

Aggregation in MongoDB: Additional Operations

Additional Aggregation Pipeline Operators

Faceted Filter

The `$facet` aggregation operation allows for a split of aggregation pipelines.

That is, `$facet` allows the user to provide a list of independent pipelines that can be run on a given document collection.

The `$facet` syntax is as follows:

```
{ $facet:
  {
    <outputField1>: [ <stage1>, <stage2>, ... ],
    <outputField2>: [ <stage1>, <stage2>, ... ],
    ...
    <outputFieldK>: [ <stage1>, <stage2>, ... ],
  }
}
```

The result of this operation is a **single** object that consists of fields: `<outputField1>`,...`<outputFieldK>` each with the result of the aggregation pipeline specified as its input.

Note. This is useful when you want to compute a variety of aggregation operations based on different groupings, or to see different slices of the dataset at the same time.

Example. This is a very simplified example based on a small collection of student grades in two courses:

```
> db.grades.find()
{ "_id" : 1, "name" : "Bob", "class" : 365, "grade" : 88 }
```

```

{ "_id" : 2, "name" : "Bob", "class" : 453, "grade" : 92 }
{ "_id" : 3, "name" : "May", "class" : 365, "grade" : 76 }
{ "_id" : 4, "name" : "May", "class" : 453, "grade" : 90 }
{ "_id" : 5, "name" : "Chris", "class" : 365, "grade" : 93 }
{ "_id" : 6, "name" : "Chris", "class" : 453, "grade" : 74 }

```

Consider a scenario where you want to return back an object that contains the following three things:

1. List of all students with a grade of 85 or above in CSC 365.
2. Average scores for students in each of the classes (as an array).
3. Highest score in CSC 453.

Each of the three requests can be executed in isolation as an aggregation pipeline:

1. To find CSC 365 students with a score of 85 or higher:

```
db.grades.aggregate({$match: {class:365, grade: {$gte: 85}}})
```

2. To find average scores in each class and format the output appropriately:

```
db.grades.aggregate({$group: {_id: "$class", avgScore:{$avg: "$grade"}}},
                    {$group: {_id:1, scores:{$push: {class:"$_id", avgScore:"$avgScore"}}}},
                    {$project: {_id:0}})
```

3. To find the highest score in CSC 453:

```
db.grades.aggregate({$group: {_id:"$class", maxGrade:{$max: "$grade"}}},
                    {$match: {_id: 453}},
                    {$project: {_id:0}} )
```

Using the `$facet` aggregation operation, we can now combine these three pipelines into as single step:

```

> db.grades.aggregate({$facet: {
  csc365_hi: [{$match: {class:365, grade: {$gte: 85}}}],
  averages: [{$group: {_id: "$class", avgScore:{$avg: "$grade"}},
              {$group: {_id:1, scores:{$push: {class:"$_id", avgScore:"$avgScore"}}}},
              {$project: {_id:0}}
            ],
  max453: [{$group: {_id:"$class", maxGrade:{$max: "$grade"}},
           {$match: {_id: 453}},
           {$project: {_id:0}}
          ]
        ]
      })
    ).pretty()

{
  "csc365_hi" : [
    {

```

```

        "_id" : 1,
        "name" : "Bob",
        "class" : 365,
        "grade" : 88
      },
      {
        "_id" : 5,
        "name" : "Chris",
        "class" : 365,
        "grade" : 93
      }
    ],
    "averages" : [
      {
        "scores" : [
          {
            "class" : 453,
            "avgScore" : 85.33333333333333
          },
          {
            "class" : 365,
            "avgScore" : 85.66666666666667
          }
        ]
      }
    ],
    "max453" : [
      {
        "maxGrade" : 92
      }
    ]
  ]
}

```

Left Join

The `$lookup` operation performs a left join of the current collection with the collection listed in the parameters of the `$lookup` on the conditions specified.

The syntax of the `$lookup` operation is

```

{
  $lookup:
  {
    from: <collection>,
    localField: <field1>,
    foreignField: <field2>,
    as: <arrayField>
  }
}

```

Here, `<collection>` specifies the collection with which to join. `<field1>` and `<field2>` complete the condition that will be checked. The condition is:

```
<coll>.<field1> = <collection>.<field2>
```

where `<coll>` is the collection on which the aggregation pipeline is run. `<arrayField>` stores the results.

This operation works as follows. For each object in the collection on which the aggregation pipeline is run, the collection `<collection>` is scanned, object-by-object. Any object from this collection that contains `<field2>` whose value is **exactly equal** to the value of the `<field1>` from the current object from the host collection is added into the `<arrayField>` as another array element.

The output object preserves all its original fields, and adds `<arrayField>` constructed as described above.

