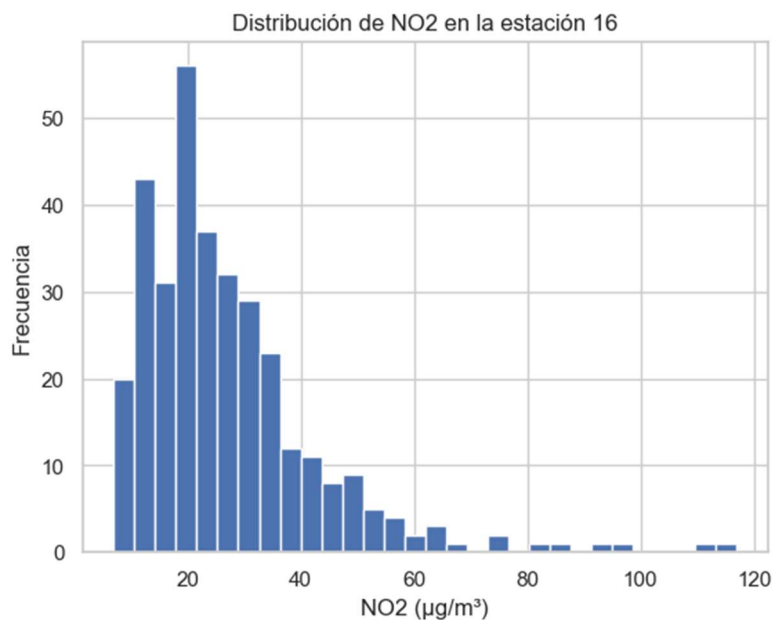# Classification Using Decision Trees

## 1. Introduction

This report describes the process carried out to apply classification models based on decision trees to the analysis of air quality data in Madrid. The Python library **scikit-learn** was used, and the focus was placed on the variable **NO2**, which was discretized into three levels: low, medium, and high. Throughout the report, the steps are detailed from data preparation to the comparison of different classification models.

## 2. Data Preparation
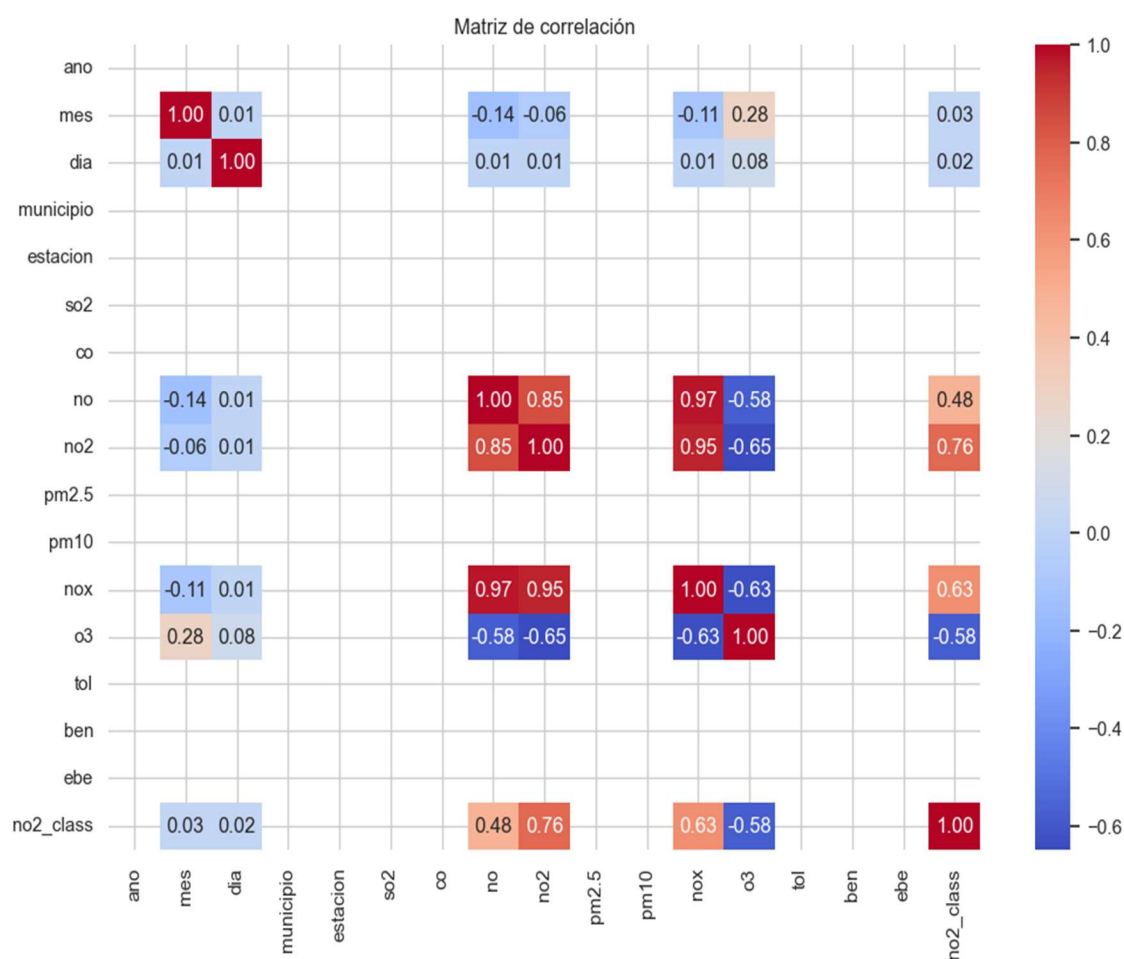
A dataset containing daily measurements of pollutants was used. **Station number 16**, which had the highest number of records (334), was selected. The target variable, **NO2**, was discretized using the **quantile-based technique** to obtain three classes. Records with missing values were removed, and the most relevant numerical columns were retained.



Distribución de NO2 en la estación 16

# 3. Correlation Analysis and Feature Selection

To select the predictor attributes, a **correlation matrix** was generated. It was observed that the variables most correlated with **NO2** were 'nox' (0.95), 'no', 'pm10', and 'co'. These were selected as the predictor features. The variable 'nox' was considered key and was later removed to evaluate its impact on model performance.



Matriz de correlación

# 4. Training and Evaluation of the Decision Tree

A **DecisionTreeClassifier** model was trained using 80% of the dataset for training and 20% for testing. The model achieved an **accuracy of 94%** and showed a good balance across the three classes. The following figures display the resulting **confusion matrix** and the **generated decision tree**.

Árbol de decisión

nox <= 22.5
gini = 0.667
samples = 267
value = [89, 91, 87]
class = medio

True

False

nox <= 21.5
gini = 0.046
samples = 84
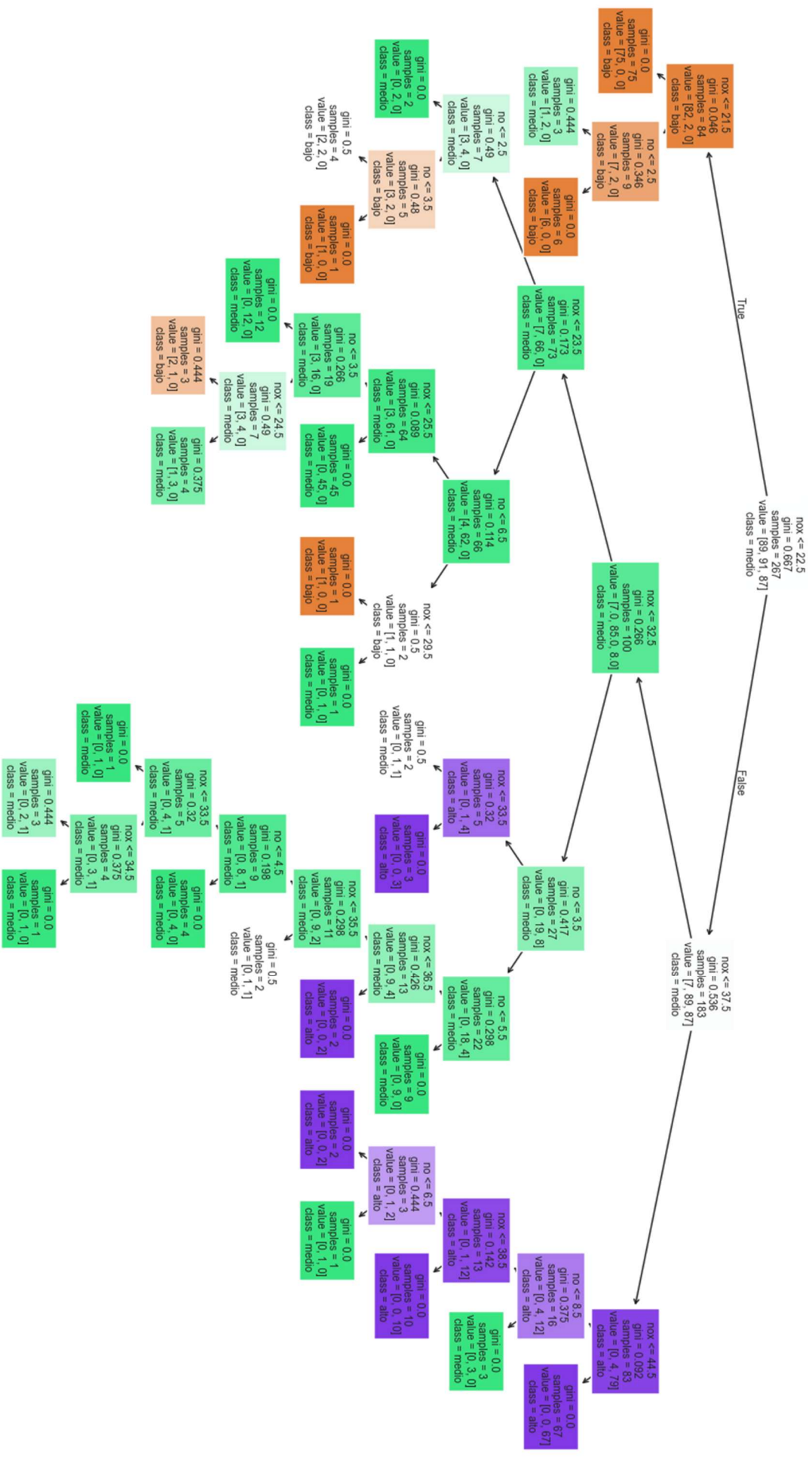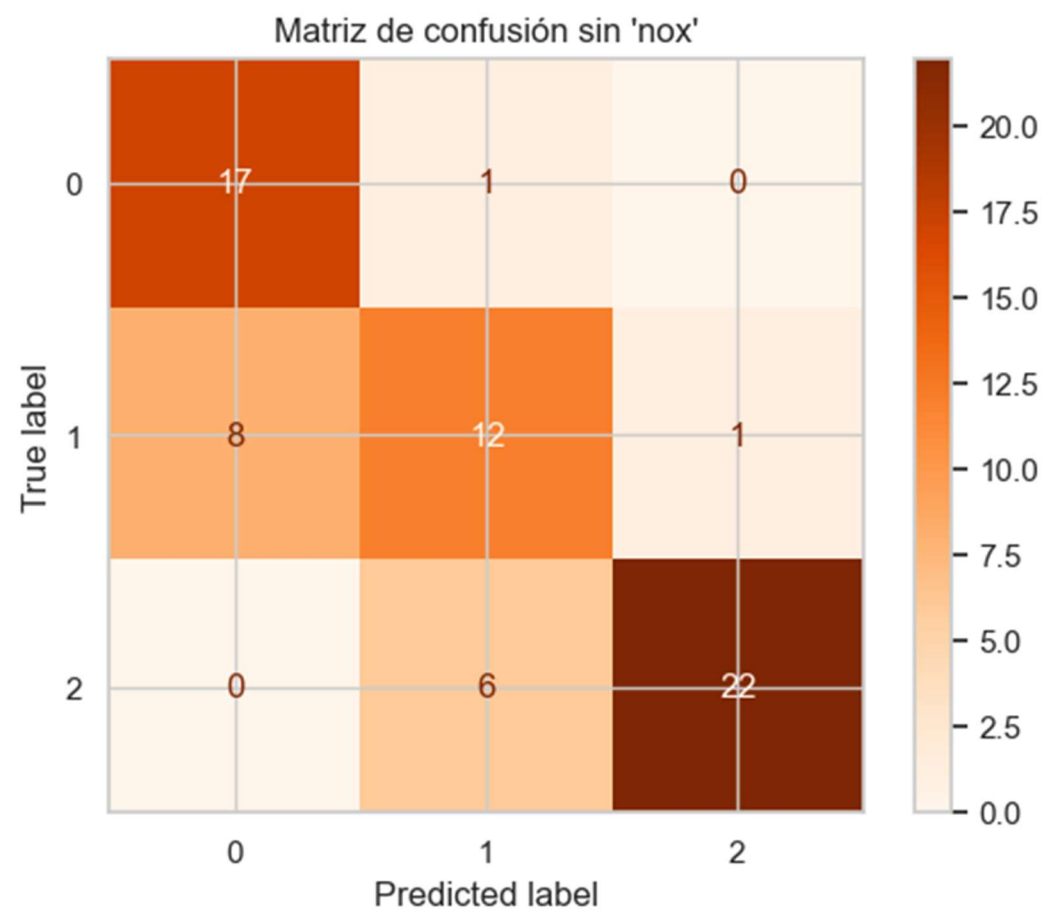value = [82, 2, 0]
class = bajo

gini = 0.0
samples = 75
value = [75, 0, 0]
class = bajo

no <= 2.5
gini = 0.346
samples = 9
value = [7, 2, 0]
class = bajo

gini = 0.0
samples = 6
value = [6, 0, 0]
class = bajo

gini = 0.444
samples = 3
value = [1, 2, 0]
class = medio

no <= 2.5
gini = 0.49
samples = 7
value = [3, 4, 0]
class = medio

gini = 0.0
samples = 2
value = [0, 2, 0]
class = medio

no <= 3.5
gini = 0.48
samples = 5
value = [3, 2, 0]
class = bajo

gini = 0.5
samples = 4
value = [2, 2, 0]
class = bajo

gini = 0.0
samples = 1
value = [1, 0, 0]
class = bajo

nox <= 23.5
gini = 0.173
samples = 73
value = [7, 66, 0]
class = medio

nox <= 25.5
gini = 0.089
samples = 64
value = [3, 61, 0]
class = medio

no <= 3.5
gini = 0.266
samples = 19
value = [3, 16, 0]
class = medio

gini = 0.0
samples = 12
value = [0, 12, 0]
class = medio

no <= 24.5
gini = 0.49
samples = 7
value = [3, 4, 0]
class = medio

gini = 0.444
samples = 3
value = [2, 1, 0]
class = bajo

gini = 0.375
samples = 4
value = [1, 3, 0]
class = medio

gini = 0.0
samples = 45
value = [0, 45, 0]
class = medio

no <= 6.5
gini = 0.114
samples = 66
value = [4, 62, 0]
class = medio

nox <= 29.5
gini = 0.5
samples = 2
value = [1, 1, 0]
class = bajo

gini = 0.0
samples = 1
value = [1, 0, 0]
class = bajo

gini = 0.0
samples = 1
value = [0, 1, 0]
class = medio

nox <= 32.5
gini = 0.266
samples = 100
value = [7, 85, 0, 8, 0]
class = medio

nox <= 33.5
gini = 0.32
samples = 5
value = [0, 1, 4]
class = alto

gini = 0.0
samples = 3
value = [0, 0, 3]
class = alto

no <= 3.5
gini = 0.417
samples = 27
value = [0, 19, 8]
class = medio

nox <= 35.5
gini = 0.298
samples = 11
value = [0, 9, 2]
class = medio

nox <= 36.5
gini = 0.426
samples = 13
value = [0, 9, 4]
class = medio

no <= 5.5
gini = 0.298
samples = 22
value = [0, 18, 4]
class = medio

gini = 0.0
samples = 9
value = [0, 9, 0]
class = medio

nox <= 33.5
gini = 0.32
samples = 9
value = [0, 8, 1]
class = medio

no <= 4.5
gini = 0.198
samples = 9
value = [0, 8, 1]
class = medio

nox <= 34.5
gini = 0.375
samples = 4
value = [0, 3, 1]
class = medio

gini = 0.0
samples = 5
value = [0, 4, 1]
class = medio

gini = 0.0
samples = 1
value = [0, 1, 0]
class = medio

gini = 0.444
samples = 3
value = [0, 2, 1]
class = medio

gini = 0.0
samples = 1
value = [0, 1, 0]
class = medio

gini = 0.5
samples = 2
value = [0, 1, 1]
class = medio

gini = 0.0
samples = 2
value = [0, 0, 2]
class = alto

gini = 0.0
samples = 2
value = [0, 0, 2]
class = alto

no <= 6.5
gini = 0.444
samples = 3
value = [0, 1, 2]
class = alto

gini = 0.0
samples = 1
value = [0, 1, 0]
class = medio

gini = 0.0
samples = 2
value = [0, 0, 2]
class = alto

nox <= 37.5
gini = 0.536
samples = 183
value = [7, 89, 87]
class = medio

no <= 8.5
gini = 0.375
samples = 16
value = [0, 4, 12]
class = alto

nox <= 38.5
gini = 0.142
samples = 13
value = [0, 1, 12]
class = alto

gini = 0.0
samples = 3
value = [0, 3, 0]
class = medio

gini = 0.0
samples = 10
value = [0, 0, 10]
class = alto

nox <= 44.5
gini = 0.092
samples = 83
value = [0, 4, 79]
class = alto

gini = 0.0
samples = 3
value = [0, 3, 0]
class = medio

gini = 0.0
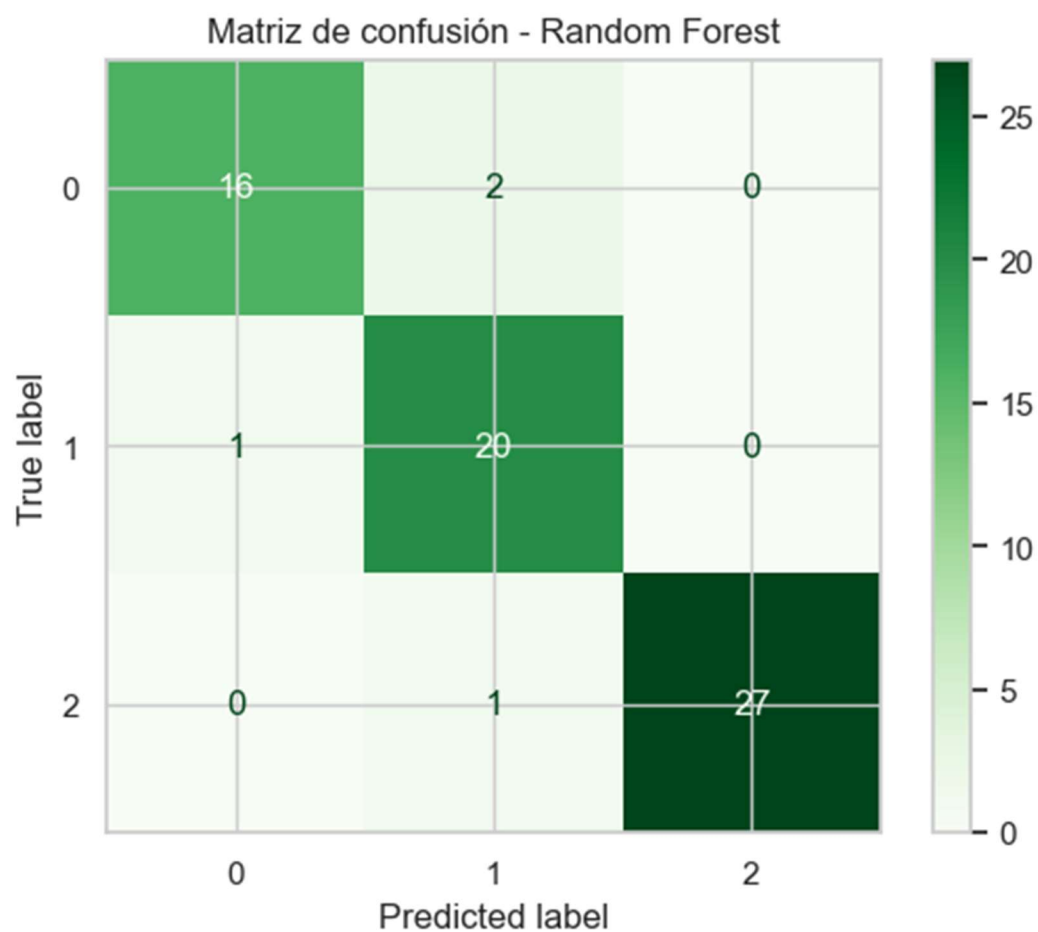samples = 67
value = [0, 0, 67]
class = alto

# 5. Evaluation After Removing the Most Correlated Variable

To analyze the impact of the most correlated variable ('nox'), it was removed from the set of features and the model was retrained. The **accuracy dropped to 76%**, demonstrating the high importance of this variable for the model. The metrics were also affected, especially for the **intermediate class**.



Matriz de confusión sin 'nox'

# 6. Comparison with Random Forest

A **RandomForestClassifier** model was also trained using the same initial set of features. The model achieved an **accuracy of 94%** as well, proving to be a robust alternative to the simple decision tree. **Random Forest** has the advantage of reducing overfitting by combining multiple decision trees.



Matriz de confusión - Random Forest

# 7. Conclusions

Throughout this analysis, it has been demonstrated that classification models based on decision trees, such as **DecisionTreeClassifier** and **RandomForestClassifier**, are effective tools for predicting air quality based on various pollutant variables. The **correlation analysis** was key in selecting the most relevant features (nox, no, pm10, co), which together enabled highly accurate predictions.

The simple decision tree model achieved **94% accuracy**, showing a good balance across the three defined **NO2 classes**. The subsequent removal of the most correlated variable (nox) resulted in a significant drop in performance (to **76%**), confirming its crucial role in the classification process.

On the other hand, the implementation of a **Random Forest** model using the same initial features also achieved high performance (**94% accuracy**), reinforcing the dataset's stability and showcasing the advantages of ensemble models in terms of **robustness and generalization**.

Overall, this experience provided hands-on application of key concepts in **supervised learning**, **data preparation**, **model evaluation**, and **result visualization**. As future improvements, it would be interesting to expand the analysis by including data from multiple stations, incorporating **meteorological variables**, or applying more advanced techniques such as **XGBoost** or **LightGBM** to compare performance.