In [ ]:
```
#**********************************************************************************************
#                                    Breast_Cancer_Prediction_Project
#**********************************************************************************************
```

In [ ]:
```
#This project is made using CSV File called Data
#This CSV file is downloaded from https://www.kaggle.com/uciml/breast-cancer-wisconsin-data
#Project by:Anubha Sharma-Data Science
#Submitted to:Jyotika
```

In [2]:
```
#1.Import Libraries
import numpy
import pandas as pd
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
```

Matplotlib is building the font cache; this may take a moment.

In [4]:
```
df=pd.read_csv("data.csv")
df.head()
print(df)
```

```
           id diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean  \
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

     smoothness_mean  compactness_mean  concavity_mean  concave points_mean  \
0            0.11840           0.27760         0.30010              0.14710
1            0.08474           0.07864         0.08690              0.07017
2            0.10960           0.15990         0.19740              0.12790
3            0.14250           0.28390         0.24140              0.10520
4            0.10030           0.13280         0.19800              0.10430
```

```
..         ...            ...            ...              ...
564      0.11100        0.11590        0.24390          0.13890
565      0.09780        0.10340        0.14400          0.09791
566      0.08455        0.10230        0.09251          0.05302
567      0.11780        0.27700        0.35140          0.15200
568      0.05263        0.04362        0.00000          0.00000

        ...  texture_worst  perimeter_worst  area_worst  smoothness_worst  \
0       ...          17.33           184.60      2019.0           0.16220
1       ...          23.41           158.80      1956.0           0.12380
2       ...          25.53           152.50      1709.0           0.14440
3       ...          26.50            98.87       567.7           0.20980
4       ...          16.67           152.20      1575.0           0.13740
..      ...            ...              ...         ...               ...
564     ...          26.40           166.10      2027.0           0.14100
565     ...          38.25           155.00      1731.0           0.11660
566     ...          34.12           126.70      1124.0           0.11390
567     ...          39.42           184.60      1821.0           0.16500
568     ...          30.37            59.16       268.6           0.08996

     compactness_worst  concavity_worst  concave points_worst  symmetry_worst  \
0              0.66560           0.7119                0.2654          0.4601
1              0.18660           0.2416                0.1860          0.2750
2              0.42450           0.4504                0.2430          0.3613
3              0.86630           0.6869                0.2575          0.6638
4              0.20500           0.4000                0.1625          0.2364
..                 ...              ...                   ...             ...
564            0.21130           0.4107                0.2216          0.2060
565            0.19220           0.3215                0.1628          0.2572
566            0.30940           0.3403                0.1418          0.2218
567            0.86810           0.9387                0.2650          0.4087
568            0.06444           0.0000                0.0000          0.2871

     fractal_dimension_worst  Unnamed: 32
0                    0.11890          NaN
1                    0.08902          NaN
2                    0.08758          NaN
3                    0.17300          NaN
4                    0.07678          NaN
..                       ...          ...
564                  0.07115          NaN
565                  0.06637          NaN
566                  0.07820          NaN
567                  0.12400          NaN
568                  0.07039          NaN

[569 rows x 33 columns]
```

In [32]:

```
df.describe()
```

Out[32]:

|  | id | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean | smoothness_mean | compactness_mean | concavity_m |
|---|---|---|---|---|---|---|---|---|---|
| **count** | 5.690000e+02 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000000 | 569.000 |
| **mean** | 3.037183e+07 | 0.372583 | 14.127292 | 19.289649 | 91.969033 | 654.889104 | 0.096360 | 0.104341 | 0.088 |
| **std** | 1.250206e+08 | 0.483918 | 3.524049 | 4.301036 | 24.298981 | 351.914129 | 0.014064 | 0.052813 | 0.079 |
| **min** | 8.670000e+03 | 0.000000 | 6.981000 | 9.710000 | 43.790000 | 143.500000 | 0.052630 | 0.019380 | 0.000 |
| **25%** | 8.692180e+05 | 0.000000 | 11.700000 | 16.170000 | 75.170000 | 420.300000 | 0.086370 | 0.064920 | 0.029 |
| **50%** | 9.060240e+05 | 0.000000 | 13.370000 | 18.840000 | 86.240000 | 551.100000 | 0.095870 | 0.092630 | 0.061 |
| **75%** | 8.813129e+06 | 1.000000 | 15.780000 | 21.800000 | 104.100000 | 782.700000 | 0.105300 | 0.130400 | 0.130 |
| **max** | 9.113205e+08 | 1.000000 | 28.110000 | 39.280000 | 188.500000 | 2501.000000 | 0.163400 | 0.345400 | 0.426 |

8 rows × 33 columns

In [34]:
```
df.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
```

In [31]:
```python
#**************************
#1.Training Model
#Here 1 stands for M
#Here 2 stands for B
#**************************

from sklearn.preprocessing import LabelEncoder
labelencoder_Y = LabelEncoder()
df.iloc[:,1]=labelencoder_Y.fit_transform(df.iloc[:,1].values)
df.head()
```

Out[31]:

| | id | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean | smoothness_mean | compactness_mean | concavity_mean | c point |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 842302 | 1 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.3001 | |
| 1 | 842517 | 1 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.0869 | |
| 2 | 84300903 | 1 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.1974 | |
| 3 | 84348301 | 1 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.2414 | |
| 4 | 84358402 | 1 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.1980 | |

5 rows × 33 columns

In [15]:
```python
#************************
# Models/ Algorithms
#************************
def models(X_train,Y_train):


    #**********************
    #logistic regression
    #**********************
    from sklearn.linear_model import LogisticRegression
    log=LogisticRegression(random_state=0)
    log.fit(X_train,Y_train)



    #***********************
    #Decision Tree
    #***********************
    from sklearn.tree import DecisionTreeClassifier
    tree=DecisionTreeClassifier(random_state=0,criterion="entropy")
    tree.fit(X_train,Y_train)


    #*********************
    #Random Forest
    #*********************
    from sklearn.ensemble import RandomForestClassifier
    forest=RandomForestClassifier(random_state=0,criterion="entropy",n_estimators=10)
    forest.fit(X_train,Y_train)

    print('[0]logistic regression accuracy:',log.score(X_train,Y_train))
    print('[1]Decision tree accuracy:',tree.score(X_train,Y_train))
    print('[2]Random forest accuracy:',forest.score(X_train,Y_train))


    return log,tree,forest
```

In [21]:
```python
model=models(X_train,Y_train)
```

```
[0]logistic regression accuracy: 0.9912087912087912
[1]Decision tree accuracy: 1.0
[2]Random forest accuracy: 0.9978021978021978
```
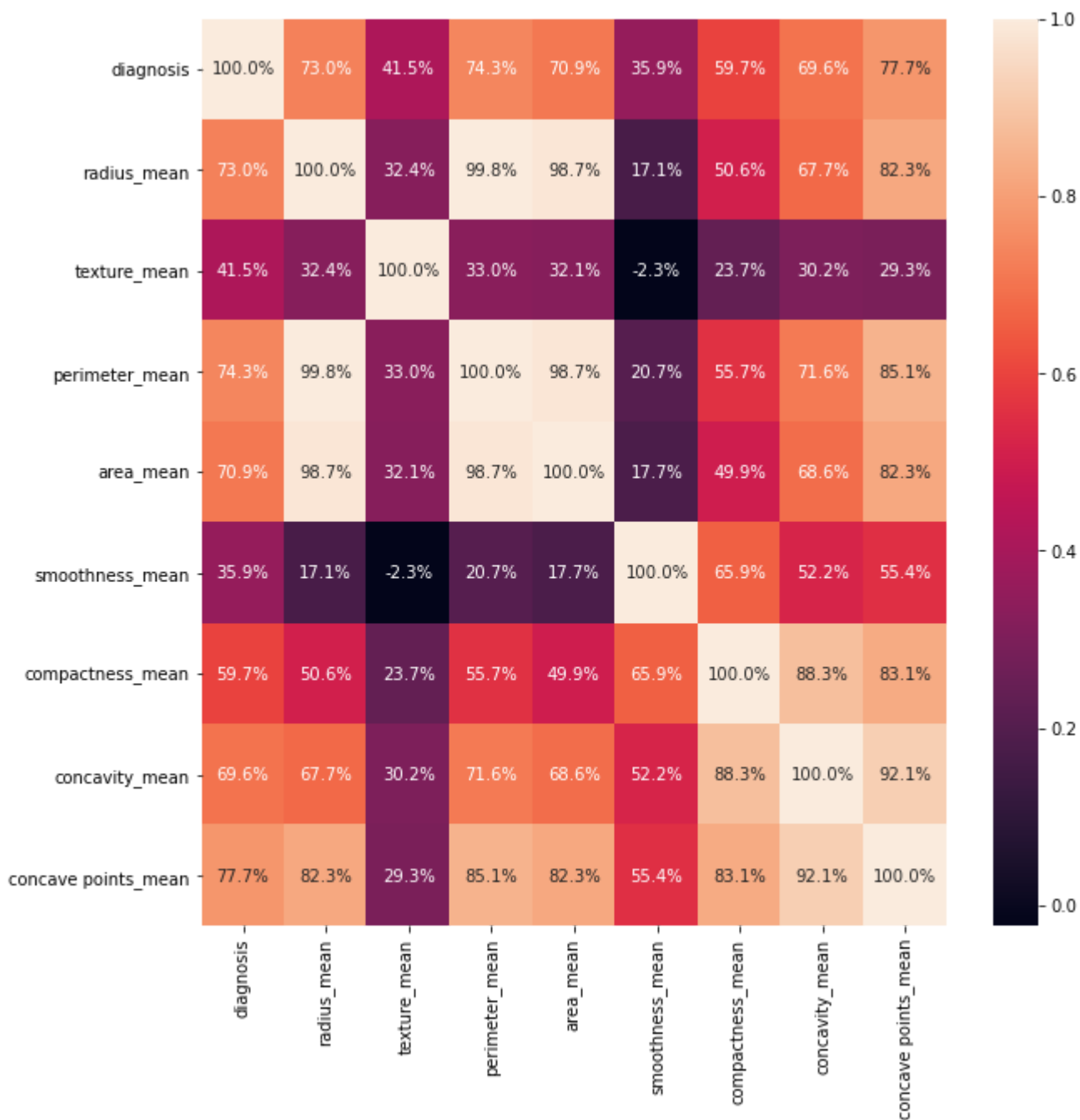
In [20]:
```python
plt.figure(figsize=(10,10))
sns.heatmap(df.iloc[:,1:10].corr(),annot=True,fmt=".01%")
```
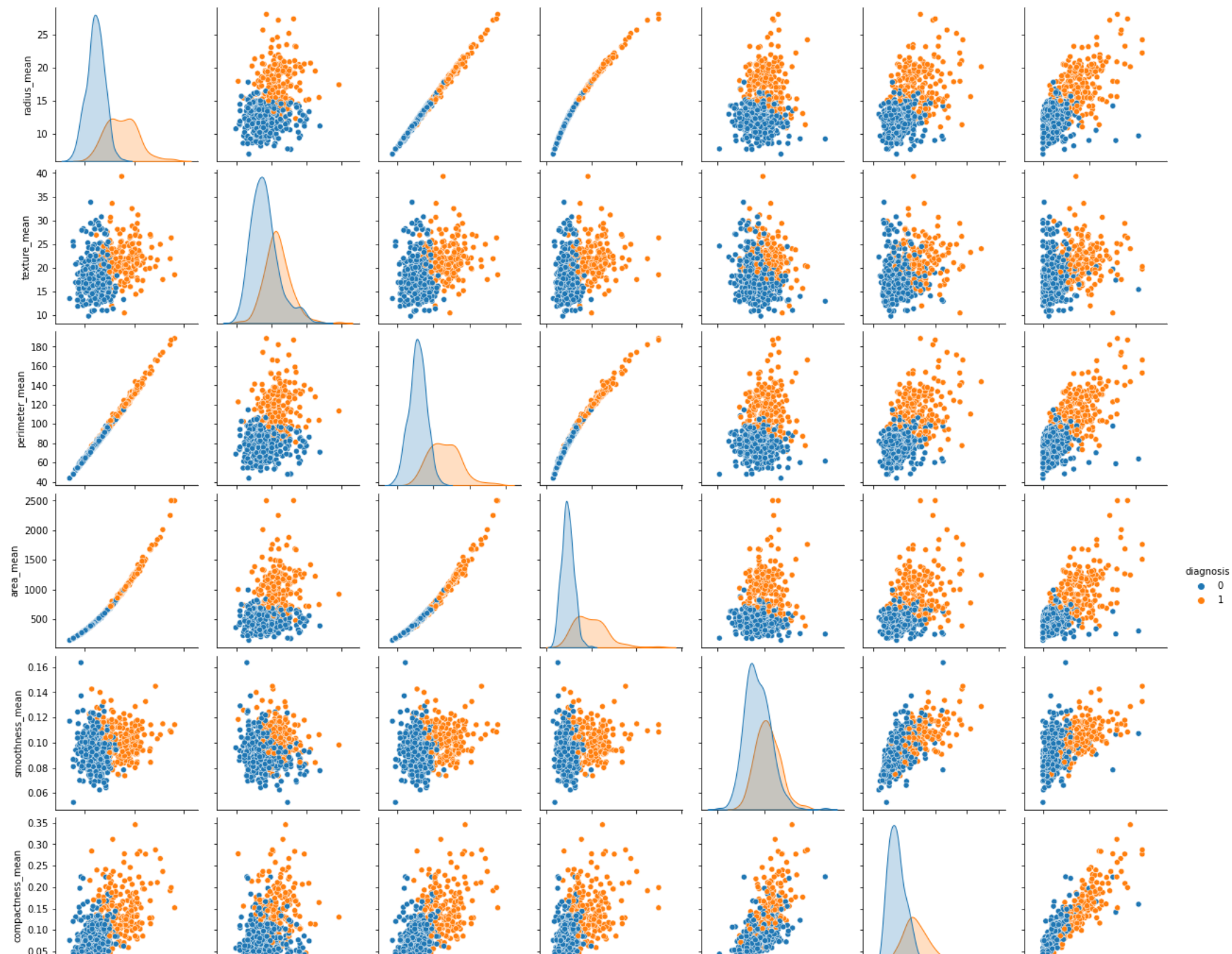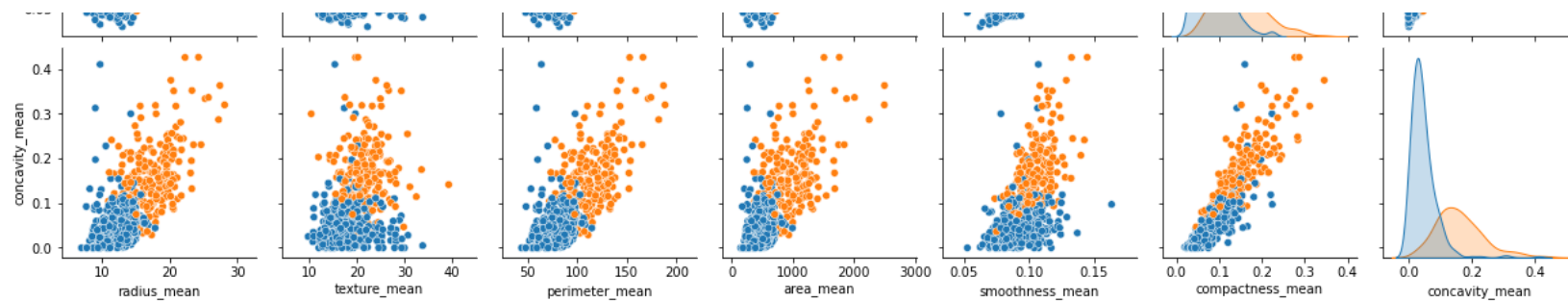
Out[20]:  <AxesSubplot:>

In [33]:
```python
sns.pairplot(df.iloc[:,1:9],hue="diagnosis")
```
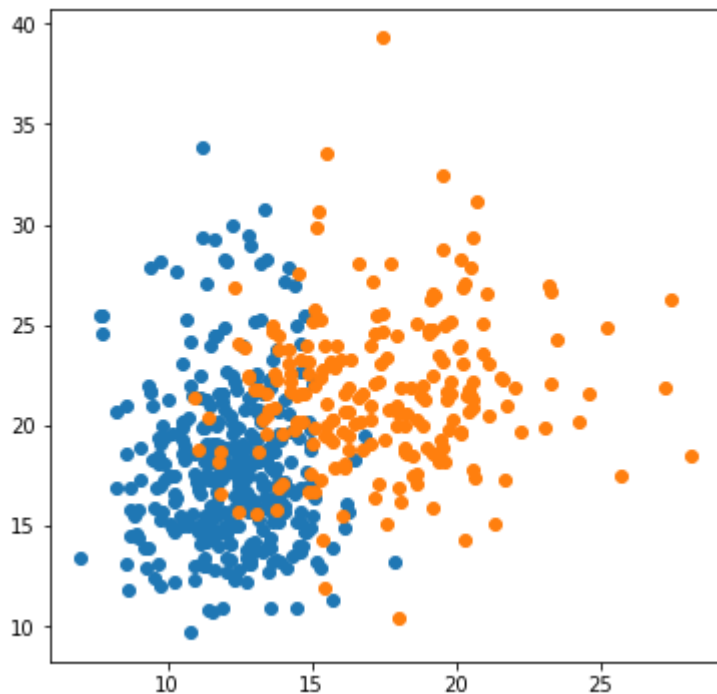
Out[33]: &lt;seaborn.axisgrid.PairGrid at 0x29bf4a19280&gt;

In [39]:
```python
plt.figure(figsize=(6,6))
plt.scatter(X[Y == 0][:, 0], X[Y == 0][:, 1],  label = '0')
plt.scatter(X[Y == 1][:, 0], X[Y == 1][:, 1],  label = '1')
```

Out[39]:  <matplotlib.collections.PathCollection at 0x29bf7330490>



In [23]:
```python
#****************************************
# Testing the models
#****************************************
```

```python
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report

for i in range(len(model)):
    print("Model",i)
    print(classification_report(Y_test,model[i].predict(X_test)))
    print('Accuracy : ',accuracy_score(Y_test,model[i].predict(X_test)))
```

```
Model 0
              precision    recall  f1-score   support

           0       0.96      0.99      0.97        67
           1       0.98      0.94      0.96        47

    accuracy                           0.96       114
   macro avg       0.97      0.96      0.96       114
weighted avg       0.97      0.96      0.96       114

Accuracy :  0.9649122807017544
Model 1
              precision    recall  f1-score   support

           0       0.94      0.96      0.95        67
           1       0.93      0.91      0.92        47

    accuracy                           0.94       114
   macro avg       0.94      0.94      0.94       114
weighted avg       0.94      0.94      0.94       114

Accuracy :  0.9385964912280702
Model 2
              precision    recall  f1-score   support

           0       0.96      1.00      0.98        67
           1       1.00      0.94      0.97        47

    accuracy                           0.97       114
   macro avg       0.98      0.97      0.97       114
weighted avg       0.97      0.97      0.97       114

Accuracy :  0.9736842105263158
```

```python
In [18]:  #*********************************************
          #Accuracy
          #*********************************************
          pred=model[2].predict(X_test)
          print('Predicted values:')
```

```
print(pred)
print('Actual values:')
print(Y_test)
```

```
Predicted values:
[1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 1 1 1 1 0 0 1 0 0 1 0 1 0 1 0 1 0 1 0
 1 0 1 0 0 0 0 0 1 0 0 0 1 1 1 1 0 0 0 0 0 0 1 1 1 0 0 1 0 1 1 1 0 0 1 0 0
 1 0 0 0 0 0 1 1 1 0 1 0 0 0 1 1 0 1 0 1 0 0 1 0 0 0 0 0 0 0 1 0 1 0 1 1 0
 1 1 0]
Actual values:
[1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 1 1 1 1 0 0 1 0 0 1 0 1 0 1 0 1 0 1 0
 1 0 1 1 0 1 0 0 1 0 0 0 1 1 1 1 0 0 0 0 0 0 1 1 1 0 0 1 0 1 1 1 0 0 1 0 1
 1 0 0 0 0 0 1 1 1 0 1 0 0 0 1 1 0 1 0 1 0 0 1 0 0 0 0 0 0 0 1 0 1 0 1 1 0
 1 1 0]
```

In [ ]:
```
#*********************************************************************************************************
#                                          End of the Project
#*********************************************************************************************************
```