

CM3400 Report

Introduction

In 2019 a new form of Coronavirus was discovered, called SARS-CoV-2. It causes mild symptoms in most people but in the elderly and those with underlying health conditions it can cause much more serious illnesses. It is highly contagious and can be transmitted from person to person. In 2020 the outbreak developed into a world-wide pandemic, more than 3 million people as of April 2021 have died after contracting the virus. I will be taking a deeper look at some of the statistics of the corona virus outbreak, as well as seeing if it is possible to model the death rate and mortality.

Dataset

The data set consists of profiles of many countries. It collates information on cases, deaths and government responses for each country in relation to the Covid-19 outbreak. The data covers a year, broken down into 3 month periods, from February 2020 to February 2021.

Load packages and data set

Firstly loading packages and data set into python and checking size of data set.

```
In [54]: import pandas as pd
import seaborn
import matplotlib.pyplot as plt
import numpy as np
from datetime import datetime
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
import sklearn.metrics as metrics
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import RandomForestRegressor

df = pd.read_csv("C:/Users/rossm/CM3400/covid.csv")
print(df.shape)
df.head(5)
```

```
Out[55]:
```

	location	continent	date	cases_per_M	deaths_per_M	stringency_index	population_density	median_age	aged_65_older	aged_70_older
0	Afghanistan	Asia	May-20	390.641	6.602	84.26	54.422	18.6	2.581	1.337
1	Albania	Europe	May-20	394.398	11.120	86.11	104.871	38.0	13.188	8.643
2	Algeria	Africa	May-20	214.202	14.868	76.85	17.348	29.1	6.211	3.857
3	Argentina	South America	May-20	372.823	11.904	90.74	16.177	31.9	11.198	7.441
4	Australia	Oceania	May-20	281.374	4.000	63.43	3.202	37.9	15.504	10.129

5 rows x 23 columns

23 variables on 525 different country/date combinations.

Statistics of individual variables

Now to look at the number of confirmed cases of coronavirus per million people in each country.

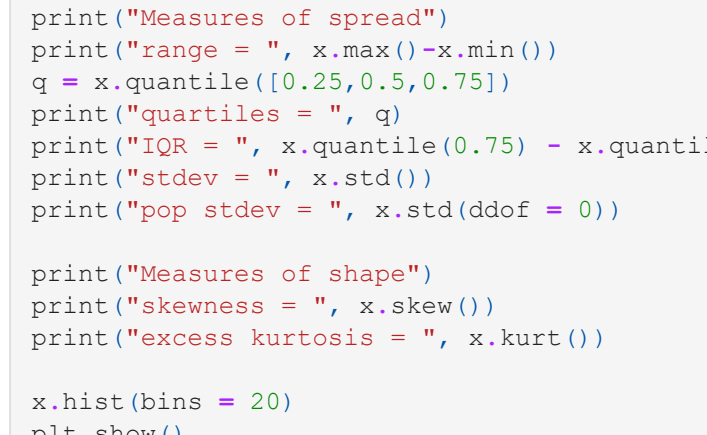
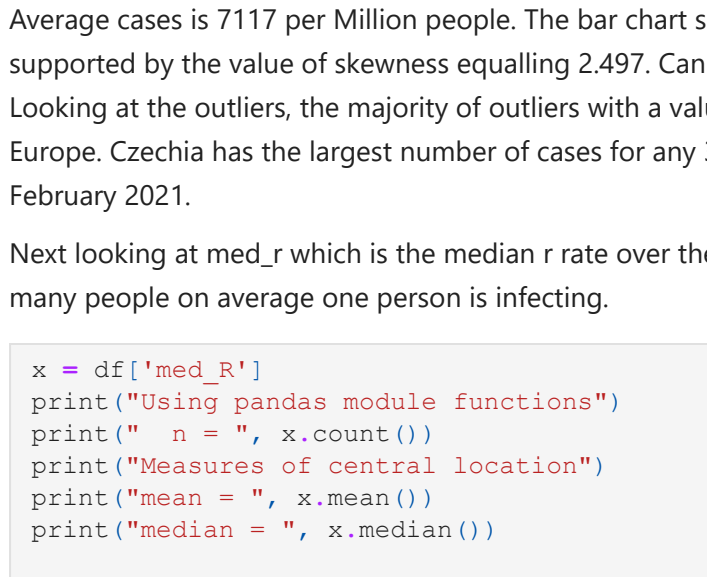
```
In [56]: x = df['cases_per_M']
print("Using pandas module functions")
print("n = ", x.count())
print("Measures of central location")
print("mean = ", x.mean())
print("median = ", x.median())

print("Measures of spread")
print("range = ", x.max()-x.min())
q = x.quantile([0.25,0.5,0.75])
print("quartiles = ", q)
print("IQR = ", x.quantile(0.75) - x.quantile(0.25))
print("stddev = ", x.std())
print("pop stddev = ", x.std(ddof = 0))

print("Measures of shape")
print("skewness = ", x.skew())
print("excess kurtosis = ", x.kurt())

x.hist(bins = 20)
plt.show()
x.plot.box()
plt.show()
```

Using pandas module functions
n = 525
Measures of central location
mean = 7117.880666666669
median = 1943.832
range = 83074.35399999999
quartiles = 0.25 322.346
0.50 1943.832
0.75 8856.786
Name: cases_per_M, dtype: float64
IQR = 8534.44
stddev = 11393.095143653863
pop stddev = 11382.239404987053
Measures of shape
skewness = 2.496126965529684
excess kurtosis = 7.327139760950313



```
In [57]: dfcases = df[df['cases_per_M'] > 40000]
dfcases.sort_values(by = ['cases_per_M'], ascending=False)
```

```
Out[57]:
```

	location	continent	date	cases_per_M	deaths_per_M	stringency_index	population_density	median_age	aged_65_older	aged_70_older
421	Czechia	Europe	Feb-21	83076.430	1434.497	69.44	137.176	43.3	19.027	11.58
495	Slovenia	Europe	Feb-21	60455.561	1205.427	69.44	102.619	44.5	19.062	12.9
429	Estonia	Europe	Feb-21	56719.026	468.136	41.67	31.033	42.7	19.452	13.4
452	Israel	Asia	Feb-21	56092.623	366.932	57.41	402.606	30.6	11.733	7.35
465	Lithuania	Europe	Feb-21	53256.246	1027.076	66.67	45.135	43.5	19.002	13.77
491	Portugal	Europe	Feb-21	50686.560	1196.661	87.96	112.371	46.2	21.502	14.92
518	United States	North America	Feb-21	47971.498	803.773	68.06	35.608	38.3	15.413	9.73
506	Sweden	Europe	Feb-21	47742.065	642.720	69.44	24.718	41.0	19.985	13.43
288	Czechia	Europe	Nov-20	46566.518	734.990	69.44	137.176	43.3	19.027	11.58
332	Luxembourg	Europe	Nov-20	44814.817	314.708	60.19	231.447	39.7	14.312	9.84
462	Lebanon	Asia	Feb-21	43292.288	652.851	92.59	594.561	31.1	8.514	5.43
498	Slovakia	Europe	Feb-21	42788.329	1434.160	71.30	113.128	41.2	15.070	9.16
270	Belgium	Europe	Nov-20	42461.099	582.417	63.89	375.564	41.8	18.571	12.84
486	Paraguay	South America	Feb-21	42360.199	679.063	74.07	55.133	29.7	7.918	5.03
461	Latvia	Europe	Feb-21	41102.172	832.891	56.48	31.212	43.9	19.754	14.13

15 rows x 23 columns

Average cases is 7117 per Million people. The bar chart shows that cases per million of people is strongly positively skewed, this is supported by the value of skewness equalling 2.497. Can see that the majority of cases is between 0 and 4000 with numerous large outliers. Looking at the outliers, the majority of outliers with a value of cases per million people greater than 40,000 come from countries inside Europe. Czechia has the largest number of cases for any 3 month period with around 83,000 cases per million from December 2020 to February 2021.

Next looking at med_r which is the median r rate over the 3 month period in each country. The r rate (reproduction number) describes how many people on average one person is infecting.

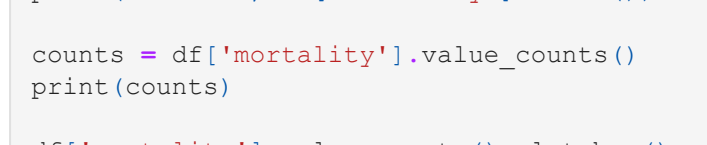
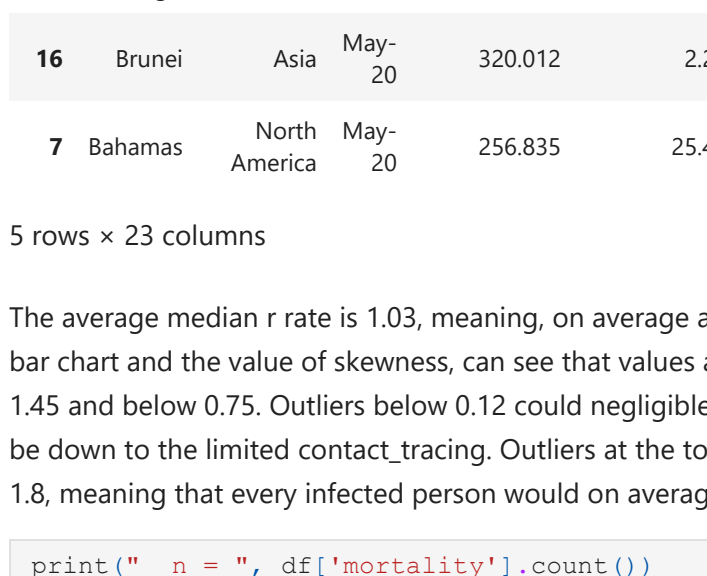
```
In [58]: x = df['med_R']
print("Using pandas module functions")
print("n = ", x.count())
print("Measures of central location")
print("mean = ", x.mean())
print("median = ", x.median())

print("Measures of spread")
print("range = ", x.max()-x.min())
q = x.quantile([0.25,0.5,0.75])
print("quartiles = ", q)
print("IQR = ", x.quantile(0.75) - x.quantile(0.25))
print("stddev = ", x.std())
print("pop stddev = ", x.std(ddof = 0))

print("Measures of shape")
print("skewness = ", x.skew())
print("excess kurtosis = ", x.kurt())

x.hist(bins = 20)
plt.show()
x.plot.box()
plt.show()
```

Using pandas module functions
n = 525
Measures of central location
mean = 1.0392190476190468
median = 1.05
range = 1.8
quartiles = 0.25 0.95
0.50 1.05
0.75 1.15
Name: med_R, dtype: float64
IQR = 0.19999999999999996
stddev = 0.22208527022986462
pop stddev = 0.2206422651284147
Measures of shape
skewness = -1.165629123222257
excess kurtosis = 0.021591273844064



```
In [59]: df.sort_values(by = ['med_R'], ascending=False).head(5)
```

```
Out[59]:
```

	location	continent	date	cases_per_M	deaths_per_M	stringency_index	population_density	median_age	aged_65_older	aged_70_older
84	Nepal	Asia	May-20	53.918	0.241	92.59	204.430	25.0	5.809	3.212
74	Malawi	Africa	May-20	14.689	0.157	60.19	197.519	18.1	2.979	1.783
112	Tajikistan	Asia	May-20	410.479	4.718	44.44	64.281	23.3	3.466	2.155
94	Peru	South America	May-20	4988.347	136.571	89.81	25.129	29.1	7.151	4.455
38	Ethiopia	Africa	May-20	10.186	0.079	80.56	104.957	19.8	3.526	2.063

5 rows x 23 columns

```
In [60]: df.sort_values(by = ['med_R'], ascending=True).head(5)
```

```
Out[60]:
```

	location	continent	date	cases_per_M	deaths_per_M	stringency_index	population_density	median_age	aged_65_older	aged_70_older
508	Tajikistan	Asia	Feb-21	116.800	0.419	16.67	64.281	23.3	3.466	2.155
147	Brunei	Asia	Aug-20	6.858	2.285	40.74	81.347	32.4	4.591	2.382
116	Trinidad and Tobago	North America	May-20	82.173	5.001	87.04	266.886	36.2	10.014	5.819
16	Brunei	Asia	May-20	320.012	2.286	52.78	81.347	32.4	4.591	2.382
7	Bahamas	North America	May-20	256.835	25.429	92.59	39.497	34.3	8.996	5.200

