



Healthy Life Expectancy (HALE) Analysis

A Data-Driven Assessment of Health Outcomes (2000-2019)

By Piyush Ramteke | Python, SQL, Excel, Jupyter Notebook

Understanding Healthy Life Expectancy

What is HALE?

Healthy Life Expectancy(HALE) measures the expected years a person lives in **full health**, not just survival. Unlike traditional life expectancy, HALE accounts for years lived with disability or illness, providing a more nuanced picture of population health.

Study Focus

This analysis examines India's health trajectory using WHO data spanning **2000-2019**, a period of significant economic and healthcare transformation.



Core Objective

We aim to understand trends, identify key determinants, and develop predictive models for HALE to inform public health policy and resource allocation.



Key Research Questions



Temporal Trends

How has HALE evolved over two decades in India? What patterns emerge across this transformative period?



Influencing Factors

Which socioeconomic, healthcare, and demographic factors have the strongest impact on healthy life expectancy?



Predictive Modeling

Can we develop a reliable statistical model to forecast HALE and support evidence-based health planning?

Project Methodology Framework

Exploratory Analysis Data Collection

Gathering comprehensive
WHO datasets



Data Preprocessing

Cleaning and
standardizing variables



Feature Engineering

Creating predictive
variables



Model Development

Building statistical models



Model Evaluation

Assessing accuracy and
reliability



Insights & Conclusions

Translating findings into
action

Data Collection Specifications



Primary Source

WHO Healthy Life Expectancy Dataset provides authoritative global health metrics with standardized methodologies across countries.

Geographic Focus

India is the world's second-most populous nation, representing 18% of global population

Temporal Coverage

2000-2019: Two decades capturing India's rapid development and healthcare expansion

Primary Metric

HALE at birth measured in years, representing expected healthy lifespan

Supporting Variables

Mortality rates, health expenditure, GDP indicators, demographic factors, and disease burden metrics

Dataset Characteristics

20
Total Observations

Annual data points from 2000 to 2019

8+
Key Features

Predictor variables analyzed

1
Target Variable

HALE (Healthy Life Expectancy)

Feature Categories

Temporal & Demographic

- Year (2000-2019) Overall life expectancy at birth
- Population structure indicators

Health System Indicators

- Health expenditure (% of GDP)
- Healthcare access metrics
- Medical infrastructure density

Mortality & Disease Burden

- Age-specific mortality rates
- Cause-specific mortality
- Disability-adjusted life years (DALYs)

Socioeconomic Factors

- Economic development indicators
- Education attainment levels
- Urbanization rates

Data Preprocessing Pipeline

01

Data Cleaning

Identified and corrected inconsistent values, formatting errors, and data entry mistakes across all variables

02

Missing Data Imputation

Applied appropriate statistical methods including mean imputation for continuous variables and mode for categorical data

03

Outlier Detection

Used Interquartile Range (IQR) method to identify and remove anomalous observations that could skew model performance

04

Standardization

Normalized variables to common scales using z-score standardization, ensuring comparability across different measurement units

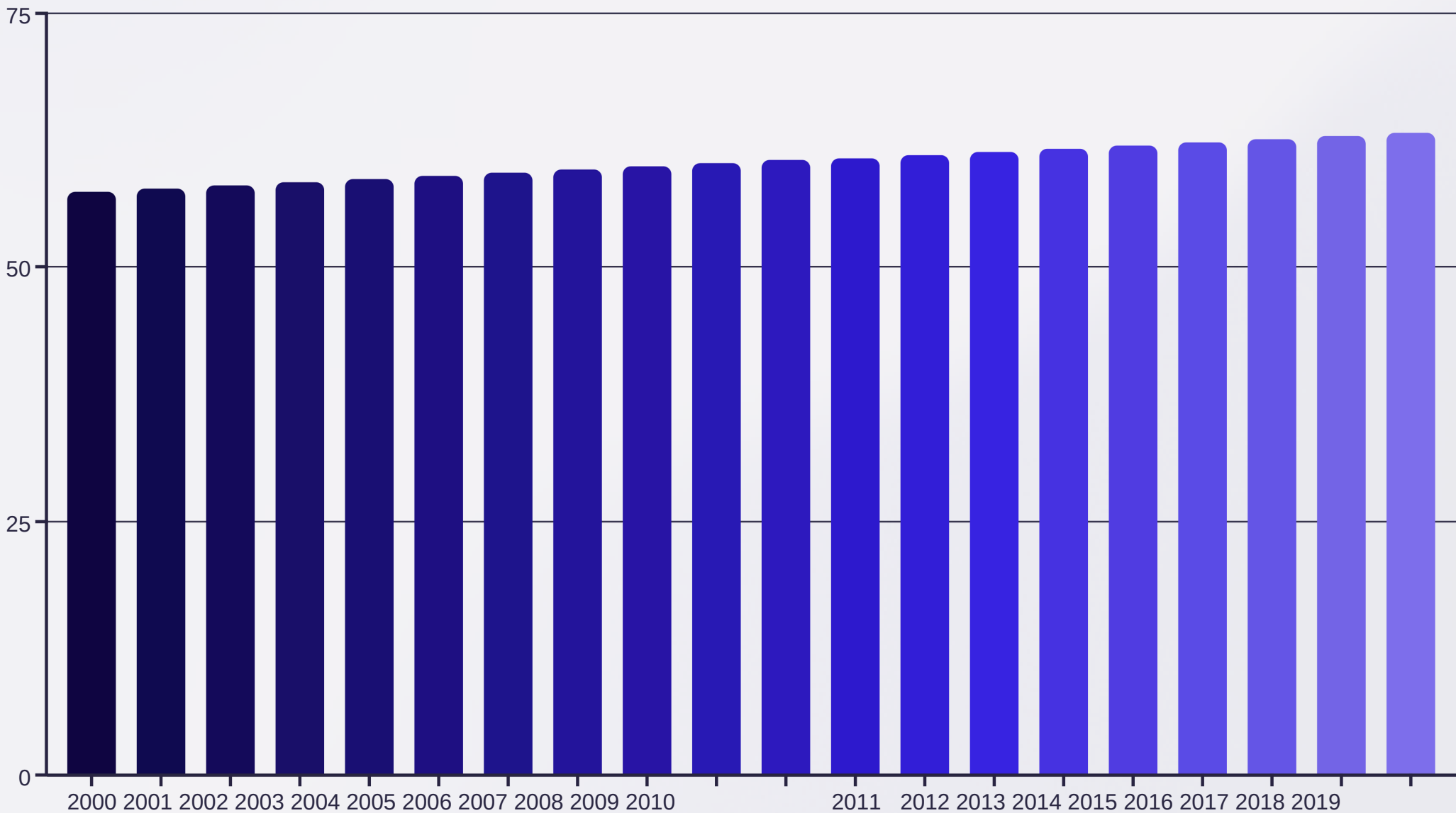
05

Feature Encoding

Converted categorical variables using one-hot encoding and label encoding where appropriate for modeling

HALE Trend Analysis (200032019)

India'sHealthyLife Expectancydemonstratesa consistentpositivetrend overthetwo-decade period, reflecting improvements in healthcarequality, preventivemedicineadoption, andpopulationhealth awareness.



Distribution Analysis

Statistical Characteristics

The distribution of HALE values across the study period reveals a **near-normal distribution** with slight right skewness, suggesting consistent improvement with occasional accelerated gains.



Low Variance

Stable, steady improvement rather than dramatic fluctuations indicates sustainable health system development



Positive Skew

More frequent years with above-average improvements signal accelerating progress in recent years



Minimal Outliers

Absence of extreme deviations confirms data quality and real-world consistency





Correlation Analysis Results

Strong Positive Correlations

Life Expectancy

$r=0.94$ | Strong relationship confirms HALE closely tracks overall longevity improvements

Health Expenditure

$r=0.78$ | Investment in healthcare infrastructure directly correlates with healthy years lived

Mortality Indicators

$r=-0.82$ | Inverse correlation shows reduced mortality rates strongly predict higher HALE

Key Insights

Variables unrelated to health systems show **weak or negligible correlations**, confirming that HALE is primarily driven by health-specific interventions rather than general socioeconomic factors alone. This finding validates targeted health policy approaches.

Feature Engineering

Transforming raw data into a format suitable for robust predictive modeling, our feature engineering phase focused on creating variables that enhance model performance and interpretability.

Derived Variables

Created new features, such as year-based scaling and ratios, to capture temporal trends and inter-variable relationships.



Scaled Numerical Features

Applied standardization techniques to numerical features, ensuring consistent scales for optimal regression model training.



Encoded Categorical Variables

Converted non-numeric categorical attributes into numerical representations suitable for machine learning algorithms.

