

COMP5338 – Advanced Data Models

Week 3: MongoDB – Aggregation Framework

Assignment Project Exam Help

Dr. Ying Zhou
School of Information Technologies

<https://powcoder.com>

Add WeChat powcoder



THE UNIVERSITY OF
SYDNEY

Outline

■ Review

■ Aggregation

- ▶ Pipeline stages
- ▶ Operators

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

COMMONWEALTH OF AUSTRALIA

Copyright Regulations 1969

WARNING

This material has been reproduced and communicated to you by or on behalf of the **University of Sydney** pursuant to Part VB of the Copyright Act 1968 (the Act).

The material in this communication may be subject to copyright under the Act. Any further reproduction or communication of this material by you may be the subject of copyright protection under the Act.

Do not remove this notice

Review

- Document Storage Systems store data as semi-structured document: XML or JSON as two dominant semi-structured formats
 - ▶ Semi-structured data is self-describing
- MongoDB is a popular document storage system that stores data as Binary representation of JSON document (BSON)
- Documents with similar structure that representing a particular type of entity are stored in the same collection
- A database is used to hold multiple collections representing related entities
- All CRUD operations (find, update, insert and delete) target single collection
- Query criteria, projection and modifier are expressed as JSON document

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Null, empty string and related operators

- Null (or null) is a special data type
 - ▶ Similar to **None**, **Null** or **Nil** in any programming language
 - ▶ It has a singleton value expressed as **null**
 - ▶ Indicating no value is given
- The interpretation of null is different depending on where it appears
- It might represents
 - ▶ The field exists, but has no value
 - ▶ The field does not exists
- This is different to given a field an empty string "" as value

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Null query example

■ Collection revisions document sample

```
{ "_id" : ObjectId("5799843ee2cbe65d76ed919b"),  
  "title" : "Hillary_Clinton",  
  "timestamp" : "2016-07-23T02:02:06Z",  
  "revid" : 731113635,  
  "user" : "BD2412",  
  "parentid" : 731113573,  
  "size" : 251742,  
  "minor" : ""}
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- We need a field to indicate if a revision is minor or not. The original schema uses a field with empty string value to indicate a minor revision; a document without this field would be a non-minor revision.

<https://docs.mongodb.com/manual/tutorial/query-for-null-fields/>

Querying for null or field existence

■ Queries

- ▶ `db.revisions.find({minor:{$exists:true}})`
 - Find all documents that a field called `minor` exists
- ▶ `db.revisions.find(where the {minor:""})`
 - Find all documents whose `minor` field has a value of "", empty string
- ▶ `db.revisions.find({minor:null})`
 - Find all documents that does not have a `minor` field or the value of `minor` field is null
- ▶ `db.revisions.find({minor:{$exists:false}})`
 - Find all documents that does not have a field called `minor`

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

It is possible to set the value to null

```
db.revisions.insertOne({title:"nulltest",  
  "timestamp" : "2018-08-14T02:02:06Z",  
  "revid" : NumberLong(7201808141159),  
  "user" : "BD2412",  
  "parentid" : 731113573,  
  "size" : NumberInt(251900),  
  "minor":null})
```

Assignment Project Exam Help

<https://powcoder.com>

```
db.revisions.insertOne({title:"nulltest",  
  "timestamp" : "2018-08-14T02:02:06Z",  
  "revid" : NumberLong(201808141157),  
  "user" : "BD2412",  
  "parentid" : NumberLong(731113573),  
  "size" : NumberInt(251800)})
```

Add WeChat powcoder

`db.revisions.find({minor:null})` would return both documents
`db.revisions.find({minor:{$exists:true}})` can differentiate the two



Aggregation

- Simple and relatively standard data analytics can be achieved through **aggregation**

- ▶ Grouping, summing up value, counting, sorting, etc
- ▶ Running on the DB engine instead of application layer

Assignment Project Exam Help

- Several options <https://powcoder.com>

- ▶ Aggregation Pipeline
- ▶ MapReduce **Add WeChat powcoder**
 - Through JavaScript Functions
 - Is able to do customized aggregations

Aggregation Pipeline

- Aggregation pipeline consists of multiple stages
 - ▶ Stages are specified using **pipeline operators** such as **\$match**, **\$group**, **\$project**, **\$sort** and so on
 - This is similar to SQL's WHERE, GROUP BY, SORT BY etc
 - Each stage is expressed as an object enclosed by curly bracket
 - ▶ Various **expressions** can be specified in each stage
 - To filter documents or to perform simple calculation on an document
 - \$substr, \$size, etc, ..
 - ▶ **\$group** stage can specify **accumulators** to perform calculation on documents with the same group key

```
db.collection.aggregate( [  
  { pipeline operator: {expression/accumulator,..., expression/accumulator} },  
  { pipeline operator: {expression/accumulator,..., expression/accumulator} },  
  ...  
] )
```

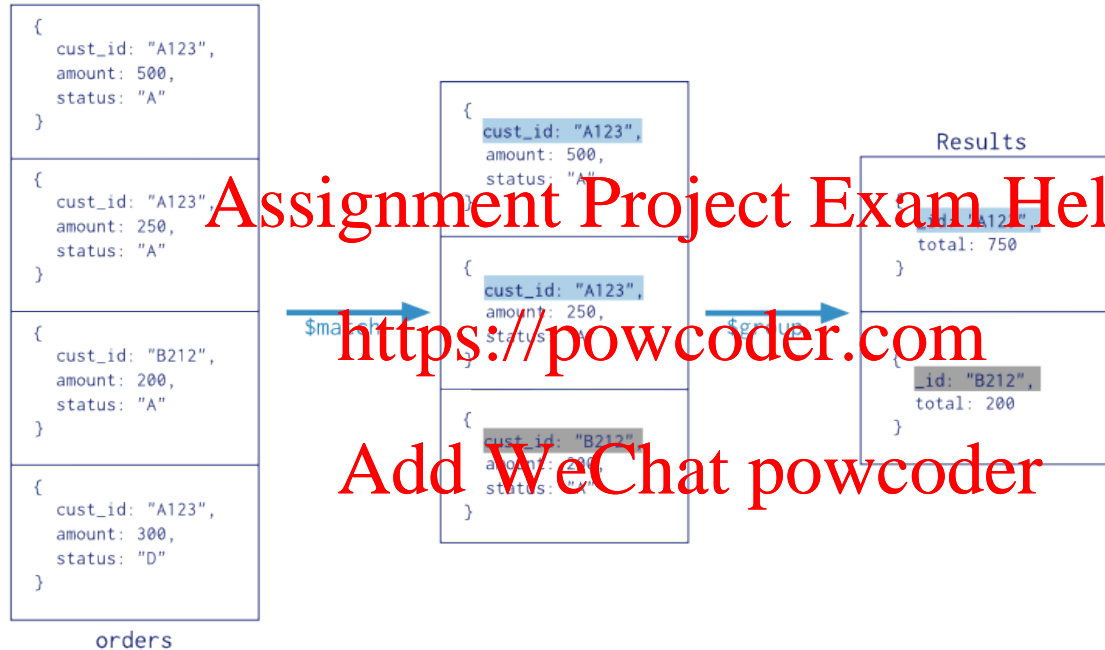
Aggregation Example

Collection
↓
db.orders.aggregate([

```

    $match stage → { $match: { status: "A" } },
    $group stage → { $group: { _id: "$cust_id", total: { $sum: "$amount" } } }
  ] )

```



```

select cust_id as _id, SUM(amount) as total
from orders
where status = "A"
group by cust_id

```

Typical aggregation stages

- \$match
- \$group
- \$project
- \$sort
- \$skip
- \$limit
- \$count
- \$sample
- \$out
- \$unwind
- \$lookup

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



\$match stage

■ \$match

- ▶ Filter the incoming documents based on given conditions

- ▶ Format:

```
{ $match: { <query> } }
```

- The query document is the same as those in the `find` query

- ▶ Example:

```
db.revisions.aggregate([ { $match: { size : { $lt: 250000 } } } ] )
```

<https://powcoder.com>
Add WeChat powcoder

Has the same effect as

```
db.revisions.find({ size : { $lt: 250000 } })
```

\$group stage

■ \$group

- ▶ Groups incoming documents by some specified expression and outputs to the next stage a document for each distinct group
 - The `_id` field of the output document has the value of the group key for each group
 - The stage can specify many other fields

```
{ $group: { _id: <expression>,
            <field1>: { accumulator: <expression> },
            ... },
```

- ▶ To specify the whole collection as a group, give `_id` field `null` value
- ▶ Use *field path* to access fields in the document and set one or many as the value of the `_id` field
 - `"$title"`, or `"$address.street"`
- ▶ There are predefined accumulators: `$sum`, `$avg`, `$first`, `$last`, etc

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

\$group stage example

- Find the earliest revision time in the whole collection

```
db.revisions.find({}, {timestamp:1, _id:0})
```

```
.sort({timestamp:1})
```

```
.limit(1)
```

```
db.revisions.aggregate([  
  {  
    $group: {  
      _id: null, earliest: {  
        $min: "$timestamp"  
      }  
    }  
  }  
])
```

- Find the earliest revision time of each page in the collection

```
db.revisions.aggregate([  
  {  
    $group: {  
      _id: "$title", earliest: {  
        $min: "$timestamp"  
      }  
    }  
  }  
])
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

field path

\$group stage example (cont'd)

- Find the number of revisions made on each page by each individual user
 - ▶ This would require grouping based on two fields: title and user
 - ▶ We need to specify these two as the `_id` field of the output document

Assignment Project Exam Help

```
db.revisions.aggregate([  
  {  
    $group: {  
      _id: {  
        title: "$title",  
        user: "$user"  
      },  
      rev_count: {  
        $sum: 1  
      }  
    }  
  }  
])
```

Composite type as _id

<https://powcoder.com>

Add WeChat powcoder

\$group by more than one field

{_id:ObjectId("..."), title: "DT", user:"A", size:123, timestamp:..., ... }
{_id:ObjectId("..."), title: "HC", user:"B", size:113, timestamp:..., ... }
{_id:ObjectId("..."), title: "DT", user:"B", size:125, timestamp:..., ... }
{_id:ObjectId("..."), title: "HC", user:"A", size:113, timestamp:..., ... }
{_id:ObjectId("..."), title: "DT", user:"A", size:125, timestamp:..., ... }

Assignment Project Exam Help

<https://powcoder.com>

```
{ $group: { _id: { title: "$title", user: "$user" },  
  rev_count: { $sum: 1 } } }
```

Add WeChat powcoder

```
{_id: {title: "DT", user:"A"}, rev_count: 2}  
{_id: {title: "HC", user:"B"}, rev_count: 1}  
{_id: {title: "DT", user:"B"}, rev_count: 1}  
{_id: {title: "HC", user:"A"}, rev_count: 1}
```


\$group examples (cont'd)

- Accumulators do not just return a single value, we can use accumulators to create an array to hold data from incoming documents

- What do the following two commands do:

```
db.revisions.aggregate([
  {$group: {_id:"$title",
    revs: {$push:{user:"$user",timestamp:"$timestamp"}}}
  ]])
```

Add WeChat powcoder

```
db.revisions.aggregate([
  {$group: {_id:"$title",rev_users: {$addToSet:"$user"}}}
  ])
```

Assignment Project Exam Help

<https://powcoder.com>

\$push accumulator

```
{_id:ObjectId("..."), title: "DT", user:"A", size:123, timestamp:..., ... }
```

```
{_id:ObjectId("..."), title: "HC", user:"B", size:113, timestamp:..., ... }
```

```
{_id:ObjectId("..."), title: "DT", user:"B", size:125, timestamp:..., ... }
```

```
{_id:ObjectId("..."), title: "HC", user:"A", size:113, timestamp:..., ... }
```

```
{_id:ObjectId("..."), title: "DT", user:"A", size:125, timestamp:..., ... }
```

Assignment Project Exam Help

<https://powcoder.com>

```
db.revisions.aggregate([
  {$group:
    {_id:"$title",
     revs:{$push:{user:"$user",timestamp:"$timestamp"}}}
  ]})
```

Add WeChat powcoder

```
{ _id: "DT",
  revs:[
    {user:"A",timestamp:...},
    {user:"B",timestamp:...},
    {user:"A",timestamp:...}
  ]
}
```

```
{ _id:"HC",
  revs:[
    {user:"A", timestamp:...},
    {user:"B", timestamp:...}
  ]
}
```

\$addToSet accumulator

```
{_id:ObjectId("..."), title: "DT", user:"A", size:123, timestamp:..., ... }  
{_id:ObjectId("..."), title: "HC", user:"B", size:113, timestamp:..., ... }  
{_id:ObjectId("..."), title: "DT", user:"B", size:125, timestamp:..., ... }  
{_id:ObjectId("..."), title: "HC", user:"A", size:113, timestamp:..., ... }  
{_id:ObjectId("..."), title: "DT", user:"A", size:125, timestamp:..., ... }
```

Assignment Project Exam Help

<https://powcoder.com>

```
db.revisions.aggregate([  
  {$group: {_id:"$title",  
    rev_users:{$addToSet:"$user"}}}  
])
```

Add WeChat powcoder

```
{_id:"DT",  
  rev_users:["A", "B"]  
}
```

```
{_id:"HC",  
  rev_users:["A", "B"]  
}
```

\$project stage

■ \$project

- ▶ Reshape the document by including/excluding field, adding new fields, resetting the value of existing field

- ▶ More powerful than the *project* argument in **find** query

- ▶ Format **Assignment Project Exam Help**

{*\$project*: {<specification(s)>}}

- ▶ The specification can be an existing field name followed by a single value indicating the inclusion or exclusion of fields
- ▶ Or it can be a field name (existing or new) followed by an expression to compute the value of the field

<field>: <expression>

- ▶ In the expression, existing field from incoming document can be accessed using field path: “*\$fieldname*”

<https://powcoder.com>

Add WeChat powcoder

\$project examples

- Find the age of each title in the collection, where the age is defined as the duration between the last and the first revision of that title, assuming the timestamp is of ISODate type

```
db.revisions.aggregate([
  { $group: { _id: "$title",
    first: { $min: "$timestamp" },
    last: { $max: "$timestamp" } } },
  { $project: { _id: 0,
    title: "$_id",
    age: { $subtract: [ "$last", "$first" ] } } }
])
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

\$group then \$project

```
{_id:ObjectId("..."), title: "DT", timestamp:"2016-07-01 00:03:46.000Z", ... }
{_id:ObjectId("..."), title: "HC", timestamp:"2016-07-01 00:55:44.000Z", ... }
{_id:ObjectId("..."), title: "DT", timestamp:"2016-07-15 12:22:35.000Z", ... }
{_id:ObjectId("..."), title: "HC", timestamp:"2016-07-28 00:03:58.000Z", ... }
{_id:ObjectId("..."), title: "DT", timestamp:"2016-07-28 00:20:19.000Z", ... }
```

Assignment Project Exam Help

<https://powcoder.com>

```
{ $group: { _id: "$title",
             first: { $min: "$timestamp" },
             last: { $max: "$timestamp" } } },
```

```
{_id:"DT", first:"2016-07-01 00:03:46.000Z", last:"2016-07-28 00:20:19.000Z"}
{_id:"HC", first:"2016-07-01 00:55:44.000Z", last:"2016-07-28 00:03:58.000Z"}
```

Add WeChat powcoder

```
{ $project: { _id: 0,
               title: "$_id",
               age: $subtract: [ "$last", "$first" ] } }
```

```
{title: "DT", age:2333793000}
```

```
{title: "HC", age:2329694000}
```

We can combine multiple operators

```
db.revisions.aggregate([
  {$group: {_id: "$title",
            first: {$min: "$timestamp"},
            last: {$max: "$timestamp"} }},
  {$project: {_id: 0,
              title: "$_id",
              age: {
                $divide: [
                  {
                    $subtract: ["$last", "$first"],
                    86400000
                  }
                ]
              },
              age_unit: {$literal: "day"}
            }
        }
  ])
```

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

\$sort, \$skip, \$limit and \$count stages

- **\$sort** stage sorts the incoming documents based on specified field(s) in ascending or descending order
 - ▶ The function and format is similar to the sort modifier in **find** query
 - ▶ { **\$sort**: { <field1>: <sort order>, <field2>: <sort order> ... } }
- **\$skip** stage skips over given number of documents
 - ▶ The function and format is similar to the skip modifier in **find** query
 - ▶ { **\$skip**: <positive integer> }
- **\$limit** stage limits the number of documents passed to the next stage
 - ▶ The function and format is similar to the limit modifier in **find** query
 - ▶ { **\$limit**: <positive integer> }
- **\$count** stage counts the number of documents passing to this stage
 - ▶ The function and format is similar to the count modifier in **find** query
 - ▶ { **\$count**: <string> }
 - ▶ String is the name of the field representing the count

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

`$sample` and `$out` stages

- The `$sample` stage randomly selects given number of documents from the previous stage
 - ▶ `{ $sample: { size: <positive integer> } }`
 - ▶ Different sampling approaches depending on the location of the stage and the size of the sample and the collection
 - ▶ May fail due to memory constraints
- The `$out` stage writes the documents in a given collection
 - ▶ should be the last one in the pipeline
 - ▶ `{ $out: "<output-collection>" }`

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

\$lookup stage

- New aggregation stages are added with major releases.
- **\$lookup** stage is added since 3.2 to perform left outer join between two collections
 - ▶ The collection already in the pipeline (maybe after a few stages)
 - ▶ Another collection (could be the same one)
- For each incoming document from the pipeline, the \$lookup stage adds a new **array field** whose elements are the matching documents from the other collection.

```
{ $lookup:
  { from: <collection to join>,
    localField: <field from the input documents>,
    foreignField: <field from the documents of the "from"
collection>,
    as: <output array field>
  }
}
```

Add WeChat powcoder

<https://powcoder.com>

Assignment Project Exam Help

\$lookup stage (cont'd)

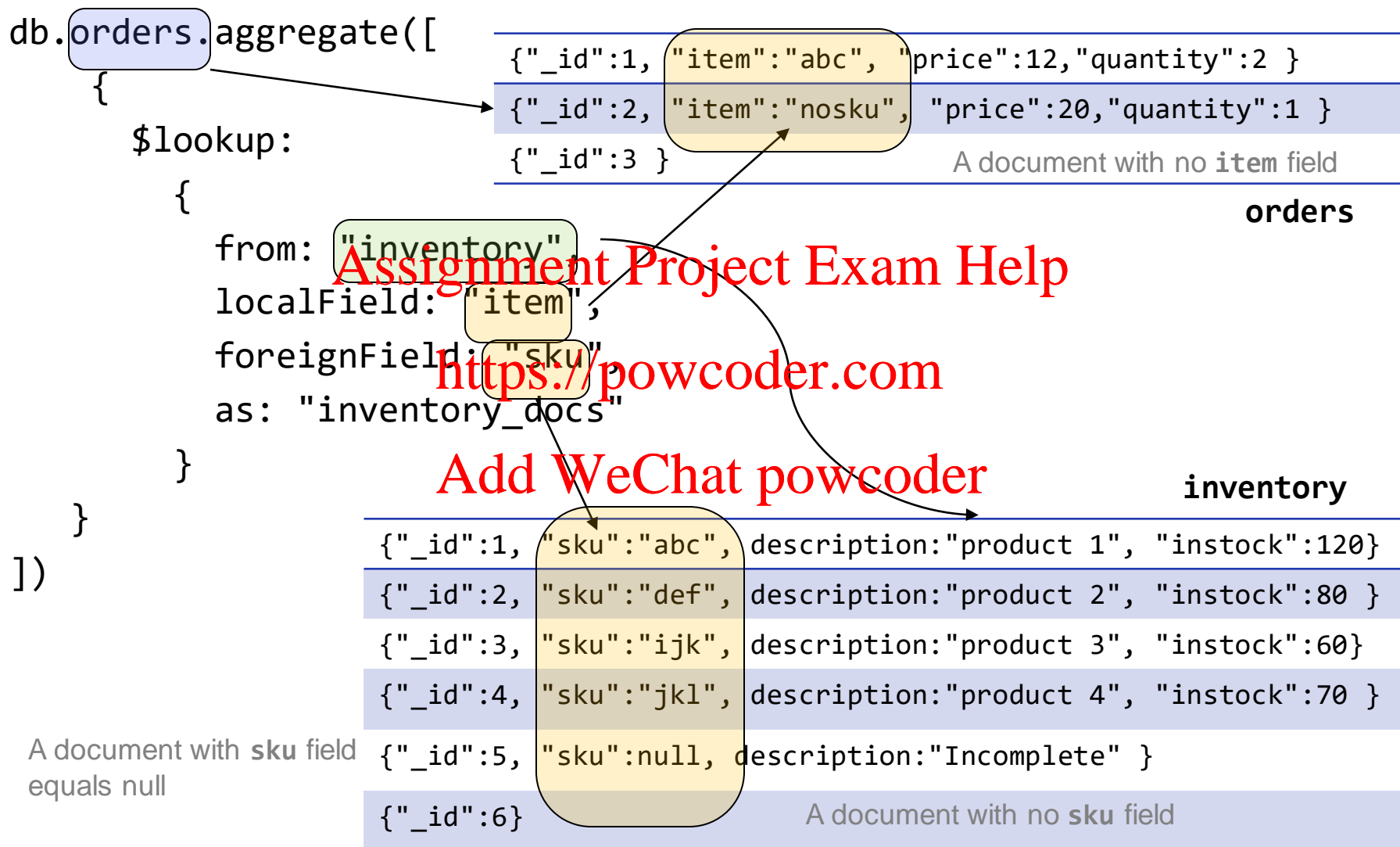
- The output of **\$lookup** stage has the same number of documents as the previous stage
- Each document is augmented with an **array field** storing matching document(s) from the other collection
- The array could contain any number of documents depending on the match, including zero
- Missing local or foreign field is treated as having **null** value

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

\$lookup stage example



https://docs.mongodb.com/manual/reference/operator/aggregation/lookup/#pipe._S_lookup

\$lookup stage example (cont'd)

```
{"_id":1, "item":"abc", "price":12,"quantity":2 }
```

```
{"_id":2, "item":"nosku", "price":20,"quantity":1 }
```

```
{"_id":3 }
```

```
{"_id":1, "sku":"abc", description:"product 1", "instock":120}
```

```
{"_id":2, "sku":"def", description:"product 2", "instock":80 }
```

```
{"_id":3, "sku":"ijk", description:"product 3", "instock":60 }
```

```
{"_id":4, "sku":"jkl", description:"product 4", "instock":70 }
```

```
{"_id":5, "sku":null, description:"Incomplete" }
```

```
{"_id":6}
```

Non exists field matches null and non exists field

```
{"_id":1, "item":"abc", "price":12,"quantity":2,  
  "inventory_docs": [  
    { "_id":1, "sku":"abc", description:"product 1", "instock":120 } ] }
```

```
{"_id":2, "item":"nosku", "price":20,"quantity":1,  
  "inventory_docs" : [ ] }
```

An empty array for no matching from other collection

```
{"_id":3, "inventory_docs" : [  
  { "_id" : 5, "sku" : null, "description" : "Incomplete" },  
  { "_id" : 6 } ] }
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Dealing with data of array type

- To aggregate (e.g. grouping) values in an array field, it is possible to flatten the array to access individual value
- **\$unwind** stage flattens an array field from the input documents to output a document for *each* element. Each output document is the input document with the value of the array field replaced by the element.
 - ▶ `{ $unwind: <field path> }`
- Behaviour
 - ▶ Input document:

```
{ "_id" : 1, "item" : "ABC1", "sizes" : [ "S", "M", "L" ] }
```
 - ▶ After **\$unwind**: **"\$sizes"**
 - ▶ Becomes 3 output documents:

```
{ "_id" : 1, "item" : "ABC1", "sizes" : "S" }  
{ "_id" : 1, "item" : "ABC1", "sizes" : "M" }  
{ "_id" : 1, "item" : "ABC1", "sizes" : "L" }
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

\$unwind example

- Find the number of items that are available in each size

```
db.inventory.aggregate( [  
  { $unwind : "$sizes" },  
  { $group:{_id: "$sizes", item_count: {$sum:1}} }  
] )
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



\$unwind then \$group

```
{ "_id" : 1, "item" : "ABC", "sizes": [ "S", "M", "L" ] }  
{ "_id" : 2, "item" : "EFG", "sizes" : [ ] }  
{ "_id" : 3, "item" : "IJK", "sizes": "M" }  
{ "_id" : 4, "item" : "LMN" }  
{ "_id" : 5, "item" : "XYZ", "sizes" : null }
```

Assignment Project Exam Help

```
{ $unwind : "$sizes" },
```

<https://powcoder.com>

```
{ "_id" : 1, "item" : "ABC", "sizes": "S" }  
{ "_id" : 1, "item" : "ABC", "sizes": "M" }  
{ "_id" : 1, "item" : "ABC", "sizes": "L" }  
{ "_id" : 3, "item" : "IJK", "sizes": "M" }
```

Add WeChat powcoder

```
{ $group:{_id: "$sizes",  
item_count: {$sum:1}}
```

```
{ "_id" : "S", "item_count": 1 }  
{ "_id" : "M", "item_count": 2 }  
{ "_id" : "L", "item_count": 1 }
```


Aggregation Operators

- A few aggregation stages allow us to add new fields or to given existing fields new values based on expression
 - ▶ In **\$group** stage we can use various *operators* or *accumulators* to compute values for new fields
 - ▶ In **\$project** stage we can use operators to compute values for new or exiting fields
- There are many predefined operators for various data types to carry out common operations in that data type
 - ▶ Arithmetic operators: **\$mod**, **\$log**, **\$sqrt**, **\$subtract**, ...
 - ▶ String operators: **\$concat**, **\$split**, **\$indexofBytes**, ...
 - ▶ Comparison operators: **\$gt**, **\$gte**, **\$lt**, **\$lte**,...
 - ▶ Set operators: **\$setEquals**, **\$setIntersection**, ...
 - ▶ Boolean operators: **\$and**, **\$or**, **\$not**, ...
 - ▶ Array operators: **\$in**, **\$size**, ..

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Aggregation vs. Query operators

- There is another set of operators that can be used in **find/update/delete** queries or the **\$match** stage of an aggregation
 - ▶ E.g. **\$gt**, **\$lt**, **\$in**, **\$all**...
- The set is smaller and are different to the operators used in **\$group** or **\$project** stage
- Some operators look the same but have different syntax and slightly different interpretation in query and in aggregation.
 - ▶ E.g. **\$gt** in query looks like
`{age: {$gt:18}}`
 - ▶ **\$gt** in **\$project** stage looks like:
`{over18: {$gt:["$age", 18]}}`

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat: powcoder

Aggregation Behaviour

- It operates on a single collection (before 3.2)
 - ▶ Join can be performed using a particular operator **\$lookup**
- It logically passes the entire collection into the pipeline
- Early filtering can improve the performance
- **\$match** and **\$sort** operator are able to use index if placed at the beginning of the pipeline

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Summary

- MongoDB stores data as BSON document
- Retrieving data from MongoDB are usually achieved through
 - ▶ **find** query
 - ▶ **aggregate** pipeline
- **find** query targets a single collection, it supports condition on any field
- **aggregate** pipeline is able to access other collection(s)
- Both provides rich set of operators
- **update/insert/delete** operation guarantees document level atomicity
- None standard query API, set of operators are growing

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

References

■ BSON types

- ▶ <https://docs.mongodb.com/manual/reference/bson-types/>

■ Aggregation Pipelines

- ▶ <https://docs.mongodb.com/manual/core/aggregation-pipeline/>

■ Aggregation operators

- ▶ <https://docs.mongodb.com/manual/reference/operator/aggregation/>

Assignment Project Exam Help

Add WeChat powcoder

