

MongoDB: BDD Orientada a Documentos

No vamos a trabajar con tablas como en SQL, sino colecciones de documentos. Un documento esta formado por varias entradas de tipo clave valor, como un JSON por ejemplo:

JSON

```
{
  id_: ObjectId("0123456789")
  created_at: ISODate("2019-06-26T00:00:00.000Z"),
  text: "Tweet que no existe.",
  user_id: "102510",
  is_quote_status: false,
  retweet_count: 0,
  display_text_range: [
    {...}
  ],
  entities: {
    hashtags: [],
    ...
  },
  in_reply_to_status_id: null
}
```

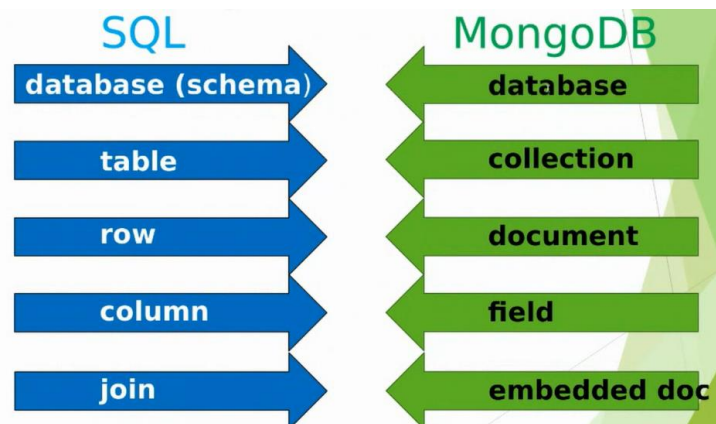
Annotations for the JSON example:

- string (points to `id_`)
- string (points to `created_at`)
- string (points to `text`)
- bool (points to `is_quote_status`)
- número (entero o decimal) (points to `retweet_count`)
- array (points to `display_text_range`)
- objeto anidado (points to `entities`)
- valores nulos (= undefined) (points to `in_reply_to_status_id: null`)

- El id es una especie de primary key para identificar a cada documento
- Los objetos anidados son las **entidades**

Aclaración: Mongo no guarda exactamente JSONs sino BSONS (Binary JSON), por el tipo de almacenamiento. EN BSON los Nulls se representan con no existir

Comparativa de como se llama en una bdd relacional vs en MongoDB



En vez de usar joins, como mongodb no tiene una forma sencilla de juntar los datos, acá vamos a tener documento embebidos.

A diferencia de SQL, no tenemos que crear ninguna tabla. Simplemente decimos que queremos utilizar una colección en particular y accedemos a ella directamente.

- `db.createCollection("conductores");`
- `db.createCollection("autos");`

Cómo armamos documentos embebidos

En mongoDB NO existe el join. Lo que hacemos es poner un documento entero dentro de otro, en lugar de tener una referencia a otro documento.



Acá estoy guardando todo el documento de conductores en dueño

Consultas en Documentos

1. Buscar documentos

La palabra clave para buscar documentos en MongoDB es **find**. Se puede usar en conjunto con sort y limit para obtener una cantidad n de documentos nada más.

```
1 db.<collection>
2   .find(<query>, <projection>)
3   .sort(<order>)
4   .limit(n)
```

Ejemplo: db.tweets.find().limit(5)

a. Búsqueda por igualdad

Buscamos los documentos que tienen cierto campo igual al que indicamos. Podemos buscar de a más de un campo:

Ejemplos:

```
3 db.tweets.find(
4   { _id: ObjectId("1143670209296785408") }
5 )
6
7 db.user.find(
8   { _id: "1143670209296785408" },
9   { _id: 0, text: 1 }
10 )
```

a. Búsqueda por and y or

Otra cosa que podemos hacer es utilizar **AND y OR**, utilizando la notación **\$and: [...], \$or: [...]**

```
1 db.tweets.find({ $or: [
2   { user_id: "Juan" },
3   { retweet_count: 2183 }
4 ] })
5
6 db.tweets.find(
7   { user_id: "Juan", retweet_count: 2183 })
8
9 db.tweets.find({ $and: [
10  { user_id: "818839458" },
11  { retweet_count: 2183 }
12 ] })
```

Acá vemos un ejemplo de or y dos notaciones posibles para el and.

b. Búsqueda por comparación: ==, >, >=, <, <=

Siguen la misma notación que el and y el or, con el \$ y el [], teniendo las siguientes palabras claves:

- ==: eq
- >: gt
- >=: gte
- <: lt
- <=: lte

```
db.tweets.find(
  { retweet_count: { $gte: 2183 } })
```

Si queremos encontrar un documento donde cierto campo este entre dos valores, debemos combinar el <= y >= con un and.

