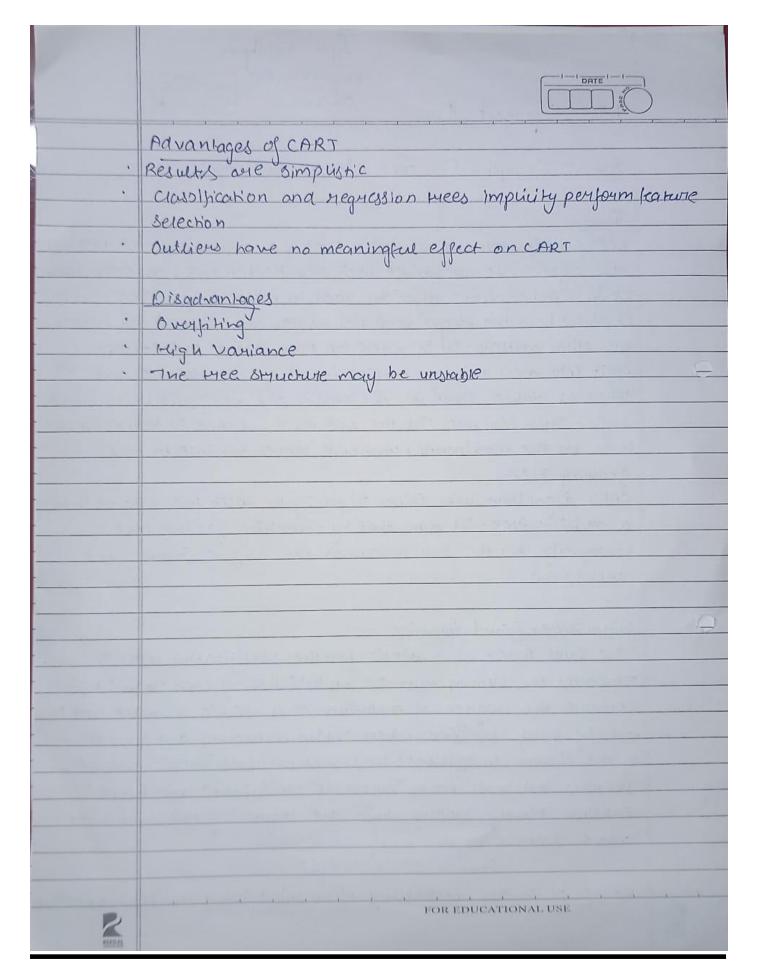
Name: Jigar Siddhpura SAPID: 60004200155

DIV: C/C2 Branch: Computer Engineering

ML - Experiment 4

	Jigar Siddhpunc
	6000421015
	ML - Experiment 4.
	Aim : To implement CART decision tore algorithm
	Theory:
	CART (classification and Reguession Tax) is a variation of the decision
	thee agonith. It can handle both classification and neguession tasks
-	It is a propredictive algorithm used in machine learning and it
	explains how the larger variables values can be predicted based
	on other matters. It is a decision tree where each fork is
	split into a pre two by considering the best attribute and
	threshold value. Furthur, the subsets are also oplit using same
	logic. This continues till the last pure sub-set is found in the
	wee on the maximum number of leaves possible in that
	growing trees.
	CART algorithm uses bring Impurity to split the dataset into
	a decision tree. It does that by searching gove the best
	homogenity for the sub nodeswith the help of Ginni index
	(Hi KAIAKON
	Ginn Index / Gini impurity
	The crini index 16 a metric for the classification tasks in CART.
	The stayes the sum of squared probabilities of each class. It
	computes the cleaner of promotion of a specific the testing
	words being classified when chosen marriaged to the
	of the fring coefficient. It works on early
	around ex our comes either succession er fatter
	conducts binary splitting only. The degree of the cini index
	way en from o to 1
	Cuni c 1 - 2 ; " (P1)
	where pi to the probability of an object being dossified to a particular
ALCOHOL:	



Slep 1:-	Jijor Siddhpura 6000 4 21015 5 DATE CALL CHIM INDEX (Tin) 2 S. Gini S. Gini S. Gini Index of Target Chim Index (Tin) 2 S. Gini S. Chim Index of Target
Step 2:-	Subset in the attribute. [COIPA] COIPA Job offer = yes Job offer = No = 9 3
	2 8 9 0 2) Cuni - Inclex (T. Garpa - L = 9, >, 8) -1 - (2) 2 - (1) 2 2 0.2194
	Gini - Irdex (T. CGPA & EST = 1-(0)2-(3)2 =0 Gini Inaen (T., CGPA & [>9, >,8, <8] = \$ x0.2194+2 x0 10 10 10 10
PEGAL PASSION	FOR EDUCATIONAL USE

<u>_</u>	DATE

	I DATE I
2)	GUNI - Index (T. & CGIPA = (> 9, (8 y = 1-/3) 2-/3) 2-0.5
	Orini Index (T1, COPA = [48]) =1-(4)2-(0)2:0
	Güni Index (T. (GPA - (139, <8) >83 . 6 x0.5+4x0.0.3
3)	Crini-Index (T, CGPA 178, <8) =1-(7) =-(2) = 0.445
	Gini Index (T, COIPA : [>9]) = 1-(3)2-(1)2,0.375
	Cini Imaex (T), CGPA (>, 8, <8 } >, 9), 6 x 0.445 +4/6 x0.32
	Sabbets Cini Index
	>9, >8 <8 0.1955 → best spurring subject
	79, 18 78 0.3 POR CONPA
	-> 8, <8 >9 O.417
	Quini (CGPA) 2 Ovini (T) - Ovini (T, COPA)
	z 0.42 -0. x 355 2 0.24.45
	Ina Interactivenes
	Interactive ness Job offer ? Mo
	yes s
	No 2 2

