

In [15]: *#Created by Rami ALmehdawi*

Import the necessary Libraries

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import sklearn as sk
from sklearn import tree
from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, confusion_matrix, classification_report, roc_
from imblearn.over_sampling import RandomOverSampler, SMOTE
```

In [16]: *# Read in the Data*

```
MW_Data = pd.read_csv("Malware-staDyn-data.csv")
```

```
In [17]: # Plots the Pandas Data Frame as a check.  
MW_Data
```

```
Out[17]:
```

	Virtual	Offset	loc	Import	Imports	var	Forwarder	UINT	LONG	BOOL	...	count_file_written	count_file_exists	count_file_deleted
0	0	0	1	0	0	0	0	0	0	0	...	0	0	0
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0
2	0	0	1	0	0	0	0	0	0	0	...	0	0	0
3	0	0	1	0	0	0	0	0	0	0	...	0	0	0
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0
...
6243	0	0	0	0	0	0	0	0	0	0	...	0	0	0
6244	0	0	58	0	0	0	0	0	0	0	...	0	4	0
6245	0	0	0	0	0	0	0	0	0	0	...	0	0	0
6246	0	0	0	0	0	0	0	0	0	0	...	0	0	0
6247	0	0	0	0	0	0	0	0	0	0	...	0	0	0

6248 rows × 1085 columns

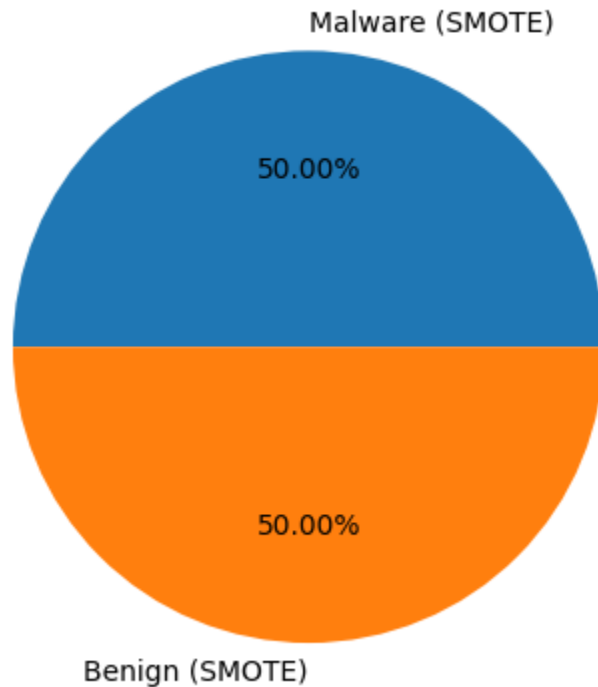
```
In [18]: # Drops all Missing values in Label and selects the Feature Column  
x = MW_Data.drop('label', axis = 1)  
y = MW_Data["label"]  
  
# Split into Training and Test Sets  
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.1, random_state=42)
```

```
In [19]: smote = SMOTE(random_state=42)  
x_train_resampled, y_train_resampled = smote.fit_resample(x_train, y_train)  
  
# Concatenate resampled data back into a dataframe  
MW_Data_resampled = pd.concat([pd.DataFrame(x_train_resampled, columns=MW_Data.columns[:-1]), pd.DataFrame(y_train_resampled, colu  
  
# Counts number of 0 or 1 occurrences within the data. Will be used in the Pie chart  
MW_Data_count = MW_Data_resampled["label"].value_counts()  
print(MW_Data_count)  
Not_Malware_count = MW_Data_count[1]  
Malware_count = MW_Data_count[0]
```

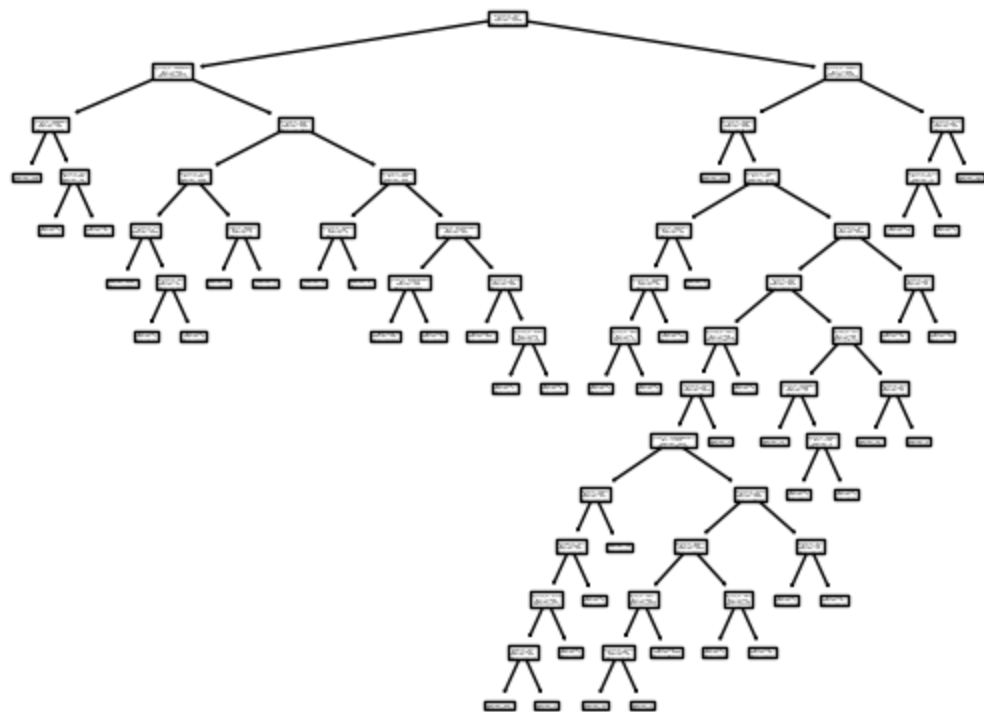
```
# Code for Pie chart.
pie_labels = "Malware (SMOTE)", "Benign (SMOTE)"
pie_sizing = [int(Malware_count), int(Not_Malware_count)]

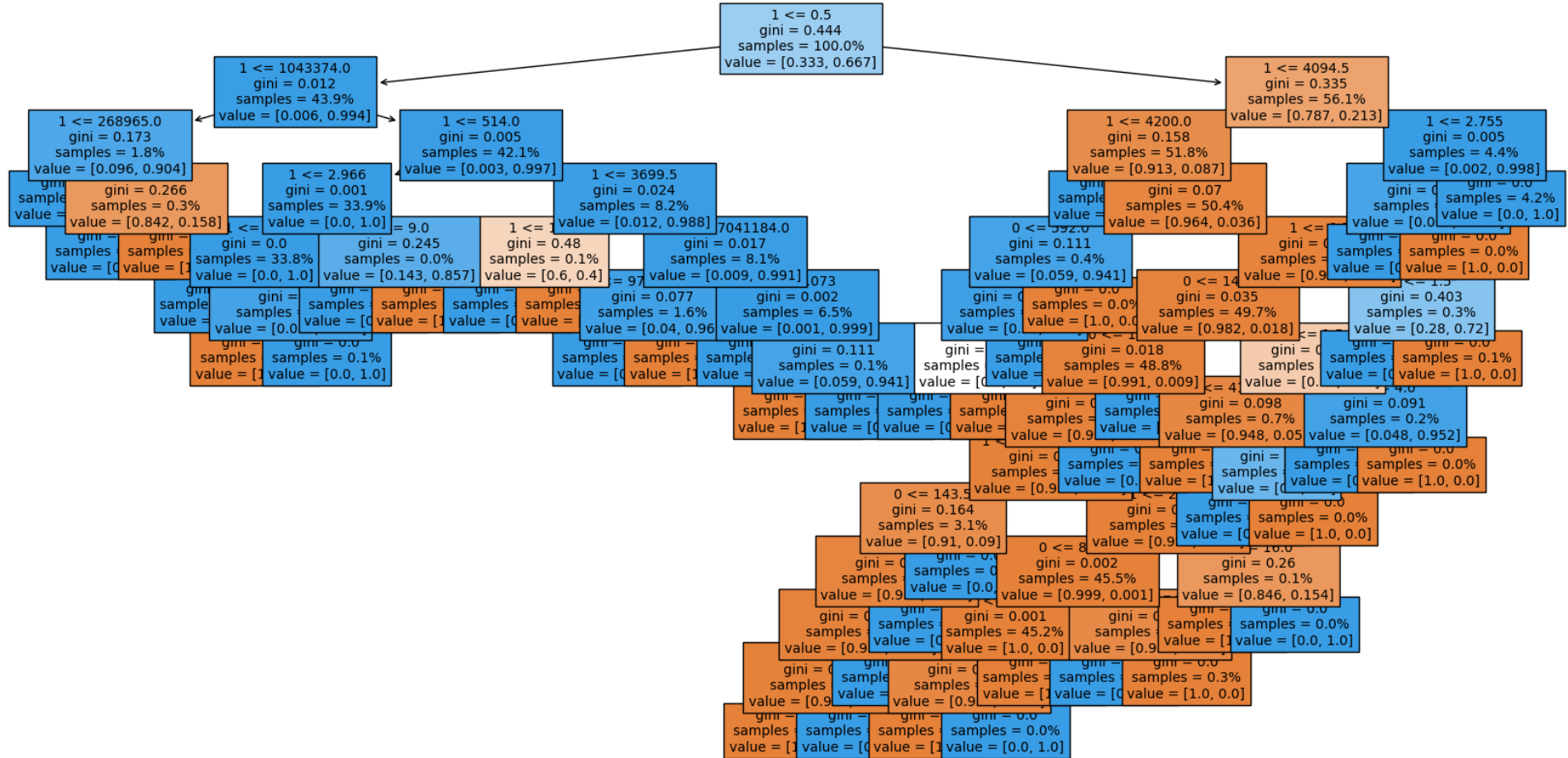
fig, ax = plt.subplots()
ax.pie(pie_sizing, labels=pie_labels, autopct='%1.2f%%')
plt.show()
```

```
label
1    5101
0    5101
Name: count, dtype: int64
```



```
In [20]: dTree3 = DecisionTreeClassifier(class_weight= {1:2, 0:1})
dTree3.fit(x_train_resampled, y_train_resampled)
tree.plot_tree(dTree3);
dTree3.fit(x_train_resampled, y_train_resampled)
plt.figure(figsize=(20, 10))
plot_tree(dTree3, filled=True, feature_names=MW_Data["label"], proportion=True, fontsize=10)
plt.show()
```





```
In [21]: y_pred = dTree3.predict(x_test)
y_pred_proba = dTree3.predict_proba(x_test)[: , 1]
```

```
In [22]: # Prints out Relevant Metrics
print("Accuracy: ", accuracy_score(y_test, y_pred))

print("AUC-ROC:", roc_auc_score(y_test, y_pred_proba))

print("F1 Score: ", f1_score(y_test, y_pred))

print("Precision:", precision_score(y_test, y_pred))

print("Recall: ", recall_score(y_test, y_pred))
```

Accuracy: 0.9952
AUC-ROC: 0.9853955727615644
F1 Score: 0.9972850678733032
Precision: 0.9963833634719711
Recall: 0.9981884057971014

