# EXPLORATORY DATA ANALYSIS

## Project: Analysis of different measures taken by different countries to control the spread of Covid-19 virus

### Team Members
Nisha Ramrakhyani (Student Id: 801208678)
Punit Mashruwala (Student Id: 801208416)
Zalak Panchal (Student Id: 801196881)

**Exploratory Data Analysis** is a process having a set of techniques for analyzing datasets, performing initial investigations on data to discover patterns, spot errors using graphical representations.
The steps involved in Exploratory Data Analysis are-
a) Importing and cleaning data and checking for errors and other special conditions.
b) Exploring Variables one at a time, visualizing distributions
c) Exploring relationships between variables two at a time, using scatter plots and other visualizations
d) Exploring multivariate reltionships using multiple regression and logistic regression.
We perform Exploratory Data Analysis to use data to answer questions and guide in decision making.
We perform Exploratory Data Analysis especially in the early stages of a project, or while working with a new dataset.

In [1]:
```python
#Importing libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from bokeh.io import output_notebook, curdoc
from bokeh.plotting import figure, show
from bokeh.models import ColumnDataSource, HoverTool, CategoricalColorMapper,
output_notebook()
```

(https://bokeh.org) BokehJS 2.2.1 successfully loaded.

In [3]:

```python
#Importing datatset
dataset = pd.read_excel('covid_19_data_2.xlsx')
print(dataset.head())
```

```
   ID  ISO    COUNTRY   REGION                           LOG_TYPE  \
0  98  AUS  Australia  Pacific  Introduction / extension of measures
1  96  AUS  Australia  Pacific  Introduction / extension of measures
2  93  AUS  Australia  Pacific  Introduction / extension of measures
3  94  AUS  Australia  Pacific  Introduction / extension of measures
4  97  AUS  Australia  Pacific  Introduction / extension of measures

                             CATEGORY  \
0  Governance and socio-economic measures
1  Governance and socio-economic measures
2                 Movement restrictions
3                 Public health measures
4                 Public health measures

                                   MEASURE TARGETED_POP_GROUP
\
0  Emergency administrative structures activated ...                NaN
1                         Economic measures                NaN
2                        Visa restrictions            checked
3           Isolation and quarantine policies            checked
4         Strengthening the public health system                NaN

                                  COMMENTS   \
0  Australian Health Sector Emergency Plan Activated
1  Implementation of an economic response to the ...
2  Citizens from China, Italy, South Korea, Iran,...
3  14 days self-quarantine, for nationals arrivin...
4                 Additional masks and funding

               NON_COMPLIANCE DATE_IMPLEMENTED               SOURCE  \
0               Not applicable       2020-02-17  Department of Health
1               Not applicable       2020-03-01  Department of Health
2  Refusal to Enter the Country       2020-03-01                  IATA
3               Not available       2020-03-01                  IATA
4               Not applicable       2020-03-12  Department of Health

         SOURCE_TYPE                                              LINK
\
0          Government  https://www.health.gov.au/news/health-alerts/n...
(https://www.health.gov.au/news/health-alerts/n...)
1          Government  https://www.health.gov.au/news/health-alerts/n...
(https://www.health.gov.au/news/health-alerts/n...)
2  Other organisations  https://www.iatatravelcentre.com/international...
(https://www.iatatravelcentre.com/international...)
3  Other organisations  https://www.iatatravelcentre.com/international...
(https://www.iatatravelcentre.com/international...)
4          Government  https://www.health.gov.au/news/health-alerts/n...
(https://www.health.gov.au/news/health-alerts/n...)

   ENTRY_DATE  covid_case_per_date  population
0 2020-03-14                 15.0    25499884
1 2020-03-14                 14.0    25499884
```

```
2 2020-03-14                    14.0    25499884
3 2020-03-14                    14.0    25499884
4 2020-03-14                    28.0    25499884
```

## Data Pre-processing

In [4]: ▶|
```python
#Selecting features that are important in our analysis
df = dataset[['COUNTRY','CATEGORY', 'MEASURE', 'COMMENTS', 'DATE_IMPLEMENTED'
print(df.head())
```

```
        COUNTRY                           CATEGORY   \
0  Australia  Governance and socio-economic measures
1  Australia  Governance and socio-economic measures
2  Australia                   Movement restrictions
3  Australia                   Public health measures
4  Australia                   Public health measures


                                       MEASURE   \
0  Emergency administrative structures activated ...
1                              Economic measures
2                              Visa restrictions
3              Isolation and quarantine policies
4            Strengthening the public health system


                                      COMMENTS DATE_IMPLEMENTED   \
0  Australian Health Sector Emergency Plan Activated        2020-02-17
1  Implementation of an economic response to the ...        2020-03-01
2  Citizens from China, Italy, South Korea, Iran,...        2020-03-01
3  14 days self-quarantine, for nationals arrivin...        2020-03-01
4                          Additional masks and funding        2020-03-12


   covid_case_per_date   population
0                 15.0   25499884
1                 14.0   25499884
2                 14.0   25499884
3                 14.0   25499884
4                 28.0   25499884
```

In [12]: ▶|
```python
#Converting datatype of column 'DATE_IMPLEMENTED' to datetime
df['DATE_IMPLEMENTED'] = pd.to_datetime(df['DATE_IMPLEMENTED'])
print(df['DATE_IMPLEMENTED'])
```

```
0        2020-02-17
1        2020-03-01
2        2020-03-01
3        2020-03-01
4        2020-03-12
            ...
4138     2020-11-02
4139     2020-11-02
4140     2020-11-02
4141     2020-11-02
4142     2020-11-02
Name: DATE_IMPLEMENTED, Length: 4107, dtype: datetime64[ns]
```

In [5]: ▶|
```python
print(df.isnull().sum())
```

```
COUNTRY              0
CATEGORY             0
MEASURE              0
COMMENTS             6
DATE_IMPLEMENTED    26
covid_case_per_date 31
population           0
dtype: int64
```

In [6]: ▶|
```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4143 entries, 0 to 4142
Data columns (total 7 columns):
 #   Column               Non-Null Count  Dtype
---  ------               --------------  -----
 0   COUNTRY              4143 non-null   object
 1   CATEGORY             4143 non-null   object
 2   MEASURE              4143 non-null   object
 3   COMMENTS             4137 non-null   object
 4   DATE_IMPLEMENTED     4117 non-null   datetime64[ns]
 5   covid_case_per_date  4112 non-null   float64
 6   population           4143 non-null   int64
dtypes: datetime64[ns](1), float64(1), int64(1), object(4)
memory usage: 226.7+ KB
```

In [7]: ▶|
```python
df.describe()
```

Out[7]:

|       | covid_case_per_date | population    |
|-------|---------------------|---------------|
| count | 4112.000000         | 4.143000e+03  |
| mean  | 6756.807636         | 1.223711e+08  |
| std   | 14039.385920        | 2.703601e+08  |
| min   | 0.000000            | 4.822233e+06  |
| 25%   | 42.000000           | 2.141325e+07  |
| 50%   | 556.000000          | 3.774215e+07  |
| 75%   | 4649.000000         | 6.788601e+07  |
| max   | 132854.000000       | 1.380004e+09  |

In [8]: ▶|
```python
sum = df['covid_case_per_date'].sum()
print(sum)
```

```
27783993.0
```

In [9]: ▶
```python
df.isnull().sum()
```

Out[9]:
```
COUNTRY               0
CATEGORY              0
MEASURE               0
COMMENTS              6
DATE_IMPLEMENTED     26
covid_case_per_date  31
population            0
dtype: int64
```

In [10]: ▶
```python
#Dealing with missing values
df = df[df['DATE_IMPLEMENTED'].notna()]
df.isnull().sum()
```

Out[10]:
```
COUNTRY               0
CATEGORY              0
MEASURE               0
COMMENTS              5
DATE_IMPLEMENTED      0
covid_case_per_date  10
population            0
dtype: int64
```

In [11]: ▶
```python
#Dealing with missing values
df = df[df['covid_case_per_date'].notna()]
df.isnull().sum()
```

Out[11]:
```
COUNTRY               0
CATEGORY              0
MEASURE               0
COMMENTS              5
DATE_IMPLEMENTED      0
covid_case_per_date   0
population            0
dtype: int64
```

In [13]:

```python
#Performing log normalization
df['log_value'] = np.log(df['covid_case_per_date'])
print(df)
```

```
         COUNTRY                              CATEGORY  \
0      Australia  Governance and socio-economic measures
1      Australia  Governance and socio-economic measures
2      Australia                   Movement restrictions
3      Australia                   Public health measures
4      Australia                   Public health measures
...          ...                                     ...
4138     Belgium                                Lockdown
4139     Belgium                       Social distancing
4140     Belgium                       Social distancing
4141     Belgium                  Public health measures
4142     Belgium                  Public health measures

                                          MEASURE  \
0      Emergency administrative structures activated ...
1                                Economic measures
2                                Visa restrictions
3                  Isolation and quarantine policies
4             Strengthening the public health system
...                                           ...
4138                                Full lockdown
4139        Closure of businesses and public services
4140                          Limit public gatherings
4141        Amendments to funeral and burial regulations
4142                          General recommendations

                                         COMMENTS DATE_IMPLEMENTED
\
0      Australian Health Sector Emergency Plan Activated      2020-02-17
1      Implementation of an economic response to the ...      2020-03-01
2      Citizens from China, Italy, South Korea, Iran,...      2020-03-01
3      14 days self-quarantine, for nationals arrivin...      2020-03-01
4                          Additional masks and funding      2020-03-12
...                                           ...             ...
4138                      Second nationwide lockdown      2020-11-02
4139               Closure of non-essential businesses      2020-11-02
4140    Mixing between households should be limited, ...      2020-11-02
4141                15 people permitted for funerals      2020-11-02
4142    Supermarkets are to remove all non-essential p...      2020-11-02

      covid_case_per_date   population   log_value
0                    15.0     25499884    2.708050
1                    14.0     25499884    2.639057
2                    14.0     25499884    2.639057
3                    14.0     25499884    2.639057
4                    28.0     25499884    3.332205
...                   ...          ...         ...
4138              11789.0     11589623    9.374922
4139              11789.0     11589623    9.374922
4140              11789.0     11589623    9.374922
4141              11789.0     11589623    9.374922
4142              11789.0     11589623    9.374922
```

```
[4107 rows x 8 columns]

C:\Users\nisha\anaconda3\lib\site-packages\pandas\core\series.py:679: Runti
meWarning: divide by zero encountered in log
  result = getattr(ufunc, method)(*inputs, **kwargs)
```

