**Group 45: Health Insurance Cost Prediction**

**Name=Bikal Bista(A43323)**

**1.Objective**

- **Goal**: Develop a machine learning model using Support Vector Machines (SVM) to predict healthcare costs for new customers.

- **Outcome**: Provide accurate and fair insurance premium estimations.

---

**2.Dataset Overview**

- **Training Data**: 2,215 instances with features like gender, marital status, area of residence, BMI, smoking status, age class, and healthcare costs.

- **Test Data**: 550 instances with only features (costs unknown).

---

**3.Methodology**

1. **Data Preprocessing**:

    o Categorical features encoded using LabelEncoder.

    o Data scaled using StandardScaler.

    o Dataset split: 80% training and 20% testing.

2. **Model Training**:

    o Algorithm: Support Vector Regressor (SVR)

    o Key Parameters: Kernel = 'rbf', C = 100, Gamma = 0.1, Epsilon = 0.1

3. **Feature Importance Analysis**:

    o Technique: Permutation Importance

    o Tool: Scikit-learn's permutation_importance function

---

**4.Model Evaluation**

- **Performance Metrics**:

    o $R^2$ Score: 0.19

o Mean Squared Error (MSE): 89,513,449.68

### 5.Observations:

    o Model shows moderate predictive performance.

    o Smoking status (fumador) is the most impactful feature.

---

## 6.Feature Importance

- **Top Influential Features**:

  1. Smoking Status (fumador)

  2. Age Class (class_etaria)

  3. Body Mass Index (imc)

- **Visualization**: The bar chart below illustrates the importance of each feature: