

```
import pandas as pd
# Read data. This data represents the cumulative known cases to date (https://covidtracking.com/about-data/faq)
url = 'https://raw.githubusercontent.com/COVID19Tracking/covid-tracking-data/master/data/states_daily_4pm_et.csv'
df = pd.read_csv(url,index_col=0,parse_dates=[0])

df.head(5)
```

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | inCuCurrently | inCuCumulative | onVentilatorCurrently | onVentilatorCumulative | recovered | hash | dateChecked | death | hospitalized | total | totalTestResults |
|------------|-------|----------|----------|---------|-----------------------|------------------------|---------------|----------------|-----------------------|------------------------|-----------|--|----------------------|-------|--------------|---------|------------------|
| | date | | | | | | | | | | | | | | | | |
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | NaN | NaN | NaN | NaN | 262.0 | d19bb8087806f8dded75fb6165d7bed1bcface44 | 2020-05-03T20:00:00Z | 9.0 | NaN | 21578.0 | 21578.0 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 403.0 | NaN | 242.0 | NaN | aa4e1102894dd9528461e4d910fab397cc40d31b | 2020-05-03T20:00:00Z | 290.0 | 1035.0 | 92500.0 | 92500.0 |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | NaN | NaN | 20.0 | 88.0 | 1999.0 | 6b15aab296af49161db39fc69eef7c6ca244994 | 2020-05-03T20:00:00Z | 76.0 | 427.0 | 52890.0 | 52890.0 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ccae2e571ac93c5836a22158f57c12313127810d | 2020-05-03T20:00:00Z | 0.0 | NaN | 57.0 | 57.0 |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 282.0 | NaN | 192.0 | NaN | 1597.0 | c687661dcf3892bcc23c38ddcb46fe8096006f1 | 2020-05-03T20:00:00Z | 362.0 | 1348.0 | 81119.0 | 81119.0 |

Double-click (or enter) to edit

```
# Drop total, postNeg, and hospitalized columns as they are redundant
# Drop other columns that will not be used
df_drop = df.drop(columns = [6, 7, 8, 9, 11, 12, 14, 15, 17, 18, 19, 20, 21, 22, 23])
df_drop = df.drop(columns = ['inCuCurrently', 'inCuCumulative',
                             'onVentilatorCurrently', 'onVentilatorCumulative',
                             'hash', 'dateChecked', 'hospitalized', 'total',
                             'postNeg', 'fips', 'deathIncrease',
                             'hospitalizedIncrease', 'negativeIncrease',
                             'positiveIncrease', 'totalTestResultsIncrease'])
df_drop.head()
```

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|
| | date | | | | | | | | |
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 290.0 | 92500.0 |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | 1999.0 | 76.0 | 52890.0 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | 0.0 | 57.0 |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 1597.0 | 362.0 | 81119.0 |

```
# Create new features
# Divide positive by totalTestResults to get positive_percent
df_drop["percent_positive"] = ""
df_drop["percent_positive"] = 100*df_drop["positive"]/df_drop["totalTestResults"]
df_drop.head()
```

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|------------------|
| | date | | | | | | | | | |
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 | 1.705441 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 290.0 | 92500.0 | 8.351351 |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | 1999.0 | 76.0 | 52890.0 | 6.487049 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | 0.0 | 57.0 | 0.000000 |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 1597.0 | 362.0 | 81119.0 | 10.651019 |

```
# Divide hospitalized by positive to get hospitalized_percent
import numpy as np
df_drop["hospitalized_percent"] = ""
df_drop["hospitalized_percent"] = np.nanmax(df_drop[['hospitalizedCurrently','hospitalizedCumulative']], axis=1)
df_drop["hospitalized_percent"] = 100*df_drop["hospitalized_percent"]/df_drop["positive"]
df_drop.head()
```

/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:3: RuntimeWarning: All-NaN axis encountered
This is separate from the ipykernel package so we can avoid doing imports until

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive | hospitalized_percent |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|------------------|----------------------|
| | date | | | | | | | | | | |
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 | 1.705441 | 3.260870 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 290.0 | 92500.0 | 8.351351 | 13.398058 |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | 1999.0 | 76.0 | 52890.0 | 6.487049 | 12.445351 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | 0.0 | 57.0 | 0.000000 | NaN |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 1597.0 | 362.0 | 81119.0 | 10.651019 | 15.601852 |

```
# Divide recovered by positive to get recovered_percent
df_drop["recovered_percent"] = ""
df_drop["recovered_percent"] = 100*df_drop["recovered"]/df_drop["positive"]
df_drop.head()
```

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive | hospitalized_percent | recovered_percent |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|------------------|----------------------|-------------------|
| | date | | | | | | | | | | | |
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 | 1.705441 | 3.260870 | 71.195652 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 290.0 | 92500.0 | 8.351351 | 13.398058 | NaN |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | 1999.0 | 76.0 | 52890.0 | 6.487049 | 12.445351 | 58.262897 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | 0.0 | 57.0 | 0.000000 | NaN | NaN |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 1597.0 | 362.0 | 81119.0 | 10.651019 | 15.601852 | 18.483796 |

```
# Divide death by positive to get death_percent
df_drop["death_percent"] = ""
df_drop["death_percent"] = 100*df_drop["death"]/df_drop["positive"]
df_drop.head()
```

```
# Fetch the latest state population data (nst-est2019-01.csv)
from google.colab import files
uploaded = files.upload()

# Choose Files | nst-est2019-01.csv
• nst-est2019-01.csv(application/vnd.ms-excel) - 676 bytes, last modified: 4/13/2020 - 100% done
Saving nst-est2019-01.csv to nst-est2019-01.csv
....

# Load latest state population data
import io
df_state_pop = pd.read_csv(io.StringIO(uploaded['nst-est2019-01.csv'].decode('utf-8')))
df_state_pop["Population"] = pd.to_numeric(df_state_pop["Population"])
df_state_pop.head()

# Add column of state populations (population) to df_drop_total_posNeg
# Need to sort rows by state using index numbering from state_list

df_drop["population"] = ""

for i in range(len(df_drop)):
    for index in range(len(df_state_pop)):
        if df_drop.iloc[i, 0] == df_state_pop.iloc[index, 0]:
            df_drop.iloc[i, 13] = df_state_pop.iloc[index, 1]

df_drop[["population"]] = df_drop["population"].apply(pd.to_numeric)

df_drop.head()

# Normalize positive to state population
df_drop["positive_norm"] = ""
df_drop["positive_norm"] = df_drop["positive"]/df_drop["population"]
df_drop.head()

# Normalize hospitalized to state population
df_drop["hospitalized_norm"] = ""
df_drop["hospitalized_norm"] = np.nanmax(df_drop[["hospitalizedCurrently", 'hospitalizedCumulative']], axis=1)
df_drop["hospitalized_norm"] = df_drop["hospitalized_norm"]/df_drop["population"]
df_drop.head()

# Normalize recovered to state population
df_drop["recovered_norm"] = ""
df_drop["recovered_norm"] = df_drop["recovered"]/df_drop["population"]
df_drop.head()

# Normalize death to state population
df_drop["death_norm"] = ""
df_drop["death_norm"] = df_drop["death"]/df_drop["population"]
df_drop.head()
```

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive | hospitalized_percent | recovered_percent | death_percent | population |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|------------------|----------------------|-------------------|---------------|------------|
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 | 1.705441 | 3.260870 | 71.195652 | 2.445652 | 731545.0 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 290.0 | 92500.0 | 8.351351 | 13.398058 | NaN | 3.754045 | 4903185.0 |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | 1999.0 | 76.0 | 52890.0 | 6.487049 | 12.445351 | 58.262897 | 2.215098 | 3017804.0 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | 0.0 | 57.0 | 0.000000 | NaN | NaN | NaN | NaN |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 1597.0 | 362.0 | 81119.0 | 10.651019 | 15.601852 | 18.483796 | 4.189815 | 7278717.0 |

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive | hospitalized_percent | recovered_percent | death_percent | population | positive_norm |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|------------------|----------------------|-------------------|---------------|------------|---------------|
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 | 1.705441 | 3.260870 | 71.195652 | 2.445652 | 731545.0 | 0.000503 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 290.0 | 92500.0 | 8.351351 | 13.398058 | NaN | 3.754045 | 4903185.0 | 0.001576 |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | 1999.0 | 76.0 | 52890.0 | 6.487049 | 12.445351 | 58.262897 | 2.215098 | 3017804.0 | 0.001137 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | 0.0 | 57.0 | 0.000000 | NaN | NaN | NaN | NaN | NaN |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 1597.0 | 362.0 | 81119.0 | 10.651019 | 15.601852 | 18.483796 | 4.189815 | 7278717.0 | 0.001187 |

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive | hospitalized_percent | recovered_percent | death_percent | population | positive_norm | hospitalized_norm |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|------------------|----------------------|-------------------|---------------|------------|---------------|-------------------|
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 | 1.705441 | 3.260870 | 71.195652 | 2.445652 | 731545.0 | 0.000503 | 0.000016 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 290.0 | 92500.0 | 8.351351 | 13.398058 | NaN | 3.754045 | 4903185.0 | 0.001576 | 0.000211 |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | 1999.0 | 76.0 | 52890.0 | 6.487049 | 12.445351 | 58.262897 | 2.215098 | 3017804.0 | 0.001137 | 0.000141 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | 0.0 | 57.0 | 0.000000 | NaN | NaN | NaN | NaN | NaN | NaN |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 1597.0 | 362.0 | 81119.0 | 10.651019 | 15.601852 | 18.483796 | 4.189815 | 7278717.0 | 0.001187 | 0.000185 |

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive | hospitalized_percent | recovered_percent | death_percent | population | positive_norm | hospitalized_norm | recovered_norm |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|------------------|----------------------|-------------------|---------------|------------|---------------|-------------------|----------------|
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 | 1.705441 | 3.260870 | 71.195652 | 2.445652 | 731545.0 | 0.000503 | 0.000016 | 0.000358 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 290.0 | 92500.0 | 8.351351 | 13.398058 | NaN | 3.754045 | 4903185.0 | 0.001576 | 0.000211 | NaN |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | 1999.0 | 76.0 | 52890.0 | 6.487049 | 12.445351 | 58.262897 | 2.215098 | 3017804.0 | 0.001137 | 0.000141 | 0.000662 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | 0.0 | 57.0 | 0.000000 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 1597.0 | 362.0 | 81119.0 | 10.651019 | 15.601852 | 18.483796 | 4.189815 | 7278717.0 | 0.001187 | 0.000185 | 0.000219 |

| | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive | hospitalized_percent | recovered_percent | death_percent | population | positive_norm | hospitalized_norm | recovered_norm | death_norm |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|------------------|----------------------|-------------------|---------------|------------|---------------|-------------------|----------------|------------|
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 | 1.705441 | 3.260870 | 71.195652 | 2.445652 | 731545.0 | 0.000503 | 0.000016 | 0.000358 | 0.000012 |
| 2020-05-03 | AL | 7725.0 | 84775.0 | NaN | NaN | 1035.0 | NaN | 290.0 | 92500.0 | 8.351351 | 13.398058 | NaN | 3.754045 | 4903185.0 | 0.001576 | 0.000211 | NaN | 0.000059 |
| 2020-05-03 | AR | 3431.0 | 49459.0 | NaN | 100.0 | 427.0 | 1999.0 | 76.0 | 52890.0 | 6.487049 | 12.445351 | 58.262897 | 2.215098 | 3017804.0 | 0.001137 | 0.000141 | 0.000662 | 0.000025 |
| 2020-05-03 | AS | 0.0 | 57.0 | NaN | NaN | NaN | NaN | 0.0 | 57.0 | 0.000000 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 2020-05-03 | AZ | 8640.0 | 72479.0 | NaN | 732.0 | 1348.0 | 1597.0 | 362.0 | 81119.0 | 10.651019 | 15.601852 | 18.483796 | 4.189815 | 7278717.0 | 0.001187 | 0.000185 | 0.000219 | 0.000050 |

```
df_state.dropna(inplace=True)
```

```
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 3321 entries, 2020-05-03 to 2020-01-22
Data columns (total 18 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   state               3321 non-null   object
 1   positive            3386 non-null   float64
 2   negative            3148 non-null   float64
 3   pending             677 non-null    float64
 4   hospitalizedCurrently 1191 non-null   float64
 5   hospitalizedCumulative 1239 non-null   float64
 6   recovered           1837 non-null   float64
 7   death              2594 non-null   float64
 8   totalTestResults    3319 non-null   float64
 9   percent_positive    3275 non-null   float64
10   hospitalized_percent 1878 non-null   float64
11   recovered_percent    1837 non-null   float64
12   death_percent       2541 non-null   float64
13   population          3125 non-null   float64
14   positive_norm       3125 non-null   float64
15   hospitalized_norm    1831 non-null   float64
16   recovered_norm      958 non-null    float64
17   death_norm         2447 non-null   float64
dtypes: float64(17), object(1)
memory usage: 573.0+ KB
```

```
# Get the unique values of 'state' column
state_list = df_state.unique()
state_list
```

```
array(['AK', 'AL', 'AR', 'AS', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE', 'FL',
       'GA', 'GU', 'HI', 'IA', 'ID', 'IL', 'IN', 'KS', 'KY', 'LA', 'MA',
       'MD', 'ME', 'MI', 'MN', 'MO', 'MP', 'MS', 'MT', 'NC', 'ND', 'NE',
       'NH', 'NJ', 'NM', 'NV', 'NY', 'OK', 'OR', 'PA', 'PR', 'RI',
       'SC', 'SD', 'TN', 'TX', 'UT', 'VA', 'VT', 'WA', 'WI', 'WV'],
      dtype=object)
```

```
#create a data frame dictionary to store the state data frames
df_state_dict = {}
for key in df_state_dict.keys():
    df_state_dict[key] = df_state[df_state.state == key]
```

```
df_state_dict['AK'].head()
```

| date | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive | hospitalized_percent | recovered_percent | death_percent | population | positive_norm | hospitalized_norm | recovered_norm | death_norm |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|-------|------------------|------------------|----------------------|-------------------|---------------|------------|---------------|-------------------|----------------|------------|
| 2020-05-03 | AK | 368.0 | 21210.0 | NaN | 12.0 | NaN | 262.0 | 9.0 | 21578.0 | 1.705441 | 3.260870 | 71.195652 | 2.445652 | 731545.0 | 0.000503 | 0.000016 | 0.000358 | 0.000012 |
| 2020-05-02 | AK | 365.0 | 21034.0 | NaN | 10.0 | NaN | 261.0 | 9.0 | 21399.0 | 1.705687 | 2.739726 | 71.506849 | 2.465753 | 731545.0 | 0.000499 | 0.000014 | 0.000357 | 0.000012 |
| 2020-05-01 | AK | 364.0 | 19961.0 | NaN | 25.0 | NaN | 254.0 | 9.0 | 20325.0 | 1.790898 | 6.868132 | 69.780220 | 2.472527 | 731545.0 | 0.000498 | 0.000034 | 0.000347 | 0.000012 |
| 2020-04-30 | AK | 355.0 | 18764.0 | NaN | 19.0 | NaN | 252.0 | 9.0 | 19119.0 | 1.856792 | 5.352113 | 70.985915 | 2.535211 | 731545.0 | 0.000485 | 0.000026 | 0.000344 | 0.000012 |
| 2020-04-29 | AK | 355.0 | 18764.0 | NaN | 14.0 | NaN | 240.0 | 9.0 | 19119.0 | 1.856792 | 3.943662 | 67.605634 | 2.535211 | 731545.0 | 0.000485 | 0.000019 | 0.000328 | 0.000012 |

```
df_state_dict['CA'].head()
```

| date | state | positive | negative | pending | hospitalizedCurrently | hospitalizedCumulative | recovered | death | totalTestResults | percent_positive | hospitalized_percent | recovered_percent | death_percent | population | positive_norm | hospitalized_norm | recovered_norm | death_norm |
|------------|-------|----------|----------|---------|-----------------------|------------------------|-----------|--------|------------------|------------------|----------------------|-------------------|---------------|------------|---------------|-------------------|----------------|------------|
| 2020-05-03 | CA | 53616.0 | 662135.0 | NaN | 4734.0 | NaN | NaN | 2215.0 | 715751.0 | 7.490873 | 8.829454 | NaN | 4.131229 | 39512223.0 | 0.001357 | 0.000120 | NaN | 0.000056 |
| 2020-05-02 | CA | 52197.0 | 634006.0 | NaN | 4722.0 | NaN | NaN | 2171.0 | 686803.0 | 7.599996 | 9.046497 | NaN | 4.159243 | 39512223.0 | 0.001321 | 0.000120 | NaN | 0.000055 |
| 2020-05-01 | CA | 50442.0 | 604543.0 | NaN | 4706.0 | NaN | NaN | 2073.0 | 654985.0 | 7.701245 | 9.329527 | NaN | 4.109671 | 39512223.0 | 0.001277 | 0.000119 | NaN | 0.000052 |
| 2020-04-30 | CA | 48917.0 | 576420.0 | NaN | 4981.0 | NaN | NaN | 1982.0 | 625337.0 | 7.822502 | 10.182554 | NaN | 4.051761 | 39512223.0 | 0.001238 | 0.000126 | NaN | 0.000050 |
| 2020-04-29 | CA | 46500.0 | 556639.0 | NaN | 5011.0 | NaN | NaN | 1887.0 | 603139.0 | 7.709666 | 10.776344 | NaN | 4.058065 | 39512223.0 | 0.001177 | 0.000127 | NaN | 0.000048 |

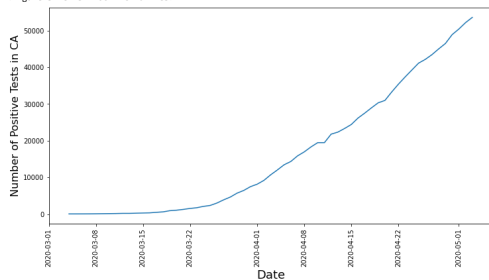
```
from matplotlib import pyplot as plt
```

```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].positive)
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Number of Positive Tests in CA', fontsize=16)
plt.show()
```

```
No handles with labels found to put in legend.
<Figure size 432x288 with 0 Axes>
```



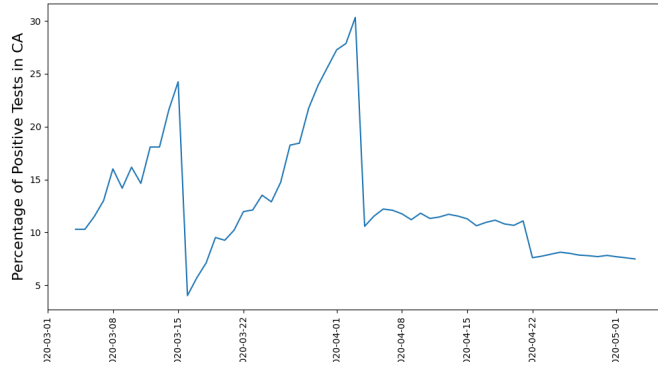
```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].percent_positive)
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Percentage of Positive Tests in CA', fontsize=16)
plt.show()
```

```
<
```

No handles with labels found to put in legend.
<Figure size 640x480 with 0 Axes>

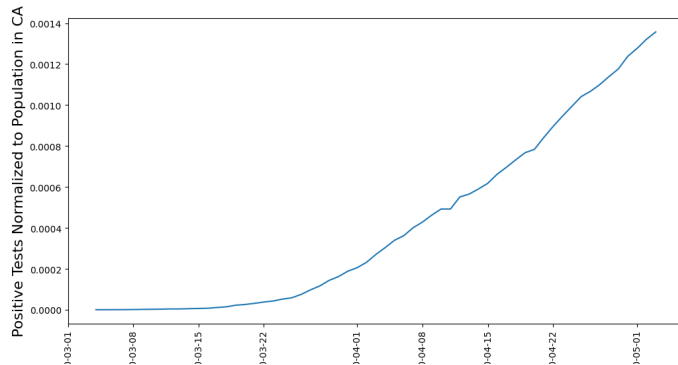


```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].positive_norm)
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Positive Tests Normalized to Population in CA', fontsize=16)
plt.show()
```

No handles with labels found to put in legend.
<Figure size 640x480 with 0 Axes>

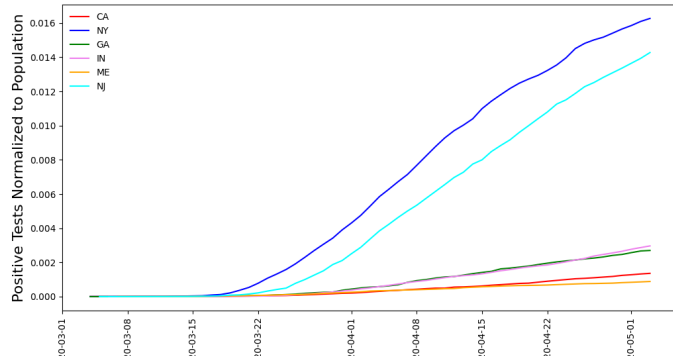


```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].positive_norm, color='red', label='CA')
plt.plot(df_state_dict['NY'].positive_norm, color='blue', label='NY')
plt.plot(df_state_dict['GA'].positive_norm, color='green', label='GA')
plt.plot(df_state_dict['IN'].positive_norm, color='violet', label='IN')
plt.plot(df_state_dict['ME'].positive_norm, color='orange', label='ME')
plt.plot(df_state_dict['NJ'].positive_norm, color='cyan', label='NJ')
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Positive Tests Normalized to Population', fontsize=16)
plt.show()
```

<Figure size 640x480 with 0 Axes>

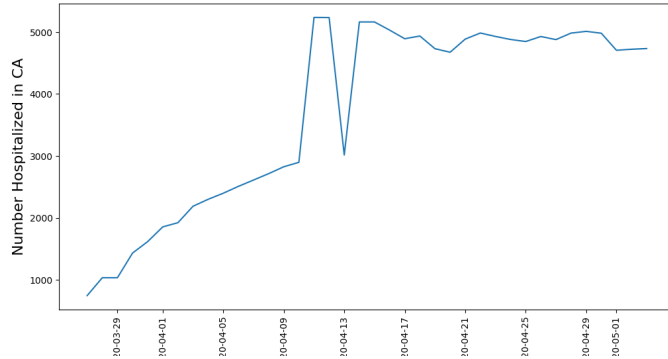


```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].hospitalizedCurrently)
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Number Hospitalized in CA', fontsize=16)
plt.show()
```

No handles with labels found to put in legend.
<Figure size 640x480 with 0 Axes>

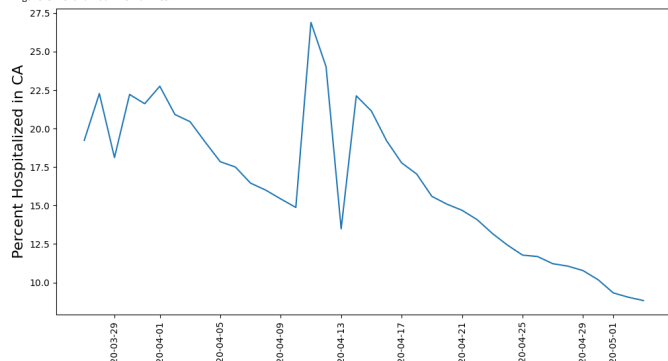


```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].hospitalized_percent)
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Percent Hospitalized in CA', fontsize=16)
plt.show()
```

No handles with labels found to put in legend.
<Figure size 640x480 with 0 Axes>

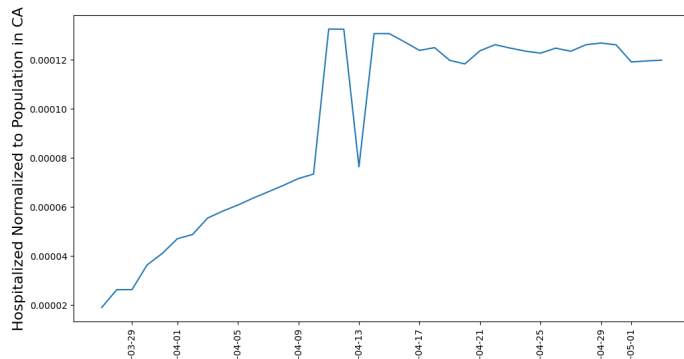


```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].hospitalized_norm)
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Hospitalized Normalized to Population in CA', fontsize=16)
plt.show()
```

No handles with labels found to put in legend.
<Figure size 640x480 with 0 Axes>



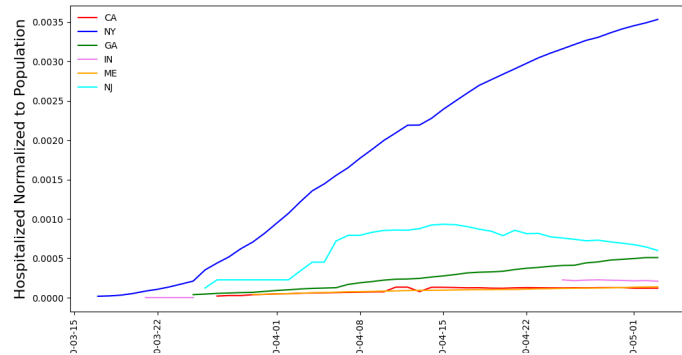
```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].hospitalized_norm, color='red', label='CA')
plt.plot(df_state_dict['NY'].hospitalized_norm, color='blue', label='NY')
plt.plot(df_state_dict['GA'].hospitalized_norm, color='green', label='GA')
plt.plot(df_state_dict['IN'].hospitalized_norm, color='violet', label='IN')
plt.plot(df_state_dict['ME'].hospitalized_norm, color='orange', label='ME')
plt.plot(df_state_dict['NJ'].hospitalized_norm, color='cyan', label='NJ')
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Hospitalized Normalized to Population', fontsize=16)
plt.show()
```

□

<Figure size 640x480 with 0 Axes>



In several states, population normalized hospitalizations plateau, although population normalized death rate continues to grow.

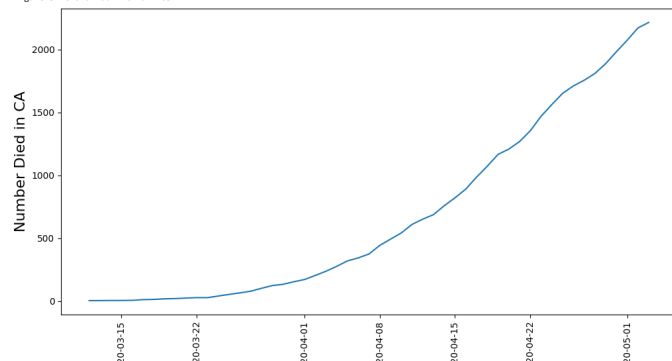
```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].death)
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Number Died in CA', fontsize=16)
plt.show()
```

No handles with labels found to put in legend.

<Figure size 640x480 with 0 Axes>



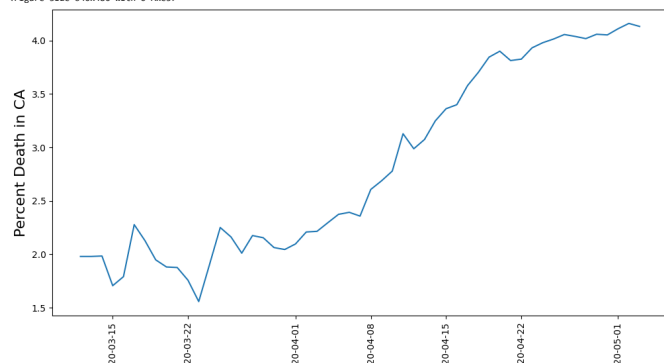
```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].death_percent)
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Percent Death in CA', fontsize=16)
plt.show()
```

No handles with labels found to put in legend.

<Figure size 640x480 with 0 Axes>



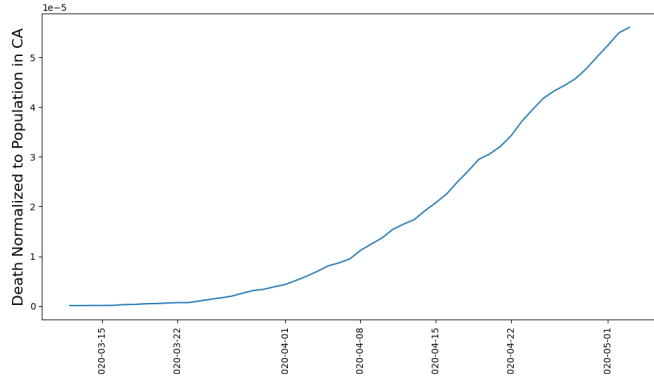
```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].death_norm)
plt.xticks(rotation='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Death Normalized to Population in CA', fontsize=16)
plt.show()
```

No handles with labels found to put in legend.

No handles with labels found to put in legend.
<Figure size 640x480 with 0 Axes>

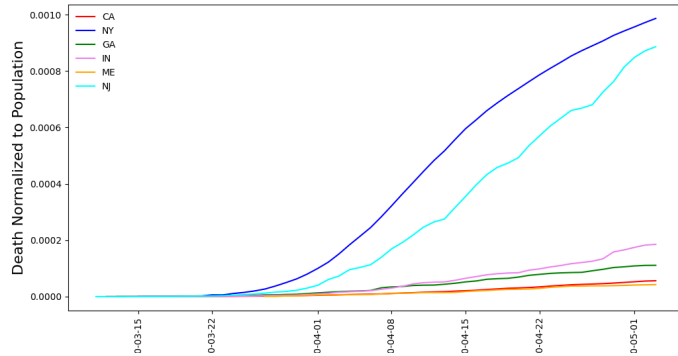


```
fig = plt.figure()
fig, ax = plt.subplots(figsize=(12, 6))
plt.rcParams.update(plt.rcParamsDefault)

plt.plot(df_state_dict['CA'].death_norm, color="red", label="CA")
plt.plot(df_state_dict['NY'].death_norm, color="blue", label="NY")
plt.plot(df_state_dict['GA'].death_norm, color="green", label="GA")
plt.plot(df_state_dict['IN'].death_norm, color="violet", label="IN")
plt.plot(df_state_dict['ME'].death_norm, color="orange", label="ME")
plt.plot(df_state_dict['NJ'].death_norm, color="cyan", label="NJ")
plt.xticks(rotate='vertical')

plt.legend(frameon=False)
plt.xlabel('Date', fontsize=18)
plt.ylabel('Death Normalized to Population', fontsize=16)
plt.show()
```

<Figure size 640x480 with 0 Axes>



Note how the population normalized death curves relate closely to population normalized positive test curves

Curve fitting done at: <http://www.xuru.org/rt/NLR.asp#CopyPaste>

```
# Fetch the parameters for each state (AexpBx^1.csv) that fit to positive_norm = a*exp(b/x)
# where x is the number of days from March 4, 2020
from google.colab import files
uploaded = files.upload()
```

Choose Files | AexpBx^1.csv
AexpBx^1.csv(application/vnd.ms-excel) - 2391 bytes, last modified: 5/3/2020 - 100% done
Saving AexpBx^1.csv to AexpBx^1.csv

```
# Load the parameters for each state (AexpBx^1.csv) that fit to positive_norm = a*exp(b/x)
import io
df_state_params = pd.read_csv(io.StringIO(uploaded['AexpBx^1.csv'].decode('utf-8')))
df_state_params.head()
```

| | State | a (10^-3) | b | fit rank | r^2 |
|---|-------|-----------|-------------|----------|----------|
| 0 | AK | 2.593040 | -75.366476 | 1.0 | 0.996906 |
| 1 | AL | 12.121593 | -111.222242 | 2.0 | 0.997430 |
| 2 | AR | 2.941186 | -75.356785 | 4.0 | 0.997586 |
| 3 | AS | NaN | NaN | NaN | NaN |
| 4 | AZ | 4.984063 | -90.295019 | 1.0 | 0.998613 |

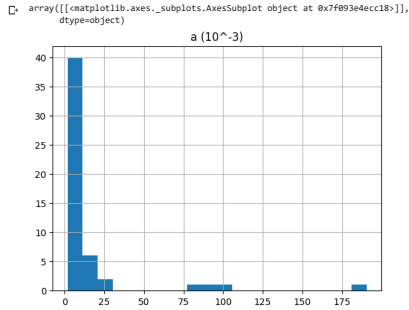
```
df_state_params.describe()
```

| | a (10^-3) | b | fit rank | r^2 |
|-------|------------|-------------|-----------|-----------|
| count | 52.000000 | 52.000000 | 52.000000 | 52.000000 |
| mean | 16.215254 | -100.951881 | 1.769231 | 0.995682 |
| std | 31.801661 | 25.545128 | 1.095720 | 0.004749 |
| min | 1.952592 | -185.986576 | 1.000000 | 0.972728 |
| 25% | 5.041013 | -116.155268 | 1.000000 | 0.995399 |
| 50% | 7.113788 | -99.476492 | 1.000000 | 0.997030 |
| 75% | 10.698133 | -80.847333 | 2.000000 | 0.998087 |
| max | 190.553218 | -49.104858 | 5.000000 | 0.999660 |

```
df_state_params.hist(column='r^2')
```

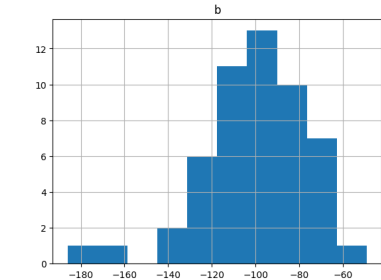
<Figure size 640x480 with 0 Axes>

```
array([[<matplotlib.axes._subplots.AxesSubplot object at 0x7f093e578940>]],
df_state_params.hist(column='a (10^-3)', bins=20)
```



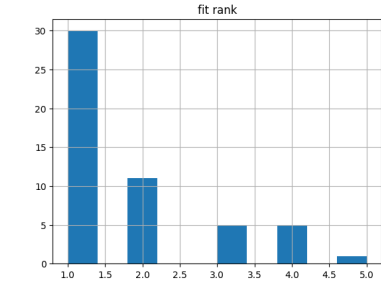
High value outliers here are NJ (fit rank 1), NY (fit rank 1), RI (fit rank 5), and SD (fit rank 4)

```
df_state_params.hist(column='b', bins=10)
array([[<matplotlib.axes._subplots.AxesSubplot object at 0x7f093e48b780>]],
dtype=object)
```



Low value outliers here are RI (fit rank 5) and SD (fit rank 4).

```
df_state_params.hist(column='fit rank')
array([[<matplotlib.axes._subplots.AxesSubplot object at 0x7f093e433ba8>]],
dtype=object)
```



The $A \cdot \exp(B/x)$ functional form works extremely well for thirty of the 52 states (57.7%).

```
# Fetch static data for each state (CovidCompleteStateData.csv)
from google.colab import files
uploaded = files.upload()

Choose Files CovidCompl_iteData.csv
CovidCompleteStateData.csv(application/vnd.ms-excel) - 82610 bytes, last modified: 4/20/2020 - 100% done
Saving CovidCompleteStateData.csv to CovidCompleteStateData.csv
```

```
# Load static data for each state (CovidCurrentStateData.csv)
import io
df_state_data = pd.read_csv(io.StringIO(uploaded['CovidCompleteStateData.csv'].decode('utf-8')))
```

| | State | Sum of NUM_Medicare_BEN | Sum of NUM_BEN_Age_Less_65 | Sum of NUM_BEN_Age_65_to_74 | Sum of NUM_BEN_Age_75_to_84 | Sum of NUM_BEN_Age_Greater_84 | Sum of NUM_Female_BEN | Sum of NUM_Male_BEN | Sum of NUM_Black_or_African_American_BEN | Sum of NUM_Asian_Pacific_Islander_BEN | Sum of NUM_Hispanic_BEN | Sum of NUM_American_IndianAlaska_Native_BEN | Sum of NUM_BEN_With |
|---------------------|-------|----------------------------|-------------------------------|--------------------------------|--------------------------------|----------------------------------|--------------------------|------------------------|---|--|----------------------------|--|------------------------|
| 0 | AK | 1820384.0 | 270970.0 | 809516.0 | 468255.0 | 175296.0 | 1034762.0 | 760009.0 | 62311.0 | 76773.0 | 46525.0 | 147917.0 | |
| 1 | AL | 10804823.0 | 2065353.0 | 4386595.0 | 2980828.0 | 1190504.0 | 6237445.0 | 4514041.0 | 1549811.0 | 30624.0 | 65500.0 | 5556.0 | |
| 2 | AR | 15892716.0 | 2818665.0 | 6370265.0 | 4555468.0 | 1848506.0 | 9275039.0 | 6507151.0 | 1334245.0 | 19642.0 | 108428.0 | 62782.0 | |
| 3 | AS | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| 4 | AZ | 10786064.0 | 886596.0 | 4861035.0 | 3377040.0 | 1294375.0 | 5944519.0 | 4747801.0 | 221183.0 | 61840.0 | 689880.0 | 179818.0 | |
| 5 rows x 16 columns | | | | | | | | | | | | | |

```
# Feature Engineering
# Land Area/Water Area
df_state_data['State Area Ratio'] = df_state_data['Land Area']/df_state_data['Water Area']
df_state_data['State Area Ratio'] = df_state_data['Land Area'].divide(df_state_data['Water Area'], fill_value=0)

# Elevation Ratio = Highest Elevation/Mean Elevation
df_state_data['Elevation Ratio'] = df_state_data['Highest Elevation']/df_state_data['Mean Elevation']
df_state_data['Elevation Ratio'] = df_state_data['Highest Elevation'].divide(df_state_data['Mean Elevation'], fill_v

# Capital Area Ratio = Capital Land Area/Capital Water Area
df_state_data['Capital Area Ratio'] = df_state_data['Capital Land Area']/df_state_data['Capital Water Area']
df_state_data['Capital Land Area'] = df_state_data['Capital Land Area'].astype(float)
df_state_data['Capital Area Ratio'] = df_state_data['Capital Land Area'].divide(df_state_data['Capital Water Area'],

# Boundaries = Number of boarding states + On Coast + Borders Another Country
df_state_data['Boundaries'] = df_state_data['Number of bordering states'] + df_state_data['On Coast'] + df_state_data

# Latitude Difference to State Capital = Latitude - Capital Latitude
df_state_data['Latitude Difference to State Capital'] = df_state_data['Latitude'] - df_state_data['Capital Latitude']

# Longitude Difference to State Capital = Capital Longitude - Longitude
df_state_data['Longitude Difference to State Capital'] = df_state_data['Capital Longitude'] - df_state_data['Longitude']

# Latitude Difference to DC = Latitude - DC Latitude
df_state_data['Latitude Difference to DC'] = df_state_data['Latitude'] - 38.904722

# Longitude Difference to DC = DC Longitude - Longitude
df_state_data['Longitude Difference to DC'] = -77.061389 - df_state_data['Longitude']

# Latitude Difference to US Center = Latitude - Center Latitude
df_state_data['Latitude Difference to Center'] = df_state_data['Latitude'] - 39.833333

# Longitude Different to US Center = Center Longitude - Longitude
```



```
df_state_data['Longitude Difference to Center'] = -98.585522 - df_state_data['Longitude']

df_state_data.head()
```

| | State | Sum of NUM_Medicare_BEN | Sum of NUM_BEN_Age_Less_65 | Sum of NUM_BEN_Age_65_to_74 | Sum of NUM_BEN_Age_75_to_84 | Sum of NUM_BEN_Age_Greater_84 | Sum of NUM_Female_BEN | Sum of NUM_Male_BEN | NUM_Black_or_African_American_BEN | Sum of NUM_Asian_Pacific_Islander_BEN | Sum of NUM_Hispanic_BEN | Sum of NUM_American_IndianAlaska_Native_BEN | Sum of NUM_BEN_With_ |
|---|-------|----------------------------|-------------------------------|--------------------------------|--------------------------------|----------------------------------|--------------------------|------------------------|-----------------------------------|--|----------------------------|--|-------------------------|
| 0 | AK | 1820384.0 | 270970.0 | 809516.0 | 468255.0 | 175296.0 | 1034762.0 | 760009.0 | | 62311.0 | 76773.0 | 46525.0 | 147917.0 |
| 1 | AL | 10804823.0 | 2065353.0 | 4386595.0 | 2980828.0 | 1190504.0 | 6237445.0 | 4514041.0 | 1549811.0 | 30624.0 | 65500.0 | | 5556.0 |
| 2 | AR | 15892716.0 | 2818665.0 | 6370265.0 | 4555468.0 | 1848506.0 | 9275039.0 | 6507151.0 | 1334245.0 | 19642.0 | 108428.0 | | 62782.0 |
| 3 | AS | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 4 | AZ | 10786064.0 | 886596.0 | 4861035.0 | 3377040.0 | 1294375.0 | 5944519.0 | 4747801.0 | | 221183.0 | 61840.0 | 689880.0 | 179818.0 |

5 rows × 126 columns

```
df_state_data.shape
```

```
(56, 126)
```

```
# Define variables for regression
df_temp1 = df_state_data.drop(df_state_data.index[[3, 12, 27, 42, 50]])
x = df_temp1.drop('State', axis = 1)
df_temp2 = df_state_params.drop(df_state_data.index[[3, 12, 27, 42, 50]])
y = df_temp2['a (10^-3)']

# Look at correlation coefficients
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', 1000)
x.corr()
```

| | Sum of NUM_Medicare_BEN | Sum of NUM_BEN_Age_Less_65 | Sum of NUM_BEN_Age_65_to_74 | Sum of NUM_BEN_Age_75_to_84 | Sum of NUM_BEN_Age_Greater_84 | Sum of NUM_Female_BEN | Sum of NUM_Male_BEN | Sum of NUM_Black_or_African_American_BEN | Sum of NUM_Asian_Pacific_Islander_BEN | Sum of NUM_Hispanic_BEN | Sum of NUM_Ameri |
|--|----------------------------|-------------------------------|--------------------------------|--------------------------------|----------------------------------|--------------------------|------------------------|---|--|----------------------------|---------------------|
| Sum of NUM_Medicare_BEN | 1.000000 | 0.961404 | 0.998624 | 0.998100 | 0.989961 | 0.999917 | 0.999897 | 0.896692 | 0.525530 | 0.893302 | |
| Sum of NUM_BEN_Age_Less_65 | 0.961404 | 1.000000 | 0.978099 | 0.969440 | 0.960650 | 0.982576 | 0.979741 | 0.926091 | 0.475021 | 0.822787 | |
| Sum of NUM_BEN_Age_65_to_74 | 0.998624 | 0.978099 | 1.000000 | 0.996374 | 0.982712 | 0.996372 | 0.986636 | 0.895722 | 0.517514 | 0.902298 | |
| Sum of NUM_BEN_Age_75_to_84 | 0.998100 | 0.969440 | 0.996374 | 1.000000 | 0.992601 | 0.997916 | 0.998296 | 0.884218 | 0.530001 | 0.899556 | |
| Sum of NUM_BEN_Age_Greater_84 | 0.989961 | 0.960650 | 0.982712 | 0.992601 | 1.000000 | 0.989606 | 0.990404 | 0.864777 | 0.561253 | 0.879739 | |
| Sum of NUM_Female_BEN | 0.999917 | 0.982576 | 0.998372 | 0.997916 | 0.989606 | 1.000000 | 0.999658 | 0.899227 | 0.523618 | 0.898962 | |
| Sum of NUM_Male_BEN | 0.999897 | 0.979741 | 0.998636 | 0.998296 | 0.990404 | 0.999658 | 1.000000 | 0.893490 | 0.527005 | 0.895526 | |
| Sum of NUM_Black_or_African_American_BEN | 0.896692 | 0.926091 | 0.895722 | 0.884218 | 0.864777 | 0.899227 | 0.893490 | 1.000000 | 0.302985 | 0.726543 | |
| Sum of NUM_Asian_Pacific_Islander_BEN | 0.525530 | 0.475021 | 0.517514 | 0.530001 | 0.561253 | 0.523618 | 0.527005 | 1.000000 | 0.000000 | 0.633875 | |
| Sum of NUM_Hispanic_BEN | 0.893302 | 0.822787 | 0.902298 | 0.899556 | 0.879739 | 0.889962 | 0.895526 | 0.726543 | 0.633875 | 1.000000 | |
| Sum of NUM_American_IndianAlaska_Native_BEN | 0.082561 | 0.059858 | 0.091513 | 0.086836 | 0.065730 | 0.082194 | 0.083230 | -0.035649 | 0.118158 | 0.130647 | |
| Sum of NUM_BEN_With_Race_Not_Elsewhere_Classified | 0.823477 | 0.774080 | 0.803783 | 0.832225 | 0.879096 | 0.821921 | 0.824850 | 0.638847 | 0.739901 | 0.734084 | |
| Sum of NUM_Non-Hispanic_White_BEN | 0.998636 | 0.978894 | 0.994391 | 0.996119 | 0.988883 | 0.997045 | 0.996745 | 0.888851 | 0.485602 | 0.865686 | |
| Sum of NUM_Minorities | 0.958442 | 0.925721 | 0.961095 | 0.957721 | 0.945115 | 0.957361 | 0.958590 | 0.868495 | 0.645767 | 0.959140 | |
| Sum of Average_Age_of_BEN | 0.682483 | 0.730359 | 0.686432 | 0.663590 | 0.637504 | 0.685316 | 0.678577 | 0.692434 | 0.132443 | 0.516973 | |
| Sum of NUM_BEN_Atrial_Fibrillation | 0.990425 | 0.969550 | 0.985604 | 0.991418 | 0.990376 | 0.990407 | 0.990816 | 0.889881 | 0.460369 | 0.851983 | |
| Sum of NUM_BEN_Asthma | 0.995532 | 0.979588 | 0.991583 | 0.992903 | 0.991762 | 0.995242 | 0.995600 | 0.893349 | 0.526980 | 0.880784 | |
| Sum of NUM_BEN_Cancer | 0.994765 | 0.972149 | 0.992903 | 0.994874 | 0.987401 | 0.994443 | 0.994921 | 0.900538 | 0.464158 | 0.883930 | |
| Sum of NUM_BEN_Heart_Failure | 0.997133 | 0.985150 | 0.995371 | 0.993915 | 0.984952 | 0.997171 | 0.996846 | 0.913138 | 0.484598 | 0.884681 | |
| Sum of NUM_BEN_Chronic_Kidney_Disease | 0.997501 | 0.980301 | 0.997095 | 0.995430 | 0.984274 | 0.997279 | 0.997616 | 0.907084 | 0.485414 | 0.893142 | |
| Sum of NUM_BEN_Chronic_Obstructive_Pulmonary_Disease | 0.986234 | 0.980624 | 0.981625 | 0.983999 | 0.978052 | 0.986989 | 0.985885 | 0.906585 | 0.429822 | 0.833395 | |
| Sum of NUM_BEN_Hyperlipidemia | 0.996237 | 0.974348 | 0.994742 | 0.996423 | 0.987588 | 0.996101 | 0.996491 | 0.903165 | 0.477347 | 0.884215 | |
| Sum of NUM_BEN_Diabetes | 0.997754 | 0.981227 | 0.996544 | 0.995687 | 0.985896 | 0.997745 | 0.997458 | 0.912776 | 0.494767 | 0.892805 | |
| Sum of NUM_BEN_Hypertension | 0.998856 | 0.982300 | 0.998079 | 0.996943 | 0.986018 | 0.998963 | 0.998633 | 0.908147 | 0.492704 | 0.886804 | |
| Sum of NUM_BEN_Ischemic_Heart_Disease | 0.994006 | 0.975145 | 0.991547 | 0.994105 | 0.985840 | 0.994115 | 0.993975 | 0.906294 | 0.457797 | 0.876921 | |
| Sum of NUM_BEN_Stroke | 0.990547 | 0.972081 | 0.988818 | 0.990024 | 0.980879 | 0.990462 | 0.990642 | 0.919103 | 0.447938 | 0.879610 | |
| Sum of PCT_MEDICARE | 0.713702 | 0.762102 | 0.716971 | 0.696228 | 0.671536 | 0.717742 | 0.709263 | 0.752793 | 0.141713 | 0.484719 | |
| % Urban Pop | 0.246412 | 0.181090 | 0.240984 | 0.259055 | 0.285836 | 0.242294 | 0.250798 | 0.181632 | 0.311802 | 0.283145 | |
| Density (Pmi2) | -0.095479 | -0.105571 | -0.096280 | -0.092015 | -0.087655 | -0.096111 | -0.095250 | -0.017993 | -0.029204 | -0.041712 | |
| Children 0-18 | 0.886252 | 0.846604 | 0.876226 | 0.888481 | 0.912717 | 0.884787 | 0.887362 | 0.723346 | 0.775958 | 0.840004 | |
| Adults 19-25 | 0.865749 | 0.826231 | 0.852680 | 0.868860 | 0.900273 | 0.864399 | 0.866800 | 0.698451 | 0.784342 | 0.808969 | |
| Adults 26-34 | 0.844661 | 0.804462 | 0.835397 | 0.852982 | 0.888196 | 0.847034 | 0.850081 | 0.667717 | 0.811427 | 0.807784 | |
| Adults 35-54 | 0.861684 | 0.820010 | 0.848035 | 0.865769 | 0.898830 | 0.860256 | 0.862945 | 0.695994 | 0.775924 | 0.803195 | |
| Adults 55-64 | 0.840536 | 0.802214 | 0.822003 | 0.845654 | 0.889137 | 0.839342 | 0.841899 | 0.678395 | 0.734919 | 0.747404 | |
| 65+ | 0.842520 | 0.796154 | 0.822919 | 0.852028 | 0.886626 | 0.841354 | 0.844588 | 0.672577 | 0.691686 | 0.734185 | |
| Latitude | -0.400391 | -0.397373 | -0.403138 | -0.407192 | -0.381151 | -0.404042 | -0.398907 | -0.449230 | -0.184344 | -0.284522 | |
| Longitude | 0.026061 | 0.029274 | 0.034115 | 0.040031 | 0.057825 | 0.040875 | 0.043076 | 0.189118 | -0.272166 | -0.097547 | |
| Land Area | 0.229013 | 0.193883 | 0.242084 | 0.230058 | 0.205886 | 0.226251 | 0.230971 | 0.128422 | 0.201398 | 0.342072 | |
| Water Area | 0.042895 | 0.056385 | 0.036723 | 0.038782 | 0.050598 | 0.042948 | 0.042530 | 0.080106 | 0.048930 | 0.044540 | |
| Mean Elevation | -0.163276 | -0.224740 | -0.147730 | -0.155029 | -0.169641 | -0.168622 | -0.157738 | -0.314621 | 0.096312 | 0.038452 | |
| Highest Elevation | -0.050881 | -0.137582 | -0.040603 | -0.049835 | -0.071530 | -0.065559 | -0.054634 | -0.233143 | 0.289714 | 0.156542 | |
| Lowest elevation | -0.354394 | -0.352655 | -0.344053 | -0.355481 | -0.373528 | -0.356287 | -0.352697 | -0.312823 | -0.524776 | -0.283582 | |
| Number of bordering states | 0.077790 | 0.135863 | 0.075964 | 0.059448 | 0.056612 | 0.079630 | 0.074603 | 0.044175 | -0.147225 | -0.075000 | |
| On Coast | 0.471024 | 0.505115 | 0.442960 | 0.461862 | 0.518306 | 0.471588 | 0.470022 | 0.512059 | 0.172256 | 0.274148 | |
| Borders Another Country | 0.358143 | 0.310618 | 0.363869 | 0.356815 | 0.359742 | 0.351935 | 0.363168 | 0.188403 | 0.423258 | 0.500772 | |
| Capital Latitude | -0.388663 | -0.393079 | -0.394070 | -0.392266 | -0.395029 | -0.392545 | -0.386636 | -0.463654 | -0.136997 | -0.269834 | |
| Capital Longitude | 0.027375 | 0.076949 | 0.015075 | 0.019608 | 0.037346 | 0.030778 | 0.023488 | 0.181245 | -0.297512 | -0.116778 | |
| Capital Land Area | 0.008902 | -0.002410 | 0.018688 | 0.009403 | -0.012432 | 0.008939 | 0.008880 | -0.020032 | -0.012972 | 0.022939 | |
| Capital Water Area | -0.087670 | -0.096352 | -0.083610 | -0.087193 | -0.092577 | -0.088397 | -0.087364 | -0.097128 | -0.020784 | -0.040196 | |
| Capital Mean Elevation | -0.194009 | -0.217624 | -0.182931 | -0.190725 | -0.206936 | -0.197080 | -0.190497 | -0.249034 | -0.121862 | -0.049687 | |
| Capital is the Largest City | -0.171080 | -0.147972 | -0.165860 | -0.173283 | -0.194417 | -0.169162 | -0.172323 | -0.133938 | -0.129514 | -0.188909 | |
| Largest City Latitude | -0.421170 | -0.421496 | -0.425109 | -0.424938 | -0.398170 | -0.424688 | -0.419431 | -0.467601 | -0.234938 | -0.317480 | |
| Largest City Longitude | 0.057094 | 0.102104 | 0.044423 | 0.050338 | 0.069739 | 0.060248 | 0.053604 | 0.201824 | -0.262204 | -0.081857 | |
| Number of Counties | 0.663716 | 0.710105 | 0.670375 | 0.654677 | 0.611985 | 0.666592 | 0.659547 | 0.684930 | 0.100824 | 0.503011 | |
| Became a State | -0.140415 | -0.200801 | -0.128869 | -0.126557 | -0.143313 | -0.144064 | -0.136075 | -0.308218 | 0.076276 | 0.035901 | |
| DaysSinceStayatHomeOrder | -0.020651 | -0.019693 | -0.030343 | -0.027347 | 0.007693 | -0.023631 | -0.018913 | -0.045827 | 0.221997 | 0.052471 | |
| DaysSinceFirstPositive | 0.368252 | 0.311941 | 0.366229 | 0.374653 | 0.390899 | 0.365580 | 0.379093 | 0.287102 | 0.259030 | 0.302889 | |
| DaysSinceTestStart | 0.290649 | 0.242592 | 0.289428 | 0.297948 | 0.312252 | 0.288578 | 0.292640 | 0.233275 | 0.191547 | 0.240467 | |
| 15-49yearsAlicauses | 0.888203 | 0.856564 | 0.874919 | 0.889982 | 0.919415 | 0.887334 | 0.889058 | 0.739853 | 0.736873 | 0.795191 | |
| 15-49yearsAsthma | 0.824682 | 0.787879 | 0.807656 | 0.860382 | 0.827220 | 0.825303 | 0.825395 | 0.667614 | 0.757078 | 0.793058 | |
| 15-49yearsChronickidneydisease | 0.918864 | 0.893772 | 0.909568 | 0.918825 | 0.935419 | 0.918655 | 0.918932 | 0.805947 | 0.715267 | 0.828201 | |
| 15-49yearsChronicobstructivepulmonarydisease | 0.896769 | 0.878089 | 0.880516 | 0.897303 | 0.928648 | 0.896643 | 0.896964 | 0.771591 | 0.635771 | 0.744901 | |
| 15-49yearsDiabetesmellitus | 0.912330 | 0.881654 | 0.900896 | 0.914260 | 0.936863 | 0.911726 | 0.912809 | 0.781996 | 0.692768 | 0.812172 | |
| 15- 49yearsinterstitiallungdiseaseandpulmonarysarcoidosis | 0.881251 | 0.864222 | 0.866766 | 0.880121 | 0.909006 | 0.881289 | 0.881043 | 0.782863 | 0.644183 | 0.738019 | |
| 15-49yearsischemicheartdisease | 0.928387 | 0.927789 | 0.916634 | 0.923405 | 0.939571 | 0.929317 | 0.927199 | 0.849388 | 0.595981 | 0.764879 | |
| 15-49yearsNeoplasms | 0.887461 | 0.860138 | 0.873067 | 0.886885 | 0.919369 | 0.886922 | 0.887968 | 0.729597 | 0.729552 | 0.858252 | |
| 15-49yearsOtherchronicrespiratorydiseases | 0.906636 | 0.885246 | 0.892415 | 0.906637 | 0.934764 | 0.906287 | 0.906852 | 0.785146 | 0.652797 | 0.775052 | |
| 15-49yearsRheumaticheartdisease | 0.903473 | 0.893269 | 0.893364 | 0.898792 | 0.916822 | 0.903373 | 0.903469 | 0.792144 | 0.690715 | 0.785070 | |
| 15-49yearsStroke | 0.919789 | 0.910295 | 0.919449 | 0.934870 | 0.919891 | 0.919891 | 0.919744 | 0.808343 | 0.702539 | 0.815118 | |
| 50-69yearsAlicauses | 0.880146 | 0.855617 | 0.863069 | 0.881923 | 0.918175 | 0.879688 | 0.880993 | 0.745024 | 0.677823 | 0.745280 | |
| 50-69yearsAsthma | 0.801803 | 0.765502 | 0.781306 | 0.805925 | 0.855811 | 0.800503 | 0.802845 | 0.641040 | 0.741394 | 0.706390 | |
| 50-69yearsChronickidneydisease | 0.917312 | 0.898401 | 0.905572 | 0.916416 | 0.938247 | 0.917261 | 0.917589 | 0.809479 | 0.675662 | 0.793827 | |
| 50-69yearsChronicobstructivepulmonarydisease | 0.879259 | 0.872843 | 0.860771 | 0.878641 | 0.912096 | 0.879671 | 0.879770 | 0.765135 | 0.542692 | 0.678077 | |
| 50-69yearsDiabetesmellitus | 0.882501 | 0.857522 | 0.865414 | 0.884673 | 0.920496 | 0.882047 | 0.883347 | 0.753933 | 0.653401 | 0.743092 | |
| 50- 69yearsinterstitiallungdiseaseandpulmonarysarcoidosis | 0.863191 | 0.840683 | 0.846169 | 0.863950 | 0.901043 | 0.862755 | 0.863919 | 0.739130 | 0.673635 | 0.725850 | |
| 50-69yearsischemicheartdisease | 0.905979 | 0.901073 | 0.890019 | 0.902683 | 0.931491 | 0.906505 | 0.905613 | 0.807307 | 0.618281 | 0.735980 | |
| 50-69yearsNeoplasms | 0.872500 | 0.853408 | 0.854035 | 0.873415 | 0.912024 | 0.872273 | 0.873220 | 0.746019 | 0.650724 | 0.719312 | |
| 50-69yearsOtherchronicrespiratorydiseases | 0.885021 | 0.875159 | 0.867604 | 0.883457 | 0.917418 | 0.885065 | 0.885343 | 0.780290 | 0.570425 | 0.701549 | |
| 50-69yearsRheumaticheartdisease | 0.892519 | 0.890373 | 0.880528 | 0.886667 | 0.908250 | 0.892709 | 0.892705 | 0.793842 | 0.640749 | 0.737680 | |
| 50-69yearsStroke | 0.907993 | 0.892290 | 0.895108 | 0.907388 | 0.930491 | 0.906246 | 0.908324 | 0.800720 | 0.65 | | |

AllAgesAllcauses0.8800030.8512930.8633990.8825920.9186340.8793560.8809770.7365530.6951840.757560

AllAgesAsthma0.8332530.7949170.8158120.8369100.8793100.8318110.8341910.6740540.7486970.754660

AllAgesChroniclkdneidisease0.9054620.8855030.8915040.9053070.9330930.9053110.9059430.7889480.6719410.774120

AllAgesChronicobstructivepulmonarydisease0.8772140.8608330.8581270.8792650.9160740.8772400.8781790.7459660.5806260.696350

AllAgesDiabetesmellitus0.8797280.8521210.8622090.8829080.9198540.8792320.8806700.7459440.6565130.743480

AllAgesInterstitiailungdiseaseandpulmonarysarcoidosis0.8539120.8260930.8357490.8567590.8972560.8533260.8550500.7145950.6811690.720490

AllAgesIschemicheartdisease0.8835350.8709540.8644600.8830690.9209410.8837200.8835840.7683830.6285180.717687

AllAgesNeoplasms0.8653250.8414500.8463250.8677260.9084290.8649320.8662730.7289200.6629080.719608

AllAgesOtherchronicrespiratorydiseases0.9035920.8859670.8884450.9029700.9329940.9033000.9040110.7856640.6228840.752887

AllAgesRheumaticheartdisease0.8803570.8752860.8661440.8750890.9040800.8801680.8808830.7673980.6475270.727462

AllAgesStroke0.8953980.8757350.8806740.8959780.9251200.8953810.8960740.7711830.6667430.752509

AllAgesTotal0.8805070.8539230.8634630.8828130.9192940.8799780.8814110.7406570.6836640.750268

Alpirtion0.8892290.8884420.8750920.8829640.9106520.8896240.8889910.7801680.6547750.729374

Highbody-massindex0.8937970.8727390.8771330.8946240.9287150.8935240.8943320.7712580.6608130.753636

Highfastingplasmaglucose0.8867950.8701240.8689090.8874170.9229850.8867790.8872610.7731980.6170420.724166

HighLDLcholesterol0.8932150.8824830.8753980.8920230.9263040.8935270.8931240.7833540.6272470.727214

Highsystolicbloodpressure0.8974530.8826310.8803460.8971310.9304640.8975470.8976700.7882830.6385190.741135

Impairedkidneyfunction0.8899340.8726930.8731730.8900340.9238350.8899740.8903210.7731150.6577980.740044

Noaccessstohandwashingfacility0.8776030.8577810.8624530.8765190.9004990.8771120.8779400.7566860.6678410.744063

Smoking0.8815790.8667260.8628310.8826120.9190200.8817620.8820990.7609420.6040360.705752

Log10Pop0.7284940.7379020.7140570.7223200.7471290.7295270.7278710.6665720.4215210.509917

DaysSinceInfection0.4225250.3730100.4197270.4312330.4439990.4207750.4244020.3653160.2630220.351935

Children0-180.1671330.1808230.1812960.1695800.1190980.1190980.1641190.1696120.0498110.148864

Allriskfactors0.8828150.8600440.8655300.8842170.9201340.8824870.8836530.7505940.6605520.737247

State Area Ratio-0.141342-0.190449-0.126323-0.134563-0.158168-0.143787-0.137979-0.260308-0.103905-0.027420

Elevation Ratio0.0203320.0073110.0295980.023691-0.0002230.0217150.0184820.1298900.0529350.012281

Capital Area Ratio-0.119284-0.151665-0.109968-0.112407-0.124967-0.120614-0.117747-0.171038-0.076449-0.052022

Boundaries0.4993560.5563930.4793300.4779600.5171250.4994870.4957220.4552720.1254650.277110

Latitude Difference to State Capital-0.268652-0.211068-0.252026-0.293417-0.327548-0.268822-0.271605-0.115573-0.406313-0.233715

Longitude Difference to State Capital-0.143646-0.133106-0.139285-0.150250-0.155036-0.143567-0.145263-0.100819-0.106975-0.108075

Latitude Difference to DC-0.400391-0.397373-0.403138-0.407192-0.381151-0.404042-0.398907-0.449230-0.184344-0.284522

Longitude Difference to DC-0.040601-0.092974-0.034115-0.040031-0.057825-0.049875-0.043076-0.1891180.2721660.097547

Latitude Difference to Center-0.400391-0.397373-0.403138-0.407192-0.381151-0.404042-0.398907-0.449230-0.184344-0.284522

Longitude Difference to Center-0.040601-0.092974-0.034115-0.040031-0.057825-0.049875-0.043076-0.1891180.2721660.097547

Note that there are many highly correlated features which need to be dropped

Create absolute value correlation matrix

corr_matrix = X.corr().abs()

Select upper triangle of correlation matrix

upper = corr_matrix.where(np.triu(np.ones(corr_matrix.shape), k=1).astype(np.bool))

Find index of feature columns with correlation greater than 0.95

to_drop = [column for column in upper.columns if any(upper[column] > 0.95)]

Drop features by index which were identified as being highly correlated

X = X.drop(X[to_drop], axis=1)

X.head()

| | Sum of NUM_Medicare_BEN | Sum of NUM_Black_or_African_American_BEN | Sum of NUM_Asian_Pacific_Islander_BEN | Sum of NUM_Hispanic_BEN | Sum of NUM_American_IndianAlaska_Native_BEN | Sum of NUM_BEN_With_Race_Not_Elsewhere_Classified | Sum of Average_Age_of_BEN | Sum of PCT_MEDICARE | % Urban Pop | Density (P/mi2) | Children 0-18 | Latitude | Longitude | Land Area |
|---|----------------------------|---|--|----------------------------|--|--|------------------------------|------------------------|-------------------|--------------------|------------------|-----------|-------------|--------------|
| 0 | 1820384.0 | 62311.0 | 76773.0 | 46525.0 | 147917.0 | 23372.0 | 996.298679 | 10.069041 | 66.0 | 1.2863 | 181405.17 | 61.370716 | -152.404419 | 570665.0 |
| 1 | 10804823.0 | 1549811.0 | 30624.0 | 65500.0 | 5556.0 | 58660.0 | 3967.220634 | 51.254704 | 59.0 | 96.9221 | 1105570.08 | 32.806671 | -86.791130 | 50644.0 |
| 2 | 15892716.0 | 1334245.0 | 19642.0 | 108428.0 | 62782.0 | 61250.0 | 3928.834167 | 94.570949 | 56.2 | 58.4030 | 686482.50 | 34.969704 | -92.373123 | 52030.0 |
| 4 | 10780604.0 | 221183.0 | 61840.0 | 689880.0 | 179818.0 | 114903.0 | 1009.367955 | 14.075942 | 89.8 | 64.9550 | 1744612.56 | 33.729759 | -111.431221 | 113595.0 |
| 5 | 42579588.0 | 2072012.0 | 3276415.0 | 5674776.0 | 113871.0 | 562214.0 | 4001.853612 | 63.398334 | 95.0 | 256.3727 | 9481941.36 | 36.116203 | -119.681564 | 155766.0 |

X.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 51 entries, 0 to 55
Data columns (total 38 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   Sum of NUM_Medicare_BEN                   51 non-null     float64
1   Sum of NUM_Black_or_African_American_BEN 51 non-null     float64
2   Sum of NUM_Asian_Pacific_Islander_BEN     51 non-null     float64
3   Sum of NUM_Hispanic_BEN                   51 non-null     float64
4   Sum of NUM_American_IndianAlaska_Native_BEN 51 non-null     float64
5   Sum of NUM_BEN_With_Race_Not_Elsewhere_Classified 51 non-null     float64
6   Sum of Average_Age_of_BEN                 51 non-null     float64
7   Sum of PCT_MEDICARE                       51 non-null     float64
8   % Urban Pop                               51 non-null     float64
9   Density (P/mi2)                           51 non-null     float64
10  Children 0-18                             51 non-null     float64
11  Latitude                                   51 non-null     float64
12  Longitude                                  51 non-null     float64
13  Land Area                                 51 non-null     float64
14  Water Area                               51 non-null     float64
15  Mean Elevation                           51 non-null     float64
16  Highest Elevation                        51 non-null     float64
17  Lowest elevation                         51 non-null     float64
18  Number of bordering states               51 non-null     float64
19  On Coast                                51 non-null     float64
20  Borders Another Country                  51 non-null     float64
21  Capital Land Area                        51 non-null     float64
22  Capital Water Area                       51 non-null     float64
23  Capital Mean Elevation                   51 non-null     float64
24  Capital is the Largest City               51 non-null     float64
25  Became a State                           51 non-null     float64
26  DaysSinceStayatHomeOrder                 51 non-null     float64
27  DaysSinceFirstPositive                    51 non-null     float64
28  DaysSinceTestStart                       51 non-null     float64
29  Log10Pop                                 51 non-null     float64
30  DaysSinceInfection                       51 non-null     float64
31  Children0-18                             51 non-null     float64
32  State Area Ratio                         51 non-null     float64
33  Elevation Ratio                         51 non-null     float64
34  Capital Area Ratio                       51 non-null     float64
35  Boundaries                              51 non-null     float64
36  Latitude Difference to State Capital      51 non-null     float64
37  Longitude Difference to State Capital     51 non-null     float64
dtypes: float64(38)
memory usage: 15.3 KB
```

X.describe()

| | Sum of NUM_Medicare_BEN | Sum of NUM_Black_or_African_American_BEN | Sum of NUM_Asian_Pacific_Islander_BEN | Sum of NUM_Hispanic_BEN | Sum of NUM_American_IndianAlaska_Native_BEN | Sum of NUM_BEN_With_Race_Not_Elsewhere_Classified | Sum of Average_Age_of_BEN | Sum of PCT_MEDICARE | % Urban Pop | Density (P/mi2) | Children 0- 18 | Latitude | Longit |
|-------|----------------------------|---|--|----------------------------|--|--|------------------------------|------------------------|----------------|--------------------|-------------------|-----------|----------|
| count | 5.100000e+01 | 5.100000e+01 | 5.100000e+01 | 5.100000e+01 | 5.100000 | 5.100000 | 5.100000 | 5.100000 | 5.100000 | 5.100000 | 5.100000e+01 | 5.100000 | 5.1000 |
| mean | 1.038431e+07 | 9.464777e+05 | 1.411691e+05 | 5.310095e+05 | 38355.372549 | 86379.019608 | 3535.191553 | 52.223011 | 74.107843 | 431.560508 | 1.486858e+06 | 39.464823 | -93.339 |
| std | 1.311026e+07 | 4.722330e+06 | 4.722330e+06 | 1.629961e+06 | 87337.002647 | 114765.665446 | 2518.178494 | 46.520870 | 14.885481 | 1647.225920 | 1.764935e+06 | 6.069546 | 19.288 |
| min | 1.655870e+05 | 2.960000e+02 | 1.660000e+02 | 4.130000e+02 | 0.000000 | 1693.000000 | 70.002893 | 0.972106 | 38.700000 | 1.286300 | 1.160123e+05 | 21.094318 | -157.498 |
| 25% | 2.252305e+06 | 5.366000e+04 | 6.445500e+03 | 3.101950e+04 | 2980.500000 | 17674.500000 | 1542.140834 | 13.385073 | 65.400000 | 50.604850 | 3.984744e+05 | 35.688955 | -102.547 |
| 50% | 6.272609e+06 | 3.156040e+05 | 2.579200e+04 | 1.042170e+05 | 7061.000000 | 3578.360041 | 51.254704 | 74.200000 | 108.049700 | 1.013513e+06 | 39.849426 | -89.616 | |
| 75% | 1.471830e+07 | 1.547566e+06 | 7.063400e+04 | 2.005865e+05 | 28506.500000 | 100449.000000 | 52.14099778 | 78.017975 | 87.550000 | 223.983100 | 1.719581e+06 | 43.041292 | -78.988 |
| max | 7.644909e+07 | 7.011107e+06 | 3.276415e+06 | 1.007620e+07 | 560433.000000 | 562214.000000 | 13644.965980 | 219.756971 | 100.000000 | 11814.541000 | 9.481941e+06 | 61.370716 | -69.3811 |

https://colab.research.google.com/drive/1FNfgMzUDDLcyH4QLNmKClyT7CFQloSRi#scrollTo=j2pO3dNu6Edz&printMode=true

11/18

```
Double-click (or enter) to edit

# Train/validate split: random 75/25% train/validate split.
from sklearn.model_selection import train_test_split

X_train, X_val, y_train, y_val = train_test_split(X, y, test_size = 0.25, random_state = 42)

X_train.shape, y_train.shape, X_val.shape, y_val.shape

[(38, 38), (38,), (13, 38), (13,)]

X_train.describe()


```

| | Sum of NUM_Medicare_BEN | Sum of NUM_Black_or_African_American_BEN | Sum of NUM_Asian_Pacific_Islander_BEN | Sum of NUM_Hispanic_BEN | Sum of NUM_American_IndianAlaska_Native_BEN | Sum of NUM_BEN_With_Race_Not_Elsewhere_Classified | Sum of Average_Age_of_BEN | Sum of PCT_MEDICARE | % Urban Pop | Density (P/mi2) | Children 0- 18 | Latitude | Longiti |
|-------|----------------------------|---|--|----------------------------|--|--|------------------------------|------------------------|----------------|--------------------|-------------------|--------------|--------------------|
| count | 3.800000e+01 | 3.800000e+01 | 3.800000e+01 | 3.800000e+01 | | 38.000000 | 38.000000 | 38.000000 | 38.000000 | 38.000000 | 3.800000e+01 | 38.000000 | 38.0000 |
| mean | 1.014125e+07 | 9.685705e+05 | 1.623107e+05 | 3.942231e+05 | 37676.868421 | | 93816.710526 | 3549.919212 | 51.835362 | 75.578947 | 509.044824 | 1.536694e+06 | 39.769502 -94.847 |
| std | 9.963253e+06 | 1.001560e+06 | 5.333709e+05 | 1.021129e+06 | 93507.556838 | | 120712.415819 | 2056.853166 | 40.781616 | 14.338667 | 1899.378670 | 1.707713e+06 | 6.596259 20.534 |
| min | 3.472690e+05 | 2.689000e+03 | 4.580000e+02 | 2.622000e+03 | 45.000000 | | 4201.000000 | 70.002893 | 0.972106 | 38.700000 | 1.286300 | 1.365427e+05 | 21.094318 -157.498 |
| 25% | 2.518838e+06 | 4.934350e+04 | 1.427175e+04 | 3.676725e+04 | 4699.500000 | | 19802.500000 | 1590.459797 | 18.987497 | 66.125000 | 43.927350 | 5.296779e+05 | 35.839934 -103.929 |
| 50% | 7.473651e+06 | 5.120990e+05 | 3.068000e+04 | 1.071920e+05 | 9615.500000 | | 62351.000000 | 3948.027401 | 53.017662 | 74.850000 | 112.688450 | 1.123149e+06 | 40.249745 -89.647 |
| 75% | 1.563758e+07 | 1.560497e+06 | 9.455175e+04 | 1.983508e+05 | 27902.250000 | | 100987.500000 | 5194.290279 | 79.721735 | 87.970000 | 249.285600 | 1.728603e+06 | 43.183955 -80.276 |
| max | 4.257959e+07 | 3.265865e+06 | 3.276415e+06 | 5.674776e+06 | 560433.000000 | | 562214.000000 | 7981.626181 | 185.232318 | 100.000000 | 11814.541000 | 9.481941e+06 | 61.370716 -69.381 |

```


# Optimizing Hyperparameters for Random Forest Regressor
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestRegressor

# Define classifier
forest = RandomForestRegressor(random_state = 1)

# Parameters to fit
max_depth = [2, 3, 4]
n_estimators = [15, 36, 37]
min_samples_split = [1.5, 2, 2.5]
min_samples_leaf = [3.5, 4, 4.5]
max_leaf_nodes = [None]
max_features = ['auto']
ccp_alpha = [0.0, 0.00025, 0.0125]
min_weight_fraction_leaf = [0.0, 0.00025, 0.0125]

hyperF = dict(n_estimators = n_estimators, max_depth = max_depth,
              min_samples_split = min_samples_split,
              min_samples_leaf = min_samples_leaf,
              max_leaf_nodes = max_leaf_nodes,
              max_features = max_features,
              ccp_alpha=ccp_alpha,
              min_weight_fraction_leaf=min_weight_fraction_leaf)

gridF = GridSearchCV(forest, hyperF, cv = 3, verbose = 10,
                    scoring='r2', return_train_score=True,
                    n_jobs = -1)
bestF = gridF.fit(X_train, y_train)

# Output best accuracy and best parameters
print('The score achieved with the best parameters = ', gridF.best_score_, '\n')
print('The parameters are:', gridF.best_params_)


```

Fitting 3 folds for each of 729 candidates, totalling 2187 fits

Parallel(n_jobs=-1):

Using backend LokyBackend with 2 concurrent workers.

Parallel(n_jobs=-1):

Done 1 tasks | elapsed: 1.3s

Parallel(n_jobs=-1):

Done 4 tasks | elapsed: 1.4s

Parallel(n_jobs=-1):

Done 9 tasks | elapsed: 1.6s

Parallel(n_jobs=-1):

Done 14 tasks | elapsed: 1.7s

Parallel(n_jobs=-1):

Batch computation too fast (0.1862s.) Setting batch_size=2.

Parallel(n_jobs=-1):

Done 22 tasks | elapsed: 1.8s

Parallel(n_jobs=-1):

Batch computation too fast (0.1818s.) Setting batch_size=4.

Parallel(n_jobs=-1):

Batch computation too fast (0.1858s.) Setting batch_size=8.

Parallel(n_jobs=-1):

Done 44 tasks | elapsed: 2.2s

Parallel(n_jobs=-1):

Done 116 tasks | elapsed: 3.9s

Parallel(n_jobs=-1):

Done 188 tasks | elapsed: 5.3s

Parallel(n_jobs=-1):

Done 276 tasks | elapsed: 6.7s

Parallel(n_jobs=-1):

Done 364 tasks | elapsed: 8.5s

Parallel(n_jobs=-1):

Done 468 tasks | elapsed: 10.6s

Parallel(n_jobs=-1):

Done 572 tasks | elapsed: 12.2s

Parallel(n_jobs=-1):

Done 692 tasks | elapsed: 14.8s

Parallel(n_jobs=-1):

Done 812 tasks | elapsed: 16.7s

Parallel(n_jobs=-1):

Done 948 tasks | elapsed: 19.6s

Parallel(n_jobs=-1):

Done 1084 tasks | elapsed: 22.1s

Parallel(n_jobs=-1):

Done 1236 tasks | elapsed: 25.0s

Parallel(n_jobs=-1):

Done 1388 tasks | elapsed: 28.3s

Parallel(n_jobs=-1):

Done 1556 tasks | elapsed: 31.0s

Parallel(n_jobs=-1):

Done 1724 tasks | elapsed: 34.5s

Parallel(n_jobs=-1):

Done 1908 tasks | elapsed: 38.2s

Parallel(n_jobs=-1):

Done 2092 tasks | elapsed: 41.9s

The score achieved with the best parameters = -7.248154911869659

The parameters are: ('ccp_alpha': 0.0, 'max_depth': 3, 'max_features': 'auto', 'max_leaf_nodes': None, 'min_samples_leaf': 4, 'min_samples_split': 2, 'min_weight_fraction_leaf': 0.0, 'n_estimators': 36)

Parallel(n_jobs=-1):

Done 2187 out of 2187 | elapsed: 43.3s finished

```


!pip install category_encoders==2.0.0

Collecting category_encoders==2.0.0
  Downloading https://files.gyrfurhustad.org/packages/6e/a1/f7a22f1d4f31be78afeb00fa78478e028a64263a3c09b5ef54e073861e/category_encoders-2.0.0-ry2-none-any.whl (87kB)
    100% |#####| 920k 2.409/s
Requirement already satisfied: patsy>=0.4.1 in /usr/local/lib/python3.6/dist-packages (from category_encoders==2.0.0) (0.5.1)
Requirement already satisfied: pandas>=0.21.1 in /usr/local/lib/python3.6/dist-packages (from category_encoders==2.0.0) (1.0.3)
Requirement already satisfied: scikit-learn>=0.20.0 in /usr/local/lib/python3.6/dist-packages (from category_encoders==2.0.0) (0.22.2.post1)
Requirement already satisfied: statsmodels>=0.6.1 in /usr/local/lib/python3.6/dist-packages (from category_encoders==2.0.0) (0.10.2)
Requirement already satisfied: scipy>=0.19.0 in /usr/local/lib/python3.6/dist-packages (from category_encoders==2.0.0) (1.4.1)
Requirement already satisfied: numpy>=1.11.3 in /usr/local/lib/python3.6/dist-packages (from category_encoders==2.0.0) (1.18.3)
Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages (from patsy>=0.4.1->category_encoders==2.0.0) (1.12.0)
Requirement already satisfied: python-dateutil>=2.6.1 in /usr/local/lib/python3.6/dist-packages (from pandas>=0.21.1->category_encoders==2.0.0) (2.8.1)
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.6/dist-packages (from pandas>=0.21.1->category_encoders==2.0.0) (2018.3)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.6/dist-packages (from scikit-learn>=0.20.0->category_encoders==2.0.0) (0.14.1)
Installing collected packages: category-encoders
Successfully installed category-encoders-2.0.0

from sklearn.ensemble import RandomForestRegressor
from sklearn.pipeline import make_pipeline
import category_encoders as ce
from sklearn.impute import SimpleImputer

pipeline = make_pipeline(
    ce.OneHotEncoder(use_cat_names=True),
    SimpleImputer(strategy='mean'),
    RandomForestRegressor(bootstrap=True, ccp_alpha=0.0, criterion='mse',
                        max_depth=3, max_features='auto', max_leaf_nodes=None,
                        max_samples=None, min_impurity_decrease=0.0,
                        min_impurity_split=None, min_samples_leaf=4,
                        min_samples_split=2, min_weight_fraction_leaf=0.0,
                        n_estimators=36, n_jobs=None, oob_score=False,
                        random_state=0, verbose=0, warm_start=False))

pipeline.fit(X_train, y_train)

# Get the model's training accuracy
print("Training Accuracy: R^2 = ", pipeline.score(X_train, y_train))

# Get the model's validation accuracy
print("Validation Accuracy: R^2 = ", pipeline.score(X_val, y_val))

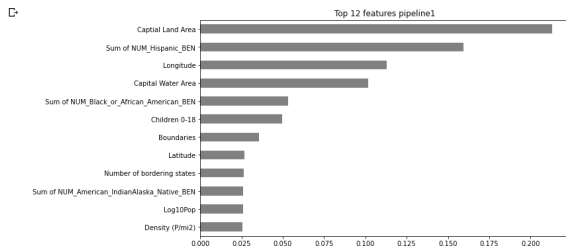
Training Accuracy: R^2 = 0.3013054634238487
Validation Accuracy: R^2 = 0.06938628867684

print('Feature Importances =')
#print(RandomForestRegressor.feature_importances_)
print(pipeline.steps[2][1].feature_importances_)

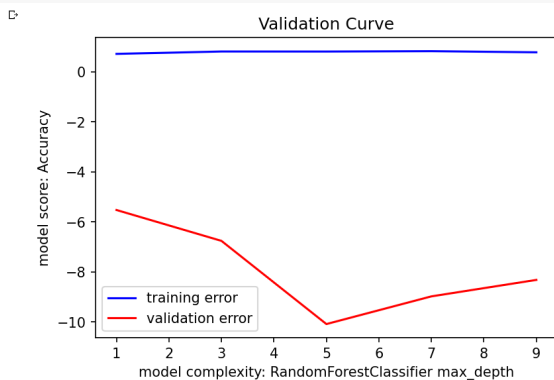

```

```
Feature Importances =  
[1.58960229e-03 5.30566745e-02 4.17418911e-03 1.59284318e-01  
2.61804152e-03 2.17672132e-02 7.81822081e-03 2.437228475e-02  
2.90624300e-03 2.56584641e-02 4.95719120e-02 2.69244094e-02  
1.12768719e-01 1.23221813e-04 4.58989637e-04 2.49941071e-03  
8.4319866e-03 0.00000000e+00 2.63169045e-02 2.37762734e-02  
0.00000000e+00 2.13822183e-01 1.81518181e-01 1.46478728e-04  
0.00000000e+00 2.3823658e-02 5.38444596e-03 0.00000000e+00  
2.61738842e-04 2.59317047e-02 0.00000000e+00 0.00000000e+00  
2.95326741e-03 0.00000000e+00 3.98608558e-03 3.56252545e-02  
3.63475050e-03 6.12666322e-03]
```

```
# Plot of feature importances from pure Random Forest Regressor  
%matplotlib inline  
import matplotlib.pyplot as plt  
# Get feature importances  
encoder = pipeline1.named_steps['onehotencoder']  
encoded = encoder.transform(X_train)  
rf = pipeline1.named_steps['randomforestregressor']  
importances1 = pd.Series(rf.feature_importances_, encoded.columns)  
# Plot feature importances  
n = 12  
plt.figure(figsize=(10,n/2))  
plt.title("Top (n) features pipeline1")  
importances1.sort_values()[::-n:].plot.barh(color='grey');
```



```
# Generate validation curves  
%matplotlib inline  
import numpy as np  
import matplotlib.pyplot as plt  
from sklearn.model_selection import validation_curve  
pipeline2 = make_pipeline(  
    ce.OrdinalEncoder(),  
    SimpleImputer(),  
    RandomForestRegressor()  
)  
  
depth = range(1, 10, 2)  
train_scores, val_scores = validation_curve(  
    pipeline2, X_train, y_train,  
    param_name='randomforestregressor__max_depth',  
    param_range=depth,  
    cv=3,  
    n_jobs=-1  
)  
  
plt.figure(dpi=150)  
plt.plot(depth, np.mean(train_scores, axis=1), color='blue', label='training error')  
plt.plot(depth, np.mean(val_scores, axis=1), color='red', label='validation error')  
plt.title('Validation Curve')  
plt.xlabel('model complexity: RandomForestClassifier max_depth')  
plt.ylabel('model score: Accuracy')  
plt.legend();
```



```
# Get drop-column importances  
column = 'Capital Land Area'  
  
pipeline3 = make_pipeline(  
    ce.OneHotEncoder(use_cat_names=True),  
    SimpleImputer(strategy='most_frequent'),  
    RandomForestRegressor(bootstrap=True, ccp_alpha=0, criterion='mse',  
        max_depth=3, max_features='auto', max_leaf_nodes=None,  
        max_samples=None, min_impurity_decrease=0.0,  
        min_impurity_split=None, min_samples_leaf=4,  
        min_samples_split=2, min_weight_fraction_leaf=0,  
        n_estimators=36, n_jobs=None, oob_score=False,  
        random_state=0, verbose=0, warm_start=False))  
  
# Fit without column  
pipeline3.fit(X_train.drop(columns=column), y_train)  
score_without = pipeline3.score(X_val.drop(columns=column), y_val)  
print(f'Validation Accuracy without {column}: {score_without}')  
  
# Fit with column  
pipeline3.fit(X_train, y_train)  
score_with = pipeline3.score(X_val, y_val)  
print(f'Validation Accuracy with {column}: {score_with}')  
  
# Compare the error with & without column  
print(f'Drop-Column Importance for {column}: {score_with - score_without}')  
  
Validation Accuracy without Capital Land Area: 0.10771758680677681  
Validation Accuracy with Capital Land Area: 0.069388620867604  
Drop-Column Importance for Capital Land Area: -0.03832895959917281
```

```
# Using Eli5 library which does not work with pipelines  
transformers = make_pipeline(  
    ce.OneHotEncoder(use_cat_names=True),  
    SimpleImputer(strategy='most_frequent')  
)  
  
X_train_transformed = transformers.fit_transform(X_train)  
X_val_transformed = transformers.transform(X_val)  
  
model1 = RandomForestRegressor(bootstrap=True, ccp_alpha=0, criterion='mse',  
    max_depth=3, max_features='auto', max_leaf_nodes=None,  
    max_samples=None, min_impurity_decrease=0.0,  
    min_impurity_split=None, min_samples_leaf=4,  
    min_impurity_split=2, min_weight_fraction_leaf=0,  
    n_estimators=36, n_jobs=None, oob_score=False,  
    random_state=0, verbose=0, warm_start=False)
```

```
modell.fit(X_train_transformed, y_train)

RandomForestRegressor(bootstrap=True, ccp_alpha=0, criterion='mse', max_depth=3,
                      max_features='auto', max_leaf_nodes=None,
                      max_samples=None, min_impurity_decrease=0.0,
                      min_impurity_split=None, min_samples_leaf=4,
                      min_samples_split=2, min_weight_fraction_leaf=0,
                      n_estimators=36, n_jobs=None, oob_score=False,
                      random_state=0, verbose=0, warm_start=False)

# Get permutation importances
! pip install eli5
from eli5.sklearn import PermutationImportance
import eli5

permuter = PermutationImportance(
    modell,
    scoring='r2',
    n_iter=2,
    random_state=42
)

permuter.fit(X_val_transformed, y_val)
feature_names = X_val.columns.tolist()

eli5.show_weights(
    permuter,
    top=None, # show permutation importances for all features
    feature_names=feature_names
)

Collecting eli5
  Downloading https://files.pythonhosted.org/packages/97/2f/c85c748f8548e468829971785347e14e45f45c6617da374711dec8cb38cc/eli5-0.10.1-py2.py3-none-any.whl (105kB)
    112kB 2.8MB/s
Requirement already satisfied: numpy>=1.9.0 in /usr/local/lib/python3.6/dist-packages (from eli5) (1.18.3)
Requirement already satisfied: scikit-learn>=0.18 in /usr/local/lib/python3.6/dist-packages (from eli5) (0.22.2.post1)
Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages (from eli5) (1.12.0)
Requirement already satisfied: graphviz in /usr/local/lib/python3.6/dist-packages (from eli5) (0.10.1)
Requirement already satisfied: attrs<16.0.0 in /usr/local/lib/python3.6/dist-packages (from eli5) (19.3.0)
Requirement already satisfied: Jinja2 in /usr/local/lib/python3.6/dist-packages (from eli5) (2.11.2)
Requirement already satisfied: tabulate>=0.7.7 in /usr/local/lib/python3.6/dist-packages (from eli5) (0.8.7)
Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-packages (from eli5) (1.4.1)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.6/dist-packages (from scikit-learn>=0.18->eli5) (0.14.1)
Requirement already satisfied: MarkupSafe>=0.23 in /usr/local/lib/python3.6/dist-packages (from Jinja2->eli5) (1.1.1)
Installing collected packages: eli5
Successfully installed eli5-0.10.1
/usr/local/lib/python3.6/dist-packages/sklearn/utils/deprecation.py:144: FutureWarning: The sklearn.metrics.scorer module is deprecated in version 0.22 and will be removed in version 0.24. The corresponding classes / functions should instead be imported from sklearn.metrics. Anythin
warnings.warn(message, FutureWarning)
/usr/local/lib/python3.6/dist-packages/sklearn/utils/deprecation.py:144: FutureWarning: The sklearn.feature_selection.base module is deprecated in version 0.22 and will be removed in version 0.24. The corresponding classes / functions should instead be imported from sklearn.feature
warnings.warn(message, FutureWarning)
Using TensorFlow backend.
      Weight      Feature
0.1883 ± 0.0054    Capital Land Area
0.0239 ± 0.2251      Longitude
0.0091 ± 0.0014      Density (Pm2)
0.0063 ± 0.0098    Sum of NUM_Black_or_African_American_BEN
0.0049 ± 0.0032    Sum of NUM_Asian_Pacific_Islander_BEN
0.0045 ± 0.0043      Latitude
0.0044 ± 0.0104    Became a State
0.0024 ± 0.0001    Capital Area Ratio
0.0022 ± 0.0019    Sum of Average_Age_of_BEN
0.0021 ± 0.0014    Boundaries
0.0015 ± 0.0025    Sum of NUM_BEN_With_Race_Not_Elsewhere_Classified
0.0007 ± 0.0016    % Urban Pop
0.0005 ± 0.0008    On Coast
0.0003 ± 0.0014    Sum of NUM_Medicare_BEN
0 ± 0.0000        Capital is the Largest City
0 ± 0.0000        DaysSinceFirstPositive
0 ± 0.0000        Borders Another Country
0 ± 0.0000        DaysSinceInfection
0 ± 0.0000        Lowest elevation
0 ± 0.0000        Children0-18
0 ± 0.0000        Elevation Ratio
-0.0002 ± 0.0021    Latitude Difference to State Capital
-0.0002 ± 0.0001    Water Area
-0.0004 ± 0.0007    DaysSinceTestStart
-0.0004 ± 0.0004    Longitude Difference to State Capital
-0.0005 ± 0.0001    Capital Mean Elevation
-0.0005 ± 0.0003    Land Area
-0.0007 ± 0.0015    Mean Elevation
-0.0008 ± 0.0002    State Area Ratio
-0.0015 ± 0.0007    Highest Elevation
-0.0038 ± 0.0423    Sum of PCT_MEDICARE
-0.0069 ± 0.0009    DaysSinceStayatHomeOrder
-0.0077 ± 0.0167    Sum of NUM_American_IndianAlaska_Native_BEN
-0.0088 ± 0.0251    Number of bordering states
-0.0148 ± 0.0124    Capital Water Area
-0.0154 ± 0.0149    Children 0-18
-0.0157 ± 0.0063    Log GDP
-0.0038 ± 0.1228    Sum of NUM_Hispanic_BEN

from sklearn.metrics import mean_squared_error, r2_score

# Coefficient of determination r2 for the training set
pipeline_score = permuter.score(X_train_transformed,y_train)
print("Coefficient of determination r2 for the training set.: ", pipeline_score)

# Coefficient of determination r2 for the validation set
pipeline_score = permuter.score(X_val_transformed,y_val)
print("Coefficient of determination r2 for the validation set.: ", pipeline_score)

# The mean squared error
y_pred = permuter.predict(X_val_transformed)
print("Mean squared error: %.2f%% mean_squared_error(y_val, y_pred))

Coefficient of determination r2 for the training set.: 0.3813954634238487
Coefficient of determination r2 for the validation set.: 0.069388620867684
Mean squared error: 538.47

# Capital Land Area continues to be of importance
# Use Importances for Feature selection
print('Shape before removing features:', X_train.shape)

Shape before removing features: (38, 38)

# Remove features of 0 importance
zero_importance = 0.0
mask = permuter.feature_importances_ > zero_importance
features1 = X_train.columns[mask]
X_train = X_train[features1]
print('Shape after removing features:', X_train.shape)

Shape after removing features: (38, 14)

# Random forest classifier with fourteen features
X_val = X_val[features1]
pipeline4 = make_pipeline(
    ce.OneHotEncoder(use_cat_names=True),
    SimpleImputer(strategy='most_frequent'),
    RandomForestRegressor(bootstrap=True, ccp_alpha=0,
                        max_depth=3, max_features='auto', max_leaf_nodes=None,
                        max_samples=None, min_impurity_decrease=0.0,
                        min_impurity_split=None, min_samples_leaf=4,
                        min_samples_split=2, min_weight_fraction_leaf=0,
                        n_estimators=36, n_jobs=None, oob_score=False,
                        random_state=0, verbose=0, warm_start=False)
)

# Fit on train, score on val
pipeline4.fit(X_train, y_train);

from sklearn.metrics import mean_squared_error, r2_score

# Coefficient of determination r2 for the training set
pipeline_score = pipeline4.score(X_train,y_train)
print("Coefficient of determination r2 for the training set.: ", pipeline_score)

# Coefficient of determination r2 for the validation set
pipeline_score = pipeline4.score(X_val,y_val)
print("Coefficient of determination r2 for the validation set.: ", pipeline_score)

# The mean squared error
y_pred = pipeline4.predict(X_val)
print("Mean squared error: %.2f%% mean_squared_error(y_val, y_pred))

Coefficient of determination r2 for the training set.: 0.26984125793319236
Coefficient of determination r2 for the validation set.: 0.11859396791756194
Mean squared error: 582.42
```

```
pipeline4.fit(X_val, y_val)

# Plot of Features
%matplotlib inline
import matplotlib.pyplot as plt

# Get feature importances
encoder = pipeline4.named_steps['onehotencoder']
encoded = encoder.transform(X_val)
rf = pipeline4.named_steps['randomforestregressor']
importances2 = pd.Series(rf.feature_importances_, encoded.columns)

# Plot feature importances
#n = 4
n = 12
plt.figure(figsize=(18,n/2))
plt.title("Top {n} features pipeline4")
importances2.sort_values()[::-n:].plot.barh(color='grey');

# Top 12 features pipeline4
Sum of NUM_Asian_Pacific_Islander_BEN
Sum of NUM_Black_or_African_American_BEN
Capital Area Ratio
% Urban Pop
Boundaries
Became a State
Longitude
Sum of Average_Age_of_BEN
Capital Land Area
Density (Pmi2)
Sum of NUM_Medicare_BEN
On Coast
```

```
!pip install pdpbox

Collecting pdpbox
  Downloading https://files.pythonhosted.org/packages/87/73/ac7da5ba1c6c83a87c412e7e7bde91a18d6ecf4474906c3e73f6f93980d49/pdpbox-0.2.0.tar.gz (57.7MB)
    57.7MB 66kB/s
Requirement already satisfied: pandas in /usr/local/lib/python3.6/dist-packages (from pdpbox) (1.0.3)
Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from pdpbox) (1.18.2)
Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-packages (from pdpbox) (1.4.1)
Requirement already satisfied: matplotlib>=2.1.2 in /usr/local/lib/python3.6/dist-packages (from pdpbox) (3.2.1)
Requirement already satisfied: joblib in /usr/local/lib/python3.6/dist-packages (from pdpbox) (0.14.1)
Requirement already satisfied: puzll in /usr/local/lib/python3.6/dist-packages (from pdpbox) (5.4.8)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.6/dist-packages (from pdpbox) (0.22.2.post1)
Requirement already satisfied: python-dateutil>=2.6.1 in /usr/local/lib/python3.6/dist-packages (from pandas->pdpbox) (2.8.1)
Requirement already satisfied: pytz>=2007.2 in /usr/local/lib/python3.6/dist-packages (from pandas->pdpbox) (2019.9)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.6/dist-packages (from matplotlib>=2.1.2->pdpbox) (1.2.0)
Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,!=2.0.1 in /usr/local/lib/python3.6/dist-packages (from matplotlib>=2.1.2->pdpbox) (2.4.7)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.6/dist-packages (from matplotlib>=2.1.2->pdpbox) (0.10.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.6/dist-packages (from python-dateutil>=2.6.1->pdpbox) (1.12.0)
Building wheels for collected packages: pdpbox
  Building wheel for pdpbox (setup.py) ... done
Created wheel for pdpbox: filename=pdpbox-0.2.0-cp36-none-any.whl size=57698722 sha256=2eb21b6a2e7bf8b4ba562bf76671ac3e8c2da28d9180574197f84c6005f80ca
Stored in directory: /root/.cache/pip/wheels/7d/08/51/63fd122b04a2c87d788464eeff94867c75bd96ae64d580a3fe
Successfully built pdpbox
Installing collected packages: pdpbox
Successfully installed pdpbox-0.2.0
```

```
model2 = RandomForestRegressor(bootstrap=True, ccp_alpha=0,
                               max_depth=3, max_features='auto', max_leaf_nodes=None,
                               max_samples=None, min_impurity_decrease=0.0,
                               min_impurity_split=None, min_samples_leaf=4,
                               min_samples_split=2, min_weight_fraction_leaf=0,
                               n_estimators=36, n_jobs=None, oob_score=False,
                               random_state=0, verbose=0, warm_start=False)

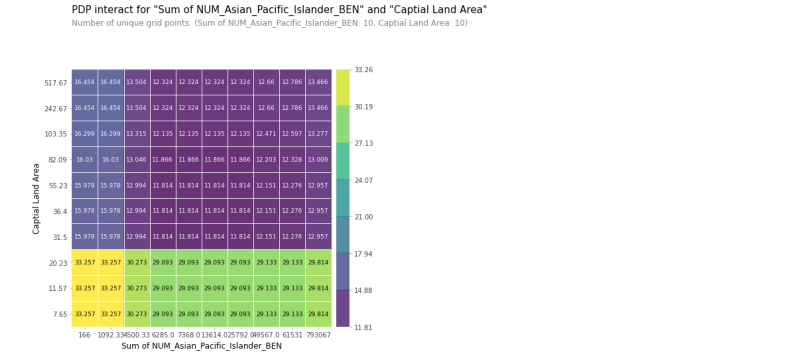
model2.fit(X_train, y_train)

RandomForestRegressor(bootstrap=True, ccp_alpha=0, criterion='mse', max_depth=3,
                       max_features='auto', max_leaf_nodes=None,
                       max_samples=None, min_impurity_decrease=0.0,
                       min_impurity_split=None, min_samples_leaf=4,
                       min_samples_split=2, min_weight_fraction_leaf=0,
                       n_estimators=36, n_jobs=None, oob_score=False,
                       random_state=0, verbose=0, warm_start=False)

# Partial Dependence Plots with 2 features
from pdpbox.pdp import pdp_interact, pdp_interact_plot
#features2 = ['Density (Pmi2)', 'Capital Land Area']
features2 = ['Sum of NUM_Asian_Pacific_Islander_BEN', 'Capital Land Area']
interaction = pdp_interact(
#
    model=model2,
    dataset=X_val,
    model_features=X_val.columns,
    features=features2
)

pdp_interact_plot(interaction, plot_type='grid', feature_names=features2);

findfont: Font family ['Arial'] not found. Falling back to DejaVu Sans.
findfont: Font family ['Arial'] not found. Falling back to DejaVu Sans.
findfont: Font family ['Arial'] not found. Falling back to DejaVu Sans.
findfont: Font family ['Arial'] not found. Falling back to DejaVu Sans.
```



```
# A two feature partial dependence plot in 3D
pdp = interaction.pdp.pivot_table(
    values='preds',
    columns=features2[0],
    index=features2[1]
)[::-1] # Slice notation to reverse index order so y axis is ascending

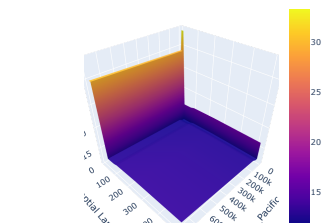
import plotly.graph_objs as go

target = 'Value of a parameter'

surface = go.Surface(x=pdp.columns,
                    y=pdp.index,
                    z=pdp.values)

layout = go.Layout(
    scene=dict(
        xaxis=dict(title=features2[0]),
        yaxis=dict(title=features2[1]),
        zaxis=dict(title=target)
    )
)

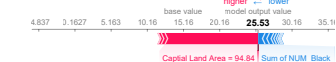
fig = go.Figure(surface, layout)
fig.show()
```



```
! pip install shap==0.23.0
! pip install -I shap
```

```
Collecting shap-0.23.0
  Downloading https://files.pythonhosted.org/packages/68/D0/8bd96821f73bdc2892a9d9ea92ca5cf2f925466537272bae8916c/shap-0.23.0.tar.gz (182KB)
    [■■■■■■■■■■] 184K 2.0MB/s
Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from shap==0.23.0) (1.18.3)
Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-packages (from shap==0.23.0) (1.4.1)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.6/dist-packages (from shap==0.23.0) (0.22.2.post1)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.6/dist-packages (from shap==0.23.0) (3.2.1)
Requirement already satisfied: pandas in /usr/local/lib/python3.6/dist-packages (from shap==0.23.0) (1.0.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.6/dist-packages (from shap==0.23.0) (4.38.0)
Requirement already satisfied: ipython in /usr/local/lib/python3.6/dist-packages (from shap==0.23.0) (5.5.0)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.6/dist-packages (from scikit-learn==0.23.0) (0.14.1)
Requirement already satisfied: python-dateutil<2.1.1 in /usr/local/lib/python3.6/dist-packages (from matplotlib==0.23.0) (2.8.1)
Requirement already satisfied: pyparsing=>=4.0.1,<4.1.2,!=4.1.2,!=4.2.0,!=4.2.0.1 in /usr/local/lib/python3.6/dist-packages (from matplotlib==0.23.0) (2.4.7)
Requirement already satisfied: kiwisolver>=0.1 in /usr/local/lib/python3.6/dist-packages (from matplotlib==0.23.0) (1.2.0)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.6/dist-packages (from matplotlib==0.23.0) (0.10.0)
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.6/dist-packages (from pandas==0.23.0) (2018.9)
Requirement already satisfied: pexpect in /usr/local/lib/python3.6/dist-packages (from ipython==shap==0.23.0) (2.1.3)
Requirement already satisfied: traitlets>=0.8 in /usr/local/lib/python3.6/dist-packages (from ipython==shap==0.23.0) (0.8.1)
Requirement already satisfied: prompt-toolkit<2.0.0,>=1.0.4 in /usr/local/lib/python3.6/dist-packages (from ipython==shap==0.23.0) (1.0.18)
Requirement already satisfied: pickleshare in /usr/local/lib/python3.6/dist-packages (from ipython==shap==0.23.0) (0.7.5)
Requirement already satisfied: setuptools>=18.5 in /usr/local/lib/python3.6/dist-packages (from ipython==shap==0.23.0) (46.1.3)
Requirement already satisfied: traitlets>=2 in /usr/local/lib/python3.6/dist-packages (from python-dateutil==>1-matplotlib==0.23.0) (1.12.0)
Requirement already satisfied: pypprocess>=0.5 in /usr/local/lib/python3.6/dist-packages (from pexpect; sys_platform == "win32">ipython==shap==0.23.0) (0.6.0)
Requirement already satisfied: ipython-genutils in /usr/local/lib/python3.6/dist-packages (from traitlets>=4.2-ipython==shap==0.23.0) (0.2.0)
Requirement already satisfied: wcwidth in /usr/local/lib/python3.6/dist-packages (from prompt-toolkit<2.0.0,>=1.0.4-ipython==shap==0.23.0) (0.1.9)
Building wheels for collected packages: shap
  Building wheel for shap (setup.py) ... done
Created wheel for shap: filename=shap-0.23.0-cp36-cp36m-linux_x86_64.whl size=235674 sha256=88888aebed63695a0bf0845e9a249aa16412ee8b20ba7b63e95ad9b3da4
Stored in directory: /root/.cache/pip/wheels/c1/2c/a4/bd1782feeb653efcd54ac2af8ade408f5778236838559
Successfully built shap
Installing collected packages: shap
Successfully installed shap-0.23.0
Collecting shap
  Downloading https://files.pythonhosted.org/packages/a8/77/504ed3e12a2ba543aac460671b0eb506caf788af2fb57bc7dc4c299845c/shap-0.35.0.tar.gz (273KB)
    [■■■■■■■■■■] 276K 2.0MB/s
Collecting numpy
  Downloading https://files.pythonhosted.org/packages/83/B7/e337e6de652fab9cf64fc72b2c0516cb43808d020ace31b82803db/numpy-1.18.4-cp36-cp36m-manylinux1_x86_64.whl (20.2MB)
    [■■■■■■■■■■] 20.2MB 1.2MB/s
Collecting scipy
  Downloading https://files.pythonhosted.org/packages/dc/29/162476404203116e7288cfb9352e9d03b7c49545dfec3539902f6a5/scipy-1.4.1-cp36-cp36m-manylinux1_x86_64.whl (26.1MB)
    [■■■■■■■■■■] 26.1MB 1.4MB/s
Collecting scikit-learn
  Downloading https://files.pythonhosted.org/packages/58/d9/312b03dafc478663e17d892fd2440b7237f6fec6ad5a484f172d777a32/scikit_learn-0.22.2.post1-cp36-cp36m-manylinux1_x86_64.whl (7.1MB)
    [■■■■■■■■■■] 7.1MB 41.3MB/s
Collecting pandas
  Downloading https://files.pythonhosted.org/packages/fb/1f/853bdcbbcf67c107888b0def7255767e4574407c6ff48081049b94943/pandas-1.0.3-cp36-cp36m-manylinux1_x86_64.whl (10.0MB)
    [■■■■■■■■■■] 10.0MB 174KB/s
Collecting tqdm>=4.25.0
  Downloading https://files.pythonhosted.org/packages/c9/40/0580128ba1be35f80c9a1cf2c24c7f809547f2d5ebc75b750919fda0/tqdm-4.46.0-py2.py3-none-any.whl (63kB)
    [■■■■■■■■■■] 71KB 8.9MB/s
Collecting joblib>=0.11
  Downloading https://files.pythonhosted.org/packages/28/5c/cfa62b65a321d4208efc0ff4c2689ef4c26261f8f64db28547f/joblib-1.0.1-py2.py3-none-any.whl (294KB)
    [■■■■■■■■■■] 290K 51.1MB/s
Collecting python-dateutil<2.6.1
  Downloading https://files.pythonhosted.org/packages/44/40/d6045d34d6ef758697a207aee09708d6306a2f7bc5413478013bc/python_dateutil-1.8.1-py2.py3-none-any.whl (227KB)
    [■■■■■■■■■■] 235K 41.7MB/s
Collecting pytz>=2017.2
  Downloading https://files.pythonhosted.org/packages/4f/ad/879454dad9688e7fad91a595744efda380781c745fd2ce2a3ad7f08880d/pytz-2020.1-py2.py3-none-any.whl (510KB)
    [■■■■■■■■■■] 512K 50.4MB/s
Collecting six>=1.5
  Downloading https://files.pythonhosted.org/packages/65/eb/197fc97bf7c390a2769695cf6ae0b7d5a282f5858802d4430cc5c1d9a/six-1.14.0-py2.py3-none-any.whl
Building wheels for collected packages: shap
  Building wheel for shap (setup.py) ... done
Created wheel for shap: filename=shap-0.23.0-cp36-cp36m-linux_x86_64.whl size=234123 sha256=5a6da0bee64fadd8c98344a30e59224018eb6e84f8d7ce5a6042e8ddca7cd
Stored in directory: /root/.cache/pip/wheels/ef/7f/f8/07558080c8f894906b3d3616d2c7d0f6b1251d5161b32cac
Successfully built shap
ERROR: google-colab 1.0.0 has requirement ipykernel<=1.2.0, but you'll have ipykernel 5.14.0 which is incompatible.
ERROR: datascience 0.10.0 has requirement folium==0.2.1, but you'll have folium 0.8.3 which is incompatible.
ERROR: convertdate 2.2.0 has requirement pytz@2020, >=2014.10, but you'll have pytz 2020.1 which is incompatible.
ERROR: abumaturus 0.1.12 has requirement iugate@0.2.7, >=0.1.25, but you'll have iugate 0.2.9 which is incompatible.
Install collected packages: shap, numpy, scipy, joblib, scikit-learn, shap, python-dateutil, pytz, pandas, tqdm, shap
Successfully installed joblib-1.0.1 numpy-1.18.4 pandas-1.0.3 python-dateutil-1.8.1 pytz-2020.1 scikit-learn-0.22.2.post1 scipy-1.4.1 shap-0.35.0 six-1.14.0 tqdm-4.46.0
WARNING: The following packages were previously imported in this runtime:
datascience,joblib,numpy,pandas,pytz,scipy,six,sklearn,tqdm
You must restart the runtime in order to use newly installed versions.
```

```
# Local Interpretation using SHAP (for prediction at State # = 4, row 32)
import shap
shap.initjs()
explainer = shap.TreeExplainer(model2)
shap_values = explainer.shap_values(X_train)
i = 32
shap_force_plot(explainer.expected_value, shap_values[i], features=X_train.loc[i], feature_names=X_train.columns)
```



```
# Find Shapley Forces across the training sample i (1 @ 0 - 37)
processor = make_pipeline(
    ce.OrdinalEncoder(),
    SimpleImputer(strategy='median')
)

X_train_processed = processor.fit_transform(X_train)
column_names = X_train.columns
shap_values_array = pd.DataFrame(columns = column_names)

for i in range(len(X_train)):
    row = X_train.iloc[[i]]
    explainer = shap.TreeExplainer(model2)
    row_processed = processor.transform(row)
    shap_values_input = explainer.shap_values(row_processed)
    shap_values_array = np.concatenate((shap_values_array, shap_values_input), axis=0)
```

```
# Create a 3D plot of force as a function of state curve displacement from mean curve and features for validation set
# A two feature partial dependence plot in 3D
import plotly.graph_objs as go
surface = go.Surface(x=volume_names,
                    y=y_train,
                    z=shap_values_array)

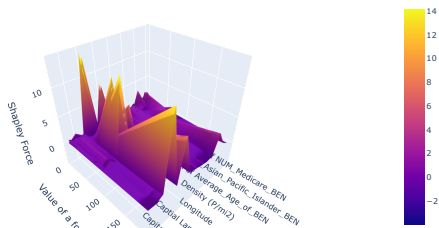
layout = go.Layout(
    scene=dict(
        xaxis=dict(title=''),
```



```

        yaxis=dict(title= 'Value of a for state'),
        zaxis=dict(title= 'Shapley Force')
    )
    fig = go.Figure(surface, layout)
    fig.show()

```



```

# Recursive Feature Elimination
from sklearn.feature_selection import RFE, f_regression
from sklearn.model_selection import StratifiedKFold

rfr = RandomForestRegressor(bootstrap=True, ccp_alpha=0,
                            max_depth=3, max_features='auto', max_leaf_nodes=None,
                            max_samples=None, min_impurity_decrease=0.0,
                            min_impurity_split=None, min_samples_leaf=4,
                            min_samples_split=2, min_weight_fraction_leaf=0,
                            n_estimators=36, n_jobs=None, oob_score=False,
                            random_state=0, verbose=0, warm_start=False)

#Selecting 7 features turns out to give maximum validation accuracy
number_selected_features = 7
rfe = RFE(rfr, n_features_to_select=number_selected_features, verbose=3)
rfe.fit(X_train,y_train)

# Fitting estimator with 14 features.
# Fitting estimator with 13 features.
# Fitting estimator with 12 features.
# Fitting estimator with 11 features.
# Fitting estimator with 10 features.
# Fitting estimator with 9 features.
# Fitting estimator with 8 features.
RFE(estimator=RandomForestRegressor(bootstrap=True, ccp_alpha=0,
                                     criterion='mse', max_depth=3,
                                     max_features='auto', max_leaf_nodes=None,
                                     max_samples=None, min_impurity_decrease=0.0,
                                     min_impurity_split=None, min_samples_leaf=4,
                                     min_samples_split=2,
                                     min_weight_fraction_leaf=0, n_estimators=36,
                                     n_jobs=None, oob_score=False,
                                     random_state=0, verbose=0,
                                     warm_start=False),
    n_features_to_select=7, step=1, verbose=3)

```

```

rfe_support = rfe.get_support()
rfe_feature = X_train.loc[:,rfe_support].columns.tolist()
print(str(len(rfe_feature)), 'selected features')

```

```

# 7 selected features

```

```

from sklearn.metrics import mean_squared_error, r2_score

# Coefficient of determination r2 for the training set
pipeline_score = rfe.score(X_train,y_train)
print('Coefficient of determination r2 for the training set.: ', pipeline_score)

# Coefficient of determination r2 for the validation set
pipeline_score = rfe.score(X_val,y_val)
print('Coefficient of determination r2 for the validation set.: ', pipeline_score)

# The mean squared error
y_pred = rfe.predict(X_val)
print('Mean squared error: %.2f'% mean_squared_error(y_val, y_pred))

# Coefficient of determination r2 for the training set.: 0.27619181813448996
# Coefficient of determination r2 for the validation set.: 0.14111881952867034
# Mean squared error: 489.58

```

```

# Retain only features with highest importance from RFE
X_train_rfe_select = X_train[rfe_feature]
X_val_rfe_select = X_val[rfe_feature]
print('Shape after removing features:', X_train_rfe_select.shape, X_val_rfe_select.shape)

# Shape after removing features: (38, 7) (13, 7)

```

```

# Random forest classifier after RFE Feature Selection on Reduced Feature Set

```

```

pipeline5 = make_pipeline(
    co.OneHotEncoder(use_cat_names=True),
    SimpleImputer(strategy='most_frequent'),
    RandomForestRegressor(bootstrap=True, ccp_alpha=0,
                        max_depth=3, max_features='auto', max_leaf_nodes=None,
                        max_samples=None, min_impurity_decrease=0.0,
                        min_impurity_split=None, min_samples_leaf=4,
                        min_samples_split=2, min_weight_fraction_leaf=0,
                        n_estimators=36, n_jobs=None, oob_score=False,
                        random_state=0, verbose=0, warm_start=False)
)

```

```

# Fit on train, score on val
pipeline5.fit(X_train_rfe_select, y_train);

```

```

# Coefficient of determination r2 for the training set
pipeline_score = pipeline5.score(X_train_rfe_select,y_train)
print('Coefficient of determination r2 for the training set.: ', pipeline_score)

# Coefficient of determination r2 for the validation set
pipeline_score = pipeline5.score(X_val_rfe_select,y_val)
print('Coefficient of determination r2 for the validation set.: ', pipeline_score)

# The mean squared error
y_pred = pipeline5.predict(X_val_rfe_select)
print('Mean squared error: %.2f'% mean_squared_error(y_val, y_pred))

# Coefficient of determination r2 for the training set.: 0.27619181813448996
# Coefficient of determination r2 for the validation set.: 0.14111881952867034
# Mean squared error: 489.58

```

```

pipeline5.fit(X_val_rfe_select, y_val)
# Plot of features
%matplotlib inline
import matplotlib.pyplot as plt

```

```

# Get feature importances
encoder = pipeline5.named_steps['onehotencoder']
encoded = encoder.transform(X_val_rfe_select)
rf = pipeline5.named_steps['randomforestregressor']
importances3 = pd.Series(rf.feature_importances_, encoded.columns)

```

```

# Plot feature importances
n = number_selected_features
plt.figure(figsize=(10,n/2))
plt.title('Top (n) features pipelines')
importances3.sort_values()[::-n:].plot.barh(color='grey');

```

