

In [30]:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import confusion_matrix, accuracy_score
from sklearn import tree
from io import StringIO
import pydotplus
from IPython.display import Image

# Load the dataset which was used throughout the whole course
data = pd.read_excel('finally_clean_data_for_plotting.xlsx')

# Show the data frame
print(data)
```

	age	sex	householdincome	howoftenwine	noofwines	income_category
0	34	2	12	10	1	\$50,000 to \$59,999
1	84	2	7	6	1	\$20,000 to \$24,999
2	29	2	13	10	1	\$60,000 to \$69,999
3	68	2	6	5	1	\$15,000 to \$19,999
4	54	2	11	9	1	\$40,000 to \$49,999
...	...	...	...	...	...	...
14556	18	2	1	9	1	Less than \$5,000
14557	18	1	1	10	1	Less than \$5,000
14558	51	1	6	6	1	\$15,000 to \$19,999
14559	21	1	1	10	2	Less than \$5,000
14560	18	2	1	10	1	Less than \$5,000

[14561 rows x 6 columns]

The dataset focusses on the correlations between age, sex, household-income (already pre-categorized according to the NESARC codebook) and wine drinking frequency with the "noofwines" variable, which represents the consumed amount of wine per occasion. In the following, the data pre-processing is performed, which includes binning and grouping of the age and frequency variable, to reduce the complexity of the resulting tree. I previously tried it without the binning and grouping first, but it did not work out.



Interpretation of the results:

The accuracy score of approximately 60.7% indicates that about 60.7% of the total predictions made by the model were correct. Further, the conclusion can be drawn that the model performs well for class 0, with a high number of correct predictions of 3536, but it fails to predict any instances for classes 1 through 9, as indicated by the zeros in those rows. This suggests that the model may be biased towards class 0 or that the other classes are underrepresented.

For better understanding, the class 0 corresponds to 1 unit, class 1 to 2 units of wine, and so on.

Now the tree can be drawn:

In [32]:

```
# Create a PNG image
png_image = graph.create_png()

# Display the image
display(Image(png_image, width=1200, height=2400))
```

