

ESCUDERO ERICA
 TRABAJO FINAL BOOTCAM DATA ENGINEERING
 DICIEMBRE 2024
 LAB “CLOUD DATAPREP”

QuickLAB

Título: Creating a Data Transformation Pipeline with Cloud Dataprep

Schedule: 1 hour 15 minutes

Cost: 5 Credits

Link:
https://www.cloudskillsboost.google/focuses/4415?catalog_rank=%7B%22rank%22%3A1%2C%22num_filters%22%3A0%2C%22has_search%22%3Atrue%7D&parent=catalog&search_id=32278924

End Lab

00:35:55

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked. [Learn more.](#)

Open Google Cloud console

Username

student-01-72edc93ba934i

Password

Lv2wsKaokzXC

Project ID

qw1klabs-gcp-01-82180e6f

- Click **Create a New Table** from the panel on the right.
- Name your table **revenue_reporting**.
- Select **Drop the Table every run**.
- Click on **Update**.
- Click **RUN**.

Once your Cloud Dataprep job is completed, refresh your BigQuery page and confirm that the output table **revenue_reporting** exists.

Note: If your job fails, try waiting a minute, pressing the back button on your browser, and running the job again with the same settings.

Click **Check my progress** to verify the objective.

✓

Verify if the Cloud Dataprep jobs output the data to BigQuery

Check my progress

Assessment Completed!

Lab instructions and tasks

GSP430

100/100

Overview

Setup and requirements

Task 1. Open Dataprep in the Google Cloud console

Task 2. Creating a BigQuery dataset

Task 3. Connecting BigQuery data to Cloud Dataprep

Task 4. Exploring ecommerce data fields with the UI

Task 5. Cleaning the data

Task 6. Enriching the data

Task 7. Running Cloud Dataprep jobs to BigQuery

Congratulations!

Captura de Pantalla del QuickLab terminado y de la tabla del transform en BQ

The image is a composite screenshot showing two parts of a Google Cloud workflow. The top part is a screenshot of the Google Cloud BigQuery console. The bottom part is a screenshot of a Google Cloud QuickLab completion screen.

Top Screenshot: Google Cloud BigQuery Console

The console shows the 'revenue_reporting' table schema. The table has the following fields:

Field name	Type	Mode	Key	Collation	Default Value	Policy Tags	Description
fullVisitorId	STRING	NULLABLE	-	-	-	-	-
channelGrouping	STRING	NULLABLE	-	-	-	-	-
time	INTEGER	NULLABLE	-	-	-	-	-
country	STRING	NULLABLE	-	-	-	-	-
city	STRING	NULLABLE	-	-	-	-	-
totalTransactionRevenue	INTEGER	NULLABLE	-	-	-	-	-
totalTransactionRevenue1	FLOAT	NULLABLE	-	-	-	-	-
transactions	INTEGER	NULLABLE	-	-	-	-	-
timeOnSite	INTEGER	NULLABLE	-	-	-	-	-
pageviews	INTEGER	NULLABLE	-	-	-	-	-
sessionQualityDim	INTEGER	NULLABLE	-	-	-	-	-
date	DATETIME	NULLABLE	-	-	-	-	-
visitId	INTEGER	NULLABLE	-	-	-	-	-
unique_session_id	STRING	NULLABLE	-	-	-	-	-
type	STRING	NULLABLE	-	-	-	-	-
productRefundAmount	STRING	NULLABLE	-	-	-	-	-
productCategory	INTEGER	NULLABLE	-	-	-	-	-

Bottom Screenshot: Cloud Dataprep Job Completion

The bottom screenshot shows a 'Creating a Data Transformation Pipeline with Cloud Dataprep' lab completion screen. It includes a timer at 00:35:55, a caution message, and a list of tasks completed. The tasks are:

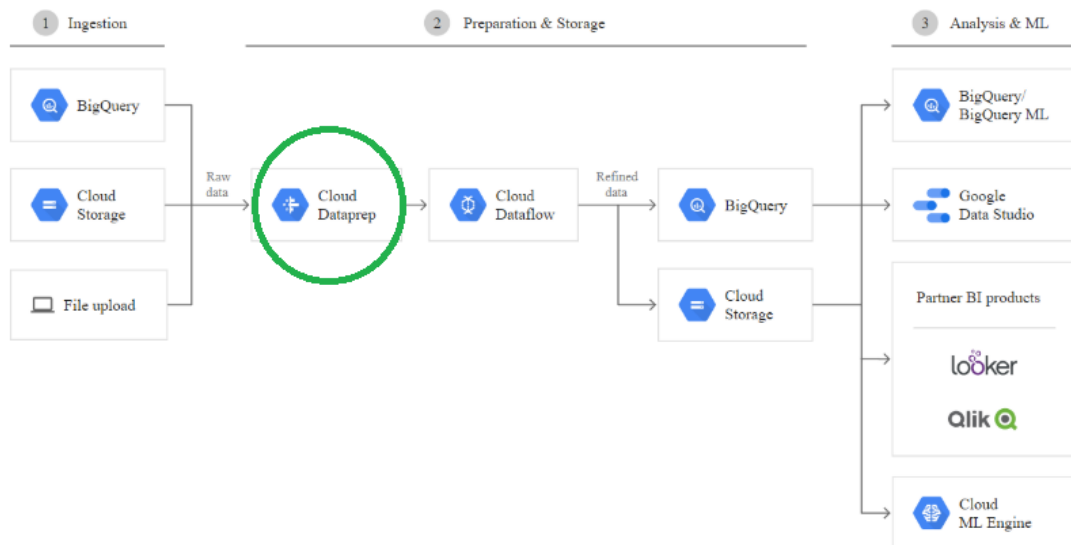
- Open Dataprep in the Google Cloud console
- Creating a BigQuery dataset
- Connecting BigQuery data to Cloud Dataprep
- Exploring ecommerce data fields with the UI
- Cleaning the data
- Enriching the data
- Running Cloud Dataprep jobs to BigQuery

The lab is marked as 'Assessment Completed!' with a green checkmark. The right sidebar shows 'Lab instructions and tasks' with a score of 100/100.

RESPUESTAS PREGUNTAS:

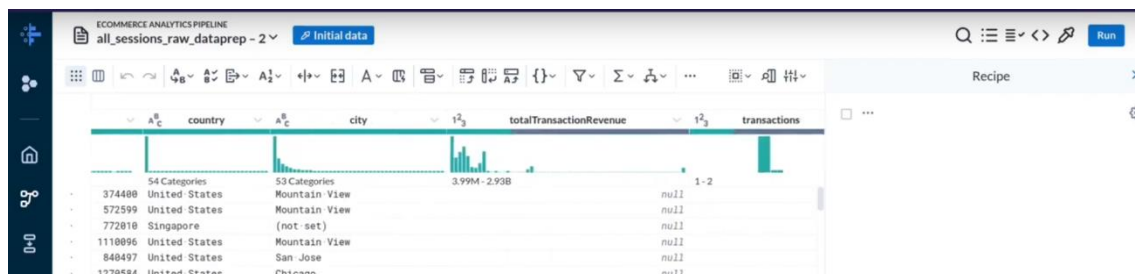
1. ¿Para qué se utiliza data prep?

Es un servicio de datos inteligente de tercero que está integrado con GCP que permite limpiar, explorar y preparar datos de forma gráfica. Permite construir un pipeline de transformación de datos, tomando los datos desde Big Query(BQ) y luego escribe los resultados de esta transformación en Big Query para luego hacer distinta análisis como Data Analytics o Science .



2. ¿Qué cosas se pueden realizar con DataPrep?

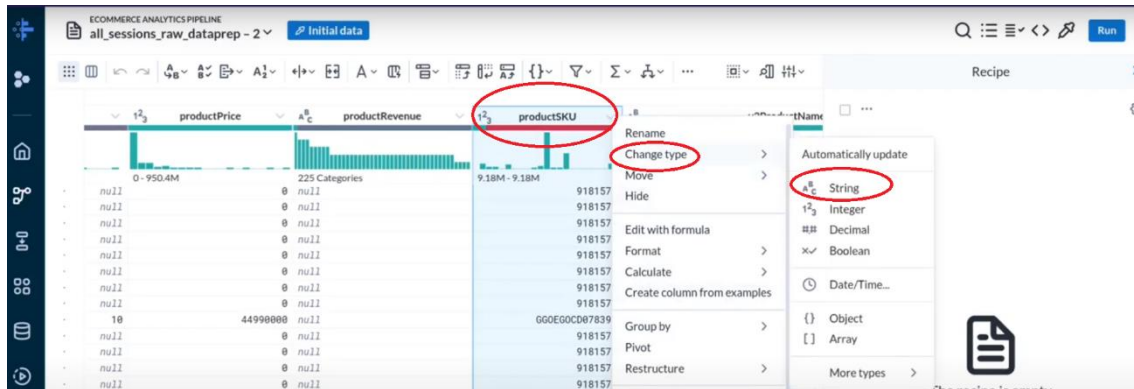
Se conecta a un Dataset en BQ y permite, de forma gráfica, explorar los campos a través de “recetas” pre-existentes mostrando cada columna y sus registros correspondientes. Al ser interactivo, pasando el mouse por arriba de la columna podemos ver la información de cada columna.



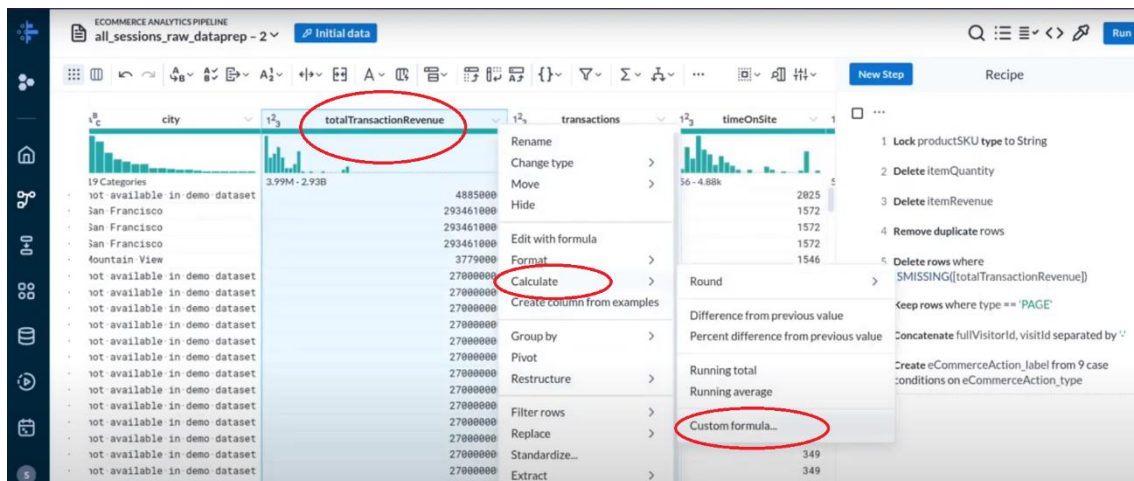
Una vez inspeccionado el dataset, se puede realizar:

- Eliminar columnas
- Cambiar el tipo de datos de las columnas
- Eliminar registros duplicados

- Crear cálculos calculados
- Y realizar distintos filtros sobre la información



Imágenes de DataPrep haciendo para transformaciones sobre los datos



3. ¿Por qué otra/s herramientas lo podrías reemplazar? Por qué?

Hay varias maneras de realizar las transformaciones de los datos, fuera de GCP.

Una es a través de Python o PySpark desarrollando pipelien de transformación a través de la ejecución de DAG's como hemos realizado en los primeros ejercicios de este TP y también se puede realizar modificaciones de datos en Power BI con la opción "transformar datos" dentro de PowerBI unan vez cargados los datos en la herramienta. Esta opción también permite de manera gráfica, hacer transformaciones sobre las columnas y los filas del dataset.

El por qué de utilizar estas opciones, depende mucho de la cantidad de datos y los recursos económicos y tipo de proyecto, si necesita escalar o no, disponibilidad de herramientas para llevar adelante el proyecto, etc.

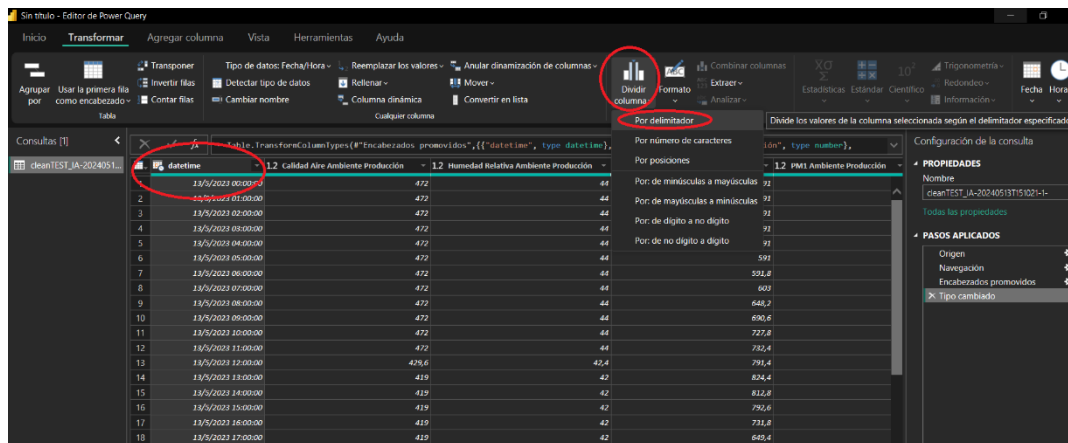


Imagen de PowerBI en el servicio de Transformar datos de power Query

4. ¿Cuáles son los casos de uso comunes de Data Prep de GCP?

Los casos comunes son aquellos negocios de venta on-line o servicios de streaming que manejan una cantidad inmensa de data que debe ser procesada por los diferentes departamentos para obtener información estratégica para la compañía. Venta de productos on-line como ser market places, supermercados, seguros, operadores celulares, video/musica streaming son negocios que generan y manejan mucha data y esta debe estar disponible, segura y fiable 7x24 hs.

5. ¿Cómo se cargan los datos en Data Prep de GCP?

Se debe configurar un dataset y una tabla determinada en BG, cargar los datos y desde dataprep se conecta a BQ-tabla para poder hacer las transformaciones. Luego el resultado de esas transformaciones se vuelve a cargar en BigQuery a través de un JOB de Dataflow creando una nueva tabla en en dataste de BQ.

6. ¿Qué tipos de datos se pueden preparar en Data Prep de GCP?

Google Cloud Dataprep está diseñada para preparar y transformar diversos tipos de datos en múltiples formatos y estructuras, facilitando su análisis y explotación. Es compatible con datos:

- estructurados,
- semiestructurados y
- no estructurados.

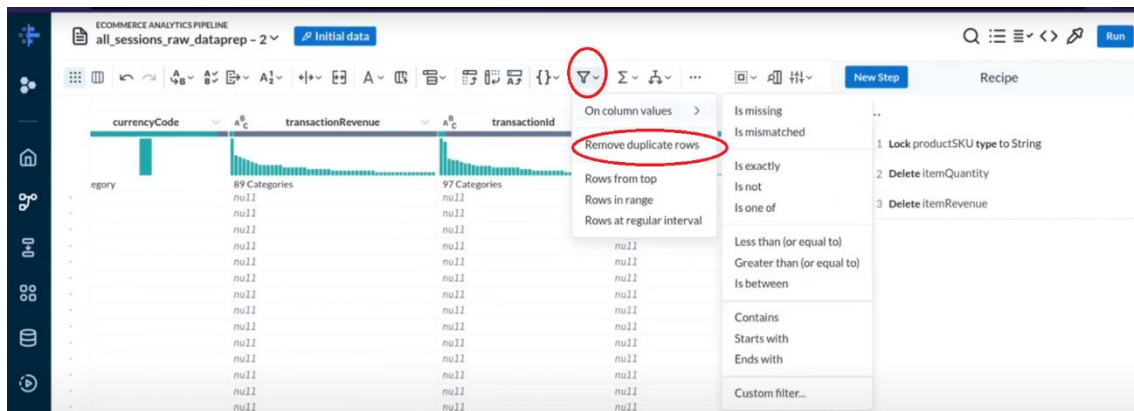
Incluye formatos comunes como CSV, JSON, Avro y Parquet, además de datos almacenados en bases de datos relacionales y sistemas de almacenamiento en la nube como BigQuery y Cloud Storage.

Dataprep maneja datos numéricos, categóricos, fechas, y texto libre, permitiendo transformaciones como la limpieza de valores faltantes, el ajuste de formatos, la eliminación de duplicados, y la corrección de inconsistencias.

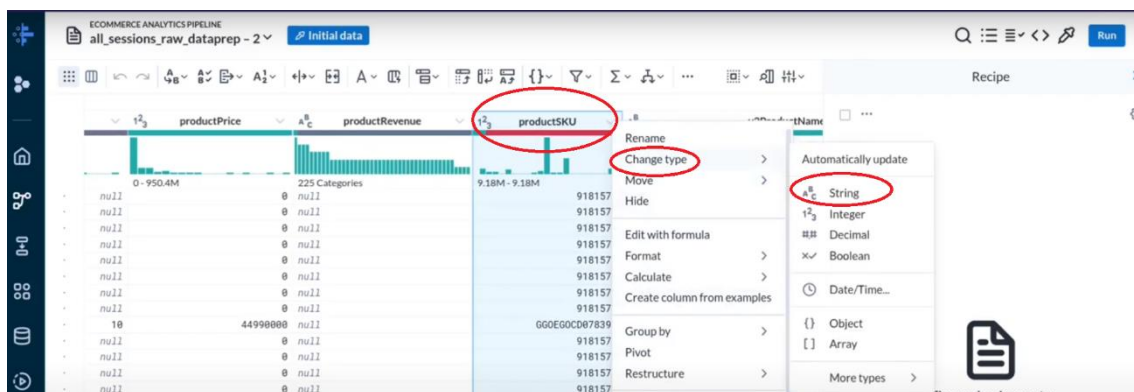
7. ¿Qué pasos se pueden seguir para limpiar y transformar datos en Data Prep de GCP?

Una vez conectado los datos, desde BQ o Cloud Storage, Dataprep nos facilita la visión de los datos y sus estadísticas sobre distribución, frecuencia y anomalías en los datos.

Se seleccionan los datos que sean nulos o duplicados con el ícono de filtrado y seleccionamos filas duplicadas o nulas y se eliminan.



Por ejemplo, para cambiar un tipo de dato a String podemos hacerlo de la siguiente manera:



Se pueden hacer operaciones con los datos de una columna, por ejemplo, el valor de una columna se puede dividir, siempre y cuando el dato sea numérico.

10. ¿Cómo se puede garantizar la calidad de los datos en Data Prep de GCP?

Dataprep permite explorar datos de manera interactiva, identificar problemas comunes como valores faltantes, duplicados, o inconsistencias, y realizar transformaciones mediante una interfaz visual. Permite ver el perfil de los datos a través de gráfico integrado proporciona estadísticas descriptivas y gráficos que destacan patrones y anomalías, lo que nos permite tomar decisiones sobre qué aspectos deben corregirse a través de la definición del pipeline de transformación.

Al ver en cada paso de la transformación como se grafican los gráficos estadísticos se puede ir supervisando el resultado de la transformación y detectar anomalías a corregir. Esto garantiza la calidad final de los datos.

Arquitectura:

El gerente de Analítica te pide realizar una arquitectura hecha en GCP que contemple el uso de esta herramienta ya que le parece muy fácil de usar y una interfaz visual que ayuda a sus desarrolladores ya que no necesitan conocer ningún lenguaje de desarrollo.

