

## 1.

(a)

Because the indicator function  $I(\bullet)$  takes value +1 if its argument is true, and -1 otherwise,  $c_3(x)=I(x<+\infty)$  always predicts +1.

And then, according to conditions of threshold classifier and indicator function, I listed the form of threshold of classifier below:

classifier	threshold		
	$(-\infty, a)$	$[a, b)$	$[b, +\infty)$
C1	-1	+1	+1
C2	+1	+1	-1
C3	+1	+1	+1
$0.1c_3$	0.1	0.2	0.1
f	0.1	2.1	0.1

Therefore, we could know the set represented by  $f=c_1+c_2+0.1c_3$  is  $\{0.1, 2.1\}$ .

(b)

According to the question, we could know that  $f>0$  when  $x$  is in intervals  $(-2, -1) \cup (1, 2)$ . Thus, I listed the form of threshold of classifier below:

classifier	threshold				
	$(-\infty, -2)$	$(-2, -1)$	$(-1, 1)$	$(1, 2)$	$(2, +\infty)$
$C1=I(x>1)$	+1	+1	-1	-1	-1
$C2=I(x<2)$					
$C3=I(x<-1)$	-1	+1	+1	+1	+1
$C4=I(x>-2)$					
f	0	2	0	0	0

Because the example model puts 0 in the case where the result is +1, I add  $C5=I(x<-\infty)$  in the end, which is always equal to -1.

classifier	threshold				
	$(-\infty, -2)$	$(-2, -1)$	$(-1, 1)$	$(1, 2)$	$(2, +\infty)$
$C1=I(x>1)$	+1	+1	-1	-1	-1
$C2=I(x<2)$					
$C3=I(x<-1)$	-1	+1	+1	+1	+1
$C4=I(x>-2)$					

$-C5 = I(x < -\infty)$	+1	+1	+1	+1	+1
f	1	3	1	1	1

Therefore, a linear combination of threshold classifiers to represent two intervals  $(-2, -1) \cup (1, 2)$  is:  $f = c1 + c2 + c3 + c4 - c5$ , where  $c1 = I(x > 1)$ ,  $c2 = I(x < 2)$ ,  $c3 = I(x < -1)$ ,  $c4 = I(x > -2)$  and  $c5 = I(x < -\infty)$ .

2.

(a)

According to the question "assume that  $Y_1, \dots, Y_m$  are identically distributed with  $\text{var}(Y_i) = \sigma^2$  for all  $i$ ", which means  $\text{var}(Y_i)$  doesn't change as the model changes and is always constant. Therefore, we could know that Wagging has the same bias as each individual model.

(b)

(b)

$$\begin{aligned} & \text{var}\left(\sum_{i=1}^m w_i Y_i\right) + 2 \sum_{i=1}^m \sum_{j=1, j \neq i}^m \text{cov}(w_i Y_i, w_j Y_j) \\ &= \sum_{i=1}^m w_i^2 \cdot \text{var}(Y_i) + 2 \sum_{i=1}^m \sum_{j=1, j \neq i}^m w_i \cdot w_j \cdot \text{cov}(Y_i, Y_j) \quad (\text{formula 1}) \end{aligned}$$

\* We have known from question:  $\text{var}(Y_i) = \sigma^2$ ,  $\text{cov}(Y_i, Y_j) = \rho \cdot \sigma^2$  and  $1 \leq i \neq j \leq m$

(formula 1) could be transformed to:  $\sum_{i=1}^m w_i^2 \cdot \sigma^2 + 2 \sum_{i=1}^m \sum_{j=1, j \neq i}^m w_i \cdot w_j \cdot \rho \cdot \sigma^2$  (formula 2)

When  $m = 2$ ,

$$\textcircled{1} (w_1, w_2) = (1, 0)$$

$$\begin{aligned} (\text{formula 2}) &= w_1^2 \cdot \sigma^2 + w_2^2 \cdot \sigma^2 + 2 \cdot w_1 \cdot w_2 \cdot \rho \cdot \sigma^2 \\ &= \sigma^2 \end{aligned}$$

$$\textcircled{2} (w_1, w_2) = \left(\frac{1}{2}, \frac{1}{2}\right)$$

$$\begin{aligned} (\text{formula 2}) &= \frac{1}{4} \sigma^2 + \frac{1}{4} \sigma^2 + \frac{1}{2} \cdot \rho \cdot \sigma^2 \\ &= \left(\frac{1}{2} + \frac{1}{2} \cdot \rho\right) \cdot \sigma^2 \\ &= \frac{1+\rho}{2} \cdot \sigma^2 \end{aligned}$$

$$\textcircled{3} (w_1, w_2) = (0, 1)$$

$$(\text{formula 2}) = 0 + \sigma^2 + 0 = \sigma^2$$

(c)

(c)

$$\text{var}(\bar{Y}) = \sum_i w_i^2 \cdot G^2 + 2 \sum_i \sum_{j \neq i} w_i \cdot w_j \cdot p \cdot G^2$$

Because we need to find weights  $w_1, \dots, w_n$  that minimize the variance  $\bar{Y} = \sum_i w_i Y_i$ , we need to calculate  $\frac{\partial \text{var}(\bar{Y})}{\partial w_i}$

$$\frac{\partial \text{var}(\bar{Y})}{\partial w_i} = 2 w_i \cdot G^2 + 2 \sum_{j \neq i} w_j \cdot p \cdot G^2 = 0$$

$$w_i + \sum_{j \neq i} w_j \cdot p = 0$$

$$w_i = -p \cdot \sum_{j \neq i} w_j$$

Now, we know that  $w_i = -p \cdot \sum_{j \neq i} w_j$

$$\text{Therefore, } \frac{w_i}{w_j} = \frac{-p \cdot \sum_{k \neq i} w_k}{-p \cdot \sum_{k \neq j} w_k} = \frac{1 - w_i}{1 - w_j}$$

Because  $w_j$  is all  $w$  except for  $w_i$ ,  $\sum_{i \neq j} w_j + w_i = 1$

$$\text{Thus, } \frac{w_i}{w_j} = \frac{1 - w_i}{1 - w_j}$$

$$w_i - w_i \cdot w_j = w_j - w_j \cdot w_i$$

$$w_i = w_j$$

Therefore, when all weights are equal to others, we could minimize  $\text{var}(\bar{Y})$

3.

(a)

(1)

```
from sklearn.datasets import fetch_california_housing
from sklearn.model_selection import train_test_split
```

```
housing = fetch_california_housing()
x = pd.DataFrame(housing.data, columns=housing.feature_names)
y = housing.target
print(x.shape)
```

The output is (20640, 8), so we could know that the value of d is 8.

(2)

```
random_seed=42
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size =
0.3,random_state=random_seed)
print(x_train.shape,x_test.shape,y_train.shape,y_test.shape)
```

I set the number of random seed is 42. And then, the outputs of x\_train, x\_test, y\_train, y\_test are (14448, 8), (6192, 8), (14448,) and (6192,) separately.

**(b)**

(1) From question(a), we could know that the shapes of x\_train and x\_test are (14448, 8) and (6192, 8) after using train\_test\_split() function. Therefore, the value of m is 8, which is equal to the value of d. The reason why this result happened is that one of parameters of train\_test\_split() function is max\_features and the default value of it is "auto", which means that max\_features=n\_features. Thus, the value of m is equal to d.

(2)

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error

rfr = RandomForestRegressor(n_estimators=100,random_state=random_seed)
rfr.fit(x_train, y_train)
y_pred_train,y_pred_test=rfr.predict(x_train),rfr.predict(x_test)
train_mse=mean_squared_error(y_train,y_pred_train) #0.03698995328644075
test_mse=mean_squared_error(y_test,y_pred_test) #0.2569942567398473
print(train_mse,test_mse)
```

I use mean\_squared\_error() function to measure accuracy. The mean squared error of an estimator measures the average of the squares of the errors—that is, the average squared difference between the estimated values and the actual value. a smaller MSE means otherwise and it is definitely the preferred and/or desired choice as it shows that your data values are dispersed closely to its central moment (mean); which is usually great. The outputs of the training and test mse are 0.03698995328644075 and 0.2569942567398473 respectively. This means that the average squared difference of train and test between the estimated values and the actual values are relatively small.

**(c)**

(1)

```
def get_correlation(rfr,test):
    base_learners=rfr.estimators_ #list of decision tree regressor
```

```

base_predictions=[]
for i in range(len(base_learners)):
    base_predictions.append(base_learners[i].predict(test))
base_predictions=np.array(base_predictions)
corr=np.corrcoef(base_predictions)

return corr

corr=get_correlation(rfr,x_test)
print(corr)

