

# Medical Insurance Premium Prediction



## ***Objectives:***

The main objective of this project is to predict medical insurance premiums for individuals based on their demographic and health-related attributes. The key objectives include:

Develop a machine learning model to predict insurance premiums accurately.

Create a user-friendly interface for individuals to input their details and obtain predicted insurance premiums.

Enhance the efficiency of insurance premium calculations for insurance companies.

## ***Benefits:***

- The project provides several benefits to different stakeholders:
  - **Individuals:** Allows individuals to estimate their insurance premiums based on their specific attributes, helping them make informed decisions about their insurance coverage.
  - **Insurance Companies:** Streamlines the insurance premium calculation process, enabling companies to provide accurate quotes to potential customers and improve customer satisfaction.
  - **Healthcare Providers:** Facilitates better financial planning by estimating insurance costs for medical procedures and treatments.
  - **Researchers:** Provides insights into the factors influencing insurance premiums and their impact on healthcare accessibility and affordability.

## ***Data Sharing Agreement:***

The dataset used for training the machine learning model was obtained from an insurance company, and a data sharing agreement was established to ensure the confidentiality and ethical use of the data. The agreement outlines the terms and conditions for accessing, storing, and analyzing the data, as well as the measures taken to protect sensitive information and comply with data privacy regulations.

Number of Columns: 7

**Attributes:** The dataset includes the following attributes:

- Age: The age of the individual.
- Gender: The gender of the individual (male or female).
- BMI: Body Mass Index, a measure of body fat based on height and weight.
- Children: The number of children the individual has.
- Smoker: Smoking status (yes or no).
- Region: The geographical region of the individual (Southeast, Southwest, Northeast, Northwest).
- Charges: The insurance charges incurred by the individual.

Columns Datatype: integer, string, float64, integer, string, string, float64

## ***Architecture:***

The project architecture consists of the following components:

**Data Collection:** Obtaining the insurance dataset from the insurance company.

**Data Preprocessing and EDA:** Cleaning the data, handling missing values, encoding categorical variables, and performing exploratory data analysis (EDA) to understand the data distribution and relationships.

- **Algorithm Selection:** The Random Forest Regressor algorithm was chosen for its ability to handle non-linear relationships and produce accurate predictions.
- **Feature Selection:** Features were selected based on their importance in predicting insurance premiums, as determined by feature importance scores.
- **Model Creation:** A Random Forest Regressor model was trained using the preprocessed dataset, with hyperparameters tuned using techniques such as grid search or random search.
- **Model Evaluation:** Model performance was evaluated using metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared score on both training and validation datasets to assess prediction accuracy and generalization.

## ***Data Preprocessing and EDA:***

- Handled missing values and outliers in the dataset.
- Encoded categorical variables using one-hot encoding.
- Visualized data distributions and relationships between variables using histograms, countplots , and correlation matrices.

## ***Model Development:***

We employ a Random Forest Regressor model for predicting insurance premiums. Random Forest is an ensemble learning algorithm that combines multiple decision trees to make predictions. We initialize the model with 100 decision trees and limit the maximum depth of each tree to 7 to prevent overfitting. The model is trained on the training dataset using the fit method.

## ***Random Forest Regressor:***

- The Random Forest Regressor model is employed for predicting insurance premiums.
- Random Forest is an ensemble learning algorithm that combines multiple decision trees to make predictions.
- It is initialized with 100 decision trees and limits the maximum depth of each tree to 7 to prevent overfitting.
- The model is trained on the training dataset using the fit method provided by the scikit-learn library.

## ***Model Evaluation:***

We evaluate the trained model using the testing dataset. Performance metrics such as R-squared score are calculated to assess the model's accuracy in predicting insurance premiums.

Additionally, we use the model to make predictions on a sample input data representing a 35-year-old male with a BMI of 35, one child, non-smoker, and residing in the Southeast region.

## ***Frequently Asked Questions (FAQs):***

**What factors influence insurance premiums?** Insurance premiums are influenced by factors such as age, gender, BMI, smoking status, number of children, and region.

**How accurate are the predictions?** The accuracy of predictions depends on the quality of input data and the performance of the machine learning model. The deployed model achieves a high level of accuracy in predicting insurance premiums.

**Is the user data secure?** Yes, user data is securely handled and processed in compliance with data privacy regulations. The data sharing agreement ensures the confidentiality and ethical use of user information.

**Can the model be updated with new data?** Yes, the model can be updated periodically with new data to improve its accuracy and relevance over time.

**How can I access the prediction tool?** The prediction tool is accessible through a user-friendly interface available on the project website. Users can input their details and obtain predicted insurance premiums instantly.

## ***Conclusion:***

The Medical Insurance Premium Prediction project provides an efficient and accurate solution for estimating insurance premiums based on individual attributes. By leveraging machine learning techniques and user-friendly interfaces, the project aims to enhance the insurance premium calculation process, benefiting individuals, insurance companies, healthcare providers, and researchers alike.