



VIT[®]

Vellore Institute of Technology

(Deemed to be University under section 3 of UGC Act, 1956)

Parallel and Distributed Computing

Review - 3

PARALLEL DISTRIBUTED ALGORITHM **FOR RAILWAY SIGNAL FAULT** **DETECTION**

*Project done under the guidance of Prof. Narayanan
Prasanth*

Members :

Arshia Shetty - 17BCE0076

Nikila GS - 17BCE0272

Priyanshu Panda - 17BCI0085

Sagar Suman Dwibedi - 17BCI0098

Anish Nagesh - 17BCI0192

DECLARATION

I hereby declare that the report entitled “PARALLEL DISTRIBUTIVE ALGORITHM FOR RAILWAY SIGNAL FAULT DETECTION” submitted by Priyanshu Panda - 17BCI0085, Sagar Suman Dwibedi - 17BCI0098, Anish Nagesh - 17BCI0192, Arshia Shetty - 17BCE0076 and Nikila GS - 17BCE0272, for the course CSE4001 Parallel and Distributed Computing (EPJ) to VIT is a record of bonafide work carried out by us under the supervision of Dr.N.Narayanan Prasanth.

I further declare that the work reported in this report has not been submitted and will not be submitted, either in part or in full, for any other courses in this institute or any other institute or university.

Place : Vellore

Date : 2/06/20

Signature of the Candidate

ACKNOWLEDGEMENT

We would like to specially thank Dr.N.Narayanan Prasanth, School of Computer Science and Engineering, whose immense guidance helped us to complete this project. We wish to express our sincere thanks to Honourable Chancellor, Dr.G.Viswanathan; respected Vice-chancellor, Dr. Anand Samuel of this prestigious college for giving us an excellent world class academic environment and facilities for pursuing our degree. We would also like to thank our friends for their timely help and suggestions rendered for the successful completion of this project. This acknowledgement would be incomplete without expressing our whole hearted thanks to our parents for their continuous support and guidance.

ABSTRACT :

To explore the data processing of high-speed railway fault signal diagnosis based on MapReduce algorithm, the partitioning of the data set has been improved and implemented using K- means clustering algorithm. After which SON algorithm and Markov models have been applied to get better and improved results. In MapReduce parallelization process, the data partition matrix T_k was stored in line segmentation, the computing load was distributed in every node of cluster, and the time consumption of mobile data matrix and the consumption of partitioned matrix were calculated.

SR. NO.	CONTENTS	PAGE NO.
	Acknowledgement	3
	Abstract	4
	Table of contents	5
	List of figures	6
	List of tables	7
	List of abbreviations	8
1	INTRODUCTION	
1.1	Objective	9
1.2	Motivation	9
2	LITERATURE SURVEY	10
3	TECHNICAL SPECIFICATION	15
4	DESIGN	15
5	PROPOSED SYSTEM	
5.1	Modules	18
5.2	Algorithm	20
6	RESULTS AND DISCUSSION	26
7	CONCLUSION	31
8	REFERENCES	32

List of Figures

FIGURE	TITLE	PAGE NO
1.	Supercomputing performance trends, data source, where MFlop/s means million floating point operations per second.	18
2.	Diagrammatic representation of K-Means algorithm	19

List of Tables

Table No.	Title	Page no.
1.	Literature survey	12

List of Abbreviations

DTMC	Discrete Time Markov Chain
TF-ID	Term Frequency – Inverse Document Frequency
HTTP	Hypertext Transfer Protocol

INTRODUCTION :

Objective:

The main objective of this project is to develop an efficient system such that we will be able to minimize the accidents that occur on a regular basis on the railway tracks. The main aim of this project is to reduce the unfortunate events that cause damage to people and the government.

Motivation:

With the increase of speed of India's high-speed railway and the continuous improvement of the railway information system, the conditions of collecting more railway running information are now available. At present, the running high-speed railway train, through the deployment of a large number of sensors, collects a variety of data. However, the traditional vibration data feature extraction and analysis technology is running on a single machine. This kind of technology, in the mass vibration data acquired by sensors, exposed the shortcomings of long processing time, various artificial intervention, and poor capability of processing big data file and so on. The emergence of cloud computing technology provides a way of thinking to solve the above problems. Map Reduce is an effective parallel computing framework of processing big data, which is one of the main models of cloud computing, and can automatically assign tasks and realize task balance. The working principle, operating mechanism and fault tolerance mechanism of Map Reduce calculation model are studied. In addition, combined with the characteristics of association rule generation algorithm, the traditional parallel algorithm is improved and the parallel optimization scheme of association rules algorithm based on Map Reduce is proposed.

LITERATURE SURVEY :

AUTHOR	PROBLEM STATEMENT	METHODOLOGY	EVALUATION METRIC	RESULT	LIMITATION
Optimization of a genetic algorithm for road traffic network division using a distributed/parallel genetic algorithm By Tomas Potuzak Published in: 2016 9th International Conference on Human System Interactions (HSI) Date:- 8 th July, 2016	Computer simulation of road traffic is an important tool for analysis and control of existing and projected road traffic networks. Such a simulation can be very time-consuming for large road traffic networks (e.g., big cities or even states), therefore many road traffic simulators should be performed in a distributed computing environment.	A distributed/parallel genetic algorithm (DGA) for the optimization of the parameters settings of the genetic algorithm for the road traffic network division. Necessary to divide the simulated road traffic network into sub-networks, whose simulations are then performed as processes on individual nodes of the distributed computer.	The parameters include:- <ol style="list-style-type: none"> 1. Number of DGA generations 2. Number of DGA parent individuals. 3. The number of mutations per DGA individual 4. Type of DGA selection 5. Type of DGA crossover. 	In order to determine, whether the new DGA settings obtained from the OGA is better than the original DGA settings, we performed a set of tests on five road traffic networks - three regular square grids, which were already used by the OGA. As expected the new DGA settings gave better results.	<ol style="list-style-type: none"> 1. OGA computation was regularly interrupted during test traffic. 2. Production of irregularly high mutations could cause negative results. 3. However, the better results are for the price of far higher DGA computation time. With the new settings, the DGA is on average 185 times slower than with the original settings.

Research of an Improved Apriori Algorithm in Data Mining Association Rules By Jiao Yabing Published In: International Journal of Computer and Communication Engineering, Vol. 2, No. 1, January 2013 [17BCI0098]	Apriori algorithm is the classic algorithm of association rules, which enumerate all of the frequent item sets. When this algorithm encountered dense data due to the large number of long patterns emerge, this algorithm's performance declined dramatically. This paper illustrates about enhancing the algorithm for better performance	In the Apriori algorithm, C_{k-1} is compared with support level once it was found. Item sets less than the support level will be pruned and L_{k-1} will come out which will connect with itself and lead to C_k . The optimized algorithm is that, before the candidate item sets C_k come out, further prune L_{k-1} , count the times of all items occurred in L_{k-1} , delete item sets with this number less than $k-1$ in L_{k-1} . In this way, the number of connecting items sets will decrease, so that the number of candidate items will decline.	This study is focused on how to solve the efficient problems of Apriori algorithm and raise another association rules mining algorithm. This has certain reference value to research and solve the issues of data expiation and information lacking.	This decreases the possibility of combination, decline the number of candidate item sets in C_k , and reduce the times to repeat the process. For large database, this algorithm can obviously save time cost and increase the efficiency of data mining. This is what Apriori algorithm do not have.	Although this process can decline the number of candidate item sets in C_k and reduce time cost of data mining, the price of it is pruning frequent item set, which could cost certain time. For dense database (such as, telecom, population census, etc.), as large amounts of long forms occur, the efficiency of this algorithm is higher than Apriori obviously
--	---	--	--	---	--

A novel fuzzy logic controller based automatic caution order system for Indian railways By K.P Ajeesh R Ashis , Mary Helna Muttath Nidheesh.K. Chandran C.R Sreedevi Published in: 2012 IEEE International Conference on Advanced Communication Control and Computing Technologies (ICACCCT) Date:- 04 October 2012	This paper focuses a new way of approach to find the solution for the artificial intelligent braking system in train using the fuzzy logic controller. The new way of approach is to reduce the man power and to automate the system.	The main function of the fuzzy logic controller used here is to automatically stop the train in each station without any manual procedure of stopping the train. The fuzzy logic controller in train gets activated about 500m from the station so that the train stops at the station smoothly and automatically. The fuzzy controller takes the decision with reference to the speed and distance of the train.	1. Crisp Input and Output Map in Matlab:- The fuzzy predictive control technology is compared with the conventional control technology. The simulation results show that the performance of train safety, comfort, parking accuracy and other performance indicators have been improved significantly by using fuzzy predictive controller.	In the simulation circuit the speed of the motor is one input of the fuzzy logic controller, another input is distance. Distance is calculated from the speed. The fuzzy logic controller output is given to the motor. This is the test circuit of the fuzzy logic controller. So the motor stops when the distance becomes zero.	<ol style="list-style-type: none"> 1. It is tedious to develop fuzzy rules and membership functions and fuzzy outputs can be interpreted in a number of ways making analysis difficult. 2. In addition, it requires lot of data and expertise to develop a fuzzy system.
---	---	---	--	--	--

<p>Comparative Analysis of Map Reduce Scheduling Algorithms</p> <p>By Sonia Sharma, Dr. <u>Parag Jain</u>.</p> <p>Published In: Jour of <u>Adv Research in Dynamical & Control Systems</u>, Vol. 11, 10- Special Issue, 2019</p> <p>[17BCI0098]</p>	<p>When Hadoop schedules reduce tasks, it neither exploits data locality nor addresses partitioning skew present in some MapReduce applications. Scheduling of mixed real-time and non-real-time applications in MapReduce environment is a challenging problem but accepts only restricted attention. First In First Out (FIFO) is the default job scheduling strategy of Hadoop, but it cannot guarantee that the job will be completed by a specific deadline.</p>	<p>In this work they have compared the various scheduling algorithms and their counterparts using MapReduce framework. They have developed FIFO, Round Robin, ETF and LATE schedulers with MapReduce compatible implementation of the same namely MR-FIFO, MR-RR, MR-ETF and MR-LATE.</p>	<p>To calculate the performance gain the one algorithm for which the value is calculated is subtracted the reference algorithm, divided by the gain by the original algorithm turnaround time. Finally, multiply it by 100 to get the percentage change.</p>	<p>The MapReduce framework along with the overview of different scheduling techniques with their applications is discussed in this paper. Different scheduling techniques to enhance the data locality, make span, efficiency, fairness and performance are discussed with special emphasis on its real time applications.</p> <p>This work they have compared the various scheduling algorithms and their counterparts using MapReduce framework and developed FIFO, Round Robin, ETF and LATE schedulers with MapReduce compatible implementation of the same namely MR-FIFO, MR-RR, MR-ETF and MR-LATE. Classic benchmarks of Max of Eigen Vectors was used to evaluate the performance of our scheduling algorithm. For each test every benchmark is run for 100-1000 tasks and the average execution time is used as the result.</p>	<p>MR-FAIR scheduler needs to be improved in using</p> <ol style="list-style-type: none"> (1) Application of Optimization algorithm to map reduce (2) Application of non-uniform task distributions. (3) Trying to enable the optimization based algorithm to schedule real time jobs.
--	---	---	--	---	---

<p>Parallel computing in railway research</p> <p>By: Qing Wu, <u>Maksym Spirvagin</u>, Colin <u>cole</u> & Tim <u>McSweeney</u>.</p> <p>Date: Dec 2018</p>	<p>Parallel computing applications in railway research as well as the enabling techniques used for the purpose.</p>	<p>one publication can use multiple enabling techniques. For example, the Domain Decomposition Method (DDM) technique has to be used with another communication technique such as the MPI to successfully conduct parallel computing</p>	<p>Iterative optimisations, such as Particle Swarm Optimisations (PSO), Genetic Algorithm (GA) and Simulated Annealing (SA), improve a group of solutions (plans, designs, etc.) iteratively in loops. Almost all iterative optimisation methods have algorithm structures that are very suitable for parallel computing, especially during the fitness assessment step</p>	<p>parallel computing applications and enabling techniques that are used in railway research were reviewed. Nine enabling techniques were reviewed and MPI, DDM and Hadoop & Apache are the top three most widely used enabling techniques. Seven major application</p>	<ol style="list-style-type: none"> (1) the engineering basis using clusters for parallel computing was arguably invented by Gene Amdahl in 1967 (2) PCs were only made available in the late 1970s (3) the MPI and <u>OpenMP</u> standards were only
--	---	--	---	---	---

DISTRIBUTED COMPUTING IN MOBILE ENVIRONMENT By Apekshit Sharma, Chandrakant Sharma	The classic models and solutions are not adapted to mobile world's hence new models; new problems and obviously new solutions should be designed. This paper advocates for the creation of a common view for mobile worlds based on the similarities between them.	The RL algorithm assumes that there is a unique token initially and utilizes the partial reversal to maintain a token oriented DAG (directed acyclic graph). A node should gain the token along the DAG to access the shared resource. The proposed algorithm in this paper is also adapted from the RL algorithm. The RL algorithm and its adaptations are sensitive to link forming and link breaking.	State: Indicates whether node i is in the ES, CS, or NCS state. Initially, state = NCS. N : The set of all nodes (neighbors) in direct wireless contact with node i . Initially, N contains all neighbors of node i . $hVector$: An array of triplets representing node i is view of height of node j , $j \in N$. Initially, $hVector[j] = \text{height}$	1. We propose a distributed algorithm for solving the h-out-of-k mutual exclusion problem for ad hoc mobile networks. 2. The proposed algorithm is sensitive to link forming and link breaking and thus is suitable for ad hoc mobile networks. 3. This paper pointed out similarities among common mobile systems	1. The nodes may not have unique node identifiers. 2. Node failures may occur. 3. Communication links can be unidirectional and not FIFO. 4. Incipient link failures may not be detectable. 5. Message delays don't obey the triangle inequality
--	--	--	--	--	--

Parallel and distributed kmeans to identify the translation initiation site of proteins By:- Laerte M. Rodrigues ; Luis E. Zárate ; Cristiane N. Nobre ; Henrique C. Freitas Published in: 2012 IEEE International Conference on Systems, Man, and Cybernetics (SMC)	Use of Parallel kmeans algorithm to enhance performance of TIS and non-TIS class balancing for databases where the imbalance is approximately 1:40.	The strategy is based on MPI (Message Passing Interface) and OpenMP (Open Multiprocessing) to exploit shared memories, and collective communications to reduce the influence of network. This way, the proposal is to achieve efficiency and scalability.	1. The number of databases . 2. The splitting of databases and clustering. 3. Stability tendency of centroids .	The processing time gain per iteration demonstrates a high speedup , indicating that the solution considerably reduces processing time and guarantees a response in a much lower time.	1. The efficiency indicates that the processing cores are being underused. 2. For small databases the small number of operations reduces the efficiency. 3. The time taken for all of the databases is directly related to the quantity of clusters that should be formed. Which cause longer processing.
---	---	---	---	--	---

Paragraph	Paragraph	Paragraph	Paragraph	Paragraph	Paragraph
<p>A quick Otsu-Kmeans algorithm for the internal pipeline detection</p> <p>By:- Tao Song ; Ji Li ; <u>Xue Han</u> ; Lei Shao</p> <p>Published in:- 2017 IEEE International Conference on Mechatronics and Automation (ICMA)</p>	<p>The traditional <u>kmeans</u> algorithm is sensitive to the initial clustering center and the slow processing speed of large data set. It is adopted to make the image partition processing and using the <u>otsu</u> algorithm to find the right clustering center.</p>	<p>Aiming at the defects of the pipeline, this paper proposes an improved <u>kmeans</u> clustering algorithm for image segmentation. First of all, according to the characteristics of image processing, the original image were divided into blocks and the sub block variance is greater than 10%</p>	<ol style="list-style-type: none"> 1. Value of Segmentation threshold . 2. <u>Gray</u> level of the image 3. RGB level of the image. 	<p>The experiment shows that the clustering effect in this paper is much better than Otsu method based on threshold and region growing method based on region, and the speed is improved than the traditional k-means algorithm. Pipeline robot has an important</p>	<ol style="list-style-type: none"> 1. RGB, the components of relevance is too large. 2. Each channel into the luminance information, easily affected by the surrounding environment. 3. The <u>color</u> of the human cognitive process is not matching, is not suitable for analysis and segmentation of <u>color</u> image.

		<p>of the original variance. Then, image segmentation is performed on this sub block by using Otsu, and the two kind of mean value is obtained as the initial clustering center of <u>kmeans</u></p>		<p>value in inspecting small pipes, and effective processing the collected images is key to reliable operation of the inspection robot.</p>	
--	--	--	--	---	--

AUTHOR	PROBLEM STATEMENT	METHODOLOGY	EVALUATION METRIC	RESULT	LIMITATION
<p>Markov Decision Process-Based Resilience Enhancement for Distribution Systems: An Approximate Dynamic Programming Approach</p> <p>Chong Wang , Member, IEEE, Ping Ju , Senior Member, IEEE, Shunbo Lei , Member, IEEE</p> <p>2020</p>	<p>Because failures in distribution systems caused by extreme weather events directly result in consumers' outages, this paper proposes a state-based decision-making model with the objective of mitigating loss of load to improve the distribution system resilience throughout the unfolding events.</p>	<p>A recursive optimization model based on Markov decision processes (MDP) is developed to make state-based actions. To overcome the curse of dimensionality caused by enormous states and actions, an approximate dynamic programming (ADP) approach based on post-decision states and iteration is used to solve the proposed MDP-based model.</p>	<p>Two test systems are used to verify the proposed model and the algorithm. The first system is the IEEE 33-bus system, and the second system is the IEEE 123-bus system. The cases are tested in MATLAB 2017a using the CPLEX 12.6 solver on computers with 3.1 GHz i5 processors and 8 GB RAMS</p>	<p>Case studies demonstrate that the state-based strategies try to make the feeders on the trajectory locate the terminal of the whole network to reduce potential loss of load, and in consequence to improve system resilience.</p>	<p>According to the decision processes, constructing the state transition tree under actions in consideration of uncertainties is one critical step to solve the MDP-based model. It is a difficult task to construct the state transition tree of the proposed model in consideration of various actions and the resulting complicated state transitions.</p>

<p>Use of Discrete and Continuous Markov Chains for System Absorbing States</p> <p><u>Tetiana Olekh</u> <u>Viktor Gogunskii</u></p> <p>2018</p>	<p>The paper shows the use of Markov chains and oriented graphs in models of gradation of state of correspondence as the degree of project perfection. In the description of these models, the decomposition of the systems under investigation into certain discrete states is performed and a scheme of transitions between these states is created.</p>	<p>A fundamental matrix for a hypothetically simulated Markov absorption chain is considered, which gives the same prediction for the future regardless of the absolute value of the elapsed time.</p>	<p>A model for assessing the degree of compliance of environmental assessments with quality criteria is presented in the form of an oriented graph. The vertices of the graph correspond to the states of degrees of conformity of ecological estimates to certain criteria, and the arcs to nonzero transition probabilities.</p>	<p>The fundamental matrix gives the same prediction for the future, regardless of the absolute value of the time elapsed from the initial moment. This property of the fundamental matrix illustrates the Markov property of the process, characterizing it as a process without aftereffect</p>	<p>The nature of the probability distribution of the initial states is determined by the conditions of the problem. For example, at an initial moment the system may be in each of the states with equal probability. The appearance of the transition matrix depends entirely on numbering the states.]</p>
---	--	--	--	--	--

TECHNICAL SPECIFICATIONS

The primary benefits of parallel computing are the time saving and the possibility of better results. For example, due to the high computing efficiency, simulation models can be developed with more details for finite element analysis (FEA) and more search iterations can be used for optimisation studies, with the outcome that better simulation results and optimisation results can be achieved.

Parallel computing helps us harness power and make the most out of it. Parallel computing uses multiple computing units to perform multiple tasks simultaneously thus reducing the amount of time taken to perform a task stored in multiple computers with no way or a time consuming way to access them. Therefore we can see that parallel computing is emerging day to day to simplify the tasks and to access information faster and efficiently.

The algorithm demands parallel computing software. On using apriori and markov model we will be able to achieve the required goal of the project. The higher level implementation of the project uses sensors on the railway tracks that collects information about the train details and updates it on the software. This gives the user an idea about how and when an error might occur.

DESIGN

There has been an ever increasing rate of the number of trains being deployed and it has become very hard for us to keep up with this rise in the number of trains and the information regarding these trains. There were many methodologies which were adopted initially to make this process of monitoring trains easier like using sensors or supercomputers alone but later it was realised that using cloud

computing using parallel distributed networks was more optimistic compared to the other methods.

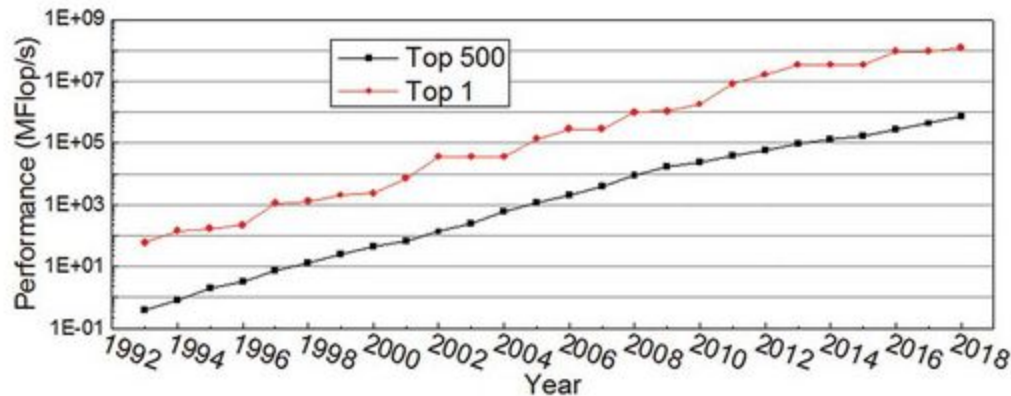


Figure 1. Supercomputing performance trends, data source, where MFlop/s means million floating point operations per second.

The above diagram shows the increase in the performance of monitoring trains over the years and it is to be duly noted that the increase in the performance was mainly due to the increase in parallel distributed networks.

This method uses K-Means algorithm to form clusters and based on these clusters formed we can monitor a train easily. The working of this algorithm in brief is given below

Step 1: Initialization

The first thing that K-Means does is randomly choose examples from the given dataset as initial centroids and that's simply it does not know yet where the center of each cluster is.

