

# Mini project 1: air quality in U.S. cities

In a way, this project is simple: you are given some data on air quality in U.S. metropolitan areas over time together with several questions of interest, and your objective is to answer the questions.

However, unlike the homeworks and labs, there is no explicit instruction provided about *how* to answer the questions or *when* exactly to begin. Thus, you will need to discern for yourself how to manipulate and summarize the data in order to answer the questions of interest, and you will need to write your own codes from scratch to obtain results. It is recommended that you examine the data, consider the questions, and plan a rough approach before you begin doing any computations.

You have some latitude for creativity: **although there are accurate answers to each question** – namely, those that are consistent with the data – **there is no singularly correct answer**. Most students will perform similar operations and obtain similar answers, but there's no specific result that must be considered to answer the questions accurately. As a result, your approaches and answers may differ from those of your classmates. If you choose to discuss your work with others, you may even find that disagreements prove to be fertile learning opportunities.

The questions can be answered using computing skills taught in class so far and basic internet searches for domain background; for this project, you may wish to refer to HW1 and Lab1 for code examples and the [EPA website on PM pollution](#) for background. However, you are also encouraged to refer to external resources (package documentation, vignettes, stackexchange, internet searches, etc.) as needed – this is the especially good idea if you find yourself thinking, "It would be really handy to do X, but I haven't seen that in class anywhere".

The broader goal of these mini projects is to cultivate your problem-solving ability in an unstructured setting. Your work will be evaluated based on the following:

- choice of method(s) used to answer questions;
- clarity of presentation;
- code style and documentation.

Please write up your results separately from your codes; codes should be included at the end of the notebook.

## Part I: Dataset

Merge the city information with the air quality data and tidy the dataset (see notes below). Write a brief description of the data.

In your *description*, answer the following questions:

- What is a CBSA (the geographic unit of measurement)?
- How many CBSAs are included in the data?
- In how many states and territories do the CBSAs reside? (*Hint: `str.split()`*)
- In which years were data values recorded?
- How many observations are recorded?
- How many variables are measured?
- Which variables are non-missing most of the time (*i.e.*, in at least 50% of instances)?
- What is PM 2.5 and why is it important?

Please write your description in narrative fashion; **please do not list answers to the questions above one by one**. A few brief paragraphs should suffice; please limit your data description to three paragraphs or less.

## Air quality data

The CBSA is a U.S. geographic area defined by the Office of Management and Budget (OMB) based on census data and referred collectively to both metropolitan statistical areas and micropolitan areas. To calculate how many CBSA's are included in the data, I calculated the length of each unique CBSA's, and the answer is 351. To calculate how many different states are included, I similar used `strsplit` function to get every names of states after the commas and then `strsplit` them accordingly to slashes between states. I then used the `set` function to count how many unique states there are, and it returned as 52. To calculate how many territories are included I used `strsplit` function to get every names states and territories. I used an if statement to find append all the territories in my territory list. I then measure the length of the list(unique) and it returned as 36.The data was recorded from 2000 to 2019, and it could be seen by just looking at how many columns of values there were.

In order to get how many variables are measured, I firstly melted the data according to years, which means letting different years become rows. Then, I pivoted the dataframe by different variables in pollutant and trend statistics. There are in total 1134 observations measured across 20 years (because the dataset has 1134 rows). It turns out there are 9 variables(some pollutants have multiple trend statistics, and there are 7 pollutants in total, but 9 variables) measured acrossing 20 years. Moreover, I then used `isna.sum()` to calculate the number of missing values for different pollutants. Note that there are in total 7020 datasets in this pivoted table, thus any sums below the half of 7020 are non-missing most of the time. It turns out that PM2.5 and O3 are the variables of nonmissing most of the time. PM2.5 is an air pollutant that is a concern for people's health when levels in air are high. PM2.5 are tiny particles in the air that reduce visibility and cause the air to appear hazy when levels are elevated. It is important because it could harm people's health

## Part II: Descriptive analysis

Focus on the PM2.5 measurements that are non-missing most of the time. Answer each of the following questions in a brief paragraph or two. Your paragraph(s) should indicate both your answer and a description of how you obtained it: **please do not include codes with your answers**.

Has PM 2.5 air pollution improved in the U.S. on the whole since 2000?

In order to calculate if pm2.5 air pollution has improved, I melted the dataset according to different years and then groupby years. The index that I looked at was weighted annual mean which was only included in pm2.5. I took the average after the groupby function in different years, and we can see a decline across 20 years. So yes, the pollution has improved from city to city in US.

Over time, has PM 2.5 pollution become more variable, less variable, or about equally variable from city to city in the U.S.?

In order to see if the pollution has become more or less variable, I used the same melted dataset from previous question. Similarly to the last question, except that I did not take the average after grouping by the melted dataframe with different years, I took the standard deviation function and we could see a decline across 20 years. So, it has become less variable

Which state has seen the greatest improvement in PM 2.5 pollution over time? Which city has seen the greatest improvement?

The way that I interpreted the word "improvement" was by looking at the difference at PM2.5 pollution between year 2019 and 2000. To look at the difference of cities between 2019 and 2000, I selected the dataframes in melted datasets(where years were the index) with years equal to 2019 and 2000. I then took out the numeric values of different years as two long vectors(or dataframe, to be more precise). I renamed these vectors with corresponding years and merged these two vectors to a dataframe together. I then added an empty column to the new dataframe to calculate the difference between two years. Lastly, I used `idxmax` to get the index of the greatest difference between two years. It turns out to be Portsmouth, OH.

To answer the questions that which state has seen the greatest improvement, I used similar methodology and the only difference is that I have to `strsplit` the column Core Based Statistical area after the comma to get different state names. I did it so and grouped by the dataset with different states and looked at their mean(since there are multiple territories within a state). Similarly I inserted a column with empty entries to the dataframe and calculated the difference between two years. I then used `idxmax()` to get the index of greatest difference between two years. The answer turns out to be WV. I used a for loop to go through all the entries without '-' because it means territory in the dataset.

Choose a location with some meaning to you (e.g. hometown, family lives there, took a vacation there, etc.). Was that location in compliance with EPA primary standards as of the most recent measurement?

I took a look at my cleaned dataset where it only included the weighted annual mean for pm2.5. The city that I chose was Cincinnati because I had high school over there. At year 2019, it had 9.2 unit of PM2.5. I looked up online and the epa standards for annual mean was 12. So yes, it has reached the epa primary standard for pm2.5 in the most recent measurement

## Extra credit: Imputation

One strategy for filling in missing values (imputation) is to use non-missing values to predict the missing ones; the success of this strategy depends in part on the strength of relationship between the variable(s) used as predictors of missing values.

Identify one other pollutant that might be a good candidate for imputation based on the PM 2.5 measurements and explain why you selected the variable you did. Can you envision any potential pitfalls to this technique?

The way that I looked this question is by looking at the correlation matrix between two variables PM2.5 and other pollutants. It turns out that O3 has the strongest correlation with PM2.5 weighted annual mean. So O3 is the most likely to be predicted. Some pitfalls about this technique is that the prediction might not be good because there are biases exist in different cities due to missing values. So the best solution is to find a way to fill up the missing values.

## Codes

```
In [1]: # packages
import numpy as np
import pandas as pd
import sklearn

# raw data
air_raw = pd.read_csv('air-quality.csv')
csba_info = pd.read_csv('csba-info.csv')
data=pd.merge(air_raw, csba_info, how = 'left', on='CBSA')
# PART I
data
```

```
Out[1]:
```

	CBSA	Pollutant	Trend Statistic	Number of Trends Sites	2000	2001	2002	2003	2004	2005	...	2011	2012	2013	2014	2015	2016	2017	2018	2019	Core Based Statistical Area
0	10100	PM10	2nd Max	1	50.000	58.000	59.000	66.000	39.000	48.000	...	29.000	62.000	66.000	36.000	43.000	65.000	40.000	49.000	35.000	Aberdeen, SD
1	10100	PM2.5	Weighted Annual Mean	1	8.600	8.600	7.900	8.400	8.100	9.000	...	7.100	7.500	7.300	6.200	6.200	5.400	5.800	6.600	5.900	Aberdeen, SD
2	10100	PM2.5	98th Percentile	1	23.000	23.000	20.000	21.000	23.000	23.000	...	18.000	23.000	22.000	17.000	14.000	14.000	13.000	22.000	18.000	Aberdeen, SD
3	10300	O3	4th Max	1	0.082	0.086	0.089	0.088	0.074	0.082	...	0.076	0.087	0.064	0.068	0.065	0.069	0.066	0.071	0.059	Adrian, MI
4	10420	CO	2nd Max	1	2.400	2.700	1.800	1.900	2.100	1.600	...	1.000	1.100	0.800	0.800	1.000	1.100	0.900	1.800	1.800	Akron, OH
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
1129	49700	NO2	Annual Mean	1	13.000	14.000	15.000	14.000	12.000	12.000	...	8.000	10.000	10.000	8.000	7.000	7.000	7.000	7.000	6.000	Yuba City, CA
1130	49700	NO2	98th Percentile	1	62.000	62.000	62.000	62.000	52.000	51.000	...	44.000	46.000	52.000	44.000	39.000	40.000	42.000	41.000	40.000	Yuba City, CA
1131	49700	O3	4th Max	2	0.081	0.077	0.090	0.085	0.076	0.075	...	0.070	0.073	0.066	0.072	0.068	0.072	0.074	0.073	0.063	Yuba City, CA
1132	49700	PM2.5	Weighted Annual Mean	1	10.600	11.900	13.100	9.500	10.000	9.500	...	8.000	6.900	8.200	9.400	9.600	8.100	9.300	10.300	8.400	Yuba City, CA
1133	49700	PM2.5	98th Percentile	1	38.000	54.000	34.000	29.000	38.000	42.000	...	37.000	24.000	25.000	25.000	31.000	22.000	32.000	37.000	27.000	Yuba City, CA
1134 rows × 25 columns																					

```
In [2]: air_raw

Out[2]:
```

	CBSA	Pollutant	Trend Statistic	Number of Trends Sites	2000	2001	2002	2003	2004	2005	...	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019
0	10100	PM10	2nd Max	1	50.000	58.000	59.000	66.000	39.000	48.000	...	46.000	29.000	62.000	66.000	36.000	43.000	65.000	40.000	49.000	35.000
1	10100	PM2.5	Weighted Annual Mean	1	8.600	8.600	7.900	8.400	8.100	9.000	...	8.700	7.100	7.500	7.300	6.200	6.200	5.400	5.800	6.600	5.900
2	10100	PM2.5	98th Percentile	1	23.000	23.000	20.000	21.000	23.000	23.000	...	27.000	18.000	23.000	22.000	17.000	14.000	14.000	13.000	22.000	18.000
3	10300	O3	4th Max	1	0.082	0.086	0.089	0.088	0.074	0.082	...	0.066	0.076	0.087	0.064	0.068	0.065	0.069	0.066	0.071	0.059
4	10420	CO	2nd Max	1	2.400	2.700	1.800	1.900	2.100	1.600	...	1.000	1.100	0.800	0.800	1.000	1.100	0.900	1.800	1.800	...
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
1129	49700	NO2	Annual Mean	1	13.000	14.000	15.000	14.000	12.000	12.000	...	8.000	8.000	10.000	10.000	8.000	7.000	7.000	7.000	7.000	6.000
1130	49700	NO2	98th Percentile	1	62.000	62.000	62.000	62.000	52.000	51.000	...	45.000	44.000	46.000	52.000	44.000	39.000	40.000	42.000	41.000	40.000
1131	49700	O3	4th Max	2	0.081	0.077	0.090	0.085	0.076	0.075	...	0.067	0.070	0.073	0.066	0.072	0.068	0.072	0.074	0.073	0.063
1132	49700	PM2.5	Weighted Annual Mean	1	10.600	11.900	13.100	9.500	10.000	9.500	...	5.900	8.000	6.900	8.200	9.400	9.600	8.100	9.300	10.300	8.400
1133	49700	PM2.5	98th Percentile	1	38.000	54.000	34.000	29.000	38.000	42.000	...	17.000	37.000	24.000	25.000	25.000	31.000	22.000	32.000	37.000	27.000
1134 rows × 24 columns																					

```
In [3]: csba_number=len(pd.unique(data['CBSA']))
csba_number #this counts how many different csba's are included

Out[3]:
351
```

```
In [4]: territory_list=data.loc[:,['Core Based Statistical Area']]
territory_name=[]
for i in range(0,len(territory_list)):
    territory_name.append(territory_list[i].split(',')[0])
len(pd.unique(territory_name))
#this counts how many different territories(cities) the dataset includes
```

```
Out[4]:
341
```

```
In [5]: state_name=[]
for i in range(0,len(territory_list)):
    state_name.append(territory_list[i].split(',')[1])
state_name_unique=[]
for i in range(0,len(state_name)):
    state_name_unique.append(state_name[i].split('-'))
state_name_unique_flat=[]
for slist in state_name_unique:
    for item in slist:
        state_name_unique_flat.append(item)
print(set(state_name_unique_flat))
len(set(state_name_unique_flat))
#this counts how many different states are included in the dataset
```

```
{'KS', 'CA', 'RI', 'D', 'AZ', 'MN', 'NM', 'MD', 'NJ', 'OR', 'NC', 'PA', 'WV', 'LA', 'HI', 'PR', 'MT', 'SC', 'MI', 'AL', 'VA', 'VT', 'NY', 'UT', 'NE', 'NV', 'OH', 'ND', 'I', 'A', 'TN', 'NA', 'FL', 'NE', 'CT', 'MS', 'MO', 'OK', 'CO', 'AR', 'DE', 'GA', 'DC', 'IL', 'TX', 'NH', 'IN', 'KY', 'WA', 'MA', 'AK', 'SD', 'WI'}
```

```
Out[5]:
52
```

```
In [6]: state_name
territory_real=[]
for i in range(0,len(state_name)):
    if(state_name[i].find('-')==1):
        territory_real.append(state_name[i])
len(pd.unique(territory_real))
```

```
Out[6]:
36
```

```
In [7]: years=['2000','2001','2002','2003','2004','2005','2006','2007','2008','2009','2010','2011','2012','2013','2014','2015','2016','2017','2018','2019']
melteddata=pd.melt(data, id_vars=['Trend Statistic','Core Based Statistical Area','CBSA','Pollutant'],
value_vars=years,
var_name='years')
melteddata.pollutant.pivot(melteddata.pollutant, pivot_table(melteddata.pollutant, index=['years', 'CBSA','Core Based Statistical Area'],columns=['Pollutant','Trend Statistic'],
melteddata.pollutant.pivot
```

```
Out[7]:
```

			Pollutant	CO	NO2	O3	PM10	PM2.5	Pb	SO2	value	
			Trend Statistic	2nd Max	98th Percentile	Annual Mean	4th Max	2nd Max	98th Percentile	Weighted Annual Mean	Max 3-Month Average	99th Percentile
years	CBSA	Core Based Statistical Area										
2000	10100	Aberdeen, SD	NaN	NaN	NaN	NaN	50.0	23.0	8.6	NaN	NaN	
10300	Adrian, MI	NaN	NaN	NaN	0.082	NaN	NaN	NaN	NaN	NaN	NaN	
10420	Akron, OH	2.4	NaN	NaN	NaN	0.085	NaN	37.0	16.2	NaN	163.0	
10500	Albany, GA	NaN	NaN	NaN	NaN	NaN	38.0	16.6	NaN	NaN	NaN	
10580	Albany-Schenectady-Troy, NY	1.1	NaN	NaN	0.070	NaN	30.0	12.4	NaN	NaN	56.0	
...	...	...	...	...	...	...	...	...	...	...	...	
2019	49340	Worcester, MA-CT	NaN	NaN	NaN	0.060	NaN	NaN	NaN	NaN	NaN	
49420	Yakima, WA	NaN	NaN	NaN	NaN	NaN	32.0	9.2	NaN	NaN	NaN	
49620	York-Hanover, PA	0.7	42.0	7.0	0.062	NaN	20.0	8.8	NaN	NaN	8.0	
49660	Youngstown-Warren-Boardman, OH-PA	NaN	NaN	NaN	0.065	31.3	19.0	7.7	NaN	NaN	5.0	
49700	Yuba City, CA	NaN	40.0	6.0	0.063	NaN	27.0	8.4	NaN	NaN	NaN	
7020 rows × 9 columns												

```
In [8]: melteddata.pollutant.pivot(isna().sum())
#note that there are in total 7020 datasets in this pivoted table, thus we divide the nasum by 7020 to get the percentage of nonmissing,
# thus PM2.5 and O3 are the variables of nonmissing most of the time
```

```
Out[8]:
```

	Pollutant	Trend Statistic	5840
value	CO	2nd Max	5688
	NO2	98th Percentile	5240
		Annual Mean	1340
	O3	4th Max	4960
	PM10	2nd Max	2740
	PM2.5	98th Percentile	2740
		Weighted Annual Mean	2740
	Pb	Max 3-Month Average	6750
	SO2	98th Percentile	5240
dtype: int64			

```
In [9]: years=['2000','2001','2002','2003','2004','2005','2006','2007','2008','2009','2010','2011','2012','2013','2014','2015','2016','2017','2018','2019']
melteddata=pd.melt(data, id_vars=['Trend Statistic','Core Based Statistical Area','CBSA','Pollutant'],
value_vars=years,
var_name='years')
melteddata[melteddata['Trend Statistic']=='Weighted Annual Mean'].groupby('years').mean().value
```

```
Out[9]:
```

years	
2000	13.857944
2001	12.688328
2002	12.352336
2003	11.853271
2004	11.642056
2005	12.479439
2006	11.360748
2007	11.573264
2008	10.625234
2009	9.671828
2010	9.490374
2011	9.638318
2012	8.973364
2013	8.798598
2014	8.660748