

In [72]:

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

In [73]:

```
df=pd.read_csv(r"C:\Users\prajapath Arjun\Downloads\insurance.csv")
print(df)
```

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.92400
1	18	male	33.770	1	no	southeast	1725.55230
2	28	male	33.000	3	no	southeast	4449.46200
3	33	male	22.705	0	no	northwest	21984.47061
4	32	male	28.880	0	no	northwest	3866.85520
...
1333	50	male	30.970	3	no	northwest	10600.54830
1334	18	female	31.920	0	no	northeast	2205.98080
1335	18	female	36.850	0	no	southeast	1629.83350
1336	21	female	25.800	0	no	southwest	2007.94500
1337	61	female	29.070	0	yes	northwest	29141.36030

[1338 rows x 7 columns]

In [74]:

```
df.head()
```

Out[74]:

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	yes	southwest	16884.92400
1	18	male	33.770	1	no	southeast	1725.55230
2	28	male	33.000	3	no	southeast	4449.46200
3	33	male	22.705	0	no	northwest	21984.47061
4	32	male	28.880	0	no	northwest	3866.85520

In [75]:

```
df.tail()
```

Out[75]:

	age	sex	bmi	children	smoker	region	charges
1333	50	male	30.97	3	no	northwest	10600.5483
1334	18	female	31.92	0	no	northeast	2205.9808
1335	18	female	36.85	0	no	southeast	1629.8335
1336	21	female	25.80	0	no	southwest	2007.9450
1337	61	female	29.07	0	yes	northwest	29141.3603

In [76]:

```
df.shape
```

Out[76]:

(1338, 7)

In [77]:

```
df.describe()
```

Out[77]:

	age	bmi	children	charges
count	1338.000000	1338.000000	1338.000000	1338.000000
mean	39.207025	30.663397	1.094918	13270.422265
std	14.049960	6.098187	1.205493	12110.011237
min	18.000000	15.960000	0.000000	1121.873900
25%	27.000000	26.296250	0.000000	4740.287150
50%	39.000000	30.400000	1.000000	9382.033000
75%	51.000000	34.693750	2.000000	16639.912515
max	64.000000	53.130000	5.000000	63770.428010

In [78]:

```
df.isnull().sum()
```

Out[78]:

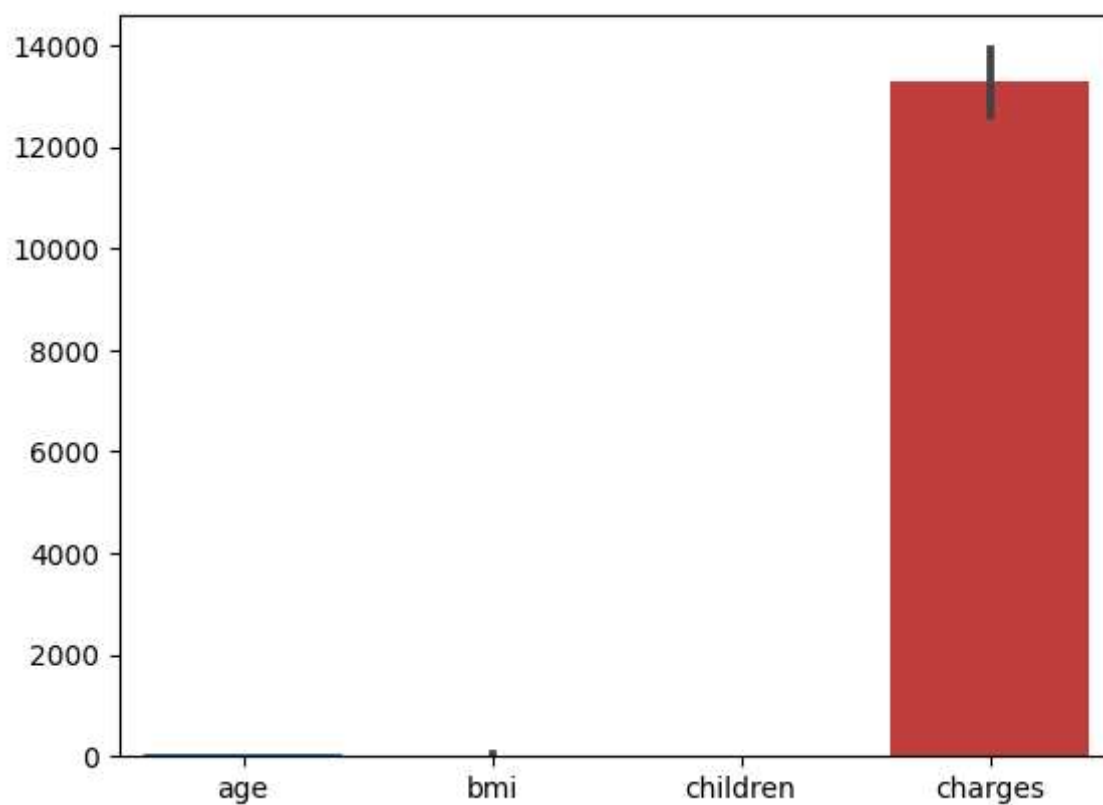
```
age      0
sex      0
bmi      0
children 0
smoker   0
region   0
charges  0
dtype: int64
```

In [79]:

```
#Exploratory Data Analysis  
sns.barplot(df)
```

Out[79]:

<Axes: >



In [80]:

```
df.columns
```

Out[80]:

```
Index(['age', 'sex', 'bmi', 'children', 'smoker', 'region', 'charges'], dtype='object')
```

In [81]:

```
smoker={"smoker":{"yes":1,"no":0}}
df=df.replace(smoker)
df
```

Out[81]:

	age	sex	bmi	children	smoker	region	charges
0	19	female	27.900	0	1	southwest	16884.92400
1	18	male	33.770	1	0	southeast	1725.55230
2	28	male	33.000	3	0	southeast	4449.46200
3	33	male	22.705	0	0	northwest	21984.47061
4	32	male	28.880	0	0	northwest	3866.85520
...
1333	50	male	30.970	3	0	northwest	10600.54830
1334	18	female	31.920	0	0	northeast	2205.98080
1335	18	female	36.850	0	0	southeast	1629.83350
1336	21	female	25.800	0	0	southwest	2007.94500
1337	61	female	29.070	0	1	northwest	29141.36030

1338 rows × 7 columns

In [82]:

```
sex={"sex":{"male":1,"female":0}}
df=df.replace(sex)
df
```

Out[82]:

	age	sex	bmi	children	smoker	region	charges
0	19	0	27.900	0	1	southwest	16884.92400
1	18	1	33.770	1	0	southeast	1725.55230
2	28	1	33.000	3	0	southeast	4449.46200
3	33	1	22.705	0	0	northwest	21984.47061
4	32	1	28.880	0	0	northwest	3866.85520
...
1333	50	1	30.970	3	0	northwest	10600.54830
1334	18	0	31.920	0	0	northeast	2205.98080
1335	18	0	36.850	0	0	southeast	1629.83350
1336	21	0	25.800	0	0	southwest	2007.94500
1337	61	0	29.070	0	1	northwest	29141.36030

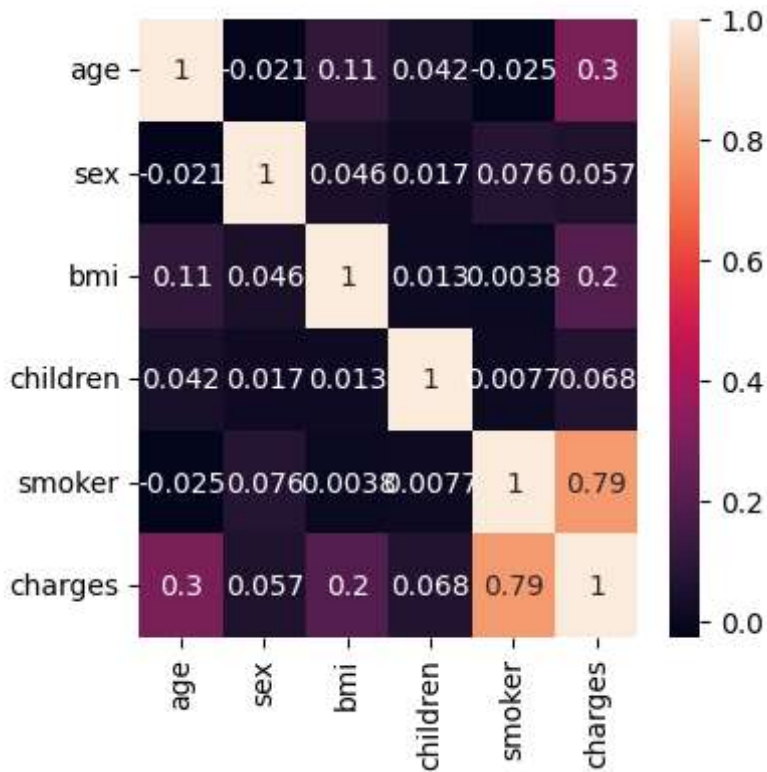
1338 rows × 7 columns

In [83]:

```
idf=df[['age', 'sex', 'bmi', 'children', 'smoker', 'charges']]
plt.figure(figsize=(4,4))
sns.heatmap(idf.corr(),annot=True)
```

Out[83]:

<Axes: >



In [84]:

```
#Training the model
X=df[['age', 'sex', 'bmi', 'children', 'smoker']]
y=df['charges']
```

In [85]:

```
#Linear Regression
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3,random_state=100)
```

In [86]:

```
from sklearn.linear_model import LinearRegression
regr=LinearRegression()
regr.fit(X_train,y_train)
print(regr.intercept_)
coeff_df=pd.DataFrame(regr.coef_,X.columns,columns=['coefficient'])
coeff_df
```

-10719.483493479494

Out[86]:

	coefficient
age	259.757578
sex	18.216925
bmi	277.903898
children	461.169867
smoker	23981.741027

In [87]:

```
score=regr.score(X_test,y_test)
print(score)
```

0.780095696440481

In [88]:

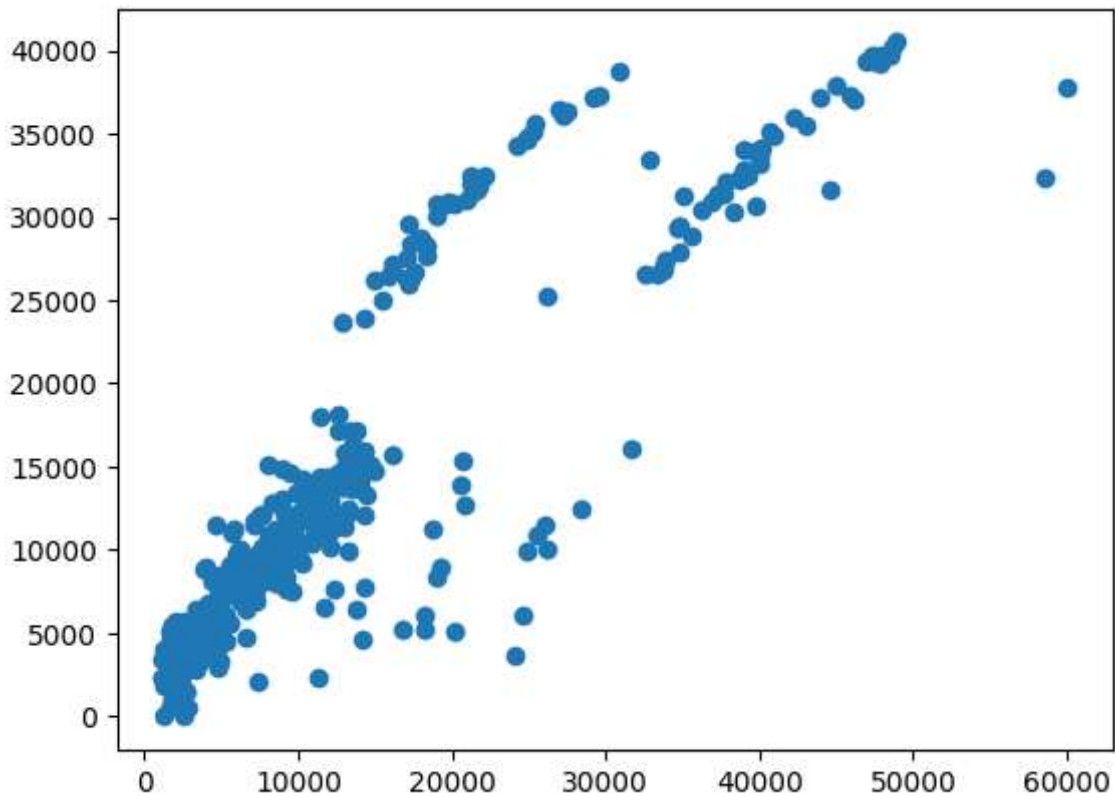
```
predictions=regr.predict(X_test)
```

In [89]:

```
plt.scatter(y_test,predictions)
```

Out[89]:

<matplotlib.collections.PathCollection at 0x2151562ce10>



In [90]:

```
x=np.array(df['smoker']).reshape(-1,1)
y=np.array(df['charges']).reshape(-1,1)
df.dropna(inplace=True)
```

In [91]:

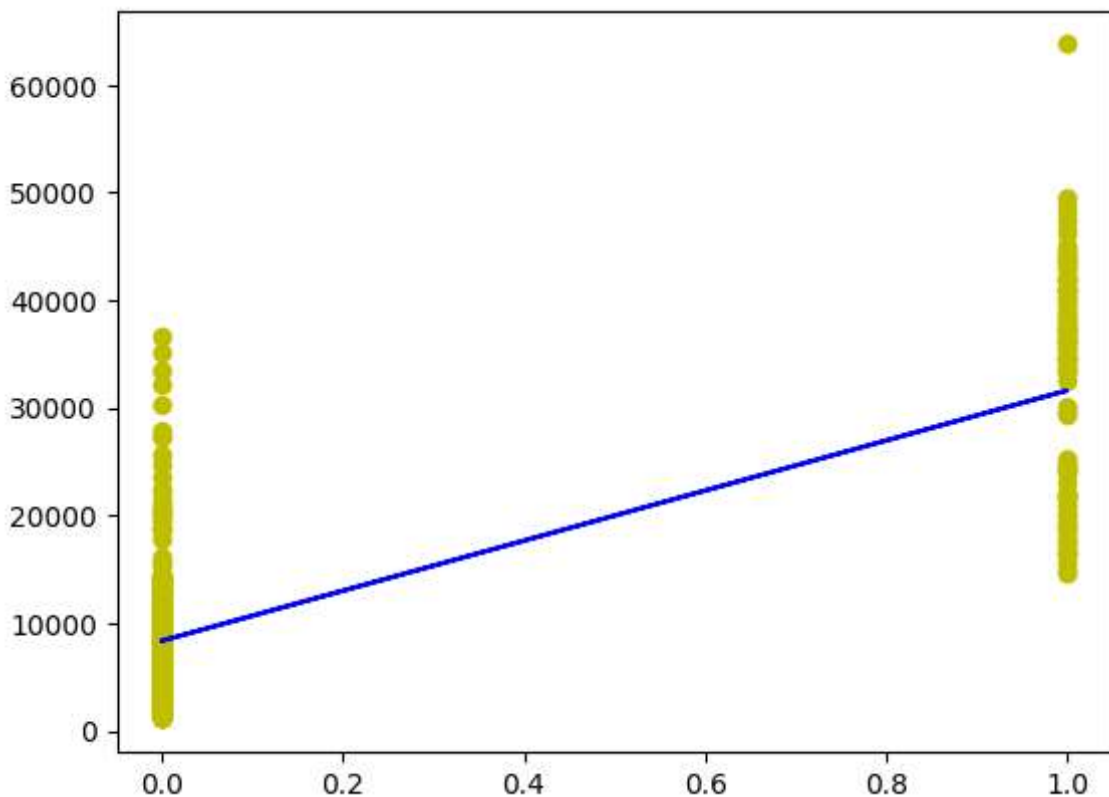
```
X_train,X_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
regr.fit(X_train,y_train)
regr.fit(X_train,y_train)
```

Out[91]:

```
LinearRegression
LinearRegression()
```

In [92]:

```
y_pred=regr.predict(X_test)
plt.scatter(X_test,y_test,color='y')
plt.plot(X_test,y_pred,color='b')
plt.show()
```



In [93]:

```
#Logistic Regression
x=np.array(df['charges']).reshape(-1,1)
y=np.array(df['smoker']).reshape(-1,1)
df.dropna(inplace=True)
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=1)
from sklearn.linear_model import LogisticRegression
lr=LogisticRegression(max_iter=10000)
```

In [94]:

```
lr.fit(x_train,y_train)
```

C:\Users\prajapath Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\utils\validation.py:1143: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

Out[94]:

```
LogisticRegression
LogisticRegression(max_iter=10000)
```


In [95]:

```
score=lr.score(x_test,y_test)
print(score)
```

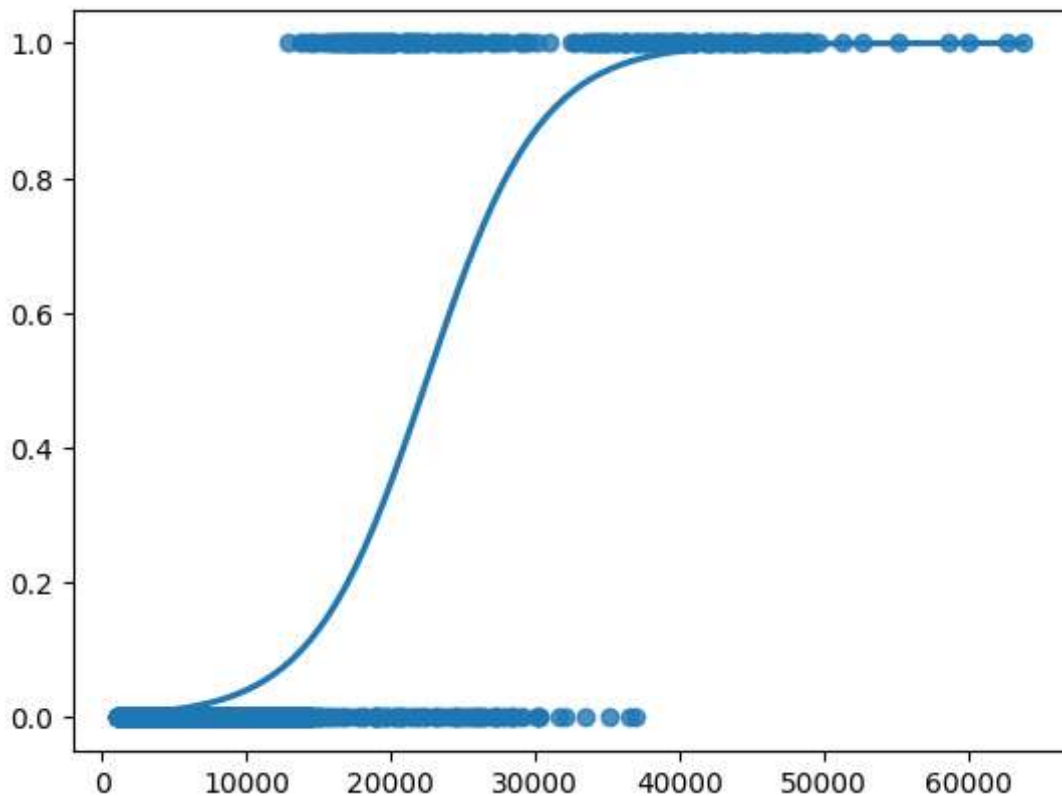
0.8930348258706468

In [96]:

```
sns.regplot(x=x,y=y,data=df,logistic=True,ci=None)
```

Out[96]:

<Axes: >



In [97]:

```
#Decision tree
from sklearn.tree import DecisionTreeClassifier
clf=DecisionTreeClassifier(random_state=0)
clf.fit(x_train,y_train)
```

Out[97]:

```
DecisionTreeClassifier
DecisionTreeClassifier(random_state=0)
```

In [98]:

```
score=clf.score(x_test,y_test)
print(score)
```

0.8880597014925373

In [99]:

```
#Random forest classifier
from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
rfc.fit(X_train,y_train)
```

C:\Users\prajapath Arjun\AppData\Local\Temp\ipykernel_6804\1232785509.py:
4: DataConversionWarning: A column-vector y was passed when a 1d array was
expected. Please change the shape of y to (n_samples,), for example using
ravel().
rfc.fit(X_train,y_train)

Out[99]:

```
▼ RandomForestClassifier
RandomForestClassifier()
```

In [100]:

```
params={'max_depth':[2,3,5,10,20], 'min_samples_leaf':[5,10,20,50,100,200], 'n_estimators'
```

In [108]:

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

Please change the shape of y to (n_samples,), for example using ravel().
estimator.fit(X_train, y_train, **fit_params)

C:\Users\prajapath Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\model_selection_validation.py:686: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

```
estimator.fit(X_train, y_train, **fit_params)
```

C:\Users\prajapath Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\model_selection_validation.py:686: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

```
estimator.fit(X_train, y_train, **fit_params)
```

C:\Users\prajapath Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\model_selection_validation.py:686: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

```
estimator.fit(X_train, y_train, **fit_params)
```

C:\Users\prajapath Arjun\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\model_selection_validation.py:686: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

In [103]:

```
grid_search.best_score_
```

Out[103]:

0.7938034188034188

In [104]:

```
rf_best=grid_search.best_estimator_  
rf_best
```

Out[104]:

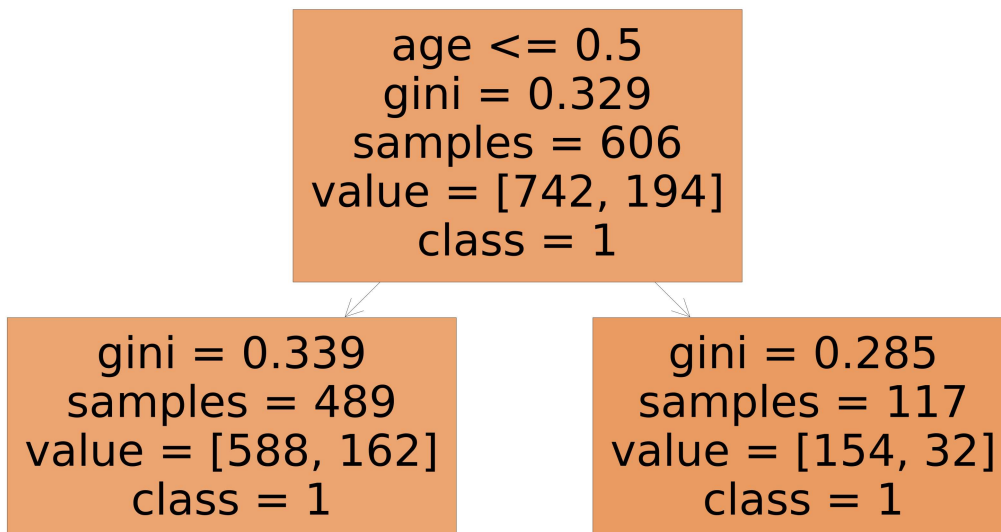
```
RandomForestClassifier  
RandomForestClassifier(max_depth=2, min_samples_leaf=5, n_estimators=10)
```

In [105]:

```
from sklearn.tree import plot_tree  
plt.figure(figsize=(80,40))  
plot_tree(rf_best.estimators_[4],feature_names=X.columns,class_names=['1','0'],filled=True)
```

Out[105]:

```
[Text(0.5, 0.75, 'age <= 0.5\ngini = 0.329\nsamples = 606\nvalue = [742, 1  
94]\nnclass = 1'),  
 Text(0.25, 0.25, 'gini = 0.339\nsamples = 489\nvalue = [588, 162]\nnclass  
= 1'),  
 Text(0.75, 0.25, 'gini = 0.285\nsamples = 117\nvalue = [154, 32]\nnclass =  
1')]
```



In [106]:

```
score=rfc.score(x_test,y_test)  
print(score)
```

0.7985074626865671

In []: