# Sidra Bibi, Ariful Islam

# Water Pump
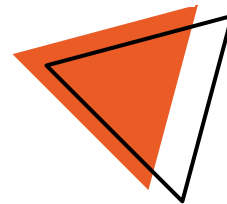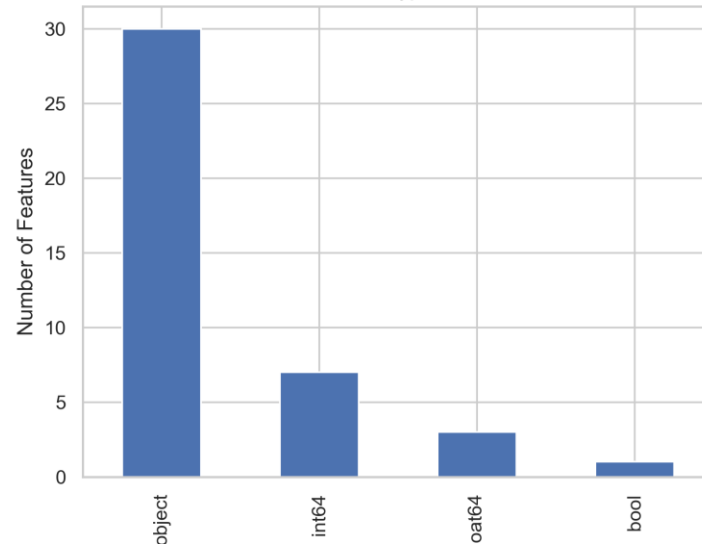
16.4.2025
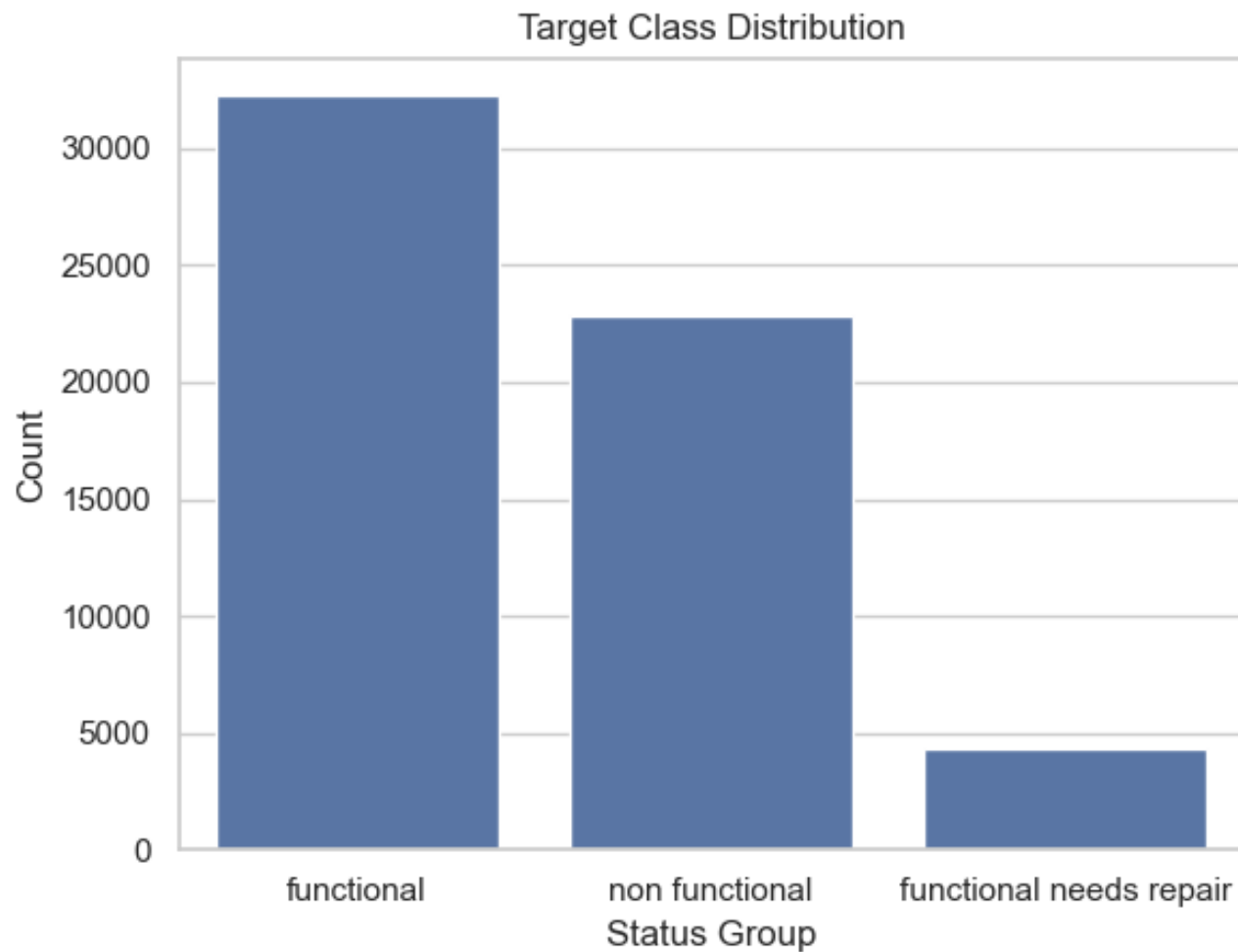
We use tech to connect human potential and opportunity with dignity & humility

# Data Types
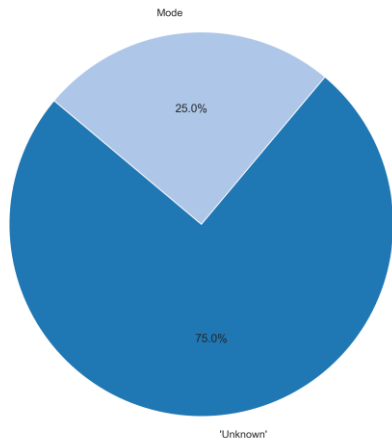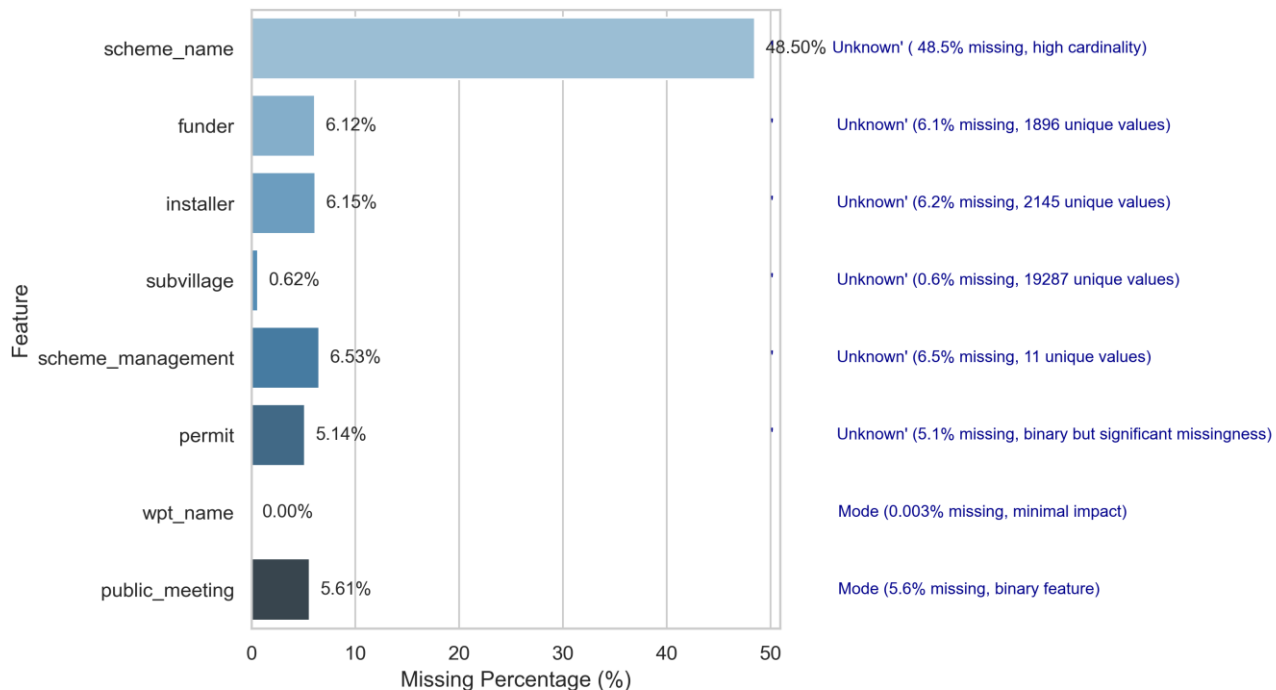


Dataset Dimensions



Count of Data Types in Dataset

# Missing values and imputation

# Proposed Encoding Techniques

| | Feature | Unique Values | Encoding Method |
|---|---|---|---|
| 0 | date_recorded | 356 | Frequency or Target Encoding |
| 1 | funder | 1896 | Frequency or Target Encoding |
| 2 | installer | 2145 | Frequency or Target Encoding |
| 3 | wpt_name | 37399 | Frequency or Target Encoding |
| 4 | basin | 9 | One-Hot Encoding |
| 5 | subvillage | 19288 | Frequency or Target Encoding |
| 6 | region | 21 | Label Encoding |
| 7 | lga | 125 | Frequency or Target Encoding |
| 8 | ward | 2092 | Frequency or Target Encoding |
| 9 | recorded_by | 1 | One-Hot Encoding |
| 10 | scheme_management | 12 | Label Encoding |
| 11 | scheme_name | 2696 | Frequency or Target Encoding |
| 12 | permit | 3 | One-Hot Encoding |
| 13 | extraction_type | 18 | Label Encoding |
| 14 | extraction_type_group | 13 | Label Encoding |
| 15 | extraction_type_class | 7 | One-Hot Encoding |
| 16 | management | 12 | Label Encoding |
| 17 | management_group | 5 | One-Hot Encoding |
| 18 | payment | 7 | One-Hot Encoding |
| 19 | payment_type | 7 | One-Hot Encoding |
| 20 | water_quality | 8 | One-Hot Encoding |
| 21 | quality_group | 6 | One-Hot Encoding |
| 22 | quantity | 5 | One-Hot Encoding |
| 23 | quantity_group | 5 | One-Hot Encoding |
| 24 | source | 10 | One-Hot Encoding |
| 25 | source_type | 7 | One-Hot Encoding |
| 26 | source_class | 3 | One-Hot Encoding |
| 27 | waterpoint_type | 7 | One-Hot Encoding |
| 28 | waterpoint_type_group | 6 | One-Hot Encoding |
| 29 | status_group | 3 | One-Hot Encoding |



Number of Unique Values per Categorical Feature

# Features Extraction from Date_recorded

- ❖ year_recorded
- ❖ month_recored
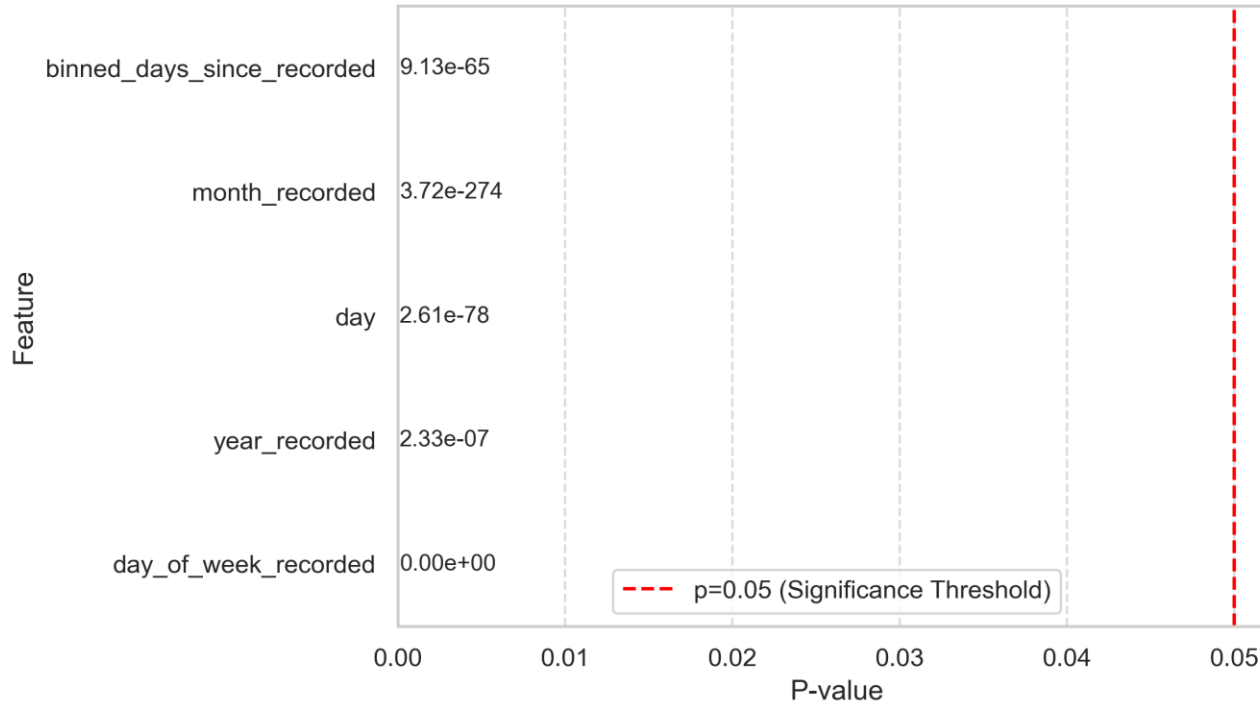- ❖ day
- ❖ day_of_week_recorded
- ❖ days_since_recorded



Correlation Heatmap of Date Features and status_group

# Date_record Features Chi-squared

| Feature | Chi² Statistic | P-value | Degrees of Freedom |
|---|---|---|---|
| binned_days_since_recorded | 1915.23 | 0.000000e+00 | 18 |
| month_recorded | 1359.44 | 3.72e-274 | 22 |
| day | 539.67 | 2.61e-78 | 60 |
| year_recorded | 321.85 | 9.13e-65 | 8 |
| day_of_week_recorded | 54.38 | 2.33e-07 | 12 |

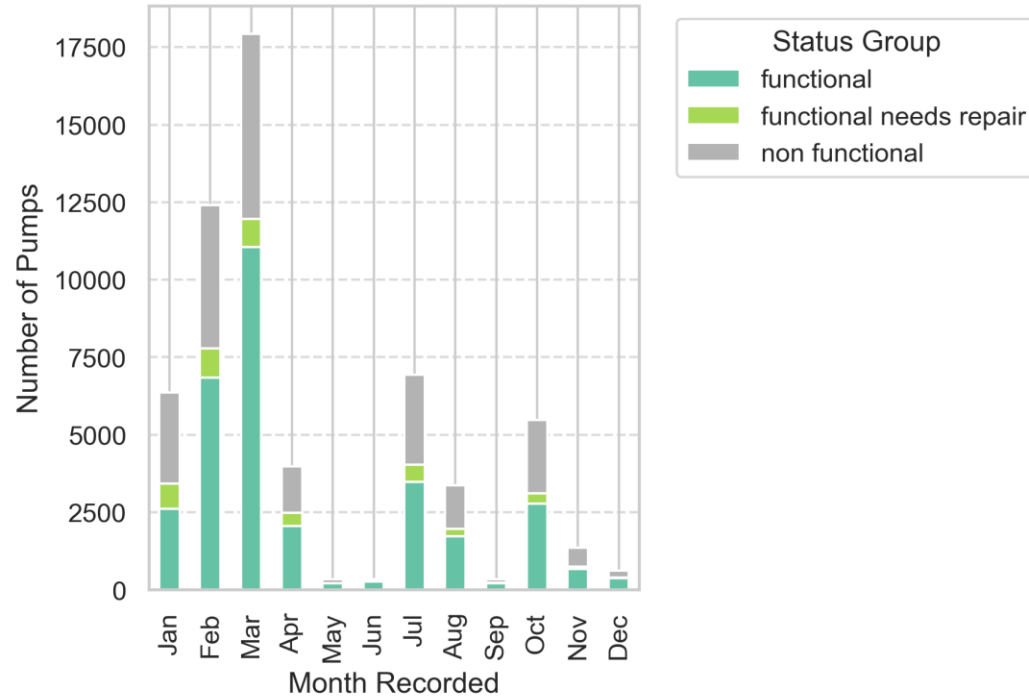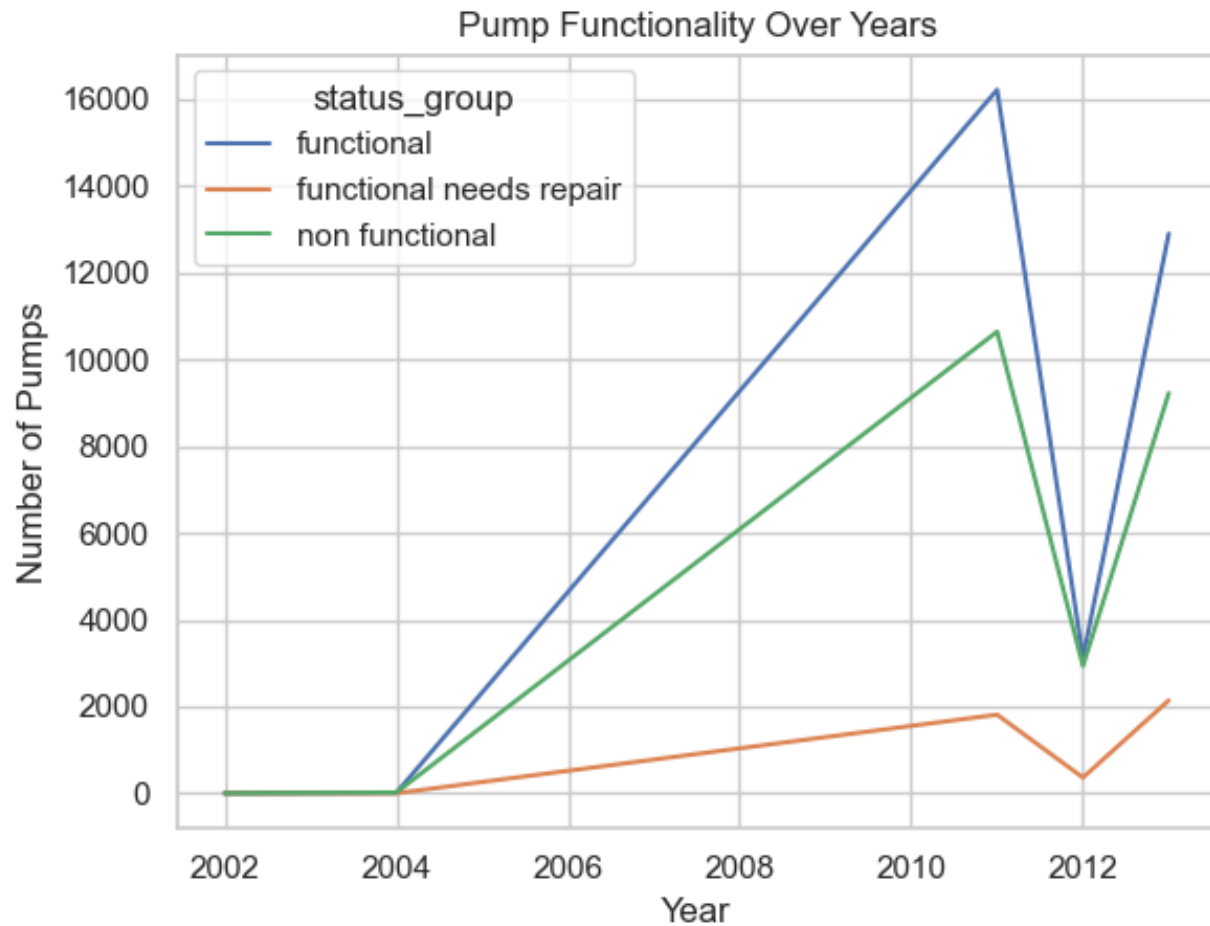# Ch2-P values



