

# Bases de Datos II

## Bases de datos basadas en documentos

Diego Sevilla Ruiz & Fernando Terroso Sáenz

Universidad de Murcia & Universidad Politécnica de Cartagena

*dsevilla@um.es, fernando.terroso@upct.es*

2024

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
- 4 Uso de MongoDB desde pymongo

# **1. Introducción a las bases de datos de documentos**

# Bases de datos Documentales

- En general, organizadas como un conjunto de **Colecciones** que contienen **documentos**:
  - Tablas relacionales  $\leftrightarrow$  Colecciones
  - Filas (*rows*)  $\leftrightarrow$  Documentos
  - No restringen la estructura de cada documento y no tienen esquema
- Cada colección es un *array* donde cada documento tiene asociado un identificador
- La base de datos puede ver el contenido del documento, y utilizar su información como parte de las búsquedas y actualizaciones
- Documentos  $\Rightarrow$  formatos jerárquicos tipo JSON o XML
- No utilizan SQL y tienen su lenguaje de consulta propio
- Utilizan normalmente **Map-Reduce** para cálculos distribuidos
- Algunas implementan otros lenguajes de consulta y procesado, como N1QL (Couchbase) y Aggregation Framework (MongoDB)
- Bases de datos de Documentos: Couchbase, MongoDB, OrientDB

## **2. Modelado de bases de datos de documentos**

El modelado de datos debe ser:

- Realizado al mayor nivel de abstracción posible
- Independiente de la tecnología subyacente

Sin embargo, en NoSQL:

- Se tiene que tener en cuenta el diseño **distribuido**
- **Optimización guiada por las consultas**

Con respecto al modelo de datos:

- Se mantienen los conceptos de entidad, relación, cardinalidades, etc.
- El modelado relacional se centra en especificar **qué datos tenemos y podemos ofrecer**
- El modelo NoSQL se centra en **optimizar qué consultas vamos a servir**
- Es “barato” **duplicar (desnormalizar)** los datos si con ello se consigue **mayor eficiencia de acceso**

# Representación de CV como tablas

Kleppmann, 2016. *Designing Data Intensive Applications*

<http://www.linkedin.com/in/williamhgates>



**Bill Gates**

Greater Seattle Area | Philanthropy

## Summary

Co-chair of the Bill & Melinda Gates Foundation. Chairman, Microsoft Corporation. Voracious reader. Avid traveler. Active blogger.

## Experience

Co-chair • Bill & Melinda Gates Foundation  
2000 – Present

Co-founder, Chairman • Microsoft  
1975 – Present

## Education

Harvard University  
1973 – 1975

Lakeside School, Seattle

## Contact Info

Blog: [thegatesnotes.com](http://thegatesnotes.com)  
Twitter: @BillGates

user_id	first_name	last_name	summary
251	Bill	Gates	Co-chair of ... blogger.
	region_id	industry_id	photo_id
	us:91	131	57817532

id	region_name
us:7	Greater Boston Area
us:91	Greater Seattle Area

id	industry_name
43	Financial Services
48	Construction
131	Philanthropy

id	user_id	job_title	organization
458	251	Co-chair	Bill & Melinda Gates F...
457	251	Co-founder, Chairman	Microsoft

id	user_id	school_name	start	end
807	251	Harvard University	1973	1975
806	251	Lakeside School, Seattle	NULL	NULL

id	user_id	type	url
155	251	blog	<a href="http://thegatesnotes.com">http://thegatesnotes.com</a>
156	251	twitter	<a href="http://twitter.com/BillGates">http://twitter.com/BillGates</a>



# Representación de relaciones

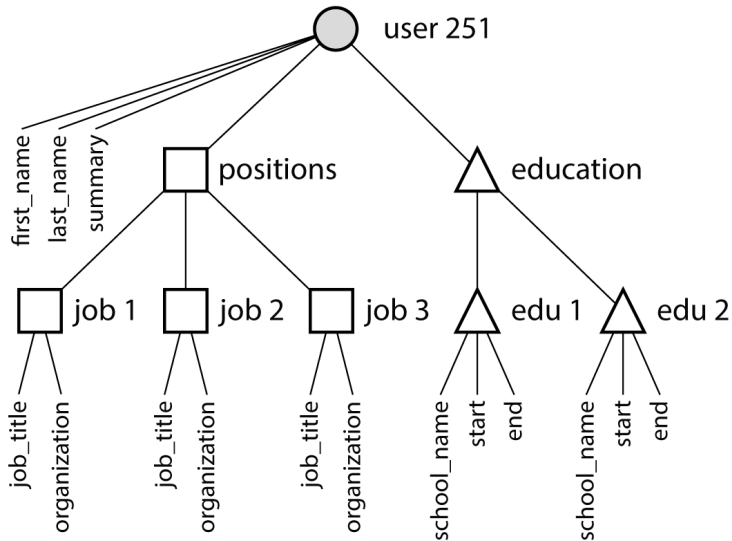
## Relaciones uno a muchos

- Las relaciones uno a muchos (por ejemplo en el CV: positions) en el modelo relacional:
  - Normalización usando varias tablas (Positions con user\_id)
    - Necesidad de más de una tabla
    - Necesidad de uso de JOIN  $\Rightarrow$  ineficiencia
  - Algunos SGBDR ofrecen la posibilidad de tener tipos de datos estructurados y campos XML o JSON. (P. ej. PostgreSQL)
    - Usualmente no se pueden usar para buscar dentro
    - No son estándar

# CV como un documento

```
1 {
2   "user_id": 251,
3   "first_name": "Bill",
4   "last_name": "Gates",
5   "summary": "Co-chair of the Bill & Melinda Gates... Active blogger.",
6   "region_id": "us:91",
7   "industry_id": 131,
8   "photo_url": "/p/7/000/253/05b/308dd6e.jpg",
9   "positions": [
10    {
11      "job_title": "Co-chair",
12      "organization": "Bill & Melinda Gates Foundation"
13    },
14    {
15      "job_title": "Co-founder, Chairman",
16      "organization": "Microsoft"
17    }
18  ],
19  "education": [
20    {
21      "school_name": "Harvard University",
22      "start": 1973,
23      "end": 1975
24    },
25    {
26      "school_name": "Lakeside School, Seattle",
27      "start": null,
28      "end": null
29    }
30  ],
31  "contact_info": {
32    "blog": "http://thegatesnotes.com",
33    "twitter": "http://twitter.com/BillGates"
```

# CV como un árbol (equivalente)



# Representación de Relaciones

## Modelo de documentos

- **Modelo de documentos**

→ analogía del **array/mapa gigante**

- **Conjunto de documentos** (objetos complejos)

- Un **identificador único**, campo *id*
- Búsqueda aleatoria eficiente por clave (**referencia**)
- Estructura jerárquica de sub-documentos contenidos → **agregación**

Más flexibilidad que el modelo relacional:  
(elección entre referencia y agregación)

# Representación de Relaciones (iii)

## Uno a muchos (ii) – NoSQL

- Relaciones **Uno a Muchos** (positions):
  - **Opción 1:** Agregando la tabla positions
  - **Opción 2:** Convertir las empresas en entidades, y utilizar una ***referencia***

# Representación de Relaciones (iii)

## Uno a muchos (ii) – NoSQL

- Relaciones **Uno a Muchos** (positions):
  - **Opción 1:** Agregando la tabla positions
  - **Opción 2:** Convertir las empresas en entidades, y utilizar una *referencia*

### ¡Modelado guiado por el acceso a datos!

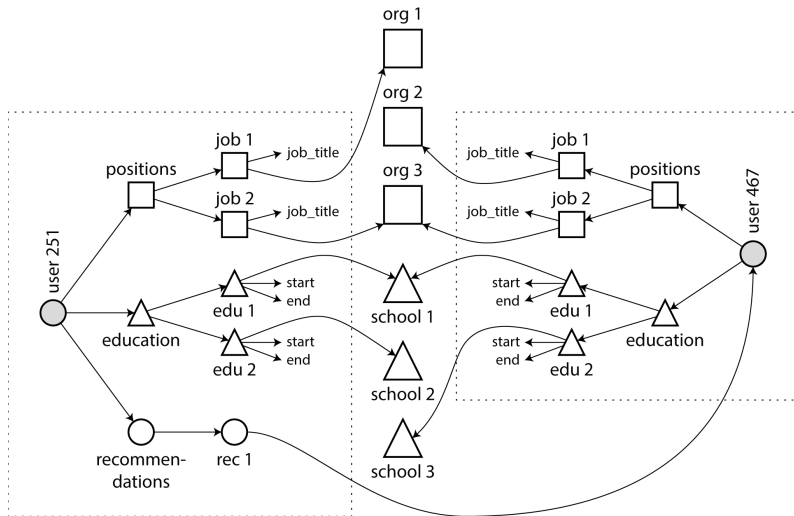
- Si los elementos “muchos” tienen una estructura sencilla ⇒ **Opción 1**
- Si los elementos “muchos” son usualmente **recuperados en una consulta** junto con el elemento “uno” ⇒ **Opción 1**
- Si los elementos “muchos” son relativamente grandes, o bien son recuperados siempre de forma separada ⇒ **Opción 2**

# Representación de Relaciones (iv)

## Muchos a uno y muchos a muchos

- Las relaciones muchos a uno y muchos a muchos:
  - Personas que viven en una región
  - Preguntas que refieren a Tags
- El modelo de documentos no aporta ventajas con respecto al modelo relacional
- Al haber muchas entidades que refieren a otra entidad, la *agregación* daría lugar a mucha **duplicación** (y a problemas de sincronización)
- ⇒ **Referencias** (sobre el ID), similar a una FK en el modelo relacional
  - **Sin embargo**, al no haber **JOINS** la aplicación tiene que hacer más de una petición a la BD

# Muchos a muchos – referencia





### **3. Introducción a MongoDB**

# Introducción

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - **Introducción**
  - Uso básico
  - Consultas MapReduce
  - Validación
  - Índices y desnormalización
- 4 Uso de MongoDB desde pymongo
  - Métodos de búsqueda
  - Métodos de inserción y actualización
  - Framework de agregación
  - Índices
  - Transacciones

# Introducción a MongoDB

- Base de datos documental
- Basado en una representación binaria de JSON, llamada BSON
- Cada base de datos se divide en un conjunto de *colecciones*
- Una colección es un conjunto arbitrariamente grande de documentos
- Cada documento contiene un campo especial `_id`, que guarda un objeto de tipo especial `ObjectID` de BSON
- Permite la definición de bases de datos, colecciones (tablas *hash* de documentos) y subcolecciones
- Permite realizar búsquedas secuenciales, programas MapReduce y un API propietario de consulta (llamado *aggregate*)
- Ofrece un *shell* que acepta JavaScript y permite realizar operaciones directamente sobre la base de datos
- También ofrece clientes en muchos lenguajes de programación

# Uso básico

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - Introducción
  - **Uso básico**
  - Consultas MapReduce
  - Validación
  - Índices y desnormalización
- 4 Uso de MongoDB desde pymongo
  - Métodos de búsqueda
  - Métodos de inserción y actualización
  - Framework de agregación
  - Índices
  - Transacciones

# Uso básico de MongoDB

- MongoDB ofrece un *shell*, aparte de las posibles conexiones remotas a través de *drivers* usando otros lenguajes de programación
- El shell acepta código JavaScript
- Veremos inicialmente su uso
- Existen una serie de objetos predefinidos:

```
1 show databases;  
   use <database>;  
3 db.<collection>; // Permite acceder a la colección
```

- Simplemente nombrando las colecciones se puede acceder a ellas para consulta y actualización
- Inserción de documentos:

```
1 db.colección.insert({< Objeto JSON >});
```

- Se puede obtener los documentos de la base de datos con:

```
1 var objeto = db.colección.findOne({ atributos : valores})  
  // ó encontrar todos los que cumplen la condición:  
3 var objs = db.colección.find({ atributos : valores})
```

- Y actualizarlos:

```
1 db.colección.update(  
  { atributos: valores}, // búsqueda  
3  objeto);              // nuevo objeto
```

- La selección del objeto se suele hacer por su `_id`, que es único

- También permite modificadores:

```
1 db.analytics.update({"url" : "www.example.com"},  
                      {"$inc" : {"pageviews" : 1}})
```

- También \$set, \$unset, \$push (para *arrays*), \$pull (para eliminar elementos de un *array* que cumplen un criterio)

- Las consultas se hacen con find:

```
db.users.find(< criterio > ,  
2           {"username" : 1, "_id" : 0})  
  
// =>  
4 {  
  "username" : "joe",  
6 }
```

- Y se pueden utilizar condicionantes para la búsqueda:

```
db.users.find({"age": {"$gte": 18, "$lte": 30}});  
2 db.users.find({"$or": [ {"age": {"$gte": 18} },  
                           {"permission": true} ] }));
```



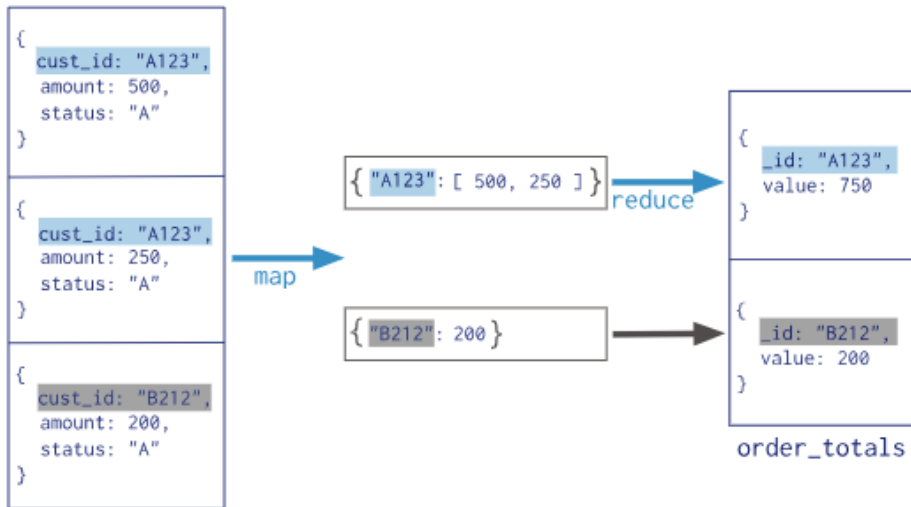
# Consultas MapReduce

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - Introducción
  - Uso básico
  - Consultas MapReduce
  - Validación
  - Índices y desnormalización
- 4 Uso de MongoDB desde pymongo
  - Métodos de búsqueda
  - Métodos de inserción y actualización
  - Framework de agregación
  - Índices
  - Transacciones

- *Map-reduce* es un paradigma de procesamiento de datos que permite condensar grandes volúmenes de datos en resultados agregados útiles.
- Una operación map-reduce consta de tres fases:
  - 1 Se aplica la función Map ( $\text{Map}(k_1, v_1) \rightarrow \text{list}(k_2, v_2)$ ) a cada documento de entrada. La función Map se aplica en paralelo a cada documento (con clave  $k_1$ ) del conjunto de datos de entrada. Esto produce una lista de pares (codificados por  $k_2$ ) para cada llamada.
  - 2 La función Reduce ( $\text{Reduce}(k_2, \text{list}(v_2)) \rightarrow \text{list}(k_3, v_3)$ ) se aplica entonces en paralelo a cada grupo.

# Consultas MapReduce

## Ejemplo de aplicación



- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - Introducción
  - Uso básico
  - Consultas MapReduce
  - **Validación**
  - Índices y desnormalización
- 4 Uso de MongoDB desde pymongo
  - Métodos de búsqueda
  - Métodos de inserción y actualización
  - Framework de agregación
  - Índices
  - Transacciones

- En sus últimas versiones MongoDB ha incorporado funcionalidades que permiten anotar y validar documentos JSON (*JSON Schema*)
- A la hora de crear una colección, podemos especificar en *JSON schema* qué configuración deben seguir los documentos almacenados en ella

# Validación del esquema – Ejemplo

- Vamos a crear una colección estudiantes y vamos a especificar su esquema asociado

```
1 db.createCollection("estudiantes", {  
  validator: {  
3    $jsonSchema: {  
      bsonType: "object",  
5      title: "Validacion del objeto Estudiante",  
      required: [ "direccion", "major", "nombre", "año" ],  
7      properties: {  
        nombre: {  
9          bsonType: "string",  
          description: "'nombre' debe de ser una cadena (string) y es  
            obligatorio"  
11        },  
      },  
    },  
  },  
}
```

## Validación del esquema – Ejemplo (cont.)

```
1      año: {  
3          bsonType: "int",  
          minimum: 2017,  
          maximum: 3017,  
5          description: "'año' debe ser un entero entre [ 2017, 3017 ] y es  
              obligatorio"  
        },  
7      gpa: {  
          bsonType: [ "double" ],  
9          description: "'gpa' debe de ser un double si el campo existe"  
        }  
11    }  
    }  
13  } )
```

## Validación del esquema – Ejemplo (cont.)

- Si ahora intentamos insertar un documento que no cumple con dicho esquema...

```
db.estudiantes.insertOne( {  
  nombre: "Alice",  
  year: Int32( 2019 ),  
  major: "Historia",  
  gpa: Int32(3),  
  direccion: {  
    city: "USA",  
    street: "Fraggel Rock"  
  }  
} )
```



# Validación del esquema – Ejemplo (cont.)

- ...obtendremos un mensaje de error pues el campo gpa se ha insertado como Int cuando debía hacerse como double

```
MongoServerError: Document failed validation
2 Additional information: {
  failingDocumentId: ObjectId("630d093a931191850b40d0a9"),
4 details: {
  operatorName: '$jsonSchema',
6 title: 'Validacion del objeto Estudiante',
  schemaRulesNotSatisfied: [
8   {
    operatorName: 'properties',
10    propertiesNotSatisfied: [
      {
12        propertyName: 'gpa',
        description: "'gpa' debe de ser un doble si el campo existe",
14        details: [
          {
16            operatorName: 'bsonType',
            specifiedAs: { bsonType: [ 'double' ] },
18            reason: 'type did not match',
            consideredValue: 3,
20            consideredType: 'int'
          }
22        ] } ] } ] } }
```

## Validación del esquema – Ejemplo (cont.)

- Sin embargo, la siguiente inserción no provoca ningún fallo

```
db.students.insertOne( {  
  2   name: "Alice",  
      year: NumberInt(2019),  
  4   major: "History",  
      gpa: Double(3.0),  
  6   address: {  
        city: "NYC",  
  8     street: "33rd Street"  
      }  
10 } )
```

# Índices y desnormalización

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - Introducción
  - Uso básico
  - Consultas MapReduce
  - Validación
  - Índices y desnormalización
- 4 Uso de MongoDB desde pymongo
  - Métodos de búsqueda
  - Métodos de inserción y actualización
  - Framework de agregación
  - Índices
  - Transacciones

# Data Access Hitting a Wall



## Current practice based on data download (FTP/GREP)

### Will not scale to the datasets of tomorrow

- You can GREP 1 MB in a second
- You can GREP 1 GB in a minute
- You can GREP 1 TB in 2 days
- You can GREP 1 PB in 3 years.
- Oh!, and 1PB ~5,000 disks
- You can FTP 1 MB in 1 sec
- You can FTP 1 GB / min (~1\$)
- ... 2 days and 1K\$
- ... 3 years and 1M\$
- At some point you need **indices** to limit search
- **parallel** data search and analysis
- This is where databases can help



[slide src: Jim Gray]

# Índices

- Sea con SQL o con datos *raw*, la búsqueda no se puede realizar secuencialmente
- Los **índices** se utilizan para acelerar esta búsqueda
- El concepto de índice está presente en casi todas las bases de datos (también es un eje principal en las NoSQL)
- En el ámbito del SQL, los índices se pueden aplicar a los valores de una o más columnas
- Se usan para acelerar la recuperación de información
- El planificador de consultas usa los índices disponibles
- Es **labor del usuario** definir los índices adecuados
- Los índices aceleran las búsquedas, pero también ocupan espacio
- En SQL se utiliza la construcción `CREATE INDEX`, `ALTER TABLE` o se incluyen en `CREATE TABLE`

# Índices (cont.)

- EXPLAIN para ver si se usa:

```
mysql> explain select * from Posts where Id=5;
2  +-----+-----+-----+-----+-----+-----+...
   | id | select_type | table | type | possible_keys | key  | ...
4  +-----+-----+-----+-----+-----+-----+...
   | 1 | SIMPLE      | Posts | ref  | Id             | Id   | ...
6  +-----+-----+-----+-----+-----+-----+...
```

- Obviamente **TAMBIÉN** se usan en las bases de datos NoSQL

# Índices (cont.)

- En general, hay tres tipos de índices:
  - 1 Basados en árboles balanceados
  - 2 Basados en tablas *hash*
  - 3 Índices *full-text*

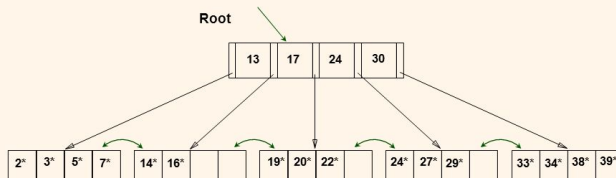
## Índices de Árboles Balanceados

- Mantienen la clave ordenada dentro de un árbol balanceado (B-Tree, B+, etc.)
- Aceleran búsquedas con operadores =, < y > y rangos (BETWEEN)
- También LIKE que no empiecen con comodín
- Puede ayudar en ORDER BY

## Example B+ Tree



- ❖ Search begins at root, and key comparisons direct it to a leaf (as in ISAM).
- ❖ Search for 5\*, 15\*, all data entries  $\geq 24^*$  ...



