

Course: CS6200

Homework 3

Name: Rushabh Shah

Github Repository: <https://github.ccs.neu.edu/rushabh0812/HW3>

Pseudo-Code for Pre-Processing:

```
method map(Key k, value v):
    name = Parse v to obtain page-name
    adjacencyList = Parse v to obtain adjacencyList
    emit(name, adjacencyList)

    for each record r in adjacencyList:
        emit(r, null);
end

method reduce(pageName p, adjacencyList[ l1, l2, l3, ...]):
    pageCountCounter = newGlobalCounter
    if( isInLinksPresent() or isOutLinksresent()):
        pageCountCounter = pageCountCounter + 1

    if(adjacencyList is null):
        emit(p, new adjacencyList[])
        return
    emit(p, adjacencyList)
```

Pseudo-Code for Page Rank:

```
method map(Key k, value v):
    alpha = Initialize alpha component of page rank
    pageCount = Extract pageCount from counters
    delta = get delta from counters
    emit(k,v)

    newPage = v.pageRank + (1-alpha)*delta/pageCount
    for each entry e in v.adjacencyList:
        emit(e, newPageRank/sizeof(adjacencyList))

    if(v.adjacencyList is empty):
        emit(dummy, newPageRank)

method reduce(Node n, List[dummy, c1, c2, c3]):
    if(n is dummy):
        for each contribution 'c' in list:
            totalContributions += c
        update delta global counter with totalContributions
        return

    Node n1 = Initialize new node
    for each entry e in List:
        if(e in adjacencyList):
            n1.adjacencyList = e
        else:
            totalContributions += e

    n1.pageRank = totalContributions
    emit(n,n1)
```

Pseudo Code for TopK records:

```
Class Mapper{
  localTopK

  setup(){
    initialize localTopK
  }

  map(... , x){
    if (x is in localTopK)
      // Adding x also evicts the now (k+1)-st record from localTopK
      localTopK.add(x)
    }

  cleanup(){
    for each x in localTopK
      emit(dummy, x)
    }

  reduce(dummy, [x1, x2, x3....]){
    initialize globalTopK

    for each record x in input list
      if (x in gloablTopK)
        // Adding x also evicts the now (k+1)-st record from globalTopK
        globalTopK.add(x)

    for each x in localTopK
      emit(NULL, x)
  }
}
```

Performance Comparison

Time For	6 m4.large machines	11 m4.large machines
Pre-Processing	1562206	801178
10 iterations of Page Rank	1184316	894631
Top-100 pages	52778	51820

Report for both configurations the amount of data transferred from Mappers to Reducers, and from Reducers to HDFS, separately for each iteration of the PageRank computation. Does it change with the cluster size? Does it change over time?

6m4.large syslog screenshot:

```
      Total megabyte-milliseconds taken by all reduce tasks
Map-Reduce Framework
  Map input records=2292111
  Map output records=68942930
  Map output bytes=3601269553
  Map output materialized bytes=1121059069
  Input split bytes=2628
  Combine input records=85621388
  Combine output records=31648992
  Reduce input groups=2292112
  Reduce shuffle bytes=1121059069
  Reduce input records=14970534
  Reduce output records=2292111
  Spilled Records=46619526
  Shuffled Maps =162
  Failed Shuffles=0
  Merged Map outputs=162
  GC time elapsed (ms)=8743
  CPU time spent (ms)=584790
  Physical memory (bytes) snapshot=16345899008
  Virtual memory (bytes) snapshot=100520697856
  Total committed heap usage (bytes)=15251013632
```

11m4.large syslog screenshot:

```
Total megabyte-milliseconds taken by all map tasks=22
Total megabyte-milliseconds taken by all reduce tasks=22
Map-Reduce Framework
  Map input records=7012253
  Map output records=68931961
  Map output bytes=3613989271
  Map output materialized bytes=1271392380
  Input split bytes=11554
  Combine input records=0
  Combine output records=0
  Reduce input groups=3150469
  Reduce shuffle bytes=1271392380
  Reduce input records=68931961
  Reduce output records=1635970
  Spilled Records=143046037
  Shuffled Maps =2014
  Failed Shuffles=0
  Merged Map outputs=2014
  GC time elapsed (ms)=122021
  CPU time spent (ms)=10696450
  Physical memory (bytes) snapshot=95290146816
  Virtual memory (bytes) snapshot=436883947520
  Total committed heap usage (bytes)=83450396672
```

Metric	6m4.large	11 m4.large
Map Input Records	2292111	7012253
Map Output Records	68942930	68931961
Reduce Input Records	14970534	68931961
Reduce Output Records	2292111	1635970
HDFS: Number of bytes read	2628	11554
HDFS: Number of bytes written	46619526	143046037

Which of the computation phases showed a good speedup? If a phase seems to show fairly poor speedup, briefly discuss possible reasons—make sure you provide concrete evidence, e.g., numbers from the log file or analytical arguments based on the algorithm’s properties.

The speedup for different phases are as follows:

1. Pre-processing: time on 6 machines/11machines = 1562206/801178= 1.94
2. Page Rank: time on 6 machines/11 machines = 1184316/894631 = 1.323
3. Top100: time on 6machines/11 machines = 52778/51820 = 1.018

Top 100 Wikipedia Pages:

Simple Dataset

0.005453741115272236,United_States_09d4
0.004104536615524893,Wikimedia_Commons_7b57
0.003404509212509808,Country
0.0022948577996755944,England
0.002271737607979354,Europe
0.002268150805290908,United_Kingdom_5ad7
0.0022629779913135284,Water
0.002171982659197757,Germany
0.0021688915379444646,France
0.002138656181574242,Earth
0.0021307553864373533,Animal
0.0020324385973417907,City
0.0018104100982807851,Week
0.0016861675135594327,Asia
0.001662584120043279,Sunday
0.0016376791725900256,Monday
0.0016211380742672439,Wednesday
0.0016067369109369653,Wiktionary
0.0016012250605147568,Money
0.0015812332480585113,Friday
0.0015751561805757337,Plant
0.001563585832571562,Saturday
0.0015432469741492245,Thursday
0.0015322597127829711,Tuesday
0.0015226660733272503,Computer
0.0015206333183474597,English_language
0.0014886861646351845,Government
0.0014854208861482282,Italy
0.0014803779262863906,India
0.0013989641624862293,Number
0.0013433209043226485,Spain
0.001324716409963496,Day
0.0013067866792937732,Japan
0.0013012872224230385,Canada
0.001259706767201987,People
0.0012330199856291776,Human
0.001199215924337052,Wikimedia_Foundation_83d9

0.00118555669038321,China
0.0011835802084830167,Australia
0.0011806317748449646,Energy
0.0011492294294854126,Sun
0.0011470167114563676,index
0.0011432737157235494,Food
0.0011309142327184433,Science
0.001110240554157543,Mathematics
0.0010524948126738657,Television
0.0010341124010405758,Capital_(city)
0.0010272504631897096,Russia
0.0010185780992265463,State
9.962916269389546E-4,Year
9.890574283683957E-4,Music
9.631823330129022E-4,Greece
9.585194243681958E-4,Language
9.565535753398897E-4,Scotland
9.427108993083259E-4,Metal
9.347069737086805E-4,Wikipedia
9.288391777659076E-4,Greek_language
9.197337782593858E-4,Planet
9.084557885467258E-4,2004
8.908705039989364E-4,Sound
8.876963803266156E-4,Religion
8.754673511143459E-4,London
8.66361357102545E-4,Africa
8.29628865930656E-4,Geography
8.275228053720926E-4,Law
8.260124904412899E-4,20th_century
8.239106344310348E-4,Liquid
8.095556887868571E-4,19th_century
8.089858293821437E-4,World
7.995000625663615E-4,Society
7.985908636151757E-4,Scientist
7.780969934950783E-4,Atom
7.660343358245363E-4,Latin
7.660326403598595E-4,History
7.599655695752877E-4,Light
7.555621797389816E-4,Sweden
7.551883781897179E-4,Poland
7.545850991653409E-4,War

7.438619807426295E-4,Culture
7.415970353280629E-4,Netherlands
7.318881731430144E-4,Building
7.186343919935827E-4,Turkey
7.165891228724882E-4,Plural
7.149167207311541E-4,God
7.088099495257791E-4,Information
6.947877783286079E-4,Centuries
6.942208027589603E-4,Chemical_element
6.890653944902166E-4,Portugal
6.78460591726409E-4,Denmark
6.704297038277953E-4,Austria
6.689640326250872E-4,Cyprus
6.662490730894218E-4,Capital_city
6.643165056288547E-4,Ocean
6.562123785555036E-4,North_America_e7c4
6.551499039158885E-4,Inhabitant
6.544220159449451E-4,Moon
6.531322745787359E-4,Species
6.518098214751361E-4,Disease
6.508095822929614E-4,Biology
6.489588261018193E-4,Book

Full Data set

0.0022996176325430014,United_States_09d4
0.0020706757462888286,2006
0.0010962546969893415,United_Kingdom_5ad7
9.503407064868011E-4,2005
7.237970401168999E-4,Biography
7.126458390750062E-4,Canada
7.096585486220868E-4,England
7.083586581322475E-4,France
6.606620815052857E-4,2004
6.046346052643152E-4,Germany
5.834674613244876E-4,Australia
5.612867583671742E-4,Geographic_coordinate_system
5.319404301287987E-4,2003
5.14646806654876E-4,India
5.048125582945026E-4,Japan

4.3327951393979037E-4,Italy
4.2664576555855884E-4,2001
4.215682707295947E-4,2002
4.134895474056176E-4,Internet_Movie_Database_7ea7
4.044529101638422E-4,Europe
3.992591069655945E-4,2000
3.8651564915375406E-4,World_War_II_d045
3.739017857762557E-4,London
3.5619955737753553E-4,English_language
3.5372697544887793E-4,Population_density
3.5325802141764664E-4,1999
3.531758918660774E-4,Spain
3.4602790020054784E-4,Record_label
3.3340079217880045E-4,Russia
3.305487328771001E-4,Race_(United_States_Census)_a07d
3.273638012116936E-4,Wiktionary
3.125647692172924E-4,Wikimedia_Commons_7b57
3.0503917113466515E-4,1998
2.9175964267938003E-4,1997
2.910829799217043E-4,Music_genre
2.885126775582817E-4,New_York_City_1428
2.876190166731202E-4,Scotland
2.73639731768505E-4,1996
2.689199997489037E-4,Sweden
2.684842086577E-4,Football_(soccer)
2.6684591485695894E-4,Television
2.5885870044762904E-4,Square_mile
2.578173018813917E-4,1995
2.565196348124615E-4,Census
2.5508563707107737E-4,California
2.5399972869177763E-4,China
2.500815847686322E-4,Netherlands
2.468565178505421E-4,New_Zealand_2311
2.463332085689414E-4,1994
2.3524152083383853E-4,1991
2.3256913386128747E-4,1993
2.3154776449134216E-4,1990
2.2986050027739592E-4,Public_domain
2.2983156791672952E-4,New_York_3da4
2.2299528678178119E-4,1992
2.2066790492250517E-4,United_States_Census_Bureau_2c85

2.190506978619827E-4,Film
2.1646293103969893E-4,Ireland
2.1628218381199082E-4,Norway
2.1580939558354207E-4,Actor
2.1512975985160752E-4,Scientific_classification
2.099371604986171E-4,Population
2.0952298720396085E-4,1989
2.076974446111385E-4,January_1
2.0523016494857353E-4,Latin
2.051695153613435E-4,1980
2.0213895815303896E-4,Brazil
2.018269394388249E-4,Mexico
2.0095768651471615E-4,Marriage
1.995367176470662E-4,1986
1.966298132821705E-4,French_language
1.9465051955724517E-4,1979
1.9432337003785546E-4,1985
1.9377404716056828E-4,1982
1.937191174994288E-4,1981
1.9249523714410635E-4,1974
1.9227148925034625E-4,Poland
1.9221462292518752E-4,Politician
1.9011369775103333E-4,South_Africa_1287
1.9009030541808163E-4,Switzerland
1.89768029544612E-4,1984
1.896208630127703E-4,1983
1.8954205245400683E-4,1987
1.8926253265456644E-4,Per_capita_income
1.8769605494478594E-4,1970
1.852752611772E-4,1988
1.8527158953722436E-4,1976
1.8509163911032442E-4,Album
1.8297151685262863E-4,Record_producer
1.8297085586158377E-4,1975
1.810133251335694E-4,1969
1.8055320354770323E-4,Paris
1.798253208501527E-4,Greece
1.7978309989739278E-4,Km²
1.797518484831652E-4,1945
1.7953735757391068E-4,1972
1.7842203913378776E-4,Soviet_Union_ad1f

1.7771562004284996E-4,1977

1.7672759870600738E-4,1978

1.7495996223508656E-4,1973