Chi-Squared Test P-values: Date Features vs. status_group

# Chi-Squared value

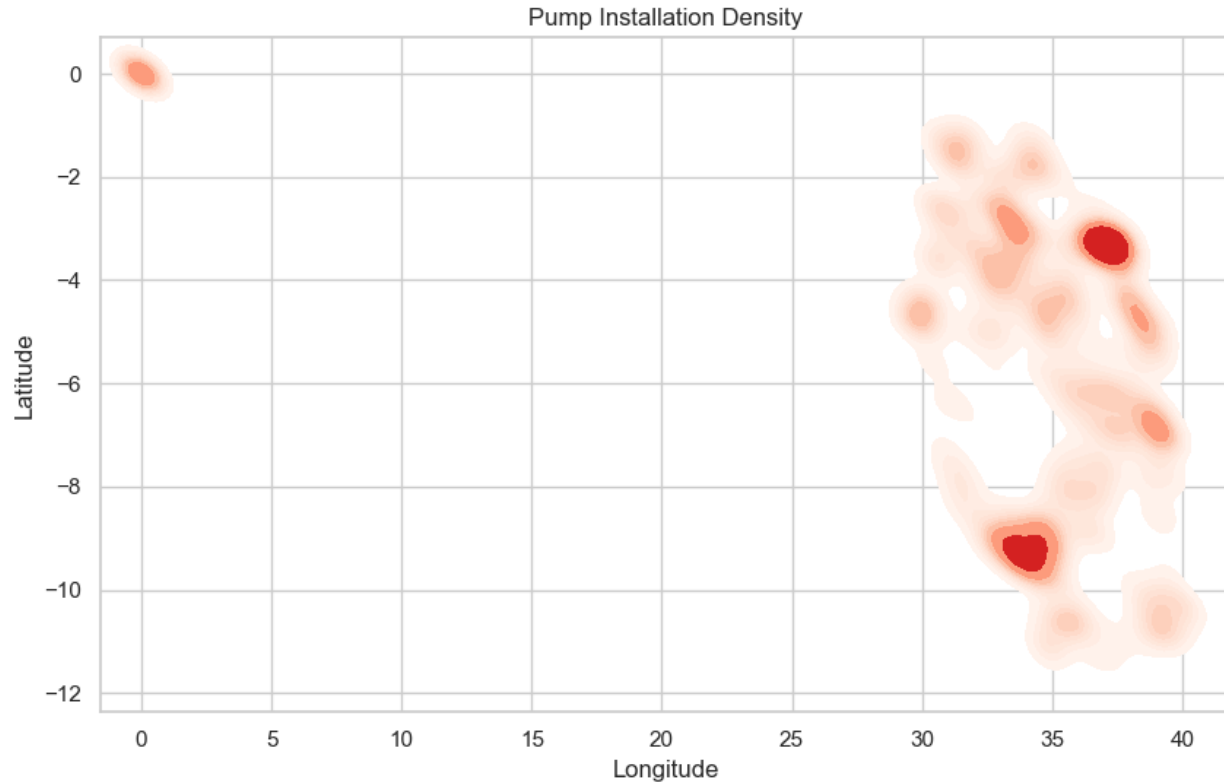| Feature | Chi² Statistic | P-value | Degrees of Freedom |
|---|---|---|---|
| binned_days_since_recorded | 1915.23 | 0.000000e+00 | 18 |
| month_recorded | 1359.44 | 3.72e-274 | 22 |
| day | 539.67 | 2.61e-78 | 60 |
| year_recorded | 321.85 | 9.13e-65 | 8 |
| day_of_week_recorded | 54.38 | 2.33e-07 | 12 |

# Monthly status



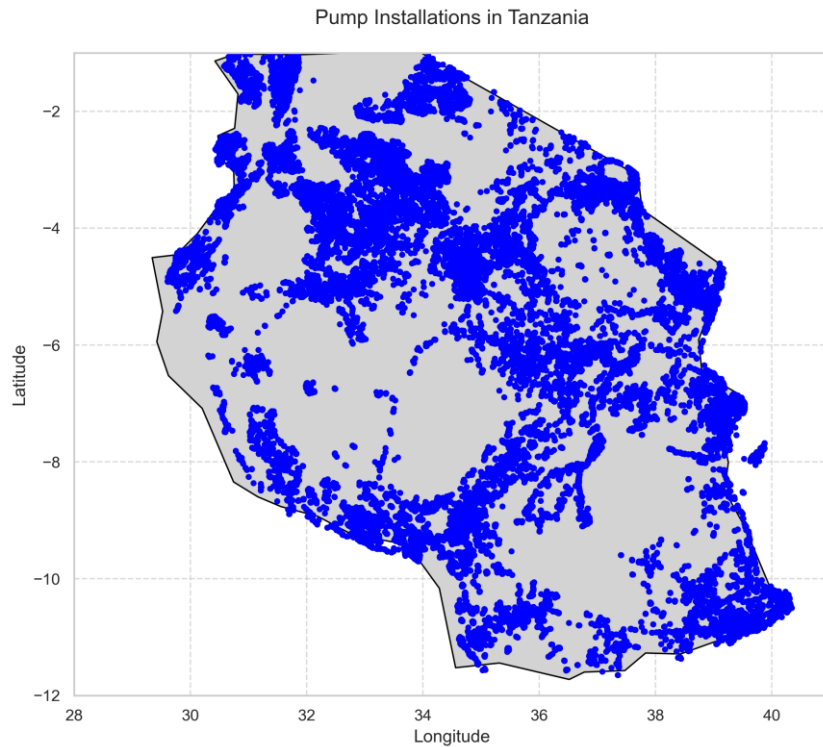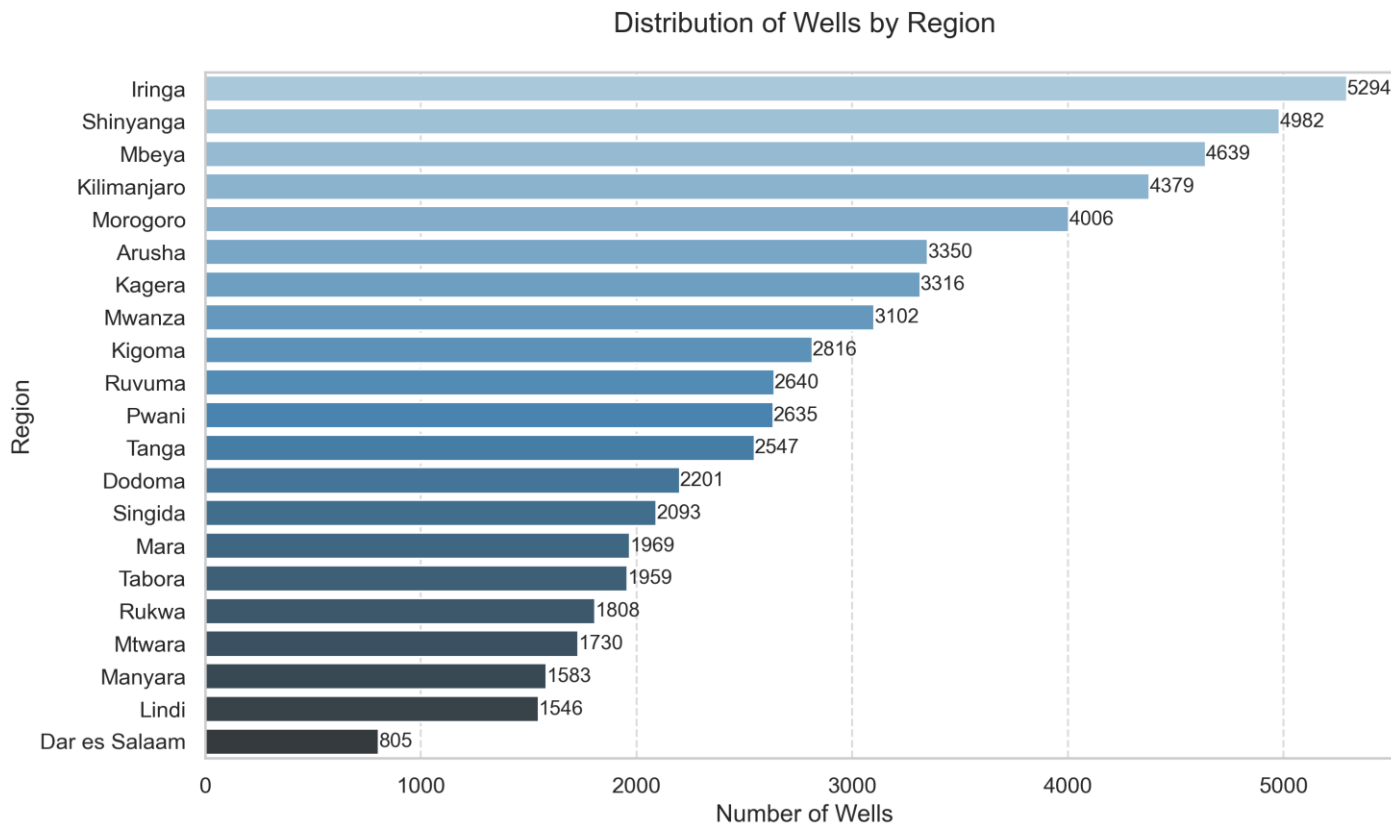Pump Status Distribution by Month Recorded

# Pump installation Density

# Water Pump in Tanzania



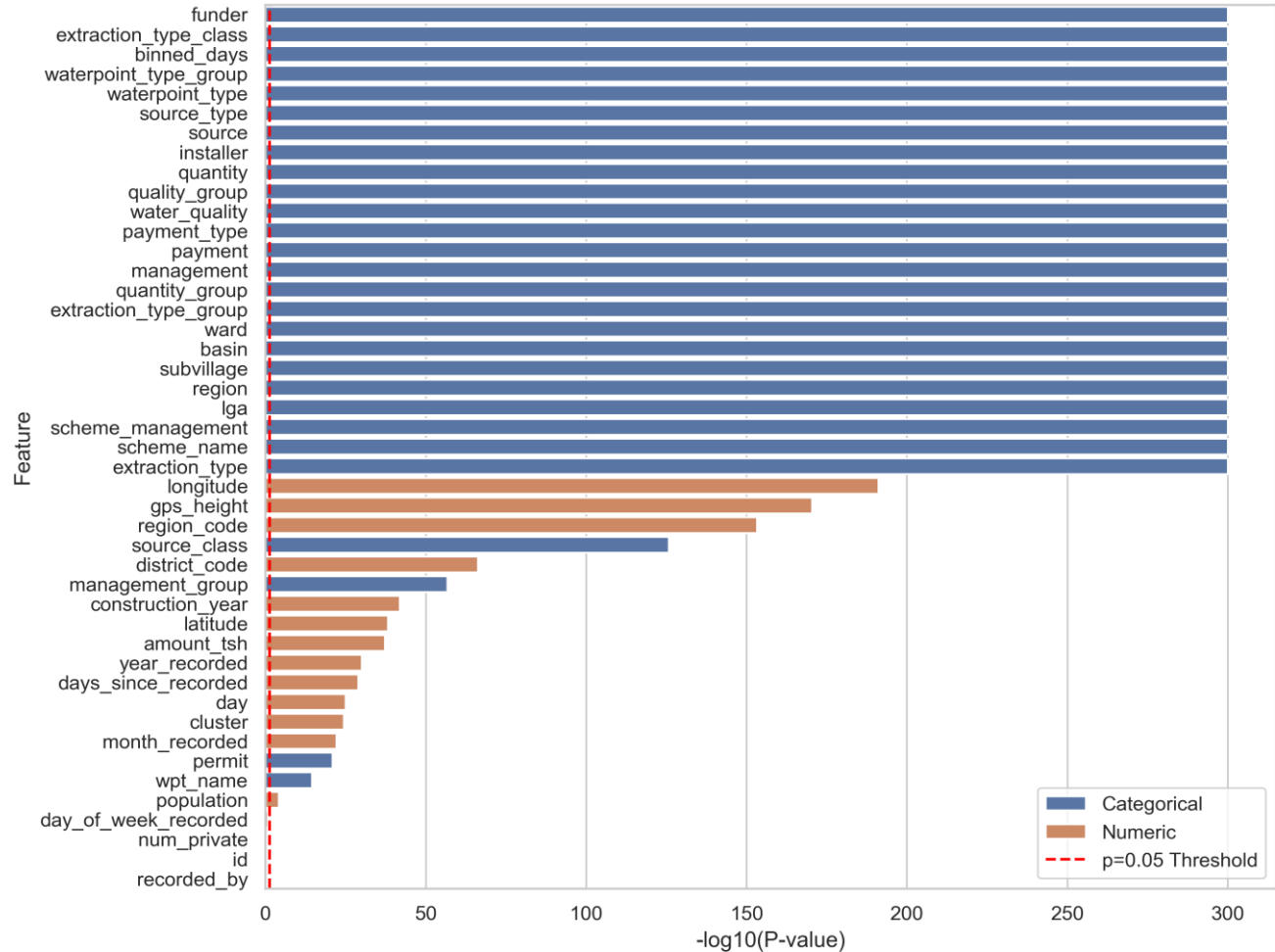Pump Installations in Tanzania

# Distrinution water pump based on region



Distribution of Wells by Region

Relationship of Features with status_group (-log10 P-values)

- ❖ Chi-squared Test
- ❖ ANOVA Test

# Thank you!