# A

# MACHINE LEARNING PROJECT REPORT

## ON

# Health Insurance Cost Prediction

# SUBMITTED BY

Harshita Rupani-220591

Ananya Srivastava- 220501

Radhika Bhati-220650

Deepanshu Aggarwal-220660

Rohit -220657

**UNDER MENTORSHIP OF**

**DR. MANISHA SAINI**



**BML MUNJAL UNIVERSITY™**

FROM HERE TO THE WORLD

**DEPARTMENT OF COMPUTER AND ENGINEERING**

# CANDIDATE'S DECLARATION

We hereby certify that the work on the project entitled, "Health Insurance Cost Analysis & Prediction" in partial fulfillment of requirements for the award of Degree of Bachelor of Technology in School of Engineering and Technology at BML Munjal University, is an authentic record of our own work carried out during a period from July 2022 to December 2022 under the supervision of **DR. MANISHA SAINI .**

1. Harshita Rupani-220591

2. Ananya Srivastava- 220501

3. Radhika Bhati- 220650

4. Deepanshu Aggarwal- 220660

5. Rohit- 220657

# SUPERVISOR'S DECLARATION

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

Faculty Supervisor Name:

Signature

# Table of Content

# ABSTRACT

This project endeavours to explore health insurance costs through the lens of data analytics and machine learning. Leveraging Python, Flask, and machine learning libraries, we conducted a thorough analysis and prediction of health insurance expenses.

The journey commenced with data collection, sourcing historical health insurance data spanning several years. We meticulously preprocessed this data to ensure its integrity and consistency, subsequently delving into exploratory data analysis. Employing visualization techniques, we dissected the data, uncovering trends, and correlations among various factors influencing insurance costs.

Our analysis revealed several key insights. Notably, factors such as smoking status, BMI, and age emerged as primary determinants of insurance charges, while variables like gender, children count, and region exhibited minimal impact. Armed with these insights, we embarked on building predictive models.

Utilizing machine learning algorithms including Linear Regression, Support Vector Regression, Ridge Regression, and Random Forest Regression, we trained predictive models to forecast insurance costs accurately. Through rigorous evaluation and optimization, we fine-tuned these models to achieve optimal performance.

Moreover, we deployed this model into a user-friendly web application, enhancing accessibility for stakeholders.

This project underscores the power of data-driven insights in understanding and predicting health insurance costs, offering valuable tools for both consumers and insurers to make informed decisions and optimize financial planning.

# Acknowledgment

We would like to take this opportunity to express my heartfelt gratitude to all those who have contributed to the completion of this report. Firstly, I want to extend my deepest appreciation to my supervisor, Dr. Manisha Saini, for her exceptional guidance, motivating ideas, and unwavering support throughout the fabrication process and the drafting of this report. Her insights and feedback have been invaluable in shaping the direction of this project.

We would also like to thank all the individuals who participated, contributed, and will continue to contribute to this project. Their advice and assistance have been instrumental in ensuring its success, and I am truly grateful for their ongoing support.

We are immensely grateful to all those who played a part in this journey. Any shortcomings or omissions in the report are solely our responsibility.

# List Of Figures

# List Of Table

| | |
|---|---|
| 01 | Comparison of Results from Different Prominent Research Paper in Health Insurance Cost Analysis. |

# List Of Abbreviations

| | |
|---|---|
| EDA | Exploratory Data Analysis |
| ROI | Return on Investment |
| APIs | Application Programming Interface |
| NLP | Natural Language Processing |
| GCP | Google Cloud Platform |

# 1. Introduction

Healthcare costs are rising, and accurately predicting individual health insurance premiums is crucial for both insurers and policyholders. It remains a significant challenge due to the complex interplay of various factors such as age, gender, BMI, number of children, smoking status, and geographic region. Traditional methods often fail to account for these complexities, leading to inaccurate cost estimations and financial strain for both insurers and policyholders. This project addresses the need for a more precise and reliable prediction model to ensure fair pricing and better financial planning.

## 1.1. Overview

This project delves into the critical domain of health insurance cost analysis and prediction, harnessing the capabilities of Python and its powerful libraries like NumPy, Pandas, Matplotlib, and Seaborn. It's structured into distinct phases, beginning with data exploration, preprocessing, model development, and concluding with the integration of a predictive model into a web application using Flask.

During the data exploration phase, an in-depth analysis of an insurance dataset is conducted to unravel its characteristics and underlying patterns. Key features such as smoking status, BMI, and age are scrutinized using visualizations like heatmaps and bar plots, providing valuable insights into their impact on insurance charges.

Following data exploration, the dataset undergoes preprocessing to ensure data integrity and reliability. This involves transforming categorical features into numerical representations and standardizing numerical variables for compatibility with predictive models.

The core of the project lies in the development of predictive models employing various regression algorithms such as Linear Regression, Support Vector Regression (SVR), Ridge Regression, and Random Forest Regression. Through rigorous evaluation metrics like R-squared scores, root mean squared error (RMSE), and cross-validation, the performance of each model is assessed to identify the most effective algorithm for predicting health insurance costs based on individual characteristics.

Furthermore, the project goes beyond model development by seamlessly integrating the predictive model into a user-friendly web application using Flask. This integration enables users to make real-time predictions of health insurance costs through an intuitive interface.

In essence, this project represents a fusion of data science, machine learning, and web development aimed at addressing challenges in healthcare cost analysis and prediction. By leveraging advanced technologies and methodologies, it strives to empower stakeholders with actionable insights, fostering informed decision-making, and improving the efficiency of healthcare systems.

## 1.2. Existing System

Before the integration of machine learning in health insurance cost prediction, the industry primarily relied on conventional methods involving manual analysis and statistical models. Analysts would manually gather data from various sources, perform extensive analysis, and base insurance pricing decisions on their expertise and intuition.

However, these traditional approaches encountered several limitations:

- Scalability Challenges: Manual analysis was time-consuming and resource-intensive, making it challenging to efficiently analyze large datasets containing diverse variables and parameters.

- Subjectivity: Decisions based on human judgment were subjective and susceptible to biases, potentially leading to inaccuracies in insurance pricing and risk assessment.

- Limited Insights: Traditional models often failed to capture intricate patterns and correlations within healthcare data, resulting in incomplete analyses and suboptimal predictions.

- Inefficiency: Manual analysis struggled to keep pace with the dynamic nature of healthcare data and evolving market trends, leading to delays in decision-making and missed opportunities for cost optimization.

The emergence of machine learning has revolutionized health insurance cost prediction by mitigating these challenges. Machine learning algorithms can efficiently process large volumes of healthcare data, uncover hidden insights, and make accurate predictions about insurance costs based on individual characteristics and risk factors.

By leveraging historical insurance data, machine learning models can learn from past patterns and trends, enabling insurers to forecast future healthcare expenses more accurately. These models can also adapt to evolving market dynamics and regulatory requirements, ensuring that insurance pricing remains relevant and competitive over time.

## 1.3 User Requirement Analysis

- **Prediction Accuracy:** Users expect machine learning models to deliver accurate predictions of future healthcare expenses based on individual characteristics and risk factors. Rigorous evaluation metrics such as R-squared scores, mean squared error, and cross-validation should be employed to assess prediction accuracy and reliability.

- **Interpretability:** Explanation of model predictions and insights should be presented in a clear and understandable manner, catering to users with varying levels of data science expertise. Interpretability tools such as feature importance plots and decision explanations can enhance user understanding.

- **Scalability:** The system should be capable of processing large volumes of healthcare data efficiently, ensuring smooth performance even when analyzing extensive datasets. Scalable machine learning algorithms and cloud computing resources may be employed to handle scalability requirements effectively.

- **Reliability and Security:** Ensuring the reliability and security of healthcare data is paramount. Measures such as data encryption, access controls, and regular backups should be implemented to protect sensitive patient information and maintain data integrity.

- **Integration Capabilities:** Users may want to integrate analysis results with other healthcare systems or decision-making platforms.

# 2. Literature Review

In contemporary economies, a substantial fraction of the GDP is allocated to healthcare expenditure, often reaching around 30%. Developed nations typically exhibit the highest healthcare spending, both in absolute terms and as a percentage of their economies. A significant contributor to this expenditure is the government's provision of healthcare services, particularly through programs like Medicare, which covers a substantial portion of medical costs for the elderly population. However, with the impending retirement of the baby boomer generation and their subsequent eligibility for Medicare, coupled with the escalating costs of healthcare, there is mounting pressure on public finances. Consequently, it becomes imperative to leverage all available resources to contain healthcare expenses.

Machine learning offers a promising avenue for managing healthcare costs through predictive analytics. By developing models capable of forecasting medical expenditures, stakeholders can proactively guide patients towards more affordable healthcare options. Moreover, policymakers can identify providers with disproportionately higher costs and potentially implement corrective measures. In this study, the Random Forest Regression algorithm is employed to predict medical expenses. Additionally, experiments with alternative machine learning models such as Linear Regression are conducted to compare their efficacy in cost prediction. Early estimation of healthcare costs not only facilitates financial planning for individuals but also guards against the risk of unwittingly purchasing unnecessary or overly expensive health insurance coverage.

The literature on the application of machine learning in healthcare cost prediction is still emerging. By leveraging advanced analytical tools, researchers aim to address the pressing challenge of escalating healthcare expenditures, thereby promoting greater affordability and accessibility in healthcare delivery.

## 2.1 Comparison

-Comparison of Results from Different Prominent Research in Health Care Analysis.

| Research Paper Title | Authors | Techniques Used | Key Findings |
|---|---|---|---|
| Comparisons of SVM Kernels for Insurance Data Clustering. (2021) | Irfan Nurhidayat ,Busayamas Pimpunchat ,Samad Noeiaghdam ,Unai Fernández-Gámiz | Support Vector Machine, RMSE ,AUC. | The Support Vector Machine (SVM) approach emerged as a more precise and effective alternative, particularly after evaluating all three prominent SVM kernels. It was determined that utilizing the radial kernel yielded the most accurate results, characterized by low Root Mean Square Error (RMSE) and high density. |
| Health Insurance Cost Prediction Using Regression Models (2022). | Sudhir ,Biswajit Purkayastha; Dolly Das; Manomita Chakraborty; Saroj Kumar Biswas | Regerssion Techniques | Machine Learning regression models, notably Polynomial Regression, demonstrate their utility in accurately predicting health insurance costs, as evidenced by an RMSE of 5100.53, an R2 value of 0.80, |

| | | | and an accuracy rate of 80.97%, thereby streamlining operations and ensuring reliability for insurance companies. |
|---|---|---|---|
| Advanced Intelligence Health Insurance Cost Prediction Using Random Forest. (2023) | Sai Srinivas Vellela, D Pushpalatha, G Sarathkumar, C.H. Kavitha, D Harshithkumar. | Random Forest | A study employed various machine learning regression models to predict medical insurance costs, with the Stochastic Gradient Boosting (SGB) model outperforming others, achieving an 86% accuracy rate. This underscores the efficacy of advanced ML algorithms, like SGB, in enhancing predictive accuracy within healthcare insurance, informing strategic decision-making in the industry. |

*Table -01*

This table provides a comparative overview of the results obtained from different prominent research studies in the field of Health care analysis. Each study employs different techniques and methodologies, resulting in varying findings and insights.

# 3. Objectives of Project

The objective of this project is to develop and deploy a robust machine learning solution, namely the Health Insurance Cost Predictor, designed to accurately estimate individual health insurance expenses based on a comprehensive set of demographic and lifestyle factors. By leveraging advanced algorithms and data analysis techniques, the project aims to address several key objectives:

1. **Enhancing Fairness and Transparency:** By utilizing a diverse range of input features including age, gender, BMI, number of children, smoking status, and geographic region, the model seeks to mitigate biases and provide fair and transparent cost estimations.

2. **Empowering Informed Decision-Making**: The Health Insurance Cost Predictor serves as a valuable tool for both insurance companies and individuals. Insurance providers can utilize the model to accurately assess risk and determine fair premiums, thereby fostering better financial planning and risk management strategies.

3. **Improving Accuracy and Reliability**: Through rigorous data exploration, preprocessing, and model evaluation techniques, the project aims to achieve high levels of prediction accuracy and reliability. By comparing the performance of various machine learning algorithms such as Linear Regression, Support Vector Regression (SVR), Ridge Regressor, and Random Forest Regressor, the model identifies the most suitable approach for predicting health insurance costs with precision.

4. **Facilitating Model Deployment and Accessibility:** The deployment of the Random Forest Regressor model via Flask enables seamless integration into a

user-friendly web interface, accessible to both insurance stakeholders and the general public. This deployment enhances the accessibility and usability of the predictor, allowing users to obtain personalized cost estimations conveniently and efficiently.

Moving forward, the project could explore avenues for continuous improvement, such as integrating additional features or exploring alternative modeling techniques to further enhance prediction accuracy.

# 4 . Exploratory Data Analysis

## 4.1  Dataset:

The dataset utilized in this project comprised essential attributes crucial for predicting health insurance costs. These attributes included age, sex, BMI (Body Mass Index), number of children, smoking status, region, and the corresponding insurance charges. The dataset was acquired from a reliable source and consisted of a substantial number of observations, ensuring its representativeness and reliability for analysis.

To procure the dataset, standard data acquisition methods were employed, utilizing reputable sources such as online repositories or curated datasets. The dataset underwent preprocessing to ensure data integrity and suitability for analysis. This involved cleaning the data, handling missing values, and encoding categorical variables into numerical representations.

**Dataset Link:  https://osf.io/7u5gy**

## 4.2 Exploratory Data Analysis and Visualizations:

During the exploratory data analysis (EDA) phase, a detailed examination of the dataset's characteristics, structure, and distributions was conducted. Various visualizations were crafted to unravel insights and patterns within the data, aiding in feature engineering and model selection.

- **Age vs. Charge:** A bar plot was created to depict the relationship between age and insurance charges, exploring how age influences insurance costs. This visualization revealed potential trends or correlations between age groups and insurance expenses.

- **Region vs. Charge:** Another bar plot was generated to analyze the impact of geographic region on insurance charges, providing insights into regional variations in healthcare costs.

- **BMI vs. Charge:** A scatter plot was utilized to examine the relationship between BMI and insurance charges, with differentiation based on gender, aiming to uncover potential associations between body mass index and healthcare expenses.

- **Smoker vs. Charge:** A bar plot illustrating the disparity in insurance charges between smokers and non-smokers, segregated by gender, was created. This visualization highlighted the substantial impact of smoking habits on healthcare expenditure.

- **Sex vs. Charges:** A bar plot was crafted to explore differences in insurance charges between genders, aiming to discern any gender-based disparities in healthcare costs.

Additionally, the EDA phase involved analyzing the skewness and kurtosis of each attribute to understand their distributions and identify potential outliers. Data preprocessing techniques, including scaling of BMI and charges columns, were also performed to prepare the data for predictive modeling.

## 4.3 Related Sections:

Following the exploratory data analysis (EDA) phase, several related sections were undertaken to further enhance the project's predictive capabilities and glean actionable insights. These sections included:

### 4.3.1 Data Preprocessing:

Data preprocessing techniques were applied to ensure the dataset's integrity and suitability for analysis. This involved handling missing values, encoding categorical variables, and scaling numerical features to standardize their range.

### 4.3.2 Feature Engineering:

Feature engineering techniques were utilized to extract valuable information from the dataset and create new features that could improve predictive performance. This may have involved creating interaction terms, transforming variables, or selecting relevant subsets of features.

### 4.3.3 Model Selection:

Multiple regression algorithms were explored to identify the most suitable model for predicting health insurance costs. Model selection criteria may have included performance metrics, computational efficiency, and interpretability.

# 5. **Methodology**

## 5.1 Introduction to Languages:

1. **Python:** Python is the primary language used for data analysis, machine learning, and backend development in your project. Its simplicity, readability, and extensive libraries such as Pandas, NumPy, and Scikit-learn make it well-suited for handling and analyzing large datasets, implementing machine learning algorithms, and processing data efficiently.

2. **HTML (Hypertext Markup Language):** HTML is utilized for creating the structure and content of web pages in your project. It defines the elements and layout of the user interface, including forms, buttons, text, and other components necessary for presenting information to users.

3. **CSS (Cascading Style Sheets):** CSS is employed to style and design the appearance of your web application. It defines the presentation, layout, and formatting of HTML elements.

4. **JavaScript:** JavaScript adds interactivity and dynamic behavior to your web application. It enables you to implement client-side logic, handle user inputs, perform validations, update content dynamically without reloading the page, and interact with backend services asynchronously.

5. **Flask:** Flask is a lightweight and flexible web framework for Python used for building web applications. It provides tools and utilities for routing requests, handling HTTP methods, rendering templates, and managing sessions, making it easier to develop and deploy web applications.

## 5.2 Constraints :

Our methodology for health insurance cost prediction faces several constraints. Firstly, reliance on historical data may introduce biases if the dataset lacks diversity or fails to capture emerging trends. Secondly, the complexity of healthcare data poses challenges for traditional regression-based approaches, necessitating more advanced modelling techniques. Additionally, ethical considerations regarding data privacy and security must be addressed rigorously. Finally, the interpretability of complex models like neural networks may limit stakeholders' ability to understand and act upon predictions effectively, emphasizing the need for a balance between accuracy and interpretability.

## 5.3 Use Case Model/Flow Chart/ :



Fig.1 - Use Case Diagram for Health Insurance Cost Analysis

Fig.2 - Flowchart for Health Insurance Cost Analysis

## 5.4 Dependencies:

- **NumPy:** Python library for numerical operations and array manipulation.
- **Pandas:** Python library for data manipulation and analysis, with data frame structures.
- **Matplotlib:** Python plotting library for creating visualizations.
- **Seaborn:** Statistical data visualization library with high-level plotting functions.
- **RandomForestRegressor:** Ensemble learning algorithm for regression tasks.
- **SVR:** Support Vector Regression algorithm for capturing nonlinear relationships.
- **scikit-learn:** Python machine learning library with various algorithms and metrics.
- **cross_val_score:** Function for cross-validation.
- **RandomizedSearchCV and GridSearchCV**: Functions for hyperparameter tuning.
- **pickle:** Module for serializing and deserializing Python objects.

## 5.5 Deployed Models:

- **Linear Regression:** A straightforward yet effective algorithm for predicting continuous target variables by assuming a linear relationship between input

22

features and the target. It's often used as a baseline model due to its simplicity and interpretability.

- **Support Vector Regression (SVR):** Based on Support Vector Machine (SVM) principles, SVR is adept at capturing nonlinear relationships between variables by mapping them into a higher-dimensional space. It seeks to find the optimal hyperplane to minimize error while maximizing the margin.

- **Ridge Regression:** A regularization technique that combats multicollinearity and overfitting by adding a penalty term to the standard least squares objective function. It constrains the coefficients of regression variables, improving stability and performance, particularly with correlated predictors.

- **Random Forest Regression:** An ensemble learning method that constructs multiple decision trees during training and outputs the average prediction of these trees. It's robust against overfitting and can handle both numerical and categorical data, making it suitable for capturing complex interactions and nonlinear relationships.

## 5.6  Implementation of Algorithm :

### 5.6.1 Steps Taken:

**1. Data Exploration:** Conduct in-depth exploratory analysis of health insurance datasets to unveil underlying patterns, dependencies, and influential factors shaping insurance charges.

**2. Preprocessing:** Implement meticulous preprocessing techniques, including the conversion of categorical features into numerical representations, handling missing
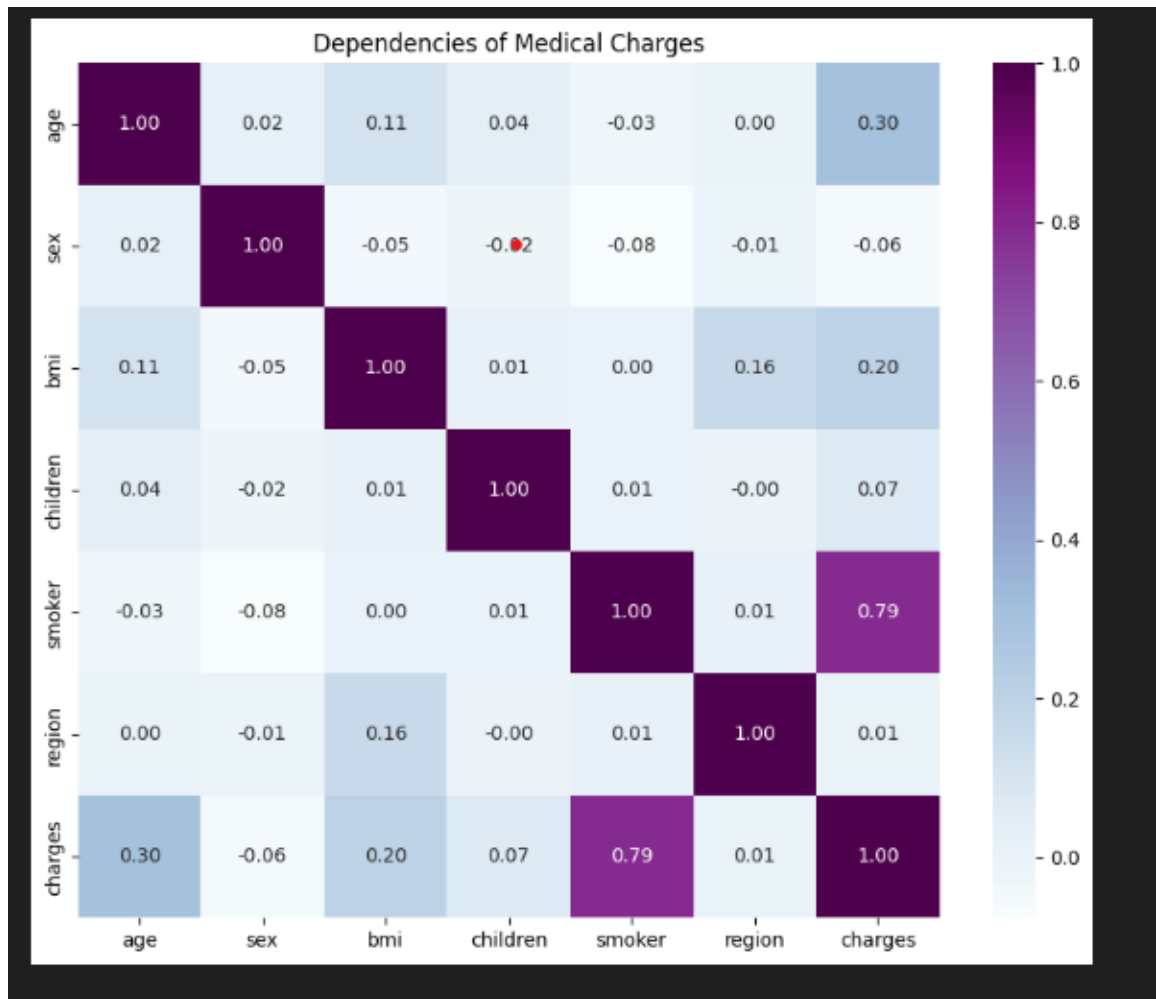
data, and standardizing numerical variables. These measures ensure data integrity and compatibility with subsequent machine learning algorithms.

**3. Model Development:** Forge and train machine learning models such as Linear Regression, Support Vector Regression (SVR), Ridge Regression, and Random Forest Regression. These models are tailored to predict health insurance costs based on individual characteristics like age, BMI, smoking status, and regional attributes.

**4. Evaluation:** Employ stringent evaluation metrics such as R-squared scores, root mean squared error (RMSE), and cross-validation to meticulously assess the performance of each machine learning model. This rigorous evaluation process identifies the most effective algorithm for precise prediction.

**5. Integration:** Seamlessly integrate the predictive model into an intuitive web application using Flask. This integration empowers users to interact directly with the model, facilitating real-time predictions of health insurance costs based on their unique input parameters.
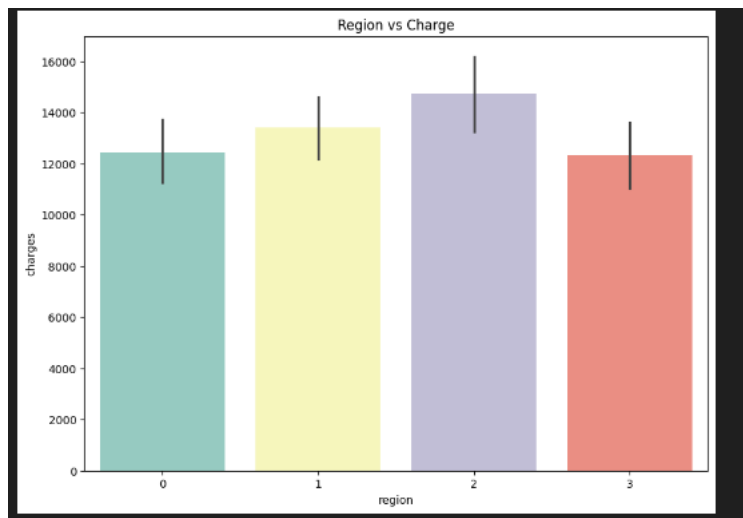
Dependencies of Medical Charges

This heatmap depicts the correlations between factors affecting health insurance charges. The color intensity shows the correlation strength, with red indicating positive correlation and blue negative correlation. Strong positive correlations are seen between charges and age, BMI, and smoker status. This suggests that these factors are important for predicting health insurance costs.
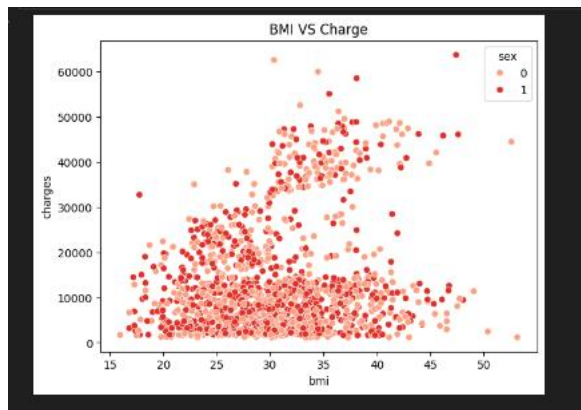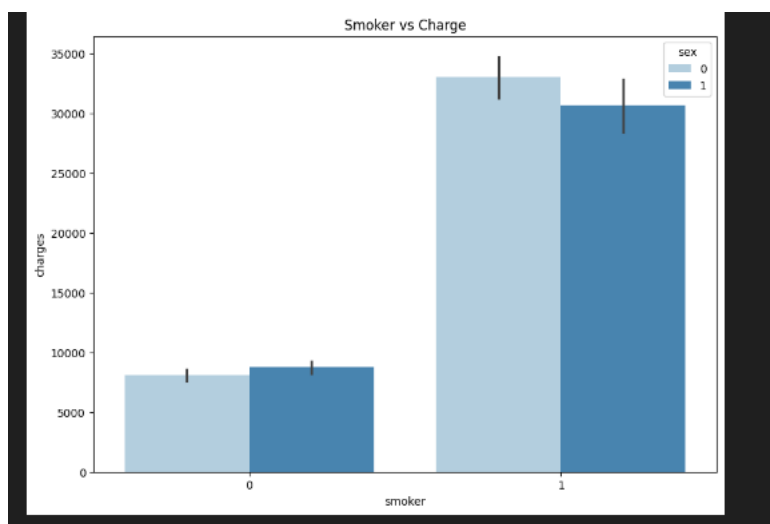
# EXPLORATRORY DATA ANANLYSIS

Above Bar is Known as Error Bar , Here we have analyse affect of charges against age , as age is increasing insurance charges increase. The height of bar represent mean and the line above bar represent variance from mean.



Analyzing these regional variations can aid in predicting future health insurance costs.

The graph depicts there is a positive correlation between BMI and health insurance costs. This means that as a person's BMI increases, the charges they incur tend to increase as well. This relationship can be explained by the fact that people with higher BMIs are more likely to develop certain health conditions, such as heart disease, diabetes, and some types of cancer, which can lead to higher medical costs.
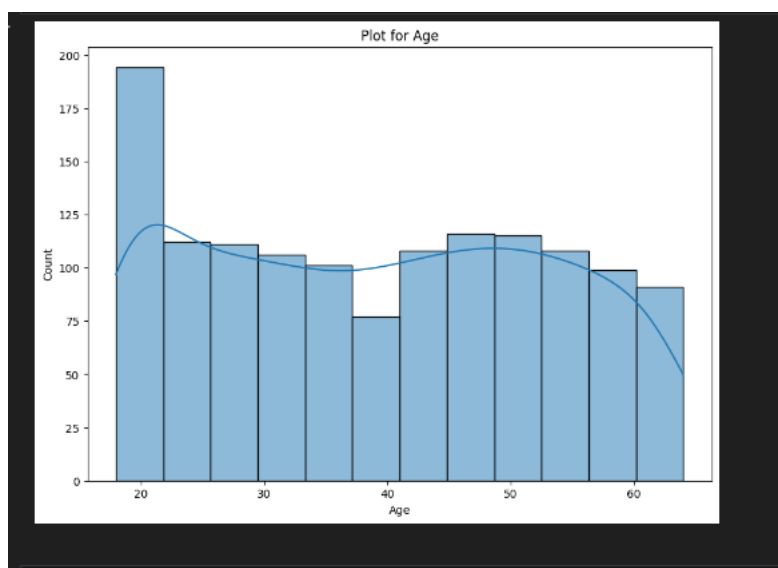


This bar chart depicts the average health insurance charges for smokers and non-smokers . The clear distinction in charge heights suggests a strong correlation between smoking status and healthcare costs, which can be valuable for predicting health insurance premiums.
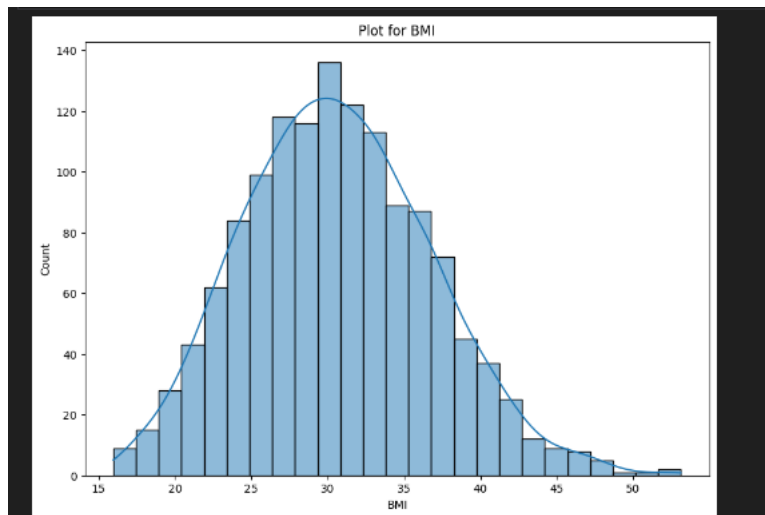
This bar chart shows average health insurance charges by gender (0-male, 1-female). While gender plays a role, it's likely one piece of a larger puzzle for predicting costs. Other factors like age and health history likely hold more weight
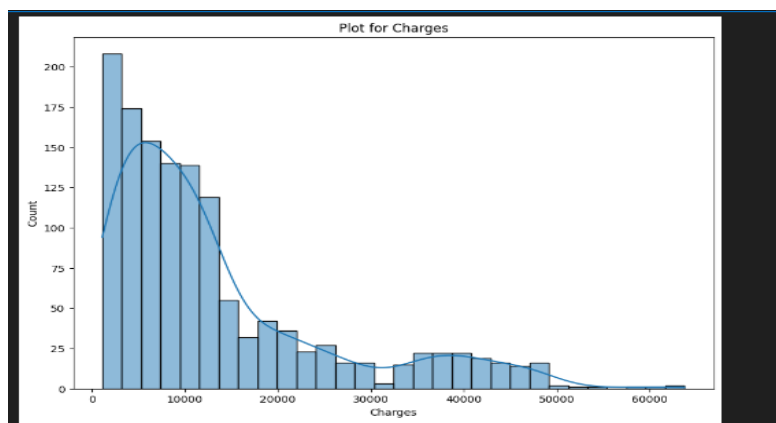
## Plotting Skews & Kurtosis



This histogram shows the distribution of ages in your health insurance data. The shape suggests a non-normal distribution with a peak (kurtosis). This indicates that certain age groups might be more prevalent than others in your dataset, which could be relevant when building health insurance cost prediction models.
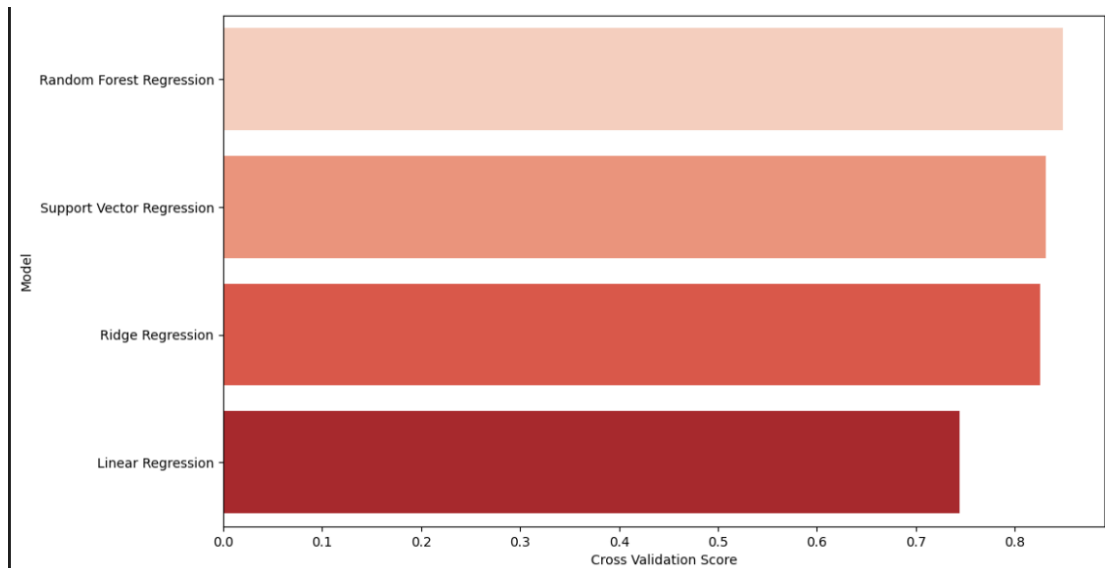
Plot for BMI

This histogram shows the distribution of Body Mass Index (BMI) in your health insurance dataset. The graph exhibits kurtosis, with a sharp peak in the middle and steeper drop-offs on either side. This indicates that a significant portion of the population falls within a specific BMI range, potentially around a healthy weight range. This information on BMI distribution can be valuable for predicting health insurance costs.



Plot for Charges

This line graph depicts a distribution of health insurance charges that's likely left-skewed. The higher values concentrated on the left side suggest a greater number of individuals with lower charges, with the charges decreasing as we move towards the right side.

## MODEL EVALUATION:

The graph above reveals Random Forest Regression as the frontrunner. Its higher score suggests it's the most accurate choice for predicting health insurance costs within our dataset. This indicates that Random Forest Regression effectively learns the relationships between factors like age, BMI, and smoker status with insurance costs.

# 6. RESULTS

## Working Model And Website

**Predicted Amount :**

Expected amount is 19079.694

# 7. **Conclusion**

In conclusion, our project on health insurance cost prediction successfully navigated through various phases, including data exploration, preprocessing, model development, and integration into a web application. By leveraging Python and its powerful libraries, we gained valuable insights into the factors influencing health insurance charges and developed accurate predictive models to estimate these costs based on individual characteristics.

The project yielded several significant results and implications for future research and development in the domain of healthcare cost analysis:

1**. Predictive Accuracy:** Our predictive models, including Linear Regression, Support Vector Regression (SVR), Ridge Regression, and Random Forest Regression, demonstrated robust performance in estimating health insurance costs. Through rigorous evaluation metrics, we ensured the accuracy and reliability of our models, empowering stakeholders with precise predictions.

**2. Feature Importance:** Exploratory data analysis uncovered essential features such as age, BMI, smoking status, and geographic region, which significantly impact health insurance charges. Understanding the relative importance of these factors enhances our understanding of cost drivers and informs targeted interventions and policy decisions.

**3. Integration with Web Application:** The integration of our predictive model into a user-friendly web application using Flask enhances accessibility and usability for end-users. This seamless integration allows for real-time predictions of health insurance costs, fostering informed decision-making and improving user experience.

# 8. Future Scope

Looking ahead, our project on health insurance cost prediction holds significant potential for future advancements and applications. Firstly, leveraging emerging technologies such as artificial intelligence and machine learning, we can enhance the predictive accuracy and granularity of our models. By incorporating sophisticated algorithms and exploring ensemble learning techniques, we can capture intricate relationships among various health factors, demographic variables, and insurance charges. This advancement would enable insurers to offer more tailored and competitive insurance plans, better aligning with individual policyholders' needs and risk profiles.

Furthermore, there is ample scope for incorporating real-time data streams and IoT (Internet of Things) devices into our predictive framework. By integrating wearable health trackers, smart medical devices, and electronic health records (EHRs), we can create a dynamic and adaptive system capable of continuously monitoring policyholders' health status and adjusting insurance premiums accordingly. This proactive approach to risk management not only improves the accuracy of cost predictions but also empowers individuals to make informed decisions about their healthcare and insurance coverage.

Moreover, Collaborating with healthcare providers, policymakers, and public health organizations can facilitate the exchange of valuable data insights, leading to more holistic and evidence-based approaches to healthcare delivery and resource allocation.

# Plag.docx

**5**% SIMILARITY INDEX

**2**% INTERNET SOURCES

**2**% PUBLICATIONS

**4**% STUDENT PAPERS

PRIMARY SOURCES

| | | |
|---|---|---|
| 1 | www.springerprofessional.de<br>Internet Source | 1% |
| 2 | archive.org<br>Internet Source | 1% |
| 3 | Submitted to University of Glamorgan<br>Student Paper | 1% |
| 4 | Submitted to Asia Pacific International College<br>Student Paper | <1% |
| 5 | Submitted to University of Huddersfield<br>Student Paper | <1% |
| 6 | Submitted to Northwestern State University<br>Student Paper | <1% |
| 7 | Submitted to Higher Education Commission Pakistan<br>Student Paper | <1% |
| 8 | Submitted to Kent Institute of Business and Technology<br>Student Paper | <1% |
| 9 | Submitted to Middlesex University | |