HW 3 - Report

Table of Contents

1.	Desig	n Discussion	2
	1.1	Pre-Processing Job	2
		Page Rank Job	
		Top-K Job	
		Data Transferred	
2.	Perfo	rmance Comparison	10
	2.1	Running Time	10
	2.2	Running Time Comparison	.11
	2.3	Top-100 Wikipedia Pages	12
3.		ergence Estimation	

1. Design Discussion

1.1 Pre-Processing Job

I have followed the same per-processing approach as described in the Assignment. I have used given SAX XML Parser. I have modified three things in the given parser:

- 1. Do not include self links
- 2. Removed the duplicate nodes from the adjacency list
- 3. Fixed the " & " issue, some pages was not parsed due to this issue.

I have created my own input file which covers all cases required in assignment and calculated manual page ranks for this input for 10 iterations. And then compared each output of my all programs to the manual calculation and they are the same.

PreProcessingJob emits the nodeld and its adjacentList for all successfully parsed documents. And the job returns the total number of pages from global counter.

Psuedo-Code

```
class Reducer
        Long localNumberOfPages;
        setup()
        {
                localNumberOfPages = 0;
        }
        reduce(nodeld, [adjacencyList1, adjacencyList2,..])
                adjacencyList = EMPTY LIST;
                // get the first non empty adjacencyList if present
                // There can be multiple documents with same nodeld in input
                // so consider the first one only
                for each val in input list do
                {
                        if(val is not EMPTY_LIST)
                        {
                                adjacencyList = val;
                                break;
                        }
                }
                // increment the number of pages
                localNumberOfPages++;
                // emit nodeld, its adjacencyList and its default page rank
                emit((nodeId, adjacencyList), NULL);
        }
       // cleanup increments the global counter by localNumberOfPages count.
        cleanup()
        {
                increment global counter with localNumberOfPages;
        }
}
```

On successful completion of above job we will return the total number of pages from global counters. Total number of pages should be same as number of reduce outputs.

1.2 Page Rank Job

For calculating page rank (computing dangling factor delta), I have used Solution -2 (Merge computation of delta in previous reduce phase) approach as per Learning module 6 - Graph Algorithms. Below are the all three approaches for calculating delta and the reasons why I have choose solution 2.

Solution-1 - Add a separate phase to each iteration to compute delta

During an iteration, first execute a Map Reduce program that computes delta. This is a simple global aggregation job, summing up Page Rank values for all dangling nodes. Then pass the newly computed delta as a parameter to modified Map Reduce program that updates all Page Ranks using new formula with delta. The drawback of this approach is that, in each iteration we are reading all the records in map to calculate the delta in global aggregation job. So here we have an overhead of reading all the inputs for each iteration.

Solution-2: Merge computation of delta in previous reduce phase. (Used Approach)

In this approach, instead of computing delta in a separate job in the beginning of iteration (i+1), it could already be computed at the end of iteration i in reduce call. In this approach we don't need the extra reading of all input records as we need in above approach. Only problem with this approach is, the page rank emitted by the reducer is not the correct one so in last iteration we need to add one more map only job to correct the page ranks calculated by the last iteration. In our program, we have a PageRankCorrectionJob which does the same. So here we can argue that we need one extra map read of all the inputs per entire job, but this can also be eliminated. In our program, We just need TopK page ranks and we can correct the output emitted by the last iteration in TopK job. But this will give the correct page rank of only top K pages and it is possible that we might need other records in future. So for the safer side, I have provided PageRankCorrectionJob to run before TopKJob.

Solution-3 - Order Inversion

The order inversion pattern can be applied to make sure each Reducer receives the old PageRank values of all dangling nodes. Each reducer can then compute delta right before executing any of the normal reduce calls. This approach has the drawback same as what we have seen in other Order Inversion algorithms. We need to have one reduce task only so we cannot have reduce task granularity. We can overcome this issue by emitting multiple copies of same dangling node's page rank such that all reducers receive exact one copy of it. But this creates map output duplication issue. We can overcome this output duplication issue in some extent by in-mapper combining. So this approach with in-mapper combining to dangling factor output in each map task and emitting that combined delta to all reducer (multiple reduce tasks) can be a good solution. But as per our problem needs we required only top k page ranks and in our solution 2 approach we can do that without using PageRankCorrectionJob which makes Solution 2 better than other two approaches.

Pseudo Code:

```
// Map processes the node with id n.
// N stores node n's current PageRank and its adjacency list
map(nid n, node N)
{
       if(isfirstIteration)
              N.pageRank = 1/|V|
       else
       {
               delta = getDeltaFromConfiguration()
              N.pageRank = alpha/|V| + (1-alpha)(delta/|V| + N.pageRank)
       }
       // Pass along the graph structure
       emit(nid n, N)
       if(|N.adjacencyList| != 0)
              //compute contributions to send along outgoing links
              p = N.pageRank / | N.adjacencyList |
              for all nid m in N.adjacencyList do
                      emit(nid m, p)
       }
       else
       {
              emit("Delta", N.pageRank)
       }
}
// In code, I have used In-mapper combining to combine delta and also used Combiner to
// combine the map outputs
// Reduce receives the node object from node m and
// the PageRank contributions for all m's inlinks
reduce(nid m, [p1, p2,...])
{
       s=0
       M=NULL
       for all p in [p1,p2,...] do
              if isNode(p) then
```

```
// The node object was found : recover graph structure M = p
else
// A PageRank contribution from an inlink was
// found : add it to the running sum
s += p1

if (m == "Delta")
    incrementGlobalCounter("Delta", s)
else
    M.pageRank = s
    emit(nid m, node M)
}
// After each iteration we will pass computed delta from global counters to next job's configuration.
```

1.3 Top-K Job

There are two approaches to find the top-K elements. First is to sort the input data and then the K largest records can easily be selected from the sorted file. This is an efficient method for very large K. Second approach is to avoid sorting if we have smaller value of K. The idea is to scan the input only once and use in-mapper combining to keep track of the top-K records in each map task. A single reduce call receives these local top-k lists and merges them into final result.

In our assignment K=100 which is a very small number, So I have used the second approach and not sorting as we do not need to transfer all input records from map to reduce, read again in reduce and write from reduce to file.

Pseudo-Code:

Pseudo code is same as the pseudo code given in Module -5: Basic Algorithms.

[2.13: Top-K records – Second Approach – find the local TopK in each Map task and merge them into global TopK in one reduce task.]

Addition to the given code, I have implemented MyTreeMap to handle the different pages with same page rank issue.

1.4 Data Transferred

Below table shows the amount of data transferred from Mappers to Reducers, and from Reducers to S3 (on AWS, job has configured to write on S3 instead of HDFS) in each phase of the entire program execution and in each iteration of page rank job. Data from syslog file of AWS EMR run of Wikipedia-full-html input execution on 11 machines (1 master and 10 workers).

Phase - Iteration	S3: Number of bytes read (S3 to Mappers)	Map output bytes	Reduce shuffle bytes (Combiner/Mapper to Reducer)	S3: Number of bytes written (Reducer to S3)
Pre-Pro.	7091101822	2088430787	1008848757	1057795115
1	1057795115	3210315684	1217336429	1180704724
2	1180704724	3214348041	1317781918	1182955026
3	1182955026	3212780522	1317996064	1181970571
4	1181970571	3213168630	1318206121	1183014094
5	1183014094	3213382361	1318341266	1183023055
6	1183023055	3213765772	1318386360	1183036262
7	1183036262	3212932819	1318206906	1181992025
8	1181992025	3212833986	1318236724	1181995559
9	1181995559	3214408272	1318455864	1183051088
10	1183051088	3213143682	1318350598	1181993998
Correction	1181993998	Map Only Job, Mapper to S3: 142752032 byte		752032 bytes
ТорК	142752032	65692	59740	3211

Below table shows the number of records transferred between mapper, combiner and reducer. Data from syslog file of AWS EMR run of Wikipedia-full-html input execution on 11 machines (1 master and 10 workers).

Iteration	Map input records	Map output records	Combine input records	Combine output records	Reduce input records	Reduce output records
Pre-Pro.	7012253	54626996	54626996	19434887	19434887	3178227
1	3178227	54718394	54718394	19556653	19556653	3178227
2	3178227	54718394	54718394	19562914	19562914	3178227
3	3178227	54718394	54718394	19559839	19559839	3178227
4	3178227	54718394	54718394	19560635	19560635	3178227
5	3178227	54718394	54718394	19561348	19561348	3178227
6	3178227	54718394	54718394	19561719	19561719	3178227
7	3178227	54718394	54718394	19560157	19560157	3178227
8	3178227	54718394	54718394	19560043	19560043	3178227
9	3178227	54718394	54718394	19562772	19562772	3178227
10	3178227	54718394	54718394	19560535	19560535	3178227
Correction 3178227 3178227			Map Onl	y Job		
ТорК	3178227	1900	0	0	1900	100

In pre-processing, the input document size is very high. But in pre-processing, we discard most of the data from document and end up to only graph representation (nodeld : Adjacency List). So output of the reducer to S3 is comparatively smaller.

In first iteration, in input file (S3 to mapper), we have only nodeld and adjacencyList, but after processing the first iteration of page rank, reducers emit the page rank of each page along with its nodeld and adjacencyList. This increases the number of bytes written to S3 from reducers for first iteration.

For all other iteration, there is no bigger difference in amount of data transferred in each step, as we are processing & passing the kind of same data and just modifying the page rank values. In each iteration, amount of data transferred from Mapper to Reducer increases because we are passing the page rank contribution to each adjacent node along with node structure entries. Though Combiner helped a lot to reduce the reduce shuffle bytes from original Map output bytes. Reducer combines the map outputs and writes it to S3 so the number of bytes reduces again.

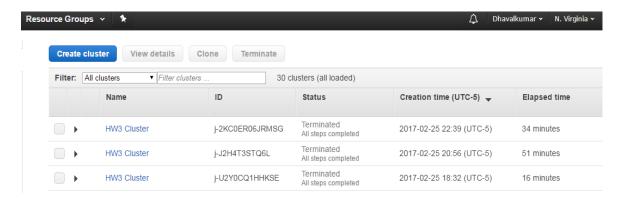
Last Correction job writes only Nodeld and PageRank values amount of data transferred to S3 from reducer drastically reduced same for TopK(only K=100 records written) job as well.

2. Performance Comparison

2.1 Running Time

I have executed program on AWS EMR for both the inputs. I have used alpha = 0.15 for my page rank calculations (Source: Wikipedia page of page rank algorithm). Though the value of alpha is not hardcoded and you can configure in Makefile.

Below is the screenshot of AWS EMR successful run:



1st row: Wikipedia-full-html input run on 11 m4.large machines (1 master and 10 workers)

2nd row: Wikipedia-full-html input run on 6 m4.large machines (1 master and 5 workers)

Below is the running time of full Wikipedia input on both configurations:

(Time in Seconds)

Job	6 m4.large Machine	11 m4.large Machines
Pre-Processing Time	1021	631
Time to Run ten Iterations Of Page Rank	1532	947
Time to Find The Top -100 Pages	40	29

^{3&}lt;sup>rd</sup> row: Wikipedia-simple-html input run on 6 m4.large machines (1 master and 5 workers)

2.2 Running Time Comparison

As we see in running times, all the computation phases shown good speed up by increasing machines as expected. With higher number of machines, our work load is divided among them and due to parallel work overall time is reduced.

If we compare the speed up of all computational phases then first two phases (PreProcessing and PageRank) has shown higher speed up then the last phase (TopK). The running time decreased to 38% in first to phases while it reduced 27% in last phase. I have also expected the same as in first two phases our MapReduce framework utilizes all the machines whereas in TopK job we are creating only on reduce task so we cannot get advantage of increase in machines in reduce work of TopK job.

Number of reduce task depends on the available nodes to MapReduce framework. As per hadoop documentation "The right number of reduces seems to be 0.95 or 1.75 multiplied by (<no. of nodes> * <no. of maximum containers per node>)." If we have higher number of machines then it will create more reduce task and we will get good parallelism which will reduce the time taken and increase the performance.

In our case for first two phases, for 11 machines, total number of reduce task is 20 whereas for 6 machines they are reduced to 10. So in first two phases, higher machine gave better performance as the load is distributed among them nicely. For last phase, number of reduce tasks are independent of machines as we are setting it manually to 1. So here we are not getting advantage of increase in machines in reduce phase and it shows lesser speed up then above two phases. Adding to

2.3 Top-100 Wikipedia Pages

Below is the top-100 Wikipedia pages with highest page rank along with their page values sorted from highest to lowest for both the simple and full datasets. I find them reasonable as I guess there will be higher in links to the pages like United_States_09d4 from the other pages which result into higher page rank of this page. And we can see other countries as well in the list this means, the other wiki pages will have links in their web page to the related country page and which will boost the page rank of the country pages.

No	Wikipedia-full-html-output-11-machines		Wikipedia-simple-html-o	utput-6-machines
•	Page Name	Page Rank	Page Name	Page Rank
		0.00262288934616		0.00518900900027
1	United_States_09d4	4360	United_States_09d4	2900
		0.00122849783940		0.00480676647470
2	2006	2600	Wikimedia_Commons_7b57	8800
		0.00120313676095		0.00394028468771
3	United_Kingdom_5ad7	6090	Country	2630
		0.00098207155818		0.00275248143611
4	Biography	9593	England	0630
		0.00091705986394		0.00268780962344
5	2005	1628	Water	6530
		0.00088020713412		0.00255408756514
6	England	7707	Animal	9160
		0.00085590507934		0.00251082408078
7	Canada	7993	City	2450
	Geo-	0.00077172622988		0.00235864709361
8	graphic_coordinate_system	9266	United_Kingdom_5ad7	2220
		0.00072502418722		0.00235040169771
9	France	5341	Germany	1510
		0.00071989680727		0.00232473485995
10	2004	9462	Earth	4630
		0.00068047669529		0.00232360794714
11	Australia	2216	France	2090
		0.00065434528014		0.00203809703716
12	Germany	1820	Europe	7730
		0.00058738697782		0.00175388421427
13	2003	8930	Wiktionary	6070
		0.00058341977400		0.00174967712175
14	India	5353	English_language	4420
		0.00058285604737		0.00173234465210
15	Japan	3971	Government	3290
	Inter-	0.00053350696600		0.00171684048471
16	net_Movie_Database_7ea7	3842	Computer	3360
17	Europe	0.00050926763789	India	0.00171317091838

		9670		4910
		0.00049145956758		0.00166738369802
18	Record_label	4876	Money	2790
		0.00048701215830		0.00155169056853
19	2001	3275	Japan	5440
		0.00048286324681		0.00152355950935
20	2002	9372	Plant	9920
		0.00047805043674		0.00150743309049
21	World_War_II_d045	7285	Italy	7990
		0.00047034299613		0.00148140734345
22	Population_density	8859	Canada	2890
22	Music gonro	0.00046721175179	Casia	0.00147112369222
23	Music_genre	0.00046466103240	Spain	3520
24	2000	1972	Food	0.00142468684896 7640
	2000	0.00044579274603	FOOU	0.00141209700626
25	Italy	5966	Human	9630
	reary	0.00043620978702	Haman	0.00139671506127
26	Wiktionary	9988	China	2910
	,	0.00043529472397		0.00138224852505
27	Wikimedia_Commons_7b57	9010	People	5760
		0.00043480265906	·	0.00132985424075
28	London	0111	Australia	0500
		0.00041850352295		0.00128443617113
29	English_language	4299	Asia	6100
		0.00040593698578		0.00127426842125
30	1999	7724	Capital_(city)	1940
		0.00036295379192		0.00126499722576
31	Spain	9990	Television	0360
		0.00035631063673		0.00126021008117
32	1998	4167	Sun	8010
22	Bussia	0.00034390662497	Number	0.00124323622892
33	Russia	8832 0.00033728202493	Number	0.00124037568145
34	1997	3117	State	4610
	1337	0.00033629712380	Juic	0.00123521166722
35	Television	0.00033023712380	Sound	1920
		0.00033462877684		0.00123254317535
36	New York City 1428	3839	Science	9430
		0.00032614642134		0.00123105663929
37	Football_(soccer)	4763	Mathematics	5570
		0.00032362786609		0.00119230462374
38	1996	1999	Metal	9440
		0.00032355337188		0.00117709258351
39	Census	6843	Year	0610
		0.00032218915558		0.00117335731376
40	Scotland	6006	2004	8480

		0.00031015464104		0.00115016588485
41	1995	7029	Language	7740
		0.00030864301289		0.00114618177921
42	China	6061	Russia	2580
		0.00030432048264		0.00112333028098
43	Population	0489	Wikipedia	8190
		0.00030405610072		0.00109856669996
44	Square_mile	0908	Religion	6040
	-	0.00030401196802	_	0.00109653914178
45	Scientific_classification	8030	19th_century	0080
		0.00030166740429		0.00108743132321
46	California	4186	Music	4420
		0.00029069115857		0.00105480073500
47	1994	5429	Scotland	6320
		0.00028762080340		0.00105370498325
48	Sweden	1212	20th_century	8870
		0.00028741610678		0.00104922273293
49	Public_domain	1891	Greece	4630
		0.00028626913927		0.00102986061318
50	Film	5013	Latin	7450
		0.00028411019015		0.00102735544285
51	Record_producer	9314	London	1320
		0.00028310205424		0.00100435725665
52	New_Zealand_2311	5546	Greek_language	0290
		0.00027888263555		0.00099901181037
53	New_York_3da4	7677	Energy	9402
		0.00027667367831		0.00098635084799
54	Netherlands	8118	World	7666
		0.00027581329226		0.00097590586513
55	Marriage	8387	Centuries	6575
		0.00027482489872		0.00094520396521
56	1993	0027	Culture	1297
l	United_States_Census_Bureau	0.00027466710694		0.00093646960342
57	_2c85	8521	History	5431
	1001	0.00027189525612	11. 1.1	0.00091452309680
58	1991	9761	Liquid	0025
	1000	0.00026832611879	Ni a tila a vil a va al a	0.00090572450764
59	1990	8734	Netherlands	8975
	1003	0.00026636924637	Dlanat	0.00090493226223
60	1992	7862	Planet	9007
61	Dolitician	0.00026489490763	Light	0.00090167635268
61	Politician	6039	Light	6388
62	Album	0.00026056445391	Society	0.00090149206214
62	Album	8753	Society	5202
62	Latin	0.00026045636147	Atom	0.00089002264065
63	Latin	8826	Atom	2958
64	Actor	0.00025833956084	Wikime-	0.00088844007077

		6043	dia_Foundation_83d9	6097
		0.00025810638770		0.00088838361057
65	Ireland	3505	Scientist	3494
		0.00025564272454		0.00088768848602
66	Per_capita_income	9733	Image	1999
		0.00025186026160		0.00088629080559
67	Studio_album	0276	Law	8419
		0.00025116500863		0.00087884516145
68	Poverty_line	4912	Geography	4898
	2	0.00024950659722		0.00087857429428
69	Km²	5139	List_of_decades	3701
		0.00024600272252	Uni-	0.000004.00450036
70	1000	0.00024689373252	form_Resource_Locator_1b	0.00086188450636
70	1989	3362 0.00024092192871	4e	3216
71	Norway	2262	Africa	0.00086056996715 2458
	NOI Way	0.00023901785842	Airica	0.00084488636788
72	Website	5252	Turkey	9011
<u> </u>	· · · · · · · · · · · · · · · · · · ·	0.00023532199065	rancy	0.00083047948823
73	1980	3687	Inhabitant	2369
		0.00022937863976		0.00082304881404
74	Animal	8423	Capital_city	3760
		0.00022920870296	· = /	0.00082151559551
75	Area	2197	Plural	0230
		0.00022703304683		0.00081372300166
76	1986	4427	Electricity	6497
		0.00022623653898		0.00079723790431
77	Personal_name	7703	Poland	5345
		0.00022611544290		0.00079712389257
78	Poland	5689	Building	2039
		0.00022568426663		0.00079465406062
79	Brazil	7636	Car	3909
90	1005	0.00022402906904	Swadan	0.00079171255623 4121
80	1985	2517 0.00022330540325	Sweden	0.00079148847053
81	1987	8049	Book	1959
- 51	1507	0.00022175338663	BOOK	0.00078693289643
82	1983	6217	Biology	1402
		0.00022109653273	01	0.00077081729454
83	1982	6380	War	8049
		0.00021938452665		0.00076816079591
84	French_language	4706	Chemical_element	9691
		0.00021934801195		0.00076093572189
85	1981	2457	God	1388
		0.00021932859347		0.00075628686441
86	1979	5849	North_America_e7c4	6695
87	1984	0.00021879019421	September_7	0.00075477818126

		2315		4075
		0.00021869239369		0.00074629735006
88	World_War_I_9429	5080	Website	0417
		0.00021857418021		0.00074266715264
89	1988	8445	Nation	0605
		0.00021801968011		0.00073971037875
90	Paris	5686	Politics	8902
		0.00021797486466		0.00073329001722
91	1974	4915	2006	5940
		0.00021567359642		0.00073223711129
92	Mexico	7163	Fish	0969
		0.00021185635773		0.00073087111762
93	19th_century	6199	Species	9328
		0.00021132429567		0.00072167441359
94	1970	9975	Mammal	4913
		0.00021086426991		0.00071780902030
95	January_1	0462	Island	3589
		0.00021070868243		0.00071710705966
96	USA_f75d	2733	Portugal	0586
		0.00020860183512		0.00071555153665
97	1975	0045	Gas	3913
		0.00020846744997		0.00071157775130
98	1976	3305	River	0913
		0.00020779805256		0.00070610750743
99	Africa	2159	Switzerland	8502
10		0.00020736244197		0.00070203049315
0	South_Africa_1287	3637	World_War_II_d045	8166

3. Convergence Estimation

I have coded to measure the convergence achieved after each iteration of page rank calculation. In each iteration, program calculates the absolute difference in old and new page rank of each node and sums it up. And at the end of iteration, the sum is divided by the number of total pages to get the average difference per node and logs it to log file.

Below is the convergence output from logs of AWS EMR execution of Wikipedia-full-html input with both configurations:

Iteration	6 m4.large Machine	11 m4.large Machines
1	0.000000184996236977	0.000000184996609594
2	0.000000076656124798	0.000000076656300012
3	0.000000030632383622	0.000000030632328128
4	0.000000013365304956	0.000000013365282971
5	0.000000006365162230	0.000000006365152615
6	0.000000003233759976	0.000000003233754581
7	0.000000001721090299	0.000000001721086772
8	0.000000000951050433	0.000000000951048336
9	0.000000000543148803	0.000000000543147593

(Note: It doesn't calculate for last iteration, as last page rank is corrected in PageRankCorrectionJob)

As we can see in above diagram, average difference in page rank in each iteration reduces which means page ranks are converging.

If we need the convergence up to specific decimal point then we can count the number of nodes whose page rank change is higher than specified margin in each iteration. This count will decrease in each iteration and when it reaches to zero, we will stop calculating the page rank.