# HOW DOES AIR QUALITY AFFECT DISEASE RATES?

GROUP 9

Daniel Tiessalo

Antoni Kajrys

Anssi Lampinen

Ignacio Wąsowicz-Peinado

Aleksei Filonov

CS-A1160 Beginner's Python for Engineers

# WHAT DO WE WANT TO FIND?

Hypothesis: High air pollution increases the prevalence of diseases

First, we plot the data against each other and use the visual method to see whether there exists any correlation

Then, we compute the correlation coefficients and their p-values

Finally, we interpret the results and their significance

# MAIN DATASET

→ US Air quality data from around 400 cities (2020)

**who_region**
**Iso3**
**country_name**
**city**
**year**
**version**
**pm10_concentration**
**pm25_concentration**
**no2_concentration**
**pm10_tempcov**
**pm25_tempcov**
**no2_tempcov**
**type_of_stations**
**reference**
**web_link**
**+ 5 other...**

drop columns and
filter by 2020

→

who_region
Iso3
country_name
**city**
year
version
**pm10_concentration**
**pm25_concentration**
**no2_concentration**
pm10_tempcov
pm25_tempcov
no2_tempcov
type_of_stations
reference
web_link
+ 5 other...

# ADDITIONAL DATASET

→ US health data from over 3000 cities (2020)

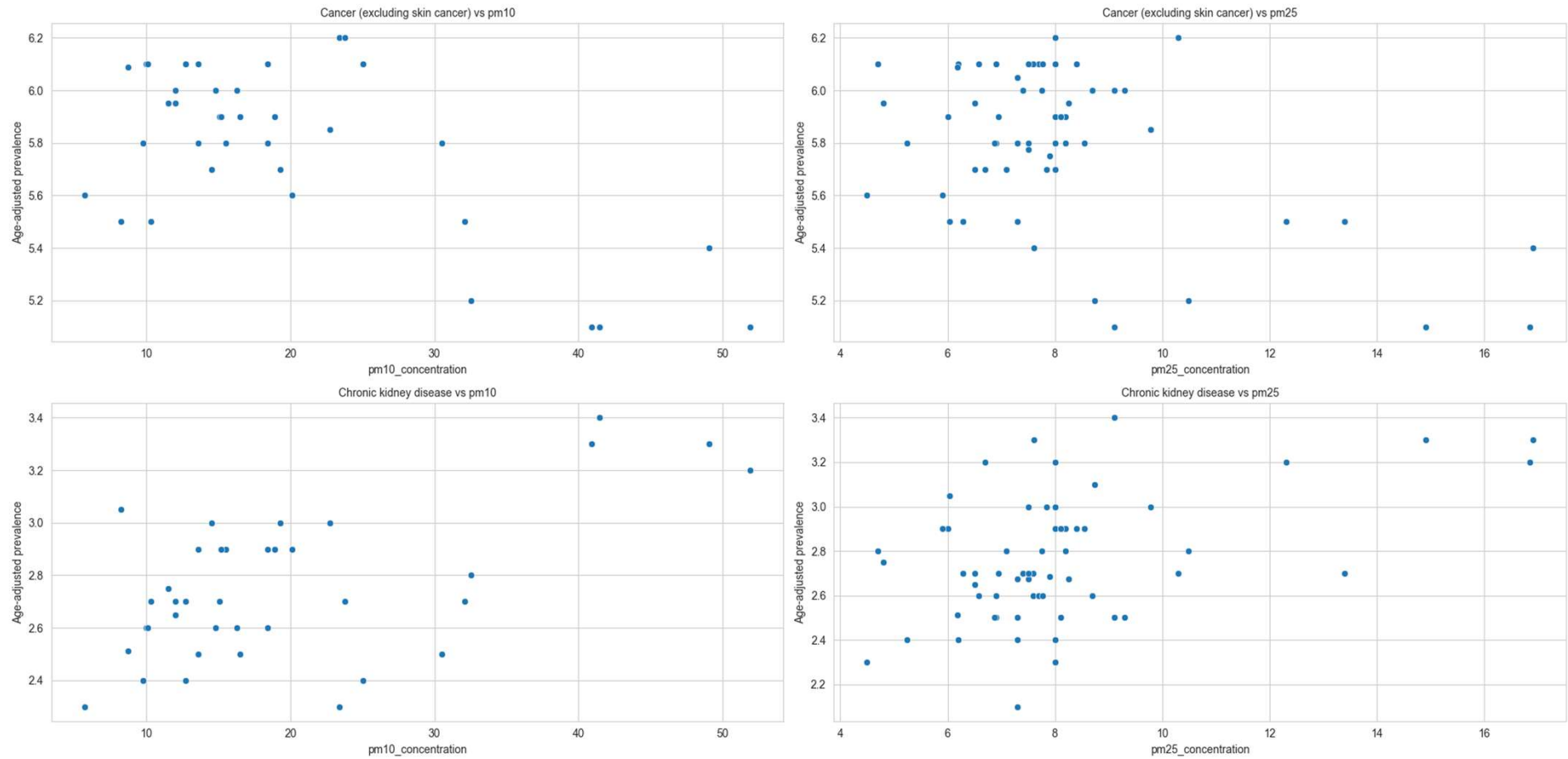| | | |
|---|---|---|
| **Year**<br>**StateAbbr**<br>**StateDesc**<br>**LocationName**<br>**DataSource**<br>**Category**<br>**Measure**<br>**Data_Value_Unit**<br>**Data_Value_Type**<br>**Data_Value**<br>**Low_Confidence_Limit**<br>**High_Confidence_Limit**<br>**TotalPopulation**<br>**LocationID**<br>**CategoryID**<br>**+ 7 other…** | drop columns and<br>filter by 2020<br>——————→ | Year<br>StateAbbr<br>StateDesc<br>**city (prev. LocationName)**<br>DataSource<br>Category<br>**Measure**<br>Data_Value_Unit<br>Data_Value_Type<br>Data_Value<br>Low_Confidence_Limit<br>High_Confidence_Limit<br>TotalPopulation<br>LocationID<br>CategoryID<br>+ 7 other… |

pivot
——————→

**Measure**
**city**
**Cancer (excluding skin cancer) among adults aged >=18 years**
**Chronic kidney disease among adults aged >=18 years**
**Chronic obstructive pulmonary disease among adults aged >=18 years**
**Coronary heart disease among adults aged >=18 years**
**Current asthma among adults aged >=18 years**
**Current smoking among adults aged >=18 years**
**+ 7 other…**

Cities in Air Pollutant - Disease Prevelance frame of reference

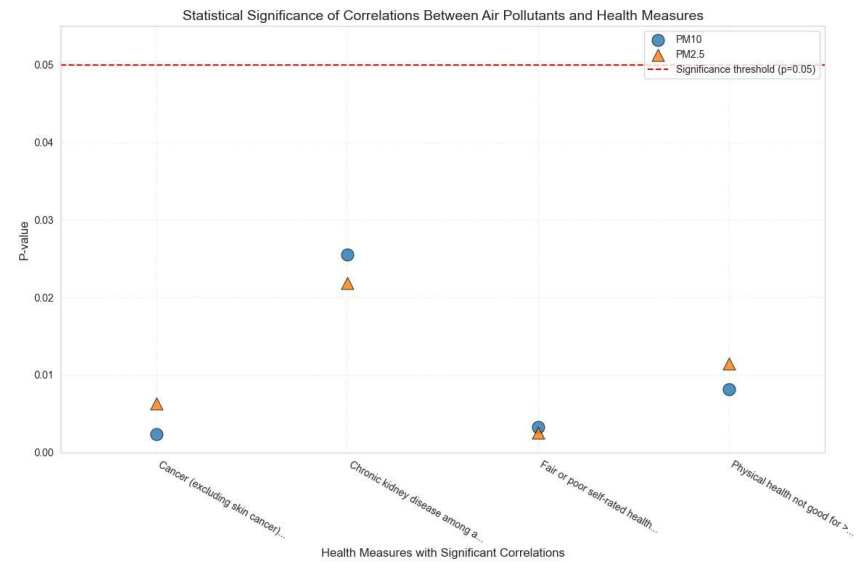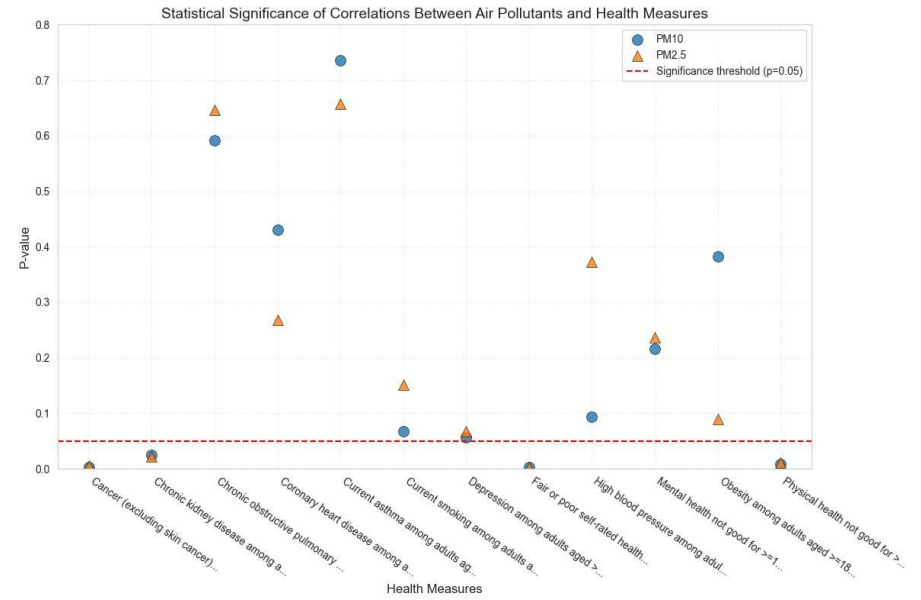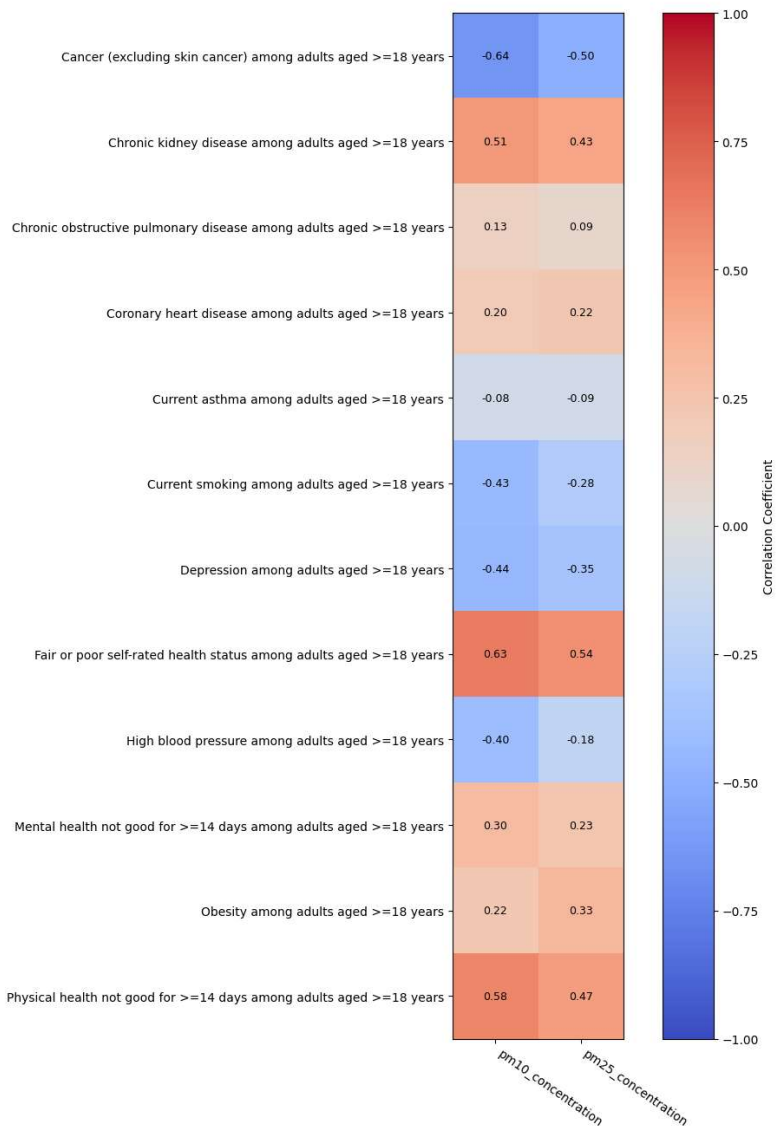CS-A1160 Beginner's Python for Engineers

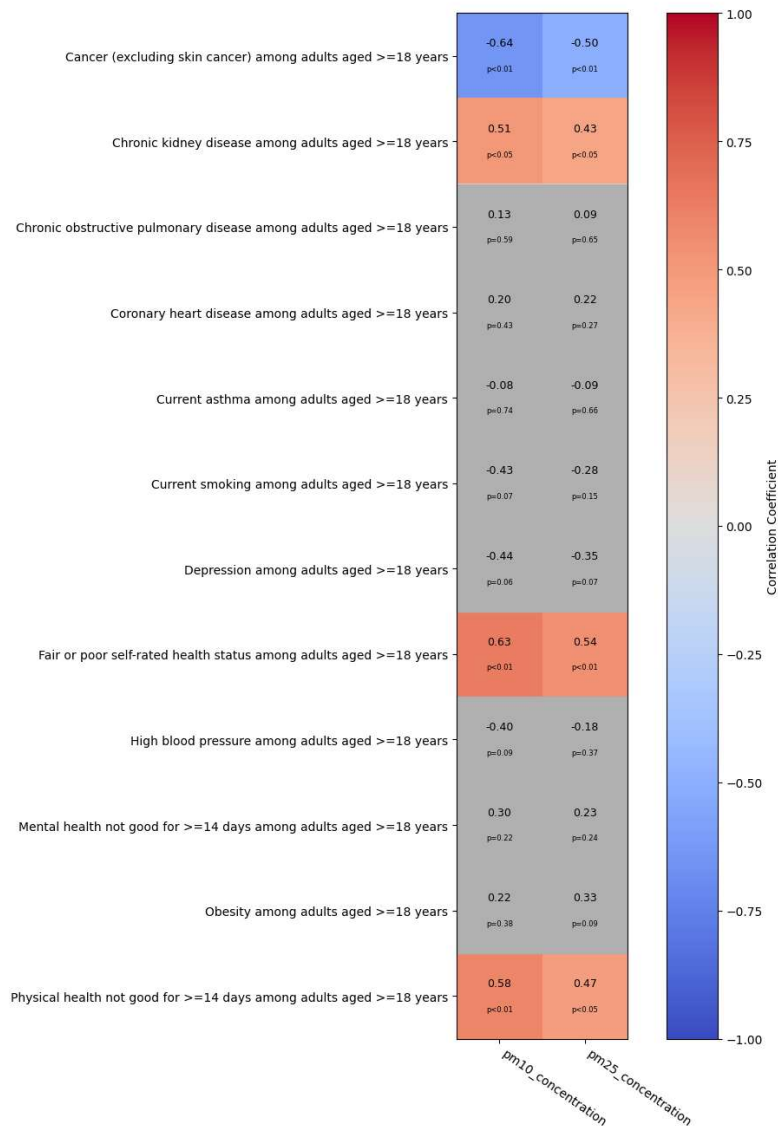# CALCULATE CORRELATIONS AND CHECK THEIR SIGNIFICANCE

- We first check if our data is normally distributed with a **Shapiro-Wilk test** → it is not

- We choose the **Spearman correlation test** because it does not require the normality of the data

- We use the **weighted test** to represent the difference in significance of different cities in the statistic based on their population

- We calculate the coefficients and their p-values

- We filter the coefficients so that only the statistically significant ones remain ($p\text{-value} < 0.05$)

Correlation between Air Pollutants and Medical Conditions

Statistical Significance of Correlations Between Air Pollutants and Health Measures

Statistical Significance of Correlations Between Air Pollutants and Health Measures

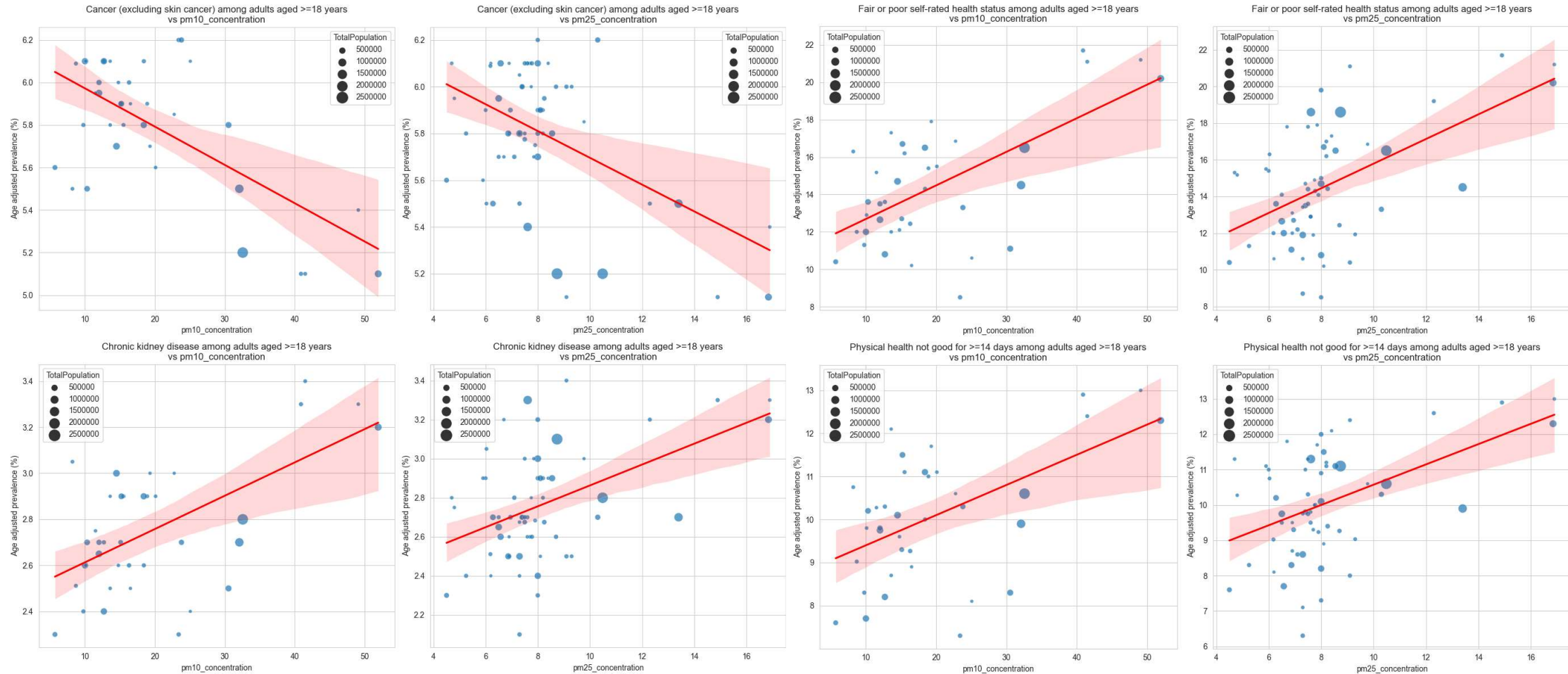Correlation between Air Pollutants Health Measures With p-values
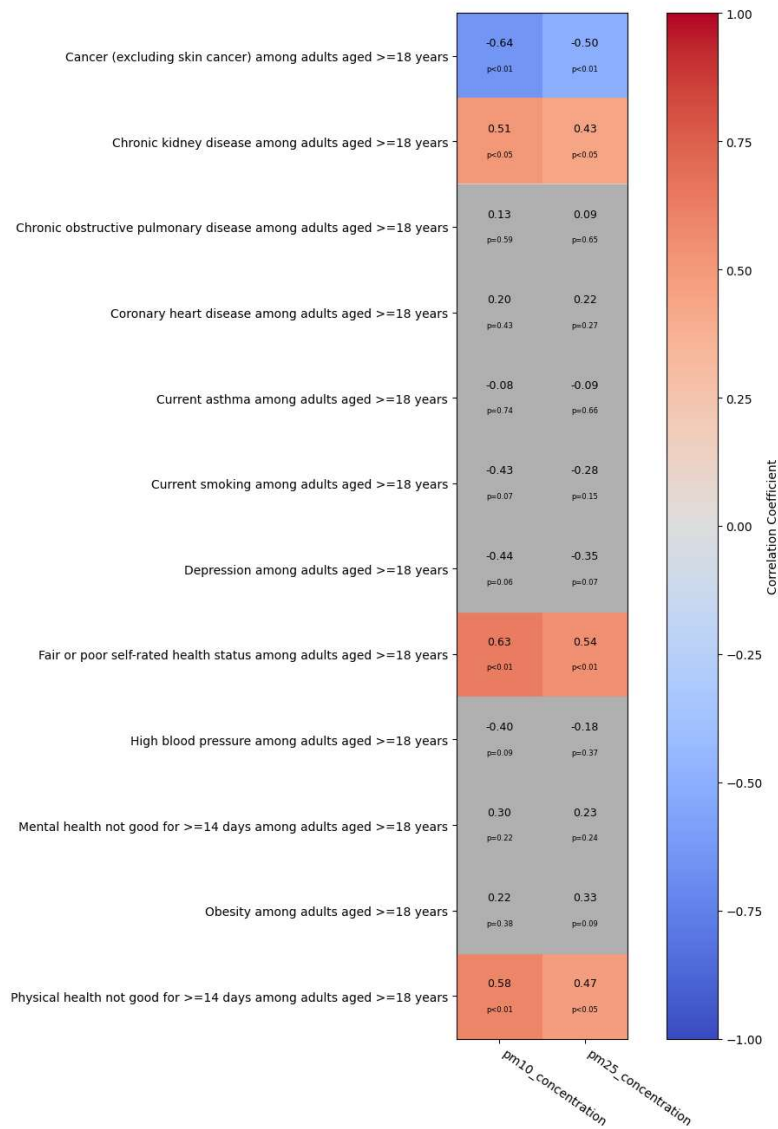
# Conclusions

- Prevalence of cancer among adults is negatively correlated with air pollution

- Chronic kidney disease is positively correlated with air pollution

- Poor self-rated health status among adults is positively correlated with air pollution

- Bad physical health is positively correlated with air pollution

- Prevalence of neither asthma nor chronic obstructive pulmonary disease have significant correlation with the pollution levels, even though both are respiratory track conditions

CS-A1160 Beginner's Python for Engineers

# BUBBLE PLOTS



CS-A1160 Beginner's Python for Engineers

Correlation between Air Pollutants Health Measures With p-values

# Conclusions

- Prevalence of cancer among adults is negatively correlated with air pollution

- Chronic kidney disease is positively correlated with air pollution

- Poor self-rated health status among adults is positively correlated with air pollution

- Bad physical health is positively correlated with air pollution

- Prevalence of neither asthma nor chronic obstructive pulmonary disease have significant correlation with the pollution levels, even though both are respiratory track conditions

CS-A1160 Beginner's Python for Engineers

# IF YOU DON'T WANT CANCER, LIVE IN CRACOW!!!

# CORRELATION IS NOT CAUSATION!

If we believed this study blindly, we would find that:

- Living in cities with high air pollution significantly reduces the chance of having cancer

- Depression is less prevalent in the cities with poor air quality

- More air pollution can decrease the prevalence of asthma and chronic obstructive pulmonary diseases

# THIS DOES NOT MEAN THAT THE STUDY IS MEANINGLESS

- Smoking negatively correlates with air pollutants in this dataset.

- The conclusions of this analysis are inspirational.

- Correlation studies usually provide a starting point instead of concrete answers.

# WE HAD SOME PROBLEMS...

WE ALSO LEARNED SOMETHING

# THANK YOU FOR YOUR ATTENTION

**Aleksei Filonov**
Presentation
Cleaning of datasets
and merging
Correlation plot

**Anssi Lampinen**
Analysis
Weighing of
statistics
Presentation

**Daniel Tiessalo**
Presentation
Cleaning of datasets
and merging
Finding datasets

**Antoni Kajrys**
Finding datasets
Statistics
Checking of significance
Seeker of methods
Presentation

**Ignacio Wąsowicz**
Statistics
Presentation
Checking of
significance