

```
In [129]: 1 import numpy as np
          2 import pandas as pd
          3 import os
```

```
In [140]: 1 # Get COVID-19 Data from John Hopkins CSSE
          2 df_confirmed = pd.read_csv('https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_2020_03_27.csv')
          3 df_deaths = pd.read_csv('https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_time_series/time_series_2020_03_27.csv')
```

```
In [215]: 1 df_confirmed.columns

Index(['Province/State', 'Country/Region', 'Lat', 'Long', '1/22/20', '1/23/20',
      '1/24/20', '1/25/20', '1/26/20', '1/27/20', '1/28/20', '1/29/20',
      '1/30/20', '1/31/20', '2/1/20', '2/2/20', '2/3/20', '2/4/20', '2/5/20',
      '2/6/20', '2/7/20', '2/8/20', '2/9/20', '2/10/20', '2/11/20', '2/12/20',
      '2/13/20', '2/14/20', '2/15/20', '2/16/20', '2/17/20', '2/18/20',
      '2/19/20', '2/20/20', '2/21/20', '2/22/20', '2/23/20', '2/24/20',
      '2/25/20', '2/26/20', '2/27/20', '2/28/20', '2/29/20', '3/1/20',
      '3/2/20', '3/3/20', '3/4/20', '3/5/20', '3/6/20', '3/7/20', '3/8/20',
      '3/9/20', '3/10/20', '3/11/20', '3/12/20', '3/13/20', '3/14/20',
      '3/15/20', '3/16/20', '3/17/20', '3/18/20', '3/19/20', '3/20/20',
      '3/21/20', '3/22/20', '3/23/20', '3/24/20', '3/25/20', '3/26/20',
      '3/27/20', '3/28/20', '3/29/20', '3/30/20', '3/31/20'],
      dtype='object')
```

```
In [207]: 1 # To Separate the training data sets 1/22/2020 - 3/18/2020
          2 dates_train = list(df_confirmed.columns[4:61])
          3 confirmed_train = pd.DataFrame(np.array(df_confirmed[df_confirmed.columns[4:61]]))
          4 deaths_train = pd.DataFrame(np.array(df_deaths[df_deaths.columns[4:61]]))
          5 p_before_train = list(df_confirmed['Province/State'])
          6 p_train = pd.DataFrame([val for val in p_before_train for i in range(len(dates_train))])
          7 c_before_train = list(df_confirmed['Country/Region'])
          8 c_train = pd.DataFrame([val for val in c_before_train for i in range(len(dates_train))])
          9 d_train = pd.DataFrame(dates_train*len(df_confirmed['Province/State']))
         10 train = pd.concat([p_train,c_train,d_train,confirmed_train,deaths_train])
         11 train.columns = ['Province_State', 'Country_Region', 'Date', 'ConfirmedCases', 'Deaths']
         12 train
         13 test.to_csv('train.csv')
```

```
In [202]: 1 # To Separate the test data sets 3/18/2020 - 3/31/2020
2 dates_test = list(df_confirmed.columns[61:])
3 confirmed_test = pd.DataFrame(np.array(df_confirmed[df_confirmed.columns[61:]]).
4 deaths_test = pd.DataFrame(np.array(df_deaths[df_deaths.columns[61:]]).
5 p_before_test = list(df_confirmed['Province/State'])
6 p_test = pd.DataFrame([val for val in p_before_test for i in range(len(
7 c_before_test = list(df_confirmed['Country/Region'])
8 c_test = pd.DataFrame([val for val in c_before_test for i in range(len(
9 c_train = [val for val in c_before_train for i in range(len(dates_test)
10 d_test = pd.DataFrame(dates_test*len(df_confirmed['Province/State']))
11 test = pd.concat([p_test,c_test,d_test,confirmed_test,deaths_test],axis
12 test.columns = ['Province_State', 'Country_Region','Date','ConfirmedCas
13 test.to_csv('test.csv')
```

```
In [208]: 1 # Load Covid Data
2 df = pd.read_csv('train.csv', sep=',')
3 df['Date'] = pd.to_datetime(df['Date'])
4 train_last_date = df.Date.unique()[-1]
5 print(f"Dataset has training data untill : {train_last_date}")
```

Dataset has training data untill : 2020-03-31T00:00:00.000000000

```
In [209]: 1 wpop = pd.read_csv('WPP2019_PopulationByAgeSex_Medium.csv')
2 country_mapper = {
3     'Iran (Islamic Republic of)' : "Iran",
4     'Bolivia (Plurinational State of)' : 'Bolivia',
5     'Brunei Darussalam' : 'Brunei',
6     'Congo' : 'Congo (Kinshasa)',
7     'Democratic Republic of the Congo' : "Congo (Brazzaville)",
8     "Côte d'Ivoire": "Cote d'Ivoire",
9     "Gambia" : "Gambia, The",
10    "Republic of Korea": "Korea, South",
11    "Republic of Moldova": "Moldova",
12    'Réunion' : "Reunion",
13    'Russian Federation' : "Russia",
14    'China, Taiwan Province of China' : "Taiwan*",
15    "United Republic of Tanzania": "Tanzania",
16    "Bahamas": "The Bahamas",
17    "Gambia": "The Gambia",
18    "United States of America (and dependencies)" : "US",
19    "Venezuela (Bolivarian Republic of)" : "Venezuela",
20    'Viet Nam' : "Vietnam"}
21 def rename_countries(x, country_dict):
22     new_name = country_dict.get(x)
23     if new_name is not None:
24         #print(x, "-->", new_name)
25         return new_name
26     else:
27         return x
28 wpop = wpop[wpop['Time']==2020].reset_index(drop=True)
29 wpop['Location'] = wpop.Location.apply(lambda x : rename_countries(x, c
30 clean_wpop = wpop[wpop['Location'].isin(df['Country_Region'].unique())]
31 population_distribution = []
32 for country, gpdf in clean_wpop.groupby("Location"):
33     aux = {f"age_{age_grp}": tot for age_grp, tot in zip(gpdf.AgeGrp, c
34     aux["Country_Region"] = country
35     population_distribution.append(aux)
36 df_pop_distrib = pd.DataFrame(population_distribution)
37 # add missing countries with median values
38 no_data = []
39 for country in df['Country_Region'].unique():
40     if country not in df_pop_distrib['Country_Region'].unique():
```

```
41         aux = df_pop_distrib.drop('Country_Region', axis=1).median(axis=1)
42         aux["Country_Region"] = country
43         no_data.append(aux)
44     df_no_data = pd.DataFrame(no_data)
45     df_pop_distrib = pd.concat([df_pop_distrib, df_no_data], axis=0)
46     # normalize features
47     norm_pop_distrib = df_pop_distrib.drop("Country_Region", axis=1).div(df_pop_distrib['total_pop'], axis=1)
48     norm_pop_distrib['total_pop'] = df_pop_distrib['total_pop']
49     norm_pop_distrib["Country_Region"] = df_pop_distrib["Country_Region"]
50     del df_pop_distrib
51     del df_no_data
52     del clean_wpop
53     del wpop
54     df = df.merge(norm_pop_distrib, on="Country_Region", how='left')
```

```
In [210]: 1 # Data From: https://ourworldindata.org/smoking#prevalence-of-smoking-a
2 smokers = pd.read_csv('share-of-adults-who-smoke.csv')
3 smokers = smokers[smokers.Year == 2016].reset_index(drop=True)
4 smokers_country_dict = {'North America' : "US",
5   'Gambia' : "The Gambia",
6   'Bahamas': "The Bahamas",
7   "'South Korea'" : "Korea, South",
8   'Papua New Guinea' : "Guinea",
9   "'Czech Republic'" : "Czechia",
10  'Congo' : "Congo (Brazzaville)"}
11 smokers['Entity'] = smokers.Entity.apply(lambda x : rename_countries(x,
12 no_datas_smoker = []
13 for country in df['Country_Region'].unique():
14     if country not in smokers.Entity.unique():
15         mean_score = smokers[['Smoking prevalence, total (ages 15+) (%)
16         mean_score['Entity'] = country
17         no_datas_smoker.append(mean_score)
18 no_data_smoker_df = pd.DataFrame(no_datas_smoker)
19 clean_smoke_data = pd.concat([smokers, no_data_smoker_df], axis=0)[['En
20 clean_smoke_data.rename(columns={"Entity": "Country_Region",
21                               "Smoking prevalence, total (ages 15+)
22 df = df.merge(clean_smoke_data, on="Country_Region", how='left')
```

/home/lechen/anaconda3/envs/lechen/lib/python3.7/site-packages/ipykernel_launcher.py:22: FutureWarning: catenation axis is not aligned. A future version of pandas will change to not sort by default.

To accept the future behavior, pass 'sort=False'.

To retain the current behavior and silence the warning, pass 'sort=True'.

```

In [211]: 1 # Add Smokers Percentages By Country
2 # Data From: https://ourworldindata.org/smoking#prevalence-of-smoking-a
3 smokers = pd.read_csv('share-of-adults-who-smoke.csv')
4 smokers = smokers[smokers.Year == 2016].reset_index(drop=True)
5 smokers_country_dict = {'North America' : "US",
6   'Gambia' : "The Gambia",
7   'Bahamas': "The Bahamas",
8   "'South Korea'" : "Korea, South",
9   'Papua New Guinea' : "Guinea",
10  "'Czech Republic'" : "Czechia",
11  'Congo' : "Congo (Brazzaville)"}
12 smokers['Entity'] = smokers.Entity.apply(lambda x : rename_countries(x,
13 no_datas_smoker = []
14 for country in df['Country_Region'].unique():
15     if country not in smokers.Entity.unique():
16         mean_score = smokers[['Smoking prevalence, total (ages 15+) (%)',
17         mean_score['Entity'] = country
18         no_datas_smoker.append(mean_score)
19 no_data_smoker_df = pd.DataFrame(no_datas_smoker)
20 clean_smoke_data = pd.concat([smokers, no_data_smoker_df], axis=0)[['Entity',
21 clean_smoke_data.rename(columns={"Entity": "Country_Region",
22                               "Smoking prevalence, total (ages 15+) (%)": "Smoking prevalence, total (ages 15+) (%)",
23 df = df.merge(clean_smoke_data, on="Country_Region", how='left')

```

/home/lechen/anaconda3/envs/lechen/lib/python3.7/site-packages/ipykernel_launcher.py:23: FutureWarning: concatenation axis is not aligned. A future version of pandas will change to not sort by default.

To accept the future behavior, pass 'sort=False'.

To retain the current behavior and silence the warning, pass 'sort=True'.

```
In [212]: 1 # Concatenate Country and Region Province
          2 def concat_country_province(country, province):
          3     if not isinstance(province, str):
          4         return country
          5     else:
          6         return country+"_"+province
          7
          8 # Concatenate region and province for training
          9 df["Country_Region"] = df[["Country_Region", "Province_State"]].apply(1
```

```

In [213]: 1 # Add Time Data information from Quarantine, Restrictions and Schools
2 # Data From: https://www.kaggle.com/koryto/countryinfo
3 country_info = pd.read_csv('covid19countryinfo.csv')
4 country_info = country_info[~country_info.country.isnull()].reset_index
5 country_info.drop([ c for c in country_info.columns if c.startswith("Ur
6 country_info.drop(columns=['pop', 'sex0', 'sex14', 'sex25', 'sex54', 's
7                 axis=1,
8                 inplace=True)
9 # Columns with dates
10 country_info["quarantine"] = pd.to_datetime(country_info["quarantine"])
11 country_info["restrictions"] = pd.to_datetime(country_info["restriction
12 country_info["schools"] = pd.to_datetime(country_info["schools"])
13 same_state = []
14 for country in df["Province_State"].unique():
15     if country in country_info.country.unique():
16         same_state.append(country)
17     else:
18         pass
19 country_to_state_country = {}
20 for state in same_state:
21     #print(state)
22     #print(df[df["Province_State"]==state]["Country/Region"].unique())
23     #print("----")
24     country_to_state_country[state] = df[df["Province_State"]==state]["
25 country_info['country'] = country_info.country.apply(lambda x : rename_
26 coutry_merge_info = country_info[["country", "density", "urbanpop", "hc
27 cols_median = ["density", "urbanpop", "hospibed", "lung", "femalelung",
28 coutry_merge_info.loc[:, cols_median] = coutry_merge_info.loc[:, cols_m
29 merged = df.merge(coutry_merge_info, left_on="Country_Region", right_on
30 merged.loc[:, cols_median] = merged.loc[:, cols_median].apply(lambda x:
31 country_dates_info = country_info[["country", "restrictions", "quaranti
32 def update_dates(a_df, col_update):
33     """
34     This creates a boolean time series with one after the start of conf
35     """
36     gpdf = a_df.groupby("Country_Region")
37     new_col = gpdf.apply(lambda df : df[col_update].notnull().cumsum())
38     a_df[col_update] = new_col
39 for col in ["restrictions", "quarantine", "schools"]:
40     print(merged.shape)

```



```

41 merged = merged.merge(country_dates_info[["country", col]],
42                        left_on=["Country_Region", "Date"],
43                        right_on=["country", col],
44                        how="left",
45                        )
46 update_dates(merged, col)
47 drop_country_cols = [x for x in merged.columns if x.startswith("country
48 merged.drop(columns=drop_country_cols, axis=1, inplace=True)

```

```
(3354, 37)
```

```
(3354, 39)
```

```
(3354, 41)
```

```

In [214]: 1 merged.to_csv('enriched_covid_19.csv', index=None)
          2 merged

```

	Unnamed: 0	Province_State	Country_Region	Date	ConfirmedCases	Fatalities	age_0-4
0	0	NaN	Afghanistan	2020-03-19	22	0	0.145717
1	1	NaN	Afghanistan	2020-03-20	24	0	0.145717
2	2	NaN	Afghanistan	2020-03-21	24	0	0.145717
3	3	NaN	Afghanistan	2020-03-22	40	1	0.145717
4	4	NaN	Afghanistan	2020-03-23	40	1	0.145717
...
3349	3349	NaN	Malawi	2020-03-27	0	0	0.152846
3350	3350	NaN	Malawi	2020-03-28	0	0	0.152846
3351	3351	NaN	Malawi	2020-03-29	0	0	0.152846
3352	3352	NaN	Malawi	2020-03-30	0	0	0.152846
3353	3353	NaN	Malawi	2020-03-31	0	0	0.152846

```
3354 rows × 39 columns
```

```
In [ ]: 1
```