LABORATORY PROGRAM - 8

Scala Program to print numbers from 1 to 100

```
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize
import re
# Initialize SparkContext and StreamingContext (batch interval = 5 sec)
sc = SparkContext("local[2]", "TextCleanStreamingApp")
ssc = StreamingContext(sc, 5)
# Set up stop words and lemmatizer
stop_words = set(stopwords.words('english'))
lemmatizer = WordNetLemmatizer()
# Connect to the socket stream
lines = ssc.socketTextStream("localhost", 9999)
# Text cleaning function
def clean_text(line):
 # Remove non-alphabetic characters and lower the case
 line = re.sub(r'[^a-zA-Z\s]', '', line)
 line = line.lower()
 # Tokenize
 tokens = word_tokenize(line)
 # Remove stop words and lemmatize
 cleaned_tokens = [lemmatizer.lemmatize(word) for word in tokens if word not in stop_words]
```

return ' '.join(cleaned_tokens)

Apply cleaning to each line
cleaned_lines = lines.map(clean_text)

Print the cleaned lines
cleaned_lines.pprint()

Start streaming
ssc.start()
ssc.awaitTermination()

OBSERVATION

```
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC: ~ Q =
```