*\* Based on the search for 15\*, we know it is not in the tree!*

Database Management Systems 3ed, R. Ramakrishnan and J. Gehrke

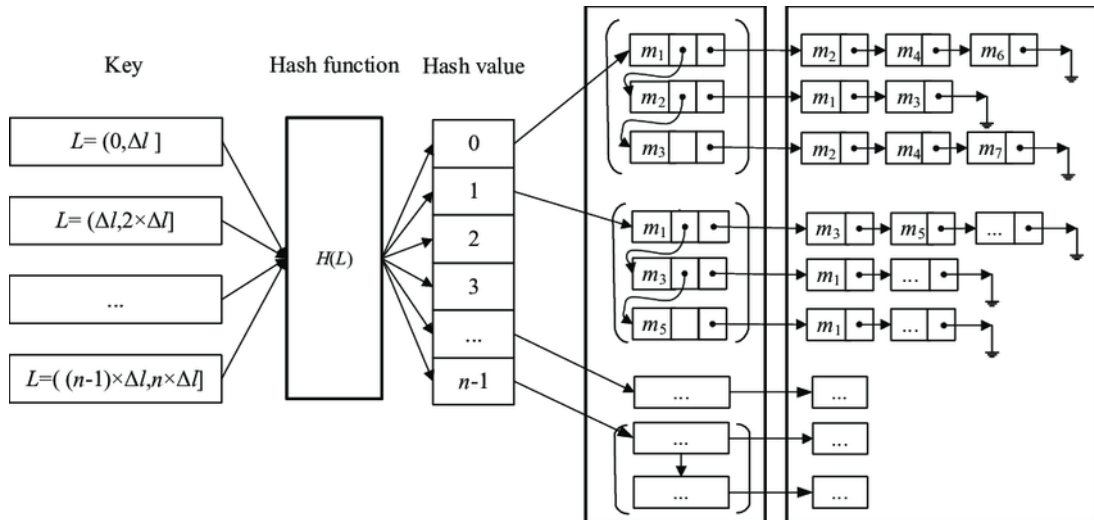
10



### Índices de tablas hash

- Sólo para consultas de igualdad (=) y desigualdad (<=>)
- Son muy rápidos
- No se pueden usar para ORDER BY
- Sólo se pueden usar con valores completos de clave (no permiten prefijos, por ejemplo)

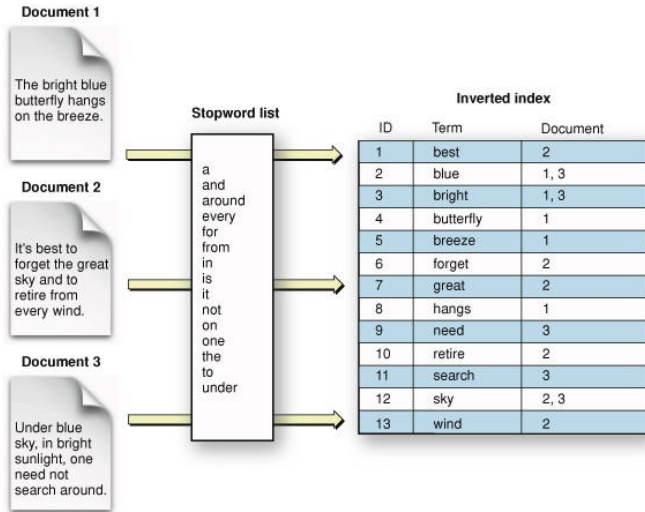
# Índices (cont.)



## Índices *full-text*

- Usados en búsquedas inversas sobre campos textuales
- Dada una palabra o conjunto de palabras, encontrar la fila o filas en las que el campo buscado (el del índice) contiene esas palabras
- Tiene limitaciones en LIKE con comodines (%)
- Por ello los SGBD proveen operadores especiales (p. ej. en MySQL “MATCH(...) AGAINST ...”)

# Índices (cont.)



- El modelo relacional basa el modelado de datos en la **normalización**, en sus distintos niveles
- Esto permite aplicar el principio de “sólo almacenar un dato en un sitio”, lo que lleva a la eficiencia de datos
- Pero a veces necesitamos también **eficiencia de acceso**, y a veces está reñido con la normalización
- El principal escollo para la velocidad es la unión (JOIN) de datos de tablas

- ¿Por qué desnormalizar?
  - **Mantener la historia** – Los datos de los clientes, por ejemplo, pueden cambiar. Necesitamos almacenar los valores de los mismos, por ejemplo, cuando emitimos una factura en el pasado. Para eso hay que copiar los datos de un cliente en ese momento
  - **Mejorar la velocidad de consulta** – A veces, para hacer consultas recurrentes, tenemos que usar uniones de varias tablas. Evitar la unión, por ejemplo, añadiendo claves ajenas a las tablas finales. O bien, como en el caso de `Posts.Tags`, replicando la información
  - **Precalcular valores necesitados de antemano** – Si hay datos que se necesitarán, como medias o totales, podemos calcularlos y almacenarlos de antemano

- Desventajas de desnormalizar:
  - **Espacio en disco** – Al igual que ocurría con los índices, la copia y duplicación de datos ocupa más espacio
  - **Anomalías de datos** – Hay que ser conscientes de que ahora un cambio en un dato puede requerir un cambio en más de una parte de la base de datos para mantener la consistencia. Los SGBDR no están preparados para mantener este tipo de consistencia. Se tienen que utilizar mecanismos de *triggers*, transacciones, etc.
  - **Ralentización de las operaciones de actualización** – Dependiendo del *ratio* entre actualizaciones y consultas, este puede ser un punto importante
  - **Más código, más dependencias, más propenso a errores** – Al necesitar más coordinación, la base de datos requiere más código, hay más dependencias ocultas, lo que puede llevar a más errores y a un mayor coste de mantenimiento

- Ejemplos vistos en Stackoverflow:
  - En la tabla `Posts` tenemos, por ejemplo, el campo `Tags`, que contiene la lista de los tags. Evita tener que hacer un `JOIN` para obtener los nombres de los tags
  - También está el campo `Score`, que aglutina todos los votos positivos recibidos (también en la tabla `Votes`)



## 4. Uso de MongoDB desde pymongo

# Métodos de búsqueda

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - Introducción
  - Uso básico
  - Consultas MapReduce
  - Validación
  - Índices y desnormalización
- 4 **Uso de MongoDB desde pymongo**
  - **Métodos de búsqueda**
  - Métodos de inserción y actualización
  - Framework de agregación
  - Índices
  - Transacciones

# Métodos de búsqueda

- El método de búsqueda principal es `find()`, que tiene muchas opciones
- En general permite especificar:
  - El filtro de búsqueda
  - Ordenación de resultados por algún campo
  - Proyección para no obtener todos los campos del documento
  - Número de resultados máximo (*limit*)
  - Número de elementos iniciales a ignorar (*skip*)
  - El tamaño del *batch*
- Como la variabilidad es muy grande, veremos ejemplos en la hoja Jupyter Notebook de prácticas y también en la documentación

- La función `find()` tiene un gran número de posibilidades para especificar la búsqueda. Se pueden utilizar cualificadores complejos como:
  - `$and`
  - `$or`
  - `$not`
- Estos calificadores unen “objetos”, no valores

# Métodos de búsqueda (cont.)

- Por otro lado, hay otros calificadores que se refieren a valores:
  - \$lt (menor)
  - \$lte (menor o igual)
  - \$gt (mayor)
  - \$gte (mayor o igual)
  - \$regex (expresión regular)

```
posts.find_one({"body": {"$regex" : "[Mm]ongo"}})
```

Condición compuesta (limitando la salida a 10 documentos):

```
1 posts.find({"$and": [{"PostTypeId": 2},  
2                       {"Id" : {"$gte" : 100}} ]})  
3     .limit(10)
```

# Métodos de búsqueda

## Ejemplos introductorios

- Encontrar el primer documento en la colección customers

```
1 myclient = pymongo.MongoClient("mongodb://localhost:27017/")
  mydb = myclient["mydatabase"]
3 mycol = mydb["customers"]
  x = mycol.find_one()
```

- Encontrar todos los documentos en la colección customers

```
for x in mycol.find():
2     print(x)
```

- También podemos indicar qué campos queremos obtener de una determinada búsqueda

```
for x in mycol.find({}, { "_id": 0, "name": 1, "address": 1 }):
2     print(x)
```

# Métodos de búsqueda

## Ejemplos básicos

- Dada la siguiente colección,

```
{ "_id": "manzanas", "cantidad": 5 }  
2 { "_id": "plátanos", "cantidad": 7 }  
  { "_id": "naranjas", "cantidad": 8 }  
4 { "_id": "aguacate", "cantidad": 14 }
```

- Si quisieramos obtener aquellos documentos con un valor en su atributo cantidad superior a 4 deberíamos ejecutar el siguiente comando

```
db.collection.find( { "cantidad": { "$gt": 4 } } )
```

- ¿Qué comando deberíamos lanzar para obtener aquellos que tengan un valor inferior o igual a 5?

# Métodos de búsqueda (cont.)

## Ejemplos básicos

- Suponiendo una colección bios cuyos documentos tienen esta forma:

```
1 {  
3   "_id" : <valor>,  
   "nombre" : { "nombre_propio" : <string>, "primer_apellido" : <string> },  
   "nacimiento" : <ISODate>,  
5   "muerte" : <ISODate>,  
   "profesion": <valor>,  
7   "trabajos": [<string>,...]  
}
```



# Métodos de búsqueda (cont.)

## Ejemplos básicos

- Una opción interesante es el filtrado por rango de fechas en caso de que tengamos atributos de ese tipo

```
2 db.bios.find( { "nacimiento": { "$gt": new Date('1940-01-01'),  
    "$lt": new Date('1960-01-01') } } )
```

- El operador `$in` nos permite buscar documentos con atributos dentro de una lista de valores

```
2 db.bios.find(  
    { "profesion": { "$in": [ 'jedi', 'samurai' ] } }  
)
```

# Métodos de búsqueda (cont.)

## Ejemplos básicos

- También podemos realizar búsquedas en campos anidados

```
1 db.bios.find( { "nombre.primer_apellido": "Skywalker" } )  
db.bios.find(  
3   { "nombre": { "nombre_propio": "Akira", "primer_apellido": "Kurosawa" } }  
  )
```

- Las búsquedas pueden realizarse también sobre campos cuyos valores son arrays

```
db.bios.find({"trabajos": "star_wars"})  
2 db.bios.find({"trabajos": { "$in": [ "star wars", "alien" ] } })  
db.bios.find({"trabajos": { "$all": [ "star wars", "the clone wars" ] } })  
4 db.bios.find({"trabajos": { "$size": 4 } }) #El array 'trabajos' debe tener 4  
   elementos
```

# Métodos de búsqueda (cont.)

## Ejemplos básicos

- El operador `$not` es útil cuando queremos buscar elementos que NO cumplan una condición

*# Buscar documentos cuyo precio no sea superior a 1,99€ o que no tengan el campo "precio"*

```
2 db.inventory.find( { "precio": { "$not": { "$gt": 1.99 } } } )
```

- En el ejemplo de arriba, ¿cuál es la diferencia respecto a usar la condición `{"$gt": 1.99}`?

# Métodos de búsqueda (cont.)

## Ejemplos básicos

- Por último, los operadores \$and, \$or y \$not pueden combinarse para definir condiciones complejas

```
db.inventory.find( {  
  "$and": [  
    {"$or":[{"cantidad":{"$lt":10}}, {"cantidad":{"$gt":50} }]},  
    {"$or":[{"vendido": true}, {"precio": {"$lt": 5}}]}  
  ]  
} )
```

- ¿Qué documentos devolverá la consulta de arriba?

# Ejercicio métodos de búsqueda

- Dada la siguiente colección llamada libros en la base de datos biblioteca

```
{
  "_id": "1",
  "titulo": "Cien años de soledad",
  "autor": "Gabriel García Márquez",
  "genero": ["Realismo mágico", "Literatura latinoamericana"],
  "anio_publicacion": 1967,
  "precio": 15.99,
  "tags": ["Ciencia ficción", "Distopía", "Política"]
},
{
  "_id": "2",
  "titulo": "El Principito",
  "autor": "Antoine de Saint-Exupéry",
  "genero": ["Fantasía", "Literatura infantil"],
  "anio_publicacion": 1943,
  "precio": 10.99,
  "tags": ["Fantasía", "Aventura", "Amistad"]
},
...
```

## Ejercicio métodos de búsqueda (cont.)

- Escribir el código Python que permita realizar las siguientes búsquedas:
  - **Ejercicio 1:** Encontrar todos los libros que sean de genero 'Fantasía' y tengan un precio mayor a 15 euros
  - **Ejercicio 2:** Encontrar todos los libros que tengan el tag 'Fantasía' y que se hayan publicado después de 1950
  - **Ejercicio 3:** Encontrar todos los libros que tengan el tag 'Fantasía' y 'Amistad' o que se hayan publicado a partir de 1982 (inclusive)

# Métodos de inserción y actualización

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - Introducción
  - Uso básico
  - Consultas MapReduce
  - Validación
  - Índices y desnormalización
- 4 Uso de MongoDB desde pymongo
  - Métodos de búsqueda
  - **Métodos de inserción y actualización**
  - Framework de agregación
  - Índices
  - Transacciones

# Métodos de inserción y actualización

- pymongo ofrece métodos para inserción y actualización:
  - `insert_one()`, `insert_many()` (batch)
  - `update_one()` - Permite actualizar un objeto con nuevos campos. El objeto se crea si se pone el parámetro `upsert` a `True`
  - `update_many()` - Permite poner nuevos valores calculados a un conjunto de objetos



