

# Life Expectancy Prediction

Thokala Homakesh, Chenna Rakesh, C.Savinay Kumar

Btech IV Year Department of Electronics and Computer Engineering,  
Sreenidhi Institute of Science and Technology, Ghatkesar, Hyderabad-501301, TS, INDIA.

V.Sowmya

Assistant Proffesor, Department of Electronics and Computer Engineering,  
Sreenidhi Institute of Science and Technology, Ghatkesar, Hyderabad-501301, TS, INDIA.

**Abstract:** The expression "Future expectancy" alludes to the number of years an individual can hope to live. Assessing the future and discovering the qualities which impact it fills in as a significant part as it is the reason for any nation to discover how great they are at keeping up great wellbeing. The venture attempts to make a model dependent on information given (WHO) to assess the future for various nations utilizing different calculations like Simple Linear Regression, Gradient Descent, Linear Regression with Polynomial Feature, Decision Tree Regression and Random Forest Regression by thinking about numerous source factors. The information offers 15 years. In this process, we have considered information of 15 years from 193 nations for examination. This paper intends to investigate the connection of the future with its connected components.

**IndexTerms** – Multiple Linear Regression, Decision tree, Random Forest.

## LINTRODUCTION

Future alludes to the measure of years an individual is anticipated to quantify. By definition, life is predicated on partner gauge of the average age that individuals from a particular populace group will be once they kick the bucket. In this cutting-edge period future of people is continuing expanding step by step. This was just conceivable by investigation of the issues and factors which influenced the future in a given climate. Be that as it may, the conditions are assorted, a few nations are created while others are as yet creating. Indeed, even today a ton of nations are as yet immature. Hence taking every one of the nations and expecting their normal future brings the significant elements into light as improvement has no closure. Wellbeing measurements give significant outcomes to improve wellbeing choices dependent on proof. An examination of the medical care area is truly fundamental to improve convenience and similarity of wellbeing insights.

The expression "Future" alludes to the number of years an individual will hope to quantify. Assessing the lifetime and discovering the qualities that outcome it is a pivotal job since it is that the reason for any nation to look out anyway reasonable, they're at keeping up physiological condition. The task attempts to frame a model upheld information given by (WHO) to guage the lifetime for diverse nations misuses different calculations like basic Linear relapse toward the, Gradient Descent, relapse toward the mean with Polynomial Feature, Decision Tree Regression and Random Forest Regression by thinking about various stockpile factors. the data offers a time span from 2000 to 2015. Among all classes of wellbeing related factors exclusively those essential components were picked that region unit extra delegate. during this task we have measured data of 15 years from 193 nations for examination. This paper plans to explore the relationship of lifetime with its associated factors.



Figure 1. World Life Expectancy Map  
Table1.Domains selected for standardized health status module

Domains indirectly assessing health	Domains directly describing health	
Discrimination/stigma	General health	Sexual activity
Participation barriers	Affect*	Fertility
Self care*	Cognition*	Hearing
Shame/embarrassment	Communication	Speech
Social functioning	Destiny	Vision
Usual activities*	Mobility*	Brushing
	Pain*	Eating
	Skin & bodily disfigurement	Digestion
	Energy/vitality	Bodily excretion

\* Domains selected for standardized health status module

## II.HEALTH STATISTICS RECORD (WHO FACTS)

### a.Population:

The worldwide populace was two.8 billion out of 1955 and is 5.8 billion presently. it'll increment by almost eighty million person's years to prevail with regards to concerning eight billion continuously 2025. In 1955, 68 percent of the overall populace lived in country regions and thirty second in metropolitan regions. In 1995 the quantitative connection was fifty fifth country and forty fifth metropolitan; by 2025 it'll be forty first rustic and 59 urbans. Consistently in 1997, concerning 365000 children were conceived, and concerning 140000 people passed on, giving a characteristic increment of concerning 22000 people day by day.

### b.Life expectancy:

Normal expectation upon entering the world in 1955 remained basically 48 ages; in 1995 it totally remained 65 years; in 2025 it'll arrive at 73 years. Continuously 2025, it's normal that no nation can take an expectation however fifty ages. In excess of fifty million people live these days in nations with an expectation of yet 45 years. More than five billion people in a single hundred twenty nations

these days have expectation of more than sixty years. Around 300,000,000 people board sixteen nations any place expectation really blurred between 1975-1995.

### c.Deaths based on age:

In 1955, 40% of all deaths were among kids beneath five years, 100 percent stood in 5–19-year-olds, twenty eighth stood among grown-ups old 20-64, and twenty one percent stood amid the over-65s. In 1995, solely 21% of all demises remained amid the under-5s, 7% among those 5-19, twenty nine percent among those 20-64, and forty three percent amongst the over-65s. By 2025, 8% of all demises are within the under-5s, 3% amid 5–19-year-olds, twenty seven percent among 20–64-year-olds and sixty three percent between the over- 65s

### d.Foremost causes of worldwide deaths:

Trendy 1997, of a world complete of 53 squillion passing's, 17.4 squillion remained a direct result of irresistible and parasitic infections; 15.2 million were a result of circulatory sicknesses; 6.3 million were a direct result of malignancy; 2.9 million were a result of metabolic interaction illnesses, in the primary constant preventive respiratory organ sickness; and 3.7 million were a result of perinatal conditions. Driving reasons for death from irresistible sicknesses were intense lower metabolic interaction diseases (3.7 million), TB (2.9 million), detachment of the guts (2.5 squillion), HIV/AIDS (2.3 squillion) and protozoal contamination (1.5-2.7 squillion).

## III.OVERVIEW OF THE PROBLEM:

Past research (as indicated by NCBI) recommends that beside financial, local, and organic components, there are a unit some elusive variables that affect person's expectation regardless of the bunch to that the individual has a place in accordance with NCBI. despite the fact that all conditions region unit steady, everybody doesn't kick the bucket at an identical age, and thus there could likewise be distinction long of life among individuals at stretches subgroups. Henceforth, the degree of expectation alone can't give the whole picture of mortality situation in a very populace. For of these reasons, it's important to take a gander at the distinction long of life reliably. Upon this current, it's moreover important to get a handle on the progressions in contrast long of life.

## IV.PROBLEM STATEMENT

To anticipate the future of various nations taking various elements (as characteristics) into thought dependent on the past record of information.

## V.USED SYSTEM PROPERTIES

**a.Computer:** Intel i5 processor 2.0 GHz

**b.Memory:** Minimum of 8GB RAM, 3gb free space for anaconda, 2 Gigabytes AMD M335 dedicated graphics card along with 2 Gigabytes of Intel graphic card

**c.Operating System:** Windows 10, Python version-3.6.5, Anaconda3 version-5.2.0

## VI.FITTING THE MODEL

As we need to predict a particular value (life expectancy which is numeric(age) and can have infinite possible values) it is a regression problem. We used Multiple Linear Regression Algorithm of various degrees (Polynomial regression)

### i.NUMEROUS LINEAR REGRESSION

Numerous linear reversions make an attempt to perfect the link among 2 or additional instructive variable star and a reply variable by fitting an reckoning to determined info. each value of the variable x is connected to a value of the variable y.

The populace regression line for *an* descriptive variable

$X_1, X_2, \dots, X_n$  is well-defined to be

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

Where,

Y is dependent variable  $\beta_0$  is the intercept

$\beta_1, \beta_2 \dots \beta_n$  are the regression coefficients for the explanatory variables i.e., independent variables.

### ii.MULTIPLE LINEAR REGRESSION (DEGREE=1)

In this model, the degree of the fitted model curve is 1. The basic formula if we exclude the random error component is

$$Y_i = \beta_0 + \beta_1 X_i$$

Y is dependent variable  $\beta_0$  is the intercept

$\beta_1$  is the regression coefficient

With the help of sklearn we used Linear Regression class and fitted the model.

```
lm=LinearRegression() #Creating a Linear regression object
lm.fit(x_train,y_train) #Fitting the model
```

### iii.CROSS VALIDATION AND REGRESSION ERROR METRICS

#### Testing the fitted model with test data as input

The output values are denormalized with the help of the following formula

$$\hat{y} = \hat{y}_{norm} \times (\max Y - \min Y) + \min Y$$

```
predict=lm.predict(x_test) #predicting the lifeexpectancy based on x_test values
predict1=predict*(max1[0]-min1[0])+min1[0] #denormalizing the output
d_y_test=y_test*(max1[0]-min1[0])+min1[0] #denormalizing the whole y_test values
print(predict1)
```

```
[[77.10917196]
 [77.01122812]
 [76.85580653]
 [76.61415495]
 [76.138172 ]
 [76.05893488]
 [76.0184429 ]
 [74.72987976]
 [74.48626363]
 [80.70732625]
 [80.65647134]
 [79.94206252]
 [80.07919142]
 [80.80077991]
 [80.67388326]
 [80.1618844 ]
 [79.88568148]
 [79.68040153]
 [79.80869756]]
```

### iv.REGRESSION ERROR METRICS

Used the following metrics:

#### i.Mean Squared Error

MSE essentially events regular square mistake of our forecasts. for every purpose, it computes sq. distinction among the guesses and therefore the mark so averages those standards. the upper this worth, the more severe the perfect is.

#### ii.Root Mean Squared Mistake

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

RMSE is just the square root of MSE.

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2} = \sqrt{MSE}$$

#### iii.Mean Absolute Fault

In MAE the fault is designed as an regular of absolute differences among the target values and the forecasts.

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

#### iv.R Squared (R<sup>2</sup>)

The constant of resolve, or R<sup>2</sup> (occasionally read as

R-two), is additional metric we may use to assess a model and it is closely connected to MSE, but has the benefit of being scale-free it doesn't substance if the output values are very huge or very minor, the R<sup>2</sup> is always profitable to be between -∞ and 1. When R<sup>2</sup> is undesirable it income that the model is worse than forecasting the mean.

$$R^2 = 1 - \frac{MSE(model)}{MSE(baseline)}$$

$$MSE(baseline) = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$$



## v. Error in multiple linear regression (degree=1)

```

mse=mean_squared_error(y_test,predict)
mae=mean_absolute_error(y_test,predict)
mse=mse*(max1[0]-min1[0])+min1[0] #denormalizing the error
mae=mae*(max1[0]-min1[0])+min1[0]

print("Mean squared error:",mse)
print("Root mean squared error:",np.sqrt(mse),"years")
print("mean absolute error:",mae)
print("The r2_score is:",r2_score(predict1,d_y_test))

```

Mean squared error: 42.17696536985433  
Root mean squared error: 6.494379521544327 years  
mean absolute error: 45.25970190455404  
The r2\_score is: 0.4250419257664805

## a. MULTIPLE POLYNOMIAL REGRESSION

Polynomial regression may be a sort of associate degree analysis (multivariate analysis) within the link among the variable quantity  $x$  then therefore the variable  $y$  is showed as an ordinal grade polynomial in  $x$ . The polynomial models are often utilized in those things wherever the link between study and informative variables is curved.

The  $k^{th}$  order polynomial model in one variable is given by

$$y = \beta_0 + \beta_1x + \beta_2x^2 + \dots + \beta_kx^k + \epsilon.$$

We fitted the model of degrees 2 to 8 to analyze which curve is most accurate for our data. We used a for loop for all the higher degree polynomial curves.

```

#Now applying polynomial regression (higher degrees) to find the best fit
for i in range(2,9):
    print("-----Degree:",i,"-----")
    x_train_poly_data=x_train
    poly_reg=PolynomialFeatures(degree=i)
    poly=x_train_poly_data.iloc[:,1].values
    pol=poly_reg.fit_transform(poly)
    l=LinearRegression()
    l.fit(pol,y_train)
    poly_x_test=poly_reg.fit_transform(x_test)
    p=l.predict(poly_x_test)
    pl=p*(max1[0]-min1[0])+min1[0]
    print('The predicted life expectancy for the degree',i)
    print(pl[1:6])
    mse_poly=mean_squared_error(y_test,p)
    mse_pol=mse_poly*(max1[0]-min1[0])+min1[0]
    mae_pol=mean_absolute_error(y_test,p)
    mae_pol=mae_pol*(max1[0]-min1[0])+min1[0]
    print("Mean squared error for degree",i,":",mse_pol)
    print("Root mean squared error for degree:",i,":",np.sqrt(mse_pol),"years")
    print("mean absolute error:",mae_pol)
    print("r2 score:",r2_score(pl,d_y_test))

    print('\n')

```

## b. Decision Tree Regression

Result Tree could be a controlled learning technique that might be pushed off for every association and Reversion issues, anyway to a great extent it's generally famous for discovering Organization issues. it's a tree-structured classifier, any place inside hubs connotes the alternatives of a dataset, branches describe the decision rules and each leaf hub addresses the outcome.

In a decision tree, there are 2 hubs, that are the ideal Node and Leaf Node. call hubs are utilized to make any call and have numerous branches, while Leaf hubs are the yield of these determinations and don't contain somewhat longer twigs. The choices or the investigate are accomplished on the possibility of alternatives of the predefined dataset. it very well may be a realistic delineation for acquiring every one of the expected responses to an issue/decision upheld given condition. it's anything but a decision tree because of, equivalent to a tree, it bounces with the premise hub, that develops extra twigs and ideas a tree-like structure. on request to mark a tree, we tend to utilize the CART algorithmic law, that perspectives for Organization and Regression Tree algorithmic.

Underneath diagram clarifies the over-all structure of a decision tree and implementation:

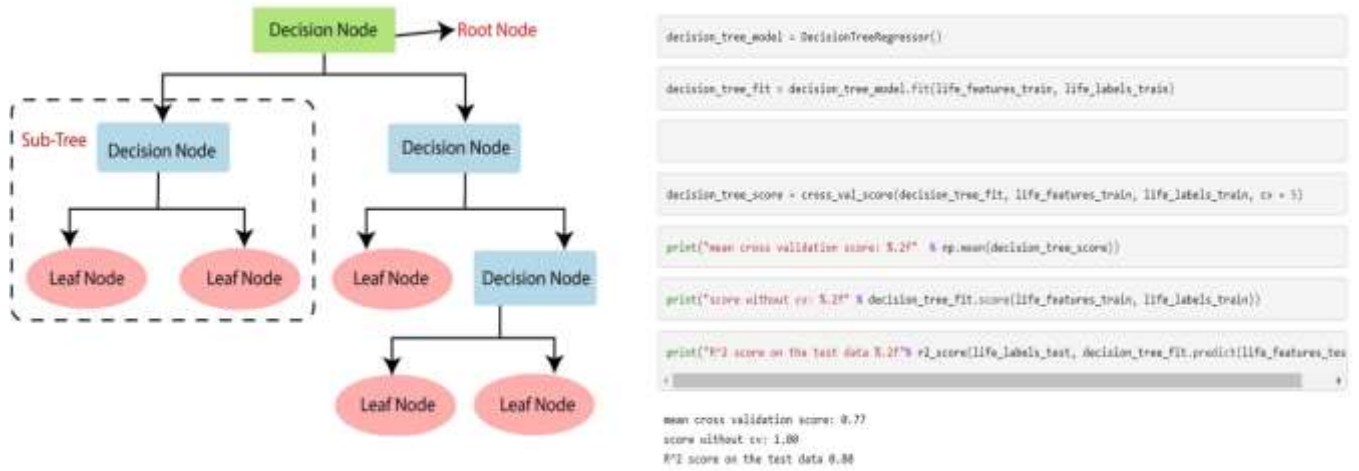


Figure 2. Structure of a decision tree

By Choice Tree Regression. Cross Proof has been performed The R sq. on the coaching knowledge is one which means that the algorithmic rule has scholarly the info by memory, with the cross proof the figure deteriorations to seventy-seven and exploitation the take a look at date we have a tendency to get eightieth. currently we have a tendency to use the algorithmic rule to forecast the values of the life\_features\_test.

**c.RANDOM FOREST REGRESSION**

Chance Forest is a machine knowledge algorithmic program that belongs to the oversight learning technique. It is often used for each Organization and Regression issues in milliliter. it's built mostly on the idea of ensemble knowledge, that could be a method of mixing numerous classifiers to unravel a posh drawback then to boost presentation of the perfect. The term proposes, "Random Forest could be a classifier that contains variety of call trees on numerous subsets of the given dataset and takes the typical to boost the prophetic accuracy of that dataset." rather than hoping happening 1 call tree, chance forest receipts the forecast after every tree then supported bulk elections of forecasts, besides it predicts ultimate result. rather than hoping on one call tree, the chance forest receipts the guess after every tree and supported bulk polls of forecasts, and it forecasts the ultimate productivity. Larger range of trees within forest ends up in developed correctness and prevents the matter of tighten.

The under diagram clarifies the occupied of the Random Forest algorithm.

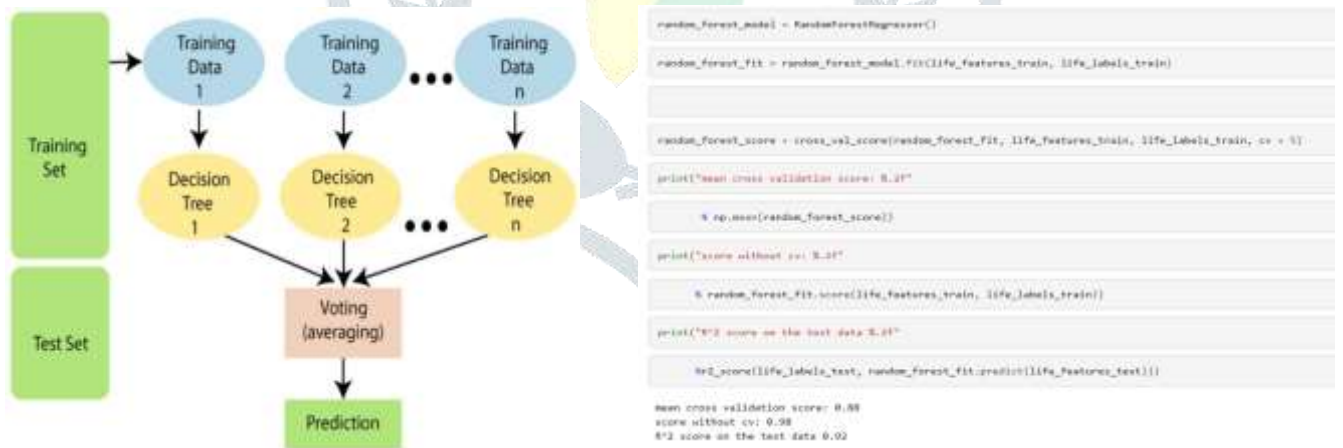


Figure 3. Structure of a Random Forest tree

**Working of Random Forest algorithm**

Chance Forest works in two-phase 1st is to form the chance forest by compounding N call tree, and additional is to create forecasts for every tree shaped within the 1st section. The operating method is explained within the below steps and diagram:

- 1: choose chance K information opinions from the coaching established.
- 2: Shape the choice trees related to the chosen information points (Subsets).
- 3: select the quantity N for call trees that you simply wish to create.
- 4: Duplication Step one & a pair.
- 5: for brand spanking novel information points, realize the calculations of every call Tree. and allot the new information points to the class that victories the bulk votes.

By using the Random Forest Regression:

The algorithms have learned 98% on the training data without cross validation and 88% with, the value is 92 % on the test data.

## VII.RESULT ANALYSIS

Algorithm	R2 Score
Linear Regression	0.82
Decision Tree	0.80
Random Forest	0.92

Table 2 Result Analysis

## VIII.CONCLUSION

In this project we are using different algorithms like Linear regression, Decision Tree and Random Forest to expect the life expectancy of different countries. In this process we consider R2 score as accuracy of life expectancy. So, we get different R2 scores for each algorithm. Out of all these algorithms, Random Forest algorithm gives high accuracy. So, we consider Random Forest algorithm to predict the life expectation of a country.

## IX.REFERENCES

- [1] <https://www.life-expectancy-who/>
- [2] <https://www.machine-learning-decision-tree-classification-algorithm/>
- [3] <https://www.machine-learning-random-forest-algorithm/>
- [4] <http://regression/Chapter12-Regression-PolynomialRegression.pdf> <https://www.datasciencecentral.com/>
- [5] <https://towardsdatascience.com/>
- [6] <https://github.com/SmartPracticeschool/>
- [7] <https://ieeexplore.ieee.org/>
- [8] <https://www.simplilearn.com/>
- [9] <https://blog.floydhub.com/best-machine-learning-books/>
- [10] <https://www.geeksforgeeks.org/best-books-to-learn-machine-learning-for-beginners-and-experts/>
- [11] internet sources and handbooks of machine learning