

---

## INFORME DE AVANCE - SPRINT 2

Aylén Sol Guzman, Edith Priscila Muniz Cantu, Maria Marcela Diaz, Jorge Andres Jola Hernandez, Franco Dylan Damian Luna Pedroso.

---

### Introducción

En concordancia con los lineamientos establecidos en la propuesta inicial, el presente Informe de Avance N°2 aborda la Ingeniería de Datos del Proyecto. Este componente se erige como un pilar fundamental en el desarrollo, dado que sienta las bases para la automatización de procesos y permite la alimentación continua del data warehouse con nuevas bases de datos. La implementación de la Ingeniería de Datos resulta crucial para garantizar la eficiencia y la actualización constante de la infraestructura, contribuyendo así al éxito y la evolución continua del proyecto.

En este segundo sprint se llevaron a cabo las siguientes acciones:

1. Utilizar un servicio de Nube
2. Llevar a cabo ETL de Datasets de Google y Yelp
3. Desarrollar Proceso Carga Incremental
4. Automatizar y crear el DataWareHouse
5. Construir diccionario de datos
6. Construcción versión inicial del Dashboard
7. Inicio de Desarrollo de productos de ML

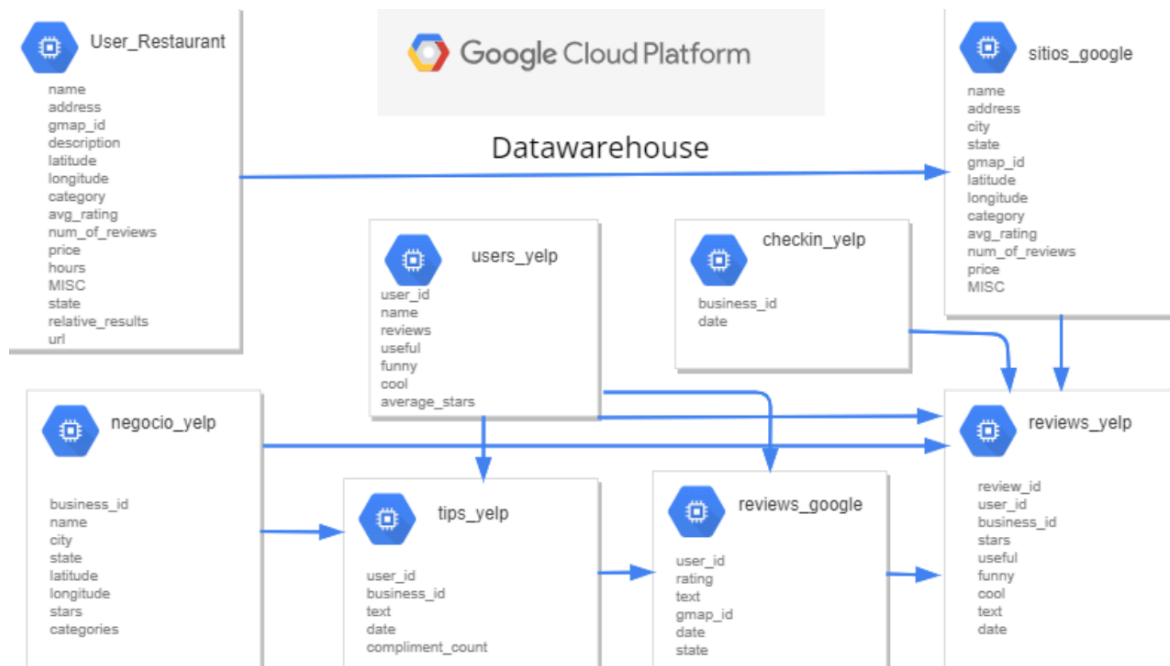
### Servicio de Nube

Para la implementación de un servicio de nube, optamos por Google Cloud debido a sus notables ventajas y eficiencias. Google Cloud ofrece una infraestructura sólida y confiable, respaldada por su capacidad para escalar según las necesidades del proyecto. Su integración con diversas herramientas y servicios facilita la gestión y el despliegue eficiente de aplicaciones y recursos. Además, la seguridad avanzada y las opciones flexibles de almacenamiento y procesamiento de datos hacen de Google Cloud una elección idónea para satisfacer nuestras necesidades de forma integral y efectiva.

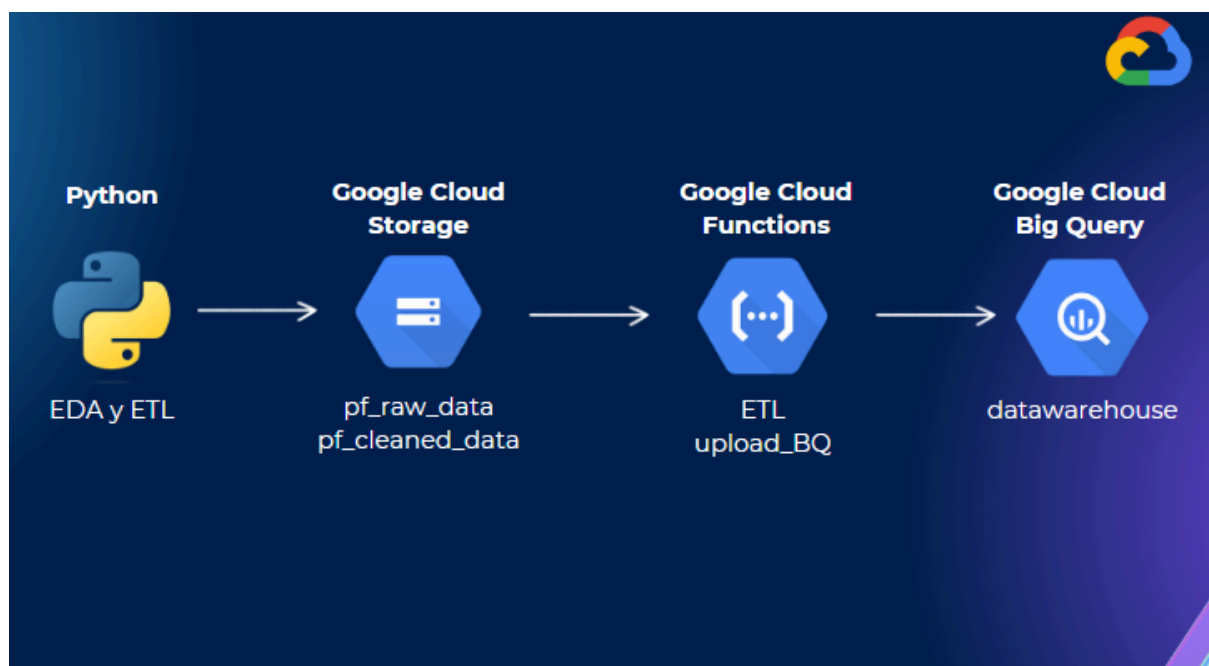
### ETL

Antes de embarcarnos en la creación del DataWareHouse y adentrarnos en el análisis de datos y Machine Learning, resulta imperativo llevar a cabo una meticulosa limpieza y transformación de los datos. En este contexto, se ha realizado inicialmente un proceso de Extract, Transform, Load (ETL) para cada uno de los conjuntos de datos que serán utilizados. Este procedimiento abarcó diversas acciones, como la estandarización de datos, eliminación de columnas innecesarias, filtrado según los estados de interés, focalización en el rubro de restaurantes, tratamiento de valores nulos, entre otras medidas.

## DataWarehouse



## Proceso de Carga Incremental



Dada la magnitud de los conjuntos de datos de Yelp y Google, que se caracterizan por ser extraordinariamente extensos, se optó por implementar una arquitectura Big Data. En esta estrategia, utilizamos Google Cloud Storage como plataforma para almacenar los datasets, organizando la información en dos buckets específicos: "pf\_raw\_data", destinado a guardar los datasets originales, y "pf\_cleaned\_data", destinado a almacenar los datasets procesados y depurados.

Para llevar a cabo la automatización de procesos, empleamos Google Cloud functions, donde se implementaron dos funciones escritas en Python:

1. ETL: Encargada de ejecutar el proceso de Extract, Transform, Load (ETL).
2. Upload\_BQ: Responsable de importar los datasets procesados al Datawarehouse.

En la creación del Datawarehouse se recurrió a Big Query, una herramienta que será aprovechada por diversas instancias como Looker y Python (Machine Learning).

## Diccionarios

### DICCIONARIO DE DATOS YELP

#### 1. Diccionario de Business

- **business\_id**: Identificador único para cada negocio en la base de datos.
- **name**: Nombre del negocio.
- **address**: Dirección del negocio.
- **city**: Ciudad donde se encuentra el negocio.
- **state**: Estado o provincia donde se encuentra el negocio.
- **postal\_code**: Código postal del área donde se encuentra el negocio.
- **latitude**: Latitud geográfica del negocio.
- **longitude**: Longitud geográfica del negocio.
- **stars**: Calificación promedio del negocio en Yelp, en una escala de 1 a 5 estrellas.
- **review\_count**: Número total de reseñas que ha recibido el negocio en Yelp.
- **categories**: Categorías o etiquetas que describen el tipo de negocio, separadas por comas.
- **Categoría**: Se refiere a la categoría objetivo encontrada en ese negocio.

#### 2. Diccionario de Check-in

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

#### 3. Diccionario de Reviews

- **user\_id (cadena)**: Identificador único del usuario que ha escrito la reseña.
- **business\_id (cadena)**: Identificador único del negocio que recibió la reseña.
- **stars (entero)**: Calificación de la reseña en términos de estrellas (1 a 5).
- **useful (entero)**: Número de votos útiles que ha recibido la reseña.
- **funny (entero)**: Número de votos de otros usuarios que encontraron la reseña divertida.
- **cool (entero)**: Número de votos de otros usuarios que encontraron un factor de fascinación en la reseña.
- **text (cadena)**: Contenido completo de la reseña escrita por el usuario.
- **Date (fecha)**: Fecha en la que se hizo la reseña.
- **Sentimiento\_score (decimal)**: Valores entre -1 y 1 que representan el sentimiento hacia el texto de la reseña.

#### 4. Diccionario de Consejos

- **user\_id**: Identificador único para cada usuario en Yelp que ha dejado un consejo en un negocio.
- **business\_id**: Identificador único para el negocio en el que se dejó el consejo.
- **text**: Contenido del consejo que el usuario dejó para el negocio.

- **Date (fecha):** Fecha en la que se hizo el consejo.
- **Sentimiento\_score (decimal):** Valores entre -1 y 1 que representan el sentimiento hacia el texto del consejo.

## 5. Diccionario de Users

- **user\_id:** Identificador único para el usuario en Yelp.
- **name:** Nombre del usuario en Yelp.
- **review\_count:** Cantidad total de reseñas que el usuario ha escrito en Yelp.
- **useful:** Cantidad total de votos "útiles" recibidos en las reseñas del usuario.
- **funny:** Cantidad total de votos "graciosos" recibidos en las reseñas del usuario.
- **cool:** Cantidad total de votos "frescos" recibidos en las reseñas del usuario.
- **average\_stars:** Número de estrellas asignado por el usuario en sus reseñas.

## DICCIONARIO DE DATOS GOOGLE

### 1. Diccionario de sitios\_google

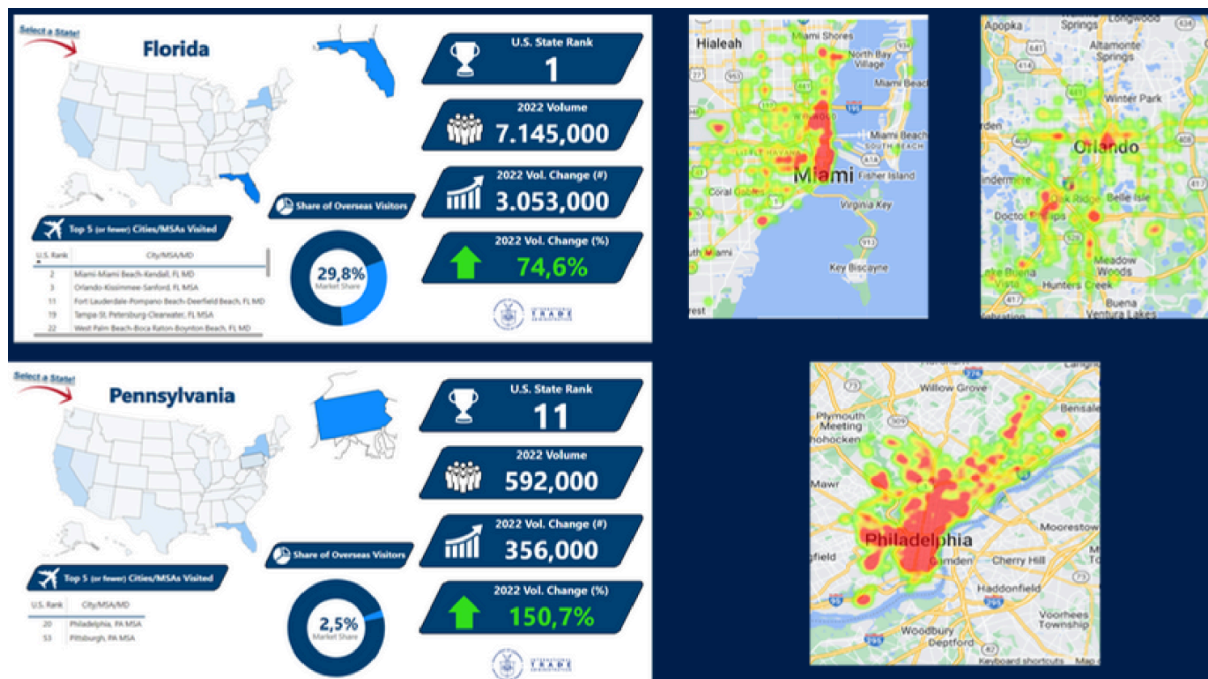
- **name:** Denominación del usuario, secuencia de caracteres que refleja el nombre del usuario.
- **address:** Ubicación del negocio, secuencia de caracteres.
- **gmap\_id:** Código global de ubicación de Google, secuencia de caracteres.
- **latitude:** Ángulo de coordenada geográfica, número decimal.
- **longitude:** Ángulo de coordenada geográfica, número decimal.
- **categories:** Clasificación del establecimiento, secuencia de caracteres.
- **avg\_rating:** Calificación promedio del establecimiento, número decimal.
- **num\_of\_reviews:** Cantidad de reseñas del establecimiento, número entero.
- **category:** Categoría clasificada como Hotel, Bar y Nightlife a la que pertenece el establecimiento, secuencia de caracteres.

### 2. Diccionario de review\_google

- **user\_id:** Código de identificación del usuario, número decimal.
- **name:** Nombre del usuario, secuencia de caracteres.
- **rating:** Calificación del usuario al establecimiento, número entero.
- **text:** Comentario acerca del establecimiento, secuencia de caracteres.
- **gmap\_id:** Código global de ubicación de Google, secuencia de caracteres.
- **date:** Fecha y hora de la reseña, fecha.
- **state:** Sigla del estado donde se ubica el establecimiento, secuencia de caracteres.
- **sentiment\_score:** Define valores entre -1 y 1 que representan el sentimiento hacia la columna text, número decimal.

## Versión inicial del Dashboard





El producto final es un modelo de recomendación desplegado en Google Cloud Vertex AI, diseñado para proporcionar recomendaciones de restaurantes basadas en el análisis de sentimientos y similitud de contenido de las reseñas. Este modelo utiliza técnicas avanzadas de procesamiento de lenguaje natural, como el análisis de sentimientos mediante NLTK, tokenización y lematización, así como la vectorización TF-IDF para representar las características de las reseñas. Además, se emplea el algoritmo de vecinos más cercanos (KNN) con métrica de similitud de coseno para encontrar restaurantes similares en función de las preferencias del usuario. El producto incluye una función que toma como entrada una ciudad específica y devuelve recomendaciones personalizadas. Este enfoque permite una adaptabilidad significativa, y se menciona una mejora adicional para incorporar la cantidad mínima de estrellas como entrada para ajustar aún más las sugerencias según las preferencias

individuales del usuario. El modelo se almacena en Google Cloud Storage y se despliega en Vertex AI Prediction, ofreciendo así una solución integral y escalable para el desarrollo de sistemas de recomendación basados en contenido.

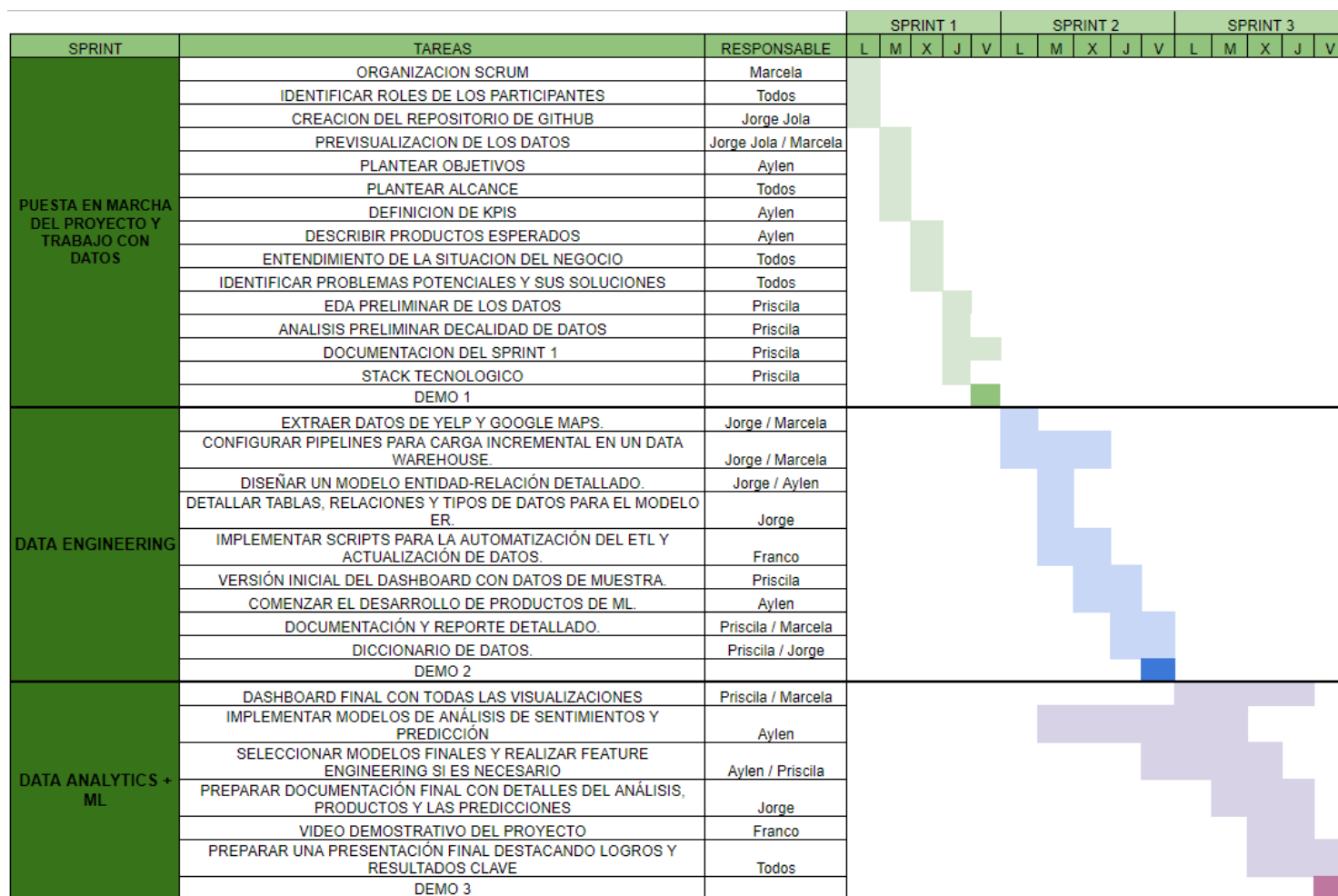


Figura x . Diagrama de Gantt para el proyecto