

# Final Report

due November 16, 2021 by 11:59 PM

Kaitlyn Lewars and Katie Meehl: The Exposure Experience

November 16, 2021

```
library(tidyverse)
library(readr)
library(scales)
library(tidymodels)
library(knitr)
library(infer)
install.packages("janitor")
library(janitor)

#Total_Air_Pollution_Death_Rate <- read_csv("~/Project Proposal/project-the-exposure-experience/data/To
#Household_Air_Pollution_Total_Deaths <- read_csv("~/Project #Proposal/project-the-exposure-experience/
#Household_Air_Pollution_Rate <- read_csv("~/Project Proposal/project-the-exposure-experience/data/Hous
#Ambient_Air_Pollution_Rate <- read_csv("~/Project Proposal/project-the-exposure-experience/data/WHO Am
#Ambient_Air_Pollution_Total_Deaths <- read_csv("~/Project Proposal/project-the-exposure-experience/dat
#Particulate_Ambient_Concentration <- read_csv("~/Project Proposal/project-the-exposure-experience/data,

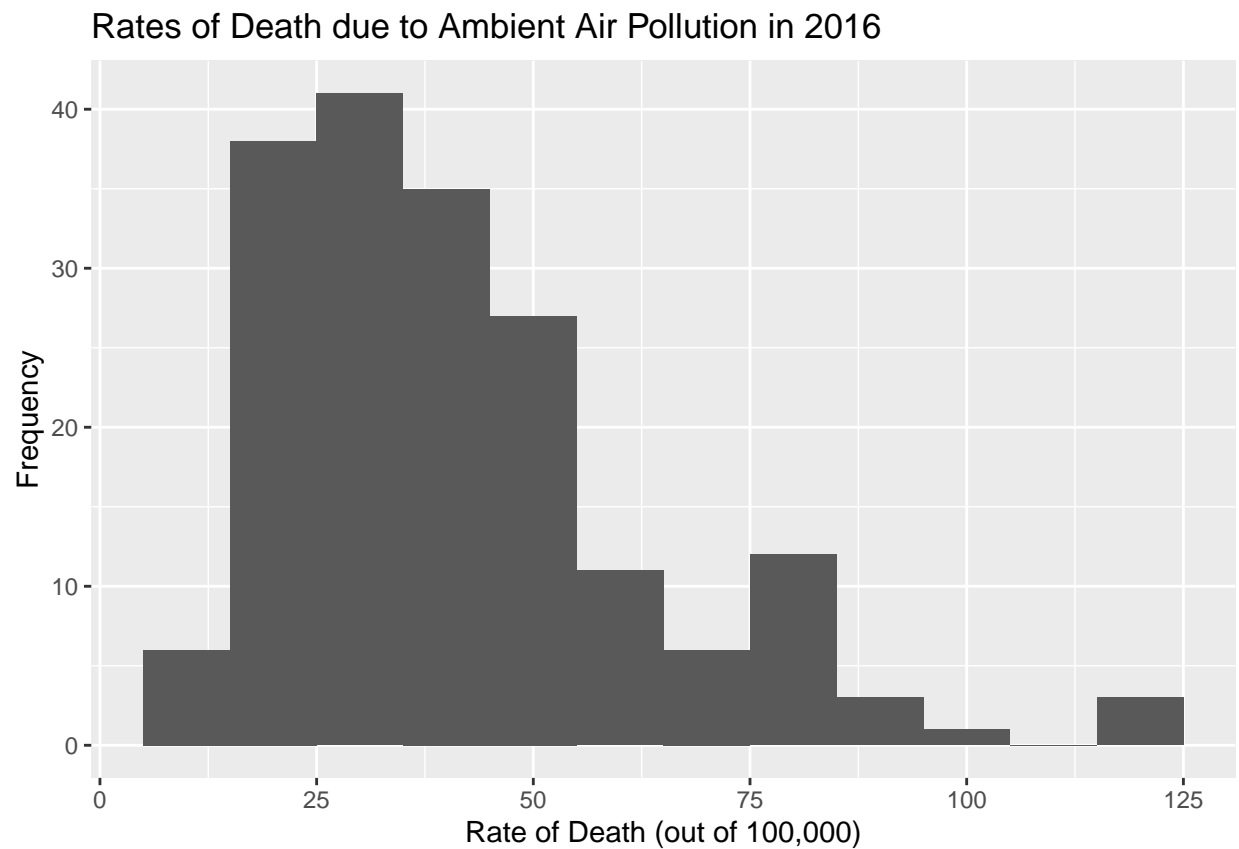
Total_Air_Pollution_Death_Rate <- read_csv("~/R/Project Proposal/data/Total Air Pollution Death Rate.csv")
Household_Air_Pollution_Total_Deaths <- read_csv("~/R/Project Proposal/data/Household Air Pollution Tot
Household_Air_Pollution_Rate <- read_csv("~/R/Project Proposal/data/Household Air Pollution Rate.csv")
Ambient_Air_Pollution_Rate <- read_csv("~/R/Project Proposal/data/WHO Ambient Air Pollution Rate.csv")
Ambient_Air_Pollution_Total_Deaths <- read_csv("~/R/Project Proposal/data/Ambient Air Pollution Total D
Particulate_Ambient_Concentration <- read_csv("~/R/Project Proposal/data/Ambient Particulate Concentrat
```

Climate change has been a recurring topic in the news in recent years as it becomes a more pressing problem. One of the important factors of climate change is air pollution. In 2017, air pollution was the 4th leading cause of mortality and the 5th leading cause of morbidity worldwide. As air pollution is a leading cause of morbidity and mortality, we thought it would be important to explore a data set investigating this problem.

In general we would like to investigate air pollution as a cause of mortality globally. There are several different types of air pollution, but we will look at household pollution and ambient matter pollution. With these two variables we will compare them to see which air pollution is the most fatal. We would also like to look into the trend of air pollution over the last 27 years. Lastly we would like to compare air pollution as a risk factor to other common risk factors. We downloaded this data from kaggle. There are several variables in this data including year, country, deaths by each type of air pollution, and deaths by other risk factors.

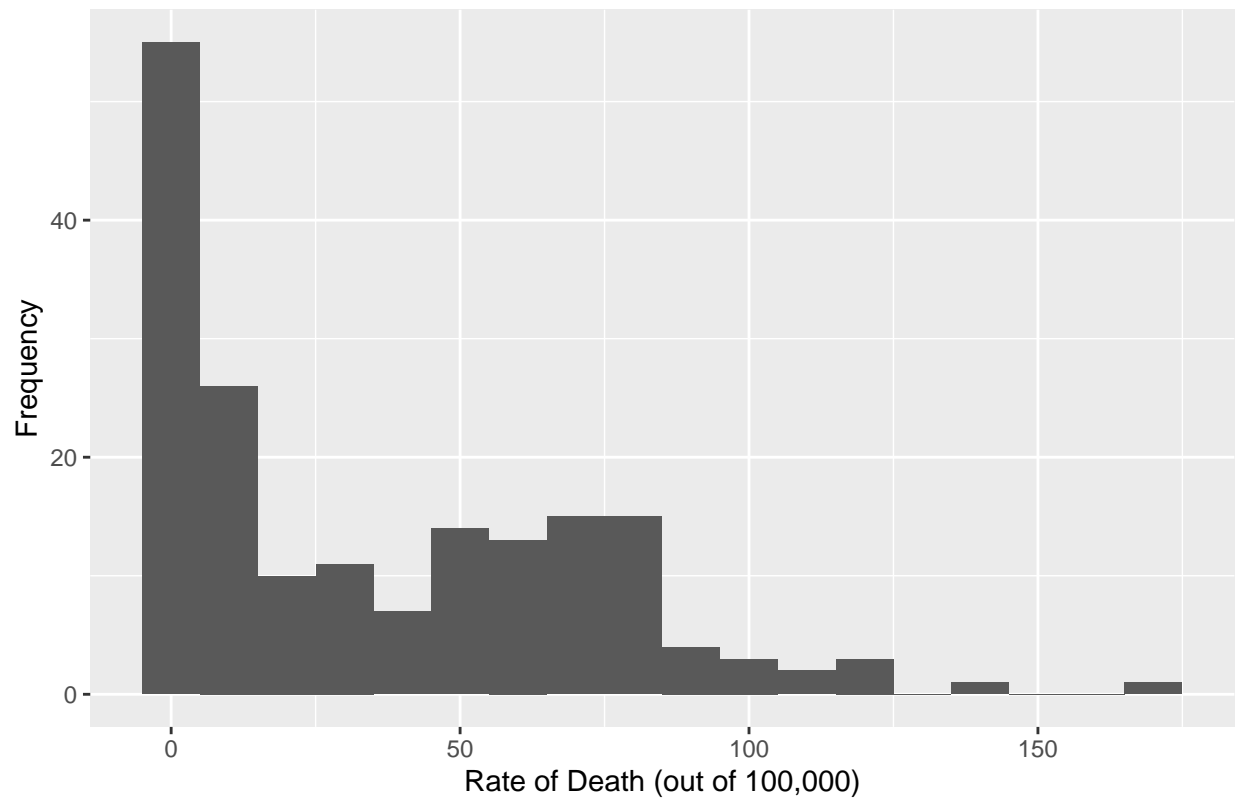
The data collection is a bit complicated. In order to estimate deaths caused by pollution they use “mathematical functions, derived from epidemiological studies from countries around the world, that relate different levels of exposure to the increased risk of death or disability from each cause, by age and sex, where applicable, estimates of population exposure to PM2.5, ozone, and household air pollution, country-specific data on underlying rates of disease and death for each pollution-linked disease, and a comprehensive set of population data, adjusted to match the UN2015 Population Prospectus and obtained from the Gridded Population of the World (GPW) database for each country” (<https://www.stateofglobalair.org/data/estimate-burden>).

```
Ambient_Air_Pollution_Rate %>%
  filter(Dim2 == "Total", Dim1 == "Both sexes") %>%
  ggplot(aes(x = FactValueNumeric)) +
  geom_histogram(binwidth = 10) +
  scale_x_continuous(labels = label_comma()) +
  labs(x = "Rate of Death (out of 100,000)",
       y = "Frequency",
       title = "Rates of Death due to Ambient Air Pollution in 2016")
```

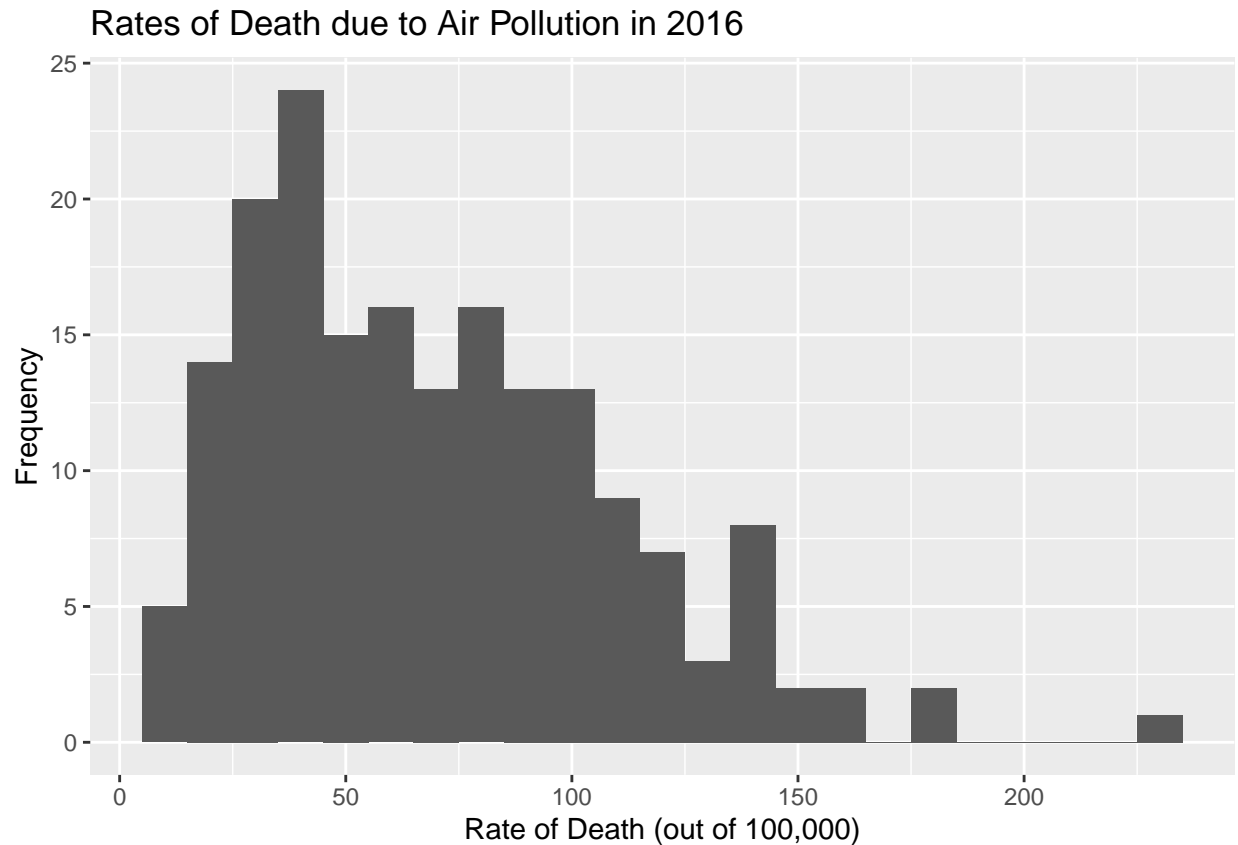


```
Household_Air_Pollution_Rate %>%
  filter(Dim2 == "Total", Dim1 == "Both sexes") %>%
  ggplot(aes(x = FactValueNumeric)) +
  geom_histogram(binwidth = 10) +
  scale_x_continuous(labels = label_comma()) +
  labs(x = "Rate of Death (out of 100,000)",
       y = "Frequency",
       title = "Rates of Death due to Household Air Pollution in 2016")
```

## Rates of Death due to Household Air Pollution in 2016



```
Total_Air_Pollution_Death_Rate %>%  
  filter(Dim2 == "Total", Dim1 == "Both sexes") %>%  
  ggplot(aes(x = FactValueNumeric)) +  
  geom_histogram(binwidth = 10) +  
  scale_x_continuous(labels = label_comma()) +  
  labs(x = "Rate of Death (out of 100,000)",  
       y = "Frequency",  
       title = "Rates of Death due to Air Pollution in 2016")
```



The four countries of interest are Democratic People's Republic of Korea, Georgia, Chad, and Bosnia and Herzegovina.

[Insert graph for N. Korea, Georgia, Chad, Bos. and Herz.]

```
# These are the possible graphs I've made starter code for. I still need to work on them to get them to work

# TRYING TO COMPARE RATES BELOW - Need to add x labels for countries
#Total_Air_Pollution_Death_Rate %>%
#filter(Dim2 == "Total", Dim1 == "Both sexes") %>%
#filter(Location %in% c("Georgia", "Chad", "Democratic People's Republic of Korea", "Bosnia and Herzegovina"))
#ggplot(aes(x = FactValueNumeric)) +
#geom_histogram() +
#labs(x = "Rate of Death (out of 100,000)",
#      y = "Frequency",
#      title = "Rates of Death due to Air Pollution in 2016")

#COMPARE COUNTRY WITHIN REGION - Need to make this work [not showing anything]
#Total_Air_Pollution_Death_Rate %>%
#filter(ParentLocation=="Europe") %>%
#filter(Dim2=="Total", Dim1=="Both Sexes") %>%
##ggplot(aes(x=FactValueNumeric)) +
#geom_histogram() +
#labs(x="Rate of Death (out of 100,000)", y="Frequency", title="Comparative Rates of Death due to Air Pollution in 2016")

# COMPARE BETWEEN SEXES
```

```

#Total_Air_Pollution_Death_Rate %>%
  #filter(Location %in% c("Georgia")) %>%
  #filter(Dim2 == "Total") %>%
  #ggplot(aes(x = FactValueNumeric, fill=Dim1)) +
  # geom_bar() +
  #scale_x_continuous(labels = label_comma()) +
  #labs(x = "Rate of Death (out of 100,000)",
        #y = "Frequency",
        #title = "Rates of Death due to Air Pollution in 2016")

# COMPARE ILLNESS CAUSED - Need to make this work [Not showing anything]
#Total_Air_Pollution_Death_Rate %>%
  #filter(Location %in% c("Georgia")) %>%
  #filter(Dim1=="Both Sexes") %>%
  #ggplot(aes(x = FactValueNumeric, fill=Dim2)) +
  #geom_bar() +
  #scale_x_continuous(labels = label_comma()) +
  #labs(x = "Rate of Death (out of 100,000)", y = "Frequency", title = "Rates of Death due to Air Pollution in 2016")

joinedambient1 <- Ambient_Air_Pollution_Rate %>%
  rename(AmbientDeathRate = FactValueNumeric) %>%
  rename(Sex = Dim1) %>%
  rename(CauseofDeath = Dim2) %>%
  select(Location, Sex, CauseofDeath, AmbientDeathRate) %>%
  left_join(Ambient_Air_Pollution_Total_Deaths) %>%
  select(Location, Sex, CauseofDeath, AmbientDeathRate, FactValueNumeric) %>%
  filter(CauseofDeath == "Total") %>%
  mutate(totalpopulation = (100000*FactValueNumeric)/AmbientDeathRate) %>%
  rename(Totaldeathsambient = FactValueNumeric)

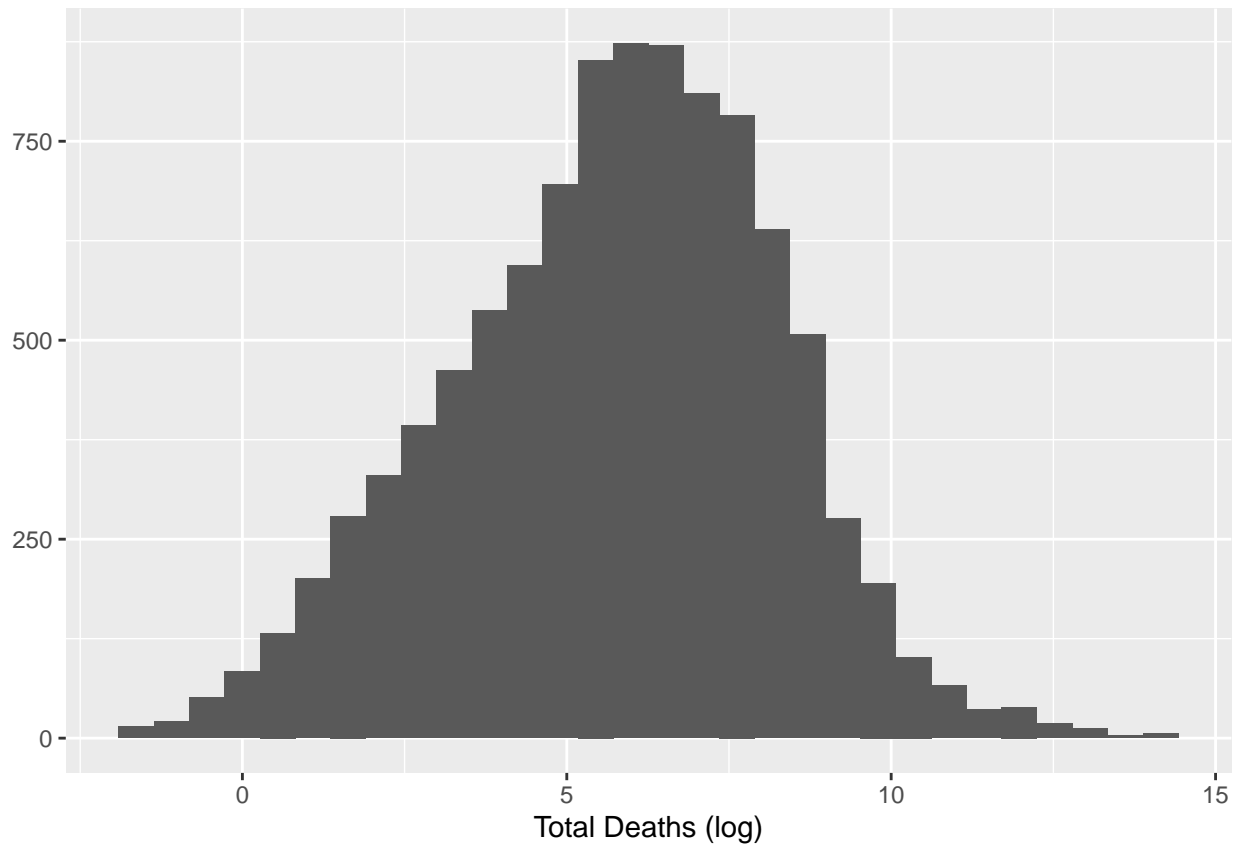
ambient <- Particulate_Ambient_Concentration %>%
  rename(AmbientAirConcentration = FactValueNumeric) %>%
  select(Location, AmbientAirConcentration) %>%
  left_join(joinedambient1, by = "Location") %>%
  select(Location, Sex, CauseofDeath, AmbientDeathRate, Totaldeathsambient, AmbientAirConcentration, totalpopulation) %>%
  filter(CauseofDeath == "Total", Sex == "Both sexes")

[Evaluate, explain, need to edit scale for both axes]

ggplot(data = ambient, aes(x = log(Totaldeathsambient))) +
  geom_histogram() +
  labs(x = "Total Deaths (log)", y = NULL)

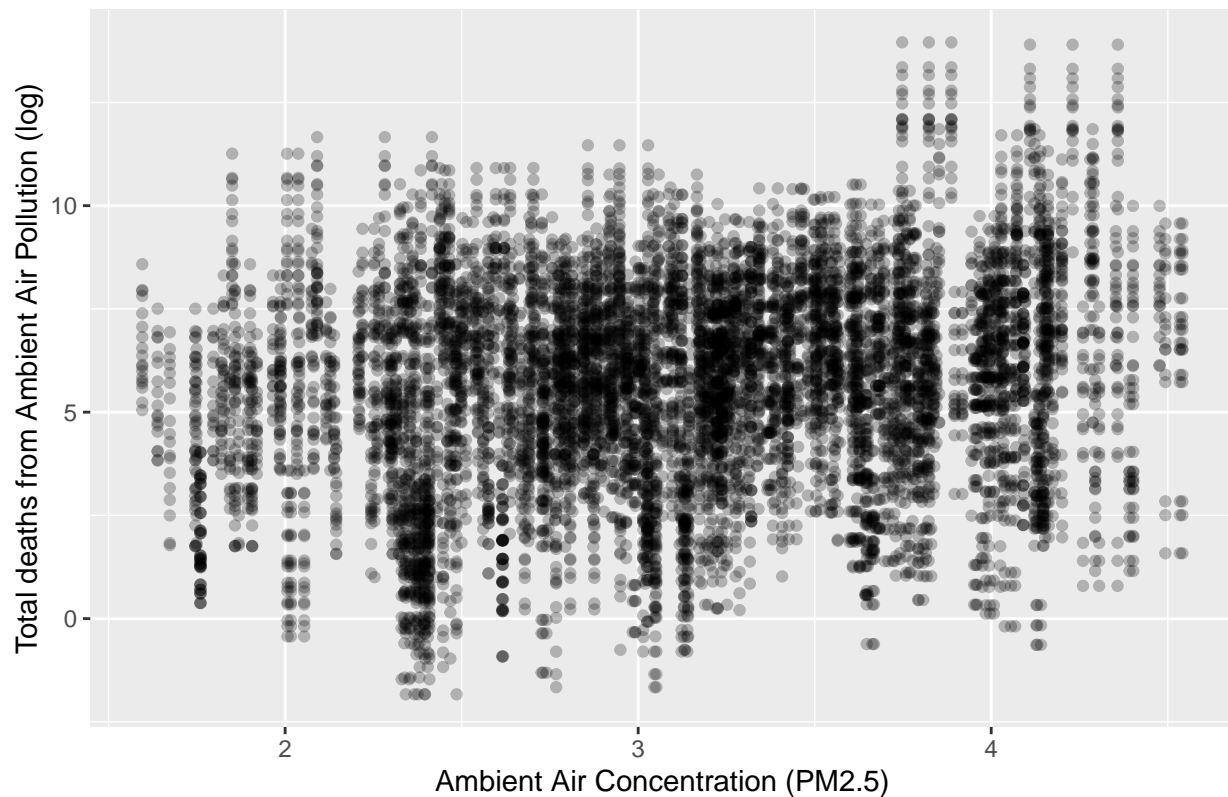
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

```



```
ggplot(data = ambient, aes(x = log(AmbientAirConcentration), y = log(Totaldeathsambient))) +  
  geom_point(alpha = 0.25) +  
  labs(title = "Deaths as a function of Ambient Air Concentration",,  
    x = "Ambient Air Concentration (PM2.5)",  
    y = "Total deaths from Ambient Air Pollution (log)")
```

## Deaths as a function of Ambient Air Concentration

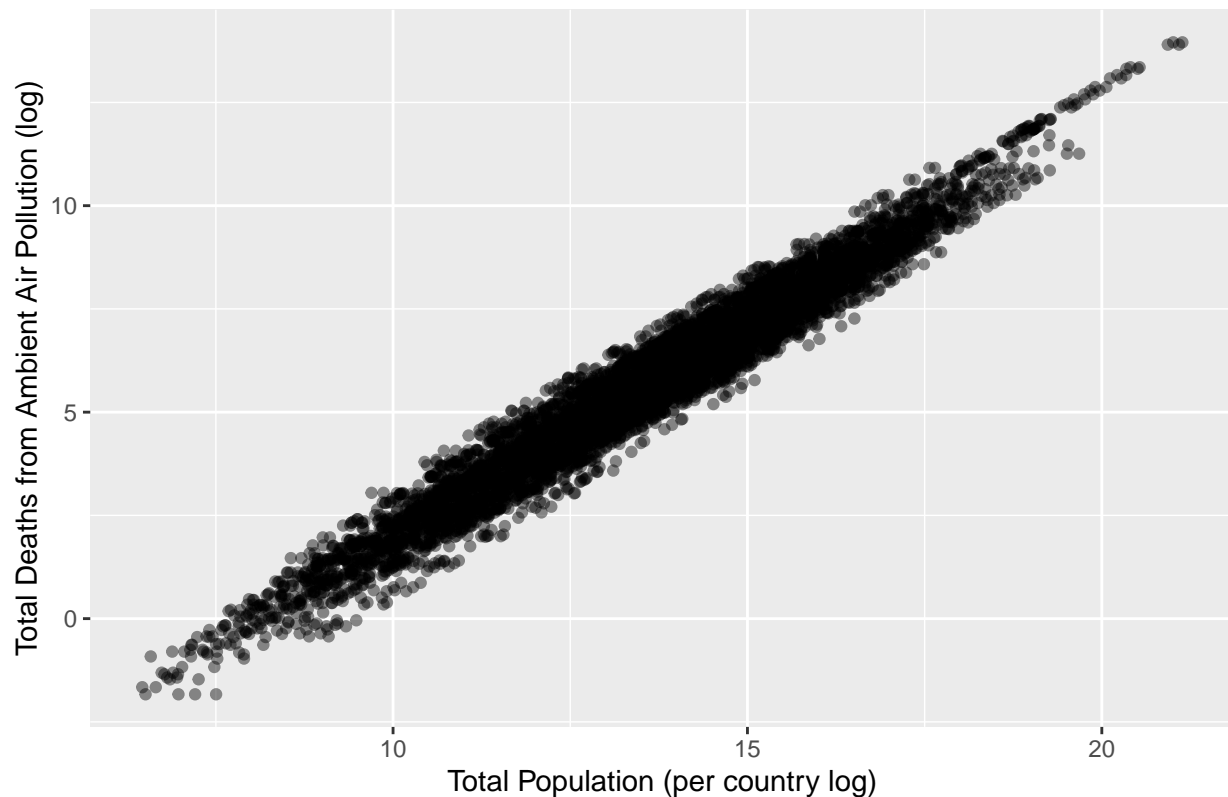


```
ambient2 <- Particulate_Ambient_Concentration %>%
  rename(AmbientAirConcentration = FactValueNumeric) %>%
  rename(Residencetype = Dim1) %>%
  select(Location, AmbientAirConcentration, Residencetype) %>%
  left_join(joinedambient1, by = "Location") %>%
  select(Location, Sex, CauseofDeath, AmbientDeathRate, Totaldeathsambient, AmbientAirConcentration, to
  filter(CauseofDeath == "Total", !Sex %in% c("Both sexes"), !Residencetype %in% c("Total"))
```

[Evaluate, explain]

```
ggplot(data = ambient2, aes(x = log(totalpopulation), y = log(Totaldeathsambient))) +
  geom_point(alpha = 0.25) +
  labs(title = "Deaths as a function of Population", ,
       x = "Total Population (per country log)",
       y = "Total Deaths from Ambient Air Pollution (log)")
```

## Deaths as a function of Population



For the next 4 sections we created two different models for a linear regression. The first one we did just based off the outdoor air concentration. The equation we got was  $\text{deaths} = -2601.0 + 274.1(\text{AirConcentration})$ . The air concentration is in PM2.5.

```
death_ambientairpol <- linear_reg() %>%
  set_engine("lm") %>%
  fit(log(Totaldeathsambient) ~ log(AmbientAirConcentration), data = ambient)
tidy(death_ambientairpol, conf.int=TRUE, exponentiate=TRUE)
```

```
## # A tibble: 2 x 7
```

##	term	estimate	std.error	statistic	p.value	conf.low	conf.high
##	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
## 1	(Intercept)	3.38	0.120	28.2	7.31e-169	3.15	3.62
## 2	log(AmbientAirConce~	0.746	0.0374	20.0	6.29e- 87	0.673	0.819

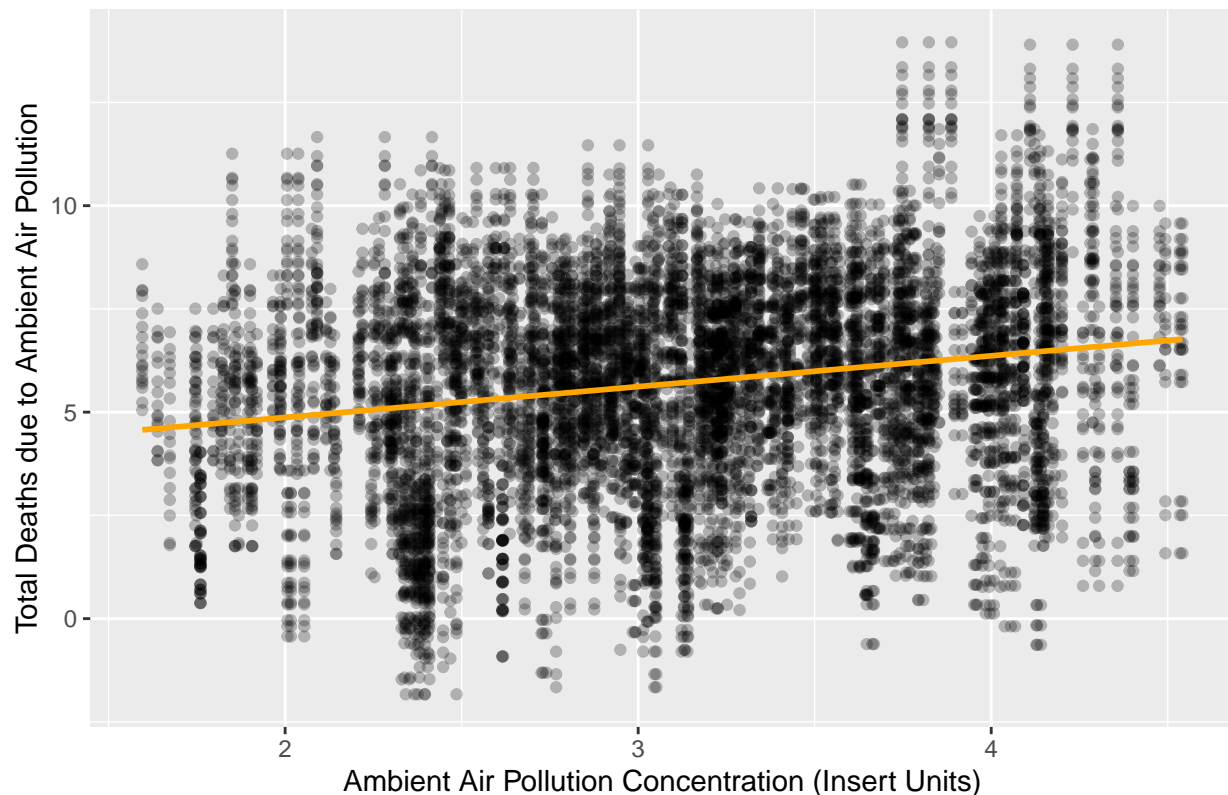
$\text{deaths} = 4.979.0 + 0.0261(\text{AirConcentration})$

[evaluate what this means]

```
ggplot(data = ambient, aes(x = log(AmbientAirConcentration), y = log(Totaldeathsambient))) +
  geom_point(alpha = 0.25) +
  geom_smooth(method = "lm", se = FALSE, color = "orange") +
  labs(title = "Deaths per Ambient Air Pollution Concentration",
       x = "Ambient Air Pollution Concentration (Insert Units)",
       y = "Total Deaths due to Ambient Air Pollution")
```



## Deaths per Ambient Air Pollution Concentration



[Add narrative about the graph - most deaths are concentrated at mid concentrations - lethal dose, what are the outliers?]

```
glance(death_ambientairpol)$r.squared
```

```
## [1] 0.03876003
```

```
death_ambientairpol_totalpopulation <- linear_reg() %>%
  set_engine("lm") %>%
  fit(log(Totaldeathsambient) ~ log(AmbientAirConcentration) + Sex + Residencetype, data = ambient2)

tidy(death_ambientairpol_totalpopulation, conf.int=TRUE, exponentiate=TRUE)
```

```
## # A tibble: 4 x 7
```

term	estimate	std.error	statistic	p.value	conf.low	conf.high
<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1 (Intercept)	3.45e+ 0	0.106	3.26e+ 1	7.91e-224	3.24	3.66
2 log(AmbientAirConc~	7.35e- 1	0.0322	2.28e+ 1	3.09e-113	0.672	0.798
3 SexMale	-1.30e-17	0.0426	-3.05e-16	1 e+ 0	-0.0835	0.0835
4 ResidencetypeUrban	-6.70e- 2	0.0427	-1.57e+ 0	1.17e- 1	-0.151	0.0167

```
glance(death_ambientairpol_totalpopulation)$r.squared
```

```
## [1] 0.03808491
```

Here we created a binomial regression model to predict the likelihood of a person dying based on the level of ambient air concentration. (will add explanation later)

```
Alive <- ambient2$totalpopulation - ambient2$Totaldeathsambient
```

```
modelagg1<-glm(cbind(round(Totaldeathsambient), round(Alive)) ~ AmbientAirConcentration, data=ambient2,
summary(modelagg1)
```

```
##
## Call:
## glm(formula = cbind(round(Totaldeathsambient), round(Alive)) ~
##     AmbientAirConcentration, family = binomial, data = ambient2)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -161.24  -12.68   -3.56    0.39   430.08
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      -8.005e+00  2.763e-04  -28975  <2e-16 ***
## AmbientAirConcentration  1.237e-02  5.419e-06   2282  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 14254994  on 13175  degrees of freedom
## Residual deviance:  9023120  on 13174  degrees of freedom
## AIC: 9122756
##
## Number of Fisher Scoring iterations: 4
```

Here we created a binomial regression model to predict the likelihood of a person dying based on the level of ambient air concentration, gender, and where they live. (will add explanation later)

```
Alive <- ambient2$totalpopulation - ambient2$Totaldeathsambient
```

```
modelagg2<-glm(cbind(round(Totaldeathsambient), round(Alive)) ~ AmbientAirConcentration + Sex + Residencetype,
summary(modelagg2)
```

```
##
## Call:
## glm(formula = cbind(round(Totaldeathsambient), round(Alive)) ~
##     AmbientAirConcentration + Sex + Residencetype, family = binomial,
##     data = ambient2)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -143.52  -12.84   -3.78    0.17   347.88
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      -8.058e+00  3.130e-04 -25742.6  <2e-16 ***
## AmbientAirConcentration  1.261e-02  5.509e-06   2288.9  <2e-16 ***
## SexMale              1.674e-01  2.427e-04    689.6  <2e-16 ***
## ResidencetypeUrban    -7.582e-02  2.458e-04   -308.4  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 14254994  on 13175  degrees of freedom
## Residual deviance:  8452356  on 13172  degrees of freedom
## AIC: 8551996
##
## Number of Fisher Scoring iterations: 4

Finally we conducted a t-test to see if the impact on mortality is different depending on the type of air
pollution (household vs. ambient)

t.test(Household_Air_Pollution_Rate$FactValueNumeric,Ambient_Air_Pollution_Rate$FactValueNumeric, var.e

##
## Welch Two Sample t-test
##
## data: Household_Air_Pollution_Rate$FactValueNumeric and Ambient_Air_Pollution_Rate$FactValueNumeric
## t = -3.9975, df = 6328.8, p-value = 6.474e-05
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -2.8883144 -0.9876038
## sample estimates:
## mean of x mean of y
##  11.87733  13.81529
```