Step 2: Cluster Assignment

Then, all the data points that are closest (similar) to a centroid will create a cluster. If we're using the Euclidean distance between data points and every centroid, a straight line is drawn between two centroids, then a perpendicular bisector divides this line into two clusters.

Step 3: Move the centroid

Now, we have new clusters, that need centers. A centroid's new value is going to be the mean of all the examples in a cluster.

We'll keep repeating step 2 and 3 until the centroids stop moving, in other words, K-means algorithm is converged.

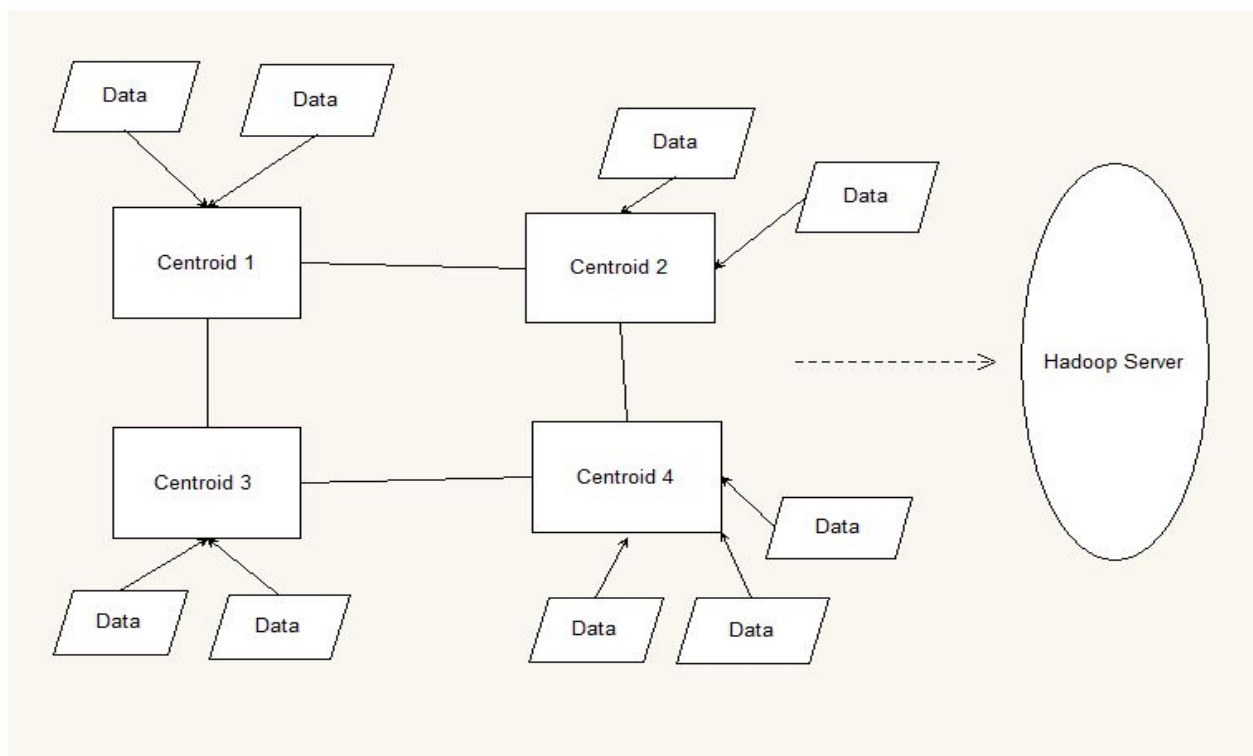


FIG 2: Diagrammatic representation of K-Means algorithm

PROPOSED SYSTEM :

MODULES :

Input Module :

This input module will consist of all the components that help pick up data from the field and supply it to the system for processing. These components can be anything like sensors on the railway tracks, computers accessed by the employees of the railway department, devices placed to check certain aspects like power and energy consumption, a passerby witnessing any mishap or even the cloud where all the data is finally stored.

This data devices are large in number and collect a large amount of data on a daily basis. It becomes really difficult to manage all of this data and classify it. For that purpose, we are using all these techniques of parallel computing for optimal processing.

Data and Signal Processing Module :

This large number of data collected by our input module needs to be processed and sorted out into useful and discardable information. For a long period of time this was being done manually and later with the use of a single system. Such large amounts of data will not even fit into a single system at any place. What we can do is store this data on the cloud and use applications like Hadoop.

In this module, we apply the MapReduce algorithm and the K-Means algorithm. The MapReduce algorithm helps us chunk out large pieces of data. It helps us take in and process that raw data into useful information. We can also use it to segregate data and see if it is useful. The K-Means algorithm is used to segregate this

processed data into different clusters. We can have n number of clusters depending on the need of the system. This is used to make a searching process easier.

The outcomes of using this module is that we can process passenger flow traffic, process high speed train noise data, signal fault diagnosis, railway equipment management and ground penetrating radar data.

Activity Forecast Module :

In this module we implement the Markov model which is based on the Markov property. The Markov property states that the future state of a system depends only on the current state and not previous states, for a randomly changing system. A railway system is not entirely predictable. We have train schedules but still they run late, get cancelled or accidents occur. This uncertainty of the system is what makes us assume that it is a constantly changing one. This module is used to do the scheduling and forecast the possibility of accidents. The predictability of this system is crucial and needs to be optimal for the perfect functioning of railway systems.

Accident Factor Analysis :

Due to unforeseen circumstances, accidents occur. These accidents can lead to loss of life and damage to infrastructure and railway equipment. The railways would do anything in their power to prevent these mishaps. To do so, we need to analyse previous accidents and study them to understand what went wrong and then fix it.

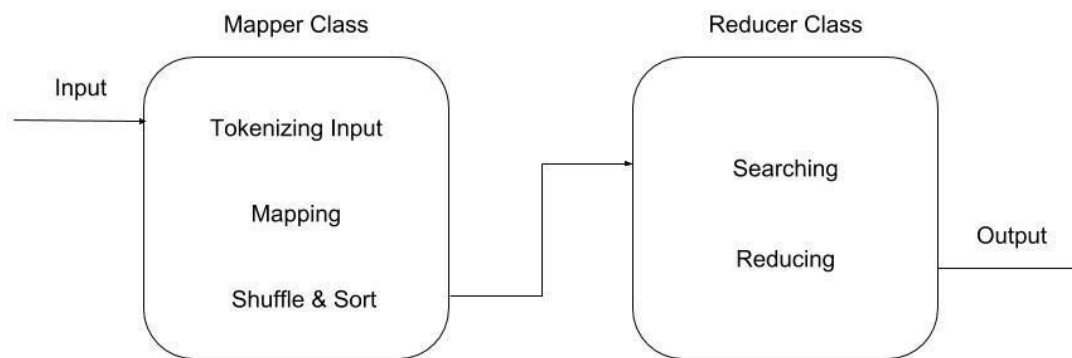
The best tool we can use in analysing large amounts of data is the Apriori algorithm. This algorithm helps us find patterns and previous occurrences in the dataset provided to us. It helps us generate a list of factors and reasoning behind the occurrence of these accidents and so we can use it to make changes to the system to prevent future accidents and irreparable damages.

ALGORITHMS :

The MapReduce algorithm contains two important tasks, namely Map and Reduce.

- The map task is done by means of Mapper Class
- The reduce task is done by means of Reducer Class.

Mapper class takes the input, tokenizes it, maps and sorts it. The output of Mapper class is used as input by Reducer class, which in turn searches matching pairs and reduces them.



MapReduce implements various mathematical algorithms to divide a task into small parts and assign them to multiple systems. In technical terms, MapReduce algorithm helps in sending the Map & Reduce tasks to appropriate servers in a cluster.

These mathematical algorithms may include the following –

•Sorting

Sorting is one of the basic MapReduce algorithms to process and analyse data.

MapReduce implements sorting algorithm to automatically sort the output key-value pairs from the mapper by their keys.

•Searching

Searching plays an important role in MapReduce algorithm. It helps in the combined phase (optional) and in the Reducer phase. Let us try to understand how Searching works with the help of an example.

•Indexing

Normally indexing is used to point to a particular data and its address. It performs batch indexing on the input files for a particular Mapper.

The indexing technique that is normally used in MapReduce is known as inverted index. Search engines like Google and Bing use inverted indexing technique. Let us try to understand how Indexing works with the help of a simple example.

•TF-IDF

TF-IDF is a text processing algorithm which is short for Term Frequency – Inverse Document Frequency. It is one of the common web analysis algorithms. Here, the term 'frequency' refers to the number of times a term appears in a document.

MAPPER CLASS

The Mapper class defines the Map job. Maps input key-value pairs to a set of intermediate key-value pairs. Maps are the individual tasks that transform the input records into intermediate records. The transformed intermediate records need not be of the same type as the input records. A given input pair may map to zero or many output pairs.

REDUCER CLASS

The `Reducer` class defines the Reduce job in MapReduce. It reduces a set of intermediate values that share a key to a smaller set of values.

Reducer implementations can access the Configuration for a job via the `JobContext.getConfiguration()` method. A Reducer has three primary phases – Shuffle, Sort, and Reduce.

- **Shuffle** – The Reducer copies the sorted output from each Mapper using HTTP across the network.
- **Sort** – The framework merge-sorts the Reducer inputs by keys (since different Mappers may have output the same key). The shuffle and sort phases occur simultaneously, i.e., while outputs are being fetched, they are merged.
- **Reduce** – In this phase the `reduce (Object, Iterable, Context)` method is called for each `<key, (collection of values)>` in the sorted inputs.

Generally MapReduce paradigm is based on sending map-reduce programs to computers where the actual data resides.

- During a MapReduce job, Hadoop sends Map and Reduce tasks to appropriate servers in the cluster.
- The framework manages all the details of data-passing like issuing tasks, verifying task completion, and copying data around the cluster between the nodes.
- Most of the computing takes place on the nodes with data on local disks that reduces the network traffic.
- After completing a given task, the cluster collects and reduces the data to form an appropriate result, and sends it back to the Hadoop server.

K-Means

Randomly choose k examples as centroids While true:

Create k clusters by assigning each example to closest centroid Compute k new centroids by averaging examples in each cluster If centroids don't change:

Break

K-means is a fast and efficient method, because the complexity of one iteration is $k \cdot n \cdot d$ where k (number of clusters), n (number of examples), and d (time of computing the Euclidean distance between 2 points).

We try different values of k , we evaluate them and we choose the best k value using the following algorithm:

Best = kMeans(points); For t in range(numTrials):

C = kMeans(points);

if dissimilarity(C) < dissimilarity(best): best = C;

return best

Dissimilarity(C) is the sum of all the variabilities of k clusters

Variability is the sum of all Euclidean distances between the centroid and each example

in the cluster.

We have used Euclidean distance to calculate the distance between every node and the centroid of every cluster to segregate the data points into different clusters.

Euclidean distance function can be defined as:

Where x and y are two vectors

MARKOV

A Markov chain is a random process with the Markov property. A random process or often called stochastic property is a mathematical object defined as a collection of random variables. A Markov chain has either discrete state space (set of possible values of the random variables) or discrete index set (often representing time) - given the fact, many variations for a Markov chain exists. Usually the term "Markov chain" is reserved for a process with a discrete set of times, that is a Discrete Time Markov chain (DTMC).

A Markov chain is represented using a probabilistic automaton (It only sounds complicated!). The changes of state of the system are called transitions. The probabilities associated with various state changes are called transition probabilities. A probabilistic automaton includes the probability of a given transition into the transition function, turning it into a transition matrix.

If the Markov chain has N possible states, the matrix will be an $N \times N$ matrix, such that entry (I, J) is the probability of transitioning from state I to state J .

Additionally, the transition matrix must be a stochastic matrix, a matrix whose entries in each row must add up to exactly 1. Why? Since each row represents its own probability distribution.

So, the model is characterized by a state space, a transition matrix describing the probabilities of particular transitions, and an initial state across the state space, given in the initial distribution.

The Rules is a sequence of pair of strings, usually presented in the form of pattern \rightarrow replacement. Each rule may be either ordinary or terminating.

Given an input string:

Check the Rules in order from top to bottom to see whether any of the patterns can be found in the input string.

If none is found, the algorithm stops.

If one (or more) is found, use the first of them to replace the leftmost occurrence of matched text in the input string with its replacement.

If the rule just applied was a terminating one, the algorithm stops. Go to step 1.

Note that after each rule application the search starts over from the first rule.

RESULTS AND DISCUSSION

```

1 import numpy as np
2 import random as rd
3 import time
4 #The statespace
5 states=["E1","E2","E3"]
6
7 #Possible sequences of events
8 transitionName=["11","12","21","22","23","31","32","33"]
9
10 #Probabilities matrix(transition matrix)
11 transitionMatrix=[[0.2,0.6,0.2],[0.1,0.6,0.3],[0.2,0.7,0.1]]
12 # sum(transitionMatrix[0])+sum(transitionMatrix[1])+sum(transitionMatrix[2])!=3:
13 print("Somewhere something went wrong transition matrix, perhaps?")
14
15 #else:
16 print("All is gonna be okay,let's move on")
17
18 #A function that implements the Markov model to forecast the state/mood
19 def activity_forecast(days,activityToday="E1"):
20     print("Start State:" + activityToday)
21
22     #shall store the sequence of the states taken.so this only has the starting state for n
23     activityList=[activityToday]
24     i=0
25
26     #No documentation available...lity of the activityList
27     prob=1
28     t0=time.time()
29     while i<days:
30         if activityToday=="E1":
31             change=rd.random.choice(transitionName[0],replace=True,p=transitionMatrix[0])
32             if change=="11":
33                 prob=prob*0.2
34                 activityList.append("E1")
35             pass
36         elif change=="12":
37             prob=prob*0.6
38             activityToday="E2"
39             activityList.append("E2")
40         else:
41             prob=prob*0.2
42             activityToday="E3"

```

Usage

Python 3.7.6 (default, Jan 8 2020, 20:23:39) [MSC v.1916 64 bit (AMD64)]
Type "copyright", "credits" or "license" for more information.

IPython 7.12.0 -- An enhanced Interactive Python.

In [1]: runfile('C:/Users/IDEAPAD/PyX/Markov.py', wdir='C:/Users/IDEAPAD/PyX')
All is gonna be okay,let's move on

Enter the Error occurred currently:E1
Start State:E1
Possible Errors:['E1', 'E3', 'E1']
End state after 2days: E1
Probability of the possible sequence of Errors: 0.04000000000000001

Time taken: 0.138229685357666 seconds

In [2]:

```

Python 3.7.6 (default, Jan 8 2020, 20:23:39) [MSC v.1916 64 bit (AMD64)]
Type "copyright", "credits" or "license" for more information.

IPython 7.12.0 -- An enhanced Interactive Python.

In [1]: runfile('C:/Users/IDEAPAD/PyX/Markov.py', wdir='C:/Users/IDEAPAD/PyX')
All is gonna be okay,let's move on

Enter the Error occurred currently:E1
Start State:E1
Possible Errors:['E1', 'E3', 'E1']
End state after 2days: E1
Probability of the possible sequence of Errors: 0.04000000000000001

Time taken: 0.138229685357666 seconds

In [2]: runfile('C:/Users/IDEAPAD/PyX/Markov.py', wdir='C:/Users/IDEAPAD/PyX')
All is gonna be okay,let's move on

Enter the Error occurred currently:E3
Start State:E3
Possible Errors:['E3', 'E3', 'E2']
End state after 2days: E2
Probability of the possible sequence of Errors: 0.06000000000000001

Time taken: 0.43747448921203613 seconds

In [3]:

```

Activate Windows
Go to Settings to activate Windows.

```

1 import time
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import MiniBatchKMeans, KMeans
5 from sklearn.metrics.pairwise import pairwise_distances_argmin
6 from sklearn.datasets.samples_generator import make_blobs
7 # Generate sample data
8 np.random.seed(0)
9 batch_size = 45
10 centers = [[1, 1], [-1, -1], [1, -1], [-1, 1]]
11 n_clusters = len(centers)
12 X, labels_true = make_blobs(n_samples=3000, centers=centers, cluster_std=0.7)
13 # Generate sample data
14 # Compute clustering with KMeans
15 k_means = KMeans(init='k-means++', n_clusters=n_clusters, n_init=10, verbose=1)
16 t0 = time.time()
17 k_means.fit(X)
18 t_batch = time.time() - t0
19 # Compute clustering with MiniBatchKMeans
20 mbk = MiniBatchKMeans(init='k-means++', n_clusters=n_clusters, batch_size=batch_size, n_init=10,
21 t0 = time.time()
22 mbk.fit(X)
23 t_mini_batch = time.time() - t0
24 # Plot result
25 fig = plt.figure(figsize=(10, 10))
26 fig.subplots_adjust(left=0.02, right=0.98, bottom=0.05, top=0.95)
27 colors = ['AF9934', 'AF9934', 'AF9934', 'AF9934']
28 # We want to have the same colors for the same cluster from the
29 # ParallelKMeans and the KMeans algorithms. Let's pair the cluster centers per
30 # cluster
31 k_means_cluster_centers = np.sort(k_means.cluster_centers_, axis=0)
32 mbk_means_cluster_centers = np.sort(mbk.cluster_centers_, axis=0)
33 k_means_labels = pairwise_distances_argmin(X, k_means_cluster_centers)
34 mbk_labels = pairwise_distances_argmin(X, mbk_means_cluster_centers)
35 order = pairwise_distances_argmin(k_means_cluster_centers, mbk_means_cluster_centers)
36 # KMeans
37 ax = fig.add_subplot(1, 3, 1)
38 for k, col in zip(range(n_clusters), colors):
39     ax = fig.add_subplot(1, 3, 1)
40     my_members = k_means_labels == k
41     cluster_center = k_means_cluster_centers[k]
42     ax.plot(X[my_members, 0], X[my_members, 1], 'w', markerfacecolor=col, marker='o')
43     ax.plot(cluster_center[0], cluster_center[1], 'o', markerfacecolor=col, markeredgecolor=col)
44     ax.set_title('KMeans')
45     ax.set_xticks(())
46     ax.set_yticks(())
47     plt.text(-3.5, 1.8, 'train time: %.2fs\ninertia: %f' % (t_batch, k_means.inertia_))
48 # MiniBatchKMeans
49 ax = fig.add_subplot(1, 3, 2)
50 for k, col in zip(range(n_clusters), colors):
51     my_members = mbk_labels == order[k]
52     cluster_center = mbk_means_cluster_centers[order[k]]
53     ax.plot(X[my_members, 0], X[my_members, 1], 'w', markerfacecolor=col, marker='o')
54     ax.plot(cluster_center[0], cluster_center[1], 'o', markerfacecolor=col, markeredgecolor=col)
55     ax.set_title('Parallelized KMeans')
56     ax.set_xticks(())
57     ax.set_yticks(())
58     plt.text(-3.5, 1.8, 'train time: %.2fs\ninertia: %f' % (t_mini_batch, mbk.inertia_))
59 # Parallelized KMeans
60 different = (mbk_labels == order)
61 # Initialize the different array to all False
62 different = (mbk_labels == order)
63 ax = fig.add_subplot(1, 3, 3)
64 for k in range(n_clusters):
65     different = ((k_means_labels == k) != (mbk_labels == order[k]))
66     identic = np.logical_not(different)
67     ax.plot(X[identic, 0], X[identic, 1], 'w', markerfacecolor='b', markers='o')
68     ax.plot(X[different, 0], X[different, 1], 'w', markerfacecolor='r', markers='o')
69     ax.set_title('Difference')
70     ax.set_xticks(())
71     ax.set_yticks(())
72     plt.suptitle('Showing the Execution time for KMeans and ParallelKMeans')
73     plt.show()

```

```

35 mbk_means_cluster_centers = np.sort(mbk.cluster_centers_, axis=0)
36 k_means_labels = pairwise_distances_argmin(X, k_means_cluster_centers)
37 mbk_labels = pairwise_distances_argmin(X, mbk_means_cluster_centers)
38 order = pairwise_distances_argmin(k_means_cluster_centers, mbk_means_cluster_centers)
39 # KMeans
40 ax = fig.add_subplot(1, 3, 1)
41 for k, col in zip(range(n_clusters), colors):
42     my_members = k_means_labels == k
43     cluster_center = k_means_cluster_centers[k]
44     ax.plot(X[my_members, 0], X[my_members, 1], 'w', markerfacecolor=col, marker='o')
45     ax.plot(cluster_center[0], cluster_center[1], 'o', markerfacecolor=col, markeredgecolor=col)
46     ax.set_title('KMeans')
47     ax.set_xticks(())
48     ax.set_yticks(())
49     plt.text(-3.5, 1.8, 'train time: %.2fs\ninertia: %f' % (t_batch, k_means.inertia_))
50 # MiniBatchKMeans
51 ax = fig.add_subplot(1, 3, 2)
52 for k, col in zip(range(n_clusters), colors):
53     my_members = mbk_labels == order[k]
54     cluster_center = mbk_means_cluster_centers[order[k]]
55     ax.plot(X[my_members, 0], X[my_members, 1], 'w', markerfacecolor=col, marker='o')
56     ax.plot(cluster_center[0], cluster_center[1], 'o', markerfacecolor=col, markeredgecolor=col)
57     ax.set_title('Parallelized KMeans')
58     ax.set_xticks(())
59     ax.set_yticks(())
60     plt.text(-3.5, 1.8, 'train time: %.2fs\ninertia: %f' % (t_mini_batch, mbk.inertia_))
61 # Parallelized KMeans
62 different = (mbk_labels == order)
63 # Initialize the different array to all False
64 different = (mbk_labels == order)
65 ax = fig.add_subplot(1, 3, 3)
66 for k in range(n_clusters):
67     different = ((k_means_labels == k) != (mbk_labels == order[k]))
68     identic = np.logical_not(different)
69     ax.plot(X[identic, 0], X[identic, 1], 'w', markerfacecolor='b', markers='o')
70     ax.plot(X[different, 0], X[different, 1], 'w', markerfacecolor='r', markers='o')
71     ax.set_title('Difference')
72     ax.set_xticks(())
73     ax.set_yticks(())
74     plt.suptitle('Showing the Execution time for KMeans and ParallelKMeans')
75     plt.show()

```

```

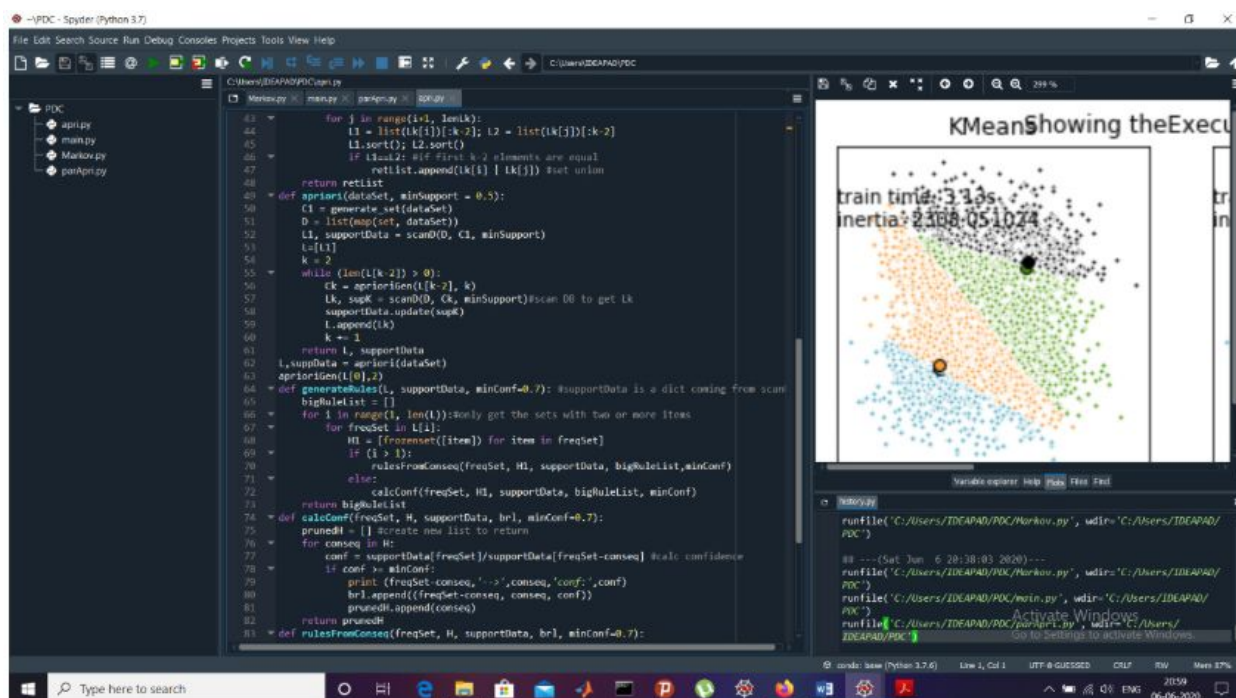
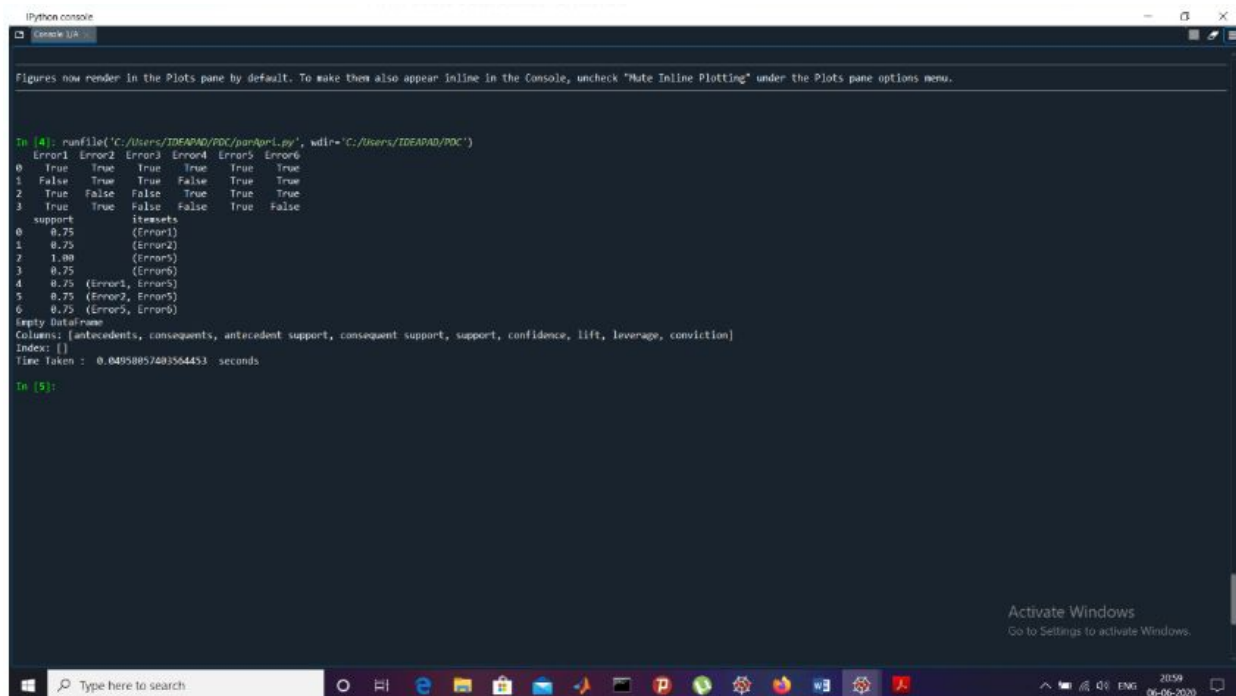
Python console
C:\Users\IA>
done sorting
end inner loop
Iteration 0, inertia 2700.062793945967
start iteration
done sorting
end inner loop
Iteration 1, inertia 2384.820079856231
start iteration
done sorting
end inner loop
Iteration 2, inertia 2322.5203163388843
start iteration
done sorting
end inner loop
Iteration 3, inertia 2318.2855756235717
start iteration
done sorting
end inner loop
Iteration 4, inertia 2388.6182529986716
start iteration
done sorting
end inner loop
Iteration 5, inertia 2388.291278578832
center shift 1.341451e-04 within tolerance 1.478063e-04
init 1/10 with method: k-means++
Inertia for init 1/10: 121.640434
init 2/10 with method: k-means++
Inertia for init 2/10: 123.549484
init 3/10 with method: k-means++
Inertia for init 3/10: 150.976098
init 4/10 with method: k-means++
Inertia for init 4/10: 173.853005
init 5/10 with method: k-means++
Inertia for init 5/10: 125.400547
init 6/10 with method: k-means++
Inertia for init 6/10: 112.546088
init 7/10 with method: k-means++
Inertia for init 7/10: 116.949463
init 8/10 with method: k-means++
Inertia for init 8/10: 118.188847
init 9/10 with method: k-means++
Inertia for init 9/10: 113.416177
init 10/10 with method: k-means++
Inertia for init 10/10: 139.504753
Minibatch iteration 1/6780: mean batch inertia: 0.907285, cwa inertia: 0.907285
Minibatch iteration 2/6780: mean batch inertia: 0.873345, cwa inertia: 0.902961
Minibatch iteration 3/6780: mean batch inertia: 0.944879, cwa inertia: 0.903877
Minibatch iteration 4/6780: mean batch inertia: 0.873345, cwa inertia: 0.902961
Type here to search
Activate Windows
Go to Settings to activate Windows.
2057
06-06-2020

