



< Volver a la semana 3

× Lecciones

este curso: Big Data Essentials: HDFS, MapReduce and Spark

Siguiente

In the following assignments you will write Hadoop Streaming applications. You will be working with Wikipedia articles dump and auxiliary datasets:

Wikipedia dataset location in HDFS: `/data/wiki/en_articles_part`

"Stop words" dataset is located in `/datasets/stop_words_en.txt` file in **local** filesystem.

Format: `article_id <tab> article_text`

Below, you can find starter template (code sample) to preprocess Wikipedia articles for MapReduce. You can use it for programming (LTI) assignments in the following week.

```
1  #!/usr/bin/env python
2
3  import sys
4  import re
5
6  reload(sys)
7  sys.setdefaultencoding('utf-8')
8
9  for line in sys.stdin:
10     try:
11         article_id, text = unicode(line.strip()).split('\t', 1)
12     except ValueError as e:
13         continue
14     text = re.sub("^\W+|\W+$", "", text, flags=re.UNICODE)
15     words = re.split("\W*\s+\W*", text, flags=re.UNICODE)
16
17     # your code goes here
18
```

Remarks:

- English Wikipedia contains a lot of characters from other languages. So, you should parse Unicode correctly.
- you need to remove punctuation marks and transform words to lowercase, to get correct quantities.

If you want to deploy the environment on your own machine, please use [bigdatateam/yarn-notebook](https://github.com/bigdatateam/yarn-notebook) Docker container.

✓ Completado

