

# Introduction to Data Science

## Assignment 1

Name: Anuska Ghosh

Course: BCA (Sec – A)

Registration Number: 2411021240012

GitHub Repository Link: [https://github.com/anuskaghosh17/IDS\\_Python](https://github.com/anuskaghosh17/IDS_Python)

### Part 1: Theoretical Understanding

#### 1. Define Data Science:

Q. What is Data Science? Discuss its key components and the CRISP-DM process.

Ans. Data Science is an interdisciplinary field that uses mathematics, statistics, computer science, and domain knowledge to extract meaningful insights from data.

It is all about asking the right questions and using data to find answers.

#### Key Components:

- a. Data Collection: Gathering raw data from various sources.
- b. Data Processing: Cleaning and organizing the data.
- c. Data Analysis: Using statistical and computational methods to derive insights.
- d. Visualization: Presenting results in an understandable format (e.g., graphs, dashboards).
- e. Decision-Making: Applying insights to solve real-world problems.

#### The CRISP-DM Process:

A structured approach to solving data science problems.

1. Business Understanding: Define objectives and questions.
2. Data Understanding: Explore the data to understand its structure.
3. Data Preparation: Clean, transform, and organize the data.
4. Modeling: Apply algorithms to analyze and predict.
5. Evaluation: Check if the model answers the problem.
6. Deployment: Implement the solution in real-world scenarios.

Q. Explain how the CRISP-DM framework is applied in solving real-world problems (e.g., predicting customer churn or recommending movies).

Ans. **Predicting Customer Churn in Telecom:-**

Problem Statement:

How can a telecom company predict which customers are likely to stop using their services?

Dataset:

Source: Telco Customer Churn Dataset. [<https://www.kaggle.com/datasets/blastchar/telco-customer-churn>]

Columns:

- Customer ID
- Demographics: Gender, age group, etc.
- Service details: Internet service, contract type, monthly charges.
- Churn (Yes/No).

CRISP-DM Process:

1. *Business Understanding:*

- Goal: Reduce churn rate by identifying at-risk customers.
- Impact: Increase revenue by targeting retention efforts.

2. *Data Understanding:*

- Analyze churn rate across demographics.
- Understand correlations between service features and churn.

3. *Data Preparation:*

- Handle missing values (e.g., impute missing charges).
- Convert categorical variables (e.g., gender, contract type) to numerical data.
- Scale numerical features like monthly charges.

4. *Modeling:*

- Train a classification model (e.g., Logistic Regression, Random Forest) to predict churn.
- Use features like contract type, monthly charges, and tenure.

5. *Evaluation:*

- Use metrics like Accuracy, Precision, Recall, and F1 Score.
- Evaluate the model on a confusion matrix to understand false positives and negatives.

6. *Deployment:*

- Provide alerts for high-risk customers to customer service teams.
- Implement proactive offers and discounts to retain customers.

## **Netflix Recommendation System:-**

### Problem Statement:

How can Netflix recommend personalized movies or TV shows to users based on their preferences?

### Dataset:

Source: MovieLens Dataset (Free public dataset). [<https://grouplens.org/datasets/movielens/>]

### Columns:

- User ID
- Movie ID
- Rating (1-5)
- Timestamp
- Movie metadata (title, genres, release year).

### CRISP-DM Process:

#### *1. Business Understanding:*

- Goal: Improve user engagement by suggesting content they're likely to enjoy.
- Impact: Increased user satisfaction and retention.

#### *2. Data Understanding:*

- Explore the dataset: Number of users, movies, and ratings.
- Analyze distribution of ratings and popular genres.

#### *3. Data Preparation:*

- Handle missing data in movie metadata.
- Transform timestamp into human-readable format.
- One-hot encode genres for analysis.

#### *4. Modeling:*

- Use Collaborative Filtering to predict user preferences:
  - (I) Find similar users and suggest movies they liked.
  - (II) Train a recommendation model (e.g., Singular Value Decomposition).
- Alternative: Content-based filtering using movie metadata.

#### *5. Evaluation:*

- Split data into training and test sets.
- Use metrics like Root Mean Square Error (RMSE) to measure model accuracy.

#### *6. Deployment:*

- Integrate the model into a recommendation engine.
- Provide real-time recommendations on the Netflix platform.

## 2. Case Study Questions:

- From the case studies in the "Module 1 Case Studies" file, answer the following:

Q. What is the main business objective of the Netflix Recommendation System?

Ans. : The main business objective of the Netflix Recommendation System is to improve user engagement by suggesting content they're likely to enjoy. This will increase the user satisfaction and retention.

## Part 2: Data Manipulation and Joins

```
[10]: import pandas as pd
[11]: students=pd.read_csv('15_Student.csv')
[12]: students
[13]:
   StudentID  Name  Marks
0         101  Alice    85
1         102   Bob    90
2         103  Charlie  88
3         104  David    92
[14]: details=pd.read_csv('15_StudentDetails.csv')
[15]: details
[16]:
   StudentID  Age  Grade
0         101   20    A
1         102   21    B
2         103   22    A
3         104   19    C
[17]:
[18]: merged_students=pd.merge(students,details,on='StudentID')
[19]: merged_students
[20]:
   StudentID  Name  Marks  Age  Grade
0         101  Alice    85.0  20.0    A
1         102   Bob    90.0  21.0    B
2         103  Charlie  88.0  22.0    A
3         104  David    92.0  19.0    C
[21]: merged_students['Name']
[22]: merged_students['Age']
[23]: merged_students['Grade']
[24]: merged_students[['Name','Age','Grade']]
[25]: merged_students[['Name','Age','Grade']].reset_index(drop=True)
[26]: merged_students[['Name','Age','Grade']].reset_index(drop=True).to_csv('merged_students_details.csv')
```

StudentID	Name	Marks	Age	Grade
101	Alice	85.0	20.0	A
102	Bob	90.0	21.0	B
103	Charlie	88.0	22.0	A
104	David	92.0	19.0	C

Types of Joins and difference between each of them:

- Inner Join: Only take the matching records from both students and details files.
- Left Join: All records from the students file (left) and matching records from the details file (right).
- Right Join: All records from the details file (right) and matching records from the students file (left).
- Outer Join: Take all the records from both students and details files.

Merged CSV File: [http://localhost:8888/lab/tree/IDS/merged\\_students\\_details.csv](http://localhost:8888/lab/tree/IDS/merged_students_details.csv)

	StudentID	Name	Marks	Age	Grade
1	101	Alice	85.0	20.0	A
2	102	Bob	90.0	21.0	B
3	103	Charlie	88.0	22.0	A
4	104	David	92.0		
5	105			19.0	C

[illegible]

