

# Practica Big Data con Mongo DB

Gisela Ruiz Renzi

## Introducción

En esta práctica se verá como R puede ser utilizado para conectarnos a una base de datos NoSQL. En particular Mongo DB. Para la realización de la práctica se utilizará Atlas un servicio en la nube gratuito para manejar bases de datos Mongo DB. MongoDB Atlas es fácil de configurar y tiene conjuntos de datos de muestra para ejemplos de R con Mongo DB. Puede cargar conjuntos de datos de muestra usando el “...” junto al botón de colecciones en la página de su clúster. No obstante, aunque se puede crear un clúster específico en Atlas, esta práctica parte de uno ya creado y que puede ser utilizado para la resolución de los ejercicios.

- Adicionalmente, se recomienda utilizar un cliente Mongo DB para conectarse a la base de datos e inspeccionar los datos contenidos. Esto será extemadamente útil para realizar las consultas. Puede considerarse Robo 3T (<https://robomongo.org/download>). Además, si se crea el propio cluster en Atlas, este tiene una interfaz amigable para inspeccionar los datos.
- Además de estas funciones, existe documentación de las colecciones y la información contenida en esta base de datos de ejemplo en <https://docs.atlas.mongodb.com/sample-data/sample-training/>
- Para la resolución de ejercicios puede consultar el Manual de Mongo DB que contiene ejemplos y explicación de la sintáxis de Mongo DB

## Paquetes R utilizados

El controlador R MongoDB preferido, mongolite, es rápido y tiene una sintaxis similar a la del shell MongoDB. Mongolite es la que se utilizará en los siguientes ejemplos. Los otros paquetes enumerados aquí no han estado tan activos en Github recientemente. Los paquetes más populares para conectar MongoDB y R son:

- **mongolite**: un controlador R MongoDB más reciente, mongolite puede realizar varias operaciones como indexación, canalizaciones de agregación, cifrado TLS y autenticación SASL, entre otras. Está basado en el paquete jsonlite para R y mongo-c-driver. Podemos instalar mongolite desde CRAN o desde RStudio (explicado en una sección posterior). RMongo: RMongo fue el primer controlador R MongoDB con una sencilla interfaz R MongoDB. Tiene una sintaxis como la del shell MongoDB. RMongo ha quedado obsoleto a partir de ahora.
- **rmongodb**: rmongodb tiene funciones para crear pipelines, manejar objetos BSON, etc. Su sintaxis es muy compleja en comparación con mongolite. Al igual que RMongo, rmongodb ha quedado obsoleto y no está disponible ni se mantiene en CRAN.

## Instalación Mongolite

Para poder usar el paquete mongolite debemos instalarlo previamente con el comando siguiente, además de importar la librería posteriormente

```
#install.packages("mongolite")  
library(mongolite)
```

## Conectarse a Mongo DB

Esta es la cadena de conexión a MongoDB en Atlas. Se podría sustituir por otro servidor o clúster en Atlas si se desea.

```
cadena_conexion = 'mongodbsrv://user01:user01@cluster0.mcblc3z.mongodb.net/test'
```

Opciones de no validación de certificados SSL. Sin esta opción puede haber error de conexión a Atlas.

```
opciones_conexion = ssl_options(weak_cert_validation = T)
```

Después de establecer la conexión a MongoDB se recupera la colección ‘trips’ usando la función mongo() en código R para obtener la colección de viajes de la base de datos sample\_training. Esta recopilación contiene datos de viajes realizados por los usuarios de un servicio de bicicletas compartidas con sede en la ciudad de Nueva York.

```
viajes = mongo(collection="trips", db="sample_training", url=cadena_conexion, options = opciones_conexion)
```

Puede verificar que su código ahora esté conectado a la colección MongoDB verificando el número total de documentos en esta base de datos. Para hacerlo, use la función count().

```
viajes$count()
```

```
## [1] 10000
```

Ahora que tiene una conexión establecida con la base de datos, podrá leer los datos de la misma para ser procesados por R.

## Cómo obtener datos en R desde MongoDB

En esta sección, aprenderá cómo recuperar datos de MongoDB y mostrar los mismos. Continuemos con trips\_collection de la sección anterior.

Puede usar la interfaz de usuario de MongoDB Atlas para ver los documentos de trip\_collection o RStudio para visualizarlos.

Obtenga cualquier documento de muestra de la colección usando el método *iterate().one()* para examinar la estructura de los datos de esta colección.

```
viajes$iterate()$one()
```

```
## $tripduration
## [1] 396
##
## $'start station id'
## [1] 317
##
## $'start station name'
## [1] "E 6 St & Avenue B"
##
## $'end station id'
## [1] 317
##
## $'end station name'
## [1] "E 6 St & Avenue B"
##
## $bikeid
## [1] 23082
##
## $usertype
## [1] "Subscriber"
##
## $'birth year'
## [1] 1981
##
## $'start station location'
## $'start station location'$type
## [1] "Point"
##
## $'start station location'$coordinates
## $'start station location'$coordinates[[1]]
## [1] -73.98185
##
## $'start station location'$coordinates[[2]]
## [1] 40.72454
##
##
##
## $'end station location'
## $'end station location'$type
## [1] "Point"
##
## $'end station location'$coordinates
## $'end station location'$coordinates[[1]]
## [1] -73.98185
##
## $'end station location'$coordinates[[2]]
## [1] 40.72454
##
##
##
## $'start time'
```

```
## [1] "2016-01-01 01:17:55 CET"
##
## $'stop time'
## [1] "2016-01-01 01:24:32 CET"
```

Ahora que conoce la estructura de los documentos, puede realizar consultas más avanzadas, como buscar los cinco viajes más largos a partir de los datos de recopilación de viajes. Y luego enumerar la duración en orden descendente.

```
viajes$find(sort = '{"tripduration" : -1}' , limit = 5)
```

```
##      tripduration start station id      start station name end station id
## 1          326222          391      Clark St & Henry St          310
## 2          279620          3165 Central Park West & W 72 St          3019
## 3          173357          3155 Lexington Ave & E 63 St          3083
## 4          152023          340  Madison St & Clinton St          2009
## 5          146099          368   Carmine St & 6 Ave          238
##      end station name bikeid  usertype birth year
## 1   State St & Smith St 18591 Subscriber    1979
## 2   NYCBS Depot - DEL 17547  Customer
## 3  Bushwick Ave & Powers St 15881  Customer
## 4  Catherine St & Monroe St 22678 Subscriber    1992
## 5  Bank St & Washington St 15553  Customer
##      start station location.type start station location.coordinates
## 1                      Point      -73.99345, 40.69760
## 2                      Point      -73.97621, 40.77579
## 3                      Point      -73.96649, 40.76440
## 4                      Point      -73.98776, 40.71269
## 5                      Point      -74.00215, 40.73039
##      end station location.type end station location.coordinates
## 1                      Point      -73.98913, 40.68927
## 2                      Point      -73.98193, 40.71663
## 3                      Point      -73.94100, 40.71248
## 4                      Point      -73.99683, 40.71117
## 5                      Point      -74.00859, 40.73620
##      start time      stop time
## 1 2016-01-01 01:58:20 2016-01-04 20:35:23
## 2 2016-01-02 17:07:26 2016-01-05 22:47:46
## 3 2016-01-02 15:25:36 2016-01-04 15:34:53
## 4 2016-01-02 00:31:59 2016-01-03 18:45:42
## 5 2016-01-01 02:59:07 2016-01-02 19:34:07
```

La consulta anterior utiliza operadores de clasificación y límite para producir este conjunto de resultados.

## Cómo analizar datos de MongoDB en R

Para analizar MongoDB con R con más detalle, puede usar el marco de agregación de MongoDB. Este marco permite a los operadores crear canalizaciones de agregación que ayudan a obtener los datos exactos con una sola consulta.

Suponga que desea verificar cuántos suscriptores realizaron viajes de una duración > 240 segundos y regresaron a la misma estación donde comenzaron. La consulta usa MongoDB [\$expr] (<https://docs.mongodb.com/manual/reference/operator/query/expr/>) para comparar dos campos en el mismo documento.

```
query = viajes$find('{"usertype":"Subscriber","tripduration":{"$gt":240},"$expr": {"$eq": [{"start stat.
```

Combinando estos operadores con algún código R, también puede ver qué tipo de usuarios son más comunes: suscriptores o clientes únicos. Para ello, se puede agrupar usuarios por tipo de usuario campo.

```
tipos_usuario = viajes$aggregate(' [{"group":{"_id":"$usertype", "Count": {"$sum":1}}}]')
```

Para comparar los resultados, puede visualizar los datos. Es conveniente convertir los datos obtenidos de mongolite en un marco de datos y usar ggplot2 para trazar.

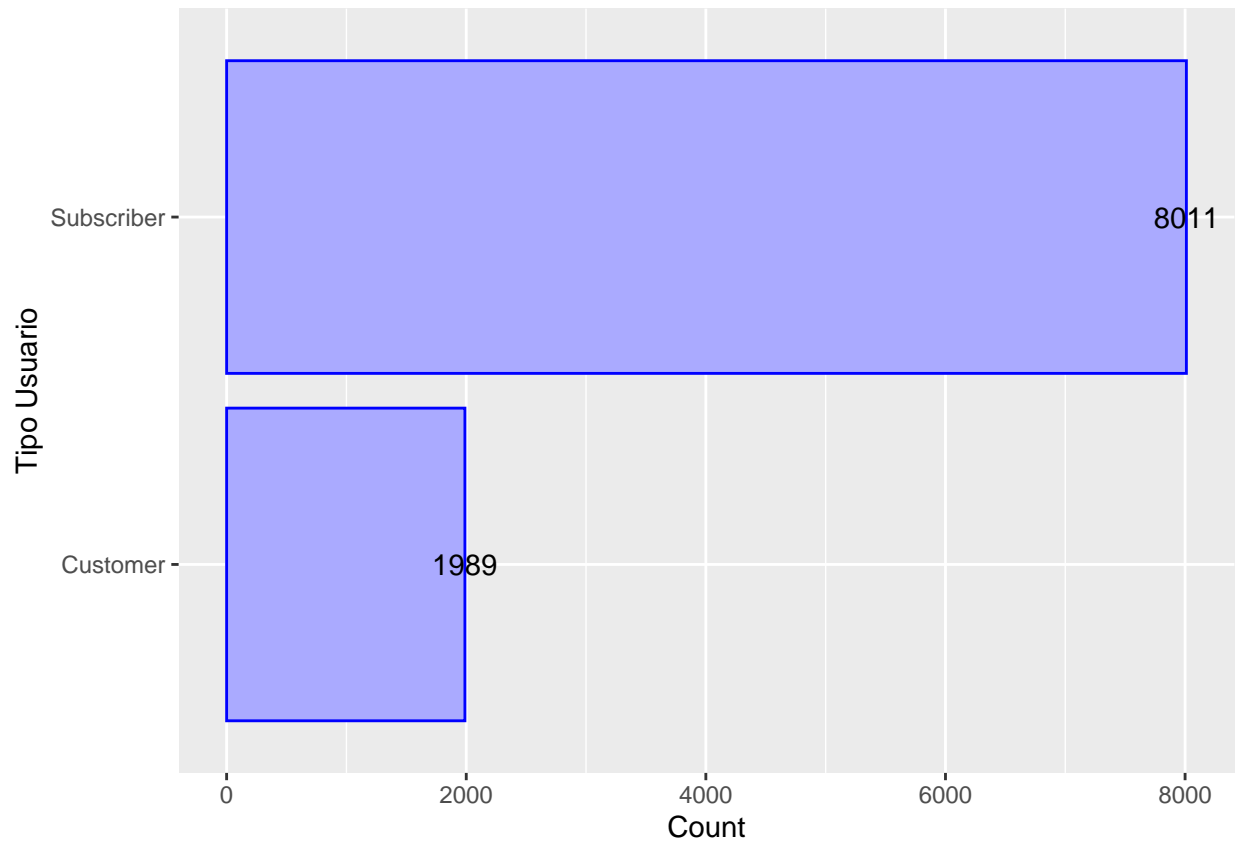
```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.3      v readr      2.1.4
## v forcats    1.0.0      v stringr    1.5.0
## v ggplot2    3.4.3      v tibble     3.2.1
## v lubridate  1.9.2      v tidyr      1.3.0
## v purrr      1.0.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```
library(lubridate)
library(ggplot2)
```

```
df <- as.data.frame(tipos_usuario)
```

```
ggplot(df,aes(x=reorder(`_id`,Count),y=Count))+
geom_bar(stat="identity",color='blue',fill='#aaaaff')+geom_text(aes(label = Count), color = "black") +c
```



## Tarea a realizar

A continuación se proponen 4 ejercicios a completar. Se recomienda completarlos en script markdown de R que se proporciona.

### Tarea 1

Exploremos otro diagrama de barras con una colección diferente - inspections. Esta recopilación contiene datos sobre las inspecciones de edificios de la ciudad de Nueva York y si pasan o no. Recupere dicha colección en R.

```
print("Completar Tarea 1")
```

```
## [1] "Completar Tarea 1"
```

```
inspecciones = mongo(collection="inspections", db="sample_training", url=cadena_conexion, options = opo
```

### Tarea 2

Suponga que desea verificar el número de empresas que no aprobaron las inspecciones en 2015 en comparación con 2016.

Si ve los datos obtenidos de la colección, notará que el campo de fecha es una Cadena. Convertirlo en tipo de fecha y luego extraer el año requerirá algún procesamiento. Pero, con la canalización de agregación de MongoDB, puede hacer todo en una sola consulta. Para manipular el campo de fecha, use el operador \$addField.

Además, agregue las deficiencias encontradas en las inspecciones por año

```
print("Completar Tarea 2")
```

```
## [1] "Completar Tarea 2"
```

```
empresas_Failed = inspecciones$aggregate ('[{"$addField": {
    "dateModificado":{"$year":{"
    "$dateFromString": {
    "dateString": "$date"}}
    }},
    { "$match": {
    "$and": [
    {"$or": [
    {"dateModificado": 2015},
    {"dateModificado": 2016}
    ]},
    {"result":"Fail"}
    ] }},
    {"$group":{"_id":"$dateModificado",
    "Count":{"$sum":1}}}
  ]')
```

### Tarea 3

Teniendo en cuenta que el resultado de la tarea anterior está agrupando los resultados por año, cree un gráfico de barras.

```
print("Completar Tarea 3")
```

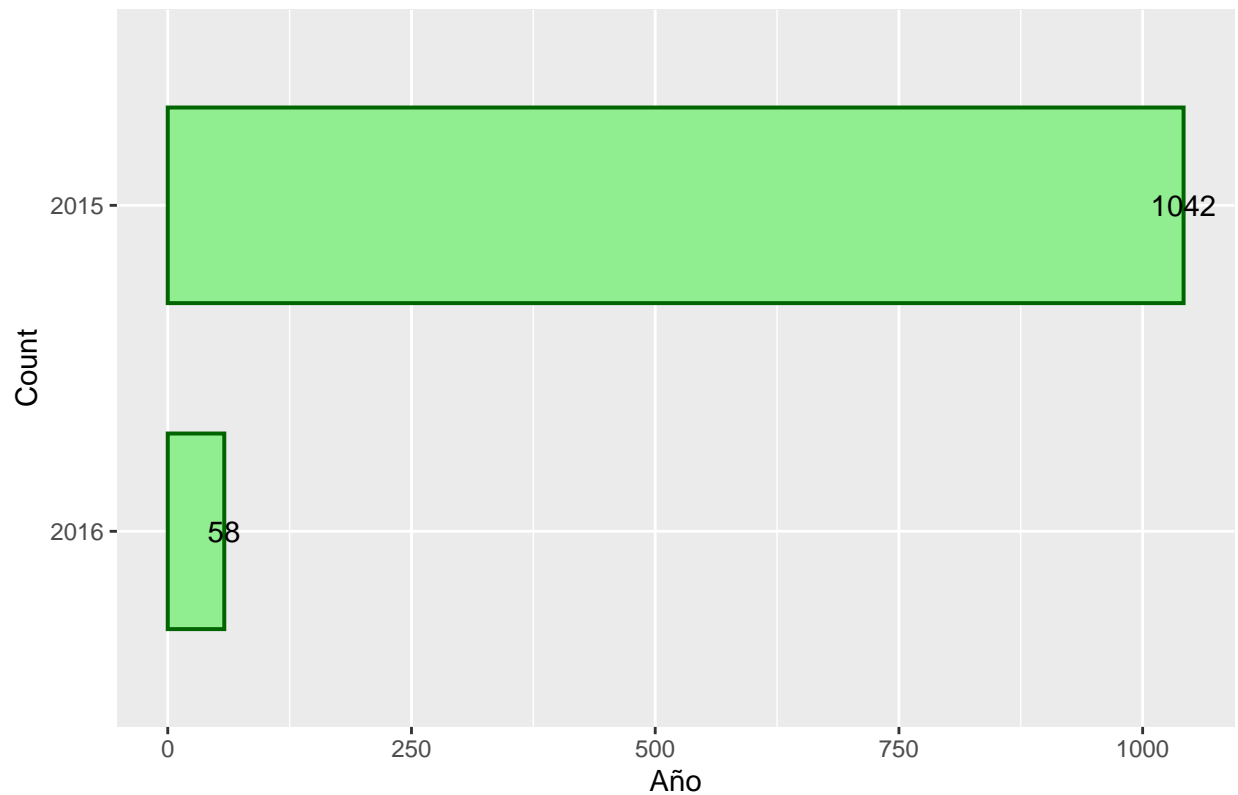
```
## [1] "Completar Tarea 3"
```

```
df_inspections_failed <- as.data.frame(empresas_Failed)
df_inspections_failed
```

```
##      _id Count
## 1 2016     58
## 2 2015    1042
```

```
ggplot(df_inspections_failed, aes(x = reorder(`_id`,Count), y = Count)) +
geom_bar(stat="identity",color='darkgreen',lwd = 0.7, fill='lightgreen', width =.6 ) +
coord_flip() +
geom_text(aes(label = Count), color = "black") +
labs(title="Empresas que no superan la inspección", y="Año", x="Count")
```

### Empresas que no superan la inspección



### Tarea 4

A continuación, se utilizará la colección 'companies', que contiene información sobre empresas, como su año de fundación y la dirección de su sede.

Supongamos que desea conocer la tendencia del número de empresas de publicidad (category\_code = 'advertising') fundadas a partir de 2000 en adelante. Para ello, utilice el operador relacional \$gt, agrupe los resultados por año de creación (founded\_year) y ordénelos para que se muestren posteriormente en un gráfico de líneas por año.

```
print("Completar Tarea 4")
```

```
## [1] "Completar Tarea 4"
```

```
empresas = mongo(collection="companies", db="sample_training", url=cadena_conexion, options = opciones_

tendencia_empr_publi=empresas$aggregate(['{"$match":
    {"category_code":"advertising",
    "founded_year":{"$gt":1999}}},
    {"$group":
    {"_id": "$founded_year", "companies_year": {"$sum": 1}}},
    {"$sort": {"_id": 1}}
    ]')
```



```
df_tendencia_empr_publi <- as.data.frame(tendencia_empr_publi)
df_tendencia_empr_publi
```

```
##      _id companies_year
## 1  2000             16
## 2  2001             13
## 3  2002             15
## 4  2003             17
## 5  2004             24
## 6  2005             32
## 7  2006             49
## 8  2007             80
## 9  2008             77
## 10 2009             27
## 11 2012              1
## 12 2013              1
```

```
ggplot(df_tendencia_empr_publi, aes(x = `_id`, y = companies_year)) +
  geom_line(color='orange', lwd = 1) +
  geom_point(size=2, color='orange') +
  geom_text(aes(label = `_id`), color = 'darkblue', size= 3.4, vjust= -0.7) +
  labs(title="Número de empresas de publicidad fundadas por año", y="Empresas fundadas", x="Año")
```

