info@stratebi.com





info@stratebi.com



CONTENIDO

1.	INTRODUCCIÓN	3
 (Requisitos Previos Creación de los Recursos	. 3
2.	PLANTEAMIENTO Y ARQUITECTURA	5
3.	ELTL CON PIPELINES	8
4.	POOLS SQL	16
5.	VISUALIZACIÓN CON POWER BI	24
6.	AUTOMATIZACIÓN CON TRIGGERS	31
7.	CONCLUSIONES	33

info@stratebi.com



1. INTRODUCCIÓN

Synapse Analytics es un servicio cloud de Azure que proporciona una serie de funcionalidades que permiten la realización de un análisis completo de los datos. Para ello, Synapse cuenta con varias capacidades funcionales y técnicas de las que podemos destacar los pipelines para implementar procesos ETL, ELT o ELTL, los notebooks de Synapse o los dataflows para realizar transformaciones y asegurar la calidad del dato, los pools de SQL para almacenar los datos a analizar a modo de Data Warehouse, entre otras muchas funcionalidades.



Figura O. Arquitectura de Synapse Analytics

Toda esta suite de utilidades permite a los usuarios realizar un proceso de analítica de datos con perspectiva end-to-end. De esta forma, una vez que el usuario haya identificado las fuentes de datos que quiere analizar, puede emplear el servicio de Synapse Analytics para realizar la extracción del origen de datos, transformar estos datos con notebooks o dataflows y almacenarlos en un destino como una base de datos Azure SQL o un pool SQL dedicado de Synapse. Nos ha parecido interesante crear un ejemplo con Azure Synapse, pues se consolida como pieza clave en las arquitecturas modernas de datos basadas en Azure. A lo largo de este artículo, daremos un paseo por Synapse Analytics con una visión end-to-end, implementando procesos ELTL (extracción, carga, transformación y carga), hasta explotar el dato con Power BI.

Requisitos Previos

Para poder implementar procesos en Synapse Analytics es necesario cumplir varios requisitos:

- Es necesario tener una suscripción de <u>Azure</u> para poder crear los recursos.
- Necesitamos tener un grupo de recursos en el que crearemos los recursos necesarios de Azure.
- Si queremos visualizar los datos una vez terminado el proceso ELTL, debemos tener instalado <u>Power BI Desktop</u>.

info@stratebi.com

strate

Creación de los Recursos

Como se ha mencionado previamente, es necesario que tengamos un grupo de recursos creado en el portal de Azure donde podamos crear el resto de los recursos necesarios. Una vez definido el grupo, debemos crear los siguientes recursos:

Name ↑↓	Type ↑↓
DedicatedDemo (stsynapsedemosv1/DedicatedDemo)	Dedicated SQL pool
stdatalakedemosv1	Storage account
Stkeyvaultdemosv1	Key vault
stsparkpooldemo (stsynapsedemosv1/stsparkpooldemo)	Apache Spark pool
stsynapsedemosv1	Synapse workspace

Figura 1. Recursos necesarios en Azure

Se pueden distinguir varios tipos de recursos entre los cuales el principal es el Synapse workspace. Este recurso será el eje principal de todo el proceso pues gestiona todas las tareas del proceso ELTL y proporciona el acceso a los datos.

Además del propio recurso de Synapse, es necesario crear previamente un ADLS Gen 2 pues es el sistema de almacenamiento soportado por Synapse y es obligatorio emplearlo cuando queremos crear un workspace de Synapse.

En cuanto a los pools de SQL y Spark, estos son creados desde Synapse por dos motivos:

- El pool de Spark se crea para otorgar la capacidad de procesamiento a los notebooks y a los dataflows. Se le pueden asignar diferentes capacidades de cómputo e incluso la habilidad de escalar entre varios nodos.
- El pool de SQL dedicado se utiliza como destino en nuestro proceso ELTL. Este pool hace las funciones de Data Warehouse y, al ser de tipo "Dedicado", solo incurre en costes cuando lo tenemos encendido.

Por último, se puede observar que se ha creado un Key Vault para almacenar algunos secretos que usa el pipeline.

info@stratebi.com



2. PLANTEAMIENTO Y ARQUITECTURA

En el ejercicio que se plantea, se desea recabar información sobre el histórico de meses anteriores para poder analizar la evolución de los precios de la luz y de los combustibles, así como ver en qué grado afectan los valores climatológicos en estas variables.

Para ello, se emplean varias APIs públicas como orígenes de datos. Estas APIs son las siguientes:

- **Red Eléctrica de España**: permite la obtención de datos sobre el precio de la electricidad, la demanda y generación energética entre otros muchos datos.
- **AEMET**: proporciona datos sobre las mediciones de los valores climatológicos diarios como la temperatura o las precipitaciones, para cada una de sus estaciones meteorológicas.
- Datos.gob.es: este portal de datos abiertos pertenece al Gobierno de España y en él se proporcionan diferentes mediciones y datos proporcionados por los principales organismos públicos. En este caso, se ha empleado una de las APIs del Ministerio de Transición Ecológica y el Reto Demográfico que permite acceder a los datos de los precios de los combustibles en las estaciones de servicio españolas.

Una vez se ha expuesto el planteamiento inicial de nuestro proceso ELTL, se presenta la arquitectura del mismo en el siguiente diagrama:

www.stratebi.com

91.788.34.10



info@stratebi.com



Figura 2. Arquitectura del proceso ELTL con Synapse

info@stratebi.com



Como se puede observar en la figura 2, en la parte de la izquierda se encuentran los orígenes de datos del proceso ELTL. Para realizar la fase de extracción de datos, se emplean pipelines de Synapse que realizan diversas llamadas a los endpoints de las APIs.

Cabe destacar que una vez se han extraído los datos del origen, se almacenan en la capa raw del ADLS Gen 2 que está conectado con el recurso de Synapse. El ADLS se puede estructurar de diferentes formas, sin embargo, gracias a la estructura de almacenamiento de ficheros jerárquica introducida en la Gen 2 de los data lakes, podemos crear varias capas en función del estado del dato. Esto se materializa en la división del data lake en tres capas:

- Raw: contiene los datos en el formato de origen.
- Clean: contiene los datos tras aplicar transformaciones y un proceso de limpieza.
- Delivery: contiene los datos que van a utilizar las aplicaciones.

Una vez se han extraído los datos y se han almacenado en la capa raw, se proceden a realizar las transformaciones con los notebooks de Synapse y se almacenan los datos en las capas clean y delivery.

En este caso, dado que los datos se encuentran a diferentes niveles de detalle (algunos se proporcionan por las APIs a nivel diario y otros a nivel horario), se ha decidido tener un pool de SQL dedicado donde almacenar los datos agregados al mismo nivel de detalle y usar este en vez de la capa delivery para que lo utilice nuestro informe de Power BI. Para realizar esta agregación, se han empleado transformaciones mediante dataflows que toman como origen la capa delivery y como destino el pool SQL dedicado.

Una vez se ha creado el Data Warehouse con el pool SQL dedicado, podemos usar el endpoint de Synapse para importar estos datos desde Power BI y realizar los informes correspondientes.

info@stratebi.com



3. **ELTL CON PIPELINES**

La implementación de los procesos ELTL con Synapse se materializa con los pipelines. Estos pipelines son una serie de actividades encadenadas una tras otra. En el ejercicio planteado, se pueden observar que hay tres fuentes diferentes y por ello, se ha decidido dividir el proceso ELTL en tres pipelines diferentes, uno para cada fuente. Cabe destacar que sobre cada origen de datos se tiene una forma diferente de realizar las peticiones. Sin embargo, todas las APIs admiten peticiones basadas en rangos de fechas por lo que, para facilitar el proceso de creación de pipelines, se van a crear unos parámetros que serán comunes en todos ellos.

Name	Туре	Default value	
year	String ~	2022	Ŵ
month	String ~	01	Ŵ

Figura 3. Parámetros de los pipelines

En primer lugar, trataremos el pipeline que extrae información de la API de la Red Eléctrica de España.





Como se puede apreciar en la figura 4, el pipeline comienza estableciendo una serie de variables. Estas variables son necesarias para establecer el rango de fechas del que se quiere extraer la información y son requisito obligatorio para las peticiones de la API de la REE. Estos valores se calculan teniendo en cuenta el parámetro que se ha mencionado previamente.

info@stratebi.com



Una vez se ha definido el rango de fechas que se quiere extraer, se ejecutan una serie de pipelines que tienen como objetivo hacer varias peticiones según la información que se quiera obtener. En este caso, se quiere extraer información de:

- Precios (€/MWh): Spot y PVPC. Son los dos tipos de precios de la electricidad del sistema eléctrico español:
 - El precio spot de electricidad es uno de los principales indicadores del mercado eléctrico y señal de su condición de adaptación entre oferta y demanda. Refleja el coste de suministrar un kilo watt hora (kWh) adicional al sistema eléctrico.
 - El PVPC (Precio de Voluntario Pequeño Consumidor) es una tarifa de discriminación horaria y precio variable. Es decir, el precio del kWh de electricidad cambia para cada hora de cada día. Por lo tanto, su precio está sometido a la volatibilidad del mercado eléctrico, según la oferta y la demanda entre las compañías generadoras de energía y las comercializadoras que venden esa energía al consumidor.
- Demanda (MWh): Real, programada y prevista.
- Generación (MWh): Datos de generación de energías renovables y no renovables.
- Intercambios (MWh): Importaciones y exportaciones de España con los países con los que existe una frontera física (Marruecos, Portugal, Francia y Andorra).

Como se puede ver en el flujo, cada uno de estos pipelines se ejecuta en paralelo, lo que permite realizar solicitudes simultáneas a la API en cuestión. Todos estos pipelines tienen una estructura similar, que se presenta en la siguiente imagen:



Figura 5. Pipeline precios de la electricidad

El pipeline se divide en dos actividades:

• Copy data: como su propio nombre indica, consiste en copiar datos de una fuente de origen a un destino. En este caso, la fuente de origen es la API de la Red Eléctrica y la fuente de destino es la capa raw de nuestro datalake.

strate

91.788.34.10

info@stratebi.com

Source dataset *	Input_REE	V Open + New 60 Preview data
	✓ Dataset properties (D
	Name	Value
	category	mercados
	lang	es
	start_date	@pipeline().parameters.start_date
	end_date	@pipeline().parameters.end_date
	widget	precios-mercados-tiempo-real
	time_trunc	hour

Figura 6. Origen de la actividad copy data

Sobre la imagen anterior, hay varios factores a tener en cuenta como el método HTTP que se emplea (GET) o el dataset que se está empleando. El dataset funciona como una interfaz entre nuestra actividad del pipeline y los datos. En este caso, el dataset está parametrizado, lo que permite usar el mismo dataset para hacer el resto de las llamadas a la API. Desde Synapse, podemos crear los integration datasets desde la pestaña "Data". En la siguiente imagen se muestra la configuración realizada para este dataset en particular:

REST REST Input_REE			REST REST	
Connection Parameters		_	Connection Parameters	
Linked service *	• API_REE ~	🖉 🎜 Test connection 🖉 Edit	+ New Delete	
Integration runtime *	AutoResolveIntegrationRuntime ~	🖉 Edit	Name	Туре
Base URL	https://apidatos.ree.es/]	category	String ~
Relative URL	@concat(dataset().lang, '/datos/', data	ා ම Preview data	lang	String ~
			start_date	String ~
			end_date	String ~
			widget	String ~
			time_trunc	String ~

Figura 7. Configuración del dataset

Este dataset es de tipo REST ya que es el servicio que vamos a emplear en la consulta a la API. Los parámetros que se han establecido son los que están definidos en la <u>documentación de la API</u>, y se emplean en la URL relativa para usar este mismo dataset en todas las peticiones a la REE.

info@stratebi.com



Por otra parte, en la actividad de copia también hay que definir el destino de los datos extraídos. Como mencionamos previamente, los datos se almacenan en primera instancia en la capa raw del ADLS Gen 2.

General	Source	Sink	Mapping	Settings	User properties			
Sink datas	set *	[📓 RawOutp	ut_Electricida	d_JSON ~	🖉 Open	+ New	Learn more 🖸
			✓ Dataset pro	operties 🕕				
			Name		Value			
			fileName		@concat(forma	tDateTime(p	ipeline().p	
			directory		@concat('raw/e	electricidad/p	precios/', fo	

Figura 8. Destino de la actividad copy data

En el caso del dataset destino también está parametrizado para poder mantener la estructura jerárquica vista en la figura 3, de forma que el nombre del fichero tenga el siguiente formato [año]-[mes].json y el directorio en el que se almacena también está parametrizado para tener separados los años en diferentes carpetas. El formato de los ficheros es de tipo JSON ya que en la capa raw se almacenan los datos tal y como vienen del origen. De esta forma se puede mantener una trazabilidad y permite detectar el origen de los posibles errores cometidos en las diferentes transformaciones.

Los datasets de Synapse requieren de un conector para su creación. Estos conectores se denominan linked services y existen varios tipos.

Name ↑↓	Type ↑↓
API_REE	REST
AzureKeyVault	Azure Key Vault
name DedicatedSynapse	Azure Synapse Analytics
stsynapsedemosv1-WorkspaceDefaultSqlServer	Azure Synapse Analytics
stsynapsedemosv1-WorkspaceDefaultStorage	Azure Data Lake Storage Gen2

Figura 9. Linked Services

info@stratebi.com



 Notebook. Los notebooks de Synapse siguen la misma filosofía que los Jupyter notebooks mezclando celdas de texto y celdas de código. En las celdas de código se pueden emplear varios lenguajes como Scala, SQL o PySpark, siendo este último el empleado en todo el ejercicio. A través de Python se pueden leer los datos de la capa raw, almacenarlos temporalmente en una estructura de datos como un dataframe, realizar las transformaciones pertinentes y posteriormente guardarlo en la capa clean o delivery del datalake.

El procedimiento es el mismo para el resto de pipelines que obtienen información de electricidad con la diferencia de que se deben cambiar los parámetros para hacer la llamada pertinente en la API.

Por último, una vez que se han almacenado los datos en las correspondientes capas, se deben agregar los datos para que estén al mismo nivel de detalle, es decir, transformar los datos para que estén a nivel diario. Si recordamos la figura 4, podemos ver que la última actividad consiste en un dataflow el cual veremos en detalle posteriormente. Este dataflow realiza una agregación para que los datos estén a nivel diario como mencionamos previamente.



Figura 10. Dataflow electricidad

En este dataflow, se utilizan varios orígenes de datos, que se corresponden con los diferentes datasets creados para los ficheros almacenados en la capa de delivery. Estos ficheros están almacenados en formato parquet. Para poder trabajar al mismo nivel de detalle, hay que realizar una agregación a los precios de la electricidad y a la demanda energética, ya que la API proporciona estos datos a nivel horario mientras que el resto está a nivel diario.

Una vez se han agregado a nivel diario, se pueden unir todos los ficheros con el fin de crear una sola tabla que contenga todos los datos de la electricidad. Esto depende del modelo que se quiera implementar posteriormente podría decidirse hacer una tabla para los precios, otra para la generación, etc. pero para este ejercicio, por simplicidad, se ha decidido usar una sola

info@stratebi.com



tabla que será una de las tablas de hechos de nuestro modelo final. Para realizar estas uniones, se utiliza la operación INNER JOIN que une los diferentes ficheros.

Por otra parte, también se realiza una operación de SELECT para eliminar las columnas duplicadas y renombrarlas, y se realiza un ALTER ROW para poder realizar una operación de tipo UPSERT en la base de datos de destino. Esto permite insertar los datos que no existían en la base de datos y actualizar los datos ya existentes.

Por último, se ha de destacar que el destino es una tabla perteneciente al pool SQL dedicado. Por ello, es necesario crear previamente la tabla con un script SQL.

```
CREATE TABLE [electricity]
(
    [year] [INT] NOT NULL,
    [month] [INT] NOT NULL,
    [day] [INT] NOT NULL,
    [real_demand] [FLOAT],
    [programmed_demand] [FLOAT],
    ...
)
```

A partir de lo anterior, también se deben crear los pipelines para extraer información sobre las condiciones climatológicas y sobre los combustibles. Estos dos pipelines son muy similares entre sí pues tienen las mismas actividades. Además, en ambos casos solo se realiza una petición a la API ya que, a diferencia de los datos proporcionados por la REE, toda la información que nos interesa se puede recabar con una sola petición. A continuación, los veremos con más detalle:



Figura 11. Pipeline de meteorología



Figura 12. Pipeline de combustibles

info@stratebi.com



Como se puede observar en las figuras 11 y 12, los pipelines constan tan solo de dos actividades que realizan el proceso ELTL.

- Notebook. En este caso, se ha preferido hacer la extracción directamente desde el notebook de Synapse debido a los formatos que devuelven las APIs. Dentro de estos notebooks, además de realizar las peticiones a las diferentes APIs, también se realizan diversas transformaciones para adaptar el formato de los datos correctamente. Además, la información se almacena en las diferentes capas del datalake según el estado del dato.
- Data Flow. De manera similar a lo que ocurre con la electricidad, los datos extraídos de las APIs tienen un nivel de detalle horario





_	FuelPrices	-	AggregatePrices		AlterRow		CoadIntoSQLPool
2	Import data from FuelParquet	+	Aggregating data by 'Provincia, year, month, day' producing columns 'biodiesel, bioetanol, gas_natural_comprimido,	+	Add expressions to alter rows	+	Export data to FuelPoolTable



Como se puede observar en las imágenes, tras obtener los datos de la capa delivery en formato parquet, se agregan los datos a nivel diario. Además, para que se pueda realizar la operación de UPSERT sobre la base de datos, se inserta un paso intermedio denominado ALTER ROW. Este paso permite actualizar las filas en caso de que se modifiquen los valores para una determinada clave. En este caso, la clave está formada por la fecha y por la provincia ya que, a diferencia de los valores de la electricidad que son a nivel nacional, los valores del precio del combustible y los datos climatológicos son diferentes dependiendo de la provincia. Por último, se define la base de datos destino, el cuál es un pool SQL dedicado. Al igual que en el dataflow de la electricidad, es necesario que se creen las tablas con sentencias SQL previo a la ejecución del dataflow.

Una vez que se han creado los tres pipelines, podemos crear otro pipeline adicional que los ejecute en paralelo.

info@stratebi.com





Figura 15. Pipeline ELTL

Este pipeline está parametrizado de forma que el usuario puede introducir el mes y el año del que quiere obtener el histórico de datos. Cada pipeline ajusta internamente los valores de los parámetros para adaptarlos al formato exigido por cada una de las APIs.

info@stratebi.com

4. **POOLS SQL**

A lo largo de la implementación hemos hecho referencia al pool SQL dedicado, pero ¿qué es realmente un pool SQL y por qué es dedicado?

Los pools SQL son el servicio cloud que proporciona Azure que cumple con las labores de un Data Warehouse tradicional. Los Data Warehouse son plataformas destinadas a almacenar datos que permiten realizar el análisis de estos. Son un elemento fundamental dentro del mundo de Data Analytics y del Business Intelligence, y los pools SQL de Synapse son la solución cloud que proporciona Microsoft.



Figura 16. Solución Big Data end-to-end

Dentro de las opciones que proporciona Synapse a la hora de crear los pools SQL, encontramos dos categorías:

 Pools SQL dedicados. Funcionan como un servidor de base de datos SQL tradicional, pero en la nube. Se puede definir su potencia y rendimiento a la hora de crearlo, aunque hay que tener en cuenta que el coste variará según la configuración que se realice.

Dedicated SQL pool details			Dedicated SQL pool details		
Name your dedicated SQL pool and choose its initial settings.			Name your dedicated SQL pool and choose its initial settings.		
Dedicated SQL pool name *	Enter dedicated SQL pool name		Dedicated SQL pool name *	Enter dedicated SQL pool name	
Performance level ①	0	DW100c	Performance level ①	O	DW1000c
Estimated price ③	Est. cost per hour 127 EUR View pricing details		Estimated price ①	Est. cost per hour 12.73 EUR View pricing details	



info@stratebi.com



Figura 17. Comparación de precios pool SQL dedicado

Como se puede observar, si se incrementa el rendimiento del pool SQL dedicado, el coste por hora también aumenta. El usuario tiene que decidir qué nivel de rendimiento va a requerir dependiendo del rendimiento que necesite. El rendimiento se mide en DWU que representan una combinación de ciertos recursos como CPU, memoria o velocidad de escritura y lectura. Puede obtener más información <u>aquí</u>.

 Pools SQL serverless. Los pools sin servidor o serverless son más sencillos de cara al usuario ya que los costes se producen por TB de datos procesado y no hay costes por horas a diferencia del pool dedicado. Sin embargo, a diferencia del pool SQL dedicado que permite trabajar con tablas relacionales, el pool serverless está enfocado en realizar consultas a los ficheros del datalake.

En nuestro caso, utilizaremos el pool SQL dedicado ya que queremos tener un modelo relacional para acceder a él desde Power BI y generar nuestro informe. De todas maneras, Synapse crea por defecto un pool serverless que permite explorar los ficheros del datalake.

Name	Туре	Status	Size
Built-in	Serverless	📀 Online	Auto
DedicatedDemo	Dedicated	Paused	DW100c

Figura 18. Pools SQL

Como hemos mencionado, el pool serverless nos permite realizar consultas a nuestros ficheros del datalake. Vamos a realizar una prueba con los datos extraídos del precio de la luz.

🗘 / 🗞 master branch 🗸 🗸 Validate all 🔶 Commit all	1 Publish
🖨 demos 🛛 🕹	
📃 New SQL script 🗸 📋 New notebook 🗸 🚯 New data flow	🗮 New integration dataset 🛛 ក Upload
$\leftarrow \rightarrow \lor \uparrow ~$ demos > delivery > electricidad > precios > 2022	
Name	 Last Modified
2022-01.parquet	27/5/2022, 10:25:29
2022-02.parquet	6/6/2022, 13:02:35
2022-03.parquet	8/6/2022, 11:16:34
2022-04.parquet	8/6/2022, 11:25:49
D 2022-05.parquet	15/6/2022, 13:58:50



www.stratebi.com

91.788.34.10



info@stratebi.com

Si seleccionamos uno de los ficheros y creamos un script SQL, podemos hacer una operación de SELECT para ver los datos que hay en dicho archivo.



Figura 20. Consulta SQL

Podemos observar que Synapse está utilizando el pool serverless (Built-in) para realizar la operación de consulta. Una vez ejecutamos la operación, podemos observar los resultados en formato de tabla e incluso como un gráfico.

