

► Haz clic aquí para obtener una pista

In [29]: # Cuenta el número de transacciones fraudulentas y no fraudulentes

plt.xticks([0, 1], ['0: No Fraudulentas', '1: Fraudulentas'])

0: No Fraudulentas

transacciones_fraud = data_limpia[data_limpia['Class'] == 1]

500

Desarrollo y evaluación de modelos

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import StandardScaler

Tamaño del conjunto de entrenamiento (X train): (226980, 30)

Tamaño del conjunto de evaluación (X_test): (56746, 30) Tamaño del conjunto de entrenamiento (y_train): (226980,)

Tamaño del conjunto de evaluación (y_test): (56746,)

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import classification_report, accuracy_score

modelo = RandomForestClassifier(max_depth=150, random_state=42)

print(f"La exactitud del modelo es: {exactitud * 100:.2f}%")

1.00

0.73

0.87

1.00

model_SVM = SVC(kernel='rbf', C=1, gamma='scale', probability=True)

Almaceno las predicciones sobre los datos de evaluación en SVM_pred

print(f"La exactitud del modelo es: {exactitud_SVM * 100:.2f}%")

Resumen de rendimiento modelo SVM Máquina de Vectores de Soporte:

1.00

0.64

0.82

1.00

analizar los resultados obtenidos por ambos modelos y extraer conclusiones.

1. Resultados del Modelo RandomForest

recall f1-score

1.00

0.78

1.00

0.89

1.00

print("Resumen de rendimiento modelo SVM Máquina de Vectores de Soporte:")

recall f1-score support

1.00

0.84

1.00

0.92

1.00

56656

56746

56746

56746

In [40]: # Pruebo el uso del modelo SVM para poder comparar que modelo funciona mejor, repetimos entreno y evaluación

56656

56746

56746

56746

Comentarios finales extraidos de el analisis realizado por Chat-GPT:

• **Precisión:** El modelo identifica correctamente el 97% de las transacciones que predice como fraudulentas.

Has aplicado y evaluado dos modelos de Machine Learning: RandomForest y SVM (Máquina de Vectores de Soporte) sobre un dataset de detección de fraudes con tarjetas de crédito, utilizando datos estandarizados. Vamos a

• F1-Score: El F1-Score de 0.84 para la clase 1 (fraudulenta) indica un buen balance entre precisión y recall, lo que es importante en contextos donde el coste de falsos negativos (transacciones fraudulentas no detectadas) es

• Recall: RandomForest tiene un mejor recall, lo que significa que detecta una mayor proporción de transacciones fraudulentas. Esto es crucial en el contexto de detección de fraudes, donde los falsos negativos (fraudes no

• Mejor equilibrio: El modelo RandomForest presenta un mejor equilibrio entre precisión y recall, lo que lo hace más adecuado para la detección de fraudes, donde es crítico identificar la mayor cantidad posible de

Para la implementación práctica, te recomendaría usar RandomForest como el modelo principal debido a su mayor fiabilidad y menor tasa de errores. Si decides utilizar SVM, asegúrate de que los datos estén correctamente

• Precisión superior: Aunque SVM tiene una precisión ligeramente superior, su menor recall y mayor número de errores falsos negativos lo hacen menos confiable como modelo principal.

• Posible complemento: SVM podría ser útil como un modelo complementario si se utilizan técnicas de ensamblado o si se busca un segundo punto de referencia en el sistema de detección.

• Recall: El modelo detecta el 73% de las transacciones realmente fraudulentas. Es decir, de las 90 transacciones fraudulentas, el modelo detecta 66 correctamente, y falla en 24.

• Matriz de confusión: El modelo solo falla en 26 casos (2 falsos positivos y 24 falsos negativos), lo que indica un excelente rendimiento general.

• Precisión: El modelo SVM es ligeramente más preciso que RandomForest, con un 98% de precisión en la detección de transacciones fraudulentas.

• Precisión: Ambos modelos son altamente precisos, pero SVM tiene una precisión ligeramente superior en la detección de transacciones fraudulentas.

• Matriz de confusión: RandomForest comete menos errores globales en comparación con SVM, lo que lo hace más confiable en la detección de fraudes.

• Menos errores: RandomForest comete menos errores en términos de falsos negativos, lo cual es preferible en la mayoría de los contextos de seguridad.

• Necesidad para SVM: Como observaste, la estandarización es crítica para el rendimiento del modelo SVM. Sin estandarización, este modelo falla drásticamente.

• RandomForest es robusto: Este modelo muestra un rendimiento consistente independientemente de la estandarización, lo que lo hace versátil y fácil de aplicar.

estandarizados y considera su uso como un modelo complementario para mejorar la precisión del sistema general de detección de fraudes.

• F1-Score: El F1-Score de 0.78 es menor que el de RandomForest, lo que sugiere un menor equilibrio entre precisión y recall.

• F1-Score: RandomForest tiene un F1-Score superior, lo que indica un mejor equilibrio general entre precisión y recall.

• Matriz de confusión: El modelo comete más errores que RandomForest, con 33 errores (1 falso positivo y 32 falsos negativos).

• Recall: Sin embargo, su capacidad para detectar transacciones fraudulentas es menor, con un 64% (el modelo no detecta 32 de las 90 transacciones fraudulentas).

90

90

Almaceno las predicciones sobre los datos de evaluación en y_pred

X_train = scaler.fit_transform(X_train)

X_test = scaler.transform(X_test)

► Haz clic aquí para obtener una pista

Crea y evalúa los modelos

Bibliotecas necesarias

modelo.fit(X_train, y_train)

Reporte de clasificación

Exactitud del modelo

Resumen de rendimiento:

0

accuracy

macro avg

weighted avg

y_pred = modelo.predict(X_test)

print("Resumen de rendimiento:")

precision

1.00

0.97

0.99

1.00

La exactitud del modelo es: 99.95%

Importo bibliotecas necesarias

Inicializar el clasificador SVM

model_SVM.fit(X_train, y_train)

Reporte de clasificación

Exactitud del modelo

accuracy

macro avg

weighted avg

[[56654

[24 [[56655]

[32

SVM_pred = model_SVM.predict(X_test)

precision

1.00

0.98

0.99

1.00

La exactitud del modelo es: 99.94%

► Haz clic aquí para obtener una pista

In [42]: from sklearn.metrics import confusion_matrix

2]

66]]

1]

58]]

Análisis de los Resultados

• Resumen de rendimiento:

■ Precision (clase 1): 0.97

■ **F1-score (clase 1):** 0.84

■ Exactitud del modelo: 99.95%

2. Resultados del Modelo SVM

■ Precision (clase 1): 0.98

■ **F1-score (clase 1):** 0.78

■ Exactitud del modelo: 99.94%

3. Comparación entre RandomForest y SVM

■ **Recall (clase 1):** 0.64

[[56655, 1], [32, 58]]

detectados) tienen un alto coste.

1. Modelo Recomendado: RandomForest:

transacciones fraudulentas.

Conclusiones Finales

2. Uso de SVM como apoyo:

3. Estandarización de datos:

Recomendación Final

Matriz de confusión:

Interpretación:

• Resumen de rendimiento:

■ **Recall (clase 1):** 0.73

[[56654, 2], [24, 66]]

Matriz de confusión:

Interpretación:

print(confusion_matrix(y_test, y_pred))

print(confusion_matrix(y_test, SVM_pred))

print(classification_report(y_test, SVM_pred))

exactitud_SVM = accuracy_score(y_test, SVM_pred)

from sklearn.svm import SVC

Entrenar el modelo

print(classification_report(y_test, y_pred))

exactitud = accuracy_score(y_test, y_pred)

In [39]: # Escribe tu código aquí

también estandarizo los datos para que funcione mejor el modelo

► Haz clic aquí para obtener una pista

In [36]: # Separa los datos de entrenamiento y evaluación

X = data_limpia.drop('Class', axis=1)

Separa del dataset

y = data_limpia['Class']

scaler = StandardScaler()

Muestra la distribución de los importes de las transacciones fraudulentas

plt.title('Distribución de los Importes de las Transacciones Fraudulentas')

sns.histplot(transacciones_fraud['Amount'], bins=50, palette="deep")

► Haz clic aquí para obtener una pista

plt.figure(figsize=(8, 6))

Mostrar el gráfico

plt.show()

250

Cantidad de Transacciones

150

50

In [32]: # Separa los datos de transacciones fraudulentas

Pongo etiquetas y título al gráfico

plt.xlabel('Importe de la Transacción') plt.ylabel('Cantidad de Transacciones')

sns.barplot(x=conteo_fraud.index, y=conteo_fraud.values, palette='deep')

plt.title('Distribución de Transacciones Fraudulentas vs No Fraudulentas')

conteo_fraud = data_limpia['Class'].value_counts()

Pregunta 1: ¿Cuántas transacciones fraudulentas hay en comparación con las no fraudulentas? (Utiliza un gráfico de barras)

Muestra la distribución de las traducciones fraudulentas con respecto de las no fraudulentas

Distribución de Transacciones Fraudulentas vs No Fraudulentas

Tipo de Transacciónes

Pregunta 2: ¿Cuál es la distribución de los importes de las transacciones fraudulentas? (Utiliza un histograma)

Distribución de los Importes de las Transacciones Fraudulentas

1000

Importe de la Transacción

Siguiendo instrucciones de Hugo Ramallo importo de la librería SKLEARN train test split,

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

Imprimo el tamaño de los datos para comprobación según se aconseja por Hugo

print("Tamaño del conjunto de entrenamiento (X_train):", X_train.shape)

print("Tamaño del conjunto de entrenamiento (y_train):", y_train.shape)

print("Tamaño del conjunto de evaluación (X_test):", X_test.shape)

print("Tamaño del conjunto de evaluación (y_test):", y_test.shape)

1500

La variable objetivo Y es class, la separ del resto de variables mediante Drop y dejo en X el resto de variables

Aqui separo los datos, uso el argumento random_state valor 42 ya que es el número comunmente utilizado, según he podidido averiguar

En la Tutoria sesión 5 Hugo comenta que estandaricemos los datos para tener mejores resultados y que el modelo funcione más eficazmente

Siguiendo la pista y tutoria de Hugo usaré en primer lugar el método RandomForest para entrenar el módelo y predict para las predicciones

Defino modelo RandomForestClassifier con los hiperparámetros indicados y lo entreno con los datos train mediante .fit()

Uso scaler para ello y estandarizo los datos de entrenamiento y test (excepto los valores objetivo y, ya que ya son 0 y 1)

2000

1: Fraudulentas

Visualiza los datos

print (conteo_fraud)

Gráfico plt.show()

283253

250000

200000

150000

100000

50000

Cantidad de Transacciones

473

Name: count, dtype: int64

Class

plt.figure(figsize=(8, 6))

Pongo etiquetas y título al gráfico

plt.ylabel('Cantidad de Transacciones')

plt.xlabel('Tipo de Transacciónes')

V7

V8

0.098698

0.085102

0.247676

0.377436

-0.270533

0.260314

0.081213

-3.807864

0.851084

0.069539

-5.433583e-

6.119264e-01

V6

V7

0.239599

-0.078803

0.791461

0.237609

0.592941

-4.918215

0.024330

-0.296827

-0.686180

1.577006

V4

V8

0.098698

0.085102

0.247676

0.377436

-0.270533

7.305334

0.294869

0.708417

0.679145

-0.414650

V5

0.001828

1.377008

-113.743307

-0.689830

-0.053468

0.612218

34.801666

V9 ...

-0.255425 ... -0.225775

-1.387024 ... -0.108300

0.363787 ...

-1.514654 ...

0.817739

1.914428

0.584800

0.432454

0.392087

0.486180

V6

-0.001139

1.331931

-26.160506

-0.769031

-0.275168

0.396792

73.301626

V9 ...

...

-0.255425 ... -0.225775

-1.387024 ... -0.108300

0.817739 ... -0.009431

-0.568671 ... -0.208254

0.363787

-1.514654 ...

0.464960

0.615375

-0.392048

-7.682956e-

3.985649e-01

V22

0.277838

-0.638672

0.771679

0.005274

0.798278

-0.559825

-0.270710

-1.015455

-0.268092

V21

-0.018307

0.247998

... -0.167716

1.943465

-0.073425

-5.540759e-01

-0.736727 ... -0.246914 -0.633753

V23

-0.110474

0.101288

0.909412

-0.190321

-0.137458

-0.026398

-0.154104

0.057504

-0.204233

-0.120794

-2.086297e-

V22

0.277838

-0.638672

0.771679

0.005274

0.798278

0.111864

0.924384

0.578229

0.800049

0.643078

V23

-0.110474

0.101288

0.909412

-0.190321

-0.137458

1.014480

0.012463

-0.037501

-0.163298

0.376777

V8

-0.000854

1.179054

-73.216718

-0.208828

0.021898

0.325704

20.007208

V24

0.066928

-0.339846

-0.689281

-1.175575

0.141267

-0.509348

-1.016226

0.640134

0.123205

0.008797

V9 ...

-0.001596 ...

1.095492 ...

-13.434066 ...

-0.644221 ...

-0.052596 ...

15.594995 ...

0.595977

V25

0.128539

0.167170

-0.327642

0.647376

-0.206010

1.436807

-0.606624

0.265745

-0.569159

-0.473649

V26

-0.189115

0.125895

-0.139097

-0.221929

0.502292

0.250034

-0.395255

-0.087371

0.546668

-0.818267

V21

-0.000371

0.723909

-34.830382

-0.228305

-0.029441

0.186194

27.202839

V21

-0.018307

0.247998

-0.009431

0.213454

0.214205

0.232045

0.265245

0.261057

V7

0.001801

1.227664

-43.557242

-0.552509

0.040859

120.589494

5.704361e-01 3.273459e-01

V24

0.066928

-0.339846

-0.689281

-1.175575

0.141267

-0.371427

-0.780055

-0.649709

-0.385050

-6.430976e-

5.971390e-01 ...

1.011592

V25

0.128539

0.167170

-0.327642

0.647376

-0.206010

-0.232794

0.750137

-0.415267

0.373205

-0.069733

V9 ...

-2.741871e-01 4.010308e-02 2.235804e-02 -5.142873e-02 ... -2.945017e-02 6.781943e-03 -1.119293e-02 4.0976

V26

-0.189115

0.125895

-0.139097

-0.221929

0.502292

0.105915

-0.257237

-0.051634

-0.384157

0.094199

1.654067e-16

01 ... -2.283949e-01

V27

0.133558

-0.008983

-0.055353

0.062723

0.219422

0.253844

0.034507

-1.206921

0.011747

0.246219

V22

-5.423504e-

1.863772e-01 5.285536e-01

7.257016e-01 6.244603e-01

-0.021053

0.014724

-0.059752

0.061458

0.215153

0.081080

0.005168

-1.085339

0.142404

0.083076

V28 Amount Clas

149.62

378.66

123.50

69.99

3.67

4.99

40.80

93.20

3.68

V23

-1.618463e-01 -3.545

2.578648e-16

1.476421e-01

V27

0.133558

-0.008983

-0.055353

0.062723

0.219422

0.943651

0.068472

0.004455

0.108821

-0.002415

V22

 0.000°

0.6237

-44.8077

-0.1617

-0.011

22.5284

-0.000015

0.724550

-10.933144

-0.542700

0.006675

0.528245

10.503090

V28 A

-0.021053

0.014724

-0.059752

0.061458

0.215153

0.823731

-0.053527

-0.026561

0.104533

0.013649

2.69