

# **Actividad 3.5**

## **Clasificación de vinos**



# Índice

Objetivo de la actividad.....	3
1. Importación de los datasets (utilizar el dataset winequality-red.csv).....	3
2. Mostrar la matriz de correlación de variables.....	4
3. Aplicar cualquier otra técnica de selección de características que consideres adecuada y justificar tu propuesta.....	5
4. Realizar una comparativa de la precisión en el entrenamiento de los diferentes modelos de Naives Bayes y KNN. Aplicando Cross Validation.....	6
5. Entrenamiento del Mejor Modelo y Matriz de Confusión.....	7
6. Comparar el resultado obtenido con el valor de calidad indicado en el dataset por medio de una matriz de confusión.....	9
NUEVOS ENFOQUES.....	10
A) Método de Clasificación con Random Forest Classifier junto con GridSearchCV para optimizar los hiperparámetro.....	10
Resultado y explicación.....	11
B) Método de Clasificación: Gradient Boosting Classifier y SMOTE.....	12
Resultado y explicación.....	13
C) Método de Clasificación: HistGradientBoostingClassifier.....	14
Resultado y explicación.....	15
7. Probar a utilizar el cuaderno con el dataset de los vinos blancos y realizar captura de los resultados obtenidos. (utilizar el dataset winequality-white.csv).....	16
Resultado y explicación.....	17
8. Conclusiones de la actividad.....	18
Repositorios.....	19

## Objetivo de la actividad

El objetivo de esta actividad es poner en práctica los conocimientos adquiridos hasta el momento para ello vamos a utilizar el siguiente dataset que contiene una serie de características físico-químicas que determina la calidad del vino en una escala de valores del 1 al 10.

Dataset es el siguiente: <https://archive.ics.uci.edu/ml/datasets/Wine+Quality>

Proyecto de partida: [https://colab.research.google.com/drive/1hwri6X-N\\_cHmpZs3l-zyK2XwRufA4EGN?usp=sharing](https://colab.research.google.com/drive/1hwri6X-N_cHmpZs3l-zyK2XwRufA4EGN?usp=sharing)

### 1. Importación de los datasets (utilizar el dataset winequality-red.csv)

#### ✓ 1. Importación del Dataset (winequality-red.csv)

Se utilizará el dataset winequality-red.csv disponible en UCI Machine Learning Repository.

Este contiene características físico-químicas de vinos tintos y su respectiva calidad en una escala de 1 a 10.

```
[1] 1 # IVAN FALCON MONZON
2 import pandas as pd # Importamos la librería Pandas para el manejo de datos
3
4 # URL del dataset de vinos tintos
5 dataset_url = "https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv"
6
7 # Cargar el dataset en un DataFrame, especificando que el separador es ';'
8 df = pd.read_csv(dataset_url, sep=';')
9
10 # Mostrar las primeras 5 filas del dataset para inspección rápida
11 print(df.head())
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	\
0	7.4	0.70	0.00	1.9	0.076	
1	7.8	0.88	0.00	2.6	0.098	
2	7.8	0.76	0.04	2.3	0.092	
3	11.2	0.28	0.56	1.9	0.075	
4	7.4	0.70	0.00	1.9	0.076	

	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	\
0	11.0	34.0	0.9978	3.51	0.56	
1	25.0	67.0	0.9968	3.20	0.68	
2	15.0	54.0	0.9970	3.26	0.65	
3	17.0	60.0	0.9980	3.16	0.58	
4	11.0	34.0	0.9978	3.51	0.56	

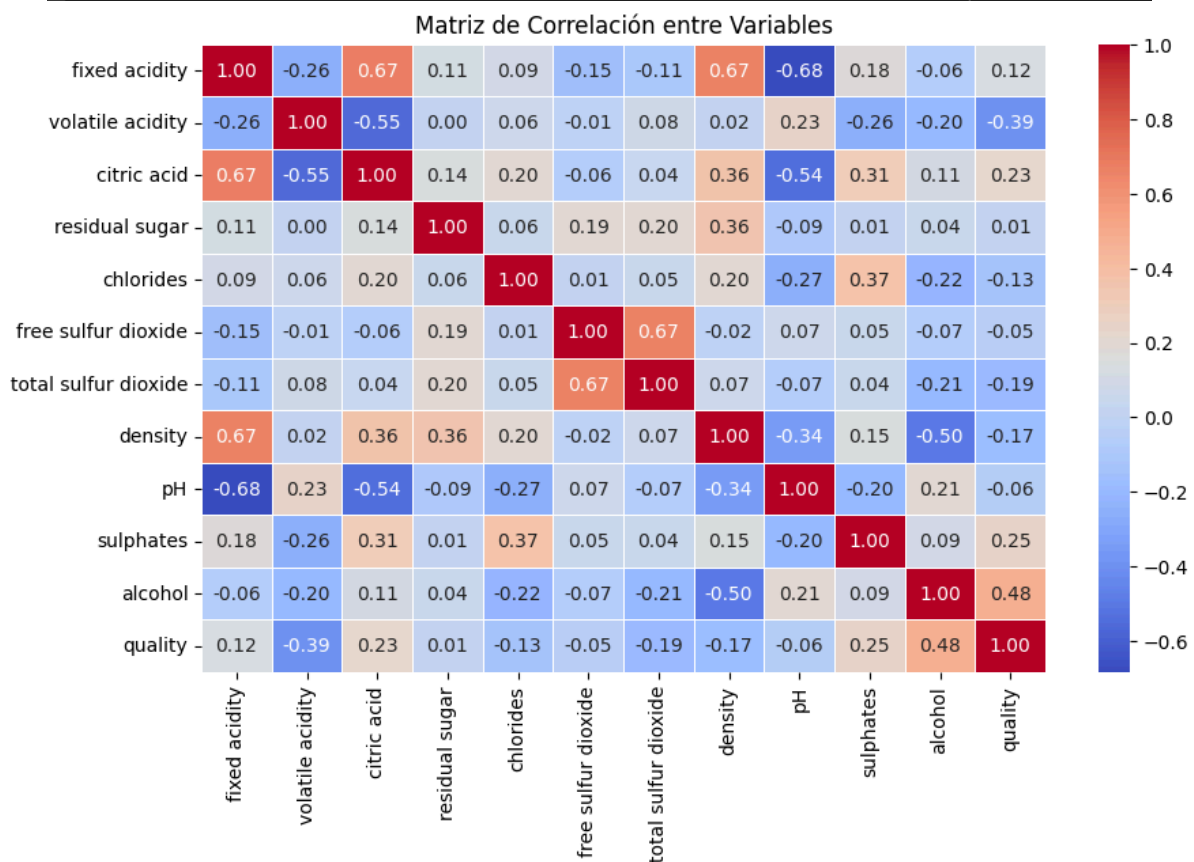
	alcohol	quality
0	9.4	5
1	9.8	5
2	9.8	5
3	9.8	6
4	9.4	5

## 2. Mostrar la matriz de correlación de variables.

### 2. Matriz de Correlación de variables

Se visualizará la matriz de correlación entre las variables para identificar relaciones fuertes.

```
[2] 1 # IVAN FALCON MONZON
2 import seaborn as sns # Importamos Seaborn para visualización de datos
3 import matplotlib.pyplot as plt # Importamos Matplotlib para manejar gráficos
4
5 # Configurar el tamaño de la figura para mejorar la visualización
6 plt.figure(figsize=(10, 6))
7
8 # Crear un mapa de calor para mostrar la correlación entre variables del DataFrame
9 sns.heatmap(
10     df.corr(), # Calcula la matriz de correlación
11     annot=True, # Muestra los valores dentro de cada celda
12     cmap='coolwarm', # Define la paleta de colores
13     fmt='.2f', # Formato de los valores numéricos (2 decimales)
14     linewidths=0.5 # Espaciado entre las celdas
15 )
16
17 # Agregar un título a la visualización
18 plt.title("Matriz de Correlación entre Variables")
19
20 # Mostrar el gráfico
21 plt.show()
```



### 3. Aplicar cualquier otra técnica de selección de características que consideres adecuada y justificar tu propuesta.

#### 3. Selección de Características

Se seleccionarán las características más relevantes para la clasificación.

Características seleccionadas:

- volatile acidity
- total sulfur dioxide
- density
- sulphates
- alcohol

```
[3] 1 # IVAN FALCON MONZON
2 from sklearn.model_selection import train_test_split # Importamos la función para dividir los datos
3
4 # Selección de características relevantes para el modelo
5 selected_features = ['volatile acidity', 'total sulfur dioxide', 'density', 'sulphates', 'alcohol']
6 X = df[selected_features] # Variables predictoras
7 y = df['quality'] # Variable objetivo (calidad del vino)
8
9 # División de los datos en conjunto de entrenamiento (80%) y prueba (20%)
10 X_train, X_test, y_train, y_test = train_test_split(
11     X, y, test_size=0.2, random_state=42
12 )
13
14 # 'random_state=42' asegura reproducibilidad en la división de datos
```

#### Clasificación por Random Forest

```
[4] 1 # IVAN FALCON MONZON
2 from sklearn.ensemble import RandomForestClassifier # Importamos el clasificador Random Forest
3 from sklearn.metrics import accuracy_score, classification_report # Métricas de evaluación
4
5 # Crear el modelo de Random Forest
6 rf_model = RandomForestClassifier(n_estimators=100, random_state=42) # 100 árboles en el bosque
7
8 # Entrenar el modelo con los datos de entrenamiento
9 rf_model.fit(X_train, y_train)
10
11 # Realizar predicciones sobre el conjunto de prueba
12 y_pred = rf_model.predict(X_test)
13
14 # Evaluación del modelo
15 accuracy = accuracy_score(y_test, y_pred) # Precisión del modelo
16 report = classification_report(y_test, y_pred) # Reporte de métricas detalladas
17
18 # Mostrar resultados
19 print(f"Precisión del modelo: {accuracy:.2f}")
20 print("Reporte de clasificación:")
21 print(report)
```

Precisión del modelo: 0.68  
Reporte de clasificación:

	precision	recall	f1-score	support
3	0.00	0.00	0.00	1
4	0.00	0.00	0.00	10
5	0.72	0.79	0.75	130
6	0.65	0.69	0.67	132
7	0.63	0.52	0.57	42
8	0.00	0.00	0.00	5
accuracy			0.68	320
macro avg	0.33	0.33	0.33	320
weighted avg	0.64	0.68	0.66	320

#### Explicación del código:

1. Importamos RandomForestClassifier de sklearn.ensemble.
2. Se crea el modelo en este caso con 100 árboles (n\_estimators=100).
3. Entrenamos el modelo con X\_train y y\_train.
4. Predicimos valores con X\_test.
5. Calculamos la precisión con accuracy\_score().
6. Mostramos el classification\_report(), que incluye precisión, recall y F1-score.

4. Realizar una comparativa de la precisión en el entrenamiento de los diferentes modelos de Naives Bayes y KNN. Aplicando Cross Validation.

#### ✓ 4. Comparación de Modelos: Naive Bayes vs KNN

Se comparará la precisión de Naive Bayes y KNN utilizando Cross Validation.

