

EDA and Stakeholder Questions Notebook

```
In [1]: 1 # Import the required libraries
        2 import pandas as pd
        3 import numpy as np
        4
        5 import scipy.stats as scs
        6
        7 import matplotlib.pyplot as plt
        8 import matplotlib.colors
        9 import seaborn as sns
        10
        11 # set up pandas to display floats in a more human friendly way
        12 pd.options.display.float_format = '{:,.2f}' format
```

```
In [2]: 1 def get_percentage_summary(label, x, y):
        2     percentage = round(x/y, 3)
        3     summary = str(round(x/y, 3)) + '% ' + label
        4     return summary
```

```
In [3]: 1 # read in the processed data
        2 df = pd.read_csv('../data/train_processed_labeled.csv')
        3 print(df.shape)
        4 df.head(3)
```

(57565, 44)

Out[3]:

| | id | amount_tsh | date_recorded | funder | gps_height | installer | longitude | latitude | wpt_name | basin | ... | source | source_type | so |
|---|-------|------------|---------------|--------------|------------|--------------|-----------|----------|-------------|---------------|-----|----------------------|----------------------|-----|
| 0 | 69572 | 6,000.00 | 2011-03-14 | roman | 1390 | roman | 34.94 | -9.86 | none | lake nyasa | ... | spring | spring | ... |
| 1 | 8776 | 0.00 | 2013-03-06 | grumeti | 1399 | grumeti | 34.70 | -2.15 | zahanati | lake victoria | ... | rainwater harvesting | rainwater harvesting | ... |
| 2 | 34310 | 25.00 | 2013-02-25 | lottery club | 686 | world vision | 37.46 | -3.82 | kwa mahundi | pangani | ... | dam | dam | ... |

3 rows x 44 columns

Processed Data Column Descriptions

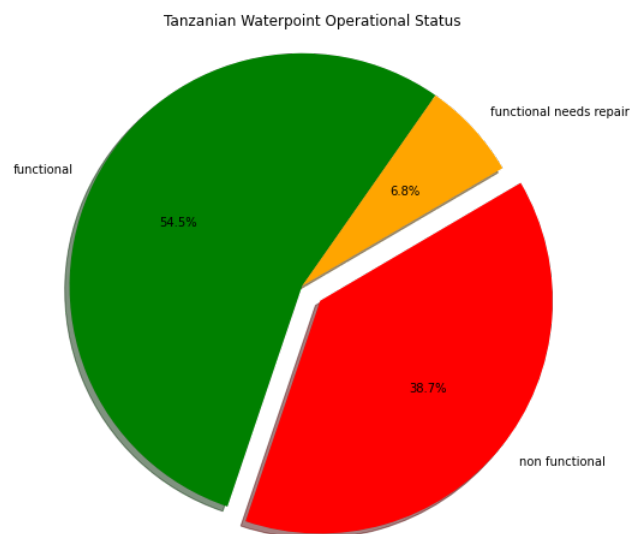
- id - Numeric identifier for the waterpoint
- amount_tsh - Total static head (amount water available to waterpoint)
- date_recorded - The date the row was entered
- funder - Who funded the well
- gps_height - Altitude of the well
- installer - Organization that installed the well
- longitude - GPS coordinate
- latitude - GPS coordinate
- wpt_name - Name of the waterpoint if there is one
- basin - Geographic water basin
- subvillage - Geographic location
- region - Geographic location, NOTE: Hierarchy is Region > LGA > Ward
- 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_type - The source of the water
- source_class - The source of the water
- waterpoint_type - The kind of waterpoint
- waterpoint_type_group - The kind of waterpoint
- recorded_year - Pulling out the year from date_recorded
- waterpoint_age - Calculate as recorded_year - construction_year
- region_with_code - Combine region and region_code. There are more region_code values than region values
- recorded_good_quality - True if quality_group == 'good', False if anything other than 'good'
- recorded_good_quantity - True if quantity_group == 'sufficient', False if anything other than 'sufficient'
- status_group - Operational status (these are the 3 classes we will attempt to predict on Test data)

Question 1: What is the operational status of waterpoints in Tanzania?

1a: What is the overall waterpoint Operational Status?

```
In [4]: 1 by_op_status = df.groupby('status_group')['id'].count()
        2 by_op_status.sort_values(ascending=False, inplace=True)
```

```
In [5]: 1 plfig = plt.figure(figsize = (8, 8))
2 my_explode = (0, 0.1, 0)
3 my_colors = ['green', 'red', 'orange']
4 plt.pie(by_op_status, labels=by_op_status.index, autopct='%1.1f%%', startangle=55, shadow =True, c
5 plt.title('Tanzanian Waterpoint Operational Status')
6 plt.axis('equal')
7 plt.rcParams.update({'font.size': 20})
8 plt.show()
```

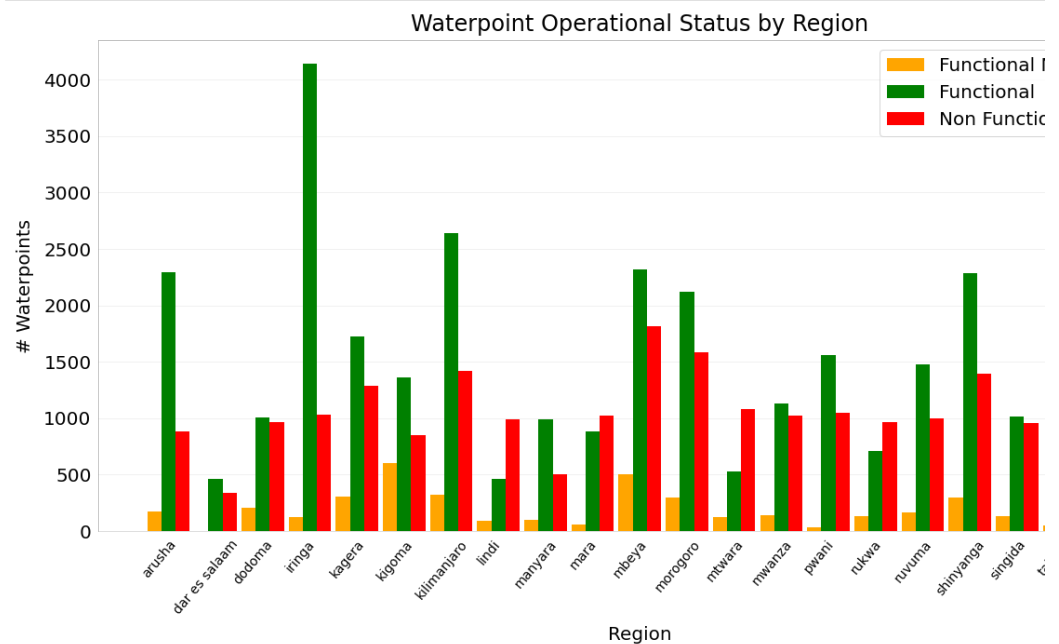


1b: What is the breakdown of waterpoint operational status by Region?

```

In [6]: 1 # Create a grouped bar chart, with region as the x-axis and status_group as the variable we're gro
2 fig, ax = plt.subplots(figsize=(18, 10))
3
4 # Our x-axis. We just want a list of numbers from zero with a value for each of the regions.
5 x = np.arange(len(df.region.unique()))
6
7 # Define bar width. We need this to offset the additional bars.
8 bar_width = 0.3
9
10 b1_series = df[df['status_group'] == 'functional needs repair'].groupby('region')['id'].count().so
11 b2_series = df[df['status_group'] == 'functional'].groupby('region')['id'].count().sort_index(0)
12 b3_series = df[df['status_group'] == 'non functional'].groupby('region')['id'].count().sort_index(
13
14
15 b1 = ax.bar(x, b1_series, width=bar_width, label='Functional Needs Repair', color='orange')
16 # Same thing, but offset the x.
17 b2 = ax.bar(x + bar_width, b2_series, width=bar_width, label='Functional', color='green')
18 # Same thing, but offset the x again
19 b3 = ax.bar(x + (bar_width*2), b3_series, width=bar_width, label='Non Functional', color='red')
20
21 # Fix the x-axes.
22 ax.set_xticks(x + bar_width / 3)
23 ax.set_xticklabels(b1_series.index, fontsize=14, rotation=50)
24
25 # Add legend.
26 ax.legend()
27
28 # Axis styling.
29 ax.spines['top'].set_visible(False)
30 ax.spines['right'].set_visible(False)
31 ax.spines['left'].set_visible(False)
32 ax.spines['bottom'].set_color('#DDDDDD')
33 ax.tick_params(bottom=False, left=False)
34 ax.set_axisbelow(True)
35 ax.yaxis.grid(True, color='EEEEEE')
36
37 # Add axis and chart labels.
38 ax.set_xlabel('Region', labelpad=10)
39 ax.set_ylabel('# Waterpoints', labelpad=10)
40 ax.set_title('Waterpoint Operational Status by Region', pad=10)
41 fig.tight_layout()
42

```



Insights:

- The 3 Regions with the highest number of Functional waterpoints are Iringa, Kilimanjaro, and Kagera.
- The 3 Regions with the lowest number of Functional waterpoints are Mtwara, Lindi, and Dar es Salaam.
- The 3 Regions with the highest number of Functional Needs Repair waterpoints are Kigoma, Mbeya, and Kilimanjaro.
- 16 of the 21 Regions have a higher number of Functional waterpoints than Non Functional waterpoints.
- 5 of the 21 Regions have a higher number of Non Functional waterpoints than Functional waterpoints. They are Lindi, Mara, Mtwara, Rukwa, and Morogoro.

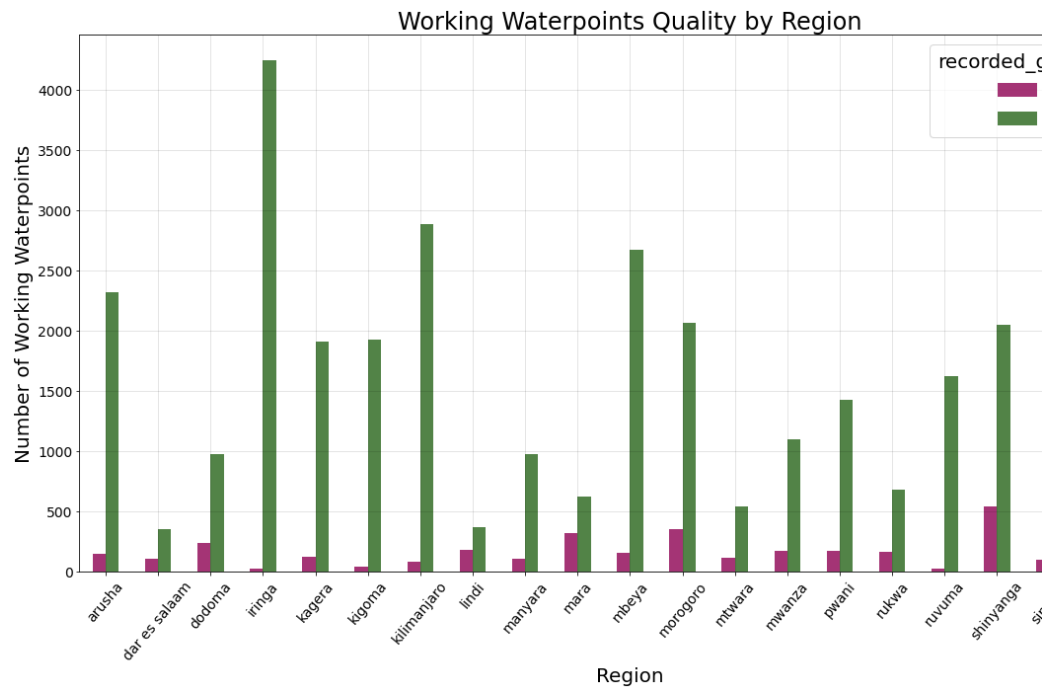
Recommendations:

- The regions of Lindi, Mara, Mtwara, Rukwa, and Tabora could potentially benefit from a water needs assessment. These regions have a higher number of Non Functional waterpoints than Functional waterpoints.
- Conduct future analysis into the Functional Needs Repair waterpoints in Kigoma, Mbeya, and Kilimanjaro.

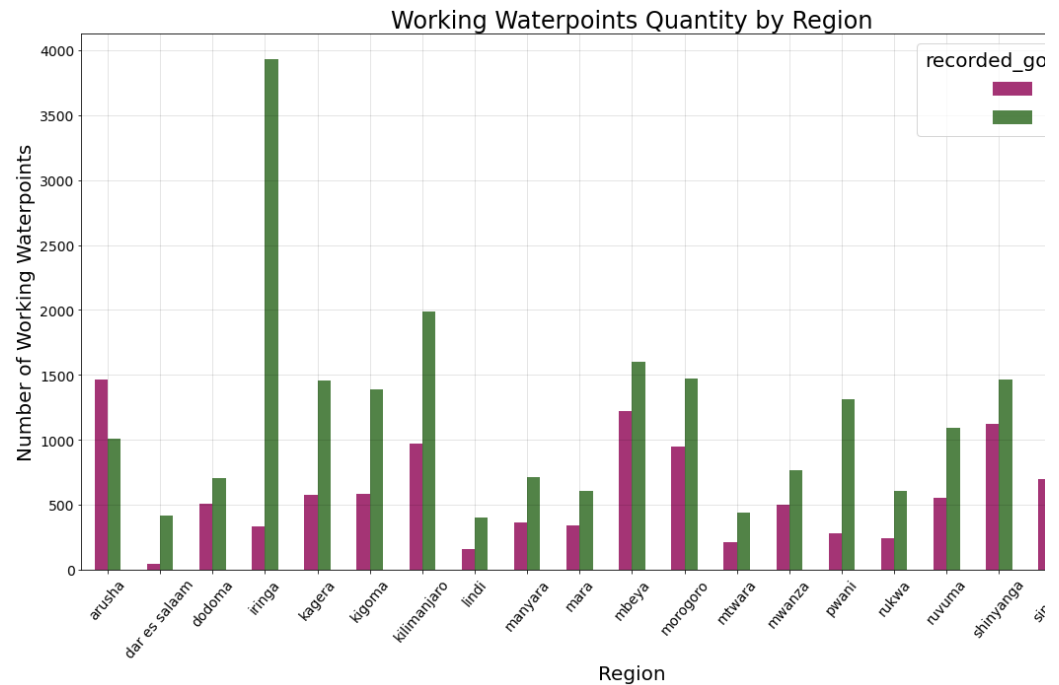
```
1 ## Question 2: What is the reported Quality and Quantity of Working (Functioning Needs Repair) waterpoints?
2 ### 2a: For All Working waterpoints, what is the breakdown of water quality by Region?
3 ### 2b: For All Working waterpoints, what is the water quantity by Region?
```

```
In [7]: 1 #Quality group
2 working_by_region_quality = df[df['status_group'] != 'non functional'].groupby(['region', 'recorded_quality']).count().reset_index()
3 #Quantity group
4 working_by_region_quantity = df[df['status_group'] != 'non functional'].groupby(['region', 'recorded_quantity']).count().reset_index()
```

```
In [8]: 1 working_by_region_quality.unstack().plot.bar(fontsize=14, rot=50, alpha = 0.80, figsize=(20,10), color = 'black', alpha = 0.1, linestyle = '-', linewidth = 1)
2 plt.grid(color = 'black', alpha = 0.1, linestyle = '-', linewidth = 1)
3 plt.xlabel('Region')
4 plt.ylabel('Number of Working Waterpoints')
5 plt.show()
```



```
In [9]: 1 working_by_region_quantity.unstack().plot.bar(fontsize=14, rot=50, alpha = 0.80, figsize=(20,10),
2         plt.grid(color = 'black', alpha = 0.1, linestyle = '-', linewidth = 1)
3         plt.xlabel('Region')
4         plt.ylabel('Number of Working Waterpoints')
5         plt.show()
```



Insights:

- All 21 Regions have a higher number of 'good quality' waterpoints than 'insufficient quality' waterpoints.
- 19 of the 21 Regions have a higher number of 'good/sufficient quantity' waterpoints than 'insufficient quantity' waterpoints.
- 2 of the 21 Regions have a higher number of 'insufficient quantity' waterpoints than 'good/sufficient quantity' waterpoints.

Recommendations:

- The regions of Indi, Mara, and Singida could potentially benefit from a water needs assesment. These regions have a higher ratio of "insufficient quality" waterpoints than the other Regions.
- The Arusha and Singida regions could potentially benefit from a water needs assesment. These are the only regions with a higher number of "insufficient quantity" waterpoints than "sufficient quantity" waterpoints.

Question 3: Is there a difference between the average age of Waterpoints by Operational Status?

3a: Difference between average age of Working (Functional and Functional Needs Repair) and Non Functional waterpoints?

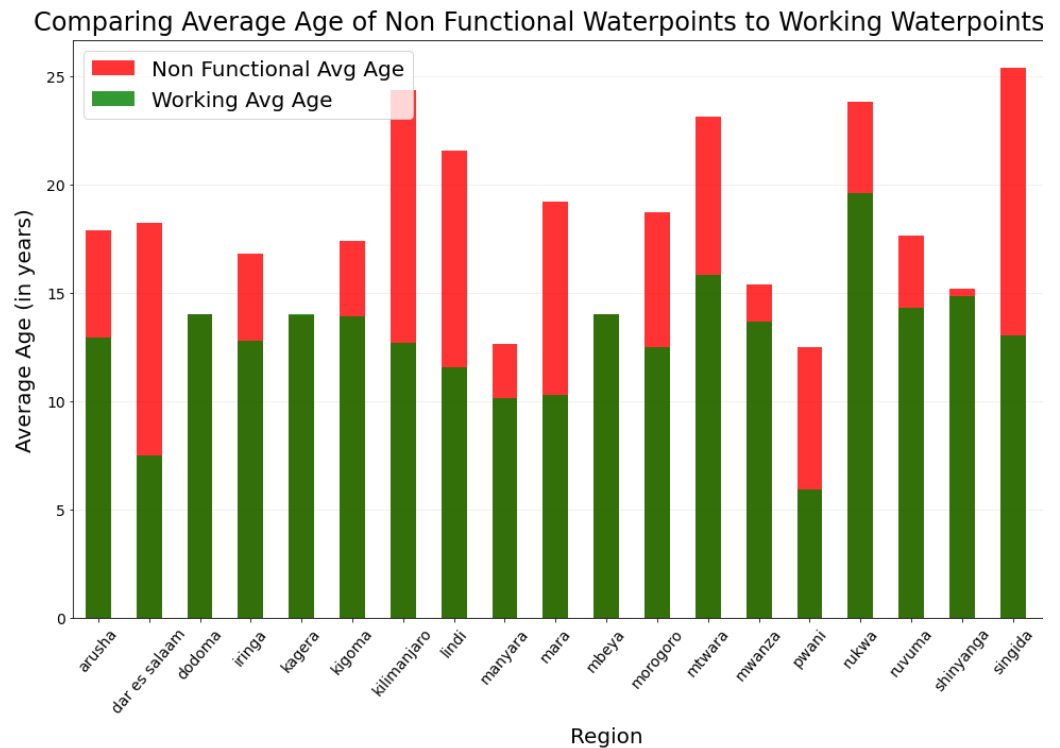
3b: Difference between average age of Functional Needs Repair and Functional waterpoints?

```
In [13]: 1 #df.groupby('region')['construction_year'].max().sort_values()
```

```
In [14]: working_by_region_mean_age = df[df['status_group'] != 'non functional'].groupby('region')['waterpoint_age'].mean()
non_functional_by_region_mean_age = df[df['status_group'] == 'non functional'].groupby('region')['waterpoint_age'].mean()
```

```
In [15]: 1 ax = non_functional_by_region_mean_age.plot(fontsize=14, kind='bar', alpha=.8, rot=50, color='red')
2         working_by_region_mean_age.plot(kind='bar', rot=50, alpha=.8, ax=ax, color="green")
3
4         ax.set_axisbelow(True)
5         ax.yaxis.grid(True, color='#EEEEEE')
6
7         # Add axis and chart labels.
8         ax.set_xlabel('Region', labelpad=10)
9         ax.set_ylabel('Average Age (in years)', labelpad=10)
10        ax.set_title('Comparing Average Age of Non Functional Waterpoints to Working Waterpoints by Region')
11        ax.legend(['Non Functional Avg Age', 'Working Avg Age'])
12
```

