

# **Tanzania Water Wells Classification**



### **Problem Overview**

The Tanzania Ministry of Water along with Taarifa, a crowd-source platform, have commissioned the development of a predictive model that is supposed to be able to predict with **water wells** are likely to fail. While much of Tanzanias population has access to basic water services, a large 39% of households still lack this basic need. An estimated 10% of preventable deaths in the country can be attributed to inadequate *wash services*. A predictive model can enable quick **predictive maintenance** on water wells and help ensure water security in many of the rural communities that are disporportionately affected by this problem.



# **Project Objectives**

The main objective of this project undertaking is to build a Classification model that can be able to classify water wells in Tanzania as functional or those that need repairs need\_repair

#### **Specific Objectives**

- 1. To conduct exploratory analysis and determine which features to include in our model
- 2. Determine the cleaning steps to be included in building the model pipeline
- 3. To build a classifiction model that can predict the status of wells with acceptable accuracy.

#### **Success Metrics**

Accuracy: 75%Recall: 80%

### The Data

The data is provided by an organization known as Taarifa in co-operation with the Tanzanian government. A detailed description can be found **here** 

#### Column Summary

- amount\_tsh Total static head (amount water available to waterpoint)
- date\_recorded The date the row was entered
- funder Who funded the well
- gps\_height Altitude of the well
- · installer Organization that installed the well
- · longitude GPS coordinate
- · latitude GPS coordinate
- wpt\_name Name of the waterpoint if there is one
- num\_private -
- · basin Geographic water basin
- · subvillage Geographic location
- region Geographic location
- region\_code Geographic location (coded)

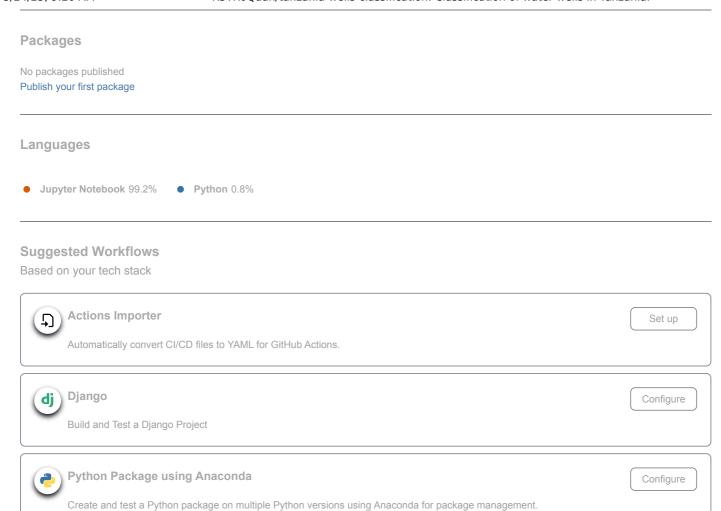
- district code Geographic location (coded)
- Iga Geographic location
- ward Geographic location
- population Population around the well
- public\_meeting True/False
- · recorded\_by Group entering this row of data
- scheme\_management Who operates the waterpoint
- · scheme\_name Who operates the waterpoint
- permit If the waterpoint is permitted
- · construction year Year the waterpoint was constructed
- extraction\_type The kind of extraction the waterpoint uses
- extraction\_type\_group The kind of extraction the waterpoint uses
- extraction\_type\_class The kind of extraction the waterpoint uses
- · management How the waterpoint is managed
- management group How the waterpoint is managed
- payment What the water costs
- payment\_type What the water costs
- water quality The quality of the water
- · quality\_group The quality of the water
- · quantity The quantity of water
- · quantity\_group The quantity of water
- · source The source of the water
- source type The source of the water
- source\_class The source of the water
- waterpoint type The kind of waterpoint
- · waterpoint\_type\_group The kind of waterpoint

## Conclusions

- The Baseline **Logistic Regression** performs well on our classification metrics i.e **accuracy** and **recall**Accuracy is a valid metric as the class imbalance is negligible as seen in the data exploration
- The **KNN model** performs similarly to the Logistic Regression. However, due to its exponentially increasing time complexity, this model has a much longer runtime The *runtime* does not justify using this model.
- Both Decision Trees and their ensemble counterpart, Random Forest performed worse than the baseline model.
- After tuning the Logistic Regression model, as it preforms best, we obtain the best parameters.
- Using a feature selector proved detrimental to the modelling process.
  - Best Model: logistic regression(tuned)

#### Releases

No releases published Create a new release



More workflows

Dismiss suggestions