



Análisis de Estadísticas de Videos en Tendencia de YouTube

Integrantes:

- Caparachin Villaverde, Nanto Gustavo (E202110812)
- García Godos Villavicencio, Jorge Daniel (E202110981)
- Huamani Franco, Ismael (E202110991)
- Poma Ludeña, Rodrigo Renato (E202111018)

INTRODUCCIÓN

Extracción e Ingesta de Información



Fuente:

Trending YouTube Video Statistics

Agente:

Importador de archivos hacia MONGO DB (ATLAS).

Base de Datos:

Servicio MongoDB alojado en la nube en **AWS**, Azure y Google Cloud.

INTRODUCCIÓN

Extracción, Análisis y Representación de los Datos



Fuente:
MongoDB con las colecciones de datos unificados

Aplicación de Análisis:
Entorno de desarrollo integrado con herramientas de análisis estadístico.

Herramienta de Visualización:
Aplicaciones interactivas de visualización de datos

METODOLOGÍA

Procedimiento realizado

- 1. Identificación de fuente de datos**
Buscamos fuentes de datos confiables y escogimos una fuente considerando la veracidad, volumen y variedad.
- 2. Dimensionamiento y creación de la base de datos**
Dimensionamos la base de datos considerando solo tres países y creamos la base de datos en su entorno PAAS creando usuarios de solo lectura y abriendo la comunicación a cualquier IP.
- 3. Obtención de datos e importación a base de datos**
Automatizamos la obtención de datos usando un script bash que descarga las fuentes de datos y los ingresa en la base de datos mongo con la librería de mongo-import.
- 4. Conexión a base de datos y extracción de data**
Utilizamos mongolite para establecer la conexión con la base de datos y extraer la información.
- 5. Comprensión y limpieza de datos**
Organizamos, categorizamos, tipificamos y unificamos las diferentes colecciones para tener un dataframe confiable en el que realizar nuestro análisis.
- 6. Análisis de outlines**
Normalizamos los datos, aplicamos diferentes técnicas como Shapiro, Anderson-Darling, Mahalanobis y Chi-Cuadrado
- 7. Exploración de datos en Shiny**
Diseñamos los componentes UI y completamos diferentes gráficos de hallazgos en la sección de server.

Extracción de los datos CSV & JSON

**Cadena de
Conexión a
MongoDB desde
RStudio**

**Extracción de
Colecciones**

**Tratar los datos
extraídos según el
tipo de fuente
asociado (CSV,
JSON).**

```
#1.1 Conexión a la base de datos
connection_string = 'mongodb+srv://read-only-user:U5PN6pQ1eirSu7e@mycluster.za173.mongodb.net/dbyoutube'

#1.2 Extrayendo la data de cada colección tipo CSV de trending
MX_trending = mongo(collection="MX_youtube_trending_data", db="dbyoutube", url=connection_string)
US_trending = mongo(collection="US_youtube_trending_data", db="dbyoutube", url=connection_string)
FR_trending = mongo(collection="FR_youtube_trending_data", db="dbyoutube", url=connection_string)

#1.3 Explorando las colecciones de la base de datos de YOUTUBE
MX_trending$run('{"listCollections": 1}')$cursor$firstBatch$name

#1.4 Extrayendo la data de cada colección tipo JSON de category
MX_category = mongo(collection="MX_category_id", db="dbyoutube", url=connection_string)
US_category = mongo(collection="US_category_id", db="dbyoutube", url=connection_string)
FR_category = mongo(collection="FR_category_id", db="dbyoutube", url=connection_string)

#1.5 Armando los archivos de trending que son tipo csv
MX <- MX_trending$find('{}')
US <- US_trending$find('{}')
FR <- FR_trending$find('{}')

#1.6 Armando los archivos de categoria que son tipo json
cat_MX <- MX_category$find('{}')
cat_US <- US_category$find('{}')
cat_FR <- FR_category$find('{}')
```

Exploración de los datos CSV

Reconocimiento de los datos en los 5 primeros registros

Identificación de una columna de tags con potencial de análisis

Hallazgo de datos propios de un video en una plataforma digital

```
> head(MX)
  video_id      publishedAt channelTitle categoryId trending_date
1 QgX0zn5nMyk 2020-08-11T15:57:47Z      Hoy          24 2020-08-12T00:00:00Z
2 zU5lgsAKIIA 2020-08-11T20:18:23Z      Mailu          22 2020-08-12T00:00:00Z
3 p7WtYeowOEY 2020-08-12T01:00:04Z SOY TRIKI TRIKI    22 2020-08-12T00:00:00Z
4 M9Pmf9AB4Mo 2020-08-11T17:00:10Z      Apex Legends    20 2020-08-12T00:00:00Z
5 eew7QJqFn7Q 2020-08-10T17:06:34Z      Rubius Z        20 2020-08-12T00:00:00Z
6 vXSv6umIWXU 2020-08-10T16:59:41Z      Genius          10 2020-08-12T00:00:00Z

tags
1 Televisa|Televisa espectaculos|programa hoy televisa|Pr
ograma Hoy|elenco programa hoy|videos programa hoy|programa hoy 2020|Andrea Legarreta|Galilea Montijo|Raúl Arañ
za|Andrea Escalona|#ConLasEstrellas|Paul Stanley|Elenco programa hoy 2020|Lambda García|Marisol González|#Telev
isaTeAcompaña|Magda Rodríguez|con permiso|unicable|pepillo origel|martha figueroa|erika buenfil|hijo de erika b
uenfil
2 lulu99|novio de lulu99|lulu y
su novio|maicol y luisa|tag del novio|novio|lulu 99|99|lulu99 y su novio canal|mailu|mai lul|pareja|novios|YOLO
aventuras|sandra cires|24 horas|esposados|con mi novio|maicol|maiking|NUEVO CANAL JUNTOS! Lulu Y Maicol <U+27
64><U+FE0F> NUESTRO PRIMER VIDEO!!|novio vs novia|MaiLu|reto|quien conoce mejor al otro|challenge|skabeche|123
spanish
3 Triki triki|Payaso triki trki|Soy triki
triki
4 Apex Legends|Apex Legends characters|new Apex Legend|Apex Legends Rampart|Apex Legends Season 6|Apex Legends
Boosted|Battle Pass|Season 6 Battle Pass|Apex Legends new season|Apex Legends game|Respawn Apex Legends|Battle
Royale game|Battle Royale|Battle Royale shooter|Apex Games|squad play|multiplayer shooter|Apex Legends PS4|Ape
x Legends Xbox|Apex Legends PC|Apex Legends Origin|Respawn Entertainment|Electronic Arts|Titanfall 2|fun battle
royale
5 elrubius|rubius|kun|kun aguero|agüero|fall guys|fall g
uys|full|battle royale|gameplay|mejor|mejores momentos|manchester|manchester city|futbol|2020|gol|goles|sergio|
kubius
6 genius|rap genius|verified|official lyrics|lyrics|lyric video|Lyric videos|pop music|hip hop|rap|ne
w pop music|jd pantoja vevo|jd pantoja hagamos las paces|jd pantoja canciones|jd pantoja letal|jd pantoja 12 19
|gbwc0
view_count likes dislikes comment_count comments_disabled ratings_disabled country
1 521751 4680 3129 2285 False False Mexico
2 989033 157522 2060 18314 False False Mexico
3 64858 6327 110 437 False False Mexico
4 2381688 146744 2794 16557 False False Mexico
5 4331471 492848 4011 10942 False False Mexico
6 403952 37242 9986 6734 False False Mexico
```

Limpieza de los datos CSV

Creamos una variable para identificar los países

```
#2.1.1 Preparamos las tablas de YOUTUBE creando un nuevo campo country: Mexico,  
MX$country = MX$country = 'Mexico'  
US$country = US$country = 'USA'  
FR$country = FR$country = 'Francia'
```

Eliminamos aquellas variables que no se considerados evaluar, tal como:
Link del video,
Descripción,
Titulo,
El código del canal.

```
#2.1.2. Eliminamos las columnas que no vamos a utilizar  
MX[, c("thumbnail_link", "description", "title", "channelId")] <- NULL #tags ar  
US[, c("thumbnail_link", "description", "title", "channelId")] <- NULL  
FR[, c("thumbnail_link", "description", "title", "channelId")] <- NULL  
#se elimina thumbnail_link porque es la direccion url no agrega valor al anali  
#se elimina description porque es la descripcion del video, contiene datos no e
```

Solo en la variable de canal (CHANNELID) se encontraron vacíos, además, no brinda un valor significativo al análisis, por lo que, no se utilizará en la evaluación

Procesamiento de los archivos JSON

Encontramos los datos de interés de los archivos JSON,

Creamos la función para extraer los datos de interés: código de categoría y descripción de la categoría

```
cat_FR$snippet$title[1] #
cbind(id=cat_FR$id[1], title=cat_FR$snippet$title[1]) #
> cat_FR$snippet$title[1]
[1] "Short Movies"
> cbind(id=cat_FR$id[1], title=cat_FR$snippet$title[1])
  id title
[1,] "18" "Short Movies"
```

```
#Inf de categoria de Mexico
vector1 = c()
vector2 = c()
for (i in 1:length(cat_MX$id)){
  vector1 = c(vector1, cat_MX$id[i])
  vector2 = c(vector2, cat_MX$snippet$title[i])
}
DFX = data.frame("id"=vector1, "title"=vector2)

#Inf de categoria de USA
vector3 = c()
vector4 = c()
for (i in 1:length(cat_US$id)){
  vector3 = c(vector3, cat_US$id[i])
  vector4 = c(vector4, cat_US$snippet$title[i])
}
DFUS = data.frame("id"=vector3, "title"=vector4)

#Inf de categoria de Francia
vector5 = c()
vector6 = c()
for (i in 1:length(cat_FR$id)){
  vector5 = c(vector1, cat_FR$id[i])
  vector6 = c(vector2, cat_FR$snippet$title[i])
}
DFFR = data.frame("id"=vector5, "title"=vector6)

#2.2.3. juntamos las 3 bases de datos
DFCAT = bind_rows(DFX,DFUS,DFFR)
```


METODOLOGÍA

Herramientas utilizadas

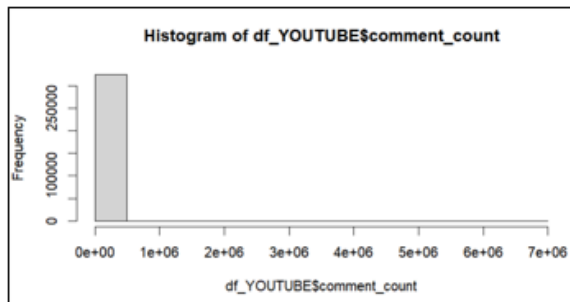
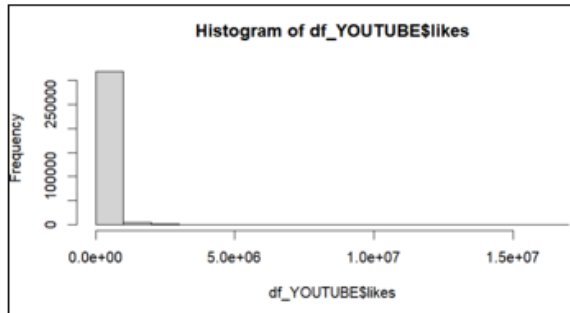
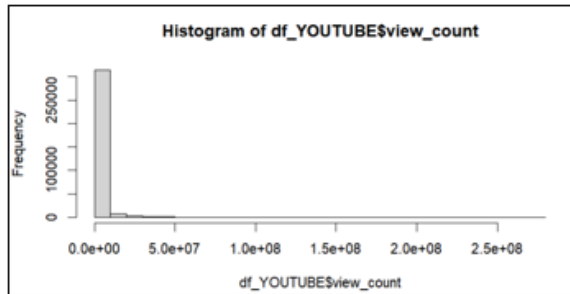
HERRAMIENTA	DETALLE	UTILIZACIÓN
Kaggle	Kaggle es un repositorio de datasets abiertos.	Obtención de nuestro dataset actualizado de tendencias de videos.
Amazon Cloud9	IDE de desarrollo de software ejecutado en EC2.	Manipulación del comprimido del dataset e importación a la base de datos.
Atlas MongoDB	Base de datos no relacional administrada por Atlas.	Base de datos de consulta de dataset.
DataBase Tools	Librería de manipulación de datos	Uso de la función mongo-import para la importación del dataset.
RStudio Cloud	IDE cloud de desarrollo de proyectos en lenguaje R.	Desarrollo de análisis de datos.
Amazon EC2	Máquina virtual multipropósito.	Ejecución de consultas con alto consumo de datos y procesamiento.

METODOLOGÍA

TÉCNICA	ETAPA DE UTILIZACIÓN
Eliminación de duplicados	Comprensión y limpieza de datos
Aplicación de Factor o Levels	Comprensión y limpieza de datos
Validación de Nulos	Comprensión y limpieza de datos ,Análisis de outliers
Validación de Ceros	Análisis de outlines
Normalización usando criterio Min-Max	Análisis de outlines
Shapiro	Análisis de outlines
Anderson-Darling	Análisis de outlines
Distancia de Mahalanobis	Análisis de outlines, Exploración de datos en Shiny
Chi cuadrado	Análisis de outlines, Exploración de datos en Shiny
Análisis de frecuencia	Exploración de datos en Shiny
Análisis de dispersión	Exploración de datos en Shiny

RESULTADOS

Análisis de dispersión



Resultado 1:

Inicialmente utilizamos 4 variables para el análisis: Likes, Views y Coments y dislikes hemos encontrado que la distribución de las variables de Likes, Views y Coments, tienen una distribución claramente asimétrica positiva.

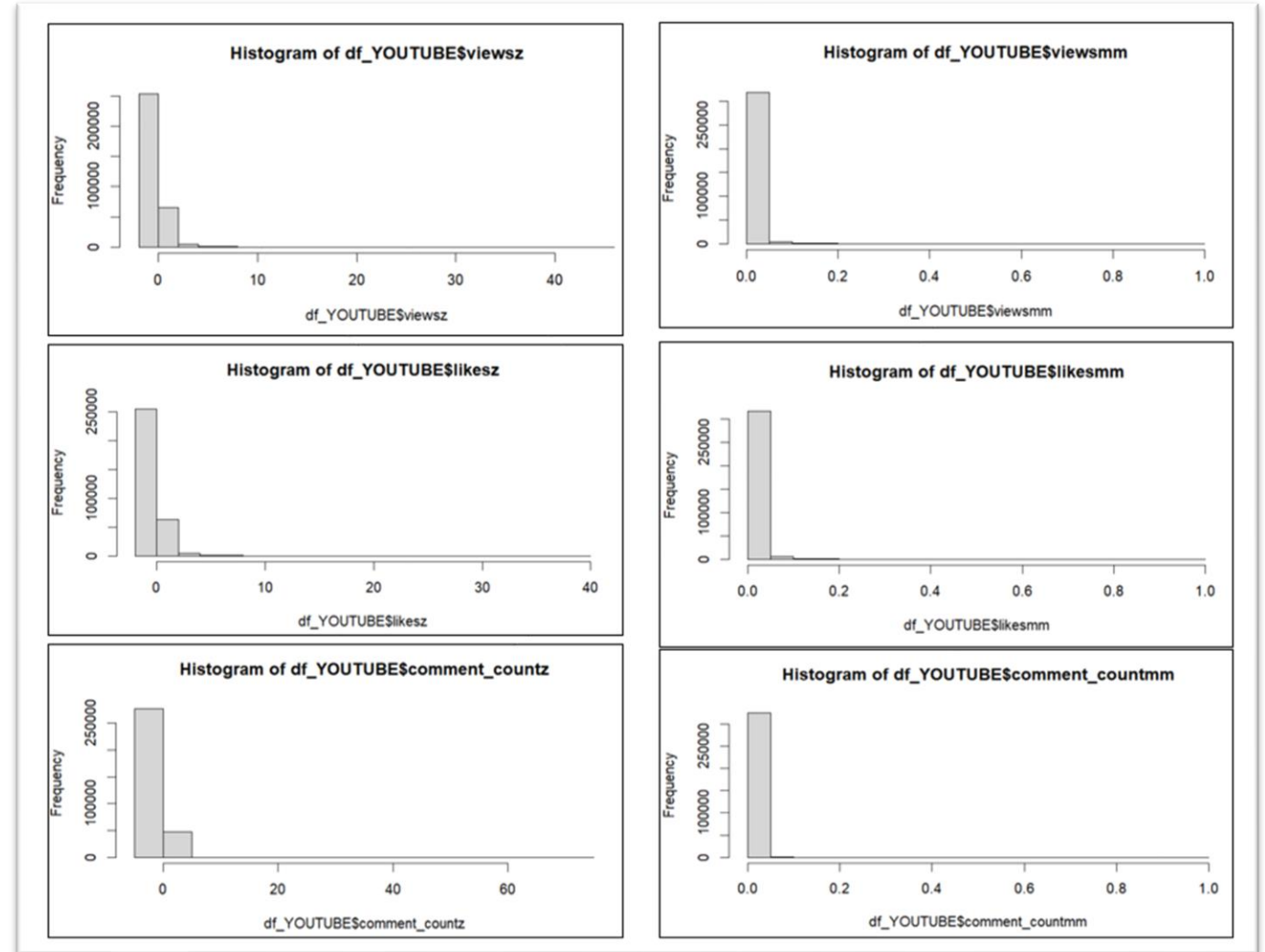
RESULTADOS

Normalización

Resultado 2:

Al ser los rangos en los cuales se distribuye cada variable son muy diferentes por lo que se procedió a normalizarlos, tanto por el método de Puntaje Z, para que conserve la distribución de los resultados, como MINMAX.

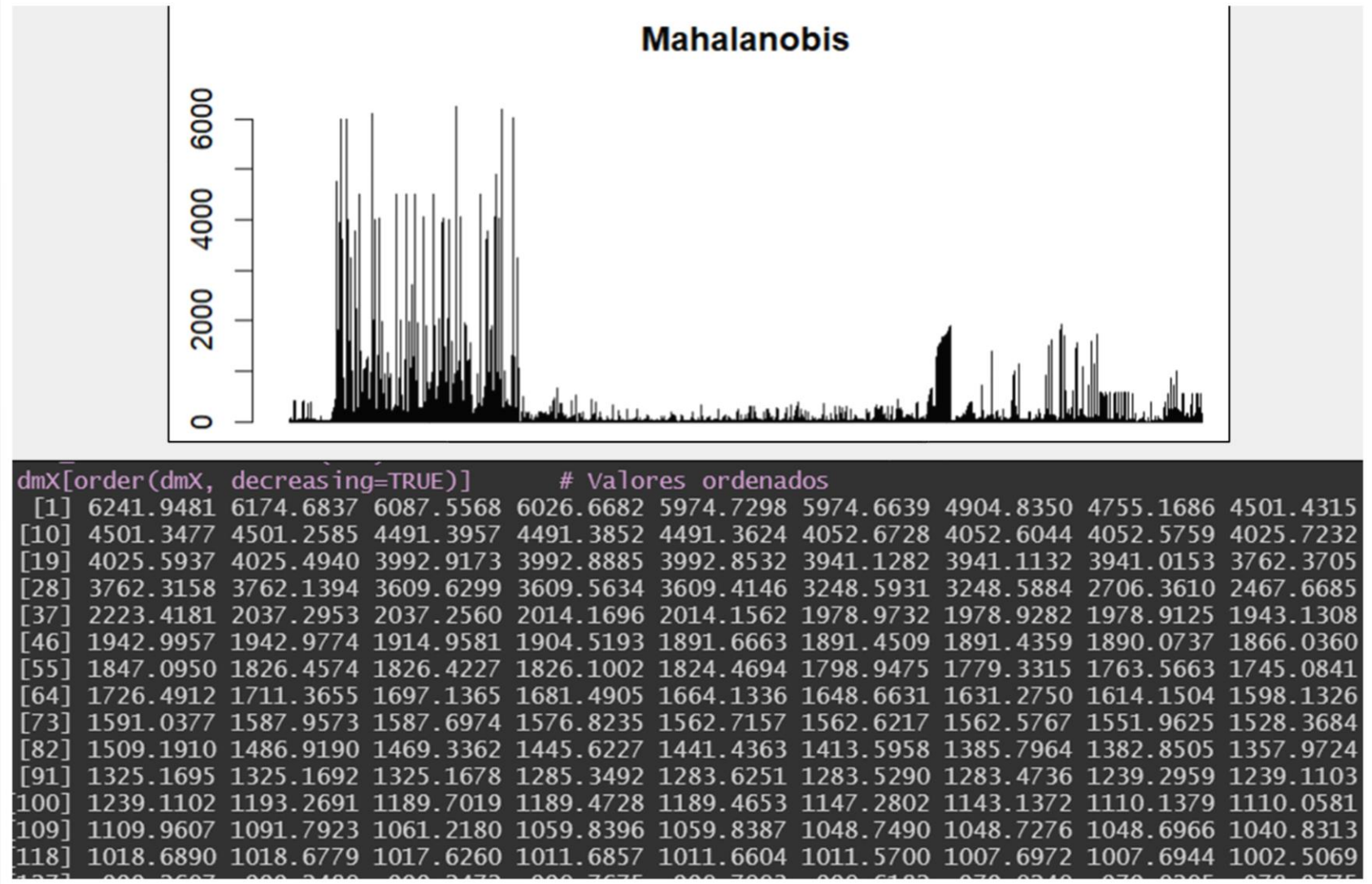
- Método de Puntaje Z
- Método de MINMAX



RESULTADOS

Resultado 3:

Hemos logrado identificar aquellos valores que en su combinación de sus variables views, likes y coments presentan un comportamiento muy diferenciado al resto, logrando valores lejos del centro alcanzando un valor de 6,241.

