

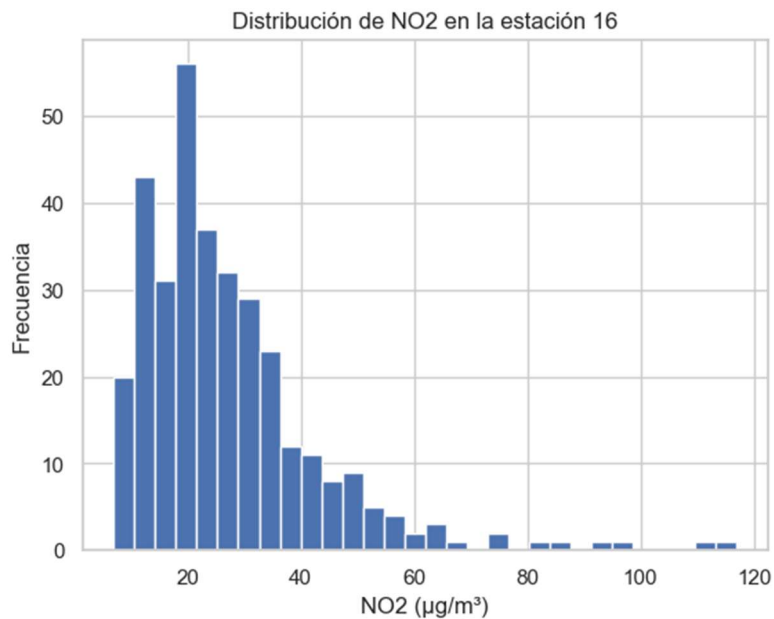
Clasificación mediante árboles de decisión

1. Introducción

Este informe describe el proceso llevado a cabo para aplicar modelos de clasificación basados en árboles de decisión al análisis de datos de calidad del aire en Madrid. Se ha utilizado la biblioteca scikit-learn de Python, y se ha trabajado con la variable NO₂, discretizada en tres niveles: bajo, medio y alto. A lo largo del informe se detallan los pasos seguidos desde la preparación de los datos hasta la comparación de distintos modelos de clasificación.

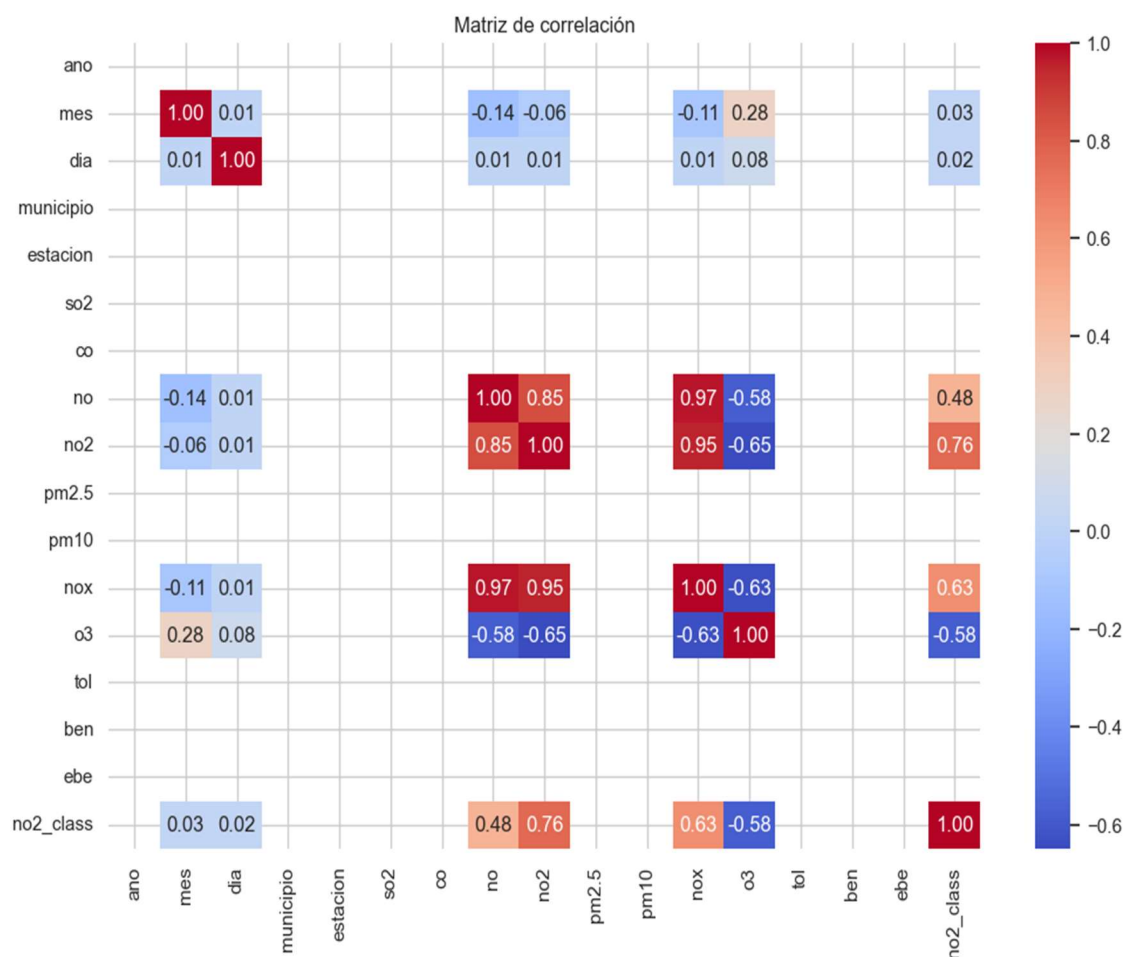
2. Preparación de los datos

Se utilizó un conjunto de datos con medidas diarias de contaminantes. Se seleccionó la estación número 16, que presentaba el mayor número de registros (334). La variable objetivo, NO₂, se discretizó usando la técnica de cuantiles para obtener tres clases. Se descartaron los registros con valores nulos y se mantuvieron las columnas numéricas más relevantes.



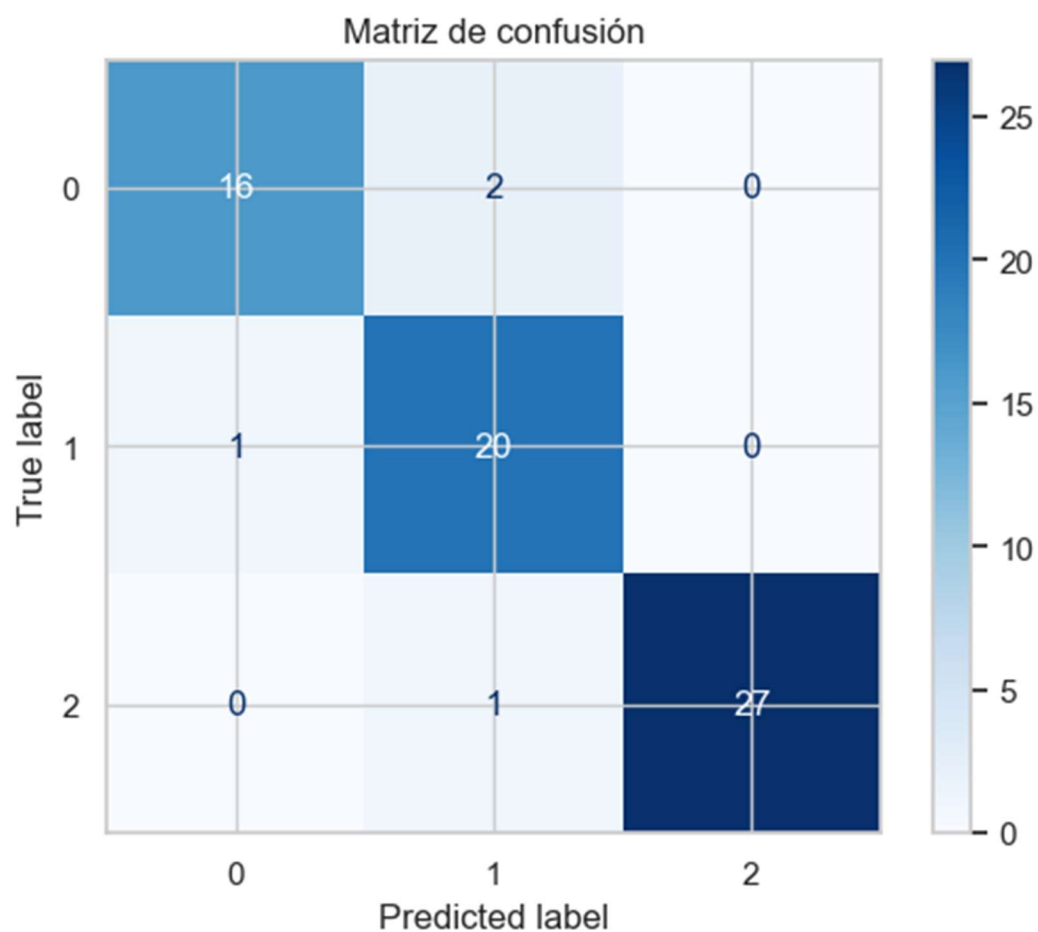
3. Análisis de correlación y selección de atributos

Para seleccionar los atributos predictores, se generó una matriz de correlación. Se observó que las variables más correlacionadas con NO2 eran 'nox' (0.95), 'no', 'pm10' y 'co'. Estas fueron seleccionadas como variables predictoras. La variable 'nox' se consideró clave y posteriormente se eliminó para evaluar su impacto en el rendimiento.

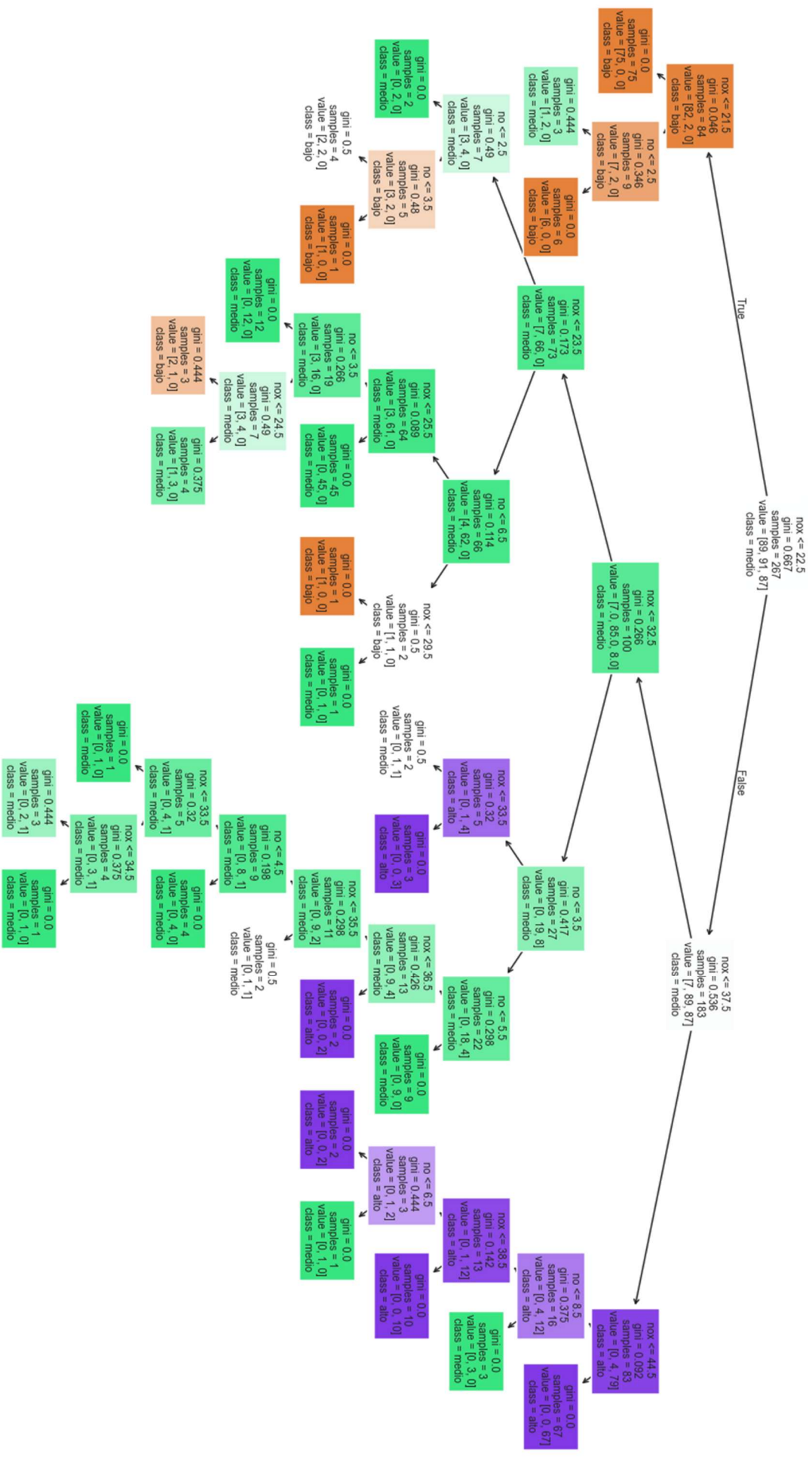


4. Entrenamiento y evaluación del árbol de decisión

Se entrenó un modelo `DecisionTreeClassifier` usando un 80% del dataset para entrenamiento y un 20% para test. El modelo obtuvo un accuracy del 94% y mostró un buen equilibrio en las tres clases. A continuación, se visualiza la matriz de confusión y el árbol de decisión generado.

