

Exploratory Data Analysis

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [2]: df = pd.read_csv('Country.csv')
```

```
In [3]: df.head()
```

```
Out[3]:
```

	id	name
0	1	Belgium
1	1729	England
2	4769	France
3	7809	Germany
4	10257	Italy

```
In [4]: df.describe()
```

```
Out[4]:
```

	id
count	11.000000
mean	12452.090909
std	8215.308472
min	1.000000
25%	6289.000000
50%	13274.000000
75%	18668.000000
max	24558.000000

```
In [5]: df.shape
```

```
Out[5]: (11, 2)
```

```
In [6]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11 entries, 0 to 10
Data columns (total 2 columns):
#   Column  Non-Null Count  Dtype
---  ---
0   id      11 non-null      int64
1   name    11 non-null      object
dtypes: int64(1), object(1)
memory usage: 304.0+ bytes

#Checking for missing null values
df.isnull().sum()

id      0
name     0
dtype: int64

#Checking for wrong entries like symbols ~,!,#,*,etc.
for col in df.columns:
    print('({}) : {}'.format(col,df[col].unique()))

id : [ 1 1729 4769 7809 10257 13274 15722 17642 19694 21518 24558]
name : ['Belgium' 'England' 'France' 'Germany' 'Italy' 'Netherlands' 'Poland'
'Portugal' 'Scotland' 'Spain' 'Switzerland']

for col in df.columns:
    df[col].replace({'?':'np.nan',inplace=True)
```

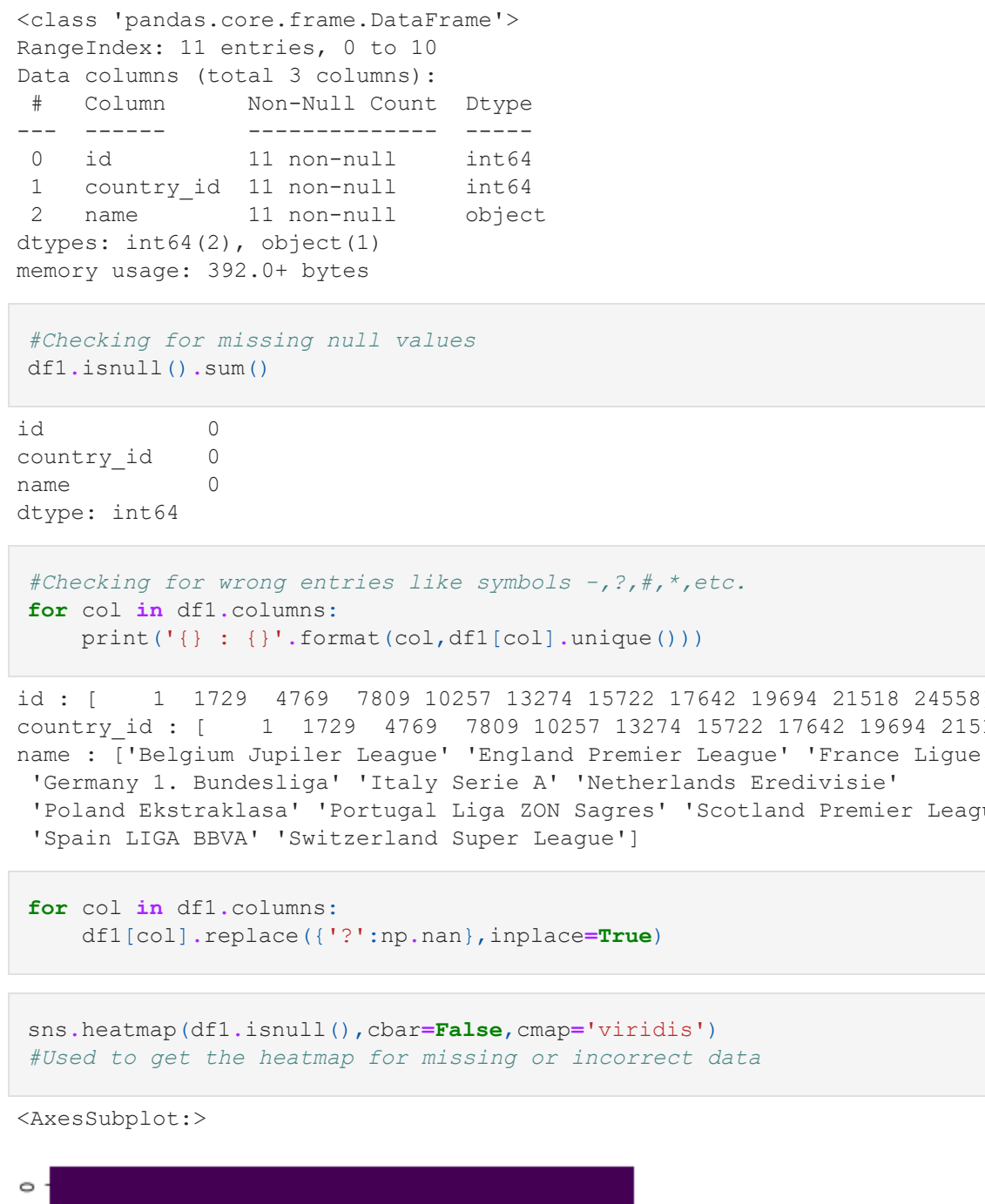
```
In [10]: sns.heatmap(df.isnull(),cbar=False,cmap='viridis')
#Used to get the heatmap for missing or incorrect data
```

