

Module 2: Assignment - 4

1. Find out the count of each word in the 'Shakespeare.txt' dataset in the 'Shakespeare.rar'

Solution:

Reading Shakespeare.txt

```
bitnami@debian:~$ hadoop fs -cat /user/bitnami/Shakespeare.txt_
```

Par. Faith sir, ha's led the drumme before the English
Tragedians: to belye him I will not, and more of his
souldiership I know not, except in that Country, he had
the honour to be the Officer at a place there called Mile-end,
to instruct for the doubling of files. I would doe the
man what honour I can, but of this I am not certaine

Cap.G. He hath out-villain'd villanie so farre, that the
raritie redeemes him

Ber. A pox on him, he's a Cat still

Int. His qualities being at this poore price, I neede
not to aske you, if Gold will corrupt him to reuolt

Par. Sir, for a Cardceue he will sell the fee-simple of
his saluation, the inheritance of it, and cut th' intaile from
all remainders, and a perpetuall succession for it perpetually

Int. What's his Brother, the other Captain Dumain?

Cap.E. Why do's he aske him of me?

Int. What's he?

Par. E'ne a Crow a'th same nest: not altogether so
great as the first in goodnesse, but greater a great deale in

Running Map Reduce Job:

```
bitnami@debian:~$ hadoop jar map_reduce_3.jar map_reduce_3/WordCount /user/bitnami/Shakespeare.txt
/user/bitnami/output_wordcount_1
2021-11-04 18:29:24,654 INFO client.DefaultNoHARMAFailoverProxyProvider: Connecting to ResourceMana
ger at /0.0.0.0:8032
Exception in thread "main" org.apache.hadoop.mapred.FileAlreadyExistsException: Output directory h
dfs://localhost:8020/user/bitnami/output_wordcount_1 already exists
    at org.apache.hadoop.mapreduce.lib.output.FileOutputFormat.checkOutputSpecs(FileOutputForm
at.java:164)
    at org.apache.hadoop.mapreduce.JobSubmitter.checkSpecs(JobSubmitter.java:277)
    at org.apache.hadoop.mapreduce.JobSubmitter.submitJobInternal(JobSubmitter.java:143)
    at org.apache.hadoop.mapreduce.Job$11.run(Job.java:1571)
    at org.apache.hadoop.mapreduce.Job$11.run(Job.java:1568)
    at java.security.AccessController.doPrivileged(Native Method)
    at javax.security.auth.Subject.doAs(Subject.java:422)
    at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1878)
    at org.apache.hadoop.mapreduce.Job.submit(Job.java:1568)
    at org.apache.hadoop.mapreduce.Job.waitForCompletion(Job.java:1589)
    at map_reduce_3.WordCount.main(WordCount.java:94)
    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
    at java.lang.reflect.Method.invoke(Method.java:498)
    at org.apache.hadoop.util.RunJar.run(RunJar.java:323)
    at org.apache.hadoop.util.RunJar.main(RunJar.java:236)
bitnami@debian:~$ hadoop jar map_reduce_3.jar map_reduce_3/WordCount /user/bitnami/Shakespeare.txt
/user/bitnami/out_wordcount_1
2021-11-04 18:29:59,595 INFO client.DefaultNoHARMAFailoverProxyProvider: Connecting to ResourceMana
ger at /0.0.0.0:8032
2021-11-04 18:30:00,034 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not
performed. Implement the Tool interface and execute your application with ToolRunner to remedy th
is.
2021-11-04 18:30:00,052 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tm
p/hadoop-yarn/staging/hadoop/.staging/job_1636027398107_0002
2021-11-04 18:30:00,327 INFO input.FileInputFormat: Total input files to process : 1
2021-11-04 18:30:00,812 INFO mapreduce.JobSubmitter: number of splits:1
2021-11-04 18:30:01,053 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1636027398107_
0002
2021-11-04 18:30:01,053 INFO mapreduce.JobSubmitter: Executing with tokens: []
2021-11-04 18:30:01,276 INFO conf.Configuration: resource-types.xml not found
2021-11-04 18:30:01,276 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2021-11-04 18:30:01,352 INFO impl.YarnClientImpl: Submitted application application_1636027398107_
0002
2021-11-04 18:30:01,415 INFO mapreduce.Job: The url to track the job: http://localhost:8088/proxy/
application_1636027398107_0002/
2021-11-04 18:30:01,416 INFO mapreduce.Job: Running job: job_1636027398107_0002
2021-11-04 18:30:08,588 INFO mapreduce.Job: Job job_1636027398107_0002 running in uber mode : fals
e
2021-11-04 18:30:08,591 INFO mapreduce.Job:  map 0% reduce 0%
2021-11-04 18:30:15,735 INFO mapreduce.Job:  map 100% reduce 0%
2021-11-04 18:30:22,797 INFO mapreduce.Job:  map 100% reduce 100%
2021-11-04 18:30:23,836 INFO mapreduce.Job: Job job_1636027398107_0002 completed successfully
```

Viewing count for each word using two reducers only output:

```
bitnami@debian:~$ hadoop fs -ls /user/bitnami/out_wordcount_1/
Found 2 items
-rw-r--r-- 1 hadoop supergroup          0 2021-11-04 18:30 /user/bitnami/out_wordcount_1/_SUCCESS
-rw-r--r-- 1 hadoop supergroup    746266 2021-11-04 18:30 /user/bitnami/out_wordcount_1/part-r-00000
```

```
bitnami@debian:~$ hadoop fs -cat /user/bitnami/out_wordcount_1/part-r-00000.....
```

```
willingly      24
willingly,      4
willingly.      1
willingly:      3
willingnesse    1
willingnesse.   1
willow  2
wills  6
wills,  4
wills.  1
wills;  1
wils  1
wils:  1
wils?  1
wilt  215
wilt,  11
wilt.  4
wilt:  2
wilt;  1
wilte  1
wimpled,      1
win  56
win,  4
win.  2
win;  1
winch,  1
winch:  1
winck  1
wincke,  1
wind  7
wind,  3
wind-pipes  1
wind-swift  1
wind:  2
winde  74
winde)  1
winde,  25
winde-obeying  1
winde-shak'd-Surge,  1
winde-shaken.  1
```

```
yrefull 1
yron  5
yron:  2
ysickle 1
ysicles,      1
yssue  4
yssue,  1
yssue:  1
yssued,  1
yssues.  1
yt  2
yut  1
z  1
zeale  11
zeale,  7
zeale:  1
zeale?  1
zeales,  1
zealous 3
zelous  2
zip  1
zo  1
zwaggerd      1
```

3. Find out the most commonly used words (Words with the count over 100 are considered common).

Solution:

Running Map Reduce Job

```
bitnami@debian:~$ hadoop jar map_reduce_4.jar map_reduce_3/WordCount /user/bitnami/Shakespeare.txt /user/bitnami/output_wordcount_5
```

```
      HDFS: Number of bytes read erasure-coded=0
Job Counters
      Killed reduce tasks=1
      Launched map tasks=1
      Launched reduce tasks=2
      Data-local map tasks=1
      Total time spent by all maps in occupied slots (ms)=9224
      Total time spent by all reduces in occupied slots (ms)=16240
      Total time spent by all map tasks (ms)=4612
      Total time spent by all reduce tasks (ms)=8120
      Total vcore-milliseconds taken by all map tasks=4612
      Total vcore-milliseconds taken by all reduce tasks=8120
      Total megabyte-milliseconds taken by all map tasks=9445376
      Total megabyte-milliseconds taken by all reduce tasks=16629760
Map-Reduce Framework
      Map input records=135281
      Map output records=817811
      Map output bytes=7680316
      Map output materialized bytes=1020341
      Input split bytes=115
      Combine input records=817811
      Combine output records=70503
      Reduce input groups=70503
      Reduce shuffle bytes=1020341
      Reduce input records=70503
      Reduce output records=806
      Spilled Records=141006
      Shuffled Maps =2
      Failed Shuffles=0
      Merged Map outputs=2
      GC time elapsed (ms)=214
      CPU time spent (ms)=6980
      Physical memory (bytes) snapshot=807919616
      Virtual memory (bytes) snapshot=10361499648
      Total committed heap usage (bytes)=645398528
      Peak Map Physical memory (bytes)=366923776
      Peak Map Virtual memory (bytes)=3451076608
      Peak Reduce Physical memory (bytes)=221237248
      Peak Reduce Virtual memory (bytes)=3456024576
Shuffle Errors
      BAD_ID=0
      CONNECTION=0
      IO_ERROR=0
      WRONG_LENGTH=0
      WRONG_MAP=0
      WRONG_REDUCE=0
File Input Format Counters
      Bytes Read=4516586
File Output Format Counters
      Bytes Written=7701
```

Viewing output for most commonly used words (count > 100)

```
bitnami@debian:~$ hadoop fs -cat /user/bitnami/output_wordcount_5/part-r-00000_
```

'Tis	523	euen	280	then,	244
'tis	692	euer	538	these	1000
1	202	eye	143	things	224
A	1513	face	173	thinke	748
Against	126	faire	464	this?	187
Ah	102	fall	196	those	411
All	327	false	177	though	356
All.	108	feare	350	thought	271
An	127	feare,	102	through	201
An.	139	fellow	111	thy	3488
And	7029	fit	123	till	383
Ant.	450	fiue	100	time	613
Art	105	follow	204	time,	222
At	203	for	5372	to	14978
Be	411	forth	231	told	198
Before	106	foule	162	truth	124
Being	101	foure	108	verie	118
Ber.	217	friends	113	vnder	178
Bru.	200	full	334	vnto	339
Brutus	101	gentle	287	vpon	1240
But	2326	goe	432	vs	1086
By	679	gone	173	vs,	266
Caesar	191	gone,	129	warrant	137
Caesar,	107	good	1984	was	1856
Cassi.	108	good,	106	way	330
Clau.	100	gracious	158	way,	112
Cleo.	210	great	688	wee	218
Court	103	ha's	160	were	1263
Did	218	had	1194	where	615
Do	277	hand	350	which	1220
Du.	108	hand,	248	while	135
Duke	444	hard	118	whom	285
Duke.	205	hath	1541	wife	139
England	111	haue	5023	wife,	100
Euen	216	haue,	104	wish	169
Exeunt.	639	he	4051	within	193
Fal.	327	he's	103	words	203
Fathers	188	he,	108	words,	105
For	1636	heare	665	worthy	148
Fortune	103	heartes	117	would	1882
France	126	heere	623	wrong	122
France,	154	heere,	163	y	133
Gentleman	122	hence	131	yong	197
Gods	238	hence,	123	you	9838
Goe	132	her	2759	you,	1043
Good	425	her,	373	you.	213
Had	137	her:	127	you:	281
Hath	274	him?	135	your	6186
Haue	410	himselfe	233	youth	103