

Milestone 3: Beyond Descriptive Stats (Dive Deeper/Go Broader)

In this milestone, you will go beyond the descriptive statistics you completed in the last milestone. This milestone is really about diving deeper to analyze your data, beyond descriptive stats. Maybe you need to analyze qualitative data or textual data to get a full picture.

	Field	Description
	amount_tsh	Total static head (amount water available to waterpoint)
	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 waterpoint 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 waterpoint is permitted
	construction_year	Year the waterpoint was constructed
	extraction_type	The kind of extraction the waterpoint uses
	extraction_type_group	The kind of extraction the waterpoint uses
	extraction_type_class	The kind of extraction the waterpoint uses
	management	How the waterpoint is managed
	management_group	How the waterpoint 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_class	The source of the water
	source_type	The source of the water
	waterpoint_type	The kind of waterpoint
	waterpoint_type_group	The kind of waterpoint
	functional	the waterpoint is operational and there are no repairs needed
	non functional	the waterpoint is not operational

Dive Deeper

Look deeper into the features you are investigating, consider:

- Relationships / Correlation, Pearson Correlation
- Linear Regression for future prediction (if the relationship is linear)
- Textual Analysis for TF-IDF (Term Frequency-Inverse Document Frequency; Row-based and column-based, stop-word removal)?

Specify 1-2 correlations you discovered. List the fields that you found to be correlated and describe what you learned from these correlations.

- Basin versus status-group - Some areas have higher non functional pumps
- Payment type versus status-group - Those who never pay for water have higher non functional pumps

Go Broader

Expand the features you are investigating. Look for connections/relationships that you may have initially missed.

- What jumps out at you now?
- Use the descriptive stats to point you to features that you may now want to consider.

What key terms did you discover in any text analysis, for whom? Any themes? If you are not analyzing text, summarize what other things you are considering in your analysis?

- Public Meeting areas has higher counts of functional and non functional pumps.
- Majority of functional pumps installed have water point permits.
- Groundwater extraction type uses gravity method.
- Majority of users never pay for water usage.
- Main water sources are groundwater.

New Metric

Create 1 or 2 new metrics to track relationships of data you discovered. Explain why you created them.

Accuracy and F1 score. We must be able to pinpoint which water pump is faulty for repairs. F1 score is used due to imbalance class (functional versus non-functional).

Import Libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
import datetime

%matplotlib inline
sns.set_style('dark')
sns.set(font_scale=1.2)

import warnings
warnings.filterwarnings('ignore')

import feature_engine.missing_data_imputers as mdi
from feature_engine.outlier_removers import Winsorizer

from sklearn.preprocessing import LabelEncoder, StandardScaler, MinMaxScaler, OneHotEncoder

pd.set_option('display.max_columns',None)
pd.set_option('display.max_row',None)
pd.set_option('display.width',1000)

np.random.seed(0)
np.set_printoptions(suppress=True)

In [2]: df = pd.read_csv('train.csv',parse_dates=['date_recorded'])

In [3]: df

Out [3]:
```

	id	amount_tsh	date_recorded	funder	gps_height	installer	longitude	latitude	wpt_name	num_private	basin	population
0	69572	6000.0	2011-03-14	Roman	1390	Roman	34.930893	-9.855322	none	0	Lake Nya	179.009983
1	8776	0.0	2013-06-03	Grumeti	1399	GRUMETI	34.688766	-2.147466	Zahanati	0	Lake Victoria	0
2	34310	25.0	2013-02-25	Lobby Club	686	World vision	37.460664	-3.821329	Kwa Mahudi	0	Pangani	0
3	67743	0.0	2013-01-28	Unioef	263	UNICEF	38.486161	-11.155289	Zahanati Nanyumbu	0	Ruvuma / Southern Coast	0
4	19728	0.0	2011-07-13	Action In A	0	Artisan	31.130847	-1.825359	Shuleni	0	Lake Victoria	0
...
59395	60739	10.0	2013-03-05	Germany Repubi	1210	CES	37.169807	-3.253847	Area Three Namba 27	0	Pangani	0
59396	27263	4700.0	2011-07-05	Cefa-njombe	1212	Cefa	35.249991	-9.070629	Kwa Yafusa Kuvula	0	Rufiji	0
59397	37057	0.0	2011-11-04	NaN	0	NaN	34.017087	-8.750434	Mashine	0	Rufiji	0
59398	31282	0.0	2011-08-03	Malec	0	Musa	35.861315	-6.378573	Mahoro	0	Rufiji	0
59399	26348	0.0	2011-03-23	World Bank	191	World	38.104048	-6.747464	Kwa Mzee Lugawa	0	Wami / Ruwu	0

59400 rows x 41 columns

Exploratory Data Analysis

```
In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 59400 entries, 0 to 59399
Data columns (total 41 columns):
#   Column              Non-Null Count  Dtype  
---  --
0   id                   59400 non-null  int64  
1   amount_tsh          59400 non-null  float64
2   date_recorded       59400 non-null  datetime64[ns]
3   funder              55765 non-null  object  
4   gps_height          59400 non-null  int64  
5   installer           55743 non-null  object  
6   longitude           59400 non-null  float64
7   latitude           59400 non-null  float64
8   wpt_name            59400 non-null  object  
9   num_private         59400 non-null  int64  
10  basin               59400 non-null  object  
11  subvillage          59029 non-null  object  
12  region              59400 non-null  object  
13  region_code         59400 non-null  int64  
14  district_code       59400 non-null  int64  
15  lga                 59400 non-null  object  
16  ward               59400 non-null  object  
17  population          59400 non-null  int64  
18  public_meeting      56566 non-null  object  
19  recorded_by         59400 non-null  object  
20  scheme_management   55523 non-null  object  
21  scheme_name         31234 non-null  object  
22  permit              56344 non-null  object  
23  construction_year   59400 non-null  int64  
24  extraction_type     59400 non-null  object  
25  extraction_type_group 59400 non-null  object  
26  extraction_type_class 59400 non-null  object  
27  management          59400 non-null  object  
28  management_group    59400 non-null  object  
29  payment             59400 non-null  object  
30  payment_type        59400 non-null  object  
31  water_quality       59400 non-null  object  
32  quality_group       59400 non-null  object  
33  quantity            59400 non-null  object  
34  quantity_group      59400 non-null  object  
35  source              59400 non-null  object  
36  source_type         59400 non-null  object  
37  source_class        59400 non-null  object  
38  waterpoint_type     59400 non-null  object  
39  waterpoint_type_group 59399 non-null  object  
40  status_group        59399 non-null  object  
dtypes: datetime64[ns](1), float64(3), int64(7), object(30)
memory usage: 18.4+ MB

In [5]: df.describe()

Out [5]:
```

	id	amount_tsh	gps_height	longitude	latitude	num_private	region_code	district_code	population
count	59400.000000	59400.000000	59400.000000	59400.000000	5.940000e+04	59400.000000	59400.000000	59400.000000	59400.000000
mean	37115.131768	317.650385	668.297239	6.507422	-2.946031e+00	0.474141	15.297003	5.629747	179.009983
std	21453.131781	299.574558	693.116350	6.507422	2.946031e+00	12.236230	17.587406	9.633649	471.482176
min	0.000000	0.000000	-90.000000	0.000000	-1.164944e+01	0.000000	1.000000	0.000000	0.000000
25%	18519.750000	0.000000	0.000000	33.060347	-8.540621e+00	0.000000	5.000000	2.000000	0.000000
50%	37061.500000	0.000000	369.000000	34.908743	-5.021597e+00	0.000000	12.000000	3.000000	25.000000
75%	56566.500000	20.000000	1919.250000	37.176387	-3.326156e+00	0.000000	17.000000	5.000000	215.000000
max	74427.000000	35000.000000	2770.000000	40.345193	-2.000000e-08	1776.000000	99.000000	80.000000	3050.000000

```
In [6]: df.columns

Out [6]:
```

Index(['id', 'amount_tsh', 'date_recorded', 'funder', 'gps_height', 'installer', 'longitude', 'latitude', 'wpt_name', 'num_private', 'basin', 'subvillage', 'region', 'region_code', 'district_code', 'lga', 'ward', 'population', 'public_meeting', 'recorded_by', 'scheme_management', 'scheme_name', 'permit', 'construction_year', 'extraction_type', 'extraction_type_group', 'extraction_type_class', 'management', 'management_group', 'payment', 'payment_type', 'water_quality', 'quality_group', 'quantity', 'quantity_group', 'source', 'source_type', 'source_class', 'waterpoint_type', 'waterpoint_type_group', 'status_group'], dtype='object')

```
In [7]: df['funder'].value_counts()

Out [7]:
```

Government Of Tanzania	9084
Danida	3114
Hexava	1292
Bwasp	1374
World Bank	1349
Babtest	1
Nyangere	1
Sinango Kihengu	1
Eze Rappa	1
Muslimu Society(shia)	1
Name: funder, Length: 1897, dtype: int64	

```
In [8]: df['installer'].value_counts()

Out [8]:
```

DWE	17402
Government	1292
RWE	1206
Commu	1060
DAMTA	1050
...	...
AFRICAN REFLECTIONS FOUNDATION	1
VC	1
go	1
Gerald Milla	1
Bonite Bottles Ltd	1
Name: installer, Length: 2145, dtype: int64	

```
In [9]: df['wpt_name'].value_counts()

Out [9]:
```

none	3563
Shuleni	1748
Zahanati	830
Makikitini	535
Kawisa	323
...	...
Kwa Bwana Ventura	1
Kwa Winiu Masafu	1
Kwa Timothy Mbembati	1
Sabasaba	1
Kwa Mama Judith	1
Name: wpt_name, Length: 37400, dtype: int64	

```
In [10]: df['num_private'].value_counts()

Out [10]:
```

0	58643
6	81
1	73
5	46
8	46
...	...
213	1
23	1
55	1
94	1
Name: num_private, Length: 65, dtype: int64	

```
In [11]: df['basin'].value_counts()

Out [11]:
```

Lake Victoria	10248
Pangani	8940
Rufiji	7976
Internal	7785
Lake Tanganyika	6432
Kasi / Ruwu	5987
Lake Nyasa	5085
Ruvuma / Southern Coast	4493
Lake Rukwa	2454
Name: basin, dtype: int64	

```
In [12]: df['subvillage'].value_counts()

Out [12]:
```

Madukani	508
Shuleni	506
Majengo	502
Kata	493
Mtakuja	262
...	...
Bunyanya	1
Maninga	1
Mkatoni	1
Majiwe	1
Heka Kati	1
Name: subvillage, Length: 19287, dtype: int64	

```
In [13]: df['region'].value_counts()

Out [13]:
```

Iringa	5294
Shinyanga	4982
Moshi	4639
Kilimanjaro	4379
Morogoro	4006
Arusha	3350
Rapera	3315
Mwanza	3102
Rigoma	2816
Ruvuma	2640
Pwani	2635
Tanga	2547
Dodoma	2201
Singida	2093
Mara	1969
Tabora	1959
Rukwa	1808
Mtwara	1730
Manyara	1583
Lindi	1546
Dar es Salaam	805
Name: region, dtype: int64	

```
In [14]: df['region_code'].value_counts()

Out [14]:
```

11	5300
17	5011
12	4639
3	4379
5	4040
18	3324
19	3047
2	3024
16	2816
4	2513
1	2201
13	2093
14	1979
20	1669
15	1608
6	1609
21	1583
80	1238
10	1025
7	917
9	805
99	423
9	390
24	326
8	300
40	1
Name: region_code, dtype: int64	

```
In [15]: df['district_code'].value_counts()

Out [15]:
```

1	12203
2	11173
3	8998
4	8999
5	4356
6	4074
7	3343
8	1043
30	995
33	874
53	745
43	505
13	391
23	293
63	195
62	109
60	63
0	23
80	12
67	6
Name: district_code, dtype: int64	

```
In [16]: df['lga'].value_counts()

Out [16]:
```

Njombe	2503
Arusha Rural	1252
Moshi Rural	1251
Bariadi	1177
Rungwe	1106
...	...
Moshi Urban	79
Kigoma Urban	71
Aruusha Urban	63
Lindi Urban	21
Nyamagana	1
Name: lga, Length: 125, dtype: int64	

```
In [17]: df['ward'].value_counts()

Out [17]:
```

Igesi	307
Tealinyi	252
Siha Kati	232
Miandu	231
Mduruma	217
...	...
Themti	1
Riviera	1
Wapiga Kaskazini	1
Mawenzi	1
Korongoni	1
Name: ward, Length: 2092, dtype: int64	

```
In [18]: df['public_meeting'].value_counts()

Out [18]:
```

True	51011
False	5055
Name: public_meeting, dtype: int64	

```
In [19]: df['recorded_by'].value_counts()

