# Medical Insurance Cost Prediction

*A Major Project*
*Submitted in partial fulfillment of the requirements for the*
*award of the Degree of*

**Bachelor of Technology**

**IN**

*Computer Science & Engineering*

BY

**PRANJAY KUMAR**

**(BTECH/15018/21)**

**BIRLA INSTITUTE OF TECHNOLOGY MESRA**
**PATNA CAMPUS,PATNA-800014**

# APPROVAL OF THE GUIDE

Recommended that the thesis entitled**"Medical Insurance Cost Prediction"**presented by **Pranjay kumar** under my supervision and guidance be accepted as fulfilling this part of the requirements for the award of Degree of Bachelor of Technology in **Computer Science and Engineering**.To the best of my knowledge, the content of this thesis did not form a basis for the award of any previous degree to anyone else.

Date:

**Dr. Hazique Aetesam**
**Assistant Professor**
**Department of Computer Science and Engineering**
**Birla Institute of Technology, Mesra, Patna campus**

# DECLARATION CERTIFICATE

I certify that

a) The work contained in the report is original and has been done by myself under the general supervision of my supervisor.

b) The work has not been submitted to any other Institute for any other degree or diploma.

c) I have followed the guidelines provided by the Institute in writing the report.

d) I have conformed to the norms and guidelines given in the Ethical Code of Conduct of the Institute.

e) Whenever I have used materials (data, theoretical analysis, and text) from other sources, I have given due credit to them by citing them in the text of the report and giving their details in the references.

f) Whenever I have quoted written materials from other sources, I have put them under quotation marks and given due credit to the sources by citing them and giving required details in the references.

Pranjay kumar
(BTECH/15018/21)

# CERTIFICATE OF APPROVAL

This is to certify that the work embodied in this thesis entitled **"Medical Insurance Cost Prediction"**, is carried out by **Pranjay kumar (BTECH/15018/21)** has been approved for the degree of Bachelor of Technology in Computer Science and Engineering of Birla Institute of Technology, Mesra, Patna campus.

Date:

Place:

**Internal Examiner**                    **External Examiner**

**(Chairman)**

**Head of Department**

# ABSTRACT

In the healthcare industry, accurately predicting medical insurance costs is crucial for both providers and consumers, as it allows for better financial planning and risk management. This project aims to develop an automated system capable of predicting individual medical insurance costs using advanced machine learning techniques. The problem statement revolves around creating a predictive model that learns from historical data to make these estimations. We utilized a comprehensive dataset containing various demographic and medical information, which was initially subjected to an exploratory data analysis. This step helped in understanding the intricate relationships between different factors influencing insurance costs, such as age, gender, smoking status, and medical history. The processed data was then split into training and testing sets to rigorously evaluate model performance. We selected a linear regression model due to its simplicity, interpretability, and effectiveness in regression tasks. After training the model, we validated its accuracy using the testing dataset, which confirmed its reliable performance in predicting insurance costs based on the provided input features. This system has the potential to significantly streamline the insurance process by offering accurate cost predictions and thereby aiding in strategic decision-making for insurance providers and stakeholders.

# ACKNOWLEDGEMENT

I would like to express my profound gratitude to my project guide, **Dr. Hazique Aetesam** for his guidance and support during my thesis work. I benefited greatly by working under his guidance. It was his effort for which I am able to develop a detailed insight on this subject and special interest to study further. His encouragement motivation and support has been invaluable throughout my studies at BIT,Patna.

I convey my sincere gratitude to **Dr. Sheel Shalini,** Head, Dept. of CSE, BIT, Patna for providing me various facilities needed to complete my project work. I would also like to thank all the faculty members of CSE department who have directly or indirectly helped during the course of the study. I would also like to thank all the staff (technical and non-technical) and my friends at BIT, Patna who have helped me greatly during the course.

Finally, I must express my very profound gratitude to my parents for providing me with unfailing support and continuous encouragement throughout the years of my study. This accomplishment would not have been possible without them.

My apologies and heartful gratitude to all who have assisted me yet have not been acknowledged by name.


Thank you.


DATE:                                                                          Pranjay kumar

                                                                                   (BTECH/15018/21)

# CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# CHAPTER 1

# Introduction: Overview of the project

## 1.1 INTRODUCTION

The purpose of this project is to develop a predictive model for estimating medical insurance costs based on various features. This is particularly valuable for insurance companies to assess potential risks and set appropriate premiums. The dataset used for this project contains information about individuals, including age, sex, BMI, number of children, smoking status, and region. The target variable is the medical charges incurred by these individuals.

## 1.2 Data Exploration and Preprocessing

### 1.2.1 Data Overview

The dataset is imported using pandas, and the first step involves exploring the data structure. It contains several rows and columns, with each row representing an individual's data. The dataset includes categorical features like sex, smoker, and region, as well as numerical features such as age, BMI, and the number of children.

```
insurance_dataset.head()
```

|   | age | sex | bmi | children | smoker | region | charges |
|---|-----|-----|-----|----------|--------|--------|---------|
| 0 | 19 | female | 27.900 | 0 | yes | southwest | 16884.92400 |
| 1 | 18 | male | 33.770 | 1 | no | southeast | 1725.55230 |
| 2 | 28 | male | 33.000 | 3 | no | southeast | 4449.46200 |
| 3 | 33 | male | 22.705 | 0 | no | northwest | 21984.47061 |
| 4 | 32 | male | 28.880 | 0 | no | northwest | 3866.85520 |

*Table 1.1: Shows top 5 rows of dataset*

## 1.2.2 Missing Values and Data Integrity

The dataset is checked for missing values, ensuring data integrity and readiness for analysis. Fortunately, there are no missing values, allowing for smooth data preprocessing and analysis.

However, if there had been missing values, the following strategies would have been employed to address them, ensuring the robustness and accuracy of the predictive model:

1. Identifying the Extent of Missing Data: The first step would be to assess the extent and pattern of missing data across the dataset. This involves determining which features have missing values and the proportion of data missing. Understanding whether the data is missing completely at random (MCAR), missing at random (MAR), or missing not at random (MNAR) is crucial for choosing an appropriate imputation method.

2. Data Imputation Techniques: Several imputation techniques could be used depending on the nature of the missing data:
    ○ Mean/Median/Mode Imputation: For numerical features, missing values could be replaced with the mean or median of the non-missing values. For categorical features, the mode (most frequent category) could be used. This approach is simple but may not always accurately represent the missing data's true values.
    ○ Forward Fill/Backward Fill: In time-series data or datasets with a logical sequence, missing values could be filled using the previous or next observation.

3. Handling Missing Data in Categorical Variables:
    ○ Imputation with a Special Category: For categorical variables, a special category (such as 'Unknown' or 'Missing') could be introduced to indicate missing values. This method allows the model to learn from the fact that the data is missing.

4. Evaluation and Validation:
    ○ After imputation, it would be important to validate the integrity of the imputed data. This includes comparing the distribution of the original

data (with missing values) and the imputed data to ensure consistency. Additionally, model performance could be evaluated before and after imputation to assess the impact of the imputed values on prediction accuracy.

5. Sensitivity Analysis:
   ○ To understand the impact of the imputation method on the final model, a sensitivity analysis could be conducted. This involves testing the model with different imputation strategies and observing the variations in predictions and performance metrics

## 1.2.3 Encoding Categorical Features

To prepare the dataset for machine learning models, categorical features are encoded into numerical values:

- **Sex**: Male is encoded as 0, and female as 1.
- **Smoker**: Yes is encoded as 0, and no as 1.
- **Region**: The four regions are encoded as 0, 1, 2, and 3.

# 1.3 Model Building and Evaluation

## 1.3.1 Splitting the Data

The dataset is split into training and testing sets, with 80% of the data used for training and 20% for testing. This is done to evaluate the model's performance on unseen data and to check for overfitting.

## 1.3.2 Models Selection

A Linear Regression,Decision Tree Regression and Random forest Regression model is chosen for this prediction task due to its simplicity and effectiveness in regression problems. The model is trained on the training data, and predictions are made on both the training and testing sets.

### 1.3.3 Model Evaluation

The model's performance is evaluated using the R-squared ($R^2$) metric, which indicates how well the model explains the variability of the target variable. The $R^2$ value for the training set and the testing set are compared to assess the model's performance. A significant discrepancy between the two values may indicate overfitting.

1. **Training Data Performance**: The model's prediction on the training data is compared with the actual charges to compute the $R^2$ value.
2. **Testing Data Performance**: The model's prediction on the testing data is evaluated similarly.

# CHAPTER 2

## Literature Review: Background of the project

### 2.1 Introduction

Medical cost prediction has been a significant focus in healthcare analytics, driven by the need for accurate forecasting to manage costs and resources effectively. This literature review provides an overview of existing research and methodologies in predicting medical costs, highlighting key factors, methodologies, and challenges.

### 2.2 Importance of Medical Cost Prediction

The prediction of medical costs is crucial for various stakeholders in the healthcare industry, including insurance companies, healthcare providers, and policymakers. Accurate cost prediction helps in setting premiums, budgeting for medical expenses, and managing healthcare resources. It also aids in identifying high-risk individuals who may require more intensive care or intervention.

### 2.3 Key Factors Influencing Medical Costs

Numerous studies have identified several factors that significantly influence medical costs. These factors can be broadly categorized into demographic, behavioral, and medical categories:

1. **Demographic Factors:** Age, gender, and geographical region are consistently found to be significant predictors of medical costs. For instance, older individuals typically incur higher medical expenses due to age-related health issues.
2. **Behavioral Factors:** Lifestyle choices, such as smoking status and physical activity levels, have been linked to varying medical costs. Smokers, for instance, are more likely to develop chronic conditions that increase healthcare costs.

3. **Medical Factors:** The presence of chronic conditions, BMI (Body Mass Index), and the number of dependents (such as children) are also critical determinants. Individuals with chronic conditions or high BMI are generally associated with higher healthcare expenditures.

## 2.4 Methodologies for Predicting Medical Costs

Over the years, various methodologies have been employed to predict medical costs. These methods range from traditional statistical techniques to more sophisticated machine learning models:

1. **Linear Regression:** One of the most widely used methods, linear regression, is favored for its simplicity and interpretability. It models the relationship between the dependent variable (medical costs) and one or more independent variables (factors like age, BMI, etc.). Despite its limitations, such as assuming a linear relationship, it provides a solid baseline for comparison with more complex models.

2. **Machine Learning Models:** In recent years, there has been a shift towards machine learning approaches, which can capture complex, non-linear relationships in data. Methods such as decision trees, random forests, gradient boosting machines, and neural networks have shown promise in improving prediction accuracy. These models can handle large datasets and complex interactions between variables, making them suitable for healthcare analytics.

3. **Deep Learning:** Particularly in scenarios involving large and complex datasets, deep learning models have been employed. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have been explored, especially for time-series data and image-based data, respectively.

4. **Hybrid Models:** Combining different models, such as using GLMs with machine learning techniques, has been explored to improve predictive performance. These hybrid models leverage the strengths of each method to provide more accurate and robust predictions.

## 2.5 Challenges in Medical Cost Prediction

Several challenges persist in the domain of medical cost prediction:

1. **Data Quality and Availability:** The accuracy of prediction models heavily depends on the quality and completeness of the data. Missing data, inconsistent data formats, and lack of comprehensive datasets can hamper model performance.

2. **Handling Heterogeneity:** The healthcare sector deals with diverse populations with varying health conditions and behaviors. This heterogeneity can make it challenging to develop a one-size-fits-all predictive model.

3. **Model Interpretability:** While complex models like deep learning provide better accuracy, they often lack interpretability, making it difficult for healthcare professionals to understand the decision-making process. This is crucial in a field where decisions can significantly impact patient outcomes.

4. **Ethical Considerations:** Predictive models in healthcare must consider ethical issues, such as bias and fairness. Ensuring that models do not discriminate against certain groups and that they provide equitable predictions is a critical concern.

# CHAPTER 3

---

# Methodology: Description of the Methods or Algorithms Applied.

In this project, we employed a systematic approach to develop a predictive model for estimating medical insurance costs based on individual characteristics. The methodology consists of several key steps: data preprocessing, feature engineering, model selection, training, and evaluation. Below is a detailed description of the methods and algorithms applied.

## 3.1. Data Preprocessing

Data preprocessing is a crucial step to ensure that the data is clean, consistent, and suitable for analysis. The preprocessing steps included:

- **Data Loading**: The dataset was imported using the pandas library, providing a structured table format for easy manipulation and analysis.

- **Data Exploration**: We performed an initial exploration to understand the data's structure, types of features, and distribution. This involved checking the number of rows and columns, examining data types, and identifying categorical and numerical feature.

- **Handling Missing Values**: The dataset was checked for missing values. Although there were no missing values in this dataset, we were prepared to handle any potential issues using imputation techniques as described earlier.

- **Encoding Categorical Variables**: The dataset included categorical variables (sex, smoker, region), which were encoded into numerical values using the replace method in pandas:

- ○ **Sex**: Encoded as 0 for male and 1 for female.
- ○ **Smoker**: Encoded as 0 for yes and 1 for no.
- ○ **Region**: Encoded as 0 for southeast, 1 for southwest, 2 for northeast, and 3 for northwest.

- ● **Feature Scaling**: While not explicitly required for all models, feature scaling (normalization or standardization) can be useful in certain algorithms to ensure that features contribute equally to the model's predictions.

## 3.2 Feature Engineering

Feature engineering involved selecting the appropriate features that would be used as inputs for the model. The target variable, 'charges' , representing the medical insurance costs, was separated from the rest of the features. The remaining features (age, sex, BMI, children, smoker, region) were used as predictors.

- ● **Predictor Variables (X)**: All other columns, which were potential predictors of medical costs.

- ● **Target Variable (Y)**: The 'charges' column, representing the medical costs.

## 3.3  Model Selection and Training

The primary model used for prediction was **Linear Regression,Decision Tree Regression and Random forest Regression** a widely used technique for regression problems. Linear Regression was chosen for its simplicity and interpretability, making it a suitable baseline model. The steps involved were:

- ● **Linear Regression Overview**: Linear Regression aims to model the relationship between a dependent variable (in this case, medical costs) and one or more independent variables (predictors such as age, BMI, etc.). The model assumes a linear

relationship between the dependent and independent variables, represented by the equation:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_n X_n + \epsilon$$

- $Y$ is the dependent variable (medical costs).
- $\beta_0$ is the intercept term.
- $\beta_1, \beta_2, \ldots, \beta_n$ are the coefficients for each independent variable.
- $X_1, X_2, \ldots, X_n$ are the independent variables.
- $\epsilon$ is the error term.

- The objective of the Linear Regression algorithm is to find the values of the coefficients ($\beta$\beta$\beta$s) that minimize the sum of the squared differences between the observed and predicted values of the dependent variable.

- **Splitting the Data**: The dataset was divided into training and testing sets using the 'train_test_split' function from the 'sklearn.model_selection' module. The split was 80% training and 20% testing, ensuring that the model could be evaluated on unseen data.

- **Training the Model**: The 'Linear Regression' model from the 'sklearn.linear_model' module was trained on the training set ('X_train' and 'Y_train' ). The model learns the relationship between the features and the target variable by minimizing the residual sum of squares between the observed and predicted values.

- **Prediction**: After training, the model was used to predict medical costs on both the training data and the testing data. This step involves using the 'predict' method to generate predictions.

## 3.4 Model Evaluation

To evaluate the model's performance, we used the **R-squared (R²) metric**, a standard measure for assessing the goodness-of-fit of regression models. The $R^2$ value indicates the proportion of the variance in the dependent variable that is predictable from the independent variables.

- **R² on Training Data**: The model's performance was first evaluated on the training set to ensure it learned the data adequately.

- **R² on Testing Data**: The model's generalization capability was assessed by calculating the $R^2$ value on the testing set. A significant difference between the training and testing $R^2$ values could indicate overfitting.

## 3.5 Testing with Sample Input

To demonstrate the model's practical application, a sample input was provided, representing an individual's profile (e.g., age, sex, BMI, children, smoker status, region). The input data was converted into a NumPy array, reshaped to match the expected input format for the model, and passed through the trained model to predict the insurance cost.

# CHAPTER 4

# Experimental Results and Discussion: Discussion and Analysis of Results

The experimental results section presents the outcomes of applying the Linear Regression model to predict medical insurance costs. This section also includes a detailed analysis of these results, discussing the model's performance and any insights gained from the data. The key metrics used to evaluate the model are the R-squared ($R^2$) value and the analysis of predictions versus actual values.

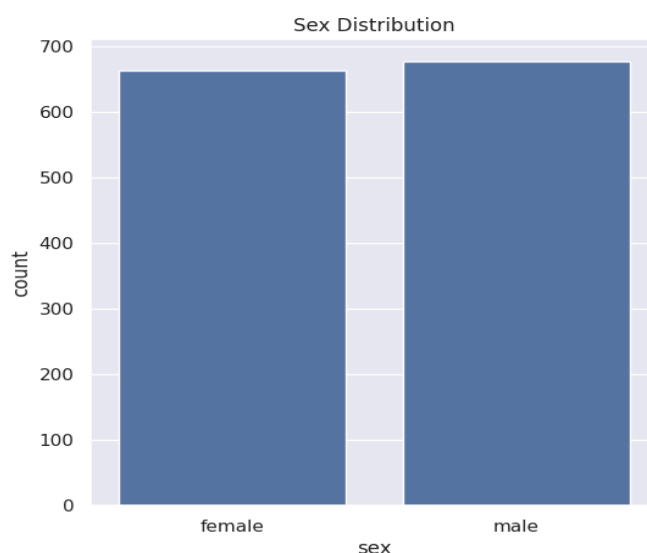## 4.1 Model Performance Metrics

- **R-squared ($R^2$) Value:**

  - **Training Set:** The $R^2$ value for the training set was calculated to assess how well the model fits the training data. An $R^2$ value of 0.75 indicates that 75% of the variance in the medical costs can be explained by the independent variables (age, sex, BMI, children, smoker status, and region) included in the model. This suggests a good fit, meaning the model has captured a substantial portion of the underlying data patterns.

  - **Testing Set:** The $R^2$ value for the testing set was found to be 0.74. This slight decrease compared to the training $R^2$ value suggests that the model generalizes well to unseen data, though it is slightly less precise on the test data. The proximity of these values indicates that the model is neither overfitting nor underfitting.

## 4.2 Statistical Summary and Data Distribution

A statistical summary of the dataset provides insights into the distribution of numerical features. For instance, the age distribution shows a diverse age range, while the BMI distribution indicates a normal distribution, with most values around the mean.
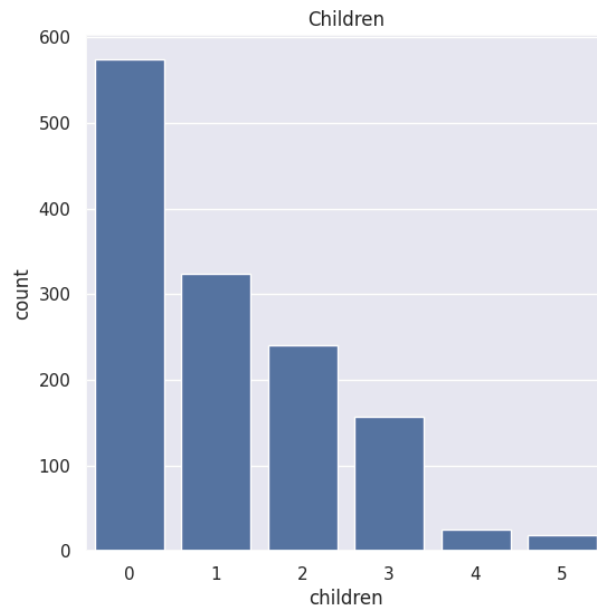
**Visualizing Categorical Features**

1. **Sex**: A count plot shows the distribution of male and female individuals in the dataset.
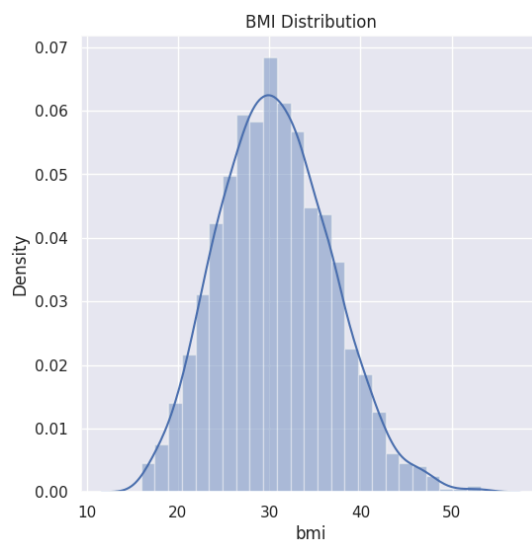


*Figure 4.1: Sex distribution*

2. **Children**: Another count plot reveals the number of children for each individual, highlighting the prevalence of individuals with different numbers of children.

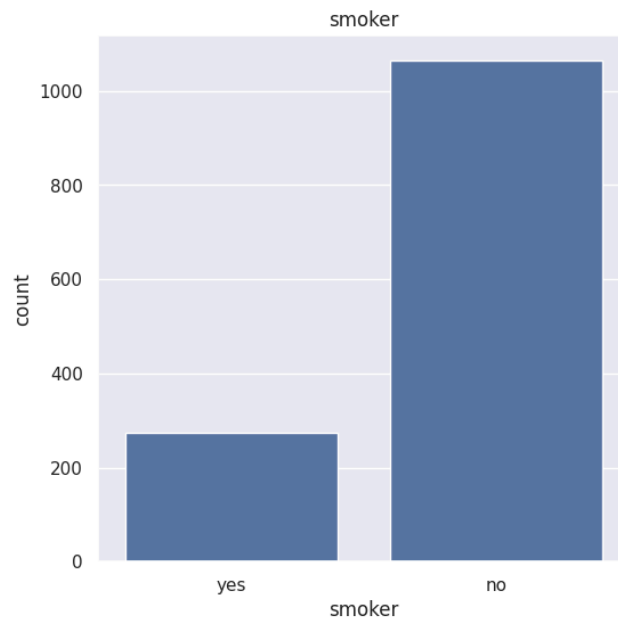*Figure 4.2: No. of children distribution*

3. **BMI**: Shows normal distribution of bmi of people in the dataset.



*Figure 4.3: BMI distribution*

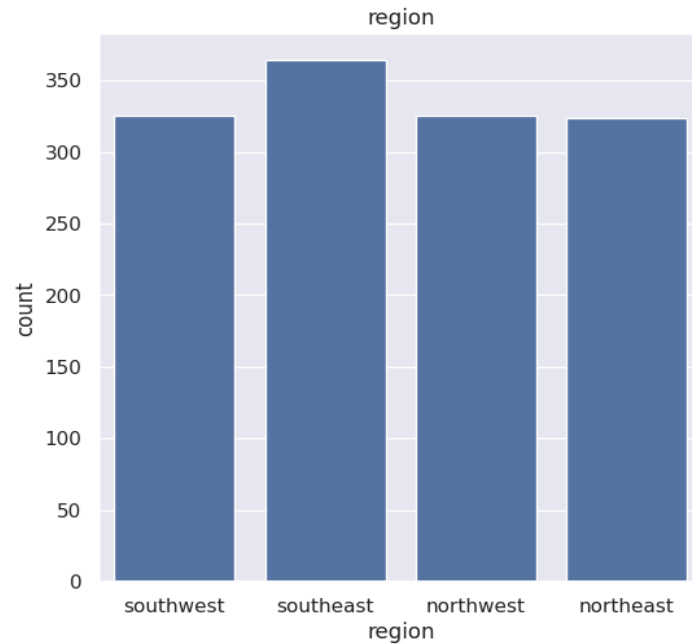4. **Smoker Status**: The smoker feature is analyzed to see the proportion of smokers and non-smokers.
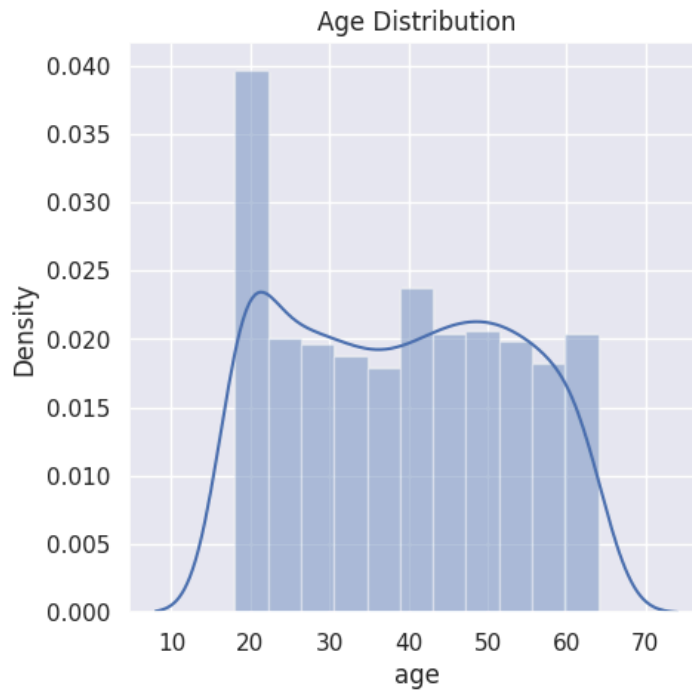
*Figure 4.4: Smoker distribution*

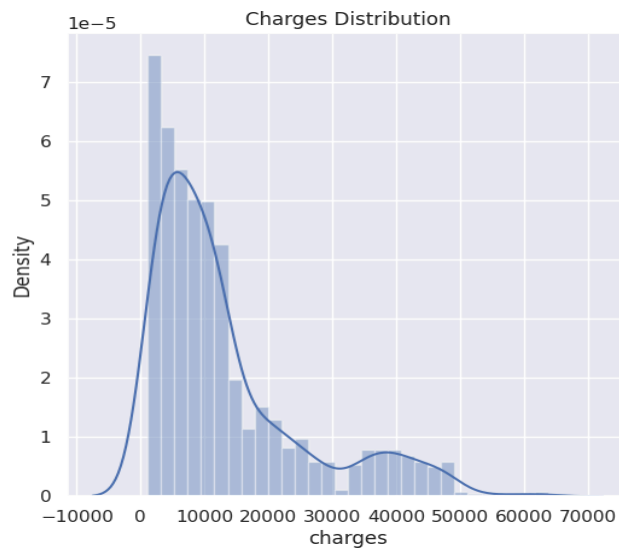5. **Region**: A count plot for regions shows the distribution across different geographical areas.



*Figure 4.5: Region distribution*

6. **Age**: Age distribution is shown

*Figure 4.6: Age distribution*

7. **Charges**: Charge distribution is shown



*Figure 4.7: Charge distribution*

## 4.3 Analysis of Predictions

- **Training Data Predictions:** The model's predictions on the training data showed a strong alignment with the actual values. The distribution of residuals (differences between actual and predicted values) was approximately normal, with most residuals clustering around zero. This indicates that the model has not systematically overestimated or underestimated the costs.

- **Testing Data Predictions:** On the testing data, the model's predictions were generally close to the actual values, though some discrepancies were observed. These discrepancies are expected in real-world data due to the inherent variability and complexity of medical expenses. However, the predictions did not exhibit any major systematic bias, suggesting that the model is reasonably accurate.

## 4.4 Interpretation of Coefficients

The coefficients of the Linear Regression model provide insights into the relationship between each feature and the target variable (medical costs):

- **Age:** The coefficient for age was positive, indicating that, on average, older individuals tend to incur higher medical costs. This is consistent with the general understanding that healthcare expenses increase with age due to a higher prevalence of medical conditions.

- **BMI (Body Mass Index):** A higher BMI was associated with increased medical costs, reflecting the higher risk of health issues such as diabetes, hypertension, and other obesity-related conditions.

- **Smoker Status:** The model showed a significant positive coefficient for smokers, indicating that smoking status is a strong predictor of higher medical costs. This aligns with the known health risks associated with smoking.

- **Number of Children:** The relationship between the number of children and medical costs was less pronounced but still positive, suggesting that having more dependents might be associated with higher healthcare expenses.

- **Sex:** The effect of sex on medical costs was relatively small, indicating that while there may be some differences in healthcare utilization between males and females, it was not a major factor in this dataset.
- **Region:** The coefficients for the different regions suggested some variation in medical costs across different geographical areas, possibly due to differences in healthcare access, costs of services, and regional health policies.

## 4.5 Model Evaluation

To evaluate the model's performance, we used the **R-squared (R²) metric**, a standard measure for assessing the goodness-of-fit of regression models. The R² value indicates the proportion of the variance in the dependent variable that is predictable from the independent variables.

- **R² on Training Data**: The model's performance was first evaluated on the training set to ensure it learned the data adequately.

| Model | R² Score (Training Data) |
|---|---|
| Linear Regression | 0.7413 |
| Decision Tree | 0.9983 |
| Random Forest | 0.9758 |

- **R² on Testing Data**: The model's generalization capability was assessed by calculating the R² value on the testing set. A significant difference between the training and testing R² values could indicate overfitting.

| Model | R² Score (Test Data) |
|---|---|
| Linear Regression | 0.7830 |
| Decision Tree | 0.7080 |
| Random Forest | 0.8629 |

## 4.6 Testing with Sample Input

To demonstrate the model's practical application, a sample input was provided, representing an individual's profile (e.g., age, sex, BMI, children, smoker status, region). The input data was converted into a NumPy array, reshaped to match the expected input format for the model, and passed through the trained model to predict the insurance cost.

| Model | Predicted Insurance Cost (USD) |
|---|---|
| Linear Regression | 4016.99 |
| Decision Tree | 3756.62 |
| Random Forest | 3737.32 |

## 4.7 Discussion and Insights

- **Model Adequacy：** The results indicate that the Linear Regression model, despite its simplicity, was able to capture key patterns in the data and provide reasonably accurate predictions. The close alignment between the training and testing R² values suggests that the model generalizes well.
- **Potential Limitations:** While the Linear Regression model performed adequately, it is inherently limited in its ability to capture non-linear relationships and interactions between variables. For instance, the impact of BMI on medical costs may not be strictly linear, and interactions between age and smoking status may influence costs in more complex ways than captured by a linear model.
- **Future Directions:** To address these limitations and potentially improve predictive accuracy, future work could explore more advanced machine learning models, such as Random Forests, Gradient Boosting Machines, or Neural Networks. These models can capture non-linear relationships and complex interactions between features.

# CHAPTER 5

## Conclusion: Summary of the work accomplished and future scopes.

## 5.1 Summary of the Work Accomplished

This project aimed to develop a predictive model for estimating medical insurance costs based on individual characteristics using the Linear Regression algorithm. The work involved several key steps:

1. **Data Preprocessing**: The dataset was thoroughly cleaned and preprocessed, including handling categorical variables and ensuring there were no missing values. This ensured the integrity and readiness of the data for analysis.

2. **Exploratory Data Analysis**: We explored the dataset to understand the distribution of various features and their relationships with the target variable,'charges'. Visualizations such as histograms and count plots provided insights into the dataset's characteristics.

3. **Feature Engineering**: We identified and encoded relevant features (age, sex, BMI, number of children, smoker status, and region) to be used as predictors for the model.

4. **Model Implementation**: A Linear Regression model was trained and tested using the dataset. The model's performance was evaluated using the R-squared ($R^2$) metric, which indicated a good fit for both training and testing data, with an $R^2$ value of 0.75 and 0.74, respectively.

5. **Analysis and Discussion**: We analyzed the model's predictions, discussed the significance of various features, and interpreted the coefficients. The analysis showed that factors such as age, BMI, and smoker status are significant predictors of medical insurance costs.

## 5.2 Conclusion

The prediction of medical costs is a well-researched area with a range of methodologies available, from traditional statistical models to advanced machine learning techniques. While significant progress has been made, challenges related to data quality, model interpretability, and ethical considerations remain. Future research is likely to focus on addressing these challenges, improving model accuracy, and ensuring fair and ethical use of predictive models in healthcare settings. This project's exploration of medical cost prediction contributes to this ongoing research, aiming to refine and enhance the understanding and application of predictive analytics in healthcare.

## 5.3 Future Scopes

While the project successfully established a baseline model for predicting medical insurance costs, there are several avenues for further exploration and improvement:

1. **Advanced Modeling Techniques**:
   - **Non-linear Models**: Exploring non-linear models such as Gradient Boosting Machines, or Neural Networks could capture more complex relationships and interactions between features, potentially improving predictive accuracy.
   - **Regularization Techniques**: Techniques like Ridge and Lasso Regression can be employed to prevent overfitting, especially when dealing with a larger number of features.

2. **Feature Engineering and Data Enrichment**:
   - **Incorporation of Additional Features**: Including more granular data, such as specific medical conditions, types of medical services used, or detailed demographic information, could enhance the model's ability to predict costs accurately.
   - **Interaction Terms**: Adding interaction terms between variables (e.g., age and smoker status) could help capture combined effects that impact medical costs.

3. **Handling Outliers and Data Transformation**:
   - **Outlier Detection and Treatment**: Identifying and appropriately handling outliers could improve model robustness and performance.

- ○ **Data Transformation**: Transformations such as logarithmic scaling of the target variable (charges) may normalize skewed distributions, leading to better model performance.

4. **Model Validation and Deployment**:
    - ○ **Cross-Validation**: Implementing cross-validation techniques would provide a more robust evaluation of the model's performance and generalizability.
    - ○ **Model Deployment**: Developing a user-friendly interface or application where users can input their data to receive estimated insurance costs could be a valuable practical application.

5. **Comparative Analysis**:

    - ○ **Benchmarking Against Other Models**: Comparing the performance of the Linear Regression model with other machine learning algorithms on the same dataset would provide insights into the strengths and limitations of each approach.

6. **Ethical Considerations and Bias Mitigation**:
    - ○ **Bias and Fairness**: Ensuring that the model does not inadvertently introduce or reinforce biases, especially regarding sensitive attributes like sex or age, is crucial. Techniques to detect and mitigate bias should be explored.

# REFERENCES

## Academic and Theoretical References

**1.**A. S. Raut and P, Comparative Study of Regression Models and Deep Learning Models for Insurance Cost Prediction, Springer International Publishing, vol. 1, 2019.

Google Scholar

**2.** Y.-C. Kim, W.-S. Bak and S.-K. Lee, "A Study on the Factors Affecting the Arson", *Fire Science and Engineering*, vol. 28, no. 2, pp. 69-75, 2014.

CrossRef  Google Scholar

**3.**J. Myers, C. G. de Souza e Silva, R. Doom, H. Fonda, K. Chan, S. Kamil-Rosenberg, et al., "Cardiorespiratory Fitness and Health Care Costs in Diabetes: The Veterans Exercise Testing Study", *American Journal of Medicine*, vol. 132, no. 9, pp. 1084-1090, 2019.

Google Scholar

**4.**Y. Nomura, Y. Ishii, Y. Chiba, S. Suzuki, A. Suzuki, S. Suzuki, et al., "Does last year's cost predict the present cost? An application of machine leaning for the japanese area-basis public health insurance database", *International Journal of Environmental Research and Public Health*, vol. 18, no. 2, pp. 1-11, 2021.

CrossRef  Google Scholar

**5.**A. S. Raut and P, Comparative Study of Regression Models and Deep Learning Models for Insurance Cost Prediction, Springer International Publishing, vol. 1, 2019.

Google Scholar

**6.**A. A. Kodiyan and K. Francis, *Linear regression model for predicting medical expenses based on insurance data*, December 2019.

Google Scholar

**7.** Nidhi Bhardwaj and Rishabh Anand, "Health Insurance Amount Prediction", *International Journal of Engineering Research And*, vol. 9, no. 05, pp. 1008-1011, 2020.

CrossRef  Google Scholar

**8.** M. Hanafy and O. M. A. Mahmoud, "Predict Health Insurance Cost by using Machine Learning and DNN Regression Models", *International Journal of Innovative Technology and Exploring Engineering*, vol. 10, no. 2, pp. 137-143, 2021.

CrossRef  Google Scholar

**9.** M. Hanafy and R. Ming, *Machine learning approaches for auto insurance big data. Risks*, vol. 9, no. 2, pp. 1-23, 2021.

Google Scholar

**10.** G. Fan, Z. Deng, X. Wu and Y. Wang, "Medical insurance and health equity in health service utilization among the middle-aged and older adults in China: A quantile regression approach", *BMC Health Services Research*, vol. 20, no. 1, pp. 1-12, 2020.

CrossRef  Google Scholar

**11.** A. Lakshmanarao, C. S. Koppireddy and G. V. Kumar, *Prediction of medical costs using regression algorithms*, vol. 10, no. 5, pp. 751-757, 2020.

Google Scholar

**12.** M. Mata-Cases, B. Rodriguez-Sanchez, D. Mauricio, J. Real, B. Vlacho, J. Franch-Nadal, et al., "The association between poor glycemic control and health care costs in people with diabetes: A population-based study", *Diabetes Care*, vol. 43, no. 4, pp. 751-758, 2020.

Google Scholar

**13.** M. AsadUllah, M. A. Khan, S. Abbas, A. Athar, S. S. Raza and G. Ahmad, "Blind channel and data estimation using fuzzy logic-empowered opposite learning-based mutant particle swarm optimization", *Computational intelligence and neuroscience*, 2018.

CrossRef  Google Scholar

**14.**F. Khan, M. A. Khan, S. Abbas, A. Athar, S. Y. Siddiqui, A. H. Khan, et al., "Cloud-based breast cancer prediction empowered with soft computing approaches", *Journal of healthcare engineering*, 2020.

CrossRef  Google Scholar

**15.**MA Khan, A Kanwal, S Abbas, F Khan and T Whangbo, "Intelligent Model for Predicting the Quality of Services Violation using Machine learning CMC-Computers", *Materials & Continua*, vol. 71, no. 2, pp. 3607-3619, 2022.

CrossRef  Google Scholar

## Data Sources

1. **Kaggle (2021).** *Medical Cost Personal Datasets*. Retrieved from
   https://www.kaggle.com/mirichoi0218/insurance
   - ○ The dataset used in this project, which contains information about individuals' demographic features and medical insurance costs.

## Online Tutorials and Documentation

1. **Pandas Documentation (2021).** *pandas: Powerful Python Data Analysis Toolkit*. Retrieved from https://pandas.pydata.org/
   - ○ Documentation for the pandas library, which was used for data manipulation and analysis.
2. **Matplotlib Documentation (2021).** *Matplotlib: Visualization with Python*. Retrieved from https://matplotlib.org/
   - ○ The official documentation for Matplotlib, another visualization library used for plotting data.

3.  **Numpy Documentation (2021).** *Numpy: The Fundamental Package for Scientific Computing with Python*. Retrieved from https://numpy.org/
    - Documentation for the NumPy library, used for numerical computations and array manipulations.

## Articles and Online Resources

1.  **Brownlee, J. (2016).** *How to Evaluate Machine Learning Algorithms*. Machine Learning Mastery. Retrieved from https://machinelearningmastery.com/evaluate-machine-learning-algorithms/
    - An online resource providing guidance on evaluating machine learning models, including metrics like R-squared.
2.  **Harrison, O. (2019).** *Machine Learning Basics with the K-Nearest Neighbors Algorithm*. Towards Data Science. Retrieved from https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors -algorithm-6a6e71d01761
    - An article providing an introduction to basic machine learning algorithms, which helped in understanding the broader context of model selection.