

Bootstrap assignment

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

Every Grader function has to return True.</p>
</div>
<div data-bbox="278 28 365 30" data-label="Section-Header>
<h3>Importing packages</h3>
</div>
<div data-bbox="232 33 745 40" data-label="Text>
<pre>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</pre>
</div>
<div data-bbox="232 43 476 49" data-label="Text>
<pre>In [2]: boston = load_boston()
x=boston.data #independent variables
y=boston.target #target variable</pre>
</div>
<div data-bbox="232 52 274 54" data-label="Text>
<pre>In [3]: x.shape</pre>
</div>
<div data-bbox="232 56 327 59" data-label="Text>
<pre>Out[3]: (506, 13)</pre>
</div>
<div data-bbox="232 62 306 64" data-label="Text>
<pre>In [4]: x[:5]</pre>
</div>
<div data-bbox="232 67 645 99" data-label="Text>
<pre>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, 3.9600e+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, 3.9600e+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,
6.9800e+00, 4.5800e+01, 6.0622e+00, 3.0000e+00, 2.2200e+02,
1.8700e+01, 3.9463e+02, 2.9400e+00],
[6.9950e-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, 3.9600e+02, 6.3300e+00]])</pre>
</div>
<div data-bbox="278 102 322 104" data-label="Section-Header>
<h3>Task 1</h3>
</div>
<div data-bbox="278 107 312 109" data-label="Section-Header>
<h4>Step - 1</h4>
</div>
<div data-bbox="285 112 373 114" data-label="Section-Header>

Creating samples

</div>
<div data-bbox="296 115 573 117" data-label="Text>
<p>Randomly create 30 samples from the whole boston data points</p>
</div>
<div data-bbox="296 116 737 120" data-label="List-Group>

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

</div>
<div data-bbox="296 121 748 127" data-label="Text>
<p>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], consider they are [5, 8, 3, 7] so our final sample will be [4, 5, 7, 8, 9, 3, 5, 8, 3, 7]</p>
</div>
<div data-bbox="285 128 378 129" data-label="Section-Header>

Create 30 samples

</div>
<div data-bbox="296 129 743 137" data-label="List-Group>

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 atleast 3 features/columns/attributes

</div>
<div data-bbox="296 136 726 140" data-label="List-Group>

Note - While selecting the random 60% datapoints from the whole data, make sure that the selected datapoints are all exclusive, repetition is not allowed.

</div>
<div data-bbox="278 143 312 145" data-label="Section-Header>
<h4>Step - 2</h4>
</div>
<div data-bbox="278 150 631 152" data-label="Text>
<p>Building High Variance Models on each of the sample and finding train MSE value</p>
</div>
<div data-bbox="285 155 720 169" data-label="List-Group>

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 $y_{pred}^i = \frac{1}{30} \sum_{k=1}^{30}$ (predicted value of x^i with k^{th} model)
Now calculate the $MSE = \frac{1}{506} \sum_{i=1}^{506} (y^i - y_{pred}^i)^2$

</div>
<div data-bbox="278 172 312 174" data-label="Section-Header>
<h4>Step - 3</h4>
</div>
<div data-bbox="285 177 413 179" data-label="Section-Header>

Calculating the OOB score

</div>
<div data-bbox="296 180 464 182" data-label="Text>
<p>Predicted house price of i^{th} data point</p>
</div>
<div data-bbox="296 181 715 183" data-label="Text>
<p> $y_{pred}^i = \frac{1}{k} \sum_{k=k}^{k=K}$ model which was built on samples not included x^i (predicted value of x^i with k^{th} model).</p>
</div>
<div data-bbox="296 182 549 189" data-label="List-Group>

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

</div>
<div data-bbox="278 192 330 194" data-label="Section-Header>
<h3>Task 2</h3>
</div>
<div data-bbox="285 197 484 199" data-label="Section-Header>

Computing CI of OOB Score and Train MSE

</div>
<div data-bbox="296 198 740 217" data-label="List-Group>

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

</div>
<div data-bbox="278 220 330 222" data-label="Section-Header>
<h3>Task 3</h3>
</div>
<div data-bbox="285 225 527 227" data-label="Section-Header>

Given a single query point predict the price of house.

</div>
<div data-bbox="278 229 778 235" data-label="Text>
<p>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.</p>
</div>
<div data-bbox="278 238 392 240" data-label="Section-Header>
<h3>A few key points</h3>
</div>
<div data-bbox="285 243 768 267" data-label="List-Group>

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.

</div>
<div data-bbox="278 270 340 272" data-label="Section-Header>
<h3>Task - 1</h3>
</div>
<div data-bbox="278 275 312 277" data-label="Section-Header>
<h4>Step - 1</h4>
</div>
<div data-bbox="285 280 373 282" data-label="Section-Header>

Creating samples

</div>
<div data-bbox="278 284 322 286" data-label="Section-Header>
<h4>Algorithm</h4>
</div>
<div data-bbox="278 287 320 290" data-label="Text>
<p></p>
</div>
<div data-bbox="285 292 446 294" data-label="Section-Header>

Write code for generating samples

</div>
<div data-bbox="227 297 783 319" data-label="Text>
<pre>In [17]: 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
selected_rows = np.random.choice(input_data.shape[0], size = 303, replace = False)
replacing_rows = np.random.choice(selected_rows.shape[0],size=203,replace = True)
selected_columns = np.random.choice(13,np.random.choice(np.arange(3,14)),replace=False)
sample_data = input_data[selected_rows[:,None],selected_columns]
target_of_sample_data = target_data[selected_rows]

#Replicating data
Replicated_sample_data = sample_data[replacing_rows]
target_of_replicated_sample_data = target_of_sample_data[replacing_rows]

#Concatenating 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 final_sample_data , final_target_data, selected_rows, selected_columns
#note please return as lists</pre>
</div>
<div data-bbox="278 322 398 324" data-label="Text>
<p>Grader function - 1 </p>
</div>
<div data-bbox="227 327 783 349" data-label="Text>
<pre>In [18]: 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)</pre>
</div>
<div data-bbox="227 352 300 354" data-label="Text>
<pre>Out[18]: True</pre>
</div>
<div data-bbox="285 357 377 359" data-label="Section-Header>

Create 30 samples

</div>
<div data-bbox="278 362 320 364" data-label="Text>
<p></p>
</div>
<div data-bbox="227 367 783 389" data-label="Text>
<pre>In [19]: # 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)</pre>
</div>
<div data-bbox="278 392 360 394" data-label="Text>
<p>Grader function - 2</p>
</div>
<div data-bbox="227 397 783 419" data-label="Text>
<pre>In [20]: def grader_30(a):
assert(len(a)==30 and len(a[0])==506)
return True
grader_30(list_input_data)</pre>
</div>
<div data-bbox="227 422 300 424" data-label="Text>
<pre>Out[20]: True</pre>
</div>
<div data-bbox="278 427 312 429" data-label="Section-Header>
<h4>Step - 2</h4>
</div>
<div data-bbox="278 432 395 434" data-label="Text>
<p>Flowchart for building tree</p>
</div>
<div data-bbox="278 437 320 439" data-label="Text>
<p></p>
</div>
<div data-bbox="285 442 469 444" data-label="Section-Header>

Write code for building regression trees

</div>
<div data-bbox="227 447 783 469" data-label="Text>
<pre>In [59]: from sklearn.tree import DecisionTreeRegressor
from tqdm import tqdm
#store all the trained models
list_of_all_models = []
for i in tqdm(range=30)):
model=DecisionTreeRegressor(max_depth=None)
model.fit(list_input_data[i],list_selected_columns[i])
list_of_all_models.append(model)</pre>
</div>
<div data-bbox="278 472 709 474" data-label="Text>
<p>100% |████████████████████| 30/30 [00:00<00:00, 326.07it/s]</p>
</div>
<div data-bbox="278 477 408 479" data-label="Text>
<p>Flowchart for calculating MSE</p>
</div>
<div data-bbox="278 482 320 484" data-label="Text>
<p></p>
</div>
<div data-bbox="278 487 786 493" data-label="Text>
<p>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.</p>
</div>
<div data-bbox="285 496 429 498" data-label="Section-Header>

Write code for calculating MSE

</div>
<div data-bbox="227 501 783 523" data-label="Text>
<pre>In [60]: pred_Y_datapoint=[] #List to store the final predicted output
for i in range(len(x)):
array_of_Y=[]
for j,model in enumerate(list_of_all_models):
pred=model.predict(x[i][list_selected_columns[j]].reshape(1,-1))
array_of_Y.append(pred)
pred_Y_datapoint.append(np.median(array_of_Y))

