

In the notebook I will show how I iterate over multiple algorithms utilizing python's sklearn library. In the "Basic" section, I will iterate though each algorithm and nothing more. Afterwards I will scale the data and then apply principle component analysis.

```
In [1]: 1 import pandas as pd
2 import numpy as np
3
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.decomposition import PCA
9 from sklearn.pipeline import Pipeline
10
11 from sklearn.ensemble import RandomForestRegressor
12 from sklearn.ensemble import GradientBoostingRegressor
13 from sklearn.linear_model import LinearRegression
14 from sklearn.tree import DecisionTreeRegressor
15 from sklearn.linear_model import SGDRegressor
16
17
18 from sklearn.model_selection import train_test_split
19
20 from sklearn.model_selection import GridSearchCV
21 from sklearn.metrics import mean_squared_error
```

```
In [2]: 1 df1 = pd.read_csv('boston_housing.csv', usecols = ['CRIM', 'ZN', 'INDUS',
2                                                         'RM', 'AGE', 'DIS', 'RA
3                                                         'B', 'LSTAT', 'target'])
```

```
In [3]: 1 df1.head()
```

Out[3]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTA
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.9
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.1
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.0
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.9
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.3

```
In [4]: 1 X = df1.drop(columns = ['target'])
        2 y = df1['target']
```

```
In [5]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30,
```

Basic

```
In [6]: 1 models = {
        2     'RandomForestRegressor' : RandomForestRegressor(max_features='sq
        3     'GradientBoostingRegressor' : GradientBoostingRegressor(),
        4     'LinearRegression' : LinearRegression(),
        5     'DecisionTreeRegressor' : DecisionTreeRegressor(),
        6     'SGDRegressor' : SGDRegressor(max_iter = 2000)
        7 }
        8
        9 hyperparameters = {
       10     'RandomForestRegressor' :
       11         {'max_depth': [6, 8, 10, 12],
       12         'n_estimators': [165, 175, 17
       13     'GradientBoostingRegressor' :
       14         {"max_depth": [2, 4, 6, 8], "mi
       15         "learning_rate": [0.001, 0.01,
       16     'LinearRegression' : {'fit_intercept' : [True, False],
       17
       18     'DecisionTreeRegressor' :
       19         {'max_depth': [2, 4, 6, 8], 'm
       20
       21     'SGDRegressor' : {'alpha' : [0.15, 0.25, 0.30, 0.35],
       22         'learning_rate' : ['constant', 'optima
       23         'penalty': ['l2', 'l1', 'elasticnet'
       24
       25 }
```

```

In [7]: 1 print('Algorithm, best hyperparameters, r-squared:')
2 i = 0
3 while i < len(list(models.values())):
4
5     grid_rf_regression = GridSearchCV(
6         estimator = list(models.values())[i],
7         param_grid = list(hyperparameters.values())[i],
8         #scoring='neg_mean_squared_error',
9         scoring='r2',
10        n_jobs=4,
11        cv = 3,
12        refit=True,
13        return_train_score=True)
14
15    grid_rf_regression.fit(X_train, y_train)
16    #display(grid_rf_regression.get_params())
17
18    print(list(models.keys())[i].upper(),':',grid_rf_regression.best_param
19
20    i += 1

```

Algorithm, best hyperparameters, r-squared:

RANDOMFORESTREGRESSOR : {'max_depth': 12, 'min_samples_leaf': 1, 'n_estimators': 180} 0.8247038461819008

GRADIENTBOOSTINGREGRESSOR : {'learning_rate': 0.1, 'max_depth': 4, 'min_samples_split': 3} 0.8431710173025984

LINEARREGRESSION : {'fit_intercept': True, 'positive': False} 0.6864013112944148

DECISIONTREEREGRESSOR : {'max_depth': 8, 'min_samples_leaf': 2} 0.7014689877766709

SGDREGRESSOR : {'alpha': 0.3, 'learning_rate': 'adaptive', 'penalty': 'l2'} -4.5017376544305714e+21

```

In [8]: 1 #display(DecisionTreeRegressor().get_params().keys())
2 #display(DecisionTreeRegressor().get_params())
3 #display(grid_rf_regression.get_params())

```

Scaling

In [9]:

```

1 scaler = StandardScaler()
2 models = {
3     'RandomForestRegressor' : Pipeline([('scaler', scaler),
4                                           ('RandomForestRegressor', Ra
5     'GradientBoostingRegressor' : Pipeline([('scaler', scaler),
6                                               ('GradientBoostingRegres
7     'LinearRegression' : Pipeline([('scaler', scaler),
8                                     ('LinearRegression', LinearRegres
9     'DecisionTreeRegressor' : Pipeline([('scaler', scaler),
10                                        ('DecisionTreeRegressor', De
11    'SGDRegressor' : Pipeline([('scaler', scaler),
12                              ('SGDRegressor', SGDRegressor())])
13    }
14
15 hyperparameters = {
16     'RandomForestRegressor' :
17         {'RandomForestRegressor__max_d
18         'RandomForestRegressor__min_s
19         'RandomForestRegressor__n_esti
20     'GradientBoostingRegressor' :
21         {"GradientBoostingRegressor__ma
22         "GradientBoostingRegressor__mi
23         "GradientBoostingRegressor__le
24     'LinearRegression' :
25         {'LinearRegression__fit_intercept'
26         'LinearRegression__positive' : [
27     'DecisionTreeRegressor' :
28         {'DecisionTreeRegressor__max_d
29         'DecisionTreeRegressor__min_s
30     'SGDRegressor' :
31         {'SGDRegressor__alpha' : [0.15, 0.25,
32         'SGDRegressor__learning_rate' : ['cons
33         'SGDRegressor__penalty' : ['l2', 'l1'
34     }

```

```

In [10]: 1 print('Algorithm, best hyperparameters, r-squared:')
2 i = 0
3 while i < len(list(models.values())):
4
5     grid_rf_regression = GridSearchCV(
6         estimator = list(models.values())[i],
7         param_grid = list(hyperparameters.values())[i],
8         #scoring='neg_mean_squared_error',
9         scoring='r2',
10        n_jobs=4,
11        cv = 3,
12        refit=True,
13        return_train_score=True)
14
15    grid_rf_regression.fit(X_train, y_train)
16    #display(grid_rf_regression.get_params())
17
18    print(list(models.keys())[i].upper(),':',grid_rf_regression.best_param
19    print(' ')
20    i += 1

```

Algorithm, best hyperparameters, r-squared:

RANDOMFORESTREGRESSOR : {'RandomForestRegressor__max_depth': 8, 'RandomForestRegressor__min_samples_leaf': 1, 'RandomForestRegressor__n_estimators': 175} 0.826227448699833

GRADIENTBOOSTINGREGRESSOR : {'GradientBoostingRegressor__learning_rate': 0.1, 'GradientBoostingRegressor__max_depth': 4, 'GradientBoostingRegressor__min_samples_split': 8} 0.843835347505636

LINEARREGRESSION : {'LinearRegression__fit_intercept': True, 'LinearRegression__positive': False} 0.6864013112944152

DECISIONTREEREGRESSOR : {'DecisionTreeRegressor__max_depth': 8, 'DecisionTreeRegressor__min_samples_leaf': 1} 0.7417983886289493

SGDREGRESSOR : {'SGDRegressor__alpha': 0.35, 'SGDRegressor__learning_rate': 'adaptive', 'SGDRegressor__penalty': None} 0.6877920559477668

Scaler and Principle Component Analysis (PCA)

In [11]:

```

1 scaler = StandardScaler()
2 pca = PCA()
3
4 models = {
5     'RandomForestRegressor' : Pipeline([('scaler', scaler),
6                                           ('pca', pca),
7                                           ('RandomForestRegressor', Ra
8     'GradientBoostingRegressor' : Pipeline([('scaler', scaler),
9                                               ('pca', pca),
10                                              ('GradientBoostingRegres
11     'LinearRegression' : Pipeline([('scaler', scaler),
12                                     ('pca', pca),
13                                     ('LinearRegression', LinearRegres
14     'DecisionTreeRegressor' : Pipeline([('scaler', scaler),
15                                         ('pca', pca),
16                                         ('DecisionTreeRegressor', De
17     'SGDRegressor' : Pipeline([('scaler', scaler),
18                               ('pca', pca),
19                               ('SGDRegressor', SGDRegressor())]),
20 }
21
22 hyperparameters = {
23     'RandomForestRegressor' :
24         {
25             'pca__n_components' : np.arange(1, 10),
26             'RandomForestRegressor__max_depth' : [3, 5, 7, 9],
27             'RandomForestRegressor__min_samples_split' : [2, 5, 10],
28             'RandomForestRegressor__n_estimators' : [100, 200, 400],
29         },
30     'GradientBoostingRegressor' :
31         { 'pca__n_components' : np.arange(1, 10),
32           'GradientBoostingRegressor__max_depth' : [3, 5, 7, 9],
33           'GradientBoostingRegressor__min_samples_split' : [2, 5, 10],
34           'GradientBoostingRegressor__learning_rate' : [0.01, 0.05, 0.1],
35         },
36     'LinearRegression' :
37         {
38             'pca__n_components' : np.arange(1, 10),
39             'LinearRegression__fit_intercept' : [True, False],
40             'LinearRegression__positive' : [True, False],
41         },
42     'DecisionTreeRegressor' :
43         {
44             'pca__n_components' : np.arange(1, 10),
45             'DecisionTreeRegressor__max_depth' : [3, 5, 7, 9],
46             'DecisionTreeRegressor__min_samples_split' : [2, 5, 10],
47         },
48     'SGDRegressor' :
49         {
50             'pca__n_components' : np.arange(1, 10),
51             'SGDRegressor__alpha' : [0.15, 0.25, 0.5],
52             'SGDRegressor__learning_rate' : ['constant', 'invscaling', 'adaptive'],
53             'SGDRegressor__penalty' : ['l2', 'l1', 'elasticnet'],
54         }

```

55	}
----	---

In [12]:

```

1 print('Algorithm, best hyperparameters, r-squared:')
2 i = 0
3 while i < len(list(models.values())):
4
5     grid_rf_regression = GridSearchCV(
6         estimator = list(models.values())[i],
7         param_grid = list(hyperparameters.values())[i],
8         #scoring='neg_mean_squared_error',
9         scoring='r2',
10        n_jobs=4,
11        cv = 3,
12        refit=True,
13        return_train_score=True)
14
15    grid_rf_regression.fit(X_train, y_train)
16    #display(grid_rf_regression.get_params())
17
18    print(list(models.keys())[i].upper(), ': ', grid_rf_regression.best_param
19    print(' ')
20    i += 1

```

Algorithm, best hyperparameters, r-squared:

RANDOMFORESTREGRESSOR : {'RandomForestRegressor__max_depth': 12, 'RandomForestRegressor__min_samples_leaf': 1, 'RandomForestRegressor__n_estimators': 165, 'pca__n_components': 12} 0.7396813003224555

GRADIENTBOOSTINGREGRESSOR : {'GradientBoostingRegressor__learning_rate': 0.1, 'GradientBoostingRegressor__max_depth': 4, 'GradientBoostingRegressor__min_samples_split': 10, 'pca__n_components': 12} 0.7943943739600443

LINEARREGRESSION : {'LinearRegression__fit_intercept': True, 'LinearRegression__n_positive': False, 'pca__n_components': 13} 0.6864013112944151

DECISIONTREEREGRESSOR : {'DecisionTreeRegressor__max_depth': 8, 'DecisionTreeRegressor__min_samples_leaf': 2, 'pca__n_components': 7} 0.5889123570638041

SGDREGRESSOR : {'SGDRegressor__alpha': 0.3, 'SGDRegressor__learning_rate': 'adaptive', 'SGDRegressor__penalty': None, 'pca__n_components': 13} 0.6874868823240415

The tree based estimators - xgboost, (not in this notebook) decision tree, and random forest - often do not benefit from scaling and pca, which is why I prefer them. Nonetheless, if you wish you use an algorithm that requires scaling or pca, you can use the code in this notebook as a guide.

Below is an example of how I might instigate a

```
In [13]: 1 pipe = Pipeline([('scaler', StandardScaler()),  
2                          ('reducer', PCA(n_components=3)),  
3                          ('regressor', RandomForestRegressor())])
```

```
In [14]: 1 pipe.fit(X_train, y_train)
```

```
Out[14]: Pipeline(steps=[('scaler', StandardScaler()), ('reducer', PCA(n_components=3)),  
                          ('regressor', RandomForestRegressor())])
```

```
In [15]: 1 print(pipe.score(X_test, y_test))
```

```
0.7161943144071752
```

```
In [ ]: 1
```