

EDS241: Assignment 1

Marie Rivers

01/16/2022

In this assignment, we use air quality data in R to investigate the relationship between PM2.5 and low birth weight in California. The data came 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 9,035 census tracts in California. Source: <https://oehha.ca.gov/calenviroscreen/report/calenviroscreen-40>

xxx...take a closer look at the data dictionary

1 Read and Clean Data

```
data <- read_xlsx(here("data", "CES4.xlsx"), sheet = "CES4.OFINAL_results") %>%
  rename(CensusTract = "Census Tract", TotalPopulation = "Total Population", CaliforniaCounty = "California County") %>%
  select(CensusTract, TotalPopulation, CaliforniaCounty, LowBirthWeight, PM25, Poverty) %>%
  mutate(LowBirthWeight = as.numeric(LowBirthWeight))
```

```
# data <- read_xlsx(here("data", "CES4.xlsx"), sheet = "CES4.OFINAL_results") %>%
#   rename(CensusTract = "Census Tract", TotalPopulation = "Total Population", CaliforniaCounty = "California County") %>%
#   select(CensusTract, TotalPopulation, CaliforniaCounty, LowBirthWeight, PM25, Poverty) %>%
#   mutate(LowBirthWeight = as.numeric(LowBirthWeight))
```

2 Question a:

What is the average concentration of PM2.5 across all census tracts in California?

```
mean_pm25 <- mean(data$PM25)
```

The average concentration of PM2.5 across all census tracts in California is 10.15 $\mu\text{g}/\text{m}^3$.

3 Question b:

What county has the highest level of poverty in California?

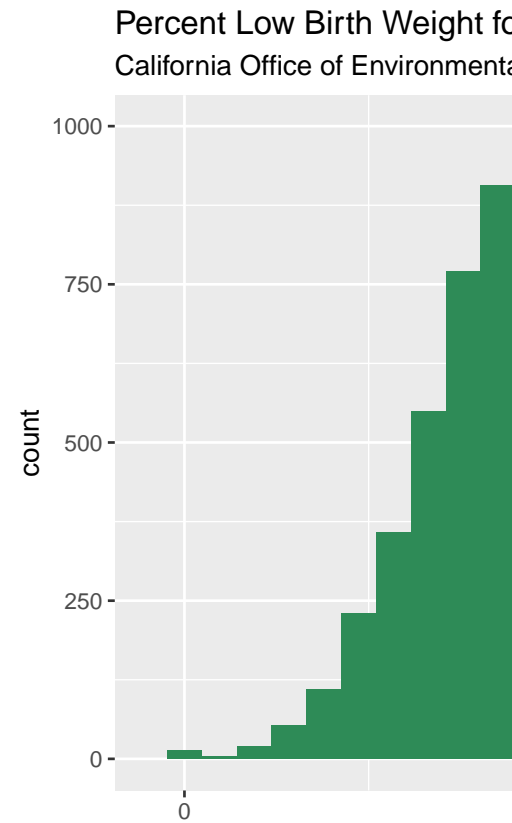
```
county_pov <- data %>%
  group_by(CaliforniaCounty) %>%
  summarise(mean_pov = mean(Poverty, na.rm = TRUE))

county_max_pov <- county_pov$CaliforniaCounty[which.max(county_pov$mean_pov)]
```

Tulare County has the highest level of poverty in California based on mean poverty for all census tracts in each county

Note: there are 75 census tracts with NA values for poverty

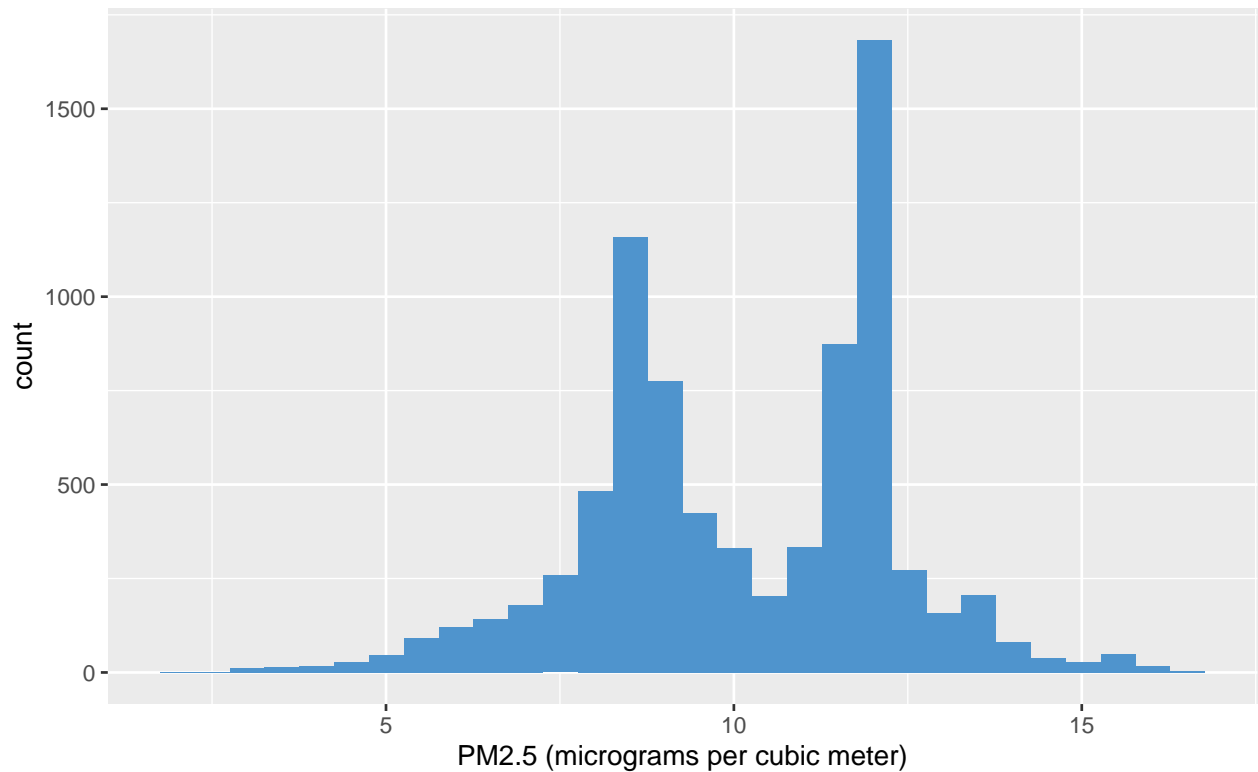
4 Question c:



Make a histogram depicting the distribution of percent low birth weight and PM2.5

Air Quality for California Census Tracts

California Office of Environmental Health Hazards Assessment (OEHHHA)



Data Source: CalEnviroScreen 4.0

xxx

Figure 1: MPG and vehicle weight

Figure 1 shows the expected negative relationship between vehicle weight and MPG.

5 Question d:

Estimate a OLS regression of `LowBirthWeight` on `PM25`. Report the estimated slope coefficient and its heteroskedasticity-robust standard error. Interpret the estimated slope coefficient. Is the effect of `PM2.5` on low birth weight statistically significant at the 5% level?

$$\text{percent low birth weight}_i = \beta_0 + \beta_1 \cdot \text{PM2.5} + \varepsilon_i$$

```
# xxx...maybe use glm instead of lm.fit
ols_model <- lm_robust(formula = LowBirthWeight ~ PM25, data = data)
summary(ols_model)
```

```
##
## Call:
## lm_robust(formula = LowBirthWeight ~ PM25, data = data)
##
## Standard error type: HC2
##
```

[illegible]

```
## (Intercept)
##      3.800988
```

```
##          PM25
## 0.1179305
```

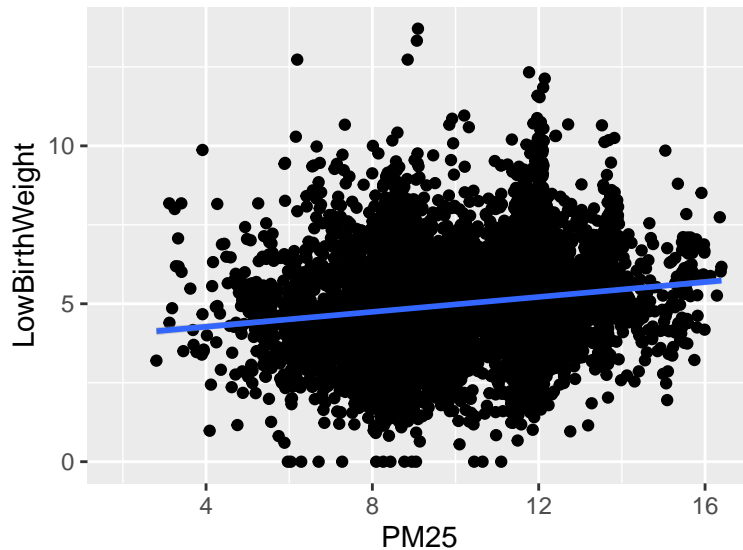
```
##          PM25
## 0.008402393
```

```
##          PM25
## 0.1014596
```

```
##          PM25
## 0.1344015
```

The estimated slope coefficient for this model is 0.118. The heteroskedasticity-robust standard error for this estimated slope coefficient is 0.008 xxx...

```
ols_model_plot <- ggplot(data = data, aes(x = PM25, y = LowBirthWeight)) +
  geom_point() +
  geom_smooth(method = lm)
ols_model_plot
```



xxx...read up on “heteroskedasticity-robust standard errors”

```
# xxx...experiment chunk

# data_rm_na <- data %>%
#   na.omit(LowBirthWeight) %>%
#   mutate(model_low_birth_wt = ols_model$fitted.values)

fit <- data.frame(ols_model$fitted.values)

fit2 <- data.frame(predict(ols_model, newdata = data, se.fit = TRUE, interval = "confidence"))
```

6 Question e:

Suppose a new air quality policy is expected to reduce PM2.5 concentration by 2 micrograms per cubic meters. Predict the new average value of low birth weight and derive its 95% confidence interval. Interpret the 95% confidence interval.

```
# data_reduced_PM <- data %>%
#   na.omit(LowBirthWeight) %>%
#   mutate(PM25_reduced = PM25 - 2) %>%
#   mutate(PM25_reduced = case_when(
#     PM25_reduced < 0 ~ 0,
#     PM25_reduced >= 0 ~ PM25_reduced)) %>%
#   mutate(test = slope_coef * PM25 + intercept) %>%
#   mutate(test_policy = slope_coef * PM25_reduced + intercept)

#mutate(model_low_birth_wt = ols_model$fitted.values)
```

```
policy_results <- data %>%
  mutate(PM25_reduced = PM25 - 2) %>%
  mutate(PM25_reduced = case_when(
    PM25_reduced < 0 ~ 0,
    PM25_reduced >= 0 ~ PM25_reduced)) %>%
  mutate(policy_low_birth_wt = slope_coef * PM25_reduced + intercept) %>%
  mutate(policy_ci_low = (ci_lower * PM25_reduced + intercept)) %>%
  mutate(policy_ci_high = (ci_upper * PM25_reduced + intercept))
```

```
mean_policy_low_birth_wt <- mean(policy_results$policy_low_birth_wt)
mean_policy_low_birth_wt
```

```
## [1] 4.762442
```

```
mean_policy_low_birth_wt_ci_low <- mean(policy_results$policy_ci_low)
mean_policy_low_birth_wt_ci_low
```

```
## [1] 4.896725
```

```
mean_policy_low_birth_wt_ci_high <- mean(policy_results$policy_ci_high)
mean_policy_low_birth_wt_ci_high
```

```
## [1] 4.896725
```

Following an air quality policy that reduces PM2.5 concentrations by 2 microns per cubic meters, the predicted new average value of percent low birth weight is 4.762‘

The 95% confidence interval for the predicted new average value of percent low birth weight ranges from 4.897 to 4.897. This means that there is a 95% chance that this interval includes the true average value of percent low birth weight based on the effect of the new policy on PM2.5 concentrations.

