

Introducción a la Ciencia de Datos

Maestría en Ciencia de Datos y Aprendizaje Automático

Facultad de Ingeniería, UdelaR

11/07/2024

Tarea 3

Joana Auriello, Pablo Molina

Introducción

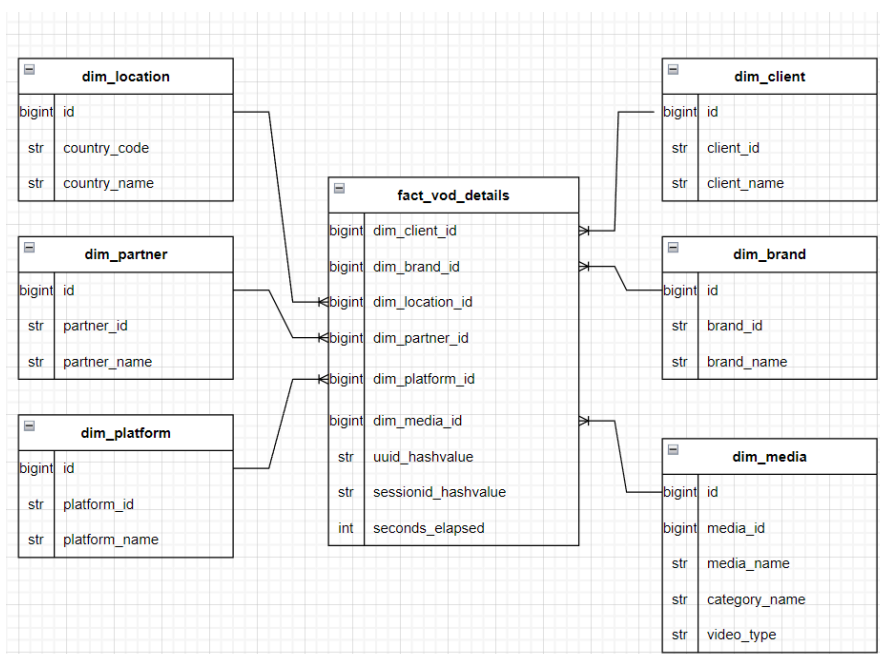
En este informe, se aborda la problemática de la segmentación de usuarios en plataformas de streaming con el objetivo de mejorar la personalización de contenido y las estrategias de retención. Utilizando datos provenientes de los logs de aplicaciones de servicios de streaming, que incluyen información sobre dispositivos, sistemas operativos, plataformas y ubicaciones geográficas. Este estudio se centra en analizar la calidad de los datos y aplicar métodos de ciencia de datos para segmentar a los usuarios basándose en sus patrones de visualización y otras características.

Datos: descripción y calidad

Los datos a analizar surgen de los logs de aplicaciones de servicios de streaming. De estos logs se extraen datos vinculados a los dispositivos, sistemas operativos, plataformas, ubicaciones geográficas, entre otros, de los usuarios que están consumiendo el contenido disponible en estas apps de streaming. Estos logs son procesados a través de pipelines ETL en la nube y se cargan a un data warehouse periódicamente.

Por razones de granularidad de los datos y privacidad, a pesar de que mantenemos una muestra representativa, en este trabajo se presenta solamente una parte del esquema completo con el que se trabaja usualmente.

Modelo de datos



En este modelo, la tabla fact_vod_details actúa como tabla de hechos y contiene los datos principales de visualización, incluyendo uuid_hashvalue (llave primaria), seconds_elapsed y sessionid_hashvalue. Las llaves foráneas en esta tabla (dim_client_id, dim_brand_id, dim_location_id, dim_partner_id, dim_platform_id, dim_media_id) se relacionan con las

tablas dimensionales correspondientes. La tabla dim_client incluye id (llave primaria), client_id y client_name. La tabla dim_brand contiene id (llave primaria), brand_id y brand_name. La tabla dim_location tiene id (llave primaria), country_code y country_name. La tabla dim_partner incluye id (llave primaria), partner_id y partner_name. La tabla dim_platform contiene id (llave primaria), platform_id y platform_name. Finalmente, la tabla dim_media incluye id (llave primaria), media_id, media_name, category_name y video_type. Estas relaciones permiten la integración de datos detallados sobre clientes, marcas, ubicaciones, sistema operativo, plataformas y títulos

A partir de este modelo se genera un dataset que contiene las siguientes columnas:

- client_id, corresponde al número de cliente, dueño del canal.
- client_name, nombre del cliente.
- brand_id, código del canal
- brand_name, nombre del canal
- country_code, código único del país.
- country_name, nombre del país.
- partner_id, código del dueño de la plataforma.
- partner_name, nombre del dueño de la plataforma
- platform_id, código del sistema operativo.
- platform_name, nombre del sistema operativo.
- media_id, código del producto que fue visto.
- media_name, nombre del producto que fue visto.
- category_name, categoría del canal.
- video_type, tipo de video visto (película, episodio, clip, tráiler)
- UUID_HASHVALUE, código para identificar usuarios, se crea a partir de atributos del dispositivo, geográficos y la IP (solo se usa cuando se tiene restricciones para compartir la información) esto se hace para poder identificar el usuario sin comprometer los datos personales.
- seconds_elapsed, tiempo total que la persona permanece en la aplicación.
- sessions, número total de sesiones, definido como count distinct de sessionid_hashvalue, siendo este último un identificador basado en uuid_hashvalue y un componente de tiempo.

Calidad de datos

Los datos cuentan con cierto tratamiento y transformaciones para su calidad, algunos de los siguientes puntos listados ya que fueron corregidos o están en proceso de corrección:

- **Inconsistencia en las Categorías:** Dado que cada cliente define sus propias categorías en su aplicación, es posible encontrar categorías similares con nombres diferentes, como "Comedia" y "Comedy". Esto puede dificultar el análisis agregado y la comparación entre usuarios. Existe una categoría universal de categorías denominadas IAB (Interactive Advertising Bureau), para solucionar este problema se puede matchear las categorías proveniente de cada brand con el estándar internacional.

- **Datos Faltantes o Incompletos:** Puede haber **registros con valores faltantes** en campos clave como `client_name`, `brand_name`, `country_name`, etc. Esto afectaría la precisión de los análisis y reportes.
- **Duplicación de Datos:** Podrían existir **duplicados** en la tabla `fact_vod_uu_details` o en las tablas dimensionales, lo que inflaría los resultados de conteos y sumas.
- **Valores Nulos o Incorrectos en Llaves Foráneas:** Las llaves foráneas en `fact_vod_uu_details` pueden tener valores nulos o incorrectos, impidiendo la correcta relación con las tablas dimensionales y afectando la integridad referencial.
- **Problemas de Integridad Referencial:** Falta de correspondencia entre las llaves primarias de las tablas dimensionales y las llaves foráneas en la tabla de hechos puede llevar a registros huérfanos o a la imposibilidad de realizar uniones correctas.
- **Formato de Datos Inconsistente:** Diferencias en el formato de datos, como fechas en diferentes formatos o códigos de país inconsistentes (USA vs US), pueden complicar la consolidación y análisis de datos.

Para abordar estos problemas, se podría implementar procesos de limpieza y estandarización de datos, como la normalización de categorías, la validación de llaves foráneas, la deduplicación de registros y la implementación de controles de calidad de datos durante la recolección y carga de datos.

Problemática a resolver con herramientas del curso

Problemas/Preguntas a Resolver

Basándonos en sus patrones de visualización y características demográficas, surgen los siguientes problemas que pueden ser resueltos con datos:

- ¿Cómo podemos segmentar a los usuarios para mejorar la personalización de contenido y las estrategias de retención?
- ¿Podemos predecir el tiempo que permanecerán en la aplicación?
- ¿Es posible automatizar el proceso de mapeo entre las categorías estándar y las utilizadas por cada aplicación?

Descripción de la implementación de la solución

Métodos

Para resolver este problema, se puede seguir un proceso estructurado que involucre varias etapas de la ciencia de datos, aplicando conceptos y herramientas presentadas en el curso.

1. Recolección y Exploración de Datos

- **Recolección:** Usar la consulta proporcionada para recolectar datos relevantes de las tablas fact_vod_uu_details, dim_client, dim_brand, dim_location, dim_partner, dim_platform y dim_media.
- **Exploración:** Realizar un análisis exploratorio de datos (EDA) para entender las distribuciones, valores atípicos y patrones generales. Visualizar datos usando histogramas, gráficos de barras y matrices de correlación para identificar relaciones importantes entre variables.

2. Limpieza y Calidad de Datos

- **Limpieza:** Tratar los valores faltantes y duplicados. Normalizar las categorías de medios para evitar inconsistencias. Validar y corregir las llaves foráneas para mantener la integridad referencial.
- **Calidad de Datos:** Implementar controles de calidad para asegurar que los datos sean precisos, completos y consistentes.

3. Segmentación y Aprendizaje Automático

- **Segmentación de Clientes:** Utilizar métodos de aprendizaje no supervisado como K-means clustering para segmentar a los usuarios en grupos basados en sus patrones de visualización y características demográficas.
- **Predicción de duración de la sesión:** Utilizando Random Forest, este algoritmo se puede utilizar tanto para problemas de clasificación como de regresión, su faceta de regresión será la que nos permita a partir de las características del usuario estimar el tiempo que este permanecerá en la aplicación.
- **Automapeo de categorías:** En este caso es posible la utilización de medidas de distancia que junto con un algoritmo de predicción nos permita comparar el parecido entre nuestras categorías y el estándar, definiendo un threshold de certeza podemos automatizar que se remplacen aquellas categorías que tienen una distancia menos o un mayor parecido con las categorías del sistema de clasificación estándar.

4. Visualización e Interpretación de Resultados

- **Visualización de Resultados:** Crear gráficos de dispersión, diagramas de cajas y gráficos de radar para visualizar los segmentos de usuarios. Usar mapas de calor para visualizar la correlación entre las variables y los clústeres formados.
- **Interpretación:** Interpretar los resultados para entender las características y comportamientos de cada segmento. Identificar patrones comunes y diferencias clave entre los segmentos.

5. Aplicaciones Prácticas

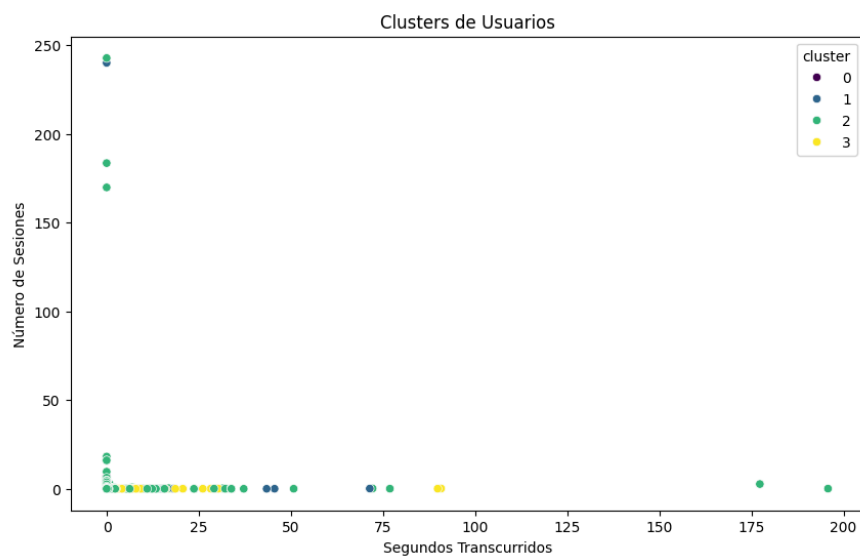
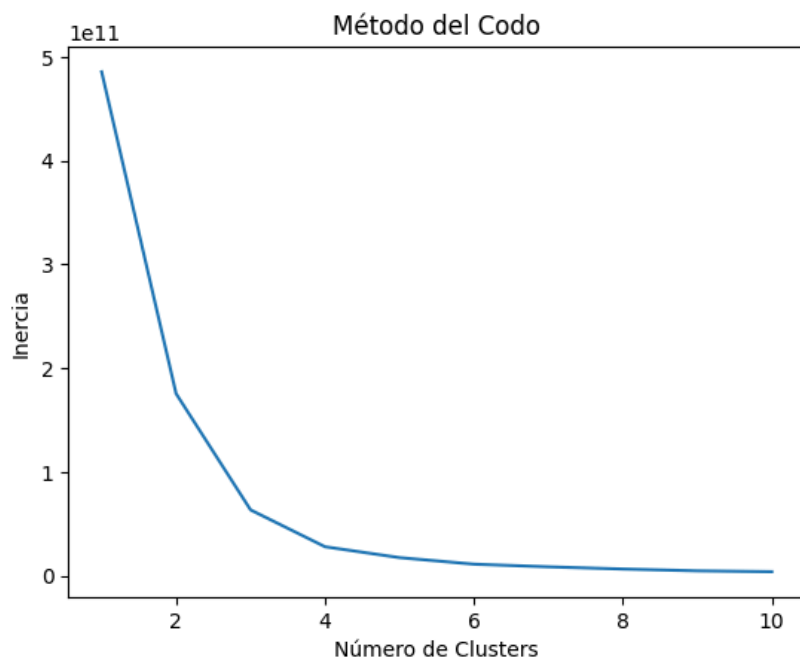
- **Personalización de Contenido:** Utilizar los segmentos de clientes para desarrollar recomendaciones personalizadas de contenido.
- **Estrategias de Retención:** Diseñar estrategias de retención específicas para cada segmento basado en sus características y comportamientos.

Anexo: Implementación

A modo de prueba buscamos lograr una prueba de concepto de cómo funcionaria el análisis de K-means en la base de datos planteada.

Visualizaciones

Seguimos el método del código para identificar que el número ideal de clusters sería 4.



1	0	1.998 967	95.4014 62	46.21 7162	1.0699 32	3985. 8927 3	5.85194 4	2.6487 42	-0.03707 9	0.00 1193
2	0	1.764 59	106.282 952	56.13 0338	1.1314 05	2497. 9164 6	12.6245 13	2.9398 94	-0.00022 4	0.00 5801
3	0	1.996 013	112.517 399	45.78 7697	1.6889 5	5584. 1902 4	23.3456 92	2.3455 88	0.248858	-0.01 6577

n=3

cluster	client_id	brand_id	country_code	partner_id	platform_id	media_id	category_name	video_type	seconds_elapsed	sessions
0	0	1.998 152	100.071 328	46.45 0271	1.2609 63	4443. 1777 5	10.8026 47	2.5729 41	0.045648	-0.00 3644
1	0	1.954 623	109.847 692	50.67 6345	1.3208 5	315.7 1921 6	26.3344 41	4.4631 35	-0.08581 5	-0.00 4428
2	0	1.771 43	105.987 627	55.51 3435	1.1095 91	2524. 3005 3	12.2873 26	2.9245 14	-0.00081 6	0.00 5634