

Analyze the captains performance in ODIs and Test matches in Cricket

This notebook gives an overview of how to read data into Spark Framework and apply some basic RDD operations

```
In [1]: # Spark Context Created.. Let's verify it  
sc
```

```
Out[1]: <pyspark.context.SparkContext at 0x7fade0596278>
```

ODI Performance Analysis ¶

```
In [4]: # Read the captains_ODI.csv file from local file system  
captains_odis = sc.textFile( "file:///home/hadoop/lab/data/captains_ODI.csv" )
```

```
In [5]: # Check the first record  
captains_odis.first()
```

```
Out[5]: 'Ponting  R T,Australia,1995-2012,230,165,51,14,124'
```

```
In [6]: # Display the first 10 records  
captains_odis.take( 10 )
```

```
Out[6]: ['Ponting   R T,Australia,1995-2012,230,165,51,14,124',  
        'Fleming   S P,New Zealand,1994-2007,218,98,106,14,105',  
        'Ranatunga  A,Sri Lanka,1982-1999,193,89,95,9,102',  
        'Dhoni     M S*,India,2004-,186,103,68,15,88',  
        'Border    A R,Australia,1979-1994,178,107,67,4,86',  
        'Azharuddin  M,India,1985-2000,174,89,77,8,96',  
        'Smith     G C,South Africa,2002-2013,149,91,51,7,74',  
        'Ganguly    S C,India,1992-2007,147,76,66,5,74',  
        'Cronje     W J,South Africa,1992-2000,140,99,37,4,74',  
        'Imran Khan,Pakistan,1974-1992,139,75,59,5,70']
```

```
In [7]: # captains_odis should be of RDD type  
type( captains_odis )
```

```
Out[7]: pyspark.rdd.RDD
```

```
In [8]: # Count the number of records  
captains_odis.count()
```

```
Out[8]: 98
```

```
In [9]: captains_odis.cache()
```

```
Out[9]: MapPartitionsRDD[5] at textFile at NativeMethodAccessorImpl.java:-2
```

```
In [10]: # The RDD does not have a header. Let's define a header and then assign it to the RDD  
fields = ("name", "country", "career", "matches", "won", "lost", "ties", "toss" )
```

```
In [11]: from collections import namedtuple
```

```
In [12]: # Lets define a tuple ( a record ) for each line. We need to iterate through each line and convert that into a record  
# The record can be defined as a namedtuple type and called Captain. Let's also link the names for the fields  
Captain = namedtuple( 'Captain', fields )
```

```
In [13]: # Define a function to parse each line and convert them into records  
def parseRecs( line ):  
    fields = line.split(",")  
    return Captain( fields[0], fields[1], fields[2], int( fields[3] ),  
                    int( fields[4] ), int(fields[5]), int(fields[6]), int(fields[7] )  
    )
```

```
In [14]: # Iterate through the data and convert them into records  
captains = captains_odis.map( lambda rec: parseRecs( rec) )
```

```
In [16]: # Now captains refer to all the records. Let's display the first 10 records.  
captains.take( 10 )
```

```
Out[16]: [Captain(name='Ponting  R T', country='Australia', career='1995-2012', matches=230,  
won=165, lost=51, ties=14, toss=124),  
Captain(name='Fleming  S P', country='New Zealand', career='1994-2007', matches=21  
8, won=98, lost=106, ties=14, toss=105),  
Captain(name='Ranatunga  A', country='Sri Lanka', career='1982-1999', matches=193,  
won=89, lost=95, ties=9, toss=102),  
Captain(name='Dhoni  M S*', country='India', career='2004-', matches=186, won=103,  
lost=68, ties=15, toss=88),  
Captain(name='Border  A R', country='Australia', career='1979-1994', matches=178,  
won=107, lost=67, ties=4, toss=86),  
Captain(name='Azharuddin  M', country='India', career='1985-2000', matches=174, wo  
n=89, lost=77, ties=8, toss=96),  
Captain(name='Smith  G C', country='South Africa', career='2002-2013', matches=14  
9, won=91, lost=51, ties=7, toss=74),  
Captain(name='Ganguly  S C', country='India', career='1992-2007', matches=147, won  
=76, lost=66, ties=5, toss=74),  
Captain(name='Cronje  W J', country='South Africa', career='1992-2000', matches=14  
0, won=99, lost=37, ties=4, toss=74),  
Captain(name='Imran Khan', country='Pakistan', career='1974-1992', matches=139, wo  
n=75, lost=59, ties=5, toss=70)]
```

```
In [17]: ## Create an sql context  
#from pyspark.sql import SQLContext  
#sqlContext = SQLContext(sc)
```

```
In [19]: #captains_df = sqlContext.createDataFrame( captains )
```

```
In [20]: #captains_df.show( 10 )
```

```
+-----+-----+-----+-----+---+-----+-----+-----+
|      name|   country|  career|matches|won|lost|ties|toss|
+-----+-----+-----+-----+---+-----+-----+-----+
| Ponting  R T| Australia|1995-2012|   230|165|  51|  14|124|
| Fleming  S P| New Zealand|1994-2007|   218|  98| 106|  14|105|
| Ranatunga  A| Sri Lanka|1982-1999|   193|  89|  95|   9|102|
| Dhoni  M S*|   India|   2004-|   186|103|  68|  15|  88|
| Border  A R| Australia|1979-1994|   178|107|  67|   4|  86|
| Azharuddin  M|   India|1985-2000|   174|  89|  77|   8|  96|
| Smith  G C| South Africa|2002-2013|   149|  91|  51|   7|  74|
| Ganguly  S C|   India|1992-2007|   147|  76|  66|   5|  74|
| Cronje  W J| South Africa|1992-2000|   140|  99|  37|   4|  74|
| Imran Khan| Pakistan|1974-1992|   139|  75|  59|   5|  70|
+-----+-----+-----+-----+---+-----+-----+
only showing top 10 rows
```

```
In [96]: # What is the type of the captains RDD
type( captains )
```

```
Out[96]: pyspark.rdd.PipelinedRDD
```

```
In [97]: # Filter only those captains that have captained for at least 100 ODI matches.
# And then we can compare the statistics of these captains
captains_100 = captains.filter( lambda rec: rec.matches > 100 )
```

```
In [98]: # How many captains have captained their country for more than 100 ODIs..
captains_100.count()
```

```
Out[98]: 16
```

```
In [99]: # Who are these captains
captains.take( 10 )
```

```
Out[99]: [Captain(name='Ponting  R T', country='Australia', career='1995-2012', matches=230,
won=165, lost=51, ties=14, toss=124),
Captain(name='Fleming  S P', country='New Zealand', career='1994-2007', matches=21
8, won=98, lost=106, ties=14, toss=105),
Captain(name='Ranatunga  A', country='Sri Lanka', career='1982-1999', matches=193,
won=89, lost=95, ties=9, toss=102),
Captain(name='Dhoni  M S*', country='India', career='2004-', matches=186, won=103,
lost=68, ties=15, toss=88),
Captain(name='Border  A R', country='Australia', career='1979-1994', matches=178,
won=107, lost=67, ties=4, toss=86),
Captain(name='Azharuddin  M', country='India', career='1985-2000', matches=174, wo
n=89, lost=77, ties=8, toss=96),
Captain(name='Smith  G C', country='South Africa', career='2002-2013', matches=14
9, won=91, lost=51, ties=7, toss=74),
Captain(name='Ganguly  S C', country='India', career='1992-2007', matches=147, won
=76, lost=66, ties=5, toss=74),
Captain(name='Cronje  W J', country='South Africa', career='1992-2000', matches=14
0, won=99, lost=37, ties=4, toss=74),
Captain(name='Imran Khan', country='Pakistan', career='1974-1992', matches=139, wo
n=75, lost=59, ties=5, toss=70)]
```

```
In [100]: # Filtering: Captains who have more wins than losses...
captains_more_wins = captains_100.filter( lambda rec: rec.won > rec.lost )
```

```
In [101]: # Captains with more wins than losses  
captains_more_wins.map( lambda rec: rec.name ).collect()
```

```
Out[101]: ['Ponting   R T',  
           'Dhoni    M S*',  
           'Border   A R',  
           'Azharuddin M',  
           'Smith     G C',  
           'Ganguly   S C',  
           'Cronje    W J',  
           'Imran Khan',  
           'Jayawardene D P M',  
           'Jayasuriya S T',  
           'Wasim Akram',  
           'Waugh     S R',  
           'Richards  I V A']
```

```
In [102]: # Captains with less wins than losses  
captains_more_losts = captains_100.filter( lambda rec: rec.won <= rec.lost )  
captains_more_losts.map( lambda rec: rec.name ).collect()
```

```
Out[102]: ['Fleming   S P', 'Ranatunga  A', 'Lara    B C']
```

```
In [103]: # Which country has played how many matches..  
