

lab 10

Write a Simple Streaming program in spark to receive text data streams on a particular port, perform basic text cleaning (like white space removal, stop word removal, lemmatization, etc), and print the cleaned text on the screen (o/p).

→ Install dependencies & download NLTK data:
`text-cleaning-stream.py`

```
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
import re
```

```
sc = SparkContext("local[*]", "Text Cleaning Stream")
```

```
scc = StreamingContext(sc, 5)
```

```
stop_words = set(stopwords.words('english'))
```

```
lemmatizer = WordNetLemmatizer()
```

```
def clean_text(line):
```

```
    line = line.lower()
```

```
    line = re.sub(r"[^a-z\s]", "", line)
```

```
    tokens = [word for word in line.split() if word not in stop_words]
```

```
    lemmatized = [lemmatizer.lemmatize(word) for word in tokens]
```

```
    return " ".join(lemmatized)
```

```
lines = scc.socketTextStream("localhost", 9999)
```

```
cleaned = lines.map(clean_text)
```

```
cleaned.pprint()
```

```
scc.start()
```

```
scc.awaitTermination()
```

`entry.getValue(i));`

o/p = The quick brown fox jumps over the lazy dog.

→ Quick brown fox jump lazy dog.