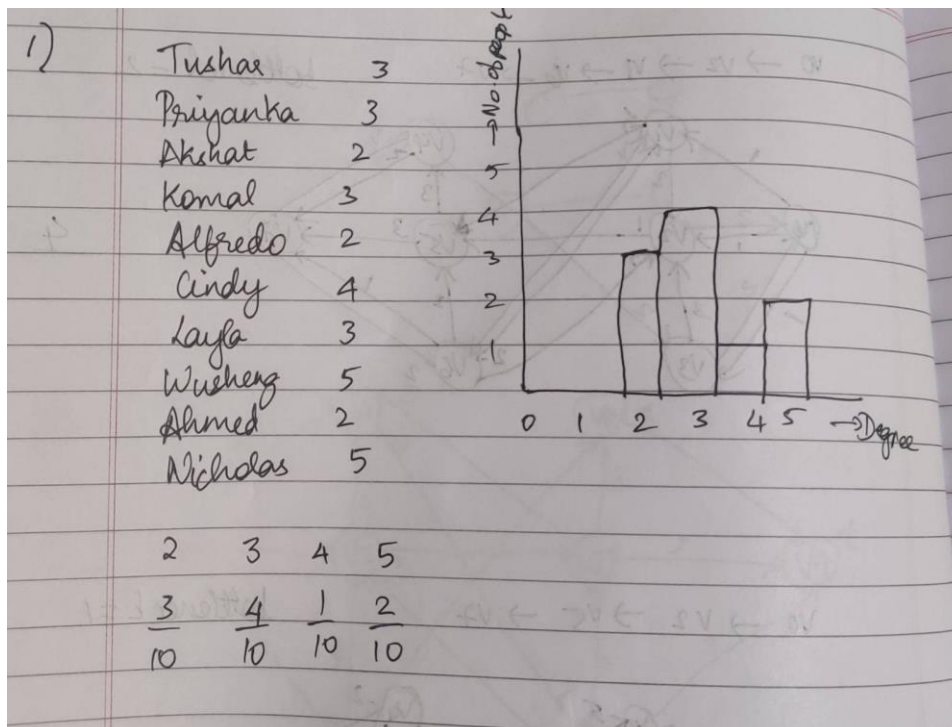
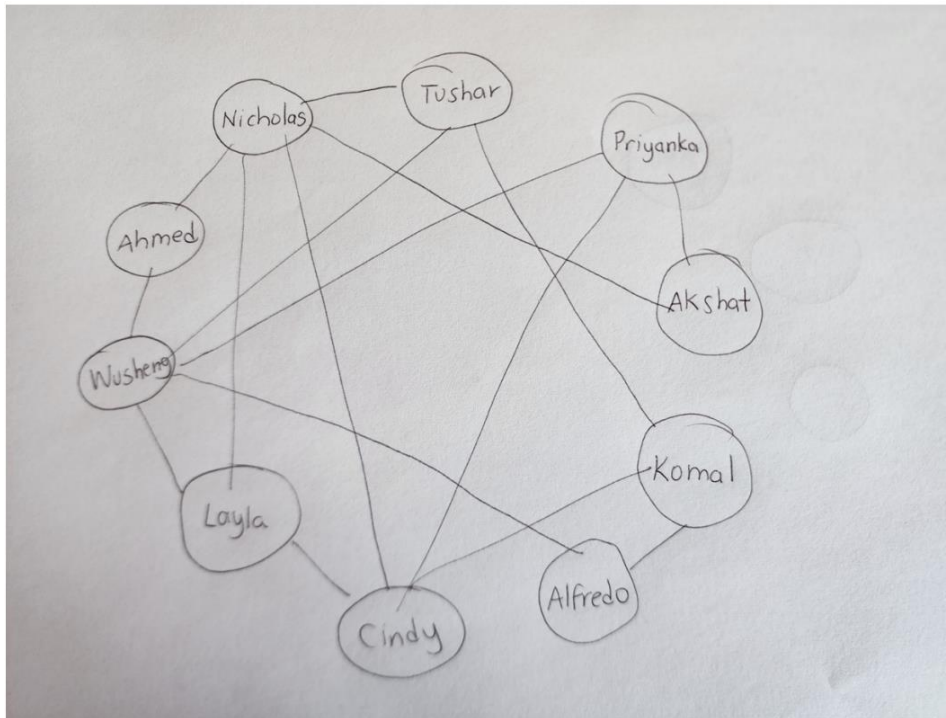
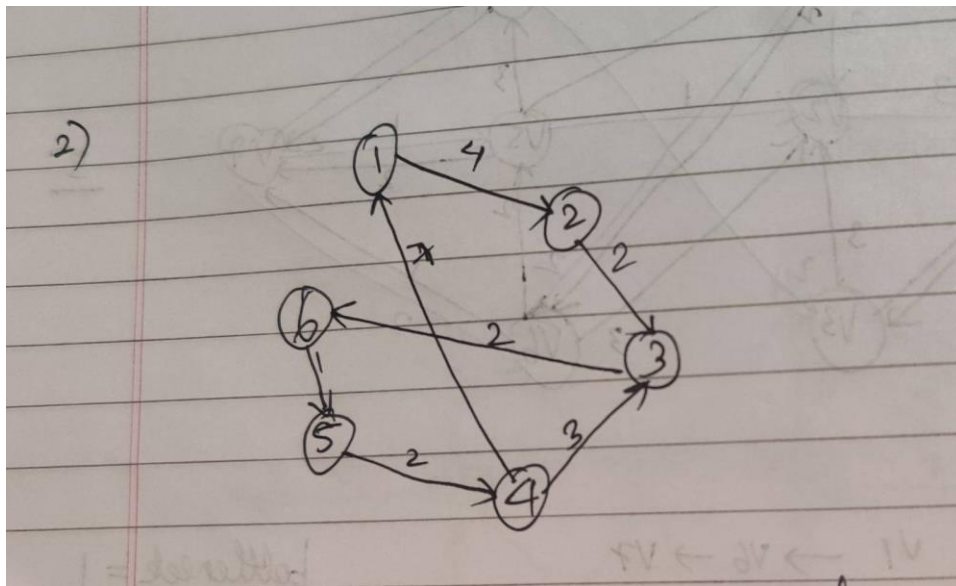


1. Given the friendship graph below, calculate and plot the degree distribution of the graph. Be sure to label the plot axes.



2. Draw the graph specified in the adjacency matrix below. Is this graph connected? If yes, is it weakly connected or strongly connected?

$$\begin{bmatrix} 0 & 4 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 2 \\ 7 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$


Strongly connected.

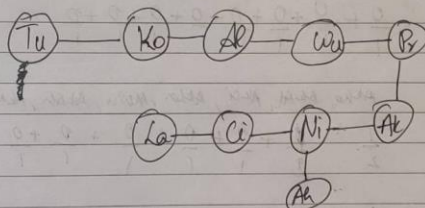
3. Use Dijkstra's or Prim's Algorithm to create a shortest path table for the friend graph from problem 1. What is the diameter of this graph? Show a minimum spanning tree of this graph as an adjacency matrix.

3) Assuming Tushar is node 1, and following the syntax from answer 1, the shortest path table is shown in the next page.

	Tu	Pr	Ak	Ko	Al	Ci	La	Wu	Ah	Ni
Tu	0	2	2	1	2	2	2	1	2	1
Pr	2	0	1	2	2	1	2	1	2	2
Ak	2	1	0	3	3	2	2	2	2	1
Ko	1	2	3	0	1	1	2	2	3	2
Al	2	2	3	1	0	2	2	1	2	3
Ci	2	1	2	1	2	0	1	2	2	1
La	2	2	2	2	2	1	0	1	2	1
Wu	1	1	2	2	1	2	1	0	1	2
Ah	2	2	2	3	2	2	2	1	0	1
Ni	1	2	1	2	3	1	1	2	1	0

The diameter of this graph is 3.

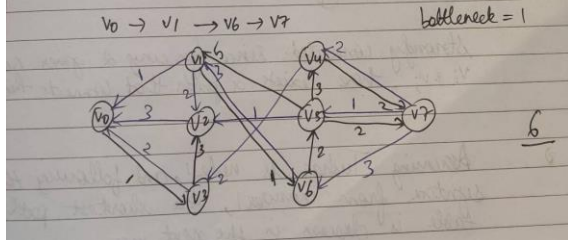
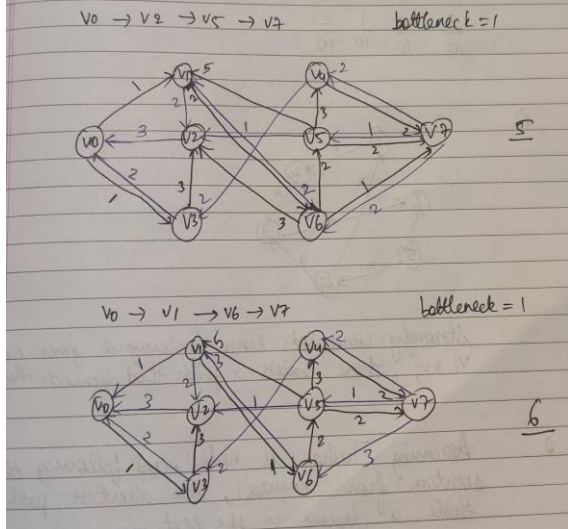
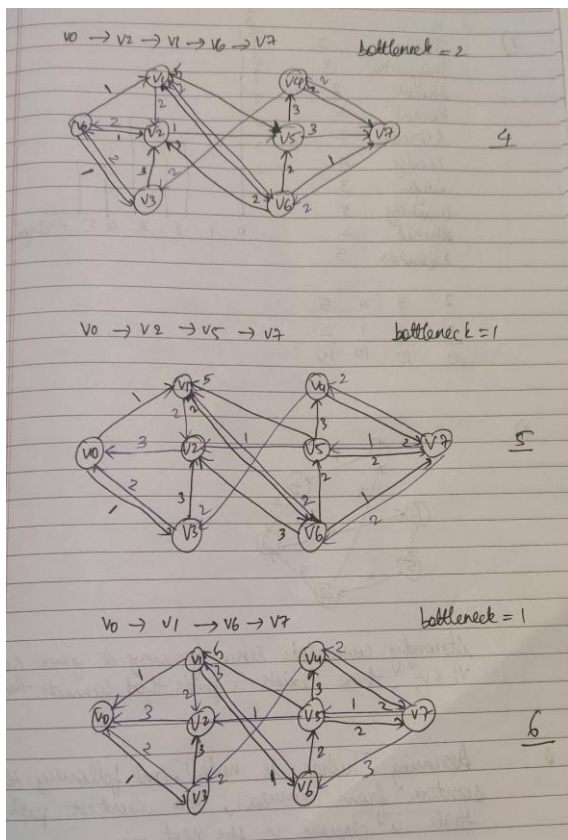
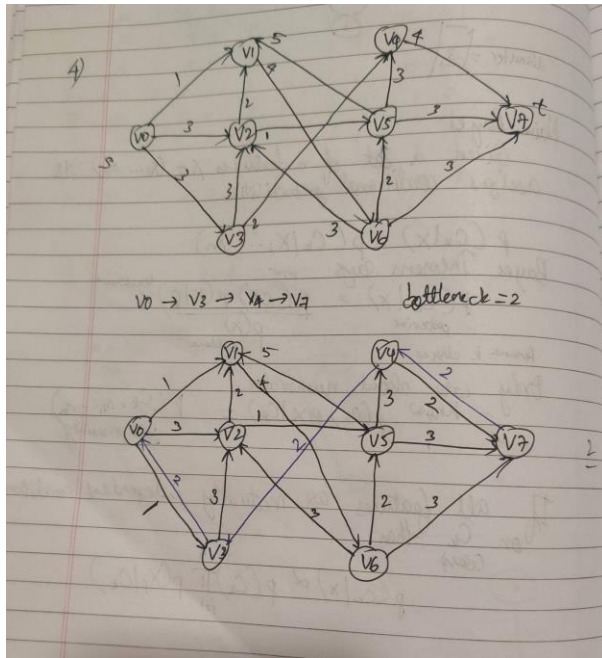
Using Prim's algorithm, MST can be defined:



Adjacency matrix	Tu	Pr	Ak	Ko	Al	Ci	La	Wu	Ah	Ni
Tu	0	0	0	1	0	0	0	0	0	0
Pr	0	0	1	0	0	0	0	1	0	0
Ak	0	1	0	0	0	0	0	0	0	1
Ko	1	0	0	0	1	0	0	0	0	0
Al	0	0	0	1	0	0	0	1	0	0
Ci	0	0	0	0	0	0	1	0	0	1
La	0	0	0	0	0	1	0	0	0	0
Wu	0	1	0	0	1	0	0	0	0	0
Ah	0	0	0	0	0	0	0	0	0	1
Ni	0	0	1	0	0	1	0	0	1	0

4. Given the following flow network from v_0 (source) to v_7 (sink), use the Ford-Fulkerson algorithm to determine the maximum flow. Provide the resulting flow and draw and label the flow network and residual network. Edge list:

$\{(v_0, v_1, 1), (v_0, v_2, 3), (v_0, v_3, 3), (v_1, v_6, 4), (v_2, v_1, 2), (v_2, v_5, 1), (v_3, v_2, 3), (v_3, v_4, 2), (v_4, v_7, 4), (v_5, v_1, 5), (v_5, v_4, 3), (v_5, v_7, 3), (v_6, v_2, 3), (v_6, v_5, 2), (v_6, v_7, 3)\}$



5. For the friendship graph in problem 1, calculate the degree centrality, betweenness centrality and closeness centrality for each node. Provide a table showing the rank of each node for each measure.

5)

	D.C - Rank	B.C - Rank	C.C - Rank
Tushar	3/5 4	1.5 - 3	0.6 - 4
Prinanka	3/5 4	0.5 - 6	0.64 - 4
Atishat	2/5 8	0.5 - 6	0.5 - 9
Komal	3/5 4	0.66 - 9	0.52 - 7
Alfredo	2/5 8	0.16 - 10	0.47 - 10
Cindy	4/5 3	1.5 - 3	0.64 - 2
Layla	3/5 4	1 - 5	0.6 - 4
Wu Sheng	5/5 1	3.5 - 1	0.69 - 1
Shruti	2/5 8	0.5 - 6	0.52 - 7
Nicholas	5/5 1	2.33 - 2	0.64 - 2

6. For the friendship graph in problem 1, what is the clustering coefficient?

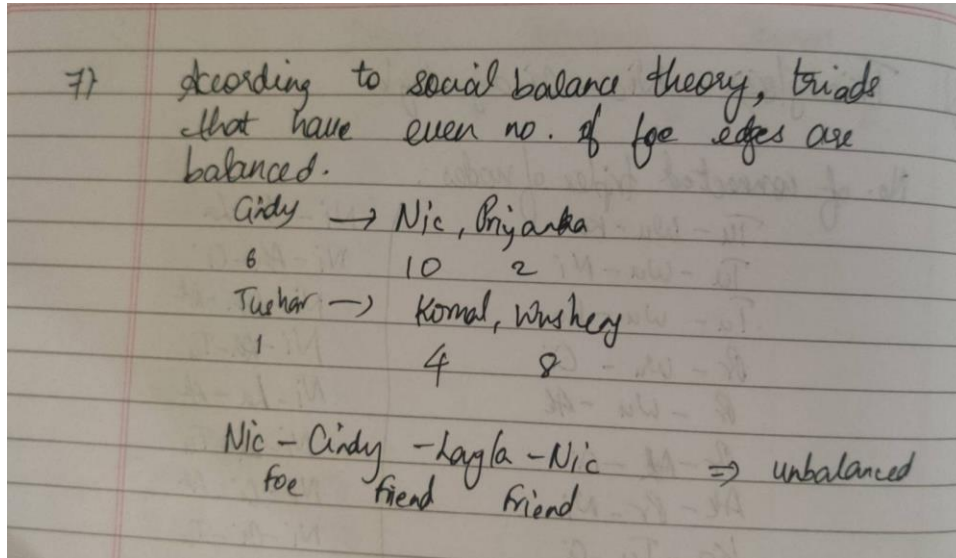
6) Triangles: Nic - Cindy - Layla

No. of connected triples of nodes:

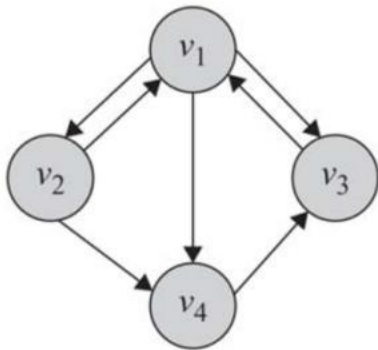
Tu - Wu - Ko	Ni - Ah - La
Tu - Wu - Ni	Ni - Ah - Ci
Tu - Wu - Ko	Ni - Ah - At
Pr - Wu - Ci	Ni - Ah - Tu
Pr - Wu - Ah	Ni - La - Ah
Pr - Ah - Ci	Ni - La - Tu
At - Pr - Ni	Ni - Ci - Ah
Ko - Tu - Ci	Ni - Ci - Tu
Ko - Ci - At	Ni - At - Tu
Ko - Tu - At	
At - Ko - Wu	
Ci - Ko - Pr	
Ci - Ko - Ni	
Ci - Ko - La	
Ci - Pr - Ni	
Ci - Pr - La	
Ci - Ni - La (x 2)	
La - Ci - Wu	
La - Ni - Wu	
Wu - La - At	
Wu - La - Pr	
Wu - La - Tu	
Wu - La - Ah	
Wu - At - Pr	
Wu - At - Tu	
Wu - At - Ah	
Wu - Pr - Tu	
Wu - Pr - Ah	
Wu - Tu - Ah	
Pr - Wu - Ni	

$\frac{1 \times 8}{41} = \frac{8}{41} = 0.097$

7. For the friendship graph in problem 1, assume that Cindy is a foe of Nicholas and Priyanka, Tushar is a foe of Wusheng and Komal and all the other edges in the graph represent friendship. According to social balance theory, is this new friend/foe graph balanced?



8. Calculate PageRank values for the graph below when • $a=1, b=0$ • $a=b=0.3$ • $a=0.85, b=1$ • $a=0, b=1$
Discuss the effects of different values of a and b for this particular problem.



Networkx's pagerank function output that doesn't factor β :

```
PageRank for each node: alpha: 1 beta: 0
Node: 1 PageRank: 0.33333260221557337
Node: 2 PageRank: 0.11111141117268056
Node: 3 PageRank: 0.33333369598839135
Node: 4 PageRank: 0.22222229062335402
PageRank for each node: alpha: 0.3 beta: 0.3
Node: 1 PageRank: 0.25876513
Node: 2 PageRank: 0.20087667499999998
Node: 3 PageRank: 0.2792185175
Node: 4 PageRank: 0.2611396775
PageRank for each node: alpha: 0.85 beta: 1
Node: 1 PageRank: 0.31409230948584155
Node: 2 PageRank: 0.1264931095706823
Node: 3 PageRank: 0.32540314534572534
Node: 4 PageRank: 0.23401143559775067
PageRank for each node: alpha: 0 beta: 1
Node: 1 PageRank: 0.25
Node: 2 PageRank: 0.25
Node: 3 PageRank: 0.25
Node: 4 PageRank: 0.25
```

Custom pagerank function that implements α and β :

```
PageRank for each node: alpha = 1 , beta = 0
Node: 1 PageRank: 0.33333333333333304
Node: 2 PageRank: 0.11111111111111127
Node: 3 PageRank: 0.3333333333333336
Node: 4 PageRank: 0.2222222222222215
PageRank for each node: alpha = 0.3 , beta = 0.3
Node: 1 PageRank: 0.25171554611551344
Node: 2 PageRank: 0.2273507066416142
Node: 3 PageRank: 0.26258067151376524
Node: 4 PageRank: 0.25835307572910704
PageRank for each node: alpha = 0.85 , beta = 1
Node: 1 PageRank: 0.25270984543118913
Node: 2 PageRank: 0.22182022457443407
Node: 3 PageRank: 0.265940267242289
Node: 4 PageRank: 0.25952966275208783
PageRank for each node: alpha = 0 , beta = 1
Node: 1 PageRank: 0.25
Node: 2 PageRank: 0.25
Node: 3 PageRank: 0.25
Node: 4 PageRank: 0.25
```

9. You have been tasked to design a classifier that decides whether students will be admitted to a CS graduate program. Applications to the program are received from students all around the world. Application contain student name, address, mobile phone number, final grade point average in undergraduate program and transcript. Describe what information from the application you would use as input to your classifier. For each piece of information, what (if anything) would need to be done to clean or transform the information into input data. For each of the transformed data input, identify the type (nominal, ordinal, interval or ration). Give at least 3 other pieces of information that would be helpful and describe why you think they would help.

Inputs:

1. Undergraduate GPA: Interval data, no transformation needed
2. Undergraduate major: Nominal data, Encode majors into numeric categories
3. Undergraduate university ranking: Ordinal data, Assign numeric ranks
4. GRE quantitative score: Interval data, no transformation needed
5. GRE verbal score: Interval data, no transformation needed
6. Number of CS courses taken: Ratio data, no transformation needed

The other fields like student name, address, and phone number need not be considered while deciding.

Additional useful inputs:

7. Letters of recommendation: Text content, extract key phrases about student qualities
8. Personal statement: Text content, analyze writing quality and stated interests
9. Years of work experience: Numeric, relevant industry experience indicates preparedness

The GPA, GRE scores, and CS coursework give quantitative measures of academic ability and preparation that would likely correlate with success in a graduate CS program.

The additional qualitative inputs help gauge interests, communication skills, recommendations, and experience - other helpful factors for evaluating applicants.

The inputs would be fed into a classification algorithm like logistic regression, SVM, or neural network to predict admissions chance.

10. You are given the following set of data

Name	City	Likes Beyoncé	In a relationship	Age	Number of concerts per year	Bought ticket to see Taylor Swift
Kate	Chicago	Yes	No	23	8	Yes
Joe	New York	No	No	36	4	Yes
Mena	New York	Yes	Yes	43	20	No
Pat	Chicago	No	Yes	19	2	No
Tim	Chicago	Yes	No	20	14	Yes
Tina	Chicago	Yes	Yes	54	7	Yes

Using entropy as a measure of purity, design a decision tree to predict whether someone bought a ticket to see Taylor Swift or not. Show how each decision node was selected

10) For splitting based on City:-
 Entropy of parent:-

\bullet - Bought Taylor swift tickets
 \times - Did not buy Taylor swift tickets

Chicago NY

$$E(\text{Parent}) = -\frac{4}{6} \log_2\left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right)$$

$$= 0.9183$$

$$E(\text{Chicago}) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right)$$

$$= 0.81125$$

$EC(NY) = -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)$
 $= 1$
 $EC(city) = \frac{4}{6} \times 0.81125 + \frac{2}{6} \times 1$
 $= 0.8741$
 $IG_{city} = 0.9183 - 0.8741$
 $= 0.0442$

For splitting based on liking Beyoncé:-

$EC(Yes) = 0.81125$
 $EC(No) = 1$
 $EC(Beyoncé) = 0.8741$
 $IG_{Beyoncé} = 0.0442$