MSE = mean_squared_error(y,pred_Y_datapoint)
print('Mean Squared Error is :',MSE)</pre>
</div>
<div data-bbox="278 526 504 528" data-label="Text>
<p>Mean Squared Error is : 9.678463534702134</p>
</div>
<div data-bbox="278 531 312 533" data-label="Section-Header>
<h4>Step - 3</h4>
</div>
<div data-bbox="278 536 436 538" data-label="Text>
<p>Flowchart for calculating OOB score</p>
</div>
<div data-bbox="278 541 320 543" data-label="Text>
<p></p>
</div>
<div data-bbox="278 546 530 559" data-label="Text>
<p>Now calculate the $OOBScore = \frac{1}{506} \sum_{i=1}^{506} (y^i - y_{pred}^i)^2$.</p>
</div>
<div data-bbox="285 562 458 564" data-label="Section-Header>

Write code for calculating OOB score

</div>
<div data-bbox="227 567 783 601" data-label="Text>
<pre>In [61]: Y_pred=[] #List to store the final predicted output
for i,data_point in enumerate(x):
list_Y_values=[] #List to store the datapoints which are not present in the 30 samples
for j,model in enumerate(list_of_all_models):
if i not in list_selected_row[j]:
pred = model.predict(data_point[list_selected_columns[j]].reshape(1,-1))
list_Y_values.append(pred)
Y_pred.append(np.median(list_Y_values))

OOB = mean_squared_error(y,Y_pred)
print('OOB score is :',OOB)

OOB score is : 23.887024138774578</pre>
</div>
<div data-bbox="278 604 330 606" data-label="Section-Header>
<h3>Task 2</h3>
</div>
<div data-bbox="227 609 783 631" data-label="Text>
<pre>In [65]: mse_35_scores=[]
oob_35_scores=[]

for i in range(35):
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)

Training 30 DT Regressors from 30 samples
list_of_all_models = []
for i in range(30):
model=DecisionTreeRegressor(max_depth=None)
model.fit(list_input_data[i],list_output_data[i])
list_of_all_models.append(model)

#Calculating the MSE score
pred_Y_datapoint=[]
for i in range(len(x)):
array_of_Y=[]
for j,model in enumerate(list_of_all_models):
pred=model.predict(x[i][list_selected_columns[j]].reshape(1,-1))
array_of_Y.append(pred)
pred_Y_datapoint.append(np.median(array_of_Y))
MSE = mean_squared_error(y,pred_Y_datapoint)
mse_35_scores.append(MSE)

#Calculating the OOB score
Y_pred=[]
for i,data_point in enumerate(x):
list_Y_values=[]
for j,model in enumerate(list_of_all_models):
if i not in list_selected_row[j]:
pred = model.predict(data_point[list_selected_columns[j]].reshape(1,-1))
list_Y_values.append(pred)
Y_pred.append(np.median(list_Y_values))
OOB = mean_squared_error(y,Y_pred)
oob_35_scores.append(OOB)</pre>
</div>
<div data-bbox="278 634 416 636" data-label="Text>
<p>Confidence Interval:</p>
</div>
<div data-bbox="278 637 783 647" data-label="Text>
<p>Confidence interval is a range of estimates for an unknown population interval,defined as an interval with a lower bound and upper bound. The 95% confidence interval is the most common and it means that there is 95% probability the true population mean will fall within the confidence interval range calculated using the samples.</p>
</div>
<div data-bbox="278 648 786 667" data-label="Text>
<p>Here we are assuming that we dont have the knowledge on population standard deviation,SE is used is to make confidence intervals of the unknown population mean. If the sampling distribution is normally distributed, the sample mean, the standard error, and the quantiles of the normal distribution can be used to calculate confidence intervals for the true population mean.</p>
</div>
<div data-bbox="278 668 681 674" data-label="Text>
<p>95% confidence interval = (sample_mean-2(standard deviation of sample)/sqrt(sample_size)), (sample_mean+2(standard deviation of sample)/sqrt(sample_size))</p>
</div>
<div data-bbox="227 677 783 699" data-label="Text>
<pre>In [72]: #Computing the Confidence Intervals for MSE scores
sample_size = len(mse_35_scores)
sample_mean = np.mean(mse_35_scores)
sample_std = np.std(mse_35_scores)
left_limit = np.round(sample_mean - 2*(sample_std/np.sqrt(sample_size)), 3)
right_limit = np.round(sample_mean + 2*(sample_std/np.sqrt(sample_size)), 3)
print('95% Confidence Interval of MSE is', left_limit, 'to', right_limit)

95% Confidence Interval of MSE is 7.929 to 8.839</pre>
</div>
<div data-bbox="227 702 783 724" data-label="Text>
<pre>In [71]: #Computing the Confidence Intervals for OOB scores
sample_size = len(oob_35_scores)
sample_mean = np.mean(oob_35_scores)
sample_std = np.std(oob_35_scores)
left_limit = np.round(sample_mean - 2*(sample_std/np.sqrt(sample_size)), 3)
right_limit = np.round(sample_mean + 2*(sample_std/np.sqrt(sample_size)), 3)
print('95% Confidence Interval of OOB is', left_limit, 'to', right_limit)

95% Confidence Interval of OOB is 23.254 to 24.26</pre>
</div>
<div data-bbox="278 727 329 729" data-label="Section-Header>
<h3>Task 3</h3>
</div>
<div data-bbox="278 732 367 734" data-label="Text>
<p>Flowchart for Task 3</p>
</div>
<div data-bbox="278 737 776 747" data-label="Text>
<p>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.</p>
</div>
<div data-bbox="278 750 320 752" data-label="Text>
<p></p>
</div>
<div data-bbox="285 755 393 757" data-label="Section-Header>

Write code for TASK 3

</div>
<div data-bbox="227 760 783 782" data-label="Text>
<pre>In [66]: 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]
xq=np.array(xq)
for i in range(30):
y_pred=[]
pred=list_of_all_models[i].predict(xq[list_selected_columns[i]].reshape(1,-1))
y_pred.append(pred)
Predict_Y_datapoint=np.median(y_pred)
print('The house price for the given query point is :',Predict_Y_datapoint)

The house price for the given query point is : 18.2</pre>
</div>
<div data-bbox="278 785 500 787" data-label="Text>
<p>Write observations for task 1, task 2, task 3 indetail</p>
</div>
<div data-bbox="278 790 318 792" data-label="Section-Header>
<h4>Task 1:</h4>
</div>
<div data-bbox="278 795 786 814" data-label="Text>
<p>As we are performing bootstrap aggregation on the given dataset, training the same set of datapoints in the samples will lead to overfitting. As we want our base learners to have high variance and low bias, but only a subset of features get picked for each model and each model is different from each other. Even if the original dataset changes, very few subsets change as the datapoints are randomly picked to form the samples.</p>
</div>
<div data-bbox="278 815 752 825" data-label="Text>
<p>Finally the aggregation of each subset of model is done to achieve the resultant model which tends to reduce variance without impacting the bias.</p>
</div>
<div data-bbox="278 828 318 830" data-label="Section-Header>
<h4>Task 2:</h4>
</div>
<div data-bbox="278 831 783 841" data-label="Text>
<p>Confidence interval is a range of estimates for an unknown population interval,defined as an interval with a lower bound and upper bound. The 95% confidence interval is the most common and it means that there is 95% probability the true population mean will fall within the confidence interval range calculated using the samples.</p>
</div>
<div data-bbox="278 842 786 861" data-label="Text>
<p>Here we are assuming that we dont have the knowledge on population standard deviation,SE is used is to make confidence intervals of the unknown population mean. If the sampling distribution is normally distributed, the sample mean, the standard error, and the quantiles of the normal distribution can be used to calculate confidence intervals for the true population mean.</p>
</div>
<div data-bbox="278 862 681 868" data-label="Text>
<p>95% confidence interval = (sample_mean-2(standard deviation of sample)/sqrt(sample_size)), (sample_mean+2(standard deviation of sample)/sqrt(sample_size))</p>
</div>
<div data-bbox="278 871 318 873" data-label="Section-Header>
<h4>Task 3:</h4>
</div>
<div data-bbox="278 874 758 884" data-label="Text>
<p>Given a query point, we are performing bagging plus column sampling to predict the real value and taking the median on all the prediction to get the final value. The runtime complexity is O(depth * N) where depth = max_depth of each tree and N is the number of base learners.</p>
</div>
</div>