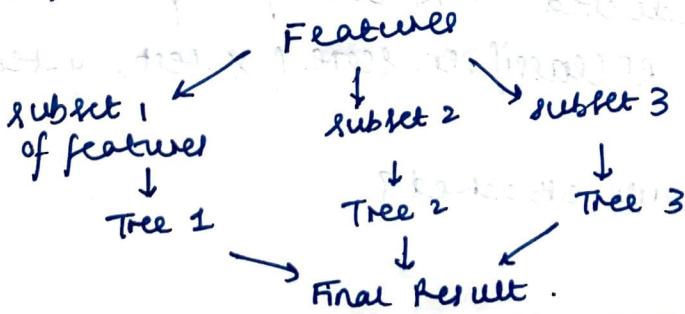


ANSWER 1

- (A) Random Forest is a type of Supervised Machine Learning Algorithm which can be used for tasks like classification & regression.
- It is an 'ensemble' learning algorithm as it aims to use multiple algorithm or multiple instances of similar algorithms in order to perform better.
 - As the name suggests, Random Forest uses multiple decision tree to train on the data. It classifies the class reported by the majority of the decision trees or while regression, it gives the mean of the predicted values by all trees.
 - Decision trees on their own are susceptible to overfitting - using multiple trees which are not correlated and use different subsets of data to train on can help alleviate this problem. low correlation between trees is done by 'feature randomness / bagging' i.e. a different subset of features is used by each tree.
 - Main parameters for Random Forests include no. of trees, no. of features and node sizes.
 - since we are training using many trees, this is time consuming, complex and requires more compute power.



(B) RFClassifier = Random Forest Classifier()
RFClassifier.fit(X-train, y-train)
y-test-predicted = RFClassifier.predict(X-test)
{code & output attached?}

(C) FEATURE IMPORTANCE

Feature Importance tells us which of the features is more relevant to the result i.e. the higher the feature importance the more relevant is that feature/variable to the result of the model.

- Formally, feature importance is computed as the 'mean and standard deviation of accumulation of the impurity decrease within each tree'. (default in sci-kit-learn)
- Another way to calculate importance is permutation importance.
- Feature importance predict which feature had the biggest impact on each prediction and hence, we can build better models by employing better data for higher impact features.

feature_importances = RFClassifier.feature_importances
plt.bar(range(X-train.shape[1]), feature_importances, align='center')
plt.show()
{code & output Attached?}

(D) y-test-predicted = RFClassifier.predict(X-test)
accuracy = RFClassifier.score(X-test, y-test)
{code & output attached?}

```
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn import metrics
import matplotlib.pyplot as plt

# Do not change anything in this code
X, y = make_classification(
    n_samples=1000,
    n_features=12,
    n_informative=4,
    n_redundant=0,
    n_repeated=0,
    n_classes=2,
    random_state=5,
    shuffle=False)
X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=y, random_state=43)

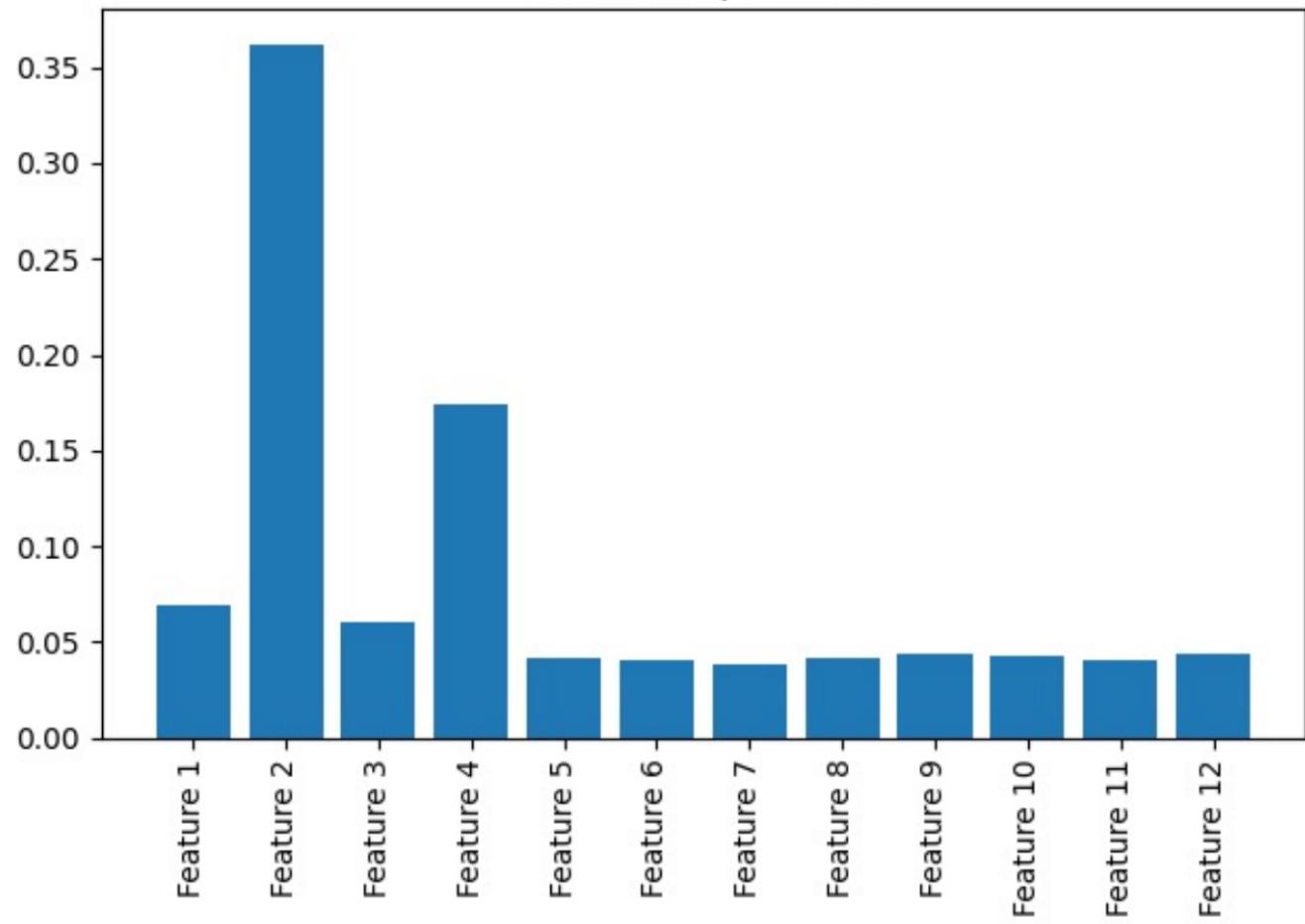
RFClassifier = RandomForestClassifier()
RFClassifier.fit(X_train, y_train)
y_test_predicted = RFClassifier.predict(X_test)

feature_importances = RFClassifier.feature_importances_
print('Feature Importances: \n')
for i in range(0, 12):
    print('Feature '+str(i+1), feature_importances[i])

plt.title('Feature Importance')
plt.bar(range(X_train.shape[1]), feature_importances, align='center')
plt.xticks(range(X_train.shape[1]), ['Feature ' + str(i) for i in
range(1,X_train.shape[1]+1)], rotation=90)
plt.tight_layout()
plt.show()

accuracy = RFClassifier.score(X_test, y_test)
print('Accuracy using score function: ', accuracy)
accuracy2 = metrics.accuracy_score(y_test, y_test_predicted)
print('Accuracy using accuracy_score function: ', accuracy2)
```

Feature Importance



The screenshot shows a dark-themed code editor interface with various icons on the left side. The main area displays a Python script named `q1.py`. The script imports necessary libraries, generates synthetic data, splits it into training and testing sets, trains a Random Forest Classifier, and prints feature importances along with accuracy metrics.

```
1  from sklearn.datasets import make_classification
2  from sklearn.model_selection import train_test_split
3  from sklearn.ensemble import RandomForestClassifier
4  from sklearn import metrics
5  import matplotlib.pyplot as plt
6
7  # Do not change anything in this code
8  X, y = make_classification(
9      n_samples=1000,
10     n_features=12,
11     n_informative=4,
12     n_redundant=0,
13     n_repeated=0,
14     n_classes=2,
15     random_state=5,
16     shuffle=False)
17 X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=y, random_state=43)
18
19 RFClassifier = RandomForestClassifier()
20 RFClassifier.fit(X_train, y_train)
21 y_test_predicted = RFClassifier.predict(X_test)
22
23 feature_importances = RFClassifier.feature_importances_
24 print('Feature Importances: \n')
25 for i in range(0, 12):
26     print('Feature '+str(i+1), feature_importances[i])
27
28 plt.title('Feature Importance')
29 plt.bar(range(X_train.shape[1]), feature_importances, align='center')
30 plt.xticks(range(X_train.shape[1]), ['Feature ' + str(i) for i in range(1,X_train.shape[1]+1)])
31 plt.tight_layout()
32 plt.show()
33
34 accuracy = RFClassifier.score(X_test, y_test)
35 print('Accuracy using score function: ', accuracy)
36 accuracy2 = metrics.accuracy_score(y_test, y_test_predicted)
37 print('Accuracy using accuracy_score function: ', accuracy2)
38
```

The terminal output shows the printed feature importances and accuracy values:

```
Feature Importances:
Feature 1 0.06978644460046514
Feature 2 0.36209875366977057
Feature 3 0.060738332739039454
Feature 4 0.17362540851119743
Feature 5 0.04138894579182502
Feature 6 0.040114503263485354
Feature 7 0.038722709628786205
Feature 8 0.04186723981214087
Feature 9 0.04436293747237699
Feature 10 0.042533356734031795
Feature 11 0.04040247839726527
Feature 12 0.04435888937961594
Accuracy using score function: 0.844
Accuracy using accuracy_score function: 0.844
```

PROBLEMS 3 OUTPUT TERMINAL

> DEBUG CONSOLE

< TERMINAL

→ Assignment3 git:(main) ✘ python3 q1.py
Feature Importances:

Feature 1 0.06978644460046514
Feature 2 0.36209875366977057
Feature 3 0.060738332739039454
Feature 4 0.17362540851119743
Feature 5 0.04138894579182502
Feature 6 0.040114503263485354
Feature 7 0.038722709628786205
Feature 8 0.04186723981214087
Feature 9 0.04436293747237699
Feature 10 0.042533356734031795
Feature 11 0.04040247839726527
Feature 12 0.04435888937961594
Accuracy using score function: 0.844
Accuracy using accuracy_score function: 0.844

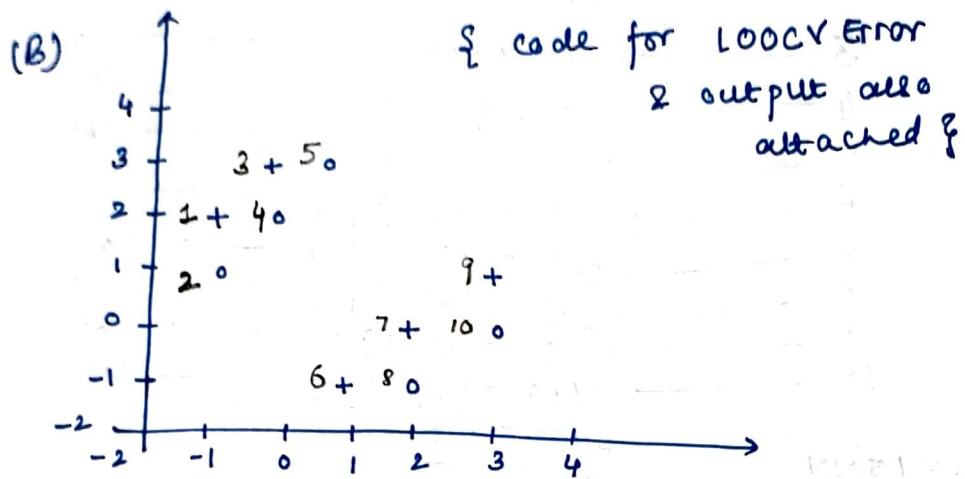
→ Assignment3 git:(main) ✘

zsh +



ANSWER 2

(A) Scode & output Attached



K = 1 / 1NN

point 1 : nearest neighbor points 2 and 4 both '0'
 $1 \rightarrow '0'$, actual '+'

point 2 : nearest neighbor point 1 '+'

$2 \rightarrow '+'$, actual '0'

point 3 : nearest is point 4/5, botn '0'
 $3 \rightarrow '0'$, actual '+'

point 4 : nearest points 1/2 botn '+'
 $4 \rightarrow '+'$, actual '0'

point 5 : nearest point 3 '+'
 $5 \rightarrow '+'$, actual '0'

point 6 : nearest point 8 '0'
 $6 \rightarrow '0'$, actual '+'

point 7 : nearest points 8/10 botn '0'
 $7 \rightarrow '0'$, actual '+'

point 8 : nearest points 6/7 both '+'
 $8 \rightarrow '+'$, actual '0'

point 9 : nearest neighbor 10 '0'
 $9 \rightarrow '0'$, actual '+'

point 10 : nearest neighbor 7/9 both '+'
 $10 \rightarrow '+'$, actual '0'

→ 1NN predicts all points wrong i.e. Error = 1.0

K = 3 / 3NN

point	Nearest Points (class)	Predicted Class
1	2(0), 3(+), 4(0)	'0'
2	1(+), 3(+), 4(0)	'+'
3	1(+), 4(0), 5(0)	'0'
4	1(+), 3(+), 5(0)	'+'
5	1(+), 3(+), 4(0)	'+'
6	7(+), 8(0), 10(0)	'0'
7	6(+), 8(0), 10(0)	'0'
8	6(+), 7(+), 9(0)	'+'
9	7(+), 8(0), 10(0)	'0'
10	9(+), 7(+), 8(0)	'+'
Error = 1.0		

K = 5 / 5NN

point	Nearest Points (class)	Predicted Class
1	2(0), 3(+), 4(0), 5(0), 6(+)	'0'
2	1(+), 3(+), 4(0), 5(0), 6(+)	'+'
3	1(+), 2(0), 4(0), 5(0), 6(+)	'0'
4	2(0), 1(+), 3(+), 5(0), 6(+)	'+'
5	...	'+'
6	...	'0'
7	...	'0'
8	...	'+'
9	...	'0'
10	...	'+'
(computed using code for higher value)		
Error = 1.0		

K = 7 / 7NN

point	Nearest Points (class)	Predicted Class
1	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'+' ✓
2	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'+' X
3	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'+' ✓
4	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'+' X
5	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'+' X
6	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'0' X
7	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'0' X
8	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'0' ✓
9	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'0' X
10	1(+), 2(0), 3(+), 4(0), 5(0), 6(+), 7(+)	'0' ✓

$$\text{Error} = \frac{6}{10} = 0.6$$

$K=9 / 9NN$

As total points are 10,

9NN class simply predicts the majority of the remaining 9 points (10 minus the point itself)

point	predicted class	
1		
2	o	
3	+	
4	o	
5	+	
6	+	
7	o	
8	o	
9	+	
10	o	

All predicted wrong
Error = 1.0.

