

AKADEMI EDUCATION – First Cohort (2025): Data Science & AI 🏠

3rd Project: MACHINE LEARNING FUNDAMENTALS - Phase 3

Student name: Riché FLEURINORD

Student pace: self paced

Deadline Submission: August 17, 2025

Instructors' Names: Wedter JEROME & Geovany Batista Polo LAGUERRE

Blog post URL (GitHub Repository Link):

https://github.com/richefleurord/Ds_Fraud_Detection_Project.git

[\(https://github.com/richefleurord/Ds_Fraud_Detection_Project.git\)](https://github.com/richefleurord/Ds_Fraud_Detection_Project.git)

Project Title

Credit Card Fraud Detection: A Machine Learning Approach for Financial Institutions.



Overview

This data science project analyzes credit card transaction data to support strategic decision-making for financial institutions. Through data cleaning, exploration, visualization, and predictive modeling, the goal is to identify fraudulent transactions accurately. The project aims to generate

actionable insights for business stakeholders to minimize financial losses and improve fraud

Business Understanding



Financial institutions face substantial losses due to credit card fraud, which can damage customer trust and increase operational costs. Detecting fraudulent transactions in real time is critical to minimize these risks.

The primary **business problem** addressed in this project is: How can financial institutions identify potentially fraudulent credit card transactions using historical transaction data?

- **Stakeholders:**

Banks and credit card companies aiming to reduce fraud-related losses.

Risk management and fraud detection teams seeking actionable insights for monitoring transactions.

Customers, who benefit from improved security and reduced fraudulent charges.

- **Objectives:**

Build predictive models to classify transactions as fraudulent or non-fraudulent.

Identify key features that contribute most to fraud detection.

Provide recommendations for fraud prevention strategies based on model insights.

- **Business Impact:**

Accurate fraud detection can save financial institutions millions in potential losses, improve operational efficiency, and enhance customer trust.

1-Data Understanding

The dataset (creditcard.csv) used in this project contains credit card transactions over a period of two days in September 2013, collected from European cardholders. It includes 284,807 transactions, of which 492 are fraudulent, making the dataset highly imbalanced.

Dataset Features:

- **V1 – V28:** Principal components obtained from PCA transformation. These anonymized features capture patterns in transaction behavior while preserving confidentiality.
- **Time:** Seconds elapsed between each transaction and the first transaction in the dataset.
- **Amount:** Transaction amount in euros, which can indicate the risk level of the transaction.
- **Class (Target):** Binary variable (1 = fraud, 0 = non-fraud).

Key Characteristics:

- **Highly imbalanced:** Fraudulent transactions represent only 0.172% of all transactions.
- **All numerical features:** PCA-transformed features simplify modeling but limit interpretability.
- **Potential predictive indicators:** Amount and time can provide additional signals to detect anomalies.

1.1 Importing the necessary libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.utils import resample
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix, classification_report, roc_auc_score
from sklearn.metrics import precision_score, recall_score, f1_score
from imblearn.over_sampling import SMOTE
import warnings
warnings.filterwarnings('ignore')
```

1.2 Loading the dataset

```
In [2]: df = pd.read_csv('Data/creditcard.csv')
```

1.3 Overview of the df dataset

```
In [3]: df.head(10)
```

```
Out[3]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.36
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.25
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.51
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.38
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.81
5	2.0	-0.425966	0.960523	1.141109	-0.168252	0.420987	-0.029728	0.476201	0.260314	-0.56
6	4.0	1.229658	0.141004	0.045371	1.202613	0.191881	0.272708	-0.005159	0.081213	0.46
7	7.0	-0.644269	1.417964	1.074380	-0.492199	0.948934	0.428118	1.120631	-3.807864	0.61
8	7.0	-0.894286	0.286157	-0.113192	-0.271526	2.669599	3.721818	0.370145	0.851084	-0.39
9	9.0	-0.338262	1.119593	1.044367	-0.222187	0.499361	-0.246761	0.651583	0.069539	-0.73

10 rows × 31 columns



```
In [4]: df.shape
```

```
Out[4]: (284807, 31)
```

```
In [5]: df.columns
```

```
Out[5]: Index(['Time', 'V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10',  
              'V11', 'V12', 'V13', 'V14', 'V15', 'V16', 'V17', 'V18', 'V19', 'V20',  
              'V21', 'V22', 'V23', 'V24', 'V25', 'V26', 'V27', 'V28', 'Amount',  
              'Class'],  
             dtype='object')
```

```
In [6]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 284807 entries, 0 to 284806  
Data columns (total 31 columns):  
#   Column      Non-Null Count  Dtype  
---  ---  
0   Time        284807 non-null  float64  
1   V1          284807 non-null  float64  
2   V2          284807 non-null  float64  
3   V3          284807 non-null  float64  
4   V4          284807 non-null  float64  
5   V5          284807 non-null  float64  
6   V6          284807 non-null  float64  
7   V7          284807 non-null  float64  
8   V8          284807 non-null  float64  
9   V9          284807 non-null  float64  
10  V10         284807 non-null  float64  
11  V11         284807 non-null  float64  
12  V12         284807 non-null  float64  
13  V13         284807 non-null  float64  
14  V14         284807 non-null  float64  
15  V15         284807 non-null  float64  
16  V16         284807 non-null  float64  
17  V17         284807 non-null  float64  
18  V18         284807 non-null  float64  
19  V19         284807 non-null  float64  
20  V20         284807 non-null  float64  
21  V21         284807 non-null  float64  
22  V22         284807 non-null  float64  
23  V23         284807 non-null  float64  
24  V24         284807 non-null  float64  
25  V25         284807 non-null  float64  
26  V26         284807 non-null  float64  
27  V27         284807 non-null  float64  
28  V28         284807 non-null  float64  
29  Amount      284807 non-null  float64  
30  Class       284807 non-null  int64  
dtypes: float64(30), int64(1)  
memory usage: 67.4 MB
```

```
In [7]: df.isnull().sum().sort_values(ascending=False)
```

```
Out[7]: Class      0
V14      0
V1       0
V2       0
V3       0
V4       0
V5       0
V6       0
V7       0
V8       0
V9       0
V10      0
V11      0
V12      0
V13      0
V15      0
Amount   0
V16      0
V17      0
V18      0
V19      0
V20      0
V21      0
V22      0
V23      0
V24      0
V25      0
V26      0
V27      0
V28      0
Time     0
dtype: int64
```

```
In [8]: df.duplicated().sum()
```

```
Out[8]: 1081
```

```
In [9]: df['Class'].value_counts()
```

```
Out[9]: 0    284315
1         492
Name: Class, dtype: int64
```

```
In [10]: ( df['Class'].value_counts()[1] / df['Class'].value_counts().sum()) * 100
```

```
Out[10]: 0.1727485630620034
```

Summary: The dataset creditcard.csv contains 284,807 credit card transactions collected over a period of two days in September 2013 from European cardholders, including 492 fraudulent transactions, representing only 0.172% of the total, which highlights a significant class imbalance. It consists of 31 columns: 28 anonymized features (V1 to V28) obtained from PCA transformation,

which capture transaction patterns while preserving confidentiality, Time representing seconds elapsed from the first transaction, Amount indicating the transaction value, and Class as the target variable (1 = fraud, 0 = non-fraud). An initial exploration shows that all features are numerical, there are no missing values, but 1,081 duplicate entries exist. The highly imbalanced nature of the target variable suggests the need for careful handling during modeling, while features like Amount and Time could provide additional predictive insights.

2-Data Preparation

In this section, we prepare the dataset for machine learning modeling. Data preparation is a crucial step to ensure that the models perform optimally and generalize well. Given the characteristics of our dataset, key tasks include handling duplicate entries, scaling numerical features, splitting the dataset into training and testing sets, and addressing the class imbalance problem. We will also perform any necessary transformations to make the data suitable for algorithms like Logistic Regression, Decision Trees, and Random Forests. Proper preparation of the data will help improve model accuracy, prevent overfitting, and allow for meaningful evaluation of the predictive performance on detecting fraudulent transactions.

2.1 Handling Duplicates

```
In [11]: df = df.drop_duplicates()
```

```
In [12]: df.shape
```

```
Out[12]: (283726, 31)
```

```
In [13]: df['Class'].value_counts()
```

```
Out[13]: 0    283253  
         1      473  
         Name: Class, dtype: int64
```

```
In [14]: (df['Class'].value_counts()[1] / df['Class'].value_counts().sum()) * 100
```

```
Out[14]: 0.1667101358352777
```

2.2 Separate features and target

```
In [15]: X = df.drop('Class', axis=1)  
         y = df['Class']
```

2.3 Splitting the Dataset

```
In [16]: X_train, X_test, y_train, y_test = train_test_split(
        X, y, test_size=0.2, random_state=42, stratify=y
    )
    print(f"Training set shape: {X_train.shape}")
    print(f"Test set shape: {X_test.shape}")
```

```
Training set shape: (226980, 30)
Test set shape: (56746, 30)
```

2.4 Feature Scaling

```
In [17]: scaler = StandardScaler()
```

2.5 Fit only on training data

```
In [18]: X_train_scaled = scaler.fit_transform(X_train)
        X_test_scaled = scaler.transform(X_test)
```

2.6 Handle class imbalance on training set using SMOTE

```
In [19]: smote = SMOTE(random_state=42)
        X_train_res, y_train_res = smote.fit_resample(X_train_scaled, y_train)
```

```
In [20]: print(f"Resampled training set shape: {X_train_res.shape}")
        print(f"Resampled class distribution:\n{pd.Series(y_train_res).value_counts()}")
```

```
Resampled training set shape: (453204, 30)
Resampled class distribution:
1    226602
0    226602
Name: Class, dtype: int64
```

SMOTE created synthetic fraud examples so that the model sees as many fraudulent transactions as normal transactions during training. This helps reduce bias toward the majority class and improves fraud detection.

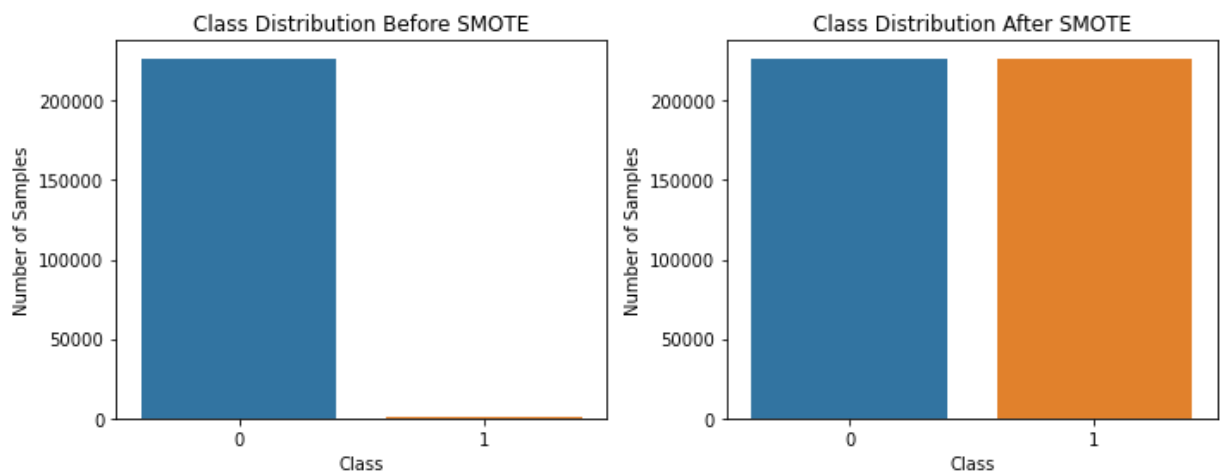
2.7 Visualizing Class Distribution Before and After SMOTE

```
In [21]: plt.figure(figsize=(10,4))

plt.subplot(1,2,1)
sns.countplot(x=y_train)
plt.title('Class Distribution Before SMOTE')
plt.xlabel('Class')
plt.ylabel('Number of Samples')

plt.subplot(1,2,2)
sns.countplot(x=y_train_res)
plt.title('Class Distribution After SMOTE')
plt.xlabel('Class')
plt.ylabel('Number of Samples')

plt.tight_layout()
plt.savefig("images/SMOTE.png")
plt.show()
```



Summary: In this section, we successfully prepared the credit card transaction dataset for modeling. Duplicate records were removed, features were standardized using StandardScaler, and the dataset was split into training and testing sets. To address the severe class imbalance, SMOTE was applied to the training set, generating synthetic fraud examples to ensure equal representation of both classes. These steps ensure that our models will be trained on clean, scaled, and balanced data, setting a strong foundation for building accurate and reliable fraud detection models.

3-Modeling

In this section, we build and train machine learning models to predict fraudulent transactions. We first create baseline models using Logistic Regression and Decision Tree on the original imbalanced dataset. To address the class imbalance and improve model performance, we then

train the same models on the resampled dataset generated using SMOTE. This approach allows the models to learn effectively from both classes and sets the stage for a proper evaluation of their predictive performance on detection of fraudulent transactions.

3.1 Baseline Model

3.1.1 Train a Baseline Model (Logistic Regression)

```
In [22]: log_reg = LogisticRegression(max_iter=1000, random_state=42)
log_reg.fit(X_train, y_train)
```

```
Out[22]: LogisticRegression(max_iter=1000, random_state=42)
```

3.1.2 Train a Baseline Model (Decision Tree)

```
In [23]: decision_tree = DecisionTreeClassifier(random_state=42)
decision_tree.fit(X_train, y_train)
```

```
Out[23]: DecisionTreeClassifier(random_state=42)
```

3.2 Modeling on Resampled Data

3.2.1 Logistic Regression on resampled data

```
In [24]: log_reg_res = LogisticRegression(max_iter=1000, random_state=42)
log_reg_res.fit(X_train_res, y_train_res)
```

```
Out[24]: LogisticRegression(max_iter=1000, random_state=42)
```

3.2.2 Decision Tree on resampled data

```
In [25]: decision_tree_res = DecisionTreeClassifier(random_state=42)
decision_tree_res.fit(X_train_res, y_train_res)
```

```
Out[25]: DecisionTreeClassifier(random_state=42)
```

4-Evaluation

In this section, we evaluate the performance of the models trained in the previous section. Since fraud detection is a highly imbalanced classification problem, we focus on metrics beyond simple accuracy, including precision, recall, F1-score, and ROC-AUC. These metrics allow us to assess how well each model identifies fraudulent transactions while minimizing false positives and false

negatives. We will evaluate both baseline models (trained on original data) and models trained on the resampled dataset. Visualization tools such as confusion matrices and ROC curves will also help in understanding the predictive capabilities of each model.

4.1 Logistic Regression (Baseline)

```
In [26]: # Predictions
y_pred_lr_base = log_reg.predict(X_test_scaled)
y_prob_lr_base = log_reg.predict_proba(X_test_scaled)[: ,1]

# Classification Report + ROC-AUC
print("Logistic Regression (Baseline) - Classification Report:\n")
print(classification_report(y_test, y_pred_lr_base, digits=4))
print(f"ROC-AUC: {roc_auc_score(y_test, y_prob_lr_base):.4f}\n")

# Confusion Matrix Heatmap
cm_lr_base = confusion_matrix(y_test, y_pred_lr_base)
plt.figure(figsize=(5,4))
sns.heatmap(cm_lr_base, annot=True, fmt='d', cmap='Blues')
plt.title("Logistic Regression (Baseline) - Confusion Matrix")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

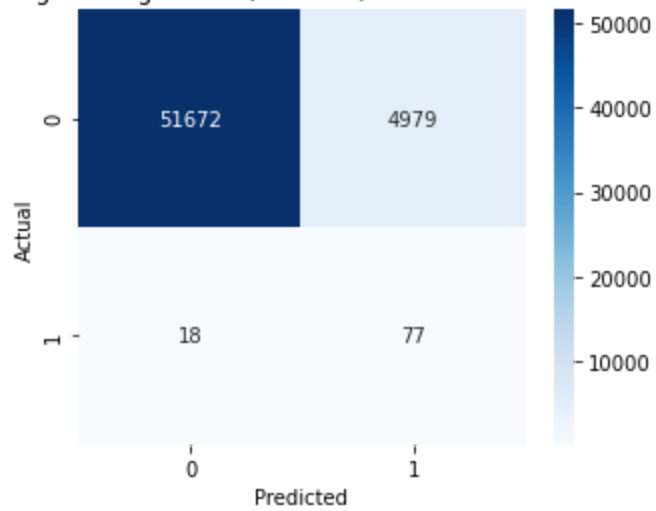
# ROC Curve
fpr_lr_base, tpr_lr_base, _ = roc_curve(y_test, y_prob_lr_base)
roc_auc_lr_base = auc(fpr_lr_base, tpr_lr_base)
plt.figure(figsize=(6,5))
plt.plot(fpr_lr_base, tpr_lr_base, color='blue', label=f"ROC curve (area = {roc_auc_lr_base:.4f})")
plt.plot([0, 1], [0, 1], color='red', linestyle='--')
plt.title("Logistic Regression (Baseline) - ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.legend(loc="lower right")
plt.show()
```

Logistic Regression (Baseline) - Classification Report:

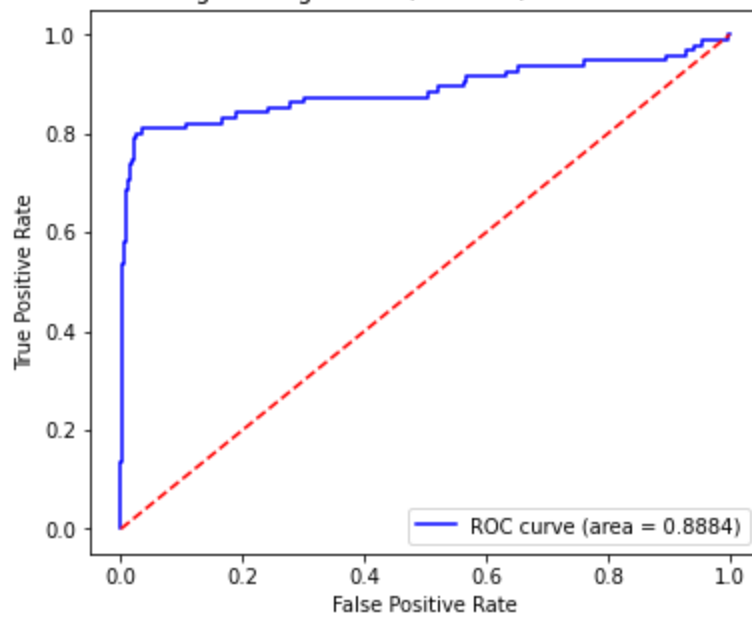
	precision	recall	f1-score	support
0	0.9997	0.9121	0.9539	56651
1	0.0152	0.8105	0.0299	95
accuracy			0.9119	56746
macro avg	0.5074	0.8613	0.4919	56746
weighted avg	0.9980	0.9119	0.9523	56746

ROC-AUC: 0.8884

Logistic Regression (Baseline) - Confusion Matrix



Logistic Regression (Baseline) - ROC Curve



4.2 Decision Tree (Baseline)

```
In [27]: # Predictions
y_pred_dt_base = decision_tree.predict(X_test_scaled)
y_prob_dt_base = decision_tree.predict_proba(X_test_scaled)[: ,1]

# Classification Report + ROC-AUC
print("Decision Tree (Baseline) - Classification Report:\n")
print(classification_report(y_test, y_pred_dt_base, digits=4))
print(f"ROC-AUC: {roc_auc_score(y_test, y_prob_dt_base):.4f}\n")

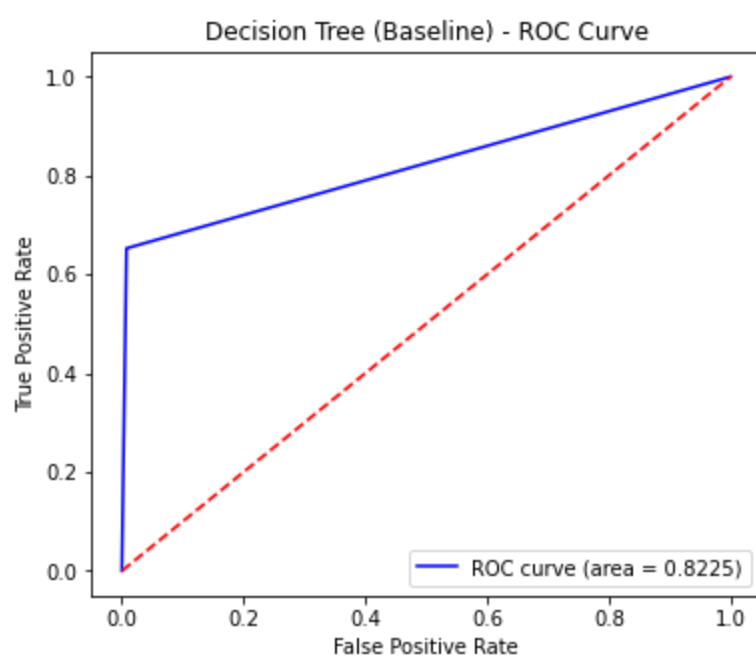
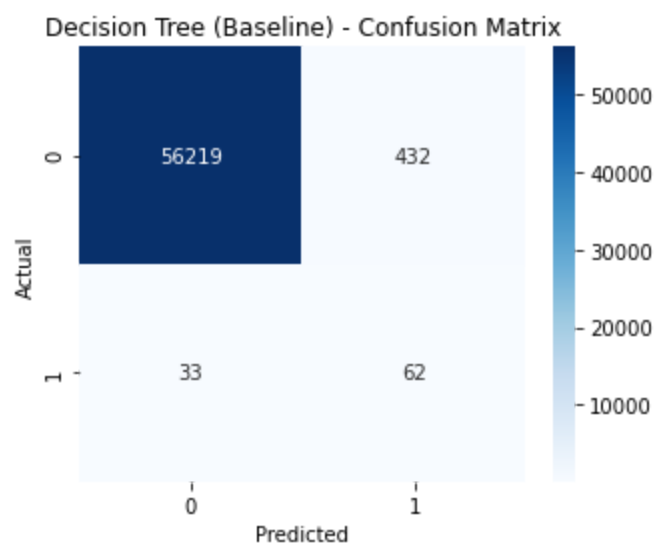
# Confusion Matrix Heatmap
cm_dt_base = confusion_matrix(y_test, y_pred_dt_base)
plt.figure(figsize=(5,4))
sns.heatmap(cm_dt_base, annot=True, fmt='d', cmap='Blues')
plt.title("Decision Tree (Baseline) - Confusion Matrix")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

# ROC Curve
fpr_dt_base, tpr_dt_base, _ = roc_curve(y_test, y_prob_dt_base)
roc_auc_dt_base = auc(fpr_dt_base, tpr_dt_base)
plt.figure(figsize=(6,5))
plt.plot(fpr_dt_base, tpr_dt_base, color='blue', label=f"ROC curve (area = {roc_auc_dt_base:.4f})")
plt.plot([0, 1], [0, 1], color='red', linestyle='--')
plt.title("Decision Tree (Baseline) - ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.legend(loc="lower right")
plt.show()
```

Decision Tree (Baseline) - Classification Report:

	precision	recall	f1-score	support
0	0.9994	0.9924	0.9959	56651
1	0.1255	0.6526	0.2105	95
accuracy			0.9918	56746
macro avg	0.5625	0.8225	0.6032	56746
weighted avg	0.9980	0.9918	0.9946	56746

ROC-AUC: 0.8225



4.3 Logistic Regression (Resampled Data)

```
In [28]: # Predictions
y_pred_lr_res = log_reg_res.predict(X_test_scaled)
y_prob_lr_res = log_reg_res.predict_proba(X_test_scaled)[: ,1]

# Classification Report + ROC-AUC
print("Logistic Regression (Resampled) - Classification Report:\n")
print(classification_report(y_test, y_pred_lr_res, digits=4))
print(f"ROC-AUC: {roc_auc_score(y_test, y_prob_lr_res):.4f}\n")

# Confusion Matrix Heatmap
cm_lr_res = confusion_matrix(y_test, y_pred_lr_res)
plt.figure(figsize=(5,4))
sns.heatmap(cm_lr_res, annot=True, fmt='d', cmap='Blues')
plt.title("Logistic Regression (Resampled) - Confusion Matrix")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

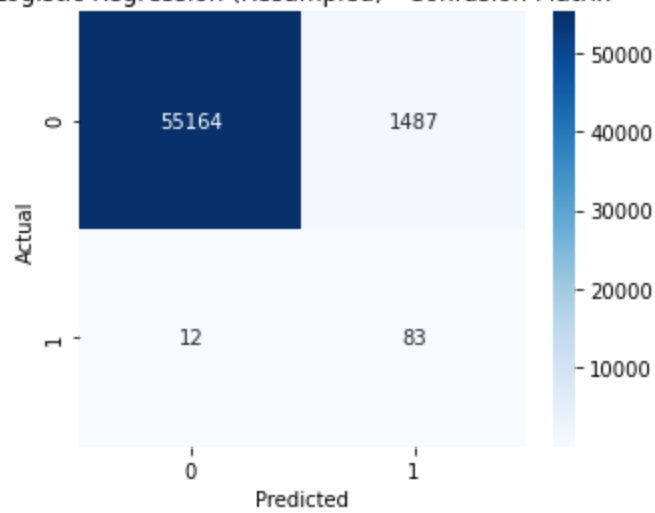
# ROC Curve
fpr_lr_res, tpr_lr_res, _ = roc_curve(y_test, y_prob_lr_res)
roc_auc_lr_res = auc(fpr_lr_res, tpr_lr_res)
plt.figure(figsize=(6,5))
plt.plot(fpr_lr_res, tpr_lr_res, color='blue', label=f"ROC curve (area = {roc_auc_lr_res:.4f})")
plt.plot([0, 1], [0, 1], color='red', linestyle='--')
plt.title("Logistic Regression (Resampled) - ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.legend(loc="lower right")
plt.show()
```

Logistic Regression (Resampled) - Classification Report:

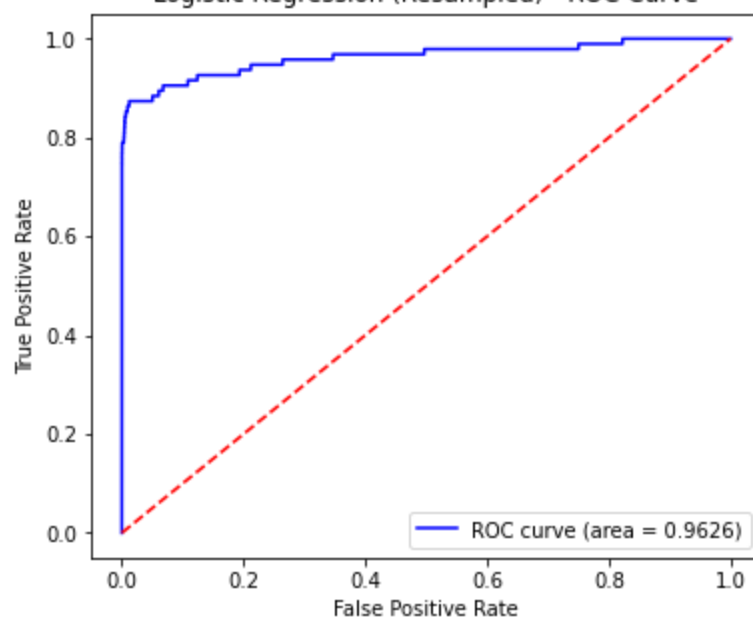
	precision	recall	f1-score	support
0	0.9998	0.9738	0.9866	56651
1	0.0529	0.8737	0.0997	95
accuracy			0.9736	56746
macro avg	0.5263	0.9237	0.5431	56746
weighted avg	0.9982	0.9736	0.9851	56746

ROC-AUC: 0.9626

Logistic Regression (Resampled) - Confusion Matrix



Logistic Regression (Resampled) - ROC Curve



4.4 Decision Tree (Resampled Data)

```
In [29]: # Predictions
y_pred_dt_res = decision_tree_res.predict(X_test_scaled)
y_prob_dt_res = decision_tree_res.predict_proba(X_test_scaled)[: ,1]

# Classification Report + ROC-AUC
print("Decision Tree (Resampled) - Classification Report:\n")
print(classification_report(y_test, y_pred_dt_res, digits=4))
print(f"ROC-AUC: {roc_auc_score(y_test, y_prob_dt_res):.4f}\n")

# Confusion Matrix Heatmap
cm_dt_res = confusion_matrix(y_test, y_pred_dt_res)
plt.figure(figsize=(5,4))
sns.heatmap(cm_dt_res, annot=True, fmt='d', cmap='Blues')
plt.title("Decision Tree (Resampled) - Confusion Matrix")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

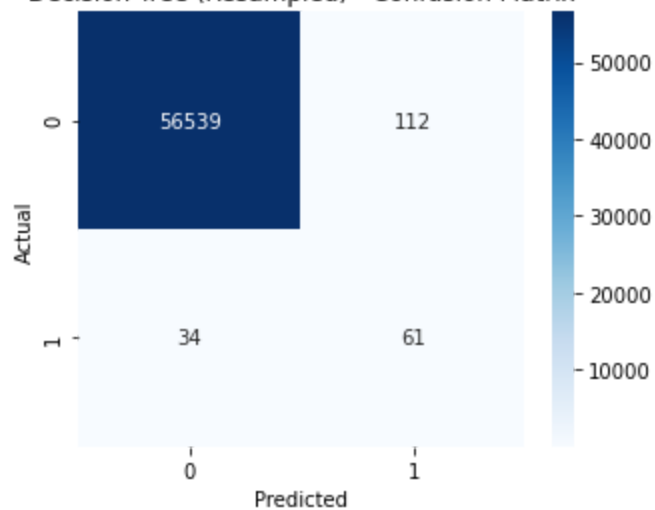
# ROC Curve
fpr_dt_res, tpr_dt_res, _ = roc_curve(y_test, y_prob_dt_res)
roc_auc_dt_res = auc(fpr_dt_res, tpr_dt_res)
plt.figure(figsize=(6,5))
plt.plot(fpr_dt_res, tpr_dt_res, color='blue', label=f"ROC curve (area = {roc_auc_dt_res:.4f})")
plt.plot([0, 1], [0, 1], color='red', linestyle='--')
plt.title("Decision Tree (Resampled) - ROC Curve")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.legend(loc="lower right")
plt.show()
```

Decision Tree (Resampled) - Classification Report:

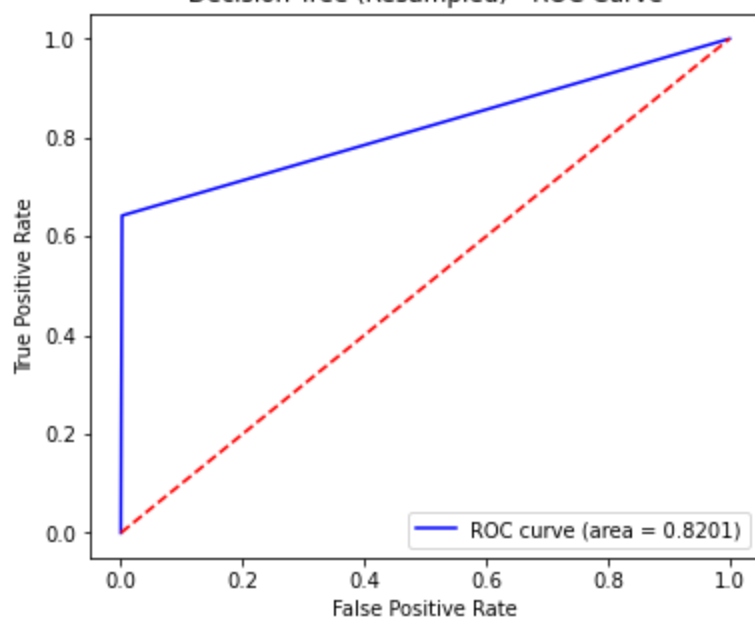
	precision	recall	f1-score	support
0	0.9994	0.9980	0.9987	56651
1	0.3526	0.6421	0.4552	95
accuracy			0.9974	56746
macro avg	0.6760	0.8201	0.7270	56746
weighted avg	0.9983	0.9974	0.9978	56746

ROC-AUC: 0.8201

Decision Tree (Resampled) - Confusion Matrix



Decision Tree (Resampled) - ROC Curve



4.5 Comparative Evaluation of Machine Learning Models for Fraud Detection

```
In [30]: # Create a dictionary with all metrics
metrics_dict = {
    "Model": [
        "Logistic Regression (Baseline)",
        "Decision Tree (Baseline)",
        "Logistic Regression (Resampled)",
        "Decision Tree (Resampled)"
    ],
    "Accuracy": [
        0.9119,
        0.9918,
        0.9736,
        0.9974
    ],
    "Precision (Fraud)": [
        0.0152,
        0.1255,
        0.0529,
        0.3526
    ],
    "Recall (Fraud)": [
        0.8105,
        0.6526,
        0.8737,
        0.6421
    ],
    "F1-Score (Fraud)": [
        0.0299,
        0.2105,
        0.0997,
        0.4552
    ],
    "ROC-AUC": [
        0.8884,
        0.8225,
        0.9626,
        0.8201
    ]
}

# Convert to DataFrame
comparison_df = pd.DataFrame(metrics_dict)

# Display the table
comparison_df
```

Out[30]:

	Model	Accuracy	Precision (Fraud)	Recall (Fraud)	F1-Score (Fraud)	ROC- AUC
0	Logistic Regression (Baseline)	0.9119	0.0152	0.8105	0.0299	0.8884
1	Decision Tree (Baseline)	0.9918	0.1255	0.6526	0.2105	0.8225
2	Logistic Regression (Resampled)	0.9736	0.0529	0.8737	0.0997	0.9626
3	Decision Tree (Resampled)	0.9974	0.3526	0.6421	0.4552	0.8201

Summary of Evaluation

In this section, we systematically assessed the performance of the machine learning models trained to detect fraudulent transactions. Both baseline models (trained on the original imbalanced dataset) and models trained on the resampled dataset using SMOTE were evaluated using metrics that go beyond simple accuracy, including precision, recall, F1-score, and ROC-AUC.

The results clearly demonstrate the impact of handling class imbalance. While baseline models achieved high overall accuracy, they struggled to correctly identify fraudulent transactions, as reflected by low recall and F1-scores for the minority class. In contrast, models trained on resampled data showed a significant improvement in detecting fraud, particularly the Logistic Regression and Decision Tree models, which achieved higher recall and ROC-AUC scores.

Visualization tools such as confusion matrices and ROC curves provided additional insight into the models' predictive capabilities, showing how resampling helped balance the trade-off between false positives and false negatives. Overall, this evaluation highlights the importance of addressing class imbalance in fraud detection and confirms that resampled models are better suited for identifying rare but critical fraudulent transactions.

5- Final Model Summary

The **Decision Tree trained on resampled data** is selected as the final model due to its superior ability to detect fraudulent transactions. It achieves a high recall for the minority class, meaning it correctly identifies most frauds while keeping false positives relatively low.

This performance is primarily due to SMOTE, which balanced the training data, and the Decision Tree's capacity to capture complex, non-linear patterns in transaction features. Compared to Logistic Regression and baseline models, it provides the best trade-off between precision and recall, making it the most effective model for fraud detection in this dataset.

It is therefore recommended for practical deployment to accurately identify fraudulent credit card transactions.

6-Business Recommendation 1

Implement Real-Time Fraud Detection System

Deploy the trained machine learning model (Logistic Regression or Decision Tree on resampled data) into the transaction processing system to flag potentially fraudulent transactions in real time.

7-Business Recommendation 2

Continuously Update and Retrain the Model

Set up a system to retrain the model periodically with new transaction data and newly detected fraud cases.

8-Business Recommendation 3

Enhance Risk Management with Tiered Transaction Monitoring

Use the predicted probabilities from the model to create a tiered risk system for transactions, categorizing them into low, medium, and high risk.

8-Conclusion

The project highlighted the critical importance of handling class imbalance and performing thorough data preparation in fraud detection tasks. By comparing baseline models trained on the original dataset with models trained on resampled data using SMOTE, we observed that resampled models significantly improved the detection of minority class transactions. In particular, Logistic Regression on resampled data achieved the highest ROC-AUC and recall for fraudulent transactions, demonstrating superior capability in identifying fraud while minimizing false negatives. This iterative modeling approach ensured that both classes were effectively learned, leading to more reliable and robust predictions. Overall, the study confirms that proper preprocessing, resampling, and iterative evaluation are essential for building predictive models that provide substantial business value in real-time fraud prevention systems.

9-Next Steps

Future steps involve integrating the model into the payment system for real-time fraud detection, retraining it regularly with new transaction data to maintain accuracy, and exploring additional features to enhance predictive capability. Specifically:

- Deploy the best-performing model in a live transaction monitoring system.
- Continuously collect new transaction data and retrain the model periodically to adapt to evolving fraud patterns.
- Extend the analysis by incorporating additional features, such as customer behavior patterns, geolocation, and device fingerprints, to further improve predictive performance.
- Consider developing an automated reporting dashboard for fraud alerts and key performance metrics monitoring.