

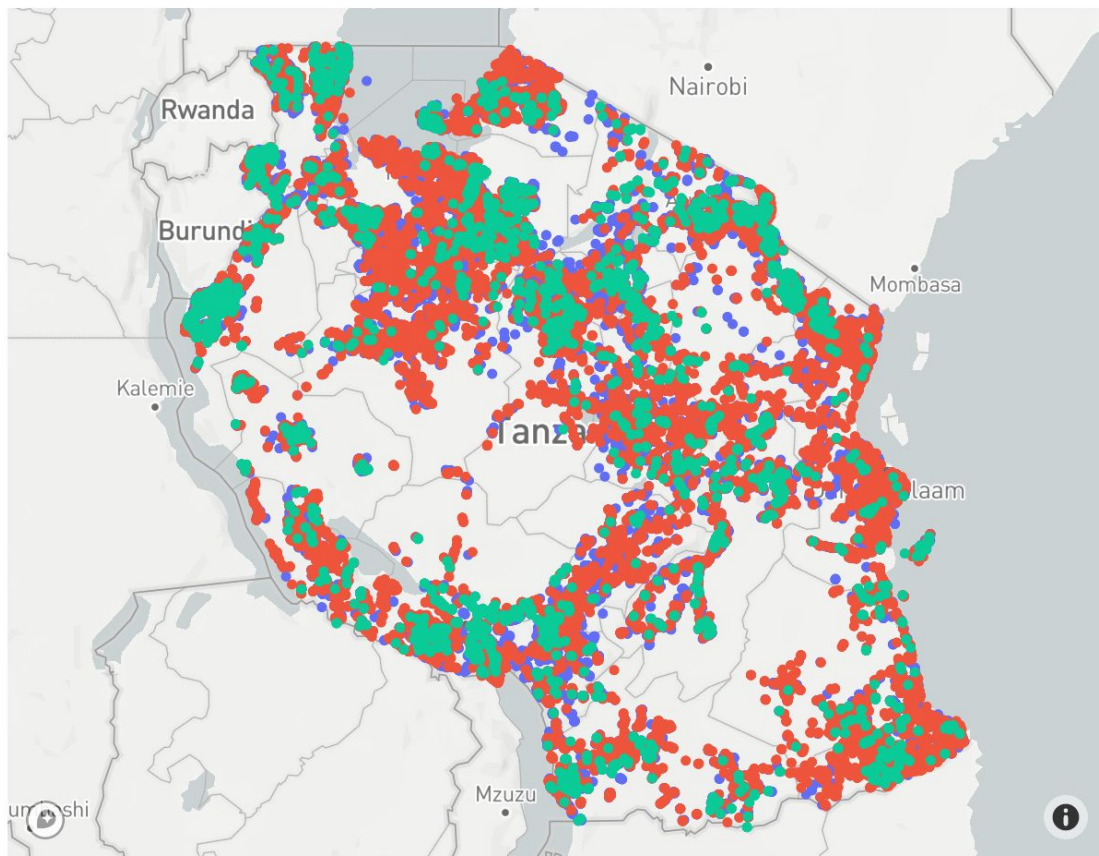
# Tanzanian Water Well Data

David Cuervo

# Background

- Competition through DrivenData
- Data collected from Taarifa and the Tanzanian Ministry of Water
- Predicting well functionality based on a number of variables about the well's age, location, and how it is managed





- status\_group=functional
- status\_group=non functional
- status\_group=functional needs repair

# Data Exploration and Cleaning

## Columns Dropped

- Went through each variable to deal with missing data and outliers
  - Dropped columns
  - Made dummy variables to replace the categorical variables in the dataset
- date\_recorded
  - wpt\_name
  - num\_private
  - region\_code
  - district\_code
  - public\_meeting
  - recorded\_by
  - scheme\_name
  - extraction\_type
  - extraction\_type\_class
  - Management
  - management\_group
  - Payment
  - water\_quality
  - quantity\_group
  - Source
  - source\_class
  - waterpoint\_type

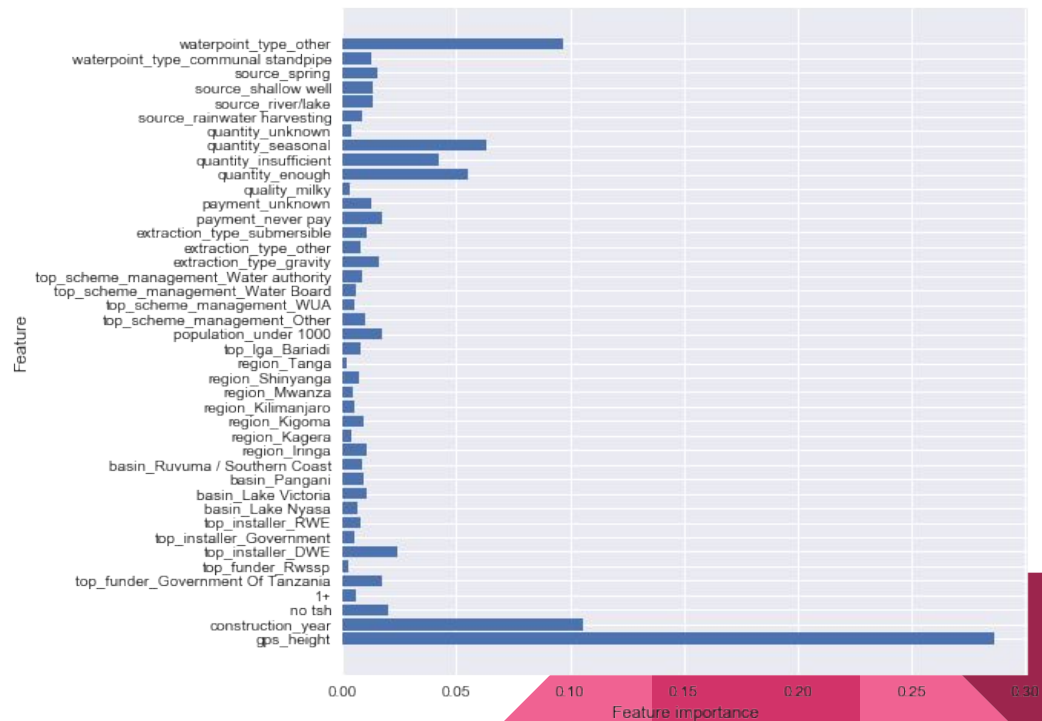
# Building a Classification Model

- Boruta algorithm for feature selection
- Used best features in 3 models:
  - Logistic regression: 72.4
  - Decision tree: **75.5**
  - Random forest: 70



# Best Features in Model

- GPS height - Water level of the well
- Year the well was constructed
- Type of well



## Next Steps

- Prioritize and anticipate maintenance on older wells and uncommon types of wells (not communal standpipes or hand pumps)
- Analyze data through a GIS software
  - Perform a Geographically Weighted Regression model



Thank you!

