

# EDS241: Assignment 1

Lewis 1

2023-02-27

The data for this assignment are taken 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. For the assignment, you will need the following variables: CensusTract, TotalPopulation, LowBirthWeight (percent of census tract births with weight less than 2500g), PM25 (ambient concentrations of PM2.5 in the census tract, in micrograms per cubic meters), Poverty (percent of population in the census tract living below twice the federal poverty line), and LinguisticIsolation (percent of households in the census tract with limited English speaking).

**Read in and clean the data** The following code loads and cleans the data.

```
# read in data, clean the column names, and select  
# variables of interest  
env_health <- read_csv(here("Assignment_1/CES4.csv")) %>%  
  clean_names() %>%  
  select(census_tract, total_population, low_birth_weight,  
         pm2_5, poverty, linguistic_isolation)
```

```
# make sure that each row is a unique census tract  
length(unique(env_health$census_tract))
```

(a) What is the average concentration of PM2.5 across all census tracts in California? [1] 8035

```
# paste a statement that includes the mean PM2.5 across all  
# census tracts in California.  
print(paste("The average concentration of PM2.5 across all census tracts in California is",  
            round(mean(env_health$pm2_5), 3)))
```

[1] “The average concentration of PM2.5 across all census tracts in California is 10.153”

```
# convert low birth weight to be numeric  
env_health$low_birth_weight <- as.numeric(env_health$low_birth_weight)  
  
# plot distribution of low birth weight  
birth_wt_hist <- ggplot(data = env_health, aes(x = low_birth_weight)) +  
  geom_histogram(stat = "bin", color = "black", fill = "steelblue",
```

```
alpha = 0.6) + theme_minimal() + labs(x = "Percent of babies born classified as low weight",
title = "The distribution of low weight babies \n as a percentage of total births")
```

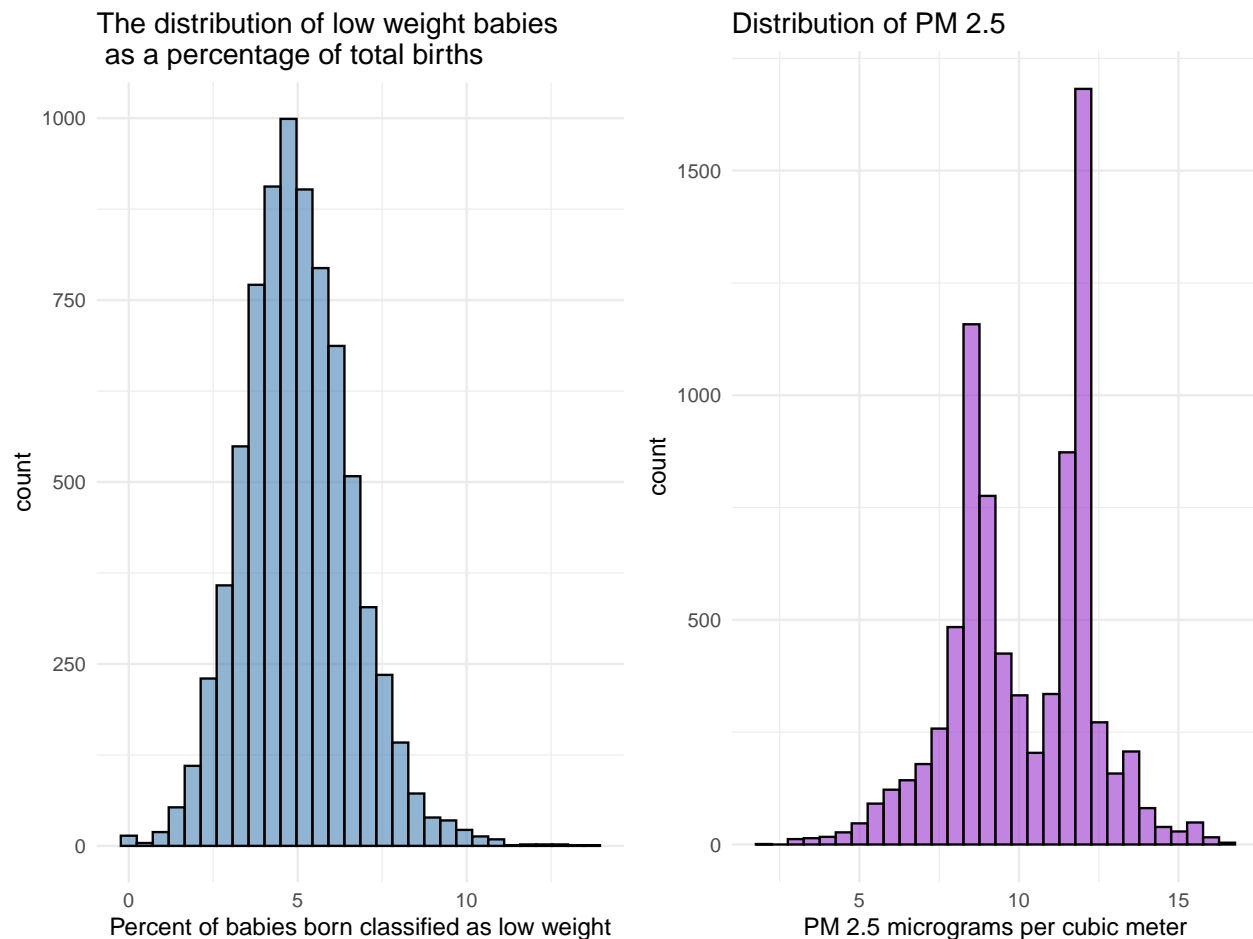
```
# plot distribution of pm2.5
```

```
pm_hist <- ggplot(data = env_health, aes(x = pm2_5)) + geom_histogram(stat = "bin",
color = "black", fill = "darkorchid", alpha = 0.6) + theme_minimal() +
labs(x = "PM 2.5 micrograms per cubic meter", title = "Distribution of PM 2.5")
```

