

Caso Técnico – Ingeniería de Datos Comercial

Implementación de Flujo ETL para Procesamiento de Entregas.

José Carlos Porras.

[Repositorio-caso-ingeneria-de-datos](#)

Objetivo

Desarrollar un proceso robusto y parametrizable para la ingesta, transformación y carga de datos históricos de entregas.

Alcance

El proyecto abarca desde la extracción de archivos planos (CSV) hasta la generación de una capa "Silver" confiable en formatos Parquet y CSV. La solución implementa una **arquitectura de configuración centralizada basada en OmegaConf y YAML**, permitiendo la parametrización dinámica de reglas de negocio (como rangos de fechas, filtros de país y factores de conversión) sin necesidad de modificar el código fuente.

Objetivo del Caso

Diseñar e implementar un flujo ETL que permita:

- Procesar datos históricos no cargados previamente.
- Filtrar información por rango de fechas y país.
- Normalizar unidades de medida.
- Clasificar tipos de entrega.
- Garantizar calidad, consistencia y estandarización de los datos.
- Generar salidas particionadas por fecha de proceso.

Flujo de Trabajo

1. Inicialización y Parametrización

- El sistema arranca leyendo el archivo config.yaml.
- Se cargan las reglas críticas: rangos de fechas, país objetivo y factores de conversión. Si el archivo de configuración no existe, el proceso se detiene.

2. Ingesta Inteligente (Extract)

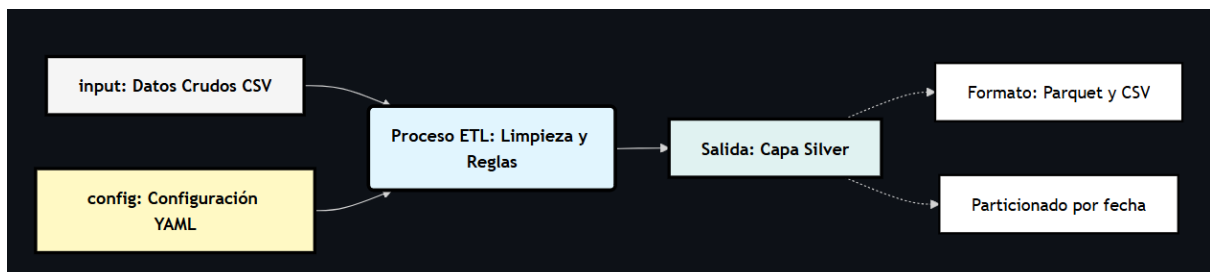
- Se conecta al archivo CSV definido.
- Se aplica un **filtrado temprano** (Push-down predicate) usando los parámetros de fecha y país. Esto asegura que solo entren al sistema los datos necesarios, optimizando el uso de memoria RAM.

3. Transformación y Calidad (Transform)

- **Normalización:** Se convierten todas las cantidades de "Cajas" a "Unidades" estándar.
- **Enriquecimiento:** Se generan columnas lógicas (entrega_rutina, entrega_bonificacion) basadas en los códigos de entrega.
- **Limpieza:** Se eliminan duplicados y se validan los tipos de datos (precios decimales, cantidades positivas).

4. Carga y Persistencia (Load)

- Se estandarizan los nombres de columnas al inglés.
- Se utiliza una estrategia híbrida con **Pandas** para escribir los resultados en disco local (Windows).
- **Resultado Final:** Se generan archivos particionados por fecha en formato **Parquet** (para consumo analítico eficiente) y se crean réplicas en **CSV** (para auditoría manual).



Parametrización del Flujo YAML

Todos los parámetros del flujo son controlados desde un archivo YAML:

- **Ruta de entrada:** Se estableció la dirección donde se debe buscar el archivo CSV para empezar a trabajar.
- **Rango de fechas:** Define el periodo de tiempo que queremos analizar.
- **País:** Se filtro para quedarnos solo con la información de una región específica (por ejemplo, "GT" para Guatemala), ignorando el resto.
- **Conversión de unidades** Se definió el valor **20** por el cual se deben multiplicar las cantidades marcadas como "**Cajas**" (**CS**) para estandarizarlas a "**Unidades**" (**ST**).
- **Tipos de entrega** Son las etiquetas que usa la empresa para diferenciar qué ventas son normales (Rutina) o (Bonificaciones).

```
1 # config/config.yaml
2 input:
3   file_path: "../data/input/data_entrega_productos.csv"
4
5 filters:
6   start_date: 20250101
7   end_date: 20250630
8   country: "GT"
9   date_format: "yyyyMMdd"
10
11 processing:
12   box_size: 20
13   delivery_types:
14     routine: ["ZPRE", "ZVE1"]
15     bonus: ["Z04", "Z05"]
16
17 output:
18   base_path: "../data/processed"
19   format: "parquet"
20   mode: "overwrite"
21   partition_col: "fecha_proceso"
22
```

Gestión de Configuración y Parametrización

Este bloque de código representa el **punto de entrada**. Utilizamos la librería **OmegaConf** para los parámetros de la lógica de procesamiento del archivo **config.yaml**.

1. **Definición de Ruta Relativa - config.yaml** usando una ruta relativa (**../**) para garantizar que el código funcione independientemente de la carpeta base del usuario.
2. **Carga de Objeto config.yaml**

```
1 import os
2 from omegaconf import OmegaConf
3
4 #YAML
5 config_path = "../config/config.yaml"
6
7 if not os.path.exists(config_path):
8     raise FileNotFoundError(f"no se encontro el file: {config_path}")
9
10 config = OmegaConf.load(config_path)
11 print(f"Archivo yaml leido correctamente")
```

Librerías utilizadas fueron:

Librería	Uso Principal
PySpark	Ingesta, Filtrado y Transformación Masiva.
OmegaConf	Lectura de reglas de negocio desde YAML.
Pandas	Escritura controlada de archivos locales (Parquet/CSV).

Tecnologías Utilizadas

- Python 3
- PySpark
- OmegaConf (YAML)
- Jupyter Notebook / Visual Studio Code
- GitHub (repositorio público)

Módulo de Extracción - `get_data`

- **Conexión:** Localiza el archivo CSV usando la ruta que definimos en el archivo de configuración (**config.yaml**).
- **Lectura Inteligente:** Utiliza **Spark** para detectar automáticamente qué tipo de dato hay en cada columna.
- **Filtrado Inmediato:** Antes de intentar hacer cálculos complejos, eliminamos todo lo que no nos sirve.

```
1
2 def get_data(spark, config):
3     #salida fecha proceso
4     #(output_path: data/processed/${fecha_proceso})
5     df = spark.read.option("header", True).option("inferSchema", True).csv(config.input.file_path)
6     df = df.withColumn("fecha_proceso", col("fecha_proceso").cast("string"))
7
8     return df.filter(
9         (col("fecha_proceso") >= config.filters.start_date) &
10        (col("fecha_proceso") <= config.filters.end_date) &
11        (col("pais") == config.filters.country)
12    )
```

Módulo de Transformación - `transform_data`

Este módulo implementa la lógica de normalización y limpieza mediante transformaciones sobre el DataFrame de Spark.

```
1
2 def transform_data(df, config):
3
4     df = df.dropDuplicates()
5
6     df = df.withColumn(
7         "cantidad_unidades",
8         when(
9             col("unidad") == "CS",
10            col("cantidad") * config.processing.box_size # llamamos al parametro box_size del config.yaml que tiene como valor 20
11        ).when(
12            col("unidad") == "ST",
13            col("cantidad")
14        ).otherwise(None) # Default null si no coincide ninguna unidad
15    )
16
17    routine_types = list(config.processing.delivery_types.routine)
18    bonus_types = list(config.processing.delivery_types.bonus)
19    all_types = routine_types + bonus_types
20
21
22    df = df.filter(
23        col("tipo_entrega").isin(all_types)
24    )
25
26
27    df = (
28        df
29        .withColumn(
30            "entrega_rutina",
31            col("tipo_entrega").isin(routine_types)
32        )
33        .withColumn(
34            "entrega_bonificacion",
35            col("tipo_entrega").isin(bonus_types)
36        )
37    )
38
39    df = df.withColumn(
40        "precio",
41        when(col("precio") > 0, col("precio"))
42        .cast(DecimalType(10, 2))
43    )
44
45    df = (
46        df
47        .dropna(subset=["cantidad_unidades"])
48        .filter(col("cantidad_unidades") > 0)
49    )
50
51
52    return df
```

1. Eliminación de Registros Duplicados.

`df = df.dropDuplicates()`

- **Operación:** Eliminación de redundancia.
- **Descripción Técnica:** Se aplica una transformación para identificar y remover filas idénticas en todo el dataset, garantizando los registros antes de iniciar el procesamiento computacional.

2. Normalización Condicional de Métricas (cantidad_unidades)

```
df = df.withColumn("cantidad_unidades",  
    when(col("unidad") == "CS", col("cantidad") * config.processing.box_size)  
    .when(col("unidad") == "ST", col("cantidad"))  
    .otherwise(None)
```

- **Condición:** (Case When).
- **Descripción Técnica:** Se estandariza la unidad de medida aplicando una lógica de ramificación:
 - **Caso CS:** Se realiza una operación aritmética de multiplicación utilizando el factor de conversión (**box_size = 20**) desde la configuración de **YAML**.
 - **Caso ST:** Se deja el valor original.
 - **(Otherwise):** Se asigna None (Null) a cualquier unidad no reconocida para facilitar su depuración o eliminación en etapas posteriores.

3. Filtrado por Pertenencia de Conjuntos (tipo_entrega)

```
all_types = routine_types + bonus_types  
df = df.filter(col("tipo_entrega").isin(all_types))
```

- **Operación:** Filtrado de tipo de entrega.
- **Descripción Técnica:** Se restringe el dataset únicamente a los registros cuyo tipo_entrega pertenezca a la lista de códigos válidos concatenados desde la configuración de **YAML**.

4. Atributos booleanos

```
.withColumn("entrega_rutina", col("tipo_entrega").isin(routine_types))  
.withColumn("entrega_bonificacion", col("tipo_entrega").isin(bonus_types))
```

- **Operación:** Generación de atributos booleanos.
- **Descripción Técnica:** Se crean dos nuevas columnas de tipo BooleanType evaluando la pertenencia del código de entrega a las sub-listas de "Rutina" o "Bonificación". Esto optimiza consultas.

5. Validación de Tipos y Redondeo de Decimales

```
when(col("precio") > 0, col("precio")).cast(DecimalType(10, 2))
```

- **Operación: Casteo de Tipo (*Type Casting*) con Redondeo.**
- **Descripción Técnica:**
 - **Validación:** Se aplica una condición para preservar solo valores estrictamente positivos (> 0). Los valores incorrectos se convierten en Null.
 - **Redondeo y Casteo:** Se fuerza la conversión a `DecimalType(10, 2)`. Esto no solo cambia el tipo de dato, sino que redondea aritméticamente los valores a 2 decimales, asegurando la precisión necesaria para cálculos monetarios y evitando errores de "punto flotante".

6. Integridad de Datos (Null)

```
.dropna(subset=["cantidad_unidades"])
```

```
.filter(col("cantidad_unidades") > 0)
```

- **Operación:** Manejo de Nulos y Restricciones Lógicas.
- **Descripción Técnica:**
 - **dropna:** Elimina los registros que resultaron en **Null** durante la normalización de unidades, limpiando datos con unidades de medida inválidas.
 - **filter:** Aplica una restricción de negocio para asegurar que no existan transacciones con cantidad cero o negativa en el dataset final.

Módulo de Carga y Estandarización - (load_data)

Esta etapa final se encarga de darle formato a los datos para su consumo ("**Capa Silver**"). No solo guardamos los datos, sino que los preparamos.

```
1
2 # Carga de datos
3 def load_data(df, config):
4     #se estandarizaron los nombres de paises
5     country_map = {
6         "GT": "Guatemala",
7         "PE": "Peru",
8         "EC": "Ecuador",
9         "SV": "El Salvador",
10        "HN": "Honduras",
11        "JM": "Jamaica"
12    }
13
14    df = df.replace(country_map, subset=["pais"])
15
16    #Estandar de nombres de columnas
17    output_df = df.select(
18        col("pais").alias("country"),
19        col("material"),
20        col("transporte").alias("transport"),
21        col("ruta").alias("route"),
22        col("precio").alias("price"),
23        col("cantidad_unidades").cast(IntegerType()).alias("unit_quantity"),
24        col("entrega_rutina").alias("routine_delivery"),
25        col("entrega_bonificacion").alias("bonus_delivery"),
26        col("fecha_proceso"),
27        current_timestamp().alias("load_date")
28    )
29
30    output_path = config.output.base_path
31    if os.path.exists(output_path):
32        shutil.rmtree(output_path)
33    os.makedirs(output_path, exist_ok=True)
34
35
36
37    pdf = output_df.toPandas() #se utilizo una conversión a Pandas para la ejecución local en Windows sin dependencias de binarios de Hadoop (winutils.exe).
38
39    pdf.to_parquet(output_path, engine='pyarrow', compression='snappy', partition_cols=['fecha_proceso'], index=False)
40
41    #CSV
42    for root, _, files in os.walk(output_path):
43        for file in files:
44            if file.endswith(".parquet"):
45                parquet_file = os.path.join(root, file)
46                csv_file = parquet_file.replace(".parquet", ".csv")
47                pd.read_parquet(parquet_file).to_csv(csv_file, index=False)
48                #print(f"Generado: {csv_file}") # validamos el csv que se creo
```

1. (Mapping) – Estandarización de Nombre de Países:

- Los códigos de país son útiles para sistemas, pero confusos para el análisis. Aplicamos un mapeo (replace) para expandirlos a sus nombres completos ("**Guatemala**", "**GT**").

2. Estándar de Nomenclatura:

- Renombramos todas las columnas del español al **Inglés** usando el estándar *snake_case* (minúsculas con guiones bajos).
- *Ejemplo:* cantidad_unidades → unit_quantity.
- *Por qué:* Esto facilita la integración con herramientas globales de BI y equipos de desarrollo internacionales.

3. Trazabilidad (Linaje del Dato):

- Agregamos la columna **load_date** con la fecha y hora exacta de ejecución (**current_timestamp**). Esto permite saber cuándo se generó cada registro.

4. (Spark + Pandas):

- Se convertimos el **DataFrame** procesado a **Pandas** para realizar la escritura física en el disco local.
- Generamos la salida en formato **Parquet** (comprimido con Snappy) particionado por fecha, y creamos copias espejo en **CSV**.

Orquestación y Ejecución Principal - Main

Aquí se inicializa el motor de Spark, se gestiona el flujo de ejecución y, lo más importante, se controlan los errores y el ciclo de vida de la aplicación.

```
1 from pyspark.sql import SparkSession
2
3 spark = SparkSession.builder \
4     .appName("ETL_CASO_INGENIERIA_DATOS") \
5     .master("local[*]") \
6     .config("spark.sql.execution.arrow.pyspark.enabled", "false") \
7     .getOrCreate()
8
9 spark.sparkContext.setLogLevel("ERROR")
10
11 try:
12     # obtener datos
13     df_raw = get_data(spark, config)
14
15     if df_raw.count() == 0:
16         print("No se encontraron registros que cumplan con los filtros")
17     else:
18         print(f"Registros encontrados: {df_raw.count()}")
19
20     #transformacion de la data y normalizacion
21     df_processed = transform_data(df_raw, config)
22     #df_processed.show(5)
23
24     #Carga de los datos
25     load_data(df_processed, config)
26
27     print("ETL REALIZADO")
28
29 except Exception as e:
30     print(f"Error durante la ejecucion: {e}")
31
32 finally:
33     spark.stop()
```

1. Inicialización del Motor (Spark Session):

- Levantamos una sesión de Spark configurada para ejecutarse en modo local (**local[*]**).

Decisión Técnica:

Se utilizó **local[*]** para simular un clúster utilizando todos los hilos del CPU.

2. Manejo de Errores - (Try-Except-Finally):

- Envolvemos todo el proceso en un bloque de seguridad. Si algo ocurre el se detecta el error, nos dice qué pasó y detiene el proceso de la ejecución.

3. Validación de Flujo (Early Stop):

- Después de leer los datos (**get_data**), hacemos un conteo rápido. Si no hay registros que cumplan con los filtros, el proceso se detiene.

Conclusión

La solución presentada cumple integralmente con los requerimientos del caso técnico mediante un flujo ETL desarrollado en **PySpark**. La arquitectura aplica principios sólidos de ingeniería de datos, garantizando flexibilidad a través de la parametrización con **OmegaConf**.

La estrategia técnica adoptada para la ejecución local en entornos Windows: se automatizó el paso de carga utilizando **Pandas** como puente de escritura. Esta decisión permitió prescindir de la configuración compleja de binarios de Hadoop (Winutils), logrando exportar los datos particionados por fecha tanto en formato **Parquet** como en **CSV** de manera eficiente. El resultado es un sistema escalable, portable y listo para ser reutilizado en distintos contextos de cliente.