

PROCESO DE ETL PRUEBA TÉCNICA

Miguel Enrique Fernández Azucena

Este documento presenta la documentación de la realización de un proceso de ETL realizado en los datos proveídos en la descripción de la prueba técnica

Programas utilizados

- SQL Server
- SQL Server Integration Services
- Visual Studio
- Visual Studio SQL Server Integration Services Projects Extension

Diagrama de archivos en Kaggle

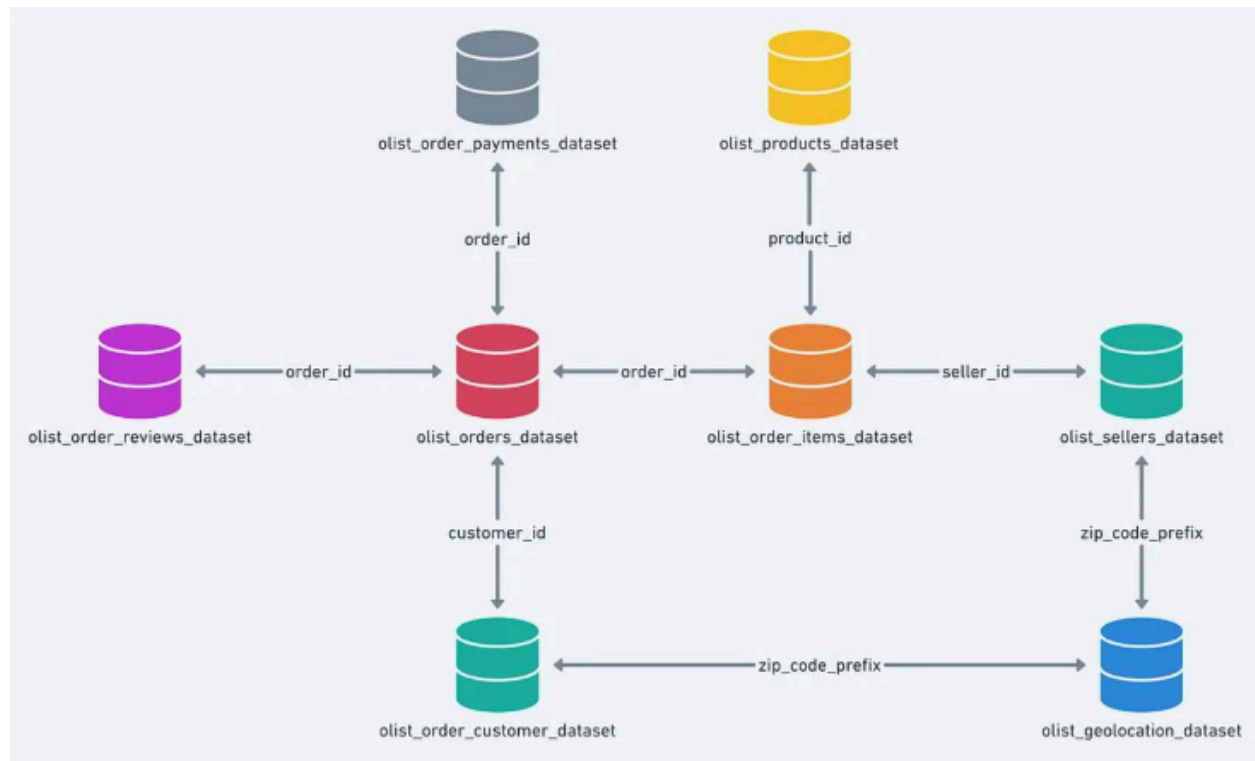


Figura 1. Diagrama de dataset Olist

El diagrama de la figura 1 describe la estructura de un dataset correspondiente a un negocio de e-commerce Brasileño. Cada uno de los iconos representa un archivo csv que contiene data correspondiente a esa tabla, es decir que el archivo **olist_order_customer_dataset** contiene información acerca de los clientes quienes hacen uso del servicio de e-commerce.

De estas tablas se utilizaron tres archivos para realizar un proceso de extracción, transformación y carga de datos a una base de datos. Los archivos seleccionados son:

- **olist_orders_dataset:** tabla que describe la información crítica de una orden, incluyendo detalles como el cliente, si la orden ya fue entregada y las diferentes fechas y horas relacionadas a la orden.
- **olist_order_items_dataset:** tabla que describe los ítems incluidos en cada orden junto con su precio nominal y su costo de envío.
- **olist_product_dataset:** tabla con información de los productos publicados en el sitio de e-commerce.

Estas tres tablas están relacionadas, siendo la tabla de order_items una tabla normalizada que granulariza la relación entre las órdenes y los ítems incluidos en cada orden. Se seleccionaron estos archivos ya que proveen información necesaria para formar perspectivas y conocimiento importante con respecto al negocio de e-commerce. De estas tablas se seleccionaron solamente ciertas columnas para su carga a la base de datos. Tanto la relación entre las tablas y las columnas seleccionadas se muestran en el siguiente diagrama entidad-relación:

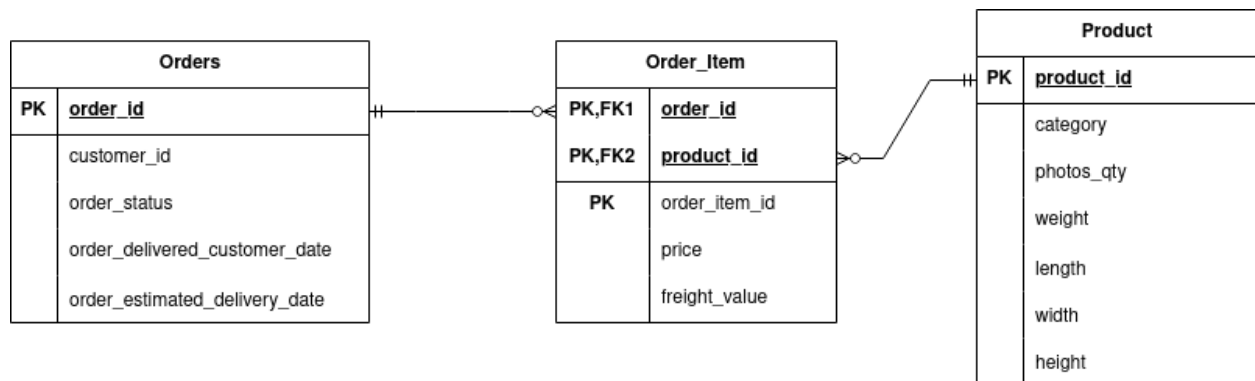


Figura 2. Diagrama entidad-relación

El proceso de ETL se realizó por medio de la herramienta de integración de SQL Server(SSIS por sus siglas en inglés). El proceso se realizó en un servidor de SQL Server local y una solución realizada en Visual Studio. Todos los archivos utilizados en el proceso y sus productos se encuentran en el repositorio de Github adjuntado a la entrega de este trabajo. Las instrucciones de creación de la base de datos se pueden ver en la figura 3.

```

CREATE DATABASE PruebaTecnica1
USE PruebaTecnica1
GO

CREATE TABLE orders(
    order_id varchar(50),
    customer_id varchar(50),
    order_status varchar(50),
    delivered_date datetime,
    estimated_delivery_date datetime,
    CONSTRAINT PK_order_id PRIMARY KEY(order_id)
);
GO

CREATE TABLE product(
    product_id varchar(50),
    category varchar(50),
    photos_qty int,
    product_weight float,
    product_length float,
    product_height float,
    product_width float,
    CONSTRAINT PK_product_id PRIMARY KEY(product_id)
);
GO

CREATE TABLE order_item(
    order_id varchar(50),
    product_id varchar(50),
    order_item_id int,
    price float,
    freight_value float
    CONSTRAINT PK_order_item PRIMARY KEY(order_id,product_id,order_item_id),
    CONSTRAINT FK_item_order_id FOREIGN KEY(order_id) REFERENCES orders(order_id),
    CONSTRAINT FK_item_product_id FOREIGN KEY(product_id) REFERENCES product(product_id)
);
GO

