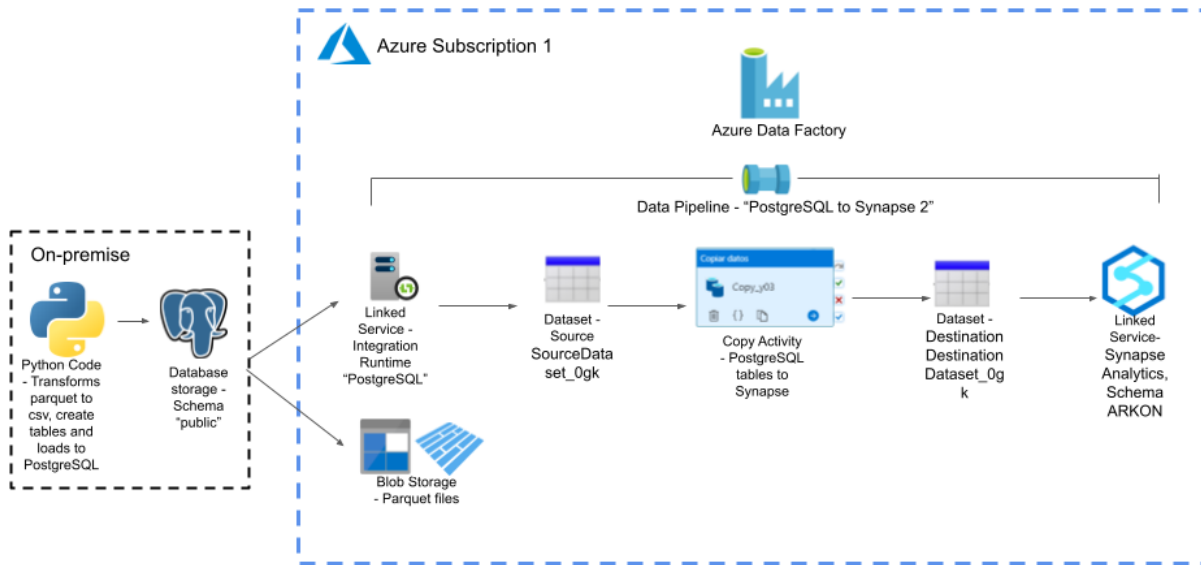


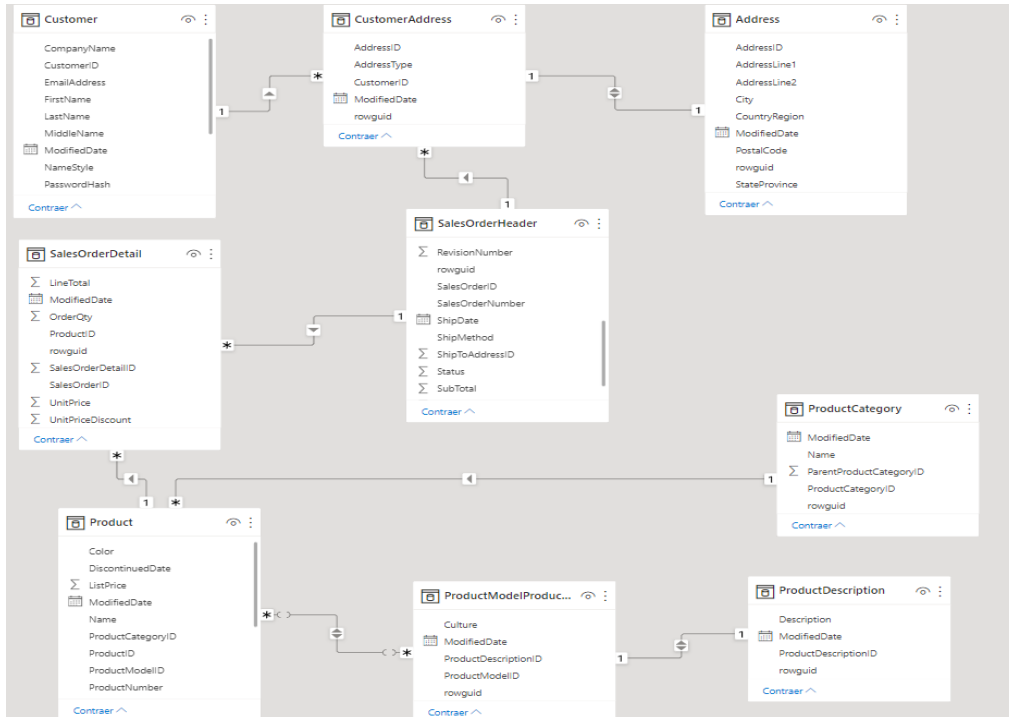
1. Diagrama de Solución

Para la solución de la prueba y la creación de un data pipeline se elaboró la siguiente propuesta de solución. Esta propuesta y el paso a paso se explica a detalle a partir de la sección 2 de este documento.



2. Modelo de Datos y DLLs

Para tener una perspectiva de la relación de tablas que serán cargadas a una base de datos PostgreSQL se cargaron primeramente los archivos en formato csv (transformados de parquet a csv con el código Python del punto 3) a PowerBI. La siguiente imagen muestra la relación entre tablas:



La creación de un servidor de PostgreSQL se realizó manualmente y se conectó a DBEAVER, que es una herramienta de administración de bases de datos. Se creó el esquema “public” y se definieron las tablas con los queries que se han agregado a anexos (al final de este documento). Las llaves para la conexión al servidor se presentan en el código Python de la sección 3.

3. Carga de datos a PostgreSQL (con Python)

El código escrito fue parte de la solución para un examen previo pero se tomó como primer paso para la solución a esta nueva prueba. El código fuente puede consultarse en el repositorio:

<https://github.com/hromolozano/Arkon-Simple-Data-Pipeline.git>

El código contiene comentarios que permiten el entendimiento del mismo pero se describe su secuencia a continuación:

1. Se importa el módulo para manipulación de datos, se determina una lista con el nombre de los archivos parquet y se lee cada documento para convertirlo a un dataframe.
2. Se creó una base de datos llamada “postgres” y los queries que crean las tablas dentro del esquema “public”. En el código se definen los nombres de los archivos de extensión txt para ser leídos y convertidos a un string. Esto resultó en una

lista donde cada elemento es el query a ejecutar en la siguiente celda de código.

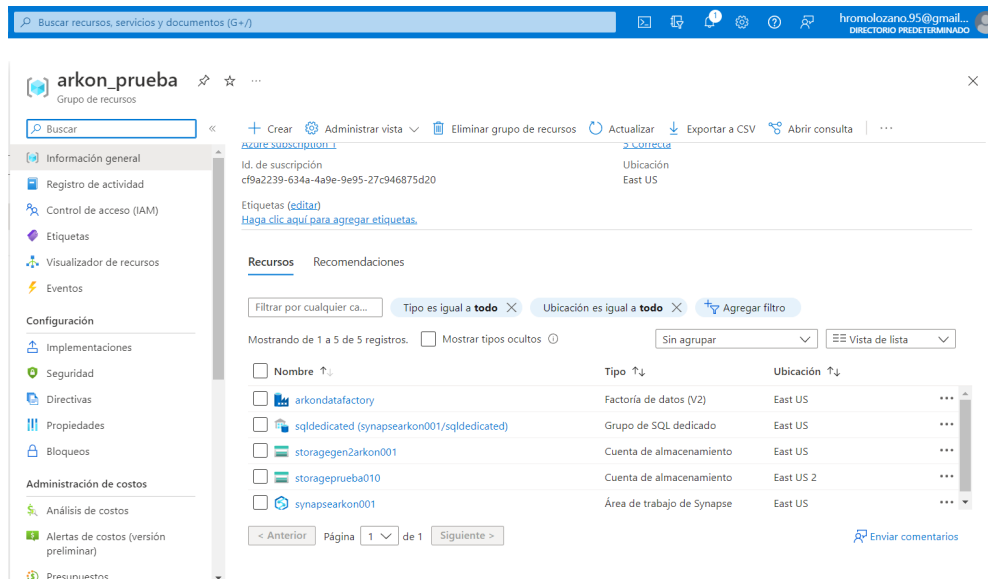
3. Se importó la librería SQLAlchemy para conectarnos desde Python a la base de datos utilizando las llaves determinadas durante la creación de la misma. Un ciclo for lee cada query y crea la tabla en la conexión a PostgreSQL. Si la tabla existe se imprime un mensaje de que la tabla ya existe.
4. Una vez creadas las tablas se inserta la información de cada archivo csv (creados en el paso 1) mediante otro ciclo for. En este caso cada ejecución de la inserción de datos reemplaza a la información existente en la tabla.

***Consideraciones para futuras correcciones y implementaciones de este código:

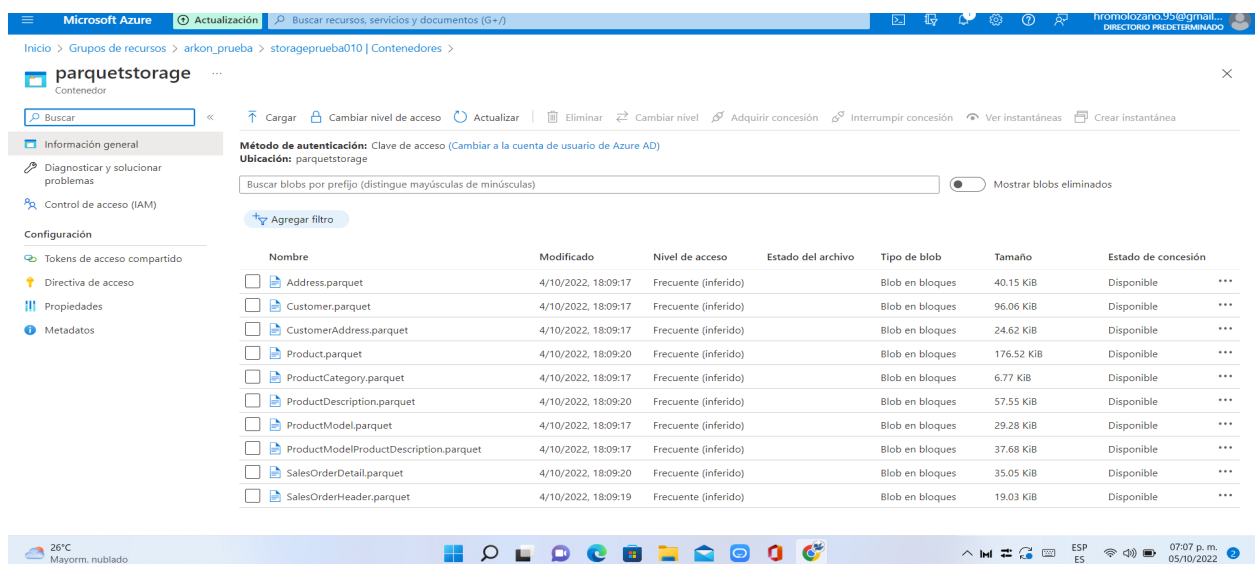
- Que el código lea desde el local path los archivos con extensión .parquet automáticamente y no mediante la determinación manual de los nombres de los archivos. El código debe ser ejecutado automáticamente cada que los archivos sean actualizados o sean agregados nuevos archivos parquet.
- Que el código también asigne el nombre de la tabla, tomando automáticamente el prefijo del archivo o los archivos leídos con la porción de código anterior.
- Que el DDL sea creado automáticamente con la lectura del csv, donde el usuario deberá pasar argumentos como las primary keys de la tabla.
- La inserción de datos debe considerar si la tabla es incremental o full, es decir, que el código permita insertar solo la información de la tabla faltante por medio de una comparación entre lo que ya está en PostgreSQL y lo que se actualizó en el nuevo archivo parquet. En el caso de este código supone una inserción full (borra en cada ejecución del código la información e inserta todo de nuevo).