Esta arquitectura debería contemplar las siguientes etapas:

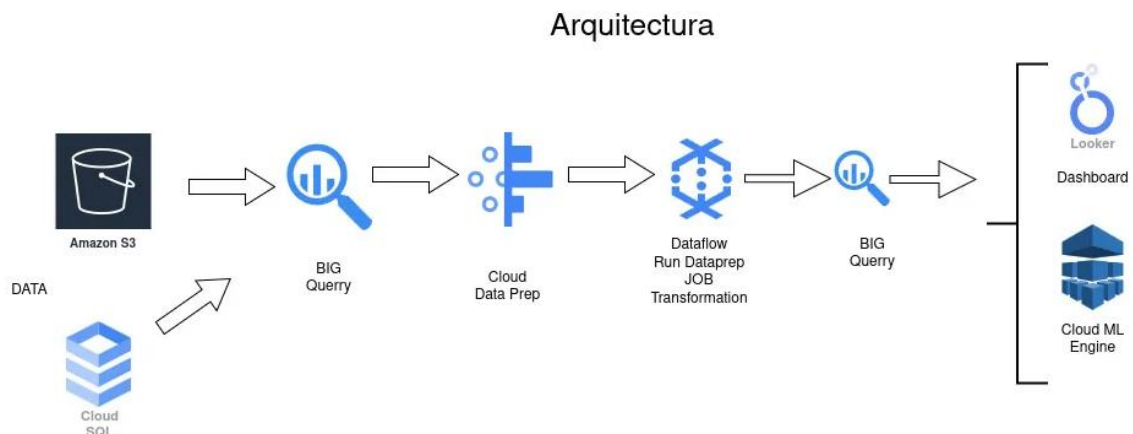
Ingesta: datos parquet almacenados en un bucket de S3 y datos de una aplicación que guarda sus datos en Cloud SQL.

Procesamiento: filtrar, limpiar y procesar datos provenientes de estas fuentes

Almacenar: almacenar los datos procesados en BigQuery

BI: herramientas para visualizar la información almacenada en el Data Warehouse

ML: Herramienta para construir un modelo de regresión lineal con la información almacenada en el Data Warehouse



Descripción de la solución propuesta:

La arquitectura diseñada para análisis de datos en GCP es una solución modular, escalable y eficiente que integra múltiples herramientas y servicios para procesar datos de diversas fuentes y obtener insights avanzados para el área que ha solicitado esta arquitectura.

El flujo comienza con la ingesta de datos desde fuentes externas como **Amazon S3** y **Cloud SQL**, que se consolidan en **BigQuery** como almacenamiento centralizado de data.

En **BigQuery**, los datos se almacenan en su forma bruta, donde están disponibles para consultas iniciales y preparación. Posteriormente, **Cloud Dataprep** se utiliza para realizar transformaciones visuales e interactivas, como limpieza de datos, cambios de formatos, creación de nuevas columnas calculadas según los requerimientos del negocio y creación de nuevas estructuras. Las reglas definidas en Dataprep se implementan mediante **Dataflow**, lo que permite procesar grandes volúmenes de datos de forma automatizada y escalable, cargando los datos transformados nuevamente en BigQuery.

Una vez que los datos están preparados, se guardan en nuevas tablas en BigQuery y permite habilitar consultas analíticas avanzadas y prepara los datos para ser consumidos por herramientas como **Looker**, que proporciona visualizaciones interactivas y dashboards dinámicos.

Adicionalmente, los datos transformados pueden alimentar **Cloud ML Engine**, donde se desarrollan y despliegan modelos de machine learning para análisis predictivo y automatización inteligente.

Esta arquitectura garantiza:

- **Eficiencia** en el manejo de grandes volúmenes de datos mediante BigQuery y Dataflow.
- **Flexibilidad** al integrar múltiples fuentes de datos externas e internas.
- **Automatización** del proceso de transformación con pipelines repetibles y confiables.
- **Escalabilidad** para soportar necesidades de datos crecientes y análisis complejos.

Es una propuesta de solución técnica robusta diseñada para optimizar el flujo de datos y habilitar análisis para la obtención de información relevante para la empresa, permitiendo a los equipos técnicos manejar datos de extremo a extremo de manera eficiente.