

Nigeria Water Quality Prediction by Machine Learning

Ranice TAN Hui Qi

Contents

1. OVERVIEW	2
2. OBJECTIVE.....	2
3. DATA	2
3.1. DATA USED.....	2
3.2. DATA PREPARATION	2
3.3. DATA QUALITY	3
a. <i>Reviewing the data and metadata</i>	3
b. <i>Checking for missing data</i>	3
c. <i>Univariate data analysis</i>	4
d. <i>Bivariate data analysis</i>	4
e. <i>Data Sampling</i>	5
4. DATA ANALYSIS.....	6
4.1. CALIBRATING BASE MODEL	6
4.2. NOMINAL LOGISTIC REGRESSION	7
4.3. STEPWISE LOGISTIC REGRESSION (FORWARD)	8
4.4. STEPWISE LOGISTIC REGRESSION (BACKWARD)	10
4.5. DECISION TREE	11
4.6. BOOTSTRAP FOREST.....	12
4.7. MODEL COMPARISON	13
5. CONCLUSION	16
6. RECOMMENDATION.....	16
7. APPENDIX	17
7.1. APPENDIX A: FIELDS IN ORIGINAL DATASET	17
7.2. APPENDIX B: DATA PREPARATION CHANGE LOG	19

1. Overview

Ensuring the availability and sustainable management of water and sanitation is one of the Sustainable Development Goal identified by the United Nations. However, providing clean potable water to rural communities in many developing countries remain challenging due to water scarcity and poor water quality. The lack of clean water poses a threat to a community's health, safety, and economy.

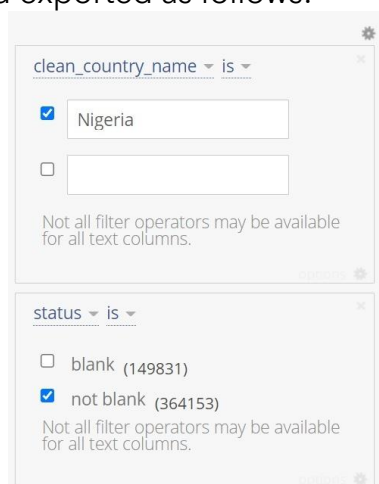
2. Objective

The objective of the analysis is to build models to predict the reliability of water points in Nigeria, so that preventive maintenance can be carried out on urgent Non-Functional systems timely. The models will also be compared to identify a model with a suitable accuracy and not overly complex.

3. Data

3.1. Data Used

The dataset is obtained from the Water Point Data Exchange (WPdx) Data Repository¹. The repository collects water point related data from rural areas at the water point or small water scheme level, based on WPdx Data Standard, and shared on a cloud-based data library. The data for 'Nigeria' was filtered and exported as follows:



The screenshot shows a filter interface with two sections. The first section is for 'clean_country_name' with a dropdown menu set to 'is'. Below it, there is a checked checkbox for 'Nigeria' and an unchecked checkbox for an empty text box. The second section is for 'status' with a dropdown menu set to 'is'. Below it, there is an unchecked checkbox for 'blank (149831)' and a checked checkbox for 'not blank (364153)'. Both sections include a note: 'Not all filter operators may be available for all text columns.'

Figure 1 Filter criterion for WPdx Data Repository

3.2. Data Preparation

The dataset was imported into JMP Pro to ensure all fields are filled, and all columns are appropriately formatted. The fields in the original dataset can be found in 17. A preliminary inspection of the data summary statistics and distribution was also conducted. The data preparation log is accessible in 19.

¹ <https://data.waterpointdata.org/dataset/Water-Point-Data-Exchange-Plus-WPdx-/eqje-vguj/data>

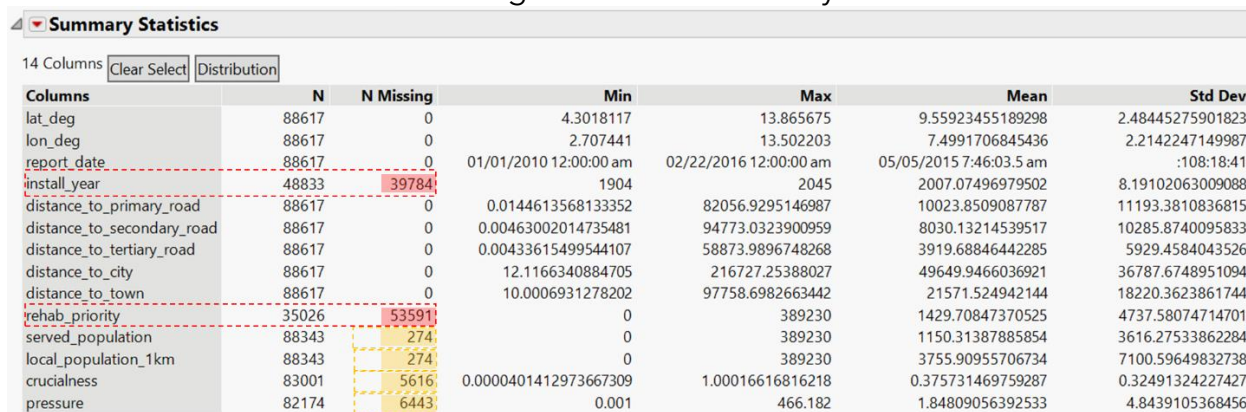
3.3. Data Quality

a. Reviewing the data and metadata

The raw data was extracted and reviewed against the WPdx data standard². The list of fields in the dataset can be found in 7.1. The metadata was first explored to ensure all data is in the correct form. Next, data that are not updated were removed from the dataset. Finally, the data that are irrelevant or redundant to the purpose of the analysis were removed. The amendments made can be found in 1-36 in the Data Change Log (19).

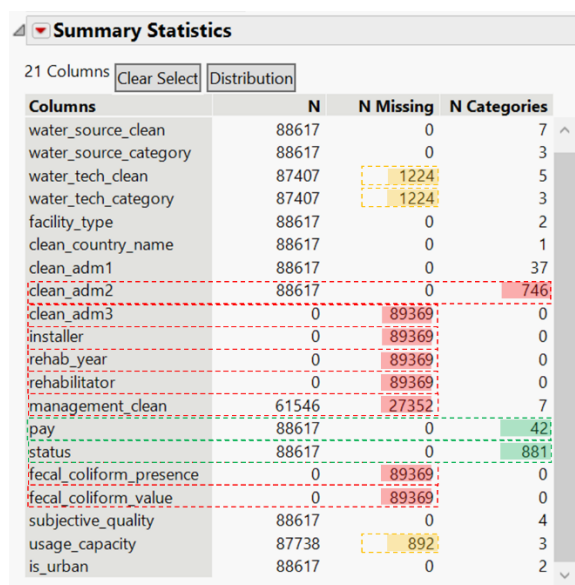
b. Checking for missing data

Column viewer was then used to generate the summary statistics as shown below:



Columns	N	N Missing	Min	Max	Mean	Std Dev
lat_deg	88617	0	4.3018117	13.865675	9.55923455189298	2.48445275901823
lon_deg	88617	0	2.707441	13.502203	7.4991706845436	2.2142247149987
report_date	88617	0	01/01/2010 12:00:00 am	02/22/2016 12:00:00 am	05/05/2015 7:46:03.5 am	:108:18:41
install_year	48833	39784	1904	2045	2007.07496979502	8.19102063009088
distance_to_primary_road	88617	0	0.0144613568133352	82056.9295146987	10023.8509087787	11193.3810836815
distance_to_secondary_road	88617	0	0.00463002014735481	94773.0323900959	8030.13214539517	10285.8740095833
distance_to_tertiary_road	88617	0	0.00433615499544107	58873.9896748268	3919.68846442285	5929.4584043526
distance_to_city	88617	0	12.1166340884705	216727.25388027	49649.9466036921	36787.6748951094
distance_to_town	88617	0	10.0006931278202	97758.6982663442	21571.524942144	18220.3623861744
rehab_priority	35026	53591	0	389230	1429.70847370525	4737.58074714701
served_population	88343	274	0	389230	1150.31387885854	3616.27533862284
local_population_1km	88343	274	0	389230	3755.90955706734	7100.59649832738
crucialness	83001	5616	0.0000401412973667309	1.00016616816218	0.375731469759287	0.32491324227427
pressure	82174	6443	0.001	466.182	1.84809056392533	4.8439105368456

Figure 2 Summary Statistics of Continuous Variables



Columns	N	N Missing	N Categories
water_source_clean	88617	0	7
water_source_category	88617	0	3
water_tech_clean	87407	1224	5
water_tech_category	87407	1224	3
facility_type	88617	0	2
clean_country_name	88617	0	1
clean_adm1	88617	0	37
clean_adm2	88617	0	746
clean_adm3	0	89369	0
installer	0	89369	0
rehab_year	0	89369	0
rehabilitator	0	89369	0
management_clean	61546	27352	7
pay	88617	0	42
status	88617	0	881
fecal_coliform_presence	0	89369	0
fecal_coliform_value	0	89369	0
subjective_quality	88617	0	4
usage_capacity	87738	892	3
is_urban	88617	0	2

Figure 3 Summary Statistics of Nominal Variables

Firstly, columns with relatively large amount of missing data (in red) were excluded (> 25000 rows) to avoid inaccurate model due to lack of data. Next, nominal variables with large number of categories were either removed (in red) or recoded (in green) based on the importance of the variables. Lastly, useful variables with rows containing

² https://www.waterpointdata.org/wp-content/uploads/2021/04/WPDx_Data_Standard.pdf

missing data (in yellow) were excluded to prevent interfering with the model results. A summary of the changes made can be found in 37-57 in the Data Change Log (19).

c. Univariate data analysis

The distribution of the remaining variables was explored and those with relatively uneven distribution were identified below. Variables with uneven distribution were recoded in 58-61 of the Data Change Log (19).



Figure 4 Uneven Distribution of variables

d. Bivariate data analysis

A pairwise correlation for continuous variables was then carried out to identify highly correlated variables that can be interchanged. Based on the results below, the variables pressure, served_population and local_population_1km seems to be relatively highly correlated with one another.

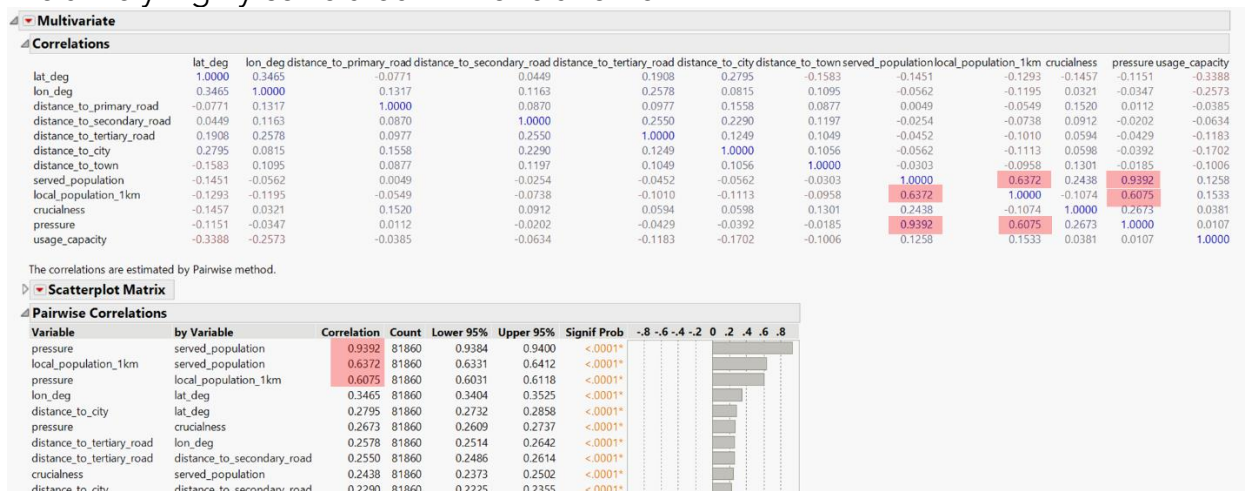


Figure 5 Pairwise Correlations of continuous variables

A principal component analysis was then conducted to identify how many variables can be excluded. Based on the eigenvalue table, the first 3 principal components have an eigenvalue of greater than 1.

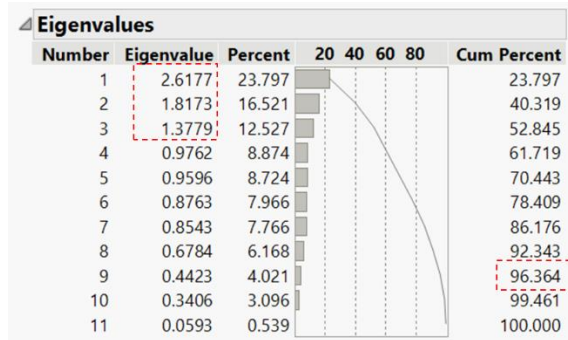


Figure 6 Eigenvalues of pairwise correlation variables

Finally, variable pressure was excluded and the analysis proceeded with 10 variables, accounting for 99% of the variance.

e. Data Sampling

A validation column was added to split the data into training, validation and test set, stratified by the STATUS recoded column.

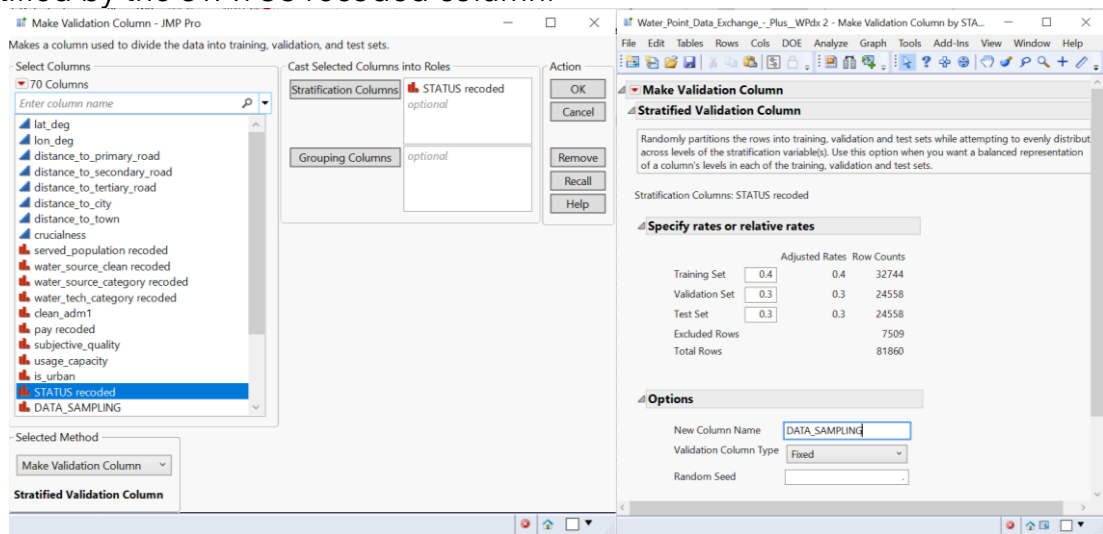


Figure 7 Making Validation Column

4. Data Analysis

4.1. Calibrating Base Model

In the dataset, there are a total of 20 predictors. 8 of them are continuous predictors and the other 12 are categorical predictors. Among the categorical predictors, only 3 is binary and the remaining 9 have more than two classes. The base nominal logistic model was generated as follows:

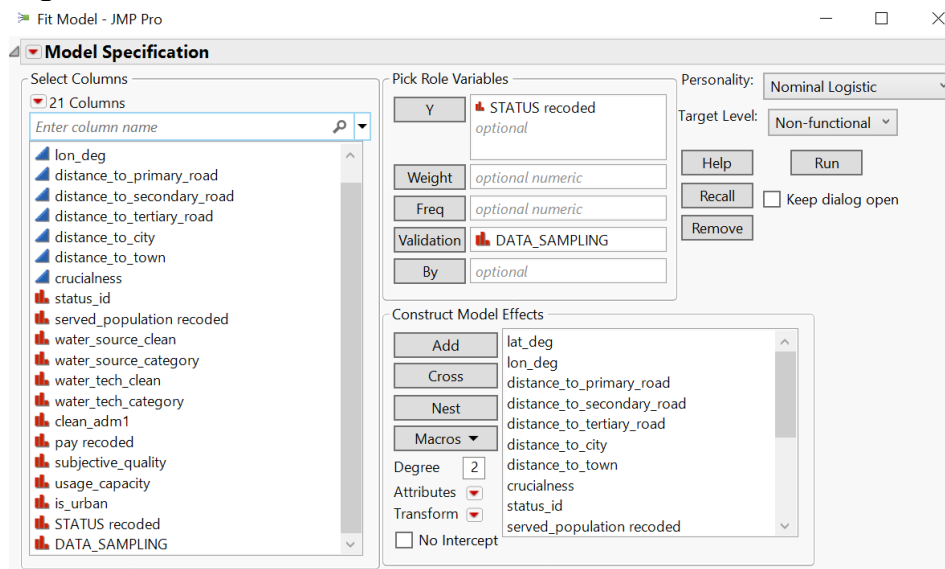


Figure 8 Model specification for base model

By interpreting the model through the parameter estimates, the status_id is marked as an unstable estimate, and the water_sources and water_techs are marked as biased or zeroed. The binary status_id variable was then excluded from the model, the water_source_clean, water_source_category and water_tech_category was recoded into binary predictors. water_tech_clean was removed as if recoded into binary, it would be a duplicate of water_tech_category.

Parameter Estimates				
Term		Estimate	Std Error	ChiSquare Prob>ChiSq
Intercept	Biased	-5.7722e-9	1652264.8	0.00 1.0000
lat_deg		1.6704e-13	72429.936	0.00 1.0000
lon_deg		1.1454e-13	74851.357	0.00 1.0000
distance_to_primary_road		2.0666e-19	3.2590273	0.00 1.0000
distance_to_secondary_road		2.017e-19	3.8027027	0.00 1.0000
distance_to_tertiary_road		-1.53e-18	6.4373221	0.00 1.0000
distance_to_city		6.325e-19	1.2224955	0.00 1.0000
distance_to_town		5.0114e-19	2.587149	0.00 1.0000
crucialness		3.5953e-14	124411.63	0.00 1.0000
status_id[No]	Unstable	31.2028948	36816.845	0.00 0.9993
served_population recoded[> 100000]		-5.02e-13	2008674.4	0.00 1.0000
served_population recoded[1 - 1000]		1.6897e-13	676560.34	0.00 1.0000
served_population recoded[1001 - 10000]		1.725e-13	675170.21	0.00 1.0000
water_source_clean[Borehole]	Biased	-6.937e-15	175919.61	0.00 1.0000
water_source_clean[Protected Shallow Well]	Biased	-6.842e-15	215360.23	0.00 1.0000
water_source_clean[Protected Spring]	Zeroed	0	0	. .
water_source_clean[Spring]	Zeroed	0	0	. .
water_tech_clean[Hand Pump]	Biased	9.7402e-13	1272044.2	0.00 1.0000
water_tech_clean[Mechanized Pump]	Biased	-2.95e-13	1193406.3	0.00 1.0000
water_tech_clean[Mechanized Pump - Diesel]	Biased	-1.748e-12	2936572.6	0.00 1.0000
water_tech_clean[Mechanized Pump - Solar]	Biased	-3.087e-13	1245664	0.00 1.0000
water_tech_category[Hand Pump]	Zeroed	0	0	. .
water_tech_category[Mechanized Pump]	Zeroed	0	0	. .
clean_adm1[Abia]		5.537e-13	389781.24	0.00 1.0000
clean_adm1[Adamawa]		-7.953e-13	693517.59	0.00 1.0000
clean_adm1[Akwa Ibom]		5.5491e-13	355069.62	0.00 1.0000
clean_adm1[Anambra]		4.4838e-13	468434.43	0.00 1.0000

Figure 9 Fit Details and Parameter Estimates of Base Model

The final dataset consists of a total of 17 predictors - 8 which are continuous and 9 which are categorical predictors. Among the 9 categorical predictors, 4 are binary and the remaining 5 have more than two classes.

4.2. Nominal Logistic Regression

Using the revised data table from calibration, the result below was achieved. The model has a misclassification rate of 0.3, a non-functional prediction rate of 0.47, and a functional prediction rate of 0.83. The false positive error is 0.52 and false negative error is 0.17.

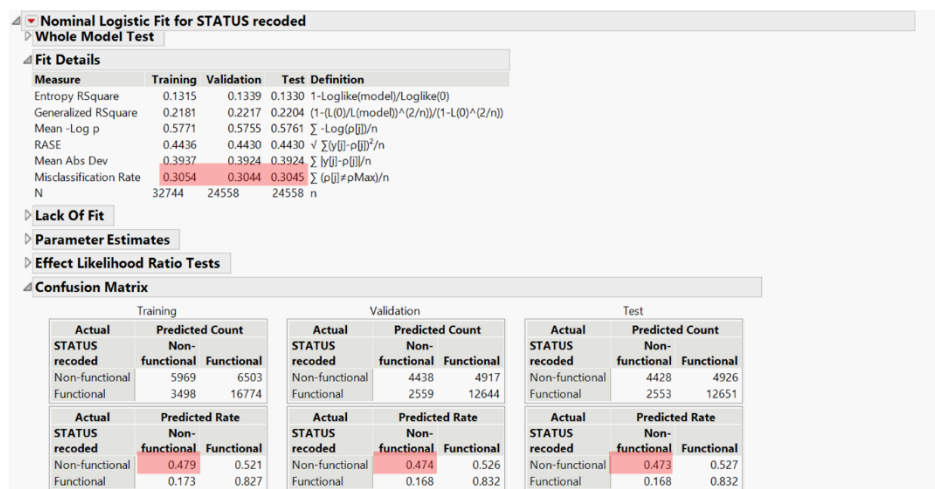


Figure 10 Results for Nominal Logistic Regression

Profit/Cost Decision Matrix

Specify Profit Matrix

Each matrix entry is the profit if you predict the response in the column when the response in the row is the actual response.

Enter values that reflect profits for correct decisions on the diagonal.
Enter values (usually negative) that reflect profits for incorrect decisions elsewhere.
Use the Undecided column to reflect profits for an alternative decision.

When you save prediction formulas, these values will be used to create best decision columns.
The best decision is the one with greatest expected profit.

Decision or Prediction

Actual \ Decision	Non-functional	Functional	Undecided
Non-functional	0	-1	.
Functional	-0.6129	0	.

To create a profit matrix for a binary response, enter a Target and Probability Threshold.
If the predicted probability exceeds the threshold, the best decision will be the target.

Target Level: Non-functional

Probability Threshold: 0.38

☐ Save to column as property.

OK Cancel

Figure 11 Specify probability threshold for non-functional level

From the effect summary, the variables that influence the model most are administrative area, quality of water, local population within 1km of the water system, payment system, crucialness of system, population served and water source.

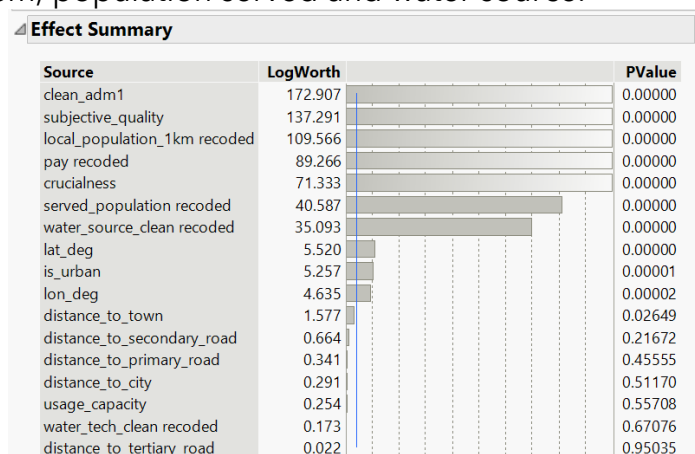


Figure 12 Effect Summary of Nominal Logistic Regression Model

From the formula of the model below, the non-functional water systems have population within 1km and served population ≤ 10000 . The sources also tend to come from boreholes, do not require payment, have unacceptable quality of water, and a low usage capacity of 250.

Lower population means lower usage of water, which may have more stagnant pipelines, resulting in an increased risk of biofilm-associated bacteria, or spread of Legionella. Hence, resulting in unacceptable quality of water. Without payment for usage of the water system, more funds will need to be obtained to carry out preventive maintenance.

+ Match(<i>clean_adm1</i>)	<table> <tr><td>"Abia"</td><td>⇒ 1.1170210786</td></tr> <tr><td>"Adamawa"</td><td>⇒ -0.143360395</td></tr> <tr><td>"Akwa Ibom"</td><td>⇒ 1.0679197633</td></tr> <tr><td>"Anambra"</td><td>⇒ 0.4021211032</td></tr> <tr><td>"Bauchi"</td><td>⇒ -0.696589969</td></tr> <tr><td>"Bayelsa"</td><td>⇒ 0.2842302442</td></tr> <tr><td>"Benue"</td><td>⇒ 0.3532064814</td></tr> <tr><td>"Borno"</td><td>⇒ -0.698458176</td></tr> <tr><td>"Cross River"</td><td>⇒ 0.9468670777</td></tr> <tr><td>"Delta"</td><td>⇒ 0.3650979272</td></tr> <tr><td>"Ebonyi"</td><td>⇒ -0.028555285</td></tr> <tr><td>"Edo"</td><td>⇒ 0.1822396568</td></tr> <tr><td>"Ekiti"</td><td>⇒ 0.4284941711</td></tr> <tr><td>"Enugu"</td><td>⇒ 0.0398419027</td></tr> <tr><td>"Federal Capital Territory"</td><td>⇒ 0.2908821834</td></tr> <tr><td>"Gombe"</td><td>⇒ 0.6959838107</td></tr> <tr><td>"Imo"</td><td>⇒ 0.5984870002</td></tr> <tr><td>"Jigawa"</td><td>⇒ -1.401581359</td></tr> <tr><td>"Kaduna"</td><td>⇒ -0.343103221</td></tr> <tr><td>"Kano"</td><td>⇒ -1.095126287</td></tr> <tr><td>"Katsina"</td><td>⇒ -0.891563864</td></tr> <tr><td>"Kebbi"</td><td>⇒ -0.359683368</td></tr> <tr><td>"Kogi"</td><td>⇒ 0.6714265927</td></tr> <tr><td>"Kwara"</td><td>⇒ -0.202576447</td></tr> <tr><td>"Lagos"</td><td>⇒ -0.23021112</td></tr> <tr><td>"Nassarawa"</td><td>⇒ 0.1019055325</td></tr> <tr><td>"Niger"</td><td>⇒ -0.212276111</td></tr> <tr><td>"Ogun"</td><td>⇒ 0.1486210324</td></tr> <tr><td>"Ondo"</td><td>⇒ 0.6103058381</td></tr> <tr><td>"Osun"</td><td>⇒ 0.0235115274</td></tr> <tr><td>"Oyo"</td><td>⇒ -0.030848953</td></tr> <tr><td>"Plateau"</td><td>⇒ 0.2187527892</td></tr> <tr><td>"Rivers"</td><td>⇒ -0.196552392</td></tr> <tr><td>"Sokoto"</td><td>⇒ -0.569751648</td></tr> <tr><td>"Taraba"</td><td>⇒ 0.8062101319</td></tr> <tr><td>"Yobe"</td><td>⇒ -1.396899071</td></tr> <tr><td>"Zamfara"</td><td>⇒ -0.855988178</td></tr> <tr><td>else</td><td>⇒ .</td></tr> </table>	"Abia"	⇒ 1.1170210786	"Adamawa"	⇒ -0.143360395	"Akwa Ibom"	⇒ 1.0679197633	"Anambra"	⇒ 0.4021211032	"Bauchi"	⇒ -0.696589969	"Bayelsa"	⇒ 0.2842302442	"Benue"	⇒ 0.3532064814	"Borno"	⇒ -0.698458176	"Cross River"	⇒ 0.9468670777	"Delta"	⇒ 0.3650979272	"Ebonyi"	⇒ -0.028555285	"Edo"	⇒ 0.1822396568	"Ekiti"	⇒ 0.4284941711	"Enugu"	⇒ 0.0398419027	"Federal Capital Territory"	⇒ 0.2908821834	"Gombe"	⇒ 0.6959838107	"Imo"	⇒ 0.5984870002	"Jigawa"	⇒ -1.401581359	"Kaduna"	⇒ -0.343103221	"Kano"	⇒ -1.095126287	"Katsina"	⇒ -0.891563864	"Kebbi"	⇒ -0.359683368	"Kogi"	⇒ 0.6714265927	"Kwara"	⇒ -0.202576447	"Lagos"	⇒ -0.23021112	"Nassarawa"	⇒ 0.1019055325	"Niger"	⇒ -0.212276111	"Ogun"	⇒ 0.1486210324	"Ondo"	⇒ 0.6103058381	"Osun"	⇒ 0.0235115274	"Oyo"	⇒ -0.030848953	"Plateau"	⇒ 0.2187527892	"Rivers"	⇒ -0.196552392	"Sokoto"	⇒ -0.569751648	"Taraba"	⇒ 0.8062101319	"Yobe"	⇒ -1.396899071	"Zamfara"	⇒ -0.855988178	else	⇒ .
"Abia"	⇒ 1.1170210786																																																																												
"Adamawa"	⇒ -0.143360395																																																																												
"Akwa Ibom"	⇒ 1.0679197633																																																																												
"Anambra"	⇒ 0.4021211032																																																																												
"Bauchi"	⇒ -0.696589969																																																																												
"Bayelsa"	⇒ 0.2842302442																																																																												
"Benue"	⇒ 0.3532064814																																																																												
"Borno"	⇒ -0.698458176																																																																												
"Cross River"	⇒ 0.9468670777																																																																												
"Delta"	⇒ 0.3650979272																																																																												
"Ebonyi"	⇒ -0.028555285																																																																												
"Edo"	⇒ 0.1822396568																																																																												
"Ekiti"	⇒ 0.4284941711																																																																												
"Enugu"	⇒ 0.0398419027																																																																												
"Federal Capital Territory"	⇒ 0.2908821834																																																																												
"Gombe"	⇒ 0.6959838107																																																																												
"Imo"	⇒ 0.5984870002																																																																												
"Jigawa"	⇒ -1.401581359																																																																												
"Kaduna"	⇒ -0.343103221																																																																												
"Kano"	⇒ -1.095126287																																																																												
"Katsina"	⇒ -0.891563864																																																																												
"Kebbi"	⇒ -0.359683368																																																																												
"Kogi"	⇒ 0.6714265927																																																																												
"Kwara"	⇒ -0.202576447																																																																												
"Lagos"	⇒ -0.23021112																																																																												
"Nassarawa"	⇒ 0.1019055325																																																																												
"Niger"	⇒ -0.212276111																																																																												
"Ogun"	⇒ 0.1486210324																																																																												
"Ondo"	⇒ 0.6103058381																																																																												
"Osun"	⇒ 0.0235115274																																																																												
"Oyo"	⇒ -0.030848953																																																																												
"Plateau"	⇒ 0.2187527892																																																																												
"Rivers"	⇒ -0.196552392																																																																												
"Sokoto"	⇒ -0.569751648																																																																												
"Taraba"	⇒ 0.8062101319																																																																												
"Yobe"	⇒ -1.396899071																																																																												
"Zamfara"	⇒ -0.855988178																																																																												
else	⇒ .																																																																												
+ 0.8806548915 * <i>crucialness</i>	<table> <tr><td>"> 100000"</td><td>⇒ -1.883335739</td></tr> <tr><td>"1 - 1000"</td><td>⇒ 1.24481025</td></tr> <tr><td>"1001 - 10000"</td><td>⇒ 0.9105313705</td></tr> <tr><td>"10001 - 100000"</td><td>⇒ -0.272005881</td></tr> <tr><td>else</td><td>⇒ .</td></tr> </table>	"> 100000"	⇒ -1.883335739	"1 - 1000"	⇒ 1.24481025	"1001 - 10000"	⇒ 0.9105313705	"10001 - 100000"	⇒ -0.272005881	else	⇒ .																																																																		
"> 100000"	⇒ -1.883335739																																																																												
"1 - 1000"	⇒ 1.24481025																																																																												
"1001 - 10000"	⇒ 0.9105313705																																																																												
"10001 - 100000"	⇒ -0.272005881																																																																												
else	⇒ .																																																																												
+ Match(<i>local_population_1km recoded</i>)	<table> <tr><td>"> 100000"</td><td>⇒ 2.0309288859</td></tr> <tr><td>"1 - 1000"</td><td>⇒ -1.271209376</td></tr> <tr><td>"1001 - 10000"</td><td>⇒ -0.818983004</td></tr> <tr><td>"10001 - 100000"</td><td>⇒ 0.0592634943</td></tr> <tr><td>else</td><td>⇒ .</td></tr> </table>	"> 100000"	⇒ 2.0309288859	"1 - 1000"	⇒ -1.271209376	"1001 - 10000"	⇒ -0.818983004	"10001 - 100000"	⇒ 0.0592634943	else	⇒ .																																																																		
"> 100000"	⇒ 2.0309288859																																																																												
"1 - 1000"	⇒ -1.271209376																																																																												
"1001 - 10000"	⇒ -0.818983004																																																																												
"10001 - 100000"	⇒ 0.0592634943																																																																												
else	⇒ .																																																																												
+ Match(<i>served_population recoded</i>)	<table> <tr><td>"> 100000"</td><td>⇒ 2.0309288859</td></tr> <tr><td>"1 - 1000"</td><td>⇒ -1.271209376</td></tr> <tr><td>"1001 - 10000"</td><td>⇒ -0.818983004</td></tr> <tr><td>"10001 - 100000"</td><td>⇒ 0.0592634943</td></tr> <tr><td>else</td><td>⇒ .</td></tr> </table>	"> 100000"	⇒ 2.0309288859	"1 - 1000"	⇒ -1.271209376	"1001 - 10000"	⇒ -0.818983004	"10001 - 100000"	⇒ 0.0592634943	else	⇒ .																																																																		
"> 100000"	⇒ 2.0309288859																																																																												
"1 - 1000"	⇒ -1.271209376																																																																												
"1001 - 10000"	⇒ -0.818983004																																																																												
"10001 - 100000"	⇒ 0.0592634943																																																																												
else	⇒ .																																																																												
+ Match(<i>water_source_clean recoded</i>)	<table> <tr><td>"Borehole"</td><td>⇒ 0.273731994</td></tr> <tr><td>"Others"</td><td>⇒ -0.273731994</td></tr> <tr><td>else</td><td>⇒ .</td></tr> </table>	"Borehole"	⇒ 0.273731994	"Others"	⇒ -0.273731994	else	⇒ .																																																																						
"Borehole"	⇒ 0.273731994																																																																												
"Others"	⇒ -0.273731994																																																																												
else	⇒ .																																																																												
+ Match(<i>pay recoded</i>)	<table> <tr><td>"No"</td><td>⇒ 0.5310284394</td></tr> <tr><td>"Yes"</td><td>⇒ -0.531028439</td></tr> <tr><td>else</td><td>⇒ .</td></tr> </table>	"No"	⇒ 0.5310284394	"Yes"	⇒ -0.531028439	else	⇒ .																																																																						
"No"	⇒ 0.5310284394																																																																												
"Yes"	⇒ -0.531028439																																																																												
else	⇒ .																																																																												
+ Match(<i>subjective_quality</i>)	<table> <tr><td>"Acceptable quality"</td><td>⇒ -0.702152481</td></tr> <tr><td>"No because of Colour"</td><td>⇒ 0.1110635428</td></tr> <tr><td>"No because of Odour"</td><td>⇒ 0.4246999554</td></tr> <tr><td>"No because of Taste"</td><td>⇒ 0.1663889825</td></tr> <tr><td>else</td><td>⇒ .</td></tr> </table>	"Acceptable quality"	⇒ -0.702152481	"No because of Colour"	⇒ 0.1110635428	"No because of Odour"	⇒ 0.4246999554	"No because of Taste"	⇒ 0.1663889825	else	⇒ .																																																																		
"Acceptable quality"	⇒ -0.702152481																																																																												
"No because of Colour"	⇒ 0.1110635428																																																																												
"No because of Odour"	⇒ 0.4246999554																																																																												
"No because of Taste"	⇒ 0.1663889825																																																																												
else	⇒ .																																																																												
+ Match(<i>usage_capacity</i>)	<table> <tr><td>250</td><td>⇒ 0.4415001739</td></tr> <tr><td>500</td><td>⇒ -0.172274156</td></tr> <tr><td>1000</td><td>⇒ -0.269226018</td></tr> <tr><td>else</td><td>⇒ .</td></tr> </table>	250	⇒ 0.4415001739	500	⇒ -0.172274156	1000	⇒ -0.269226018	else	⇒ .																																																																				
250	⇒ 0.4415001739																																																																												
500	⇒ -0.172274156																																																																												
1000	⇒ -0.269226018																																																																												
else	⇒ .																																																																												

Figure 13 Formulas of variables that impact the reliability of water system

4.3. Stepwise Logistic Regression (Forward)

A forward stepwise logistic regression was also carried out as follows. Minimum BIC was selected as the stopping rule due to the large dataset (~80000 rows) acquired.

4.4. Stepwise Logistic Regression (Backward)

A backward stepwise logistic regression was also carried out as follows. Minimum BIC was selected as the stopping rule due to the large dataset (~80000 rows) acquired.

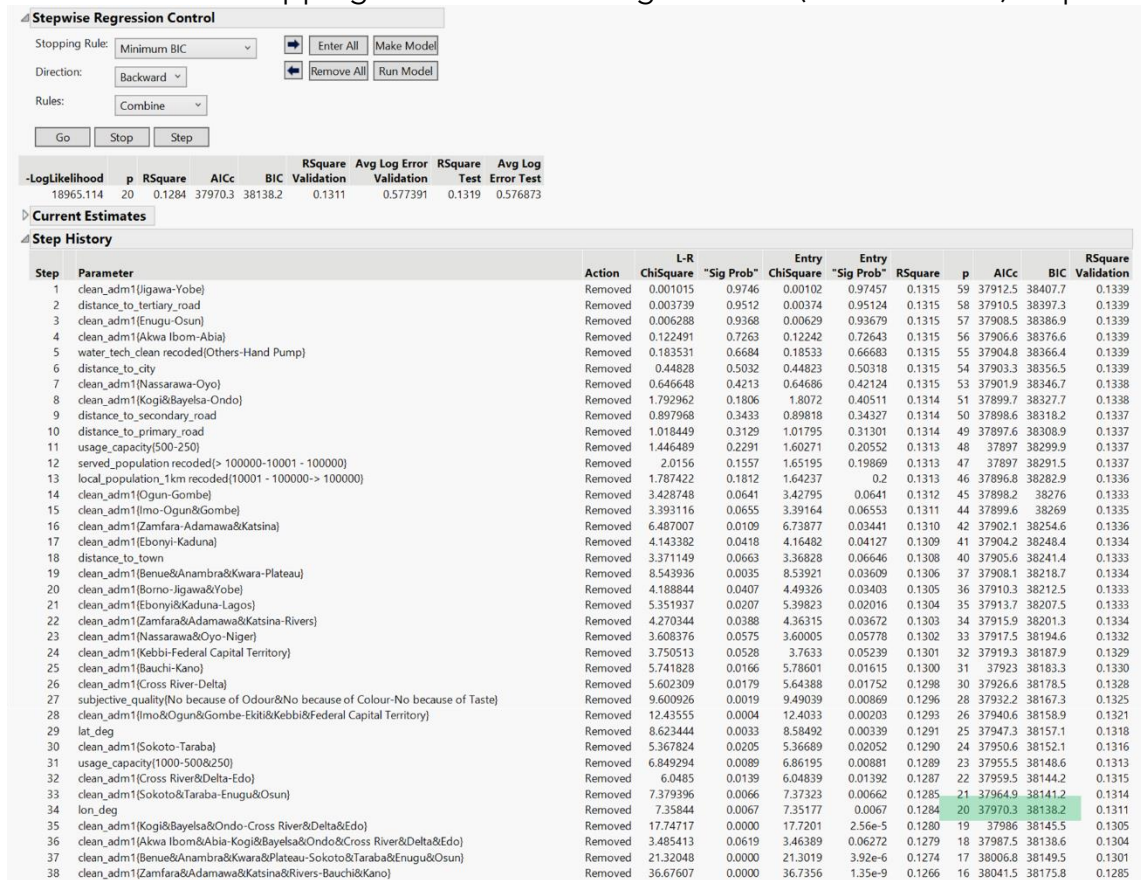


Figure 17 Stepwise (Backward) Regression Control

The model has a misclassification rate of 0.31, a non-functional prediction rate of 0.47, and a functional prediction rate of 0.83. The false positive error is 0.53 and false negative error is 0.17.

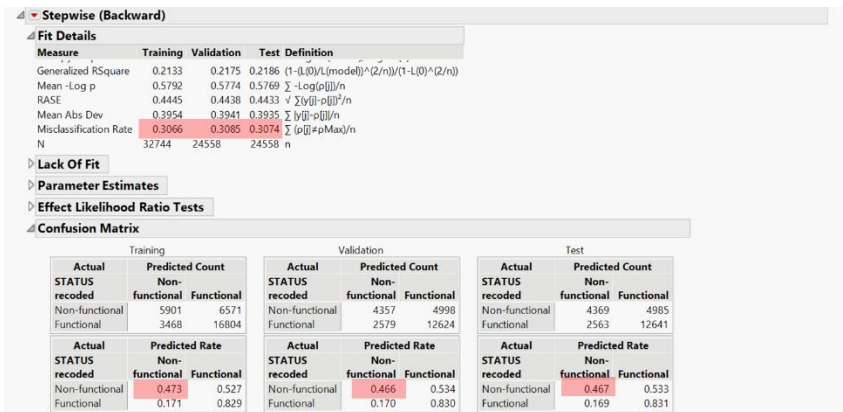


Figure 18 Results for Stepwise Logistic Regression (Backward)

From the step history, the variables that influence the model most are local population within 1km, administrative area and whether the system is in an urban area. The variables are like some of the variables in the previous logistic regression. They are also directionally similar for predicting non-functionality.

while other aspects such as microorganism and bacterial growth cannot be assessed through this.

Leaf Report
Response Prob

Leaf Label	Functional	.2 .4 .6 .8	Non-functional	.2 .4 .6 .8
&crucialness<0.0912698413&crucialness>=0.0200892857&local_population_1km recoded(1 - 1000)&is_urban(True)^&clean_adm1(Ebonyi, ...	0.0174		0.9826	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)^&water_source_clean recoded(Borehole)&clean_ad...	0.0267		0.9733	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)^&water_source_clean recoded(Borehole)&clean_ad...	0.0346		0.9654	
&crucialness>=0.0912698413&subjective_quality(No because of Odour, No because of Taste, No because of Colour)^&pay recoded(No)&...	0.0739		0.9261	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)&clean_adm1(Abia, Kogi, Ondo, Gombe, Kebbi, Ogu...	0.0973		0.9027	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)&clean_adm1(Edo, Oyo, Benue, Anambra, Ekiti, Osu...	0.0995		0.9005	
clean_adm1(Akwa Ibom, Abia, Kogi, Ondo, Cross River, Bayelsa, Delta, Imo, Ogun, Kebbi, Ekiti, Edo, Anambra, Gombe, Benue, Sokoto, Plateau, ...	0.1220		0.8780	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)^&water_source_clean recoded(Borehole)&clean_ad...	0.1412		0.8588	
&crucialness>=0.0912698413&subjective_quality(No because of Odour, No because of Taste, No because of Colour)^&pay recoded(No)&...	0.1992		0.8008	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)^&water_source_clean recoded(Borehole)&local_po...	0.2026		0.7974	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)&clean_adm1(Edo, Oyo, Benue, Anambra, Ekiti, Osu...	0.2103		0.7897	
&crucialness>=0.0912698413&subjective_quality(No because of Odour, No because of Taste, No because of Colour)&clean_adm1(Akwa Ibo...	0.2368		0.7632	
clean_adm1(Akwa Ibom, Abia, Kogi, Ondo, Cross River, Bayelsa, Delta, Imo, Ogun, Kebbi, Ekiti, Edo, Anambra, Gombe, Benue, Sokoto, Plateau, ...	0.2595		0.7405	
&crucialness>=0.0912698413&subjective_quality(No because of Odour, No because of Taste, No because of Colour)&clean_adm1(Akwa Ibo...	0.2610		0.7390	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)&clean_adm1(Abia, Kogi, Ondo, Gombe, Kebbi, Ogu...	0.2730		0.7270	
&crucialness>=0.0912698413&subjective_quality(No because of Odour, No because of Taste, No because of Colour)&clean_adm1(Akwa Ibo...	0.2731		0.7269	
&crucialness<0.0912698413&crucialness>=0.0200892857&local_population_1km recoded(1 - 1000)&is_urban(True)&clean_adm1(Abia, Keb...	0.2745		0.7255	
&crucialness<0.0912698413&crucialness>=0.0200892857&local_population_1km recoded(1 - 1000)&is_urban(True)^&clean_adm1(Plateau, ...	0.2767		0.7233	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)&clean_adm1(Abia, Kogi, Ondo, Gombe, Kebbi, Ogu...	0.2939		0.7061	
&crucialness>=0.0314220374&subjective_quality(Acceptable quality)^&water_tech_clean recoded(Hand Pump)&served_population recode...	0.3105		0.6895	
&crucialness>=0.0912698413&subjective_quality(No because of Odour, No because of Taste, No because of Colour)&clean_adm1(Akwa Ibo...	0.3123		0.6877	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)^&water_source_clean recoded(Borehole)&clean_ad...	0.3412		0.6588	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)&clean_adm1(Abia, Kogi, Ondo, Gombe, Kebbi, Ogu...	0.3521		0.6479	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)^&water_source_clean recoded(Borehole)&local_po...	0.3616		0.6384	
&crucialness>=0.0912698413&subjective_quality(Acceptable quality)&pay recoded(No)^&water_source_clean recoded(Other)&water_tech...	0.3690		0.6310	
&crucialness>=0.0912698413&subjective_quality(No because of Odour, No because of Taste, No because of Colour)&clean_adm1(Edo, Kwa...	0.3817		0.6183	

Figure 21 Decision Tree Leaf Report

4.6. Bootstrap Forest

A bootstrap forest was carried out with the results as follows. The resulting model is a decision tree with >600 splits. The model has a misclassification rate of 0.25, a non-functional prediction rate of 0.59, and a functional prediction rate of 0.84. The false positive error is 0.4 and false negative error is 0.15.

Overall Statistics

Measure	Training	Validation	Test	Definition
Entropy RSquare	0.2444	0.2005	0.1999	1-Loglike(model)/Loglike(0)
Generalized RSquare	0.3772	0.3182	0.3174	$(1 - (L(0)/L(model))^{(2/n)}) / (1 - (L(0))^{(2/n)})$
Mean -Log p	0.5021	0.5313	0.5316	$\sum -\log(p_{ij})/n$
RASE	0.4089	0.4237	0.4235	$\sqrt{\sum (y_{ij} - p_{ij})^2/n}$
Mean Abs Dev	0.3573	0.3699	0.3691	$\sum y_{ij} - p_{ij} /n$
Misclassification Rate	0.2471	0.2739	0.2738	$\sum (p_{ij} \neq p_{Max})/n$
N	32744	24558	24558	n

Confusion Matrix

Training			Validation			Test		
Actual	Predicted Count		Actual	Predicted Count		Actual	Predicted Count	
STATUS recoded	Functional	Non-functional	STATUS recoded	Functional	Non-functional	STATUS recoded	Functional	Non-functional
Functional	17206	3066	Functional	12594	2609	Functional	12671	2533
Non-functional	5025	7447	Non-functional	4117	5238	Non-functional	4192	5162
Actual	Predicted Rate		Actual	Predicted Rate		Actual	Predicted Rate	
STATUS recoded	Functional	Non-functional	STATUS recoded	Functional	Non-functional	STATUS recoded	Functional	Non-functional
Functional	0.849	0.151	Functional	0.828	0.172	Functional	0.833	0.167
Non-functional	0.403	0.597	Non-functional	0.440	0.560	Non-functional	0.448	0.552

Figure 22 Bootstrap Forest Model Results

From the prediction profiler, water systems with low crucialness, high local population, require payment and other water sources have higher probability of non-functionality. Low crucialness may mean lower priority for maintenance or rehabilitation. High local population may translate to higher usage, and if facility is not designed for it, may lead to insufficiency. As for water sources, it may be because wells and springs may require higher maintenance than boreholes.



Figure 23 Bootstrap Model Prediction Profiler

4.7. Model Comparison

The table below compares the results from the different models generated. The decision tree and bootstrap forest models had lower misclassification rates between 0.25 – 0.30, higher prediction of true positives of 0.55 – 0.62 and higher area under ROC curves of 0.77 to 0.83. The decision tree, however, has higher Type I error of 0.23.

Table 1 Comparison of results from all models

Model		Misclassification Rate	Positive: Non-Functional		Negative: Functional		AUC
			True Positive	False Positive	True Negative	False Negative	
Nominal Logistic Regression	TR	0.31	0.48	0.17	0.83	0.52	0.74
	VA	0.30	0.47	0.17	0.83	0.53	0.74
	TE	0.30	0.47	0.17	0.83	0.53	0.74
Stepwise (Forward) Log. Reg.	TR	0.31	0.47	0.17	0.83	0.53	0.74
	VA	0.31	0.47	0.17	0.83	0.53	0.74
	TE	0.31	0.47	0.17	0.83	0.53	0.74
Stepwise (Backward) Log. Reg.	TR	0.31	0.47	0.17	0.83	0.53	0.74
	VA	0.31	0.47	0.17	0.83	0.53	0.74
	TE	0.31	0.47	0.17	0.83	0.53	0.74
Decision Tree	TR	0.29	0.62	0.23	0.77	0.38	0.77
	VA	0.30	0.61	0.23	0.77	0.39	0.76
	TE	0.29	0.60	0.23	0.77	0.40	0.76
Bootstrap Forest	TR	0.25	0.60	0.15	0.85	0.40	0.83
	VA	0.27	0.56	0.17	0.83	0.44	0.79
	TE	0.27	0.55	0.17	0.83	0.45	0.79

In the model comparison, the profit matrix was specified by updating the probability threshold for non-functional levels to 0.38, estimated from the percentage of non-functionality in the dataset.

Profit/Cost Decision Matrix

Specify Profit Matrix

Each matrix entry is the profit if you predict the response in the column when the response in the row is the actual response.

Enter values that reflect profits for correct decisions on the diagonal.
Enter values (usually negative) that reflect profits for incorrect decisions elsewhere.
Use the Undecided column to reflect profits for an alternative decision.

When you save prediction formulas, these values will be used to create best decision columns.
The best decision is the one with greatest expected profit.

Decision or Prediction

	Functional	Non-functional	Undecided
Actual Functional	0	-0.6129	.
Actual Non-functional	-1	0	.

To create a profit matrix for a binary response, enter a Target and Probability Threshold.
If the predicted probability exceeds the threshold, the best decision will be the target.

Target Level: Non-functional

Probability Threshold: 0.38 Set

☐ Save to column as property.

OK Cancel

Figure 24 Probability Threshold Setting

After updating the profit matrix, the prediction for non-functionality improves for all the models. However, the Type I and Type II errors also increases. In general, based on the performance metrics, the decision tree and bootstrap forest models have higher sensitivity and accuracy than the logistic regression models. Some of the specificity and precision are on par.

		Training		Validation		Testing	
		Rates		Rates		Rates	
Method	Actual	Non-functional	Functional	Non-functional	Functional	Non-functional	Functional
Predictor Fit Nominal Logistic	Non-functional	0.705	0.295	0.698	0.302	0.703	0.297
	Functional	0.343	0.657	0.339	0.661	0.335	0.665
Stepwise (Forward)	Non-functional	0.699	0.301	0.694	0.306	0.701	0.299
	Functional	0.343	0.657	0.342	0.658	0.335	0.665
Stepwise (Backward)	Non-functional	0.702	0.298	0.699	0.301	0.705	0.295
	Functional	0.344	0.656	0.341	0.659	0.334	0.666
Predictor Partition	Non-functional	0.754	0.246	0.742	0.258	0.741	0.259
	Functional	0.338	0.662	0.341	0.659	0.336	0.664
Predictor Bootstrap Forest	Non-functional	0.790	0.210	0.757	0.243	0.761	0.239
	Functional	0.308	0.692	0.327	0.673	0.321	0.679

	Method	TP	FN	TN	FP	Sensitivity	Specificity	Precision	Accuracy	F1	MCC	Profit
Training	Predictor Fit Nominal Logistic	8790	3682	13328	6944	0.7048	0.6575	0.5587	0.6755	0.6233	0.3521	-0.242
	Predictor Stepwise (Forward)	8715	3757	13325	6947	0.6988	0.6573	0.5564	0.6731	0.6195	0.3462	-0.245
	Predictor Stepwise (Backward)	8752	3720	13307	6965	0.7017	0.6564	0.5568	0.6737	0.621	0.3481	-0.244
	Predictor Partition	9404	3068	13416	6856	0.7540	0.6618	0.5784	0.6969	0.6546	0.4038	-0.222
	Predictor Bootstrap Forest	9857	2615	14025	6247	0.7903	0.6918	0.6121	0.7294	0.6899	0.4684	-0.197
Validation	Predictor Fit Nominal Logistic	6534	2821	10049	5154	0.6985	0.6610	0.559	0.6753	0.621	0.3495	-0.244
	Predictor Stepwise (Forward)	6496	2859	10011	5192	0.6944	0.6585	0.5558	0.6722	0.6174	0.3431	-0.246
	Predictor Stepwise (Backward)	6537	2818	10015	5188	0.6988	0.6588	0.5575	0.674	0.6202	0.3476	-0.244
	Predictor Partition	6943	2412	10014	5189	0.7422	0.6587	0.5723	0.6905	0.6463	0.3894	-0.228
	Predictor Bootstrap Forest	7083	2272	10225	4978	0.7571	0.6726	0.5873	0.7048	0.6615	0.4174	-0.217
Testing	Predictor Fit Nominal Logistic	6576	2778	10107	5097	0.7030	0.6648	0.5634	0.6793	0.6255	0.3576	-0.24
	Predictor Stepwise (Forward)	6556	2798	10113	5091	0.7009	0.6652	0.5629	0.6788	0.6244	0.356	-0.241
	Predictor Stepwise (Backward)	6592	2762	10120	5084	0.7047	0.6656	0.5646	0.6805	0.6269	0.3601	-0.239
	Predictor Partition	6928	2426	10091	5113	0.7406	0.6637	0.5754	0.693	0.6476	0.3928	-0.226
	Predictor Bootstrap Forest	7122	2232	10328	4876	0.7614	0.6793	0.5936	0.7106	0.6671	0.4281	-0.213

Figure 25 Results after modifying probability threshold

Between the decision tree and bootstrap forest, the latter seems to outperform the decision tree in all aspects. However, the bootstrap forest's non-functionality prediction for the training set, 0.79, is much higher than the validation and training sets of 0.76. This may be a sign of overfitting in the bootstrap forest model.

The variables that contribute to the model predictions are shown in the table below. Most of the models show that a low local population, free water system, higher

crucialness, urban areas and various administrative areas (Abia, Akwa Ibom, Cross River, Gombe, Kogi, Ondo, Teraba) tend to have higher contribution to non-functionality. The bootstrap forest model tends to contradict the parameters for the other models.

Higher crucialness and urban areas tend to signify higher usage of water. Low local population may also suggest lower population density, and hence the water system may be more spread out. This may result in stagnation in certain parts of the system, resulting in higher risk of bacterial growth. A water system with no payment required also means funds may need to be acquired elsewhere for maintenance of the system, which may be more challenging than a system that requires payment.

	Variable that contributes to non-functionality prediction				
Variables	Log. Reg.	Stepwise (Forward)	Stepwise (Backward)	Decision Tree	Bootstrap Forest
Subjective quality of water	No	No		Mixed	
Local population within 1km	<10000	<10000	<10000	<1000	>10000
Payment required	No	No		No	Yes
Crucialness of system	Higher	Higher		Higher	Lower
Population served	<10000	<10000			
Water sources	Borehole	Borehole		Borehole	Others
Distance to town		Higher			
Urban area		True	True	True	
Administrative Area	Various	Various	Various	Various	Various

Comparing complexity, the decision tree model is less complex with 101 splits, compared to the bootstrap forest with over 600 splits. Hence, balancing performance metrics, complexity, and making sense of variables, the decision tree model seems to be the most suitable out of the 5 models.

5. Conclusion

In conclusion, the analysis has churned out several models to predict the non-functionality of the water systems in Nigeria. Out of all the models, the Decision Tree was deemed best in terms of performance, complexity, and variable sensibility.

The variables that affect the model are the local population within 1km, payment system, crucialness of system, water sources, whether the system is in an urban area and various administrative area.

6. Recommendation

Based on the analysis of the models, the following recommendations are generated:

- a. Priority to be put into urban areas with low population but high crucialness of system

Based on the model, these areas seem to signal higher non-functionality in water systems. By prioritising the rehabilitation or maintenance of these systems, any reliability issues can be addressed in a timely manner. Hence, ensuring that clean and safe water can be delivered in a reliable and sustainable manner throughout the country.

- b. Suggest introducing fair pricing for usage of water system

A water system that does not require pricing may result in misuse of systems, and hence increase its risk of non-functionality. However, a fair pricing needs to be introduced and govern to ensure that the supplier's do not exploit the market for an essential commodity.

- c. Subjective quality of water needs to be accompanied by scientific quantitative analysis

The subjective quality of water alone may be insufficient to determine the usability of the water as it only covers colour, odour, and taste. Other components that affect the quality of the water such as fecal matter, pH, microorganism and bacterial growth cannot be assessed through this.

- d. Further investigation into administrative areas required

Various administrative areas (Abia, Akwa Ibom, Cross River, Gombe, Kogi, Ondo, Teraba) seem to have higher contribution to non-functionality. Further analysis may need to be carried out for these areas to identify the root cause (such as low funding etc.) and apply the corrective measures to ensure that the area is accessible to safe water.

7. Appendix

7.1. Appendix A: Fields in Original Dataset

Name	Label	Role	Level
row_id	Row number	ID	Continuous
source	Name of organization collecting the data record	Text	Nominal
lat_deg	Latitude of Water system	Input	Continuous
lon_deg	Longitude of Water system	Input	Continuous
report_date	Date that data was collected	Input	Continuous
status_id	Identify if any water is available on the day of visit	Text	Nominal
water_source_clean	Describe the water source	Text	Nominal
water_source_category	Describe the water source	Text	Nominal
water_tech_clean	Describe the system being used to transport the water	Text	Nominal
water_tech_category	Describe the system being used to transport the water	Text	Nominal
facility_type	Identify if the system is improved	Text	Nominal
clean_country_name	Name of country	Text	Nominal
clean_adm1	Name of primary administrative division	Text	Nominal
clean_adm2	Name of secondary administrative division	Text	Nominal
clean_adm3	Name of tertiary administrative division	Text	Nominal
install_year	4-digit installation year	Input	Continuous
installer	Name of entity that installed the system	Text	Nominal
rehab_year	4-digit of most recent major rehabilitation	Input	Nominal
rehabilitator	Name of entity that completed the rehabilitation	Text	Nominal
management_clean	Name of entity that manages the water point	Text	Nominal
pay	Payment amount and basis	Text	Nominal
status	Physical condition of water point	Text	Nominal
fecal_coliform_presence	Results of e.coli or water quality test	Input	Nominal
fecal_coliform_value	Results of e.coli or water quality test	Input	Nominal
subjective_quality	Information of the water including taste, appearance, odour	Input	Nominal
activity_id	Unique ID for the specific water point infrastructure	ID	Nominal
scheme_id	Identifier for a small, piped scheme that connects multiple water points	ID	Nominal
wpdx_id	Identifier	ID	Nominal
notes	Additional information	Text	Nominal
orig_lnk	Link to the data record for a specific water point or full data set	URL	Nominal
photo_lnk	URL of a photograph of the water system	URL	Nominal
country_id	Identifier for country name	ID	Nominal
data_lnk	Link to dataset	URL	Nominal
public_data_source	Link to dataset	URL	Nominal
distance_to_primary_road	Distance to primary road	Input	Continuous
distance_to_secondary_road	Distance to secondary road	Input	Continuous
distance_to_tertiary_road	Distance to tertiary road	Input	Continuous
distance_to_city	Distance to city	Input	Continuous
distance_to_town	Distance to town	Input	Continuous
rehab_priority	Priority of rehabilitation	Input	Continuous
served_population	Population served by water system	Input	Continuous
local_population_1km	Local population within 1km of water system	Input	Continuous

crucialness	Crucialness of water point (i.e. are there alternative water points nearby?)	Input	Continuous
pressure	Pressure measurement of water system (i.e. it the waterpoint over or under-utilised?)	Input	Continuous
usage_capacity	Recommended maximum users per water point	Input	Continuous
is_urban	Is in an urban area defined by EU Human Settlement Database	Text	Nominal
latest_record	Whether data contains latest information	Text	Nominal
location_id	Location ID of water system	ID	Continuous
cluster_size	Number of times location appears on dataset	Input	Continuous
lat_deg_original	Latitude of water system	Input	Nominal
lon_deg_original	Longitude of water system	Input	Nominal
water_source	Describe the water source	Text	Nominal
water_tech	Describe the system being used to transport the water	Text	Nominal
country_name	Name of country	Text	Nominal
adm1	Name of primary administrative division	Text	Nominal
adm2	Name of secondary administrative division	Text	Nominal
adm3	Name of tertiary administrative division	Text	Nominal
clean_country_id	Country ID	ID	Nominal
management	Entity managing water system	Text	Nominal
created_timestamp	Data creation timestamp	Input	Continuous
updated_timestamp	Data update timestamp	Input	Continuous
New Georeferenced Column	Georeferenced Latitude and Longitude combined	Text	Nominal
lat_lon_deg	Latitude and Longitude combined	Text	Nominal
count	Number of records of this type	Input	Continuous
converted	List of columns by WPDx from original data source	Text	Nominal

7.2. Appendix B: Data Preparation Change Log

Item	Variable Name	Issue	Action
1	row_id	Data characterised as continuous	Change to nominal
2	usage_capacity	Data characterised as continuous	Change to nominal
3	latest_record	Some data not updated, may not be representative of actual	Hide and exclude 752 rows with 'latest_record' = False
4	row_id	Row number irrelevant to analysis	Hide and exclude variable
5	source	Data source irrelevant to analysis	Hide and exclude variable
6	report_Date	Data source irrelevant to analysis	Hide and exclude variable
7	facility_type	Data source irrelevant to analysis	Hide and exclude variable
8	clean_country_name	Data source irrelevant to analysis	Hide and exclude variable
9	activity_id	Water point ID irrelevant to analysis	Hide and exclude variable
10	scheme_id	Water connection ID irrelevant to analysis	Hide and exclude variable
11	wpx_id	WPDx ID irrelevant to analysis	Hide and exclude variable
12	notes	Additional information irrelevant to analysis	Hide and exclude variable
13	orig_lnk	URL link irrelevant to analysis	Hide and exclude variable
14	photo_lnk	URL link irrelevant to analysis	Hide and exclude variable
15	country_id	Country ID irrelevant to analysis	Hide and exclude variable
16	data_lnk	URL link irrelevant to analysis	Hide and exclude variable
17	public_data_source	URL link irrelevant to analysis	Hide and exclude variable
18	latest_record	Record latestness irrelevant to analysis	Hide and exclude variable
19	location_id	Duplicate information, able to use latitude and longitude data	Hide and exclude variable
20	cluster_size	Number of times record appears irrelevant to analysis	Hide and exclude variable
21	lat_deg_original	Duplicate information	Hide and exclude variable
22	lon_deg_original	Duplicate information	Hide and exclude variable
23	water_source	Duplicate information	Hide and exclude variable
24	water_tech	Duplicate information	Hide and exclude variable
25	country_name	Country name irrelevant to analysis	Hide and exclude variable
26	adm1	Duplicate information, clean data available	Hide and exclude variable
27	adm2	Duplicate information, clean data available	Hide and exclude variable
28	adm3	Duplicate information, clean data available	Hide and exclude variable
29	clean_country_id	Country ID irrelevant to analysis	Hide and exclude variable
30	management	Management entity irrelevant to analysis	Hide and exclude variable
31	created_timestamp	Data creation time irrelevant to analysis	Hide and exclude variable
32	updated_timestamp	Data updated time irrelevant to analysis	Hide and exclude variable
33	New Georeferenced Column	Duplicate information, contains latitude and longitude data	Hide and exclude variable
34	lat_lon_deg	Duplicate information, contains latitude and longitude data	Hide and exclude variable
35	count	Number of records of this type irrelevant to analysis	Hide and exclude variable
36	converted	List of columns converted irrelevant to analysis	Hide and exclude variable
37	install_year	More than 30% of data (39784 rows) is missing	Hide and exclude variable

38	rehab_priority	More than 30% of data (53591 rows) is missing	Hide and exclude variable
39	clean_adm2	Contains too many categories (746)	Hide and exclude variable, use clean_adm1 instead
40	clean_adm3	Data is missing	Hide and exclude variable
41	installer	Data is missing	Hide and exclude variable
42	rehab_year	Data is missing	Hide and exclude variable
43	rehabilitator	Data is missing	Hide and exclude variable
44	management_clean	More than 30% of data (27352 rows) is missing	Hide and exclude variable
45	fecal_coliform_presence	Data is missing	Hide and exclude variable
46	fecal_coliform_value	Data is missing	Hide and exclude variable
47	pay	Contains too many categories (42)	Recode into 'pay recoded' - 2 categories (Yes and No)
48	pay	Duplicate with pay recoded	Hide and exclude variable
49	status	Contains too many categories (881)	Recode into 'STATUS recoded' 2 categories (Functional and Non-Functional)
50	status	Duplicate with STATUS recoded	Hide and exclude variable
51	served_population	Contains 274 rows with missing data, may affect model analysis	Hide and exclude rows with missing data
52	local_population_1km	Contains 274 rows with missing data, may affect model analysis	Hide and exclude rows with missing data
53	crucialness	Contains 5616 rows with missing data, may affect model analysis	Hide and exclude rows with missing data
54	pressure	Contains 6443 rows with missing data, may affect model analysis	Hide and exclude rows with missing data
55	water_tech_clean	Contains 1224 rows with missing data, may affect model analysis	Hide and exclude rows with missing data
56	water_tech_category	Contains 1224 rows with missing data, may affect model analysis	Hide and exclude rows with missing data
57	usage_capacity	Contains 892 rows with missing data, may affect model analysis	Hide and exclude rows with missing data
58	served_population	Distribution is uneven with > 99% water points serving < 20000 population	Recode the columns into 'served_population_recoded' 5 categories - '0', '1-1000, 1001-10000, 10001 - 100000 and > 100000)
59	served_population	Duplicate with recoded column	Hide and exclude variable
60	local_population_1km	Distribution is uneven with > 97.5% water points serving < 20000 population	Recode the columns into 'local_population_1km_recoded' 5 categories - '0', '1-1000, 1001-10000, 10001 - 100000 and > 100000)
61	local_population_1km	Duplicate with recoded column	Hide and exclude variable
62	pressure	Correlated with served_population, may serve as redundancy	Hide and exclude variable
63	status_id	Binary predictor unstable in model	Hide and exclude variable
65	water_source_clean	Biasness in model may be due to uneven category distribution	Recode 2 categories - 'Borehole' and 'Others'
65	water_source_clean	Duplicate with water_source_clean	Hide and exclude variable
66	water_source_category	Similar to water_source_clean data, resulting in zeroed error	Hide and exclude variable
67	water_tech_clean	Biasness in model may be due to uneven category distribution	Recode 2 categories - 'Handpump' and 'Others'
68	water_tech_clean	Duplicate with water_tech_clean_recoded	Hide and exclude variable

69	water_tech_category	Similar to water_tech_clean data, resulting in zeroed error	Hide and exclude variable
----	---------------------	--	---------------------------