

Analytic Implementation of Big Data in Cricket using Hadoop Framework

Bhawna¹, Seema Baghla²,

¹Student M. Tech. (Comp. Engg.), ²Assistant Professor (Comp. Engg.)
Yadavindra College of Engineering, Punjabi University Guru Kashi Campus,
Talwandi Sabo, Bathinda, Punjab, India

Email: bhawnagarg232@gmail.com¹, seemabaghla.ycoe1@gmail.com²

Abstract-Big data creates a new way in the cricket to predict the results of a game that which team is going to win the game and also help in to improve the performance of the team. In this paper, data mining is performed over the self-created database using the Apache hadoop framework. With the use of mined data, a team can analyse the strategy followed by the opponent team and predict how they are going to play in upcoming matches. The popularity of a particular player can also be analysed by the number of followers on twitter. Players dependent parameters are considered based upon their batting and bowling data to extract the best insights and to classify the batsman and bowler into different categories. In experimental work, analysis of parameters is done using hadoop and hive for all major cricket teams.

Key-words- Cricket team prediction, Apache hadoop, Mapreduce algorithm, HCatalog tool, Beeswax Hive tool.

1 INTRODUCTION

Big data is a term used to refer the data sets that are too big to be handled using the existing database management tools [11]. Only those technologies that can present vast quantities of structured and unstructured data quickly in the right context are the ones that will provide genuine insights. Not only will these insights enhance our knowledge but also lead to proactive decision making which can redefine the future of all our ventures. Internationally, sports like baseball and basketball have already adopted statistical analysis to enhance their performance. Today cricket has been added to the list of technologically enhanced sports with the adoption of the advanced big data analytics. Big data creates a new way in the field of cricket to predict the outcome of a game by analysing the numerous amounts of data generated in every match of cricket. This data can be generated from the various wearable devices like fitness trackers that present the real-time stat of each player including the speed of batsman hitting the ball and the speed of bowlers throwing the ball, their heart rate and the acceleration. According to the data generated coaches can analyse the performance of each batsman and the bowler.

2 LITERATURE REVIEW

Gandomi and Haider (2015) [1] explained the concept of analytics to gain valuable and valid insights from big data. They reviewed various

analytics techniques for video, text, social media data, audio and predictive analytics. **Che et al. (2013)** [2] discussed the platforms for managing and processing big data as well as efforts anticipated on the big data mining. They pointed out the privacy issue which is collapsed by big data. They proposed and expressed the garbage mining which is a serious issue that is not realized and addressed by other in the era of big data. **Bhosale and Gadekar (2014)** [3] described the challenges faced during the processing of big data like heterogeneity, lack of structure, and privacy and these challenges must be addressed for the efficient and fast processing of big data. They use the hadoop as the solution for the processing of big data. **Dean and Ghemawat (2004)** [4] explained the implementation of MapReduce that processes many terabytes of data on a large cluster of commodity machines. **Khan et al. (2014)** [5] created the data-life cycle using the terminology and technologies of big data. This cycle enhances the efficiency of data management by converting raw data to published data. This paper also discussed the issue of data integrity. **Rein and Memmert (2016)** [6] introduced the big data in the field of elite soccer and also assist in other sports science field as enormous data become an extensive phenomenon. **Shweta and Garg (2013)** [7] used Weka to find out the best association rule. Here author consider three association rule algorithms and compare them: Apriori Association

