

Hadoop - Wordcount:

LAB - G : (20/5/28)

short-all.sh

```
hadoop map.py → #!/usr/bin/env python3  
import sys  
def mapper():  
    for line in sys.stdin:  
        line = line.strip()  
        if not line:  
            continue  
        words = line.split()  
        for word in words:  
            print(f'{word}\t1')  
    } --name-- = "main":  
    mapper()
```

hadoop red.py → #!/usr/bin/env python3

c_word = None

c_count = 0

for line in sys.stdin:

line = line.strip()

if not line:

continue

word, count = line.split(' \t', 1)

try:

count = int(count)

except ValueError:

continue

if c_word == word:

c_count += count

else:

c_word:

print(f'{c_word}\t{c_count}')

c_word = word

c_count = count

if c_word:

print(f'{c_word}\t{c_count}')

(some commands or next top ten words)

Hadoop - Average Temperature

start-dfs.sh

start-yarn.sh

jps

hdfs --daemon start namenode

hdfs namenode -format

start-dfs.sh

hdfs dfs -mkdir /weatherdata

~~hdfs dfs~~

git clone https://github.com/

hdfs dfs -put hadoopbook/input/ncdc/all/* /weatherdata

mapper.py

namemapper.py

chmod +x mapper.py

name reducer.py

chmod +x reducer.py

hadoop jar jar-path \

input/weatherdata \

output/weather_output \

-mapper mapper.py

-reducer reducer.py

hdfs dfs -cat /weather-output/

part-0000

1901 46.70

1902 21.66

#!/usr/bin/env python3

import sys

for line in sys.stdin:

year = line[5:9]

temp = line[87:92]

if temp == "9999":

print(f"Year {year} Temp {temp}")

reducer.py

import sys

c-year = None

temp_sum = 0

count = 0

for line in sys.stdin:

year, temp = line.split()

split('t')

temp = int(temp)

if c-year == year:

temp_sum += temp

count += 1

else:

if c-year :

print(f"Year {c-year} Temp {temp}")

temp_sum / count

c-year = year

temp_sum = temp

count = 1

if c-year:

print(f"Year {c-year} Temp {temp}")

Hadoop → Map Reduce program Hadoop → Mean Max Temp (about 800 commands or previous program)

mapper.py:

#!/usr/bin/env python3

import sys

for line in sys.stdin:

 year = line[5:19]

 month = line[19:21]

 temp = line[87:92]

 if temp == "+9999":

 print(f"{{year}}-{month}{{temp}}")

reducer.py:

import sys

c_year = None

temp_sum = 0

count = 0

for line in sys.stdin:

 year_month, temp = line.split("-").split("t")

 temp = int(temp)

 if c_year == year_month:

 temp_sum += temp

 count += 1

 else:

 c_year

 print(f"{{c_year}} {{temp_sum / count :.2f}}")

 c_year = year_month

 temp_sum = temp

 count = 1

 else:

 c_year

 print(f"{{c_year}} {{temp_sum / count :.2f}}")

Output:

1901-01 - 25.88

1901-02 - 91.72

1901-03 - 22.54

Hadoop → Map Reduce program to list top 10 max words

start-all.sh

hadoop.map.py

```
#!/usr/bin/env python3
import sys
import re
for line in sys.stdin:
    words = re.findall(r'\w+', line.lower())
    for word in words:
        print(f'{word}\t1')
```

hadoop.red.py

```
#!/usr/bin/env python3
import sys
```

```
from collections import defaultdict
```

```
count_map = defaultdict(int)
```

```
for line in sys.stdin:
    words, count = line.strip().split('\t')
    count_map[word] += int(count)
```

```
sorted_words = sorted(count_map.items(), key=lambda x: (-x[1], x[0]))
```

```
for word, count in sorted_words[:10]:
    print(f'{word}\t{count}'")
```

chmod +x map.py

chmod +x red.py

hdfs dfs -mkdir /bdall

hdfs dfs -ls /

(make a input file (inp.txt))

hdfs dfs -put inp.txt /bdall

hadoop jar /home/hadoop/hadoop/share/hadoop/tools/lib/

hadoop-streaming-3.3.6.jar -input /bdall/input.txt -output /bdall/output -mapper /home/hadoop/mop.py

-reducer /home/hadoop/red.py -file /home/hadoop/mop.py

-file /home/hadoop/red.py

hadoop fs -ls /bdall ("if/fibonacci") string

hadoop fs -ls /bdall/output

hadoop fs -cat /bdall/output/part-0000

fibonacci (true) numbers
(true) fibonacci = done true

if(fib(0) == 0) {
 fib(0) = 0; return 0;

(true) true = + [true] false true

fibonacci(1) = 1; (true) false true both 02 = fibon(1) 02
((0) & (1)) : 0

: [0]: fibon(0) both 02 = fibon(0), fibon(1)
(true) 3: fibon(1) both 02

if (fib(0) == 0) {
 fib(0) = 0; return 0;

((0) & (0)) : 0
fib(0) both 02

((true) 0) : 0
fib(0) both 02

((0) & (1)) : 0
fib(1) both 02

Scala: → To print numbers from 1 to 100 using for loop
In Ubuntu user.

\$ sudo apt install scala

\$ nano prints.scala

object prints {

```
def main(args: Array[String]): Unit = {
    for (i <= 1 to 100) {
        print(i)
    }
}
```

ctrl + o, enter, ctrl + x
\$ scala prints.scala

\$ scala prints
1 2 3 ... 100

- 3610-1-202
Habmo

object prints {

def main(args: Array[String]): Unit =

for (i <= 1 to 100) {

print(i)

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

}

<

Spark : count words > 4

sudo apt update

wget https://archive.apache.org/dist/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz

tar -xvzf spark-3.5.1-bin-hadoop3.tgz

ls /opt | grep spark

cd /opt

ls /opt | grep spark

nono ~/.bashrc

→ add these values on top
 export SPARK_HOME = /opt/spark
 export PATH = \$SPARK_HOME/bin:\$PATH
 export JAVA_HOME = /usr/lib/jvm/java-11-openjdk-amd64

source ~/.bashrc

spark-shell

New window

mkdir ~/sparks

nono wordcount.py

```
→ from pyspark import SparkContext
sc = SparkContext("local", "wordCountApp")
text_rdd = sc.textFile("input.txt")
words_rdd = text_rdd.flatMap(lambda line: line.split())
pairs_rdd = words_rdd.map(lambda word: (word, 1))
words_count_rdd = pairs_rdd.reduceByKey(lambda a, b: a+b)
words_count_rdd = words_count_rdd.filter(lambda x: x[1] > 4)
filter_rdd = words_count_rdd.collect()
for word, count in filter_rdd:
    print(f"Word {word} : {count}")
filter_rdd.saveAsTextFile("output_directory")
```

new input.txt (your text)

(choose) spark-submit wordcount.py input.txt

output:

The: 6

lorem: 6

ipsum: 7

: 12.000 words

sentences and words

- tricking (po. nrof words pro topic)

- koffit. Ip. nrof words pro topic

((topic))² ("no" "do", "it", "et", "kno", "an", "ai", "st", "st") + 98 = Hercegovačka

3. alien know = pmr : (pmr: know) 95% normal job

(("2")) x if ((pmr <= ("pmr")) alienknow.w { w 910}

(("2")) x if ((pmr <= ("pmr")) alienknow.w { w 910})

know (a - 910)

? C = (pmr: on)) floor - testFloor know

((("testFloor" - 15 - 0.5)) 940) else

(("12 / 10" - 0.5))

((("12 / 10" - 0.5)) 940) else

((("12 / 10" - 0.5)) 940) else

((("12 / 10" - 0.5)) 940) else

modifFloor.210 (?) = result know

((("testFloor" - 15 - 0.5)) 940) else

((("testFloor" - 15 - 0.5)) 940) else

((("12 / 10" - 0.5)) 940) else

((("12 / 10" - 0.5)) 940) else

= know) or

((("12 / 10" - 0.5)) 940) else

((("12 / 10" - 0.5)) 940) else

= know) or

open-ended question : Text cleaning & Print cleaned text (Spark)

~~Mapper.py:~~

~~File used to implement python~~

~~Map~~

import org.apache.spark.sql.functions -

import org.apache.spark.sql.types -

val stopWords = Set("the", "is", "in", "and", "to", "a", "of", "on", "that")

def lemmatize(word: String): String = word match {

case w if w.endsWith("ing") => w.stripSuffix("ing")

case w if w.endsWith("ed") => w.stripSuffix("s")

case _ => word

}

val cleanText = udf((line: String) => {

line.toLowerCase

- replaceAll("[^a-zA-Z]+", "")

- split(" ")

- filter(w => w.nonEmpty && !stopWords.contains(w))

- map(lemmatize)

- mkString(" ")

)

val lines = spark.readStream

- format("socket")

- option("host", "localhost")

- option("port", 9999)

, load()

val cleaned =

lines.withColumn("cleaned",

cleanText(col("value")))

val query = cleaned.writeStream

• outputMode ("opend")
• format ("console")
• short()

query.awaitTermination()

}

as
22/5/25