

CSCI 5408

**DATA MANAGEMENT AND
WAREHOUSING**

Assignment - 2

Gitlab Link: https://git.cs.dal.ca/jspatel/csci5408_s24_b00982253_jay_patel

References

Problem 1A.....	3
Algorithm.....	3
FlowChart.....	5
Data Cleaning/Transformation.....	6
MongoDB Connection.....	8
Output.....	9
Problem 1B.....	11
Algorithm.....	11
Program.....	14
Problem 2.....	18
Entities.....	18
Insert Data.....	19
Relationships.....	22
Final Output Graph.....	24
Problem 3.....	25
Algorithm.....	25
MongoDB Connection.....	30
Output.....	30
References.....	34

Problem 1A

Algorithm

Initialize MongoDB Connection:

- Create a MongoClient instance using a connection URI.
- Connect to the ReuterDb database and the News collection.

Read and Parse the News File:

- Create an instance of FileOperations using a singleton pattern.
- Read the raw data from the specified file (reut2-014.sgm) into a StringBuilder.

Extract and Clean Data:

- Use regular expressions to extract individual news articles (<REUTERS> tags) from the raw data.
- For each extracted article, extract the TITLE and BODY tags.
- Clean the extracted TITLE and BODY content by removing special characters and trimming whitespace.

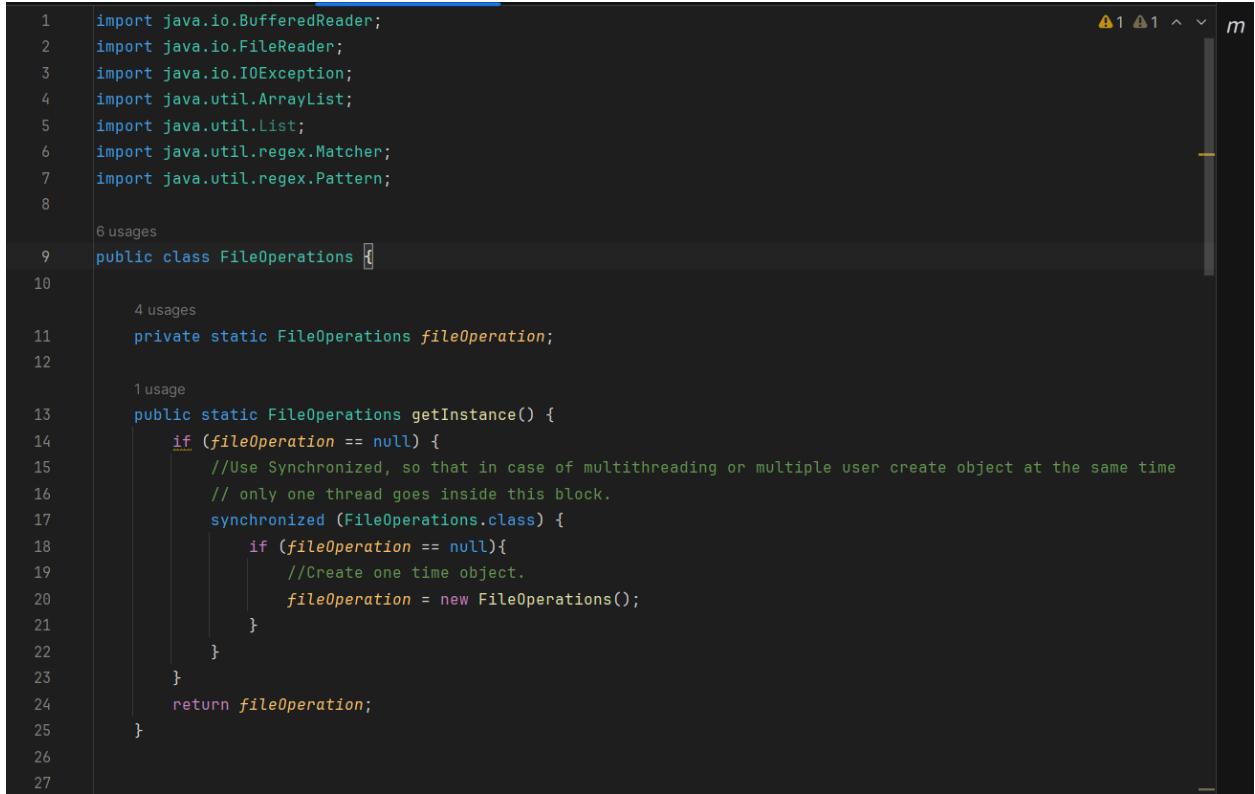
Convert Data to MongoDB Documents:

- Convert the cleaned TITLE and BODY content into News objects.
- Create a list of Document objects from the News objects.

Insert Data into MongoDB:

- Insert the list of Document objects into the News collection in MongoDB.
- Print a success message upon successful insertion.

The main method reads news articles from the specified file (reut2-014.sgm) using the FileOperations class, which implements a singleton pattern to ensure only one instance is created.



The screenshot shows a Java code editor with a dark theme. The code is a Singleton pattern implementation for a `FileOperations` class. The code includes imports for `java.io.BufferedReader`, `java.io.FileReader`, `java.io.IOException`, `java.util.ArrayList`, `java.util.List`, `java.util.regex.Matcher`, and `java.util.regex.Pattern`. The class has a private static field `fileOperation` and a static method `getInstance()` that returns it. The `getInstance()` method uses a synchronized block to ensure thread safety. The code editor shows 6 usages of the class and 4 usages of the field.

```
1 import java.io.BufferedReader;
2 import java.io.FileReader;
3 import java.io.IOException;
4 import java.util.ArrayList;
5 import java.util.List;
6 import java.util.regex.Matcher;
7 import java.util.regex.Pattern;
8
9 public class FileOperations {
10
11     private static FileOperations fileOperation;
12
13     public static FileOperations getInstance() {
14         if (fileOperation == null) {
15             //Use Synchronized, so that in case of multithreading or multiple user create object at the same time
16             // only one thread goes inside this block.
17             synchronized (FileOperations.class) {
18                 if (fileOperation == null){
19                     //Create one time object.
20                     fileOperation = new FileOperations();
21                 }
22             }
23         }
24         return fileOperation;
25     }
26
27 }
```

Figure 1: Singleton `FileOperations` class.

The `getNews` method in `FileOperations` reads the file content into a `StringBuilder` and then uses regular expressions to extract the content within `<REUTERS>` tags. Each extracted article's `TITLE` and `BODY` tags are further parsed, cleaned, and stored in `News` objects. Once the articles are parsed and cleaned, the `insertDataToDB` method converts these `News` objects into `MongoDB Document` objects and inserts them into the `News` collection in the `ReuterDb` database. The design encapsulates file reading, data extraction, cleaning, and database insertion in a structured and efficient manner.

FlowChart

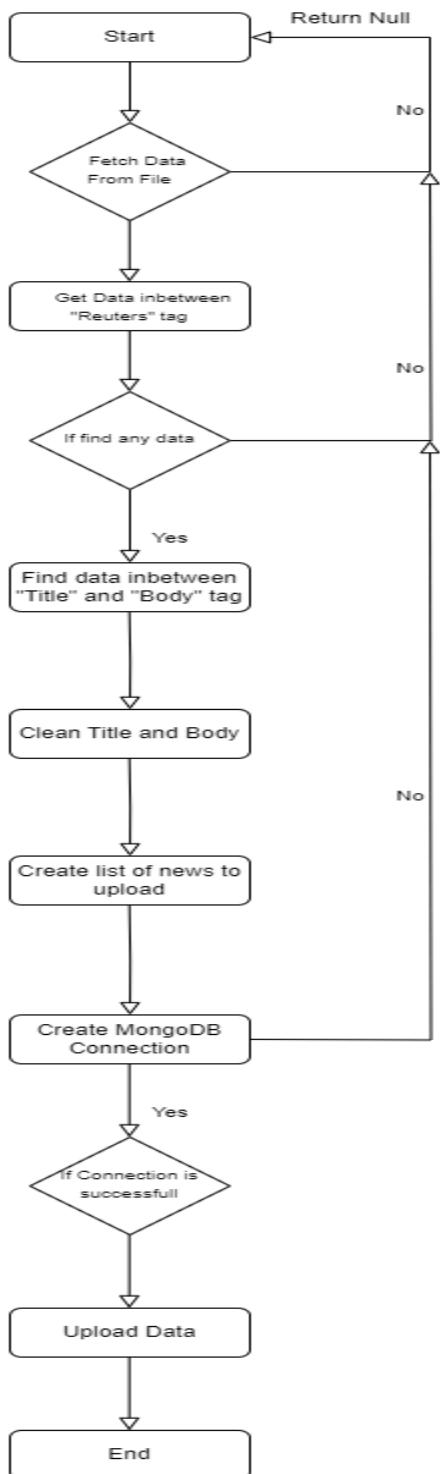


Figure 2: FlowChart of Problem 1A.

Data Cleaning/Transformation

1. Extraction of Relevant Data

- I used <REUTERS.*?>(.*)</REUTERS> regex pattern matches the entire content enclosed within <REUTERS> and </REUTERS> tags. The .*? makes the match non-greedy, ensuring that it captures the shortest possible string that fits the pattern.

```
private List<String> getInfoFromRawData(StringBuilder content) {  
    List<String> results = new ArrayList<>();  
    String regex = "<REUTERS.*?>(.*)</REUTERS>";  
    Pattern pattern = Pattern.compile(regex);  
    Matcher matcher = pattern.matcher(content);  
    while (matcher.find()) {  
        results.add(matcher.group(1));  
    }  
    return results;  
}
```

Figure 3: Extract data between “REUTERS” tag.

2. Parsing Title and Body

- I used <TITLE.*?>(.*)</TITLE> this pattern extracts content between <TITLE> and </TITLE> tags, similarly using .*? for a non-greedy match.
- Moreover, I used the pattern <BODY.*?>(.*)</BODY> to extract content between <BODY> and </BODY> tags, similarly using .*? for a non-greedy match.

```
private String getTitleAndBody(String content, String tagName) {  
    String regex = "<" + tagName + ".*?>(.*)</" + tagName + ">";  
    Pattern pattern = Pattern.compile(regex);  
    Matcher matcher = pattern.matcher(content);  
    if (matcher.find()) {  
        return matcher.group(1);  
    }  
    return null;  
}
```

Figure 4: Extract title and body from each news.

3. Cleaning Data

- Removing Special Characters:
 - Used this pattern [^a-zA-Z0-9\s] for matching any character that is not a letter, digit, or whitespace. It removes all such characters to ensure that only alphanumeric characters and spaces are retained.

- Replacing HTML Entities:
 - < and > is replaced with an empty string.
 - This step removes specific HTML entities that could appear in the text.
- Trimming Whitespace:
 - I used this pattern \\s+ which matches one or more whitespace characters (spaces, tabs, newlines) and is replaced with a single space to consolidate multiple spaces into one.
 - I also used trim() to remove white spaces at from both ends.

```
private String cleanTitleAndBody(String content) {
    if(content == null)
        return content;
    return content.replaceAll( regex: "[^a-zA-Z0-9\\s]", replacement: "" )
        .replace( target: "&lt;", replacement: "" ).replace( target: ">", replacement: "" )
        .replaceAll( regex: "\\s+", replacement: " " ).trim();
}
```

Figure 5: Remove special characters and unnecessary spaces.

4. Transforming Data into Objects

- Conversion to News Objects:
 - The cleaned TITLE and BODY strings are used to create News objects.
 - Each News object holds the title and body of the news article in a structured format.

```
private List<News> filterData(List<String> info){

    List<News> response = new ArrayList<>();

    for (String currInfo : info) {
        News currNew = new News();
        String title = getTitleAndBody(currInfo, tagName: "TITLE");
        String body = getTitleAndBody(currInfo, tagName: "BODY");
        if (title != null || body != null) {
            currNew.setTitle(cleanTitleAndBody(title));
            currNew.setBody(cleanTitleAndBody(body));
            response.add(currNew);
        }
    }

    return response;
}
```

Figure 6: Create News from filtered data.

5. Preparation for MongoDB Insertion

- Conversion to MongoDB Document Objects:
- Document Creation:
 - Each news object is converted into a document with fields, such as title and body. This transformation prepares the data for insertion into the MongoDB collection.

```
private static void insertDataToDB(List<News> response){  
  
    List<Document> data = new ArrayList<>();  
  
    for (News currNew : response){  
        Document doc = new Document("title", currNew.getTitle())  
            .append("body", currNew.getBody());  
        data.add(doc);  
    }  
  
    collection.insertMany(data);  
    System.out.println("Data Inserted Successfully To MongoDB Database!!");  
}  
}
```

Figure 7: Insert data into MongoDB database.

MongoDB Connection

```
public class ReutReader {  
  
    2 usages  
    private static MongoClient mongoClient;  
    2 usages  
    private static MongoCollection<Document> collection;  
    2 usages  
    private static FileOperations fileOperation;  
  
    public static void main(String[] args) {  
        fileOperation = FileOperations.getInstance();  
        List<News> response = fileOperation.getNews( filepath, "reut2-014.sgm");  
  
        ConnectMongoDB();  
        MongoDatabase database = mongoClient.getDatabase( databaseName: "ReuterDb");  
        collection = database.getCollection( "News");  
        insertDataToDB(response);  
    }  
  
    1 usage  
    private static void ConnectMongoDB(){  
        MongoClientURI uri = new MongoClientURI("mongodb+srv://root:root@default-cluster.xn53b0c.mongodb.net/?retryWrites=true&w=majority&appName=default-cluster");  
        mongoClient = new MongoClient(uri);  
    }  
}
```

Figure 8: MongoDB database connection string and collection name.

- The `ConnectMongoDB` method establishes a connection to a MongoDB database using the `MongoClient` class with a connection URI. This URI specifies the MongoDB cluster's address, authentication credentials, and connection options.

Output

- Here I create database with name “ReuterDb” and collection name is “News”.

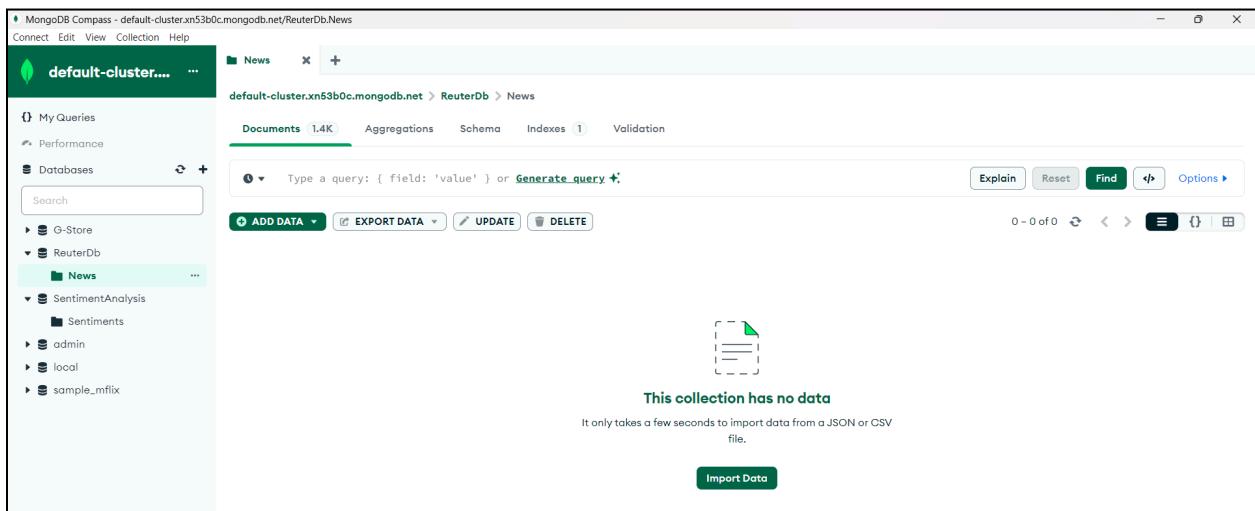


Figure 9: Before we run Java program.

```

for (News currNew : response) {
    Document doc = new Document("title", currNew.getTitle())
        .append("body", currNew.getBody());
    data.add(doc);
}
collection.insertMany(data);
System.out.println("Data Inserted Successfully To MongoDB Database!!!");
}
}

```

Aug 01, 2024 2:51:11 PM com.mongodb.diagnostics.logging.JULLogger log

INFO: Monitor thread successfully connected to server with description ServerDescription{address=ac-oelinyp-shard-00-01.xn53b0c.mongodb.net:27017, type=REPLICA_SET_SECONDARY, state=LIVE}

Aug 01, 2024 2:51:11 PM com.mongodb.diagnostics.logging.JULLogger log

INFO: Monitor thread successfully connected to server with description ServerDescription{address=ac-oelinyp-shard-00-00.xn53b0c.mongodb.net:27017, type=REPLICA_SET_SECONDARY, state=LIVE}

INFO: Setting max election id to 7fffffff000000000001ec from replica set primary ac-oelinyp-shard-00-02.xn53b0c.mongodb.net:27017

Aug 01, 2024 2:51:11 PM com.mongodb.diagnostics.logging.JULLogger log

INFO: Setting max set version to 130 from replica set primary ac-oelinyp-shard-00-02.xn53b0c.mongodb.net:27017

Aug 01, 2024 2:51:11 PM com.mongodb.diagnostics.logging.JULLogger log

INFO: Discovered replica set primary ac-oelinyp-shard-00-02.xn53b0c.mongodb.net:27017

Aug 01, 2024 2:51:12 PM com.mongodb.diagnostics.logging.JULLogger log

INFO: Opened connection [connectionId{localValue:4, serverValue:684707}] to ac-oelinyp-shard-00-02.xn53b0c.mongodb.net:27017

Data Inserted Successfully To MongoDB Database!

Process finished with exit code 0

Figure 10: Run Java Program.

MongoDB Compass - default-cluster.xn53b0c.mongodb.net/ReuterDb.News

Connect Edit View Collection Help

default-cluster.xn53b0c.mongodb.net > ReuterDb > News

Documents 1.4K Aggregations Schema Indexes 1 Validation

Type a query: { field: 'value' } or [Generate query](#)

[Explain](#) [Reset](#) [Find](#) [Options](#)

[ADD DATA](#) [EXPORT DATA](#) [UPDATE](#) [DELETE](#)

1–20 of 698

Document 1:
`_id: ObjectId("66abcb0f7388451ac777314b")
title: "JOHANNESBURG GOLD SHARES CLOSE MIXED TO FIRMER"
body: "Gold share prices closed mixed to slightly firmer in quiet and cautious..."`

Document 2:
`_id: ObjectId("66abcb0f7388451ac7773149")
title: "FEEDER CATTLE FUTURES SET NEW HIGHS TURN MIXED"
body: "Feeder cattle futures advanced 0.26 to 0.30 cent at the start and posted ..."`

Document 3:
`_id: ObjectId("66abcb0f7388451ac777314a")
title: "CBT CORN SPREADS"
body: "1000 hrs cdt MONTHS LAST DIFFERENCE JulMay 314 over 312 over DecJul 12..."`

Document 4:
`_id: ObjectId("66abcb0f7388451ac777314b")
title: "TEXAS PANHANDLE OKLA FEEDLOT ROUNDUP USDA"
body: "Cattle in the panhandle area Monday were 0.50 to 1.50 dlr higher. Trading ..."`

Document 5:
`_id: ObjectId("66abcb0f7388451ac777314c")
title: "MIDWEST GRAIN FUTURES 1100 EDT"
body: "MINNEAPOLIS WHEAT MAY7 284 34 UP 1 14 JUL7 280 14 OFF 12 SEPT 277 UP 1..."`

Document 6:
`_id: ObjectId("66abcb0f7388451ac777314d")
title: "VOLCKER PUSHES SPENDING CUTS OVER TRADE BILL"`

> MONGOSH

Figure 11: After successfully run the java program, data on MongoDB

Problem 1B

Algorithm

Initialize Spark Session

- Create a SparkSession instance with the application name "A2".

```
String appName = "A2";
SparkSession sparkSession = SparkSession.builder().appName(appName).getOrCreate();
```

Figure 12: Start spark session with “A2” app name.

Read Input Data

- Use the FileOperations class to read the main data file ("reut2-014.sgm") and the stop words file ("StopWords.txt") using the getData and getStopWords methods, respectively.
- Collect the stop words into a List<String>.

```
1 usage
public JavaRDD<String> getData(SparkSession sparkSession, String filePath){
    Dataset<String> fileData = sparkSession.read().option("multiline", false).textFile(filePath);
    //Print File to Console
    System.out.println("Content of File: ");
    fileData.show();
    return fileData.javaRDD();
}

1 usage
public JavaRDD<String> getStopWords(SparkSession sparkSession, String filePath){
    Dataset<String> fileData = sparkSession.read().option("multiline", false).textFile(filePath);
    //Print Stop Words to Console
    System.out.println("Stop Words: ");
    fileData.show();
    return fileData.javaRDD();
}
```

Figure 13: Fetch data from “reut2-014.sgm”, and “StopWords.txt”.

Clean and Tokenize Text

- Clean Text:
 - For each line in the main data file, replace all non-alphabetic characters (except spaces) with a space.
 - Convert the text to lowercase.
 - Trim leading and trailing whitespace.

- Split the cleaned text into individual words by whitespace.

```
JavaRDD<String> cleanedText = news.flatMap(line -> Arrays.asList(
    line.replaceAll( regex: "[^a-zA-Z\\s]", replacement: " ")
        .toLowerCase(Locale.ROOT)
        .trim()
        .split( regex: "\\s+")
).iterator());
```

Figure 14: Remove unnecessary spaces and convert text to lowercase.

Filter Words

- Remove Empty Strings:
 - Filter out any empty strings.
- Remove Stop Words:
 - Filter out words that are in the stop words list.
- Ensure Alphabetic Words:
 - Filter to keep only alphabetic words.

```
JavaRDD<String> words = cleanedText
    .filter(s -> !s.isEmpty())
    .filter(s -> !stopWordList.contains(s))
    .filter(s -> s.matches( regex: "[a-z]+"));

words = words.map(s -> s.replaceAll( regex: "<|>|[^\w]+", replacement: ""));
```

Figure 15: Compare text with stop words.

Additional Cleanup (if necessary)

- Remove Special Characters:
 - Further, clean up words by removing specific special characters.

```
words = words.map(s -> s.replaceAll( regex: "<|>|[^\w]+", replacement: ""));
```

Figure 16: Remove special characters.

Count Word Frequencies

- Map each word to a key-value pair where the key is the word and the value is 1.
- Aggregate the counts for each word using reduceByKey.

```
JavaPairRDD<String, Integer> wordCounts = words.mapToPair(word -> new Tuple2<>(word, 1))
    .reduceByKey((a, b) -> a + b);
```

Figure 17: Create pair for each word and set count default count to 1.

Output Word Frequencies:

- Collect and print the word counts in the format (word, count).

```
wordCounts.collect().foreach(  
    wordCount -> System.out.println("(" + wordCount._1() + ", " + wordCount._2() + ")");
```

Figure 18: Print each word on console.

Find and Print Maximum and Minimum Frequency Words.

- Find Maximum Frequency Word:
 - Swap the key-value pairs to facilitate sorting by frequency.
 - Sort by frequency in descending order and take the top result.
- Find Minimum Frequency Word:
 - Swap the key-value pairs to facilitate sorting by frequency.
 - Sort by frequency in ascending order and take the top result.
- Print the results for both maximum and minimum frequency words in the format (word, frequency).

```
Tuple2<Integer, String> maxFreqWord = wordCounts.mapToPair(Tuple2::swap) JavaPairRDD<Integer, String>  
    .sortByKey( ascending: false)  
    .take( num: 1) List<Tuple2<...>>  
    .get(0);  
  
Tuple2<Integer, String> minFreqWord = wordCounts.mapToPair(Tuple2::swap) JavaPairRDD<Integer, String>  
    .sortByKey( ascending: true)  
    .take( num: 1) List<Tuple2<...>>  
    .get(0);  
  
System.out.println(  
    "Word with maximum frequency : (" + maxFreqWord._2() + " , " + maxFreqWord._1() + ")");  
  
System.out.println(  
    "Word with minimum frequency : (" + minFreqWord._2() + " , " + minFreqWord._1() + ")");
```

Figure 19: Calculate the Max. and Min. Frequency of words.

Stop Spark Session

- Close the SparkSession.

```
sparkSession.stop();
```

Figure 20: Close spark session.

Program

- I created the cluster using the same method and configuration as I learned and followed in Lab 6 - BIG DATA: HADOOP AND APACHE SPARK, and I didn't put any screenshots for the cluster creation as the assignment didn't require that.
 - I create a cluster with the name “spark-freqcount”.
 - Upload .jar, StopWords.txt, and reut2-014.sgm files to cluster.
 - Here I show that the files are successfully uploaded and we can see them on the cluster using the “ls” command.

```
ssh.cloud.google.com/v2/ssh/projects/cloud-428816/zones/us-central1-a/instances/spark-freqcount-m?authuser=1&hl=en_US&projectN... — □ ×
ssh.cloud.google.com/v2/ssh/projects/cloud-428816/zones/us-central1-a/instances/spark-freqcount-m?authuser=1&hl=en_US&projec...
SSH-in-browser UPLOAD FILE DOWNLOAD FILE
[ ] [ ] [ ] [ ] [ ]
jayapatel41120@spark-freqcount-m:~$ ls
IO2-1.0-SNAPSHOT.jar StopWords.txt reut2-014.sgm
jayapatel41120@spark-freqcount-m:~$ █
```

Figure 21: Upload files on cluster:

- Now I run the program using the command “spark-submit –class Program.FreqCount IO2-1.0-SNAPSHOT.jar”.

```
ssh.cloud.google.com/v2/ssh/projects/cloud-428816/zones/us-central1-a/instances/spark-freqcount-m?authuser=1&hl=en_US&projectNumber=287287505722&useAdminProxy=true&pageViewId=AEEF84B02-A62C-4453-A1CE-3AE586C9BC95 - Google Chrome
ssh.cloud.google.com/v2/ssh/projects/cloud-428816/zones/us-central1-a/instances/spark-freqcount-m?authuser=1&hl=en_US&projectNumber=287287505722&useAdminProxy=true&pageViewId=AEEF84B02-A62C-4453-A1CE-3AE586C9BC95 - Google Chrome

SSH-in-browser
UPLOAD FILE DOWNLOAD FILE

jyat Patel@l1120s:spark-freqcount-m: $ ls
IO2-1.0-SNAPSHOT.jar StopWords.txt reut2-014.sgm
jyat Patel@l1120s:spark-freqcount-m: $ spark-submit --class Program.FreqCount IO2-1.0-SNAPSHOT.jar
24/08/01 16:39:40 INFO SparkEnv: Registering MapOutputTracker
24/08/01 16:39:40 INFO SparkEnv: Registering BlockManagerMaster
24/08/01 16:39:40 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
24/08/01 16:39:40 INFO SparkEnv: Registering OutputCommitCoordinator
24/08/01 16:39:41 INFO RMFailoverProxyProvider: Connecting to ResourceManager at spark-freqcount-m.us-central1-a.c.cloud-428816.internal/10.128.15.194:8032
24/08/01 16:39:41 INFO AdminProxy: Registered with Application Master server at spark-freqcount-m.us-central1-a.c.cloud-428816.internal/10.128.15.194:10200
24/08/01 16:39:43 INFO Configuration: resource-types.xml not found
24/08/01 16:39:43 INFO ResourceUtils: Unable to find 'resource-types.xml'
24/08/01 16:39:45 INFO YarnClientImpl: Submitted application application_1722528378513_0006
24/08/01 16:39:45 INFO DefaultHARMFailoverProxyProvider: Connecting to ResourceManager at spark-freqcount-m.us-central1-a.c.cloud-428816.internal/10.128.15.194:8030
24/08/01 16:39:46 INFO MetricsSystemConfig: Probing properties from hadoop-metrics2.properties
24/08/01 16:39:46 INFO MetricsSystemImpl: Scheduled Metric snapshot period at 10 seconds().
24/08/01 16:39:46 INFO MetricsSystemImpl: google-hadoop file-system metrics system started
24/08/01 16:39:47 INFO GlobalStorageStatistics: Detected potential high latency for operation op_get_file_status. latencyMs=503; previousMaxLatencyMs=0; operationCount=1; context=gs://data
```

Figure 22:Run java program on cluster.

- Here in the code, I print the content of “reut2-014.sgm”, and “StopWords.txt”, and also display words with the highest and lowest frequency.
 - Here I take all the stopwords from the GitHub repo sebleier/NLTK's list of english stopwords.

```

Content of File:
+-----+
|          value|
+-----+
|<!DOCTYPE lewis S...|
|<REUTERS TOPICS="...|
|<DATE> 7-APR-1987...|
|  <TOPICS></TOPICS>|
|  <PLACES></PLACES>|
|  <PEOPLE></PEOPLE>|
|  <ORGS></ORGS>|
|<EXCHANGES></EXCH...|
|<COMPANIES></COMP...|
|          <UNKNOWN>|
|          &#5;&#5;&#5;E|
|&#22;&#22;&#1;f11...|
|b f BC-JOHANNESBU...|
|          <TEXT>&#2;|
|<TITLE>JOHANNESBU...|
|<DATELINE>    JOH...|
|slightly firmer i...|
|reaction to a ret...|
|dlrs and a firmer...|
|  Heavyweight V...|
+-----+
only showing top 20 rows

```

Figure 23: Print content of file on console.

```

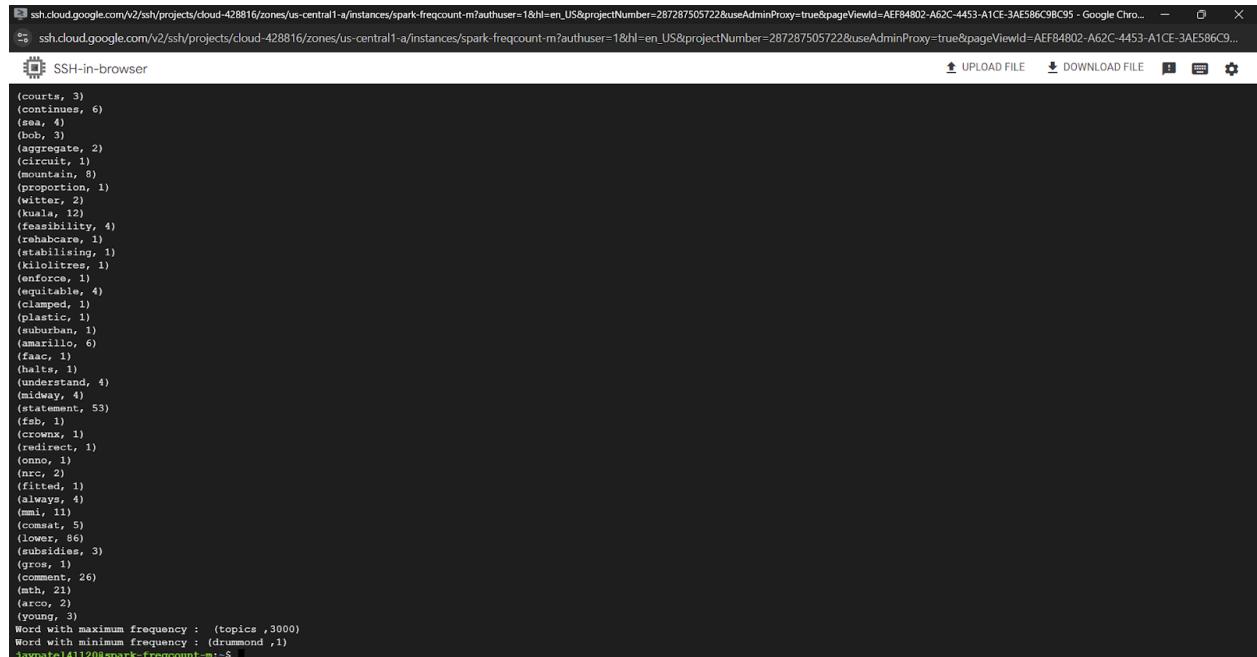
Stop Words:
+-----+
|      value|
+-----+
|      i|
|      me|
|      my|
|      myself|
|      we|
|      our|
|      ours|
|      ourselves|
|      you|
|      your|
|      yours|
|      yourself|
|      yourselves|
|      he|
|      him|
|      his|
|      himself|
|      she|
|      her|
|      hers|
+-----+
only showing top 20 rows

```

Figure 24: Print stop words on console.

```
(francs, 25)
(bone, 2)
(constraints, 3)
(drummond, 1)
(pentoxide, 1)
(yun, 1)
(brewery, 2)
(secure, 4)
(ogren, 4)
(dayton, 7)
(fred, 1)
(preventing, 1)
(end, 58)
(leeeway, 2)
(barring, 1)
(financials, 1)
(sumatran, 2)
(tough, 7)
(pig, 1)
(specia, 1)
(clients, 2)
(recurrence, 1)
(koido, 2)
(gnma, 2)
(less, 24)
(supporting, 5)
(mch, 15)
(paso, 1)
(supporter, 1)
(irresponsible, 1)
(adequately, 1)
(smooth, 1)
(plus, 19)
(grant, 5)
(different, 2)
(shape, 1)
(sarasota, 1)
(completing, 2)
(phillips, 2)
(jeddah, 1)
```

Figure 25: Show random words and its frequency.



The screenshot shows a terminal window titled "SSH-in-browser" with the URL "ssh.cloud.google.com/v2/ssh/projects/cloud-428816/zones/us-central1-a/instances/spark-freqcount-m?authuser=1&hl=en_US&projectNumber=287287505722&useAdminProxy=true&pageViewId=AEC84802-A62C-4453-A1CE-3AE586C9BC95 - Google Chrome". The window displays a list of words and their frequencies, starting with "courts, 3) (continues, 6) (ree, 4) (bob, 3) (aggregate, 2) (circuit, 1) (mountain, 8) (proportion, 1) (witter, 2) (kuala, 12) (feasibility, 4) (rehabcare, 1) (stabilising, 1) (kilometres, 1) (enforce, 1) (available, 4) (clamped, 1) (plastic, 1) (suburban, 1) (amarillo, 6) (fac, 1) (halts, 1) (understand, 4) (midway, 4) (statement, 53) (tsp, 1) (crownx, 1) (radiant, 1) (com, 1) (nrc, 2) (fitted, 1) (always, 4) (nmi, 11) (comsat, 5) (lower, 86) (subsidies, 3) (gros, 1) (comment, 26) (nth, 21) (arc, 2) (young, 3) Word with maximum frequency : (topics ,3000) Word with minimum frequency : (drummond ,1) jaypathel41120@spark-freqcount-m:~\$

Figure 26: Show last lines of output

```
Word with maximum frequency : (topics ,3000)
Word with minimum frequency : (drummond ,1)
jaypate141120@spark-freqcount-m:~$ █
```

Figure 27:Show Max. and Min. Frequency.

- The word with the Max. frequency is “topics” with a frequency of 3000.
- The word with the Min. frequency is “drummond” with a frequency of 1.

Problem 2

Entities

- From Assignment 1, I have four entities as follows:

Resources

- Resource_ID
- Name
- Type

Park

- Park_ID
- Name
- Location
- Area
- Type
- Opening Hours

Site

- Site_ID
- Site Number
- Location
- Capacity
- Rate
- Status

Reservation

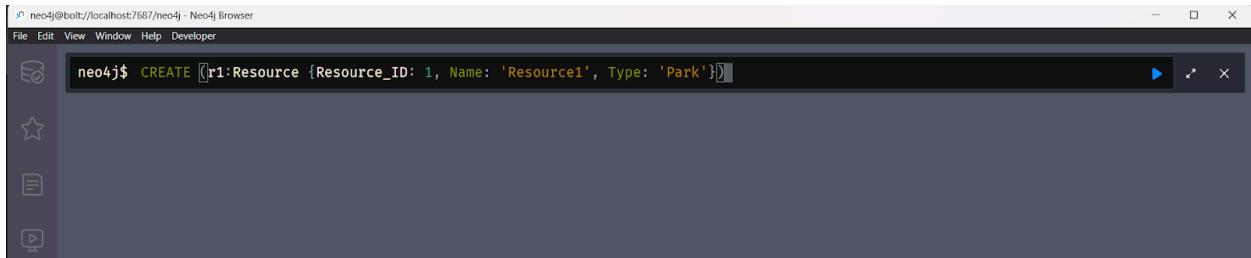
- Reservation_ID
- Date
- Time
- Status
- Payment Status

The relationships are:

- Resources have Parks.
- Parks have Sites.
- Sites are reserved by Reservations.

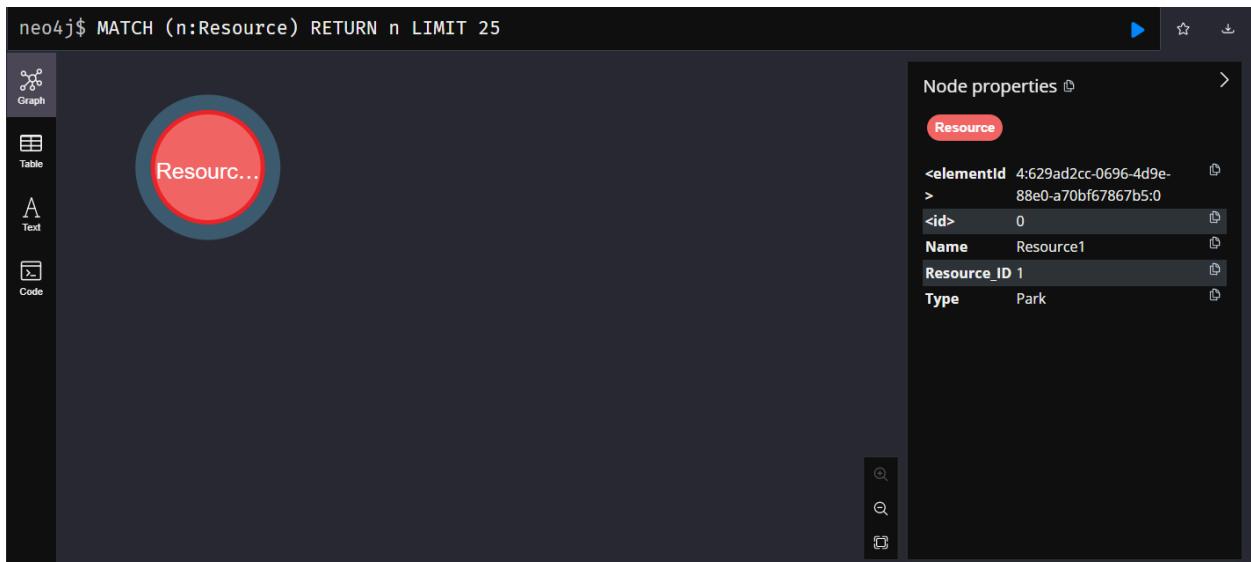
Insert Data

- Now first I insert data for all the four entities and than establish appropriate relationships.
- So here I insert data for resource entity.
- The Cypher query creates a node labeled Resource with properties Resource_ID set to 1, Name set to 'Resource1', and Type set to 'Park'.



A screenshot of the Neo4j Browser interface. The title bar says "neo4j@bolt://localhost:7687/neo4j - Neo4j Browser". The menu bar includes File, Edit, View, Window, Help, and Developer. The main area shows a Cypher command: "neo4j\$ CREATE (r1:Resource {Resource_ID: 1, Name: 'Resource1', Type: 'Park'})". Below the command is a toolbar with icons for Graph, Table, Text, and Code. The status bar at the bottom right shows "neo4j\$".

Figure 28: Insert data into resource entity.



A screenshot of the Neo4j Browser interface. The title bar says "neo4j\$ MATCH (n:Resource) RETURN n LIMIT 25". The main area displays a single red circular node labeled "Resource...". To the right is a panel titled "Node properties" showing the following data for the node:

Resource
<elementId> 4:629ad2cc-0696-4d9e-> 88e0-a70bf67867b5:0
<id> 0
Name Resource1
Resource_ID 1
Type Park

The left sidebar has tabs for Graph, Table, Text, and Code, with "Graph" selected. The status bar at the bottom right shows "neo4j\$".

Figure 29: Show resource entity.