Ejemplo:

```
collection.find(
{
  "$and": [
    {
      "retweet_count": { "$gt": 1000 },
      "retweet_count": { "$lt": 2000 },
      "user.location": "Frias, Santiago del Estero"
    }
  ]
},
```

c. Búsqueda de documentos embebidos

```

16 db.tweets.find(
17   { retweet_count: { $gte: 2183 } })
18
19 db.tweets.find(
20   {"user.name": "Juan"},
21   {"name": "$user.name"}
22 )
23
24
25 db.tweets.find(
26   {"retweet_count": { $gt: "$fa" } }
27 )

```

Para buscar en un campo de un documento embebido, debemos usar ña notación: **“doc.field”**

d. Búsqueda con regex

```

29 db.tweets.find( {"user.name": /Juan/ } ) \
30
31 db.tweets.find( {"user.name": { $regex: "Juan" } })
32
33 db.tweets.find( {"user.screen_name": /juan/i } )

```

e. Búsqueda según in o nin (not in): palabra clave \$in o \$nin

```

35 db.tweets.find( {"user.name":{$in:["Juan","Laura"]}} )
36
37 db.tweets.find( {"user.name": {$nin: ["Pablo"]}} )

```

f. Búsqueda por arrays

Algún elemento en el array: palabra clave \$in

```

3 db.tweets.find({ activities: "swimming" })
4
5 ▼ db.tweets.find({"entities.hashtags":
6 ▼   { $elemMatch:
7     {"text": {$in: ['futbol', 'copaamerica']}}
8   }
9 })

```

Todos los elementos del array: palabra clave \$all

```

db.tweets.find(
  {activities: {$all: ["swimming", "boxing"]}}
)

```

Por tamaño de array: palabra clave \$size

```

db.tweets.find( {activities: {$size: 3} } )

```

El find puede ser seguido por count, lo que colapsa la consulta en un valor numérico, por un limit, que trae una cantidad n de resultados, y también se pueden ordenar.

- Count: **db.<collection>.find().count()**
- Sort: **db.<collection>.find().sort({fieldname: order})**

Ejemplos:

```
db.tweets.find().count()

db.tweets.find().sort( {created_at: 1} )

db.tweets.find().sort( {retweet_count: -1} )
```

2. Agregaciones de documentos

En MongoDB, se refiere a las agregaciones como un **pipeline**. Es como una tubería de procesamiento de datos en la que los documentos de una colección fluyen a través de una **serie de etapas secuenciales**

Cada pipeline de agregación consta de una o mas etapas. Cada etapa toma como entrada un conjunto de documentos y produce como salida otro conjunto de docs, que se pasa a la siguiente etapa en la pipeline.

Ejemplo de una pipeline:



Notación

```
db.collection.aggregate(
  [
    { $stage1 },
    { $stage2 },
    ...
    { $stageN },
  ]
)
```

agregation pipeline

Las etapas más comunes incluyen:

- **\$match**: filtra los docs que coinciden con ciertos criterios.

Permite realizar una query al igual que find()

Notación: { **\$match**: { <query> } }

Ejemplos:

```
db.tweets.aggregate( [ { $match: { "user.name": "Juan" } } ] )

db.tweets.aggregate( [ { $match: {
  "entities.hashtags.text": { $eq: "futbol" }
} } ] )
```

- **\$proyect**: redefine la estructura de los documentos, seleccionando campos específicos y/o añadiendo nuevos campos

Notación: { \$project: { <field>: <1, 0 or expression> ... } }

Ejemplos:

```
db.tweets.aggregate( [ { $project: { user: 1 } } ] )
```

Acá estoy seleccionando el campo específico user.

Y utilizando **\$project** en conjunto con **\$addFields** podemos crear campos nuevos

Notación:

{ \$project: { <field>: <1, 0 or expression> ... } }

{ \$addFields: { <field>: <expression> ... } }

Ejemplos:

```
db.tweets.aggregate( [ { $project: { hashtag_count:
  { $size: "$entities.hashtags" }
} })

db.tweets.aggregate( [ { $addFields: { hashtag_count:
  { $size: "$entities.hashtags" }
} })
```

- **\$limit** y **\$skip**: limita la cantidad de docs que pasan a la siguiente etapa y skip omite un número especificado de documentos

Ejemplo

```
db.tweets.aggregate(
  [
    { $limit: 1 }
  ]
)

db.tweets.aggregate(
  [
    { $skip: 32 }
  ]
)
```

- **\$sort**: ordena los docs según un campo especificado

```
db.country.aggregate(
  [
    { $sort: { area: -1 } }
  ]
)
```

- **\$group**: agrupa docs por un campo especificado, y puede realizar operaciones de agregación (suma, avg, etc.)

Notación:

{ \$group: {

 _id: <expression> ,

 <field>: <expression>

```
}}
```

Y se le pueden aplicar acumuladores: **avg, sum, max, min, first, last**

El **id de agrupación** es **por que** voy a agrupar. Al colocar el \$ lo que hace es usar ese campo, sino me lo toma como valor literal.

Es decir, si yo escribo `_id: "user_name"` me va a buscar un `_id` que sea `user_name`. En cambio `_id: "$user_name"` agrupa por nombre de usuario.

Ejemplo:

```
31 db.country.aggregate( [ { $group: {  
32   _id: "$user.name",  
33   avg_retweeters: {  
34     $avg: "$retweet_count"  
35   }  
36 } ] )
```

- **\$unwind:** "Deconstruye" una lista de elementos en sus elementos individuales

Notación:

```
db.country.aggregate([  
  { $unwind: <field> }  
])
```

- **\$lookup:** simula operaciones de **join** con otras colecciones

3. ABM de documentos

3.1 Insertar documentos

Podemos insertar un elemento o muchos, con One o Many.

Notación:

```
db.<collection>.insertOne( <document> )  
db.<collection>.insertMany( [ <document1>, <document2>, ... ] )
```

Ejemplo:

```
3 db.user.insertOne(  
4   {user_id: "999999", ...}  
5 ) -> ID
```

3.2 Actualizar documentos

La palabra clave para la actualización de documentos es Update. También podemos actualizar uno o muchos documentos a la vez.

Podemos agregar nuevos atributos (como si agregáramos una nueva columna en SQL) simplemente usando un **db.autos.update(...)**

En la sintaxis del update vamos a tener un filtro, entonces puede realizar la actualización sobre los cuales aplicamos el filtro

Es por ello, que para estas BDD, **no necesariamente todos los objetos necesitan tener los mismos datos**, porque algunos documentos van a poder tener un atributo q yo agregue, pero otros no.

Notación:

db.<collection>.updateOne(<filter>, <update>, {upsert: <boolean>})

db.<collection>.updateMany(<filter>, <update>, {upsert: <boolean>})

Ejemplo:

```
db.user.updateOne(
  {"user.name": /Juan/i}, <documento_completo>
)
```

Palabras claves

- **\$set y \$unset:** Agrega o elimina un campo, sin afectar los otros del documento

Ejemplos:

```
db.tweets.updateOne( {id_: ID}, {$set: {is_fake: true}} )

db.tweets.updateMany(
  {is_fake: true},
  {
    $set: {is_real: false},
    $unset: {is_fake: 1}
  }
)
```

- **\$inc:** modifica un campo numérico en la cantidad deseada

Ejemplos:

```
db.tweets.updateOne( {id_: ID}, {$inc: {retweet_count: 10}} )
```

3.3 Eliminar documentos

La palabra clave es **delete**, y también puedo eliminar uno o muchos documentos, con deleteOne o deleteMany.

Notación:

db.<collection>.deleteOne(<filter>)

db.<collection>.deleteMany(<filter>)

Ejemplo:

```
db.tweets.deleteOne(
  {is_real: false}
)
```


Ejemplos de queries en MongoDB Compass

En la IDE que vamos a usar tenemos el siguiente formato:

The screenshot shows the MongoDB Compass query builder interface. At the top, there is a filter icon, a dropdown menu, and buttons for 'Explain', 'Reset', 'Find', and 'Options'. Below this, the 'Project' field is set to '{}'. The 'Sort' field is set to '{ field: -1 } or [['field', -1]]'. The 'Max Time MS' field is set to '60000'. The 'Collation' field is set to '{ locale: 'simple' }'. The 'Skip' field is set to '0' and the 'Limit' field is set to '0'. The 'Index Hint' field is set to '{ field: -1 }'.

