

Name	PRAVEENAA KULANDHAIVEL
ASU ID Number	1225814683

CSE 472: Social Media Mining

Homework II - Network Models and Data Mining

Prof. Huan Liu
 Due at 2022 Sept 21, 11:59 PM

This is an **individual** homework assignment. Please submit a digital copy of this homework to **GradeScope**. For your solutions, even when not explicitly asked you are supposed to concisely justify your answers.

1. [Network Models]

- a. Assuming that we are interested in a dense random graph, what should we choose as the value of p? How does the value of p affect the sparseness? Where p defines the probability of forming edges.

Expected number of edges connected to a node, $c \geq (n-1)p$
 higher value of p \Rightarrow more edges \Rightarrow dense graph
 highest value of p = 1, there will be $(n-1)$ edges connected.

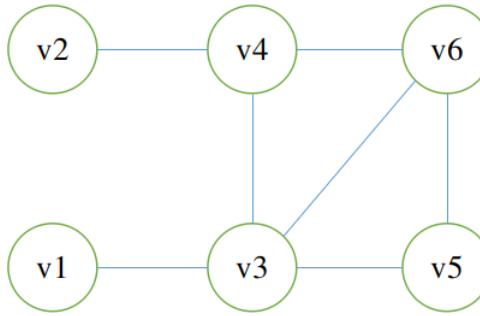
- b. We can make a simple random graph model of a network with clustering or transitivity as follows. We take n vertices and go through each distinct trio of three vertices, of which there are $\binom{n}{3}$, and with independent probability p we connect the members of the trio together using three edges to form a triangle, where $p = \frac{c}{\binom{n-1}{2}}$ with c constant. Show that the mean degree of a vertex in this network is $2c$.

Total number of pairs forming triangles - $\binom{n-1}{2}$

Probability for each triangle $\Rightarrow \frac{c}{\binom{n-1}{2}}$

\Rightarrow Therefore, each triangle contributes 2 edges for degree making the average $2c$.

- c. In a citation network, each node represents a paper, and an edge exists between two nodes if one paper cites another. When a new paper is published (i.e., a new node added to the graph), that paper cites previous papers with probability *proportional to their degrees*, hence, following the **preferential attachment model**. Suppose at each time-step a new paper is submitted to the following citation network. At $t=0$, we have 6 papers of $V_0 = \{v_1, v_2, v_3, v_4, v_5, v_6\}$, please calculate the probability of the new node connecting to each node. At $t=1$, v_7 joins and attaches to $n (= 4)$ papers with the highest probabilities, please update each node's probability; and do the same for $t=2$, when v_8 joins. Leave your answers as common fractions (e.g., $\frac{1}{6}, \frac{4}{5}$).



t	x	v_1	v_2	v_3	v_4	v_5	v_6	v_7	v_8
$t=0$	$P(v_x)$	$\frac{1}{14}$	$\frac{1}{14}$	$\frac{4}{14}$	$\frac{3}{14}$	$\frac{2}{14}$	$\frac{3}{14}$	NA	NA
$t=1$	$P(v_x)$	$\frac{1}{22}$	$\frac{1}{22}$	$\frac{5}{22}$	$\frac{4}{22}$	$\frac{3}{22}$	$\frac{4}{22}$	$\frac{4}{22}$	NA
$t=2$	$P(v_x)$	$\frac{1}{30}$	$\frac{1}{30}$	$\frac{6}{30}$	$\frac{5}{30}$	$\frac{3}{30}$	$\frac{5}{30}$	$\frac{5}{30}$	$\frac{4}{30}$

Note: NA in the table represents the situation when the new node cannot connect to the other nodes in the network.

Algorithm 1: Preferential Attachment Algorithm

Input : Graph $G(V_0, E_0)$, where $|V_0| = n_0$ and $d_v \geq 1 \forall v \in V_0$, number of expected connections $n \leq n_0$, time to run the algorithm t

Output: A scale-free network

```

1 //Initial graph with  $n_0$  nodes with degrees at least 1
2  $G(V, E) = G(V_0, E_0)$ 
3 for 1 to  $t$ 
4    $V = V \cup \{v_i\}$       //add new node  $v_i$ 
5   while  $d_i \neq n$ 
6     Connect  $v_i$  to a random node  $v_j \in V$ ,  $i \neq j$  (i.e.,  $E = E \cup \{e(v_i, v_j)\}$ ) with probability  $P(v_j) = \frac{d_j}{\sum_k d_k}$ 
7 return  $G(V, E)$ 

```

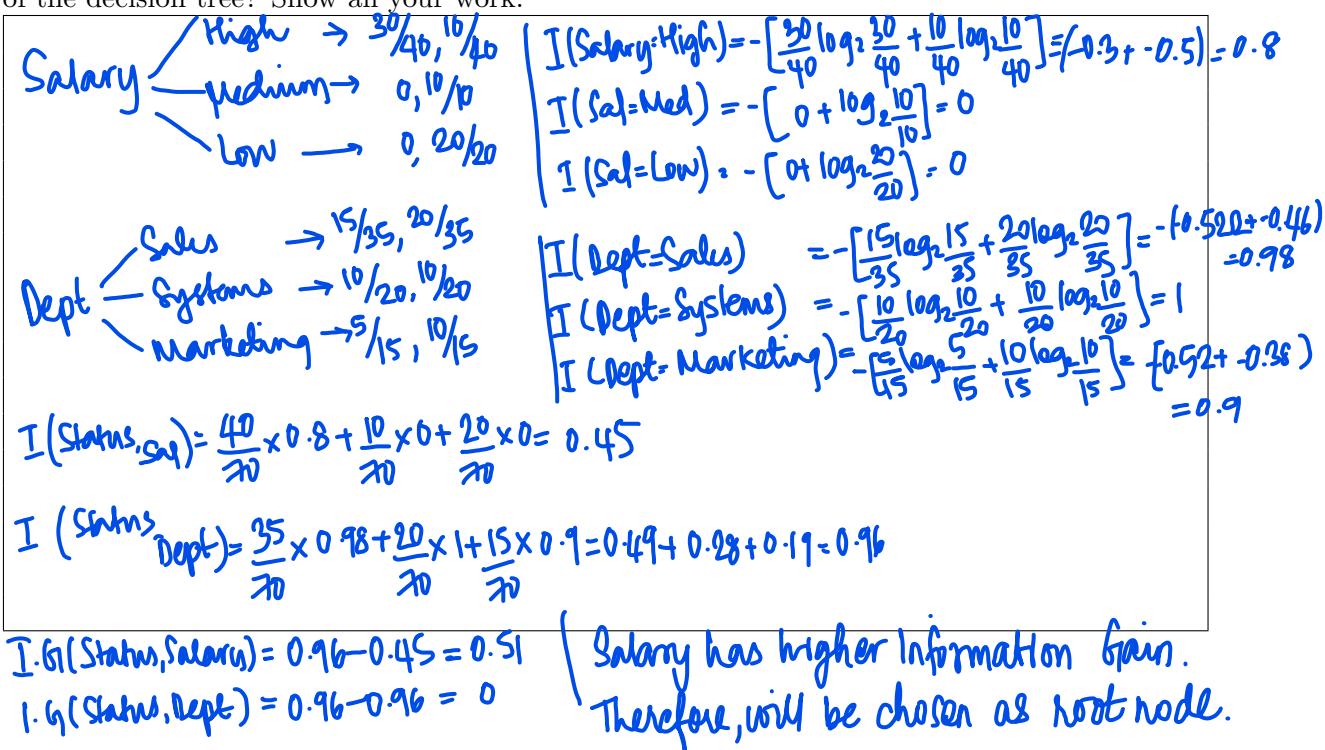
2. [Data Mining] Consider the given dataset from an employee database. For a given entry row, *count* column represents the number of data tuples having the values for *department*, *salary*, and *status* given in that row. For example, there are 15 instances with values of (*department* = sales, *salary* = high, *status*=senior). Let *status* be the class label attribute, answer the following questions. Please use 2 as the default value for base of all logs.

Department	Salary	Status	Count
Sales	High	Senior	15
Sales	Low	Junior	20
Systems	Medium	Junior	10
Systems	High	Senior	10
Marketing	High	Junior	10
Marketing	High	Senior	5

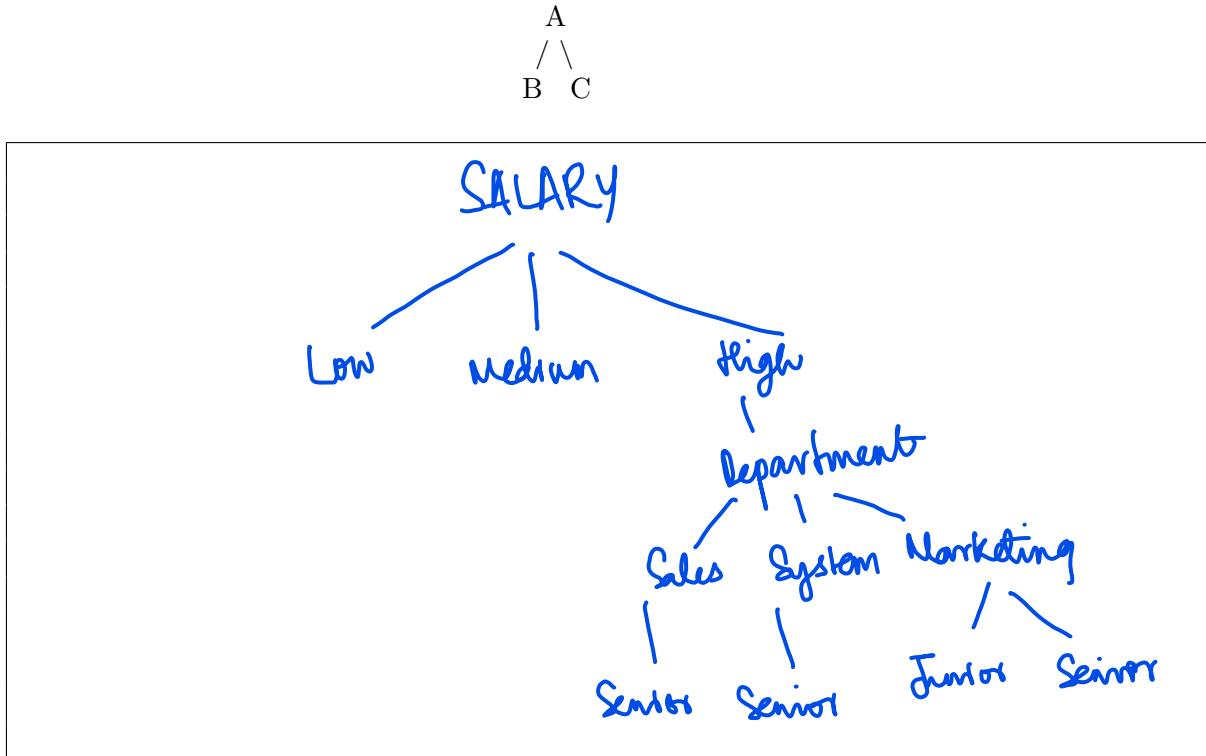
- a. What is the value for the $H(\text{Status})$? Where $H(x)$ defines the entropy of x .

$$H(\text{Status}) = -\frac{30}{70} \log_2 \frac{30}{70} - \frac{40}{70} \log_2 \frac{40}{70} \Rightarrow -0.42 \times 1.2 - 0.57 \times 0.8 \\ = 0.504 + 0.456 \\ = 0.96$$

- b. Based on the Information Gain values, which feature is the most probable to be the root node of the decision tree? Show all your work.



- c. Draw the final decision tree. An example of how to draw the tree on the text box:



- d. Given a data instance having the values “Sales” and “High” for the attributes *department* and *salary*, respectively, what would a Naive Bayesian classification of the class attribute *Status* for the instance be? Detail all your calculations.

$$\begin{aligned}
 P(\text{Junior} | \text{Sales} \& \text{High}) &= \frac{P(\text{Sales}, \text{High} | \text{Junior}) P(\text{Junior})}{P(\text{Sales}, \text{High})} \\
 &= \frac{(20/40) \times (10/40) \times (10/40)}{P(\text{Sales}, \text{High})} = 0.07 \times \left(\frac{1}{P(\text{Sales}, \text{High})} \right) \\
 P(\text{Senior} | \text{Sales} \& \text{High}) &= \frac{P(\text{Sales}, \text{High} | \text{Senior}) P(\text{Senior})}{P(\text{Sales}, \text{High})} \\
 &= \frac{(15/30) \times (30/30) \times (30/40)}{P(\text{Sales}, \text{High})} = 0.21 \times \left(\frac{1}{P(\text{Sales}, \text{High})} \right)
 \end{aligned}$$

Therefore, Value for Senior > Value of Junior. Status would be Senior