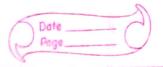
Name: UID: class:	Vivek Hemantsi 2019110046 B.E ETRX ((cc)	Date		
			China China	3 4 0 - 1947 1	-2.6	
	Assignment - 1 - 12-11-2022					
	Air Traffic		h ii- 1 .	14211		
	Letus considu		vervation	recorded in	na	
	database	(3) (4) PM	1100	7		
	· A =	Day, se	,			
				with 20.		
	Samo Frach	= [On Time	Late: V	ery Later Co	ancelled	
	To find most blooky statistics for all or					
	To find most likely classification for any other unseen instance:					
	Week Day, Winter, High, None, 3?					
11\	ell = [count Chale) × dount (Bathar)]					
	· Conditional Probability: PCY/21-,20)					
	P(x1/y).P(x2/y)					
	PCAIB) = PCAOB) _ PCBOA) KANY)					
	p(H1). P(Xe)-P(Mn)					
	= PCBIA)·PCA)					
	015 - 031 / [PCB) + 008] 7 Els 10					
			Class			
	Attribute	om+ Luz	o Late !	Very	'Concelled	
		Time		Late		
322	Weekday	9/14 = 0.64	1/2=0.5	3/3=1	0/1=0	
		2/14=0.14	1/2=0.5	0/3=0	111 = 1	
200	. /	1114 = 0.07	0/2=0	013=0	011 = 0	
	Holiday	2/14 = 0.14	012=0	013=0	0/1 = 0	
	11007249	•				
	Spring	4114 = 0.29	0/2=0	0/3=0	0/1=0	
1 8		6/14 = 0.43	0/2=0	0/3=0	0/1 = 0	
	Autumn	2/14 = 0.14	0/2=0	\$13 = 0·33		
V	Winter	2/14=0.14	2/2=1	2/3 = 0.67	0/1=0	

	6
Date	
Page	()
0	

						1
	None	5/14= 0.36	0/2=0	0/3 = 0	011=0	
Fog	High	4114=0.29	1/2 = 0.5	113 = 0.33	1/1=1	-
U	Normal	5/14 = 0.36	1/2=0.5	2/3=,0.67	0/1=0	1
	None	5/14 = 0.36	1/2=0.5	1/3 = 0.33	0/1=0	
air	Slight	8/14 = 0.57	012=0	0/3=0	0/1=0	
N	Heavy	1/14=0:07	1/2=0.5	2/3=0.67	1/1=1	
,	Prior	14/2-0-70	2/20= 0.10	3/20=0.15	1120=0.05	
	Probability	120		الما المعتبد		_
	Weekday	- Winter-	High-Non	e - 333		
Case 1: Class = On-Time: Prior Probability (On-time)						
				ekday Con-		
		x (Winter (On-Time) - season				
	* (Foo High (on-time) - Fog					
	x None (On-Time) - Rain					
	$= 0.70 \times 0.64 \times 0.14 \times 0.29 \times 0.36$ = 0.0065 Prior Probability (Lake)					
	= 0.0065 Prior Probability Clase) Case 2: Class = Late: 0.10 x 0.5 x 1 x 0.5 x 0.5					
	= 0.0125					
	Prior Probability (Very Late)					
	Case 3: (Class = Ve		0.15×1×0		0
			0109			
			pric	or probability ((Cancelled)	
	Case 4: class = Cancelled: 0.05x0x0x1x0					
	= 0.000					
	Highest Probability Occurs for the case 2 (Late).					
	Therefore, Case 2 is the strongest;					
ώ.	Hence correct classification is late					
	So, When the day is Weekday, Season is Winter, Fog is High, Rain is None, Class mostly					
	Fog i. liketo	s High, Ró Late	unis Mon	e, Class	mostly	
			ACTIVITIES OF			
IF				Scanned with CamScanr	ner	

1. Expected Frequencies The expected eij frequency which can be computed as ejj = count (A=ai) x count (B=bi) where N= no. of data tuples. count (A=ai) = no of tuples having value ai for A count (B=bj) = no. of tuples having value Now, ell = [count Cmale) x count (fiction)]/N ell = [300 x 450] / 1500 · : e12 = [1200 x 450] /1500 = 360 e12 = [300 × 1050] / 1500 = 210 : e22 = [1200 × 1050] /1500 = 360

	-12 m	000	ing is it to	Hobelett.				
1 1 2 11	X-= 5	& (0ij	- جي ا	is augoris				
	$\chi^2 = \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} \frac{Coij - e_{jj}}{e_{ij}}$							
	Oij = Observed Prequency							
	eij = Expected frequency							
			•					
	men are r	no. of nows	eno. of col	lums				
Co.	- 5500 F - 1015 -	05) + -(0	2 - 02 2). =	J	Н			
		104	Q 23.	,				
	preferred read	Male	Female	Total				
i i	Hickion 1111	250 (96)		450				
į.	Non-Fiction	50 (216)	1000 C840)	1050				
i.	Total	300	1200	1.1509				
	Ho: Preferred Reading & gender are independent							
	of each o	ther.	(0)	-1				
Jar' C	Stone		of FIFT					
į i	Ha: Preferre	d Reading	(gender	are triot				
. 23	indonanda	nt each c	Har	toid - ada				
	Maepenae	Cach		1))				
	the dear of the day are (1) is 1)							
		(1) 1/5 = 5	7 70 1000		1			



Frequencies must equal the total observed frequency

$$\frac{1}{30} = \frac{(250-90)^2 + (50-210)^2 + (200-360)^2}{4(200-360)^2}$$

= 284.44 + 121.90 + 71.11 + 30.48

For this given tables allows bounded ist

-					
		Male	Female	Sum (now)	
	fiction	250 C90)	200 (360)	450	
	Non-Action	50 (210)	1000 C840	1050	
	sum ((01.)	300	1200	1500	

the degree of freedom are (2-1) (2-1) = 1,

For 1 degree of freedom.

value 12 with degree of Freedom 1 and 0.01 significance level from the standard statistical table is 6635 or 10.828 Ctaken from the table of upper 1. points of 1/2)

Our receive value is above this value. Therefore we can reject the hypothesis that gender & preferred reading are independent and

correlated for given group of people.