

▼ Bootstrap assignment

There will be some functions that start with the word "grader" ex: grader_sampples(), grader_30().. etc, you should not change those function definition.

Every Grader function has to return True.

Importing packages

```
1 import numpy as np # importing numpy for numerical computation
2 from sklearn.datasets import load_boston # here we are using sklearn's boston
3 from sklearn.metrics import mean_squared_error # importing mean_squared_error
```

```
1 boston = load_boston()
2 x=boston.data #independent variables
3 y=boston.target #target variable
```

```
1 x.shape
```

```
↳ (506, 13)
```

```
1 x[:5]
```

```
↳ array([[6.3200e-03, 1.8000e+01, 2.3100e+00, 0.0000e+00, 5.3800e-01,
          6.5750e+00, 6.5200e+01, 4.0900e+00, 1.0000e+00, 2.9600e+02,
          1.5300e+01, 3.9690e+02, 4.9800e+00],
         [2.7310e-02, 0.0000e+00, 7.0700e+00, 0.0000e+00, 4.6900e-01,
          6.4210e+00, 7.8900e+01, 4.9671e+00, 2.0000e+00, 2.4200e+02,
          1.7800e+01, 3.9690e+02, 9.1400e+00],
         [2.7290e-02, 0.0000e+00, 7.0700e+00, 0.0000e+00, 4.6900e-01,
          7.1850e+00, 6.1100e+01, 4.9671e+00, 2.0000e+00, 2.4200e+02,
          1.7800e+01, 3.9283e+02, 4.0300e+00],
         [3.2370e-02, 0.0000e+00, 2.1800e+00, 0.0000e+00, 4.5800e-01,
          6.9980e+00, 4.5800e+01, 6.0622e+00, 3.0000e+00, 2.2200e+02,
          1.8700e+01, 3.9463e+02, 2.9400e+00],
         [6.9050e-02, 0.0000e+00, 2.1800e+00, 0.0000e+00, 4.5800e-01,
          7.1470e+00, 5.4200e+01, 6.0622e+00, 3.0000e+00, 2.2200e+02,
          1.8700e+01, 3.9690e+02, 5.3300e+00]])
```

▼ Task 1

Step - 1

- **Creating samples**

Randomly create 30 samples from the whole boston data points

- Creating each sample: Consider any random 303(60% of 506) data points from whole data set and then replicate any 203 points from the sampled points

For better understanding of this procedure let's check this 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 [4, 5, 7, 8, 9, 3], consider they are [5, 8, 3, 7] so our final sample will be [4, 5, 7, 8, 9, 3, 5, 8, 3, 7]

- **Create 30 samples**

- Note that as a part of the Bagging when you are taking the random samples **make sure each of the sample will have different set of columns**
Ex: 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... Make sure each sample will have at least 3 features/[columns/attributes](#)

Step - 2

Building High Variance Models on each of the sample and finding train MSE value

- Build a regression trees on each of 30 samples.
- Computed the predicted values of each data point(506 data points) in your corpus.
- Predicted house price of i^{th} data point

$$y_{pred}^i = \frac{1}{30} \sum_{k=1}^{30} (\text{predicted value of } x^i \text{ with } k^{th} \text{ model})$$

- Now calculate the $MSE = \frac{1}{506} \sum_{i=1}^{506} (y^i - y_{pred}^i)^2$

Step - 3

- **Calculating the OOB score**

- Predicted house price of i^{th} data point

$$y_{pred}^i = \frac{1}{k} \sum_{k=\text{model which was built on samples not included } x^i} (\text{predicted value of } x^i \text{ with } k)$$

- Now calculate the $OOBScore = \frac{1}{506} \sum_{i=1}^{506} (y^i - y_{pred}^i)^2$.

▼ Task 2

- **Computing CI of OOB Score and Train MSE**

- Repeat Task 1 for 35 times, and for each iteration store the Train MSE and OOB score

- After this we will have 35 Train MSE values and 35 OOB scores
- using these 35 values (assume like a sample) find the confidence intervals of MSE and OOB Score
- you need to report CI of MSE and CI of OOB Score
- Note: Refer the Central_Limit_theorem.ipynb to check how to find the confidence interval

▼ Task 3

- **Given a single query point predict the price of house.**

Consider $x_q = [0.18, 20.0, 5.00, 0.0, 0.421, 5.60, 72.2, 7.95, 7.0, 30.0, 19.1, 372.13, 18.60]$ Predict the house price for this point as mentioned in the step 2 of Task 1.

▼ Task - 1

Step - 1

- **Creating samples**

Algorithm

Pesudo Code for generating Sample

```
def generating_samples(input_data, target_data):  
    Selecting_rows <--- Getting 303 random row indices from the input_data  
    Replaing_rows <--- Extracting 206 random row indices from the "Selecting_rows"  
    Selecting_columns<--- Getting from 3 to 13 random column indices  
    sample_data<--- input_data[Selecting_rows[:,None],Selecting_columns]  
    target_of_sample_data <--- target_data[Selecting_rows]  
    #Replicating Data  
    Replicated_sample_data <--- sample_data [Replaing_rows]  
    target_of_Replicated_sample_data<--- target_data[Replaing_rows]  
    # Concatinating data  
    final_sample_data <--- perform vertical stack on sample_data, Replicated_sample_data  
    final_target_data<--- perform vertical stack on target_of_sample_data.reshape(-1,1), target_of_Replicated_sample_data.reshape(-1,1)  
    return final_sample_data, final_target_data, Selecting_rows, Selecting_columns
```

- **Write code for generating samples**

```
1  def generating_samples(input_data, target_data):  
2  
3      '''In this function, we will write code for generating 30 samples '''  
4      # you can use random.choice to generate random indices without replacement  
5      # Please have a look at this link https://docs.scipy.org/doc/numpy-1.16.0/,  
6      # Please follow above pseudo code for generating samples  
7  
8      #print(np.floor(0.6*len(input_data)))  
9      select_rows = np.random.choice(len(input_data),size = int(np.floor(0.6*len(input_data))))  
10     #print(select_rows,select_rows.shape)  
11     replace_rows = np.random.choice((select_rows),size = int(np.ceil(0.4*len(input_data))))  
12     #print(replace_rows,replace_rows.shape)  
13  
14     no_col = 8  
15     select_cols = np.random.choice(len(input_data[0,:]),size =no_col, replace = True)  
16     #print(select_cols)  
17     sam = np.meshgrid(select_cols,select_rows)  
18     #print((sam[0].shape))  
19     #print(sam[1].shape)  
20  
21     sam_data = input_data[sam[1],sam[0]]  
22     #print(sam_data.shape)  
23  
24     target_sam_data = target_data[select_rows]  
25     #print(target_sam_data.shape)  
26  
27     sam = np.meshgrid(select_cols,replace_rows)  
28     #print((sam[0].shape))  
29     #print(sam[1].shape)
```

```

30     replicate_data = input_data[sam[1],sam[0]]
31     #print(replicate_data.shape)
32
33     target_rep_data = target_data[replace_rows]
34     #print(target_rep_data.shape)
35
36     final_sam_data = np.concatenate((sam_data,replicate_data),axis = 0)
37     final_target_data = np.concatenate((target_sam_data,target_rep_data))
38
39     #print(final_sam_data.shape,final_target_data.shape,select_rows.shape,select_columns.shape)
40
41     return list(final_sam_data) , list(final_target_data),list(select_rows),list(select_columns)
42     #note please return as lists

```

Grader function - 1

```

1  def grader_samples(a,b,c,d):
2      length = (len(a)==506 and len(b)==506)
3      sampled = (len(a)-len(set([str(i) for i in a]))==203)
4      rows_length = (len(c)==303)
5      column_length= (len(d)>=3)
6      assert(length and sampled and rows_length and column_length)
7      return True
8  a,b,c,d = generating_samples(x, y)
9  grader_samples(a,b,c,d)

```

☞ True

- **Create 30 samples**

Run this code 30 times, so that you will 30 samples, and store them in a lists as shown below:

```

list_input_data=[]
list_output_data=[]
list_selected_row=[]
list_selected_columns=[]

for i in range(0,30):
    a,b,c,d=generating_sample(input_data,target_data)
    list_input_data.append(a)
    list_output_data.append(b)
    list_selected_row.append(c)
    list_selected_columns.append(d)

```

```

1  # Use generating_samples function to create 30 samples
2  # store these created samples in a list

```

```

3  list_input_data =[]
4  list_output_data =[]
5  list_selected_rows= []
6  list_selected_columns=[]
7  no_base_models = 30
8  count = 0
9  i=0
10 while (i != (no_base_models+count)):
11     a,b,c,d = generating_samples(x, y)
12     d = sorted(d)
13     i = i+1
14     if ((d in list_selected_columns)):
15         #print(count)
16         count = count+1
17         continue
18     list_selected_columns.append(d)
19     list_input_data.append(a)
20     list_output_data.append(b)
21     list_selected_rows.append(c)
22
23

```

```

1  print(len(list_selected_columns),len(list_selected_rows))
2  print(len(list_input_data),len(list_output_data))

```

```

↳ 30 30
   30 30

```

Grader function - 2

```

1  def grader_30(a):
2      assert(len(a)==30 and len(a[0])==506)
3      return True
4  grader_30(list_input_data)

```

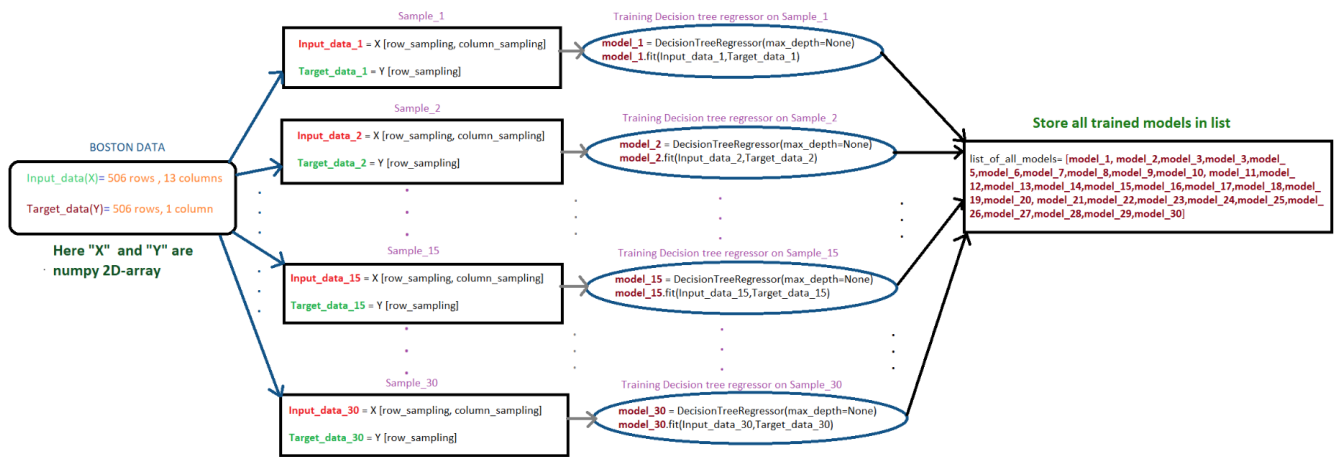
```

↳ True

```

Step - 2

Flowchart for building tree

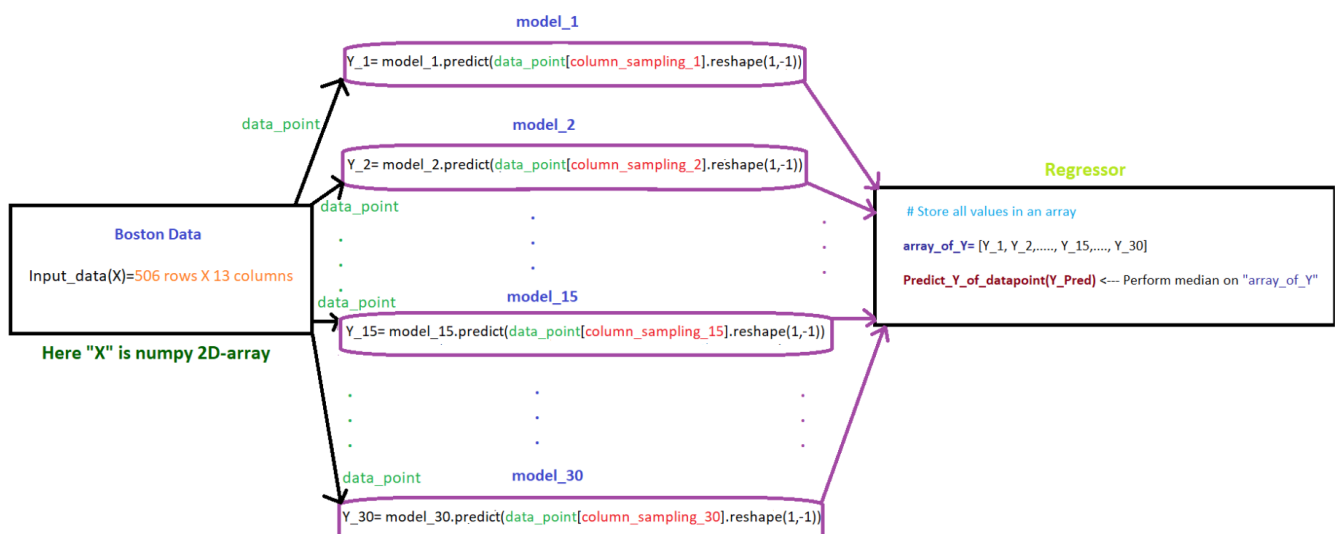


- Write code for building regression trees

```

1 from sklearn.tree import DecisionTreeRegressor
2
3 list_all_models = []
4 for i in range(no_base_models):
5     input_data = list_input_data[i]
6     target_data = list_output_data[i]
7     model = DecisionTreeRegressor()
8     model.fit(input_data, target_data)
9     list_all_models.append(model)
10
11 #print(list_all_models)
  
```

Flowchart for calculating MSE



After getting predicted_y for each data point, we can use sklearn's mean_squared_error to calculate the MSE between predicted_y and actual_y.

- **Write code for calculating MSE**

```

1 import numpy as np
2 from sklearn.metrics import mean_squared_error
3
4 y_pred = []
5 for i in range(no_base_models):
6     sam_cols = list_selected_columns[i]
7     sam_data = x[:,sam_cols]
8     #print(sam_data.shape)
9     model = list_all_models[i]
10    y = model.predict(sam_data)
11    #print(y.shape)
12    y_pred.append(y)
13 y_pred = np.array(y_pred)
14 print((y_pred.shape))
15 y_pred= np.median(y_pred, axis= 0)
16 print((y_pred.shape))
17 mse = mean_squared_error(y,y_pred)
18 print(mse)

```

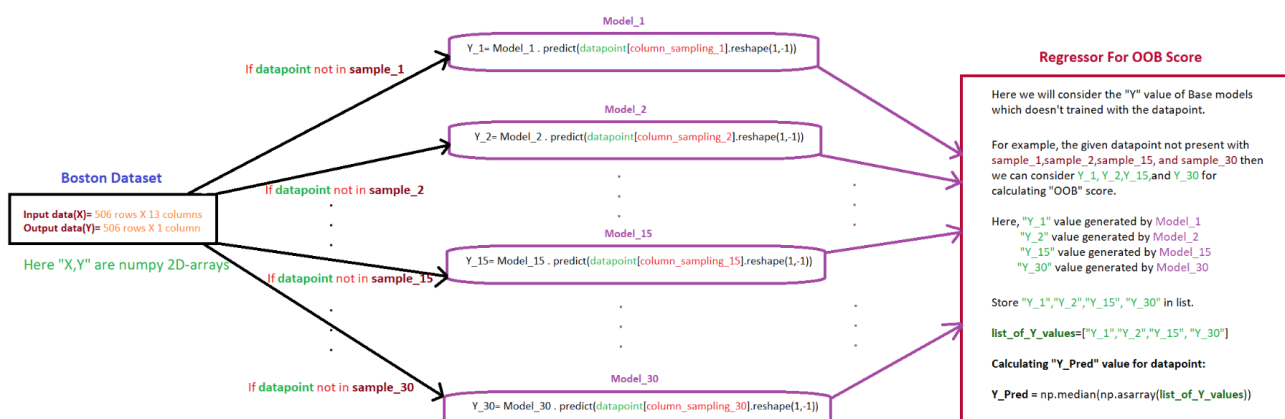
```

↳ (30, 506)
   (506,)
   22.11139328063241

```

Step - 3

Flowchart for calculating OOB score



Now calculate the $OOB Score = \frac{1}{506} \sum_{i=1}^{506} (y^i - y_{pred}^i)^2$.

- Write code for calculating OOB score

```
1 y_pred = []
2 for i in range(len(x)):
3     sam=[]
4     for j in range(no_base_models):
5         if i not in list_selected_rows[j]:
6             model = list_all_models[j]
7             sam.append(model.predict(x[i,list_selected_columns[j]].reshape(1,-1)))
8     y_pred.append(np.median(sam))
9
10 #print(len(y_pred))
11 oobscore = mean_squared_error(y,y_pred)
12 print(oobscore)
13
14
```

➤ 48.836062252964425

```
1 print(no_base_models,mse,oobscore)
```

➤ 30 22.11139328063241 48.836062252964425

▼ Task 2

```
1 Train_mse = []
2 OOb_score = []
3 no_iter = 35
4 for count2 in range(no_iter):
5     list_input_data=[]
6     list_output_data=[]
7     list_selected_rows= []
8     list_selected_columns=[]
9     no_base_models = 30
10    count = 0
11    i=0
12    while (i != (no_base_models+count)):
13        a,b,c,d = generating_samples(x, y)
14        d = sorted(d)
15        i = i+1
16        if ((d in list_selected_columns)):
17            #print(count)
18            count = count+1
19            continue
20        list_selected_columns.append(d)
21        list_input_data.append(a)
22        list_output_data.append(b)
```

```

22     list_output_data.append(y)
23     list_selected_rows.append(c)
24
25     #print(len(list_input_data), len(list_selected_columns ))
26     list_all_models = []
27     i=0
28     for i in range(no_base_models):
29         input_data = list_input_data[i]
30         target_data = list_output_data[i]
31         model = DecisionTreeRegressor()
32         model.fit(input_data, target_data)
33         list_all_models.append(model)
34
35     y_pred = []
36     for i in range(no_base_models):
37         sam_cols = list_selected_columns[i]
38         sam_data = x[:,sam_cols]
39         #print(sam_data.shape)
40         model = list_all_models[i]
41         y_pred.append(model.predict(sam_data))
42
43     y_pred = np.array(y_pred)
44     #print((y_pred.shape))
45     y_pred= np.median(y_pred, axis= 0)
46     #print((y_pred.shape))
47     mse = mean_squared_error(y,y_pred)
48     #print(mse)
49     Train_mse.append(mse)
50
51     y_pred = []
52     for i in range(len(x)):
53         sam=[]
54         for j in range(no_base_models):
55             if i not in list_selected_rows[j]:
56                 model = list_all_models[j]
57                 sam.append(model.predict(x[i,list_selected_columns[j]].reshape(1,-1)))
58         y_pred.append(np.median(sam))
59
60     #print(len(y_pred))
61     oobscore = mean_squared_error(y,y_pred)
62     print(mse,oobscore)
63     OOb_score.append(oobscore)
64
65     print(len(Train_mse))
66     print(len(OOb_score))

```



```

34.98829051383399 45.772623517786506
27.723399209486146 29.881210474308304
27.935177865612715 34.208211462450606
25.881111660079156 29.504160079051466
19.98754940711454 21.600968379446552
24.246689723319996 41.32165019762818
15.200286561264864 20.85303359683799
7.379609683794486 14.544293478260863
17.458710474308205 20.80408102766797
26.312895256917127 37.78571146245058
14.318824110671793 20.11350296442683
26.29342885375504 32.893542490118634
20.154298418972207 26.444392292490043
21.782677865612577 33.84576086956515
27.77843379446646 27.244426877470392
23.41337944664018 23.297524703557258
27.819061264822103 28.64650691699599
23.01041996047444 22.844693675889417
23.303913043478186 22.502801383399188
33.51211956521746 47.2355583003952
27.935177865612626 29.55178853754942
13.776704545454422 20.667826086956385
30.73194664031629 44.32474802371543
16.82976778656112 21.768246047430765
23.508285573122603 24.676210474308377
27.93517786561269 28.218463438735213
18.09698616600785 25.15606719367582

```

```

1 print((Train_mse))
2 print((00b_score))

```

```

[34.98829051383399, 27.723399209486146, 27.935177865612715, 25.88111166007915
[45.772623517786506, 29.881210474308304, 34.208211462450606, 29.5041600790514

```

```

-----
>>>

```

```

1 x_mean = np.mean(Train_mse)
2 x_std = np.std(Train_mse)
3 #print(x_mean, x_std)
4 # for 95% confidence interval
5 upper_lim = x_mean-2*(x_std/np.sqrt(len(Train_mse)))
6 lower_lim = x_mean+2*(x_std/np.sqrt(len(Train_mse)))
7 #print(upper_lim,lower_lim)
8 print('95% confidence interval of mean square error is [{0},{1}]'.format(upper_lim,lower_lim))

```

```

95% confidence interval of mean square error is [20.92001236793118,26.1007883

```

```

1 x_mean = np.mean(00b_score)
2 x_std = np.std(00b_score)
3 #print(x_mean, x_std)
4 # for 95% confidence interval
5 upper_lim = x_mean-2*(x_std/np.sqrt(len(00b_score)))
6 lower_lim = x_mean+2*(x_std/np.sqrt(len(00b_score)))
7 #print(upper_lim,lower_lim)
8 print('95% confidence interval for Out of Bag(00B) score is [{0},{1}]'.format(upper_lim,lower_lim))

```

```

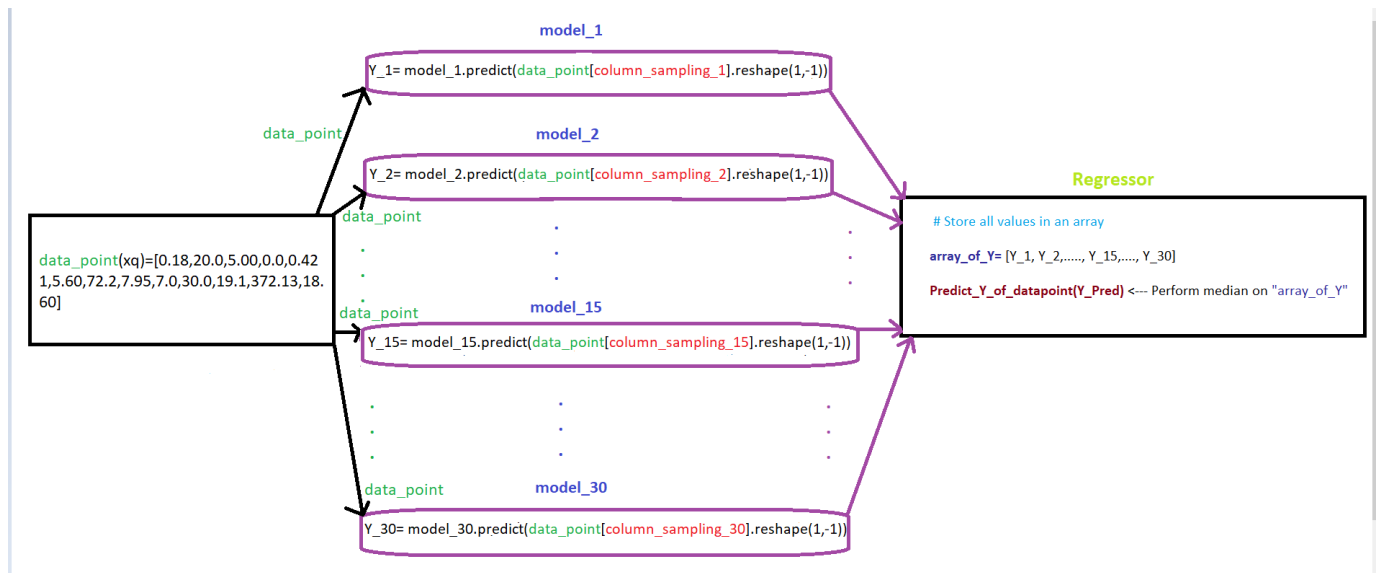
95% confidence interval for Out of Bag(00B) score is [26.580092596300947,32.5

```

▼ Task 3

Flowchart for Task 3

Hint: We created 30 models by using 30 samples in TASK-1. Here, we need send query point "xq" to 30 models and perform the regression on the output generated by 30 models.



- **Write code for TASK 3**

```

1  xq= np.array([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
2  ✓ for i in range(no_base_models):
3      data = xq[list_selected_columns[i]].reshape(1, -1)
4      model = list_all_models[i]
5      y_pred.append(model.predict(data))
6  target = np.median(y_pred)
7  print(target)
  
```

📄 21.9

Write observations for task 1, task 2, task 3 in detail

Task1

1. By increasing the number of base learners in the Random Forest, Training MSE and OOB_score can be reduced as the more models learn different aspects of the data.
2. we can observe that by decreasing the no of features sampled Training_mse and OOB_Score decreased to some extent and later increased.(Conclusion is to keep the no of features as reasonably well so that good learning happens by the different base learners.)

Task2

1. As the number of iterations increases, the sample size of train_mse and oob_scores increases and there by 95% confidence interval range decreases by a lot extent.

Task3

1. As the no of base learners increases and the no of sampled features are reasonably well, the prediction of the input test point will be accurate and the test mse will be less.