Out [19]:
```

GeoData Consultants Ltd	59400
Name: recorded_by, dtype: int64	

```
In [20]: df['scheme_management'].value_counts()

Out [20]:
```

VWC	36793
MUG	5206
Water authority	2153
MUA	2883
Water Board	3148
Parastatal	1680
Private operator	1063
Company	1061
Other	766
SWC	97
Trust	72
None	1
Name: scheme_management, dtype: int64	

```
In [21]: df['scheme_name'].value_counts()

Out [21]:
```

K	682
None	644
Borehole	546
Chalinsie wate	405
M	400
...	...
Mwigimbii piped scheme	212
Kakonko/Mbizi gravity water supply	1
Pwani water supply	1
Mwadi wa maji wa natalawe	1
QUICK WINDS	1
Name: scheme_name, Length: 2696, dtype: int64	

```
In [22]: df['permit'].value_counts()

Out [22]:
```

True	38852
False	17492
Name: permit, dtype: int64	

```
In [23]: df['construction_year'].value_counts()

Out [23]:
```

0	20709
2010	2645
2008	2613
2009	2533
2000	2091
2007	1587
2006	1471
2003	1286
2011	1256
2004	1123
2012	1084
2002	1075
1978	1037
1995	1014
2005	1011
1999	979
1998	966
1990	954
1985	945
1980	811
1984	779
1982	744
1984	738
1972	708
1974	676
1987	644
1992	640
1993	608
2001	540
1998	521
1983	488
1975	437
1986	434
1976	414
1970	411
1991	324
1989	316
1987	302
1981	238
1977	202
1979	192
1973	184
2013	176
1971	145
1960	102
1967	88
1983	85
1968	77
1969	59
1964	40
1962	30
1961	21
1965	19
1966	17
Name: construction_year, dtype: int64	

```
In [24]: df['extraction_type'].value_counts()

Out [24]:
```

gravity	26780
nira/tanira	8154
other	6430
submersible	4764
swm 80	3470
mono	2865
india mark ii	2600
afridev	1770
kab	1415
other - rope pump	451
other - swm 81	229
windmill	117
india mark iii	98
cemo	90
other - play pump	85
walimal	48
climax	32
other - mkulima/shinyanga	2
Name: extraction_type, dtype: int64	

```
In [25]: df['extraction_type_group'].value_counts()

Out [25]:
```

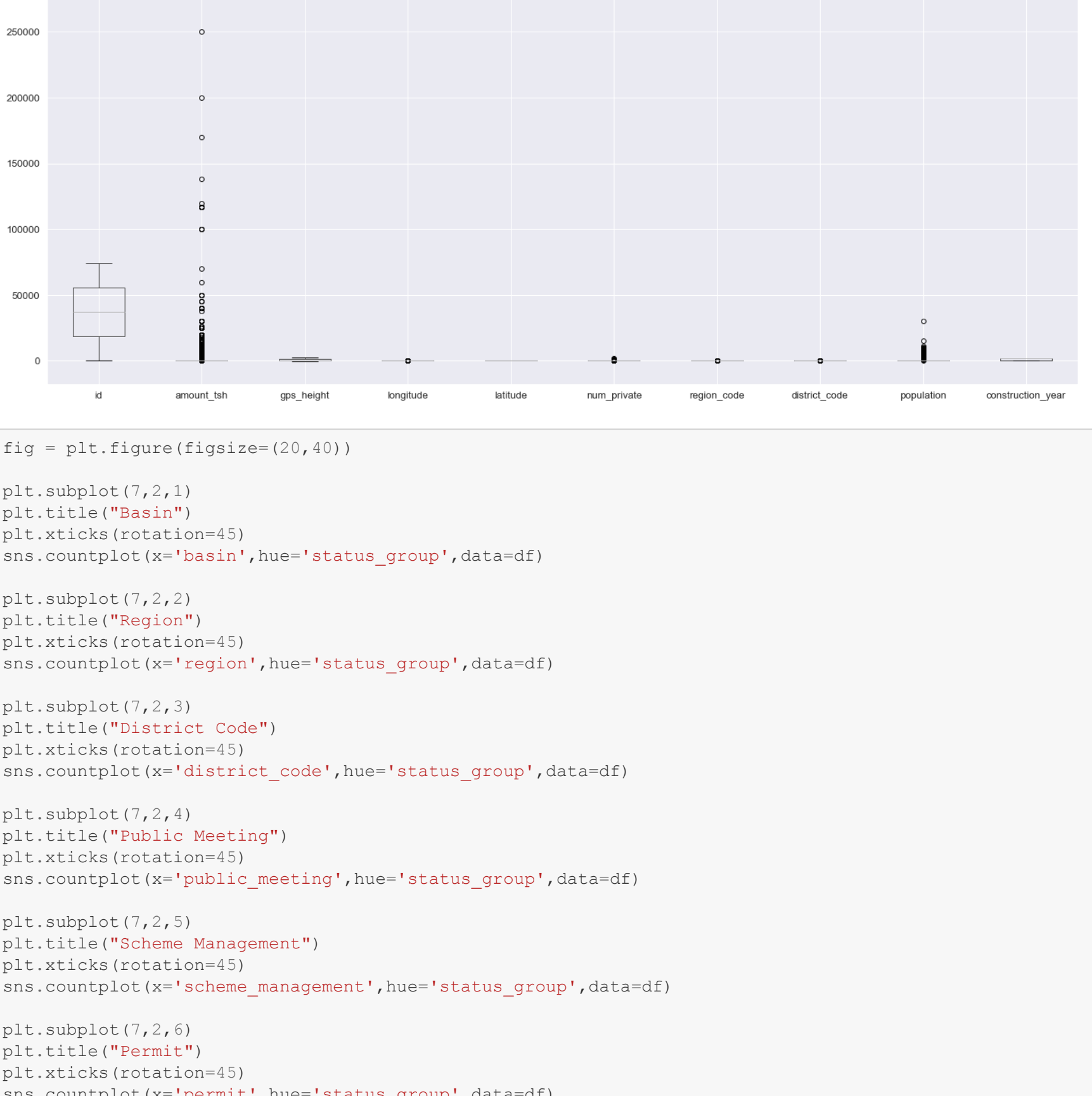
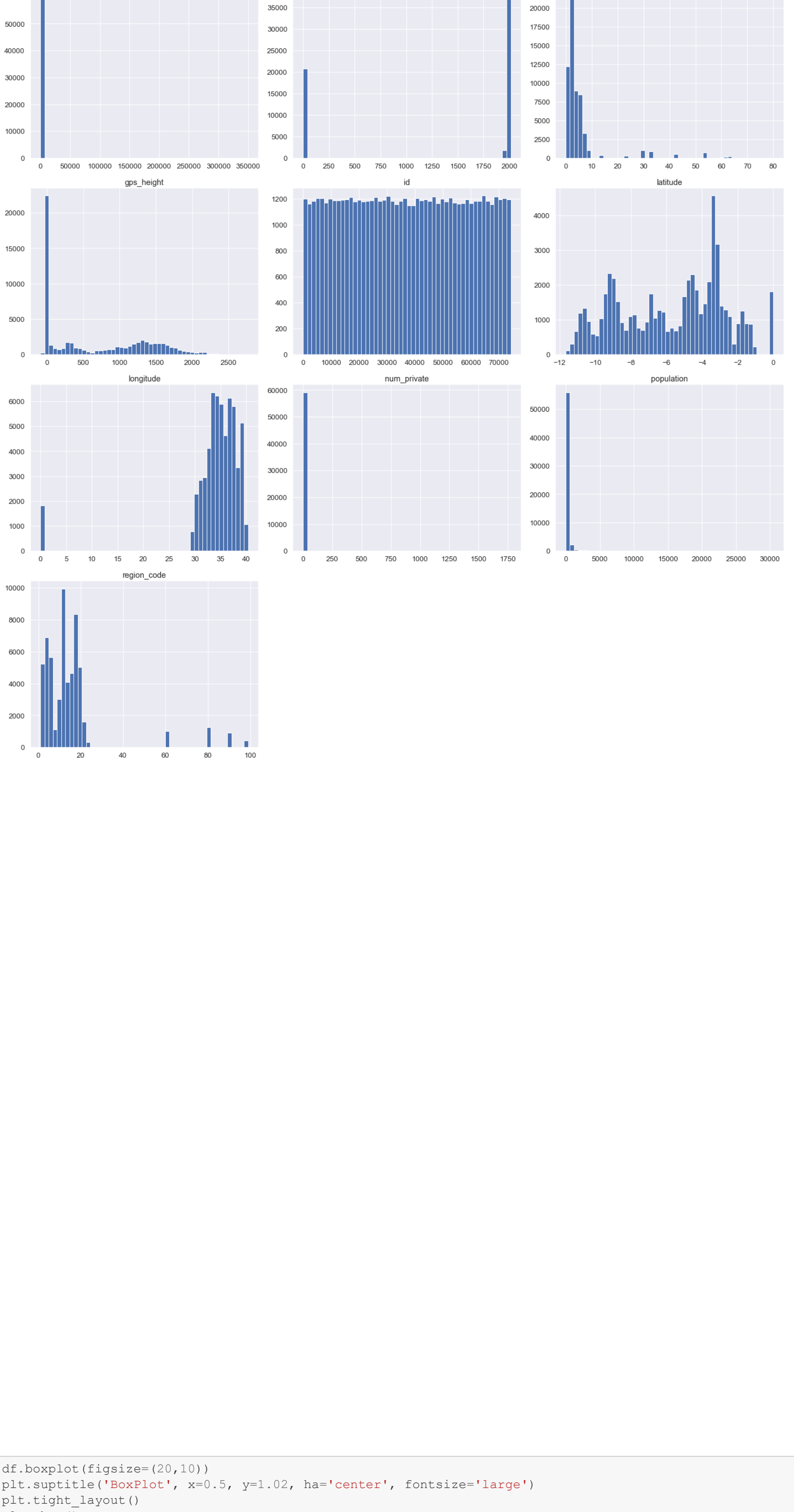
gravity	26780
nira/tanira	8154
other	6430
submersible	6179
swm 80	3670
mono	2865
india mark ii	2400
afridev	1770
rope pump	451
other handpump	364
other motorpump	122
wind-powered	117
india mark iii	98
Name: extraction_type_group, dtype: int64	

