

# The Air We Breathe: Air Quality and Health Outcomes in Kentucky

Jianing Gong<sup>a</sup>, Madhvi Malhotra<sup>b</sup>, Maggie Sullivan<sup>c</sup>, and Sanha Tahir<sup>d</sup>

<sup>a</sup>All authors are M.S. candidates in Georgetown University's Data Science for Public Policy Program

This manuscript was compiled on December 6, 2022

1 Much research has been conducted on the relationship between air pollution and health outcomes. For our final project, we dived deeper  
2 into this relationship in a single state in the United States: Kentucky. With 2019 data from the U.S. Centers for Disease Control and Preven-  
3 tion (CDC) and the U.S. Environmental Protection Agency (EPA), we conducted an exploratory analysis with data visualization and linear  
4 regression on overall, respiratory, and mental health outcomes. Due to a very limited sample size (potential political reasons discussed), our  
5 results did not show a statistically significant relationship, but we provided recommendations for further research.

Air Quality | Health Outcomes | Sulfur Dioxide | Policy

1 **A**ir!, an invisible friend we need for survival and breathe it every second of our everyday lives. However, we are living  
2 in times where we are also breathing various invisible particles with it. Sometimes this can be harmful and would  
3 result in many unknown diseases. As per the survey by WHO (2021)<sup>\*</sup>, 58 % of premature deaths related to outdoor  
4 air pollution were due to ischemic heart disease and stroke, while 18 % were due to chronic obstructive pulmonary  
5 disease and acute lower respiratory infection, respectively, and 6 % cause of death was due to lung cancer. Not only  
6 that, according to the American Psychological Association, polluted air is harmful to people's brains as well. When  
7 children are exposed to highly polluted air, their cognitive abilities may be damaged. Furthermore, there are studies  
8 that have concluded that exposure to certain PM can lead to oxidative stress-mediated inflammation in the brain.<sup>†</sup>. Other studies have found supporting evidence that these pollutants are connected to negative mental health outcomes.<sup>‡</sup>

10 Common air pollutants with known health impacts were first regulated as "Criteria Pollutants" by the 1970  
11 Clean Air Act, which established health-based National Ambient Air Quality Standards (NAAQS). The six criteria  
12 pollutants are carbon monoxide(CO), ground-level ozone, lead, nitrogen dioxide(NO<sub>2</sub>), particulate matter(PM), and  
13 sulfur dioxide(SO<sub>2</sub>).

14 Out of these, Sulfur dioxide(SO<sub>2</sub>) has a particularly significant presence in Kentucky, which is one the largest  
15 coal-producing states in the United States. In fact, in 2016, Kentucky ranked as the fourth-highest coal producer in  
16 the United States, producing 42.9 million tons<sup>§</sup>. SO<sub>2</sub> is a direct product of the combustion of coal and the smelting  
17 of sulfur-containing ores. It is a colorless, soluble gas with a characteristic pungent smell that forms sulphuric acid  
18 when combined with water. SO<sub>2</sub> can also be directly credited for the creation of respiratory symptoms in healthy  
19 patients and those with underlying lung disease. Sulfur dioxide can cause respiratory problems such as bronchitis and  
20 can irritate your nose, throat, and lungs. It may cause coughing, wheezing, phlegm, and asthma attacks. The effects  
21 are worse when you are exercising. Sulfur dioxide has been linked to cardiovascular disease. Groups that are most  
22 sensitive to sulfur dioxide include children, adults with lung diseases, and asthmatics.

23 However, studies have not demonstrated dose-dependent health risk responses to increased exposure to these  
24 pollutants, except at high concentrations. Furthermore, many studies examining the effects of ambient-level exposure  
25 to NO<sub>2</sub>, SO<sub>2</sub>, and CO have failed to find associations with adverse health outcomes. This research aimed to find out  
26 the correlation between air quality and health status and is based on the 2019 data collected in Kentucky. The re-  
27 search analyzes 10 counties and over 13,000 data, running a regression model and showing results in a map visualization.

\* Ambient (outdoor) air pollution, WHO,[https://www.who.int/news-room/fact-sheets/detail/ambient-\(outdoor\)-air-quality-and-health](https://www.who.int/news-room/fact-sheets/detail/ambient-(outdoor)-air-quality-and-health)

† Particulate matter, oxidative stress, and neurotoxicity, Sheba MohanKumara, Arezoo Campbell, Michelle Block, Bellina Veronesi, 2007

‡ Growing evidence for the impact of air pollution on depression, Ali, Naureen Akber and Khoja, Adeel, 2019

§ [https://eecc.ky.gov/Energy/Coal%20Facts%20%20Annual%20Editions/Kentucky%20Coal%20Facts%20-%202017th%20Edition%20\(2017\).pdf](https://eecc.ky.gov/Energy/Coal%20Facts%20%20Annual%20Editions/Kentucky%20Coal%20Facts%20-%202017th%20Edition%20(2017).pdf)

32 **Data**

33 **Data Sources.** Our data consisted of three main datasets:

- 34 1. Annual\_aqi\_by\_county\_2019.csv (referred to as “AQI” through the remainder of this paper),  
35 2. Daily\_42401\_2019\_SO2.csv (referred to as “Sulfur”), and  
36 3. PLACES\_Local\_Data\_for\_Better\_Health\_Census\_Tract\_Data\_2021\_release.csv (referred to as “Health”)

37 We downloaded the AQI and the Sulfur data from the U.S. Environmental Protection Agency (EPA) and the  
38 Health data from the Centers for Disease Control and Prevention (CDC). These two U.S. Government agencies  
39 regularly collect and publish these official datasets as part of their mandate. As a result, this data is freely available,  
40 relatively free of errors, and consistently available across many years.

41 We subset all datasets to 2019 data for the state of Kentucky. The 2021 Health data release is the most up-to-date  
42 data available on the CDC’s website and contains data collected in 2018 and 2019. More recent data AQI and Sulfur  
43 data are available on the EPA website; however, we downloaded 2019 data to be able to compare to the Health data in  
44 the same year. According to the CDC<sup>¶</sup> website, this Health data is based on “the Behavioral Risk Factor Surveillance  
45 System (BRFSS), Census 2010 population counts, annual Census county population estimates, and the Census  
46 American Community Survey (ACS) estimates.” According to the EPA website<sup>||</sup>, the AQI and Sulfur datasets contain  
47 both daily and annual data because of the Clean Air Act, which “requires that state, local, and tribal air pollution  
48 control agencies monitor the air for ambient levels of certain pollutants. In addition to the required monitoring, many  
49 agencies perform additional and/or voluntary monitoring of substances and meteorological parameters.”

51

52 **AQI.** For the Air Quality data, we have used the Annual and Daily AQI data. The annual AQI data has the pollutant  
53 summary; all metrics (mean, maxes, percentiles) are based on the daily maximum 8-hour average, not the individual  
54 sample values. We have data for days with Carbon Dioxide (C O), Nitrogen (NO2), Ozone, PM2.5, and PM 10. The  
55 daily summary files show the AQI by county and have one record per day per county. For our visualization, we have  
56 calculated the monthly mean of the AQI data.

57 **Sulfur.** For the Sulfur data, we have used the Daily AQI Summary data. The file contains SO2 readings collected  
58 from different sites and associated metadata. The readings themselves include units of measure, and arithmetic  
59 mean/maximum/minimum values of observed counts. The associated metadata includes site numbers, site geolocations,  
60 method collection and local dates. These latter meta details were critical in helping us identify feasible methods to  
61 rollup data for appropriate analysis. The daily file shows the AQI by county, and may have multiple records based on  
62 different sample collection methods per day, per site.

63 **Health.** For the Health data, in line with our research questions, we chose to focus on respiratory and mental health-  
64 related health measures for each of the four health categories in the dataset. Our final list of health categories and  
65 the corresponding measures (and measure ids) are below:

66

- 67 • Health Outcomes: current asthma (CASTHMA), Chronic obstructive pulmonary disease among adults aged  
68 >=18 years’ (COPD), Cancer (excluding skin cancer) among adults aged >=18 years (CANCER), and Depression  
69 among adults aged >=18 years (DEPRESSION)
- 70 • Prevention: Visits to doctor for a routine checkup within the past year among adults aged (CHECKUP),  
71 Current lack of health insurance among adults aged 18-64 years (ACCESS2), and Taking medicine for high  
72 blood pressure control among adults aged >=18 years with high blood pressure (BPMED)
- 73 • Risk Behaviors: Current smoking among adults aged >=18 years (CSMOKING) and No leisure-time physical  
74 activity among adults aged >=18 years (LPA)
- 75 • Health Status: Physical health not good for >=14 days among adults aged >=18 years (PHLTH), Mental  
76 health not good for >=14 days among adults aged >=18 years (MHLTH), Fair or poor self-rated health status  
77 among adults aged >=18 years’ (GHLTH)

78

<sup>¶</sup>Methodology, CDC, <https://www.cdc.gov/places/methodology/index.html>, 2022

<sup>||</sup>About AQS Data, EPA, [https://aqs.epa.gov/aqsweb/documents/about\\_aqs\\_data.html](https://aqs.epa.gov/aqsweb/documents/about_aqs_data.html), 2022

79 We excluded all other health measures that were not relevant to our current study, such as “Arthritis among adults  
80 aged  $\geq 18$  years” and “Chronic kidney disease among adults aged  $\geq 18$  years.” All 120 Kentucky counties were  
81 represented in the data, with each county containing anywhere from 1 to 190 prevalence values for the focus health  
82 measure. For our analysis, we took each county’s mean and median prevalence values for each of our focus measures,  
83 so there was one corresponding crude prevalence value for each county. In addition to the county, the dataset included  
84 location data by CountyFIPS, Geolocation, and LocationId.

85  
86 We focused on the county name and the latitude and longitude values in the Geolocation column. For our  
87 final analytic sample, we dropped columns with duplicate information and other less-important to only keep the  
88 following column values: ‘index’, ‘Year’, ‘CountyName’, ‘CountyFIPS’, ‘Category’, ‘Measure’, ‘Data\_Value\_Unit’,  
89 ‘Data\_Value\_Type’, ‘Data\_Value’, ‘TotalPopulation’, ‘Geolocation’, ‘LocationID’, ‘MeasureId’. See the image below  
90 for more information about the original dataset post subset to 2019 and Kentucky.

91 **Limitations.** Our data sources and research approach do have limitations that may affect the results and conclusion.  
92 For example, the health data only includes results from adults (age 18+). Therefore, we cannot accurately isolate  
93 whether people at certain ages (especially children) experience the health impacts of air pollution differently. Also,  
94 the health data on cancer does not disaggregate on the type (only noting that skin cancer is excluded), so we cannot  
95 isolate specific cancers like lung cancer. We also lack detailed demographic data (such as race or gender) to better  
96 understand the population in the health surveys. With only one year of data, it is also difficult to accurately identify  
97 trends in any of the datasets.

98 Finally, our sample size is extremely limited in Kentucky. Although we have health data for all Kentucky counties,  
99 the AQI data was only collected in 27 counties, while Sulfur data was collected in even fewer counties, at 10. The  
100 Kentucky state government website provides some reason for this: “The locations of the monitoring stations are  
101 selected in accordance with U.S. EPA regulatory requirements and are generally established near populous areas or  
102 pollutant sources. Each year the station locations are reviewed to ensure that adequate coverage is being provided.”

103 Although coal production (and consequently coal mining employment) in Kentucky has been declining for years,  
104 the coal industry is still integral to Kentucky culture and politics both in rhetoric and campaign contributions. As  
105 one Kentucky professor shared, “Like bourbon and horse racing, coal is an integral part of Kentucky’s identity. If  
106 you’re going to win an election, you need to be supportive of the industry.”\*\* For example, in a heated senate race in  
107 2016, a GOP-associated non-profit reportedly helped incumbent Mitch McConnell conduct a 12-week, 750,000 digital  
108 campaign which explicitly aimed “to educate Kentuckians on the disastrous policies of the Obama Administration  
109 when it comes to the Commonwealth’s coal-based economy, and on legislation aimed at stopping those policies.” Coal  
110 is still king in Kentucky.<sup>††</sup>

113

## 114 Methods

115 **Data Cleaning.** The data acquired from the sources was in excellent condition and did not require extensive cleaning  
116 at the elemental level.

117 For the purposes of our explorations, analysis, and visualizations, however, certain modifications needed to be  
118 made and are detailed below:

119 • **Health:** We subset the data using a custom function to our focus measures:

- 120 – Current asthma (CASTHMA),  
121 – Chronic obstructive pulmonary disease among adults aged  $\geq 18$  years’ (COPD),  
122 – Cancer (excluding skin cancer) among adults aged  $\geq 18$  years (CANCER),  
123 – Depression among adults aged  $\geq 18$  years (DEPRESSION)  
124 – Visits to doctor for a routine checkup within the past year among adults aged (CHECKUP)

\*\*Coal-Fired Politics: In Kentucky Senate Race, Bitter Rivals Woo a Dying Industry, Katherine Bagley, <https://insideclimateneWS.org/news/08102014-coal-fired-politics-kentucky-senate-race-bitter-rivals-woo-dying-industry/>, 2014

††Trump, McConnell promised to bring back coal. But Kentucky lost 2,700 jobs on their watch, <https://www.courier-journal.com/in-depth/news/politics/elections/kentucky/2020/10/14/election-2020-kentucky-coals-influence-nowhere-found/5927720002/>

- 125           – Current lack of health insurance among adults aged 18-64 years (ACCESS2)  
 126           – Taking medicine for high blood pressure control among adults aged  $\geq 18$  years with high blood pressure  
 127           (BPMED)

128           We did this by aggregating on the basis of category and Measure.  
 129

- 130           • **Air Quality:** This was a relatively clean and aggregated data, by virtue of it being an annual file data.  
 131           • **Sulfur Data:** While most counties only had one site for sulfur data collection, the counties of Jefferson and  
 132           Henderson had two and three sites. Each site would then have records for all days of the year and for different  
 133           data collection methods. Using thresholds set by the EPA, the AQI values were then used to categorize days as:

- 134           – Good days (if  $AQI < 50$ )  
 135           – Moderate Days ( $50 < AQI < 100$ )  
 136           – Unhealthy Days for Sensitive Groups ( $100 < AQI < 150$ )  
 137           – Unhealthy Days ( $150 < AQI$ )

138           Our analysis then led us to group by county name and site number, aggregating for general, maximum, and  
 139           minimum AQI values as well as aggregating to find the total days for each type in the entire year of 2019.  
 140

141           See [0\\_Final Data Cleaning.ipynb](#) for all data cleaning code.

142           **Data Merging.** When we began our data merging and before we had condensed our datasets, our final air-health-sulfur  
 143           dataset was confusingly large, with over a million rows.  
 144

145           Upon deeper analysis, we realized that it was because of the duplication of County Names in both data sets; each  
 146           county would have multiple records of health measures and multiple records of sulfur data reports. In this instance of  
 147           our journey, we realized the need for condensing data. By this, we mean that each county should only have 1 record.  
 148

149           While this would have made our coding exercises much less arduous, we understand that it would have been an  
 150           oversimplification. It would have condensed different measures of health into a single number and readings from  
 151           different sites into one.  
 152

153           As such, we decided to group the SO2 using County Name and Site Number, and similar methods for the health  
 154           dataset. These methods maintained the integrity of the individual readings while generalizing just enough for us to  
 155           gain valuable insights.  
 156

157           After this, it was not challenging to merge the three datasets on the name of the Kentucky county. As this was  
 158           government-issued, there was a great deal of standardization across the County Names.  
 159

160           This post-condensation merging left us with 6,900 rows; a massive reduction in the volume of the merged dataset  
 161           and allowing for more efficient analysis down the line.  
 162

163           However, a county is indeed a large area for us to make any conclusive analyses. With this motivation, we set  
 164           out to calculate the distance between regex-cleaned geolocations to calculate the distance between the sites of data  
 165           collection, using the following Haversine formula:  
 166

$$167 \quad d = 2 * r * \arcsin * \sqrt{(\sin^2((x_2 - x_1)/2)) + (\cos(x_1) * \cos(x_2)) * (\sin^2((y_2 - y_1)/2))}$$

168           where:  
 169

170            $x_i$  = latitude of point i  
 171            $y_i$  = longitude of point i  
 172

173 It is pertinent to note that while both health and sulfur data had geolocations, these were stored in a 'point',  
174 string format. For further analysis, longitude and latitude were extracted using regex methods. We then refined the  
175 location information to precisely locate collection points within 15 miles of each other.

176  
177 Prior to regression, we also identified that each county in the health data set might have up to multiple different  
178 data values for each measure. Therefore we grouped by county and took the mean measure of the data values for  
179 each focus measure. We pivoted this data to a wide dataframe where each column had the average data value. This  
180 allowed us to have one unique value for each county for each focus measure. See [1\\_Final Data Merging.ipynb](#) for all  
181 data cleaning codes.

182 **Regression.** When running a simple linear regression on the data, created two new dataframes: one by merging the  
183 AQI and Health data alone and the other by merging Sulfur and Health data alone. This allowed us to maintain  
184 clean dataframes for two different sets of potential independent variables and test multiple dependent variables.

185  
186 To do so more easily, we created four functions (two for each dataset). These functions created a custom scatterplot  
187 with a regression line and then ran a linear regression model and printed the summary. See [4\\_Regression.ipynb](#)  
188 for these functions. We then used these to test seven dependent variables related to health (all focus measures and  
189 Health Outcomes and Health Status).

190

191 **Results.** We explored many possible questions with our data. In the regression analysis, likely due to our small sample  
192 size of Kentucky counties, we found common trends in our scatterplots but no statistically significant relationships.  
193 For example, fewer people reported depression as the percentage of moderate+ AQI days increased. However, the  
194 air quality data in most cases actually resulted in an opposite trend of what we would expect (i.e., health outcomes  
195 improved in relation to a higher number of moderate+ AQI days). In one example, the average percentage of people  
196 reporting poor health actually decreased as the percentage of moderate+ AQI days increased. In contrast, sulfur  
197 data generally, but this was primarily due to an outlier county. These relationships, again, were not statistically  
198 significant but should be explored with a larger sample size. See [4\\_Regression.ipynb](#) for all regression plots and  
199 model summaries and below, we record some of our more significant investigations:

200

- 201 • **Figure 1** shows the relationship between the number of non-good AQI days and depression. We found out that  
202 as the percentage of moderate or worse air quality days increases, fewer people report cancer. More precisely,  
203 with each 0.1 increase in the proportion of moderate+ AQI days, the portion of people who report depression  
204 goes by 7 percent.
- 205 • **Figure 2** on the other hand, shows a positive relationship between mean SO2 observations and the percentage  
206 of people reporting poor health; an increase of 1 point in the mean SO2 observations increases the portion of  
207 people reporting poor health with 0.15
- 208 • **Figure 3** show that for every 0.1 percentage increase in the portion of days with moderate or unhealthy AQI,  
209 there is a 9 percent decrease in people reporting poor health.

210 As noted above, we failed to reject the null hypothesis for these regression results i.e. this result does not hold  
211 statistical significance, which is generally counter-intuitive to our expectations. While our suggested reasons for this  
212 insignificance have already been discussed, they might be worth recapping:

- 213 • Small data sample, both in terms of the number of counties and in terms of the time for which this analysis was  
214 being done
- 215 • Political context may discourage data collection of sensitive measures

216 We also went on to investigate the variation in reported monthly average AQI across the year:

- 217 • **Figure 4** shows a heat map of the Monthly Average AQI for all the county's in Kentucky. The higher the AQI  
218 value, the greater the level of air pollution and the greater the health concern. We can see that the AQI is lower  
219 during colder months except for Jefferson, where the AQI is consistently higher as compared to other counties.  
220 The grey boxes show the missing AQI for the County.

221 In order to understand the relationship between the sites for sulfur data collection and health measures data  
222 collection, we went on to use the folium package to map both datasets. The interactive visualizations provided [here](#),  
223 will allow for deeper insight into each data point, but here are a few highlights:

- 224
- **Figure 5** shows the locations of health data collection points with associated measures as hover points.  
225 Superimposed on this map are the sites for sulfur data collection, surrounded by a 15-mile, crimson circle. The  
226 figure illustrates that generally, health measures are being collected within 15 miles of the where sulfur data is  
227 being collected. Having said that, there does not seem to be a relation between sulfur collection sites and the  
228 health measures being collected in that area.
  - **Figure 6** shows a magnification of this map, in the busy city of Louisville, which has three different sites for  
229 sulfur data collection. On the other hand, there is a vast number of sites for health data collection, which are  
230 almost entirely clustered towards the eastern side of the city. The mean AQI values for these sites are between 1  
231 and 3, which fails to make any clear association between the sites that are selected for sulfur data collection and  
232 health measure data collection. However, it is interesting to note that there are more health measures being  
233 recorded for this urban city (a total of 6,840) than in other areas e.g. Lexington, another urban area that has 2  
234 sulfur collection sites and 1092 records of health measures, or Paducah, a smaller area with only one sulfur  
235 collection point and 204 records of health measures.
  - **Figure 7** shows a similar magnification of Mammoth National Park. While we only have one site for sulfur  
236 data collection, there are 3 sites for health measure collections (with 12 measures being recorded in each) in  
237 the busy city of Louisville, which has three different sites for sulfur data collection. The sites for health data  
238 collection points are clustered towards the eastern side of the city. The mean AQI values for these sites are  
239 between 1 and 3, which fails to make any clear association between the sites that are selected for sulfur data  
240 collection and health measure data collection.

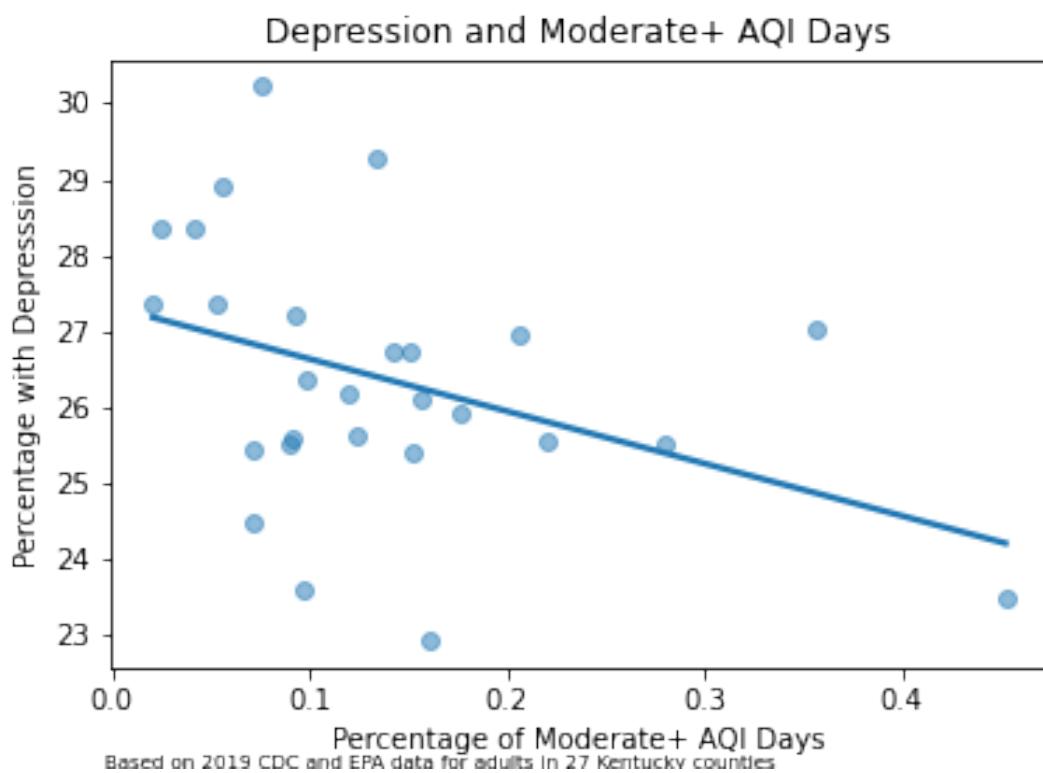
241

242

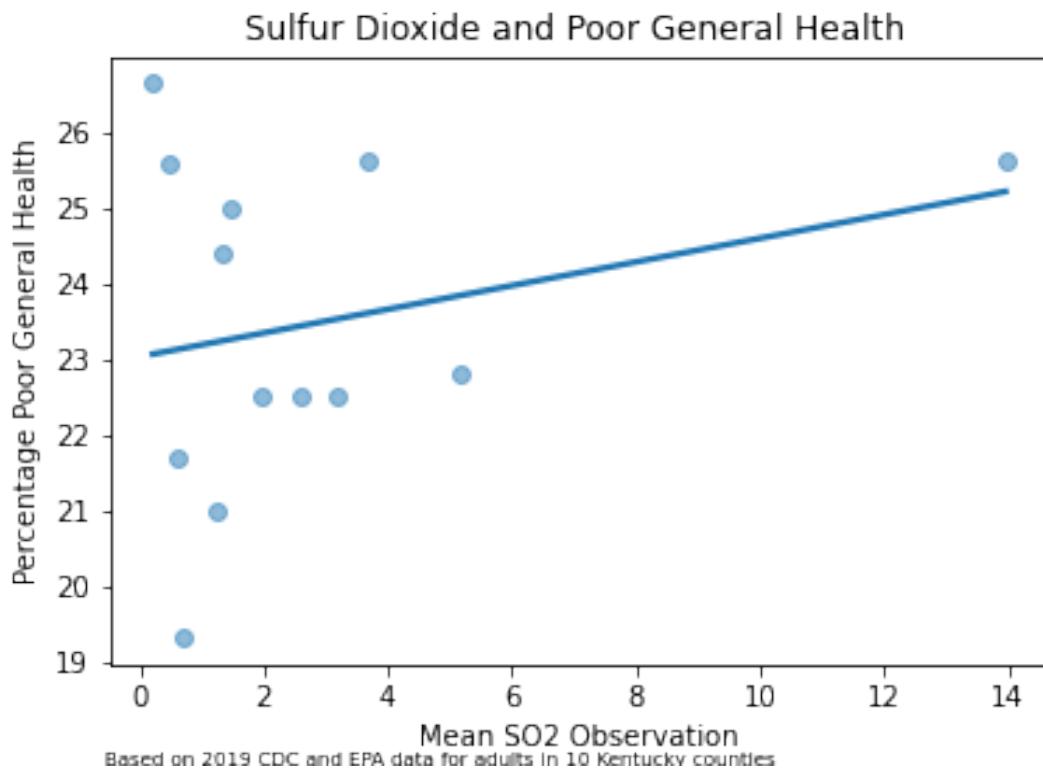
243 **Summary.** Our findings do not support a statistically significant relationship between air quality and health outcomes  
244 in Kentucky. However, as previously discussed, there are many limitations to this research exercise. The primary  
245 limitations of our analysis include the lack of AQI and Sulfur data for all 120 counties in Kentucky and the study of  
246 data from only one year. For future analysis, we recommend a number of ways to expand on our research. First, this  
247 could be adapted to a longitudinal study to better understand trends over time in health outcomes and air quality  
248 levels. As many health effects are not seen for years, this longitudinal study might demonstrate a stronger impact  
249 than our study of a sample from one year. For another example, one could conduct a national analysis of AQI, sulfur,  
250 and health data, with dummy variables for coal-producing states and non-coal-producing states. This would help  
251 isolate the influence of coal production on health outcomes.

252

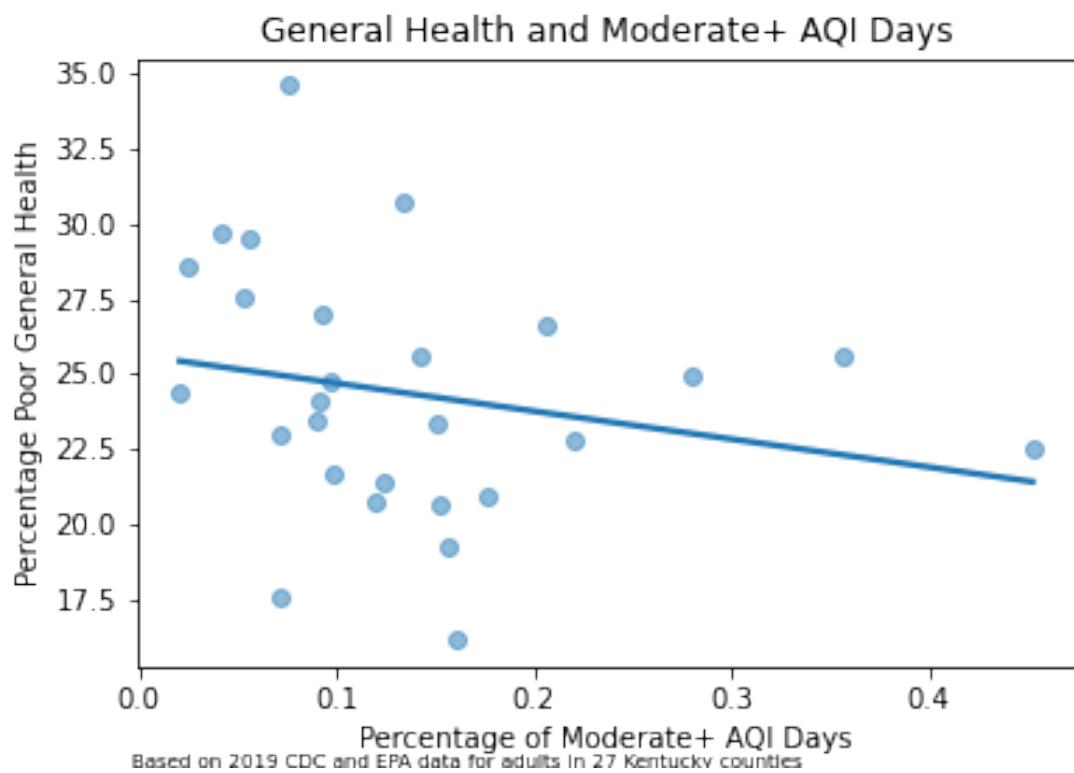
253 In addition to providing a basis for future research, our work has policy implications as it highlights a major  
254 data gap in air quality across Kentucky counties. This research could be used as an argument to expand air quality  
255 measures across additional Kentucky counties. However, due to the political influence of coal mining and power  
256 companies in Kentucky, this is likely to face strong pushback from the private sector as well as government officials  
257 with strong ties (particularly financial ties from campaign donations) from these private sector companies.



**Fig. 1.** Regression result for Depression and Moderate+ AQI Days in Kentucky. Further statistical modeling did not find this relationship significant.



**Fig. 2.** Regression result for Sulfur Dioxide and Poor General Health in Kentucky. Further statistical modeling did not find this relationship significant.



**Fig. 3.** Regression result for Poor General Health and Moderate+ AQI Days in Kentucky. Further statistical modeling did not find this relationship significant.

## Heat Map of Monthly Average AQI for Kentucky

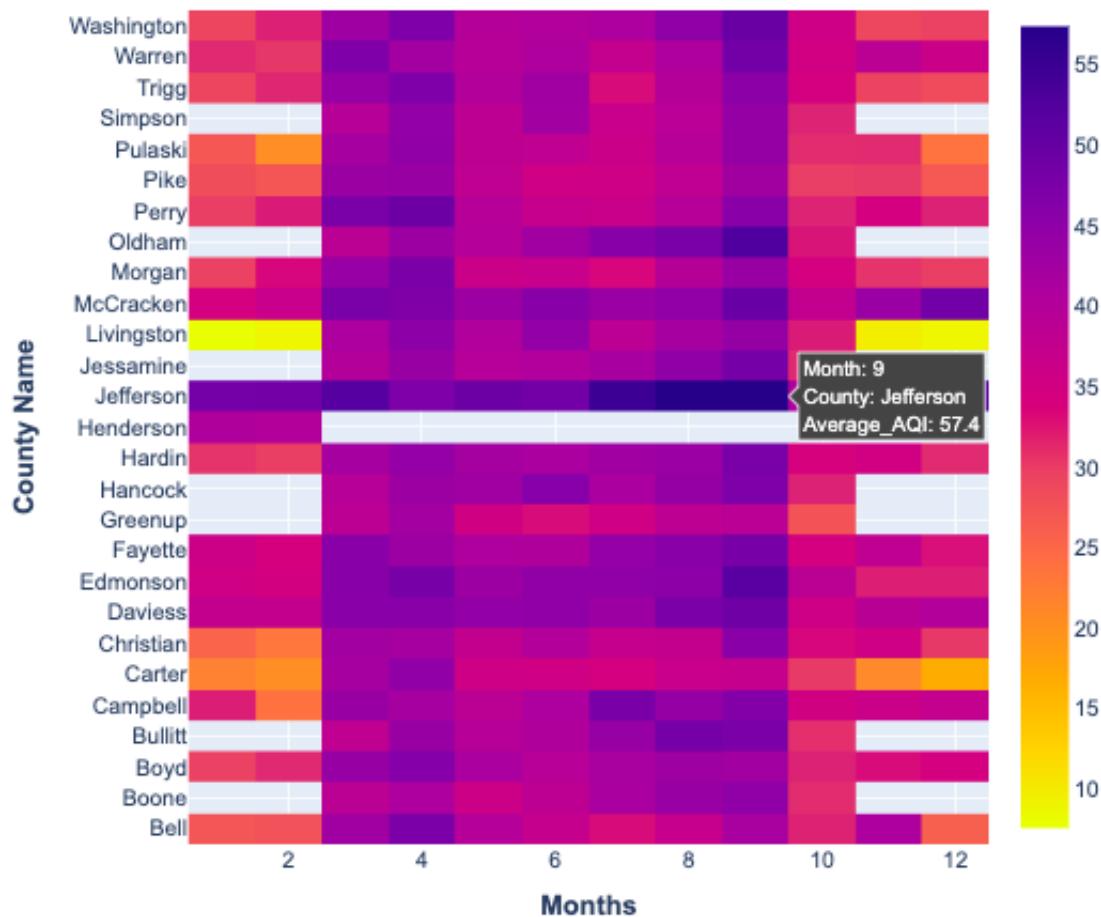


Fig. 4. Heat Map of Monthly Average AQI for Kentucky

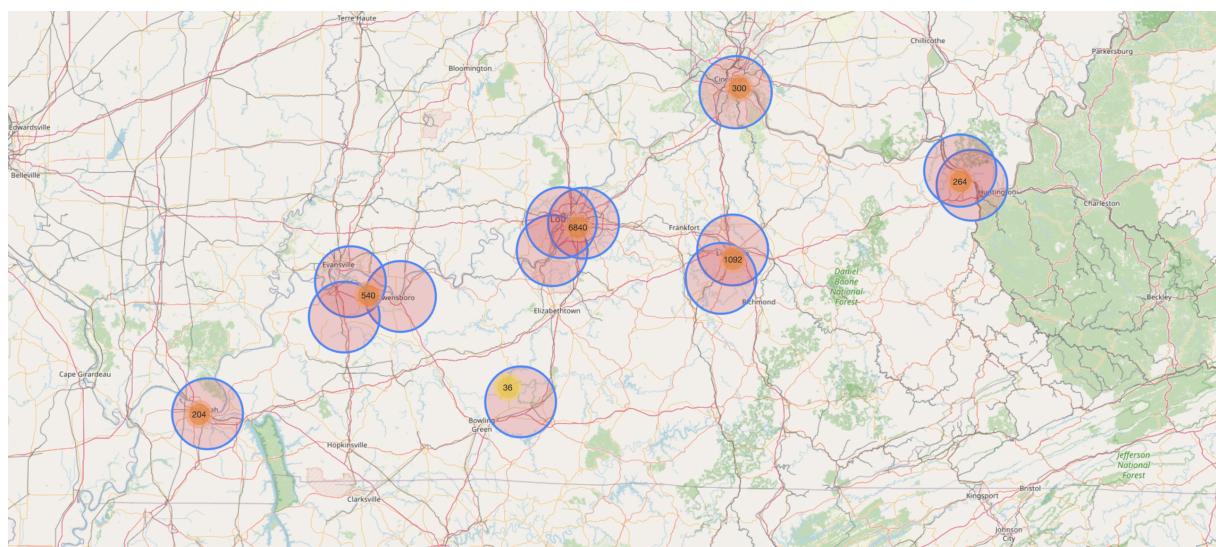
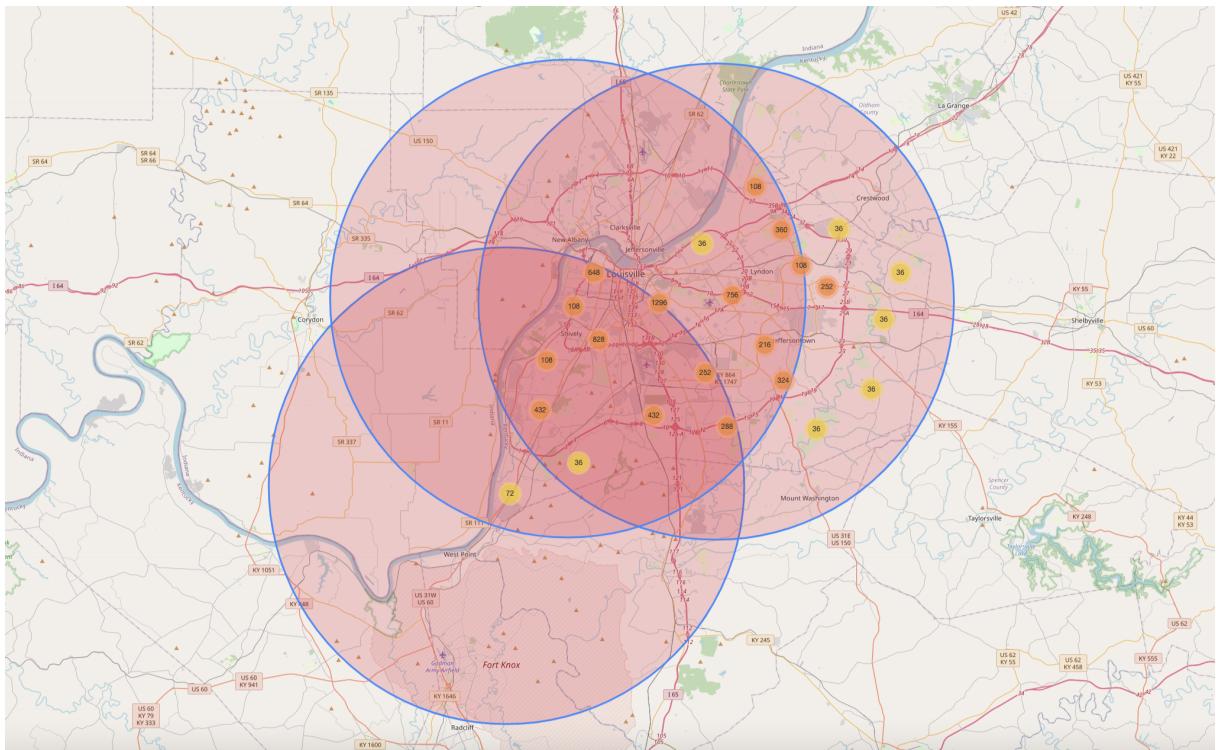
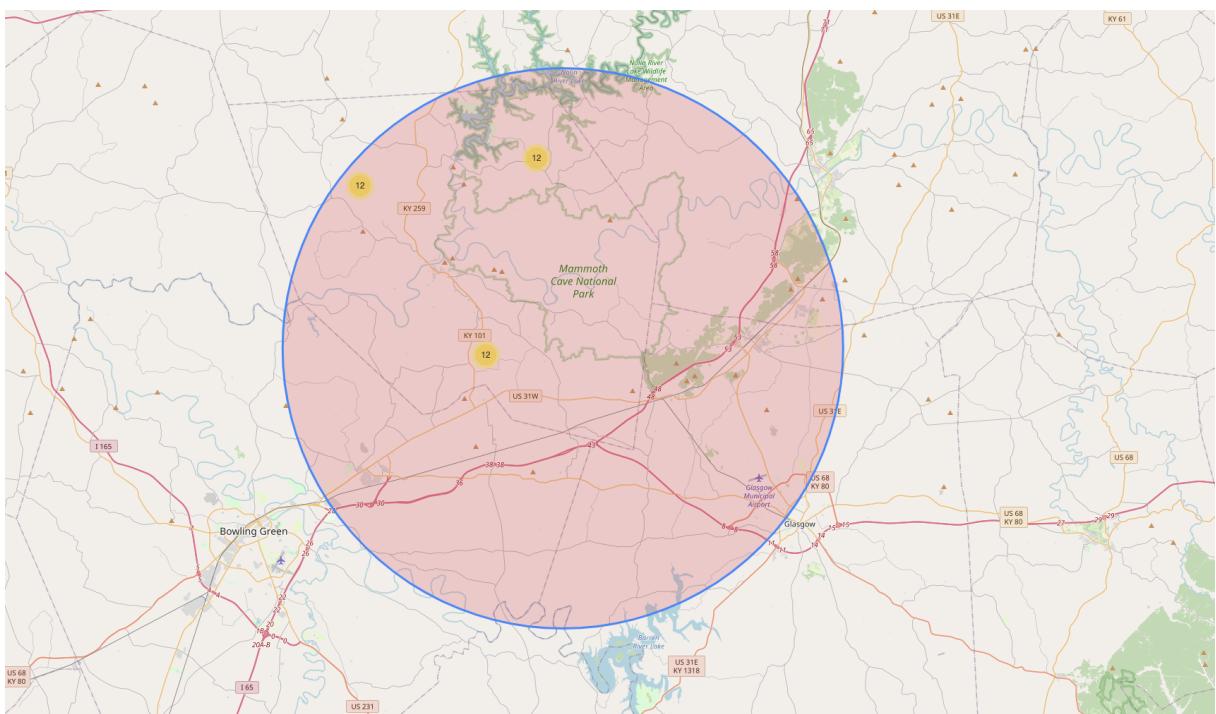


Fig. 5. Comparing data collection points to sulfur and health



**Fig. 6.** Comparing data collection points to sulfur and health: zoomed-in view of Louisville, the largest city in Kentucky.



**Fig. 7.** Comparing data collection points to sulfur and health: zoomed-in view of Mammoth National Park