# Health Insurance Premium Prediction Using Machine Learning Models: A Comparative Study of Implementations

Pruthviraj Santo

Student ID: SCFU123010

B.Tech (CSE) - V Sem

MIT Vishwaprayag University

December 5, 2025

# Contents

## Abstract

This report presents a comprehensive project on predicting health insurance premiums using machine learning regression models implemented from scratch. The dataset used is the Insurance dataset, featuring attributes such as age, sex, BMI, children, smoker status, and region. Three models are developed: Linear Regression using the normal equation, Gradient Boosting Regressor with decision trees, and Random Forest Regressor with bagging and feature randomness. The models are trained on the full dataset and evaluated using manual train-test splits. Performance metrics like Mean Absolute Error (MAE) and $R^2$ score are computed. The implementation includes robust input handling for interactive predictions. All code is hosted on GitHub at `https://github.com/MIDSTAN/Health-Insurance-Premium-Prediction`. This project demonstrates the feasibility of scratch implementations for educational purposes, achieving competitive performance comparable to scikit-learn counterparts.

Keywords: Machine Learning, Regression, Insurance Premium Prediction, Scratch Implementation, Gradient Boosting, Random Forest

# Chapter 1

# Introduction

## 1.1 Background

Health insurance premium prediction is a critical application in the insurance industry, enabling accurate pricing based on individual risk factors. Traditional actuarial methods rely on statistical models, but machine learning offers enhanced predictive power through data-driven approaches. This project focuses on implementing three regression models from scratch to predict medical expenses (premiums) using a publicly available dataset.

The Insurance dataset, sourced from Kaggle, contains 1338 records with features: age (numerical), sex (binary: 0-female, 1-male), BMI (numerical), children (numerical), smoker (binary: 0-no, 1-yes), and region (categorical: 1-4). The target is 'expenses' (continuous).

## 1.2 Objectives

- Implement Linear Regression, Gradient Boosting, and Random Forest regressors from scratch.

- Train models on the dataset and evaluate using MAE and $R^2$.

- Develop an interactive testing script for user predictions.

- Compare model performances and discuss insights.

- Host the project on GitHub with a comprehensive README.

## 1.3 Scope and Limitations

The implementations are simplified for educational purposes. No hyperparameter tuning or cross-validation is performed. Evaluation uses a manual 80/20 split. Future work

could include regularization and ensemble improvements.

## 1.4   Project Structure

The project includes training scripts for each model, a unified testing script, and utilities. All files are in `/home/midstan/Documents/Health Insurance Premium/Model/Training/` and `Testing/`.

# Chapter 2

# Literature Review

## 2.1 Overview of Regression Models

Linear Regression [?] assumes a linear relationship between features and target, solved via ordinary least squares. It serves as a baseline.

Gradient Boosting [?] builds an ensemble of weak learners (decision trees) sequentially, fitting to residuals. It excels in handling non-linearities.

Random Forest [?] uses bagging with random feature subsets, reducing variance through averaging multiple decision trees.

## 2.2 Prior Work on Insurance Prediction

Studies like [?] use ensemble methods on similar datasets, reporting $R^2$ scores up to 0.87. Scratch implementations are rare but valuable for understanding internals [?].

## 2.3 Gap Analysis

Most works use scikit-learn; this project emphasizes from-scratch coding to demystify algorithms.

# Chapter 3

# Methodology

## 3.1 Dataset Description

The dataset is loaded via pandas:

```python
df = pd.read_csv(data_path)
X = df[['age', 'sex', 'bmi', 'children', 'smoker', 'region']].
    values
y = df['expenses'].values
```

Features are pre-encoded (numerical). No missing values or scaling applied.

## 3.2 Model Implementations

### 3.2.1 Linear Regression from Scratch

Uses the normal equation:

$$\theta = (X^T X)^{-1} X^T y \tag{3.1}$$

where $X$ includes a bias column. Prediction: $\hat{y} = X\theta$.

Code snippet:

```python
def fit(self, X, y):
    # Add bias term (intercept column of ones)
    X_b = np.c_[np.ones((X.shape[0], 1)), X]
    # Normal equation: theta = (X^T X)^-1 X^T y
    theta_best = np.linalg.inv(X_b.T.dot(X_b)).dot(X_b.T).dot(y)
    self.intercept_ = theta_best[0]
    self.coef_ = theta_best[1:]
    return self

def predict(self, X):
```

```
11      # Add bias term
12      X_b = np.c_[np.ones((X.shape[0], 1)), X]
13      # Predictions: X_b * theta
14      return X_b.dot(np.r_[self.intercept_, self.coef_])
```

### 3.2.2  Gradient Boosting Regressor

Initial prediction: mean of $y$. Iteratively fits trees to residuals, updating predictions with learning rate $\eta$:

$$F_m(x) = F_{m-1}(x) + \eta h_m(x) \tag{3.2}$$

Early stopping if MAE ¡ tol.

