Applying Bootstrap Aggregated Random Forest using custom functions.

```
In [1]: import numpy as np
        import random as random
        from scipy.sparse import vstack
        # Here we are using sklearn's boston dataset
        from sklearn.datasets import load boston
        from sklearn.metrics import mean squared error
In [2]: boston = load boston()
        x=boston.data #independent variables
        y=boston.target #target variable
```

In [3]: x.shape

Out[3]: (506, 13)

In [4]: x[0]

Creating samples: Randomly create 30 samples from the whole boston data points.

consider they are [5,8,3,7] so our final sample will be [4,5,7,8,9,3,5,8,3,7].

then replicate any 203 points from the sampled points.

Out[4]: array([6.320e-03, 1.800e+01, 2.310e+00, 0.000e+00, 5.380e-01, 6.575e+00, 6.520e+01, 4.090e+00, 1.000e+00, 2.960e+02, 1.530e+01, 3.969e+02,

For example: Assume we have 10 data points [1,2,3,4,5,6,7,8,9,10], first we take 6 data points randomly, consider we have selected [4,5,7,8,9,3] now we will replicate 4 points from this sample,

Creating each sample: Considering any random 303(60% of 506) data points from whole data set and

Also, As a part of Bagging when we are taking the random samples we would also make sure that each of our sample will have different set of columns.

there in each sample. Similarly we will create 30 samples like this.

For example: Assume we have 10 columns[1,2,3,4,5,6,7,8,9,10] for the first sample we will select [3,4,5,9,1,2] and for the second sample [7,9,1,4,5,6,2] and so on. Atleast 3 feautres/columns would be

In [5]: # In this function we will write code for generating 30 samples def generating samples(input data, target data):

selected rows=np.random.choice(len(input data), 303, replace=False) # Selecting 303 unique random d

replicating rows=np.random.choice(selected rows, 203) # Replicating 203 points from 303 points whic

https://numpy.org/doc/stable/reference/random/generated/numpy.random.choice.html

num=random.randint(3,12) # Selecting random number for columns with minimun 3 columns.

ata points. #print(len(selected rows))

#print(len(replicating rows))

h were generated above.

selected columns=[]

Using np.random.choice for genrating random sample.

4.980e+00])

```
#print(num)
            selected columns=random.sample(range(1,13), num) # Selecting random columns.
            #print(selected columns)
            sample data=input data[selected rows[:, None], selected columns] # Adding the rows and columns to a
         variable.
            #print(sample data)
            target sample data=target data[selected rows] # Adding output variable rows similar to as sample da
        ta.
            #print(len(target sample data))
            replicated sample data=input data[replicating rows[:, None], selected columns] # Adding repeating ro
        ws and columns.
            #print(replicated_sample_data.shape)
            target replicate sample data=target data[replicating rows] # Adding output variable repeating rows
         similar to as sample data.
            #print(len(target replicate sample data))
            final sample=np.vstack((sample data, replicated sample data)) # Combining the data.
            final target=np.vstack((target sample data.reshape(-1,1), target replicate sample data.reshape(-1,1
        )))  # Combining output data.
            #print(final sample.shape)
            #print(final target.shape)
            return final sample, final target, selected rows, selected columns
        Creating 30 samples.
In [7]:
        # Use generating samples function to create 30 samples.
        # Storing these samples in a list.
        list input data =[]
        list output data =[]
        list selected row= []
        list selected columns=[]
        for i in range (0,30):
            a,b,c,d=generating samples(x,y)
```

from tqdm import tqdm

In [9]: **from sklearn.tree import** DecisionTreeRegressor

for i in tqdm(range(len(x))):

pred.append(p)

print('Mean Square Error: ',mse)

Mean Square Error: 0.07789946264658806

for j in range(len(list models)):

model=list models[j]

columns=list selected columns[j]

Taking median of data point from each model

Now we are calculating Out of Bag Score (OOB score).

trained using the same data point for which we are predicting the value.

#Skipping the data points from which models was trained.

Repeating all the above steps again to obtain different MSE and OOB scores.

p = model.predict(data[columns].reshape(1,-1))

p = model.predict(data[columns].reshape(1,-1))

data=x[i] pred=[]

100%|

00:00, 259.34it/s]

predictions=[]

data=x[i] pred=[]

#print (pred)

for i in tqdm(range(len(x))):

for j in range(len(list models)):

rows=list selected row[j]

model=list models[j]

pred.append(p)

mse=mean squared error(y,predictions)

In [13]: # Repeating all above steps 35 times.

for sample in tqdm(range(35)):

list input data =[] list output data =[] list selected row= [] list selected columns=[]

for i in range (0,30):

list models=[]

predictions=[]

for i in range(len(x)):

a,b,c,d=generating samples(x,y)

for i in range(len(list input data)):

list models.append(model)

model=list models[j]

pred.append(p)

model=DecisionTreeRegressor(max depth=None)

columns=list selected columns[j]

predictions.append(np.median(pred))

model.fit(list_input_data[i], list_output_data[i])

list_input_data.append(a) list output data.append(b) list selected row.append(c) list selected columns.append(d)

list mse=[] list oob=[]

predictions.append(np.median(pred))

if i not in rows:

columns=list selected columns[j]

In [10]:

list models=[]

list input data.append(a) list output data.append(b) list selected row.append(c) list selected columns.append(d)

```
for i in tqdm(range(len(list_input_data))):
             model=DecisionTreeRegressor(max depth=None)
             model.fit(list input data[i], list output data[i])
             list models.append(model)
         100%|
                                                                                               | 30/30 [00:00<
         00:00, 191.27it/s]
         Now we are predicting the value of each data point with our 30 models and storing the median of
         prediction of 30 models. Then calculating the Mean Sqaure Error from median of predictions.
In [11]: predictions=[]
         # Iterating over each data point.
```

Now we are training 30 Decision Tree's with high variance for our sample data.

Making prediction on each data points with 30 models we created above.

Training 30 different high variance models for our sample data.

predictions.append(np.median(pred)) #print(predictions) mse=mean squared error(y,predictions) # Computing Mean squre error

This is similar to above step, the only difference is that we will not make predition if our model was

In [12]: # Performing the same above steps but this time not to include the the points from which model was trai

| 506/506 [00:01<

| 506/506 [00:00<

35/35 [01:35

print('OOB Score: ',mse)

00:00, 537.75it/s] OOB Score: 14.080603198656924

data=x[i] pred=[] for j in range(len(list models)):

p = model.predict(data[columns].reshape(1,-1))

#print(predictions) mse=mean squared error(y,predictions) list mse.append(mse) predictions=[] for i in range(len(x)): data=x[i] pred=[] for j in range(len(list models)): columns=list selected columns[j] rows=list selected row[j] model=list models[j] if i not in rows: p = model.predict(data[columns].reshape(1,-1)) pred.append(p) #print(pred) predictions.append(np.median(pred))

In [14]: # Computing Confidence Intervals.

https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.norm.html import scipy.stats as st # Creating 95% confidence interval for MSE

In [15]: # Repeating the same steps as done in Task 1. But this time just to predict for a single query point. xq = [0.18, 20.0, 5.00, 0.0, 0.421, 5.60, 72.2, 7.95, 7.0, 30.0, 19.1, 372.13, 18.60]pt=np.array(xq) for i in range(len(list models)): cols=list selected columns[i] model=list models[i] pred=model.predict(pt[cols].reshape(1,-1))

Using our 30 models to predict the price of house from a new data point.

ans.append(pred) print('Price of the house of new data point:',np.median(ans)) Price of the house of new data point: 18.9

interval mse=st.norm.interval(alpha=0.95, loc=np.mean(list mse), scale=st.sem(list mse)) print('Confidence Interval for Mean Square Errors:',interval mse)

https://www.statology.org/confidence-intervals-python/

mse=mean squared_error(y,predictions) list oob.append(mse) 100%| <00:00, 2.72s/it] Now we are computing Confidence Intervals of MSE and OOB scores.

Creating 95% confidence interval for MSE interval oob=st.norm.interval(alpha=0.95, loc=np.mean(list oob), scale=st.sem(list oob)) print('Confidence Interval for OOB scores:',interval oob) Confidence Interval for Mean Square Errors: (0.08597271699380007, 0.16566371321125348) Confidence Interval for OOB scores: (13.513575196555944, 14.537933193188813)