4. Carga de archivos Parquet a Blob Storage

Después de haber creado la base de datos se pensó en una solución dentro de los servicios de la nube de Azure. Se creó una cuenta gratuita y se crearon los recursos necesarios dentro del Azure Portal. El grupo de recursos se nombró "arkon_prueba" y contiene los servicios de un Data Factory, un Blob storage y un Synapse Analytics que incluye un storage. Además se creó un SQL Pool dedicado para recibir la información ingestada.



Para ingestar los archivos de formato parquet se creó un contenedor dentro del servicio de Blob storage donde se cargó manualmente cada archivo.

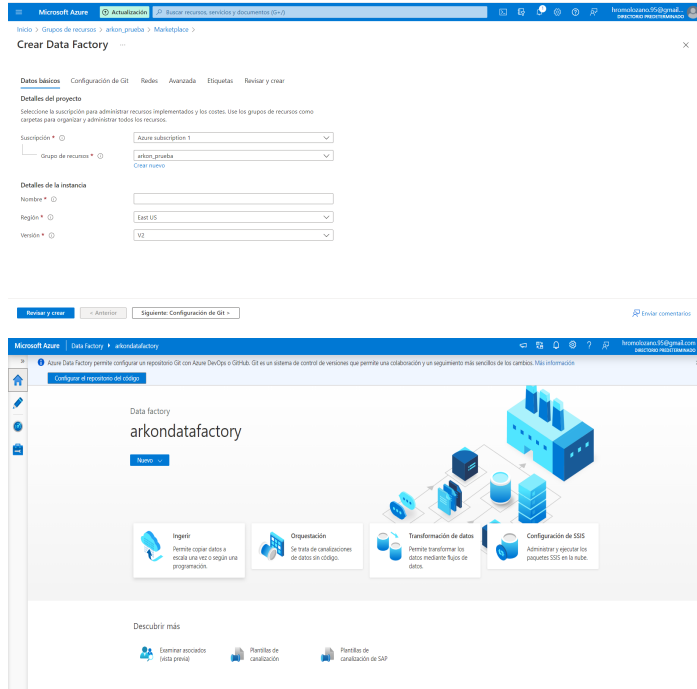


***Consideraciones para futuras implementaciones de esta parte de la solución:

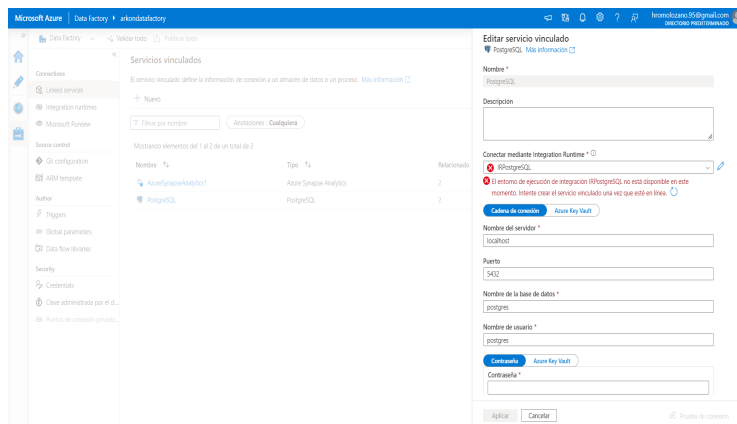
- Que el código Python pueda conectarse a la nube y que cargue el archivo parquet a este contenedor cada que los archivos son actualizados o son agregados al local path. Si esto se realiza es posible que se pueda eliminar la base de datos PostgreSQL y la ingesta a Synapse sea desde Blob storage.

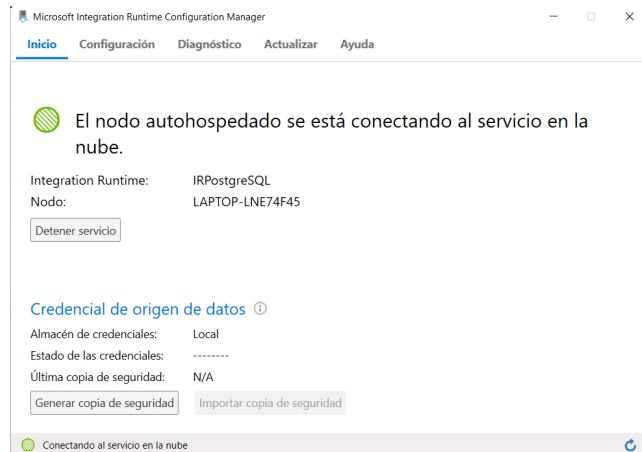
5. Pipeline Azure Data Factory a Synapse Analytics

Se creó el servicio de Data Factory y posteriormente se crearon los linked services que permitirán la conexión a la fuente de datos y a la instancia final donde serán ingestados.

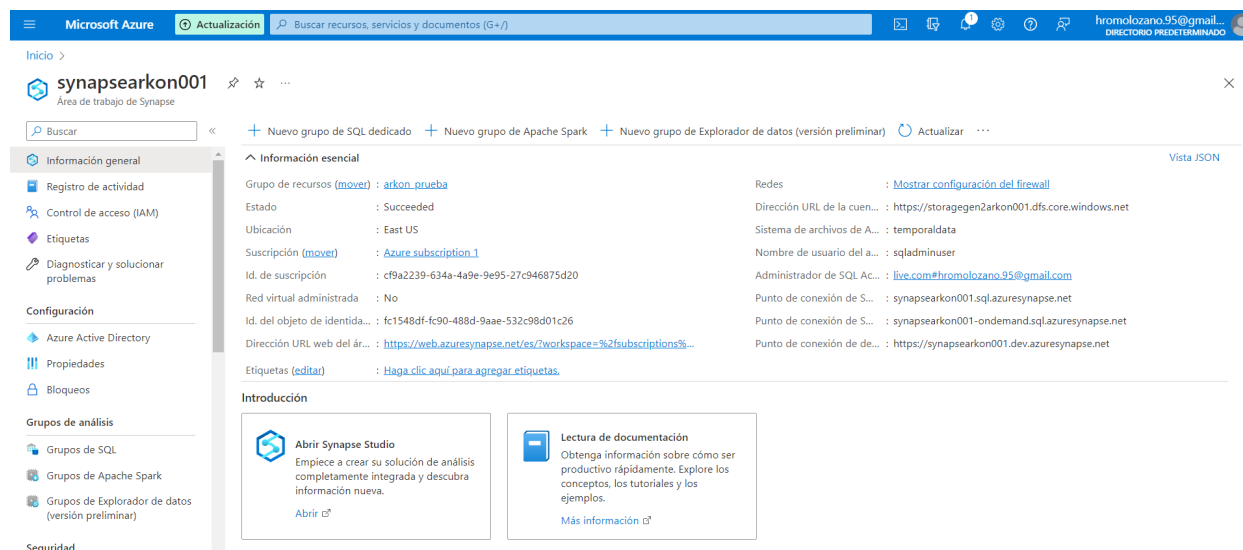


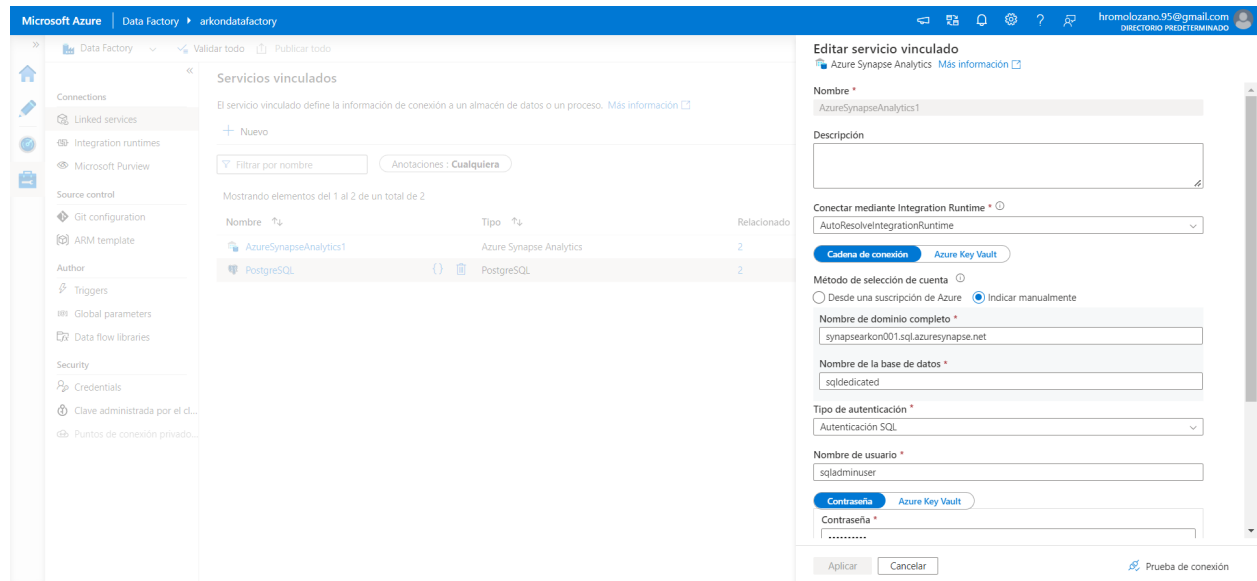
La primera conexión fue poder enlazar el servidor de PostgreSQL (en mi computadora local) con la nube. Esto se realizó con la aplicación de Integration Runtime. El procedimiento para esto es crear el linked service primero y después ingresar las keys a la aplicación descargada de Integration Runtime. Esto conecta nuestra base de datos para ser la fuente en el pipeline de Data Factory.



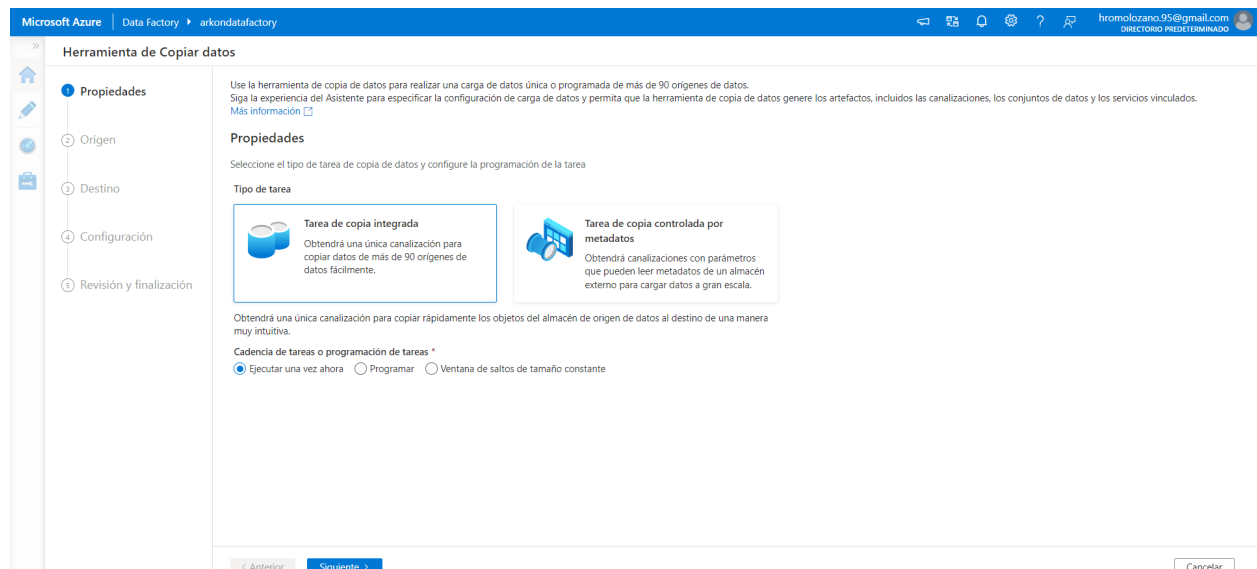


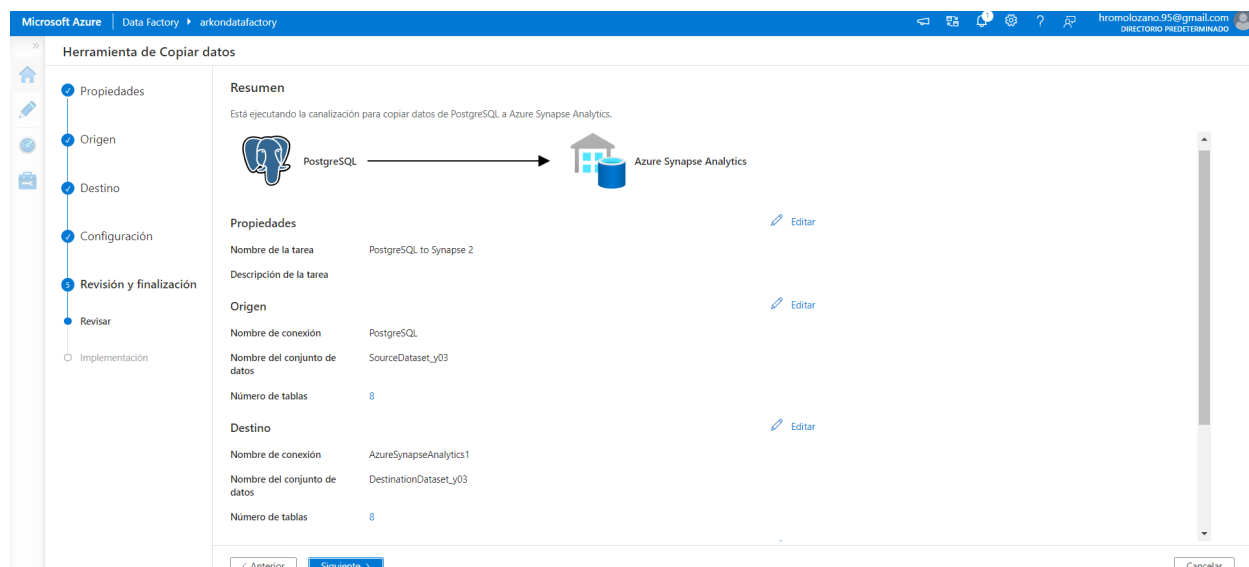
Posteriormente se creó el servicio de Synapse Analytics y se creó el linked service para esta instancia. Para el linked service también se creó un SQL pool. Dentro de esta nueva base de datos dedicada se creó el esquema “ARKON”, desde una conexión de DBeaver y no desde los scripts de SQL.



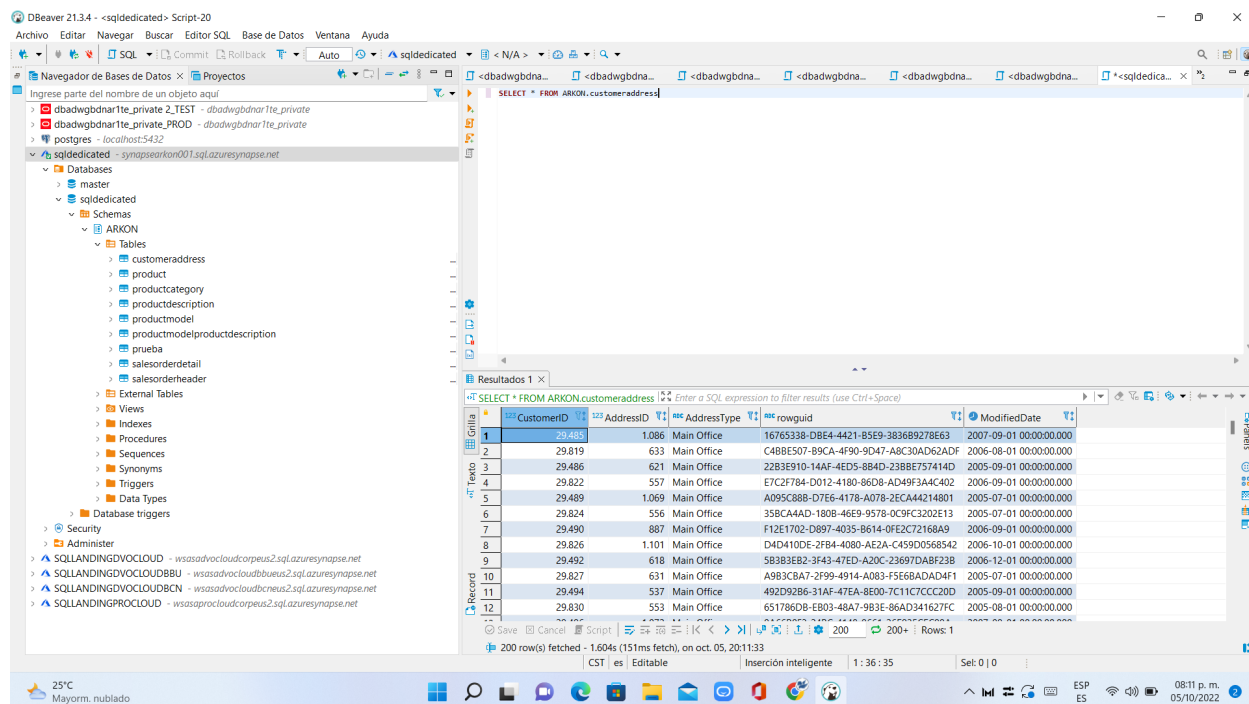


Una vez creada esta infraestructura se tiene acceso a la opción de Ingesta de datos desde el Azure Data Factory Studio. Aquí podemos configurar paso a paso la copia de datos de una fuente a nuestra base de datos final. Este proceso nos guía de una manera más práctica y sencilla para elegir nuestros servicios vinculados y la creación de los datasets. Este pipeline se ejecutó de manera inmediata pero puede ser programado para que se haga una vez al día una vez que con el código Python hemos actualizado la información (en un escenario productivo).





Terminando este proceso y la ejecución se validó la información para concluir que el data está ingestado en Synapse.



6. Anexos

Se anexan los queries para la creación de tablas del modelo de datos presentado:


```
CREATE TABLE IF NOT EXISTS public.address
(
  AddressID integer PRIMARY KEY,
  AddressLine1 varchar(4000),
  AddressLine2 varchar(4000),
  City varchar(4000),
  StateProvince varchar(4000),
  CountryRegion varchar(4000),
  PostalCode integer,
  rowguid varchar(4000),
  ModifiedDate date
)
```

```
CREATE TABLE public.customer
(
  CustomerID integer PRIMARY KEY,
  NameStyle boolean,
  Title varchar(4000),
  FirstName varchar(4000),
  MiddleName varchar(4000),
  LastName varchar(4000),
  Suffix varchar(4000),
  CompanyName varchar(4000),
  SalesPerson varchar(4000),
  EmailAddress varchar(4000),
  Phone integer,
  PasswordHash varchar(4000),
  PasswordSalt varchar(4000),
  rowguid varchar(4000),
  ModifiedDate date
)
```

```
CREATE TABLE public.customeraddress
(
  CustomerID integer,
  AddressID integer,
  AddressType varchar(4000),
  rowguid varchar(4000),
  ModifiedDate date,
  PRIMARY KEY (CustomerID, AddressID)
)
```

```
CREATE TABLE public.product
(
  ProductID integer,
  Name varchar(4000),
  ProductNumber varchar(4000),
  Color varchar(4000),
  StandardCost float,
  ListPrice float,
  Size varchar(6),
  Weight float,
  ProductCategoryID integer,
  ProductModelID integer,
  SellStartDate date,
  SellEndDate date,
  DiscontinuedDate date,
  ThumbNailPhoto varchar(4000),
  ThumbnailPhotoFileName varchar(4000),
  rowguid varchar(4000),
  ModifiedDate date,
  PRIMARY KEY (ProductID, ProductCategoryID, ProductModelID)
)
```

```
CREATE TABLE public.productcategory
(
  ProductCategoryID integer PRIMARY KEY,
  ParentProductCategoryID integer,
  Name varchar(4000),
  rowguid varchar(4000),
  ModifiedDate date
)
```

```
CREATE TABLE public.productdescription
(
  ProductDescriptionID integer PRIMARY KEY,
  Description varchar(4000),
  rowguid varchar(4000),
  ModifiedDate date
)
```

```
CREATE TABLE public.productmodelproductdescription
(
  ProductModelID integer,
  ProductDescriptionID integer,
  Culture varchar(12),
```

```
rowguid varchar(4000),  
ModifiedDate date,  
PRIMARY KEY (ProductModelID, ProductDescriptionID)  
)
```

```
CREATE TABLE public.salesorderdetail  
(  
SalesOrderID integer,  
SalesOrderDetailID integer,  
OrderQty integer,  
ProductID integer,  
UnitPrice float,  
UnitPriceDiscount float,  
LineTotal float,  
rowguid varchar(4000),  
ModifiedDate date,  
PRIMARY KEY (SalesOrderID, ProductID)  
)
```

```
CREATE TABLE public.salesorderheader  
(  
SalesOrderID integer,  
RevisionNumber integer,  
OrderDate date,  
DueDate date,  
ShipDate date,  
Status integer,  
OnlineOrderFlag boolean,  
SalesOrderNumber varchar(4000),  
PurchaseOrderNumber varchar(4000),  
AccountNumber varchar(4000),  
CustomerID integer,  
ShipToAddressID integer,  
BillToAddressID integer,  
ShipMethod varchar(4000),  
CreditCardApprovalCode varchar(4000),  
SubTotal float,  
TaxAmt float,  
Freight float,  
TotalDue float,  
Comment varchar(4000),  
rowguid varchar(4000),  
ModifiedDate date,  
PRIMARY KEY (CustomerID, SalesOrderID)
```

)

Create a POSTGRESQL

Create Python Code for data loading to database

Create an Azure Data Factory Pipeline
Linked Service to Postgresql
Create Blob Storage

Creación de grupo de recursos

Creación de Azure Data Factory

Creación de Synapse Analytics
Lake Storage Gen 2

Storage Account

Blob Storage for parquet files

The screenshot shows the Microsoft Azure portal interface for a storage account named 'parquetstorage'. The left sidebar contains navigation links for 'Información general', 'Configuración', 'Propiedades', and 'Metadatos'. The main content area displays a table of blobs with the following columns: Nombre, Modificado, Nivel de acceso, Estado del archivo, Tipo de blob, Tamaño, and Estado de concesión. The table lists several parquet files, including Address.parquet, Customer.parquet, CustomerAddress.parquet, Product.parquet, ProductCategory.parquet, ProductDescription.parquet, ProductModel.parquet, ProductModelProductDescription.parquet, SalesOrderDetail.parquet, and SalesOrderHeader.parquet. The status of each blob is 'Disponibile'.

Nombre	Modificado	Nivel de acceso	Estado del archivo	Tipo de blob	Tamaño	Estado de concesión
Address.parquet	4/10/2022, 18:09:17	Frecuente (inferido)		Blob en bloques	40.15 KiB	Disponibile
Customer.parquet	4/10/2022, 18:09:17	Frecuente (inferido)		Blob en bloques	96.06 KiB	Disponibile
CustomerAddress.parquet	4/10/2022, 18:09:17	Frecuente (inferido)		Blob en bloques	24.62 KiB	Disponibile
Product.parquet	4/10/2022, 18:09:20	Frecuente (inferido)		Blob en bloques	176.52 KiB	Disponibile
ProductCategory.parquet	4/10/2022, 18:09:17	Frecuente (inferido)		Blob en bloques	6.77 KiB	Disponibile
ProductDescription.parquet	4/10/2022, 18:09:20	Frecuente (inferido)		Blob en bloques	57.55 KiB	Disponibile
ProductModel.parquet	4/10/2022, 18:09:17	Frecuente (inferido)		Blob en bloques	29.28 KiB	Disponibile
ProductModelProductDescription.parquet	4/10/2022, 18:09:17	Frecuente (inferido)		Blob en bloques	37.68 KiB	Disponibile
SalesOrderDetail.parquet	4/10/2022, 18:09:20	Frecuente (inferido)		Blob en bloques	35.05 KiB	Disponibile
SalesOrderHeader.parquet	4/10/2022, 18:09:19	Frecuente (inferido)		Blob en bloques	19.03 KiB	Disponibile

Connect PostgreSQL server to Azure Integration Runtime Service Linked Service to Integration Runtime (server and database keys to PostgreSQL)

The screenshot shows the Microsoft Azure Data Factory portal. The left sidebar contains navigation options: Home, Connections, Linked services, Integration runtimes, Microsoft Purview, Source control, Git configuration, ARM template, Author, Triggers, Global parameters, Data flow libraries, Security, Credentials, and Clave administrada por el d... (Administered by the d...). The main content area is titled 'Servicios vinculados' (Linked services) and shows a message: 'No hay servicio vinculado para r...' (There is no linked service for r...). The right sidebar is titled 'Nuevo servicio vinculado' (New linked service) and shows the configuration form for a PostgreSQL linked service. The form includes fields for Nombre (Name), Descripción (Description), and a dropdown for 'Conectar mediante Integration Runtime' (Connect via Integration Runtime) set to 'IRPostgreSQL'. Below this, there are fields for 'Nombre del servidor' (Server name) set to 'localhost', 'Puerto' (Port) set to '5432', 'Nombre de la base de datos' (Database name) set to 'postgres', and 'Nombre de usuario' (Username) set to 'postgres'. There are also buttons for 'Contraseña' (Password) and 'Azure Key Vault'. At the bottom, there are buttons for 'Crear' (Create), 'Atrás' (Back), 'Prueba de conexión' (Test connection), and 'Cancelar' (Cancel). A green checkmark indicates 'Conexión correcta' (Correct connection).

The screenshot shows the Microsoft Azure Data Factory portal. The left sidebar is the same as the previous screenshot. The main content area is titled 'Entornos de ejecución de integración' (Integration execution environments) and shows a table with the following data:

Nombre	Tipo	Subtipo	Estado	Relacionado	Región	Versión
AutoResolveIntegrationR...	Azure	Pública	En ejecución	0	Resolver automáticamente	---
IRPostgreSQL	Autohospedado	---	En ejecución	1	---	5.21.8273.2

Propuesta_Pri x | Untitled docu x | Linked service x | Copy and tran x | Azure Data Fa x | arkondatfacti x | arkondatfacti x | #8. Azure Dat x | (1) WhatsApp x | +

adf.azure.com/es/management/datalinkedservices?factory=%2Fsubscriptions%2Fc9a2239-634a-4a9e-9e95-27c946875d20%2FresourceGroups%2Farkon_prueba%2Fproviders...

OACs Cursos Importantes Click Up Arkon Abraxas Group - Ca... Medium Python Directorio de Servic... Knights-Watch GitHub Hector iETL con Azure Dat... Azure DIAGRAM

Microsoft Azure | Data Factory | arkondatfactory

Validar todo Publicar todo

Connections

- Linked services
- Integration runtimes
- Microsoft Purview

Source control

- Git configuration
- ARM template

Author

- Triggers
- Global parameters
- Data flow libraries

Security

- Credentials
- Clave administrada por el d...
- Puntos de conexión privado...

Servicios vinculados

El servicio vinculado define la información de conexión a un almacén de datos o un proceso. Más información

+ Nuevo

Filtrar por nombre Anotaciones: Cualquiera

Mostrando elementos del 1 al 1 de un total de 1

Nombre	Tipo	Relacionado	Anotaciones
PostgreSQL	PostgreSQL	1	

25°C Nublado

Data Pipeline Creation

Final

Propue x | Azure x | arkond x | arkond x | synaps x | Azure x | Untitle x | Whats x | Quick x | Conne x | Crea x | Enfor x | +

adf.azure.com/es/monitoring/pipeline/runs?factory=%2Fsubscriptions%2Fc9a2239-634a-4a9e-9e95-27c946875d20%2FresourceGroups%2Farkon_prueba%2Fproviders%2...

OACs Cursos Importantes Click Up Arkon Abraxas Group - Ca... Medium Python Directorio de Servic... Knights-Watch GitHub Hector iETL con Azure Dat... Azure DIAGRAM

Microsoft Azure | Data Factory | arkondatfactory

Dashboard

- Runs
- Pipeline runs
- Trigger runs
- Runtimes & sessions
- Integration runtimes
- Data flow debug
- Notifications
- Alerts & metrics

Ejecuciones de canalización

Desencadenado Depurar Volver a ejecutar Cancelar opciones Actualizar Editar columnas Lista Gantt

Filtrar por id. o nombre de... Hora local: Últimas 24 horas Nombre de canalización: PostgreSQL to Synapse 2 Estado: Todo

Ejecuciones: Ejecuciones más recientes Desencadenadas por: Todo Agregar filtro

Mostrando elementos del 1 al 1

Nombre de canalización	Inicio de la ejecución	Fin de la ejecución	Duración	Desencadenadas por	Estado	Error	Ejecutar	Parámetros	Anotaci
PostgreSQL to Synapse 2	Oct 5, 2022, 12:49:38 pm	--	00:00:14	Desencadenador ma...	En curso		Original		

Última actualización hace 0 minutos

25°C Nublado

