

## **1. Data cleaning including missing values, outliers and multi-collinearity.**

1. Demographic information: Age, gender, location, and other demographic information can be relevant indicators of fraud.
2. Purchase behaviour: The amount and frequency of purchases, as well as the type of products or services being purchased, can indicate whether a customer is fraudulent.
3. Payment information: The type of payment method used, such as credit card or PayPal, can be a useful indicator of fraud. Additionally, the billing address, card expiration date, and CVV number can also provide information about the likelihood of fraud.
4. Device information: The type of device used for a transaction, such as a desktop computer or mobile phone, can be an indicator of fraud. Additionally, the IP address and geolocation information can also provide information about the likelihood of fraud.
5. Historical data: Fraudulent customers often have a pattern of behaviour that can be detected over time. By analyzing past transactions and behaviour, it is possible to identify patterns that are indicative of fraud.

## **2. Describe your fraud detection model in elaboration.**

Fraud detection is a challenging task that requires a combination of machine learning algorithms and domain expertise to identify fraudulent transactions. In my fraud detection model, I use both decision trees and random forests to identify fraudulent customers.

Decision trees are a type of machine learning algorithm that works by recursively splitting the data into subsets based on the values of the input features. The goal of the algorithm is to create a tree structure that separates fraudulent data from non-fraudulent data. At each node of the tree, a decision is made based on the value of a feature that provides the best separation of the data. The final prediction is made by following the path from the root of the tree to a leaf node.

Random forests are an extension of decision trees that create multiple trees and combine their predictions to make a final prediction. In a random forest, each tree is trained on a different subset of the data, and the features are selected randomly at each node. The final prediction is made by taking the average of the predictions made by each tree.

To implement my fraud detection model, I did the follow these steps:

1. Collect and pre-process the data: The first step would be to collect the data and pre-process it to remove missing values, outliers, and irrelevant features.
2. Split the data into training and test sets: The data would then be split into training and test sets to evaluate the performance of the model.
3. Train the decision tree and random forest models: The decision tree and random forest models would then be trained on the training data. The hyperparameters of the models would be tuned to achieve the best performance.
4. Evaluate the models on the test data: The performance of the models would then be evaluated on the test data using metrics such as accuracy, precision, recall, and F1-score.
5. Compare the results: The results of the decision tree and random forest models would then be compared to determine which model performs better.

6. Deploy the model: The final step would be to deploy the model in a production environment and monitor its performance regularly to ensure that it continues to perform well over time.

In conclusion, decision trees and random forests are powerful algorithms that can be used to detect fraud in a variety of applications. By combining their predictions, we can build a robust fraud detection system that is capable of identifying fraudulent transactions with high accuracy.

### **3. How did you select variables to be included in the model?**

Use the correlation heatmap and VIF data. Simply look for any traits that are highly connected with one another, and then discard any that are less correlated with the Fraud Attribute.

### **4. Demonstrate the performance of the model by using the best set of tools.**

To demonstrate the performance of the decision tree and random forest models, we can use several tools and techniques. Here are some of the most common tools and techniques used to evaluate the performance of these models:

1. Confusion Matrix: A confusion matrix is a table that summarizes the performance of a classifier by comparing the predicted labels to the true labels. The confusion matrix provides important metrics such as accuracy, precision, recall, and F1-score, which can be used to evaluate the performance of the models.
2. ROC Curve: A receiver operating characteristic (ROC) curve is a graphical representation of the performance of a classifier that plots the true positive rate against the false positive rate. The ROC curve can be used to visualize the trade-off between the sensitivity (recall) and specificity of the classifier, and it is often used to evaluate the performance of binary classification models.
3. Precision-Recall Curve: The precision-recall curve is a graphical representation of the performance of a classifier that plots precision against the recall. The precision-recall curve can be used to visualize the trade-off between the precision and recall of the classifier, and it is often used to evaluate the performance of imbalanced classification problems, such as fraud detection.
4. Cross-Validation: Cross-validation is a technique that involves dividing the data into multiple folds, training the model on each fold, and evaluating its performance on the remaining data. This technique can be used to assess the generalization performance of the models, and to tune the hyperparameters of the models to achieve the best performance.

By using these tools and techniques, we can compare the performance of the decision tree and random forest models and determine which model is better suited for the task of fraud detection. Additionally, these tools can also be used to visualize the results of the models and to identify the features that are most important for predicting fraudulent transactions.

## **5. What are the key factors that predict fraudulent customers?**

The key factors that predict fraudulent customers can vary depending on the specific problem and data set. However, some common factors that can contribute to fraudulent behaviour include:

1. Demographic information: Information about the customer's age, gender, income, and location can provide insight into the likelihood of fraudulent behaviour. For example, certain age groups or geographic locations may be more likely to engage in fraudulent activities.
2. Spending patterns: Changes in spending patterns, such as unusual or large purchases, can indicate fraudulent activity. For example, a sudden increase in spending on luxury items may indicate the use of a stolen credit card.
3. Payment method: The type of payment method used, such as a credit card or bank transfer, can also be a predictor of fraudulent behaviour. For example, certain payment methods may be more susceptible to fraud.
4. Account information: Information about the customer's account, such as the length of time the account has been open or the number of recent account changes, can also provide insight into the likelihood of fraudulent behaviour.
5. Interactions with the system: Information about the customer's interactions with the system, such as the frequency and type of transactions, can also be used to predict fraudulent behaviour. For example, frequent small transactions may indicate the use of a stolen credit card.
6. External sources: Information from external sources, such as social media or public records, can also be used to predict fraudulent behaviour. For example, a history of criminal behaviour may indicate a higher likelihood of fraudulent activity.

These are some of the common factors that can be used to predict fraudulent behaviour. However, the most important factors will depend on the specific problem and data set and may require domain expertise and feature engineering to identify.

## **6. Do these factors make sense? If yes, How? If not, How not?**

Yes, these factors make sense as predictors of fraudulent behaviour. The reasoning behind why each factor is relevant is mentioned above. In general, these factors make sense as predictors of fraudulent behaviour because they provide information about the customer's behaviour and background, which can be used to identify potential fraud. However, it is important to note that these factors may not be relevant or predictive in all cases, and the most important factors will depend on the specific problem and data set.

## **7. What kind of prevention should be adopted while the company update its infrastructure?**

When updating a company's infrastructure, it's important to consider the potential security implications and take steps to prevent fraud and other security incidents. Here are some key preventive measures that a company can adopt while updating its infrastructure:

1. Conduct a risk assessment: Before making any changes to the infrastructure, it is important to conduct a thorough risk assessment to identify potential security

weaknesses and threats. This will help to prioritize the necessary updates and ensure that the most critical issues are addressed first.

2. Implement security by design: When designing and implementing the updates to the infrastructure, it is important to incorporate security measures from the start. This includes using secure protocols and encryption, implementing access controls and authentication, and designing systems with resilience in mind to ensure that they can withstand potential security incidents.
3. Update software and hardware: Ensure that all software and hardware components are up-to-date and secure. This includes updating operating systems, firewalls, and other security software, as well as replacing any outdated hardware.
4. Train employees: Provide employees with training on the updated infrastructure, including how to use it securely, how to identify and report security incidents, and the company's policies and procedures for protecting sensitive information.
5. Regularly monitor and test systems: Regularly monitor and test the updated infrastructure to identify and address any potential security weaknesses. This can include conducting penetration tests, reviewing logs, and conducting regular security audits.
6. Develop a response plan: Develop a comprehensive response plan in case of a security incident, including procedures for reporting incidents, communicating with affected customers, and restoring systems and data.

By following these preventive measures, a company can ensure that its infrastructure is secure and protected against potential security incidents while updating its systems.

#### **8. Assuming these actions have been implemented, how would you determine if they work?**

To determine if the preventive measures implemented to secure a company's infrastructure are effective, you can use the following methods:

1. Vulnerability scans and penetration testing: Regularly conducting vulnerability scans and penetration tests can help identify any potential security weaknesses and vulnerabilities in the infrastructure. This information can then be used to evaluate the effectiveness of the security measures and make any necessary changes.
2. Incident response and reporting: Monitoring the number and types of security incidents and how they are handled can provide insight into the effectiveness of the response plan and the overall security posture of the infrastructure.
3. Compliance and security audits: Conducting regular security audits and compliance assessments can help ensure that the infrastructure is in compliance with industry standards and regulations and that the preventive measures are being effectively implemented.

4. User feedback and surveys: Collecting user feedback and conducting surveys can help gauge the overall level of security awareness and understanding among employees and identify any areas where additional training or education may be necessary.
5. Monitoring system logs: Regularly monitoring system logs can help detect any unusual or suspicious activity and provide insight into the effectiveness of the access controls and authentication measures.

By using these methods to evaluate the effectiveness of the preventive measures, a company can ensure that its infrastructure is secure and protected against potential security incidents. It is also important to regularly review and update the preventive measures as needed to keep up with changes in technology and the evolving threat landscape.