Rule, Predictive Apriori Association Rule and Tertius Association Rule. Author compares the result of these three algorithms which shows that Apriori Association algorithm performs better than any other algorithm. **Anjali et al. (2015) [8]** explained the use of big data analysis in the field of sports. They use the real time data of cricket tweets to get the accurate results. For the analysis of data, map reduce algorithm is used. They suggest to use the big data analysis in other sports also like football, baseball etc. **Agarwal et al. (2017) [9]** explained the use of statistical modeling approach in the field of cricket to predict the best appropriate team to be lined up for a specific match. **Manikandan and Ravi (2014) [10]** explained that the use of big data tools like mapreduce algorithm over hadoop and hdfs helps the organization in taking best business **Mukherjee and Shaw (2016) [11]** discussed the practical and theoretical challenges that are hindering in the development of big data analytics. **Phaneendra and Reddy (2013) [12]** explained the hadoop architecture consisting HDFS to handle big data systems. The authors also focused on various challenges faced by enterprises while handling big data. **Ahlawat et al. (2016) [13]** defined the ecosystem, developed various models and categorized the elements on the basis of big data. They are able to identify various techniques and technologies for big data analytics. They analysed the areas for the usage of big data analytics. The various data forms are also discussed in this paper.

3 METHODOLOGY

The role played by big data for obtaining the objectives are elaborated with the working of MapReduce algorithm and Hortonworks Sandbox 2.2.0. Initially the pre-processed and the compiled (csv) data files are uploaded in hadoop based Hortonworks platform using file browser tool. These files are then accessed using HCatalog tool to form the tables. To obtain the desired results, scripts and queries are written on Beewax Hive tool. Twenty- four parameters are considered based upon the batting and bowling data of a player. Table 1 shows the description of the batting and bowling attributes of cricket.

Table 1 Cricket attributes with description

Attribute	Description
Player	Name of the players
Matchtype	Specify the type of match like "ODI" or "T20"
Country	Name of the countries

Startingyear	Players started to play the match
Last year	Specify the year till the player play the match
Matches	How many matches are played by the player
Innings	divisions of a match during which one team takes its turn to bat
Notout	Player will be not out if he comes out to bat in an innings
Runs	Runs scored by the player
Strikerate	
i. Batting strikerate	i. how frequently a batsman achieves the primary goal of batting, namely scoring runs
ii. Bowling strikerate	ii. how frequently a bowler achieves the primary goal of bowling, namely taking wickets
Highestscore	High number of scores made by the player in specific matchtype
Centuries	Score of 100 attain by the batsman
Halfcenturies	Score of 50 attain by the batsman
Ducks	When the player is dismissed without facing a ball
Overs	an over consists of six consecutive balls bowled by a <i>single</i> bowler
Maidenovers	A maiden over is one in which no runs are scored
Wickets	Number of wickets taken by the bowler
Economy	It is the average number of runs agreed for each over bowled
Fourwickethaul	Refers to a bowler taking four wickets in a single innings
Fivewickethaul	Refers to a bowler taking wickets in a single innings
Catches	Number of times bowler dismissed the batsman by catching a ball
Stumpings	Method to dismiss the player
Ballsfaced	Number of balls faced by the player
Average	
i. Batting average	i. Average is the total number of runs that they have scored divided by the number of times they have been out.
ii. Bowling average	ii. It is the ratio of runs conceded per wickets taken.

4 RESULTS AND DISCUSSIONS

The process is initiated by collecting the data on the basis of some attributes as shown in table 1 related to the field of cricket, to enhance the winning chances of a team in the matches of “ODI” and “T20”. Figure 1 represents the database of cricketers based upon their batting performances.

Player	Matchtype	Country	Starting	Lastyear	Matches	Innings	Notout	Runs	Highest	Average	Ballsfaced	Strike	centuries	halfcentury	Ducks
1	Abdullah	ODI	2010	2010	2	1	0	3	3	16	18.75	0	0	0	0
2	Afraz Zazi	ODI	2014	2017	17	16	1	264	60	17.6	488	54.09	0	2	3
3	Afraz Zazi	ODI	2013	2013	1	1	0	9	9	9	21	42.85	0	0	0
4	Afraz Zazi	ODI	2010	2015	12	8	4	15	14	6.25	38	65.78	0	0	1
5	Afraz Zazi	ODI	2012	2018	9	3	2	1	1	1	5	20	0	0	1
6	Afraz Zazi	ODI	2009	2009	1	1	0	2	2	2	10	20	0	0	0
7	Afraz Zazi	ODI	2012	2017	31	16	9	23	7	3.28	94	24.46	0	0	1
8	Amir Ham	ODI	2013	2017	31	9	5	40	21	10	38	105.26	0	0	1
9	Amir Ham	ODI	2009	2018	86	81	6	1608	101	21.44	2580	62.32	1	7	5
10	Asghar S	ODI	2010	2018	51	46	4	507	62	21.59	853	106.33	0	3	2
11	Asghar S	ODI	2009	2010	3	0	0	0	0	0	0	0	0	0	0
12	Dawlat Z	ODI	2010	2010	2	1	1	2	2	2	3	66.66	0	0	0
13	Dawlat Z	ODI	2011	2018	72	69	25	500	47	20.83	624	80.12	0	0	5
14	Dawlat Z	ODI	2012	2017	33	16	7	68	13	7.55	42	161.9	0	0	3
15	Fareed A	ODI	2014	2017	5	2	2	1	1	1	6	16.66	0	0	0
16	Fareed A	ODI	2016	2017	7	0	0	0	0	0	0	0	0	0	0
17	Fareed A	ODI	2011	2018	41	34	5	615	82	21.2	851	72.26	0	4	5
18	Gulbedin	ODI	2012	2018	38	31	9	432	56	19.63	357	121	0	1	2
19	Hamid H	ODI	2009	2016	32	17	4	92	17	7.07	171	53.8	0	0	3
20	Hamid H	ODI	2010	2016	22	9	6	50	22	16.66	46	108.69	0	0	0
21	Hamid H	ODI	2012	2012	4	1	1	0	0	0	0	0	0	0	0
22	Hazratul	ODI	2013	2017	17	17	3	322	72	23	556	57.91	0	1	2
23	Hazratul	ODI	2013	2013	1	1	1	1	1	1	1	100	0	0	0
24	Hazratul	ODI	2009	2009	2	1	1	23	23	0	11	209.09	0	0	0
25	Hazratul	ODI	2016	2016	1	1	0	18	18	18	24	75	0	0	0
26	Hazratul	ODI	2017	2018	10	10	1	170	54	18.88	246	69.1	0	2	2
27	Ispahull	ODI	2010	2015	5	3	1	7	7	3.5	27	25.92	0	0	1
28	Ispahull	ODI	2012	2012	4	1	1	0	0	0	3	0	0	0	0
29	Javed Ahr	ODI	2010	2018	39	36	0	920	81	25.55	1161	79.24	0	7	1
30	Javed Ahr	ODI	2012	2017	3	2	0	6	6	3	7	85.71	0	0	1
31	Kacim J	ODI	2017	2017	1	1	0	9	9	9	9	27.27	0	0	0

Fig. 1. Batting database of cricket

Player	Matchtype	Country	Starting	Lastyear	Matches	Innings	Notout	Runs	Highest	Average	Ballsfaced	Strike	centuries	halfcentury	Ducks
1	Abdullah	ODI	2010	2010	2	1	0	3	3	16	18.75	0	0	0	0
2	Afraz Zazi	ODI	2014	2017	17	16	1	264	60	17.6	488	54.09	0	2	3
3	Afraz Zazi	ODI	2013	2013	1	1	0	9	9	9	21	42.85	0	0	0
4	Afraz Zazi	ODI	2010	2015	12	8	4	15	14	6.25	38	65.78	0	0	1
5	Afraz Zazi	ODI	2012	2018	9	3	2	1	1	1	5	20	0	0	1
6	Afraz Zazi	ODI	2009	2009	1	1	0	2	2	2	10	20	0	0	0
7	Afraz Zazi	ODI	2012	2017	31	16	9	23	7	3.28	94	24.46	0	0	1
8	Amir Ham	ODI	2013	2017	31	9	5	40	21	10	38	105.26	0	0	1
9	Amir Ham	ODI	2009	2018	86	81	6	1608	101	21.44	2580	62.32	1	7	5
10	Asghar S	ODI	2010	2018	51	46	4	507	62	21.59	853	106.33	0	3	2
11	Asghar S	ODI	2009	2010	3	0	0	0	0	0	0	0	0	0	0
12	Dawlat Z	ODI	2010	2010	2	1	1	2	2	2	3	66.66	0	0	0
13	Dawlat Z	ODI	2011	2018	72	69	25	500	47	20.83	624	80.12	0	0	5
14	Dawlat Z	ODI	2012	2017	33	16	7	68	13	7.55	42	161.9	0	0	3
15	Fareed A	ODI	2014	2017	5	2	2	1	1	1	6	16.66	0	0	0
16	Fareed A	ODI	2016	2017	7	0	0	0	0	0	0	0	0	0	0
17	Fareed A	ODI	2011	2018	41	34	5	615	82	21.2	851	72.26	0	4	5
18	Gulbedin	ODI	2012	2018	38	31	9	432	56	19.63	357	121	0	1	2
19	Hamid H	ODI	2009	2016	32	17	4	92	17	7.07	171	53.8	0	0	3
20	Hamid H	ODI	2010	2016	22	9	6	50	22	16.66	46	108.69	0	0	0
21	Hamid H	ODI	2012	2012	4	1	1	0	0	0	0	0	0	0	0
22	Hazratul	ODI	2013	2017	17	17	3	322	72	23	556	57.91	0	1	2
23	Hazratul	ODI	2013	2013	1	1	1	1	1	1	1	100	0	0	0
24	Hazratul	ODI	2009	2009	2	1	1	23	23	0	11	209.09	0	0	0
25	Hazratul	ODI	2016	2016	1	1	0	18	18	18	24	75	0	0	0
26	Hazratul	ODI	2017	2018	10	10	1	170	54	18.88	246	69.1	0	2	2
27	Ispahull	ODI	2010	2015	5	3	1	7	7	3.5	27	25.92	0	0	1
28	Ispahull	ODI	2012	2012	4	1	1	0	0	0	3	0	0	0	0
29	Javed Ahr	ODI	2010	2018	39	36	0	920	81	25.55	1161	79.24	0	7	1
30	Javed Ahr	ODI	2012	2017	3	2	0	6	6	3	7	85.71	0	0	1
31	Kacim J	ODI	2017	2017	1	1	0	9	9	9	9	27.27	0	0	0

Fig. 2. Bowling database of cricket

Figure 2 represents the dataset of cricketers based upon their bowling performances.

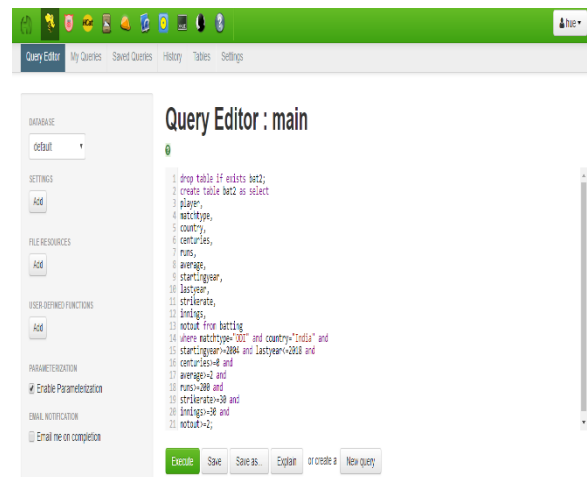


Fig. 3. Script 1 written in hive tool

Figure 3 shows how the script is written to mine the data. The resultant output of script 1 is shown in figure 4 which is in the tabular form.

player	matchtype	country	centuries	runs	average	startingyear	lastyear	strike	innings	notout	
0	AM Rahane	ODI	India	3	2962	55.26	2011	2018	78.63	67	3
1	AT Rayudu	ODI	India	2	1650	50.23	2013	2016	76.28	30	9
2	EI Kumar	ODI	India	0	379	14.57	2012	2018	76.25	40	14
3	IK Pathan	ODI	India	0	1544	23.39	2004	2012	79.54	67	21
4	KD Karthik	ODI	India	0	1496	29.92	2004	2017	72.79	67	17
5	MS Dhoni	ODI	India	9	5793	51.0	2004	2018	67.54	259	77
6	P Kumar	ODI	India	0	292	13.9	2007	2012	88.21	30	12
7	R Ashwin	ODI	India	0	675	16.07	2010	2017	66.96	61	19
8	RA Jadeja	ODI	India	0	1914	31.37	2009	2017	65.29	63	32
9	RG Sharma	ODI	India	17	6594	44.55	2007	2018	66.96	174	26

Fig. 4. Results of script 1 in tabular

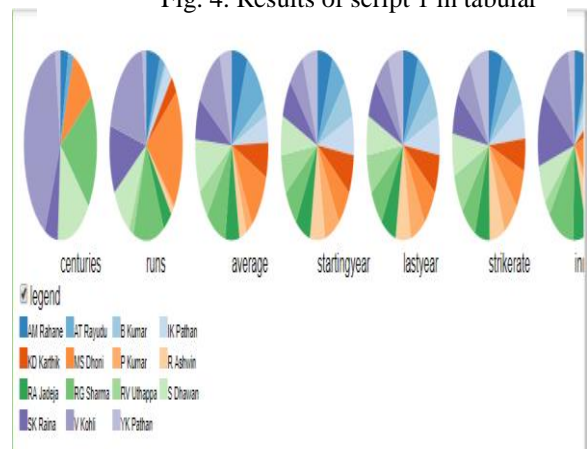


Fig. 5. Results of script 1 in pie chart

Figure 5 represents the result of script 1 in pie chart form.

Figure 6 represents the script 2 written in the beewax hive tool to mine the data from bowling database of cricket.

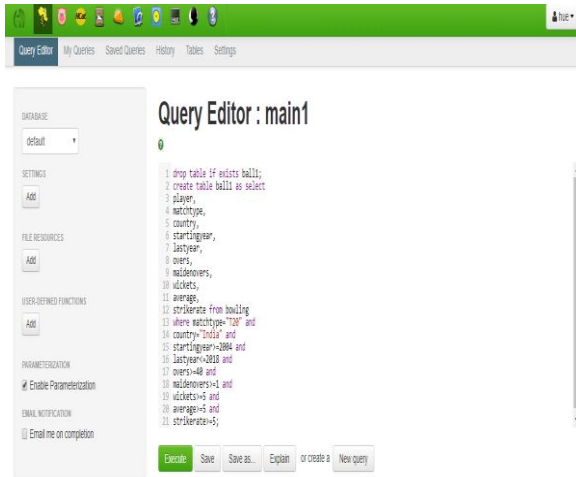


Fig. 6. Script 2 written in hive tool

Figure 7 shows the result of script 2 in the tabular form.

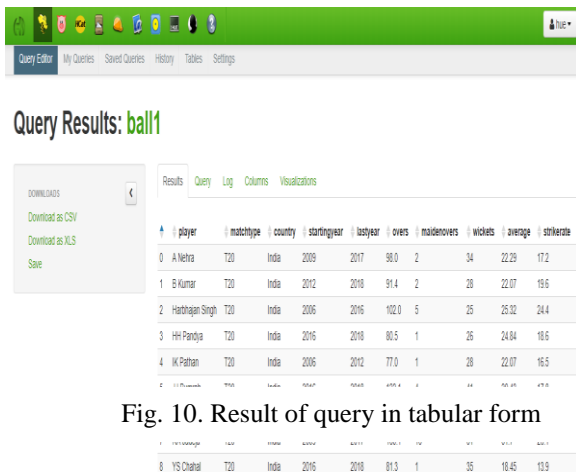


Fig. 10. Result of query in tabular form

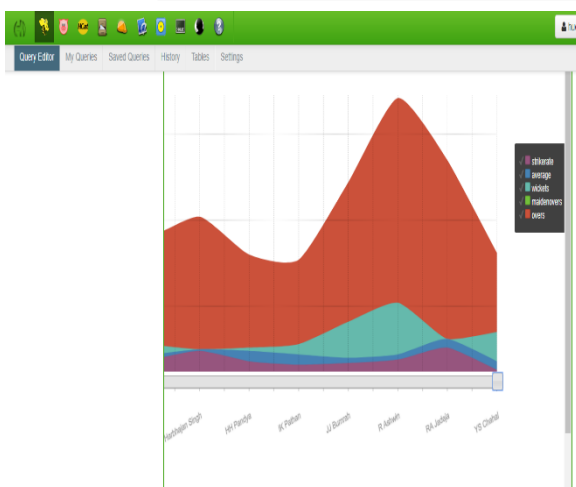


Fig. 8. Results of script 2 in area covered

Figure 9 shows the query written in hive tool to analysis the popularity of players on social media site like twitter based upon follower > 3000000.

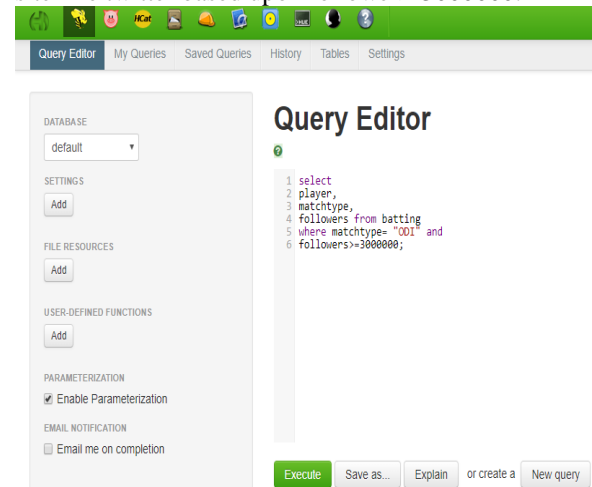


Fig. 9. Query written in hive tool

The screenshot shows the 'Query Results: Unsaved Query' interface. It displays a table with the following data:

player	matchtype	followers
0 AH Rahane	ODI	4000000
1 G Gambhir	ODI	8300000
2 MS Dhoni	ODI	7100000
3 R Ashwin	ODI	8900000
4 R Sharma	ODI	12100000
5 S Chavara	ODI	3000000
6 SK Raina	ODI	14700000
7 SR Tendulkar	ODI	26500000
8 V Kohli	ODI	25700000
9 V Sehwag	ODI	17500000

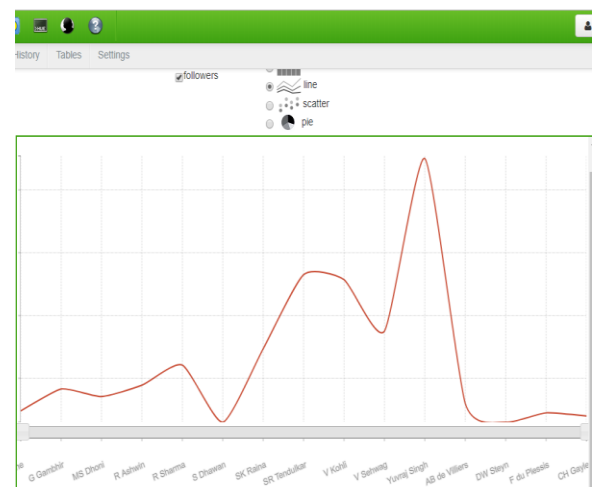


Fig. 11. Result of query in line form and figure 11 shows the same resultant output in line graph form.

5 Conclusions and future scope

Even though the use of Big Data analytics has been a game-changer in the world of sports, it should be for informative purpose only. And for professional cricketers these results can take on a whole new meaning. If technologies behind the data are customized for cricket teams they can pave the way towards a higher level of performances. Similar systems could easily condense years of data and help teams make better decisions when it comes to team selection and game strategy. By combining data driven insights with years of experience and practice, even the most experienced professionals can make better decisions. If a team stops depending on its skills and gut instincts, the whole excitement of the game would vanish. So, there should be a distinctive gap between sports and Big Data analytics.

So, it can be concluded that

- Teams should effectively consume data to derive insights.
- Predicting the entire outcome of the game is difficult, but technology can assist in doing so.
- Human decision making still dominates technology as most teams have similar data but performance varies.

In future, with the help of Machine Learning algorithms can be used to identify complex yet meaningful patterns in the data, which then allows us to predict or classify future instances or events. We can use data from the first innings, such as the number of deliveries bowled, wickets left, runs scored per deliveries faced and partnership for the last wicket, and compare that against total runs scored. Machine learning techniques like SVM, Neural Network, Random Forest can be used to create a model from the historical first innings data, considering the teams playing the match. The same model can be used to predict the second innings which is interrupted by rain. This will give a more accurate prediction than the D/L (Duckworth-Lewis system) method which is the most popular use of mathematics in cricket, as we are using a lot of historical data and all relevant variables.

REFERENCES

- [1] A. Gandomi, M. Haider (2015), "Beyond the hype: Big data concepts, methods, and analytics", *International Journal of Information Management*, 35(2), pp. 137-144.
- [2] D. Che, M. Safran, Z. Peng (2013), "From big data to big data mining: Challenges, issues, and opportunities", *International Conference on database systems for Advanced Applications*, 7827, pp. 1-15.
- [3] H. S. Bhosale, D. P. Gadekar (2014), "A review paper on big data and Hadoop" *International Journal of Scientific and Research Publications*, 4(10), pp. 1-7.
- [4] J. Dean, S. Ghemawat (2004), "MapReduce: Simplified data processing on large clusters", *Communications of the ACM*, 53(1), pp. 1-13.
- [5] N. Khan, I. Yaqoob, I. A. T. Hashem, Z. Inayat, W. K. M. Ali, M. Alam, M. Shiraz, A. Gani (2014), "Big data: Survey, technologies, opportunities, and challenges", *The Scientific World Journal*, pp. 1-18.
- [6] R. Rein, D. Memmert (2016), "Big data and tactical analysis in elite soccer: Future challenges and opportunities for sports science", *SpringerPlus*, 5, pp. 1-13.
- [7] Shweta, K. Garg (2013), "Mining efficient association rules through Apriori algorithm using attributes and comparative analysis of various association rule algorithms", *International Journal of Advance Research in Computer Science and Technology*, 3(6), pp. 306-312.
- [8] S. Anjali, V. Aswini, M. Abirami (2015), "Predictive analysis with cricket tweets using big data", *International Journal of Scientific & Engineering Research*, 6(10), pp. 78-83.
- [9] S. Agarwal, L. Yadav, S. Mehta (2017), "Cricket team prediction with Hadoop: Statistical Modeling Approach", *Procedia Computer Science*, 122, pp. 525-532.
- [10] S.G Manikandan, S. Ravi (2014), "Big data analysis using apache Hadoop", *IEEE International Conference on IT Convergence and Security*, pp. 1-4.
- [11] S. Mukherjee, R. Shaw (2016), "Big data- Concepts, applications, challenges and future scope", *International Journal of Advance Research in Computer and Communication Engineering*, 5(2), pp. 66-73.
- [12] S.V. Phaneendra, E.M. Reddy (2013), "Big data- Solutions for RDBMS problems- A Survey", *International Journal of Advance Research in Computer and Communication Engineering*, 2(9), pp. 3686-3691.
- [13] T. Ahlawat, Dr. R. K. Rambola (2016), "Literature review on big data", *International Journal of Advancement in Engineering Technology, Management & Applied Science*, 3(5), pp. 21-30.