

```
In [1]: 1 ## Setting up notebook
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import warnings
7 warnings.filterwarnings('ignore')
8
9 #pd.set_option('display.max_rows', None)
```

Problem 1 [25 marks]: Eurovision Winners Data Set

1A

```
In [2]: 1 df = pd.read_csv(r'/Users/elliotttoates/OneDrive - University of Exeter/Data Science-Elliott's MacBook Pro/P
2
3 display (df)
4 print()
5 print('Column Headings:',list (df.keys()))
6
7
```

	Year	Date	Host City	Winner	Song	Performer	Points	Margin	Runner-up
0	1957	3-Mar	Frankfurt	Netherlands	"Net als toen"	Corry Brokken	31	14	France
1	1958	12-Mar	Hilversum	France	"Dors, mon amour"	Andr_ Claveau	27	3	Switzerland
2	1959	11-Mar	Cannes	Netherlands	"Een beetje"	Teddy Scholten	21	5	United Kingdom
3	1960	29-Mar	London	France	"Tom Pillibi"	Jacqueline Boyer	32	7	United Kingdom
4	1961	18-Mar	Cannes	Luxembourg	"Nous les amoureux"	Jean-Claude Pascal	31	7	United Kingdom
...
60	2014	10-May	Copenhagen	Austria	"Rise Like a Phoenix"	Conchita Wurst	290	52	Netherlands
61	2015	23-May	Vienna	Sweden	"Heroes"	Mïns Zelmerl_w	365	62	Russia
62	2016	14-May	Stockholm	Ukraine	"1944"	Jamala	534	23	Australia
63	2017	13-May	Kiev	Portugal	"Amar pelos dois"	Salvador Sobral	758	143	Bulgaria
64	2018	18-May	Lisbon	Israel	"Toy"	Netta	529	93	Cyprus

65 rows x 9 columns

Column Headings: ['Year', 'Date', 'Host City', 'Winner', 'Song', 'Performer', 'Points', 'Margin', 'Runner-up']

1B

```
In [3]: 1 display(df.loc[df.Year.duplicated(), :]) #extracting rows with duplicate values in year collumn
2 display(df.iloc[[12]])
```

	Year	Date	Host City	Winner	Song	Performer	Points	Margin	Runner-up
13	1969	29-Mar	Madrid	United Kingdom	"Boom Bang-a-Bang"	Lulu	18	No runner-up	No runner-up
14	1969	29-Mar	Madrid	Netherlands	"De troubadour"	Lenny Kuhr	18	No runner-up	No runner-up
15	1969	29-Mar	Madrid	France	"Un jour, un enfant"	Frida Boccara	18	No runner-up	No runner-up
	Year	Date	Host City	Winner	Song	Performer	Points	Margin	Runner-up
12	1969	29-Mar	Madrid	Spain	"Vivo cantando"	Salom_	18	No runner-up	No runner-up

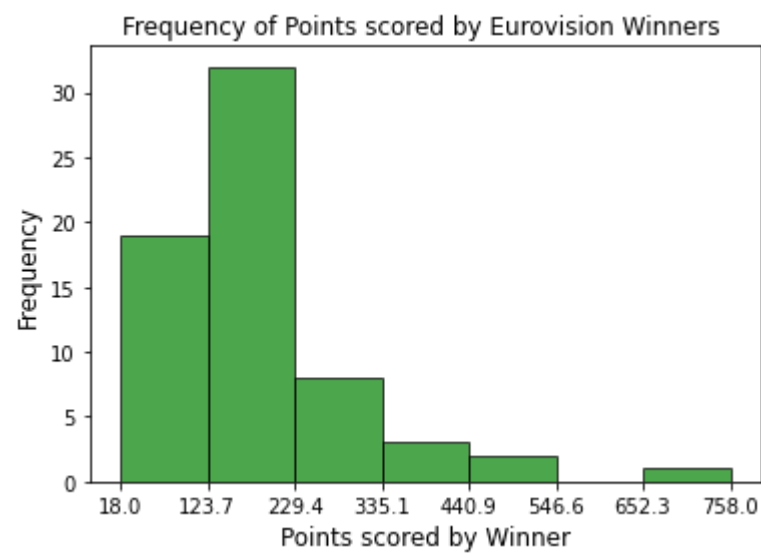
The 3 Additional rows can be found by searching for rows that have duplicate values for in the "Year" column of the dataframe. It is apparent that in 1969 the data set has 4 different winners of which the last 3 are considered duplicates. After having done research online it the additional rows are valid as there was a four-way tie in the 1969 Eurovision')

1C

```
In [4]: 1 limitDF = df.loc[:,['Year','Points','Margin']]
```

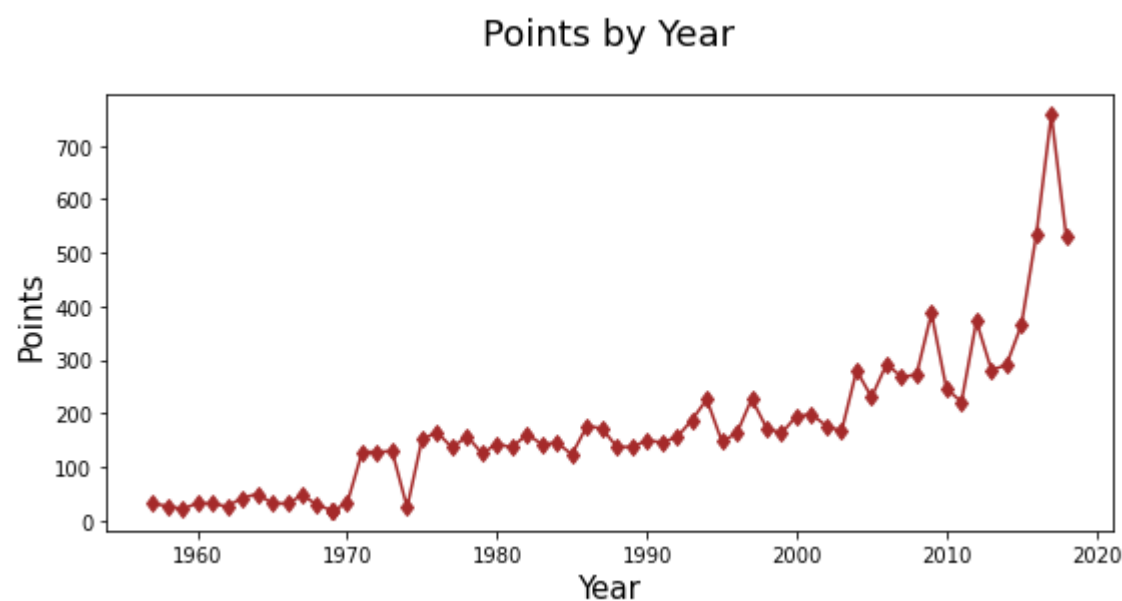
1D

```
In [5]: 1 %matplotlib inline
2
3 plt.title('Frequency of Points scored by Eurovision Winners ')
4 plt.xlabel('Points scored by Winner', size = 'large')
5
6 plt.ylabel('Frequency', size = 'larger')
7
8 plt.xticks(np.linspace(18,758,8))
9
10
11 plt.hist(df[ 'Points' ], bins=7,alpha=0.7,color='green',edgecolor='black');
```



1E

```
In [8]: 1 #part 1E
2 %matplotlib inline
3 plt.figure(figsize=(9,4))
4 plt.title('Points by Year', size = 18, pad = 25)
5 plt.xlabel('Year', size = 15)
6 plt.ylabel('Points', size = 15)
7
8 plt.plot(df[ 'Year' ],df[ 'Points' ],color = 'brown', marker = 'd' );
9
10
```



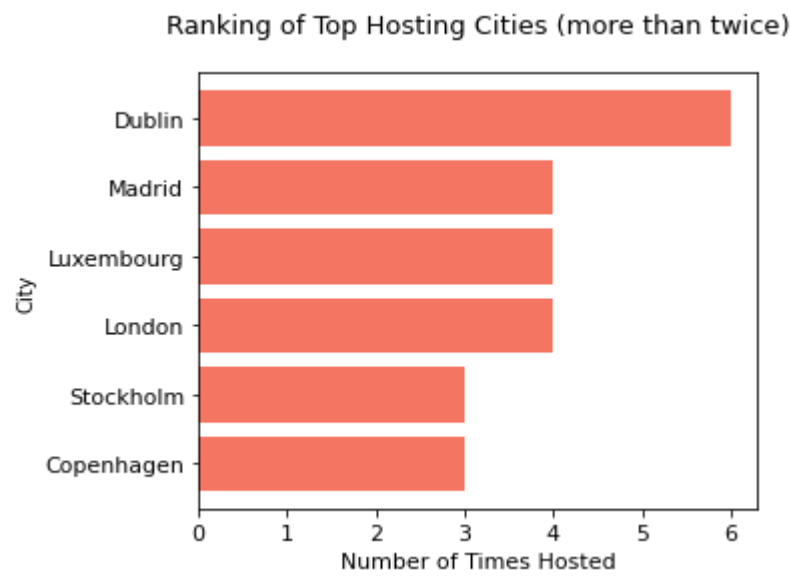
Rise in points : explanation

A plausible explanation for the rise in points that respective eurovision winners receive is that it is due to increasing amount of participants while the points that each participant allocates to others in a given year has remained the same.

With the data set provided, one cannot confirm this explanation as it doesnt include data about the amount of participants however upon researching further its apparent that this is the explanations as the number of participants has risen from 7 countries in 1956 to records of 43 participants in 2008.

1F

```
In [9]: 1 df2 = pd.DataFrame(df[ 'Host City' ].value_counts())
2
3 df2 = df2.iloc[0:6]
4 df2 = df2.reindex([ 'Copenhagen', 'Stockholm', 'London', 'Luxembourg', 'Madrid', 'Dublin' ])
5
6 df2.plot.barh(color='#F47562', legend = False, width = 0.8, figsize= (5,4))
7
8 plt.title('Ranking of Top Hosting Cities (more than twice)',size = 13,pad = 20)
9 plt.ylabel('City',size = 10)
10 plt.xlabel('Number of Times Hosted',size = 11)
11 plt.tick_params(axis = 'both', labelsiz = 11)
12
13
14
```



1G

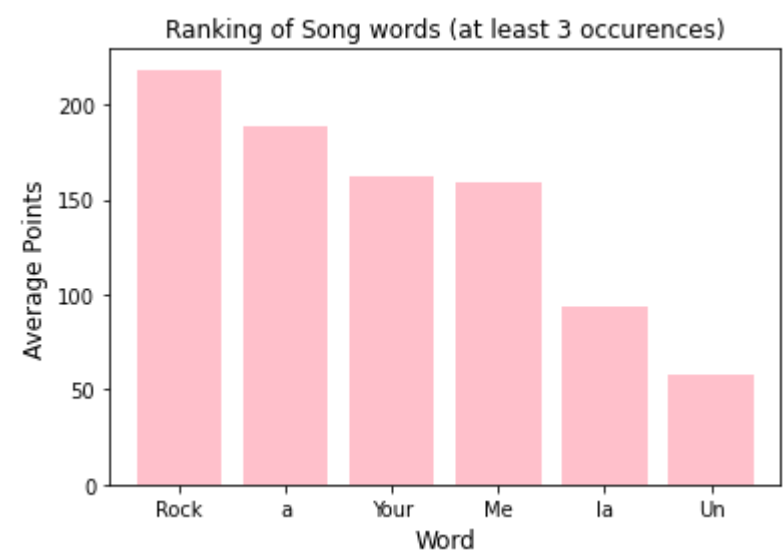
```
In [10]: 1 #part a
2 def function (row):
3     df[ 'Song' ] = df.Song.str.replace('[_()''",-]', '')
4     n = row.name
5
6     wordlist = df.loc[n, 'Song'].split(' ')
7     points = [df.loc[n, 'Points']] * len(wordlist)
8
9     return(list(zip(wordlist, points)))
```

```
In [11]: 1 function(df.loc[0,:])
```

```
Out[11]: [('Net', 31), ('als', 31), ('toen', 31)]
```

```
In [12]: 1 #part b
2 tuple_list = ((df.apply(function,axis = 1)).explode()).tolist()
3 #part c
4 df_words = pd.DataFrame(tuple_list, columns = [ 'Word', 'Points' ])
5 #part D
6 df_words= df_words.groupby('Word').filter(lambda x: len(x) > 2)
7 #part e
8 df_words = (df_words.groupby('Word',as_index= False).mean())
```

```
In [13]: 1 #part f
2 df_words = df_words.sort_values (by='Points',ascending = False)
3
4 %matplotlib inline
5
6 df_words.plot.bar(x='Word',y = 'Points',rot =0,color = 'pink',width = 0.8,legend = False)
7
8 plt.title('Ranking of Song words (at least 3 occurences)')
9 plt.xlabel('Word', size = 'large')
10 plt.ylabel('Average Points', size = 'larger');
```



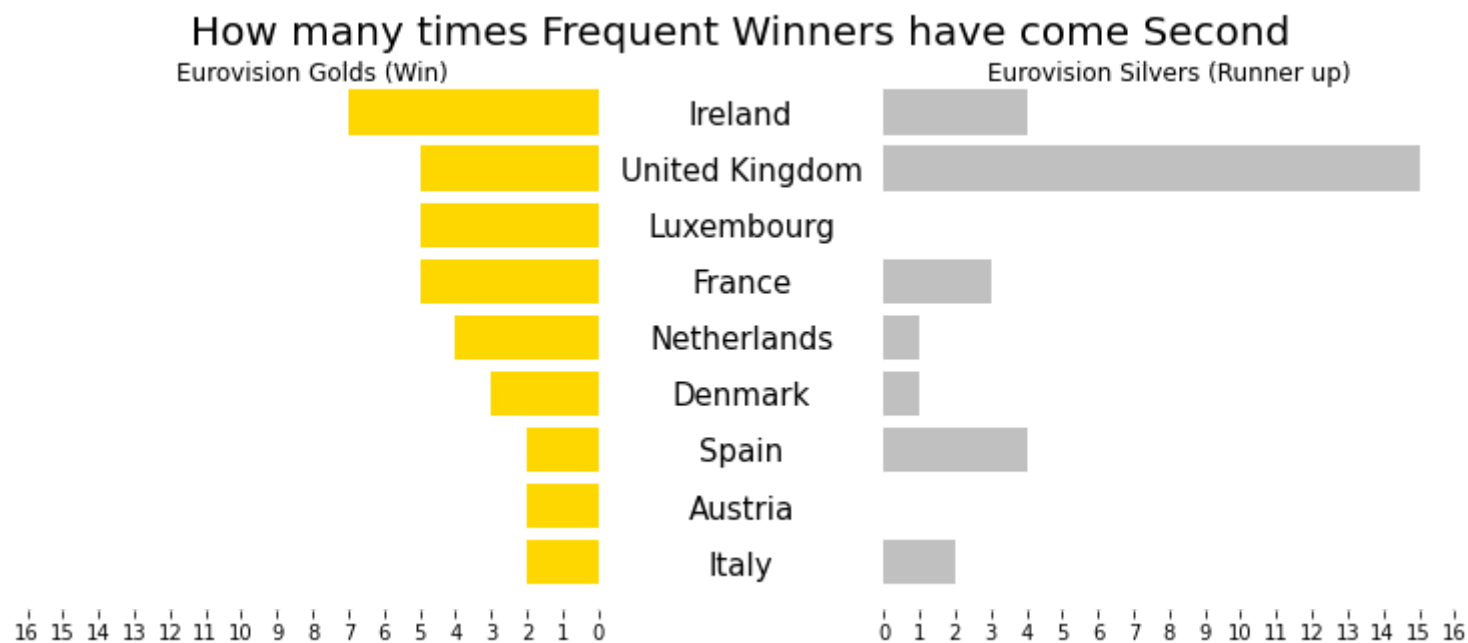
1H

```
In [14]: 1 df['Total Wins'] = df['Winner'].map(df['Winner'].value_counts())
2 df ['Total 2nd Place'] = df['Winner'].map(df['Runner-up'].value_counts())
3
4
5 df['Total 2nd Place'] = df['Total 2nd Place'].fillna(0)
6 df['Total 2nd Place'] = df['Total 2nd Place'].astype('int')
7
8
9 df = df.drop_duplicates (subset = 'Winner', ignore_index = True)
10
11 df['Country'] = df['Winner']
12
13 df = df[['Country','Total Wins', 'Total 2nd Place']]
14
15 df = df.head (9)
16 df = df.sort_values (by =['Total Wins'] ,ascending = True)
17 df
```

Out[14]:

	Country	Total Wins	Total 2nd Place
4	Italy	2	2
5	Austria	2	0
7	Spain	2	4
3	Denmark	3	1
0	Netherlands	4	1
1	France	5	3
2	Luxembourg	5	0
6	United Kingdom	5	15
8	Ireland	7	4

```
In [15]: 1 Country = df['Country']
2 num_country = len(Country)
3 Wins = df['Total Wins']
4 Runner_up = df['Total 2nd Place']
5 pos = np.arange(num_country) + .5
6
7 fig, (ax_left, ax_right) = plt.subplots(ncols=2, figsize = (11,5))
8 ax_left.barh(pos, Wins, align='center', facecolor='gold')
9 ax_left.set_yticks([])
10
11 ax_right.barh(pos, Runner_up, align='center', facecolor = 'silver')
12 ax_right.set_yticks(pos)
13 ax_right.set_yticklabels(Country, ha='center', x=-0.225, size = 15)
14
15 ax_right.set_xlim((0,16))
16 ax_left.set_xlim((16,0))
17
18 plt.subplots_adjust(left = 0.05,right = 0.95,wspace = 0.5)
19
20 for i in [ax_left,ax_right]:
21
22     i.spines["top"].set_visible(False)
23     i.spines["right"].set_visible(False)
24     i.spines["bottom"].set_visible(False)
25     i.spines["left"].set_visible(False)
26     i.tick_params(left = False)
27
28     i.set_xticks(range(0,17,1))
29
30 plt.suptitle('How many times Frequent Winners have come Second',size =20, y=0.95)
31 ax_left.set_title('Eurovision Golds (Win)',y = 0.95)
32 ax_right.set_title('Eurovision Silvers (Runner up)',y = 0.95);
```



The extra analysis above is a 'tornado' graph that offers insight into the amount of times countries that have one eurovision have also come second place.

One can observe from the figure above that the amount of times a country has come second is not indicative of the amount of wins. I chose the tornado style to show this as it is capable of showing this general assymetry between the wins and 2nd places well.

Problem 2 [25 marks]: UEFA Euro 2020 Scores Data Set

2A

```
In [16]: 1 df = pd.read_csv(r'/Users/elliottoates/OneDrive - University of Exeter/Data Science-Elliott's MacBook Pro/P
2
3 display (df)
4 print()
5 print('Column Headings:',list (df.keys()))
```

40	Group stage: Matchday 1	15.06.2021	False	False	False	Hungary	Portugal	0	3
41	Group stage: Matchday 1	14.06.2021	False	False	False	Spain	Sweden	0	0
42	Group stage: Matchday 1	14.06.2021	False	False	False	Poland	Slovakia	1	2

2B

```
In [17]: 1 unique_countries = pd.unique(df[['team_name_home', 'team_name_away']].values.ravel())
2 print('list of Unique Teams:',unique_countries)
```

list of Unique Teams: [' Italy ' ' England ' ' Denmark ' ' Spain ' ' Ukraine ' ' Czech Republic ' ' Belgium ' ' Switzerland ' ' Sweden ' ' Germany ' ' France ' ' Croatia ' ' Portugal ' ' Netherlands ' ' Austria ' ' Wales ' ' Hungary ' ' Poland ' ' Slovakia ' ' Scotland ' ' Finland ' ' Russia ' ' North Macedonia ' ' Turkey ']

```
In [18]: 1 df['team_name_home'] = df['team_name_home'].str.strip()
2 df['team_name_away'] = df['team_name_away'].str.strip()
3 unique_countries = pd.unique(df[['team_name_home', 'team_name_away']].values.ravel())
4 print('list of Unique Teams:',unique_countries)
```

list of Unique Teams: ['Italy' 'England' 'Denmark' 'Spain' 'Ukraine' 'Czech Republic' 'Belgium' 'Switzerland' 'Sweden' 'Germany' 'France' 'Croatia' 'Portugal' 'Netherlands' 'Austria' 'Wales' 'Hungary' 'Poland' 'Slovakia' 'Scotland' 'Finland' 'Russia' 'North Macedonia' 'Turkey']

2C

```
In [19]: 1 df_G = df.iloc [15:]
```

2D

```
In [20]: 1 def calcScore(row):
2     n = row.name
3     h_points=1
4     a_points=1
5     if df_G.loc[n,'team_home_score']>(df_G.loc[n,'team_away_score']):
6         h_points=3
7         a_points-=1
8     elif df_G.loc[n,'team_home_score']<(df_G.loc[n,'team_away_score']):
9         a_points=3
10        h_points-=1
11    return(h_points,a_points)
12
```

```
In [21]: 1 calcScore(df_G.loc[17,:])
```

Out[21]: (3, 0)

2E

In [22]:

1	df_G['team_home_points'] = (df_G.apply(calcScore, axis = 1, result_type = 'expand'))[0]
2	df_G['team_away_points'] = (df_G.apply(calcScore, axis = 1, result_type = 'expand'))[1]
3	df_G
41	stage: 14.06.2021 False False False Spain Sweden 0 0 Matchday 1
42	Group stage: 14.06.2021 False False False Poland Slovakia 1 2 Matchday 1
43	Group stage: 14.06.2021 False False False Scotland Czech Republic 0 2 Matchday 1
44	Group stage: 13.06.2021 False False False Netherlands Ukraine 3 2 Matchday 1
45	Group stage: 13.06.2021 False False False Austria North Macedonia 3 1 Matchday 1

2F

In [23]:

1	def GetGroupTeams(Team):
2	group = [Team]
3	away_opp = list(df_G.loc[df_G['team_name_home'] == Team, 'team_name_away'])
4	home_opp = list(df_G.loc[df_G['team_name_away'] == Team, 'team_name_home'])
5	group += away_opp + home_opp
6	return group
7	

In [24]:

1	GetGroupTeams('England')
---	--------------------------

Out[24]: ['England', 'Scotland', 'Croatia', 'Czech Republic']

2G

In [25]:

1	def GetGroupRows(Team):
2	group_results = df_G[df_G['team_name_home'].isin(GetGroupTeams(Team))]
3	group_results = group_results.reset_index(drop=True)
4	return group_results

In [26]:

1	GetGroupRows('England')
---	-------------------------

Out[26]:

	stage	date	pens	pens_home_score	pens_away_score	team_name_home	team_name_away	team_home_score	team_away_score	poss
0	Group stage: Matchday 3	22.06.2021	False	False	False	Croatia	Scotland	3	1	
1	Group stage: Matchday 3	22.06.2021	False	False	False	Czech Republic	England	0	1	
2	Group stage: Matchday 2	18.06.2021	False	False	False	England	Scotland	0	0	
3	Group stage: Matchday 2	18.06.2021	False	False	False	Croatia	Czech Republic	1	1	
4	Group stage: Matchday 1	14.06.2021	False	False	False	Scotland	Czech Republic	0	2	
5	Group stage: Matchday 1	13.06.2021	False	False	False	England	Croatia	1	0	

6 rows × 22 columns

2H

```
In [27]: 1 def GetTeamStats (Team):
2         Matches= df_G[df_G.isin([Team]).any(axis=1)]
3
4         Points = Matches.loc[Matches['team_name_home'] == Team, 'team_home_points'].sum()\
5         +Matches.loc[Matches['team_name_away'] == Team, 'team_away_points'].sum()
6
7         Scored = Matches.loc[Matches['team_name_home'] == Team, 'team_home_score'].sum()\
8         +Matches.loc[Matches['team_name_away'] == Team, 'team_away_score'].sum()
9
10        Conceded = Matches.loc[Matches['team_name_home'] == Team, 'team_away_score'].sum()\
11        +Matches.loc[Matches['team_name_away'] == Team, 'team_home_score'].sum()
12
13        Stats = {'Points':Points, 'Goals Scored': Scored, 'Goals Conceded': Conceded}
14
15        return Stats
```

```
In [28]: 1 GetTeamStats ('England')
```

Out[28]: {'Points': 7, 'Goals Scored': 2, 'Goals Conceded': 0}

2I

```
In [29]: 1 def GetGroupTable(Team):
2         group = GetGroupTeams(Team)
3         table = pd.DataFrame(GetTeamStats(Team),index=[])
4         for i in group :
5             stats = pd.DataFrame(GetTeamStats(i),index =[i])
6             table = table.append(stats)
7
8         table['Goal Difference'] = table['Goals Scored'] - table['Goals Conceded']
9         table = table.sort_values(['Points', 'Goal Difference', 'Goals Scored'],ascending = False)
10        return table
```

```
In [30]: 1 GetGroupTable('England')
2
```

Out[30]:

	Points	Goals Scored	Goals Conceded	Goal Difference
England	7	2	0	2
Croatia	4	4	3	1
Czech Republic	4	3	2	1
Scotland	1	1	5	-4

```
In [31]: 1 GetGroupTable('Italy')
```

Out[31]:

	Points	Goals Scored	Goals Conceded	Goal Difference
Italy	9	7	0	7
Wales	4	3	2	1
Switzerland	4	4	5	-1
Turkey	0	1	8	-7

2J

```
In [32]: 1 def GetAllGroupTables():
2
3         for y in unique_countries:
4             dft=pd.DataFrame(sorted(GetGroupTeams(y)))
5
6         for i in unique_countries:
7             dft[i]=pd.DataFrame(sorted(GetGroupTeams(i)))
8             dft=dft.T.drop_duplicates().T
9
10        for x in dft.iloc[1]:
11            print(GetGroupTable(x))
```


In [33]:

1	GetAllGroupTables()					
	SCOTLAND	1	1	5	-4	
		Points	Goals Scored	Goals Conceded	Goal Difference	
	Belgium	9	7	1	6	
	Denmark	3	5	4	1	
	Finland	3	1	3	-2	
	Russia	3	2	7	-5	
		Points	Goals Scored	Goals Conceded	Goal Difference	
	Sweden	7	4	2	2	
	Spain	5	6	1	5	
	Slovakia	3	2	7	-5	
	Poland	1	4	6	-2	
		Points	Goals Scored	Goals Conceded	Goal Difference	
	Netherlands	9	8	2	6	
	Austria	6	4	3	1	
	Ukraine	3	4	5	-1	
	North Macedonia	0	2	8	-6	
		Points	Goals Scored	Goals Conceded	Goal Difference	
	France	5	4	3	1	
	Portugal	4	7	6	1	
	Germany	4	6	5	1	

K : extra analysis

In [34]:

```
1 df['possession_away'] = df['possession_away'].str.strip(' %').astype(int)
2 df['possession_home'] = df['possession_home'].str.strip(' %').astype(int)
3 df['duels_won_home'] = df['duels_won_home'].str.strip(' %').astype(int)
4 df['duels_won_away'] = df['duels_won_away'].str.strip(' %').astype(int)
5
6 def GetAllTeamStats (Team):
7
8     Matches = df[df.isin([Team]).any(axis=1)]
9     Matches = Matches.iloc [1:]
10
11     AvgScored = (df.loc[df['team_name_home'] == Team, 'team_home_score'].sum()\
12 +df.loc[df['team_name_away'] == Team, 'team_away_score'].sum())
13
14
15     AvgConceded = (df.loc[df['team_name_home'] == Team, 'team_away_score'].sum()\
16 +df.loc[df['team_name_away'] == Team, 'team_home_score'].sum())
17
18
19     AvgPossession = (df.loc[df['team_name_home'] == Team, ('possession_home')].sum()\
20 +df.loc[df['team_name_away'] == Team, ('possession_away')].sum())
21
22
23     AvgSONT = (df.loc[df['team_name_home'] == Team, ('shots_on_target_home')].sum()\
24 +df.loc[df['team_name_away'] == Team, ('shots_on_target_away')].sum())
25
26
27     AvgSONTConceded = (df.loc[df['team_name_home'] == Team, ('shots_on_target_away')].sum()\
28 +df.loc[df['team_name_away'] == Team, ('shots_on_target_home')].sum())
29
30
31     AvgDuelsWon = (df.loc[df['team_name_home'] == Team, ('duels_won_home')].sum()\
32 +df.loc[df['team_name_away'] == Team, ('duels_won_away')].sum())
33
34     AllStats = {'Team': Team,
35                 'Goals scored': AvgScored,
36                 'Goals conceded': AvgConceded,
37                 'Mean Possession(%)': AvgPossession,
38                 'Shots on Target': AvgSONT,
39                 'Shots on Target Conceded': AvgSONTConceded,
40                 'Mean Duels Won(%)' : AvgDuelsWon
41             }
42
43     return AllStats
```

```
In [35]: 1 df_finalists = pd.DataFrame([GetAllTeamStats ('England'),GetAllTeamStats ('Italy')])
2 factors = ['Goals scored','Goals conceded','Mean Possession(%)','Shots on Target','Shots on Target Conceded
3
4 df = df_finalists
5
6 new_max = 100
7 new_min = 0
8 new_range = new_max - new_min
9
10 for factor in factors:
11     max_val = df_finalists[factor].max()
12     min_val = df_finalists[factor].min()
13     val_range = max_val - min_val
14     df_finalists[factor + '_adj'] = df_finalists[factor].apply(
15         lambda x: (((x/ max_val) )))
16
17 df_finalists = df_finalists.loc[:,['Team','Goals scored_adj','Goals conceded_adj','Mean Possession(%)_adj',
18
19
20 df_finalists.set_index('Team',inplace=True)
21 df_finalists
22
```

Out[35]:

	Goals scored_adj	Goals conceded_adj	Mean Possession(%)_adj	Shots on Target_adj	Shots on Target Conceded_adj	Mean Duels Won(%)_adj
Team						
England	0.846154	0.5	0.929688	0.777778	1.000000	1.00000
Italy	1.000000	1.0	1.000000	1.000000	0.842105	0.93956

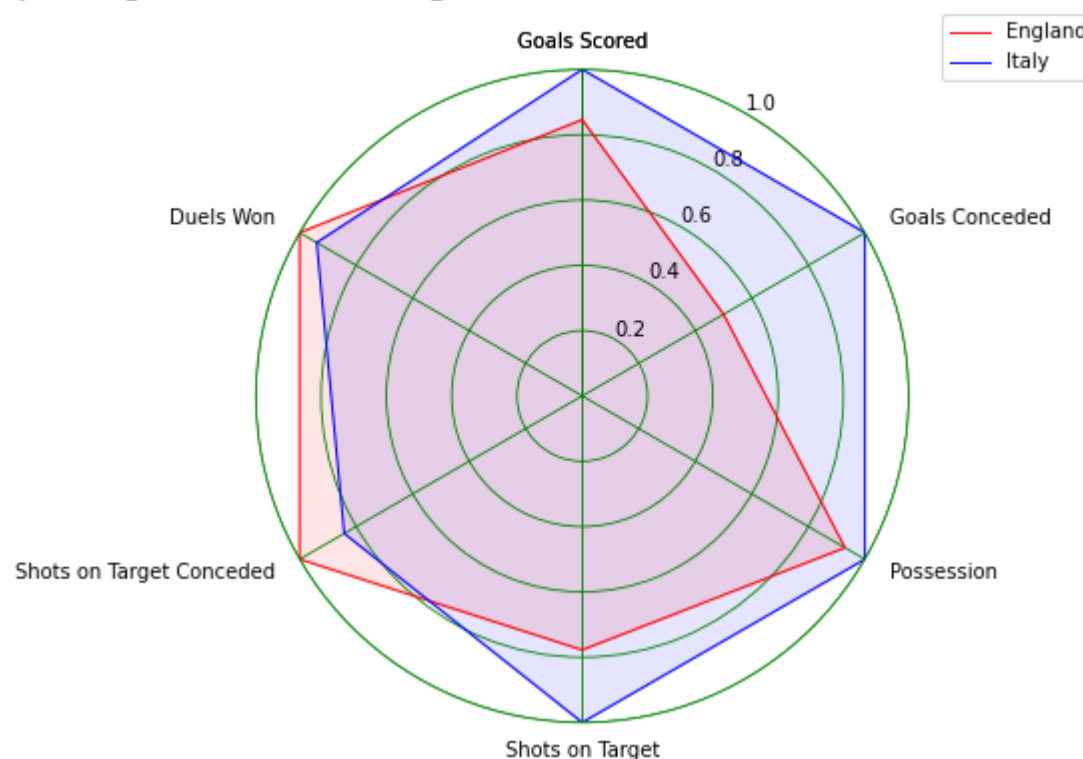
Figure

```

In [36]: 1 labels = ['', 'Goals Conceded', 'Possession', 'Shots on Target', 'Shots on Target Conceded', 'Duels Won', 'Goals Scored']
2 num_vars = len(labels)
3
4 angles = np.linspace(0, 2 * np.pi, num_vars, endpoint=False).tolist()
5 angles += angles[:1]
6
7 fig, ax = plt.subplots(figsize=(6, 6), subplot_kw=dict(polar=True))
8
9 #England
10 values = df_finalists.loc['England'].tolist()
11 values += values[:1]
12 ax.plot(angles, values, color='red', linewidth=1, label= 'England')
13 ax.fill(angles, values, color='red', alpha=0.1)
14
15 #Italy
16 values = df_finalists.loc['Italy'].tolist()
17 values += values[:1]
18 ax.plot(angles, values, color='blue', linewidth=1, label= 'Italy')
19 ax.fill(angles, values, color='blue', alpha=0.1)
20
21
22 #Grid
23 ax.set_theta_offset(np.pi / 2)
24 ax.set_theta_direction(-1)
25 ax.set_thetagrids(np.degrees(angles), labels)
26 ax.grid(color='green')
27 ax.spines['polar'].set_color('green')
28
29 #Labels
30 for label, angle in zip(ax.get_xticklabels(), angles):
31     if angle in (0, np.pi):
32         label.set_horizontalalignment('center')
33     elif 0 < angle < np.pi:
34         label.set_horizontalalignment('left')
35     else:
36         label.set_horizontalalignment('right')
37 ax.set_ylim(0, 1)
38 ax.set_rlabel_position(180 / num_vars)
39
40 #title
41 ax.set_title('Comparing relative strengths in the Euro 2020 finalists campaign', y=1.1, Size = 20)
42
43 ax.legend(loc='upper right', bbox_to_anchor=(1.3, 1.1));

```

Comparing relative strengths in the Euro 2020 finalists campaign



The above figure is a radar chart that compares the relative strengths and weaknesses of the Italy and European Teams in their campaign to the Final.

Data is collected from the group and knockout-stages (excluding the final).

The different attributes have all been scaled to show the relative weakness of either team. For example with goals scored, Italy has a value of 1, and England has one of ~0.84 which is the amount of goals England scored in the run-up to the final relative to the amount Italy scored.

A graph like this could be used in preparation of the final to help either team focus on aspects of the game that the other team is significantly stronger in. Although in the case of this final, both teams had in most metrics a similar campaign. One must note that for goals conceded England has a lower score but due to the nature of the metric this is actually "better".

Problem 3 [25 marks]: COVID-19

3A

```
In [37]: 1 df = pd.read_csv(r'/Users/elliotttoates/OneDrive - University of Exeter/Data Science-Elliott's MacBook Pro/P
2 display (df)
3 print('Column Headings:',list (df.keys()))
```

	iso_code	continent	location	date	total_cases	new_cases	new_cases_smoothed	total_deaths	new_deaths	new_deaths_smoothed	...
0	AFG	Asia	Afghanistan	2020-02-24	5.0	5.0	NaN	NaN	NaN	NaN	...
1	AFG	Asia	Afghanistan	2020-02-25	5.0	0.0	NaN	NaN	NaN	NaN	...
2	AFG	Asia	Afghanistan	2020-02-26	5.0	0.0	NaN	NaN	NaN	NaN	...
3	AFG	Asia	Afghanistan	2020-02-27	5.0	0.0	NaN	NaN	NaN	NaN	...
4	AFG	Asia	Afghanistan	2020-02-28	5.0	0.0	NaN	NaN	NaN	NaN	...
...
168159	ZWE	Africa	Zimbabwe	2022-03-09	240343.0	0.0	405.714	5400.0	0.0	0.571	...
168160	ZWE	Africa	Zimbabwe	2022-03-10	241548.0	1205.0	577.857	5408.0	8.0	1.714	...
168161	ZWE	Africa	Zimbabwe	2022-03-11	241548.0	0.0	401.286	5408.0	0.0	1.571	...
168162	ZWE	Africa	Zimbabwe	2022-03-12	242069.0	521.0	435.714	5412.0	4.0	2.143	...
168163	ZWE	Africa	Zimbabwe	2022-03-13	242515.0	446.0	472.286	5414.0	2.0	2.143	...

168164 rows × 67 columns

Column Headings: ['iso_code', 'continent', 'location', 'date', 'total_cases', 'new_cases', 'new_cases_smoothed', 'total_deaths', 'new_deaths', 'new_deaths_smoothed', 'total_cases_per_million', 'new_cases_per_million', 'new_cases_smoothed_per_million', 'total_deaths_per_million', 'new_deaths_per_million', 'new_deaths_smoothed_per_million', 'reproduction_rate', 'icu_patients', 'icu_patients_per_million', 'hosp_patients', 'hosp_patients_per_million', 'weekly_icu_admissions', 'weekly_icu_admissions_per_million', 'weekly_hosp_admissions', 'weekly_hosp_admissions_per_million', 'total_tests', 'new_tests', 'total_tests_per_thousand', 'new_tests_per_thousand', 'new_tests_smoothed', 'new_tests_smoothed_per_thousand', 'positive_rate', 'tests_per_case', 'tests_units', 'total_vaccinations', 'people_vaccinated', 'people_fully_vaccinated', 'total_boosters', 'new_vaccinations', 'new_vaccinations_smoothed', 'total_vaccinations_per_hundred', 'people_vaccinated_per_hundred', 'people_fully_vaccinated_per_hundred', 'total_boosters_per_hundred', 'new_vaccinations_smoothed_per_million', 'new_people_vaccinated_smoothed', 'new_people_vaccinated_smoothed_per_hundred', 'stringency_index', 'population', 'population_density', 'median_age', 'aged_65_old', 'aged_70_old', 'gdp_per_capita', 'extreme_poverty', 'cardiovasc_death_rate', 'diabetes_prevalence', 'female_smokers', 'male_smokers', 'handwashing_facilities', 'hospital_beds_per_thousand', 'life_expectancy', 'human_development_index', 'excess_mortality_cumulative_absolute', 'excess_mortality_cumulative', 'excess_mortality_cumulative_per_million']

3B

```
In [38]: 1 print (df['date'].unique())

['2020-02-24' '2020-02-25' '2020-02-26' '2020-02-27' '2020-02-28'
'2020-02-29' '2020-03-01' '2020-03-02' '2020-03-03' '2020-03-04'
'2020-03-05' '2020-03-06' '2020-03-07' '2020-03-08' '2020-03-09'
'2020-03-10' '2020-03-11' '2020-03-12' '2020-03-13' '2020-03-14'
'2020-03-15' '2020-03-16' '2020-03-17' '2020-03-18' '2020-03-19'
'2020-03-20' '2020-03-21' '2020-03-22' '2020-03-23' '2020-03-24'
'2020-03-25' '2020-03-26' '2020-03-27' '2020-03-28' '2020-03-29'
'2020-03-30' '2020-03-31' '2020-04-01' '2020-04-02' '2020-04-03'
'2020-04-04' '2020-04-05' '2020-04-06' '2020-04-07' '2020-04-08'
'2020-04-09' '2020-04-10' '2020-04-11' '2020-04-12' '2020-04-13'
'2020-04-14' '2020-04-15' '2020-04-16' '2020-04-17' '2020-04-18'
'2020-04-19' '2020-04-20' '2020-04-21' '2020-04-22' '2020-04-23'
'2020-04-24' '2020-04-25' '2020-04-26' '2020-04-27' '2020-04-28'
'2020-04-29' '2020-04-30' '2020-05-01' '2020-05-02' '2020-05-03'
'2020-05-04' '2020-05-05' '2020-05-06' '2020-05-07' '2020-05-08'
'2020-05-09' '2020-05-10' '2020-05-11' '2020-05-12' '2020-05-13'
'2020-05-14' '2020-05-15' '2020-05-16' '2020-05-17' '2020-05-18'
'2020-05-19' '2020-05-20' '2020-05-21' '2020-05-22' '2020-05-23'
'2020-05-24' '2020-05-25' '2020-05-26' '2020-05-27' '2020-05-28'
'2020-05-29' '2020-05-30' '2020-05-31' '2020-06-01' '2020-06-02']
```

3C

```
In [39]: 1 print ('number of unique countries in data set:',len(df['iso_code'].unique()))
```

number of unique countries in data set: 238

3D

```
In [40]: 1 df = df[df['iso_code'].str.contains("OWID_") == False]
2 df = df.reset_index()

In [41]: 1 print ('number of unique countries in data set:',len(df['iso_code'].unique()))

number of unique countries in data set: 223
```

3E

```
In [42]: 1 df_v = df[df['new_vaccinations'].notna()]
2 df_v = df_v.reset_index()
3
4 earliest_date = df_v['date'].min()
5 print (earliest_date)

2020-12-03
```

3F

```
In [43]: 1 df_v['total_vaccinations_per_million'] = 10000 * df_v['total_vaccinations_per_hundred']
```

3G

```
In [44]: 1 df_v['date'] = df_v['date'].apply(pd.to_datetime)
2 df_v['DaysSince3Dec20'] = (df_v['date']-df_v['date'].min()).dt.days
3
```

3H

```
In [45]: 1 min_deaths = df_v.groupby('iso_code')['total_deaths_per_million'].min() #earliest date deaths /mill
2
3 max_deaths = df_v.groupby('iso_code')['total_deaths_per_million'].max() #latest date deaths/mill
4
5 total_vax = df_v.groupby('iso_code')['total_vaccinations_per_hundred'].max() #latest date vax/100
6
7 df_country =pd.DataFrame ({'total_vaccinations_per_hundred': total_vax, 'total_deaths_per_million(earliest date)': min_deaths, 'total_deaths_per_million(latest date)': max_deaths})
8 df_country['TotalDeathPerMSinceVac'] = (df_country.iloc[:, 2])-(df_country.iloc[:, 1])
9 df_country
10
```

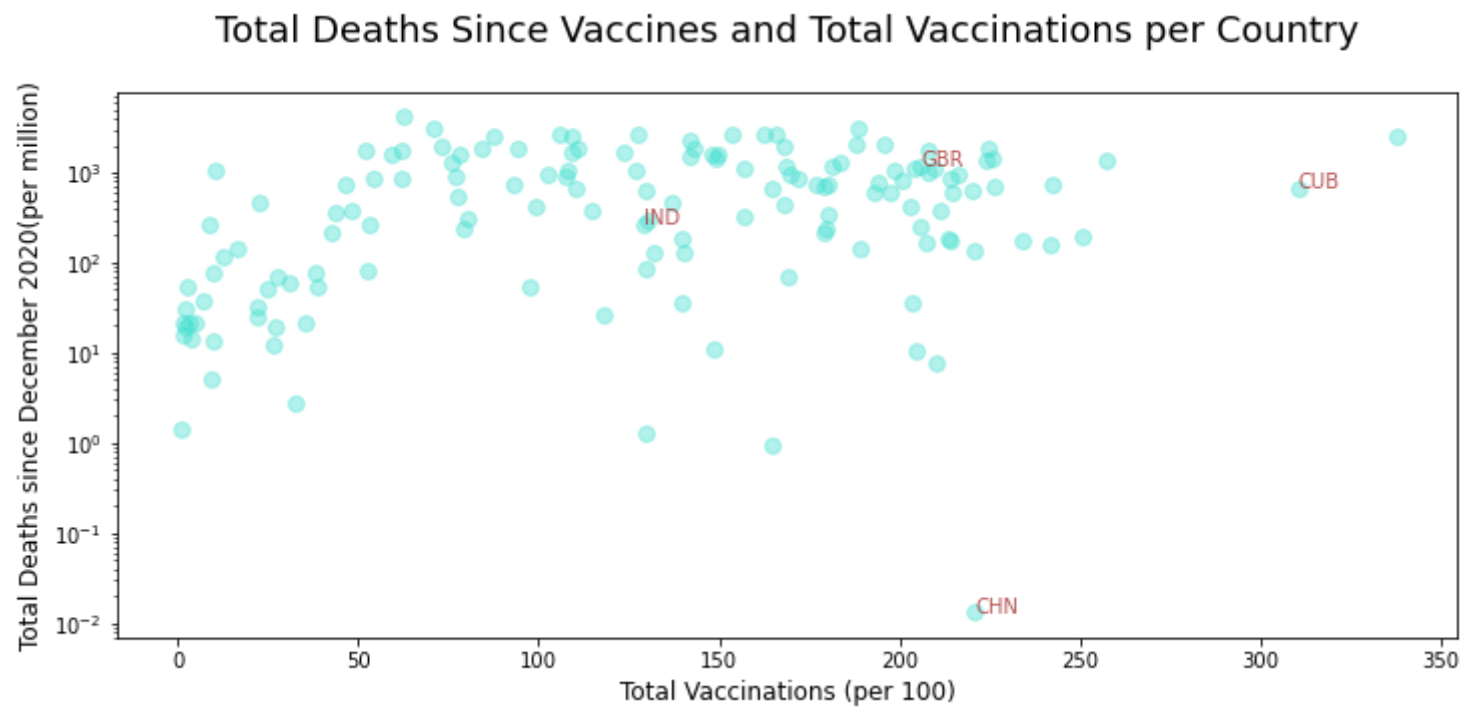
Out[45]:

	total_vaccinations_per_hundred	total_deaths_per_million(earliest date)	total_deaths_per_million(latest date)	TotalDeathPerMSinceVac
iso_code				
ABW	157.01	858.249	1977.704	1119.455
AFG	12.76	72.021	185.840	113.819
AIA	18.26	NaN	NaN	NaN
ALB	93.05	437.184	1191.465	754.281
AND	24.65	1512.527	1564.237	51.710
...
VCT	62.13	98.860	952.646	853.786
VNM	202.99	0.357	419.247	418.890
ZAF	54.16	813.747	1660.704	846.957
ZMB	16.31	65.008	209.295	144.287
ZWE	52.95	96.474	358.597	262.123

171 rows x 4 columns

3I

```
In [46]: 1 %matplotlib inline
2
3 plt.figure(figsize=(12,5))
4 plt.title('Total Deaths Since Vaccines and Total Vaccinations per Country', size = 18, pad = 25)
5 plt.ylabel('Total Deaths since December 2020(per million)', size = 12)
6 plt.xlabel('Total Vaccinations (per 100)', size = 12)
7
8
9 ax = plt.gca()
10 ax.set_yscale('log')
11
12 ax.text(220.84,0.013, "CHN", color = 'brown',alpha = 0.8) #CHN
13 ax.text(129.14,260.925, "IND", color = 'brown',alpha = 0.8) #IND
14 ax.text(205.72,1183.704, "GBR", color = 'brown',alpha = 0.8)
15 ax.text(310.42,665.960, "CUB", color = 'brown',alpha = 0.8)
16
17 plt.scatter(df_country ['total_vaccinations_per_hundred'],df_country['TotalDeathPerMSinceVac'], alpha=0.4, c
18
```



3J

```
In [47]: 1 df_GBR = df_v[df_v['iso_code'].str.contains('GBR')]
2 df_GBR.drop(columns=['level_0','index'],
3 axis=1,
4 inplace=True)
5 df_GBR = df_GBR.reset_index(drop=True)
6
7 df_GBR.shape #Checking that its 424 rows
```

Out[47]: (424, 69)

3K

```
In [48]: 1 df_GBR['total_vaccinations_per_hundred_smoothed'] = df_GBR['total_vaccinations_per_hundred'].rolling(window=
```

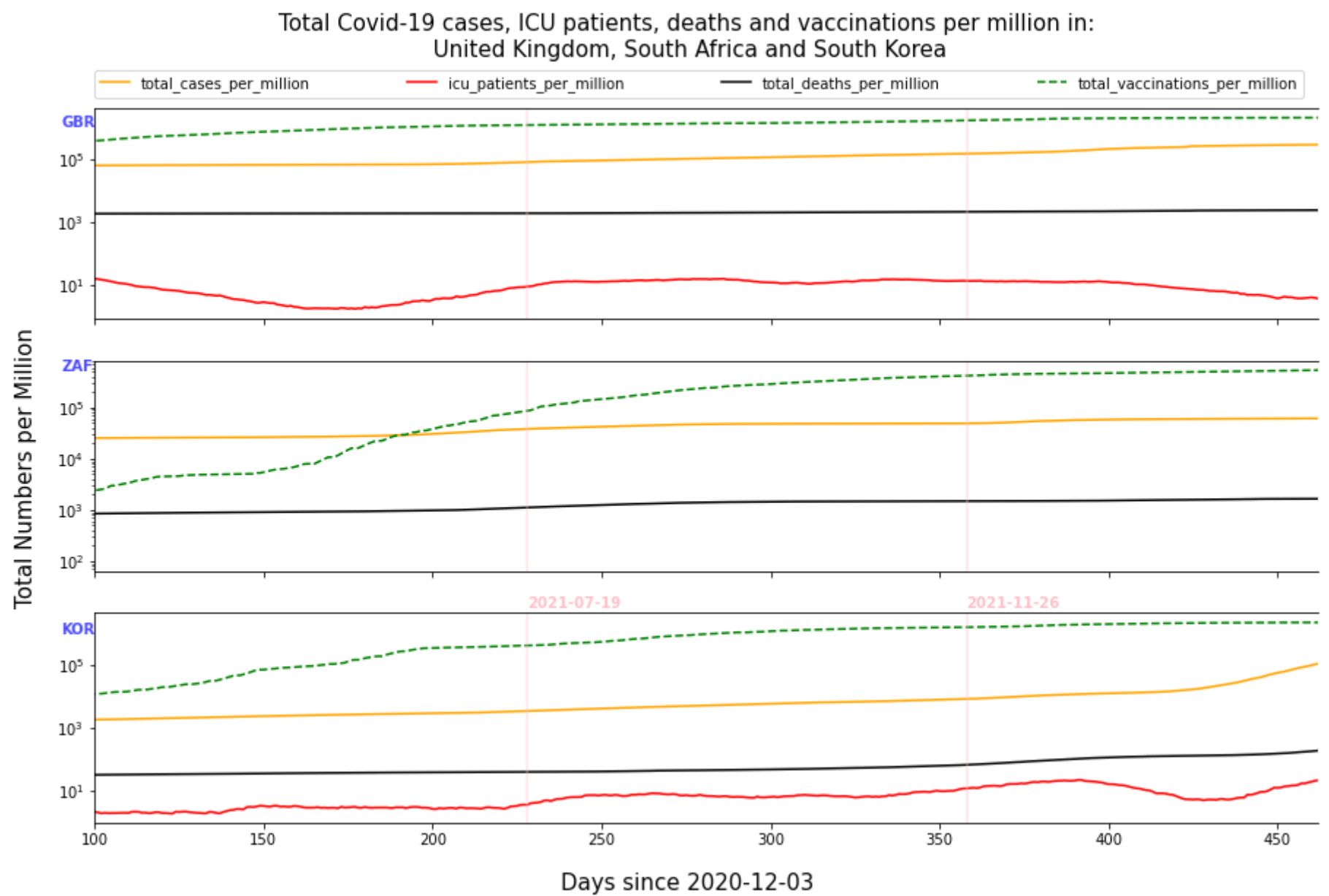
3L

In [49]:

```
1  #ZAF
2  df_ZAF = df_v[df_v['iso_code'].str.contains('ZAF')]
3  df_ZAF.drop(columns=['level_0', 'index'],
4               axis=1,
5               inplace=True)
6  df_ZAF = df_ZAF.reset_index(drop=True)
7
8  #KOR
9  df_KOR = df_v[df_v['iso_code'].str.contains('KOR')]
10 df_KOR.drop(columns=['level_0', 'index'],
11              axis=1,
12              inplace=True)
13 df_KOR = df_KOR.reset_index(drop=True)
```


In [50]:

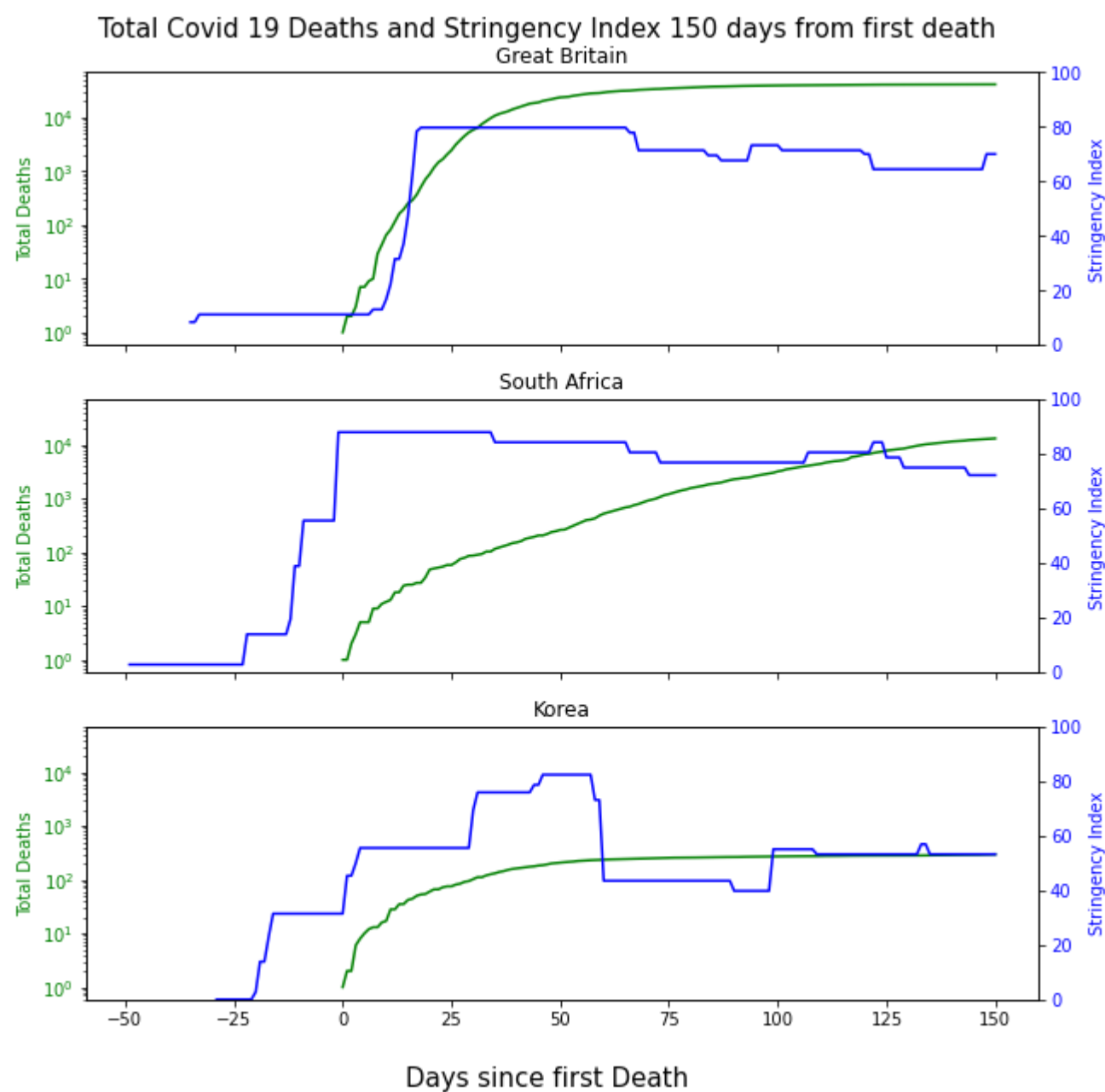
```
1 #figure
2 fig, axes = plt.subplots(nrows=3, ncols=1,figsize=(15,9),sharex = True)
3
4 plt.suptitle('Total Covid-19 cases, ICU patients, deaths and vaccinations per million in:\n United Kingdom,
5 fig.supxlabel('Days since 2020-12-03',size = 15, y = 0.05)
6 fig.supylabel('Total Numbers per Million',size = 15, x = 0.075)
7
8 #DATA
9 x = 'DaysSince3Dec20'
10 y1 = 'total_cases_per_million'
11 y2 = 'icu_patients_per_million'
12 y3 = 'total_deaths_per_million'
13 y4 = 'total_vaccinations_per_million'
14
15 data = [df_GBR,df_ZAF,df_KOR]
16
17 #subplot - axes
18 for z,i in enumerate(data) :
19     i.plot(x=x,y =[y1,y2,y3] ,ax=axes[z],legend =False,logy=True,color = ['orange','red','black','green'])
20     i.plot(x=x,y =y4 ,ax=axes[z],legend =False,logy=True,color = 'green',linestyle = '--')
21     axes[z].set_xlim(left=100,right =462)
22     axes[z].set(xlabel=None)
23     axes[z].axvline(x=228,ymin = 0, ymax = 1,color = 'pink',alpha = 0.4 )
24     axes[z].axvline(x=358,ymin = 0, ymax = 1,color = 'pink',alpha = 0.4 )
25     axes[z].spines["right"].set_visible(False)
26
27 #titles
28 axes[0].legend(ncol = 4, loc = (0,1.05),mode = 'expand')
29
30 axes[0].text(x=90,y=1000000,s = 'GBR', color = 'blue',weight='demi',alpha = 0.7)
31 axes[1].text(x=90,y=500000,s = 'ZAF', color = 'blue',weight='demi',alpha = 0.7)
32 axes[2].text(x=90,y=1000000,s = 'KOR', color = 'blue',weight='demi',alpha = 0.7)
33
34 #specific dates
35 axes[2].text(x=228,y=7000000,s = '2021-07-19', color = 'pink',weight='demi',alpha = 1)
36 axes[2].text(x=358,y=7000000,s = '2021-11-26', color = 'pink',weight='demi',alpha = 1);
37
38
```



Extra analysis


```
In [51]: 1 def make_dfs (x):
2         y = df[df['iso_code'].str.contains(x)]
3         y['date'] = y['date'].apply(pd.to_datetime)
4         y['Days Since 1st Death'] =(y['date']-(y.loc[df['total_deaths'] == 1, 'date'].iloc[0])).dt.days
5         y = y[y['Days Since 1st Death']<=150]
6         y = y[['iso_code','total_deaths','stringency_index','Days Since 1st Death']]
7         y = y[y['stringency_index'].notna()]
8         y = y.reset_index(drop=True)
9         return y
```

```
In [52]: 1 fig,(ax,ax1,ax2) =plt.subplots (3,1,figsize = (10,10), sharex = True,sharey = True)
2
3 plt.suptitle('Total Covid 19 Deaths and Stringency Index 150 days from first death',size=15,weight = 'ultra')
4 fig.supxlabel('Days since first Death',size = 15, y = 0.05)
5
6 isocodes = ['GBR','ZAF','KOR']
7 axes = [ax,ax1,ax2]
8
9 ax.set (title = 'Great Britain')
10 ax1.set (title = 'South Africa')
11 ax2.set (title = 'Korea')
12
13
14 for i,z in zip (isocodes,axes):
15
16     axa=z.twinx()
17     x = make_dfs(i) ['Days Since 1st Death']
18     y1 = make_dfs(i)['total_deaths']
19     y2 = make_dfs(i) ['stringency_index']
20
21     z.plot(x, y1, color='g')
22     axa.plot(x, y2, color='b')
23
24     z.set_ylabel( 'Total Deaths',color = 'g')
25     z.set(yscale = 'log')
26
27     axa.set_ylabel('Stringency Index', color = 'b')
28     axa.set_ylim(0,100)
29
30     def color_y_axis(ax, color):
31         for t in ax.get_yticklabels():
32             t.set_color(color)
33     color_y_axis(z, 'g')
34     color_y_axis(axa, 'b')
```



The Figure above provides insight into variation in government responses to the pandemic. Although most countries followed the same trend of getting stricter, I wanted to visualize variation between different responses.

Source: Our World in Data

The rate of which governments implemented stricter measures played a critical role in preventing covid deaths. The figure above compares the response of the British, South-Korean and South-African governments.

It is clear that the South Korean government who raised measures significantly even before the first death where then better able to combat the pandemic over the next 150 days, while in the UK the increase in the stringency of responses lagged the growth in deaths