```

Figura 3. Instrucciones de creación de las tablas en la base de datos

Estas columnas se ven reflejadas dentro del proyecto de integración realizado en Visual Studio. Como parte del proyecto se definen los archivos csv como fuentes de datos de archivo plano, como parte de este proceso los datos se mapean con los siguientes tipos.

- **Orders**
 - **Order_id:** string
 - **Customer_id:** string
 - **Order_status:** string
 - **Order_purchase_timestamp:** database timestamp [DB_DBTIMESTAMP]
 - **Order_approved_at:** database timestamp [DB_DBTIMESTAMP]

- **Order_delivered_carrier_date:** database timestamp [DB_DBTIMESTAMP]
- **Order_delivered_customer_date:** database timestamp [DB_DBTIMESTAMP]
- **Order_estimated_delivery_date:** database timestamp [DB_DBTIMESTAMP]
- **Order_item**
 - **Order_id:** string
 - **Order_item_id:** integer
 - **Product_id:** string
 - **Seller_id:** string
 - **Shipping_limit_date:** database timestamp [DB_DBTIMESTAMP]
 - **Price:** float
 - **Freight_value:** float
- **Product**
 - **Product_id:** string
 - **Product_category_name:** string
 - **Product_name_length:** integer
 - **Product_description_length:** integer
 - **Product_photos_qty:** integer
 - **Product_weight_g:** float
 - **Product_length_cm:** float
 - **Product_height_cm:** float
 - **Product_width_cm:** float

De esta forma Visual Studio extrae los datos de los archivos para que estos se carguen a la base de datos. En el siguiente paso se seleccionan las columnas de interés y se cargan en la base de datos.

El proceso de carga se debe de hacer en dos pasos, ya que si se busca paralelizar la carga de la tabla order_item junto con las demás tablas se produce un problema de dependencia, por ello se realiza primero la carga de las tablas de orders y product, por último se cargan los valores de la tabla order_item. El primer paso de carga se encuentra dentro del primer data flow mostrado en la figura 4.

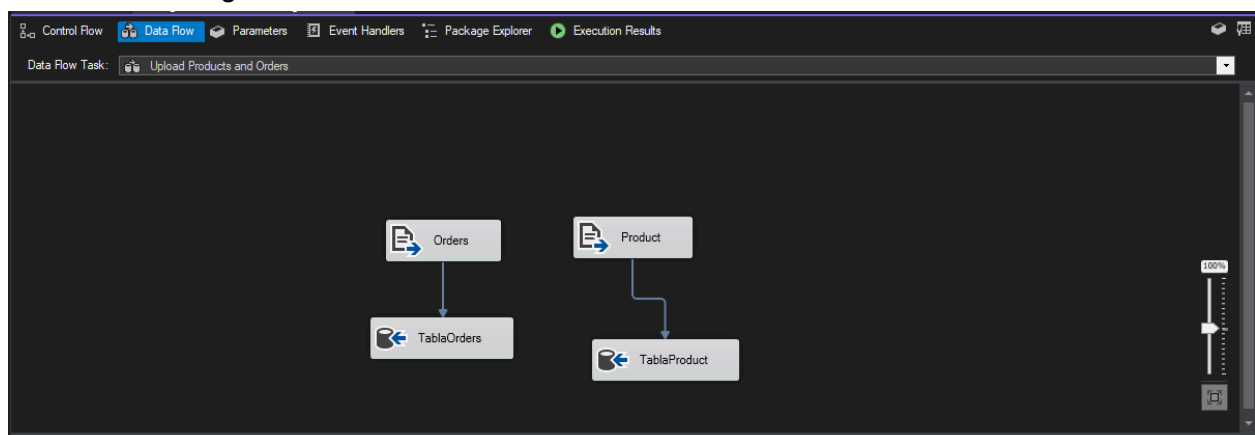


Figura 4. Primer proceso de data flow para subir los datos de productos y orders

Una vez ha terminado ese proceso se procede a cargar los datos de la tabla order_item indicado en su proceso de data flow mostrado en la figura 5.

