

Documentación Ingeniería de Datos

Extracción y carga de datos

El objetivo era extraer los datos de la pagina NYC Taxi & Limousine Commission, del periodo de 2018 hasta 2023 para Yellow taxi, Green taxi y For-Hire Vehicle para Green y Yellow taxi 2018 a 2023 y para For-Hire Vehicle los años 2022 a 2023



Para lograr este objetivo decidimos usar un enfoque de web scraping utilizando la biblioteca de Python BeautifulSoup para extraer los nuevos enlaces disponibles en la pagina creamos el siguiente código:

```
# Configura la variable de entorno GOOGLE
1 # Configura la variable de entorno GOOGLE_APPLICATION_CREDENTIALS con la ruta al archivo de credenciales JSON
2 os.environ["GOOGLE_APPLICATION_CREDENTIALS"] = "/Workspace/Users/ccfontanillar@unadvirtual.edu.co/project-test-419903-c2e3452605db.json"
3
4 def get_filtered_urls(url):
5     """
6     Obtiene las URLs filtradas de la página web especificada.
7
8     Args:
9     url (str): URL de la página web.
10
11     Returns:
12     list: Lista de URLs filtradas.
13     """
14     filtered_urls = []
15
16     response = requests.get(url)
17     if response.status_code == 200:
18         soup = BeautifulSoup(response.text, 'html.parser')
19         links = soup.find_all('a', href=True)
20
21         for link in links:
22             href = link.get('href')
23             absolute_url = urljoin(url, href)
24             # Filtrar los enlaces por los que contienen "yellow_tripdata", "green_tripdata" o "fhvhv_tripdata" y el año 2024
25             if any(keyword in absolute_url for keyword in ["yellow_tripdata", "green_tripdata", "fhvhv_tripdata"]) and "2024" in absolute_url:
26                 filtered_urls.append(absolute_url)
27
28     return filtered_urls
29
30 def process_and_upload_trip_data(filtered_urls):
31     """
32     Lee cada URL en un DataFrame y lo convierte en formato Parquet antes de subirlo al Cloud Storage.
33
34     Args:
```

```

30 def process_and_upload_trip_data(filtered_urls):
31     """
32     Lee cada URL en un DataFrame y lo convierte en formato Parquet antes de subirlo al Cloud Storage.
33
34     Args:
35         filtered_urls (list): Lista de URLs filtradas.
36
37     Returns:
38         None
39     """
40     try:
41         # Crea un cliente de Cloud Storage
42         storage_client = storage.Client()
43
44         # Nombre de tu bucket en GCS
45         bucket_name = 'upload-scraping-gcp'
46
47         # Obtiene el bucket
48         bucket = storage_client.get_bucket(bucket_name)
49
50         # Recorre las URLs filtradas, lee cada una en un DataFrame y la guarda en un archivo Parquet en GCS
51         for url in filtered_urls:
52             try:
53                 # Leer el archivo Parquet en un DataFrame de Pandas
54                 df = pd.read_parquet(url)
55
56                 # Generar un nombre de archivo único para el Parquet
57                 parquet_file_name = os.path.basename(url)
58
59                 # Subir el archivo Parquet a GCS
60                 blob = bucket.blob(f"import/raw/{parquet_file_name}")
61                 blob.upload_from_string(df.to_parquet(), content_type='application/octet-stream')
62                 print(f"Archivo Parquet subido a Cloud Storage: gs://{bucket_name}/import/raw/{parquet_file_name}")
63
64             except Exception as e:
65                 print(f"Error al procesar {url}: {e}")
66
67     except Exception as e:
68         print(f"Error en el proceso de carga a Cloud Storage: {e}")
69

```

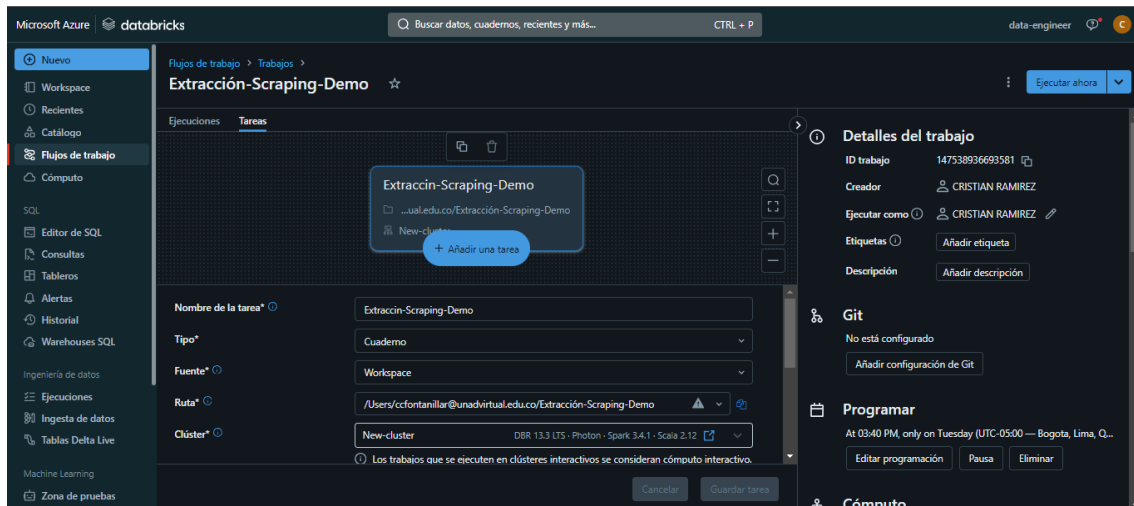
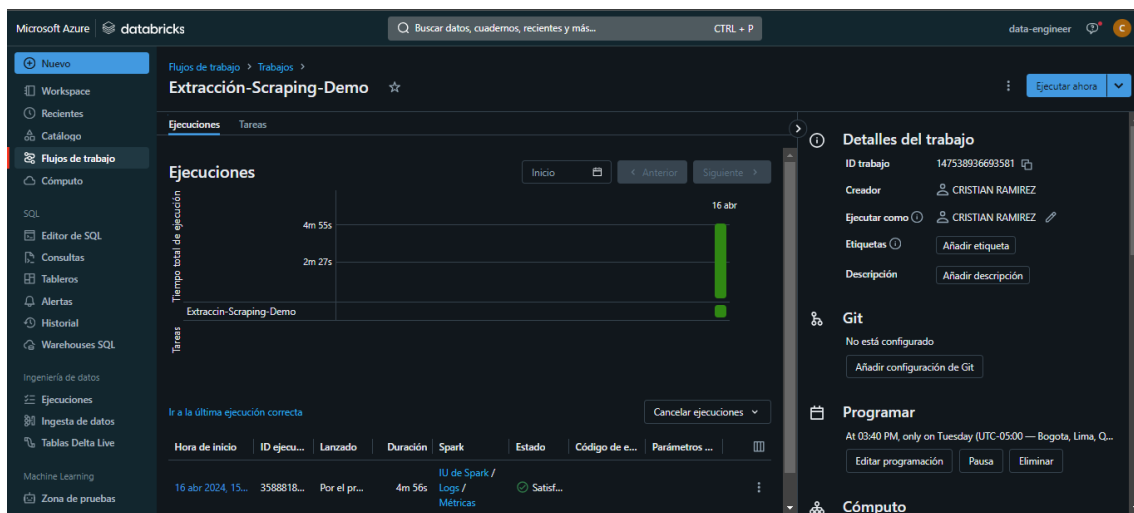
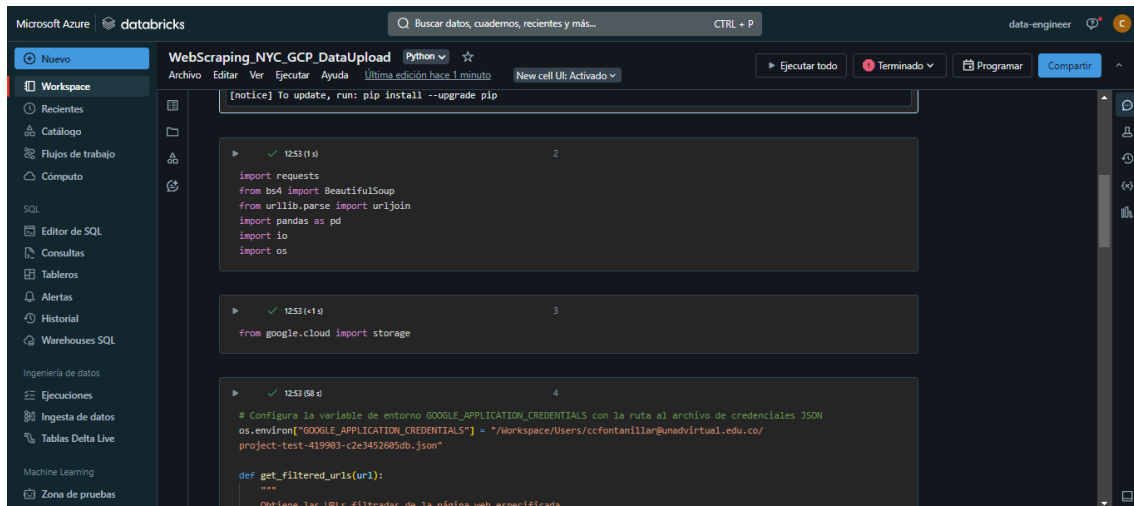
Este código consiste en dos funciones principales: `get_filtered_urls(url)` y `process_and_upload_trip_data(filtered_urls)`.

La función `get_filtered_urls(url)` toma una URL como entrada y utiliza bibliotecas como `requests` y `BeautifulSoup` para obtener el contenido HTML de la página web y analizarlo en busca de enlaces. Luego, filtra estos enlaces según ciertos criterios, como la presencia de palabras clave específicas y el año 2024 en la URL, y devuelve una lista de URLs que cumplen con estos criterios.

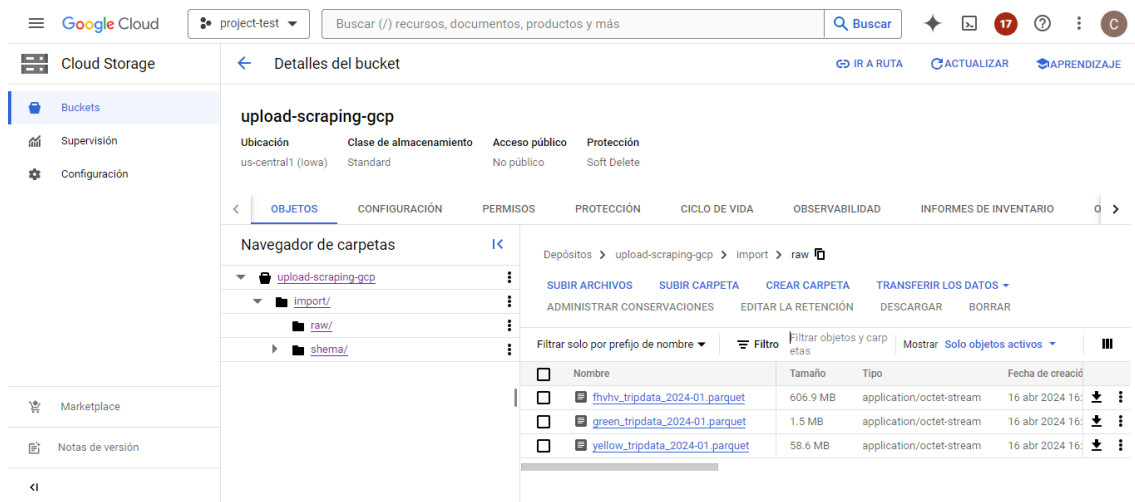
La función `process_and_upload_trip_data(filtered_urls)` toma la lista de URLs filtradas como entrada. Itera sobre cada URL, lee su contenido como un DataFrame de Pandas (asumiendo que son archivos Parquet), y luego los sube al Google Cloud Storage en una ubicación específica dentro de un bucket especificado. Esta función maneja excepciones para capturar cualquier error que pueda ocurrir durante el proceso y muestra mensajes de error apropiados.

Además de estas funciones, el código incluye un bloque principal que define la URL del sitio web del cual se desean extraer los datos. Luego, utiliza las funciones mencionadas anteriormente para realizar el proceso completo de obtención, procesamiento y carga de datos.

Después de crear la lógica creamos un Workspace en Databricks donde creamos un nuevo notebook y se hizo las configuraciones necesarias para la producción del código.



Tarea programada



Archivos subidos a GCS

Justificación del uso de la nube

La decisión de utilizar los servicios en la nube de Google Cloud Platform (GCP), específicamente Google Cloud Storage, Cloud Functions y BigQuery, dentro de nuestro proyecto de Data Science ha sido fundamentada en una serie de consideraciones estratégicas y prácticas.

En primer lugar, la flexibilidad y escalabilidad inherentes a la nube son aspectos importantes para un proyecto de esta naturaleza. Al optar por servicios en la nube, podemos adaptar fácilmente nuestros recursos según las demandas cambiantes del proyecto, sin incurrir en costos adicionales significativos por infraestructura subutilizada o la necesidad de adquirir hardware adicional. Esta flexibilidad nos permite concentrar nuestros esfuerzos en la eficiencia y la optimización de los procesos de ingeniería de datos, en lugar de preocuparnos por la gestión y mantenimiento de la infraestructura subyacente.

Aprovechamos el servicio de Google Cloud Storage como almacenamiento temporal de datos no estructurados y crudos para después aplicar un preprocesamiento que permitiera aplicar ETL lo mismo como destino de los datos de scraping.

Además, el uso de Cloud Functions nos permite automatizar el procesamiento de datos en tiempo real de manera eficiente. Al recibir archivos Parquet en nuestro almacenamiento, podemos activar eventos que desencadenen procesos de transformación (ETL) y carga de datos en BigQuery de manera ágil y eficaz. Esto garantiza una actualización constante de nuestra base de datos, lo que es fundamental para la toma de decisiones informada y el análisis en tiempo real.

Por otro lado, la elección de BigQuery como nuestro almacén de datos principal se debe a su capacidad para gestionar conjuntos de datos a gran escala con un rendimiento excepcional. Con la capacidad de crear y manipular dataframes como el proyecto_taxis, que integra diversas tablas relacionales como climas, economía de autos, modelos aprobados, entre otros, podemos realizar consultas complejas de manera eficiente y

rápida. Esto es fundamental para la generación de insights y la construcción de modelos de machine learning, ya que nos permite acceder y analizar grandes volúmenes de datos de forma oportuna y efectiva.

Además, al aprovechar las herramientas y servicios integrados de GCP, podemos construir pipelines de datos robustos y confiables que automatizan el flujo de datos desde su origen hasta su destino final. Esto no solo aumenta la eficiencia operativa, sino que también garantiza la integridad y la calidad de los datos, lo que es crucial para la generación de informes precisos y la toma de decisiones fundamentadas.