

Taller: Flujo de datos de datos SQL (Oracle) a NoSQL (MongoDB)

Objetivo:

Al finalizar este taller, los participantes serán capaces de realizar un flujo de datos de Oracle a MongoDB utilizando un programa de Google Colab.

Requisitos:

- Cuenta de Google
- Conocimientos básicos de Python
- Una base de datos de Oracle
- Una base de datos de MongoDB Atlas

Metodología:

- En grupos de trabajo, realicen uno de los cuatro ejercicios asignados.
- Construya la solución basado en el código entregado.

Ejercicios:

Ejercicio 1: CSV a MongoDB

Desde Python con el siguiente código se sube la base seleccionada Noticias:

Se genera la conexión:

```
1 !pip install pymongo
2 !pip install --upgrade pymongo
3
Requirement already satisfied: pymongo in /usr/local/lib/python3.10/dist-packages (4.6.0)
Requirement already satisfied: dnspython<3.0.0,>=1.16.0 in /usr/local/lib/python3.10/dist-packages (from pymongo) (2.4.2)
Requirement already satisfied: pymongo in /usr/local/lib/python3.10/dist-packages (4.6.0)
Requirement already satisfied: dnspython<3.0.0,>=1.16.0 in /usr/local/lib/python3.10/dist-packages (from pymongo) (2.4.2)

1 from pymongo.mongo_client import MongoClient
2 from pymongo.server_api import ServerApi
3
4 uri = "mongodb+srv://Adri_30:Melek@cluster0.etcfqhm.mongodb.net/?retryWrites=true&w=majority"
5
6 # Create a new client and connect to the server
7 client = MongoClient(uri, server_api=ServerApi('1'))
8
9 # Send a ping to confirm a successful connection
10 try:
11     client.admin.command('ping')
12     print("Pinged your deployment. You successfully connected to MongoDB!")
13 except Exception as e:
14     print(e)

Pinged your deployment. You successfully connected to MongoDB!
```

Se genera la base de datos.csv y se sube a python:

✓

[47]

```
1 # Seleccionar una base de datos especifica
2 db = client["Noticias_"]
3
4 # Ahora puedes usar 'db' para referirte a tu base de datos
5 coleccion = db["Noticias_work"]
6
```

✓

01

```
1 # Intentando determinar el delimitador correcto del archivo CSV
2 import csv
3 import pandas as pd
4
5 # Ruta al archivo CSV
6 ruta_archivo_csv = '/content/Noticias_bases.csv'
7
8 # Cargar el archivo CSV con el delimitador detectado
9 df = pd.read_csv(ruta_archivo_csv, delimiter=';')
10
11 # Mostrar las primeras filas del DataFrame
12 df.head()
13
```

	Unnamed: 0	Unnamed: 1	Unnamed: 2	Unnamed: 3	Unnamed: 4	Unnamed: 5	Unnamed: 6	Unnamed: 7	Unnamed: 8	Unnamed: 9	...	Unnamed: 16	Unnamed: 17	Unnamed: 18	Unnamed: 19
0	NaN	Observatorio de Noticias en Tecnologia	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN
1	NaN	#	Fecha Actualización	Fecha de Noticia	Nombre Estudiante	Link Noticia	Titulo	Resumen (No mayor a 50 palabras)	Categorias (Maximo 4 palabras claves que descr...	NaN	...	NaN	NaN	NaN	NaN
2	NaN	1	17/11/2023	10/5/2023	Maria Camila Prada	https://news.un.org/es/story/2023/05/1520892	Las tecnologías digitales son una herramienta ...	El artículo de la ONU informa sobre la crisis ...	Crisis humanitaria en YemenIn Falta de fondos...	NaN	...	NaN	NaN	NaN	NaN

Este código delimita el archivo CSV, se genera automáticamente antes de cargar los datos en un DataFrame.

✓

01

```
1 # Leer las primeras líneas del archivo para determinar el delimitador
2 with open(ruta_archivo_csv, 'r', encoding='utf-8') as archivo:
3     muestra = archivo.read(1024)
4
5 delimitador = csv.Sniffer().sniff(muestra).delimiter
6
7 # Cargar el archivo CSV con el delimitador detectado
8 df = pd.read_csv(ruta_archivo_csv, delimiter=delimitador)
9
10 # Mostrar las primeras filas del DataFrame
11 df.head()
12
```

	Unnamed: 0	Unnamed: 1	Unnamed: 2	Unnamed: 3	Unnamed: 4	Unnamed: 5	Unnamed: 6	Unnamed: 7	Unnamed: 8	Unnamed: 9	...	Unnamed: 16	Unnamed: 17	Unnamed: 18	Unnamed: 19
0	NaN	Observatorio de Noticias en Tecnologia	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN
1	NaN	#	Fecha Actualización	Fecha de Noticia	Nombre Estudiante	Link Noticia	Titulo	Resumen (No mayor a 50 palabras)	Categorias (Maximo 4 palabras claves que descr...	NaN	...	NaN	NaN	NaN	NaN
2	NaN	1	17/11/2023	10/5/2023	Maria Camila Prada	https://news.un.org/es/story/2023/05/1520892	Las tecnologías digitales son una herramienta ...	El artículo de la ONU informa sobre la crisis ...	Crisis humanitaria en YemenIn Falta de fondos...	NaN	...	NaN	NaN	NaN	NaN
3	NaN	2	24/08/2023	05-06-23	Alexander Victoria Garcia	https://elpais.com/tecnologia/2023-06-06/ibm-a...	IBM anuncia su primer centro de datos cuántico...	La computación cuántica encuentra viabilidad a...	IBM, Computacion Cuantica, red	NaN	...	NaN	NaN	NaN	NaN
4	NaN	3	05/09/2023	28/06/2023	Ana Maria Cruz Pacheco	https://www.bbc.com/news/science-environment-5...	What is net zero and how are the UK and other ...	YnEl Reino Unido podría quedarse rezagado en	Captura y almacenamiento de carbono, gases de ...	NaN	...	NaN	NaN	NaN	NaN

Se genera visualización del data frame:

```
1 # Convertir el DataFrame de pandas a una lista de diccionarios
2 datos_para_insertar = df.to_dict(orient='records')
3
4 # Ahora puedes usar 'insert_many' con esta lista de diccionarios
5 coleccion.insert_many(datos_para_insertar)
6
```

```
InsertManyResult([ObjectId('65637afc377dc27eb0e2cbe9'), ObjectId('65637afc377dc27eb0e2cbea'), ObjectId('65637afc377dc27eb0e2cbeb'), ObjectId('65637afc377dc27eb0e2cbec'),
ObjectId('65637afc377dc27eb0e2cbef'), ObjectId('65637afc377dc27eb0e2cbf0'), ObjectId('65637afc377dc27eb0e2cbf1'),
ObjectId('65637afc377dc27eb0e2cbf2'), ObjectId('65637afc377dc27eb0e2cbf3'), ObjectId('65637afc377dc27eb0e2cbf4'), ObjectId('65637afc377dc27eb0e2cbf5'), ObjectId('65637afc377dc27eb0e2cbf6'),
ObjectId('65637afc377dc27eb0e2cbf7'), ObjectId('65637afc377dc27eb0e2cbf8'), ObjectId('65637afc377dc27eb0e2cbf9'), ObjectId('65637afc377dc27eb0e2cbfa'), ObjectId('65637afc377dc27eb0e2cbfb'),
ObjectId('65637afc377dc27eb0e2cbfc'), ObjectId('65637afc377dc27eb0e2cbfd'), ObjectId('65637afc377dc27eb0e2cbfe'), ObjectId('65637afc377dc27eb0e2cbff'), ObjectId('65637afc377dc27eb0e2cc00'),
ObjectId('65637afc377dc27eb0e2cc01'), ObjectId('65637afc377dc27eb0e2cc02'), ObjectId('65637afc377dc27eb0e2cc03'), ObjectId('65637afc377dc27eb0e2cc04'), ObjectId('65637afc377dc27eb0e2cc05'),
ObjectId('65637afc377dc27eb0e2cc06'), ObjectId('65637afc377dc27eb0e2cc07'), ObjectId('65637afc377dc27eb0e2cc08'), ObjectId('65637afc377dc27eb0e2cc09'), ObjectId('65637afc377dc27eb0e2cc0a'),
ObjectId('65637afc377dc27eb0e2cc0b'), ObjectId('65637afc377dc27eb0e2cc0c'), ObjectId('65637afc377dc27eb0e2cc0d'), ObjectId('65637afc377dc27eb0e2cc0e'), ObjectId('65637afc377dc27eb0e2cc0f'),
ObjectId('65637afc377dc27eb0e2cc10'), ObjectId('65637afc377dc27eb0e2cc11'), ObjectId('65637afc377dc27eb0e2cc12'), ObjectId('65637afc377dc27eb0e2cc13'), ObjectId('65637afc377dc27eb0e2cc14'),
ObjectId('65637afc377dc27eb0e2cc15'), ObjectId('65637afc377dc27eb0e2cc16'), ObjectId('65637afc377dc27eb0e2cc17'), ObjectId('65637afc377dc27eb0e2cc18'), ObjectId('65637afc377dc27eb0e2cc19'),
ObjectId('65637afc377dc27eb0e2cc1a'), ObjectId('65637afc377dc27eb0e2cc1b'), ObjectId('65637afc377dc27eb0e2cc1c'), ObjectId('65637afc377dc27eb0e2cc1d'), ObjectId('65637afc377dc27eb0e2cc1e'),
ObjectId('65637afc377dc27eb0e2cc1f'), ObjectId('65637afc377dc27eb0e2cc20'), ObjectId('65637afc377dc27eb0e2cc21'), ObjectId('65637afc377dc27eb0e2cc22'), ObjectId('65637afc377dc27eb0e2cc23'),
ObjectId('65637afc377dc27eb0e2cc24'), ObjectId('65637afc377dc27eb0e2cc25'), ObjectId('65637afc377dc27eb0e2cc26'), ObjectId('65637afc377dc27eb0e2cc27'), ObjectId('65637afc377dc27eb0e2cc28'),
ObjectId('65637afc377dc27eb0e2cc29'), ObjectId('65637afc377dc27eb0e2cc2a'), ObjectId('65637afc377dc27eb0e2cc2b'), ObjectId('65637afc377dc27eb0e2cc2c'), ObjectId('65637afc377dc27eb0e2cc2d'),
ObjectId('65637afc377dc27eb0e2cc2e'), ObjectId('65637afc377dc27eb0e2cc2f'), ObjectId('65637afc377dc27eb0e2cc30'), ObjectId('65637afc377dc27eb0e2cc31'), ObjectId('65637afc377dc27eb0e2cc32'),
ObjectId('65637afc377dc27eb0e2cc33'), ObjectId('65637afc377dc27eb0e2cc34'), ObjectId('65637afc377dc27eb0e2cc35'), ObjectId('65637afc377dc27eb0e2cc36'), ObjectId('65637afc377dc27eb0e2cc37'),
ObjectId('65637afc377dc27eb0e2cc38'), ObjectId('65637afc377dc27eb0e2cc39'), ObjectId('65637afc377dc27eb0e2cc40'), ObjectId('65637afc377dc27eb0e2cc41'),
ObjectId('65637afc377dc27eb0e2cc42'), ObjectId('65637afc377dc27eb0e2cc43'), ObjectId('65637afc377dc27eb0e2cc44'), ObjectId('65637afc377dc27eb0e2cc45'), ObjectId('65637afc377dc27eb0e2cc46'),
ObjectId('65637afc377dc27eb0e2cc47'), ObjectId('65637afc377dc27eb0e2cc48'), ObjectId('65637afc377dc27eb0e2cc49'), ObjectId('65637afc377dc27eb0e2cc4a'), ObjectId('65637afc377dc27eb0e2cc4b'),
ObjectId('65637afc377dc27eb0e2cc4c'), ObjectId('65637afc377dc27eb0e2cc4d'), ObjectId('65637afc377dc27eb0e2cc4e'), ObjectId('65637afc377dc27eb0e2cc4f'), ObjectId('65637afc377dc27eb0e2cc50'),
ObjectId('65637afc377dc27eb0e2cc51'), ObjectId('65637afc377dc27eb0e2cc52'), ObjectId('65637afc377dc27eb0e2cc53'), ObjectId('65637afc377dc27eb0e2cc54'), ObjectId('65637afc377dc27eb0e2cc55'),
ObjectId('65637afc377dc27eb0e2cc56'), ObjectId('65637afc377dc27eb0e2cc57'), ObjectId('65637afc377dc27eb0e2cc58'), ObjectId('65637afc377dc27eb0e2cc59'), ObjectId('65637afc377dc27eb0e2cc60'),
ObjectId('65637afc377dc27eb0e2cc61'), ObjectId('65637afc377dc27eb0e2cc62'), ObjectId('65637afc377dc27eb0e2cc63'), ObjectId('65637afc377dc27eb0e2cc64'),
ObjectId('65637afc377dc27eb0e2cc65'), ObjectId('65637afc377dc27eb0e2cc66'), ObjectId('65637afc377dc27eb0e2cc67'), ObjectId('65637afc377dc27eb0e2cc68'), ObjectId('65637afc377dc27eb0e2cc69'),
ObjectId('65637afc377dc27eb0e2cc6a'), ObjectId('65637afc377dc27eb0e2cc6b'), ObjectId('65637afc377dc27eb0e2cc6c'), ObjectId('65637afc377dc27eb0e2cc6d'), ObjectId('65637afc377dc27eb0e2cc6e'),
ObjectId('65637afc377dc27eb0e2cc6f'), ObjectId('65637afc377dc27eb0e2cc70'), ObjectId('65637afc377dc27eb0e2cc71'), ObjectId('65637afc377dc27eb0e2cc72'), ObjectId('65637afc377dc27eb0e2cc73'),
ObjectId('65637afc377dc27eb0e2cc74'), ObjectId('65637afc377dc27eb0e2cc75'), ObjectId('65637afc377dc27eb0e2cc76'), ObjectId('65637afc377dc27eb0e2cc77'), ObjectId('65637afc377dc27eb0e2cc78'),
ObjectId('65637afc377dc27eb0e2cc79'), ObjectId('65637afc377dc27eb0e2cc7a'), ObjectId('65637afc377dc27eb0e2cc7b'), ObjectId('65637afc377dc27eb0e2cc7c'), ObjectId('65637afc377dc27eb0e2cc7d'),
ObjectId('65637afc377dc27eb0e2cc7e'), ObjectId('65637afc377dc27eb0e2cc7f'), ObjectId('65637afc377dc27eb0e2cc80'), ObjectId('65637afc377dc27eb0e2cc81'), ObjectId('65637afc377dc27eb0e2cc82'),
ObjectId('65637afc377dc27eb0e2cc83'), ObjectId('65637afc377dc27eb0e2cc84'), ObjectId('65637afc377dc27eb0e2cc85'), ObjectId('65637afc377dc27eb0e2cc86'), ObjectId('65637afc377dc27eb0e2cc87'),
ObjectId('65637afc377dc27eb0e2cc88'), ObjectId('65637afc377dc27eb0e2cc89'), ObjectId('65637afc377dc27eb0e2cc8a'), ObjectId('65637afc377dc27eb0e2cc8b'), ObjectId('65637afc377dc27eb0e2cc8c'),
```

Se verifica en Mongo, y se muestra la data cargada:

Atlas

Adriana's Or...

Access Manager

Billing

All ClustersGet Help +Adriana

Project 0

Data Services

App Services

Charts

Overview

Cluster0

6.0.11

AWS N. Virginia (us-east-1)

Deployment

Database

Data Lake

Services

Device Sync

Triggers

Data API

Data Federation

Search

Stream Processing

Security

Quickstart

Backup

Database Access

Network Access

Advanced

Goto

Overview

Real Time

Metrics

Collections

Search

Profiler

Performance Advisor

Online Archive

Cmd Line Tools

DATABASES: 11

COLLECTIONS: 25

Visualize Your Data

Refresh

Create Database

Search Namespaces

Noticies..

Noticies_work

Find

Indexes

Schema Anti-Patterns

Aggregation

Search Indexes

Filter

Type a query: { field: 'value' }

Reset

Apply

Options

QUERY RESULTS: 1-20 OF MANY

_id

ObjectId('65637afc377dc27eb0e2cbe9')

Unnamed: 0

NaN

Unnamed: 1

"Observatorio de Noticias en Tecnologia"

Unnamed: 2

NaN

Unnamed: 3

NaN

Unnamed: 4

NaN

Unnamed: 5

NaN

Unnamed: 6

NaN

Unnamed: 7

NaN

Unnamed: 8

NaN

Unnamed: 9

NaN

Unnamed: 10

NaN

Unnamed: 11

NaN

Unnamed: 12

NaN

Unnamed: 13

NaN

Unnamed: 14

NaN