The output of all the pairwise correlations between the test set predictions of
the 100 trees is:
[[1.          0.75771861 0.76903759 ... 0.74834924 0.75885034 0.77774916]
 [0.75771861 1.          0.75988204 ... 0.74944788 0.75795134 0.76498828]
 [0.76903759 0.75988204 1.          ... 0.75258947 0.7667629  0.77247691]
 ...
 [0.74834924 0.74944788 0.75258947 ... 1.          0.74320362 0.76258888]
 [0.75885034 0.75795134 0.7667629  ... 0.74320362 1.          0.77001229]
 [0.77774916 0.76498828 0.77247691 ... 0.76258888 0.77001229 1.          ]]

```

The details of output could be seen in the “correlation.csv”.

(2)

```

avg_correlation=[]
avg_correlation.append(np.mean(corr))
print(avg_correlation)

```

The average of all the pairwise correlations between the test set predictions of the 100 trees is:

```
[0.7652796942329071]
```

**(d)**

```

train_MSE=[]
test_MSE=[]
avg_corr=[]

```

for m in range(1,9):

```

    rfr1=RandomForestRegressor(n_estimators=100,random_state=random_seed,m
ax_features=m)
    rfr1.fit(x_train,y_train)
    y_train_pred=rfr1.predict(x_train)
    train_MSE.append(mean_squared_error(y_train_pred,y_train))

```

```

y_test_pred=rfr1.predict(x_test)
test_MSE.append(mean_squared_error(y_test_pred,y_test))

corr=get_correlation(rfr1,x_test)
avg_corr.append(np.mean(corr))

result={"training accuracy":train_MSE,"testing accuracy":test_MSE,"average correlation":avg_corr}
result_df=pd.DataFrame(result)
result_df.to_csv("result.csv")

```

I tabulated the training and test accuracies, and the average correlations for all m values, which is listed below:

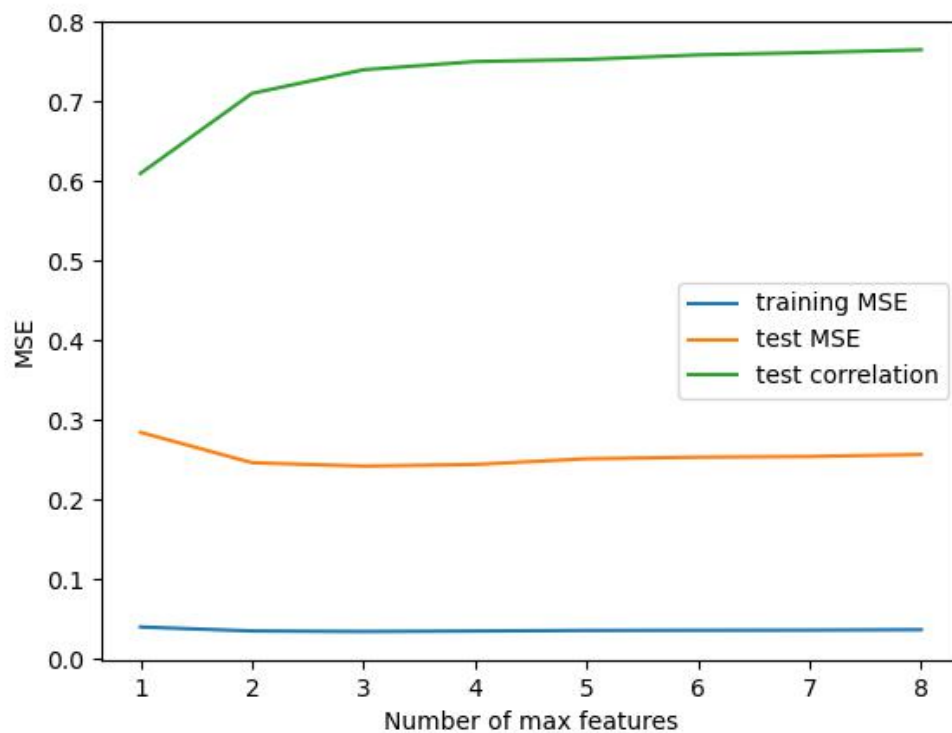
	training accuracy	testing accuracy	average correlation
0	0.040372228	0.284834419	0.610075047
1	0.035390846	0.24668487	0.710542084
2	0.034710454	0.242433586	0.740238506
3	0.035305724	0.244514728	0.750476323
4	0.036032001	0.251657464	0.753193766
5	0.036273677	0.253629474	0.758850978
6	0.036443317	0.254527552	0.761802971
7	0.036989953	0.256994257	0.765279694

```

plt.plot(range(1,9),train_MSE)
plt.plot(range(1,9),test_MSE)
plt.plot(range(1,9),avg_corr)
plt.legend(["training MSE","test MSE","test correlation"])
plt.xlabel("Number of max features")
plt.ylabel("MSE")
plt.show()

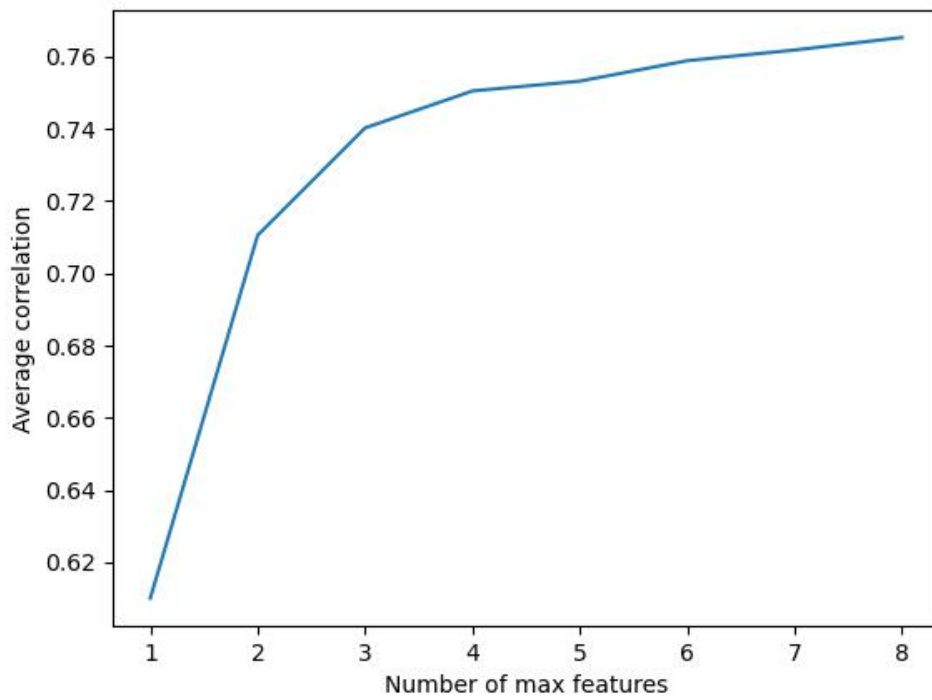
```

I plot the training and test accuracies against m in a single figure, which is listed below:

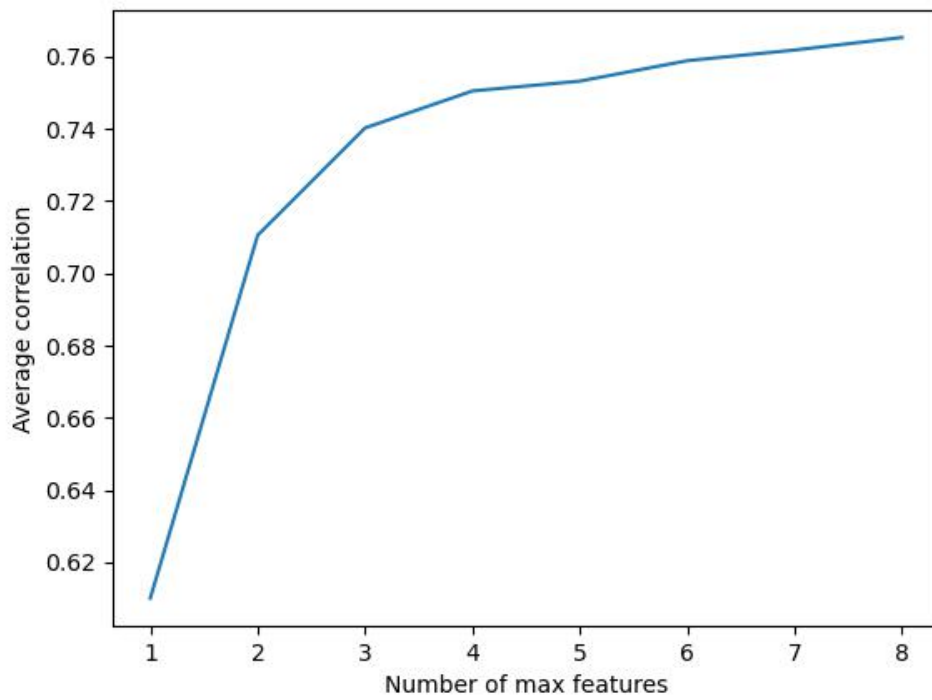


```
plt.plot(range(1,9),avg_corr)
plt.xlabel("Number of max features")
plt.ylabel("Average correlation")
plt.show()
```

I plot the average correlation against  $m$  in another figure, which is listed below:



(e)



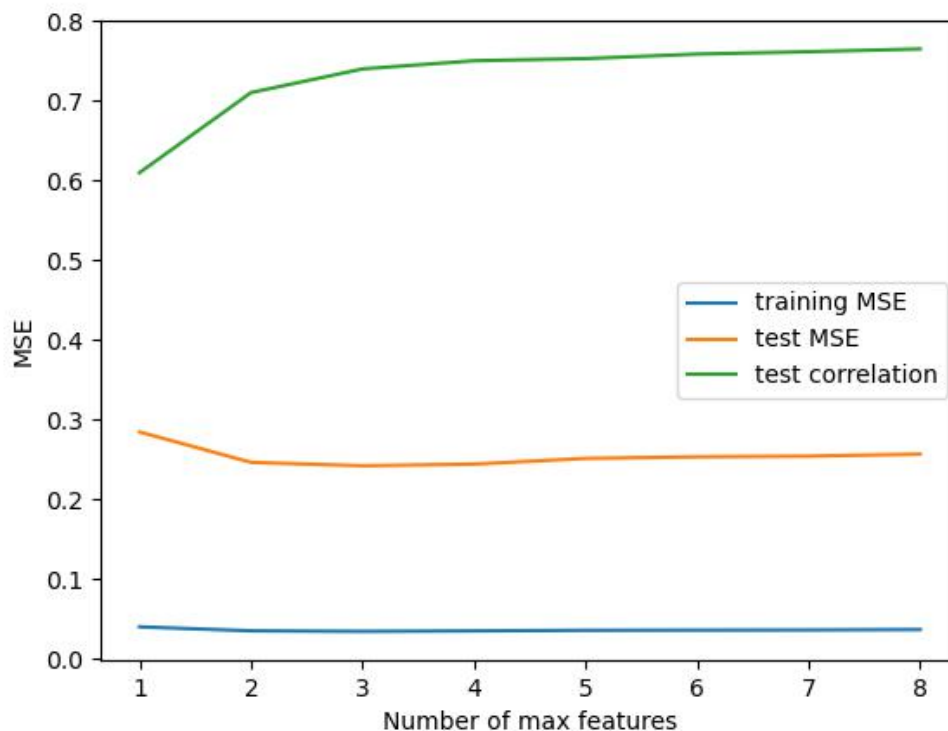
From the plot of the average correlation against  $m$  I listed above, we could find that



the average correlation increases with the increase of  $m$ . If we use  $m=1$  to build 8 decision trees, these 8 decision trees will most likely use different features. Thus, these 8 trained decision trees are completely different, which means that the correlation between them is very low. If we use  $m=8$  to build 8 decision trees, these 8 decision trees are very similar to others, because some training datasets overlap when training them.

(f)

As far as I am concerned, I think the claim of this data scientist is false.



The formula of expected prediction error (EPE):

$$\overbrace{\mathbb{E}((\tilde{Y} - Y)^2)}^{\text{expected prediction error}} = \overbrace{\mathbb{E}((\tilde{Y} - \mathbb{E}(\tilde{Y}))^2)}^{\text{variance}} + \overbrace{(\mathbb{E}(\tilde{Y}) - \mathbb{E}(Y))^2}^{\text{bias (squared)}} + \overbrace{\mathbb{E}((Y - \mathbb{E}(Y))^2)}^{\text{irreducible noise}}$$

According to the plot of the training and test accuracies against  $m$  I listed above, we could find that when  $m=1$ , the correlation and variance between machine learners is the lowest, but the bias is the highest. However, according to the formula of EPE, when we select a appropriate model, we have to think about not only the relatively low variance but also the relatively low bias.

4.

(a)

```
from sklearn.datasets import load_digits
```

```
from sklearn.model_selection import train_test_split
```

```
digits = load_digits()
```

```
random_seed = 999
```

```
X_train, X_test, y_train, y_test = train_test_split(digits.data, digits.target,  
test_size=0.3, random_state=random_seed)
```

**(b)**

First, we should find out the index of OOB data from sample. And then, we could find the corresponding OOB data from the dataset according to these indexes. Thirdly, we use these OOB data to calculate probability of predicted result. Then, we construct a matrix, where each row is the predicted value of each data and each column is the predicted result of the current data in each model. This matrix contains a fixed amount of data. After that, after training with each machine learner, we will find out the probability of predicted result from the corresponding OOB data, and then we add these probabilities into the matrix according to the index of OOB data. Therefore, the amount of memory used should be at most within a constant factor of that for the original dataset size, although we introduce additional variables or data structures if needed. And then, the non-zero data in the matrix represents the current OOB value. Calculate accuracy based on current OOB data and predicted OOB data. After that, we use 1 minus accuracy to get OOB error. Finally, when the number of model is at least 10, we could stop training when the mean of oob error of the current model is greater than that of its previous five models.

**(c)**

I implement the fit function according to the method and thought which I have already described in question(b). I listed code below:

```
from sklearn.metrics import accuracy_score
```

```
def fit(self, X, y, random_state=None):
```

```
    if random_state:
```

```
        np.random.seed(random_state)
```

```
    self.best_n = 0
```

```
    probs_oob = None
```

```
    oob_pred_total = np.zeros((X.shape[0], 10))
```

```
    for i in range(self.n_estimators):
```

```
        estimator = clone(self.base_estimator_)
```

```
        # construct a bootstrap sample
```

```
        bst_index = np.random.choice(range(X.shape[0]), X.shape[0],
```

```

replace=True)
    oob_index = np.setdiff1d(range(X.shape[0]), bst_index)
    bst_X = X[bst_index, :] #training data features
    bst_y = y[bst_index] #training data target

    # train on bootstrap sample
    estimator.fit(bst_X, bst_y) #trained bootstrap sample
    self.estimators_.append(estimator) #list of decision tree regressor

    # compute OOB error
    oob_X = X[oob_index, :]
    # oob_y = y[oob_index]
    oob_pred = estimator.predict_proba(oob_X) #predict result of y

    for j, index in enumerate(oob_index):
        oob_pred_total[index] += oob_pred[j]

    current_oob = list(set(np.nonzero(oob_pred_total)[0]))
    # print(len(current_oob))
    current_oob_pred = np.argmax(oob_pred_total[current_oob,:],
axis=1)
    oob_error = 1 - accuracy_score(y[current_oob], current_oob_pred) #
replace ... with your code
    self.oob_errors_.append(oob_error)

    # save the OOB error and the new model
    self.oob_errors_.append(oob_error)
    self.estimators_.append(estimator)

    # stop early if smoothed OOB error increases (for the purpose of
    # this problem, we don't stop training when the criterion is
    # fulfilled, but simply set self.best_n to (i+1)).
    if (self.best_n == 0) and (i >= 10 and
np.mean(self.oob_errors_[i:i-5:-1]) >= np.mean(self.oob_errors_[i-5:i-10:-1])): #
replace OOB criterion with your code
        self.best_n = (i+1)

```

(d)

```

def errors(self, X, y):
    """
    Parameters
    -----
    X: an input array of shape (n_sample, n_features)
    y: an array of shape (n_sample,) containing the classes for the input

```

examples

Returns

-----

error\_rates: an array of shape (n\_estimators,), with the error\_rates[i] being the error rate of the ensemble consisting of the first (i+1) models.

"""

error\_rates = []

# compute all the required error rates

estimator\_predictions = []

for i, estimator in enumerate(self.estimators\_):

    prob\_pred = estimator.predict\_proba(X)

    if i == 0:

        estimator\_predictions.append(prob\_pred)

    else:

        estimator\_predictions.append(prob\_pred

+

estimator\_predictions[i - 1])

for i in range(len(self.estimators\_)):

    class\_pred = np.argmax(estimator\_predictions[i], axis=1)

    accuracy = accuracy\_score(y, class\_pred)

    error = 1 - accuracy

    error\_rates.append(error)

return error\_rates

**(e)**

(1)

base\_estimator = DecisionTreeClassifier(random\_state=random\_seed)

bagging\_estimator = OOBaggingClassifier(base\_estimator, 100)

bagging\_estimator.fit(X\_train, y\_train, random\_seed)

print(bagging\_estimator.best\_n)

The output of the number of basis models chosen by the OOB error method is 69.

(2)

test\_error=bagging\_estimator.errors(X\_test,y\_test)

oob\_error=bagging\_estimator.oob\_errors\_

plt.plot(range(1,201),test\_error)

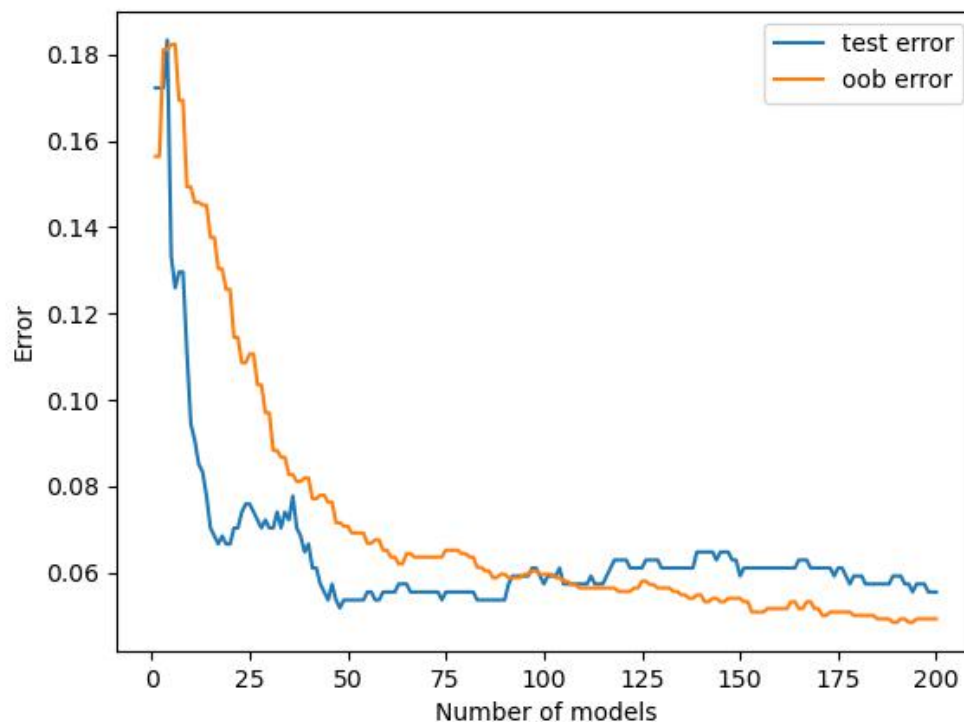
plt.plot(range(1,201),oob\_error)

plt.legend(["test error", "oob error"])

plt.xlabel("Number of models")

```
plt.ylabel("Error")
plt.show()
```

The plot of the OOB error and the test error against the number of models (from 1 to 200) is:

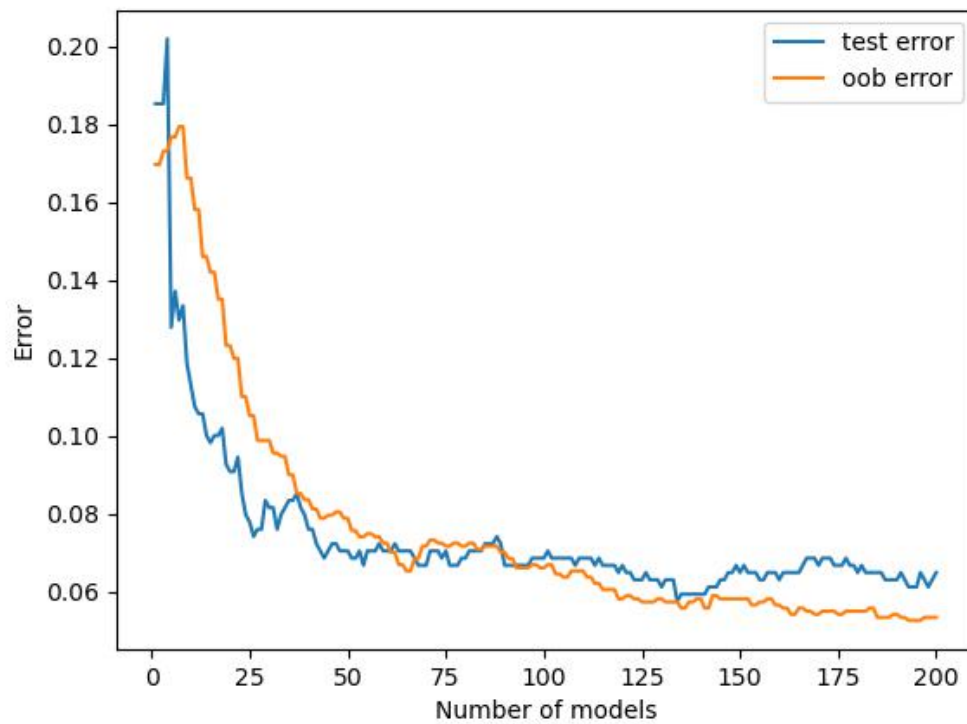


From the plot I listed above, we could know that the trends of test error and oob error are very similar because their trends are both down before the number of model is 100 and oob error always greater than test error. However, after  $m=100$ , the trend of test error is leveling off, but the trend of oob error is still down and oob error always lower than test error.

(f)

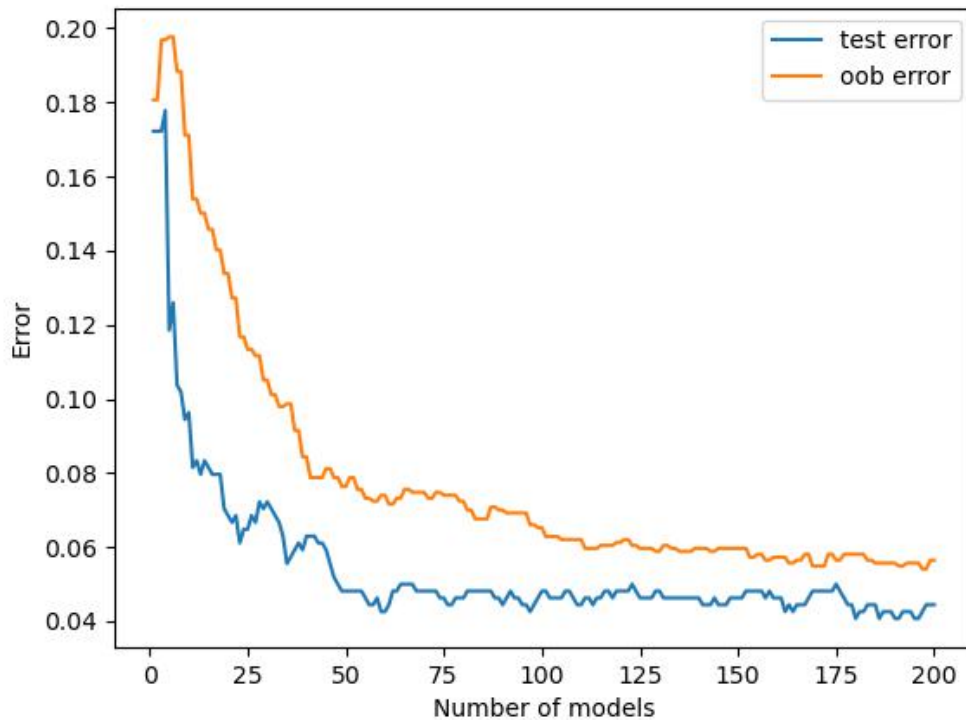
(1) When I set random seeds equal to 111, the number of basis models chosen by the OOB error method is 70.

The plot of the OOB error and the test error against the number of models (from 1 to 200) is:



(2) When I set random seeds equal to 555, the number of basis models chosen by the OOB error method is 66.

The plot of the OOB error and the test error against the number of models (from 1 to 200) is:



Compared three different outputs when random seeds equal to 111, 555 and 999 respectively, we could find that the range of best\_n is roughly from 50 to 70.

The trends of test error and oob error when random seeds equal to 111 and 999 are very similar, they are both down before the number of model is 100 and oob error always greater than test error. However, after  $m=100$ , the trend of test error is leveling off, but the trend of oob error is still down and oob error always lower than test error.

But when random seeds equal to 555, the trends of test error and oob error are both down and oob error always greater than test error.