



Course:	Introduction to Data Science	CourseCode:	DS 2001
Program:	BS(DS)	Semester:	Fall 2023
Duration:	1 Hour	Total Marks:	50
Paper Date:	02-10-2023	Page(s):	6
Section:	BS (DS) A, B, C	Section:	RNC 2A
Exam:	Mid I	Roll No:	

**Instructions:** Answer in the space provided. You can ask for rough sheets, but they will not be graded or marked. In case of confusion or ambiguity make a reasonable assumption. Questions during exam are not allowed.

Question#1:

10x4 = 40 Marks

The dataset represents a sample of employee performance evaluation data, containing various attributes related to individual employees within an organization. It includes information such as employee IDs, department affiliations, ages, genders, years of experience, performance ratings, joining dates, and salaries. Each row corresponds to a unique employee, and the dataset provides insights into factors affecting employee performance and compensation.

Employee_ID	Department	Age	Gender	Experience (Years)	Rating (1-5)	Joining Date	Salary
E001	Sales	35	Male	8	4	2020-06-15	60000
E002	HR	28	Female	4	3	2021-01-20	55000
E003	Engineering	42	Male	15	5	2019-03-10	75000
E004	Marketing	31	NULL	0	4	2020-11-05	32000
E005	Sales	29	Male	7	NULL	2021-09-18	58000
E006	Engineering	36	Male	10	4	2020-04-25	70000
E006	Sales	36	Male	8	4	20-04-2020	70000

Answer the following questions:

a) What is the type of each feature?

Employee-ID = string  
 Department = string  
 Age = integer  
 Experience = integer  
 Rating = integer  
 Joining Date = string

Salary = integer

b) Identify at least three quality issues with this data.

- i) Null values: there are two null/missing values in ~~gender~~ "gender" column, row 4 and "Rating" column, row 5 and negative value in col 5 row 6
- ii) Index error: In last two rows, employee ID has same value "E006"
- iii) Wrong Date format: in joining date column last row "20-04-2020" has wrong format or rest of the data

c) Is there a correlation between years of experience and salary? If so, what is the nature of this correlation?

Correlation is that they have direct relationship between them. are in wrong format

As experience increases from 0 to 4, salary increases highly. But from 4 years onwards, salary increases at lesser rate comparatively.

d) Can you figure out imbalance distribution in any of the features?

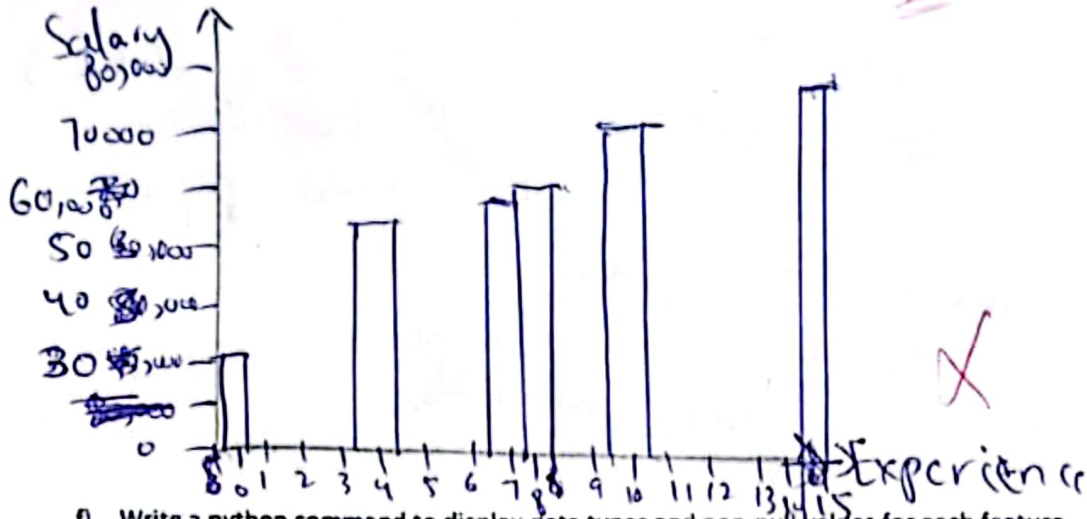
Age: There is '311' value which is an outlier and impossible outcome

Experience: "-8" in last row, experience is always positive.

Salary: There is gap between values



e) Create a histogram of salaries. Identify the type of distribution.



f) Write a python command to display data types and non-null values for each feature.

`df.info()`

g) Write python code to group the data by "Gender" and calculate the average age for each gender.

`ge - gender = df.groupby("Gender")["Age"].n`

h) Write python code to Determine the number of male and female employees in the dataset. 4

~~male\_no = df.groupby("Gender").size.unstack()~~  
~~male\_no = df.groupby("Gender").nunique()~~  
~~male\_no = df.groupby("Gender")~~  
 male\_no = 0  
 female\_no = 0  
 for i in range(df.shape[0]):  
 if (df["Gender"].iloc[i] == "Male"):  
 male\_no = male\_no + 1  
 elif (df["Gender"].iloc[i] == "Female"):  
 female\_no = female\_no + 1

i) Write a python code to Group the data by "Department" and calculate the average salary for each department. 4

avg salary = df.groupby("Department")["Salary"].  
 mean()

j) Write a python code to calculate the mean, median, and standard deviation of the "Salary" column. 3

smean = df["Salary"].mean()  
 smedian = df["Salary"].median()  
 listsmode = df["Salary"].mode()  
 smode = listsmode[0]

a) What are the key challenges in data cleaning, and how do you address them?

The key challenges in data cleansing are:-

- Filling missing values.
- Dropping desired data which is not useful.
- Replacing the ~~data~~ missing values with mean, a median etc.
- Detecting outliers.
- Removing or replacing outliers.

b) Why is it important to identify outliers in a dataset, and what methods can be used for outlier detection?

→ It is important to identify outliers in a dataset.

because data having outliers in it will not give accurate results. Data will produce unusual results.

→ Outlier detection method is following:-

$$Q_1 = 0.25 \times N \quad \therefore \text{where } n \text{ is sample number.}$$

$$Q_3 = 0.75 \times N$$

$$\rightarrow \text{upper bound} = Q_3 + 1.5 \text{ IQR} \quad \therefore \text{IQR} = Q_3 - Q_1$$

$$\rightarrow \text{lower bound} = Q_1 - 1.5 \text{ IQR.}$$

If entries of dataset, are greater than upper bound & lower than lower bound, then the dataset

FAST School of Computing will have outliers in it.