

PROYECTO MACHINE LEARNING

MODELO DE PREDICCION :

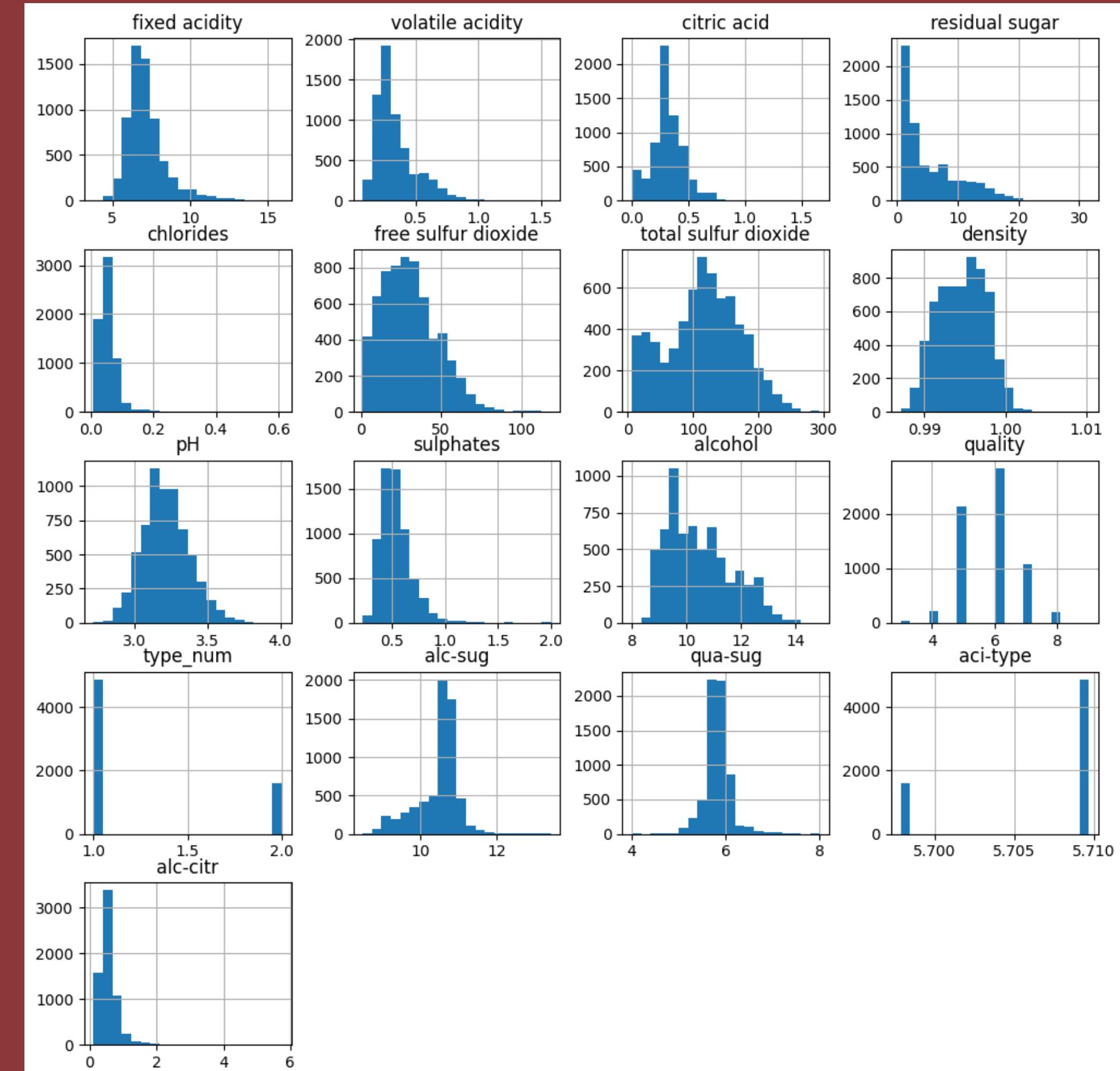
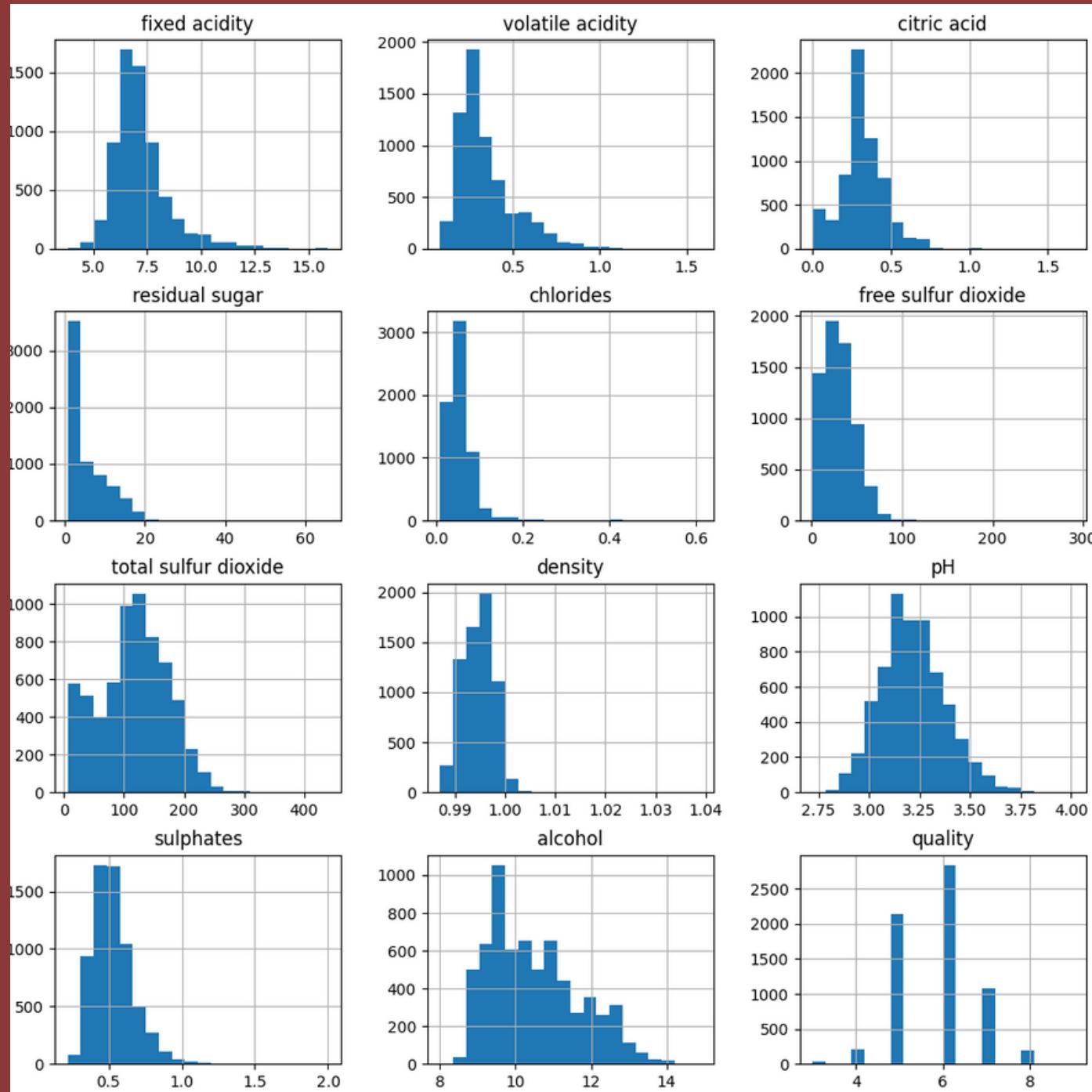
CALIDAD DEL VINO



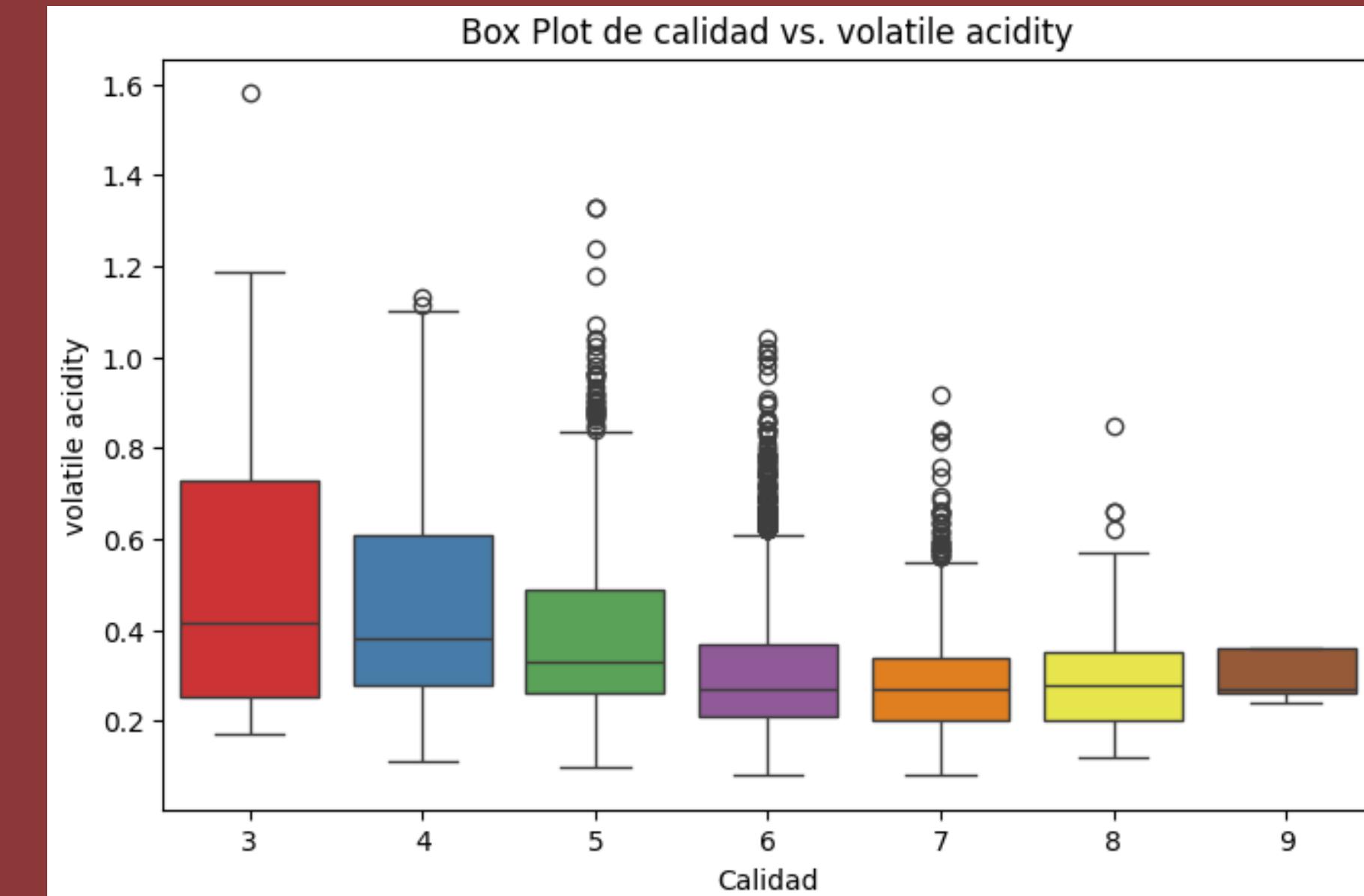
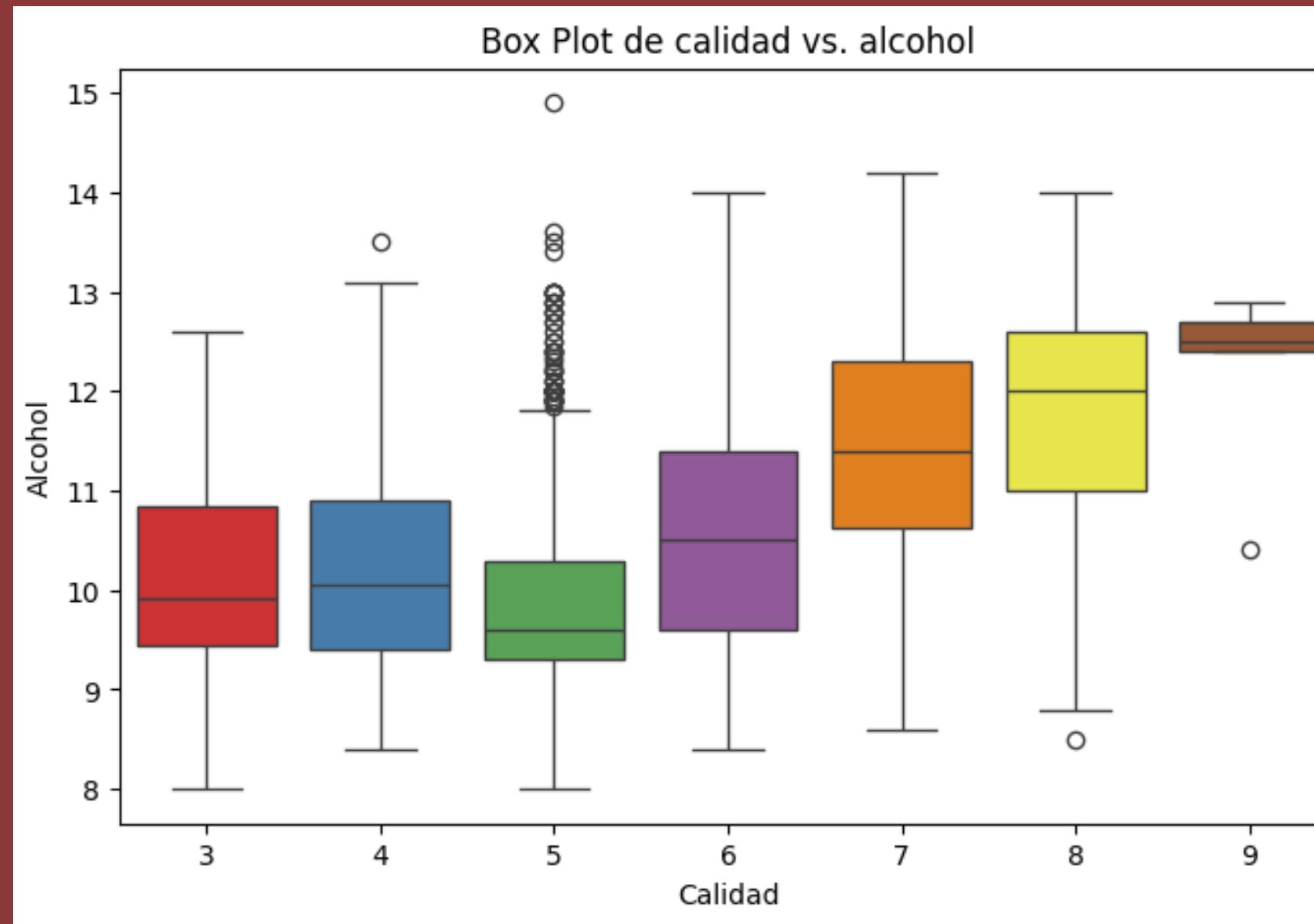
Pasos realizados para llegar al modelo final:

- Obtención de datos
- Procesamiento de datos
- Elección del target
- Estudio de sus variables y correlaciones
- Entrenamiento de nuestros modelos
- Evaluación de los datos finales
- Generación de modelo con posibilidad de actualización

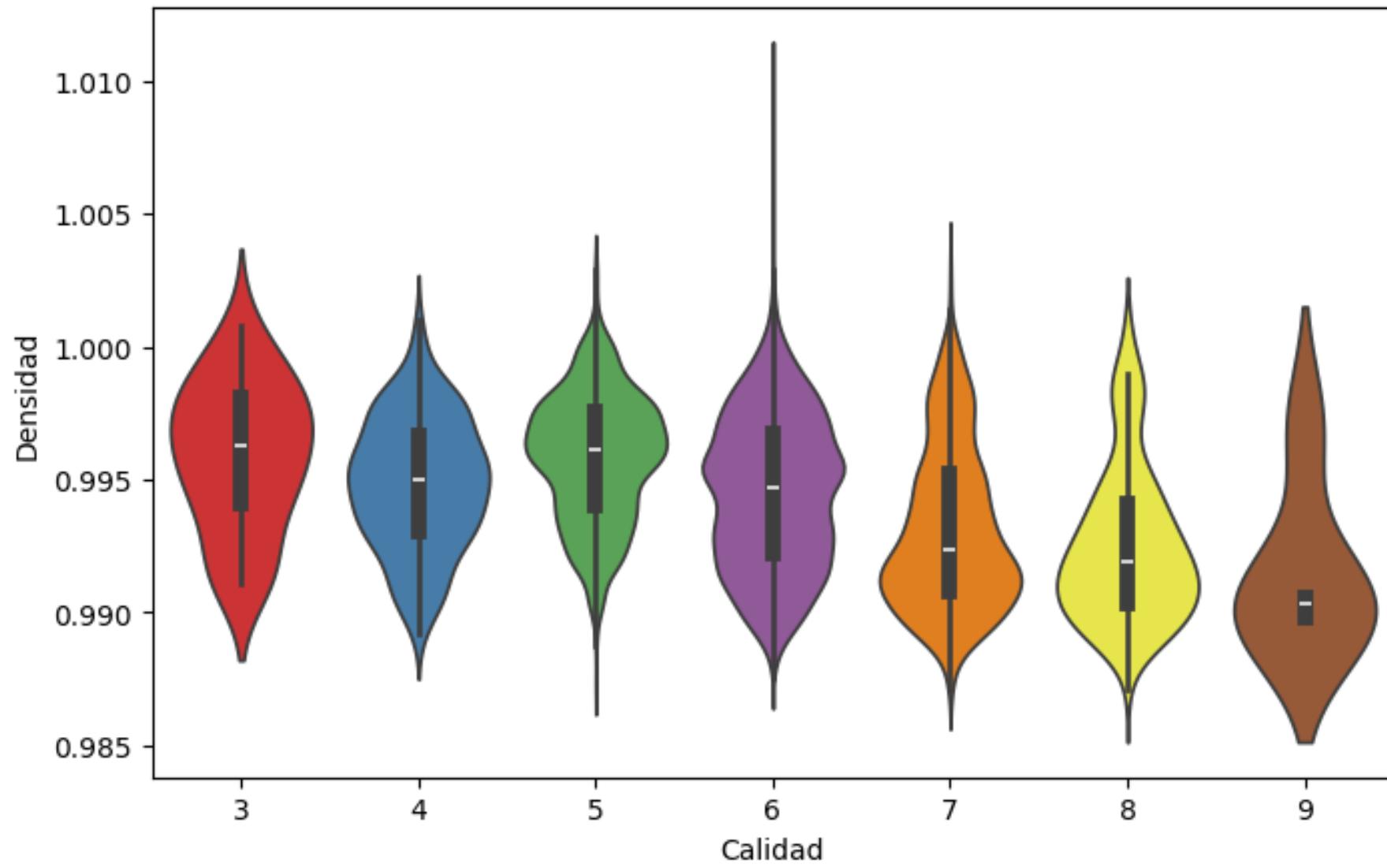
Analisis de distribución de nuestras variables :



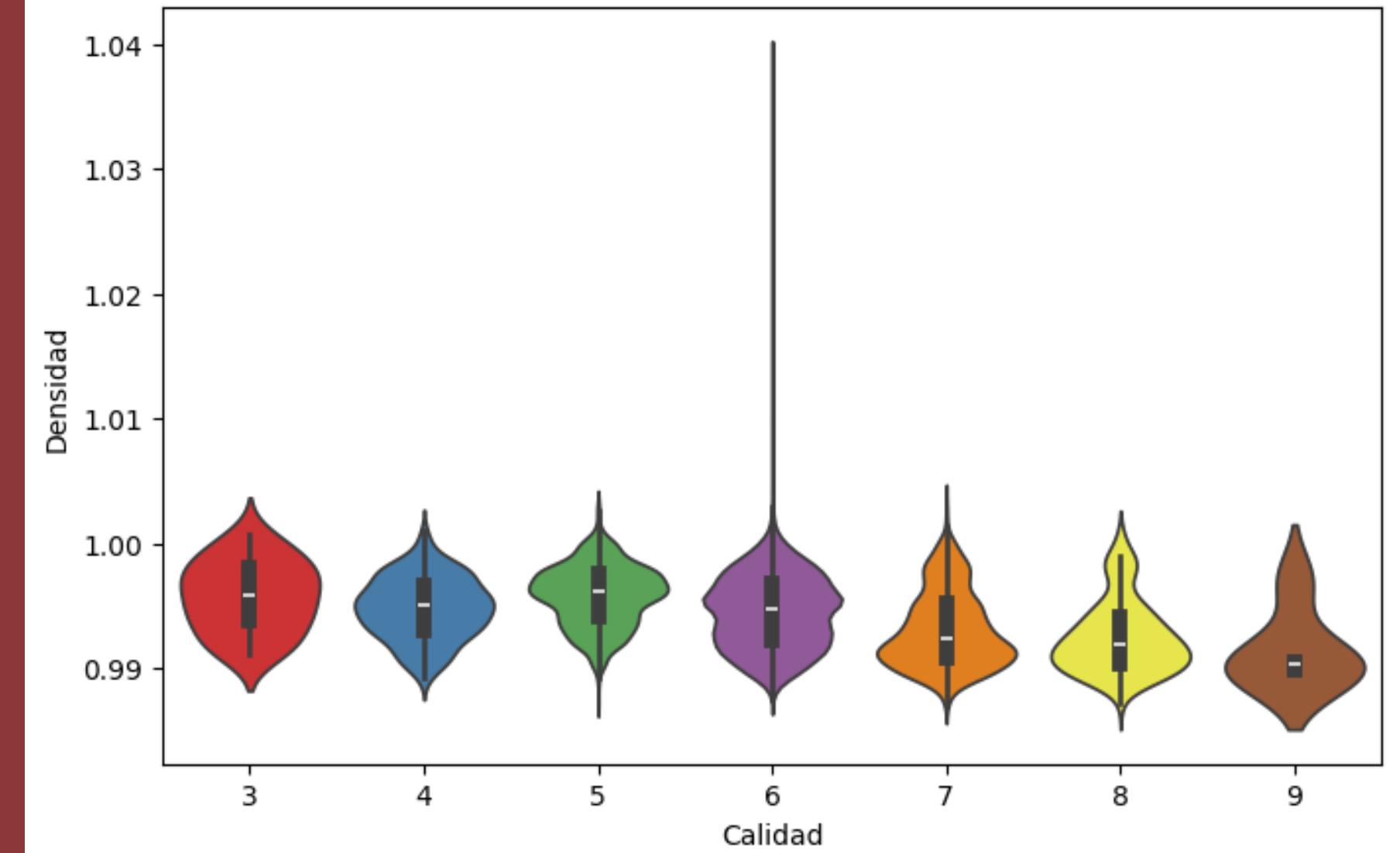
Analizamos la correlación de las variables con nuestro target :



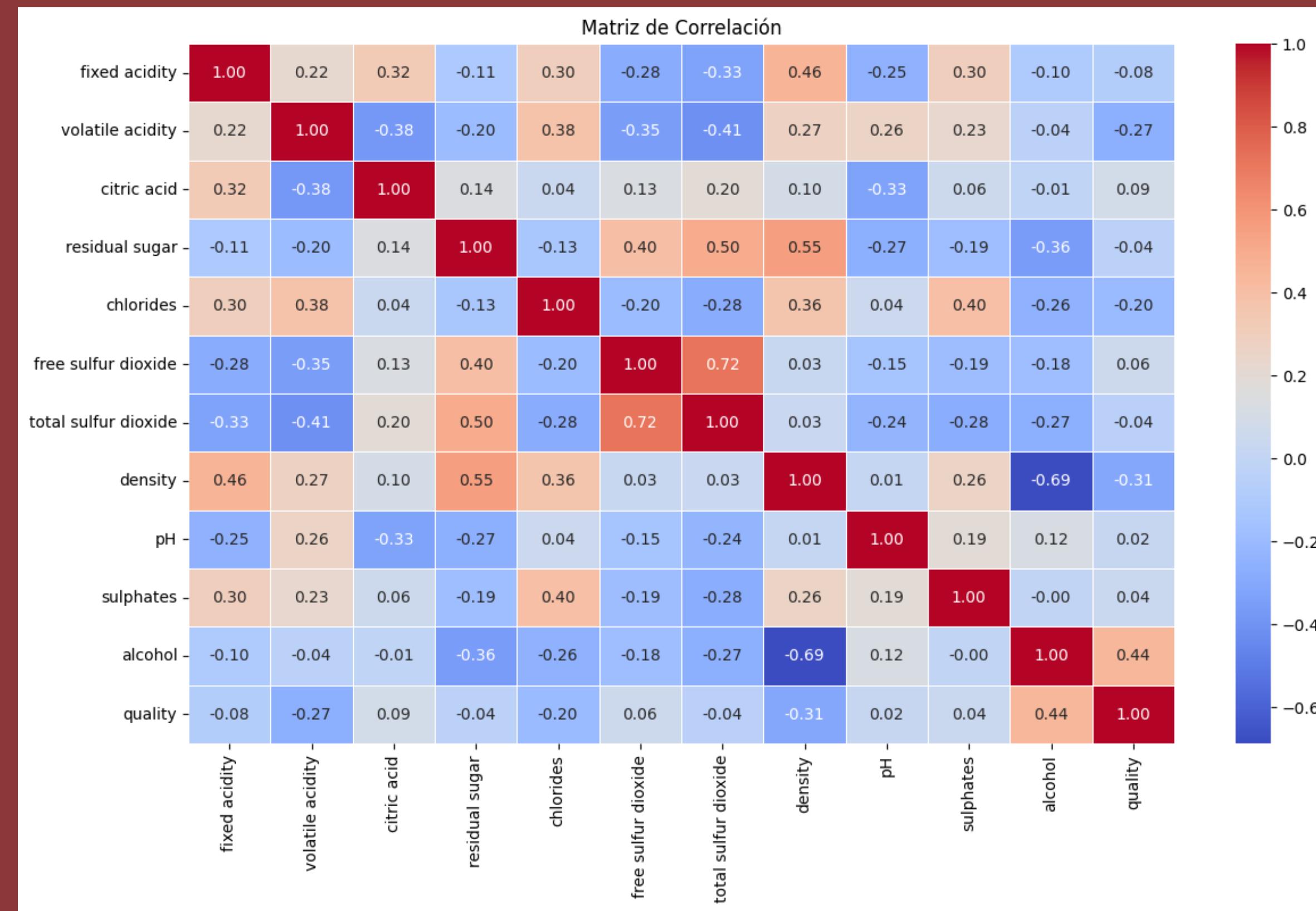
Violin Plot de calidad vs. densidad



Violin Plot de calidad vs. densidad

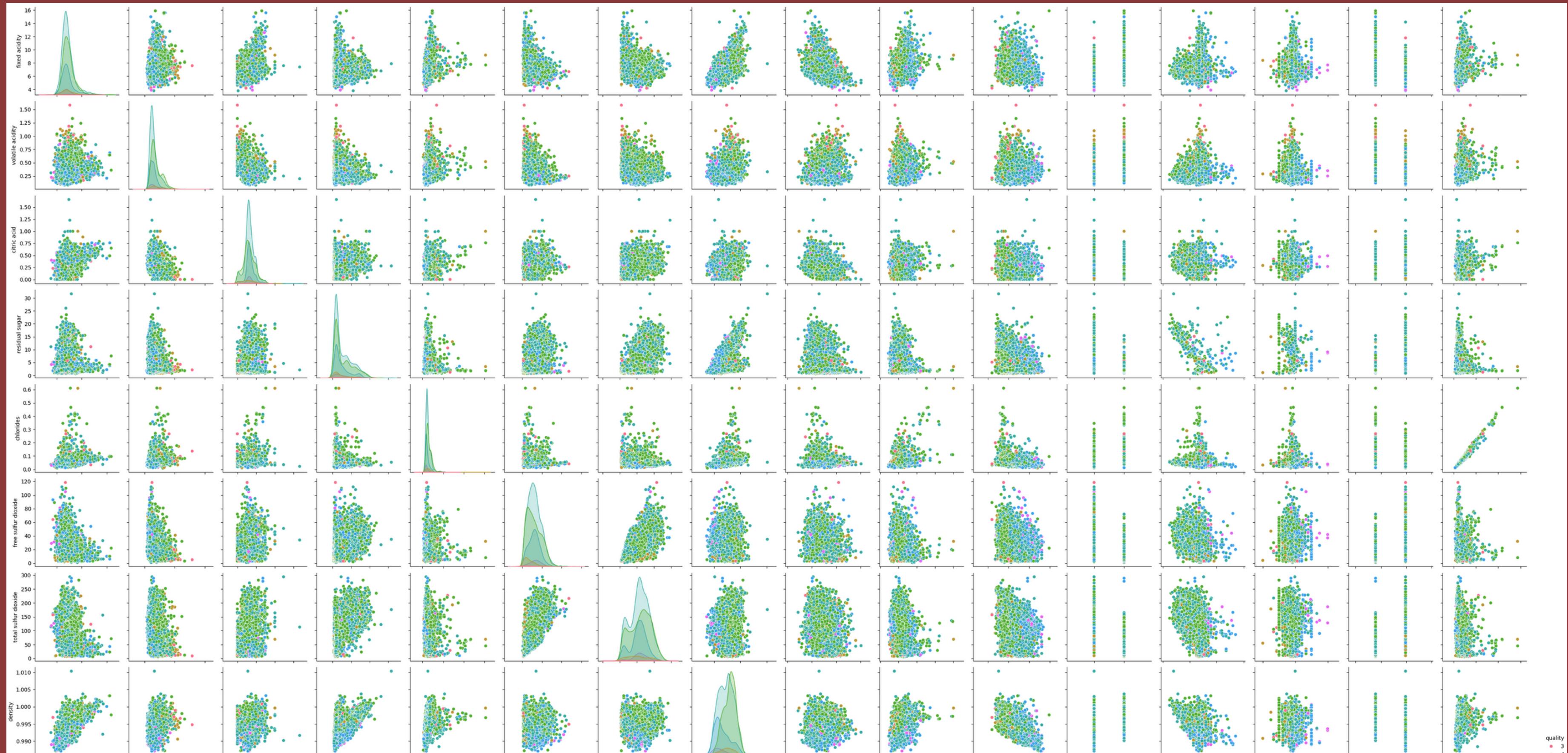


Al tener baja la correlación entre nuestras variables, buscamos relacionarlas entre ellas para conseguir un mayor apoyo a nuestro entrenamiento

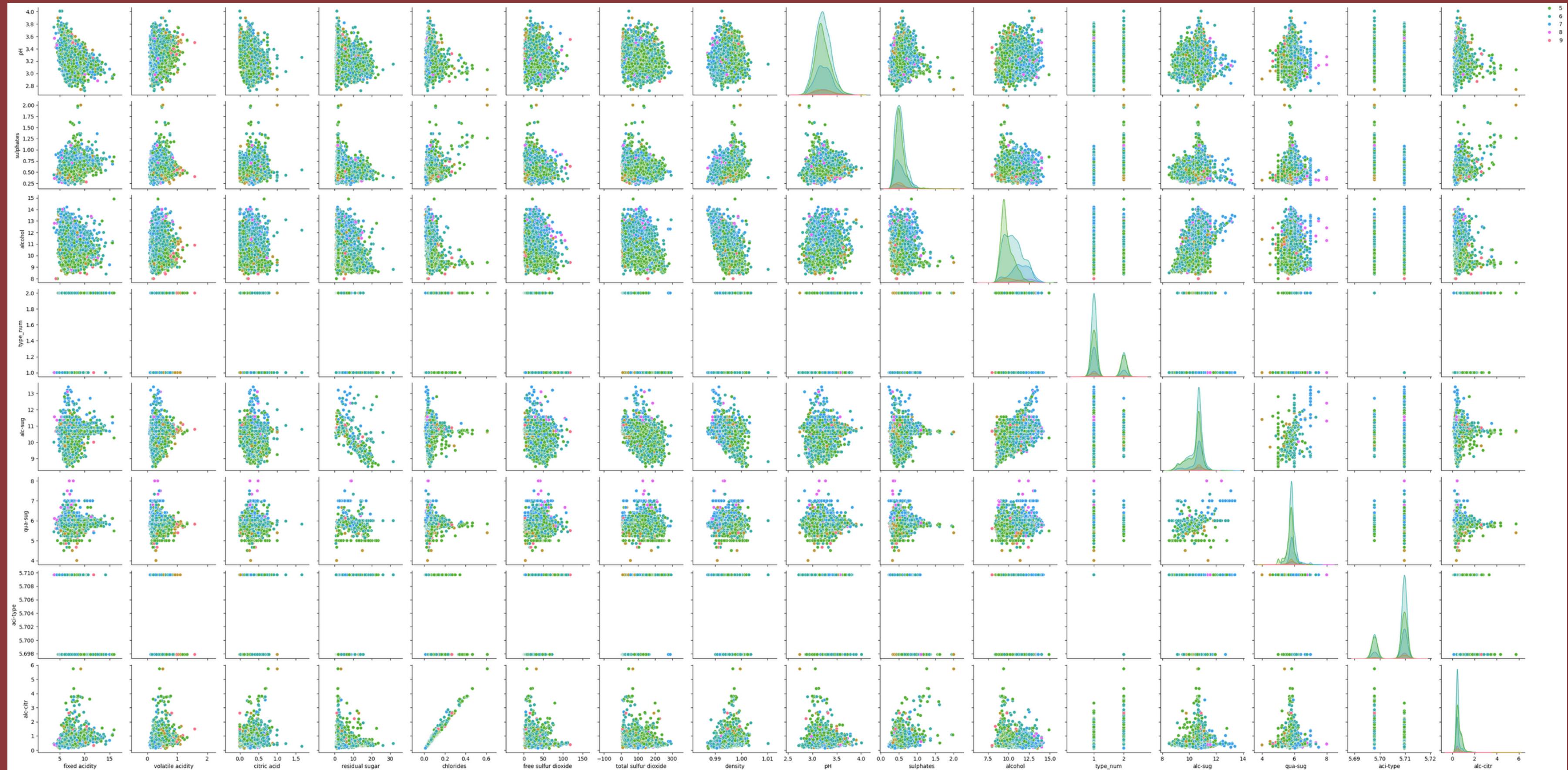


Tras realizar el Feature engineering:

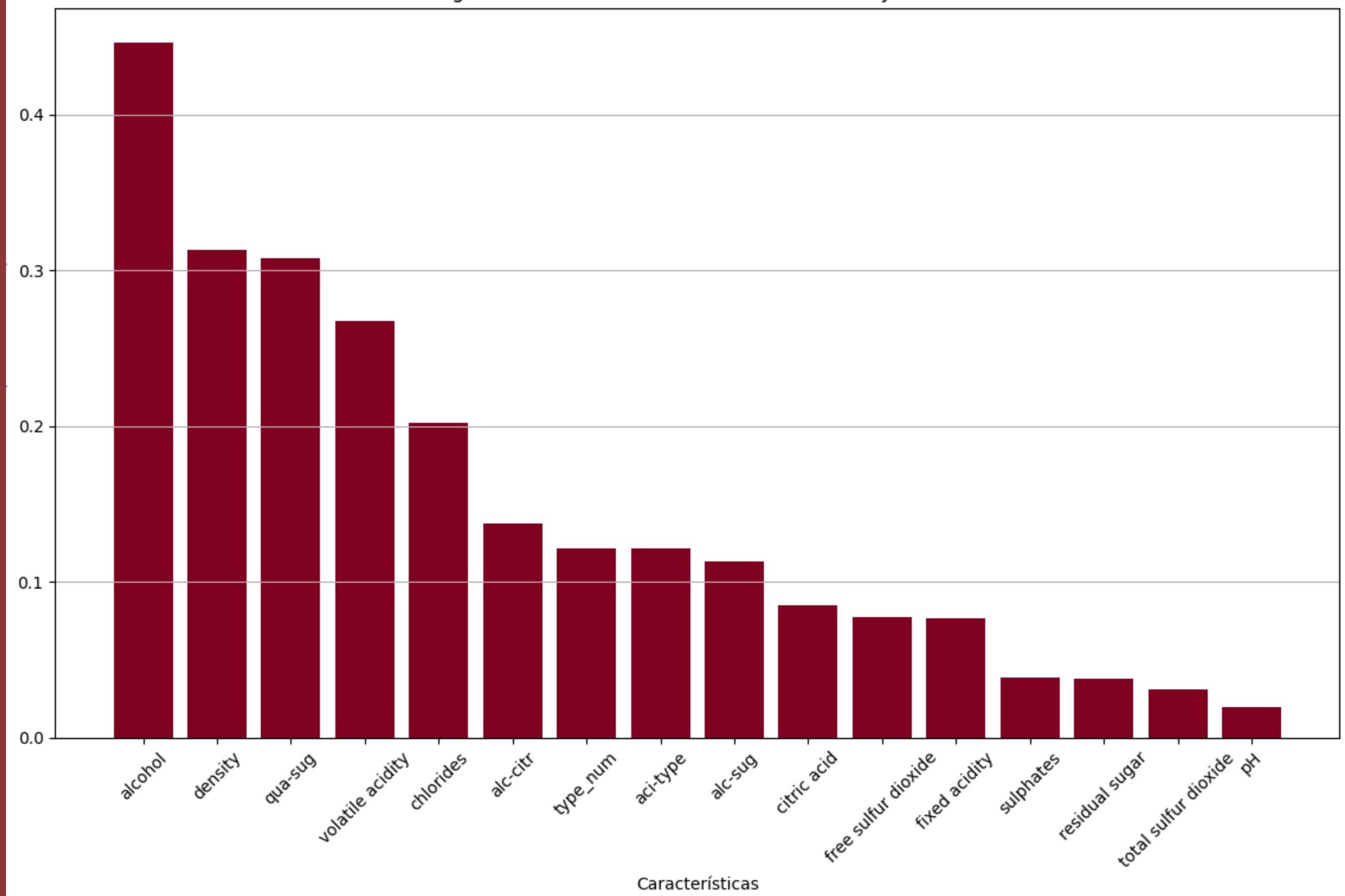
| Matriz de Correlación | | | | | | | | | | | | | | | | | |
|-----------------------|---------------|------------------|-------------|----------------|-----------|----------------|----------------------|---------|-------|-----------|---------|---------|----------|---------|---------|----------|----------|
| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality | type_num | alc-sug | qua-sug | aci-type | alc-citr |
| fixed acidity | 1.00 | 0.22 | 0.32 | -0.12 | 0.30 | -0.29 | -0.33 | 0.47 | -0.25 | 0.30 | -0.10 | -0.08 | 0.49 | 0.09 | 0.02 | -0.49 | 0.31 |
| volatile acidity | 0.22 | 1.00 | -0.38 | -0.21 | 0.38 | -0.36 | -0.42 | 0.27 | 0.26 | 0.23 | -0.04 | -0.27 | 0.65 | 0.16 | -0.03 | -0.65 | 0.39 |
| citric acid | 0.32 | -0.38 | 1.00 | 0.14 | 0.04 | 0.14 | 0.20 | 0.09 | -0.33 | 0.06 | -0.01 | 0.08 | -0.19 | -0.09 | -0.01 | 0.19 | 0.04 |
| residual sugar | -0.12 | -0.21 | 0.14 | 1.00 | -0.13 | 0.43 | 0.51 | 0.54 | -0.27 | -0.19 | -0.37 | -0.04 | -0.35 | -0.78 | -0.12 | 0.35 | -0.20 |
| chlorides | 0.30 | 0.38 | 0.04 | -0.13 | 1.00 | -0.20 | -0.28 | 0.37 | 0.04 | 0.40 | -0.26 | -0.20 | 0.51 | 0.06 | -0.02 | -0.51 | 0.98 |
| sulfur dioxide | -0.29 | -0.36 | 0.14 | 0.43 | -0.20 | 1.00 | 0.72 | 0.03 | -0.16 | -0.19 | -0.18 | 0.08 | -0.48 | -0.35 | -0.05 | 0.48 | -0.25 |
| total sulfur dioxide | -0.33 | -0.42 | 0.20 | 0.51 | -0.28 | 0.72 | 1.00 | 0.03 | -0.24 | -0.28 | -0.27 | -0.03 | -0.70 | -0.42 | -0.08 | 0.70 | -0.35 |
| density | 0.47 | 0.27 | 0.09 | 0.54 | 0.37 | 0.03 | 0.03 | 1.00 | 0.01 | 0.26 | -0.70 | -0.31 | 0.40 | -0.52 | -0.13 | -0.40 | 0.29 |
| pH | -0.25 | 0.26 | -0.33 | -0.27 | 0.04 | -0.16 | -0.24 | 0.01 | 1.00 | 0.19 | 0.12 | 0.02 | 0.33 | 0.20 | 0.02 | -0.33 | 0.08 |
| sulphates | 0.30 | 0.23 | 0.06 | -0.19 | 0.40 | -0.19 | -0.28 | 0.26 | 0.19 | 1.00 | -0.00 | 0.04 | 0.49 | 0.10 | 0.02 | -0.49 | 0.42 |
| alcohol | -0.10 | -0.04 | -0.01 | -0.37 | -0.26 | -0.18 | -0.27 | -0.70 | 0.12 | -0.00 | 1.00 | 0.45 | -0.03 | 0.47 | 0.17 | 0.03 | -0.10 |
| quality | -0.08 | -0.27 | 0.08 | -0.04 | -0.20 | 0.08 | -0.03 | -0.31 | 0.02 | 0.04 | 0.45 | 1.00 | -0.12 | 0.11 | 0.31 | 0.12 | -0.14 |
| type_num | 0.49 | 0.65 | -0.19 | -0.35 | 0.51 | -0.48 | -0.70 | 0.40 | 0.33 | 0.49 | -0.03 | -0.12 | 1.00 | 0.24 | 0.02 | -1.00 | 0.55 |
| alc-sug | 0.09 | 0.16 | -0.09 | -0.78 | 0.06 | -0.35 | -0.42 | -0.52 | 0.20 | 0.10 | 0.47 | 0.11 | 0.24 | 1.00 | 0.37 | -0.24 | 0.14 |
| qua-sug | 0.02 | -0.03 | -0.01 | -0.12 | -0.02 | -0.05 | -0.08 | -0.13 | 0.02 | 0.02 | 0.17 | 0.31 | 0.02 | 0.37 | 1.00 | -0.02 | 0.01 |
| aci-type | -0.49 | -0.65 | 0.19 | 0.35 | -0.51 | 0.48 | 0.70 | -0.40 | -0.33 | -0.49 | 0.03 | 0.12 | -1.00 | -0.24 | -0.02 | 1.00 | -0.55 |
| alc-citr | 0.31 | 0.39 | 0.04 | -0.20 | 0.98 | -0.25 | -0.35 | 0.29 | 0.08 | 0.42 | -0.10 | -0.14 | 0.55 | 0.14 | 0.01 | -0.55 | 1.00 |



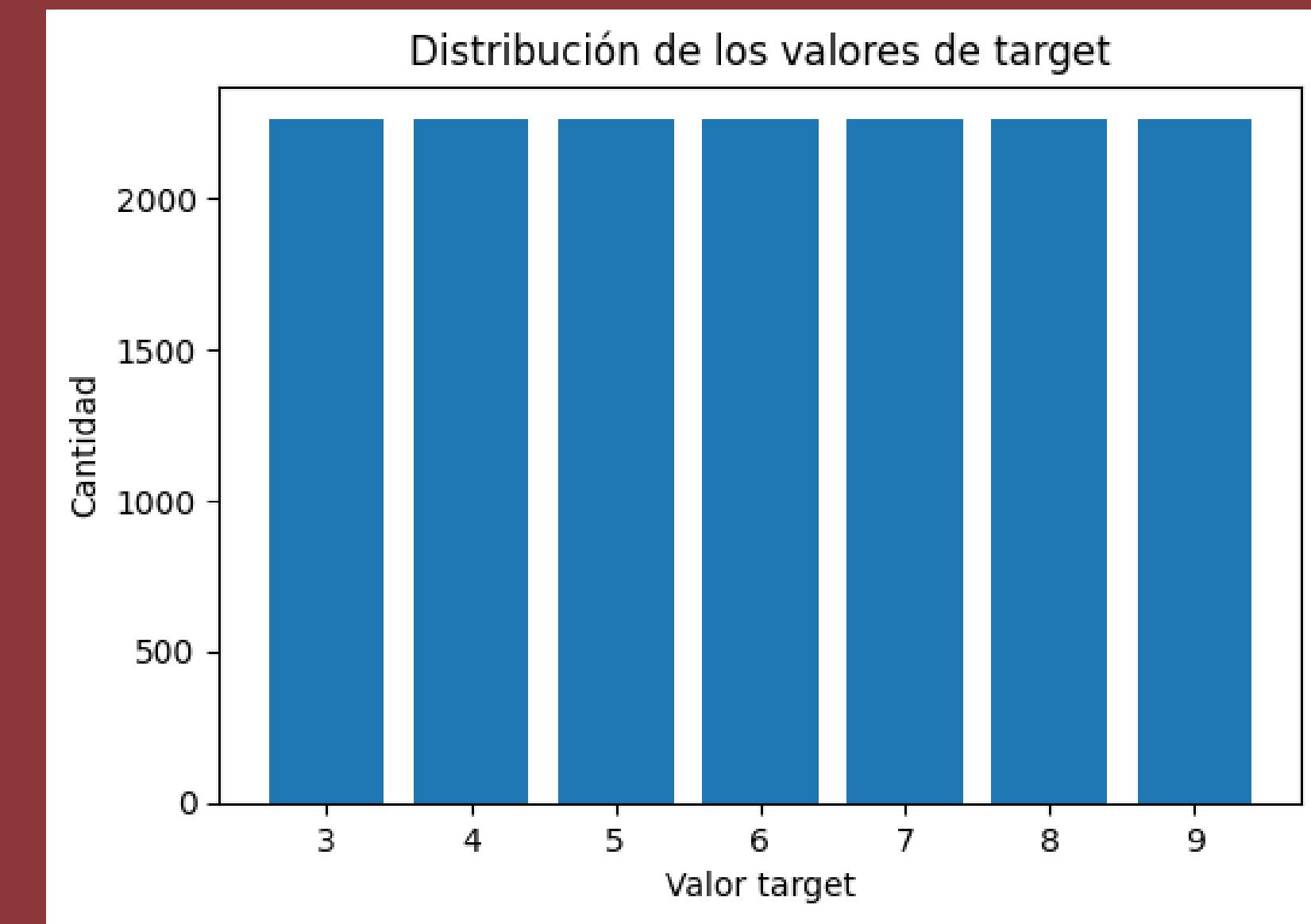
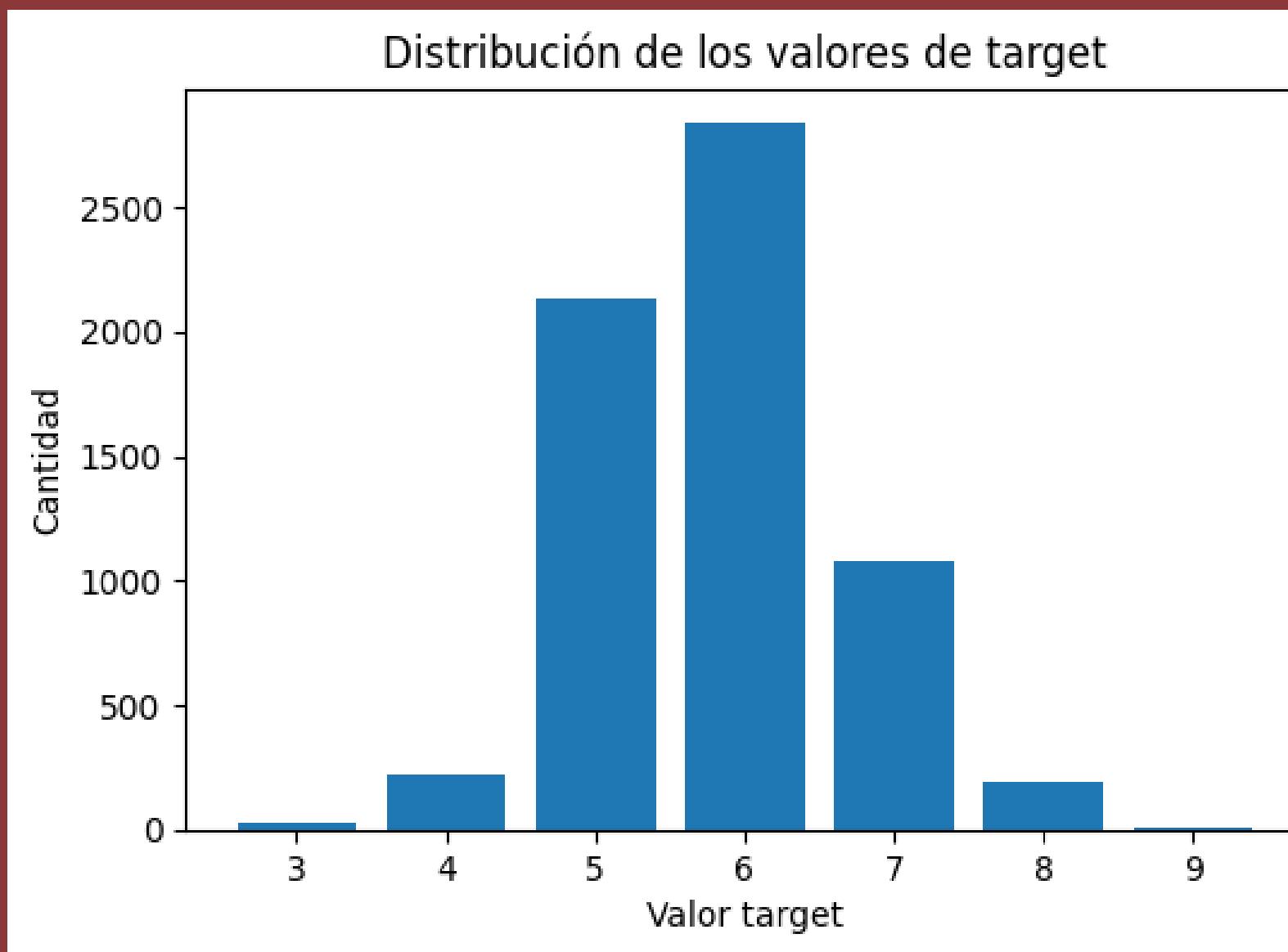
1



Magnitud de Correlación entre Características y Calidad



Evaluamos haciendo un RandomOverSampler para que pueda aprender bien de todos nuestros target



Comenzamos con el entrenamiento de nuestros modelos :

En primer lugar afronte mi problema como clasificación, para que me diera una clasificación del 3 al 9 en función de su calidad.

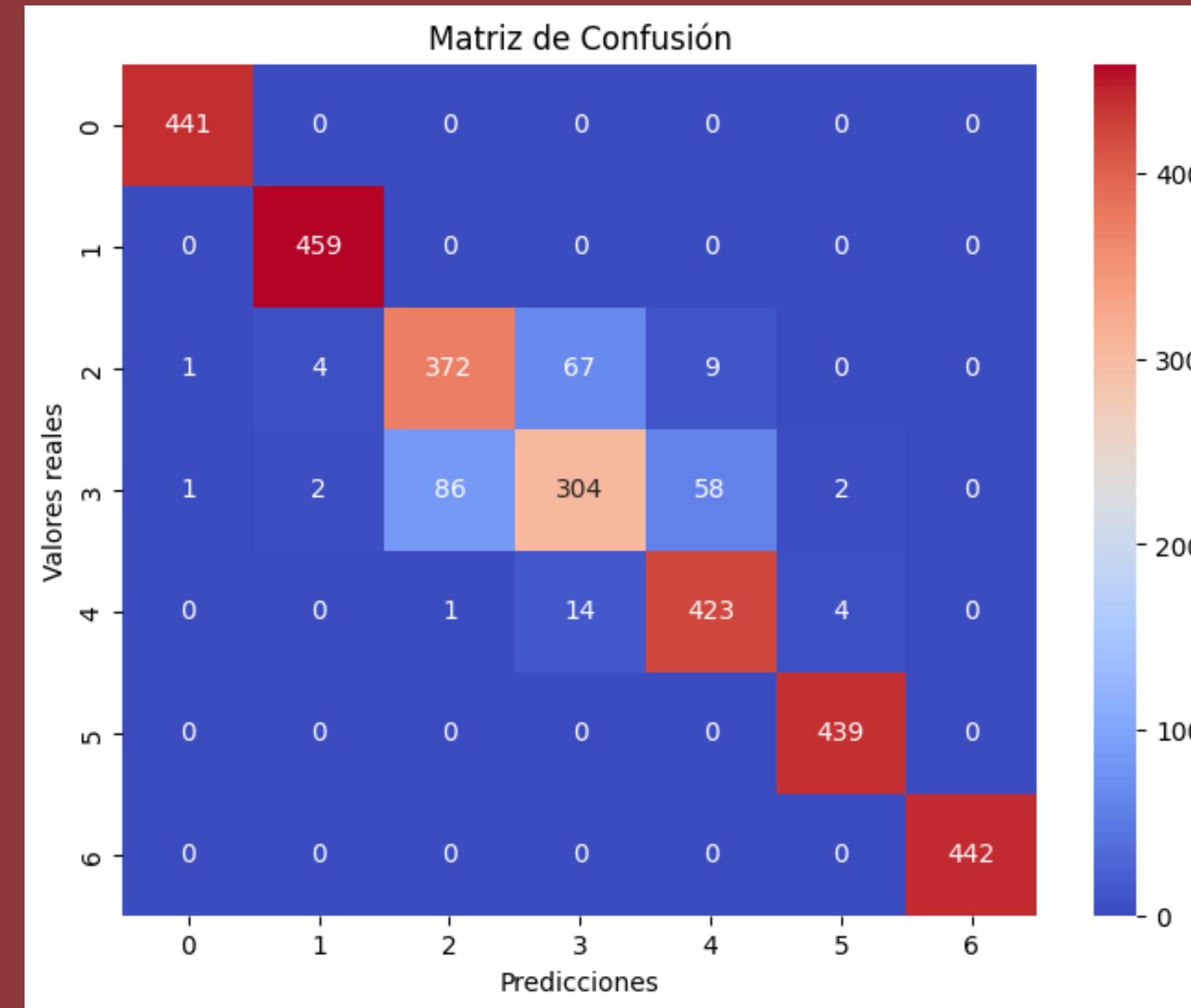
Una vez hecho esto , entendí que seria mas aconsejable obtener un valor real, mas que una aproximación, por lo que después pase a los modelos de regresión.

A continuación os enseñare mis modelos de clasificación y sus resultados obtenidos

Modelos de clasificación

Gradient Boosting

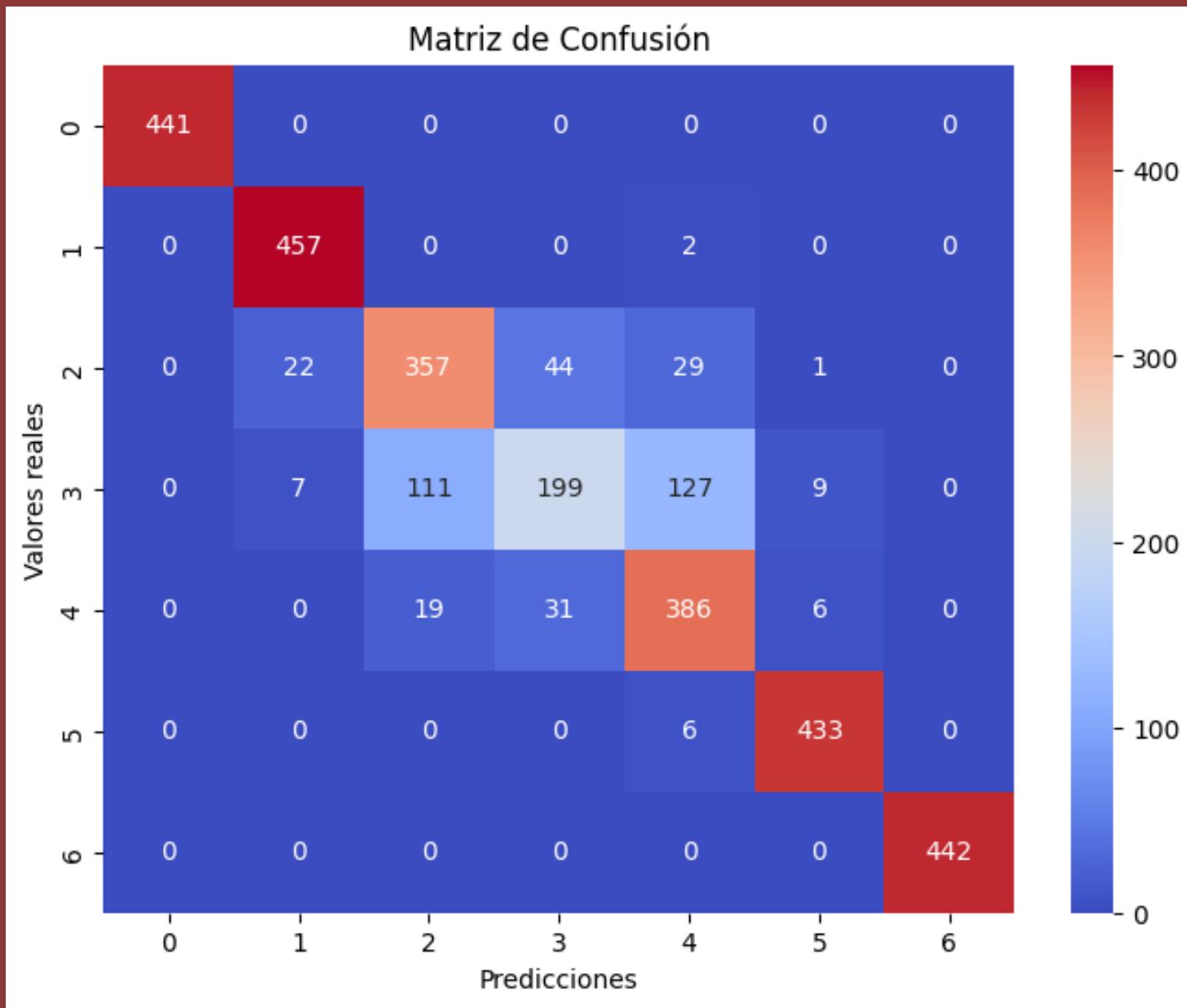
El mejor modelo de clasificación



Tasa de aciertos :
92.04%

classifier_learning_rate: 0.2,
'classifier_max_depth': 6,
'classifier_n_estimators': 150

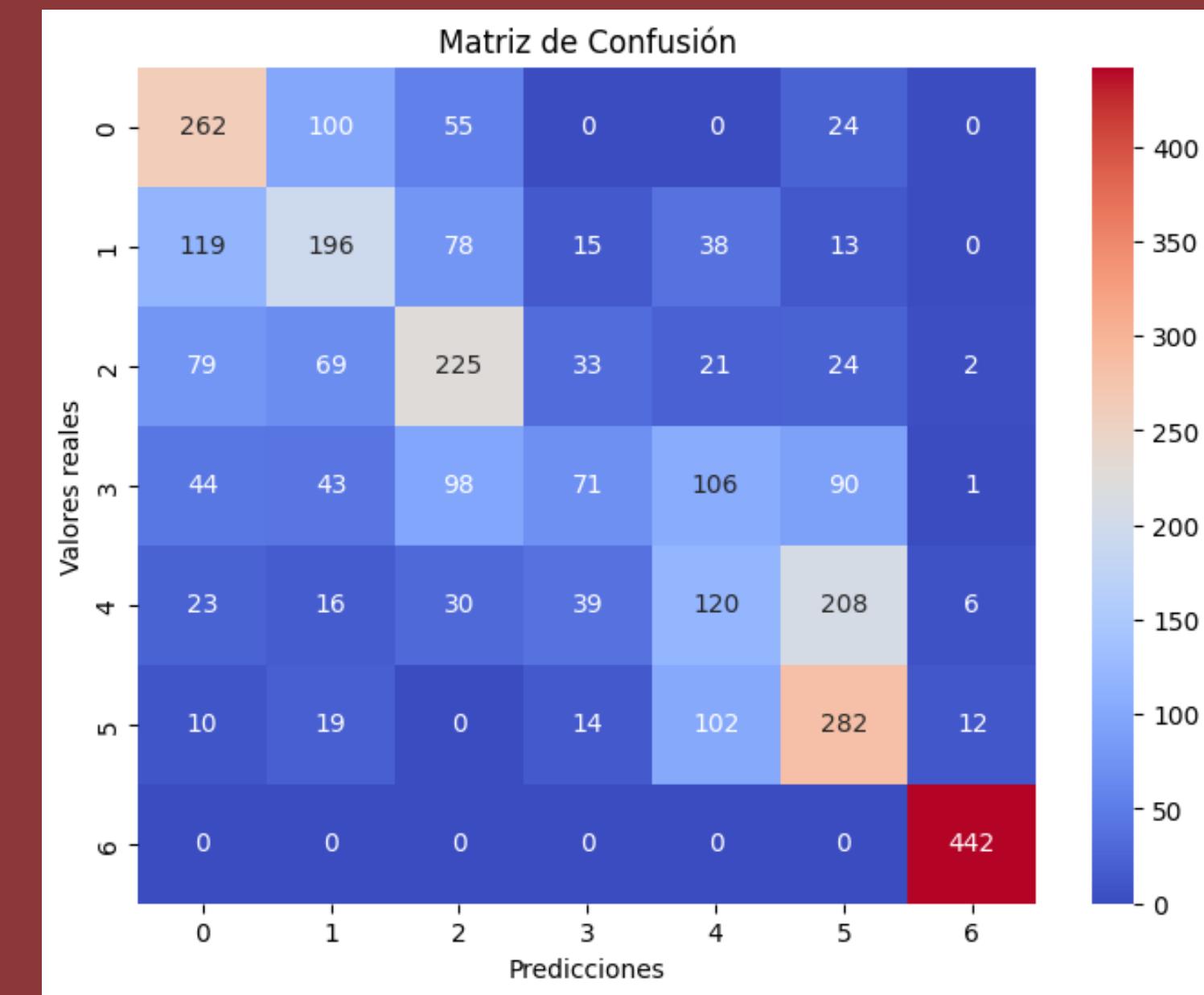
Random Forest



Tasa de aciertos 86.77%

classifier': RandomForestClassifier(),
'classifier_max_depth': 10,
'classifier_max_features': 1,
'classifier_n_estimators': 200

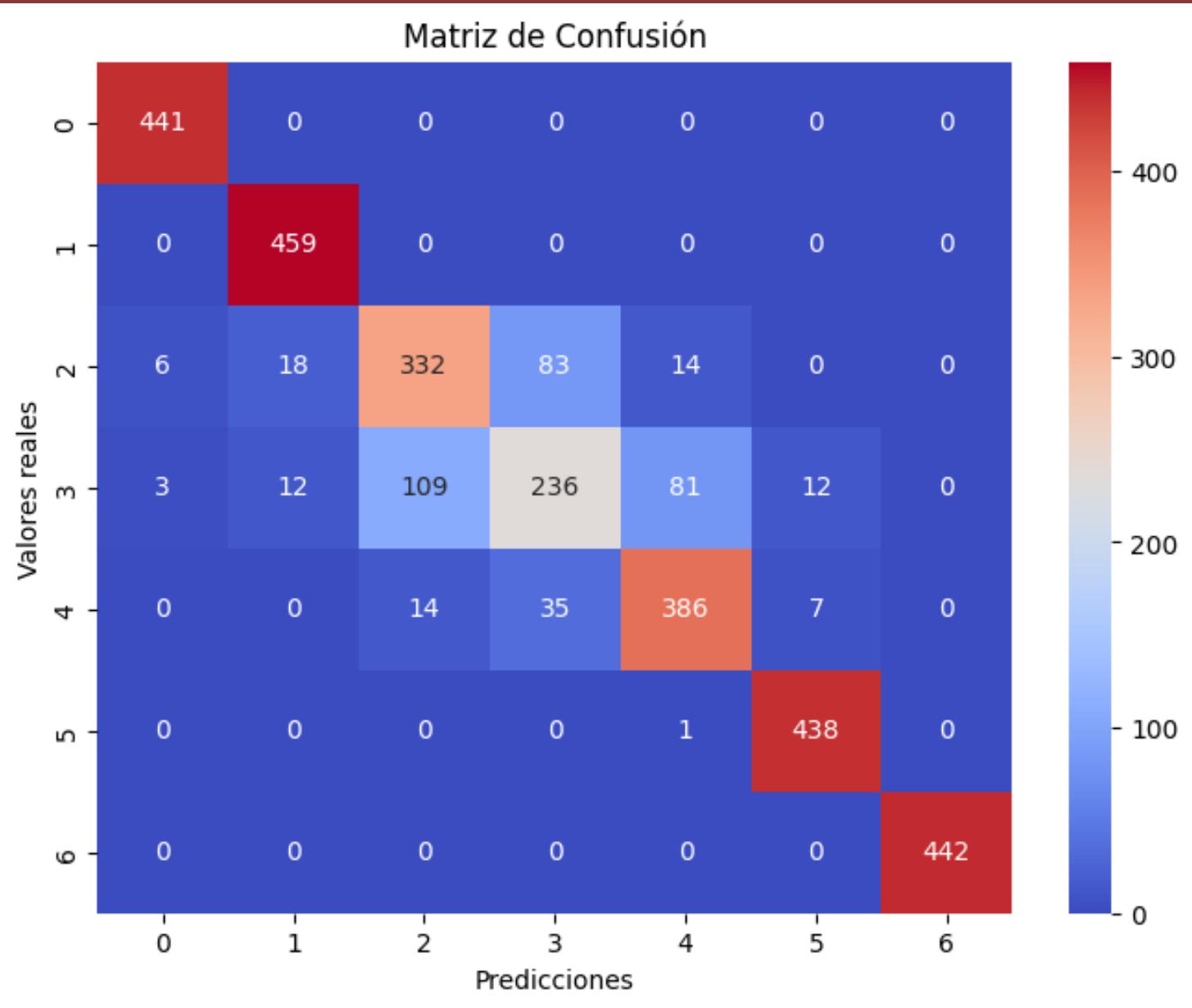
Logistic Regression



Tasa de aciertos : 51.07%

classifier_C': 10,
classifier_max_iter: 100,
classifier_penalty: 'l1'
classifier_solver: 'liblinear'

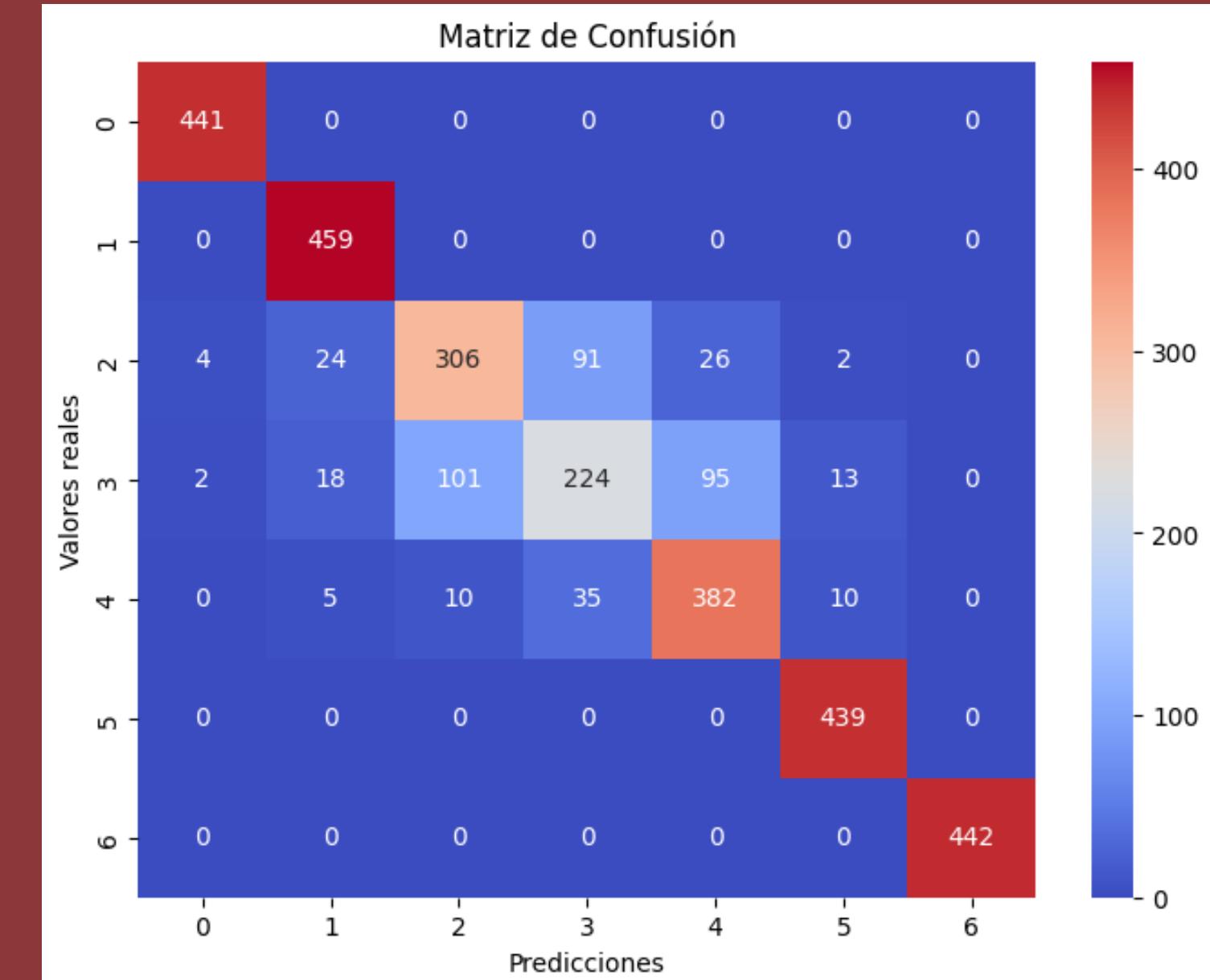
SVC



Tasa de Aciertos 87.38%

'classifier': SVC(),
'classifier__C': 200

KNN



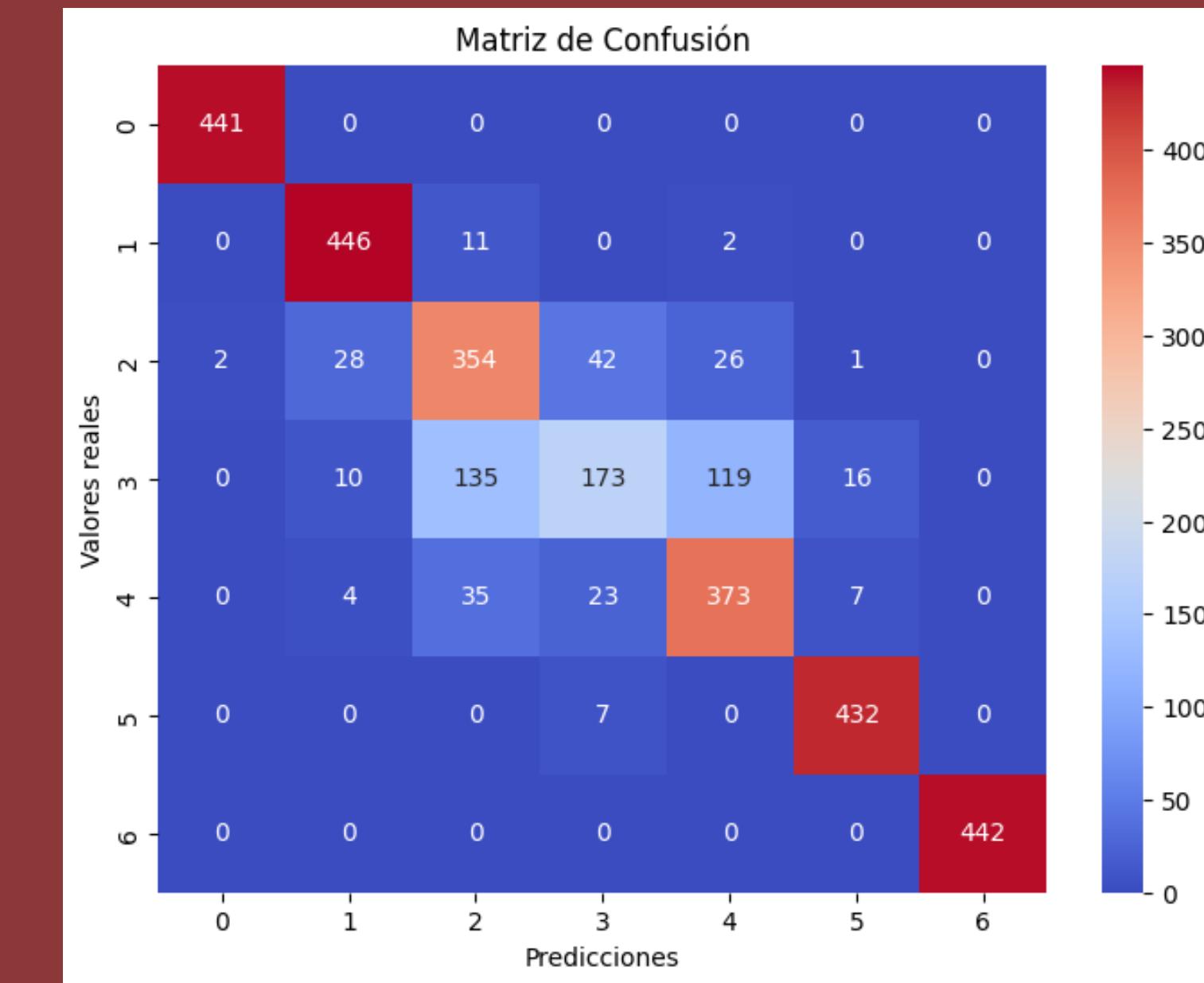
Tasa de aciertos 86.07%

'classifier': KNeighborsClassifier(),
'classifier__n_neighbors': 3

No supervisado PCA con Random Forest

Tasa de acierto 85.04%

```
classifier_max_depth': 10,  
'classifier_min_samples_leaf': 2,  
'classifier_n_estimators': 200,  
'pca_n_components': 10,  
'scaler': StandardScaler()
```

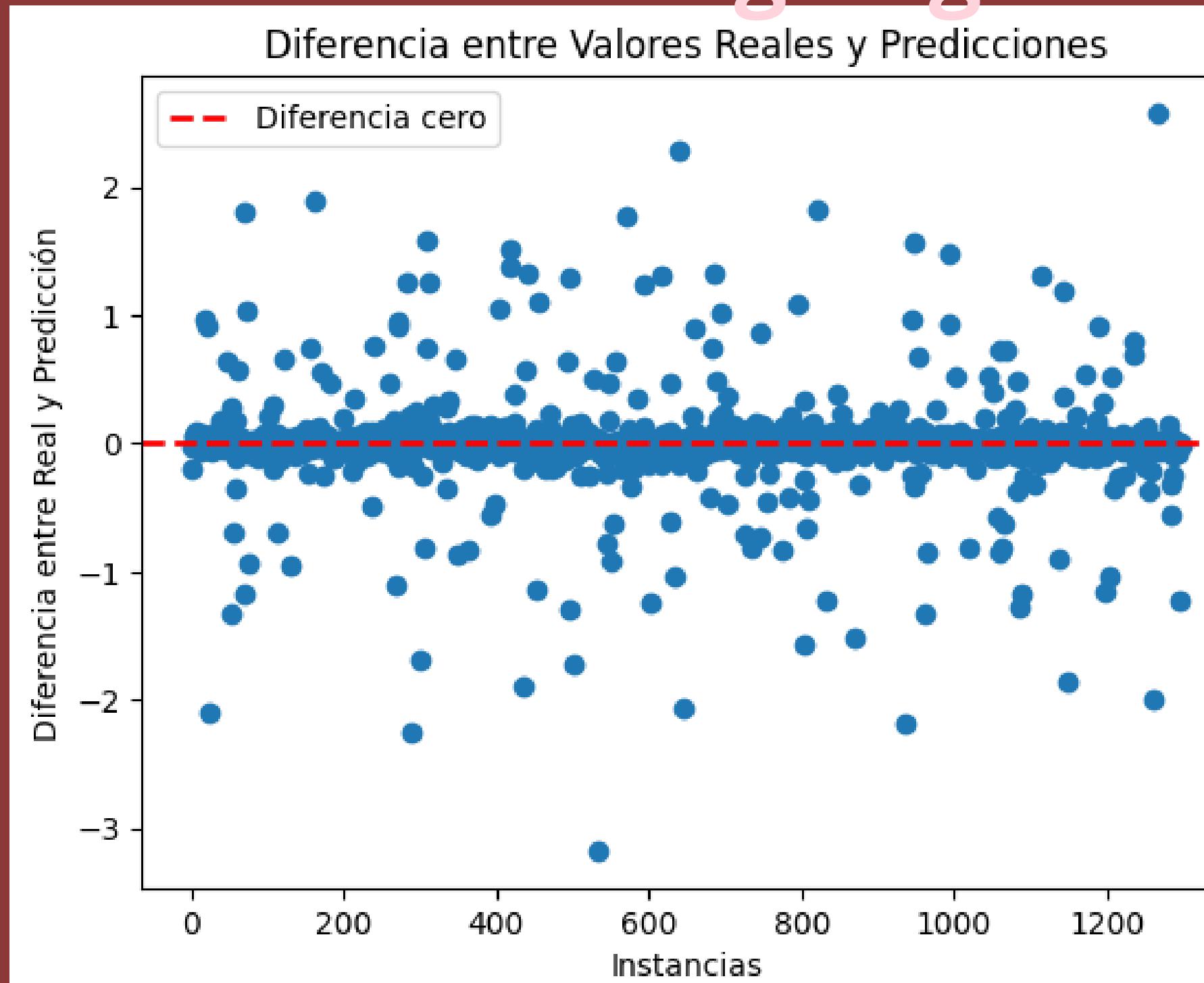


Regression

Una vez probados los modelos de clasificación, decidí probar con los modelos de regresión para ver otra forma para afrontar mi modelo.

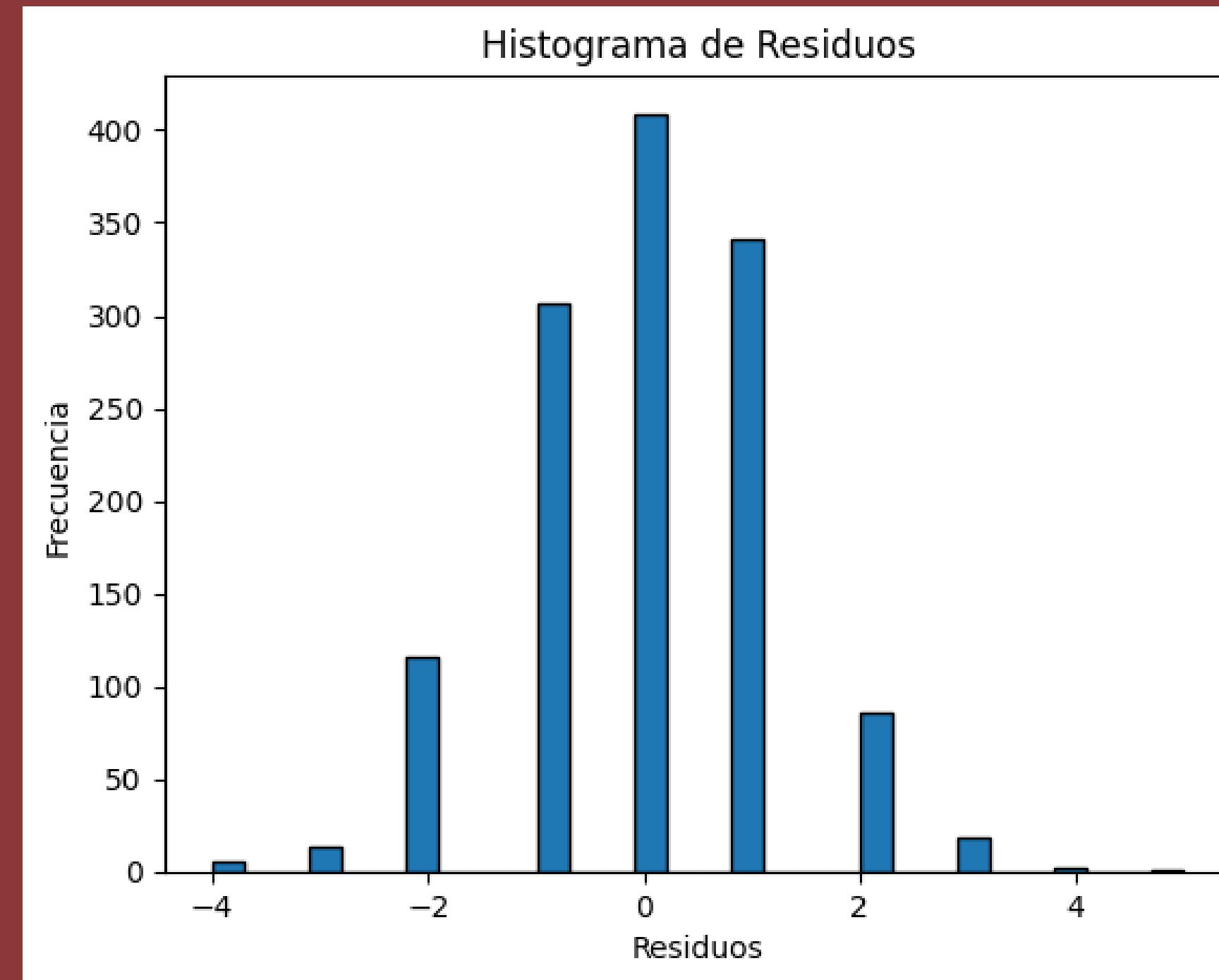
A continuación paso a enseñaros mi mejor modelo de regresión y el modelo final

Gradient Boosting Regression



Mean Absolute Error (MAE): 0.1437630517740377
Mean Absolute Percentage Error (MAPE): 2.622878504956207%
Mean Squared Error (MSE): 0.13138993696312906
Root Mean Squared Error (RMSE): 0.3624774985611232

Decission Tree Regression



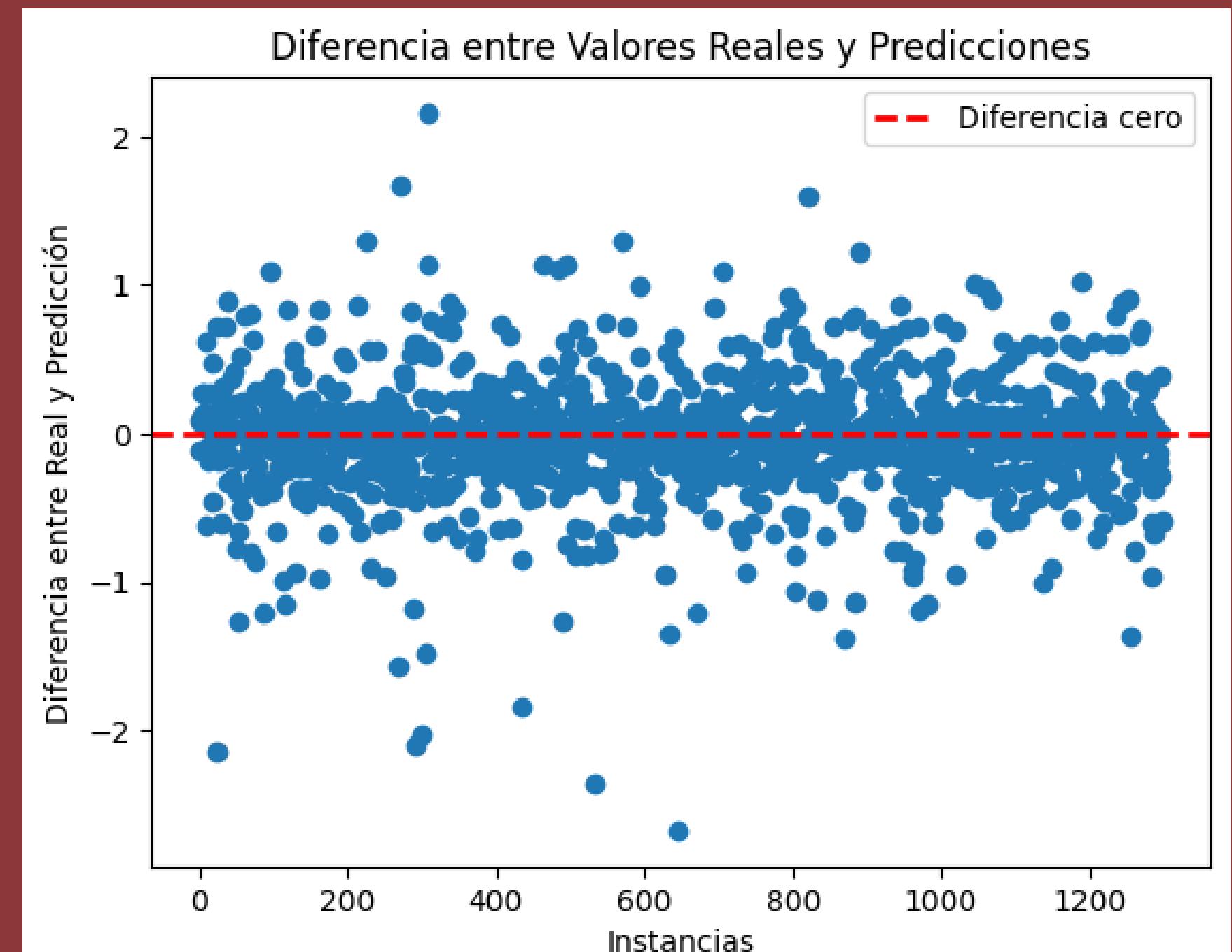
Mean Absolute Error (MAE): 0.9090909090909091

Mean Absolute Percentage Error (MAPE): 15.890105657054809%

Mean Squared Error (MSE): 1.448382126348228

Root Mean Squared Error (RMSE): 1.2034874849154966

Random Forest Regression



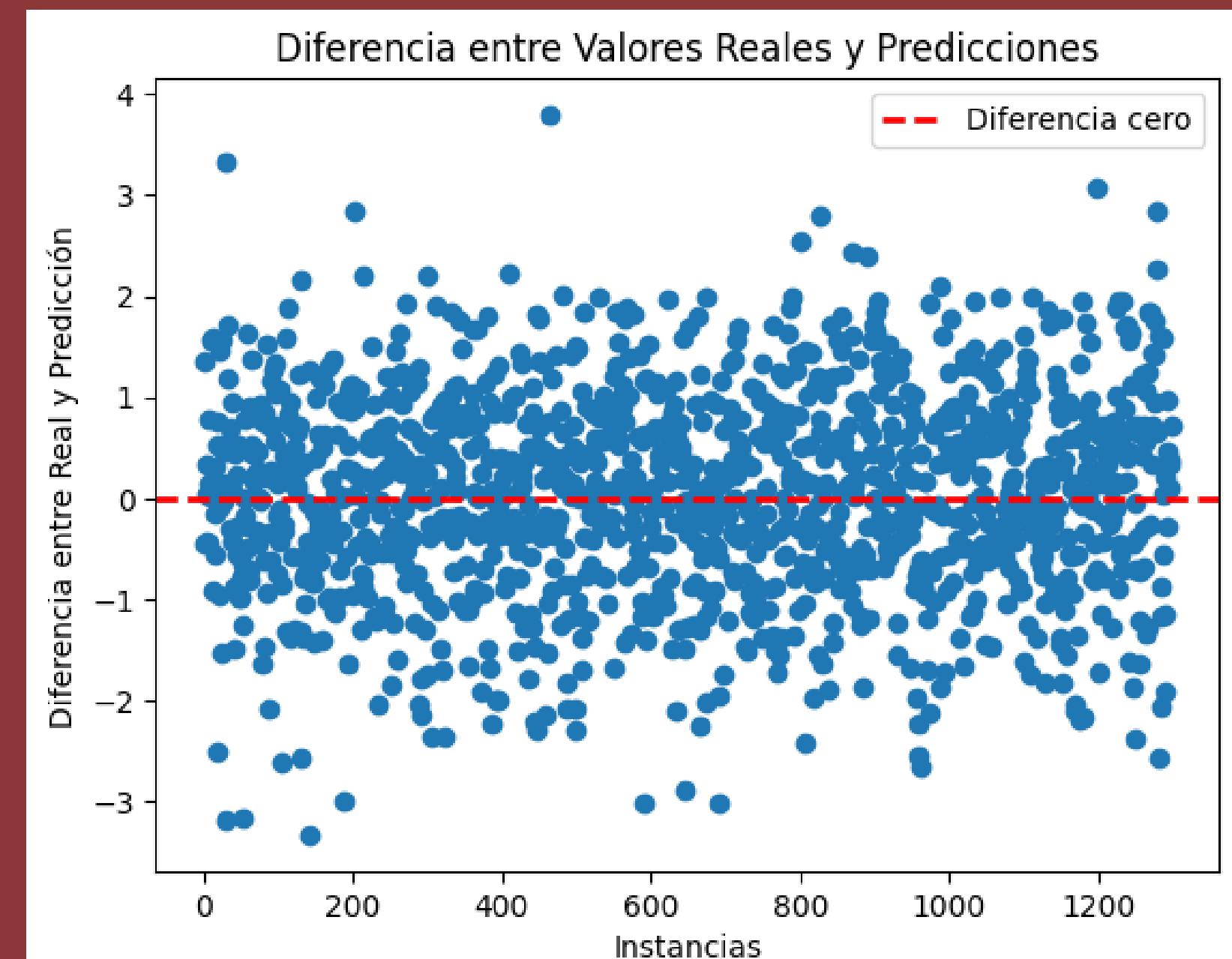
Mean Absolute Error (MAE): 0.26286238245014526

Mean Absolute Percentage Error (MAPE): 4.7258278744264395%

Mean Squared Error (MSE): 0.1605248141611384

Root Mean Squared Error (RMSE): 0.40065548063284795

LASSO



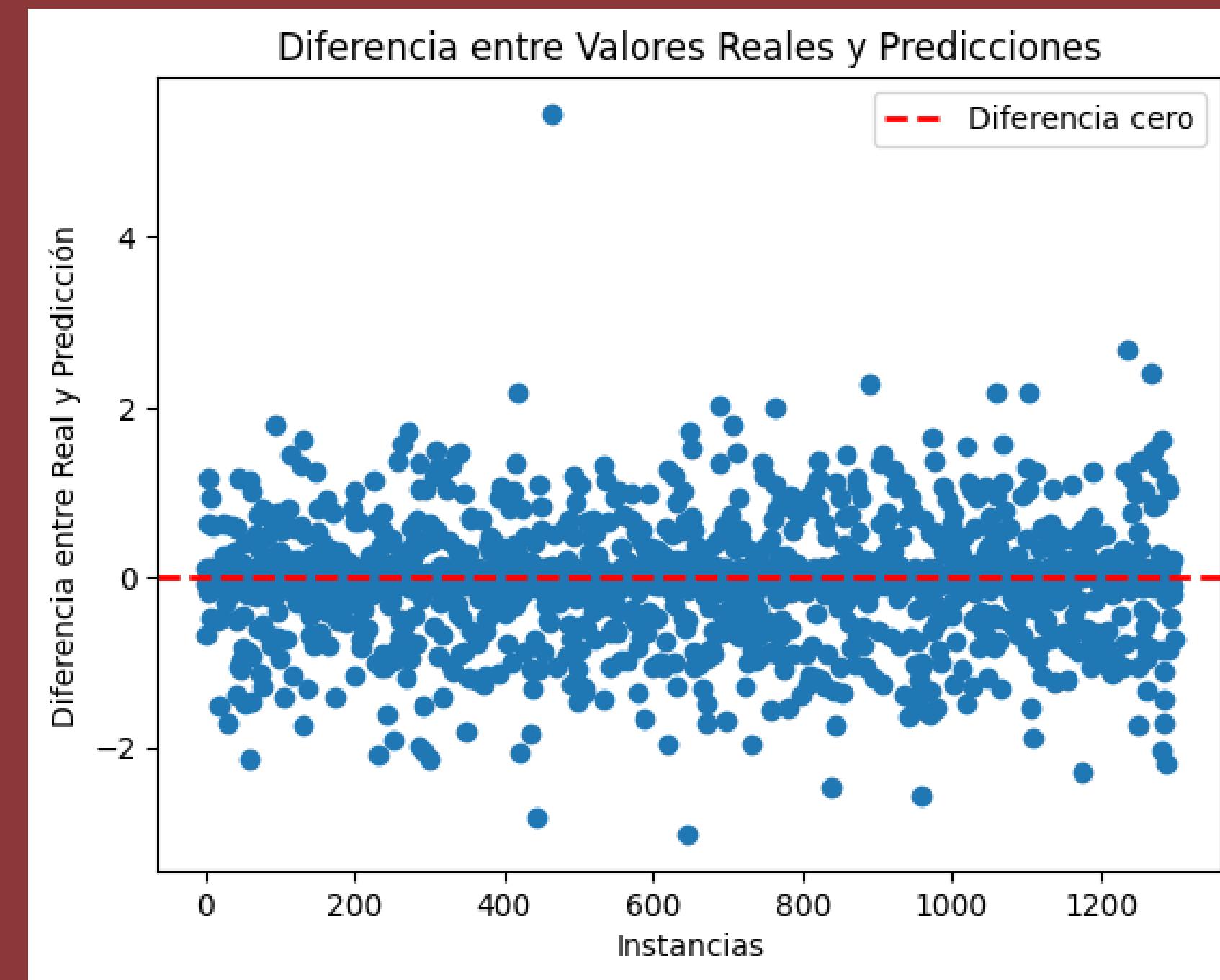
Mean Absolute Error (MAE): 0.7998389355733486

Mean Absolute Percentage Error (MAPE): 14.21285908986947%

Mean Squared Error (MSE): 1.0139882637452864

Root Mean Squared Error (RMSE): 1.0069698425202644

SVC



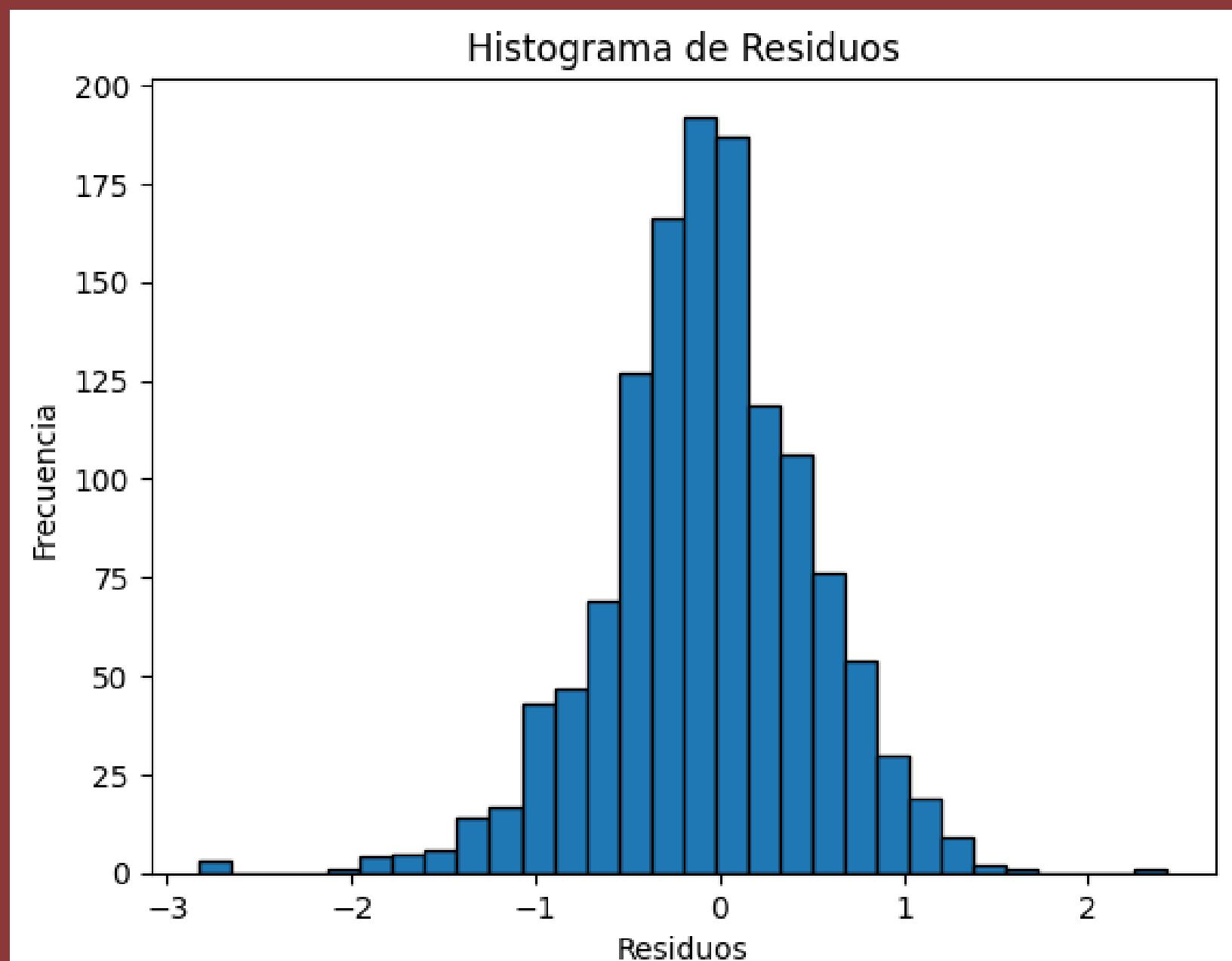
Mean Absolute Error (MAE): 0.4950699727764762

Mean Absolute Percentage Error (MAPE): 8.726739301198808%

Mean Squared Error (MSE): 0.5072335220166678

Root Mean Squared Error (RMSE): 0.7122032870021507

PCA + Random Fores



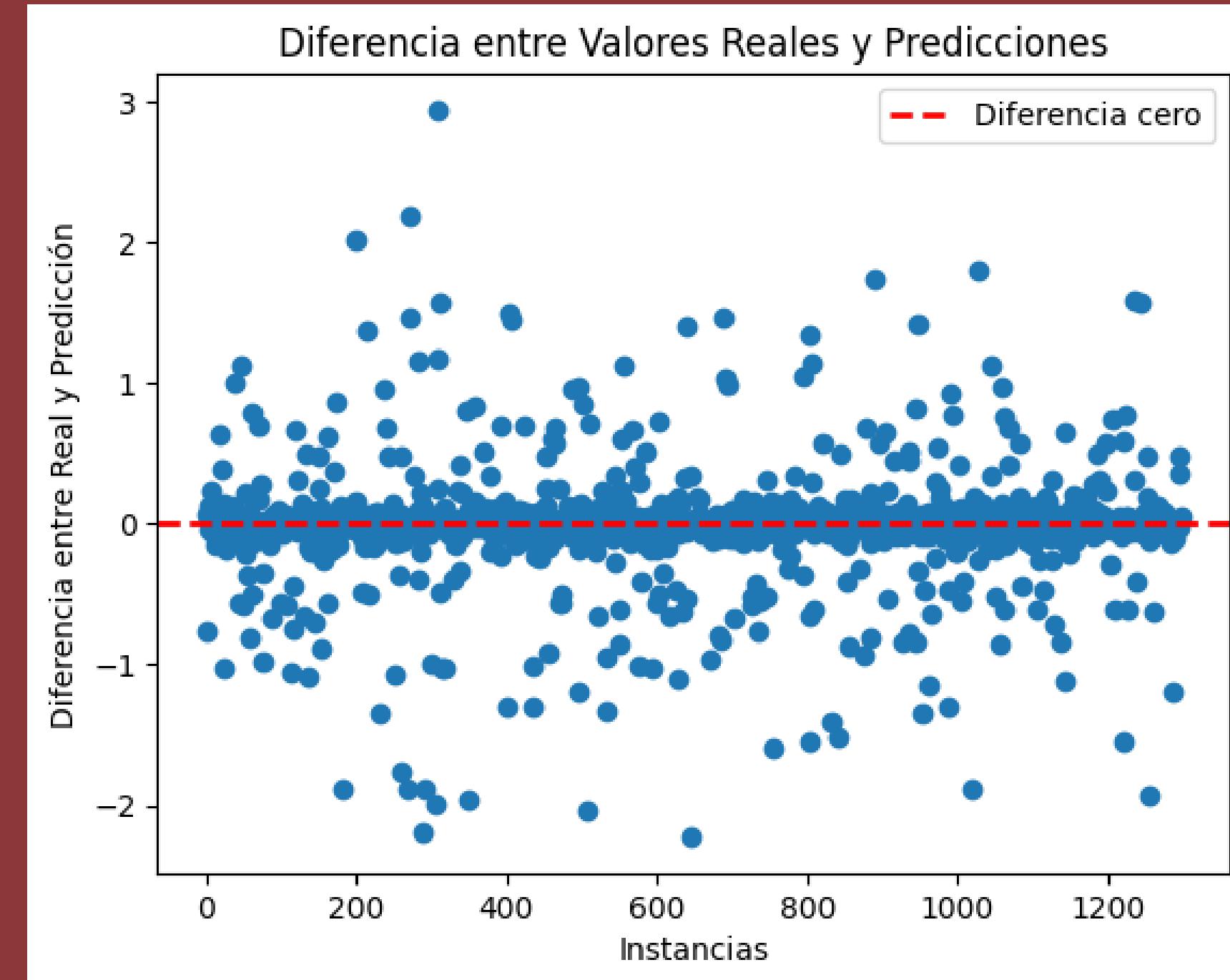
Mean Absolute Error (MAE): 0.4299380569533648

Mean Absolute Percentage Error (MAPE): 7.6386932085148755%

Mean Squared Error (MSE): 0.3241427107536108

Root Mean Squared Error (RMSE): 0.5693353236482089

PCA + Gradient Boosting



Mean Absolute Error (MAE): 0.18381300106680273

Mean Absolute Percentage Error (MAPE): 3.317348100547596%

Mean Squared Error (MSE): 0.15903451492197246

Root Mean Squared Error (RMSE): 0.39879131751076585

Para finalizar, hemos creado un Streamlit, donde podemos comprobar el funcionamiento de nuestro modelo, insertandole los parametros que queramos comprobar y que nos diga su predicción.

Paso a enseñaros a utilizarlo

