

Logistic Regression

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 %matplotlib inline
```

Types of Variable

```
1 cercan_df.info()

<class 'pandas.core.frame.DataFrame'>
Index: 835 entries, 0 to 857
Data columns (total 36 columns):
 #   Column                                     Non-Null Count  Dtype
---  -
 0   Age                                       835 non-null    int64
 1   Number_of_sexual_partners               835 non-null    float64
 2   First_sexual_intercourse                 835 non-null    float64
 3   Num_of_pregnancies                       835 non-null    float64
 4   Smokes                                   835 non-null    float64
 5   Smokes_(years)                          835 non-null    float64
 6   Smokes_(packs/year)                     835 non-null    float64
 7   Hormonal_Contraceptives                 835 non-null    float64
 8   Hormonal_Contraceptives_(years)         835 non-null    float64
 9   IUD                                       835 non-null    float64
10   IUD_(years)                             835 non-null    float64
11   STDs                                     835 non-null    float64
12   STDs_(number)                           835 non-null    float64
13   STDs:condylomatosis                     835 non-null    float64
14   STDs:cervical_condylomatosis            835 non-null    float64
15   STDs:vaginal_condylomatosis             835 non-null    float64
16   STDs:vulvo-perineal_condylomatosis      835 non-null    float64
17   STDs:syphilis                           835 non-null    float64
18   STDs:pelvic_inflammatory_disease        835 non-null    float64
19   STDs:genital_herpes                     835 non-null    float64
20   STDs:molluscum_contagiosum              835 non-null    float64
21   STDs:AIDS                               835 non-null    float64
22   STDs:HIV                                 835 non-null    float64
23   STDs:Hepatitis_B                        835 non-null    float64
24   STDs:HPV                                835 non-null    float64
25   STDs:_Number_of_diagnosis                835 non-null    int64
26   STDs:_Time_since_first_diagnosis         835 non-null    float64
27   STDs:_Time_since_last_diagnosis          835 non-null    float64
28   Dx:Cancer                               835 non-null    int64
29   Dx:CIN                                   835 non-null    int64
30   Dx:HPV                                   835 non-null    int64
31   Dx                                        835 non-null    int64
32   Hinselmann                              835 non-null    int64
33   Schiller                                 835 non-null    int64
34   Cytology                                835 non-null    int64
35   Biopsy                                   835 non-null    int64
dtypes: float64(26), int64(10)
memory usage: 241.4 KB
```

```
1 # find the categorical variables
2 cat = [i for i in cercan_df.columns if cercan_df[i].dtype=='O']
3 print(f'There are {len(cat)} categorical variables\n')
4 print('The categorical variables are: ', cat)
```

There are 0 categorical variables

The categorical variables are: []

Explore Numerical Variables

```
1 # find the numerical variables
2 num = [i for i in cercan_df.columns if cercan_df[i].dtype != 'O']
3 print(f'There are {len(num)} numerical variables\n')
4 print('The numerical variables are:', num)
```

There are 36 numerical variables

The numerical variables are: ['Age', 'Number_of_sexual_partners', 'First_sexual_intercourse', 'Num_of_pregnancies', 'Smokes', 'Smokes_(years)', 'Smokes_(packs/year)', 'Hormonal_Contraceptives', 'Hormonal_Contraceptives_(years)', 'IUD', 'IUD_(years)', 'STDs', 'STDs_(number)', 'STDs:condylomatosis', 'STDs:cervical_condylomatosis', 'STDs:vaginal_condylomatosis', 'STDs:vulvo-perineal_condylomatosis', 'STDs:syphilis', 'STDs:pelvic_inflammatory_disease', 'STDs:genital_herpes', 'STDs:molluscum_contagiosum', 'STDs:AIDS', 'STDs:HIV', 'STDs:Hepatitis_B', 'STDs:HPV', 'STDs:_Number_of_diagnosis', 'STDs:_Time_since_first_diagnosis', 'STDs:_Time_since_last_diagnosis', 'Dx:Cancer', 'Dx:CIN', 'Dx:HPV', 'Dx', 'Hinselmann', 'Schiller', 'Cytology', 'Biopsy']

```
1 # checking if the length of the numerical is the same length as the dataframe
2 col = [i for i in cercan_df.columns]
3 len(col)

36
```

Outliers in numerical values

```
1 # creating a function to check all of the columns that has outliers
2 def outliers(data):
```

```

3  outliers=[]
4  for i in data.columns:
5      low_bound = data[i].quantile(0.25) - (1.5 * (data[i].quantile(0.75) - data[i].quantile(0.25)))
6      upper_bound = data[i].quantile(0.75) + (1.5 * (data[i].quantile(0.75) - data[i].quantile(0.25)))
7      if ((data[i] < low_bound) | (data[i] > upper_bound)).any():
8          outliers.append(i)
9  return outliers

```

```

1 # checking all the columns in the dataframe
2 outliers(cercan_df)

```

```

['Age',
 'Number_of_sexual_partners',
 'First_sexual_intercourse',
 'Num_of_pregnancies',
 'Smokes',
 'Smokes_(years)',
 'Smokes_(packs/year)',
 'Hormonal_Contraceptives_(years)',
 'IUD',
 'IUD_(years)',
 'STDs',
 'STDs_(number)',
 'STDs:condylomatosis',
 'STDs:vaginal_condylomatosis',
 'STDs:vulvo-perineal_condylomatosis',
 'STDs:syphilis',
 'STDs:pelvic_inflammatory_disease',
 'STDs:genital_herpes',
 'STDs:molluscum_contagiosum',
 'STDs:HIV',
 'STDs:Hepatitis_B',
 'STDs:HPV',
 'STDs:_Number_of_diagnosis',
 'STDs:_Time_since_first_diagnosis',
 'STDs:_Time_since_last_diagnosis',
 'Dx:Cancer',
 'Dx:CIN',
 'Dx:HPV',
 'Dx',
 'Hinselmann',
 'Schiller',
 'Cytology',
 'Biopsy']

```

```

1 # checking to see which ones are booleans as to not include in the boxplot
2 cercan_df.describe(), 2

```

min	13.000000	1.000000	10.000000
25%	21.000000	2.000000	15.000000
50%	26.000000	2.000000	17.000000
75%	32.000000	3.000000	18.000000
max	84.000000	28.000000	32.000000

	Num_of_pregnancies	Smokes	Smokes_(years)	Smokes_(packs/year)	\
count	835.000000	835.000000	835.000000	835.000000	
mean	2.283832	0.147305	1.234329	0.458571	
std	1.408152	0.354623	4.111264	2.239363	
min	0.000000	0.000000	0.000000	0.000000	
25%	1.000000	0.000000	0.000000	0.000000	
50%	2.000000	0.000000	0.000000	0.000000	
75%	3.000000	0.000000	0.000000	0.000000	
max	11.000000	1.000000	37.000000	37.000000	

	Hormonal_Contraceptives	Hormonal_Contraceptives_(years)	IUD	\
count	835.000000	835.000000	835.000000	
mean	0.651639	2.080520	0.099401	
std	0.446366	3.601364	0.299379	
min	0.000000	0.000000	0.000000	
25%	0.000000	0.000000	0.000000	
50%	1.000000	0.500000	0.000000	
75%	1.000000	3.000000	0.000000	
max	1.000000	30.000000	1.000000	

	... STDs:_Time_since_first_diagnosis	STDs:_Time_since_last_diagnosis	\
count	...	835.000000	835.000000
mean	...	4.182036	3.239521
std	...	1.809358	1.843420
min	...	1.000000	1.000000
25%	...	4.000000	3.000000
50%	...	4.000000	3.000000
75%	...	4.000000	3.000000
max	...	22.000000	22.000000

std	0.145519	0.105520	0.145519	0.107182	0.200518	0.288020
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

	Citology	Biopsy
count	835.000000	835.000000
mean	0.051497	0.064671
std	0.221142	0.246091
min	0.000000	0.000000
25%	0.000000	0.000000
50%	0.000000	0.000000
75%	0.000000	0.000000
max	1.000000	1.000000

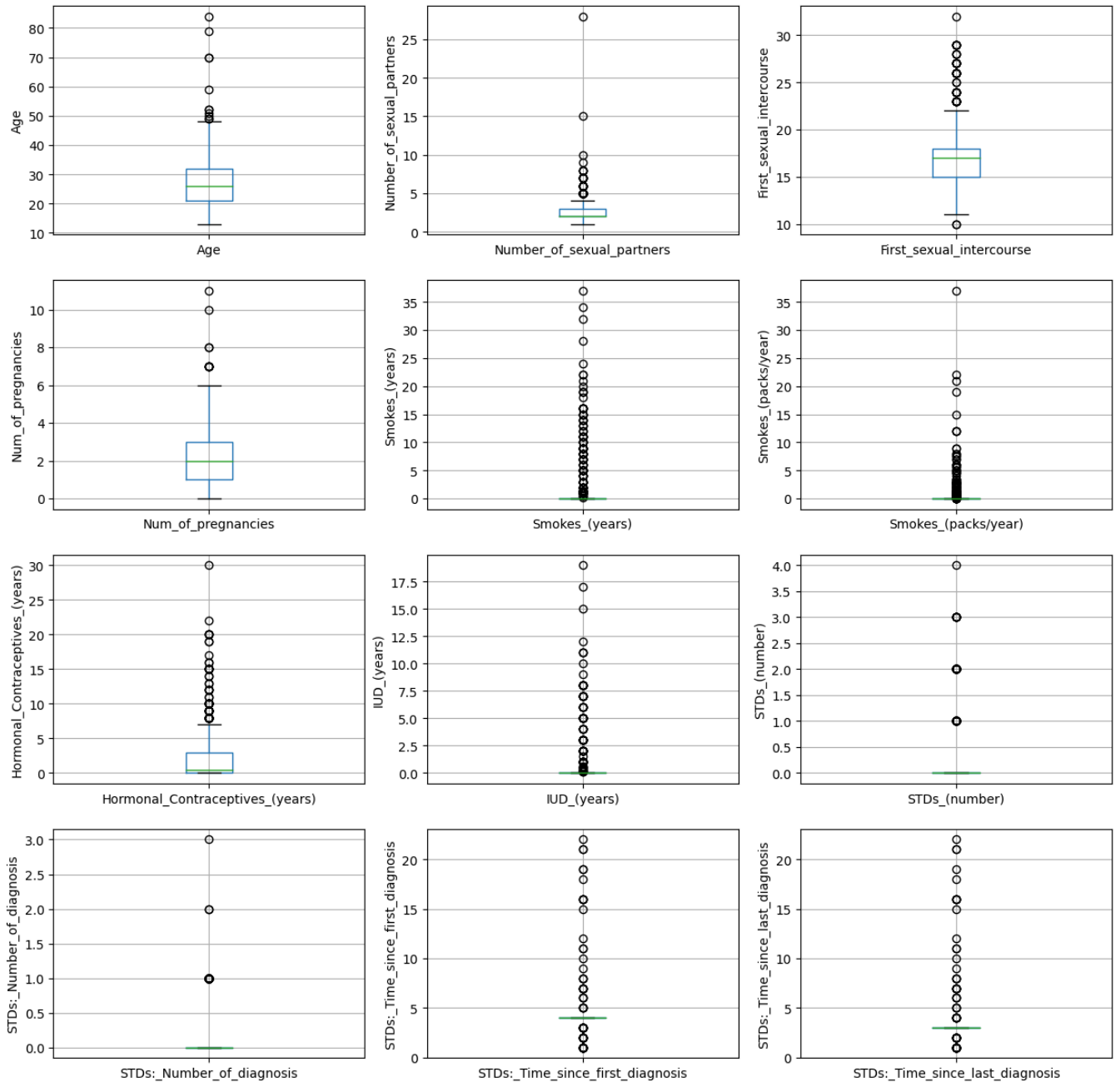
[8 rows x 36 columns],
2)

```
1 cercan_df['Smokes_(years)'].value_counts()
```

Smokes_(years)	
0.000000	712
1.266973	15
5.000000	9
9.000000	9
1.000000	8
3.000000	7
2.000000	7
16.000000	6
7.000000	6
8.000000	6
11.000000	5
4.000000	5
10.000000	5
14.000000	4
15.000000	4
6.000000	4
13.000000	3
0.500000	3
19.000000	3
12.000000	3
22.000000	2
32.000000	1
20.000000	1
28.000000	1
24.000000	1
18.000000	1
34.000000	1
21.000000	1
37.000000	1
0.160000	1
Name: count, dtype: int64	

```
1 # creating a boxplot for each column that is not boolean and has an outlier
2 plt.figure(figsize=(15,15))
3
4 plt.subplot(4, 3, 1)
5 fig = cercan_df.boxplot(column='Age')
6 fig.set_title('')
7 fig.set_ylabel('Age')
8
9 plt.subplot(4, 3, 2)
10 fig = cercan_df.boxplot(column='Number_of_sexual_partners')
11 fig.set_title('')
12 fig.set_ylabel('Number_of_sexual_partners')
13
14 plt.subplot(4, 3, 3)
15 fig = cercan_df.boxplot(column='First_sexual_intercourse')
16 fig.set_title('')
17 fig.set_ylabel('First_sexual_intercourse')
18
19 plt.subplot(4, 3, 4)
20 fig = cercan_df.boxplot(column='Num_of_pregnancies')
21 fig.set_title('')
22 fig.set_ylabel('Num_of_pregnancies')
23
24 plt.subplot(4, 3, 5)
25 fig = cercan_df.boxplot(column='Smokes_(years)')
26 fig.set_title('')
27 fig.set_ylabel('Smokes_(years)')
28
29 plt.subplot(4, 3, 6)
30 fig = cercan_df.boxplot(column='Smokes_(packs/year)')
31 fig.set_title('')
32 fig.set_ylabel('Smokes_(packs/year)')
33
34 plt.subplot(4, 3, 7)
35 fig = cercan_df.boxplot(column='Hormonal_Contraceptives_(years)')
36 fig.set_title('')
37 fig.set_ylabel('Hormonal_Contraceptives_(years)')
38
39 plt.subplot(4, 3, 8)
40 fig = cercan_df.boxplot(column='IUD_(years)')
41 fig.set_title('')
42 fig.set_ylabel('IUD_(years)')
43
44 plt.subplot(4, 3, 9)
45 fig = cercan_df.boxplot(column='STDs_(number)')
46 fig.set_title('')
47 fig.set_ylabel('STDs_(number)')
48
49 plt.subplot(4, 3, 10)
50 fig = cercan_df.boxplot(column='STDs:_Number_of_diagnosis')
51 fig.set_title('')
52 fig.set_ylabel('STDs:_Number_of_diagnosis')
53
54 plt.subplot(4, 3, 11)
55 fig = cercan_df.boxplot(column='STDs:_Time_since_first_diagnosis')
56 fig.set_title('')
57 fig.set_ylabel('STDs:_Time_since_first_diagnosis')
58
59 plt.subplot(4, 3, 12)
60 fig = cercan_df.boxplot(column='STDs:_Time_since_last_diagnosis')
61 fig.set_title('')
62 fig.set_ylabel('STDs:_Time_since_last_diagnosis')
```

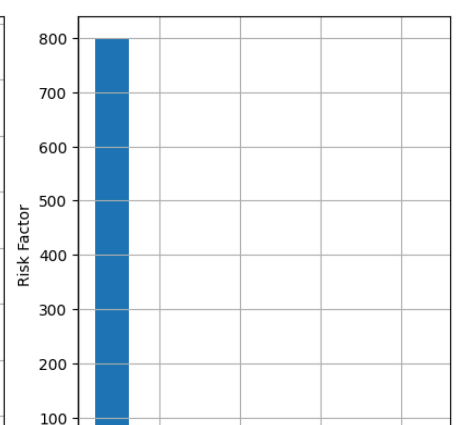
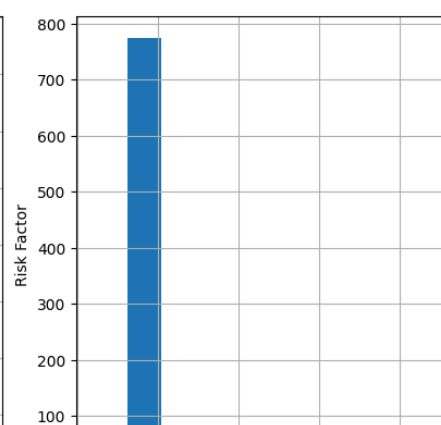
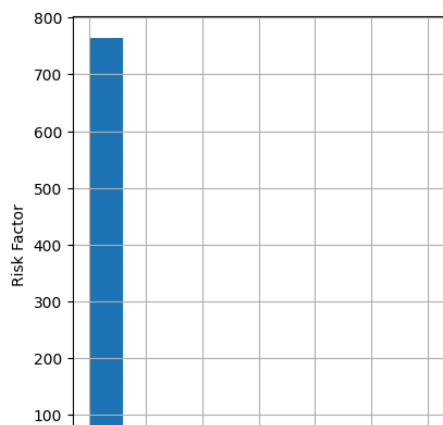
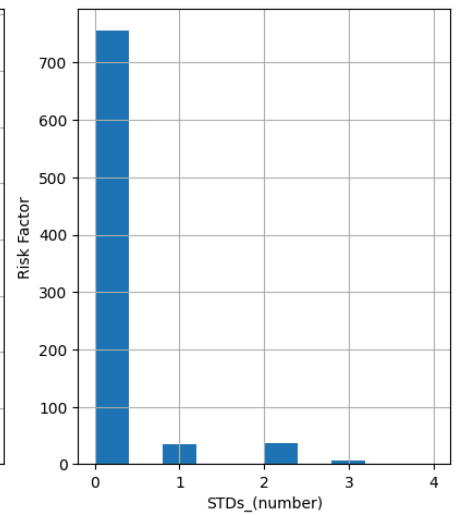
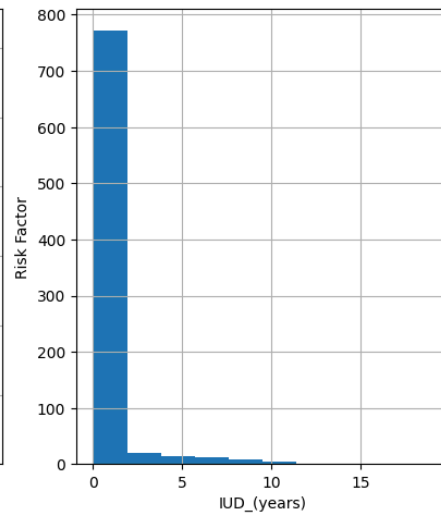
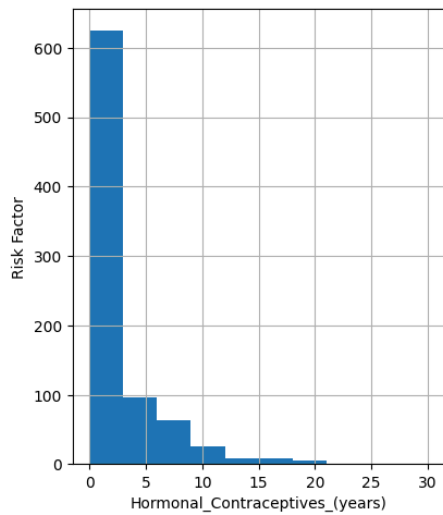
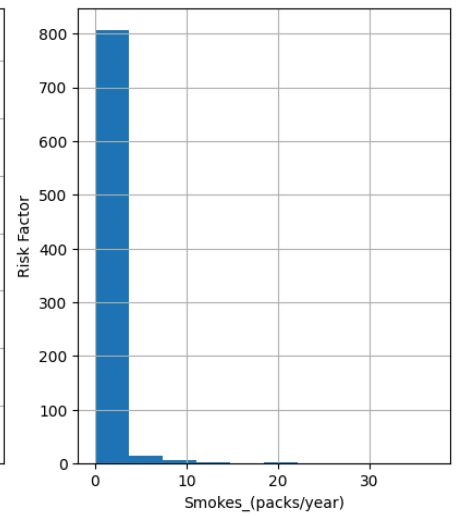
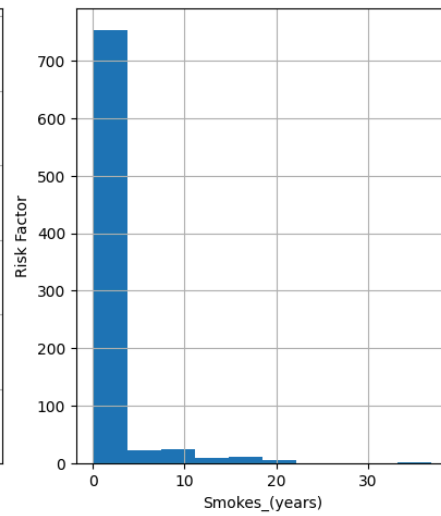
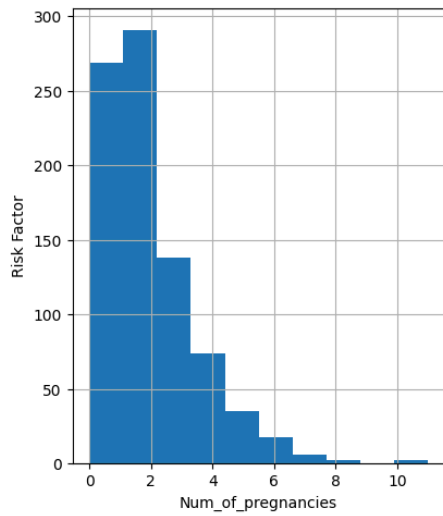
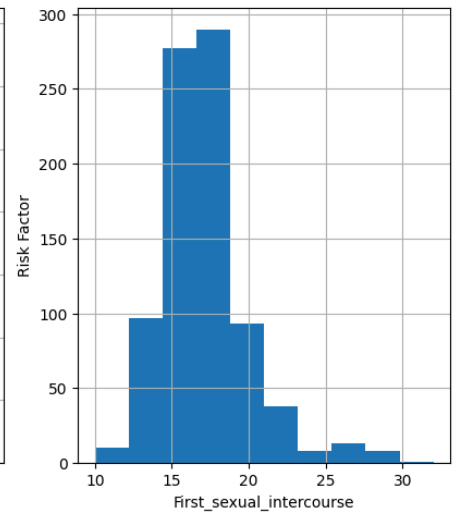
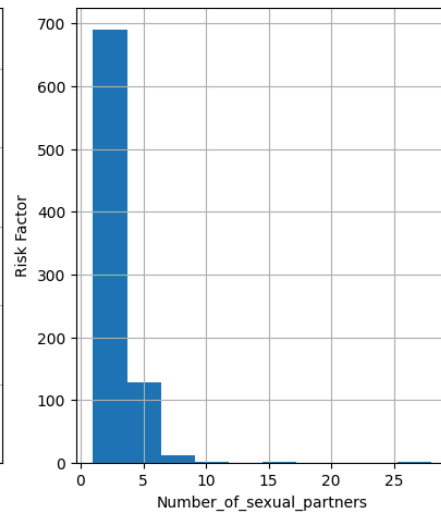
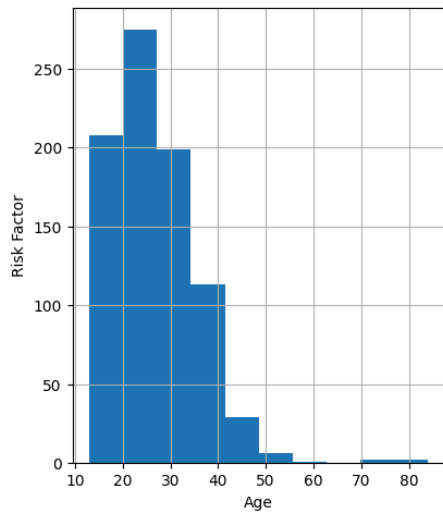
```
Text(0, 0.5, 'STDs:_Time_since_last_diagnosis')
```

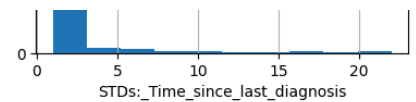
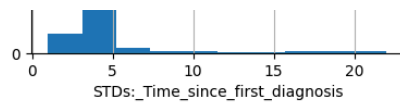
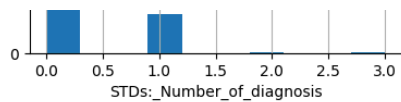


✓ Check the distribution of variables

```
1 # creating a boxplot for each column that is not boolean and has an outlier
2 # plotting the histogram to check the outliers
3 plt.figure(figsize=(15,25))
4
5 plt.subplot(4, 3, 1)
6 fig = cercan_df['Age'].hist(bins=10)
7 fig.set_xlabel('Age')
8 fig.set_ylabel('Risk Factor')
9
10 plt.subplot(4, 3, 2)
11 fig = cercan_df['Number_of_sexual_partners'].hist(bins=10)
12 fig.set_xlabel('Number_of_sexual_partners')
13 fig.set_ylabel('Risk Factor')
14
15 plt.subplot(4, 3, 3)
16 fig = cercan_df['First_sexual_intercourse'].hist(bins=10)
17 fig.set_xlabel('First_sexual_intercourse')
18 fig.set_ylabel('Risk Factor')
19
20 plt.subplot(4, 3, 4)
21 fig = cercan_df['Num_of_pregnancies'].hist(bins=10)
22 fig.set_xlabel('Num_of_pregnancies')
23 fig.set_ylabel('Risk Factor')
24
25
26 plt.subplot(4, 3, 5)
27 fig = cercan_df['Smokes_(years)'].hist(bins=10)
28 fig.set_xlabel('Smokes_(years)')
29 fig.set_ylabel('Risk Factor')
30
31 plt.subplot(4, 3, 6)
32 fig = cercan_df['Smokes_(packs/year)'].hist(bins=10)
33 fig.set_xlabel('Smokes_(packs/year)')
34 fig.set_ylabel('Risk Factor')
35
36 plt.subplot(4, 3, 7)
37 fig = cercan_df['Hormonal_Contraceptives_(years)'].hist(bins=10)
38 fig.set_xlabel('Hormonal_Contraceptives_(years)')
39 fig.set_ylabel('Risk Factor')
40
41 plt.subplot(4, 3, 8)
42 fig = cercan_df['IUD_(years)'].hist(bins=10)
43 fig.set_xlabel('IUD_(years)')
44 fig.set_ylabel('Risk Factor')
45
46 plt.subplot(4, 3, 9)
47 fig = cercan_df['STDs_(number)'].hist(bins=10)
48 fig.set_xlabel('STDs_(number)')
49 fig.set_ylabel('Risk Factor')
50
51 plt.subplot(4, 3, 10)
52 fig = cercan_df['STDs:_Number_of_diagnosis'].hist(bins=10)
53 fig.set_xlabel('STDs:_Number_of_diagnosis')
54 fig.set_ylabel('Risk Factor')
55
56 plt.subplot(4, 3, 11)
57 fig = cercan_df['STDs:_Time_since_first_diagnosis'].hist(bins=10)
58 fig.set_xlabel('STDs:_Time_since_first_diagnosis')
59 fig.set_ylabel('Risk Factor')
60
61 plt.subplot(4, 3, 12)
62 fig = cercan_df['STDs:_Time_since_last_diagnosis'].hist(bins=10)
63 fig.set_xlabel('STDs:_Time_since_last_diagnosis')
64 fig.set_ylabel('Risk Factor')
```

```
Text(0, 0.5, 'Risk Factor')
```





Since all outliers are skewed, interquartile range would be used to find the outliers.

```
1 low_bound = cercan_df['Age'].quantile(0.25) - (1.5 * (cercan_df['Age'].quantile(0.75) - cercan_df['Age'].quantile(0.25)))
2 upper_bound = cercan_df['Age'].quantile(0.75) + (1.5 * (cercan_df['Age'].quantile(0.75) - cercan_df['Age'].quantile(0.25)))
3 print(f'Age outliers are value < {low_bound} or > {upper_bound}')
```

Age outliers are value < 4.5 or > 48.5

```
1 low_bound = cercan_df['Number_of_sexual_partners'].quantile(0.25) - (1.5 * (cercan_df['Number_of_sexual_partners'].quantile(0.75) - cercan_df['Number_of_sexual_partners'].quantile(0.25)))
2 upper_bound = cercan_df['Number_of_sexual_partners'].quantile(0.75) + (1.5 * (cercan_df['Number_of_sexual_partners'].quantile(0.75) - cercan_df['Number_of_sexual_partners'].quantile(0.25)))
3 print(f'Number_of_sexual_partners outliers are value < {low_bound} or > {upper_bound}')
```

Number_of_sexual_partners outliers are value < 0.5 or > 4.5

```
1 low_bound = cercan_df['First_sexual_intercourse'].quantile(0.25) - (1.5 * (cercan_df['First_sexual_intercourse'].quantile(0.75) - cercan_df['First_sexual_intercourse'].quantile(0.25)))
2 upper_bound = cercan_df['First_sexual_intercourse'].quantile(0.75) + (1.5 * (cercan_df['First_sexual_intercourse'].quantile(0.75) - cercan_df['First_sexual_intercourse'].quantile(0.25)))
3 print(f'First_sexual_intercourse outliers are value < {low_bound} or > {upper_bound}')
```

First_sexual_intercourse outliers are value < 10.5 or > 22.5

```
1 low_bound = cercan_df['Num_of_pregnancies'].quantile(0.25) - (1.5 * (cercan_df['Num_of_pregnancies'].quantile(0.75) - cercan_df['Num_of_pregnancies'].quantile(0.25)))
2 upper_bound = cercan_df['Num_of_pregnancies'].quantile(0.75) + (1.5 * (cercan_df['Num_of_pregnancies'].quantile(0.75) - cercan_df['Num_of_pregnancies'].quantile(0.25)))
3 print(f'Num_of_pregnancies outliers are value < {low_bound} or > {upper_bound}')
```

Num_of_pregnancies outliers are value < -2.0 or > 6.0

```
1 low_bound = cercan_df['Smokes_(years)'].quantile(0.25) - (1.5 * (cercan_df['Smokes_(years)'].quantile(0.75) - cercan_df['Smokes_(years)'].quantile(0.25)))
2 upper_bound = cercan_df['Smokes_(years)'].quantile(0.75) + (1.5 * (cercan_df['Smokes_(years)'].quantile(0.75) - cercan_df['Smokes_(years)'].quantile(0.25)))
3 print(f'Smokes_(years) outliers are value < {low_bound} or > {upper_bound}')
```

Smokes_(years) outliers are value < 0.0 or > 0.0

```
1 cercan_df['Smokes_(years)'].value_counts()
```

```
Smokes_(years)
0.000000    712
1.266973     15
5.000000      9
9.000000      9
1.000000      8
3.000000      7
2.000000      7
16.000000      6
7.000000      6
8.000000      6
11.000000      5
4.000000      5
10.000000      5
14.000000      4
15.000000      4
6.000000      4
13.000000      3
0.500000      3
19.000000      3
12.000000      3
22.000000      2
32.000000      1
20.000000      1
28.000000      1
24.000000      1
```

```

18.000000    1
34.000000    1
21.000000    1
37.000000    1
0.160000     1
Name: count, dtype: int64

```

```

1 low_bound = cercan_df['Smokes_(packs/year)'].quantile(0.25) - (1.5 * (cercan_df['Smokes_(packs/year)'].quantile(0.75) - cercan_df['Sr
2 upper_bound = cercan_df['Smokes_(packs/year)'].quantile(0.75) + (1.5 * (cercan_df['Smokes_(packs/year)'].quantile(0.75) - cercan_df[
3 print(f'Smokes_(packs/year) outliers are value < {low_bound} or > {upper_bound}')

```

```
Smokes_(packs/year) outliers are value < 0.0 or > 0.0
```

```
1 cercan_df['Smokes_(packs/year)'].value_counts()
```

```

Smokes_(packs/year)
0.000000    712
0.513202     18
1.000000     6
3.000000     5
2.000000     4
...
7.500000     1
37.000000     1
2.250000     1
0.003000     1
0.300000     1
Name: count, Length: 62, dtype: int64

```

```

1 low_bound = cercan_df['Hormonal_Contraceptives_(years)'].quantile(0.25) - (1.5 * (cercan_df['Hormonal_Contraceptives_(years)'].quantil
2 upper_bound = cercan_df['Hormonal_Contraceptives_(years)'].quantile(0.75) + (1.5 * (cercan_df['Hormonal_Contraceptives_(years)'].quant
3 print(f'Hormonal_Contraceptives_(years) outliers are value < {low_bound} or > {upper_bound}')

```

```
Hormonal_Contraceptives_(years) outliers are value < -4.5 or > 7.5
```

```

1 low_bound = cercan_df['IUD_(years)'].quantile(0.25) - (1.5 * (cercan_df['IUD_(years)'].quantile(0.75) - cercan_df['IUD_(years)'].quant
2 upper_bound = cercan_df['IUD_(years)'].quantile(0.75) + (1.5 * (cercan_df['IUD_(years)'].quantile(0.75) - cercan_df['IUD_(years)'].qua
3 print(f'IUD_(years) outliers are value < {low_bound} or > {upper_bound}')

```

```
IUD_(years) outliers are value < 0.0 or > 0.0
```

```
1 cercan_df['IUD_(years)'].value_counts()
```

```

IUD_(years)
0.00    752
3.00     11
2.00     10
5.00      9
1.00      8
8.00      7
7.00      7
6.00      5
4.00      5
11.00     3
0.08      2
0.50      2
0.33      1
9.00      1
0.41      1
0.16      1
0.91      1
1.50      1
10.00     1
12.00     1
15.00     1
0.25      1
17.00     1
19.00     1
0.58      1
0.17      1
Name: count, dtype: int64

```

```

1 low_bound = cercan_df['STDs_(number)'].quantile(0.25) - (1.5 * (cercan_df['STDs_(number)'].quantile(0.75) - cercan_df['STDs_(number)
2 upper_bound = cercan_df['STDs_(number)'].quantile(0.75) + (1.5 * (cercan_df['STDs_(number)'].quantile(0.75) - cercan_df['STDs_(number)
3 print(f'STDs_(number) outliers are value < {low_bound} or > {upper_bound}')

```

```
STDs_(number) outliers are value < 0.0 or > 0.0
```

```
1 cercan_df['STDs_(number)'].value_counts()
```

```

STDs_(number)
0.0    756

```

```

2.0    37
1.0    34
3.0     7
4.0     1
Name: count, dtype: int64

```

```

1 low_bound = cercan_df['STDs:_Number_of_diagnosis'].quantile(0.25) - (1.5 * (cercan_df['STDs:_Number_of_diagnosis'].quantile(0.75) - (
2 upper_bound = cercan_df['STDs:_Number_of_diagnosis'].quantile(0.75) + (1.5 * (cercan_df['STDs:_Number_of_diagnosis'].quantile(0.75) - (
3 print(f'STDs:_Number_of_diagnosis outliers are value < {low_bound} or > {upper_bound}'))

```

```
STDs:_Number_of_diagnosis outliers are value < 0.0 or > 0.0
```

```
1 cercan_df['STDs:_Number_of_diagnosis'].value_counts()
```

```

STDs:_Number_of_diagnosis
0    764
1     68
2      2
3      1
Name: count, dtype: int64

```

```

1 low_bound = cercan_df['STDs:_Time_since_first_diagnosis'].quantile(0.25) - (1.5 * (cercan_df['STDs:_Time_since_first_diagnosis'].quar
2 upper_bound = cercan_df['STDs:_Time_since_first_diagnosis'].quantile(0.75) + (1.5 * (cercan_df['STDs:_Time_since_first_diagnosis'].qi
3 print(f'STDs:_Time_since_first_diagnosis outliers are value < {low_bound} or > {upper_bound}'))

```

```
STDs:_Time_since_first_diagnosis outliers are value < 4.0 or > 4.0
```

```

1 low_bound = cercan_df['STDs:_Time_since_last_diagnosis'].quantile(0.25) - (1.5 * (cercan_df['STDs:_Time_since_last_diagnosis'].quant:
2 upper_bound = cercan_df['STDs:_Time_since_last_diagnosis'].quantile(0.75) + (1.5 * (cercan_df['STDs:_Time_since_last_diagnosis'].quar
3 print(f'STDs:_Time_since_last_diagnosis outliers are value < {low_bound} or > {upper_bound}'))

```

```
STDs:_Time_since_last_diagnosis outliers are value < 3.0 or > 3.0
```

✓ Declare feature vector and target variable

```

1 # since y does not have a dataframe then there is not target variable
2 X = cercan_df.drop(['Biopsy'], axis=1)
3 y = cercan_df['Biopsy']

```

✓ Split data into separate training and test set

```

1 # split X and y into training and testing sets
2 from sklearn.model_selection import train_test_split
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

```

```

1 # check the shape of X_train and X_test
2 X_train.shape, X_test.shape

```

```
((668, 35), (167, 35))
```

✓ Engineering outliers in numerical values

```

1 def max_value(df3, variable, top):
2     return np.where(df3[variable]>top, top, df3[variable])
3
4 for df3 in [X_train, X_test]:
5     df3['Age'] = max_value(df3, 'Age', 48.5)
6     df3['Number_of_sexual_partners'] = max_value(df3, 'Number_of_sexual_partners', 4.5)
7     df3['First_sexual_intercourse'] = max_value(df3, 'First_sexual_intercourse', 22.5)
8     df3['Num_of_pregnancies'] = max_value(df3, 'Num_of_pregnancies', 6.0)
9     df3['Hormonal_Contraceptives_(years)'] = max_value(df3, 'Hormonal_Contraceptives_(years)', 7.5)
10    df3['STDs:_Time_since_first_diagnosis'] = max_value(df3, 'STDs:_Time_since_first_diagnosis', 4.0)
11    df3['STDs:_Time_since_last_diagnosis'] = max_value(df3, 'STDs:_Time_since_last_diagnosis', 3.0)

```

```
1 X_train['Age'].max(), X_test['Age'].max()
```

```
(48.5, 48.5)
```

```
1 X_train['Number_of_sexual_partners'].max(), X_test['Number_of_sexual_partners'].max()
```

```
(4.5, 4.5)

1 X_train['First_sexual_intercourse'].max(), X_test['First_sexual_intercourse'].max()

(22.5, 22.5)

1 X_train['Num_of_pregnancies'].max(), X_test['Num_of_pregnancies'].max()

(6.0, 6.0)

1 X_train['Hormonal_Contraceptives_(years)'].max(), X_test['Hormonal_Contraceptives_(years)'].max()

(7.5, 7.5)

1 X_train['STDs:_Time_since_first_diagnosis'].max(), X_test['STDs:_Time_since_first_diagnosis'].max()

(4.0, 4.0)

1 X_train['STDs:_Time_since_last_diagnosis'].max(), X_test['STDs:_Time_since_last_diagnosis'].max()

(3.0, 3.0)
```

Feature Scaling

```
1 X_train.describe()
```

	Age	Number_of_sexual_partners	First_sexual_intercourse	Num_of_pregnancies	Smokes	Smokes_(years)	Smokes_(packs)
count	668.000000	668.000000	668.000000	668.000000	668.000000	668.000000	668
mean	26.738024	2.439371	16.855539	2.232036	0.157186	1.267709	0
std	7.613179	1.096057	2.289344	1.282470	0.364248	4.021316	2
min	13.000000	1.000000	10.000000	0.000000	0.000000	0.000000	0
25%	21.000000	2.000000	15.000000	1.000000	0.000000	0.000000	0
50%	26.000000	2.000000	17.000000	2.000000	0.000000	0.000000	0
75%	32.000000	3.000000	18.000000	3.000000	0.000000	0.000000	0
max	48.500000	4.500000	22.500000	6.000000	1.000000	37.000000	37

8 rows × 35 columns

```
1 cols = X_train.columns

1 from sklearn.preprocessing import MinMaxScaler
2
3 scaler = MinMaxScaler()
4 X_train = scaler.fit_transform(X_train)
5 X_test = scaler.fit_transform(X_test)

1 X_train = pd.DataFrame(X_train, columns=[cols])

1 X_test = pd.DataFrame(X_test, columns=[cols])

1 X_train.describe()
```

8 rows \times 35 columns

```
1 logreg.predict_proba(X_test)[:,-1]

array([0.0681819 , 0.01347467, 0.00869753, 0.01290119, 0.02662565,
        0.01586956, 0.23995319, 0.01001325, 0.01516459, 0.01886381,
        0.0375028 , 0.01333394, 0.00685647, 0.01404528, 0.53764177,
        0.01875839, 0.00988367, 0.01690609, 0.01573696, 0.01899051,
        0.00938461, 0.01714533, 0.06345496, 0.01551824, 0.01470809,
        0.04212003, 0.01579026, 0.01583826, 0.00429639, 0.67613157,
        0.06014656, 0.01566792, 0.30161164, 0.01495033, 0.02193021,
        0.01657418, 0.01285653, 0.01468255, 0.00987957, 0.01576251,
        0.01814364, 0.01586245, 0.01789926, 0.01523051, 0.00858158,
        0.01565034, 0.33374702, 0.01651829, 0.01331129, 0.01298081,
        0.01688045, 0.01531596, 0.01696703, 0.01597227, 0.0056862 ,
        0.06579735, 0.01734199, 0.01569216, 0.01404457, 0.01789167,
        0.0121981 , 0.02314671, 0.01317484, 0.01818465, 0.01780952,
        0.00938555, 0.02032783, 0.03116828, 0.01515568, 0.01683743,
        0.5686354 , 0.01427848, 0.02018063, 0.01906333, 0.05380697,
        0.01429906, 0.01584917, 0.4604003 , 0.01485665, 0.01323249,
        0.01530071, 0.48906646, 0.01372073, 0.01375002, 0.00771065,
        0.02069225, 0.01422335, 0.01643524, 0.01726076, 0.01768023,
        0.00393108, 0.01948758, 0.01545286, 0.01266695, 0.00708049,
        0.01509622, 0.01387173, 0.01441253, 0.00923305, 0.01347648,
        0.79265512, 0.02026134, 0.01956075, 0.01569309, 0.01634629,
        0.01830688, 0.59323342, 0.0149568 , 0.01692373, 0.0191292 ,
        0.01307902, 0.01068693, 0.0224706 , 0.01831756, 0.75915319,
        0.01333159, 0.65157887, 0.02036183, 0.01523858, 0.01130142,
        0.00992905, 0.0036863 , 0.01626376, 0.0168602 , 0.01441408,
        0.01684998, 0.0153332 , 0.01546251, 0.0118241 , 0.01695569,
        0.01394412, 0.6441652 , 0.00887825, 0.0162521 , 0.01162517,
        0.01750155, 0.78262516, 0.01443275, 0.00977053, 0.01562444,
        0.01615672, 0.01489604, 0.01127947, 0.00549581, 0.00951773,
        0.01538071, 0.0144317 , 0.01054269, 0.01151979, 0.01306861,
        0.04473572, 0.06128548, 0.01222499, 0.00920906, 0.01399183,
        0.00929528, 0.01551541, 0.4956282 , 0.01021991, 0.25699498,
        0.01296524, 0.01812096, 0.01607635, 0.01659872, 0.02094492,
        0.01500768, 0.01210755])
```

```
1 from sklearn.metrics import accuracy_score
2
3 print('Model accuracy score: {0:0.4f}'.format(accuracy_score(y_test, y_pred_test)))

Model accuracy score: 0.9461
```

[illegible]

Training-set accuracy score: 0.9671

macro avg	0.76	0.74	0.75	167
weighted avg	0.94	0.95	0.94	167

✓ Classification accuracy

```
1 TP = cm[0,0]
2 TN = cm[1,1]
3 FP = cm[0,1]
4 FN = cm[1,0]

1 # print classification accuracy
2 classification_accuracy = (TP + TN) / float(TP + TN + FP + FN)
3 print('Classification accuracy:{0:0.4f}'.format(classification_accuracy))

Classification accuracy:0.9461
```

✓ Classification error

```
1 # print classification error
2 classification_error = (FP + FN) / float(TP + TN + FP + FN)
3 print('Classification error:{0:0.4f}'.format(classification_error))

Classification error:0.0539
```

✓ Precision

```
1 # print precision score
2 precision = TP / float(TP + FP)
3 print('Precision:{0:0.4f}'.format(precision))

Precision:0.9745
```

✓ Recall

```
1 recall = TP / float(TP + FN)
2 print('Recall or Sensitivity: {0:0.4f}'.format(recall))

Recall or Sensitivity: 0.9684
```

✓ True Positive Rate

```
1 tpr = TP / float(TP + FN)
2 print('True Positive Rate: {0:0.4f}'.format(tpr))

True Positive Rate: 0.9684
```

✓ False Positive Rate

```
1 fpr = FP / float(FP + FN)
2 print('False Positive Rate: {0:0.4f}'.format(fpr))

False Positive Rate: 0.4444
```

✓ Specificity

```
1 specificity = TN / (TN + FP)
2 print('Specificity: {0:0.4f}'.format(specificity))

Specificity: 0.5556
```

✓ Adjusting the threshold level


```

1 # print the first 10 predicted probabilities of two classes- 0 and 1
2 y_pred_prob = logreg.predict_proba(X_test)[0:10]
3 y_pred_prob

```

```

array([[0.9318181 , 0.0681819 ],
       [0.98652533, 0.01347467],
       [0.99130247, 0.00869753],
       [0.98709881, 0.01290119],
       [0.97337435, 0.02662565],
       [0.98413044, 0.01586956],
       [0.76004681, 0.23995319],
       [0.98998675, 0.01001325],
       [0.98483541, 0.01516459],
       [0.98113619, 0.01886381]])

```

```

1 # store the probabilities in dataframe
2 y_pred_prob = pd.DataFrame(data=y_pred_prob, columns=['Prob of - No biopsy (0)', 'Prob of - Biopsy (1)'])
3 y_pred_prob

```

	Prob of - No biopsy (0)	Prob of - Biopsy (1)
0	0.931818	0.068182
1	0.986525	0.013475
2	0.991302	0.008698
3	0.987099	0.012901
4	0.973374	0.026626
5	0.984130	0.015870
6	0.760047	0.239953
7	0.989987	0.010013
8	0.984835	0.015165
9	0.981136	0.018864

Next steps: [View recommended plots](#)

```

1 # print the first 10 predicted probabilities for class 1 - Probability of biopsy
2 logreg.predict_proba(X_test)[0:10, 1]

```

```

array([0.0681819 , 0.01347467, 0.00869753, 0.01290119, 0.02662565,
       0.01586956, 0.23995319, 0.01001325, 0.01516459, 0.01886381])

```

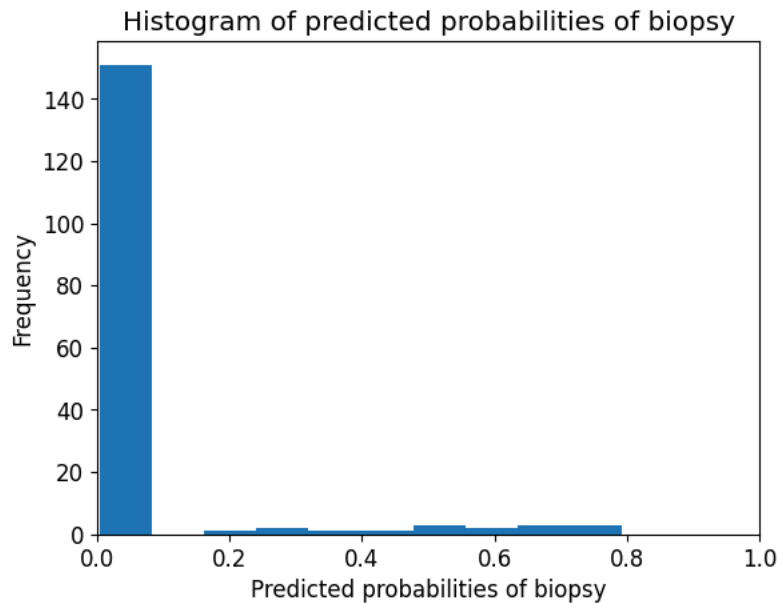
```
1 y_pred1 = logreg.predict_proba(X_test)[: ,1]
```

```

1 # plot histogram of predicted probabilities
2 # adjust the font size
3 plt.rcParams['font.size'] = 12
4
5 # plot histogram with 10 bins
6 plt.hist(y_pred1, bins=10)
7
8 # set the title of predicted probabilities
9 plt.title('Histogram of predicted probabilities of biopsy')
10
11 # set the x-axis limit
12 plt.xlim(0,1)
13
14 # set the title
15 plt.xlabel('Predicted probabilities of biopsy')
16 plt.ylabel('Frequency')

```

Text(0, 0.5, 'Frequency')



Observations:

- The histogram is highly positive skewed
- The first column tells us that there are approximately 160 observation with probability between 0.0 and 0.1
- The small number of observation predict the possibility of undergoing biopsy.
- Majority of observation predicts the possibility of not undergoing biopsy.

✓ Lower the threshold

```

1 from sklearn.preprocessing import binarize
2
3 for i in range(1,5):
4     cm1=0
5     y_pred1 = logreg.predict_proba(X_test)[: ,1]
6     y_pred1= y_pred1.reshape(-1,1)
7     y_pred2= binarize(y_pred1, threshold=i/10)
8     y_pred2 = np.where(y_pred2 == 1, 1, 0)
9     cm1 = confusion_matrix(y_test, y_pred2)
10    print ('With',i/10, 'threshold the Confusion Matrix is ', '\n\n', cm1, '\n\n',
11           'with', cm1[0,0]+cm1[1,1], 'correct predictions, ', '\n\n',
12           cm1[0,1], 'Type I errors( False Positives), ', '\n\n',
13           cm1[1,0], 'Type II errors( False Negatives), ', '\n\n',
14           'Accuracy score: ', (accuracy_score(y_test, y_pred2)), '\n\n',
15           'Sensitivity: ',cm1[1,1]/(float(cm1[1,1]+cm1[1,0])), '\n\n',
16           'Specificity: ',cm1[0,0]/(float(cm1[0,0]+cm1[0,1])), '\n\n',
17           '=====', '\n\n')

```

With 0.1 threshold the Confusion Matrix is

```
[[149  8]
 [ 2  8]]
```

with 157 correct predictions,

8 Type I errors(False Positives),

2 Type II errors(False Negatives),

Accuracy score: 0.9401197604790419

Sensitivity: 0.8

Specificity: 0.9490445859872612

=====

With 0.2 threshold the Confusion Matrix is

```
[[149  8]
 [ 2  8]]
```

```

with 157 correct predictions,

8 Type I errors( False Positives),

2 Type II errors( False Negatives),

Accuracy score:  0.9401197604790419

Sensitivity:  0.8

Specificity:  0.9490445859872612

=====

```

With 0.3 threshold the Confusion Matrix is

```

[[150  7]
 [ 3  7]]

```

```

with 157 correct predictions,

7 Type I errors( False Positives),

3 Type II errors( False Negatives),

Accuracy score:  0.9401197604790419

Sensitivity:  0.7

Specificity:  0.9554140127388535

=====

```

Comments

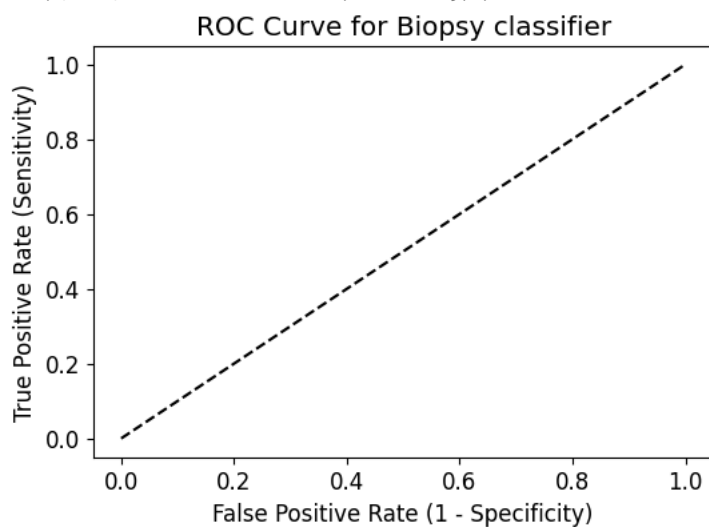
✓ ROC - AUC

```

1 from sklearn.metrics import roc_curve
2
3 fpr, tpr, thresholds = roc_curve(y_test, y_pred1, pos_label = 'Yes')
4
5 plt.figure(figsize=(6,4))
6
7 plt.plot(fpr,tpr, linewidth=2)
8 plt.plot([0,1], [0,1], 'k--')
9 plt.rcParams['font.size'] =12
10 plt.title('ROC Curve for Biopsy classifier')
11 plt.xlabel('False Positive Rate (1 - Specificity)')
12 plt.ylabel('True Positive Rate (Sensitivity)')

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_ranking.py:1029: UndefinedMe
warnings.warn(
Text(0, 0.5, 'True Positive Rate (Sensitivity)')

```



```

1 # compute ROC AUC
2 from sklearn.metrics import roc_auc_score
3
4 ROC_AUC = roc_auc_score(y_test, y_pred1)
5 print('ROC AUC : {:.4f}'.format(ROC_AUC))

ROC AUC : 0.9350

1 # calculate cross-validated ROC AUC
2 from sklearn.model_selection import cross_val_score
3
4 Cross_validated_ROC_AUC = cross_val_score(logreg, X_train, y_train, cv=5, scoring='roc_auc').mean()
5 print('Cross validated ROC AUC : {:.4f}'.format(Cross_validated_ROC_AUC))

Cross validated ROC AUC : 0.9405

```

✓ k-Fold Cross Validation

```

1 # Applying 5-Fold Cross Validation
2 from sklearn.model_selection import cross_val_score
3
4 scores = cross_val_score(logreg, X_train, y_train, cv = 5, scoring='accuracy')
5 print('Cross-validation scores:{}'.format(scores))

Cross-validation scores:[0.96268657 0.96268657 0.94029851 0.96240602 0.94736842]

1 # compute Average cross-validation score
2 print('Average cross-validation score: {:.4f}'.format(scores.mean()))

Average cross-validation score: 0.9551

```

✓ Hyperparameter Optimization using GridSearch CV

```

1 # Hyperparameter Optimization using GridSearchCV
2 from sklearn.model_selection import GridSearchCV
3
4 parameters = [{'penalty': ['l1', 'l2']],
5               {'C':[1, 10, 100, 1000]}]
6
7 grid_search = GridSearchCV(estimator = logreg,
8                             param_grid = parameters,
9                             scoring = 'accuracy',
10                            cv = 5,
11                            verbose=0)
12
13 grid_search.fit(X_train, y_train)

```

```
1 # examine the best model
2 # best score achieved during the GridSearchCV
3 print('GridSearch CV best score : {:.4f}\n\n'.format(grid_search.best_score_))
4
5 # print parameters that give the best results
6 print('Parameters that give the best results : ', '\n\n', (grid_search.best_params_))
7
8 # print estimator that was chosen by the GridSearch
9 print('\n\nEstimator that was chosen by the search : ', '\n\n', (grid_search.best_estimator_))
```

Parameters that give the best results :

Estimator that was chosen by the search :

[/usr/local/lib/python3.10/nltk-packages/sklearn/model_selection/_search.py,334](https://userwarning.github.io/python3.10/nltk-packages/sklearn/model_selection/_search.py,334). UserWarning: One or more of the test scores are non

Gridsearch CV score on test set: 0.9281

.....

1. The logistic regression model accuracy score is 0.9461. So, the model does a very good job in predicting whether a patient having the aforementioned illness would undergo biopsy.
2. Small number of observations predict that patients would undergo biopsy. Majority of observations predict that patients would not undergo biopsy.
3. The model shows no signs of overfitting.
4. ROC AUC of our model approaches towards 1. So, we can conclude that our classifier does a good job in predicting whether a patient having the aforementioned illness would undergo biopsy.
5. Our original model test accuracy is 0.9461 while GridSearch CV accuracy is 0.9281. We can see that GridSearch CV decreased the performance for this particular model.