

20/05/25 Lab 9: - Scala program to print numbers from 1 to 100 using a for loop

```
scala> for (i <- 1 to 100) {  
  println(i)  
}
```

1

2

3

4

...

...

→ Using RDD and FlatMap count how many times each word appears in a file and write out a list of words whose count is strictly greater than 4 using Spark

```
scala> val textFile = sc.textFile("/home/  
bmcecece/text.txt")
```

```
scala> val words = textFile.flatMap(line  
=> line.split("\\W+"))
```

```
scala> val wordPairs = words.map(word  
=> (word.toLowerCase, 1))
```

```
scala> val wordCounts = wordPairs.  
  reduceByKey(_+_)
```



```
scala> val frequentWords = wordCounts.filter
  {case (_, count) => count > 4}
```

```
scala> frequentWords.collect {l} foreach
  (println)
```

```
(hello, 5)
```

```
(hi, 6)
```

→ For a given Text file, create a MapReduce program, to sort the content in an alphabetic order listing only top 10 maximum occurrences of words (Hadoop program)

→ Mapper.py

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

```
for line in sys.stdin:
    line = line.strip().lower()
```

```
words = re.findall(r'[a-z]{1,16}', line)
```

```
for word in words:
```

```
    print(f"{word} {1 + 1}")
```

Reducer.py

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

```
import sys
```

```
from collections import defaultdict
```

```
word-count = defaultdict(int)
```

```
for line in sys.stdin:
```

```
    line = line.strip()
```

```
    if not line:
```

```
        continue
```

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

```
    word-count[word] += int(count)
```

```
top-words = sorted(word-count.items(),
```

```
    key = lambda n: (-n[1], n[0]))[:10]
```

```
top-words = sorted(top-words,
```

```
    key = lambda n: n[1])
```

```
for word, count in top-words:
```

```
    print(f"{word}\t{count}")
```

→ `hadoop jar /home/bmccge/had.jar \`

`-input /home/bmccge/inp.txt \`

`-output /home/bmccge/output-dir \`

`-mapper mapper.py \`

`-reducer reducer.py \`

`-file mapper.py \`

`-file reducer.py`

25/5/2015