

IPL 2016 Match Prediction

Big Data | 11-11-2016

| SNo | Name              | USN          | Class/Section |
|-----|-------------------|--------------|---------------|
| 1   | Rakshith Gowda VR | 1PI14CS079   | CSE H         |
| 2   | Sunil Pai G       | 1PI14CS115   | CSE H         |
| 3   | Varun L Pai       | 1PI14CS126   | CSE H         |
| 4   | Pradyumna Sripad  | 01FB14ECS152 | CSE C         |

### Introduction

This project is to build a prediction model that can calculate the team scores of any IPL 2016 match, given the batting and bowling order. It involves creating a training model using Player Profiles and Player V Player ball data to calibrate the team score.

#### ALGORITHM/DESIGN

Step I)

- Player V Player data was obtained from cricsheet.org. This data contained Per ball summary of every match in yaml format. We have converted it into CSV files for the ease of reading and manipulating data.
- Player Profiles were scraped from cricketarchive.com. This data consists of two CSV files, one for bowlers and the other for batsmen and the files contain a player's statistics.

```
#webscrape players profiles using BeautifulSoup in python
    #website:cricketarchive.com
    #algorithm
    teams<-list of all the links of each team played in ipl
 6
    for every_team_link in teams:
       players<-urlopen(every team link)
        soup<-BeautifulSoup(players)
9
        allplayers<-soup.find_all href to players
        for each_player_link in allplayers:
            player profile<-urlopen(each player link)
            player_soup<-BeautifulSoup(player_profile)
15
            profile<-find_player_profile_tables()
16
            for every profile table in profile:
                if(check t20 format(every profile table)):
                    player_details<-get_player_details(every_profile_table)
18
19
                    if (batsman_profile(player_details)):
20
                        add entry into t20 batsman profile(player details)
                    else add entry into t20 bowler profile (player details)
21
```

Fig: algorithm to scrape the data.

|    | A                               | В       | С       | D        | E    | F          | G       | Н         | -        | J           | K             | L         |
|----|---------------------------------|---------|---------|----------|------|------------|---------|-----------|----------|-------------|---------------|-----------|
| 1  | Name                            | Matches | Innings | Not_Outs | Runs | High_Score | Average | No_Of_100 | No_Of_50 | Strike_Rate | Catches_Taken | Stumpings |
| 2  | Aakash Shyamlal Chopra          |         | 19      | 1        | 334  | 72*        | 18.55   | 0         | 1        | 91.25       | 4             | NA        |
| 3  | Aaron James Finch               | 162     | 158     | 15       | 4945 | 15         | 34.58   | 2         | 37       | 133.97      | 64            | NA        |
| 4  | Aavishkar Madhav Salvi          | 19      | 2       | 2        | 16   | 10*        | NA      | 0         | 0        | 69.56       | 5             | NA        |
| 5  | Abhimanyu Mithun                | 45      | 22      | 8        | 91   | 12*        | 6.5     | 0         | 0        | 126.38      | 15            | NA        |
| 6  | Abhinav Kumar                   | 2       | 2       | 1        | 56   | 55*        | 56      | 0         | 1        | 186.66      | 0             | NA        |
| 7  | Abhinav Mukund                  | 27      | 26      | 5        | 609  | 67*        | 29      | 0         | 2        | 111.53      | 15            | NA        |
| 8  | Abhinav Mukund                  | 27      | 26      | 5        | 609  | 67*        | 29      | 0         | 2        | 111.53      | 15            | NA        |
| 9  | Abhishek Arunkumar Jhunjhunwala | 38      | 32      | 6        | 472  | 54*        | 18.15   | 0         | 2        | 103.96      | 20            | NA        |
| 10 | Abhishek Mohan Nayar            | 88      | 75      | 18       | 1193 | 7          | 20.92   | 0         | 2        | 122.35      | 21            | .NA       |

Fig2: batsmen profiles snippet.

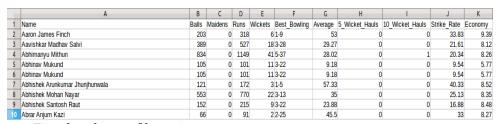


Fig3: bowler profiles snippet.

#### Step II)

- 10 Batsman clusters are created on the Batsman Player Profiles using the Spark MLLib K-Means Clustering function. Clustering is done based on the player's Batting Average and Strike Rate.
- 10 Bowler clusters are created on the Bowler Player Profiles using the Spark MLLib K-Means Clustering function. Clustering is done based on Balls Bowled, Wickets Taken and Economy.

```
val Batvectors = allBatDF.rdd.map(r => Vectors.dense( r.getDouble(1),r.getDouble(2)))

Batvectors.cache()
val BatkMeansModel = KMeans.train(Batvectors, 10, 100)
val Batpredictions = BatrowsRDD.map(r => (r_1, BatkMeansModel.predict(Vectors.dense(r_2,r_3))))
val BatpredDF = Batpredictions.toDF("Batsman_name", "CLUSTER")
```

Fig4: Snippet of the scala code used to cluster the batsmen. Similar code is used for bowlers.

|    | A                       | В        | С             | D       |
|----|-------------------------|----------|---------------|---------|
| 1  | Batsman_name            | avgScore | avgStrikeRate | CLUSTER |
| 2  | David Andrew Miller     | 34.83    | 137.39        | 3       |
| 3  | Kuldeep Yadav           | 22       | 122.22        | 0       |
| 4  | Ravinder Singh Bopara   | 27.48    | 119.59        | 0       |
| 5  | Harmeet Singh Bansal    | 4.4      | 110           | 8       |
| 6  | Rohan Sunil Gavaskar    | 16.14    | 99.12         | 8       |
| 7  | Samuel Badree           | 12.72    | 86.95         | 2       |
| 8  | Musavir Asmat Khote     | 29.4     | 137.38        | 3       |
| 9  | Sanjay Bapusaheb Bangar | 15.24    | 124.5         | 4       |
| 10 | Vijay Hari Zol          | 27.09    | 122.63        | 0       |

Fig5: Snippet of clustered batsmen file.

|    | A                       | В     | С       | D       | E       |
|----|-------------------------|-------|---------|---------|---------|
| 1  | Bowler_name             | Balls | Wickets | Economy | CLUSTER |
| 2  | David Andrew Miller     | 18    | 0       | 10.33   | 3       |
| 3  | Kuldeep Yadav           | 520   | 37      | 7.2     | 1       |
| 4  | Ravinder Singh Bopara   | 3262  | 165     | 7.35    | 4       |
| 5  | Harmeet Singh Bansal    | 594   | 27      | 7.8     | 1       |
| 6  | Rohan Sunil Gavaskar    | 87    | 3       | 7.72    | 3       |
| 7  | Samuel Badree           | 3135  | 146     | 5.61    | 7       |
| 8  | Musavir Asmat Khote     | 132   | 6       | 8.5     | 3       |
| 9  | Sanjay Bapusaheb Bangar | 578   | 31      | 7.35    | 1       |
| 10 | Aiden Craig Blizzard    | 6     | 0       | 10      | 3       |

Fig6: Snippet of clustered bowlers file.

#### Step III)

-Player vs Player Match Data is used to find what run a batsman may score for every ball bowled by a particular bowler. Aggregated Cluster vs Cluster probability is used in case the batsman has not faced the bowler. Algorithm used for achieving this is as given below in the figure.

```
1 #algorithm to generate player vs player probability
 2 #input ->ipl season 3-8 match data, clustered player profiles
 3 #assumption ->if player has not played->probability is zero for every type of score against any bowler
 5 no of zero=empty dict()
6 no of one=empty dict()
 7 no of two=empty dict()
 8 no of three=empty dict()
9 no of four=empty dict()
10 no_of_five=empty_dict()
11 no of six=empty dict()
12 no of seven=empty dict()
13 no of outs=empty dict()
14 type of run=[no of zero,no of one,no of two,no of three,no of four,no of five,no of six,no of seven,no of outs]
16 all matches=list of all match files
17 for every match file in all matches:
18
       for every row in every match file:
19
           batsman_name<-every_row['Batsman']
20
           bowler name<-every row['Bowler']
21
           run<-every_row['runs']
22
            type_of_run[run][batsman_name][bowler_name]+=1
23
24 for every_batsman in player_profiles:
25
        for every bowler in bowler profiles:
            for every type of score in range(0,9):#8 being out.
26
               if(check_if_count_of_runs_is_zero(type_of_run,every_type_of_score,every_batsman,every_bowler)):
27
28
                    names=get_list_of_all_the batsmen_in_the_corresponding_cluster
29
                    for every name in names:
                       total_count+=get_the_count_of_runs_against_the_same_bowler(type_of_run,every_type_of_score,every_name,every_bowler)
30
31
32
                    avg count=ceil(total count/names.length)
                    type_of_run[every_type_of_score][every_batsman][every_bowler]=avg_count
34
            probab=compute_probability_of_each_type_of_run.#tuple of probability values
36
37
            file_write(every_batsman,every_bowler,probab)
```

Fig7:algorithm to generate player vs player statistics using the cluster information.

This algorithm generates a CSV file for every batsman .These CSV files contain the player probabilities against all the bowlers that played in IPL.An example is shown on the next page.

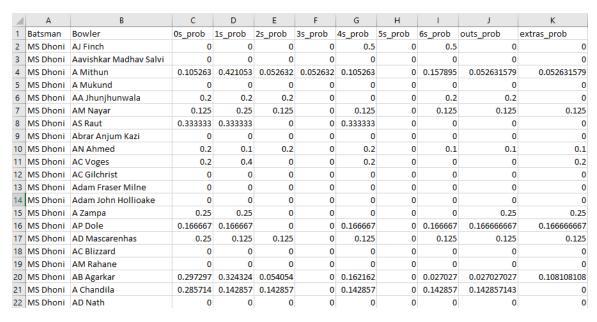
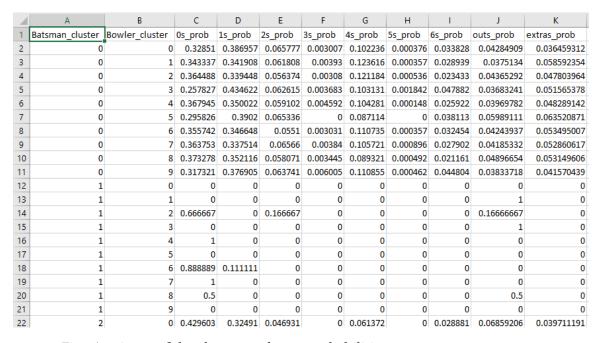


Fig8:Snippet of the player vs player information for MS Dhoni

If the player vs player probability turns out to be zero then it is likely that the batsman has not played against that bowler. In these cases, we use the cluster vs cluster probabilities which looks like the following.



Figo: A snippet of the cluster vs cluster probabilities.

#### Step IV)

-This part involves creating a model to simulate the IPL match and predict the winner based on their scores.

We have made an attempt to create two different models and the difference between the two is in choosing the possible score that for a batsman facing the ball -one model generates random number between o and 1 to pick the score according to the cumulative range probabilities that we compute every time, and the other model picks the score which has the highest probability. The rest of both algorithms work in an identical way.

#### Working:

- -Both of our score predicting algorithm are sensitive to the batting and the bowling order of each side which should be fed to the model manually by creating a CSV. Here we explain the working of the algorithms for team 1 i.e team which bats first. The same applies for team 2.
- -Once both batting and bowling order of the opposite sides is read, we create a list which contains the bowler name with his wicket taking probability (total number of wickets/number of balls he bowled) and also a list of all the batsmen from team 1 with their scoring probabilities and out's probability against the bowlers of team 2. From this we compute the probability of batsman getting not out.
- -If in case the probability of scoring is zero for a batsman that means it is likely that he has not faced the bowler. In this case we use the cluster vs cluster probabilities computed in step 3 which results in a list which looks like < Batsman Name, Bowler Name, oRun, oRun, 1Run, 1Run, 2Run, 2Run, 3Run, 3Run, 4Run, 4Run, 5Run, 5Run, 6Run, 6Run, Not Out Probability, Out Probability> and for the bowlers < Bowler Name, Wicket Probability, Wicket Probability >
- -One can notice that we are duplicating the probabilities for each type of score. The first one is used to get the cumulative probability and the second one is for resetting and incrementing it. Let the second probability be called as increment constant.
- -Now that we have the two lists, we run a loop for 20 overs by setting the Strike and Non-strike batsman to 1 and 2 respectively at the beginning. We use the bowling order specified by the user to choose the bowler for every over and iterate 6 times, once per each ball.

- -From the beginning of the over we check for batsman getting out by comparing if bowler's wicket probability is higher than the probability of the strike batsman not getting out. If out, increment wicket by 1, ball by 1, and change Strike to Wicket+2(new batsman number) and also divide the bowler's wicket probability by a factors say 5. If there is no wicket, find out which run is being scored.
- -To find which run is scored, In the first model we use the cumulative probabilities which we find by summing over the probabilities of every score. These cumulative probabilities are divided by the total cumulative sum hence resulting numbers in range o-1. Hence we get a fixed range for every score. Now we generate a random number between o and 1 and check in which range this number falls. Choose the score accordingly.

In the second model we choose that run with the highest probability being scored.

- Now that we know which run is scored ,reset its probability to the original i.e increment constant and the probability of all the other runs that were not scored is incremented by their respective increment constant except for Not out probability . Also ,we increase bowler's wicket taking probability by its increment constant which ensures that every run not scored gains an increased probability to score in the next turn.
- Add the run scored to total runs. Increment ball by 1. If run scored was odd, swap strike and non-strike batsman. At the very end check if total wickets are less than 10. If so then break out, else increment over count and continue looping until 20 overs get completed.
- -We carry out the same for team 2 .Finally we check the total scores. Whichever team scores highest is declared as the winner.

#### **EXPERIMENTAL RESULTS**

#### Model 1:-

We ran the code for model 1 15 times each on three random matches of IPL 2016 namely Match 1(MI Vs RPS), Match 5(KKR Vs MI) and Match 29(RPS Vs MI).

- -For match 1, out of 15, the model predicted that Rising Pune Giants would win 11 times and Mumbai would win 4 times. Actual result of the match -RPS won.
- -For match 5, out of 15, the model predicted that Mumbai would win 11 times and kkr would win only 4 times . Actual result of the match -MI won. The predicted scores as well as results are depicted in the image below.

-For match 29 ,out of 15,the model predicted that Mumbai would win 12 times and Pune would win only 4 times. Actual result of the match -MI won.

Looking at the frequency of wins per each time in a match ,from above, we can say that our model works very well in predicting the winner of the game and the scores predicted are also comparable with the original scores.

|                      |          | KKR Vs MI   | Match 5 |            |        |
|----------------------|----------|-------------|---------|------------|--------|
|                      |          |             |         |            |        |
|                      | KKR Runs | KKR Wickets | MI Runs | MI Wickets | Winner |
|                      | 184      | 10          | 213     | 6          | MI     |
|                      | 212      | 10          | 243     | 6          | MI     |
|                      | 241      | 10          | 243     | 6          | MI     |
|                      | 242      | 10          | 225     | 5          | KKR    |
|                      | 193      | 10          | 229     | 6          | MI     |
|                      | 197      | 10          | 236     | 6          | MI     |
|                      | 190      | 10          | 239     | 5          | MI     |
|                      | 206      | 10          | 204     | 6          | KKR    |
|                      | 198      | 10          | 203     | 6          | MI     |
|                      | 188      | 10          | 227     | 6          | MI     |
|                      | 177      | 10          | 249     | 6          | MI     |
|                      | 176      | 10          | 220     | 6          | MI     |
|                      | 215      | 10          | 222     | 6          | MI     |
|                      | 221      | 10          | 210     | 6          | KKR    |
|                      | 228      | 10          | 225     | 6          | KKR    |
| Actual scores/Result | 187      | 5           | 188     | 4          | MI     |

Figure 10: Results and scores predicted by model 1 of match KKR vs MI.

Model 2:-

We ran the model 2 code on 8 different matches in IPL 2016. The analysis is depicted in the table.

| Match No. | Match      | Innings 1 | Runs 1 | Wickets 1 | Innings 2 | Runs 2 | Wickets 2 | Winner | actual Innings 1 | actual wickets | actual score | actual wickets | actual winner |
|-----------|------------|-----------|--------|-----------|-----------|--------|-----------|--------|------------------|----------------|--------------|----------------|---------------|
| 1         | MI vs RPS  | MI        | 116    | 9         | RPS       | 117    | 4         | RPS    | 121              | 8              | 126          | 1              | RPS           |
| 2         | DD vs KKR  | DD        | 78     | 10        | KKR       | 124    | 4         | KKR    | 98               | 10             | 99           | 1              | KKR           |
| 3         | KXIP vs GL | KXIP      | 76     | 10        | GL        | 125    | 5         | GL     | 161              | 6              | 162          | 5              | GL            |
| 4         | RCB vs SRH | RCB       | 94     | 10        | SRH       | 134    | 4         | SRH    | 227              | 4              | 182          | 6              | RCB           |
| 5         | KKR vs MI  | KKR       | 100    | 10        | MI        | 115    | 6         | MI     | 187              | 5              | 188          | 4              | MI            |
| 6         | RPS vs GL  | RPS       | 95     | 9         | GL        | 120    | 2         | GL     | 163              | 5              | 164          | 3              | GL            |
| 7         | KXIP vs DD | KXIP      | 96     | 9         | DD        | 157    | 3         | DD     | 111              | 9              | 113          | 2              | DD            |
| 29        | RPS vs MI  | RPS       | 93     | 10        | MI        | 146    | 3         | MI     | 159              | 5              | 161          | 2              | MI            |

From the above table it is evident that Model 2 predicts right 11 times out of the 12 random matches we chose although the scores are comparable only to a certain extent. From this we can conclude that , since model's successes ratio being very high , can be used to predict the outcome of an IPL match in 2016.

#### **FUTURE ENHANCEMENTS**

- Creating a user-friendly GUI.
- Using data from every IPL 2016 match predicted played to further enhance the training model to include changes in form and consecutively predicting IPL 2017.
- Creating a model which takes into account of the type of pitch, weather conditions and other external factors that effect a match to predict scores.
- Make the probabilities sensitive to the current "form" of the player to get accurate results.

#### **REFERENCES**

- Spark MLlibDocumentation :

https://spark.apache.org/mllib/

- BeautifulSoup Documentation:

https://www.crummy.com/software/BeautifulSoup/bs4/doc/

- Player profile data source:

http://cricketarchive.com/

- IPL match data source:

http://cricsheet.org/

- Stack Overflow

http://stackoverflow.com/

## EVALUATIONS (Leave this for the faculty)

| Date | Evaluator | Comments | Score |
|------|-----------|----------|-------|
|      |           |          |       |
|      |           |          |       |
|      |           |          |       |
|      |           |          |       |
|      |           |          |       |
|      |           |          |       |
|      |           |          |       |

# CHECKLIST

| SNo | Item             | Status   |
|-----|------------------|--|
| 1.  | Source code      | complete   |
|     | documented       |  |
| 2   | Source code      | incomplete   |
|     | uploaded to      |  |
|     | CCBD server      |  |
| 3   | Recorded video   | incomplete   |
|     | of demo          |  |
| 4   | Instructions for | complete   |
|     | building and     |  |
|     | running the      |  |
|     | code. Your code  |  |
|     | must be usable   |  |
|     | out of the box.  |  |
| 5   | Dataset used for | <u>Data set formats</u> :                                      |
|     | project          | 1)Match data:  |
|     | uploaded. Please | Batsman name(onstrike),ball number,bowler name,runs scored.    |
|     | include a        | 2) <u>Player profile data</u> :                                |
|     | description of   | a) <u>Batsmen</u> :  |
|     | the dataset      | name,matches,innings,not outs,runs,high                        |
|     | format. This     | score,average,no_of_100,no_of_50,strike rate,catches,stumping. |
|     | includes input   | <u>b)bowler</u> :  |
|     | file format.     | name,balls,maidens,runs,wickets,best bowling,average, 5        |
|     |                  | wickets,10 wickets,strike rate,economy.                        |