



Máster en Data Science. URJC
Esta es la asignatura

Este es el título.
AUTOR 1, AUTOR 2

1 de abril de 2018

Índice

Lista de figuras	3
Glosario	4
1. MongoDB	4
1.1. Diseño	4
1.2. Carga de los datos.	5
2. Análisis	6
2.1. Listado de todas las publicaciones de un autor determinado.	6
2.2. Número de publicaciones de un autor determinado.	6
2.3. Número de artículos en revistas para el año 2017.	6
Conclusiones	7
Bibliografía	8
ANEXOS	8
A. Parseador de XML a JSON	8

Lista de Figuras

1. A la izquierda, la query sin índice, a la deracha, aplicando un índice sobre la columna year. 6

1. MongoDB

La primera parte de la práctica estará enfocada en la Base de datos (BBDD) no relacional MongoDB. Esta se trate de la BBDD más importante actualmente. En el momento de escribir este texto, se encuentra en el 5to de BBDD más utilizadas según la web <https://db-engines.com/en/ranking> y solo por detrás de las BBDD relacionales.

MongoDB tiene las siguientes características:

- Es schemas less, es decir, no es preciso definir un esquema de datos. Aunque no por ello, tenemos que dejar de saber como son almacenados los datos.
- Usa el formato BSON para almacenar los datos, haciendo que sea facil su uso desde lenguajes de programación como Python y JavaScript.
- Permite realizar agregaciones. Incluye el framework aggregation el cuál le permite hacer operaciones realmente potentes.
- Podemos crear índices (incluyendo índices parciales) de cualquiera de sus campos, lo que acelera notablemente las búsquedas.
- Su característica de *sharding* le permite ser muy escalable.
- Cuenta con una gran comunidad de desarrolladores detrás lo cual es una garantía de futuro.
- y muchas más.

Para la realización de esta parte de la práctica vamos a partir de tres documentos JavaScript Object Notation (JSON), los cuales podemos ver la forma de obtenerlos en el anexo A.

Cada uno de estos documentos contiene la información sobre los tres tipos de publicaciones que vamos a usar:

- Articles.
- Inproceedings.
- Incollections.

1.1. Diseño

Aunque se dice que MongoDB se trata de una BBDD *Schemasless*, eso no quita de definir una correcta estructura para almacenar los datos la cual nos permita que las consultas que se hagan sean tanto más sencillas y más rápidas. También es importante el saber definir correctamente los índices que vamos a aplicar. En este sentido hay que tener en cuenta que el uso de índices aumenta el rendimiento de nuestras consultas de búsqueda en la BBDD, pero resulta perjudiciales a la hora de hacer inserciones, además de ocupar espacio en memoria.

También es importante a la hora de definir la forma en la que se guardarán los datos el hecho de como gestionar las relaciones entre documentos. A grandes ragos, y sin entrar en muchos detalles, existen dos estrategias:

- Utilizar referencias entre documentos. Es decir, que uno de los campos de un documento almacene un identificador que sirva para identificar unívocamente a otro documento. Esta relacion puede darse en ambos sentidos. La principal ventaja es el ahorro de espacio al no tener elementos duplicados, así como la facilidad de gestionar una posible actualización de los datos. Por contra, las queries de búsqueda serán más complejas al tener que también hacer una búsqueda por el identificador del segundo documento.

- El segundo método consiste en insertar el segundo documento dentro del primero como si de un campo más se tratase. Esto Facilita las queries de búsqueda ya que recuperamos ambos documentos a la vez, pero tiene la pega de que muy probablemente dupliquemos datos lo cuál se traduce en una mayor cantidad de datos a almacenar y hace más complejo la actualización.

El hecho de optar por una u otra estrategia dependerá de las operaciones que vayamos a realizar sobre nuestra BBDD.

Para el caso que nos ocupa, hemos decidido crear una colección por cada tipo de publicación que vamos a trabajar. También vamos a crear una cuarta colección dentro de la cual cada documento representará un autor. Además, cada uno de estos podocmentos almacenara una información mínima de las publicaciones de dicho autor con el fin de facilitar las queries futuras. Además, se considera que la información de un libro, y la relación de este con sus autores, no debería sufrir cambios (salvo casos excepcionales), por lo que el tener estos datos embebidos dentro de otros no supondrá un gran problema a la hora de gestionar las actualizaciones.

