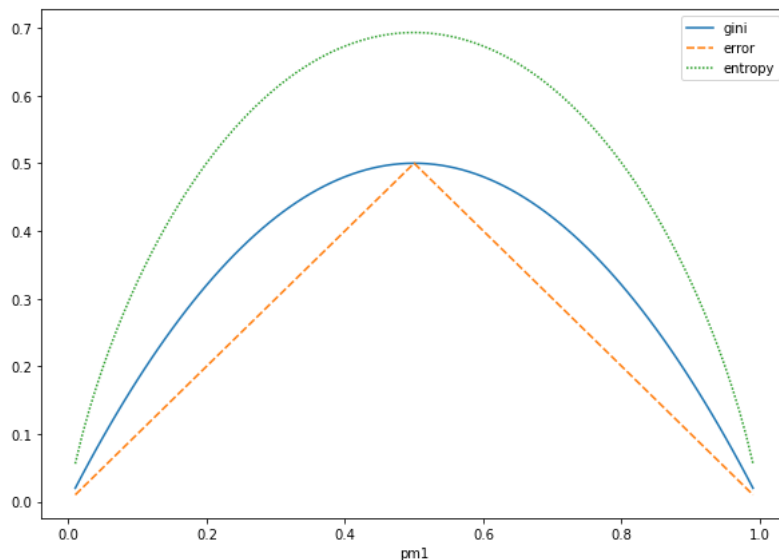


Youssef Mahmoud 905854027

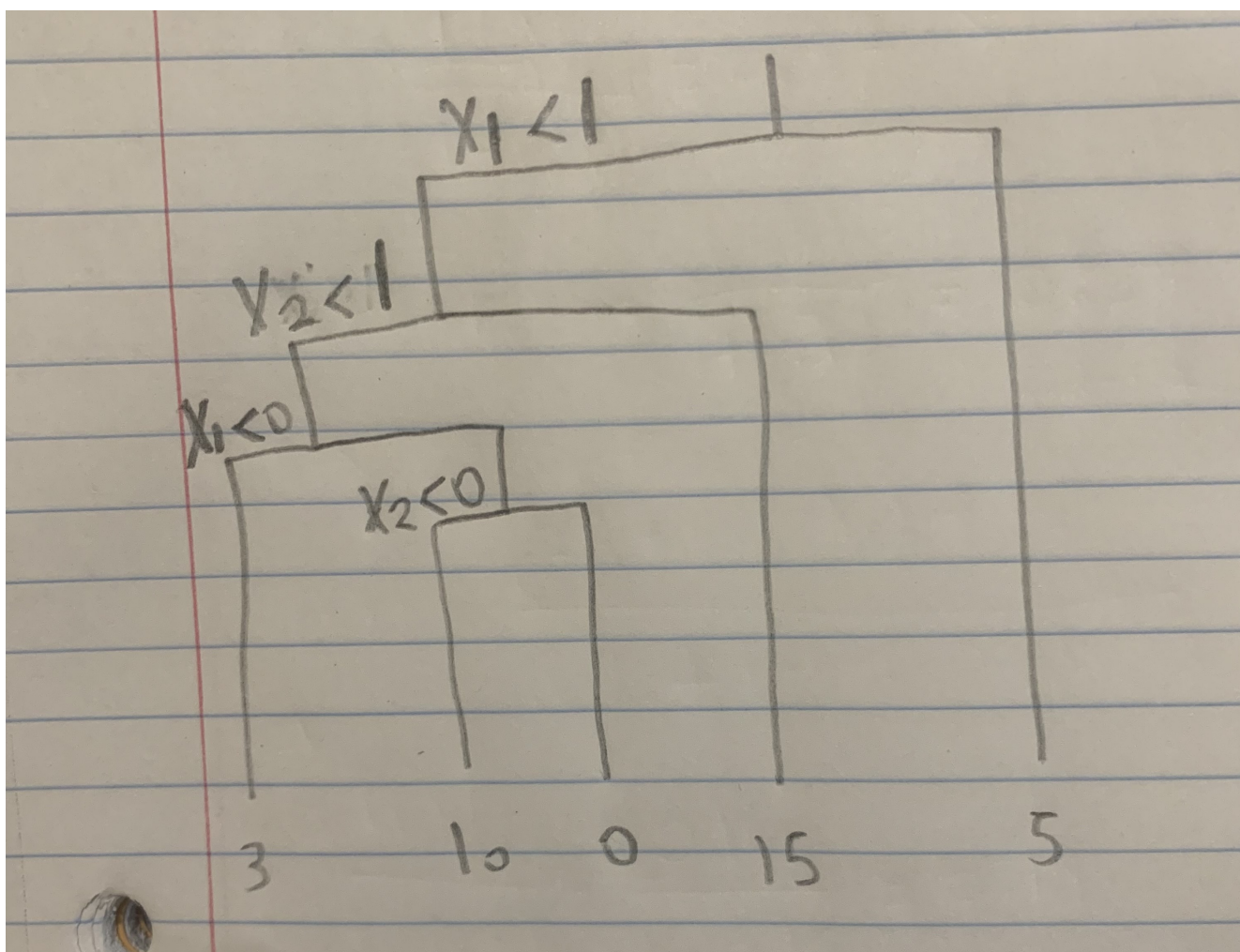
```
In [140]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib as mpl
4 import matplotlib.pyplot as plt
5 import graphviz
6 from sklearn import tree
7 from sklearn.model_selection import cross_val_score
8 import sklearn
9 from sklearn.metrics import accuracy_score, roc_auc_score
10 from sklearn.ensemble import RandomForestClassifier
11 from sklearn.model_selection import train_test_split
12 from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier, export_graphviz, plot_tree
13 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
14 from sklearn.metrics import confusion_matrix, mean_squared_error
15 import seaborn as sns
16 from sklearn.pipeline import Pipeline
17 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
18 from sklearn.ensemble import AdaBoostRegressor
19 from sklearn.tree import DecisionTreeRegressor
20 import xgboost as XGB
21 from sklearn.preprocessing import OneHotEncoder
22 from sklearn.impute import SimpleImputer
23 from sklearn.compose import ColumnTransformer
24 from sklearn.pipeline import Pipeline
25 %matplotlib inline
```

