## Lab 3_4

## CIS612 Big Data and Parallel data processing systems

**Semi-Structured Data Processing with MongoDB**

Create a database with two Collections from the given Yelp Business.json and Review.json data files in the semi-structured database server MongoDB
For this lab, you can use 5 JSON files in YelpDataSets.zip from the Yelp site (data set in 2017) in the CIS612 Lab section. Or you can directly download from the Yelp site at https://www.yelp.com/dataset

| | | | |
|---|---|---|---|
| yelp_academic_dataset_business | 19-01-2022 17:35 | JSON Source File | 1,16,078 KB |
| yelp_academic_dataset_review | 19-01-2022 17:51 | JSON Source File | 52,16,669 ... |

Here is yelpdata sets required

Creating Two Collections in MongoDB and Write a Data Pipelining to Retrieval Analytic Information Using Aggregate Pipelining and Join
1. Import Business.json and Review.json Data files from the Yelp site into MongoDB to create a Collection from each json file named business and review respectively.
2. Find the following information Using MongoDB Aggregation Pipelining:
Q1: For those business whose category is either "Fast Food" or "Restaurants", count the number of the business by each "city" and "stars".
"by each city and stars" means Group By city, stars in SQL.
Q2_1: For those business whose category is either "Fast Food" or "Restaurants" and review_count > 10 and stars >= 4, find all the reviews with stars >= 4 by performing Join with two Collections business and review.

MongoDB query to retrieve the data from JSON files

bash

mongoimport --db yelpDB --collection business --file D:\labb3\yelp_academic_dataset_business.json

mongoimport --db yelpDB --collection review --file D:\labb3\yelp_academic_dataset_review.json

```
Microsoft Windows [Version 10.0.22631.4317]
(c) Microsoft Corporation. All rights reserved.

D:\labb3>mongoimport --db yelpDB --collection business --file D:\labb3\yelp_academic_dataset_business.json
2024-10-27T11:33:43.224-0400    connected to: mongodb://localhost/
2024-10-27T11:33:46.225-0400    [#####...............] yelpDB.business    24.2MB/113MB (21.4%)
2024-10-27T11:33:49.224-0400    [##########..........] yelpDB.business    47.8MB/113MB (42.2%)
2024-10-27T11:33:52.224-0400    [###############.....] yelpDB.business    71.7MB/113MB (63.2%)
2024-10-27T11:33:55.229-0400    [##################..] yelpDB.business    96.2MB/113MB (84.8%)
2024-10-27T11:33:57.304-0400    [####################] yelpDB.business    113MB/113MB (100.0%)
2024-10-27T11:33:57.311-0400    150346 document(s) imported successfully. 0 document(s) failed to import.

D:\labb3>mongoimport --db yelpDB --collection review --file D:\labb3\yelp_academic_dataset_review.json
2024-10-27T11:34:09.464-0400    connected to: mongodb://localhost/
2024-10-27T11:34:12.465-0400    [....................] yelpDB.review    23.9MB/4.98GB (0.5%)
2024-10-27T11:34:15.466-0400    [....................] yelpDB.review    47.8MB/4.98GB (0.9%)
2024-10-27T11:34:18.465-0400    [....................] yelpDB.review    70.7MB/4.98GB (1.4%)
2024-10-27T11:34:21.466-0400    [....................] yelpDB.review    95.6MB/4.98GB (1.9%)
2024-10-27T11:34:24.464-0400    [....................] yelpDB.review    111MB/4.98GB (2.2%)
2024-10-27T11:34:27.465-0400    [....................] yelpDB.review    136MB/4.98GB (2.7%)
2024-10-27T11:34:30.465-0400    [....................] yelpDB.review    159MB/4.98GB (3.1%)
2024-10-27T11:34:33.464-0400    [....................] yelpDB.review    185MB/4.98GB (3.6%)
2024-10-27T11:34:36.465-0400    [....................] yelpDB.review    212MB/4.98GB (4.2%)
2024-10-27T11:34:39.467-0400    [#...................] yelpDB.review    238MB/4.98GB (4.7%)
2024-10-27T11:34:42.465-0400    [#...................] yelpDB.review    261MB/4.98GB (5.1%)
2024-10-27T11:34:45.465-0400    [#...................] yelpDB.review    287MB/4.98GB (5.6%)
2024-10-27T11:34:48.465-0400    [#...................] yelpDB.review    312MB/4.98GB (6.1%)
2024-10-27T11:34:51.465-0400    [#...................] yelpDB.review    338MB/4.98GB (6.6%)
2024-10-27T11:34:54.465-0400    [#...................] yelpDB.review    365MB/4.98GB (7.2%)
2024-10-27T11:34:57.465-0400    [#...................] yelpDB.review    389MB/4.98GB (7.6%)
2024-10-27T11:35:00.465-0400    [#...................] yelpDB.review    415MB/4.98GB (8.1%)
2024-10-27T11:35:03.464-0400    [##..................] yelpDB.review    438MB/4.98GB (8.6%)
2024-10-27T11:35:06.465-0400    [##..................] yelpDB.review    464MB/4.98GB (9.1%)
2024-10-27T11:35:09.465-0400    [##..................] yelpDB.review    490MB/4.98GB (9.6%)
2024-10-27T11:35:12.465-0400    [##..................] yelpDB.review    517MB/4.98GB (10.1%)
2024-10-27T11:35:15.464-0400    [##..................] yelpDB.review    542MB/4.98GB (10.6%)
2024-10-27T11:35:18.465-0400    [##..................] yelpDB.review    567MB/4.98GB (11.1%)
2024-10-27T11:35:21.464-0400    [##..................] yelpDB.review    591MB/4.98GB (11.6%)
```



MongoDB Compass - localhost:27017/yelpDB.business

Connections  Edit  View  Collection  Help

**Compass**

{} My Queries

CONNECTIONS (1)

Search connections

▼ 🖥 localhost:27017
  ▶ 🗄 admin
  ▶ 🗄 config
  ▶ 🗄 local
  ▼ 🗄 yelpDB
    📁 business
    📁 review

🌢 Welcome    📁 business    +

localhost:27017 > yelpDB > business

Documents 150.3K    Aggregations    Schema    Indexes 1    Validation

Type a query: { field: 'value' } or  Generate query +

[Explain]  [Reset]  [Find]

⊕ ADD DATA ▼    📝 EXPORT DATA ▼    ✏ UPDATE    🗑 DELETE    25 ▼  1 – 25 of 150346

```
_id: ObjectId('671e5d57d3d230eb4502af07')
business_id : "Pns2l4eNsfO8kk83dixA6A"
name : "Abby Rappoport, LAC, CMQ"
address : "1616 Chapala St, Ste 2"
city : "Santa Barbara"
state : "CA"
postal_code : "93101"
latitude : 34.4266787
longitude : -119.7111968
stars : 5
review_count : 7
is_open : 0
▸ attributes : Object
categories : "Doctors, Traditional Chinese Medicine, Naturopathic/Holistic, Acupunct…"
hours : null

_id: ObjectId('671e5d57d3d230eb4502af08')
business_id : "qkRM_2X51Yqxk3btlwAQIg"
name : "Temple Beth-El"
address : "400 Pasadena Ave S"
city : "St. Petersburg"
state : "FL"
postal_code : "33707"
```



Connections  Edit  View  Collection  Help

**Compass**

{} My Queries

CONNECTIONS (1)

Search connections

▼ 🖥 localhost:27017
  ▶ 🗄 admin
  ▶ 🗄 config
  ▶ 🗄 local
  ▼ 🗄 yelpDB
    📁 business
    📁 review

🌢 Welcome    📁 review    +

localhost:27017 > yelpDB > review

Documents 7.0M    Aggregations    Schema    Indexes 1    Validation

Type a query: { field: 'value' } or  Generate query +

[Explain]  [Reset]  [Find]  [</>]

⊕ ADD DATA ▼    📝 EXPORT DATA ▼    ✏ UPDATE    🗑 DELETE    25 ▼  1 – 25 of 6990280

```
_id: ObjectId('671e5d716b3706e1d4a7b1bd')
review_id : "KU_O5udG6zpxOg-VcAEodg"
user_id : "mh_-eMZ6K5RLWhZyISBhwA"
business_id : "XQfwVwDr-v0ZS3_CbbE5Xw"
stars : 3
useful : 0
funny : 0
cool : 0
text : "If you decide to eat here, just be aware it is going to take about 2 h…"
date : "2018-07-07 22:09:11"

_id: ObjectId('671e5d716b3706e1d4a7b1be')
review_id : "ZKvDG2sBvHVdF5oBNUOpAQ"
user_id : "wSTuiTk-sKNdcFyprzZAjg"
business_id : "B5X5oSG3SfvQGtKEGQ1tSQ"
stars : 3
useful : 1
funny : 1
cool : 0
text : "This easter instead of going to Lopez Lake we went to Los Padres Natio…"
date : "2016-03-30 22:46:33"

_id: ObjectId('671e5d716b3706e1d4a7b1bf')
review_id : "_ZeMknuYdlQcUqng_Im3yg"
```

Importing the Business.josn and Review.json into MongoDB server

**Python code:**

```python
from flask import Flask, jsonify, request, make_response
from pymongo import MongoClient
import pandas as pd

app = Flask(__name__)
client = MongoClient("mongodb://localhost:27017/")
db = client["yelpDB"]
business_collection = db["business"]
review_collection = db["review"]

@app.route('/count_businesses', methods=['GET'])
def count_businesses():
    pipeline = [
        {"$match": {"categories": {"$in": ["Fast Food", "Restaurants"]}}},
        {"$group": {"_id": {"city": "$city", "stars": "$stars"}, "count": {"$sum": 1}}},
        {"$sort": {"_id.city": 1, "_id.stars": 1}}
    ]
    result = list(business_collection.aggregate(pipeline))
    if request.args.get('format') == 'csv':
        return convert_to_csv(result, 'count_businesses')
    return jsonify(result)

@app.route('/high_rated_reviews', methods=['GET'])
def high_rated_reviews():
    pipeline = [
        {"$match": {"categories": {"$in": ["Fast Food", "Restaurants"]}, "review_count": {"$gt": 10}, "stars":
{"$gte": 4}}},
        {"$lookup": {
            "from": "review",
            "localField": "business_id",
            "foreignField": "business_id",
            "as": "high_reviews"
        }},
        {"$unwind": "$high_reviews"},
        {"$match": {"high_reviews.stars": {"$gte": 4}}},
            {"$project": {"review_id": "$high_reviews.review_id", "business_id": "$business_id", "stars":
"$high_reviews.stars", "review_text": "$high_reviews.text"}}
    ]
    result = list(business_collection.aggregate(pipeline))
    for document in result:
        if "_id" in document:
            document["_id"] = str(document["_id"])
        if "review_id" in document:
            document["review_id"] = str(document["review_id"])
```

```python
            if "business_id" in document:
                document["business_id"] = str(document["business_id"])
        if request.args.get('format') == 'csv':
            return convert_to_csv(result, 'high_rated_reviews')
        return jsonify(result)


@app.route('/low_rated_reviews', methods=['GET'])
def low_rated_reviews():
    pipeline = [
        {"$match": {"categories": {"$in": ["Fast Food", "Restaurants"]}, "review_count": {"$gt": 10}, "stars":
{"$lt": 2}}},
        {"$lookup": {
            "from": "review",
            "localField": "business_id",
            "foreignField": "business_id",
            "as": "low_reviews"
        }},
        {"$unwind": "$low_reviews"},
        {"$match": {"low_reviews.stars": {"$lt": 2}}},
            {"$project": {"review_id": "$low_reviews.review_id", "business_id": "$business_id", "stars":
"$low_reviews.stars", "review_text": "$low_reviews.text"}}
    ]
    result = list(business_collection.aggregate(pipeline))
    for document in result:
        if "_id" in document:
            document["_id"] = str(document["_id"])
        if "review_id" in document:
            document["review_id"] = str(document["review_id"])
        if "business_id" in document:
            document["business_id"] = str(document["business_id"])
        if request.args.get('format') == 'csv':
            return convert_to_csv(result, 'low_rated_reviews')
        return jsonify(result)


def convert_to_csv(data, filename):
    df = pd.DataFrame(data)
    print(df.head())  # Debug: print the first few rows of the DataFrame
    response = make_response(df.to_csv(index=False))
    response.headers["Content-Disposition"] = f"attachment; filename={filename}.csv"
    response.headers["Content-Type"] = "text/csv"
    return response


if __name__ == '__main__':
    app.run()
```

# Getting the files in CSV format:

## Count_business:



## High rated reviews:



## Low Rated reviews: