

# Code the Game – Stage 1 Report

Team ID: CO1051

## Complete Code and Other Files:

### 1. Python-3 Codes:

```
commentary_extract.py x
1 import re
2 import pandas as pd
3
4 file = open('test.txt','r')
5
6 data = []
7 data = file.readlines(-1)
8 #print(data)
9 pattern = r"([0-9]\.[0-9][0-9]pm)"
10 time = []
11 commentary = []
12
13 #pdb.set_trace()
14 for i in range(len(data)):
15     ext = re.split(pattern,data[i])
16     for j in range(len(ext)):
17         if(re.match(pattern,ext[j])):
18             time.append(ext[j])
19             commentary.append(ext[j+1])
20
21 df = pd.DataFrame(data = {"Time":time,"Commentary":commentary})
22 df.to_csv("./Commentary.csv", sep=',',index=False)
23
```

```
commentary_extract.py  end_of_over.py  Final_PdfT
1 import pandas as pd
2 import re
3
4 file = open('test.txt','r')
5 data = []
6 data = file.readlines(-1)
7 end = []
8 for i in range(21):
9     j=0
10    small_data = data[i]
11    #print(i)
12    #print(data[i])
13    while(j<100):
14        pattern = r"(MCID=\\"+str(j)+"\\")(.?)(</P>)"
15        ext_data = re.search(pattern, small_data)
16        if ext_data is None:
17            break;
18        #print(ext_data.group(2))
19        if('|' in ext_data.group(2)):
20            end.append(ext_data.group(2))
21        j = j+1
22
23    end1=[]
24    j=0
25    i=0
26    while i<(len(end)):
27        if(len(end[i])<=10):
28            end1.append(end[i]+end[i+1])
29            i=i+1
30        else:
31            end1.append(end[i])
32
33        j=j+1
34        i=i+1
35
```

```

commentary_extract.py  end_of_over.py  Final_PdfToCsv.py  Final_with_sentin
28     end1.append(end[i]+end[i+1])
29     i=i+1
30     else:
31         end1.append(end[i])
32
33     j=j+1
34     i=i+1
35
36 # for i in range(len(end1)):
37 #     print(end1[i])
38
39 df = pd.DataFrame(data = {"Over":[] , "Runs Scored":[] ,"Wickets": [], "Total":{}})
40 # print(df.columns)
41
42 for i in range(20):
43     string = ""
44     count = 0
45     str2 = ""
46     x=0
47     string = end1[i]
48     for j in string:
49         # print(j)
50         if(j == '|' and count == 0):
51             df["Over"][i] = (str2)
52             str2 = ""
53             count = count+1
54         elif(j == 'R' and count == 1 ):
55             df["Runs Scored"][i] = (str2)
56             str2 = ""
57             count = count+1
58         elif(j == "s" and count ==2):
59             str2 = ""
60             # count = count+1
61         elif( j == "W" and count ==2):
62             df["Wickets"][i] = str2
63             str2 = ""
64             x = 1
65         elif(j == "|" and count == 2):
66             # if(x==0):

```

```

47     string = end1[i]
48     for j in string:
49         # print(j)
50         if(j == '|' and count == 0):
51             df["Over"][i] = (str2)
52             str2 = ""
53             count = count+1
54         elif(j == 'R' and count == 1 ):
55             df["Runs Scored"][i] = (str2)
56             str2 = ""
57             count = count+1
58         elif(j == "s" and count ==2):
59             str2 = ""
60             # count = count+1
61         elif( j == "W" and count ==2):
62             df["Wickets"][i] = str2
63             str2 = ""
64             x = 1
65         elif(j == "|" and count == 2):
66             # if(x==0):
67             #     df["Wickets"][i] = 0
68             str2 = ""
69             count = count+1
70         elif(j == ":" and count ==3):
71             str2 = ""
72             count = count+1
73         else:
74             str2 = str2+j
75     df["Total"][i] = (str2)
76
77 df1 = pd.DataFrame(data = {"Overs" : df.Over , "Runs": df["Runs Scored"],"Wickets" : df.Wickets , "Total" : df.Total})
78 #dataframe.to_csv("end_of_over_scores.csv", sep = ',' , index = False)
79 df1.Wickets.fillna(0 , inplace = True)
80 df1.to_csv("./EOO_Summary.csv", sep=',',index=False)

```

```

commentary_extract.py x end_of_over.py x Final_PdfToCsv.py Final_with_sentiment.py
1 import re
2 import pandas as pd
3
4 file = open('test.txt','r')
5 data = []
6 data = file.readlines(-1)
7
8 over = []
9 score = []
10 commentry = []
11 string = []
12 flag=1
13 var=0
14 overcount=0
15 commOnprevPage=0
16 test = ""
17
18 for i in range(21):
19     j=0
20     small_data = data[i]
21     while(j<100):
22         pattern = r"(MCID=\""+str(j)+r"\">)(.*?)(</P>)"
23         ext_data = re.search(pattern, small_data)
24         string = []
25         if ext_data is None:
26             break;
27         else:
28             var=0
29             if(len(ext_data.group(2)) >5 and overcount>0):
30                 if(flag==1):
31                     while(1):
32                         if(len(ext_data.group(2)) <=5 and ext_data.group(2)!=""):
33                             if j==0 and commOnprevPage==1 and len(ext_data.group(2))<=5:
34                                 commentry.append([test])
35                                 print("111111=====",test)
36                                 test=""
37                                 commOnprevPage=0
38                             elif commOnprevPage==1 and len(ext_data.group(2))<=5:
39                                 commentry.append([test])
40                                 print("3333333333=====",test)
41                                 test=""
42                                 commOnprevPage=0

```



```

commentary_extract.py x end_of_over.py x Final_PdfToCsv.py Final_with_sentiment.py x score_vs_ball.py
91         None
92         elif("." in ext_data.group(2) and len(ext_data.group(2))<=5):
93             flag=1
94             #j = j-1
95             break
96             j = j+1
97             pattern = r"(MCID="\'+str(j)+r"\">)(.*?)(</P>)"
98             ext_data = re.search(pattern, small_data)
99             if ext_data is None:
100                 break;
101             #commentry.append(ext_data.group(2))
102             if(ext_data is None):
103                 break
104             if(len(ext_data.group(2))<=5 and ext_data.group(2)!=" and ext_data.group(2)!=" "):
105                 if(j==0 or j==1 and commOnprevPage==1):
106                     commOnprevPage=0
107                     commentry.append([test])
108                     # print("44444444===== ",test)
109                     test=""
110                 if("." not in ext_data.group(2)):
111                     if(len(ext_data.group(2))==1 and "D" not in ext_data.group(2) and "O" not in ext_data.group(2)):
112                         score.append(ext_data.group(2))
113                     elif("lb" not in ext_data.group(2) and "w" not in ext_data.group(2) and "nb" not in ext_data.group(2)):
114                         None
115                     else:
116                         score.append(ext_data.group(2))
117                 else:
118                     x = len(ext_data.group(2))
119                     overcount=1
120                     over.append((ext_data.group(2))[:x-1])
121             j=j+1
122
123 print(len(score))
124 print(len((over)))
125 print((len(commentry)))
126 for i in range(len(commentry)):
127     x = "".join(commentry[i])
128     commentry[i] = x
129 df = pd.DataFrame(data={"Overs": over, "ScoreBoard": score,"Commentry": commentry})
130 df.to_csv("./Final_CSV.csv", sep=',', index=False)
131

```

```

commentary_extract.py  x  end_of_over.py  x  Final_PdfToCsv.py  ●  Final_with_sentiment.py  ●
1  import json
2  import pandas as pd
3  from watson_developer_cloud import NaturalLanguageUnderstandingV1
4  import watson_developer_cloud.natural_language_understanding.features.v1 \
5  as Features
6
7  sentiment = []
8  score = []
9  df = pd.read_csv('Final_CSV.csv')
10 over = df.Overs
11 ScoreBoard = df.ScoreBoard
12 Commentry = df.Commentry
13 Score = df.ScoreBoard
14 Overs = df.Overs
15
16 ScoreNew = []
17 CommentryNew = []
18 OversNew = []
19 for i in range(len(df)):
20     if(Score[i] == 'W' or Score[i]=='6' or Score[i]=='4'):
21         ScoreNew.append(Score[i])
22         OversNew.append(Overs[i])
23         CommentryNew.append(Commentry[i])
24     if('\n' in Score[i]):
25         if('w' in Score[i+1] or '4' in Score[i+1] or '6' in Score[i+1] or Score[i] >= '2'):
26             ScoreNew.append(Score[i])
27             OversNew.append(Overs[i])
28             CommentryNew.append(Commentry[i])
29     if(len(ScoreNew)>=50):
30         break
31 if(len(ScoreNew)>=50):
32     exit()
33 #print(len(OversNew) , len(ScoreNew) , len(CommentryNew))
34 natural_language_understanding = NaturalLanguageUnderstandingV1(
35     username="*****",
36     password="*****",
37     version="2017-02-27")
38
39 for i in Commentry:
40     response = natural_language_understanding.analyze(
41         text = i,
42         features=[

```

```

commentary_extract.py  x  end_of_over.py  x  Final_PdfToCsv.py  ●  Final_with_sentiment.py  ●  score_vs_ball.py
34 natural_language_understanding = NaturalLanguageUnderstandingV1(
35     username="*****",
36     password="*****",
37     version="2017-02-27")
38
39 for i in Commentry:
40     response = natural_language_understanding.analyze(
41         text = i,
42         features=[
43             Features.Sentiment(
44             )
45         ]
46     )
47
48     res = response['sentiment']['document']['score']
49     if res < 0:
50         res = -1*res
51     score.append(res)
52 sover = []
53 scomm = []
54 sscore = []
55 ssenti = []
56 df1 = pd.DataFrame(data = {'Over' : Overs, 'Score' : Score, 'Commentry' : Commentry, 'Sentiment' : score})
57 df1 = df1.sort_values(['Sentiment'],ascending = [0])
58 df1.to_csv("./Senti_Score.csv", sep=',',index=False)
59 # print(len(Overs) , len(Score) , len(Commentry))
60 # print(len(df1))
61 # print(df1.head())
62 remain = 50 - len(OversNew)
63 count = 0
64 df1 = pd.read_csv('Senti_Score.csv')
65 # print(len(df1.Over))
66 # print(df1.Over[0] , df1.Score[0] , df1.Commentry[0])
67 for i in range(len(df1)):
68     if count>=remain:
69         break
70     if df1.Commentry[i] not in CommentryNew:
71         # print(df1.Over[i])
72         sover.append(df1.Over[i])
73         scomm.append(df1.Commentry[i])
74         sscore.append(df1.Score[i])
75         ssenti.append(df1.Sentiment[i])

```

```

commentary_extract.py x end_of_over.py x Final_PdfToCsv.py Final_with_sentiment.py score_vs_ball.p
46 )
47
48 res = response['sentiment']['document']['score']
49 if res < 0:
50     res = -1*res
51     score.append(res)
52 sover = []
53 scomm = []
54 sscore = []
55 ssenti = []
56 df1 = pd.DataFrame(data = {'Over' : Overs, 'Score' : Score, 'Commentry' : Commentry, 'Sentiment' : score})
57 df1 = df1.sort_values(['Sentiment'], ascending = [0])
58 df1.to_csv("./Senti_Score.csv", sep=',', index=False)
59 # print(len(Overs) , len(Score) , len(Commentry))
60 # print(len(df1))
61 # print(df1.head())
62 remain = 50 - len(OversNew)
63 count = 0
64 df1 = pd.read_csv('Senti_Score.csv')
65 # print(len(df1.Over))
66 # print(df1.Over[0] , df1.Score[0] , df1.Commentry[0])
67 for i in range(len(df1)):
68     if count >= remain:
69         break
70     if df1.Commentry[i] not in CommentryNew:
71         # print(df1.Over[i])
72         sover.append(df1.Over[i])
73         scomm.append(df1.Commentry[i])
74         sscore.append(df1.Score[i])
75         ssenti.append(df1.Sentiment[i])
76         count = count+1
77 OversNew = OversNew + sover
78 CommentryNew = CommentryNew + scomm
79 ScoreNew = ScoreNew + sscore
80
81 # print(len(OversNew) , len(ScoreNew) , len(CommentryNew))
82 df = pd.DataFrame(data = {'Over' : OversNew, 'Score' : ScoreNew , 'Commentry' : CommentryNew})
83 df = df.sort_values(['Over'])
84 df.to_csv("./Top50Events.csv", sep=',', index=False)
85

```

```

commentary_extract.py x end_of_over.py x Final_PdfToCsv.py Final_with_sentiment.py sco
1 import pandas as pd
2 df = pd.read_csv('Final_CSV.csv')
3
4 x = df['Overs'].values
5 y = df['ScoreBoard'].values
6
7 for i in range(len(y)):
8     y[i] = y[i][0]
9     if (y[i] == 'W'):
10         y[i] = -1
11     else:
12         y[i] = int(y[i])
13
14
15 import matplotlib.pyplot as plt
16 ind = np.arange(min(x)-0.1, max(x)+1, 1)
17 ind = ind.reshape((len(ind), 1))
18 get_ipython().magic('matplotlib inline')
19 plt.figure(figsize=(20,10))
20 plt.plot(x,y,color='blue')
21 plt.xticks(ind)
22 plt.title('Score v/s Ball Analysis', fontsize = 40)
23 plt.xlabel('Balls in Overs 1-20', fontsize = 30)
24 plt.ylabel('Runs', fontsize = 30)
25 plt.show()
26
27 df1 = pd.read_csv('EOO_Summary.csv')
28 x = df1['Overs'].values
29 y = df1['Runs'].values
30
31 ind = np.arange(min(x), max(x)+1, 1)
32 ind = ind.reshape((len(ind), 1))
33
34 width = 0.5
35
36 fig, ax = plt.subplots(figsize=(20,10))
37 rectsl = ax.bar(ind, y, width, color='r')
38 plt.xticks(ind)
39 plt.title('Overs vs Runs Analysis', fontsize = 40)
40 plt.xlabel('Over', fontsize = 30)
41 plt.ylabel('Runs', fontsize = 30)
42

```



## 2). Python-2 Codes:

```
PdfMinerInstaller.py x
1 #####STEPS TO INSTALL PDFMINER PACKAGE FOR PDF TO TXT CONVERTOR#####
2 # 1). First unzip the tar file pdfminer-20140328.tar
3 # 2). Then use "cd pdfminer-20140328" command to go inside the pdfminer-20140328 directory
4 # 3). Then run this given script it will install the PdfMiner Package on to your system
5
6 import os
7 os.system(python setup.py install)
```

```
PdfMinerInstaller.py x PdfToTxt.py x
1 import os
2 os.system(pdf2txt.py -o test.txt -t tag MivsKKR.pdf)
```

### 3). R Codes:

```

1 library("tm")
2 library("SnowballC")
3 library("wordcloud2")
4 library("RColorBrewer")
5 df = read.csv("Final_CSV.csv")
6 df1 = df['Commentry']
7 for (i in df1)
8 {
9   i <- Corpus(VectorSource(i))
10  i <- tm_map(i, removeWords, stopwords("english"))
11  i <- tm_map(i, stemDocument)
12  dtm <- TermDocumentMatrix(i)
13  m <- as.matrix(dtm)
14  v <- sort(rowSums(m), decreasing=TRUE)
15  d <- data.frame(word = names(v), freq=v)
16 }
17 wordcloud2(data = d)
18
19
```

### Aspects taken into account for selecting top events –

We have followed the data driven approach for this natural language processing task. We have used quantitative methods to discover relations. After converting the data to structured form, we have gone through methods like word frequency counting, ranking by means of the Term Frequency – Inverse Document Frequency metric, n-grams, and clustering. We effectively selected the events corresponding to balls which makes the game biased towards change like balls with sixes, fours, wickets, no-balls etc. After that we evaluated the sentiment score of commentary corresponding to each ball and selected the balls which gave us the highest score. We have chosen this strategy so that we may be able to analyze any special events that have occurred with a particular delivery.

We used IBM Watson Natural Language Understanding API for different tasks like categorization of content into five-level classification hierarchy. (However there is no need of such categorization but to finally provide the output to user (who will use the GUI ) with full perfection.) , identifying high level concepts that are not directly referenced in the text, identified the entities mentioned in the content, this includes people, places, events and other types of entities and also done the emotional analysis of each entity. We searched the relevant keywords from the document so that particular analysis can be done on those keywords. Identification of semantic roles is done by parsing the text into subject-object and action form, and identified entities and keywords that are subjects and objects for an action.

## **We are looking forward to work on –**

We have selected the top 50 events till now. Now we have the challenging task to make the best team. For this we will go through the corresponding commentary of top events and we will apply named entity recognition on that to extract the name of players from those events. After that we will focus on the keywords around that entity. (We will use regex parsing for that and will prepare keywords on our own in a separate file which will include “bat”, “bowl”, “sixes”, “fours”, “field”, “keeping”, “stumping”, etc. because ultimately we need the finalized team which have best players from every segment. After that we will evaluate the parameters like average strike rate, bowling economy, catches taken, stumping made, run out, and for bowlers we will consider the parameters like no-ball, Yorker, short pitch, good-length, long pitch. This will give use features corresponding to each player. After that we will calculate the score by allocating weights to each parameter and players with best score will be selected. We will use Watson API for drawing insights from the data and applying Natural Language processing concepts.

For doing this task we will create the GUI using flask in python which have the option to select multiple pdf's file, after that we will convert those pdf into text. Then different parsing methods will be applied that will produce three different csv files and covert unstructured data into structured form. One will contain the commentary data, another will contain the runs/wicket on each ball along with commentary corresponding to that ball, and one will contain the end of over summary. After that we will integrate the second csv file with IBM Watson natural language tool kit API, Watson discovery API to again select the top 50 events and follow the strategy as discussed above, the second and third csv file will be used to plot the visualization and directly analyzer can get aware with the insights of the data.

## **Google Drive Link Containing All the Data:**

<https://drive.google.com/open?id=0B3VdPmUhlac1TENTDJwOWxlQVU>

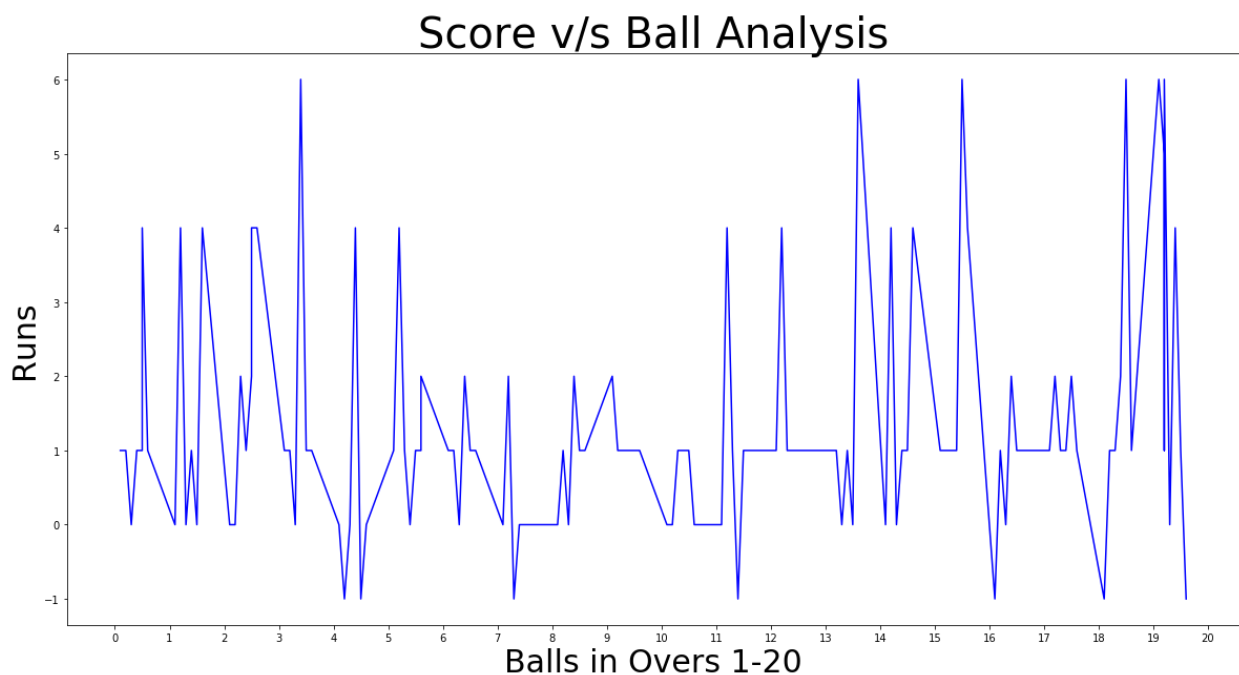
## **Conclusion:**

We concluded from the results that after doing this analysis we are now able to generate the top 50 major events in the match, and now using this data we are able to predict what is status of each player in the given match and with the help of these stats and more data which will be given to us in the upcoming stage we will able to predict who were the star players during the

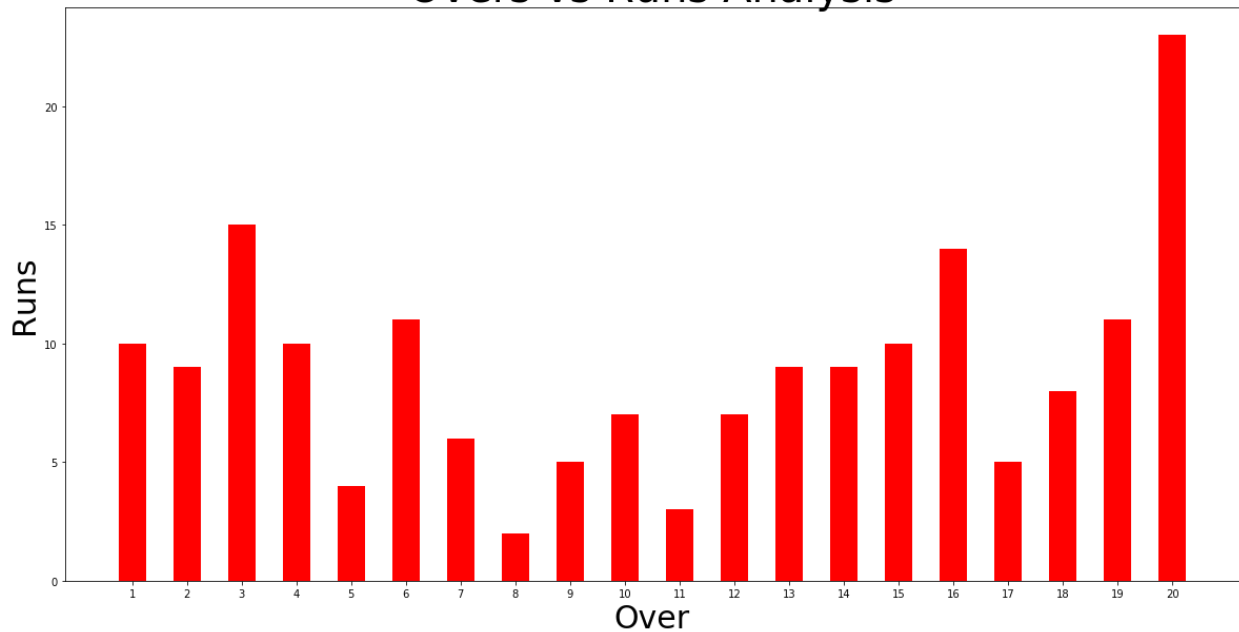
season and this will also help us getting the best all-rounders, best openers, best bowlers, best keeper and best fielders during the season. So all this done we will choose the top 11 players which will form the best team and will surely win any of the match in the next season.

## Alternate Approaches –

This is the alternative approach that we can use for the major event extraction. Unlike data driven approach which we have used to extract the 50 major events from the data, it is often on patterns that express rules representing expert knowledge. It is inherently based on linguistic and lexicographic knowledge, as well as existing human knowledge regarding the contents of the text that is to be processed. In this approach we could have divided the commentary by means of three keywords (score, ball, batsman) that together describe the major events. For example let's take the instance if a batsman score 6 runs in the last ball of the last over it is considered as a major event rather than the ball of any middle over. Also the batsman that is already playing good throughout the series, if that batsman gets out then it will be the major event. So to find these relations we need the expert knowledge and create the rules manually. This approach is known as knowledge driven approach which we could have use for the event extraction.



## Overs vs Runs Analysis



## Word Cloud

