

COVID-19 Data - Data Cleaning with SQL

```
In [1]: import sqlite3
import numpy as np
import pandas as pd
```

```
In [2]: # Create SQL Engine, Connection, and Cursor
# and connect to database file.
connection = sqlite3.connect('covid_large_dataset.db')
cursor = connection.cursor()
```

```
In [3]: # Query data from sqlite3 database for dates, states, and count totals per state/date
command1 = """
SELECT DISTINCT
    case_month,
    res_state,
    count(res_state) AS state_total
FROM
    covid_data
WHERE
    case_month IS NOT NULL AND res_state IS NOT NULL
GROUP BY
    case_month, res_state
ORDER BY
    case_month, res_state;
"""
```

```
In [4]: # Import filter data from SQL to DataFrame.
# Execute command and read into DataFrame.
df = pd.read_sql(sql=command1, con=connection)
```

```
In [5]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1459 entries, 0 to 1458
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   case_month  1459 non-null  object
 1   res_state   1459 non-null  object
 2   state_total 1459 non-null  int64
dtypes: int64(1), object(2)
memory usage: 34.3+ KB
```

```
In [6]: df.head()
```

Out[6]:

	case_month	res_state	state_total
0	2020-01	AL	103
1	2020-01	AR	17
2	2020-01	AZ	63
3	2020-01	CA	389
4	2020-01	CO	85

In [7]: *# Create a list of sorted dates and a list of sorted states
for use as index and columns of new table.*

```
dates = df['case_month'].value_counts()
states = df['res_state'].value_counts()
```

```
dates.sort_index(inplace=True)
states.sort_index(inplace=True)
```

```
dates_list = list(dates.index)
states_list = list(states.index)
```

In [8]: *# Create the new table.*

```
df2 = pd.DataFrame(index=states_list, columns=dates_list)
```

Assign values to each date, state pair in new table.

```
for index, row in df.iterrows():
    state = row["res_state"]
    date = row["case_month"]
    total = row["state_total"]
    df2[date][state] = total
```

In [9]: df2.head(10)

Out[9]:

	2020-01	2020-02	2020-03	2020-04	2020-05	2020-06	2020-07	2020-08	2020-09	2020-10
AK	NaN	NaN	239	137	146	614	2293	2104	3068	936
AL	103	59	2677	6208	11427	21971	49743	37556	27988	3794
AR	17	18	1156	2742	5354	16458	21307	18628	21805	2896
AZ	63	57	2679	7243	16190	71744	86892	22845	16826	3341
CA	389	489	19293	46525	71232	174261	282735	151891	101378	12494
CO	85	98	6793	12128	9287	7091	14650	9801	14372	4331
CT	17	36	3766	9257	3629	29767	4007	3505	4564	1634
DC	NaN	36	1232	3305	2175	566	4361	712	633	79
DE	NaN	NaN	300	4242	4384	1929	3134	2397	2921	427
FL	132	290	16081	20817	21705	81505	135395	149549	111662	12942

10 rows × 28 columns

In []: