

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('cars.csv')
df
```

	Unnamed: 0	model	year	price	transmission	mileage	fuelType	tax	mpg	engineSize	Make
0	0	T-Roc	2019	25000	Automatic	13904	Diesel	145	49.6	2.0	VW
1	1	T-Roc	2019	26883	Automatic	4562	Diesel	145	49.6	2.0	VW
2	2	T-Roc	2019	20000	Manual	7414	Diesel	145	50.4	2.0	VW
3	3	T-Roc	2019	33492	Automatic	4825	Petrol	145	32.5	2.0	VW
4	4	T-Roc	2019	22900	Semi-Auto	6500	Petrol	150	39.8	1.5	VW
...
99182	10663	A3	2020	16999	Manual	4018	Petrol	145	49.6	1.0	Audi
99183	10664	A3	2020	16999	Manual	1978	Petrol	150	49.6	1.0	Audi
99184	10665	A3	2020	17199	Manual	609	Petrol	150	49.6	1.0	Audi
99185	10666	Q3	2017	19499	Automatic	8646	Petrol	150	47.9	1.4	Audi
99186	10667	Q3	2016	15999	Manual	11855	Petrol	150	47.9	1.4	Audi

In [3]:

```
df.columns
```

Out[3]:

```
Index(['Unnamed: 0', 'model', 'year', 'price', 'transmission', 'mileage',
      'fuelType', 'tax', 'mpg', 'engineSize', 'Make'],
      dtype='object')
```

In [4]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 99187 entries, 0 to 99186
Data columns (total 11 columns):
 #   Column          Non-Null Count  Dtype  
---  -
 0   Unnamed: 0      99187 non-null  int64  
 1   model           99187 non-null  object  
 2   year            99187 non-null  int64  
 3   price           99187 non-null  int64  
 4   transmission    99187 non-null  object  
 5   mileage         99187 non-null  int64  
 6   fuelType        99187 non-null  object  
 7   tax             99187 non-null  int64  
 8   mpg             99187 non-null  float64 
 9   engineSize      99187 non-null  float64 
10  Make            99187 non-null  object  
dtypes: float64(2), int64(5), object(4)
memory usage: 8.3+ MB
```

In [8]:

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

Out[8]:

```
Petrol      54928
Diesel      40928
Hybrid       3078
Other         247
Electric         6
Name: fuelType, dtype: int64
```

In [9]:

```
x=df[['Unnamed: 0', 'year', 'price', 'mileage', 'tax']]
y=df['fuelType']
```

In [*]:

```
d={"fuelType":{"Petrol":1, 'Diesel':2, 'Hybrid':3, 'Other':4, 'Electric':5}}
df=df.replace(df)
print(df)
```

In [24]:

```
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 [25]:

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

Out[25]:

```
RandomForestClassifier()
```

Depth of Tree

In [26]:

```
parameters={"max_depth":[1,2,3,4,5],"min_samples_leaf":[5,23,45,76,78],'n_estimators':[10,20,30,40,50]}
```

Cross Validate

In [27]:

```
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[27]:

```
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 [28]:

```
grid_search.best_score_
```

Out[28]:

```
0.6242959549411162
```

In [29]:

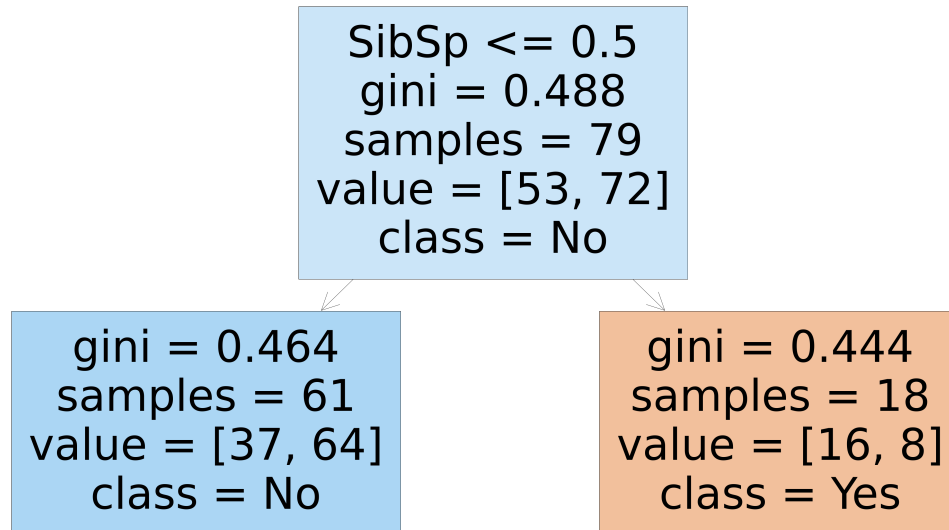
```
rfc_best=grid_search.best_estimator_
```

In [30]:

```
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[30]:

```
[Text(2232.0, 1630.8000000000002, 'SibSp <= 0.5\n'gini = 0.488\n'nsamples = 79\n'nvalue = [53, 72]\n'nclass = No'),
 Text(1116.0, 543.5999999999999, 'gini = 0.464\n'nsamples = 61\n'nvalue = [37, 64]\n'nclass = No'),
 Text(3348.0, 543.5999999999999, 'gini = 0.444\n'nsamples = 18\n'nvalue = [16, 8]\n'nclass = Yes')]
```



In []: