



# PIPELINES

MACHINE LEARNING

VICTOR CABALLERO

JESÚS GONZÁLEZ

VIRIDIANA ESPINOSA

SUPERVISADO / NO SUPERVISADO



# Supervisado

Data Set: Credit\_npo

*Librerías*

```
1 from sklearn.compose import ColumnTransformer
2 from sklearn.impute import SimpleImputer
3 from sklearn.pipeline import Pipeline
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.model_selection import GridSearchCV, train_test_split
6 from sklearn.metrics import classification_report, confusion_matrix
7 import pandas as pd
8 from sklearn.preprocessing import StandardScaler
9 from imblearn.over_sampling import SMOTE
10 from imblearn.pipeline import Pipeline as ImbPipeline
```

[1]





# Target - Características

## Carga el conjunto de datos



```
[72] 1 # Cargar el conjunto de datos de crédito
2 file_path_credit_npo = './data/credit_npo.csv'
3 credit_data = pd.read_csv(file_path_credit_npo)
4
5 # Revisar las primeras filas del conjunto de datos
6 credit_data.head()
```

Python

	SeriousDlqin2yrs	RevolvingUtilizationOfUnsecuredLines	age	NumberOfTime30-59DaysPastDueNotWorse	DebtRatio	MonthlyIncome	NumberOfOpenCreditLinesAndLoans	NumberOfTimes90DaysLate	NumberRealEstate
0	0	0.081892	37	0	0.070709	5656.0		12	1
1	0	0.023413	74	0	0.209197	4870.0		9	0
2	0	0.000000	43	0	0.080784	5000.0		2	0
3	0	0.492754	44	0	0.412735	7333.0		4	0
4	0	1.000000	63	0	0.000000	8333.0		3	0

## Definición de la variable objetivo y las características



```
[73] 1 # Definir la variable objetivo
2 y = credit_data['SeriousDlqin2yrs']
3
4 # Seleccionar las características (features) para el modelo
5 X = credit_data.drop(columns=['SeriousDlqin2yrs'])

[74] 1 # Dividir los datos en conjuntos de entrenamiento y prueba
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Python

# Procesamiento de datos

```
1 # Redefinir la configuración del preprocesamiento
2 # Imputar valores faltantes y escalar las características numéricas
3 num_pipeline = Pipeline(
4     steps=[
5         ("imputer", SimpleImputer(strategy="mean")), # Imputar valores faltantes
6         ("scaler", StandardScaler()) # Escalar las características
7     ]
8 )
9
10 # Definir el transformador de columnas
11 preprocessing = ColumnTransformer(
12     transformers=[
13         ("num", num_pipeline, [
14             "RevolvingUtilizationOfUnsecuredLines", "age",
15             "NumberOfTime30-59DaysPastDueNotWorse",
16             "DebtRatio", "MonthlyIncome",
17             "NumberOfOpenCreditLinesAndLoans",
18             "NumberOfTimes90DaysLate",
19             "NumberRealEstateLoansOrLines",
20             "NumberOfTime60-89DaysPastDueNotWorse",
21             "NumberOfDependents"
22        ]), # Características numéricas
23     ],
24     remainder="drop" # Descarta las columnas que no se procesan
25 )
```

## Missings:

- MonthlyIncome
- NumberOfDependents



# Definir el modelo y los parámetros para GridSearch

```
1 # Definir el modelo
2 model = RandomForestClassifier(random_state=42)
3
4 # Crear la tubería con SMOTE, el preprocesamiento y el modelo
5 pipeline = ImbPipeline(steps=[
6     ('preprocessing', preprocessing), # Primero, preprocesar los datos
7     ('smote', SMOTE(random_state=42)), # Aplicar SMOTE
8     ('model', model) # Modelo de RandomForest
9 ])
10
11 # Definir los parámetros a buscar en GridSearch
12 param_grid = {
13     'model__n_estimators': [10, 50, 100, 200], # Probar con más o menos árboles
14     'model__max_features': ['auto', 'sqrt'], # Considerar diferentes divisiones
15     'model__max_depth': [5, 10, 20, None], # Probar más niveles
16     'model__min_samples_split': [2, 5, 10, 20], # Probar diferentes cortes
17     'model__min_samples_leaf': [1, 2, 4, 6], # Probar diferentes tamaños de hojas
18     'model__class_weight': ['balanced', None] # Ajustar pesos de clase
19 }
20
21 # Ajustar GridSearchCV con validación cruzada
22 grid_search = GridSearchCV(pipeline, param_grid, cv=5, scoring='f1', verbose=2, n_jobs=-1)
23 grid_search.fit(X_train, y_train) # Ajustar el modelo a los datos de entrenamiento
```

Python

```
Fitting 5 folds for each of 1024 candidates, totalling 5120 fits
[CV] END model__class_weight=balanced, model__max_depth=5, model__max_features=auto, model__min_samples_leaf=1, model__min_samples_split=2, model__n_estimators=10; total time= 0.0s
[CV] END model__class_weight=balanced, model__max_depth=5, model__max_features=auto, model__min_samples_leaf=1, model__min_samples_split=2, model__n_estimators=10; total time= 0.0s
[CV] END model__class_weight=balanced, model__max_depth=5, model__max_features=auto, model__min_samples_leaf=1, model__min_samples_split=2, model__n_estimators=10; total time= 0.0s
[CV] END model__class_weight=balanced, model__max_depth=5, model__max_features=auto, model__min_samples_leaf=1, model__min_samples_split=2, model__n_estimators=10; total time= 0.0s
[CV] END model__class_weight=balanced, model__max_depth=5, model__max_features=auto, model__min_samples_leaf=1, model__min_samples_split=2, model__n_estimators=10; total time= 0.0s
warnings.warn(some_fits_failed_message, FitFailedWarning)
/Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages/numpy/core.py:2820: RuntimeWarning: invalid value encountered in cast
    _data = np.array(data, dtype=dtype, copy=copy,
/Library/Frameworks/Python.framework/Versions/3.12/lib/python3.12/site-packages/sklearn/model_selection/_search.py:1102: UserWarning: One or more of the test scores are non-finite: [
nan      nan      nan ... 0.42615475 0.42749987 0.4217664 ]
    warnings.warn(
```



# Mejores parámetros

```
1 # Imprimir los mejores parámetros encontrados  
2 print("Mejores parámetros:", grid_search.best_params_)  
3 print("\nMejor puntuación de validación:", grid_search.best_score_)  
4
```

```
Mejores parámetros: {'model__class_weight': 'balanced', 'model__max_depth': 20, 'model__max_features': 'sqrt', 'model__min_samples_leaf': 6, 'model__min_samples_split': 20, 'model__n_estimators': 100}  
Mejor puntuación de validación: 0.43133781550968076
```

```
1 # Hacer predicciones sobre el conjunto de prueba utilizando los mejores parámetros  
2 y_pred = grid_search.predict(X_test)  
3  
4 # Evaluar el rendimiento del modelo  
5 conf_matrix = confusion_matrix(y_test, y_pred)  
6 class_report = classification_report(y_test, y_pred)  
7  
8 # Imprimir la matriz de confusión y el informe de clasificación  
9 print("Matriz de Confusión:\n", conf_matrix)  
10 print("\nInforme de Clasificación:\n", class_report)
```

[79]

```
... Matriz de Confusión:  
[[2208 137]  
 [ 89  74]]
```

Informe de Clasificación:

	precision	recall	f1-score	support
0	0.96	0.94	0.95	2345
1	0.35	0.45	0.40	163
accuracy			0.91	2508
macro avg	0.66	0.70	0.67	2508
weighted avg	0.92	0.91	0.92	2508



No  
Supervisados





# No Supervisado

Data Set: Bebidas energéticas

## Librerías

```
1 # Importación de bibliotecas
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 from sklearn.pipeline import Pipeline, TransformerMixin
7 from sklearn.preprocessing import StandardScaler, PolynomialFeatures
8 from sklearn.decomposition import PCA
9 from sklearn.cluster import KMeans
10 from sklearn.metrics import silhouette_score, classification_report, confusion_matrix
11 from sklearn.linear_model import LogisticRegression
12 from sklearn.model_selection import GridSearchCV, train_test_split, KFold
13 from sklearn.base import BaseEstimator, ClusterMixin
14
15
16 # Configuración de estilo para gráficos
17 sns.set(style="whitegrid")
18 plt.rcParams['figure.figsize'] = (12, 6)
19
```





# Carga y exploración de los datos

01

Carga del archivo CV



```
1 # Cargar el archivo CSV y dividir en columnas
2 data = pd.read_csv('./data/empowering_drinks.csv', sep='|')
3 data.head(5)
4
```

	Azúcares	Vitamína	Cafeína	Ácido Cítrico	Taurina
0	1.518613		0.232053	1.034819	1.013009
1	0.246290		-0.827996	0.733629	0.965242
2	0.196879		1.109334	1.215533	1.395148
3	1.691550		0.487926	1.466525	2.334574
4	0.295700		1.840403	0.663351	1.186068

02

Estimador KMeans con validación cruzada

```
1 # Definición de un estimador para KMeans con validación cruzada
2 class KMeansWithSilhouette(BaseEstimator, ClusterMixin):
3     def __init__(self, n_clusters=3):
4         self.n_clusters = n_clusters
5         self.kmeans = KMeans(n_clusters=self.n_clusters, random_state=42)
6
7     def fit(self, X, y=None):
8         self.kmeans.fit(X)
9         return self
10
11    def predict(self, X):
12        return self.kmeans.predict(X)
13
14    def score(self, X, y=None): # Aceptar 'y' como argumento adicional
15        labels = self.kmeans.predict(X)
16        return silhouette_score(X, labels)
17
```

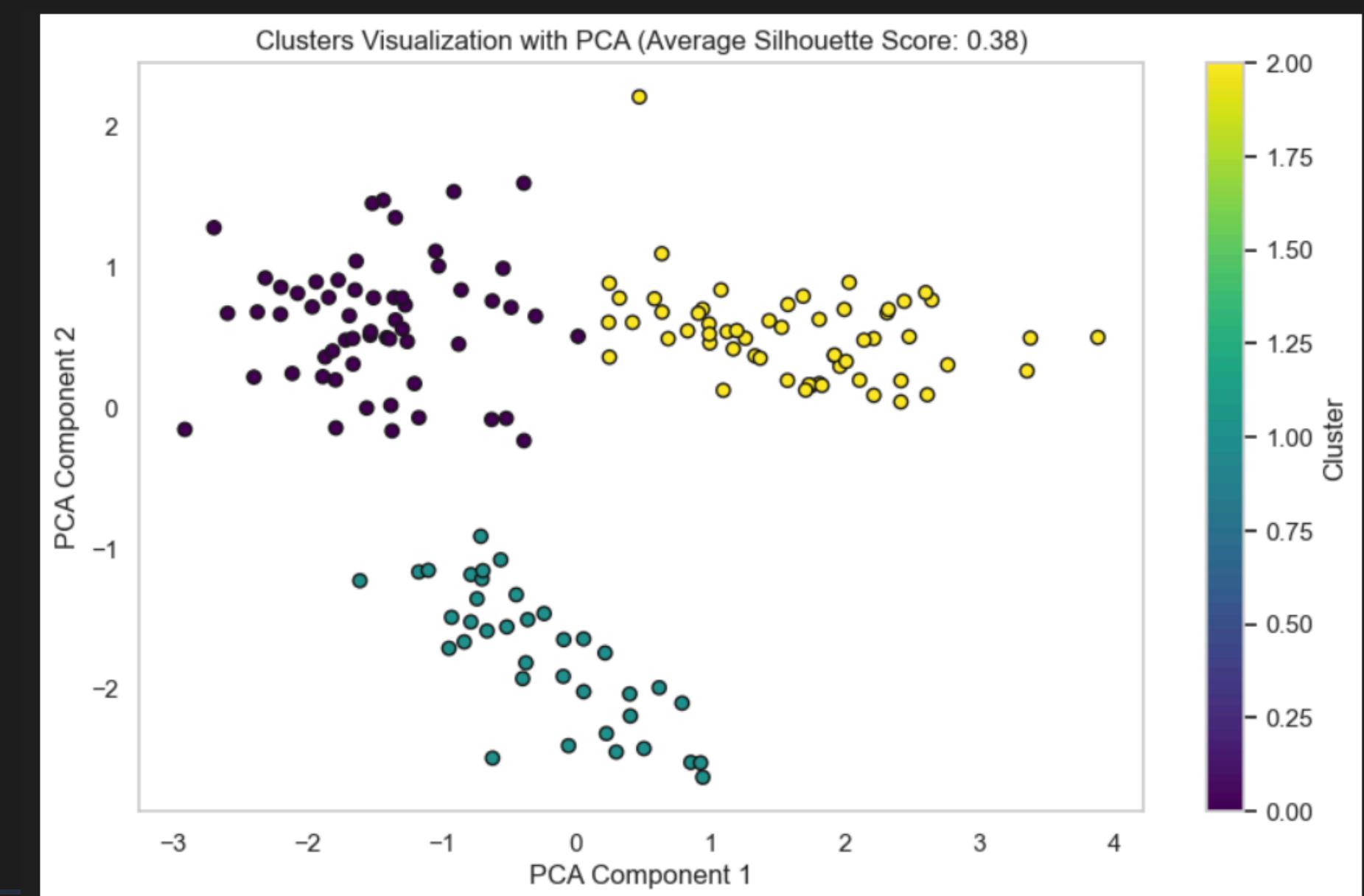
03

Preprocesamiento y Clustering

```
1 # Preprocesamiento y Clustering
2 X = data.values # Convertir a numpy array para procesar
3
4 # Crear el pipeline para clustering
5 pipeline_clustering = Pipeline(steps=[
6     ('scaler', StandardScaler()),
7     ('kmeans', KMeansWithSilhouette(n_clusters=3)) # Usar KMeans con validación cruzada
8 ])
9
10 # Validación cruzada para el modelo KMeans
11 kf = KFold(n_splits=5, shuffle=True, random_state=42)
12 silhouette_scores = []
13
14 for train_index, test_index in kf.split(X):
15     X_train, X_test = X[train_index], X[test_index]
16     pipeline_clustering.fit(X_train) # Ajustar solo en el conjunto de entrenamiento
17     score = pipeline_clustering.score(X_test) # Calcular el silhouette score
18     silhouette_scores.append(score)
19
20 average_silhouette_score = np.mean(silhouette_scores)
21
```

# Predictión Clusters / Etiquetas / PCA visualización Clusters

```
1 # Preprocesamiento y Clustering
2 X = data.values # Convertir a numpy array para procesar
3
4 # Crear el pipeline para clustering
5 pipeline_clustering = Pipeline(steps=[  
6     ('scaler', StandardScaler()),  
7     ('kmeans', KMeansWithSilhouette(n_clusters=3)) # Usar KMeans con validación cruzada  
8 ])
9
10 # Validación cruzada para el modelo KMeans
11 kf = KFold(n_splits=5, shuffle=True, random_state=42)
12 silhouette_scores = []
13
14 for train_index, test_index in kf.split(X):
15     X_train, X_test = X[train_index], X[test_index]
16     pipeline_clustering.fit(X_train) # Ajustar solo en el conjunto de entrenamiento
17     score = pipeline_clustering.score(X_test) # Calcular el silhouette score
18     silhouette_scores.append(score)
19
20 average_silhouette_score = np.mean(silhouette_scores)
21
```



# Preparación de datos



```
1 # Preparar datos para el modelo predictivo
2 X_features = data.drop(columns=['Cluster'])
3 y_labels = data['Cluster']
4
5 # Dividir los datos en conjuntos de entrenamiento y prueba
6 X_train, X_test, y_train, y_test = train_test_split(X_features, y_labels, test_size=0.2, random_state=42)
7
```

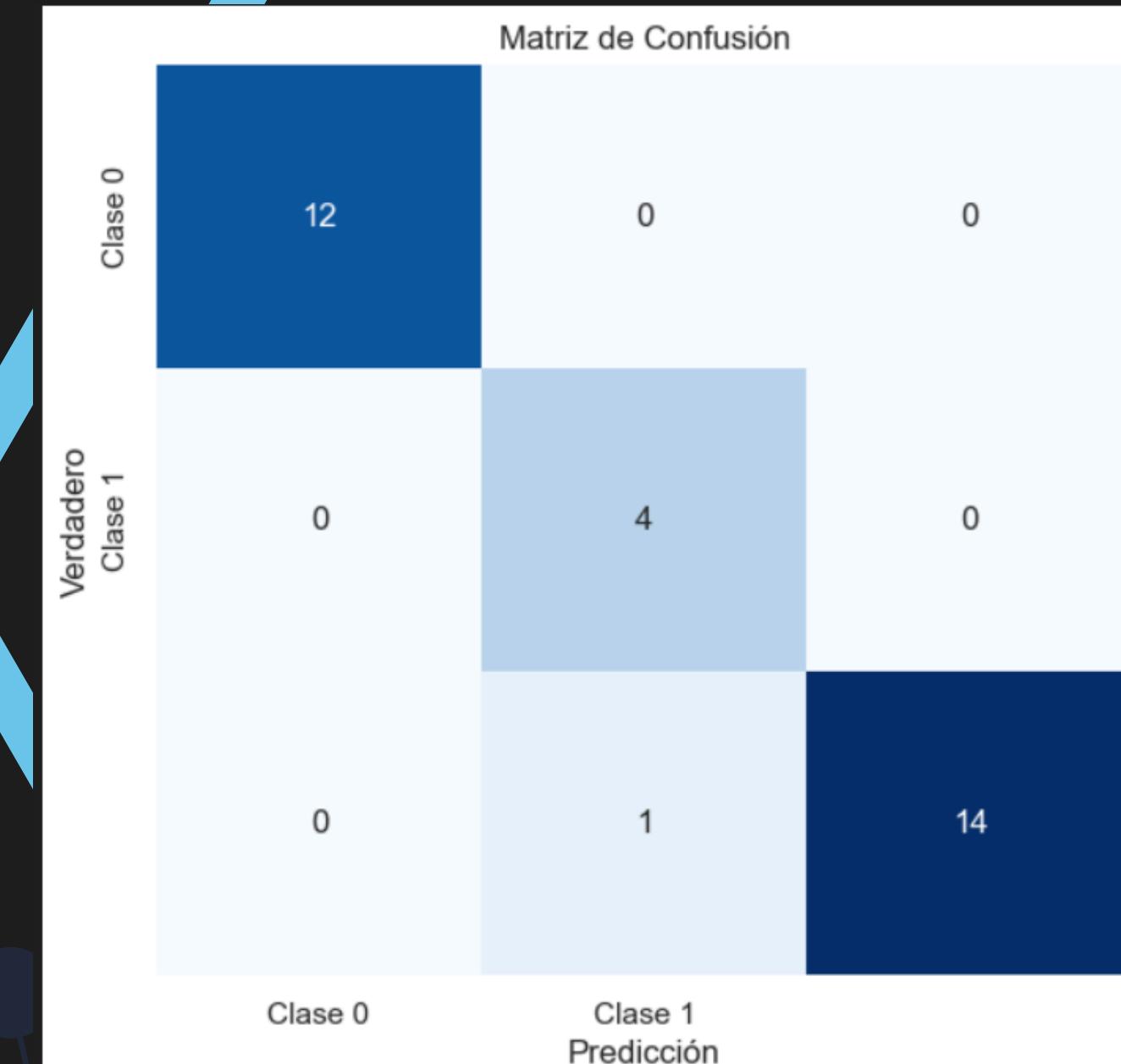
```
1 # Crear el pipeline con transformaciones y modelo
2 pipeline_model = Pipeline(steps=[
3     ('polynomial_features', PolynomialFeatures(degree=2, include_bias=False)),
4     ('scaler', StandardScaler()),
5     ('logistic_regression', LogisticRegression(max_iter=1000))
6 ])
7
8 # Configurar GridSearchCV para ajuste de hiperparámetros
9 param_grid = {
10     'logistic_regression__C': [0.01, 0.1, 1, 10, 100],
11     'logistic_regression__solver': ['liblinear', 'saga']
12 }
13 grid_search = GridSearchCV(pipeline_model, param_grid, cv=5, scoring='accuracy')
14
15 # Ajustar el modelo a los datos de entrenamiento
16 grid_search.fit(X_train, y_train)
17
18 # Mejor parámetros y puntuación
19 best_params = grid_search.best_params_
20 best_score = grid_search.best_score_
```

# Evaluación Modelo y Matriz de Confusión

```
1 # Evaluar el modelo en el conjunto de prueba
2 y_pred = grid_search.predict(X_test)
3
4 # Reporte de clasificación
5 report = classification_report(y_test, y_pred)
6 conf_matrix = confusion_matrix(y_test, y_pred)
7
8 best_params, best_score, report, conf_matrix, average_silhouette_score
9
```

```
{'logistic_regression_C': 10, 'logistic_regression_solver': 'liblinear'},
1.0,
precision      recall   f1-score support\n0          1.00    1.00     1.00    12\n1          0.80    0.89     0.89     4\n2          0.89    0.89     0.89     2
```

```
1 plt.figure(figsize=(8, 6))
2 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels=['Clase 0', 'Clase 1'], yticklabels=['Clase 0', 'Clase 1'])
3 plt.ylabel('Verdadero')
4 plt.xlabel('Predicción')
5 plt.title('Matriz de Confusión')
6 plt.show()
```



# ¡Gracias!