(b) Make a histogram depicting the distribution of percent low birth weight and PM2.5.  
Figure 1: Low birth weight and PM 2.5 Distributions

```
# plot the histograms in one chart
```

```
gridExtra::grid.arrange(birth_wt_hist, pm_hist, ncol = 2)
```



(c) Estimate an OLS regression of LowBirthWeight on PM2.5. Report the estimated slope coefficient and its heteroskedasticity-robust standard error. Interpret the estimated slope coefficient. Is the effect of PM25 on LowBirthWeight statistically significant at the 5% level?  
[1] "The estimated slope coefficient for PM 2.5 is 0.11793 and the corresponding standard error for this value is 0.0084"

```
# summary table for the model
summary(birth_weight_pm25_mod)
```

```
##
```

[illegible]

```
# finding the mean pm2.5 value if it was 2 micrograms per
# cubic meter lower than before
pm2_5_lower <- mean(env_health$pm2_5 - 2)

# use predict() to find the average value and confidence
# interval values given this new mean PM2.5
predict(birth_weight_pm25_mod, newdata = list(pm2_5 = pm2_5_lower),
       se.fit = TRUE, interval = "confidence")
```

```
## $fit
##          fit          lwr          upr
## [1,] 4.76244 4.712526 4.812354
##
## $se.fit
##          1
## 0.02546284
```

Our model predicts that the new average percentage of babies born with a low birth weight would be 4.76%.

```
# adding poverty variable to the model
birth_wt_pm25_poverty <- lm_robust(low_birth_weight ~ pm2_5 +
  poverty, data = env_health, se_type = "HC1", alpha = 0.05)

# summarizing the new model
summary(birth_wt_pm25_poverty)
```

[illegible]

The estimated coefficient for PM 2.5 decreases in this new model from 0.1179 when just PM 2.5 was in the model to 0.05911 here. It appears as if poverty was causing omitted variable bias in the first model. Poverty is positively associated with low birth weights. Based on how the coefficient for PM2.5 changed, poverty must be positively associated with PM2.5 as well.

```
# adding an indicator variable for linguistic isolation
env_full <- env_health %>%
  mutate(linguistic_isolation = as.numeric(linguistic_isolation)) %>%
  mutate(ling_iso_over_6point9 = ifelse(linguistic_isolation >
    6.9, 1, 0))
```

