

## **Preprocessing Steps and Rationale**

The dataset underwent essential preprocessing steps to ensure data quality and model efficiency. Missing values were identified and removed to prevent inconsistencies. Categorical variables, if any, were encoded using one-hot encoding to convert them into numerical format. To standardize the spectral reflectance data, feature scaling was applied using `StandardScaler`, ensuring that all wavelength bands had a consistent range and distribution. This step was crucial for enhancing the performance of machine learning models, particularly those sensitive to feature magnitudes.

## **Insights from Dimensionality Reduction**

To explore the data distribution and identify patterns, dimensionality reduction techniques such as Principal Component Analysis (PCA) and t-SNE were implemented. PCA helped reduce the dataset's complexity while retaining significant variance, with the first two principal components explaining a substantial portion of the total variance. The visualization of PCA projections provided insight into the separability of data points based on spectral reflectance. Additionally, t-SNE was used for nonlinear dimensionality reduction, revealing potential clustering tendencies that might assist in distinguishing different sample characteristics.

## **Model Selection, Training, and Evaluation**

For the predictive task, XGBoost was selected due to its superior handling of structured tabular data and its ability to capture complex relationships. The dataset was split into training (80%) and testing (20%) sets to evaluate generalization performance. Hyperparameter tuning was conducted using `GridSearchCV`, optimizing parameters such as the number of estimators, learning rate, and max depth to enhance model performance. The model was evaluated using regression metrics, including Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and  $R^2$  Score. The results demonstrated that XGBoost effectively predicted DON concentration with relatively low error rates and a high  $R^2$  score, indicating a strong correlation between actual and predicted values.

## **Key Findings and Suggestions for Improvement**

The study showed that spectral reflectance data can effectively predict DON concentration, and machine learning models like XGBoost perform well in this context. The PCA and t-SNE visualizations confirmed that certain spectral bands contribute significantly to the prediction. However, further improvements can be made by experimenting with additional feature selection techniques to retain only the most relevant wavelengths. Incorporating more advanced deep learning architectures, such as convolutional neural networks (CNNs), could further enhance performance by capturing intricate spectral patterns. Additionally, increasing the dataset size or using data augmentation techniques may improve model generalization and robustness. Future work could also explore alternative ensemble methods to refine predictions further.