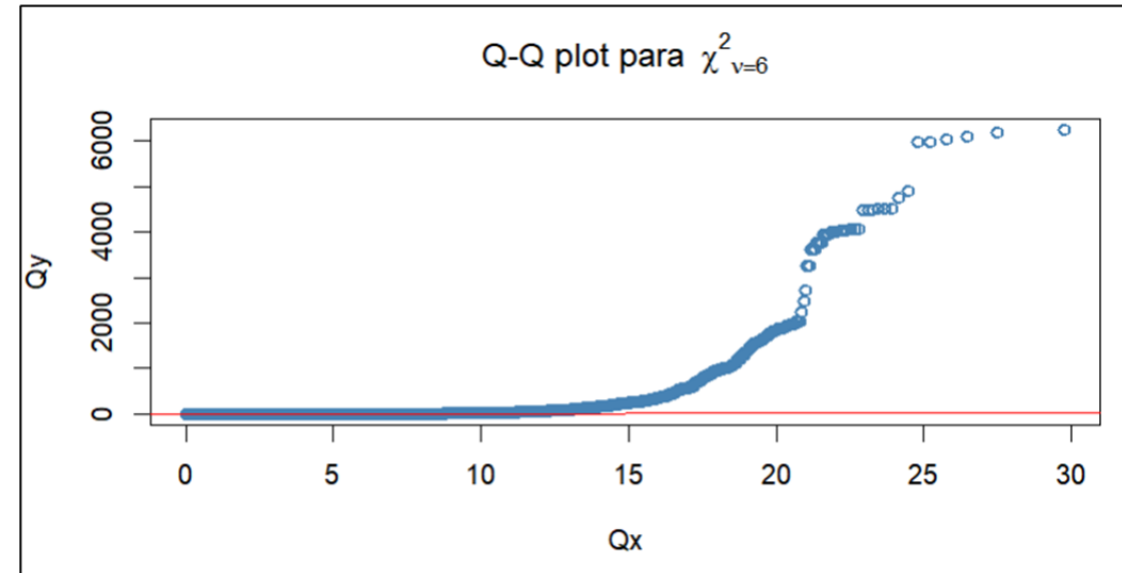
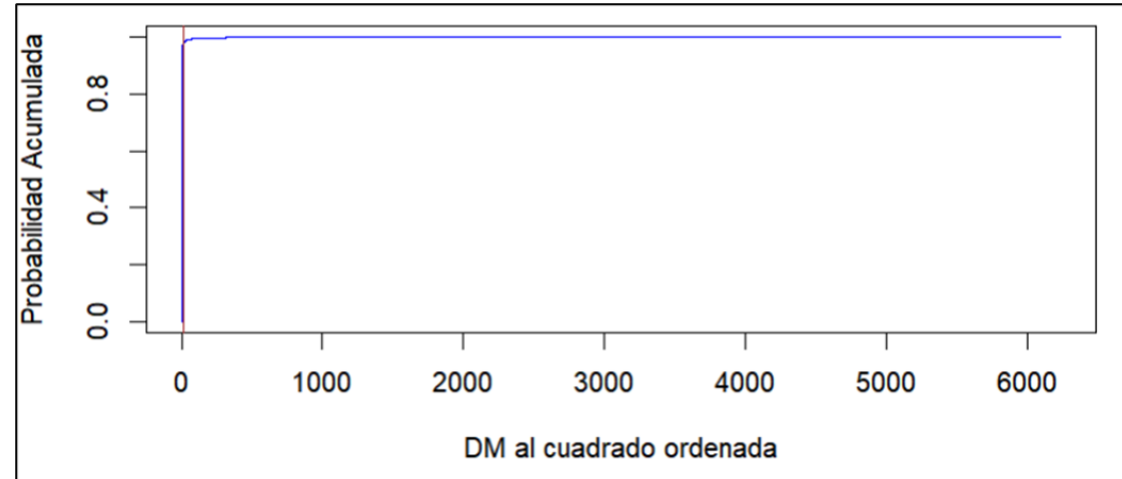


RESULTADOS

Resultado 4:

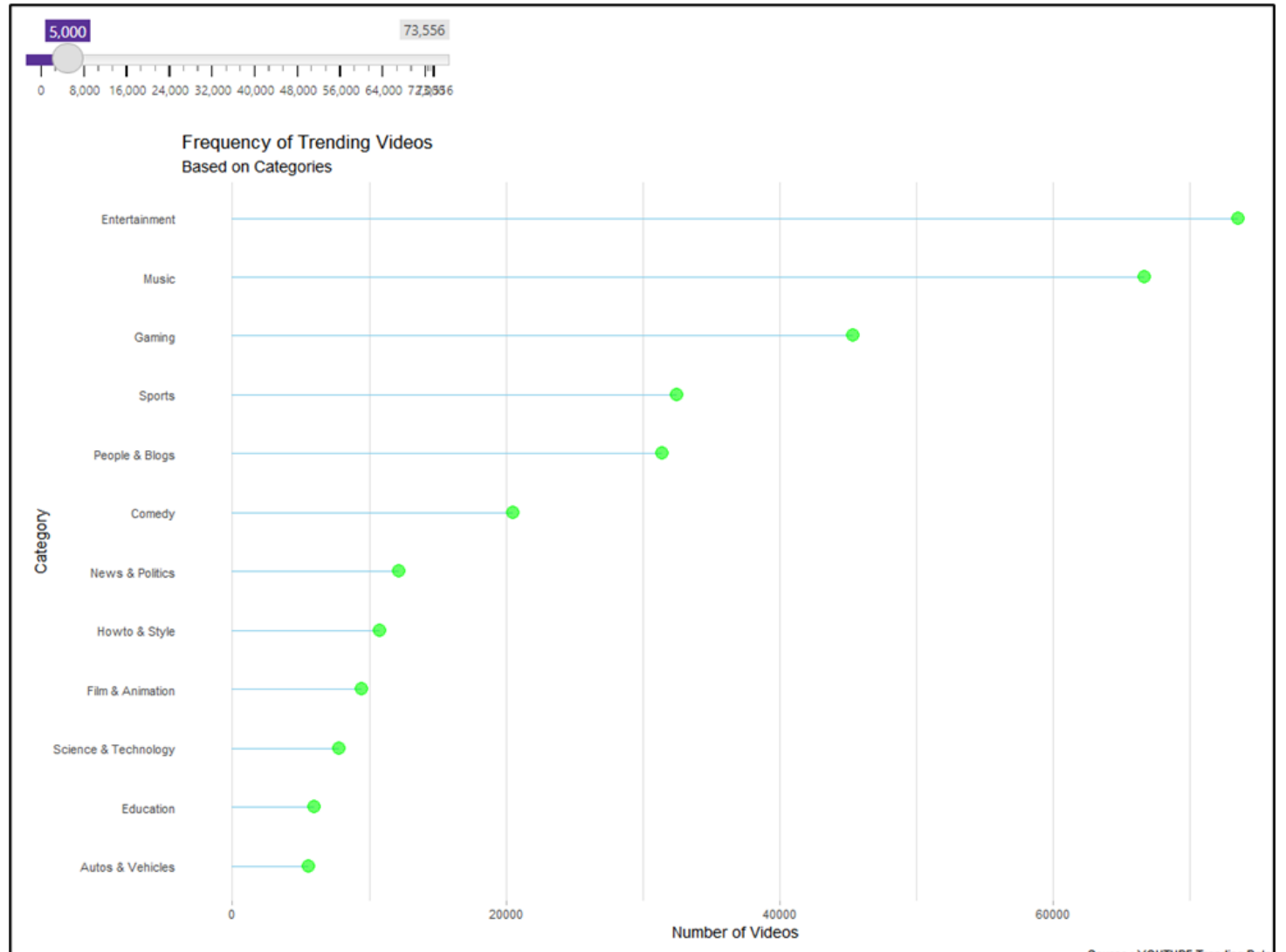
El umbral encontrado a través de chicuadro es de 16.26. Los datos que superan ese valor de dispersión son 5,532; lo cual representa el 1,70%.

Sin embargo, si la evaluación es mediante la gráfica de cuantil – cuantil, observamos una mayor dispersión en relación al valor de chicuadro . Los resultados se presentan a continuacion.



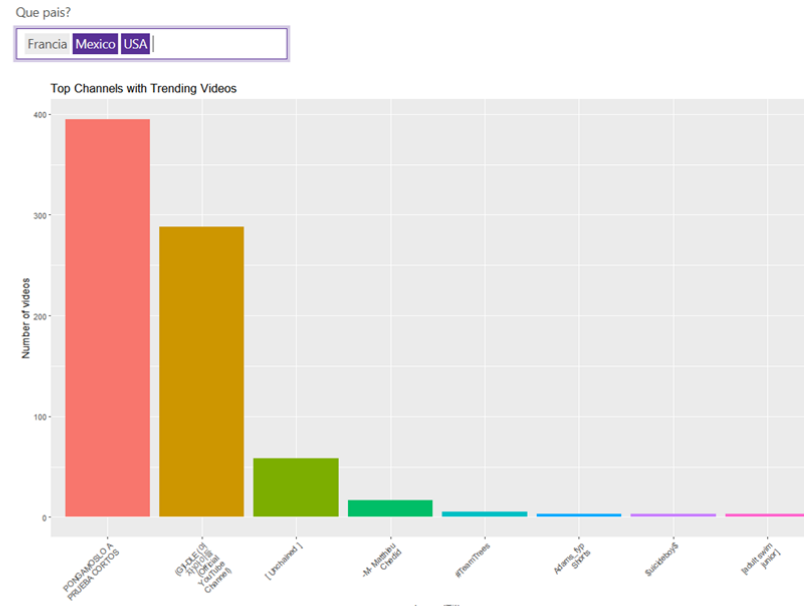
RESULTADOS

Resultado 5: Encontramos en el análisis una mayor inclinación por la generación de videos en las categorías de Entretenimiento y Música, del otro lado del análisis se encuentran Vehículos, Educación, y Ciencia y tecnología.

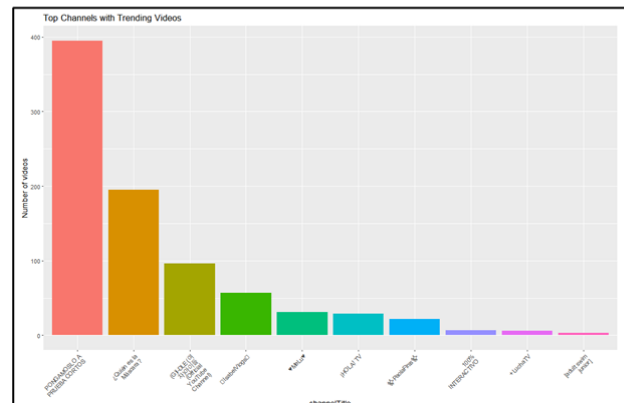


RESULTADOS

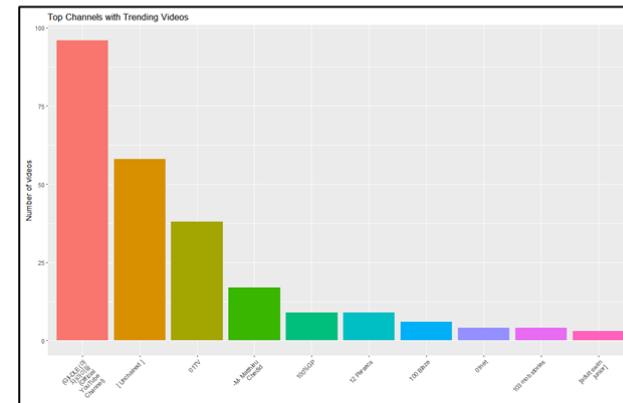
Resultado 6: Es interesante, encontrar que entre los canales que más contenido sube, está en primer lugar “Pongámoslo a prueba cortos”. Lo que sorprende es un canal chino que hace contenido para cada uno de los países, posicionándose como uno de los que más generen en cada uno de ellos.



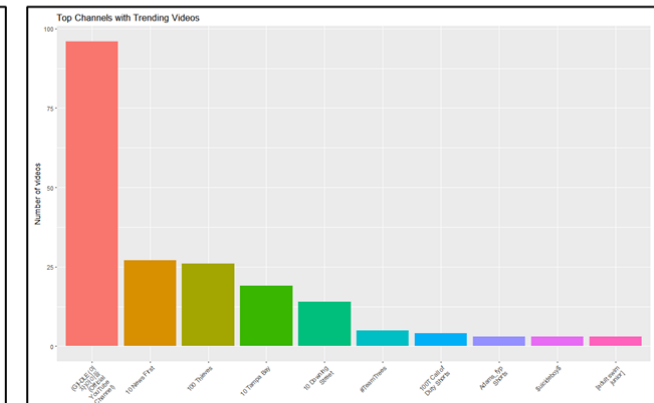
Francia



USA



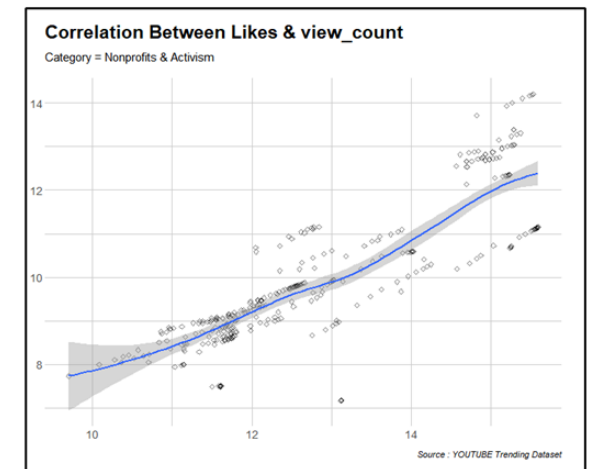
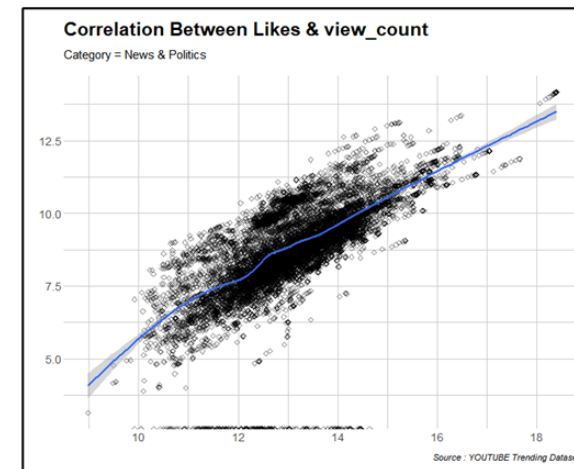
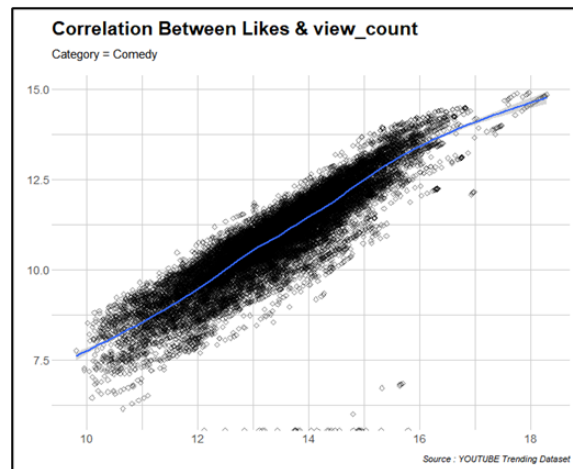
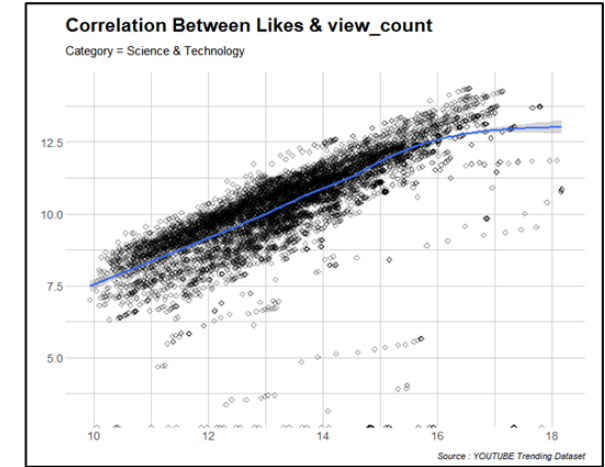
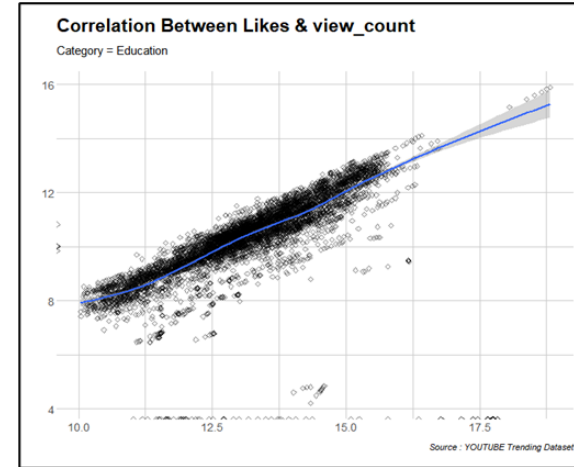
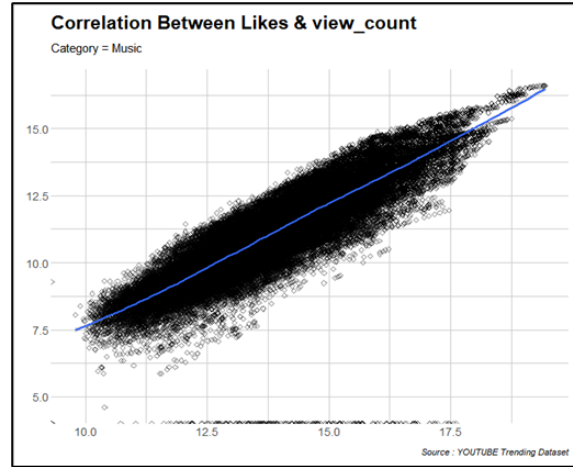
Mexico



RESULTADOS

Resultado 7: Por último, hicimos el ejercicio para evaluar la correlación entre el número de reproducciones y LIKES.

- De primera vista, vemos que las categorías que podríamos llamar divertidas, como “Música” y “Comedy”, tienen una distribución casi homogénea en estas dos variables. No pasa de forma tan notoria en “Educación” y “News Politycs”. Menos aún en “Ciencias y Tecnología”, y “Sin fines de lucro y activismo”.
- Probablemente, estos resultados pueden ayudar a ahondar en la forma psicología del comportamiento de las personas. Que si bien en general, todos tienen una misma dirección en cosas no serias, en temas que son más serios, existe ligera diferencia en la opinión de cada persona



CONCLUSIONES

- La herramienta de MongoDB nos permitió extraer los datos de manera segura y organizada utilizando usuarios de solo lectura con la finalidad de la inmutabilidad de la base de datos.
- Hemos comprobado la facilidad de usar una instancia virtualizada como EC2, la misma que ha ayudado al traslado de altos volúmenes de datos sin interrupción y con la posibilidad de programarse como una tarea repetitiva.
- Hemos constatado la capacidad de R-Studio sobre su capacidad para la conexión a una base de datos como servicio, como es MongoDB, la misma que después de la exploración y procesamiento de la información también ayudo a su presentación de información mediante Shiny Apps.
- Hemos encontrado también una limitación en la conexión a la base de datos como servicio, cuando se trata de la extracción de datos pesados, ya que en oportunidades la carga exigió bastante recurso computacional. Revisar la cantidad de conexiones máximas hacia la BD, es necesario.
- Hemos encontrando que un canal chino genera una cantidad significativa de videos en diferentes países tales como USA, México y Francia. De forma muy diferente al comportamiento de otros canales que tienen mayor presencia en ciertos países.
- Menos del 0.01% podemos decir que son videos que han logrado una variación representativa al resto. Si lo vemos en números, existe una probabilidad de 0.01%. de que si generas un video alcance un gran éxito.

BIBLIOGRAFÍA

A continuación, mostramos los recursos que nos permitieron completar este trabajo

- A. N. PETTITT, A two-sample Anderson-Darling rank statistic, *Biometrika*, Volume 63, Issue 1, 1976, Pages 161–168, <https://doi.org/10.1093/biomet/63.1.161>
- Emad-Eldin, A. A., & Öztürk, A. (1988). A modified one-sample QQ plot and a test for normality. *Journal of Statistical Computation and Simulation*, 29(1), 1-15.
- G. Fasano, A. Franceschini, A multidimensional version of the Kolmogorov–Smirnov test, *Monthly Notices of the Royal Astronomical Society*, Volume 225, Issue 1, March 1987, Pages 155–170, <https://doi.org/10.1093/mnras/225.1.155>
- Hanusz, Z., & Tarasińska, J. (2015). Normalization of the Kolmogorov–Smirnov and Shapiro–Wilk tests of normality. *Biometrical Letters*, 52(2), 85-93.
- Lee, R., Qian, M., & Shao, Y. (2014). On rotational robustness of Shapiro-Wilk type tests for multivariate normality. *Open Journal of Statistics*, 4(11), 964.
- Peter A. W. Lewis. (1961). Distribution of the Anderson-Darling Statistic. *The Annals of Mathematical Statistics*, 32(4), 1118–1124. <http://www.jstor.org/stable/2237910>
- Ronald L. Iman (1982) Graphs for use with the Lilliefors Test for Normal and Exponential Distributions, *The American Statistician*, 36:2, 109-112, DOI: 10.1080/00031305.1982.10482799
- Walfish, S. (2006). A review of statistical outlier methods. *Pharmaceutical technology*, 30(11), 82.



Gracias