| Sr.No | Practical Name | Date | Sign |
|-------|----------------|------|------|
| 1 | Write a program to demonstrate bitwise operations. | | |
| 2 | Implement Page Rank Algorithm. | | |
| 3 | Implement Dynamic programming algorithm for computing the edit distance between string s1 and s2 (Levenshtein Distance). | | |
| 4 | Write a program to Compute Similarity between two text documents. | | |
| 5 | Implement a map reduce program to count words and ignore case in Hadoop environment | | |
| 6 | Implement IR system using lucene | | |
| 7 | Write a program for Pre-processing of a Text Document: stop word removal. | | |
| 8 | Write a program to implement simple web crawler. | | |

# Practical 1

```
In [5]: a = 60                # 60 = 0011 1100
        b = 13                # 13 = 0000 1101
        c = 0

        c = a & b;            # 12 = 0000 1100
        print ("Line 1 - Value of c is ", c)

        c = a | b;            # 61 = 0011 1101
        print ("Line 2 - Value of c is ", c)

        c = a ^ b;            # 49 = 0011 0001
        print ("Line 3 - Value of c is ", c)

        c = ~a;               # -61 = 1100 0011
        print ("Line 4 - Value of c is ", c)

        c = a << 2;           # 240 = 1111 0000
        print ("Line 5 - Value of c is ", c)

        c = a >> 2;           # 15 = 0000 1111
        print ("Line 6 - Value of c is ", c)

Line 1 - Value of c is  12
Line 2 - Value of c is  61
Line 3 - Value of c is  49
Line 4 - Value of c is  -61
Line 5 - Value of c is  240
Line 6 - Value of c is  15
```

# Practical 2

## Simple Page Rank

```java
package MyPageRank;
import java.util.Arrays;
import java.util.Scanner;

public class MyPageRank {
    static int path[][];
    static int nodes;
    static int outbounds[];
    static void calculatePageRank(){
        double tempRank[]=new double[nodes];
        //Giving equal rank to all the pages or nodes
        //double damping=0.85;
        for(int i=0;i<nodes;i++){
            tempRank[i]=1/(float)nodes;
            System.out.println("Initial page rank for page "+(i+1)+" : "+tempRank[i]);
        }
        double newPageRank[]=new double[nodes];
        Arrays.fill(newPageRank, 0);
        int itr=1;

        while(itr<=2){
        for(int i=0;i<nodes;i++){
            for(int j=0;j<nodes;j++){
                if(path[j][i]==1){
                    newPageRank[i]=newPageRank[i]+tempRank[j]/(float)outbounds[j];
                }

            }

            //if you want to include damping or teleportation factor
            //newPageRank[i]=(1-damping)+damping*(newPageRank[i]);

        System.out.println("For iteration "+itr+" page rank for node : "+(i+1)+" is = "+newPageRank[i]);
        }
        for(int i=0;i<nodes;i++){
            tempRank[i]=newPageRank[i];
        }
            Arrays.fill(newPageRank,0);
            itr++;
        }
    }

    public static void main(String arg[]){
        Scanner sc=new Scanner(System.in);
        System.out.println("Enter number of nodes :");
        nodes=sc.nextInt();
        path=new int[nodes][nodes];
        outbounds=new int[nodes];
        Arrays.fill(outbounds,0);
        System.out.println("Enter the graph : ");
        for(int i=0;i<nodes;i++){
            for(int j=0;j<nodes;j++){
                path[i][j]=sc.nextInt();
                if(i==j){
                    path[i][j]=0;
                }
```

```
57                    //calculate outbound links for all
58                    if(path[i][j]==1){
59                        outbounds[i]+=1;
60                    }
61                }
62            }
63        calculatePageRank();
64        }
65
66    }
```

run:
Enter number of nodes :
4
Enter the graph :
1 0 1 1
0 0 1 0
1 1 0 0
1 0 1 0
Initial page rank for page 1 : 0.25
Initial page rank for page 2 : 0.25
Initial page rank for page 3 : 0.25
Initial page rank for page 4 : 0.25
For iteration 1 page rank for node : 1 is = 0.25
For iteration 1 page rank for node : 2 is = 0.125
For iteration 1 page rank for node : 3 is = 0.5
For iteration 1 page rank for node : 4 is = 0.125
For iteration 2 page rank for node : 1 is = 0.3125
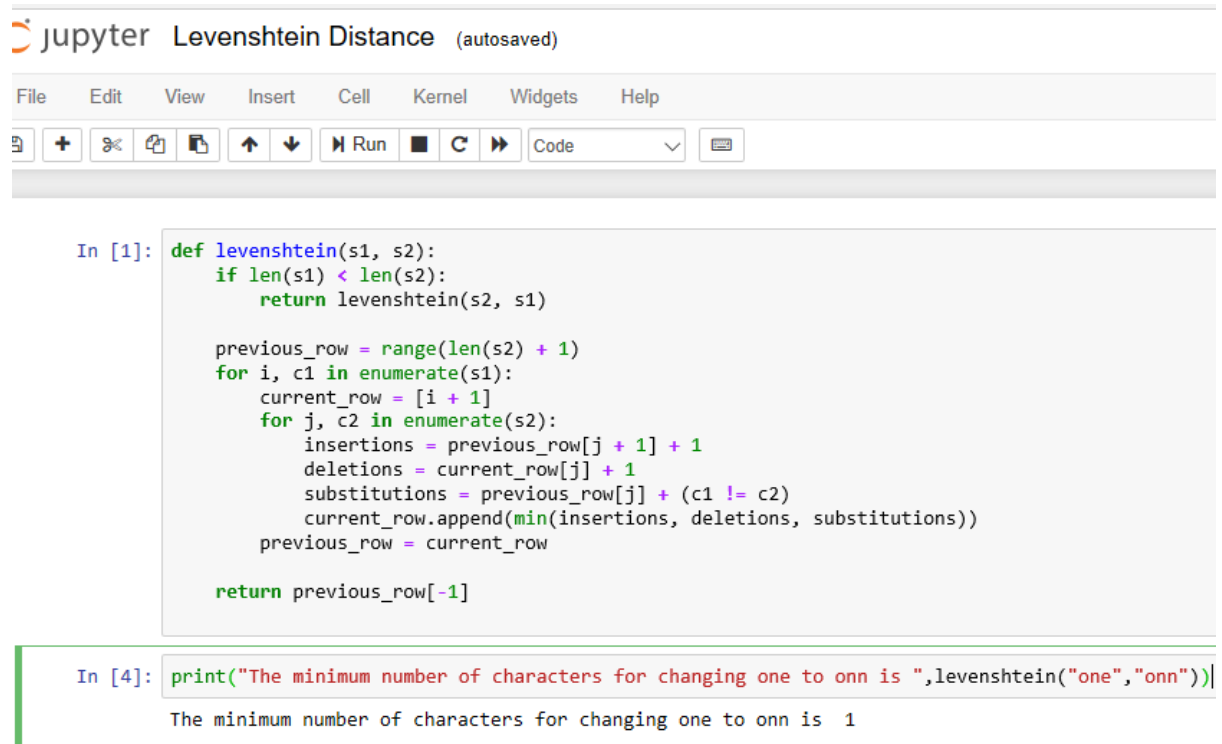For iteration 2 page rank for node : 2 is = 0.25
For iteration 2 page rank for node : 3 is = 0.3125
For iteration 2 page rank for node : 4 is = 0.125
BUILD SUCCESSFUL (total time: 43 seconds)

# Practical 3

Comupute levenshtein distance

Jupyter Levenshtein Distance (autosaved)

File    Edit    View    Insert    Cell    Kernel    Widgets    Help

Code

```python
In [1]: def levenshtein(s1, s2):
    if len(s1) < len(s2):
        return levenshtein(s2, s1)

    previous_row = range(len(s2) + 1)
    for i, c1 in enumerate(s1):
        current_row = [i + 1]
        for j, c2 in enumerate(s2):
            insertions = previous_row[j + 1] + 1
            deletions = current_row[j] + 1
            substitutions = previous_row[j] + (c1 != c2)
            current_row.append(min(insertions, deletions, substitutions))
        previous_row = current_row

    return previous_row[-1]
```

```python
In [4]: print("The minimum number of characters for changing one to onn is ",levenshtein("one","onn"))

The minimum number of characters for changing one to onn is  1
```

# Practical 4

Coumpute Cosine Similarity

Jupyter  Doc1.txt✔  a few seconds ago        Jupyter  Doc2.txt✔  a few seconds ago

File    Edit    View    Language              File    Edit    View    Language

1  Apples are tasty,Apples are yummy       1  Apples are yummy, Apples are red

Jupyter  Cosine Similarity  (autosaved)

File    Edit    View    Insert    Cell    Kernel    Widgets    Help

💾  +  ✂  ⎘  ⎘  ↑  ↓  ▶ Run  ■  C  ⏩    Code    ⌨

In [2]:

```python
from sklearn.feature_extraction.text import CountVectorizer
import pandas as pd
from sklearn.metrics.pairwise import cosine_similarity

f = open('Doc1.txt')
doc1 = str(f.read())
#doc1="Apples are tasty,Apples are yummy"
f = open('Doc2.txt')
doc2 = str(f.read())
#doc2="Apples are yummy, Apples are red"

documents = [doc1, doc2]

count_vectorizer = CountVectorizer()
sparse_matrix = count_vectorizer.fit_transform(documents)

#print(sparse_matrix)

doc_term_matrix = sparse_matrix.todense()

#print(doc_term_matrix)
df = pd.DataFrame(doc_term_matrix,
                  columns=count_vectorizer.get_feature_names(),
                  index=['doc1', 'doc2'])
print(df)

print(cosine_similarity(df[0:1], df))
```
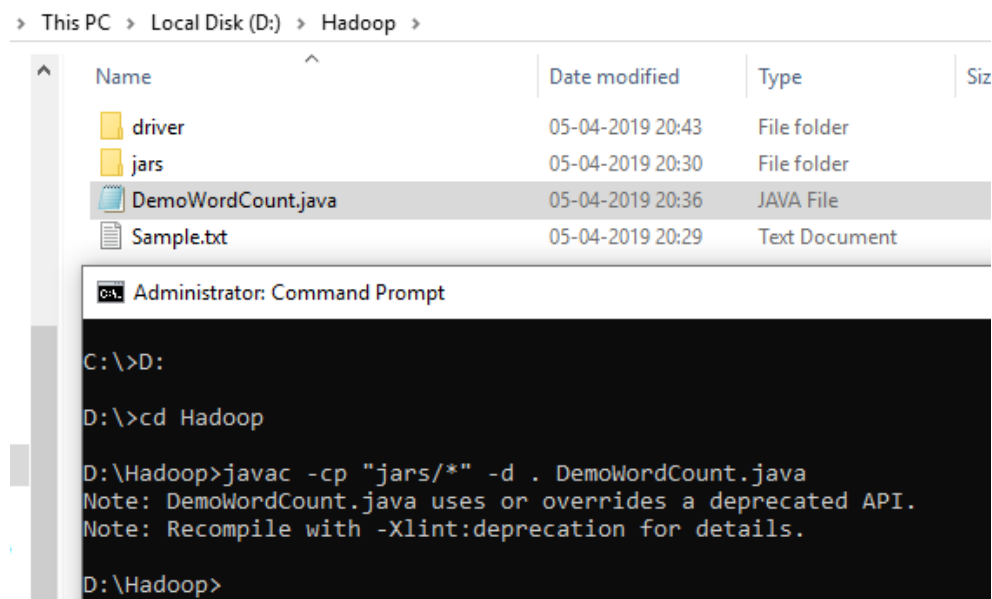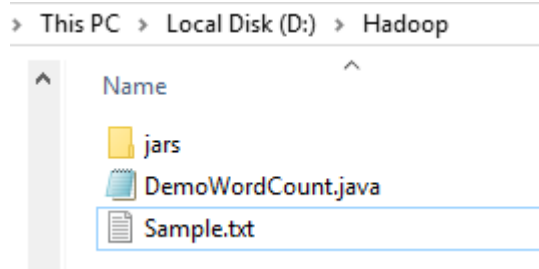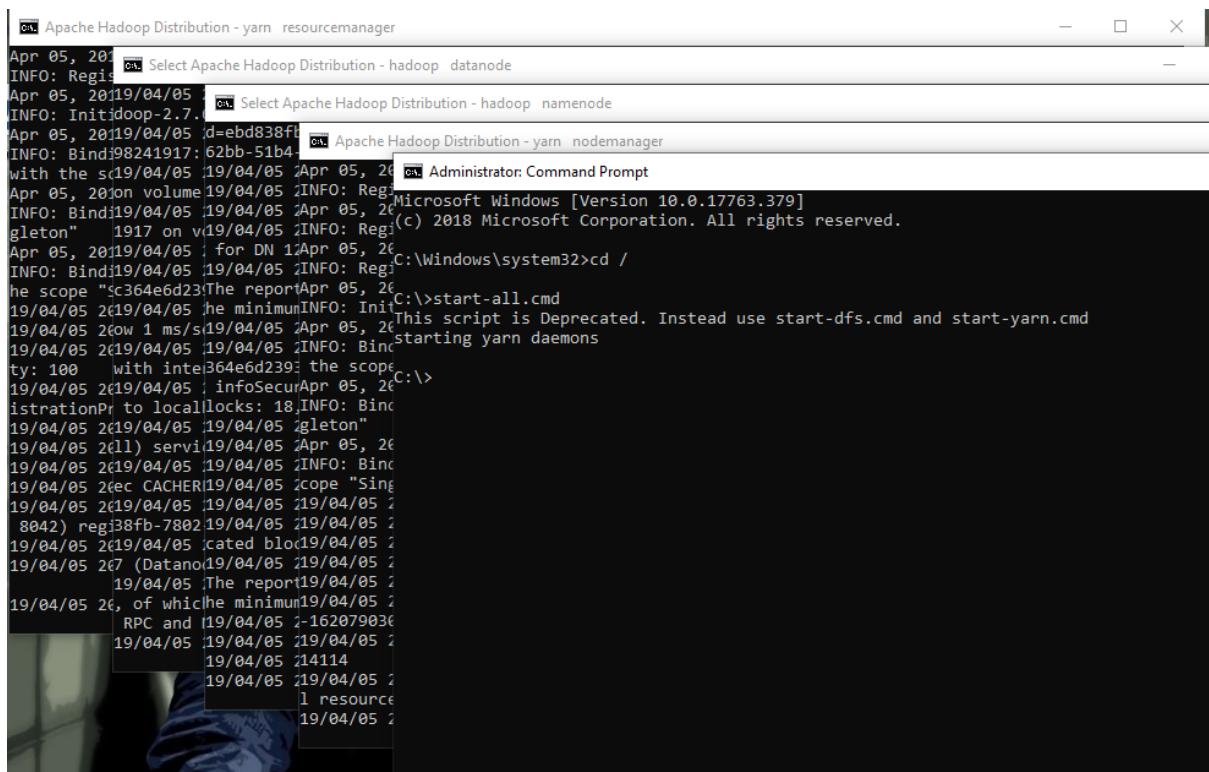
```
      apples  are  red  tasty  yummy
doc1       2    2    0      1      1
doc2       2    2    1      0      1
[[1.  0.9]]
```

