# Predective Analytics for Personalized Insurance Pricing

Vaishnavi Sanjeevkumar Kulkarni
Dept. of Industrial & Systems
Engineering
University at Buffalo – State University
of NewYork
Buffalo, NY, USA
vk52@buffalo.edu

*Abstract*— **This project applies machine learning techniques to predict personalized health insurance charges based on features like age, BMI, smoker status, and region. XGBoost outperformed Random Forest and Linear Regression in accuracy, achieving the lowest RMSE and highest R² score by capturing complex, non-linear relationships. Key predictors such as smoker status and BMI were identified, aligning with industry knowledge. While XGBoost and Random Forest provided superior performance, their lack of interpretability and computational complexity were limitations. Future work will focus on improving model transparency, incorporating real-time data, and addressing fairness to enhance personalized and equitable insurance pricing.**

*Keywords - Machine learning, XGBoost, Random Forest, Linear Regression, insurance pricing, personalized pricing, predictive modelling, feature importance, BMI, smoker status, model interpretability, RMSE, R² score, fairness, real-time data, health insurance, non-linear relationships.*

## I. INTRODUCTION

As industries become increasingly data-driven, predictive analytics becomes increasingly important, especially in the insurance industry. As the market becomes more competitive, health insurance companies must assess client risk accurately and customize premium pricing. Traditionally, pricing is based on broad demographic and health factors, which is inefficient and inequitable.

The project "Predictive Analytics for Personalized Insurance Pricing" aims to develop a machine learning model that predicts insurance charges with greater precision. In order to create an accurate and fair pricing model, we will use advanced machine learning techniques.

Using ensemble methods such as XGBoost, we aim to capture complex relationships between variables such as age, BMI, and lifestyle habits. The objective of this research is to improve the effectiveness of insurance pricing models and ensure equitable outcomes across diverse client demographics.

The implications of this work go beyond insurance; they can also be applied to other sectors such as healthcare, finance, and retail, where personalized risk assessment plays an important role. As this project promotes innovation in predictive analytics and ethical standards, it aims to improve decision-making across multiple domains.

## II. LITERATURE SURVEY

### A. "Advanced Artificial Inteligence techniques for Predictive Analytics in ife Insurance: Enhancing Risk Assesment and Pricing Accuracy"

In traditional life insurance pricing methodologies, risk categorization and static demographic factors are relying on wide categories of risk. Through the use of AI-powered predictive analytics, research shows that personalized pricing structures can be achieved by incorporating dynamic health data and behavioural patterns.

Researchers have demonstrated that machine learning algorithms, such as Random Forests and Gradient Boosting Machines, are capable of constructing accurate risk profiles from diverse data sources. By integrating real-time health data through wearable devices, premium adjustments can be made continuously.

It is important, however, to consider algorithmic bias and regulatory compliance in existing literature. For pricing decisions to be transparent and fair, robust governance frameworks and explainable AI techniques should be implemented. In order to maintain regulatory frameworks, industry stakeholders and regulatory bodies must continue to collaborate.

### B. "Advanced Artificial Inteligence Models for Predictive Analytics Insurance: Techniques, Applications and Real world Case Studies."

Insurers are transforming the way they price insurance through the integration of AI-powered predictive analytics. In today's dynamic risk landscape, traditional methods of insurance pricing relying on broad categories and historical data have proven inadequate. Data from telematics, wearables, and social media can be analyzed using AI models to provide more accurate risk assessments.

Machine learning algorithms, particularly decision trees, SVMs, and random forests, excel at identifying complex patterns and relationships within data. These technologies enable dynamic pricing models that adjust premiums based on real-time behaviour and risk factors. Implementation benefits include improved profitability through accurate risk assessment, enhanced customer experience through

personalized pricing, and optimized pricing structures that reflect individual risk profiles.

Case studies demonstrate successful implementations, with companies like Progressive Insurance utilizing telematics data for usage-based insurance programs. However, challenges remain regarding data privacy and the need for transparent, explainable AI models.

## III. CHALLENGES

While advancements in predictive analytics for health insurance are significant, several knowledge gaps and challenges persist:

### A. Model Interpretability and Explainability

Many complex machine learning models lack transparency, making it difficult for stakeholders to understand premium determinations. There is a need for standardized methods to enhance model interpretability without sacrificing accuracy.

### B. Data Quality and Availability

Predictive modelling depends on high-quality data. Insurers often face issues with data silos, missing information, and inconsistencies. More research is needed on data integration and handling missing data effectively.

### C. Ethical Considerations and Fairness

The risk of bias in predictive models can reinforce existing inequalities. While interest in fairness-aware machine learning is growing, best practices for implementing fairness constraints are still underexplored.