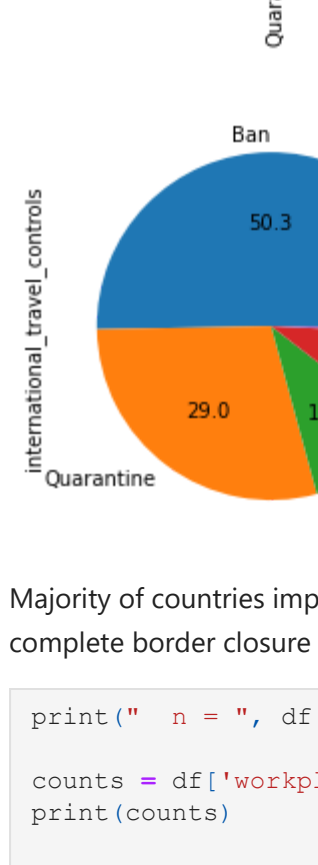
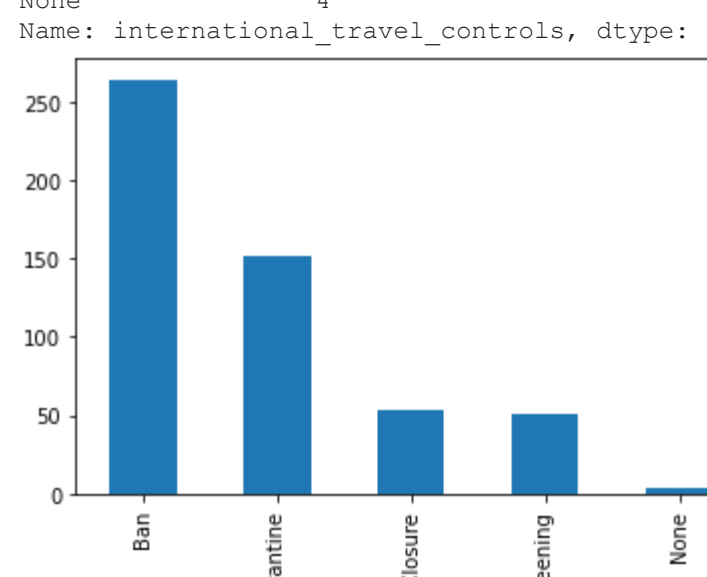
5 rows x 23 columns

The average stringency index is 1.03, meaning, on average across the year, 1 person was infecting just above 1 other person. From looking at bar chart and the value of skewness, can see that values are slightly negatively skewed. Most values between 1 and 1.25. Some outliers over 1.45 and below 0.75. Outliers below 0.12 could be negligible as there were cases in those countries during that period, so the low r rate could be down to the limited contact tracing. Outliers at the top include Nepal which had the largest median R value for any 3 month period of 1.8, meaning that every infected person would on average infect another 1.8 people.

```
In [61]: print("n = ", df['mortality'].count())
counts = df['mortality'].value_counts()
print(counts)

df['mortality'].value_counts().plot.bar()
plt.show()
df['mortality'].value_counts().plot.pie(autopct = '%1.1f')
plt.show()
```

n = 525
Ban 304
High 221
Name: mortality, dtype: int64

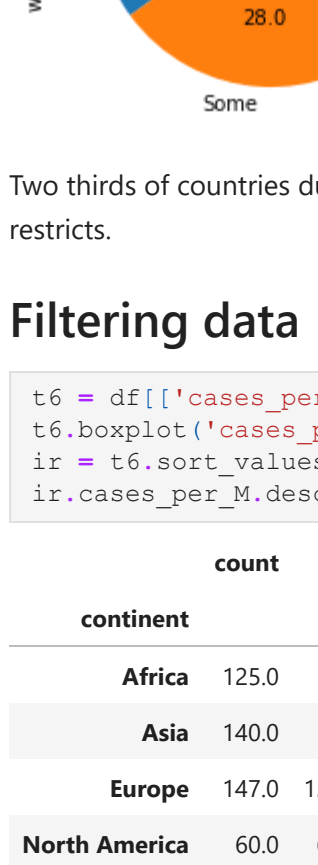
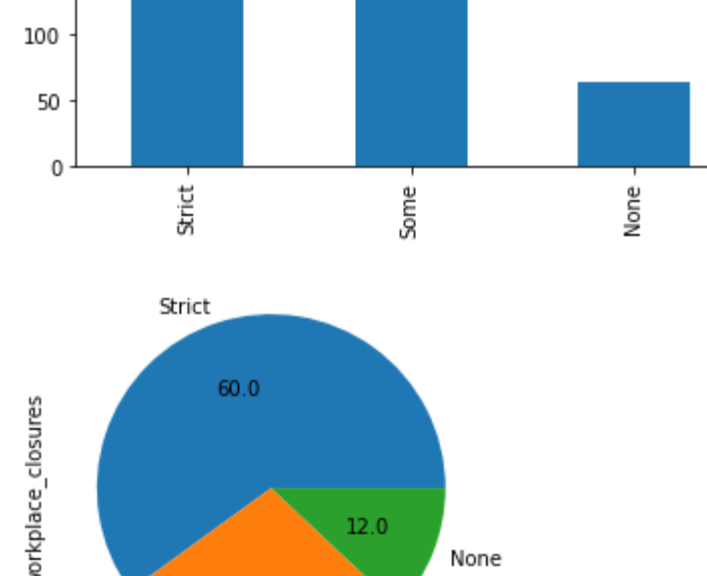