El primer set de { } es el campo por el que queremos filtrar. Por ejemplo, podemos buscar un documento por su `_id`.

En Project es un **idem del Select en SQL**, donde digo con que campos del docs me quiero quedar cuando me devuelva los resultados

Ejemplo:

The screenshot shows the MongoDB Compass query builder interface for a specific query. The breadcrumb path is 'tweets.cbukni.mongodb.net > tweets > tweets'. The 'Documents' tab is selected, showing '100.0K' documents. The query is set to '{_id: "1143929112639279104"}'. The 'Project' field is set to '{text: true}'. The 'Sort' field is set to '{ field: -1 } or [['field', -1]]'. The 'Max Time MS' field is set to '60000'. The 'Collation' field is set to '{ locale: 'simple' }'. The 'Skip' field is set to '0' and the 'Limit' field is set to '0'. The 'Index Hint' field is set to '{ field: -1 }'. Below the query builder, there is an 'EXPORT DATA' button and a pagination bar showing '1 - 1 of 1'. The result is displayed as a JSON object: '{_id: "1143929112639279104", text: "Maravilloso! https://t.co/tl90vBYM4G"}'.

A menos que se indique lo contrario, el `_id` me lo devuelve siempre. Para que no aparezca, el Project hay que poner `{_id: 0}`

Ejercicios del Taller

1. Hallar los tweets del usuario con user id '818839458'.

```
db.tweets.find(  
  {user_id: "818839458"}  
)
```

2. Hallar aquellos tweets que tengan más de 500000 retweets.

```
db.tweets.find(  
  {retweet_count: {$gte: 500000}}  
)
```

Lo que ponemos en la IDE es lo siguiente:

```
🔍 {retweet_count: {$gte: 500000}}
```

3. Mostrar la cantidad de retweets de los tweets que se hayan hecho desde Argentina o Brasil.

```
collection.find(  
  {$or:[  
    {"place.country": "Argentina"},  
    {"place.country": "Brasil"}  
  ]},  
  {"retweet_count": 1}  
)
```

4. Hallar los usuarios que tengan tweets con 200000 o más retweets y sean en idioma español.

```
collection.find(  
  {$and:[  
    {"retweet_count": {$gte : 200000}},  
    {"lang": "es"}  
  ]},  
  {"user": 1}  
)
```

```
🔍 {retweet_count: {$gte : 200000},  
  lang: "es"}
```

Formato en la IDE: Project {"user": 1, "_id": 0}

5. Mostrar la cantidad de retweets para los tweets que no se hayan hecho en Argentina ni Brasil, pero si tengan un lugar definido y sean en español.

```
collection.find(  
  {$and:[  
    {"place.country": {$nin: ["Argentina", "Brasil"]}},  
    {"place.country": {$ne: null}},  
    {"lang": "es"}  
  ]},  
  {"retweet_count": 1}  
)
```

También puedo hacer directamente:

```
collection.find(  
  {$and: [  
    {"place.country": {$nin: ["Argentina", "Brasil", null]}},  
    {"lang": "es"}  
  ]},  
  {"retweet_count": 1}  
)
```

6. Mostrar los screen name de aquellos usuarios que tengan “Juan” como parte de su nombre.

```
collection.find(  
  {"user.name": {$regex: "Juan"}},  
  {"user.screen_name": 1}  
)
```

7. Mostrar de los 10 tweets con más retweets, su usuario y la cantidad de retweets.

```
collection.find(
  {},
  {user: 1, retweet_count: 1}
).sort(retweet_count: -1)
.limit(10)
```

Consultas de Agregación

En la pestaña de Aggregations vamos a ir agregando stages, cada uno con un Stage de aggregation

Stage 1: \$match

```
1 /**
2  * query: The query in MQL.
3  */
4 {
5   "user.name": /juan/i
6 }
```

Output after \$match stage (Sample of 10 documents)

```
{ "_id": "1143970625905733633",
  "contributors": null,
  "cooccurrence_checked": true,
  "coordinates": null,
  "created_at": Object,
  "display_text_range": Array (2),
  "entities": Object,
  "favorite_count": 0,
  "full_text": "@sylvyala Que hermosa !!; FEL8Z" }
{ "_id": "1143883926710226944",
  "contributors": null,
  "cooccurrence_checked": true,
  "coordinates": null,
  "created_at": Object,
  "display_text_range": Array (2),
  "entities": Object,
  "favorite_count": 1,
  "full_text": "@doloceb Che gallina sos cagón, le" }
```

+ Add Stage

Y esto me va a ir devolviendo algunos documentos

Puedo agregar un 2do Stage de proyección por ejemplo:

Stage 2: \$project

```
1 /**
2  * specifications: The fields to
3  * include or exclude.
4  */
5 {
6   "user": 1
7 }
```

Output after \$project stage (Sample of 10 documents)

```
{ "_id": "1143970625905733633",
  "user": Object }
{ "_id": "1143883926710226944",
  "user": Object }
{ "_id": "1143900570975113217",
  "user": Object }
```

Lo mas interesante que podemos hacer los las Aggregations es **agrupar con \$group**

1. Mostrar de los 10 tweets con más retweets, su usuario y la cantidad de retweets. Ordenar la salida de forma ascendente.

```
1 collection.aggregate([
2   {
3     $project: {
4       user: 1,
5       retweet_count: 1
6     }
7   },
8   {
9     $sort: { retweet_count: -1 }
10  },
11  {$limit: 10}
12 ])
```

2. Encontrar los 10 hashtags más usados.

```

1 ▼ collection.aggregate([
2 ▼ {
3   $unwind:
4 ▼   {
5     path: "$entities.hashtags",
6     preserveNullAndEmptyArrays: false
7   }
8 },
9 ▼ {
10   $group:
11 ▼   {
12     _id: "$entities.hashtags.text",
13     cantUsos: {
14       $count: {}
15     }
16   }
17 },
18 ▼ {
19   $sort:
20 ▼   {
21     cantUsos: -1
22   }
23 },
24 ▼ {
25   $limit:
26     10
27 }
28 ])

```

3. Por cada usuario obtener una lista de ids de tweets y el largo de la misma.

```

1 ▼ collection.aggregate([
2 ▼ {
3 ▼   $group: {
4     _id: "$user",
5     lista_de_tweets: {
6       $push: "$_id"
7     },
8     largo: {
9       $count: {}
10    }
11  }
12 }
13 ])

```

El operador push es como un append, lo que hace es armarme una lista.

4. Hallar para cada intervalo de una hora cuantos tweets realizo cada usuario.

```

1 ▾ collection.aggregate([
2 ▾ {
3     $match:
4         //extraigo la hora de la fecha
5     {
6         created_at: {
7             $gt: ISODate("2019-06-26T01:00:00.000")
8         }
9     }
10 },
11 ▾ {
12     $project: //me quedo con el usuario y la hora del tweet
13     {
14         user: 1,
15         hora: {
16             $hour: "$created_at"
17         }
18     }
19 },
20 ▾ {
21     $group: //agrupo por usuario y por hora. Hago un count y
22             //me queda la cant de tweets por hora por usuario
23     {
24         _id: {
25             user_name: "$user_name",
26             hora: "$hora"
27         },
28         cantidad: {
29             $count: {}
30         }
31     }
32 }
33 ])
```

5. Hallar la máxima cantidad de retweets totales que tuvo algún usuario.

```

1 ▾ collection.aggregate([
2 ▾ {
3     $group: {
4         _id: "$user",
5         cantidad_retweets: {
6             $sum: "$retweet_count"
7         }
8     }
9 },
10 ▾ {
11     $sort:
12     {
13         cantidad_retweets: -1
14     }
15 },
16 ▾ {
17     $limit:
18         //dice algun usuario, me quedo con 1
19         1
20 },
21 ▾ {
22     $project: {
23         cantidad_retweets: 1
24     }
25 }
26 ])
```

6. Encontrar a los 5 usuarios más mencionados. (les hicieron @)

7. Hallar la cantidad de retweets promedio para los tweets que se hayan hecho desde Argentina y aquellos que no.