

---

# Predicting Faulty Water Pumps in Tanzania

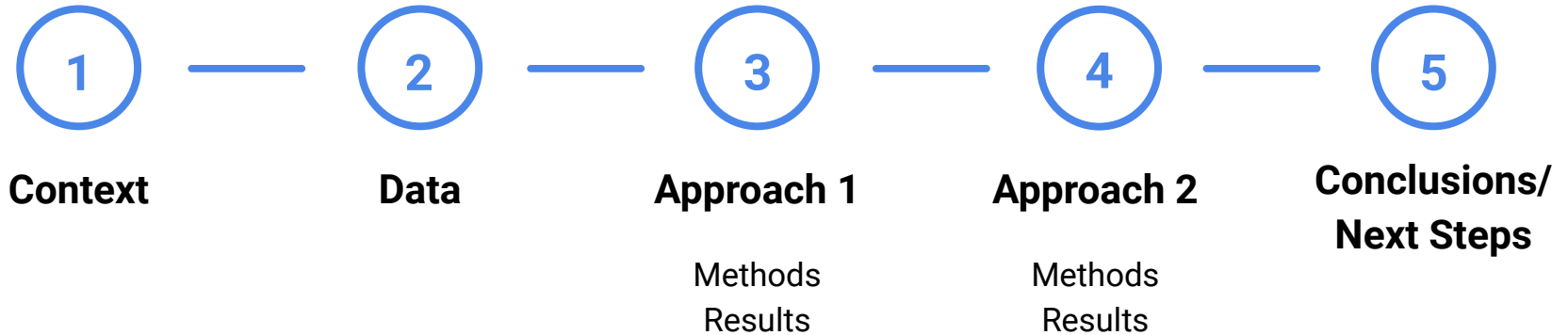
by Kaitlyn Zeichick



Source: Tobias Hayes

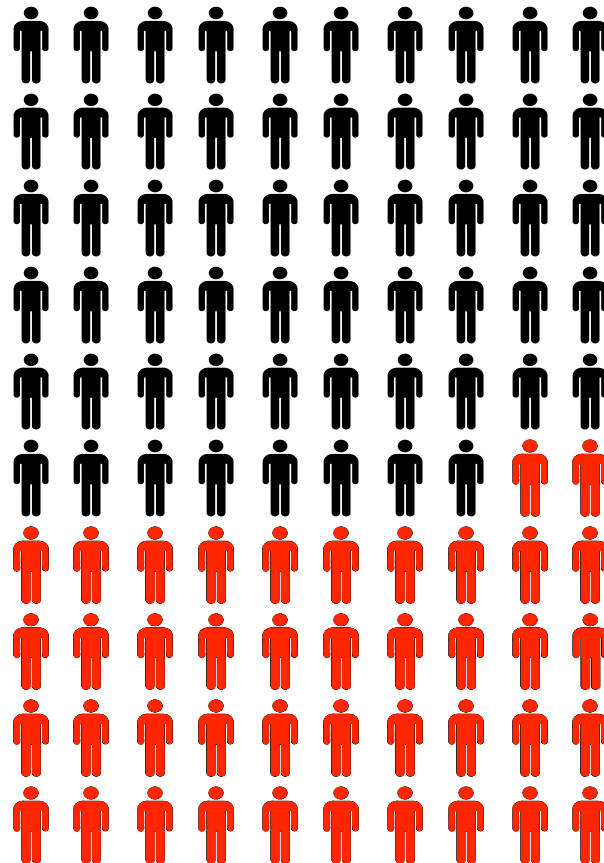


# Outline



# Tanzania Water Crisis

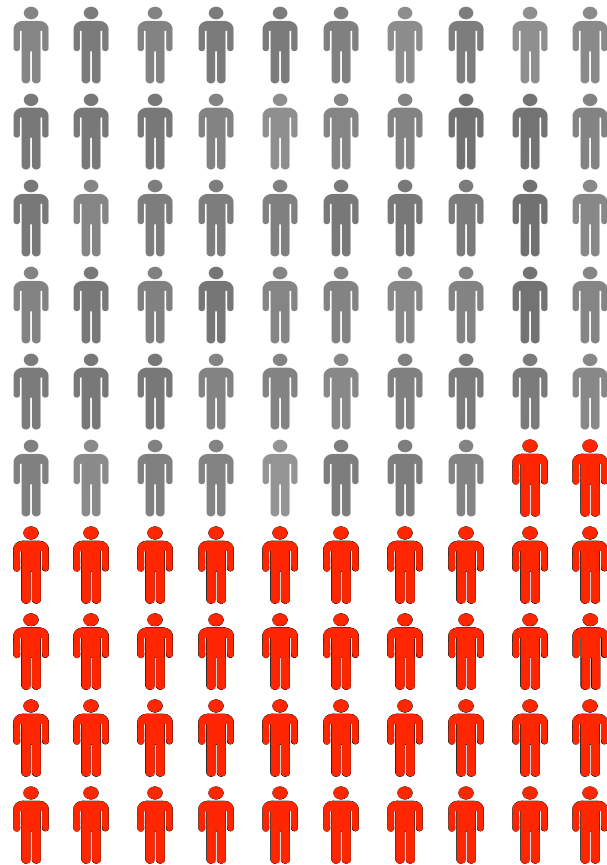
- 24 million people without basic access to safe water



# Tanzania Water Crisis



=





## Consequences

- Malaria
- Cholera





## Solution : Water Pumps!

- Pros:
  - Clean
  - Cheap
- Cons:
  - Easily broken





## Data



Water  
Management



Location



Water  
Usage



Water  
Related



Pump  
Features

**Functional**

---

**Non-Functional**



## Two Problem Statements

- 1) WHERE? Locate already existing faulty water pumps.
- 2) WHY? Identify features that are highly correlated with faulty pumps.



## Random Forest: Choosing a Metric

Consequence of many  
False Negatives:



Consequence of many  
False Positives:





## Random Forest: Results

Accuracy

79%

Recall

80%

# Logistic Regression

Why and how to water pumps break?





## Logistic Regression

Accuracy

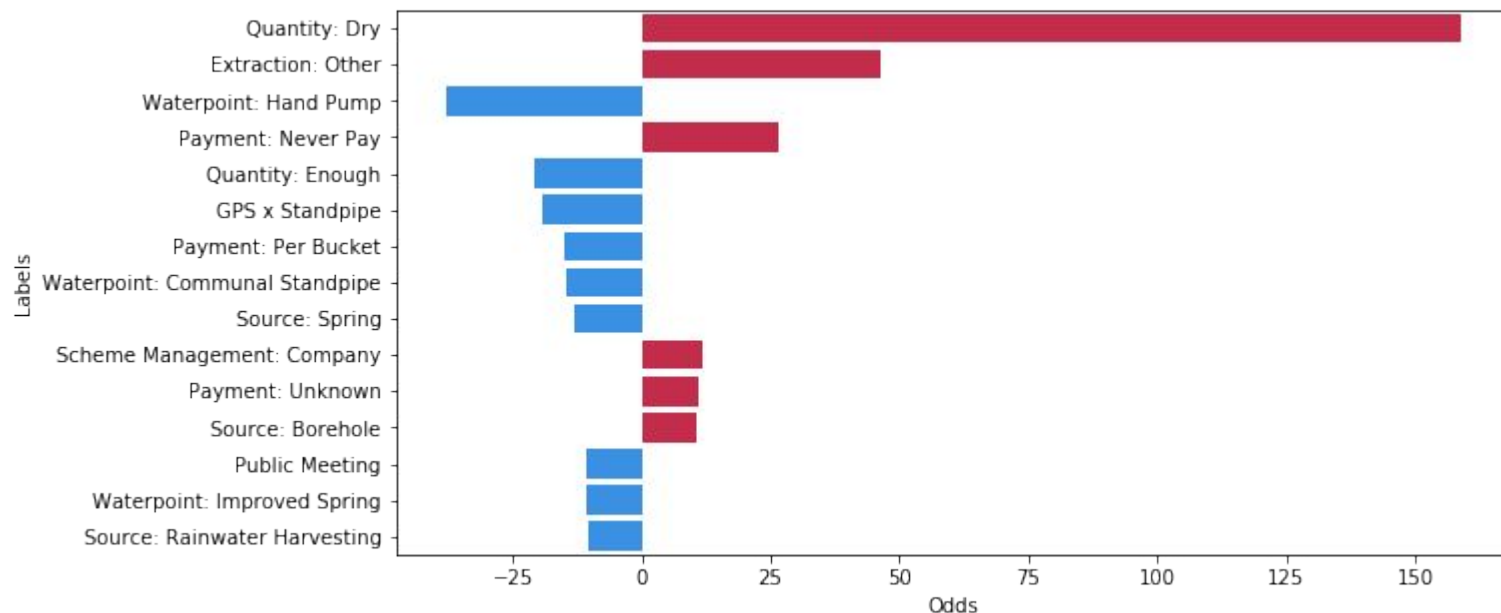
70%

Recall

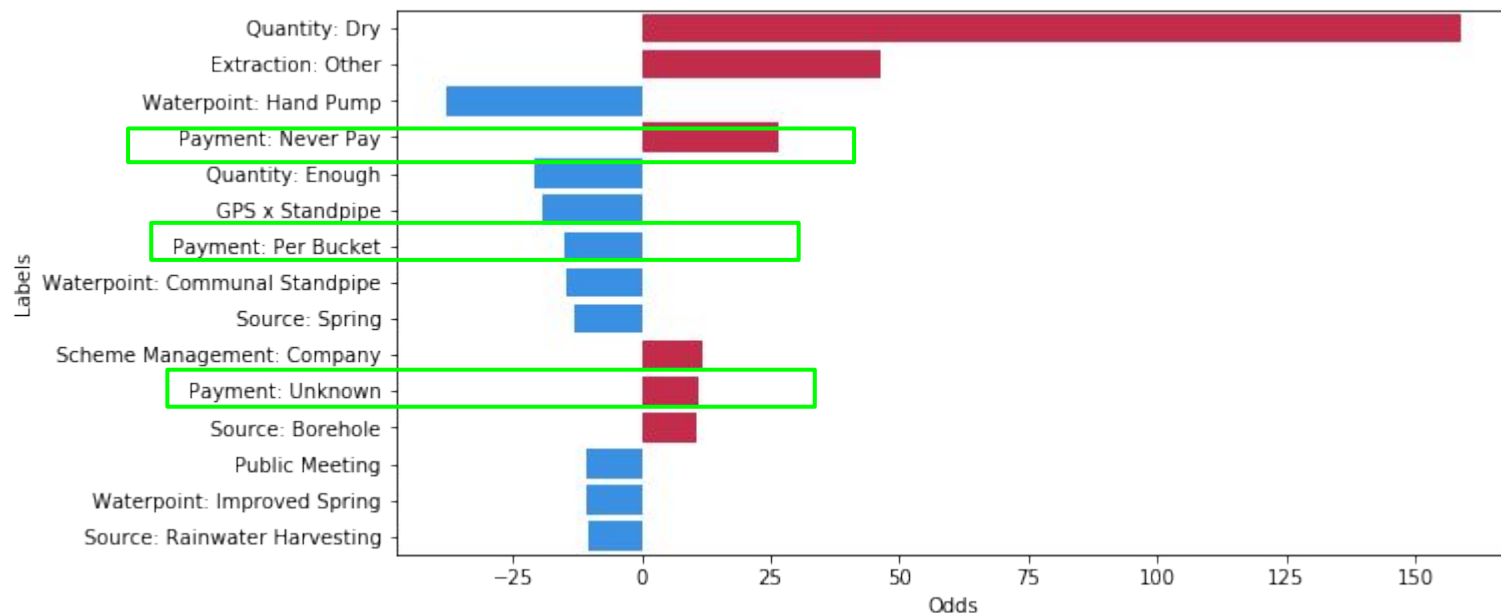
68%



## Logistic Regression: Top 15 Features



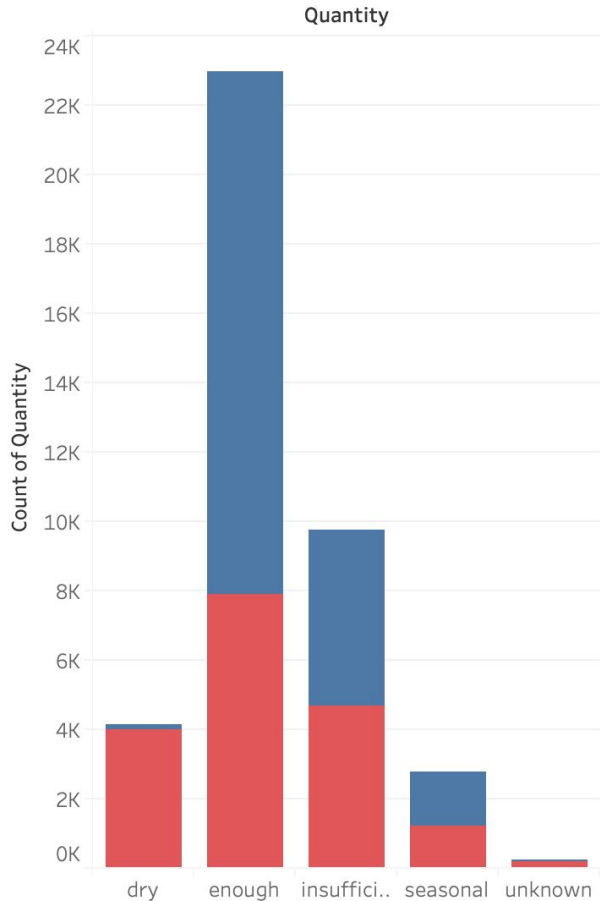
## Logistic Regression: Top 15 Features





## Logistic Regression: Top Feature

- Increases odds of pump being non-functional by 159%
- Explanation:
  - Dry running
  - Defined as non-functional





## Conclusions

- Identify broken water pumps:
  - Random Forest - 79% accuracy
- Identify important features:
  - Logistic Regression - 70% accuracy
  - Quantity: Dry



## Future work

- Improve logistic regression model
  - More feature engineering
- Improve random forest model
  - Ensemble methods
- Interpretation



**Thank you!**  
**Questions?**



# Waterpoint Type Group

cattle trough

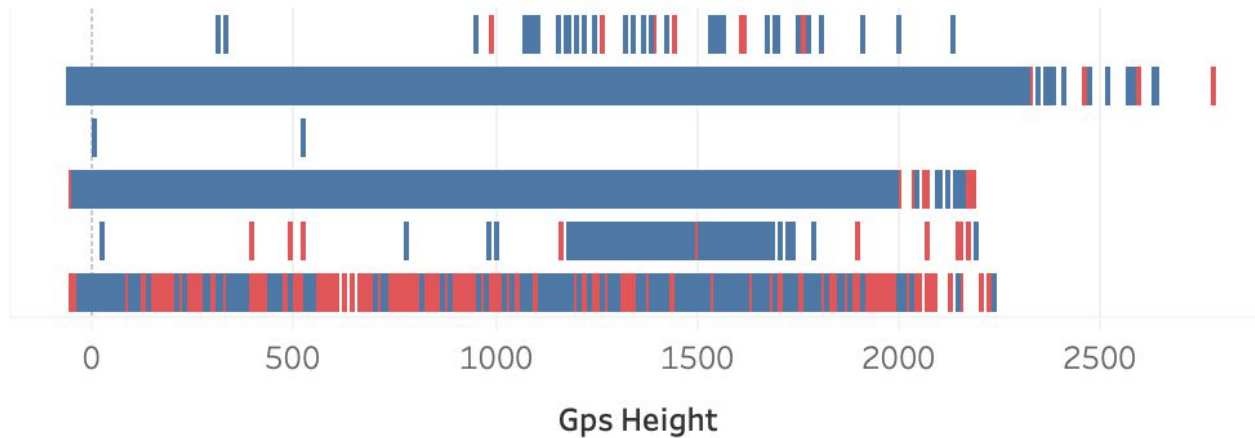
communal standpipe

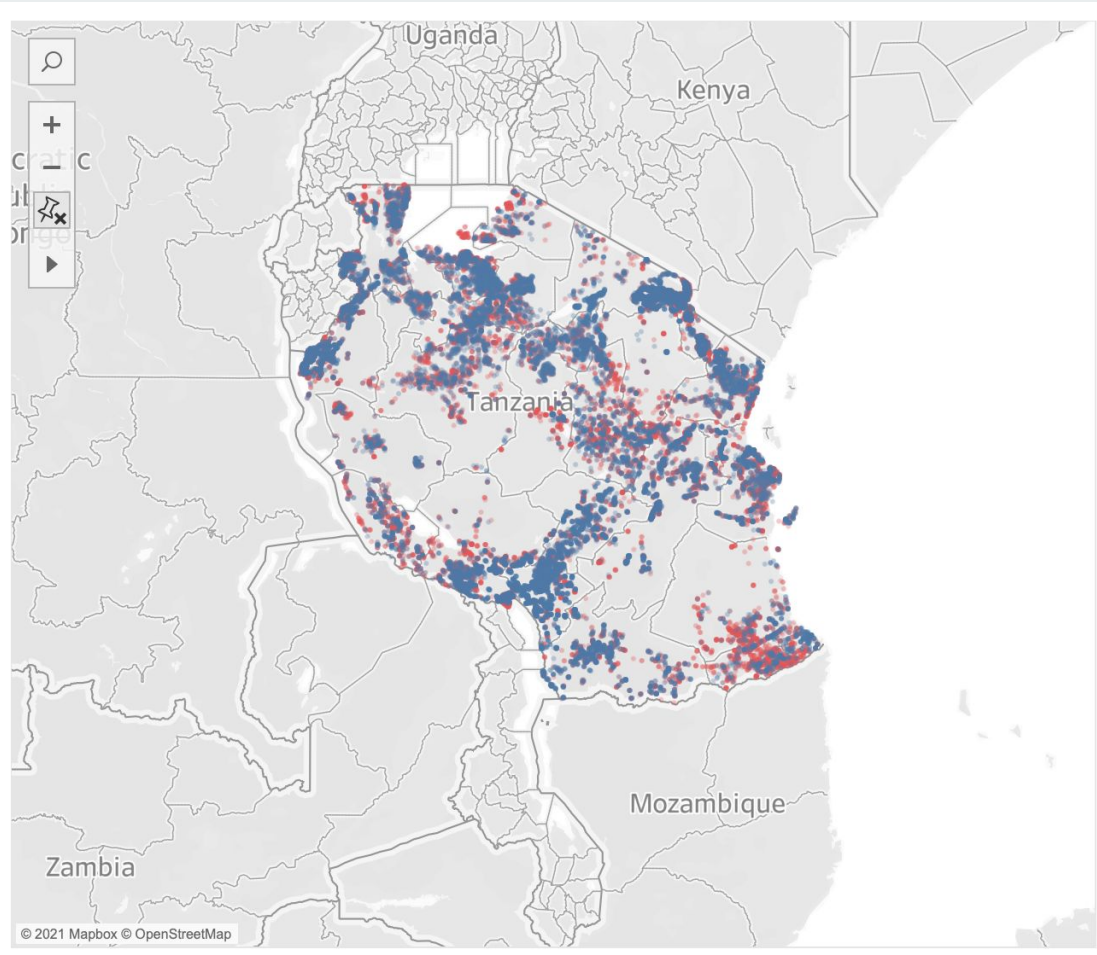
dam

hand pump

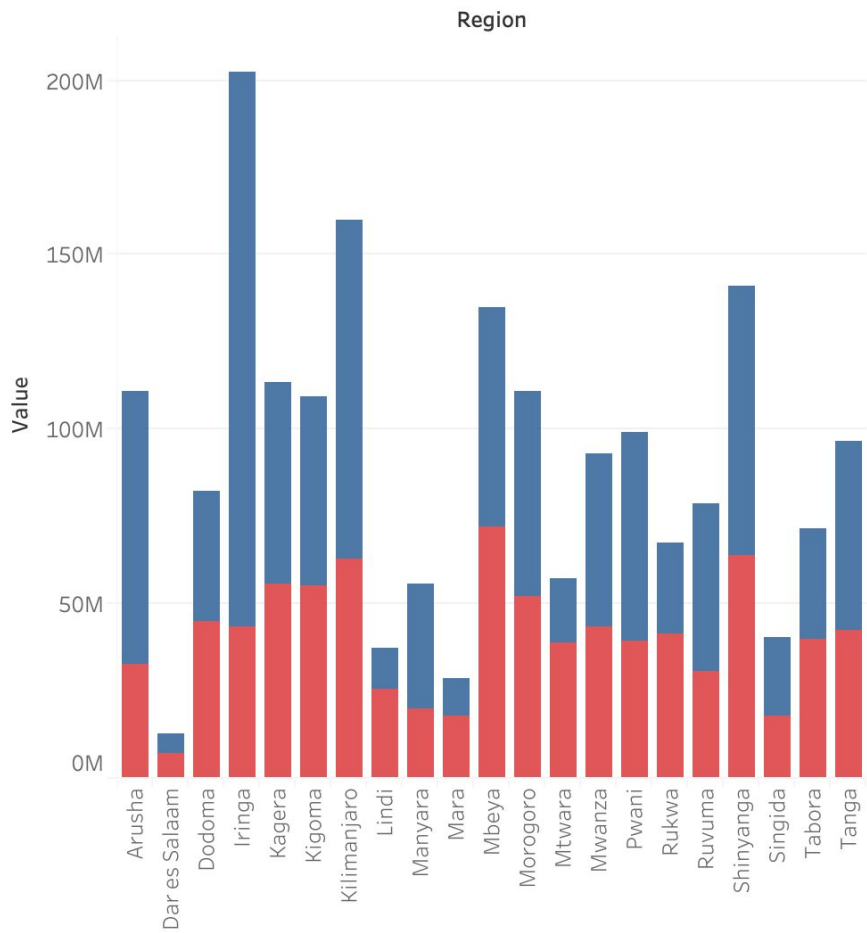
improved spring

other









Visualizing Important Features