```
In [26]: df['extraction_type_class'].value_counts()

Out [26]:
```

gravity	26780
handpump	16456
other	6430
submersible	6179
motorpump	2987
rope pump	451


```
In [41]: df.hist(bins=50, figsize=(20,20))
plt.suptitle('Feature Distribution', x=0.5, y=1.02, ha='center')
plt.tight_layout()
plt.show()
```



```

sns.pairplot(df)

plt.subplot(7,2,7)
plt.title("Extraction")
plt.xticks(rotation=45)
sns.countplot(x="extraction_type_class",hue='status_group',data=df)

plt.subplot(7,2,8)
plt.title("Management Grouping")
plt.xticks(rotation=45)
sns.countplot(x="management_group",hue='status_group',data=df)

plt.subplot(7,2,9)
plt.title("Payment Type")
plt.xticks(rotation=45)
sns.countplot(x="payment_type",hue='status_group',data=df)

plt.subplot(7,2,10)
plt.title("Quality Grouping")
plt.xticks(rotation=45)
sns.countplot(x="quality_group",hue='status_group',data=df)

plt.subplot(7,2,11)
plt.title("Quantity Grouping")
plt.xticks(rotation=45)
sns.countplot(x="quantity_group",hue='status_group',data=df)

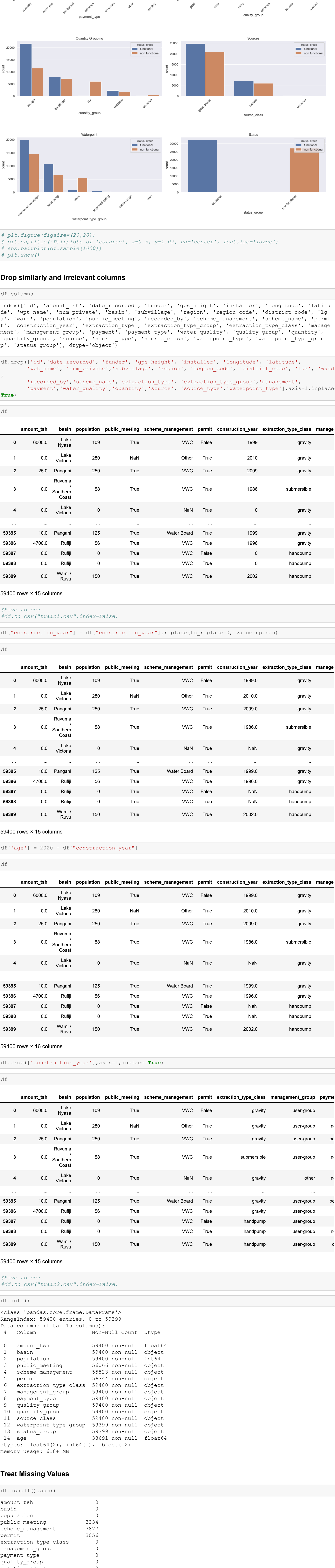
plt.subplot(7,2,12)
plt.title("Sources")
plt.xticks(rotation=45)
sns.countplot(x="source_class",hue='status_group',data=df)

plt.subplot(7,2,13)
plt.title("Waterpoint")
plt.xticks(rotation=45)
sns.countplot(x="waterpoint_type_group",hue='status_group',data=df)

plt.subplot(7,2,14)
plt.title("Status")
plt.xticks(rotation=45)
sns.countplot(x="status_group",hue='status_group',data=df)

plt.tight_layout()
plt.show()

```



```

quantity_group      0
source_class         0
waterpoint_type_group 1
status_group         1
age                 20709
dtype: int64

```

In [58]: `df.dropna(inplace=True)`

```
df.isnull().sum()
amount_tsh      0
has_invoice      0
```

E
 F
 F
 S
 F

extraction_type_class	0
management_group	0
payment_type	0
quality_group	0
quantity_group	0
source_class	0
waterpoint_type_group	0
status_group	0

```
age                                0
dtype: int64

In [60]: df

Out[60]:
```

	amount_tsh	basin	population	public_meeting	scheme_management
0	2000000	Lake	100	True	1000000

0	0000.0	Nyssa	109	True	VW0
2	25.0	Pangani	250	True	VW0
3	0.0	Ruvuma Southern Coast	58	True	VW0
5	20.0	Pangani	1	True	VW0

	10	0.0	Wami / Ruvu	345	True	Private operator
...
59391	0.0	Pangani	210	True	Water authority	
59394	500.0	Wami / Ruvu	89	True	VWC	
59395	10.0	Pangani	126	True	Water Board	

id	name	age	sex	status	weight	height
59399	Wami / Rufu	700	True	True	150	150
59396	Wami / Rufu	4700	True	True	56	56
59399	Wami / Rufu	0	True	True	150	150

32514 rows x 7 columns

```
df.groupby('name').agg({'age': 'sum', 'weight': 'sum'})
```

```
In [61]: df.reset_index(drop=True, inplace=True)
```

```
In [62]: df
```

```
Out[62]:
```

	amount_tsh	basin	population	public_meeting	scheme_management
0	6000.0	Lake Nyasa	109	True	VWC

1	25.0	Pangani	250	True	WFO
2	0.0	Ruvuma / Southern Coast	58	True	WFO
3	20.0	Pangani	1	True	WFO
4	0.0	Wami /	345	True	Private operators

32509	0.0	Pangani	210	True	Water authority
32510	500.0	Wami / Ruvu	89	True	WVO
32511	10.0	Pangani	125	True	Water Board
32512	1200.0	Pangani	66	True	Water Board

32512	4700.0	Ruip	50	True	VWC
32513	0.0	Wami / Ruvi	150	True	VWC

32514 rows x 15 columns

```
In [63]: #Save to csv
         #df.to_csv("train3.csv",index=False)
```

Correlation

```
In [64]: df.corr()
```

```
Out[64]:
```

	amount_tsh	population	age
--	------------	------------	-----

```

In [65]: plt.figure(figsize=(16,5))
          sns.heatmap(df.corr(), cmap='coolwarm', annot=True, fmt='.2f',

```

```
plt.show()
```

	low	high
low	1.00	-0.01
high	-0.01	1.00

population	-0.01	1.00
age	-0.01	-0.02

```
In [66]: sns.pairplot(df[['amount_tsh', 'population', 'age']].sample(
plt.suptitle('Pairplots Of features', x=0.5, y=1.02, ha='center',
plt.show())
```

Pairplots of features

amount fish species

The figure contains three plots. The first is a scatter plot of 'age' (y-axis, 0-30) against 'amount_tsh' (x-axis, 0-40,000). The second is a scatter plot of 'age' (y-axis, 0-30) against 'population' (x-axis, 0-5,000). The third is a histogram of 'age' with a y-axis from 0 to 20 and an x-axis from 0 to 40. The histogram shows a distribution peaking around age 10-15.

Treat Duplicate Values

```
In [67]: df.duplicated(keep='first').sum()
Out[67]: 9658
```

```
In [68]: df[df.duplicated(keep=False)] #Check duplicate values
```

```
Out[68]:
```

	amount_tsh	basin	population	public_meeting	scheme_managemen
0	6000.0	Lake Nyasa	109	True	WV
1	25.0	Pangani	250	True	WV
6	0.0	Dogoni	4	True	Water

6	0.0	Panglima	1	True	Water Dis
7	0.0	Lake Tanganyika	200	True	WV
9	0.0	Rufiji	50	True	WV
...	
32493	0.0	Lake Rukwa	1	False	WV

32497	2000.0	Lake Nyssa	0	True	Water Bo
32503	6.0	Pangani	1	True	Water Bo
32509	0.0	Pangani	210	True	Water autho
32511	10.0	Pangani	125	True	Water Bo

13039 rows × 15 columns

Treat Outliers

```
In [69]: df.columns
```

```
Out[69]: Index(['amount_tsh', 'basin', 'population', 'public_meeting
```

```
In [70]: df.describe()
```

	amount_tsh	population	age
count	32514.000000	32514.000000	32514.000000

mean	515.297933	257.588639	23.270776
std	3235.720694	548.275049	12.604989
min	0.000000	0.000000	7.000000
25%	0.000000	25.000000	12.000000
50%	0.000000	130.000000	20.000000

	75%	250.000000	300.000000	34.000000
	max	250000.000000	305000.000000	60.000000

```
In [71]: windsorizer = Winsorizer(distribution='skewed',tail='right',on='')

In [72]: windsorizer.fit(df)
```

```
Out[72]: Winsorizer(distribution='skewed', fold=3.0,
                  variables=['amount_tsh', 'population'])

In [73]: df2 = winsorizer.transform(df)

In [74]: df2
```

	amount_tsh	basin	population	public_meeting	scheme_management
0	1000.0	Lake Nyasa	109.0	True	VWOC
1	25.0	Pangani	250.0	True	VWOC
2	0.0	Ruvuma	58.0	True	VWOC

		Southern Coast				
3	20.0	Pangani	1.0	True	VWC	
4	0.0	Wami / Ruvu	345.0	True	Private operator	
32599	0.0	Pangani	210.0	True	Water authority	

32510	500.0	Wami / Ruvu	89.0	True	VWC
32511	10.0	Pangani	125.0	True	Water Board
32512	1000.0	Rufiji	56.0	True	VWC
32513	0.0	Wami / Ruvu	150.0	True	VWC

```
32514 rows × 15 columns
```

```
In [75]: df2.describe()
```

```
Out[75]:
```

	amount_tsh	population	age
count	32514.000000	32514.000000	32514.000000

mean	195.524912	217.319893	23.270776
std	339.119814	263.161043	12.604989
min	0.000000	0.000000	7.000000
25%	0.000000	25.000000	12.000000
50%	0.000000	130.000000	20.000000
75%	350.000000	300.000000	24.000000

	75%	250.000000	500.000000	750.000000
max	1000.000000	1125.000000	60.000000	

```
In [76]: windsorizer.left_tail_caps_
Out[76]: {}

In [77]: windsorizer.right_tail_caps_
Out[77]: {}
```

```
In [77]: df2.to_csv('train4.csv', sep=';',
Out[77]: ('amount_tsh': 1000.0, 'population': 1125.0)

In [78]: #Save to csv
         #df2.to_csv("train4.csv", index=False)

In [79]: df2.columns
```

```
Out[79]: Index(['amount_tsh', 'basin', 'population', 'public_meeting',
       'type_class', 'management_group', 'payment_type', 'quality',
       'a', 'waterpoint_type_group', 'status_group', 'age'], dtype=
object)

In [80]: #Rearrange columns
df2 = df2[['amount_tsh', 'age', 'population', 'basin', 'public_meeting',
           'extraction_type_class', 'management_group', 'payment_type',
           'a', 'waterpoint_type_group', 'status_group', 'quality']]
```

```
y = {
    'source_class', 'waterpoint_type_group', 'status'
```