14-12-2020

## Chapter-4
## Collecting Text Data from Web

Objectives:-
- Extract and process data from web pages
- Describe different kinds of semi-structured data, such as Json and XML
- Extract real time data using Application programming interfaces
- Extract data from various file format.

=> <u>Semi Structured Data</u>

1) <u>JSON</u>
   Java Script object notion

- Data is stored as key-value pair

- Datatype of values can be :
  - A String
  - A number
  - Another JSON object
  - An array
  - A Boolean
  - Null.

16-12-2020

2) XML

Extensible markup language.

- Human readable and extensible
- May on may not have a declaration, but if there's a declaration, it has to be the starting line
- The declaration has three parts → Version, encoding, standalone.

- Can be represented as a tree called XML tree.

⇒ Using API's to retrieve real time data.

Application programming interface ~~some~~

Some website do not provide their data directly. Instead they provide API's using which we can extract data from them.