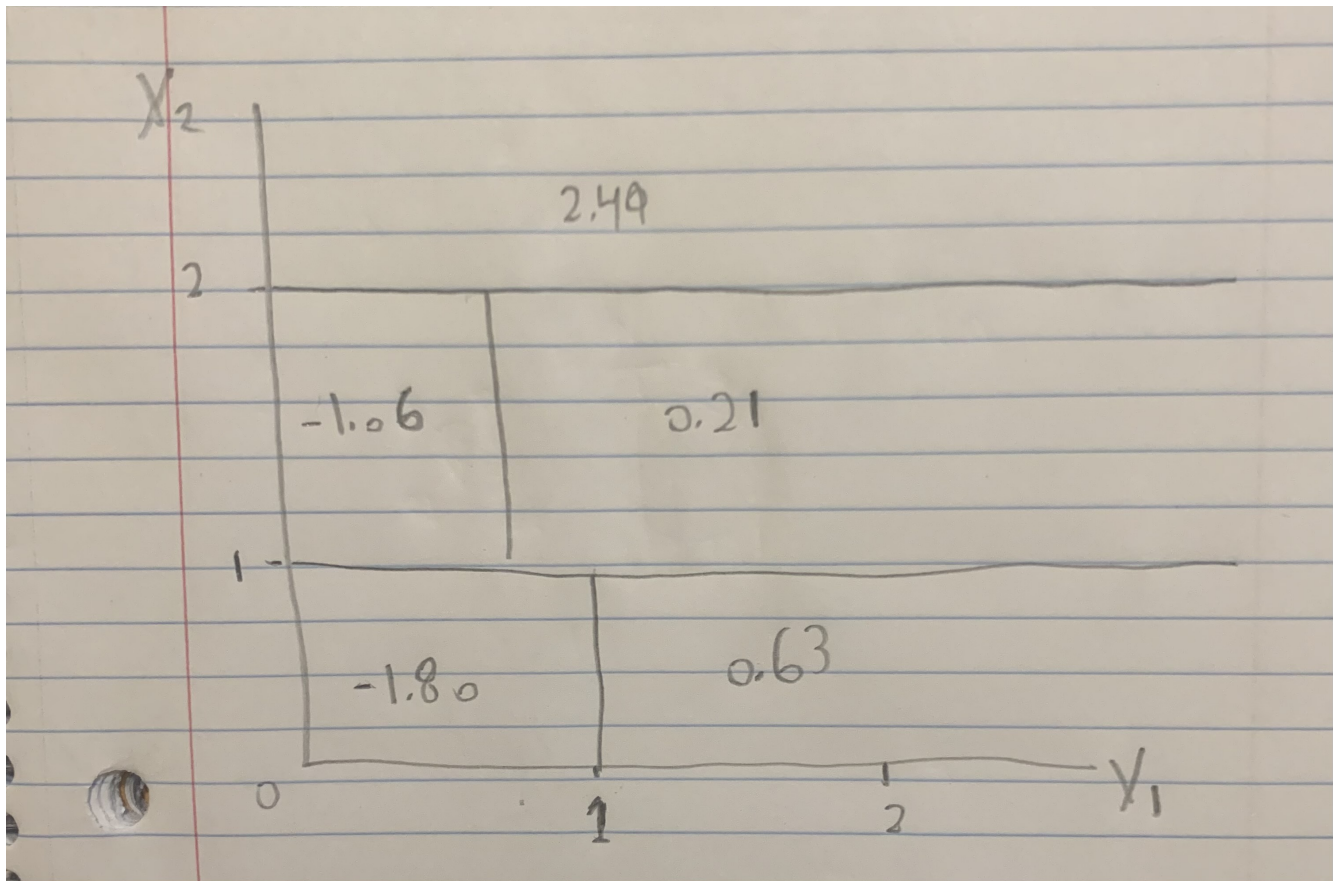
### 8.4.3

```
In [2]: 1 pml = np.arange(0.01, 1, 0.01)
2 pm2 = 1 - pml
3 gini = (pml * (1 - pml)) + (pm2 * (1 - pm2))
4
5
6 err=[]
7 for p in pml:
8     if 0.5 <= p:
9         err += [1 - p]
10    if p < 0.5:
11        err += [p]
12
13
14 entropy = -((pml * np.log(pml)) + (pm2 * np.log(pm2)))
15
16
17
18
19 df = pd.DataFrame(np.stack([pml, gini, err, entropy], axis=1),
20                  columns=['pml', 'gini', 'error', 'entropy']).set_index('pml')
21 plt.figure(figsize=(10, 7))
22 sns.lineplot(data=df);
```



## 8.4.4





### 8.4.5

```
In [3]: 1 x_boot = np.array([0.1,0.15,0.2,0.2,0.55,0.6,0.6,0.65,0.7,0.75])
2
3 def majority_clf(votes):
4     pro_votes = (0.5 < votes).sum()
5     majority_is_pro = (len(votes)/2) < pro_votes
6     return majority_is_pro, pro_votes
7
8 def avg_clf(votes):
9     avg = np.mean(votes)
10    return (0.5 < avg, avg)
11
12 print('Is red (by Majority): {}, votes_red={}'.format(*majority_clf(x_boot)))
13 print('Is red (by Average) : {}, avg={}'.format(*avg_clf(x_boot)))
```

Is red (by Majority): True, votes\_red=6  
Is red (by Average) : False, avg=0.45

### 8.3.Boston

```
In [4]: 1 df1 = pd.read_csv("desktop/housing.csv")
```

In [5]: 1 df1

Out[5]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV
0	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90	5.33	36.2
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
501	0.06263	0.0	11.93	0	0.573	6.593	69.1	2.4786	1	273	21.0	391.99	9.67	22.4
502	0.04527	0.0	11.93	0	0.573	6.120	76.7	2.2875	1	273	21.0	396.90	9.08	20.6
503	0.06076	0.0	11.93	0	0.573	6.976	91.0	2.1675	1	273	21.0	396.90	5.64	23.9
504	0.10959	0.0	11.93	0	0.573	6.794	89.3	2.3889	1	273	21.0	393.45	6.48	22.0
505	0.04741	0.0	11.93	0	0.573	6.030	80.8	2.5050	1	273	21.0	396.90	7.88	11.9

506 rows × 14 columns

In [6]: 1 df1.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 506 entries, 0 to 505
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  ------  -
0   CRIM        506 non-null    float64
1   ZN          506 non-null    float64
2   INDUS       506 non-null    float64
3   CHAS        506 non-null    int64
4   NOX         506 non-null    float64
5   RM          506 non-null    float64
6   AGE         506 non-null    float64
7   DIS         506 non-null    float64
8   RAD         506 non-null    int64
9   TAX         506 non-null    int64
10  PTRATIO     506 non-null    float64
11  B           506 non-null    float64
12  LSTAT       506 non-null    float64
13  MEDV       506 non-null    float64
dtypes: float64(11), int64(3)
memory usage: 55.5 KB
```

```
In [7]: 1 df1_train, df1_test = train_test_split(
2       df1,
3       train_size = 0.5,
4       random_state = 425, # seed
5       )
6 df1_test.shape
```

Out[7]: (253, 14)

```
In [8]: 1 features = ['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT']
2 X_train = df1_train[features]
3 Y_train = np.log(df1_train.MEDV)
4
5 X_test = df1_test[features]
6 Y_test = np.log(df1_test.MEDV)
```

## Regression Tree

```
In [9]: 1 regtree_mod = DecisionTreeRegressor(random_state = 425)
```

```
In [10]: 1 pipe = Pipeline(steps = [
2       ("model", regtree_mod)
3       ])
4 pipe
```

Out[10]: Pipeline(steps=[('model', DecisionTreeRegressor(random\_state=425))])

```
In [11]: 1 ccp_alpha_grid = np.linspace(start = 0.0, stop = 0.1, num = 100)
2 tuned_parameters = {
3     "model__ccp_alpha": ccp_alpha_grid
4 }
5 tuned_parameters
```

```
Out[11]: {'model__ccp_alpha': array([0.          , 0.0010101 , 0.0020202 , 0.0030303 , 0.0040404 ,
0.00505051, 0.00606061, 0.00707071, 0.00808081, 0.00909091,
0.01010101, 0.01111111, 0.01212121, 0.01313131, 0.01414141,
0.01515152, 0.01616162, 0.01717172, 0.01818182, 0.01919192,
0.02020202, 0.02121212, 0.02222222, 0.02323232, 0.02424242,
0.02525253, 0.02626263, 0.02727273, 0.02828283, 0.02929293,
0.03030303, 0.03131313, 0.03232323, 0.03333333, 0.03434343,
0.03535354, 0.03636364, 0.03737374, 0.03838384, 0.03939394,
0.04040404, 0.04141414, 0.04242424, 0.04343434, 0.04444444,
0.04545455, 0.04646465, 0.04747475, 0.04848485, 0.04949495,
0.05050505, 0.05151515, 0.05252525, 0.05353535, 0.05454545,
0.05555556, 0.05656566, 0.05757576, 0.05858586, 0.05959596,
0.06060606, 0.06161616, 0.06262626, 0.06363636, 0.06464646,
0.06565657, 0.06666667, 0.06767677, 0.06868687, 0.06969697,
0.07070707, 0.07171717, 0.07272727, 0.07373737, 0.07474747,
0.07575758, 0.07676768, 0.07777778, 0.07878788, 0.07979798,
0.08080808, 0.08181818, 0.08282828, 0.08383838, 0.08484848,
0.08585859, 0.08686869, 0.08787879, 0.08888889, 0.08989899,
0.09090909, 0.09191919, 0.09292929, 0.09393939, 0.09494949,
0.0959596 , 0.0969697 , 0.0979798 , 0.0989899 , 0.1          ]))
```

```
In [12]: 1 n_folds = 6
2 search = GridSearchCV(
3     pipe,
4     tuned_parameters,
5     cv = n_folds,
6     scoring = "neg_root_mean_squared_error",
7     # Refit the best model on the whole data set
8     refit = True
9 )
```

