

# Actividad 13

Jose Manuel Enriquez

29 de marzo del 2025

## 1 Introduction

### 1.1 ¿Qué es un árbol de decisión?

Random Forest es un algoritmo de aprendizaje supervisado basado en la combinación de múltiples árboles de decisión para mejorar la precisión y reducir el sobreajuste. Funciona mediante un proceso llamado Bootstrap Aggregating (Bagging), donde se generan subconjuntos de datos aleatorios con reemplazo, y se entrena un árbol de decisión en cada uno. Además, en cada división del árbol, solo se considera un subconjunto aleatorio de características, lo que reduce la correlación entre los árboles y mejora la generalización del modelo. Para la clasificación, el resultado final se obtiene por votación mayoritaria, mientras que en regresión se toma el promedio de las predicciones. Sus principales aplicaciones incluyen la detección de fraudes en finanzas, el diagnóstico de enfermedades en medicina, los sistemas de recomendación en plataformas de contenido, la predicción del clima y la segmentación de imágenes en visión computacional. Entre sus ventajas destacan su capacidad para reducir el sobreajuste, su robustez ante datos ruidosos y su flexibilidad para tareas de clasificación y regresión. Sin embargo, puede ser computacionalmente costoso y menos interpretable que un árbol de decisión individual.

## 2 Metodología

Este ejercicio fue realizado en una Jupyter Notebook y Python 3.12.4 Lo primero que se hizo fue importar todas las librerías necesarias.

Primero se importan todas las librerías necesarias

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 from sklearn.metrics import confusion_matrix
7 from sklearn.metrics import classification_report
8 from sklearn.model_selection import train_test_split
```

```

9 from sklearn.linear_model import LogisticRegression
10 from sklearn.decomposition import PCA
11 from sklearn.tree import DecisionTreeClassifier
12
13 from pylab import rcParams
14
15 !pip install imbalanced-learn
16
17 from imblearn.under_sampling import NearMiss
18 from imblearn.over_sampling import RandomOverSampler
19 from imblearn.combine import SMOTETomek
20 from imblearn.ensemble import BalancedBaggingClassifier
21
22 from collections import Counter
23
24 #set up graphic style in this case I am using the color scheme from
   xkcd.com
25 rcParams['figure.figsize'] = 14, 8.7 # Golden Mean
26 LABELS = ["Normal", "Fraud"]
27 #col_list = ["cerulean", "scarlet"]# https://xkcd.com/color/rgb/
28 #sns.set(style='white', font_scale=1.75, palette=sns.xkcd_palette(
   col_list))
29
30 %matplotlib inline

```

Luego importamos los datos

```

1 df = pd.read_csv("creditcard.csv")
2 df.head(n=5)

```

Comparamos los datos que tienen 0 y 1

```

1 pd.value_counts(df['Class'], sort = True) #class comparison 0=
   Normal 1=Fraud

```

Creamos dos dataframes para los dos casos

```

1 normal_df = df[df.Class == 0] #registros normales
2 fraud_df = df[df.Class == 1] #casos de fraude

```

Se crean las variables de los datos y hacemos la primera predicción

```

1 y = df['Class']
2 X = df.drop('Class', axis=1)
3 X_train, X_test, y_train, y_test = train_test_split(X, y,
   train_size=0.7)
4
5 def mostrar_resultados(y_test, pred_y):
6     conf_matrix = confusion_matrix(y_test, pred_y)
7     plt.figure(figsize=(8, 8))
8     sns.heatmap(conf_matrix, xticklabels=LABELS, yticklabels=LABELS
   , annot=True, fmt="d");
9     plt.title("Confusion matrix")
10    plt.ylabel('True class')
11    plt.xlabel('Predicted class')
12    plt.show()
13    print(classification_report(y_test, pred_y))

```

```

14
15 def run_model_balanced(X_train, X_test, y_train, y_test):
16     clf = LogisticRegression(C=1.0,penalty='l2',random_state=1,
17                             solver="newton-cg",class_weight="balanced")
18     clf.fit(X_train, y_train)
19     return clf
20
21 model = run_model_balanced(X_train, X_test, y_train, y_test)
22
23 pred_y = model.predict(X_test)
24 mostrar_resultados(y_test, pred_y)

```

Y se obtuvo la siguiente gráfica:

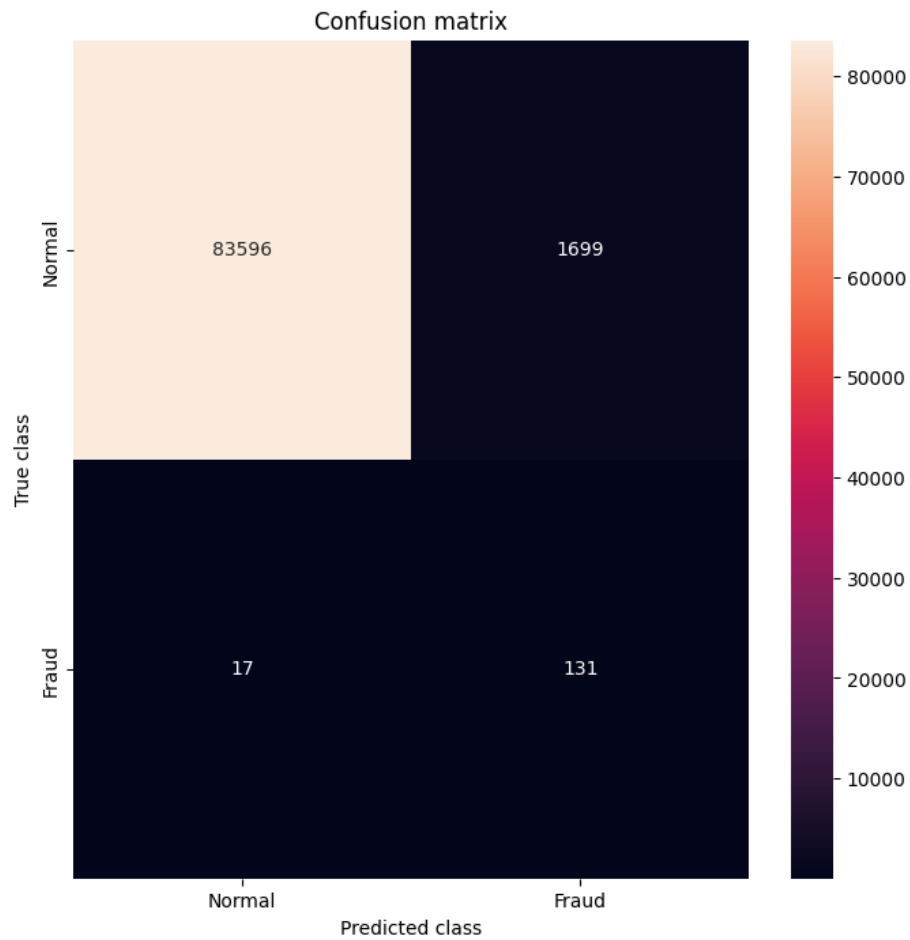


Figure 1: Confussion Matrix 1

Luego se crea un modelo con 100 arboles

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 # Crear el modelo con 100 arboles
4 model = RandomForestClassifier(n_estimators=100,
5                               bootstrap = True, verbose=2,
6                               max_features = 'sqrt')
7 # entrenar!
8 model.fit(X_train, y_train)
9
10 pred_y = model.predict(X_test)
11 mostrar_resultados(y_test, pred_y)
```

Y ahora se obtuvo la siguiente gráfica

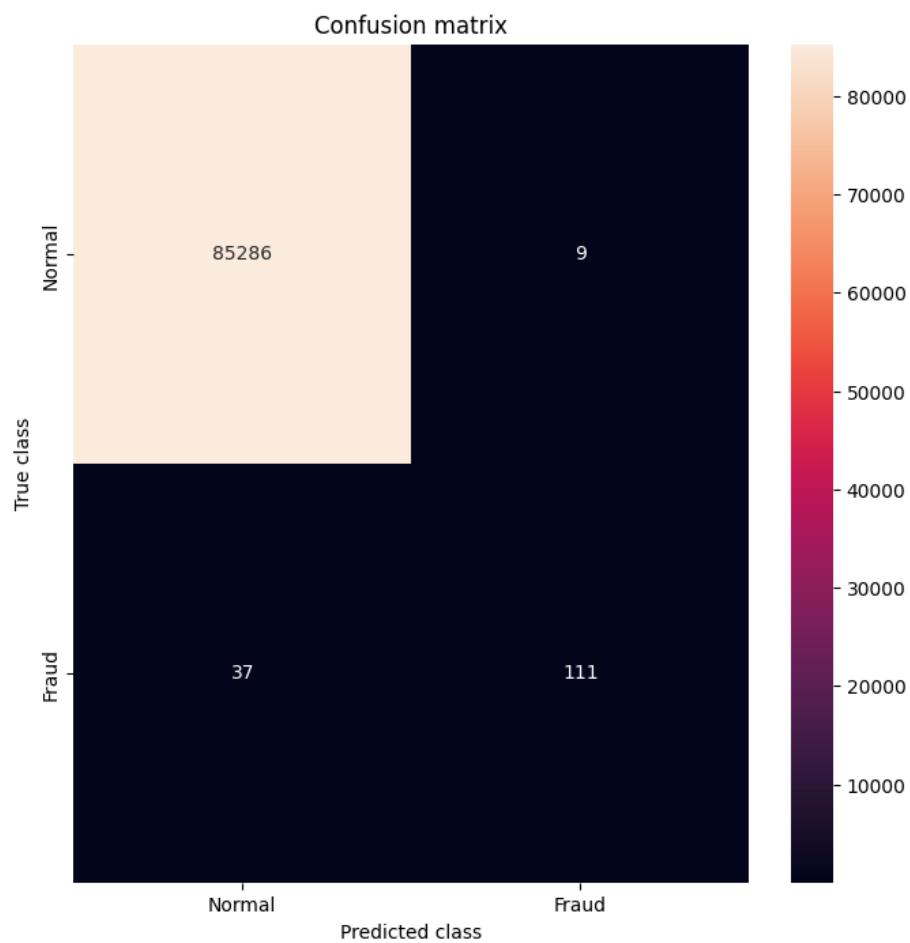


Figure 2: Confussion Matrix 2

### **3 Resultados**

Se obtuvo que la calificación de ROC AUC es de 0.95

### **4 Conclusiones**

Me pareció muy interesante realizar predicciones con random forest con Python ya que es de las primeras veces en las cuales implemento esta herramienta.