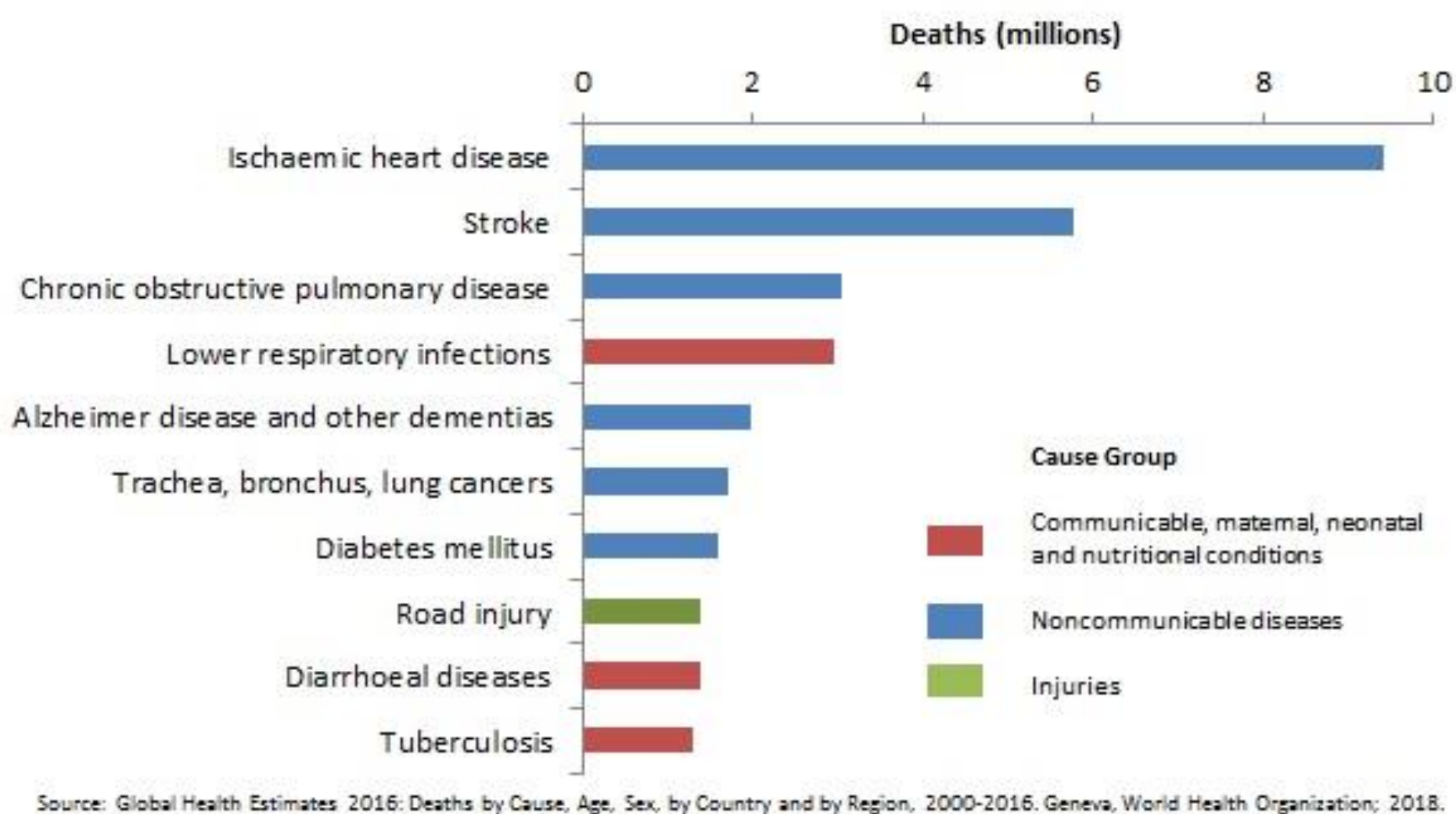


# Visualization of relationship between chronic diseases and preventions in 500 US Cities

Bu Qianqian - Myles Lefkovitz - Salim Noorallah Ladak - Tri Nguyen

## MOTIVATION

Top 10 global causes of deaths, 2016



Source: Global Health Estimates 2016: Deaths by Cause, Age, Sex, by Country and by Region, 2000-2016. Geneva, World Health Organization; 2018.

Diseases and health conditions represent 9 of the 10 most common causes of death globally. Of those 10 causes, chronic diseases (non-communicable diseases) represent 6 (blue bars). Understanding of related factors is of utmost importance to effective public health planning.

## Chronic Diseases

**Data Source.** Published by the CDC (Center for Disease Control and Prevention)

- Contain prevalence of 13 diseases, 9 prevention practices and 5 unhealthy behaviors in 500 cities at **National, State, City, and Census Tract levels.**
- CSV format, 800,000 rows of data/235MB.
- Contain geographic data for each area.

## THE DATA

**Goals.** Study which preventions and behaviors best predict a particular disease at state and national levels.

**Processing.** Data is cleaned and reformatted in Python, using Pandas and Numpy libraries.

## ANALYSIS

### Questions

- Interpretable Relationship**
- Between health outcomes and preventions/unhealthy behaviours?
  - Prediction accuracy?
  - Which model works best?
  - What are the key features?

### Experiment Design

- ML with Scikit-Learn**
- Multiple model comparison analysis
  - Top feature selection at national and state levels
  - Generate output file for visualization

### Algorithms Evaluation

- Multiple Regression Models**
- Linear regression
  - Ridge regression
  - Lasso regression
  - Support vector regression (SVR)

### Results

- Linear Relationships**
- Found top 5 predictors for each health outcome
  - Formulated SVR model with hyperparameter tuning
  - Test accuracies averaged 0.89 at national level

## VISUALIZATION

## INTERACTIVE MAPS

Developed in Tableau.  
Show most prevalent diseases in each city/state.  
Most related factors to each disease.  
Calculation at **state** and **national** level.

## RESULTS

**Best predicting model:** Support Vector Regression  
**Most prevalent diseases:** high cholesterol and high blood pressure.

**Most important factor:** blood pressure medication

**Visualization** enabling users to explore the findings and make comparison across diseases and cities/states.

