

## **Department of Computer Science and Engineering (Data Science)**

### **Machine Learning – IV**

#### **Experiment 1**

**Name: Umang Kirit Lodaya**

**SAP ID: 60009200032**

**Batch: D11**

#### **Aim:**

To implement and understand the Word Count algorithm using the MapReduce programming model, a fundamental application for processing and analyzing large text datasets in a parallel and distributed manner.

#### **Theory:**

##### **Introduction to MapReduce:**

- MapReduce is a programming model designed for processing and generating large datasets in parallel across distributed clusters.
- It consists of two main phases: Map and Reduce.

##### **Word Count Overview:**

- **Map Step:**
  - Mapper Function:
    - For each input record (line of text), the mapper function emits key-value pairs.
    - The key is the word, and the value is the count (usually 1).
- **Shuffling and Sorting:**
  - The MapReduce framework automatically shuffles and sorts the emitted key-value pairs, grouping them by key.
- **Reduce Step:**
  - Reducer Function:
    - For each group of values with the same key, the reducer function sums the counts to get the total occurrences of each word.

## **Step-by-Step Implementation:**

- **Step 1: Mapper Function**
  - Implement the mapper function to tokenize words and emit key-value pairs.
- **Step 2: Shuffling and Sorting**
  - Understand the automatic shuffling and sorting mechanism of the MapReduce framework.
- **Step 3: Reducer Function**
  - Implement the reducer function to aggregate word counts.
- **Step 4: Input Data Formatting**
  - Prepare input data in a format suitable for MapReduce processing (e.g., a collection of text documents).
- **Step 5: Experimentation with Input Sizes**
  - Experiment with different sizes of input text datasets to observe the scalability of the MapReduce approach.

## **Implementation Tips:**

- Handle special cases like punctuation, case sensitivity, and stemming appropriately.
- Optimize the mapper and reducer functions for efficient word counting.

## **Lab Experiments to be Performed in This Session:**

Execute the Word Count using Map Reduce on a dataset to gain insights into its functionality and operation.

**NAME: UMANG KIRIT LODAYA**

**SAP ID: 60009200032**

**BATCH: K1 / D11**

**AIM: WORD COUNT USING MAP REDUCE**

```
In [1]: import os
```

```
In [2]: PATH = '/content/drive/MyDrive/Colab Notebooks/ML - IV'
```

```
In [3]: FILES = []
        for x in os.listdir(PATH):
            if x.startswith("file"):
                FILES.append(x)
```

```
In [4]: FILES
```

```
Out[4]: ['file.txt', 'file2.txt']
```

```
In [5]: FILE = []
        for fil in FILES:
            with open(PATH + f'/{fil}', 'r') as f:
                FILE.extend(f.readlines())
```

```
In [6]: F = []
        for i in range(len(FILE)):
            sent = FILE[i].lower().split()
            new = []
            for word in sent:
                n = ''
                for char in word:
                    if char.isalnum():
                        n += char

                new.append(n)

            F.extend(new)
```

```
In [7]: with open(PATH + '/mapper.txt', 'w') as f:
        for word in F:
            f.write(f"{word},1\n")
```

**SAMPLE OUTPUT OF MAPPER**

mapper.txt ×

```
1 the,1
2 first,1
3 sentence,1
4 establishes,1
5 the,1
6 key,1
7 idea,1
8 of,1
9 balance,1
10 theory,1
11 the,1
12 next,1
13 sentence,1
14 begins,1
15 with,1
16 balance,1
17 theory,1
18 and,1
19 ends,1
20 with,1
21 social,1
22 ties,1
23 which,1
24 is,1
25 the,1
26 focus,1
27 of,1
28 the,1
```

```
In [8]: MAP_OUT = None
with open(PATH + '/mapper.txt', 'r') as f:
    MAP_OUT = f.readlines()

SHUFFLE = sorted(MAP_OUT)

with open(PATH + '/shuffle.txt', 'w') as f:
    for word in SHUFFLE:
        f.write(f"{word}")
```

**SAMPLE OUTPUT OF SHUFFLER**

```
mapper.txt  shuffle.txt X
1  a,1
2  a,1
3  a,1
4  a,1
5  a,1
6  a,1
7  about,1
8  and,1
9  and,1
10 and,1
11 and,1
12 awkward,1
13 balance,1
14 balance,1
15 balance,1
16 balance,1
17 be,1
18 before,1
19 before,1
20 begins,1
21 begins,1
22 benefits,1
23 but,1
24 but,1
25 choose,1
26 choppy,1
27 cliques,1
28 coherent,1
29 cohesion,1
30 comes.1
```

