

Mini-project #3: MongoDB Aggregate Pipeline Queries- Nicholas Gincley & Lauren Kahrs

- Show the aggregation pipeline in screen captures. Show 2-3 stages in each screenshot, and include the first 1-2 sample results.
- Only show the requested fields in the final answer; there should not be any nested fields (arrays or objects).

Part 1: Querying Provided Collections

1. Use *restaurants* and *nyTheaters*. Display the restaurant street, borough, and zipcode along with the theater street where the restaurant zipcode matches the theater zipcode. Sort the results by borough and zipcode (both in ascending order). The solution has 6,811 documents.

Getting query1Restaurants view

25359 Documents in the Collection

Select an operator to construct expressions used in the aggregation pipeline stages. [Learn more](#)

Preview of Documents in the Collection

```
{ "_id": ObjectId("6066193ec11a1832dcc02f4d"),  
  "address": Object,  
  "borough": "Queens",  
  "cuisine": "American",  
  "grades": Array,  
    "name": "Brunos On The Boulevard",  
    "restaurant_id": "40356151"  
}
```

Output after \$project stage (Sample of 20 documents)

```
1 //**  
2 * specifications: The fields to  
3 * include or exclude.  
4 */  
5 {  
6   id: 0,  
7   borough: 1,  
8   restaurant_street: "$address.street",  
9   zipcode: "$address.zipcode"  
10 }
```

```
{ "borough": "Queens",  
  "restaurant_street": "Astoria Boulevard",  
  "zipcode": "11369"  
}  
  
{ "borough": "Brooklyn",  
  "restaurant_street": "Flatbush Avenue",  
  "zipcode": "11225"  
}
```

Getting query1Theaters view

81 Documents in the Collection

Select an operator to construct expressions used in the aggregation pipeline stages. [Learn more](#)

Preview of Documents in the Collection

```
{
  "_id": ObjectId("59a47286cf9a3a73e51e78d"),
  "theaterId": 1115,
  "location": Object
}
```

```
{
  "_id": ObjectId("59a47287cf9a3a73e51ea18"),
  "theaterId": 2851,
  "location": Object
}
```

\$project

Output after \$project stage (Sample of 20 documents)

```
1 // **
2 * specifications: The fields to
3 * include or exclude.
4 */
5 {
6   _id: 0,
7   zipcode: "$location.address.zipcode",
8   street: "$location.address.street1"
9 }
```

```
{
  "zipcode": "11354",
  "street": "13187 40th Road c300"
}
```

```
{
  "zipcode": "10941",
  "street": "1 N Galleria Dr"
}
```

Getting query

25359 Documents in the Collection

Select an operator to construct expressions used in the aggregation pipeline stages. [Learn more](#)

Preview of Documents in the Collection

```
{
  "borough": "Queens",
  "restaurant_street": "Astoria Boulevard",
  "zipcode": "11369"
}
```

```
{
  "borough": "Brooklyn",
  "restaurant_street": "Flatbush Avenue",
  "zipcode": "11225"
}
```

\$lookup

Output after \$lookup stage (Sample of 20 documents)

```
1 // **
2 * from: The target collection.
3 * localField: The local join field.
4 * foreignField: The target join field.
5 * as: The name for the results.
6 * pipeline: The pipeline to run on the joined collection.
7 * let: Optional variables to use in the pipeline function.
8 */
9 {
10   from: 'query1Theaters',
11   localField: 'zipcode',
12   foreignField: 'zipcode',
13   as: 'theaters'
14 }
```

```
{
  "borough": "Queens",
  "restaurant_street": "Astoria Boulevard",
  "zipcode": "11369",
  "theaters": Array
}
```

```
{
  "borough": "Brooklyn",
  "restaurant_street": "Flatbush Avenue",
  "zipcode": "11225",
  "theaters": Array
}
```

☰

Smatch

☒

+

```
1 ▾ /**  
2  * query: The query in MQL.  
3  */  
4 ▾ {  
5   theaters: {$ne: []}  
6 }  

```

Output after Smatch stage ① (Sample of 20 documents)

borough: "Queens"
restaurant_street: "Guardia Airport Parking"
zipcode: "11371"
▶ theaters: Array

borough: "Manhattan"
restaurant_street: "West 55 Street"
zipcode: "10019"
▶ theaters: Array

☰

Sunwind

☒

+

```
1 ▾ /**  
2  * path: Path to the array field.  
3  * includeArrayIndex: Optional name for index.  
4  * preserveNullAndEmptyArrays: Optional  
5  * toggle to unwind null and empty values.  
6  */  
7 ▾ {  
8   path: "$theaters",  
9   preserveNullAndEmptyArrays: true  
10 }  

```

Output after Sunwind stage ① (Sample of 20 documents)

borough: "Queens"
restaurant_street: "Guardia Airport Parking"
zipcode: "11371"
▶ theaters: Object

borough: "Queens"
restaurant_street: "Guardia Airport Parking"
zipcode: "11371"
▶ theaters: Object

☰

Sproject

☒

+

```
1 ▾ /**  
2  * specifications: The fields to  
3  * include or exclude.  
4  */  
5 ▾ {  
6   _id: 0,  
7   borough: 1,  
8   zipcode: 1,  
9   restaurant_street: 1,  
10  theater_street: "$theaters.street"  
11 }  

```

Output after Sproject stage ① (Sample of 20 documents)

borough: "Queens"
restaurant_street: "Guardia Airport Parking"
zipcode: "11371"
theater_street: "LaGuardia Airport"

borough: "Queens"
restaurant_street: "Guardia Airport Parking"
zipcode: "11371"
theater_street: "LaGuardia Airport"

☰

Scount

☒

+

```
1 ▾ /**  
2  * Provide the field name for the count.  
3  */  
4 'restaurant_street'
```

Output after Scount stage ① (Sample of 1 document)

restaurant_street: 6811

The screenshot displays a data transformation interface with two stages. The top stage is labeled '\$sort' and shows a sample of 14 documents. The bottom stage is labeled '\$count' and shows a sample of 1 document.

Stage 1: \$sort

Input (Sample of 14 documents):

```
1 /*  
2  * Provide any number of field/order pairs.  
3  */  
4 {  
5    property_type: -1  
6  }
```

Output (Sample of 14 documents):

```
count:1  
average_price:2450.00  
minimum_cleaning_fee:599.00  
maximum_cleaning_fee:599.00  
suburb:"Manhattan"  
property_type:"Townhouse"
```

```
count:3  
average_price:225.00  
minimum_cleaning_fee:75.00  
maximum_cleaning_fee:250.00  
suburb:"Manhattan"  
property_type:"Loft"
```

Stage 2: \$count

Input (Sample of 1 document):

```
1 /*  
2  * Provide the field name for the count.  
3  */  
4 {  
5    'property_type':  
6  }
```

Output (Sample of 1 document):

```
property_type:14
```

Part 2: Preparing Your Own Collection

1. Find a data set at <https://opendata.cityofnewyork.us/data/>. You should choose a data set that contains either a zipcode or borough. These fields may be named something different, but the data should match up to one of these fields from the *restaurants* collection. There are over 2,000 data sets available.

<https://data.cityofnewyork.us/Business/Consumer-Services-Mediated-Complaints/nre2-6m2s/data>

2. Give a high-level description of what is contained in the data set (one or two sentences.) This should not be a list of the fields; describe it meaningfully and concisely.

This dataset contains information complaints about businesses in the US (majority in the NYC area) that have been mediated by the DCA (Department of Consumer Affairs) and it's outcome, satisfaction of the complainer and any restitution.

3. List how many instances (e.g., rows) and how many fields (e.g., columns) are in the data set. Give the URL for the data set.

There are 2506 rows and 17 columns in this table.

<https://data.cityofnewyork.us/Business/Consumer-Services-Mediated-Complaints/nre2-6m2s/data>

4. Load the data into BigQuery and create a subset of the data with 1,000-2,000 instances.

Done (If you wanna do this too I downloaded the data from the link above as a csv then in BigQuery click create new table, then select import data from the dropdown list at the

top, select csv as the type on the right, select the file, tell it to auto create the schema then create.)

5. Describe in English in what ways you modified the data set (to reduce the number of instances and fields, if you choose to reduce the number of fields.) Give the BigQuery query or queries that you use to reduce the data set.

First I removed all rows not pertaining the New York.

```
delete from tactical-curve-258016.gincline_dw.Customer_Complaints  
  
where Business_State != 'NY';
```

This removed 460 rows leaving us with 2046 rows.

Then I removed all entries where the city is null (as this is what we will turn into Borough)

```
delete from tactical-curve-258016.gincline_dw.Customer_Complaints  
  
where Business_City is null;
```

This removed another 59 rows leaving us with 1987 rows

Then I removed many of the unnecessary columns in the dataset

```
ALTER TABLE tactical-curve-258016.gincline_dw.Customer_Complaints  
  
DROP COLUMN Longitude,  
  
DROP COLUMN Latitude,  
  
DROP COLUMN Building_Address_Unit,  
  
DROP COLUMN Business_Building,  
  
DROP COLUMN Mediation_Close_Date,  
  
DROP COLUMN Mediation_Start_Date,  
  
DROP COLUMN Business_State;
```

We dropped state because now we know all the businesses are in New York.

Next we changed all instances where the business city is New York to it's borough of Manhattan along with normalizing so other city names and reformatting them to later match the restaurant table.

```
update tactical-curve-258016.ginclene_dw.Customer_Complaints  
  
set Business_City = "Manhattan"  
  
where Business_City = "NEW YORK";
```

```
update tactical-curve-258016.ginclene_dw.Customer_Complaints  
  
set Business_City = "Queens"  
  
where Business_City = "QUEENS VLG"  
  
or Business_City = "QUEENS VILLAGE"  
  
or Business_City = "QUEENS";
```

```
update tactical-curve-258016.ginclene_dw.Customer_Complaints  
  
set Business_City = "Bronx"  
  
where Business_City = "BRONX";
```

```
update tactical-curve-258016.ginclene_dw.Customer_Complaints  
  
set Business_City = "Brooklyn"  
  
where Business_City = "BROOKLYN";
```

```
update tactical-curve-258016.ginclene_dw.Customer_Complaints  
  
set Business_City = "Staten Island"  
  
where Business_City = "STATEN ISLAND";
```

Finally we will now remove all entries that are not 1 of our 5 boroughs.

```
delete from tactical-curve-258016.ginclene_dw.Customer_Complaints
```

```
where Business_City != "Manhattan" AND  
  
Business_City != "Queens" AND  
  
Business_City != "Staten Island" AND  
  
Business_City != "Brooklyn" AND  
  
Business_City != "Bronx";
```

This leaves us with a total of 1269 rows in our dataset

6. Export the data to either JSON or CSV and give a Google Drive link to the final data set.

<https://drive.google.com/file/d/1nCNGvHpRPufAsR6n88Mvw0Ak7lK4bpFd/view?usp=sharing>

Part 3: Creating Data Demands and Solutions

Create two data demands and solutions that meet the given criteria. Your data demands and solutions should be presented in the same style as the data demands and solutions in Part 1.

- Using your collection prepared in Part 2, write a data demand and solution that requires using \$group, \$match, \$project, and \$sort stages.

Give each complaint type, how many of each complaint, average restitution for each complaint, and each industry associated with the complaint. Leave out any null values.

1269 Documents in the Collection

C

Select an operator to construct expressions used in the aggregation pipeline stages. [Learn more](#)

\$match

☒

+

```
1 // **  
2 * query: The query in MQL.  
3 */  
4 {  
5   "Satisfaction": "Yes",  
6   "Complaint_Type": {$ne: null}  
7 }
```

Preview of Documents in the Collection

```
_id: ObjectId("607340a93d4ab00908f20960")  
Business_Name: "CENTURY 21 DEPARTMENT STORE LLC"  
Industry: "Retail Store - 820"  
Complaint_Type: "Non-Delivery of Goods - N01"  
Complaint_Result: "Out of Business - OOB"  
Satisfaction: "NA"  
Restitution: 0  
Business_Street: "CORTLANDT ST"  
Business_City: "Manhattan"
```

```
_id: ObjectId("607340a93d4ab00908f20961")  
Business_Name: "BROOKLYN MOTOR CARS INC."  
Industry: "Secondhand Dealer Auto - 005"  
Complaint_Type: "Breach of Contract - B03"  
Complaint_Result: "Out of Business - OOB"  
Satisfaction: "NA"  
Restitution: 0  
Business_Street: "5TH AVE"  
Business_City: "Brooklyn"
```

Output after \$match stage ⓘ (Sample of 20 documents)

```
_id: ObjectId("607340a93d4ab00908f20967")  
Business_Name: "A-TWINAUTOMOTIVE EXPERTS CORP"  
Industry: "Tow Truck Company - 124"  
Complaint_Type: "Outstanding Judgment - J01"  
Complaint_Result: "Agency Collected Judgement - ACJ"  
Satisfaction: "Yes"  
Restitution: 1480  
Business_Street: "E 81ST ST"  
Business_City: "Brooklyn"
```

```
_id: ObjectId("607340a93d4ab00908f20969")  
Business_Name: "ROSSMANN REPAIR GROUP INC."  
Industry: "Electronic & Home Appliance Service Dealer -"  
Complaint_Type: "Non-Delivery of Service - N02"  
Complaint_Result: "Cash Amount - AMT"  
Satisfaction: "Yes"  
Restitution: 950  
Business_Street: "W 27TH ST"  
Business_City: "Manhattan"
```


\$match

1 /**
2 * query: The query in MQL.
3 */
4 {
5 "Industry": /^Restaurant/
6 }
7

Output after \$match stage (Sample of 20 documents)

1 /**
2 * specifications: The fields to
3 * include or exclude.
4 */
5 {
6 _id: 0,
7 Business_Name: 1,
8 Business_Zip: 1,
9 Business_Street: 1,
10 complaint tvoe: 1
11

Output after \$project stage (Sample of 20 documents)

Satisfaction: "Yes"
Restitution: 49
Business_Street: "ADAM CLAYTON POWELL JR BLVD"
Business_City: "Manhattan"
Business_Zip: "10039"
Complainant_Zip: "10451"
Business_Name: "KING BARKA JAMAICAN AND AMERICAN F

Complaint_Type: "Surcharge/Overcharge - S02"
Complaint_Result: "Cash Amount - AMT"
Satisfaction: "Yes"
Complainant_Zip: "10604"
_id: ObjectId("6078d0b3ffef133040909a10")
Industry: "Restaurant - 818"
Business_Street: "DYCKMAN ST"
Business_City: "Manhattan"

\$lookup

1 /**
2 * from: The target collection.
3 * localField: The local join field.
4 * foreignField: The target join field.
5 * as: The name for the results.
6 * pipeline: The pipeline to run on the joined collec
7 * let: Optional variables to use in the pipeline fie
8 */
9 {
10 from: 'restaurants',
11 localField: 'Business_Zip',
12 foreignField: 'address.zipcode',
13 as: 'zip'
14 }
15

Output after \$lookup stage (Sample of 1 document)

Business_Street: "ADAM CLAYTON POWELL JR BLVD"
Business_Zip: "10039"
zip: Array
Business_Name: "KING BARKA JAMAICAN AND AMERICAN RES
INC"
Complaint_Type: "Misrepresentation - M01"

\$match

1 /**
2 * query: The query in MQL.
3 */
4 {
5 zip: { \$ne: [] }
6 }
7

Output after \$match stage (Sample of 1 document)

Business_Zip: "10039"
zip: Array
Business_Name: "KING BARKA JAMAICAN AND AMERICAN RES
INC"
Complaint_Type: "Misrepresentation - M01"
Business_Street: "ADAM CLAYTON POWELL JR BLVD"

\$unwind

1 /**
2 * path: path to the array field.
3 * includeArrayIndex: Optional name for index.
4 * preserveNullAndEmptyArrays: Optional
5 * toggle to unwind null and empty values.
6 */
7 {
8 path: "\$zip",
9 preserveNullAndEmptyArrays: true
10 }
11

Output after \$unwind stage (Sample of 20 documents)

Business_Name: "KING BARKA JAMAICAN AND AMERICAN RES
INC"
Complaint_Type: "Misrepresentation - M01"
Business_Street: "ADAM CLAYTON POWELL JR BLVD"
Business_Zip: "10039"
zip: Object

Business_Zip: "10039"
zip: Object
Business_Name: "KING BARKA JAMAICAN AND AMERICAN RES
INC"
Complaint_Type: "Misrepresentation - M01"
Business_Street: "ADAM CLAYTON POWELL JR BLVD"

Business
Business
zip: Obje
Business
Complain

