

Pump it Up: Data Mining the Water Table



Deepali Sharma
January, 2023



Tanzania begins water rationing due to drought



- **Stakeholder:**
 - Ministry of water; Government of Tanzania
- **Business Problem:**
 - Predict the pumps that are functional/faulty/need repairs
 - Proper Allocation of resources based on pump status
 - Find other factors that are responsible for non-functioning pumps

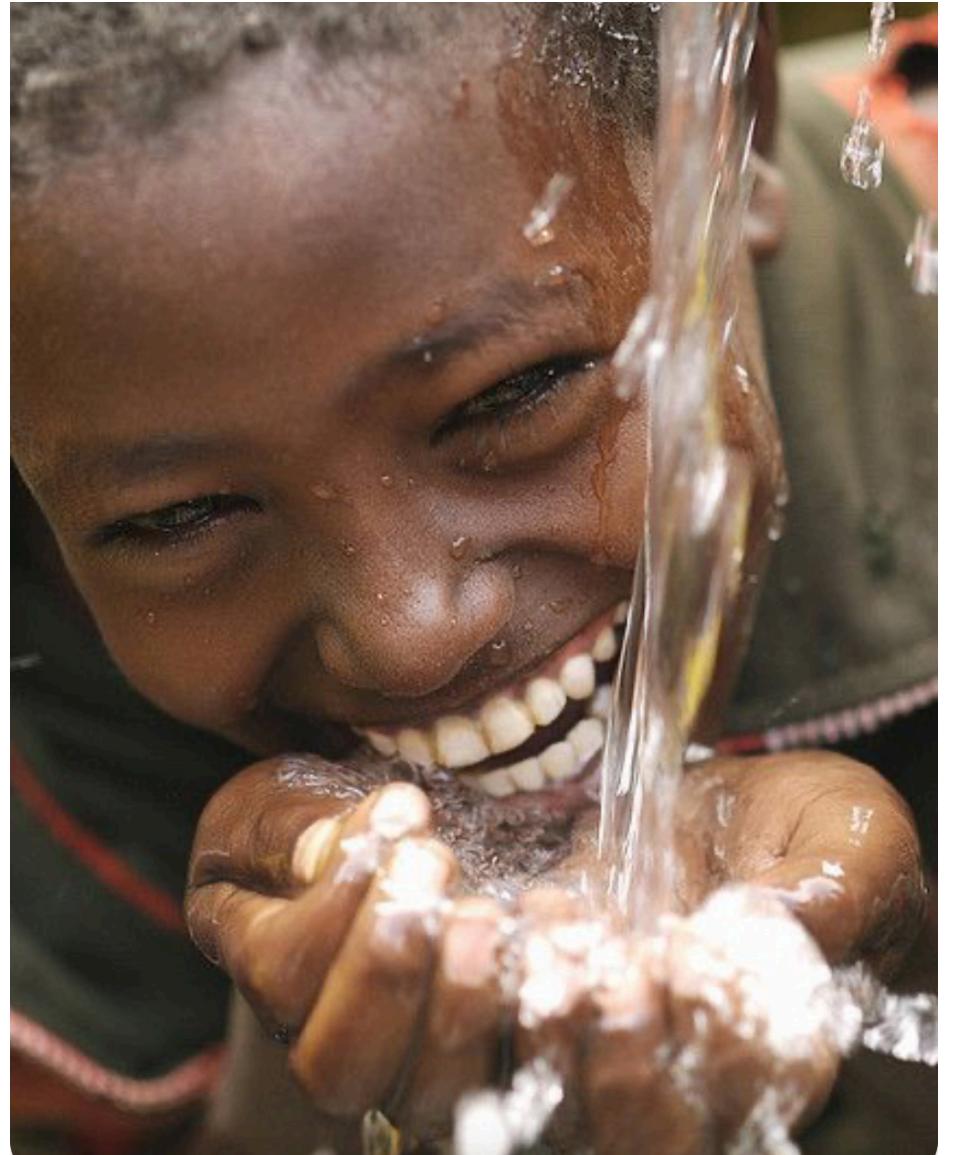
<https://www.africanews.com/2022/10/28/tanzania-begins-water-rationing-due-to-drought//>

Data:



- Data collected by **Taarifa** and Tanzanian Ministry of Water
 - Information contains geographical locations, quality of water, water **extraction** method, installed/funded by etc.etc

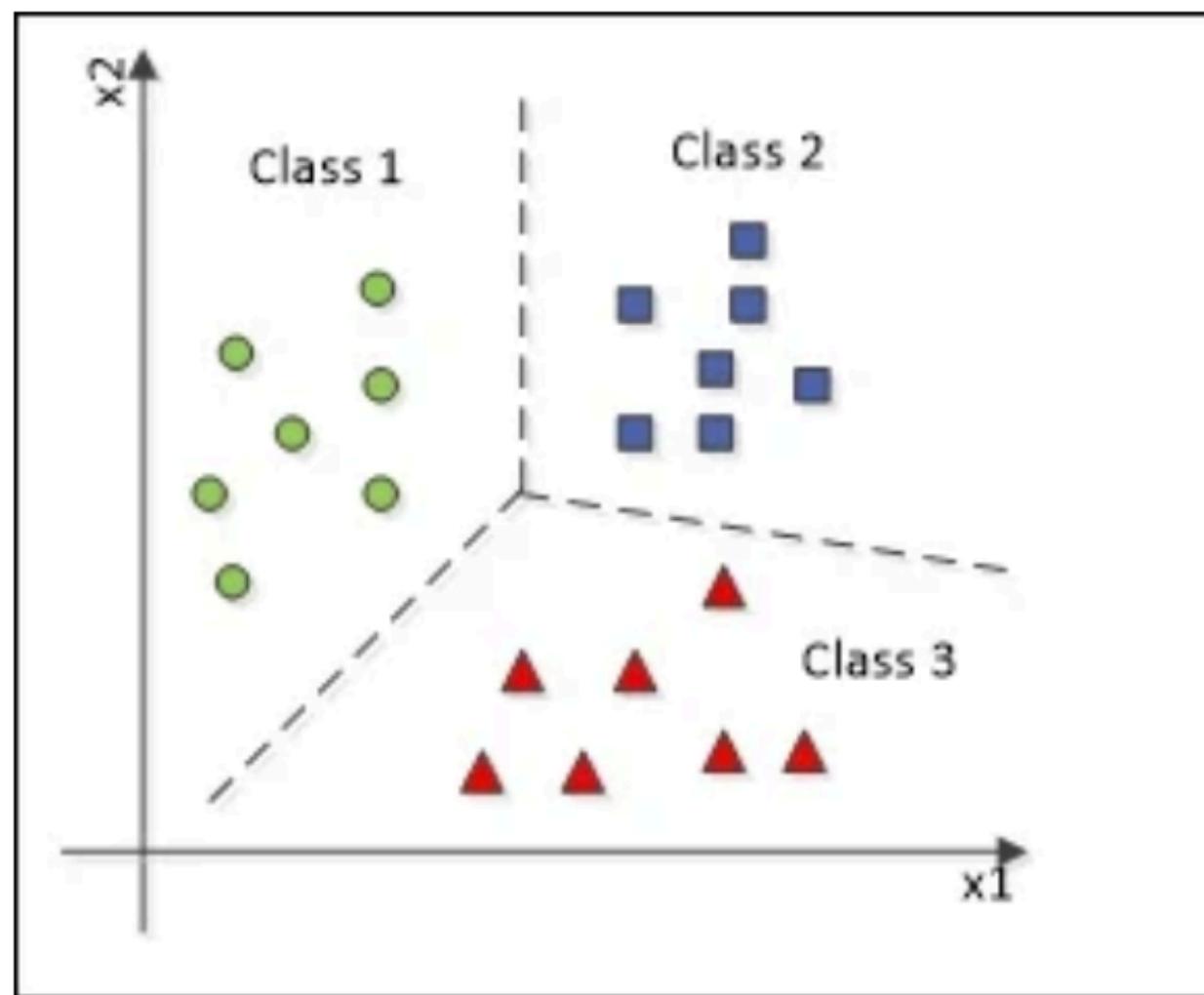
Goal:



- Goal is to build a **classifier** that would predict the status of water pump utilizing the available information
 - appropriately **redistribute crews** and funds to fix non-functioning pumps
 - where **not to build pumps**, which **extraction techniques** to use etc.

Modeling Algorithm

- Clean and pre-process the data
- Find the model that best describes the data
- Check the model against validation and test data
- Find the most important features that affect pump functionality
- Use the knowledge gained to make recommendations that will improve the pump status



Classification Problem

Results from the Best Model:

Model Performance:

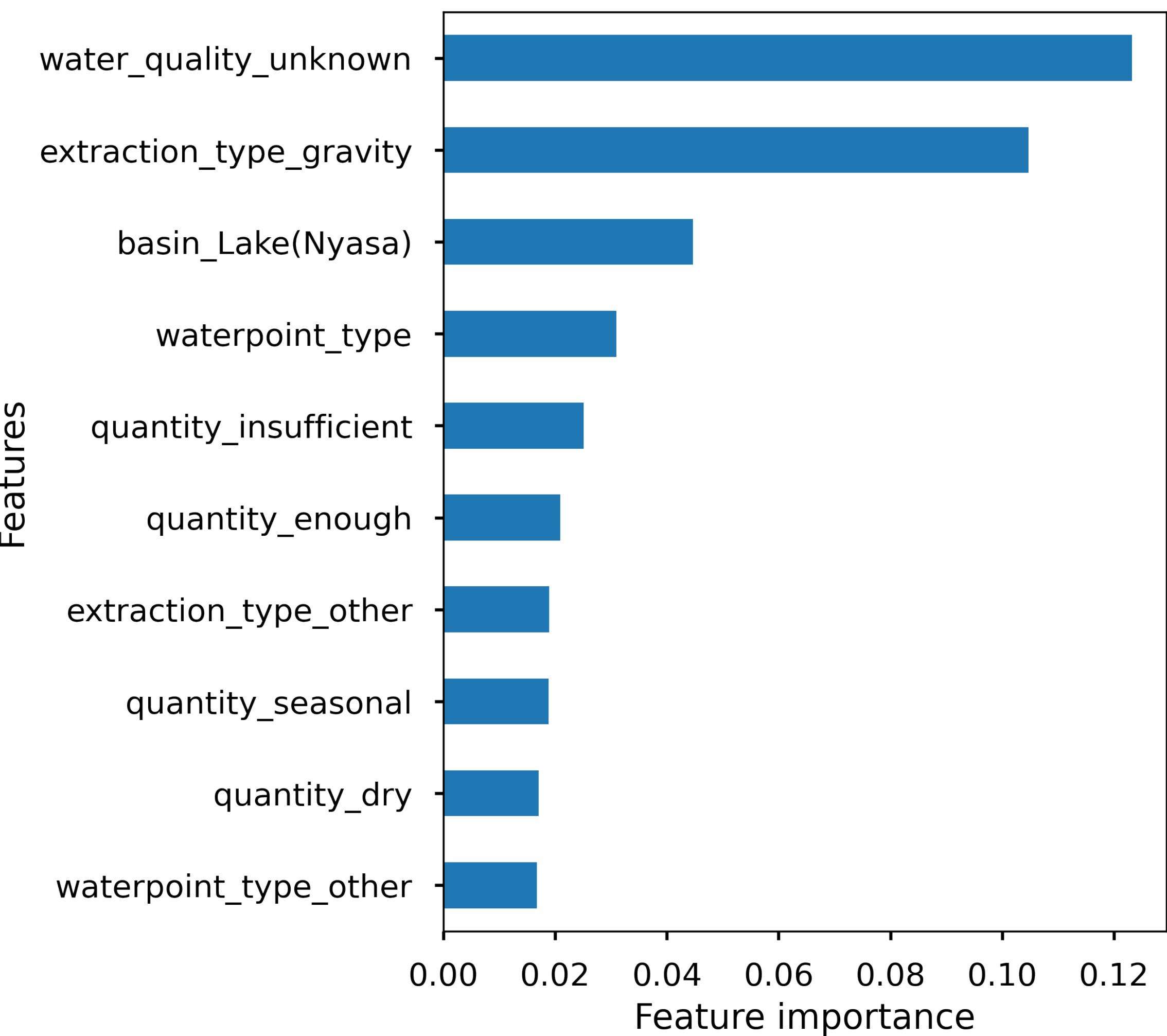
- ROC_AUC:
 - **overall** performance between TP and FP rates
- ACCURACY:
 - **correct** predictions
- RECALL:
 - **correct true predictions** for a given class
(TP+FN)
- PRECISION:
 - correct true predictions out of all predictions
(TP+FP)

	TRAIN	TEST
ROC_AUC	91%	88%
ACCURACY	80%	78%
RECALL	78%	78%
PRECISION	80%	78%

Results from the Best Model:

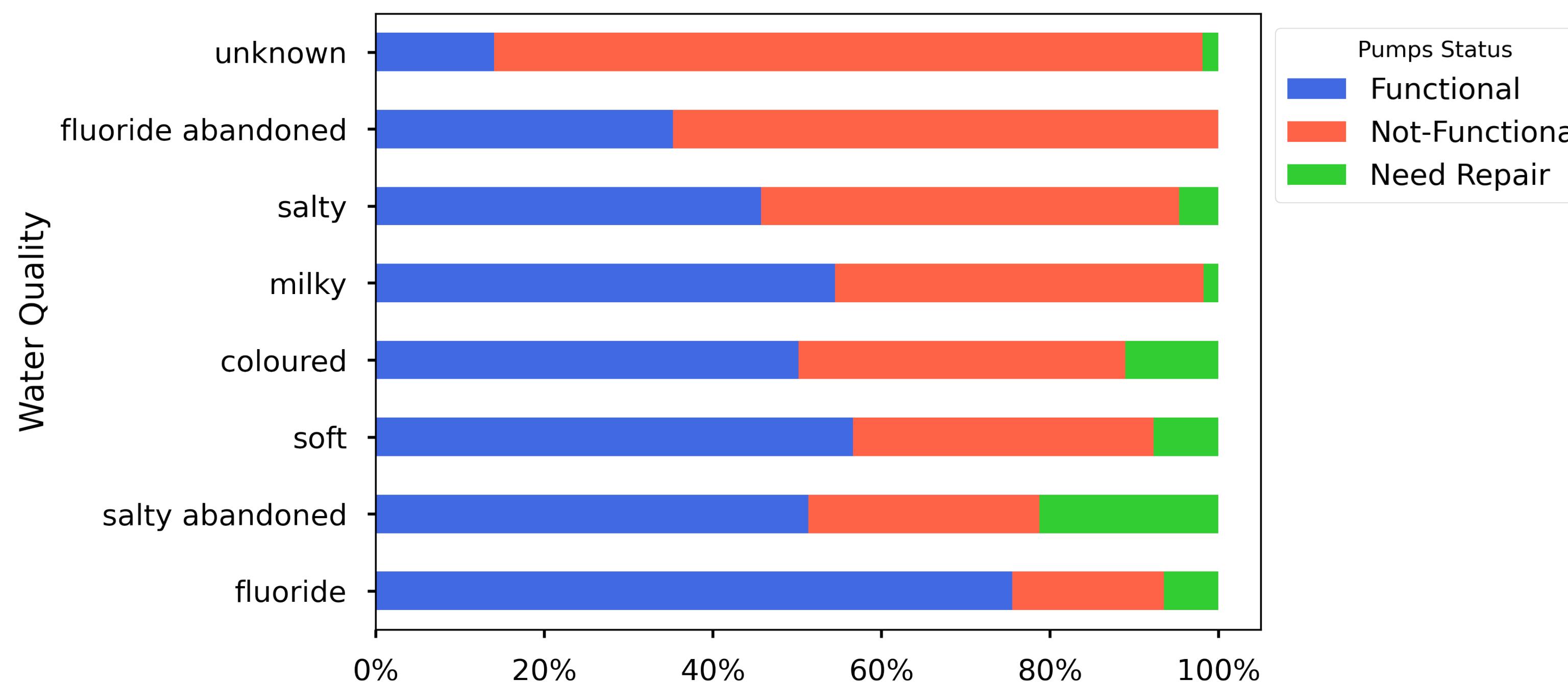
- Top 3 Features
 - Water Quality
 - Extraction Type
 - WaterPoint Type

**Relative Importance of Top 10 Features
for Predicting Water Pump Status**



Water Quality: the top feature

Relationship between Water quality and Pump status



- **salty water and unknown** type water pumps fall in the category of non-functional/ need repairs category

Salt Water Pumps

Salt water pumps are commonly equipped with materials such as stainless steel, aluminum, and thermoplastics, that perform better than cast iron in corrosive environments. Absolute Water Pumps suggests thorough flushing of your pump's inside and out with fresh clean non-salty and non-corrosive water after every use in a salt water or other corrosives environment.

