

Open in app ↗

Get unlimited access



Search Medium



Published in Towards Data Science

You have **2** free member-only stories left this month. [Upgrade for unlimited access.](#)



George Seif

Follow

Apr 29, 2019 · 3 min read · ✨ · 🎧 Listen



Save



The easy way to work with CSV, JSON, and XML in Python

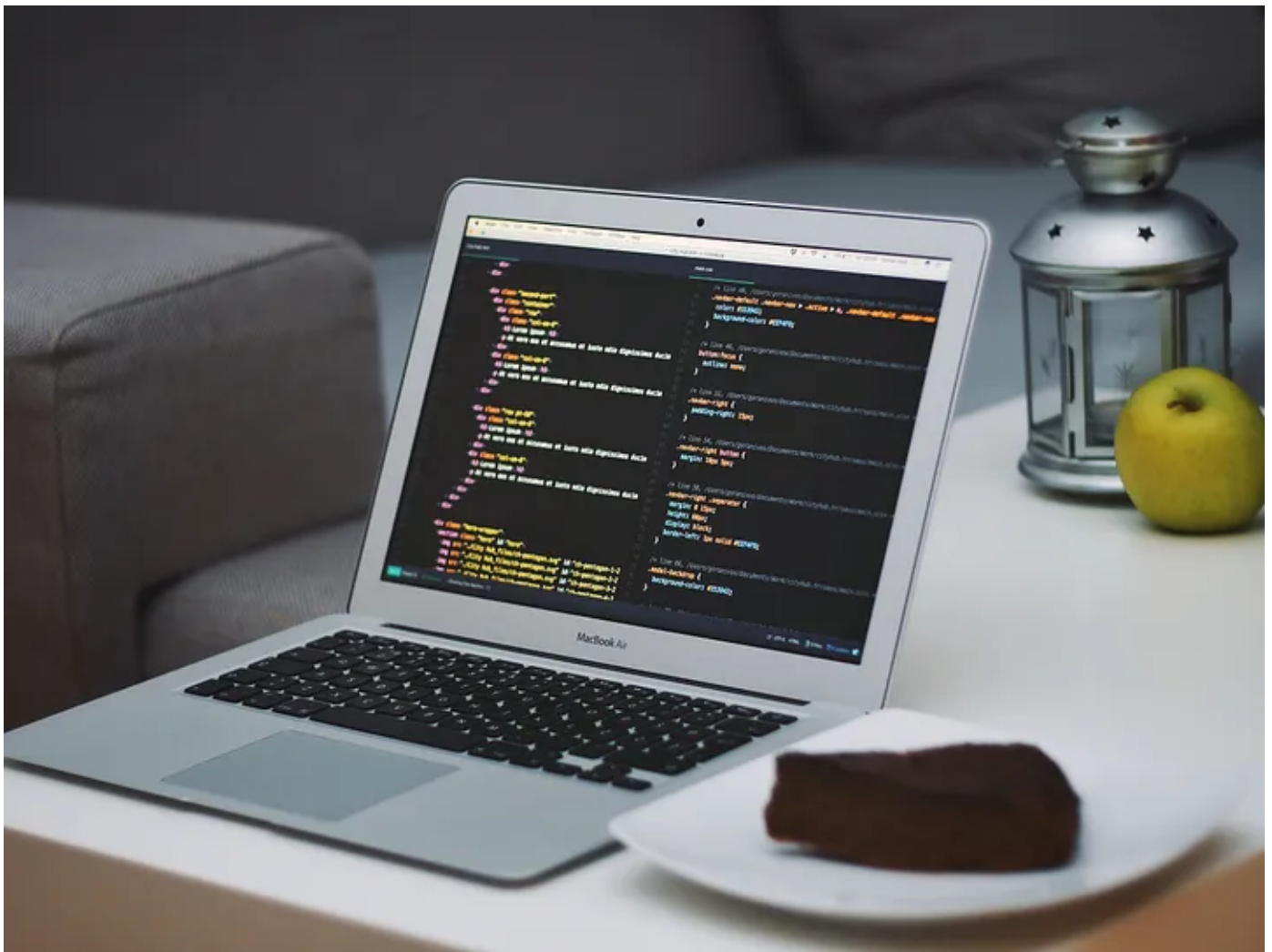


1.1K



1





Want to be inspired? Come join my *[Super Quotes newsletter](#)*. 😎

Python's superior flexibility and ease of use are what make it one of the most popular programming language, especially for Data Scientists. A big part of that is how simple it is to work with large datasets.

Every technology company today is building up a data strategy. They've all realised that having the right data: insightful, clean, and as much of it as possible, gives them a key competitive advantage. Data, if used effectively, can offer deep, beneath the surface insights that can't be discovered anywhere else.

Over the years, the list of possible formats that you can store your data in has grown significantly. But, there are 3 that dominate in their everyday usage: CSV, JSON, and XML. In this article, I'm going to share with you the easiest ways to work with these 3 popular data formats in Python!

CSV Data

A CSV file is the most common way to store your data. You'll find that most of the data coming from Kaggle competitions is stored in this way. We can do both read and write of a CSV using the built-in Python csv library. Usually, we'll read the data into a list of lists.

Check out the code below. When we run `csv.reader()` all of our CSV data becomes accessible. The `csvreader.next()` function reads a single line from the CSV; every time you call it, it moves to the next line. We can also loop through every row of the csv using a for-loop as with `for row in csvreader`. Make sure that you have the same number of columns in each row, otherwise, you'll likely end up running into some errors when working with your list of lists.

```
1  import csv
2
3  filename = "my_data.csv"
4
5  fields = []
6  rows = []
7
8  # Reading csv file
9  with open(filename, 'r') as csvfile:
10     # Creating a csv reader object
11     csvreader = csv.reader(csvfile)
12
13     # Extracting field names in the first row
14     fields = csvreader.next()
15
16     # Extracting each data row one by one
17     for row in csvreader:
18         rows.append(row)
19
20 # Printing out the first 5 rows
21 for row in rows[:5]:
22     print(row)
```

Writing to CSV in Python is just as easy. Set up your field names in a single list, and your data in a list of lists. This time we'll create a `writer()` object and use it to write our data to file very similarly to how we did the reading.

```
1  import csv
2
3  # Field names
4  fields = ['Name', 'Goals', 'Assists', 'Shots']
5
6  # Rows of data in the csv file
7  rows = [ ['Emily', '12', '18', '112'],
8           ['Katie', '8', '24', '96'],
9           ['John', '16', '9', '101'],
10          ['Mike', '3', '14', '82']]
11
12  filename = "soccer.csv"
13
14  # Writing to csv file
15  with open(filename, 'w+') as csvfile:
16      # Creating a csv writer object
17      csvwriter = csv.writer(csvfile)
18
19      # Writing the fields
20      csvwriter.writerow(fields)
21
22      # Writing the data rows
23      csvwriter.writerows(rows)
```

csv_python_2.py hosted with ❤ by GitHub

[view raw](#)

Of course, installing the wonderful Pandas library will make working with your data far easier once you've read it into a variable. Reading from CSV is a single line as is writing it back to file!

```
1  import pandas as pd
2
3  filename = "my_data.csv"
4
5  # Read in the data
6  data = pd.read_csv(filename)
7
8  # Print the first 5 rows
9  print(data.head(5))
10
11 # Write the data to file
12 data.to_csv("new_data.csv", sep=",", index=False)
```

csv_python_3.py hosted with ❤ by GitHub

[view raw](#)

We can even use Pandas to convert from CSV to a list of dictionaries with a quick one-liner. Once you have the data formatted as a list of dictionaries, we'll use the `dicttoxml` library to convert it to XML format. We'll also save it to file as a JSON!

```

1  import pandas as pd
2  from dicttoxml import dicttoxml
3  import json
4
5  # Building our dataframe
6  data = {'Name': ['Emily', 'Katie', 'John', 'Mike'],
7         'Goals': [12, 8, 16, 3],
8         'Assists': [18, 24, 9, 14],
9         'Shots': [112, 96, 101, 82]
10        }
11
12  df = pd.DataFrame(data, columns=data.keys())
13
14  # Converting the dataframe to a dictionary
15  # Then save it to file
16  data_dict = df.to_dict(orient="records")
17  with open('output.json', "w+") as f:
18      json.dump(data_dict, f, indent=4)
19
20  # Converting the dataframe to XML
21  # Then save it to file
22  xml_data = dicttoxml(data_dict).decode()
23  with open("output.xml", "w+") as f:
24      f.write(xml_data)

```

csv_python_4.py hosted with ❤ by GitHub

[view raw](#)

JSON Data

JSON provides a clean and easily readable format because it maintains a dictionary-style structure. Just like CSV, Python has a built-in module for JSON that makes reading and writing super easy! When we read in the CSV, it will become a dictionary. We then write that dictionary to file.

```
1  import json
2  import pandas as pd
3
4  # Read the data from file
5  # We now have a Python dictionary
6  with open('data.json') as f:
7      data_listofdict = json.load(f)
8
9  # We can do the same thing with pandas
10 data_df = pd.read_json('data.json', orient='records')
11
12 # We can write a dictionary to JSON like so
13 # Use 'indent' and 'sort_keys' to make the JSON
14 # file look nice
15 with open('new_data.json', 'w+') as json_file:
16     json.dump(data_listofdict, json_file, indent=4, sort_keys=True)
17
18 # And again the same thing with pandas
19 export = data_df.to_json('new_data.json', orient='records')
```

json_python_1.py hosted with ❤ by GitHub

[view raw](#)

And as we saw before, once we have our data you can easily convert to CSV via pandas or use the built-in Python CSV module. When converting to XML, the `dicttoxml` library is always our friend.

```
1 import json
2 import pandas as pd
3 import csv
4
5 # Read the data from file
6 # We now have a Python dictionary
7 with open('data.json') as f:
8     data_listofdict = json.load(f)
9
10 # Writing a list of dicts to CSV
11 keys = data_listofdict[0].keys()
12 with open('saved_data.csv', 'wb') as output_file:
13     dict_writer = csv.DictWriter(output_file, keys)
14     dict_writer.writeheader()
15     dict_writer.writerows(data_listofdict)
```

json_python_2.py hosted with ❤ by GitHub

[view raw](#)

XML Data

XML is a bit of a different beast from CSV and JSON. Generally, CSV and JSON are widely used due to their simplicity. They're both easy and fast to read, write, and interpret as a human. There's no extra work involved and parsing a JSON or CSV is very lightweight.

XML on the other hand tends to be a bit heavier. You're sending more data, which means you need more bandwidth, more storage space, and more run time. But XML does come with a few extra features over JSON and CSV: you can use namespaces to build and share standard structures, better representation for inheritance, and an industry standardised way of representing your data with XML schema, DTD, etc.

To read in the XML data, we'll use Python's built-in XML module with sub-module ElementTree. From there, we can convert the ElementTree object to a dictionary using the `xmltodict` library. Once we have a dictionary, we can convert to CSV, JSON, or Pandas Dataframe like we saw above!


```
1  import xml.etree.ElementTree as ET
2  import xmltodict
3  import json
4
5  tree = ET.parse('output.xml')
6  xml_data = tree.getroot()
7
8  xmlstr = ET.tostring(xml_data, encoding='utf8', method='xml')
9
10
11 data_dict = dict(xmltodict.parse(xmlstr))
12
13 print(data_dict)
14
15 with open('new_data_2.json', 'w+') as json_file:
16     json.dump(data_dict, json_file, indent=4, sort_keys=True)
```

xml_python_1.py hosted with ❤ by GitHub

[view raw](#)

Like to learn?

Follow me on [twitter](#) where I post all about the latest and greatest AI, Technology, and Science! Connect with me on [LinkedIn](#) too!

Artificial Intelligence

Technology

Innovation

Programming

Coding

Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

Emails will be sent to humera.shaziya@gmail.com. [Not you?](#)



Get this newsletter