Notes Syllabus Question Papers Results and Many more... Available @

www.binils.com

tion at	Reg. N	o.:		
	Ougstion Par	non Codo . 70	0079	
	Question Pa	per Code : 10	0014	
B.E./	B.Tech. DEGREE EXAMINA	ATIONS, NOVEMB	ER/DECEMBER 2022	
	This	rd Semester		
	Computer Sci	ence and Engineerir	ng	
	CS 3352 – FOUNDA	TIONS OF DATA S	CIENCE	
(Common	to: Computer and Communi	cation Engineering	/ Information Technol	ogy)
	(Regu	lations 2021)		
Time: Thr			Maximum: 100 m	arks
		ALL questions. $(10 \times 2 = 20 \text{ marks})$		
2. List	an overview of common er tions to be employed.		data and which clean	sing
(c) fa (g) ne	sify the below list of data family size (d) academic et worth (dollars) (h) third-p ef note on them.	major (e) sexual	preference (f) IQ s	core
4. Diffe:	erentiate discrete and contin	uous variables.		
5. What	t is a percentile rank? Give a	an example.		
	sider Helen sent 10 greeting rds, what is the kind of relati			back
7. List t	the attributes of a Numpy ar	ray. Give an examp	le for it.	
A-40,	ate a data frame with key ar 0, C-5, B-10, C-10. Find the s group.			
9. What	it is the purpose of errorbar f	unction in Matplotli	b? Give an example.	
10. Show Code.	wcase 3-dimensional drawing.	ig in Matplotlib wi	th corresponding Py	thon

Notes Syllabus Question Papers Results and Many more... Available @

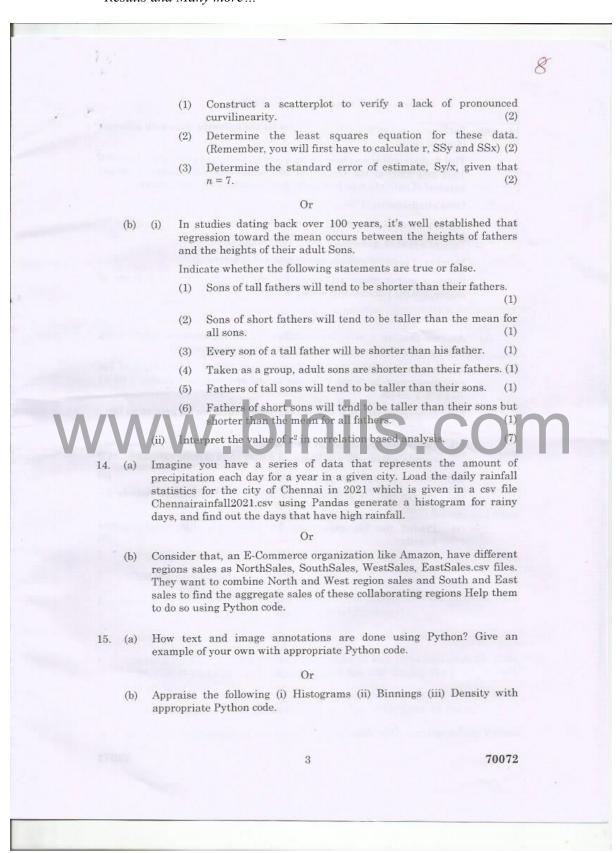
www.binils.com

		DADT D (5 12 - C7 1)
		PART B — $(5 \times 13 = 65 \text{ marks})$
.11.	(a)	Examine the different facets of data with the challenges in their processing.
		Or Or
	(b)	Explore the various steps associated with data science process and explain any three steps of it with suitable diagrams and example.
12.	(a)	Demonstrate the different types of variables used in data analysis with an example for each.
		Or
	(b)	The number of friends reported by Facebook users is summarized in the following frequency distribution.
		FRIENDS f
		400 – above 2
		350 – 399 5
		300 – 349 12
		250 – 299 17
		200 - 249 23 $150 - 199$ 49
		100 = 149 27
		50 - 99 29
		0-49 36
VV	VA	VV Total 200
VV	VAV	(i) What is the shape of this distribution?
VV		(i) What is the shape of this distribution? (ii) Find the relative frequencies.
VV		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300–349.
VV		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300–349. (iv) Convert to a histogram.
VV	di p	(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300–349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display?
13.	(a)	(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300-349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7)
13.		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300-349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter
13.		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300-349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my
		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300-349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my neighborhood: Drivers (X) Cars (Y) 5 4
		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300-349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my neighborhood: Drivers (X) Cars (Y) 5 4 5 3
		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300-349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my neighborhood: Drivers (X) Cars (Y) 5 4 5 3 2 2
		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300-349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my neighborhood: Drivers (X) Cars (Y) 5 4 5 3 2 2 2 2
		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300-349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my neighborhood: Drivers (X) Cars (Y) 5 4 5 3 2 2 2 2 3 2 3 2
		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300-349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my neighborhood: Drivers (X) Cars (Y) 5 4 5 3 2 2 2 3 2 1 1 1
		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300–349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my neighborhood: Drivers (X) Cars (Y) 5 4 5 3 2 2 2 2 3 2 1 1 1 2 2
		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300–349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my neighborhood: Drivers (X) Cars (Y) 5 4 5 3 2 2 2 2 3 2 1 1 1 2 2
		(i) What is the shape of this distribution? (ii) Find the relative frequencies. (iii) Find the approximate percentile rank of the interval 300–349. (iv) Convert to a histogram. (v) Why would it not be possible to convert to a stem and leaf display? (i) Categorize the different types of relationships using Scatter plots. (7) (ii) Each of the following pairs represents the number of licensed drivers (X) and the number of cars (Y) for seven houses in my neighborhood: Drivers (X) Cars (Y) 5 4 5 3 2 2 2 2 3 2 1 1 1 2 2

Notes Syllabus Question Papers Results and Many more...

Available @

www.binils.com



Notes Syllabus Question Papers Results and Many more...

www.binils.com

Available @

PART C — $(1 \times 15 = 15 \text{ marks})$

16. (a) Perform an exploratory data analysis for the following data with different types of plots:

The dataset contains cases from a study that was conducted between 1958 and 1970 at the University of Chicago's Billings Hospital on the survival of patients who had undergone surgery for breast cancer.

Data attributes:-

Age of patient at the time of operation (numerical)

Patient's year of operation (year - 1900, numerical)

Number of positive axillary nodes detected (numerical)

Survival status (class attribute) 1 = the patient survived 5 years or longer, 2 = the patient died within 5 year

Or

(b) Assume that an r of - .80 describes the strong negative relationship between years of heavy smoking (X) and life expectancy (Y).

Assume, furthermore, that the distributions of heavy smoking and life expectancy each have the following means and sums of squares: 5 60 35 70~x~y~X~Y~SS~SS

(i) Determine the least squares regress on equation for predicting life expectancy from years of heavy smoking (3

Determine the standard error of estimate, Sy/x, assuming that the correlation of 80 was based on n = 50 pairs of observations. (3)

(iii) Supply a rough interpretation of Sy/x.

- (3)
- (iv) Predict the life expectancy for John, who has smoked heavily for 8 years. (3)
- (v) Predict the life expectancy for Katie, who has never smoked heavily.
 (3)

70072