

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5
6 from csc665 import features, metrics
7
8 from sklearn.model_selection import train_test_split
9 from sklearn.ensemble import RandomForestRegressor
```

/Users/kalininalex/miniconda3/envs/py35_intel/lib/python3.5/site-packages/sklearn/ensemble/weight_boosting.py:29: DeprecationWarning: numpy.core.umath_tests is an internal NumPy module and should not be imported. It will be removed in a future NumPy release.
from numpy.core.umath_tests import inner1d

```
In [2]: 1 csv_df = pd.read_csv("../Melbourne_housing_FULL.csv")
```

```
In [3]: 1 X, y = features.preprocess_ver_1(csv_df)
```

```
In [4]: 1 RANDOM_STATE = 10
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_si
```

```
In [5]: 1 X_train.shape, X_test.shape
```

```
Out[5]: ((7109, 20), (1778, 20))
```

```
In [6]: 1 rf = RandomForestRegressor(
2     n_estimators=100,
3     random_state=RANDOM_STATE,
4     n_jobs=-1)
```

```
In [7]: 1 rf.fit(X_train, y_train)
```

```
Out[7]: RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
                                max_features='auto', max_leaf_nodes=None,
                                min_impurity_split=1e-07, min_samples_leaf=1,
                                min_samples_split=2, min_weight_fraction_leaf=0.0,
                                n_estimators=100, n_jobs=-1, oob_score=False, random_state=10,
                                verbose=0, warm_start=False)
```

```
In [8]: 1 y_train_pred = rf.predict(X_train)
```

```
In [9]: 1 metrics.r2_score(y_train_pred, y_train)
```

```
Out[9]: 0.9742581071101623
```

```
In [10]: 1 metrics.r2_score(rf.predict(X_test), y_test)
```

```
Out[10]: 0.8285220493562198
```

```
In [11]: 1 from sklearn.tree import DecisionTreeRegressor
```

```
In [12]: 1 dt = DecisionTreeRegressor()
```

```
In [13]: 1 dt.fit(X_train, y_train)
```

```
Out[13]: DecisionTreeRegressor(criterion='mse', max_depth=None, max_features=None,
                                max_leaf_nodes=None, min_impurity_split=1e-07,
                                min_samples_leaf=1, min_samples_split=2,
                                min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                                splitter='best')
```

```
In [14]: 1 metrics.r2_score(dt.predict(X_train), y_train)
```

```
Out[14]: 1.0
```

```
In [15]: 1 metrics.r2_score(dt.predict(X_test), y_test)
```

```
Out[15]: 0.6664240158986743
```

```
In [16]: 1 N = X_train.shape[0]
          2 N
```

```
Out[16]: 7109
```

```
In [17]: 1 X_train.shape
```

```
Out[17]: (7109, 20)
```

```
In [18]: 1 indices = np.random.randint(0, 7109, 700)
```

```
In [19]: 1 X_train.iloc[indices].shape
```

```
Out[19]: (700, 20)
```

```
In [20]: 1 np.random.seed(RANDOM_STATE)
2 n = 200
3 trees = []
4 n_samples = X_train.shape[0]
5
6 for i in range(n):
7     tree = DecisionTreeRegressor()
8     indices = np.random.randint(0, n_samples, int(n_samples * 0.1))
9     _ = tree.fit(X_train.iloc[indices, :], y_train[indices])
10    trees.append(tree)
```

▼ Test Score

```
In [21]: 1 y_test.shape
```

```
Out[21]: (1778,)
```

```
In [22]: 1 y_pred_list = []
2 for t in trees:
3     y_pred_list.append(t.predict(X_test))
4 y_pred_list = np.array(y_pred_list)
```

```
In [23]: 1 y_pred_list.shape
```

```
Out[23]: (200, 1778)
```

```
In [24]: 1 y_pred_sampled = y_pred_list.mean(axis=0)
```

```
In [25]: 1 metrics.r2_score(y_pred_sampled, y_test)
```

```
Out[25]: 0.8028782437450729
```

```
In [ ]: 1
```