La siguiente tabla muestra los datos de las colecciones con las que vamos a trabajar:

Colección	Información	Indices
Articles	<ul style="list-style-type: none">▪ <code>_id</code>▪ <code>authors</code>: Array con los nombres de los autores.	<ul style="list-style-type: none">▪ <code>_id</code>
Incollections	<ul style="list-style-type: none">▪ <code>_id</code>▪ <code>authors</code>: Array con los nombres de los autores.	<ul style="list-style-type: none">▪ <code>_id</code>
Inproceedings	<ul style="list-style-type: none">▪ <code>_id</code>▪ <code>authors</code>: Array con los nombres de los autores.	<ul style="list-style-type: none">▪ <code>_id</code>
Authors	<ul style="list-style-type: none">▪ <code>_id</code>▪ <code>authors</code>: Array con los nombres de los autores.	<ul style="list-style-type: none">▪ <code>_id</code>

1.2. Carga de los datos.

Los datos han sido cargados usando la herramienta **mongoimport**. Esta herramienta se ejecuta desde la línea de comandos y se instala junto con el propio MongoDB. Hay que tener en cuenta que, en caso de no existir la base de datos o la colección donde se indica que se deben cargar los datos, es el propio MongoDB el encargado de crearlas.

Para cargar nuestros datos se han ejecutado los siguientes comandos:

```
mongoimport --db=dblp --collection=articles articles/articles.json
mongoimport --db=dblp --collection=incollections incollections/incollections.json
mongoimport --db=dblp --collection=Inproceedings Inproceedings/Inproceedings.json
```

Al cargar los datos de esta forma, los *array* de autores y de ee han sido cargados como *arrays* de objetos, los cuales tienen un campo llamado **__VALUE** que contiene la cadena de texto. Es por ello que aplicamos un preprocesamiento para convertir estos campos en *arrays* formados por cadenas de texto. El código utilizado para ello:

```
db.articles.aggregate([{$addFields: {author: '$author.__VALUE', title: '$title.__VALUE',
ee: '$ee.__VALUE'}}], {$out: "articles"})

db.incollections.aggregate([{$addFields: {author: '$author.__VALUE', ee: '$ee.__VALUE'}},
{$out: "incollections"})
```

En el caso de los *inproceedings*, también necesitamos cambiar el tipo del campo año ya que no ha sido reconocido como numérico al realizar la carga.

```

db.inproceedings.find().forEach(function(obj){
  db.inproceedings.update(
    {"_id": obj._id, 'year': {$exists : true}},
    {$set: {"year": NumberInt(obj.year)}}
  })
})

```

2. Análisis

2.1. Listado de todas las publicaciones de un autor determinado.

```

db.authors.aggregate([
  { $match : { _id : "Chin-Wang Tao" } },
  { $project: {publication: {$concatArrays:
    ["$incollections.title", "$articles.title", "$inproceedings.title"]}}}
  }
] )

```

2.2. Número de publicaciones de un autor determinado.

```

db.authors.aggregate([
  { $match : { _id : "Chin-Wang Tao" } },
  { $project: {publication: {$concatArrays:
    ["$incollections.title", "$articles.title", "$inproceedings.title"]}}}
  },
  { $project: {number_of_publications: {$size: "$publications"}}}
] )

```

2.3. Número de artículos en revistas para el año 2017.

```
db.articles.find({year: 2017}).count()
```

En esta simple query podemos ver la importancia que adquieren los índices en cuanto a tiempo de ejecución:

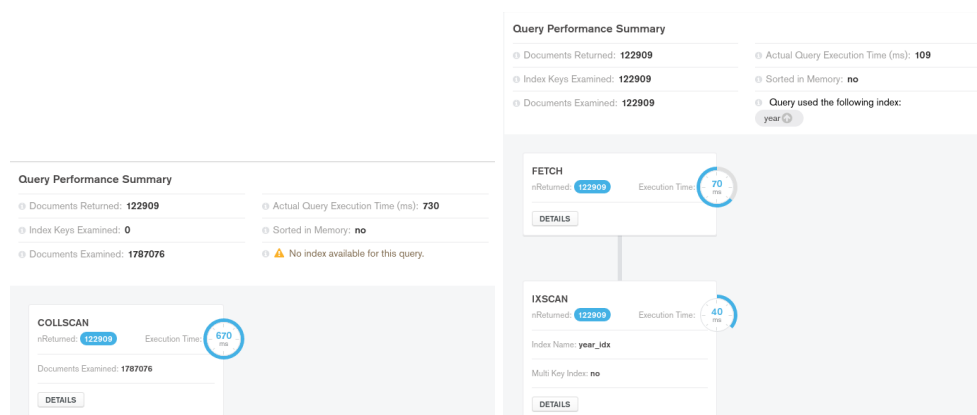


Figura 1: A la izquierda, la query sin índice, a la derecha, aplicando un índice sobre la columna year.

El tiempo de respuesta es aproximadamente 5 veces inferior al hacer uso del índice.

Conclusiones

Aún ninguna.

ANEXOS

A. Parseador de XML a JSON

En construcción.