

PUMP IT UP

TERNARY CLASSIFICATION MODEL

Bao Tram Duong

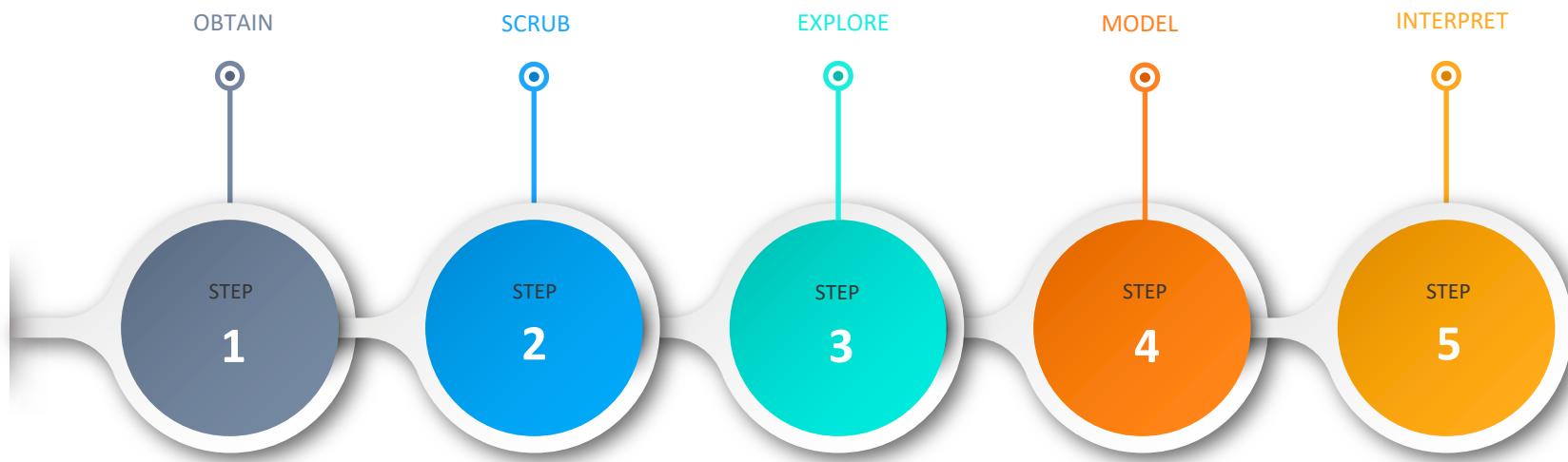


Tanzania, the largest country in East Africa, is suffering from a water crisis

- 4 million people lack access to an improved source of safe water
- 30 million people lack access to improved sanitation
- Water-borne illnesses, such as malaria and cholera, account for over half of the diseases affecting the population



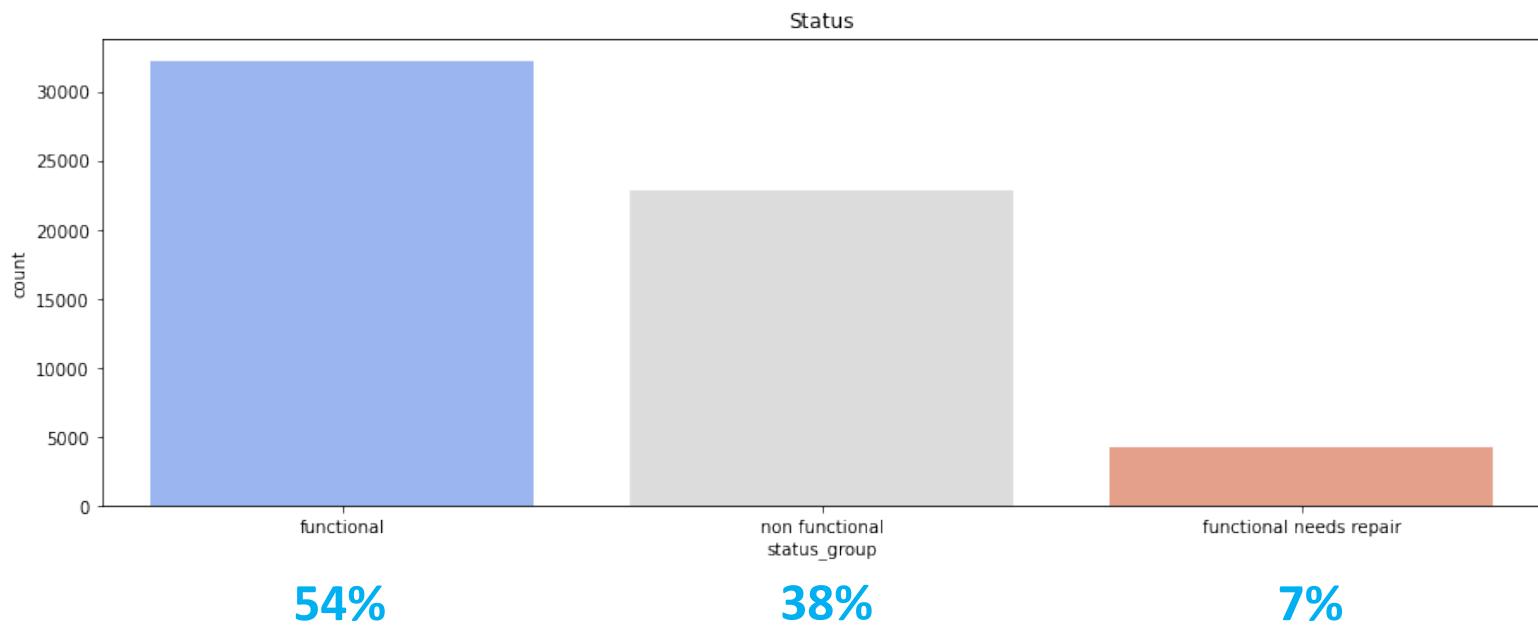
OSEMN



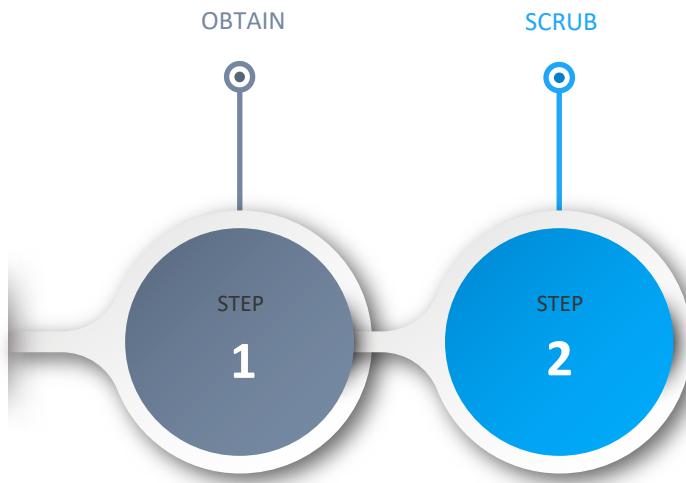


Using data provided by **The Tanzanian Water Ministry** and **Taarifa, DrivenData** began a competition to solve this problem by building a classification system to predict whether a given water source is working correctly.

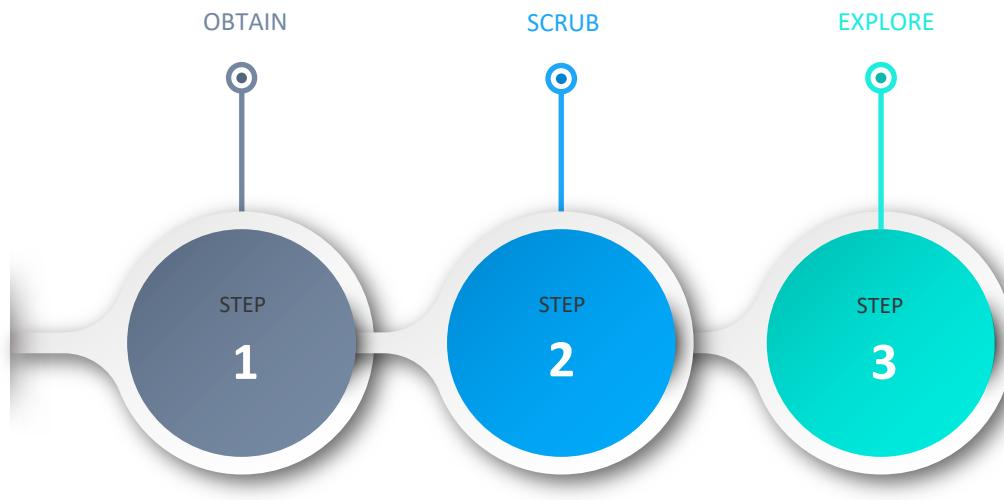
- 59,400 water points
- 40 features



Challenges

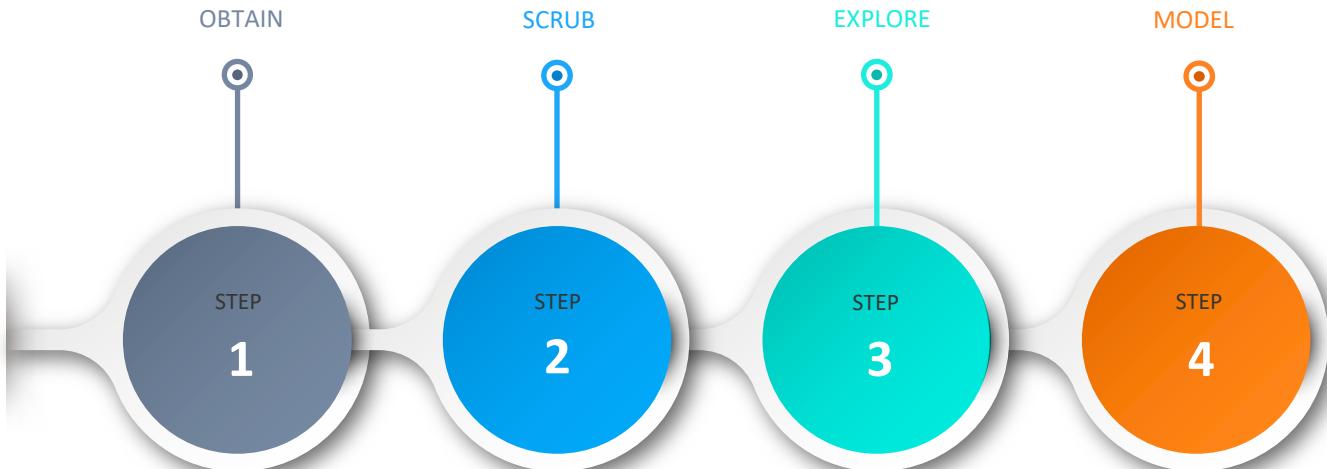


1. Replace **missing values** with mean, median, or classify them as 'other'
2. Remove **redundant features**
3. Fix **misspellings and variations**
4. Select for the 20 most common values with features that have **high number of unique values** and categorize the rest as 'other'
5. Correct **class imbalance** with SMOTE (Synthetic Minority Over-Sampling Technique)

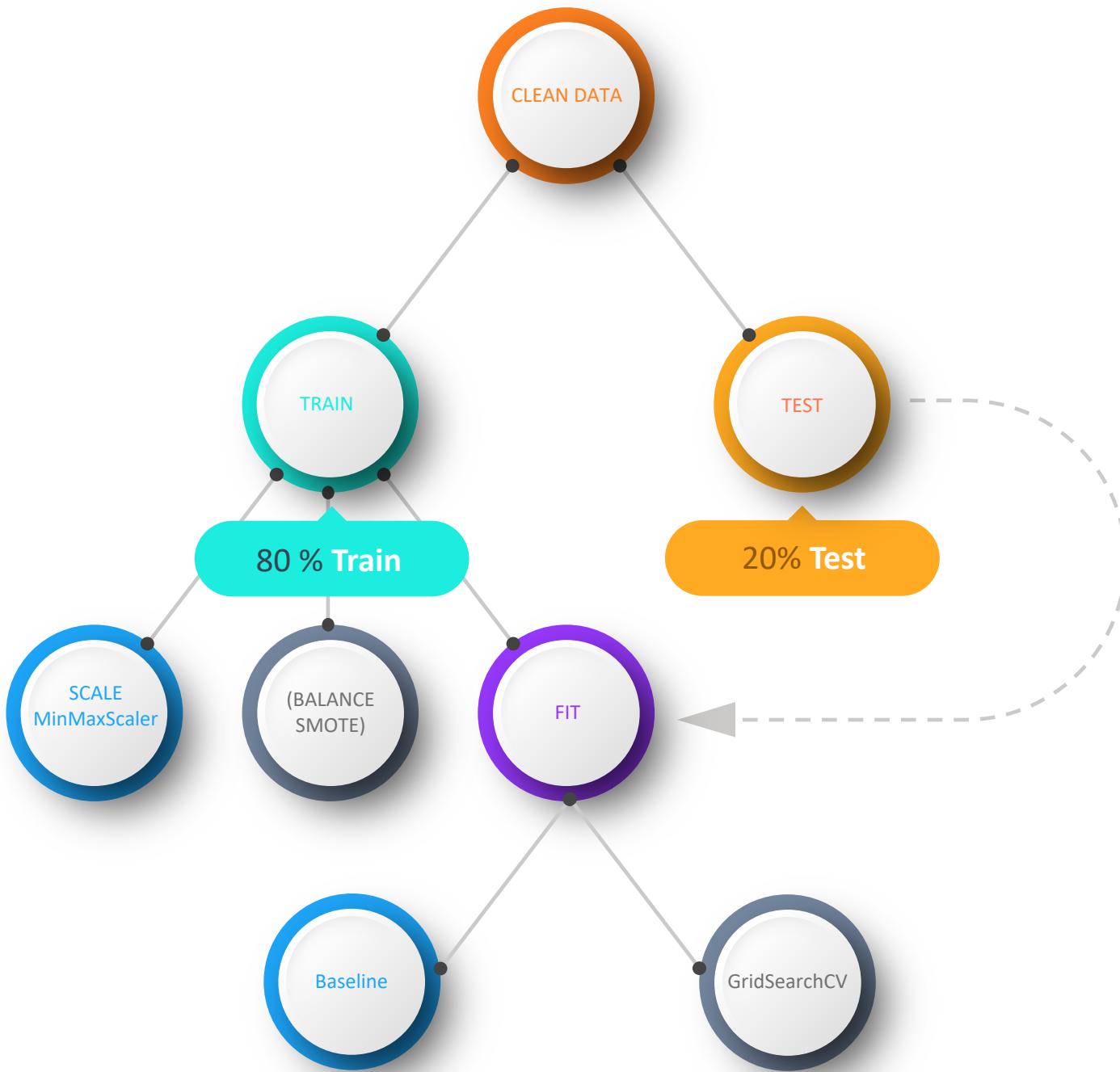


Recommendations

1. **Focus on sustainability:** early preventative strategy rather than letting things go broken
2. **Decentralized management:** we need to restructure authority so that there is a system of co-responsibility between the central, regional and local levels.
3. **Improved payment system:**
 - A local payment system should be put in place so that the user-group can be independently responsible for their own water points
 - Direct funding from international donors to village-level should also be implemented instead of having to go through the long bureaucratic process where money get lost along the way between ministry and district level.



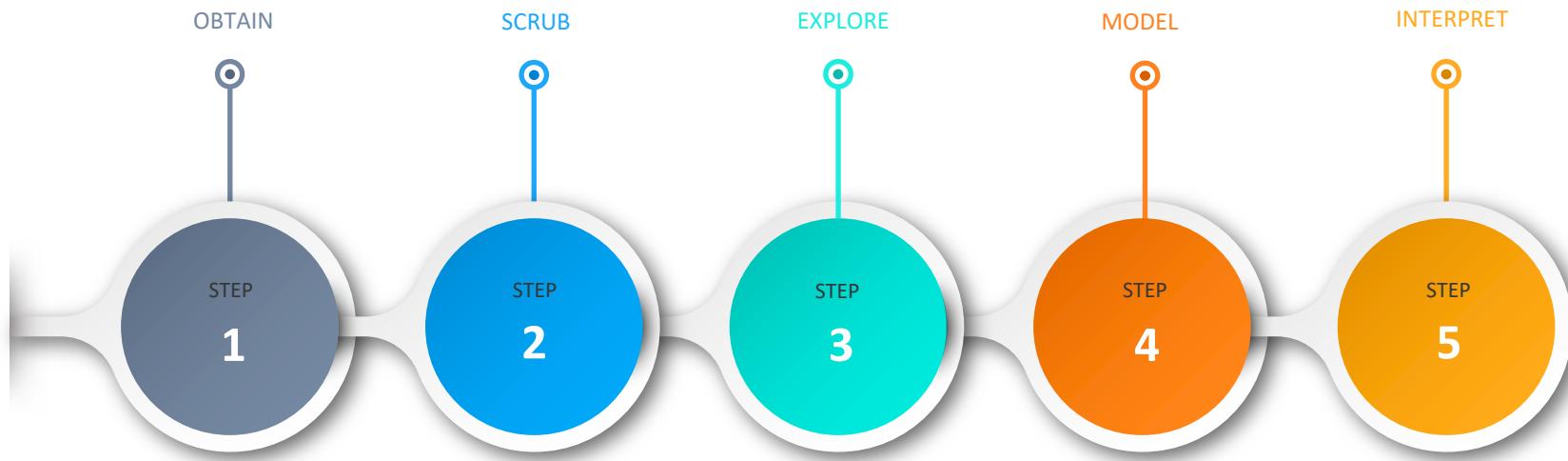
1. First, we convert the labels from strings to integers:
 - Non-functional = 0
 - Functional = 1
 - Functional-needs-repair = 2
2. Dummies encode all categorical variables



Best Model

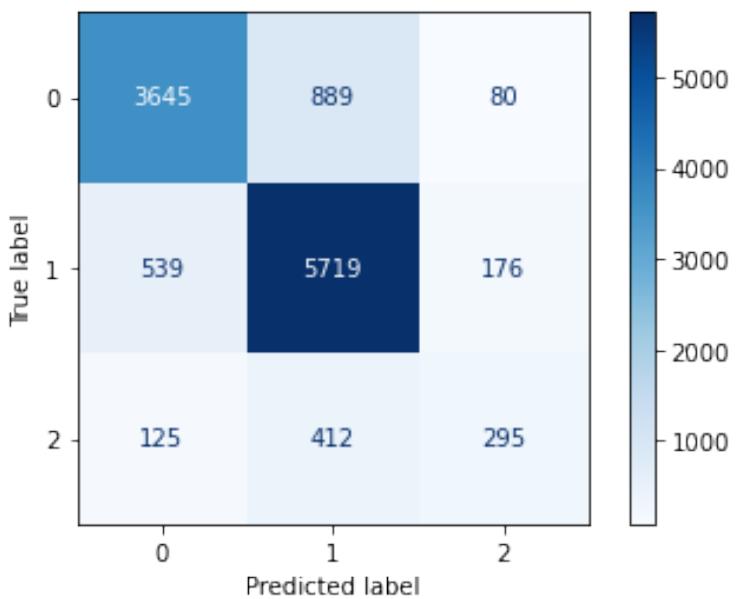
Gradient Boosting Classifier is a group of machine learning algorithms that combine many weak learning models together to create a strong predictive model.

The weak learner used is **Decision Tree Classifier**, which is the simplest tree-based method.

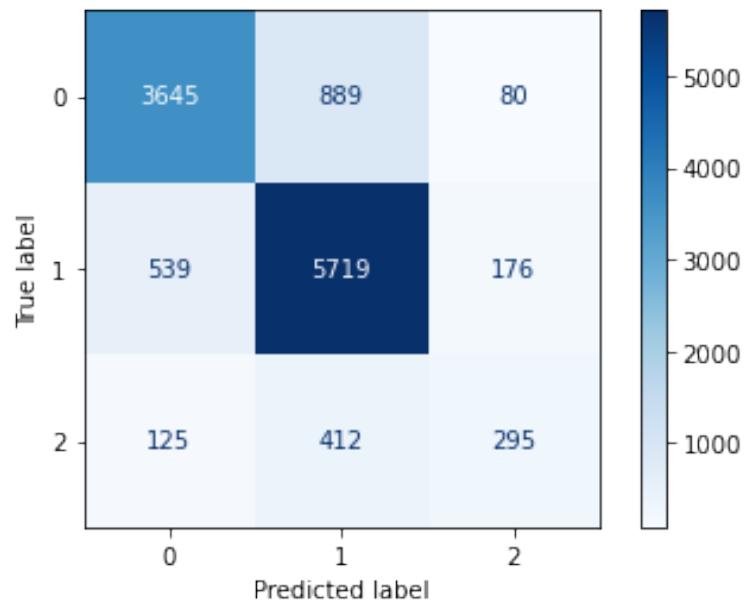


Accuracy 81.30%

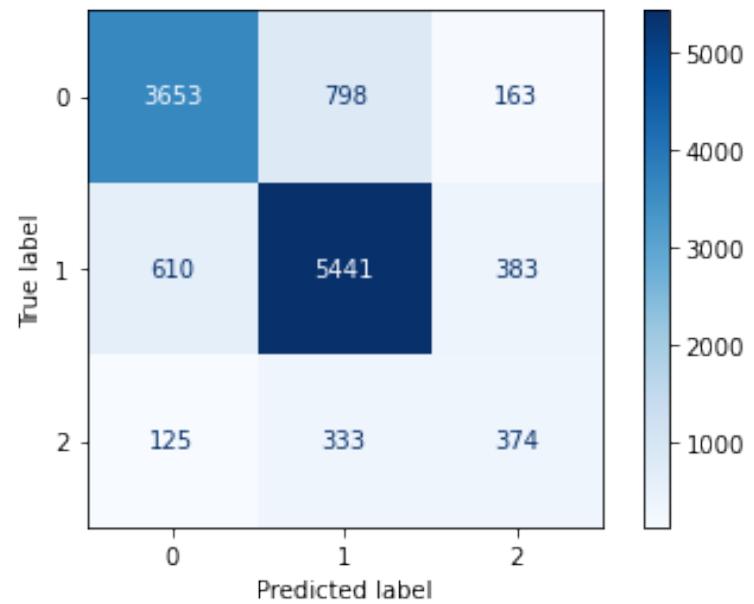
- Train accuracy: 99.33%
- Test accuracy: 81.30%



		Model: Gradient Boosting			
		precision	recall	f1-score	support
	0	0.85	0.79	0.82	4614
	1	0.81	0.89	0.85	6434
	2	0.54	0.35	0.43	832
		accuracy		0.81	11880
		macro avg	0.73	0.68	11880
		weighted avg	0.81	0.81	11880



Imbalanced Gradient Boost Model
81.30%



Balanced Gradient Boost Model
79.86%

Future Work

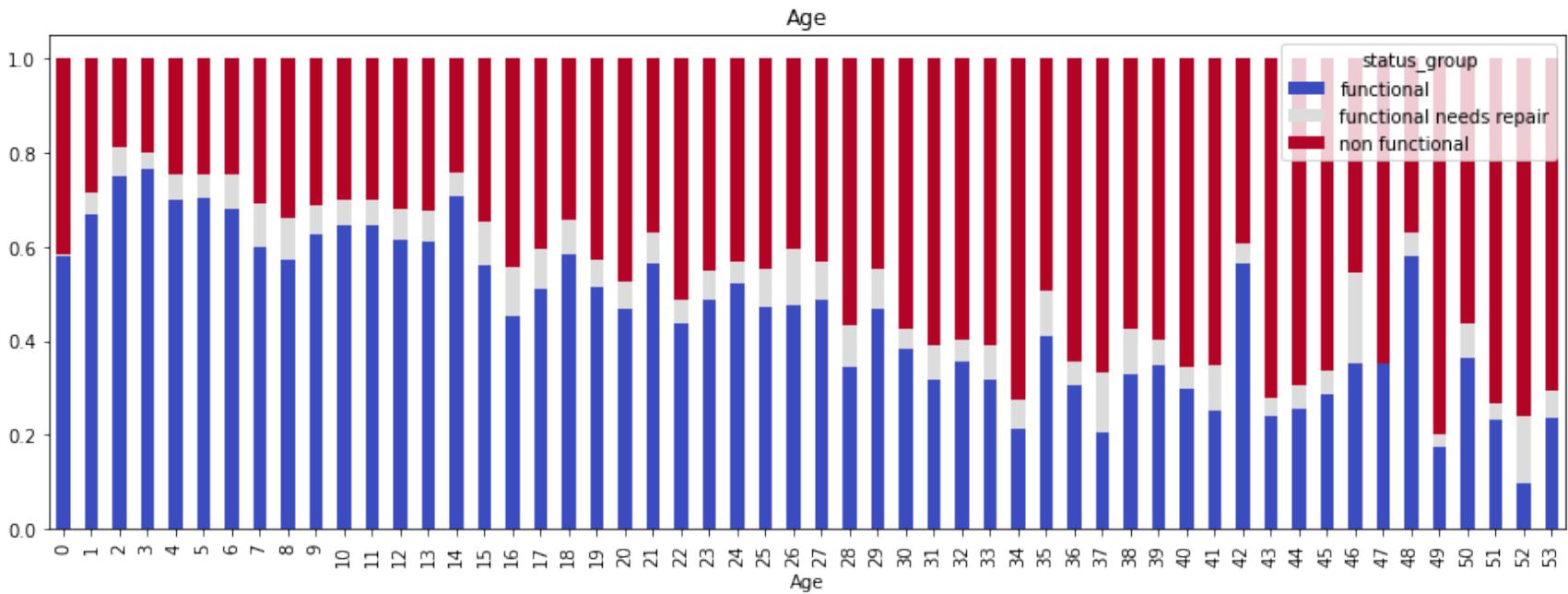
1. Since correcting class imbalance did not improve the model, we can try model stacking i.e build a binary classification between functional vs non-functional and another binary classification between functional vs. functional needs repair.
2. Try more parameters tuning with more and wider range of options
3. Work to reduce overfit while maintaining and/or improving accuracy score

THANK YOU

APPENDIX

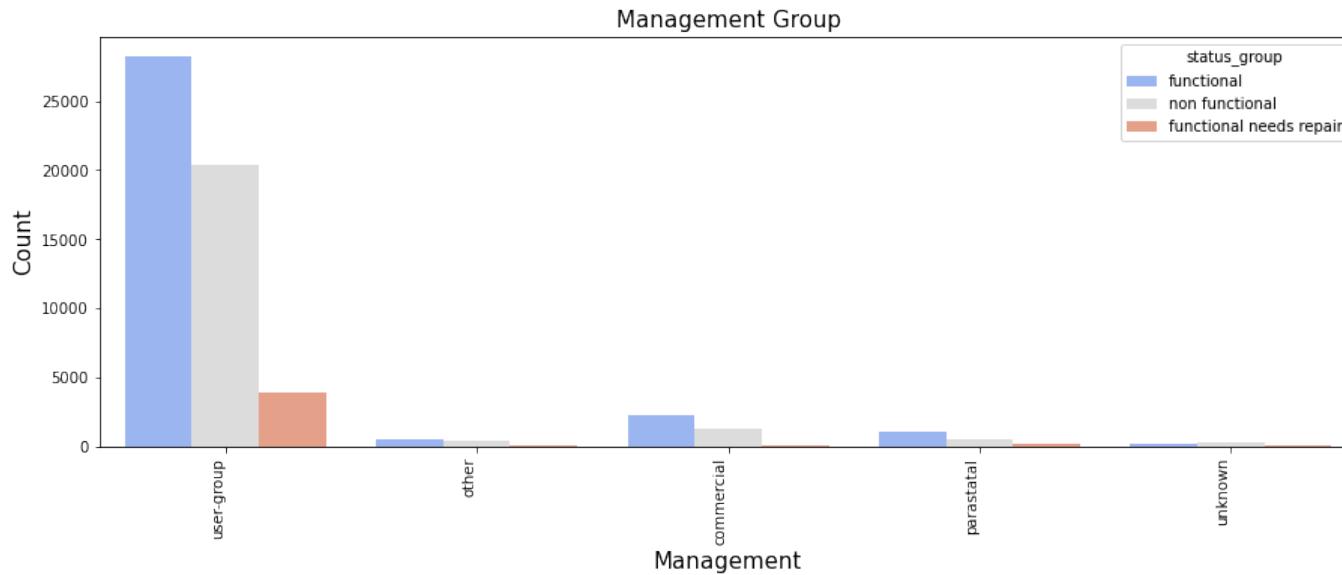
Problems

1. **Sustainability:** Regardless of hundreds of millions of dollars over budget and years past the original deadline of the Water Sector Development Program (WSDP), local government and communities find themselves unable to raise the money to fix and maintain their water points.
2. **Power struggle:** Full responsibility for operating, maintaining and sustaining water points is done at the village level. However, disbursement of funds and report of functionality must follow a long bureaucratic process all the way from the village, to the district, and, finally, to the Ministry of Water. The problem found is not only the miscommunication but also the power struggle around roles, responsibilities, and accountability between many different levels of government.



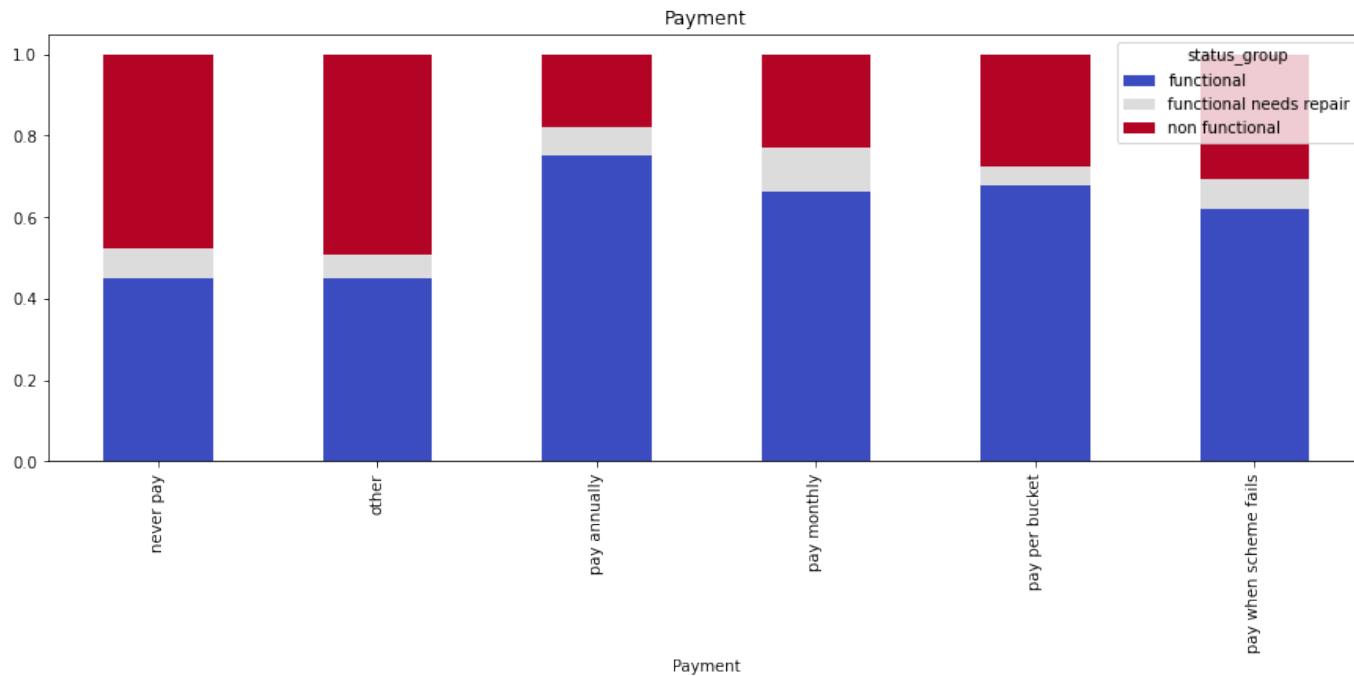
- Within the first year, 30% of water points become non-functional.
- Only 54% of water points are working 15 years after installation.

Problem: Regardless of hundreds of millions of dollars over budget and years past the original deadline of the Water Sector Development Programme (WSDP), local government and communities find themselves unable to raise the money to fix and maintain their water points.



- Under NAWAPO (National Water Policy), user groups are to take the full responsibility for operating, maintaining and sustaining water points at the village level.
- Data is published by the ministry which are based on the coverage reported by district, are not reliable.
- Disbursement of funds must follow a long bureaucratic process of accountability, requiring upwards (vertical) reporting at each level of government, all the way from the village, to the district, and, finally, to the Ministry of Water.

Problem: Power struggle around roles, responsibilities, and accountability between different levels of government



- If the water point management charges money, the more likely that it is better maintained and kept functional.
- Regular payment, used for regular maintenance and upkeep, is a better approach for preventative treatments rather than trying to secure a large amount of fund for when the system breaks down.

Balanced

	Model	Accuracy	CV	Precision	Recall	F1 Score	MAE	MSE	RMSE	AUC	Bias	Variance
0	Decision Tree	76.01	0.72	0.65	0.67	0.66	0.27	0.33	0.57	-	0.02	0.39
1	Logistic Regression	65.31	0.74	0.60	0.67	0.58	0.42	0.55	0.74	0.82	0.24	0.55
2	KNN	75.66	0.72	0.64	0.68	0.66	0.28	0.34	0.59	-	0.05	0.40
3	Bagged Tree	78.21	0.76	0.67	0.69	0.68	0.24	0.30	0.54	-	0.01	0.37
4	Random Forest	79.37	0.77	0.69	0.70	0.69	0.23	0.28	0.53	-	0.02	0.37
5	Gradient Boost	79.86	0.78	0.69	0.70	0.70	0.22	0.27	0.52	-	0.03	0.37
6	ADABOost	63.01	0.72	0.56	0.61	0.55	0.44	0.59	0.77	-	0.21	0.53
7	XGBoost	77.03	0.78	0.67	0.71	0.68	0.26	0.32	0.57	-	0.09	0.41
8	SVM	72.50	0.76	0.64	0.71	0.64	0.32	0.42	0.65	-	0.17	0.49

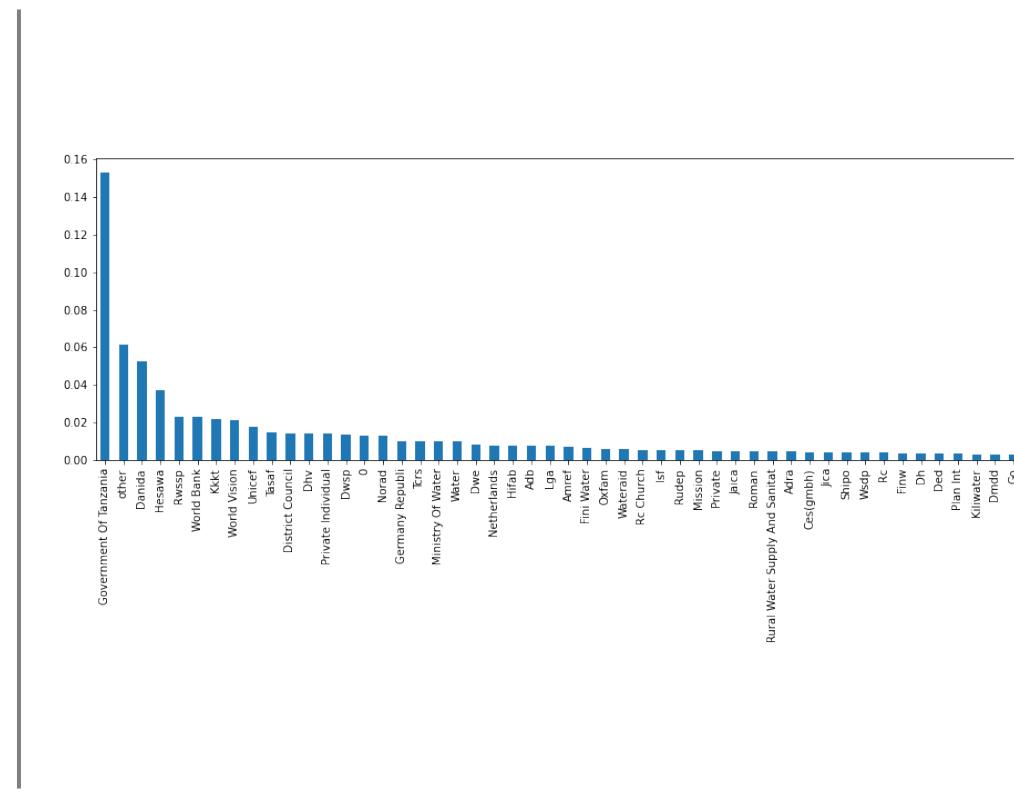
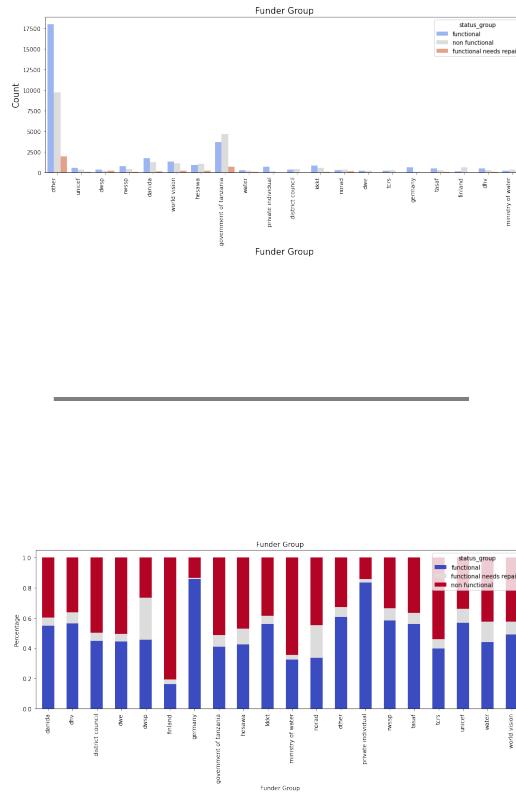
Imbalanced

	Model	Accuracy	CV	Precision	Recall	F1 Score	MAE	MSE	RMSE	AUC	Bias	Variance
0	Imbalance Decision Tree	75.96	0.72	0.65	0.66	0.65	0.27	0.32	0.57	-	0.01	0.36
1	Imbalance Logistic Regression	74.87	0.74	0.68	0.56	0.58	0.27	0.30	0.55	0.81	0.01	0.25
2	Imbalance KNN	78.96	0.73	0.70	0.66	0.68	0.23	0.27	0.52	-	0.01	0.32
3	Imbalance Bagged Tree	79.16	0.76	0.69	0.66	0.67	0.23	0.27	0.52	-	-0.01	0.32
4	Imbalance Random Forest	80.37	0.77	0.71	0.67	0.69	0.22	0.26	0.51	-	0.00	0.32
5	Imbalance Gradient Boost	81.30	0.78	0.73	0.68	0.70	0.20	0.24	0.49	-	0.00	0.31
6	Imbalance ADABoost	71.80	0.72	0.60	0.51	0.50	0.30	0.32	0.57	-	0.02	0.22
7	Imbalance XGBoost	80.34	0.78	0.76	0.63	0.66	0.21	0.24	0.49	-	0.00	0.27
8	Imbalance SVM	78.54	0.76	0.75	0.60	0.63	0.23	0.26	0.51	-	0.02	0.25

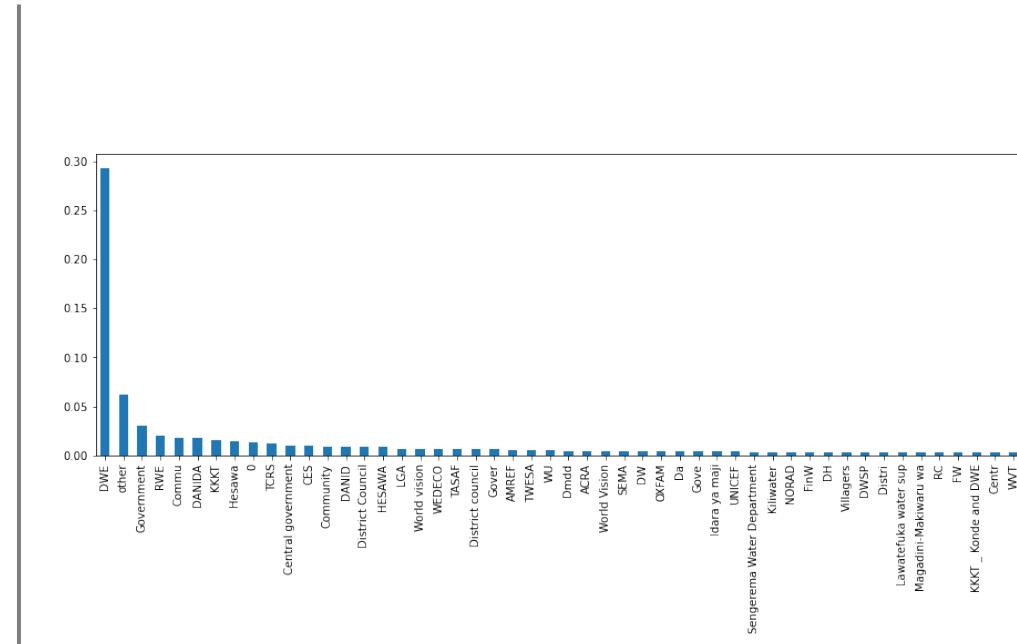
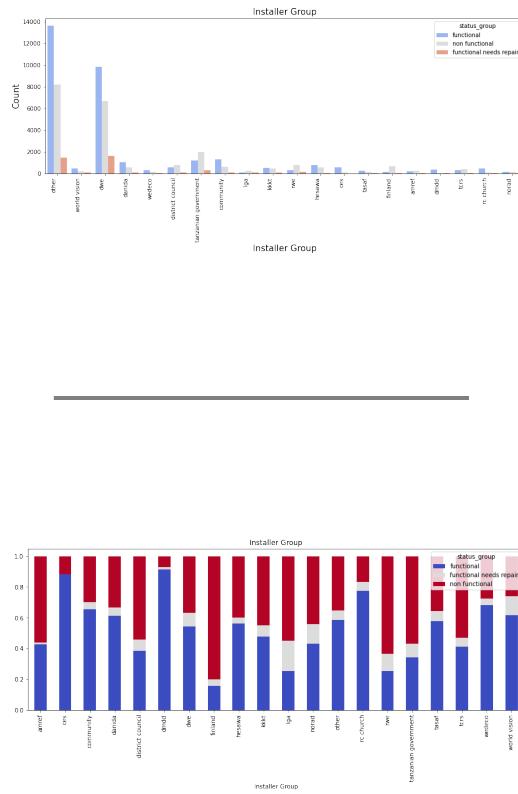
Metrics

- **Accuracy:** Out of all the classes, how much we predicted correctly. It should be high as possible.
- **Precision:** What proportion of positive identifications was actually correct? i.e. $TP / (TP + FP)$. It should be high as possible.
- **Recall:** What proportion of positive identifications was actually correct? i.e. $TP / (TP + FN)$. It should be high as possible.
- **F1:** It is difficult to compare two models with low precision and high recall or vice versa. It is often convenient to combine precision and recall into a single metric called the F1 score. We want to maximize f1 score as well.
- **Cross validation:** this should be as close as possible with the model's RMSE or else we overfit our model.
- **Macro average** is the average of precision/recall/f1-score.
- **Weighted average** is just the weighted average of precision/recall/f1-score.

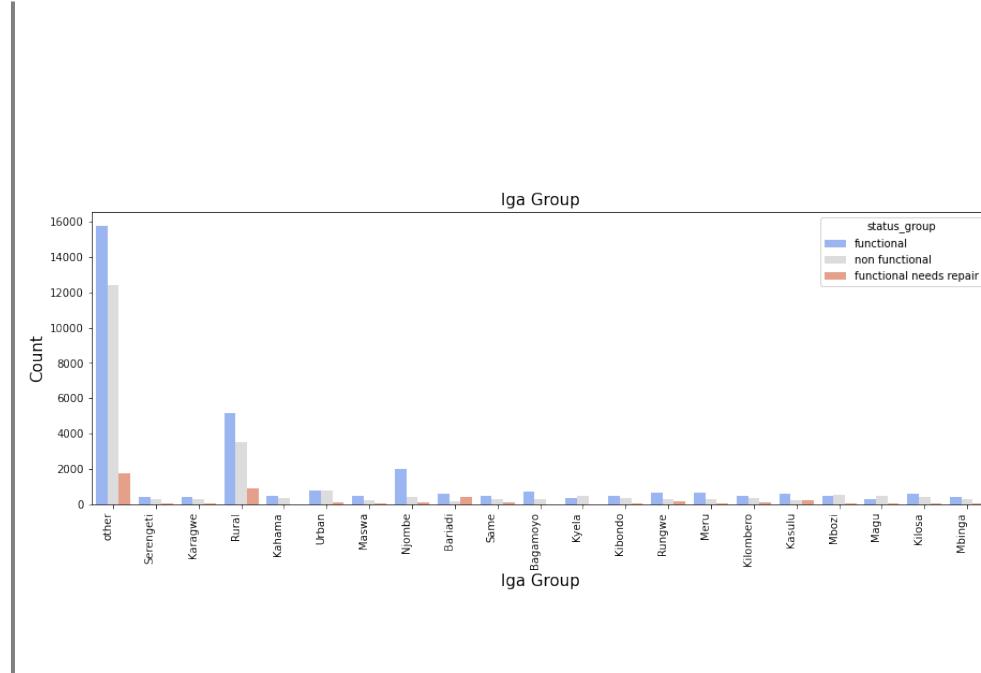
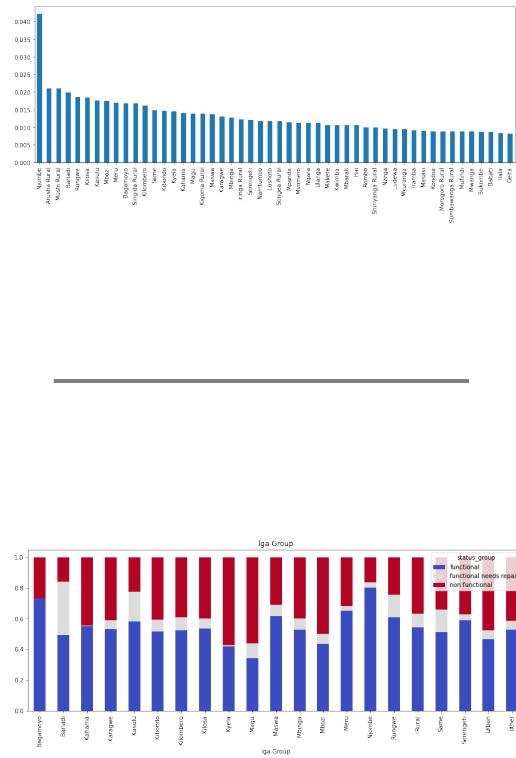
funder



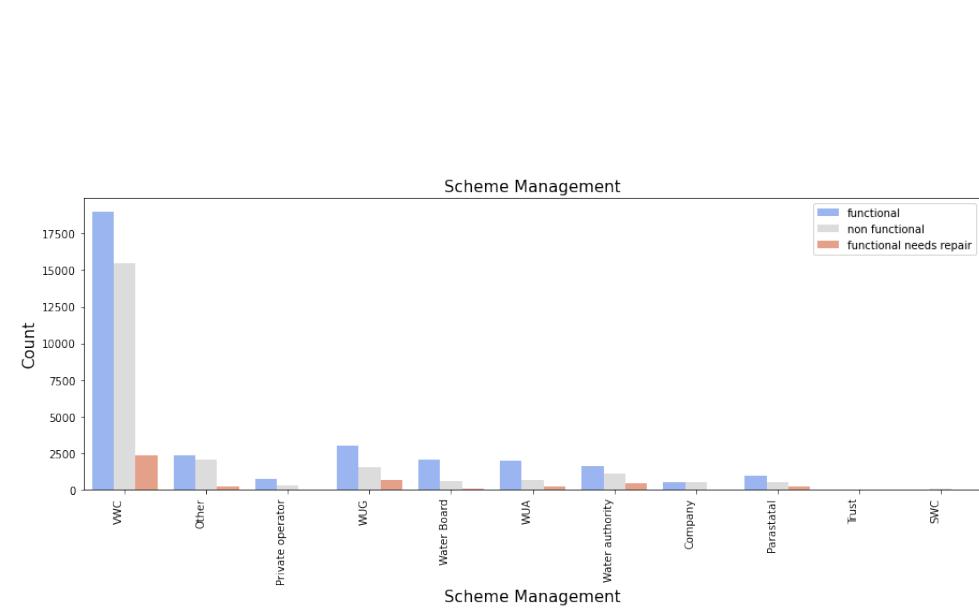
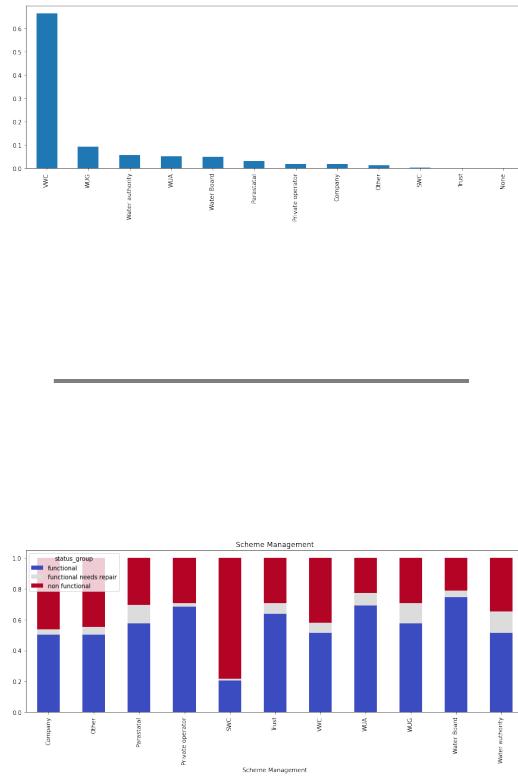
installer



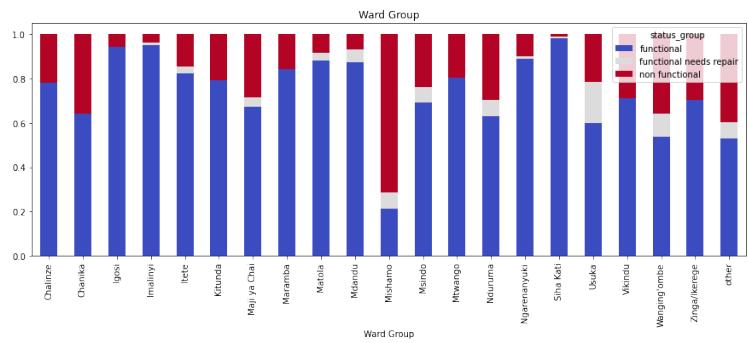
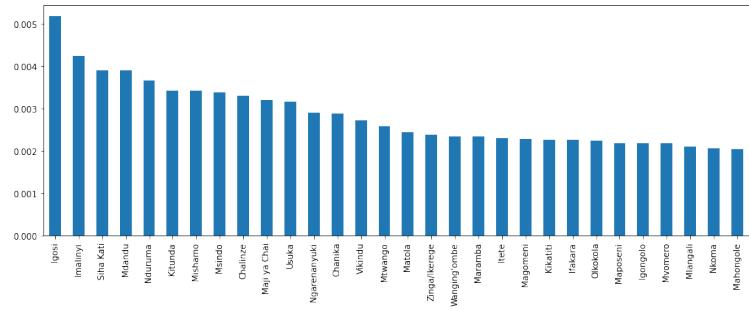
Iga



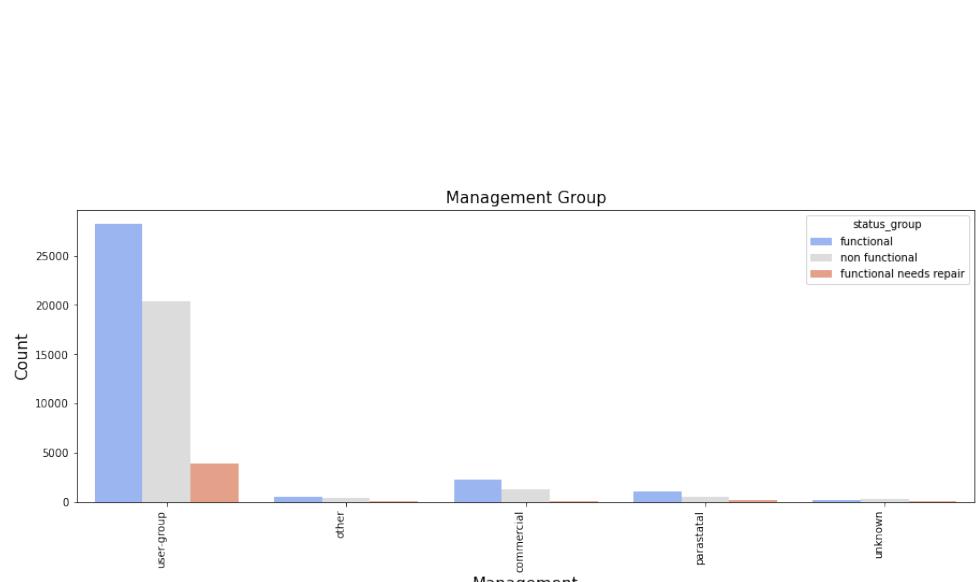
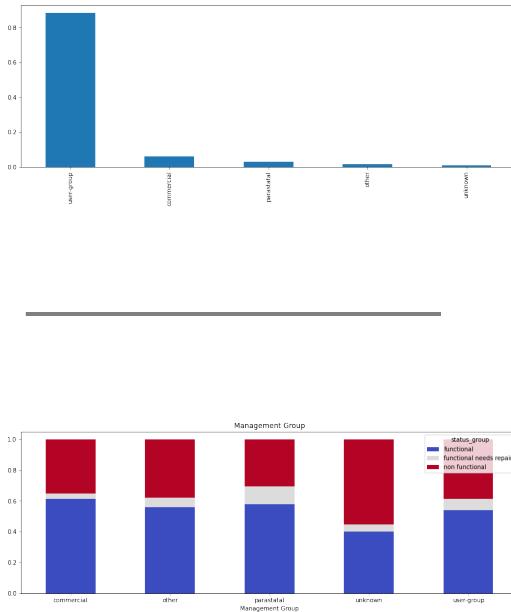
scheme management



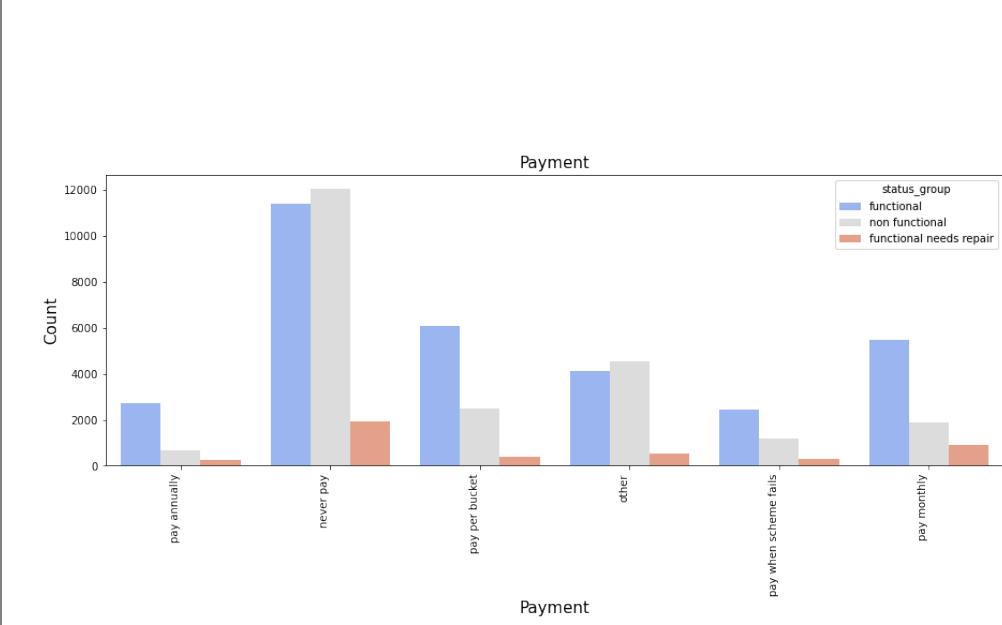
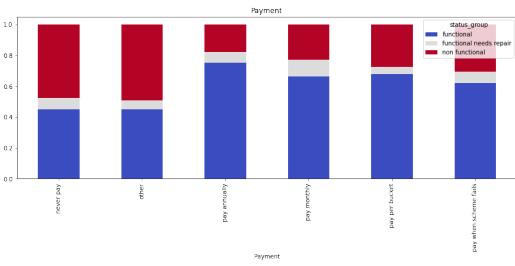
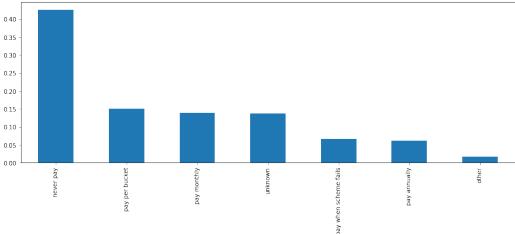
ward



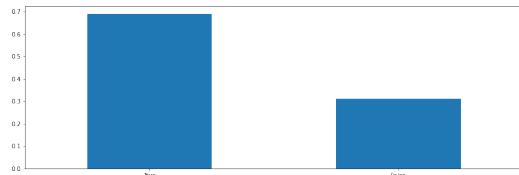
management group

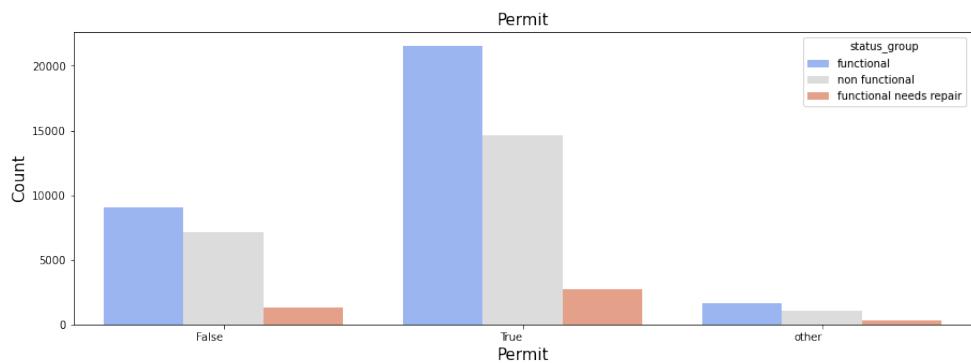
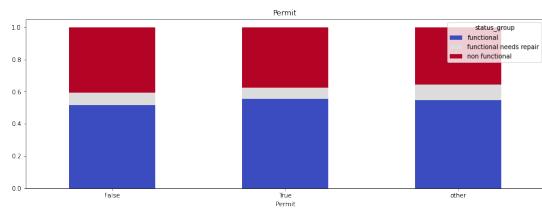


payment

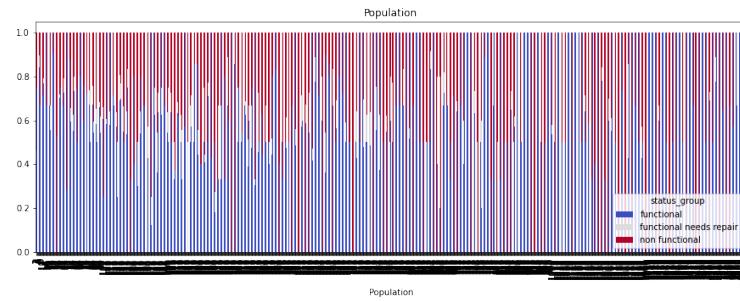
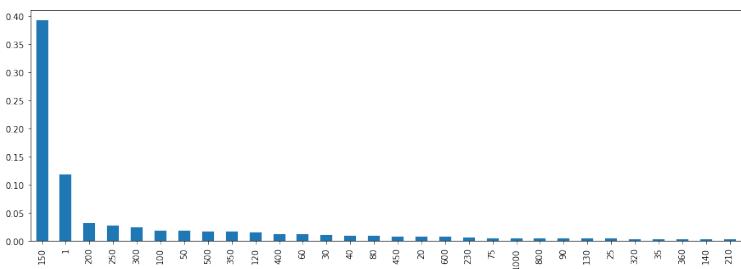


permit

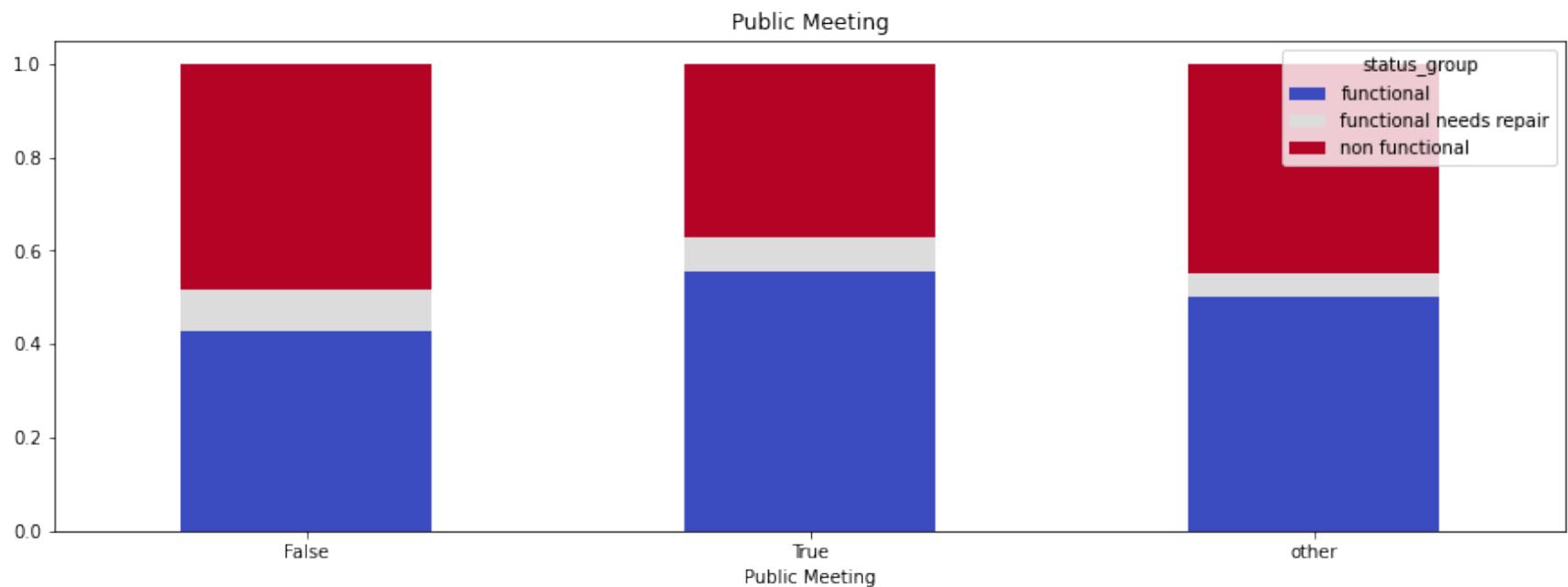




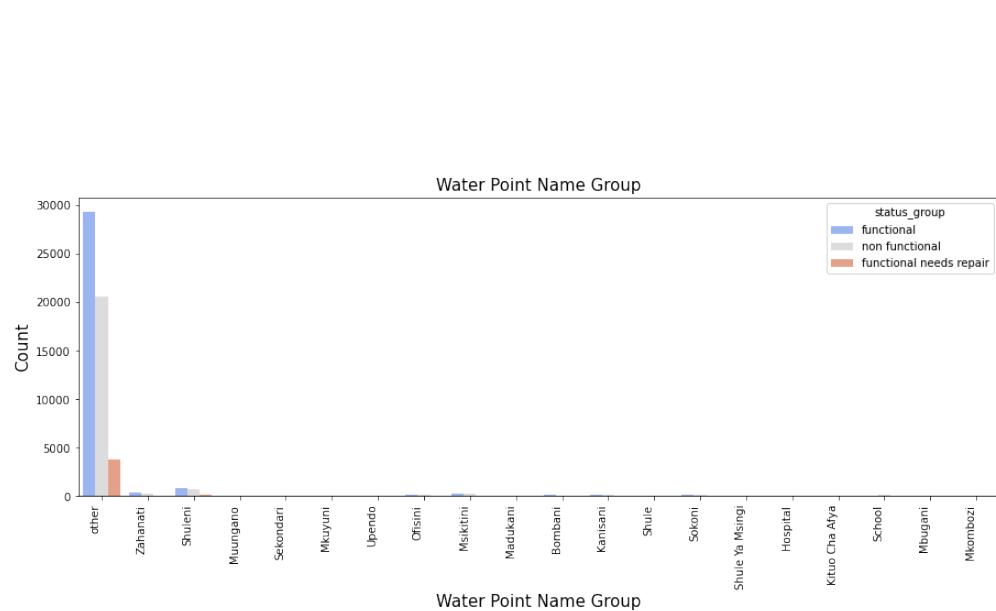
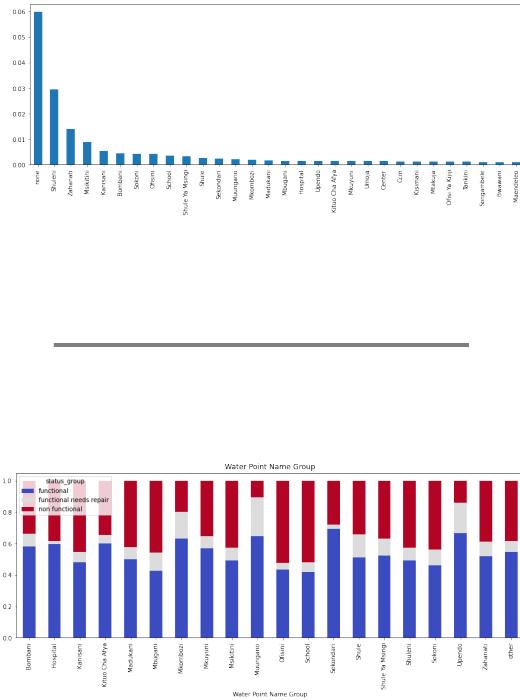
population



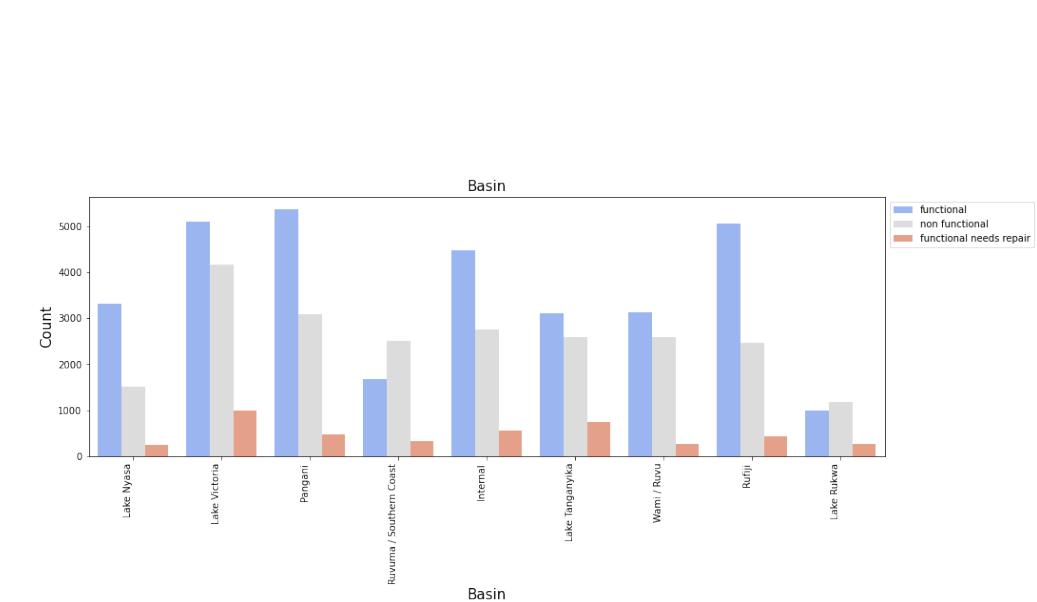
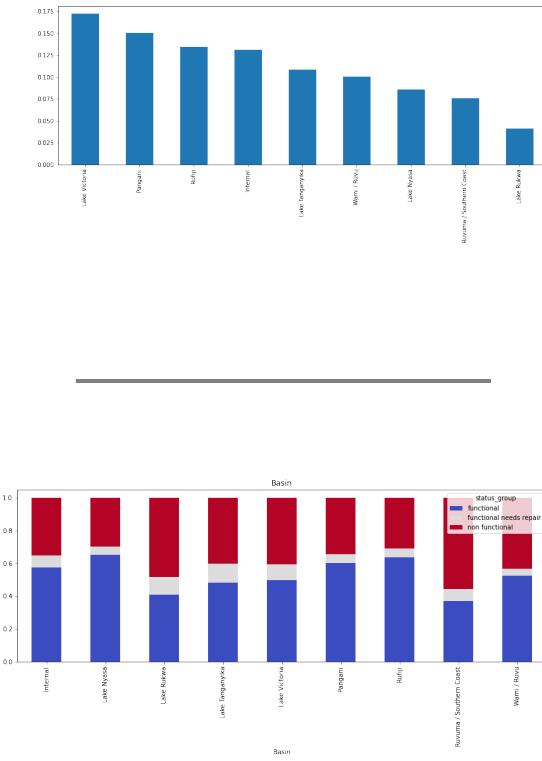
public meeting



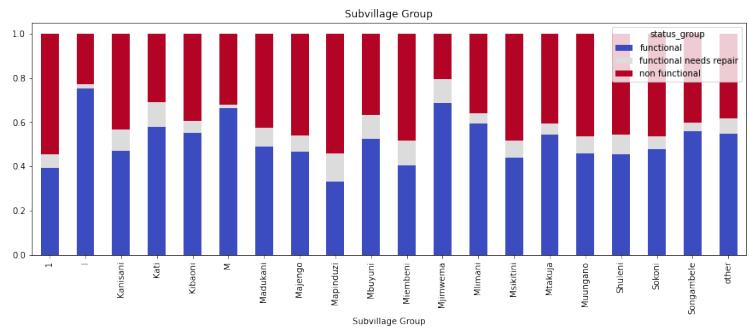
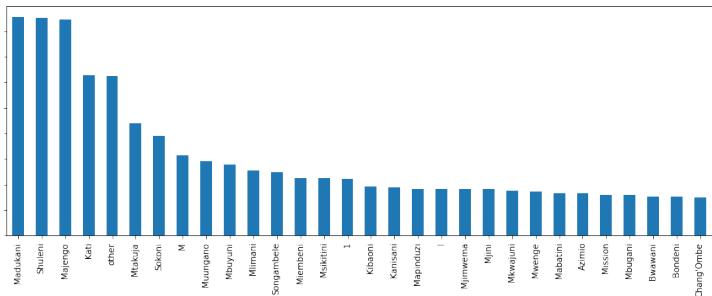
water point names



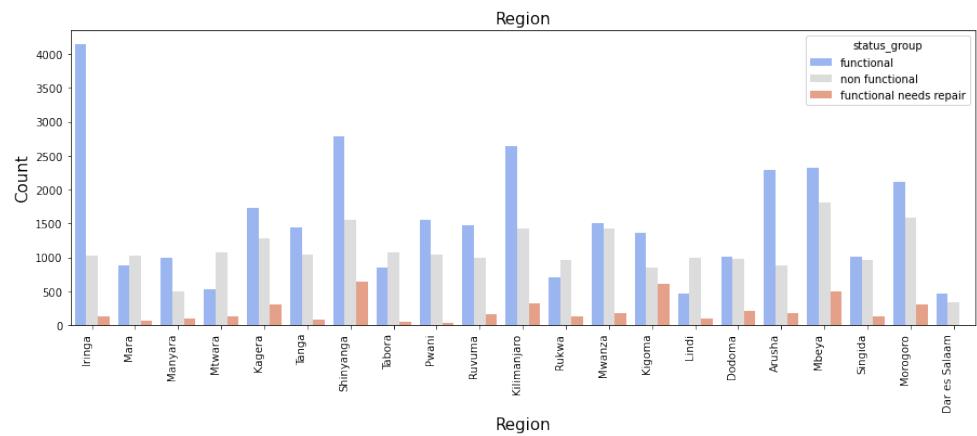
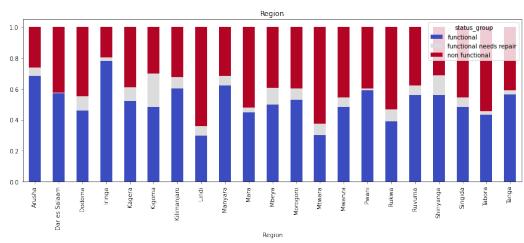
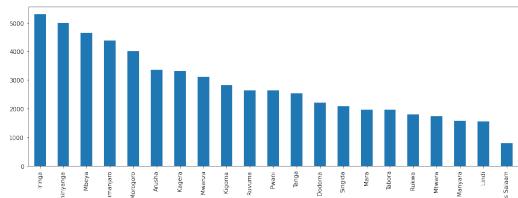
basin



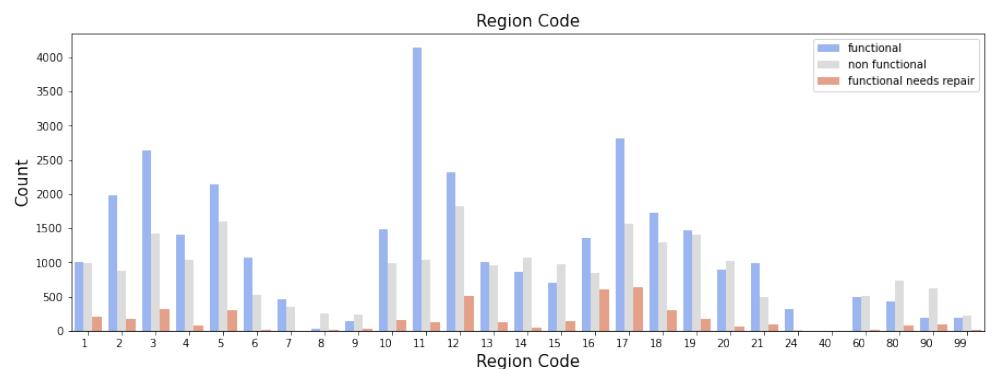
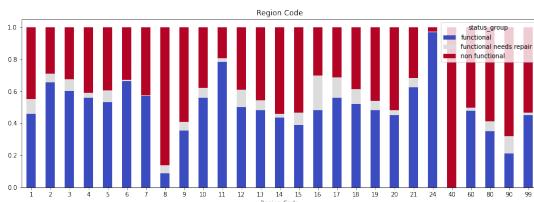
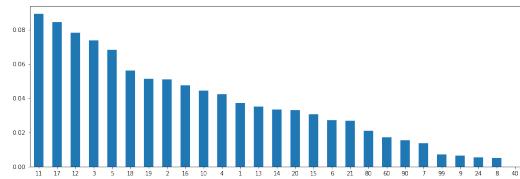
subvillage



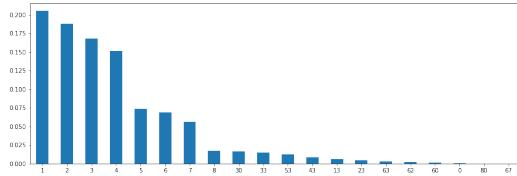
region



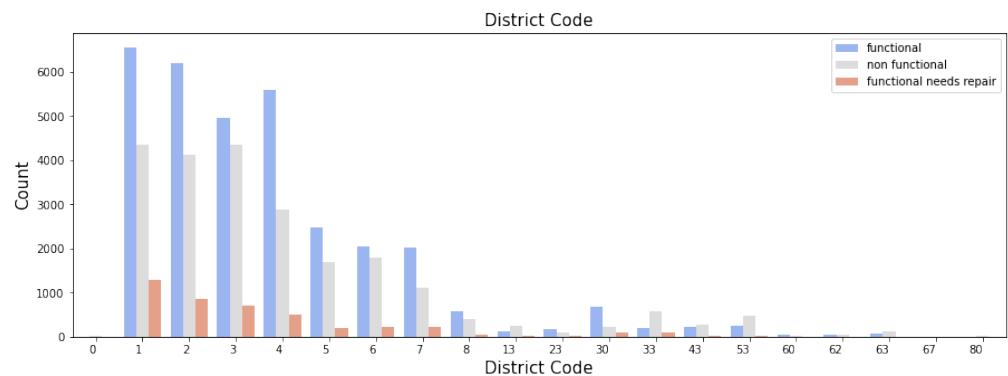
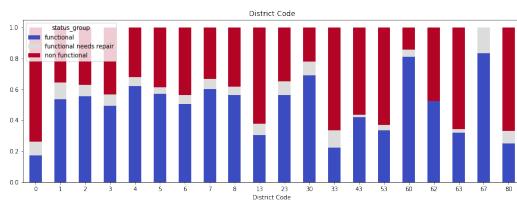
region code



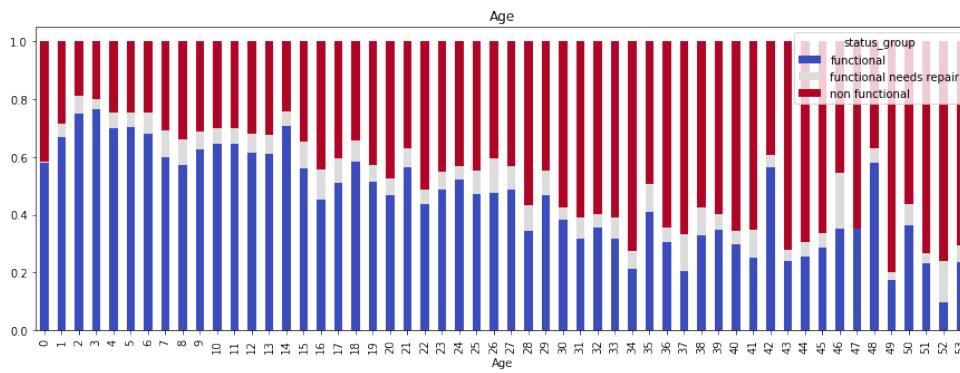
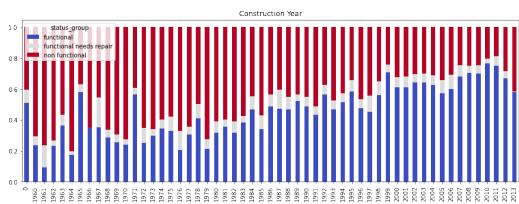
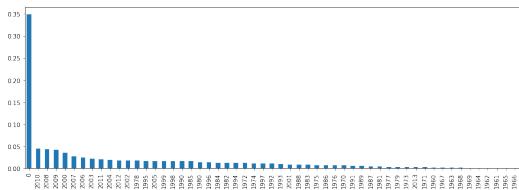
district code



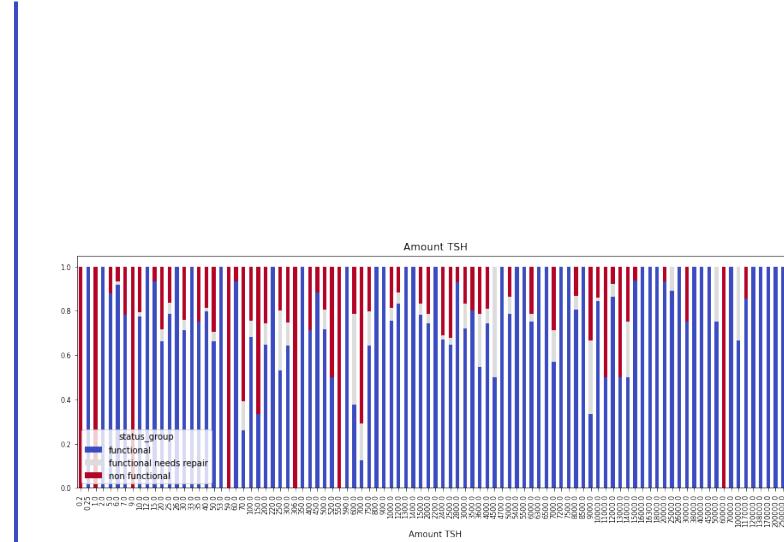
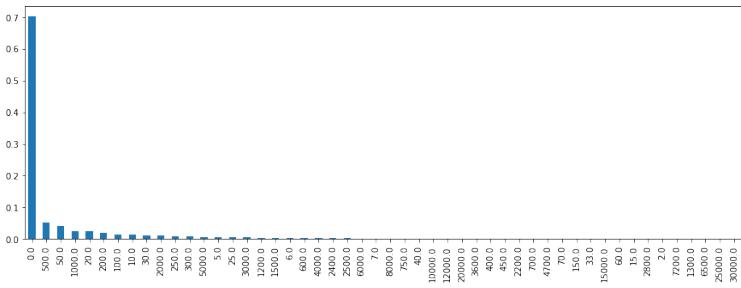
—



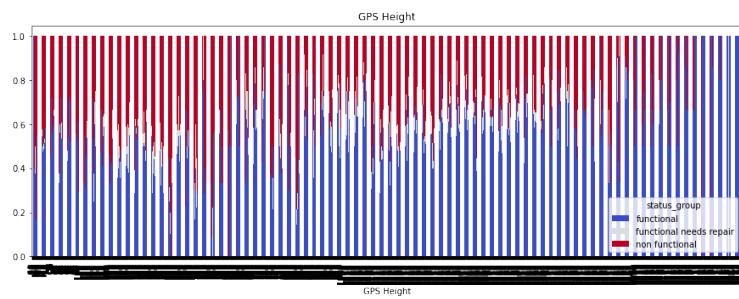
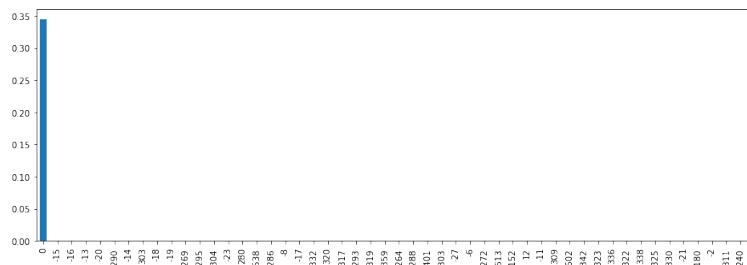
age



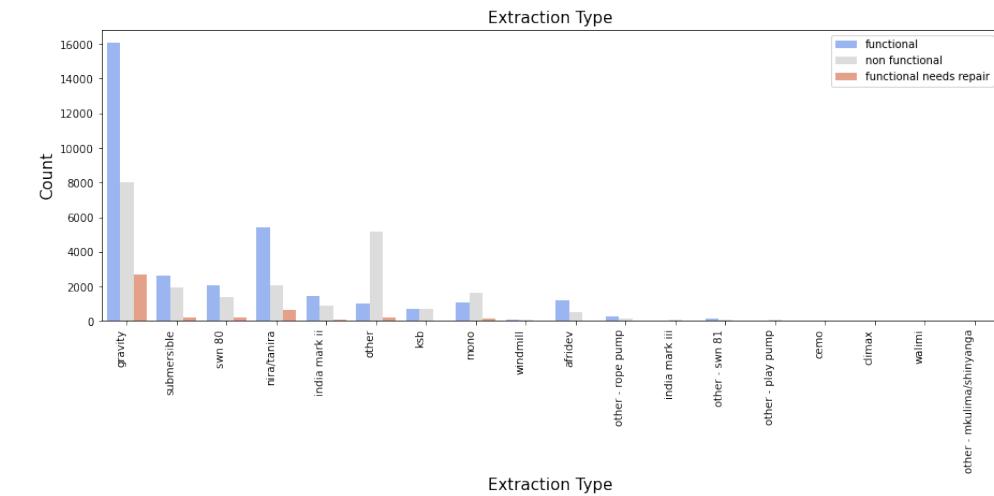
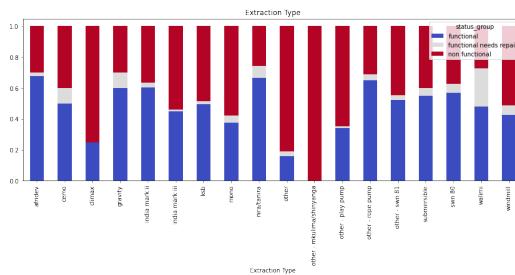
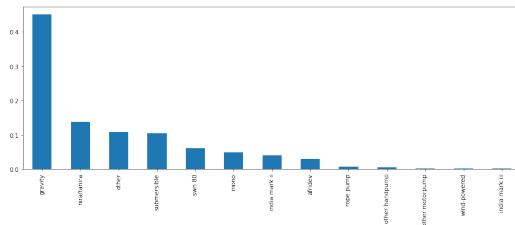
amount tsh



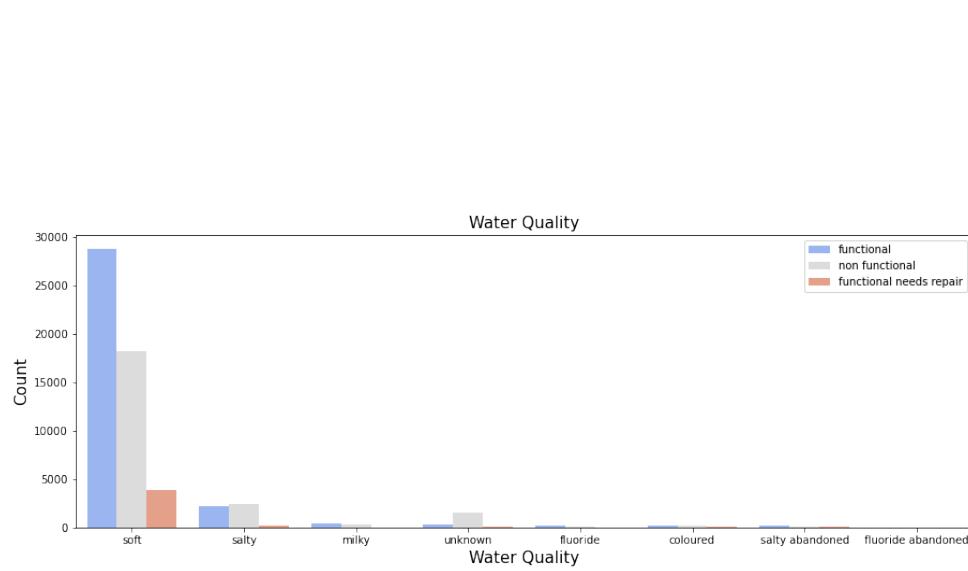
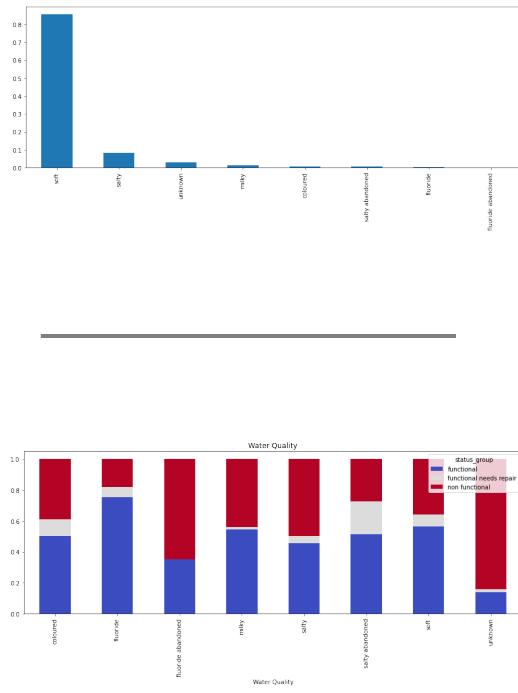
gps height



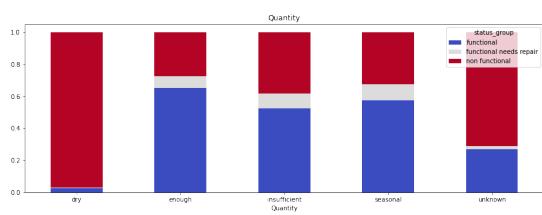
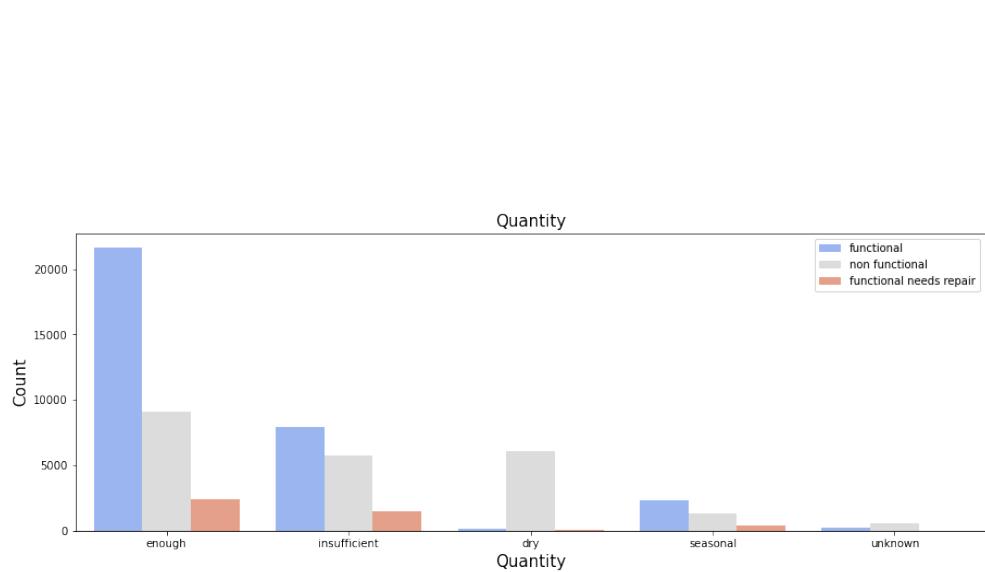
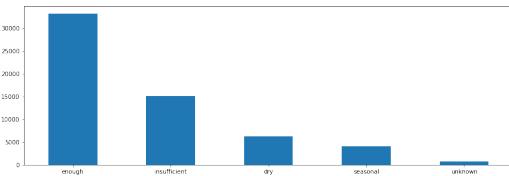
extraction type



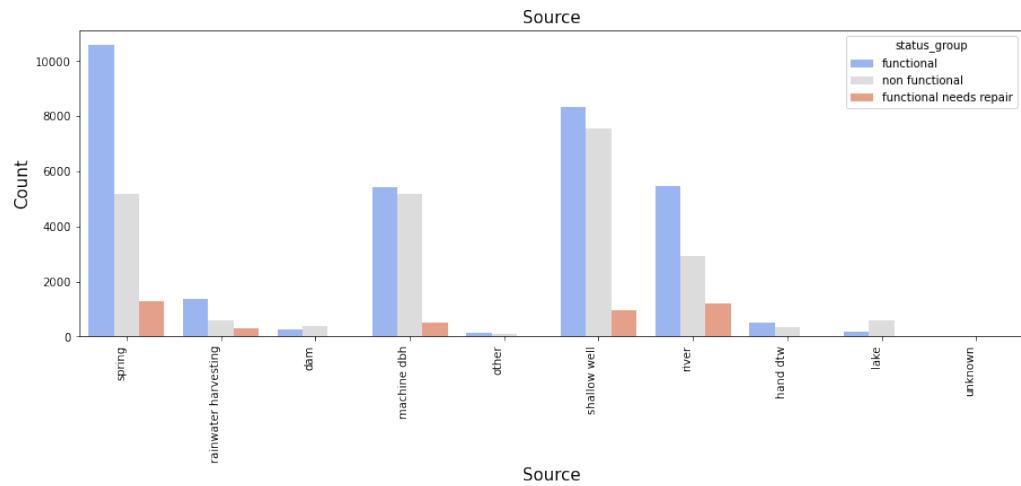
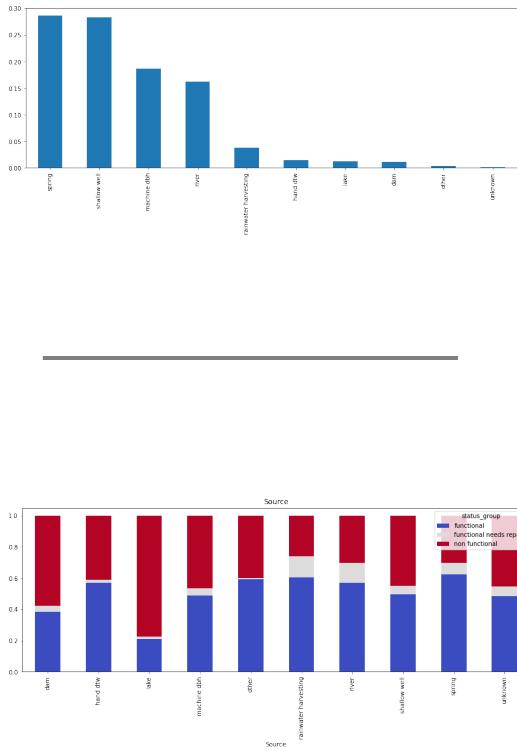
water quality



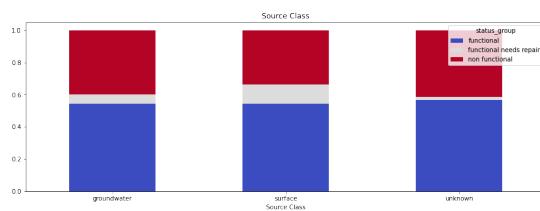
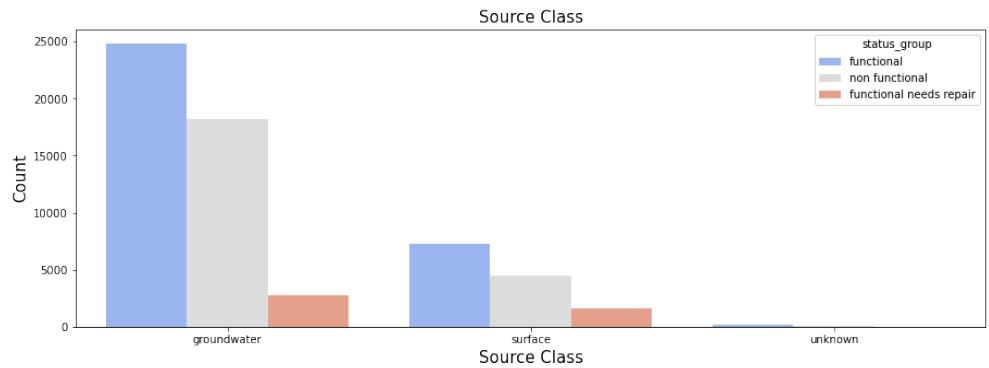
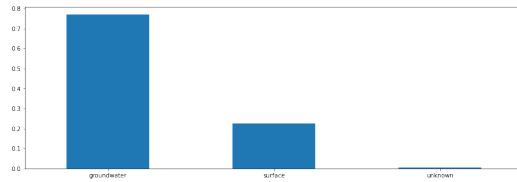
quantity



source



source class



water point type

