

Adult Census Income Prediction

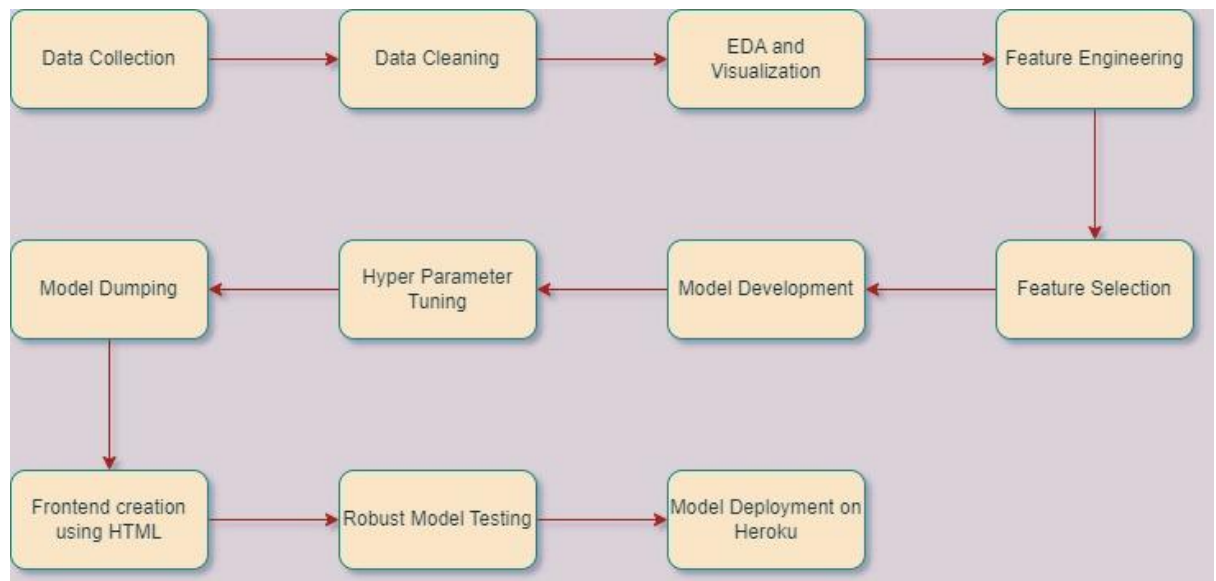
Architecture



Written By	Gurjeet Singh
Document Version	0.1
Description	Architecture
Date Issued	14 November, 2022

1. Architecture

Proposed methodology



2. Architecture Description

2.1 Data Collection

The dataset named Adult Census Income is available in Kaggle and UCRepository. This data was extracted from the 1994 census bureau dataset by Ronny Kohavi and Barry Becker (Data Mining and Visualization, Silicon Graphics). The dataset contains 15 columns and 32k+ rows. The dataset includes age, workclass, fnlwgt, education, education-num, marital-status, occupation, relationship, race, sex, capital-gain, capital-loss, hours-per-week, country, salary. This is given in the comma separated value format (.csv). The prediction task is to determine whether a person makes over \$50K a year or not.

2.2 Exploratory Data Analysis

Exploratory Data Analysis (EDA) is an approach to analyse the data using visual techniques. It is used to discover trends, patterns, or to check assumptions with the help of statistical summary and graphical representations.

In this step, we find relation between our dependent and independent features and also treat any outliers for our data to get prepared to be fed to our ML model.

2.3 Data pre-processing

In this step object data types are converted to int or float according to need and also correlation between different features are checked.

2.4 Feature Engineering

In this step categorical features are encoded using one hot encoder or label encoder and new features are created to be passed onto our model.

2.5 Feature Selection

Feature Selection is the method of reducing the input variable to our model by using only relevant data and getting rid of noise in data. It is the process of automatically choosing relevant features for our machine learning model based on the type of problem you are trying to solve.

2.6 Model Development

After cleaning the data and completing the feature engineering and selection, we split data into training and testing sets implement various classification algorithms like Logistic Regression, Random Forest, Gradient Boosting, Decision Tree, SVM, XGBoost, Catboost, and KNN to make the best possible model that can be made for accurate and correct prediction.

2.7 Hyperparameter Tuning

Hyperparameter tuning consists of finding a set of optimal hyperparameter values for a learning algorithm while applying this optimized algorithm to any data set. That combination of hyperparameters maximizes the model's performance, minimizing a predefined loss function to produce better results with fewer errors.

2.8 Model Dumping

After the training is completed a best model is selected among all the models by choosing an accuracy parameter according to our need and dumped the model in a pickle file format.

2.9 Frontend using HTML

In this step, a frontend is created using HTML to receive input from the user and connected to a flask backend server to do the predictions.

2.10 Robust Model Testing

The model is tested by giving an input data to check for its accuracy and response time.

2.11 Model Deployment on Heroku

The model is deployed on Heroku using a corn job on the server to keep the server and server code running forever.