- Insert data for Park entity.
- The Cypher query creates a node labeled 'Park' with properties: 'Park_ID' set to 1, 'Type' set to 'National Park', 'Location' set to 'Location1', 'Area' set to '5000 sq km', 'Name' set to 'Park1', and 'Opening_Hours' set to '9 AM - 5 PM'.

neo4j\$ CREATE [p1:Park {Park_ID: 1, Type: 'National Park', Location: 'Spring Garden Road Halifax', Area: '5000 sq km', Name: 'Spring Garden', Opening_Hours: '9 AM - 5 PM'}]

Figure 30: Insert data into park entity.

neo4j\$ MATCH (n:Park) RETURN n LIMIT 25

Node properties

Park
<elementId> 4:629ad2cc-0696-4d9e->
<> 88e0-a70bf67867b5:1
<id> 1
Area 5000 sq km
Location Spring Garden Road Halifax
Name Spring Garden
Opening_Hours 9 AM - 5 PM
Park_ID 1
Type National Park

Figure 31: Show park entity.

- Insert data for Site entity.
- The Cypher query creates a node labeled 'Site' with properties: 'Site_ID' set to 1, 'Site_Number' set to 101, 'Location' set to 'LocationA', 'Capacity' set to 50, 'Rate' set to 100, and 'Status' set to 'Available'.

neo4j\$ CREATE [s1:Site {Site_ID: 1, Site_Number: 101, Location: 'Left Corner', Capacity: 50, Rate: 100, Status: 'Available'}]

Figure 32: Insert data into site entity.

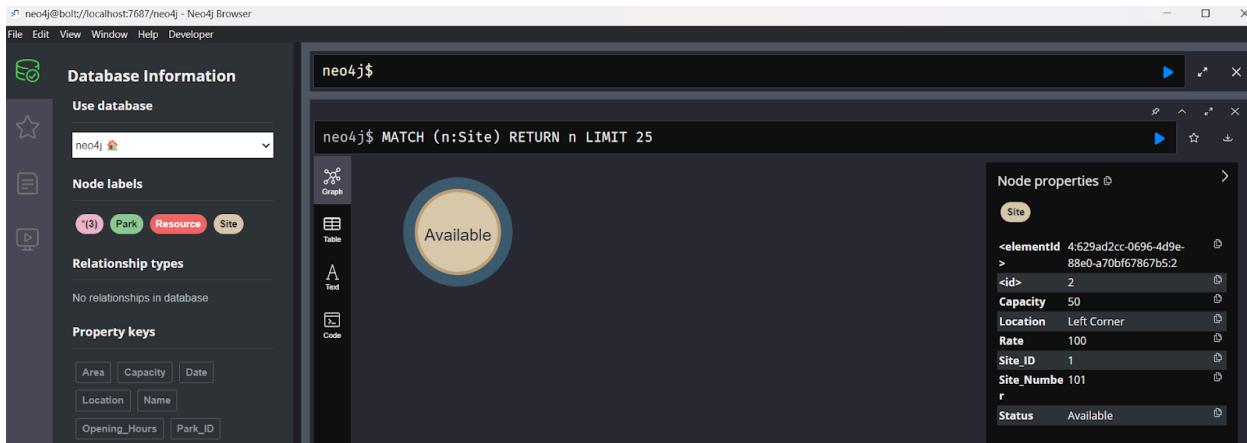


Figure 33: Show site entity.

- Insert data into Reservation entity.
- The Cypher query creates a node labeled 'Reservation' with properties: 'Reservation_ID' set to 1, 'Date' set to '2024-08-01', 'Time' set to '10:00', 'Status' set to 'Confirmed', and 'Payment_Status' set to 'Paid'.

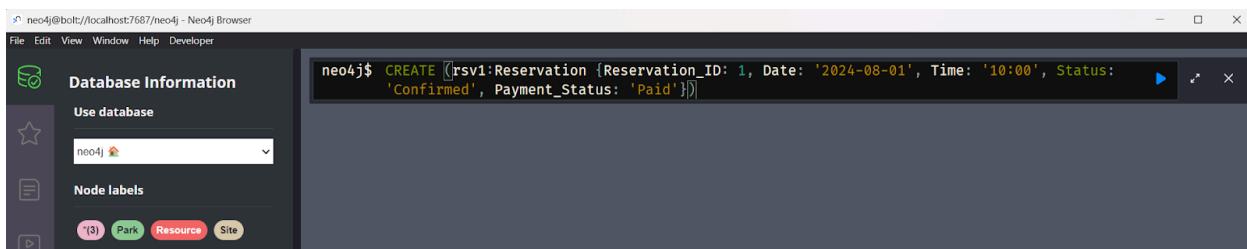


Figure 34: Insert into reservation entity.

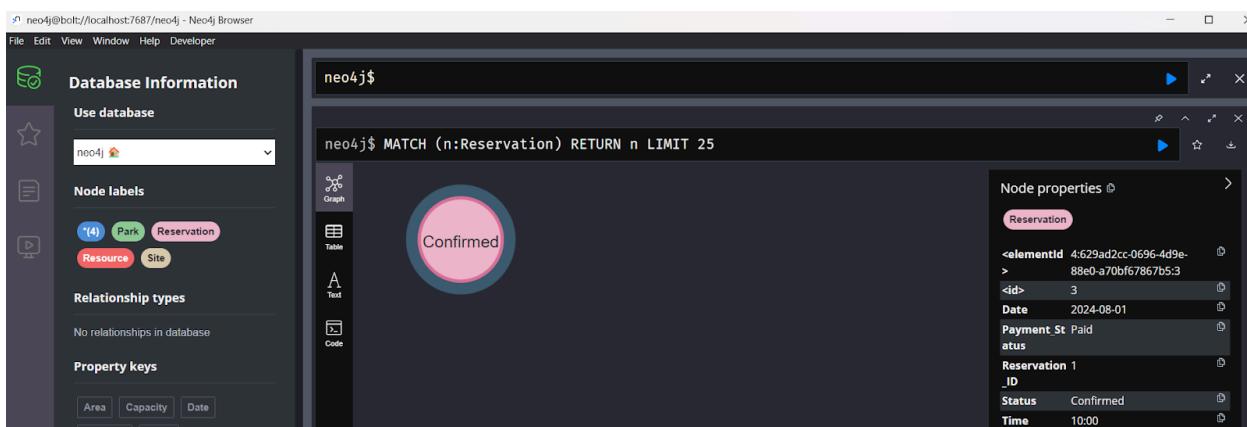
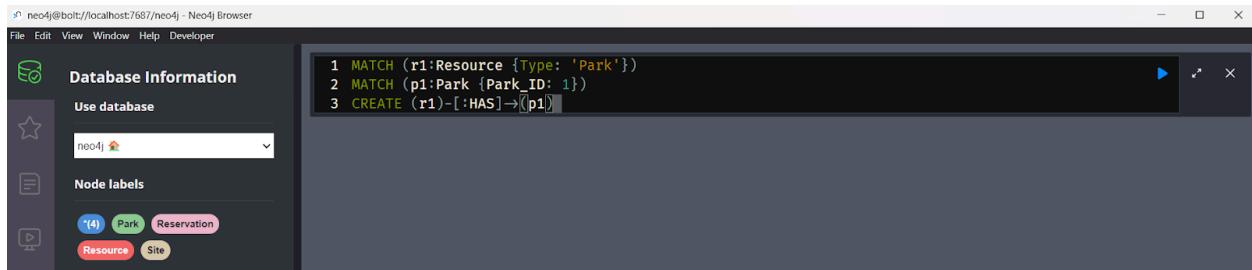


Figure 35: Show reservation entity.

Relationships

- Establish relationship between resource and park.
- So resources is handle by nova scotia.
- Resource entity has so many type like park, lake, and natural resources.
- Here we only consider park.



The screenshot shows the Neo4j Browser interface. On the left, the 'Database Information' sidebar shows the database 'neo4j' selected. Under 'Node labels', there are four nodes: Park (4), Reservation, Resource (1), and Site. Under 'Relationship types', there is one relationship type: HAS (1). In the main query editor area, the following Cypher code is written:

```
1 MATCH (r1:Resource {Type: 'Park'})  
2 MATCH (p1:Park {Park_ID: 1})  
3 CREATE (r1)-[:HAS]-(p1)
```

Figure 36: Establish relationship between resource and park.

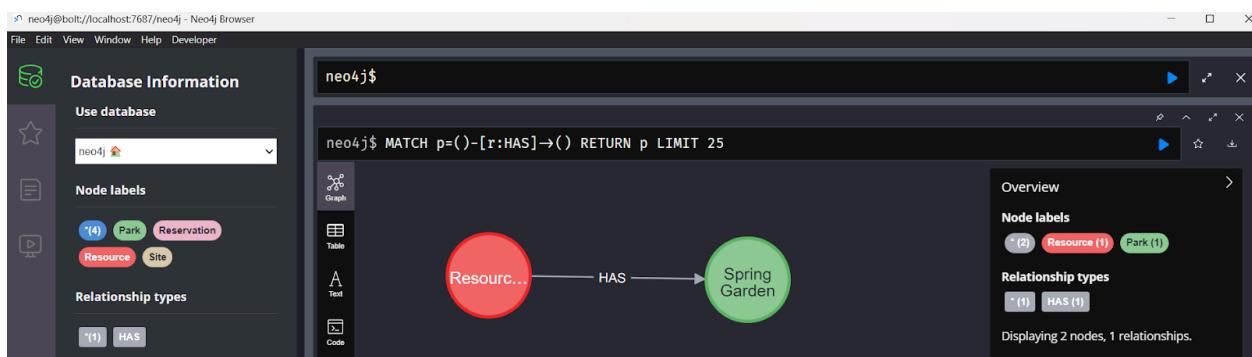
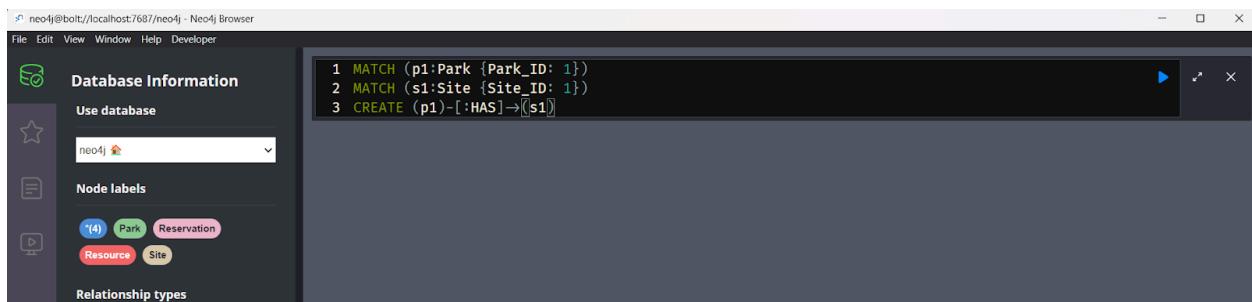


Figure 37: Show relationship between resource and park.