	DATE 1—1				
	Guini Index (T, Interactiveness (fyes, Not) = 6 (0.28)+4 (0.5) c				
	Guini_ Index (7, Interactiveness - [yesis) = 1- (5)2-(1)2,0-28				
	Orini Index (7, Inletactiveness, INOG): 1-(2)(1)2,0.5				
	A (limi / Interactiveness) = Crinil T) - Crinil T, Interactivenessy = 0.42-0.368				
	3 0.027				
	Pradical knowledge Job offer : No				
	Very good &				
	Average 1 2				
	Subseld Crini- Index				
	(very good, Grood) Averlage 0.3054 -> best- spuitting				
	(very good, average) around 0.40 subset joy				
0	(crood, Average) vory good 0. @ 3750 practical knowledge				
	Downi (Practical knowledge) e Grini(T) - Grini (T, Practical knowledge) = 0.42 - 0.3054				
	2 0.1146				
	Communication stills				
	Comm Sixily Job off or 40				
	6100 d				
	Moderate 3 6				
	1 ² 00 × 0 2				
2	FOR EDUCATIONAL USE				

			DATE 1-1	, ,		
	A Guini (comm & su) + Guini(T) - Guini (T, comm skiu)					
	1	2 6.42 -0	1755			
	Aldribute	Gilni Index	O(MIN)			
	COPP	0:1955				
	Interactiveness		0.025			
	Practical Know					
		0.1322		10		
	Cap	A	W. B. BURNER			
	39,7,8/	<8				
	29, 2,8/ ≥9,2,8/	NO	Territoria de la constitución de			
lep 4	Attribute	Gini Index	DGi ni			
1	Interactiveness	0.056				
To the second	Practical know	0.125	0.0934			
	Comm skicu	0	0.2184			
		Maria de la companya		0		
		- 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1				
R	FOR EDUCATIONAL USE					

```
import pandas as pd
import numpy as np

def variable_count(att):
    types = pd.unique(att)
    no_of_types = len(types)
    counts = att.value_counts() # count of each unique attr
    return no_of_types, counts, types
```

```
lef gini_of_attribute(no_of_types, counts, rows, cla, types, att1, cl):
   gini_a = 0
   type_cl_count = 0
   type_count = 0
   gini = []
   div_index = 0
   if no_of_types == 2:
       for i in range(len(types)):
           temp = df.loc[df[att1.name] == types[i]]
           type_count = len(temp)
           p = 1
           for j in range(len(cla)):
               temp = df.loc[(df[att1.name] == types[i]) & (df[cl.name] == cla[j])]
               type_cl_count = len(temp)
               p -= pow((type_cl_count/type_count), 2)
           gini_a += (type_count/rows) * p
   elif no_of_types > 2:
       for i in range(no_of_types):
           temp1 = df.loc[df[att1.name] == types[i]]
           temp2 = df.loc[df[att1.name] != types[i]]
           type_count1 = len(temp1)
           type_count2 = len(temp2)
           p1 = 1
           p2 = 1
           for j in range(len(cla)):
               temp3 = df.loc[(df[att1.name] == types[i]) & (df[cl.name] == cla[j])]
               type_cl_count1 = len(temp3)
               p1 -= pow((type_cl_count1/type_count1), 2)
               temp4 = df.loc[(df[att1.name] != types[i]) & (df[cl.name] == cla[j])]
               type_cl_count2 = len(temp4)
               p2 -= pow((type_cl_count2/type_count2), 2)
           gini.append((type_count1/rows) * p1 + (type_count2/rows) * p2)
           gini_a = min(gini)
           div_index = gini.index(gini_a)
   return gini_a, div_index
```

```
df = pd.read_csv('/content/gdrive/MyDrive/ML/CART.csv')
col = list(df.columns.values.tolist())
```

```
cl = df.iloc[:, -1]
# cla = [yes,no]
no_of_types, counts, cla = variable_count(cl)
rows = len(cl)
gini = 1 - (pow((counts[0]/rows), 2) + pow((counts[1]/rows), 2))
print(gini)
```

```
gini_a = []
div = []
t = []
att = len(df.columns) - 1
import pandas as pd
from sklearn import tree
rom sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix, accuracy_score
def plot_decision_tree():
 df['Age'] = df['Age'].apply(lambda x: 1 if x == 'youth' else (2 if x == 'middle' else 3))
 df['Income'] = df['Income'].apply(lambda x: 1 if x == 'low' else (2 if x == 'medium' else 3))
 df['Student'] = df['Student'].apply(lambda x: 1 if x == 'no' else 2)
 df['Credit_Rating'] = df['Credit_Rating'].apply(lambda x: 1 if x == 'fair' else 2)
 df['Buys_Computer'] = df['Buys_Computer'].apply(lambda x: 1 if x == 'no' else 2)
 X = df.iloc[:, 0:3]
 y = df.iloc[:, -1]
 clf = tree.DecisionTreeClassifier()
 clf.fit(X, y)
 tree.plot_tree(clf)
for i in range(att):
   att1 = df.iloc[:, i]
   no_of_types, counts, types = variable_count(att1)
   t.append(types)
   gini_a1, div_index = gini_of_attribute(no_of_types, counts, rows, cla, types, att1, cl)
   gini_a.append(gini_a1)
   div.append(div_index)
delta_gini = <mark>list(map(lambda</mark> item : gini - item, gini_a))
print(delta_gini) # highest delta gini wala lenge
index = delta_gini.<mark>index(max</mark>(delta_gini))
print("\n")
print(col[index], "is the root variable")
  [0.10204081632653056, 0.01632653061224476, 0.09183673469387743, 0.030612244897959162]
```

Age is the root variable

