# Health Outcomes v. Warehouse Location

Something.

# Intro: Diesel Trucks are Bad For Babies


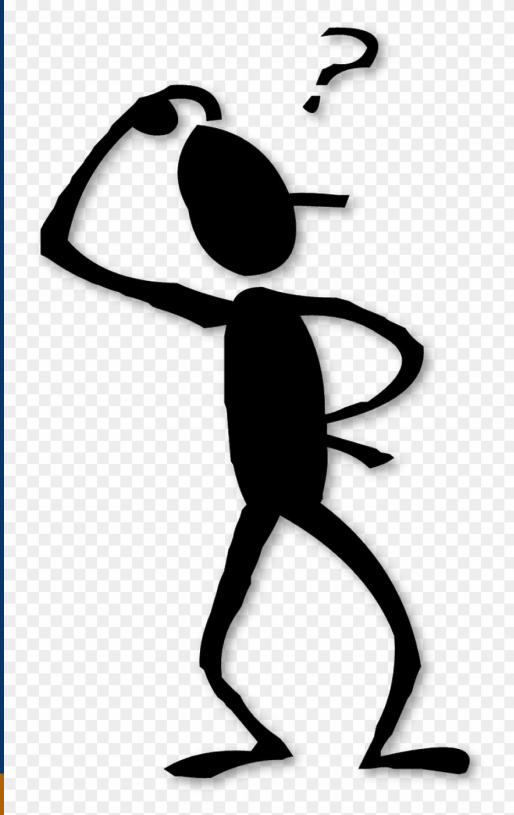
FEATURES DROPPED          SURPRISING FEATURES

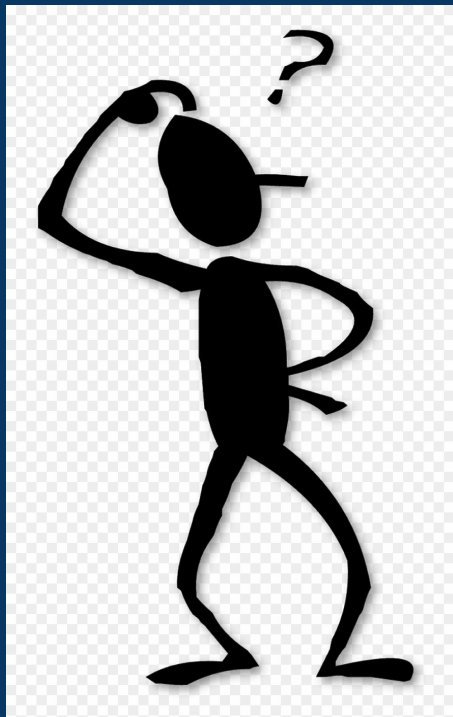**BLURB**                              **MANIPULATION**

# THE PROBLEM STATEMENT SLIDE

- **What is the (quantifiable) effect of increased warehouse presence in California in the last decade on emergency healthcare?**
- **How well do the CalEnviroScreen scores reflect emergency healthcare counts?**
- **What indicators from the CalEnviroScreen dataset best determine the number of emergency healthcare visits?**

# Aggregated Modeling With Additional Features



As Data Scientists in *OEHHA*, we are tasked with developing models aggregating the four time-points from each report with additional information on warehouse density to assess primary mitigating factors addressing negative health outcomes.

# Data Source:
# California EnviroScreen reports



## California Office of Environmental Health Hazard Assessment

From the **California Office of Environmental Health Hazard Assessment**

https://oehha.ca.gov/calenviroscreen

A series of four datasets and reports, published 20–, 20–, 20–, and 20–, with pollution, basic health, and socioeconomic measurements for each of California's zip codes or census tracts.

These measurements are compiled into a small number of summary scores, including a broad California EnviroScreen score indicating the regions with the most pressing needs.

# The CalEnviroScreen Model

EnviroScreen-specific "scores" are derived from measurements.
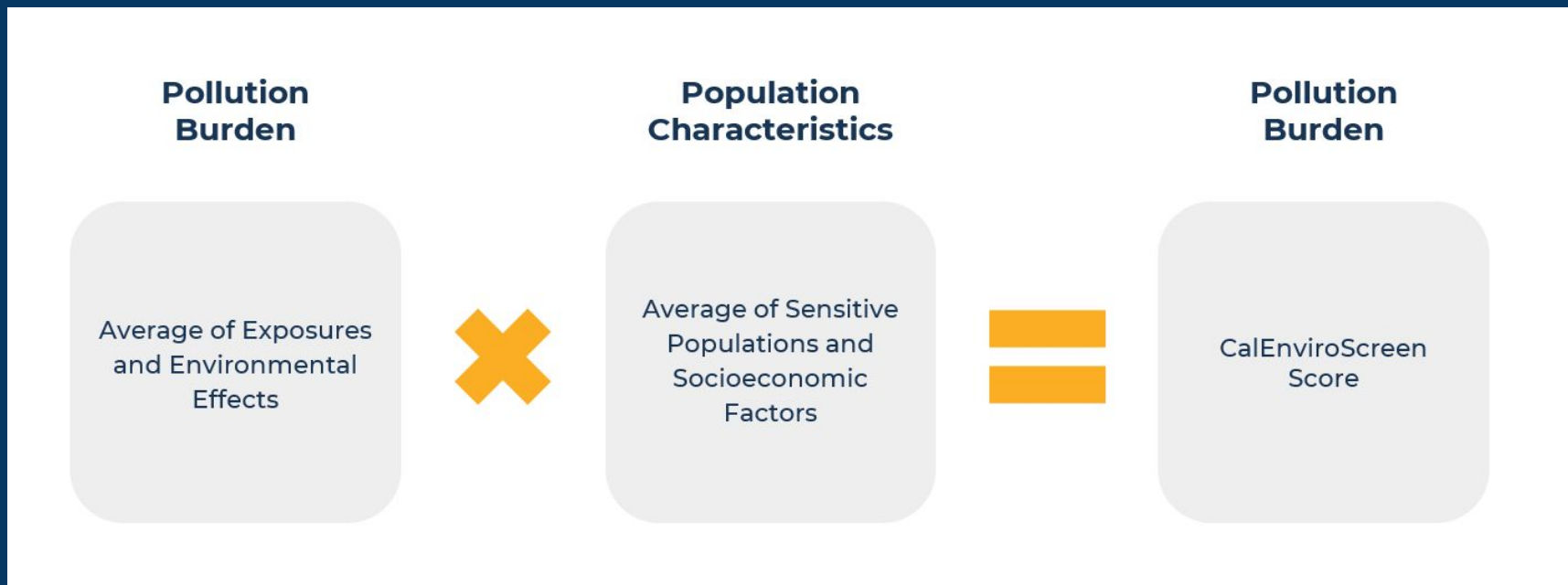
- Pollution Burden Score
  - Exposures
    - Ozone concentrations
    - Particulate matter emissions and concentrations (diesel, PM2.5)
    - Drinking water contaminants, lead risk
    - Toxic releases from facilities, pesticide use
    - Traffic density
  - Environmental Effects
    - Solid waste, sites
    - Groundwater threats and impaired water body count

# The CalEnviroScreen Model

EnviroScreen-specific "scores" are derived from measurements, also included in the dataset. Impact weights are determined by the CalEPA.
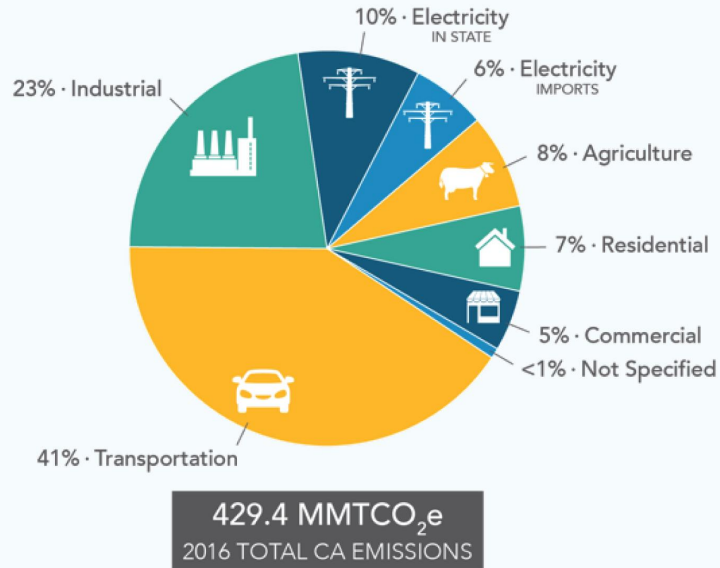
- Population characteristics
  - Sensitive population
    - Asthma
    - Cardiovascular disease
    - Low birth weight infants
  - Socioeconomic factors
    - Educational attainment
    - Housing burdened low income households
    - Linguistic isolation
    - Poverty
    - Unemployment

# The CalEnviroScreen Model

**Pollution Burden**

Average of Exposures and Environmental Effects

**✕**

**Population Characteristics**

Average of Sensitive Populations and Socioeconomic Factors

**=**

**Pollution Burden**

CalEnviroScreen Score

# Motivation: Assessing Effect of Emission Sources



10% · Electricity IN STATE
6% · Electricity IMPORTS
8% · Agriculture
7% · Residential
5% · Commercial
<1% · Not Specified
41% · Transportation
23% · Industrial

429.4 MMTCO$_2$e
2016 TOTAL CA EMISSIONS

**SOURCE:** CA Air Resources Board, 2018 GHG Emission Inventory, July 2018
https://www.arb.ca.gov/cc/inventory/data/data.htm

# Motivation: Emissions Exceptions

**Your vehicle does not need a smog inspection if your:**

- Gasoline-powered vehicle is a 1975 year model or older (This includes motorcycles and trailers.)

- Diesel-powered vehicle is a 1997 and older year model OR with a Gross Vehicle Weight of more than 14,000 pounds.

- Powered by natural gas and weighs more than 14,000 pounds.

- An electric vehicle.

- Gasoline-powered and less than eight model-years old.

*SOURCE: CA DMV*
dmv.ca.gov/portal/vehicle-registration/smog-inspections/



*1995 Freightliner for Sale in San Rafael, Marin County, San Francisco Bay Area, CA*

# Motivation: Landscape Changes



SOURCE:
riversidewarehouses.com/listings/1020-prosperity-way-beaumont-ca-92223

# EDA with with health, pollution, and Poverty(SB)



Asthma vs Pollution with refrence to Disadvantaged community

**BLURB**

-
-
-
-
-
-
-
-
-

**MANIPULATION**

12

# EDA visuals cont. (marshall)



Poverty vs Asthma

BLURB

MANIPULATION

# Most important features for health cont.

**BLURB**

**MANIPULATION**

# EDA : Census data info – warehouse counts

Broad statistics on warehouse counts

**MANIPULATION**                    SURPRISING FEATURES

Joined onto
CAES data by tract
or zip

Changes with time

Joined onto
CAES data by tract
or zip

BLURB

Time by zip — board warehouse business changes.
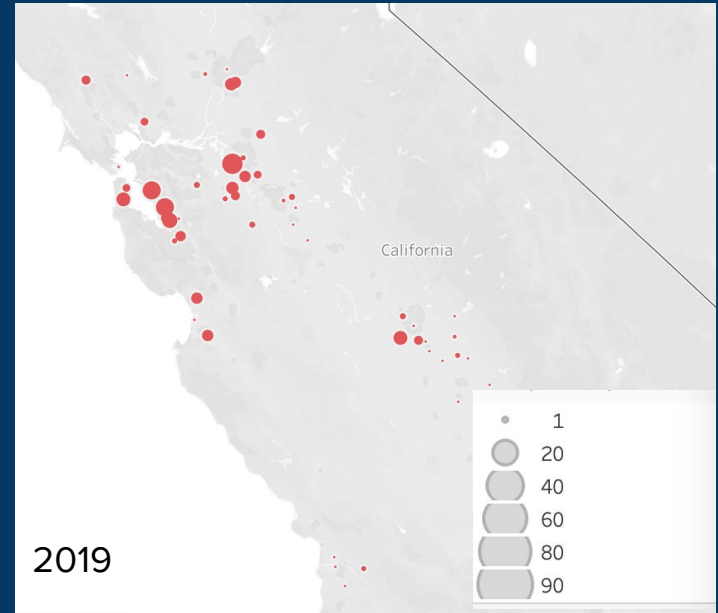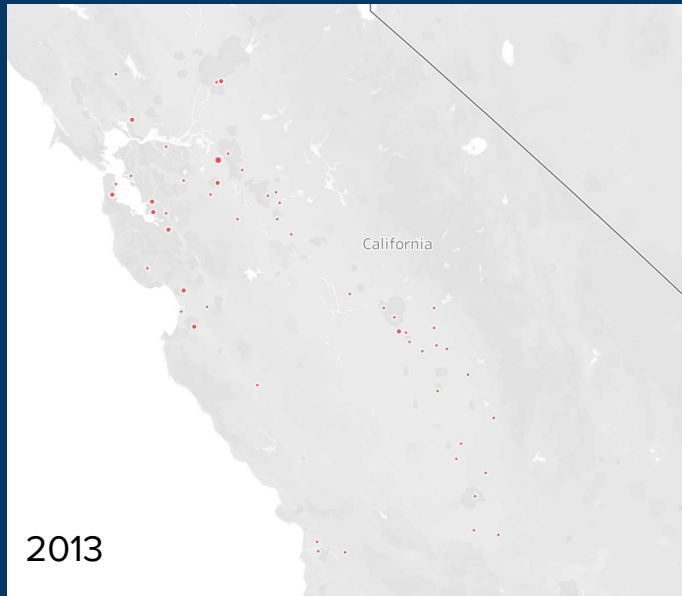
Or, just time with california as a whole.
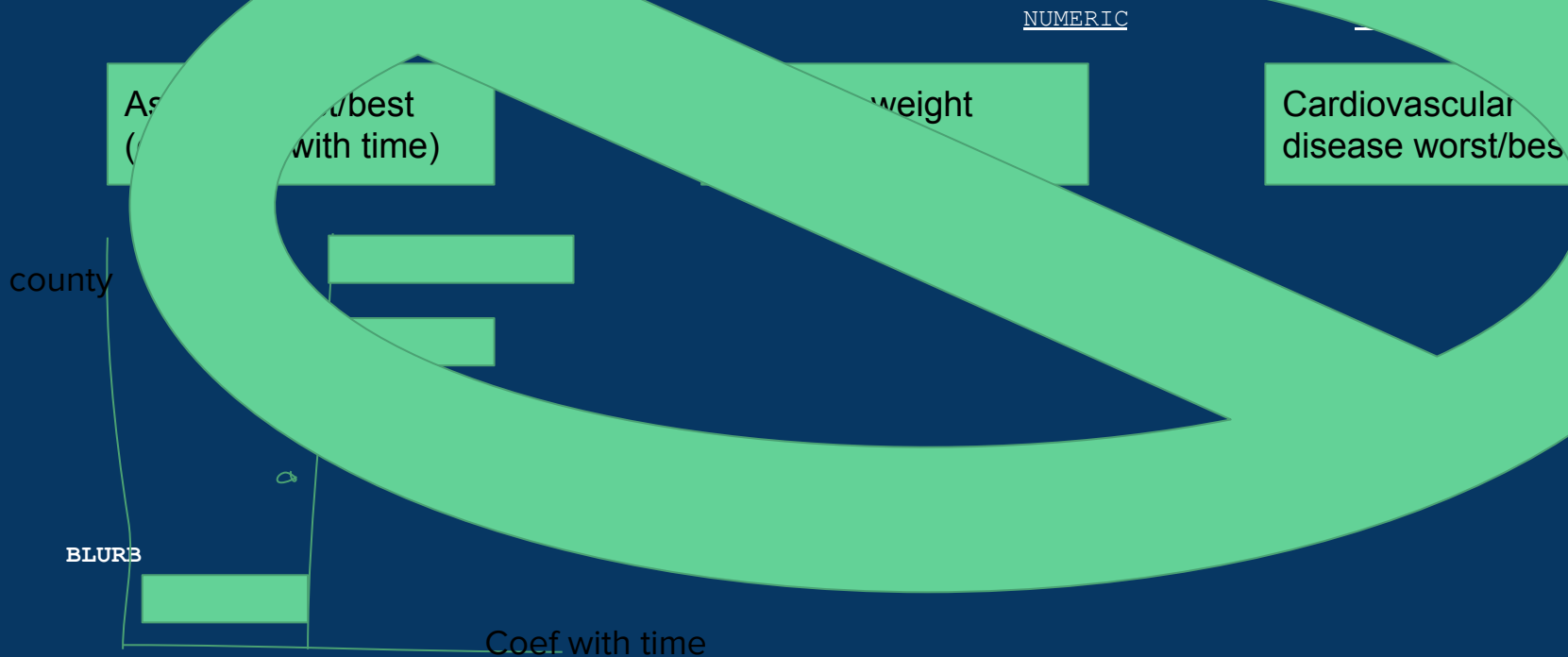
MAP

# Time by zip — board warehouse business changes.



2013



2019

Legend:
- 1
- 20
- 40
- 60
- 80
- 90

# Time by zip — board warehouse business changes.



2013



2019

|  | 1 |
|---|---|
|  | 20 |
|  | 40 |
|  | 60 |
|  | 80 |
|  | 90 |

# Time by zip — wh... ? Health or pollution...

NUMERIC

As... /best
(... with time)

...weight

Cardiovascular
disease worst/bes...

county

**BLURB**

Coef with time

These highlight regions that may be trouble soon, or rapidly improving. (eda should show bad/good *currently*)

# Time by zip — what are the biggest changers? Health or pollution

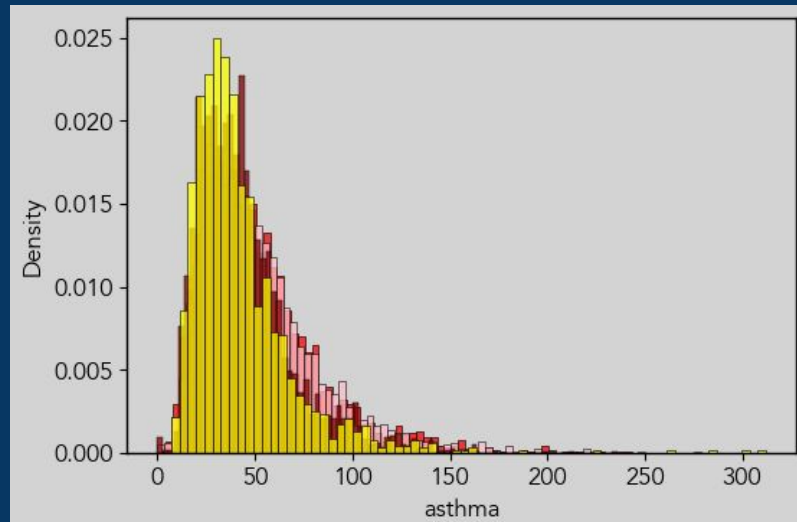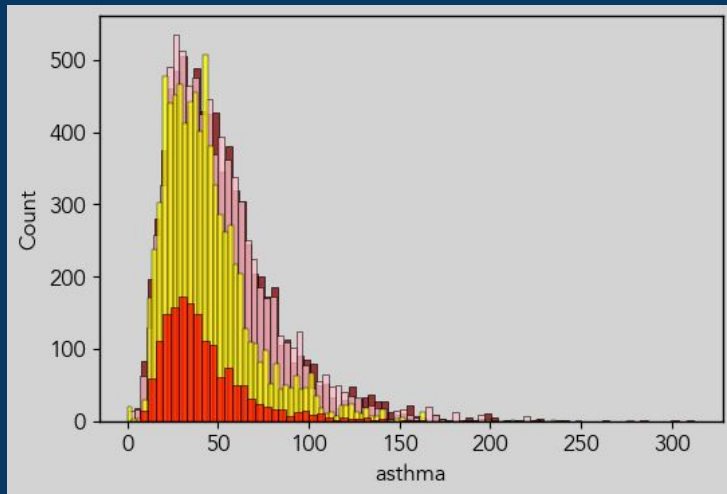Only fitting four values for each county: caes 1, 2, 3, 4 years.

Colored map

These highlight regions that may be trouble soon, or rapidly improving. (eda should show bad/good *currently*)
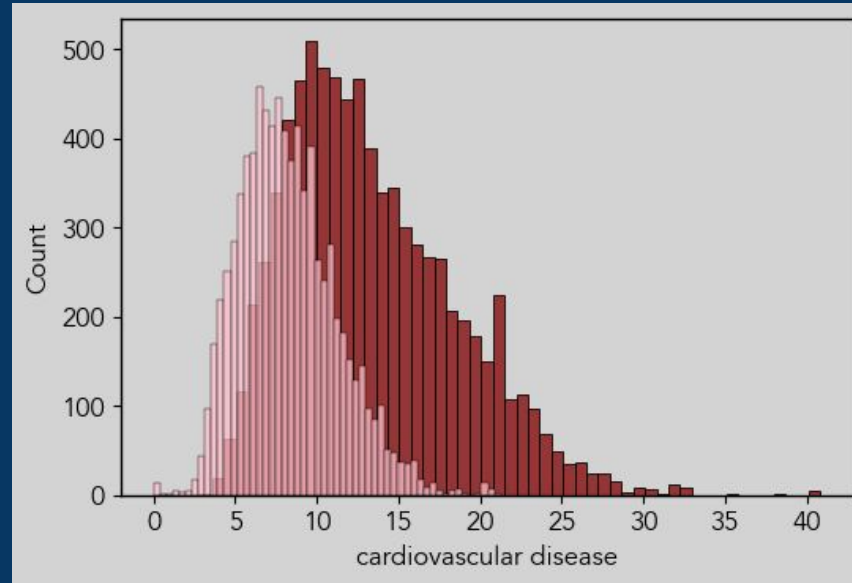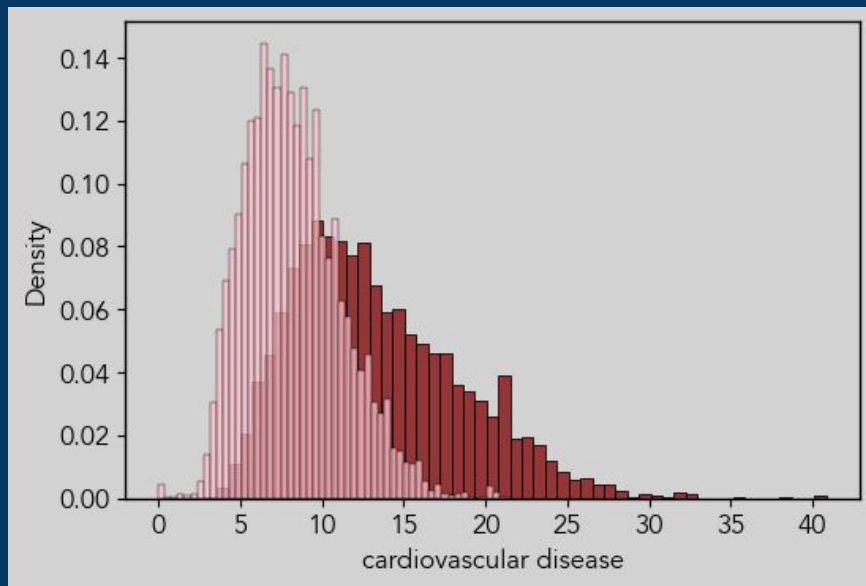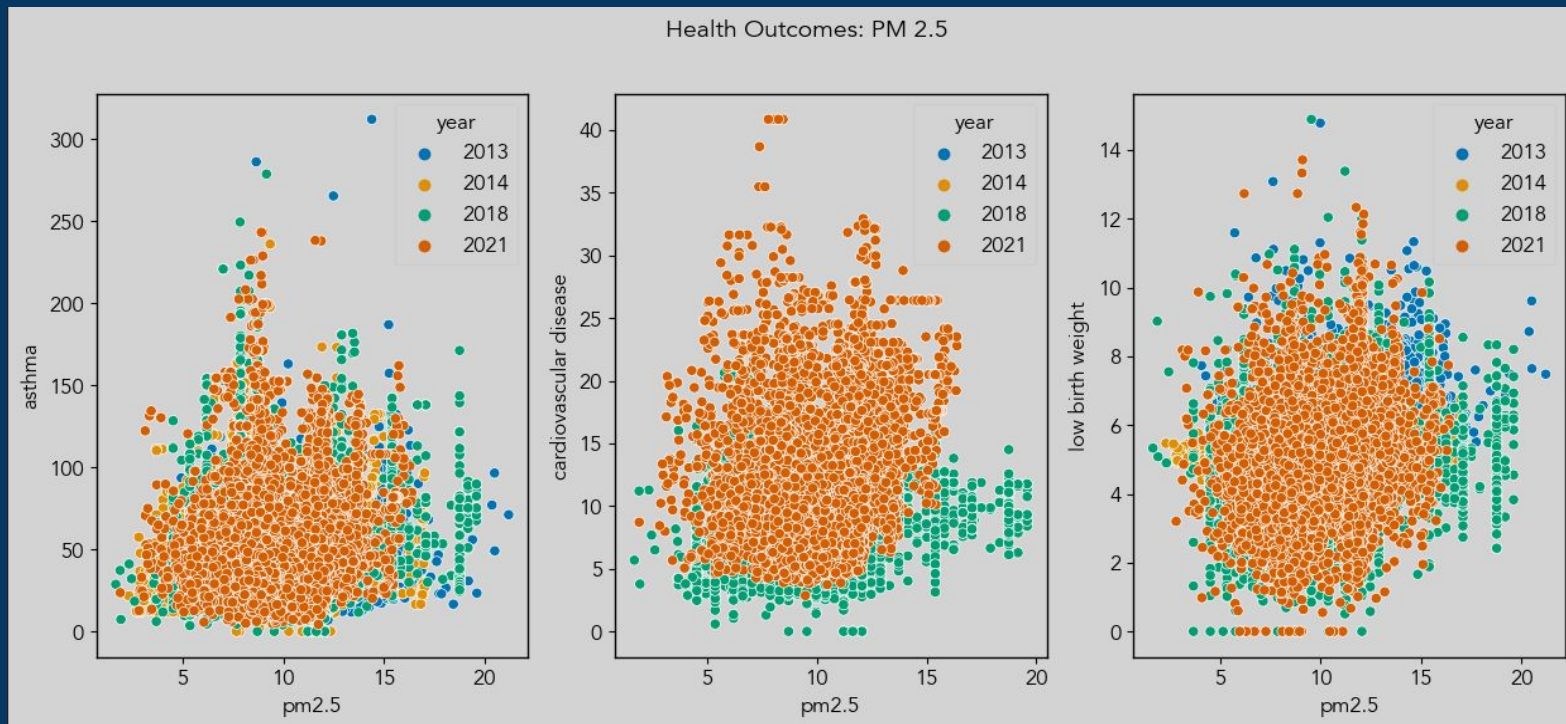
# EDA: Low-Birth Weight

# EDA: Asthma

# EDA: Cardiovascular

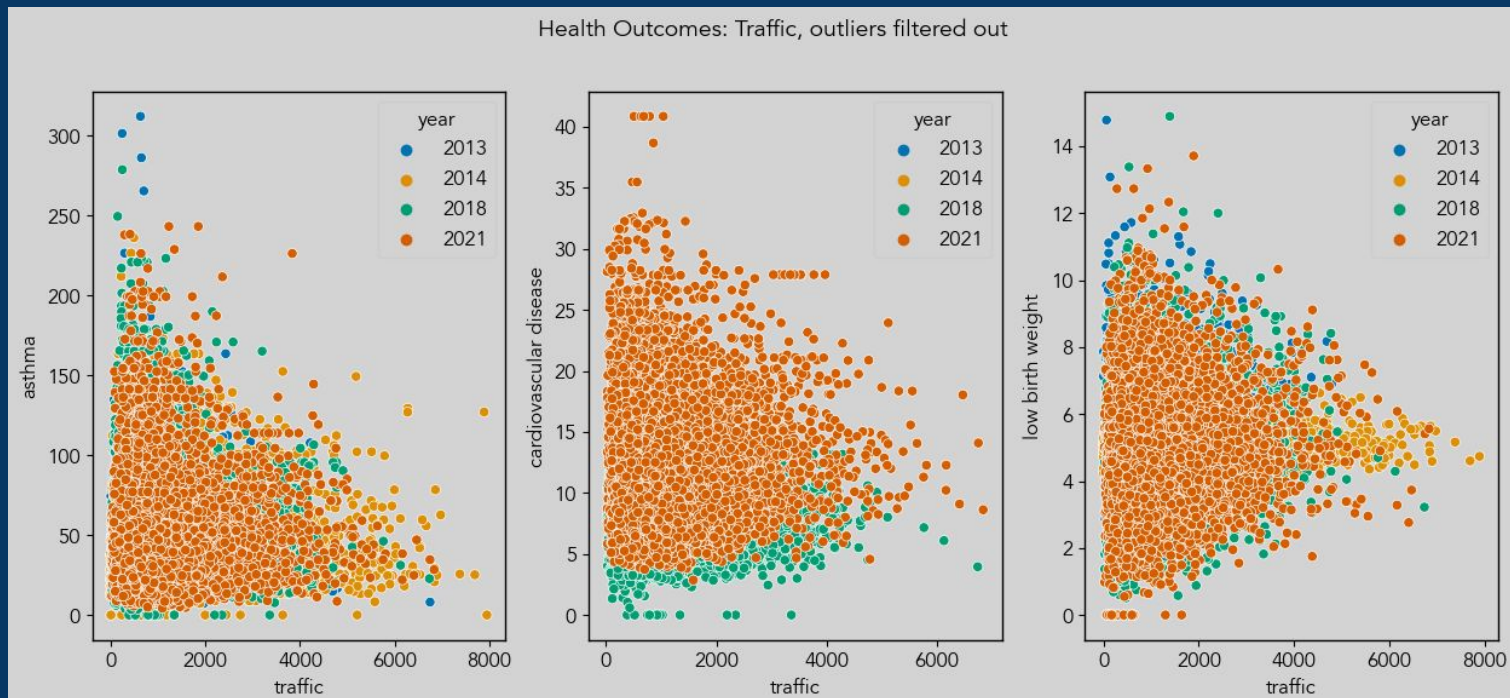# EDA: Health Outcomes, PM 2.5



Health Outcomes: PM 2.5

# EDA: Health Traffic Volume



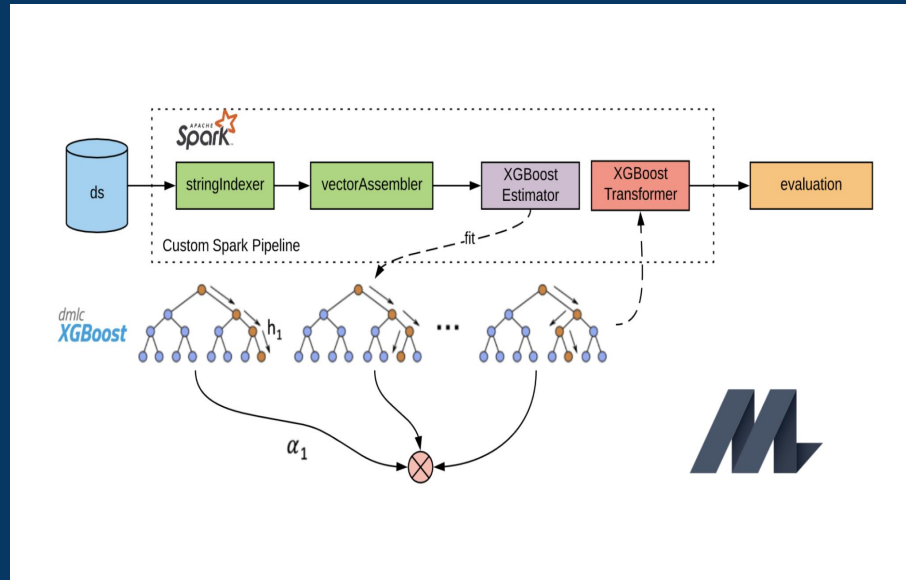Health Outcomes: Traffic, outliers filtered out

# EDA: Health Outcomes, Diesel PM



Health Outcomes: Diesel PM, outliers filtered

# XGboost, scaled, highest correlated features for Asthma

**BLURB**

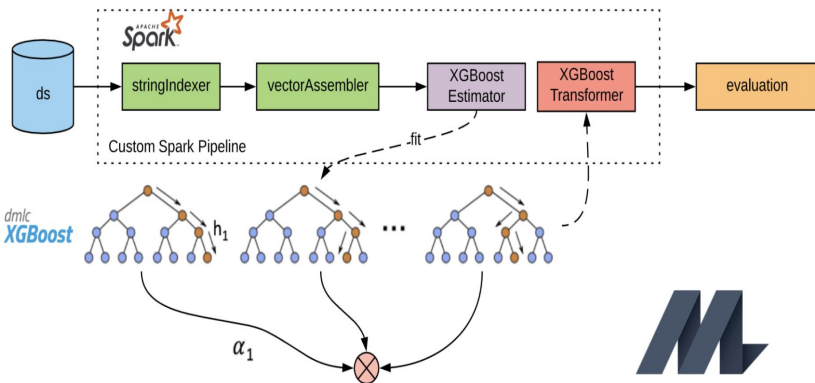- **XGBoost** can be used directly for **regression predictive modeling**.

-

**FINAL METRICS**

Train Accuracy:0.9472151826696329

Test Accuracy:0.7639090300436409

RMSE score:14.376141

# XGboost GS CV fit to best params Asthma

**BLURB**
-

**FINAL METRICS**
Train Accuracy:0.9472151826696329
Test Accuracy:0.7639090300436409
RMSE score:14.376141

# Model 2(marshall) Random Forest Reg



NUMERIC                    CATEGORICAL

**BLURB**                              **FINAL METRICS**

# Linear model : time and space only

NUMERIC                    CATEGORICAL

Year
Latitude
longitude

**BLURB**

**FINAL METRICS**
**Asthma -** (0.042, 0.060)
Low birth weight - (0.007, 0.010)
Cardiovascular disease - (0.38, 0.38)

# Linear model: CAES score features only



NUMERIC                    CATEGORICAL

**BLURB**

Do this to evaluate the CAES scores. "Have they done
the feature engineering for us already?"

**FINAL METRICS**

# Linear model: "Selected columns" (pollution and industry info. Most of the features.)

**BLURB**
**No year. No space. No CAES specific scores.**

**FINAL METRICS**
**r^2**
**Asthma –** 0.593
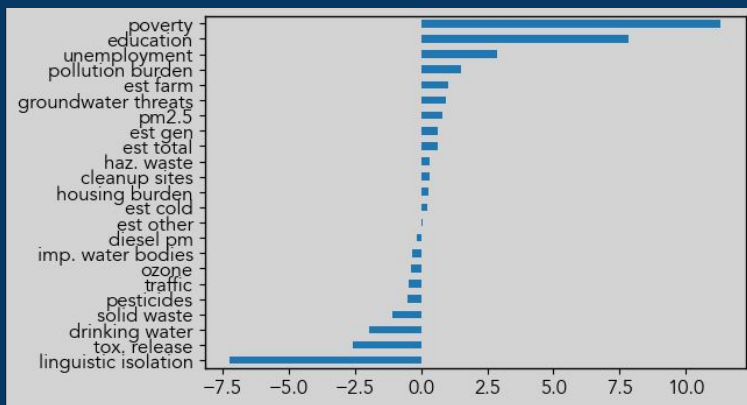Low birth weight - 0.385
Cardiovascular disease - 0.497

# More on linear "selected features."

Which features were most influential on the linear scale? (not svd/PCR, but just relative to scaled data. This is coefficient relative to scaled data.)

**FINAL METRICS, R^2:**
**Asthma** - 0.593
**Low birth weight** - 0.385
**Cardiovascular disease** - 0.497

# Model: SVR

**Epsilon-Support Vector Regression**

*regularization: L2, C = 1*

*Feature Importances: really highlights*
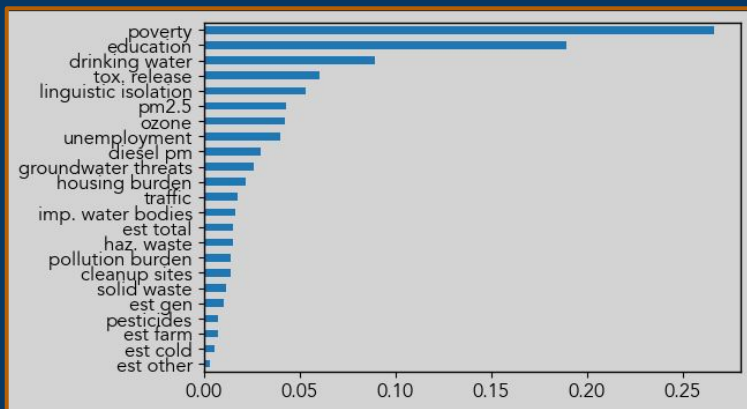
**FINAL METRICS**

# Model: Random Forest Regression

*n-estimators = 100*     *max_leaf_nodes = 10*
*max_depth = 10*         *max_features : auto*

NUMERIC                              CATEGORICAL



*different importances:*                    FINAL METRICS

Twenty breeds of dogs

keep this slide

DOGS
DOGS
DOGS