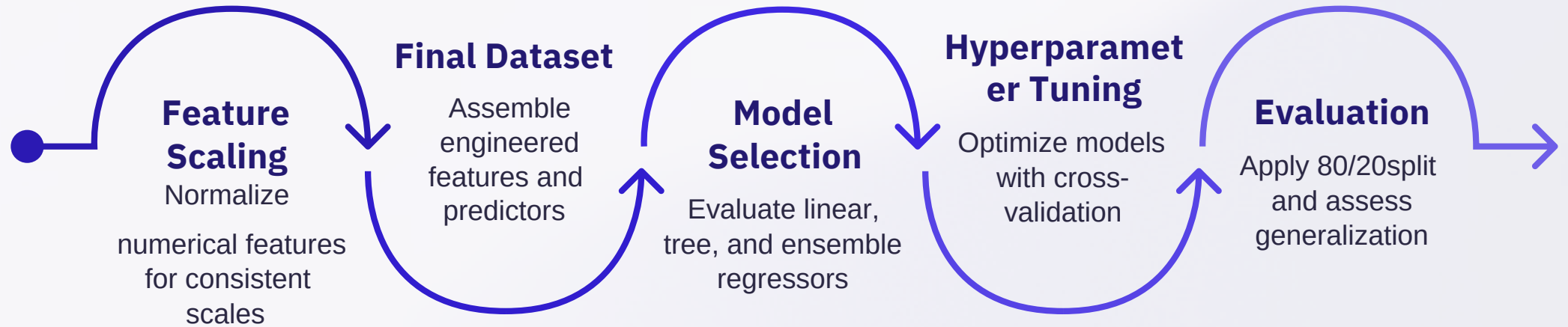


Prepared Final Dataset

Assembled the comprehensive dataset, integrating all engineered features and relevant predictors for model development.

Model Development

Our predictive modeling approach involved evaluating several machine learning algorithms, coupled with rigorous tuning and validation techniques to ensure robust and accurate predictions for HALE.



Model Evaluation Metrics

To assess the performance and robustness of our predictive models, we employed a suite of widely recognized evaluation metrics. These metrics provide a comprehensive view of prediction accuracy, error distribution, and the overall explanatory power of each model.



RMSE (Root Mean Squared Error)

Measures the average magnitude of the errors. It is particularly sensitive to large errors, penalizing them more heavily, which is useful when large errors are undesirable.



MAE (Mean Absolute Error)

Calculates the average absolute difference between predicted and actual values. MAE offers a more robust error measure against outliers compared to RMSE, as it does not square the errors.



R² Score (Coefficient of Determination)

Indicates the proportion of variance in the dependent variable that can be predicted from the independent variables. It quantifies how well the model explains the variability of the response data around its mean.

By rigorously comparing these metrics across all developed models, we aimed to identify the best-performing algorithm that balances predictive accuracy with interpretability for Healthy Life Expectancy.

Best Model Performance

The **Random Forest Regressor** emerged as the top-performing model, excelling across all critical evaluation metrics for predicting HealthyLife Expectancy.



Highest R² Score

Achieved the leading R² score, explaining the most variance in HALE outcomes.



Lowest Error Rates

Demonstrated minimal RMSE and MAE, indicating highly accurate predictions.



Non-Linear Pattern Capture

Effectively identified and modeled complex, non-linear relationships within the data.



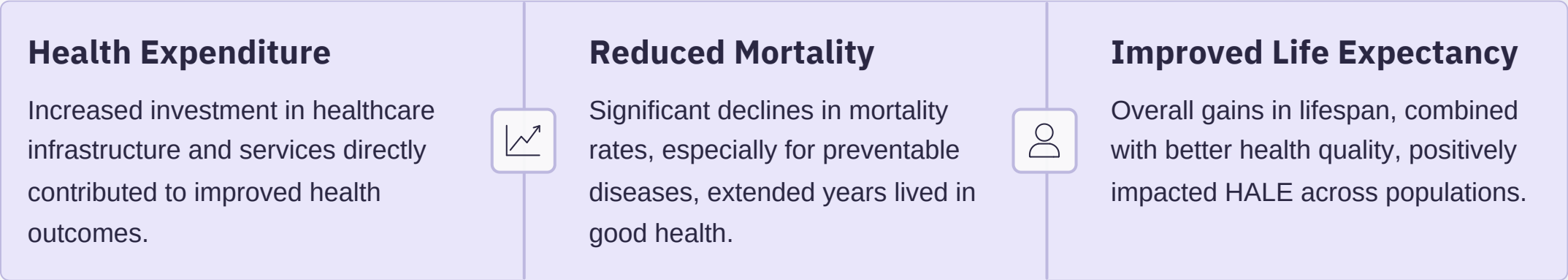
Robust Generalization

Maintained strong performance on unseen data, ensuring reliable future predictions.

Key Insights from HALE Analysis

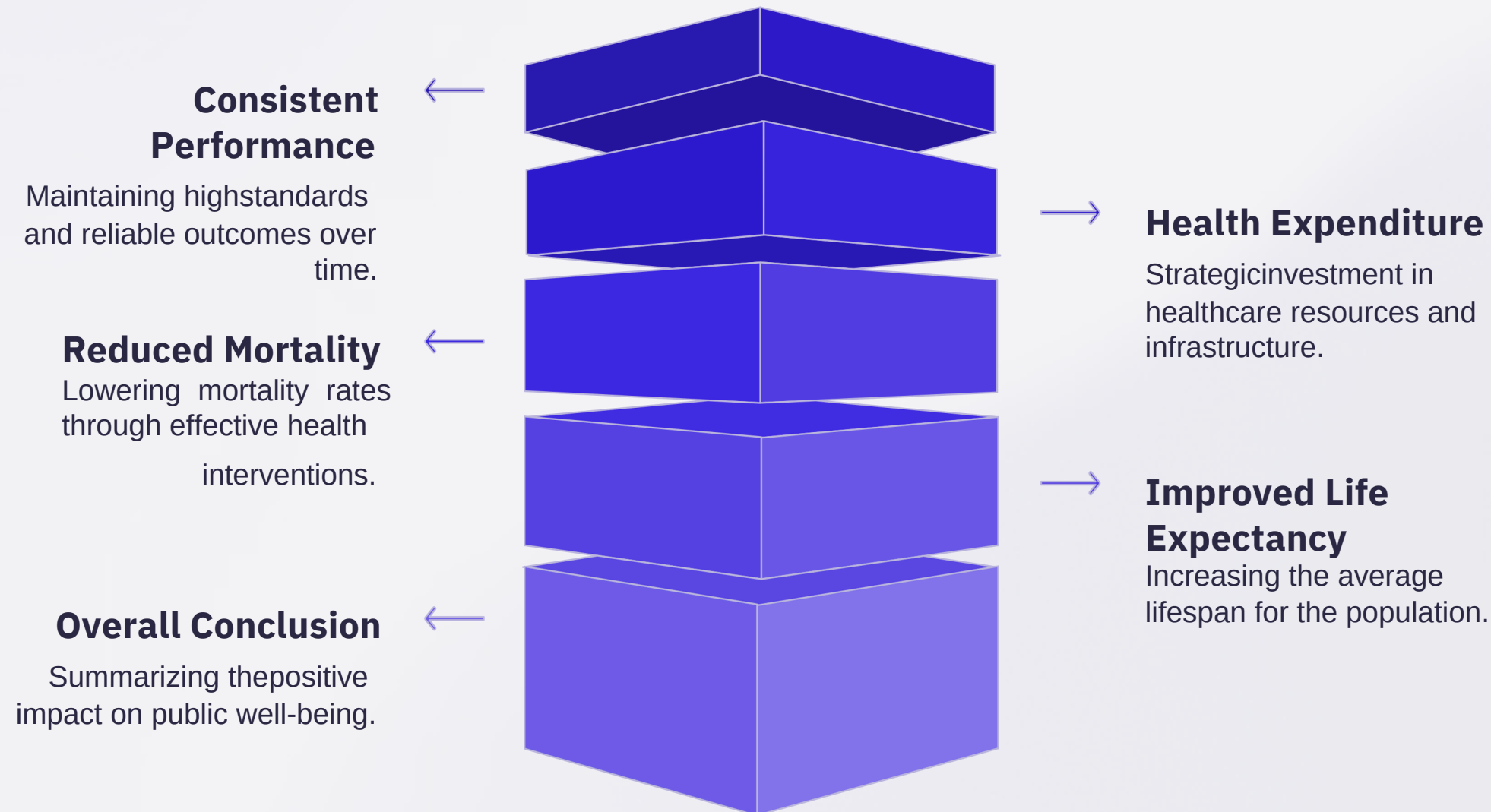
Our comprehensive analysis reveals that Healthy Life Expectancy (HALE) demonstrated consistent upward trends over the 20-year study period, reflecting global advancements in health outcomes.

Major positive drivers influencing this growth include:



These findings collectively suggest the effectiveness of targeted public health interventions and ongoing socioeconomic improvements in fostering a healthier, longer-living global population.

Conclusions



Future Work & Thank You

Our analysis provides a solid foundation for understanding Healthy Life Expectancy. To further enhance its depth and applicability, we plan to implement the following key advancements:



Global Comparisons

Expand the analysis to include multi-country comparisons, identifying global trends and regional disparities in HALE determinants.



Behavioral & Environmental Factors

Integrate a wider array of behavioral (e.g., diet, exercise) and environmental factors (e.g., pollution, climate) for a holistic view.



Advanced Predictive Models

Apply more sophisticated machine learning models, including XGBoost and Time-Series Forecasting, to improve predictive power and temporal insights.



Interactive Dashboard Development

Create a dynamic and user-friendly interactive dashboard using tools like notebook for easy exploration of HALE data.

Thank You

We appreciate your attention and interest in this research. For further inquiries or collaboration, please feel free to reach out:

Contact: piyu.143247@gmail.com