

**DBMS Assignment 5**

**Sarah Groark**

**Nov 20, 2024**

Query 1 = Over how many years was the unemployment data collected?

```
> db.unemployment.distinct("Year")
< [
  1990, 1991, 1992, 1993, 1994,
  1995, 1996, 1997, 1998, 1999,
  2000, 2001, 2002, 2003, 2004,
  2005, 2006, 2007, 2008, 2009,
  2010, 2011, 2012, 2013, 2014,
  2015, 2016
]
> db.unemployment.distinct("Year").length
< 27
```

This query takes the distinct aggregate and processes the “Year” field in the schema to output all of the distinct years that are recorded in the data. As a result, the query output represents all of the years, in a list format, that data was recorded. Finally, the second query outputs the length of this list to return the range of years that unemployment data was collected. It should be noted that these two queries do not both have to be executed to get the final answer of 27, but was done so for visualization purposes.

Query 2 = How many states were reported on in this dataset?

```
> db.unemployment.distinct("State")
< [
  'Alabama', 'Arizona', 'Arkansas',
  'California', 'Colorado', 'Connecticut',
  'Delaware', 'Hawaii', 'Idaho',
  'Illinois', 'Indiana', 'Iowa',
  'Kansas', 'Kentucky', 'Louisiana',
  'Maine', 'Maryland', 'Massachusetts',
  'Michigan', 'Minnesota', 'Mississippi',
  'Missouri', 'Montana', 'Nebraska',
  'Nevada', 'New Hampshire', 'New Jersey',
  'New Mexico', 'New York', 'North Carolina',
  'North Dakota', 'Ohio', 'Oklahoma',
  'Oregon', 'Pennsylvania', 'Rhode Island',
  'South Carolina', 'South Dakota', 'Tennessee',
  'Texas', 'Utah', 'Vermont',
  'Virginia', 'Washington', 'West Virginia',
  'Wisconsin', 'Wyoming'
]
> db.unemployment.distinct("State").length
< 47
```

This query aggregates the distinct states from the schema and outputs the list of the unique states that data has been recorded under. The next query calculates the length of this list to output the amount of states that were reported on in the dataset. Similar to above, the second query can be calculated independently of the first to retrieve the answer.

Query 3 = What does this query compute?

```
db.unemployment.find({Rate : {$lt: 1.0}}).count()
```

```
> db.unemployment.find({Rate : {$lt: 1.0}}).count()
```

```
< 657
```

This query returns the number of instances that a state county had an unemployment rate less than 1% (during any of the recorded months and years).

Query 4 =

```
> db.unemployment.find({Rate : {$gt: 10.0}}, {County:1, State:1, _id:0})
```

```
< {
  State: 'Mississippi',
  County: 'Kemper County'
}
{
  State: 'Mississippi',
  County: 'Jefferson County'
}
{
  State: 'Mississippi',
  County: 'Sharkey County'
}
{
  State: 'Mississippi',
  County: 'Tunica County'
}
```

The query outputs the county and associated state where the rate is higher than 10%. It does so by filtering by the Rate field using \$gt: 10.0 to retrieve the documents that meet this condition. The second part of the query projects onto the State and County fields to output a neat result.

Query 5 = Calculate the average unemployment rate across all states

```
1 ▾ [
2 ▾ {
3   $group:
4   /**
5    * _id: null.
6    * AverageRate: computed average unemployment rate.
7    */
8   {
9     _id: null,
10    AverageRate: {
11      $avg: "$Rate"
12    }
13  },
14  {
15    $project:
16    /**
17     * specifications: The fields to
18     * include or exclude.
19     */
20    {
21      _id: 0,
22      AverageRate: 1
23    }
24  }
25 }
26 ]
```

OUTPUT:

```
{
  "AverageRate": 5.577907
}
```

This query occurs in two stages: a group and project stage. The group stage calculates the average unemployment rate across all of the documents (thus, all of the states) and stores it as AverageRate. The second project stage excludes the \_id field and solely outputs the AverageRate, representing the average unemployment rate across all of the states reported.

Query 6 = Find all **counties** with an unemployment rate between **5% and 8%**

```
1 ▾ [
2 ▾ {
3   $match:
4   /**
5    * filter Rates between 5 and 8%.
6    */
7   {
8     Rate: {
9       $gt: 5.0,
10      $lt: 8.0
11    }
12  },
13  {
14    $group:
15    /**
16     * _id: The id of the group - County field.
17     */
18    {
19      _id: "$County"
20    }
21  },
22  {
23    $project:
24    /**
25     * Fields to include: County.
26     */
27    {
28      County: "$_id",
29      _id: 0
30    }
31  }
32 }
33 ]
```

#### PIPELINE OUTPUT

Sample of 10 documents

OUTPUT OPTIONS ▾

County : "Vanderburgh County"

County : "Allegan County"

County : "Northumberland County"

County : "Washakie County"

County : "Belmont County"

County : "Miller County"

This query occurs in three different stages: \$match, \$group, and \$project. The match stage is responsible for filtering the documents to retrieve those that meet the criteria – an unemployment rate between 5 and 8%, this is done through the expression: Rate{ \$gt: 5.0, \$lt: 8.0}. Then, the group stage asserts the County field as the \_id, so that the counties are grouped together. The project stage then includes the County field in the output and excludes the \_id field to maintain neatness.

Query 7 = Find the state with the **highest unemployment rate**. Hint. Use { \$limit: 1 }

```
1  [
2    {
3      $project:
4        /**
5         * specifications: The fields to
6         * include or exclude.
7         */
8        {
9          State: 1,
10         Rate: 1,
11         _id: 0
12       },
13     },
14     {
15       $sort:
16         /**
17          * Provide any number of field/order pairs.
18          */
19         {
20           Rate: -1
21         }
22     },
23     {
24       $limit:
25         /**
26          * Provide the number of documents to limit.
27          */
28         1
29     }
30   ]
```

OUTPUT:

```
{
  "State": "Colorado",
  "Rate": 58.4
}
```

This query projects the State and Rate fields of the documents, then sorts them by the Rate in ascending order (highest to lowest). The final stage, \$limit, takes the top document in the pipeline output, which represents the highest unemployment rate.

Query 8 = Count how many **counties** have an unemployment rate above **5%**.

```
1  [
2    {
3      $match:
4        /**
5         * query: matches with instances of
6         * unemployment rates over 5.0%.
7         */
8      {
9        Rate: {
10          $gt: 5.0
11        }
12      },
13    },
14    {
15      $group:
16        /**
17         * _id: The id of the group (Count)
18         */
19      {
20        _id: "$County"
21      }
22    },
23    {
24      $count:
25        /**
26         * Count number of counties.
27         */
28        "count"
29    }
30  ]
```

OUTPUT:

```
{
  "Count": 1697
}
```

This query's pipeline first matches the documents that meet the criteria of having an unemployment rate above 5%. The group stage groups the current pipeline output by the County field to access these counties that are above 5%. Then, the \$count stage counts the amount of these instances.

9 = Calculate the average unemployment rate per state by year.

```
1  ▾ [
2  ▾ {
3    $group:
4    ▾ /**
5      * query: The query in MQL.
6      */
7    {
8      ▾ _id: {
9        state: "$State",
10       year: "$Year"
11      },
12      avg_unemployment_rate: {
13        $avg: "$Rate"
14      }
15    },
16  },
17  ▾ {
18    $project:
19    ▾ /**
20      * specifications: The fields to
21      * include or exclude.
22      */
23    {
24      ▾ _id: 0,
25        state: "$_id.state",
26        year: "$_id.year",
27        avg_unemployment_rate:
28          "$avg_unemployment_rate"
29      }
30    },
31  ]
```

```
state : "Alabama"
year : 2014
avg_unemployment_rate : 7.917164179104478
```

```
state : "Arizona"
year : 2014
avg_unemployment_rate : 9.470555555555556
```

```
state : "Arkansas"
year : 2014
avg_unemployment_rate : 6.968
```

```
state : "California"
year : 2014
avg_unemployment_rate : 8.820258620689655
```

```
state : "Colorado"
year : 2014
avg_unemployment_rate : 5.198958333333334
```



This query groups the `_id` field by state and year as well as calculates the average unemployment rate as specified by the `_id` field. Next, the project stage includes the state, year, and `avg_unemployment_rate` fields in the output.

10 = (Extra Credit) For each state, calculate the total unemployment rate across all counties (sum of all county rates).

```
1 ▾ [
2 ▾ {
3   $group:
4   /**
5    * _id: The id of the group (State field).
6    * fieldN: sum of unemployment rates.
7    */
8   {
9     _id: "$State",
10    total_unemployment_rate: {
11      $sum: "$Rate"
12    }
13  },
14 },
15 {
16   $project:
17   /**
18    * include total_unemployment_rate
19    * and state fields.
20    */
21   {
22     total_unemployment_rate: 1,
23     _id: 0,
24     state: "$_id"
25   }
26 }
27 ]
```

total\_unemployment\_rate : 553.4  
state : "Delaware"

total\_unemployment\_rate : 3533  
state : "Wyoming"

total\_unemployment\_rate : 8776.9  
state : "Montana"

total\_unemployment\_rate : 13668.3  
state : "Iowa"

This query takes the group stage to group by the State field and also calculates the sum of these rates by state. The project field includes the state and total\_unemployment\_rate fields in the output, but excludes the \_id field.

11 = (Extra Credit) The same as Query 10 but for states with data from 2015 onward

```
1 ▾ [
2 ▾ {
3   $match:
4   /**
5    * Criteria: years from 2015 and forward
6    */
7   {
8     Year: {
9       $gte: 2015
10    }
11  },
12 },
13 {
14   $group:
15   /**
16    * _id: The id of the group (State).
17    * total_unemployment_rate: summation of the rates by state.
18    */
19   {
20     _id: "$State",
21     total_unemployment_rate: {
22       $sum: "$Rate"
23     }
24  },
25 },
26 {
27   $project:
28   /**
29    * excludes _id, includes state and total_unemployment_rate.
30    */
31   {
32     _id: 0,
33     state: "$_id",
34     total_unemployment_rate: 1
35   }
36 }
37 ]
```

total\_unemployment\_rate : 3078.1  
state : "Maryland"

total\_unemployment\_rate : 3377.5  
state : "Utah"

total\_unemployment\_rate : 9605.4  
state : "Arkansas"

total\_unemployment\_rate : 10408.2  
state : "Indiana"

total\_unemployment\_rate : 4298.2  
state : "North Dakota"

This query does the same as #10 but includes a match stage at the beginning of the pipeline to filter documents with data from the year 2015 and onward.