

Dig Data Management Assignment 2

信科四 郑元嘉 1800920541

Task1 - Word Count Task on 新概念英语第二册.txt Based on Hadoop MapReduce

In this task, we count the words in 新概念英语第二册.txt applying Hadoop MapReduce framework.

How I Implement

In this task, we use Python to go through it. Hadoop provides a streaming method to make it possible to run MapReduce tasks for all programming languages.

Here shows the simplistic codes of mapper,

```
1  #!/usr/bin/python3.5
2  import sys
3  from nltk.tokenize import RegexpTokenizer
4
5  tokenizer = RegexpTokenizer(r'\w+')
6
7  for line in sys.stdin:
8      words = tokenizer.tokenize(line.strip())
9
10     for word in words:
11         print ('%s\t%s' % (word, 1))
```

It's quite simple and straightforward. We can design any parsing pattern to build up our key-value pairs. We simply use library 'nltk' in Python to tokenize English sentences.

```
1  #!/usr/bin/python3.5
2  import sys
3
4  current_word, current_count = None, 0
5
6  for line in sys.stdin:
7      word, count = line.strip().split('\t', 1)
8      count = int(count)
9
10     if current_word == word:
11         current_count += count
12     else:
13         if current_count:
14             print('%s\t%s' % (current_word, current_count))
15             current_count, current_word = count, word
16 if current_count:
17     print('%s\t%s' % (current_word, current_count))
```

We need to treat the key-value pairs as string and do customized parsing here, and then we do a simple counting task. Note that the standard input key-value pairs are sorted by key before passed into reducers, it helps simplify the logic of counting.

Result of Word Count

Top 10

1	the	819
2	a	418
3	to	390
4	and	278
5	of	274
6	I	266
7	was	229
8	in	219
9	he	167
10	it	154

It makes sense to note that they're almost stopwords.

Last 10 (English only)

1	added	1
2	actual	1
3	actresses	1
4	actress	1
5	Across	1
6	acquire	1
7	accustomed	1
8	Accurate	1
9	accurate	1
10	account	1