For splitting based on relationship status:-

$EC(Yes) = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right)$
 $= 0.3899 + 0.5283$
 $= 0.9183$
 $EC(No) = -\frac{3}{3} \log_2\left(\frac{3}{3}\right) = 0$
 $\therefore IG_{relationship} = \frac{1}{6} \times 0.9183 - \left(\frac{3}{6} \times 0.9183 + \frac{3}{6} \times 0\right)$
 $IG_{relationship} = 0.45913 \rightarrow \text{Maximum}$

for splitting based on age:-

Since it is a continuous numerical value the avg. of pairs of values was taken
 19 \rightarrow 19.5
 20 \rightarrow 21.5
 23 \rightarrow 29.5
 36 \rightarrow 39.5
 43 \rightarrow 48.5
 54 \rightarrow 48.5
 ↑ In sorted order

splitting based on avg. age range

For ≤ 19.5

$EC(Y) = -1 \log_2(1) = 0$
 $EC(N) = -\frac{1}{5} \log_2\left(\frac{1}{5}\right) - \frac{4}{5} \log_2\left(\frac{4}{5}\right)$
 $= 0.4642 + 0.2535$
 $= 0.72172$
 $E(\leq 19.5) = 0 + \frac{5}{6} \times 0.72172 = 0.60143$
 $IG_{\leq 19.5} = 0.31687$

For ≤ 21.5

$EC(Y) = 0.81125$
 $EC(N) = 0.81125$
 $EC(\leq 21.5) = 0.8741$
 $IG_{\leq 21.5} = 0.0442$

For ≤ 29.5

$EC(Y) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right)$
 $= 0.9183$
 $EC(N) = 0.9183$
 $EC(\leq 29.5) = \frac{3}{6} (0.9183) + \frac{3}{6} (0.9183)$
 $= 0.9183$
 $IG_{\leq 29.5} = 0$

For ≤ 39.5

$EC(Y) = 1$
 $EC(N) = 0.81125$
 $EC(\leq 39.5) = 0.8741$
 $IG_{\leq 39.5} = 0.0442$

For ≤ 48.5

$EC(Y) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$
 $= 0.5283 + 0.4421$
 $= 0.9704$
 $EC(N) = 0$
 $EC(\leq 48.5) = 0.80905$
 $IG_{\leq 48.5} = 0.1092$

For splitting based on No. of concert:-
 \therefore it is a continuous, we take avg. values.

2 \rightarrow 3
 4 \rightarrow 5.5
 7 \rightarrow 7.5
 8 \rightarrow 11
 14 \rightarrow 17
 20

For ≤ 3

$EC(Y) = 0$
 $EC(N) = 0.72172$
 $EC(\leq 3) = 0.60143$
 $IG_{\leq 3} = 0.31687$

For ≤ 5.5

$EC(Y) = 0.81125$
 $EC(N) = 0.81125$
 $EC(\leq 5.5) = 0.8741$
 $IG_{\leq 5.5} = 0.0442$

For ≤ 7.5

$EC(Y) = 0.9183$
 $EC(N) = 0.9183$
 $EC(\leq 7.5) = 0.9183$
 $IG_{\leq 7.5} = 0$

For ≤ 11

$IG_{\leq 11} = 0.0442$

For ≤ 17

$IG_{\leq 17} = 0.31687$

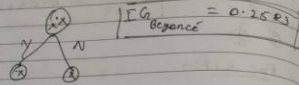
Out of all the Information gain values, we get the max for In a relationship.

Based on this, we consider:-
 Mena
 Pat
 Tina

For 2nd iteration:-

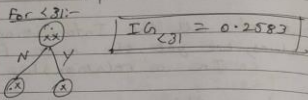
$EC(\text{parent}) = -1 \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right)$
 $= 0.81125$
 $EC(\text{city}) = -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)$
 $= 1$
 $EC(NY) = 0$
 $EC(\text{city}) = 0.5 \times 0.81125 + 0.5 \times 1$
 $= 0.655625$
 $IG_{\text{city}} = 0.252675$

For liking Beyonce :-

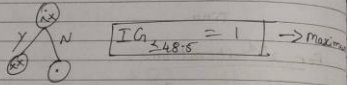


for age (taking avg) :-

19 → 31
43 → 48.5
54

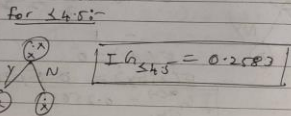


For < 48.5 :-

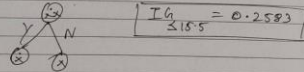


for No. of concerts :- (taking avg)

2 → 4.5
7 → 15.5
20



For < 15.5 :-



The 2nd split is done using Age

Final Decision tree :-

