



# Air pollution in Berlin

Using machine learning algorithms to model atmospheric aerosol particles in Berlin

# Air pollution

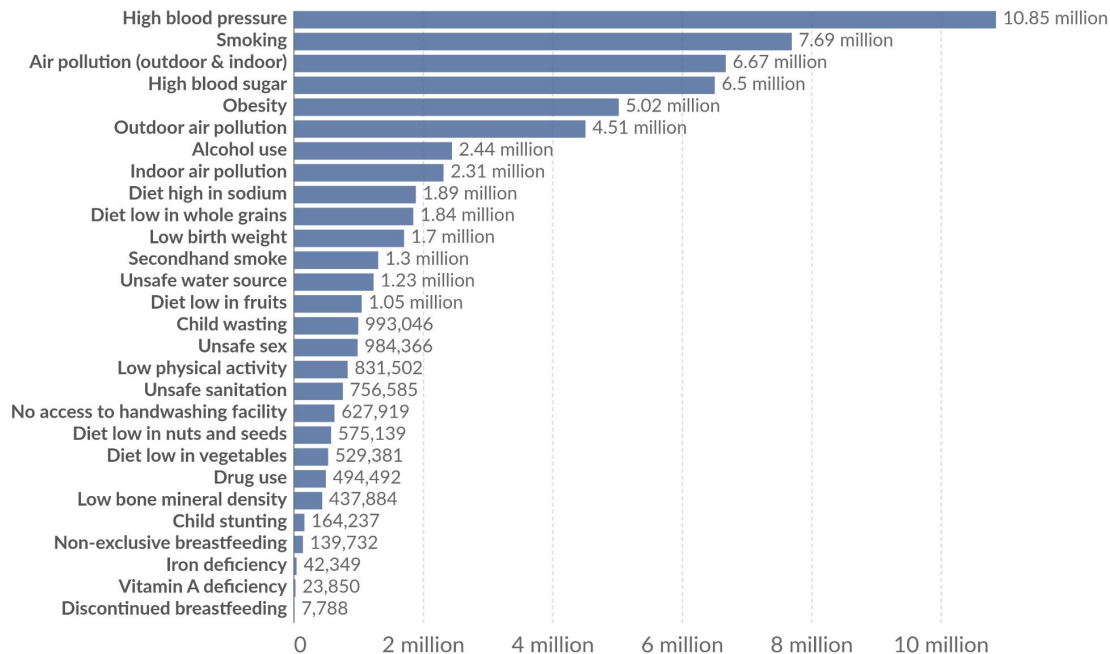
#3 death cause globally

- 01 | stroke
- 02 | heart disease
- 03 | lung cancer
- 04 | respiratory diseases (asthma)
- 05 | infertility

## Number of deaths by risk factor, World, 2019

Total annual number of deaths by risk factor, measured across all age groups and both sexes.

Our World  
in Data



Source: IHME, Global Burden of Disease (2019)

OurWorldInData.org/causes-of-death • CC BY

# PM 2.5 concentration

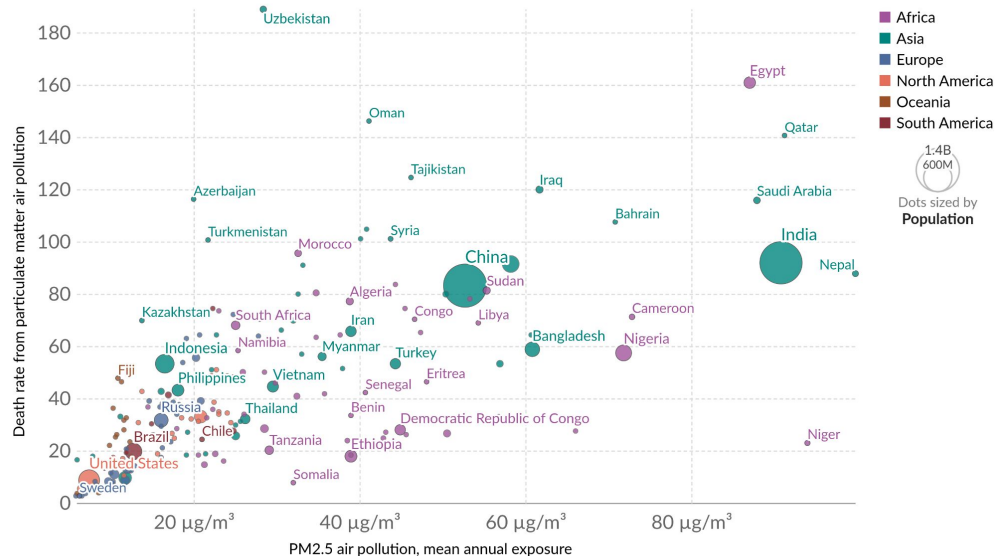
## Particulate matter

- Atmospheric aerosol particles with diameter  $<2.5\text{ }\mu\text{m}$
- Microscopic particles suspended in the air
- Respirable (penetrate into lungs)

## Death rate from particulate matter air pollution vs PM2.5 concentration, 2017

Age-standardized death rate from particulate matter (PM2.5) exposure per 100,000 people versus the average mean annual exposure to particulate matter smaller than 2.5 microns (PM2.5), measured in micrograms per cubic meter.

Our World  
in Data



Source: IHME, Global Burden of Disease (2019); Brauer et al. (2017) via World Bank

OurWorldInData.org/air-pollution/ • CC BY

# WHO guidelines

Berlin: worst air quality in Germany

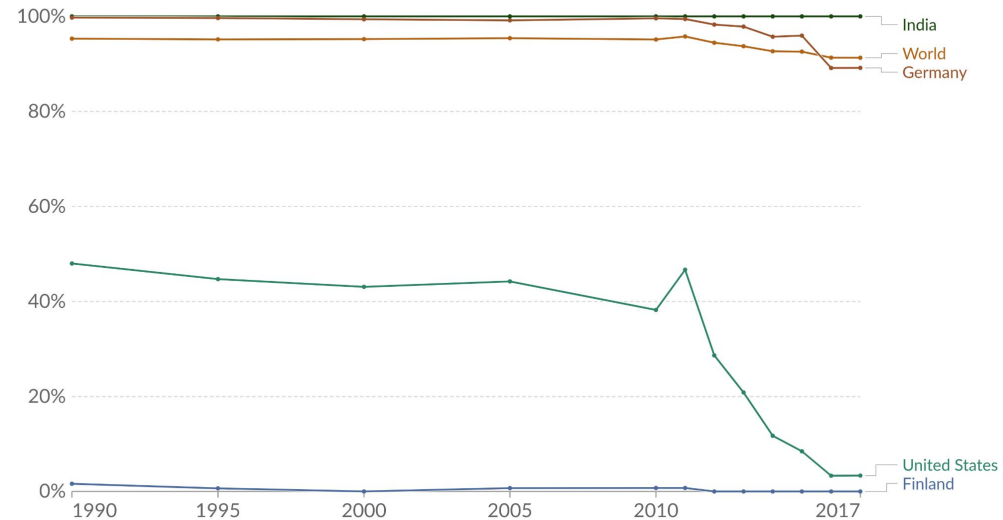


Source: Federico Gambarini/dpa

## Share of the population exposed to air pollution levels above WHO guidelines, 1990 to 2017

Our World in Data

The share of the population exposed to outdoor concentrations of particulate matter (PM<sub>2.5</sub>) that exceed the WHO guideline value of 10 micrograms per cubic meter per year. 10µg/m<sup>3</sup> represents the lower range of WHO recommendations for air pollution exposure over which adverse health effects are observed.



Source: Brauer et al. (2017) via World Bank

OurWorldInData.org/outdoor-air-pollution • CC BY



# Influencing factors



## Topography

Mountains, buildings, narrow streets,  
vegetation

## Weather

Temperature: summer vs. winter  
Wind speed: dilution / accumulation  
Wind direction: transport from source  
Precipitation



## Emission source

Natural vs. anthropogenic  
Distance from source:  
Traffic vs. volcanic eruption, Sahara dust





