

# AIML Project

## GROUP 6:

9 - JAIDEEP SINGH HUNDAL  
10 - SHREYA JAGADALE  
25 - MRINAAL PALIWAL  
41 - PALLAVI SAWANT

In [1]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.tree import DecisionTreeClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.svm import SVC
from sklearn.linear_model import LinearRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import KFold, LeaveOneOut, ShuffleSplit, StratifiedKFold
from sklearn.model_selection import cross_val_score
from sklearn.naive_bayes import BernoulliNB
from sklearn.naive_bayes import GaussianNB
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import mean_squared_error, mean_absolute_error
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.preprocessing import StandardScaler
from sklearn.tree import DecisionTreeClassifier
from sklearn.tree import DecisionTreeRegressor
from sklearn import tree
from sklearn.datasets import make_blobs
from sklearn.metrics import silhouette_score
from sklearn.cluster import KMeans
from sklearn import datasets
from sklearn import decomposition
from sklearn.feature_selection import SelectPercentile, SelectKBest
import warnings
warnings.filterwarnings('ignore') #ignoring warnings
```

In [2]:

```
data = pd.read_csv('C:\\Users\\Shreyash\\Desktop\\data.csv')
data
```

Out[2]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes:
0	842302	M	17.99	10.38	122.80	1001.0	
1	842517	M	20.57	17.77	132.90	1326.0	
2	84300903	M	19.69	21.25	130.00	1203.0	
3	84348301	M	11.42	20.38	77.58	386.1	
4	84358402	M	20.29	14.34	135.10	1297.0	
...	...	...	...	...	...	...	
564	926424	M	21.56	22.39	142.00	1479.0	
565	926682	M	20.13	28.25	131.20	1261.0	
566	926954	M	16.60	28.08	108.30	858.1	
567	927241	M	20.60	29.33	140.10	1265.0	
568	92751	B	7.76	24.54	47.92	181.0	

569 rows × 33 columns



# Data Cleaning

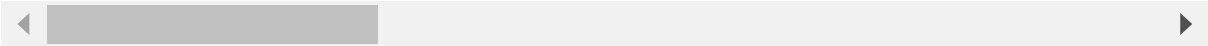
In [3]:

```
data.head()
```

Out[3]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_1
0	842302	M	17.99	10.38	122.80	1001.0	0.1
1	842517	M	20.57	17.77	132.90	1326.0	0.0
2	84300903	M	19.69	21.25	130.00	1203.0	0.1
3	84348301	M	11.42	20.38	77.58	386.1	0.1
4	84358402	M	20.29	14.34	135.10	1297.0	0.1

5 rows × 33 columns



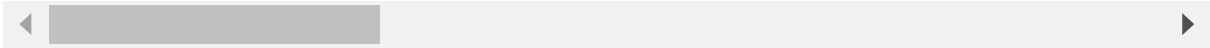
In [4]:

```
data.head(10)
```

Out[4]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_l
0	842302	M	17.99	10.38	122.80	1001.0	0.1
1	842517	M	20.57	17.77	132.90	1326.0	0.0
2	84300903	M	19.69	21.25	130.00	1203.0	0.1
3	84348301	M	11.42	20.38	77.58	386.1	0.1
4	84358402	M	20.29	14.34	135.10	1297.0	0.1
5	843786	M	12.45	15.70	82.57	477.1	0.1
6	844359	M	18.25	19.98	119.60	1040.0	0.0
7	84458202	M	13.71	20.83	90.20	577.9	0.1
8	844981	M	13.00	21.82	87.50	519.8	0.1
9	84501001	M	12.46	24.04	83.97	475.9	0.1

10 rows × 33 columns



In [5]:

```
data.info()
```

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

```
RangeIndex: 569 entries, 0 to 568
```

```
Data columns (total 33 columns):
```

#	Column	Non-Null Count	Dtype
0	id	569 non-null	int64
1	diagnosis	569 non-null	object
2	radius_mean	569 non-null	float64
3	texture_mean	569 non-null	float64
4	perimeter_mean	569 non-null	float64
5	area_mean	569 non-null	float64
6	smoothness_mean	569 non-null	float64
7	compactness_mean	569 non-null	float64
8	concavity_mean	569 non-null	float64
9	concave points_mean	569 non-null	float64
10	symmetry_mean	569 non-null	float64
11	fractal_dimension_mean	569 non-null	float64
12	radius_se	569 non-null	float64
13	texture_se	569 non-null	float64
14	perimeter_se	569 non-null	float64
15	area_se	569 non-null	float64
16	smoothness_se	569 non-null	float64
17	compactness_se	569 non-null	float64
18	concavity_se	569 non-null	float64
19	concave points_se	569 non-null	float64
20	symmetry_se	569 non-null	float64
21	fractal_dimension_se	569 non-null	float64
22	radius_worst	569 non-null	float64
23	texture_worst	569 non-null	float64
24	perimeter_worst	569 non-null	float64
25	area_worst	569 non-null	float64
26	smoothness_worst	569 non-null	float64
27	compactness_worst	569 non-null	float64
28	concavity_worst	569 non-null	float64
29	concave points_worst	569 non-null	float64
30	symmetry_worst	569 non-null	float64
31	fractal_dimension_worst	569 non-null	float64
32	Unnamed: 32	0 non-null	float64

```
dtypes: float64(31), int64(1), object(1)
```

```
memory usage: 146.8+ KB
```

In [6]:

```
data.shape
```

Out[6]:

```
(569, 33)
```

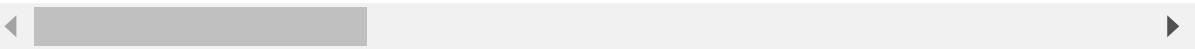
In [7]:

```
data.isnull()
```

Out[7]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_m
0	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False
...	...	...	...	...	...	...	...
564	False	False	False	False	False	False	False
565	False	False	False	False	False	False	False
566	False	False	False	False	False	False	False
567	False	False	False	False	False	False	False
568	False	False	False	False	False	False	False

569 rows × 33 columns

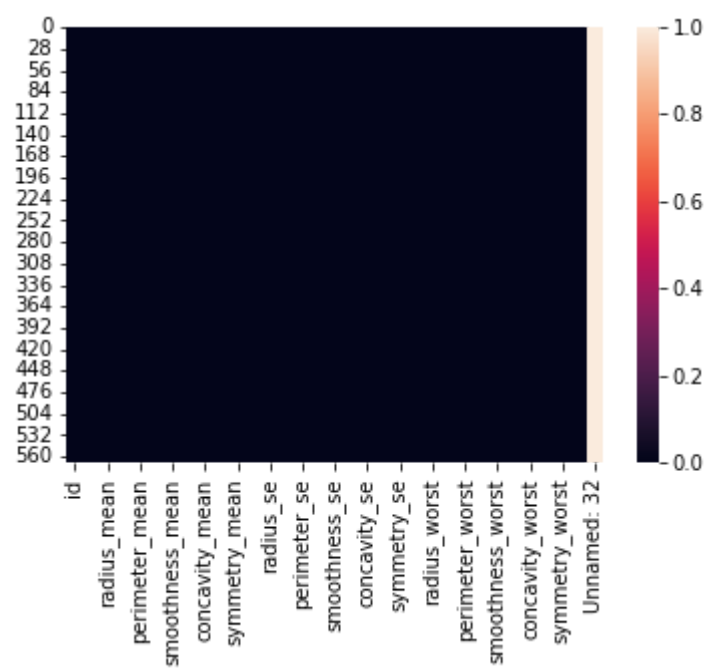


In [8]:

```
sns.heatmap(data.isnull())
```

Out[8]:

<AxesSubplot:>

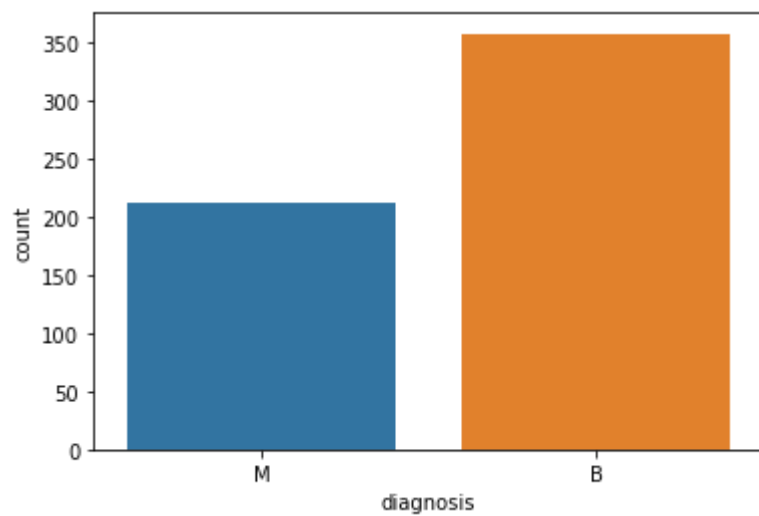


In [9]:

```
sns.countplot(x='diagnosis', data=data)
```

Out[9]:

<AxesSubplot:xlabel='diagnosis', ylabel='count'>



In [10]:

```
data.diagnosis.value_counts()
```

Out[10]:

```
B    357
M    212
Name: diagnosis, dtype: int64
```

In [11]:

```
data.dtypes
```

Out[11]:

```
id                int64
diagnosis         object
radius_mean       float64
texture_mean      float64
perimeter_mean    float64
area_mean         float64
smoothness_mean   float64
compactness_mean  float64
concavity_mean    float64
concave points_mean float64
symmetry_mean     float64
fractal_dimension_mean float64
radius_se         float64
texture_se        float64
perimeter_se      float64
area_se           float64
smoothness_se     float64
compactness_se    float64
concavity_se      float64
concave points_se float64
symmetry_se       float64
fractal_dimension_se float64
radius_worst      float64
texture_worst     float64
perimeter_worst   float64
area_worst        float64
smoothness_worst  float64
compactness_worst float64
concavity_worst   float64
concave points_worst float64
symmetry_worst    float64
fractal_dimension_worst float64
Unnamed: 32       float64
dtype: object
```

In [12]:

```
data['diagnosis'].unique()
```

Out[12]:

```
array(['M', 'B'], dtype=object)
```

In [13]:

```
channel_map = {'M':0, 'B':1}
```

In [14]:

```
data['diagnosis'] = data['diagnosis'].map(channel_map)
```

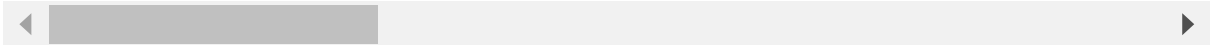
In [15]:

data

Out[15]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes:
0	842302	0	17.99	10.38	122.80	1001.0	
1	842517	0	20.57	17.77	132.90	1326.0	
2	84300903	0	19.69	21.25	130.00	1203.0	
3	84348301	0	11.42	20.38	77.58	386.1	
4	84358402	0	20.29	14.34	135.10	1297.0	
...	...	...	...	...	...	...	
564	926424	0	21.56	22.39	142.00	1479.0	
565	926682	0	20.13	28.25	131.20	1261.0	
566	926954	0	16.60	28.08	108.30	858.1	
567	927241	0	20.60	29.33	140.10	1265.0	
568	92751	1	7.76	24.54	47.92	181.0	

569 rows × 33 columns





In [16]:

```
data.info()
```

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

```
RangeIndex: 569 entries, 0 to 568
```

```
Data columns (total 33 columns):
```

#	Column	Non-Null Count	Dtype
0	id	569 non-null	int64
1	diagnosis	569 non-null	int64
2	radius_mean	569 non-null	float64
3	texture_mean	569 non-null	float64
4	perimeter_mean	569 non-null	float64
5	area_mean	569 non-null	float64
6	smoothness_mean	569 non-null	float64
7	compactness_mean	569 non-null	float64
8	concavity_mean	569 non-null	float64
9	concave points_mean	569 non-null	float64
10	symmetry_mean	569 non-null	float64
11	fractal_dimension_mean	569 non-null	float64
12	radius_se	569 non-null	float64
13	texture_se	569 non-null	float64
14	perimeter_se	569 non-null	float64
15	area_se	569 non-null	float64
16	smoothness_se	569 non-null	float64
17	compactness_se	569 non-null	float64
18	concavity_se	569 non-null	float64
19	concave points_se	569 non-null	float64
20	symmetry_se	569 non-null	float64
21	fractal_dimension_se	569 non-null	float64
22	radius_worst	569 non-null	float64
23	texture_worst	569 non-null	float64
24	perimeter_worst	569 non-null	float64
25	area_worst	569 non-null	float64
26	smoothness_worst	569 non-null	float64
27	compactness_worst	569 non-null	float64
28	concavity_worst	569 non-null	float64
29	concave points_worst	569 non-null	float64
30	symmetry_worst	569 non-null	float64
31	fractal_dimension_worst	569 non-null	float64
32	Unnamed: 32	0 non-null	float64

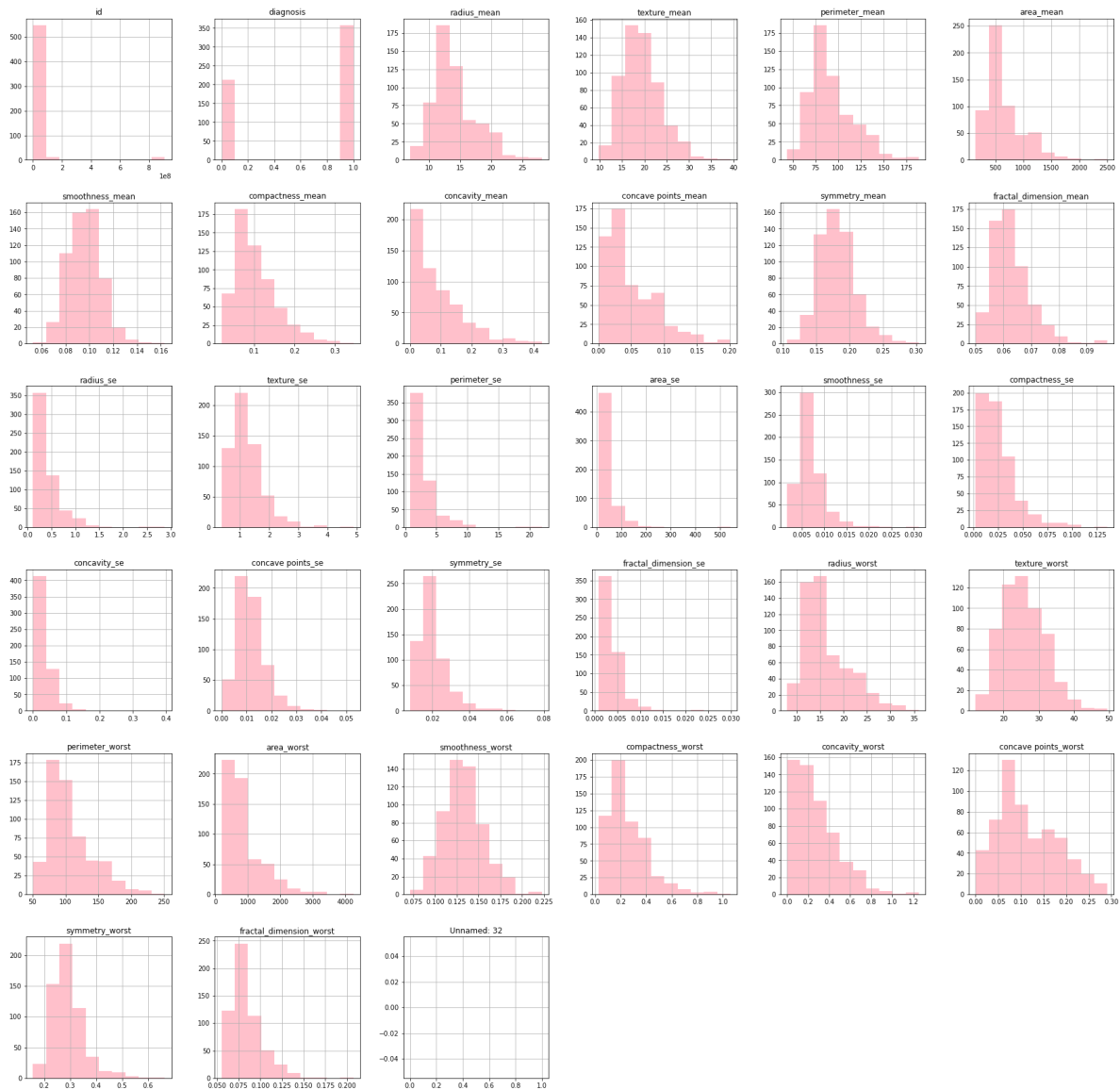
```
dtypes: float64(31), int64(2)
```

```
memory usage: 146.8 KB
```

## Data Visualization

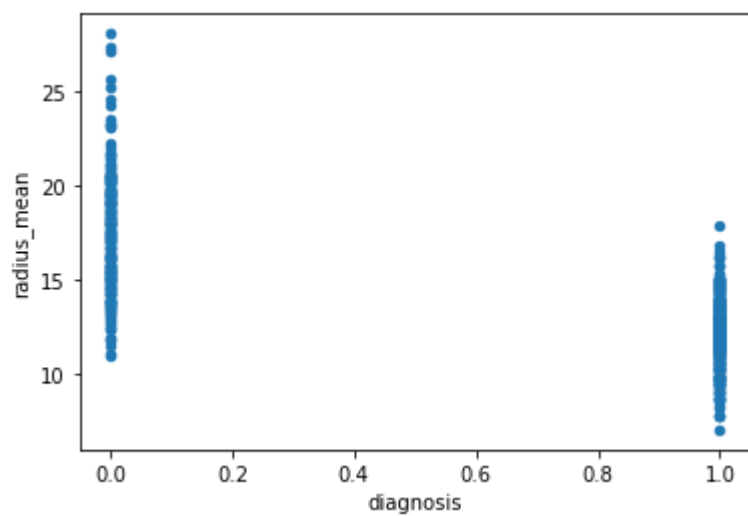
In [17]:

```
#plot the histograms for each feature:  
data.hist(figsize = (30,30), color = 'pink')  
plt.show()
```



In [18]:

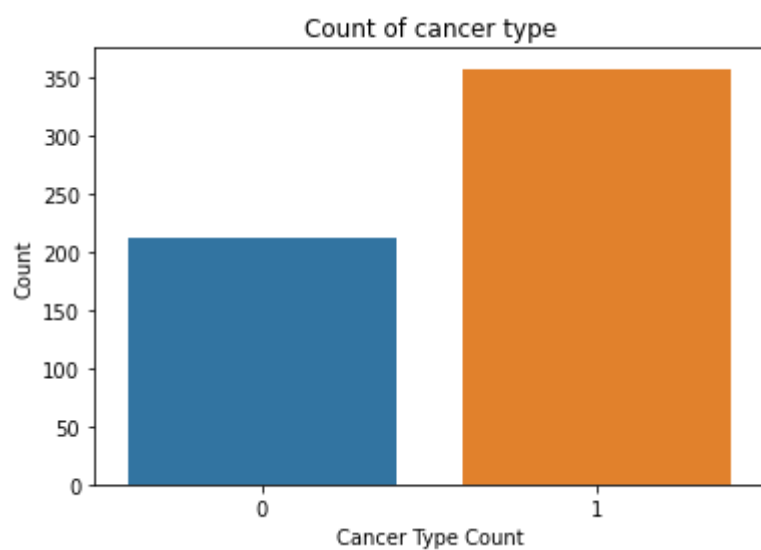
```
#scatterplot  
data.plot.scatter(x='diagnosis',y='radius_mean');
```



In [19]:

```
# Analyzing the target variable
```

```
plt.title('Count of cancer type')  
sns.countplot(data['diagnosis'])  
plt.xlabel('Cancer Type Count')  
plt.ylabel('Count')  
plt.show()
```



In [20]:

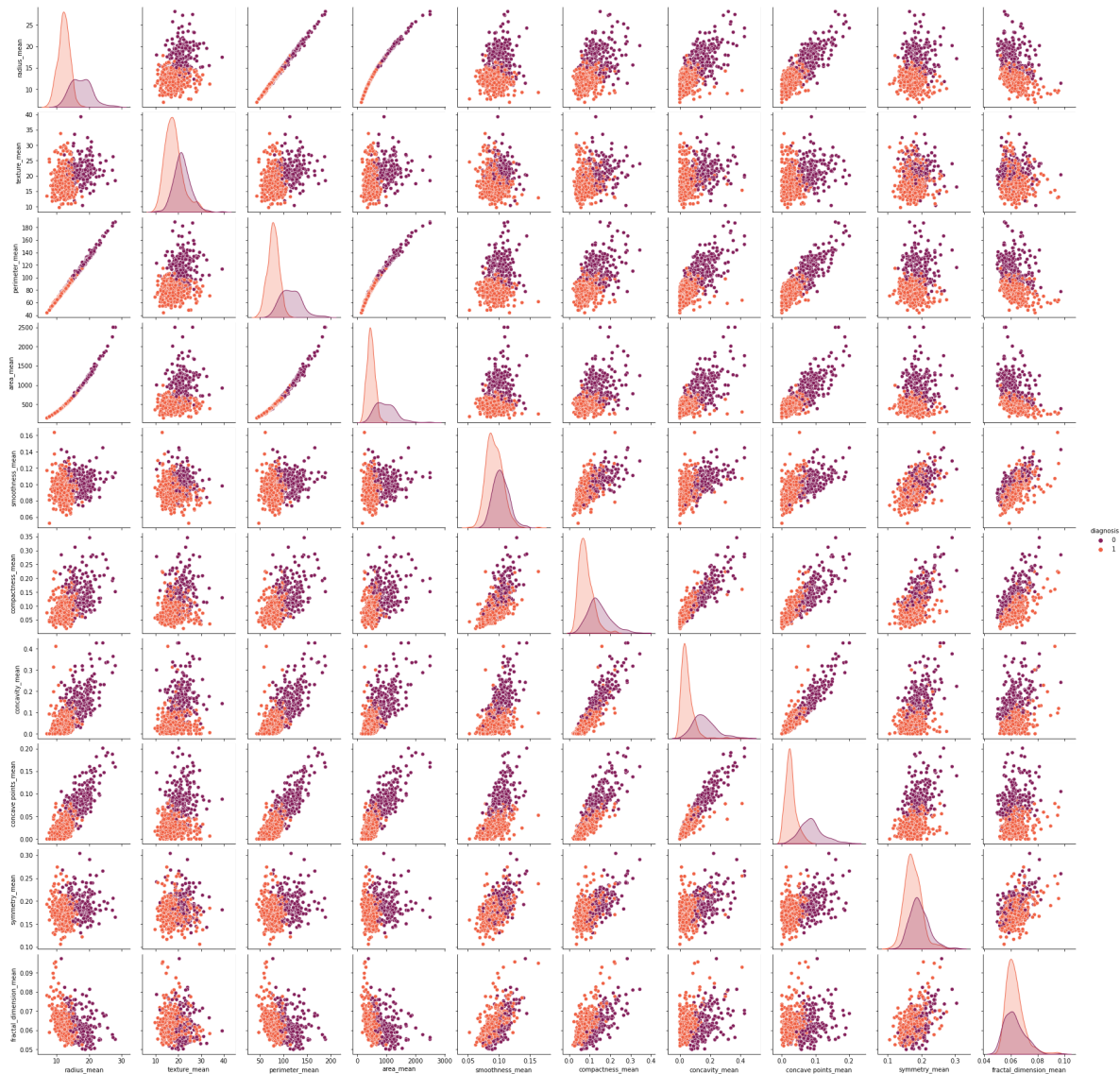
*#generate a scatter plot with the following columns:*

```
columns = ['diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean', 'smoothness_mean', 'compactness_mean', 'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean']
```

```
sns.pairplot(data=data[columns], hue="diagnosis", palette='rocket')
```

Out[20]:

<seaborn.axisgrid.PairGrid at 0x26220609ee0>



In [21]:

```
cor = data.corr()  
cor
```

Out[21]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	are
id	1.000000	-0.039769	0.074626	0.099770	0.073159	0
diagnosis	-0.039769	1.000000	-0.730029	-0.415185	-0.742636	-0
radius_mean	0.074626	-0.730029	1.000000	0.323782	0.997855	0
texture_mean	0.099770	-0.415185	0.323782	1.000000	0.329533	0
perimeter_mean	0.073159	-0.742636	0.997855	0.329533	1.000000	0
area_mean	0.096893	-0.708984	0.987357	0.321086	0.986507	1
smoothness_mean	-0.012968	-0.358560	0.170581	-0.023389	0.207278	0
compactness_mean	0.000096	-0.596534	0.506124	0.236702	0.556936	0
concavity_mean	0.050080	-0.696360	0.676764	0.302418	0.716136	0
concave points_mean	0.044158	-0.776614	0.822529	0.293464	0.850977	0
symmetry_mean	-0.022114	-0.330499	0.147741	0.071401	0.183027	0
fractal_dimension_mean	-0.052511	0.012838	-0.311631	-0.076437	-0.261477	-0
radius_se	0.143048	-0.567134	0.679090	0.275869	0.691765	0
texture_se	-0.007526	0.008303	-0.097317	0.386358	-0.086761	-0
perimeter_se	0.137331	-0.556141	0.674172	0.281673	0.693135	0
area_se	0.177742	-0.548236	0.735864	0.259845	0.744983	0
smoothness_se	0.096781	0.067016	-0.222600	0.006614	-0.202694	-0
compactness_se	0.033961	-0.292999	0.206000	0.191975	0.250744	0
concavity_se	0.055239	-0.253730	0.194204	0.143293	0.228082	0
concave points_se	0.078768	-0.408042	0.376169	0.163851	0.407217	0
symmetry_se	-0.017306	0.006522	-0.104321	0.009127	-0.081629	-0
fractal_dimension_se	0.025725	-0.077972	-0.042641	0.054458	-0.005523	-0
radius_worst	0.082405	-0.776454	0.969539	0.352573	0.969476	0
texture_worst	0.064720	-0.456903	0.297008	0.912045	0.303038	0
perimeter_worst	0.079986	-0.782914	0.965137	0.358040	0.970387	0
area_worst	0.107187	-0.733825	0.941082	0.343546	0.941550	0
smoothness_worst	0.010338	-0.421465	0.119616	0.077503	0.150549	0
compactness_worst	-0.002968	-0.590998	0.413463	0.277830	0.455774	0
concavity_worst	0.023203	-0.659610	0.526911	0.301025	0.563879	0
concave points_worst	0.035174	-0.793566	0.744214	0.295316	0.771241	0
symmetry_worst	-0.044224	-0.416294	0.163953	0.105008	0.189115	0
fractal_dimension_worst	-0.029866	-0.323872	0.007066	0.119205	0.051019	0
Unnamed: 32	NaN	NaN	NaN	NaN	NaN	

33 rows × 33 columns

In [22]:

```
cor.shape
```

Out[22]:

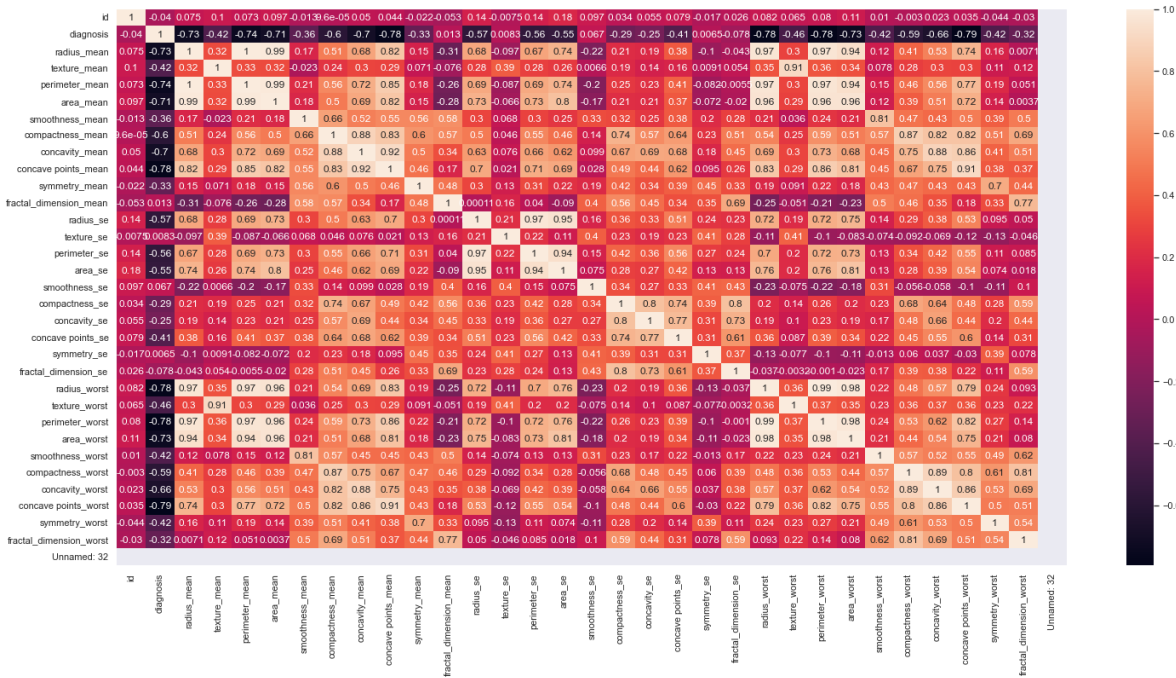
(33, 33)

In [23]:

```
sns.set(rc = {'figure.figsize':(25,12)})
sns.heatmap(cor,annot=True)
```

Out[23]:

<AxesSubplot:>



From above correlation we can see that 'concave points\_worst' has the highest absolute correlation with diagnosis

In [24]:

```
srx = data['concave points_worst']
srx.shape
```

Out[24]:

(569,)

In [25]:

```
srxx = srx.values.reshape(-1,1)
srxx.shape
```

Out[25]:

(569, 1)

## We will use 'concave points\_worst' as feature for simple LR

In [26]:

```
y = data['diagnosis']
y.shape
```

Out[26]:

(569,)

In [27]:

```
ydt = y.values.reshape(-1,1)
ydt.shape
```

Out[27]:

(569, 1)

In [28]:

```
xtr, xts, ytr, yts = train_test_split(srxx,ydt, test_size = 0.1, random_state = 0)
print("Size of training set:", xtr.shape)
print("Size of testing set:", xts.shape)
```

Size of training set: (512, 1)

Size of testing set: (57, 1)

## 9 - JAIDEEP SINGH

### Multivariate Logistic Regression

In [62]:

```
LR=LogisticRegression()
```

In [63]:

```
LR.fit(xtr,ytr)
```

Out[63]:

```
LogisticRegression()
```

In [64]:

```
y_pred = LR.predict(xts)
```

In [65]:

```
acc = mean_squared_error(yts, y_pred)
acc
```

Out[65]:

```
0.04678362573099415
```

In [66]:

```
LR.score(xts, yts)
```

Out[66]:

```
0.9532163742690059
```

In [67]:

```
df=pd.DataFrame({'Actual': yts.flatten(), 'Predicted': y_pred.flatten()})
df
```

Out[67]:

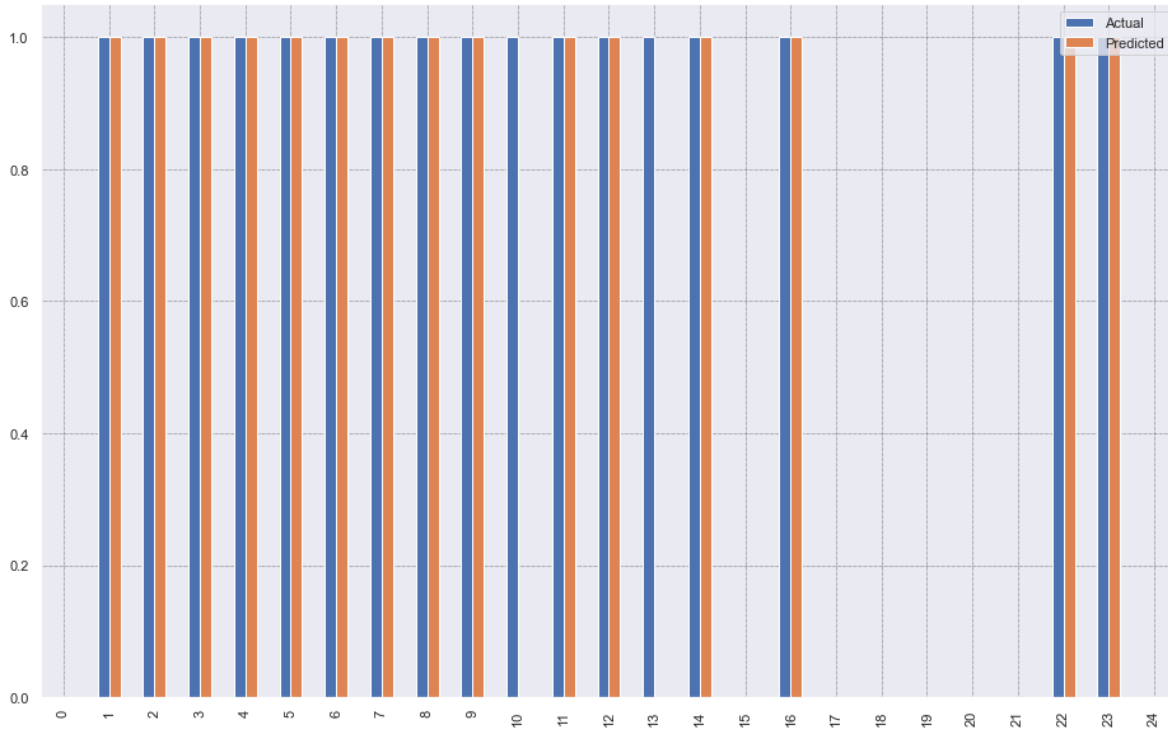
	Actual	Predicted
0	0	0
1	1	1
2	1	1
3	1	1
4	1	1
...	...	...
166	0	0
167	0	0
168	1	1
169	1	1
170	1	1

171 rows × 2 columns



In [68]:

```
df1=df.head(25)
df1.plot(kind='bar',figsize=(16,10))
plt.grid(which='major',linestyle='-',linewidth='0.5',color='green')
plt.grid(which='major',linestyle=':',linewidth='0.5',color='black')
plt.show()
```



In [69]:

```
print('Mean Absolute Error:',mean_absolute_error(yts,y_pred))
print('Mean Squared Error:',mean_squared_error(yts,y_pred))
print('Root Mean Squared Error:',np.sqrt(mean_squared_error(yts,y_pred)))
```

Mean Absolute Error: 0.04678362573099415  
Mean Squared Error: 0.04678362573099415  
Root Mean Squared Error: 0.21629522817435004

## 25 - MRINAAL PALIWAL

### Naive-Bayes

In [70]:

```
nb = BernoulliNB()
gnb = GaussianNB()
mnb = MultinomialNB()
```

In [71]:

```
nb.fit(xtr,ytr)
gnb.fit(xtr,ytr)
mnb.fit(xtr,ytr)
```

Out[71]:

MultinomialNB()

## BernoulliNB

In [72]:

```
ypred = nb.predict(xts)
```

In [73]:

```
accuracy_score(yts,ypred)
```

Out[73]:

0.631578947368421

In [74]:

```
print(classification_report(yts,ypred))
```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	63
1	0.63	1.00	0.77	108
accuracy			0.63	171
macro avg	0.32	0.50	0.39	171
weighted avg	0.40	0.63	0.49	171

In [75]:

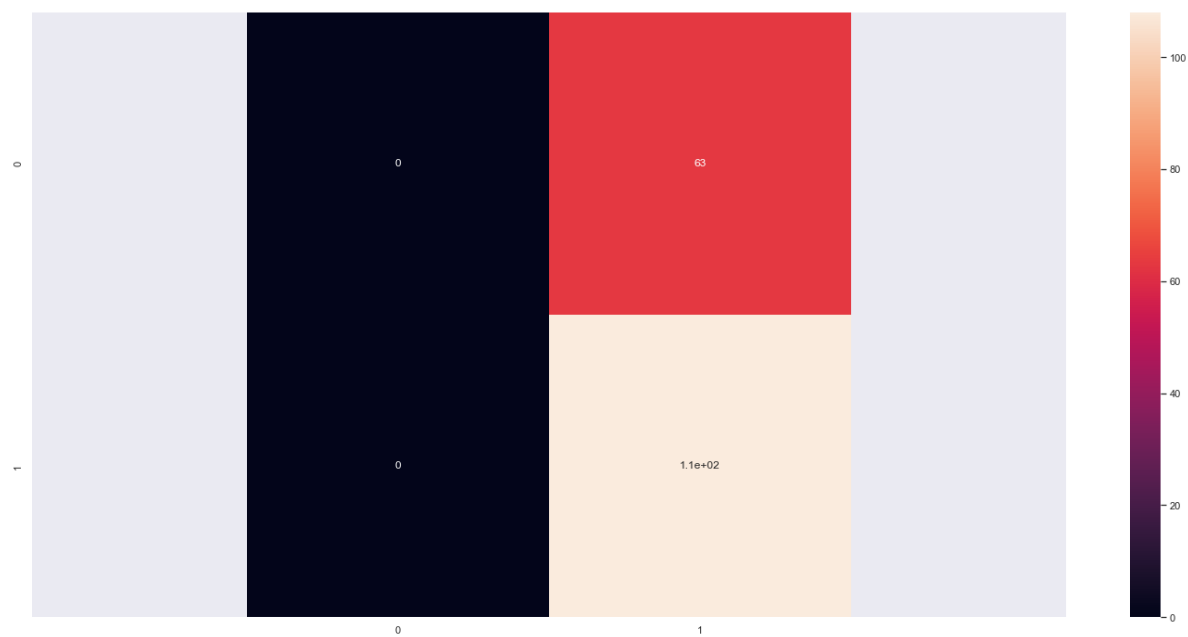
```
cf = confusion_matrix(yts,ypred)
cf
```

Out[75]:

```
array([[ 0, 63],
       [ 0, 108]], dtype=int64)
```

In [76]:

```
sns.heatmap (cf,annot=True)  
plt.axis('equal')  
plt.show()
```



In [77]:

```
nb.score(xts,yts)
```

Out[77]:

0.631578947368421

In [78]:

```
df=pd.DataFrame({'Actual': yts.flatten(),'Predicted': y_pred.flatten()})
df
```

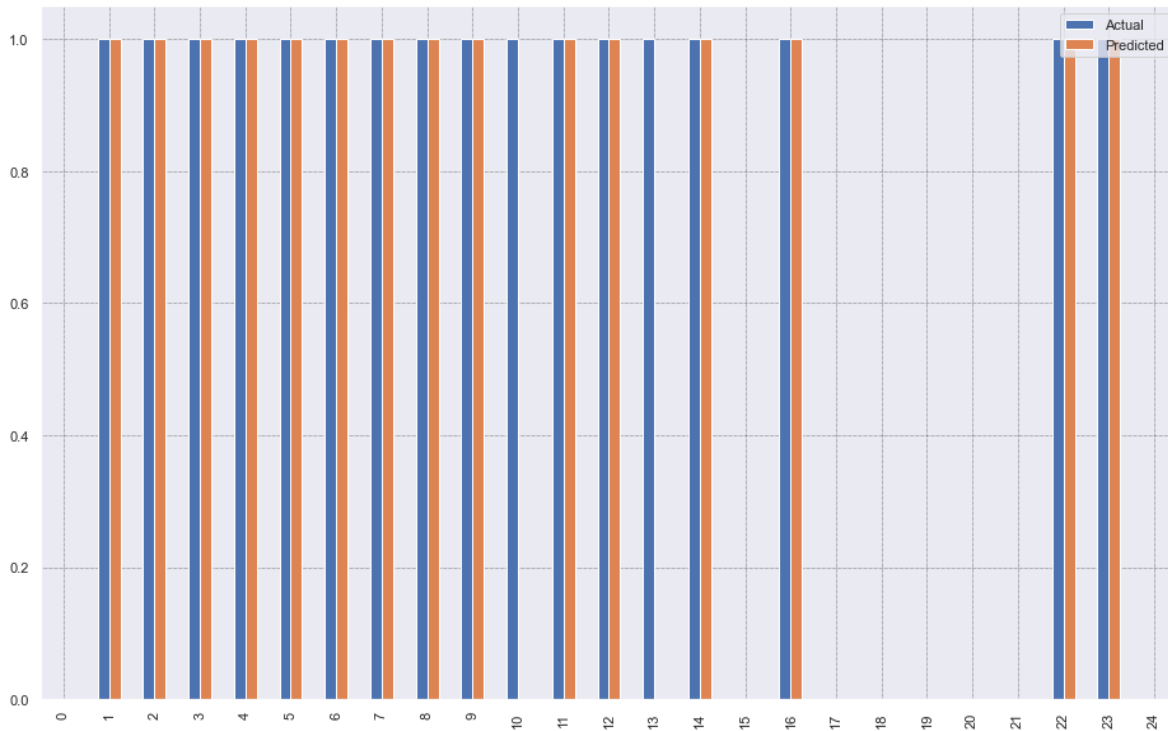
Out[78]:

	Actual	Predicted
0	0	0
1	1	1
2	1	1
3	1	1
4	1	1
...	...	...
166	0	0
167	0	0
168	1	1
169	1	1
170	1	1

171 rows × 2 columns

In [79]:

```
df1=df.head(25)
df1.plot(kind='bar',figsize=(16,10))
plt.grid(which='major',linestyle='-',linewidth='0.5',color='green')
plt.grid(which='major',linestyle=':',linewidth='0.5',color='black')
plt.show()
```



## GaussianNB

In [80]:

```
ypred = gnb.predict(xts)
```

In [81]:

```
accuracy_score(yts,ypred)
```

Out[81]:

0.9239766081871345

In [82]:

```
print(classification_report(yts,ypred))
```

	precision	recall	f1-score	support
0	0.89	0.90	0.90	63
1	0.94	0.94	0.94	108
accuracy			0.92	171
macro avg	0.92	0.92	0.92	171
weighted avg	0.92	0.92	0.92	171

In [83]:

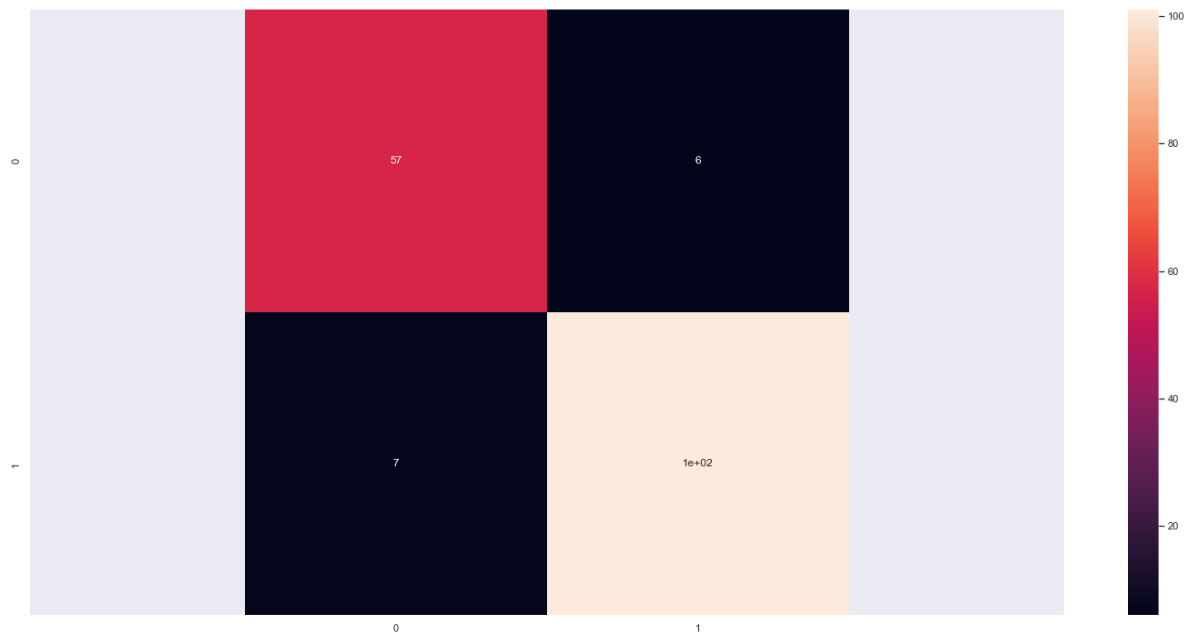
```
cf = confusion_matrix(yts,ypred)
cf
```

Out[83]:

```
array([[ 57,   6],
       [  7, 101]], dtype=int64)
```

In [84]:

```
sns.heatmap (cf,annot=True)
plt.axis('equal')
plt.show()
```



In [85]:

```
gnb.score(xts,yts)
```

Out[85]:

```
0.9239766081871345
```

In [86]:

```
df=pd.DataFrame({'Actual': yts.flatten(), 'Predicted': y_pred.flatten()})  
df
```

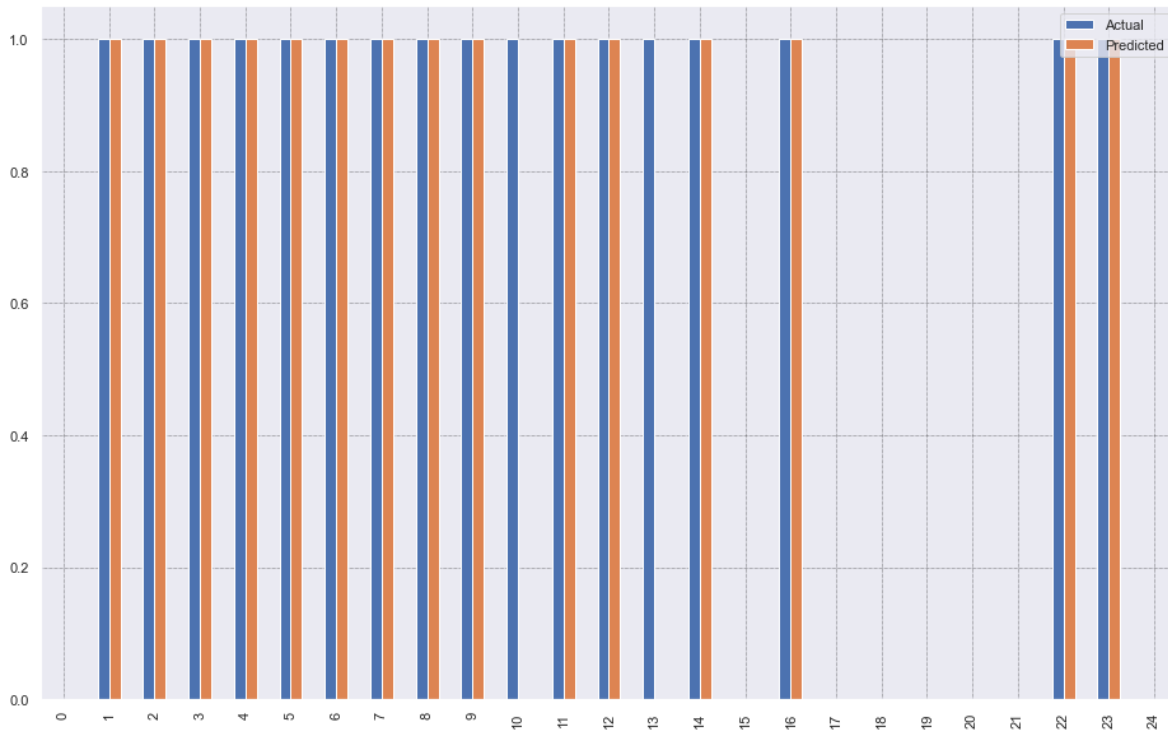
Out[86]:

	Actual	Predicted
0	0	0
1	1	1
2	1	1
3	1	1
4	1	1
...	...	...
166	0	0
167	0	0
168	1	1
169	1	1
170	1	1

171 rows × 2 columns

In [87]:

```
df1=df.head(25)
df1.plot(kind='bar',figsize=(16,10))
plt.grid(which='major',linestyle='-',linewidth='0.5',color='green')
plt.grid(which='major',linestyle=':',linewidth='0.5',color='black')
plt.show()
```



## MultinomialNB

In [88]:

```
ypred = mnbn.predict(xts)
```

In [89]:

```
accuracy_score(yts,ypred)
```

Out[89]:

0.9005847953216374



In [90]:

```
print(classification_report(yts,ypred))
```

	precision	recall	f1-score	support
0	0.96	0.76	0.85	63
1	0.88	0.98	0.93	108
accuracy			0.90	171
macro avg	0.92	0.87	0.89	171
weighted avg	0.91	0.90	0.90	171

In [91]:

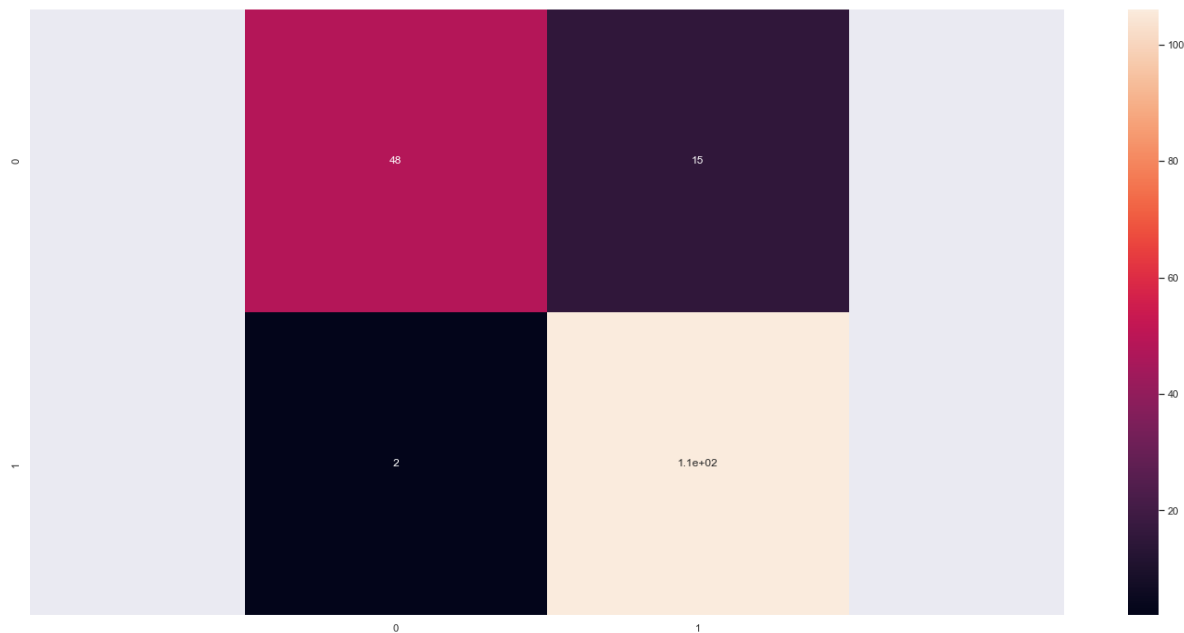
```
cf = confusion_matrix(yts,ypred)
cf
```

Out[91]:

```
array([[ 48,  15],
       [  2, 106]], dtype=int64)
```

In [92]:

```
sns.heatmap (cf,annot=True)
plt.axis('equal')
plt.show()
```



In [93]:

```
mnb.score(xts,yts)
```

Out[93]:

```
0.9005847953216374
```

In [94]:

```
df=pd.DataFrame({'Actual': yts.flatten(),'Predicted': y_pred.flatten()})
df
```

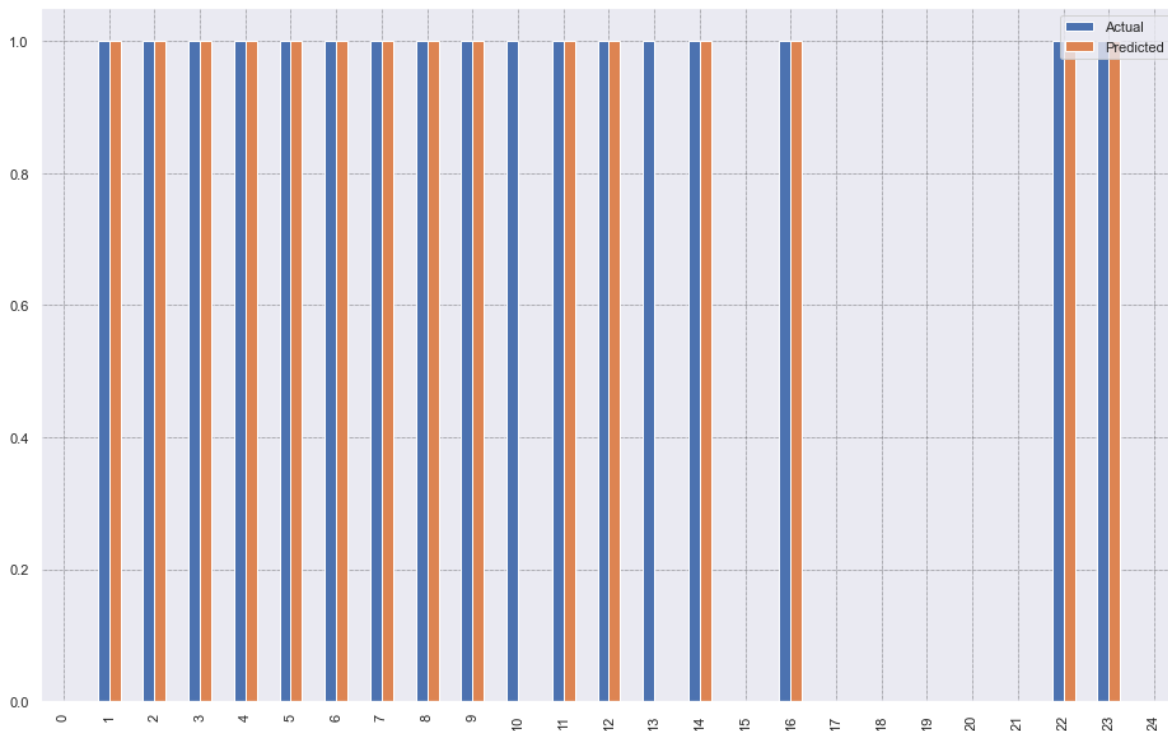
Out[94]:

	Actual	Predicted
0	0	0
1	1	1
2	1	1
3	1	1
4	1	1
...	...	...
166	0	0
167	0	0
168	1	1
169	1	1
170	1	1

171 rows × 2 columns

In [95]:

```
df1=df.head(25)
df1.plot(kind='bar',figsize=(16,10))
plt.grid(which='major',linestyle='-',linewidth='0.5',color='green')
plt.grid(which='major',linestyle=':',linewidth='0.5',color='black')
plt.show()
```



## Cross Validation

In [104]:

```
# k-fold cross validation technique
kf = KFold(n_splits=5, random_state=1,shuffle=True)
# stratified kfold cross validation technique
skf = StratifiedKFold(n_splits=5)
# LeaveOneOut cross validation technique
loocv = LeaveOneOut()
# shuffle split cross validation technique
shvc = ShuffleSplit()
```

In [105]:

```
# evaluating the data sets with kfold and DecisionTreeClassifier
dst = DecisionTreeClassifier()
scores = cross_val_score(dst,feat,classes,scoring='accuracy',cv=kf)
print('Accuracy using Decision Tree: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using Decision Tree: 93.147027%  
[0.94736842 0.92982456 0.90350877 0.93859649 0.9380531 ]

In [106]:

```
# evaluating the data sets with kfold and Naive Bayes Classifier
nb = GaussianNB()
scores = cross_val_score(nb,feat,classes,scoring='accuracy',cv=kf)
print('Accuracy using GaussianNB: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using GaussianNB: 93.851886%  
[0.94736842 0.93859649 0.9122807 0.93859649 0.95575221]

In [107]:

```
# evaluating the data sets with kfold and SVM
svm = SVC()
scores = cross_val_score(svm,feat,classes,scoring='accuracy',cv=kf)
print('Accuracy using SVC: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using SVC: 91.386431%  
[0.90350877 0.92982456 0.88596491 0.94736842 0.90265487]

In [108]:

```
# evaluating the data sets with kfold and KNN Classifier
knn = KNeighborsClassifier()
scores = cross_val_score(knn,feat,classes,scoring='accuracy',cv=kf)
print('Accuracy using KNN: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using KNN: 92.268281%  
[0.93859649 0.89473684 0.88596491 0.96491228 0.92920354]

In [109]:

```
# evaluating the data sets with Stratified KFold and DecisionTree
scores = cross_val_score(dst,feat,classes,scoring='accuracy',cv=skf)
print('Accuracy using Decision Tree: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using Decision Tree: 90.504580%  
[0.90350877 0.92105263 0.90350877 0.92105263 0.87610619]

In [110]:

```
# evaluating the data sets Stratified KFold and Naive Bayes Classifier
scores = cross_val_score(nb,feat,classes,scoring='accuracy',cv=skf)
print('Accuracy using GaussianNB: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using GaussianNB: 93.851886%  
[0.92105263 0.92105263 0.94736842 0.94736842 0.95575221]

In [111]:

```
# evaluating the data sets Stratified KFold and SVM
scores = cross_val_score(svm,feat,classes,scoring='accuracy',cv=skf)
print('Accuracy using SVM: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using SVM: 91.217202%

```
[0.85087719 0.89473684 0.92982456 0.94736842 0.9380531 ]
```

In [112]:

```
# evaluating the data sets Stratified KFold and KNN
scores = cross_val_score(knn,feat,classes,scoring='accuracy',cv=skf)
print('Accuracy using KNN: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using KNN: 92.794597%

```
[0.88596491 0.93859649 0.93859649 0.94736842 0.92920354]
```

In [113]:

```
# evaluating the data sets with LeaveOneOut and DecisionTree
scores = cross_val_score(dst,feat,classes,scoring='accuracy',cv=loo cv)
print('Accuracy using GaussianNB: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using GaussianNB: 92.267135%

[illegible]

In [114]:

```
# evaluating the data sets LeaveOneOut and Naive Bayes Classifier
scores = cross_val_score(nb,feat,classes,scoring='accuracy',cv=looocv)
print('Accuracy using GaussianNB: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using GaussianNB: 93.848858%

```
[1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 0. 1. 1. 0. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1.
 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 0. 1. 0. 1. 1. 1. 1.
 1. 1. 1. 0. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 0. 1. 0. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 0.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 0. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]
```

In [115]:

```
# evaluating the data sets LeaveOneOut and SVM
scores = cross_val_score(svm,feat,classes,scoring='accuracy',cv=looocv)
print('Accuracy using SVM: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using SVM: 91.212654%

[illegible]

In [116]:

```
# evaluating the data sets LeaveOneOut and KNN
scores = cross_val_score(knn,feat,classes,scoring='accuracy',cv=looocv)
print('Accuracy using KNN: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using KNN: 93.321617%

```
[1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 0. 1. 0. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 0. 0. 1. 1. 1.
 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1.
 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 0.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 0. 1. 1. 1. 1.
 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1.
 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 0.
 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 0. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]
```

In [117]:

```
# evaluating the data sets with ShuffleSplit and DecisionTree
scores = cross_val_score(dst,feat,classes,scoring='accuracy',cv=shvc)
print('Accuracy using GaussianNB: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using GaussianNB: 90.701754%

```
[0.94736842 0.84210526 0.92982456 0.96491228 0.89473684 0.89473684
 0.89473684 0.89473684 0.89473684 0.9122807 ]
```

In [118]:

```
# evaluating the data sets with ShuffleSplit and Naive Bayes Classifier
scores = cross_val_score(nb,feat,classes,scoring='accuracy',cv=shvc)
print('Accuracy using GaussianNB: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using GaussianNB: 95.087719%

```
[0.96491228 0.92982456 0.98245614 0.94736842 1.          0.94736842
 0.96491228 0.89473684 1.          0.87719298]
```



In [119]:

```
# evaluating the data sets with ShuffleSplit and SVM
scores = cross_val_score(svm,feat,classes,scoring='accuracy',cv=shvc)
print('Accuracy using GaussianNB: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using GaussianNB: 91.403509%

```
[0.94736842 0.87719298 0.94736842 0.85964912 0.9122807  0.87719298
 0.94736842 0.98245614 0.9122807  0.87719298]
```

In [120]:

```
# evaluating the data sets with ShuffleSplit and KNN
scores = cross_val_score(knn,feat,classes,scoring='accuracy',cv=shvc)
print('Accuracy using KNN: %2f%%'%(scores.mean()*100))
print(scores)
```

Accuracy using KNN: 93.859649%

```
[0.98245614 1.          0.87719298 0.94736842 0.98245614 0.9122807
 0.9122807  0.92982456 0.92982456 0.9122807 ]
```

## 41 - PALLAVI SAWANT

### PreProcessing

#### Define Target Data

Diagnosis column will be our target data  
If the cancer is Benign, it will be 0  
If the cancer is Malignant, it will be 1

In [21]:

```
#split the dataset into training and testing sets  
features = df.drop(['diagnosis'],axis=1).values  
classes=df['diagnosis'].values
```

In [22]:

```
feat_train, feat_test, class_train, class_test = train_test_split(features, classes,  
                                                                    test_size=0.2, random_state=9)
```

In [23]:

```
print('features train shape: ', feat_train.shape)  
print('classes train shape: ', class_train.shape)  
print('features test shape: ', feat_test.shape)  
print('classes test shape: ', class_test.shape)
```

```
features train shape: (455, 30)  
classes train shape: (455,)  
features test shape: (114, 30)  
classes test shape: (114,)
```

#### Decision Tree Classifier

Criterion=Gini

In [24]:

```
#Training
dectree=DecisionTreeClassifier(criterion='gini')

dectree.fit(feat_train,class_train)
```

Out[24]:

DecisionTreeClassifier()

In [25]:

```
#predict target values
pred=dectree.predict(feat_test)
print(pred)
```

```
['M' 'B' 'M' 'B' 'B' 'M' 'M' 'M' 'M' 'M' 'M' 'B' 'B' 'B' 'M' 'M' 'B' 'B'
'B' 'B' 'M' 'B' 'B' 'M' 'B' 'B' 'B' 'B' 'B' 'M' 'M' 'B' 'M' 'M' 'M' 'M'
'M' 'B' 'M' 'B' 'M' 'B' 'B' 'B' 'B' 'B' 'M' 'M' 'B' 'M' 'M' 'B' 'B' 'B'
'M' 'B' 'B' 'M' 'B' 'B' 'M' 'M' 'M' 'B' 'B' 'B' 'M' 'B' 'M' 'B' 'B' 'B'
'B' 'B' 'M' 'B' 'B' 'B' 'B' 'B' 'B' 'M' 'B' 'M' 'B' 'B' 'B' 'B' 'M' 'B'
'B' 'B' 'B' 'B' 'B' 'M' 'B' 'B' 'B' 'M' 'B' 'M' 'B' 'B' 'B' 'B' 'M' 'B'
'B' 'B' 'B' 'B' 'B' 'M']
```

In [26]:

```
#confusion matrix and accuracy
print("Accuracy",accuracy_score(class_test,pred))
print("Classification Report\n",classification_report(class_test,pred))
print("Confusion Matrix\n",confusion_matrix(class_test,pred))
```

Accuracy 0.956140350877193

Classification Report

	precision	recall	f1-score	support
B	0.97	0.96	0.97	74
M	0.93	0.95	0.94	40
accuracy			0.96	114
macro avg	0.95	0.95	0.95	114
weighted avg	0.96	0.96	0.96	114

Confusion Matrix

```
[[71  3]
 [ 2 38]]
```

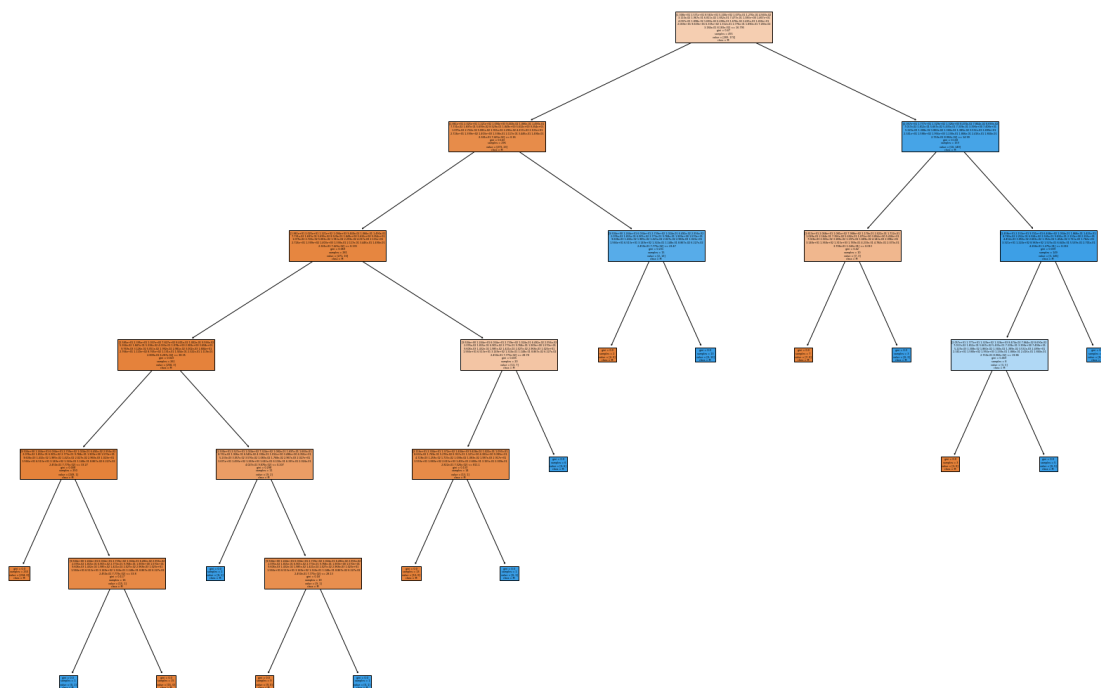
In [27]:

```
#Testing
pred=dectree.predict(feat_test)
print("Accuracy:",metrics.accuracy_score(class_test,pred))
```

Accuracy: 0.956140350877193

In [28]:

```
from sklearn import tree
fig=plt.figure(figsize=(30,20))
_=tree.plot_tree(dectree,feature_names=features,class_names=classes,filled=True)
```



In [29]:

```
text_representation=tree.export_text(dectree)
print(text_representation)
```

```
|--- feature_20 <= 16.80
|   |--- feature_27 <= 0.16
|   |   |--- feature_27 <= 0.13
|   |   |   |--- feature_13 <= 38.35
|   |   |   |   |--- feature_21 <= 33.27
|   |   |   |   |   |--- class: B
|   |   |   |   |--- feature_21 > 33.27
|   |   |   |   |   |--- feature_21 <= 33.80
|   |   |   |   |   |   |--- class: M
|   |   |   |   |   |--- feature_21 > 33.80
|   |   |   |   |   |   |--- class: B
|   |   |   |--- feature_13 > 38.35
|   |   |   |   |--- feature_28 <= 0.21
|   |   |   |   |   |--- class: M
|   |   |   |   |--- feature_28 > 0.21
|   |   |   |   |   |--- feature_21 <= 28.13
|   |   |   |   |   |   |--- class: B
|   |   |   |   |   |--- feature_21 > 28.13
|   |   |   |   |   |   |--- class: M
|   |   |   |--- feature_27 > 0.13
|   |   |   |   |--- feature_21 <= 28.78
|   |   |   |   |   |--- feature_23 <= 811.10
|   |   |   |   |   |   |--- class: B
|   |   |   |   |   |--- feature_23 > 811.10
|   |   |   |   |   |   |--- class: M
|   |   |   |   |--- feature_21 > 28.78
|   |   |   |   |   |--- class: M
|   |   |--- feature_27 > 0.16
|   |   |   |--- feature_21 <= 23.47
|   |   |   |   |--- class: B
|   |   |   |--- feature_21 > 23.47
|   |   |   |   |--- class: M
|--- feature_20 > 16.80
|   |--- feature_1 <= 14.99
|   |   |--- feature_17 <= 0.01
|   |   |   |--- class: B
|   |   |--- feature_17 > 0.01
|   |   |   |--- class: M
|   |--- feature_1 > 14.99
|   |   |--- feature_26 <= 0.22
|   |   |   |--- feature_1 <= 19.86
|   |   |   |   |--- class: B
|   |   |   |--- feature_1 > 19.86
|   |   |   |   |--- class: M
|   |   |--- feature_26 > 0.22
|   |   |   |--- class: M
```

Criterion=Entropy

In [30]:

```
#Training
dectree=DecisionTreeClassifier(criterion='entropy')
dectree.fit(feet_train,class_train)
```

Out[30]:

```
DecisionTreeClassifier(criterion='entropy')
```

In [31]:

```
#confusion matrix and accuracy
pred=dectree.predict(feet_test)
print(pred)
print("Accuracy",accuracy_score(class_test,pred))
print("Classification Report\n",classification_report(class_test,pred))
print("Confusion Matrix\n",confusion_matrix(class_test,pred))
```

```
['B' 'B' 'M' 'B' 'B' 'M' 'M' 'M' 'M' 'M' 'M' 'B' 'B' 'B' 'M' 'M' 'B' 'B'
 'B' 'B' 'M' 'B' 'B' 'M' 'M' 'B' 'B' 'B' 'B' 'M' 'M' 'B' 'M' 'B' 'M' 'M'
 'M' 'B' 'M' 'B' 'B' 'B' 'B' 'B' 'B' 'B' 'M' 'M' 'B' 'M' 'B' 'B' 'B' 'B'
 'M' 'B' 'B' 'M' 'B' 'B' 'M' 'M' 'M' 'B' 'B' 'B' 'M' 'B' 'M' 'B' 'B' 'B'
 'B' 'B' 'M' 'B' 'B' 'B' 'B' 'B' 'B' 'B' 'B' 'M' 'B' 'B' 'B' 'B' 'M' 'B'
 'B' 'B' 'B' 'B' 'B' 'M' 'B' 'B' 'B' 'M' 'B' 'M' 'B' 'B' 'B' 'B' 'M' 'B'
 'B' 'B' 'B' 'B' 'B' 'M']
```

Accuracy 0.9210526315789473

Classification Report

	precision	recall	f1-score	support
B	0.92	0.96	0.94	74
M	0.92	0.85	0.88	40
accuracy			0.92	114
macro avg	0.92	0.90	0.91	114
weighted avg	0.92	0.92	0.92	114

Confusion Matrix

```
[[71 3]
 [ 6 34]]
```

In [32]:

```
text_representation = tree.export_text(dectree)
print(text_representation)
```

```
|--- feature_22 <= 105.95
|   |--- feature_27 <= 0.13
|   |   |--- feature_10 <= 0.64
|   |   |   |--- feature_21 <= 33.27
|   |   |   |   |--- class: B
|   |   |   |--- feature_21 > 33.27
|   |   |   |   |--- feature_21 <= 33.80
|   |   |   |   |   |--- class: M
|   |   |   |   |--- feature_21 > 33.80
|   |   |   |   |   |--- class: B
|   |   |--- feature_10 > 0.64
|   |   |   |--- feature_24 <= 0.11
|   |   |   |   |--- class: B
|   |   |   |--- feature_24 > 0.11
|   |   |   |   |--- class: M
|   |--- feature_27 > 0.13
|   |   |--- feature_1 <= 20.30
|   |   |   |--- feature_28 <= 0.36
|   |   |   |   |--- class: B
|   |   |   |--- feature_28 > 0.36
|   |   |   |   |--- feature_22 <= 78.63
|   |   |   |   |   |--- class: B
|   |   |   |   |--- feature_22 > 78.63
|   |   |   |   |   |--- class: M
|   |   |--- feature_1 > 20.30
|   |   |   |--- class: M
|--- feature_22 > 105.95
|   |--- feature_22 <= 117.45
|   |   |--- feature_21 <= 27.46
|   |   |   |--- feature_24 <= 0.14
|   |   |   |   |--- class: B
|   |   |   |--- feature_24 > 0.14
|   |   |   |   |--- feature_24 <= 0.15
|   |   |   |   |   |--- class: M
|   |   |   |   |--- feature_24 > 0.15
|   |   |   |   |   |--- class: B
|   |   |--- feature_21 > 27.46
|   |   |   |--- feature_4 <= 0.09
|   |   |   |   |--- feature_9 <= 0.06
|   |   |   |   |   |--- class: M
|   |   |   |   |--- feature_9 > 0.06
|   |   |   |   |   |--- class: B
|   |   |   |--- feature_4 > 0.09
|   |   |   |   |--- class: M
|   |--- feature_22 > 117.45
|   |   |--- feature_27 <= 0.09
|   |   |   |--- feature_29 <= 0.06
|   |   |   |   |--- class: B
|   |   |   |--- feature_29 > 0.06
|   |   |   |   |--- class: M
|   |--- feature_27 > 0.09
|   |   |--- class: M
```

# KNN Classifier

In [33]:

```
feat_train, feat_test, class_train, class_test = train_test_split(features,
                                                                    classes, test_size=0.2, random_state=9)
```

In [34]:

```
knn=KNeighborsClassifier(n_neighbors=4)
knn.fit(feat_train,class_train)
```

Out[34]:

```
KNeighborsClassifier(n_neighbors=4)
```

In [35]:

```
pred=knn.predict(feat_test)
print("Accuracy:",metrics.accuracy_score(class_test,pred))
```

Accuracy: 0.9210526315789473

In [36]:

```
neighbors=np.arange(1,9)
train_accuracy=np.empty(len(neighbors))
test_accuracy=np.empty(len(neighbors))
for i,k in enumerate(neighbors):
    #setup as knn vlassifier with k neighbors
    knn1=KNeighborsClassifier(n_neighbors=k)
    #fit the model
    knn1.fit(feat_train,class_train)
    pred=knn1.predict(feat_test)
    #compute accuracy on the training set
    train_accuracy[i]=knn1.score(feat_train,class_train)
    #compute accuracy on the test set
    test_accuracy[i]=knn1.score(feat_test,class_test)
    print("Accuracy:",i,metrics.accuracy_score(class_test,pred))
print("train_accuracy\n",train_accuracy)
print("test_accuracy\n",test_accuracy)
```

Accuracy: 0 0.9122807017543859

Accuracy: 1 0.9298245614035088

Accuracy: 2 0.9385964912280702

Accuracy: 3 0.9210526315789473

Accuracy: 4 0.9385964912280702

Accuracy: 5 0.9298245614035088

Accuracy: 6 0.9298245614035088

Accuracy: 7 0.9122807017543859

train\_accuracy

```
[1.          0.94725275  0.95824176  0.94285714  0.94725275  0.94505495
 0.94065934  0.93406593]
```

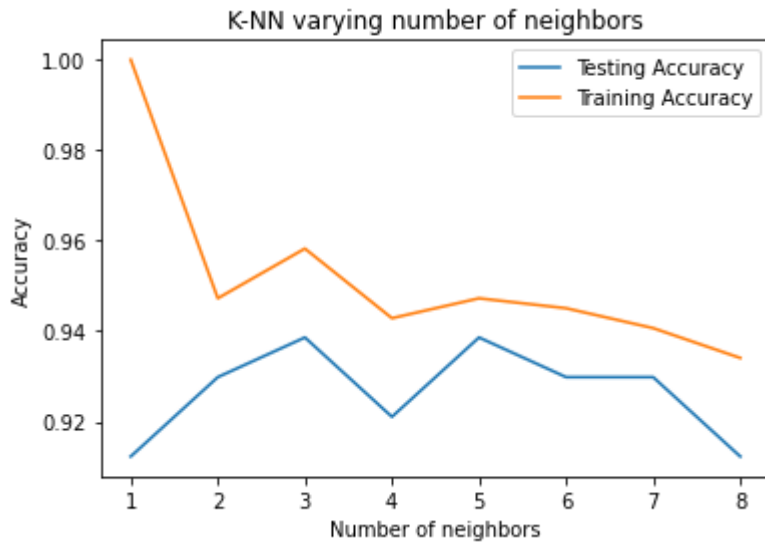
test\_accuracy

```
[0.9122807  0.92982456 0.93859649 0.92105263 0.93859649 0.92982456
 0.92982456 0.9122807 ]
```



In [37]:

```
plt.title('K-NN varying number of neighbors')
plt.plot(neighbors, test_accuracy, label='Testing Accuracy')
plt.plot(neighbors, train_accuracy, label='Training Accuracy')
plt.legend()
plt.xlabel('Number of neighbors')
plt.ylabel('Accuracy')
plt.show()
```



Conclusion: From above graph we see training accuracy is more than that of testing accuracy

## 10 - SHREYA JAGADALE

### Support Vector Machine

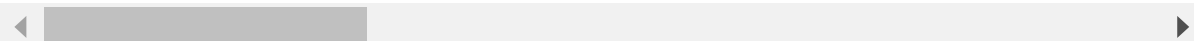
In [38]:

```
df.head()
```

Out[38]:

	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	com
0	M	17.99	10.38	122.80	1001.0	0.11840	
1	M	20.57	17.77	132.90	1326.0	0.08474	
2	M	19.69	21.25	130.00	1203.0	0.10960	
3	M	11.42	20.38	77.58	386.1	0.14250	
4	M	20.29	14.34	135.10	1297.0	0.10030	

5 rows × 31 columns



In [39]:

```
df['diagnosis'].unique()
```

Out[39]:

```
array(['M', 'B'], dtype=object)
```

In [40]:

```
df['diagnosis'] = df['diagnosis'].map({'M':1, 'B':0})
```

In [41]:

```
classess=df.diagnosis  
classess
```

Out[41]:

```
0      1  
1      1  
2      1  
3      1  
4      1  
..  
564    1  
565    1  
566    1  
567    1  
568    0
```

Name: diagnosis, Length: 569, dtype: int64

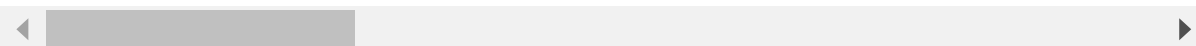
In [42]:

```
featuress = df.drop(['diagnosis'],axis=1)
featuress
```

Out[42]:

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_
0	17.99	10.38	122.80	1001.0	0.11840	0.
1	20.57	17.77	132.90	1326.0	0.08474	0.
2	19.69	21.25	130.00	1203.0	0.10960	0.
3	11.42	20.38	77.58	386.1	0.14250	0.
4	20.29	14.34	135.10	1297.0	0.10030	0.
...	...	...	...	...	...	...
564	21.56	22.39	142.00	1479.0	0.11100	0.
565	20.13	28.25	131.20	1261.0	0.09780	0.
566	16.60	28.08	108.30	858.1	0.08455	0.
567	20.60	29.33	140.10	1265.0	0.11780	0.
568	7.76	24.54	47.92	181.0	0.05263	0.

569 rows × 30 columns



In [43]:

```
from sklearn import preprocessing
#get col names
names = featuress.columns
#create scaler object
scaler = preprocessing.StandardScaler()
#fit data on the scaler object
scaled_df = scaler.fit_transform(featuress)
featuress = pd.DataFrame(scaled_df,columns=names)
```

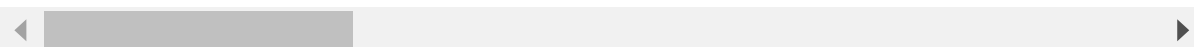
In [44]:

```
featuress
```

Out[44]:

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_
0	1.097064	-2.073335	1.269934	0.984375	1.568466	3.2
1	1.829821	-0.353632	1.685955	1.908708	-0.826962	-0.4
2	1.579888	0.456187	1.566503	1.558884	0.942210	1.0
3	-0.768909	0.253732	-0.592687	-0.764464	3.283553	3.4
4	1.750297	-1.151816	1.776573	1.826229	0.280372	0.5
...	...	...	...	...	...	...
564	2.110995	0.721473	2.060786	2.343856	1.041842	0.2
565	1.704854	2.085134	1.615931	1.723842	0.102458	-0.0
566	0.702284	2.045574	0.672676	0.577953	-0.840484	-0.0
567	1.838341	2.336457	1.982524	1.735218	1.525767	3.2
568	-1.808401	1.221792	-1.814389	-1.347789	-3.112085	-1.1

569 rows × 30 columns



In [45]:

```
feat_train, feat_test, class_train, class_test = train_test_split(featuress,
                                                                    classess, train_size=0.9, random_state=100)
```

In [46]:

```
svclassifier=SVC(kernel='linear')
svclassifier.fit(feat_train,class_train)
```

Out[46]:

```
SVC(kernel='linear')
```

In [47]:

```
print('features train shape: ', feat_train.shape)
print('classes train shape: ', class_train.shape)
print('features test shape: ', feat_test.shape)
print('classes test shape: ', class_test.shape)
```

```
features train shape: (512, 30)
classes train shape: (512,)
features test shape: (57, 30)
classes test shape: (57,)
```

In [48]:

```
pred=svclassifier.predict(feat_test)
```

In [49]:

```
accuracy_score(class_test,pred)
```

Out[49]:

0.9473684210526315

In [50]:

```
print(classification_report(class_test,pred))
```

	precision	recall	f1-score	support
0	0.92	1.00	0.96	35
1	1.00	0.86	0.93	22
accuracy			0.95	57
macro avg	0.96	0.93	0.94	57
weighted avg	0.95	0.95	0.95	57

In [51]:

```
cf=confusion_matrix(class_test,pred)
cf
```

Out[51]:

```
array([[35,  0],
       [ 3, 19]], dtype=int64)
```

In [52]:

```
kernels = ['linear','rbf','poly']
for kernel in kernels:
    sv = SVC(kernel=kernel).fit(feat_train,class_train)
    pred=sv.predict(feat_test)
    print("Accuracy:("+kernel+")", accuracy_score(class_test,pred))
```

```
Accuracy:(linear) 0.9473684210526315
Accuracy:(rbf) 0.9649122807017544
Accuracy:(poly) 0.8947368421052632
```

In [53]:

```
gammas = [0.1,1,10,100]
for gamma in gammas:
    sv = SVC(kernel='rbf', gamma=gamma).fit(feet_train,class_train)
    pred=sv.predict(feet_test)
    print("Accuracy:", gamma , ""), accuracy_score(class_test,pred))
```

```
Accuracy:( 0.1 ) 0.9473684210526315
Accuracy:( 1 ) 0.6140350877192983
Accuracy:( 10 ) 0.6140350877192983
Accuracy:( 100 ) 0.6140350877192983
```

In [54]:

```
degrees = [0,1,2,3,4,5,20]
for degree in degrees:
    sv = SVC(kernel='poly', degree=degree).fit(feet_train,class_train)
    pred=sv.predict(feet_test)
    print("Accuracy:", degree , "):", accuracy_score(class_test,pred))
```

```
Accuracy:( 0 ): 0.6140350877192983
Accuracy:( 1 ): 0.9473684210526315
Accuracy:( 2 ): 0.8771929824561403
Accuracy:( 3 ): 0.8947368421052632
Accuracy:( 4 ): 0.8596491228070176
Accuracy:( 5 ): 0.8596491228070176
Accuracy:( 20 ): 0.7543859649122807
```

## Principal Component Analysis

In [57]:

```
pca = decomposition.PCA(n_components=3)
pca.fit(featuress)
X1 = pca.transform(featuress)
```

In [58]:

```
print(cancer.data.shape)
print(featuress.shape)
print(X1.shape)
```

```
(569, 30)
(569, 30)
(569, 3)
```

In [59]:

```
from sklearn.model_selection import train_test_split
(train_feat,test_feat,train_classes,test_classes)= train_test_split(features,
                                                                    classes,train_size=0.5,random_state=6)
dectree = DecisionTreeClassifier()
dectree.fit(train_feat,train_classes)
```

Out[59]:

DecisionTreeClassifier()

In [60]:

```
from sklearn import metrics
pred=dectree.predict(test_feat)
print("Accuracy:",metrics.accuracy_score(test_classes,pred))
```

Accuracy: 0.9157894736842105

In [61]:

```
from sklearn.model_selection import train_test_split
(train_feat,test_feat,train_classes,test_classes)= train_test_split(X1,classes,
                                                                    train_size=0.5,random_state=6)
dectree = DecisionTreeClassifier()
dectree.fit(train_feat,train_classes)
```

Out[61]:

DecisionTreeClassifier()

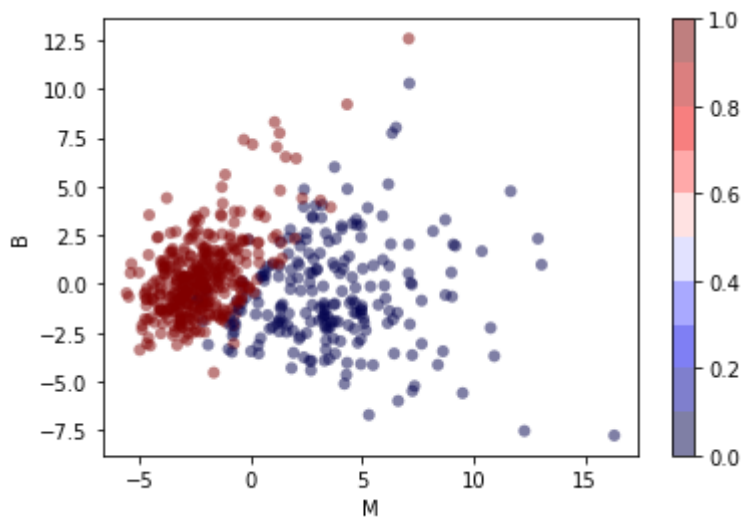
In [62]:

```
from sklearn import metrics
pred=dectree.predict(test_feat)
print("Accuracy:",metrics.accuracy_score(test_classes,pred))
```

Accuracy: 0.9298245614035088

In [63]:

```
plt.scatter(X1[:,0],X1[:,1],c=cancer.target,edgecolor='none',
            alpha=0.5,cmap=plt.cm.get_cmap('seismic',10))
plt.xlabel('M')
plt.ylabel('B')
plt.colorbar();
```



## Select K Percentile and K Best

In [102]:

```
X_train,X_test,y_train,y_test = train_test_split(featuress,
            classess,test_size=0.9,random_state=100)
```

In [103]:

```
select = SelectPercentile(percentile=20)
select.fit(X_train,y_train)
```

Out[103]:

```
SelectPercentile(percentile=20)
```

In [104]:

```
X_train_selected = select.transform(X_train)
print("X_train.shape: {}".format(X_train.shape))
print("X_train_selected.shape: {}".format(X_train_selected.shape))
```

```
X_train.shape: (56, 30)
X_train_selected.shape: (56, 6)
```



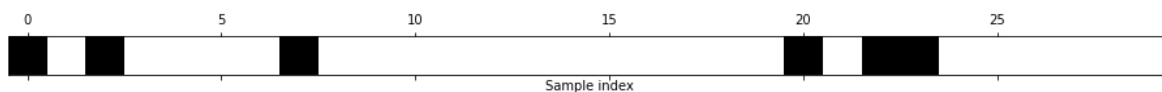
In [105]:

```
mask = select.get_support()
print(mask)
plt.matshow(mask.reshape(1,-1),cmap='gray_r')
plt.xlabel("Sample index")
plt.yticks(())
```

```
[ True False  True False False False False  True False False False False
 False False False False False False False False  True False  True  True
 False False False False False False]
```

Out[105]:

([], [])



In [106]:

```
from sklearn.tree import DecisionTreeClassifier
X_test_selected = select.transform(X_test)
lr = DecisionTreeClassifier()
lr.fit(X_train,y_train)
print("Score with all features: {:.3f}".format(lr.score(X_test,y_test)))
lr.fit(X_train_selected,y_train)
print("Score with all features: {:.3f}".format(lr.score(X_test_selected,y_test)))
```

Score with all features: 0.854

Score with all features: 0.903

In [107]:

```
select=SelectKBest(k=2)
select.fit(X_train,y_train)
```

Out[107]:

SelectKBest(k=2)

In [108]:

```
X_train_selected = select.transform(X_train)
print("X_train.shape: {}".format(X_train.shape))
print("X_train_selected.shape: {}".format(X_train_selected.shape))
```

X\_train.shape: (56, 30)

X\_train\_selected.shape: (56, 2)

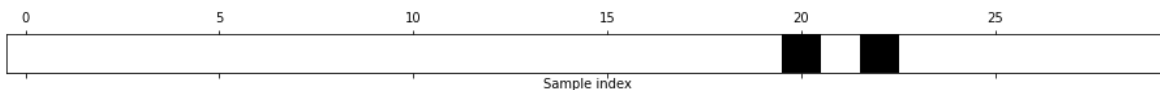
In [109]:

```
mask = select.get_support()
print(mask)
plt.matshow(mask.reshape(1,-1),cmap='gray_r')
plt.xlabel("Sample index")
plt.yticks(())
```

```
[False False False False False False False False False False False False
 False False False False False False False False  True False  True False
 False False False False False False]
```

Out[109]:

([], [])



## Feature Selection-Model Based

In [86]:

```
from sklearn.feature_selection import SelectFromModel
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier
select = SelectFromModel(RandomForestClassifier(n_estimators=10,
                                                random_state=100),threshold="median")
#select = SelectFromModel(SVC(kernel='linear'))
```

In [87]:

```
select.fit(X_train,y_train)
X_train_l1 = select.transform(X_train)
print("X_train.shape: {}".format(X_train.shape))
print("X_train_l1.shape: {}".format(X_train_l1.shape))
```

X\_train.shape: (56, 30)

X\_train\_l1.shape: (56, 15)

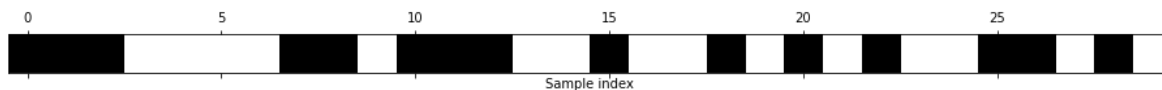
In [88]:

```
mask = select.get_support()
print(mask)
plt.matshow(mask.reshape(1,-1),cmap='gray_r')
plt.xlabel("Sample index")
plt.yticks(())
```

```
[ True  True  True False False False False  True  True False  True  True
  True False False  True False False  True False  True False  True False
 False  True  True False  True False]
```

Out[88]:

([], [])



In [89]:

```
X_test_l1 = select.transform(X_test)
score = SVC().fit(X_train,y_train).score(X_test,y_test)
print("Test Score: {:.3f}".format(score))
score = SVC().fit(X_train_l1,y_train).score(X_test_l1,y_test)
print("Test Score: {:.3f}".format(score))
```

Test Score: 0.942

Test Score: 0.930

## Iterative Feature Selection

In [120]:

```
from sklearn.feature_selection import RFE
select = RFE(RandomForestClassifier(n_estimators=20,
                                   random_state=100),n_features_to_select=30)
select.fit(X_train,y_train)
```

Out[120]:

```
RFE(estimator=RandomForestClassifier(n_estimators=20, random_state=100),
    n_features_to_select=30)
```

```
mask = select.get_support() print(mask) plt.matshow(mask.reshape(1,-1),cmap='gray_r') plt.xlabel("Sample
index") plt.yticks(())
```

In [121]:

```
from sklearn.linear_model import LogisticRegression
X_train_rfe =select.transform(X_train)
X_test_rfe = select.transform(X_test)
score = LogisticRegression().fit(X_train_rfe,y_train).score(X_test_rfe,y_test)
print("Test score: {:.3f}".format(score))
print("Test score: {:.3f}".format(select.score(X_test,y_test)))
```

Test score: 0.975

Test score: 0.932

## Top 5 Accuracy:

Random Forest=0.97 

SVM kernel(rbf)= 0.96

Multivariate Logistic Regression=0.95

Decision Tree Classifier entropy=0.95

KNN shvc=0.93