

Welcome to

Big Data & Hadoop

Session

Session 3 - Adv. Map Reduce



+91-9538998962

sandeep@knowbigdata.com

WELCOME - KNOWBIGDATA

- Interact - Ask Questions
- Real Life Project
- Lifetime access of content
- Quizzes & Certification Test
- Class Recording
- 10 x (3hr class)
- Cluster Access
- Socio-Pro Visibility
- 24x7 support
- Mock Interviews

ABOUT ME

2014	KnowBigData	Founded
2014	Amazon	Built High Throughput Systems for Amazon.com site using in-house NoSql.
2012		
2012	InMobi	Built Recommender after churning 200 TB
2011	tBits Global	Founded tBits Global Built an enterprise grade Document Management System
2006	D.E.Shaw	Built the big data systems before the term was coined
2002	IIT Roorkee	Finished B.Tech somehow.
2002		



COURSE CONTENT

I	Understanding BigData, Hadoop Architecture
II	Environment Overview, MapReduce Basics
➔ III	Adv MapReduce & Testing
IV	Analytics using Pig
V	Analytics using Hive
VI	NoSQL, HBASE
VII	Oozie, Mahout,
VIII	Zookeeper, Apache Storm
IX	Apache Flume, Apache Spark
X	YARN, Big Data Sets & Project Assignment

TODAY'S CLASS

- **Streaming Job**
- **Description**
- **Visualization**
- **Hands ON**
- **Discussion on Problems**
- **Limitations**
- **Testing**
- **Assignments**

A Hadoop Library which makes it possible to use *any binary* as mapper or reducer

Example

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming.jar  
-input sgiri/wordcount/input/  
-output mylogin/output/  
-mapper 'sed "s/ /\n/g"  
-reducer "/usr/bin/uniq -c"
```


sa re sa ga



replace space with new line
sed 's/ /\n/g'



sa
re
sa
ga



ga



uniq -c



ga 1

re



uniq -c



re 1

sa
sa



uniq -c



sa 2

ga 1
re 1
sa 2



Ship a script

```
#mycmd.sh - clean up further
```

```
#!/bin/bash
```

```
sed "s/ /\n/g"|sed "s/[^a-zA-Z0-9]//g"|tr "A-Z" "a-z"
```

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming.jar
```

```
-input sgiri/wordcount/input/
```

```
-output sgiri/wordcount/output7/
```

```
-mapper ./mycmd.sh
```

```
-reducer "/usr/bin/uniq -c"
```

```
-file mycmd.sh
```

STREAMING JOB - HANDS-ON



I. Frequencies of letters [a-z] - Do you need Map/Reduce?

I. Frequencies of letters [a-z] - Do you need Map/Reduce?

Without MR Approach:

- Create an integer array A of 26 size.
- Scan the text character by character
- Increase $A[0]$ for 'a' and $A[25]$ for 'z'
 - and others in between

I. Frequencies of letters [a-z] - Do you need Map/Reduce?

Without MR Approach:

- Create an integer array A of 26 size.
- Scan the text, character by character
- Increase A[0] for 'a' and a[25] for 'z'
 - and others in between

Problems?

1. Although memory/RAM will not be a limitation
2. The Network or Disk IO will be bottle neck
3. The CPU will be bottle

I. Frequencies of letters [a-z] - Do you need Map/Reduce?

MR Approach:

Mapper

- Does the same thing as previous example
- Prints the array at the end

Reducer

- Sums up the values at the end

Mapper1	Mapper2
a 234	a 12
b 23028409	b 122
c 328782	c 90
d 37637	d 22
...26 line	...26 line

a [234, 12 ...]
b [23028409, 122 ...]

2. Find anagrams in a huge text. An anagram is basically a different arrangement of letters in a word.

Input:

“the cat act in tic tac toe.”

Output:

cat, tac, act

MAP / REDUCE

Problem 2: Anagram

*“the cat act in tic
tac toe.”*

Mapper

Forms Key
by Sorting
Chars and value is
actual word

eht	<i>the</i>
act	<i>cat</i>
act	<i>act</i>
in	<i>in</i>
cit	<i>tic</i>
act	<i>tac</i>
eot	<i>toe</i>

act	act
act	cat
act	tac
cit	tic
eht	the
eot	toe
in	in

act	act, cat, tac
cit	tic
eht	the
eot	toe
in	in

Reducer

Simply
prints distinct
values for a key

act, cat, tac
in
the
tic
toe

3a. A file contains the DNA sequence of people. Find all the people who have same DNAs.

Input:

“User1 ACGT”

“User2 TGCA”

“User3 ACG”

“User4 ACGT”

“User5 ACG”

“User6 AGCT”

Output:

User1, User4

User2

User3, User 5

User6

```
select dna, concat(users) from mytable group by dna
```


3b. A file contains the DNA sequence of people. Find all the people who have same or mirror image of DNAs.

Input:

“User1 ACGT”

“User2 TGCA”

“User3 ACG”

“User4 ACGT”

“User5 ACG”

“User6 AGCT”

Output:

User1, User2, User4

User3, User 5

User6

MAP / REDUCE

Problem 3: DNA

User1 ACGT
User2 TGCA
User3 ACG
User4 ACGT
User5 ACG
User6 AGCT
User7 TCGA

Mapper

Key:
Smaller (DNA ,
Reverse)
Value: User Id

ACGT USER1
ACGT User2
ACG User3
ACGT User4
ACG User5
AGCT User6
AGCT User7

$\min(\text{ACGT}, \text{TCGA})$
=ACGT

ACG User3
ACG User5
ACGT USER1
ACGT User2
ACGT User4
AGCT User6
AGCT User7

ACG User3,User5
ACGT USER1,User2, User 4
AGCT User6,User7

Reducer

Simply
prints Users

User3,User5
USER1,User2, User 4
User6,User7

4. In an unusual democracy, everyone is not equal. The vote count is a function of worth of the voter. Though everyone is voting for each other.

As example, if A with a worth of 5 and B with a worth of 1 are voting for C, the vote count of C would be 6.

You are given a list of people with their value of vote. You are also given another list describing who voted for who all.

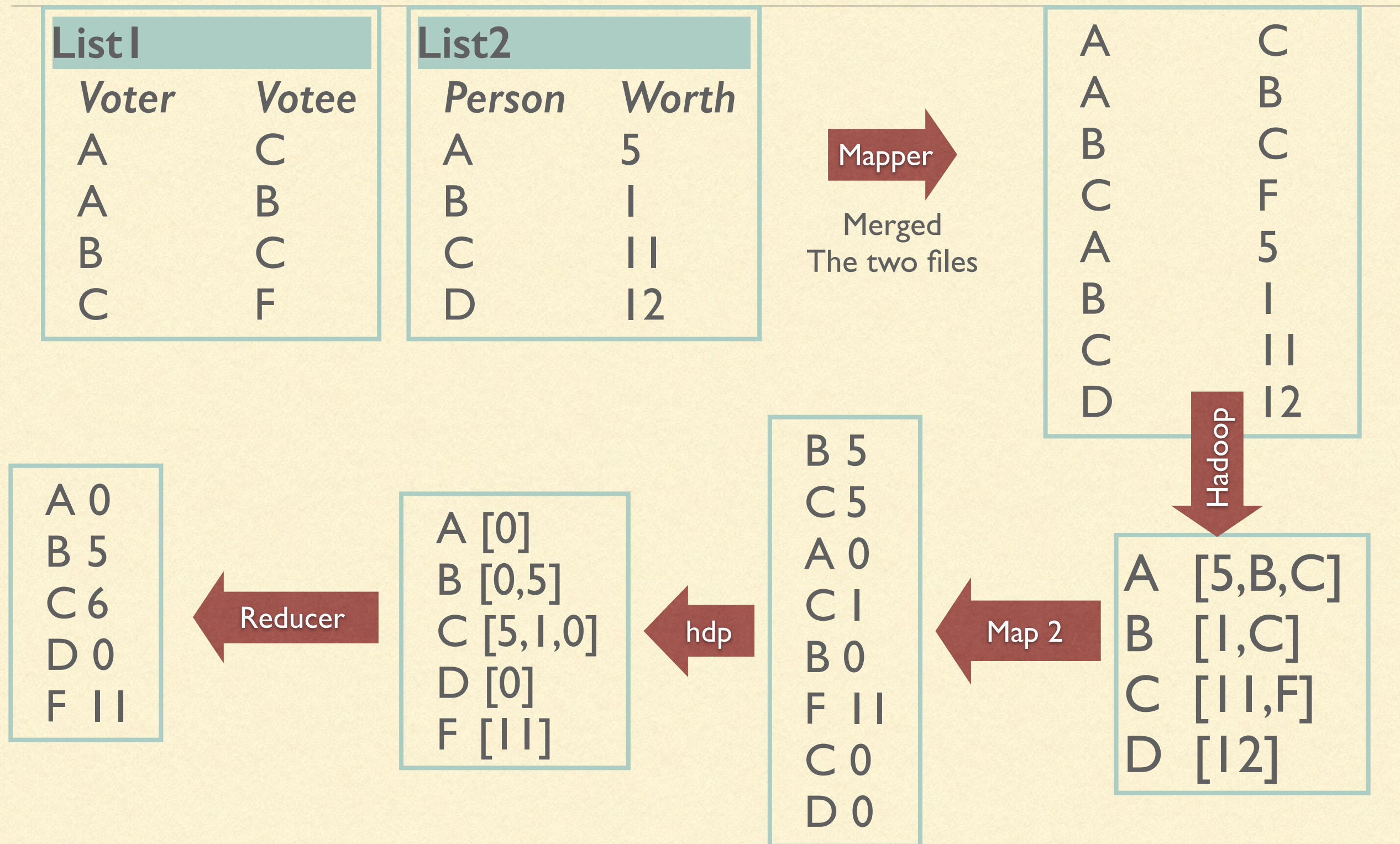
Find out what is the vote count of everyone?

List1		List2		Result	
Voter	Votee	Person	Worth	Person	VoteCount
A	C	A	5	A	0
B	C	B	1	B	5
C	F	C	11	C	6
A	B			F	11



MAP / REDUCE

Problem 4: Voting



MAP / REDUCE

Limitation / When Not to Use?

- If the job can be done by a single machine in reasonable time
- Computation depends on previously computed values.
 - e.g. Fibonacci Series
- Full-text indexing or ad hoc searching
- Algorithms depend on shared global state

Testing

1. First test on very small data
2. Separately Test Mapper and Reducer
3. Streaming Job's Mapper could be tested with simple unix command:
 1. `cat inputfile | mymapper | sort | myreducer >> outputfile`
4. Use MRUnit (Java MR Session, Next class)
5. To test predictions, you may want to test the part

MAP / REDUCE - JAVA

```
public class StubMapper extends Mapper<LongWritable, Text, Text, IntWritable> {

    @Override
    public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {

        /*
         * TODO implement
         */

    }
}
```

```
public class StubReducer extends Reducer<Text, IntWritable, Text, DoubleWritable> {

    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
        throws IOException, InterruptedException {

        /*
         * TODO implement
         */

    }
}
```

```
public class StubDriver {
    public static void main(String[] args) throws Exception {
        Job job = new Job();
        job.setMapperClass(StubMapper.class);
        job.setReducerClass(StubReducer.class);
    }
}
```

MAP / REDUCE

Assignment - Problem I

Code M/R for all the problems in your favourite language.

Assignment - Problem 2

Based on the content from a very large text archive, formulate the next words recommendation.

For each word, prepare a top 5 recommendations of the word that would go next.

happy birthday, newyear, marriage

how are, do, did

...

Download the content from wikipedia using this:

<http://www.evanjones.ca/software/wikipedia2text.html>



Big Data & Hadoop

Thank you.



+91-9538998962

sandeep@knowbigdata.com