Over the course of the year most countries had a low mortality rate during the 3 month period. Low mortality rate of 57.9% compared to a high mortality rate of 42.1%.

```
In [62]: print("n = ", df['international_travel_controls'].count())
counts = df['international_travel_controls'].value_counts()
print(counts)

df['international_travel_controls'].value_counts().plot.bar()
plt.show()
df['international_travel_controls'].value_counts().plot.pie(autopct = '%1.1f')
plt.show()
```

n = 525
Ban 264
Quarantine 152
BorderClosure 54
Screening 51
None 4
Name: international_travel_controls, dtype: int64

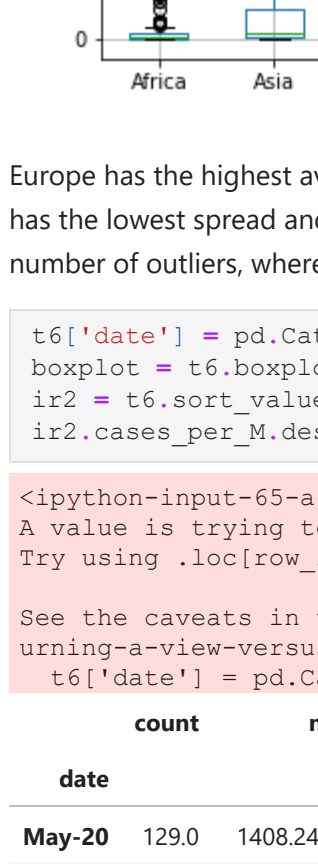
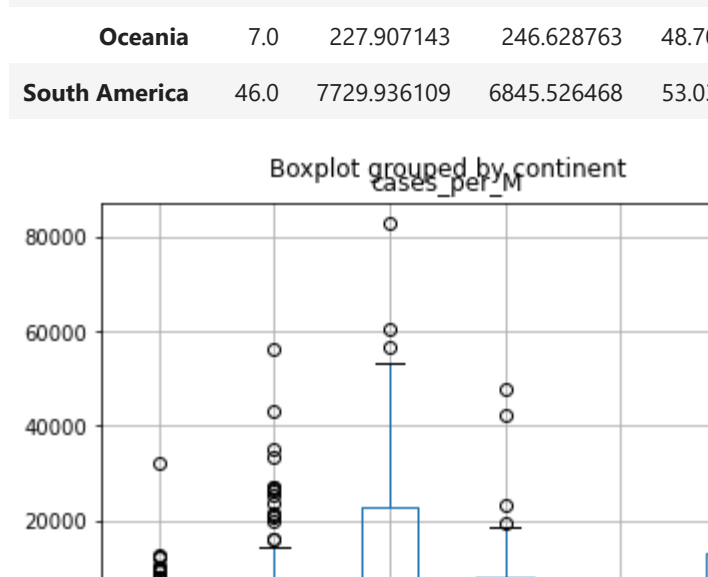


Majority of countries imposed bans on countries, that were considered to be high risk, from international travel. Whereas 10.3% imposed complete border closure with no travel in or out. Less than 1% imposed no travel restrictions during the 3 month period.

```
In [63]: print("n = ", df['workplace_closures'].count())
counts = df['workplace_closures'].value_counts()
print(counts)

df['workplace_closures'].value_counts().plot.bar()
plt.show()
df['workplace_closures'].value_counts().plot.pie(autopct = '%1.1f')
plt.show()
```

n = 525
Strict 315
Some 147
None 63
Name: workplace_closures, dtype: int64



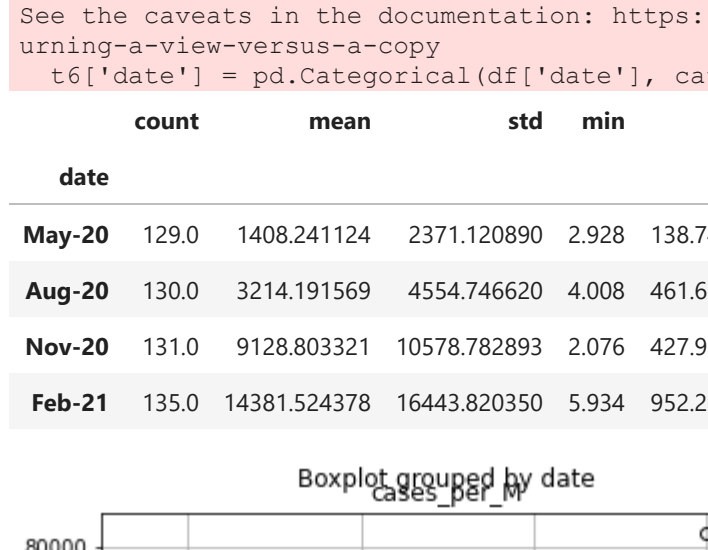
Two thirds of countries during a 3 month period imposed strict workplace closures. Whilst 12% imposed none, with 28% imposing some restrictions.

Filtering data

```
In [64]: t6 = df[['cases_per_M', 'continent', 'date']]
t6.plot('cases_per_M', by = 'continent')
t6 = t6.sort_values('cases_per_M').groupby('continent')
t6.describe()
```

```
Out[64]:
```

	count	mean	std	min	25%	50%	75%	max
Africa	1250	1765.336880	4000.299981	8.095	115.71800	322.346	1165.90200	32245.271
Asia	1400	5239.409607	9251.070900	2.076	237.38000	1315.724	6336.67325	56092.623
Europe	1470	13718.374891	15752.497741	17.955	1492.50350	6246.246	22788.93400	83076.430
North America	600	6605.486000	9205.816195	82.173	619.64500	3543.546	7969.69900	47971.498
Oceania	70	227.907143	246.628763	48.706	70.71400	82.078	296.52800	730.082
South America	460	7729.936109	6845.526468	53.032	1682.02575	6294.498	13202.76675	24782.046



Europe has the highest average cases per million people (13,700), as well as the highest value of values. Oceania has the lowest spread and average (227), could be down to there only being data for seven 3-month periods. Largest and Africa have a number of outliers, whereas South America has none.

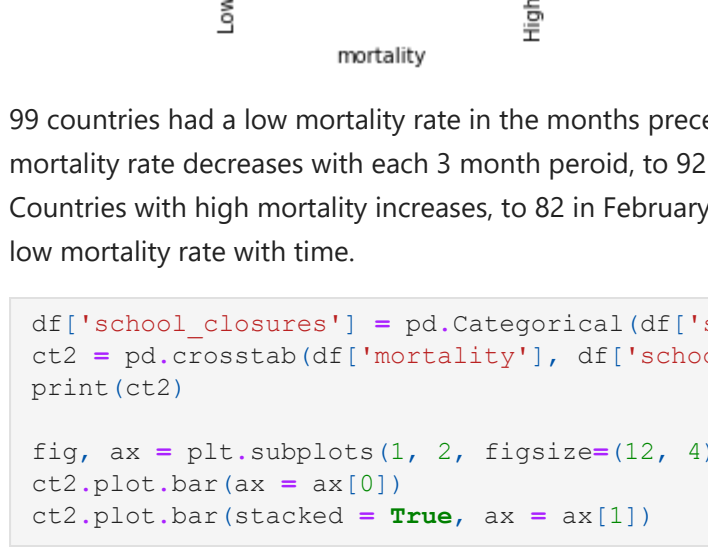
```
In [65]: t6['date'] = pd.Categorical(df['date'], categories=['May-20', 'Aug-20', 'Nov-20', 'Feb-21'])
boxplot = t6.boxplot('cases_per_M', by = 'date')
t6 = t6.sort_values('cases_per_M').groupby('date')
t6.describe()
```

ipython-input-65-a8d1ab3c3de1: SettingWithCopyWarning: A value is being set on a copy of a slice from 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
t6['date'] = pd.Categorical(df['date'], categories=['May-20', 'Aug-20', 'Nov-20', 'Feb-21'])
```

```
Out[65]:
```

	count	mean	std	min	25%	50%	75%	max
May-20	1290	1408.241124	2371.120690	2.928	138.74200	387.9010	1972.37800	19752.105
Aug-20	1300	3214.191560	4554.746620	4.008	461.65925	1279.7225	3825.04475	23844.855
Nov-20	1310	9128.803321	10578.782893	2.076	427.99350	4769.3960	14077.58450	46566.518
Feb-21	1350	14381.524378	16443.820359	5.934	952.28600	8541.4000	24568.58800	83076.430



The average number of cases per million people increases in with time. Averages cases from February 2020 to May 2020 was 1408 per million people. Increased to 14381 per million people for December 2020 to February 2021. Three months preceding February 2021 has the largest spread of values. Lots of outliers in August 2020 where supposedly some countries were hit by an outbreak early than others, by February 2021 the virus has spread more evenly across the world.

Cross-Tabulation

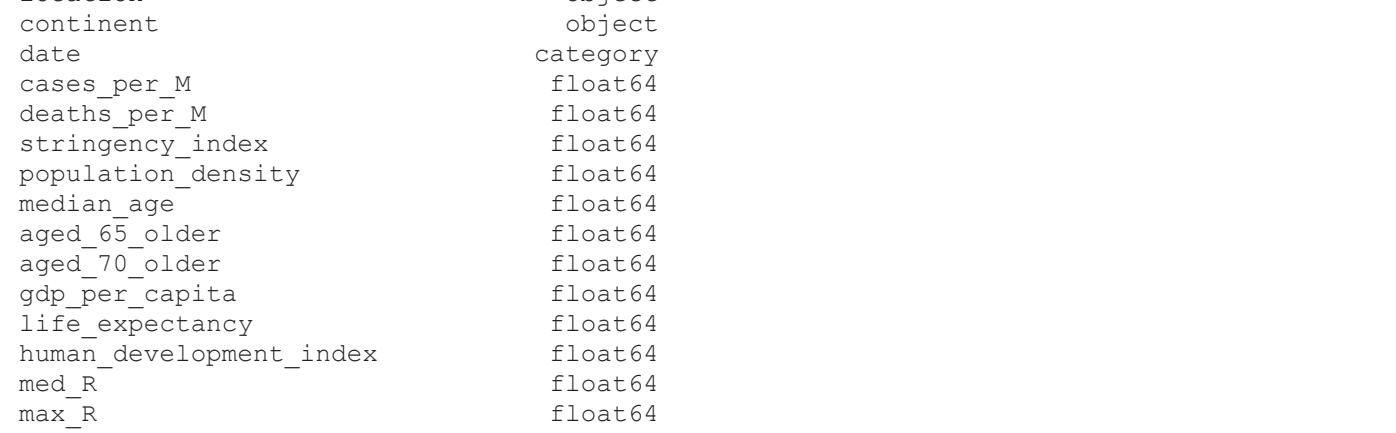
```
In [66]: df['date'] = pd.Categorical(df['date'], categories=['May-20', 'Aug-20', 'Nov-20', 'Feb-21'])
df['mortality'] = pd.Categorical(df['mortality'], categories=['Low', 'High'])
ct = pd.crosstab(df['mortality'], df['date'])
print(ct)
```

```
fig, ax = plt.subplots(1, 2, figsize=(12, 4))
ct1.plot.bar(ax = ax[0])
ct2.plot.bar(stacked = True, ax = ax[1])
```

	date	May-20	Aug-20	Nov-20	Feb-21
Low	mortality	99	92	60	53
High	mortality	30	38	71	82

```
Out[66]:
```

casesSubplot: xlabel='mortality'



92 countries had a low mortality rate in the months preceding May 2020, whereas only 30 had a high mortality rate. Countries with low mortality rate decreases with each 3 month period, to 92 in August 2020, then 60 in November 2020 and finally 53 in February 2021.

Countries with high mortality increases, to 82 in February 2021. Shows trend in increase in high mortality rate with time and decrease in low mortality rate with time.

```
In [67]: df['school_closures'] = pd.Categorical(df['school_closures'], categories=['Minimal', 'Some', 'Required'])
print(df['school_closures'].value_counts())
df['school_closures'].value_counts().plot.bar()
plt.show()
```

```
fig, ax = plt.subplots(1, 2, figsize=(12, 4))
ct1.plot.bar(ax = ax[0])
ct2.plot.bar(stacked = True, ax = ax[1])
```

	date	May-20	Aug-20	Nov-20	Feb-21
Low	mortality	85	137	82	101
High	mortality	37	83	81	82

```
Out[67]:
```

casesSubplot: xlabel='mortality'



When mortality rate is high largest percentage of governments impose school closures, 101 out of 221, with only 37 out of 221 imposing minimal closures.

Data Processing

Creating a new dataset variable.

```
In [68]: df['date_time'] = pd.to_datetime(df['date'], format = '%b-%y')
df.dtypes
```

```
Out[68]:
```

	location	continent	date	cases_per_M	deaths_per_M	stringency_index	population_density	median_age	aged_65_older	aged_70_older
0	Afghanistan	Asia	May-20	390.641	6.602	84.26	54.422	18.6	2.581	1.337
1	Albania	Europe	May-20	394.398	11.120	86.11	104.871	38.0	13.188	8.643
2	Algeria	Africa	May-20	214.202	14.868	76.85	17.348	29.1	6.211	3.857
3	Argentina	South America	May-20	372.823	11.904	90.74	16.177	31.9	11.198	7.441
4	Australia	Oceania	May-20	281.374	4.000	63.43	3.202	37.9	15.504	10.129

5 rows x 25 columns

Removing all missing values from data set before modelling.

```
In [70]: print("Show count of missing values/n", df.isna().all())
df.dropna(inplace = True)
print("Show count of missing values/n", df.isna().all())
```

Show count of missing values

	location	continent	date	cases_per_M	deaths_per_M	stringency_index	population_density	median_age	aged_65_older	aged_70_older	gdp_per_capita	life_expectancy	human_development_index	med_R	max_R	workplace_closures	contact_tracing	resting_policy	cancel_public_events	facial_coverings	international_travel_controls	school_closures	mortality	date_time
0	Afghanistan	Asia	May-20	390.641	6.602	84.26	54.422	18.6	2.581	1.337	11.344	75.2	0.755883	1.039219	1.039219	True	True	True	True	True	True	True	True	2020-05-01
1	Albania	Europe	May-20	394.398	11.120	86.11	104.871	38.0	13.188	8.643	10.544	78.9	0.850815	1.039219	1.039219	True	True	True	True	True	True	True	True	2020-05-01
2	Algeria	Africa	May-20	214.202	14.868	76.85	17.348	29.1	6.211	3.857	10.544	78.9	0.850815	1.039219	1.039219	True	True	True	True	True	True	True	True	2020-05-01
3	Argentina	South America	May-20	372.823	11.904	90.74	16.177	31.9	11.198	7.441	10.544	78.9	0.850815	1.039219	1.039219	True	True	True	True	True	True	True	True	2020-05-01
4	Australia	Oceania	May-20	281.374	4.000	63.43	3.202	37.9	15.504	10.129	10.544	78.9	0.8508											


```
[75]: y = df['deaths_per_M']

X = df.drop(columns= ['deaths_per_M','location','mortality','mortality_rate','datetime','median_age'])
X = pd.get_dummies(X)

X_train, X_test, y_train, y_test = train_test_split(X, y, train_size = 0.5, random_state=42)

rf = RandomForestRegressor(random_state=0)
RFmodel = rf.fit(X_train,y_train)
y_pred = RFmodel.predict(X_test)

print("\nRandom Forest")
print("R^2 = ", metrics.r2_score(y_test, y_pred))
print("MSE = ", metrics.mean_squared_error(y_test, y_pred))
print("MAE = ", metrics.mean_absolute_error(y_test, y_pred))

Random Forest
R^2 = 0.6821563765439793
MSE = 15008.744068883914
MAE = 63.57223638783272
```

Model explains 69.0% of the variability in the data.

There for random forest regression can account for more of the variability in the data than the linear regression model.

Classification models

Creating a decision tree classifier model of deaths per million people against all other variables. Again splitting data 50/50 into training and test sets.

```
In [76]: y = df['mortality']

X = df.drop(columns= ['deaths_per_M','location','mortality','mortality_rate','datetime','median_age'])
X = pd.get_dummies(X)

scaler = preprocessing.StandardScaler()
X = scaler.fit_transform(X)

X_train, X_test, y_train, y_test = train_test_split(X, y, train_size = 0.5, random_state=42)

tree = DecisionTreeClassifier(random_state = 0)
model = tree.fit(X_train,y_train)
y_pred = model.predict(X_test)
print("decision tree")
print("accuracy (on test set) = ", metrics.accuracy_score(y_test, y_pred))
print("Confusion Matrix = \n", metrics.confusion_matrix(y_test, y_pred))

decision tree
accuracy (on test set) = 0.8212927756653993
Confusion Matrix =
[[ 82  28]
 [ 19 134]]

Model accurately predicts 82.9% of mortality classifications. 19 were wrongly predicted as low. 29 wrongly predicted as high.
```

```
In [77]: y = df['mortality']

X = df.drop(columns= ['deaths_per_M','location','mortality','mortality_rate','datetime','median_age'])
X = pd.get_dummies(X)

scaler = preprocessing.StandardScaler()
X = scaler.fit_transform(X)

X_train, X_test, y_train, y_test = train_test_split(X, y, train_size = 0.5, random_state=0)

knn = KNeighborsClassifier(n_neighbors = 10)
kNNmodel = knn.fit(X_train,y_train)
y_pred = kNNmodel.predict(X_test)
print("kNN")
print("accuracy (on test set) = ", metrics.accuracy_score(y_test, y_pred))
print("Confusion Matrix = \n", metrics.confusion_matrix(y_test, y_pred))

kNN
accuracy (on test set) = 0.7604562737642585
Confusion Matrix =
[[ 87  16]
 [ 19 133]]

Model accurately predicts 76.0% of mortality classifications. 19 were wrongly predicted as low. 16 wrongly predicted as high.
```

Comparing the two models, the decision tree model is more accurate than the K nearest neighbours model.