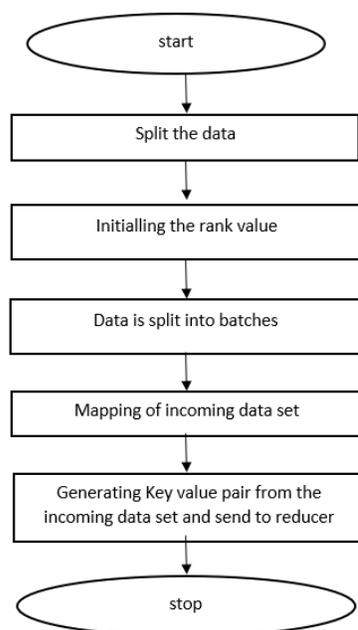


FIG 2: Flowchart of Mapper Phase

In mapper module, the important concept is generating $\langle \text{key,value} \rangle$ pairs. These pairs are stored in an intermediate file which will be accessed by reducer. Figure 3 shows flowchart for mapper Phase.

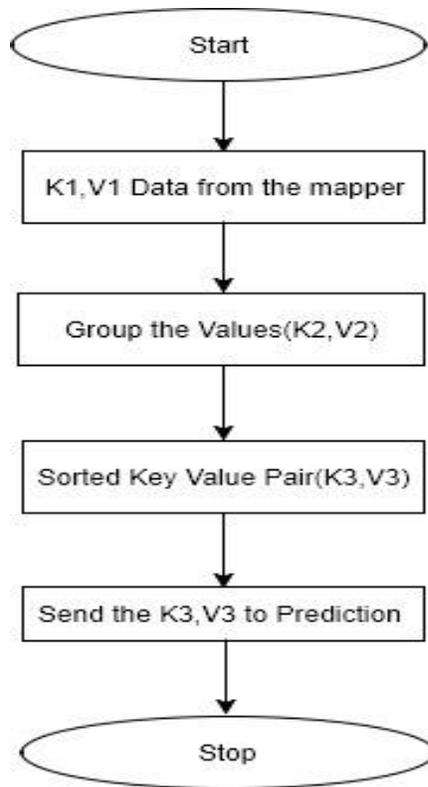


FIG 3: Flowchart of Reducer Phase

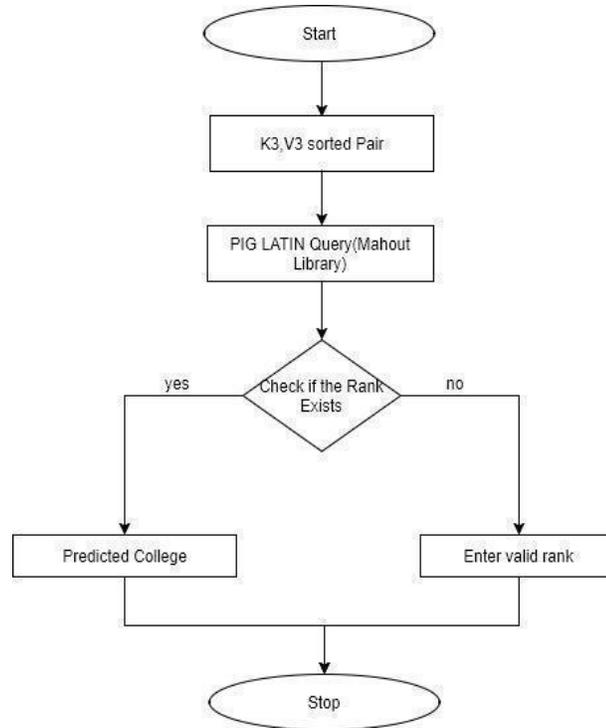


Fig 4 Flowchart of Prediction phase

In mapper module, the important concept is generating <key, value> pairs. These pairs are stored in an intermediate file which will be accessed by reducer. The major work of reducer is to collect data from mapper and then sort the desired result. Figure 5.4 shows flowchart for reducer. The reducer first collects the intermediate data from all mappers [6]. It calculates the incoming data from these intermediate data. This will be the updated or create the new row in the database. The data of whole process will be in sorted order and these values will be the input for prediction. The final output from the reducer Figure 5.4 shows flowchart for prediction i.e. the (key,value) pair (K3, V3) is sent to the PIG Latin query which makes use of the mahout library. This query performs the basic check that searches whether the rank is present in the CSV file. If it is present, then the probable/predicted college will be displayed as output. If it's not present, then an alert pops up saying "Please enter a valid rank".

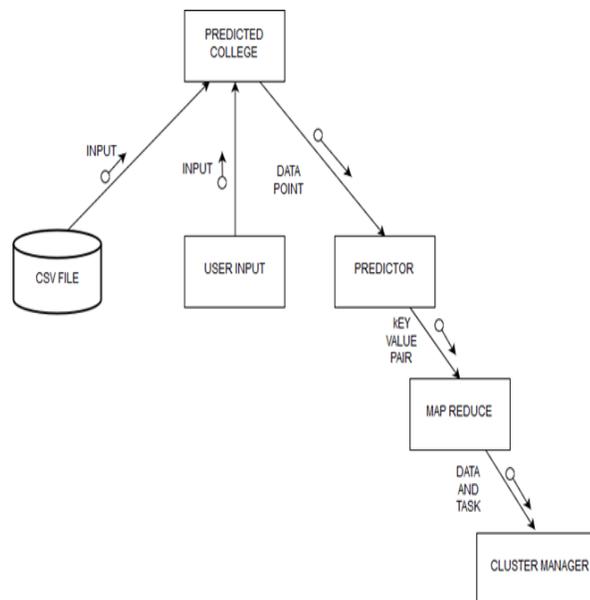


Fig 5 Structural Chart of Seat Matrix Prediction


```

Success!
Job Stats (time in seconds):
JobId Alias Feature Outputs
job_local1268283861_0004 ket_ket_short3,newone3 MAP_ONLY file:/tmp/temp523829163/tmp1279997182,

Input(s):
Successfully read records from: "/home/cloudera/Desktop/kes.txt"

Output(s):
Successfully stored records in: "file:/tmp/temp523829163/tmp1279997182"

Job DAG:
job_local1268283861_0004

2017-12-24 23:46:48,899 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2017-12-24 23:46:48,899 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2017-12-24 23:46:48,899 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2017-12-24 23:46:48,899 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2017-12-24 23:46:48,189 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2017-12-24 23:46:48,189 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(3490,"College :9 - P E S Institute of Technology,Bangalore.")

```

Fig 9 Predicted College in Command Prompt

The below command shows how to store the above result into a file in HDFS:

```
store newone3 into '/user/cloudera/ketoutput' using PigStorage('\t');
```

All the results that are stored in this manner can be used for comparisons whenever required. The contents of the resultant file will be stored with the name "part-m-00000". The below Fig. 10 shows a view when the file is opened.

Metrics used to verify the results are Rank and college. Its approximately accurate for the year 2014 after verifying the 2014 CET [7] allotment sheet. But for the year 2016 it was not accurate. Need to include Category and branch to get the best accurate result.

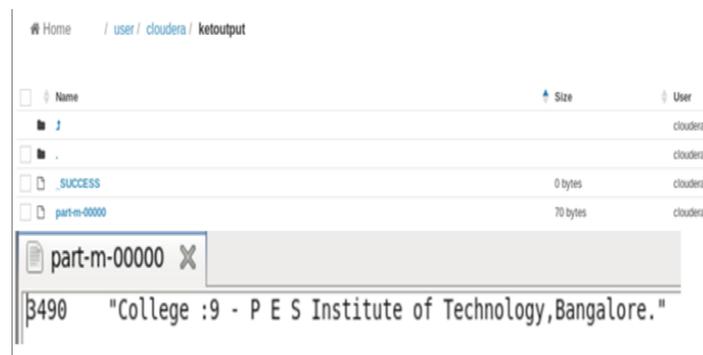


Fig 10 Predicted College stored in HBase and HDFS

Conclusion and Future Scope

Big data is a term for massive data sets having large, more varied and complex structure with the difficulties of storing, analyzing and visualizing for further processes or results. Predictive analytics is an enabler of big data: Businesses collect vast amounts of real-time customer data and predictive analytics uses this historical data, combined with customer insight, to predict future events. Predictive analytics enable organizations to use big data (both stored and real-time) to move from a historical view to a forward-looking perspective of the customer. Advantages 1) Predicting the college based on rank is made easier. 2) Time consumed in prediction is drastically reduced. 3) Predicted results can be stored for future use. 3) Result matrix for the years to come can be paced along with the existing data in HDFS. Limitations and future work 1) As of now Web UI is not made available to the user to key in the input. Prediction of college can be made by taking candidates belonging to reserved categories into consideration. The main conclusions of the study may be presented in a short Conclusion Section. In this section, the author(s) should also briefly discuss the limitations of the research and Future Scope for improvement.

REFERENCES

- [1] Sanjay Kapoor , Deepak Kumar , Sandeep Jain , Brijesh Khandelwal , Anil Mittal, Manoj Kumar "An optimized e- Allotment algorithm for admissions in educational institutes," International Conference on Technology Enhanced Education (ICTEE), IEEE, 2012, pp.17-25.
- [2] Rotsnarani Sethy, "Big Data Analysis using Hadoop: A Survey", International Journal of Advanced Research in Computer Science and Software Engineering, Volume 5, Issue 7, July 2015, pp 1153- 1157

- [3] SHIVAM AGARWAL, "DATA MINING CONCEPTS AND TECHNIQUES", International Conference on Machine Intelligence and Research Advancement, 2013, pp 203-207
- [4] B.N. Lakshmi, G.H. Raghunandhan, "A Conceptual Overview of Data Mining," National Conference on Innovations in Emerging Technology-2011, pp 27-33
- [5] V.K.Deepa, J. Remy R. Geetha, "Rapid Development of Applications in Data Mining", International Conference on Green High Performance Computing March 14-15, 2013 pp 13-16
- [6] Jinlong Wang, Jing Liu, Russell Higgs, Li Zhou, "The Application of Data Mining Technology to Big Data", IEEE International Conference on Computational Science and Engineering (CSE) and IEEE International Conference on Embedded and Ubiquitous Computing (EUC) 21-24 July 2017 pp 284-289
- [7] <http://kea.kar.nic.in/>

AUTHORS PROFILE

Mr. Ravikiran Devadiga pursued Bachelor of Engineering from University of VTU, Belgaum in 2012 and M.Tech from VTU University in year 2019. He is currently working as Programmer in Department of Computer Sciences, RNSIT since 2012.



Dr Poonam G is Currently Working as Associate Prof in RVCE Bangalore. Her area of interest in Research are Distributed Clustering, MapReduce And Machine Learning.

