

Insurance Premium Prediction

Sahil & Siddhi

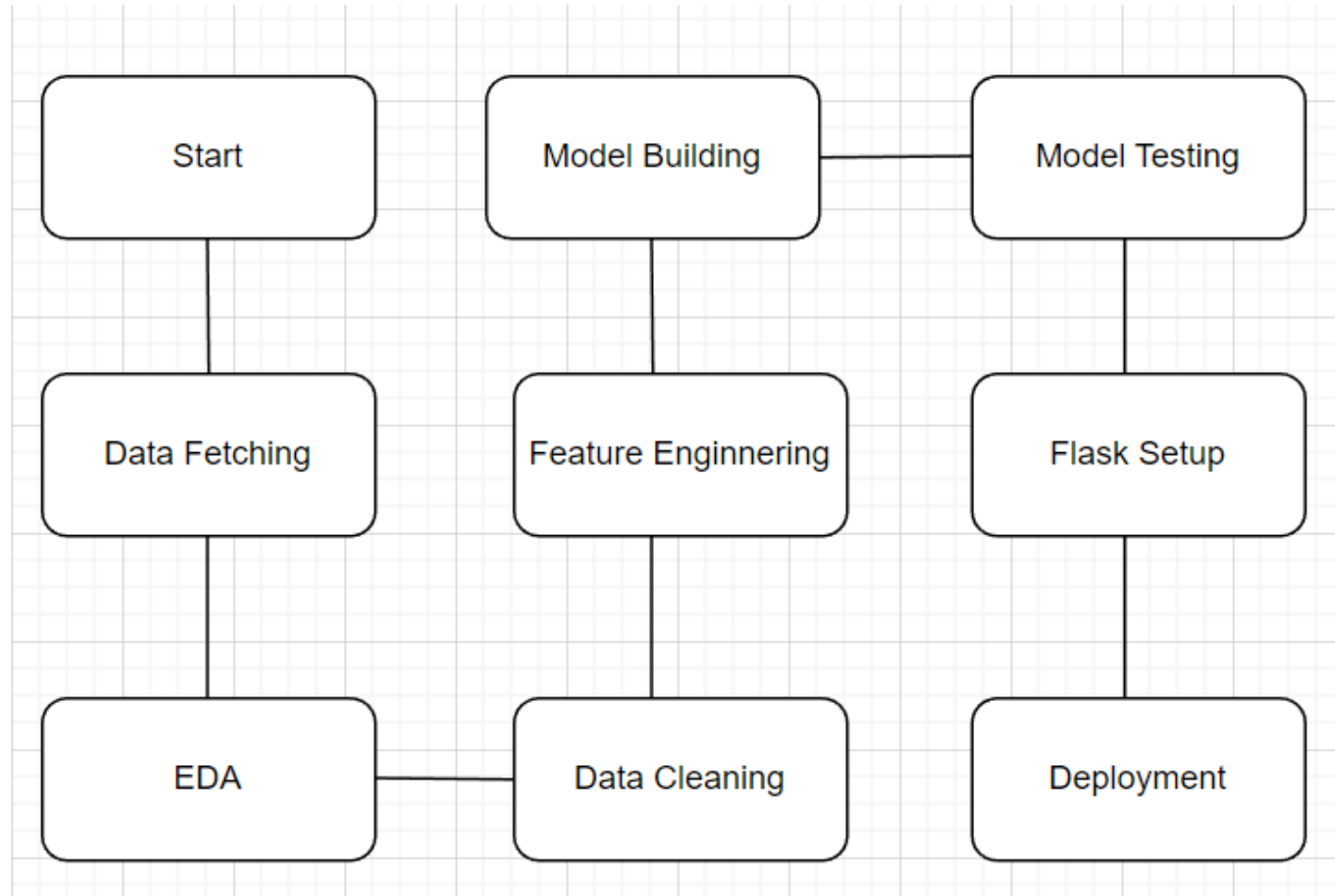
Objective :

The goal of this project is to give an estimate of how much they need on their individual health situation and Build a solution that should able to predict the premium of the personal for health insurance.

Benefits :

- Gets idea about how much amount required annually according to their own of health status.
- This can help a person in focusing more on the health aspect of an insurance.
- Help in giving premium of health insurance.

Architecture



Data Collection and validation

- The dataset was taken from the Kaggle competition page.
- Data type of columns – Validating the data type of the columns if wrong, then it was corrected.
- Null values in columns – Validating the column in the dataset have null values or missing information.

Model Training

➤ Data Pre-processing:

- Performing EDA to get insights of the data like identifying distribution, outliers etc.
- Check any null values present in the dataset. If present then impute those null values.
- Encode the categorical features/columns.
- Perform Standard Scalar to scale down values.

Model Selection

After pre-processing and model training, we find the best model for premium prediction. The model is trained on multiple regression algorithms like Linear Regression, Decision Trees, Random Forest, Gradient Boosting, Extreme Gradient Boosting and K-Nearest Neighbors (KNN). After prediction we will find accuracy of those predictions using evaluation metrics like RMSE (Root mean squared error) and `r2_score` (R-squared).

Predictions

Then all the trained models were used for validating test set.

We perform pre-processing techniques on it.

The best RMSE and r^2 score model were saved for developing API for prediction of premium.

Dataset Information:

DATASET LINK: [LINK](#)

kaggle

Create

Home

Competitions

Datasets

Models

Code

Discussions

Learn

More

Search

Sign In

Register

Insurance Premium Prediction

Data Card

Code (91)

Discussion (0)

101


New Notebook

Download (14 kB)

features (sex, smoker and region) that were converted into factors with numerical value designated for each level.

The purposes of this exercise to look into different features to observe their relationship,

age



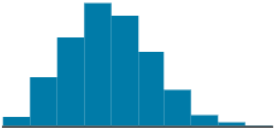
1864

sex

male51%


female49%

bmi



1653.1

children



0

19	female	27.9	0
18	male	33.8	1
28	male	33.0	3
33	male	22.7	0
32	male	28.9	0
31	female	25.7	0

We use cookies on Kaggle to deliver our services, analyze web traffic, and improve your experience on the site. By using Kaggle, you agree to our use of cookies.

Got it

Learn more

Web Interface

Welcome To Insurance Premium Predictor

Enter Age Of The Customer

28

Select Gender Of The Customer

female

Enter BMI Value Of The Customer

33.8

Enter Number Of The Children

0

Select Type Of The Customer(Smoker)

no

Select Region Of The Customer

southeast

Predict Premium

Premium: 4615.2748107641255

Activate Windows

Go to Settings to activate Windows

Welcome To Insurance Premium Predictor

Enter Age Of The Customer

Enter Weight Of The Customer

Select Gender Of The Customer

female

Enter BMI Value Of The Customer

Enter BMI Value Of The Customer

Enter Number Of The Children

Enter Number Of The Children

Select Type Of The Customer(Smoker)

no

Select Region Of The Customer

northeast

Predict Premium

Q & A

Q1) What is the source data?

The source of the data is Kaggle. The data is in the form of 'csv' file.

Q2) What was the type of the data?

The data was combination of categorical and numerical values.

Q3) What's the complete flow you followed in this project?

Refer the 3rd slide for better understanding

Q4) What techniques were you using for data pre-processing?

Visualizing relation of independent variables with each other and dependent variable.

Checking distribution of Continuous variables.

Checking any null values present in the dataset.

Converting categorical data into numeric values.

Scaling the data.

Q5) How training was done or what models were used?

Before training the model the dataset is divided into training set and testing/validation set.

The scaling was performed of training and validation set.

The categorical columns were converted into numeric values.

Algorithms like Linear Regression, Decision Trees, Random Forest, Gradient Boosting, KNN, and Extreme Gradient Boosting were used for model training and based on RMSE & r2_score the Gradient boosting model is saved for Validation.

Q6) How prediction was done?

On the basis of trained model, the prediction was performed. We also created API interface for estimating cost of premium on the basis of personal health information/status.

Q7) What are the different stages of deployment?

When the model is ready we deploy it in Heroku platform.