Results Messages							
View Table Chart	\mapsto Export results \vee						
✓ Search							
рурс	spot	datetime	year	month	day		
257.91	206	2022-05-01T00:00:00.000+02:00	2022	05	01		
255.88	199.7	2022-05-01T01:00:00.000+02:00	2022	05	01		
257.39	197.56	2022-05-01T02:00:00.000+02:00	2022	05	01		
256.58	192.2	2022-05-01T03:00:00.000+02:00	2022	05	01		
264.22	197.43	2022-05-01T04:00:00.000+02:00	2022	05	01		
269.51	201.32	2022-05-01T05:00:00.000+02:00	2022	05	01		
261.56	195	2022-05-01T06:00:00.000+02:00	2022	05	01		
262.05	194.95	2022-05-01T07:00:00.000+02:00	2022	05	01		

Figura 21. Resultados en formato de tabla



Figura 22. Resultados en formato gráfica

```
www.stratebi.com
```

info@stratebi.com

La gráfica que proporciona Synapse es interactiva, pudiendo añadir o quitar las series que queremos mostrar e incluso permite destacar los valores cuando pasamos el cursor por encima como se muestra en el siguiente ejemplo.



Figura 23. Filtrando la gráfica

De esta forma, podemos realizar un análisis exhaustivo de nuestros datos. Además de todo esto, el pool SQL sin servidor nos permite analizar no solo uno sino varios ficheros de una carpeta. Esto nos permite, por ejemplo, realizar una consulta para obtener todos los datos de los precios de la electricidad del 2022 que están almacenados en la carpeta 2022. Para ello, debemos modificar ligeramente la consulta SQL.



Figura 24. Consulta SQL para todos los ficheros de una carpeta

Como se ha podido ver en los resultados, los datos de los precios de la luz se mantienen a nivel horario en el datalake. Sin embargo, para hacer nuestro informe de Power BI queremos que todos los datos estén a nivel horario. Para ello, mediante los dataflows se agregan estos datos y se almacenan en el pool dedicado.

El pool dedicado está enfocado en trabajar con tablas SQL y no con los ficheros del ADLS. Por este motivo, este pool permite la realización de operaciones CRUD sobre las tablas SQL. Existe una gran diferencia en el funcionamiento de este pool respecto al pool serverless, ya que cuando estamos empleando el pool SQL dedicado, tenemos que encenderlo y apagarlo cuando





info@stratebi.com

dejemos de usarlo. Esto lo podemos hacer desde el portal de Azure o desde el workspace de Synapse, salvo que su inicio y pausa se automatice dentro del pipeline

r	DedicatedDemo (stsynapsedemosv1/DedicatedDemo) &									
»	Pause 🗹 Scale	🔊 Restore 🕂 New restore point 📋 Delete 🛛 🗹 Open in Synapse Studio								
	↑ Essentials									
	Resource group (move)	: <u>stratebi-demo-v1</u>								
	Status	: Online								
	Location	: West Europe								

Figura 25. Pool SQL dedicado desde el portal de Azure

Name		Туре	Status	Size
Built-in		Serverless	🕑 Online	Auto
DedicatedDemo	11 🗹 🖓	Dedicated	🕑 Online	DW100c

Figura 26. Pool SQL dedicado desde Synapse

Una vez encendido el pool dedicado, podemos empezar a utilizarlo. En la siguiente imagen, se puede observar cómo están las tablas creadas con los datos obtenidos tras haber realizado el proceso ELTL con los pipelines.

Data	+ *	* **
Workspace	Linked	
\bigtriangledown Filter resources by name		
Lake database		1
 SQL database 		1
🔺 🔒 DedicatedDemo (SC)L)	
Tables		
👌 🖪 dbo.electricit	у	
🕨 📰 dbo.fuel		
🕨 🔚 dbo.weather		
External tables		
External resource	es	
▷ 🖹 Views		
Programmability		
Schemas		
Security		

info@stratebi.com



Figura 27. Pool SQL dedicado

Al igual que con el pool serverless, se pueden realizar consultas sobre las diferentes tablas:



Figura 28. Operaciones con el pool SQL dedicado

Hemos visto cómo podemos emplear los diferentes tipos de pools SQL que ofrece Azure con Synapse. Sin embargo, el detalle de tener que encender y apagar nuestro pool dedicado cuando hacemos el proceso ELTL con los pipelines puede llegarnos a causar que incurramos en costes innecesarios por no apagar el pool cuando terminan de ejecutar las actividades del pipeline. Para ello, vamos a realizar algunas modificaciones en nuestro pipeline para que se inicie el pool antes de ejecutar el ELTL y para que se pause una vez se haya terminado.



Figura 29. Pipeline ELTL con control del pool SQL

Como podemos observar en la imagen, hay una actividad que ejecuta el pipeline que realiza los tres procesos ELTL de las diferentes APIs, y otras dos actividades de tipo "Web" que controlan el pool SQL. Estas actividades realizan llamadas a la API REST de Azure.



info@stratebi.com

La URL que debemos emplear para iniciar el pool SQL es la siguiente:

https://management.azure.com/subscriptions/{subscriptionid}/resourceGroups/{resource-groupname}/providers/Microsoft.Synapse/workspaces/{workspacename}/sqlPools/{database-name}/resume?api-version=2019-06-01-preview

Para la actividad de parar el pool SQL es similar:

https://management.azure.com/subscriptions/{subscriptionid}/resourceGroups/{resource-groupname}/providers/Microsoft.Synapse/workspaces/{workspacename}/sqlPools/{database-name}/pause?api-version=2019-06-01-preview

Además, debemos realizar algunas configuraciones en la actividad para que Synapse pueda tener permisos suficientes a la hora de modificar el estado del pool SQL dedicado.

General Settings Use	er properties
URL	@concat('https://management.azure.c
Method *	POST ~
Headers	+ New
Body	Resume SQL pool
	1
Datasets	+ Add dataset reference
Linked services	+ Add linked service reference
Integration runtime * ①	📀 AutoResolveIntegrationRuntime 🛛 🖉 Edit
Disable certificate validation	
HTTP request timeout \bigcirc	00:01:00
Authentication	System Assigned Managed Identity \sim
Resource *	https://management.azure.com

info@stratebi.com



Figura 30. Configuración de las actividades web

También se deben asignar permisos de "Contributor" desde el portal de Azure al propio service principal de Synapse. Esto se puede realizar en el portal web en el apartado de "Access Control (IAM)":

Synapse workspace	sv1
✓ Search (Ctrl+/)	«
S Overview	^
Activity log	
Å Access control (IAM)	
🗳 Tags	
Diagnose and solve problems	

stsynapsedemosv1 assignments - stsynapsedemosv1

Assignments for the selected user, group, service principal, or managed identity at this scope or inherited to this scope.

ho Search by assignment name or description								
Role assignments (2) (i)	Description	C	Commenciation and	Condition				
Kole	Description	Scope	Group assignment	Condition				
Contributor	Grants full access to manage all re	This resource		None				

Figura 31. Permisos del workspace de Synapse en el portal de Azure

Una vez otorgados los permisos, se pueden ejecutar correctamente las actividades web.

info@stratebi.com



5. VISUALIZACIÓN CON POWER BI

Una vez se han definido los pipelines asociados al proceso ELTL, podemos visualizar los datos mediante un informe de Power BI. Para ello, desde Power BI Desktop, seleccionamos la opción de "Obtener datos" y buscamos la opción de Synapse:

Buscar	Azure	
Todo	Azure SQL Database	~
Archivo	Azure Synapse Analytics SQL	
Base de datos	Permite importar datos de Azuro Apalveis Sonicos Permite importar datos de Azure Synapse Analytics SQL	
Power Platform	azure Database for PostgreSQL	
Azure	Azure Blob Storage	
Servicios en línea	Almacenamiento de tablas de Azure	
Otras	🧝 Azure Cosmos DB	
	💥 Azure Data Explorer (Kusto)	
	Azure Data Lake Storage Gen2	
	Azure Data Lake Storage Gen1	
	🔹 Azure HDInsight (HDFS)	
	Azure HDInsight Spark	
	🧬 HDInsight Interactive Query	
	Azure Cost Management	
	🤪 Azure Databricks	
	Área de trabajo de Azure Synapse Analytics (beta)	•

Obtener datos

Figura 32. Orígenes de datos en Power BI

Una vez hemos elegido esta opción, previo a seguir con la configuración del origen de datos, es necesario que se inicie el pool SQL dedicado ya que, de lo contrario, la conexión entre Power BI Desktop y Synapse será fallida. Con el pool iniciado, se debe configurar el acceso al servidor desde Power BI.

stsvna	r () psedemosv1.sql.azuresvnapse.net
Base de	datos (opcional)
Dedica	atedDemo
Modo (Conectividad de datos 🕞
Impo	ortar
O Direc	ctQuery
> Opcio	nes avanzadas

Base de datos SQL Server

info@stratebi.com



Figura 33. Configuración conexión Synapse SQL

La dirección del servidor se puede localizar desde el portal de Azure:



Figura 34. Configuración del endpoint de Synapse

En este caso, debemos utilizar el endpoint dedicado, sobre el que se han creado todas las tablas resultantes del proceso ELTL. Tras identificarnos con nuestro usuario, podremos cargar las tablas de nuestro modelo:

	Q	electricity	/			
pciones de presentación 💌	Ŀ	Vista previa	descargada	el miércoles		
stsynapsedemosy1.sgl.azuresynapse.net; DedicatedDemo [3]		year	month	day	real_demand	pro
		2022	2	16	4305316	
		2022	1	3	3926446	
		2022	4	6	4230731	
✓ 🔠 weather		2022	2	5	3863123	
		2022	1	29	3997083	
		2022	2	2	4368309	
		2022	3	26	3609222	
		2022	4	2	3642430	
		2022	4	9	3524936	
		2022	4	20	3909975	
		2022	3	7	4173965	
		2022	2	10	4340698	
		2022	1	13	4478691	
		1 Los da a límite	tos de la vist es de tamañ	a previa se h	ian truncado debio	lo

info@stratebi.com



Figura 35. Tablas de Synapse en Power BI

Una vez hemos cargado los datos, podemos realizar las transformaciones que consideremos pertinentes para nuestro modelo como la obtención de dimensiones para nuestra tabla de hecho. Es importante tener en cuenta que, dado que el modo de carga establecido para Power BI es "Import", una vez se hayan cargado los datos, se puede apagar el pool SQL dedicado para no incurrir en costes innecesarios.

Con el modelo preparado, podemos realizar nuestro informe. A modo demostrativo, se muestran los siguientes ejemplos:



Figura 36. Análisis electricidad en Power BI



info@stratebi.com

Combustib	les					Provir Todas	ncia ~	All	Last 3
iasoleo por Date Gasoleo A © Gasoleo .0 .5 .5 .0 .0 .0 .0 .0 .0 .0 .0 .0 .0	B • Gasoleo	premium	,	Gasolina pr 2,0	or Date	ſ		Gaso Gaso Gaso Gaso Gaso	ilina 98 E5 ilina 98 E10 ilina 95 E5 Pre ilina 95 E10 ilina 95 E5
ene 2022 f	feb 2022	mar 20	22 abr 20	1,4	feb 2022 n	nar 2022	abr 2022		
ene 2022 1 Fop 10 provincias co Provincia	feb 2022 In gasóleo A Gasoleo A	mar 20 A más barat Gasoleo B	22 abr 20 to Gasoleo premium	22 1,4 ene 2022 Top 10 pro Provincia	feb 2022 n vincias con gase Gasolina 98	nar 2022 Dina 98 E5 Gasolina	abr 2022 más barata Gasolina 95	Gasolina	Gasolina
ene 2022 1 Top 10 provincias co Provincia MELILLA	feb 2022 In gasóleo A Gasoleo A 1,24	mar 20 A más barat Gasoleo B	22 abr 20 to Gasoleo premium 1,28	22 1,4 ene 2022 Top 10 pro Provincia	feb 2022 n vincias con gase Gasolina 98 E5	olina 98 E5 Gasolina 98 E10	abr 2022 más barata Gasolina 95 E5 Prem.	Gasolina 95 E10	Gasolina 95 E5
ene 2022 1 op 10 provincias co Provincia MELILLA LAS PALMAS	feb 2022 In gasóleo A Gasoleo A 1,24 1,25	mar 20 A más barat Gasoleo B	22 abr 20 to Gasoleo premium 1,28 1,36	22 1.4 ene 2022 Top 10 pro Provincia CEUTA	feb 2022 n vincias con gase Gasolina 98 E5 1,40	aar 2022 blina 98 E5 Gasolina 98 E10	abr 2022 más barata Gasolina 95 E5 Prem. 1,41	Gasolina 95 E10	Gasolina 95 E5 1,37
ene 2022 1 Top 10 provincias co Provincia MELILLA LAS PALMAS STA, CRUZ DE TENERIFE	feb 2022 In gasóleo A Gasoleo A 1,24 1,25 1,28	mar 20 A más barat Gasoleo B	22 abr 20 to Gasoleo premium 1,28 1,36 1,39	22 1,4 ene 2022 Top 10 pro Provincia CEUTA LAS PALMAS	reb 2022 n vincias con gase Gasolina 98 5 1,40 1,43	olina 98 E5 Gasolina 98 E10	abr 2022 más barata Gasolina 95 E5 Prem. 1,41 1,31	Gasolina 95 E10	Gasolina 95 E5 1,37 1,31
ene 2022 i op 10 provincias co Provincia MELILLA AS PALMAS STA. CRUZ DE TENERIFE CEUTA	feb 2022 In gasóleo A Gasoleo A 1,24 1,25 1,28 1,34	mar 20 A más bara Gasoleo B	22 abr 20 to Gasoleo premium 1,28 1,36 1,39 1,38	22 Top 10 pro Provincia CEUTA LAS PAL/MAS STA. CRUZ D	feb 2022 n vincias con gas Gasolina 98 E5 1,40 1,43 1,44	olina 98 E5 Gasolina 98 E10 1,44	abr 2022 más barata Gasolina 95 E5 Prem. 1,41 1,31 1,29	Gasolina 95 E10	Gasolina 95 E5 1,37 1,31 1,32
ene 2022 1 op 10 provincias co Provincia WELILLA LAS PALLMAS STA. CRUZ DE TENERIFE CEUTA TERUEL	feb 2022 n gasóleo A Gasoleo A 1,24 1,25 1,28 1,34 1,58	mar 20 A más barar Gasoleo B 1,12	22 abr 20 to Gasoleo premium 1,28 1,36 1,39 1,38 1,68	22 Top 10 pro Provincia CEUTA LAS PALMAS STA, CRUZ D TENERIFE	feb 2022 n vincias con gas Gasolina 98 E5 1,40 1,43 1,43 E 1,44	ar 2022 Dina 98 E5 Gasolina 98 E10 1,44	abr 2022 más barata Gasolina 95 E5 Prem. 1,41 1,31 1,29 1 76	Gasolina 95 E10	Gasolina 95 E5 1,37 1,31 1,32 1,66
ene 2022 1 op 10 provincias co Provincia WELILLA LAS PALMAS STA. CRUZ DE TENERIFE CEUTA TERUEL LLIDIA	feb 2022 m gasóleo A Gasoleo A 1,24 1,25 1,28 1,34 1,58 1,59	mar 20 A más barat Gasoleo B 1,12 1,12	22 abr 20 Co Gasoleo premium 1,28 1,36 1,39 1,38 1,38 1,69	1,4 ene 2022 Top 10 pro Provincia CEUTA LAS PALMAS STA. CRU2 D TENERIFE HUELVA	feb 2022 n vincias con gas Gasolina 98 E5 1,40 1,43 1,43 E 1,44 1,79 1,79	nar 2022 Dina 98 E5 Gasolina 98 E10 1,44 1,69	abr 2022 más barata Gasolina 95 E5 Prem. 1,41 1,31 1,29 1,76	Gasolina 95 E10 1,83	Gasolina 95 E5 1,37 1,31 1,32 1,66
ene 2022 1 op 10 provincias co Provincia WELILLA LAS PALMAS STA. CRUZ DE TENERIFE CEUTA TERUEL LEIDA LLEIDA ALMERIA	feb 2022 Gasoleo A Gasoleo A 1,24 1,25 1,28 1,34 1,58 1,59 1,59	mar 20 A más baran Gasoleo B 1,12 1,17 1,19	22 abr 20 Co Gasoleo premium 1,28 1,36 1,39 1,38 1,68 1,69 1,69	1,4 ene 2022 Top 10 pro Provincia CEUTA LAS PALMAS STA. CRUZ HUEUA CIUDAD REAL CORDOR REAL	reb 2022 n n vincias con gas Gasolina 98 E5 1,40 1,43 PE 1,44 1,79 L 1,79 1 70	nar 2022 Dina 98 E5 Gasolina 98 E10 1,44 1,69	abr 2022 más barata Gasolina 95 E5 Prem. 1,41 1,31 1,29 1,76 1,73 1,73 1,73	Gasolina 95 E10 1,83	Gasolina 95 E5 1,37 1,31 1,32 1,66 1,66 1,66
ene 2022 f op 10 provincias co Provincia WELILLA LAS PALMAS STA. CRUZ DE TENERIFE CEUTA TERUEL LLEIDA LLEIDA LLEIDA ALMERIA	feb 2022 on gasóleo A Gasoleo A 1,24 1,25 1,28 1,34 1,59 1,59 1,60	mar 20 A más barat Gasoleo B 1,12 1,17 1,19 1,19	22 abr 20 to Gasoleo premium 1,28 1,36 1,39 1,38 1,68 1,69 1,69 1,69 1,69	1,4 ene 2022 Top 10 pro Provincia CEUTA LAS PAL/MAS STA. CRUZ HUELVA CIUDAD REAL CIUDAD REAL CIUDAD REAL CIUDAD REAL CIUDAD REAL	reb 2022 n vincias con gase 5 6 6 7 8 7 8 7 8 7 8 7 8 8 7 8 8 7 8 8 7 8 8 8 7 8 8 8 8 8 7 8	aar 2022 blina 98 E5 Gasolina 98 E10 1,44 1,69	abr 2022 más barata Gasolina 95 E5 Prem. 1,41 1,31 1,29 1,76 1,73 1,73 1,73 1,73 1,73	Gasolina 95 E10 1,83 1,69	Gasolina 95 E5 1,37 1,31 1,32 1,66 1,66 1,66 1,66
ene 2022 i pp 10 provincias co Provincia WELILLA LAS PALMAS STA. CRUZ DE TENERIFE LEUTA LEUDA ALMERIA ALLENCIA WURCIA	reb 2022 regasóleo J Casoleo A 1,24 1,25 1,28 1,34 1,58 1,59 1,59 1,60 1,60	mar 20 A más barai Gasoleo B 1,12 1,17 1,19 1,19 1,19	22 abr 20 Co Gasoleo premium 1,28 1,36 1,39 1,38 1,68 1,69 1,69 1,69 1,69 1,68	22 Top 10 pro Provincia CEUTA LAS PALMAS STA. CRUZ D TENERIFE HUELVA CIUDAD REAL CORDOBA MURCIA JAFN	reb 2022 n vinctas con gas gasolina 98 E5 1,40 1,43 E 1,44 1,79 1,79 1,79 1,80 4 so	nar 2022 blina 98 E5 Gasolina 98 E10 1,44 1,69	abr 2022 más barata Gasolina 95 E5 Prem. 1,41 1,31 1,29 1,76 1,73 1,73 1,73 1,73	Gasolina 95 E10 1,83 1,69 1,87	Gasolina 95 E5 1,37 1,31 1,32 1,66 1,66 1,66 1,66

