

Chad Bradford

Peyton Lambourne

Anshi Mathur

James Segovia

Developing Models to Predict Life Expectancy Using Data from the WHO

PROJECT 2 PRESENTATION

March 31, 2025

OUR GOALS:

Our Group Project 2 focuses on harnessing our acquired skill set to assess and clean a rich data set that enables a build of strong predictive models for life expectancy.

Hypothesis: Health, social, and/or economic factors can be used to model life expectancy outcomes with high predictability.

01

IDENTIFY AN INFORMING DATA SET

- Relevance
- Data set quality
- Size and balance
- Feature diversity (without leakage)
- Current
- Ethics

02

ASSESS AND CLEAN THE DATA SET

- Routine cleaning
- Identified and dropped leaking features
- Used lag features to model temporal trends
- Imputed missing values

03

MODEL BUILDING AND TESTING

- Model build and testing
 - Classification
 - Regression

Classify Life Expectancy

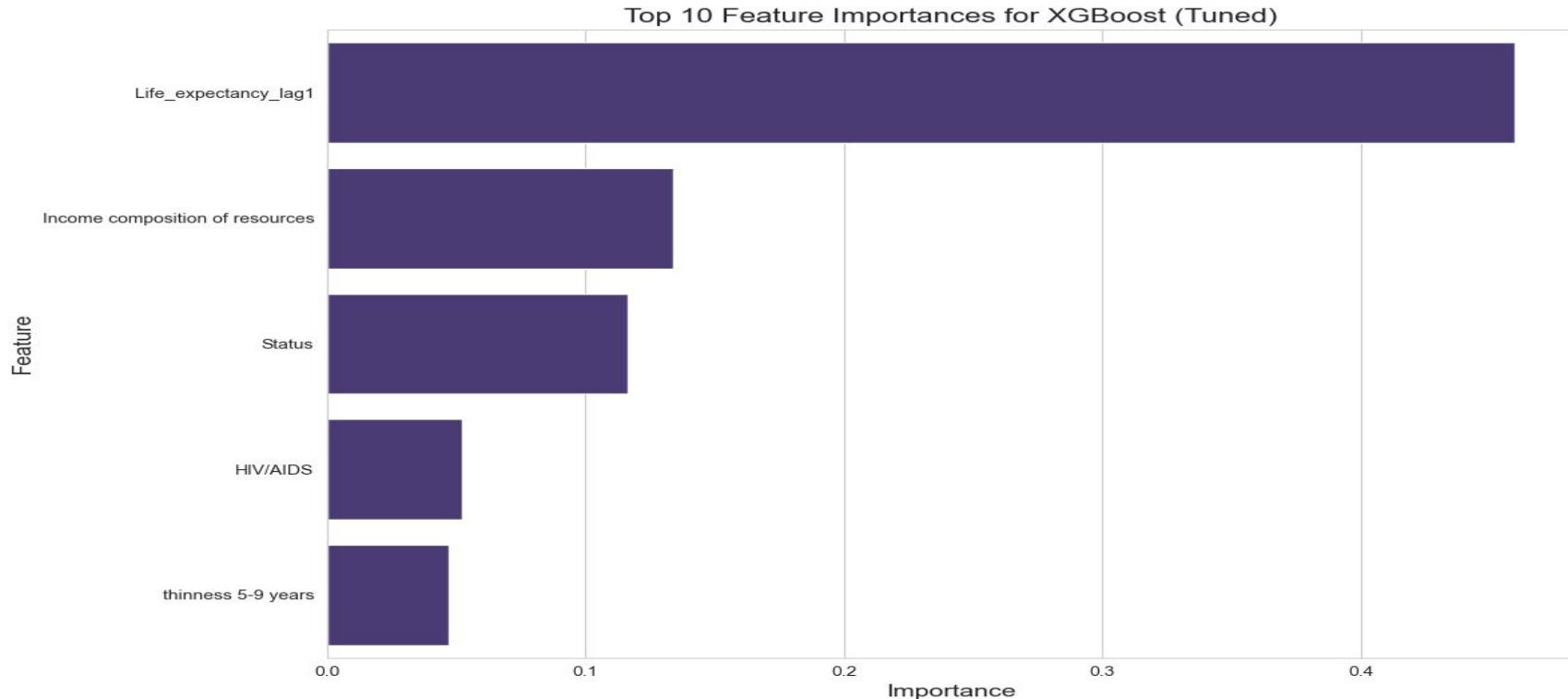
Let's get the boring stuff out of the way





The graph made me do it...

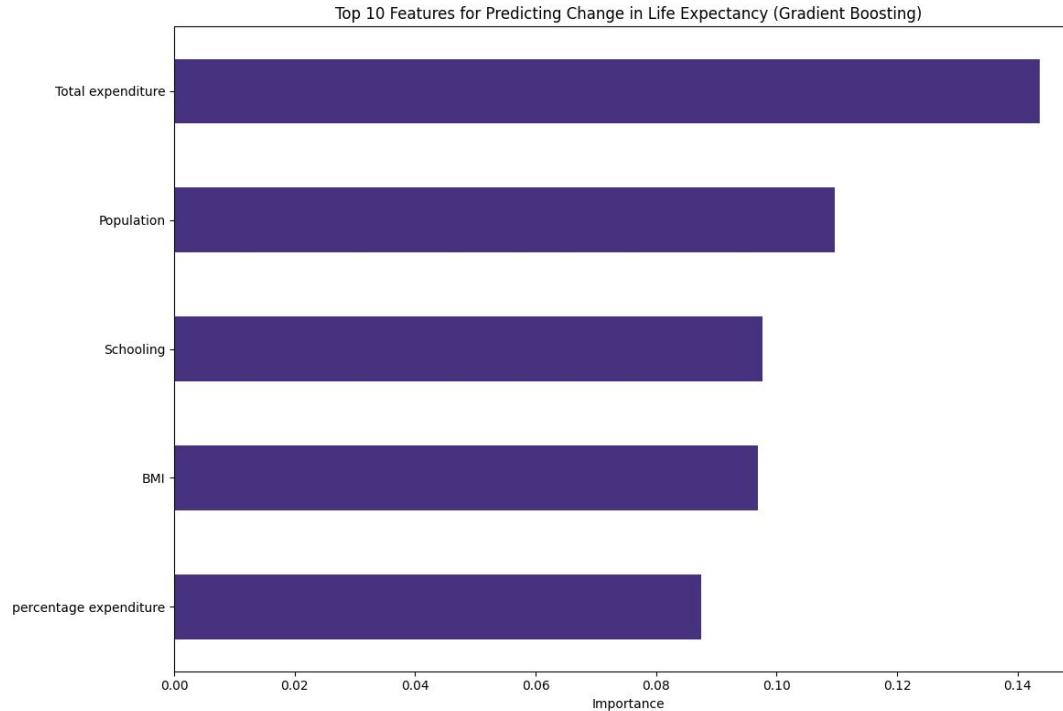
Wow!! Who could have predicted this?



GET ME OFF THIS PLANE

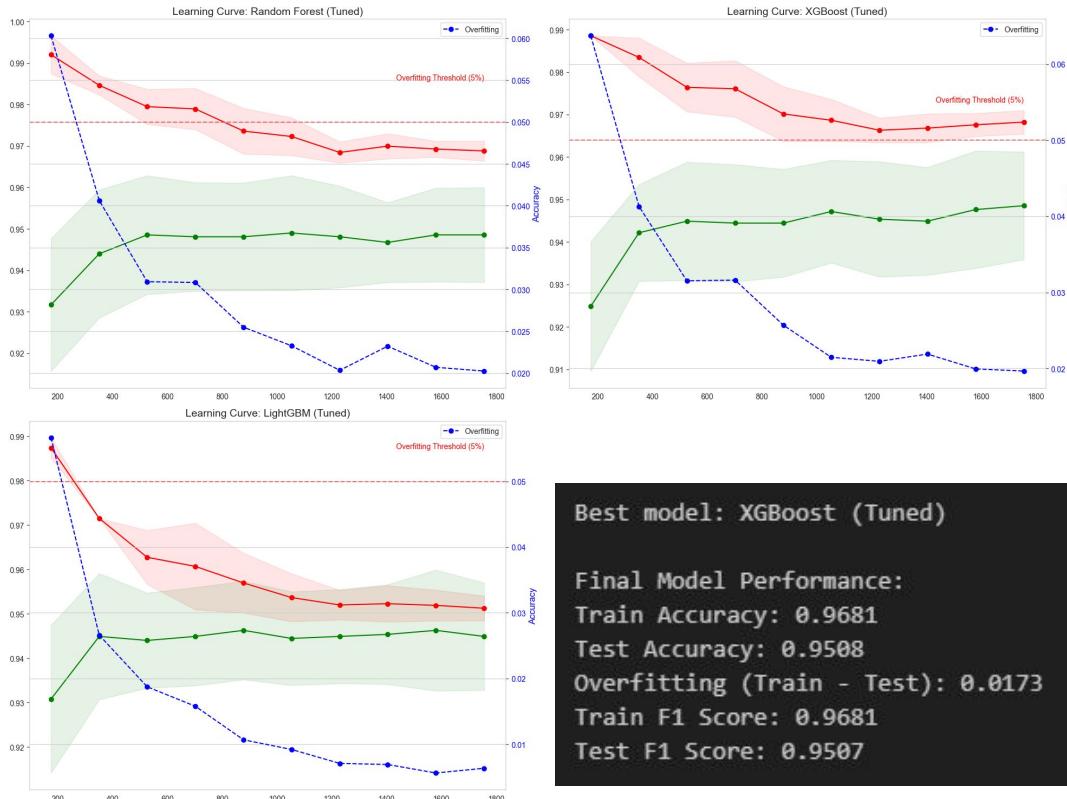
Delta!!!

I made the spirit airlines of machine learning code



So what does this mean?

Sub text go brrr





Regression Analysis

The Problem Statement: What are the factors that impact Life Expectancy?

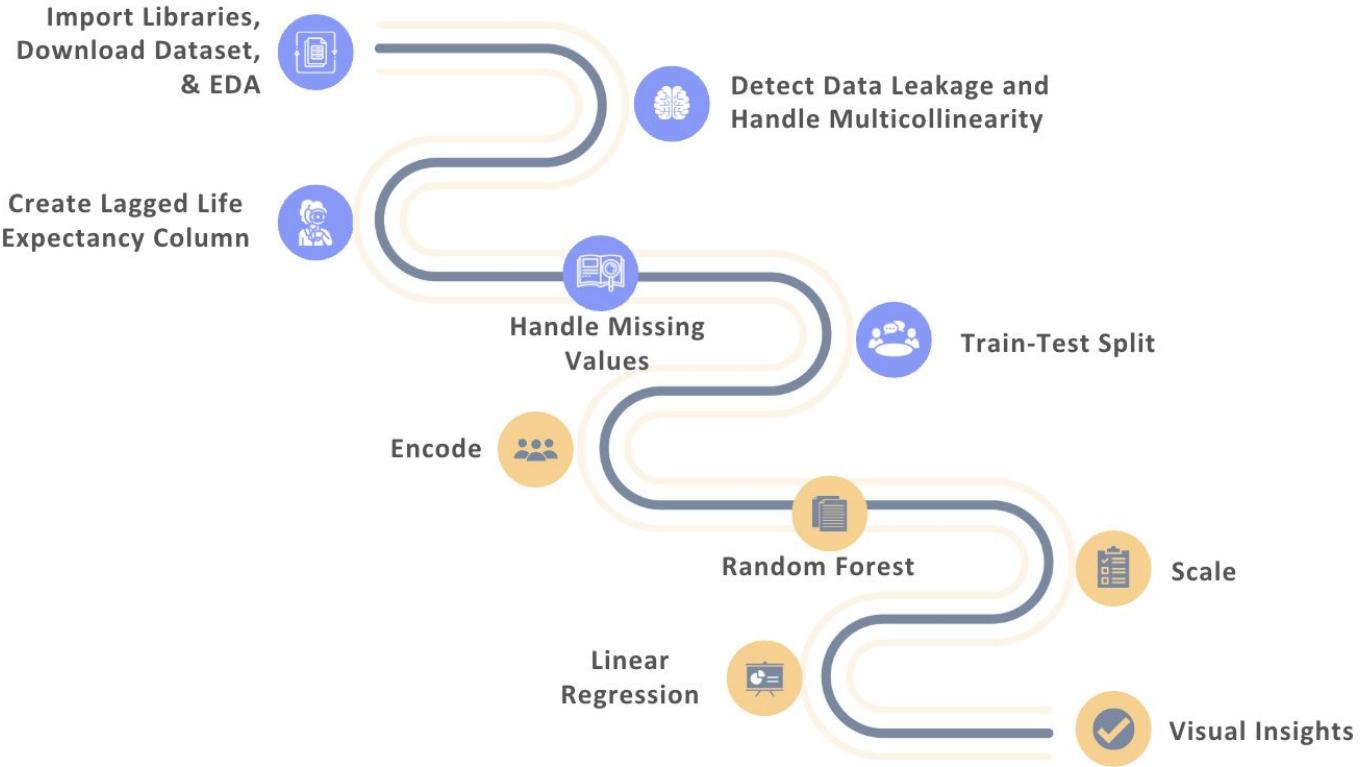
Goal: Predicting life expectancy using a variety of features.

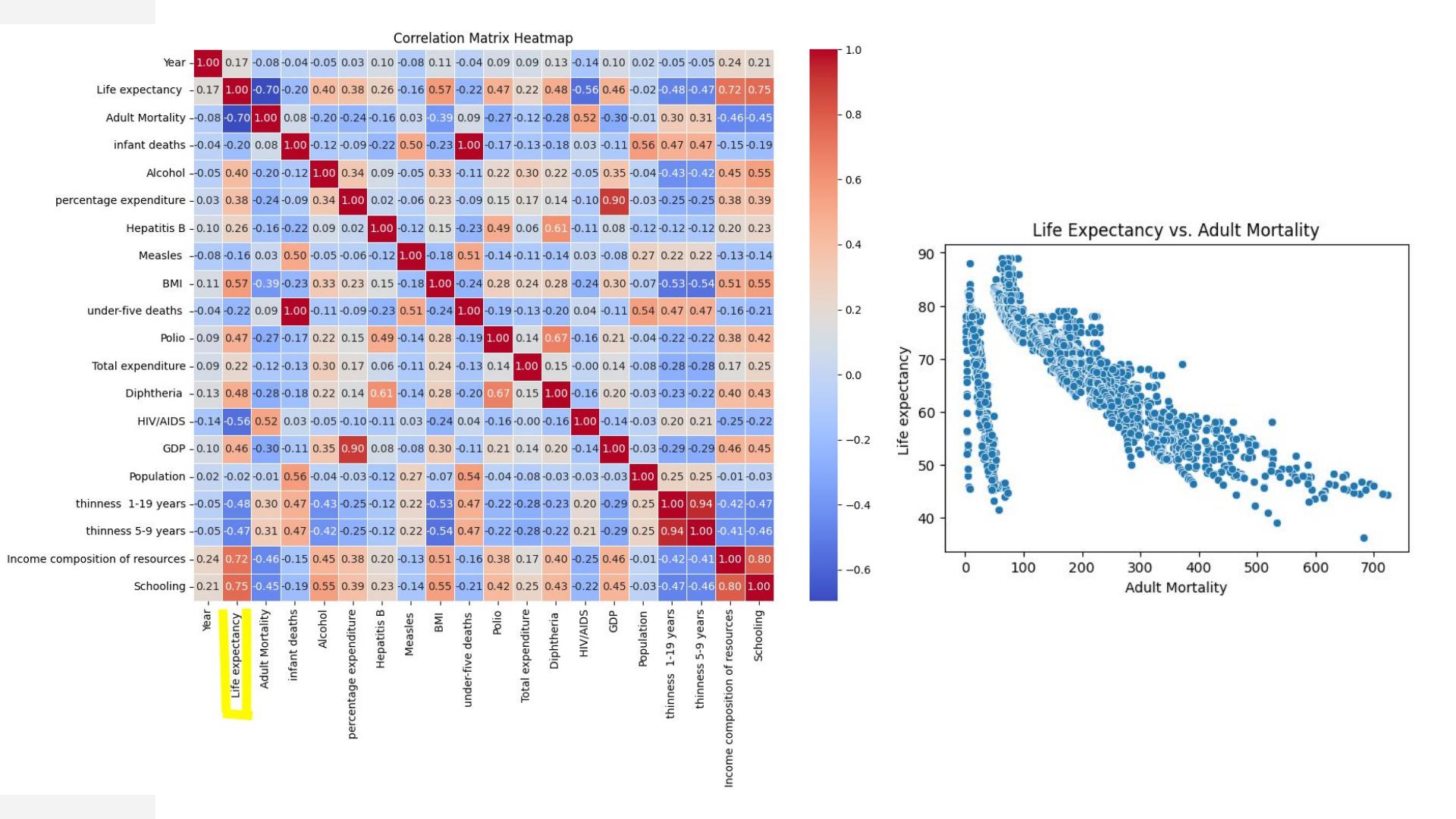
Our Approach: Linear Regression and Random Forest.

Why Regression?: A regression model is ideal for predicting continuous outcomes, such as life expectancy, from various numeric and categorical predictors.



Process Workflow



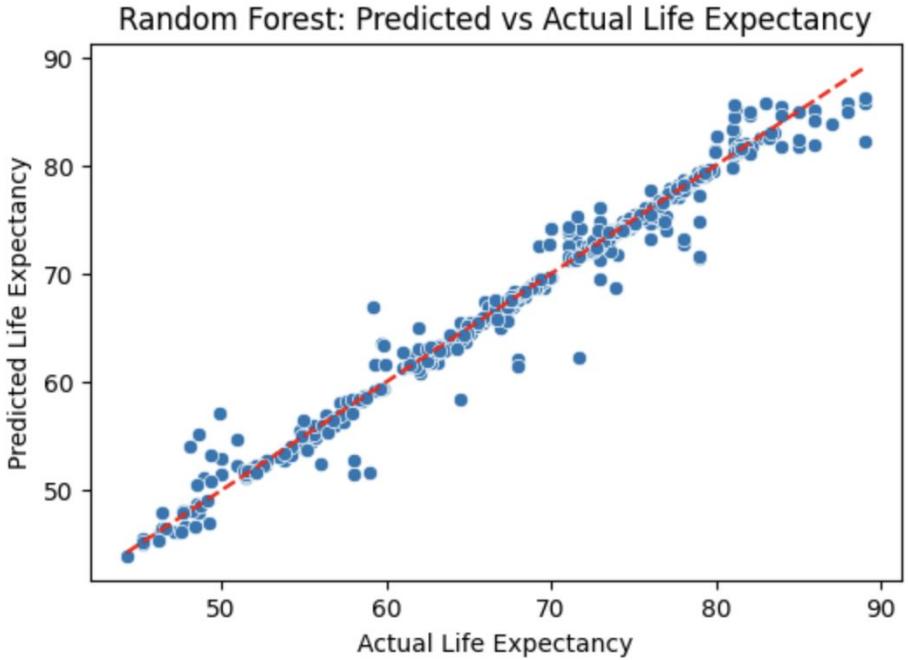


Results

Linear Regression R² Score:	0.9697
Linear Regression MSE:	2.8684
Random Forest R² Score:	0.9742
Random Forest MSE:	2.4434



