-: LAB : 09 :-

**A** Write a Scala program to print numbers from 1 to 100 using for loop

=> Spark-shell

Scala > object PrintNumbers {
    def main (args: Array[String]): Unit = {
        // Loop from 1 to 100 and print each
        //number
        for (i <- 1 to 100) {
            println(i)
        }
    }
}

defined object PrintNumbers

scala> PrintNumbers. main (Array())

**# Output:**
1
2
3
.
.
.
.
.
.
.
100

**☆** Using RDD and FlatMap, count how many times each word time appears in a file and write out a list of words whose count is strictly greater than $\xi$ using Spark.

=>

  textfile (sc.textfile)

val data = sc. textFile("sparkdata. txt")
data. collect;
val splitdata = data. flatmap (line => line
split (" "));

split data . collect;

val mapdata = splitdata .map (word => (word,1)).

mapdata . collect;

val reduce data = mapdata . reduce ByKey (+);

reduce data. collect;

# output :-

Hello: 5
Tree: 4 //

☆ Hadoop program (open ended)

→ mapper.py

```
[usr/bin/env python3
import sys
import re

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

    words = re.findall(r'[bra-z]#[b:(-)]
    for word in words:
        print(f'{word}\t1')
```

→ Reducer.py

```
[usr/bin/env python3

import sys
from collections improve default dict

word-count.default dist(int)

for line in sys.std in:
    line = line.strip()
    if not line
        continue
    word count= line .split('t')
    word.count[word]+= int(count)
```

top_words = sorted (word. count items
key = contine x : (-x(1), x(0))) [:10)

top word count is top words:
print (f" {word}it {count}")

→ hadop jar /home/ bmsre/ badyor
- mapper /mapper. py `
- reducer reducer. py `
- file mapper. py
- file reducer - py.

d-4 20/6/25