### D. Dynamic Risk Assessment:

Health risks evolve rapidly, yet many models do not account for these changes. Research is needed to develop adaptive techniques that can incorporate new information and update pricing dynamically.

### E. Integration of Diverse Data Sources::

Combining varied data sources (e.g., health records, lifestyle data) presents technical and ethical challenges. Comprehensive frameworks for integrating heterogeneous data types are necessary.

### F. Regulatory and Compliance Issues::

Navigating varying regulations while implementing advanced analytics can be challenging. More investigation is required into the impact of regulatory compliance on model development and deployment.

## IV. METHODOLOGY

To develop the predictive model for personalized insurance pricing, a systematic approach consisting of the following steps.

2. *Data Preprocessing*
   2.1. *Dataset Overview*
      - The dataset consists of 1,338 rows and 7 columns, including features such as age, BMI, smoking status, and insurance charges (the target variable).
      - Features were categorized as:
        - Numeric features: age, BMI, children, and charges.

- Categorical Features: sex, smoker and region.

2.2. *Feature Encoding*
   - Categorical Features were encoded using OneHotEncoder:
     - Sex: Converted into a binary variable (e.g. 0 for male, 1 for female).
     - Smoker: Converted into binary variable (e.g. 0 for non-smoker and 1 for smoker).
     - Region: One-hot encoded into four binary variables (southwest, southeast, northwest and northeast).
   - Numeric features were standardized using StandardScaler to ensure uniform scales.

2.3. *Handling Correlations*
   - Correlation analysis was conducted to examine relationships among numeric feature. It was observed that:
     - Smoker had the strongest positive correlation with charges.
     - BMI and age also had moderate correlations with charges.

2.4. *Data Splitting*
   - The dataset was split into training (80%) and testing (20%) sets to ensure that the models were evaluated on unseen data.
   - Stratified splitting was applied where necessary to ensure class balance for binary features such as smoker.

3. *Model Development*
   Three models were developed and evaluated:
   3.1. *Linear Regression:*
      - A simple interpretable model that assumes a linear relationship between the predictors and the target.
      - Assumptions such as homoscedasticity and normality of residuals were validated through diagnostic plots.
   3.2. *Random Forest:*
      - An ensemble learning method that combines multiple decision trees.
      - Hyperparameters such as n_estimators (number of trees), max_depth, and min_samples_split were tuned using GridSearchCV.
   3.3. *XGBoost:*
      - A gradient boosting framework that builds trees iteratively to minimize the loss function.
      - Regularization terms were applied to prevent overfitting. Parameters such as learning_rate, max_depth, and n_estimators were optimized using cross validation.

4. *Model Evaluation*
   4.1. *Performance Metrics:*
      - Root Mean Squared Error (RMSE): Measures average prediction error, lower is better.
      - $R^2$ Score: Represents the proportion of variance explained by the model, higher is better.
   4.2. *Validation Methods:*

- Cross-Validation: Used during hyperparameter tuning to ensure model robustness.
- Residual Analysis: Performed to validate key assumptions for linear regression.
- Feature Importance: Extracted from Ensemble model (Random Forest and XGBoost) to identify the most significant predictors.

## V. ASSUMPTIONS OF THE MODELS

Each machine learning model operates under specific assumptions that ensure the validity and reliability of its results. Below are the key assumptions for the model used:

*1. Linear Regression*
  i.   Linearity: Assumes a linear relationship between predictors ($X$) and the target ($y$):
$$y = X\beta + \varepsilon$$
  ii.  Independent Errors (i.i.d.): Residuals ($\varepsilon$) are independent and identically distributed.
  iii. Homoscedasticity: Residual variance remains constant across all levels of predictors.
  iv.  Normality of Errors: Residuals follow a normal distribution.
  v.   No Multicollinearity: Predictors are not highly correlated.

*2. Random Forest*
  i.   Independence Between Trees: Each tree is trained on a unique bootstrap sample, ensuring diversity.
  ii.  No Assumptions About Input Data: Can handle non-linear relationships and does not require specific data distributions.
  iii. Sufficient Data: Requires enough samples for robust decision tree construction.

*3. XGBoost:*
  i.   Additive Model: Builds predictions iteratively but adding weak learners:
$$y = \sum_{i=1}^{n} f_i(x)$$
  ii.  Independent Observations: Assumes data points are independent.
  iii. No Distribution Assumption: Input data does not need to follow a specific distribution.
  iv.  Feature Relevance: Requires meaningful features for accurate predictions.

## VI. NOTATIONS AND THEIR DIMENSIONALITIES

This section formalizes the mathematical notations used in the project.

*1. Input Variables:*
  i.   $x_i$: Feature vector for the $i$-th individual,
$$x_i \in \mathbb{R}^d$$
*where* $d = 7$(e.g., age, BMI, smoker etc.).
  ii.  $X$: Feature matrix for all $n = 1338$ samples,
$X \in \mathbb{R}^{n \times d}$.
*2. Output Variables:*
  i.   $y_i$: Target variable (insurance charges) for the $i$-th individual, $y_i \in \mathbb{R}$.
  ii.  $y$: Vector of insurance charges for all samples,
$y \in \mathbb{R}^n$
*3. Model Parameters:*

  i.   Linear Regression:
- $\beta \in \mathbb{R}^d$: Coefficients of the linear model.
- $\varepsilon \in \mathbb{R}^n$: Residual error.
  ii.  Random Forest and XGBoost:
- $f(x)$: Aggregated prediction function,
$$f(x) \in \mathbb{R}$$
- $\theta_k$: Parameters of the $k$-th tree.

*1. Optimization Goals:*
  i.   Linear regression: Minimize residual sum of squares:
$$\min_{\beta} \|y - X\beta\|^2$$
  ii.  Random Forest: Minimize impurity at each split.
  iii. XGBoost: Minimize a custom loss function with regularization:
$$L(\theta) = \sum_{i=1}^{n} l(y_i, \hat{y}_i) + \Omega(\theta)$$

## VII. MODEL FORMULATION AND ESTIMATORS

*1. Linear Regression*
**Formulation:**
Predicts the target $y_i$ as a linear combination of the input feature $x_i$:
$$y = X\beta + \varepsilon$$
Where,
$y \in \mathbb{R}^n, X \in \mathbb{R}^{n \times d}, \beta \in \mathbb{R}^d, \varepsilon \in \mathbb{R}^n$.
**Estimator:**
Minimizes residual sum of square:
$$\hat{\beta} = (X^T X)^{-1} X^T y$$

*2. Random Forest*
**Formulation:**
An ensemble of $T$ decision trees, predicting as:
$$\hat{y}_i = \frac{1}{T} \sum_{t=1}^{T} f_t(x_i)$$
Where $f_t(x_i)$ is the prediction from tree t.
**Estimator:**
Uses recursive portioning to minimize impurity Final prediction is the average of all tree outputs.

*3. XGBoost*
**Formulation:**
A gradient-boosted ensemble that predicts as:
$$\hat{y}_i = \sum_{t=1}^{T} f_t(x_i)$$
Minimizes an objective function:
$$L(\theta) = \sum_{i=1}^{n} l(y_i, \hat{y}_i) + \Omega(\theta)$$
Where $l(y_i, \hat{y}_i)$ is the loss function and $\Omega(\theta)$ is the regularization term.
**Estimator:**
Optimizes $L(\theta)$ iteratively using gradient descent.

## VIII. ALGORITHM FOR OPTIMIZATION

*1. Linear Regression:*
**Goal:** Minimize residual sum of squares:
$$\hat{\beta} = arg \min_{\beta} \|y - X\beta\|^2$$
**Optimization:** Ordinary Least Squares (OLS):

$$\hat{\beta} = (X^TX)^{-1}X^Ty$$

2. *Random Forest:*
   **Goal:** Build multiple decision trees and minimize impurity at each split.
   **Optimization:** Recursive portioning
   - Bootstrap sampling creates diverse datasets for each tree.
   - Evaluate feature splits by impurity reduction
   - Select the best split recursively until stopping criteria
   - Final prediction is average of tree outputs.

3. *XGBoost:*
   **Goal:** Iteratively build trees to minimize an objective function:
   $$L(\theta) = \sum_{i=1}^{n} l(y_i, \hat{y}_i) + \Omega(\theta)$$

   **Optimization:**
   - Compute gradients of the loss function.
   - Add a tree to correct errors from previous predictions.
   - Regularize with $\Omega(\theta)$ to control model complexity.
   - Update predictions
   $$\hat{y}_i^{new} = \hat{y}_i + \eta \cdot g_i$$

## IX. DETAILED DATASET DESCRIPTION

The dataset used in this project is a health insurance dataset that contains various attributes related to individuals and their corresponding insurance charges. The objective of this project is to use this dataset to develop a predictive model for estimating personalized insurance pricing.

1. *Dataset Overview*
   The dataset consists of 1,388 observations and 7 features. These features include both numeric and categorical variables as well as the target variables (insurance charges).
   **Features:**
   - Age: The age of the individual (Numeric)
   - Sex: The gender of the individual (Categorical)
   - BMI: Body mass Index (BMI) of the individua (Numeric)
   - Children: The number of children/dependents covered by the insurance (Numeric)
   - Smoker: Whether the individual smokes (Categorical)
   - Region: Geographical region of the individual (Categorical)
   - Charges: The target variable representing the insurance charges for the individual (Numeric)

2. *Data Types and Structure*
   - Numeric Variables: age, BMI, children, and charges.
   - Categorical Variables: sex, smoker, and region.
   The dataset is balanced in terms of the number of entries across the categorical variables, but there is a skewed distribution for the target variable charges, as it exhibits a right-skew, with a few individuals having very high charges.

3. *Target Variable: charges*

The target variable charges represent the insurance charges billed to each individual. The distribution of insurance charges is skewed, meaning most individuals have relatively low charges, while a few have much higher charges.

4. *Correlation Analysis:*
   Since the dataset contains both numeric and categorical variables, the correlation analysis is performed on numeric variables only. We compute the correlation matrix for age, BMI, children and charges.
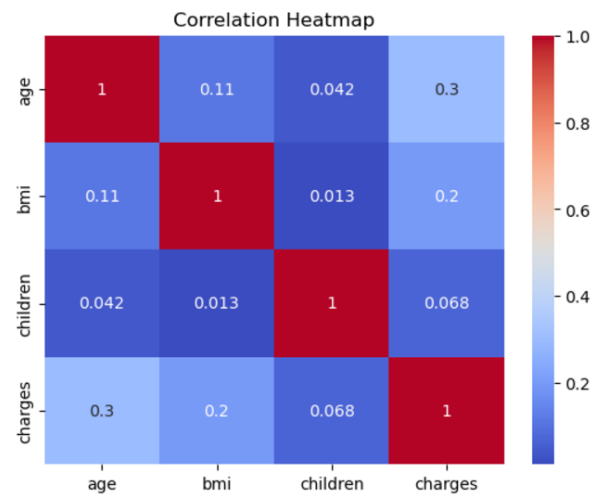   **Code Snippet:**

```python
# Load dataset
data = pd.read_csv("C:/Users/Vaishnavi Kulkarni/Documents/IE 500 Statistical Machine Learning/Health_insurance.csv")

# Overview
print(data.info())
print(data.describe())
print(data.head())

# Filter numeric columns for correlation analysis
numeric_data = data.select_dtypes(includes[np.number])

# Correlation heatmap
sns.heatmap(numeric_data.corr(), annot=True, cmap="coolwarm")
plt.title("Correlation Heatmap")
plt.show()
```

**Correlation Heatmap:**



**Insights:**
- Age and Charges: There is a moderate positive correlation between ang and insurance charges. As people age, their insurance premiums tend to increase.
- BMI and Charges: A strong positive correlation exists between BMI and charges, suggesting that individuals with higher BMI tend to have higher insurance charges.
- Children and Charges: A slight positive correlation exists between the number of children and charges. Individuals with more children tend to have a slightly higher premiums, likely due to increased dependents.

5. *Data Preprocessing*
   **Encoding Categorical Features:**
   - Sex: The sex variable is encoded as a binary feature (0 for male, 1 for female).
   - Smoker: The smoker variable is encoded as a binary (0 for non-smoker and 1 for smoker)
   - Region: The region variable is encoded using one-hot encoding, creating binary features for each region (e.g., southwest, southeast, northwest, northeast)
   **Scaling:**

The numeric features (age, BMI, and children) are scaled using StandardScaler to ensure they are on same scale for model training.

6. *Data Splitting*

The dataset is split into two parts:
- Training Set: 80% of the data used for training the models.
- Testing Set: 20% of the data used to evaluate the models.

## X. TRAINING-VALIDATION-TESTING PARTITION

In machine learning, dividing the dataset into training, validation, and testing sets is essential for model evaluation and generalization. Below is the summary of the dataset portioning for this project.

1. *Purpose of Each Partition*
- Training Set: Used to train the model and optimize its parameter.
- Validation Set: Used for tuning hyperparameters and model selection.
- Testing Set: Used to evaluate the model's performance on unseen data.

2. *Dataset Partitioning Strategy*

For this project, the data is partitioned into the following proportions:
- Training Set: 80% of the data is used for training the model.
- Testing Set: 20% of the data is used to test the model.

No separate validation set is used in this case because we will use the cross-validation technique and hyperparameter tuning during model training.

3. *Code Snippet for Data Splitting*

```
# Split the data into features (X) and target (y)
X = data.drop(columns=['charges'])  # Features (excluding 'charges')
y = data['charges']  # Target (insurance charges)

# Define the preprocessing for numeric and categorical data
numeric_features = ['age', 'bmi', 'children']
categorical_features = ['sex', 'smoker', 'region']

# Apply transformations to each feature
preprocessor = ColumnTransformer(
    transformers=[
        ('num', StandardScaler(), numeric_features),
        ('cat', OneHotEncoder(drop='first'), categorical_features)
    ])

# Create the pipeline
pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('model', RandomForestRegressor(random_state=42))
])

# Train-Test Split (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Output dataset shapes
print(f"Training set size: {X_train.shape[0]} samples")
print(f"Testing set size: {X_test.shape[0]} samples")

Training set size: 1070 samples
Testing set size: 268 samples
```

4. *Code Snippet for Validation through Cross-Validation*

```
# Perform cross-validation
cv_scores = cross_val_score(pipeline, X, y, cv=5, scoring='neg_mean_squared_error')

# Print RMSE scores from cross-validation
print(f"Cross-validation RMSE scores: {-cv_scores}")
print(f"Average RMSE across folds: {-cv_scores.mean()}")

Cross-validation RMSE scores: [22446373.51382653 28945532.90278041 19056449.32502298 25496557.79877909
22535802.45603657]
Average RMSE across folds: 23696143.199289113
```

## XI. MODELING PERFORMANCE AND LIMITATIONS

1. *Performance Summary*

| Model | RMSE | $R^2$ Score | Comments |
|---|---|---|---|
| Linear Regression | 5796.28 | 0.78 | Fast, simple, interpretable but may underperform on complex data |

Predictive Analytics for Personalized Insurance Pricing

| Model | RMSE | $R^2$ Score | Comments |
|---|---|---|---|
| Random forest | 4567.96 | 0.87 | Good for capturing non-linear relationships, but lacks interpretability |
| XGBoost | 4517.29 | 0.87 | High performance, handles missing data well, but complex and computationally intensive. |

2. *Performance Comparison*
- Code Snippet:

```
models = ['Linear Regression', 'Random Forest', 'XGBoost']
rmse_scores = [5796.28, 4567.96, 4517.29]
r2_scores = [0.78, 0.87, 0.87]

# Create the figure and axis
fig, ax1 = plt.subplots(figsize=(10, 6))

# Bar plot for RMSE
ax1.bar(models, rmse_scores, color='skyblue', label='RMSE', alpha=0.7)
ax1.set_ylabel('RMSE', color='blue')
ax1.tick_params(axis='y', labelcolor='blue')

# Add a secondary y-axis for R²
ax2 = ax1.twinx()
ax2.plot(models, r2_scores, color='red', marker='o', label='R² Score', linewidth=2)
ax2.set_ylabel('R² Score', color='red')
ax2.tick_params(axis='y', labelcolor='red')

# Add titles and legend
plt.title("Model Performance Comparison")
fig.tight_layout()
plt.show()
```
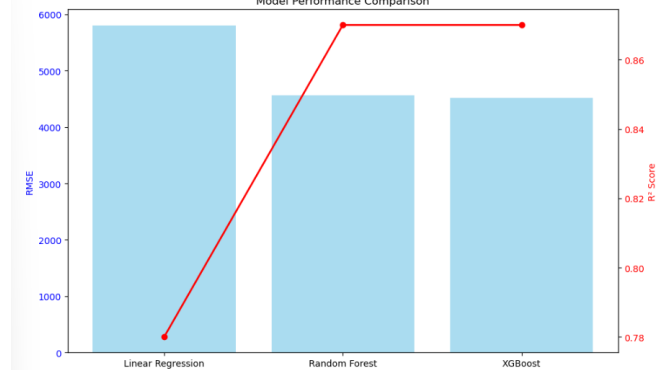
- Model Performance Comparison Graph:



3. *Limitations of the Models*

**Linear Regression:**
- Assumes linear relationships may not capture more complex patterns.
- Sensitive to outliers.
- Assumes homoscedasticity and normality of residuals, which might not always hold.

**Random Forest:**
- Lack of interpretability - difficult to explain individual predictions.
- Can be computationally expensive, especially with large datasets.
- Tuning the hyperparameters is necessary for optimal performance.

**XGBoost:**
- Requires careful tuning of hyperparameters.

- Computationally expensive, especially for large datasets.
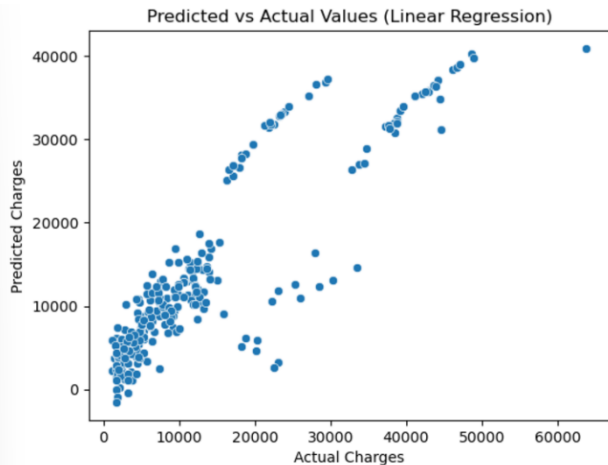- Less interpretable compared to simpler models like Linear Regression.

### XII. VALIDATION OF ASSUMPTIONS

1. *Linear Regression Assumptions*
   - Linearity: The relationship between features and target is expected to be approximately linear.
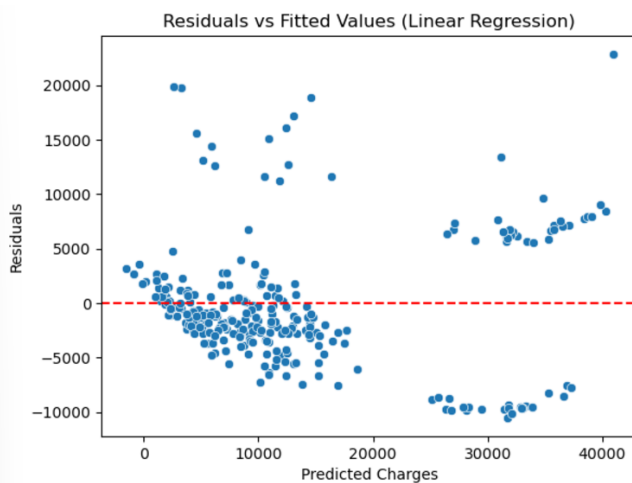     **Code Snippet and Visualization:**

```
sns.scatterplot(x=y_test, y=y_pred_lr)
plt.title('Predicted vs Actual Values (Linear Regression)')
plt.xlabel('Actual Charges')
plt.ylabel('Predicted Charges')
plt.show()
```


Predicted vs Actual Values (Linear Regression)

   - Independence of Errors (i.i.d.): Independence of residuals can be validated by plotting residuals vs fitted values, The residuals should appear randomly scattered without any discernable pattern.
     **Code Snippet and Visualization:**

```
residuals_lr = y_test - y_pred_lr
sns.scatterplot(x=y_pred_lr, y=residuals_lr)
plt.axhline(0, color='red', linestyle='--')
plt.title('Residuals vs Fitted Values (Linear Regression)')
plt.xlabel('Predicted Charges')
plt.ylabel('Residuals')
plt.show()
```
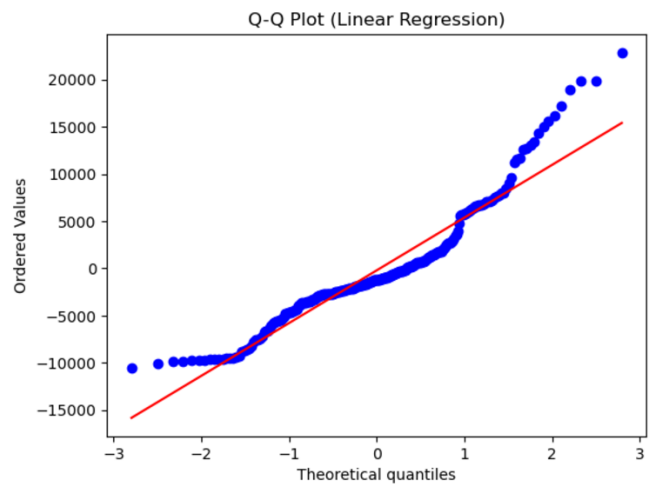

Residuals vs Fitted Values (Linear Regression)

   - Normality of Errors: Q-Q plot is used to check if the residuals are normally distributed. If the points roughly lie along the line, then the assumption holds.
     **Code Snippet and Visualization**

```
import scipy.stats as stats
stats.probplot(residuals_lr, dist="norm", plot=plt)
plt.title('Q-Q Plot (Linear Regression)')
plt.show()
```
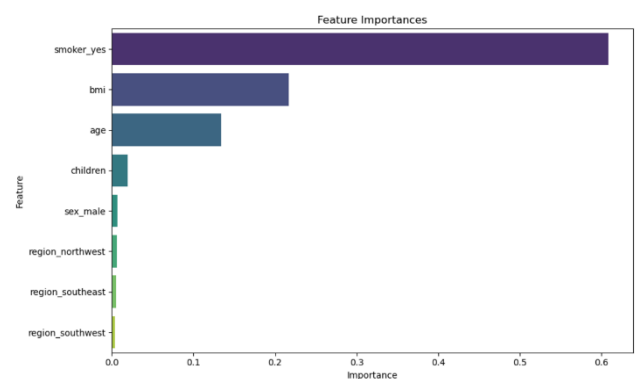

Q-Q Plot (Linear Regression)

2. *Random Forest Assumptions*
   - Feature Importance:
     **Code Snippet and Visualization:**

```
rf_pipeline.fit(X_train, y_train)
importances = rf_pipeline.named_steps['regressor'].feature_importances_
feature_names = numerical_features + list(preprocessor.named_transformers_['cat'].get_feature_names_out())
# Create a DataFrame for better handling
importances_df = pd.DataFrame({
    'Feature': feature_names,
    'Importance': importances
}).sort_values(by='Importance', ascending=False)

# Plot the feature importances
plt.figure(figsize=(10, 6))
sns.barplot(
    x='Importance',
    y='Feature',
    data=importances_df,
    palette='viridis'
)
plt.title('Feature Importances')
plt.xlabel('Importance')
plt.ylabel('Feature')
plt.tight_layout()
plt.show()
```


Feature Importances

3. *XGBoost Assumptions*
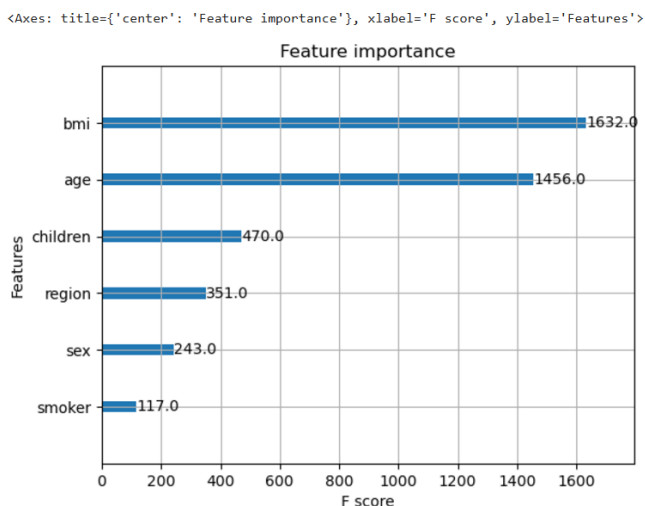   - Feature Importance:
   **Code Snippet and Visualization:**

```python
from sklearn.preprocessing import LabelEncoder
from xgboost import XGBRegressor, plot_importance

# Encode categorical columns
label_encoders = {}
for col in ['sex', 'smoker', 'region']:
    label_encoders[col] = LabelEncoder()
    X_train[col] = label_encoders[col].fit_transform(X_train[col])
    X_test[col] = label_encoders[col].transform(X_test[col])

# Train the XGBRegressor model
xgb_model = XGBRegressor(
    objective='reg:squarederror',
    n_estimators=100,
    learning_rate=0.1,
    max_depth=6,
    random_state=42
)
xgb_model.fit(X_train, y_train)

# Plot feature importance
plot_importance(xgb_model, importance_type='weight', max_num_features=10)
```

```
<Axes: title={'center': 'Feature importance'}, xlabel='F score', ylabel='Features'>
```



Feature importance chart — bmi: 1632.0, age: 1456.0, children: 470.0, region: 351.0, sex: 243.0, smoker: 117.0

## XIII. MODELING OUTCOMES

Let's discuss the outcomes, significance and limitations of three models.

1. *Linear Regression*
   1.1. *Outcome:*
      - **RMSE:** The highest among all models, reflecting poor predictive accuracy.
      - $R^2$ **Score:** Moderate, indicating limited explanatory power for the variance in charges.
   1.2. *Justification:*
      Linear Regression assumes a linear relationship between features and the target variable. While this assumption makes the model highly interpretable, it is overly simplistic for this dataset. On-linear relationships, such as exponential increase in charges for smokers or high-BMI individuals are not captured effectively. Additionally, interactions between variables such as age and smoking status are ignored.
   1.3. *Significance:*

- Linear Regression clearly shows the direct contribution of each feature which is valuable in contexts like regulatory compliance or policy explanations.
- The model is computationally efficient and easy to implement.
   1.4. *Limitations:*
      - The model's simplicity leads to underfitting, especially when dealing with nonlinear relationships.
      - Linear Regression is sensitive to outliers, such as individuals with extremely high charges, which skew predictions.

2. *Random Forest*
   2.1. *Outcome:*
      - **RMSE:** Significantly lower than Linear Regression, reflecting improved accuracy in capturing complex relationships.
      - $R^2$ **Score:** Higher, indicating that more variance in insurance charges is explained.
   2.2 *Justification:*
      Random Forest, an ensemble learning method, builds multiple decision trees and aggregate their predictions. This allows the model to capture non-linear relationships and interactions between variables, such as the combined impact of smoking and BMI on charges. The feature importance analysis highlighted that smoker status and BMI were key predictors, aligning with domain knowledge.
      The model's robustness comes from random sampling of both features and data points, reducing overfitting while maintaining flexibility to adapt to data complexities.
   2.3 *Significance:*
      - Random Forest captures non-linear patterns effectively, improving predictive accuracy.
      - Offers insights into which features drive predictions.
   2.4 *Limitations:*
      - The model is "black box", making it challenging to explain individual predictions without additional tools like SHAP.
      - Training a Random Forest can be resource-intensive, especially with a large number of trees.