# Métodos de inserción y actualización

## Ejemplos básicos de inserción

- A la hora de insertar un elemento podemos obtener su `_id` asignado

```
1 micoleccion = db["clientes"]
midict = { "nombre": "Luke", "direccion": "Tatooine, 83" }
3 x = mycol.insert_one(midict)
print(x.inserted_id)
```

- Para insertar un conjunto de documentos a la vez deberemos incluirlos en un array JSON

```
milista = [
2     { "nombre": "Amy", "direccion": "Apple st 652"},
    { "nombre": "Hannah", "direccion": "Mountain 21"},
4     { "nombre": "Michael", "direccion": "Valley 345"}
]
6 x = micoleccion.insert_many(milista)
#imprimimos us _id asignados
8 print(x.inserted_ids)
```

# Métodos de inserción y actualización

## Ejemplos básicos de actualización

- A la hora de actualizar un documento debemos especificar un criterio de búsqueda y los campos y valores a actualizar

```
micoleccion = db["clientes"]  
2 mifiltro = { "direccion": "Valley 345" }  
nuevosvalores = { "$set": { "direccion": "Rue del percebe, 13" } }  
4 micoleccion.update_one(mifiltro, nuevosvalores)
```

- Para actualizar múltiples documentos la estructura es la misma

```
micoleccion = db["clientes"]  
2 mifiltro = { "direccion": { "$regex": "^S" } }  
nuevosvalores = { "$set": { "nombre": "Rey" } }  
4 x = micoleccion.update_many(mifiltro, nuevosvalores)  
print(x.modified_count, "documentos actualizados.")
```

# Índices, explain()

- La llamada `.explain()` a cualquier búsqueda muestra el plan de ejecución
- Se puede crear un índice si la búsqueda por ese campo va a ser crítica
- Se pueden crear más índices:
  - ASCENDING
  - DESCENDING
  - HASHED
  - y geoespaciales

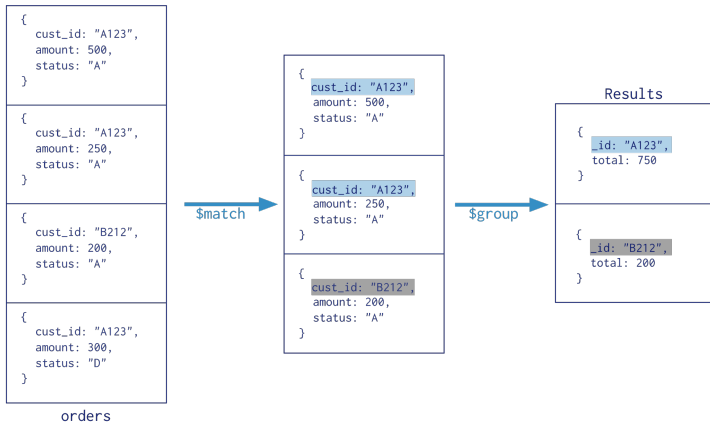
# Framework de agregación

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - Introducción
  - Uso básico
  - Consultas MapReduce
  - Validación
  - Índices y desnormalización
- 4 Uso de MongoDB desde pymongo
  - Métodos de búsqueda
  - Métodos de inserción y actualización
  - **Framework de agregación**
  - Índices
  - Transacciones

- El framework de agregación de MongoDB es una de sus características más potentes
- Permite escribir expresiones, divididas en una serie de etapas, que realizan operaciones como agregaciones, transformaciones y uniones en los datos de las bases de datos MongoDB
- Esto le permite realizar cálculos y análisis a través de documentos y colecciones dentro de su base de datos MongoDB

# Framework de agregación

Collection  
↓  
db.orders.aggregate( [  
  \$match stage → { \$match: { status: "A" } },  
  \$group stage → { \$group: { \_id: "\$cust\_id", total: { \$sum: "\$amount" } } }  
] )



# Framework de agregación

Framework de agregación:

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

```
1 posts.aggregate( [  
    {'$project' : { 'Id' : 1 }},  
3    {'$limit': 20}  
    ])
```

# Framework de agregación (cont.)

```
by_length = posts.aggregate( [  
2   {'$project': {  
    'Body' : {'$ifNull' : ['$Body', '' ]}  
4   }},  
    {'$project' : {  
6     'id' : {'$strLenBytes': '$Body'},  
     'value' : {'$literal' : 1}  
8   }},  
    {'$group' : {  
10    '_id' : '$id',  
    'count' : {'$sum' : '$value'}  
12  }},  
    {'$sort' : { '_id' : 1}}  
14  ])
```



## Framework de agregación (cont.)

Se pueden añadir etapas. En particular, por ejemplo, se puede **filtrar** inicialmente añadiendo al principio:

```
{ '$match': { 'Body' : { '$regex': 'HBase' } } }
```

La construcción `$lookup` permite el acceso a otra (o a la misma) colección. Es equivalente a un **JOIN**

## Framework de agregación (cont.)

P. ej., listar los posts con Score mayor o igual a 40 junto con el usuario que lo ha hecho:

```
1 posts.aggregate( [  
    {'$match': { 'Score' : {'$gte': 40}}},  
3    {'$lookup': {  
        'from': "users",  
5        'localField': "OwnerId",  
        'foreignField': "Id",  
7        'as': "owner"}  
    }  
9 ])
```

# Framework de agregación

## Colección de ejemplo

- Vamos a suponer que contamos con una colección de películas (movies) que contienen documentos que tienen esta estructura:

```
1 {
2   '_id': ObjectId('573a1392f29313caabcbdb497'),
3   'awards': {'nominations': 7,
4             'text': 'Won 1 Oscar. Another 2 wins & 7 nominations.',
5             'wins': 3},
6   'cast': ['Janet Gaynor', 'Fredric March', 'Adolphe Menjou', 'May Robson'],
7   'countries': ['USA'],
8   'directors': ['William A. Wellman', 'Jack Conway'],
9   'fullplot': 'Esther Blodgett is just...
10  'genres': ['Drama'],
11  'imdb': {'id': 29606, 'rating': 7.7, 'votes': 5005},
12  'languages': ['English'],
13  'lastupdated': '2015-09-01 00:55:54.333000000',
14  'plot': 'A young woman ...',
15  'poster': 'https://m.media-amazon.com/images/M/MV5....jpg',
16  'rated': 'NOT RATED',
17  'released': datetime.datetime(1937, 4, 27, 0, 0),
18  'runtime': 111,
19  'title': 'A Star Is Born',
20  'tomatoes': {'critic': {'meter': 100, 'numReviews': 11, 'rating': 7.4},
21              'dvd': datetime.datetime(2004, 11, 16, 0, 0),
```

# Framework de agregación (cont.)

## Colección de ejemplo

```
23     'fresh': 11,  
    'lastUpdated': datetime.datetime(2015, 8, 26, 18, 58, 34),  
    'production': 'Image Entertainment Inc.',  
25     'rotten': 0,  
    'viewer': {'meter': 79, 'numReviews': 2526, 'rating': 3.6},  
27     'website': 'http://www.vcientertainment.com/Film-Categories?product_id=73'},  
    'type': 'movie',  
29     'writers': ['Dorothy Parker (screen play)',  
    'Alan Campbell (screen play)',  
31     'Robert Carson (screen play)',  
    'William A. Wellman (from a story by)',  
33     'Robert Carson (from a story by)'],  
    'year': 1937}
```

## Framework de agregación – Ejemplo (i)

- En primer lugar vamos a realizar una primera búsqueda sobre dicha colección filtrando por el nombre de la película (title) y ordenando por su año (year):

```
movie_collection = db["movies"]
2 pipeline = [
    {
4         "$match": {
            "title": "A Star Is Born"
6         },
    },
8     {
        "$sort": {
10         "year": pymongo.ASCENDING
        },
12     },
]
14 results = movie_collection.aggregate(pipeline)
```

## Framework de agregación – Ejemplo (i) (cont.)

- Si ahora quisieramos limitar la salida a un único documento...

```
pipeline = [  
    {  
        "$match": {  
            "title": "A Star Is Born"  
        }  
    },  
    {  
        "$sort": {  
            "year": pymongo.ASCENDING  
        }  
    },  
    { "$limit": 1 }  
]  
results = movie_collection.aggregate(pipeline)
```

## Framework de agregación – Ejemplo (ii)

- Vamos a suponer ahora que existe otra colección llamada `comments` cuyos documentos tienen este formato

```
{  
  '_id': ObjectId('5a9427648b0beebe69579d3'),  
  'movie_id': ObjectId('573a1390f29313caabcd4217'),  
  'date': datetime.datetime(1983, 4, 27, 20, 39, 15),  
  'email': 'cameron_duran@fakegmail.com',  
  'name': 'Cameron Duran',  
  'text': '...'}  
2  
4  
6
```

## Framework de agregación – Ejemplo (ii) (cont.)

- Esto permitiría añadir un campo `related_comments` a cada película con sus críticas asociadas

```
1 stage_lookup_comments = {  
    "$lookup": {  
3         "from": "comments",  
        "localField": "_id",  
5         "foreignField": "movie_id",  
        "as": "related_comments",  
7     }  
    }  
9 # Limitamos a los primeros 5 documentos  
    stage_limit_5 = { "$limit": 5 }  
11 pipeline = [stage_lookup_comments, stage_limit_5]  
    results = movie_collection.aggregate(pipeline)
```



## Framework de agregación – Ejemplo (iii)

- Una funcionalidad muy interesante es la agrupación de documentos con el comando `$group`
- Vamos a usar dicho comando para contar el número de películas por año, ordenadas de forma ascendente

```
stage_group_year = {  
  2   "$group": {  
      "_id": "$year",  
      4   # contamos el número de películas en el grupo  
      "movie_count": { "$sum": 1 },  
  6   }  
}  
8   stage_match_years = {  
    "$match": {  
      10   "year": {  
          "$type": "number",  
      12   }  
    }
```

## Framework de agregación – Ejemplo (iii) (cont.)

```
    }  
14 }  
    stage_sort_year_ascending = {  
16     "$sort": {"_id": pymongo.ASCENDING}  
    }  
18 pipeline = [stage_match_years, stage_group_year, stage_sort_year_ascending]  
    results = movie_collection.aggregate(pipeline)
```

# Índices

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - Introducción
  - Uso básico
  - Consultas MapReduce
  - Validación
  - Índices y desnormalización
- 4 Uso de MongoDB desde pymongo
  - Métodos de búsqueda
  - Métodos de inserción y actualización
  - Framework de agregación
  - **Índices**
  - Transacciones

# Creación de índices

- En MongoDB, la creación de índices se realiza mediante el método `create_index`

```
db.collection.create_index([(<key and index type specification>)], <options> )
```

- En el siguiente ejemplo se crea un índice con una sola clave (nombre) en orden descendente

```
collection.create_index([("nombre", pymongo.DESCENDING)])
```

- Podemos listar los índices que tenemos creados en una colección con

```
collection.index_information()
```

- En el caso de las colecciones con una alta proporción de escritura por lectura, los índices son caros porque cada inserción y actualización también debe actualizar cualquier índice.

# Nombrado y borrado de índices

- Al crear un índice, puede darle un nombre personalizado.
- Ayuda a distinguir los distintos índices de la colección.
- Para especificar el nombre del índice, incluya la opción de nombre al crear el índice:

```
1 db.collection.create_index(  
    { "<campo>": "<valor>" },  
3    { "name": "<NombreDelIndice>" }  
    )
```

- Para borrar un índice podemos hacerlo con el método `.drop_index`

```
collection.drop_index("<nombreDelIndice>")
```

# Tipos de índices en MongoDB

## Índice de campo único (*single field index*)

- Recogen y ordenan los datos de un solo campo de cada documento de una colección
- Por defecto, todas las colecciones tienen un índice de este tipo sobre su campo “\_id”
- Es posible crear un *single field index* en cualquier campo de un documento (incluyendo los campos de documentos incrustados)
- Cuando se crea un índice, se especifica
  - El campo sobre el que crear el índice
  - El orden de ordenación de los valores indexados, ascendente (1) o descendente (-1)

```
1 db.<collection>.create_index( { "<campo>": <tipoDeOrdenacion> } )
```

- Por ejemplo, si el departamento de recursos humanos necesita a menudo buscar empleados por su apellido. Puede crear un índice en el apellido de empleado para mejorar el rendimiento de esa consulta

```
1 db.empleados.create_index( { "apellido": 1 } )
```

# Tipos de índices en MongoDB

## Índices compuestos (*compound indexes*)

- Los índices compuestos recopilan y ordenan datos de dos o más campos de cada documento de una colección
- Los datos se agrupan por el primer campo del índice y, a continuación, por cada campo subsiguiente
- Para su creación se sigue un formato parecido al de los *single field index*

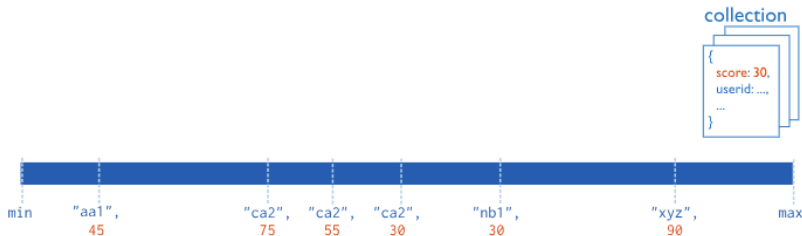
```
1 db.<collection>.create_index( {  
    <campo1>: <tipoDeOrdenacion>,  
3    <campo2>: <tipoDeOrdenacion>,  
    ...  
5    <campoN>: <tipoDeOrdenacion>  
} )
```

# Tipos de índices en MongoDB

## Índices compuestos (*compound indexes*) (ii)

- Ejemplo: si quisieramos crear un índice sobre los campos `userid` (de forma ascendente) y `score` (de forma descendente) de los documentos de una colección `usuarios` deberíamos ejecutar

```
db.usuarios.create_index({"userid": 1, "score": -1})
```



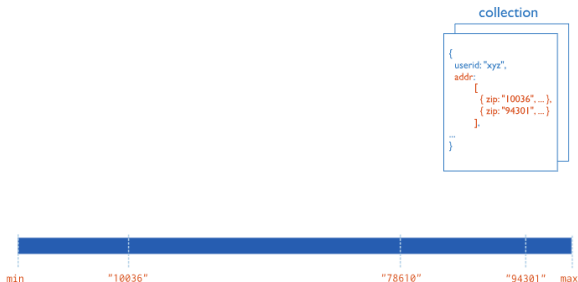


# Tipos de índices en MongoDB

## Índices multiclave (*multikey indexes*)

- Los índices multiclave recopilan y ordenan datos de campos que contienen un array como valor asociado
- Su creación es igual a los índices anteriores pero MongoDB establece automáticamente que ese índice sea multiclave

```
1 db.empleados.create_index( { "addr.zip": 1 } )
```



# Tipos de índices en MongoDB

## Índices de texto (*text indexes*)

- Los índices de texto permiten realizar búsquedas sobre campos de tipo `string`
- Los índices de texto mejoran el rendimiento cuando se buscan palabras o frases específicas dentro del contenido de una cadena
- Una colección sólo puede tener un índice de texto, pero ese índice puede abarcar varios campos

```
1 db.<coleccion>.create_index( {  
    <campo1>: "text",  
3    <campo2>: "text",  
    ...  
5    <campoN>: "text"  
} )
```

# Tipos de índices en MongoDB

## Índices de texto (*text indexes*) – Ejemplo

- Vamos a suponer que creamos una colección `blog` con los siguientes documentos

```
db.blog.insert_many( [  
  2   {  
      _id: 1,  
      "contenido": "Esta mañana me he comado una taza de café.",  
      "sobre": "bebida",  
      "palabras_clave": [ "café" ]  
  6   },  
  8   {  
      _id: 2,  
      "contenido": "¿A quién le apetece un helado de chocolate de postre?",  
      about: "comida",  
      "palabras_clave": [ "encuesta" ]  
 12   },  
 14   {
```

# Tipos de índices en MongoDB (cont.)

## Índices de texto (*text indexes*) – Ejemplo

```
16   _id: 3,  
    "contenido": "Mis sabores favoritos son la fresa y el chocolate",  
    "sobre": "helado",  
18   "palabras_clave": [ "comida", "postre" ]  
    }  
20 ] )
```

- Puesto que vamos a querer realizar muchas búsquedas sobre el campo contenido, vamos a definir un índice sobre dicho campo

```
db.blog.create_index( { "content": "text" } )
```

# Tipos de índices en MongoDB (cont.)

## Índices de texto (*text indexes*) – Ejemplo

- Esto permite realizar, por ejemplo, una búsqueda para recuperar los documentos que contengan la palabra *café* en el campo contenido

```
1 db.blog.find(  
  {  
3    "$text": { "$search": "café" }  
  }  
5 )
```

# Tipos de índices en MongoDB

## Indices hash (*hash indexes*)

- Recogen y almacenan los *hash* de los valores del campo indexado
- El campo que usemos para crear una clave hash de fragmento debe tener una cardinalidad alta, es decir, un gran número de valores diferentes
- Permiten la fragmentación en diferentes clusters (*sharding*) mediante claves hash

```
1 db.<coleccion>.create_index(  
  {  
3     "<campo>": "hashed"  
  }  
5 )
```

# La importancia de los índices en la eficiencia

- Los índices son cruciales para la eficiencia
- Y hay que conocer cuáles usar y si realmente una búsqueda los usará
- Recordemos la agregación sencilla de antes:

```
1 posts.aggregate( [  
    {'$match': { 'Score' : {'$gte': 40}}},  
3    {'$lookup': {  
        'from': "users",  
5        'localField': "OwnerUserId",  
        'foreignField': "Id",  
7        'as': "owner"}  
    }  
9 ])
```

- Nótese cómo se busca al usuario por el campo “Id” (ojo, no “\_id”)

# La importancia de los índices en la eficiencia (cont.)

- Tiempo de ejecución sin índice: **24 segundos**
- Defino un índice para el campo “Id” (valdría también *hash* porque se hace una búsqueda exacta):

```
1 db.users.create_index({"Id": 1})  
  > 'Id_1'
```

- Ejecuto de nuevo: **67,2 milisegundos**



# La importancia de los índices en la eficiencia (cont.)

- En pymongo hay que hacerlo de forma especial para saber que se ha usado el índice:

```
db.command('aggregate', 'posts', pipeline=pipeline, explain=True)
2 ...
'winningPlan': {'queryPlan': {'stage': 'EQ_LOOKUP',
4   'foreignCollection': 'stackoverflow.users',
   'localField': 'OwnerUserId',
6   'foreignField': 'Id',
   'asField': 'owner',
8   'strategy': 'IndexedLoopJoin',
   'indexName': 'Id_1',
10  'indexKeyPattern': {'Id': 1},
   'inputStage': {'stage': 'COLLSCAN',
12  'planNodeId': 1,
   'filter': {'Score': {'$gte': 40}},
14  'direction': 'forward'}}},
...
```

# Transacciones

- 1 Introducción a las bases de datos de documentos
- 2 Modelado de bases de datos de documentos
- 3 Introducción a MongoDB
  - Introducción
  - Uso básico
  - Consultas MapReduce
  - Validación
  - Índices y desnormalización
- 4 Uso de MongoDB desde pymongo
  - Métodos de búsqueda
  - Métodos de inserción y actualización
  - Framework de agregación
  - Índices
  - **Transacciones**

# Transacciones en MongoDB

- En MongoDB, una operación en un solo documento es atómica
- El uso de documentos incrustados permite capturar relaciones entre datos en una sola estructura de documento
  - Elimina la necesidad de transacciones distribuidas para muchos casos
- Para situaciones que requieren atomicidad de lecturas y escrituras en varios documentos, MongoDB soporta transacciones distribuidas
- Pueden ser utilizadas en múltiples operaciones, colecciones, bases de datos y documentos

# Transacciones en MongoDB

- En MongoDB ofrece una TransactionAPI que a su vez hace uso de la callback API
- Esta API hace la siguiente secuencia de operaciones:
  - 1 Comienza una transacción
  - 2 Ejecuta las operaciones indicadas
  - 3 Hace *commit* del resultado (o aborta en caso de error)

# Transacciones en MongoDB – Ejemplo

- 1. Nos conectamos a la BD

```
1 client = MongoClient(MONGODB_URI)
```

- 2. Definimos el *callback* que especifica la secuencia de operaciones a realizar

```
1 def callback(sesion):  
    coleccion_1 = session.client.db.foo  
3    coleccion_2 = session.client.db.bar  
    # Importante:: Se debe de pasar el parámetro de sesión a los operadores.  
5    coleccion_1.insert_one({"abc": 1}, session=sesion)  
    coleccion_2.insert_one({"xyz": 999}, session=sesion)
```

- 3. Iniciamos la sesión de cliente y ejecutamos la transacción (el bloque `with` ejecuta el *commit*)

```
with client.start_session() as session:  
    session.with_transaction(callback)  
client.close()
```

- (si algún paso falla, se puede llamar siempre a `session.abort_transaction()`, por ejemplo en un `catch` de una excepción)