

240ページの「Googleニュースをスクレイピングする」についての補足

2018年10月16日（追記あり）

Googleニュースのページの仕様変更により、この節で紹介している原著のコードではヘッドライン記事のURLを集められなくなりました。

そこで、日経トレンドネットのページから、記事のURLをスクレイピングできるコード（プログラム）を用意いたしました。記事のURLを単純にスクレイピングする `scraping1.py` と、重複するURLを取り除いてスクレイピングする `scraping2.py` の2種類がございます。

240ページで紹介しているBeautifulSoupのインストールを済ませれば、これらのプログラムを動かします。プログラムの概要については、プログラム中のコメントをご参照ください。

【2020年6月4日追記】

上記の日経トレンドネットが日経クロストrendに移行したことにもない、サイトの構造が変わりました。そこで、コードを以下のように修正いたしました。

1. URL を変更：
【旧】 `https://trendy.nikkeibp.co.jp/news/`
【新】 `https://xtrend.nikkei.com/atcl/contents/new/`
2. URL パス条件を変更：
【旧】 `if 'atcl/news'`
【新】 `if 'atcl/contents'`
3. (おまけ)リダイレクト対策として：
【旧】 `urljoin(self.site, url)`
【新】 `urljoin(r.url, url)`

以上により、記事の URL をスクレイピングできるコードについては、`chap20-crawler-nikkeibp-trendy-news.py` と、重複する URL を取り除いてスクレイピングする `chap20-crawler-nikkeibp-trendy-news-nodup.py` の2種類をご用意しています。

以上、補足いたします。