$\therefore K=7$ leads to the minimum LOOC Validation Error.

$$\begin{cases} K = 1, 3, 5, 9 & \text{Error} = 1.0 \\ K = 7 & \text{Error} = 0.6 \end{cases}$$

```

import numpy as np

# preparing input
x_coords = [-1, -1, 0, 0, 1, 1, 2, 2, 3, 3]
y_coords = [2, 1, 3, 2, 3, -1, 0, -1, 1, 0]
points = [list(x) for x in zip(x_coords, y_coords)]
classification_input = ['+', 'o', '+', 'o', 'o', '+', '+', 'o', '+', 'o']

def euclidean_dist(x1, y1, x2, y2):
    # returns euclidean distance between (x1, y1) and (x2, y2)
    return np.sqrt((x1-x2)*(x1-x2) + (y1-y2)*(y1-y2));

def nearest_neighbors(x_coords, y_coords, classification, x, y):
    # returns a list of points sorted in ascending order by distance from (x,y)
    distances = []
    for i in range(0, len(x_coords)):
        distances.append((euclidean_dist(x, y, x_coords[i], y_coords[i]),
                           x_coords[i], y_coords[i], classification[i]))
    # append a tuple to the list of form (euclidean_dist(x, y, xi, yi), xi,
    # yi, class of (xi, yi))
    distances.sort()
    # sort the list by euclidean_dist i.e. first element of the tuple (default)
    return distances

def get_knn(x_coords, y_coords, classification, x, y, k):
    # returns 'k' nearest neighbors of (x,y)
    return nearest_neighbors(x_coords, y_coords, classification, x, y)[0: k];

def knn_prediction(x_coords, y_coords, classification, x, y, k):
    # checks class of k nearest neighbours and returns the majority class as
    prediction for that point
    knn = get_knn(x_coords, y_coords, classification, x, y, k)
    # print(x, y, knn)
    PLUSes = 0
    Os = 0
    for j in range(0, k):
        if(knn[j][3]=='o'):
            Os += 1;
        else:
            PLUSes += 1
    if(Os>PLUSes):
        return 'o'
    elif(PLUSes>Os):
        return '+'
    else:
        return '='

def knn(x_coords, y_coords, classification, k):

```

```

# returns array of predicted values for all the points in the input set
prediction = []
for i in range(0, len(x_coords)):
    prediction.append(knn_prediction(x_coords, y_coords, classification,
x_coords[i], y_coords[i], k))
return prediction

def getL00CVError(x_coords, y_coords, classification, k):
    wrong_prediction = 0
    for i in range(0, len(x_coords)):
        new_x_coords = x_coords[:i]+x_coords[i+1:]
        new_y_coords = y_coords[:i]+y_coords[i+1:]
        new_classification = classification[:i]+classification[i+1:]
        element_prediction = knn_prediction(new_x_coords, new_y_coords,
new_classification, x_coords[i], y_coords[i], k)
        # print(x_coords[i], y_coords[i], element_prediction, classification[i])
        if(element_prediction != classification[i]):
            wrong_prediction += 1
    return wrong_prediction/len(x_coords)

# Running classification on the training set for an example output]
print('Actual Classes:', classification_input)
predictionk1 = knn(x_coords, y_coords, classification_input, 1)
print('Predicted Classes:')
print('k=1 ', predictionk1)
predictionk3 = knn(x_coords, y_coords, classification_input, 3)
print('k=3 ', predictionk3)
predictionk5 = knn(x_coords, y_coords, classification_input, 5)
print('k=5 ', predictionk5)
predictionk7 = knn(x_coords, y_coords, classification_input, 7)
print('k=7 ', predictionk7)
predictionk9 = knn(x_coords, y_coords, classification_input, 3)
print('k=9 ', predictionk9)

print('\nL00CV Errors: ')
print('k=1 ', getL00CVError(x_coords, y_coords, classification_input, 1))
print('k=3 ', getL00CVError(x_coords, y_coords, classification_input, 3))
print('k=5 ', getL00CVError(x_coords, y_coords, classification_input, 5))
print('k=7 ', getL00CVError(x_coords, y_coords, classification_input, 7))
print('k=9 ', getL00CVError(x_coords, y_coords, classification_input, 9))

```

q2.py 1, M X q1.py M

PROBLEMS 3 OUTPUT TERMINAL

q2.py > ...

```
1 import numpy as np
2
3 # preparing input
4 x_coords = [-1, -1, 0, 0, 1, 1, 2, 2, 3, 3]
5 y_coords = [2, 1, 3, 2, 3, -1, 0, -1, 1, 0]
6 points = [list(x) for x in zip(x_coords, y_coords)]
7 classification_input = ['+', 'o', '+', 'o', 'o', '+', '+', 'o', '+', 'o']
8
9 def euclidean_dist(x1, y1, x2, y2):
10     # returns euclidean distance between (x1, y1) and (x2, y2)
11     return np.sqrt((x1-x2)*(x1-x2) + (y1-y2)*(y1-y2));
12
13 def nearest_neighbors(x_coords, y_coords, classification, x, y):
14     # returns a list of points sorted in ascending order by distance from (x,y)
15     distances = []
16     for i in range(0, len(x_coords)):
17         distances.append((euclidean_dist(x, y, x_coords[i], y_coords[i]), x_coords[i], y_coords[i]))
18     # append a tuple to the list of form (euclidean_dist(x, y, xi, yi), xi, yi)
19     distances.sort()
20     # sort the list by euclidean_dist i.e. first element of the tuple (default)
21     return distances
22
23 def get_knn(x_coords, y_coords, classification, x, y, k):
24     # returns 'k' nearest neighbors of (x,y)
25     return nearest_neighbors(x_coords, y_coords, classification, x, y)[0: k];
26
27 def knn_prediction(x_coords, y_coords, classification, x, y, k):
28     # checks class of k nearest neighbours and returns the majority class as prediction
29     knn = get_knn(x_coords, y_coords, classification, x, y, k)
30     # print(x, y, knn)
31     PLUSes = 0
32     Os = 0
33     for j in range(0, k):
34         if(knn[j][3]=='o'):
35             Os += 1;
36         else:
37             PLUSes += 1
38     if(Os>PLUSes):
39         return 'o'
40     elif(PLUSes>Os):
41         return '+'
42     else:
43         return '='
```

> DEBUG CONSOLE

TERMINAL

zsh +

```
→ Assignment3 git:(main) ✘ python3 q2.py
Actual Classes: ['+', 'o', '+', 'o', 'o', '+', '+', 'o', '+', 'o']
Predicted Classes:
k=1 ['+', 'o', '+', 'o', 'o', '+', '+', 'o', '+', 'o']
k=3 ['o', 'o', 'o', '+', 'o', '+', 'o', '+', '+', '+']
k=5 ['o', '+', 'o', 'o', 'o', '+', '+', '+', '+', '+']
k=7 ['+', '+', '+', '+', 'o', 'o', 'o', 'o', 'o', 'o']
k=9 ['o', 'o', 'o', '+', 'o', '+', 'o', '+', '+', '+']

L00CV Errors:
k=1 1.0
k=3 1.0
k=5 1.0
k=7 0.6
k=9 1.0
→ Assignment3 git:(main) ✘
```

q2.py 1, M X q1.py M

PROBLEMS 3 OUTPUT TERMINAL

q2.py > ...

```
42     else:
43         return '='
44
45 def knn(x_coords, y_coords, classification, k):
46     # returns array of predicted values for all the points in the input set
47     prediction = []
48     for i in range(0, len(x_coords)):
49         prediction.append(knn_prediction(x_coords, y_coords, classification, x_coords[i], y_coords[i], k))
50     return prediction
51
52 def getL00CVError(x_coords, y_coords, classification, k):
53     wrong_prediction = 0
54     for i in range(0, len(x_coords)):
55         new_x_coords = x_coords[:i]+x_coords[i+1:]
56         new_y_coords = y_coords[:i]+y_coords[i+1:]
57         new_classification = classification[:i]+classification[i+1:]
58         element_prediction = knn_prediction(new_x_coords, new_y_coords, new_classification, x_coords[i], y_coords[i], k)
59         # print(x_coords[i], y_coords[i], element_prediction, classification[i])
60         if(element_prediction != classification[i]):
61             wrong_prediction += 1
62     return wrong_prediction/len(x_coords)
63
64
65 # Running classification on the training set for an example output]
66 print('Actual Classes:', classification_input)
67 predictionk1 = knn(x_coords, y_coords, classification_input, 1)
68 print('Predicted Classes:')
69 print('k=1 ', predictionk1)
70 predictionk3 = knn(x_coords, y_coords, classification_input, 3)
71 print('k=3 ', predictionk3)
72 predictionk5 = knn(x_coords, y_coords, classification_input, 5)
73 print('k=5 ', predictionk5)
74 predictionk7 = knn(x_coords, y_coords, classification_input, 7)
75 print('k=7 ', predictionk7)
76 predictionk9 = knn(x_coords, y_coords, classification_input, 9)
77 print('k=9 ', predictionk9)
78
79 print('\nL00CV Errors: ')
80 print('k=1 ', getL00CVError(x_coords, y_coords, classification_input, 1))
81 print('k=3 ', getL00CVError(x_coords, y_coords, classification_input, 3))
82 print('k=5 ', getL00CVError(x_coords, y_coords, classification_input, 5))
83 print('k=7 ', getL00CVError(x_coords, y_coords, classification_input, 7))
84 print('k=9 ', getL00CVError(x_coords, y_coords, classification_input, 9))
```

> DEBUG CONSOLE

TERMINAL

Assignment3 git:(main) ✘ python3 q2.py

Actual Classes: ['+', 'o', '+', 'o', 'o', '+', '+', 'o', '+', 'o']

Predicted Classes:

k=1 ['+', 'o', '+', 'o', 'o', '+', '+', 'o', '+', 'o']
k=3 ['o', 'o', 'o', '+', 'o', '+', 'o', '+', '+', '+']
k=5 ['o', '+', 'o', 'o', 'o', '+', '+', '+', '+']
k=7 ['+', '+', '+', '+', 'o', 'o', 'o', 'o', 'o']
k=9 ['o', 'o', 'o', '+', 'o', '+', '+', '+', '+']

L00CV Errors:

k=1 1.0
k=3 1.0
k=5 1.0
k=7 0.6
k=9 1.0

Assignment3 git:(main) ✘

Ln 68, Col 27 Spaces: 4 UTF-8 LF {} Python Prettier

ANSWER 3

	A ₁	A ₂	A ₃	output y
x ₁	1	0	0	0
x ₂	1	0	1	0
x ₃	0	1	0	0
x ₄	1	1	1	1
x ₅	1	1	0	1

We want to choose the attribute to split in such a way that remaining entropy is minimized.

i.e. we compute Remainder (A_i) for each A_i at split and select one that gives least value.

$$\Rightarrow \text{Remainder}(A_i) = \sum_{k=1}^{2^d} \frac{p_k + n_k}{p+n} B\left(\frac{p_k}{p_k + n_k}\right)$$

K=1 corresponds to A_i=0, K=2 to A_i=1

$$\text{where } B(q) = -(q \log_2 q + (1-q) \log_2 (1-q))$$

P_k is no. of positive exs and

n_k is no. of neg exs for kth A_i

first split

$$\text{Remainder}(A_{i1}) = \frac{4}{5} \left(-\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \right)$$

$$+ \frac{1}{5} \left(-0 - \frac{1}{1} \log_2 \left(\frac{1}{1} \right) \right) =$$

$$\frac{4}{5} \left(\frac{1}{2} + \frac{1}{2} \right) + 0 = 0.8$$

$$\text{Remainder}(A_2) = \frac{3}{5} \left(-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right)$$

$$+ \frac{2}{5} \left(-0 - \frac{2}{2} \log_2 \left(\frac{2}{2} \right) \right)$$

$$= \frac{3}{5} \left(2 \log_2 \left(\frac{3}{2} \right) + \log_2 (3) \right) = 0.55.$$

$$\begin{aligned}
 \text{Remainder}(A_3) &= \frac{2}{5} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) \right) \\
 &\quad + \frac{3}{5} \left(-\frac{1}{3} \log_2 \left(\frac{1}{3}\right) - \frac{2}{3} \log_2 \left(\frac{2}{3}\right) \right) \\
 &= \frac{1}{5} (1+1) + \frac{1}{5} (\log_2 3 + 2 \log_2 \frac{3}{2}) \\
 &= 0.9509 = 0.95
 \end{aligned}$$

Hence, A_2 minimizes remaining entropy.

choose A_2 to split.

All with $A_2 = 0$ are output $y = 0$.
so, we are left with x_3, x_4, x_5 .

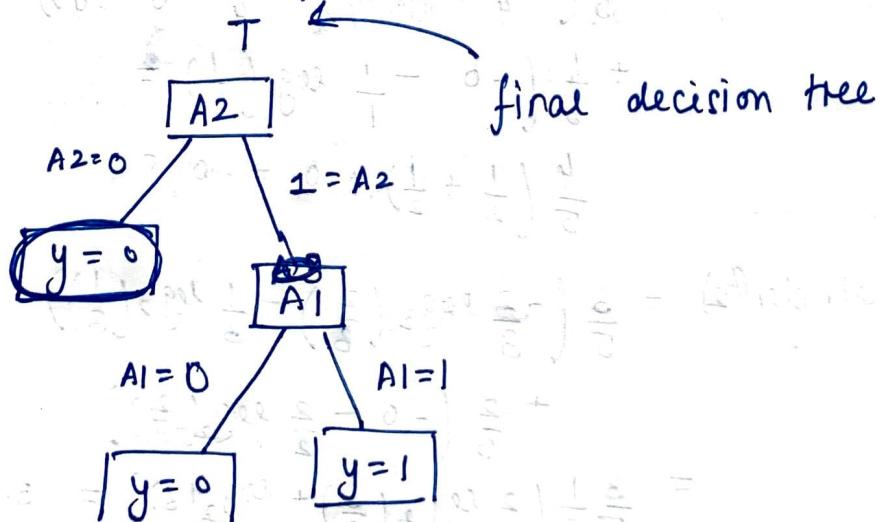
and A_1 and A_3 .

$$\begin{aligned}
 \text{Remainder}(A_1) &= \frac{2}{3} \left(-\frac{1}{2} \log_2 \left(\frac{2}{2}\right) + 0 \right) \\
 &\quad + \frac{1}{3} \left(-0 - \frac{1}{1} \log_2 \left(\frac{1}{1}\right) \right) = 0
 \end{aligned}$$

$$\begin{aligned}
 \text{Remainder}(A_3) &= \frac{1}{3} \left(-\frac{1}{1} \log_2 \left(\frac{1}{1}\right) \right) + \frac{2}{3} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) \right) \\
 &= \frac{1}{3} (0) + \frac{2}{3} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) \right) \\
 &= \frac{1}{3} (0) + \frac{2}{3} \left(-\frac{1}{2} + \frac{1}{2} \right) = \frac{2}{3} = 0.67
 \end{aligned}$$

Hence, A_3 minimizes remaining entropy to 0.

i.e. we split wrt A_3 now.



ANSWER 4

(a) fitted value at $x=7$ is $\hat{y}(7)$

$$\hat{y} = 18 - 0.4x$$

$$\hat{y}(7) = 18 - 0.4(7) = 18 - 2.8 = 15.2$$

\therefore fitted value at $x=7$ is 15.2

(b) residual at $x=3$ and $y=20$
is $20 - \hat{y}(3)$

$$= 20 - 18 + 0.4(3)$$

$$= 2 + 0.4(3) = 3.2$$