```
In [13]: 1 search.fit(X_train, Y_train)
```

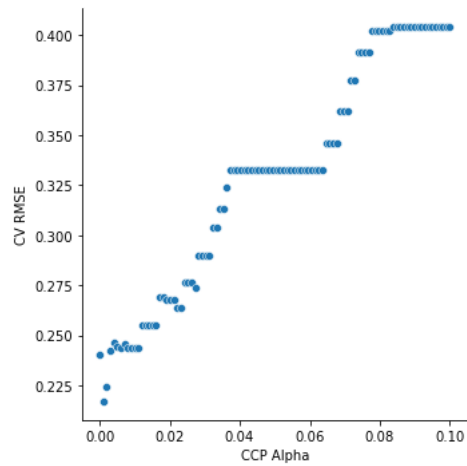
```
Out[13]: GridSearchCV(cv=6,
    estimator=Pipeline(steps=[('model',
    DecisionTreeRegressor(random_state=425))]),
    param_grid={'model__ccp_alpha': array([0.          , 0.0010101 , 0.0020202 , 0.0030303 , 0.0040404 ,
0.00505051, 0.00606061, 0.00707071, 0.00808081, 0.00909091,
0.01010101, 0.01111111, 0.01212121, 0.01313131, 0.01414141,
0.01515152, 0.01616162, 0.01717172, 0.01818182, 0.01919192,
0.02020202...
0.07070707, 0.07171717, 0.07272727, 0.07373737, 0.07474747,
0.07575758, 0.07676768, 0.07777778, 0.07878788, 0.07979798,
0.08080808, 0.08181818, 0.08282828, 0.08383838, 0.08484848,
0.08585859, 0.08686869, 0.08787879, 0.08888889, 0.08989899,
0.09090909, 0.09191919, 0.09292929, 0.09393939, 0.09494949,
0.0959596 , 0.0969697 , 0.0979798 , 0.0989899 , 0.1          ]))},
    scoring='neg_root_mean_squared_error')
```

```

In [14]: 1 cv_res = pd.DataFrame({
2         "ccp_alpha": np.array(search.cv_results_["param_model__ccp_alpha"]),
3         "rmse": -search.cv_results_["mean_test_score"]
4     })
5
6 plt.figure()
7 sns.relplot(
8     data = cv_res,
9     x = "ccp_alpha",
10    y = "rmse"
11 ).set(
12     xlabel = "CCP Alpha",
13     ylabel = "CV RMSE"
14 );
15 plt.show()

```

<Figure size 432x288 with 0 Axes>



```

In [15]: 1 -search.best_score_
2

```

Out[15]: 0.2170777866126541

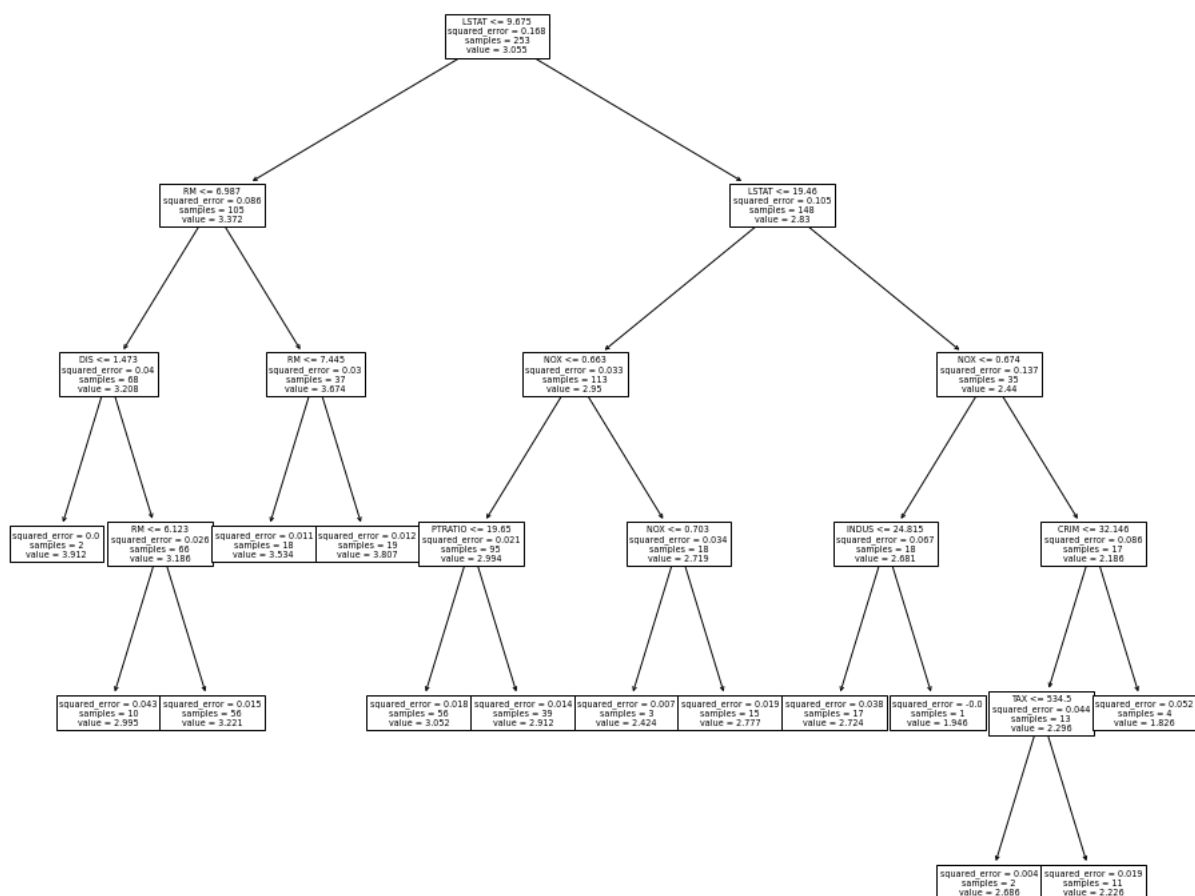
```

In [16]: 1 search.best_estimator_

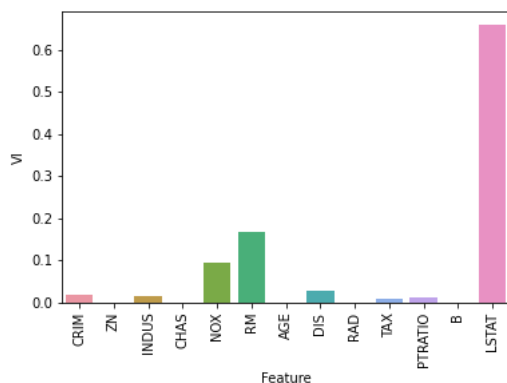
```

Out[16]: Pipeline(steps=[('model',  
DecisionTreeRegressor(ccp\_alpha=0.00101010101010101,  
random\_state=425))])

```
In [17]: 1 plt.figure(figsize=(16,14))
2 sklearn.tree.plot_tree(
3     search.best_estimator_['model'],
4     feature_names = features
5 );
6 plt.show()
```



```
In [18]: 1 vi_df = pd.DataFrame({
2         "feature": features,
3         "vi": search.best_estimator_['model'].feature_importances_
4     })
5
6 plt.figure()
7 sns.barplot(
8     data = vi_df,
9     x = "feature",
10    y = "vi"
11 ).set(
12     xlabel = "Feature",
13     ylabel = "VI"
14 );
15 plt.xticks(rotation = 90);
16 plt.show()
```



```
In [19]: 1 mean_squared_error(
2         y_test,
3         search.best_estimator_.predict(X_test),
4         squared = False
5     )
```

Out[19]: 0.21790720977320982

## Random Forest

```
In [20]: 1 rf_mod = RandomForestRegressor(
2         # Number of trees
3         n_estimators = 100,
4         criterion = 'squared_error',
5         # Number of features to use in each split
6         max_features = 'sqrt',
7         oob_score = True,
8         random_state = 425
9     )
```

```
In [21]: 1 pipe = Pipeline(steps = [
2         ("model", rf_mod)
3     ])
4 pipe
```

Out[21]: Pipeline(steps=[('model',  
RandomForestRegressor(max\_features='sqrt', oob\_score=True,  
random\_state=425))])

