# Insurance Cost Prediction

Dataset:
https://drive.google.com/file/d/1NBk1TFkK4NeKdodR2DxIdBp2Mk1mh4AS/view?usp=drive_link

## Problem Statement

Insurance companies need to accurately predict the cost of health insurance for individuals to set premiums appropriately. However, traditional methods of cost prediction often rely on broad actuarial tables and historical averages, which may not account for the nuanced differences among individuals. By leveraging machine learning techniques, insurers can predict more accurately the insurance costs tailored to individual profiles, leading to more competitive pricing and better risk management.

## Insurance Cost Prediction need

The primary need for this project arises from the challenges insurers face in pricing policies accurately while remaining competitive in the market. Inaccurate predictions can lead to losses for insurers and unfairly high premiums for policyholders. By implementing a machine learning model, insurers can:
- Enhance Precision in Pricing: Use individual data points to determine premiums that reflect actual risk more closely than generic estimates.
- Increase Competitiveness: Offer rates that are attractive to consumers while ensuring that the pricing is sustainable for the insurer.
- Improve Customer Satisfaction: Fair and transparent pricing based on personal health data can increase trust and satisfaction among policyholders.
- Enable Personalized Offerings: Create customized insurance packages based on predicted costs, which can cater more directly to the needs and preferences of individuals.
- Risk Assessment: Insurers can use the model to refine their risk assessment processes, identifying key factors that influence costs most significantly.
- Policy Development: The insights gained from the model can inform the development of new insurance products or adjustments to existing ones.
- Strategic Decision Making: Predictive analytics can aid in broader strategic decisions, such as entering new markets or adjusting policy terms based on risk predictions.
- Customer Engagement: Insights from the model can be used in customer engagement initiatives, such as personalized marketing and tailored advice for policyholders.

## Data description

The dataset comprises the following 11 attributes:
1. Age: Numeric, ranging from 18 to 66 years.
2. Diabetes: Binary (0 or 1), where 1 indicates the presence of diabetes.

3. BloodPressureProblems: Binary (0 or 1), indicating the presence of blood pressure-related issues.
4. AnyTransplants: Binary (0 or 1), where 1 indicates the person has had a transplant.
5. AnyChronicDiseases: Binary (0 or 1), indicating the presence of any chronic diseases.
6. Height: Numeric, measured in centimeters, ranging from 145 cm to 188 cm.
7. Weight: Numeric, measured in kilograms, ranging from 51 kg to 132 kg.
8. KnownAllergies: Binary (0 or 1), where 1 indicates known allergies.
9. HistoryOfCancerInFamily: Binary (0 or 1), indicating a family history of cancer.
10. NumberOfMajorSurgeries: Numeric, counting the number of major surgeries, ranging from 0 to 3 surgeries.
11. PremiumPrice: Numeric, representing the premium price in currency, ranging from 15,000 to 40,000.

# Block 1 Tableau Visualization

The core idea is to harness the power of visual analytics to dissect the complex relationships between individual health profiles and insurance costs. By visualizing these relationships, insurers can gain a clearer understanding of the factors that most significantly influence premium prices, allowing for more accurate and equitable pricing models.

## Goal

The goal of this dashboard is threefold:
1. Visualization of Key Data Points: To clearly display the distribution and impact of various health-related factors on insurance premiums.
2. Predictive Analysis: To apply statistical tools within Tableau to predict insurance costs based on individual risk factors, aiding in better premium setting.
3. Strategic Insight Generation: To enable insurance companies to derive strategic insights that can influence policy adjustments, risk assessments, and customer segmentation.

## Dashboard Components

1. Summary Statistics Dashboard:
   ● Key Metrics Display: Show average premium, average age, and distributions of key health conditions.
   ● Count of Individuals by Health Conditions: Bar or pie charts showing the count of individuals with diabetes, blood pressure problems, transplants, chronic diseases, allergies, and cancer history.
2. Premium Pricing Analysis Dashboard:
   ● Premium Distribution: Histogram or boxplot displaying the distribution of premium prices.

- Premiums by Age Group and Health Factors: Line charts or bar charts showing average premiums across different age groups segmented by health conditions like diabetes, blood pressure issues, etc.
- Correlation Heatmap: Display correlations between all numerical factors including age, height, weight, number of surgeries, and premium price.

3. Risk Factors Analysis Dashboard:
   - Surgical Impact on Premiums: Analysis of how the number of major surgeries correlates with premium costs.
   - Impact of Chronic Conditions: Stacked bar charts showing premium variations with the presence of chronic diseases, transplants, and other health issues.
   - Allergies and Family History Influence: Explore how known allergies and a history of cancer in the family impact premium prices.

4. Demographic Insights Dashboard:
   - Premiums by Height and Weight (BMI): Scatter plot analyzing the relationship between body metrics and premium costs.
   - Geographical Analysis: If region data is available or can be inferred, provide a geographical breakdown of premium costs.

## Interactive Features

- Filters and Sliders: Allow users to filter by age, weight, height, and health conditions to explore specific subsets of data.
- Drill-down Capabilities: Enable clicking on a specific chart element to view more detailed data or related visualizations.
- Tooltips: Provide detailed information when hovering over data points or bars in charts.

## Insights Gathering

- Develop Predictive Insights: Utilize regression analysis visuals to predict premium prices based on input variables.
- Identify Risk Profiles: Highlight risk profiles that typically lead to higher premiums and suggest mitigating strategies.
- Policy Recommendations: Based on the visualized data, offer recommendations for insurance policy adjustments or new product development.

# Block 2: EDA and Hypothesis Testing for Insurance Cost Prediction

For the Insurance Cost Prediction project, EDA will involve visualizing distributions, identifying outliers, and exploring correlations between different variables like age, diabetes status, weight, and premium costs. This analysis aims to unearth significant predictors of insurance costs and understand the demographic and health-related characteristics that most influence premium pricing.

Hypothesis testing will be used to formally test assumptions such as:
- Are premium costs significantly higher for smokers compared to non-smokers?
- Does the presence of chronic diseases lead to higher insurance premiums?
- Is there a significant difference in premium costs based on the number of major surgeries a person has had?

These statistical tests will help validate whether the observed patterns in the data are statistically significant or occur by chance, thereby reinforcing the robustness of subsequent predictive modeling.

## Ideas for EDA and Hypothesis Testing

1. Distribution Analysis:
   - Plot any and all visualization that could not be made in tableau.
2. Correlation Analysis:
   - Generate a correlation matrix or heatmap to visualize the relationships between all numerical variables.
   - Focus on the correlation between premium prices and potential predictors to identify strong associations.
3. Outlier Detection:
   - Identify outliers in key variables using IQR (Interquartile Range) method or Z-scores.
   - Assess the impact of outliers on the overall distribution and consider strategies for handling them.
4. Hypothesis Testing:
   - T-tests/ANOVA: Use these tests to compare the means of premium prices across different groups defined by categorical variables (e.g., smokers vs. non-smokers, number of surgeries).
   - Chi-square tests: Evaluate the association between two categorical variables (e.g., presence of chronic disease and history of cancer in family).
   - Regression Analysis: Apply linear regression to test hypotheses about the impact of various predictors on premium prices.

# Block 3: ML Modeling

1. Data Preprocessing:
   - Handling Missing Values: Although initial data checks may not show missing values, always prepare to implement strategies for handling them.
   - Feature Engineering: Create new features that might improve model performance, such as Body Mass Index (BMI) from height and weight.
   - Scaling and Encoding: Apply appropriate scaling to numerical features and encoding to categorical features to prepare the data for machine learning algorithms.

2. Model Selection:
   - Linear Regression: Start with a simple model to establish a baseline for prediction accuracy.
   - Tree-based Models: Implement models like Decision Trees, Random Forests, and Gradient Boosting Machines for their ability to handle non-linear relationships and feature importance analysis.
   - Neural Networks: Explore more complex models like neural networks if the initial models show promising results but require more flexibility in capturing interactions.
3. Model Evaluation and Validation:
   - Cross-Validation: Use techniques like k-fold cross-validation to ensure that the model performs well across different subsets of the dataset.
   - Performance Metrics: Depending on the business objective, use metrics like RMSE (Root Mean Squared Error), MAE (Mean Absolute Error), or $R^2$ (Coefficient of Determination) to evaluate model performance.
   - Confidence Intervals/Prediction Intervals: Provide these intervals along with predictions to give users an idea of prediction reliability.
4. Interpretability and Explainability:
   - Feature Importance: Use techniques like permutation importance in tree-based models or SHAP values to explain the influence of each feature on the prediction.
   - Model Insights: Translate the model's findings into actionable business insights, such as identifying risk factors that significantly increase insurance costs, which can be used for targeted interventions.

# Block 4: Web-Based Calculator for Estimating Insurance Premiums

The final phase of the Insurance Cost Prediction project involves deploying the developed machine learning model into a practical, user-friendly application. This application will serve as a web-based calculator for insurance agents or customers, enabling them to estimate insurance premiums based on individual data inputs. Deployment will be carried out through a simple Flask API or a Streamlit application, both of which are popular frameworks for deploying data science projects due to their ease of use and flexibility.

## Deployment Objectives
- Accessibility: Make the predictive model accessible to users without the need for technical background, directly via a web interface.
- Real-time Estimations: Provide real-time insurance cost predictions as users input or modify their data.
- User-Friendly Interface: Ensure the interface is intuitive, guiding users smoothly through the process of entering their data and receiving their premium estimations.

- Showcasing your Project: it is very important for you to be able to showcase something to the recruiters. It will also help you showcase some of your MLOps skills

## Deployment Process

1. API Development with Flask:
   - Flask Setup: Create a Flask application that serves as the backend. The Flask app will handle requests from the front end, process them using the machine learning model, and return the premium predictions.
   - Endpoint Creation: Design endpoints that receive user inputs in the form of JSON and return the estimated premiums.
   - Model Integration: Integrate the trained machine learning model into the Flask application, ensuring it can access and process the input data effectively.
2. Streamlit Application:
   - Streamlit Setup: Develop a Streamlit application creating the user interface.
   - User Input Forms: Use Streamlit widgets to create forms where users can input their data (age, BMI, health conditions, etc.).
   - Model Invocation: Set up the application to use the model to predict premiums based on user inputs and display the results directly in the application.

## Note:

- The suggestions/Ideas provided above are intended to assist you. The primary aim is to offer guidance on what aspects can be analyzed. If no valuable insights can be derived from a particular analysis, feel free to skip it.