

```
In [1]:
```

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

```
In [2]:
```

```
df=pd.read_csv('Test2_dataset.csv')
```

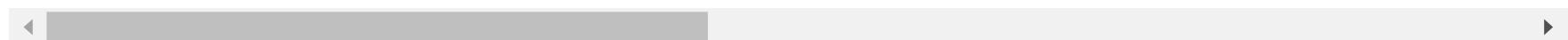
```
In [3]:
```

```
df.head()
```

```
Out[3]:
```

	Pid	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_nr
0	8670	M	15.46	19.48	101.70	748.9	0.10920	0.12230	0.14660	0.08
1	8913	B	12.89	13.12	81.89	515.9	0.06955	0.03729	0.02260	0.01
2	8915	B	14.96	19.10	97.03	687.3	0.08992	0.09823	0.05940	0.04
3	9047	B	12.94	16.17	83.18	507.6	0.09879	0.08836	0.03296	0.02
4	85715	M	13.17	18.66	85.98	534.6	0.11580	0.12310	0.12260	0.07

5 rows × 32 columns



EDA

```
In [4]:
```

```
#Understanding Data
df.describe()
```

```
Out[4]:
```

	Pid	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_nr
count	5.690000e+02	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	566.000000	566.0
mean	3.037183e+07	14.127292	19.289649	91.969033	654.889104	0.096360	0.104341	0.089270	0.0
std	1.250206e+08	3.524049	4.301036	24.298981	351.914129	0.014064	0.052813	0.079667	0.0

	Pid	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	co points_
min	8.670000e+03	6.981000	9.710000	43.790000	143.500000	0.052630	0.019380	0.000000	0.0
25%	8.692180e+05	11.700000	16.170000	75.170000	420.300000	0.086370	0.064920	0.029680	0.0
50%	9.060240e+05	13.370000	18.840000	86.240000	551.100000	0.095870	0.092630	0.061680	0.0
75%	8.813129e+06	15.780000	21.800000	104.100000	782.700000	0.105300	0.130400	0.131600	0.0
max	9.113205e+08	28.110000	39.280000	188.500000	2501.000000	0.163400	0.345400	0.426800	0.2

8 rows × 31 columns

◀	▶
---	---

In [5]:

`df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 32 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Pid              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   566 non-null    float64 
 9   concave_points_mean 566 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    566 non-null    float64 
 19  concave_points_se 567 non-null    float64 
 20  symmetry_se     569 non-null    float64
```

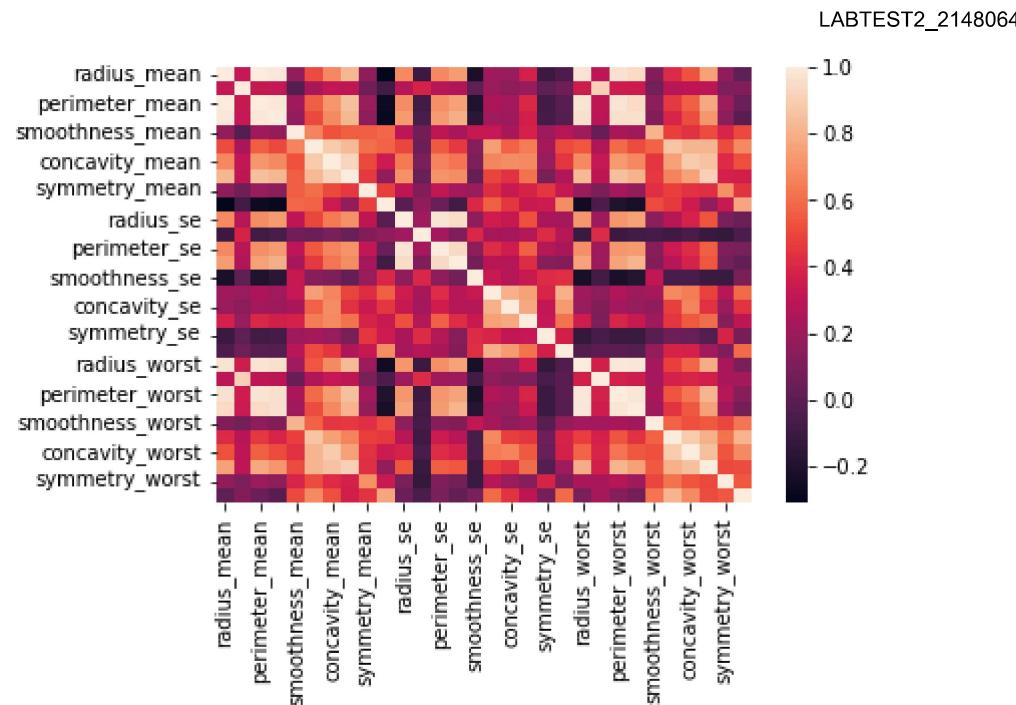
```
21 fractal_dimension_se      569 non-null   float64
22 radius_worst              569 non-null   float64
23 texture_worst              546 non-null   float64
24 perimeter_worst             569 non-null   float64
25 area_worst                  569 non-null   float64
26 smoothness_worst             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
dtypes: float64(30), int64(1), object(1)
memory usage: 142.4+ KB
```

```
In [6]: #Deleting unimportant columns
df.drop('Pid',axis=1,inplace=True)

#Deleting since unique and no impact on data
```

```
In [7]: #Relationship b/w variables
sns.heatmap(df.corr())
```

```
Out[7]: <AxesSubplot:>
```



In [8]:

```
#Univariate Analysis

l=list(df.loc[:, df.dtypes != object].columns)
len(l)

# plot all the columns present in list l together using subplot of dimension (2,3).

c=0
plt.figure(figsize=(15,10))
for i in l:
    c=c+1
    plt.subplot(6,7, c)
    plt.title(i)
    sns.distplot(df[i])
plt.show()
```

C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

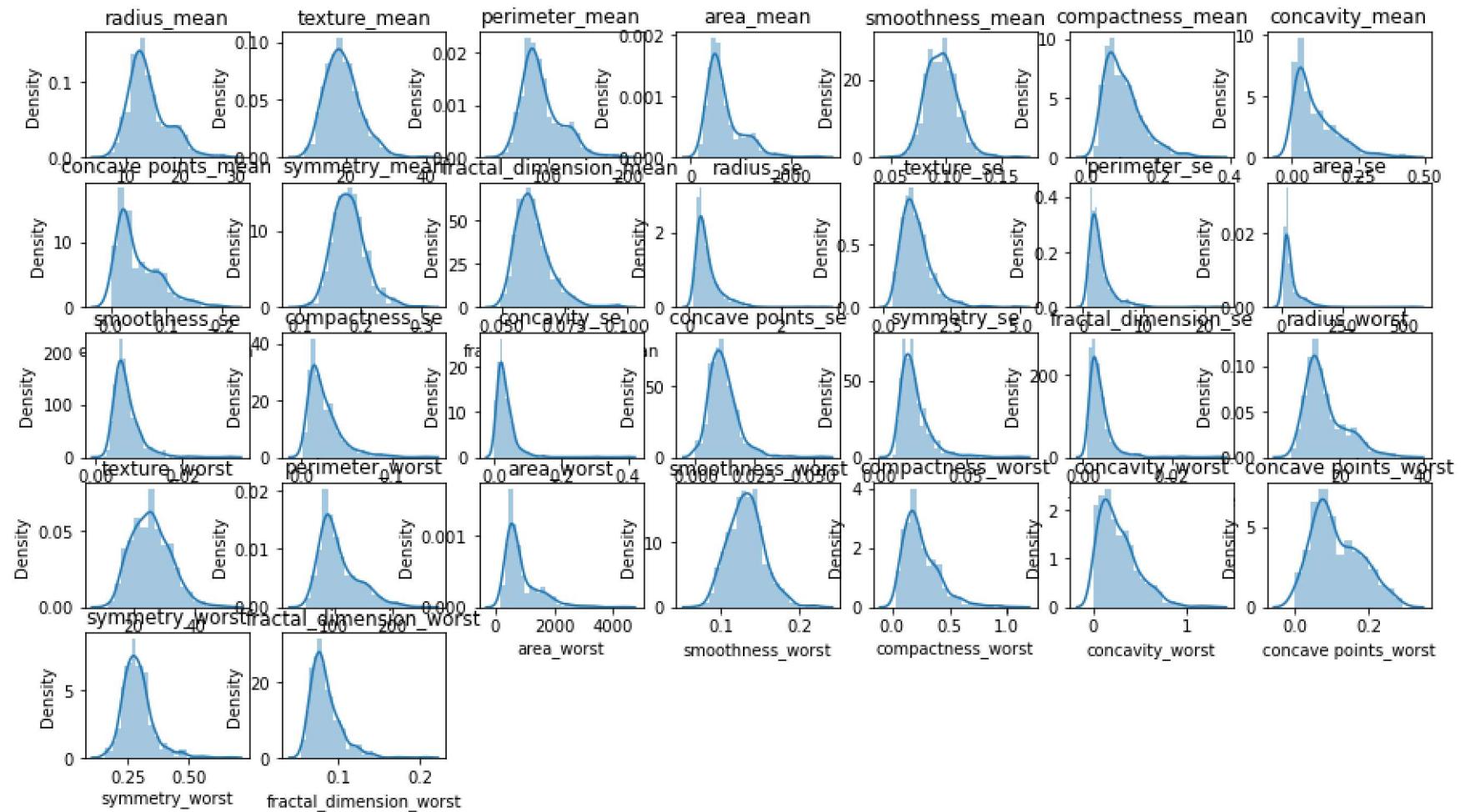
```
warnings.warn(msg, FutureWarning)
```

C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with

```
similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
```

```
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
    warnings.warn(msg, FutureWarning)
```



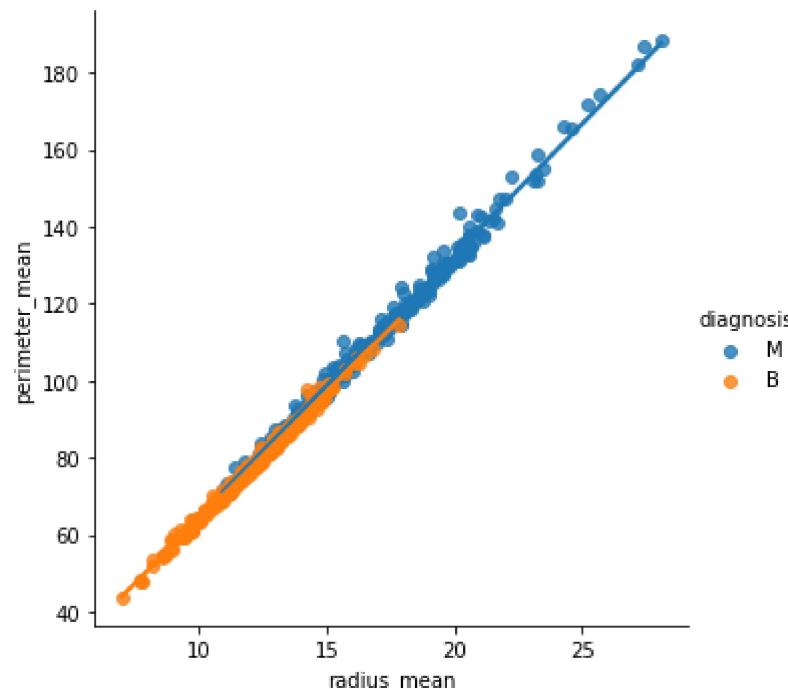
We can observe from the graphs that most of the variables are having a normal distribution, some variables such as `texture_worst` and `parameter_worst` are having a negatively-skewed distribution.

In [59]:

```
sns.lmplot('radius_mean', 'perimeter_mean', hue='diagnosis', data=df)
```

Out[59]:

```
<seaborn.axisgrid.FacetGrid at 0x1fe7e22d220>
```



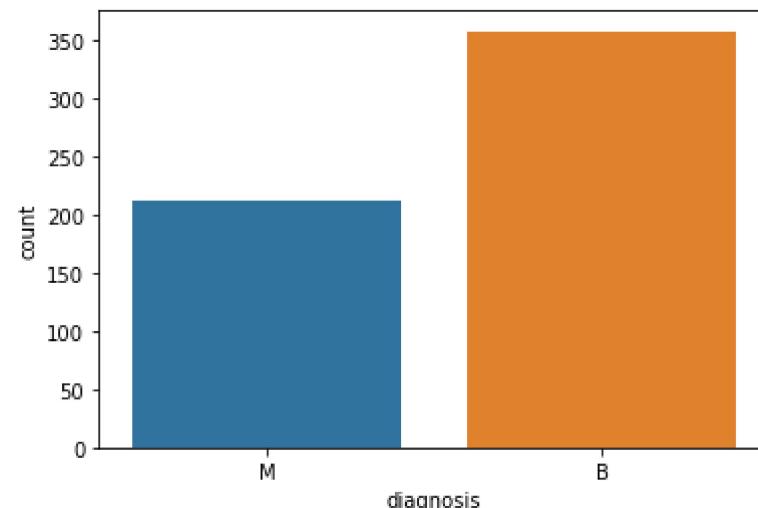
From the data we can observe a strong correlation between `perimeter_mean` and `radius_mean` and we can also observe that M have a higher value for both and F has comparatively smaller values.Hence we can eliminate one out of the 2 variables as they both are highly correlated.

```
In [21]: sns.countplot(df['diagnosis'])
```

```
C:\Users\Admin\Anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
```

```
    warnings.warn(
```

```
Out[21]: <AxesSubplot:xlabel='diagnosis', ylabel='count'>
```



```
In [10]: df1=pd.get_dummies(df,drop_first=True)  
df1.head()
```

```
Out[10]:
```

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean
0	15.46	19.48	101.70	748.9	0.10920	0.12230	0.14660	0.08087	0.19
1	12.89	13.12	81.89	515.9	0.06955	0.03729	0.02260	0.01171	0.13
2	14.96	19.10	97.03	687.3	0.08992	0.09823	0.05940	0.04819	0.18
3	12.94	16.17	83.18	507.6	0.09879	0.08836	0.03296	0.02390	0.17
4	13.17	18.66	85.98	534.6	0.11580	0.12310	0.12260	0.07340	0.21

5 rows × 31 columns

```
In [ ]:
```

```
In [12]: df1.isnull().sum()
```

```
Out[12]: radius_mean
```

0

```
texture_mean          0
perimeter_mean        0
area_mean             0
smoothness_mean       0
compactness_mean      0
concavity_mean        3
concave points_mean   3
symmetry_mean         0
fractal_dimension_mean 0
radius_se              0
texture_se              0
perimeter_se           0
area_se                0
smoothness_se          0
compactness_se          0
concavity_se            3
concave points_se      2
symmetry_se             0
fractal_dimension_se   0
radius_worst            0
texture_worst           23
perimeter_worst         0
area_worst              0
smoothness_worst        0
compactness_worst        0
concavity_worst          0
concave points_worst    0
symmetry_worst           0
fractal_dimension_worst 0
diagnosis_M             0
dtype: int64
```

Data Cleaning

Since texture_worst has highest number of missing values and it's graph is negatively skewed we've imputed the missing values with median.

In [13]:

```
df1.fillna(df.median(), inplace=True)
```

```
C:\Users\Admin\AppData\Local\Temp\ipykernel_20720/435454688.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError. Select only valid
```

```
columns before calling the reduction.  
df1.fillna(df.median(), inplace=True)
```

In []:

In []:

In [32]:

```
import math  
#Since the classes are unbalanced we will apply oversampling  
  
def _oversample_positives(df, target):  
    """ Oversample the minority classes to match  
    the majority class.  
  
    :param df: pandas dataframe - input df.  
    :param target: string - classification target column.  
  
    :return: pandas datframe - oversampled version  
    """  
  
    class_count = df[target].value_counts()  
  
    print("Before oversampling:\n %s" % class_count)  
  
    for i in range(1, len(class_count)):  
        df_i = df[df[target] == i]  
        oversampling_factor_i = class_count[0] / float(class_count[i])  
        print(len(df_i))  
        print("Oversampling factor for class %i: %s" %(i, str(oversampling_factor_i)))  
  
        # Integer part of oversampling  
        df = df.append(  
            [df_i] * int(math.floor(oversampling_factor_i) - 1),  
            ignore_index=False)  
  
        # Float part of oversampling  
        df = df.append(  
            [df_i.sample(frac=oversampling_factor_i % 1)],  
            ignore_index=False)  
  
    print("After oversampling:\n %s" % df[target].value_counts())  
    print("Shape after oversampling: %s" % str(df.shape))
```

```
    return df
df_oversampled = _oversample_positives(df1, 'diagnosis_M')
```

Before oversampling:

```
0    357
1    212
Name: diagnosis_M, dtype: int64
212
```

Oversampling factor for class 1: 1.6839622641509433

After oversampling:

```
1    357
0    357
Name: diagnosis_M, dtype: int64
Shape after oversampling: (714, 31)
```

In [33]:

```
#Storing Target variable in y and independent variables in x

X = df_oversampled.drop(["diagnosis_M"],axis=1)
Y = df_oversampled["diagnosis_M"]
```

In [34]:

```
# Scaling the Data : MinMax Scaler

#import MinMaxScaler from sklearn.preprocessing
from sklearn.preprocessing import MinMaxScaler

#storing MinMaxScalar in scaler
scaler = MinMaxScaler()

#fitting and transforming X
X = scaler.fit_transform(X)
```

In [35]:

```
# Train Test Split

#import train_test_split from sklearn.model_selection
from sklearn.model_selection import train_test_split

#splitting the data with size=0.2
x_train,x_test,y_train,y_test =train_test_split(X,Y,test_size = 0.2)
```

Dimensionality Reduction

```
In [36]: from sklearn.decomposition import PCA  
# Make an instance of the Model  
pca = PCA(.95)
```

```
In [37]: pca.fit(x_train)
```

```
Out[37]: PCA(n_components=0.95)
```

```
In [38]: x_train = pca.transform(x_train)  
x_test = pca.transform(x_test)
```

Model Building

```
In [39]: from sklearn.linear_model import LogisticRegression
```

```
In [40]: model = LogisticRegression(random_state=0)
```

```
In [41]: model.fit(x_train, y_train)
```

```
Out[41]: LogisticRegression(random_state=0)
```

```
In [43]: y_pred=model.predict(x_test)
```

```
In [47]: print('Accuracy of logistic regression classifier on test set: {:.2f}'.format(model.score(x_test, y_test)))
```

```
Accuracy of logistic regression classifier on test set: 0.98
```

Visualization and Interpretation

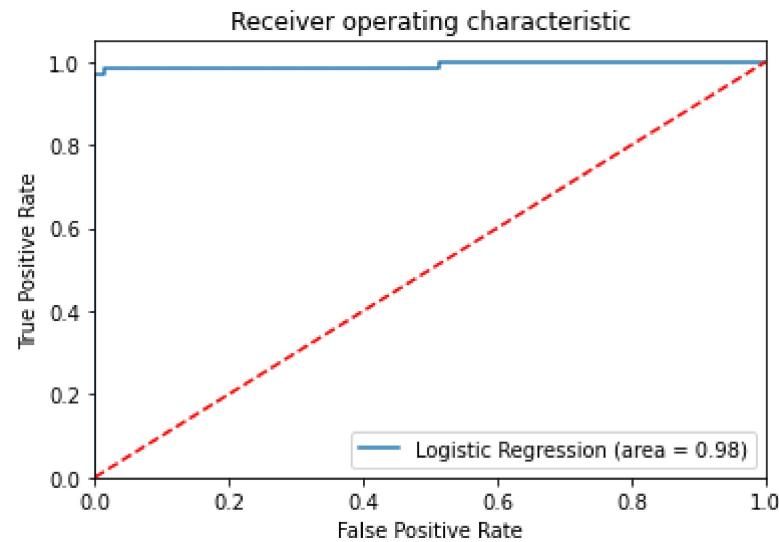
```
In [48]: from sklearn.metrics import confusion_matrix
confusion_matrix = confusion_matrix(y_test, y_pred)
print(confusion_matrix)
```

```
[[74  0]
 [ 3 66]]
```

```
In [49]: from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.96	1.00	0.98	74
1	1.00	0.96	0.98	69
accuracy			0.98	143
macro avg	0.98	0.98	0.98	143
weighted avg	0.98	0.98	0.98	143

```
In [51]: #ROC Curve
from sklearn.metrics import roc_auc_score
from sklearn.metrics import roc_curve
logit_roc_auc = roc_auc_score(y_test, model.predict(x_test))
fpr, tpr, thresholds = roc_curve(y_test, model.predict_proba(x_test)[:,1])
plt.figure()
plt.plot(fpr, tpr, label='Logistic Regression (area = %0.2f)' % logit_roc_auc)
plt.plot([0, 1], [0, 1], 'r--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic')
plt.legend(loc="lower right")
plt.savefig('Log_ROC')
plt.show()
```



The optimum values for FPR and TPR are approximately(0.02,0.97)