

P8130 Final Project

Abstract

Introduction (brief context and background of the problem)

Methods (data description and statistical methods)

Results

Conclusions/Discussion

```
library(tidyverse)
```

Read in dataset

```
cdi = read_csv("./cdi.csv") %>%
  janitor::clean_names()

## no missing value
cdi %>%
  select(everything()) %>%
  summarise_all(funs(sum(is.na(.)))) %>%
  knitr::kable()
```

| id | cty | state | area | pop | pop18 | pop65 | docs | beds | crimes | hsgrad | bagrad | poverty | unemp | pcincom | totalinc | region |
|----|-----|-------|------|-----|-------|-------|------|------|--------|--------|--------|---------|-------|---------|----------|--------|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Data cleaning

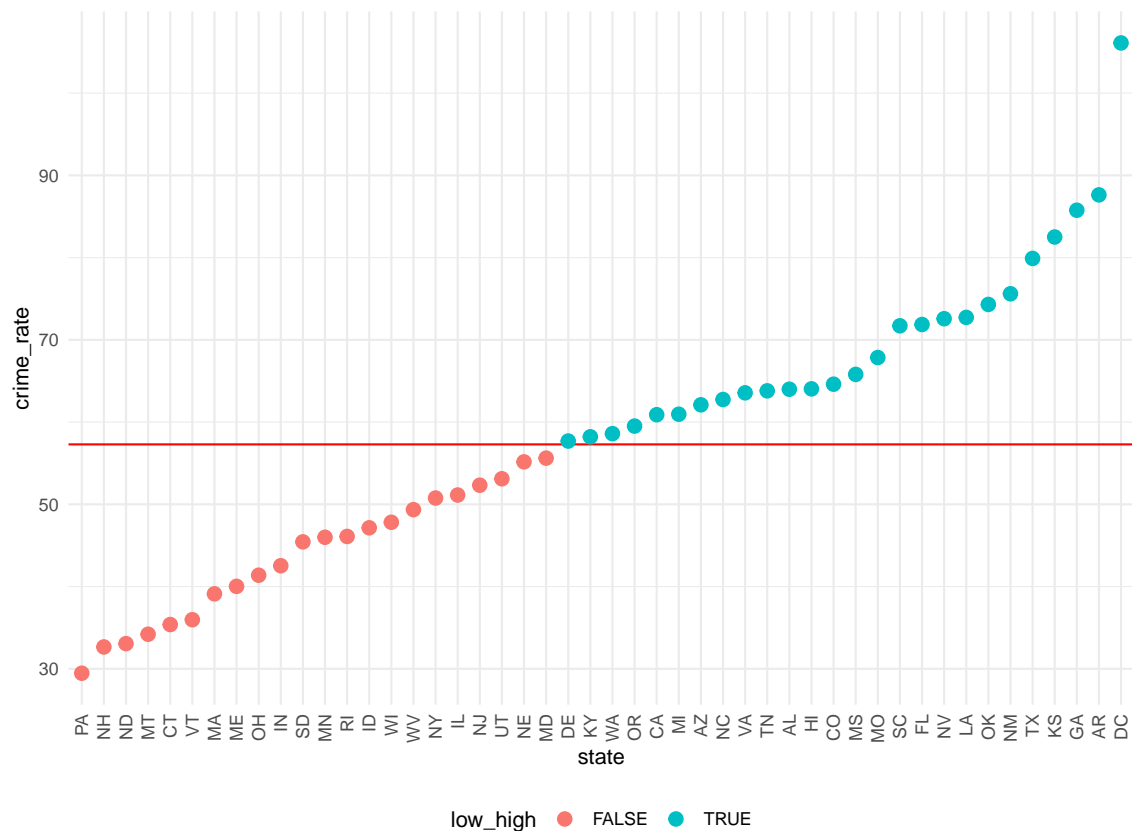
```
cdi =
  cdi %>%
  mutate(crm_1000 = crimes/pop*1000, # as indicated by the project prompt
         docs_rate_1000 = docs/pop*1000, # every 1000 people how many doctors
         beds_docs = beds/docs,
         region = factor(region)) %>%
  select(-id, -cty, -crimes)
cdi

## # A tibble: 440 x 17
##   state area   pop pop18 pop65 docs  beds hsgrad bagrad poverty unemp
##   <chr> <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 CA    4060 8863164 32.1  9.7 23677 27700  70    22.3  11.6    8
## 2 IL     946 5105067 29.2 12.4 15153 21550 73.4  22.8  11.1   7.2
## 3 TX    1729 2818199 31.3  7.1  7553 12449 74.9  25.4  12.5   5.7
## 4 CA    4205 2498016 33.5 10.9  5905  6179 81.9  25.3   8.1   6.1
## 5 CA     790 2410556 32.6  9.2  6062  6369 81.2  27.8   5.2   4.8
## 6 NY     71 2300664 28.3 12.4  4861  8942 63.7  16.6  19.5   9.5
## 7 AZ    9204 2122101 29.2 12.5  4320  6104 81.5  22.1   8.8   4.9
```

```
## 8 MI      614 2111687 27.4 12.5 3823 9490 70 13.7 16.9 10
## 9 FL      1945 1937094 27.1 13.9 6274 8840 65 18.8 14.2 8.7
## 10 TX      880 1852810 32.6 8.2 4718 6934 77.1 26.3 10.4 6.1
## # ... with 430 more rows, and 6 more variables: pcincome <dbl>, totalinc <dbl>,
## #   region <fct>, crm_1000 <dbl>, docs_rate_1000 <dbl>, beds_docs <dbl>
```

```
mean_crm = mean(cdi$crm_1000)
cdi_state = cdi %>%
  group_by(state) %>%
  summarize(crime_rate = mean(crm_1000)) %>%
  mutate(low_high = ifelse(crime_rate > mean_crm, TRUE, FALSE))
```

```
cdi_state %>%
  mutate(state = fct_reorder(state, crime_rate)) %>%
  ggplot(aes(x = state, y = crime_rate)) +
  geom_hline(yintercept = mean_crm, color = "red") +
  geom_point(aes(color = low_high), size = 3) +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1))
```



Data Exploration

```
## summary statistics
knitr::kable(summary(cdi))
```

| | state | area | pop | pop18 | pop65 | docs | beds | hsgrad | bagrad | poverty | unemp | pcincome | totalinc | region | crm_1000 | docs_rate_1000 | beds_docs | |
|--------|----------|---------|--------|--------|--------|---------|---------|--------|--------|---------|--------|----------|----------|--------|----------|----------------|-----------|---------|
| Length | Min. | Min. | Min. | Min. | Min. | Min. | Min. | Min. | Min. | Min. | Min. | Min. | Min. | 1:103 | Min. | Min. | Min. | |
| | : | : | : | 16.40 | : | : | : | 46.60 | : | : | : | : | : | : | : | : | 0.07969 | |
| | 15.0 | 100043 | | 3.000 | 39.0 | 92.0 | | 8.10 | 1.400 | 2.200 | 8899 | 1141 | | 4.601 | 0.3559 | | | |
| Class | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 1st | 2:108 | 1st | 1st | 1st | |
| | Qu.: | Qu.: | Qu.: | 26.20 | Qu.: | Qu.: | Qu.: | 73.98 | 15.28 | Qu.: | Qu.: | 16.018 | | Qu.: | Qu.: | Qu.: | 1.34565 | |
| | 451.2 | 139027 | | 9.875 | 182.8 | 390.8 | | | | 5.300 | 5.100 | | 2311 | | 38.102 | 1.2127 | | |
| Mode | Median | Median | Median | Median | Median | Median | Median | Median | Median | Median | Median | Median | Median | 3:152 | Median | Median | Median | |
| | : | : | : | 28.10 | 11.750 | : | : | 77.70 | 19.70 | : | : | 17759 | | : | : | : | 1.83419 | |
| | 656.5 | 217280 | | | 401.0 | 755.0 | | | | 7.900 | 6.200 | | 3857 | | 52.429 | 1.7509 | | |
| NA | Mean | Mean | Mean | Mean | Mean | Mean | Mean | Mean | Mean | Mean | Mean | Mean | Mean | 4: | Mean | Mean | Mean | |
| | : | : | : | 28.57 | 12.170 | : | : | 77.56 | 21.08 | : | : | 18561 | | 77 | : | : | 1.97855 | |
| | 1041.4 | 393011 | | | 988.0 | 1458.6 | | | | 8.721 | 6.597 | | 7869 | | 57.286 | 2.1230 | | |
| NA | 3rd | 3rd | 3rd | 3rd | 3rd | 3rd | 3rd | 3rd | 3rd | 3rd | 3rd | 3rd | 3rd | NA | 3rd | 3rd | 3rd | |
| | Qu.: | Qu.: | Qu.: | 30.02 | 13.025 | Qu.: | Qu.: | 82.40 | 25.32 | Qu.: | Qu.: | 20.900 | Qu.: | 20.270 | Qu.: | Qu.: | Qu.: | 2.42710 |
| | 946.8 | 436064 | | | 1036.0 | 1575.8 | | | | | 7.500 | | 8654 | | 72.597 | 2.4915 | | |
| NA | Max. | Max. | Max. | Max. | Max. | Max. | Max. | Max. | Max. | Max. | Max. | Max. | Max. | NA | Max. | Max. | Max. | |
| | :20062.3 | 8863164 | | 49.70 | 33.800 | 23677.2 | 27700.9 | 2.90 | 52.30 | 36.300 | 21.300 | 37541 | 184230 | | :295.98 | 17.0375 | 7.41667 | |

cdi

```
## # A tibble: 440 x 17
##   state area    pop pop18 pop65 docs  beds hsgrad bagrad poverty unemp
##   <chr> <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 CA    4060 8863164  32.1   9.7 23677 27700   70    22.3   11.6    8
## 2 IL     946 5105067  29.2  12.4 15153 21550   73.4   22.8   11.1   7.2
## 3 TX    1729 2818199  31.3   7.1  7553 12449   74.9   25.4   12.5   5.7
## 4 CA    4205 2498016  33.5  10.9  5905  6179   81.9   25.3    8.1   6.1
## 5 CA     790 2410556  32.6   9.2  6062  6369   81.2   27.8    5.2   4.8
## 6 NY     71 2300664  28.3  12.4  4861  8942   63.7   16.6   19.5   9.5
## 7 AZ   9204 2122101  29.2  12.5  4320  6104   81.5   22.1    8.8   4.9
## 8 MI     614 2111687  27.4  12.5  3823  9490    70   13.7   16.9  10
## 9 FL   1945 1937094  27.1  13.9  6274  8840    65   18.8   14.2   8.7
## 10 TX    880 1852810  32.6   8.2  4718  6934   77.1   26.3   10.4   6.1
## # ... with 430 more rows, and 6 more variables: pcincome <dbl>, totalinc <dbl>,
## #   region <fct>, crm_1000 <dbl>, docs_rate_1000 <dbl>, beds_docs <dbl>
```