

EDS241: Assignment 1

Paloma Cartwright

2022-01-15

The data for this assignment come from CalEnviroScreen 4.0, a mapping and data tool produced by the California Office of Environmental Health Hazards Assessment (OEHHA). The data are compiled and constructed from a variety of sources and cover all 8,035 census tracts in California. Source: <https://oehha.ca.gov/calenviroscreen/report/calenviroscreen-40>.

The full data are contained in the file CES4.xls, which is available on Gauchospace (note that the Excel file has three “tabs” or “sheets”). The data is in the tab “CES4.0FINAL_results” and “Data Dictionary” contains the definition of the variables.

1 Clean and plot data

The following code loads and cleans the data.

```
# Load data and clean names
ces4_raw <- read_excel("CES4.xlsx", sheet = "CES4.0FINAL_results") %>%
  na_if("NA") %>%
  clean_names()

# Select columns wanted data
ces4 <- ces4_raw %>%
  select("census_tract", "total_population",
         "california_county", "low_birth_weight", "pm2_5", "poverty")

ces4_clean <- na.omit(ces4)
```

The following code chunks answer the questions provided in the assignment.

(a)

```
# Avg PM2.5 Conc

avg_pm25_conc <- mean(ces4$pm2_5)
```

The average PM2.5 concentration across all census tracts in California is 10.1526999.

(b)

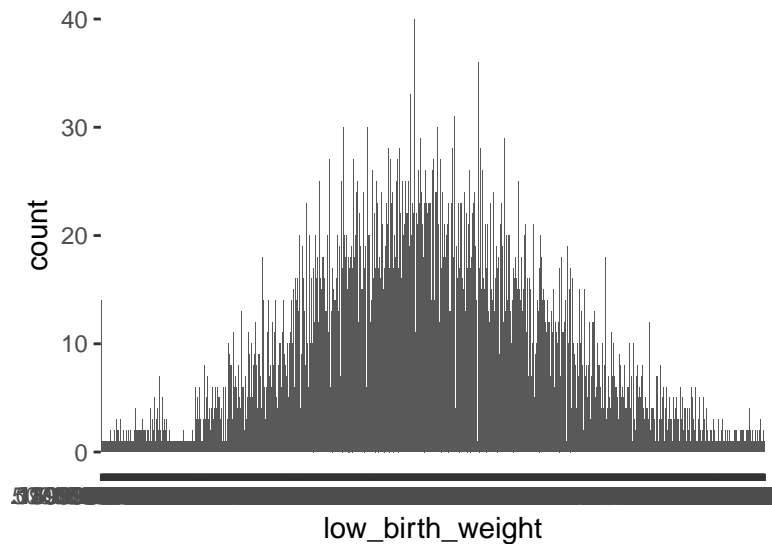
```
# highest poverty

highest_poverty <- ces4 %>%
  group_by(california_county) %>%
  summarize(avg_poverty = na.omit(mean(poverty))) %>%
  arrange(desc(avg_poverty))
```

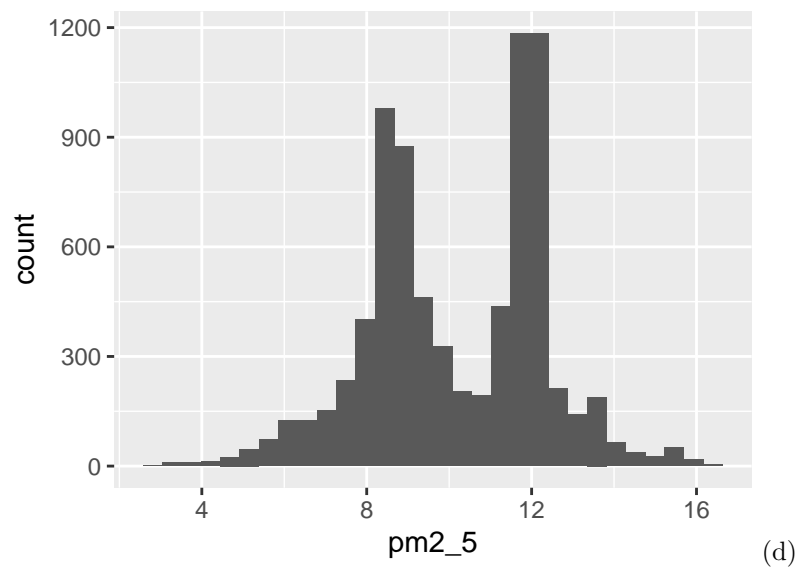
The county with the highest level of poverty is: .

(c) The histograms are shown below:

```
ggplot(data = ces4_clean, aes(x = low_birth_weight)) +  
  geom_bar()
```



```
ggplot(data = ces4_clean, aes(x = pm2_5)) +  
  geom_histogram()
```



(d)

```
model <- lm_robust(low_birth_weight ~ pm2_5, data = ces4)  
summary(model)
```

```
##  
## Call:  
## lm_robust(formula = low_birth_weight ~ pm2_5, data = ces4)  
##  
## Standard error type: HC2  
##  
## Coefficients:
```

```
##           Estimate Std. Error t value Pr(>|t|) CI Lower CI Upper DF
## (Intercept)  3.8010    0.088583  42.91 0.000e+00  3.6273  3.9746 7806
## pm2_5        0.1179    0.008402  14.04 3.256e-44  0.1015  0.1344 7806
##
## Multiple R-squared:  0.02499 ,    Adjusted R-squared:  0.02486
## F-statistic:  197 on 1 and 7806 DF,  p-value: < 2.2e-16
```

The slope coefficient for the regression of Low Birth Rate on PM2.5 is 0.1179305. The heteroskedasticity-robust standard error is 0.0084024. The slope coefficient tells us that for every 1 microgram per cubic meter increase in PM2.5 concentration, there is a 0.1179305 increase in the percentage of census tract births with weight less than 2500g. This is statistically significant at the 5% level.

(e)
(f)

```
model2 <- lm_robust(low_birth_weight ~ pm2_5 + poverty, data = ces4)
summary(model2)
```

```
##
## Call:
## lm_robust(formula = low_birth_weight ~ pm2_5 + poverty, data = ces4)
##
## Standard error type:  HC2
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|) CI Lower CI Upper DF
## (Intercept)  3.54374    0.084733  41.823 0.000e+00  3.37764  3.70984 7802
## pm2_5        0.05911    0.008293   7.127 1.116e-12  0.04285  0.07536 7802
## poverty      0.02744    0.001002  27.374 1.287e-157  0.02547  0.02940 7802
##
## Multiple R-squared:  0.1169 ,    Adjusted R-squared:  0.1167
## F-statistic: 494.8 on 2 and 7802 DF,  p-value: < 2.2e-16
```

The estimated coefficient for poverty is 0.0274353. This tells us that holding PM2.5 constant, for every one percent increase in the population of the census tract living below twice the federal poverty line there is a 0.0274353 increase in the percentage of census tract births with weight less than 2500g. The coefficient for PM2.5 increases which displays that when poverty is considered in the regression, the PM2.5 value becomes more prominent. These variables are correlated. In the model that omits poverty, the regression coefficient of PM2.5 captures the effect of poverty on low birth weights.

(g)

```
linhyp <- linearHypothesis(model2, c("pm2_5 = 0", "poverty=0"),
                             white.adjust = "hc2")
linhyp
```

Res.Df	Df	Chisq	Pr(>Chisq)
7.8e+03			
7.8e+03	2	990	1.29e-215

With the p-value being $1.2877564 \times 10^{-215}$, we can reject the null hypothesis that the effect of PM2.5 is equal to the effect of poverty.