

SURVEY ON WOMEN PARTICIPATION IN NIGERIA POLITICS

Abstract

A huge level of violence affects women in Nigerian politics and several research have found that men and women experience electoral violence differently, albeit men still participate actively in elections. This research seeks to examine how family dynamics influences women's political participation in the face of violence. We employed a survey to ask **791 respondents** if they would allow female relatives to participate in violent elections. This was analysed logistic regression model to test the variables and examine how electoral violence affects women's political participation in elections. Python programming language was used for the implementation of the model. Findings reveal that electoral violence affects women's participation in political activities irrespective of their educational qualification. The RandomForest model employed for predictive experiment gives **accuracy is 94.4%, precision is 94.6%, recall is 94.4%, and F1-score is 94.5%**.

Importing necessary libraries

```
In [1]: 1 # Installing the Libraries with the specified version.
2 # uncomment and run the following line if Google Colab is being used
3 # !pip install scikit-learn==1.2.2 seaborn==0.13.1 matplotlib==3.7.1 numpy==1.25.2 pandas==1.5.3 imbalanced-Learn==0.10.1 xgboost==2.0.1
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
```

```
In [2]: 1 # Installing the Libraries with the specified version.
2 # uncomment and run the following lines if Jupyter Notebook is being used
3 # !pip install scikit-learn
4 # !pip install seaborn
5 # !pip install matplotlib
6 # !pip install numpy
7 # !pip install pandas
8 # !pip install imblearn
9 # !pip install xgboost -q --user
10 # !pip install --upgrade -q threadpoolctl
11 # !pip install scikit-plot
```

Note: After running the above cell, kindly restart the notebook kernel and run all cells sequentially from the start again.

```
In [3]: 1 # Import necessary libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import warnings
7
8 # To help with reading and manipulation of data
9 pd.set_option("display.max_columns", None) # Remove Limit for displayed columns
10 pd.set_option("display.max_rows", 200) # Set limit for displayed rows
11
12 # To split the data
13 from sklearn.model_selection import train_test_split
14
15 # To impute missing values
16 from sklearn.impute import SimpleImputer, KNNImputer
17
18 # To preprocess the data for modelling
19 from sklearn.preprocessing import StandardScaler, OrdinalEncoder, OneHotEncoder
20
21 # To build a Random Forest Classifier
22 from sklearn.ensemble import RandomForestClassifier
23 from sklearn.model_selection import cross_val_score
24
25 # To tune a model
26 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
27
28 # To create a pipeline for production
29 from sklearn.pipeline import Pipeline, make_pipeline
30 from sklearn.compose import ColumnTransformer, make_column_selector as selector
31
32 # To get different performance metrics
33 from sklearn.metrics import (
34     classification_report,
35     confusion_matrix,
36     recall_score,
37     accuracy_score,
38     precision_score,
39     f1_score,
40     roc_curve, # This was previously duplicated
41     auc,
42 )
43
44 # Import itertools for ROC visualization
45 from itertools import cycle
46 from sklearn.preprocessing import label_binarize
47
48 # Import necessary libraries for performance visualization
49 from scikitplot.metrics import plot_confusion_matrix
50
51 # Ignore warnings
52 warnings.filterwarnings("ignore")
```

Loading Data

```
In [4]: 1 churn_dataset = pd.read_csv("survey_response.csv")
```

```
In [5]: 1 # Checking the number of rows and columns in the data
2 churn_dataset.shape
```

Out[5]: (791, 14)

- The dataset has 791 rows and 14 columns

Data Overview

- Observations
- Sanity checks

```
In [6]: 1 # Let's create a copy of the data
2 data = churn_dataset.copy()
```

```
In [7]: 1 # Let's view the first 5 rows of the data
2 data.head()
```

Out[7]:

	Timestamp	Gender	Work Sector	Educational Qualification	Age range	Do you have a permanent voters card?	Are you likely to vote when there is electoral violence around you?	Are you likely to prevent a "female" loved one from going to vote after violence occurs?	Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?	In your opinion, does violence impact the confidence of women in engaging in political activities?	Did you vote in 2023 General Elections?	If No, why not?	Have you ever witnessed any form of electoral violence during elections?	Have you ever witnessed any form of harassment on social media?
0	04/04/2024 17:57	Female	Formal Sector (9-5 jobs, Professionals, Hybrid...)	HND, B.Sc.	18-30	No	No	No	No	No	No	No PVC	No	No
1	04/04/2024 18:05	Female	Formal Sector (9-5 jobs, Professionals, Hybrid...)	HND, B.Sc.	18-30	No	No	Yes	Yes	Yes	No	Others	Yes	Yes
2	04/04/2024 18:06	Female	Informal Sector (Artisans, Traders)	HND, B.Sc.	31-40	Yes	No	Yes	Yes	Yes	No	Unavailable (distance, health issues)	No	No
3	04/04/2024 18:08	Male	Formal Sector (9-5 jobs, Professionals, Hybrid...)	HND, B.Sc.	18-30	No	No	Yes	Yes	Yes	No	No PVC	Yes	No
4	04/04/2024 18:10	Female	Informal Sector (Artisans, Traders)	HND, B.Sc.	18-30	Yes	No	Yes	Yes	Yes	Yes	Others	Yes	Yes

```
In [8]: 1 # Let's view the last 5 rows of the data
2 data.tail()
```

Out[8]:

	Timestamp	Gender	Work Sector	Educational Qualification	Age range	Do you have a permanent voters card?	Are you likely to vote when there is electoral violence around you?	Are you likely to prevent a "female" loved one from going to vote after violence occurs?	Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?	In your opinion, does violence impact the confidence of women in engaging in political activities?	Did you vote in 2023 General Elections?	If No, why not?	Have you ever witnessed any form of electoral violence during elections?	Have you ever witnessed any form of harassment on social media?
786	NaN	Female	Formal Sector (9-5 jobs, Professionals, Hybrid...)	SSCE and below	31-40	Yes	No	Yes	Uncertain	Yes	Yes	Others	Uncertain	Yes
787	NaN	Female	Informal Sector (Artisans, Traders)	HND, B.Sc.	60 and above	No	No	Yes	No	No	No	Work (Journalist, Health officials, Security a...)	Yes	No
788	NaN	Male	Formal Sector (9-5 jobs, Professionals, Hybrid...)	Postgraduate	31-40	No	Uncertain	Yes	Yes	Yes	No	No PVC	Yes	No
789	NaN	Male	Informal Sector (Artisans, Traders)	HND, B.Sc.	51-60	Yes	No	Yes	Uncertain	Uncertain	No	Others	No	No
790	NaN	Female	Formal Sector (9-5 jobs, Professionals, Hybrid...)	HND, B.Sc.	18-30	Yes	Yes	Yes	Yes	Yes	Yes	Others	Yes	No

```
In [9]: 1 # Timestamp consists of uniques ID for clients and hence will not add value to the modeling  
2 data.drop(["Timestamp"], axis=1, inplace=True)
```

```
In [10]: 1 data.columns.tolist()
```

```
Out[10]: ['Gender',  
          'Work Sector',  
          'Educational Qualification',  
          'Age range',  
          'Do you have a permanent voters card?',  
          'Are you likely to vote when there is electoral violence around you?',  
          'Are you likely to prevent a "female" loved one from going to vote after violence occurs?',  
          'Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?',  
          'In your opinion, does violence impact the confidence of women in engaging in political activities?',  
          'Did you vote in 2023 General Elections?',  
          'If No, why not?',  
          'Have you ever witnessed any form of electoral violence during elections?',  
          'Have you ever witnessed any form of harassment on social media?']
```

Data Type Conversions

```
In [11]: 1 # Let's view the statistical summary of the numerical columns in the data  
2 data.describe().T
```

```
Out[11]:
```

	count	unique	top	freq
Gender	791	2	Female	533
Work Sector	791	2	Formal Sector (9-5 jobs, Professionals, Hybrid...)	492
Educational Qualification	791	7	HND, B.Sc.	438
Age range	791	6	18-30	380
Do you have a permanent voters card?	791	2	Yes	524
Are you likely to vote when there is electoral violence around you?	791	3	No	604
Are you likely to prevent a "female" loved one from going to vote after violence occurs?	791	3	Yes	691
Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?	791	3	Yes	603
In your opinion, does violence impact the confidence of women in engaging in political activities?	791	3	Yes	552
Did you vote in 2023 General Elections?	791	2	No	586
If No, why not?	791	5	Others	273
Have you ever witnessed any form of electoral violence during elections?	791	3	Yes	395
Have you ever witnessed any form of harassment on social media?	791	3	Yes	402

Observations:

- Dataset: the data has no missing value
- Data type: all column has object has the datatype
- Columns: all the columns have categorical values

```
In [12]: 1 for i in data.describe(include=["object"]).columns:  
2     print("Unique values in", i, "are :")  
3     print(data[i].value_counts())  
4     print("*" * 50)  
5     print("*" * 50)
```

```

Unique values in Gender are :
Female      533
Male        258
Name: Gender, dtype: int64
*****
***** Unique values in Work Sector are :
Formal Sector (9-5 jobs, Professionals, Hybrid jobs)    492
Informal Sector (Artisans, Traders)                      299
Name: Work Sector, dtype: int64
*****
***** Unique values in Educational Qualification are :
HND, B.Sc.          438
Postgraduate        145
SSCE and below      140
ND, NCE             39
Mbbis in view       12
Btech               9
Undergraduate       8
Name: Educational Qualification, dtype: int64
*****
***** Unique values in Age range are :
18-30              380
31-40              173
51-60              116
60 and above       71
41-50              50
51-61              1
Name: Age range, dtype: int64
*****
***** Unique values in Do you have a permanent voters card? are :
Yes      524
No       267
Name: Do you have a permanent voters card?, dtype: int64
*****
***** Unique values in Are you likely to vote when there is electoral violence around you? are :
No       604
Yes     140
Uncertain   47
Name: Are you likely to vote when there is electoral violence around you?, dtype: int64
*****
***** Unique values in Are you likely to prevent a "female" loved one from going to vote after violence occurs? are :
Yes      691
No       91
Uncertain   9
Name: Are you likely to prevent a "female" loved one from going to vote after violence occurs?, dtype: int64
*****
***** Unique values in Do you believe that violence deters women from participating in political activities such as rallies, elections, and camp
aigns? are :
Yes      603
No       132
Uncertain   56
Name: Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?, dt
pe: int64
*****
***** Unique values in In your opinion, does violence impact the confidence of women in engaging in political activities? are :
Yes      552
No       193
Uncertain   46
Name: In your opinion, does violence impact the confidence of women in engaging in political activities?, dtype: int64
*****
***** Unique values in Did you vote in 2023 General Elections? are :
No      586
Yes     205
Name: Did you vote in 2023 General Elections?, dtype: int64
*****
***** Unique values in If No, why not? are :
Others            273
No PVC            218
Unavailable (distance, health issues)                  159
Electoral Violence                     83
Work (Journalist, Health officials, Security agents, Electoral officers)  58
Name: If No, why not?, dtype: int64
*****
***** Unique values in Have you ever witnessed any form of electoral violence during elections? are :
Yes      395
No       343
Uncertain   53
Name: Have you ever witnessed any form of electoral violence during elections?, dtype: int64
*****
***** Unique values in Have you ever witnessed any form of harassment on social media? are :
Yes      402
No       364
Uncertain   25
Name: Have you ever witnessed any form of harassment on social media?, dtype: int64
*****

```

Lets regroup the Education Qualification in to four categories

- Tertiary (college or university)
- Post-Secondary/Vocational
- SSCE and Below
- Postgraduate

```
In [13]: 1 data["Educational Qualification"].replace("Mbbs in view", "Tertiary (college or university)", inplace=True)
2 data["Educational Qualification"].replace("Btech", "Tertiary (college or university)", inplace=True)
3 data["Educational Qualification"].replace("HND, B.Sc.", "Tertiary (college or university)", inplace=True)
4 data["Educational Qualification"].replace("ND, NCE", "Post-Secondary/Vocational", inplace=True)
5 data["Educational Qualification"].replace("Undergraduate ", "Tertiary (college or university)", inplace=True)
```

```
In [14]: 1 print("*" * 50)
2 print("Unique values in Educational Qualification are :")
3 print(data["Educational Qualification"].value_counts())
4 print("*" * 50)
```

```
*****
Unique values in Educational Qualification are :
Tertiary (college or university)    467
Postgraduate                      145
SSCE and below                     140
Post-Secondary/Vocational          39
Name: Educational Qualification, dtype: int64
*****
```

Observation

- there is no missing value

Exploratory Data Analysis (EDA)

- EDA is an important part of this project inorder to reveal hidden information from the data.
- It is important to investigate and understand the data better before building a model with it.

Some of the Questions Answered through the EDA:

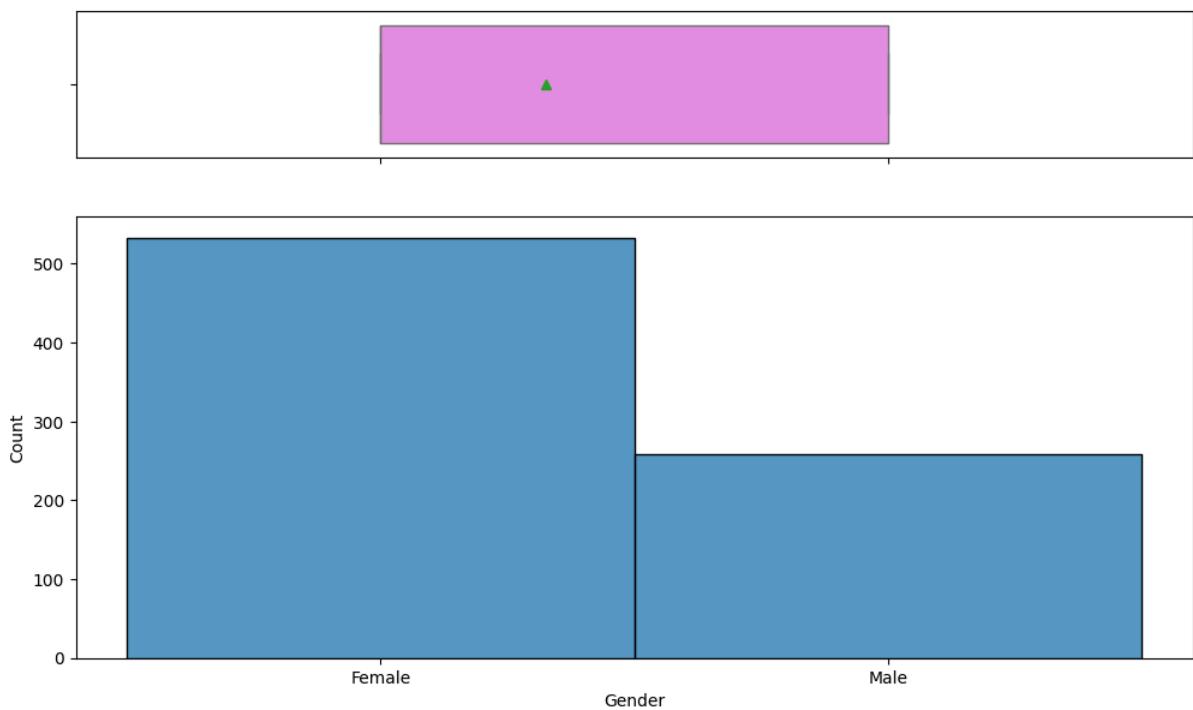
1. How is the gender distributed?
2. What is the distribution of the Educational Qualification?
3. What is the distribution of the permanent voters card?
4. How does the change in likelyhood to vote when there is electoral violence vary by the gender
5. How does the witness of any form of electoral violence during elections vary by the gender
6. How does the witnessing any form of harassment on social media vary by the gender
7. What are the attributes that have a strong correlation with each other?

Univariate analysis

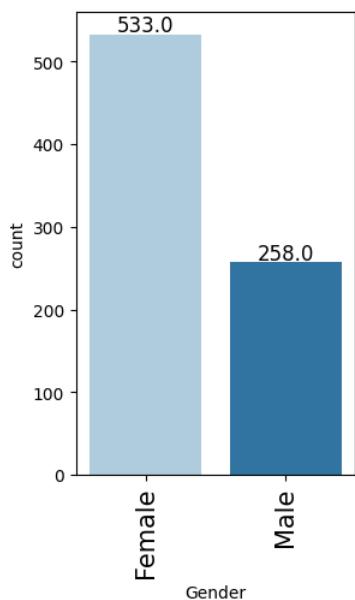
```
In [15]: 1 # function to plot a boxplot and a histogram along the same scale.
2 def histogram_boxplot(data, feature, figsize=(12, 7), kde=False, bins=None):
3     """
4     Boxplot and histogram combined
5
6     data: dataframe
7     feature: dataframe column
8     figsize: size of figure (default (12,7))
9     kde: whether to show density curve (default False)
10    bins: number of bins for histogram (default None)
11    """
12    f2, (ax_box2, ax_hist2) = plt.subplots(
13        nrows=2, # Number of rows of the subplot grid= 2
14        sharex=True, # x-axis will be shared among all subplots
15        gridspec_kw={"height_ratios": (0.25, 0.75)},
16        figsize=figsize,
17    ) # creating the 2 subplots
18    sns.boxplot(
19        data=data, x=feature, ax=ax_box2, showmeans=True, color="violet"
20    ) # boxplot will be created and a star will indicate the mean value of the column
21    sns.histplot(
22        data=data, x=feature, kde=kde, ax=ax_hist2, bins=bins, palette="winter"
23    ) if bins else sns.histplot(
24        data=data, x=feature, kde=kde, ax=ax_hist2
25    )
26
27 # function to create labeled barplots
28 def labeled_barplot(data, feature, perc=False, n=None):
29     """
30     Barplot with percentage at the top
31
32     data: dataframe
33     feature: dataframe column
34     perc: whether to display percentages instead of count (default is False)
35     n: displays the top n category levels (default is None, i.e., display all levels)
36     """
37
38     total = len(data[feature]) # Length of the column
39     count = data[feature].nunique()
40     if n is None:
41         plt.figure(figsize=(count + 1, 5))
42     else:
43         plt.figure(figsize=(n + 1, 5))
44
45     plt.xticks(rotation=90, fontsize=15)
46     ax = sns.countplot(
47         data=data,
48         x=feature,
49         palette="Paired",
50         order=data[feature].value_counts().index[:n].sort_values(),
51     )
52
53     for p in ax.patches:
54         if perc == True:
55             label = "{:.1f}%".format(
56                 100 * p.get_height() / total
57             ) # percentage of each class of the category
58         else:
59             label = p.get_height() # count of each Level of the category
60
61         x = p.get_x() + p.get_width() / 2 # width of the plot
62         y = p.get_height() # height of the plot
63
64         ax.annotate(
65             label,
66             (x, y),
67             ha="center",
68             va="center",
69             size=12,
70             xytext=(0, 5),
71             textcoords="offset points",
72         ) # annotate the percentage
73
74     plt.show() # show the plot
```

Observations on Gender

```
In [16]: 1 histogram_boxplot(data, "Gender")
```



```
In [17]: 1 labeled_barplot(data, "Gender")
```

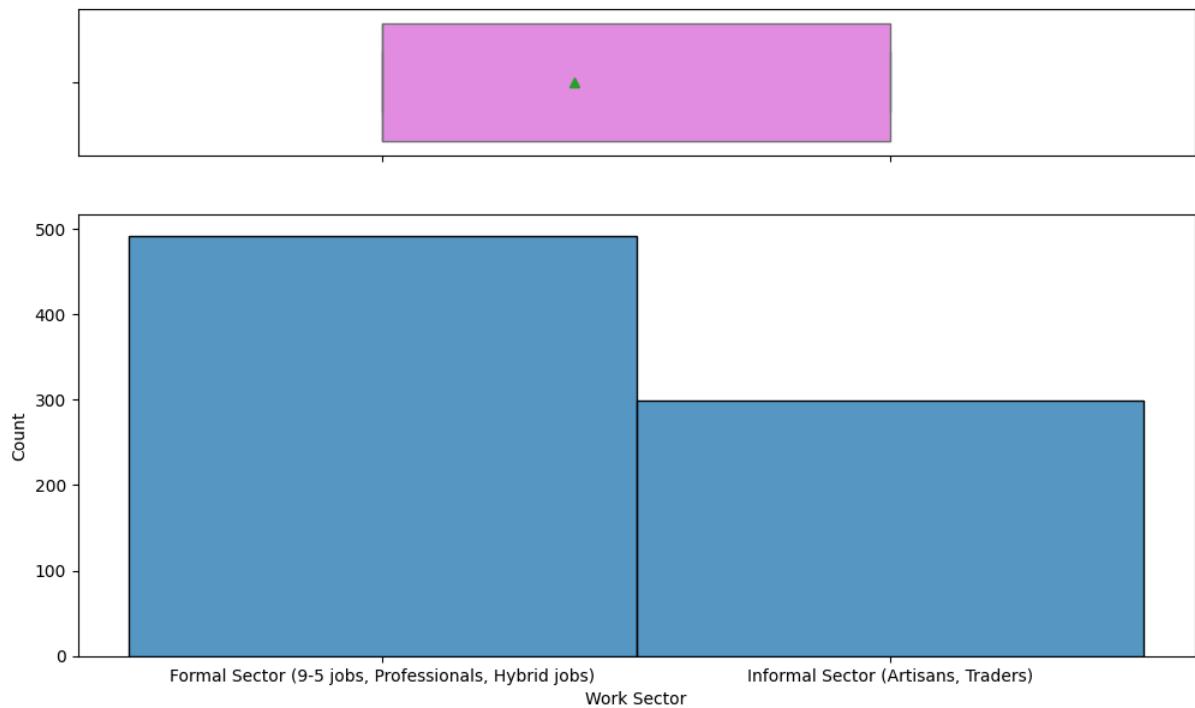


Observation

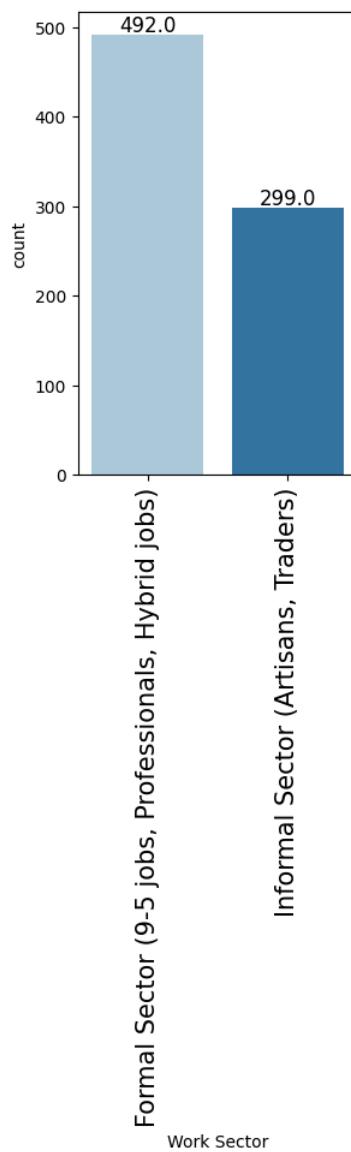
The are more female respondent than male

Observations on Work Sector

```
In [18]: 1 histogram_boxplot(data, "Work Sector")
```



```
In [19]: 1 labeled_barplot(data, "Work Sector")
```

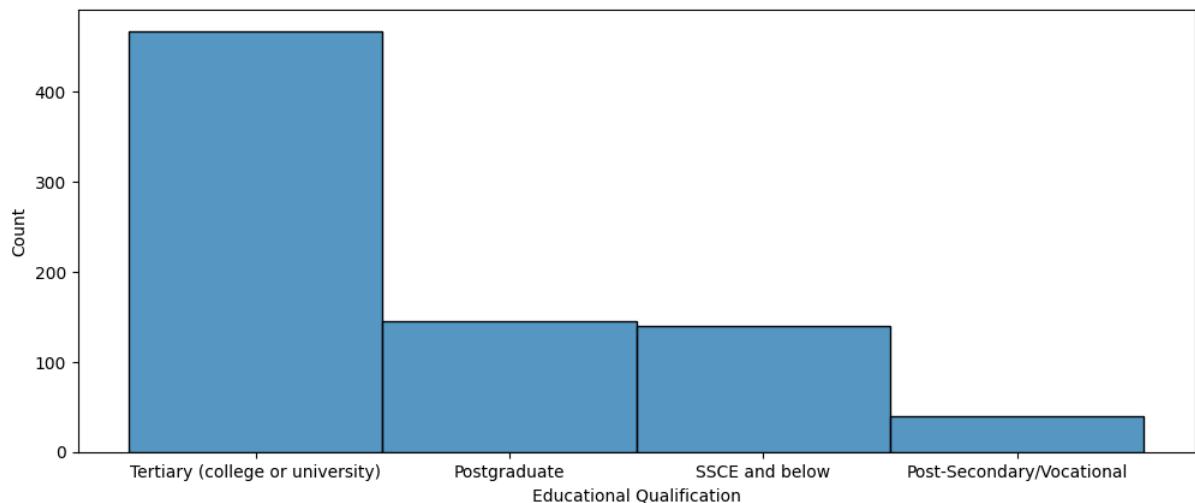
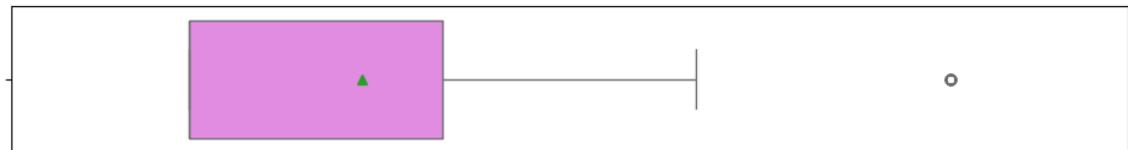


Observation

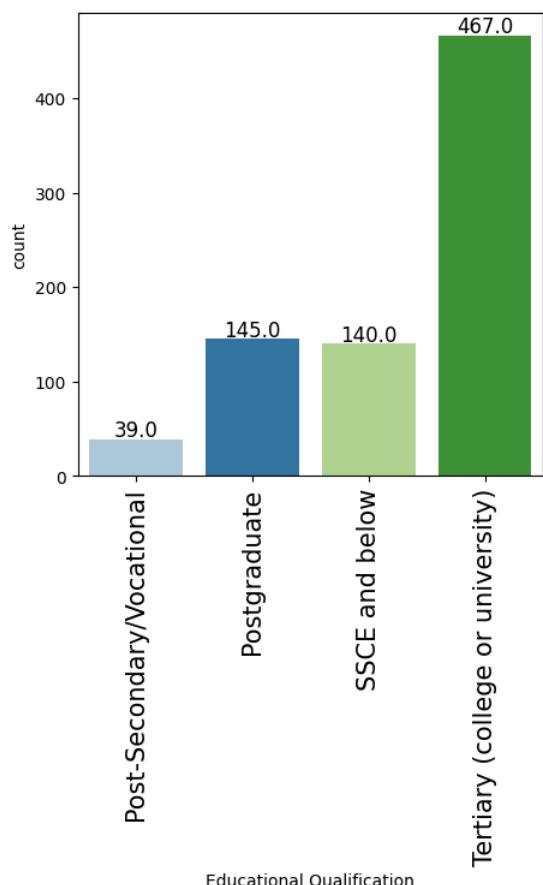
- Most of the respondent work in formal sector. That is the respondent is either works,
 1. on a 9:00AM to 5:00Pm schedule or
 2. as a professionals
 3. on jobs with hybrid mode
- Few Informal sectors (such as Artisans and Traders) responded.

Observations on Educational Qualification

```
In [20]: 1 histogram_boxplot(data, "Educational Qualification")
```



```
In [21]: 1 labeled_barplot(data, "Educational Qualification")
```

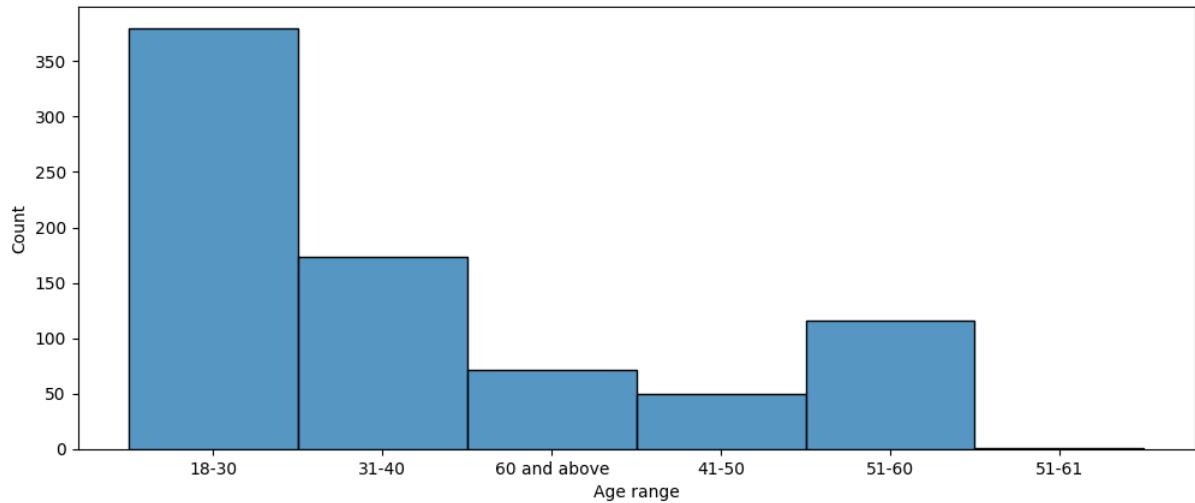
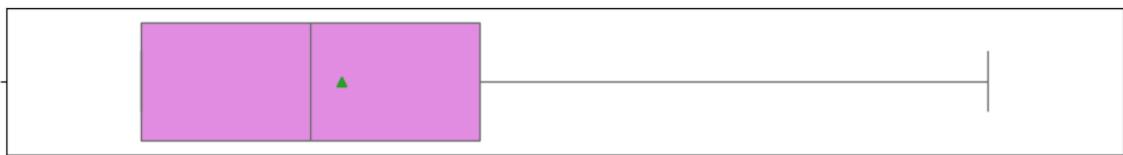


Observation

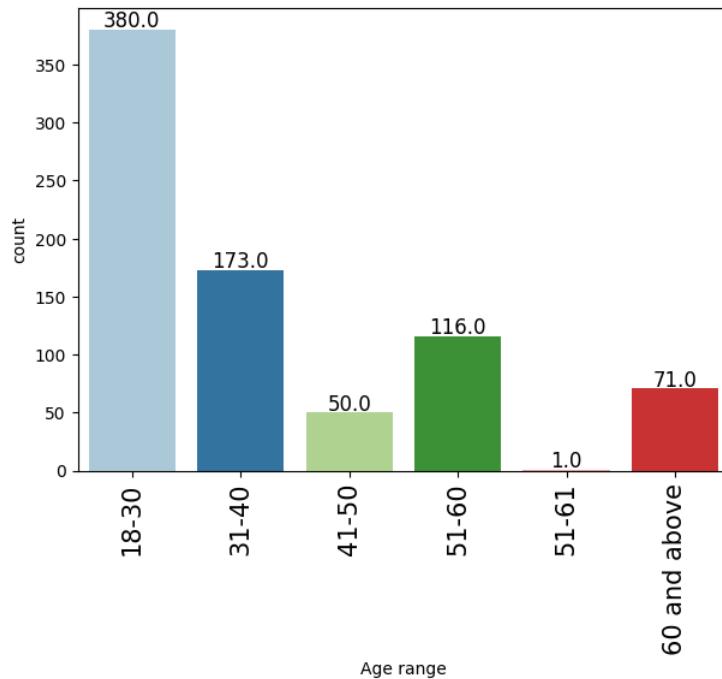
- More than 400 respondent holds HND or BSC Degree.
- Less than 200 respondent hold postgraduate degree.
- The respondent with ND or NCE or MBBS in View or BTech or Undergraduate are all less than 100.

Observations on Age Range

```
In [22]: 1 histogram_boxplot(data, "Age range")
```



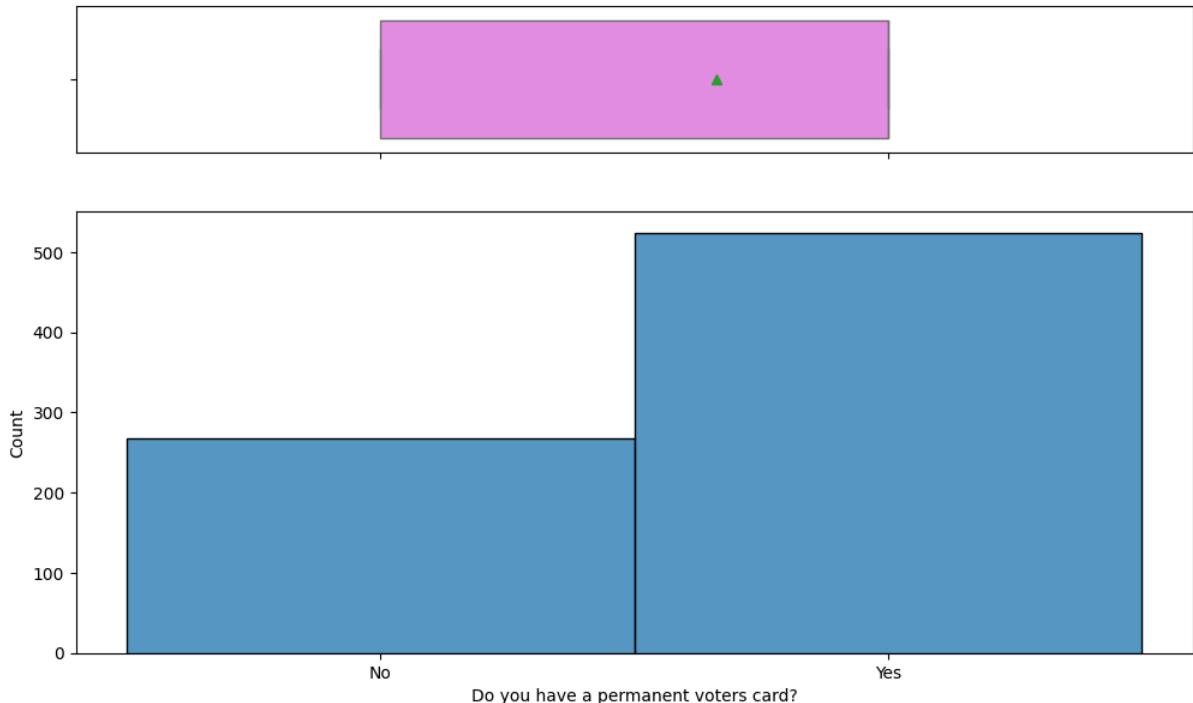
```
In [23]: 1 labeled_barplot(data, "Age range")
```



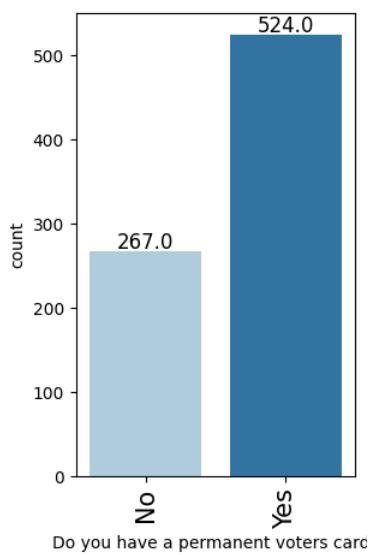
- 18-30 is the age range with the highest number of respondent follow by 31-40.
- 41-50 is the age range with the lowest number of respondent.

Observations on 'Respondent with permanent voters card?'

```
In [24]: 1 histogram_boxplot(data, "Do you have a permanent voters card?")
```



```
In [25]: 1 labeled_barplot(data, "Do you have a permanent voters card?")
```

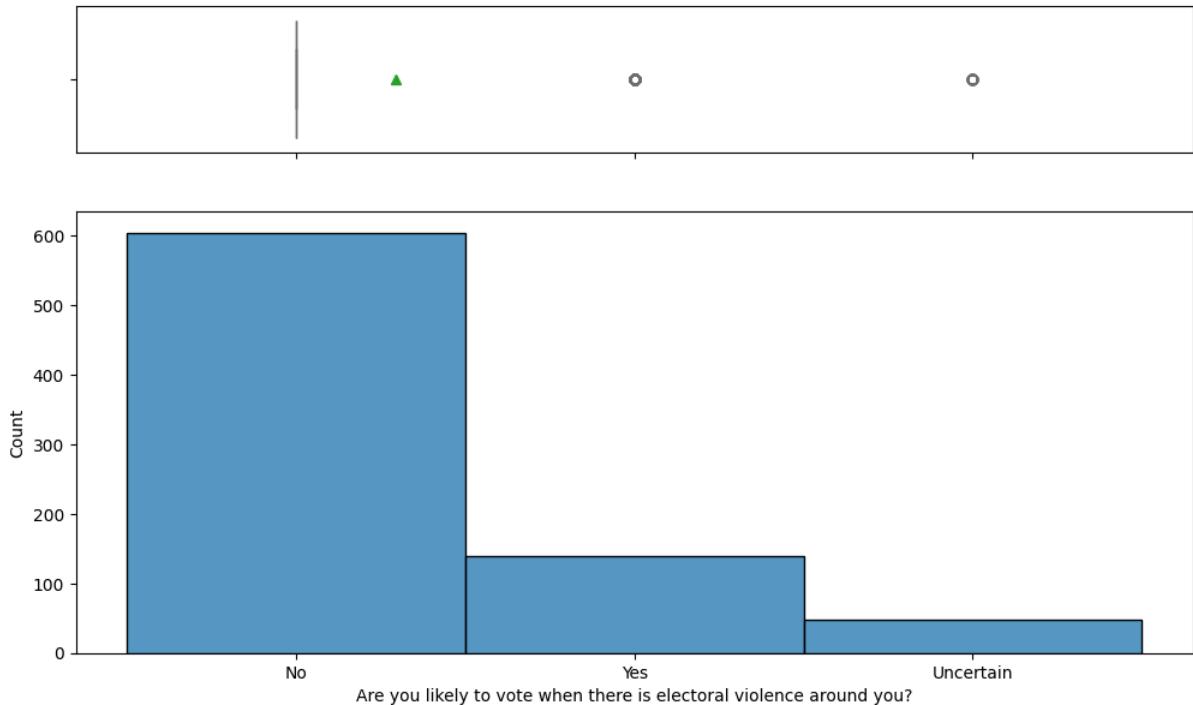


Observation

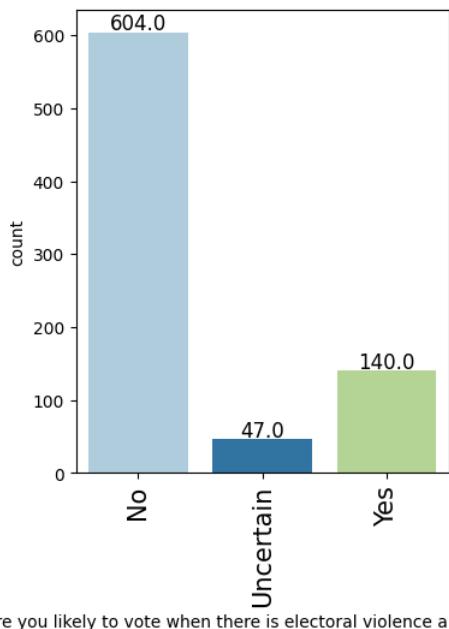
- About 70% of the respondent have permanent voters card

Observations on Respondent that are likely to vote when there is electoral violence around you?

```
In [26]: 1 histogram_boxplot(data, "Are you likely to vote when there is electoral violence around you?")
```



```
In [27]: 1 labeled_barplot(data, "Are you likely to vote when there is electoral violence around you?")
```



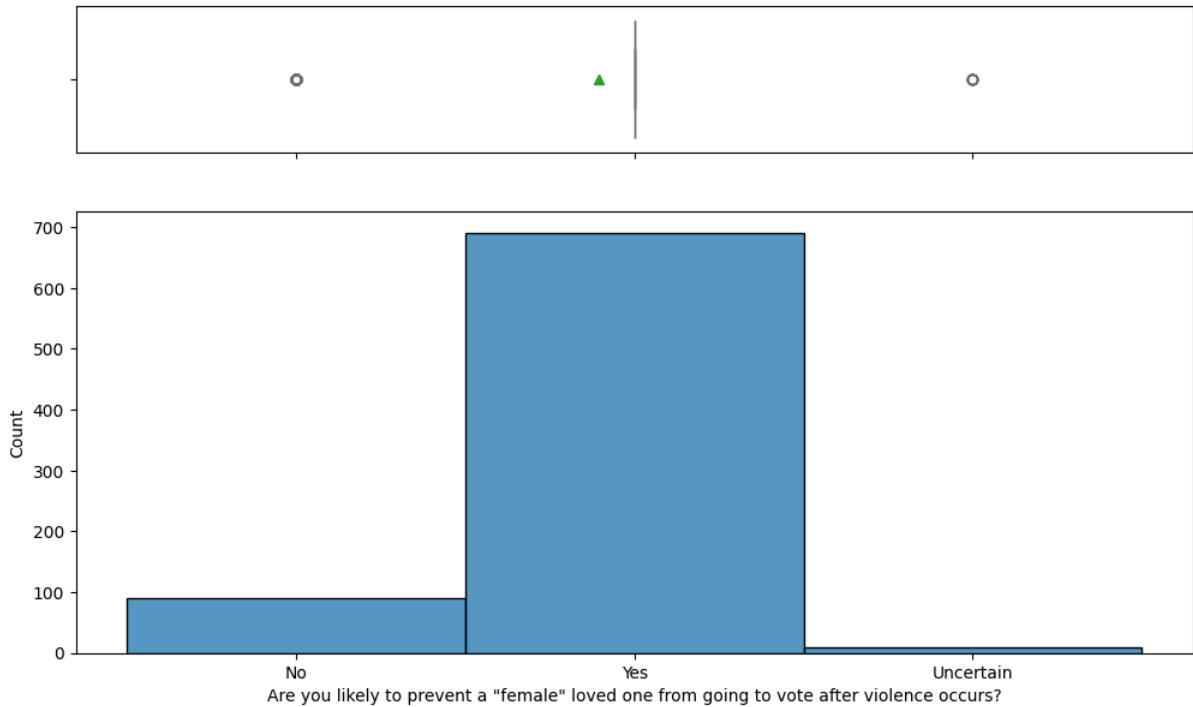
re you likely to vote when there is electoral violence arc

Observation

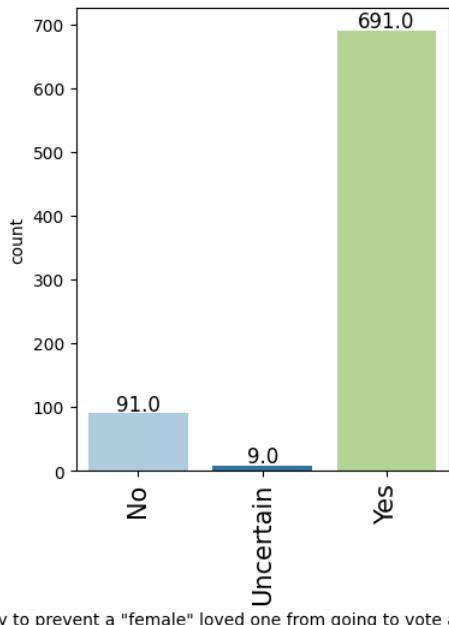
- around 80% of the respondent are not likely to vote when there is electoral violence around.

Observations on Respondent that are likely to prevent a "female" loved one from going to vote after violence occurs?

```
In [28]: 1 histogram_boxplot(data, 'Are you likely to prevent a "female" loved one from going to vote after violence occurs?')
```



```
In [29]: 1 labeled_barplot(data, 'Are you likely to prevent a "female" loved one from going to vote after violence occurs?')
```

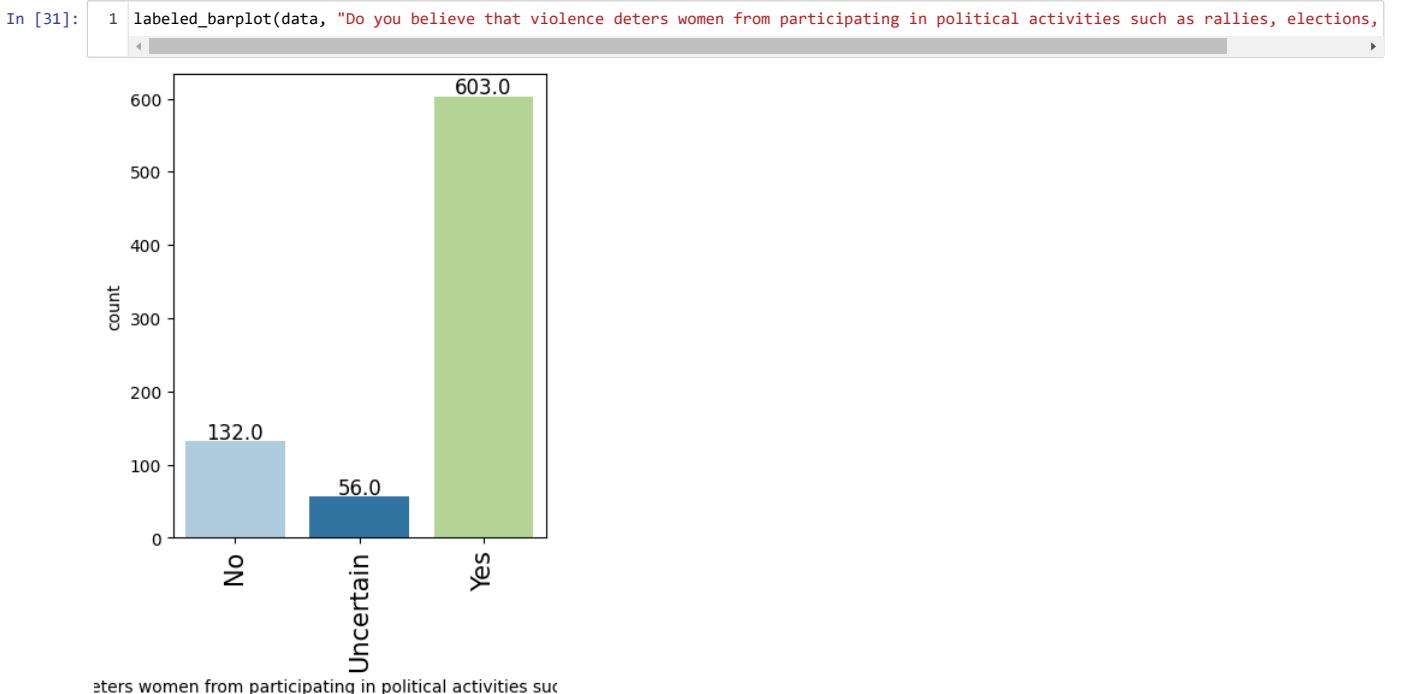
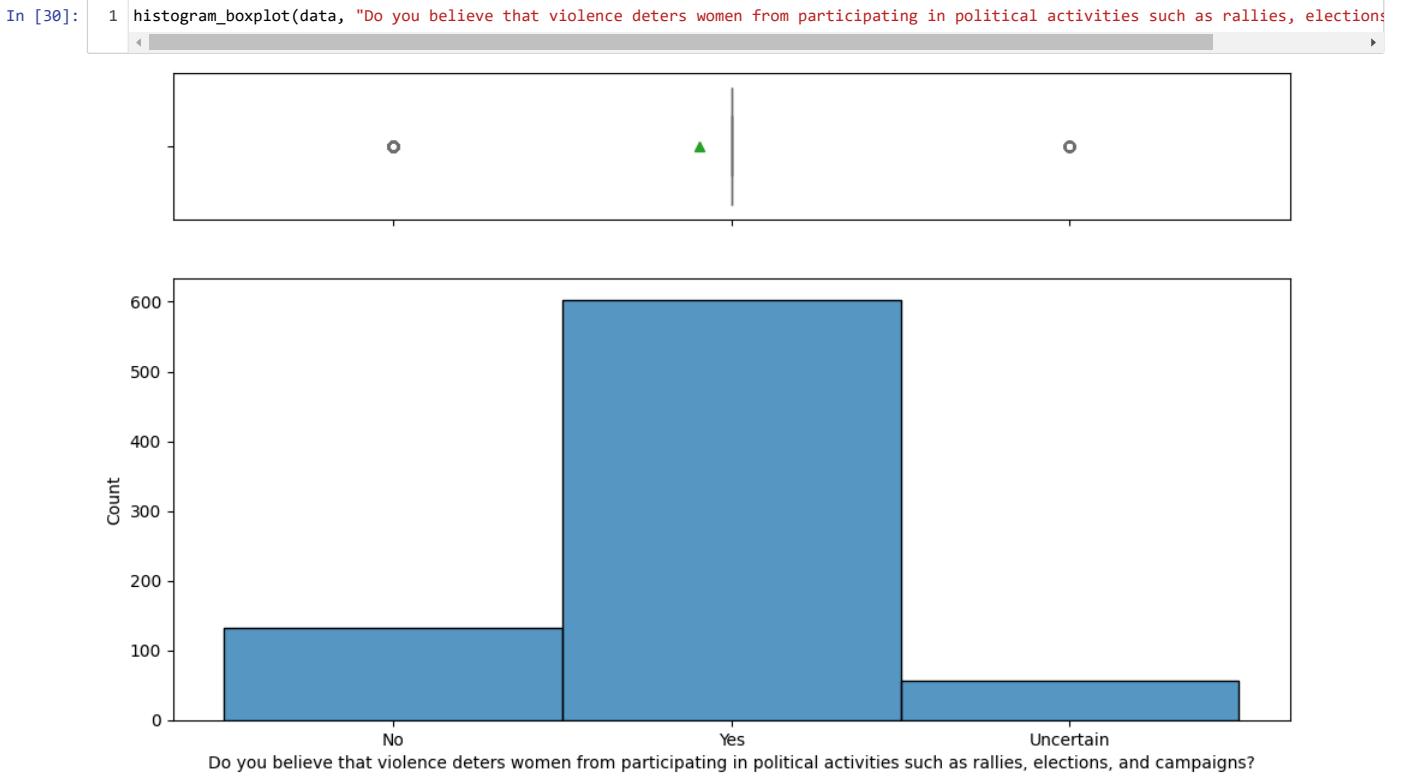


y to prevent a "female" loved one from going to vote a!

Observation

- 90% of the respondent indicate that they will prevent female loved ones from going to vote after violence occurs.

Observations on Respondent that believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?

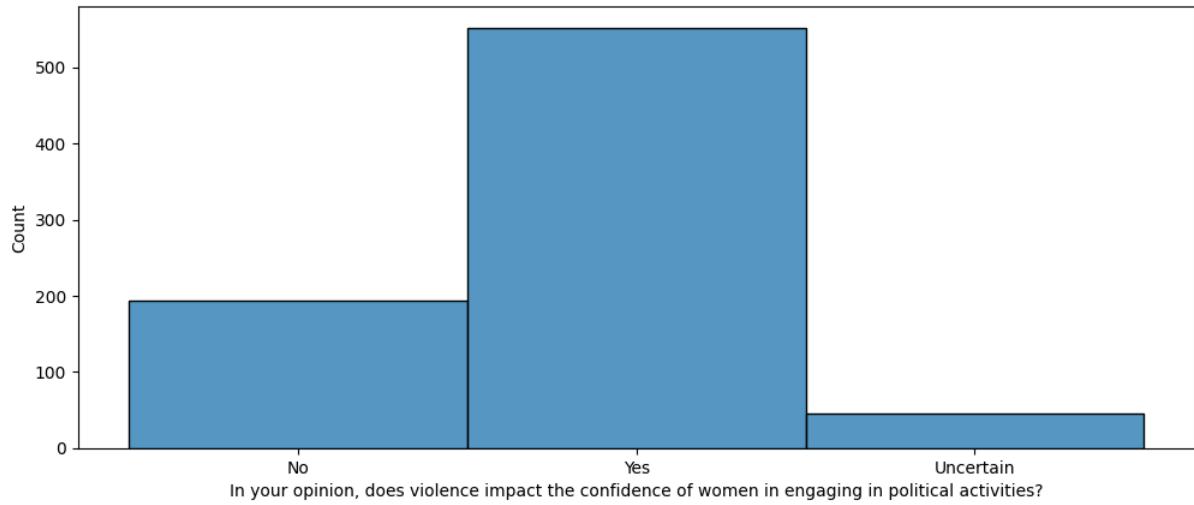
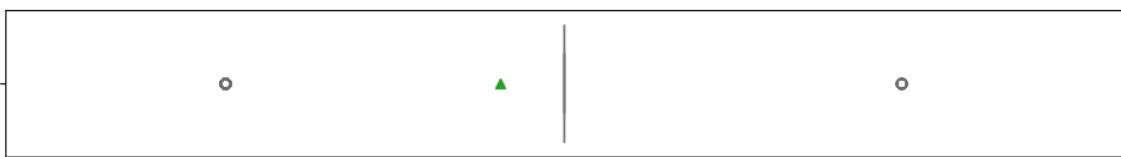


Observation

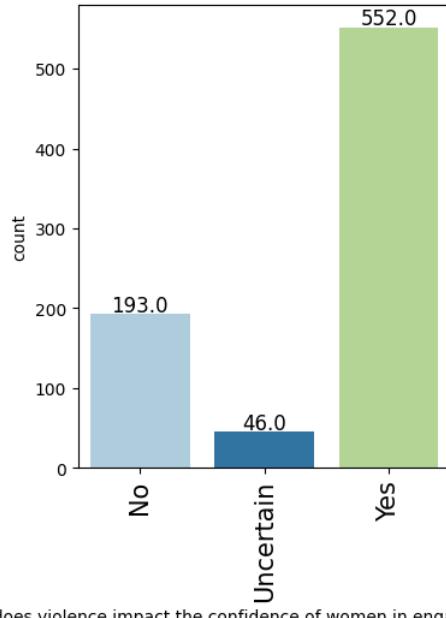
- Majority of the respondent believe that violence deters women/female from participating in political activities such as rallies, elections and campaigns

Observations on Respondent who have the opinoin that violence impact the confidence of women in engaging in political activities?

```
In [32]: 1 histogram_boxplot(data, "In your opinion, does violence impact the confidence of women in engaging in political activitites?")
```



```
In [33]: 1 labeled_barplot(data, "In your opinion, does violence impact the confidence of women in engaging in political activitites?")
```



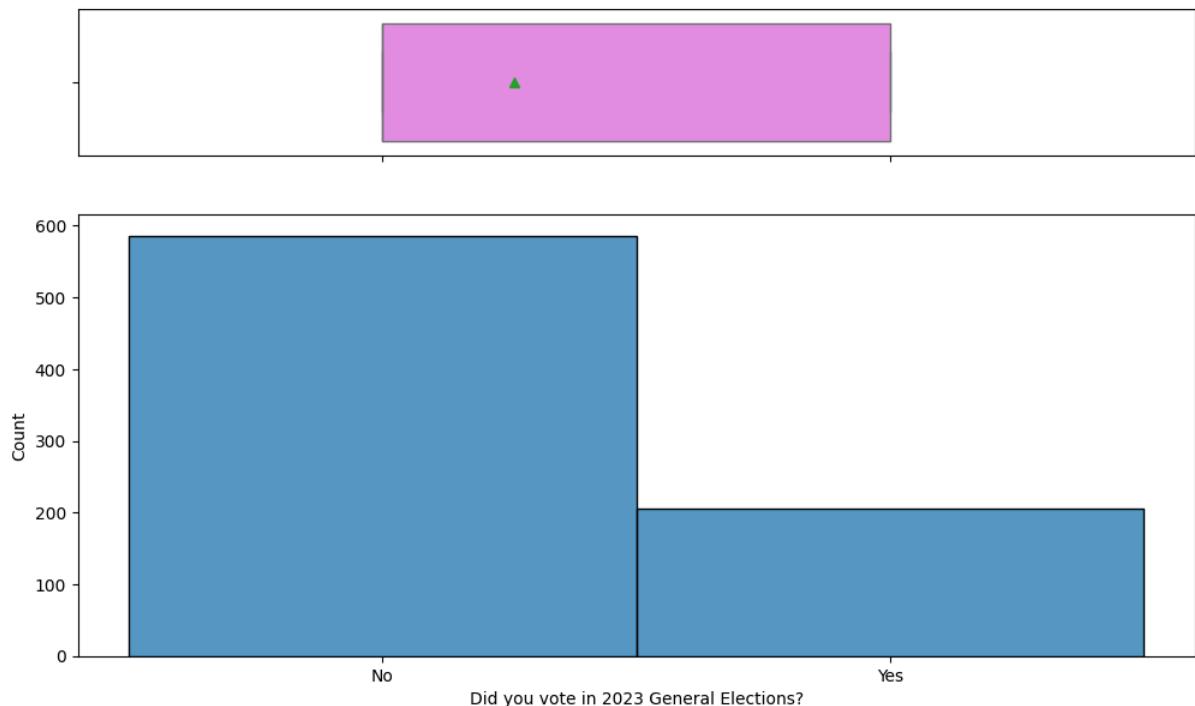
Does violence impact the confidence of women in engag

Observation

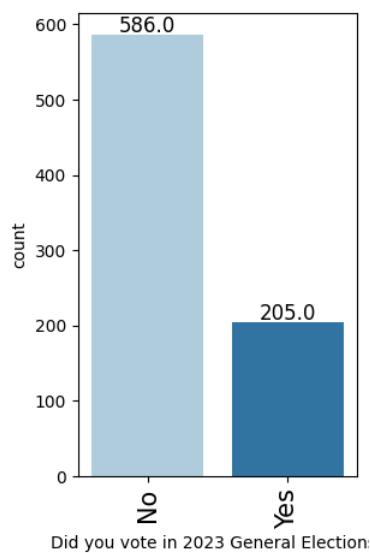
- Majority of the respondent beleive violence impact the confidence of women in engaging in political activities

Observations on Respondent that vote in 2023 General Elections?

```
In [34]: 1 histogram_boxplot(data, "Did you vote in 2023 General Elections?")
```



```
In [35]: 1 labeled_barplot(data, "Did you vote in 2023 General Elections?")
```

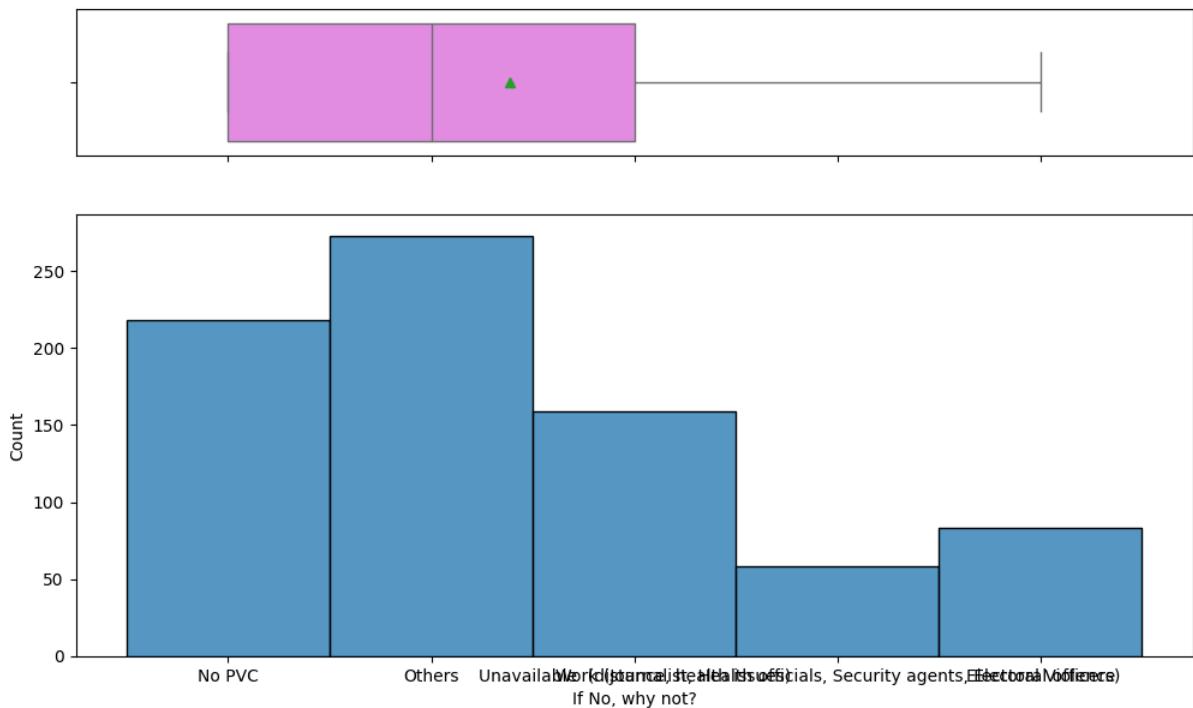


Observation

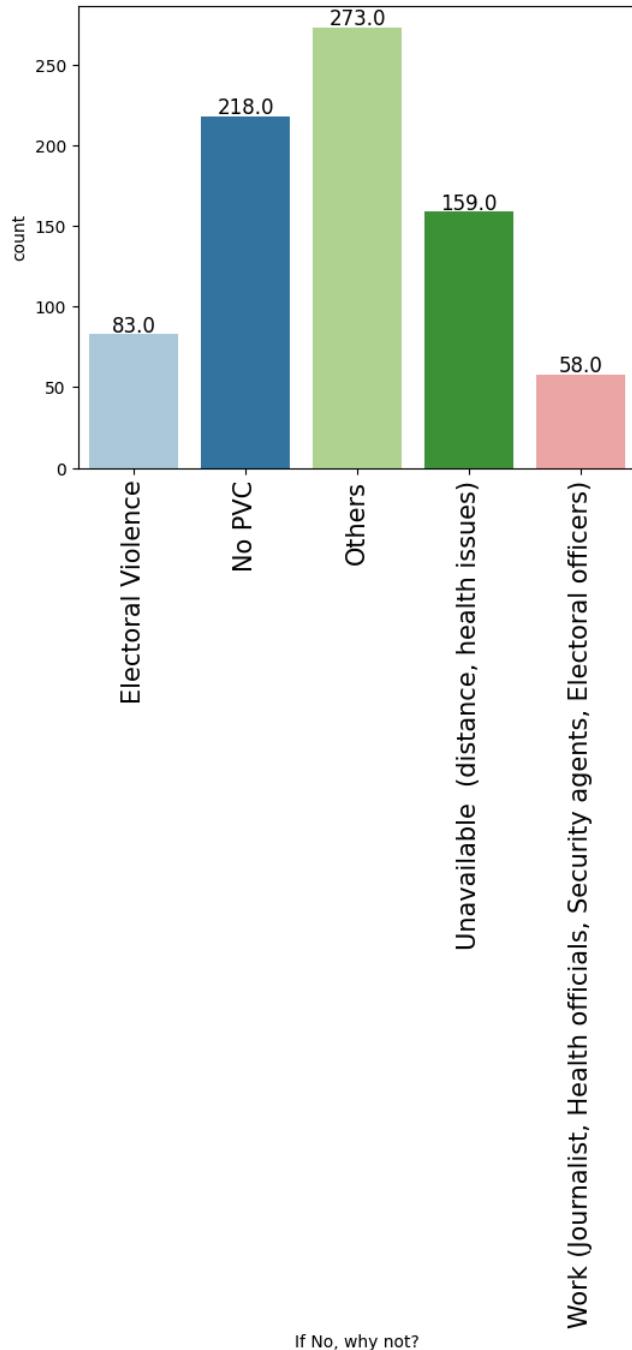
- Less than 300 respondent did not vote in 2023 General Election

Observation on Why Some Respondent did Not Vote in 2023 general election

In [36]: 1 histogram_boxplot(data, "If No, why not?")



```
In [37]: 1 labeled_barplot(data, "If No, why not?")
```

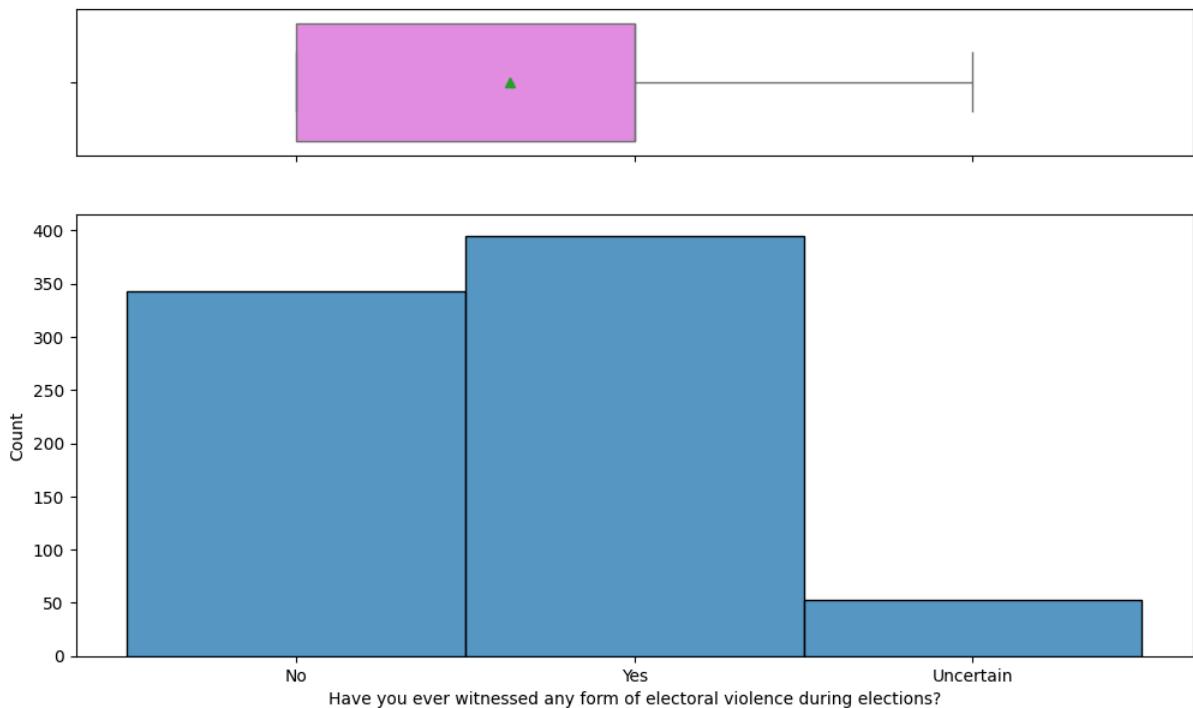


Observation

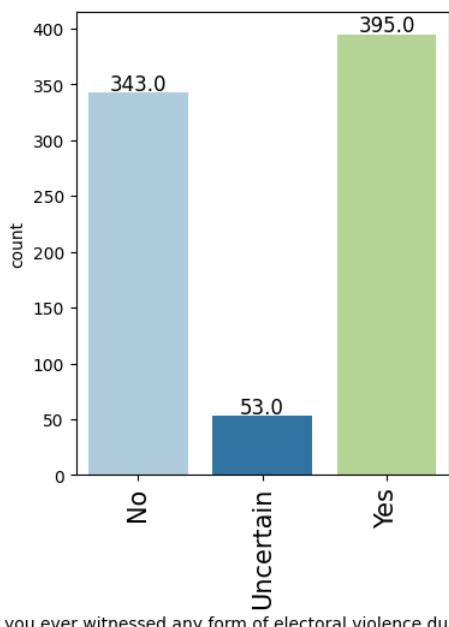
- Majority of the respondent that did not vote are influenced by
 1. Other reason best known to them
 2. Lack of Permanent Voters Card (PVC) and
 3. Unavailability (such as distance, health issues)
- Less than 100 respondent did not vote due to Electoral Violence
- Work (Journalist, Health officials, Security agents, Electoral officers) is reason for the least number of respondent that did not vote

Observation on Respondent who has ever witnessed any form of electoral violence during elections?

```
In [38]: 1 histogram_boxplot(data, "Have you ever witnessed any form of electoral violence during elections?")
```



```
In [39]: 1 labeled_barplot(data, "Have you ever witnessed any form of electoral violence during elections?")
```



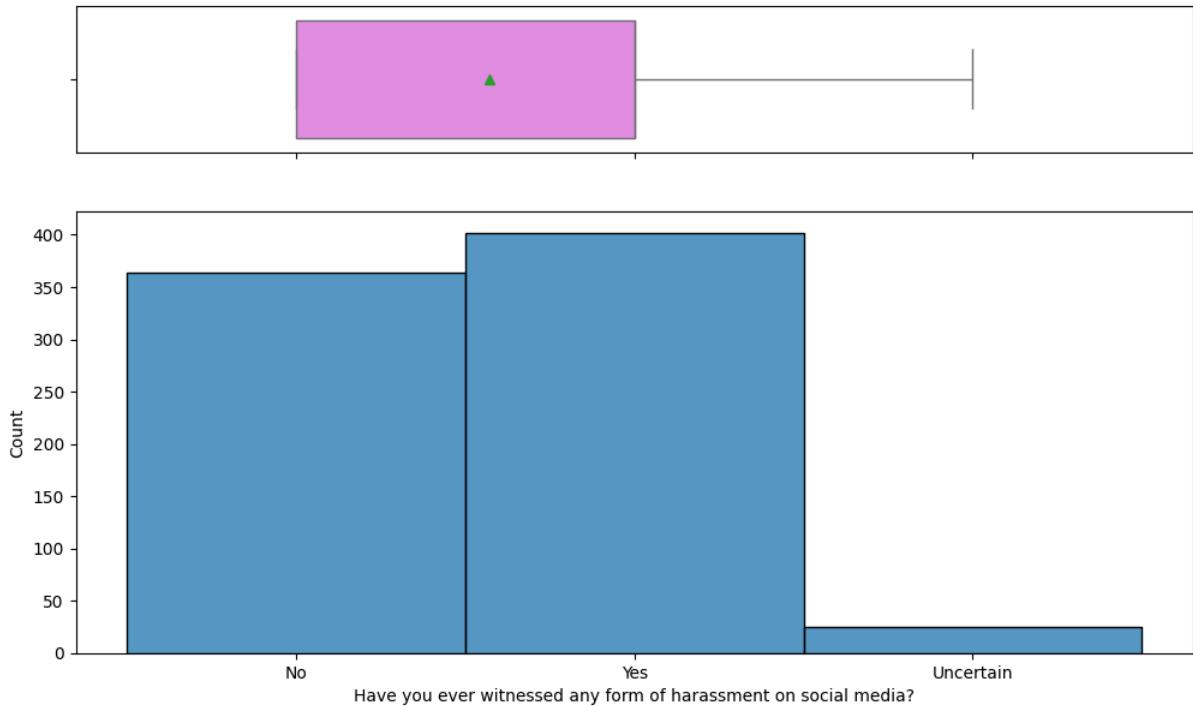
: you ever witnessed any form of electoral violence duri

Observation

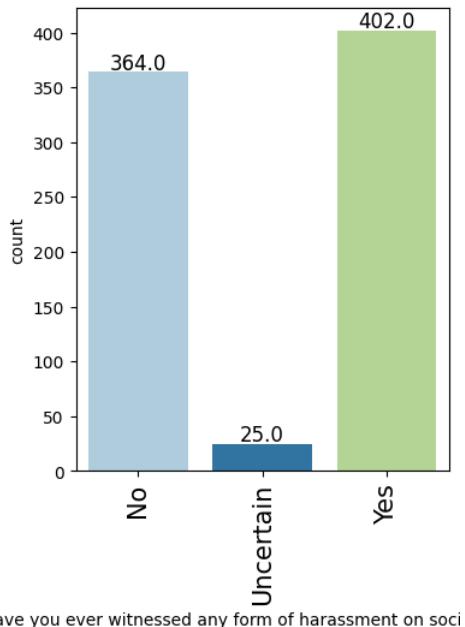
- Majority od the respondent has witness at least a particular rform of electoral violence

Observation on Respondent who has ever witnessed any form of harassment on social media?

```
In [40]: 1 histogram_boxplot(data, "Have you ever witnessed any form of harassment on social media?")
```



```
In [41]: 1 labeled_barplot(data, "Have you ever witnessed any form of harassment on social media?")
```



Have you ever witnessed any form of harassment on social media?

Observation

- Majority of the respondent has witnessed at least a particular form of harrasment on socia media

Bivariate Analysis

```
In [42]: 1 for i in data.describe(include=["object"]).columns:  
2     print("Unique values in", i, "are :")  
3     print(data[i].value_counts())  
4     print("*" * 50)  
5     print("*" * 50)
```

Unique values in Gender are :

Female	533
Male	258

Name: Gender, dtype: int64

Unique values in Work Sector are :

Formal Sector (9-5 jobs, Professionals, Hybrid jobs)	492
Informal Sector (Artisans, Traders)	299

Name: Work Sector, dtype: int64

Unique values in Educational Qualification are :

Tertiary (college or university)	467
Postgraduate	145
SSCE and below	140
Post-Secondary/Vocational	39

Name: Educational Qualification, dtype: int64

Unique values in Age range are :

18-30	380
31-40	173
51-60	116
60 and above	71
41-50	50
51-61	1

Name: Age range, dtype: int64

Unique values in Do you have a permanent voters card? are :

Yes	524
No	267

Name: Do you have a permanent voters card?, dtype: int64

Unique values in Are you likely to vote when there is electoral violence around you? are :

No	604
Yes	140
Uncertain	47

Name: Are you likely to vote when there is electoral violence around you?, dtype: int64

Unique values in Are you likely to prevent a "female" loved one from going to vote after violence occurs? are :

Yes	691
No	91
Uncertain	9

Name: Are you likely to prevent a "female" loved one from going to vote after violence occurs?, dtype: int64

Unique values in Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns? are :

Yes	603
No	132
Uncertain	56

Name: Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?, dtype: int64

Unique values in In your opinion, does violence impact the confidence of women in engaging in political activities? are :

Yes	552
No	193
Uncertain	46

Name: In your opinion, does violence impact the confidence of women in engaging in political activities?, dtype: int64

Unique values in Did you vote in 2023 General Elections? are :

No	586
Yes	205

Name: Did you vote in 2023 General Elections?, dtype: int64

Unique values in If No, why not? are :

Others	273
No PVC	218
Unavailable (distance, health issues)	159
Electoral Violence	83
Work (Journalist, Health officials, Security agents, Electoral officers)	58

Name: If No, why not?, dtype: int64

Unique values in Have you ever witnessed any form of electoral violence during elections? are :

Yes	395
No	343
Uncertain	53

Name: Have you ever witnessed any form of electoral violence during elections?, dtype: int64

Unique values in Have you ever witnessed any form of harassment on social media? are :

Yes	402
No	364
Uncertain	25

Name: Have you ever witnessed any form of harassment on social media?, dtype: int64

```
In [43]: 1 data.columns
```

```
Out[43]: Index(['Gender', 'Work Sector', 'Educational Qualification', 'Age range',
   'Do you have a permanent voters card?',
   'Are you likely to vote when there is electoral violence around you?',
   'Are you likely to prevent a "female" loved one from going to vote after violence occurs?',
   'Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?',
   'In your opinion, does violence impact the confidence of women in engaging in political activities?',
   'Did you vote in 2023 General Elections?', 'If No, why not?',
   'Have you ever witnessed any form of electoral violence during elections?',
   'Have you ever witnessed any form of harassment on social media?'],
  dtype='object')
```

```
In [44]: 1 ## Encoding values of attributes into numerical value respectively, for analysis.
2
3 dataCorr = data.copy()
4
5 dataCorr["Gender"].replace("Male", 1, inplace=True)
6 dataCorr["Gender"].replace("Female", 0, inplace=True)
7
8 dataCorr["Work Sector"].replace("Formal Sector (9-5 jobs, Professionals, Hybrid jobs)", 1, inplace=True)
9 dataCorr["Work Sector"].replace("Informal Sector (Artisans, Traders)", 0, inplace=True)
10
11 dataCorr["Do you have a permanent voters card?"].replace("Yes", 1, inplace=True)
12 dataCorr["Do you have a permanent voters card?"].replace("No", 0, inplace=True)
13
14 dataCorr["Did you vote in 2023 General Elections?"].replace("Yes", 1, inplace=True)
15 dataCorr["Did you vote in 2023 General Elections?"].replace("No", 0, inplace=True)
16
17 dataCorr["Educational Qualification"].replace("Postgraduate", 20, inplace=True)
18 dataCorr["Educational Qualification"].replace("Tertiary (college or university)", 15, inplace=True)
19 dataCorr["Educational Qualification"].replace("Post-Secondary/Vocational", 10, inplace=True)
20 dataCorr["Educational Qualification"].replace("SSCE and below", 5, inplace=True)
21 dataCorr["Educational Qualification"].astype('float64')
22
23 dataCorr["Age range"].replace("18-30", 24, inplace=True)
24 dataCorr["Age range"].replace("31-40", 35.5, inplace=True)
25 dataCorr["Age range"].replace("51-60", 55.5, inplace=True)
26 dataCorr["Age range"].replace("41-50", 45.5, inplace=True)
27 dataCorr["Age range"].replace("51-61", 56, inplace=True)
28 dataCorr["Age range"].replace("60 and above", 80, inplace=True)
29 dataCorr["Age range"].astype('float64')
30
31 dataCorr["Are you likely to vote when there is electoral violence around you?"].replace("No", 0, inplace=True)
32 dataCorr["Are you likely to vote when there is electoral violence around you?"].replace("Yes", 1, inplace=True)
33 dataCorr["Are you likely to vote when there is electoral violence around you?"].replace("Uncertain", 0.5, inplace=True)
34 dataCorr["Are you likely to vote when there is electoral violence around you?"].astype('float64')
35
36 dataCorr["Are you likely to prevent a "female" loved one from going to vote after violence occurs?"].replace("No", 0, inplace=True)
37 dataCorr["Are you likely to prevent a "female" loved one from going to vote after violence occurs?"].replace("Yes", 1, inplace=True)
38 dataCorr["Are you likely to prevent a "female" loved one from going to vote after violence occurs?"].replace("Uncertain", 0.5, inplace=True)
39 dataCorr["Are you likely to prevent a "female" loved one from going to vote after violence occurs?"].astype('float64')
40
41 dataCorr["Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns"]
42 dataCorr["Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns"]
43 dataCorr["Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns"]
44 dataCorr["Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns"]
45
46 dataCorr["If No, why not?"].replace("Others", 0.3, inplace=True)
47 dataCorr["If No, why not?"].replace("No PVC", 0, inplace=True)
48 dataCorr["If No, why not?"].replace("Unavailable (distance, health issues)", 0.5, inplace=True)
49 dataCorr["If No, why not?"].replace("Electoral Violence", 1, inplace=True)
50 dataCorr["If No, why not?"].replace("Work (Journalist, Health officials, Security agents, Electoral officers)", 0.5, inplace=True)
51 dataCorr["If No, why not?"].astype('float64')
52
53 dataCorr["Have you ever witnessed any form of electoral violence during elections?"].replace("Yes", 1, inplace=True)
54 dataCorr["Have you ever witnessed any form of electoral violence during elections?"].replace("No", 0, inplace=True)
55 dataCorr["Have you ever witnessed any form of electoral violence during elections?"].replace("Uncertain", 0.5, inplace=True)
56 dataCorr["Have you ever witnessed any form of electoral violence during elections?"].astype('float64')
57
58 dataCorr["Have you ever witnessed any form of harassment on social media?"].replace("Yes", 1, inplace=True)
59 dataCorr["Have you ever witnessed any form of harassment on social media?"].replace("No", 0, inplace=True)
60 dataCorr["Have you ever witnessed any form of harassment on social media?"].replace("Uncertain", 0.5, inplace=True)
61 dataCorr["Have you ever witnessed any form of harassment on social media?"].astype('float64')
62
63 dataCorr["In your opinion, does violence impact the confidence of women in engaging in political activities?"].replace("Yes", 1, inplace=True)
64 dataCorr["In your opinion, does violence impact the confidence of women in engaging in political activities?"].replace("No", 0, inplace=True)
65 dataCorr["In your opinion, does violence impact the confidence of women in engaging in political activities?"].replace("Uncertain", 0.5, inplace=True)
66 dataCorr["In your opinion, does violence impact the confidence of women in engaging in political activities?"].astype('float64')
```

```
Out[44]: 0    0.0
1    1.0
2    1.0
3    1.0
4    1.0
...
786   1.0
787   0.0
788   1.0
789   0.5
790   1.0
Name: In your opinion, does violence impact the confidence of women in engaging in political activities?, Length: 791, dtype: float64
```

```
In [45]: caCbrr.rename(columns = {'Educational Qualification':'Edu. Qlf.'}, inplace = True)
caCbrr.rename(columns = {'Do you have a permanent voters card?':'PVC'}, inplace = True)
caCbrr.rename(columns = {'Did you vote in 2023 General Elections?':'Vote in 2023 Gen. Elec.'}, inplace = True)
caCbrr.rename(columns = {'Are you likely to vote when there is electoral violence around you?':'Vote During Elec. Vio.'}, inplace = True)
caCbrr.rename(columns = {'Are you likely to prevent a "female" loved one from going to vote after violence occurs?':'Allow Female Vote During Elec. Vio.'}, inplace = True)
caCbrr.rename(columns = {'Do you believe that violence deters women from participating in political activities such as rallies, elections, and other forms of political engagement?':'Violence Deters Women From Parti.'}, inplace = True)
caCbrr.rename(columns = {'In your opinion, does violence impact the confidence of women in engaging in political activities?':'Violence Impact Women Confi. In Parti.'}, inplace = True)
caCbrr.rename(columns = {'Have you ever witnessed any form of electoral violence during elections?':'Witnessed any Elec. Vio.'}, inplace = True)
caCbrr.rename(columns = {'Have you ever witnessed any form of harassment on social media?':'Witnessed Haras. Social Media'}, inplace = True)
10
caCbrr.columns
```

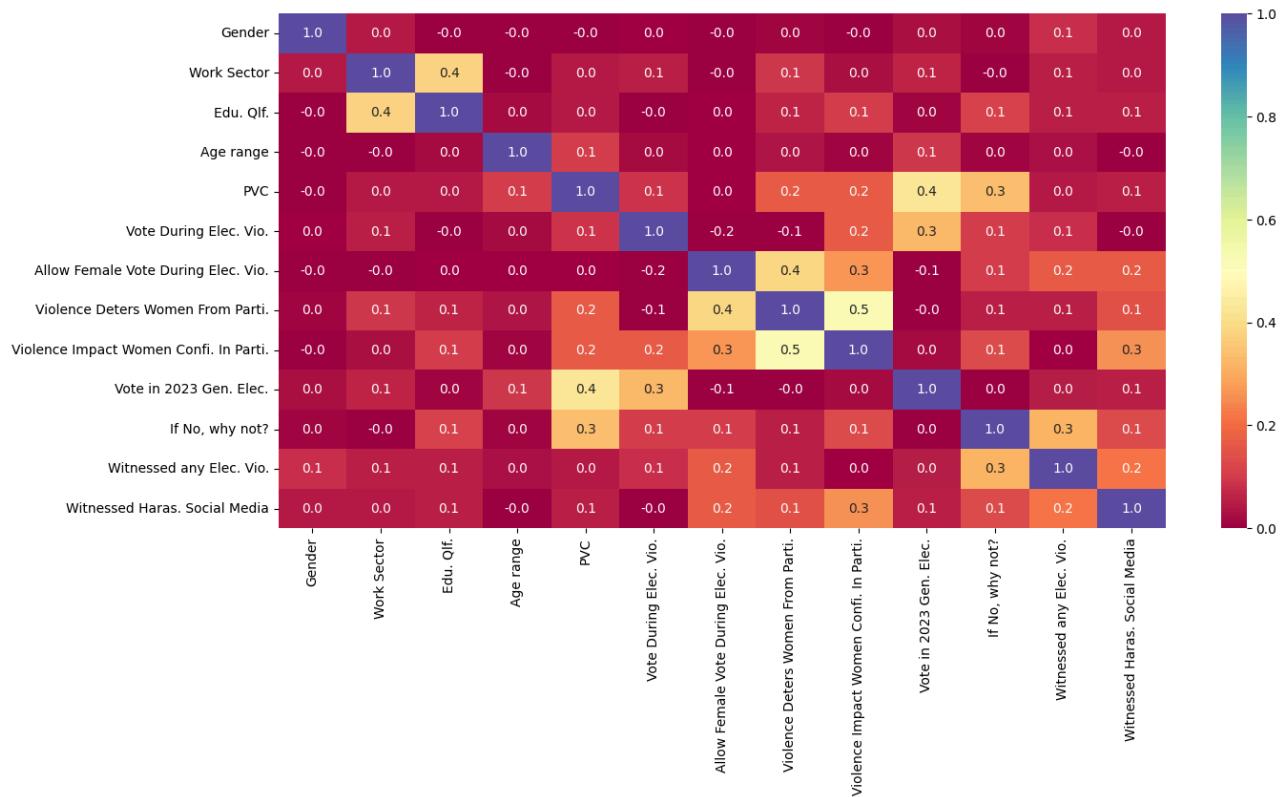
```
Out[45]: Index(['Gender', 'Work Sector', 'Edu. Qlf.', 'Age range', 'PVC',
       'Vote During Elec. Vio.', 'Allow Female Vote During Elec. Vio.',
       'Violence Deters Women From Parti.',
       'Violence Impact Women Confi. In Parti.', 'Vote in 2023 Gen. Elec.',
       'If No, why not?', 'Witnessed any Elec. Vio.',
       'Witnessed Haras. Social Media'],
      dtype='object')
```

```
In [46]: 1 dataCorr.corr(method='pearson', min_periods=0, numeric_only = True)
```

Out[46]:

	Gender	Work Sector	Edu. Qlf.	Age range	PVC	Vote During Elec. Vio.	Allow Female Vote During Elec. Vio.	Violence Deters Women From Parti.	Violence Impact Women Confi. In Parti.	Vote in 2023 Gen. Elec.	If No, why not?	Witnessed any Elec. Vio.	Witnessed Haras. Social Media
Gender	1.000000	0.041848	-0.006414	-0.011497	-0.016611	0.004688	-0.036502	0.015561	-0.022303	0.025450	0.008553	0.078462	0.040072
Work Sector	0.041848	1.000000	0.381546	-0.009647	0.044513	0.056062	-0.021084	0.092999	0.026451	0.062421	-0.040700	0.053180	0.040755
Edu. Qlf.	-0.006414	0.381546	1.000000	0.022666	0.046042	-0.005317	0.004132	0.062870	0.104492	0.009453	0.115801	0.057553	0.052234
Age range	-0.011497	-0.009647	0.022666	1.000000	0.103473	0.017472	0.007416	0.033838	0.010143	0.090162	0.009245	0.029101	-0.034038
PVC	-0.016611	0.044513	0.046042	0.103473	1.000000	0.087849	0.006357	0.162179	0.162623	0.422200	0.338936	0.043147	0.053929
Vote During Elec. Vio.	0.004688	0.056062	-0.005317	0.017472	0.087849	1.000000	-0.150338	-0.085752	0.163322	0.322217	0.103107	0.080689	-0.014504
Allow Female Vote During Elec. Vio.	-0.036502	-0.021084	0.004132	0.007416	0.006357	-0.150338	1.000000	0.385850	0.265218	-0.078540	0.103252	0.170516	0.168421
Violence Deters Women From Parti.	0.015561	0.092999	0.062870	0.033838	0.162179	-0.085752	0.385850	1.000000	0.517594	-0.007867	0.057761	0.057173	0.144856
Violence Impact Women Confi. In Parti.	-0.022303	0.026451	0.104492	0.010143	0.162623	0.163322	0.265218	0.517594	1.000000	0.020044	0.135969	0.006728	0.259027
Vote in 2023 Gen. Elec.	0.025450	0.062421	0.009453	0.090162	0.422200	0.322217	-0.078540	-0.007867	0.020044	1.000000	0.001425	0.049470	0.056218
If No, why not?	0.008553	-0.040700	0.115801	0.009245	0.338936	0.103107	0.103252	0.057761	0.135969	0.001425	1.000000	0.319391	0.132555
Witnessed any Elec. Vio.	0.078462	0.053180	0.057553	0.029101	0.043147	0.080689	0.170516	0.057173	0.006728	0.049470	0.319391	1.000000	0.218226
Witnessed Haras. Social Media	0.040072	0.040755	0.052234	-0.034038	0.053929	-0.014504	0.168421	0.144856	0.259027	0.056218	0.132555	0.218226	1.000000

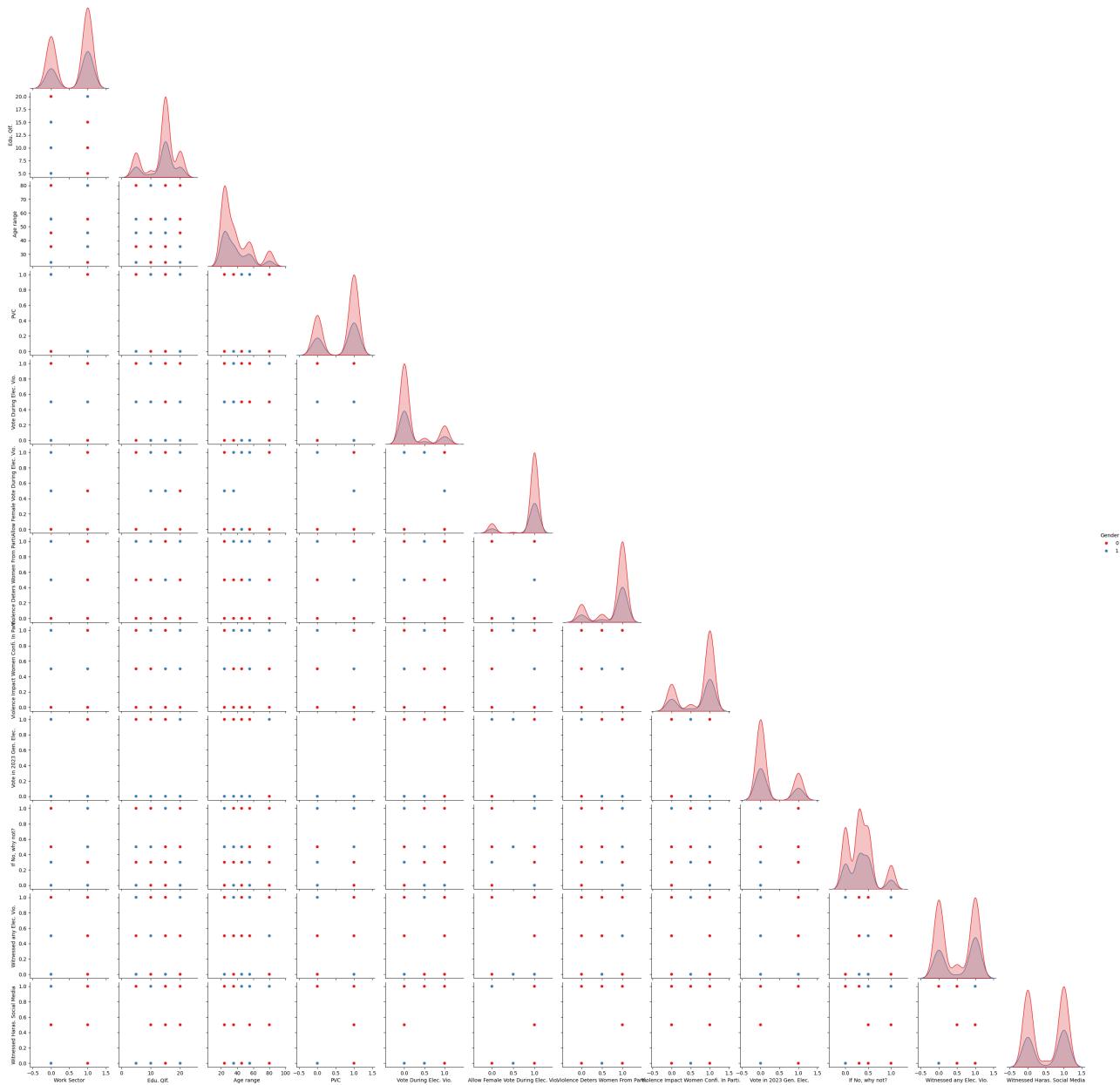
```
In [47]: 1 plt.figure(figsize=(15, 7))
2 sns.heatmap(dataCorr.corr(method='pearson'), annot=True, vmin=-0, vmax=1, fmt=".1f", cmap="Spectral")
3 plt.show()
```



Observation

- The education qualification is well correlated to the work sector of the respondent
- vote is 2023 election is well correlated with have a PVC and vote during election violence but slightly correlated to having witnessed harrassment on social media
- violence deter women from participating in electoral activities is highly correlated to allowing female to vote during electoral violence
- witnessing any form of electoral violence is also correlated to Why respondent do not vote in the 2023 general election ("If No, why Not")
- Gender shows a slight correlation with 'witnessed any Electoral Violence'

```
In [48]: 1 sns.set_palette(sns.color_palette("Set1", 8))
2 sns.pairplot(dataCorr, hue="Gender", corner=True)
# plot_kws={'Line_kws':{'color':'red'}}
4
5 plt.show()
```

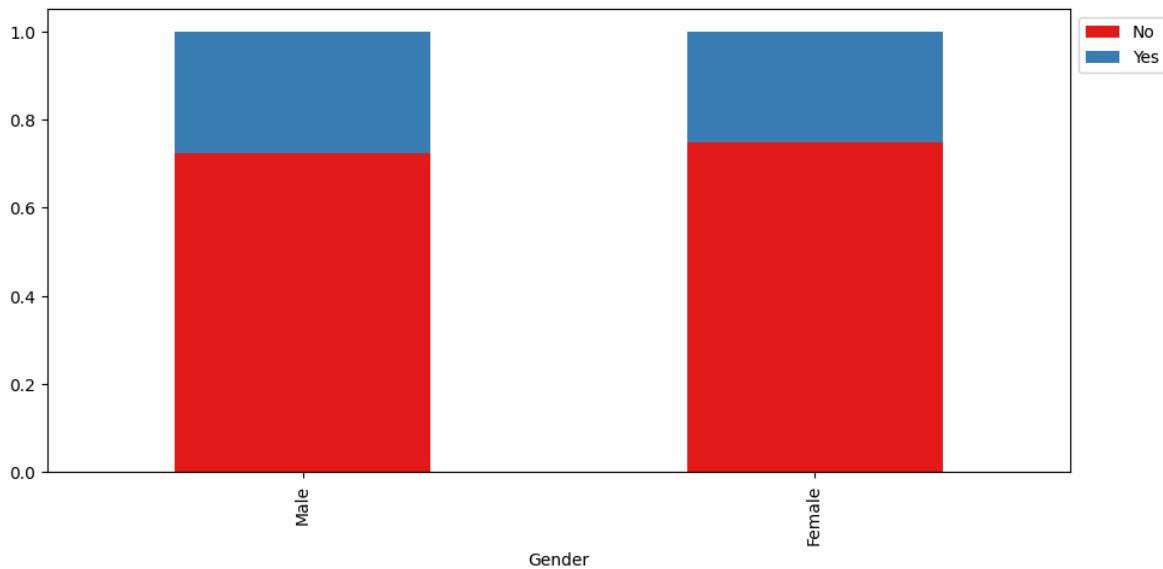


```
In [49]: 1 # function to plot stacked bar chart
2
3
4 def stacked_barplot(data, predictor, target):
5     """
6     Print the category counts and plot a stacked bar chart
7
8     data: dataframe
9     predictor: independent variable
10    target: target variable
11    """
12    count = data[predictor].nunique()
13    sorter = data[target].value_counts().index[-1]
14    tab1 = pd.crosstab(data[predictor], data[target], margins=True).sort_values(
15        by=sorter, ascending=False
16    )
17    print(tab1)
18    print("-" * 120)
19    tab = pd.crosstab(data[predictor], data[target], normalize="index").sort_values(
20        by=sorter, ascending=False
21    )
22    tab.plot(kind="bar", stacked=True, figsize=(count + 9, 5))
23    plt.legend(
24        loc="lower left", frameon=False,
25    )
26    plt.legend(loc="upper left", bbox_to_anchor=(1, 1))
27    plt.show()
```

Gender vs Did you vote in 2023 General Elections?

```
In [50]: 1 stacked_barplot(data, "Gender", "Did you vote in 2023 General Elections?")
```

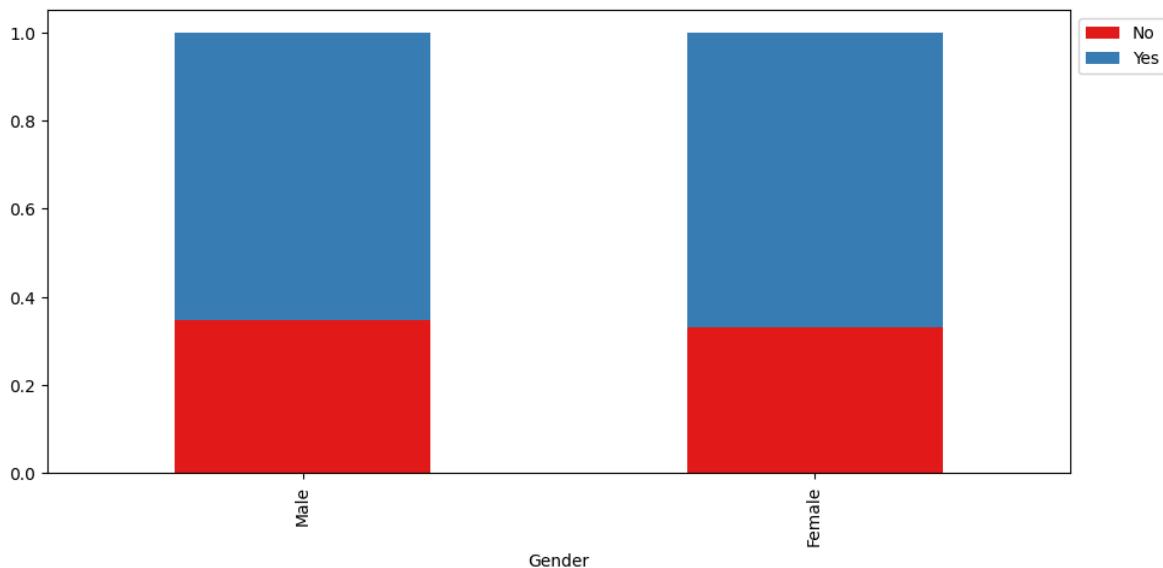
Did you vote in 2023 General Elections?	No	Yes	All
Gender			
All	586	205	791
Female	399	134	533
Male	187	71	258



Gender vs Do you have a permanent voters card?

```
In [51]: 1 stacked_barplot(data, "Gender", "Do you have a permanent voters card?")
```

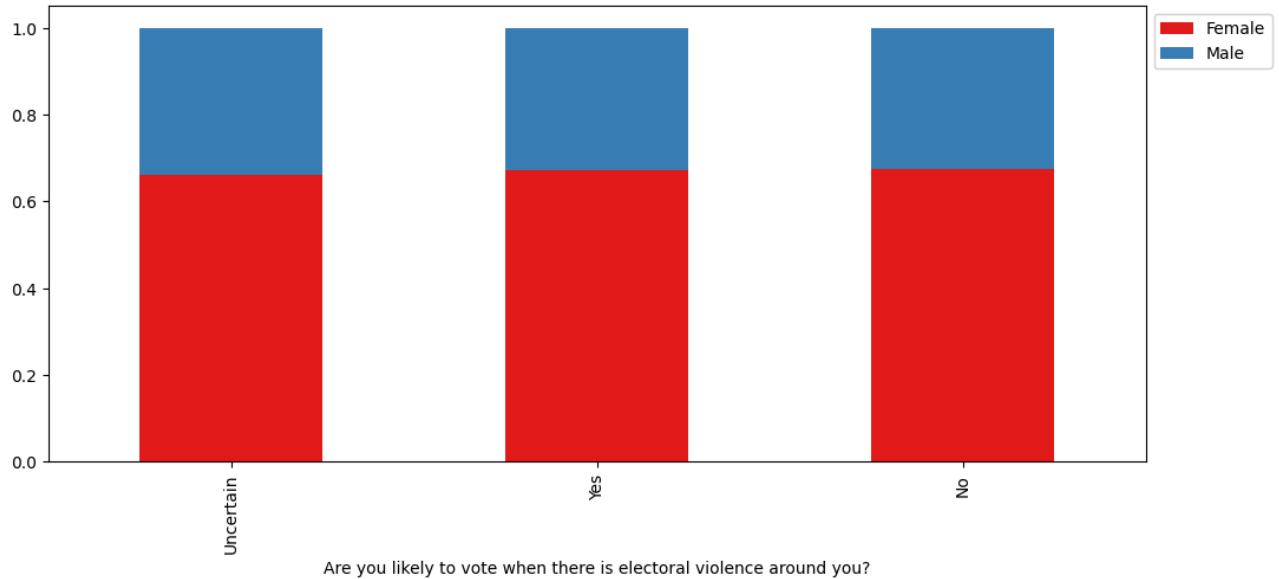
Do you have a permanent voters card?	No	Yes	All
Gender			
All	267	524	791
Female	177	356	533
Male	90	168	258



Gender vs Are you likely to vote when there is electoral violence around you?

```
In [52]: 1 stacked_barplot(data, "Are you likely to vote when there is electoral violence around you?", "Gender")
```

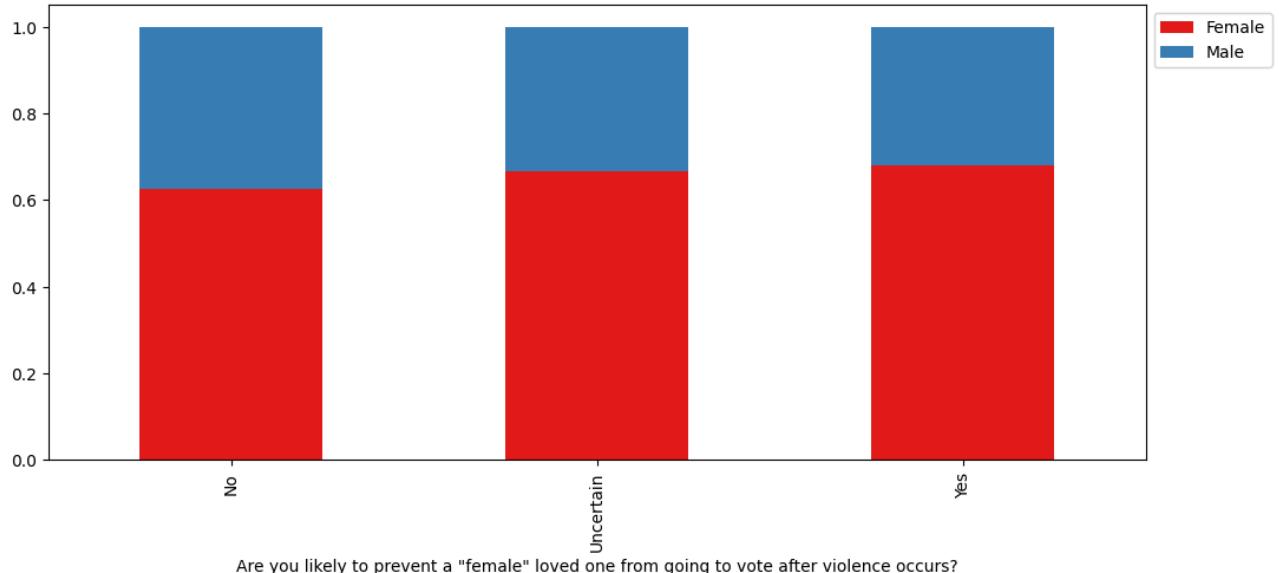
Gender	Female	Male	All
Are you likely to vote when there is electoral violence around you?			
All	533	258	791
No	408	196	604
Yes	94	46	140
Uncertain	31	16	47



Gender vs Are you likely to prevent a "female" loved one from going to vote after violence occurs?

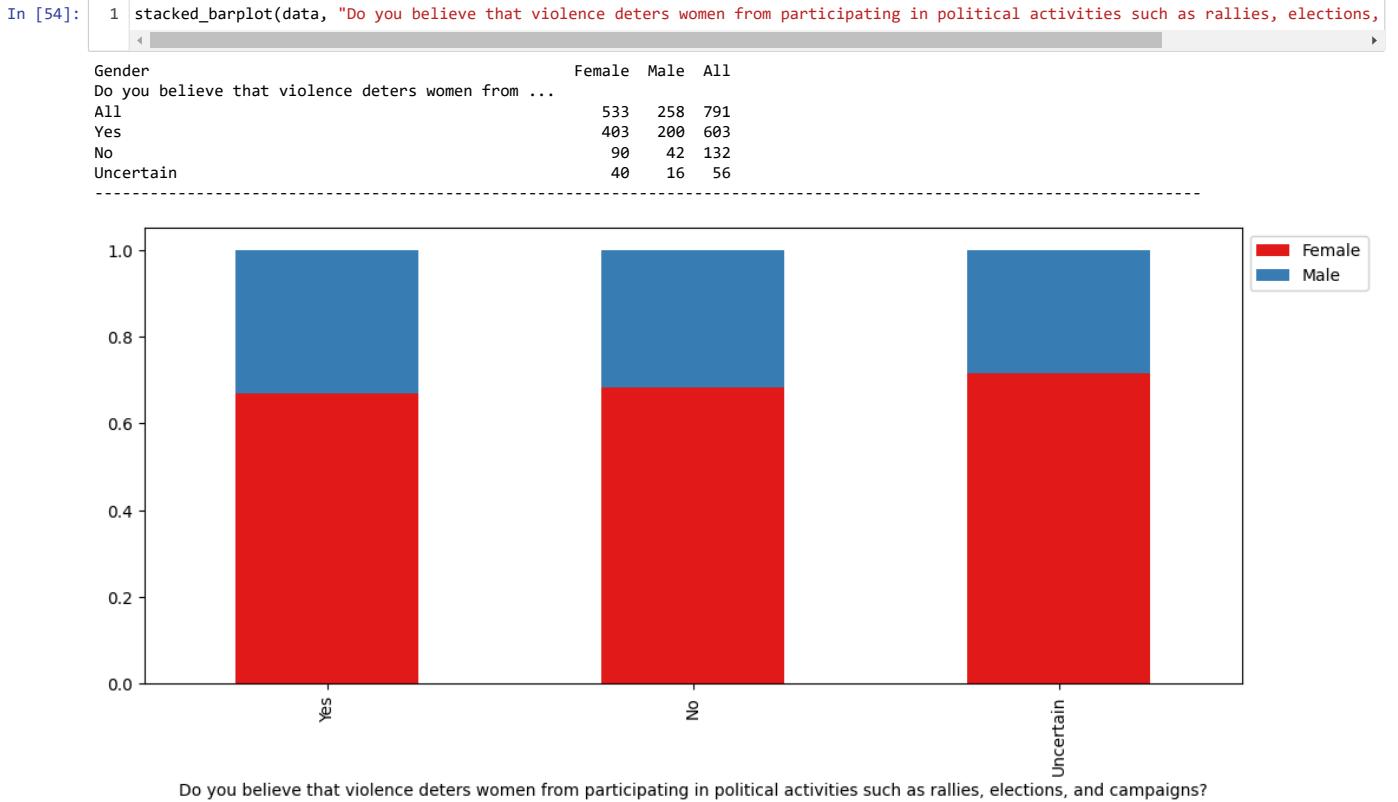
```
In [53]: 1 stacked_barplot(data, 'Are you likely to prevent a "female" loved one from going to vote after violence occurs?', "Gender")
```

Gender	Female	Male	All
Are you likely to prevent a "female" loved one from going to vote after violence occurs?			
All	533	258	791
Yes	470	221	691
No	57	34	91
Uncertain	6	3	9

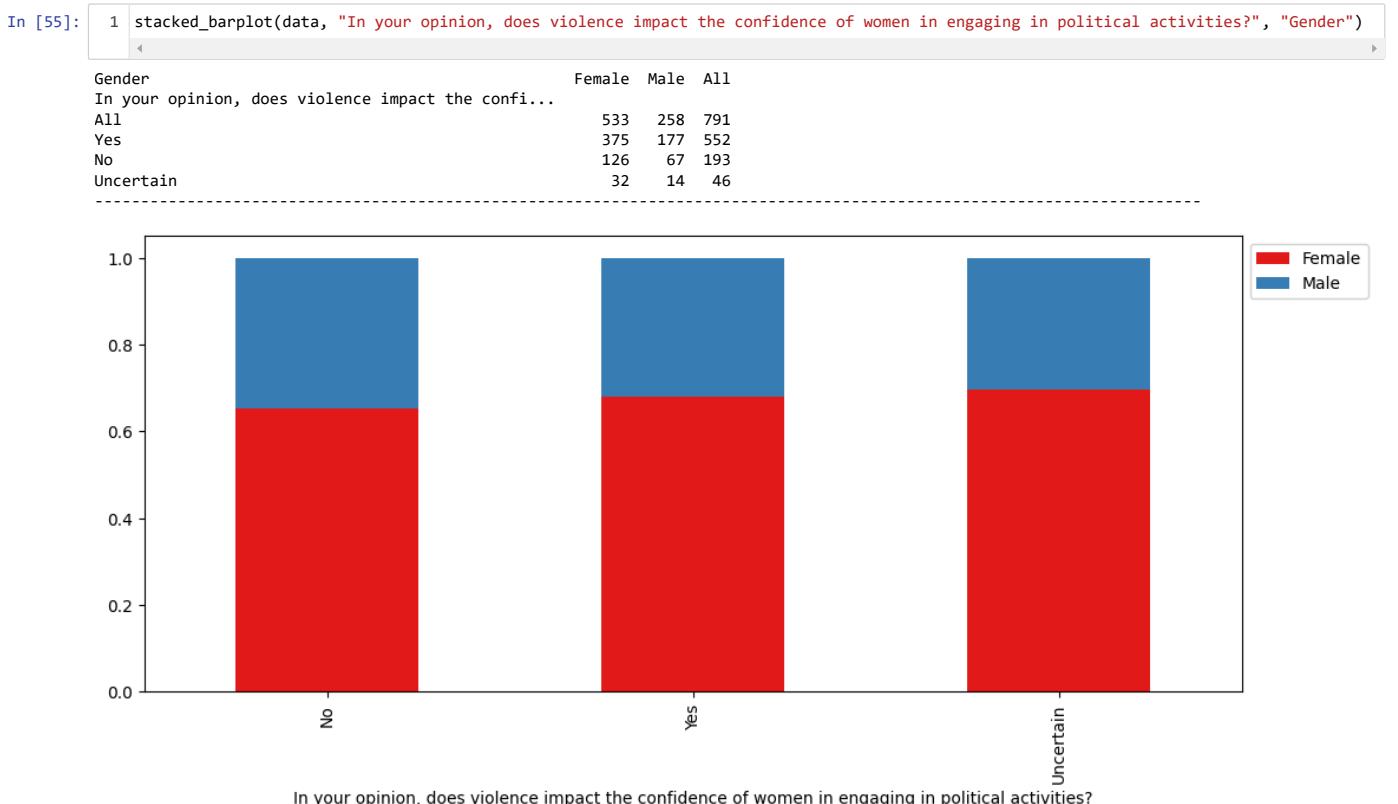


- The customers from two extreme income groups - Earning less than 40K and Earning more than 120k+ are the ones attriting the most.

Gender vs Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?



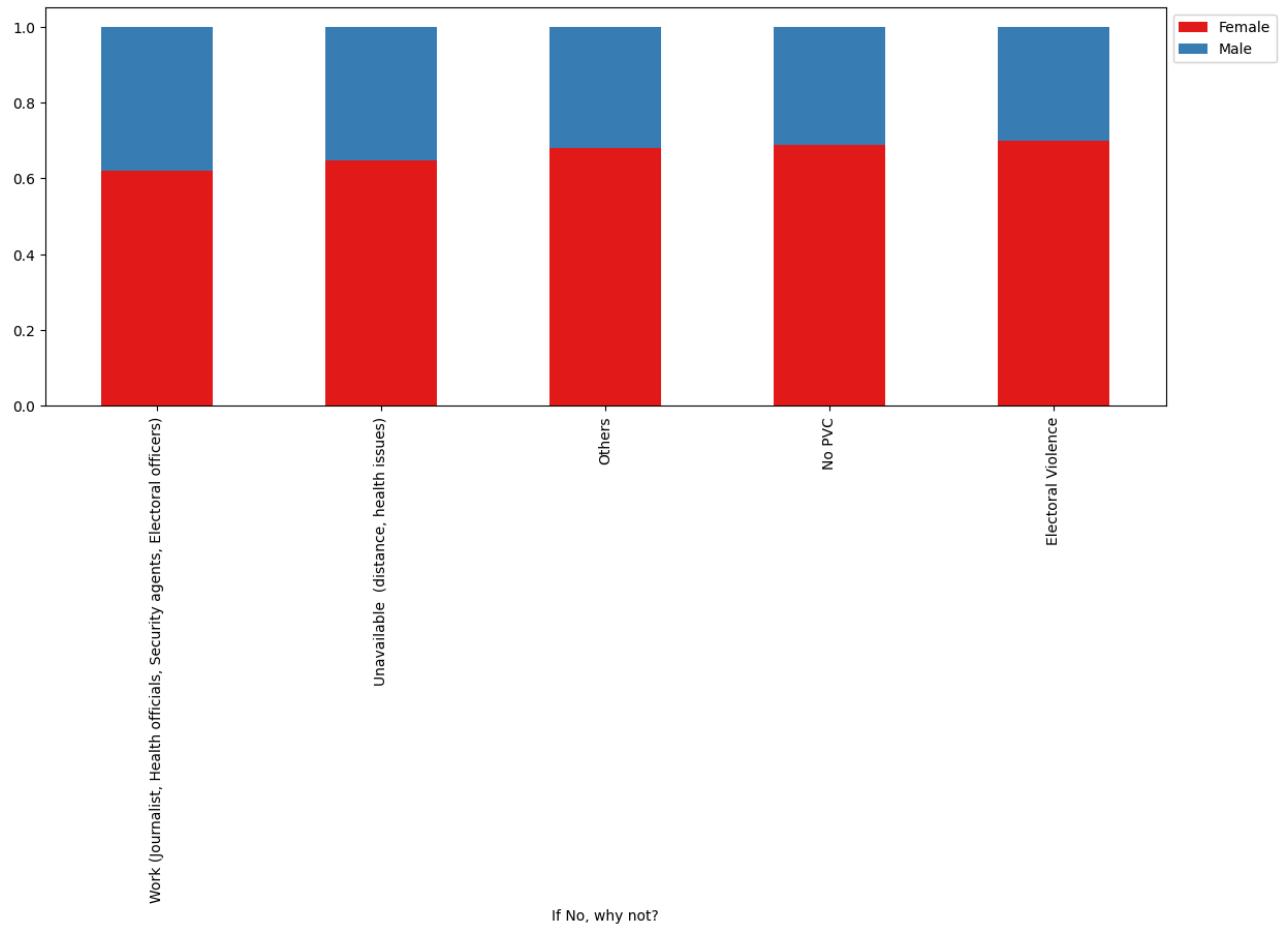
Gender vs In your opinion, does violence impact the confidence of women in engaging in political activities?



Gender vs If No, why not?

```
In [56]: 1 stacked_barplot(data, "If No, why not?", "Gender")
```

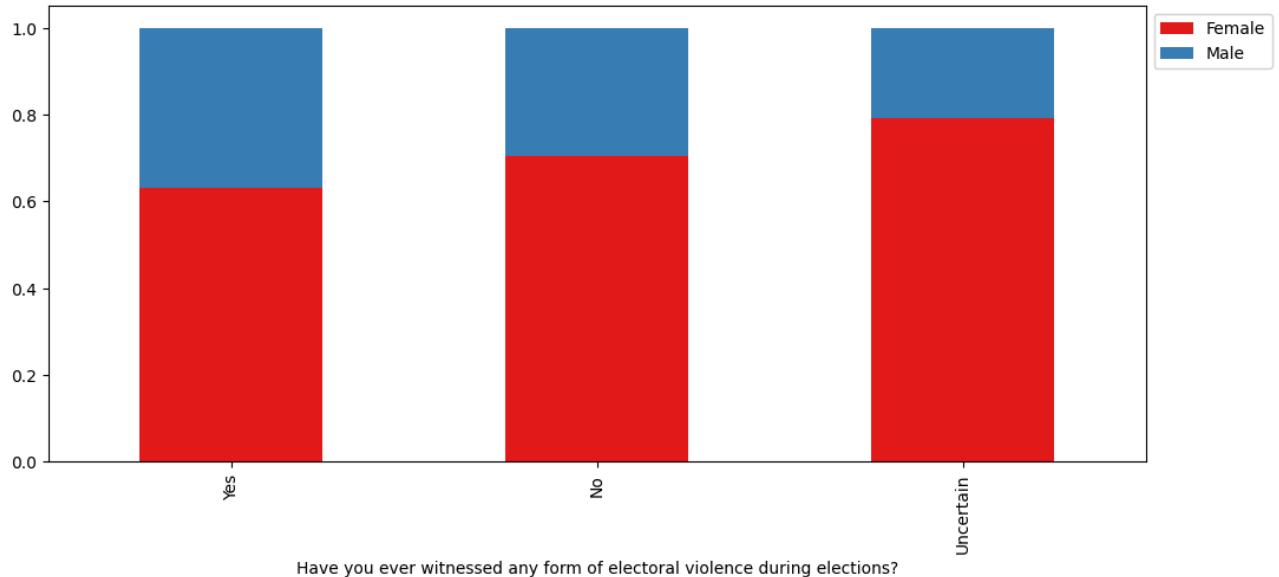
	Female	Male	All
Gender			
If No, why not?			
All	533	258	791
Others	186	87	273
No PVC	150	68	218
Unavailable (distance, health issues)	103	56	159
Electoral Violence	58	25	83
Work (Journalist, Health officials, Security ag...	36	22	58



Gender vs Have you ever witnessed any form of electoral violence during elections?

```
In [57]: 1 stacked_barplot(data, "Have you ever witnessed any form of electoral violence during elections?", "Gender")
```

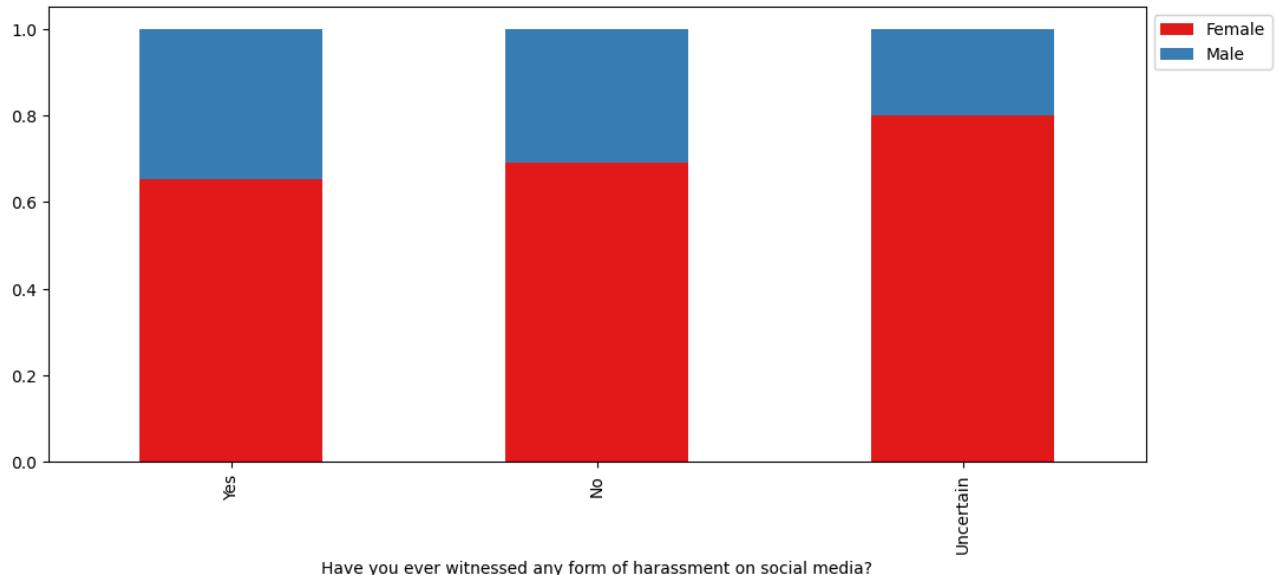
Gender	Female	Male	All
Have you ever witnessed any form of electoral violence during elections?			
All	533	258	791
Yes	249	146	395
No	242	101	343
Uncertain	42	11	53



Gender vs Have you ever witnessed any form of harassment on social media?

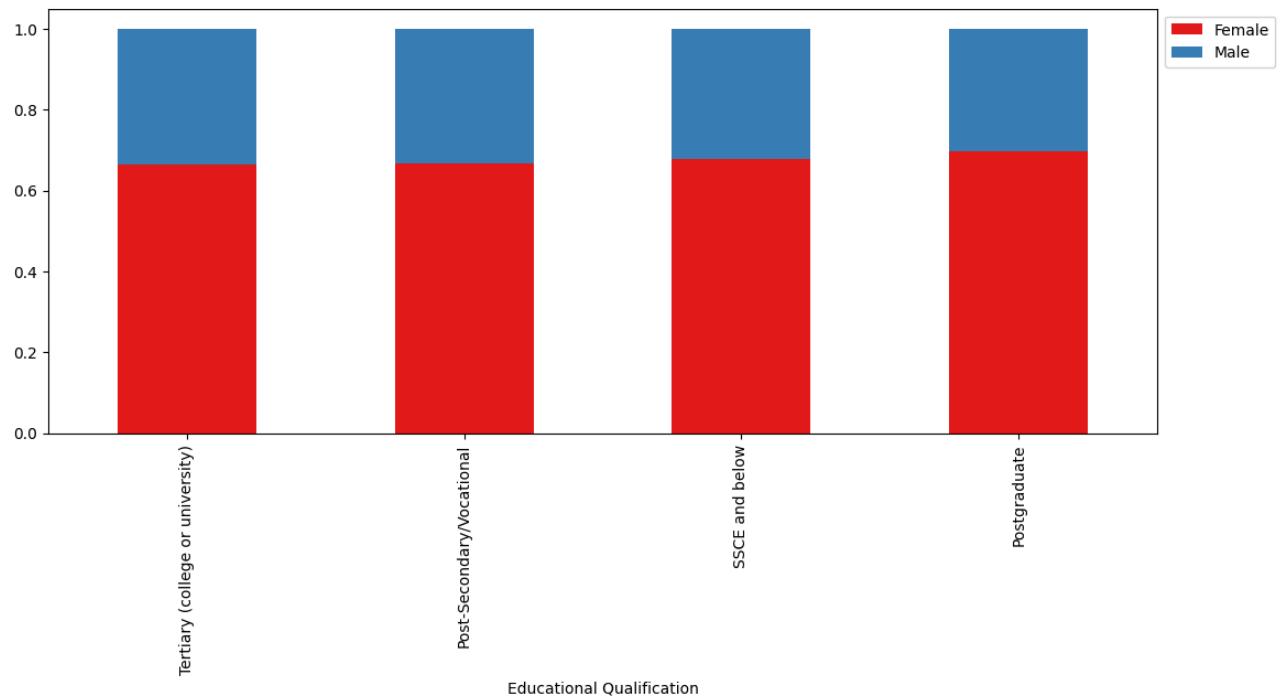
```
In [58]: 1 stacked_barplot(data, "Have you ever witnessed any form of harassment on social media?", "Gender")
```

Gender	Female	Male	All
Have you ever witnessed any form of harassment ...			
All	533	258	791
Yes	262	140	402
No	251	113	364
Uncertain	20	5	25



```
In [59]: 1 stacked_barplot(data, "Educational Qualification", "Gender")
```

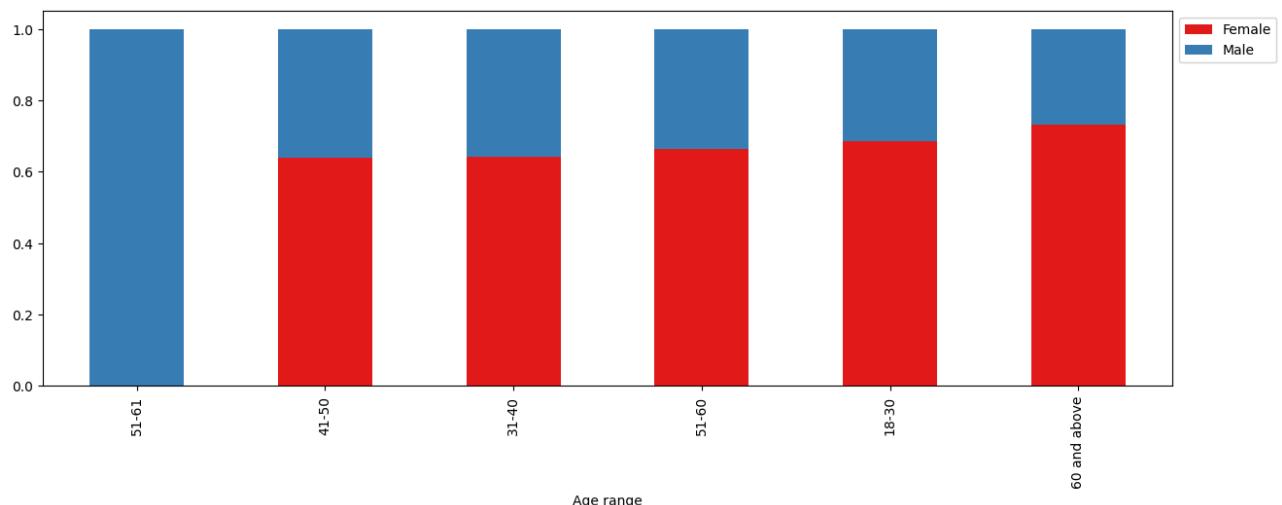
Gender	Female	Male	All
Educational Qualification			
All	533	258	791
Tertiary (college or university)	311	156	467
SSCE and below	95	45	140
Postgraduate	101	44	145
Post-Secondary/Vocational	26	13	39



Gender vs Age range

```
In [60]: 1 stacked_barplot(data, "Age range", "Gender")
```

Gender	Female	Male	All
Age range			
All	533	258	791
18-30	261	119	380
31-40	111	62	173
51-60	77	39	116
60 and above	52	19	71
41-50	32	18	50
51-61	0	1	1



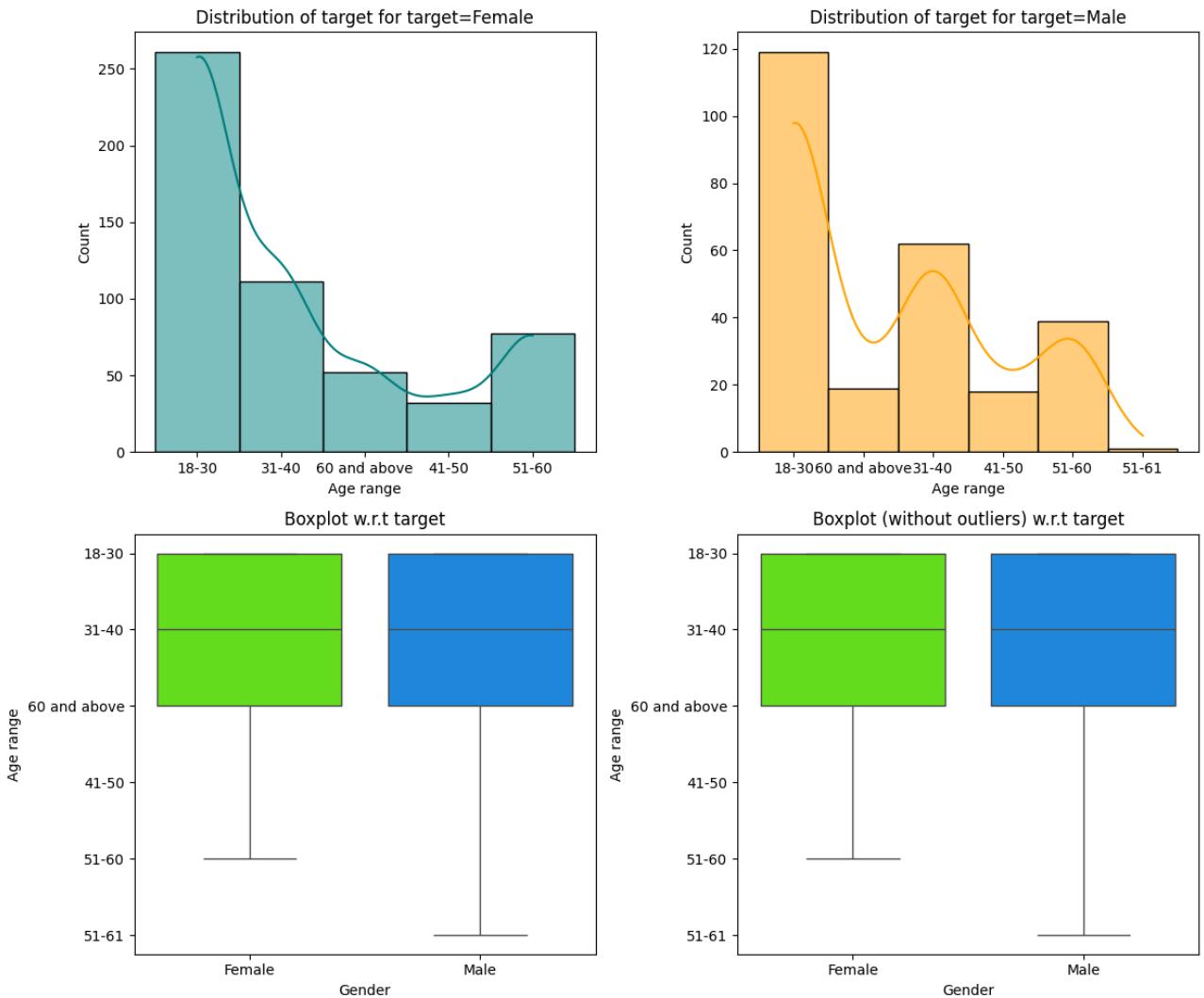
```

In [61]: 1  ### Function to plot distributions
2
3
4 def distribution_plot_wrt_target(data, predictor, target):
5
6     fig, axs = plt.subplots(2, 2, figsize=(12, 10))
7
8     target_uniq = data[target].unique()
9
10    axs[0, 0].set_title("Distribution of target for target=" + str(target_uniq[0]))
11    sns.histplot(
12        data=data[data[target] == target_uniq[0]],
13        x=predictor,
14        kde=True,
15        ax=axs[0, 0],
16        color="teal",
17    )
18
19    axs[0, 1].set_title("Distribution of target for target=" + str(target_uniq[1]))
20    sns.histplot(
21        data=data[data[target] == target_uniq[1]],
22        x=predictor,
23        kde=True,
24        ax=axs[0, 1],
25        color="orange",
26    )
27
28    axs[1, 0].set_title("Boxplot w.r.t target")
29    sns.boxplot(data=data, x=target, y=predictor, ax=axs[1, 0], palette="gist_rainbow")
30
31    axs[1, 1].set_title("Boxplot (without outliers) w.r.t target")
32    sns.boxplot(
33        data=data,
34        x=target,
35        y=predictor,
36        ax=axs[1, 1],
37        showfliers=False,
38        palette="gist_rainbow",
39    )
40
41    plt.tight_layout()
42    plt.show()

```

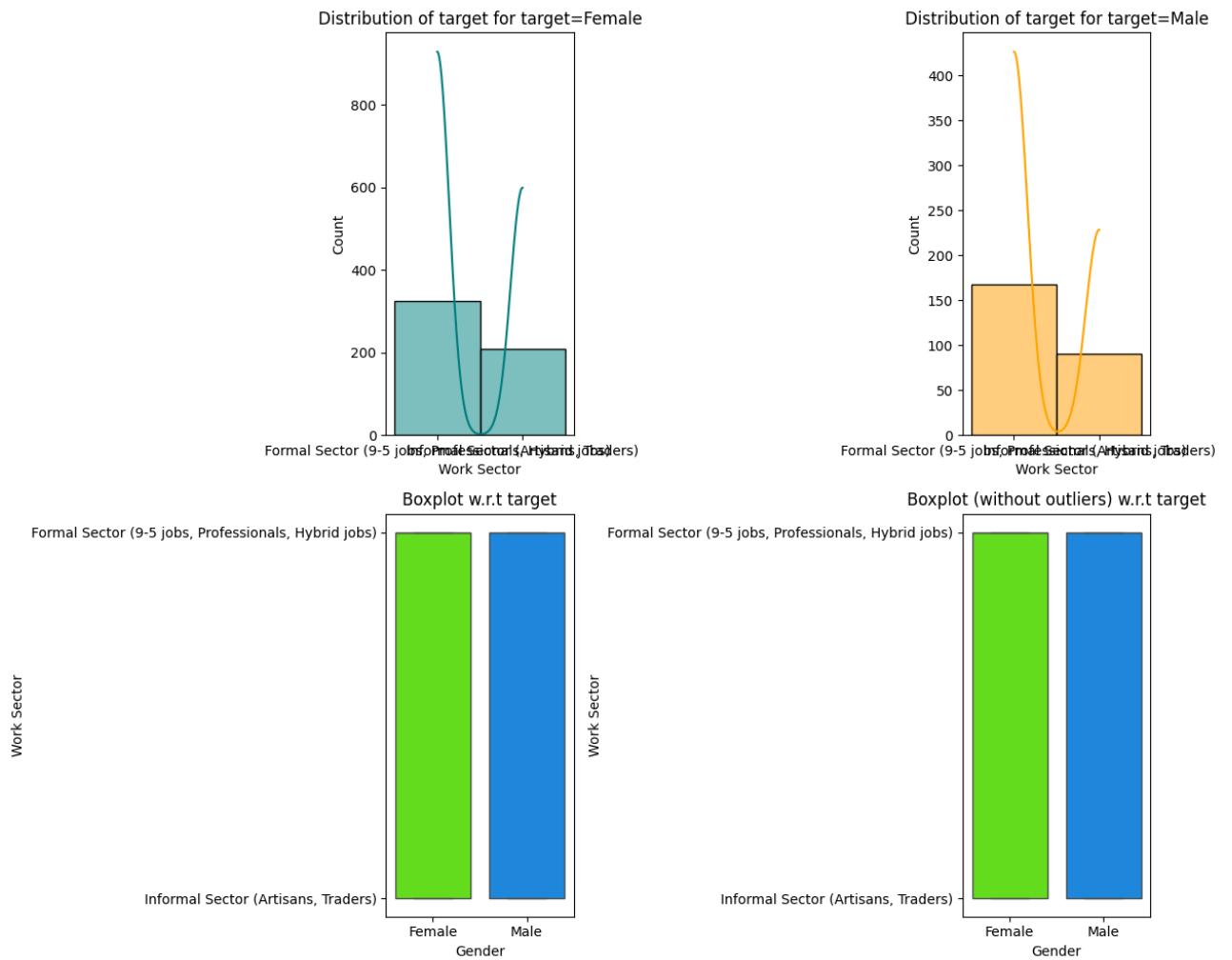
Gender vs Age range

```
In [62]: 1 distribution_plot_wrt_target(data, "Age range", "Gender")
```



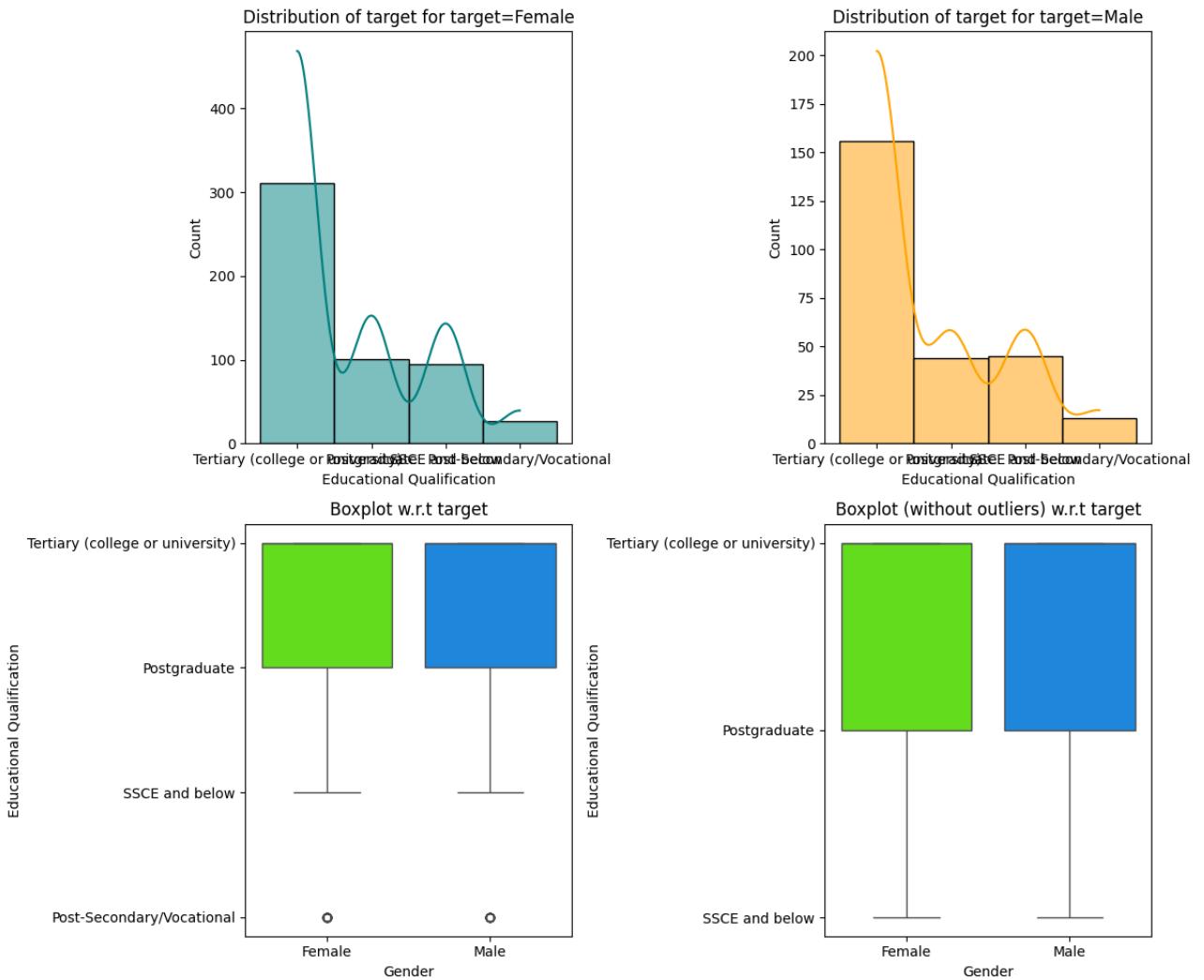
Gender vs Work Sector

```
In [63]: 1 distribution_plot_wrt_target(data, "Work Sector", "Gender")
```



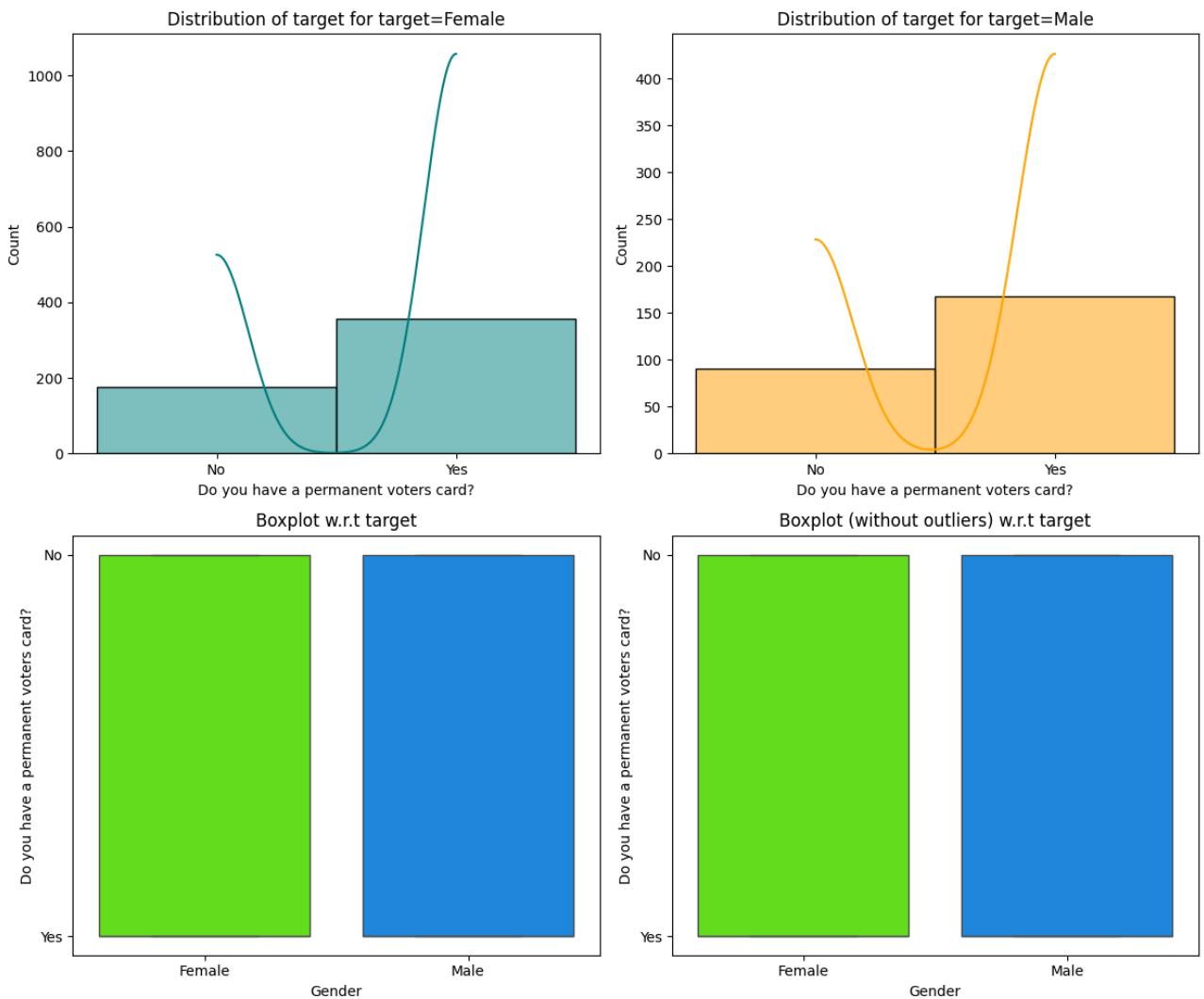
Gender vs Educational Qualification

```
In [64]: 1 distribution_plot_wrt_target(data, "Educational Qualification", "Gender")
```



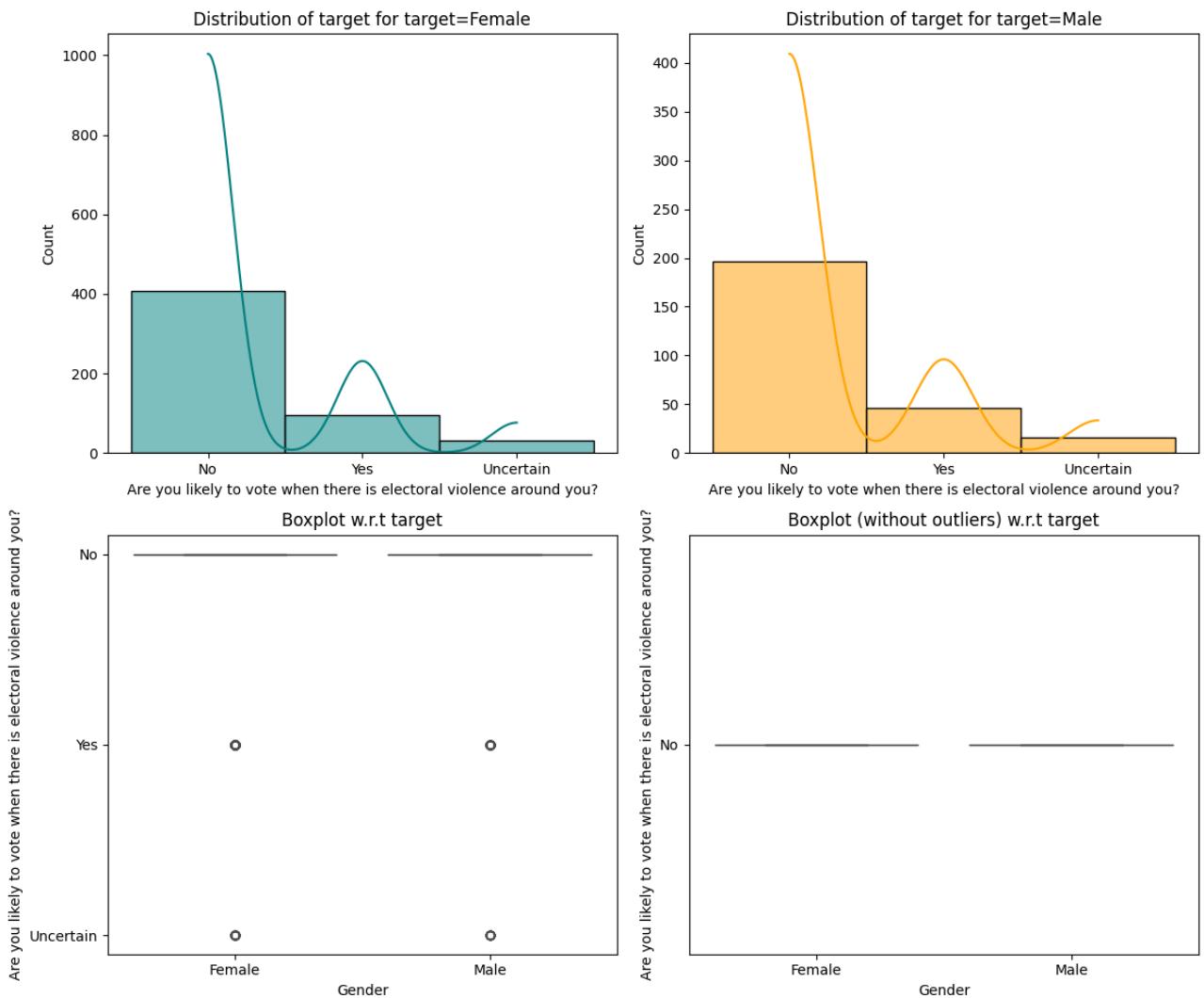
Gender vs Do you have a permanent voters card?

```
In [65]: 1 distribution_plot_wrt_target(data, "Do you have a permanent voters card?", "Gender")
```

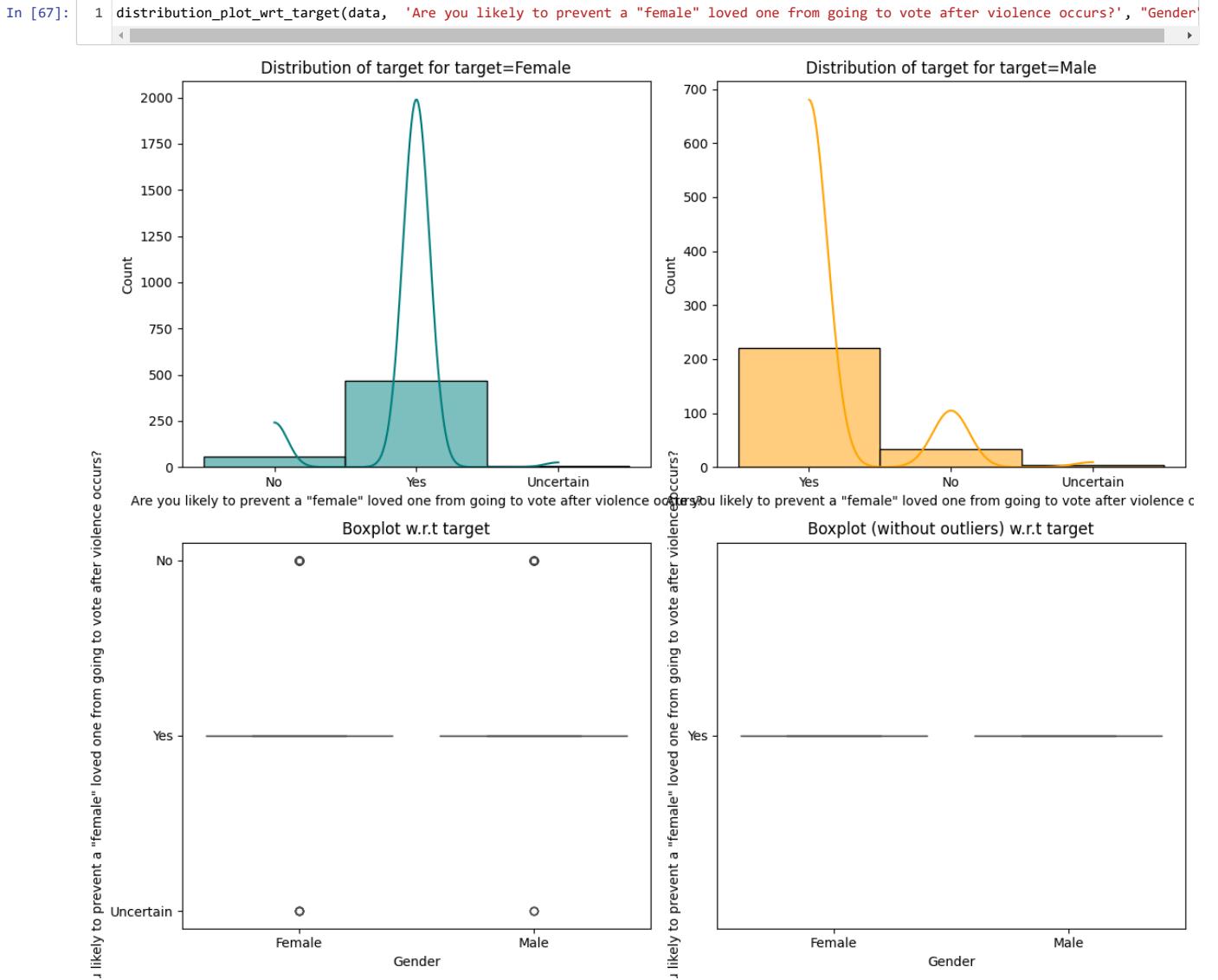


Gender vs Are you likely to vote when there is electoral violence around you?

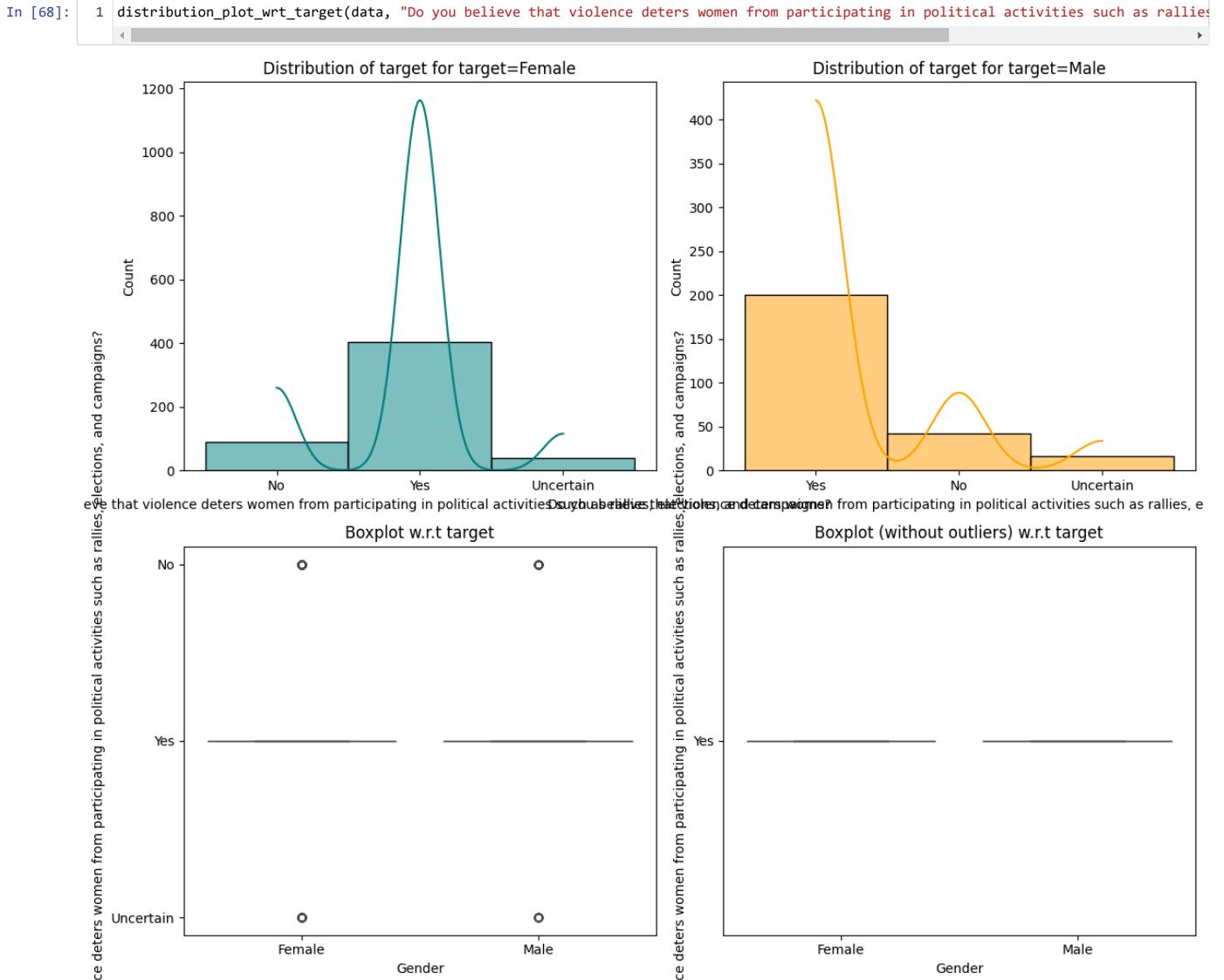
```
In [66]: 1 distribution_plot_wrt_target(data, "Are you likely to vote when there is electoral violence around you?", "Gender")
```



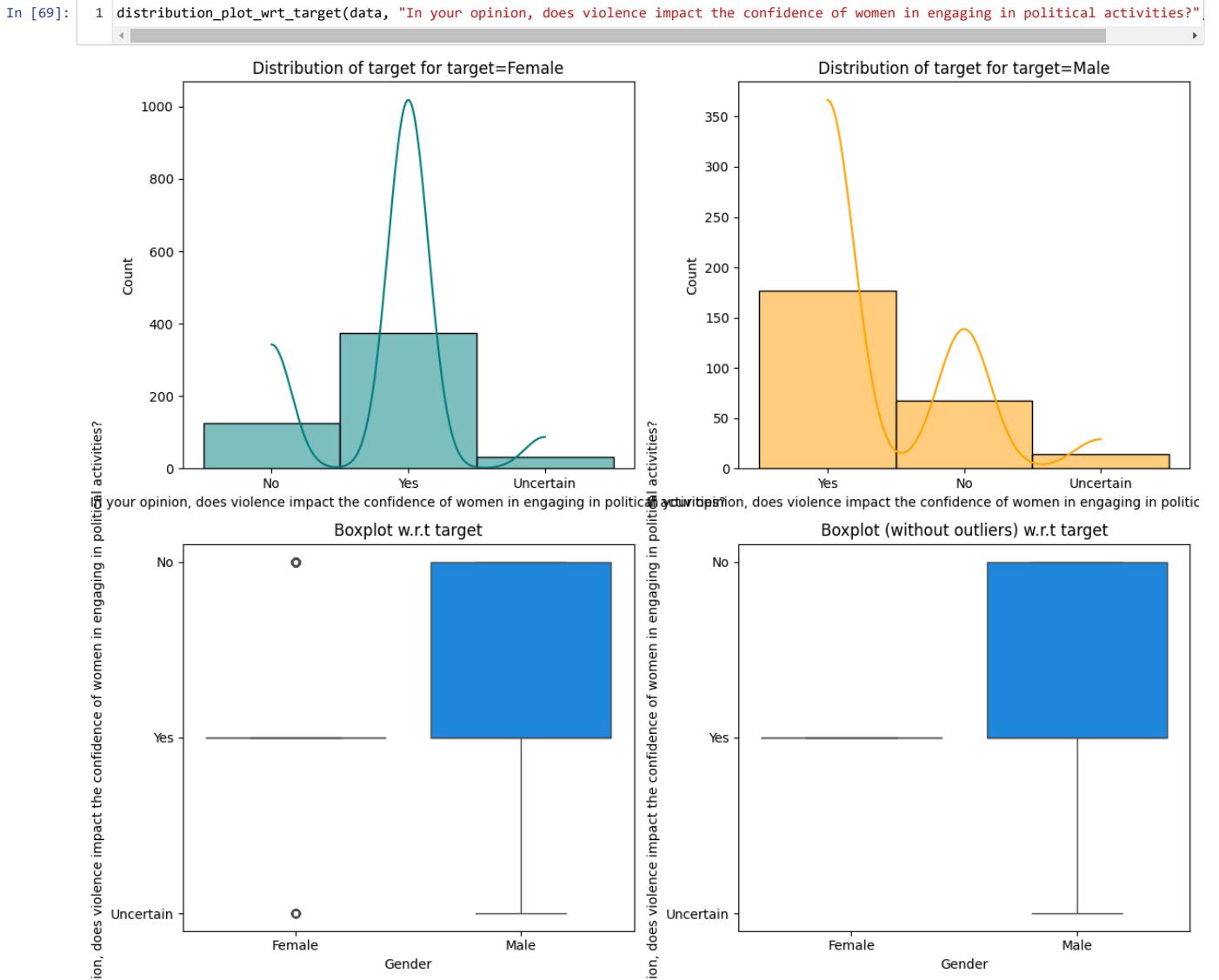
Gender vs Are you likely to prevent a "female" loved one from going to vote after violence occurs?



Gender vs Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?

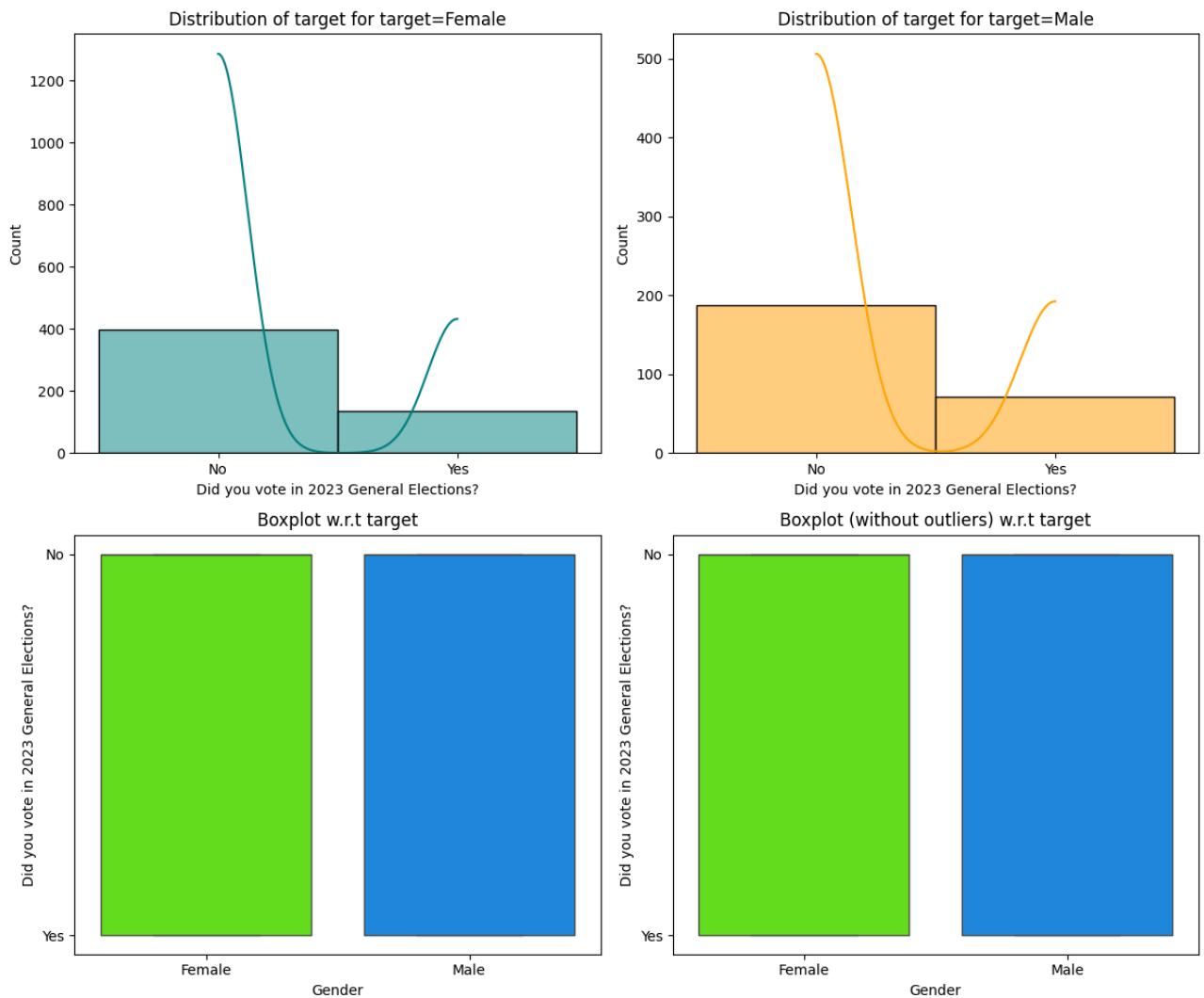


Gender vs In your opinion, does violence impact the confidence of women in engaging in political activities?



