# Fraud detection in insurance claims using the Insurance Fraud

## Problem Description :

Insurance fraud is a growing problem that costs the industry billions of dollars every year. The objective of this project is to build a machine learning model that can accurately detect fraudulent insurance claims, based on data about the claimant and the claim.

The dataset for this project is the Insurance Fraud dataset, which contains information on over 50,000 insurance claims, including details about the claimant, the policy, and the claim itself. The dataset was obtained from Kaggle (https://www.kaggle.com/roshansharma/insurance-claim).

The deliverables for this project are:

- A trained machine learning model that can accurately detect fraudulent insurance claims
- A report detailing the methodology used to build the model and the results obtained
- A Jupyter notebook containing the code used to build and evaluate the model

**Dataset Description**

The Insurance Fraud dataset contains information on 50,000 insurance claims, including details about the claimant, the policy, and the claim itself. The dataset includes the following columns:

- Policy ID: A unique identifier for the policy
- Claim ID: A unique identifier for the claim
- Claim Amount: The amount of the claim
- Policy Holder's Gender: The gender of the policy holder
- Policy Holder's Age: The age of the policy holder
- Number of Vehicles: The number of vehicles covered by the policy
- Total Claim Amount: The total amount claimed on the policy
- Fraudulent: Whether the claim was fraudulent or not (0 = not fraudulent, 1 = fraudulent)

**Background Information**

Insurance fraud is a growing problem that affects both insurance companies and policyholders. Fraudulent claims can lead to higher premiums for everyone, and can even result in insurance companies going bankrupt.

Detecting fraudulent insurance claims is a challenging problem, as fraudsters often go to great lengths to hide their activities. Machine learning models have been shown to be effective at detecting fraudulent insurance claims, by identifying patterns and anomalies in the data that may be indicative of fraud.

# Possible Framework:

### 1. Problem Understanding

- Define the problem and the objectives of the project
- Understand the business context and the stakeholders' needs
- Determine the key performance indicators (KPIs) that will be used to evaluate the model

### 2. Data Collection and Understanding

- Identify the sources of data for the project
- Collect the data and perform an initial data quality check
- Understand the features and their relationships with the target variable
- Identify any data preprocessing or feature engineering steps that need to be performed

### 3. Data Preparation

- Perform data preprocessing and feature engineering steps identified in the previous step
- Split the data into training, validation, and test sets
- Normalize or standardize the features to make them comparable
- Address any issues related to class imbalance or missing data

### 4. Model Development

- Select appropriate machine learning algorithms to train and test the model
- Train and optimize the model using the training and validation sets
- Evaluate the model's performance using the test set
- Interpret the model's results and identify any limitations or areas for improvement

### 5. Model Deployment and Monitoring

- Deploy the model in a production environment
- Monitor the model's performance and retrain the model if necessary
- Implement processes to ensure the model is updated and maintained over time

**6. Documentation**

- Document the project and its findings, including the data sources, methodology, and results
- Share the project with stakeholders and other relevant parties
- Provide recommendations for future work and improvements

**7. Review and Iteration**

- Review the project and identify areas for improvement
- Iterate on the project and incorporate any changes or new information

# Code Explanation :

Here is the simple explanation for the code which is provided in the code.py file.

This code is designed to identify fraudulent insurance claims using a dataset called Insurance Fraud. It uses a machine learning algorithm called Random Forest Classifier, which is commonly used for classification tasks like fraud detection. The dataset contains information about insurance claims, including features such as claim amount, policy holder age, and the type of claim being made.

**Requirements to Run the Code**

To run this code, you will need Python 3 and several libraries, including pandas, numpy, matplotlib, and scikit-learn. You can install these libraries using pip, the package installer for Python. The dataset used in this code can be downloaded from Kaggle.

**Section 1: Importing the Libraries**

The first section of the code imports the necessary libraries, including pandas, numpy, matplotlib, and scikit-learn. These libraries are essential for manipulating data, performing calculations, and building machine learning models.

**Section 2: Loading the Dataset**

The second section of the code loads the dataset into a pandas dataframe. This dataset contains information about insurance claims, including whether the claim was fraudulent or not. The data is split into two parts: training data and test data.

**Section 3: Preprocessing the Data**

The third section of the code preprocesses the data to prepare it for machine learning. This includes handling missing values, encoding categorical variables, and scaling numerical variables.

**Section 4: Building the Random Forest Classifier Model**

The fourth section of the code builds a Random Forest Classifier model to predict whether a claim is fraudulent or not. The model is trained on the training data and evaluated on the test data.

**Section 5: Evaluating the Model**

The fifth section of the code evaluates the performance of the Random Forest Classifier model. This includes calculating the accuracy, precision, and recall of the model.

# Future Work :

1. **Feature Engineering:** Feature engineering is an important aspect of building a predictive model, and it involves creating new features from the existing ones. For example, in this dataset, we could create features like the time taken to file a claim, the number of previous claims made by a customer, the amount claimed, and the total premium paid, among others. We could also use external data sources, such as weather data, to create new features.
2. **Model Selection:** In this project, we used a Random Forest Classifier to build the model. However, there are many other classification algorithms that we could try, such as logistic regression, decision trees, and support vector machines. We could use techniques like grid search and cross-validation to determine which algorithm and hyperparameters perform best on our data.
3. **Ensembling:** Ensemble models are a combination of multiple models that are used to make more accurate predictions. We could explore ensemble methods like bagging, boosting, and stacking to improve the accuracy of our predictions.
4. **Data Sampling:** Imbalanced data can be a challenge in fraud detection. To address this issue, we could explore different data sampling techniques such as oversampling, undersampling, and SMOTE to balance the dataset and improve the accuracy of our predictions.
5. **Interpretability:** In fraud detection, it is important to be able to explain why a claim was flagged as fraudulent. We could explore techniques like SHAP values, LIME, and partial dependence plots to explain the predictions made by our model.

**Step-by-Step Guide:**

1. **Feature Engineering:** Identify potential features that could be created from the existing data or external data sources. Create new features and add them to the dataset.
2. **Model Selection:** Choose a set of classification algorithms to evaluate on the dataset. Use techniques like grid search and cross-validation to determine which algorithm and hyperparameters perform best on the data.
3. **Ensembling:** Combine multiple models to create an ensemble model. Use techniques like bagging, boosting, and stacking to improve the accuracy of the predictions.
4. **Data Sampling:** Use techniques like oversampling, undersampling, and SMOTE to address the issue of imbalanced data.
5. **Interpretability:** Use techniques like SHAP values, LIME, and partial dependence plots to explain the predictions made by the model.

# Exercise :

**Try to answers the following questions by yourself to check your understanding for this project. If stuck, detailed answers for the questions are also provided.**

1. **How can you adjust the threshold for the classification model?**
   You can adjust the threshold for the classification model by changing the value of the **threshold** variable in the code. The threshold determines the minimum probability required for a prediction to be classified as fraudulent.

2. **Can you try using a different classification algorithm for this fraud detection problem?**
   Yes, you can try using a different classification algorithm for this fraud detection problem. Some popular algorithms include decision trees, random forests, and support vector machines. To do this, you would need to import the appropriate algorithm from a machine learning library such as scikit-learn and modify the code accordingly.

3. **How can you improve the performance of the model?**
   You can improve the performance of the model by using a larger and more diverse dataset for training, performing feature engineering to extract more useful features from the data, optimizing the hyperparameters of the classification algorithm, and using more advanced techniques such as ensemble learning.

4. **What if you have imbalanced classes in the dataset? How can you handle this?**
   If you have imbalanced classes in the dataset, you can handle this by using techniques such as undersampling or oversampling to balance the classes, or by using algorithms that are specifically designed to handle imbalanced datasets such as SMOTE (Synthetic Minority Over-sampling Technique).

5. **What if you want to deploy this model in a web application? How can you do this?**
   If you want to deploy this model in a web application, you can use a web framework such as Flask or Django to create a REST API that accepts input data and returns the model predictions. You would need to serialize the trained model to a file format such as pickle or joblib so that it can be loaded by the web application. You would also need to handle any security concerns such as user authentication and data encryption to protect the application from fraudsters.