Assignment 3

# PROBLEM1-TASK2

Manraj Singh

B00877934

# Problem-1

## Task-2

- o Pseudocode for data extraction:
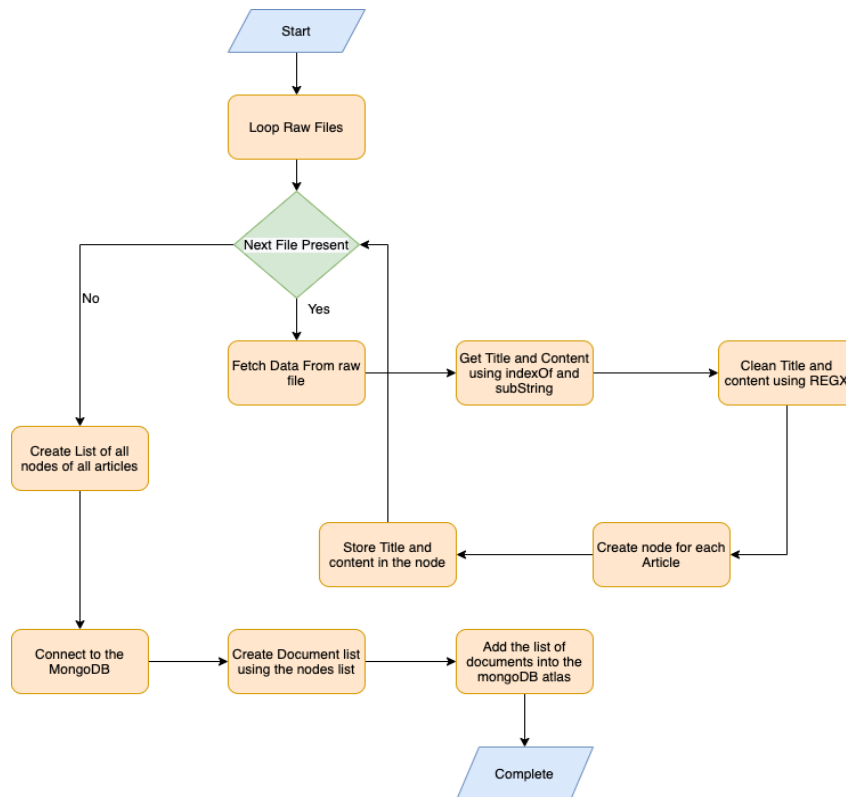    - Recursive delete all the pre-stored files in the Output/ directory.
    - Generate the API link for all the keywords using link [4] :
      *https://newsapi.org/v2/everything?q=<KeywordName>&sortBy=popularity&api Key=6e7ddac0bd2e44ec82aa90073f939e16&language=en&pageSize=90*
    - Replace the <keywordName> in ablove mentioned link with the keywords mentioned in the problem task 2.
    - Hit the new generated urls using HttpURLConnection and InputStreamReader and BufferedReader.
    - Fetch all the articles from each API hit in the loop.
    - To get the articles only from the raw data calculated startPoint and endpoint of the new SubString
        - StartPoint = rawData.indexOf("/"articles/":[");
        - EndPoint  = rawData.lastIndexOf("]");
        - Create substring from startPoint and endpoint
    - Create folder Structure to store raw files. Path : "Output/RawArticlesData/"
    - Break the data into group of 5 articles.
    - Write each group of 5 articles in a file inside the created folder structure.
    - Create new file with an incremented name *"RawArticlesData_<number++>_"* for further articles entry.



*Figure 1: RawArticlesData.json files [1][3]*

- As shown in *figure:1 [1][3]* 125 raw files generated with 5 articles each.
- To get the articles segregated from the raw data I used the index of "{" and "}" brackets.

- FlowChart for MongoDB connection and uploading cleaned data:[2]



Figure 2: Mongo Connection flowchart and data cleaning [5]

- Data is fetched from the raw files containing 5 articles each.
- Title and content fields are fetched for further processing using indexOf and substring.
- To clean the data REGX is used:
  - Regex variable = "[^0-9a-zA-Z:,\\s?!()\\/\\.]+" is used to remove all the special characters and emoticons.
  - Regex variable = "\\<.*?\\>"  is used to remove all the URLs and Http codes.
  - Regex variable = "(\r\n|\n)" is used to remove the special characters like "\r" or "\n"
- Nodes are created one per each article to contain the cleaned title and content data for each article.
- These nodes are stored in a list of nodes for further use.
- A json file is maintained containing all the the nodes data in json format (./Output/MongoArticlesProcessed/ MongoArticlesProcessed.json).
- MongoDB client connection is made.

- MongoDB Database named : "MyMongoNews" is made.
- MongoCollection named : "articles" is made
- Mongo DB connection is created using MongoDB atlas credentials.
- Documents list is created which contains Documents one per each node in the node list.
- Documents contain title and content in the json form.
- This Documents list is added to the MongoDB



*Figure 3: MongoDB Atlas View after adding all the articles [2]*



*Figure 4: MongoDB Data entry in local file (./Output/MongoArticlesProcessed/*
*MongoArticlesProcessed.json) [3]*

# References

[1] 35.226.207.33. 2021. *Spark Master at spark://assignment3-data-ms.us-central1-a.c.tribal-quasar-316422.internal:7077*. [online] Available at: <http://35.226.207.33:8080/> [Accessed 5 July 2021].

[2] Cloud.mongodb.com. 2021. *Cloud: MongoDB Cloud*. [online] Available at: <https://cloud.mongodb.com/v2/60da206dda3bb271226e1ead#metrics/replicaSet/60da21bb39e2ea2b93467554/explorer/myMongoNews/articles/find> [Accessed 5 July 2021].

[3] Console.cloud.google.com. 2021. *Google Cloud Platform*. [online] Available at: <https://console.cloud.google.com/compute/instancesDetail/zones/us-central1-a/instances/assignment3-data-ms?project=tribal-quasar-316422&rif_reserved> [Accessed 5 July 2021].

[4] Newsapi.org. 2021. *News API – Search News and Blog Articles on the Web*. [online] Available at: <https://newsapi.org/> [Accessed 5 July 2021].

[5] App.diagrams.net. 2021. *Flowchart Maker & Online Diagram Software*. [online] Available at: <https://app.diagrams.net/> [Accessed 5 July 2021].