



DEPARTAMENTO  
DE COMPUTACION

Facultad de Ciencias Exactas y Naturales - UBA

## Reentrega - Trabajo Práctico #2

12 de Diciembre de 2016

Base de Datos

Integrante	LU	Correo electrónico
Pedro Rodriguez	197/12	pedro3110.jim@gmail.com
Lucas Tavoraro Ortiz	322/12	tavo92@gmail.com



**Facultad de Ciencias Exactas y Naturales**  
Universidad de Buenos Aires

Ciudad Universitaria - Pabellón I

Intendente Güiraldes 2160 - C1428EGA

Ciudad Autónoma de Buenos Aires - Argentina

Tel/Fax: (54 11) 4576-3359

<http://exactas.uba.ar>

# Índice

<b>1. Correcciones Realizadas</b>	<b>3</b>
<b>2. Introducción</b>	<b>4</b>
<b>3. Documentos para Mongo</b>	<b>4</b>
3.1. Documento 1 y consultas asociadas . . . . .	4
3.1.1. Insert . . . . .	4
3.1.2. Mantenimiento . . . . .	4
3.1.3. Map-Reduce para mayor número de crímenes cometidos por alguna persona . . . . .	5
3.1.4. Map-Reduce para el número de crímenes promedio cometidos por criminales . . . . .	5
3.2. Documento 2 y consultas asociadas . . . . .	5
3.2.1. Insert . . . . .	6
3.2.2. Mantenimiento . . . . .	6
3.2.3. Map-Reduce para las 10 ciudades con mayor cantidad de crímenes . . . . .	6
3.2.4. Map-Reduce para cantidad de crímenes por localidad y año . . . . .	7
3.3. Documento 3 y consultas asociadas . . . . .	7
3.3.1. Insert . . . . .	8
3.3.2. Mantenimiento . . . . .	8
3.3.3. Map-Reduce para personas que fueron testigos en la mayor cantidad de casos . . . . .	8
3.4. Documento 4 y consultas asociadas . . . . .	9
3.4.1. Insert . . . . .	10
3.4.2. Mantenimiento . . . . .	10
3.4.3. Map-Reduce para casos con mayor número de involucrados . . . . .	10
3.4.4. Map-Reduce para cantidad total de evidencia por caso . . . . .	11
<b>4. Conclusiones</b>	<b>11</b>

## 1. Correcciones Realizadas

Para esta reentrega, efectuamos varias correcciones:

- cambiamos la estructura del TP: lo dividimos en 4 partes, según los 4 documentos que decidimos crear para cumplir con las queries que el gobierno quiere responder rápidamente, y con las queries que se desean resolver con el map-reduce.
- en 'Ciudades con mayor número de crímenes' cambiamos **dni** por **idCaso**, para considerar el caso en que un mismo crimen es cometido por varias personas como un único crimen
- en 'Personas como testigos en casos', buscamos las personas que se han visto involucradas como testigos en el mayor número de casos en lugar de buscar la cantidad de personas que se han visto involucradas como testigos en el mayor número de casos.
- eliminamos el documento que habíamos pensado para resolver la query sobre 'la cantidad de crímenes de cada tipo que ocurrieron en los últimos 45 días'. Para responder esta query, agregamos más datos a otro de los documentos que habíamos creado.
- cambiamos el método de resolución de todas las queries con map-reduce, utilizando en la mayoría de los casos dos map-reduce en lugar de uno solo, simplificando así enormemente la complejidad de las queries.
- llenamos una base de datos de MongoDB con documentos de ejemplo para probar las queries map-reduce.
- agregamos los JSONSchema correspondientes a cada documento

## 2. Introducción

En el presente trabajo llevamos nuestro proyecto a MongoDB, analizaremos maneras de representarlo para distintas queries con sus ventajas y desventajas y escribiremos el código de cada map-reduce.

## 3. Documentos para Mongo

A continuación vamos a definir el documento para Mongo en cada una de las queries dadas y argumentar la elección de dicho documento contando, a alto nivel, como serían los inserts, mantenimiento y búsqueda con map-reduce.

### 3.1. Documento 1 y consultas asociadas

Elegimos diseñar el documento **doc1**. A continuación exponemos un ejemplo:

```
1      [
2          {
3              dniCulpable : 37206752,
4              idCaso : 2
5          },
6      ]
```

Guardaremos en cada documento el **dniCulpable** de la persona que fue encontrada culpable del crimen **idCaso**. Para poder crear este documento tomamos del DER original los siguientes atributos: de la entidad **Personas**, el atributo **dni**, y de la entidad **CasoCriminal**, el atributo **idCaso**.

Con este documento podemos responder las siguientes queries que el gobierno quiere responder frecuentemente:

- cuál es el número promedio de crímenes cometidos por personas que ya han sido encontradas culpables de algún crimen?
- cuál es el mayor número de crímenes cometido por alguna persona?

El JSONSchema correspondiente sería:

```
1  {
2      "type" : object,
3      "properties": {
4          "dniCulpable" : {"type" : integer},
5          "idCaso" : {"type" : integer}
6      }
7  }
```

#### 3.1.1. Insert

Cada vez que una persona es hallada culpable de un nuevo caso, creamos un nuevo documento con su dni y el caso en cuestión. Esto es una operación muy rápida pues no depende de otros valores.

#### 3.1.2. Mantenimiento

Las ventajas de esta manera de representar la información están en su mantenimiento. Cada documento ocupa poco espacio y, al no depender de nada más, es fácil de mantener y escalar.

En un principio analizamos la posibilidad de tener un arreglo de casos por cada dni, pero esto requería mantener una estructura que no sabemos cuanto puede crecer en el tiempo y puede generar problemas al alcanzar el límite de disco. Por eso decidimos utilizar algo más sencillo, pero que podamos garantizar fácilmente su mantenimiento y escalabilidad.

### 3.1.3. Map-Reduce para mayor número de crímenes cometidos por alguna persona

Para resolver la consulta empleamos dos map-reduce.

En el primero, map11 mapea a cada criminal, un array de longitud igual a la cantidad de crímenes que este cometió. Luego, el reduce11 realiza una agregación, calculando cuántos crímenes cometió cada criminal.

En el segundo, map12 mapea a una misma clave todos los pares *< criminal, cantidadCrímenes >*. Y el reduce 12, se ocupa de buscar cuál fue la máxima cantidad de crímenes cometida por algún criminal, y devolver dicho valor.

```
1 map11 = function () {
2   emit( this . dniCulpable , 1 );
3 };
4 var reduce11 = function(key,values) {
5   return Array.sum( values );
6 };
7 var map12 = function () {
8   emit( 1, this.value );
9 };
10 var reduce12 = function(key, values) {
11   max = values[0];
12   for (var i=0; i<values.length; i++) {
13     if(values[i] > max) {
14       max = values[i];
15     }
16   }
17   return max;
18 };
19 db . doc1 . mapReduce (map11, reduce11, {out: "temp1"} );
20 db . temp1 . mapReduce (map12, reduce12, {out: "ej1"});
```

### 3.1.4. Map-Reduce para el número de crímenes promedio cometidos por criminales

Para resolver esta segunda consulta, también empleamos dos map-reduce.

Para el primer map-reduce, utilizamos el mismo que utilizamos en el ejercicio anterior: usamos la función map11 y reduce11, para obtener para cada criminal, la cantidad de crímenes que cometió.

Para el segundo map-reduce, utilizamos la misma función de mapeo map12 del ejercicio anterior, para mapear a una misma clave todas las cantidades de crímenes cometidas por los distintos criminales. Luego, usamos la función reduce21 para calcular el promedio de estos valores.

```
1 var reduce21 = function(key, values) {
2   sum = 0;
3   cnt = 0;
4   for (var i=0; i<values.length; i++) {
5     cnt = cnt + 1;
6     sum = sum + values[i];
7   }
8   return sum / cnt;
9 };
10 db . doc1 . mapReduce (map11, reduce11, {out: "temp2"} );
11 db . temp2 . mapReduce (map12, reduce21, {out: "ej2"} );
```

## 3.2. Documento 2 y consultas asociadas

Elegimos el documento **doc2**. A continuación, un ejemplo:

```
1 [
2   {
```

```

3         idCiudad: 4,
4         idCrimen: 12423,
5         anio: 2001,
6         fecha: '2001-10-11',
7         idCategoriaCrimen: 2
8     }
9 ]

```

Guardaremos en cada documento, una referencia a la ciudad en **idCiudad**, un crimen como **idCrimen** y el **anio** en que dicho crimen fue cometido, y por otra parte la **fecha** en la que este fue cometido. Para crear este documento, tomamos del DER original los siguientes atributos: de la entidad **Ciudad**, el atributo **idCiudad**, de la entidad **CasoCriminal** los atributos **idCaso**, y **fechaOcurrio**. También, de la entidad **CategoríaCaso**, el atributo **idCategoría**.

Con este documento podemos responder las siguientes queries que el gobierno quiere responder frecuentemente:

- qué cantidad, de cada tipo de crimen ocurrieron en los últimos 45 días?
- cuales son las 10 ciudades con mayor número de crímenes?

El JSONSchema correspondiente sería:

```

1 {
2     "type" : object,
3     "properties" : {
4         "idCiudad" : {"type" : integer},
5         "idCrimen" : {"type" : integer},
6         "anio" : {"type" : integer},
7         "fecha" : {"type" : Date},
8         "idCategoriaCrimen" : {"type" : integer}
9     }
10 }

```

### 3.2.1. Insert

Cada vez que se reporte un nuevo crimen, creamos un documento nuevo con la información.

### 3.2.2. Mantenimiento

Es similar a la anterior query analizada.

### 3.2.3. Map-Reduce para las 10 ciudades con mayor cantidad de crímenes

Para resolver esta consulta, utilizamos la función `map31`, que devuelve para ciudad, un vector con un 1 por cada crimen cometido en dicha ciudad. Luego, la función `reduce31` se ocupa de calcular la suma total de este vector para cada ciudad. Para obtener sólo las 10 ciudades con mayor cantidad de crímenes, debemos ordenar el documento devuelto por este map-reduce, y quedarnos sólo con las primeras 10 entradas.

```

1 var map31 = function() {
2     emit(this . idCiudad, 1);
3 };
4 var reduce31 = function(key, values) {
5     return Array.sum(values);
6 };
7 db . doc2 . mapReduce(map31, reduce31, {out:"temp3"} );
8 db.temp3.find().sort({value: -1}).limit(10);

```

### 3.2.4. Map-Reduce para cantidad de crímenes por localidad y año

En este caso, resolvemos la consulta con un único map-reduce. En el `map71`, lo que hacemos es recorrer todo el documento, y mapear cada crimen según la clave doble `< ciudad, anio >`. Esta función de mapeo va a devolver un arreglo con un 1 por cada crimen cometido en una tupla `< ciudad, anio >`. Luego, el `reduce71` lo único que debe hacer es sumar los valores de este arreglo, para obtener la cantidad de crímenes obtenidos en cada par de localidad y año posibles.

```
1 var map71 = function() {
2   emit(
3     {
4       ciudad: this . idCiudad ,
5       anio: this . anio
6     }, 1 );
7 };
8 var reduce71 = function(key, values) {
9   return Array.sum(values);
10 };
11 db . doc2 . mapReduce(map71, reduce71, {out:"ej7"} );
```

### 3.3. Documento 3 y consultas asociadas

Construimos el documento **doc3**.

Para este punto, decidimos probar con otro modelo de documento, que tiene otras ventajas y desventajas respecto de los utilizados anteriormente, que marcaremos más adelante.

Una instancia de ejemplo para este documento podría ser:

```
1 [
2   {
3     dniTestigo: 27893435,
4     idCasosInvolucradaComoTestigo:
5       [
6         {idCaso: 1},
7         {idCaso: 4},
8         {idCaso: 6}
9       ]
10  }
11 ]
```

Guardaremos en cada documento, una referencia al testigo con **dniTestigo**, y para cada uno de estos, un array **idCasosInvolucradaComoTestigo**, que tiene en cada entrada un **idCaso** con los casos en los que el testigo estuvo involucrado. Para crear este documento, tomamos del DER original los siguientes atributos: de la entidad **Personas**, el atributo **dni**, y de la entidad **CasoCriminal**, el atributo **idCaso**.

Con este documento podemos responder las siguientes queries que el gobierno quiere responder frecuentemente:

- cuales son las personas que se han visto involucradas como testigos en el mayor número de casos?

El JSONSchema correspondiente sería:

```
1 {
2   "type" : object,
3   "properties": {
4     "dniTestigo" : {"type" : integer},
5     "idCasosInvolucradaComoTestigo" : {
6       "type" : array,
7       "casos": {
8         "type" : object,
9         "properties": {
```

```

10         "idCaso" : {"type" : integer}
11     }
12 }
13 }
14 }
15 }

```

### 3.3.1. Insert

El Insert es mas costoso que en los anteriores pues debemos buscar el documento que contenga el dni del testigo del nuevo caso. Notemos que en caso de no existir un documento debemos crear uno nuevo. Una vez que lo tenemos, debemos agregar un nuevo caso al arreglo.

### 3.3.2. Mantenimiento

Mencionamos este caso anteriormente. La desventaja que tiene es que va a ir creciendo de tamaño cada documento y eventualmente puede llenar el disco. La solución que se nos ocurre es que, en caso de pasar eso, creamos un nuevo documento para ese testigo en donde a partir de ese momento iremos agregando los casos. En el map-reduce asumiremos que tenemos este diseño ya que es más interesante.

### 3.3.3. Map-Reduce para personas que fueron testigos en la mayor cantidad de casos

Para resolver esta consulta, empleamos dos map-reduce.

En el primero, usamos map41 para agrupar en una misma clave todos los pares  $\langle \text{testigo}, \text{cantCasosTestigoInvolucrado} \rangle$ . Luego, reduce41 calcula la máxima cantidad de casos en los que algún testigo estuvo involucrado, y devuelve un objeto con esa cantidad, *valMax* y con la lista entera de pares  $\langle \text{testigo}, \text{cantCasosTestigoInvolucrado} \rangle$ .

En el segundo mapeo, con map42 emitimos para una misma clave, de la lista de pares, solamente los testigos para los cuales la cantidad de casos en los que están involucrados es igual a la cantidad máxima registrada (*valMax*). Finalmente, el reduce41, se ocupa de devolver todos los testigos filtrados por map42.

```

1  var map41 = function() {
2      emit (1, {
3          testigo: this . dniTestigo,
4          valor: this . idCasosInvolucradoComoTestigo . length
5      });
6  };
7  var reduce41 = function(key, values) {
8      id_max = values[0].testigo;
9      val_max = values[0].valor;
10     for(var i=1; i<values.length; i++) {
11         if(val_max < values[i].valor) {
12             id_max = values[i].testigo;
13             val_max = values[i].valor;
14         }
15     }
16     return {
17         valMax: val_max,
18         valor: values
19     };
20 };
21 var map42 = function() {
22     for(var i=0; i<this . value . valor . length; i++) {
23         if(this . value . valor[i] . valor == this . value . valMax) {
24             emit(1,
25                 this . value . valor[i] . testigo
26             );
27         }

```



```

28     }
29 };
30 var reduce42 = function(key, values) {
31     const result = {
32         list: []
33     };
34     values.forEach((value, index) => result.list = result.list.concat(
35         value));
36     return result;
37 };
38 db . doc3 . mapReduce(map41, reduce41, {out:"temp4"} );
39 db . temp4 . mapReduce(map42, reduce42, {out:"ej4"} );

```

### 3.4. Documento 4 y consultas asociadas

Finalmente, pensamos otro diseño más, para resolver las consultas restantes. Llamamos a este **doc4**. A continuación exponemos una instancia de ejemplo.

```

1  [
2      {
3          idCaso: 4,
4          {
5              idEvidencias:
6                  [
7                      {idEvidencia: 3},
8                      {idEvidencia: 4}
9                  ],
10             dniInvolucrados:
11                 [
12                     {dniInvolucrado: 37206752}
13                 ]
14             }
15         }
16     ]

```

Guardamos en cada documento un **idCaso**, y para cada uno de estos tendremos un arreglo con las distintas evidencias en **idEvidencias** y otro con los involucrados en **dniInvolucrados**. Para poder generar este documento, tomamos del DER original los siguientes atributos: de la entidad **CasoCriminal**, el atributo **idCaso**, de la entidad **Evidencia**, el atributo **idEvidencia**, y de la entidad **Personas**, el atributo **dni**.

Con este documento podemos responder las siguientes queries que el gobierno quiere responder frecuentemente:

- cuáles son los casos en los que se han visto involucradas (cualquier rol) el mayor número de personas?
- cuáles son los casos con mayor número de evidencias?

El JSONSchema correspondiente sería:

```

1  {
2      "type" : object,
3      "properties": {
4          "idCaso" : {"type" : integer},
5          "idEvidencias" : {
6              "type" : array,
7              "evidencias": {
8                  "type" : object,
9                  "properties": {
10                      "idEvidencia" : {"type" : integer}

```

```

11         }
12     }
13 },
14     "dniInvolucrados" : {
15         "type" : array,
16         "involucrados": {
17             "type" : object,
18             "properties": {
19                 "dniInvolucrado" : { "type" :
20                                     integer}
21             }
22         }
23     }
24 }

```

### 3.4.1. Insert

A lo largo del tiempo, irán apareciendo evidencias nuevas y gente involucrada. Como en la anterior querie, debemos actualizar el documento agregando dicha información según qué caso se trate.

### 3.4.2. Mantenimiento

A diferencia del anterior, en este caso podemos pensar que es realmente muy difícil que este documento ocupe demasiado espacio en disco, pues el orden de magnitud de las evidencias y los involucrados no parece ser tanto. En caso de llevar acabo un sistema así, siempre se puede consultar un limite para poder tomar una decisión de diseño de estas características. Además de esta diferencia, en caso de llenarse el espacio reservado para que crezca el documento en disco, podemos utilizar una técnica similar a la ya expuesta en la querie anterior.

### 3.4.3. Map-Reduce para casos con mayor número de involucrados

Para resolver esta consulta, empleamos dos map-reduce.

En el primero, la función map51 mapea sobre una misma clave, todos los pares  $\langle idCaso, cantidadInvolucrados \rangle$  que se pueden obtener del documento. Luego, reduce51 se ocupa de detectar cuál es la máxima cantidad de involucrados, y devuelve un objeto con esa máxima cantidad *cantMax* junto con el vector con todos los pares extraídos del documento.

Luego, en map52 se agrupa en una misma clave todos los casos cuya cantidad de involucrados es igual a la máxima registrada. Por último, el reduce52 se encarga de simplemente devolver un arreglo con todos los valores resultantes del filtro aplicado por map52.

```

1 var map51 = function() {
2     emit(1, {
3         idCaso: this . idCaso,
4         cantidad: this . dniInvolucrados . length
5     });
6 };
7 var reduce51 = function(key, values) {
8     id_max = values[0].idCaso;
9     val_max = values[0].cantidad;
10    for(var i=1; i<values.length; i++) {
11        if(val_max < values[i].cantidad) {
12            id_max = values[i].idCaso;
13            val_max = values[i].cantidad;
14        }
15    }
16    return {

```

```

17     cantMax: val_max,
18     valor: values
19   };
20 };
21 var map52 = function() {
22   for(var i=0; i<this . value . valor . length; i++) {
23     if(this . value . valor[i] . cantidad == this . value . cantMax) {
24       emit(1, this . value . valor[i] . idCaso);
25     }
26   }
27 };
28 var reduce52 = function(key, values) {
29   const result = {
30     list: []
31   };
32   values.forEach((value, index) => result.list = result.list.concat(
33     value));
34   return result;
35 };
36 db . doc4 . mapReduce(map51, reduce51, {out:"temp5"} );
37 db . temp5 . mapReduce(map52, reduce52, {out:"ej5"} );

```

#### 3.4.4. Map-Reduce para cantidad total de evidencia por caso

Para resolver esta consulta, empleamos un único map-reduce. Map61 emite para cada caso, la cantidad de evidencias relacionadas con este. El reduce, simplemente devuelve para cada clave, el valor emitido por map61.

```

1 var map61 = function() {
2   emit( this . idCaso, this . idEvidencias . length );
3 };
4 var reduce61 = function(key, values) {
5   const result = {
6     list: []
7   };
8   values.forEach((value, index) => result.list = result.list.concat(
9     value));
10  return result;
11 };
12 db . doc4 . mapReduce(map61, reduce61, {out:"ej6"} );

```

## 4. Conclusiones

En el trabajo realizamos distintos tipos de documentos en donde pudimos ver que cada uno tiene ciertas ventajas y desventajas. Nos pareció enriquecedor poder llevar adelante las ideas que tuvimos para poder realizar un desarrollo de como funciona cada una. A modo de cierre, podemos concluir que el diseño del documento es algo que necesita pensarse bien ya que define mucho de todo lo demás. En particular es importante poder saber, en casos reales, como se va a comportar el documento, es decir, cada cuanto tenemos inserts, que recursos tenemos, si podemos realizar operaciones de mantenimiento frecuentemente, que queries y con que frecuencia recibiremos, etc. Con esta información clara, se puede tomar una buena decisión en cuanto al diseño para poder elegir la que tenga las ventajas que necesita nuestro problema.