

NAME \Rightarrow Kamal deep

Roll NO \Rightarrow 2501660006

Course name \Rightarrow BCA Cyber security

Topic \Rightarrow Fundamentals of Data
collection and cleaning

Signature \Rightarrow

Assignment \Rightarrow 2

- 1) Learning Objectives
- 2) Identify different types and sources of data.
- 3) Data can be categorized by its sources (Internal, External, Primary or secondary) and by its type (quantitative or qualitative). Internal sources are from within an organization, while external sources are from outside.

Types of Data :-

- (i) Quantitative Data : Numerical data can be measured and calculated, such as sales figures or employee headcount
- (ii) Qualitative Data : Non-numerical data that describes qualities or characteristics, such as customer feedback, observation, or interview transcripts.

Sources of Data

- (i) Internal sources : Data that originates and is collected from within an organization, like sales records, company reports or data from public database.

Primary data collection methods like surveys, interviews and observations.

Service interviews and observations are qualitative data collection methods used to gather in-depth insights. Interviews involve direct conversation to collect verbal responses, while observations involve ~~and~~ watching and decoding behaviour in a natural setting. Both are used to understand experiences, perceptions, and behaviours, often in the context of social science research or service quality analysis.

Service interviews

- What it is: A conversation between a researcher and one or more participants to gather detailed information through oral questions and responses.
- How it works: Researchers ask open-ended questions to understand a person's beliefs, attitudes and experiences.
ex ⇒ Asking customers about their experience with a new service by having them describe their feelings and the steps they took to use it.

3) Perform basic data cleaning including handling missing values and duplicates

3) Basic data cleaning involves identifying and correcting errors, specifically by handling missing values (through deletion or imputation) and removing duplicates (by deleting redundant entries). A systematic process first assesses data quality, then removes duplicates and irrelevant data, fixes structural errors, handles missing data, and finally validates the cleaned data.

use data integrity and ethical handling during processing. Implement technical and organizational measures such as access controls, encryption, validation, and regular audits, while also adhering to ethical principles like transparency, privacy, and minimizing data collection.

Assignments tasks

Task 1: Design a short survey (5-10 questions) to collect feedback about student satisfaction at your college.

Your feedback is invaluable in helping us improve our certification programs. This survey will take approximately 3-5 minutes to complete. All responses are confidential.

Section 1: General Program Experience.

1) Which certification program did you complete or are you currently enrolled in?

- [list Program A]
- [list Program B]
- [list Program C]

Other (Please specify)

2) How satisfied were you with the overall quality of the certification program?

- very satisfied
- satisfied
- neutral
- dissatisfied
- very dissatisfied

5) Effectively do you feel ~ ~ ~ - field?

- extremely effective
- very effective
- moderately effective
- slightly effective

4) How would you rate the instructor's knowledge of the subject matter?

- Excellent
- Good
- Average
- Poor
- very Poor

5) How satisfied are you with the accessibility and support provided by faculty and staff?

- very satisfied
- satisfied
- Neutral
- Dissatisfied
- very Dissatisfied

6) Where the online resources, library materials, ~ ~ ~ ?

- Yes completely
- Mostly
- Neutral
- Not really
- No, not at all

Wk2: collect 10-15 responses (real or simulated) and store them in a spreadsheet.

You can either use an online form (like google forms or microsoft forms) to collect real responses automatically or manually enter simulated data into an existing spreadsheet file.

method 1: using online forms

(Recommended for real responses)

Steps using google forms:

1. Create the form
2. Link to a spreadsheet.
3. Share the form
4. Collect Responses
5. Review the spreadsheet

Using microsoft forms:

1. Create the Forms
2. Link to Excel
3. Share the Form.

method 2: manually entering simulated data

- 1) Open a spreadsheet
- 2) Setup Headers
- 3) Enter simulated Responses
- 4) Save the file

Task 3: Identify and Fix issues such as missing data, duplicates entries, or inconsistent formatting.

A systematic process of data cleaning (or data cleansing) is required. This involves parsing that data to detect errors and then applying various techniques to correct them, often using tools like spreadsheets, SQL, or programming libraries.

1: Identifying the issues (Data Profiling)
missing Data

Duplicate Entries

- Inconsistent Formatting

Duplicate entries

- Remove duplicates
- consolidate records

Step 2. Fix the Issues

Tools and Best Practices

missing data

- Remove Records
- Impute values
- Predictive Modeling

- spread sheets
- programming language
- specialized software

Task 4: write a short note on ethical considerations while collecting personal data.

Ethical data collection requires prioritizing informed consent, privacy, and transparency. Key considerations include clearly explaining how data will be used, securing the data against breaches, only collecting necessary information and using it for the purpose it was originally intended, while also being accountable for any misuse.

Informed consent.
voluntary and clear
no coercion.
Right to withdraw

Transparency and accountability
• openness
• Accountability

Privacy and security
Protect data
minimize risk
Confidentiality

Responsible less

- Purpose limitation
- Avoid bias.