

Project Name: Evaluation and Prediction of Recurrence of Thyroid Cancer

Institution Name: Shaheed Sukhdev College of Business Studies

Guidance Under: Dr. Rishi Ranjan Sahay

Course Name: Certificate course on Data Analytics & Business Intelligence

Made By: Kriti Khurana & Kritika Mittal

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

df=pd.read_csv("Thyroid_Diff.csv")

df.head()
```

	Age	Gender	Smoking	Hx	Smoking	Hx	Radiothreapy	Thyroid	Function	\
0	27	F	No		No		No		Euthyroid	
1	34	F	No		Yes		No		Euthyroid	
2	30	F	No		No		No		Euthyroid	
3	62	F	No		No		No		Euthyroid	
4	62	F	No		No		No		Euthyroid	

	Physical Examination	Adenopathy	Pathology	Focality
Risk \				
0	Single nodular goiter-left	No	Micropapillary	Uni-Focal
Low				
1	Multinodular goiter	No	Micropapillary	Uni-Focal
Low				
2	Single nodular goiter-right	No	Micropapillary	Uni-Focal
Low				
3	Single nodular goiter-right	No	Micropapillary	Uni-Focal
Low				
4	Multinodular goiter	No	Micropapillary	Multi-Focal
Low				

	T	N	M	Stage	Response	Recurred
0	T1a	N0	M0	I	Indeterminate	No
1	T1a	N0	M0	I	Excellent	No

2	T1a	N0	M0	I	Excellent	No
3	T1a	N0	M0	I	Excellent	No
4	T1a	N0	M0	I	Excellent	No

```
df.shape
```

```
(383, 17)
```

```
df.columns
```

```
Index(['Age', 'Gender', 'Smoking', 'Hx Smoking', 'Hx Radiothreapy',
      'Thyroid Function', 'Physical Examination', 'Adenopathy',
      'Pathology',
      'Focality', 'Risk', 'T', 'N', 'M', 'Stage', 'Response',
      'Recurred'],
      dtype='object')
```

```
df.tail()
```

	Age	Gender	Smoking	Hx Smoking	Hx Radiothreapy	Thyroid
378	72	M	Yes	Yes	Yes	Euthyroid
379	81	M	Yes	No	Yes	Euthyroid
380	72	M	Yes	Yes	No	Euthyroid
381	61	M	Yes	Yes	Yes	Clinical Hyperthyroidism
382	67	M	Yes	No	No	Euthyroid

	Physical Examination	Adenopathy	Pathology	Focality
378	Single nodular goiter-right	Right	Papillary	Uni-Focal
379	Multinodular goiter	Extensive	Papillary	Multi-Focal
380	Multinodular goiter	Bilateral	Papillary	Multi-Focal
381	Multinodular goiter	Extensive	Hurthel cell	Multi-Focal
382	Multinodular goiter	Bilateral	Papillary	Multi-Focal

	T	N	M	Stage	Response	Recurred
378	T4b	N1b	M1	IVB	Biochemical Incomplete	Yes
379	T4b	N1b	M1	IVB	Structural Incomplete	Yes
380	T4b	N1b	M1	IVB	Structural Incomplete	Yes
381	T4b	N1b	M0	IVA	Structural Incomplete	Yes
382	T4b	N1b	M0	IVA	Structural Incomplete	Yes

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 383 entries, 0 to 382
```

```
Data columns (total 17 columns):
```

#	Column	Non-Null Count	Dtype
0	Age	383 non-null	int64
1	Gender	383 non-null	object
2	Smoking	383 non-null	object
3	Hx Smoking	383 non-null	object
4	Hx Radiothreapy	383 non-null	object
5	Thyroid Function	383 non-null	object
6	Physical Examination	383 non-null	object
7	Adenopathy	383 non-null	object
8	Pathology	383 non-null	object
9	Focality	383 non-null	object
10	Risk	383 non-null	object
11	T	383 non-null	object
12	N	383 non-null	object
13	M	383 non-null	object
14	Stage	383 non-null	object
15	Response	383 non-null	object
16	Recurred	383 non-null	object

```
dtypes: int64(1), object(16)
```

```
memory usage: 51.0+ KB
```

```
df.describe()
```

	Age
count	383.000000
mean	40.866841
std	15.134494
min	15.000000
25%	29.000000
50%	37.000000
75%	51.000000
max	82.000000

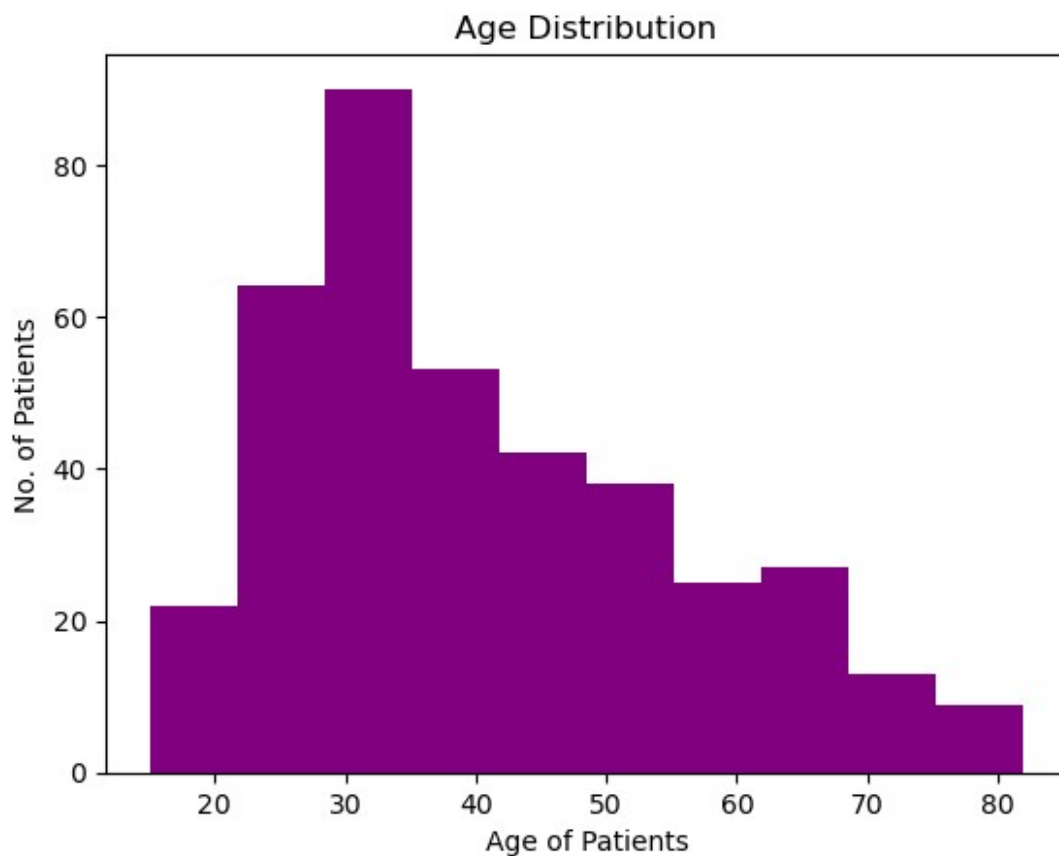
```
df.isnull().sum()
```

Age	0
Gender	0
Smoking	0
Hx Smoking	0
Hx Radiothreapy	0
Thyroid Function	0
Physical Examination	0
Adenopathy	0
Pathology	0

```
Focality          0
Risk              0
T                0
N                0
M                0
Stage            0
Response          0
Recurred         0
dtype: int64
```

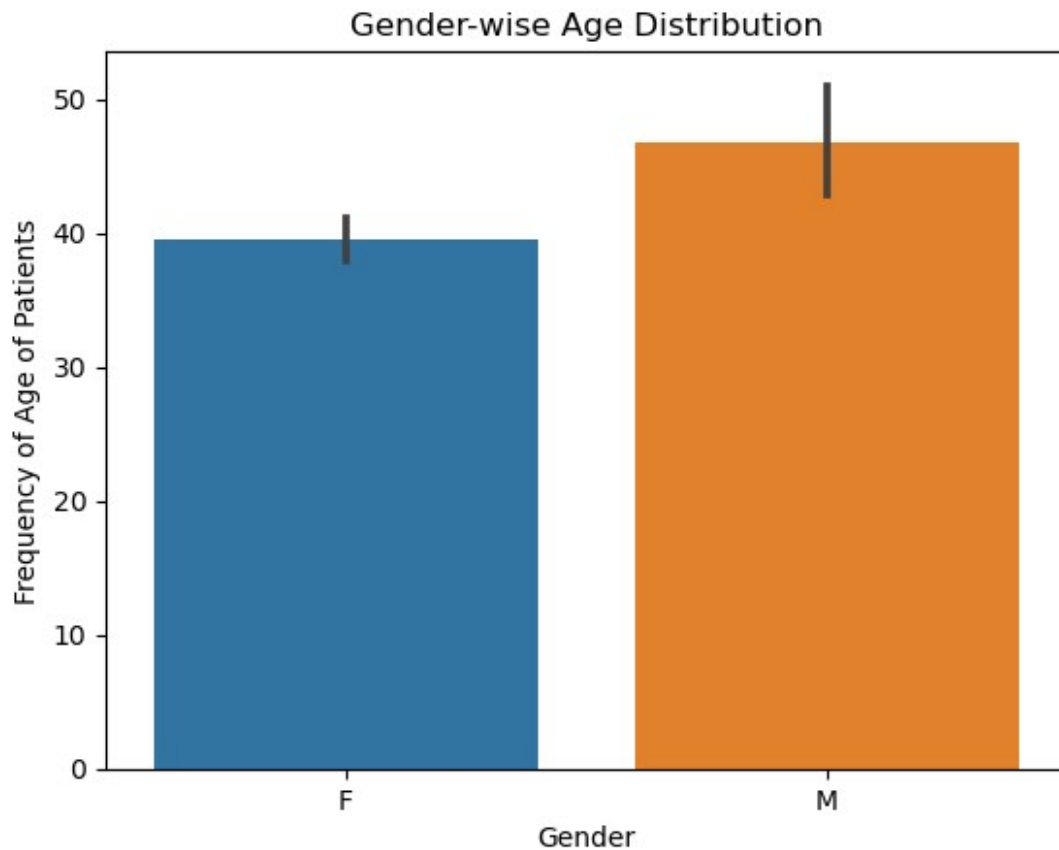
EXPLORATORY DATA ANALYSIS (EDA)

```
# 1. Age Distribution Graph
plt.hist(df['Age'],bins=10, color='purple')
plt.title('Age Distribution')
plt.xlabel('Age of Patients')
plt.ylabel('No. of Patients')
plt.savefig('insight1.png')
plt.show()
```

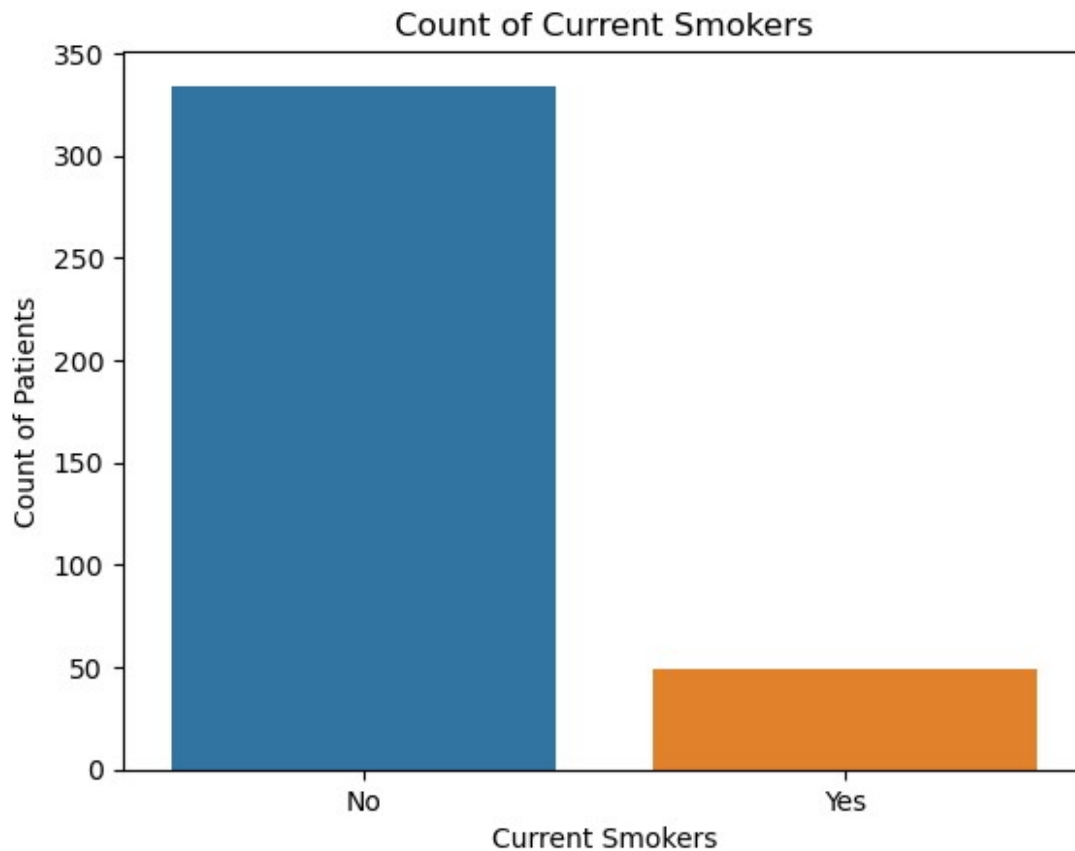


```
# 2. Gender-wise Age Distribution
sns.barplot(x='Gender',y='Age',data=df)
```

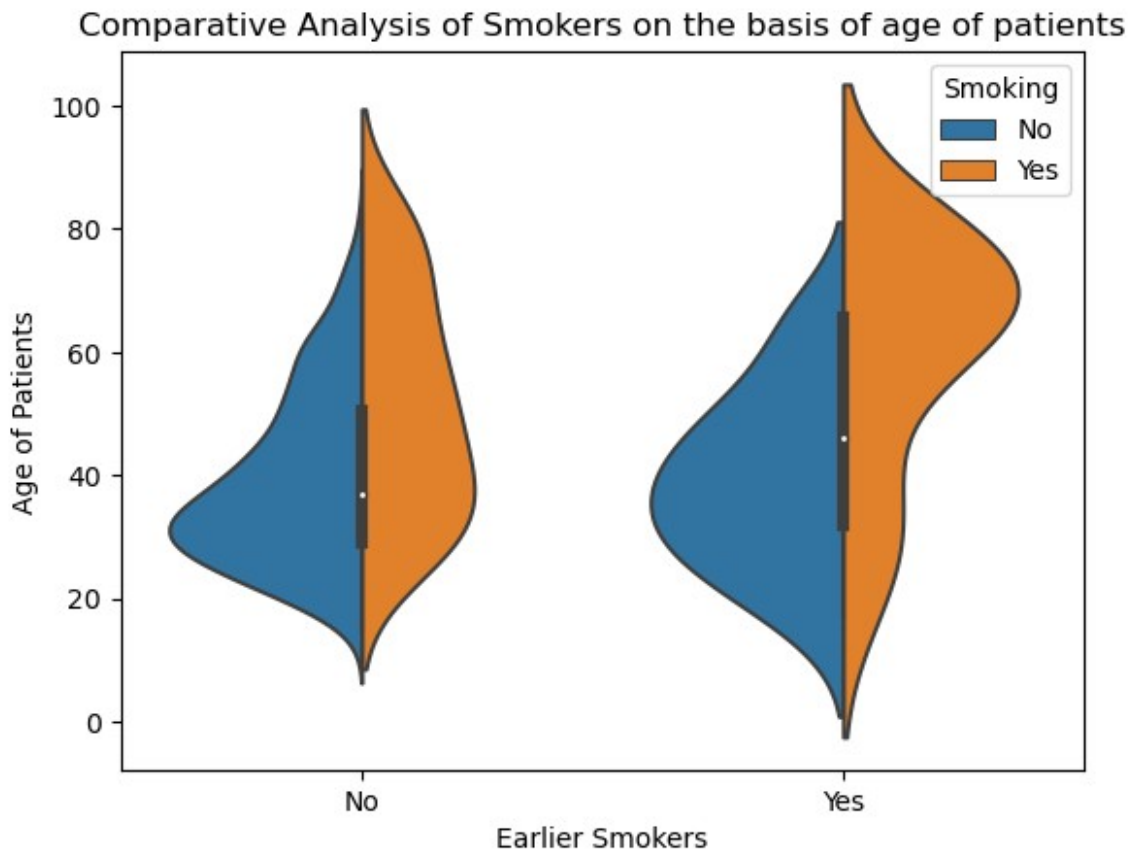
```
plt.title('Gender-wise Age Distribution')
plt.xlabel('Gender')
plt.ylabel('Frequency of Age of Patients')
plt.savefig('insight2.png')
plt.show()
```



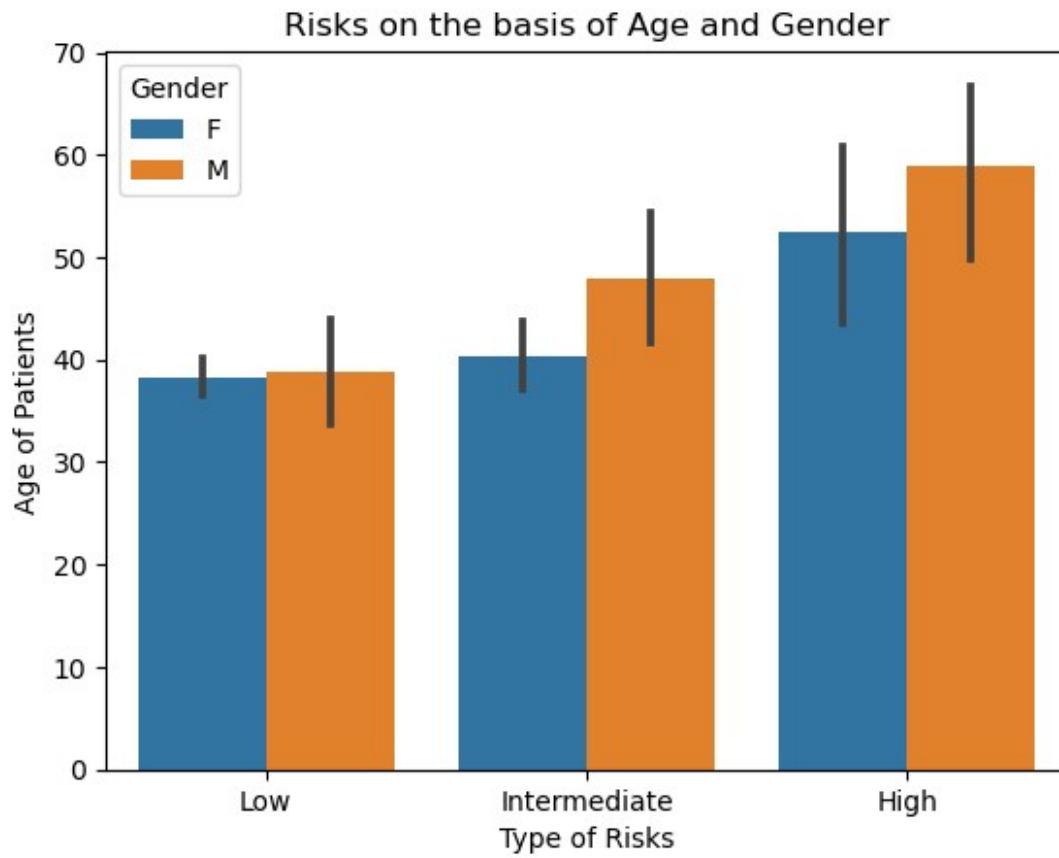
```
# 3. Count of Current Smokers
sns.countplot(x='Smoking',data=df)
plt.title('Count of Current Smokers')
plt.xlabel('Current Smokers')
plt.ylabel('Count of Patients')
plt.savefig('insight3.png')
plt.show()
```



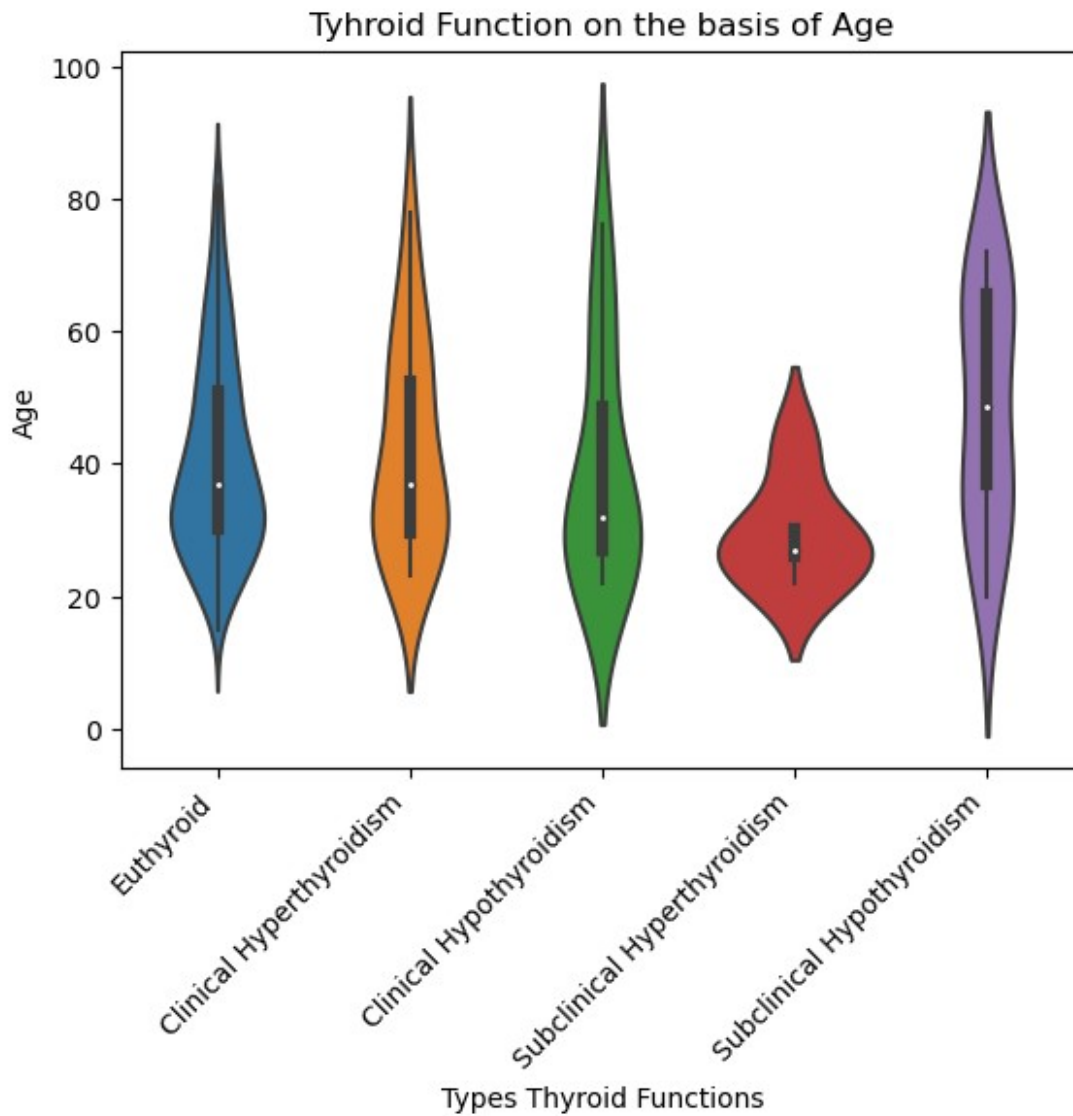
```
# 4. Comparative Analysis of Smokers on the basis of age of patients
sns.violinplot(x='Hx
Smoking',y='Age',data=df,hue='Smoking',split=True)
plt.title('Comparative Analysis of Smokers on the basis of age of
patients')
plt.xlabel('Earlier Smokers')
plt.ylabel('Age of Patients')
plt.savefig('insight4.png')
plt.show()
```



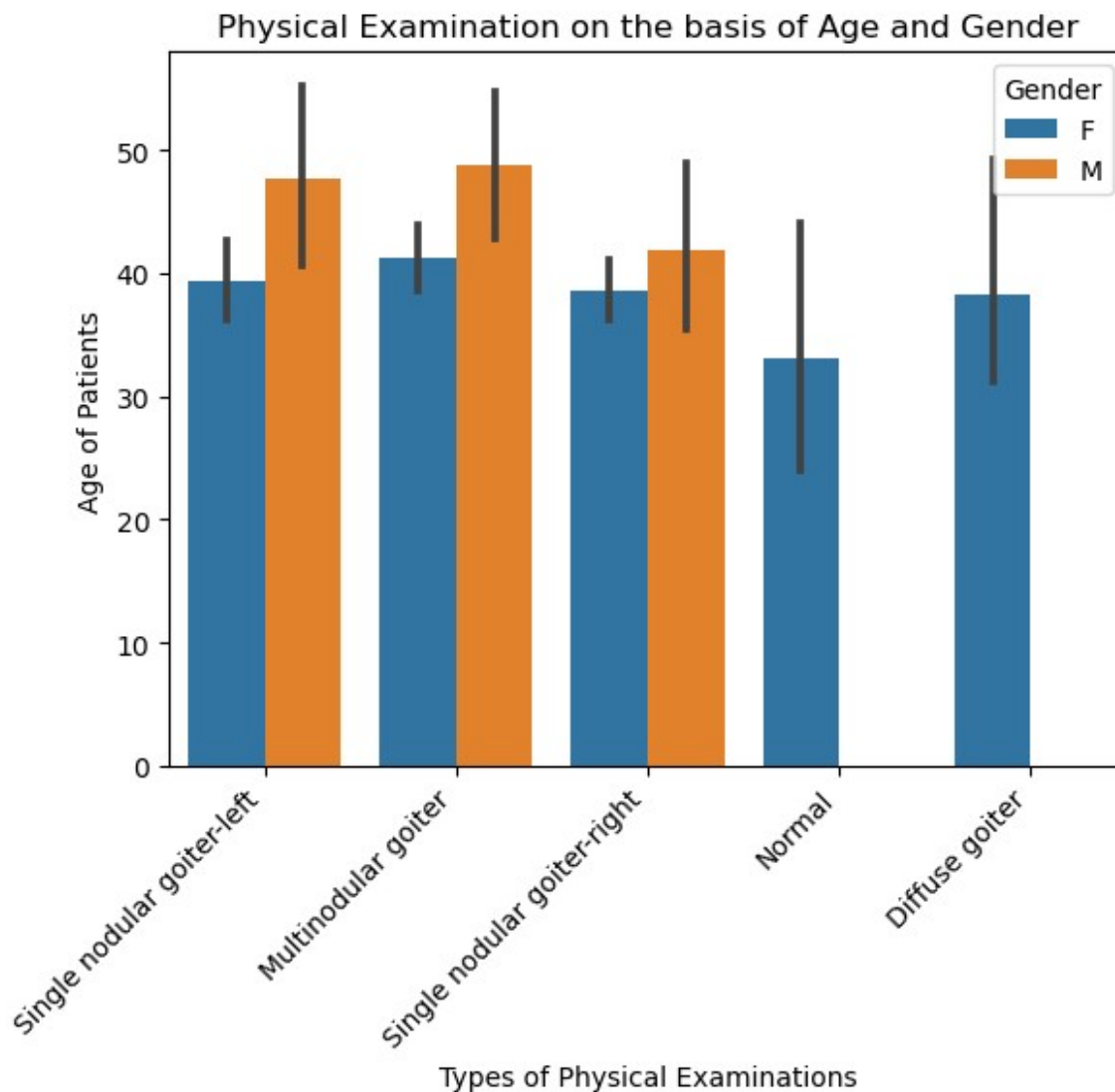
```
# 5. Risk on the basis of Age and Gender
sns.barplot(x='Risk',y='Age',data=df,hue='Gender')
plt.title('Risks on the basis of Age and Gender')
plt.xlabel('Type of Risks')
plt.ylabel('Age of Patients')
plt.savefig('insight5.png')
plt.show()
```



```
# 6. Thyroid Function on the basis of Age
sns.violinplot(x='Thyroid Function',y='Age', data=df)
plt.title('Thyroid Function on the basis of Age')
plt.xlabel('Types Thyroid Functions')
plt.ylabel('Age')
plt.xticks(rotation=45,ha='right')
plt.savefig('insight6.png')
plt.show()
```

```
# 7. Physical Examination on the basis of Age and Gender
sns.barplot(x='Physical Examination',y='Age',data=df, hue='Gender')
plt.title('Physical Examination on the basis of Age and Gender')
plt.xticks(rotation=45,ha='right')
plt.xlabel('Types of Physical Examinations')
plt.ylabel('Age of Patients')
plt.savefig('insight7.png')
plt.show()
```



CLASSIFICATION

```
# Identify categorical and continuous variables
categorical = ['Gender', 'Smoking', 'Hx Smoking', 'Hx
Radiothreapy', 'Thyroid Function', 'Physical Examination',
'Adenopathy', 'Pathology',
'Focality', 'Risk', 'T', 'N', 'M', 'Stage', 'Response',
'Recurred']
continuous = ['Age']

# Convert categorical variables to dummy variables
df = pd.get_dummies(df, columns=categorical, drop_first=True)

df
```

Age	Gender_M	Smoking_Yes	Hx Smoking_Yes	Hx Radiothreapy_Yes	\

0	27	False	False	False	False
1	34	False	False	True	False
2	30	False	False	False	False
3	62	False	False	False	False
4	62	False	False	False	False
..
378	72	True	True	True	True
379	81	True	True	False	True
380	72	True	True	True	False
381	61	True	True	True	True
382	67	True	True	False	False

Thyroid Function_Clinical Hypothyroidism		Thyroid
Function_Euthyroid \		
0		False
True		
1		False
True		
2		False
True		
3		False
True		
4		False
True		
..		...
...		
378		False
True		
379		False
True		
380		False
True		
381		False
False		
382		False
True		

Thyroid Function_Subclinical Hyperthyroidism \	
0	False

1	False
2	False
3	False
4	False
..	...
378	False
379	False
380	False
381	False
382	False

Thyroid Function_Subclinical Hypothyroidism \	
0	False
1	False
2	False
3	False
4	False
..	...
378	False
379	False
380	False
381	False
382	False

Physical Examination_Multinodular goiter ...		N_N1b	M_M1
Stage_II \			
0	False	...	False
False			False
1	True	...	False
False			False
2	False	...	False
False			False
3	False	...	False
False			False
4	True	...	False
False			False
..
...			
378	False	...	True
False			True
379	True	...	True
False			True
380	True	...	True
False			True
381	True	...	True
False			False
382	True	...	True
False			False

Stage_III	Stage_IVA	Stage_IVB	Response_Excellent \
-----------	-----------	-----------	----------------------

0	False	False	False	False
1	False	False	False	True
2	False	False	False	True
3	False	False	False	True
4	False	False	False	True
...
378	False	False	True	False
379	False	False	True	False
380	False	False	True	False
381	False	True	False	False
382	False	True	False	False

	Response_Indeterminate Recurred_Yes	Response_Structural Incomplete
0	True	False
False		
1	False	False
False		
2	False	False
False		
3	False	False
False		
4	False	False
False		
...
...		
378	False	False
True		
379	False	True
True		
380	False	True
True		
381	False	True
True		
382	False	True
True		

[383 rows x 41 columns]

changing the data type for getting the interger values

df = df.astype(int)

df.info()

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 383 entries, 0 to 382

Data columns (total 41 columns):

#	Column	Non-Null Count
Dtype		
---	-----	-----

0	Age	383 non-null
int32		
1	Gender_M	383 non-null
int32		
2	Smoking_Yes	383 non-null
int32		
3	Hx Smoking_Yes	383 non-null
int32		
4	Hx Radiothreapy_Yes	383 non-null
int32		
5	Thyroid Function_Clinical Hypothyroidism	383 non-null
int32		
6	Thyroid Function_Euthyroid	383 non-null
int32		
7	Thyroid Function_Subclinical Hyperthyroidism	383 non-null
int32		
8	Thyroid Function_Subclinical Hypothyroidism	383 non-null
int32		
9	Physical Examination_Multinodular goiter	383 non-null
int32		
10	Physical Examination_Normal	383 non-null
int32		
11	Physical Examination_Single nodular goiter-left	383 non-null
int32		
12	Physical Examination_Single nodular goiter-right	383 non-null
int32		
13	Adenopathy_Extensive	383 non-null
int32		
14	Adenopathy_Left	383 non-null
int32		
15	Adenopathy_No	383 non-null
int32		
16	Adenopathy_Posterior	383 non-null
int32		
17	Adenopathy_Right	383 non-null
int32		
18	Pathology_Hurthel cell	383 non-null
int32		
19	Pathology_Micropapillary	383 non-null
int32		
20	Pathology_Papillary	383 non-null
int32		
21	Focality_Uni-Focal	383 non-null
int32		
22	Risk_Intermediate	383 non-null
int32		
23	Risk_Low	383 non-null
int32		
24	T_T1b	383 non-null

```

int32
 25  T_T2                                383 non-null
int32
 26  T_T3a                              383 non-null
int32
 27  T_T3b                              383 non-null
int32
 28  T_T4a                              383 non-null
int32
 29  T_T4b                              383 non-null
int32
 30  N_N1a                              383 non-null
int32
 31  N_N1b                              383 non-null
int32
 32  M_M1                               383 non-null
int32
 33  Stage_II                           383 non-null
int32
 34  Stage_III                           383 non-null
int32
 35  Stage_IVA                           383 non-null
int32
 36  Stage_IVB                           383 non-null
int32
 37  Response_Excellent                  383 non-null
int32
 38  Response_Indeterminate              383 non-null
int32
 39  Response_Structural Incomplete      383 non-null
int32
 40  Recurred_Yes                        383 non-null
int32
dtypes: int32(41)
memory usage: 61.5 KB

```

```

#splitting the dataset into training and testing
from sklearn.model_selection import train_test_split
# Splitting the dataset into training and testing sets
x = df.drop('Recurred_Yes', axis=1)
y = df['Recurred_Yes']
x_train, x_test, y_train, y_test = train_test_split(x, y,
test_size=0.2, random_state=20)

```

```
x_train
```

	Age	Gender_M	Smoking_Yes	Hx	Smoking_Yes	Hx
Radiothreapy_Yes \						
0	27	0	0		0	0

228	20	1	0	0	0
373	31	1	1	0	1
346	32	1	0	1	0
148	33	0	0	0	0
..
331	51	0	0	0	0
218	48	1	0	0	0
223	56	0	0	0	0
271	45	0	0	0	0
355	32	0	0	0	0

Thyroid Function_Clinical Hypothyroidism Thyroid
Function_Euthyroid \

0	0
1	
228	0
1	
373	0
1	
346	0
1	
148	0
1	
..	...
...	
331	0
1	
218	0
1	
223	0
1	
271	0
1	
355	0
1	

Thyroid Function_Subclinical Hyperthyroidism \	
0	0
228	0
373	0

346	0					
148	0					
..	...					
331	0					
218	0					
223	0					
271	0					
355	0					
Thyroid Function_Subclinical Hypothyroidism \						
0	0					
228	0					
373	0					
346	0					
148	0					
..	...					
331	0					
218	0					
223	0					
271	0					
355	0					
Physical Examination_Multinodular goiter ... N_N1a N_N1b M_M1						
\						
0	0	...	0	0	0	0
228	0	...	0	1	0	
373	0	...	0	1	1	
346	1	...	0	1	0	
148	1	...	0	0	0	
..	
331	0	...	0	0	0	
218	0	...	1	0	0	
223	0	...	0	0	0	
271	0	...	0	0	0	
355	0	...	0	0	0	
Stage_II Stage_III Stage_IVA Stage_IVB Response_Excellent \						
0	0	0	0	0	0	
228	0	0	0	0	0	
373	1	0	0	0	0	

346	0	0	0	0	0
148	0	0	0	0	1
..
331	0	0	0	0	0
218	0	0	0	0	0
223	0	0	0	0	0
271	0	0	0	0	1
355	0	0	0	0	1

	Response_Indeterminate	Response_Structural Incomplete
0	1	0
228	0	1
373	0	1
346	0	1
148	0	0
..
331	0	1
218	1	0
223	0	1
271	0	0
355	0	0

[306 rows x 40 columns]

y_train

0	0
228	1
373	1
346	1
148	0
..	..
331	1
218	0
223	1
271	0
355	0

Name: Recurred_Yes, Length: 306, dtype: int32

x_test

	Age	Gender_M	Smoking_Yes	Hx Smoking_Yes	Hx
Radiothreapy_Yes \					
303	73	0	0	0	0
145	29	0	0	0	0
160	28	0	0	0	0
174	50	0	0	0	0

239	33	1	0	0	0
..
312	27	1	0	0	0
115	37	0	0	0	0
104	33	0	0	0	0
254	31	1	1	1	0
322	63	1	1	0	0

Thyroid Function_Clinical Hypothyroidism Thyroid
Function_Euthyroid \

303	0
1	
145	0
1	
160	1
0	
174	0
1	
239	0
1	
..	...
...	
312	0
1	
115	0
1	
104	1
0	
254	0
1	
322	0
1	

Thyroid Function_Subclinical Hyperthyroidism \

303	0
145	0
160	0
174	0
239	0
..	...
312	0
115	0
104	0

254	0
322	0
Thyroid Function_Subclinical Hypothyroidism \	
303	0
145	0
160	0
174	0
239	0
..	...
312	0
115	0
104	0
254	0
322	0
Physical Examination_Multinodular goiter ... N_N1a N_N1b M_M1	
\	
303	1 ... 0 0 0
145	0 ... 0 0 0
160	0 ... 0 0 0
174	1 ... 0 0 0
239	1 ... 0 0 0
..
312	1 ... 0 1 0
115	0 ... 0 0 0
104	0 ... 0 0 0
254	0 ... 0 0 0
322	0 ... 0 1 0
Stage_II Stage_III Stage_IVA Stage_IVB Response_Excellent \	
303	1 0 0 0 0
145	0 0 0 0 1
160	0 0 0 0 1
174	0 0 0 0 0
239	0 0 0 0 0
..
312	0 0 0 0 0
115	0 0 0 0 1
104	0 0 0 0 1

254	0	0	0	0	0
322	1	0	0	0	0

	Response_Indeterminate	Response_Structural Incomplete
303	0	1
145	0	0
160	0	0
174	1	0
239	0	1
..
312	0	1
115	0	0
104	0	0
254	1	0
322	0	1

[77 rows x 40 columns]

y_test

303	1
145	0
160	0
174	0
239	1
..	
312	1
115	0
104	0
254	0
322	1

Name: Recurred_Yes, Length: 77, dtype: int32

```
from sklearn.metrics import accuracy_score, precision_score,
recall_score, f1_score, confusion_matrix
```

Logistic Regression

```
from sklearn.linear_model import LogisticRegression
```

```
model1=LogisticRegression()
```

```
model1.fit(x_train,y_train)
```

```
C:\Users\KRITI\anaconda3\Lib\site-packages\sklearn\linear_model\
_logistic.py:460: ConvergenceWarning: lbfgs failed to converge
(status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>
Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result(
LogisticRegression()

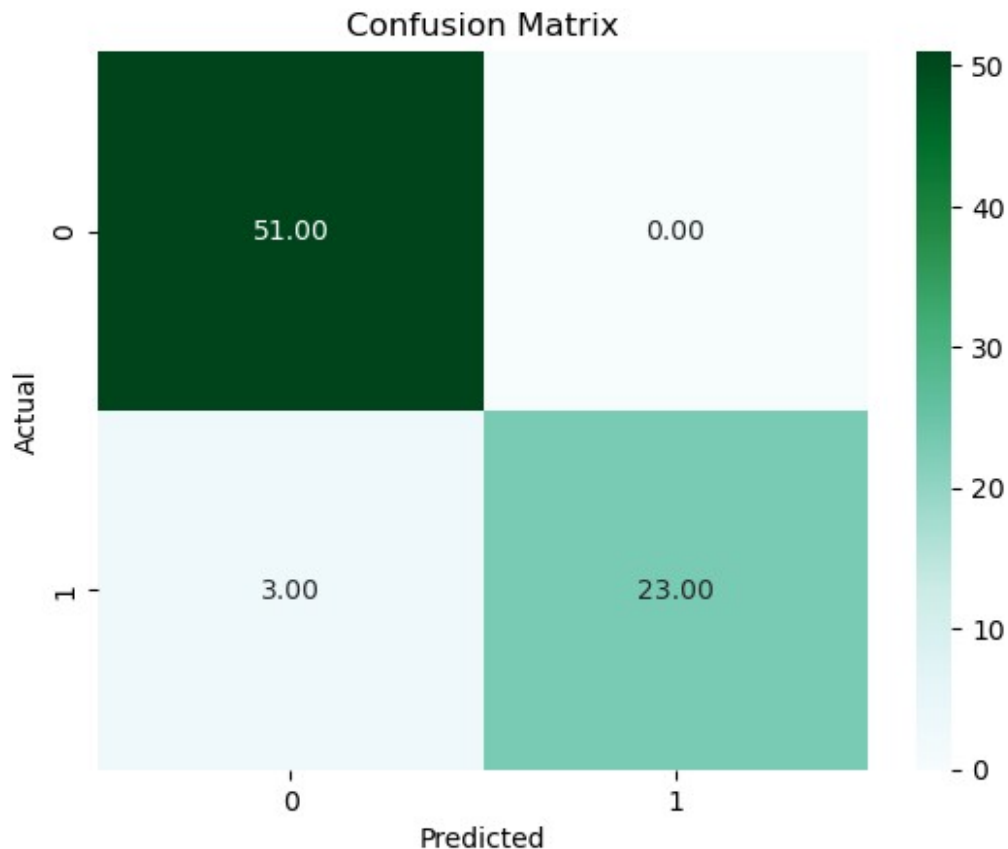
y_pred = model1.predict(x_test)
y_pred

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

cm = confusion_matrix(y_test, y_pred)
print(cm)
print("Accuracy score is",accuracy_score(y_test, y_pred))
print("Precision score is",precision_score(y_test, y_pred))
print("Recall score is",recall_score(y_test, y_pred))
print("F1 score is",f1_score(y_test, y_pred))

[[51  0]
 [ 3 23]]
Accuracy score is 0.961038961038961
Precision score is 1.0
Recall score is 0.8846153846153846
F1 score is 0.9387755102040816

# Plot the confusion matrix as a heatmap
sns.heatmap(cm, annot=True, fmt=".2f", cmap="BuGn")
plt.title('Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.savefig('ConfusionLogistic.png')
plt.show()
```



K-Nearest Neighbors Classifier

```
from sklearn.preprocessing import StandardScaler

# Feature scaling
scaler = StandardScaler()
x_train = scaler.fit_transform(x_train)
x_test = scaler.transform(x_test)

from sklearn.neighbors import KNeighborsClassifier
model2 = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
model2.fit(x_train, y_train)

KNeighborsClassifier()

y_pred = model2.predict(x_test)

y_pred
array([0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0,
0,
1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0,
0,
```

```
0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0,
1,
0, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1])
```

```
cm = confusion_matrix(y_test, y_pred)
print(cm)
print("Accuracy score is",accuracy_score(y_test, y_pred))
print("Precision score is",precision_score(y_test, y_pred))
print("Recall score is",recall_score(y_test, y_pred))
print("F1 score is",f1_score(y_test, y_pred))
```

```
[[47  4]
 [ 4 22]]
```

Accuracy score is 0.8961038961038961

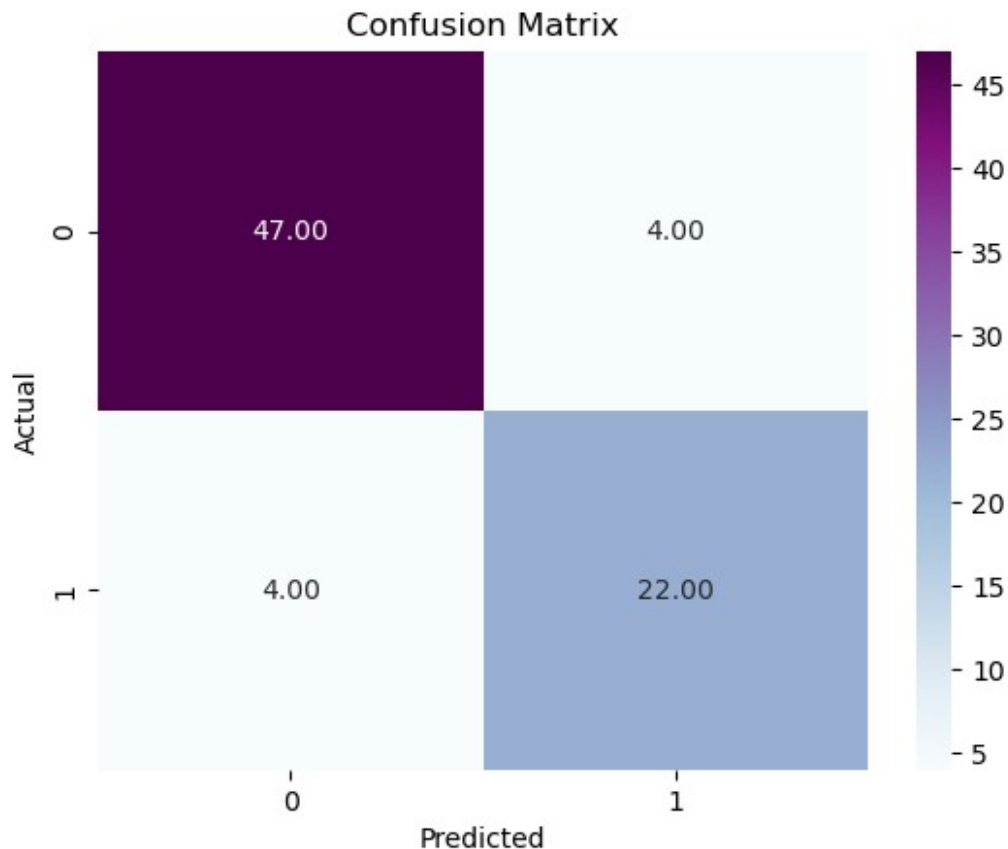
Precision score is 0.8461538461538461

Recall score is 0.8461538461538461

F1 score is 0.8461538461538461

Plot the confusion matrix as a heatmap

```
sns.heatmap(cm, annot=True, fmt=".2f", cmap="BuPu")
plt.title('Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.savefig('ConfusionKNN.png')
plt.show()
```

Support Vector Machine

```
from sklearn.svm import SVC

model3 = SVC(kernel = 'rbf', random_state = 20)
model3.fit(x_train, y_train)

SVC(random_state=20)

y_pred = model3.predict(x_test)

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

cm = confusion_matrix(y_test, y_pred)
print(cm)
print("Accuracy score is", accuracy_score(y_test, y_pred))
```

```

print("Precision score is",precision_score(y_test, y_pred))
print("Recall score is",recall_score(y_test, y_pred))
print("F1 score is",f1_score(y_test, y_pred))

[[49  2]
 [ 1 25]]
Accuracy score is 0.961038961038961
Precision score is 0.9259259259259259
Recall score is 0.9615384615384616
F1 score is 0.9433962264150944

# Plot the confusion matrix as a heatmap
sns.heatmap(cm, annot=True, fmt=".2f", cmap="GnBu")
plt.title('Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.savefig('ConfusionSVM.png')
plt.show()

```