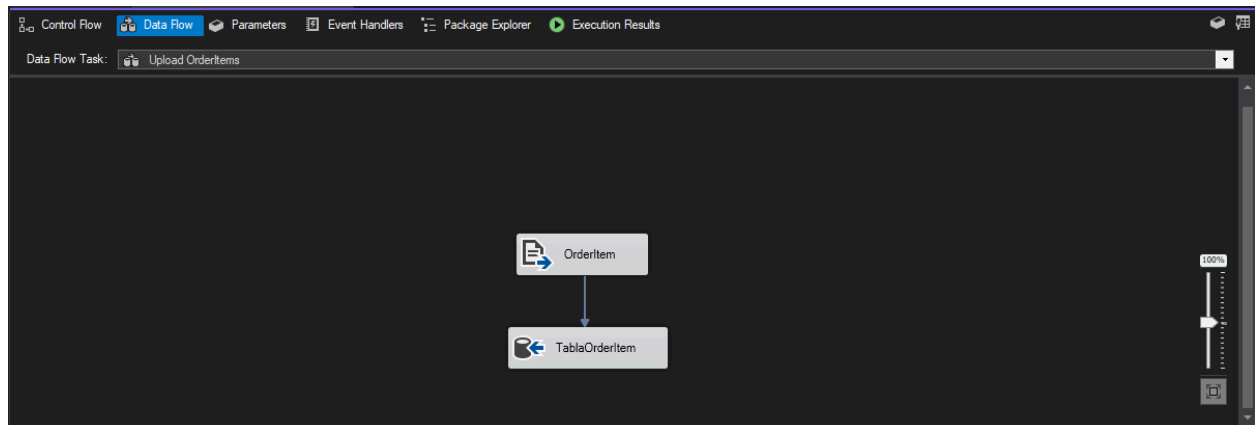


Figura 5. Segundo proceso de dataflow para la carga de datos a tabla order_item.

En la figura 6 se muestra el flujo general del proceso ETL en Visual Studio, donde el primer Dataflow titulado 'Carga de Product y Orders' corresponde a la figura 4 y el segundo data flow corresponde a la carga de los datos de la tabla order_item mostrada en la figura 5.

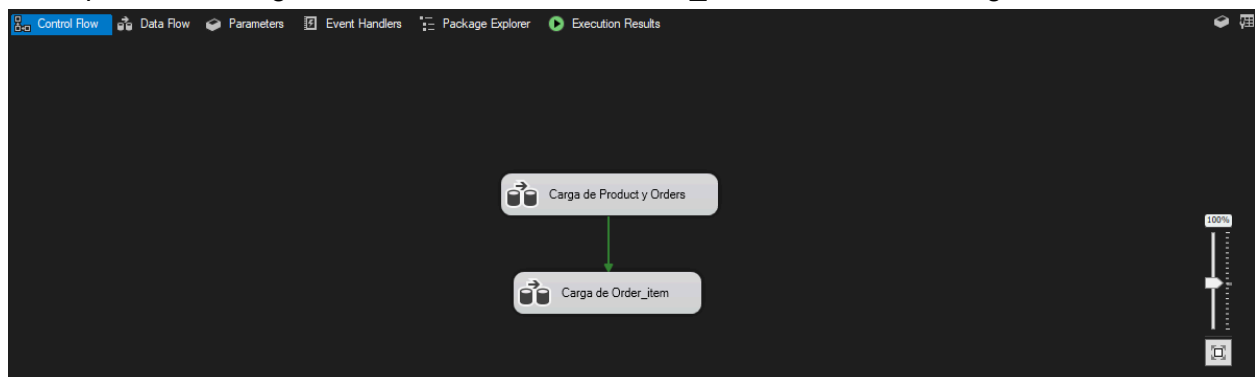


Figura 6. Flujo general de proceso ETL en Visual Studio.

Los resultados se pueden ver cargados en la base de datos por medio de las queries ilustradas en las figuras 7, 8 y 9.

Query executed successfully. DESKTOP-17UUQAC (16.0 RTM) | DESKTOP-17UUQAC\miken ... | PruebaTecnica1 | 00:00:00 | 32,951 rows

product_id	category	photos_qty	product_weight	product_length	product_height	product_width
1	perfumaria	6	300	20	16	16
2	automotivo	4	1225	55	10	26
3	cama_mesa_banho	2	300	45	15	35
4	utilidades_domesticas	3	550	19	24	12
5	relogios_presentes	4	250	22	11	15
6	automotivo	1	100	16	15	16
7	cool_stuff	1	700	25	5	15
8	consoles_games	1	600	30	20	20
9	cama_mesa_banho	1	6000	40	4	30
10	moveis_decoracao	3	600	26	8	22
11	cama_mesa_banho	1	1800	47	21	41
12	beleza_saude	6	300	30	15	15
13	fashion_calçados	3	1850	36	37	16
14	informatica_acessorios	6	650	25	10	15
15	utilidades_domesticas	2	750	30	5	35
16	cool_stuff	1	355	22	16	12
17	moveis_decoracao	3	2600	105	3	70
18	cool_stuff	3	600	40	30	40
19	brinquedos	5	350	45	30	30
20	pet_shop	1	200	20	11	11
21	utilidades_domesticas	3	150	20	10	15

Figura 7. Dados carregados em tabela product

Query executed successfully. DESKTOP-17UUQAC (16.0 RTM) | DESKTOP-17UUQAC\miken ... | PruebaTecnica1 | 00:00:00 | 99,441 rows

order_id	customer_id	order_status	delivered_date	estimated_delivery_date
1	"3ce43f183e68e07877b285a838db11a"	delivered	2017-09-20 23:43:48.000	2017-09-29 00:00:00.000
2	f6dd3ec061db4e3987629fe6b26e5c0e	delivered	2017-05-12 16:04:24.000	2017-05-15 00:00:00.000
3	"6489ae5e433f3693df5ad4372dab6d3"	delivered	2018-01-22 13:19:16.000	2018-02-05 00:00:00.000
4	d4eb9395c8c0431ee92fce09860c5a06	delivered	2018-08-14 13:32:39.000	2018-08-20 00:00:00.000
5	"58dbd0b2d70206bf40e62cd34e84d795"	delivered	2017-03-01 16:42:31.000	2017-03-17 00:00:00.000
6	"816cbea969fe5b689b39cfc97a506742"	delivered	2017-05-22 13:44:35.000	2017-06-06 00:00:00.000
7	"32e2e6ab09e778d99bf2e0ecd4898718"	delivered	2017-12-18 22:03:38.000	2018-01-04 00:00:00.000
8	"9ed5e522dd9d85b4af4a077526d8117"	delivered	2018-07-09 14:04:07.000	2018-07-25 00:00:00.000
9	"16150771df4776261284213b89c304e"	delivered	2018-03-29 18:17:31.000	2018-03-29 00:00:00.000
10	"351d3cb2cee3c7d0af6616c82df21d3"	delivered	2018-07-04 17:28:31.000	2018-07-23 00:00:00.000
11	c6fc061d8f6b1e2b2eac259bac71a49	delivered	2018-03-29 00:04:19.000	2018-04-09 00:00:00.000
12	"6a899e55865de6549a58d2c6845e5604"	delivered	2018-08-07 13:56:52.000	2018-08-07 00:00:00.000
13	"5d178120c29c61748ea95bac23cb8f25"	delivered	2018-07-31 01:04:15.000	2018-08-22 00:00:00.000
14	"2355af7c75e7c98b43a87b2a7210dc5"	delivered	2018-02-26 13:55:22.000	2018-03-06 00:00:00.000
15	"2a30c97668e81df7c17a8b14447eeeba"	delivered	2018-08-22 12:02:27.000	2018-08-28 00:00:00.000
16	"8a250edc40ebc5c3940ebc940f16a7eb"	delivered	2018-04-30 17:54:25.000	2018-05-09 00:00:00.000
17	fff5169e583fd07ac9fec88962f189d	delivered	2018-05-18 16:46:31.000	2018-05-22 00:00:00.000
18	"3773b0cf1a6fbd29233ea1c1b573c4f22"	delivered	2017-08-19 15:22:17.000	2017-09-01 00:00:00.000
19	"2b01d668726fb0b751c55918c043b7b"	delivered	2018-03-12 18:46:34.000	2018-03-19 00:00:00.000
20	"98884e672b5ca30573b99628"	delivered	2018-03-27 14:51:47.000	2018-04-05 00:00:00.000
21	"6a3b26c9f2707d258605e27bef19f488"	delivered	2017-12-09 17:27:23.000	2017-12-07 00:00:00.000

Figura 8. Dados carregados em tabela orders

order_id	product_id	order_item_id	price	freight_value
"00010242fe8c5a6d1ba2dd792cb16214"	"4244733e06e7ecb4970a6e2683c13e61"	1	58.9000015258789	13.289999961853
"00018772f0320c557190d7a144bddd3"	e5f2d52b802189ee658865ca93d83a8f	1	239.899993896484	19.9300003051758
"000229ec398224ef6ca0657da4fc703e"	c777355d18b72b67abbeef9df44d0fd	1	199	17.8700008392334
"00024acbcd0a6daa1e931b038114c75"	"7634da152a4610f1595efa32f14722fc"	1	12.9899997711182	12.789999961853
"00042b26cf59d7ce69dfabb4e55b4fd9"	ac6c362306f830de03045865e4e10089	1	199.899993896484	18.1399993896484
"00048cc3ae777c65dbb7d2a0634bc1ea"	ef92defde845ab8450f9d70c526ef70f	1	21.8999996185303	12.6899995803833
"00054e8431b9d7675808bcb819b4a32"	"8d4f2bb7e93e6710a28f34a83ee7d28"	1	19.8999996185303	11.8500003814697
"00057fe39319847cbb9d288c5617fa6"	"557d850972a7d6f792fd18ae1400d9b6"	1	810	70.75
"0005a1a1728c9d785b8e2b08b904576c"	"310ae3c140f94b03219ad0adc3c778"	1	145.949996948242	11.6499996185303
"0005f50442cb953dcd1d21e1fb923495"	"4539b0e1091c278df193e5a1d63b39f"	1	53.9900016784668	11.3999996185303
"00061f2a7bc09da83e415a52dc8a4af1"	d63c1011f49d98b976c352955b1c4bea	1	59.9900016784668	8.88000011444092
"00063b381e2406b52ad429470734ebd5"	f177554ea93259a5b282f24e33f65ab6	1	45	12.9799995422363
"0006ec9db01a64e59a68b2c340bf65a7"	"99a4788cb24856965c36a24e339b6058"	1	74	23.3199996948242
"000828aa423d2a3f00fcb17cd7d8719"	"368c6c730842d78016ad823897a372db"	1	49.9000015258789	13.3699998855591
"000828aa423d2a3f00fcb17cd7d8719"	"368c6c730842d78016ad823897a372db"	2	49.9000015258789	13.3699998855591
"0009792311464db532ff769bf7b182ae"	"8cab8abac59158715e0d70a36c807415"	1	99.9000015258789	27.6499996185303
"0009c9a17916a706d71784483a5d643"	"3f27ac8e699df3d300ec4a5d8c5cf0b2"	1	639	11.3400001525879
"000aed2e25dbad2f9dbd70584c5a2ded"	"4fa33915031a8cde03dd0d3e8fb27f01"	1	144	8.77000045776367
"000c3e6612759851cc3bb4b83257986"	b50c950aba0dcead2c48032a690ce817	1	99	13.710000038147
"000e562887b1f2006d75e0be9558292e"	"5ed9eaf534f6936b51d0b6c5e4d5c2e9"	1	25	16.1100006103516
"000e63d38ae8c00b0bcb5a30573b99628"	"553a0e75904d3116a072507a3635d2877"	1	47.9000015258789	8.88000011444092

Figura 9. Datos cargados en la tabla order_item

Finalmente se realiza el despliegue del paquete que resulta del proceso de ETL, de esta forma se facilita el acceso al proceso automático de ETL ya que se puede ejecutar desde el mismo servidor de base de datos. Desde Visual Studio se despliega el paquete al catálogo de servicios de integración de SQL Server. En la figura 10 se muestra este proceso realizado de forma satisfactoria.

Action	Result
Loading project	Passed
Connecting to destination server	Passed
Changing protection level	Passed
Deploying project	Passed

Figura 10. Despliegue de paquete a servicio de integración dentro de servidor SQL

La ejecución de este paquete se puede realizar de la siguiente manera. Se expande el folder de Integration Services Catalog dentro del programa de manejo de SQL Server (SSMS por sus siglas en inglés). Se expanden las carpetas de 'SSISDB', 'PruebaTecnica1', 'Projects', 'PruebaTecnica' y 'Package' hasta llegar al propio archivo de package.

Luego hacer click derecho en el archivo Package y seleccionar ejecutar tal como lo muestra la figura 11.

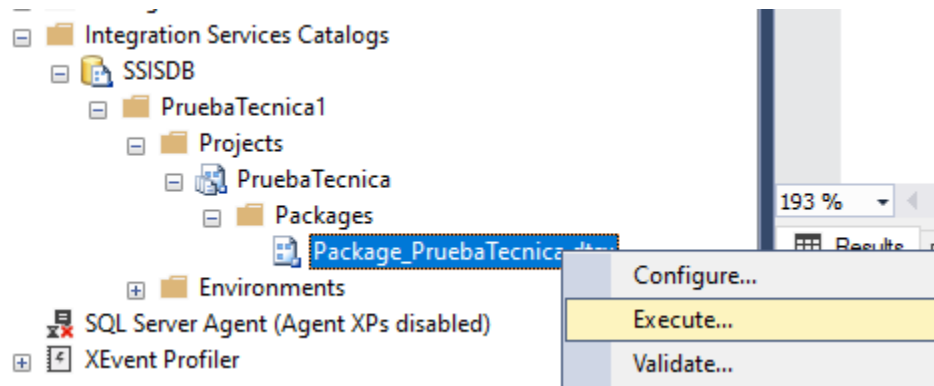


Figura 11. Ejecución de proceso ETL desde SSMS

Fuentes

1. Fuente de datos de servicio brasileño de e-commerce, Kaggle.
<https://www.kaggle.com/datasets/olistbr/brazilian-ecommerce>
2. Proyectos de ETL con SSIS y Visual Studio, Microsoft.
<https://learn.microsoft.com/en-us/sql/integration-services/ssis-how-to-create-an-etl-package?view=sql-server-ver16>
3. Instalación de SSIS, Microsoft.
<https://learn.microsoft.com/en-us/sql/integration-services/install-windows/install-integration-services?view=sql-server-ver16>
4. Tipos de datos de SSIS en Visual Studio, Microsoft.
<https://learn.microsoft.com/en-us/sql/integration-services/data-flow/integration-services-data-types?view=sql-server-ver16>