7 Question f:

Add the variable Poverty as an explanatory variable to the regression in (d). Interpret the estimated coefficient on Poverty. What happens to the estimated coefficient on PM2.5, compared to the regression in (d). Explain.

8 Question g:

From the regression in (f), test the null hypothesis that the effect of PM2.5 is equal to the effect of Poverty.

xxx...delete everything below here

9 Clean and plot data

The following code loads and cleans the data.

```

# Load data

data("mtcars")
raw_data <- mtcars

# Clean data

## Add model names as a column
## [this is just an example manipulation, I rarely assign rownames to a column]

clean_data <- tibble::rownames_to_column(raw_data, "model")

```

The code chunk below shows how to produce a scatter plot of MPG against weight.

```

# Plot 1

plot_1 <- ggplot(clean_data, aes(y=mpg, x = wt))+
  geom_point()+
  theme_cowplot(14)+
  labs(x = "Weight (1000 lbs)", y = "Miles per gallon")

```

Figure 1: MPG and vehicle weight

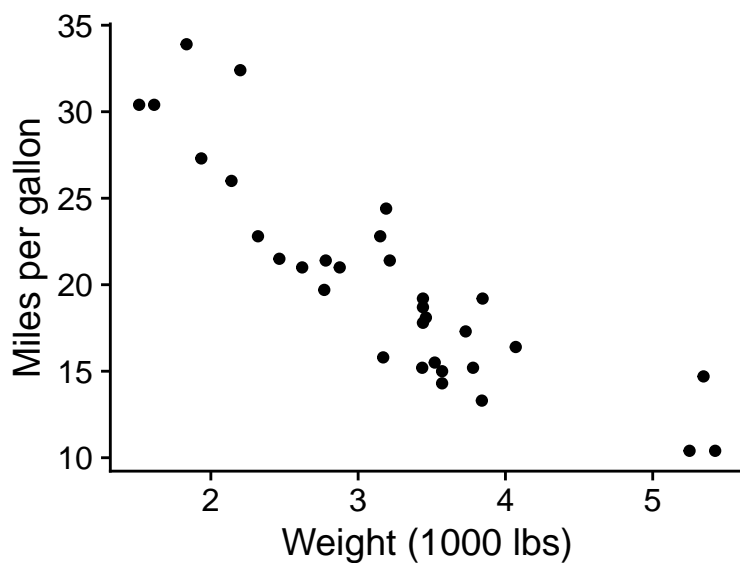


Figure 1 shows the expected negative relationship between vehicle weight and MPG.

10 Run and interpret regression models

In order to more formally analyze the relationship between MPG, vehicle weight, and cylinders we estimate the following regression:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i \quad (1)$$

where Y_i is MPG for vehicle model i , X_{1i} is the vehicle weight, X_{2i} is the number of cylinders in the engine, and u_i the regression error term. We will consider a regression including only vehicle weight, and a regression including vehicle weight and number of cylinders.

In R, we run the following code:

```
model_1 <- lm(mpg ~ wt, data=clean_data)
model_2 <- lm(mpg ~ wt + cyl, data=clean_data)
```

Table 1 shows the estimated coefficients from estimating equation (1).

Table 1: MPG and vehicle weight

	MPG	
	(1)	(2)
Weight (1000 lbs)	-5.344*** (0.559)	-3.191*** (0.757)
Cylinders		-1.508*** (0.415)
Observations	32	32
R ²	0.753	0.830
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01	

In model (1), the estimated β_1 coefficient implies that a 1000 pound increase in vehicle weight reduces miles per gallon by 5.3 miles. Adding the number of cylinders in model (2) reduces $\hat{\beta}_1$ from -5.3 to -3.2.