```
Out[10]: <AxesSubplot>
```



```
In [11]: plt.figure(figsize=(10,10))
sns.heatmap(df.corr(),cbar=True,annot=True,cmap='Blues')
```

```
Out[11]: <AxesSubplot>
```



```
In [16]: df1 = pd.read_csv('League.csv')
```

```
In [17]: df1.head()
```

```
Out[17]:
```

	id	country_id	name
0	1	1	Belgium Jupiler League
1	1729	1729	England Premier League
2	4769	4769	France Ligue 1
3	7809	7809	Germany 1. Bundesliga
4	10257	10257	Italy Serie A

```
In [18]: df1.describe()
```

```
Out[18]:
```

	id	country_id
count	11.000000	11.000000
mean	12452.090909	12452.090909
std	8215.308472	8215.308472
min	1.000000	1.000000
25%	6289.000000	6289.000000
50%	13274.000000	13274.000000
75%	18668.000000	18668.000000
max	24558.000000	24558.000000

```
In [19]: df1.shape
```

```
Out[19]: (11, 3)
```

```
In [20]: df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11 entries, 0 to 10
Data columns (total 3 columns):
#   Column  Non-Null Count  Dtype
---  ---
0   id      11 non-null      int64
1   country_id  11 non-null      int64
2   name    11 non-null      object
dtypes: int64(2), object(1)
memory usage: 392.0+ bytes

#Checking for missing null values
df1.isnull().sum()

id      0
country_id  0
name     0
dtype: int64

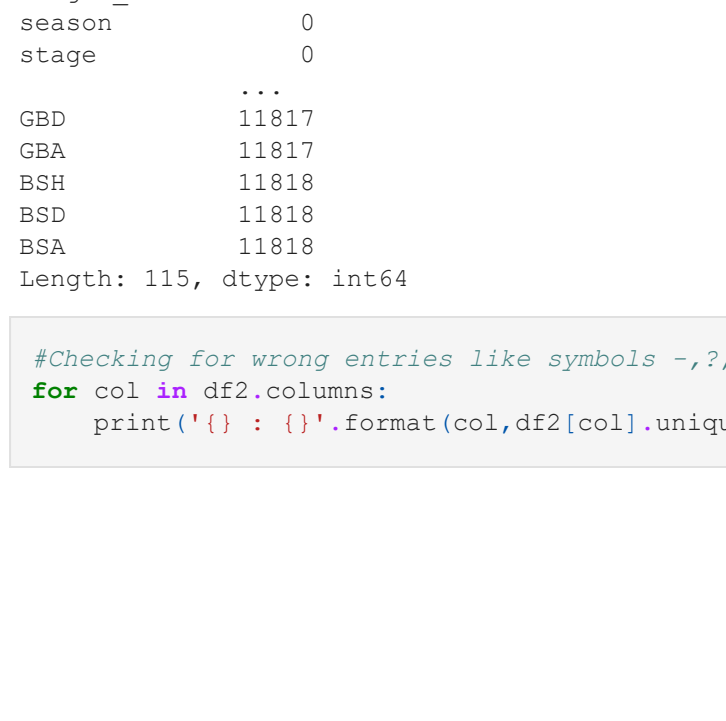
#Checking for wrong entries like symbols ~,!,#,*,etc.
for col in df1.columns:
    print('({}) : {}'.format(col,df1[col].unique()))

id : [ 1 1729 4769 7809 10257 13274 15722 17642 19694 21518 24558]
country_id : [ 1 1729 4769 7809 10257 13274 15722 17642 19694 21518 24558]
name : ['Belgium Jupiler League' 'England Premier League' 'France Ligue 1'
'Germany 1. Bundesliga' 'Italy Serie A' 'Netherlands Eredivisie'
'Poland Ekstraklasa' 'Portugal Liga ZON Sagres' 'Scotland Premier League'
'Spain Liga BBVA' 'Switzerland Super League']

for col in df1.columns:
    df1[col].replace({'?':'np.nan',inplace=True)
```

```
In [24]: sns.heatmap(df1.isnull(),cbar=False,cmap='viridis')
#Used to get the heatmap for missing or incorrect data
```

```
Out[24]: <AxesSubplot>
```



```
In [25]: plt.figure(figsize=(10,10))
sns.heatmap(df1.corr(),cbar=True,annot=True,cmap='Blues')
```

```
Out[25]: <AxesSubplot>
```



```
In [15]: df2 = pd.read_csv('Match.csv')
```

```
In [16]: df2.head()
```

```
Out[16]:
```

	id	country_id	league_id	season	stage	date	match_api_id	home_team_api_id	away_team_api_id	home_team_goal	...	SJA	VCH	VCI
0	1	1	1	2008/2009	1	17-08-2008 00:00	492473	9987	9993	1	...	4.00	1.65	3.4
1	2	1	1	2008/2009	1	16-08-2008 00:00	492474	10000	9994	0	...	3.80	2.00	3.2
2	3	1	1	2008/2009	1	16-08-2008 00:00	492475	9984	8635	0	...	2.50	2.35	3.2
3	4	1	1	2008/2009	1	17-08-2008 00:00	492476	9991	9998	5	...	7.50	1.45	3.7
4	5	1	1	2008/2009	1	16-08-2008 00:00	492477	7947	9985	1	...	1.73	4.50	3.4

5 rows x 115 columns

```
In [28]: df2.describe()
```

```
Out[28]:
```

	id	country_id	league_id	stage	match_api_id	home_team_api_id	away_team_api_id	home_team_goal	away_team_goal
count	25979.000000	25979.000000	25979.000000	25979.000000	2.597900e+04	25979.000000	25979.000000	25979.000000	25979.000000
mean	12990.000000	11738.630317	11738.630317	18.242773	1.195429e+06	9984.371993	9984.475115	1.544594	1.544594
std	7499.635658	7553.936759	7553.936759	10.407354	4.946279e+05	14087.453758	14087.445135	1.297158	1.297158
min	1.000000	1.000000	1.000000	1.000000	4.831290e+05	1601.000000	1601.000000	0.000000	0.000000
25%	6495.500000	4769.000000	4769.000000	9.000000	7.684365e+05	8475.000000	8475.000000	1.000000	0.000000
50%	12990.000000	10257.000000	10257.000000	18.000000	1.147511e+06	8697.000000	8697.000000	1.000000	1.000000
75%	19484.500000	17642.000000	17642.000000	27.000000	1.709852e+06	9925.000000	9925.000000	2.000000	2.000000
max	25979.000000	24558.000000	24558.000000	38.000000	2.216672e+06	274581.000000	274581.000000	10.000000	9.000000

8 rows x 105 columns

```
In [29]: df2.shape
```

```
Out[29]: (25979, 115)
```

```
In [30]: df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 25979 entries, 0 to 25978
Columns: 115 entries, id to BSA
dtypes: float64(96), int64(19), object(10)
memory usage: 22.8+ MB

#Checking for missing null values
df2.isnull().sum()

id      0
country_id  0
league_id  0
season    0
stage     0
...
GBS      11817
GBR      11817
BSR      11818
BSD      11818
BSA      11818
Length: 115, dtype: int64

#Checking for wrong entries like symbols ~,!,#,*,etc.
for col in df2.columns:
    print('({}) : {}'.format(col,df2[col].unique()))
```


[illegible]


```
id ~ | 1 1 2 .. 1456 1457 1458]
team_fifa_api_id ~ | 6 434 171 614 47 1901 650 245 1861 229 1530
111082 11989 147 112513 1 814 2 39 448 240 100409
37 1906 241 1848 896 32 21 675 1887 1889
1734 234 88 1926 1998 160 189 4 59
22 23 111376 1943 1896 190 378 1796 210 1842
647 1961 112409 110364 450 78 1750 110915 5 192
231 1867 110747 1799 110502 635 242 110569 180 181
182 467 110500 1824 468 1746 162 452 10020 2013
1111271 1971 79 31 100626 1915 236 1903 286
897 246 110374 110636 100632 25 111657 144 673 110556 674
1860 110316 1888 110744 100632 110565 111429 110832 322 1805
165 62 28 184 485 80 100634 100879 166 81
100225 1952 11229 44 620 110745 15 29 1822 82
111083 100081 110724 472 1862 46 1738 1739 347 873
111091 1871 95 10018 64 1853 239 65 9 1844
201 2007 217 897 66 169 573 63 10 11
1893 219 68 12 69 1747 70 1900 83 111560
1904 1891 1823 71 100737 1910 13 72 1792 112225
171 477 1564 111540 100741 479 100087 1892 10030 1843
73 50 1754 1914 100741 200 111086 11087 110746 1570
1790 247 15 456 86 480 1793 449 243 457
573 203 379 74 744 1905 1818 1902 546 52
631 874 1819 1837 111974 58 1913 34 324 481
110770 111092 226 17 100646 670 237 459 48 898
108003 110229 10804 680 232 1806 136 1960 260 1715
54 18 1809 294 1908 55 1895 461 110456 462
206 36 483 1909 1887 665 100651 110913 1795 38
109 19 682 1917 1907 1873 175 110 435 1742
900 110749 244 112512 150505]
team_api_id ~ | 9930 8485 8576 8564 10215 10217 8593 9865 8635 8121
6322 108893 9912 158085 8625 10252 8524 8315 9906 9406
8583 10229 8634 9976 9931 8178 9823 9993 10211 9807
6493 9772 8658 8655 8483 8613 9911 9857 8559 8827
9789 9788 4170 8678 10164 9858 8521 8191 7819 8529
7788 8344 208931 8530 9910 9925 9984 9880 8455 8533
8342 7869 2186 9826 8262 8526 9883 9836 8284 9938
8457 8372 9776 9810 10268 8534 8388 8558 7842 8351
8668 4087 10218 8596 8722 6631 8674 9773 9908 9824
10243 10235 8535 8194 8358 9891 9879 9987 10233 9991
8303 6391 9764 8569 6433 8019 8020 7878 9356 9835
8357 9747 9790 8429 9904 9860 9791 10278 8177 10251
8226 8667 8234 8636 8066 1957 9885 8350 8295 8597
10212 8552 8550 8549 8829 9898 10249 8348 9927 274501
8030 8673 8197 6421 8588 9860 9791 10278 8177 10251
8244 9994 8689 10199 9748 9905 9864 8661 8456 10260
10212 8552 8550 8549 8829 9898 10249 8348 9927 274501
9996 9761 10214 8481 9830 9809 8464 10261 9831 9850
6269 8165 8388 8242 2033 8573 8371 1773 6403 8460
5540 9847 10167 8426 6413 4547 9818 8028 8033 8023
8031 2183 8462 8640 10172 8696 8548 10370 9798 8603
8633 8560 8479 8690 9837 9851 7841 10219 8551 9803
9993 8686 8649 1605 9852 9882 7943 7794 10228 10189
9777 8302 10179 8025 9874 8466 8614 9986 9768 9869
9875 10190 9800 8152 8467 9997 9985 10194 8472 10003
10865 10291 9804 8986 8941 10242 8611 8600 9771 10267
9873 10281 9876 10269 10205 8277 7844 10238 9839 8475
9817 8697 8659 8654 10001 8024 8528 8525 10265 8721
8602 7955 9868 10192 8021 8394 8027 10000]
date ~ ['2010-02-22 00:00:00' *2014-09-19 00:00:00' *2015-09-10 00:00:00'
*2011-02-22 00:00:00' *2012-02-22 00:00:00' *2013-09-20 00:00:00']
buildUpPlaySpeed ~ | 60 52 47 70 58 62 59 65 45 48 30 53 38 56 40 31 35 55 42 46 50 23 41 39