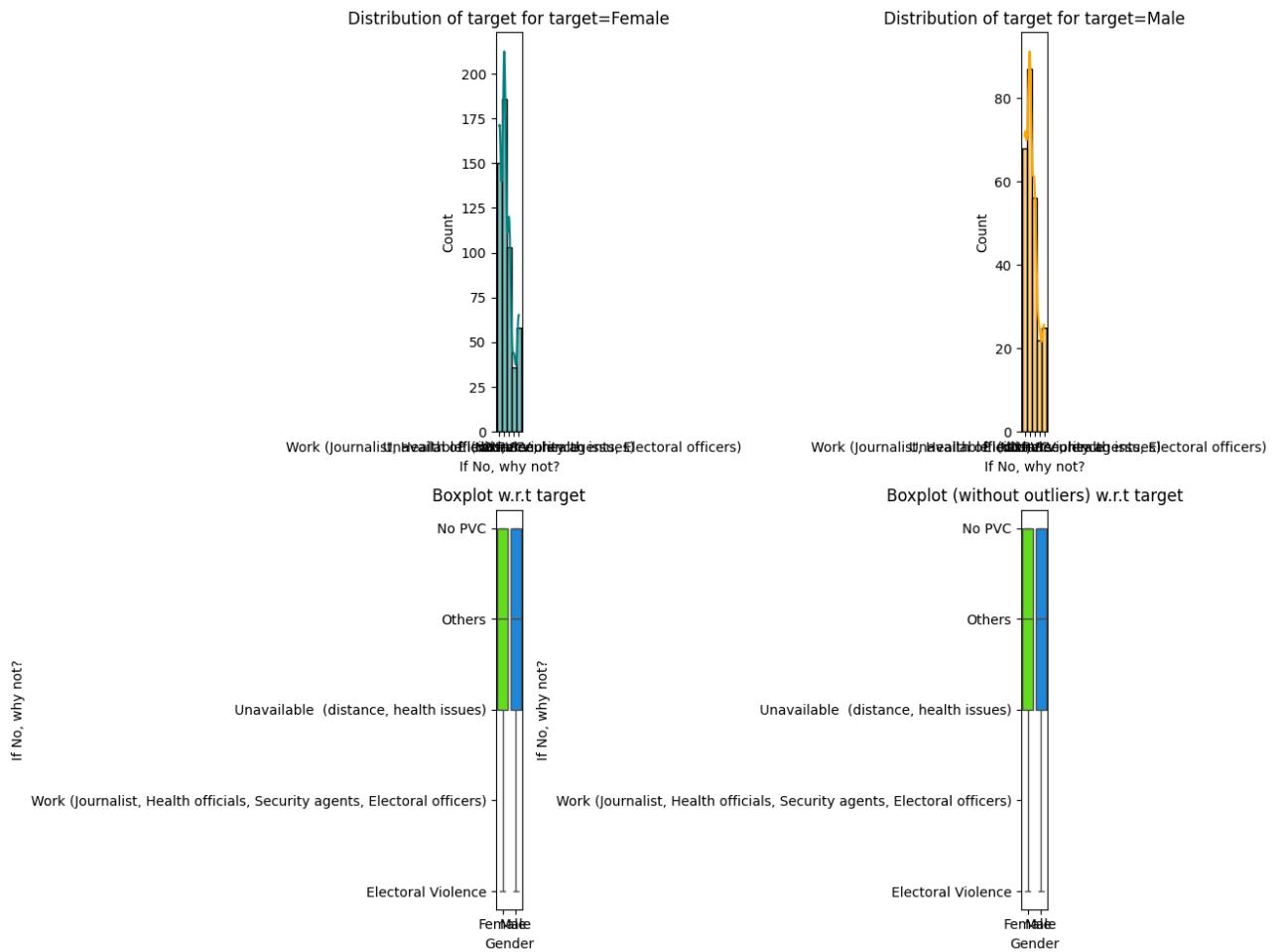
Gender vs Did you vote in 2023 General Elections?

```
In [70]: 1 distribution_plot_wrt_target(data, "Did you vote in 2023 General Elections?", "Gender")
```

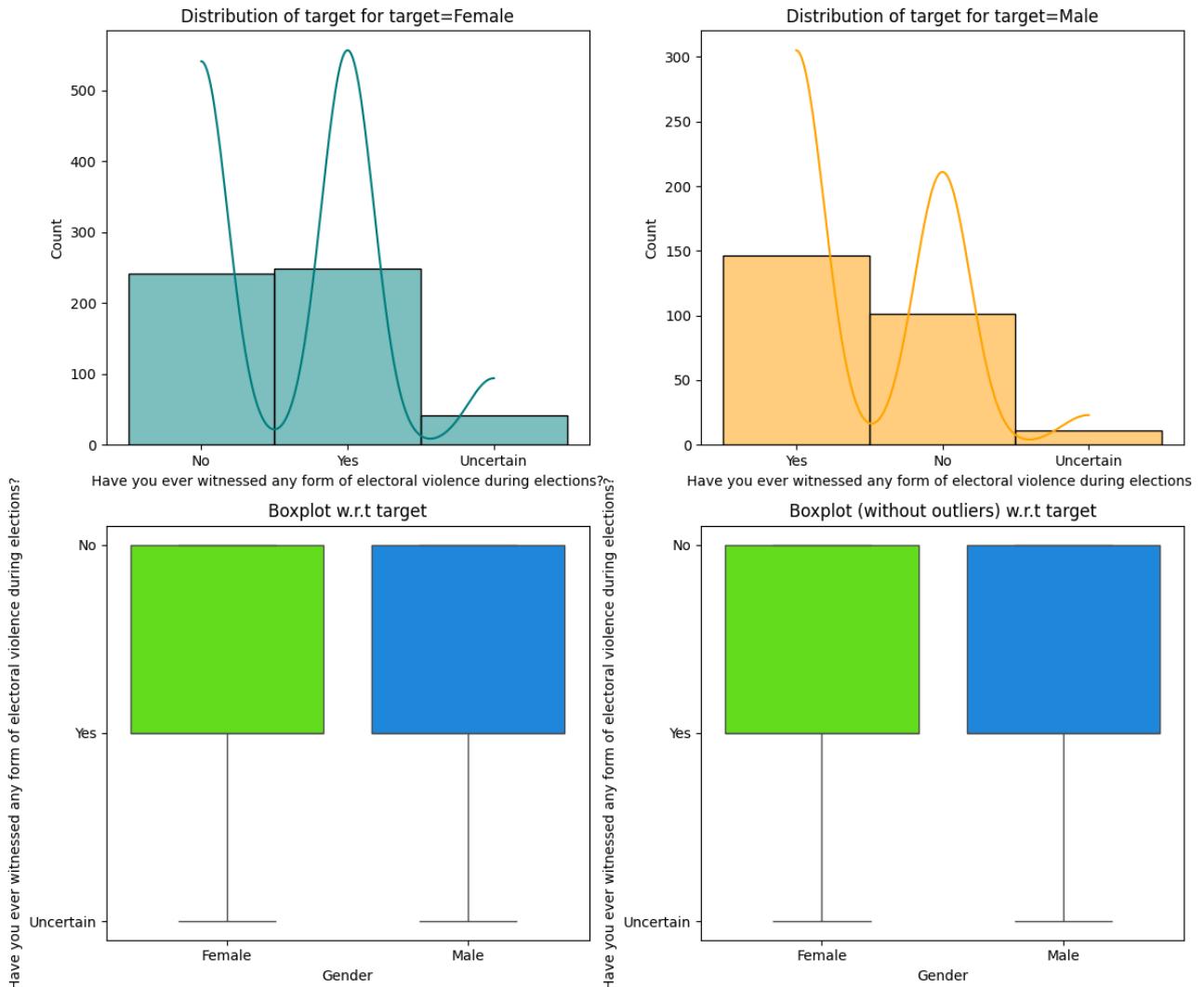


Gender vs If No, why not?

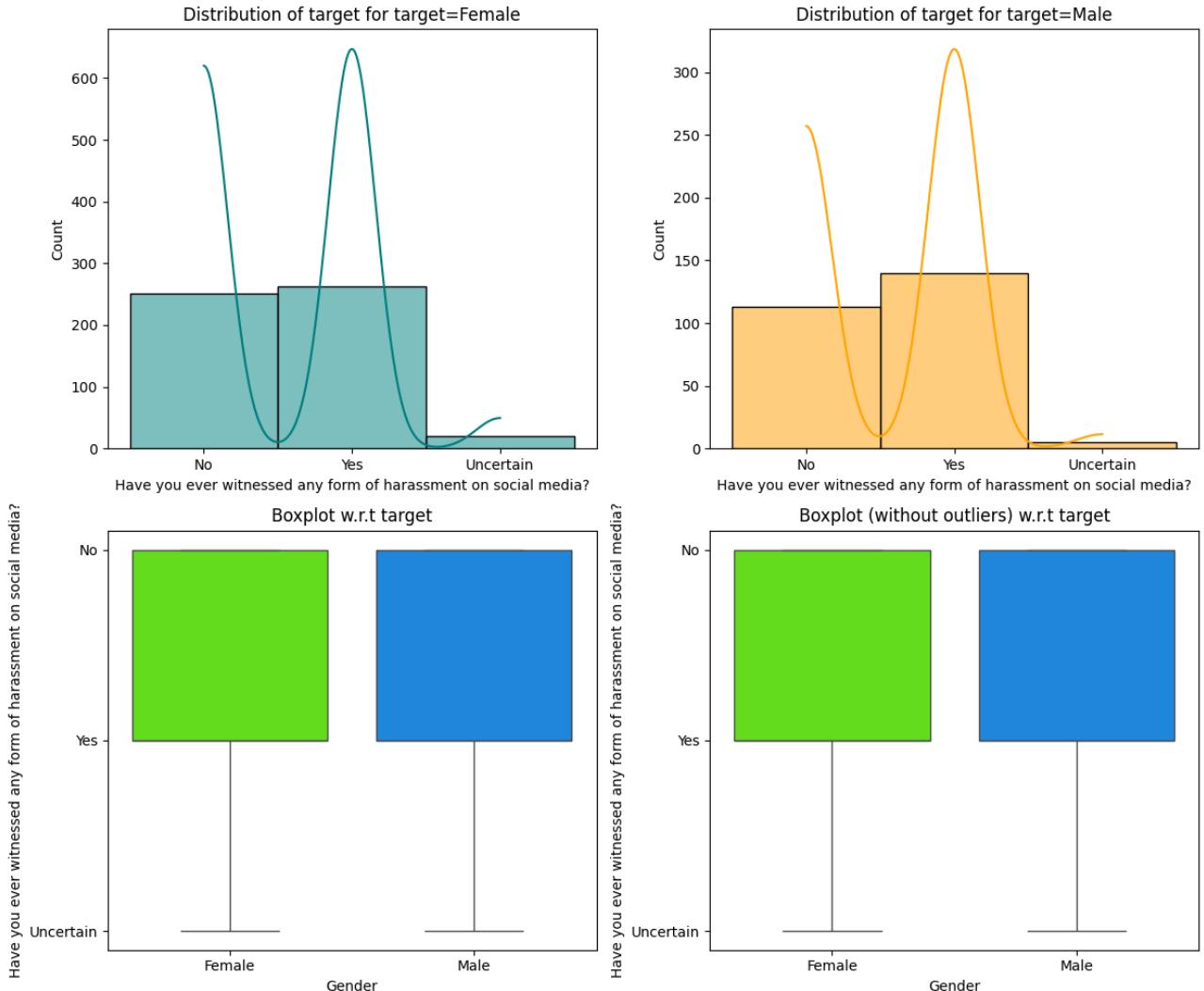
```
In [71]: 1 distribution_plot_wrt_target(data, "If No, why not?", "Gender")
```



```
In [72]: 1 distribution_plot_wrt_target(data, "Have you ever witnessed any form of electoral violence during elections?", "Gender")
```



```
In [73]: 1 distribution_plot_wrt_target(data, "Have you ever witnessed any form of harassment on social media?", "Gender")
```



Multivariate Analysis

```
In [74]: 1 def line_plot( cols):
2     num_plots = len(cols)
3     rows = (num_plots + 2) // 3 # Determine the number of rows needed
4
5     plt.figure(figsize=(23, rows * 4)) # Adjust the figure size based on the number of rows
6
7     for i, variable in enumerate(cols):
8         plt.subplot(rows, 3, i + 1)
9
10    # Count the occurrences of "Yes" and "No" for the given variable, grouped by Educational Qualification and Gender
11    counts = data.groupby(['Educational Qualification', 'Gender', variable]).size().reset_index(name='Count')
12    # counts = data.groupby(['Edu. Qlf.', 'Gender', variable]).size().reset_index(name='Count')
13
14    # Plot the data with different markers for "Yes" and "No"
15    sns.lineplot(data=counts, x="Educational Qualification", y='Count', hue="Gender", style=variable, markers=True, dashes=False)
16
17    plt.ylabel(f"{variable}") # Set y-axis Label
18    plt.title(variable)
19    plt.tight_layout()
20
21    plt.show()
```

In [75]: 1 dataCorr

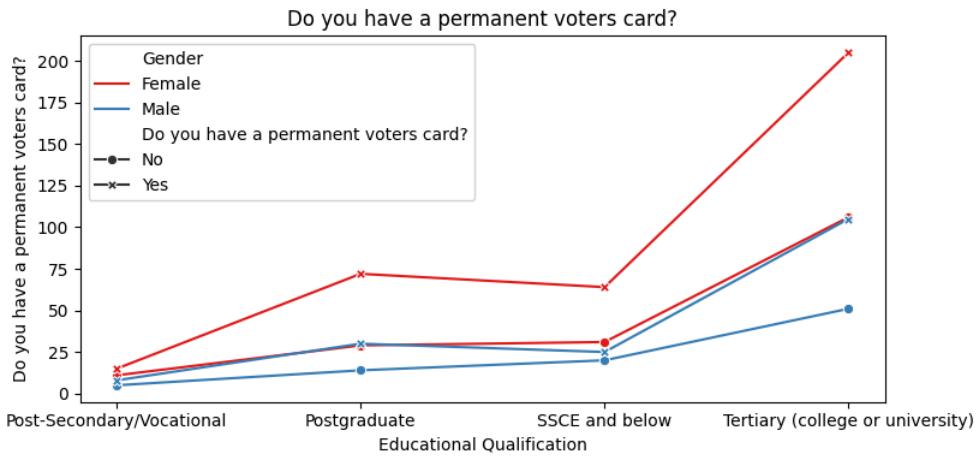
Out[75]:

	Gender	Work Sector	Edu. Qlf.	Age range	PVC	Vote During Elec. Vio.	Allow Female Vote During Elec. Vio.	Violence Deters Women From Parti.	Violence Impact Women Conf. In Parti.	Vote in 2023 Gen. Elec.	If No, why not?	Witnessed any Elec. Vio.	Witnessed Haras. Social Media
0	0	1	15	24.0	0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.0
1	0	1	15	24.0	0	0.0	1.0	1.0	1.0	0	0.3	1.0	1.0
2	0	0	15	35.5	1	0.0	1.0	1.0	1.0	0	0.5	0.0	0.0
3	1	1	15	24.0	0	0.0	1.0	1.0	1.0	0	0.0	1.0	0.0
4	0	0	15	24.0	1	0.0	1.0	1.0	1.0	1	0.3	1.0	1.0
...
786	0	1	5	35.5	1	0.0	1.0	0.5	1.0	1	0.3	0.5	1.0
787	0	0	15	80.0	0	0.0	1.0	0.0	0.0	0	0.5	1.0	0.0
788	1	1	20	35.5	0	0.5	1.0	1.0	1.0	0	0.0	1.0	0.0
789	1	0	15	55.5	1	0.0	1.0	0.5	0.5	0	0.3	0.0	0.0
790	0	1	15	24.0	1	1.0	1.0	1.0	1.0	1	0.3	1.0	0.0

791 rows × 13 columns

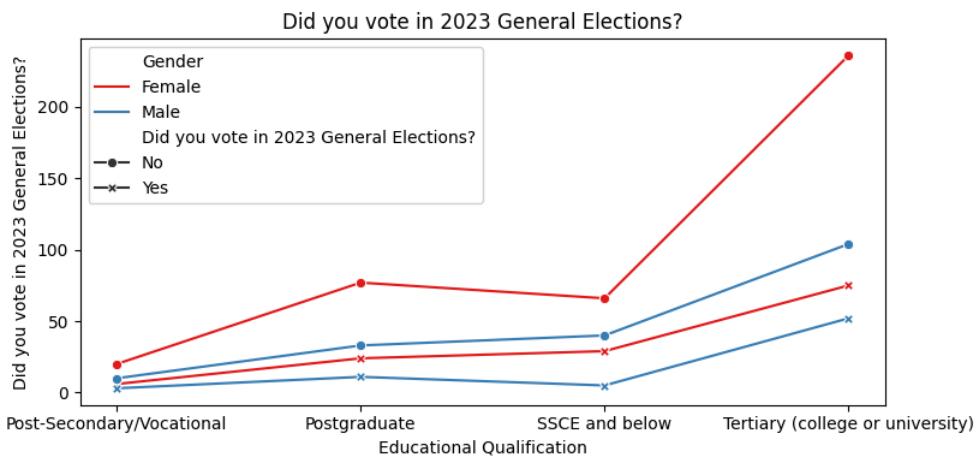
Gender vs Education Qualification vs Do you have a permanent voters card?

```
In [76]: 1 cols = data[
2     ["Do you have a permanent voters card?"]
3 ].columns.tolist()
4 line_plot(cols)
```



Gender vs Education Qualification vs Did you vote in 2023 General Elections?

```
In [77]: 1 cols = data[
2     ["Did you vote in 2023 General Elections?"]
3 ].columns.tolist()
4 line_plot(cols)
```

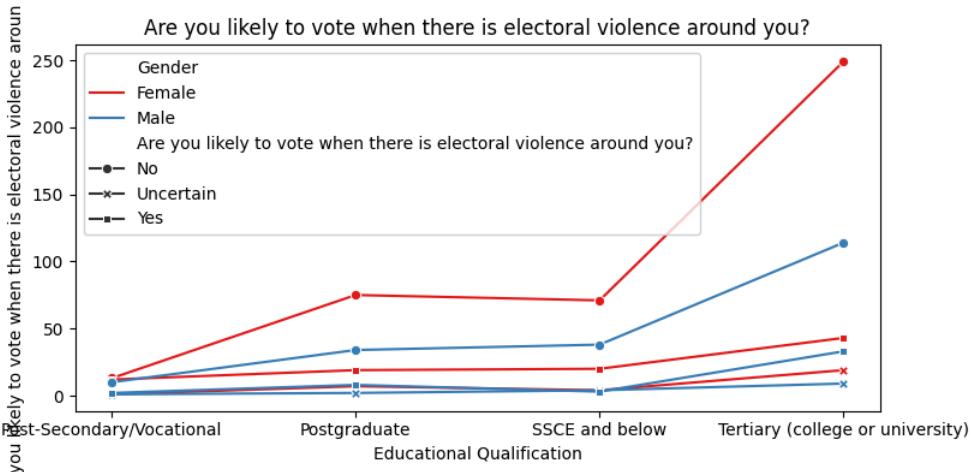


Observation

- Majority of the females that do not vote in the 2023 General Election have Tertiary Education Qualification
- Majority of the female that have PVC have Tertiary Education Qualification

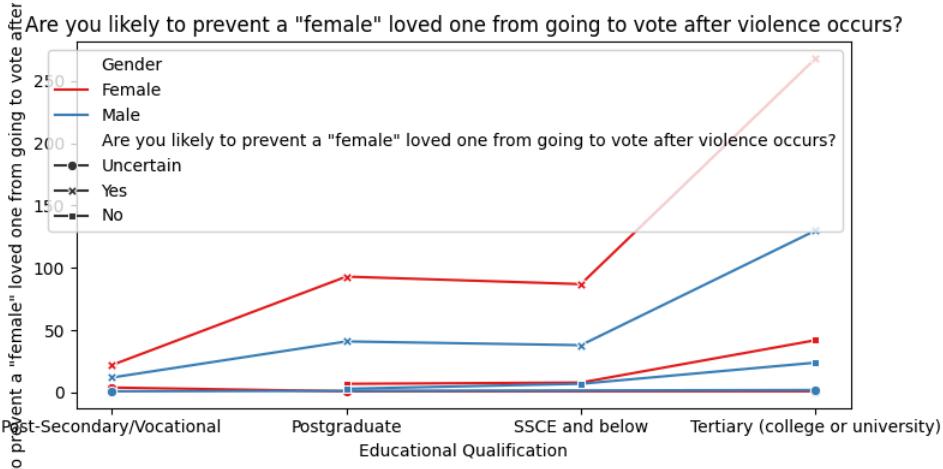
Gender vs Education Qualification vs Are you likely to vote when there is electoral violence around you?

```
In [78]: 1 cols = data[
2     ["Are you likely to vote when there is electoral violence around you?"]
3 ].columns.tolist()
4 line_plot(cols)
```



Gender vs Education Qualification vs Are you likely to prevent a "female" loved one from going to vote after violence occurs?

```
In [79]: 1 cols = data[
2     ['Are you likely to prevent a "female" loved one from going to vote after violence occurs?']
3 ].columns.tolist()
4 line_plot(cols)
```



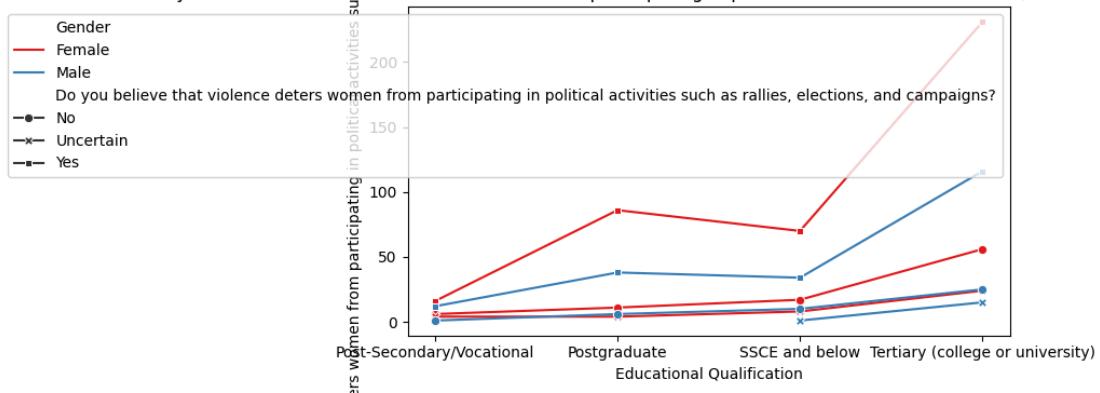
Observation

- Majority of the females that their education level is Tertiary (college or university) indicate they are not likely to vote when there is electoral violence around them
- Although more male shows that are not likely to vote when there is electoral violence around them
- Females with Post-Secondary/Vocational qualification still show willingness to vote despite electoral violence around them
- only Females with Post-Secondary/Vocational are uncertain to prevent a "female" loved one's from going to vote after violence occurs

Gender vs Education Qualification vs Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?

```
In [80]: 1 cols = data[
2     ["Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?"]
3     ]
4 ].columns.tolist()
5 line_plot(cols)
6
```

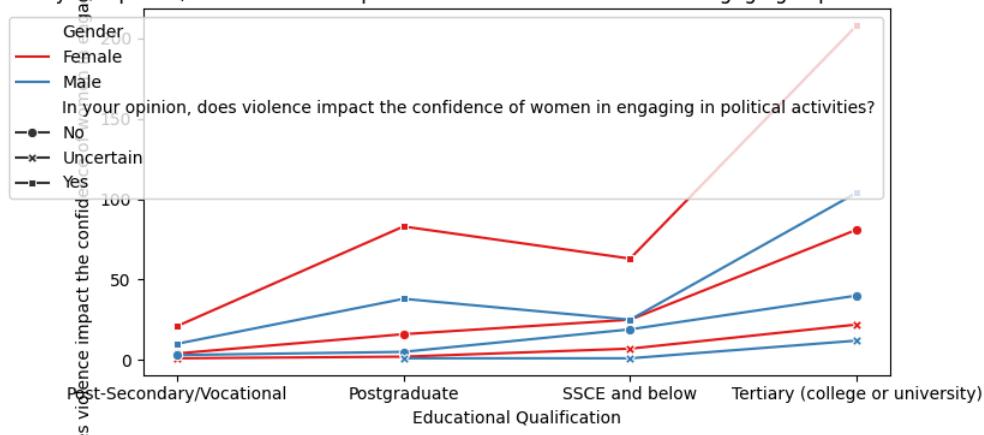
Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?



Gender vs Education Qualification vs In your opinion, does violence impact the confidence of women in engaging in political activities?

```
In [81]: 1 cols = data[
2     ['In your opinion, does violence impact the confidence of women in engaging in political activities?']
3     ]
4 ].columns.tolist()
5 line_plot(cols)
6
```

In your opinion, does violence impact the confidence of women in engaging in political activities?

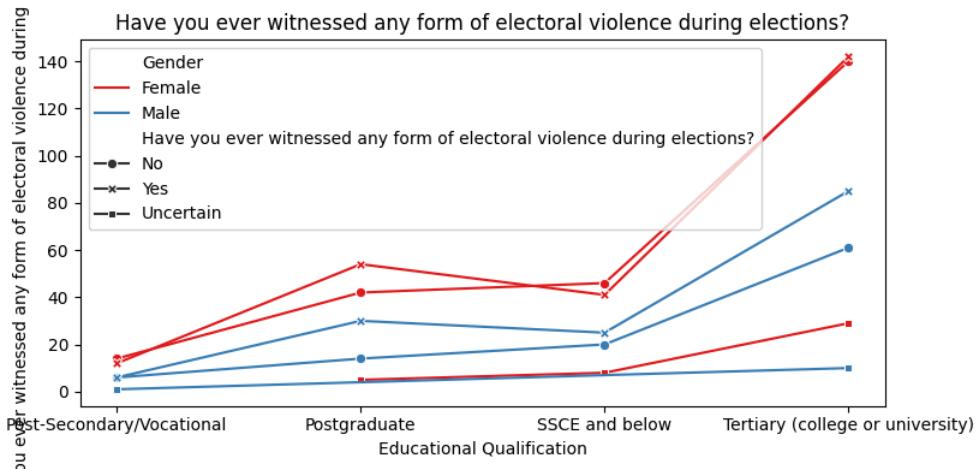


Observation

- Male respondent who have Tertiary (college or university) qualification are **uncertain** if violence deters women from participating in practical activities such as as rallies, elections, and campaigns
- Male respondent who have SSCE and Below claims that violence **does not** deter women from participating in practical activities such as as rallies, elections, and campaigns
- All the female respondent, irrespective of their educational level claims that **agree** that violence deters women from participating in practical activities such as as rallies, elections, and campaigns

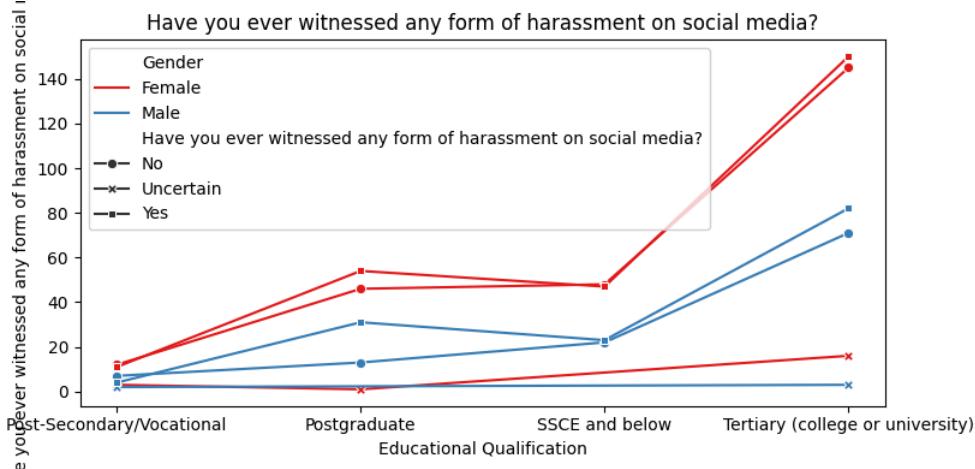
Gender vs Education Qualification vs Have you ever witnessed any form of electoral violence during elections?

```
In [82]: 1 cols = data[
2     ["Have you ever witnessed any form of electoral violence during elections?"]
3 ].columns.tolist()
4 line_plot(cols)
```



Gender vs Education Qualification vs Have you ever witnessed any form of harassment on social media?

```
In [83]: 1 cols = data[
2     ['Have you ever witnessed any form of harassment on social media?']
3 ].columns.tolist()
4 line_plot(cols)
```



Observation

- Male respondent who have BTech and undergraduate are the lowest among all the respondent that has ever witnessed any form of electoral violence during elections
- All the female respondent, irrespective of their educational level claims that **agree** that they have witnessed any form of electoral violence during elections with females with BTech being the highest.
- Male respondent who have BTech and Mbbs in view are the lowest among all the respondent that has ever witnessed any form of harrassment on social media
- All the female respondent, irrespective of their educational level claims that **agree** that they have witnessed any form of harrassment on social media with females with undergraduate being the highest follow by Mbbs in view.

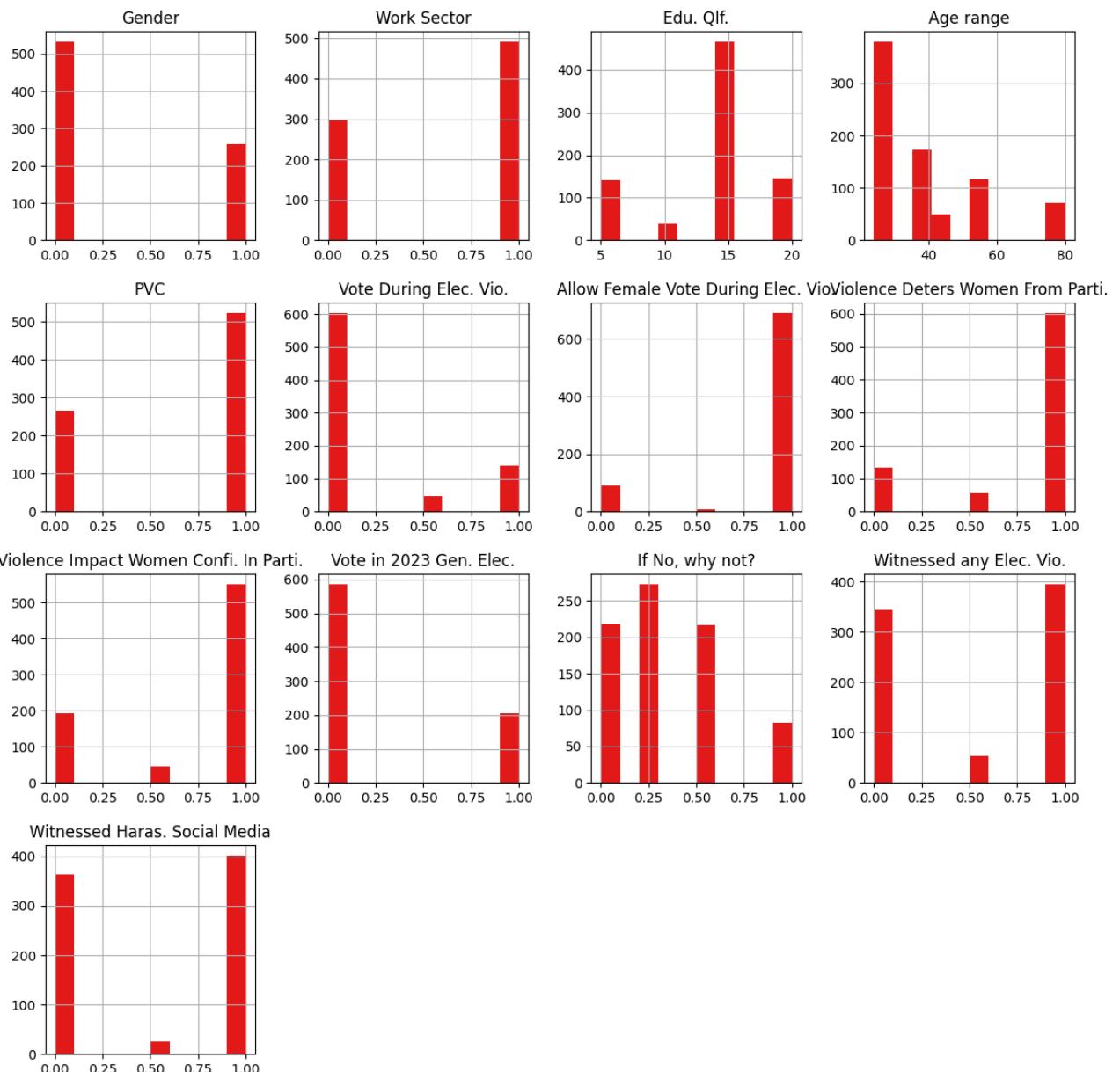
Summary of EDA

- The dataset has total of 533 Female and 258 Male respondent
- There are no missing values in the dataset
- High level of correlation exist between some of the variables.
- There are also outliers

Data Cleaning

- Timestamp column contains uniques ID for response. This column has been dropped.
- The target variable is encoded to numeric.

```
In [84]: 1 # creating histograms
2 dataCorr.hist(figsize=(14, 14))
3 plt.show()
```



Let's find the percentage of outliers, in each column of the data, using IQR.

```
In [85]: 1 Q1 = dataCorr.quantile(0.25) # To find the 25th percentile and 75th percentile.
2 Q3 = dataCorr.quantile(0.75)
3
4 IQR = Q3 - Q1 # Inter Quantile Range (75th percentile - 25th percentile)
5
6 lower = (
7     Q1 - 1.5 * IQR
8 ) # Finding Lower and upper bounds for all values. All values outside these bounds are outliers
9 upper = Q3 + 1.5 * IQR
```

```
In [86]: 1 lower
```

```
Out[86]: Gender           -1.50
Work Sector        -1.50
Edu. Qlf.          15.00
Age range          -8.25
PVC                -1.50
Vote During Elec. Vio. 0.00
Allow Female Vote During Elec. Vio. 1.00
Violence Deters Women From Parti. 1.00
Violence Impact Women Confi. In Parti. -0.25
Vote in 2023 Gen. Elec. -1.50
If No, why not?    -0.75
Witnessed any Elec. Vio. -1.50
Witnessed Haras. Social Media -1.50
dtype: float64
```

```
In [87]: 1 upper
```

```
Out[87]: Gender           2.50
Work Sector        2.50
Edu. Qlf.          15.00
Age range          77.75
PVC                2.50
Vote During Elec. Vio.    0.00
Allow Female Vote During Elec. Vio. 1.00
Violence Deters Women From Parti. 1.00
Violence Impact Women Confi. In Parti. 1.75
Vote in 2023 Gen. Elec. 2.50
If No, why not?      1.25
Witnessed any Elec. Vio. 2.50
Witnessed Haras. Social Media 2.50
dtype: float64
```

```
In [88]: 1 (
2     (dataCorr.select_dtypes(include=["float64", "int64"]) < lower)
3     | (dataCorr.select_dtypes(include=["float64", "int64"]) > upper)
4 ).sum() / len(data) * 100
```

```
Out[88]: Gender           0.000000
Work Sector        0.000000
Edu. Qlf.          40.960809
Age range          8.975980
PVC                0.000000
Vote During Elec. Vio.    23.640961
Allow Female Vote During Elec. Vio. 12.642225
Violence Deters Women From Parti. 23.767383
Violence Impact Women Confi. In Parti. 0.000000
Vote in 2023 Gen. Elec. 0.000000
If No, why not?      0.000000
Witnessed any Elec. Vio. 0.000000
Witnessed Haras. Social Media 0.000000
dtype: float64
```

Observation

- After identifying outliers, we can decide whether to remove/treat them or not. It depends on one's approach, here we are not going to treat them as there will be outliers in real case scenario (in Education qualification, Vote During Elec. Vio., Allow Female Vote During Elec. Vio., etc) and we would want our model to learn the underlying pattern for such customers.

```
In [89]: 1 data1=[]
2 data1 = dataCorr.copy()
```

```
In [90]: 1 imputer = SimpleImputer(strategy="most_frequent")
```

```
In [91]: 1 X = data1.drop(['Allow Female Vote During Elec. Vio.'], axis=1)
2 y = data1['Allow Female Vote During Elec. Vio.']
```

```
In [92]: 1 # from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)
3
```

```
In [93]: 1 reqd_col_for_impute = ["Edu. Qlf.", "Gender", "Vote During Elec. Vio."]
2
3 # reqd_col_for_impute = ["Educational Qualification", "Gender", "Are you likely to vote when there is electoral violence around you?"]
```

Model Building

Model evaluation criterion

The nature of predictions made by the classification model will translate as follows:

- True positives (TP) are failures correctly predicted by the model.
- False negatives (FN) are real failures in a generator where there is no detection by model.
- False positives (FP) are failure detections in a generator where there is no failure.
- Accuracy
- Precision
- Recall
- F1-Score
- ROC-Curve

Which metric to optimize?

- We need to choose the metric which will ensure that the maximum number of generator failures are predicted correctly by the model.
- We would want Recall to be maximized as greater the Recall, the higher the chances of minimizing false negatives.
- We want to minimize false negatives because if a model predicts that a machine will have no failure when there will be a failure, it will increase the maintenance cost.

Let's define a function to output different metrics (including recall) on the train and test set and a function to show confusion matrix so that we do not have to use the same code repetitively while evaluating models.

Model with original data

```
In [94]: 1 # Define bin edges based on the known ranges for the relevant classes
2 bins = np.linspace(0, 1, 3) # Adjust based on your data
3 y_train_binned = np.digitize(y_train, bins) - 1
4 y_test_binned = np.digitize(y_test, bins) - 1
5
6 # Reassign your labels
7 y_train = y_train_binned.copy()
8 y_test = y_test_binned.copy()
9
```

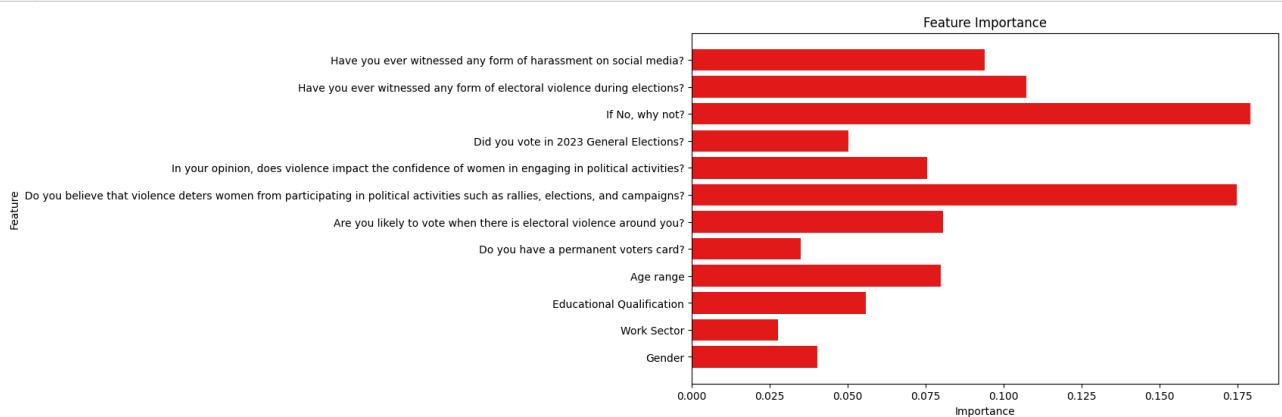
```
In [95]: 1 # Initialize the Random Forest model
2 rf_model = RandomForestClassifier(random_state=1)
3
4 # Fit the model on the training data
5 rf_model.fit(X_train, y_train)
6
7 # Predictions on the training data
8 y_pred_rf_train = rf_model.predict(X_train)
9
10 # Calculate training performance
11 rf_train_accuracy = accuracy_score(y_train, y_pred_rf_train)
12 rf_train_recall = recall_score(y_train, y_pred_rf_train, average='weighted')
13
14 # Print training performance metrics
15 print("Training Performance of Random Forest Model:")
16 print(f"Random Forest Training Accuracy: {rf_train_accuracy:.3f}")
17 print(f"Random Forest Training Recall: {rf_train_recall:.3f}")
18
```

Training Performance of Random Forest Model:
Random Forest Training Accuracy: 0.988
Random Forest Training Recall: 0.988

```
In [96]: 1 # Cross-validated metrics
2 cv_scores_accuracy = cross_val_score(rf_model, X_train, y_train, cv=10, scoring='accuracy')
3 cv_scores_recall = cross_val_score(rf_model, X_train, y_train, cv=10, scoring='recall_weighted')
4 cv_scores_precision = cross_val_score(rf_model, X_train, y_train, cv=10, scoring='precision_weighted')
5 cv_scores_f1 = cross_val_score(rf_model, X_train, y_train, cv=10, scoring='f1_weighted')
6 cv_scores_roc_auc = cross_val_score(rf_model, X_train, y_train, cv=10, scoring='roc_auc')
7 cv_scores_log_loss = cross_val_score(rf_model, X_train, y_train, cv=10, scoring='neg_log_loss')
8
9 # Average metrics
10 average_cv_accuracy = cv_scores_accuracy.mean()
11 average_cv_recall = cv_scores_recall.mean()
12 average_cv_precision = cv_scores_precision.mean()
13 average_cv_f1 = cv_scores_f1.mean()
14 average_cv_roc_auc = cv_scores_roc_auc.mean()
15 average_cv_log_loss = -cv_scores_log_loss.mean() # Negate to get positive Log Loss
16
17 # Print results
18 print("10-fold Cross validation Training Performance of Random Forest Model:")
19 print(f"Average CV Accuracy: {average_cv_accuracy:.3f}")
20 print(f"Average CV Recall: {average_cv_recall:.3f}")
21 print(f"Average CV Precision: {average_cv_precision:.3f}")
22 print(f"Average CV F1 Score: {average_cv_f1:.3f}")
23 # print(f"Average CV ROC AUC: {average_cv_roc_auc:.3f}")
24 # print(f"Average CV Log Loss: {average_cv_log_loss:.3f}")
25
```

10-fold Cross validation Training Performance of Random Forest Model:
Average CV Accuracy: 0.961
Average CV Recall: 0.961
Average CV Precision: 0.961
Average CV F1 Score: 0.960

```
In [97]: 1 # Extract feature importances
2 importances = rf_model.feature_importances_
3
4 # Use the columns from X_train as feature names, then map them to the original column names from `data`
5 short_to_full_names = {
6     'Gender': 'Gender',
7     'Work Sector': 'Work Sector',
8     'Edu. Qlf.': 'Educational Qualification',
9     'Age range': 'Age range',
10    'PVC': 'Do you have a permanent voters card?',
11    'Vote During Elec. Vio.': 'Are you likely to vote when there is electoral violence around you?',
12    'Violence Deters Women From Parti.': 'Do you believe that violence deters women from participating in political activities such as rallies, elections, and campaigns?',
13    'Violence Impact Women Confi. In Parti.': 'In your opinion, does violence impact the confidence of women in engaging in political activities such as rallies, elections, and campaigns?',
14    'Vote in 2023 Gen. Elec.': 'Did you vote in 2023 General Elections?',
15    'If No, why not?': 'If No, why not?',
16    'Witnessed any Elec. Vio.': 'Have you ever witnessed any form of electoral violence during elections?',
17    'Witnessed Haras. Social Media': 'Have you ever witnessed any form of harassment on social media?'
18 }
19
20 # Match the names in X_train with the full column names from data
21 feature_names_full = [short_to_full_names[name] for name in X_train.columns]
22
23
24 # Plot feature importance
25 plt.figure(figsize=(10, 6))
26 plt.barh(feature_names_full, importances)
27 plt.title('Feature Importance')
28 plt.xlabel('Importance')
29 plt.ylabel('Feature')
30 plt.show()
31
```



Performance Analysis on test data

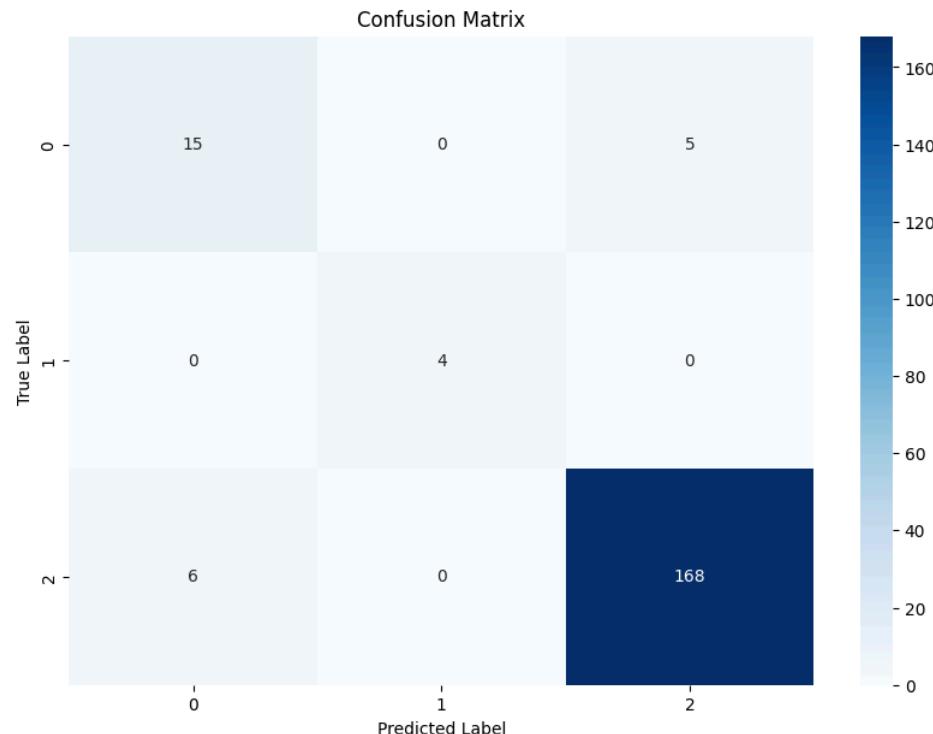
- RandomForest classification model for predicting participation of Female in Election process

```
In [98]: 1 # Predict on the test dataset
2 y_pred_test = rf_model.predict(X_test)
3
4 # Ensure predictions are integers (within the expected range)
5 y_pred_test = np.clip(y_pred_test, 0, 3) # Adjust this according to your binning
6
7 # Calculate evaluation metrics for the test set
8 test_accuracy = accuracy_score(y_test_binned, y_pred_test)
9 test_recall = recall_score(y_test_binned, y_pred_test, average='weighted')
10 test_precision = precision_score(y_test_binned, y_pred_test, average='weighted')
11 test_f1 = f1_score(y_test_binned, y_pred_test, average='weighted')
12
13 print("\nTest Performance Value:")
14 print(f"Random Forest Test Accuracy: {test_accuracy:.3f}")
15 print(f"Random Forest Test Recall: {test_recall:.3f}")
16 print(f"Random Forest Test Precision: {test_precision:.3f}")
17 print(f"Random Forest Test F1 Score: {test_f1:.3f}")
18
19 # Generate a classification report
20 report = classification_report(y_test, y_pred_test, target_names=['Class 0', 'Class 1', 'Class 2'], output_dict=False)
21
22 # Print the report
23 print("\n Random Forest Performance Report:")
24 print(report)
25
```

Test Performance Value:
 Random Forest Test Accuracy: 0.944
 Random Forest Test Recall: 0.944
 Random Forest Test Precision: 0.946
 Random Forest Test F1 Score: 0.945

	precision	recall	f1-score	support
Class 0	0.71	0.75	0.73	20
Class 1	1.00	1.00	1.00	4
Class 2	0.97	0.97	0.97	174
accuracy			0.94	198
macro avg	0.90	0.91	0.90	198
weighted avg	0.95	0.94	0.95	198

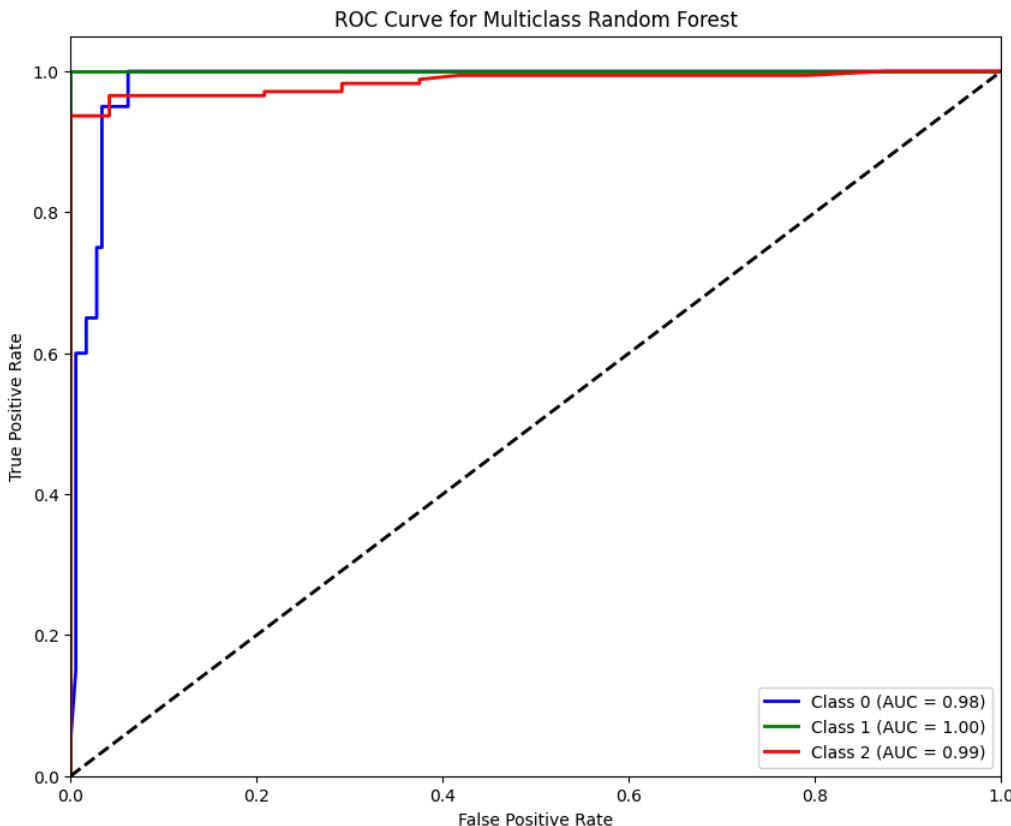
```
In [99]: 1 # Confusion matrix
2 cnf_matrix = confusion_matrix(y_test_binned, y_pred_test)
3
4 # Plot confusion matrix
5 plt.figure(figsize=(10, 7))
6 sns.heatmap(cnf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels=np.unique(y_test_binned), yticklabels=np.unique(y_test_binned))
7 plt.title("Confusion Matrix")
8 plt.xlabel("Predicted Label")
9 plt.ylabel("True Label")
10 plt.show()
11
```



```

In [100]: 1 # Convert Labels to binary format
2 y_test_bin = label_binarize(y_test, classes=[0, 1, 2]) # Adjust the classes accordingly (for 3 classes)
3 n_classes = y_test_bin.shape[1]
4
5
6 # Predict probabilities for each class
7 y_score = rf_model.predict_proba(X_test)
8
9
10 # Compute ROC curve and ROC area for each class
11 fpr = dict()
12 tpr = dict()
13 roc_auc = dict()
14
15 for i in range(n_classes):
16     fpr[i], tpr[i], _ = roc_curve(y_test_bin[:, i], y_score[:, i])
17     roc_auc[i] = auc(fpr[i], tpr[i])
18
19 # Plotting the ROC curves for each class
20 plt.figure(figsize=(10, 8))
21 colors = cycle(['blue', 'green', 'red'])
22
23 for i, color in zip(range(n_classes), colors):
24     plt.plot(fpr[i], tpr[i], color=color, lw=2,
25               label='Class {} (AUC = {:.2f})'.format(i, roc_auc[i]))
26
27 plt.plot([0, 1], [0, 1], 'k--', lw=2) # Diagonal Line
28 plt.xlim([0.0, 1.0])
29 plt.ylim([0.0, 1.05])
30 plt.xlabel('False Positive Rate')
31 plt.ylabel('True Positive Rate')
32 plt.title('ROC Curve for Multiclass Random Forest')
33 plt.legend(loc="lower right")
34 plt.show()
35

```



DESCRIPTION OF THE WORK SO FAR

- **STEP 1: EXPLORATORY DATA ANALYSIS (EDA)**
 - the EDA analysis provide proper understanding of the data
 - show existence of outliers
- **STEP 2: SPLITTING OF DATA**
 - Split data into Train and Test Set
- **STEP 3: APPLY MINMAX**
 - Apply MinMax Scaler on Train Set
 - Apply MinMax Scaler on Test Set
- **STEP 4: TRAIN MODEL**
 - Train the RandomForest Model on the Train Set
- **STEP 5: EVALUATE**
 - Evaluate the Model performance on the Test Set
 - metric include Accuracy, Recall, F1-score, Precision and ROC Curve
 - the model recorded the following performance
 1. **Accuracy of 94.4%**,
 2. **F1-score of 94.5%**,
 3. **Recall of 94.4%** and
 4. **Precision 94.6%**

