



CAPSTONE PROJECT: PUMP IT UP

Data Analysis with R

Tadaaki Sun
Foundations of Data Science

Capstone Project Report

Introduction

Tanzania is one of the poor nations around the world who is currently suffering from issues that are related to the water situation. Approximately a third of their country is dry and rarely receive rain, where they have minimal access to water that is clean and drinkable. There are several water points that provide water to the community of Tanzania, however, they are usually far away from villages and hard to retrieve. To solve the many arising issues related to water, the Tanzanian Ministry of Water has decided to hire TDS (Tadaaki's Data Science) Consulting, to have assistance in predicting which water pumps in their country are either, non-functional, functional, or functional but needs repair.

Through several meetings with the Tanzanian Ministry of Water, TDS has observed that the Ministry of Water are interested in understanding the core reason behind the fault of water points in Tanzania. To specifically find the reason for default, any waterpoints that are functional but needs repair will be counted as non-functional for simplicity.

As the Tanzanian Ministry of Water and Taarifa thrive to provide water throughout the Tanzanian community, we would like to propose a recommendation based on our analysis using data analysis tools. The recommendation will then be used to tackle the root issue of the water issue and plan to prevent in the future by providing clean and accessible water points for the Tanzanian community.

Objective

The summarize the objective of this resport, we have identified the general and specific objectives being covered:

- **General Objective**

Use statistical methods to predict the root cause of non-functional water points, which will further help in improving the maintenance operations in water points in Tanzania to ensure that clean, and portable water is available to the Tanzanian communities.

- **Specific Objectives:**

- Apply exploratory data analysis and visualization techniques to gain insight into the relationships between the different variables in the dataset.
- Determine which variables are the most associated with the different status of the water points.
- Create a predictive model that can be used to predict the status of water points (functional, non-functional, and functional but need repair).
- Create a list of actionable recommendations from the analysis of the data and the results of the model.
- Provide deliverables (report, presentation slides, and code for the model) to assist

The code, report, and presentation slide which include the detailed analysis and predictive model used to identify the recommended solutions for the Tanzanian Ministry of Water. We hope the documents can be a guide to predict more issues and bring a great community in Tanzania.

Taarifa Dataset

The Tanzanian Ministry of Water has partnered with Taarifa, who currently is working on an Innovation Project in Tanzania, who provided the data such as the status of the pump, when they were installed, how they are managed, etc. They have gathered them from hand held devices, and papers to be gathered in an excel sheet.

The data consists about 40 features and 59,400 observations as a training set, and a further 18,500 observations for the test set, that may assist us in creating a predictive model for finding the non-functional waterpoints in Tanzania. The model will be used in a test dataset provided to measure the effectiveness (the full details of features are in the Appendix 1).

In addition to the features above, each water point id will have information on which status group (functional, functional but needs repair, or non-functional) it belongs to. This data will be the dependent variable which will be tested with the other features to predict the reason the pumps are not working.

For our own purposes, we will treat any waterpoint status of “functional but needs repair” as “non-functional” status for simplicity in identifying the reason for default. Furthermore, we will use only the training set to create a predictive model to check the efficacy of the model for future predictions.

The rest of the variables are the independent variables used to investigate what the root cause of defaults of the waterpoints. Most of the data are categorical variables, where there are nominal and ordinal, and continuous variables. Some of the nominal categorical variables include the installer/funder of the waterpoint, information of region/village/etc., and source types. Ordinal categorical variables include the availability of public meetings, permits, etc. Lastly, there are continuous variables such as the amount of water heads, population, and construction years. With a great amount of data for each feature, a thorough preliminary data analysis was done to better understand the data.

Base Line Performance

To identify a feature which is the root of the cause for the faulty water points, the final predictive model will be classified as either functional or non-functional. The amount of water points that were functional but needs repair was only a fractional amount in the whole training dataset.

Once the status is changed, we have a result where there were 32229 functional and 27171 non-functional waterpoints. Therefore, the current data will have approximately 54% (32229/59400) accuracy if we did not have a model. This is not a good predictor as there is only a half chance of guessing the outcome. Therefore, we will create a new model that will have a higher accuracy compared to 54% to better identify the non-functional waterpoint.

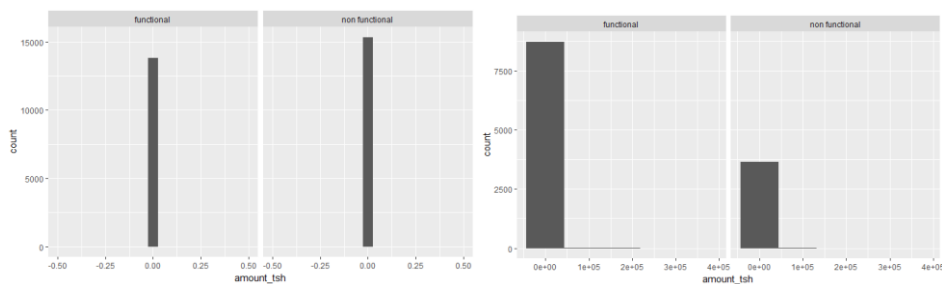
We will then split the data into a training and testing set by having 70% of randomly selected data in the training and 30% in the test set. The new datasets also calculated a result of 54% (22582/41581) accuracy for the baseline performance.

Preliminary Data Analysis

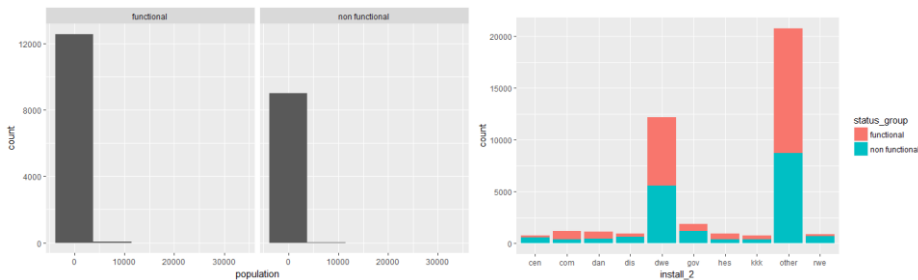
Once the baseline performance is set, the next step is to examine the independent variables. All features were carefully examined by utilizing various visualization techniques such as bar plots and histograms to make sense of variables such as quantity, quality group, water point type. Furthermore, any missing plots that may cause the graphs to skew can be identified at an early stage. We will adjust the data to have a better distribution to create variability in each feature. The utilization of various R packages can assist in creating better visual diagrams and apply some statistical classification methods to have a better prediction with the data. For example, the ggplot2 package can will be utilized to create bar charts and histograms that will give us insights into features that are currently non-functional.

With the first preliminary exploration, there were a lot of missing data as the information were coming from various sources such as handheld devices to paper. There are various levels of features which provided similar results and there is a need to categorize them as duplicates. Most of the operation data had thousands of high levels and are hard to predict the main cause. The training data had informed us that approximately half were functional and non-functional and a fraction being functional but needed repair. Finally, there were few old water point information and a lot of recently implemented pumps. These findings are all limitations when creating the predictive model and would need to be carefully assessed.

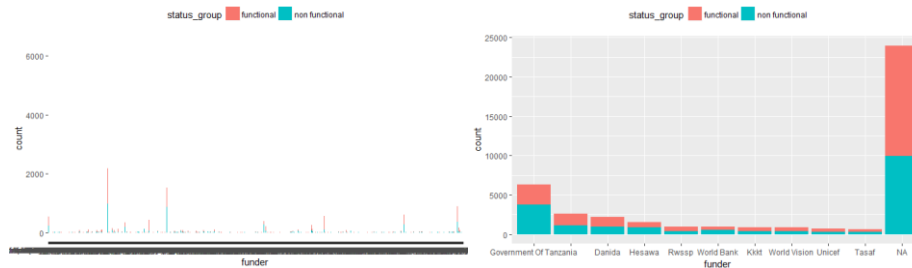
By quickly observing the data in the preliminary data analysis, there were many variables where the information was missing or logically questionable about the results. The dataset provided by Taarifa does a good job in tracking most of the variables with minimal NA's, however, there were many 0's and unknowns that still hinder us from providing insights. For example, there is one feature, "amount_tsh", which is the amount of water around the waterpoint. There are nearly 30,000 waterpoints, yet half were functional or non-functional (left hand diagram below where water amount equals 0). In a logical perspective, a waterpoint without water would mean that it can not provide any water to the community around, and therefore should be characterized as non-functional.



A further analysis was done to check the status for waterpoints that at least have more than 0 water (right hand diagram on the previous page). By creating a histogram for waterpoints that had water, two-thirds of them were functional and one-third was non-functional. There is a good case where the amount of water could be an amazing identifier for the faults, however, there were too many waterpoints that did not have water, but functional. Therefore, we decided to remove this and any other feature (for example population, indicated in the graph below) that had a lot of poorly recorded data to find a feature with correct data.



Furthermore, features where majority of the results were “other” or “NA” were removed, as we are not sure the exact reason why the water pump is at fault. The installer feature had one majority data, despite the “other’s”, however, the rest of the data were too small. The funder feature had over 2000 levels, where the majority were NA’s once grouped. It is unfortunate that these variables are ignored.



Risks and Limitations

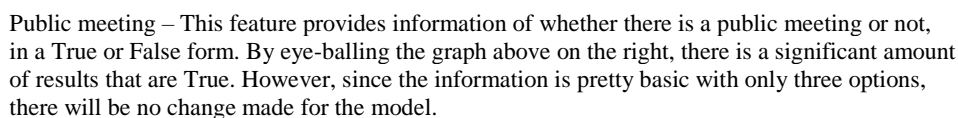
Before we move into the exploratory data analysis, we will need to identify the changes we will make to the data, and what types of risks and limitations there are. Based on the preliminary data analysis, there are few risks and limitations with the data:

1. The missing values will either be omitted by logic and assumptions – there may be 0’s and NA (not available data)’s resulting from mistakes and no proper methods of collecting certain data
2. The duplicate features will be removed – features that are strongly correlated does not add value in the model

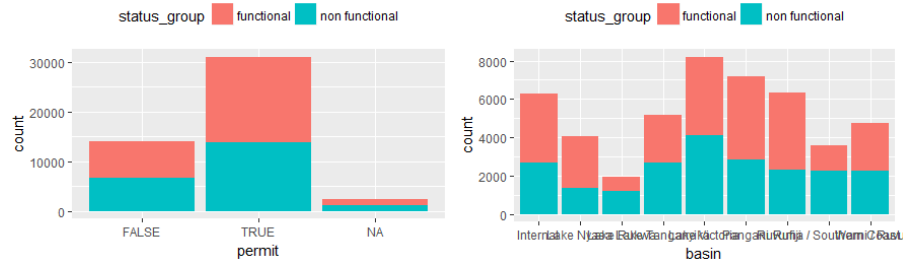
- Once the data is cleaned and ready for further analysis, various classification models (logistic regression, random forest, generalized boosting regression model, XGBoost, SVM) will be considered for use to find the best model for this case. The best model will have a significant increase compared to the baseline performance. A feature with higher coefficients would have a strong relationship with the dependant variable, and therefore would be the cause for the non-functional waterpoints.

After visualizing and going through all features that have a significant amount of data for the feature, we decided to compare them to reduce the amount of noise when creating the model. The features that had the most significant data with minimal missing fields were region, public meeting, permit, basin, payment type, management group, extraction type group, waterpoint type, GPS height, quantity, water quality, scheme management, and source. The following features has minimal amount of missing data and so we decided they were significant enough to contribute to the model.

Commented [TS1]: Talk about the relationship between your target variable (Functional - non Functional) with the other variables

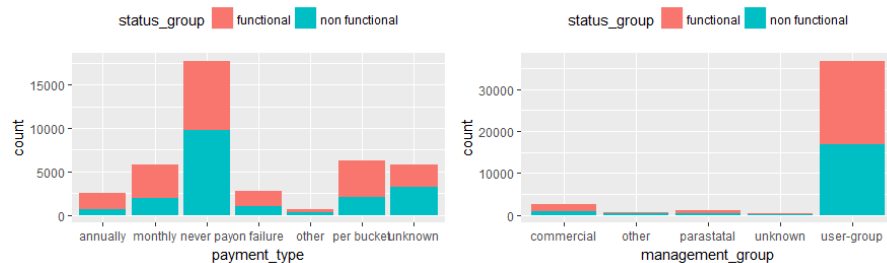


Permit – This feature informs us if the waterpoint is permitted to operate. The graph on the bottom left shows that the information is straightforward and does not need additional modification.



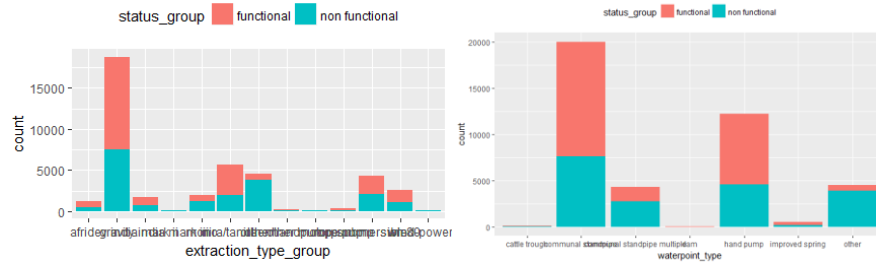
Basin – The feature gives us the geographic water basin location. The bar graph on the top right shows the variability of the data and will be a good feature for the predictive model to identify the significance of each water basin.

Payment type – This feature identifies how the waterpoint is being paid. There is a good amount of types within the feature, however, the bottom left graph shows that there is a significant amount of those who never pay. We can assume here that waterpoints that are never paid can lead poor maintenance updates, leading to non-functional waterpoints. Since the amount was obviously different from other types, we will create a new feature that groups all features other than “never pay” to cut the noise level for this feature. The new feature will be called “paytype0” and will be used in the predictive model.



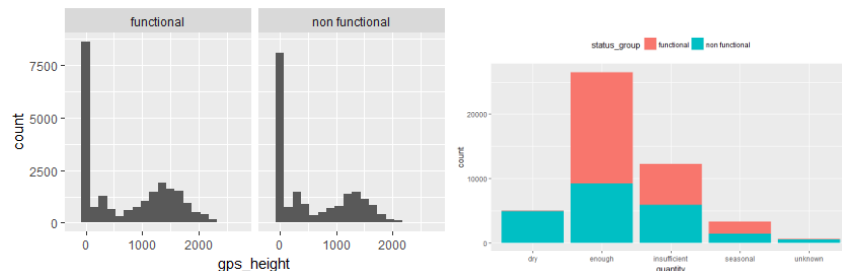
Management group This feature informs us about the organization who is managing the waterpoints. The graph above on the right clearly shows that there is one management group type, “user-group”, that dominates. Similar to the payment type earlier, we will group the rest of the types and decide whether or not the management group is a “user-group” to use in the predictive model as “management0”.

Extraction type group – This feature gives us what kind of extraction the waterpoint uses. The second extraction group type, “gravity”, seems to be the dominating factor for this feature, and therefore the rest will be grouped to create a new feature that measures the significance of this particular type.



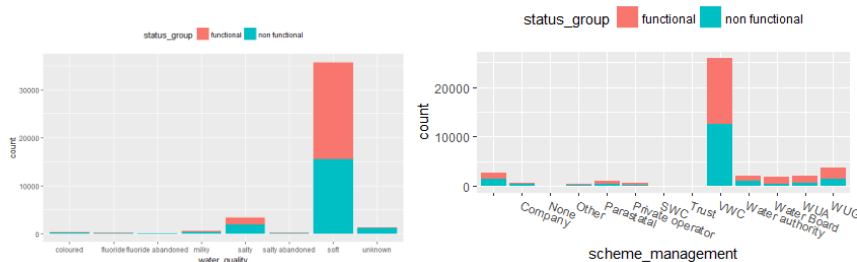
Waterpoint type – This feature explains the various types there are for a waterpoint. There are few types which dominates the waterpoint types. Therefore, we will group the most insignificant data with “other” group to decrease the noise level.

GPS height – This feature is the altitude of the well. As you can see in the graph below on the left-hand side, most of them had a value of 0 which meant that there is no proper pump installed. Since this may be a sign for a dysfunctional waterpoint, we have grouped any data that is greater than 0 to be “other”, and called it the “gps0” feature.



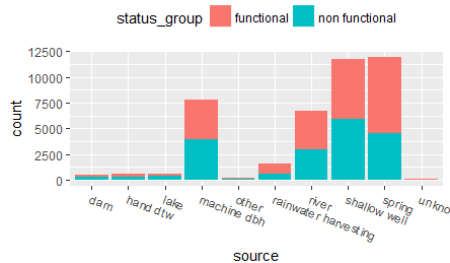
Quantity – This feature provides the information of how much water is within the area of the waterpoint. The graph on the top right hand side is an example of a feature that has enough variability where no change was made. With eye-balling we can clearly see that any waterpoint that is dry relate to a non-functional waterpoint and we will test through the model if it shows the significance.

Water quality – This feature provides the information of the quality of the water. From the bar graph on the left-hand side below, we see that the water quality type “soft” has the most amount of data and needs modification. We have created a new feature where we group the information based on whether it was soft or not. We can clearly identify if this was significant rather than having one outlier skewing the result, if the model was run with no adjustments.



Scheme management – This feature informs of the organization who operates the waterpoint. The graph on the top right-hand side shows that organization, VWC, controls the majority of the operation. Therefore, we will create a new feature that will group everything other than “VWC” to be “other” and use the feature for the predictive model.

Source – This feature is the source of the water for the waterpoint. The bar graph below has few types such as spring, shallow well, and river that has majority of the source types. Therefore we will reduce the level of types and reduce the noise of the feature to predict the non-functional waterpoints with the more significant data types.



With the features decided, we carefully looked at the model that makes the most sense in predicting the accuracy of non-functional waterpoints. In the beginning, the methods considered were logistic regression, random forest, SVM, or XGBoost. After carefully examining how different predictors calculates the accuracy, logistic regression was the simplest algorithm that can predict accurately with lowering the noise of the features.

The first step was to remove features such as latitude and longitude, which are the coordination points in a map. For a logistic regression model, it would see the numbers as significant, when there is no meaning in the numbers specifically. After removing them, any outliers/ data with high levels were grouped so that the noise of the model is reduced. Then the features were used in the logistic regression function. The features with less significance (higher p-value), were removed from the model and resulted with a prediction model to be tested with the current training and test data sets.

Final Results

As all features were modified and ready for final analysis, we have considered various models to best predict the waterpoints. We have mentioned earlier that we have considered logistic regression, random forest, generalized boosting regression model, XGBoost, and SVM. With the preliminary and explorative data analysis, we have realized that random forest, generalized boosting regression model, XGBoost, and SVM do not serve the purpose for our prediction. The above models will use classification trees to group the features into various groups. The findings will be useful to predict the pathways to a faulty waterpoint, but it is not what we are looking for.

Therefore, we decided to use the logistic regression model as it best predicts categorical value. For simplicity, we have decided to use the binomial logistic regression and treat “functional but needs repair” status groups to “non-functional” as we would like to predict whether the waterpoint is functional (0) or not (1).

The logistic regression model that consists of region, public meeting, basin, permit, group of payment type, group of extraction type, group of GPS height, waterpoint type, quantity, group of scheme management and source, has predicted insignificant and significant features in the model (see full logistic regression summary in Appendix 2). Water quality and management group was omitted as it was very insignificant and the updated information are the results.

Amongst the many features that were significant quantity had the most significant results with all types with a p-value less than $2e-16$ and the coefficients ranging from 2-4 being the most significant as it is greater than 1. The next few significant features included the waterpoint type and some of the regions where they also have significant p-values less than $2e-16$.

The predictions of the training and test set resulted with an accuracy of 72% and 69% respectively. This model predicted the status groups of the waterpoint with more than 18% accuracy compared to the baseline performance. Not only it provided us with a better prediction, but the coefficients of the features, informs us that quantity, region, and waterpoint type are a good predictor to why the waterpoints are at fault. With this predictive model, we would like to inform Taarifa and the Tanzanian Ministry of Water the next steps.

Recommendation

For our recommendation, we would like to propose a solution that will LIT up the community of Tanzania. We would like you to learn(L) from your mistakes, improve(I) on the current data, and track(T) new data to improve the model and create a better data hygiene to figure out insights.

Taarifa and the Tanzanian Ministry of Water will learn from the current model, where you will locate any waterpoint that are dry, regions that produced high coefficients via the predictive model, and waterpoint type that was classified as “other”. These were the most significant reasons for the default of the waterpoints. The following areas will be a starting point find out the faults and further recorded.

Commented [TS2]: Talk more about logistic regression and the specific findings based on numbers

Improving on the current data is also a critical task as there are useful features to include in the model. I have introduced the water mount near the waterpoint which is a feature that could have been used, if the amount of water was recorded. Population is another feature to be considered in the next model as well. There is no meaning in having a waterpoint where there are 0 population. This information needs to be properly collected to be used in the model. Lastly, the funders and installers will create a lot of noise in the model with the amount of high levels. The government can use the model to predict which of them creates the best waterpoint type that will last long to supply the citizens.

To expand on learning and improving the features and models, there is also the need of tracking new information that will make positive additions for the future. While accessing the waterpoints that need to be fixed, I would recommend implementing new data features that will make a difference in future models. The recommendation is to record the failure reasons of the waterpoints, track the remove data to find the average lifetime use of a waterpoint, and track the shortest distance for closest and alternative waterpoints.

In conclusion, there may be a lot of improvements to finding a solution for the non-functional waterpoints. However, this report includes the first analysis of the key problem that is causing the issue, and further recommend the next steps to quickly learn more about the waterpoints. With a better data hygiene and keeping in record of the data with few mediums, the data would be clean and more efficient. Please take the time to LIT up your country by learning from these insights, improving the features, and tracking new information that will help you quickly identify the next faulty waterpoint and have a brighter future.

Appendix 1

Full Data Features

- amoun_tsh - Total static head (amount water available to water point)
- date_recorded - The date the row was entered
- funder - Who funded the well
- gps_height - Altitude of the well
- installer - Organization that installed the well
- longitude - GPS coordinate
- latitude - GPS coordinate
- wpt_name - Name of the water point if there is one
- num_private -
- basin - Geographic water basin
- subvillage - Geographic location
- region - Geographic location
- region_code - Geographic location (coded)
- district_code - Geographic location (coded)
- lga - Geographic location
- ward - Geographic location
- population - Population around the well
- public_meeting - True/False
- recorded_by - Group entering this row of data
- scheme_management - Who operates the waterpoint
- scheme_name - Who operates the waterpoint
- permit - If the water point is permitted
- construction_year - Year the water point was constructed
- extraction_type - The kind of extraction the water point uses
- extraction_type_group - The kind of extraction the water point uses
- extraction_type_class - The kind of extraction the water point uses
- management - How the water point is managed
- management_group - How the water point is managed
- payment - What the water costs
- payment_type - What the water costs
- water_quality - The quality of the water
- quality_group - The quality of the water
- quantity - The quantity of water
- quantity_group - The quantity of water
- source - The source of the water
- source_type - The source of the water
- source_class - The source of the water
- waterpoint_type - The kind of water point
- waterpoint_type_group - The kind of water point

Appendix 2

```
Call:
glm(formula = status_group ~ region + public_meeting + basin +
    permit + paytype0 + exttype0 + waterpoint_type + gps0 + quantity +
    schman0 + source, family = binomial(link = "logit"), data = train)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-3.5456	-0.8660	-0.4973	0.9517	2.6108

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	3.092745	0.151580	20.403	< 2e-16	***
regionDar es Salaam	0.709091	0.216332	3.278	0.001046	**
regionDodoma	0.781119	0.139231	5.610	2.02e-08	***
regionIringa	0.597443	0.134971	4.426	9.58e-06	***
regionKagera	0.670653	0.147821	4.537	5.71e-06	***
regionKigoma	1.530135	0.124646	12.276	< 2e-16	***
regionKilimanjaro	0.677025	0.072811	9.298	< 2e-16	***
regionLindi	1.812724	0.192048	9.439	< 2e-16	***
regionManyara	0.251624	0.099806	2.521	0.011698	*
regionMara	1.273491	0.147268	8.647	< 2e-16	***
regionMbeya	1.978519	0.156677	12.628	< 2e-16	***
regionMorogoro	1.374668	0.127402	10.790	< 2e-16	***
regionMtwara	1.623765	0.195510	8.305	< 2e-16	***
regionMwanza	0.498003	0.147328	3.380	0.000724	***
regionPwani	0.911553	0.135606	6.722	1.79e-11	***
regionRukwa	1.589968	0.133351	11.923	< 2e-16	***
regionRuvuma	1.283748	0.166147	7.727	1.10e-14	***
regionShinyanga	1.151707	0.134194	8.582	< 2e-16	***
regionSingida	0.794021	0.121545	6.533	6.46e-11	***
regionTabora	1.215902	0.143261	8.487	< 2e-16	***
regionTanga	-0.097934	0.082861	-1.182	0.237242	
public_meetingTRUE	-0.261154	0.045312	-5.763	8.24e-09	***
basinLake Nyasa	-1.584814	0.123537	-12.829	< 2e-16	***
basinLake Rukwa	0.023793	0.110280	0.216	0.829185	
basinLake Tanganyika	0.009776	0.080103	0.122	0.902869	
basinLake Victoria	0.295193	0.082142	3.594	0.000326	***
basinPangani	0.630663	0.086914	7.256	3.98e-13	***
basinRufiji	-0.337658	0.097323	-3.469	0.000522	***
basinRuvuma / Southern Coast	-0.401226	0.162073	-2.476	0.013302	*

Appendix 2 (continued)

basinWami / Ruvu	-0.314483	0.097204	-3.235	0.001215	**
permitTRUE	-0.186615	0.031664	-5.894	3.78e-09	***
paytype0	0.831793	0.028512	29.174	< 2e-16	***
exttype0	-0.416338	0.054626	-7.622	2.51e-14	***
waterpoint_typecommunal standpipe multiple	0.779068	0.044983	17.319	< 2e-16	***
waterpoint_typehand pump	-0.858593	0.051711	-16.604	< 2e-16	***
waterpoint_typeother	1.154954	0.050356	22.936	< 2e-16	***
gps0	-0.367755	0.084155	-4.370	1.24e-05	***
quantityenough	-4.494721	0.104640	-42.954	< 2e-16	***
quantityinsufficient	-3.927134	0.106461	-36.888	< 2e-16	***
quantityseasonal	-4.335746	0.113826	-38.091	< 2e-16	***
quantityunknown	-2.364238	0.197216	-11.988	< 2e-16	***
schman0	0.386636	0.031762	12.173	< 2e-16	***
sourceother	-0.046968	0.061434	-0.765	0.444557	
sourceciver	-0.009966	0.062502	-0.159	0.873309	
sourceshallow well	0.205192	0.048479	4.233	2.31e-05	***
sourcespring	-0.307593	0.064360	-4.779	1.76e-06	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 51342 on 37261 degrees of freedom
Residual deviance: 39104 on 37216 degrees of freedom
(4319 observations deleted due to missingness)
AIC: 39196

Number of Fisher Scoring iterations: 6