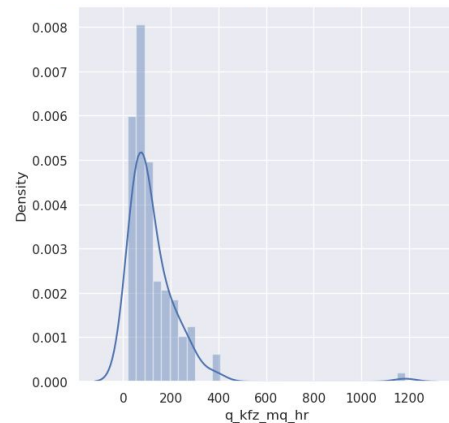
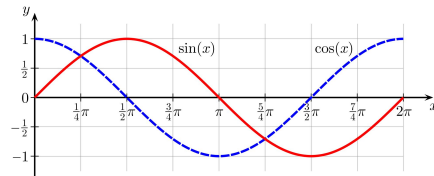
# Data: hourly time series

- **Air quality**
  - PM2.5 concentrations at measurement station in traffic (Mariendorfer Damm)
    - <https://luftdaten.berlin.de>
- **Features**
  - **Weather** data: air temperature, humidity, wind speed, wind direction, precipitation (Alexanderplatz)
    - <https://opendata.dwd.de>
  - **Traffic** data: quantity and velocity of cars per street (Mariendorfer Damm & distribution over other streets)
    - <https://api.viz.berlin.de/daten/verkehrsdetektion>
  - **Time** features: hour, weekday, month



# Preparation steps

- Clean data
  - Measurement errors
- Fill gaps
  - Not straightforward for e.g. cyclical data



- Feature engineering
  - Transform cyclical features: wind, time
    - sine & cosine
  - Lagged features of traffic data
  - Add traffic distribution stats as features:
    - mean, sd, min, max, skewness, kurtosis
- Feature Selection
  - Feature importance





# Models and results

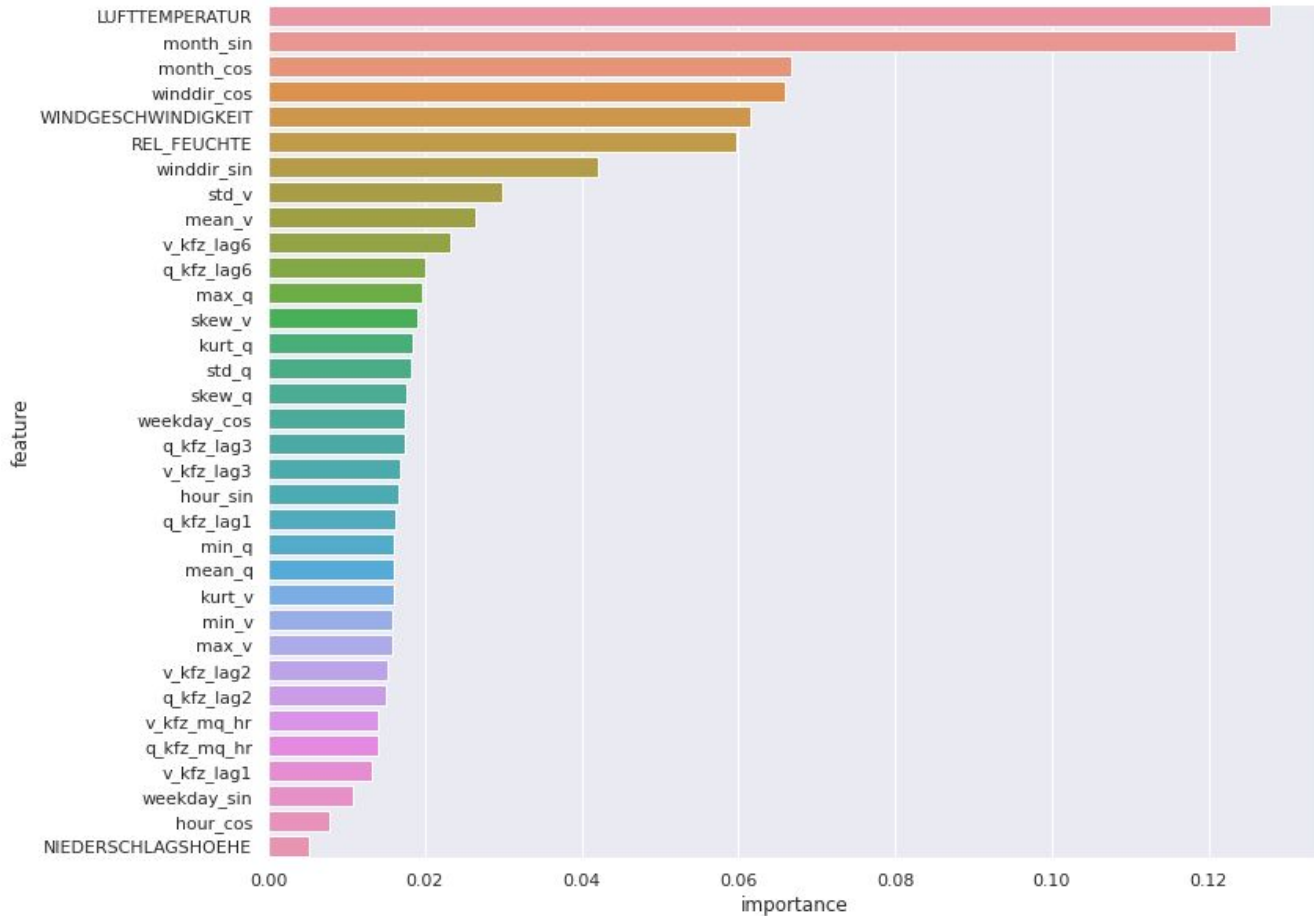
- Supervised learning – regression – time series
- GridSearchCV for all models
  - K-fold cross-validation
  - Hyperparameter tuning
- Linear Regression
  - Accuracy: train 0.29, test 0.28
- Linear regression with polynomial features
  - Accuracy: train 0.64, test 0.49
- Random Forest Regression
  - Accuracy: train 0.69, test 0.61