Example. In addition to the `grades` collection shown in the previous example, consider the following collection showing instructors of different classes:

```
> db.prof.find()
{ "_id" : 1, "course" : 365, "name" : "Alex" }
{ "_id" : 2, "course" : 369, "name" : "Alex" }
{ "_id" : 3, "course" : 453, "name" : "Phil" }
{ "_id" : 4, "course" : 307, "name" : "Davide" }
{ "_id" : 5, "course" : 466, "name" : "Foaad" }
{ "_id" : 6, "course" : 357, "name" : "Clint" }
```

The following `$lookup` operation adds the list of students taking each course to each object from the `profs` collection for which there are students in the `grades` collection.

```
> db.prof.aggregate({$lookup: {
...           from: "grades",
...           localField: "course",
...           foreignField: "class",
...           as: "roster"
...         }
...       })
{ "_id" : 1, "course" : 365, "name" : "Alex",
  "roster" : [ { "_id" : 1, "name" : "Bob", "class" : 365, "grade" : 88 },
               { "_id" : 3, "name" : "May", "class" : 365, "grade" : 76 },
               { "_id" : 5, "name" : "Chris", "class" : 365, "grade" : 93 } ] }
{ "_id" : 2, "course" : 369, "name" : "Alex", "roster" : [ ] }
{ "_id" : 3, "course" : 453, "name" : "Phil",
  "roster" : [ { "_id" : 2, "name" : "Bob", "class" : 453, "grade" : 92 },
               { "_id" : 4, "name" : "May", "class" : 453, "grade" : 90 },
               { "_id" : 6, "name" : "Chris", "class" : 453, "grade" : 74 } ] }
{ "_id" : 4, "course" : 307, "name" : "Davide", "roster" : [ ] }
{ "_id" : 5, "course" : 466, "name" : "Foaad", "roster" : [ ] }
{ "_id" : 6, "course" : 357, "name" : "Clint", "roster" : [ ] }
```

Note. `$lookup` performs, what is known as **left outer join**: the objects from the host collection that are not joined with any objects from the foreign collection are retained in the output with an empty array for the "join" key-value pair.

Bucketing

The **bucketing** operation is the version of **grouping** that uses a continuous variable to break the data into separate groups/buckets. Essentially, rather than combining the objects based on the same value of a specific key, **bucketing** combines them based on the value of a key falling in a specific range (bucket).

In MongoDB aggregation pipelines bucketing is performed using the **\$bucket** operator. The syntax of a **\$bucket** step is show below:

```
{
  $bucket: {
    groupBy: <expression>,
    boundaries: [ <lowerbound1>, <lowerbound2>, ... ],
    default: <literal>,
    output: {
      <output1>: { <$accumulator expression> },
      ...
      <outputN>: { <$accumulator expression> }
    }
  }
}
```

Here:

- **<expression>** value of the **groupBy** key is the value that will be bucketed. This is usually a referenece to a key storing the value, but it can also be any other numeric expression (see **\$project** operation documentation for numeric expressions).
- The array of numeric values supplied for the **boundaries** key stores the breakdown of the space into buckets. Given the K values

[n_1 , n_2 , ..., n_K]

there will be $K-1$ buckets of the form: $[n_1, n_2), [n_2, n_3), \dots, [n_{K-1}, n_K]$. Each bucket $[n_i, n_{i+1}]$ will be represented in the output by a single object with the key **_id**: n_i . Values that are smaller than n_1 or greater than n_K fall outside of the bucket ranges...

- ... and are handled by the **default** key. Its value is a literal that will be used as the unique identifier of the default bucket, which is reserved for accumulating information about objects that do not fall into any of the other buckets.
- Finally, the **output** key describes how the output objects look like. The key-value pairs inside this key are formed the same way as all key-value pairs except for **_id** are formed in **\$group** operation.

Example. Let's count how many grades in ranges 75-85, 85-95 and 95-100 are in the `grades` collection.

```
> db.grades.aggregate({$bucket: { groupBy:"$grade",
...                               boundaries: [75,85,95,100],
...                               default: "less than 75",
...                               output: {num: {$sum: 1}}
...                               }
...                               })
{ "_id" : 75, "num" : 1 }
{ "_id" : 85, "num" : 4 }
{ "_id" : "less than 75", "num" : 1 }
```

Special Version of Projection

\$addField. The `$addField` operation works like a projection operation which includes *all* attributes from the input document and adds additional attributes to the document. The syntax of the command is

```
{ $addField: { <newField>: <expression>, ...,
               <newField>: <expression> } }
```

Here `<newField>` represents the names of the new fields to be added to the output documents, while `<expression>` has the same syntax as the expressions used in the `$project` operation.

Example. Consider the following simple list of courses with the sizes of each section specified.

```
> db.classes.find()
{ "_id" : 1, "class" : "CSC 369", "roster" : 28 }
{ "_id" : 2, "class" : "CSC 445", "roster" : 34 }
{ "_id" : 3, "class" : "CSC 466", "roster" : 17 }
{ "_id" : 4, "class" : "CSC 357", "name" : "Systems Programming", "sections" : [ "01", "03", "05", "07" ], "roster" : 34 }
{ "_id" : 5, "class" : "CSC 101", "roster" : 100 }
{ "_id" : 6, "class" : "CSC 480", "name" : "AI", "roster" : 28 }
{ "_id" : 7, "roster" : 34, "class" : "CSC 202" }
{ "_id" : 8, "class" : "CSC 202", "sections" : [ "01", "02", "03", "04" ], "roster" : [ 20, 20, 30, 30 ] }
```

The following aggregation pipeline adds a new key to each document, specifying whether the course in question has large sections.

```
> db.classes.aggregate({$addField:
                        {sectionSize:
                          {$cond: [{$gte: ["$roster", 30]}, "large", "small"]}}})
{ "_id" : 1, "class" : "CSC 369", "roster" : 28, "sectionSize" : "small" }
{ "_id" : 2, "class" : "CSC 445", "roster" : 34, "sectionSize" : "large" }
{ "_id" : 3, "class" : "CSC 466", "roster" : 17, "sectionSize" : "small" }
{ "_id" : 4, "class" : "CSC 357", "name" : "Systems Programming", "sections" : [ "01", "03", "05", "07" ],
  "roster" : [ 34, 20, 32, 25 ], "sectionSize" : "large" }
{ "_id" : 5, "class" : "CSC 101", "roster" : 100, "sectionSize" : "large" }
```

```
{ "_id" : 6, "class" : "CSC 480", "name" : "AI", "roster" : 28, "sectionSize" : "small" }
{ "_id" : 7, "roster" : 34, "class" : "CSC 202", "sectionSize" : "large" }
{ "_id" : 8, "class" : "CSC 202", "sections" : [ "01", "02", "03", "04" ], "roster" : [ 20, 20, 30, 30 ],
  "sectionSize" : "large" }
```

Additional Operations

\$sample. The `$sample` operation returns a random sample of documents from the collection. The syntax of the aggregation pipeline command is:

```
{$sample: {size: <positive integer>}}
```

Here, `<positive integer>` represents the number of documents to put in the sample. Note, this operation samples with replacement, so a document from the original collection can be placed in the sample multiple times.

Example. In the example below, we ask for a sample of size three twice in a row and observe different documents returned.

```
> db.classes.aggregate({$sample: {size: 3}})
{ "_id" : 4, "class" : "CSC 357", "name" : "Systems Programming",
  "sections" : [ "01", "03", "05", "07" ], "roster" : [ 34, 20, 32, 25 ] }
{ "_id" : 5, "class" : "CSC 101", "roster" : 100 }
{ "_id" : 2, "class" : "CSC 445", "roster" : 34 }
> db.classes.aggregate({$sample: {size: 3}})
{ "_id" : 7, "roster" : 34, "class" : "CSC 202" }
{ "_id" : 1, "class" : "CSC 369", "roster" : 28 }
{ "_id" : 8, "class" : "CSC 202", "sections" : [ "01", "02", "03", "04" ],
  "roster" : [ 20, 20, 30, 30 ] }
```

\$count. This operation returns an object with a single key-value pair. The key name is provided as the input to the operation. The value is the number of documents produced on the previous stage of the aggregation pipeline. The format of the aggregation pipeline document for `$count` is as follows:

```
{$count: <string>}
```

Here, `<string>` is the name of the key in the output document.

example. The example below runs an aggregation pipeline that counts how many courses have sections that have less than 25 students in them.

```
> db.classes.aggregate({$match: {"roster": {$lt: 25}}})
{ "_id" : 3, "class" : "CSC 466", "roster" : 17 }
{ "_id" : 4, "class" : "CSC 357", "name" : "Systems Programming",
  "sections" : [ "01", "03", "05", "07" ], "roster" : [ 34, 20, 32, 25 ] }
{ "_id" : 8, "class" : "CSC 202", "sections" : [ "01", "02", "03", "04" ],
  "roster" : [ 20, 20, 30, 30 ] }
```

```
> db.classes.aggregate({$match: {"roster": {$lt: 25}},
                        {$count: "numberSmallClasses"})
```

```
{ "numberSmallClasses" : 3 }
```

\$out. The `$out` pipeline aggregation command can only show up as the last command of the pipeline. It directs the result of the previous stage of the pipeline to be inserted into a given new collection. The format of the command is simple:

```
{ $out: "<output-collection>" }
```

Here, "`<output-collection>`" is a string representing the name of the new collection into which the result of the aggregation pipeline will be inserted.

Example. Here is an example of creating a new collection `classNames` consisting of only course names.

```
> db.classes.aggregate({$project: {"_id":0, "class":1}},
                        {$out: "classNames"})
```

```
> db.classNames.find()
```

```
{ "_id" : ObjectId("5c4eb2502209cf8793d2fd04"), "class" : "CSC 369" }
{ "_id" : ObjectId("5c4eb2502209cf8793d2fd05"), "class" : "CSC 445" }
{ "_id" : ObjectId("5c4eb2502209cf8793d2fd06"), "class" : "CSC 466" }
{ "_id" : ObjectId("5c4eb2502209cf8793d2fd07"), "class" : "CSC 357" }
{ "_id" : ObjectId("5c4eb2502209cf8793d2fd08"), "class" : "CSC 101" }
{ "_id" : ObjectId("5c4eb2502209cf8793d2fd09"), "class" : "CSC 480" }
{ "_id" : ObjectId("5c4eb2502209cf8793d2fd0a"), "class" : "CSC 202" }
{ "_id" : ObjectId("5c4eb2502209cf8793d2fd0b"), "class" : "CSC 202" }
```