Relies on scikit-learn's DecisionTreeRegressor for weak learners (hybrid approach).

Code snippet:

```
1  def fit(self, X, y):
2      # Initialize predictions with the mean (base model)
3      self.initial_prediction = np.mean(y)
4      predictions = np.full_like(y, self.initial_prediction, dtype=
           np.float64)
5      # Compute initial residuals
6      residuals = y - predictions
7      # Build trees iteratively
8      for i in range(self.n_estimators):
9          # Fit a decision tree to the current residuals (negative
               gradient for regression)
10         tree = DecisionTreeRegressor(max_depth=self.max_depth,
               random_state=42)
11         tree.fit(X, residuals)
12         # Predict residuals with the tree
13         tree_pred = tree.predict(X)
14         # Update predictions: add learning_rate * tree prediction
15         predictions += self.learning_rate * tree_pred
16         # Update residuals
17         residuals = y - predictions
18         # Store the tree
19         self.trees.append(tree)
20         # Early stopping: check if mean absolute residual is
               below tolerance
21         mean_abs_residual = np.mean(np.abs(residuals))
22         if mean_abs_residual < self.tol:
```

```
23                    print(f"Early␣stopping␣at␣iteration␣{i+1}:␣mean␣abs␣
                          residual␣=␣{mean_abs_residual:.2f}")
24                    break
25          print(f"Training␣completed.␣Final␣mean␣abs␣residual:␣{np.mean
                (np.abs(residuals)):.2f}")
26          print(f"Number␣of␣trees␣built:␣{len(self.trees)}")
27          return self
28
29   def predict(self, X):
30          # Start with initial prediction
31          predictions = np.full(X.shape[0], self.initial_prediction,
                dtype=np.float64)
32          # Add contributions from all trees
33          for tree in self.trees:
34              predictions += self.learning_rate * tree.predict(X)
35          return predictions
```

### 3.2.3   Random Forest Regressor from Scratch

Includes custom DecisionTreeRegressor using MSE for splits:

$$MSE = \frac{1}{n} \sum (y_i - \bar{y})^2 \tag{3.3}$$

Best split minimizes weighted child MSE. Bagging via bootstrap samples; random features $(\sqrt{n_{feat}})$.

Code snippet (Decision Tree fit):

```
1    def mse(y: np.ndarray) -> float:
2        """
3    ␣␣␣␣Mean␣Squared␣Error␣for␣a␣set␣of␣targets.
4    ␣␣␣␣"""
5        return np.mean((y - np.mean(y)) ** 2)
6
7    def best_split(X: np.ndarray, y: np.ndarray, feature_idxs: List[
         int]) -> Tuple[int, float, float, np.ndarray, np.ndarray]:
8        """
9    ␣␣␣␣Find␣the␣best␣split:␣feature␣and␣threshold␣minimizing␣MSE.
10   ␣␣␣␣"""
11       best_idx, best_thresh, best_mse, best_left, best_right = None
             , None, float('inf'), None, None
12       current_mse = mse(y)
13
```

```python
14      for idx in feature_idxs:
15          thresholds = np.unique(X[:, idx])
16          for thresh in thresholds:
17              left_mask = X[:, idx] <= thresh
18              right_mask = ~left_mask
19              if np.sum(left_mask) == 0 or np.sum(right_mask) == 0:
20                  continue
21              left_y, right_y = y[left_mask], y[right_mask]
22              weighted_mse = (len(left_y) * mse(left_y) + len(
                    right_y) * mse(right_y)) / len(y)
23              if weighted_mse < best_mse:
24                  best_idx = idx
25                  best_thresh = thresh
26                  best_mse = weighted_mse
27                  best_left = left_y
28                  best_right = right_y
29
30      return best_idx, best_thresh, best_mse, best_left, best_right
31
32  def build_tree(X: np.ndarray, y: np.ndarray, feature_idxs: List[
        int], max_depth: int = None, min_samples_split: int = 2, depth
        : int = 0) -> Node:
33      """
34      Recursively build the decision tree.
35      """
36      if len(y) < min_samples_split or (max_depth is not None and
            depth >= max_depth):
37          return Node(value=np.mean(y))
38
39      idx, thresh, _, _, _ = best_split(X, y, feature_idxs)
40      if idx is None:
41          return Node(value=np.mean(y))
42
43      left_mask = X[:, idx] <= thresh
44      right_mask = ~left_mask
45
46      left = build_tree(X[left_mask], y[left_mask], feature_idxs,
            max_depth, min_samples_split, depth + 1)
47      right = build_tree(X[right_mask], y[right_mask], feature_idxs
            , max_depth, min_samples_split, depth + 1)
48
```

```
49        return Node(idx, thresh, left, right)
```

For forest:

```python
def fit(self, X: np.ndarray, y: np.ndarray):
    np.random.seed(self.random_state)
    n_features = X.shape[1]
    for _ in range(self.n_estimators):
        X_boot, y_boot = self._bootstrap_sample(X, y)
        feature_idxs = self._get_feature_idxs(n_features)
        tree = DecisionTreeRegressorScratch(
            max_depth=self.max_depth,
            min_samples_split=self.min_samples_split,
            random_state=np.random.randint(0, 1000)  # Vary seed
                per tree
        )
        tree.feature_idxs = feature_idxs  # Assign random
            features
        tree.fit(X_boot, y_boot)
        self.trees.append(tree)
        self.feature_subsets.append(feature_idxs)
    return self

def predict(self, X: np.ndarray) -> np.ndarray:
    predictions = np.array([tree.predict(X) for tree in self.
        trees])
    return np.mean(predictions, axis=0)
```

## 3.3   Evaluation Metrics

MAE: $\frac{1}{n} \sum |y_i - \hat{y}_i|$

$R^2$: $1 - \frac{SS_{res}}{SS_{tot}}$

Manual split: 80/20 with fixed seed.

# Chapter 4

# Implementation and Demo

## 4.1 Training Scripts

Three separate scripts train and pickle models: - `linearRegression_scratch.py`: Saves `lr_insurance_model_scratch.pkl`, $R^2 \approx 0.7505$. - `gradientBoosting.py`: Saves `gb_insurance_model` early stops at low residual. - `randomForest_scratch.py`: Saves `rf_insurance_model_scratch.pkl`, $R^2 \approx 0.8452$.

Example output for Linear Regression:

```
Model saved to lr_insurance_model_scratch.pkl
Training R² score: 0.7505
Intercept: 1018.32
Coefficients: [ 273.119  -64.774  290.084  473.981 12350.934 -132.568]
```

## 4.2 Testing Script

Unified `test.py` loads models, evaluates on test set, and runs interactive mode with validation.

Evaluation table (sample):

| Model | MAE | $R^2$ Score |
|-------|---------|--------|
| GB | 2564.12 | 0.8621 |
| RF | 2828.11 | 0.8503 |
| LR | 4068.09 | 0.7537 |

Table 4.1: Model Performance on Test Set

Interactive demo: User inputs validated (e.g., age 18-100). Outputs predictions, e.g.:

```
GB: $1723.45
RF: $1698.76
```

```
LR: $1654.32
```

To demo: Run `python test.py` in terminal. Handles empty/invalid inputs gracefully.

## 4.3   GitHub Repository

Project uploaded to `https://github.com/MIDSTAN/Health-Insurance-Premium-Prediction`.
README.md includes: - Installation: `pip install pandas numpy scikit-learn` -
Usage: Run training scripts, then test.py. - Dataset link: Kaggle Insurance. - Screenshots
of outputs. - License: MIT.

# Chapter 5

# Results and Discussion

## 5.1   Performance Analysis

GB outperforms with lowest MAE ($\approx$2564) and highest $R^2$ ($\approx$0.86), due to sequential error correction. RF follows closely ($R^2$ $\approx$0.85), benefiting from ensemble diversity. LR lags ($R^2$ $\approx$0.75) as it assumes linearity, missing interactions like smoker-age.

## 5.2   Insights

- Smoker feature dominates coefficients ($\approx$12350 in LR), highlighting risk. - Scratch RF is slower but educational; real-world use scikit-learn. - Overfitting likely on full train; test split mitigates.

## 5.3   Challenges

- RF split finding: O(n log n) per node; optimized for small data. - Pickling custom classes: Requires definitions in loader.

# Chapter 6

# Conclusion and Future Work

## 6.1 Summary

This project successfully implements and compares scratch ML models for insurance prediction, achieving strong metrics. Interactive demo enhances usability.

## 6.2 Future Enhancements

- Add Ridge/Lasso for LR. - Full GB from scratch (custom trees). - Hyperparameter optimization via grid search. - Deploy as web app (Flask/Streamlit). - Include presentation slides (10-12): Title, Intro, Methods, Results, Demo, Conclusion.

## 6.3 Acknowledgments

Thanks to xAI Grok for code assistance.

# Chapter 7

# Full Code Listings

## 7.1 Linear Regression Scratch

```python
import pandas as pd
import numpy as np
import pickle
import os

# Path to the dataset
data_path = "/home/midstan/Documents/Health Insurance Premium/
    Model/Converting Dataset for Training/Dataset/
    insurance_converted.csv"

# Load the dataset
df = pd.read_csv(data_path)

# Prepare features (X) and target (y)
X = df[['age', 'sex', 'bmi', 'children', 'smoker', 'region']].
    values
y = df['expenses'].values

class LinearRegressionScratch:
    """
    Linear Regression implemented from scratch using the normal
    equation.
    """
    def __init__(self):
        self.coef_ = None
        self.intercept_ = None

```

```python
24      def fit(self, X, y):
25          # Add bias term (intercept column of ones)
26          X_b = np.c_[np.ones((X.shape[0], 1)), X]
27          # Normal equation: theta = (X^T X)^-1 X^T y
28          theta_best = np.linalg.inv(X_b.T.dot(X_b)).dot(X_b.T).dot
                (y)
29          self.intercept_ = theta_best[0]
30          self.coef_ = theta_best[1:]
31          return self
32
33      def predict(self, X):
34          # Add bias term
35          X_b = np.c_[np.ones((X.shape[0], 1)), X]
36          # Predictions: X_b * theta
37          return X_b.dot(np.r_[self.intercept_, self.coef_])
38
39  # Instantiate the model
40  lr_model = LinearRegressionScratch()
41
42  # Train the model on the full dataset
43  lr_model.fit(X, y)
44
45  # Compute R  score manually
46  y_pred = lr_model.predict(X)
47  ss_res = np.sum((y - y_pred) ** 2)
48  ss_tot = np.sum((y - np.mean(y)) ** 2)
49  r2_score = 1 - (ss_res / ss_tot)
50
51  # Save the trained model to a file for later use
52  model_path = 'lr_insurance_model_scratch.pkl'
53  with open(model_path, 'wb') as f:
54      pickle.dump(lr_model, f)
55
56  print(f"Model saved to {model_path}")
57  print(f"Training R  score: {r2_score:.4f}")
58  print(f"Intercept: {lr_model.intercept_:.2f}")
59  print(f"Coefficients: {lr_model.coef_}")
```

## 7.2   Gradient Boosting

```python
import pandas as pd
import numpy as np
from sklearn.tree import DecisionTreeRegressor
import pickle
import os

# Path to the dataset
data_path = "/home/midstan/Documents/Health Insurance Premium/
    Model/Converting Dataset for Training/Dataset/
    insurance_converted.csv"

# Load the dataset
df = pd.read_csv(data_path)

# Prepare features (X) and target (y)
X = df[['age', 'sex', 'bmi', 'children', 'smoker', 'region']].
    values
y = df['expenses'].values

class GradientBoostingRegressor:
    """
    A simple implementation of Gradient Boosting for Regression.
    - Base model: Mean of the target values.
    - Weak learners: Decision Trees fitted to negative gradients
    (residuals).
    - Stops early if mean absolute residual is below a tolerance.
    """
    def __init__(self, n_estimators=100, learning_rate=0.1,
        max_depth=3, tol=1.0):
        """
        :param n_estimators: Maximum number of trees.
        :param learning_rate: Shrinkage factor for each tree's
    contribution.
        :param max_depth: Maximum depth of each decision tree.
        :param tol: Tolerance for mean absolute residual to stop
    early.
        """
        self.n_estimators = n_estimators
        self.learning_rate = learning_rate
        self.max_depth = max_depth
```

```
34          self.tol = tol # Tolerance for residuals (adjusted for
                  expenses scale ~thousands)
35          self.trees = [] # List to store decision trees
36          self.initial_prediction = None # Mean as base prediction
37
38      def fit(self, X, y):
39          # Initialize predictions with the mean (base model)
40          self.initial_prediction = np.mean(y)
41          predictions = np.full_like(y, self.initial_prediction,
                  dtype=np.float64)
42          # Compute initial residuals
43          residuals = y - predictions
44          # Build trees iteratively
45          for i in range(self.n_estimators):
46              # Fit a decision tree to the current residuals (
                      negative gradient for regression)
47              tree = DecisionTreeRegressor(max_depth=self.max_depth
                      , random_state=42)
48              tree.fit(X, residuals)
49              # Predict residuals with the tree
50              tree_pred = tree.predict(X)
51              # Update predictions: add learning_rate * tree
                      prediction
52              predictions += self.learning_rate * tree_pred
53              # Update residuals
54              residuals = y - predictions
55              # Store the tree
56              self.trees.append(tree)
57              # Early stopping: check if mean absolute residual is
                      below tolerance
58              mean_abs_residual = np.mean(np.abs(residuals))
59              if mean_abs_residual < self.tol:
60                  print(f"Early␣stopping␣at␣iteration␣{i+1}:␣mean␣
                          abs␣residual␣=␣{mean_abs_residual:.2f}")
61                  break
62          print(f"Training␣completed.␣Final␣mean␣abs␣residual:␣{np.
                  mean(np.abs(residuals)):.2f}")
63          print(f"Number␣of␣trees␣built:␣{len(self.trees)}")
64          return self
65
66      def predict(self, X):
```

```
67          # Start with initial prediction
68          predictions = np.full(X.shape[0], self.initial_prediction
               , dtype=np.float64)
69          # Add contributions from all trees
70          for tree in self.trees:
71              predictions += self.learning_rate * tree.predict(X)
72          return predictions
73
74 # Instantiate the model (adjust hyperparameters as needed)
75 gb_model = GradientBoostingRegressor(
76     n_estimators=100000,
77     learning_rate=0.1,
78     max_depth=3,
79     tol=1.0 # Stop if mean abs residual < 1 (small relative to
           expenses)
80 )
81
82 # Train the model on the full dataset
83 gb_model.fit(X, y)
84
85 # Save the trained model to a file for later use
86 model_path = 'gb_insurance_model.pkl'
87 with open(model_path, 'wb') as f:
88     pickle.dump(gb_model, f)
89 print(f"Model saved to {model_path}")
90
91 # Example: How to load and predict after training
92 # (Uncomment to test if you have sample input data)
93 # Load the model
94 with open(model_path, 'rb') as f:
95     loaded_model = pickle.load(f)
96 # Example new input (must match feature order and types)
97 new_X = np.array([[18,1,33.8,1,0,3]]) # Sample from your data
98 prediction = loaded_model.predict(new_X)
99 print(f"Predicted expenses: {prediction[0]:.2f}")
```

