	Assignment is below at the end  • https://scikit-learn.org/stable/modules/tree.html • https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html • https://scikit-learn.org/stable/modules/generated/sklearn.tree.plot_tree.html
In [1]:	<pre>import seaborn as sns import matplotlib.pyplot as plt %matplotlib inline plt.rcParams['figure.figsize'] = (20, 6) plt.rcParams['font.size'] = 14 import pandas as pd  df = pd.read_csv('/data/adult.data', index_col=False)</pre>
<pre>In [3]: In [4]: Out[4]:</pre>	
<pre>In [5]: Out[5]: In [6]:</pre>	<pre>df.head()</pre>
Out[6]:  In [7]:  In [8]:	<pre>Index(['age', 'workclass', 'fnlwgt', 'education', 'education-num',</pre>
In [9]: Out[9]:	non_num_columns = ['workelass', 'education', 'marital-status',
In [10]: Out[10]:	32560 1 0 32561 rows × 2 columns  dummies.shape (32561, 2)
<pre>In [11]: Out[11]:</pre>	
In [12]: Out[12]: In [13]:	<pre>sex = onehot.transform(df[transform_columns]) sex</pre>
Out[13]: In [14]: Out[14]:	array([[0, 1.],
<pre>In [15]: Out[15]:</pre>	<pre>In addition to OneHot encoding there is Ordinal Encoding  enc = preprocessing.OrdinalEncoder() enc.fit(df[["salary"]]) salary = enc.transform(df[["salary"]]) salary  array([[0.],</pre>
<pre>In [16]: Out[16]: In [17]:</pre>	[0.], [0.], [1.]]) enc.categories_[0]
	<pre># transformed = pd.get_dummies(df[transform_columns])  onehot = preprocessing.OneHotEncoder(handle_unknown="infrequent_if_exist", sparse=False).fit(df[transform_columns])  enc = preprocessing.OrdinalEncoder()  enc.fit(df[["salary"]])  transformed = onehot.transform(df[transform_columns])  new_cols = list(onehot.categories_[0].flatten())  df_trans = pd.DataFrame(transformed, columns=new_cols)  x = pd.concat(</pre>
In [18]: Out[18]:	/Users/kevinbrogan/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1. 2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.     warnings.warn(  x.head()
In [19]:	<pre>3 53 234721     7</pre>
<pre>In [20]: Out[20]: In [21]:</pre>	<pre>axis=1,)  xt["salary"] = enc.fit_transform(golden[["salary"]])  xt.salary.value_counts()  0.0    12435 1.0    3846 Name: salary, dtype: int64  enc.categories_</pre>
Out[21]:  In [22]:  In [23]:	<pre>[array([' &lt;=50K.', ' &gt;50K.'], dtype=object)]  from sklearn.tree import DecisionTreeClassifier from sklearn.ensemble import RandomForestClassifier from sklearn.ensemble import GradientBoostingClassifier  Choose the model of your preference: DecisionTree or RandomForest  model = RandomForestClassifier(criterion='entropy')</pre>
	DecisionTreeClassifier(criterion='entropy')  model.treenode_count
Out[26]:  In [27]:  Out[27]:  In [28]:	<pre>list(zip(x.drop(['fnlwgt','salary'], axis=1).columns, model.feature_importances_))  [('age', 0.3224696797613446),     ('education-num', 0.16039663827087378),     ('capital-gain', 0.22779810532320635),     ('capital-loss', 0.07851624804931023),     ('hours-per-week', 0.1555059695817829),     ('Female', 0.001074996873666915),     ('Male', 0.054238362139815366)]</pre>
Out[28]:  In [29]:  Out[29]:	
<pre>In [30]: Out[30]: In [31]: Out[31]:</pre>	
In [32]:	<pre>'capital-gain', 'capital-loss', 'hours-per-week', ' Female', ' Male']  predictions = model.predict(xt.drop(['fnlwgt','salary'], axis=1)) predictionsx = model.predict(x.drop(['fnlwgt','salary'], axis=1))  from sklearn.metrics import (     accuracy_score,     classification_report,     confusion_matrix, auc, roc_curve</pre>
<pre>In [34]: Out[34]: In [35]: Out[35]: In [36]:</pre>	accuracy_score(xt.salary, predictions)
Out[36]: In [37]:	print(classification_report(xt.salary, predictions))  precision recall f1-score support  0.0 0.86 0.92 0.89 12435 1.0 0.66 0.50 0.57 3846  accuracy macro avg 0.76 0.71 0.73 16281
In [38]: In [39]:	<pre>weighted avg</pre>
Out[39]:  In [40]:  Out[40]:  In [41]:	<pre>confusion_matrix(x.salary, predictionsx) array([[24097, 623],</pre>
In [42]:	1.0 0.89 0.65 0.75 7841  accuracy 0.90 32561 macro avg 0.89 0.81 0.84 32561 weighted avg 0.90 0.90 0.89 32561  print(classification_report(x.salary, predictionsx))  precision recall f1-score support  0.0 0.90 0.97 0.93 24720 1.0 0.89 0.65 0.75 7841
In [43]: Out[43]:	<b>0</b> 39 77516 13 2174 0 40 0.0 0.0 1.0
	1 50 83311 13 0 0 13 0.0 0.0 1.0  2 38 215646 9 0 0 0 40 0.0 0.0 1.0  3 53 234721 7 0 0 0 40 0.0 0.0 1.0  4 28 338409 13 0 0 40 0.0 1.0 0.0  x = x.drop([' Female', ' Male', 'fnlwgt'], axis=1) xt = xt.drop([' Female', ' Male', 'fnlwgt'], axis=1)  x.head()  # xt.head()
Out[45]:	age         education-num         capital-gain         capital-loss         hours-per-week         salary           0         39         13         2174         0         40         0.0           1         50         13         0         0         13         0.0           2         38         9         0         0         40         0.0           3         53         7         0         0         40         0.0           4         28         13         0         0         40         0.0
<pre>In [46]: In [47]: Out[47]: In [48]:</pre>	RandomForestClassifier(criterion='entropy')
Out[48]: In [49]: In [50]:	
<pre>In [52]: Out[52]: In [53]:</pre>	<pre>[ 3198 648]] confusion_matrix(xt.salary, predictions_RFC) array([[11570, 865],</pre>
In [54]:	accuracy
In [55]:	
<pre>In [56]: In [57]: Out[57]: In [58]:</pre>	<pre></pre>
<pre>In [59]: In [60]: In [61]: Out[61]:</pre>	<pre>x = pd.concat([x, df_trans],</pre>
<pre>In [62]: In [63]: Out[63]: In [64]:</pre>	<pre>predictions = modelRFC_2.predict(xt.drop(['salary'], axis=1))  confusion_matrix(xt.salary, predictions)  array([[11408, 1027],</pre>
In [65]:	<pre>Second Iteration, marital-status + race  transform_cols = ['race'] onehot = preprocessing.OneHotEncoder(handle_unknown="infrequent_if_exist", sparse=False).fit(df[transform_cols])  /Users/kevinbrogan/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1. 2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.     warnings.warn(</pre>
<pre>In [66]: Out[66]: In [67]:</pre>	<pre>onehot.categories_ [array([' Amer-Indian-Eskimo', ' Asian-Pac-Islander', ' Black', ' Other',</pre>
<pre>In [68]: In [69]: Out[69]:</pre>	
In [70]: In [71]: Out[71]: In [72]:	<pre>RandomForestClassifier(criterion='entropy')  predictions_iter2 = modelRFC_3.predict(xt.drop(['salary'], axis=1))  confusion_matrix(xt.salary, predictions_iter2)  array([[11402, 1033],</pre>
	precision recall f1-score support  0.0 0.88 0.92 0.90 12435 1.0 0.69 0.60 0.64 3846  accuracy 0.84 16281 macro avg 0.79 0.76 0.77 16281 weighted avg 0.84 0.84 0.84 16281  Third Iteration, marital-status + race + workclass
<pre>In [73]: In [74]: Out[74]:</pre>	transform_cols = ['workclass'] onehot = preprocessing.OneHotEncoder(handle_unknown="infrequent_if_exist", sparse=False).fit(df[transform_cols])  /Users/kevinbrogan/anaconda3/lib/python3.10/site-packages/sklearn/preprocessing/_encoders.py:828: FutureWarning: `sparse` was renamed to `sparse_output` in version 1. 2 and will be removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default value.  warnings.warn(  onehot.categories_  [array(4/-2/  Redeval gov/  Local gov/  Nover worked/  Redeval gov/  Redeval gov/  Nover worked/  Redeval gov/
In [75]: In [76]:	
<pre>In [77]: Out[77]: In [78]: In [79]:</pre>	<pre>modelRFC_4 = RandomForestClassifier(criterion='entropy') modelRFC_4.fit(x.drop(['salary'], axis=1), x.salary)</pre>
In [79]: Out[79]: In [80]:	array([[11340, 1095],
In [ ]:	weighted avg 0.83 0.84 0.83 16281