

PREDICTING LIFE EXPECTANCY USING MACHINE LEARNING

By : SOHAM R. SAHARE

PREDICTING LIFE EXPECTANCY USING MACHINE **LEARNING**

Project Summary :

A typical Regression Machine Learning project leverages historical data to predict insights into the future. This problem statement is aimed at predicting Life Expectancy rate of a country given various features. Life expectancy is a statistical measure of the average time a human being is expected to live, Life expectancy depends on various factors: Regional variations, Economic Circumstances, Sex Differences, Mental Illnesses, Physical Illnesses, Education, Year of their birth and other demographic factors. This problem statement provides a way to predict average life expectancy of people living in a country when various factors such as year, GDP, education, alcohol intake of people in the country, expenditure on healthcare system and some specific disease related deaths that happened in the country are given.

The output algorithms have been used to test if they can maintain their accuracy in predicting the life expectancy for data they have not been trained. Some of the algorithms that can be possibly used are:

- Linear Regression
- Random Forest Regression
- Decision Trees

And by Hyperparameter Tuning increasing the accuracy of predictive model.

Project Requirement :

This project fundamentally aims in predicting the life expectancy. The primary requirement of the project is the suitable dataset which will aid the prediction. The dataset will provide various details like kind of diseases leading to the death. By using supervised machine learning techniques, we can extract a model that will be able to predict the life expectancy of future years.

Functional Requirement :

Create a data model present on the database. The data set are made available to the public to the purpose of health data analysis. It is related to the different countries depending on the different countries while finding the data set in different countries might be difficult and hence some countries are excluded from the final data set.

Technical Requirement :

The merged data set by using the databases in the .csv formats from Kaggle. Datasets need to be integrated into the Python IDE.

Software Requirements :

- Python IDE
- IBM Cloud
- IBM Watson
- IBM Node-Red

Project Deliverables :

The project is scheduled for 29 days from 15th May 2020 to 14th June 2020.

Project Team :

The project is done individually by Soham R. Sahare.

Phases of the development

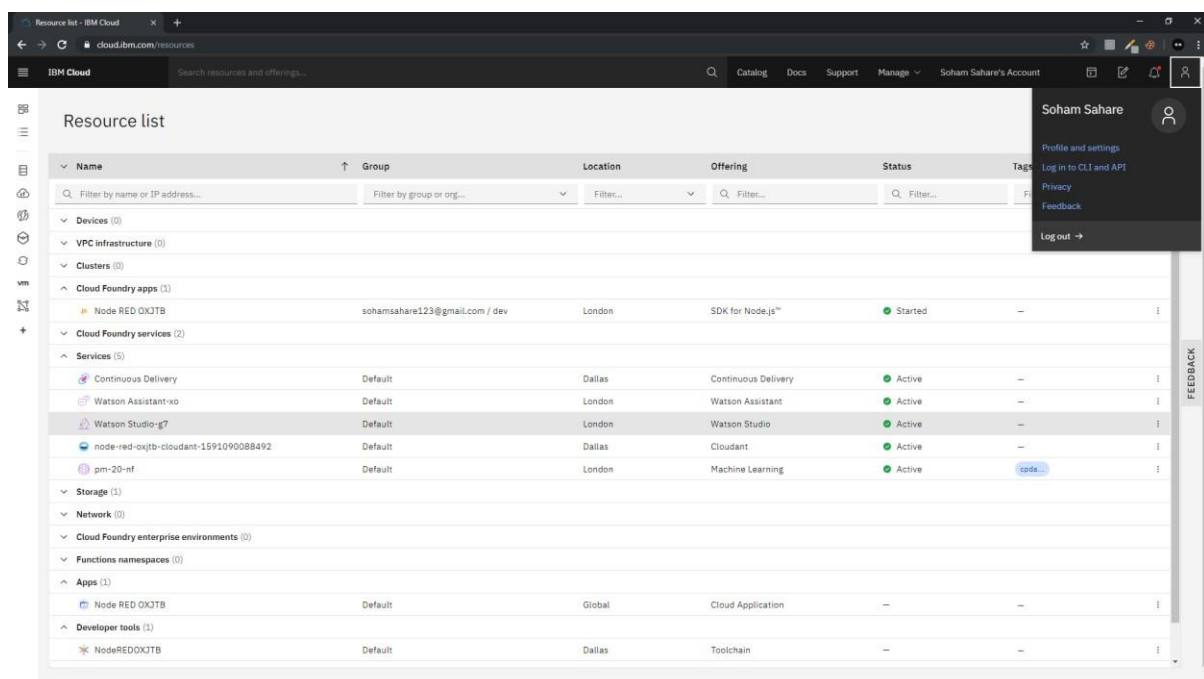
- 1) Project Planning & Kickoff :
 - Writing Project Scope, Schedule, Team & Deliverables for the project.
 - Setup the Development Environment – Creating GitHub account, Slack account, working with zoho writer and Understanding GitHub.
- 2) Explore IBM Cloud Platform :
 - Creating IBM account
 - Creating a Node-Red Starter Application
- 3) Explore IBM Watson Services :
 - Exploring IBM Watson usecases.
 - Exploring IBM Watson for Machine Learning.
- 4) Introduction To Watson Studio :
 - Building a Machine Learning model in IBM Watson Studio.
 - Automating Machine Learning Model.

5) Predicting Life Expectancy with Python :

- Collecting Dataset from [Kaggle](#).
- Creating Necessary IBM Cloud Services.
- Watson Studio Project
- Creating Machine Learning Service
- Creating a Jupyter Notebook in IBM Watson Studio and Import the dataset.
- Building a Machine Learning Model And Create Endpoints For Node-RED Integration.
- Build Node-RED Flow To Integrate ML Services

Some Screenshots while Project Development :

1) Resource list of all the apps and services developed in IBM Cloud.



4) Accuracy of different Machine Learning Models.

The screenshot displays a Jupyter notebook interface within IBM Watson Studio. The notebook is titled 'Predicting-Life-Expectancy-usin...' and is located in the 'My projects / Predicting-Life-Expectancy-usin...' workspace. The code is written in Python and compares the performance of three machine learning models: Linear Regression, Random Forest, and Decision Tree. The accuracy is measured using the R-squared score on a test set.

```

In [21]: model = LinearRegression(fit_intercept = True, normalize = True).fit(X_train, y_train)
         model.score(X_test, y_test)
Out[21]: 0.862909472654433

In [22]: random = RandomForestRegressor()
         random.fit(X_train, y_train)
         random.score(X_test, y_test)
Out[22]: 0.9587094816026089

In [23]: tree = DecisionTreeRegressor().fit(X_train, y_train)
         tree.score(X_test, y_test)
Out[23]: 0.912581611587117

In [24]: param_grid = [
         {'n_estimators': [10, 25],
          'max_features': [5, 10],
          'max_depth': [10, 50, None],
          'bootstrap': [True, False]}
         ]

         grid_search_forest = GridSearchCV(random, param_grid, cv=10, scoring='neg_mean_squared_error')
         grid_search_forest.fit(X_train, y_train)
Out[24]: GridSearchCV(cv=10, error_score='raise-deprecating',
          estimator=RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
          max_features='auto', max_leaf_nodes=None,
          min_impurity_decrease=0.0, min_impurity_split=None,
          min_samples_leaf=1, min_samples_split=2,
          min_weight_fraction_leaf=0.0, n_estimators=10, n_jobs=None,
          oob_score=False, random_state=None, verbose=0, warm_start=False),
          fit_params=None, iid='warn', n_jobs=None,
          param_grid=[{'n_estimators': [10, 25], 'max_features': [5, 10], 'max_depth': [10, 50, None], 'bootstrap': [True, False]}],
          pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
          scoring='neg_mean_squared_error', verbose=0)

In [25]: grid_search_forest.best_params_
Out[25]: {'bootstrap': False, 'max_depth': None, 'max_features': 10, 'n_estimators': 25}

In [30]: Random = RandomForestRegressor(bootstrap = False, max_depth = 50, max_features = 5, n_estimators = 25).fit(X_train, y_train)
         Random.score(X_test, y_test)
Out[30]: 0.9647692936688555

In [40]: from watson_machine_learning_client import WatsonMachineLearningAPIClient
  
```

The results show that the Random Forest model achieved the highest accuracy with an R-squared score of approximately 0.96. The Decision Tree model followed with an R-squared score of approximately 0.91, and the Linear Regression model had the lowest accuracy with an R-squared score of approximately 0.86.

5) Node-Red Deployed.

The screenshot shows a web browser window displaying a Node-RED deployment interface. The interface is titled "Machine Learning Model" and shows a "Prediction : 75.10589150562382". Below the prediction, there is a list of "Inputs" with their corresponding values. At the bottom, there are "SUMMIT" and "CANCEL" buttons.

Inputs	Value
Year	2008
Adult Mortality	1
Infant Deaths	1
Alcohol	5.61
Percentage Expenditure	36.622
Hepatitis B	99
Measles	0
BMI	52.6
Under five deaths	1
Polio	99
Total Expenditure	5.87
Diphtheria	99
ENR/AIDS	0.1
GDP	37.539
Population	2947314
chinness 1-19 years	1.6
chinness 5-9 years	1.6
chinness 1-19 years	1.6
chinness 5-9 years	1.6
Income composition of resource	0.713
Schooling	12

Developing

SUMMIT CANCEL

Conclusion :

From this project we can predict Life Expectancy with 96.47 % accuracy. Using Hyper parameterized Random Forest Regressor.

Important Links :

Node-Red URL : <https://node-red-oxjtb.eu-gb.mybluemix.net/ui/#!/0?socketid=YdWb57ES7FMjDH76AAAH>

Node-Red Editor URL : <https://node-red-oxjtb.eu-gb.mybluemix.net/red/#flow/a1e8fade.b01388>

Deployed video
: https://drive.google.com/file/d/1zj4NXGgM2DkH_eSawpeNB6SQiyef26p-/view?usp=sharing

Dataset : <https://www.kaggle.com/kumarajarshi/life-expectancy-who>