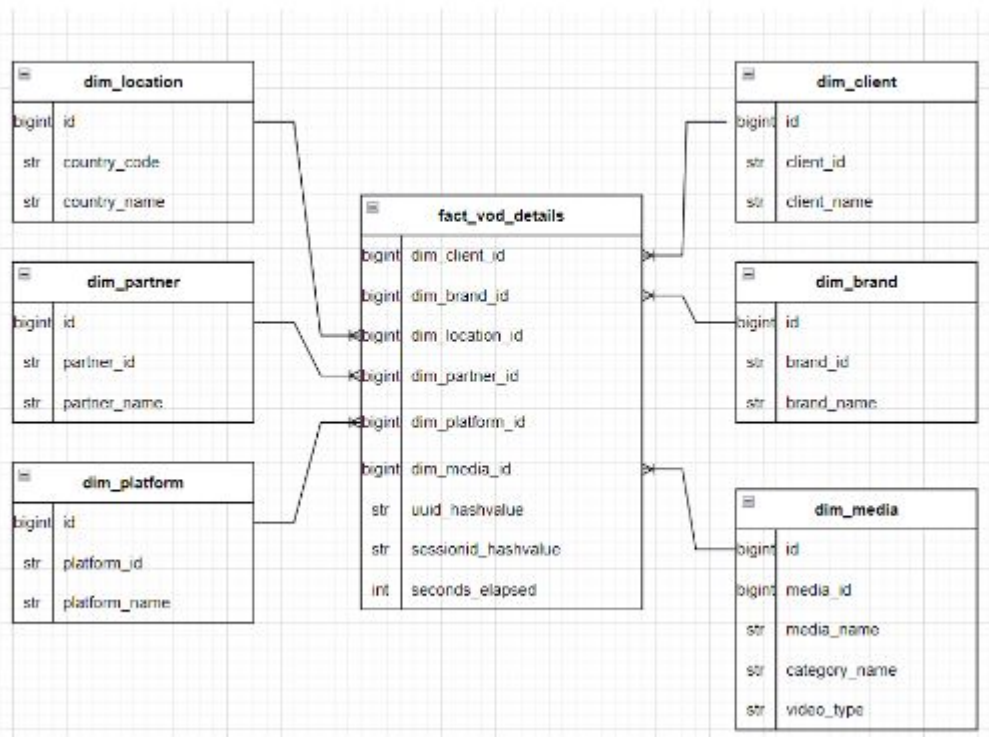




Introducción a la ciencia de datos - Tarea 3

Grupo 20 Auriello - Molina
11/07/2024

Datos



Calidad de datos



Problemas resueltos:

- **Inconsistencia en las Categorías**
- **Datos Faltantes o Incompletos**
- **Duplicación de Datos**
- **Valores Nulos o Incorrectos en Llaves Foráneas**
- **Problemas de Integridad Referencial** (Falta de correspondencia entre las llaves primarias de las tablas dimensionales)
- **Formato de Datos Inconsistente (US vs USA)**

Problemática 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?

Soluciones propuestas



- **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.