

Level 1: コーパスを用意し、pd.DataFrame として読み込もう。

コーパスの説明

・コーパスの概要

Livedoor ニュースコーパス

展開すると livedoor ニュースの項目ごとに文書(text)がまとめられており、その中でも dokujo-tsushin というカテゴリを採択。そのままでは文書数 852, 総文字数 1368059 とデータが多かった為、処理しやすいように上から 20 文書、総文字数(text)29772 とした。

・コーパスの取得手順

出典: <https://www.rondhuit.com/download.html#news%20corpus>

ldcc-20140209.tar.gz をダウンロード。

・pd.DataFrame.head() の出力結果。

```
(base) asahinatarou@talol rep3 % /opt/miniconda3/bin/python /Users/taro/DataMining/rep3/level1.py
URL date head text
0 http://news.livedoor.com/article/detail/4778030/ 2010-05-22T14:30:00+0900 友人代表のスピーチ、独女はどうこなしている？ もうすぐジューン・ブライドと呼ば
れる6月、独女の中には自分の式はまだなに呼ばれてばかり...
1 http://news.livedoor.com/article/detail/4778031/ 2010-05-21T14:30:00+0900 ネットで断ち切れない元カレとの縁 携帯電話が普及する以前、恋人への連絡ソ
ールは一般電話が普通だった。恋人と別れたら、手帳に書く...
2 http://news.livedoor.com/article/detail/4782522/ 2010-05-23T11:00:00+0900 相次ぐ芸能人の“すっぴん”披露 その時、独女の心境は？ 「男性はやっぱり、女性の“すっ
ぴん”が大好きなんですかね」と不満そうに話すのは、出版関係で働...
3 http://news.livedoor.com/article/detail/4788357/ 2010-05-25T14:00:00+0900 ムダな抵抗！？ ヒップの加齢による変化は「たわむ-下がる-内に流
れる」、バストは「そげる-たわむ-外に流れる...
4 http://news.livedoor.com/article/detail/4788362/ 2010-05-26T14:00:00+0900 税金を払うのは私たちなんですけど！ 6月から支給される子ども手当だが、当初
は子ども一人当たり月額2万6000円が支給されるはずだ...
```

文書数、総文字数

```
(base) asahinatarou@talol rep3 % /opt/miniconda3/bin/python /Users/taro/DataMining/rep3/level1.py
文書数: 20
総文字数: 29772
```

工夫した点。

まず、テキストファイルがデータフレームに適した状態ではなかった(画像1)ため、URL, date, head, text に分け、df という名の DataFrame とした。その際に URL, date, head は一行ずつ格納することができたが、text は非常に文章量が多いので別途 data という変数を用いて処理した。また、読者に読みやすくするためか改行や空白が所々あったのでそれらの処理も行った。

感想

前処理にかなりの時間が掛かってしまったが、文書内容が非常に面白いので ok.

画像 1