3 *XGBoost*
   3.1 *Outcomes:*
   - **RMSE:** The lowest among all the models, demonstrating superior predictive performance.
   - $R^2$ **Score:** The highest, explaining the largest proportion of variance in insurance charges.
   3.2 *Justification:*
      XGBoost uses gradient boosting which iteratively refines predictions by focusing on correcting the largest errors from previous iterations. This process allows XGBoost to handle both non-linear relationships and iterations more effectively than Random Forest. Additionally, regularization prevents overfitting, making the model both accurate and robust.

Feature importance analysis in XGBoost confirmed the dominance of smoker status, BMI and age as key drivers of insurance charges. The ability to tune hyperparameters such as learning rate and tree depth further optimizes its performance.

### 3.3 Significance:

- XGBoost provides the best accuracy and reliability, making it the most suitable model for this task.
- Effective for handling complex datasets with non-linear interactions and noise.

### 3.4 Limitations:

- XGBoost requires careful tuning of multiple hyperparameters, which can be time-consuming.
- Despite feature importance analysis, explaining individual predictions is challenging without additional tools.

## XIV. CONCLUSION

This project demonstrated the effectiveness of machine learning models in predicting personalized insurance charges, with XGBoost emerging as the most accurate, achieving the lowest RMSE and highest $R^2$ score. Its ability to capture non-linear relationships and interactions, particularly between critical features like smoker status, BMI, and age, set it apart from other models. Random Forest also performed well, effectively modelling complex patterns, but it lacked the precision of XGBoost. In contrast, Linear Regression, while simple and interpretable, struggled with the dataset's non-linear relationships, leading to lower accuracy. The analysis confirmed smoker status and BMI as the most influential predictors of insurance charges, aligning with industry expectations. However, ensemble models like XGBoost and Random Forest presented challenges in interpretability and required significant computational resources. Future improvements could focus on incorporating more diverse features, such as real-time health data, and addressing fairness concerns to ensure equitable and transparent insurance pricing systems.

## XV. FUTURE SCOPE

1. *Enhancing Model Interpretability*
   While ensemble models like XGBoost and Random Forest deliver high accuracy, their lack of transparency makes them challenging to use in contexts requiring regulatory compliance or user trust. Future work will focus on integrating techniques such as SHAP to provide clear insights into the decision-making process. Additionally, simpler rule-based models or hybrid approaches could be explored to balance accuracy and interpretability.

2. *Expanding Feature Diversity*
   The dataset could be enriched with more granular features to improve predictive power. For example, incorporating real-time health data from wearable devices (e.g., heart rate, sleep patterns) and lifestyle information (e.g., exercise frequency) could enable dynamic risk assessments. Adding socio-economic factors, such as income or education, would also help capture broader determinants of insurance costs, making the models more robust and fairer.

3. *Real-Time Data Integration*
   Leveraging real-time data could revolutionize insurance pricing by enabling dynamic and personalized adjustments based on changes in individual health metrics or behaviours. Future models could integrate APIs for wearable devices and health monitoring apps, allowing for continuous updates to pricing models. This would create a more personalized and fair insurance pricing system, where charges reflect current health conditions and lifestyle choices.

4. *Addressing Bias and Ensuring Fairness*
   Future efforts should prioritize implementing fairness-aware machine learning algorithms to mitigate potential biases in pricing models. Techniques like adversarial debiasing or fair representation learning can ensure that predictions are not influenced by sensitive attributes like gender or ethnicity. Addressing these biases is crucial for maintaining fairness and regulatory compliance in automated pricing systems.

5. *Scaling and Cross-Domain Applications*
   Scalability is another key focus for future work. The current models could be optimized for large-scale datasets, making them applicable to broader insurance markets. Moreover, the methodologies developed in this project can extend beyond insurance to other domains such as healthcare (e.g., predicting treatment costs) and finance (e.g., credit scoring). Cross-domain applications could further benefit from transfer learning techniques to adapt models across different industries efficiently.

## REFERENCES

[1] B. P. Kasaraneni, "Advanced Artificial Intelligence Techniques for Predictive Analytics in Life Insurance: Enhancing Risk Assessment and Pricing Accuracy," *Distributed Learning and Broad Applications in Scientific Research*, vol. 5, pp. 547-588, 2019.
In **IEEE format**, the reference for your paper would be:
[2] K. K. Kondapaka, "Advanced Artificial Intelligence Models for Predictive Analytics in Insurance: Techniques, Applications, and Real-World Case Studies," *Australian Journal of Machine Learning Research & Applications*, vol. 1, no. 1, pp. 244-290, 2021.