In [14]: ▶| `country_array = df.COUNTRY.unique()`
`country_array`

Out[14]: ```
array(['Australia', 'France', 'Germany', 'India', 'Italy',
       'United States', 'New Zealand', 'Canada', 'Norway', 'Sweden',
       'United Kingdom', 'Mexico', 'Singapore', 'Spain', 'Sri Lanka',
       'Belgium'], dtype=object)
```

# Performing EDA

**1) What was the number of cases after taking the different measures for each country?**

In [15]: ▶|
```python
# Below graphs show the number of covid cases for each country after each mea
# On Hovering over the graph, we can see the different measures taken in the
for i in range(len(country_array)):
    plt.figure(figsize=(10,5))
    new_df = df[df["COUNTRY"] == country_array[i]]
    new_df = new_df.drop_duplicates(subset=['DATE_IMPLEMENTED'])
    new_df = new_df.sort_values(by="DATE_IMPLEMENTED")
    source = ColumnDataSource(new_df)
    plot = figure(plot_width=300, plot_height=300)
    source = ColumnDataSource(data=dict(
        x=new_df['DATE_IMPLEMENTED'],
        y=new_df['covid_case_per_date'],
        desc=new_df['MEASURE'],
    ))
    p = figure(title=country_array[i],x_axis_type='datetime', x_axis_label="C
            plot_height=400, plot_width=700,
            tools=[HoverTool(tooltips='@desc')],toolbar_location=None)
    p.line('x', 'y', source=source)
    show(p)
```

**Australia**



*Cases on given date in : Australia*

**France**

Cases on given date in : France

### Germany



Cases on given date in : Germany

### India

Cases on given date in : India

### Italy



Number of Covid Cases

Cases on given date in : Italy

### United States



Number of Covid Cases

1/2020      3/2020      5/2020      7/2020      9/2020

*Cases on given date in : United States*



**New Zealand**

*Cases on given date in : New Zealand*



**Canada**

*Cases on given date in : Canada*

## Norway



Cases on given date in : Norway

## Sweden



Cases on given date in : Sweden

## United Kingdom

Cases on given date in : United Kingdom



**Mexico**

Cases on given date in : Mexico



**Singapore**

*Cases on given date in : Singapore*



**Spain**

*Cases on given date in : Spain*



**Sri Lanka**

*Cases on given date in : Sri Lanka*

**Belgium**



*Cases on given date in : Belgium*

```
<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>
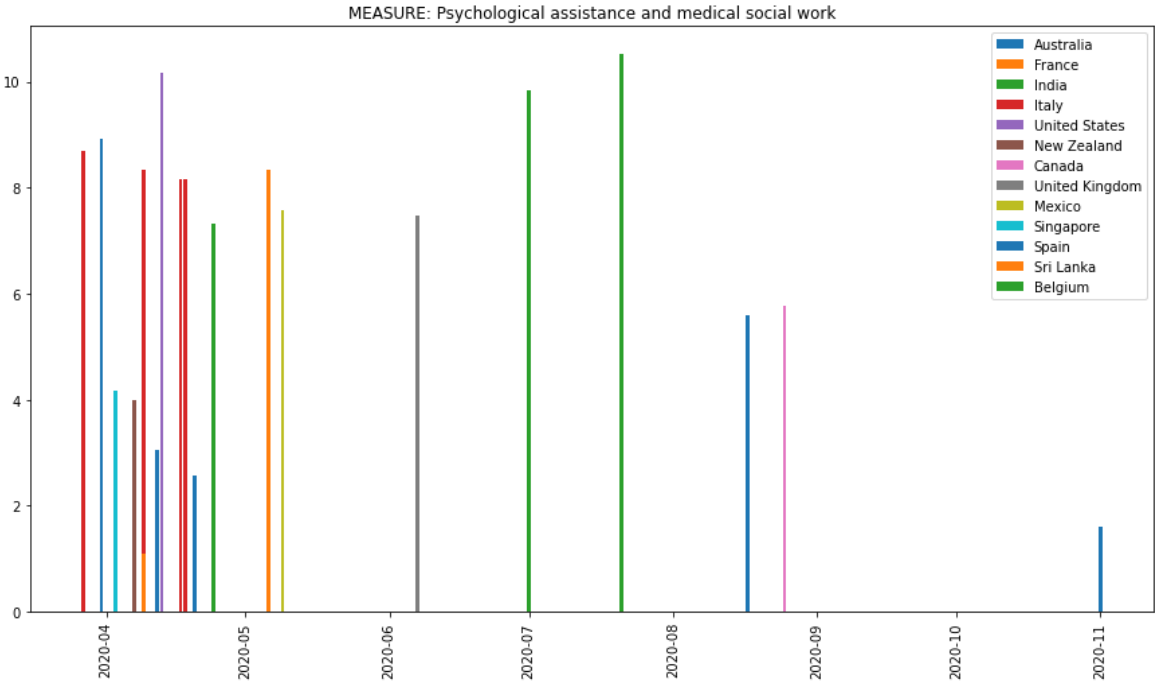
<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>
```

```
<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>

<Figure size 720x360 with 0 Axes>
```

In [16]: ▶|
```python
for i in range(len(country_array)):
    new_df = df[df["COUNTRY"] == country_array[i]]
    new_df = new_df.drop_duplicates(subset=['DATE_IMPLEMENTED'])
    new_df = new_df.sort_values(by="DATE_IMPLEMENTED")
    if i == 0:
        df0 = new_df
    else:
        df0 = df0.append(new_df)
```

In [17]: ▶|
```python
sum = df0['covid_case_per_date'].sum()
print(sum)
```

```
11102926.0
```

**2) What was the impact of each measure on the number of Covid-19 cases in the different countries?**

In [23]: ▶| 
```python
#Below are the scatterplots that show the impact of each measure on the numbe
import seaborn as sns
plt.tight_layout()
g = sns.FacetGrid(df0, col="MEASURE", hue="COUNTRY", height=3, aspect=3, col_
g.map_dataframe(sns.scatterplot, x="DATE_IMPLEMENTED", y="covid_case_per_date
g.set_axis_labels("DATE_IMPLEMENTED", "Covid_case_per_date")
g.add_legend()
```

Out[23]: <seaborn.axisgrid.FacetGrid at 0x7f03a72c3748>

<Figure size 432x288 with 0 Axes>

In [24]:

```python
# These bar-graphs show the impact of each measure on the number of Covid-19
# We have plotted the log-value of covid cases on the y-axis to normalize the
Measure_array=df0["MEASURE"].unique()
for i in range(len(Measure_array)):
    new_df = df0[df0['MEASURE'] == Measure_array[i]]
    Country_array=new_df["COUNTRY"].unique()
    if len(Country_array) > 1:
        plt.figure(figsize=(15,8))
        for j in range(len(Country_array)):
            df1 = df0[df0["MEASURE"]== Measure_array[i]]
            df1 = df1[df1["COUNTRY"] == Country_array[j]]
            if len(df1) > 0:
                plt.bar(df1['DATE_IMPLEMENTED'],df1['log_value'],label=Countr
                plt.xticks(rotation=90)
        plt.legend()
        title="MEASURE: "+Measure_array[i]
        plt.title(title)
        plt.show()
```



MEASURE: Emergency administrative structures activated or established

MEASURE: Economic measures



MEASURE: Strengthening the public health system

### MEASURE: Limit public gatherings



### MEASURE: Awareness campaigns



### MEASURE: Isolation and quarantine policies

MEASURE: General recommendations



MEASURE: Testing policy

### MEASURE: Surveillance and monitoring

Legend:
- Australia
- France
- Germany
- India
- Italy
- United States
- New Zealand
- Canada
- Norway
- United Kingdom
- Singapore
- Spain
- Sri Lanka

### MEASURE: Domestic travel restrictions

Legend:
- Australia
- France
- India
- Italy
- New Zealand
- Norway
- Mexico
- Spain
- Sri Lanka
- Belgium

### MEASURE: Visa restrictions

Legend:
- Australia
- France
- Germany
- India
- Italy
- United States
- New Zealand
- Canada
- Norway
- Sweden
- United Kingdom
- Singapore
- Sri Lanka

MEASURE: Border closure



MEASURE: Psychological assistance and medical social work

MEASURE: Other public health measures enforced



MEASURE: Closure of businesses and public services

MEASURE: Schools closure



MEASURE: Additional health/documents requirements upon arrival

MEASURE: International flights suspension



MEASURE: State of emergency declared

MEASURE: Partial lockdown



MEASURE: Amendments to funeral and burial regulations



MEASURE: Border checks

MEASURE: Requirement to wear protective gear in public



MEASURE: Military deployment



MEASURE: Limit product imports/exports

MEASURE: Changes in prison-related policies



MEASURE: Health screenings in airports and border crossings

MEASURE: Full lockdown

Determining the Countries with more than 15000 new Covid-19 cases in a day and the measures taken by them

```
In [18]:  ▶  df0['max_covid_cases'] = np.where(df0['covid_case_per_date']>= 15000, True, F
              max_covid_cases_countries = df0[df0['max_covid_cases'] == True]
              max_covid_cases_countries['MEASURE'].unique()
```

```
Out[18]:  array(['Isolation and quarantine policies', 'Testing policy',
                 'General recommendations',
                 'Strengthening the public health system',
                 'Other public health measures enforced',
                 'Psychological assistance and medical social work',
                 'Awareness campaigns', 'Economic measures',
                 'Limit product imports/exports',
                 'Emergency administrative structures activated or established',
                 'Additional health/documents requirements upon arrival',
                 'Mass population testing', 'Surveillance and monitoring',
                 'Visa restrictions', 'Closure of businesses and public services',
                 'Border closure', 'Military deployment', 'Limit public gatherings',
                 'Schools closure', 'Requirement to wear protective gear in public'],
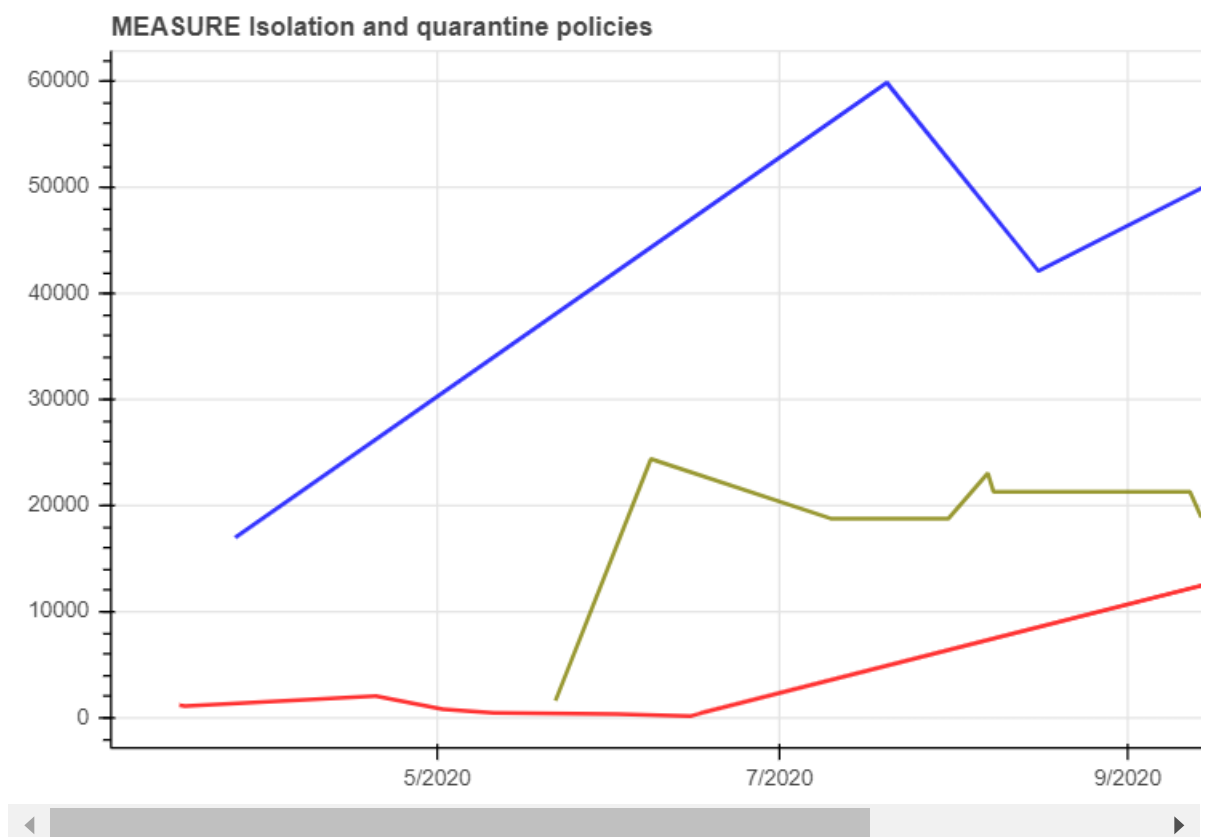                dtype=object)
```

In [19]:
```python
Measure_array = max_covid_cases_countries["MEASURE"].unique()
Measure_array
```

Out[19]:
```
array(['Isolation and quarantine policies', 'Testing policy',
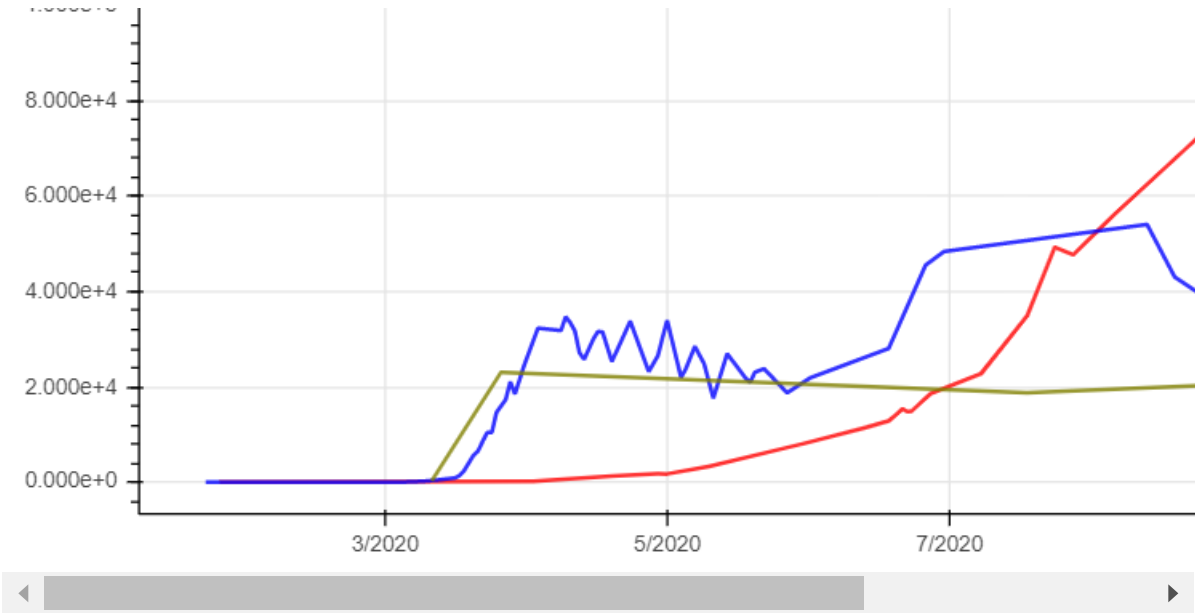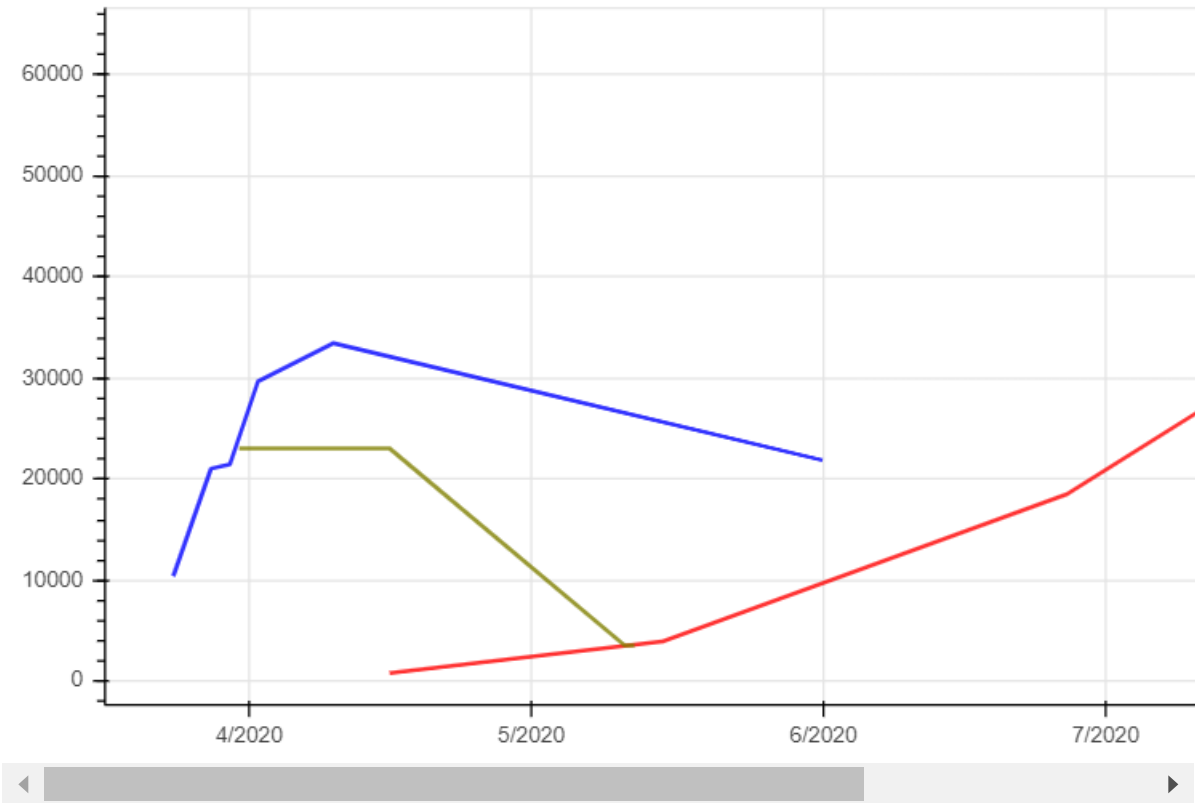       'General recommendations',
       'Strengthening the public health system',
       'Other public health measures enforced',
       'Psychological assistance and medical social work',
       'Awareness campaigns', 'Economic measures',
       'Limit product imports/exports',
       'Emergency administrative structures activated or established',
       'Additional health/documents requirements upon arrival',
       'Mass population testing', 'Surveillance and monitoring',
       'Visa restrictions', 'Closure of businesses and public services',
       'Border closure', 'Military deployment', 'Limit public gatherings',
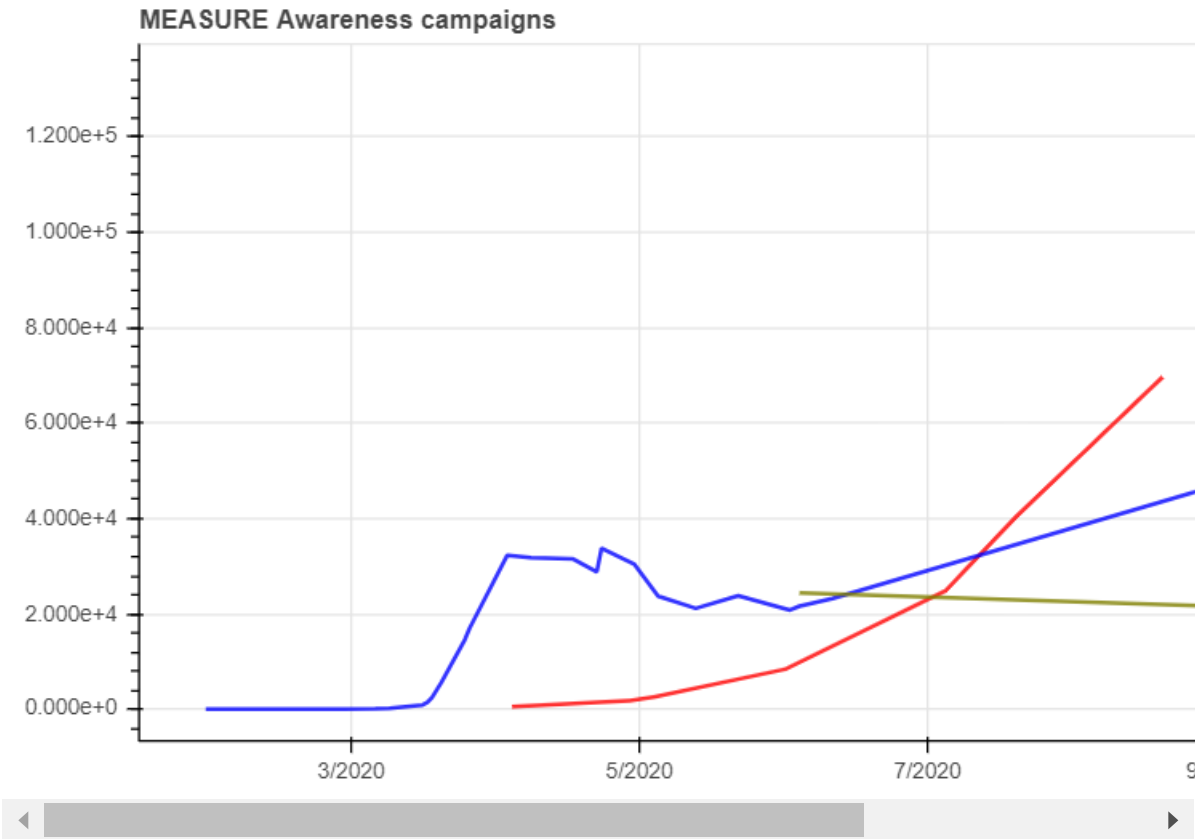       'Schools closure', 'Requirement to wear protective gear in public'],
      dtype=object)
```

**3) Countries with maximum Covid-19 cases, the measures taken by them and the impact of those measures**

In [20]: ▶|

```python
# Below graphs show the countries with maximum Covid-19 cases, the measures t
# On hovering over the graph, we can see the number of new covid cases on a p
Measure_array = max_covid_cases_countries["MEASURE"].unique()
mapper_color=['red','blue','olive','green','black','purple','maroon']


for i in range(len(Measure_array)):
    p = figure(plot_width=800, plot_height=400, x_axis_type="datetime",toolba
    p.title.text = "MEASURE "+Measure_array[i]
    new_df = max_covid_cases_countries[max_covid_cases_countries['MEASURE'] =
    Country_array = new_df["COUNTRY"].unique()
    if len(Country_array) > 1:
        for j in range(len(Country_array)):
            df1 = df[df["MEASURE"]== Measure_array[i]]
            df1 = df1[df1["COUNTRY"] == Country_array[j]]
            df1 = df1.drop_duplicates(subset=['DATE_IMPLEMENTED'])
            df1 = df1.sort_values(by="DATE_IMPLEMENTED")
            source = ColumnDataSource(data=dict(
                x=df1['DATE_IMPLEMENTED'],
                y=df1['covid_case_per_date'],
                desc=df1['covid_case_per_date']
                ))
            p.line('x','y', line_width=2, alpha=0.8,
                    legend_label=Country_array[j],source=source,color=mapper_c
            p.add_tools(HoverTool(tooltips='@desc'))
        show(p)
        p.legend.location = "top_left"
        p.legend.click_policy="hide"
```



MEASURE Isolation and quarantine policies



MEASURE Testing policy

MEASURE General recommendations



MEASURE Strengthening the public health system

8.000e+4

6.000e+4

4.000e+4

2.000e+4

0.000e+0

3/2020          5/2020          7/2020

### MEASURE Other public health measures enforced



60000

50000

40000

30000

20000

10000

0

4/2020          5/2020          6/2020          7/2020

### MEASURE Psychological assistance and medical social work



40000

30000

## MEASURE Awareness campaigns



## MEASURE Economic measures

## MEASURE Emergency administrative structures activated or established



## MEASURE Surveillance and monitoring

**MEASURE Visa restrictions**



**MEASURE Closure of businesses and public services**



**MEASURE Limit public gatherings**

**4) Which countries took the maximum amount of measures?**

In [28]:
```python
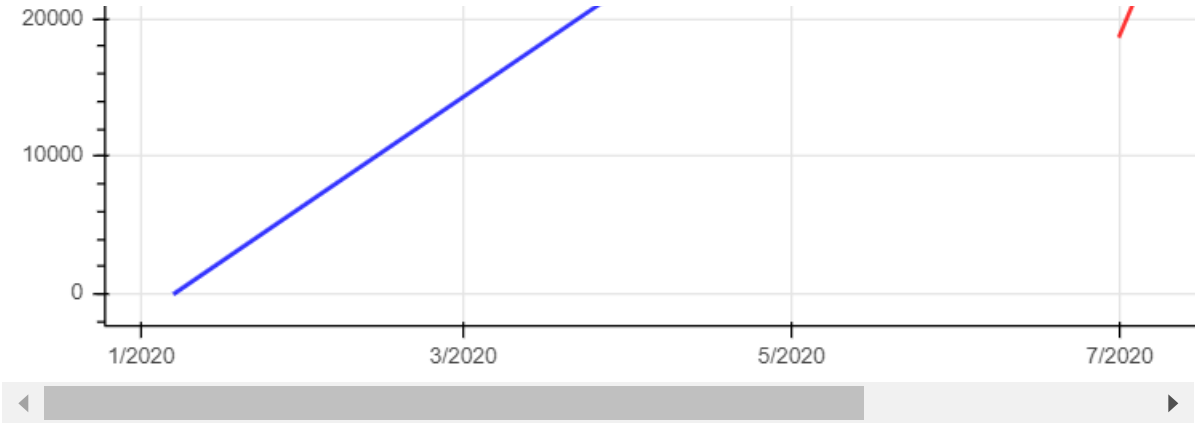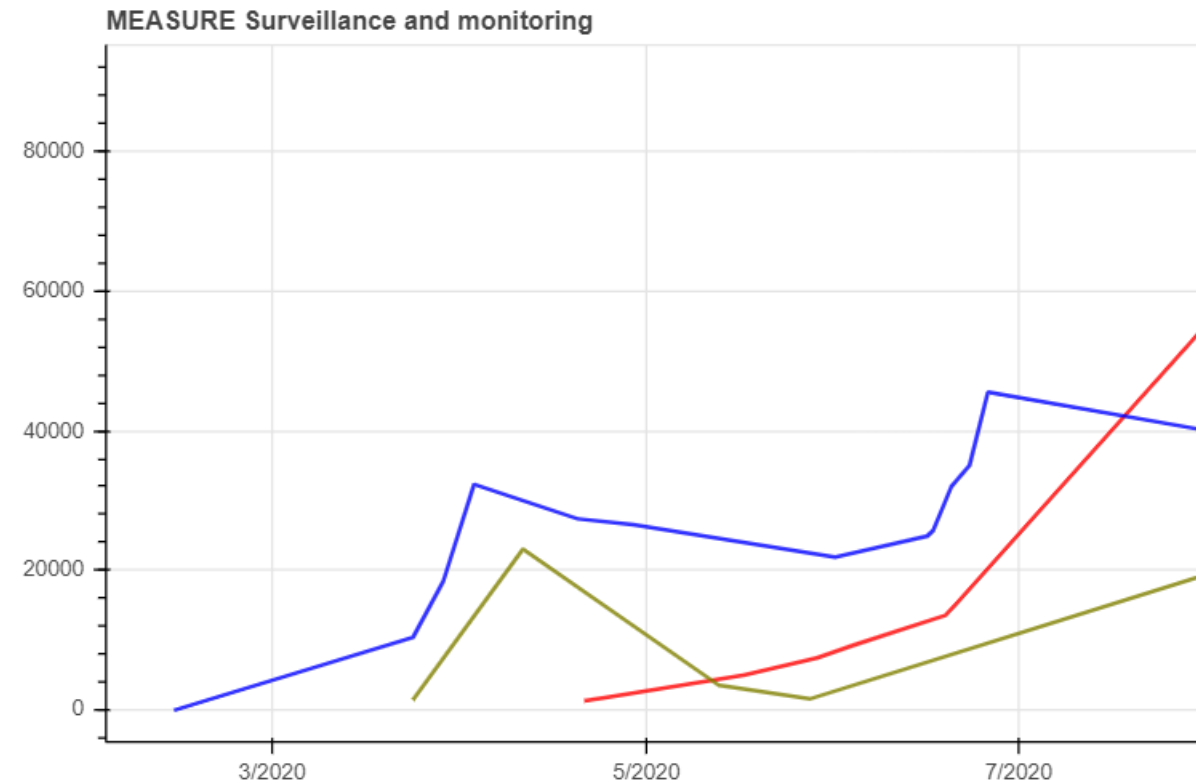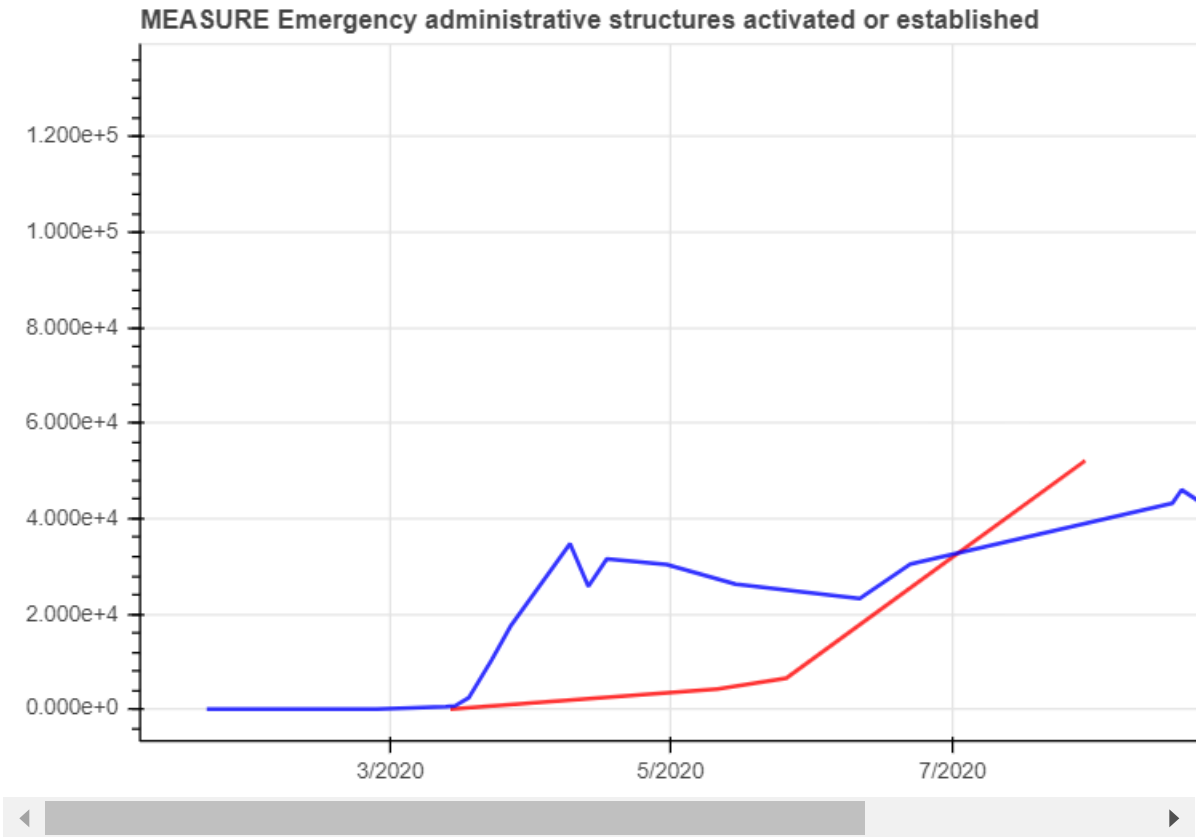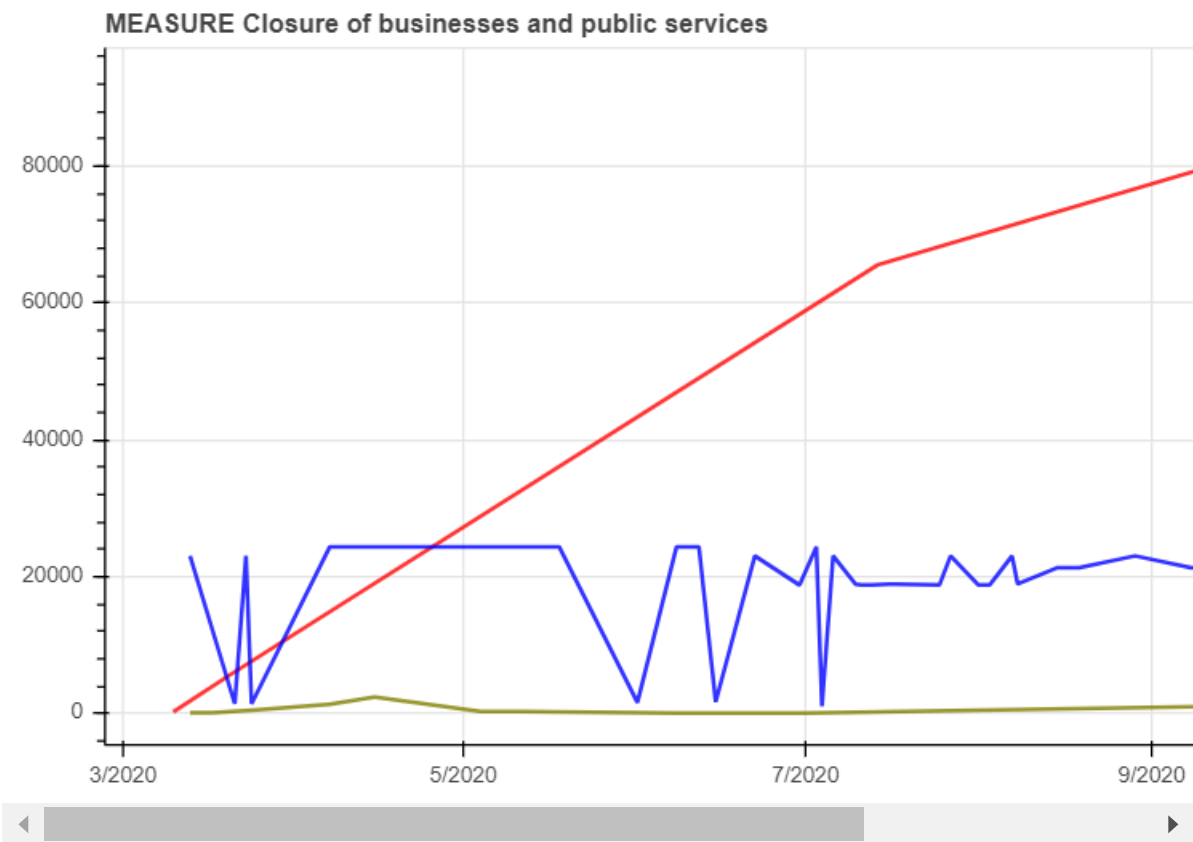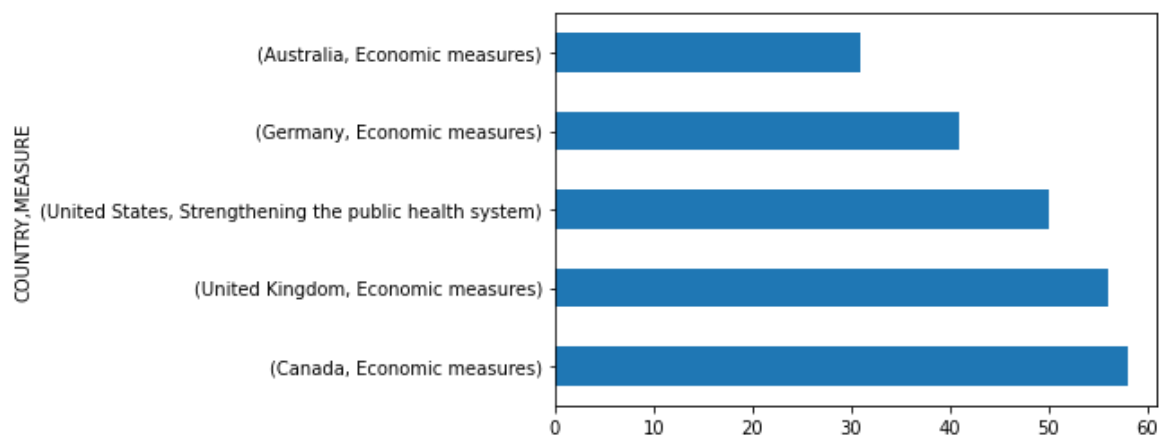#Countries that took maximum amount of measures
sorting=df0.groupby(['COUNTRY','MEASURE']).size()
sorting.sort_values(inplace=True, ascending=False)
countries_max_measures = sorting.head()
countries_max_measures
countries_max_measures.plot.barh(stacked=True)
```

Out[28]: <matplotlib.axes._subplots.AxesSubplot at 0x7f039ff8d9e8>



**5) Which countries opted for a complete lockdown?**

In [29]: ▶

```python
#Countries that opted for a complete lockdown
comp_lockdown=df.loc[df['MEASURE'] == 'Full lockdown']
full_country_array = comp_lockdown.COUNTRY.unique()
plt.figure(figsize=(10,8))
for i in range(len(full_country_array)):
    country_df = comp_lockdown[comp_lockdown["COUNTRY"] == full_country_array
    plt.plot(country_df['DATE_IMPLEMENTED'],country_df['covid_case_per_date']
    plt.legend()
    title="Full Lockdown"
    plt.title(title)
plt.show()
```



Full Lockdown

**6) What is the distribution of covid cases across the different measures taken by different countries?**

In [31]:  ▶|   ```python
#This is the boxplot of measures vs the log-value of number of Covid-cases
plt.figure(figsize=(30,10))
plt.xticks(rotation=90)
sns.set_theme(style="whitegrid")
sns.boxplot(df0['MEASURE'],df0['log_value'])
```

/usr/local/lib/python3.6/dist-packages/seaborn/_decorators.py:43: FutureWar
ning: Pass the following variables as keyword args: x, y. From version 0.1
2, the only valid positional argument will be `data`, and passing other arg
uments without an explicit keyword will result in an error or misinterpreta
tion.
  FutureWarning

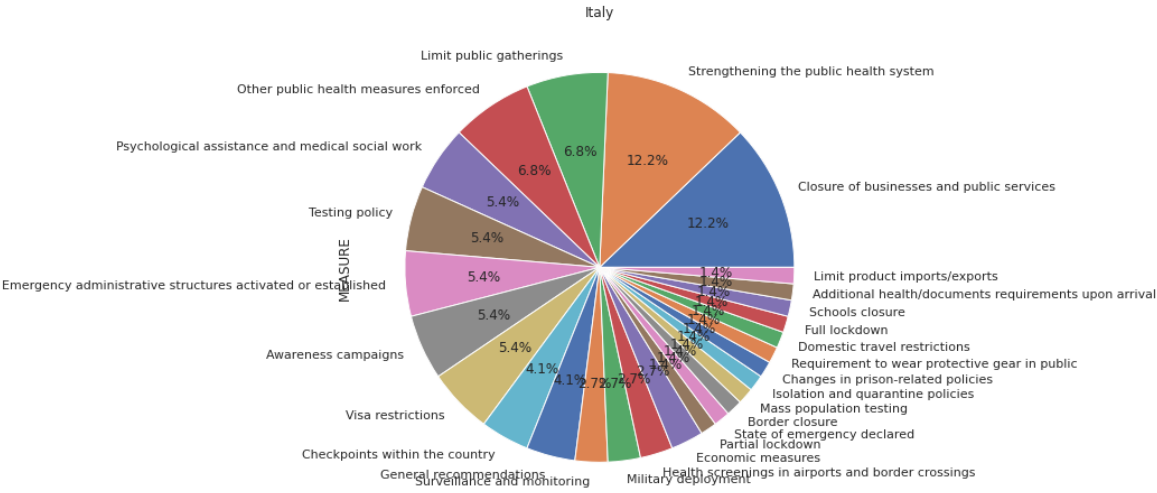Out[31]:  <matplotlib.axes._subplots.AxesSubplot at 0x7f03a0309550>



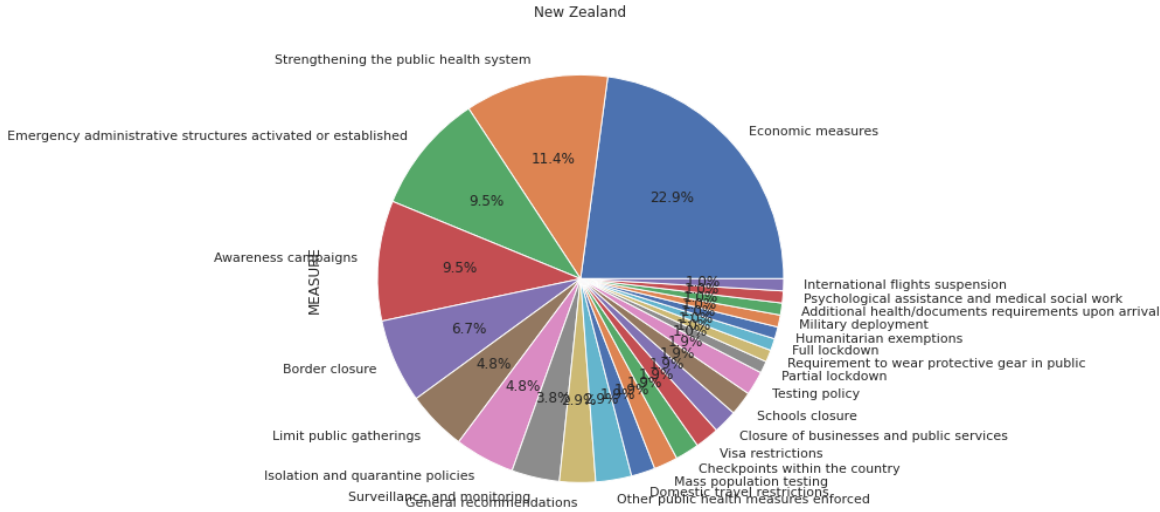**7) What is the distribution of measures taken by different countries?**
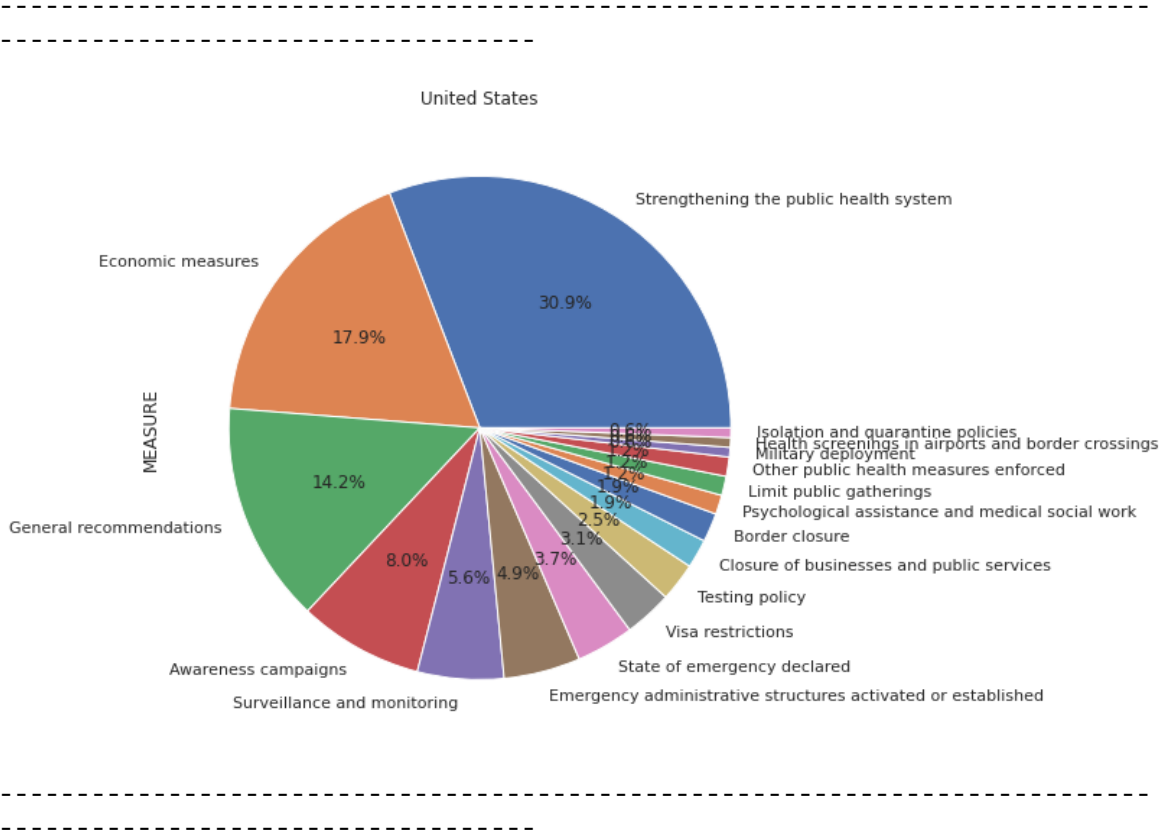
In [32]:

```python
# These pie-charts show amount of different measures taken by different count
country_array = df0['COUNTRY'].unique()
for i in range(len(country_array)):
    countries = df0[df0['COUNTRY'] == country_array[i]]
    plt.figure(figsize=(8,8))
    countries['MEASURE'].value_counts().plot.pie(autopct="%1.1f%%")
    plt.title(country_array[i])
    plt.show()
    print("-------------------------------------------------------------
```
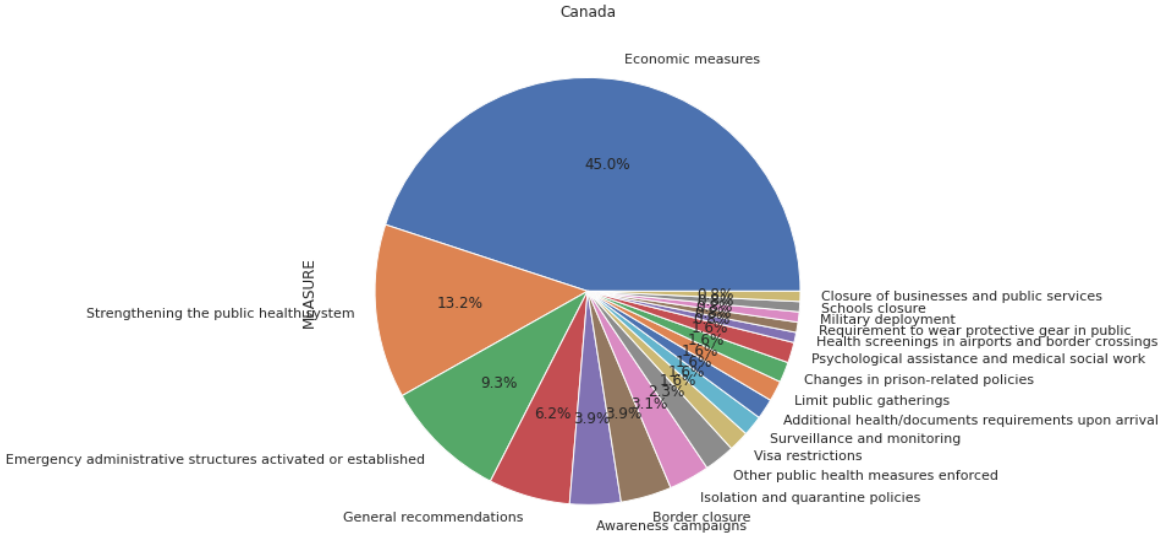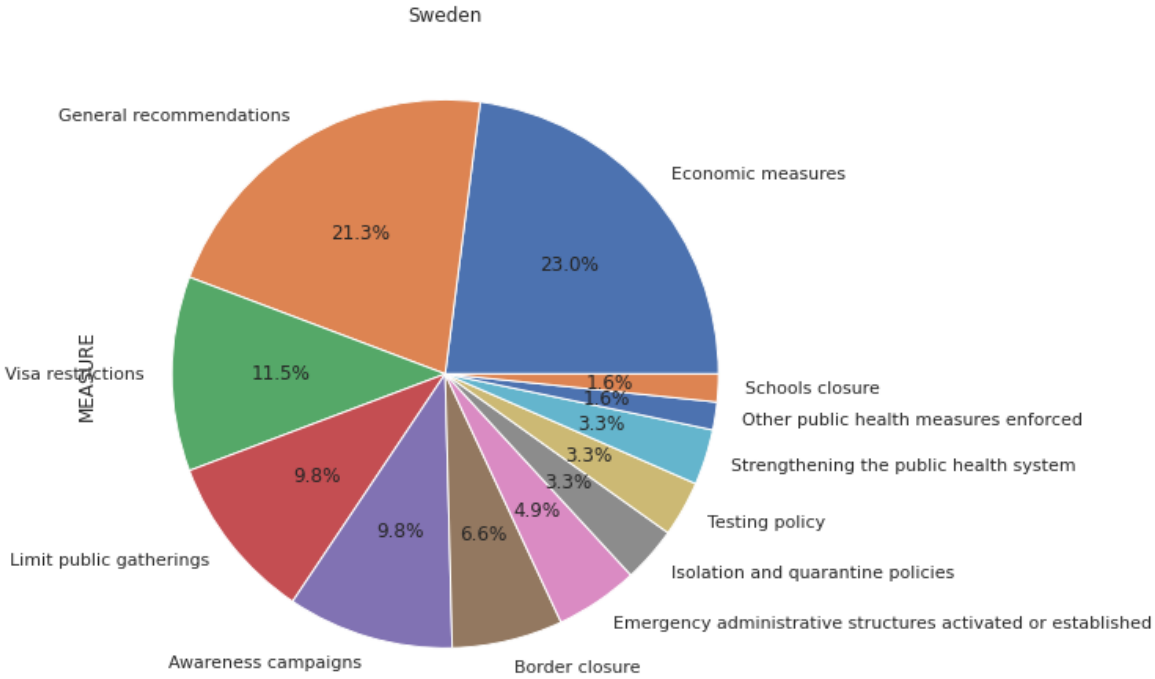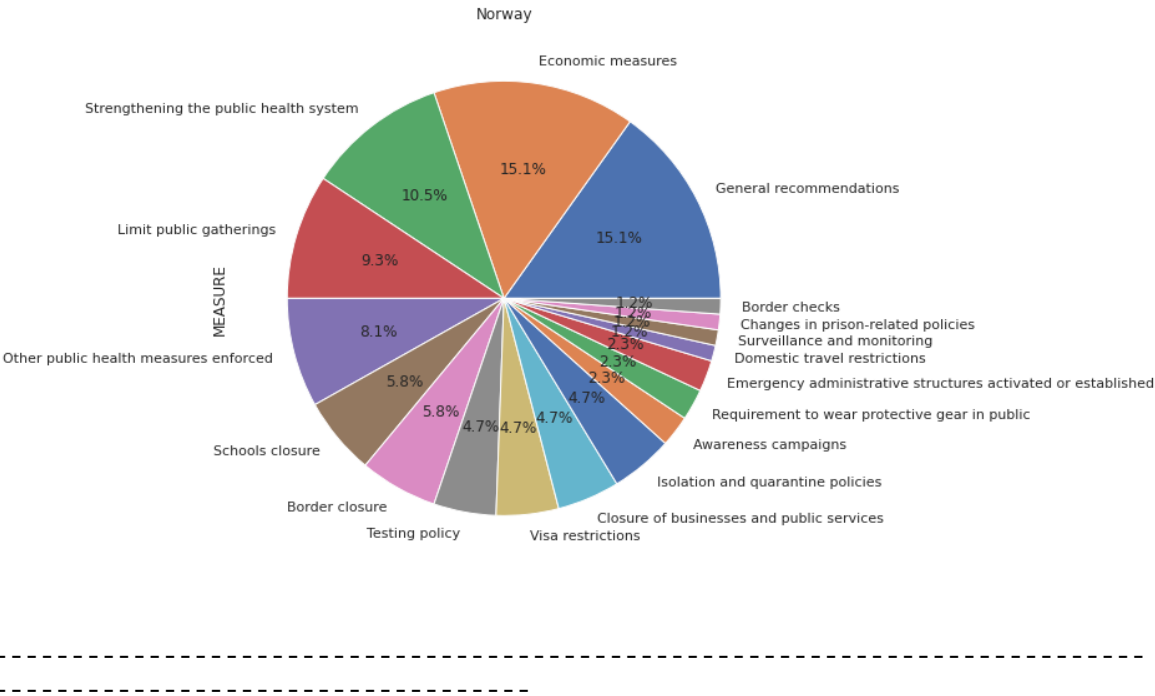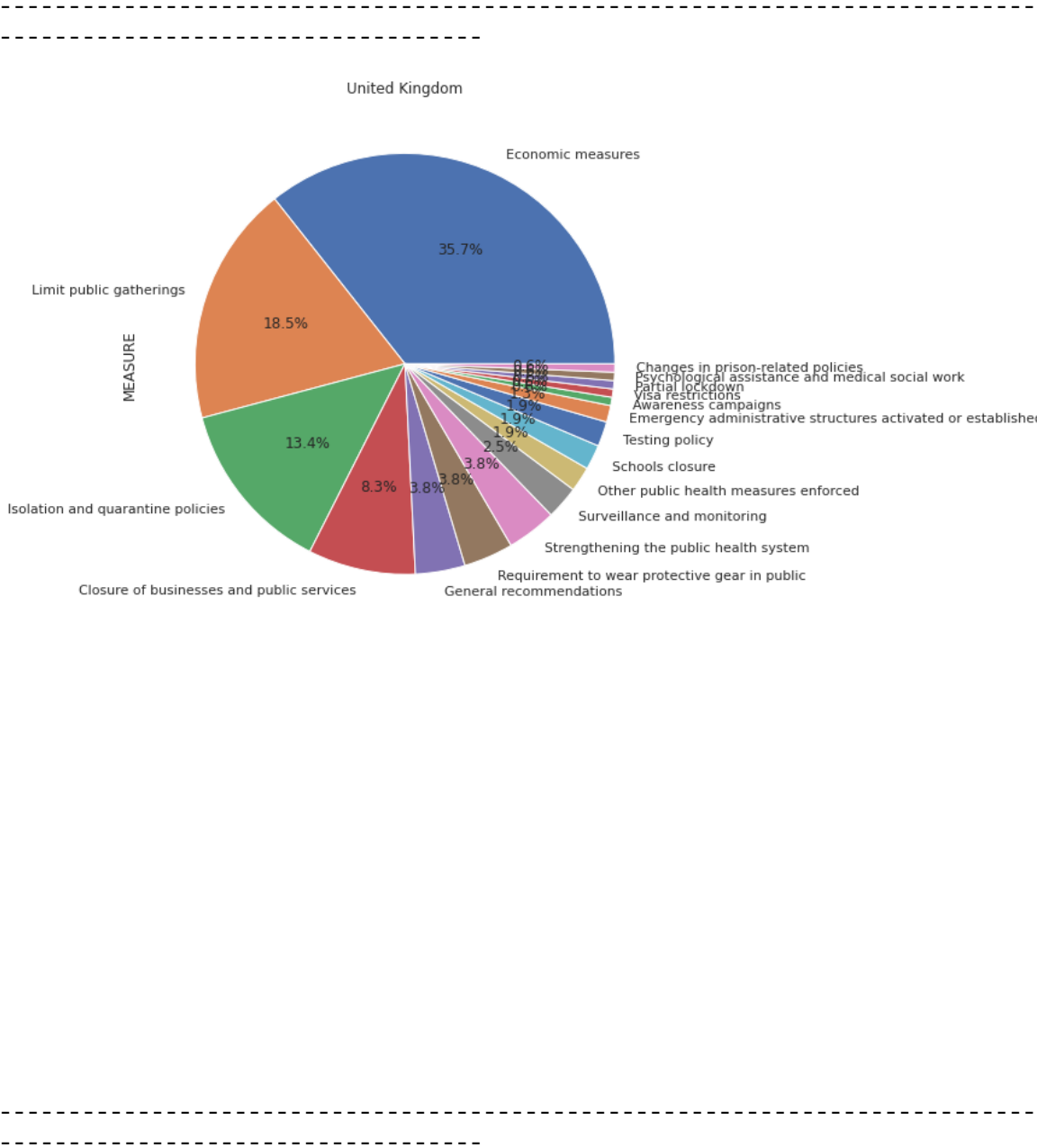


```
-------------------------------------------------------------------------
----------------------------------
```

## France



------------------------------------------------------------------------

## Germany



------------------------------------------------------------------------

India



Italy

------------------------------------------------------------------------
------------------------------------

### United States



------------------------------------------------------------------------
------------------------------------

### New Zealand

Canada

## Norway



## Sweden

------------------------------------------------------------------------
-----------------------------------

United Kingdom

## Mexico



## Singapore
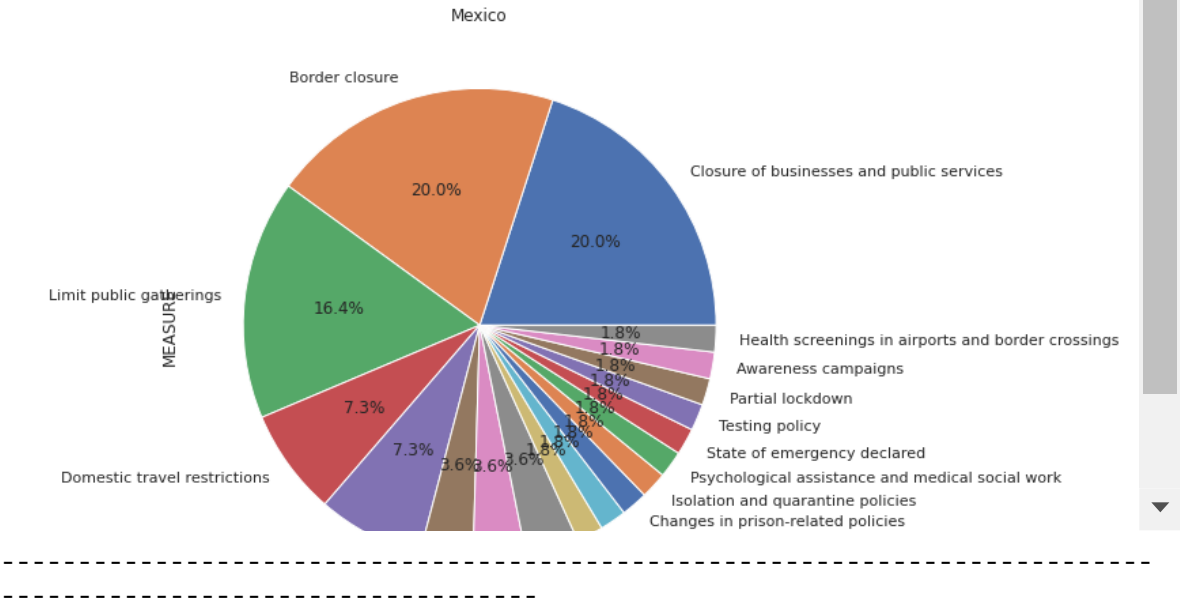
## Spain



## Sri Lanka

Belgium

```
In [ ]:   print(df0.groupby('COUNTRY').size())
```
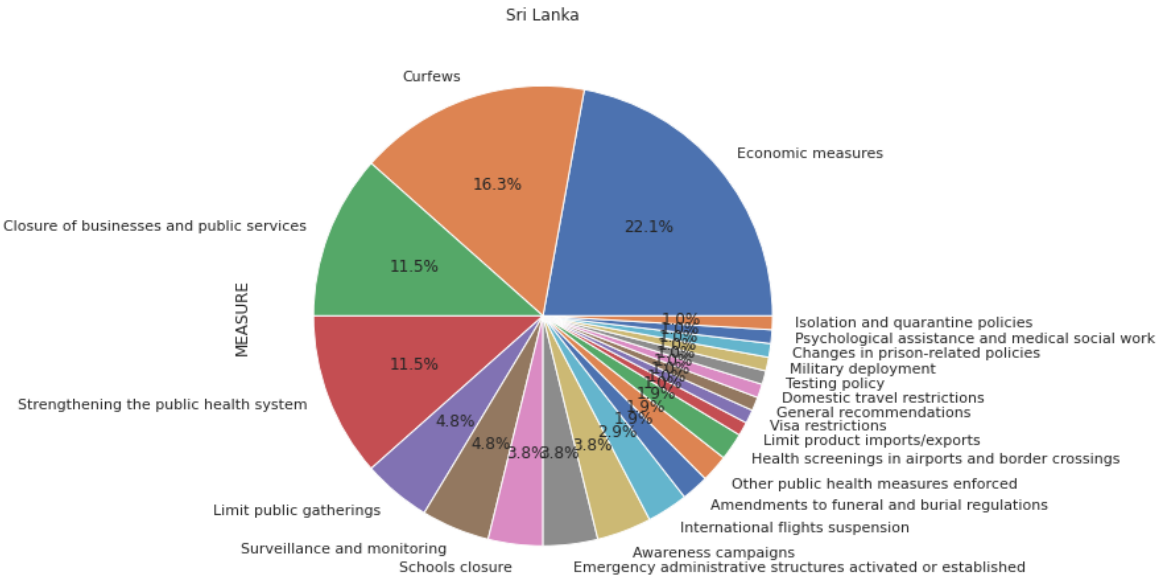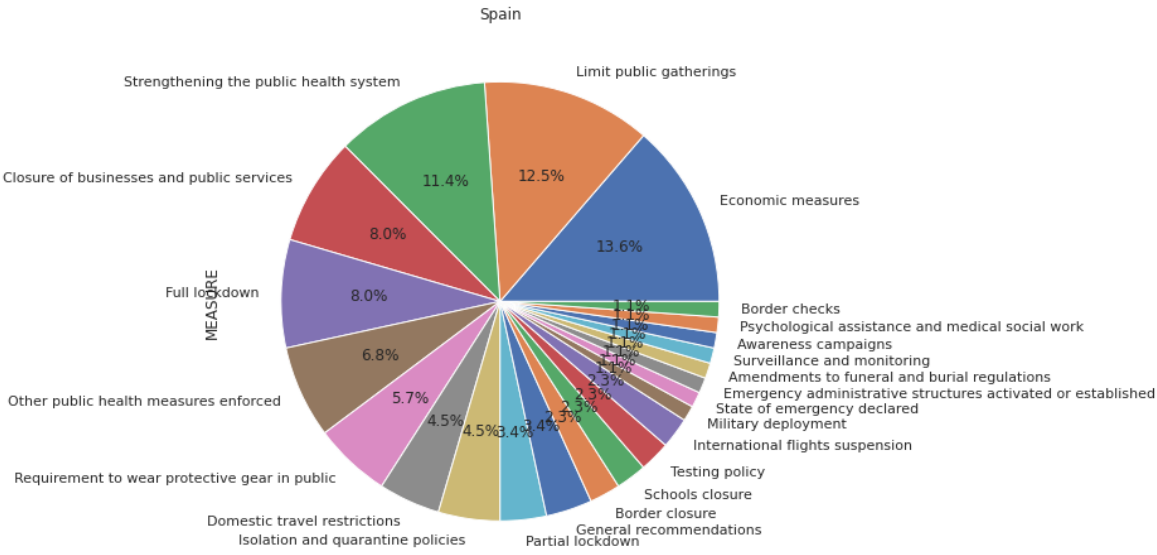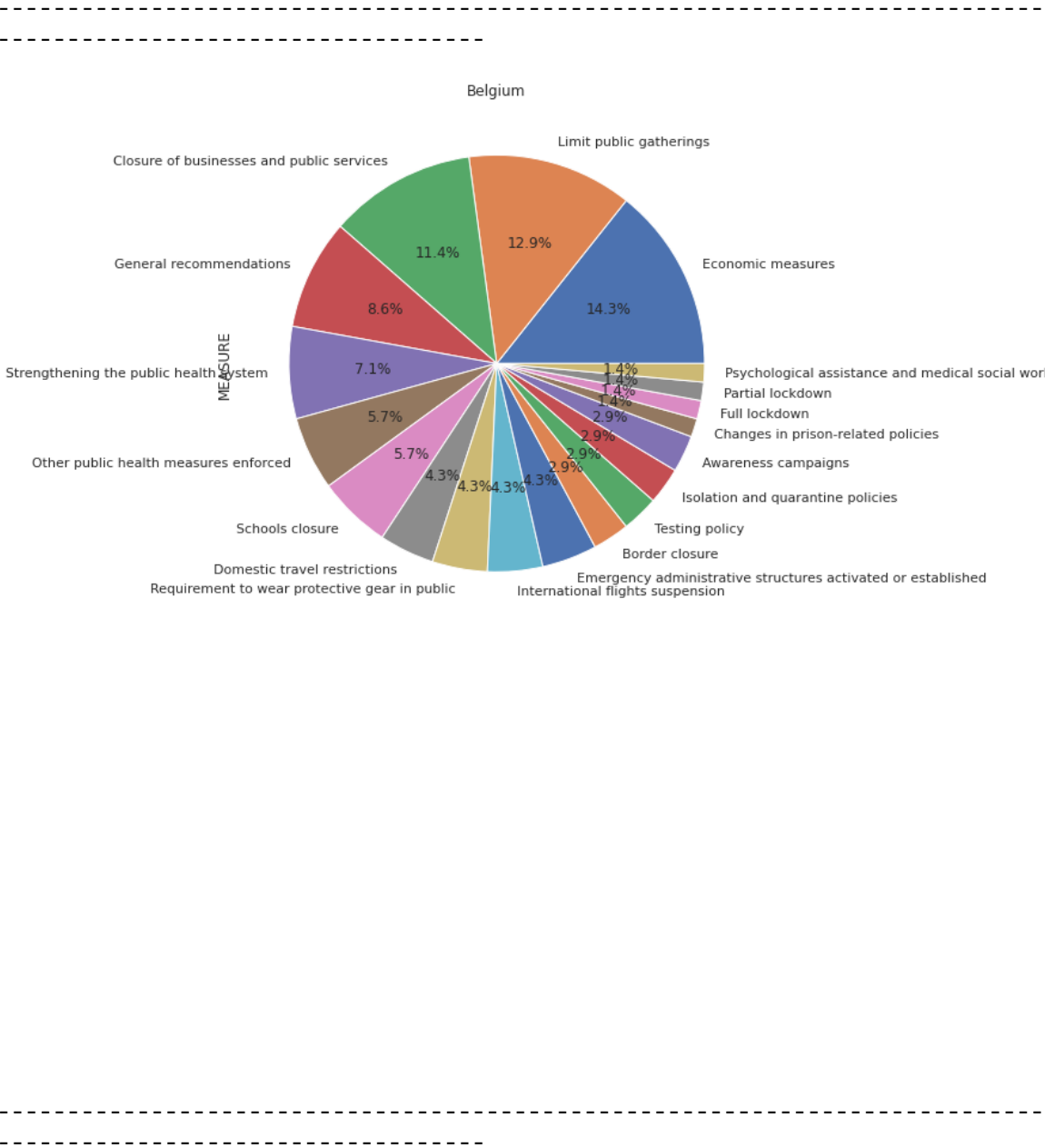
```
COUNTRY
Australia         136
Belgium            70
Canada            129
France             76
Germany           124
India              91
Italy              74
Mexico             55
New Zealand       105
Norway             86
Singapore         108
Spain              88
Sri Lanka         104
Sweden             61
United Kingdom    157
United States     162
dtype: int64
```

**8) What is the count of different measures taken by each countries?**

In [33]:  ▶|  `#This count plot shows the count of different measures taken by each country`
`plt.figure(figsize=(20,15))`
`sns.countplot(x='COUNTRY', hue='MEASURE', data=df0)`

Out[33]:  `<matplotlib.axes._subplots.AxesSubplot at 0x7f03a26a1860>`