

```
In [11]: import numpy as np
import pandas as pd
import matplotlib
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
import gc
import datetime
from sklearn.ensemble import AdaBoostRegressor
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import train_test_split
from sklearn.model_selection import learning_curve
from sklearn.metrics import SCORERS
warnings.simplefilter('ignore')
matplotlib.rcParams['figure.dpi'] = 100
sns.set()

In [15]: # Extract a minibatch from the original training dataset
def TrTeSplit(X_tr, y_tr, size):
    X_trspl, X_tespl, y_trspl, y_tespl = train_test_split(X_tr, y_tr, test_size=size, random_state=0)
    X_trspl = X_trspl.reset_index(drop = True)
    X_tespl = X_tespl.reset_index(drop = True)
    y_trspl = y_trspl.reset_index(drop = True)
    y_tespl = y_tespl.reset_index(drop = True)

    return X_trspl, X_tespl, y_trspl, y_tespl

# Parameter Tunning using cross validation
def model_param_select(X, y, nfolds):
    n_estimators = [45, 50, 55]
    learning_rate = [0.5, 1, 1.5]
    #the range is being reduced since after several runs
    #I have removed some redundant range value that will yield longer run time
    param_grid = {'n_estimators': n_estimators, 'learning_rate': learning_rate}
    grid_search = GridSearchCV(AdaBoostRegressor(base_estimator = DecisionTreeRegressor(max_depth=20)), param_grid, cv=nfolds)
    grid_search.fit(X, y)
    grid_search.best_params_
    return grid_search.best_params_

def RMSE(y_test, y_pred):
    rss=((y_test-y_pred)**2).sum()
    mse=np.mean((y_test-y_pred)**2)
    rmse = np.sqrt(np.mean((y_test-y_pred)**2))
    return rmse

def DataPre(new_train):
    y_tr = new_train['meter_reading']
    X_tr = new_train.drop(columns = ['meter_reading', 'timestamp','Unnamed: 0'])
    X_trspl, X_tespl, y_trspl, y_tespl = TrTeSplit(X_tr, y_tr, 0.33)
    del X_tr, y_tr
    return X_trspl, X_tespl, y_trspl, y_tespl

def LearningVisual(model, X_trspl, y_trspl, X_tespl, y_tespl, percentSize):
    # Visualize learning curves
    plt.figure()
    train_sizes, train_scores, test_scores = \
        learning_curve(model, X_trspl, y_trspl, train_sizes=np.linspace(0.1, 1, percentSize),
                        scoring="neg_root_mean_squared_error", cv=10)
    train_sizes, val_train_scores, val_test_scores = \
        learning_curve(model, X_tespl, y_tespl, train_sizes=np.linspace(0.1, 1, percentSize),
                        scoring="neg_root_mean_squared_error", cv=10)
    plt.plot(train_sizes, -test_scores.mean(1), 'o-', color="r",label="train set")
    plt.plot(train_sizes, -val_test_scores.mean(1), 'o-', color='b', label="validation set")
    plt.xlabel("Dataset size")
    plt.ylabel("Root Mean Squared Error")
    plt.title('Learning curves')
    plt.legend(loc="best")

    plt.show()
```

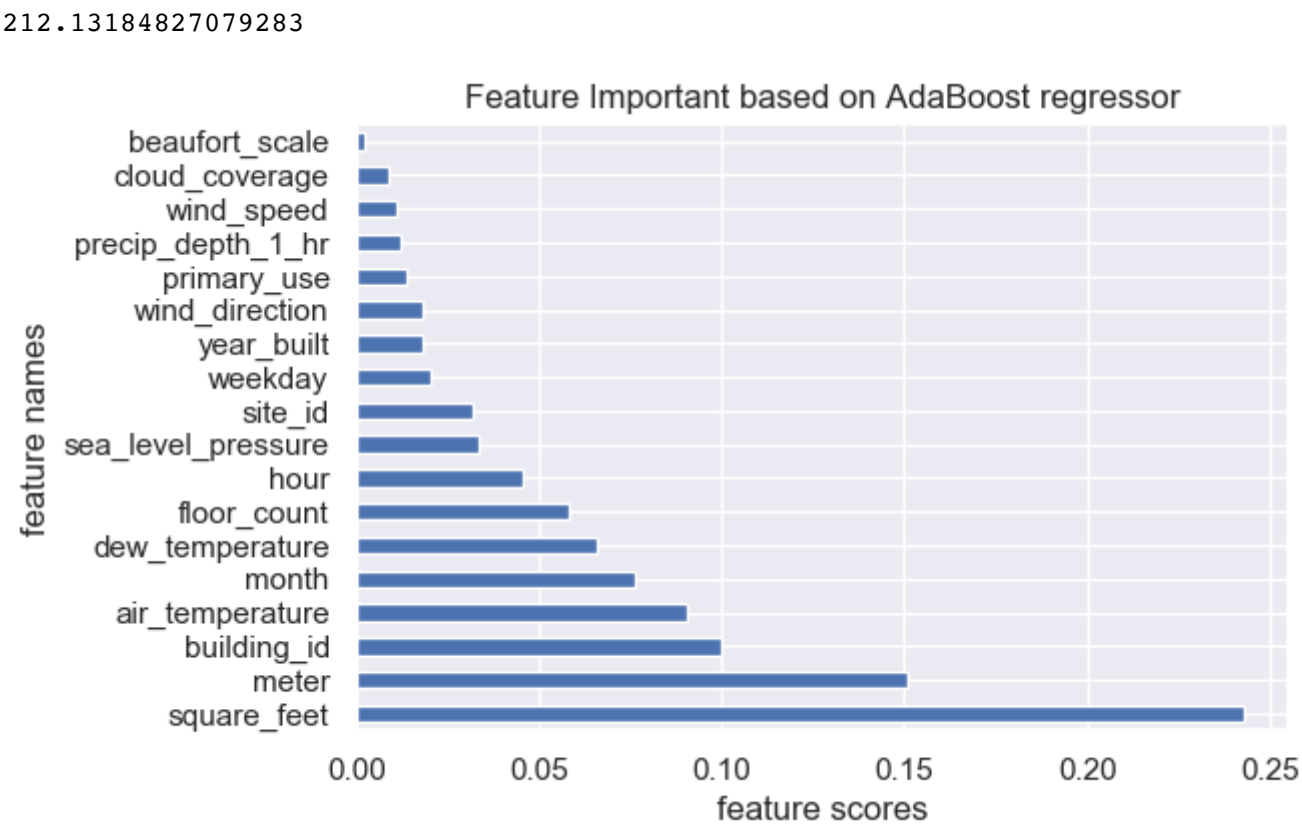
```
In [16]: # Load Data
new_train = pd.read_csv('sample_train_tiny.csv')
X_trspl, X_tespl, y_trspl, y_tespl = DataPre(new_train)
del new_train
```

```
In [17]: #Implement Adaboost model with the choosen parameter

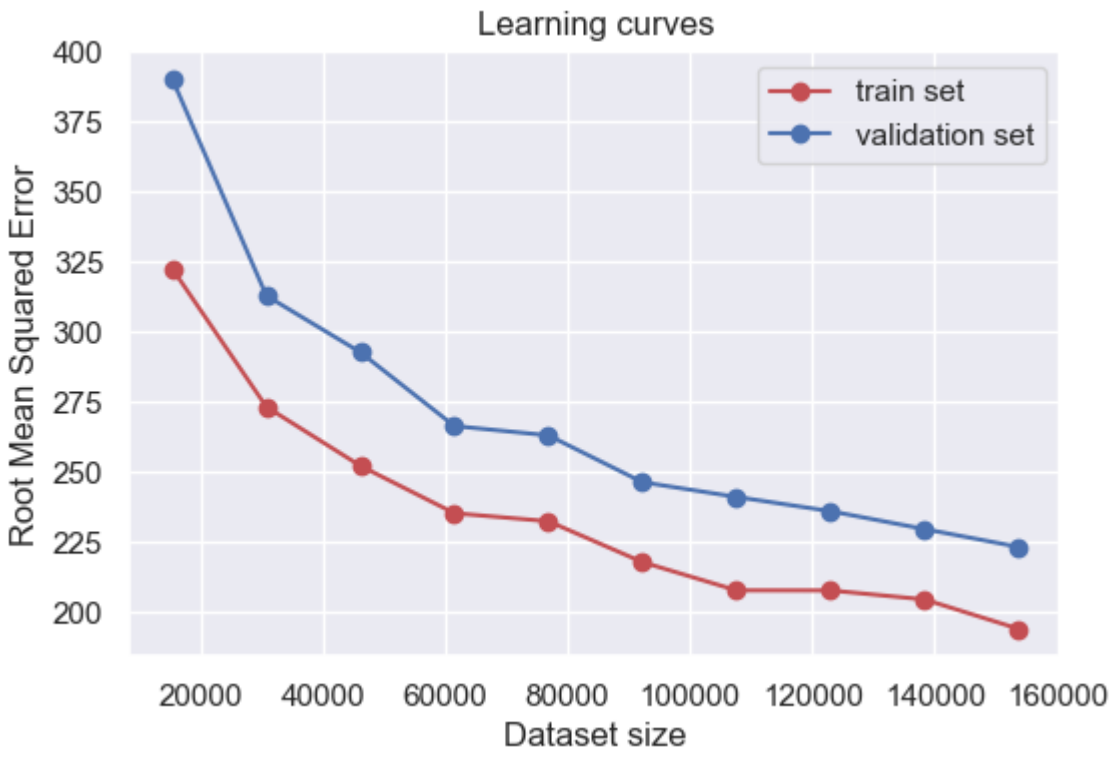
model = AdaBoostRegressor(base_estimator = DecisionTreeRegressor(max_depth=20), n_estimators = 50, learning_rate = 1.0)
model.fit(X_trspl, y_trspl)
y_pred = model.predict(X_tespl)

rmse = RMSE(y_tespl, y_pred)
print(rmse)

feat_importances = pd.Series(model.feature_importances_, index=X_trspl.columns)
feat_importances.nlargest(20).plot(kind='barh')
plt.title('Feature Important based on AdaBoost regressor')
plt.xlabel('feature scores')
plt.ylabel('feature names')
plt.show()
```



In [14]: LearningVisual(model, X\_trspl, y\_trspl, X\_tespl, y\_tespl, 10)



In [ ]: X\_train = pd.read\_csv(r'Dataset/new\_train.csv')

```
In [ ]: y_train = X_train['meter_reading']
X_train = X_train.drop(columns = ['meter_reading', 'timestamp'])
#X_test = X_test.drop(columns = ['meter_reading', 'timestamp'])
```

In [8]: X\_trspl, X\_tespl, y\_trspl, y\_tespl = TrTeSplit(X\_train, y\_train, 0.33)

```
In [ ]: #Implement Tree model with the choosen parameter

#numbers_sizes = (i*10*exp for exp in range(2, 9) for i in range(1, 10))
#for i in numbers_sizes:
model = AdaBoostRegressor(base_estimator = DecisionTreeRegressor(max_depth=20), n_estimators = 50, learning_rate = 1.0)
model.fit(X_trspl, y_trspl)
y_pred = model.predict(X_tespl)

rmse = RMSE(y_tespl, y_pred)
print(rmse)

feat_importances = pd.Series(model.feature_importances_, index=X_trspl.columns)
feat_importances.nlargest(20).plot(kind='barh')
plt.title('Feature Important based on Decision Tree regressor')
plt.xlabel('feature scores')
plt.ylabel('feature names')
plt.show()
```

In [ ]: