

---

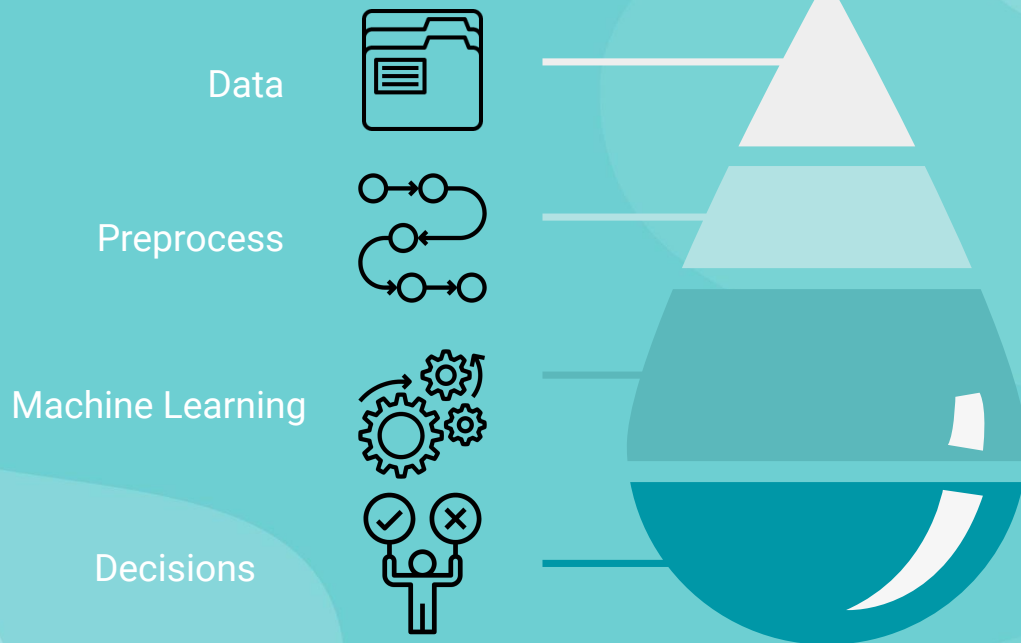
# **PUMP IT UP! DATA MINING THE WATER TABLE**

Data Science for Social Good

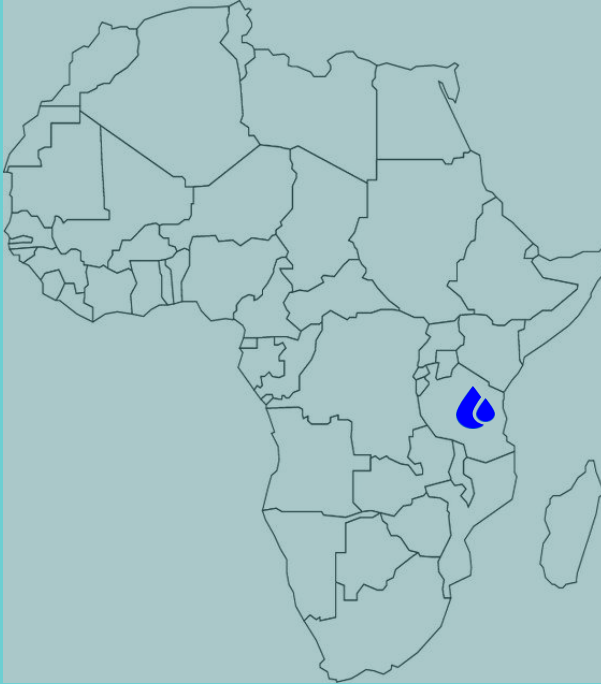
---

**HOW CAN WE APPROACH THE WATER SCARCITY IN TANZANIA ?**

# MACHINE LEARNING TECHNIQUES CAN HELP OUT



# TANZANIA WATER AND SANITATION CRISIS



ABOUT 59  
MILLION



40 % URBAN  
60 % RURAL

50 MIN ROUND TRIP  
TO ACCESS WATER  
POINTS



25 MILLION LACK OF SAFE  
ACCESS TO SAFE WATER

NO PROPER SANITATION → WATERBONE DISEASES

# WHAT DATA DO I HAVE AND WHAT ARE THE NEXT STEPS ?

## STARTING POINT



59400 WATERPOINTS  
41 FEATURES:  
GEOGRAPHICAL  
TECHNICAL  
DEMOGRAPHIC

## STEPS



DATA CLEANING  
DATA EXPLORATION  
FEATURE ENGINEERING AND MACHINE  
LEARNING

## TARGET

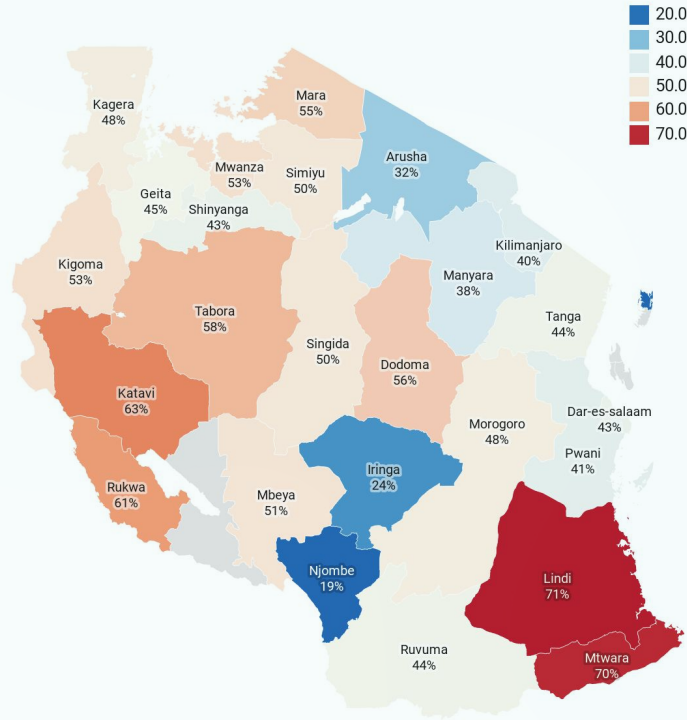


WATER PUMP FUNCTIONALITY  
# FUNCTIONAL  
# NON FUNCTIONAL  
# FUNCTIONAL - NEEDS REPAIRS

# MOST REGIONS HAVE >50 % NON FUNCTIONAL WELLS!!

## Proportion of non functional pumps

Lindi and Mtwara = Highest defect rate

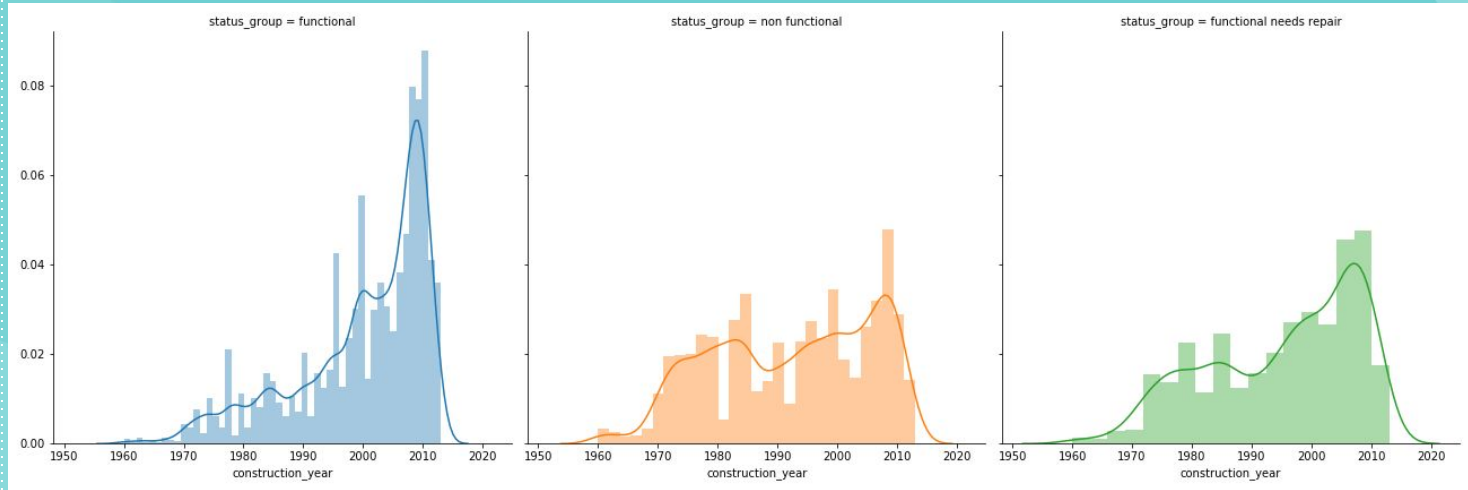


Created with Datawrapper

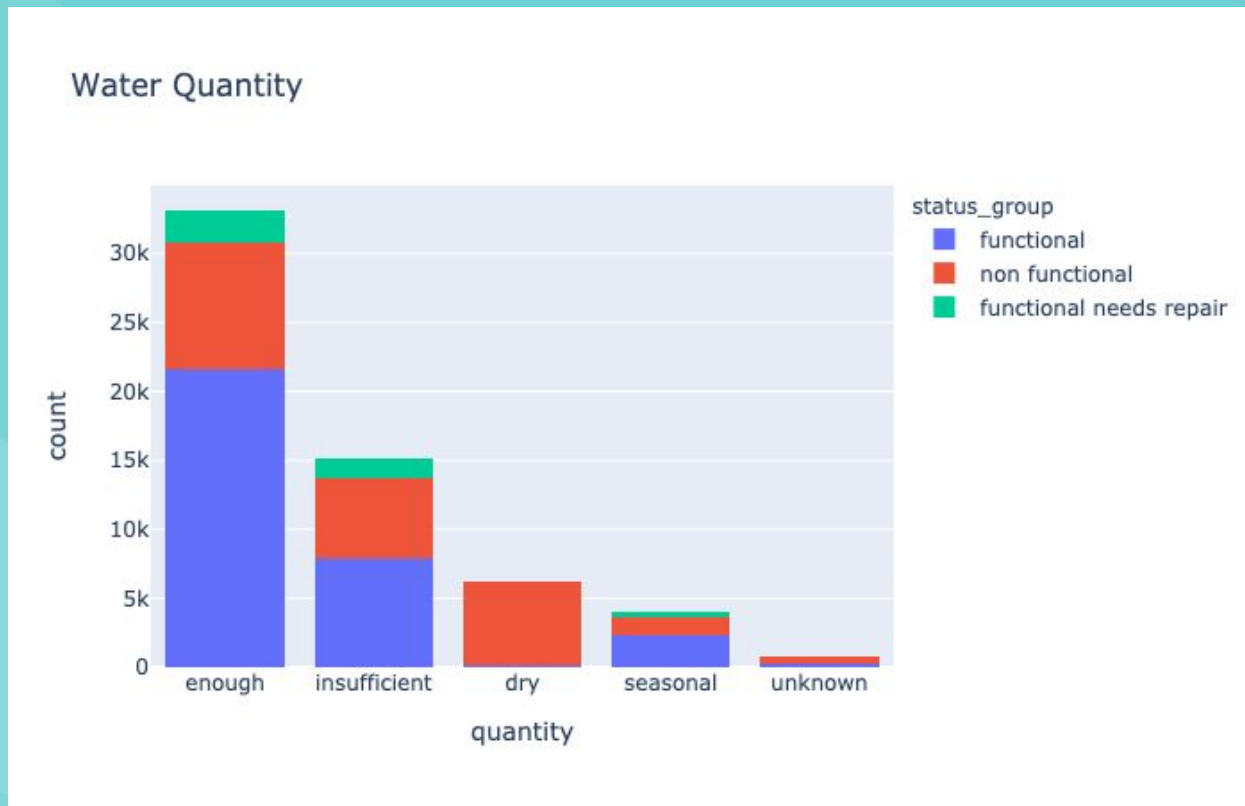
LINDI AND MTWARA LESS FUNCTIONAL  
REGIONS

NJOMBE IRINGA MOST FUNCTIONAL  
REGIONS

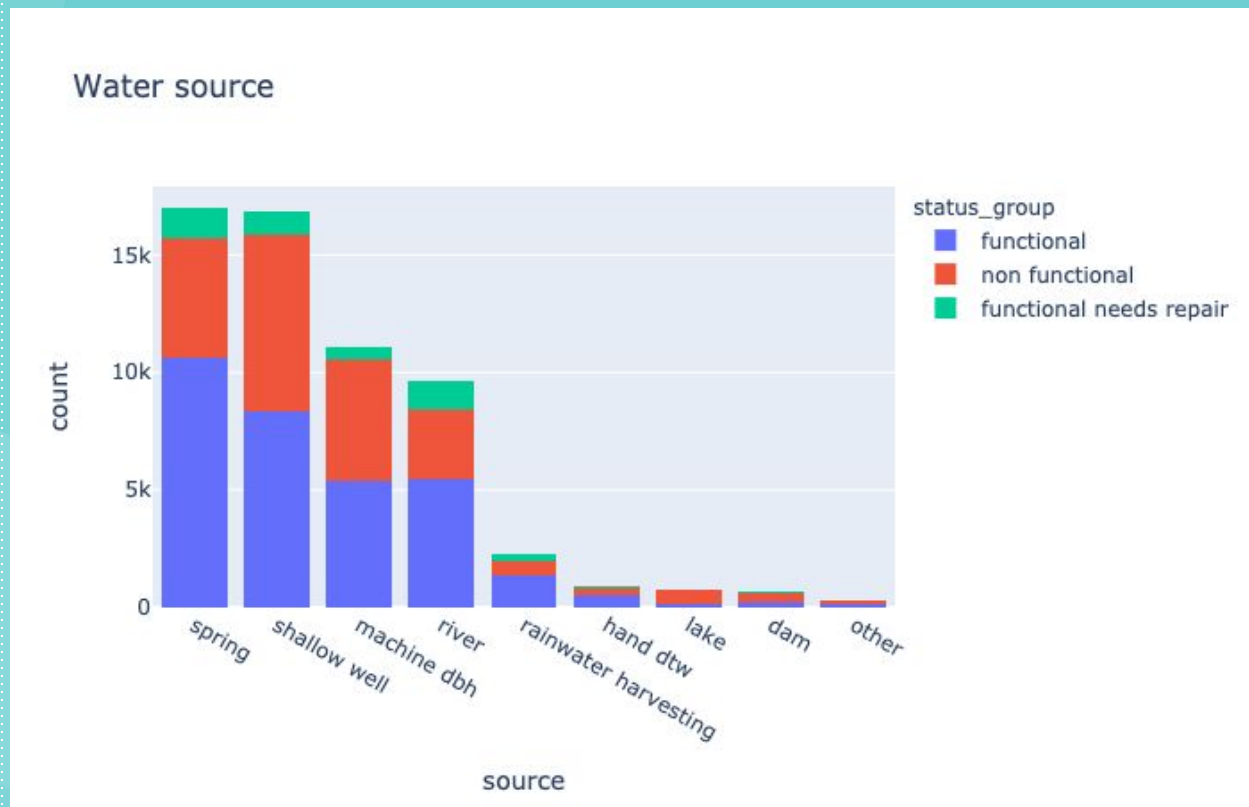
# WATER PUMP AGE DEFINES FUNCTIONALITY



# SUFFICIENT WATER DEFINES FUNCTIONALITY

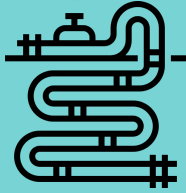


# WATER SOURCE DEFINES FUNCTIONALITY





# PREDICTIVE MODELLING METHODOLOGY



PREPROCESSING PIPELINE

LOGISTIC REGRESSION

K-NEAREST NEIGHBORS

DECISION TREE

RANDOM FOREST

ADABOOST

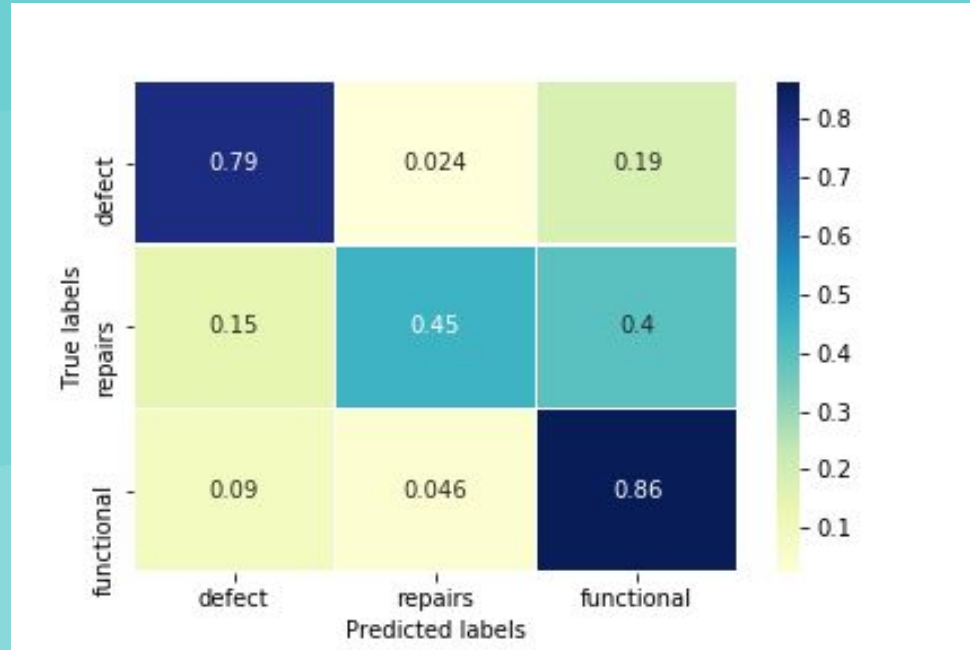
XGBOOST

NEURAL NETWORK



RANDOMIZED GRID SEARCH

## XGBOOST: BEST MODEL WITH 80 % PRECISION



## **AND NOW ?**

**FUNDING OF PROJECTS IN REGIONS  
WHICH HAVE LESS FUNCTIONAL WELLS**

**PAY ATTENTION TO PARAMETERS WHICH  
DETERMINE FUNCTIONALITY OF WELLS**

**OPTIMIZE IMPUTING STRATEGY OF  
MISSING VALUES**

**APPLY/OPTIMIZE NEURAL NETWORK  
MODEL**

**CREATE AN APPLICATION**

LAST BUT NOT LEAST

ANY QUESTIONS?

THANK YOU FOR YOUR ATTENTION



[github.com/roble-chris](https://github.com/roble-chris)



[www.linkedin.com/in/christianrobledo90](https://www.linkedin.com/in/christianrobledo90)