69 66 75 25 67 63 64 57 68 43 24 36 61 73 37 51 44 49 71 74 76 54 32 80
34 72 29 78 33 26 28 20 77]
buildUpPlaySpeedClass ~ | ['balanced'] *fast' *slow'
buildUpPlayGribbling ~ | nan 48. 41. 64. 57. 70. 53. 47. 40. 43. 46. 61. 49. 66. 51. 32. 37. 45.
52. 50. 38. 55. 35. 63. 30. 29. 34. 24. 39. 31. 60. 44. 36. 56. 54. 33.
59. 58. 42. 69. 62. 67. 65. 77. 26. 68. 71. 26. 71. 74.]
buildUpPlayGribblingClass ~ | ['little' 'normal' 'lots']
buildUpPlayPassing ~ | 50 56 54 70 52 62 45 53 40 30 44 35 47 55 66 38 33 39 65 41 51 37 48 58
32 29 26 59 72 36 69 49 31 46 34 25 57 22 28 60 79 75 64 73 61 49 27 74
67 42 68 77 71 63 24 20 23 80]
buildUpPlayPassingClass ~ | ['mixed' 'long' 'short']
buildUpPlayPositioningClass ~ | ['organised' *free form']
chanceCreationPassing ~ | 60 54 70 53 45 40 56 51 65 48 55 50 66 30 67 39 58 57 68 35 52 33 49 41
28 44 63 38 36 61 47 59 37 77 34 21 43 69 32 42 72 46 62 71 64 73 31 80
76 291
chanceCreationPassingClass ~ | ['normal' *risky' *safe']
chanceCreationCrossing ~ | 65 63 70 48 50 68 72 35 34 38 45 60 20 53 36 57 51 67 55 52 64 54 62 40
33 59 44 30 47 49 78 73 59 74 56 25 24 31 66 46 37 73 61 58 43 26 42 77
39 41 76 27 80 23 75 32]
chanceCreationCrossingClass ~ | ['normal' *lots' *little']
chanceCreationShooting ~ | 35 64 70 52 57 63 50 75 69 60 65 66 44 38 67 46 39 30 40 47 48 35 43 42
56 34 51 59 37 36 79 54 49 72 53 68 61 22 41 59 45 76 73 62 80 78 33 32
71 28 31 23 77 24 27 74 29]
chanceCreationShootingClass ~ | ['normal' *lots' *little']
chanceCreationPositioningClass ~ | ['organised' *free form']
defencePressure ~ | 50 47 60 40 42 41 49 30 43 38 48 58 52 39 46 57 65 68 64 70 37 53 44 25
36 59 51 43 55 54 59 33 34 67 66 61 72 32 23 27 23 24 63 31 29 56 62]
defencePressureClass ~ | ['medium' *deep' *high']
defenceAggression ~ | 55 44 70 47 40 42 45 35 50 49 57 30 38 62 58 65 53 68 43 48 32 54 37 39
41 66 63 60 34 61 51 46 72 52 59 56 33 69 71 64 27 28 36 24 31 29]
defenceAggressionClass ~ | ['press' *double' *contain']
defenceTeamWidth ~ | 45 54 70 52 60 63 30 50 53 49 65 61 62 64 59 67 51 35 36 58 57 48 37 55
68 42 56 41 40 66 47 38 46 43 39 44 33 69 32 34 29 31 73]
defenceTeamWidthClass ~ | ['normal' *wide' *narrow']
defenceDefenderLineClass ~ | ['cover' *offside trap']
```

```
In [123]: for col in df6.columns:
df6[col].replace({'':np.nan},inplace=True)

In [124]: sns.heatmap(df6.isnull(),cbar=False,cmap=‘viridis’)
#Used to get the heatmap for missing or incorrect data

Out[124]: <AxesSubplot>
```



```
In [125]: #Once again, the way to deal with this amount of missing data is tuple deletion to be done later
plt.figure(figsize=(10,10))
sns.heatmap(df6.corr(),cbar=True,annot=True,cmap=‘Blues’)

Out[125]: <AxesSubplot>
```



```
In [130]: #Since there is one column with a large amount of missing data, we deem it to be insignificant and delete that
df6.drop(['buildUpPlayGribbling'],axis=1,inplace=True)

In [131]: sns.heatmap(df6.isnull(),cbar=False,cmap=‘viridis’)
#Used to get the heatmap for missing or incorrect data

Out[131]: <AxesSubplot>
```



This concludes EDA for all tables involved in dataset