```

```

Python console
C:\Users\IA>
Minibatch iteration 2/6780: mean batch inertia: 0.751354, cwa inertia: 0.902009
Minibatch iteration 3/6780: mean batch inertia: 0.848779, cwa inertia: 0.902961
Minibatch iteration 4/6780: mean batch inertia: 0.873345, cwa inertia: 0.902961
Minibatch iteration 5/6780: mean batch inertia: 0.799538, cwa inertia: 0.898660
Minibatch iteration 6/6780: mean batch inertia: 0.733222, cwa inertia: 0.893713
Minibatch iteration 7/6780: mean batch inertia: 0.712878, cwa inertia: 0.888290
Minibatch iteration 8/6780: mean batch inertia: 0.708779, cwa inertia: 0.882566
Minibatch iteration 9/6780: mean batch inertia: 0.696883, cwa inertia: 0.877071
Minibatch iteration 10/6780: mean batch inertia: 0.664984, cwa inertia: 0.870590
Minibatch iteration 11/6780: mean batch inertia: 0.687137, cwa inertia: 0.865089
Minibatch iteration 12/6780: mean batch inertia: 0.599674, cwa inertia: 0.857129
Minibatch iteration 13/6780: mean batch inertia: 0.924680, cwa inertia: 0.859152
Minibatch iteration 14/6780: mean batch inertia: 0.708779, cwa inertia: 0.857000
Minibatch iteration 15/6780: mean batch inertia: 0.989223, cwa inertia: 0.858653
Minibatch iteration 16/6780: mean batch inertia: 0.873356, cwa inertia: 0.859094
Minibatch iteration 17/6780: mean batch inertia: 0.641411, cwa inertia: 0.852566
Minibatch iteration 18/6780: mean batch inertia: 0.719605, cwa inertia: 0.848578
Minibatch iteration 19/6780: mean batch inertia: 0.707739, cwa inertia: 0.844235
Minibatch iteration 20/6780: mean batch inertia: 0.545386, cwa inertia: 0.835272
Minibatch iteration 21/6780: mean batch inertia: 0.789142, cwa inertia: 0.833889
Minibatch iteration 22/6780: mean batch inertia: 0.795621, cwa inertia: 0.832741
Minibatch iteration 23/6780: mean batch inertia: 0.796889, cwa inertia: 0.831665
Minibatch iteration 24/6780: mean batch inertia: 0.847895, cwa inertia: 0.832152
Minibatch iteration 25/6780: mean batch inertia: 0.770980, cwa inertia: 0.830315
Minibatch iteration 26/6780: mean batch inertia: 0.676289, cwa inertia: 0.825694
Minibatch iteration 27/6780: mean batch inertia: 0.754161, cwa inertia: 0.823548
Minibatch iteration 28/6780: mean batch inertia: 0.891640, cwa inertia: 0.825585
Minibatch iteration 29/6780: mean batch inertia: 0.891254, cwa inertia: 0.822554
Minibatch iteration 30/6780: mean batch inertia: 0.578257, cwa inertia: 0.819818
Minibatch iteration 31/6780: mean batch inertia: 0.747693, cwa inertia: 0.817674
Minibatch iteration 32/6780: mean batch inertia: 0.551889, cwa inertia: 0.809702
Minibatch iteration 33/6780: mean batch inertia: 0.899605, cwa inertia: 0.812397
Minibatch iteration 34/6780: mean batch inertia: 0.678435, cwa inertia: 0.808148
Minibatch iteration 35/6780: mean batch inertia: 0.687932, cwa inertia: 0.804521
Minibatch iteration 36/6780: mean batch inertia: 0.670149, cwa inertia: 0.800494
Minibatch iteration 37/6780: mean batch inertia: 0.744299, cwa inertia: 0.798808
Minibatch iteration 38/6780: mean batch inertia: 0.752676, cwa inertia: 0.797425
Minibatch iteration 39/6780: mean batch inertia: 0.788671, cwa inertia: 0.797162
Minibatch iteration 40/6780: mean batch inertia: 0.787611, cwa inertia: 0.794477
Minibatch iteration 41/6780: mean batch inertia: 0.577928, cwa inertia: 0.787982
Minibatch iteration 42/6780: mean batch inertia: 0.967456, cwa inertia: 0.793365
Minibatch iteration 43/6780: mean batch inertia: 0.717536, cwa inertia: 0.791091
Minibatch iteration 44/6780: mean batch inertia: 0.729917, cwa inertia: 0.789256
Minibatch iteration 45/6780: mean batch inertia: 0.914597, cwa inertia: 0.793813
Minibatch iteration 46/6780: mean batch inertia: 0.621974, cwa inertia: 0.787886
Minibatch iteration 47/6780: mean batch inertia: 0.785907, cwa inertia: 0.787826
Minibatch iteration 48/6780: mean batch inertia: 0.781489, cwa inertia: 0.787634
Minibatch iteration 49/6780: mean batch inertia: 0.666305, cwa inertia: 0.783995
Type here to search
Activate Windows
Go to Settings to activate Windows.
2058
06-06-2020

```

CONCLUSION

The primary benefits of parallel computing are the time saving and the possibility of better results. For example, due to the high computing efficiency, simulation models can be developed with more details for finite element analysis (FEA) and more search iterations can be used for optimisation studies, with the outcome that better simulation results and optimisation results can be achieved.

Parallel computing helps us harness power and make the most out of it. Parallel computing uses multiple computing units to perform multiple tasks simultaneously thus reducing the amount of time taken to perform a task stored in multiple computers with no way or a time consuming way to access them. Therefore we can see that parallel computing is emerging day to day to simplify the tasks and to access information faster and efficiently.

We successfully implemented all the algorithms that we intended to and parallelised them as well. We hope to implement the project on a larger scale some day if given the opportunity to.

REFERENCES

1. L. M. Rodrigues, L. E. Zárate, C. N. Nobre and H. C. Freitas, "Parallel and distributed kmeans to identify the translation initiation site of proteins," *2012 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, Seoul, 2012, pp. 1639-1645, doi: 10.1109/ICSMC.2012.6377972.
2. T. Song, J. Li, X. Han and L. Shao, "A quick Otsu-Kmeans algorithm for the internal pipeline detection," *2017 IEEE International Conference on Mechatronics and Automation (ICMA)*, Takamatsu, 2017, pp. 193-197, doi: 10.1109/ICMA.2017.8015812.
3. Martin L. Puterman. 1994. Markov Decision Processes: Discrete Stochastic Dynamic Programming (1st. ed.). John Wiley & Sons, Inc., USA.
4. T. Olekh and V. Gogunskii, "Use of Discrete and Continuous Markov Chains for System Absorbing States," *2019 IEEE International Conference on Advanced Trends in Information Theory (ATIT)*, Kyiv, Ukraine, 2019, pp. 518-521, doi: 10.1109/ATIT49449.2019.9030457.
5. Potuzak, Tomas. "Optimization of a genetic algorithm for road traffic network division using a distributed/parallel genetic algorithm." In *2016 9th International Conference on Human System Interactions (HSI)*, pp. 21-27. IEEE, 2016.
6. Yabing, J. (2013). Research of an improved apriori algorithm in data mining association rules. *International Journal of Computer and Communication Engineering*, 2(1), 25.
7. Ajeesh, K. P., Ashis, R., Muttath, M. H., Chandran, N. K., Sreedevi, C. R., & Praveen, R. P. (2012, August). A novel fuzzy logic controller based

- automatic caution order system for Indian railways. In *2012 IEEE International Conference on Advanced Communication Control and Computing Technologies (ICACCCT)* (pp. 248-253). IEEE.
8. Patel, Hiral M. "A comparative analysis of MapReduce scheduling algorithms for Hadoop." *Int. J. Innov. Emerg. Res. Eng* 2.2 (2015).
 9. Wu, Q., Spiryagin, M., Cole, C., & McSweeney, T. (2018). Parallel computing in railway research. *International Journal of Rail Transportation*, 1-24.
 10. Sharma, Apekshit, and Chandrakant Sharma. "DISTRIBUTED COMPUTING IN MOBILE ENVIRONMENT."