

1 Introduction

PageRank is an algorithm proposed by Google to rank the searching results. It can also be used to measure the importance of the website pages. With billions of website pages in the internet, the efficiency of PAGERANK becomes crucial. Spark is a very popular BigData framework which can distribute data processing tasks across multiple nodes to speed up processing tasks on very large data sets. In this project, we use spark to implement the pagerank algorithm to speed it up on Berkeley-Stanford web graph dataset.

2 Berkeley-Stanford Web Graph

Berkeley-Stanford web graph was crawled in Dec 2002 which has 685230 website pages and 7600595 links. The first four lines of this dataset is the dataset description and each of the rest lines can be split into two columns. The first column defines the start node of the link and the second is the end node. The data cleaning is easy: remove the first four lines.

3 Page Rank Algorithm

PAGERANK algorithm assumes that: 1.important web pages get are more likely to be pointed to by other pages. 2. web pages are important if they are linked by important websites. We use w to represent a website, and define s_w as the rank score of the website, N as the number of the websites, d_w as the number of websites w' pointed by website w . PAGERANK is an iterative algorithm which starts by initializing $s_w = \frac{1}{N}$. Then, we do the following updating until it is converged.

$$s_w = \gamma \cdot \frac{1}{N} + (1 - \gamma) \cdot \sum_{w': w' \rightarrow w} \frac{s_{w'}}{d_{w'}}, \text{ where } \gamma = 0.15 \quad (1)$$

4 Parallel Computation

4.1 Dataset to RDD

To utilize the parallel computation, our first step is to load the dataset to RDD:

```
1 filepath = "/home/zhiyuan/Documents/BigData/FinalProj/web-BerkStan_clean.txt"
2 from_to = sc.textFile(filepath)
```

4.2 Extract the Links

Each line consists of two nodes separated by *tab* and the first node points to the second node. We extract the links by use function *lambda*. After *groupByKey()*, for each started node, we have a list of ended node pointed by the started node. We can think *links* has this data structure: $[(w' : [w_1, w_2...]), ...]$

```
1 links = from_to.map(lambda link: (link.split("\t")[0], link.split("\t")[1])).
  groupByKey()
```

4.3 Initialize s_w

$sw_init = 1.0/685230$ initialize the s_w for each webpage w' . We use variable *ranks* to store all the s_w . We can think the *ranks* has this data structure: $[(w', sw_init)]$.

```
1 ranks = links.map(lambda web_neighbors: (web_neighbors[0], sw_init))
```

4.4 Calculate $\sum_{w'} \frac{s_{w'}}{d_{w'}}$

We can think `links.join(ranks)` has this kind of data structure: $[(w', ([w_1, w_2, \dots], sw))]$. Function `cal_add_item` calculate the $\sum_{w'} \frac{s_{w'}}{d_{w'}}$ for each w in $[(w', ([w_1, w_2, \dots], sw))]$.

```
1 def cal_add_item(urls, sw):
2     dw = len(urls)
3     for url in urls:
4         yield url, sw / dw
5
6 toadd = links.join(ranks).flatMap(
7     lambda from_to_sw: cal_add_item(from_to_sw[1][0], from_to_sw[1][1]))
```

4.5 Update s_w

`Ranks` has this kind of data structure: $[(w, s_w)]$.

```
1 ranks = toadd.reduceByKey(add).mapValues(lambda rank: rank*0.85 + 0.15*sw_init)
```

4.6 Convergence

Remember `Ranks` has this kind of data structure: $[(w, s_w)]$. We can join the old `Ranks` and new `Ranks`, then only take out the values which have this kind of data structure $[(oldvalue, newvalue)]$, finally we calculate the $\sum_w |(s_{wold} - s_{wnew})|$

```
1 e = old_ranks.join(ranks).values().map(lambda old_new: abs(old_new[0] - old_new[1])
2     ).reduce(lambda x, y: x + y)
```

We can set the threshold, so that once the $e \leq threshold$, we can stop the iteration. In this program, I just do 20 iterations and print the e for each iteration.

5 Results

We can save the `Ranks` by:

```
1 ranks.saveAsTextFile("/home/zhiyuan/Documents/BigData/FinalProj/iteration10")
```

At each iteration, we plot the $\sum_w |(s_{wold} - s_{wnew})|$ to check whether pagerank algorithm is converged or not.

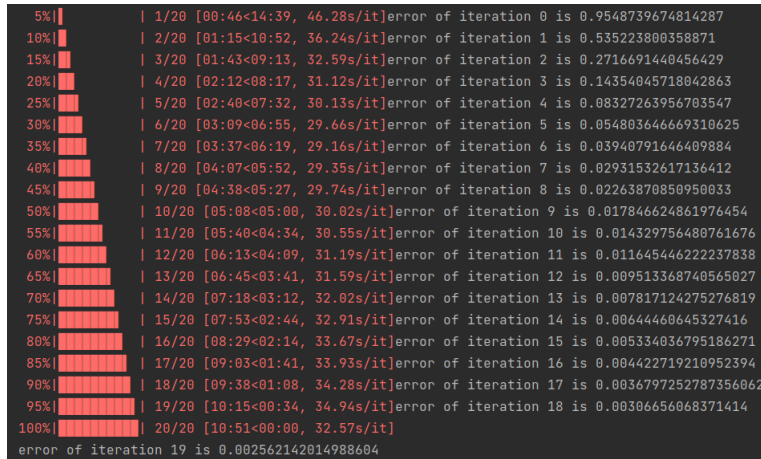


Figure 1: Print the $\sum_w |(s_{wold} - s_{wnew})|$ at each iteration

6 Conclusion

From 2, we can see, the error decreases with the iteration number increasing. From iteration 14, the decreasing almost stops, so, we can think it is converged at iteration 14.

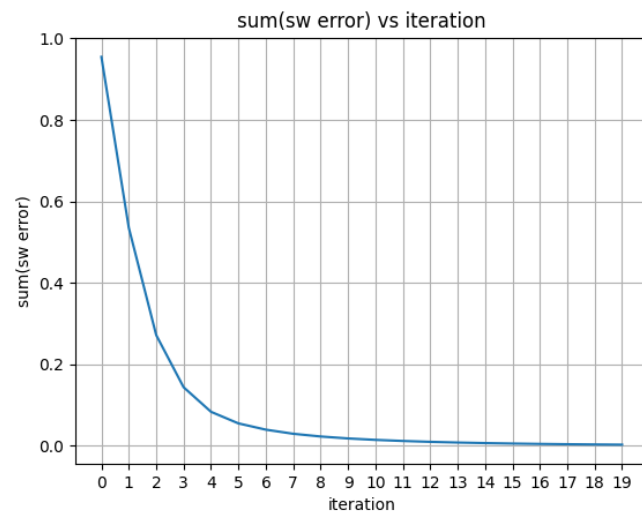


Figure 2: Plot $\sum_w |(s_{wold} - s_{wnew})|$ vs iteration number