Out[15]: <matplotlib.legend.Legend at 0x7f2671cd2048>



```
In [16]: 1 # Get the difference between the average age of all working (functional and functional needs repair)
2 # and non functional waterpoints by Region
3 means_zipped = zip(non_functional_by_region_mean_age.values, working_by_region_mean_age.values)
4 age_diffs = []
5
6 for item in means_zipped:
7     diff_non_to_functional = round(item[0] - item[1], 3)
8     age_diffs.append(diff_non_to_functional)
9
10 labeled_age_diffs = zip(working_by_region_mean_age.index, age_diffs)
11
12 for tup in labeled_age_diffs:
13     print(tup[0].capitalize(), ': Diff between average Working and Non-Functional waterpoints: ', '
```

Arusha : Diff between average Working and Non-Functional waterpoints: 4.945 years
Dar es salaam : Diff between average Working and Non-Functional waterpoints: 10.753 years
Dodoma : Diff between average Working and Non-Functional waterpoints: 0.0 years
Iringa : Diff between average Working and Non-Functional waterpoints: 4.004 years
Kagera : Diff between average Working and Non-Functional waterpoints: -0.013 years
Kigoma : Diff between average Working and Non-Functional waterpoints: 3.436 years
Kilimanjaro : Diff between average Working and Non-Functional waterpoints: 11.698 years
Lindi : Diff between average Working and Non-Functional waterpoints: 10.02 years
Manyara : Diff between average Working and Non-Functional waterpoints: 2.467 years
Mara : Diff between average Working and Non-Functional waterpoints: 8.913 years
Mbeya : Diff between average Working and Non-Functional waterpoints: -0.001 years
Morogoro : Diff between average Working and Non-Functional waterpoints: 6.227 years
Mtwara : Diff between average Working and Non-Functional waterpoints: 7.275 years
Mwanza : Diff between average Working and Non-Functional waterpoints: 1.732 years
Pwani : Diff between average Working and Non-Functional waterpoints: 6.551 years
Rukwa : Diff between average Working and Non-Functional waterpoints: 4.22 years
Ruvuma : Diff between average Working and Non-Functional waterpoints: 3.319 years
Shinyanga : Diff between average Working and Non-Functional waterpoints: 0.37 years
Singida : Diff between average Working and Non-Functional waterpoints: 12.359 years
Tabora : Diff between average Working and Non-Functional waterpoints: -0.125 years
Tanga : Diff between average Working and Non-Functional waterpoints: 2.492 years

Insights (Funcional vs Non Funcional):

The average age of Non Functional waterpoints is greater than the average age of Working waterpoints in 17 out of 21 Regions. Average age same in 4 out of 21 Regions. This measurement supports intuition that older waterpoints are more likely to be non functional.

Recommendation:

Regions that have Working waterpoints with average age close to the average age of Non Functional waterpoints should consider increasing monitoring/maintenance of Working waterpoints and/or new waterpoint installation.

Q3b: Difference between average age of Functional Needs Repair and Functional waterpoints?

```
In [17]: 1 functional_by_region_mean_age = df[df['status_group'] == 'functional'].groupby('region')['waterpoi
2 functional_needs_repair_by_region_mean_age = df[df['status_group'] == 'functional needs repair'].g
```



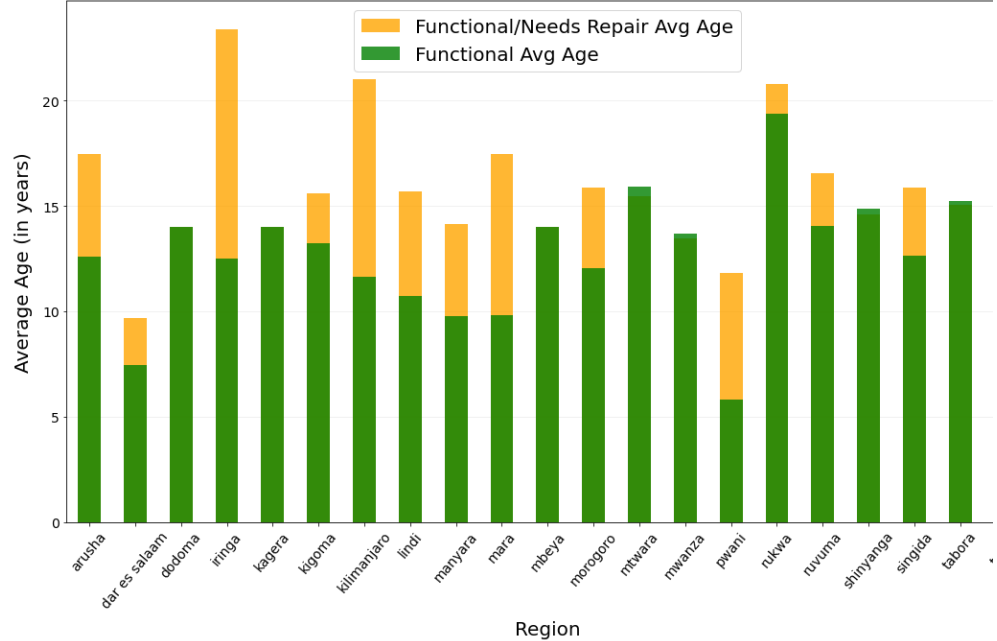
```

In [18]: 1 ax = functional_needs_repair_by_region_mean_age.plot(fontsize=14, kind='bar', alpha=.8, rot=50, co
2 functional_by_region_mean_age.plot(kind='bar', rot=50, alpha=.8, ax=ax, color="green")
3
4 ax.set_axisbelow(True)
5 ax.yaxis.grid(True, color='EEEEEE')
6
7 # Add axis and chart labels.
8 ax.set_xlabel('Region', labelpad=10)
9 ax.set_ylabel('Average Age (in years)', labelpad=10)
10 ax.set_title('Comparing Average Age of Functional/Needs Repair Waterpoints to Functional Waterpoint
11 ax.legend(['Functional/Needs Repair Avg Age', 'Functional Avg Age'])

```

Out[18]: <matplotlib.legend.Legend at 0x7f2671bd84e0>

Comparing Average Age of Functional/Needs Repair Waterpoints to Functional Waterpoint



```

In [19]: 1 # Get the difference between the average age of all Functional needs repair and Functional waterpo
2 means_zipped = zip(functional_needs_repair_by_region_mean_age.values, functional_by_region_mean_age
3 age_diffs = []
4
5 for item in means_zipped:
6     diff_non_to_functional = round(item[0] - item[1], 2)
7     age_diffs.append(diff_non_to_functional)
8
9 functioning_labeled_age_diffs = zip(functional_needs_repair_by_region_mean_age.index, age_diffs)
10
11 for tup in functioning_labeled_age_diffs:
12     print(tup[0].capitalize(), ': Diff between average age of Functional Needs Repair and Function.

```

```

Arusha : Diff between average age of Functional Needs Repair and Functional waterpoints:  4.9 years
Dar es salaam : Diff between average age of Functional Needs Repair and Functional waterpoints:  2.2 y
Dodoma : Diff between average age of Functional Needs Repair and Functional waterpoints:  0.0 years
Iringa : Diff between average age of Functional Needs Repair and Functional waterpoints:  10.88 years
Kagera : Diff between average age of Functional Needs Repair and Functional waterpoints:  0.0 years
Kigoma : Diff between average age of Functional Needs Repair and Functional waterpoints:  2.37 years
Kilimanjaro : Diff between average age of Functional Needs Repair and Functional waterpoints:  9.36 ye
Lindi : Diff between average age of Functional Needs Repair and Functional waterpoints:  4.97 years
Manyara : Diff between average age of Functional Needs Repair and Functional waterpoints:  4.37 years
Mara : Diff between average age of Functional Needs Repair and Functional waterpoints:  7.64 years
Mbeya : Diff between average age of Functional Needs Repair and Functional waterpoints:  0.0 years
Morogoro : Diff between average age of Functional Needs Repair and Functional waterpoints:  3.83 years
Mtwara : Diff between average age of Functional Needs Repair and Functional waterpoints:  -0.46 years
Mwanza : Diff between average age of Functional Needs Repair and Functional waterpoints:  -0.26 years
Pwani : Diff between average age of Functional Needs Repair and Functional waterpoints:  6.03 years
Rukwa : Diff between average age of Functional Needs Repair and Functional waterpoints:  1.41 years
Ruvuma : Diff between average age of Functional Needs Repair and Functional waterpoints:  2.5 years
Shinyanga : Diff between average age of Functional Needs Repair and Functional waterpoints:  -0.26 yea
Singida : Diff between average age of Functional Needs Repair and Functional waterpoints:  3.19 years
Tabora : Diff between average age of Functional Needs Repair and Functional waterpoints:  -0.17 years
Tanga : Diff between average age of Functional Needs Repair and Functional waterpoints:  6.68 years

```

Insights (Functional Needs Repair/Functional):

The average age of Functional Needs Repair waterpoints is greater than Functional waterpoints in 14 out of 21 Regions. The average is the same in 7 out of 21 Regions. This measurement supports intuition that waterpoints in need of repair would tend to be older.

Recommendation:

Regions that have Functional waterpoints with average age approaching to average age of waterpoints that Need Repair should consider increased monitoring/preventative maintenance of waterpoints.

Future Work questions**For All Working waterpoints (Functioning and Functioning Needs Repair):**

- what are payment types by Region?
- what are the recorded populations served by Region?
- what are the source_classes or source_type?
- who are the installers?
- who are the management groups?

For ALL waterpoints

- who are the management groups Grouped By Region, then Status?
- who are the installers Grouped By Region, then Status?

For Non Functioning waterpoints:

- who are the installers?
- who are the management groups?

Additional questions:

How many people use waterpoints? Entire country sum, breakdown by region, water basin.

For lower QUALITY waterpoints, Defined as quality_group anything other than 'good' and/or quantity_group listed as 'insufficient', what is the location? Where are they located (region, waterbasin, lat/long)

For lower QUANTITY waterpoints, Defined as quantity_group listed as 'insufficient', what is the population count, where are they located (region, lat/long)

What do we know about the payment types for waterpoints? Breakdown by country, region water basin, population bins

What do we know about waterpoint age? Country wide, water basin wide, region wide? Min, Max, Median, Mean, Dist?

How up-to-date is the data for waterpoints by region? Are there trends in missing values? Are some regions missing more data than others?

What do we know about waterpoints of unknown age? What is the population served by waterpoints of unknown age? How many, where are they located (region, waterbasin, lat/long)

What do we know about the waterpoint_type_group? Breakdown by country, region, water basin, bin by population. (communal standpipe 'h' 'improved spring' 'cattle trough' 'dam')

What is the breakdown of Orgs that perform management for the waterpoints? Any managed that don't have a permit or permit status unknown? management column not scheme_management (scheme indicates the funding mechanism, I think) Is there a difference between scheme_management entities for waterpoints? If so, what does that look like? 13 total: 'vwc' 'wug' 'other' 'private operator' 'water board' 'wua' 'company' 'parastatal' 'unknown' 'other - school' 'trust'

What is the breakdown of extraction type /extraction type class by country, region, water basin.

In []:

1