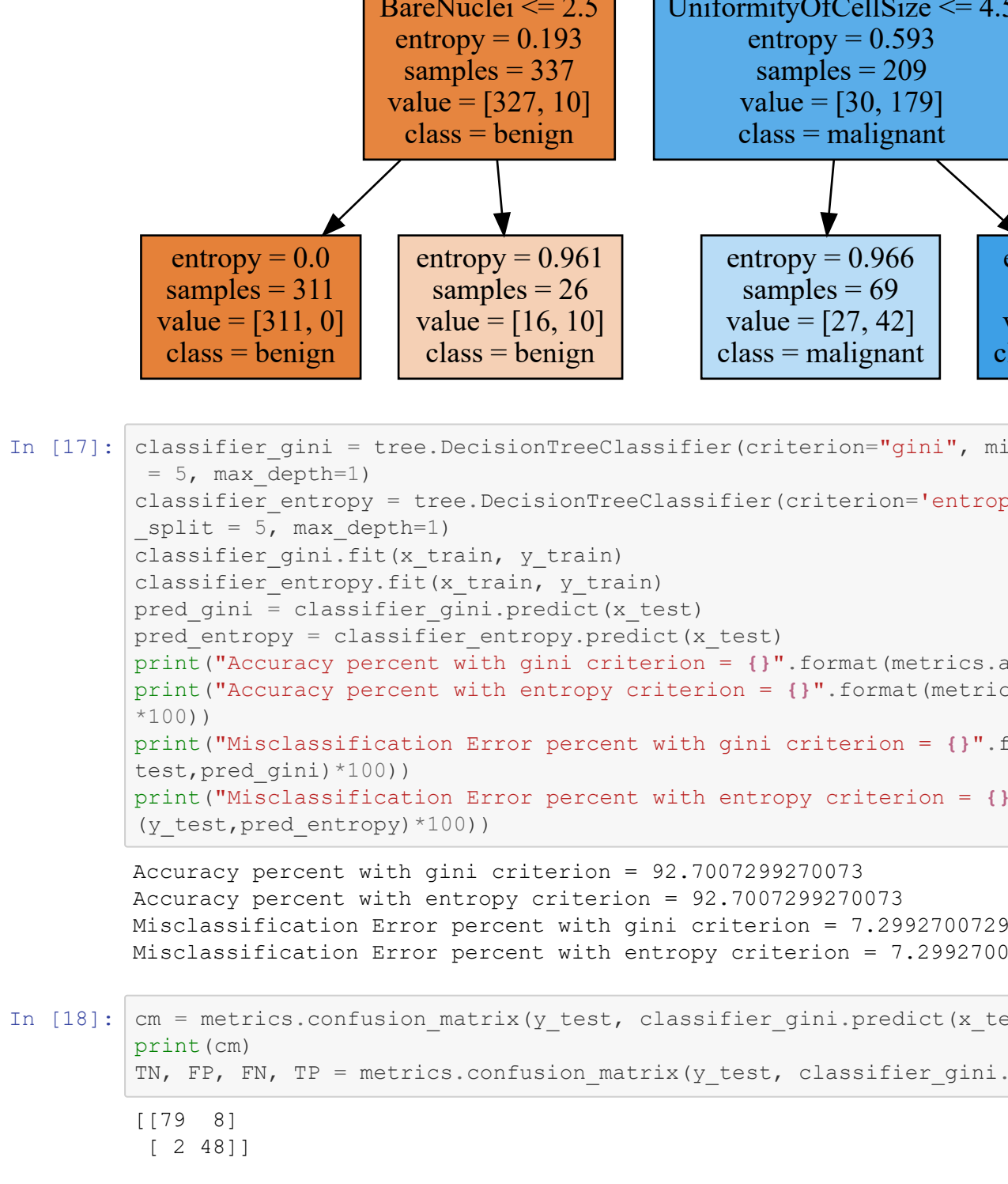



```
In [16]: # Graph with crit="entropy" and max_depth = 2
classifier_entropy = tree.DecisionTreeClassifier(criterion="entropy", min_samples_leaf = 2, min_samples_split = 5, max_depth = 2)
clf_entropy.fit(x_train, y_train)

dot_data = tree.export_graphviz(classifier_entropy,
                                feature_names=feature3,
                                class_names=['benign', 'malignant'],
                                filled=True)

graph = graphviz.Source(dot_data, format="png")
graph
```



```
In [17]: classifier_gini = tree.DecisionTreeClassifier(criterion="gini", min_samples_leaf = 2, min_samples_split = 5, max_depth=1)
classifier_entropy = tree.DecisionTreeClassifier(criterion="entropy", min_samples_leaf = 2, min_samples_split = 5, max_depth=1)
classifier_gini.fit(x_train, y_train)
classifier_entropy.fit(x_train, y_train)
pred_gini = classifier_gini.predict(x_test)
pred_entropy = classifier_entropy.predict(x_test)
print("Accuracy percent with gini criterion = {}".format(metrics.accuracy_score(y_test, pred_gini)*100))
print("Accuracy percent with entropy criterion = {}".format(metrics.accuracy_score(y_test, pred_entropy)*100))
print("Misclassification Error percent with gini criterion = {}".format(100 - metrics.accuracy_score(y_test, pred_gini)*100))
print("Misclassification Error percent with entropy criterion = {}".format(100 - metrics.accuracy_score(y_test, pred_entropy)*100))

Accuracy percent with gini criterion = 92.7007299270073
Accuracy percent with entropy criterion = 92.7007299270073
Misclassification Error percent with gini criterion = 7.299270072992698
Misclassification Error percent with entropy criterion = 7.299270072992698
```

```
In [18]: cm = metrics.confusion_matrix(y_test, classifier_gini.predict(x_test))
print(cm)
TN, FP, FN, TP = metrics.confusion_matrix(y_test, classifier_gini.predict(x_test)).ravel()
```

```
In [19]: print("TP = {}".format(TP))
print("FP = {}".format(FP))
print("TN = {}".format(TN))
print("FN = {}".format(FN))

# calculate accuracy
print((TP+TN)/float(TP+TN+FP+FN))
print(metrics.accuracy_score(y_test, classifier_gini.predict(x_test)))

TP = 48
FP = 8
FN = 2
TN = 79
0.927007299270073
0.927007299270073
```

We get same accuracy in both the case, manually and using function.

```
In [20]: # calculate Misclassification Error
print((FP+FN)/float(TP+TN+FP+FN))
print(1 - metrics.accuracy_score(y_test, classifier_gini.predict(x_test)))

0.072992700729927
0.07299270072992703
```

We get same Misclassification Error in both the case, manually and using function.

Conclusion:

- gini index at first split = 0.453
- entropy at first split = 0.931
- Misclassification Error at first split with gini criterion = 0.07299270072992698
- Misclassification Error at first split with entropy criterion = 0.07299270072992698
- Information gain

= parent entropy - average entropy of the children
= 0.931 - ((337/546)*0.193 + (209/546)*0.593)
= 0.585

- "Uniformity Of Cell Shape" feature is selected at first split.
- Decision boundary is <= 2.5

Problem 3

```
In [21]: from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
```

```
In [22]: bc_df3 = "kdbc.data"
column_names3 = ["id", "Diagnosis",
                 "radius_mean", "texture_mean", "perimeter_mean", "area_mean", "smoothness_mean",
                 "compactness_mean", "concavity_mean", "concave points_mean", "symmetry_mean", "fractal_dimension_mean",
                 "radius_se", "texture_se", "perimeter_se", "area_se", "smoothness_se",
                 "compactness_se", "concavity_se", "concave points_se", "symmetry_se", "fractal_dimension_se",
                 "radius_worst", "texture_worst", "perimeter_worst", "area_worst", "smoothness_worst",
                 "compactness_worst", "concavity_worst", "concave points_worst", "symmetry_worst", "fractal_dimension_worst"]
bc_df3 = pd.read_csv(bc_file3, names=column_names3)
bc_df3
```

Out [22]:

	id	Diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27780	0.30010	0.303
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08890	0.096
2	8430003	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0.191
3	8434801	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.241
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.198
...
554	926424	M	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.12490	0.241
555	926882	M	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.144
556	926954	M	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.095
557	927241	M	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.351
558	927551	B	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.000

569 rows x 11 columns

```
In [23]: bc_df3['Diagnosis'] = bc_df3['Diagnosis'].apply(lambda x: 1 if x == 'M' else 0)
bc_df3.drop(['id'], axis=1, inplace=True)
bc_df3
```

Out [23]:

	Diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean
0	1	17.99	10.38	122.80	1001.0	0.11840	0.27780	0.30010	0
1	1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08890	0
2	1	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0
3	1	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0
4	1	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0
...
554	1	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.12490	0
555	1	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0
556	1	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0
557	1	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0
558	0	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0

569 rows x 10 columns

```
In [24]: # split the original data
feature3 = list(bc_df3.columns[1:])
x3 = bc_df3[feature3]
y3 = bc_df3['Diagnosis']
x3_train, x3_test, y3_train, y3_test = model_selection.train_test_split(x3, y3,
                                                                          test_size=0.2,
                                                                          random_state=0)
```

```
In [25]: # plot a graph of the original data with criterion="gini" and max_depth = 2
clf_gini = tree.DecisionTreeClassifier(criterion="gini", min_samples_leaf = 3, min_samples_split = 5, max_depth = 2)
clf_gini.fit(x3_train, y3_train)

dot_data = tree.export_graphviz(clf_gini, feature_names=feature3,
                                class_names=['benign', 'malignant'],
                                filled=True)

# Draw graph
graph = graphviz.Source(dot_data, format="png")
graph
```



```
In [26]: e = y3_test
p = clf_gini.predict(x3_test)
```

```
In [27]: # Precision, Recall and F1 Score for the original data
print("Precision, Recall and F1 Score with macro = {}".format(metrics.precision_recall_fscore_support(e, p, average='macro')))
print("Precision, Recall and F1 Score with micro = {}".format(metrics.precision_recall_fscore_support(e, p, average='micro')))
print("Precision, Recall and F1 Score with weighted = {}".format(metrics.precision_recall_fscore_support(e, p, average='weighted')))
print(metrics.classification_report(e, p))

Precision, Recall and F1 Score with macro = (0.9671497584541062, 0.9606224198158145, 0.9635549872122762, None)
Precision, Recall and F1 Score with micro = (0.9649122807017544, 0.9649122807017544, 0.9649122807017544, None)
Precision, Recall and F1 Score with weighted = (0.9652851936604796, 0.9649122807017544, 0.9647888903845291, None)

precision    recall    f1-score   support

0         0.96         0.99         0.97         67
1         0.98         0.94         0.96         47

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

```
In [28]: cm = metrics.confusion_matrix(e, p)
print(cm)
TN, FP, FN, TP = metrics.confusion_matrix(e, p).ravel()
```

```
In [29]: print("TP = {}".format(TP))
print("FP = {}".format(FP))
print("TN = {}".format(TN))
print("FN = {}".format(FN))
print("TPR = {}".format(TP/float(TP+FP)))
print("FPR = {}".format(FP/float(TN+FP)))

TP = 44
FP = 1
FN = 3
TN = 66
TPR = 0.9777777777777777
FPR = 0.014925373134328358
```

```
In [30]: # Scale the data
scaler = StandardScaler()
scaler.fit(bc_df3.values)
scaled_data = scaler.transform(bc_df3.values)
```

```
In [31]: # Apply PCA for n_components=1
pca_1 = PCA(n_components=1)
pca_1.fit(scaled_data)
pca_1.transformed = pca_1.transform(scaled_data)
```

```
In [32]: pca_df_1 = pd.DataFrame(pca_1.transformed, columns = ['PC1'])
pca_df_1
```

Out [32]:

	PC1
0	9.225770
1	2.656802
2	5.892492
3	7.135401
4	4.129423
...	...
554	6.593983
555	4.024833
556	1.530077
557	10.405008
558	-5.504862

569 rows x 1 columns

```
In [33]: feature3 = list(pca_df_1.columns)
x3 = pca_df_1[feature3]
y3 = bc_df3['Diagnosis']
x3_train, x3_test, y3_train, y3_test = model_selection.train_test_split(x3, y3,
                                                                          test_size=0.2,
                                                                          random_state=0)
```

```
In [34]: clf_gini = tree.DecisionTreeClassifier(criterion="gini", min_samples_leaf = 3, min_samples_split = 5, max_depth = 2)
clf_gini.fit(x3_train, y3_train)

dot_data = tree.export_graphviz(clf_gini, feature_names=feature3,
                                class_names=['benign', 'malignant'],
                                filled=True)

graph = graphviz.Source(dot_data, format="png")
graph
```



```
In [35]: e = y3_test
p = clf_gini.predict(x3_test)
```

```
In [36]: print("Precision, Recall and F1 Score with macro = {}".format(metrics.precision_recall_fscore_support(e, p, average='macro')))
print("Precision, Recall and F1 Score with micro = {}".format(metrics.precision_recall_fscore_support(e, p, average='micro')))
print("Precision, Recall and F1 Score with weighted = {}".format(metrics.precision_recall_fscore_support(e, p, average='weighted')))
print(metrics.classification_report(e, p))

Precision, Recall and F1 Score with macro = (0.8812034739454093, 0.8902826294061607, 0.8839558374442095, None)
Precision, Recall and F1 Score with micro = (0.8859649122807017, 0.8859649122807017, 0.8859649122807017, None)
Precision, Recall and F1 Score with weighted = (0.890726350615994, 0.8859649122807017, 0.8866346038928657, None)

precision    recall    f1-score   support

0         0.94         0.87         0.90         67
1         0.83         0.91         0.87         47

accuracy         0.88         0.89         0.88         114
macro avg        0.89         0.89         0.89         114
weighted avg     0.89         0.89         0.89         114
```

```
In [37]: cm = metrics.confusion_matrix(e, p)
print(cm)
TN, FP, FN, TP = metrics.confusion_matrix(e, p).ravel()
```

```
In [38]: print("TP = {}".format(TP))
print("FP = {}".format(FP))
print("TN = {}".format(TN))
print("FN = {}".format(FN))
print("TPR = {}".format(TP/float(TP+FP)))
print("FPR = {}".format(FP/float(TN+FP)))

TP = 43
FP = 9
FN = 4
TN = 58
TPR = 0.8269230769230769
FPR = 0.13432835820895522
```

```
In [39]: pca_2 = PCA(n_components=2)
pca_2.fit(scaled_data)
pca_2.transformed = pca_2.transform(scaled_data)
```

```
In [40]: pca_df_2 = pd.DataFrame(pca_2.transformed, columns = ['PC1', 'PC2'])
pca_df_2
```

Out [40]:

	PC1	PC2
0	9.225770	2.116196
1	2.656802	-3.784776
2	5.892492	-1.005579
3	7.135401	10.318716
4	4.129423	-1.905579
...
554	6.593983	-3.454947
555	4.024833	-1.958008
556	1.530077	-1.958071
557	10.405008	1.849078
558	-5.504862	-0.768348

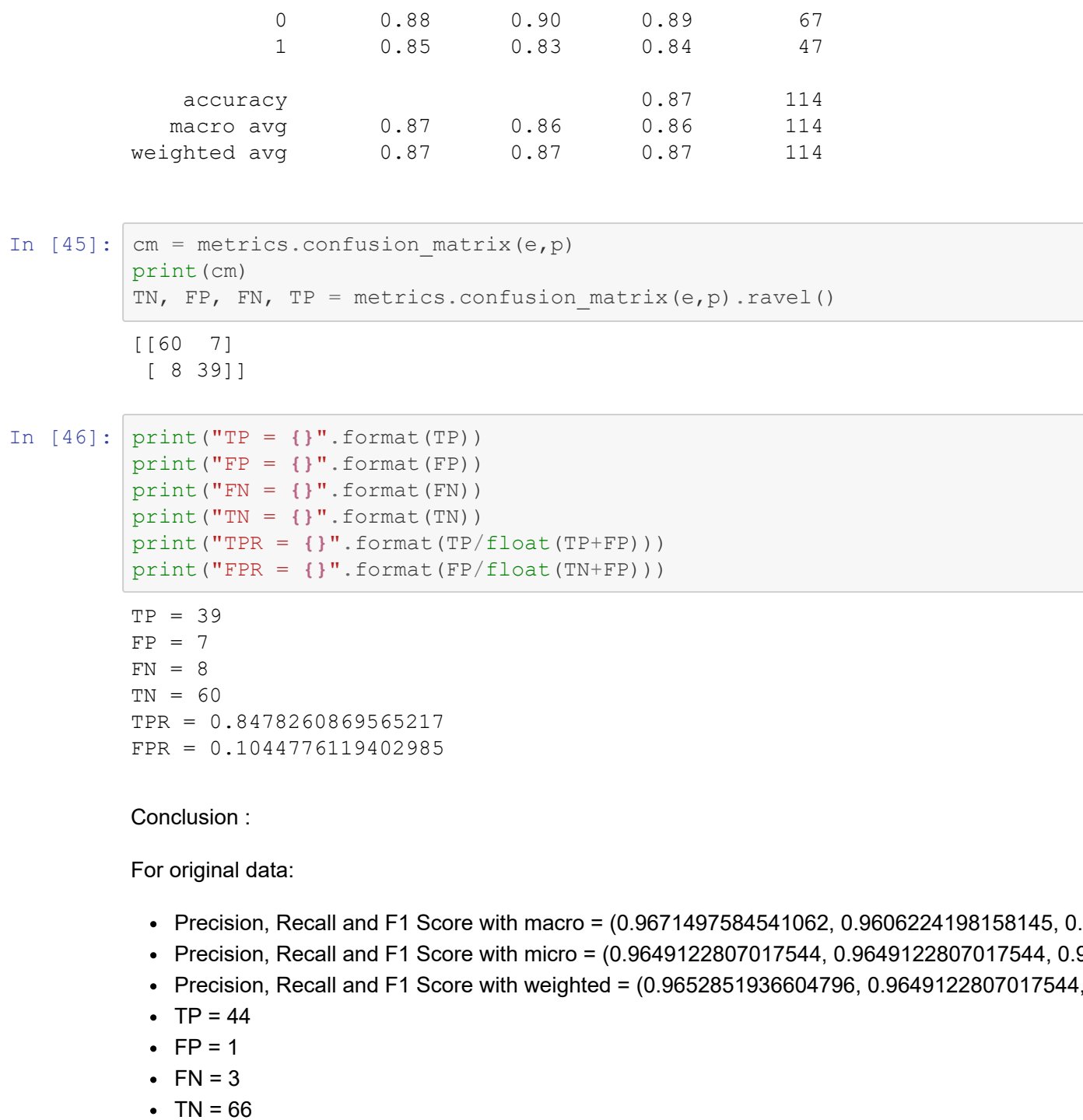
569 rows x 2 columns

```
In [41]: feature3 = list(pca_df_2.columns)
x3 = pca_df_2[feature3]
y3 = bc_df3['Diagnosis']
x3_train, x3_test, y3_train, y3_test = model_selection.train_test_split(x3, y3,
                                                                          test_size=0.2,
                                                                          random_state=0)
```

```
In [42]: clf_gini = tree.DecisionTreeClassifier(criterion="gini", min_samples_leaf = 3, min_samples_split = 5, max_depth = 2)
clf_gini.fit(x3_train, y3_train)

dot_data = tree.export_graphviz(clf_gini, feature_names=feature3,
                                class_names=['benign', 'malignant'],
                                filled=True)

graph = graphviz.Source(dot_data, format="png")
graph
```



```
In [43]: e = y3_test
p = clf_gini.predict(x3_test)
```

```
In [44]: print("Precision, Recall and F1 Score with macro = {}".format(metrics.precision_recall_fscore_support(e, p, average='macro')))
print("Precision, Recall and F1 Score with micro = {}".format(metrics.precision_recall_fscore_support(e, p, average='micro')))
print("Precision, Recall and F1 Score with weighted = {}".format(metrics.precision_recall_fscore_support(e, p, average='weighted')))
print(metrics.classification_report(e, p))

Precision, Recall and F1 Score with macro = (0.8650895140664961, 0.8626548110511274, 0.8637992831541219, None)
Precision, Recall and F1 Score with micro = (0.868421052631579, 0.868421052631579, 0.868421052631579, None)
Precision, Recall and F1 Score with weighted = (0.868421052631579, 0.868421052631579, 0.8682009683707477, None)

precision    recall    f1-score   support

0         0.88         0.90         0.89         67
1         0.85         0.83         0.84         47

accuracy         0.87         0.86         0.86         114
macro avg        0.87         0.86         0.86         114
weighted avg     0.87         0.87         0.87         114
```

```
In [45]: cm = metrics.confusion_matrix(e, p)
print(cm)
TN, FP, FN, TP = metrics.confusion_matrix(e, p).ravel()
```

```
In [46]: print("TP = {}".format(TP))
print("FP = {}".format(FP))
print("TN = {}".format(TN))
print("FN = {}".format(FN))
print("TPR = {}".format(TP/float(TP+FP)))
print("FPR = {}".format(FP/float(TN+FP)))

TP = 39
FP = 7
FN = 8
TN = 60
TPR = 0.8478260869565217
FPR = 0.1044776119402985
```

Conclusion:

For original data:

- Precision, Recall and F1 Score with macro = (0.9671497584541062, 0.9606224198158145, 0.9635549872122762, None)
- Precision, Recall and F1 Score with micro = (0.9649122807017544, 0.9649122807017544, 0.9649122807017544, None)
- Precision, Recall and F1 Score with weighted = (0.9652851936604796, 0.9649122807017544, 0.9647888903845291, None)
- TP = 44
- FP = 1
- FN = 3
- TN = 66
- TPR = 0.9777777777777777
- FPR = 0.014925373134328358

After applying PCA with n_component=1:

- Precision, Recall and F1 Score with macro = (0.8812034739454093, 0.8902826294061607, 0.8839558374442095, None)
- Precision, Recall and F1 Score with micro = (0.8859649122807017, 0.8859649122807017, 0.8859649122807017, None)
- Precision, Recall and F1 Score with weighted = (0.890726350615994, 0.8859649122807017, 0.8866346038928657, None)
- TP = 43
- FP = 9
- FN = 4
- TN = 58
- TPR = 0.8269230769230769
- FPR = 0.13432835820895522

After applying PCA with n_component=2:

- Precision, Recall and F1 Score with macro = (0.8650895140664961, 0.8626548110511274, 0.8637992831541219, None)
- Precision, Recall and F1 Score with micro = (0.868421052631579, 0.868421052631579, 0.868421052631579, None)
- Precision, Recall and F1 Score with weighted = (0.868421052631579, 0.868421052631579, 0.8682009683707477, None)
- TP = 39
- FP = 7
- FN = 8
- TN = 60
- TPR = 0.8478260869565217
- FPR = 0.1044776119402985

Decision trees work with continuous variables same as PCA works.

For example, consider age as the target variable. So, you compute the variance of age and it comes out to be x. Next, decision tree looks at various splits and calculates the total weighted variance of each of these splits. It chooses the split which provides the minimum variance (Greedy approach).

TPR value is increasing and FPR value is decreasing as the number of components retained by PCA also increased.

Problem 4

```
In [47]: A = np.random.normal(5, 2, 500)
B = np.random.normal(-5, 2, 500)
```

```
In [48]: C1 = np.repeat(0, 500)
C2 = np.repeat(1, 500)
```

```
In [49]: dfA = pd.DataFrame(zip(A, C
```