

## CIS400 Assignment 1 report

Huahao Shang  
[hushang@syr.edu](mailto:hushang@syr.edu)  
SUID: hushang

1.

- (a) Three files are successfully read using read\_csv and read\_xls and three files don't have NA data.
- (b) There are 3140 rows after the merge
- (c) There are 3093 rows and 23 columns

2.

(a)

```
# fips:discrete
contry-x: categorical
state_x: categorical
state_code_x: categorical
male: discrete
female: discrete
median_age: discrete
population: discrete
female_percentage: discrete
lat_x: continuous
long_x: continuous
Areaname: categorical
LND010200D: continuous
county_y: categorical
state_y: categorical
lat_y: continuous
Long_y: continuous
date: categorical
cases: discrete
stste_code_y: categorical
deaths: discrete
populatin-density: discrete
case-ratio: discrete
```

(b)

The histogram are not very surprised.

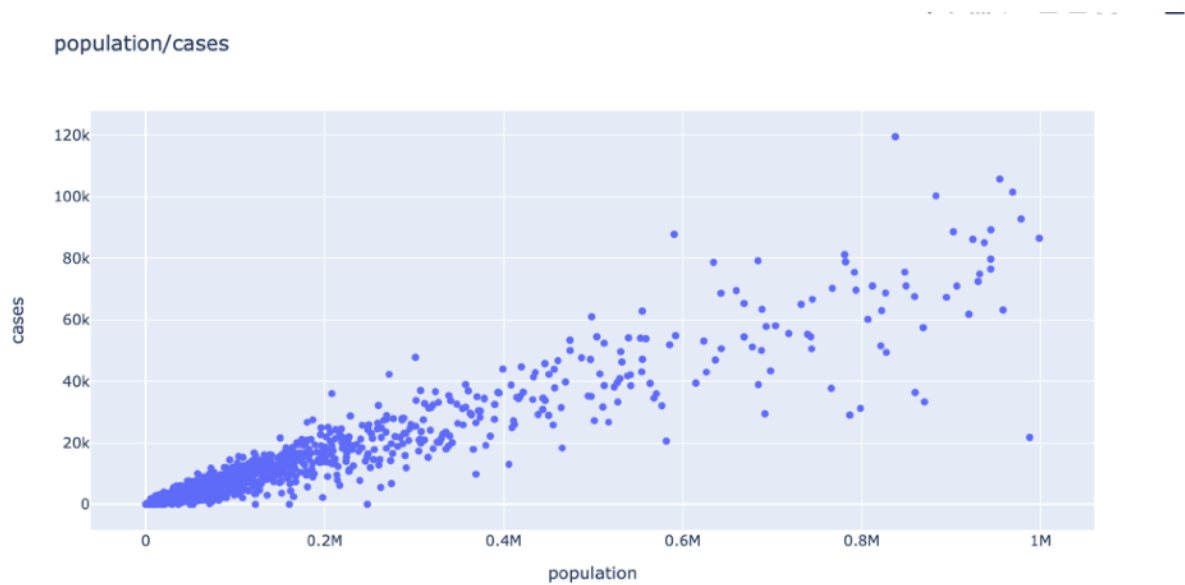
'male', 'female', 'median\_age', 'population', 'female\_percentage', 'cases', 'deaths', 'population\_density', 'case\_ratio' are obtained. The data and the graph are in a expected looking.

(c)

For heat map, I expected the cases, case\_ratio and death are related to the population and population density. And Land area, longitude and latitude are not correlated to any of the cases and death numbers.

The real heat map turns out as expected.

(d).



The trend is that cases will grow as population getting larger.

3

(a).

Train data length: 2072

Test data length: 1021

(b)

$R^2$  is 0.90

Beta 1 is 0.07 This indicates it will be a smooth slope

Beta 0 is 248.21

(c)

MAE is 1418.24

The error seems fine, from the shown pred\_cases and real cases, they are not have to much difference.

(d)

```
print(mean_absolute_error)
```

```
it[13]:
```

	pred_cases	cases
448	6,880.76	7273
667	1,545.67	1287
3139	808.99	617
230	479.63	97
3093	3,672.06	3834

Mean absolute error is  
1418.240890654271

```
141: # define function to import via libraries
```

Plot of predicted and actual

