

In [1]:

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

In [2]:

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

age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp	diabetes	totChol	sysBl
39	4.0	0	0.0	0.0	0	0	0	195.0	106.
46	2.0	0	0.0	0.0	0	0	0	250.0	121.
48	1.0	1	20.0	0.0	0	0	0	245.0	127.
61	3.0	1	30.0	0.0	0	1	0	225.0	150.
46	3.0	1	23.0	0.0	0	0	0	285.0	130.
...
50	1.0	1	1.0	0.0	0	1	0	313.0	179.
51	3.0	1	43.0	0.0	0	0	0	207.0	126.
48	2.0	1	20.0	NaN	0	0	0	248.0	131.
44	1.0	1	15.0	0.0	0	0	0	210.0	126.
52	2.0	0	0.0	0.0	0	0	0	269.0	133.

In [3]:

```
df.columns
```

Out[3]:

```
Index(['male', 'age', 'education', 'currentSmoker', 'cigsPerDay', 'BPMed
s',
      'prevalentStroke', 'prevalentHyp', 'diabetes', 'totChol', 'sysBP',
      'diaBP', 'BMI', 'heartRate', 'glucose', 'TenYearCHD'],
      dtype='object')
```

In [4]:

`df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4238 entries, 0 to 4237
Data columns (total 16 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   male                  4238 non-null   int64
 1   age                   4238 non-null   int64
 2   education             4133 non-null   float64
 3   currentSmoker         4238 non-null   int64
 4   cigsPerDay            4209 non-null   float64
 5   BPMeds                4185 non-null   float64
 6   prevalentStroke       4238 non-null   int64
 7   prevalentHyp          4238 non-null   int64
 8   diabetes              4238 non-null   int64
 9   totChol               4188 non-null   float64
10   sysBP                 4238 non-null   float64
11   diaBP                 4238 non-null   float64
12   BMI                   4219 non-null   float64
13   heartRate             4237 non-null   float64
14   glucose               3850 non-null   float64
15   TenYearCHD            4238 non-null   int64
dtypes: float64(9), int64(7)
memory usage: 529.9 KB
```

`df['_glucose'].value_counts()`

In [6]:

```
x=df[['male', 'age', 'currentSmoker', 'prevalentStroke', 'prevalentHyp', 'diabetes']]
y=df['TenYearCHD']
```

`d={"Verified":{"True":1,'False ':2}} df=df.replace(df) print(df)`

In [7]:

```
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.70)
```

In [8]:

```
from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

Out[8]:

`RandomForestClassifier()`

Depth of Tree

In [9]:

```
parameters={"max_depth":[1,2,3,4,5],"min_samples_leaf":[5,23,45,76,78],'n_estimators':[10
```

Cross Validate

In [10]:

```
from sklearn.model_selection import GridSearchCV
grid_search=GridSearchCV(estimator=rfc,param_grid=parameters,cv=2,scoring="accuracy")
grid_search.fit(x_train,y_train)
```

Out[10]:

```
GridSearchCV(cv=2, estimator=RandomForestClassifier(),
             param_grid={'max_depth': [1, 2, 3, 4, 5],
                          'min_samples_leaf': [5, 23, 45, 76, 78],
                          'n_estimators': [10, 23, 45, 65, 7]}},
             scoring='accuracy')
```

Score

In [11]:

```
grid_search.best_score_
```

Out[11]:

```
0.8544458475709404
```

In [12]:

```
rfc_best=grid_search.best_estimator_
```

In [13]:

```
from sklearn.tree import plot_tree
plt.figure(figsize=(80,40))
plot_tree(rfc_best.estimators_[5],feature_names=x.columns,class_names=[ 'Yes', 'No'],filled
```

Out[13]:

[illegible]