

Chapter 14: Revisiting AI Project Cycle, Data Collection, Data Access

Introduction

Artificial Intelligence (AI) has revolutionized the way machines interact with data and make decisions. In earlier chapters, we were introduced to the AI Project Cycle—a structured approach to developing AI-based solutions. In this chapter, we **revisit the AI Project Cycle** with a focus on two crucial stages: **Data Collection** and **Data Access**. Understanding how to gather, handle, and use data is foundational to any AI project. Without quality data, AI models cannot learn or make intelligent predictions.

14.1 The AI Project Cycle – A Quick Recap

The **AI Project Cycle** includes the following stages:

1. **Problem Scoping** Identify and define the problem you want to solve.
2. **Data Acquisition / Collection** Gather relevant data required to train your AI model.
3. **Data Exploration** Understand the nature, patterns, and structure of the data.
4. **Modelling** Build and train an AI model using the data.
5. **Evaluation** Assess the performance of the model using metrics.

Note: In this chapter, our main focus is **Data Collection (Stage 2)** and **Data Access**—how data is sourced, types of data, and legal considerations.

14.2 Data Collection

What is Data Collection?

Data Collection is the process of gathering information from various sources to be used for training AI models. It is the **second and one of the most important stages** in the AI Project Cycle.

Why is Data Collection Important?

- AI models learn patterns from data.
- Better data = Better learning = More accurate predictions.
- Poor data can lead to **biased** or **inaccurate** models.

Types of Data:

Type	Description	Example
Structured Data	Well-organized in tables or databases	Excel files, CSVs
Unstructured Data	Not organized in pre-defined format	Images, videos, texts, audio
Semi-Structured	Partially organized	JSON files, XML documents

Sources of Data:

1. Primary Data

- Collected directly by the user or organization.
- Tools: Surveys, interviews, sensors, observations.

2. Secondary Data

- Collected by others and reused.
- Sources: Government portals, research websites, public datasets.

Data Collection Tools and Platforms:

- Google Forms
 - Microsoft Excel / Google Sheets
 - APIs (Application Programming Interfaces)
 - Mobile apps/sensors
 - Kaggle, UCI Machine Learning Repository
-

14.3 Data Access

Once data is collected, it needs to be **accessed, managed, and stored securely** for further use in model training.

Methods of Data Access:

Method	Description
Local Files	Stored on your device (e.g., .csv, .xlsx)
Cloud Storage	Data stored on cloud platforms (Google Drive, Dropbox)
Databases	Structured data stored in DBMS like MySQL, MongoDB
APIs	Data accessed programmatically from websites or services
Web Scraping	Automated extraction of data from websites (with permission)

⚠ Always ensure legal compliance and permissions when accessing data.

14.4 Legal and Ethical Considerations

AI projects deal with real-world data that can sometimes include personal or sensitive information. It's important to handle such data ethically.

Key Principles:

1. **Data Privacy** Do not share personal or sensitive data without consent.
2. **Data Ownership** Ensure you have the right to use the data.
3. **Bias and Fairness** Avoid using data that may be biased towards a particular group.
4. **Copyright Laws** Respect copyrights when using text, image, or other media data.

Legal Frameworks to Know:

- GDPR (General Data Protection Regulation – EU)
 - IT Act (India)
 - Data Protection Bill (India – upcoming regulation)
-

14.5 Quality of Data: Garbage In, Garbage Out

The performance of an AI model depends heavily on the **quality** of data. If bad data is used, the model will give inaccurate predictions.

Good Data Characteristics:

- Relevant
 - Accurate
 - Complete
 - Clean (free of errors or duplicates)
 - Diverse (to avoid bias)
-

14.6 Hands-On Activity Ideas (Optional for Teachers/Students)

1. **Create a Survey Form** using Google Forms to collect data on students' daily screen time.
 2. **Access a public dataset** from data.gov.in and identify whether it's structured or unstructured.
 3. **Use an API** like OpenWeatherMap to fetch live temperature data.
-

Summary

In this chapter, we revisited the **AI Project Cycle** with a focus on **Data Collection** and **Data Access**—two essential components of building effective AI solutions. We explored various **types and sources of data**, discussed **tools for collecting data**, and learned how to **access data** using different methods such as cloud storage, databases, and APIs. We also covered **legal and ethical responsibilities** associated with data usage. Remember, **data is the foundation** of any AI project—its quality, availability, and responsible handling determine the success of your AI model.