- Establish relationship between park and site.
- One park has many site on different locations.



The screenshot shows the Neo4j Browser interface. On the left, the 'Database Information' sidebar shows the database 'neo4j' selected. Under 'Node labels', there are four nodes: Park (4), Reservation, Resource (1), and Site. Under 'Relationship types', there is one relationship type: HAS (1). In the main query editor area, the following Cypher code is written:

```
1 MATCH (p1:Park {Park_ID: 1})  
2 MATCH (s1:Site {Site_ID: 1})  
3 CREATE (p1)-[:HAS]-(s1)
```

Figure 38: Establish relationship between park and site.

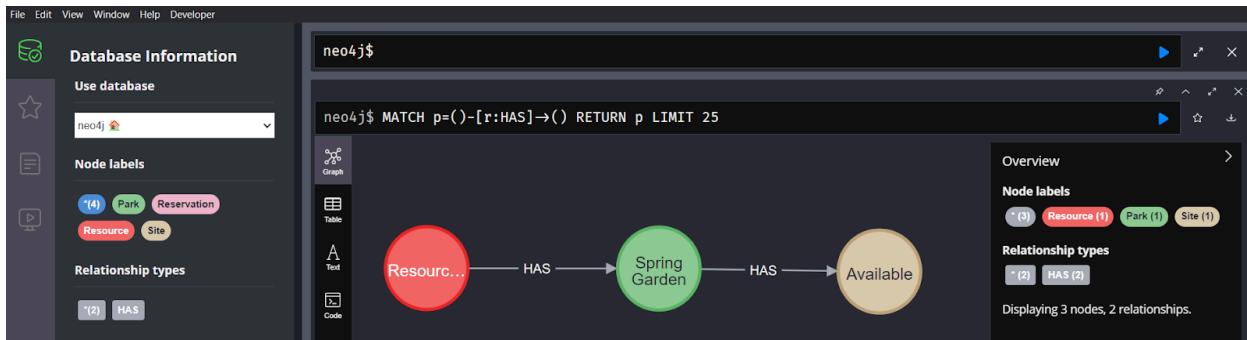


Figure 39: Show relationship between park and site.

- Establish relationship between site and reservation
- One site can be reserved by many people, but only one at a time.

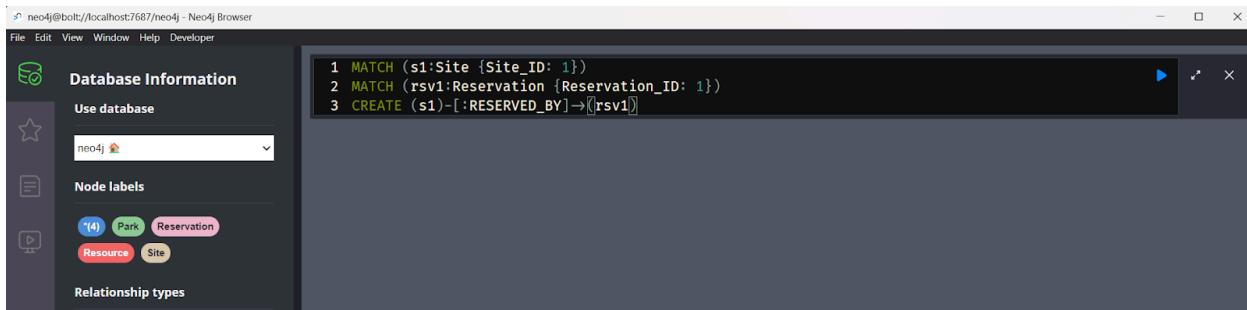


Figure 40: Establish relationship between site and reservation.

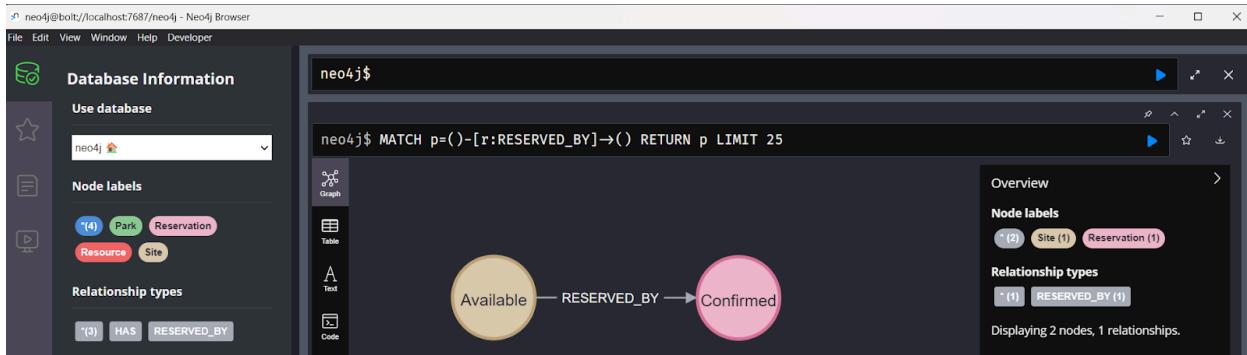


Figure 41: Show relationship between site and reservation.

Final Output Graph

- Here is the entire graph that we get from all the data we inserted and relationships we established.
- It is clearly show in graph that all the four entities are related to each other with proper relationships.

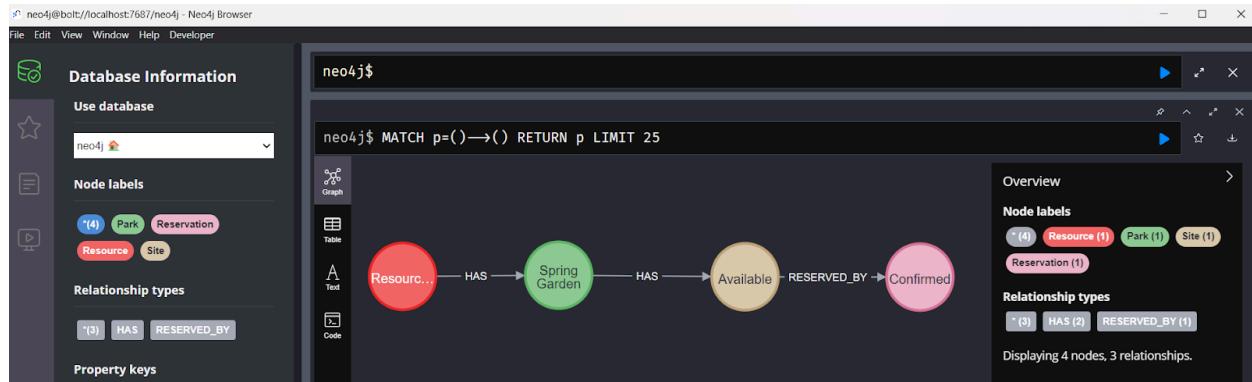


Figure 42: Show overall graph after final step.

Problem 3

Algorithm

Initialize MongoDB Connection

- Create a MongoClient and connect to MongoDB using a predefined URI.

Read and Process Data

- Read BOW Data:
 - Use FileOperations to read data from the file reut2-014.sgm.
 - Extract relevant information and filter it to get a list of string arrays.

```
1 usage
private StringBuilder readFile(String filepath){
    StringBuilder content = new StringBuilder();
    try (BufferedReader br = new BufferedReader(new FileReader(filepath))) {
        String line;
        while ((line = br.readLine()) != null) {
            content.append(line);
        }
        return content;
    } catch (IOException e) {
        System.out.println(e.getMessage());
    }
    return null;
}

1 usage
private List<String> getInfoFromRawData(StringBuilder content) {
    List<String> results = new ArrayList<>();
    String regex = "<REUTERS.*?>(.*)</REUTERS>";
    Pattern pattern = Pattern.compile(regex);
    Matcher matcher = pattern.matcher(content);
    while (matcher.find()) {
        results.add(matcher.group(1));
    }
    return results;
}
```

Figure 43: Fetch data from “reut2-014.sgm” and read data inbetween “REUTERS” tag.

```

1 usage
private List<String[]> filterData(List<String> info){
    List<String[]> response = new ArrayList<>();

    for (String currInfo : info) {
        String title = getTitleAndBody(currInfo, tagName: "TITLE");
        String filteredData = cleanTitleAndBody(title);
        if (filteredData != null){
            String[] words = filteredData.split( regex: " " );
            response.add(words);
        }
    }

    return response;
}

1 usage
private String getTitleAndBody(String content, String tagName) {
    String regex = "<" + tagName + ".*?>(.*)</" + tagName + ">";
    Pattern pattern = Pattern.compile(regex);
    Matcher matcher = pattern.matcher(content);
    if (matcher.find()) {
        return matcher.group(1);
    }
    return null;
}

```

Figure 44: Extract title from news.

```

1 usage
private String cleanTitleAndBody(String content) {
    if(content == null)
        return content;
    return content.replaceAll( regex: "[^a-zA-Z0-9\\s]", replacement: "" )
        .replace( target: "&lt;", replacement: "" ).replace( target: ">", replacement: "" )
        .replaceAll( regex: "\\s+", replacement: " " ).trim();
}

```

Figure 45: Remove special characters and unnecessary spaces.

- Read Words:
 - Retrieve lists of negative and positive words from negativeWords.txt and positiveWords.txt.

```

2 usages
public List<String> getWords(String filepath){
    List<String> words = new ArrayList<>();
    try (BufferedReader br = new BufferedReader(new FileReader(filepath))) {
        String line;
        while ((line = br.readLine()) != null) {
            words.add(line);
        }
        return words;
    } catch (IOException e) {
        System.out.println(e.getMessage());
    }
    return null;
}

```

Figure 46: Fetch positiveWords and negativeWords from files.

Calculate Sentiments

- For each data entry:
 - Count occurrences of positive and negative words.
 - Determine sentiment polarity (positive, negative, or neutral).
 - Create a TitleModel object with the sentiment details.

```

1 usage
public List<TitleModel> calculateSentiments(List<String[]> data, List<String> negativeWords, List<String> positiveWords) {
    List<TitleModel> response = new ArrayList<>();
    int index = 0;

    for (String[] curr : data) {
        int positive = 0;
        int negative = 0;
        List<String> curPositiveWords = new ArrayList<>();
        List<String> curNegativeWords = new ArrayList<>();

        index++;

        for (String str : curr) {
            if (positiveWords.contains(str.trim().toLowerCase())) {
                positive++;
                curPositiveWords.add(str);
            } else if (negativeWords.contains(str.trim().toLowerCase())) {
                negative++;
                curNegativeWords.add(str);
            }
        }

        TitleModel titleModel = createModel(index, curr, positive, negative, curPositiveWords, curNegativeWords);

        response.add(titleModel);
    }

    return response;
}

```

Figure 47: Iterate through each new and calculate its sentiment.

```

1 usage
private static String convertArrayToString(String[] titleArray) {
    StringBuilder stringBuilder = new StringBuilder();
    for (String str : titleArray) {
        stringBuilder.append(str);
    }
    return stringBuilder.toString();
}

1 usage
private static Type findPolarity(int positive, int negative) {
    if (positive > negative) {
        return Type.POSITIVE;
    } else if (negative > positive) {
        return Type.NEGATIVE;
    } else {
        return Type.NEUTRAL;
    }
}

```

Figure 48: Convert title array to string and find its polarity.

```

1 usage
private TitleModel createModel(int index, String[] title, int positive, int negative, List<String> curPositiveWords, List<String> curNegativeWords){
    TitleModel titleModel = new TitleModel();
    titleModel.setId(index);
    titleModel.setTitle(convertArrayToString(title));
    titleModel.setNegativeMatch(curNegativeWords);
    titleModel.setPositiveMatch(curPositiveWords);

    Type polarity = findPolarity(positive, negative);

    if (polarity.equals(Type.POSITIVE)) {
        titleModel.setPolarity(Type.POSITIVE);
        titleModel.setMatch(curPositiveWords);
        titleModel.setScore("+ " + positive);
    } else if (polarity.equals(Type.NEGATIVE)) {
        titleModel.setPolarity(Type.NEGATIVE);
        titleModel.setMatch(curNegativeWords);
        titleModel.setScore("- " + negative);
    } else {
        titleModel.setPolarity(Type.NEUTRAL);
        List<String> neutralList = new ArrayList<>();
        neutralList.addAll(curNegativeWords);
        neutralList.addAll(curPositiveWords);
        titleModel.setMatch(neutralList);
        titleModel.setScore("0");
    }

    return titleModel;
}

```

Figure 49: Create news model after calculating polarity.

Generate Excel Report

- Create an Excel workbook with a sheet labeled "Sentiments".
- Populate the sheet with the sentiment data from TitleModel objects.
- Format the workbook with appropriate headers and data rows.

```

1 usage
private static Workbook createExcel(List<TitleModel> sentiments) {
    Workbook workbook = new XSSFWorkbook();
    Sheet sheet = workbook.createSheet( s: "Sentiments");

    String[] headers = {"ID", "Title", "Match", "Positive Match", "Negative Match", "Score", "Polarity"};
    Row headerRow = sheet.createRow( i: 0);
    for (int i = 0; i < headers.length; i++) {
        Cell cell = headerRow.createCell(i);
        cell.setCellValue(headers[i]);
    }

    for (TitleModel sentiment : sentiments) {
        Row row = sheet.createRow(sentiment.getId());
        row.createCell( i: 0).setCellValue(sentiment.getId());
        row.createCell( i: 1).setCellValue(sentiment.getTitle());
        row.createCell( i: 2).setCellValue(String.join( delimiter: " ", " ", sentiment.getMatch()));
        row.createCell( i: 3).setCellValue(String.join( delimiter: " ", " ", sentiment.getPositiveMatch()));
        row.createCell( i: 4).setCellValue(String.join( delimiter: " ", " ", sentiment.getNegativeMatch()));
        row.createCell( i: 5).setCellValue(sentiment.getScore());
        row.createCell( i: 6).setCellValue(sentiment.getPolarity().toString());
    }
}

return workbook;
}

```

Figure 50: Convert list of news to one workbook.

Upload Data to MongoDB

- Connect to the SentimentAnalysis database and access the Sentiments collection.
- Convert the Excel data into MongoDB documents.
- Insert the documents into the MongoDB collection.

```

1 usage
private static void uploadExcelSheet(Workbook workbook) {
    List<Document> documents = new ArrayList<>();
    Sheet sheet = workbook.getSheetAt( i: 0);
    for (int i = 1; i <= sheet.getLastRowNum(); i++) {
        Row row = sheet.getRow(i);
        if (row != null) {
            Document doc = new Document()
                .append("ID", (int) row.getCell( i: 0).getNumericCellValue())
                .append("Title", row.getCell( i: 1).getStringCellValue())
                .append("Match", row.getCell( i: 2).getStringCellValue())
                .append("Positive Match", row.getCell( i: 3).getStringCellValue())
                .append("Negative Match", row.getCell( i: 4).getStringCellValue())
                .append("Score", row.getCell( i: 5).getStringCellValue())
                .append("Polarity", row.getCell( i: 6).getStringCellValue());

            documents.add(doc);
        }
    }
    collection.insertMany(documents);
    System.out.println("Data Inserted Successfully To MongoDB Database!!!");
}

```

Figure 51: Upload data to MongoDB.

MongoDB Connection

```
public class BOW {  
  
    2 usages  
    private static MongoClient mongoClient;  
    2 usages  
    private static MongoCollection<Document> collection;  
    2 usages  
    private static Calculations calculations;  
  
    1 usage  
    private static void connectMongoDB() {  
        MongoClientURI uri = new MongoClientURI("mongodb+srv://root:root@default-cluster.xn53b0c.mongodb.net/?retryWrites=true&w=majority&appName=default-cluster");  
        mongoClient = new MongoClient(uri);  
    }  
  
    public static void main(String[] args) {  
        FileOperations fileOperation = FileOperations.getInstance();  
        List<String[]> data = fileOperation.getBOW();  
        List<String> negativeWords = fileOperation.getWords( filepath: "negativeWords.txt" );  
        List<String> positiveWords = fileOperation.getWords( filepath: "positiveWords.txt" );  
        calculations = new Calculations();  
  
        List<TitleModel> response = calculations.calculateSentiments(data, negativeWords, positiveWords);  
        Workbook workbook = createExcel(response);  
  
        connectMongoDB();  
        MongoDatabase database = mongoClient.getDatabase( databaseName: "SentimentAnalysis" );  
        collection = database.getCollection( &quot;Sentiments" );  
        uploadExcelSheet(workbook);  
    }  
}
```

Figure 52: Connect to MongoDB database.

Output

- Here I upload data to MongoDB database instead of create excel file and store in local files.
- I create database name “SentimentAnalysis” and collection name “Sentiments”.

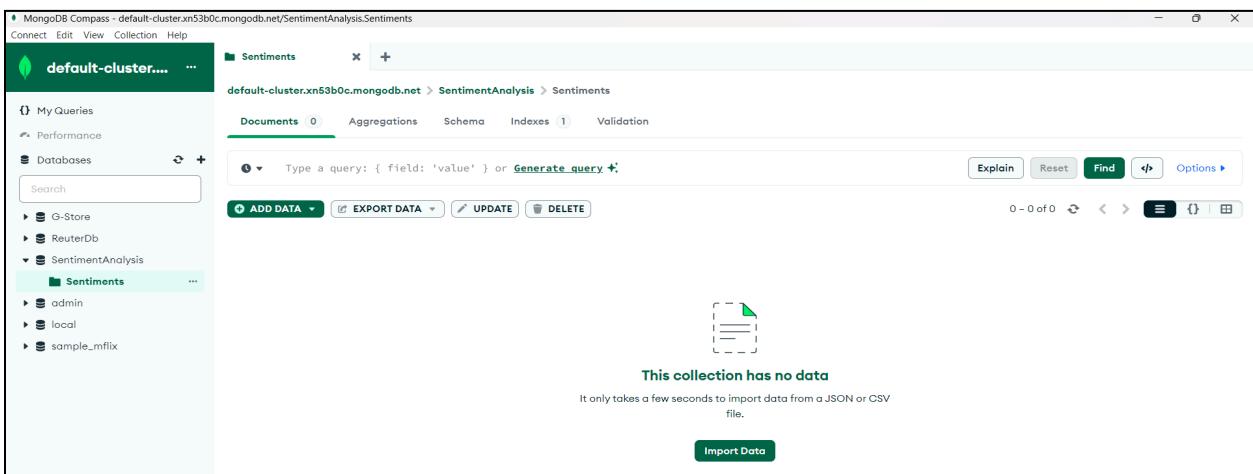


Figure 53: MongoDB database before I run Java program.

```

Aug 01, 2024 5:32:40 PM com.mongodb.diagnostics.logging.JULLogger log
INFO: Setting max set version to 130 from replica set primary ac-oeilnyp-shard-00-02.xn53b0c.mongodb.net:27017
Aug 01, 2024 5:32:40 PM com.mongodb.diagnostics.logging.JULLogger log
INFO: Discovered replica set primary ac-oeilnyp-shard-00-02.xn53b0c.mongodb.net:27017
Aug 01, 2024 5:32:41 PM com.mongodb.diagnostics.logging.JULLogger log
INFO: Opened connection [connectionId{localValue:4, serverValue:708428}] to ac-oeilnyp-shard-00-02.xn53b0c.mongodb.net:27017
Data Inserted Successfully To MongoDB Database!!

Process finished with exit code 0

```

Figure 54: Run Java Program.

ID	Title String	Match String	Positive Match String	Negative Match String	Score String
1	"JOHANNESBURGGOLDSHARESCL...	"GOLD, FIRMER"	"GOLD, FIRMER"	"	"+2"
2	"FEEDERCATTLEFUTURESETTNE...	"	"	"	"0"
3	"CBTCORNSPREADS"	"	"	"	"0"
4	"TEXASPANHANDLEWOKLAFEEDL...	"	"	"	"0"
5	"MIDWESTGRAINFUTURES1108E...	"	"	"	"0"
6	"VOLCKERPUSSHESPENDINGCUT...	"	"	"	"0"
7	"MIDWESTGRAINFUTURES1101E...	"	"	"	"0"
8	"CBTSOYBEANSPREADS"	"	"	"	"0"
9	"SIOUXFALLSCATTELEUP050DLR...	"FALLS"	"	"FALLS"	"-1"
10	"CBTSOYBEANOLSPREADS"	"	"	"	"0"
11	"FLASHSIOUTXFALLSHOGSUP175...	"FALLS, HOGS"	"	"FALLS, HOGS"	"-2"
12	"CBTSOYBEANMEALSPREADS"	"	"	"	"0"
13	"AMERICANSPORTSADVISORSIT...	"	"	"	"0"
14	"NYCOFFEEFUTURES1SLIGHTLYL...	"	"	"	"0"
15	"THATCHERFIRMASSHARESPREMO...	"	"	"	"0"

Figure 55: Data On MondoDB (Show ID, Title, Match, Positive Match, Negative Match, Score).

Title String	Match String	Positive Match String	Negative Match String	Score String	Polarity String
"JOHANNESBURGGOLDSHARESCL...	"GOLD, FIRMER"	"GOLD, FIRMER"	"	"+2"	"POSITIVE"
"FEEDERCATTLEFUTURESETTNE...	"	"	"	"0"	"NEUTRAL"
"CBTCORNSPREADS"	"	"	"	"0"	"NEUTRAL"
"TEXASPANHANDLEWOKLAFEEDL...	"	"	"	"0"	"NEUTRAL"
"MIDWESTGRAINFUTURES1108E...	"	"	"	"0"	"NEUTRAL"
"VOLCKERPUSSHESPENDINGCUT...	"	"	"	"0"	"NEUTRAL"
"MIDWESTGRAINFUTURES1101E...	"	"	"	"0"	"NEUTRAL"
"CBTSOYBEANSPREADS"	"	"	"	"0"	"NEUTRAL"
"SIOUXFALLSCATTELEUP050DLR...	"FALLS"	"	"FALLS"	"-1"	"NEGATIVE"
"CBTSOYBEANOLSPREADS"	"	"	"	"0"	"NEUTRAL"
"FLASHSIOUTXFALLSHOGSUP175...	"FALLS, HOGS"	"	"FALLS, HOGS"	"-2"	"NEGATIVE"
"CBTSOYBEANMEALSPREADS"	"	"	"	"0"	"NEUTRAL"
"AMERICANSPORTSADVISORSIT...	"	"	"	"0"	"NEUTRAL"
"NYCOFFEEFUTURES1SLIGHTLYL...	"	"	"	"0"	"NEUTRAL"
"THATCHERFIRMASSHARESPREMO...	"	"	"	"0"	"NEUTRAL"

Figure 56: Data On MondoDB (Polarity).

- Here I create the excel sheet and than upload data to mongoDB, I understand that if we have to upload data than no need to cerate data in form of excel sheet we can directly upload data to mongoDB in form documents.
- I create workbook because if I didn't want to upload data on mongoDB than I just have to just write few lines of code that is FileOutputStream to store excel sheet in local files.

MongoDB Compass - default-cluster.xn53b0c.mongodb.net/SentimentAnalysis.Sentiments

default-cluster.xn53b0c.mongodb.net > SentimentAnalysis > Sentiments

Documents 0 Aggregations Schema Indexes 1 Validation

Generate query Explain Reset Find Options

ADD DATA EXPORT DATA UPDATE DELETE

1 – 20 of 100

ID	Title String	Match String	Positive Match String	Negative Match String	Score String
1	"JOHANNESBURGOLDSHARESCL...	"GOLD, FIRMER"	"GOLD, FIRMER"	"	"+2"
2	"GOLDAANDSILVERCLOSEOEFHIG...	"GOLD"	"GOLD"	"	"+1"
3	"SWEDENERICSSONWINSUSORD...	"WINS"	"WINS"	"	"+1"
4	"CBTWHEATFUTURESTOPFIRME...	"FIRMER"	"FIRMER"	"	"+1"
5	"PORKBELLYFUTURESSARTHING...	"EASE"	"EASE"	"	"+1"
6	"CBTSOYBEANFUTURESTOPENSIL...	"FIRMER"	"FIRMER"	"	"+1"
7	"DOLLARREMAINSINNARROWRAL...	"QUIET"	"QUIET"	"	"+1"
8	"NEWBEDFORDINSTITUTIONFOR...	"SAVINGS"	"SAVINGS"	"	"+1"
9	"SPOTTINSLIGHTLYLOWERONEU...	"FREE"	"FREE"	"	"+1"
10	"USENERGYFUTURESSSTEADYBUT...	"STEADY, QUIET"	"STEADY, QUIET"	"	"+2"
11	"NYCOTTONFUTURESSSTEADYLINE...	"STEADY"	"STEADY"	"	"+1"
12	"UKGRAINFUTURESCLOSESTEADL...	"STEADY, FIRMER"	"STEADY, FIRMER"	"	"+2"
13	"SPOTLARDALONGGREASESTEAL...	"STEADY"	"STEADY"	"	"+1"
14	"CALTONLTCHOFFERINGDECLAR...	"EFFECTIVE"	"EFFECTIVE"	"	"+1"
15	"SUPERRITEFOODSINCINGSRFIS...	"SUPER"	"SUPER"	"	"+1"

Figure 57: Show only positive sentiments news.

MongoDB Compass - default-cluster.xn53b0c.mongodb.net/SentimentAnalysis.Sentiments

default-cluster.xn53b0c.mongodb.net > SentimentAnalysis > Sentiments

Documents 0 Aggregations Schema Indexes 1 Validation

Generate query Explain Reset Find Options

ADD DATA EXPORT DATA UPDATE DELETE

1 – 20 of 106

ID	Title String	Match String	Positive Match String	Negative Match String	Score String
1	"SIOUXFALLSCATTLEUP050DLR...	"FALLS"	"	"FALLS"	"-1"
2	"FLASHSIOUXFALLSHOGSUP175...	"FALLS, HOGS"	"	"FALLS, HOGS"	"-2"
3	"FLASHSPAUHLHOGSUP1092080...	"HOGS"	"	"HOGS"	"-1"
4	"VOLCKERSAYSRESTRICTIVEMO...	"RESTRICTIVE, HURT"	"	"RESTRICTIVE, HURT"	"-2"
5	"WORLD SUGAR FUTURESTUMBLED...	"TUMBLE"	"	"TUMBLE"	"-1"
6	"FLASHMOAHAHOGSUP1800DLRUS...	"HOGS"	"	"HOGS"	"-1"
7	"MISSOURIDIRECTHOGSUP1801...	"HOGS"	"	"HOGS"	"-1"
8	"VOLCKERSEESTIGHTPOLICYHUL...	"HURTING"	"	"HURTING"	"-1"
9	"MIDWAYAIRLINESIMOWYMARC...	"FALLS"	"	"FALLS"	"-1"
10	"AIDC ISSUES AUSTRALIA DOLL...	"ISSUES"	"	"ISSUES"	"-1"
11	"SIOUXCITYHOGSSTEADYUP125...	"HOGS"	"	"HOGS"	"-1"
12	"LONESTARLCLCESSESFLATUSC...	"NONE"	"	"NONE"	"-1"
13	"ST JOSEPHHOGSUP180150DLRU...	"HOGS"	"	"HOGS"	"-1"
14	"IaSoMinndirecthogsetima...	"hogs"	"	"hogs"	"-1"
15	"OHIODIRECTHOGSUP1800DLRUS...	"HOGS"	"	"HOGS"	"-1"

Figure 58: Show only negative sentiments news.

MongoDB Compass - default-cluster.xn53b0c.mongodb.net/SentimentAnalysis.Sentiments

default-cluster.xn53b0c.mongodb.net > SentimentAnalysis > Sentiments

Documents 0 Aggregations Schema Indexes 1 Validation

Generate query Explain Reset Find Options

ADD DATA EXPORT DATA UPDATE DELETE 1 – 20 of 492

ID	Title	Match	Positive Match	Negative Match	Score
1	"FEEDERCATTLEFUTURESSETTNE...	"	"	"	"0"
2	"CBTCORNSPREADS"	"	"	"	"0"
3	"TEXASPANHANDLEOKLAFEEDL...	"	"	"	"0"
4	"MIDWESTGRAINFUTURES1108E...	"	"	"	"0"
5	"VOLCKERPUSHESSPENDINGCUT...	"	"	"	"0"
6	"MIDWESTGRAINFUTURES1101E..."	"	"	"	"0"
7	"CBTSOYBEANSPREADS"	"	"	"	"0"
8	"CBTSOYBEANOILSPREADS"	"	"	"	"0"
9	"CBTSOYBEANMEALSPREADS"	"	"	"	"0"
10	"AMERICANSPORTSADVISORS1t..."	"	"	"	"0"
11	"NYCOFFEEFUTURESLIGHTLYH..."	"	"	"	"0"
12	"THATCHERFIRNASPRESSUREMO..."	"	"	"	"0"
13	"PORTUGUESEAIRLINECONFIRM..."	"	"	"	"0"
14	"DOLLARCLOSESLITTLECHANGE..."	"	"	"	"0"
15	"LIVEHOGFUTURESHIGHEREARL..."	"	"	"	"0"

Figure 59: Show only neutral sentiments news.

- To upload data to MongoDB, I chose the following fields:
 - ID
 - Title
 - Match
 - Positive Match
 - Negative Match
 - Score
 - Polarity
- After analyzing the sentiments for all news articles, I found that:
 - A total of 698 news articles are present.
 - 100 news articles have positive sentiments.
 - 106 news articles have negative sentiments.
 - 492 news articles have neutral sentiments.

References

- [1] diagrams.net, "Free Online Diagram Editor", <https://app.diagrams.net/>, [Accessed July 31, 2024].
- [2] S. Bleier, "Stop Words List", <https://gist.github.com/sebleier/554280>, [Accessed July 31, 2024].
- [3] Neo4j, "Cypher Query Language,
<https://neo4j.com/docs/cyphermanual/current/introduction/>, [Accessed July 31, 2024].
- [4] M. Kulakowski, "Negative Words List",
<https://gist.github.com/mkulakowski2/4289441>, [Accessed July 31, 2024].
- [5] M. Kulakowski, "Positive Words List",
<https://gist.github.com/mkulakowski2/4289437>, [Accessed July 31, 2024]