In this assignment you will use Spark to compute various statistics for word pairs. At the same time, you will learn some simple techniques of natural language processing.

Dataset location: */data/wiki/en_articles*

Format: *article_id <tab> article_text*

While parsing the articles, do not forget about Unicode (even though this is an English Wikipedia dump, there are many characters from other languages), remove punctuation marks and transform words to lowercase to get the correct quantities. Here is a starting snippet:

```python
1   #! /usr/bin/env python
2
3   from pyspark import SparkConf, SparkContext
4   sc = SparkContext(conf=SparkConf().setAppName("MyApp").setMaster("local[2]"))
5
6   import re
7
8   def parse_article(line):
9       try:
10          article_id, text = unicode(line.rstrip()).split('\t', 1)
11      except ValueError as e:
12          return []
13      text = re.sub("^\W+|\W+$", "", text, flags=re.UNICODE)
14      words = re.split("\W*\s+\W*", text, flags=re.UNICODE)
15      return words
16
17  wiki = sc.textFile("/data/wiki/en_articles_part1/articles-part", 16).map
        (parse_article)
18  result = wiki.take(1)[0]
19
20  for word in result[:50]:
21      print word
```

You can use this code as a starter template. Now proceed to LTI assignments.

If you want to deploy the environment on your own machine, please use bigdatateam/spark-course1 Docker container.

Marcar como completo

👍  👎  🏳