```
In [22]: 1 B_grid = [50, 100, 150, 200, 250, 300]
2 m_grid = ['sqrt', 'log2', 1.0] # max_features = 1.0 uses all features
3 tuned_parameters = {
4     "model__n_estimators": B_grid,
5     "model__max_features": m_grid
6 }
7 tuned_parameters
```

Out[22]: {'model\_\_n\_estimators': [50, 100, 150, 200, 250, 300],  
'model\_\_max\_features': ['sqrt', 'log2', 1.0]}



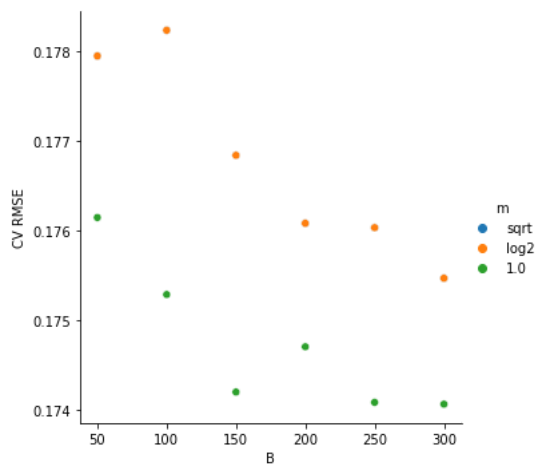
```
In [23]: 1 # Set up CV
2 n_folds = 6
3 search = GridSearchCV(
4     pipe,
5     tuned_parameters,
6     cv = n_folds,
7     scoring = "neg_root_mean_squared_error",
8     # Refit the best model on the whole data set
9     refit = True
10 )
```

```
In [24]: 1 search.fit(X_train,Y_train)
```

```
Out[24]: GridSearchCV(cv=6,
                      estimator=Pipeline(steps=[('model',
                                                  RandomForestRegressor(max_features='sqrt',
                                                                      oob_score=True,
                                                                      random_state=425))]),
                      param_grid={'model__max_features': ['sqrt', 'log2', 1.0],
                                  'model__n_estimators': [50, 100, 150, 200, 250, 300]},
                      scoring='neg_root_mean_squared_error')
```

```
In [25]: 1 cv_res = pd.DataFrame({
2     "B": np.array(search.cv_results_["param_model__n_estimators"]),
3     "rmse": -search.cv_results_["mean_test_score"],
4     "m": search.cv_results_["param_model__max_features"]
5 })
6
7 plt.figure()
8 sns.relplot(
9     # kind = "line",
10    data = cv_res,
11    x = "B",
12    y = "rmse",
13    hue = "m",
14    ).set(
15        xlabel = "B",
16        ylabel = "CV RMSE"
17    );
18 plt.show()
```

<Figure size 432x288 with 0 Axes>



```
In [26]: 1 -search.best_score_
2
```

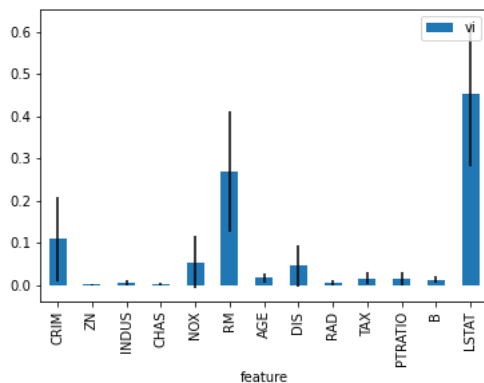
```
Out[26]: 0.17406330114074656
```

```
In [27]: 1 search.best_estimator_
2
```

```
Out[27]: Pipeline(steps=[('model',
                          RandomForestRegressor(max_features=1.0, n_estimators=300,
                                              oob_score=True, random_state=425))])
```

```
In [28]: 1 vi_df = pd.DataFrame({
2         "feature": features,
3         "vi": search.best_estimator_['model'].feature_importances_,
4         "vi_std": np.std([tree.feature_importances_ for tree in search.best_estimator_['model'].estimators_], axis = 0)
5     })
6
7 plt.figure()
8 vi_df.plot.bar(x = "feature", y = "vi", yerr = "vi_std")
9 plt.xticks(rotation = 90);
10 plt.show()
```

<Figure size 432x288 with 0 Axes>



```
In [29]: 1 mean_squared_error(
2         y_test,
3         search.best_estimator_.predict(X_test),
4         squared = False
5     )
```

Out[29]: 0.153092179332972

## Boosting

```
In [30]: 1 bst_mod = AdaBoostRegressor(
2         # Default base estimator is DecisionTreeRegressor with max_depth = 3
3         base_estimator = DecisionTreeRegressor(max_depth = 3),
4         # Number of trees (to be tuned)
5         n_estimators = 50,
6         # Learning rate (to be tuned)
7         learning_rate = 1.0,
8         random_state = 425
9     )
```

```
In [31]: 1 pipe = Pipeline(steps = [
2         ("model", bst_mod)
3     ])
4 pipe
```

Out[31]: Pipeline(steps=[('model',  
AdaBoostRegressor(base\_estimator=DecisionTreeRegressor(max\_depth=3),  
random\_state=425))])

```
In [32]: 1 d_grid = [
2         DecisionTreeRegressor(max_depth = 1),
3         DecisionTreeRegressor(max_depth = 2),
4         DecisionTreeRegressor(max_depth = 3),
5         DecisionTreeRegressor(max_depth = 4)
6     ]
7 B_grid = [50, 100, 150, 200, 250, 300, 350, 400]
8 lambda_grid = [0.2, 0.4, 0.6, 0.8, 1.0]
9 tuned_parameters = {
10     "model__base_estimator": d_grid,
11     "model__n_estimators": B_grid,
12     "model__learning_rate": lambda_grid
13 }
14 tuned_parameters
```

```
Out[32]: {'model__base_estimator': [DecisionTreeRegressor(max_depth=1),
DecisionTreeRegressor(max_depth=2),
DecisionTreeRegressor(max_depth=3),
DecisionTreeRegressor(max_depth=4)],
'model__n_estimators': [50, 100, 150, 200, 250, 300, 350, 400],
'model__learning_rate': [0.2, 0.4, 0.6, 0.8, 1.0]}
```

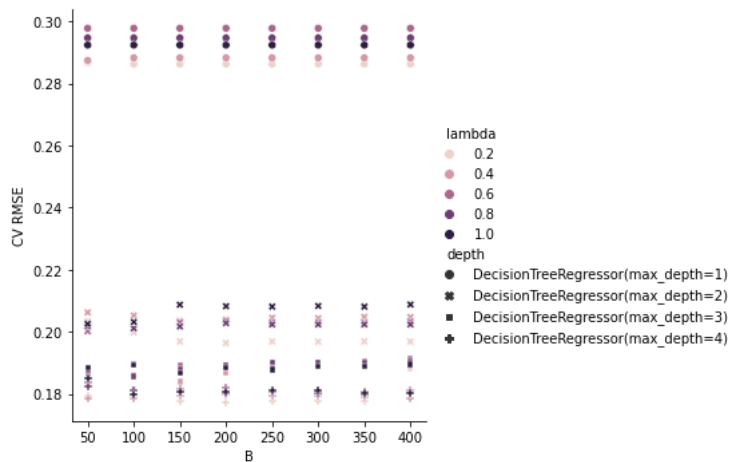
```
In [33]: 1 n_folds = 6
2 search = GridSearchCV(
3     pipe,
4     tuned_parameters,
5     cv = n_folds,
6     scoring = "neg_root_mean_squared_error",
7     # Refit the best model on the whole data set
8     refit = True
9 )
```

```
In [34]: 1 search.fit(X_train,Y_train)
2
```

```
Out[34]: GridSearchCV(cv=6,
    estimator=Pipeline(steps=[('model',
        AdaBoostRegressor(base_estimator=DecisionTreeRegressor(max_depth=3),
            random_state=425))])),
    param_grid={'model__base_estimator': [DecisionTreeRegressor(max_depth=1),
        DecisionTreeRegressor(max_depth=2),
        DecisionTreeRegressor(max_depth=3),
        DecisionTreeRegressor(max_depth=4)],
        'model__learning_rate': [0.2, 0.4, 0.6, 0.8, 1.0],
        'model__n_estimators': [50, 100, 150, 200, 250, 300,
            350, 400]},
    scoring='neg_root_mean_squared_error')
```

```
In [35]: 1 cv_res = pd.DataFrame({
2         "B": np.array(search.cv_results_["param_model__n_estimators"]),
3         "rmse": -search.cv_results_["mean_test_score"],
4         "lambda": search.cv_results_["param_model__learning_rate"],
5         "depth": search.cv_results_["param_model__base_estimator"],
6     })
7
8 plt.figure()
9 sns.relplot(
10     # kind = "line",
11     data = cv_res,
12     x = "B",
13     y = "rmse",
14     hue = "lambda",
15     style = "depth"
16 ).set(
17     xlabel = "B",
18     ylabel = "CV RMSE"
19 );
20 plt.show()
```

<Figure size 432x288 with 0 Axes>



```
In [36]: 1 -search.best_score_
2
```

Out[36]: 0.17728787774924967

```
In [37]: 1 search.best_estimator_
2
```

Out[37]: Pipeline(steps=[('model',  
AdaBoostRegressor(base\_estimator=DecisionTreeRegressor(max\_depth=4),  
learning\_rate=0.2, n\_estimators=200,  
random\_state=425))])

```
In [38]: 1 mean_squared_error(
2     Y_test,
3     search.best_estimator_.predict(X_test),
4     squared = False
5 )
```

Out[38]: 0.15874890293823835

Random forest performed best.

## 8.4.Carseat

```
In [39]: 1 df2 = pd.read_csv("Desktop/Cars.csv")
```

In [40]: 1 df2

Out[40]:

	Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	Education	Urban	US
0	9.50	138	73	11	276	120	Bad	42	17	Yes	Yes
1	11.22	111	48	16	260	83	Good	65	10	Yes	Yes
2	10.06	113	35	10	269	80	Medium	59	12	Yes	Yes
3	7.40	117	100	4	466	97	Medium	55	14	Yes	Yes
4	4.15	141	64	3	340	128	Bad	38	13	Yes	No
...	...	...	...	...	...	...	...	...	...	...	...
395	12.57	138	108	17	203	128	Good	33	14	Yes	Yes
396	6.14	139	23	3	37	120	Medium	55	11	No	Yes
397	7.41	162	26	12	368	159	Medium	40	18	Yes	Yes
398	5.94	100	79	7	284	95	Bad	50	12	Yes	Yes
399	9.71	134	37	0	27	120	Good	49	16	Yes	Yes

400 rows x 11 columns

In [41]: 1 df2.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 11 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Sales           400 non-null    float64
1   CompPrice       400 non-null    int64
2   Income          400 non-null    int64
3   Advertising     400 non-null    int64
4   Population      400 non-null    int64
5   Price          400 non-null    int64
6   ShelveLoc      400 non-null    object
7   Age            400 non-null    int64
8   Education       400 non-null    int64
9   Urban          400 non-null    object
10  US             400 non-null    object
dtypes: float64(1), int64(7), object(3)
memory usage: 34.5+ KB
```

In [42]: 1 df2['Sales1'] = df2.Sales.map(lambda x: 1 if x>8 else 0)

In [43]: 1 df2

Out[43]:

	Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	Education	Urban	US	Sales1
0	9.50	138	73	11	276	120	Bad	42	17	Yes	Yes	1
1	11.22	111	48	16	260	83	Good	65	10	Yes	Yes	1
2	10.06	113	35	10	269	80	Medium	59	12	Yes	Yes	1
3	7.40	117	100	4	466	97	Medium	55	14	Yes	Yes	0
4	4.15	141	64	3	340	128	Bad	38	13	Yes	No	0
...	...	...	...	...	...	...	...	...	...	...	...	...
395	12.57	138	108	17	203	128	Good	33	14	Yes	Yes	1
396	6.14	139	23	3	37	120	Medium	55	11	No	Yes	0
397	7.41	162	26	12	368	159	Medium	40	18	Yes	Yes	0
398	5.94	100	79	7	284	95	Bad	50	12	Yes	Yes	0
399	9.71	134	37	0	27	120	Good	49	16	Yes	Yes	1

400 rows x 12 columns

```
In [44]: 1 X1 = df2.drop(['Sales', 'Sales1'], axis=1)
2 Y1 = df2['Sales1']
3 X1_train, X1_test, Y1_train, Y1_test = train_test_split(X1, Y1, test_size=0.25, random_state=101)
```

```
In [45]: 1 num_features = ['CompPrice', 'Income', 'Advertising', 'Population', 'Price', 'Age', 'Education']
2 cat_features = ['Urban', 'US', 'ShelveLoc']
3 features = np.concatenate([num_features, cat_features])
```

```
In [46]: 1 categorical_tf = Pipeline(steps = [
2     ("cat_impute", SimpleImputer(strategy = 'most_frequent')),
3     ("encoder", OneHotEncoder())
4 ])
5
6 # Transformer for continuous variables
7 numeric_tf = Pipeline(steps = [
8     ("num_impute", SimpleImputer(strategy = 'mean')),
9 ])
10
11 # Column transformer
12 col_tf = ColumnTransformer(transformers = [
13     ('num', numeric_tf, num_features),
14     ('cat', categorical_tf, cat_features)
15 ])
```

```
In [47]: 1 classtree_mod = DecisionTreeClassifier(
2     criterion = 'gini',
3     random_state = 425
4 )
```

```
In [48]: 1 pipe = Pipeline(steps = [
2     ("col_tf", col_tf),
3     ("model", classtree_mod)
4 ])
5 pipe
```

```
Out[48]: Pipeline(steps=[('col_tf',
                           ColumnTransformer(transformers=[('num',
                                                              Pipeline(steps=[('num_impute',
                                                                                      SimpleImputer()),
                                                              ['CompPrice', 'Income',
                                                              'Advertising', 'Population',
                                                              'Price', 'Age',
                                                              'Education']),
                           ('cat',
                              Pipeline(steps=[('cat_impute',
                                                  SimpleImputer(strategy='most_frequent')),
                                                  ('encoder',
                                                    OneHotEncoder()))],
                              ['Urban', 'US',
                              'ShelveLoc'])])),
                           ('model', DecisionTreeClassifier(random_state=425))])
```

```
In [49]: 1 ccp_alpha_grid = np.linspace(start = 0.0, stop = 0.05, num = 100)
2 tuned_parameters = {
3     "model__ccp_alpha": ccp_alpha_grid
4 }
5 tuned_parameters
```

```
Out[49]: {'model__ccp_alpha': array([0.          , 0.00050505, 0.0010101 , 0.00151515, 0.0020202 ,
    0.00252525, 0.0030303 , 0.00353535, 0.0040404 , 0.00454545,
    0.00505051, 0.00555556, 0.00606061, 0.00656566, 0.00707071,
    0.00757576, 0.00808081, 0.00858586, 0.00909091, 0.00959596,
    0.01010101, 0.01060606, 0.01111111, 0.01161616, 0.01212121,
    0.01262626, 0.01313131, 0.01363636, 0.01414141, 0.01464646,
    0.01515152, 0.01565657, 0.01616162, 0.01666667, 0.01717172,
    0.01767677, 0.01818182, 0.01868687, 0.01919192, 0.01969697,
    0.02020202, 0.02070707, 0.02121212, 0.02171717, 0.02222222,
    0.02272727, 0.02323232, 0.02373737, 0.02424242, 0.02474747,
    0.02525253, 0.02575758, 0.02626263, 0.02676768, 0.02727273,
    0.02777778, 0.02828283, 0.02878788, 0.02929293, 0.02979798,
    0.03030303, 0.03080808, 0.03131313, 0.03181818, 0.03232323,
    0.03282828, 0.03333333, 0.03383838, 0.03434343, 0.03484848,
    0.03535354, 0.03585859, 0.03636364, 0.03686869, 0.03737374,
    0.03787879, 0.03838384, 0.03888889, 0.03939394, 0.03989899,
    0.04040404, 0.04090909, 0.04141414, 0.04191919, 0.04242424,
    0.04292929, 0.04343434, 0.04393939, 0.04444444, 0.04494949,
    0.04545455, 0.04595959, 0.04646465, 0.04696969, 0.04747475,
    0.04797979, 0.04848485, 0.04898989, 0.04949495, 0.05       ]))
```

```
In [50]: 1 n_folds = 5
2 search = GridSearchCV(
3     pipe,
4     tuned_parameters,
5     cv = n_folds,
6     scoring = "roc_auc",
7     # Refit the best model on the whole data set
8     refit = True
9 )
```

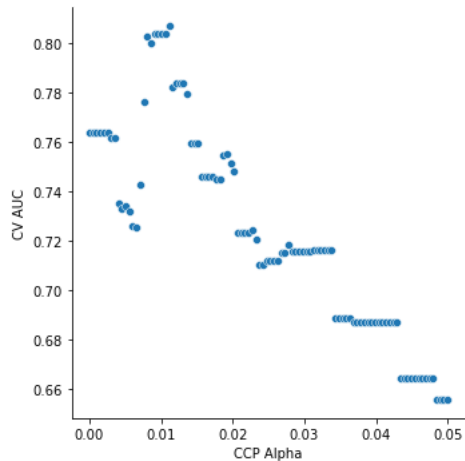
```
In [51]: 1 search.fit(X1_train,Y1_train)
2
```

```
Out[51]: GridSearchCV(cv=5,
                    estimator=Pipeline(steps=[('col_tf',
                                              ColumnTransformer(transformers=[('num',
                                          Pipeline(steps=[('num_impute',
                                                            SimpleImputer()))],
                                          ['CompPrice',
                                           'Income',
                                           'Advertising',
                                           'Population',
                                           'Price',
                                           'Age',
                                           'Education']),
                                          ('cat',
                                           Pipeline(steps=[('cat_impute',
                                                            SimpleImputer(strategy='mos
t_frequent'))],
                                          ('encoder',
                                           OneHotEncoder()))],
                                          ['Urban',
                                           'US',
                                           'ShelveLo...

0.03282828, 0.03333333, 0.03383838, 0.03434343, 0.03484848,
0.03535354, 0.03585859, 0.03636364, 0.03686869, 0.03737374,
0.03787879, 0.03838384, 0.03888889, 0.03939394, 0.03989899,
0.04040404, 0.04090909, 0.04141414, 0.04191919, 0.04242424,
0.04292929, 0.04343434, 0.04393939, 0.04444444, 0.04494949,
0.04545455, 0.04595959, 0.04646465, 0.04696969, 0.04747475,
0.04797979, 0.04848485, 0.04898989, 0.04949495, 0.05      ])),
                    scoring='roc_auc')
```

```
In [52]: 1 cv_res = pd.DataFrame({
2         "ccp_alpha": np.array(search.cv_results_["param_model__ccp_alpha"]),
3         "auc": search.cv_results_["mean_test_score"]
4     })
5
6 plt.figure()
7 sns.relplot(
8     # kind = "line",
9     data = cv_res,
10    x = "ccp_alpha",
11    y = "auc"
12 ).set(
13     xlabel = "CCP Alpha",
14     ylabel = "CV AUC"
15 );
16 plt.show()
```

<Figure size 432x288 with 0 Axes>



```
In [53]: 1 search.best_score_
2
```

Out[53]: 0.806889592760181

```
In [57]: 1 accuracy_score(
2         Y1_test,
3         search.best_estimator_.predict(X1_test)
4     )
```

Out[57]: 0.77

```
In [58]: 1 search.best_estimator_
2
```

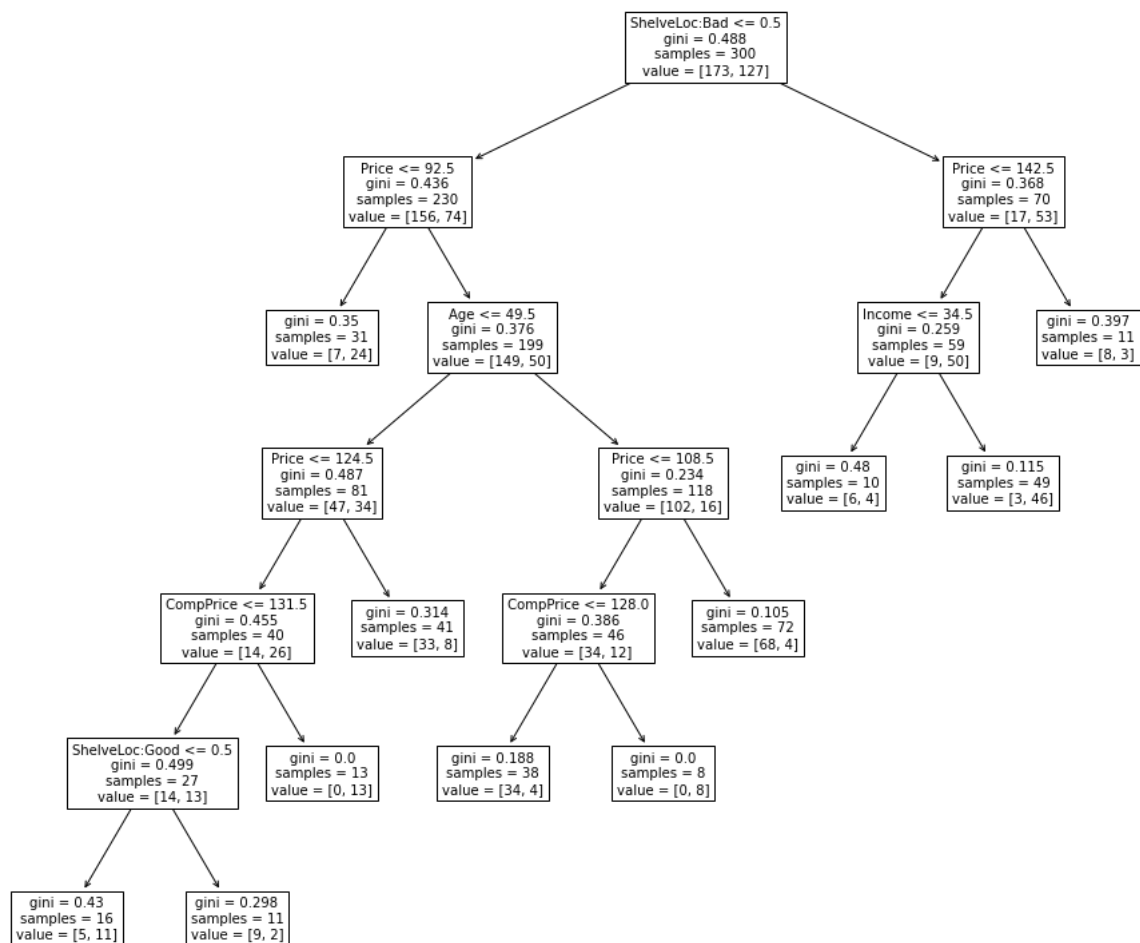
```
Out[58]: Pipeline(steps=[('col_tf',
                           ColumnTransformer(transformers=[('num',
                                                             Pipeline(steps=[('num_impute',
                                                                 SimpleImputer()))],
                                                             ['CompPrice', 'Income',
                                                                 'Advertising', 'Population',
                                                                 'Price', 'Age',
                                                                 'Education']),
                                                             ('cat',
                                                              Pipeline(steps=[('cat_impute',
                                                                 SimpleImputer(strategy='most_frequent')),
                                                                 ('encoder',
                                                                  OneHotEncoder()))],
                                                             ['Urban', 'US',
                                                                 'ShelveLoc']))],
                           ('model',
                            DecisionTreeClassifier(ccp_alpha=0.011111111111111112,
                                                    random_state=425))])
```



```

In [61]: 1 features = np.concatenate([
2         features[:-3],
3         ['ShelveLoc:Bad', 'ShelveLoc:Good', 'ShelveLoc:Medium'],
4         ['Urban:Yes', 'Urban:No'],
5         ['US:Yes', 'US:No']
6     ])
7
8 plt.figure(figsize=(16,14))
9 plot_tree(
10     search.best_estimator_['model'],
11     feature_names = features
12 );
13 plt.show()

```



```

In [68]: 1 df3 = df2.copy()

```

In [69]: 1 df3

Out[69]:

	Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	Education	Urban	US	Sales1
0	9.50	138	73	11	276	120	Bad	42	17	Yes	Yes	1
1	11.22	111	48	16	260	83	Good	65	10	Yes	Yes	1
2	10.06	113	35	10	269	80	Medium	59	12	Yes	Yes	1
3	7.40	117	100	4	466	97	Medium	55	14	Yes	Yes	0
4	4.15	141	64	3	340	128	Bad	38	13	Yes	No	0
...	...	...	...	...	...	...	...	...	...	...	...	...
395	12.57	138	108	17	203	128	Good	33	14	Yes	Yes	1
396	6.14	139	23	3	37	120	Medium	55	11	No	Yes	0
397	7.41	162	26	12	368	159	Medium	40	18	Yes	Yes	0
398	5.94	100	79	7	284	95	Bad	50	12	Yes	Yes	0
399	9.71	134	37	0	27	120	Good	49	16	Yes	Yes	1

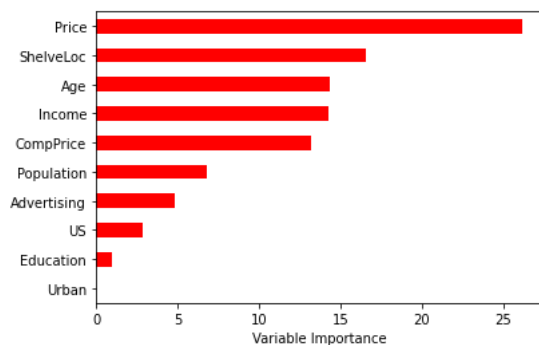
400 rows × 12 columns

```
In [70]: 1 df3['ShelveLoc'] = df3['ShelveLoc'].map({'Bad':0, 'Medium':1, 'Good':2})
2 df3['Urban'] = df3['Urban'].map({'No':0, 'Yes':1})
3 df3['US'] = df3['US'].map({'No':0, 'Yes':1})
```

```
In [93]: 1 X2 = df3.drop(['Sales', 'Sales1'], axis=1)
2 Y2 = df3['Sales1']
3 X2_train, X2_test, Y2_train, Y2_test = train_test_split(X2, Y2, test_size=0.25, random_state=101)
```

```
In [138]: 1 mod1 = classtree_mod.fit(X2_train, Y2_train)
```

```
In [139]: 1 Importance = pd.DataFrame({'Importance':mod1.feature_importances_*100}, index=X2.columns)
2 Importance.sort_values(by='Importance', axis=0, ascending=True).plot(kind='barh', color='r', )
3 plt.xlabel('Variable Importance')
4 plt.gca().legend_ = None
```



```
In [74]: 1 roc_auc_score(
2     Y1_test,
3     search.best_estimator_.predict_proba(X1_test)[: , 1]
4 )
```

Out[74]: 0.7867867867867868

```
In [77]: 1 accuracy_score(
2     Y1_test,
3     search.best_estimator_.predict(X1_test)
4 )
```

Out[77]: 0.77

## Random Forest Classifier

```
In [80]: 1 rf_mod = RandomForestClassifier(
2         # Number of trees
3         n_estimators = 100,
4         criterion = 'gini',
5         # Number of features to use in each split
6         max_features = 'sqrt',
7         oob_score = True,
8         random_state = 425
9     )
```

```
In [81]: 1 pipe = Pipeline(steps = [
2         ("col_tf", col_tf),
3         ("model", rf_mod)
4     ])
5 pipe
```

```
Out[81]: Pipeline(steps=[('col_tf',
                           ColumnTransformer(transformers=[('num',
                                                             Pipeline(steps=[('num_impute',
                                                                 SimpleImputer()))],
                                                             ['CompPrice', 'Income',
                                                                 'Advertising', 'Population',
                                                                 'Price', 'Age',
                                                                 'Education'])),
                           ('cat',
                             Pipeline(steps=[('cat_impute',
                                                                 SimpleImputer(strategy='most_frequent')),
                                                                 ('encoder',
                                                                  OneHotEncoder()))],
                             ['Urban', 'US',
                              'ShelveLoc'])])),
              ('model',
               RandomForestClassifier(max_features='sqrt', oob_score=True,
                                     random_state=425))])
```

```
In [82]: 1 B_grid = [50, 100, 150, 200, 250, 300]
2 m_grid = ['sqrt', 'log2', 1.0] # max_features = 1.0 uses all features
3 tuned_parameters = {
4     "model__n_estimators": B_grid,
5     "model__max_features": m_grid
6 }
7 tuned_parameters
```

```
Out[82]: {'model__n_estimators': [50, 100, 150, 200, 250, 300],
          'model__max_features': ['sqrt', 'log2', 1.0]}
```

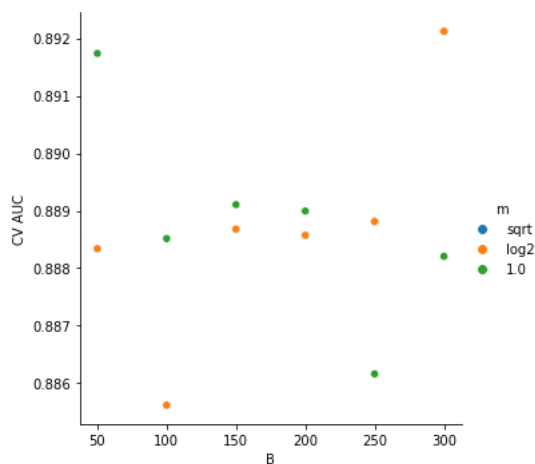
```
In [83]: 1 n_folds = 5
2 search = GridSearchCV(
3     pipe,
4     tuned_parameters,
5     cv = n_folds,
6     scoring = "roc_auc",
7     # Refit the best model on the whole data set
8     refit = True
9 )
```

```
In [84]: 1 search.fit(X1_train, Y1_train)
2
```

```
Out[84]: GridSearchCV(cv=5,
                      estimator=Pipeline(steps=[('col_tf',
                                                  ColumnTransformer(transformers=[('num',
                                                                                      Pipeline(steps=[('num_impute',
                                                                                          SimpleImputer()))],
                                                                                      ['CompPrice',
                                                                                       'Income',
                                                                                       'Advertising',
                                                                                       'Population',
                                                                                       'Price',
                                                                                       'Age',
                                                                                       'Education'])),
                                                  ('cat',
                                                                                      Pipeline(steps=[('cat_impute',
                                                                                          SimpleImputer(strategy='mos
                                                                                          t_frequent'))],
                                                                                      ('encoder',
                                                                                          OneHotEncoder()))],
                                                                                      ['Urban',
                                                                                       'US',
                                                                                       'ShelveLoc'])))],
                      ('model',
                      RandomForestClassifier(max_features='sqrt',
                                             oob_score=True,
                                             random_state=425))),
          param_grid={'model__max_features': ['sqrt', 'log2', 1.0],
                      'model__n_estimators': [50, 100, 150, 200, 250, 300]},
          scoring='roc_auc')
```

```
In [85]: 1 cv_res = pd.DataFrame({
2         "B": np.array(search.cv_results_["param_model__n_estimators"]),
3         "auc": search.cv_results_["mean_test_score"],
4         "m": search.cv_results_["param_model__max_features"]
5     })
6
7 plt.figure()
8 sns.relplot(
9     # kind = "line",
10    data = cv_res,
11    x = "B",
12    y = "auc",
13    hue = "m"
14 ).set(
15     xlabel = "B",
16     ylabel = "CV AUC"
17 );
18 plt.show()
```

<Figure size 432x288 with 0 Axes>



```
In [86]: 1 search.best_score_
2
```

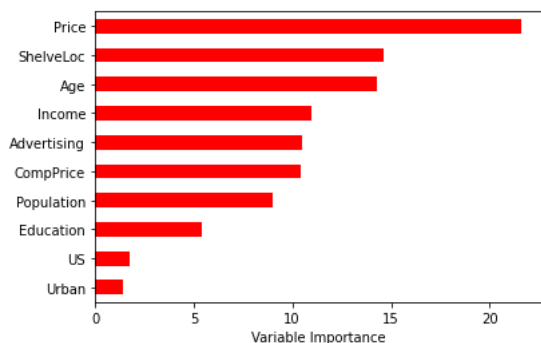
```
Out[86]: 0.8921239819004525
```

```
In [88]: 1 accuracy_score(
2         Y1_train,
3         search.best_estimator_.predict(X1_train)
4     )
```

Out[88]: 1.0

```
In [136]: 1 mod2 = rf_mod.fit(X2_train , Y2_train)
```

```
In [137]: 1 Importance = pd.DataFrame({'Importance':mod2.feature_importances_*100}, index=X2.columns)
2 Importance.sort_values(by='Importance', axis=0, ascending=True).plot(kind='barh', color='r', )
3 plt.xlabel('Variable Importance')
4 plt.gca().legend_ = None
```



```
In [100]: 1 roc_auc_score(
2         Y1_test,
3         search.best_estimator_.predict_proba(X1_test)[: , 1]
4     )
```

Out[100]: 0.8983268983268984

```
In [101]: 1 accuracy_score(
2         Y1_test,
3         search.best_estimator_.predict(X1_test)
4     )
5
6
```

Out[101]: 0.83

## Boosting

```
In [106]: 1 bst_mod = AdaBoostClassifier(
2         base_estimator = DecisionTreeClassifier(max_depth = 3),
3         # Number of trees (to be tuned)
4         n_estimators = 50,
5         # Learning rate (to be tuned)
6         learning_rate = 1.0,
7         random_state = 425
8     )
```

```
In [107]: 1 pipe = Pipeline(steps = [
2         ("col_tf", col_tf),
3         ("model", bst_mod)
4         ])
5 pipe
```

```
Out[107]: Pipeline(steps=[('col_tf',
                           ColumnTransformer(transformers=[('num',
                                                             Pipeline(steps=[('num_impute',
                                                                 SimpleImputer()))],
                                                             ['CompPrice', 'Income',
                                                                 'Advertising', 'Population',
                                                                 'Price', 'Age',
                                                                 'Education']),
                                                             ('cat',
                                                              Pipeline(steps=[('cat_impute',
                                                                 SimpleImputer(strategy='most_frequent')),
                                                                 ('encoder',
                                                                  OneHotEncoder()))],
                                                             ['Urban', 'US',
                                                                 'ShelveLoc'])])),
                           ('model',
                            AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_depth=3),
                                                random_state=425))])
```

```
In [111]: 1 d_grid = [
2     DecisionTreeClassifier(max_depth = 1),
3     DecisionTreeClassifier(max_depth = 2),
4     DecisionTreeClassifier(max_depth = 3),
5     DecisionTreeClassifier(max_depth = 4)
6 ]
7 B_grid = np.linspace(10, 100, 10).astype(int)
8 lambda_grid = [0.2, 0.4, 0.6, 0.8, 1.0]
9 tuned_parameters = {
10     "model__base_estimator": d_grid,
11     "model__n_estimators": B_grid,
12     "model__learning_rate": lambda_grid
13 }
14 tuned_parameters
```

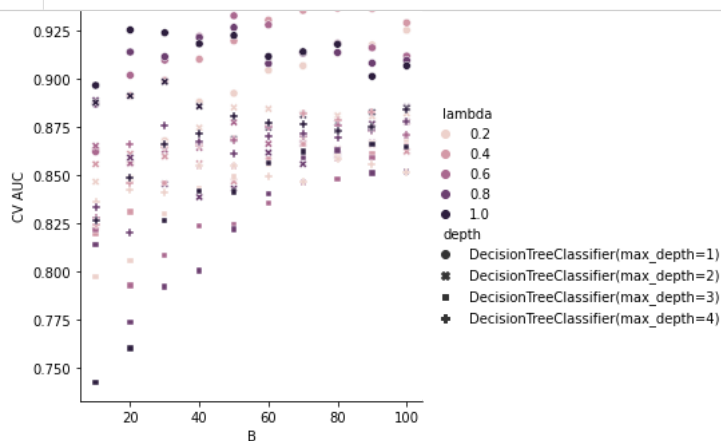
```
Out[111]: {'model__base_estimator': [DecisionTreeClassifier(max_depth=1),
DecisionTreeClassifier(max_depth=2),
DecisionTreeClassifier(max_depth=3),
DecisionTreeClassifier(max_depth=4)],
'model__n_estimators': array([ 10,  20,  30,  40,  50,  60,  70,  80,  90, 100]),
'model__learning_rate': [0.2, 0.4, 0.6, 0.8, 1.0]}
```

```
In [112]: 1 n_folds = 5
2 search = GridSearchCV(
3     pipe,
4     tuned_parameters,
5     cv = n_folds,
6     scoring = "roc_auc",
7     # Refit the best model on the whole data set
8     refit = True
9 )
```

```
In [113]: 1 search.fit(X1_train,Y1_train)
```

```
Out[113]: GridSearchCV(cv=5,
                      estimator=Pipeline(steps=[('col_tf',
                                                ColumnTransformer(transformers=[('num',
                                                                                      Pipeline(steps=[('num_impute',
                                                                                          SimpleImputer()))],
                                                                                      ['CompPrice',
                                                                                      'Income',
                                                                                      'Advertising',
                                                                                      'Population',
                                                                                      'Price',
                                                                                      'Age',
                                                                                      'Education'])),
                                                ('cat',
                                                                                      Pipeline(steps=[('cat_impute',
                                                                                          SimpleImputer(strategy='mos
t_frequent'))],
                                                                                      ('encoder',
                                                                                          OneHotEncoder()))],
                                                                                      ['Urban',
                                                                                      'US',
                                                                                      'ShelveLo...
                                                                                      AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_depth=3),
                                                                                          random_state=425))]),
                      param_grid={'model__base_estimator': [DecisionTreeClassifier(max_depth=1),
                                                            DecisionTreeClassifier(max_depth=2),
                                                            DecisionTreeClassifier(max_depth=3),
                                                            DecisionTreeClassifier(max_depth=4)],
                                'model__learning_rate': [0.2, 0.4, 0.6, 0.8, 1.0],
                                'model__n_estimators': array([ 10, 20, 30, 40, 50, 60, 70, 80, 90, 100])},
                      scoring='roc_auc')
```

```
In [115]: 1 cv_res = pd.DataFrame({
2         "B": np.array(search.cv_results_["param_model__n_estimators"]),
3         "auc": search.cv_results_["mean_test_score"],
4         "lambda": search.cv_results_["param_model__learning_rate"],
5         "depth": search.cv_results_["param_model__base_estimator"],
6     })
7
8 plt.figure()
9 sns.relplot(
10     # kind = "line",
11     data = cv_res,
12     x = "B",
13     y = "auc",
14     hue = "lambda",
15     style = "depth"
16 ).set(
17     xlabel = "B",
18     ylabel = "CV AUC"
19 );
20 plt.show()
```



```
In [125]: 1 search.best_score_
```

```
Out[125]: 0.9366921784098257
```

```
In [117]: 1 accuracy_score(
2         y1_train,
3         search.best_estimator_.predict(X1_train)
4         )
```

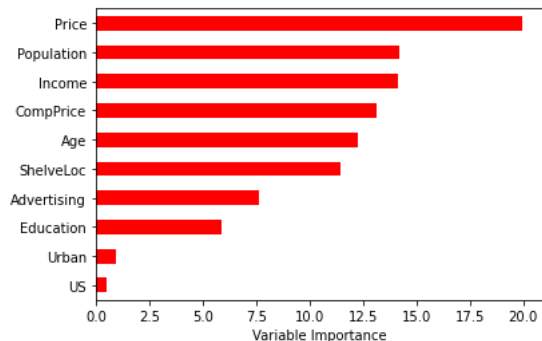
Out[117]: 0.94

```
In [126]: 1 search.best_estimator_
```

```
Out[126]: Pipeline(steps=[('col_tf',
                           ColumnTransformer(transformers=[('num',
                                                             Pipeline(steps=[('num_impute',
                                                                 SimpleImputer()))],
                                                             ['CompPrice', 'Income',
                                                                 'Advertising', 'Population',
                                                                 'Price', 'Age',
                                                                 'Education']),
                                                           ('cat',
                                                            Pipeline(steps=[('cat_impute',
                                                                 SimpleImputer(strategy='most_frequent')),
                                                                 ('encoder',
                                                                  OneHotEncoder()))],
                                                           ['Urban', 'US',
                                                            'ShelveLoc'])])),
                          ('model',
                           AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_depth=1),
                                                learning_rate=0.4, n_estimators=80,
                                                random_state=425))])
```

```
In [128]: 1 mod3 = bst_mod.fit(X2_train,Y2_train)
```

```
In [135]: 1 Importance = pd.DataFrame({'Importance':mod3.feature_importances_*100, index=X2.columns)
2         Importance.sort_values(by='Importance', axis=0, ascending=False).plot(kind='barh', color='r', )
3         plt.xlabel('Variable Importance')
4         plt.gca().legend_ = None
```



```
In [127]: 1 roc_auc_score(
2         y1_test,
3         search.best_estimator_.predict_proba(X1_test)[: , 1]
4         )
```

Out[127]: 0.9223509223509224

```
In [119]: 1 accuracy_score(
2         y1_test,
3         search.best_estimator_.predict(X1_test)
4         )
```

Out[119]: 0.86

Boosting performed best.