

Title

Omar Sherif Fouda (Data Science Assignment 2)

Abstract

In this document, the assignment is divided into two parts: (part 1: matches results data)/(part 2: covid 19 data). Part 1 is depending on exercise 1, so I have the steps for the first exercise (codes) then the discussion and conclusion section in which I will be analyzing the data. Part 2 I haven't done from part 1; it is depending on exercise 2, so I have the steps for the second exercise (codes) then comes the discussion and conclusion section in which I will also have my analyses.

Steps

Exercise 1

```
In [1]: import pandas as pd

In [2]: df=pd.read_csv("results.csv")

In [3]: df.head(4)
```

	date	home_team	away_team	home_score	away_score	tournament	city	country	neutral
0	1872-11-30	Scotland	England	0	0	Friendly	Glasgow	Scotland	False
1	1873-03-08	England	Scotland	4	2	Friendly	London	England	False
2	1874-03-07	Scotland	England	2	1	Friendly	Glasgow	Scotland	False
3	1875-03-06	England	Scotland	2	2	Friendly	London	England	False

Adding a column win/lose/draw: We first create a list of conditions

```
In [4]: x=df['home_score']>df['away_score']
conditions =
(x>0),
(x==0),
(x<0)
]

In [5]: values= ['win','lose','draw']

In [6]: import numpy as np

In [7]: df['result']= np.select(conditions, values)

In [8]: df.head(3)
```

	date	home_team	away_team	home_score	away_score	tournament	city	country	neutral	result
0	1872-11-30	Scotland	England	0	0	Friendly	Glasgow	Scotland	False	draw
1	1873-03-08	England	Scotland	4	2	Friendly	London	England	False	lose
2	1874-03-07	Scotland	England	2	1	Friendly	Glasgow	Scotland	False	lose

```
In [9]: df['result'].value_counts()

Out[9]:
lose    21009
win     12224
draw     9955
Name: result, dtype: int64
```

```
In [10]: df['result'].value_counts(normalize=True)

Out[10]:
lose    0.486455
win     0.283042
draw     0.230503
Name: result, dtype: float64
```

```
In [11]: x=df['result'].value_counts()

In [12]: x=np.array(x)

In [13]: x

Out[13]:
array([21009, 12224, 9955])
```

```
In [14]: N=x.sum()

In [15]: N

Out[15]:
43188
```

```
In [16]: import statsmodels.api as sm
from statsmodels.stats.proportion import proportion_confint

In [17]: ci_win=proportion_confint(count=x[1],nobs=N,alpha=(1-.95))

In [18]: ci_win

Out[18]:
(0.2787935599544235, 0.28729011526083115)
```

```
In [19]: ci_lose=proportion_confint(count=x[0],nobs=N,alpha=(1-.95))

In [20]: ci_lose

Out[20]:
(0.481474905950987, 0.49116843552218753)
```

```
In [21]: ci_draw=proportion_confint(count=x[2],nobs=N,alpha=(1-.95))

In [22]: ci_draw

Out[22]:
(0.2265318471530234, 0.23447584016752862)
```

```
In [23]: df['country'].unique()

Out[23]:
array(['Scotland', 'England', 'Wales', 'Ireland', 'United States',
       'Uruguay', 'Austria', 'Hungary', 'Argentina', 'Belgium', 'France',
       'Netherlands', 'Poland', 'Switzerland', 'Sweden', 'Germany',
       'Italy', 'Chile', 'Norway', 'Finland', 'Luxembourg', 'Russia',
       'Denmark', 'Spain', 'Japan', 'Brazil', 'Paraguay', 'China PR',
       'Canada', 'Estonia', 'Guatemala', 'Czechoslovakia', 'Poland',
       'Yugoslavia', 'New Zealand', 'Romania', 'Latvia', 'Portugal',
       'Northern Ireland', 'Australia', 'Lithuania', 'Turkey', 'Mexico',
       'Kazakhstan', 'Soviet Union', 'Haiti', 'Philippines', 'Bulgaria',
       'Jamaica', 'Kenya', 'Czech Republic', 'Peru', 'Honduras',
       'Bhutan', 'Guayana', 'Uganda', 'El Salvador', 'Barbados',
       'Irish Free State', 'Trinidad and Tobago', 'Greenland', 'Cape Verde',
       'Egypt', 'Dominica', 'Gosadeloupe', 'Palestine',
       'Netherlands Guayana', 'French Guiana', 'Cuba', 'Colombia',
       'Saint Kitts and Nevis', 'Belize', 'Panama', 'Bolivia and Moravia',
       'Slovakia', 'Manchuria', 'Croatia', 'Costa Rica', 'Afghanistan',
       'Martinique', 'Southern Rhodesia', 'Iceland', 'Albania',
       'Madagascar', 'Northern Rhodesia', 'Tanganyika', 'Bosnia and Herzegovina',
       'French Somaliland', 'Belgian Congo', 'Mauritius', 'Hong Kong',
       'Vietnam', 'Macau', 'Republic of Ireland', 'Ethiopia', 'Burkina Faso',
       'Puerto Rico', 'Reunion', 'Israel', 'Sierra Leone', 'Mozambique',
       'Bolivia', 'Gold Coast', 'South Africa', 'Netherlands Antilles',
       'New Caledonia', 'Fiji', 'Nigeria', 'Vanuatu', 'Ceylon',
       'French Polynesia', 'Gambia', 'Singapore', 'Portuguese Guinea',
       'German DR', 'New Hebrides', 'Borneo', 'Samarang', 'Cambodia',
       'India', 'Labanon', 'Pakistan', 'Malaya', 'South Korea',
       'Vietnam Republic', 'Togo', 'Indonesia', 'Sudan', 'Mali', 'Syria',
       'Tunisia', 'Nyassaland', 'Ghana', 'Morocco', 'United Arab Republic',
       'North Korea', 'Bahamas', 'Guinea-Bissau', 'Mali Federation',
       'Mali', 'Vietnam DR', 'Cyprus', 'Iraq', 'Saint Lucia', 'Senegal',
       'Libya', 'Gabon', 'Thailand', 'Comoro', 'Tanzania', 'Grenada',
       'Guinea', 'Central African Republic', 'Cameroon', 'Algeria',
       'Seychelles', 'Ivory Coast', 'Laos', 'Cuba', 'Angola',
       'Jordan', 'Zambia', 'Saint Vincent and the Grenadines', 'Bermuda',
       'Niger', 'Malawi', 'DR Congo', 'Upper Volta', 'Taiwan', 'Guyana',
       'Mauritania', 'Sierra Leone', 'Saudi Arabia', 'Eswatini', 'Mozambique',
       'Papua New Guinea', 'Bahrain', 'Lesotho', 'Somalia', 'Saire',
       'Sri Lanka', 'Antigua and Barbuda', 'Faroe Islands', 'Qatar',
       'Sri Lanka', 'Burundi', 'Guam', 'Chad', 'Angola',
       'Dominican Republic', 'Seychelles', 'Sao Tome and Principe',
       'Botswana', 'Benin', 'Rwanda', 'Bangladesh',
       'United Arab Emirates', 'Zimbabwe', 'Oman', 'Equatorial Guinea',
       'Cape Verde', 'Liechtenstein', 'Nepal', 'Greenland',
       'Western Samoa', 'Belize', 'Brunei', 'Djibouti', 'Burkina Faso',
       'Yemen AR', 'Anguilla', 'Curaçao', 'Cayman Islands', 'Monaco',
       'Solomon Islands', 'Saint Martin', 'Nauru', 'Saint Kitts',
       'San Marino', 'Slovenia', 'Moldova', 'Ukraine', 'Mazakhstan',
       'Kazakhstan', 'Afghanistan', 'Turkmenistan', 'Georgia',
       'Kyrgyzstan', 'Armenia', 'Belarus', 'Guernsey', 'Azerbaijan',
       'North Macedonia', 'Jersey', 'Montserrat', 'Gibraltar', 'Myanmar',
       'Bosnia and Herzegovina', 'Tonga', 'Andorra', 'Yemen',
       'United States Virgin Islands', 'Palau', 'Cook Islands',
       'British Virgin Islands', 'Eritrea', 'Comoros', 'Micronesia',
       'Maldives', 'Laos', 'Isle of Man', 'Samoa',
       'Serbia and Montenegro', 'Mayotte', 'Bangladesh', 'Northern Cyprus',
       'Sri Lanka', 'Montenegro', 'Northern Mariana Islands', 'Timor',
       'North and Central Islands', 'South Sudan', 'Kosovo', 'East Timor',
       'Taiwan', dtype=object)
```

```
In [24]: dfegy=df[df['country']=='Egypt']

In [25]: dfegy.head(1)
```

	date	home_team	away_team	home_score	away_score	tournament	city	country	neutral	result
1463	2/19/1932	Egypt	Hungary	0	0	Friendly	Cairo	Egypt	False	draw
1661	3/16/1934	Egypt	Israel	7	1	FIFA World Cup qualification	Cairo	Egypt	False	lose
1895	6/19/1936	Egypt	Greece	3	1	Friendly	Cairo	Egypt	False	lose
2927	12/24/1948	Egypt	Norway	1	1	Friendly	Cairo	Egypt	False	draw
3080	2/17/1950	Egypt	Greece	2	0	Friendly	Cairo	Egypt	False	lose

Create a variable type of the match 'friendly/No friendly' and 'home/away match'

```
In [26]: conditions = [
dfegy['tournament']=='Friendly',
dfegy['tournament']=='Friendly'
]

In [27]: values=['Friendly','Official']

In [28]: dfegy['typematch']= np.select(conditions, values)

/var/folders/gv/41kfxw90v7gdm2ntzqf4hz40000gn/71pykernel.71081/2597923672.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
dfegy['typematch']= np.select(conditions, values)
```

```
In [29]: dfegy['typematch'].value_counts()

Out[29]:
Official    226
Friendly    181
Name: typematch, dtype: int64
```

Comparing the probability of win/lose/draw by type of match

```
In [30]: x=npd.crosstab(dfegy['typematch'],dfegy['result'],margins=True)
x

Out[30]:
      result draw lose win  All
typematch
Friendly    37    99  45  181
Official    40   139  47  226
All         77   238  92  407
```

```
In [31]: x=np.array(x)

Out[31]:
array([[ 37,  99,  45, 181],
       [ 40, 139,  47, 226],
       [ 77, 238,  92, 407]])
```

```
In [32]: CI_egywin_friendly=proportion_confint(count=x[0,2],nobs=x[0,3],alpha=(1-.95))
CI_egywin_friendly

Out[32]:
(0.183652804868489, 0.31158476521246)
```

```
In [33]: CI_egywin_official=proportion_confint(count=x[1,2],nobs=x[1,3],alpha=(1-.95))
CI_egywin_official

Out[33]:
(0.1550517855722795, 0.260874179675435)
```

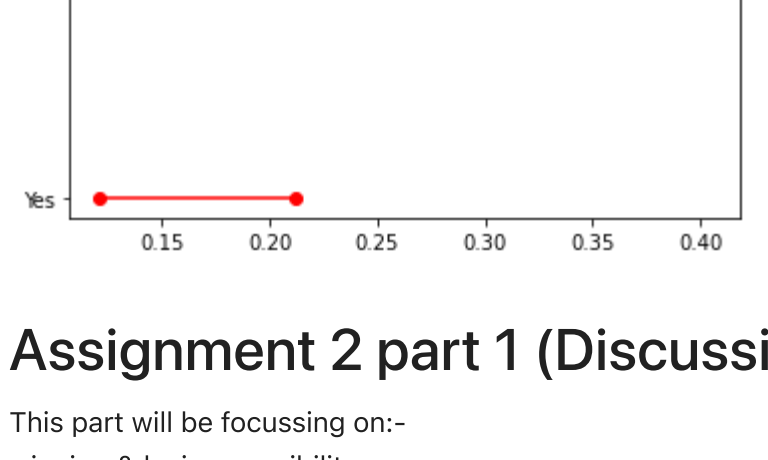
Plotting the Confidence intervals

```
In [34]: ci_egywin = {}
ci_egywin['typematch'] = ['Friendly','Official']
ci_egywin['lb'] = [CI_egywin_friendly[0],CI_egywin_official[0]]
ci_egywin['ub'] = [CI_egywin_friendly[1],CI_egywin_official[1]]
df_ci = pd.DataFrame(ci_egywin)

Out[34]:
  typematch    lb    ub
0  Friendly  0.185653  0.311585
1  Official  0.155052  0.260877
```

```
In [35]: import matplotlib.pyplot as plt
for lb,ub,y in zip(df_ci['lb'],df_ci['ub'],range(len(df_ci))):
    plt.plot((lb,ub),(y,y),'r')
plt.xticks(range(len(df_ci)),list(df_ci['typematch']))

Out[35]:
(<matplotlib.axis.YTick at 0x7f83b1448640>,
 <matplotlib.axis.YTick at 0x7f83b1437e00>,
 (Text(0, 0, 'Friendly'), Text(0, 1, 'Official')))
```



Let's now check out the impact of fans on the match results. We will compare the probability of winning between home matches and away matches.

```
In [36]: dfegy['home']=dfegy['home_team']=='Egypt'

/var/folders/gv/41kfxw90v7gdm2ntzqf4hz40000gn/71pykernel.71081/217501764.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
dfegy['home']=dfegy['home_team']=='Egypt'
```

```
In [37]: dfegy['home'].value_counts()

Out[37]:
True     258
False    149
Name: home, dtype: int64
```

```
In [38]: x=npd.crosstab(dfegy['home'],dfegy['result'],margins=True)
x

Out[38]:
      result draw lose win  All
home
False    33    67  49  149
True     44   171  43  258
All      77   238  92  407
```

```
In [39]: dfegy.head(4)
```

	date	home_team	away_team	home_score	away_score	tournament	city	country	neutral	result	typematch	home
1463	2/19/1932	Egypt	Hungary	0	0	Friendly	Cairo	Egypt	False	draw	Friendly	True
1661	3/16/1934	Egypt	Israel	7	1	FIFA World Cup qualification	Cairo	Egypt	False	lose	Official	True
1895	6/19/1936	Egypt	Greece	3	1	Friendly	Cairo	Egypt	False	lose	Friendly	True
2927	12/24/1948	Egypt	Norway	1	1	Friendly	Cairo	Egypt	False	draw	Friendly	True

```
In [40]: x=np.array(x)

Out[40]:
array([[ 33,  67,  49, 149],
       [ 44, 171,  43, 258],
       [ 77, 238,  92, 407]])
```

```
In [41]: CI_egywin_home=proportion_confint(count=x[1,2],nobs=x[1,3],alpha=(1-.95))
CI_egywin_home

Out[41]:
(0.12119174183927744, 0.21214159149405587)
```

```
In [42]: CI_egywin_away=proportion_confint(count=x[0,2],nobs=x[0,3],alpha=(1-.95))
CI_egywin_away

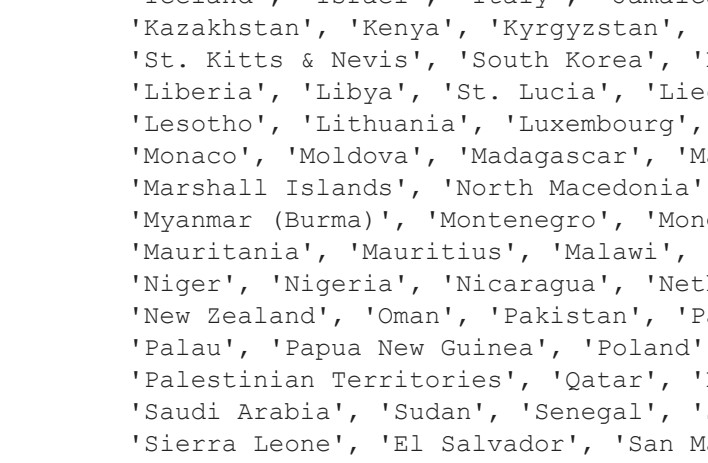
Out[42]:
(0.2534251314484364, 0.4042529395852553)
```

```
In [43]: ci_egywin = {}
ci_egywin['home'] = ['Yes','No']
ci_egywin['lb'] = [CI_egywin_home[0],CI_egywin_away[0]]
ci_egywin['ub'] = [CI_egywin_home[1],CI_egywin_away[1]]
df_ci = pd.DataFrame(ci_egywin)

Out[43]:
  home    lb    ub
0  Yes  0.121192  0.212142
1  No   0.253425  0.404293
```

```
In [44]: for lb,ub,y in zip(df_ci['lb'],df_ci['ub'],range(len(df_ci))):
    plt.plot((lb,ub),(y,y),'r')
plt.xticks(range(len(df_ci)),list(df_ci['home']))

Out[44]:
(<matplotlib.axis.YTick at 0x7f83b15542d0>,
 <matplotlib.axis.YTick at 0x7f83b154db20>,
 (Text(0, 0, 'Yes'), Text(0, 1, 'No')))
```



Assignment 2 part 1 (Discussion & Conclusion)

This winning is focusing on-
winning & losing possibility
difference and changes in home and away matches
presenting confidence interval and graphs

win and lose possibility

```
In [45]: x=npd.crosstab(dfegy['typematch'],dfegy['result'],margins=True)
x

Out[45]:
      result draw lose win  All
typematch
Friendly    37    99  45  181
Official    40   139  47  226
All         77   238  92  407
```

This table shows the stats of egypt in all games.
Differing them by the type of the competition.
20% is the winning probability in official games.
25% is the winning probability in the friendly matches.
Egypt plays better in the friendly games.

home and away stats

```
In [46]: x=npd.crosstab(dfegy['home'],dfegy['result'],margins=True)
x

Out[46]:
      result draw lose win  All
home
False    33    67  49  149
True     44   171  43  258
All      77   238  92  407
```

This table shows the stats of egypt in all games.
Differing them by home and away games (true: home/false: away).
obviously, egypt's probability of winning in the away matches is much higher.

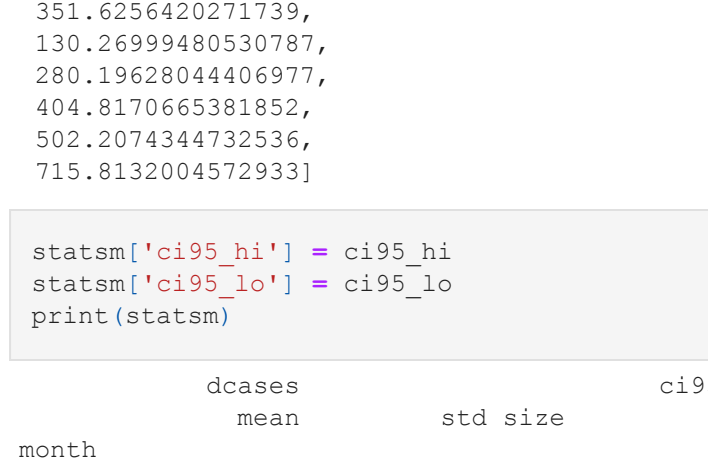
confidence interval and graphs

```
In [47]: ci_egywin = {}
ci_egywin['typematch'] = ['Friendly','Official']
ci_egywin['lb'] = [CI_egywin_friendly[0],CI_egywin_official[0]]
ci_egywin['ub'] = [CI_egywin_friendly[1],CI_egywin_official[1]]
df_ci = pd.DataFrame(ci_egywin)

Out[47]:
  typematch    lb    ub
0  Friendly  0.185653  0.311585
1  Official  0.155052  0.260877
```

```
In [48]: import matplotlib.pyplot as plt
for lb,ub,y in zip(df_ci['lb'],df_ci['ub'],range(len(df_ci))):
    plt.plot((lb,ub),(y,y),'r')
plt.xticks(range(len(df_ci)),list(df_ci['typematch']))

Out[48]:
(<matplotlib.axis.YTick at 0x7f83b15542d0>,
 <matplotlib.axis.YTick at 0x7f83b154db20>,
 (Text(0, 0, 'Yes'), Text(0, 1, 'No')))
```



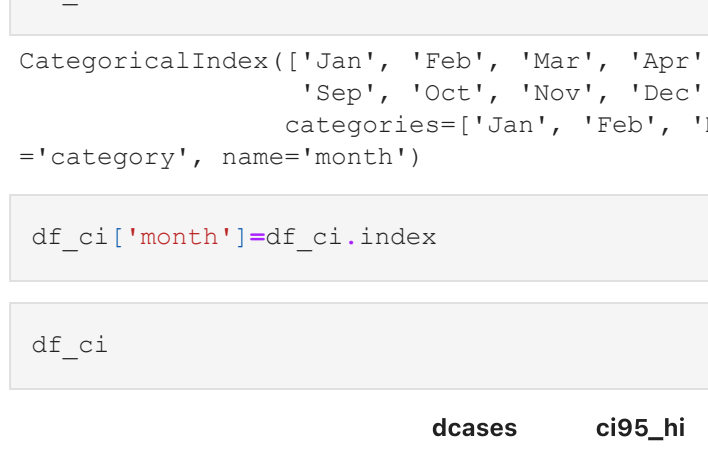
Here is the confidence interval of the probability of winning official and friendly matches and under it a graph representing it.
As shown if we have had more matches still egypt has the advantage in the friendly matches.

```
In [49]: ci_egywin = {}
ci_egywin['home'] = ['Yes','No']
ci_egywin['lb'] = [CI_egywin_home[0],CI_egywin_away[0]]
ci_egywin['ub'] = [CI_egywin_home[1],CI_egywin_away[1]]
df_ci = pd.DataFrame(ci_egywin)

Out[49]:
  home    lb    ub
0  Yes  0.121192  0.212142
1  No   0.253425  0.404293
```

```
In [50]: for lb,ub,y in zip(df_ci['lb'],df_ci['ub'],range(len(df_ci))):
    plt.plot((lb,ub),(y,y),'r')
plt.xticks(range(len(df_ci)),list(df_ci['home']))

Out[50]:
(<matplotlib.axis.YTick at 0x7f83b1636340>,
 <matplotlib.axis.YTick at 0x7f83b163db00>,
 (Text(0, 0, 'Yes'), Text(0, 1, 'No')))
```



Here is the confidence interval of the probability of winning home and away matches and under it a graph representing it.
As shown if we have had more matches still egypt has the advantage in the away matches.

Exercise 2

```
In [51]: import pandas as pd

In [52]: df=pd.read_csv("covid_data.csv",encoding='latin-1')

In [53]: df.head()
```

	date	iso3c	country	income	region	continent	dcases	deaths	population	weekdays	month
0	2020-02-24	AFG	Afghanistan	Low income	South Asia	Asia	5	0	38041754	Mon	Feb
1	2020-02-25	AFG	Afghanistan	Low income	South Asia	Asia	0	0	38041754	Tue	Feb
2	2020-02-26	AFG	Afghanistan	Low income	South Asia	Asia	0	0	38041754	Wed	Feb
3	2020-02-27	AFG	Afghanistan	Low income	South Asia	Asia	0	0	38041754	Thu	Feb
4	2020-02-28	AFG	Afghanistan	Low income	South Asia	Asia	0	0	38041754	Fri	Feb

```
In [54]: df['country'].unique()

Out[54]:
array(['Afghanistan', 'Angola', 'Albania', 'Andorra',
       'United Arab Emirates', 'Argentina', 'Armenia',
       'Antigua & Barbuda', 'Australia', 'Austria', 'Azerbaijan',
       'Burundi', 'Belgium', 'Benin', 'Burkina Faso', 'Bangladesh',
       'Bulgaria', 'Bahrain', 'Bhutan', 'Botswana', 'Bosnia & Herzegovina',
       'Belarus', 'Belize', 'Bolivia', 'Brazil', 'Barbados', 'Brunei',
       'Bhutan', 'Botswana', 'Central African Republic', 'Canada',
       'Switzerland', 'Chile', 'China', 'Cote d'Ivoire', 'Cameroon',
       'Congo - Kinshasa', 'Congo - Brazzaville', 'Colombia', 'Comoros',
       'Cape Verde', 'Costa Rica', 'Cuba', 'Cyprus', 'Czechia', 'Germany',
       'Guinea', 'Guinea-Bissau', 'Guatemala', 'Guinea', 'Guyana', 'Honduras', 'Hungary',
       'Ecuador', 'Egypt', 'El Salvador', 'Ethiopia', 'Ecuador', 'Finland',
       'Fiji', 'France', 'Gabon', 'United Kingdom', 'Georgia', 'Ghana',
       'Guinea', 'Gambia', 'Guinea-Bissau', 'Equatorial Guinea', 'Greenland',
       'Grenada', 'Guatemala', 'Guyana', 'Honduras', 'Croatia', 'Haiti',
       'Hungary', 'Indonesia', 'India', 'Ireland', 'Iran', 'Iraq',
       'Israel', 'Israel', 'Italy', 'Jamaica', 'Jordan', 'Japan',
       'Kazakhstan', 'Kenya', 'Kyrgyzstan', 'Cambodia', 'Kiribati',
       'St. Kitts & Nevis', 'South Korea', 'Kuwait', 'Laos', 'Lebanon',
       'Liberia', 'Libya', 'Liechtenstein', 'Lithuania', 'Luxembourg', 'Latvia', 'Mali',
       'Lesotho', 'Lithuania', 'Luxembourg', 'Latvia', 'Morocco',
       'Monaco', 'Moldova', 'Madagascar', 'Maldives', 'Mexico',
       'Marshall Islands', 'North Macedonia', 'Mali', 'Malta',
       'Myanmar (Burma)', 'Mauritius', 'Mongolia', 'Mozambique',
       'Mauritania', 'Mauritius', 'Malawi', 'Malaysia', 'Mambila',
       'Niger', 'Nigeria', 'Nicaragua', 'Netherlands', 'New Zealand', 'New Zealand', 'Oman', 'Pakistan', 'Panama', 'Peru', 'Philippines',
       'Palau', 'Papua New Guinea', 'Poland', 'Portugal', 'Paraguay',
       'Palestinian Territories', 'Poland', 'Romania',
       'Saudi Arabia', 'Sudan', 'Senegal', 'Singapore', 'Solomon Islands',
       'Sierra Leone', 'El Salvador', 'San Marino', 'Somalia', 'Serbia',
       'South Sudan', 'Sao Tome & Principe', 'Suriname', 'Slovakia',
       'Slovenia', 'Sweden', 'Eswatini', 'Seychelles', 'Syria', 'Chad',
       'Togo', 'Thailand', 'Tajikistan', 'Timor-Leste',
       'Trinidad and Tobago', 'Tunisia', 'Turkey', 'Tanzania', 'Uganda',
       'Ukraine', 'Ukraine', 'United States', 'Uzbekistan',
       'St. Vincent & the Grenadines', 'Venezuela', 'Vietnam', 'Vanuatu',
       'Samoa', 'Yemen', 'South Africa', 'Zambia', 'Zimbabwe',
       dtype=object)
```

```
In [55]: from pandas.api.types import CategoricalDtype
cases='cat',cat='name','True','Wed','Thu'
cat_type = CategoricalDtype(categories=cases, ordered=True)
df['weekdays']=df['weekdays'].astype(cat_type)
```

```
In [56]: from pandas.api.types import CategoricalDtype
cases='cat',cat='name','Apr','May','Jun','Jul','Aug','Sep','Oct','Nov','Dec'
cat_type = CategoricalDtype(categories=cases, ordered=True)
df['month']=df['month'].astype(cat_type)
```

```
In [57]: import numpy as np

In [58]: dfegy=df[df['country']=='Egypt']

In [59]: stats=dfegy.groupby('weekdays').agg(['dcases': (np.mean, np.std, np.size)])

In [60]: stats
```

	dcases	mean	std	size
weekdays				
Fri	567.161616	428.532849	99	
Sat	558.80122	422.818365	98	
Sun	545.520408	422.358748	98	
Mon	561.46839	442.13749	98	
Tue	566.153061	418.125460	98	
Wed	561.479592	406.357812	98	
Thu	567.683673	410.020004	98	

```
In [61]: stats=dfegy.groupby('month').agg(['dcases': (np.mean, np.std, np.size)])
stats
```

	dcases	mean	std	size
month				
Jan	899.645161	251.335970	31	
Feb	374.409091	288.087868	44	
Mar	329.290323	310.292849	62	
Apr	504.133333	354.264665	60	
May	879.774194	346.475245	62	
Jun	1032.633333	454.674216	60	
Jul	483.661290	437.595552	62	
Aug	145.806452	60.683224	62	
Sep	339.033333	225.855584	60	
Oct	497.580645	362.315487	62	
Nov	598.616667	370.082496	60	
Dec	792.709677	300.346870	62	

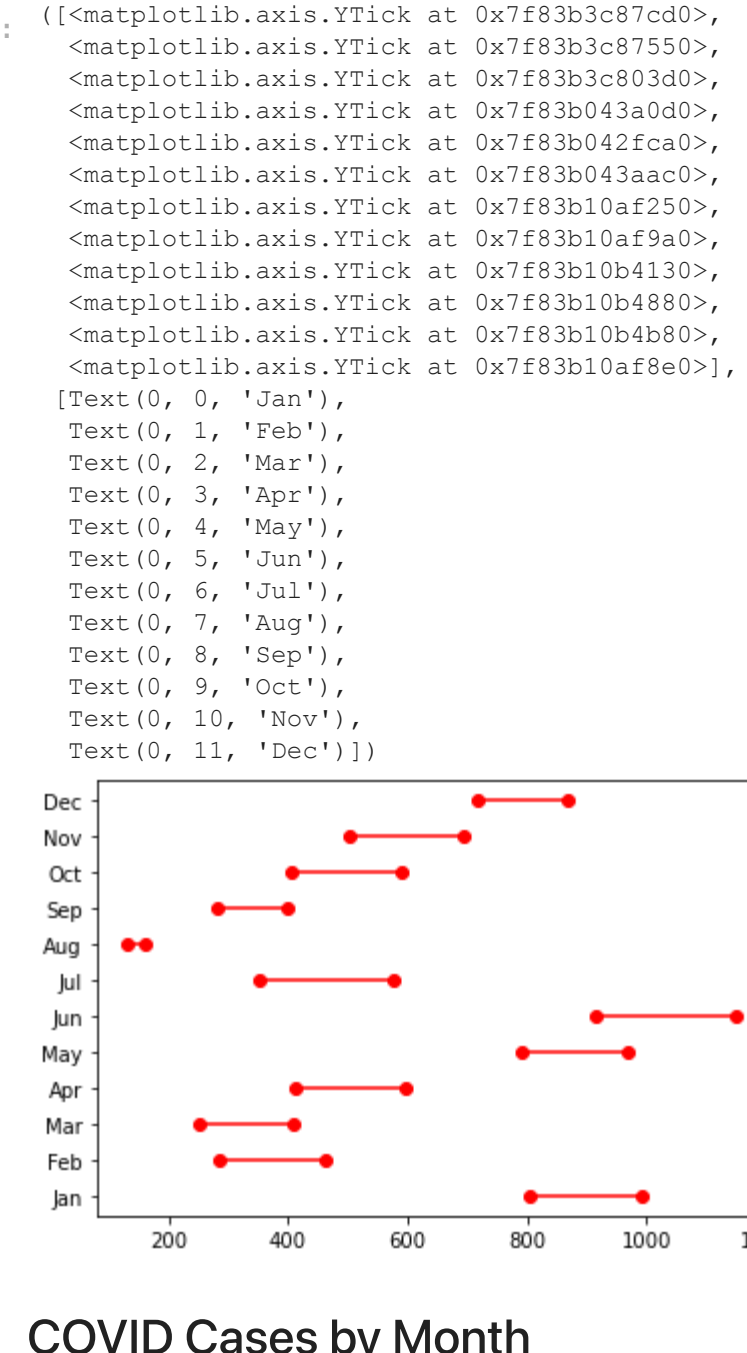
```
In [62]: import numpy as np

In [63]: statm.index

Out[63]:
CategoricalIndex(['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug',
                  'Sep', 'Oct', 'Nov', 'Dec'],
                  dtype='object', name='month')
```

```
In [64]: ci95_hi = []
ci95_lo = []

In [65]: for i in statm.index:
    n
```

COVIDCases by Month

```
In [77]: import pandas as pd

In [78]: df=pd.read_csv('covid_data.csv',encoding='latin-1')

In [79]: from pandas.api.types import CategoricalDtype
cdate=['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sep', 'Oct', 'Nov', 'Dec']
cat_type = CategoricalDtype(categories=cdate, ordered=True)
df['month'] = df['month'].astype(cat_type)

In [80]: dfegy=df[df['country']=='Egypt']

In [81]: import numpy as np

In [82]: stats=dfegy.groupby('month').agg(['dcasees': [np.mean, np.std, np.size]])

In [83]: stats=stats.reset_index()

In [84]: stats.shape

Out[84]: (12, 3)

In [85]: stats.columns

Out[85]: MultiIndex([('dcasees', 'mean'),
              ('dcasees', 'std'),
              ('dcasees', 'size')])

In [86]: stats.columns=['mean','std','size']

In [87]: stats.columns

Out[87]: Index(['mean', 'std', 'size'], dtype='object')

In [88]: def get_ci_lb(x, alpha=0.05):
    sample_size=len(x)
    sample_mean=np.mean(x)
    sample_size=len(x)
    margin_of_error = t.ppf(1 - alpha/2, sample_size-1)*sample_size/np.sqrt(sample_size-1)
    return sample_mean + margin_of_error

In [89]: x=dfegy['dcasees']

In [90]: from scipy.stats import norm,t

In [91]: get_ci_lb(x)

Out[91]: 529.794591276133

In [92]: def get_ci_lb(x, alpha=0.05):
    sample_size=len(x)
    sample_mean=np.mean(x)
    sample_size=len(x)
    margin_of_error = t.ppf(1 - alpha/2, sample_size-1)*sample_size/np.sqrt(sample_size-1)
    return sample_mean + margin_of_error

In [93]: get_ci_lb(x)

Out[93]: 592.49449184819

In [94]: ci_dcasees=stats.groupby('month').agg(['dcasees': [np.mean, np.std, np.size, get_ci_lb, get_ci_ub]])
stats

Out[94]:
```

dcasees						
month	mean	std	size	get_ci_lb	get_ci_ub	
Jan	899.645161	251.335970	31	807.454383	991.835939	
Feb	374.409091	288.087868	44	286.822374	461.995808	
Mar	329.290323	310.292849	62	250.490604	408.089951	
Apr	504.133333	354.246465	60	412.61768	595.644899	
May	879.774194	346.475245	62	791.785956	967.762431	
Jun	1032.633333	454.674216	60	915.176529	1150.088138	
Jul	463.661290	437.595552	62	352.532828	574.789752	
Aug	145.806452	60.683224	62	130.395708	161.217105	
Sep	339.033333	225.855634	60	280.688649	397.378017	
Oct	497.580645	362.321547	62	405.568201	589.593089	
Nov	598.616667	370.082496	60	510.014220	694.219113	
Dec	792.709677	300.346870	62	716.435854	868.983501	

```
In [95]: statsm.index

Out[95]: CategoricalIndex(['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug',
                  'Sep', 'Oct', 'Nov', 'Dec'],
                  categories=['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', ...], ordered=True, dtype=
                  <category>, name='month', freq=None)

In [96]: statsm.columns=['mean','std','size','lb','ub']

In [97]: statsm['month']=statsm.index

In [98]: statsm

Out[98]:
```

month	mean	std	size	lb	ub	month
Jan	899.645161	251.335970	31	807.454383	991.835939	Jan
Feb	374.409091	288.087868	44	286.822374	461.995808	Feb
Mar	329.290323	310.292849	62	250.490604	408.089951	Mar
Apr	504.133333	354.246465	60	412.61768	595.644899	Apr
May	879.774194	346.475245	62	791.785956	967.762431	May
Jun	1032.633333	454.674216	60	915.176529	1150.088138	Jun
Jul	463.661290	437.595552	62	352.532828	574.789752	Jul
Aug	145.806452	60.683224	62	130.395708	161.217105	Aug
Sep	339.033333	225.855634	60	280.688649	397.378017	Sep
Oct	497.580645	362.321547	62	405.568201	589.593089	Oct
Nov	598.616667	370.082496	60	510.014220	694.219113	Nov
Dec	792.709677	300.346870	62	716.435854	868.983501	Dec

```
In [99]: import matplotlib.pyplot as plt

In [100]: plt.plot('month', 'mean', data=statsm, marker='o', color='black', markersize=4, linewidth=1, linestyle='--')
plt.plot('month', 'mean', data=statsm, marker='o', color='black', markersize=4, linewidth=1, linestyle='--')
plt.xlabel('month')
plt.ylabel('Average Daily COVID CASES')
plt.show()
```



We will do the same graph now but only on the year 2021. We will first create the column 'year'. It works as follows:

1. We first tell to Python that the 'date' column is composed of dates object
2. We can then tell to Python to extract the year from the date

```
In [101]: df['date']=[0]

In [102]: '2020-02-24'

In [103]: df['date']= pd.to_datetime(df['date'],format='%Y-%m-%d')

In [104]: df['date']

Out[104]: Timestamp('2020-02-24 00:00:00')

In [105]: df['year']=[0]

In [106]: df['year']

Out[106]: 2020

I compute now, by country, by year, by month the following statistics on the daily COVID cases: mean, std, size, CI(95%) LB, and UB.
```

```
In [106]: statsdcasees=df.groupby(['country','year','month']).agg(['dcasees': [np.mean, np.std, np.size, get_ci_lb, get_ci_ub]])

In [107]: statsdcasees

Out[107]:
```

dcasees						
country	year	month	mean	std	size	get_ci_lb
Afghanistan	2020	Jan	NaN	NaN	NaN	NaN
		Feb	0.833333	2.047241	6.0	-1.308818
		Mar	5.258065	10.871883	31.0	1.270225
		Apr	55.366667	40.385627	30.0	40.286426
		May	430.741935	266.692078	31.0	332.918491
Zimbabwe	2021	Aug	513.322581	386.841948	31.0	371.427809
		Sep	201.566667	135.119789	30.0	151.112108
		Oct	69.580645	58.035492	31.0	48.293055
		Nov	54.933333	82.622087	30.0	24.081739
		Dec	2536.548387	2572.199964	31.0	1593.057823

4488 rows x 5 columns

Transform the data index to columns

```
In [108]: statsdcasees=statsdcasees.reset_index()

In [109]: statsdcasees

Out[109]:
```

country	year	month	mean	std	size	get_ci_lb	get_ci_ub
Afghanistan	2020	Jan	NaN	NaN	NaN	NaN	NaN
		Feb	0.833333	2.047241	6.0	-1.308818	2.975485
		Mar	5.258065	10.871883	31.0	1.270225	9.245904
		Apr	55.366667	40.385627	30.0	40.286426	70.446908
		May	430.741935	266.692078	31.0	332.918491	528.565379
Zimbabwe	2021	Aug	513.322581	386.841948	31.0	371.427809	655.217353
		Sep	201.566667	135.119789	30.0	151.112108	252.021225
		Oct	69.580645	58.035492	31.0	48.293055	90.868235
		Nov	54.933333	82.622087	30.0	24.081739	85.784928
		Dec	2536.548387	2572.199964	31.0	1593.057823	3480.038951

4488 rows x 8 columns

Filter now the Egypt data on 2021

```
In [110]: statsEgy=statsdcasees[(statsdcasees['country']=='Egypt') & (statsdcasees['year']==2021)]

In [111]: statsEgy

Out[111]:
```

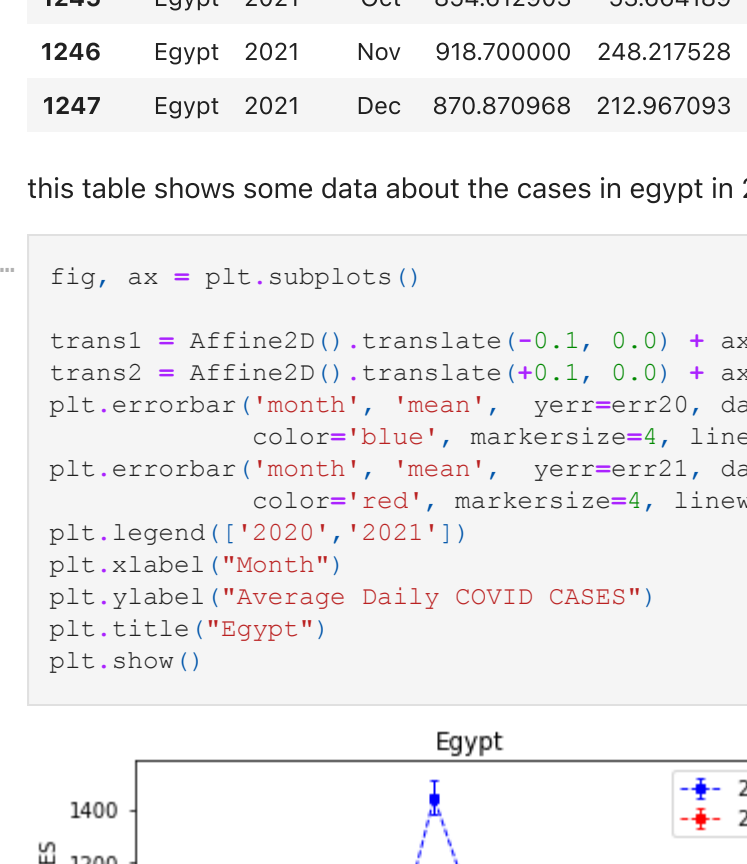
country	year	month	mean	std	size	get_ci_lb	get_ci_ub
Egypt	2021	Jan	899.645161	251.335970	31.0	807.454383	991.835939
		Feb	588.321429	40.988659	28.0	572.427298	604.215559
		Mar	635.709677	36.309038	31.0	622.391415	649.027940
		Apr	847.366667	92.416405	30.0	812.857814	861.875520
		May	1132.193548	55.087457	31.0	1111.987306	1152.399791
		Jun	621.066667	197.770841	30.0	547.219281	694.915513
		Jul	96.129032	61.659139	31.0	73.512277	118.745787
		Aug	134.806452	73.635326	31.0	107.796796	161.816107
		Sep	536.100000	151.276488	30.0	479.612431	592.587569
		Oct	854.612903	53.664189	31.0	834.928720	874.297087
		Nov	918.700000	248.217528	30.0	826.014052	1011.385948
		Dec	870.870968	212.967093	31.0	792.754007	948.987928

```
In [112]: MultiIndex(['country',
                  ('year', ''),
                  ('month', ''),
                  ('dcasees', 'mean'),
                  ('dcasees', 'std'),
                  ('dcasees', 'size'),
                  ('dcasees', 'get_ci_lb'),
                  ('dcasees', 'get_ci_ub')])

In [113]: statsEgy.columns=['country','year','month','mean','std','size','lb','ub']

Draw then the average daily COVID cases by month in Egypt.
```

```
In [114]: plt.plot('month', 'mean', data=statsEgy, marker='o', color='black', markersize=4, linewidth=1, linestyle='--')
plt.plot('month', 'mean', data=statsEgy, marker='o', color='black', markersize=4, linewidth=1, linestyle='--')
plt.xlabel('month')
plt.ylabel('Average Daily COVID CASES')
plt.title('Egypt')
plt.show()
```



Let's compare now between 2020 and 2021

```
In [118]: statsEgy21=statsdcasees[(statsdcasees['country']=='Egypt') & (statsdcasees['year']==2021)]
statsEgy21

Out[118]:
```

country	year	month	mean	std	size	get_ci_lb	get_ci_ub
Egypt	2021	Jan	899.645161	251.335970	31.0	807.454383	991.835939
		Feb	588.321429	40.988659	28.0	572.427298	604.215559
		Mar	635.709677	36.309038	31.0	622.391415	649.027940
		Apr	847.366667	92.416405	30.0	812.857814	861.875520
		May	1132.193548	55.087457	31.0	1111.987306	1152.399791
		Jun	621.066667	197.770841	30.0	547.219281	694.915513
		Jul	96.129032	61.659139	31.0	73.512277	118.745787
		Aug	134.806452	73.635326	31.0	107.796796	161.816107
		Sep	536.100000	151.276488	30.0	479.612431	592.587569
		Oct	854.612903	53.664189	31.0	834.928720	874.297087
		Nov	918.700000	248.217528	30.0	826.014052	1011.385948
		Dec	870.870968	212.967093	31.0	792.754007	948.987928

```
In [119]: statsEgy20=statsdcasees[(statsdcasees['country']=='Egypt') & (statsdcasees['year']==2020)]
statsEgy20

Out[119]:
```

country	year	month	mean	std	size	get_ci_lb	get_ci_ub
Egypt	2020	Jan	NaN	NaN	NaN	NaN	NaN
		Feb	0.062500	0.250000	16.0	-0.070716	0.195716
		Mar	22.870968	20.228267	31.0	15.451179	30.290756
		Apr	160.900000	54.996144	30.0	140.364102	181.435898
		May	627.354839	330.723908	31.0	506.044330	748.665347
		Jun	1444.200000	176.114853	30.0	1376.437633	1509.962367
		Jul	831.193548	326.183427	31.0	711.548501	950.835959
		Aug	156.806452	42.607839	31.0	141.777770	172.435133
		Sep	141.966667	23.389481	30.0	133.232891	158.700442
		Oct	140.548387	24.650678	31.0	131.506445	149.590329
		Nov	278.533333	71.137426	30.0	251.970182	305.096485
		Dec	714.548387	354.179036	31.0	584.634468	844.462306

this table shows some data about the cases in egypt in 2020

```
In [130]: statsEgy=statsdcasees[(statsdcasees['country']=='Egypt') & (statsdcasees['year']==2021)]
statsEgy

Out[130]:
```

country	year	month	mean	std	size	get_ci_lb	get_ci_ub
Egypt	2021	Jan	899.645161	251.335970	31.0	807.454383	991.835939
		Feb	588.321429	40.988659	28.0	572.427298	604.215559
		Mar	635.709677	36.309038	31.0	622.391415	649.027940
		Apr	847.366667	92.416405	30.0	812.857814	861.875520
		May	1132.193548	55.087457	31.0	1111.987306	1152.399791
		Jun	621.066667	197.770841	30.0	547.219281	694.915513
		Jul	96.129032	61.659139	31.0	73.512277	118.745787
		Aug	134.806452	73.635326	31.0	107.796796	161.816107
		Sep	536.100000	151.276488	30.0	479.612431	592.587569
		Oct	854.612903	53.664189	31.0	834.928720	874.297087
		Nov	918.700000	248.217528	30.0	826.014052	1011.385948
		Dec	870.870968	212.967093	31.0	792.754007	948.987928

this table shows some data about the cases in egypt in 2021

```
In [131]: fig, ax = plt.subplots()

trans1 = Affine2D().translate(-0.1, 0.0) + ax.transData
trans2 = Affine2D().translate(-0.1, 0.0) + ax.transData
plt.errorbar('month', 'mean', yerr=err2, data=statsEgy20, marker='s', capsize=2,
            color='blue', markersize=4, linewidth=1, linestyle='--', transform=trans1)
plt.errorbar('month', 'mean', yerr=err2, data=statsEgy21, marker='s', capsize=2,
            color='red', markersize=4, linewidth=1, linestyle='--', transform=trans2)
plt.legend(['2020', '2021'])
plt.xlabel('Month')
plt.ylabel('Average Daily COVID CASES')
plt.title('Egypt')
plt.show()
```


This graph shows the difference between 2020 & 2

		ddeaths	
month	mean	std	size
Jan	54.354839	3.638208	31
Feb	31.181818	24.163467	44
Mar	21.822581	20.717083	62
Apr	28.166667	17.826058	60
May	37.483871	20.947821	62
Jun	51.116667	23.279504	60
Jul	35.596774	27.745010	62
Aug	13.354839	7.503216	62
Sep	18.400000	7.962199	60
Oct	26.709677	16.823468	62
Nov	36.783333	27.679211	60
Dec	36.435484	15.919491	62

As shown, Egypt had less cases than Belgium and that could mean one of three things:-
Egypt is doing better in controlling the virus
Europe has more of this virus (based on the region)
Egypt is not collecting enough data and that is the most reliable option because of the difference in the population
However, Egypt always had more death numbers and that means that the health organization in egypt is bad; by this information we know that egypt is not doing better in controlling the virus

Thank You!