```
[5] 1 # IVAN FALCON MONZON
2 from sklearn.naive_bayes import GaussianNB # Importamos el modelo Naive Bayes Gaussiano
3 from sklearn.neighbors import KNeighborsClassifier # Importamos K-Nearest Neighbors (KNN)
4 from sklearn.model_selection import cross_val_score # Importamos Cross Validation para evaluación
5
6 # Definimos los modelos
7 nb_model = GaussianNB() # Modelo Naive Bayes Gaussiano
8 knn_model = KNeighborsClassifier(n_neighbors=5) # Modelo KNN con 5 vecinos
9
10 # Evaluación con validación cruzada de 5 folds (cv=5)
11 nb_scores = cross_val_score(nb_model, X, y, cv=5) # Evaluamos Naive Bayes
12 knn_scores = cross_val_score(knn_model, X, y, cv=5) # Evaluamos KNN
13
14 # Mostramos la precisión promedio de cada modelo
15 print(f"Precisión promedio Naive Bayes: {nb_scores.mean():.2f}")
16 print(f"Precisión promedio KNN: {knn_scores.mean():.2f}")
```

➡ Precisión promedio Naive Bayes: 0.58  
Precisión promedio KNN: 0.50

##### Resultados obtenidos:

- Precisión Naive Bayes: 0.58
- Precisión KNN: 0.50

## 5. Entrenamiento del Mejor Modelo y Matriz de Confusión

### 5. Entrenamiento del Mejor Modelo y Matriz de Confusión

Se entrenará el modelo que obtuvo mejor rendimiento (En mi caso Naive Bayes).

```
1 # IVAN FALCON MONZON
2 from sklearn.metrics import confusion_matrix, classification_report # Importamos métricas de evaluación
3
4 # Entrenar el modelo seleccionado (Naive Bayes, ya que tuvo mejor desempeño)
5 best_model = GaussianNB()
6 best_model.fit(X_train, y_train) # Entrenamos el modelo con los datos de entrenamiento
7
8 # Realizamos predicciones en el conjunto de prueba
9 y_pred = best_model.predict(X_test)
10
11 # Matriz de confusión para evaluar el desempeño del modelo
12 conf_matrix = confusion_matrix(y_test, y_pred)
13 print("Matriz de Confusión:")
14 print(conf_matrix)
15
16 # Reporte detallado con métricas de clasificación (precisión, recall, f1-score)
17 print("Reporte de Clasificación:")
18 print(classification_report(y_test, y_pred))
```



Matriz de Confusión:

```
[[ 0  0  1  0  0  0]
 [ 1  1  8  0  0  0]
 [ 0  0 99 31  0  0]
 [ 0  1 49 65 17  0]
 [ 0  0  2 31  9  0]
 [ 0  0  0  1  4  0]]
```

Reporte de Clasificación:

	precision	recall	f1-score	support
3	0.00	0.00	0.00	1
4	0.50	0.10	0.17	10
5	0.62	0.76	0.69	130
6	0.51	0.49	0.50	132
7	0.30	0.21	0.25	42
8	0.00	0.00	0.00	5
accuracy			0.54	320
macro avg	0.32	0.26	0.27	320
weighted avg	0.52	0.54	0.52	320

## Guardar el modelo para próximos modelos

```
[7] 1 # IVAN FALCON MONZON
2 import joblib # Importamos joblib para guardar y cargar modelos
3
4 # Guardar el modelo entrenado en un archivo .pkl
5 modelo_path = 'modelo_vinos.pkl' # Nombre del archivo donde se guardará el modelo
6 joblib.dump(best_model, modelo_path)
7
8 # Verificar que el modelo se ha guardado correctamente
9 print(f"Modelo guardado exitosamente en '{modelo_path}')
```

⇒ Modelo guardado exitosamente en 'modelo\_vinos.pkl'

## Cargar el modelo guardado

```
[8] 1 # IVAN FALCON MONZON
2 import joblib # Importamos joblib para cargar el modelo guardado
3
4 # Cargar el modelo previamente guardado
5 modelo_path = 'modelo_vinos.pkl'
6 loaded_model = joblib.load(modelo_path)
7
8 print(f"Modelo cargado exitosamente desde '{modelo_path}')
```

```
9
10 # Verificación opcional: realizar una predicción con los datos de prueba
11 sample_prediction = loaded_model.predict(X_test[:5]) # Predecir las primeras 5 muestras
12 print("Ejemplo de predicciones con el modelo cargado:", sample_prediction)
```

⇒ Modelo cargado exitosamente desde 'modelo\_vinos.pkl'  
Ejemplo de predicciones con el modelo cargado: [5 5 5 5 6]



## 6. Comparar el resultado obtenido con el valor de calidad indicado en el dataset por medio de una matriz de confusión

### 6. Comparación con la Calidad del Dataset

Se evaluará si la predicción del modelo coincide con los valores reales de calidad.

```
[9] 1 # IVAN FALCON MONZON
    2 from sklearn.metrics import confusion_matrix, classification_report # Importamos las métricas necesarias
    3
    4 # Aplicar el modelo cargado a los datos de prueba
    5 y_pred = loaded_model.predict(X_test) # Realizamos predicciones sobre el conjunto de prueba
    6
    7 # Generar la Matriz de Confusión
    8 conf_matrix = confusion_matrix(y_test, y_pred)
    9 print("Matriz de Confusión:")
   10 print(conf_matrix)
   11
   12 # Generar el Reporte de Clasificación (precisión, recall, F1-score, etc.)
   13 print("Reporte de Clasificación:")
   14 print(classification_report(y_test, y_pred))
```

```
Matriz de Confusión:
[[ 0  0  1  0  0  0]
 [ 1  1  8  0  0  0]
 [ 0  0 99 31  0  0]
 [ 0  1 49 65 17  0]
 [ 0  0  2 31  9  0]
 [ 0  0  0  1  4  0]]
Reporte de Clasificación:
              precision    recall  f1-score   support

     3         0.00         0.00         0.00         1
     4         0.50         0.10         0.17         10
     5         0.62         0.76         0.69        130
     6         0.51         0.49         0.50        132
     7         0.30         0.21         0.25         42
     8         0.00         0.00         0.00          5

 accuracy          0.54          0.54          0.54        320
 macro avg         0.32         0.26         0.27        320
 weighted avg         0.52         0.54         0.52        320
```

```
[10] 1 # IVAN FALCON MONZON
      2 # Evaluación final del modelo usando la precisión en el conjunto de prueba
      3 accuracy = loaded_model.score(X_test, y_test) # Calculamos la precisión final
      4 print(f"Precisión final del modelo: {accuracy:.4f}")
```

```
Precisión final del modelo: 0.5437
```

## NUEVOS ENFOQUES

### ▼ NUEVOS ENFOQUES

En este apartado entrenaré con diferentes modelos y métodos para que de un mejor resultado en la predicción.

A) Método de Clasificación con Random Forest Classifier junto con GridSearchCV para optimizar los hiperparámetro

### ▼ A) Método de Clasificación con Random Forest Classifier junto con GridSearchCV para optimizar los hiperparámetros

```
[11] 1 # IVAN FALCON MONZON
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import joblib # Para guardar y cargar el modelo
6 from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV
7 from sklearn.ensemble import RandomForestClassifier
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.metrics import confusion_matrix, classification_report
10
11 # Cargar el dataset de vinos tintos
12 dataset_url = "https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv"
13 df = pd.read_csv(dataset_url, sep=';')
14
15 # Selección de características relevantes del dataset
16 selected_features = ['volatile acidity', 'total sulfur dioxide', 'density', 'sulphates', 'alcohol']
17 X = df[selected_features] # Variables predictoras
18 y = df['quality'] # Variable objetivo (calidad del vino)
19
20 # Normalización de los datos para mejorar el rendimiento del modelo
21 scaler = StandardScaler()
22 X_scaled = scaler.fit_transform(X)
23
24 # División en conjunto de entrenamiento y conjunto de prueba (80%-20%)
25 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)
26
27 # Definir el modelo Random Forest con búsqueda de hiperparámetros utilizando GridSearchCV
28 tuned_parameters = {
29     'n_estimators': [50, 100, 200], # Número de árboles en el bosque
30     'max_depth': [None, 10, 20], # Profundidad máxima de los árboles
31     'min_samples_split': [2, 5, 10] # Mínimo número de muestras para dividir un nodo
32 }
33 rf_model = GridSearchCV(RandomForestClassifier(random_state=42), tuned_parameters, cv=5, scoring='accuracy')
34
35 # Entrenar el modelo con los mejores parámetros encontrados a través de GridSearchCV
36 rf_model.fit(X_train, y_train)
37 print(f"Mejores parámetros encontrados: {rf_model.best_params_}") # Imprimir los mejores parámetros
38
39 # Guardar el modelo optimizado en un archivo
40 joblib.dump(rf_model.best_estimator_, 'modelo_vinos_mejorado.pkl')
41 print("Modelo mejorado guardado exitosamente en 'modelo_vinos_mejorado.pkl'")
42
43 # Cargar el modelo optimizado desde el archivo
44 loaded_model = joblib.load('modelo_vinos_mejorado.pkl')
45 print("Modelo mejorado cargado exitosamente.")
46
47 # Realizar predicciones en el conjunto de prueba
48 y_pred = loaded_model.predict(X_test)
49
50 # Generar la Matriz de Confusión para evaluar las predicciones
51 conf_matrix = confusion_matrix(y_test, y_pred)
52 print("\nMatriz de Confusión:")
53 print(conf_matrix)
54
55 # Generar el Reporte de Clasificación, que muestra precisión, recall, f1-score, etc.
56 print("\nReporte de Clasificación:")
57 print(classification_report(y_test, y_pred))
58
59 # Evaluación final del modelo: Precisión en el conjunto de prueba
60 accuracy = loaded_model.score(X_test, y_test)
61 print(f"\nPrecisión final del modelo mejorado: {accuracy:.4f}")
```

## Resultado y explicación

```
➡ Mejores parámetros encontrados: {'max_depth': 20, 'min_samples_split': 5, 'n_estimators': 100}
Modelo mejorado guardado exitosamente en 'modelo_vinos_mejorado.pkl'
Modelo mejorado cargado exitosamente.

Matriz de Confusión:
[[ 0  0  1  0  0  0]
 [ 0  0  6  4  0  0]
 [ 0  0 103 27  0  0]
 [ 0  0 33 90  9  0]
 [ 0  0  1 19 22  0]
 [ 0  0  0  0  5  0]]

Reporte de Clasificación:
              precision    recall  f1-score   support

     3         0.00         0.00         0.00         1
     4         0.00         0.00         0.00        10
     5         0.72         0.79         0.75       130
     6         0.64         0.68         0.66       132
     7         0.61         0.52         0.56        42
     8         0.00         0.00         0.00         5

 accuracy          0.67         0.67         0.67       320
 macro avg         0.33         0.33         0.33       320
 weighted avg         0.64         0.67         0.65       320

Precisión final del modelo mejorado: 0.6719
```

### Explicación del código (A)

1. Cargar y preparar los datos: Se carga el dataset, se seleccionan las 3 características relevantes y se normalizan los datos.
2. División de los datos: Se divide el dataset en entrenamiento y prueba (80%/20%).
3. Optimización de hiperparámetros: Se utiliza GridSearchCV para encontrar los mejores parámetros del modelo Random Forest.
4. Entrenamiento y evaluación: El modelo optimizado se entrena, se guarda, se carga y se evalúa utilizando la matriz de confusión, el reporte de clasificación y la precisión.

## B) Método de Clasificación: Gradient Boosting Classifier y SMOTE

### ▼ B) Método de Clasificación: Gradient Boosting Classifier y SMOTE

El código utiliza el **Gradient Boosting Classifier** junto con **SMOTE** para abordar el desbalanceo en las clases y **GridSearchCV** para optimizar los hiperparámetros.

```
[ ] 1 # IVAN FALCON MONZON
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import joblib
6 from imblearn.over_sampling import SMOTE # Para balancear las clases con SMOTE
7 from sklearn.model_selection import train_test_split, GridSearchCV
8 from sklearn.ensemble import GradientBoostingClassifier # Clasificador de Gradient Boosting
9 from sklearn.preprocessing import StandardScaler # Para normalizar los datos
10 from sklearn.metrics import confusion_matrix, classification_report # Para evaluar el modelo
11
12 # Cargar el dataset de vinos tintos desde la URL
13 dataset_url = "https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv"
14 df = pd.read_csv(dataset_url, sep=";")
15
16 # Selección de características relevantes
17 selected_features = ['volatile acidity', 'total sulfur dioxide', 'density', 'sulphates', 'alcohol']
18 X = df[selected_features] # Variables predictoras
19 y = df['quality'] # Variable objetivo (calidad del vino)
20
21 # Normalización de los datos para mejorar el rendimiento del modelo
22 scaler = StandardScaler()
23 X_scaled = scaler.fit_transform(X)
24
25 # Manejo del desbalanceo con SMOTE (Generar muestras sintéticas para las clases minoritarias)
26 smote = SMOTE(random_state=42)
27 X_resampled, y_resampled = smote.fit_resample(X_scaled, y)
28
29 # División de los datos en conjunto de entrenamiento (80%) y prueba (20%)
30 X_train, X_test, y_train, y_test = train_test_split(X_resampled, y_resampled, test_size=0.2, random_state=42)
31
32 # Definir el modelo Gradient Boosting con búsqueda de hiperparámetros utilizando GridSearchCV
33 tuned_parameters = {
34     'n_estimators': [100, 200], # Número de árboles en el modelo
35     'learning_rate': [0.05, 0.1], # Tasa de aprendizaje
36     'max_depth': [3, 5] # Profundidad máxima de los árboles
37 }
38 gb_model = GridSearchCV(GradientBoostingClassifier(random_state=42), tuned_parameters, cv=5, scoring='accuracy')
39
40 # Entrenar el modelo con los mejores parámetros encontrados a través de GridSearchCV
41 gb_model.fit(X_train, y_train)
42 print(f"Mejores parámetros encontrados: {gb_model.best_params_}") # Imprimir los mejores parámetros
43
44 # Guardar el modelo optimizado en un archivo
45 joblib.dump(gb_model.best_estimator_, 'modelo_vinos_optimizado.pkl')
46 print("Modelo optimizado guardado exitosamente en 'modelo_vinos_optimizado.pkl'")
47
48 # Cargar el modelo optimizado desde el archivo guardado
49 loaded_model = joblib.load('modelo_vinos_optimizado.pkl')
50 print("Modelo optimizado cargado exitosamente.")
51
52 # Realizar predicciones con el modelo cargado sobre el conjunto de prueba
53 y_pred = loaded_model.predict(X_test)
54
55 # Generar la Matriz de Confusión para evaluar las predicciones
56 conf_matrix = confusion_matrix(y_test, y_pred)
57 print("\nMatriz de Confusión:")
58 print(conf_matrix)
59
60 # Generar el Reporte de Clasificación, que incluye precisión, recall, f1-score, etc.
61 print("\nReporte de Clasificación:")
62 print(classification_report(y_test, y_pred))
63
64 # Evaluación final del modelo: Precisión en el conjunto de prueba
65 accuracy = loaded_model.score(X_test, y_test)
66 print(f"\nPrecisión final del modelo optimizado: {accuracy:.4f}")
```

## Resultado y explicación

➡ Mejores parámetros encontrados: {'learning\_rate': 0.05, 'max\_depth': 5, 'n\_estimators': 200}  
Modelo optimizado guardado exitosamente en 'modelo\_vinos\_optimizado.pkl'  
Modelo optimizado cargado exitosamente.

Matriz de Confusión:

```
[[128  0  4  0  0  0]
 [ 4 117  3  4  3  0]
 [ 3 10 109 23  5  0]
 [ 1  9  39  90 12  4]
 [ 0  0  3 11  96  6]
 [ 0  0  0  0  2 132]]
```

Reporte de Clasificación:

	precision	recall	f1-score	support
3	0.94	0.97	0.96	132
4	0.86	0.89	0.88	131
5	0.69	0.73	0.71	150
6	0.70	0.58	0.64	155
7	0.81	0.83	0.82	116
8	0.93	0.99	0.96	134
accuracy		0.82		818
macro avg	0.82	0.83	0.83	818
weighted avg	0.82	0.82	0.82	818

Precisión final del modelo optimizado: 0.8215

### Explicación del flujo del código (B:

1. Carga y preparación de datos: Se carga el dataset, se seleccionan las características relevantes y se normalizan los datos.
2. Manejo de desbalanceo: Se usa SMOTE para generar ejemplos sintéticos y equilibrar las clases en el dataset.
3. División de los datos: Se divide el dataset en entrenamiento y prueba (80%/20%).
4. Optimización de hiperparámetros: Se usa GridSearchCV para encontrar los mejores parámetros para el modelo Gradient Boosting.
5. Evaluación del modelo: Se calcula la precisión, la matriz de confusión y el reporte de clasificación para evaluar el rendimiento del modelo.

## C) Método de Clasificación: HistGradientBoostingClassifier

### ▼ C) Método de Clasificación: HistGradientBoostingClassifier

El código utiliza el HistGradientBoostingClassifier para realizar la clasificación de vinos basándose en varias características seleccionadas. Este modelo es una variante más eficiente y rápida del Gradient Boosting que es particularmente útil cuando se manejan grandes volúmenes de datos o cuando el conjunto de datos contiene muchas características.

```
[ ] 1 # IVAN FALCON MONZON
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import joblib
6 from imblearn.over_sampling import SMOTE # Para balancear las clases con SMOTE
7 from sklearn.model_selection import train_test_split, GridSearchCV # Para dividir y ajustar los hiperparámetros
8 from sklearn.ensemble import HistGradientBoostingClassifier # Clasificador HistGradientBoosting
9 from sklearn.preprocessing import StandardScaler # Para normalizar los datos
10 from sklearn.metrics import confusion_matrix, classification_report # Para evaluar el modelo
11
12 # Cargar el dataset de vinos tintos desde la URL
13 dataset_url = "https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv"
14 df = pd.read_csv(dataset_url, sep=";")
15
16 # Selección de características relevantes para la predicción
17 selected_features = ['volatile acidity', 'total sulfur dioxide', 'density', 'sulphates', 'alcohol']
18 X = df[selected_features] # Variables predictoras
19 y = df['quality'] # Variable objetivo (calidad del vino)
20
21 # Normalización de los datos para mejorar el rendimiento del modelo
22 scaler = StandardScaler()
23 X_scaled = scaler.fit_transform(X)
24
25 # Manejo del desbalanceo con SMOTE (Generar muestras sintéticas para las clases minoritarias)
26 smote = SMOTE(random_state=42)
27 X_resampled, y_resampled = smote.fit_resample(X_scaled, y)
28
29 # División de los datos en conjunto de entrenamiento (80%) y prueba (20%)
30 X_train, X_test, y_train, y_test = train_test_split(X_resampled, y_resampled, test_size=0.2, random_state=42)
31
32 # Definir el modelo HistGradientBoosting con búsqueda de hiperparámetros utilizando GridSearchCV
33 tuned_parameters = {
34     'learning_rate': [0.01, 0.05, 0.1], # Tasa de aprendizaje
35     'max_iter': [100, 200, 300], # Número de iteraciones (número de árboles)
36     'max_depth': [3, 5, 7] # Profundidad máxima de los árboles
37 }
38
39 # Inicializar y optimizar el modelo utilizando GridSearchCV
40 gb_model = GridSearchCV(HistGradientBoostingClassifier(random_state=42), tuned_parameters, cv=5, scoring='accuracy')
41
42 # Entrenar el modelo con los mejores parámetros encontrados a través de GridSearchCV
43 gb_model.fit(X_train, y_train)
44 print(f"Mejores parámetros encontrados: {gb_model.best_params_}") # Imprimir los mejores parámetros
45
46 # Guardar el modelo optimizado en un archivo
47 joblib.dump(gb_model.best_estimator_, 'modelo_vinos_avanzado.pkl')
48 print("Modelo avanzado guardado exitosamente en 'modelo_vinos_avanzado.pkl'")
49
50 # Cargar el modelo optimizado desde el archivo guardado
51 loaded_model = joblib.load('modelo_vinos_avanzado.pkl')
52 print("Modelo avanzado cargado exitosamente.")
53
54 # Realizar predicciones con el modelo cargado sobre el conjunto de prueba
55 y_pred = loaded_model.predict(X_test)
56
57 # Generar la Matriz de Confusión para evaluar las predicciones
58 conf_matrix = confusion_matrix(y_test, y_pred)
59 print("\nMatriz de Confusión:")
60 print(conf_matrix)
61
62 # Generar el Reporte de Clasificación, que incluye precisión, recall, f1-score, etc.
63 print("\nReporte de Clasificación:")
64 print(classification_report(y_test, y_pred))
65
66 # Evaluación final del modelo: Precisión en el conjunto de prueba
67 accuracy = loaded_model.score(X_test, y_test)
68 print(f"\nPrecisión final del modelo avanzado: {accuracy:.4f}")
```

## Resultado y explicación

Mejores parámetros encontrados: {'learning\_rate': 0.1, 'max\_depth': 7, 'max\_iter': 300}  
Modelo avanzado guardado exitosamente en 'modelo\_vinos\_avanzado.pkl'  
Modelo avanzado cargado exitosamente.

Matriz de Confusión:

```
[[132  0  0  0  0  0]
 [ 3 122  4  1  1  0]
 [ 4  5 110 26  5  0]
 [ 2  5 37  95 11  5]
 [ 0  0  3 14  92  7]
 [ 0  0  0  0  1 133]]
```

Reporte de Clasificación:

	precision	recall	f1-score	support
3	0.94	1.00	0.97	132
4	0.92	0.93	0.93	131
5	0.71	0.73	0.72	150
6	0.70	0.61	0.65	155
7	0.84	0.79	0.81	116
8	0.92	0.99	0.95	134
accuracy			0.84	818
macro avg	0.84	0.84	0.84	818
weighted avg	0.83	0.84	0.83	818

Precisión final del modelo avanzado: 0.8362

### Resumen del flujo del código (C:

1. Carga y preparación de datos: Se carga el dataset y se seleccionan las características que se usarán para la predicción.
2. Normalización de los datos: Se normalizan los datos para garantizar que todas las características tengan la misma escala y mejorar el rendimiento del modelo.
3. Manejo del desbalanceo: Se usa SMOTE para crear ejemplos sintéticos y equilibrar las clases en el conjunto de datos.
4. Entrenamiento del modelo: Se utiliza HistGradientBoostingClassifier con GridSearchCV para encontrar los mejores parámetros.
5. Evaluación del modelo: Se evalúa el rendimiento del modelo en el conjunto de prueba mediante la matriz de confusión y el reporte de clasificación, y se calcula la precisión final.

## 7. Probar a utilizar el cuaderno con el dataset de los vinos blancos y realizar captura de los resultados obtenidos. (utilizar el dataset winequality-white.csv)

### 7. Aplicación al Dataset de Vinos Blancos

Se probará el modelo con winequality-white.csv con el último método que pruebo en el anterior punto: **Método de Clasificación: HistGradientBoostingClassifier**

```
[ ] 1 # IVAN FALCON MONZON
2 import pandas as pd
3 import joblib
4 from imblearn.over_sampling import SMOTE # Para balancear las clases con SMOTE
5 from sklearn.model_selection import train_test_split, GridSearchCV # Para dividir y ajustar los hiperparámetros
6 from sklearn.ensemble import HistGradientBoostingClassifier # Clasificador HistGradientBoosting
7 from sklearn.preprocessing import StandardScaler # Para normalizar los datos
8 from sklearn.metrics import confusion_matrix, classification_report # Para evaluar el modelo
9
10 # Cargar el dataset de vinos blancos desde la URL
11 dataset_url_white = "https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-white.csv"
12 df_white = pd.read_csv(dataset_url_white, sep=',')
13
14 # Selección de características relevantes para la predicción
15 selected_features = ['volatile acidity', 'total sulfur dioxide', 'density', 'sulphates', 'alcohol']
16 X_white = df_white[selected_features] # Variables predictoras
17 y_white = df_white['quality'] # Variable objetivo (calidad del vino)
18
19 # Normalización de los datos para mejorar el rendimiento del modelo
20 scaler = StandardScaler()
21 X_white_scaled = scaler.fit_transform(X_white)
22
23 # Manejo del desbalanceo con SMOTE (Generar muestras sintéticas para las clases minoritarias)
24 smote = SMOTE(random_state=42, k_neighbors=3)
25 X_resampled_white, y_resampled_white = smote.fit_resample(X_white_scaled, y_white)
26
27 # División de los datos en conjunto de entrenamiento (80%) y prueba (20%)
28 X_train_white, X_test_white, y_train_white, y_test_white = train_test_split(X_resampled_white, y_resampled_white, test_size=0.2, random_state=42)
29
30 # Definir el modelo HistGradientBoosting con búsqueda de hiperparámetros utilizando GridSearchCV
31 tuned_parameters = {
32     'learning_rate': [0.01, 0.05, 0.1], # Tasa de aprendizaje
33     'max_iter': [100, 200, 300], # Número de iteraciones (número de árboles)
34     'max_depth': [3, 5, 7] # Profundidad máxima de los árboles
35 }
36
37 # Inicializar y optimizar el modelo utilizando GridSearchCV
38 gb_model_white = GridSearchCV(HistGradientBoostingClassifier(random_state=42), tuned_parameters, cv=5, scoring='accuracy')
39
40 # Entrenar el modelo con los mejores parámetros encontrados a través de GridSearchCV
41 gb_model_white.fit(X_train_white, y_train_white)
42 print(f"Mejores parámetros encontrados para vinos blancos: {gb_model_white.best_params_}") # Imprimir los mejores parámetros
43
44 # Guardar el modelo optimizado
45 joblib.dump(gb_model_white.best_estimator_, 'modelo_vinos_blanco_avanzado.pkl')
46 print("Modelo avanzado para vinos blancos guardado exitosamente en 'modelo_vinos_blanco_avanzado.pkl'")
47
48 # Cargar el modelo optimizado
49 loaded_model_white = joblib.load('modelo_vinos_blanco_avanzado.pkl')
50 print("Modelo avanzado para vinos blancos cargado exitosamente.")
51
52 # Realizar predicciones con el modelo cargado sobre el conjunto de prueba
53 y_pred_white = loaded_model_white.predict(X_test_white)
54
55 # Generar la Matriz de Confusión para evaluar las predicciones
56 conf_matrix_white = confusion_matrix(y_test_white, y_pred_white)
57 print("\nMatriz de Confusión para Vinos Blancos:")
58 print(conf_matrix_white)
59
60 # Generar el Reporte de Clasificación, que incluye precisión, recall, f1-score, etc.
61 print("\nReporte de Clasificación para Vinos Blancos:")
62 print(classification_report(y_test_white, y_pred_white))
63
64 # Evaluación final del modelo: Precisión en el conjunto de prueba
65 accuracy_white = loaded_model_white.score(X_test_white, y_test_white)
66 print(f"\nPrecisión final del modelo avanzado para vinos blancos: {accuracy_white:.4f}")
```



## Resultado y explicación

Mejores parámetros encontrados para vinos blancos: {'learning\_rate': 0.1, 'max\_depth': 7, 'max\_iter': 300}  
Modelo avanzado para vinos blancos guardado exitosamente en 'modelo\_vinos\_blanco\_avanzado.pkl'

Modelo avanzado para vinos blancos cargado exitosamente.

Matriz de Confusión para Vinos Blancos:

```
[[425  0  1  2  0  0  0]
 [ 0 415  9  5  4  1  0]
 [ 3 22 335 79 27  2  1]
 [ 0 10 67 311 50  6  2]
 [ 4  1 15 52 351  9  0]
 [ 0  0  1  6  7 439  0]
 [ 0  0  0  0  0  0 416]]
```

Reporte de Clasificación para Vinos Blancos:

	precision	recall	f1-score	support
3	0.98	0.99	0.99	428
4	0.93	0.96	0.94	434
5	0.78	0.71	0.75	469
6	0.68	0.70	0.69	446
7	0.80	0.81	0.81	432
8	0.96	0.97	0.96	453
9	0.99	1.00	1.00	416
accuracy			0.87	3078
macro avg	0.88	0.88	0.88	3078
weighted avg	0.87	0.87	0.87	3078

Precisión final del modelo avanzado para vinos blancos: 0.8746

### Explicación del flujo del código:

1. Cargar el Dataset: Se carga el archivo CSV de vinos blancos.
2. Seleccionar Características y Normalizar: Se seleccionan las características relevantes y se normalizan con StandardScaler.
3. Balanceo de Clases: Se aplica SMOTE para balancear las clases del dataset.
4. División en Entrenamiento y Prueba: Se divide el dataset balanceado en entrenamiento y prueba.
5. Entrenar el Modelo: Se entrena un modelo HistGradientBoostingClassifier con búsqueda de hiperparámetros (GridSearchCV).
6. Guardar y Cargar el Modelo: Se guarda el modelo optimizado y luego se carga.
7. Predicciones y Evaluación: Se realizan predicciones, y se genera la matriz de confusión y el reporte de clasificación para evaluar el rendimiento del modelo.

## 8. Conclusiones de la actividad

### A) Método de Clasificación con Random Forest Classifier junto con GridSearchCV para optimizar los hiperparámetros

- Precisión: 0.6719
- Conclusión: El modelo de Random Forest con búsqueda de hiperparámetros ha alcanzado una precisión de 67%. Aunque el rendimiento no es bajo, es el menos preciso en comparación con los otros modelos. Esto podría ser debido a la complejidad de los datos.

### B) Método de Clasificación: Gradient Boosting Classifier y SMOTE

- Precisión: 0.8215
- Conclusión: Este modelo muestra una mejora con una precisión del 82%. El uso de SMOTE para balancear las clases minoritarias ha sido positivo, ayudando al modelo a generalizar mejor en las clases menos representadas. Esto sugiere que los datos del conjunto de vinos tienen un desbalance que se puede manejar bien con esta técnica.

### C) Método de Clasificación: Hist Gradient Boosting Classifier

- Precisión: 0.8362
- Conclusión: Este modelo tiene el mejor rendimiento hasta ahora con una precisión del 83.6%. HistGradientBoostingClassifier es más eficiente que el Gradient Boosting tradicional y parece ir bien con este problema, posiblemente debido a su capacidad para manejar de manera más eficiente grandes volúmenes de datos y el desbalance.

## 7. Aplicación al Dataset de Vinos Blancos

- Precisión: 0.8746
- Conclusión: El modelo de HistGradientBoostingClassifier ha alcanzado una precisión en los vinos blancos de 87.46%. Este rendim

### Conclusión General:

- El modelo HistGradientBoostingClassifier ha demostrado ser el más efectivo y robusto para estos datos, con un buen rendimiento en ambos datasets.
- El uso de SMOTE y técnicas de optimización como GridSearchCV han tenido un impacto positivo en los modelos, especialmente en el Gradient Boosting y HistGradientBoosting.
- Los resultados también indican que la aplicación de estos modelos a diferentes tipos de vinos (tintos y blancos) proporciona mejoras significativas en precisión, lo que demuestra que la técnica es adaptable y eficaz.

En resumen, Hist Gradient Boosting Classifier es la opción más recomendada para este tipo de clasificación, con una buena capacidad de adaptación a diferentes datasets y un rendimiento superior a los demás métodos.

# Repositorios

-Google Colab: <https://colab.research.google.com/drive/1wstA4BDdBcrbCLebw1eb4zKeYKj4tiCVL?usp=sharing>

-Github: [https://github.com/IvanFalconMonzon/SNS\\_ACT3\\_5\\_IvanFalconMonzon.git](https://github.com/IvanFalconMonzon/SNS_ACT3_5_IvanFalconMonzon.git)