

Supervised ML Classifiers for Predicting Water Wells Condition in Tanzania

Leveraging Data to Improve Water Access

- **Student:** Daniel Mwaka
- **Student Pace:** DSF-FT12
- **Phase:** 3
- **Instructor:** Samuel Karu

Problem Statement

- Access to clean water is a critical challenge in Tanzania.
- Many established water points are in disrepair or non-functional, impacting millions leading to water scarcity.
- The lack of an effective, data-driven framework for predicting a well's functional status is a leverageable improvement area.
- **The Goal:** To recommend an evidence-based supervised ML classification model for predicting the functional condition of water wells in Tanzania.



Business Understanding

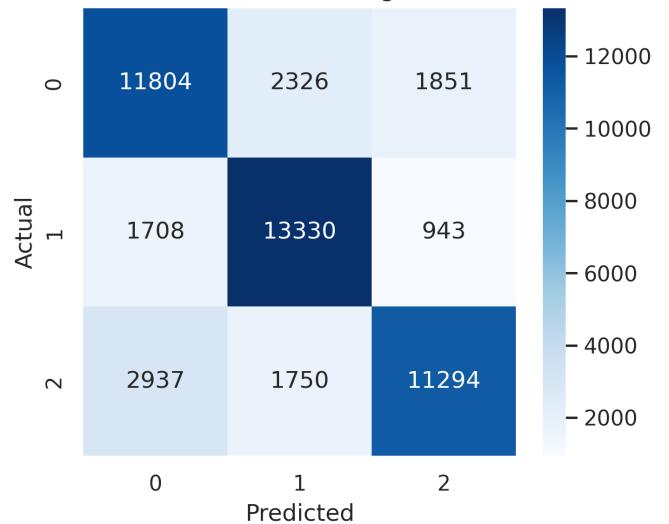
- **Business Problem:** Predicting the condition of water wells can enable stakeholders to prioritize maintenance, and allocate resources efficiently to support water security.
- **Core Question:** Can the operational status of water-wells be accurately predicted using available data and machine learning models?
- **Project Scope:**
 - Examine features related to a Water-well's function to identify potential predictor variables.
 - Build multiple classifiers, tune them, evaluate, and compare their respective prediction performance to determine the best fit model.
 - Verify the reliability of the best-fit model by using it to predict the target variable for a previously unseen dataset (testdata.csv).
 - Deduce feasible business recommendations and viable next steps.

Data Understanding

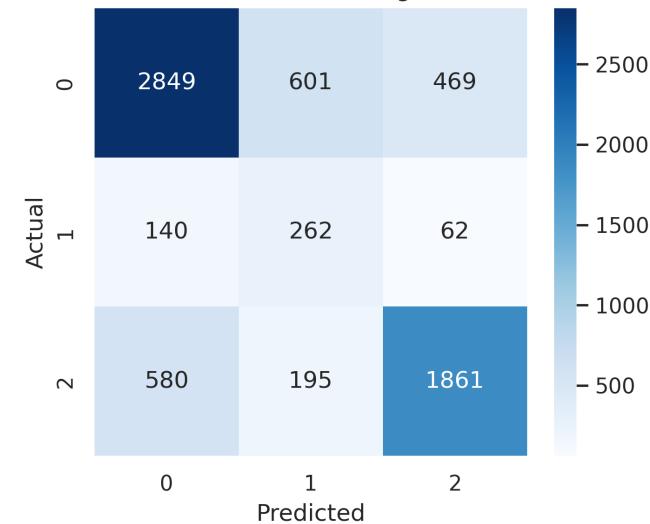
- **Datasets Source:** Kaggle competition dataset on Tanzanian Water Wells from
<https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/data/>
- **Datasets:**
 - trainingset.csv (40 columns and 59,400 entries).
 - trainingsetlabels.csv (2 columns and 59,400 entries).
 - Testdata.csv (40 columns and 14,850 entries).
- **Target Variable:** Water-well Status (Functional, Functional needs repair, or Non-functional).
- **Selected Predictor Features:** gps height, population, well age, basin, region, permit, extraction type, management group, payment type, water quality, water quantity, source type, and water-point type.

Gradient Boosting Classifier

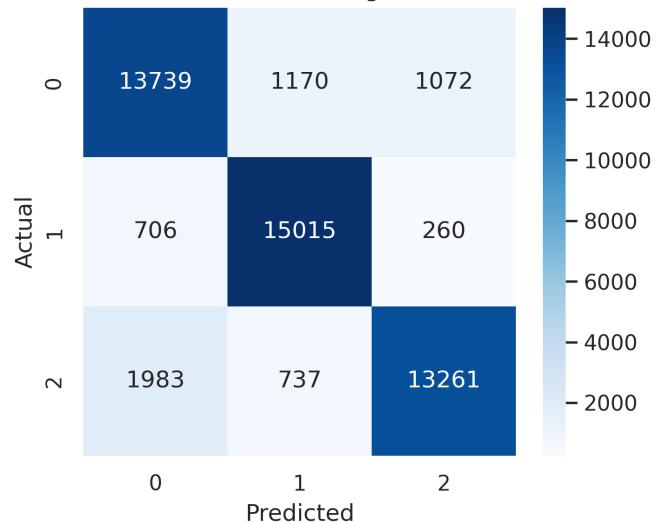
Untuned Gradient Boosting (Train Set)



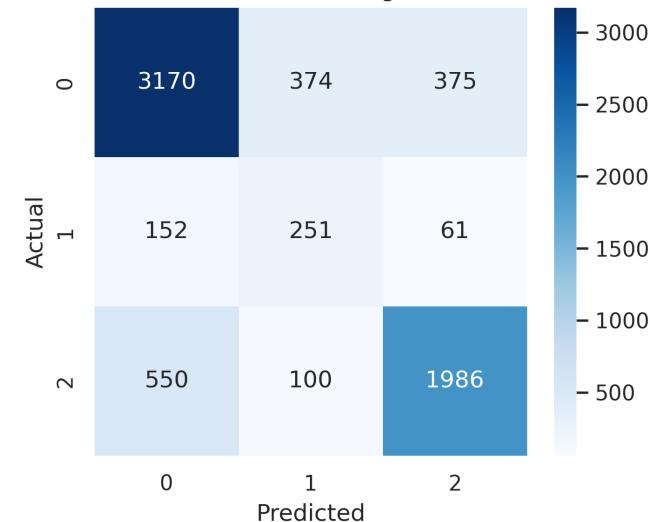
Untuned Gradient Boosting (Test Set)



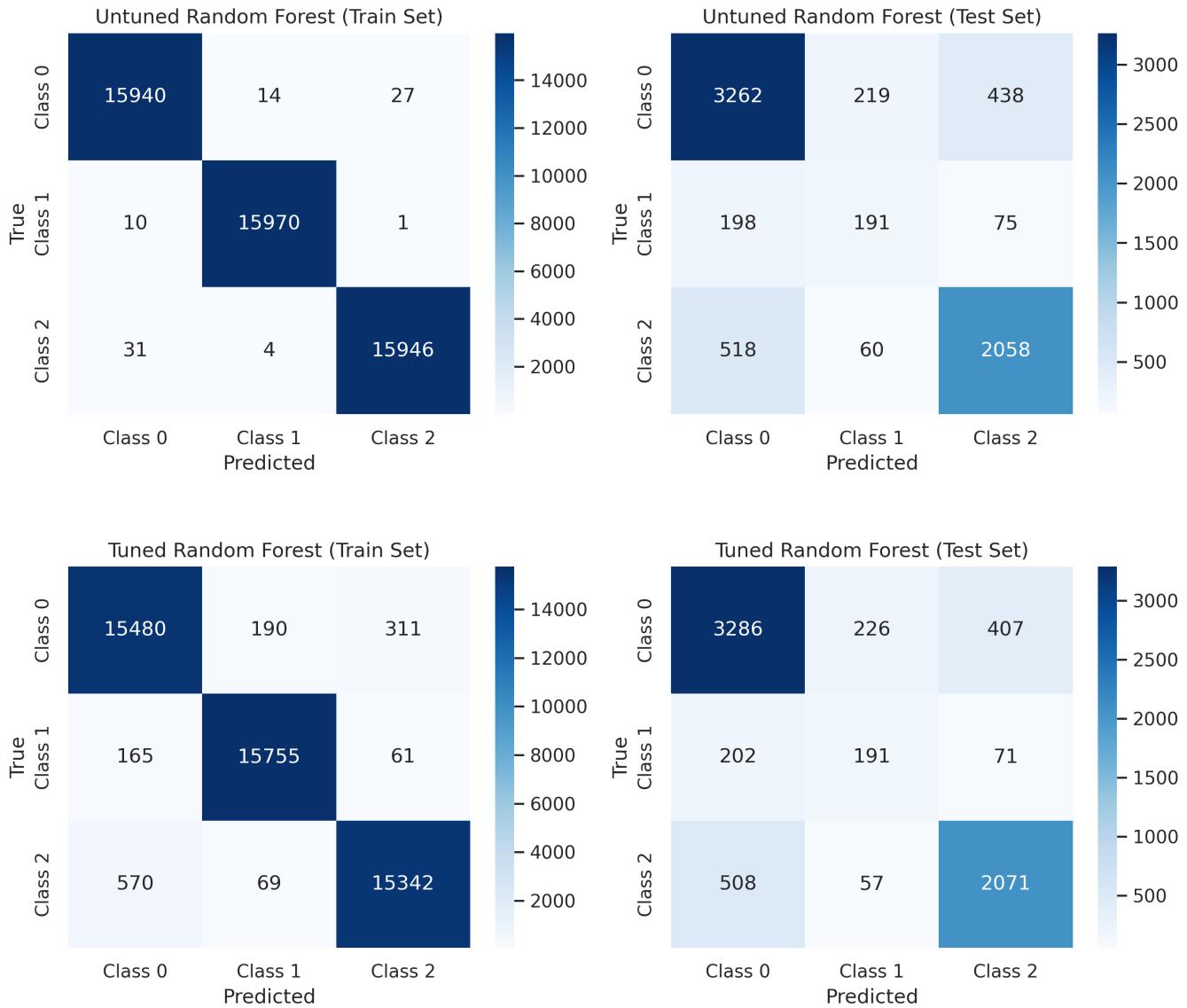
Tuned Gradient Boosting (Train Set)



Tuned Gradient Boosting (Test Set)



Random Forest Classifier

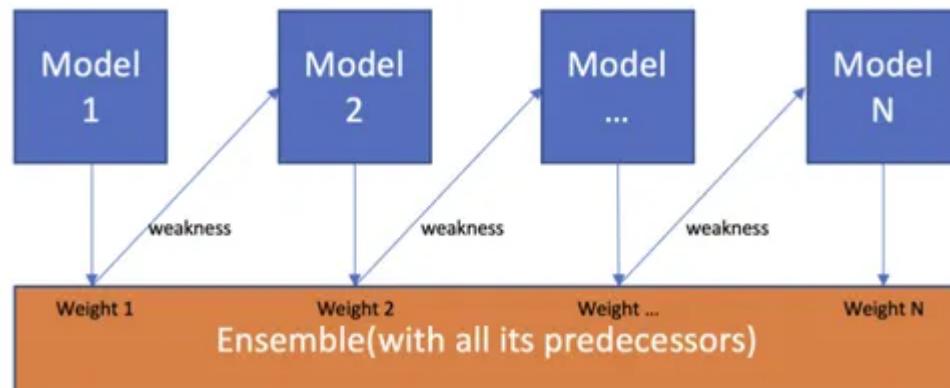


Model Performance Comparison

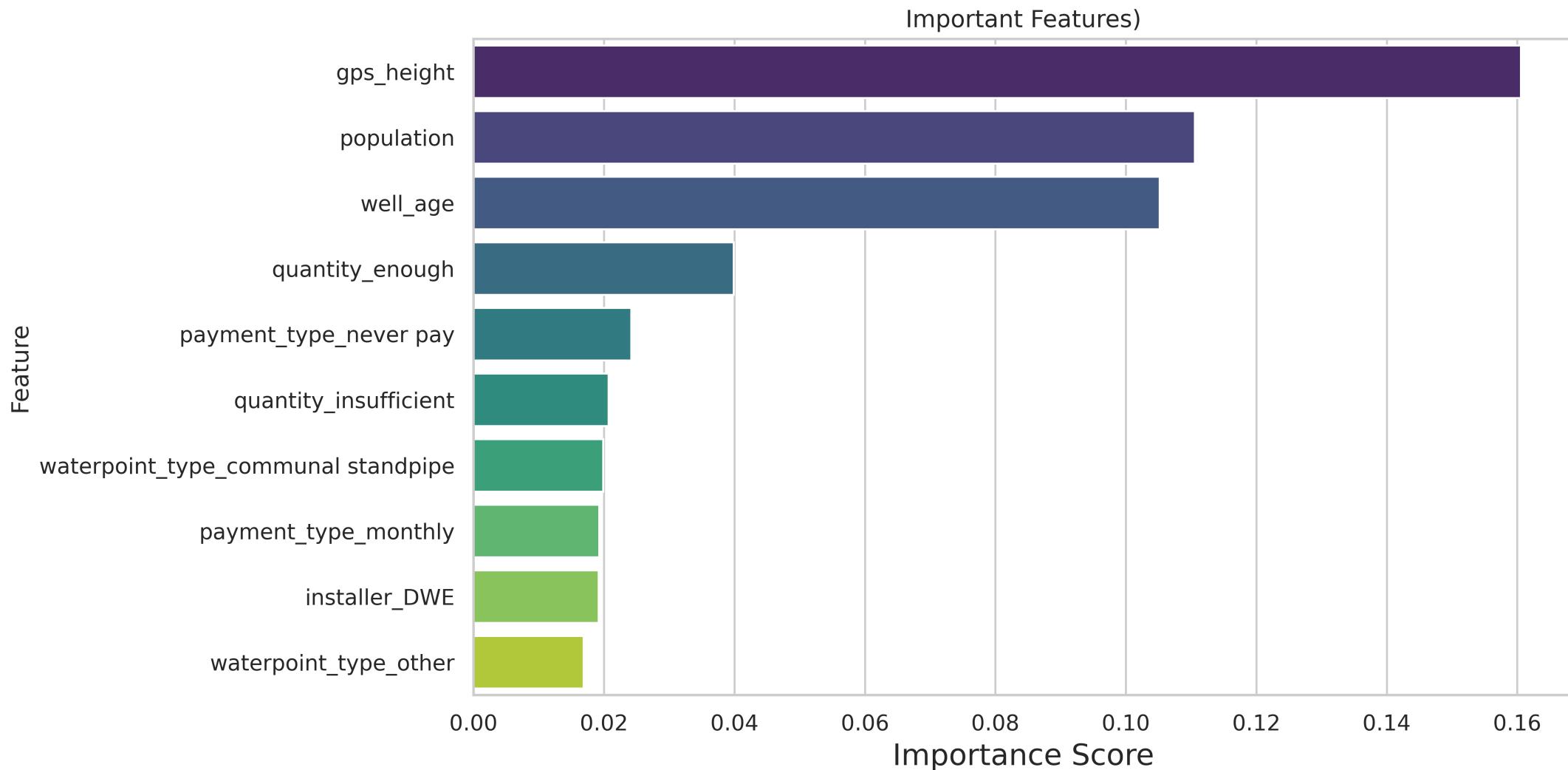
	Model	Train Accuracy	Test Accuracy	Train F1-score	Test F1-score	Train ROC-AUC	Test ROC-AUC
0	Tuned Decision Tree	0.998269	0.752244	0.998269	0.756041	0.999997	0.755808
1	Tuned Gradient Boosting	0.876353	0.770338	0.876001	0.777184	0.971252	0.889499
2	Tuned Random Forest	0.971508	0.790426	0.971518	0.790532	0.998519	0.897939

Selected Model for Deployment

- The ***tuned Gradient Boosting Classifier*** is selected for deployment.
- Reasons for Selection:
 - Highest prediction accuracy in testdata.csv dataset (71.06%).
 - Strikes an effective balance between precision and recall (highest f1-score).
 - Generalizes well to unseen data (minimal gap between train-set and test-set).



Feature Importance Insights



Business Recommendations

- **Prioritize Maintenance:** Use model predictions to identify high-risk wells and allocate maintenance resources efficiently.
- **Improve Data Collection:** Enhance data quality, particularly for key features influencing well condition, such as management, payment types, and environmental factors.
- **Stakeholder Engagement:** Share insights with local authorities and NGOs to inform decision-making, schedule maintenance routines, and find patterns on factors impacting long-term functionality.
- **Inform New Infrastructure:** Use insights on factors contributing to failure to inform the design and construction frameworks for new groundwater projects.

Next Steps

- **Model Deployment:** Integrate the Tuned Gradient Boosting Classifier into a user-friendly dashboard for real-time predictions.
- **Integration into Planning:** Use model predictions to optimize maintenance schedules and target interventions.
- **Pilot Targeted Interventions:** Use the model to pilot interventions in high-risk areas/well types and measure the impact.
- **Feature Expansion:** Incorporate additional data sources (e.g., weather, usage patterns) to enhance model accuracy and update the model regularly with new data to maintain performance and relevance.

