

2.1 Data Ingestion: Read the dataset. Do the descriptive statistics and do null value condition check, check for duplicates and outliers and write an inference on it. Perform Univariate and Bivariate Analysis and Multivariate Analysis.

The data set has a 1473 entries and 2 float, one integer and 7 object are present in our contraceptive.in Wife's age has null data.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1473 entries, 0 to 1472
Data columns (total 10 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Wife_age                             1402 non-null   float64
1   Wife_education                       1473 non-null   object
2   Husband_education                   1473 non-null   object
3   No_of_children_born                 1452 non-null   float64
4   Wife_religion                       1473 non-null   object
5   Wife_Working                        1473 non-null   object
6   Husband_Occupation                  1473 non-null   int64
7   Standard_of_living_index            1473 non-null   object
8   Media_exposure                      1473 non-null   object
9   Contraceptive_method_used           1473 non-null   object
dtypes: float64(2), int64(1), object(7)
memory usage: 115.2+ KB
```

Wife age and No.of children born features has 71 and 21 features respectively.

```
Wife_age                71
Wife_education           0
Husband_education        0
No_of_children_born     21
Wife_religion            0
Wife_Working            0
Husband_Occupation       0
Standard_of_living_index 0
Media_exposure           0
Contraceptive_method_used 0
dtype: int64
```

we can impute the features through mode and median and check the column whether the missing values presented are not.

```
Wife_age
25.0
No_of_children_born
2.0

Wife_age                0
Wife_education           0
Husband_education        0
No_of_children_born     0
Wife_religion            0
Wife_Working            0
Husband_Occupation       0
Standard_of_living_index 0
Media_exposure           0
Contraceptive_method_used 0
dtype: int64
```

Summary statistics:

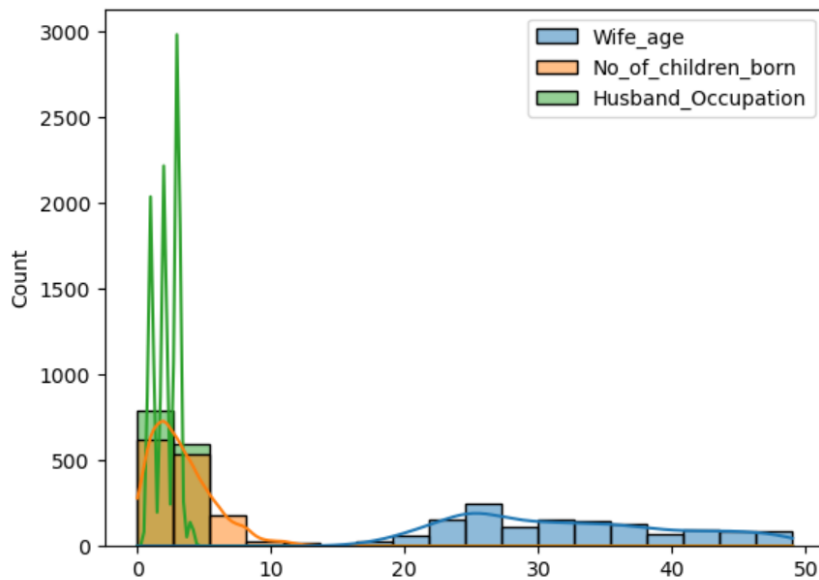
We can see the minimum and maximum age is 31 and 49. The other feature are ordinal categories.

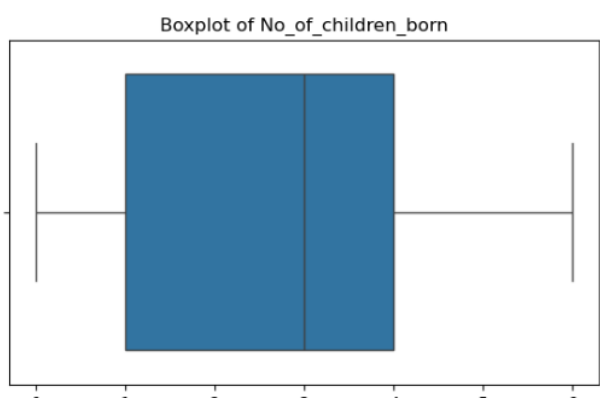
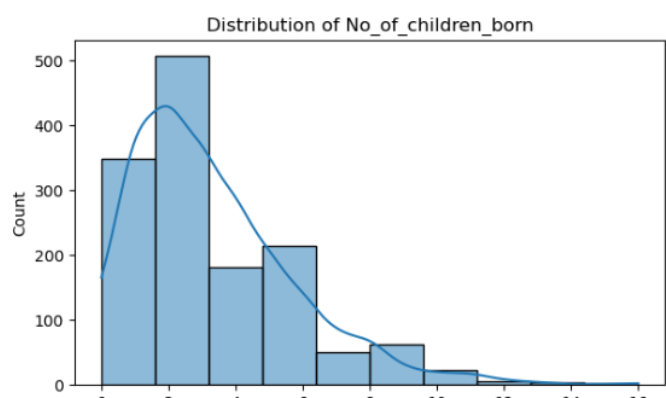
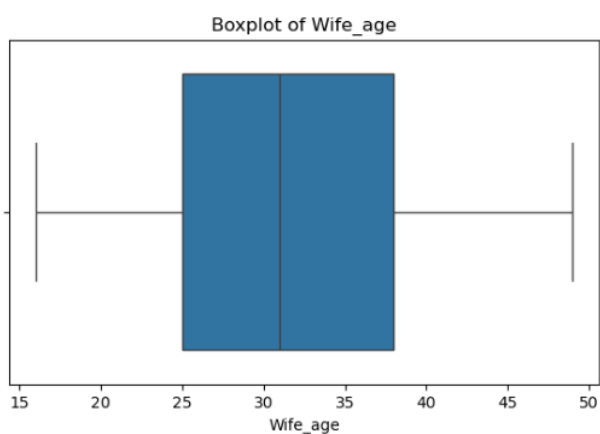
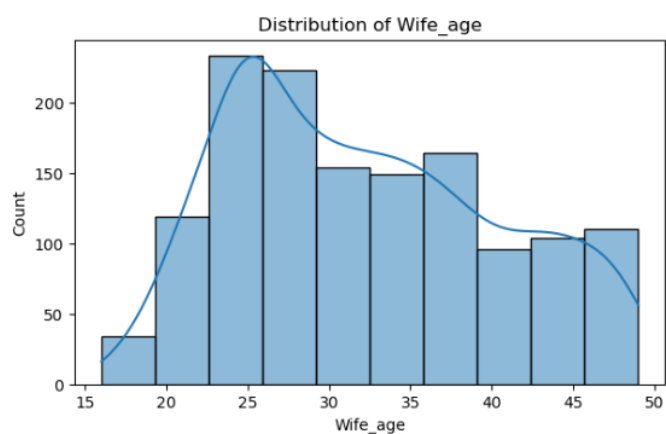
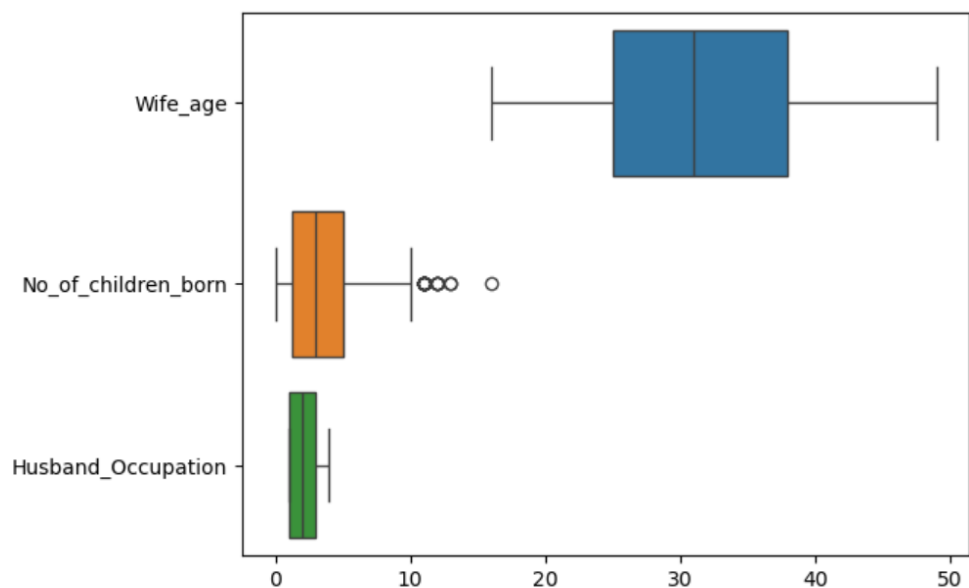
	count	mean	std	min	25%	50%	75%	max
Wife_age	1473.0	32.239647	8.235759	16.0	25.0	31.0	38.0	49.0
No_of_children_born	1473.0	3.236253	2.352985	0.0	1.0	3.0	4.0	16.0
Husband_Occupation	1473.0	2.137814	0.864857	1.0	1.0	2.0	3.0	4.0

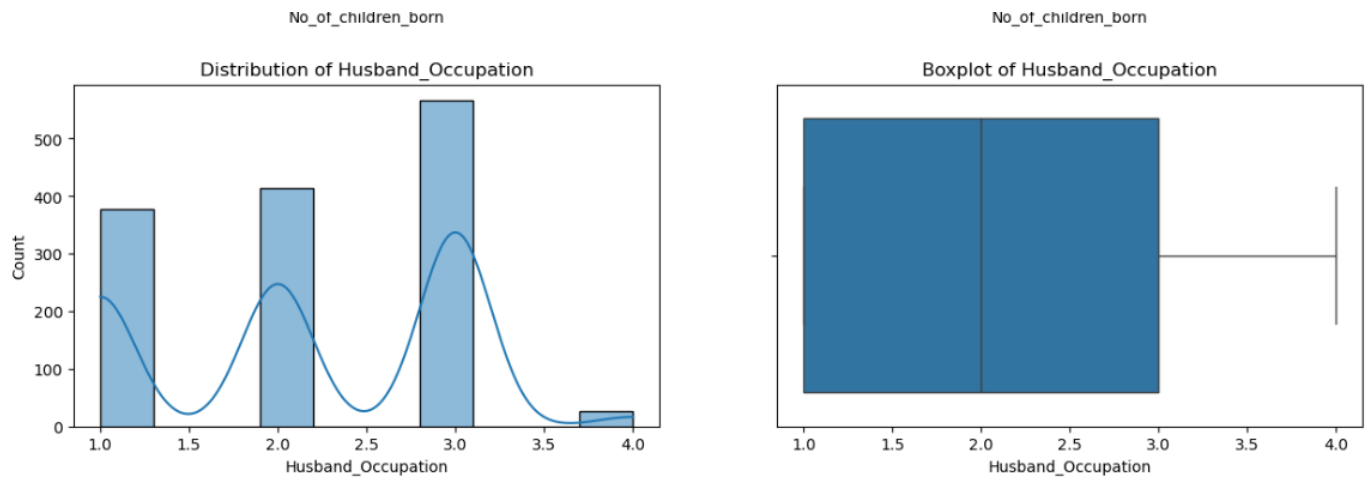
UNIVARIATE ANALYSIS:

Boxplot and Histogram Analysis:

The wife age, no.of children born and husband occupation column have numerical and right skewed distribution and the outliers are detected in children born features.



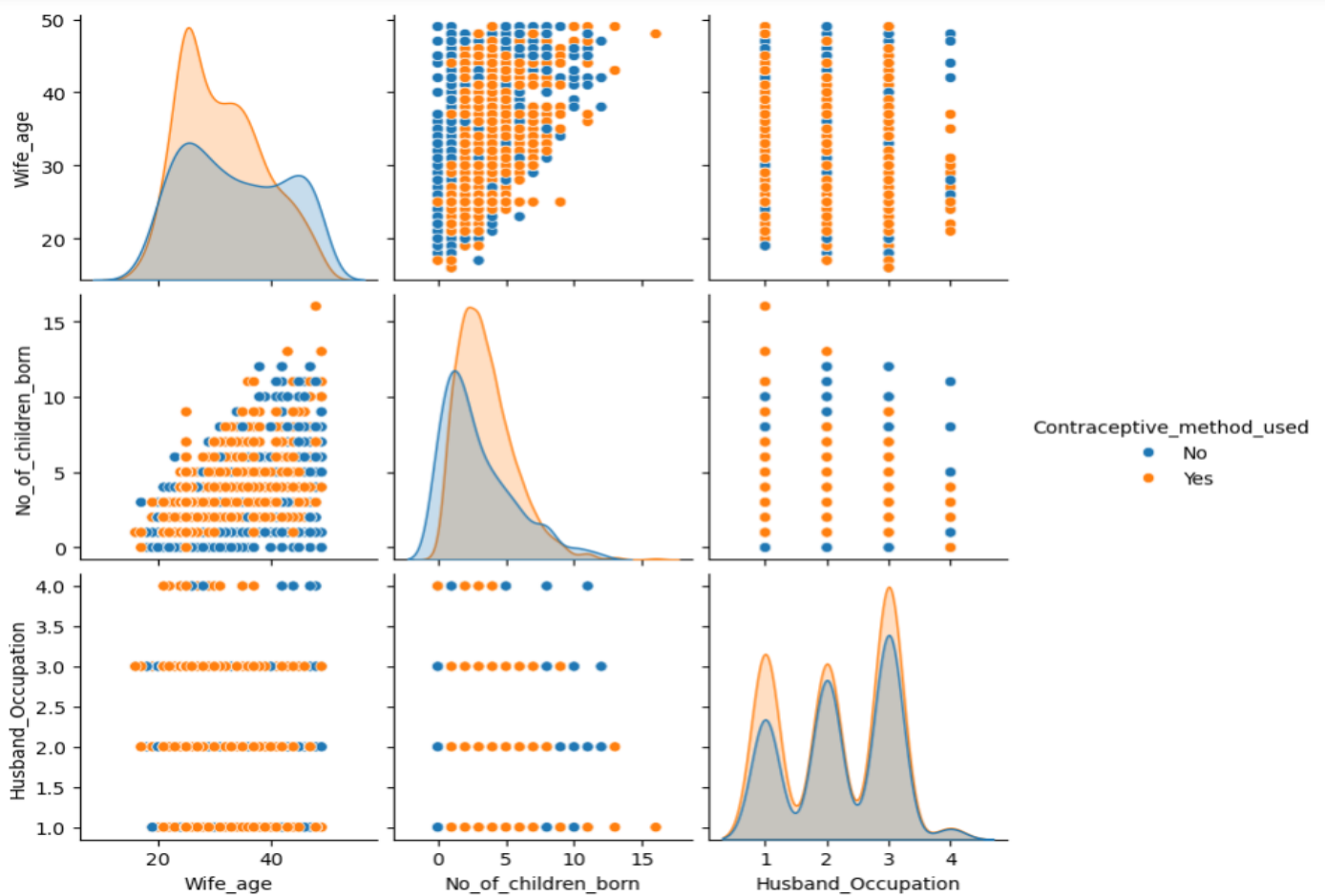




The feature of wife age has normally distributed and children born are right skewed it seems to be positive distribution

Scatter Plot:

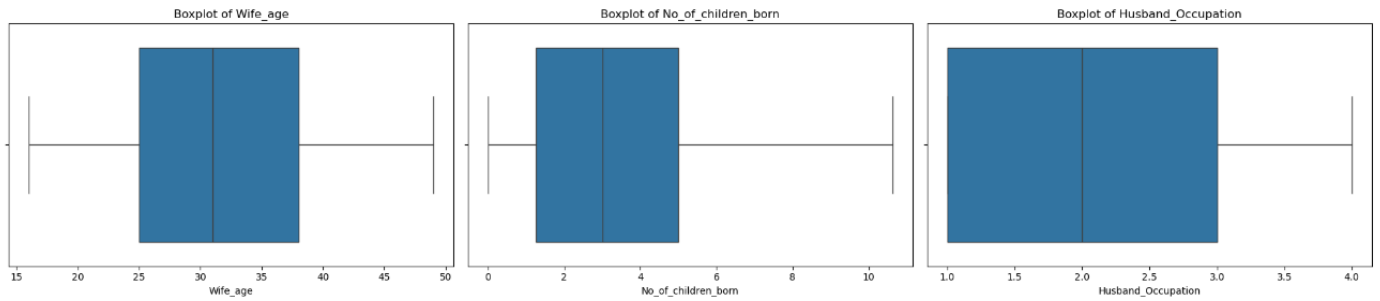
In this scatterplot, the more contraceptive used in wife_age VS no_of_children_born and husband_occupation vs wife_age and there is a less contraceptive are not used



After Removal of Outlier:

Proportion of outliers in Wife_age: 0.00%
Proportion of outliers in No_of_children_born: 0.00%
Proportion of outliers in Husband_Occupation: 0.00%

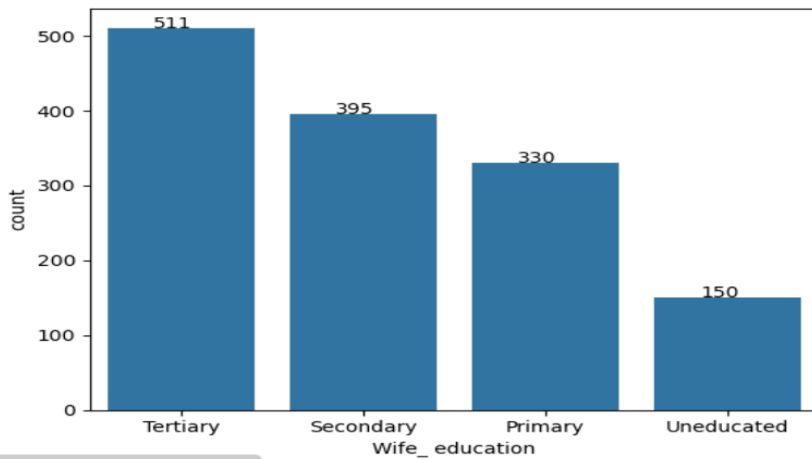
Outlier Analysis



Compare the Features through count plot:

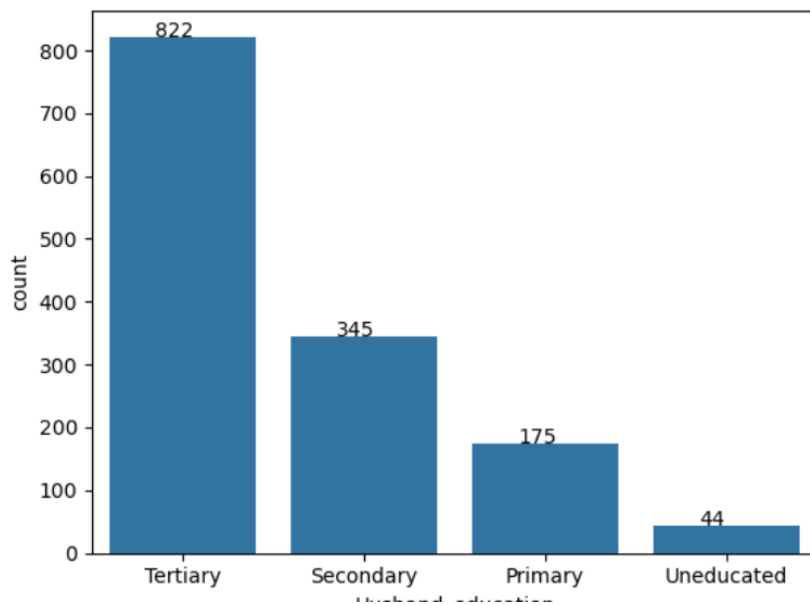
Most of the education are high because the uneducated column are low 150 and tertiary are high(511), in this data has women's are more educated.

Wife Education:



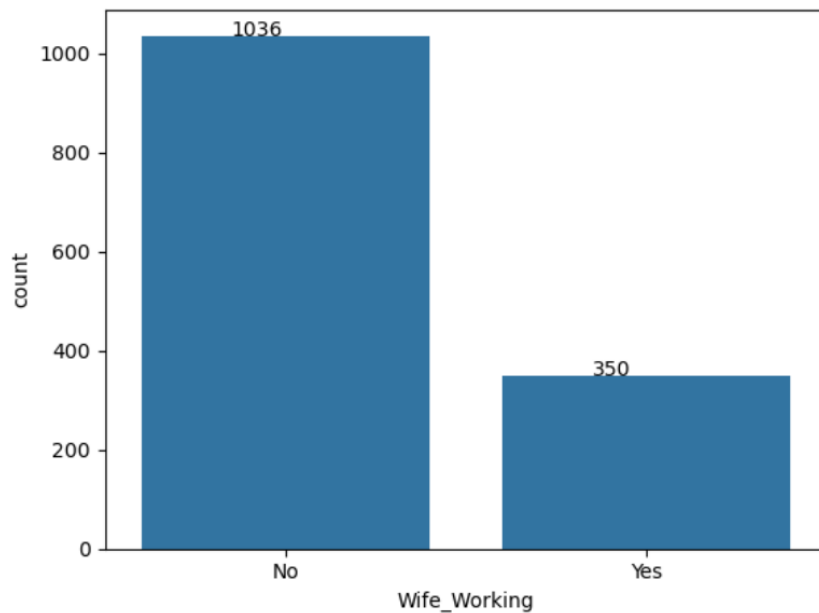
Husband Education:

Mostly are educated in husband education(Tertiary) and there is a low in uneducated terms.



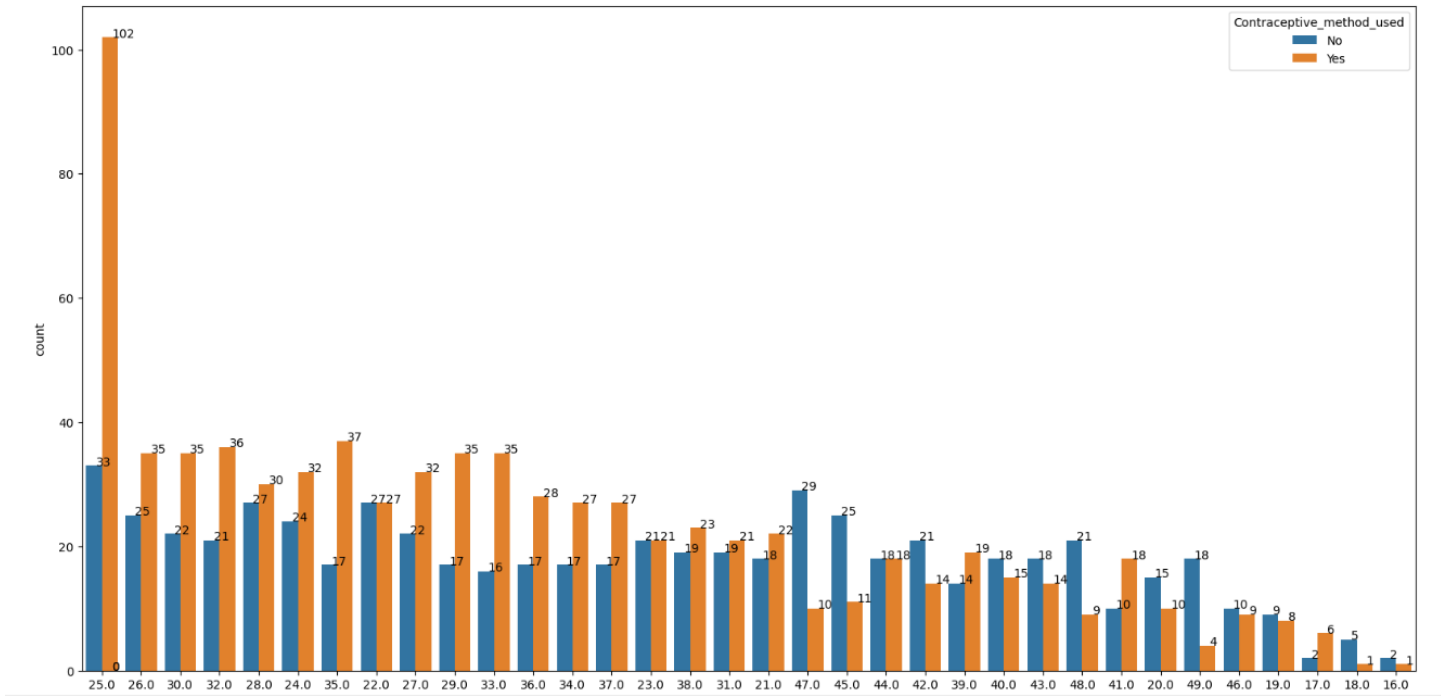
Wife Working:

The working terms are not educated is the highest factor in contraceptive data.

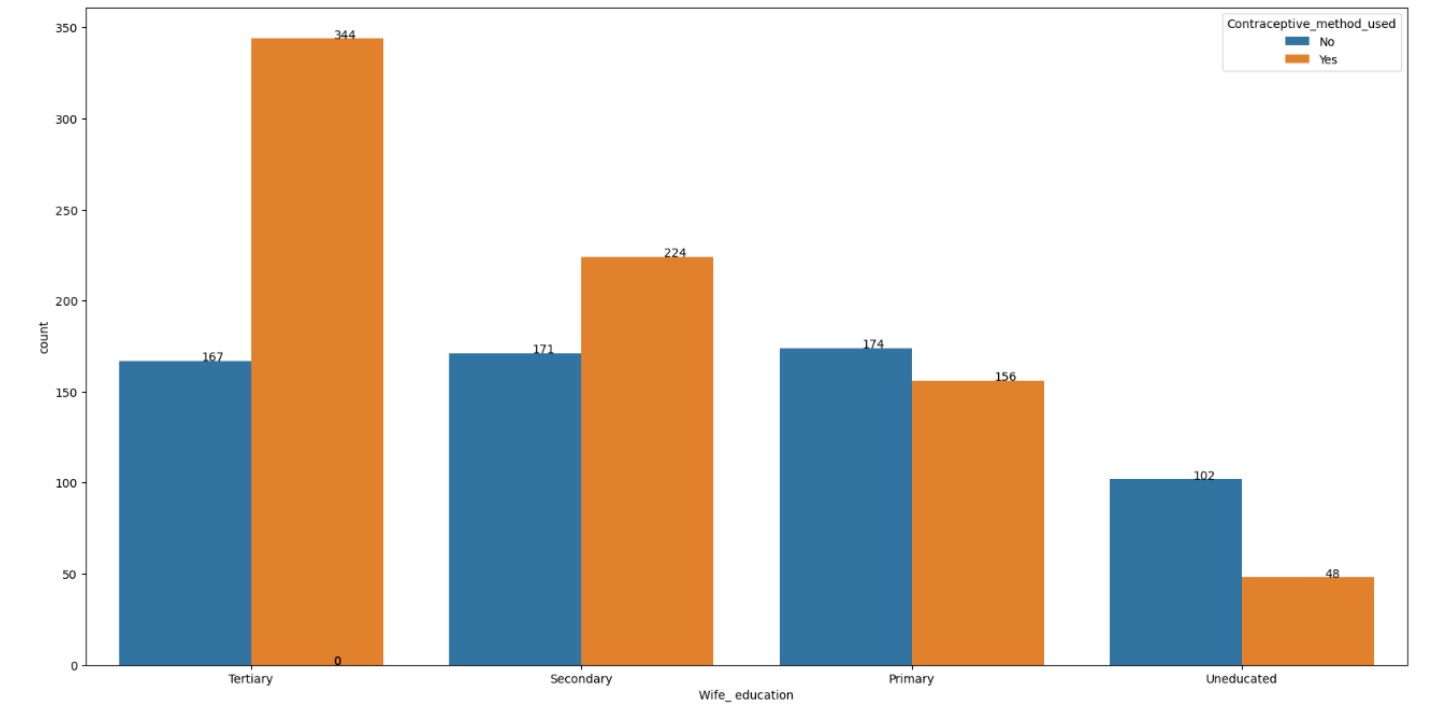


Describes the contraceptive prevalent through Countplot:

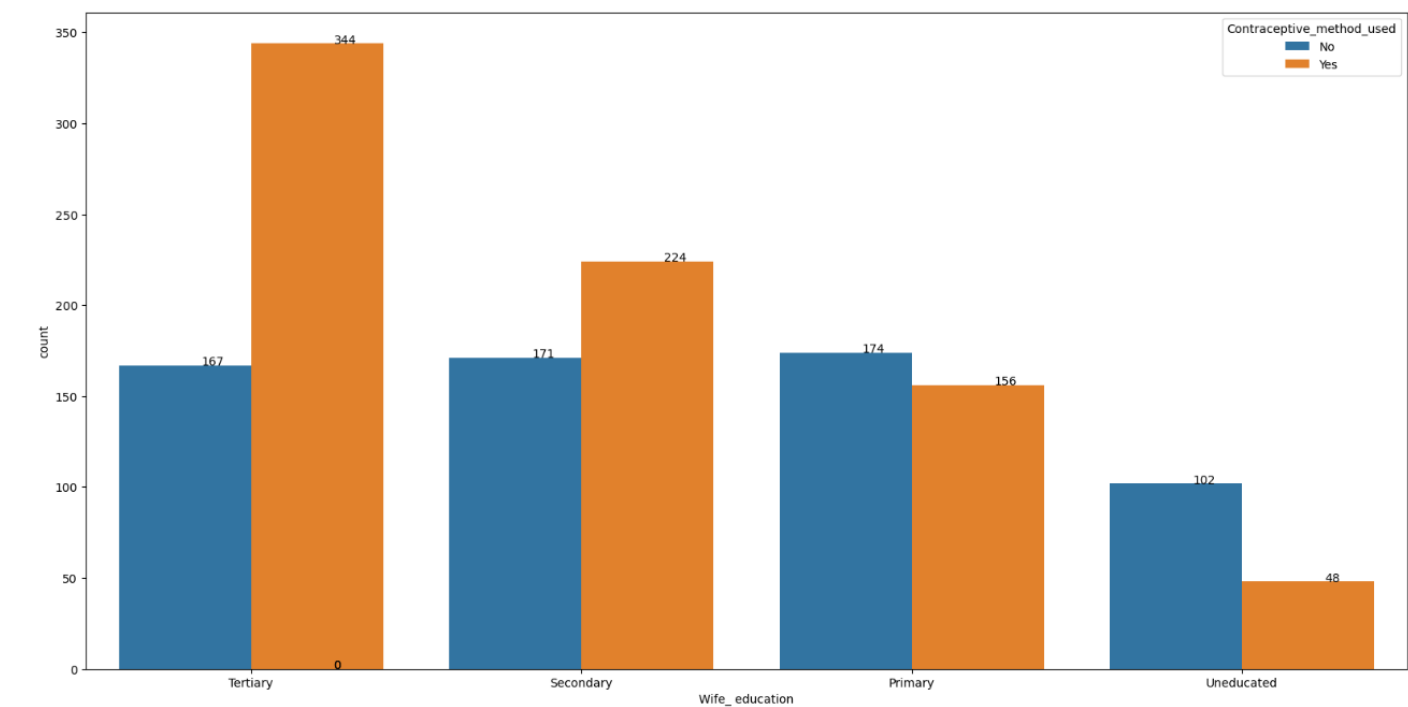
In this countplot, most the womens are contraceptive method in graphs and few charts shows not used the contraceptive method.



Wife Education:



Husband Occupation:



2.2 Do not scale the data. Encode the data (having string values) for Modelling. Data Split: Split the data into train and test (70:30). Apply Logistic Regression and LDA (linear discriminant analysis) and CART.

The dataset has imputed that the 0 referred as contraceptive used and 1 referred as contraceptive not used for encoding features and therefore highly used the contraceptive method.

	Wife_age	No_of_children_born	Husband_Occupation	Contraceptive_method_used	Wife_education_Secondary	Wife_education_Tertiary	Wife_education_Uneducated	Hu
0	24.0	3.0	2.0	1	0	0	0	
1	45.0	10.0	3.0	1	0	0	1	
2	43.0	7.0	3.0	1	0	0	0	
3	42.0	9.0	3.0	1	1	0	0	
4	36.0	8.0	3.0	1	1	0	0	
...	
96	41.0	1.0	2.0	1	0	0	0	
97	28.0	5.0	2.0	1	0	1	0	
98	17.0	3.0	3.0	1	0	0	0	
99	27.0	5.0	2.0	1	0	0	0	
100	27.0	4.0	3.0	1	0	0	1	

100 rows × 16 columns

```
0    772
1    614
Name: Contraceptive_method_used, dtype: int64
```


LDA :

LDA is used to find a linear combination of features that separates two or more classes object or events. It's also used in dimensionality reduction. It maximizes the components axis for class separation and it assumes the normal distribution of the input variable. The different transformation can be applied to the data to make it normal (from positive or negative to normal distribution).

Train Model Score:

The train and test score has good performance and it indicates the model is not over fitting and under fitting

68.1%

Test Model Score:

67.1%

Prediction Test:

```
array([0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0,
       0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1,
       1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1,
       1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0,
       0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
       1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
       1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1,
       1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0,
       0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1,
       0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
       0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
       1, 1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1,
       0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0,
       0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0,
       0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1,
       0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1,
       0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0,
       0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1,
       0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0])
```

Logistic regression:

The Logistic regression are used to predict the categorical independent class. The importance is to calculate the features importance by taking average of the absolute values of the coefficients across all classes.

Train Model Score:

In this model has a good performance presented in this model. Same as LDA

68.55%

Test Model Score:

67.54%

```
array([[0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0,
       0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1,
       1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1,
       1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0,
       0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0,
       1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
       1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1,
       0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0,
       0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1,
       0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
       0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
       1, 1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1,
       0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0,
       0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0,
       0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1,
       0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1,
       0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0,
       0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1,
       0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0])
```

2.3 Performance Metrics: Check the performance of Predictions on Train and Test sets using Accuracy, Confusion Matrix, Plot ROC curve and get ROC_AUC score for each model Final Model: Compare Both the models and write inference which model is best/optimized.

Linear Discriminant Analysis:

Confusion_Matrix:

190 true positive and 90 true negative – These are all correct predictions.

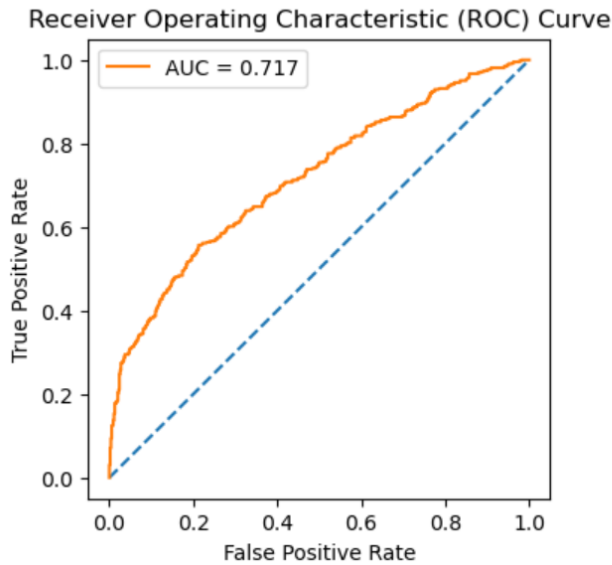
100 and 36 are false positive and False negative – These are all incorrect predictions

```
[[190  36]
 [100  90]]
```

AUC – ROC Curve:

This model has fair classification between True positive and false positive(71%)

AUC: 0.717



Classification Report:

In this report, the model **84% correctly** predict the womens are **contraceptive used** and the f1 score (74%) is weighted balance between precision and recall

	precision	recall	f1-score	support
0	0.66	0.84	0.74	226
1	0.71	0.47	0.57	190
accuracy			0.67	416
macro avg	0.68	0.66	0.65	416
weighted avg	0.68	0.67	0.66	416

Logistic regression:

[illegible]

Confusion Matrix:

190 and 90 are correct prediction and 100 and 36 are incorrect predictions.

```
[[190  36]
 [100  90]]
```

Classification Report:

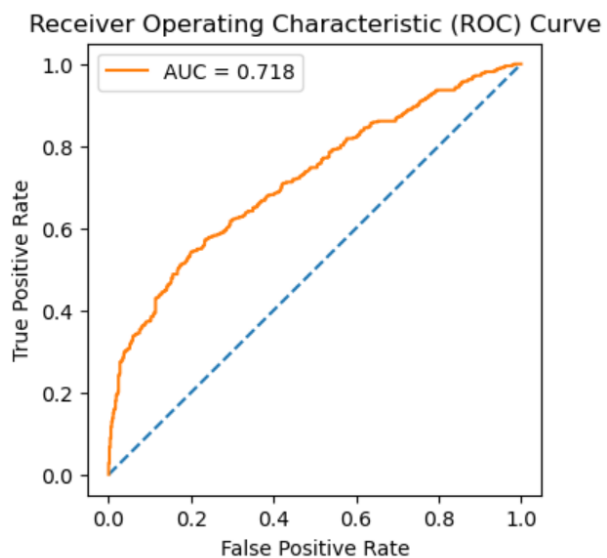
The Logistic model 84% correctly predict as a contraceptive are used. The F1 score also proved the 74% weighted balance between precision and recall

	precision	recall	f1-score	support
0	0.66	0.84	0.74	226
1	0.72	0.48	0.57	190
accuracy			0.68	416
macro avg	0.69	0.66	0.66	416
weighted avg	0.68	0.68	0.66	416

ROC-AUC Curve:

The ROC-AUC Curve determines the fair classification between positive and negative rate 71 percentage.

AUC: 0.718



2.4 Inference: Basis on these predictions, what are the insights and recommendations.

The logistic and Lda suggest that the most women are used contraceptive method and it can mainly prevent from pregnancy. But it also reduce a risk of ovarian cancer and protect against acute pelvic inflammatory disease. It increases the risk of cardiovascular disease.

The 84% predicted as a women are used the contraceptive method. so the women were not pregnant.