**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.</b>** Importing packages In [1]: import numpy as np # importing numpy for numerical computation from sklearn.datasets import load\_boston # here we are using sklearn's boston dataset from sklearn.metrics import mean\_squared\_error # importing mean\_squared\_error metric 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[:5] Out[4]: 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 lets check this examples, 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], consder 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, 7, 8, 7, 8, 7, 8, 7, 8, 7, 8, 7, 8, 7, 8, 7, 8, 7, 8, 7, 8, 7, 8, 7, 8, 9, 1, 2] and for the second sample [7, 9, 1, 4, 5, 6, 2] and so on... Make sure each sample will have atleast 3 feautres/columns/attributes Note - While selecting the random 60% datapoints from the whole data, make sure that the selected datapoints are all exclusive, repetition is not allowed. 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  $y^i_{pred} = \frac{1}{30} \sum_{k=1}^{30} (\text{predicted value of } x^i \text{ with } k^{th} \text{ model})$ • Now calculate the  $MSE = rac{1}{506} \sum_{i=1}^{506} (y^i - y^i_{pred})^2$ Step - 3 · Calculating the OOB score • Predicted house price of  $i^{th}$  data point  $y^i_{pred} = \frac{1}{k} \sum_{\mathbf{k} = \text{model which was buit on samples not included } x^i$  (predicted value of  $x^i$  with  $k^{th}$  model).
• Now calculate the  $OOBScore = \frac{1}{506} \sum_{i=1}^{506} (y^i - y^i_{pred})^2$ . 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 intravels 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 intravel Task 3 Given a single query point predict the price of house. Consider 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] Predict the house price for this point as mentioned in the step 2 of Task 1. A few key points • Remember that the datapoints used for calculating MSE score contain some datapoints that were initially used while training the base learners (the 60% sampling). This makes these datapoints partially seen (i.e. the datapoints used for calculating the MSE score are a mixture of seen and unseen data). Whereas, the datapoints used for calculating OOB score have only the unseen data. This makes these datapoints completely unseen and therefore appropriate for testing the model's performance on unseen data. Given the information above, if your logic is correct, the calculated MSE score should be less than the OOB score. • The MSE score must lie between 0 and 10. • The OOB score must lie between 10 and 35. • The difference between the left nad right confidence-interval values must not be more than 10. Make sure this is true for both MSE and OOB confidence-interval values. Task - 1 Step - 1 Creating samples **Algorithm alt** text • Write code for generating samples In [5]:  $from \ sklearn.tree \ import \ Decision Tree Regressor$ import pandas as pd from math import sqrt import numpy as np def generating\_samples(input\_data, target\_data): '''In this function, we will write code for generating 30 samples ''' # you can use random.choice to generate random indices without replacement # Please have a look at this link https://docs.scipy.org/doc/numpy-1.16.0/reference/generated/numpy.random.choice.html for more details # Please follow above pseudo code for generating samples # return sampled\_input\_data , sampled\_target\_data, selected\_rows, selected\_columns #note please return as lists Selecting\_rows=np.random.choice(np.arange(input\_data.shape[0]),size=303,replace=False) Replacing\_rows=np.random.choice(np.arange(Selecting\_rows.shape[0]),size=203,replace=False) Selecting\_columns=np.random.choice(np.arange(input\_data.shape[1]), size=np.random.randint(4,13), replace=False) sample\_data=input\_data[Selecting\_rows[:,None],Selecting\_columns] target\_of\_sample\_data=target\_data[Selecting\_rows] #Replacted data Replicated\_sample\_data=sample\_data[Replacing\_rows] target\_of\_replicated\_sample\_data=target\_of\_sample\_data[Replacing\_rows] #concantacing the data final\_sample\_data=np.vstack([sample\_data, Replicated\_sample\_data]) final\_target\_data=np.vstack([target\_of\_sample\_data.reshape(-1,1)], target\_of\_replicated\_sample\_data.reshape(-1,1)]) return list(final\_sample\_data), list(final\_target\_data), list(Selecting\_rows), list(Selecting\_columns) In [ ]: **Grader function - 1 </fongt>** In [7]: def grader\_samples(a,b,c,d): length = (len(a)==506 and len(b)==506)sampled = (len(a)-len(set([str(i) for i in a]))==203) $rows_length = (len(c)==303)$ column\_length= (len(d)>=3) assert(length and sampled and rows\_length and column\_length) return True  $a,b,c,d = generating_samples(x, y)$ grader\_samples(a, b, c, d) Out[7]: True Create 30 samples **alt** text **# Use generating\_samples function to create 30 samples** # store these created samples in a list list\_input\_data =[] list\_output\_data =[] list\_selected\_row= [] list\_selected\_columns=[] for i in range(0,30):  $a,b,c,d = generating_samples(x, y)$ list\_input\_data.append(a) list\_output\_data.append(b) list\_selected\_row.append(c) list\_selected\_columns.append(d) **Grader function - 2** def grader\_30(a): assert(len(a)==30 and len(a[0])==506) grader\_30(list\_input\_data) Out[9]: True Step - 2 Flowchart for building tree **alt** text Write code for building regression trees In [10]: list\_of\_all\_models=[] def Tree(input\_data, target\_data): clf=DecisionTreeRegressor(max\_depth=None) clf.fit(input\_data, target\_data) return clf for i in range(0,30):# Forming 30 differnet base trees X=list\_input\_data[i] Y=list\_output\_data[i] T=Tree(X,Y) list\_of\_all\_models.append(T) Flowchart for calculating MSE **alt** text After getting predicted\_y for each data point, we can use sklearns mean\_squared\_error to calculate the MSE between predicted\_y and actual\_y. Write code for calculating MSE In [11]: list\_of\_all\_predcts=[] for i in range(0,506): temp\_y\_pred=[] for j in range(0,30): pred=x[i,list\_selected\_columns[j]].reshape(1,-1) y\_pred=list\_of\_all\_models[j].predict(pred)#predcit y\_value from model j which dont have that point temp\_y\_pred.append(y\_pred) list\_of\_all\_predcts.append((1/30)\*np.sum(np.array(temp\_y\_pred))) MSE=(1/506)\*np.sum((y-list\_of\_all\_predcts)\*\*2) print(f'Mean Square ERROR {MSE}' ) Mean Square ERROR 2.0956678902552084 Step - 3 Flowchart for calculating OOB score **alt** text Now calculate the  $OOBScore = \frac{1}{506} \sum_{i=1}^{506} (y^i - y^i_{pred})^2$ . Write code for calculating OOB score In [12]: #input\_data should be a dataframe #point\_check should be a array with len of columns #columns should be indexes of columns that are present def check\_point(df,pc,col):#Return True if the row is found in that input\_data set df=pd.DataFrame(df) a=np.array(pc[col]) return (df==a).all(1).any() In [13]: list\_of\_new\_y\_pred=[] for i in range(0,506): some\_model\_predcts=[]#used to store the moddel outputs for a point i from dataset X found\_tree=False#Just to check whether a datapoint is found a tree without that datapoint k=0 for j in range(0,30): a=check\_point(list\_input\_data[j],x[i],list\_selected\_columns[j])#check whether a point exists in the jth model if a:#If exists pass pass else: #else pred the value for that datapoint found\_tree=True pred=x[i,list\_selected\_columns[j]].reshape(1,-1) y\_pred=list\_of\_all\_models[j].predict(pred)#predcit y\_value from model j which dont have that point some\_model\_predcts.append(y\_pred)#stroing all the predicted values for ith point list\_of\_new\_y\_pred.append((1/k)\*np.sum(np.array(some\_model\_predcts))) In [14]: OOB=(1/506)\*np.sum((np.array(y)-np.array(list\_of\_new\_y\_pred))\*\*2) print(f'The OOBScore is {OOB}') The OOBScore is 12.288839332109914 Task 2 In [15]: #Repeating task1 For 35 times MSE\_35=[] 00B\_35=[] for  $\underline{in}$  range(0,35): **# Use generating\_samples function to create 30 samples** # store these created samples in a list list\_input\_data\_35 =[] list\_output\_data\_35 =[] list\_selected\_row\_35= [] list\_selected\_columns\_35=[] for i in range(0,30):  $a,b,c,d = generating_samples(x, y)$ list\_input\_data\_35.append(a) list\_output\_data\_35.append(b) list\_selected\_row\_35.append(c) list\_selected\_columns\_35.append(d) list\_of\_all\_models\_35=[] def Tree(input\_data, target\_data): clf=DecisionTreeRegressor(max\_depth=None) clf.fit(input\_data, target\_data) return clf for i in range(0,30): X=list\_input\_data\_35[i] Y=list\_output\_data\_35[i] T=Tree(X,Y) list\_of\_all\_models\_35.append(T) list\_of\_all\_predcts\_35=[] for i in range(0,506): temp\_y\_pred\_35=[] for j in range(0,30): pred=x[i,list\_selected\_columns\_35[j]].reshape(1,-1) y\_pred=list\_of\_all\_models\_35[j].predict(pred)#predcit y\_value from model j which dont have that point temp\_y\_pred\_35.append(y\_pred) list\_of\_all\_predcts\_35.append((1/30)\*np.sum(np.array(temp\_y\_pred\_35)))  $MSE=(1/506)*np.sum((np.array(y)-np.array(list_of_all_predcts_35))**2)$ #print(f'The MSE score is {MSE}') MSE\_35.append(MSE) list\_of\_new\_y\_pred\_35=[] for i in range(0,506): some\_model\_predcts\_35=[]#used to store the moddel outputs for a point i from dataset X found tree=False k=0 for j in range(0,30): a=check\_point(list\_input\_data\_35[j],x[i],list\_selected\_columns\_35[j])#check whether a point exists in the jth model if a:#If exists pass else:#else pred the value for that datapoint found\_tree=True pred=x[i,list\_selected\_columns\_35[j]].reshape(1,-1) y\_pred=list\_of\_all\_models\_35[j].predict(pred)#predcit y\_value from model j which dont have that point some\_model\_predcts\_35.append(y\_pred)#stroing all the predicted values for ith point  $list_of_new_y_pred_35.append((1/k)*np.sum(np.array(some_model_predcts_35)))$ #find mean of the points and store in a new array 00B=(1/506)\*np.sum((np.array(y)-np.array(list\_of\_new\_y\_pred\_35))\*\*2) #print(f'The 00B score is {00B}') 00B\_35.append(00B) print(f'{\_}', end=' ') 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 In [16]: print(00B\_35) print(MSE\_35) [13.851426226748801, 13.448073027656783, 12.558260261587503, 13.672198242439741, 11.884333794890965, 14.067577550053254, 12.649400054228744, 13.641431683249365, 14.363030086444452, 14.277909113694674, 15.456461855139315, 14.44697205495984, 13.489906754108379, 15.092412901532821, 11.605218582644417, 13.885719970579938, 13.84268656940532, 14.307102281613947, 1 4.642076539555324, 16.88380249374791, 14.696349530746028, 14.151338012295108, 13.423881447757719, 13.164984463342053, 15.279193139103945, 12.6425369128406, 13.567132715542039, 15.3  $56227055581536, \ 15.351175802113401, \ 14.267427883373157, \ 11.984782599173311, \ 13.704549773912065, \ 14.70050516610204, \ 12.828463287961071, \ 13.606294762288092]$ [2.25068308629776, 2.190505307910018, 2.5246168900982404, 2.291803683721269, 2.04580164690382, 2.3123905357927095, 2.0121528546332885, 2.758933383527988, 2.46655966231156, 2.582642 7768018343, 2.8019691279005845, 2.463394397935833, 2.3613526053042495, 2.5328084114087726, 1.7465541501976274, 2.252641168713654, 2.1606913665641923, 2.7757870297784604, 2.80189559 7228321, 2.600298145873888, 2.4585574771690326, 2.4981203107672947, 2.3821162139930343, 2.2889391749621817, 2.7737661187114613, 2.281520138544569, 2.2155091128678084, 2.50255868031 4806, 2.427236209925339, 2.2020905717806087, 1.8268674700553502, 2.2974526127824126, 2.8648773425322815, 2.294777592348606, 2.2581529016005457] In [17]: # 00B\_35=[14.404254247090917, 11.356516121696368, 13.231594747137741, 16.010466855943655, 13.756024396771558, 14.092364041828846, 12.000842528238277, 12.228644015500443, 12.7244094 # MSE\_35=[2.4604760679442155, 1.9940492319054686, 2.215918712800015, 2.81720723641689, 2.2205558123505575, 2.2496287623544964, 1.8501277231050526, 2.133561784170205, 2.283687479842 In [18]: def find\_CI(li):#Find confidence interval for the given sampled data m=np.mean(li)se=np.std(li)/sqrt(len(li)) CI=[m-(2\*se), m+(2\*se)]return CI print(f'Confidence Interval of MSE is {find\_CI(MSE\_35)}') print(f'Confidence interval of 00B is {find\_CI(00B\_35)}') Confidence Interval of MSE is [2.2969856011565746, 2.4747871849725342] Confidence interval of 00B is [13.536496023269269, 14.280123553668659] 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. **alt** text Write code for TASK 3 In [19]: **#Predicting the Point using task1** temp\_y\_pred\_=[] xq= np.array([0.18,20.0,5.00,0.0,0.421,5.60,72.2,7.95,7.0,30.0,19.1,372.13,18.60]) for j in range(0,30): pred=xq[list\_selected\_columns[j]].reshape(1,-1) y\_pred=list\_of\_all\_models[j].predict(pred)#predcit y\_value from model j which dont have that point temp\_y\_pred\_.append((y\_pred)) print((1/30)\*np.sum(np.asarray(temp\_y\_pred\_))) 19.32333333333333 Write observations for task 1, task 2, task 3 indetail **#Observations** #For Task-1 #1.In task one we have build a random forest with 30-base trees #2. In this task-1 we have calculated the MSE and 00B for the given RandomForest #In this task we have repeated task-1 for 35 times and got 35 different values of MSE and 00B #With that sample we have calculated the CI #we have predicted the value for a given query point using the First RandomForest model we built in this.