# Practical 5

## Simple map reduce program

This PC › Local Disk (D:) › Hadoop

Name

- driver
- jars
- democount.jar
- DemoWordCount.java
- Sample.txt

```
D:\Hadoop>javac -cp "jars/*" -d . DemoWordCount.java
Note: DemoWordCount.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.

D:\Hadoop>jar -cvf democount.jar driver
added manifest
adding: driver/(in = 0) (out= 0)(stored 0%)
adding: driver/DemoWordCount$Map.class(in = 1760) (out= 752)(deflated 57%)
adding: driver/DemoWordCount$Reduce.class(in = 1669) (out= 707)(deflated 57%)
adding: driver/DemoWordCount.class(in = 1851) (out= 922)(deflated 50%)

D:\Hadoop>
```

```
D:\Hadoop>hadoop fs -mkdir /inputcount
```

```
D:\Hadoop>hadoop fs -put Sample.txt /inputcount
```

Administrator: Command Prompt

```
D:\Hadoop>hadoop jar democount.jar driver.DemoWordCount /inputcount/Sample.txt /outputcount
19/04/05 21:20:28 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
19/04/05 21:20:30 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement
ol interface and execute your application with ToolRunner to remedy this.
19/04/05 21:20:31 INFO input.FileInputFormat: Total input paths to process : 1
19/04/05 21:20:31 INFO mapreduce.JobSubmitter: number of splits:1
19/04/05 21:20:32 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1554478113696_0002
19/04/05 21:20:32 INFO impl.YarnClientImpl: Submitted application application_1554478113696_0002
19/04/05 21:20:32 INFO mapreduce.Job: The url to track the job: http://DESKTOP-JBVMG56:8088/proxy/application_155
696_0002/
19/04/05 21:20:32 INFO mapreduce.Job: Running job: job_1554478113696_0002
19/04/05 21:20:55 INFO mapreduce.Job: Job job_1554478113696_0002 running in uber mode : false
19/04/05 21:20:56 INFO mapreduce.Job:  map 0% reduce 0%
19/04/05 21:21:11 INFO mapreduce.Job:  map 100% reduce 0%
19/04/05 21:21:33 INFO mapreduce.Job:  map 100% reduce 100%
19/04/05 21:21:35 INFO mapreduce.Job: Job job_1554478113696_0002 completed successfully
19/04/05 21:21:35 INFO mapreduce.Job: Counters: 49
        File System Counters
                FILE: Number of bytes read=431
                FILE: Number of bytes written=248303
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=315
                HDFS: Number of bytes written=234
                HDFS: Number of read operations=6
                HDFS: Number of large read operations=0
```

```
D:\Hadoop>hadoop fs -ls /outputcount
Found 2 items
-rw-r--r--   1 Manasvi supergroup          0 2019-04-05 21:21 /outputcount/_SUCCESS
-rw-r--r--   1 Manasvi supergroup        234 2019-04-05 21:21 /outputcount/part-r-00000

D:\Hadoop>hadoop fs -cat /outputcount/part-r-00000
Program 1
able    2
all     1
and     2
are     1
because 1
```

```java
package driver;
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.fs.Path;

public class DemoWordCount
{

public static class Map extends Mapper<LongWritable,Text,Text,IntWritable>
{
public void map(LongWritable key,Text value,Context ctx) throws IOException,InterruptedException
        {
        Text word=new Text();
        IntWritable one=new IntWritable(1);

        String line=value.toString();

        StringTokenizer tokenizer=new StringTokenizer(line);

        while(tokenizer.hasMoreTokens())
                {
                        String w=tokenizer.nextToken();
                        word.set(w.toLowerCase());
                        ctx.write(word,one);
                }


        }
}


public static class Reduce extends Reducer<Text,IntWritable,Text,IntWritable>
{
        public void reduce(Text key,Iterable<IntWritable> values,Context ctx) throws IOException,InterruptedException
        {
                int sum=0;

                for(IntWritable x:values)
                {
                        int i=x.get();
                        sum=sum+i;
                }
                IntWritable finalCount=new IntWritable(sum);

                ctx.write(key,finalCount);
        }
}
```

```java
public static void main(String[] args) throws Exception
    {
            Configuration conf= new Configuration();

            Job job = new Job(conf,"My Word Count Program");
            job.setJarByClass(DemoWordCount.class);
            job.setMapperClass(Map.class);
            job.setReducerClass(Reduce.class);
            job.setOutputKeyClass(Text.class);
            job.setOutputValueClass(IntWritable.class);
            job.setInputFormatClass(TextInputFormat.class);
            job.setOutputFormatClass(TextOutputFormat.class);

            Path outputPath = new Path(args[1]);
            FileInputFormat.addInputPath(job, new Path(args[0]));
            FileOutputFormat.setOutputPath(job, new Path(args[1]));
            outputPath.getFileSystem(conf).delete(outputPath);

            //exiting the job only if the flag value becomes false
            System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}
```

# Practical 6

## IR system using lucene

```java
1   package UsingLucene;
2   import java.io.File;
3   import java.io.FileReader;
4   import java.io.IOException;
5   import org.apache.lucene.analysis.Analyzer;
6   import org.apache.lucene.analysis.standard.StandardAnalyzer;
7   import org.apache.lucene.util.Version;
8   import org.apache.lucene.*;
9   import org.apache.lucene.document.Document;
10  import org.apache.lucene.document.StoredField;
11  import org.apache.lucene.document.TextField;
12  import org.apache.lucene.index.DirectoryReader;
13  import org.apache.lucene.index.IndexReader;
14  import org.apache.lucene.index.IndexWriter;
15  import org.apache.lucene.index.IndexWriterConfig;
16  import org.apache.lucene.search.IndexSearcher;
17  import org.apache.lucene.search.Query;
18  import org.apache.lucene.search.ScoreDoc;
19  import org.apache.lucene.search.TopDocs;
20  import org.apache.lucene.store.Directory;
21  import org.apache.lucene.store.FSDirectory;
22  import org.apache.lucene.util.QueryBuilder;
23  public class UsingLucene {
24
25      static int index(File i, File d) throws IOException{
26
27
28          Analyzer analyzer = new StandardAnalyzer(Version.LUCENE_46);
29          IndexWriterConfig config=new IndexWriterConfig(Version.LUCENE_46,analyzer);
30
31          IndexWriter iWrite = new IndexWriter(FSDirectory.open(i),config);
32
33          File []files=d.listFiles();
34          for(File f:files){
35              System.out.println(f.getName());
36              Document doc=new Document();
37              doc.add(new TextField("content",new FileReader(f)));
38              doc.add(new StoredField("fileName",f.getCanonicalPath()));
39              iWrite.addDocument(doc);
40          }
41          int indexes=iWrite.maxDoc();
42          iWrite.close();
43          return indexes;
44      }
45
46      static void search(File f,String q) throws IOException{
47
48
49          Directory directory = FSDirectory.open(f);
50
51          IndexReader  indexReader  = DirectoryReader.open(directory);
```

```
53          IndexSearcher searcher = new IndexSearcher(indexReader);
54
55          Analyzer analyzer = new StandardAnalyzer(Version.LUCENE_46);
56
57          QueryBuilder builder = new QueryBuilder(analyzer);
58          Query query = builder.createPhraseQuery("content", q);
59
60          TopDocs topDocs =searcher.search(query, 3);
61          ScoreDoc[] hits = topDocs.scoreDocs;
62
            for (int i = 0; i < hits.length; i++) {
64              int docId = hits[i].doc;
65              Document d = searcher.doc(docId);
66              System.out.println(d.get("fileName") + " Score :"+hits[i].score);
67                                  }
68          System.out.println("Found " + hits.length);
69      }
70
71      public static void main(String []args) throws Exception{
72          File indexFile = new File("D:\\Books To Refer\\pracs\\pracs\\Lucene Indexes");
73          File dirFile = new File("D:\\Books To Refer\\pracs\\pracs\\ToBeIndexed");
74          //int numIndexes = index(indexFile,dirFile);
75          //System.out.println("Total files indexed " + numIndexes);
76          search(indexFile,"Tom and jerry");
77      }
78  }
79
```

**Information1.txt - Notepad**

File Edit Format View Help

Tom and jerry are running around
Tom fell on the ground and jerry made the laughing sound
Tom and jerry are playing around
Tom and jerry are chased by bruno

**information2.txt - Notepad**

File Edit Format View Help

Tom is dirnking millk
Jerry is eating chesse
Tom and jerry are dinning silently

**information3.txt - Notepad**

File Edit Format View Help

Tom is studying
jerry is singing
Tom and jerry are about to fight.
Tom and jerry are walking

```
run:
D:\Books To Refer\pracs\pracs\ToBeIndexed\Information1.txt Score :0.5397748
D:\Books To Refer\pracs\pracs\ToBeIndexed\information3.txt Score :0.5036848
D:\Books To Refer\pracs\pracs\ToBeIndexed\information2.txt Score :0.4451987
Found 3
BUILD SUCCESSFUL (total time: 0 seconds)
```

## Practical 7

Program to remove stop words

Jupyter  Sample.txt✔ 9 minutes ago

File    Edit    View    Language

```
1  Apples are tasty,Apples are yummy
2
```

Jupyter  Stopword Removal Last Checkpoint: 9 minutes ago  (unsaved changes)

File    Edit    View    Insert    Cell    Kernel    Widgets    Help

💾  +  ✂  📋  📋  ↑  ↓  ▶ Run  ■  C  ⏩    Code    ⌄    ⌨

In [2]:
```python
import io
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
#import nltk
#nltk.download('stopwords')

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

file1 = open("Sample.txt")
line = file1.read()    #to stream the file content
words = line.split()
for r in words:
    if not r in stop_words:
        appendFile = open('FilteredSample.txt','a')
        appendFile.write(" "+r)
        appendFile.close()
file2=open("FilteredSample.txt")
print(file2.read())
```

```
 Apples tasty,Apples yummy
```

Jupyter  FilteredSample.txt✔ a minute ago

File    Edit    View    Language

```
1    Apples tasty,Apples yummy
```

# Practical 8

## Program for web crawler

```java
import org.jsoup.Jsoup;
import org.jsoup.nodes.Document;
import org.jsoup.nodes.Element;
import org.jsoup.select.Elements;
import java.io.IOException;
import java.util.HashSet;
public class WebCrawler {
    private HashSet<String> links;
    public WebCrawler() {
        links = new HashSet<String>();
    }
    public void getPageLinks(String URL) throws Exception {
        //4. Check if you have already crawled the URLs
        //(we are intentionally not checking for duplicate content in this example)
        if (!links.contains(URL)) {
            //4. (i) If not add it to the index
            if (links.add(URL)) {
                System.out.println(URL);
            }
            //2. Fetch the HTML code
            Document document = Jsoup.connect(URL).get();
            //3. Parse the HTML to extract links to other URLs
            Elements linksOnPage = document.select("a[href]");

            //5. For each extracted URL... go back to Step 4.
            for (Element page : linksOnPage) {
                getPageLinks(page.attr("abs:href"));
            }
        }
    }
    public static void main(String[] args) throws Exception {
        //1. Pick a URL from the frontier
        new WebCrawler().getPageLinks("https://www.google.com/");
    }
}
```

```
https://www.google.com/
https://mail.google.com/mail/?tab=wm
https://www.google.co.in/imghp?hl=en&tab=wi
https://www.google.co.in/intl/en/about/products?tab=ih
https://about.google/
https://about.google/products/
https://about.google/stories/
https://about.google/stories/pedalingforpeace/
http://www.facebook.com/sharer.php?u=https://www.google.co.in/intl/ALL_in
https://www.facebook.com/recover/initiate/?ars=facebook_login
https://www.facebook.com/login/identify/?ctx=recover&ars=facebook_login#
https://www.facebook.com/
https://www.facebook.com/#
```