countries = captains.map( lambda rec: ( rec.country , rec.matches) )
```

```
In [104]: countries.take( 10 )
```

```
Out[104]: [('Australia', 230),  
            ('New Zealand', 218),  
            ('Sri Lanka', 193),  
            ('India', 186),  
            ('Australia', 178),  
            ('India', 174),  
            ('South Africa', 149),  
            ('India', 147),  
            ('South Africa', 140),  
            ('Pakistan', 139)]
```

```
In [105]: # Aggregate by countries  
matches_countries = countries.reduceByKey( lambda a, b: a + b )
```

```
In [106]: matches_countries.take( 20 )
```

```
Out[106]: [('Pakistan', 781),  
            ('South Africa', 463),  
            ('Ireland', 93),  
            ('Australia', 832),  
            ('West Indies', 658),  
            ('Kenya', 114),  
            ('India', 770),  
            ('Bermuda', 31),  
            ('Netherlands', 31),  
            ('Afghanistan', 50),  
            ('England', 554),  
            ('Canada', 27),  
            ('Zimbabwe', 394),  
            ('Sri Lanka', 710),  
            ('Bangladesh', 251),  
            ('New Zealand', 608)]
```



```
In [107]: # Sort the countries by the number of matches they played. Sort by names...(sort by  
key)  
matches_countries.sortByKey().collect()
```

```
Out[107]: [('Afghanistan', 50),  
            ('Australia', 832),  
            ('Bangladesh', 251),  
            ('Bermuda', 31),  
            ('Canada', 27),  
            ('England', 554),  
            ('India', 770),  
            ('Ireland', 93),  
            ('Kenya', 114),  
            ('Netherlands', 31),  
            ('New Zealand', 608),  
            ('Pakistan', 781),  
            ('South Africa', 463),  
            ('Sri Lanka', 710),  
            ('West Indies', 658),  
            ('Zimbabwe', 394)]
```

```
In [108]: # Sort the countries by the number of matches they played by descending order.. by n
          ames (key)
          matches_countries.sortByKey( ascending = False ).collect()
```

```
Out[108]: [('Zimbabwe', 394),
            ('West Indies', 658),
            ('Sri Lanka', 710),
            ('South Africa', 463),
            ('Pakistan', 781),
            ('New Zealand', 608),
            ('Netherlands', 31),
            ('Kenya', 114),
            ('Ireland', 93),
            ('India', 770),
            ('England', 554),
            ('Canada', 27),
            ('Bermuda', 31),
            ('Bangladesh', 251),
            ('Australia', 832),
            ('Afghanistan', 50)]
```

```
In [109]: # Sort by values.. default is ascending...  
matches_countries.sortBy( lambda rec: rec[1] ).collect()
```

```
Out[109]: [('Canada', 27),  
            ('Bermuda', 31),  
            ('Netherlands', 31),  
            ('Afghanistan', 50),  
            ('Ireland', 93),  
            ('Kenya', 114),  
            ('Bangladesh', 251),  
            ('Zimbabwe', 394),  
            ('South Africa', 463),  
            ('England', 554),  
            ('New Zealand', 608),  
            ('West Indies', 658),  
            ('Sri Lanka', 710),  
            ('India', 770),  
            ('Pakistan', 781),  
            ('Australia', 832)]
```

```
In [110]: # Sort by value by descending  
matches_countries.sortBy( lambda rec: rec[1], ascending = False ).collect()
```

```
Out[110]: [('Australia', 832),  
            ('Pakistan', 781),  
            ('India', 770),  
            ('Sri Lanka', 710),  
            ('West Indies', 658),  
            ('New Zealand', 608),  
            ('England', 554),  
            ('South Africa', 463),  
            ('Zimbabwe', 394),  
            ('Bangladesh', 251),  
            ('Kenya', 114),  
            ('Ireland', 93),  
            ('Afghanistan', 50),  
            ('Bermuda', 31),  
            ('Netherlands', 31),  
            ('Canada', 27)]
```

```
In [111]: # Captains by percentage of wins
captains_100_percent_wins = captains_100.map(
    lambda rec: ( rec.name, round( rec.won/rec.matches, 2 ) ) )
# Sort by percentage wins
captains_100_percent_wins.sortBy(
    lambda rec: rec[1], ascending = False ).collect()
```

```
Out[111]: [('Ponting R T', 0.72),
('Cronje W J', 0.71),
('Richards I V A', 0.64),
('Waugh S R', 0.63),
('Smith G C', 0.61),
('Wasim Akram', 0.61),
('Border A R', 0.6),
('Jayasuriya S T', 0.56),
('Dhoni M S*', 0.55),
('Jayawardene D P M', 0.55),
('Imran Khan', 0.54),
('Ganguly S C', 0.52),
('Azharuddin M', 0.51),
('Lara B C', 0.47),
('Ranatunga A', 0.46),
('Fleming S P', 0.45)]
```

```
In [112]: ## Filter countries which have played more than hundred matches
```

```
In [113]: ## Sort counties by the percentage of matches they won
```

```
In [114]: ## Find the top10 lucky captains in terms of TOSS win
```

Test Performance Analysis

```
In [115]: # Read the Captains_Test.csv file
captains_tests = sc.textFile( "file:///home/hadoop/lab/data/captains_Test.csv" )
```

```
In [116]: # Parse the records  
captains_tests_recs = captains_tests.map( lambda rec: parseRecs( rec ) )
```

```
In [117]: # Display the first 10 records  
captains_tests_recs.take( 10 )
```

```
Out[117]: [Captain(name='Smith  G C', country='South Africa', career='2002-2014', matches=10  
9, won=53, lost=29, ties=27, toss=58),  
Captain(name='Border  A R', country='Australia', career='1978-1994', matches=93, w  
on=32, lost=22, ties=38, toss=47),  
Captain(name='Fleming  S P', country='New Zealand', career='1994-2008', matches=8  
0, won=28, lost=27, ties=25, toss=38),  
Captain(name='Ponting  R T', country='Australia', career='1995-2012', matches=77,  
won=48, lost=16, ties=13, toss=37),  
Captain(name='Lloyd  C H', country='West Indies', career='1966-1984', matches=74,  
won=36, lost=12, ties=26, toss=35),  
Captain(name='Dhoni  M S*', country='India', career='2005-', matches=60, won=27, l  
ost=18, ties=15, toss=27),  
Captain(name='Waugh  S R', country='Australia', career='1985-2004', matches=57, wo  
n=41, lost=9, ties=7, toss=31),  
Captain(name='Ranatunga  A', country='Sri Lanka', career='1982-2000', matches=56,  
won=12, lost=19, ties=25, toss=29),  
Captain(name='Atherton  M A', country='England', career='1989-2001', matches=54, w  
on=13, lost=21, ties=20, toss=23),  
Captain(name='Cronje  W J', country='South Africa', career='1992-2000', matches=5  
3, won=27, lost=11, ties=15, toss=22)]
```

```
In [118]: # Filter the captains who have captained for more than 100 tests  
captains_tests_100 = captains_tests_recs.filter( lambda rec: rec.matches > 100 )
```

```
In [119]: ## How many captains?  
captains_tests_100.take( 10 )
```

```
Out[119]: [Captain(name='Smith  G C', country='South Africa', career='2002-2014', matches=10  
9, won=53, lost=29, ties=27, toss=58)]
```

```
In [120]: # Filter the captains who have captained for more than 50 tests  
captains_tests_50 = captains_tests_recs.filter( lambda rec: rec.matches > 50 )
```

```
In [121]: captains_tests_50.take( 10 )
```

```
Out[121]: [Captain(name='Smith  G C', country='South Africa', career='2002-2014', matches=10  
9, won=53, lost=29, ties=27, toss=58),  
Captain(name='Border  A R', country='Australia', career='1978-1994', matches=93, w  
on=32, lost=22, ties=38, toss=47),  
Captain(name='Fleming  S P', country='New Zealand', career='1994-2008', matches=8  
0, won=28, lost=27, ties=25, toss=38),  
Captain(name='Ponting  R T', country='Australia', career='1995-2012', matches=77,  
won=48, lost=16, ties=13, toss=37),  
Captain(name='Lloyd  C H', country='West Indies', career='1966-1984', matches=74,  
won=36, lost=12, ties=26, toss=35),  
Captain(name='Dhoni  M S*', country='India', career='2005-', matches=60, won=27, l  
ost=18, ties=15, toss=27),  
Captain(name='Waugh  S R', country='Australia', career='1985-2004', matches=57, wo  
n=41, lost=9, ties=7, toss=31),  
Captain(name='Ranatunga  A', country='Sri Lanka', career='1982-2000', matches=56,  
won=12, lost=19, ties=25, toss=29),  
Captain(name='Atherton  M A', country='England', career='1989-2001', matches=54, w  
on=13, lost=21, ties=20, toss=23),  
Captain(name='Cronje  W J', country='South Africa', career='1992-2000', matches=5  
3, won=27, lost=11, ties=15, toss=22)]
```

```
In [122]: # Sort the captains by percentage of wins  
captain_top = captains_tests_50.map(  
    lambda rec: ( rec.name, round( rec.won/rec.matches, 2 ) ) ).sortBy(  
    lambda rec: rec[1], ascending = False )
```

```
In [123]: captain_top.collect()
```

```
Out[123]: [('Waugh S R', 0.72),  
            ('Ponting R T', 0.62),  
            ('Cronje W J', 0.51),  
            ('Vaughan M P', 0.51),  
            ('Smith G C', 0.49),  
            ('Lloyd C H', 0.49),  
            ('Dhoni M S*', 0.45),  
            ('Fleming S P', 0.35),  
            ('Border A R', 0.34),  
            ('Atherton M A', 0.24),  
            ('Ranatunga A', 0.21)]
```

Merge both ODI and Test Performance by Various Cricket Captains

```
In [124]: # Lets join both ODI and Test captaincy details. Default is inner join...  
all_time_best_captains = captains_100_percent_wins.join( captain_top )
```

```
In [125]: all_time_best_captains.collect()
```

```
Out[125]: [('Smith G C', (0.61, 0.49)),  
            ('Border A R', (0.6, 0.34)),  
            ('Fleming S P', (0.45, 0.35)),  
            ('Cronje W J', (0.71, 0.51)),  
            ('Ponting R T', (0.72, 0.62)),  
            ('Ranatunga A', (0.46, 0.21)),  
            ('Dhoni M S*', (0.55, 0.45)),  
            ('Waugh S R', (0.63, 0.72))]
```



```
In [126]: ## Best by test match wins  
all_time_best_captains.sortBy( lambda rec: rec[1][1], ascending = False ).collect()
```

```
Out[126]: [('Waugh S R', (0.63, 0.72)),  
           ('Ponting R T', (0.72, 0.62)),  
           ('Cronje W J', (0.71, 0.51)),  
           ('Smith G C', (0.61, 0.49)),  
           ('Dhoni M S*', (0.55, 0.45)),  
           ('Fleming S P', (0.45, 0.35)),  
           ('Border A R', (0.6, 0.34)),  
           ('Ranatunga A', (0.46, 0.21))]
```

```
In [127]: ## Best by ODI match wins  
all_time_best_captains.sortBy( lambda rec: rec[1][0], ascending = False ).collect()
```

```
Out[127]: [('Ponting R T', (0.72, 0.62)),  
           ('Cronje W J', (0.71, 0.51)),  
           ('Waugh S R', (0.63, 0.72)),  
           ('Smith G C', (0.61, 0.49)),  
           ('Border A R', (0.6, 0.34)),  
           ('Dhoni M S*', (0.55, 0.45)),  
           ('Ranatunga A', (0.46, 0.21)),  
           ('Fleming S P', (0.45, 0.35))]
```

```
In [128]: ## Now let's flatten the structure and store the results into a file...  
best_captains = all_time_best_captains.map( lambda rec: ( rec[0], rec[1][0], rec[1]  
[1] ) )  
best_captains.take( 10 )
```

```
Out[128]: [('Smith G C', 0.61, 0.49),  
           ('Border A R', 0.6, 0.34),  
           ('Fleming S P', 0.45, 0.35),  
           ('Cronje W J', 0.71, 0.51),  
           ('Ponting R T', 0.72, 0.62),  
           ('Ranatunga A', 0.46, 0.21),  
           ('Dhoni M S*', 0.55, 0.45),  
           ('Waugh S R', 0.63, 0.72)]
```

```
In [130]: best_captains.saveAsTextFile( "file:///home/hadoop/sparklab/best_captains_0.csv")
```

```
In [131]: ## Consolidate into one partition  
best_captains.repartition( 1 ).saveAsTextFile( "file:///home/hadoop/sparklab/best_ca  
ptains_1.csv")
```