

# P8130 Final Project

## Abstract

## Introduction (brief context and background of the problem)

## Methods (data description and statistical methods)

## Results

## Conclusions/Discussion

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5      v purrr   0.3.4
## v tibble  3.1.5      v dplyr   1.0.7
## v tidyr   1.1.3      v stringr 1.4.0
## v readr   2.0.1      v forcats 0.5.1
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
library(pastecs)
```

```
##
```

```
## Attaching package: 'pastecs'
```

```
## The following objects are masked from 'package:dplyr':
```

```
##
```

```
##      first, last
```

```
## The following object is masked from 'package:tidyr':
```

```
##
```

```
##      extract
```

## Read in dataset

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

```
## Rows: 440 Columns: 17

## -- Column specification -----
## Delimiter: ","
## chr (2): cty, state
## dbl (15): id, area, pop, pop18, pop65, docs, beds, crimes, hsgrad, bagrad, p...

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
cdi
```

```
## # A tibble: 440 x 17
##       id cty    state area    pop pop18 pop65 docs  beds crimes hsgrad bagrad
##   <dbl> <chr>  <chr> <dbl>  <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1     1  Los_An~ CA    4060 8.86e6 32.1   9.7 23677 27700 688936 70    22.3
## 2     2   Cook  IL     946 5.11e6 29.2  12.4 15153 21550 436936 73.4  22.8
## 3     3  Harris TX    1729 2.82e6 31.3   7.1 7553 12449 253526 74.9  25.4
## 4     4  San_Di~ CA    4205 2.50e6 33.5  10.9 5905 6179 173821 81.9  25.3
## 5     5   Orange CA     790 2.41e6 32.6   9.2 6062 6369 144524 81.2  27.8
## 6     6   Kings NY      71 2.30e6 28.3  12.4 4861 8942 680966 63.7  16.6
## 7     7  Marico~ AZ   9204 2.12e6 29.2  12.5 4320 6104 177593 81.5  22.1
## 8     8   Wayne MI     614 2.11e6 27.4  12.5 3823 9490 193978 70    13.7
## 9     9    Dade FL   1945 1.94e6 27.1  13.9 6274 8840 244725 65    18.8
## 10    10 Dallas TX     880 1.85e6 32.6   8.2 4718 6934 214258 77.1  26.3
## # ... with 430 more rows, and 5 more variables: poverty <dbl>, unemp <dbl>,
## #   pcincome <dbl>, totalinc <dbl>, region <dbl>
```

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

```
## Warning: 'funs()' was deprecated in dplyr 0.8.0.
## Please use a list of either functions or lambdas:
##
##   # Simple named list:
##   list(mean = mean, median = median)
##
##   # Auto named with 'tibble::lst()':
##   tibble::lst(mean, median)
##
##   # Using lambdas
##   list(~ mean(., trim = .2), ~ median(., na.rm = TRUE))
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was generated.
```

```
## # A tibble: 1 x 17
##       id cty state area    pop pop18 pop65 docs  beds crimes hsgrad bagrad
##   <int> <int> <int> <int> <int> <int> <int> <int> <int> <int> <int> <int>
## 1     0     0     0     0     0     0     0     0     0     0     0     0
## # ... with 5 more variables: poverty <int>, unemp <int>, pcincome <int>,
## #   totalinc <int>, region <int>
```

## 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)
cdi
```

```
## # A tibble: 440 x 20
##       id cty      state area      pop pop18 pop65 docs  beds crimes hsgrad bagrad
##   <dbl> <chr>   <chr> <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1     1  Los_An~ CA      4060 8.86e6 32.1   9.7 23677 27700 688936 70    22.3
## 2     2   Cook    IL       946 5.11e6 29.2  12.4 15153 21550 436936 73.4  22.8
## 3     3  Harris  TX      1729 2.82e6 31.3   7.1 7553 12449 253526 74.9  25.4
## 4     4  San_Di~ CA      4205 2.50e6 33.5  10.9 5905 6179 173821 81.9  25.3
## 5     5  Orange  CA       790 2.41e6 32.6   9.2 6062 6369 144524 81.2  27.8
## 6     6  Kings   NY       71 2.30e6 28.3  12.4 4861 8942 680966 63.7  16.6
## 7     7  Marico~ AZ     9204 2.12e6 29.2  12.5 4320 6104 177593 81.5  22.1
## 8     8  Wayne  MI       614 2.11e6 27.4  12.5 3823 9490 193978 70    13.7
## 9     9   Dade   FL     1945 1.94e6 27.1  13.9 6274 8840 244725 65    18.8
## 10    10 Dallas TX       880 1.85e6 32.6   8.2 4718 6934 214258 77.1  26.3
## # ... with 430 more rows, and 8 more variables: poverty <dbl>, unemp <dbl>,
## #   pcincome <dbl>, totalinc <dbl>, region <dbl>, crm_1000 <dbl>,
## #   docs_rate_1000 <dbl>, beds_docs <dbl>
```

## Data Exploration

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

id	cty	state	area	pop	pop18	pop65	docs	beds	crimes	hsgrad	bagrad	poverty	unemp	pcincome	totalinc	region	crm_1000	docs_rate_1000	beds_docs
Min.	Length	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.	Min.
:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:
1.0			15.0	100043	3.000	39.0	92.0	563	8.10	1.400	2.200	8899	1141	4.60	10.3559				
1st	Class	Class	1st	1st	1st	1st	1st	1st	1st	1st	1st	1st	1st	1st	1st	1st	1st	1st	1st
Qu.:1	1st	1st	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1	Qu.:1
ac-	ac-	451.2	139027	9.875	182.8	390.8	6220	5.300	5.100	2311	38.10	2.127							
ter	ter																		
Median	Mode	Mode	Median	Median	Median	Median	Median	Median	Median	Median	Median	Median	Median	Median	Median	Median	Median	Median	Median
:220.5	char-:	char-:	:	:28.10	11.750	:	:	:77.70	19.70	:	:17759	:3.000	:	:1.834	19				
ac-	ac-	656.5	217280	401.0	755.0	11820	7.900	6.200	3857	52.42	9.7509								
ter	ter																		
Mean	NA	NA	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean	Mean
:220.5	:	:	:28.57	12.170	:	:	:77.56	21.08	:	:18561	:2.461	:	:1.978	55					
			1041.3	93011	988.0	1458.2	7112	8.72	16.597	7869	57.28	2.1230							
3rd	NA	NA	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd	3rd
Qu.:3	3rd	3rd	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3	Qu.:3
330.2			946.8	436064	1036.0	1575.2	6280	7.500	8654	72.59	2.4915								

id	cty	state	area	pop	pop18	pop65	docs	beds	crimes	hsgrad	bagrad	poverty	unemp	pcincome	totalinc	region	crm_1000	beds_1000
Max.	NA	NA	Max.	Max.	Max.	Max.	Max.	Max.	Max.	Max.	Max.	Max.	Max.	Max.	Max.	Max.	Max.	Max.
:440.0			:2006288631491.70	33.800236772077006889362.90	52.3036.30021.30037541184230.000295.937.037741667													

```
cdi
```

```
## # A tibble: 440 x 20
##       id cty      state  area    pop pop18 pop65  docs  beds crimes hsgrad bagrad
##   <dbl> <chr>   <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1     1   1 Los_An~  CA    4060 8.86e6 32.1   9.7 23677 27700 688936   70    22.3
## 2     2   2 Cook    IL     946 5.11e6 29.2  12.4 15153 21550 436936  73.4   22.8
## 3     3   3 Harris  TX    1729 2.82e6 31.3   7.1  7553 12449 253526  74.9   25.4
## 4     4   4 San_Di~  CA    4205 2.50e6 33.5  10.9  5905  6179 173821  81.9   25.3
## 5     5   5 Orange  CA     790 2.41e6 32.6   9.2  6062  6369 144524  81.2   27.8
## 6     6   6 Kings   NY     71 2.30e6 28.3  12.4  4861  8942 680966  63.7   16.6
## 7     7   7 Marico~  AZ    9204 2.12e6 29.2  12.5  4320  6104 177593  81.5   22.1
## 8     8   8 Wayne   MI     614 2.11e6 27.4  12.5  3823  9490 193978   70    13.7
## 9     9   9 Dade    FL    1945 1.94e6 27.1  13.9  6274  8840 244725   65    18.8
## 10    10  10 Dallas  TX     880 1.85e6 32.6   8.2  4718  6934 214258  77.1   26.3
## # ... with 430 more rows, and 8 more variables: poverty <dbl>, unemp <dbl>,
## #   pcincome <dbl>, totalinc <dbl>, region <dbl>, crm_1000 <dbl>,
## #   docs_rate_1000 <dbl>, beds_docs <dbl>
```