# 2. Handling Raw Data in Different Formats

**Lan Du**

## 2. Handling Raw Data in Different Formats
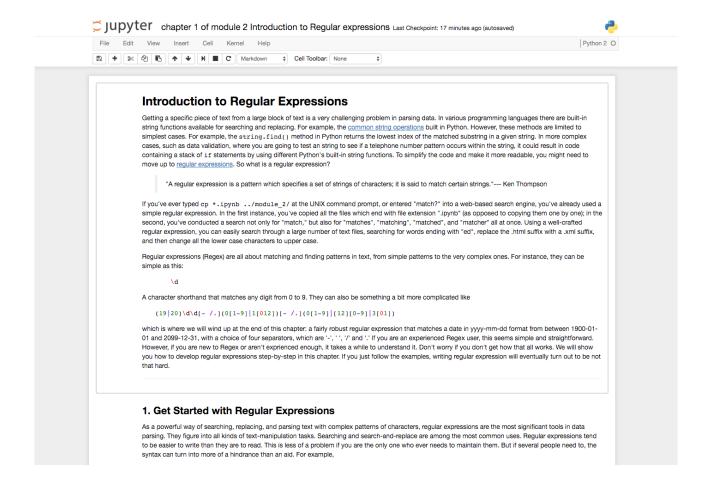
Lan Du

Generated by [Alexandria](https://www.alexandriarepository.org) (https://www.alexandriarepository.org) on July 22, 2016 at 10:23 am AEST

# Contents

# 1
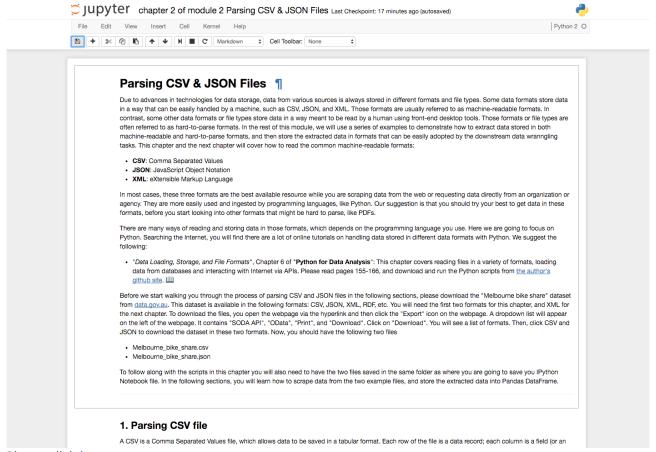# Introduction to Regular Expressions

This chapter introduces the fundamentals of regular expressions under the umbrella of Python. You will learn by example how to develop regular expressions for dealing with street addresses, Roman numerals, phone numbers, and dates. Unlike Module 1, the text and Python code are all included in one JuPyteR Notebook file, as shown below.

**Jupyter**  chapter 1 of module 2 Introduction to Regular expressions  Last Checkpoint: 17 minutes ago (autosaved)

Python 2

| File | Edit | View | Insert | Cell | Kernel | Help |

Markdown    Cell Toolbar: None

## Introduction to Regular Expressions

Getting a specific piece of text from a large block of text is a very challenging problem in parsing data. In various programming languages there are built-in string functions available for searching and replacing. For example, the common string operations built in Python. However, these methods are limited to simplest cases. For example, the `string.find()` method in Python returns the lowest index of the matched substring in a given string. In more complex cases, such as data validation, where you are going to test an string to see if a telephone number pattern occurs within the string, it could result in code containing a stack of `if` statements by using different Python's built-in string functions. To simplify the code and make it more readable, you might need to move up to regular expressions. So what is a regular expression?

> "A regular expression is a pattern which specifies a set of strings of characters; it is said to match certain strings."--- Ken Thompson

If you've ever typed `cp *.ipynb ../module_2/` at the UNIX command prompt, or entered "match?" into a web-based search engine, you've already used a simple regular expression. In the first instance, you've copied all the files which end with file extension ".ipynb" (as opposed to copying them one by one); in the second, you've conducted a search not only for "match," but also for "matches", "matching", "matched", and "matcher" all at once. Using a well-crafted regular expression, you can easily search through a large number of text files, searching for words ending with "ed", replace the .html suffix with a .xml suffix, and then change all the lower case characters to upper case.

Regular expressions (Regex) are all about matching and finding patterns in text, from simple patterns to the very complex ones. For instance, they can be simple as this:

`\d`

A character shorthand that matches any digit from 0 to 9. They can also be something a bit more complicated like

`(19|20)\d\d[- /.](0[1-9]|1[012])[- /.](0[1-9]|[12][0-9]|3[01])`

which is where we will wind up at the end of this chapter: a fairly robust regular expression that matches a date in yyyy-mm-dd format from between 1900-01-01 and 2099-12-31, with a choice of four separators, which are '-', ' ', '/' and '.' If you are an experienced Regex user, this seems simple and straightforward. However, if you are new to Regex or aren't experienced enough, it takes a while to understand it. Don't worry if you don't get how that all works. We will show you how to develop regular expressions step-by-step in this chapter. If you just follow the examples, writing regular expression will eventually turn out to be not that hard.

## 1. Get Started with Regular Expressions

As a powerful way of searching, replacing, and parsing text with complex patterns of characters, regular expressions are the most significant tools in data parsing. They figure into all kinds of text-manipulation tasks. Searching and search-and-replace are among the most common uses. Regular expressions tend to be easier to write than they are to read. This is less of a problem if you are the only one who ever needs to maintain them. But if several people need to, the syntax can turn into more of a hindrance than an aid. For example,

Please click here (https://www.alexandriarepository.org/wp-content/uploads/20151201030200/chapter-1-of-module-2.zip) to download the zip file that contains one notebook file. After unzipping the file, load the notebook with your Jupyter Notebook APP or with the command line. In order finish this chapter, you should read the text, run the code cell-by-cell and observe the output. You should also write Python code to finish the three exercises listed at the end of the notebook.

# 2
# Parsing CSV & JSON Files

This chapter introduces how to scrape data stored in CSV (Comma Separated Values) and JSON (JavaScript Object Notation) formats. You will learn to use Pandas' built-in function to read CSV and JSON files, and transform the loaded data into a well tabulated form. Similar to Chapter 1, the text content and Python code are all included in an JuPyteR (iPython) Notebook file, as shown below.
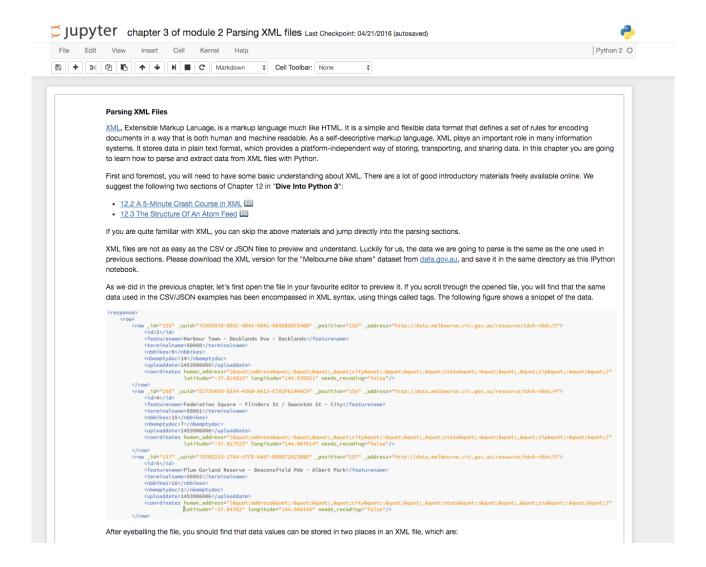


Please click here (https://www.alexandriarepository.org/wp-content/uploads/20151201042543/Chapter-2-of-module-2.zip) to download the zip file that contains the following list of files:

- chapter 2 of module 2 Parsing CSV & JSON Files.ipynb: the JuPyteR notebook containing all the Python code.
- csv1.png: an image showing what a CSV file looks like.
- elevations.json: a JSON dump used by the notebook
- json20.png: an image showing what a JSON file looks like.
- Melbourne_bike_share.csv: the CSV file to be parsed.
- Melbourne_bike_share.json: the JSON file to be parsed.

All the files listed above should be stored in the same folder so that you can run the notebook. Please read the text, run the code cell-by-cell and observe the output. You should also need to write Python code to finish the two exercises listed at the end of the notebook.

# 3
# Parsing XML files

This chapter discusses a number of ways of scraping data from XML files. You will learn the fundamental structure of XML format, and the use of different Python libraries to load/explore XML files and extract data from them. A screenshot of the notebook is shown below:
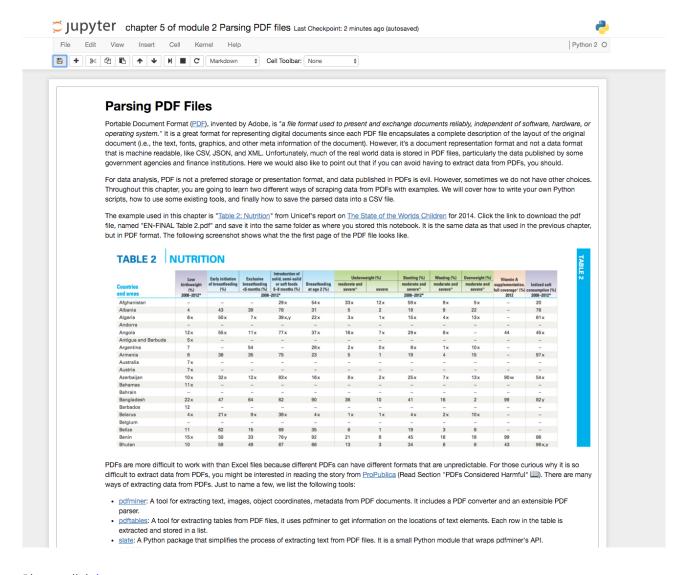


Please click [here](https://www.alexandriarepository.org/wp-content/uploads/20151201030200/Chapter-3-of-module-2.zip) (https://www.alexandriarepository.org/wp-content/uploads/20151201030200/Chapter-3-of-module-2.zip) to download the zip file that contains the following list of files:

- chapter 3 of module 2 Parsing XML files.ipynb: an IPython notebook containing all the Python code.
- Melbourne_bike_share.xml: an XML file to be parsed.
- xml_example.png: an image showing what am XML file looks like.

All the files listed above should be stored in the same folder so that you can run the notebook. Please read the text, run the code cell-by-cell and observe the output. You should also need to write Python code to finish the two exercises listed at the end of the notebook.

# 4
# Parsing PDF files

You might have found that extracting data from CSV, JSON, XML and even Excel files are not that hard. Then, how about PDF files? Can you extract tabular data stored in PDF files without scratching your head? PDF is a document representation format, rather than a proper data format. Unfortunately, scraping data from PDFs could not be avoidable since organisations, like government agencies and finance institutions prefer to use PDFs to release their data. This chapter uses an running example to show you the process of scraping simple tabular data from PDF pages. Here is a screenshot of the IPython notebook file:



Please click here (https://www.alexandriarepository.org/wp-content/uploads/20160502233906/Chapter-5-of-module-21.zip) to download the zip file that contains the following list of files:

- chapter 5 of module 2 Parsing PDF files.ipynb: an IPython notebook containing all the Python code.
- EN_FINAL_Table 2.pdf: a pdf file to be parsed.
- en_final_table_2.txt: a text file given by pdf2txt.py.
- en_final_table_2_1.csv: an output file.
- en_final_table_2_2.csv: am output file.
- EN_FINAL_Table_2_page_1.png: an image showing what a PDF page looks like. It is a screenshot of the first page of EN_FINAL_Table 2.pdf.

All the files listed above should be stored in the same folder so that you can run the notebook. Please read the text, run the code cell-by-cell and observe the output. You should also need to write Python code to finish the one exercise listed at the end of the notebook.