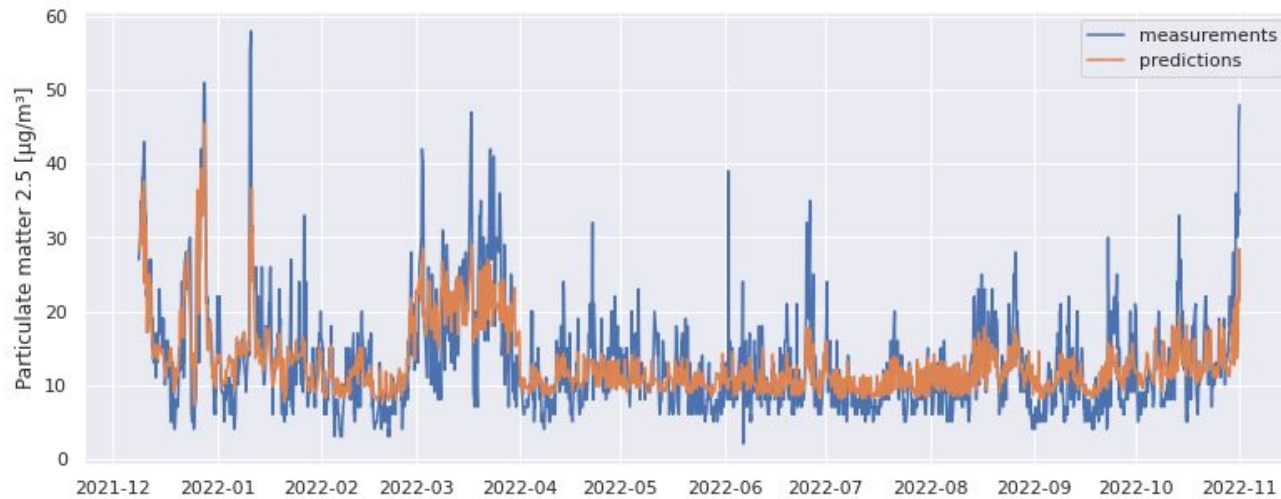
# Feature importance

drop features:  
importance < 0.02





# Results





# Implications

- Weather features most important
  - nonlinearity
- Traffic: speed and speed variability seems more important than number of cars
  - speed limits would improve air quality
- Limits of machine learning for environmental observational data
  - often have gaps, discontinuities, not long enough
  - need for relevant features (often not available)
- Given the limits of the data: RandomForest proved quite powerful





# Many thanks!