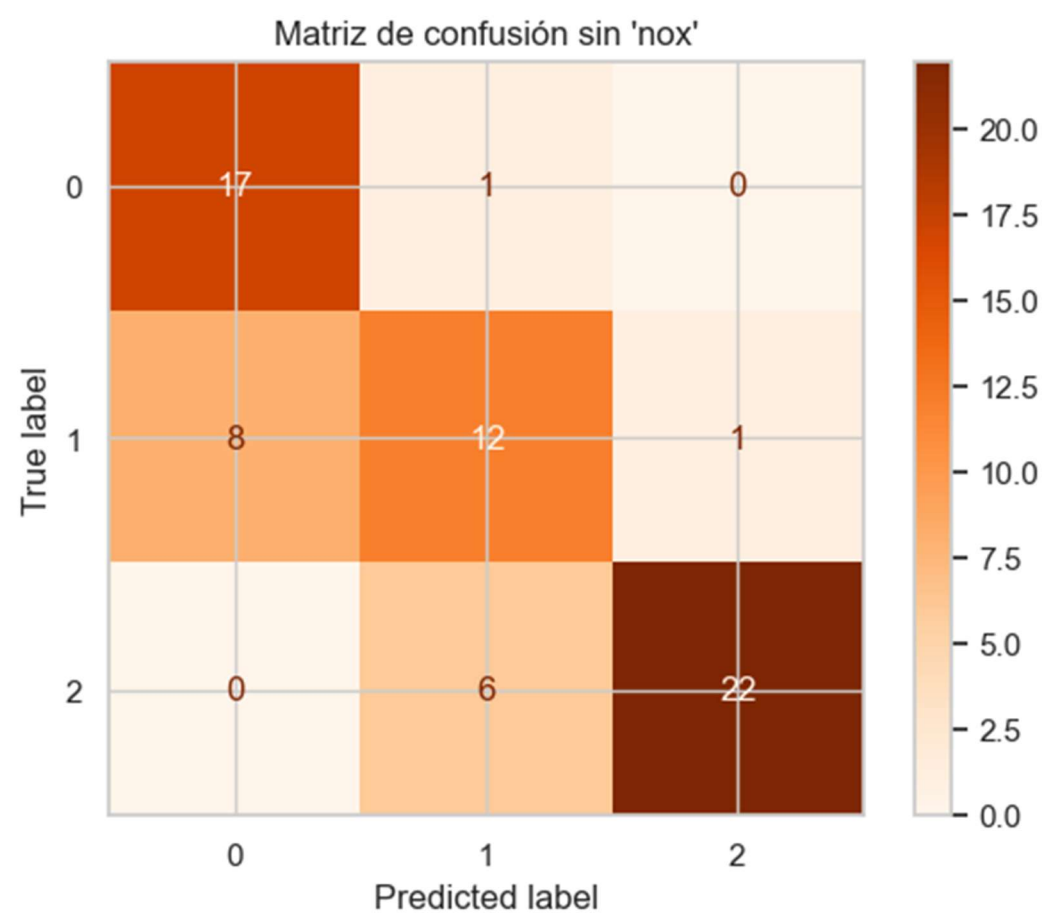


Árbol de decisión



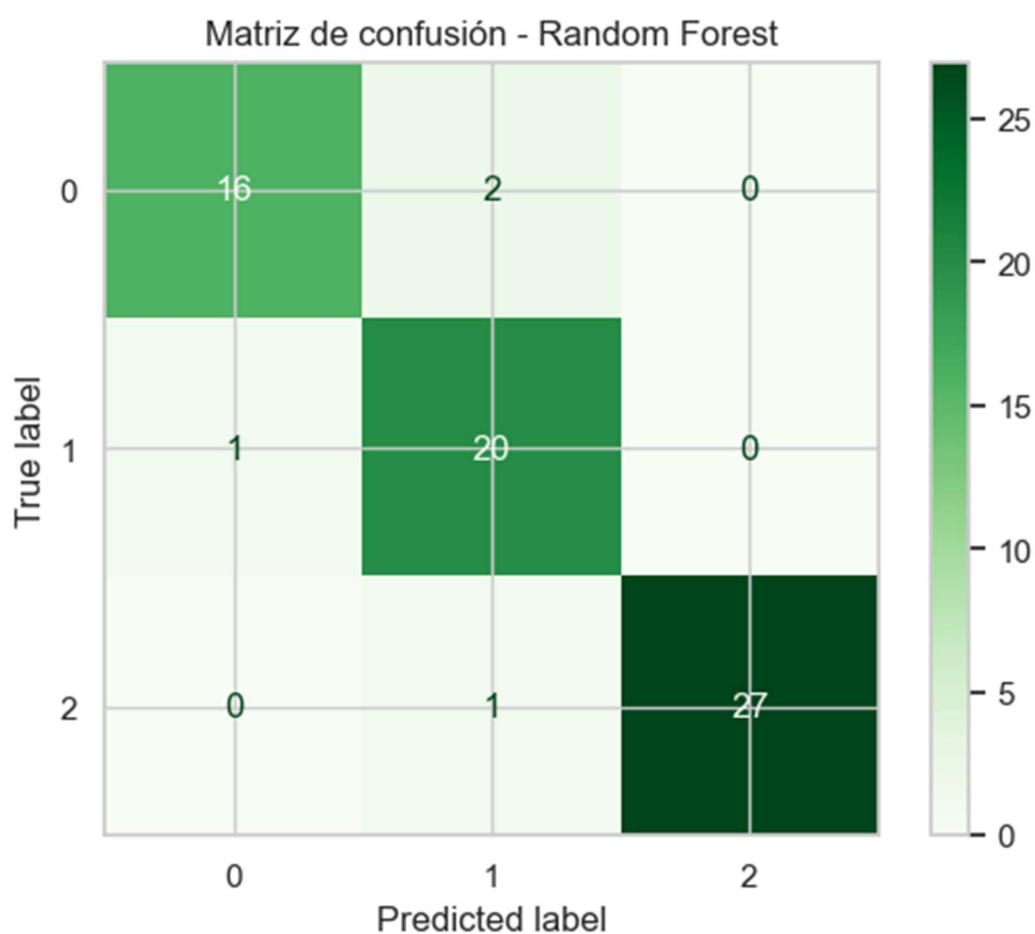
5. Evaluación tras eliminar la variable más correlacionada

Para analizar el impacto de la variable más correlacionada ('nox'), se eliminó del conjunto de atributos y se volvió a entrenar el modelo. El accuracy descendió a 76%, lo cual demuestra la alta relevancia de esta variable para el modelo. Las métricas también se vieron afectadas, especialmente en la clase intermedia.



6. Comparación con Random Forest

Se entrenó también un modelo RandomForestClassifier con los mismos atributos iniciales. El modelo obtuvo también un accuracy del 94%, demostrando ser una alternativa robusta al árbol de decisión simple. Random Forest tiene la ventaja de reducir el sobreajuste gracias a la combinación de múltiples árboles.



7. Conclusiones

A lo largo de este análisis se ha demostrado que los modelos de clasificación basados en árboles de decisión, como `DecisionTreeClassifier` y `RandomForestClassifier`, son herramientas eficaces para predecir la calidad del aire en función de diferentes variables contaminantes. El análisis de correlaciones ha sido clave para seleccionar los atributos más relevantes (`nox`, `no`, `pm10`, `co`) que, en conjunto, permiten realizar predicciones con un alto grado de precisión.

El modelo de árbol de decisión simple logró un **94% de accuracy**, mostrando un buen equilibrio en las tres clases de `NO2` definidas. La posterior eliminación de la variable más correlacionada (`nox`) provocó una caída notable en el rendimiento del modelo (hasta el 76%), lo que confirma su papel fundamental en el proceso de clasificación. Por otro lado, la implementación de un modelo `Random Forest` con los mismos atributos iniciales demostró un rendimiento igualmente alto (94%), reforzando la estabilidad del conjunto de datos y mostrando las ventajas de los modelos ensamblados en términos de robustez y generalización.

En conjunto, esta experiencia ha permitido poner en práctica conceptos fundamentales del aprendizaje supervisado, la preparación de datos, la evaluación de modelos y la visualización de resultados. Como futuras líneas de mejora, sería interesante ampliar el análisis incluyendo datos de varias estaciones, incorporar variables meteorológicas, o aplicar técnicas más avanzadas como `XGBoost` o `LightGBM` para comparar el rendimiento.