## 7.3   Random Forest Scratch

```
1 import pandas as pd
2 import numpy as np
```

```python
import pickle
import os
from typing import Tuple, List

# Path to the dataset
data_path = "/home/midstan/Documents/Health Insurance Premium/
    Model/Converting Dataset for Training/Dataset/
    insurance_converted.csv"

# Load the dataset
df = pd.read_csv(data_path)

# Prepare features (X) and target (y)
X = df[['age', 'sex', 'bmi', 'children', 'smoker', 'region']].
    values
y = df['expenses'].values

class Node:
    """
    Node in the decision tree.
    """
    def __init__(self, feature_idx=None, threshold=None, left=
        None, right=None, value=None):
        self.feature_idx = feature_idx
        self.threshold = threshold
        self.left = left
        self.right = right
        self.value = value  # Leaf node value (mean of targets)

def mse(y: np.ndarray) -> float:
    """
    Mean Squared Error for a set of targets.
    """
    return np.mean((y - np.mean(y)) ** 2)

def best_split(X: np.ndarray, y: np.ndarray, feature_idxs: List[
    int]) -> Tuple[int, float, float, np.ndarray, np.ndarray]:
    """
    Find the best split: feature and threshold minimizing MSE.
    """
```

```
38      best_idx , best_thresh , best_mse , best_left , best_right = None
            , None , float('inf'), None , None
39      current_mse = mse(y)
40
41      for idx in feature_idxs:
42          thresholds = np.unique(X[:, idx])
43          for thresh in thresholds:
44              left_mask = X[:, idx] <= thresh
45              right_mask = ~left_mask
46              if np.sum(left_mask) == 0 or np.sum(right_mask) == 0:
47                  continue
48              left_y , right_y = y[left_mask], y[right_mask]
49              weighted_mse = (len(left_y) * mse(left_y) + len(
                    right_y) * mse(right_y)) / len(y)
50              if weighted_mse < best_mse:
51                  best_idx = idx
52                  best_thresh = thresh
53                  best_mse = weighted_mse
54                  best_left = left_y
55                  best_right = right_y
56
57      return best_idx , best_thresh , best_mse , best_left , best_right
58
59  def build_tree(X: np.ndarray , y: np.ndarray , feature_idxs: List[
        int], max_depth: int = None , min_samples_split: int = 2, depth
        : int = 0) -> Node:
60      """
61      Recursively build the decision tree.
62      """
63      if len(y) < min_samples_split or (max_depth is not None and
            depth >= max_depth):
64          return Node(value=np.mean(y))
65
66      idx , thresh , _, _, _ = best_split(X, y, feature_idxs)
67      if idx is None:
68          return Node(value=np.mean(y))
69
70      left_mask = X[:, idx] <= thresh
71      right_mask = ~left_mask
72
```

```
73        left = build_tree(X[left_mask], y[left_mask], feature_idxs,
              max_depth, min_samples_split, depth + 1)
74        right = build_tree(X[right_mask], y[right_mask], feature_idxs
              , max_depth, min_samples_split, depth + 1)
75
76        return Node(idx, thresh, left, right)
77
78   def predict_tree(node: Node, x: np.ndarray) -> float:
79       """
80       Predict with a single tree.
81       """
82       if node.value is not None:
83           return node.value
84       if x[node.feature_idx] <= node.threshold:
85           return predict_tree(node.left, x)
86       else:
87           return predict_tree(node.right, x)
88
89   class DecisionTreeRegressorScratch:
90       """
91       Decision Tree Regressor from scratch (CART-like, MSE
           criterion).
92       """
93       def __init__(self, max_depth: int = None, min_samples_split:
              int = 2, random_state: int = 42):
94           self.max_depth = max_depth
95           self.min_samples_split = min_samples_split
96           self.random_state = random_state
97           self.root = None
98           self.feature_idxs = None
99
100      def fit(self, X: np.ndarray, y: np.ndarray):
101          np.random.seed(self.random_state)
102          self.n_features = X.shape[1]
103          self.feature_idxs = list(range(self.n_features))
104          self.root = build_tree(X, y, self.feature_idxs, self.
                 max_depth, self.min_samples_split)
105          return self
106
107      def predict(self, X: np.ndarray) -> np.ndarray:
108          return np.array([predict_tree(self.root, x) for x in X])
```

```
109
110  class RandomForestRegressorScratch:
111      """
112      Random Forest Regressor from scratch: Ensemble of decision
         trees with bagging and random features.
113      """
114      def __init__(self, n_estimators: int = 100, max_depth: int =
             3, min_samples_split: int = 2, max_features: str = 'sqrt',
              random_state: int = 42):
115          self.n_estimators = n_estimators
116          self.max_depth = max_depth
117          self.min_samples_split = min_samples_split
118          self.max_features = max_features
119          self.random_state = random_state
120          self.trees = []
121          self.feature_subsets = []
122
123      def _bootstrap_sample(self, X: np.ndarray, y: np.ndarray) ->
             Tuple[np.ndarray, np.ndarray]:
124          n_samples = X.shape[0]
125          idxs = np.random.choice(n_samples, n_samples, replace=
                 True)
126          return X[idxs], y[idxs]
127
128      def _get_feature_idxs(self, n_features: int) -> List[int]:
129          if self.max_features == 'sqrt':
130              return np.random.choice(n_features, max(1, int(np.
                     sqrt(n_features))), replace=False).tolist()
131          elif self.max_features == 'log2':
132              return np.random.choice(n_features, max(1, int(np.
                     log2(n_features))), replace=False).tolist()
133          else:
134              raise ValueError("max_features must be 'sqrt' or '
                     log2'")
135
136      def fit(self, X: np.ndarray, y: np.ndarray):
137          np.random.seed(self.random_state)
138          n_features = X.shape[1]
139          for _ in range(self.n_estimators):
140              X_boot, y_boot = self._bootstrap_sample(X, y)
141              feature_idxs = self._get_feature_idxs(n_features)
```

```
142              tree = DecisionTreeRegressorScratch(
143                  max_depth=self.max_depth,
144                  min_samples_split=self.min_samples_split,
145                  random_state=np.random.randint(0, 1000)  # Vary
                         seed per tree
146              )
147              tree.feature_idxs = feature_idxs  # Assign random
                    features
148              tree.fit(X_boot, y_boot)
149              self.trees.append(tree)
150              self.feature_subsets.append(feature_idxs)
151          return self
152
153      def predict(self, X: np.ndarray) -> np.ndarray:
154          predictions = np.array([tree.predict(X) for tree in self.
                 trees])
155          return np.mean(predictions, axis=0)
156
157  # Instantiate the model (hyperparameters to roughly match
         original: 100 trees, depth 3, sqrt features)
158  rf_model = RandomForestRegressorScratch(
159      n_estimators=100,
160      max_depth=3,
161      min_samples_split=2,
162      max_features='sqrt',
163      random_state=42
164  )
165
166  # Train the model on the full dataset
167  rf_model.fit(X, y)
168
169  # Compute R  score manually
170  y_pred = rf_model.predict(X)
171  ss_res = np.sum((y - y_pred) ** 2)
172  ss_tot = np.sum((y - np.mean(y)) ** 2)
173  r2_score = 1 - (ss_res / ss_tot)
174
175  # Save the trained model to a file for later use
176  model_path = 'rf_insurance_model_scratch.pkl'
177  with open(model_path, 'wb') as f:
178      pickle.dump(rf_model, f)
```

```
179
180 print(f"Model␣saved␣to␣{model_path}")
181 print(f"Training␣R␣␣score:␣{r2_score:.4f}")
182 print(f"Number␣of␣trees:␣{len(rf_model.trees)}")
```

## 7.4    Test Script

```
1  import numpy as np
2  import pickle
3  import pandas as pd
4  from sklearn.tree import DecisionTreeRegressor # Only for GB
      prediction (loaded trees)
5  import os
6  from typing import Tuple, List
7
8  # Path to the dataset (for evaluation)
9  data_path = "/home/midstan/Documents/Health␣Insurance␣Premium/
      Model/Converting␣Dataset␣for␣Training/Dataset/
      insurance_converted.csv"
10
11 # Define the GradientBoostingRegressor class (must match the one
      used during training)
12 class GradientBoostingRegressor:
13     # ... (full class definition as above)
14
15 # Define LinearRegressionScratch class
16 class LinearRegressionScratch:
17     # ... (full class definition as above)
18
19 # Define Random Forest Scratch classes
20 class Node:
21     # ... (full class definition as above)
22
23 # ... (all other functions and classes as in training script)
24
25 # Manual train_test_split function
26 def manual_train_test_split(X, y, test_size=0.2, random_state=42)
      :
27     np.random.seed(random_state)
28     indices = np.arange(X.shape[0])
```

```
29        np.random.shuffle(indices)
30        split = int((1 - test_size) * len(indices))
31        train_idx = indices[:split]
32        test_idx = indices[split:]
33        return X[train_idx], X[test_idx], y[train_idx], y[test_idx]
34
35   # Manual MAE and R2 functions
36   def manual_mean_absolute_error(y_true, y_pred):
37        return np.mean(np.abs(y_true - y_pred))
38
39   def manual_r2_score(y_true, y_pred):
40        ss_res = np.sum((y_true - y_pred) ** 2)
41        ss_tot = np.sum((y_true - np.mean(y_true)) ** 2)
42        return 1 - (ss_res / ss_tot)
43
44   # Model paths (full paths for scratch models; adjust GB if needed
         )
45   model_paths = {
46        'gb': '/home/midstan/Documents/Health Insurance Premium/
            gb_insurance_model.pkl',
47        'rf': '/home/midstan/Documents/Health Insurance Premium/
            rf_insurance_model_scratch.pkl',
48        'lr': '/home/midstan/Documents/Health Insurance Premium/
            lr_insurance_model_scratch.pkl'
49   }
50
51   # Load models
52   models = {}
53   loaded_successfully = []
54   for name, path in model_paths.items():
55        try:
56            with open(path, 'rb') as f:
57                if name == 'gb':
58                    models[name] = pickle.load(f)
59                elif name == 'rf':
60                    models[name] = pickle.load(f)
61                elif name == 'lr':
62                    models[name] = pickle.load(f)
63            print(f"{name.upper()} model loaded successfully!")
64            loaded_successfully.append(name)
65        except FileNotFoundError:
```

```
66          print(f"Error:␣Model␣file␣'{path}'␣not␣found.␣Skipping␣{
                name}␣model.")
67      except Exception as e:
68          print(f"Error␣loading␣{name}␣model:␣{e}.␣Skipping.")
69
70  if not loaded_successfully:
71      print("No␣models␣loaded.␣Exiting.")
72      exit(1)
73
74  # Evaluation on test set
75  print("\n" + "="*50)
76  print("MODEL␣PERFORMANCE␣COMPARISON␣ON␣TEST␣SET␣(80/20␣split)")
77  print("="*50)
78  try:
79      # Load data for evaluation
80      df = pd.read_csv(data_path)
81      X = df[['age', 'sex', 'bmi', 'children', 'smoker', 'region'
             ]].values
82      y = df['expenses'].values
83      X_train, X_test, y_train, y_test = manual_train_test_split(X,
              y, test_size=0.2, random_state=42)
84      print("Model\t\tMAE\t\tR␣Score")
85      print("-" * 40)
86      for name in loaded_successfully:
87          model = models[name]
88          y_pred = model.predict(X_test)
89          mae = manual_mean_absolute_error(y_test, y_pred)
90          r2 = manual_r2_score(y_test, y_pred)
91          print(f"{name}\t\t{mae:.2f}\t\t{r2:.4f}")
92      print("Evaluation␣completed␣successfully!")
93  except Exception as e:
94      print(f"Error␣during␣evaluation:␣{e}")
95      import traceback
96      traceback.print_exc()  # Print full traceback for debugging
97
98  print("\n" + "="*50)
99  print("INTERACTIVE␣PREDICTION␣MODE")
100 print("="*50)
101
102 def get_user_input():
103     """
```

```python
      Collects user input for the features in the correct order: [
          age , sex , bmi , children , smoker , region]
      with robust error handling for invalid inputs.
      """
      print("\nEnter the following details for insurance premium 
          prediction:")

      # Helper function to get validated int input
      def get_valid_int(prompt , min_val=None , max_val=None):
          while True:
              try:
                  value = input(prompt).strip()
                  if not value:
                      print("Input cannot be empty. Please try 
                          again.")
                      continue
                  int_val = int(value)
                  if min_val is not None and int_val < min_val:
                      print(f"Value must be at least {min_val}. 
                          Please try again.")
                      continue
                  if max_val is not None and int_val > max_val:
                      print(f"Value must be at most {max_val}. 
                          Please try again.")
                      continue
                  return int_val
              except ValueError:
                  print("Invalid integer value. Please enter a 
                      valid number.")

      # Helper function to get validated float input
      def get_valid_float(prompt , min_val=None , max_val=None):
          while True:
              try:
                  value = input(prompt).strip()
                  if not value:
                      print("Input cannot be empty. Please try 
                          again.")
                      continue
                  float_val = float(value)
                  if min_val is not None and float_val < min_val:
```

```python
138                     print(f"Value must be at least {min_val}. 
                            Please try again.")
139                     continue
140                 if max_val is not None and float_val > max_val:
141                     print(f"Value must be at most {max_val}. 
                            Please try again.")
142                     continue
143                 return float_val
144             except ValueError:
145                 print("Invalid float value. Please enter a valid 
                        number.")
146
147     age = get_valid_int("Age (18-100): ", min_val=18, max_val
            =100)
148     sex = get_valid_int("Sex (0 for female, 1 for male): ", 
            min_val=0, max_val=1)
149     bmi = get_valid_float("BMI (10-60): ", min_val=10, max_val
            =60)
150     children = get_valid_int("Number of children (0-5): ", 
            min_val=0, max_val=5)
151     smoker = get_valid_int("Smoker (0 for no, 1 for yes): ", 
            min_val=0, max_val=1)
152     region = get_valid_int("Region (1-4): ", min_val=1, max_val
            =4)
153
154     return np.array([[age, sex, bmi, children, smoker, region]])
155
156 # Interactive prediction loop
157 while True:
158     new_X = get_user_input()
159     # Make predictions with all loaded models
160     print("\nPredicted insurance expenses from all models:")
161     print("-" * 40)
162     for name in loaded_successfully:
163         model = models[name]
164         prediction = model.predict(new_X)
165         print(f"{name.upper()}: ${prediction[0]:.2f}")
166     # Ask if user wants to predict another
167     while True:
168         continue_choice = input("\nPredict another? (y/n): ").
                lower().strip()
```

```python
169            if continue_choice in ['y', 'yes']:
170                break
171            elif continue_choice in ['n', 'no']:
172                print("Goodbye!")
173                exit(0)
174            else:
175                print("Please enter 'y' for yes or 'n' for no.")
```

# Chapter 8

# Presentation Slides Outline

- Slide 1: Title and Agenda

- Slide 2: Problem Statement

- Slide 3: Dataset Overview

- Slide 4: Linear Regression Details

- Slide 5: Gradient Boosting

- Slide 6: Random Forest

- Slide 7: Implementation Flow

- Slide 8: Results Table

- Slide 9: Demo Screenshots

- Slide 10: Challenges and Insights

- Slide 11: Conclusion

- Slide 12: Q&A

Note: Slides prepared in PowerPoint/Google Slides with code snippets and plots.