Assignment Project Exam Help

Parsing Raw Data

https://powcoder.com

Faculty of Information Technology
Monash University, Australia

Add WeChat powcoder

FIT5196 Week 3

- Easy-to-parse (machine-readable) formats:
  - CSV: Comma Separated Values
  - JSON: JavaScript Object Notation
  - XML: eXtensible Markup Language
- Hard-to-parse formats:
  - Excel
  - PDF: Portable Document Format

- RDF: Resource Description Framework
  - A standard model for data interchange on the Web.
  - RDF has features that facilitate data merging even if the underlying schemas differ
  - RDF supports the evolution of schemas over time without requiring all the data consumers to be changed.
  - RDF is ideal for storing graph data, such as Knowledge Graphs.
  - Python libs
    - `http://rdflib.readthedocs.io/en/3.4.0/intro_to_graphs.html`
- HDF5 : Hierarchical Data Format
  - HDF5 contains an internal file system-like node structure
  - HDF5 can stores multiple datasets and supports metadata
  - HDF5 is a good choice for efficiently read and write large datasets.
  - Python libs
    - PyTables
    - h5py
    - pandas.HDFStore()

# CSV: Comma Separated Values



- TSV: Tab Separated Values
- Software: Microsoft Excel, Open Office Calc, and Google Spreadsheets.

# CSV: tools

- Pandas functions for reading tabular data
  - read_csv(): Read delimited data from a file, URL, or file-like object. Use comma as default delimiter.
  - read_table(): Read delimited data from a file, URL, or file-like object. Use tab as default delimiter.
  - read_fwf(): read data in fixed-width column format, i.e., no delimiters.
  - read_clipboard() // Version of read_table that reads data from the clipboard. Useful for converting tables from web pages.

**Outline**

MONASH University

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# XML: Extensible Markup Language[1]

- XML is a software- and hardware-independent tool for storing and transporting data.
  - It simplifies data sharing and platform changes — no need to worry about issues of exchanging data between incompatible systems
  - It simplifies data transport — XML stores data in plain text format
  - It simplifies data availability — With XML, data can be available to all kinds of "reading machines"
- XML was designed to be both human- and machine-readable.

_____

[1]Materials in the following 4 slides are based on
http://www.w3schools.com/xml/default.asp

# XML: DOM tree

```xml
<?xml version="1.0" encoding="UTF-8"?>
<bookstore>
  <book category="cooking">
    <title lang="en">Everyday Italian</title>
    <author>Giada De Laurentiis</author>
    <year>2005</year>
    <price>30.00</price>
  </book>
  <book category="children">
    <title lang="en">Harry Potter</title>
    <author>J K. Rowling</author>
    <year>2005</year>
    <price>29.99</price>
  </book>
  <book category="web">
    <title lang="en">XQuery Kick Start</title>
    <author>James McGovern</author>
    <author>Per Bothner</author>
    <author>Kurt Cagle</author>
    <author>James Linn</author>
    <author>Vaidyanathan Nagarajan</author>
    <year>2003</year>
    <price>49.99</price>
  </book>
  <book category="web" cover="paperback">
    <title lang="en">Learning XML</title>
    <author>Erik T. Ray</author>
    <year>2003</year>
    <price>39.95</price>
  </book>
</bookstore>
```

- According to the DOM (Document Object Model), everything in an XML document is a node.
- The DOM says:
  - The entire document is a document node
  - Every XML element is an element node
  - The text in the XML elements are text nodes
  - Every attribute is an attribute node
  - Comments are comment nodes

# XML: DOM tree

```xml
<?xml version="1.0" encoding="UTF-8"?>
<bookstore>
    <book category="cooking">
        <title lang="en">Everyday Italian</title>
        <author>Giada De Laurentiis</author>
        <year>2005</year>
        <price>30.00</price>
    </book>
    <book category="children">
        <title lang="en">Harry Potter</title>
        <author>J K. Rowling</author>
        <year>2005</year>
        <price>29.99</price>
    </book>
    <book category="web">
        <title lang="en">XQuery Kick Start</title>
        <author>James McGovern</author>
        <author>Per Bothner</author>
        <author>Kurt Cagle</author>
        <author>James Linn</author>
        <author>Vaidyanathan Nagarajan</author>
        <year>2003</year>
        <price>49.99</price>
    </book>
    <book category="web" cover="paperback">
        <title lang="en">Learning XML</title>
        <author>Erik T. Ray</author>
        <year>2003</year>
        <price>39.95</price>
    </book>
</bookstore>
```
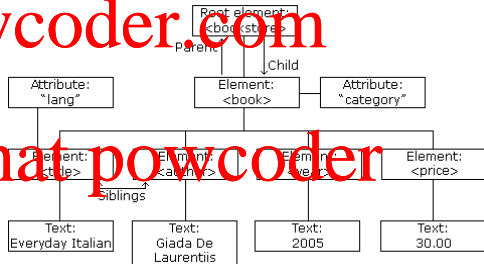
- According to the DOM (Document Object Model), everything in an XML document is a node.

# XML: tools

- ElementTree
  - https://docs.python.org/3/library/xml.etree.elementtree.html
  - Python's built-in XML parser.
- lxml:
  - http://lxml.de/
  - Strong performance in parsing very large files
- BeautifulSoup
  - https://www.crummy.com/software/BeautifulSoup/bs4/doc/
  - A Python library for pulling data out of HTML and XML files
  - Works with your favourite parser, e.g., html.parser and lxml-xml

Demonstration with Jupyter notebook.

# JSON: JavaScript Object Notation

- JSON: one of the most commonly used formats for transferring data between web services and other applications via HTTP.
- JSON is completely language independent but uses conventions that are familiar to programmers of the C-family of languages.
- JSON is built on two structures[2]:
  - A collection of name/value pairs. In various languages, this is realised as an object, record, struct, dictionary, hash table, keyed list, or associative array.
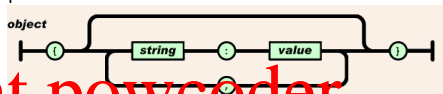  - An ordered list of values. In most languages, this is realised as an array, vector, list, or sequence.

---

[2]Materials on the following 3 slides are based on http://www.json.org/

# JSON: Structure

```json
{
  "firstName": "John",
  "lastName": "Smith",
  "isAlive": true,
  "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021-3100"
  },
  "phoneNumbers": [
    {
      "type": "home",
      "number": "212 555-1234"
    },
    {
      "type": "office",
      "number": "646 555-4567"
    },
    {
      "type": "mobile",
      "number": "123 456-7890"
    }
  ],
  "children": [],
  "spouse": null
}
```

- Three basic elements:
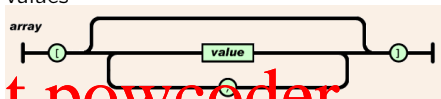  - Object: an unordered set of name/value pairs

# JSON: Structure

```
{
  "firstName": "John",
  "lastName": "Smith",
  "isAlive": true,
  "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021-3100"
  },
  "phoneNumbers": [
    {
      "type": "home",
      "number": "212 555-1234"
    },
    {
      "type": "office",
      "number": "646 555-4567"
    },
    {
      "type": "mobile",
      "number": "123 456-7890"
    }
  ],
  "children": [],
  "spouse": null
}
```

- Three basic elements:
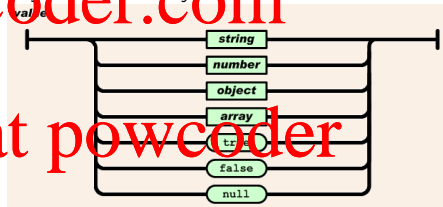  - Array: an array is an ordered collection of values

# JSON: Structure

```
{
  "firstName": "John",
  "lastName": "Smith",
  "isAlive": true,
  "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021-3100"
  },
  "phoneNumbers": [
    {
      "type": "home",
      "number": "212 555-1234"
    },
    {
      "type": "office",
      "number": "646 555-4567"
    },
    {
      "type": "mobile",
      "number": "123 456-7890"
    }
  ],
  "children": [],
  "spouse": null
}
```

- Three basic elements:
  - ▸ Value: a string in double quotes, or a number, or true or false or null, or an object or an array.

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# JSON v.s. XML

```xml
<?xml version="1.0" encoding="UTF-8"?>
<bookstore>
  <book category="cooking">
    <title lang="en">Everyday Italian</title>
    <author>Giada De Laurentiis</author>
    <year>2005</year>
    <price>30.00</price>
  </book>
  <book category="children">
    <title lang="en">Harry Potter</title>
    <author>J K. Rowling</author>
    <year>2005</year>
    <price>29.99</price>
  </book>
  <book category="web">
    <title lang="en">XQuery Kick Start</title>
    <author>James McGovern</author>
    <author>Per Bothner</author>
    <author>Kurt Cagle</author>
    <author>James Linn</author>
    <author>Vaidyanathan Nagarajan</author>
    <year>2003</year>
    <price>49.99</price>
  </book>
  <book category="web" cover="paperback">
    <title lang="en">Learning XML</title>
```

```json
{
  "firstName": "John",
  "lastName": "Smith",
  "isAlive": true,
  "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021-3100"
  },
  "phoneNumbers": [
    {
      "type": "home",
      "number": "212 555-1234"
    },
    {
      "type": "office",
      "number": "646 555-4567"
    },
    {
      "type": "mobile",
      "number": "123 456-7890"
```

# JSON: tools

- json: a built-in Python library used to parse JSON files.
- pandas json functions:
  - read_json(): Convert a JSON string to pandas object.
  - json_normalize(): "Normalise" semi-structured JSON data into a flat table

# PDF: Portable Document Format

- A file format used to present and exchange documents
  - ▶ "looks really do matter" from Adobe
  - ▶ PDF can contains text, image, link, button, form field, audio and video.
    PDF file encapsulates a complete description of the layout information, e.g.,
    fonts, graphics, and other meta information of the document.
- Not a data format

# PDF: An example

MONASH University



TABLE 2 | NUTRITION

# PDF: Parsing Tools

- pdfminer: A tool for extracting text, images, object coordinates, metadata from PDF documents.
- pdftable: A tool for extracting tables from PDF files; it uses pdfminer to get information on the locations of text elements.
- slate: A small Python module that wraps pdfminer's API.
- Tabula: A simple tool for extracting data tables out of PDF files

# Summary: what to do this week

1. Download, run and read the notebooks provided in Moodle, and also read the recommended reading materials associated with each notebook.
2. Try to finish the exercises in each chapter, and post your findings and experience in the discussion forum.
3. Attend tutorial 3 in the following week.
   - Pandas, Excel files
4. **Assessment 1** released.