

Mycotoxin Prediction Report

KATTA BHAVANA

24.03.2025

B.TECH 3RD YEAR

INTRODUCTION

Problem Statement:

The objective of this project is to predict the concentration of vomitoxin (**DON**) in corn samples using hyperspectral imaging data.

Goal:

Develop a robust machine learning model to predict DON levels based on spectral reflectance values.

Approach:

The pipeline involves data preprocessing, exploratory data analysis (**EDA**), model training, hyperparameter tuning, evaluation, and interpretation.

DATA PREPROCESSING

Dataset Overview:

448 spectral reflectance values representing different wavelengths
Target variable: vomitoxin_ppb (DON concentration).

Data Cleaning:

Missing values were filled using median imputation.

Feature Selection:

Only numeric columns were used.
Non-relevant columns like hsi_id were removed.

Normalization:

Features were normalized using StandardScaler for better model performance.

MODEL TRAINING

Model Selection:

Random Forest Regressor was chosen for its interpretability and robustness with high-dimensional data.

Hyperparameter Tuning:

Optuna was used for automated hyperparameter tuning.

Best Parameters:

n_estimators: 177

max_depth: 20

Data Splitting:

80% Training, 20% Testing using `train_test_split`.

```
[I 2025-03-24 11:20:12,598] A new study created in memory with name: no-name-5ac58f43-bfc2-430a-b215-02ddceb6091e
[I 2025-03-24 11:20:19,722] Trial 0 finished with value: 142961285.93544254 and parameters: {'n_estimators': 81, 'max_depth': 21}. Best is trial 0 with value: 142961285.93544254.
[I 2025-03-24 11:20:34,850] Trial 1 finished with value: 126913276.34136614 and parameters: {'n_estimators': 179, 'max_depth': 17}. Best is trial 1 with value: 126913276.34136614.
[I 2025-03-24 11:20:50,148] Trial 2 finished with value: 109645561.07006422 and parameters: {'n_estimators': 177, 'max_depth': 20}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:21:01,697] Trial 3 finished with value: 124173007.13542625 and parameters: {'n_estimators': 284, 'max_depth': 5}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:21:07,414] Trial 4 finished with value: 122934041.23323953 and parameters: {'n_estimators': 138, 'max_depth': 5}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:21:27,380] Trial 5 finished with value: 118507136.22026053 and parameters: {'n_estimators': 226, 'max_depth': 25}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:21:48,358] Trial 6 finished with value: 130033427.46156336 and parameters: {'n_estimators': 215, 'max_depth': 24}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:22:08,725] Trial 7 finished with value: 113169133.32256608 and parameters: {'n_estimators': 233, 'max_depth': 28}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:22:15,604] Trial 8 finished with value: 139171088.66479492 and parameters: {'n_estimators': 81, 'max_depth': 16}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:22:26,274] Trial 9 finished with value: 126157835.71047144 and parameters: {'n_estimators': 265, 'max_depth': 5}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:22:36,817] Trial 10 finished with value: 127831888.88651796 and parameters: {'n_estimators': 138, 'max_depth': 12}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:22:56,013] Trial 11 finished with value: 120274296.44174959 and parameters: {'n_estimators': 220, 'max_depth': 28}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:23:11,140] Trial 12 finished with value: 115162954.24986956 and parameters: {'n_estimators': 170, 'max_depth': 30}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:23:33,299] Trial 13 finished with value: 115450929.20330128 and parameters: {'n_estimators': 256, 'max_depth': 22}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:23:47,553] Trial 14 finished with value: 122714817.78037949 and parameters: {'n_estimators': 182, 'max_depth': 13}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:23:58,111] Trial 15 finished with value: 122949723.60085997 and parameters: {'n_estimators': 123, 'max_depth': 27}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:24:18,839] Trial 16 finished with value: 120936301.53751706 and parameters: {'n_estimators': 241, 'max_depth': 20}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:24:44,843] Trial 17 finished with value: 128365022.43300276 and parameters: {'n_estimators': 299, 'max_depth': 30}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:25:00,901] Trial 18 finished with value: 113461355.02718346 and parameters: {'n_estimators': 198, 'max_depth': 13}. Best is trial 2 with value: 109645561.07006422.
[I 2025-03-24 11:25:11,408] Trial 19 finished with value: 120669851.28666121 and parameters: {'n_estimators': 163, 'max_depth': 9}. Best is trial 2 with value: 109645561.07006422.
Best parameters: {'n_estimators': 177, 'max_depth': 20}
RandomForestRegressor
RandomForestRegressor(max_depth=20, n_estimators=177)
```

MODEL EVALUATION

Metrics Used:

Mean Absolute Error (MAE): Measures average prediction error.

Root Mean Squared Error (RMSE): Evaluates overall prediction accuracy.

R² Score: Measures how well the model explains variability in the data.

Results:

After training with the best hyperparameters:

MAE: 3847.11 ppb

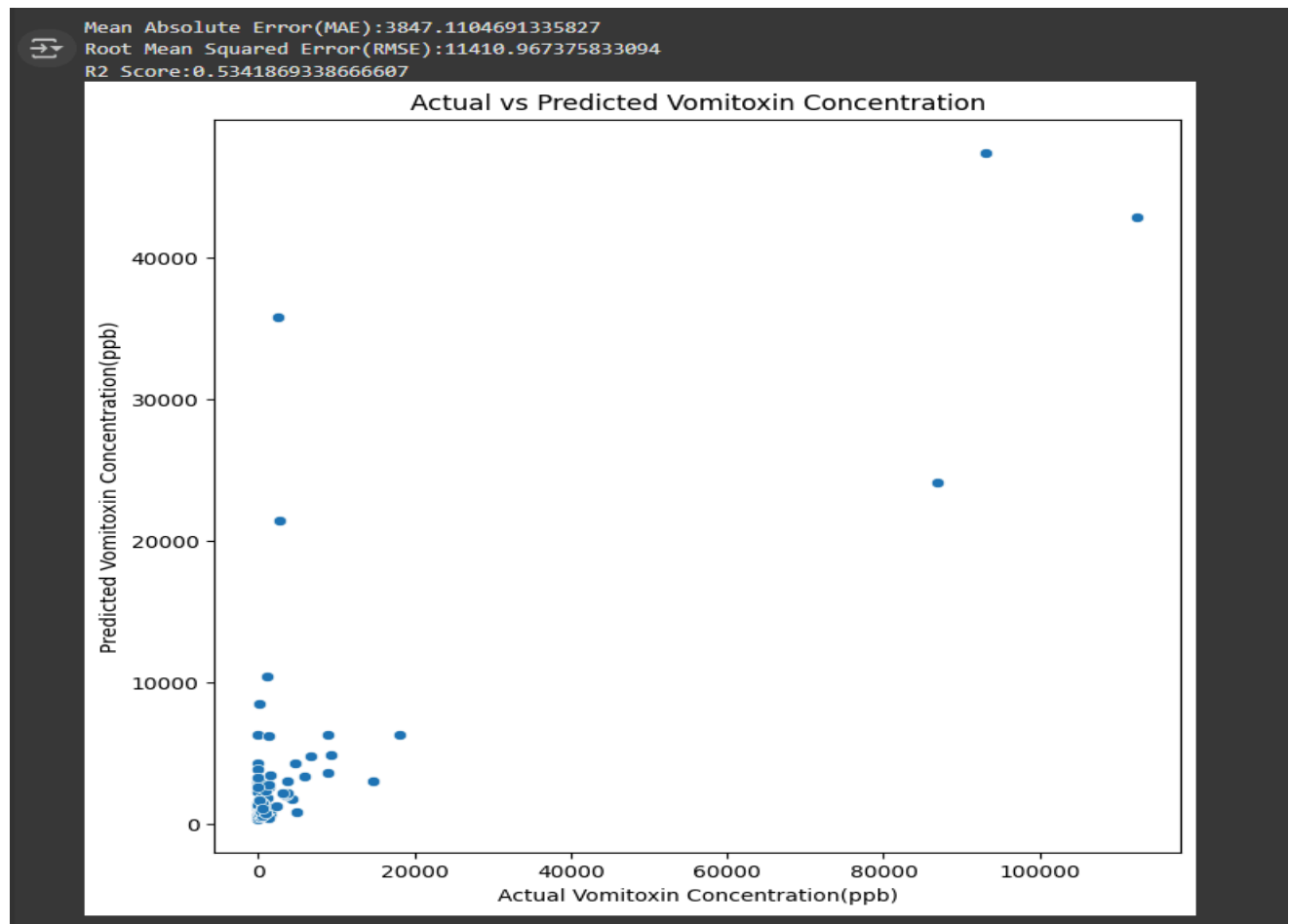
RMSE: 11,410.97 ppb

R² Score: 0.534

These results indicate the model has moderate predictive capability. While the model captures some patterns in the data, further improvements are possible.

Visualization:

A scatter plot comparing actual vs. predicted values showed a moderate correlation, indicating the model's ability to capture general trends.



MODEL INTERPRETABILITY

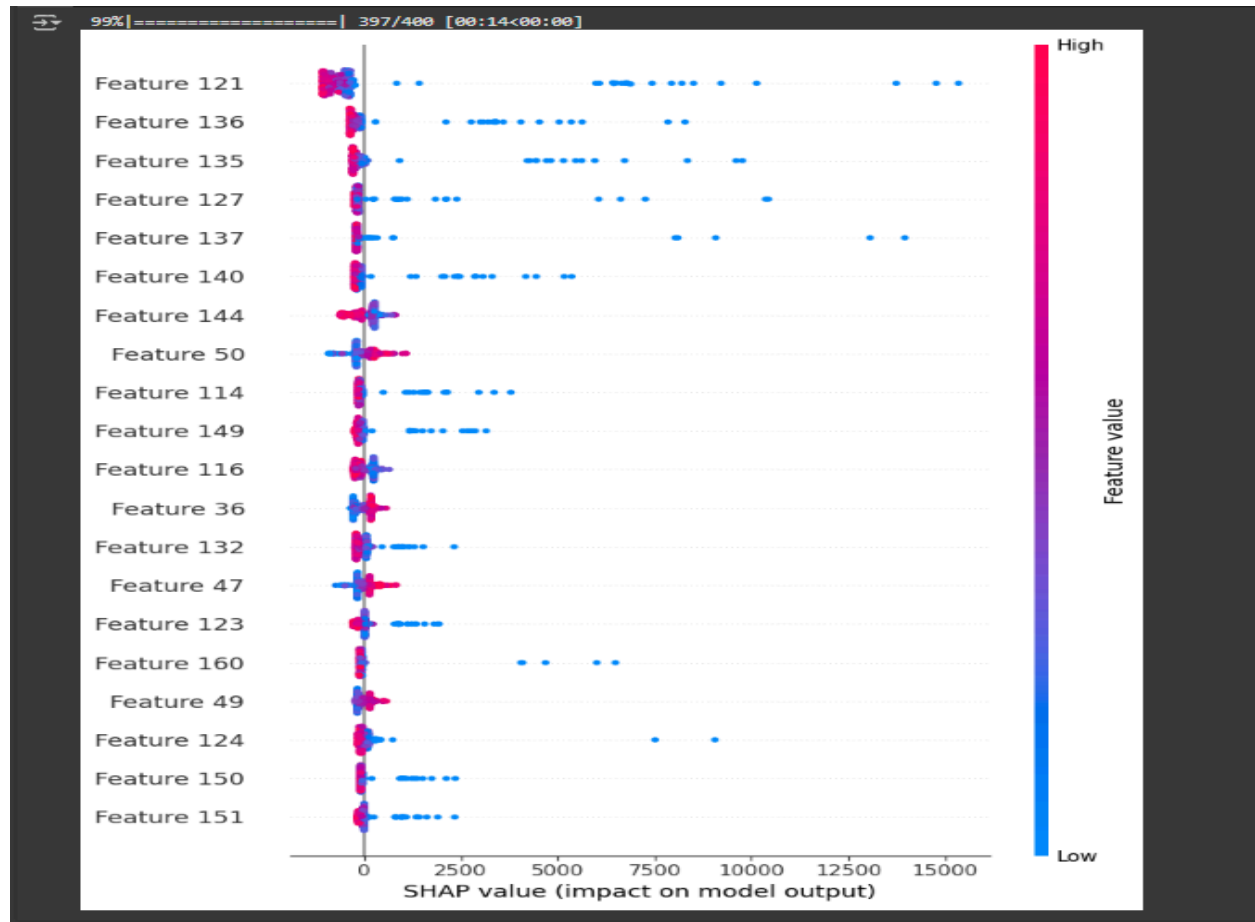
SHAP (SHapley Additive exPlanations) was used to interpret feature importance.

Key Insights from SHAP Plots:

Certain spectral bands contributed significantly to the prediction.

The model relied heavily on specific wavelengths that may indicate early signs of contamination.

Some bands had lower importance, suggesting potential redundancy in the data. Further analysis of these wavelengths could enhance domain knowledge and aid in refining the model.



CONCLUSION

The Random Forest Regressor achieved a reasonable performance with an R^2 score of **0.534**.

The model shows promise in predicting vomitoxin concentration but has room for improvement.

Future improvements can include:

- Trying ensemble models like XGBoost or LightGBM.
- Performing further feature engineering to create meaningful spectral indices.
- Collecting more data or reducing noise in the reflectance measurements.

REFERENCES

1. **Scikit-Learn:** For model training and evaluation.
2. **Optuna:** For hyperparameter tuning.
3. **SHAP:** For interpretability.
4. **Matplotlib & Seaborn:** For data visualization.