# Analyzing financial data for fraud detection using the Enron email

## Problem Description :

The Enron scandal was one of the biggest corporate scandals in history that resulted in the company's bankruptcy and several top executives being indicted for fraud. In this project, we will be using the Enron email dataset to identify fraudulent activities by analyzing financial data.

**Objectives**: The main objective of this project is to detect fraudulent activities in the financial data using machine learning algorithms. Specifically, we will be working on the following tasks:

- Exploratory data analysis of the financial dataset

- Feature engineering to extract relevant features from the data

- Building a machine learning model to identify fraudulent activities

- Evaluating the performance of the model using appropriate metrics

**Dataset**: The dataset used for this project is the Enron Email Dataset, which consists of email messages and financial data of Enron employees. This dataset was made public after the Enron scandal and can be found on Kaggle. The financial dataset contains information on the company's earnings, stock prices, and executive compensation.

**Deliverables**: The deliverables for this project include the following:

- A report detailing the data analysis, feature engineering, and model building process

- A machine learning model capable of identifying fraudulent activities in the financial data

- A presentation summarizing the findings and recommendations for further action

Background Information: In the late 1990s, Enron was one of the largest energy companies in the world, with revenues of over $100 billion. However, the company's executives engaged in fraudulent activities, such as hiding debts and inflating earnings, to make the company appear more profitable than it was. This eventually led to the company's downfall and bankruptcy. The Enron scandal resulted in several top executives being indicted for fraud and other charges, including the CEO and CFO. The scandal led to new regulations being put in place to prevent similar corporate scandals in the future.

# Framework:

1.**Data Loading and Preprocessing**

•        Load the dataset into a pandas dataframe.

•        Preprocess the data by handling missing values, outliers, and irrelevant features.

2.**Feature Engineering**

•        Create new features from existing features to capture relevant information

•        Remove redundant or irrelevant features

3.**Exploratory Data Analysis**

•        Visualize the distribution of features

•        Identify correlations and patterns in the data

•        Conduct hypothesis testing to identify significant differences between groups.

4.**Model Selection**

•        Choose appropriate machine learning algorithms for the problem

•        Split the dataset into training and testing sets

•        Evaluate the performance of different algorithms using metrics such as accuracy, precision, recall, and F1-score.

5.**Hyperparameter Tuning and Model Optimization**

•        Fine-tune the hyperparameters of the chosen model using cross-validation

•        Optimize the model by selecting the best set of hyperparameters based on performance.

6.**Model Evaluation**

•        Evaluate the final model on the test set

•        Analyze the performance of the model using metrics such as confusion matrix, ROC curve, and AUC score.

7.**Deployment and Future Work**

•        Deploy the model into a production environment

- Continuously monitor the model's performance and update it as necessary


- Explore additional feature engineering techniques and machine learning algorithms to improve the model's performance.

# Code Explanation :

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

**Section 1: Importing Libraries**

In this section, we import the necessary libraries for the project such as pandas, numpy, matplotlib, seaborn, and sklearn.

**Section 2: Data Loading**

In this section, we load the Enron email dataset using the pandas library. The dataset contains information about email exchanges among Enron employees.

**Section 3: Data Cleaning and Preparation**

In this section, we clean the dataset by removing unnecessary columns and rows with missing values. We also engineer new features such as the number of emails sent and received by each employee.

**Section 4: Exploratory Data Analysis**

In this section, we analyze the cleaned dataset to gain insights into the data using visualizations and statistical methods. We explore the distribution of various features and look for patterns or anomalies that might indicate fraudulent behavior.

**Section 5: Feature Selection and Model Building**

In this section, we select the relevant features and build a machine learning model to detect fraud in the dataset. We use various algorithms such as Random Forest Classifier, Decision Tree Classifier, and Logistic Regression to build our models.

**Section 6: Model Evaluation and Fine-tuning**

In this section, we evaluate the performance of our models using metrics such as accuracy, precision, recall, and F1-score. We also fine-tune our models by adjusting hyperparameters to improve their performance.

# Future Work:

 Analyzing financial data for fraud detection using the Enron Email

This project has focused on exploring and analyzing the Enron Email dataset to identify potential fraudulent activities. However, there is still a lot of work that can be done to improve the accuracy and efficiency of the fraud detection system. Here are some potential future work that can be done:

**1.Feature Engineering**: Feature engineering is the process of selecting and transforming the most important features that contribute the most to the target variable. In this project, we have used some basic features such as email address, total payments, and total stock value. However, there are many other features that can be extracted from the Enron email dataset, such as email content, email frequency, and email recipient. These features can help to increase the accuracy of the model.

**2.Data preprocessing and cleaning:** Preprocessing and cleaning the data is a crucial step in any machine learning project. In this project, we have done some basic data cleaning such as removing NaN values and scaling the data. However, there are many other preprocessing techniques that can be applied, such as outlier detection and removal, data normalization, and data transformation.

3**.Model Selection and Tuning**: In this project, we have used a Random Forest Classifier for fraud detection. However, there are many other algorithms that can be used such as Decision Trees, SVM, and Neural Networks. Each algorithm has its own strengths and weaknesses, and selecting the right algorithm can greatly improve the accuracy of the model. Additionally, hyperparameter tuning can also improve the performance of the model.

**4.Building a Real-time Fraud Detection System:** The current implementation of the model requires manual input of data and model retraining. Building a real-time fraud detection system can help to identify fraudulent activities in real-time. This can be achieved by deploying the model on a cloud-based platform and integrating it with the email server.

**5.Testing on a larger dataset:** The Enron Email dataset is relatively small, with only 146 records. Testing the model on a larger dataset can help to improve its accuracy and generalizability. A larger dataset can be obtained by scraping email data from other companies or using publicly available email datasets.

**Step-by-Step Guide to Implement Future Work:**

1.      **Feature Engineering:**

•       Explore the dataset to identify potential new features that can contribute to the model's accuracy

•       Extract the new features from the dataset and preprocess them if necessary

•       Add the new features to the model and evaluate its performance

2**.**      **Data preprocessing and cleaning:**

•       Use outlier detection and removal techniques to remove any anomalies from the data

•       Normalize or transform the data to improve its distribution

•       Use feature scaling techniques to scale the data to a uniform range

3**.**      **Model Selection and Tuning:**

•       Evaluate the performance of various algorithms such as Decision Trees, SVM, and Neural Networks

•       Choose the best algorithm based on its performance and efficiency

- Tune the hyperparameters of the chosen algorithm to improve its accuracy

4**.** **Building a Real-time Fraud Detection System:**

- Deploy the model on a cloud-based platform such as AWS or Google Cloud

- Integrate the model with the email server to detect fraudulent activities in real-time

- Implement a notification system to alert the relevant personnel when fraudulent activities are detected

5. **Testing on a larger dataset:**

- Obtain a larger dataset by scraping email data from other companies or using publicly available email datasets

- Preprocess and clean the dataset as necessary

- Test the model on the larger dataset and evaluate its performance

By implementing these future work, the accuracy and efficiency of the fraud detection system can be greatly improved, leading to better identification and prevention of fraudulent activities.

# 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.  **What other feature engineering techniques could you use to improve the performance of the classification models?**
    Answer: Some additional feature engineering techniques include creating interaction variables, adding polynomial features, or including domain-specific features.

2.  **How would you tune the hyperparameters of the models to improve their performance?**
    Answer: We can use grid search, randomized search or Bayesian optimization to tune the hyperparameters of the models.

3.  **How would you evaluate the performance of the models on imbalanced data?**
    Answer: We can use metrics like precision, recall, F1-score, ROC-AUC score, and PR-AUC score to evaluate the performance of the models on imbalanced data. We can also use techniques like oversampling or undersampling to balance the data before training the models.

4.  **What are some other classification algorithms that could be used for fraud detection and how do they compare to the ones used in this project?**
    Answer: Some other classification algorithms that could be used for fraud detection include SVM, Naïve Bayes, Decision Trees, Random Forests, Gradient Boosting, and Neural Networks. The performance of these algorithms can vary depending on the nature of the data and the specific problem being addressed.

**5.      How would you deploy the model in a production environment?**
Answer: We can deploy the model in a production environment using a RESTful API, which can be called from a web or mobile application. We can also use cloud-based services like Amazon SageMaker or Microsoft Azure Machine Learning to deploy the model in a scalable and cost-effective manner.