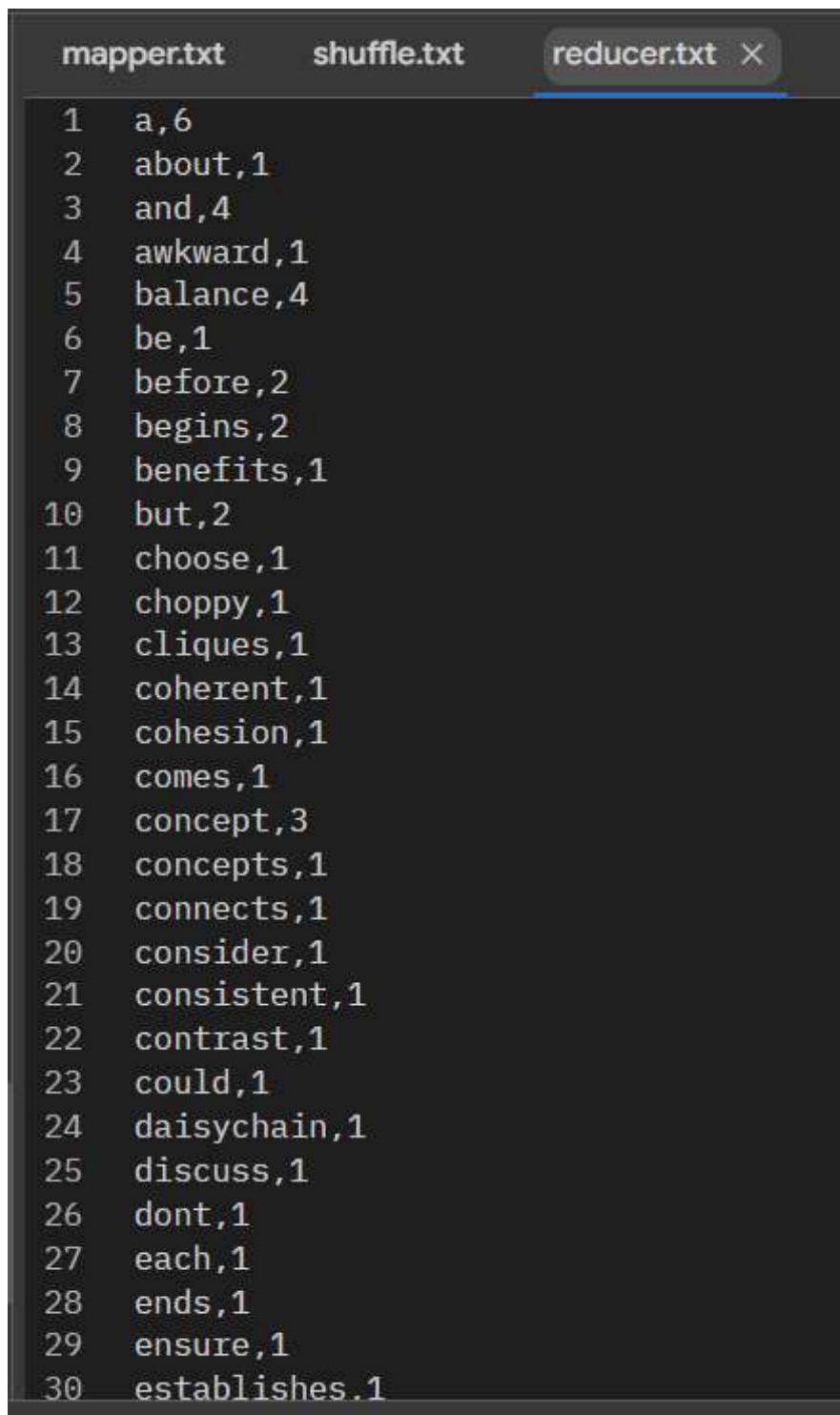
```
In [9]: SHUFFLE = None
with open(PATH + '/shuffle.txt', 'r') as f:
    SHUFFLE = f.readlines()

freq = {}

for word in SHUFFLE:
    word = word[:-1]
    word, count = word.split(',')
    count = int(count)
    freq[word] = freq.get(word, 0) + 1
```

```
In [10]: with open(PATH + '/reducer.txt', 'w') as f:
    for key in freq:
        f.write(f"{key},{freq[key]}\n")
```

## SAMPLE OUTPUT OF REDUCER



```
mapper.txt  shuffle.txt  reducer.txt X
1  a,6
2  about,1
3  and,4
4  awkward,1
5  balance,4
6  be,1
7  before,2
8  begins,2
9  benefits,1
10 but,2
11 choose,1
12 choppy,1
13 cliques,1
14 coherent,1
15 cohesion,1
16 comes,1
17 concept,3
18 concepts,1
19 connects,1
20 consider,1
21 consistent,1
22 contrast,1
23 could,1
24 daisychain,1
25 discuss,1
26 dont,1
27 each,1
28 ends,1
29 ensure,1
30 establishes.1
```