Items 1-12 of 427

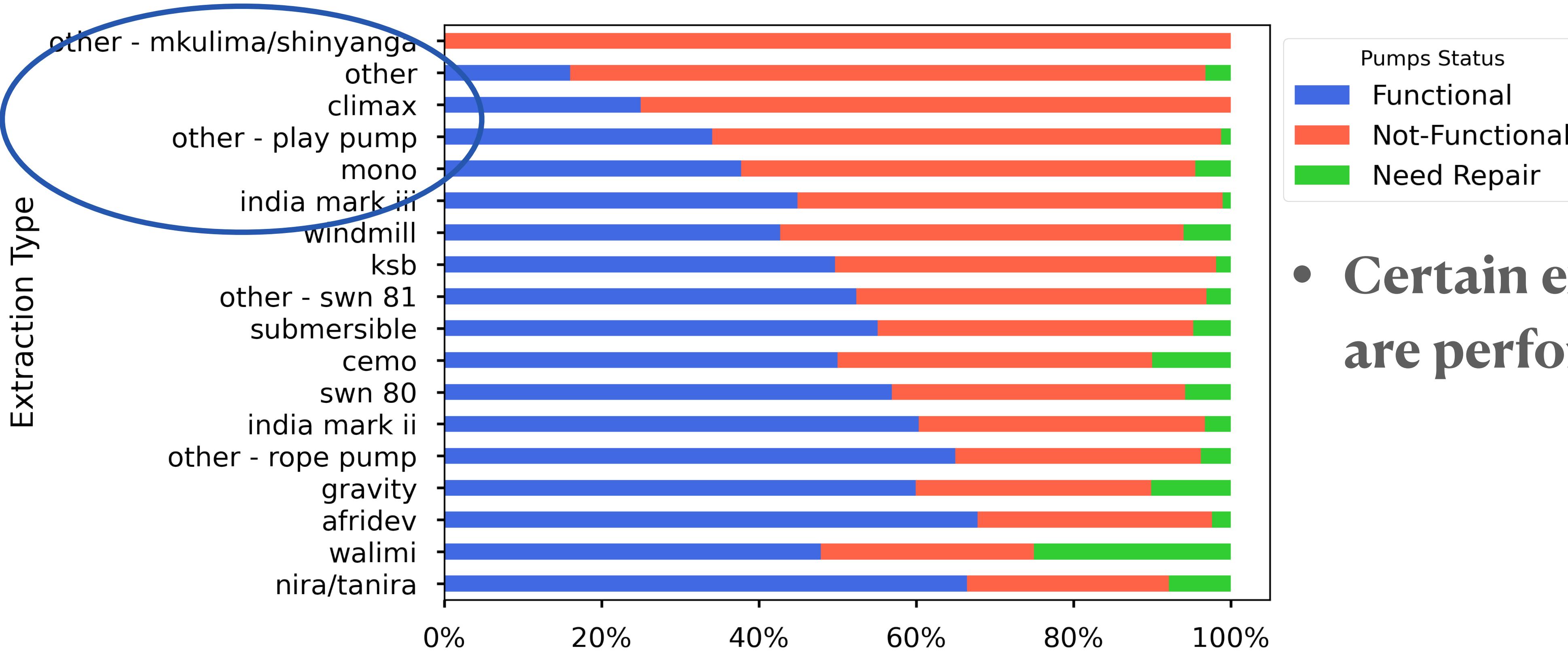
Sort By Special Price ▾



- **prioritize fixing water pumps with appropriate type e.g. aluminum pumps for salty locations!**

Extraction type: other important feature

Relationship between Extraction Method and Pump status

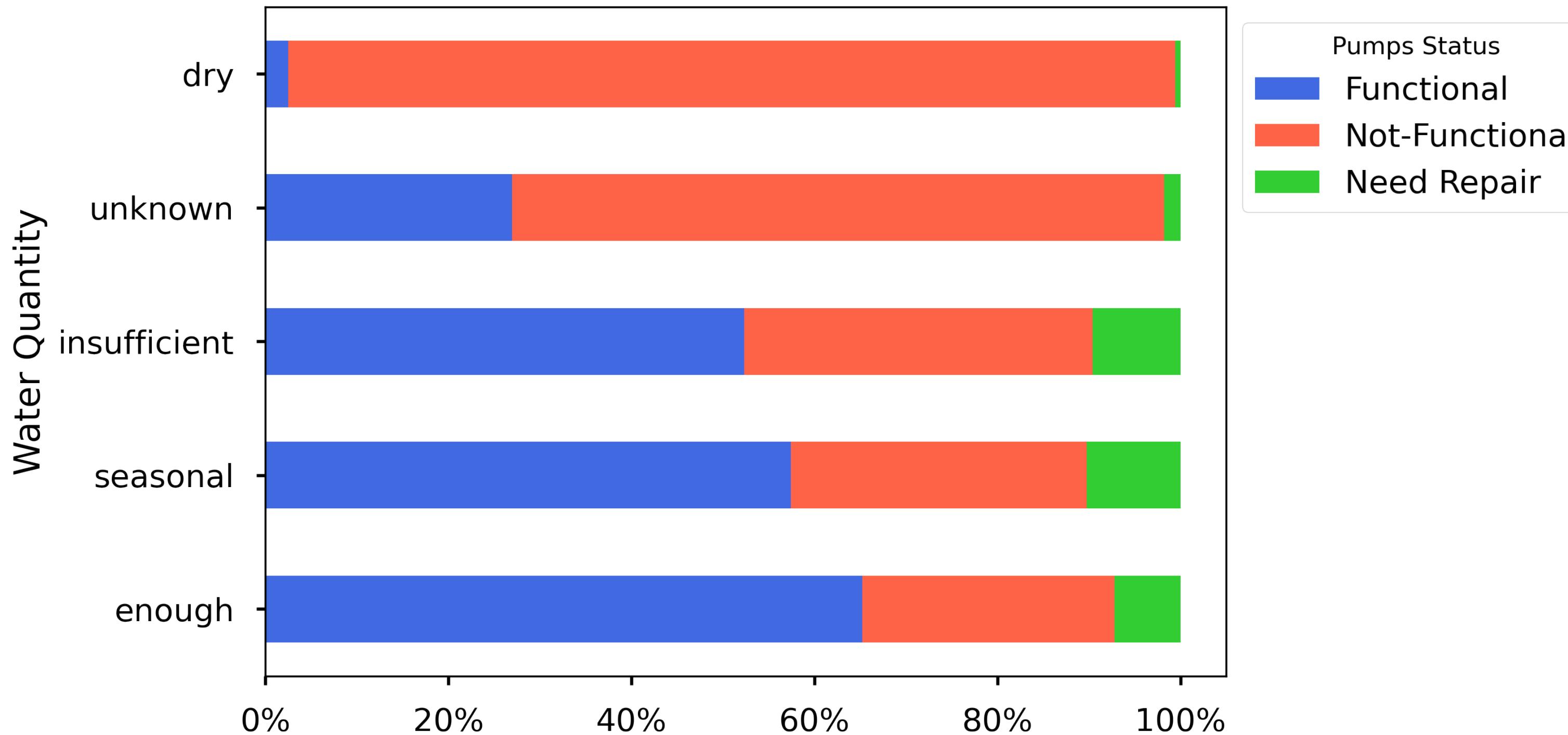


- Certain extraction types are performing worse!

- Only invest in the extraction techniques that are robust!

Water quantity: another important feature

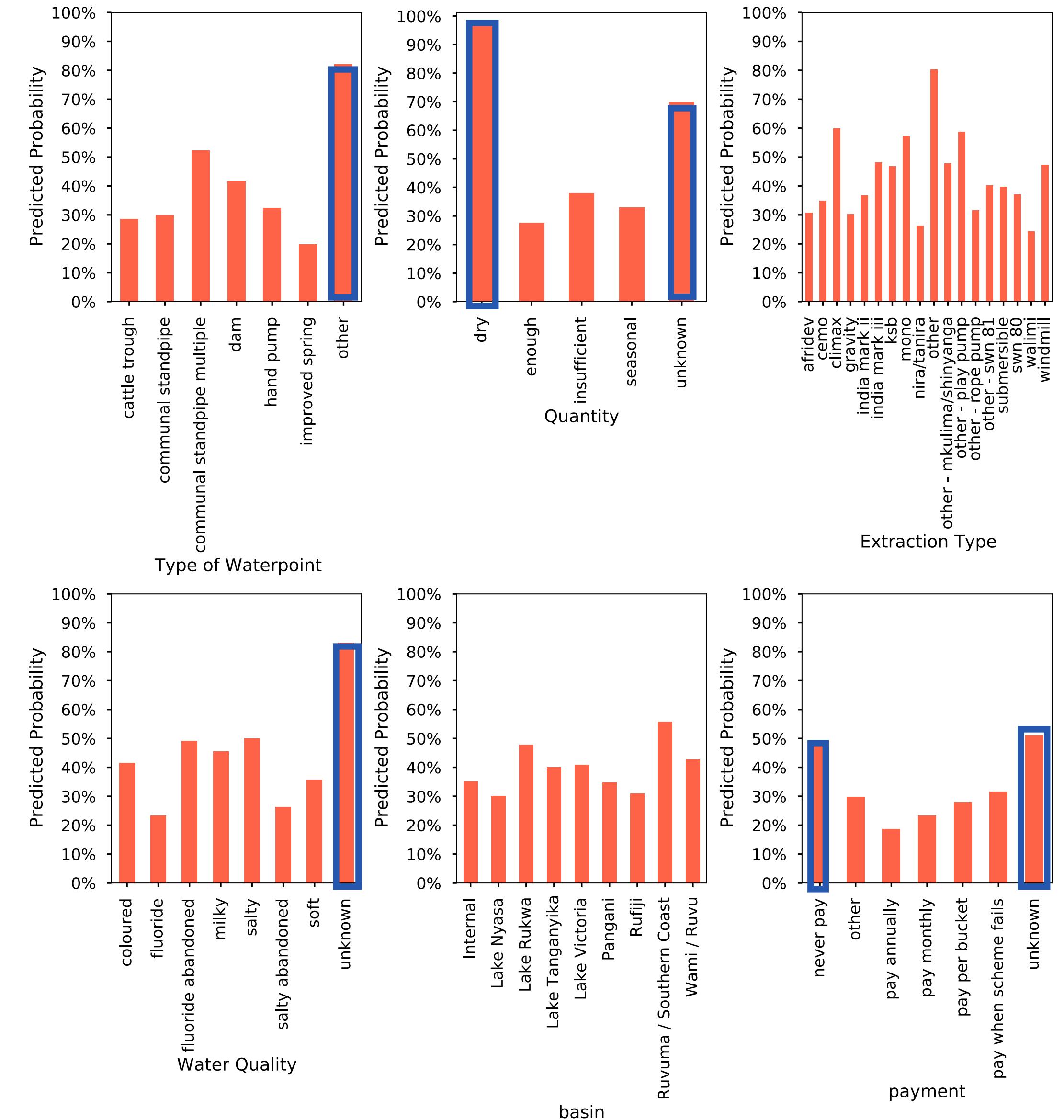
Relationship between Quantity and Pump status



- Pumps in **dry region** are quoted as **non-functional** or just labeled as such.
- Large fraction of **unknown pumps** are also non-functional
- Investigate the pumps in **dry** and **unknown** categories!

Predicted Probabilities

**Predicted Probability of Pumps that are not-functional
in Relation to Most Important Features**



- Predicted probabilities for non-functional pumps are higher among these categories
 - Pumps types:
 - **other**
 - Quantity:
 - **dry or unknown quantity**
 - Water quality:
 - **unknown** followed by **salty**
 - Payment:
 - **never pay** or **unknown**

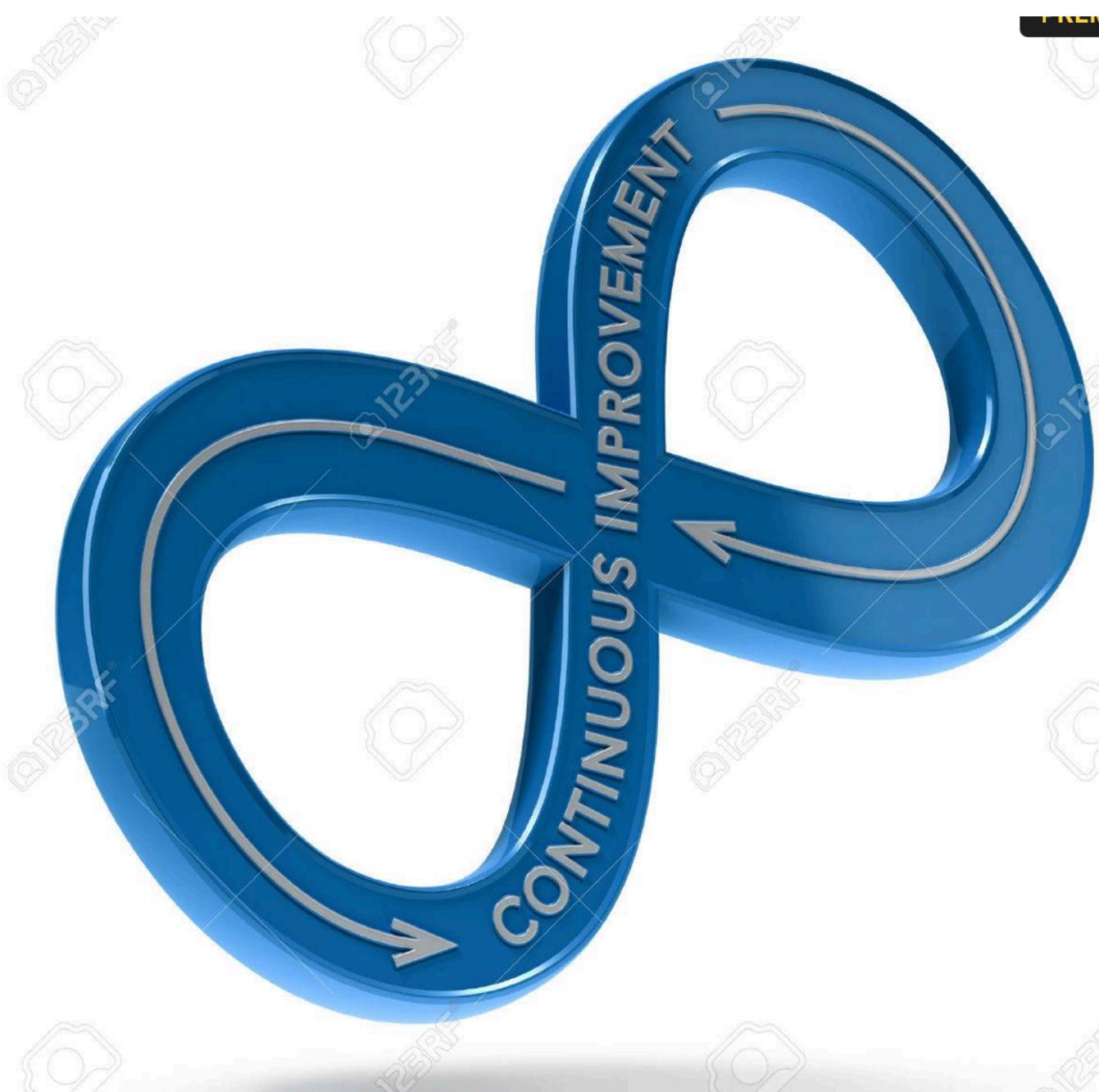
Some Recommendations:



Recommendations

- **Reallocate funds to replace pumps at salty locations.**
- **Budgets for R&D to find out the right extraction type for a given location.**
- **Why for some pumps payments are never paid or unknown.**
- Look into the water pumps that are in **dry areas** and or where the **quantity of water is not known**

Limitations and Improvements



- Clever strategy to clean up some of the columns
- Optimize GridSearchCV running time
- Check other **classification models like CNN** to see their performance.
- Explore the **dependence on the scoring metric** used in the GridSearchCV
- Eliminate or quantify effect of parameters such as installer, funder, scheme_name on the model performance

THANK YOU!



- **Deepali Sharma:** email:(deeps.sharma@gmail.com, deepali@rcf.rhic.bnl.gov)
- **LinkedIn:** <https://www.linkedin.com/in/deepali-sharma-a83a126/>
- **GitHub:** <https://github.com/deepssharma>