



# PROYECTO FINAL

GOOGLE + YELP – Grupo 8



## INFORME DE AVANCE 2 – DATA ENGINEERING

### Introducción

De acuerdo a lo establecido en la propuesta inicial, en este informe de avance N°2 se presenta la Ingeniería de Datos del Proyecto.

Este es un pilar fundamental en el desarrollo, ya que se desarrollan las bases para que los procesos sean automáticos y el dataware house pueda ser alimentado constantemente con nuevas bases de datos.

### Objetivos

A continuación se establecen los objetivos propuestos para este segundo informe de avance:

- a) ETL de Datasets de Google Bussiness y Yelp
- b) Diccionario de datos
- c) Arquitectura Propuesta y Diagrama
- d) Creación del DataWare House
- e) Automatizar el DataWare House
- f) Desarrollar Proceso Carga Incremental.
- g) Utilizar un servicio de Nube

#### A) ETL de las bases de Google Bussiness y YELP!

Antes de la creación del DataWare House y de comenzar con el proceso de Analytics y Machine Learning, es fundamental y necesario, realizar una correcta Limpieza y Transformación de los datos.

En este sentido, en primer lugar se ha realizado un ETL de cada uno de los datasets con los que se trabajará.

Esto incluyó entre otras acciones, eliminación la estandarización de datos, eliminación de columnas, filtrado a los estados con los que se trabajará, filtrado de datos al rubro bar, nightlife y hoteles, procesamiento de nulls en caso de encontrarlos, creación de nuevas columnas, etc.

#### B) Diccionario De Datos

Una vez realizado el ETL y con una versión final de cada uno de los datasets, se procedió a elaborar los diccionarios de datos

Los mismos se presentan a continuación:



## INFORME DE AVANCE 2 – DATA ENGINEERING

### DICCIONARIO DE DATOS YELP

#### i. Diccionario de Bussines

- **business\_id**: Un identificador único para cada negocio en la base de datos.
- **name**: El nombre del negocio.
- **address**: La dirección del negocio.
- **city**: La ciudad donde se encuentra el negocio.
- **state**: El estado o provincia donde se encuentra el negocio.
- **postal\_code**: El código postal del área donde se encuentra el negocio.
- **latitude**: La latitud geográfica del negocio.
- **longitude**: La longitud geográfica del negocio.
- **stars**: La calificación promedio del negocio en Yelp, en una escala de 1 a 5 estrellas.
- **review\_count**: El número total de reseñas que ha recibido el negocio en Yelp.
- **categories**: Las categorías o etiquetas que describen el tipo de negocio, separadas por comas.
- **Categoría**: se refiere a la categoría objetivo encontrada en ese negocio.

#### ii. Diccionario de checkin

- **business\_id**: Es un identificador único para cada negocio en Yelp.
- **num\_visitas**: Representa número de veces que se registró un check-ins (visita) en el negocio.

#### iii. Diccionario de review

- **user\_id** (cadena): Es el identificador único del usuario que ha escrito la reseña.
- Cada usuario tiene un **id\_usuario** único que lo distingue de otros usuarios.
- **business\_id** (cadena): Es el identificador único del negocio que recibió la reseña. Cada negocio tiene una Identificación del negocio única que lo distingue de otros negocios.
- **stars** (entero): Es la calificación de la reseña en términos de estrellas. Puede ser un valor entero entre 1 y 5, donde 1 es la calificación más baja y 5 es la calificación más alta.
- **useful** (entero): Representa el número de votos útiles que ha recibido la reseña de otros usuarios. Los usuarios pueden votar si una reseña les resultó útil o no.
- **funny** (entero): Representa el número de votos de otros usuarios que encontraron la reseña divertida. Los usuarios pueden votar si una reseña les resultó divertida o no.
- **cool** (entero): Representa el número de votos de otros usuarios que encontraron un cierto factor de fascinación, estilo o innovación que atrae la atención y genera una reacción positiva.
- **text** (cadena): Es el contenido de la reseña escrita por el usuario. Aquí se encuentra el texto completo de la reseña que describe la experiencia del usuario con el negocio.



## INFORME DE AVANCE 2 – DATA ENGINEERING

- Date(fecha): representa a la fecha en la que se hizo la reseña.
- Sentimiento\_score (decimal): contiene valores entre -1 y 1. Estos valores representan el sentimiento a la columna text. Un valor cercano a 1 indica un sentimiento positivo fuerte, mientras que un valor cercano a -1 indica un sentimiento negativo fuerte. Los valores cercanos a 0 sugieren un sentimiento neutro o ambiguo en relación con el texto correspondiente.

### iv. Diccionario de tips

- user\_id: Es un identificador único para cada usuario en Yelp que ha dejado un consejo (tip) en un negocio.
- business\_id: Es un identificador único para el negocio en el que se dejó el consejo.
- text: Es el contenido del consejo que el usuario dejó para el negocio. Es un comentario o sugerencia que el usuario comparte.
- Date(fecha): representa a la fecha en la que se hizo la tips.
- Sentimiento\_score (decimal): contiene valores entre -1 y 1. Estos valores representan el sentimiento a la columna text. Un valor cercano a 1 indica un sentimiento positivo fuerte, mientras que un valor cercano a -1 indica un sentimiento negativo fuerte. Los valores cercanos a 0 sugieren un sentimiento neutro o ambiguo en relación con el texto correspondiente.

### v. Diccionario de user

- user\_id: Es un identificador único para el usuario. Es una cadena de texto que identifica al usuario en Yelp.
- name: Es el nombre del usuario. Es una cadena de texto que representa el nombre del usuario en Yelp.
- review\_count: Es la cantidad total de reseñas que el usuario ha escrito en Yelp. Es un número entero que indica cuántas reseñas ha dejado el usuario.
- useful: Representa la cantidad total de votos "útiles" que el usuario ha recibido en sus reseñas por parte de otros usuarios. Es un número entero que indica cuántos usuarios han encontrado útiles las reseñas del usuario.
- funny: Representa la cantidad total de votos "graciosos" que el usuario ha recibido en sus reseñas por parte de otros usuarios. Es un número entero que indica cuántos usuarios han encontrado graciosas las reseñas del usuario.
- cool: Representa la cantidad total de votos "frescos" que el usuario ha recibido en sus reseñas por parte de otros usuarios. Es un número entero que indica cuántos usuarios han encontrado un cierto factor de fascinación, estilo o innovación que atrae la atención y genera una reacción positiva.
- average\_stars: Representa el número de estrellas asigna por el usuario en sus reseñas.



## INFORME DE AVANCE 2 – DATA ENGINEERING

### DICCIONARIO DE DATOS GOOGLE BUSSINESS

#### DICCIONARIO DIM\_SITIOS\_GOOGLE

- **name:** nombre del usuario, cadena de texto que representa el nombre del usuario  
texto
- **address:** La dirección del negocio, texto
- **gmap\_id:** código de ubicación global de Google, texto
- **latitude:** ángulo de coordenada geográfica, numero decimal
- **longititude:** ángulo de coordenada geográfica, numero decimal
- **categories:** clasificación del establecimiento, texto
- **avg\_rating:** promedio de votación del establecimiento, numero decimal
- **num\_of\_reviews:** número de reseñas del establecimiento, numero entero
- **categoría:** categoría filtrada por Hotel, Bar y Nightlife a la que pertenece el establecimiento, texto

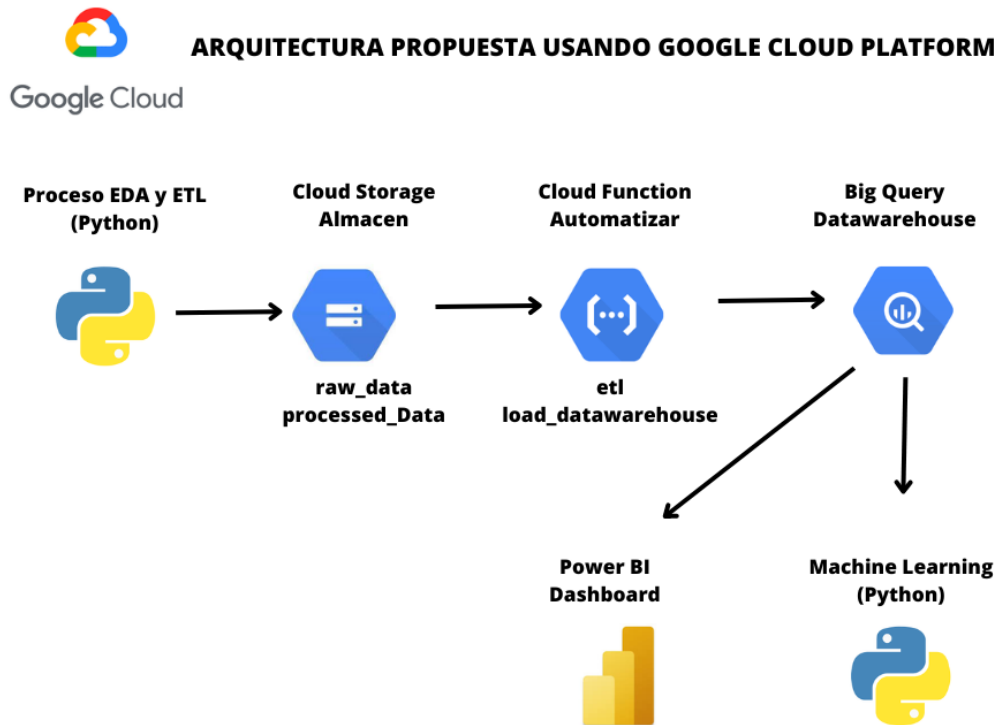
#### DICCIONARIO REVIEW-GOOGLE

- **user\_id:** código del usuario, numero decimal
- **name:** nombre del usuario, texto
- **rating:** votación del usuario al establecimiento, numero entero
- **text:** comentario sobre el establecimiento, texto
- **gmap\_id:** código de ubicación global de Google, texto
- **date:** fecha y hora de la reseña, fecha
- **estado:** sigla del estado donde se encuentra el establecimiento, texto
- **sentimiento\_score:** contiene valores entre -1 y 1. Estos valores representan el sentimiento a la columna text, numero decimal



## INFORME DE AVANCE 2 – DATA ENGINEERING

### C) Arquitectura Propuesta y Diagrama



Por la dimensión de los datasets de YELP y GOOGLE BUSINESS que son demasiados grandes, se eligió la arquitectura Big Data.

La tecnología elegida ha sido GOOGLE CLOUD PLATFORM.

Usamos CLOUD STORARE para almacenar los datasets. Tenemos dos buckets:

- RAW\_DATA: usado para guardar los datasets originales
- PROCESSED\_DATA: usado para guardar los datasets procesado y limpios

Para automatizar usamos CLOUD FUNCTIONS, tenemos dos funciones con lenguaje Python:

- ETL: Realiza el proceso de ETL
- LOAD\_DATWAREHOUSE: importa los datasets procesados al DATAWAREHOUSE

Para crear el DATAWAREHOUSE usamos BIG QUERY, el cual será usado por:

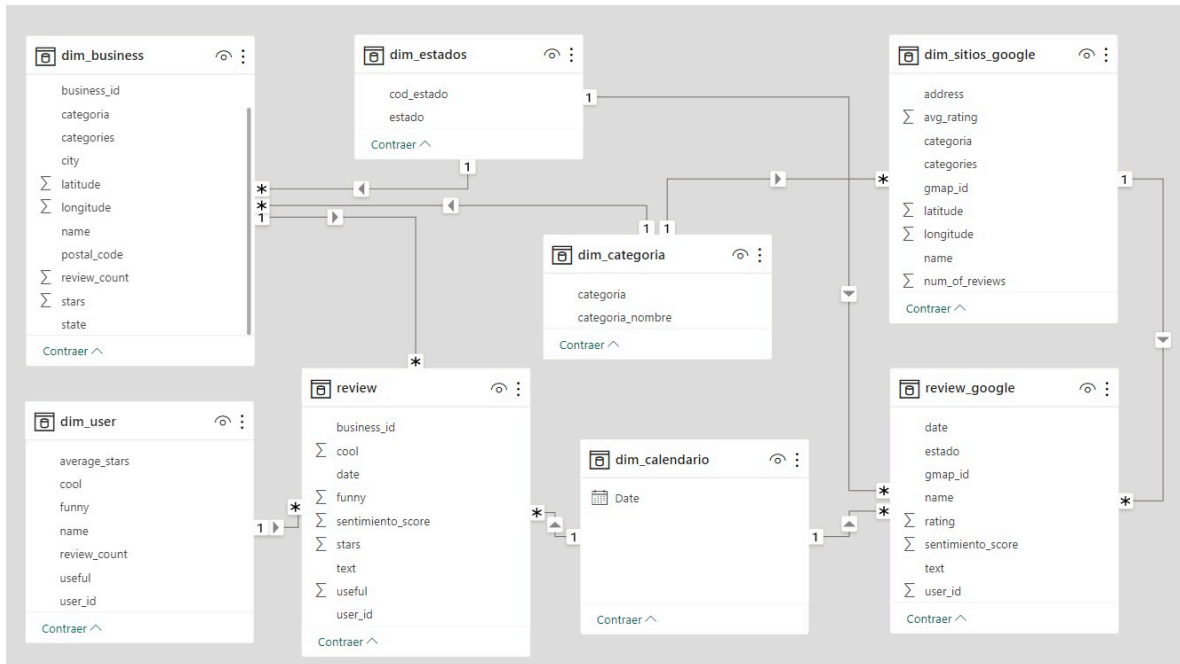
- POWER BI: para diseñar el dashboard
- MACHINE LEARNING: con Python realizamos un sistema de recomendación



## INFORME DE AVANCE 2 – DATA ENGINEERING

### D) Creación del DataWare House

Se presenta a continuación el Diseño del Data Warehouse.



El Datawarehouse diseñado, posee lo siguiente:

- Tablas Dimensiones:
  - o Dim\_bussiness
  - o Dim\_user
  - o Dim\_estados
  - o Dim\_categoria
  - o Dim\_calendario
  - o Dim\_sitios\_google
- Tablas hechos:
  - o Reviews
  - o Review\_google

### E) AUTOMATIZAR EL DW

Para automatizar usamos el DW, usamos CLOUD FUNCTIONS, tenemos dos funciones con lenguaje Python:

- ETL: Realiza el proceso de ETL

