

**Práctica Base de Datos no Convencionales  
(Máster Data Science)**

**Beatriz Visitación, Natalia Alonso  
y Susana Albarrán Félix**

# Contenidos

1	Introducción .....	3
2	Captura, procesamiento de los datos y carga de los datos en MongoDB. ....	4
3	Análisis de los datos. ....	5
3.1	Consulta 1.....	5
3.2	Consulta 2.....	6
3.3	Consulta 3.....	6
3.4	Consulta 4.....	7
3.5	Consulta 5.....	7
3.6	Consulta 6.....	10
3.7	Consulta 7.....	10
3.8	Consulta 8.....	11
3.9	Consulta 9.....	12
3.10	Consulta 10.....	13

# 1 Introducción

GESTIÓN DEL DOCUMENTO			
<b>Autor</b>	Beatriz Visitación, Natalia Alonso y Susana Albarrán	<b>Fecha</b>	04/05/2020
<b>Nombre fichero</b>	MDS_Memoria_Visitacion_Alonso_Y_Albarra_MongoDB.pdf	<b>Versión</b>	1.0
<b>Descripción:</b>  En esta memoria, se incluyen todos los pasos que se han seguido para realizar la práctica de Bases de Datos no Convencionales con MongoDB.  Para la realización de la práctica se ha usado el Lenguaje de programación Python (3.7) y MongoDB. Las consultas se han realizado en Robo 3T.			

En esta práctica se partirá de una fuente de datos real, exactamente de la base de datos **DBLP Computer Science Bibliography**, considerada como la mayor recopilación existente de referencias bibliográficas académicas centrada en la informática (en cualquiera de sus variantes). En particular, almacena los datos relativos a revistas científicas y congresos académicos sobre informática, en muchos casos remontándose hasta publicaciones de los años 60. Se caracteriza por un elevado grado de interdependencia interna. Además, el tamaño de esta fuente (en la forma de un único fichero) hará difícil su manejo, mostrando (aun a pequeña escala) algunas de las dificultades habituales en el procesamiento de esta información.

Consta de tres partes fundamentales: captura, procesamiento de datos y almacenamiento, y análisis de los datos.

El objetivo de este trabajo es utilizar esta información para realizar consultas en MongoDB y así poder obtener una serie de medidas que permitirán evaluar a los investigadores en informática.

## 2 Captura, procesamiento de los datos y carga de los datos en MongoDB.

La captura y procesamiento de datos consiste en descargar el fichero dbpl.xml <https://dblp.uni-trier.de/xml/> y a continuación procesarlo para convertirlo a formato JSON.

En nuestro caso se ha procesado el fichero dbpl.xml desde Python. Estos son los pasos que se han seguido:

- Convertir el fichero dbpl.xml a formato .json (Código fuente: **1.Convertir xml\_json.py**). Mediante este código se genera el fichero '**jsontdata.json**'.
- Se realiza un filtrado sobre el fichero '**jsontdata.json**' (Código fuente: **2.read\_json.py**) para sólo trabajar con los tres tipos de publicaciones: artículos de revista (article), artículos de congresos (inproceedings) y artículos de libros (incollecction). Mediante este código se generan tres ficheros: **incollecctions.json**, **article.json**, **inproceedings.json**.

Una vez convertidos al formato necesario los datos se almacenan en una base de datos MongoDB (**db.collection\_publication**). Para realizar este paso se usa el código fuente **3.crear\_db\_inserMany.py**.

Este último script para cargar los datos en MongoDB, lee los ficheros incollecctions.json, article.json e inproceedings.json, y para cada uno de ellos se recorre toda la lista de publicaciones extrayendo los campos **autor**, **title** y **year**. En algunos casos el campo autor es un array, ya que una publicación puede tener varios autores. Adicionalmente, se crea un nuevo campo **type** que identificará a los 3 tipos de publicaciones con los que se va a trabajar. También se ha observado que algunos de estos campos eran compuestos y tenían un subcampo "**#text**", en estos casos se ha seleccionado sólo el contenido de dicho campo. Estas son las instrucciones que se usan para la creación de la base de datos y la inserción de los distintos documentos:

```
conection = pymongo.MongoClient()
db = conection.dblp
collection = db.collection_publication
collection.insert_many(lista_incollecctions)
collection.insert_many(lista_inproceedings)
collection.insert_many(lista_articles)
conection.close()
```

La estructura de cada documento de la colección creada es la siguiente:

Key	Value	Type
▼ (1) ObjectId("5eb16563f8aa40adc0b578cb")	{ 5 fields }	Object
_id	ObjectId("5eb16563f8aa40adc0b578cb")	ObjectId
type	incollecction	String
> authors	[ 1 element ]	Array
title	Ear Shape for Biometric Identification.	String
year	2011	Int32

### 3 Análisis de los datos.

En esta etapa se realizan consultas sobre la base de datos almacenada en MongoDB (**db.collection\_publication**), usando el lenguaje de consulta propio de esta Base de Datos.

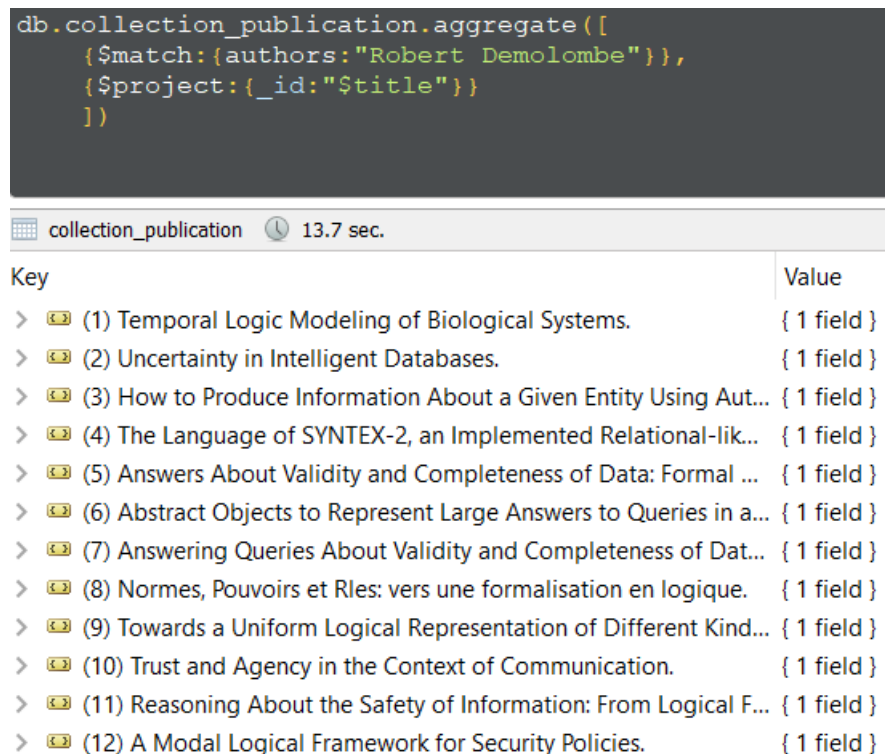
Antes de realizar las consultas se ha creado un índice para reducir el tiempo de ejecución de algunas de ellas. Se ha decidido indexar el campo **autor**, ya que es uno de los más usados (*db.collection\_publication.createIndex({authors:1})*). Adicionalmente, se probaron otros índices, pero se concluyó que el rendimiento de las consultas no mejoraba.

#### 3.1 Consulta 1.

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

En este caso se ha hecho un match por el autor: Robert Demolombe

```
db.collection_publication.aggregate([
  {$match:{authors:"Robert Demolombe"}},
  {$project:{_id:"$title"}}
])
```



```
db.collection_publication.aggregate([
  {$match:{authors:"Robert Demolombe"}},
  {$project:{_id:"$title"}}
])
```

Key	Value
> (1) Temporal Logic Modeling of Biological Systems.	{ 1 field }
> (2) Uncertainty in Intelligent Databases.	{ 1 field }
> (3) How to Produce Information About a Given Entity Using Aut...	{ 1 field }
> (4) The Language of SYNTAX-2, an Implemented Relational-lik...	{ 1 field }
> (5) Answers About Validity and Completeness of Data: Formal ...	{ 1 field }
> (6) Abstract Objects to Represent Large Answers to Queries in a...	{ 1 field }
> (7) Answering Queries About Validity and Completeness of Dat...	{ 1 field }
> (8) Normes, Pouvoirs et Rles: vers une formalisation en logique.	{ 1 field }
> (9) Towards a Uniform Logical Representation of Different Kind...	{ 1 field }
> (10) Trust and Agency in the Context of Communication.	{ 1 field }
> (11) Reasoning About the Safety of Information: From Logical F...	{ 1 field }
> (12) A Modal Logical Framework for Security Policies.	{ 1 field }

### 3.2 Consulta 2.

#### Número de publicaciones de un autor determinado.

Se usa el mismo autor que en la consulta 1: Robert Demolombe

```
db.collection_publication.aggregate([
  {$match:{authors:"Robert Demolombe"}},
  {$count:"Numero_publicaciones"},
  {$project:{_id:"$Numero_publicaciones"}}
])
```

```
db.collection_publication.aggregate([
  {$match:{authors:"Robert Demolombe"}},
  {$count:"Numero_publicaciones"},
  {$project:{_id:"$Numero_publicaciones"}}
])
```

collection\_publication 8.55 sec.

Key	Value
> (1) 90	{ 1 field }

### 3.3 Consulta 3.

#### Número de artículos en revista para el año 2018.

```
db.collection_publication.aggregate([
  {$match:{type:"article"}},
  {$match:{year:2018}},
  {$count:"Numero_articles_2018"},
  {$project:{_id:"$Numero_articles_2018"}}
])
```

```
db.collection_publication.aggregate([
  {$match:{type:"article"}},
  {$match:{year:2018}},
  {$count:"Numero_articles_2018"},
  {$project:{_id:"$Numero_articles_2018"}}
])
```

collection\_publication 6.89 sec.

Key	Value
> (1) 179219	{ 1 field }

### 3.4 Consulta 4.

**Número de autores ocasionales, es decir, que tengan menos de 5 publicaciones en total.**

```
db.collection_publication.aggregate([
  {$unwind:"$authors"},
  {$sortByCount:"$authors"},
  {$match:{count: {$lt: 5}}},
  {$count:"Número_autores_con_menos_de_5_publicaciones"},
  {$project:{_id:"$Número_autores_con_menos_de_5_publicaciones"}}
],{allowDiskUse:true})
```

```
db.collection_publication.aggregate([
  {$unwind:"$authors"},
  {$sortByCount:"$authors"},
  {$match:{count: {$lt: 5}}},
  {$count:"Número_autores_con_menos_de_5_publicaciones"},
  {$project:{_id:"$Número_autores_con_menos_de_5_publicaciones"}}
],{allowDiskUse:true})
```

collection\_publication 79.8 sec.

Key	Value
> (1) 1946102	{ 1 field }

### 3.5 Consulta 5.

**Número de artículos de revista (article) y número de artículos en congresos (inproceedings) de los diez autores con más publicaciones totales.**

- Se crea una colección nueva **"Top10\_autores"**:

```
db.collection_publication.aggregate([
  {$unwind:"$authors"},
  {$sortByCount:"$authors"},
  {$limit:10},
  {$out:"Top10_autores"}
],{allowDiskUse:true})
```

- Se crea otra colección ("**Autores\_articles**") donde se obtiene el número de autores para artículos de revista:

```
db.collection_publication.aggregate([
  {$match:{type:"article"}},
  {$unwind:"$authors"},
  {$sortByCount:"$authors"},
  {$out:"Autores_articles"}
],{allowDiskUse:true})
```

- Se unen las dos tablas anteriores:

```
db.Top10_autores.aggregate([
  {$lookup:{
    from : "Autores_articles",
    localField : "_id",
    foreignField : "_id",
    as : "Número_artículos"}},
  {$unwind:"$Número_artículos"}
],{allowDiskUse:true})
```

```
db.Top10_autores.aggregate([
  {$lookup:{
    from : "Autores_articles",
    localField : "_id",
    foreignField : "_id",
    as : "Número_artículos"}},
  {$unwind:"$Número_artículos"}
],{allowDiskUse:true})
```

Top10\_autores 0.006 sec.

Key	Value
> (1) H. Vincent Poor	{ 3 fields }
> (2) Mohamed-Slim Alouini	{ 3 fields }
> (3) Philip S. Yu	{ 3 fields }
> (4) Wei Zhang	{ 3 fields }
> (5) Wei Wang	{ 3 fields }
> (6) Lajos Hanzo	{ 3 fields }
> (7) Wei Li	{ 3 fields }
> (8) Yu Zhang	{ 3 fields }
> (9) Wen Gao 0001	{ 3 fields }
> (10) Victor C. M. Leung	{ 3 fields }



**/\* y número de artículos en congresos (inproceedings) de los diez autores con más publicaciones totales. \*/**

- Se crea otra colección ("**Autores\_inproceeding**") donde se obtiene el número de autores para artículos de congresos:

```
db.collection_publication.aggregate([
  {$match:{type:"inproceeding"}},
  {$unwind:"$authors"},
  {$sortByCount:"$authors"},
  {$out:"Autores_inproceeding"}
],{allowDiskUse:true})
```

- Se unen las dos tablas anteriores:

```
db.Top10_autores.aggregate([
  {$lookup:{
    from : "Autores_inproceeding",
    localField : "_id",
    foreignField : "_id",
    as : "Número_inproceeding"}},
  {$unwind:"$Número_inproceeding"}
],{allowDiskUse:true})
```

```
db.Top10_autores.aggregate([
  {$lookup:{
    from : "Autores_inproceeding",
    localField : "_id",
    foreignField : "_id",
    as : "Número_inproceeding"}},
  {$unwind:"$Número_inproceeding"}
],{allowDiskUse:true})
```

Top10\_autores 0.002 sec.

Key	Value
> (1) H. Vincent Poor	{ 3 fields }
> (2) Mohamed-Slim Alouini	{ 3 fields }
> (3) Philip S. Yu	{ 3 fields }
> (4) Wei Zhang	{ 3 fields }
> (5) Wei Wang	{ 3 fields }
> (6) Lajos Hanzo	{ 3 fields }
> (7) Wei Li	{ 3 fields }
> (8) Yu Zhang	{ 3 fields }
> (9) Wen Gao 0001	{ 3 fields }
> (10) Victor C. M. Leung	{ 3 fields }

### 3.6 Consulta 6.

**Número medio de autores de todas las publicaciones que tenga en su conjunto de datos.**

```
db.collection_publication.aggregate([
  {$unwind:"$authors"},
  {$sortByCount:"$_id"},
  {$group:{
    _id: null,
    Media: {$avg: "$count"}}},
  {$project: {_id: "$Media"}}
],{allowDiskUse:true})
```

```
db.collection_publication.aggregate([
  {$unwind:"$authors"},
  {$sortByCount:"$_id"},
  {$group:{
    _id: null,
    Media: {$avg: "$count"}}},
  {$project: {_id: "$Media"}}
],{allowDiskUse:true})
```

collection\_publication 62.1 sec.

Key	Value
> (1) 3.07986722835064	{ 1 field }

### 3.7 Consulta 7.

**Listado de coautores de un autor (Se denomina coautor a cualquier persona que haya firmado una publicación).**

```
db.collection_publication.aggregate([
  {$match:{authors:"Robert Demolombe"}},
  {$project: {"Co-autores" : "$authors"}},
  {$unwind:"$Co-autores"},
  {$sortByCount:"$Co-autores"},
  {$match:{_id : {$ne : "Robert Demolombe"}}}
])
```

```
db.collection_publication.aggregate([
  {$match:{authors:"Robert Demolombe"}},
  {$project: {"Co-autores" : "$authors"}},
  {$unwind:"$Co-autores"},
  {$sortByCount:"$Co-autores"},
  {$match:{_id : {$ne : "Robert Demolombe"}}},
])
```

collection\_publication 3.15 sec.

Key	Value
> (1) Luis Farias del Cerro	{ 2 fields }
> (2) Andrew J. I. Jones	{ 2 fields }
> (3) Frdric Cuppens	{ 2 fields }
> (4) Vincent Louis	{ 2 fields }
> (5) Emiliano Lorini	{ 2 fields }
> (6) Maria del Pilar Pozos Parra	{ 2 fields }
> (7) Jos Carmo	{ 2 fields }
> (8) Naji Obeid	{ 2 fields }
> (9) Andreas Herzig	{ 2 fields }
> (10) Arantza Illarramendi	{ 2 fields }
> (11) Laurence Cholvy	{ 2 fields }
> (12) Jos Miguel Blanco	{ 2 fields }

### 3.8 Consulta 8.

**Edad de los 5 autores con un periodo de publicaciones más largo (Se considera la edad de un autor el número de años transcurridos desde la fecha de su primera publicación hasta la última registrada).**

```
db.collection_publication.aggregate([
  {$unwind:"$authors"},
  {$group: {
    _id: "$authors",
    Minimo: { $min: "$year" },
    Maximo: { $max: "$year" } }},
  {$addFields:{Edad: { $subtract: ["$Maximo", "$Minimo"] } }},
  {$sort: {"Edad" : -1}},
  {$limit: 5}
],{allowDiskUse:true})
```

```
db.collection_publication.aggregate([
  {$unwind:"$authors"},
  {$group: {
    _id: "$authors",
    Minimo: { $min: "$year" },
    Maximo: { $max: "$year" } }},
  {$addFields:{Edad: { $subtract: ["$Maximo", "$Minimo"] } }},
  {$sort: {"Edad" : -1}},
  {$limit: 5}
],{allowDiskUse:true})
```

collection_publication 88.7 sec.	
Key	Value
▼ (1) Alan M. Turing	{ 4 fields }
_id	Alan M. Turing
Minimo	1937
Maximo	2012
Edad	75
> (2) Rudolf Carnap	{ 4 fields }
> (3) Claude E. Shannon	{ 4 fields }
> (4) David Nelson	{ 4 fields }
> (5) Martin Davis	{ 4 fields }

### 3.9 Consulta 9.

**Número de autores novatos, es decir, que tengan una Edad menor de 5 años. (Se considera la edad de un autor el número de años transcurridos desde la fecha de su primera publicación hasta la última registrada).**

```
db.collection_publication.aggregate([
  {$unwind:"$authors"},
  {$group: {
    _id: "$authors",
    Minimo: { $min: "$year" },
    Maximo: { $max: "$year" } }},
  {$addFields:{ Edad: { $subtract: ["$Maximo", "$Minimo"] } }},
  {$match:{Edad: {$lt: 5}}},
  {$count : "Número_autores_novatos"},
  {$project : {_id: "$Número_autores_novatos"}}
],{allowDiskUse:true})
```

```
db.collection_publication.aggregate([
  {$unwind:"$authors"},
  {$group: {
    _id: "$authors",
    Minimo: { $min: "$year" },
    Maximo: { $max: "$year" } }},
  {$addFields:{ Edad: { $subtract: ["$Maximo", "$Minimo"] },
  {$match:{Edad: { $lt: 5 } }},
  {$count : "Número_autores_novatos"},
  {$project : { _id: "$Número_autores_novatos" } }
  ],{allowDiskUse:true}})
```

collection\_publication 90.1 sec.

Key	Value
> (1) 1858893	{ 1 field }

### 3.10 Consulta 10.

**Porcentaje de publicaciones en revistas con respecto al total de publicaciones.**

```
db.collection_publication.aggregate([
  {$facet:
    {
      "Numero_articulos": [{$match:{type:"article"}}, {$count:"Num_articles"}],
      "Numero_total": [ {$count : "Num_total"}]
    }
  },
  {$unwind:"$Numero_articulos"},
  {$unwind:"$Numero_total"},
  {$addFields:{ Porcentaje: { $divide: ["$Numero_articulos.Num_articles",
"$Numero_total.Num_total"] }}}
  ],{allowDiskUse:true})
```

```
db.collection_publication.aggregate([
  {$facet:
    {
      "Numero_articulos": [{$match:{type:"article"}}, {$count:"Num_articles"}],
      "Numero_total": [ {$count : "Num_total"}]
    }
  },
  {$unwind:"$Numero_articulos"},
  {$unwind:"$Numero_total"},
  {$addFields:{ Porcentaje: { $divide: ["$Numero_articulos.Num_articles", "$Numero_total.Num_total"] }}}
  ],{allowDiskUse:true})
```

collection\_publication 13.1 sec.

Key	Value	Type
<ul style="list-style-type: none"> <li>(1)           <ul style="list-style-type: none"> <li>Numero_articulos               <ul style="list-style-type: none"> <li>Num_articles</li> </ul> </li> <li>Numero_total               <ul style="list-style-type: none"> <li>Num_total</li> </ul> </li> <li>Porcentaje</li> </ul> </li> </ul>	{ 3 fields } { 1 field } 2228904 { 1 field } 4857061 0.45889973381022	Object Object Int32 Object Int32 Double