

1. Overview

Based on the descriptive and exploratory analysis done in notebook 00_data_understanding, this Python Script will work on 2 models: logistic and decision tree classifier, we will chose the best model based on the one that has better evaluation metrics. We will then improve the chosen model with tuned hyperparameters.

2. Data Understanding

2.1 Data Description

This notebook will use the dataset: df_data_processed excel sheet created in the previous notebook: 01_data_preprocessing

2.2 Import Necessary Libraries

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5 import seaborn as sns
6 from sklearn.exceptions import ConvergenceWarning
7
8
9 from sklearn.linear_model import LogisticRegression
10 from sklearn.tree import DecisionTreeClassifier
11 from sklearn.metrics import roc_curve, auc, confusion_matrix
12 from sklearn.model_selection import GridSearchCV
13 from sklearn.metrics import make_scorer, roc_auc_score, recall_score
14
15 import pickle
16 import warnings
17 warnings.filterwarnings('ignore', category=ConvergenceWarning)
18 warnings.simplefilter('ignore')
```

3. Code

3.1 Import the database

```
In [2]: 1 df = pd.read_excel('df_data_processed.xlsx')
        2 df.head()
```

Out[2]:

	amount_tsh	gps_height	population	basin	region	extraction_type_class	payment_type
0	-0.084999	2.053863	-0.041306	-0.540016	-0.633090	-0.521411	-0.897587
1	-0.100621	-0.965049	-0.379739	-0.540016	0.555492	-0.463637	0.771866
2	-0.100621	-0.965049	-0.379739	1.471270	0.131062	2.617222	0.771866
3	-0.100621	-0.965049	-0.379739	-1.053126	0.131062	-0.521411	-1.330306
4	-0.006889	0.511216	-0.125914	0.697368	0.135714	2.617222	-0.641415

5 rows × 36 columns

```
In [3]: 1 df.shape
```

Out[3]: (59400, 36)

3.2 Import the database

```
In [4]: 1 df_train = df[df['is_test']==0]
        2 df_test = df[df['is_test']==1]
```

```
In [5]: 1 y_train = df_train['status_group']
        2 X_train = df_train.drop(['status_group', 'is_test'], axis=1)
        3
        4 y_test = df_test['status_group']
        5 X_test = df_test.drop(['status_group', 'is_test'], axis=1)
```

3.3 Baseline model creations

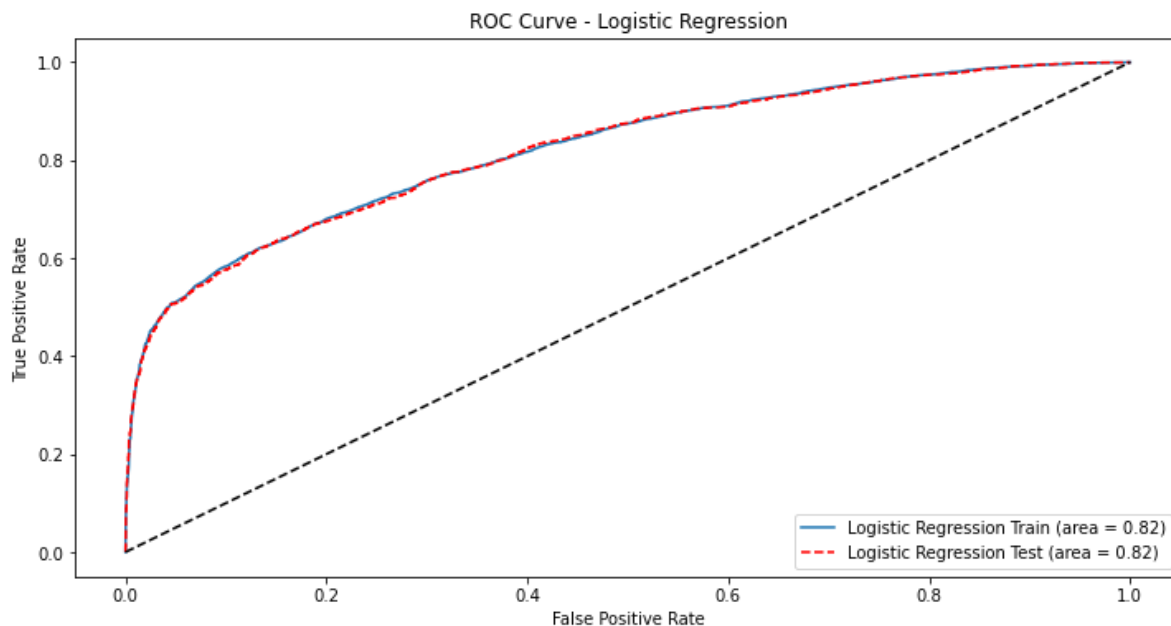
3.3.1 Logistic regression

In [6]:

```

1  # Initialize the Logistic Regression model
2  log_reg = LogisticRegression()
3
4  # Fit the model to the training data
5  log_reg.fit(X_train, y_train)
6
7  # Predict probabilities on the training and test set
8  y_pred_prob_log_reg_train = log_reg.predict_proba(X_train)[:, 1] # Traini
9  y_pred_prob_log_reg_test = log_reg.predict_proba(X_test)[:, 1] # Test pro
10
11 # Compute ROC curve and AUC for training data
12 fpr_log_reg_train, tpr_log_reg_train, _ = roc_curve(y_train, y_pred_prob_1
13 auc_log_reg_train = auc(fpr_log_reg_train, tpr_log_reg_train)
14
15 # Compute ROC curve and AUC for test data
16 fpr_log_reg_test, tpr_log_reg_test, _ = roc_curve(y_test, y_pred_prob_log_
17 auc_log_reg_test = auc(fpr_log_reg_test, tpr_log_reg_test)
18
19 # Plotting ROC Curves
20 plt.figure(figsize=(12, 6))
21 plt.plot(fpr_log_reg_train, tpr_log_reg_train, label='Logistic Regression
22 plt.plot(fpr_log_reg_test, tpr_log_reg_test, color='red', linestyle='--',
23 plt.plot([0, 1], [0, 1], 'k--')
24 plt.xlabel('False Positive Rate')
25 plt.ylabel('True Positive Rate')
26 plt.title('ROC Curve - Logistic Regression')
27 plt.legend(loc="lower right")
28 plt.show()

```



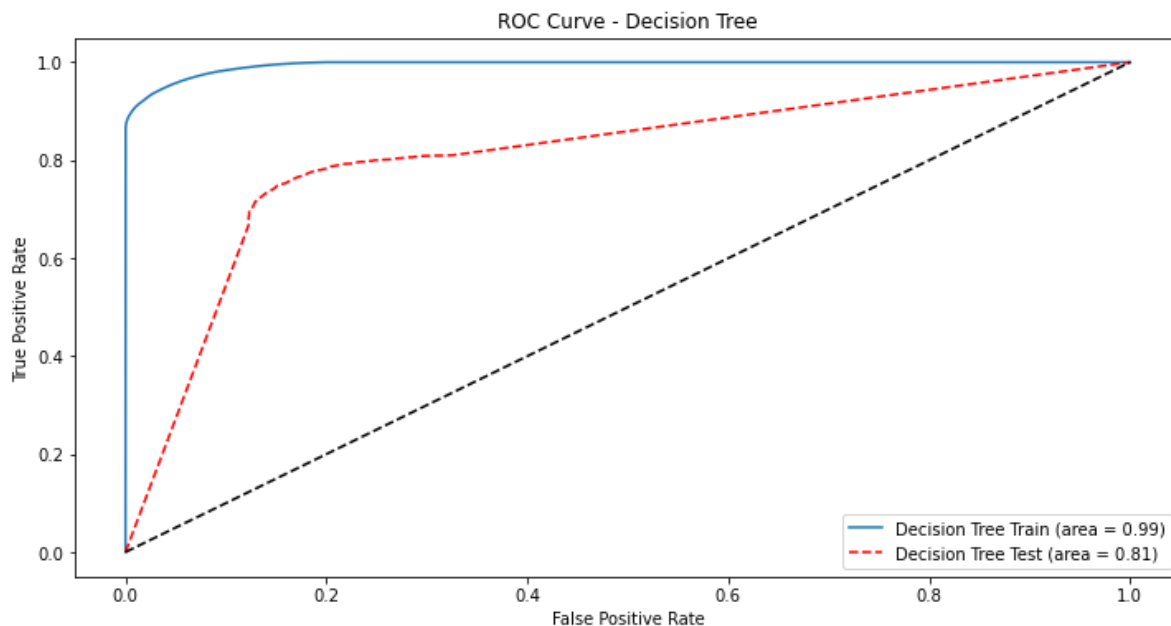
3.3.2 Decision Tree

In [7]:

```

1  # Initialize the Decision Tree model
2  decision_tree = DecisionTreeClassifier()
3
4  # Fit the model to the training data
5  decision_tree.fit(X_train, y_train)
6
7  # Predict probabilities on the training and test set
8  y_pred_prob_tree_train = decision_tree.predict_proba(X_train)[:, 1] # Tr
9  y_pred_prob_tree_test = decision_tree.predict_proba(X_test)[:, 1] # Test
10
11 # Compute ROC curve and AUC for training data
12 fpr_tree_train, tpr_tree_train, _ = roc_curve(y_train, y_pred_prob_tree_tr
13 auc_tree_train = auc(fpr_tree_train, tpr_tree_train)
14
15 # Compute ROC curve and AUC for test data
16 fpr_tree_test, tpr_tree_test, _ = roc_curve(y_test, y_pred_prob_tree_test)
17 auc_tree_test = auc(fpr_tree_test, tpr_tree_test)
18
19 # Plotting ROC Curves
20 plt.figure(figsize=(12, 6))
21 plt.plot(fpr_tree_train, tpr_tree_train, label='Decision Tree Train (area
22 plt.plot(fpr_tree_test, tpr_tree_test, color='red', linestyle='--', label=
23 plt.plot([0, 1], [0, 1], 'k--')
24 plt.xlabel('False Positive Rate')
25 plt.ylabel('True Positive Rate')
26 plt.title('ROC Curve - Decision Tree')
27 plt.legend(loc="lower right")
28 plt.show()

```



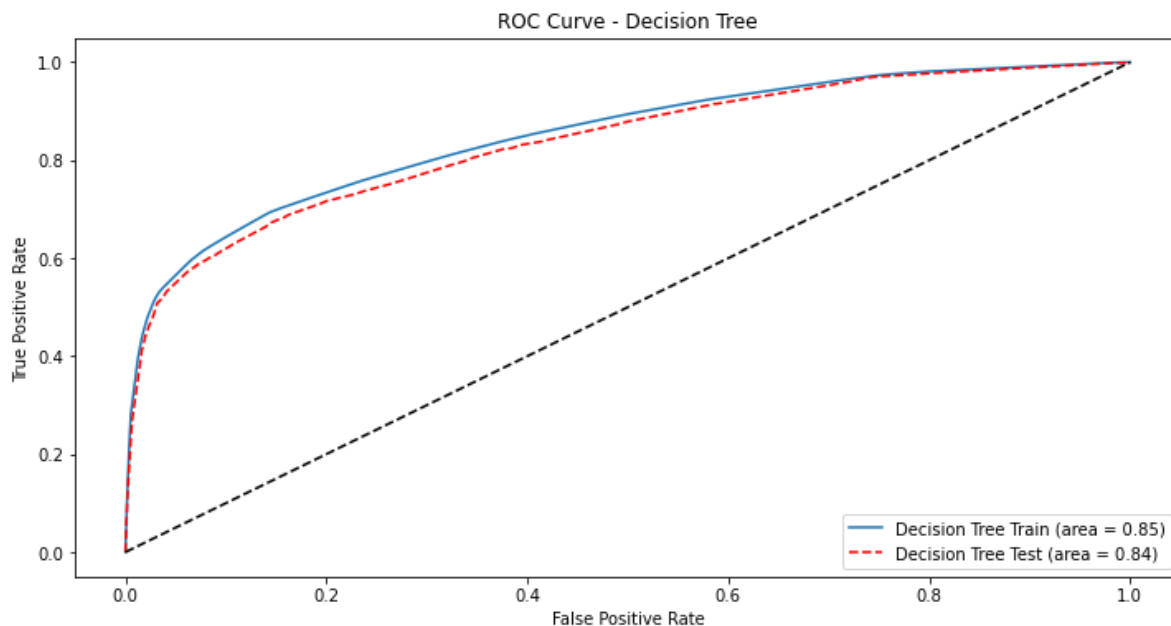
checking max_depth to mitigate overfitting

In [8]:

```

1  # Initialize the Decision Tree model
2  decision_tree = DecisionTreeClassifier(max_depth=7)
3
4  # Fit the model to the training data
5  decision_tree.fit(X_train, y_train)
6
7  # Predict probabilities on the training and test set
8  y_pred_prob_tree_train = decision_tree.predict_proba(X_train)[: , 1] # Train
9  y_pred_prob_tree_test = decision_tree.predict_proba(X_test)[: , 1] # Test
10
11 # Compute ROC curve and AUC for training data
12 fpr_tree_train, tpr_tree_train, _ = roc_curve(y_train, y_pred_prob_tree_train)
13 auc_tree_train = auc(fpr_tree_train, tpr_tree_train)
14
15 # Compute ROC curve and AUC for test data
16 fpr_tree_test, tpr_tree_test, _ = roc_curve(y_test, y_pred_prob_tree_test)
17 auc_tree_test = auc(fpr_tree_test, tpr_tree_test)
18
19 # Plotting ROC Curves
20 plt.figure(figsize=(12, 6))
21 plt.plot(fpr_tree_train, tpr_tree_train, label='Decision Tree Train (area = 0.85)')
22 plt.plot(fpr_tree_test, tpr_tree_test, color='red', linestyle='--', label='Decision Tree Test (area = 0.84)')
23 plt.plot([0, 1], [0, 1], 'k--')
24 plt.xlabel('False Positive Rate')
25 plt.ylabel('True Positive Rate')
26 plt.title('ROC Curve - Decision Tree')
27 plt.legend(loc="lower right")
28 plt.show()

```



3.4 Hyper tuning

3.4.1 Decision Tree Classifier

We are going to do hyper parameter tuning with Decision Tree classifier and the Logistic regression and we will keep the model that gives the best results

The code below is commented as it takes an approximated time of 20 minutes for it to run. However, in the following cell you can see that the best_tree is saved in a pickle

```
In [9]: 1 # Initialize the Decision Tree model
2 decision_tree = DecisionTreeClassifier(class_weight="balanced")
3
4 # Define the parameter grid to search
5 param_grid = {
6     'max_depth': range(8, 13), # Explore depths from 7 to 11
7     'min_samples_split': range(3, 7, 2), # Minimum number of samples required
8     'min_samples_leaf': range(2, 5), # Minimum number of samples required
9     'max_features': ['auto', 'log2', None] # Number of features to consider
10 }
11
12 # Define the scoring function using AUC
13 scorer = make_scorer(recall_score, average='binary')
14
15 # Setup the grid search with cross-validation
16 grid_search = GridSearchCV(estimator=decision_tree, param_grid=param_grid,
17
18 # Fit grid search on the training data
19 grid_search.fit(X_train, y_train)
20
21 # Find the best model
22 best_tree = grid_search.best_estimator_
```

```
In [10]: 1 # Output the best parameter combination and the corresponding score
2 print("Best parameters found:", grid_search.best_params_)
3 print("Best Recall achieved:", grid_search.best_score_)
4
5 # Optional: Evaluate the best model on the test set
6 y_pred_proba_best_tree = best_tree.predict_proba(X_test)[: , 1]
7
8 # Let's apply a threshold to the probabilities of y_pred_proba_best_tree to
9 y_pred_dt = np.where(y_pred_proba_best_tree >= 0.40, 1, 0)
10
11 test_recall = recall_score(y_test, y_pred_dt)
12
13 print("Test Recall of best model:", test_recall)
```

```
Best parameters found: {'max_depth': 12, 'max_features': None, 'min_samples_
eaf': 4, 'min_samples_split': 3}
Best Recall achieved: 0.7383857002960346
Test Recall of best model: 0.7900262467191601
```

```
In [11]: 1 # Save the best_tree in a pickle
2 pickle.dump(best_tree, open(f"model_objects/best_tree.pkl", 'wb'))
```

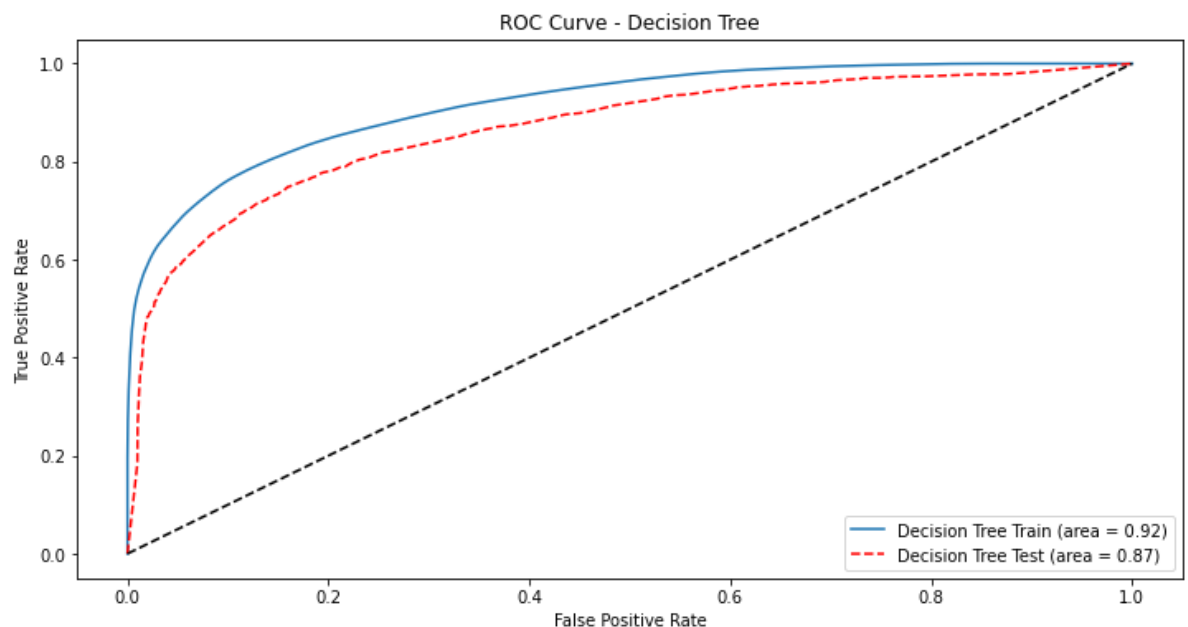
Let's do the curve ROC and see the values AUC with the values for this Decision TreeClassifier

In [12]:

```

1  # Predict probabilities on the training and test set
2  y_pred_prob_tree_train = best_tree.predict_proba(X_train)[: , 1] # Trainir
3  y_pred_proba_best_tree = best_tree.predict_proba(X_test)[: , 1] # Test pro
4
5  # Compute ROC curve and AUC for training data
6  fpr_tree_train, tpr_tree_train, _ = roc_curve(y_train, y_pred_prob_tree_tr
7  auc_tree_train = auc(fpr_tree_train, tpr_tree_train)
8
9  # Compute ROC curve and AUC for test data
10 fpr_tree_test, tpr_tree_test, _ = roc_curve(y_test, y_pred_proba_best_tree
11 auc_tree_test = auc(fpr_tree_test, tpr_tree_test)
12
13 # Plotting ROC Curves
14 plt.figure(figsize=(12, 6))
15 plt.plot(fpr_tree_train, tpr_tree_train, label='Decision Tree Train (area
16 plt.plot(fpr_tree_test, tpr_tree_test, color='red', linestyle='--', label=
17 plt.plot([0, 1], [0, 1], 'k--')
18 plt.xlabel('False Positive Rate')
19 plt.ylabel('True Positive Rate')
20 plt.title('ROC Curve - Decision Tree')
21 plt.legend(loc="lower right")
22 plt.show()

```



3.4.2 logistic regression

We are going to comment the cell below as it takes an approximate time of 20 minutes for it to run.

```
In [13]: 1 # Initialize the Logistic Regression model
2 logistic_regression = LogisticRegression()
3
4 # Define the parameter grid to search
5 param_grid = {
6     'C': [0.01, 0.1, 1, 10], # Inverse of regularization strength
7     'solver': ['newton-cg', 'lbfgs', 'liblinear'], # Algorithm to use in
8     'max_iter': [100, 200], # Maximum number of iterations taken for the
9 }
10
11 # Define the scoring function using AUC
12 scorer = make_scorer(recall_score, average='binary')
13
14 # Setup the grid search with cross-validation
15 grid_search = GridSearchCV(estimator=logistic_regression, param_grid=param_grid)
16
17 # Fit grid search on the training data
18 grid_search.fit(X_train, y_train)
19
20 # Find the best model
21 best_log_reg = grid_search.best_estimator_
22
23
```

```
In [14]: 1 # Output the best parameter combination and the corresponding score
2 print("Best parameters found:", grid_search.best_params_)
3 print("Best Recall achieved:", grid_search.best_score_)
4
5 # Let's apply a threshold to the probabilities of y_pred_prob_log_reg_test
6 y_pred_prob_log_reg_test = best_log_reg.predict_proba(X_test)[:, 1] # Test probabilities
7 y_pred_lr = np.where(y_pred_prob_log_reg_test >= 0.40, 1, 0)
8
9 test_recall_lr = recall_score(y_test, y_pred_lr)
10
11 print("Test Recall of best model:", test_recall_lr)
```

```
Best parameters found: {'C': 0.1, 'max_iter': 100, 'solver': 'newton-cg'}
Best Recall achieved: 0.5513373780124042
Test Recall of best model: 0.6233595800524935
```

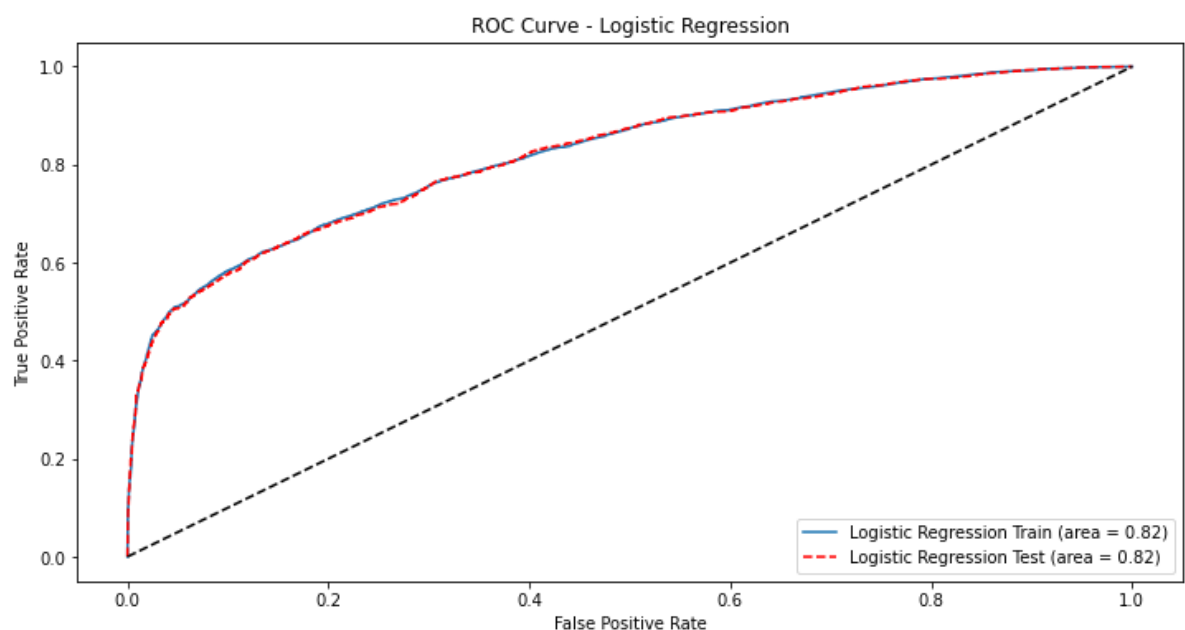
Let's do the curve ROC and see the values AUC with the values for this Logistic Regressor

In [15]:

```

1  # Predict probabilities on the training and test set using the Logistic Re
2  y_pred_prob_log_reg_train = best_log_reg.predict_proba(X_train)[: , 1] # 1
3  y_pred_prob_log_reg_test = best_log_reg.predict_proba(X_test)[: , 1] # Tes
4
5  # Compute ROC curve and AUC for training data
6  fpr_log_reg_train, tpr_log_reg_train, _ = roc_curve(y_train, y_pred_prob_1
7  auc_log_reg_train = auc(fpr_log_reg_train, tpr_log_reg_train)
8
9  # Compute ROC curve and AUC for test data
10 fpr_log_reg_test, tpr_log_reg_test, _ = roc_curve(y_test, y_pred_prob_log_
11 auc_log_reg_test = auc(fpr_log_reg_test, tpr_log_reg_test)
12
13 # Plotting ROC Curves
14 plt.figure(figsize=(12, 6))
15 plt.plot(fpr_log_reg_train, tpr_log_reg_train, label='Logistic Regression
16 plt.plot(fpr_log_reg_test, tpr_log_reg_test, color='red', linestyle='--',
17 plt.plot([0, 1], [0, 1], 'k--')
18 plt.xlabel('False Positive Rate')
19 plt.ylabel('True Positive Rate')
20 plt.title('ROC Curve - Logistic Regression')
21 plt.legend(loc="lower right")
22 plt.show()
23

```

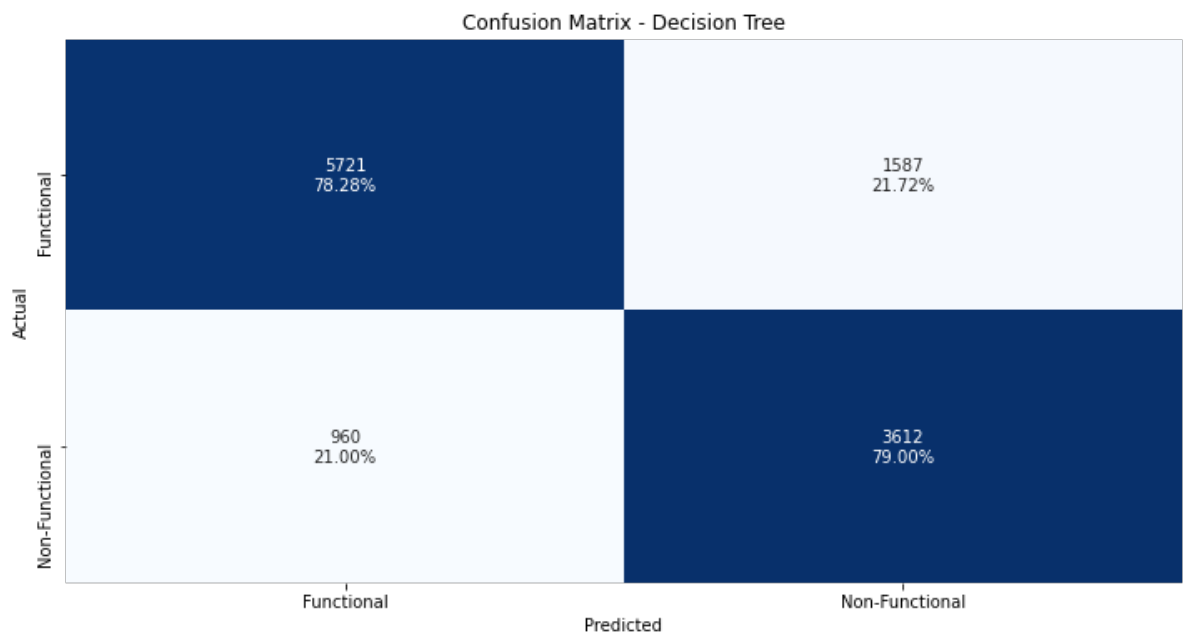


3.5 Confusion matrix

3.5.1 Decision Tree Classifier

```
In [21]: 1 # Confusion Matrix for Decision Tree
          2 cm_tree = confusion_matrix(y_test, y_pred_dt)
```

```
In [22]: 1 # Normalize the confusion matrix by row (actual class)
          2 cm_tree_normalized = cm_tree.astype('float') / cm_tree.sum(axis=1)[:, np.newaxis]
          3
          4 # Create labels for each cell
          5 labels = np.array(["{0}\n{1:.2%}".format(value, percentage) for value, percentage in
          6                     zip(cm_tree_normalized.flatten(), cm_tree_normalized.sum(axis=1).flatten())])
          7
          8 # Plotting the Confusion Matrix for Decision Tree
          9 plt.figure(figsize=(12, 6))
         10 sns.heatmap(cm_tree_normalized, annot=labels, fmt='', cmap='Blues', xticklabels=2, yticklabels=2)
         11 plt.xlabel('Predicted')
         12 plt.ylabel('Actual')
         13 plt.title('Confusion Matrix - Decision Tree')
         14 plt.show()
```



False Negatives (FN): 21.19%

- Impact: A high rate of false negatives means that a significant proportion of the positive class (e.g., non-functional pumps) is being misclassified as negative (e.g., functional pumps). This could lead to serious issues in the business context, as non-functional pumps that are not identified will not receive the necessary maintenance or repairs, leading to prolonged downtimes and possibly affecting the service quality and user satisfaction.
- Business Problem Impact: This could result in increased downtime for the pumps, higher maintenance costs over time, and a negative impact on customer satisfaction due to unreliable water supply.

False Positives (FP): 21.61%

- Impact: A relatively low rate of false positives indicates that only a small proportion of the negative class (e.g., functional pumps) is being misclassified as positive (e.g., non-

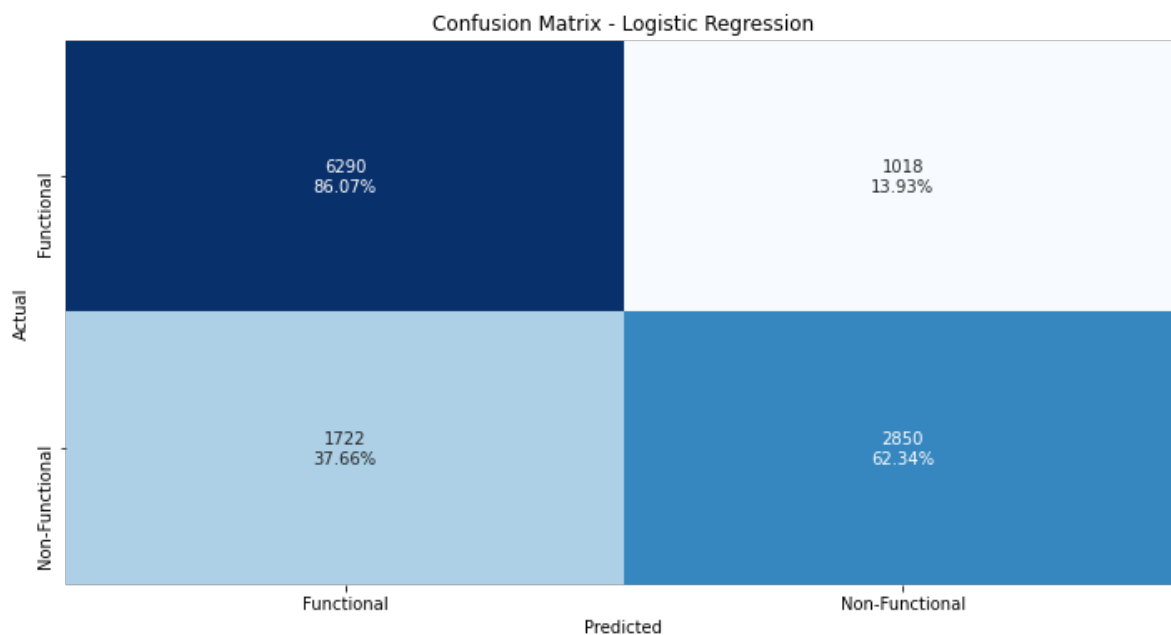
functional pumps). While this is less severe compared to false negatives, it still leads to unnecessary maintenance actions being taken on functional pumps.

- **Business Problem Impact:** This could lead to inefficient allocation of resources, where time and effort are spent on checking or repairing pumps that are actually functional. This can increase operational costs and divert attention from genuinely non-functional pumps that need repairs.

3.5.2 Logistic Regression

```
In [23]: 1 # Confusion Matrix for Logistic Regression
          2 cm_log_reg = confusion_matrix(y_test, y_pred_lr)
```

```
In [24]: 1 # Normalize the confusion matrix by row (actual class)
          2 cm_log_reg_normalized = cm_log_reg.astype('float') / cm_log_reg.sum(axis=1)
          3
          4 # Create labels for each cell
          5 labels = np.array(["{0}\n{1:.2%}".format(value, percentage) for value, percentage in
          6                      zip(cm_log_reg_normalized.flatten(), cm_log_reg_normalized.sum(axis=1).flatten())])
          7
          8 # Plotting the Confusion Matrix for Logistic Regression
          9 plt.figure(figsize=(12, 6))
         10 sns.heatmap(cm_log_reg_normalized, annot=labels, fmt='', cmap='Blues',
         11             xticklabels=['Functional', 'Non-Functional'], yticklabels=['Functional', 'Non-Functional'])
         12 plt.xlabel('Predicted')
         13 plt.ylabel('Actual')
         14 plt.title('Confusion Matrix - Logistic Regression')
         15 plt.show()
```



The logistic regression model has a False negatives of 37.66%, which is greater than the one of the Decision Tree model that has a False negative percentage of 21.19%. This is the most critical metric that we want to ensure is very small because the False Negative percentage

represents the risk of undetected non-functional pumps. Predicting a pump as functional when in reality it turns out to be non-functional could be fatal for certain communities.

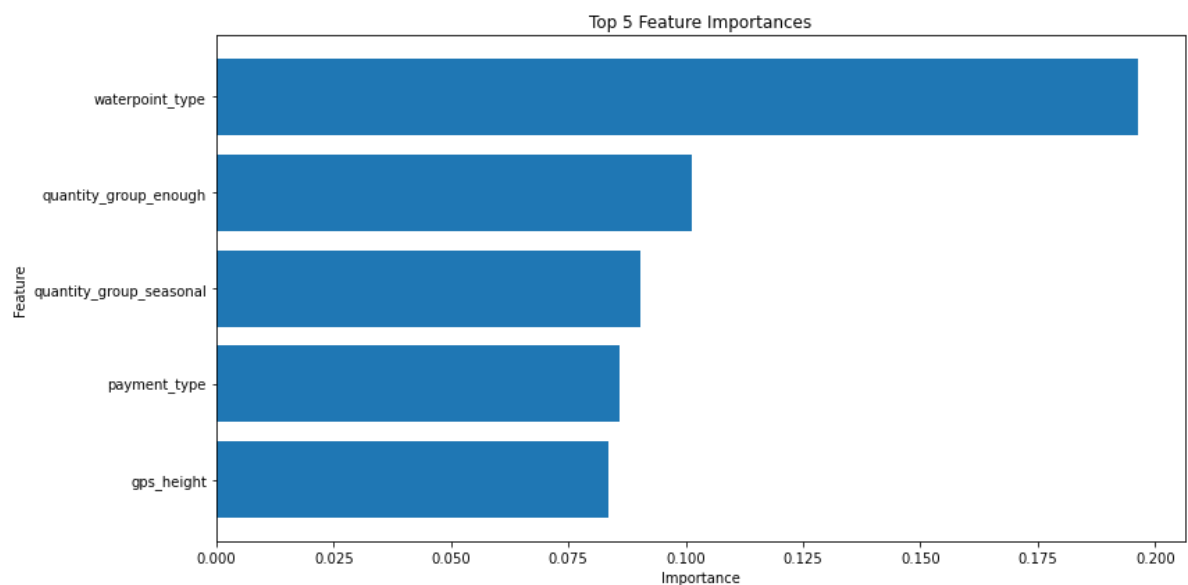
In all, considering that the Logistic Regression model has a higher False negative than the Decision Tree model, we decide to use the Decision Tree Model Classifier

4. Feature importance

We are now going to execute a feature importance code to be able to see the level of importance of all variables when doing the predictions

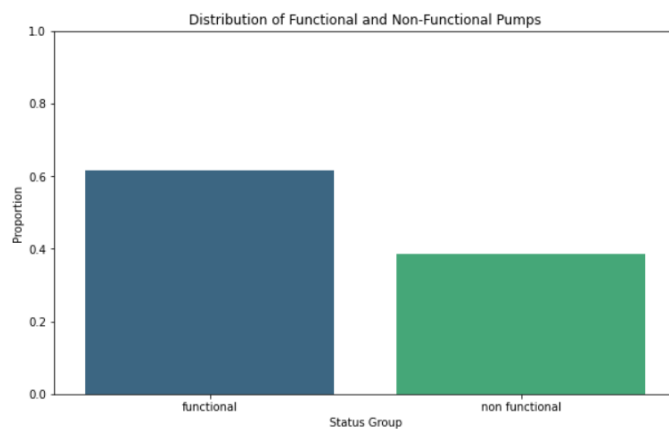
In [25]:

```
1  # Obtain the most important features affecting the status of a pump
2  importances = best_tree.feature_importances_
3
4  # Obtener los nombres de las características
5  feature_names = X_train.columns
6
7  # Create a bar graph for the importance of the characteristics
8  # Order importances in descending order
9  indexes = np.argsort(importances)[::-1]
10
11 # Get the top 5 important features
12 top_indexes = indexes[:5]
13
14 plt.figure(figsize=(12, 6))
15 plt.title("Top 5 Feature Importances")
16 plt.barh(range(5), importances[top_indexes], align="center")
17 plt.yticks(range(5), feature_names[top_indexes])
18 plt.xlabel("Importance")
19 plt.ylabel("Feature")
20 plt.tight_layout()
21 plt.gca().invert_yaxis() # Invert the y-axis to have the most important f
22 plt.show()
```



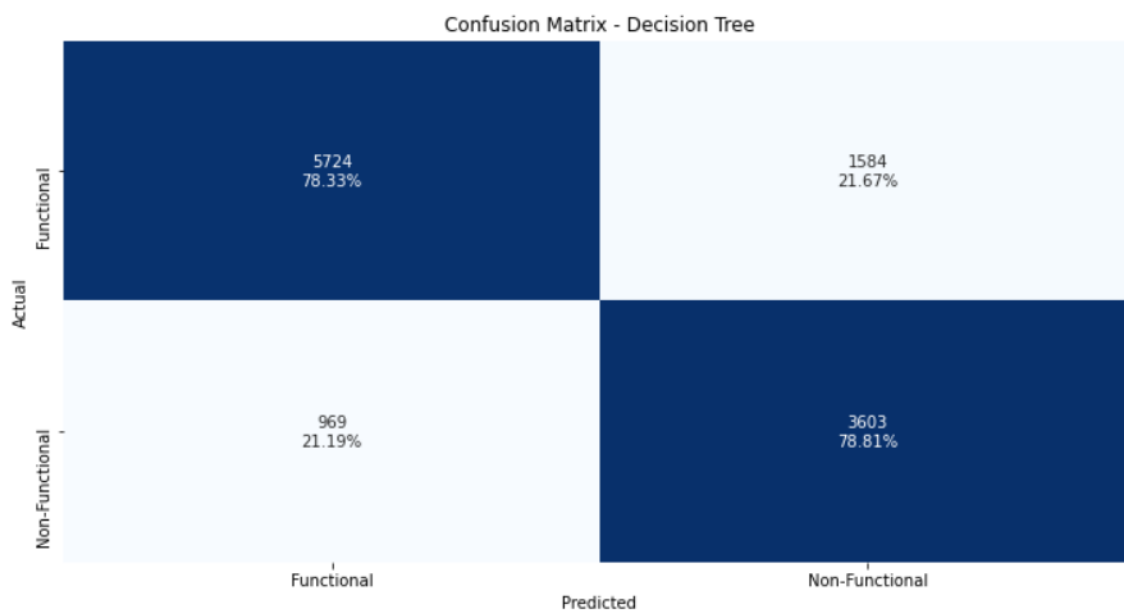
5. Conclusion

Considering the distribution of the dependent variable



As we can see there is not an imbalance problem even though the majority of pumps are functional.

Diving into the model results, let's begin by looking into the confusion matrix



The confusion matrix indicates that the model has a relatively low rate of false negatives (21.19%). The false positive rate (21.67%) is relatively low, meaning fewer resources will be wasted on unnecessary maintenance. However, the primary concern should be reducing the false negative rate to ensure that non-functional pumps are correctly identified and repaired promptly. The result is now 21.19%, which is somewhat low and satisfactory, but further progress should be made to reduce this even further.

Based on the metrics, the best Recall score is obtained with a Decision Tree Classifier. Moreover, the AUC for this model is of 0.87 for the test. In the case of the Logistic Regression model, the recall score was worse even and it had an AUC score (of 0.82).

The variables that are most important and that permit us to best descriminate are:

1. waterpoint_type
2. quantity_group
3. payment_type

We are interested in these 3 variables because they are the ones that have the most influence when determining whether a pump is functional or non-functional.

Here we will show the contingency tables for each variable divided into functional, functional with repairs, and non functional pumps:

1.

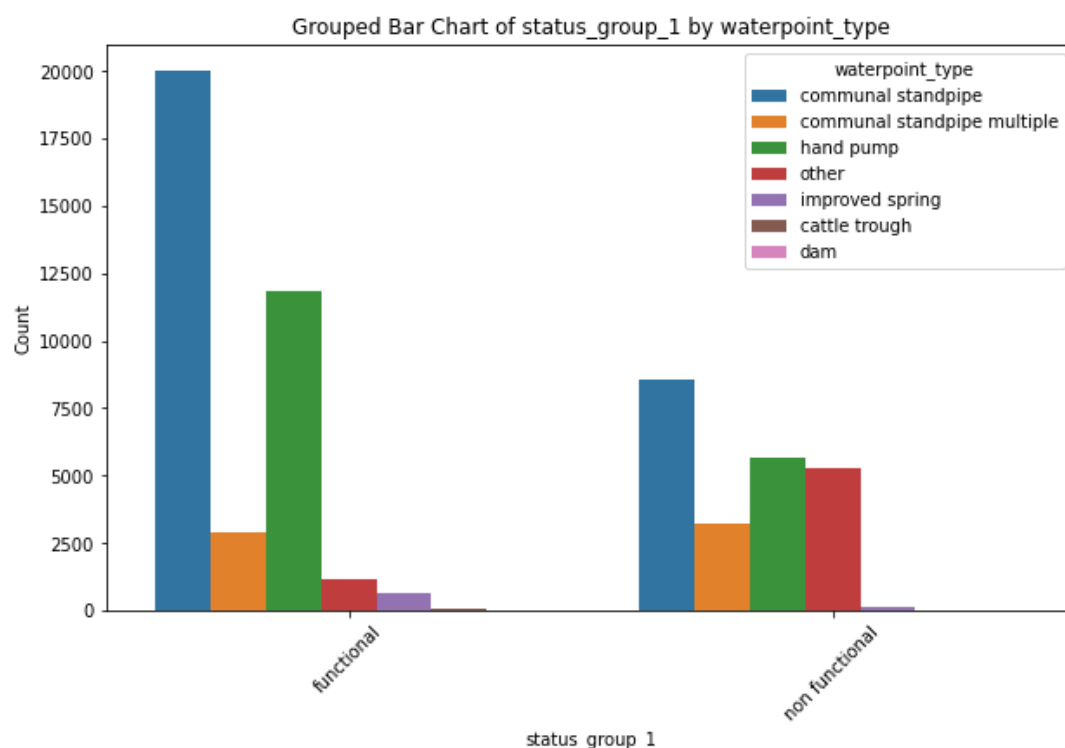
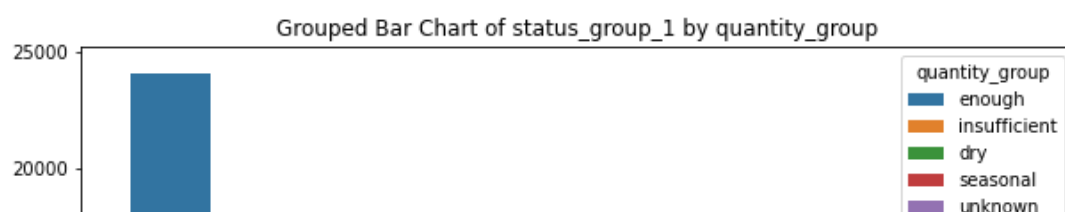


Table for waterpoint_type:

waterpoint_type	cattle trough	communal standpipe	communal standpipe multiple	dam	hand pump	improved spring	other
status_group_1							
functional	0.24%	54.64%	7.88%	0.02%	32.35%	1.77%	3.10%
non functional	0.13%	37.40%	14.11%	0.00%	24.77%	0.60%	22.99%

2.



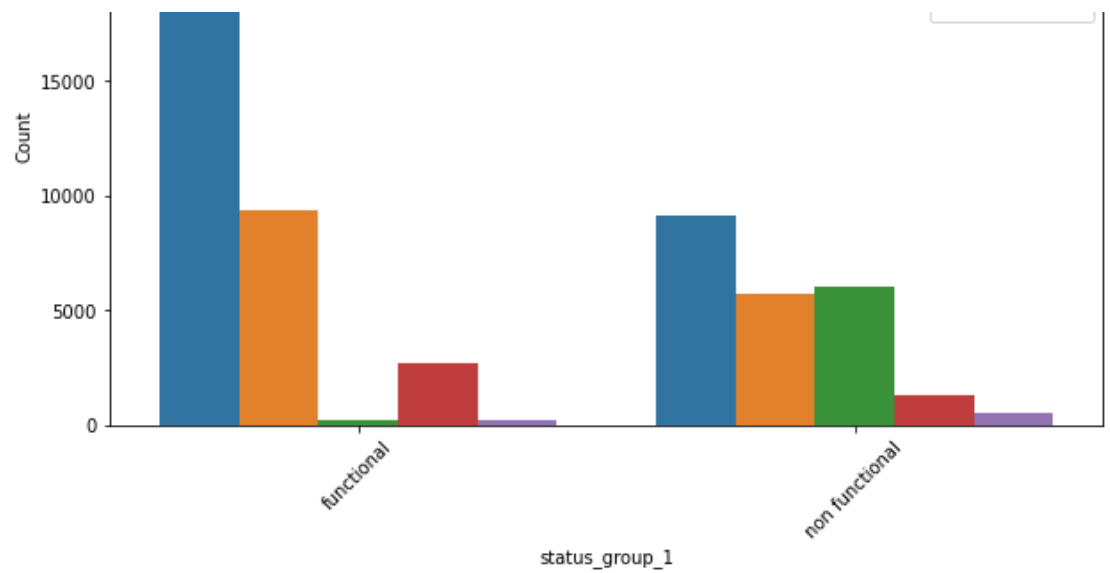


Table for quantity_group:

quantity_group	dry	enough	insufficient	seasonal	unknown
status_group_1					
functional	0.53%	65.75%	25.61%	7.49%	0.62%
non functional	26.52%	40.04%	25.25%	5.74%	2.46%

3.

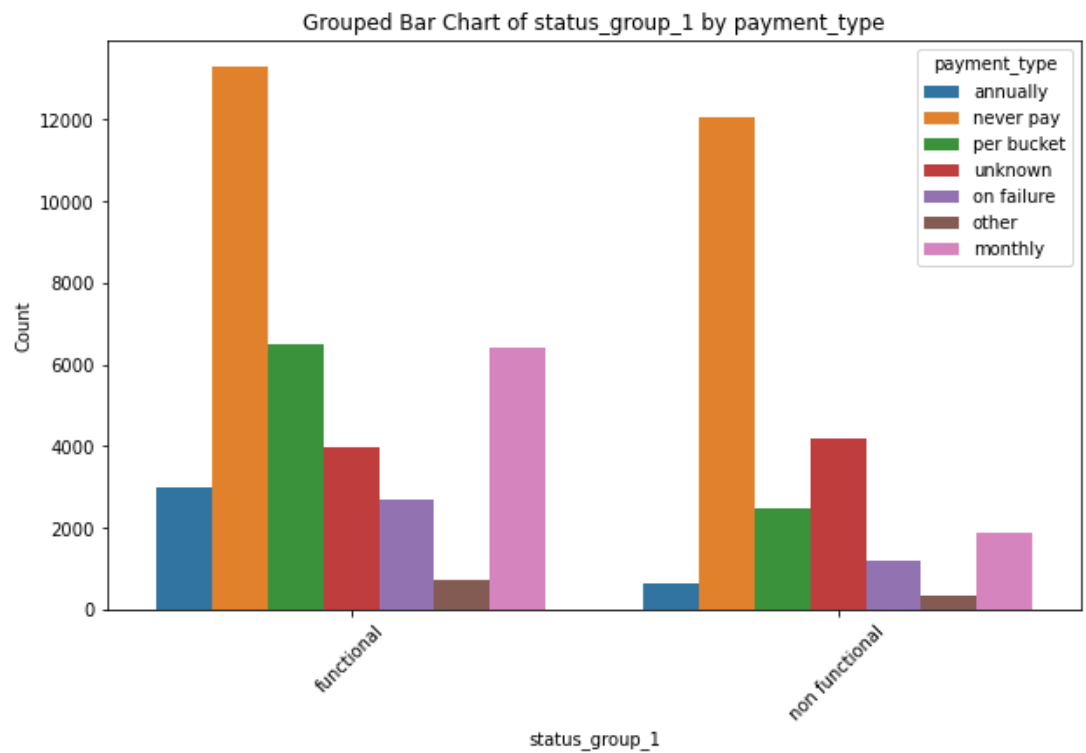


Table for payment_type:

payment_type	annually	monthly	never pay	on failure	other	per bucket	unknown
--------------	----------	---------	-----------	------------	-------	------------	---------

status_group_1							
functional	8.17%	17.52%	36.32%	7.40%	1.99%	17.77%	10.83%
non functional	2.87%	8.29%	52.85%	5.29%	1.42%	10.89%	18.39%

6. Recommendations

1. Considering that most of the functional pumps have monthly payment plans or a per bucket, the Tanzanian government can consider modifying the existing payment plans of those pumps where the payments are different from those payment types, so that the chance of the pump being functional can be increased.
2. Considering that almost none of the functional pumps are dry, it is possible to verify which pumps are dry as a proxy variable to know if they are functional or not and thus focus efforts on repairing them.
3. Considering that non-functional pumps have in most cases a waterpoint_type different from cattle trough, communal standpipe, communal standpipe multiple, dam, hand pump and improved spring, it is possible to verify which pumps do not have these waterpoint_types as a proxy variable to know if they are functional or not and thus focus efforts on repairing them.