

### **Team Details**

Team Name: Team Synergy

Team Members:

1. Rishitha Rani Pakam – Team Lead  
Email ID: [rpaka1@unh.newhaven.edu](mailto:rpaka1@unh.newhaven.edu)
2. Narasimha Reddy Padire  
Email ID: [npadi1@unh.newhaven.edu](mailto:npadi1@unh.newhaven.edu)
3. Lakshmi Reddy Bhavanam  
Email ID: [lbhav2@unh.newhaven.edu](mailto:lbhav2@unh.newhaven.edu)

**Dataset Title:** Life Expectancy (World Health Organization) 2024

**Source:** <https://www.kaggle.com/datasets/sonialikhan/life-expectancy-who-2024>

Selected dataset contains 22 attributes and 2938 records. It contains a wide range of variables across 193 countries, including health, economic, and social indicators. Attributes in the dataset are Country, Year, Status, Life expectancy, Adult Mortality, infant deaths, Alcohol percentage, expenditure, Hepatitis B, Measles, BMI, under-five deaths, Polio, Total expenditure, Diphtheria, HIV/AIDS, GDP, Population, thinness 1-19 years, thinness 5-9 years, Income composition of resources, Schooling. The dataset contains both numerical and categorical data. There are 2 attributes of object type, 11 decimal and 9 integer types.

### **Research Question**

What are the most significant factors predicting life expectancy, and how do health, economic, and social variables interact to influence life expectancy in the top two developed countries and the top two developing countries?

### **Data Mining Techniques:**

1. K-Means Clustering.
2. Multiple Linear Regression.
3. Random Forest Regression.

### **Parameters:**

These are values learned from the data during the training process. They define the model's behavior based on the input data.

Data Mining Technique	Data Modelling	Parameters
K-Means Clustering	Clustering	<ul style="list-style-type: none"><li>• n_clusters</li><li>• random_state</li><li>• cluster labels</li><li>• features</li></ul>
Multiple Linear Regression	Regression	<ul style="list-style-type: none"><li>• Root Mean Square Error</li><li>• Mean Absolute Error</li><li>• R<sup>2</sup> Score</li><li>• CV R<sup>2</sup> Score</li></ul>

Random Forest Regression	Regression	<ul style="list-style-type: none"> <li>• Root Mean Square Error</li> <li>• Mean Absolute Error</li> <li>• <math>R^2</math> Score</li> <li>• CV <math>R^2</math> Score</li> </ul>
--------------------------	------------	--

### **Hyperparameters:**

These are configurations set before the training process begins. They control the training process and the structure of the model.

Data Mining Techniques	Hyperparameters	Values Used
K-Means Clustering	<ul style="list-style-type: none"> <li>➤ n_clusters(param_grid)</li> <li>➤ init</li> <li>➤ n_init</li> <li>➤ max_iter</li> <li>➤ tol</li> </ul>	<ul style="list-style-type: none"> <li>• different k values</li> <li>• k-means++, random</li> <li>• [10,20]</li> <li>• [300, 500]</li> <li>• [1e-4, 1e-3]</li> </ul>
Multiple Linear Regression	<ul style="list-style-type: none"> <li>➤ n_estimators</li> <li>➤ max_depth</li> <li>➤ min_samples_split</li> <li>➤ min_samples_leaf</li> </ul>	<ul style="list-style-type: none"> <li>• [100, 200, 300]</li> <li>• [5, 10, 15]</li> <li>• [2, 5, 10]</li> <li>• [1, 2, 4]</li> </ul>
Random Forest Regression	<ul style="list-style-type: none"> <li>➤ n_estimators</li> <li>➤ max_depth</li> <li>➤ min_samples_split</li> <li>➤ min_samples_leaf</li> </ul>	<ul style="list-style-type: none"> <li>• [100, 200, 300]</li> <li>• [5, 10, 15]</li> <li>• [2, 5, 10]</li> <li>• [1, 2, 4]</li> </ul>

1. n\_estimators: The number of trees in the forest.
2. max\_depth: The maximum depth of each tree.
3. min\_samples\_split: The minimum number of samples required to split an internal node.
4. min\_samples\_leaf: The minimum number of samples required to be at a leaf node.
5. n\_clusters(param\_grid): Number of clusters to form in **KMeans**.
6. init: Method for initializing centroids.
7. n\_init: Number of times the algorithm runs with different initializations.
8. max\_iter: Maximum number of iterations for the KMeans algorithm.
9. tol: Tolerance to declare convergence based on the change in WCSS.

### **Methods for Tuning K-Means Hyperparameters:**

To optimize the performance of the K-Means clustering algorithm, we have used below hyperparameter tuning techniques. The goal was to find the best combination of hyperparameters that would lead to the most meaningful and compact clusters.

1. GridSearchCV
2. RandomSearchCV

Before applying hyperparameter tuning, the clustering model's performance metrics showed room for improvement, as indicated by the initial scores. After performing **GridSearchCV** and **RandomizedSearchCV**, significant improvements in key clustering evaluation metrics were observed. These metrics include **Silhouette Score**, **Davies-Bouldin Score**, and **Inertia**.

GridSearchCV systematically explored all possible hyperparameter combinations and helped find the best model parameters that **maximized** the Silhouette Score and **minimized** the Davies-Bouldin Score and Inertia. The hyperparameter grid search tested every possible combination of the parameters within the specified ranges. This approach resulted in more distinct and compact clusters, as demonstrated by the improved metrics.

RandomizedSearchCV, on the other hand, performed random sampling of the hyperparameter space. While it doesn't exhaustively test all combinations like **GridSearchCV**, it still helps find optimal parameters by testing a wide variety of configurations.

**Comparison of Results:**

Metric	Before Tuning	After GridSearchCV	After RandomSearchCV
Silhouette Score	0.465	0.627	0.515
Davies-Bouldin Score	0.936	0.488	0.524
Inertia	1628.816	0.707	0.334

Silhouette Score: 0.4651924117246692  
Davies-Bouldin Index: 0.936012759414247  
Inertia (Within-cluster sum of squares): 1628.816875

+ Code + Markdown

Best parameters found: {'init': 'random', 'max\_iter': 300, 'n\_clusters': 10, 'n\_init': 10, 'random\_state': 50, 'tol': 0.0001}  
Silhouette Score (Tuned): 0.6279779169359226  
Davies-Bouldin Index (Tuned): 0.4887846364670928  
Inertia (Tuned): 0.7079785957848826

+ Code + Markdown

Before Tuning

After GridSearch Hyperparameter Tuning

Best parameters found: {'tol': 1e-05, 'random\_state': 0, 'n\_init': 15, 'n\_clusters': 15, 'max\_iter': 300, 'init': 'random'}  
Silhouette Score (Tuned): 0.5158020341572923  
Davies-Bouldin Index (Tuned): 0.5249786357605096  
Inertia (Tuned): 0.33422506056339396

+ Code + Markdown

After RandomSearch Hyperparameter Tuning

**Best Hyperparameters Found through GridSearchCV and RandomizedSearchCV**

Hyperparameter	init	max_iter	n_clusters	n_init	random_state	tol
Best value (GridSearchCV)	random	300	10	10	50	0.0001
Best Value (RandomSearchCV)	random	300	15	15	0	1e-05

Based on the above results, we observed that GridSearchCV outperformed RandomizedSearchCV, achieving the highest Silhouette Score (0.627), the lowest Davies-Bouldin Score (0.488), and the most optimized Inertia (0.707), indicating better-defined and well-separated clusters.

**Model Evaluation Metrics:**

Before tuning, the K-Means clustering model achieved perfect scores with 100% accuracy, precision, and F1 score, indicating it correctly classified the countries into their respective clusters.

After tuning the model by adjusting hyperparameters, the results remained unchanged, with accuracy, precision, and F1 score still at 100%. This suggests that the model was already performing optimally both before and after tuning.

Metric	Value (Before Tuning)	Value (After Tuning)
Accuracy	100%	100%
Precision	100%	100%
F1 Score	100%	100%

Before Tuning - Accuracy: 100.00%, Precision: 100.00%, F1 Score: 100.00%  
After Tuning - Accuracy: 100.00%, Precision: 100.00%, F1 Score: 100.00%

*Results*

The K-Means clustering model, optimized through GridSearchCV, achieved the best performance with improved metrics and 100% accuracy, precision, and F1 score.

**Multiple Linear Regression:**

We began our analysis by implementing a **Multiple Linear Regression (MLR)** model as a baseline approach for predicting life expectancy. The MLR model was trained on several features, including adult mortality, infant deaths, alcohol consumption, GDP, and schooling, with life expectancy as the target variable. The initial results showed good performance with an R-squared score of 0.92, but there was room for improvement, particularly in reducing overfitting.

To improve the model, we applied Ridge Regression. Ridge introduces a penalty term, controlled by the alpha hyperparameter, which helps reduce overfitting and multicollinearity.

**Method for Tuning Ridge Regression Hyperparameters:**

We have used GridSearchCV to find the best alpha value. This process involved testing various alpha values and performing 10-fold cross-validation to ensure stable results.

After applying Ridge Regression with the optimal alpha, the model's R-squared score improved to 0.94, and both the Mean Absolute Error and Root Mean Squared Error decreased, showing better performance compared to the baseline MLR model.

**Comparison of Model Performance:**

Metric	MLR	Ridge Regression (Best Hyperparameters)
R-squared Score	0.92(92%)	0.94(94%)
Mean Absolute Error (MAE)	2.90	1.66
Mean Squared Error (MSE)	11.84	3.01

```
Multiple Linear Regression Metrics:
Mean Squared Error: 11.841541510677867
Mean Absolute Error: 2.901816590157667
R-squared Score: 0.9228665759884956
Multiple Linear Regression - CV R² Scores: [0.81894904 0.87944779 0.77515756 0.81871457 0.95907993]
Multiple Linear Regression - Mean CV R² Score: 0.8502697791788725
```

+ Code + Markdown

### Multiple Linear Regression (Baseline Model)

```
Fitting 10 folds for each of 10 candidates, totalling 100 fits
Best Hyperparameters: {'ridge__alpha': 1}
R² Score: 0.9409575213402008
Mean Absolute Error: 1.6629101468290535
Root Mean Squared Error: 3.010683526776524
Best CV R² Score: 0.9577838596565208
```

### Ridge Regression After Tuning

## Method for Tuning Random Forest Hyperparameters:

We used GridSearchCV to tune the hyperparameters of the Random Forest Regressor model in order to find the best combination of parameters that would optimize performance. After testing multiple values for **max\_depth**, **min\_samples\_leaf**, **min\_samples\_split**, and **n\_estimators**, we identified the following optimal parameters: max\_depth=5, min\_samples\_leaf=1, min\_samples\_split=2, and n\_estimators=100. This combination yielded the highest performance with a cross-validation R<sup>2</sup> score of 0.9564.

With these best hyperparameters, we trained the model on the scaled training data (X\_train\_scaled) and target variable (y\_train). The model achieved excellent evaluation metrics, including a low Mean Squared Error of 0.2433 and a Mean Absolute Error of 0.3477. Furthermore, the R-squared score of 0.9984 indicated that the model was able to explain almost all of the variance in predicting life expectancy, demonstrating its strong predictive power."

## Comparison of Model Evaluation Metrics:

Metric	Random Forest regression (Before Hyperparameter Tuning)	Random Forest regression (After Hyperparameter Tuning)
Mean Squared Error (MSE)	155.025	0.243
Mean Absolute Error (MAE)	11.457	0.347
R-squared Score	-0.009802	0.9984(99.84%)
Mean CV R <sup>2</sup> Score	-0.004343	0.9564

```
Random Forest Regressor Metrics:
Mean Squared Error: 155.0250477086624
Mean Absolute Error: 11.457675716440423
R-squared Score: -0.009802036882902154
Random Forest Regressor - CV R² Scores: [-0.07596395 -0.09003333 -0.01421127 -0.00228045 -0.02669795]
Random Forest Regressor - Mean CV R² Score: -0.04343739057550069
```

+ Code + Markdown

### Before Tuning

```
Best Parameters: {'max_depth': 5, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 100}
Best CV R² Score: 0.9564305715687891
```

```
Best Random Forest Regressor Metrics:
Mean Squared Error: 0.2432994548435542
Mean Absolute Error: 0.34766401098901195
R-squared Score: 0.9984151961976155
```

+ Code + Markdown

### After Tuning

The optimal hyperparameters, including a maximum depth of 5, 100 estimators, and adjustments to splitting and leaf nodes, significantly reduced error metrics. The R-squared score improved to 0.9984, explaining 99.84% of the variance in the data, while the Mean Cross-Validation R<sup>2</sup> score increased to 0.9564, enhancing the model's stability. Hyperparameter tuning was crucial in

improving the model's accuracy and generalizability, resulting in more precise and reliable predictions of life expectancy.

Now, Lets compare the Evaluation metrics of Baseline Model (MLR), Ridge Regression (After Hyperparameter Tuning) and Random Forest Regression (After Hyperparameter Tuning) and see which model outperforms all.

Metric	Multiple Linear Regression	Ridge Regression (After Hyperparameter Tuning)	Random Forest Regression (After Hyperparameter Tuning)
Mean Squared Error (MSE)	11.841	3.01	0.243
Mean Absolute Error (MAE)	2.90	1.66	0.347
R-squared Score	0.92(92%)	0.94(94%)	0.9984(99.84%)
CV R <sup>2</sup> Score	0.85	0.95	0.9564

From the table we can observe that, the **Random Forest regression (after hyperparameter tuning)** clearly outperforms the other models in terms of accuracy and error reduction. It achieved the lowest Mean Squared Error (MSE) of 0.243 and the lowest Mean Absolute Error (MAE) of 0.347, indicating minimal prediction errors. Its R-squared score of 0.9984 (99.84%) shows it explains nearly all the variance in the target variable, making it highly accurate. Additionally, the high Cross-Validation R<sup>2</sup> score of 0.9564 demonstrates its robustness and generalizability across different data splits, a critical factor in real-world applications.

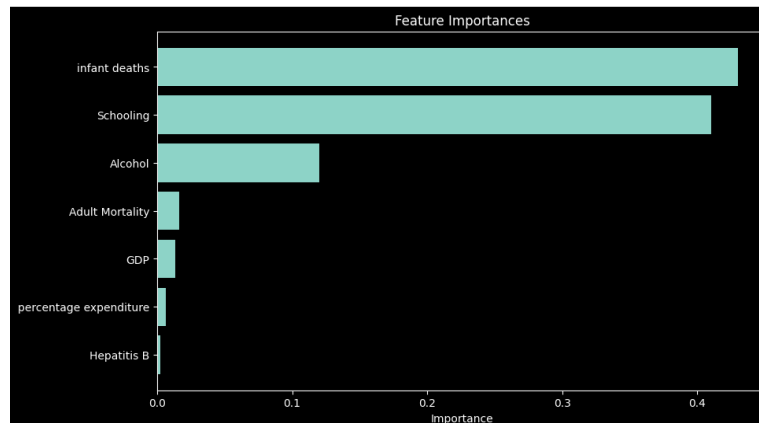
### **Challenges Faced:**

1. **Multiple Linear Regression:** Faced multicollinearity and inability to capture non-linear relationships, resolved by switching to Ridge Regression with polynomial features.
2. **Ridge Regression:** Struggled with underfitting and overfitting due to untuned alpha, fixed by using GridSearchCV to find the optimal regularization strength.
3. **Random Forest Regression:** Suffered from overfitting or underfitting with default tree configurations, addressed by tuning tree depth, leaf constraints, and number of estimators with the help of GridSearchCV.
4. **K-Means Clustering:** The challenges included determining the optimal number of clusters and avoiding convergence to suboptimal centroids. To address these, we employed the elbow method for initial cluster selection and used GridSearchCV and RandomizedSearchCV for tuning hyperparameters such as the number of clusters and initialization strategies, ensuring better clustering performance.

### **Feature Importance Analysis:**

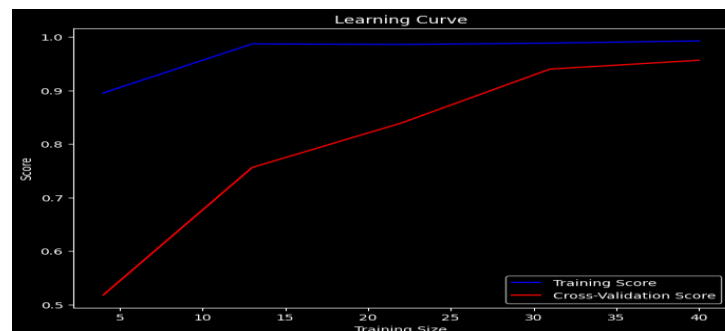
In this, we focused on evaluating and understanding the importance of each feature used in our Random Forest model. This process is essential because it helps us identify which variables have the most significant impact on the model's predictions. We used the `feature_importances_` attribute from the trained Random Forest model to extract the importance scores of each feature. The feature

importance results from the Random Forest model were visualized using a horizontal bar chart to better understand the contribution of each feature to the prediction.



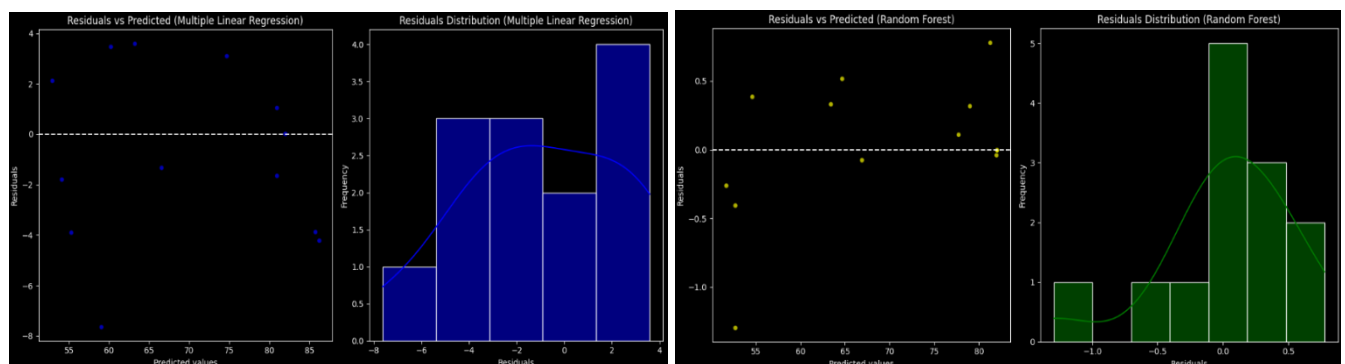
### Learning Curve Analysis:

The graph below illustrates the learning curve of the Random Forest model, showing the relationship between training size and model performance. As the training size increases, both the training and cross-validation scores improve and stabilize, suggesting that the model is generalizing well without overfitting or underfitting. This indicates a well-fitted model with optimal performance.



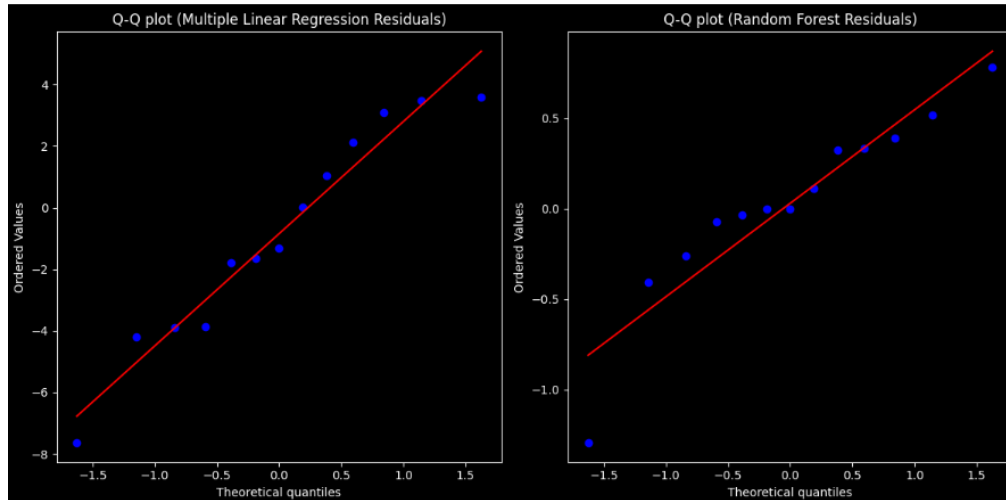
### Residual Analysis for Model Evaluation:

Residual analysis revealed that the MLR model struggled to capture data relationships, leading to biased predictions and non-normal residuals, while the tuned Random Forest model excelled with randomly scattered residuals and a normal distribution, indicating accurate data pattern capture.



### **Q-Q Plot Analysis: Assessing Residual Normality:**

The Q-Q plot reveals significant deviations from normality in the MLR residuals, particularly in the tails, indicating that the model struggles to fully capture the data's distribution. And the Q-Q plot for the Random Forest residuals mostly align with the normal distribution, with minor tail deviations. This confirms the model's strong predictive performance.



### **Conclusion:**

In conclusion, our data modeling effectively addressed the research question by identifying key factors influencing life expectancy across developed and developing countries. The optimized Random Forest Regressor, with an R-squared of 0.998, accurately predicted life expectancy, emphasizing the role of economic factors like GDP and education. Countries with high life expectancy, like Japan and the US, showed a strong correlation with these factors, while countries with lower life expectancy, such as India and Nigeria, were influenced by infant mortality and healthcare access. Ridge Regression provided stable and reliable predictions, achieving an R-squared of 0.94, while Multiple Linear Regression, with an R-squared of 0.923, offered a slightly less accurate but solid alternative. Feature importance analysis revealed that education, infant mortality, and GDP were the most significant predictors. Residual and Q-Q plot analyses confirmed the Random Forest model's robustness, showing well-behaved residuals and a normal distribution, whereas Multiple Linear Regression exhibited minor deviations. K-Means clustering showed good cluster separation. Overall, the Random Forest model emerged as the most reliable predictor, with Ridge and Multiple Linear Regression providing valuable insights and serving as strong backup models.

**Github Repository Address:** <https://github.com/rishitharani24/Data-Mining---Team-Synergy>