

Pythonではじめる クローリング・スクレイピング

ここからはPythonを使ってクローリング・スクレイピングを行います。Pythonの開発環境構築と基本文法から、クローリング・スクレイピングの一連の過程をスクリプトにするまでを解説します。

2.1 Pythonを使うメリット

クローリング・スクレイピングにPythonを使うメリットとして、第1章のまとめで述べたように汎用的なプログラミング言語であることを挙げました。さらに次のメリットがあります。

- 言語自体の特性
- 強力なライブラリ
- スクレイピング後の処理との親和性

2.1.1 言語自体の特性

1点目のメリットはPythonの言語そのものの特性です。

Pythonは教育用に使われることもある読みやすく書きやすい言語です。一度書かれたプログラムは、その後何度も他の人(未来の自分を含む)に読まれることになるため、読みやすさは重要です。

PythonはBattery Included(電池付属)と言われています。電池つき電化製品のように、豊富な標準ライブラリが付属しており、インストール後すぐに使いはじめられることの比喩です。ここでは、Pythonの標準ライブラリだけでクローリング・スクレイピングを行います。後の章でより強力なサードパーティライブラリも活用しますが、手軽に使いはじめられるPythonの良さがわかるはずです。

さらに、複数のWebサイトから高速にデータを取得するためには、非同期処理が有効です。PythonにはTwistedやTornadoなどの非同期処理のためのフレームワークが存在し、Python 3.4からはasyncioと呼ばれる非同期処理のための標準ライブラリもあります。非同期処理の分野においてはNode.jsやGo言語が有名ですが、Pythonでも手軽に扱えます。

2.1.2 強力なサードパーティライブラリの存在

2点目は豊富なサードパーティライブラリの存在です。Python Package Index (PyPI)^{*1}には、世界中の開発者が数多くのライブラリを公開しており、簡単に使うことができます。特に、第3章で紹介するBeautiful Soupやlxmlは有名なスクレイピングライブラリですし、第6章で紹介するScrapyは強力なクローリング・スクレイピングフレームワークです。このような強力なライブラリによって巨人の肩の上に乗り、短いプログラムを書くだけで素早くクローリング・スクレイピングを行えます。

2.1.3 スクレイピング後の処理との親和性

3点目はクローリング・スクレイピングでデータを取得した後、データ分析などの処理を行う際にもPythonが強力な武器になる点です。データ分析においてもPythonには優秀なライブラリが揃っているため、1つの言語を修得するだけで大抵のことは実現できてしまいます。

Pythonでは数値計算や科学技術計算の分野で古くからNumPyやSciPyといったライブラリが有名で、これらをベースとしたデータ分析用のライブラリが存在します。例えばpandas(5.3.2参照)は、NumPyをベースとしてデータの前処理(欠損値や表記ゆれの処理)や集計を簡単に行えるライブラリです。またmatplotlib(5.3.3参照)は数値データをグラフで可視化できます。データ分析の分野ではR言語が有名ですが、これらのライブラリでPythonでも同様の分析が行えます。他にも様々なライブラリがあります。

2.2 Pythonのインストールと実行

Python 3をインストールし、仮想環境内で実行します。

2.2.1 Python 2とPython 3

Pythonには大きく分けて、Python 2系とPython 3系の2つのバージョンがあります。Pythonは後方互換性を重視する言語ですが、2008年にリリースされたバージョン3.0において後方互換性を崩す大きな変更が入りました。そのためライブラリの3系への対応が進まず、2系も並行して開発されてきた経緯があります。

近年では多くのライブラリが3系に対応し、Python 3を問題なく使える状況になってきています。3

*1 <https://pypi.python.org/pypi>

系は2系に比べて、よりわかりやすい文法への変更、Unicodeサポートの強化、標準ライブラリの整理、非同期I/Oのサポートなど様々な改善が行われています。2系のサポートは2020年で打ち切られる予定になっていることもあり、今から使いはじめるのであれば3系を使うのがオススメです。本書ではPython 3のみを使用し、Python 3.5.1 (OS X) と Python 3.4.3 (Ubuntu) を対象として解説します。

2.2.2 パッケージマネージャーによるPython 3のインストール

OS XではHomebrew (1.2.1参照) で、UbuntuではAPTでインストールします*2*3*4。

```
$ brew install python3 # OS Xの場合
```

```
$ sudo apt-get install -y python3 python3.4-venv # Ubuntuの場合
```

インストールに成功すると、python3コマンドでバージョンを確認できます。

```
$ python3 -V
Python 3.5.1
```

2.2.3 仮想環境 (venv) の使用

近年のプログラミング言語では**仮想環境 (Virtual Environment)** と呼ばれる、ランタイムやライブラリを環境(用途)ごとに分離できる仕組みが多く使われています。

● 仮想環境とは

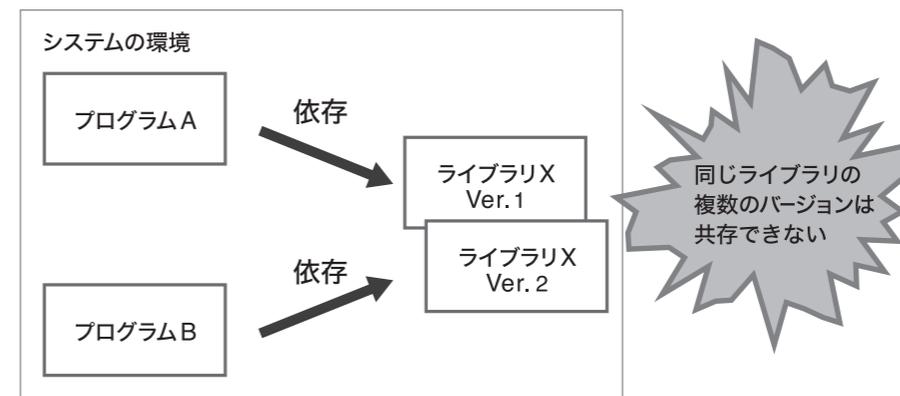
例えば、1台のコンピュータで2つの異なるプログラムAとBを書いているとしましょう。2つのプログラムはXというライブラリに依存しており、プログラムAではライブラリXのバージョン1に、プログラムBではライブラリXのバージョン2に依存しているとします。Pythonでは1つの環境に同じライブラリの複数バージョンをインストールできないため、インストールできるのはXのバージョン1か2のいずれかのみです。ライブラリXのバージョン1と2に互換性がない場合、プログラムAとBのどちらかは正常に動作しなくなってしまいます(図2.1)。

*2 本書で使用するプラットフォームでは、Python 2系と3系はコマンド名やライブラリのインストール場所が異なるので共存できません。3系をインストールしても、既存の2系が使えなくなることはありません。

*3 Pythonをインストールするためのツールとして、pyenv (<https://github.com/yyuu/pyenv>) やpythonz (<https://github.com/saghul/pythonz>) などが存在します。しかしPythonはメジャーバージョンが同じであれば基本的に後方互換があるので、マイナーバージョンが異なるPythonを使い分けたいといったこだわりがない限り、利用する必要はないでしょう。

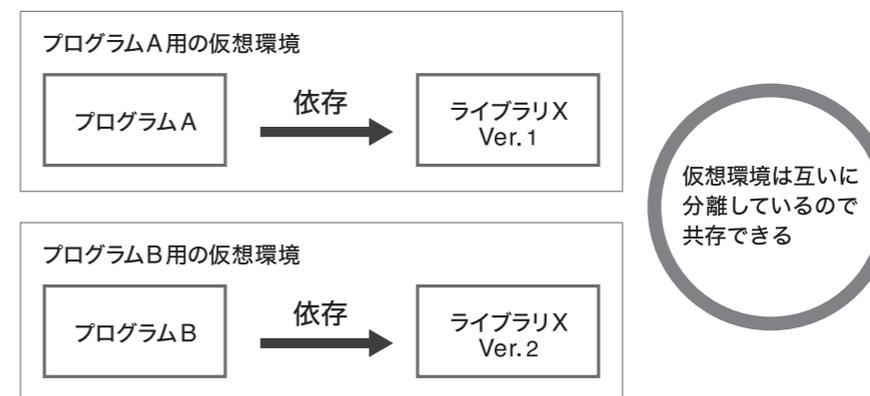
*4 仮想環境を利用するために、Ubuntuではvenvパッケージをインストールする必要があります。本書の解説には用いていませんがUbuntu 16.04ではpython3.4-venvパッケージの代わりにpython3-venvパッケージをインストールします。

▼ 図2.1 仮想環境がない時



仮想環境を使うと、このような事態を避けられます。プログラムA用の仮想環境にXのバージョン1をインストールし、プログラムB用の仮想環境にXのバージョン2をインストールすることで、互いに干渉することなくプログラムを開発できます(図2.2)。

▼ 図2.2 仮想環境がある時



Pythonでは**venv**という標準モジュールで仮想環境を利用できます*5。

仮想環境内ではpython3コマンドのようにバージョン番号がついたコマンドを使わなくても、pythonコマンドで仮想環境に紐付けられたPythonランタイムを起動できるメリットもあります。また、仮想

*5 Pythonでは仮想環境を使うためにvirtualenv (<https://virtualenv.pypa.io/en/latest/>) というサードパーティのツールが広く使われてきました。Python 3.3以降はvenvを使えばvirtualenvを使う必要はありません。

環境は通常の開発で広く使われているため、覚えておいて損はないでしょう。

2

● 仮想環境の使い方

次のコマンドで、仮想環境を作成します*6*7。-m オプションは指定したモジュールをスクリプトとして実行することを意味します。

```
$ python3 -m venv scraping
```

カレントディレクトリに `scraping` ディレクトリが作成されます。このディレクトリの名前が仮想環境の名前になります。ディレクトリの名前は自由に設定できます*8。

```
$ ls scraping/
bin      include  lib      pyvenv.cfg
```

仮想環境に入るためには、`.` (ドット) コマンドで `activate` スクリプトを実行します。

```
$ . scraping/bin/activate # ドットと引数の間にスペースを入れる。
```

`.` コマンドは引数で指定したファイルを読み込み、現在のシェルで実行するコマンドです。Bash や Zsh などでは `source` コマンドも同じ意味を持ちます。

仮想環境に入ると、シェルのプロンプトの先頭に `(scraping)` と表示されるようになります。

```
(scraping) $
```

仮想環境内では `python` コマンドで Python 3 が実行できます。

```
(scraping) $ python -V
Python 3.5.1
```

`which` コマンドで `python` コマンドのパスを確認すると、`scraping/bin` ディレクトリ内の `python` コマンドを指していることがわかります。

```
(scraping) $ which python
/path/to/scraping/bin/python
```

- *6 Windows の VirtualBox 上で仮想マシンを実行している場合、共有フォルダとしてマウントされているディレクトリ (Appendix の手順のデフォルトでは `/vagrant/`) 内に仮想環境を作成するとシンボリックリンクの作成に失敗します。仮想環境はホームディレクトリ (`~`) など、共有フォルダ以外のディレクトリに作成してください。
- *7 `venv` モジュール導入当初は、`pyenv` というコマンドで仮想環境を作成していましたが、このコマンドは Python 3.6 以降では `Deprecated` になる予定です。
- *8 仮想環境を作成した後にディレクトリをリネームしたり移動したりすると、正常に動作しなくなるので注意してください。

2

`deactivate` コマンドで仮想環境から抜けることができます。仮想環境から抜けるとシェルのプロンプトから `(scraping)` の文字がなくなり、元に戻ります。

```
(scraping) $ deactivate
```

仮想環境自体が不要になった場合はディレクトリをまるごと削除します。

以降の説明では、特に断りのない限り仮想環境内でコマンドを実行します。書かれている通りに実行しても結果が異なったり予期せぬエラーが発生する場合は、まず仮想環境内で Python 3 を使用しているか確認してください。

2.2.4 インタラクティブシェルの使用

`python` コマンドを引数なしで実行すると、インタラクティブシェルが起動します。Python のコードを対話的に実行できるので、ライブラリの使い方の確認などに便利です。

```
(scraping) $ python
Python 3.5.1 (default, Apr 18 2016, 03:49:24)
[GCC 4.2.1 Compatible Apple LLVM 7.0.2 (clang-700.1.81)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

1~3 行目にはバージョンなどの情報と簡単な説明が表示され、4 行目に `>>>` が表示されて入力を受け付ける状態になります。この `>>>` をプロンプトと呼びます。プロンプトのある行に Python の式を入力して Enter を押すと、次の行にその式の値が表示されます。

```
>>> 1 + 1
2
```

プロンプトが表示されている時に `[Ctrl]+D` を押すか、`exit()` と入力して Enter を押すとインタラクティブシェルを終了できます。

入力中の場合は、`[Ctrl]+C` を押すと入力中の文字列がリセットされて、プロンプトが表示されます。通常のシェルと同じように、上下の矢印キーで以前の入力を表示したり、`[Ctrl]+R` で以前の入力からインクリメンタルサーチすることも可能です。

Python のオンラインマニュアルでもインタラクティブシェル形式の解説が多くあります。新しい機能を使うときなどはインタラクティブシェルで試す習慣をつけると良いでしょう。

2.3 Pythonの基礎知識

本書を読み進め、実際にクローリング・スクレイピングする上で必要な、Pythonの文法をはじめとした基礎知識を解説します。最初は流し読みして、後で不明点があったときに戻ってきても良いでしょう。

2.3.1 スクリプトファイルの実行と構成

Pythonスクリプトファイルの実行方法と構成を解説します。

● Pythonスクリプトファイルの実行

Pythonのスクリプトは.pyという拡張子のファイルに保存します。テキストエディターを使って、リスト2.1の中身を持つファイルをgreet.pyという名前で作成してください。ファイルのエンコーディングはUTF-8を使います。

▼ リスト2.1 greet.py — Pythonスクリプトの例

```
import sys

def greet(name):
    print('Hello, {}!'.format(name))

if len(sys.argv) > 1:
    name = sys.argv[1]
    greet(name)
else:
    greet('world')
```

Pythonのスクリプトファイルを実行するには、ファイルパスをpythonコマンドの引数に渡します。ファイルを保存したら、仮想環境内で次のコマンドを実行して、「Hello, world!」と表示されることを確認します。

```
(scraping) $ python greet.py
Hello, world!
```

この例では引数に文字列を与えて実行すると、表示が変わります。

```
(scraping) $ python greet.py Guido
Hello, Guido!
```

● Pythonスクリプトの構成

Pythonを使ったことがない方でも、他の言語の経験があればある程度リスト2.1のコードの意味するところが理解できるのではないのでしょうか。このコードにコメントをつけるとリスト2.2のようになります。行の#以降はコメントです。

▼ リスト2.2 greet_with_comments.py — Pythonスクリプトの例 (コメントつき)

```
import sys # import文でsysモジュールを読み込む。

# def文でgreet()関数を定義する。インデントされている行が関数の中身を表す。
def greet(name):
    print('Hello, {}!'.format(name)) # 組み込み関数print()は文字列を出力する。

# if文でもインデントが範囲を表す。sys.argvはコマンドライン引数のリストを表す変数。
if len(sys.argv) > 1:
    # if文の条件が真のとき
    name = sys.argv[1] # 変数は定義せずに代入できる。
    greet(name) # greet()関数を呼び出す。
else:
    # if文の条件が偽のとき
    greet('world') # greet()関数を呼び出す。
```

Pythonのスクリプトは上から順に実行されます。関数は事前に定義されている必要があります。基本的に1行に1文だけを書き、行末にセミコロンなどの記号は不要です。

Pythonでは**インデント**が大きな意味を持ちます。次のようにブロックをインデントで表します。

```
if a == 1:
    print('a is 1')
else:
    print('a is not 1')
```

C言語やJavaScriptでは、if文などでブロックを表すときに{}(波括弧)を使い、多くの場合は読みやすいようにブロック内をインデントします。

```
if (a == 1) {
    printf("a is 1\n");
} else {
    printf("a is not 1\n");
}
```

しかし、ブロック内をインデントしなくても意味は変わりません。

第 6 章 フレームワーク Scrapy

6

ここまでは、個々のライブラリを組み合わせることでクロール・スクレイピングを行ってきました。様々な Web サイトを対象にクローラーのプログラムを書いていると、同じような処理を繰り返し書くことに気づくかもしれません。この手間を省くために使えるのが、クロール・スクレイピングのためのフレームワーク、Scrapy です。Scrapy を使うと、どんな Web サイトでも使える共通処理をフレームワークに任せて、ユーザーは個々の Web サイトごとに異なる処理だけを書けばよくなります。

短いコードで効率的にクロール・スクレイピングできるので、様々なサイトからデータを抜き出したい場合や、継続的にクロールを行いたい場合には、学習コストを払うだけの価値があるでしょう。

6.1 Scrapy の概要

Scrapy^{*1}はクロール・スクレイピングのための Python のフレームワークです。豊富な機能が備わっており、ユーザーはページからデータを抜き出すという本質的な作業に集中できます。

- Web ページからのリンクの抽出
- robots.txt の取得と拒否されているページのクロール防止
- XML Sitemap の取得とリンクの抽出
- ドメインごと / IP アドレスごとのクロール時間間隔の調整
- 複数のクロール先の並行処理
- 重複する URL のクロール防止
- エラー時の回数制限付きのリトライ
- クローラーのデーモン化とジョブの管理

これまでは個々の機能を持つライブラリを自分の書いたプログラムから呼び出して利用してきました。フレームワークである Scrapy では、流儀にしたがってプログラムを書き、それらを Scrapy が呼び出して実行します。新しく流儀を覚える必要はありますが、一度覚えてしまえば面倒な処理をフレームワークに任せられるため、手軽にクロール・スクレイピングができます。

*1 <https://scrapy.org/> 本書ではバージョン 1.1.0 を使用します。

6

Scrapy はイベント駆動型のネットワークプログラミングのフレームワークである Twisted^{*2} をベースにしており、Web サイトからのダウンロード処理は非同期に実行されます。このため、ダウンロードを待つ間にもスクレイピングなど別の処理を実行でき、効率よくクロール・スクレイピングできます。Scrapy のような、スクレイピングのための多機能なフレームワークは他の言語ではあまり見かけず^{*3}、これらの用途に Python を使う大きな理由の一つだと言えます。

Scrapy は長らく Python 3 では動作しませんでしたでしたが、2016 年 5 月にリリースされたバージョン 1.1 から Python 3 に対応しました。

6.1.1 Scrapy のインストール

Scrapy をインストールします。Ubuntu では OpenSSL の開発用パッケージが必要なので一緒にインストールします。Scrapy は lxml に依存しているため、libxml2 と libxslt の開発用パッケージも必要です (3.3.2 参照)^{*4}。

```
(scraping) $ sudo apt-get install -y libssl-dev libffi-dev # Ubuntu の場合
```

```
(scraping) $ pip install scrapy
```

インストールに成功すると scrapy コマンドが使えるようになります。

```
(scraping) $ scrapy version
Scrapy 1.1.0
```

6.1.2 Spider の実行

Scrapy を使うのに、主に作成するのが Spider というクラスです。対象の Web サイトごとに Spider を作成し、クロールの設定や、スクレイピングの処理を記述します。

まずは簡単な Spider を実行してみましょう。リスト 6.1 は、Scrapy のサイトに掲載されているサンプルコード^{*5}にコメントを加えたものです。

*2 <https://twistedmatrix.com/>

*3 近いものとして Ruby の Anemone (<http://anemone.rubyforge.org/>) がありますが、Scrapy よりも薄いフレームワークで、残念ながら 2012 年で開発が止まっています。

*4 Scrapy のような C 拡張ライブラリに依存するライブラリをインストールする場合、依存ライブラリのビルドに失敗してもインストール済みになってしまい、正常に動作しないことがあります。その場合、ビルド失敗の原因を取り除き、一度アンインストールするとインストールし直せます。また、インストール中に問題が発生する場合は `pip install -U pip` で pip 自体をアップグレードすると問題が解決することもあります。

*5 <https://scrapy.org/> のページ中央

▼ リスト6.1 mspider.py — Scrapinghub社のブログから投稿のタイトルを取得するSpider

```
import scrapy

class BlogSpider(scrapy.Spider):
    name = 'blogspider' # Spiderの名前。
    # クローリングを開始するURLのリスト。
    start_urls = ['https://blog.scrapinghub.com']

    def parse(self, response):
        """
        トップページからカテゴリページへのリンクを抜き出してたどる。
        """
        for url in response.css('ul li a::attr("href")').re('.*category/.*'):
            yield scrapy.Request(response.urljoin(url), self.parse_titles)

    def parse_titles(self, response):
        """
        カテゴリページからそのカテゴリの投稿のタイトルをすべて抜き出す。
        """
        for post_title in response.css('div.entries > ul > li a::text').extract():
            yield {'title': post_title}
```

このSpiderは、ScrapyをメンテナンスしているScrapinghub社のブログをクローリングします。一覧・詳細パターンWebサイトの処理するSpiderで、ブログのトップページ(一覧ページ)からカテゴリページ(詳細ページ)へのリンクを抽出してたどり、カテゴリページからそのカテゴリに含まれるすべての投稿のタイトルを取得します。リンクをたどる流れを図で表します。

```
/ (トップページ)
├──>/category/autoscraping/ (カテゴリページ)
├──>/category/professional-services/
├──>/category/tools/
└──>...
```

scrapyコマンドのrunspiderサブコマンド(以降ではscrapy runspiderコマンドと表記します)の引数にファイルパスを指定して実行します。ログが表示され、数秒程度でクローリングが完了します。

```
(scraping) $ scrapy runspider mspider.py -o items.json
2016-05-25 19:00:16 [scrapy] INFO: Scrapy 1.1.0 started (bot: scrapybot)
2016-05-25 19:00:16 [scrapy] INFO: Overridden settings: {'FEED_URI': 'items.json', 'FEED_FORMAT': 'json'}
...
2016-05-25 19:00:16 [scrapy] INFO: Spider opened
2016-05-25 19:00:16 [scrapy] INFO: Crawled 0 pages (at 0 pages/min), scraped 0 items (at 0 items/min)
2016-05-25 19:00:17 [scrapy] DEBUG: Crawled (200) <GET https://blog.scrapinghub.com> (referer: None)
2016-05-25 19:00:18 [scrapy] DEBUG: Crawled (200) <GET https://blog.scrapinghub.com/category/autoscraping/> (referer: https://blog.scrapinghub.com)
```

```
2016-05-25 19:00:18 [scrapy] DEBUG: Scraped from <200 https://blog.scrapinghub.com/category/autoscraping/>
{'title': 'Announcing Portia, the Open Source Visual Web\0a0Scraper! on'}
2016-05-25 19:00:18 [scrapy] DEBUG: Scraped from <200 https://blog.scrapinghub.com/category/autoscraping/>
{'title': 'Introducing Dash on'}
...
```

作成されたitems.jsonファイルの中身を見ると、投稿のタイトルがJSON Lines形式で取得できていることがわかります。**JSON Lines**^{*6}形式とは各行にJSONオブジェクトを持つテキスト形式で、ファイルへの追記が容易な点が特徴です。

```
(scraping) $ cat items.json
{"title": "Announcing Portia, the Open Source Visual Web\0a0Scraper! on"}
{"title": "Introducing Dash on"}
{"title": "Spiders activity graphs on"}
{"title": "Autoscraping casts a wider\0a0net on"}
{"title": "Scrapy Tips from the Pros May 2016\0a0Edition on"}
...
```

ここでは完成したものを実行するだけでしたが、以降でSpiderを少しずつ作りながら、その動作を解説していきます。

6.2 Spiderの作成と実行

Yahoo!ニュースを対象に、基礎的なSpiderを作成します。

- Yahoo!ニュース
<http://news.yahoo.co.jp/>

作成するSpiderは、Yahoo!ニュースのトップページに表示されているトピックスの一覧(図6.1)から個別のトピックスへのリンクをたどり、トピックスのタイトルと本文を抽出するというものです。

*6 <http://jsonlines.org/>