■ Visual Insights



- **Strong positive correlation**
- **Some instances of outliers**
- **Model Effectiveness**

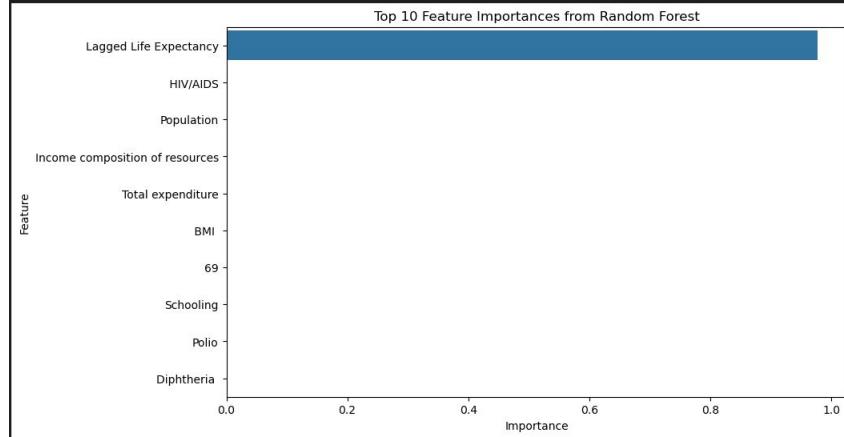




Visual Insights *Continued*

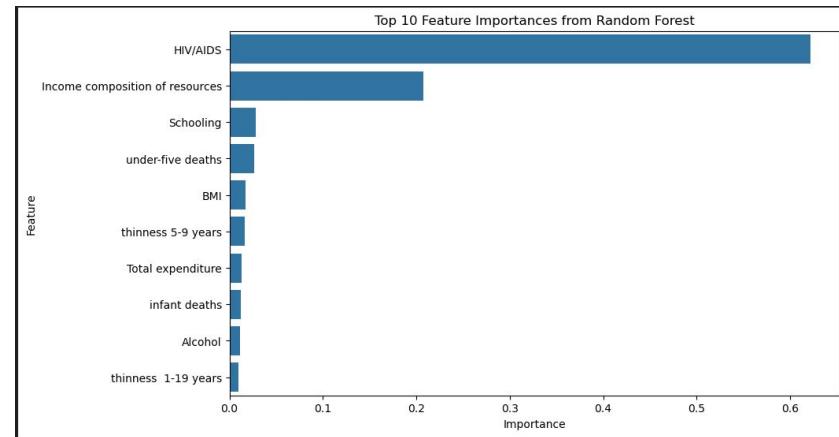
- Key Implications
- Potential Contributing Factors





Further Exploration

- Overfitting
- Additional Evaluation and Tuning





Model	R ² Score	MSE
Linear Regression	0.679854	30.138812
Ridge Regression	0.680977	30.033117
Lasso Regression	0.614680	36.274393
ElasticNet Regression	0.628472	34.975971
Random Forest	0.933913	6.221531
Gradient Boosting	0.896896	9.706326
Support Vector Regression	0.741046	24.378221
Decision Tree	0.877399	11.541767
K-Neighbors	0.842791	14.799808

SVR Tuning Results:
R²: 0.8192
MSE: 17.0193

ANALYSIS

- Random Forest still performs well
- SVR improved from .741 to .892 with hyperparameter tuning



THANK YOU

March 31, 2025