# How to conduct a big data analysis on air pollution and health?

Mathematical Institute of the Serbian Academy of Sciences and Arts
Mart 29, 2022
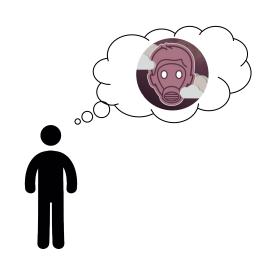
Ana Trisovic, Harvard University

# How to conduct a big data analysis on air pollution and health?

How to conduct a big data analysis on air pollution and health?

# Analysis design

- Secondary data analysis - using data from existing data sources, integrating it and applying study design
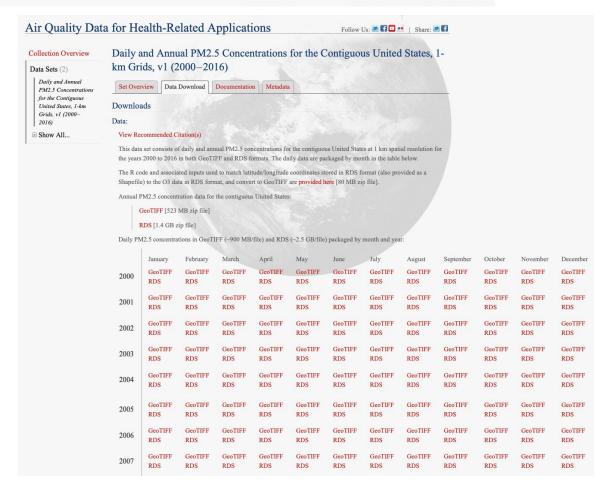- Hypothesis-driven - used to answer presupposed hypothesis or a research question

# Exposure - hypothesized cause of the disease

| |
|---|
| Nitrogen dioxide $NO_2$ |
| Ozone $O_3$ |
| Particulate Matter $PM_{2.5}$ |

Air Quality Data for Health-Related Applications

Follow Us: | Share:

## Daily and Annual PM2.5 Concentrations for the Contiguous United States, 1-km Grids, v1 (2000–2016)

| Set Overview | Data Download | Documentation | Metadata |

### Downloads

Data:

View Recommended Citation(s)

This data set consists of daily and annual PM2.5 concentrations for the contiguous United States at 1 km spatial resolution for the years 2000 to 2016 in both GeoTIFF and RDS formats. The daily data are packaged by month in the table below.

The R code and associated inputs used to match latitude/longitude coordinates stored in RDS format (also provided as a Shapefile) to the O3 data in RDS format, and convert to GeoTIFF are provided here [80 MB zip file].

Annual PM2.5 concentration data for the contiguous United States:

　GeoTIFF [523 MB zip file]

　RDS [1.4 GB zip file]

Daily PM2.5 concentrations in GeoTIFF (~900 MB/file) and RDS (~2.5 GB/file) packaged by month and year:

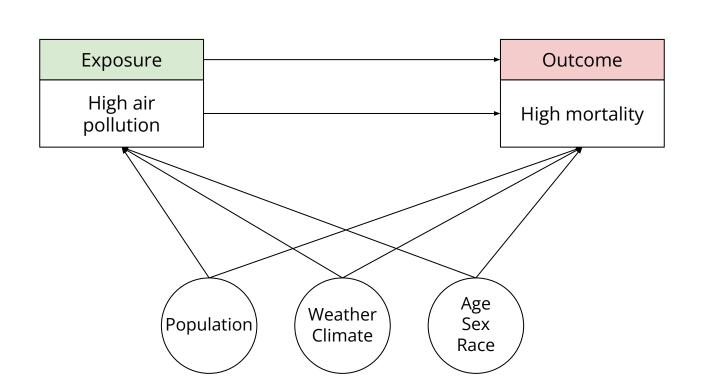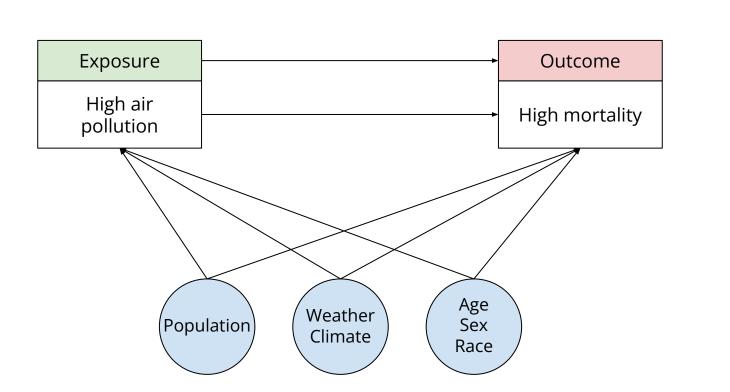| | January | February | March | April | May | June | July | August | September | October | November | December |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2000 | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS |
| 2001 | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS |
| 2002 | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS |
| 2003 | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS |
| 2004 | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS |
| 2005 | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS |
| 2006 | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS |
| 2007 | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS | GeoTIFF RDS |

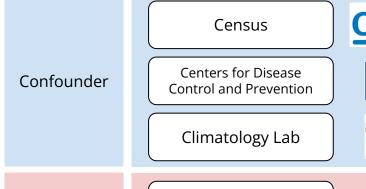# Outcome - hypothesized to have a causal relationship with exposure

- Medicare administrative data, also known as health services utilization data, are collected by the Centers for Medicare and Medicaid Services (CMS) and derived from reimbursement information or the payment of bills.
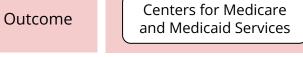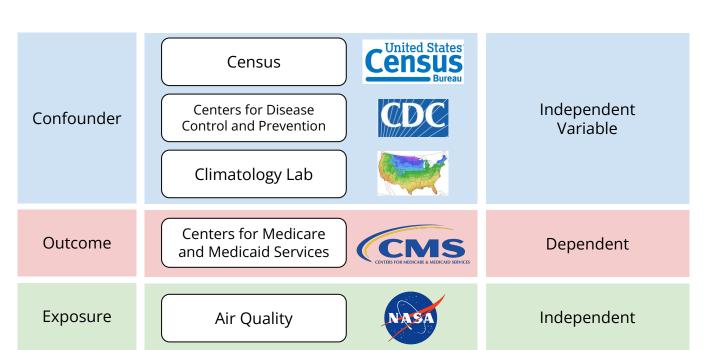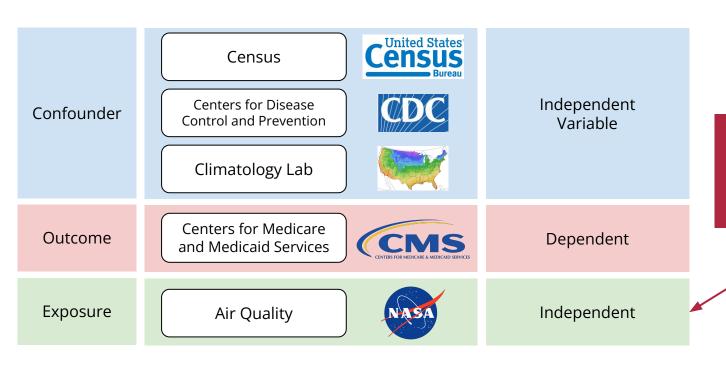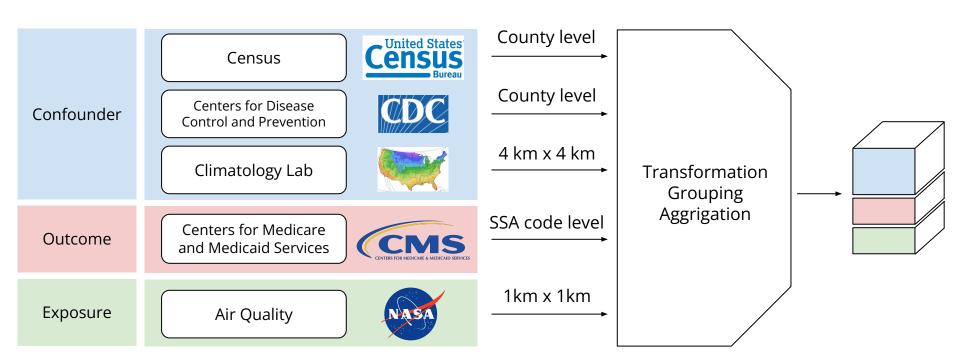  - Medical diagnosis
  - Mortality

```
┌─────────────────┐                              ┌─────────────────┐
│    Exposure     │─────────────────────────────▶│    Outcome      │
└─────────────────┘                              └─────────────────┘
```

```
  ┌─────────────────┐                              ┌─────────────────┐
  │    Exposure     │ ───────────────────────────▶ │     Outcome     │
  ├─────────────────┤                              ├─────────────────┤
  │    High air     │ ───────────────────────────▶ │  High mortality │
  │    pollution    │                              │                 │
  └─────────────────┘                              └─────────────────┘
            ▲                                                ▲
             ╲                                             ╱
              ╲                                          ╱
               ╲                                      ╱
                ╲           ┌──────────┐           ╱
                 ╲────────  │Population│  ────────╱
                            └──────────┘
```

# Troubles with data

- Temporal and spatial resolution
  - Different in different datasets
- Missing data
  - More-or-less every dataset
- Inconsistent data
  - I.e., single person having inconsistent information on age/sex/race
- Badly formatted data
  - I.e., medical, satellite

# Missing data

- Exclude records with missing values if they measure subpopulation, exposure or outcome
  - If there are few records with missing values (<5% of records)
- Don't exclude if they measure a confounder
- If there are over 5% of the records with missing values on subpopulation, exposure, outcome or important confounder - rethink study design

# Data sources

Summary

- Hypothesis-driven epidemiological data analyses require data for exposures, outcomes and confounders
- Secondary data analysis often require data cleaning, transformation and aggregation

A

B

C

A.0 2011
A.1 2012
A.2 2013
A.3 2014
A.4 2015
A.5 2016
A.6 2017

1 km

1 km

1 km

1 km

time

# NetCDF format

- Dimensions
- Variables
- Data
- Metadata

netCDF

# NetCDF format

- Dimensions
- Variables
- Data
- Metadata

# NetCDF format

- Dimensions
- Variables
- Data
- Metadata

# NetCDF format

- Dimensions
- Variables
- Data
- Metadata



```
VI_terra

<xarray.Dataset>
Dimensions:                                (lat: 134, lon: 182, time: 345)
Coordinates:
  * time                                   (time) object 2004-12-18 00:00:00 ... 2019-12-19 00:00:00
  * lat                                    (lat) float64 21.75 21.75 ... 21.2
  * lon                                    (lon) float64 -158.3 ... -157.6
Data variables:
    crs                                    int8 ...
    _500m_16_days_EVI                      (time, lat, lon) float32 ...
    _500m_16_days_MIR_reflectance          (time, lat, lon) float32 ...
    _500m_16_days_NDVI                     (time, lat, lon) float32 ...
    _500m_16_days_NIR_reflectance          (time, lat, lon) float32 ...
    _500m_16_days_VI_Quality               (time, lat, lon) float64 ...
    _500m_16_days_blue_reflectance         (time, lat, lon) float32 ...
    _500m_16_days_composite_day_of_the_year (time, lat, lon) float32 ...
    _500m_16_days_pixel_reliability        (time, lat, lon) float64 ...
    _500m_16_days_red_reflectance          (time, lat, lon) float32 ...
    _500m_16_days_relative_azimuth_angle   (time, lat, lon) float32 ...
    _500m_16_days_sun_zenith_angle         (time, lat, lon) float32 ...
    _500m_16_days_view_zenith_angle        (time, lat, lon) float32 ...
Attributes:
    title:        MOD13A1.006 for aid0001
    Conventions:  CF-1.6
    institution:  Land Processes Distributed Active Archive Center (LP DAAC)
    source:       AppEEARS v2.40
    references:   See README.txt
```

# Data analysis toolbox

# Data analysis toolbox

# Data analysis toolbox
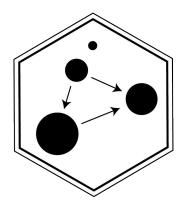
# Data analysis toolbox



https://fasrc.github.io/CRE/

https://fasrc.github.io/CausalGPS/

# Data processing and analysis

Summary

- Big data analysis can be conducted using solely free and open source software!
- NetCDF file format is great when working with high-dimensional geospatial datasets

How to conduct a big data analysis on air pollution and health?

# Capturing the data pipeline

- Review
- Verification
- Collaboration

"An article about computational science in a scientific publication is **not** the scholarship itself, it is merely **advertising** of the scholarship. The actual scholarship is the software, [data] ... and set of instructions which generated the figures. "

~ Prof Claerbout

Buckheit. 1995. Wavelab and reproducible research.

# Automation

- Fast(er) analysis
- Troubleshooting
- Reuse
- Education and training

```
JOB A A.submit
JOB B B.submit
JOB C C.submit
PARENT A CHILD B
PARENT B CHILD C
```

```
JOB     A   A.submit
JOB     B   B.submit
JOB     C   C.submit
PARENT  A CHILD  B
PARENT  B CHILD  C
```

```
[ant746@rce6-5:~/shared_space/ci3_ant746/mcbs-medpar-mbsf (master) $ cat workflow.dag
JOB A A.submit
JOB B B.submit
JOB C C.submit

PARENT A CHILD B
PARENT B CHILD C
[ant746@rce6-5:~/shared_space/ci3_ant746/mcbs-medpar-mbsf (master) $ condor_submit_dag workflow.dag

-----------------------------------------------------------------
File for submitting this DAG to Condor          : workflow.dag.condor.sub
Log of DAGMan debugging messages                : workflow.dag.dagman.out
Log of Condor library output                    : workflow.dag.lib.out
Log of Condor library error messages            : workflow.dag.lib.err
Log of the life of condor_dagman itself         : workflow.dag.dagman.log

Submitting job(s).
1 job(s) submitted to cluster 84471.
-----------------------------------------------------------------
```
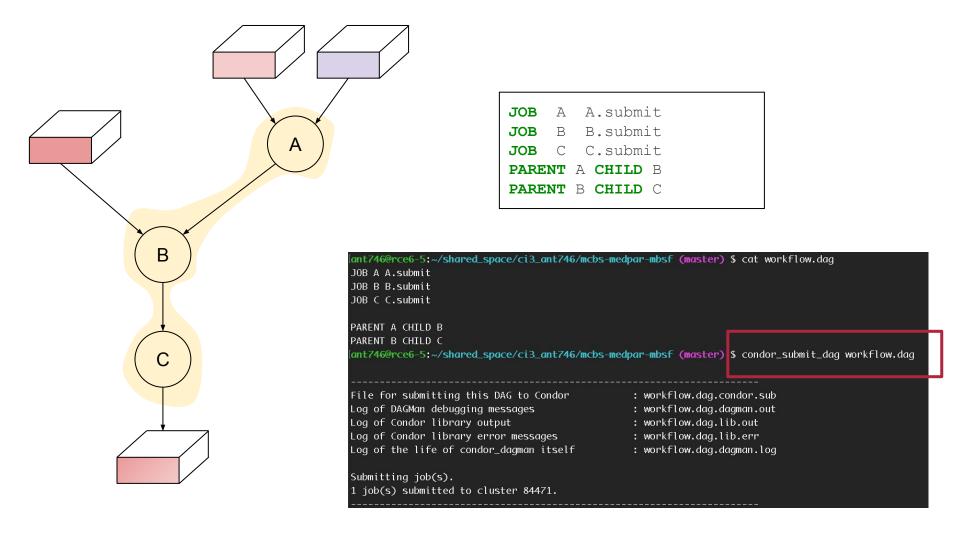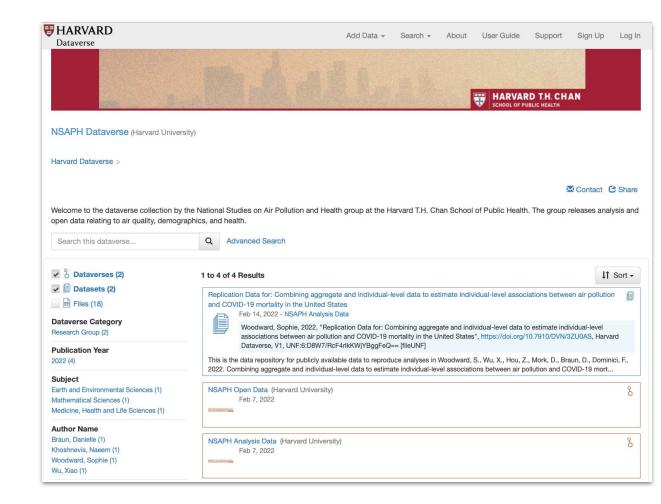
# Workflow engines

- A free and open-source software platform to archive, share, and cite research data
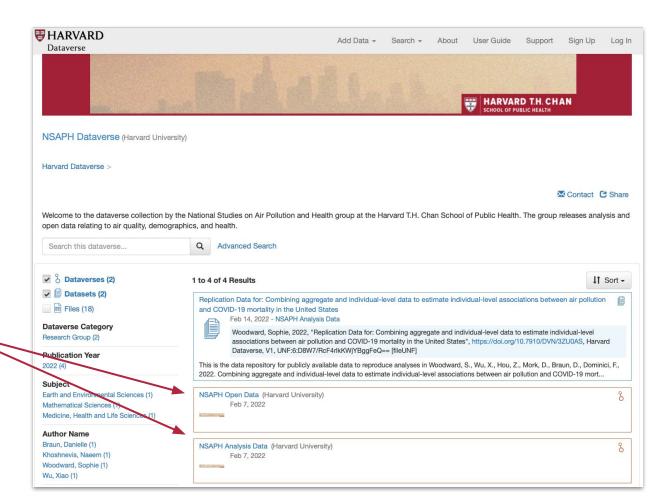  - Focus on data sharing and making data available

# 78 institutions around the globe run Dataverse installations as their official data repository



dataverse.org

# Data sharing

# Data sharing



Data collections

# Data sharing

- Data should be licensed
- Metadata
- It should be complete
- It should be shared in a (free, open) machine-readable format

# Dissemination

Summary

- Sharing of data, code and computational processes is necessary due to the requirements of policy makers, journals, funding agencies.

Email: anatrisovic@g.harvard.edu
GitHub & Twitter: atrisovic