Figura 37. Análisis del precio del combustible en Power BI



Figura 38. Análisis del clima en Power BI

Por otra parte, una vez hemos creado el informe lo podemos publicar en un área de trabajo o workspace.

info@stratebi.com





Figura 39. Workspace de Power BI

Gracias a la API del servicio de Power BI, podemos actualizar el dataset con nuevos datos. Esto permite actualizar el dataset desde el pipeline que teníamos en Synapse con una Web Activity.



Figura 40. Pipeline actualización dataset

A continuación, vamos a analizar en detalle el pipeline. En primer lugar, se encuentra la web activity que, como hemos mencionado, se encarga de realizar una petición a la API de Power BI. En la siguiente imagen, se puede observar la configuración de la actividad:



info@stratebi.com

General Settings U	ser properties
URL	https://api.powerbi.com/v1.0/myorg/g
Method *	POST ~
Headers	+ New
Body	Refresh dataset
Datasets	+ Add dataset reference
Linked services	+ Add linked service reference
Integration runtime * 🛈	AutoResolveIntegrationRuntime ~
Disable certificate validation	on 🗌
HTTP request timeout ${\scriptstyle \bigcirc}$	00:01:00
Authentication	System Assigned Managed Identity \sim
Resource *	https://analysis.windows.net/powerbi/api
> Advanced	

Figura 41. Configuración de la actividad

La URL que se emplea es la siguiente:

https://api.powerbi.com/v1.0/myorg/groups/{group_id}/datasets/{dataset_id}/refreshes
Tanto el group_id como el dataset_id se pueden obtener desde el servicio de Power BI.

S https://app.powerbi.com/groups/<group_id>/datasets/<dataset_id>/details

Figura 42. Dataset id y group id en Power BI Service

Hay que tener en cuenta que para que Synapse pueda actualizar el dataset de Power BI, el service principal de Synapse debe tener asignados los permisos necesarios al workspace de Power BI. Para ello, debemos incluir el service principal de Synapse en uno de los grupos de seguridad de Active Directory y otorgar permiso de *Contributor* al workspace donde se encuentra el dataset.

Una vez configurada la actualización desde Synapse, nos queda un último paso que realizar. Si recordamos nuestro pipeline, después del ELTL, se apaga el pool SQL dedicado. La actividad de refresco del dataset, se debe añadir al pipeline principal antes de apagar el pool SQL. Por otra parte, como se observa en la figura 40, es conveniente añadir un bucle que espere hasta que el dataset esté actualizado para que no haya problemas con la actualización si apagamos el pool SQL antes. Este bucle tiene la siguiente forma:

www.stratebi.com

91.788.34.10



info@stratebi.com



Figura 43. Bucle de espera actualización del dataset

En este bucle, se espera un tiempo determinado, establecido de forma arbitraria y tras ese periodo, se obtiene el estado de la actualización. Esta acción se repite hasta que se haya completado la actualización. Una vez configurada la actualización del dataset podemos incluir la actividad en el pipeline principal.

		Execute Pipeline	
Web		Refresh Power BI	
		dataset	
SQL poor resume	Execute Pipeline		Web
			 SQL pool pause

Figura 44. Pipeline de Synapse

info@stratebi.com



6. **AUTOMATIZACIÓN CON TRIGGERS**

Por último, vamos a automatizar el pipeline para que se ejecute cada mes obteniendo los datos del mes anterior. Para lograrlo, se van a emplear los triggers. Existen diferentes tipos de triggers que se activan por eventos personalizados, eventos relacionados con el almacenamiento o eventos temporales. En este caso, el trigger va a ser temporal y se va a activar el día 1 de cada mes a las 12:00. En la siguiente imagen, se puede observar cómo es la configuración de este trigger:

Edit	trigger							
Name	*							
Mon	thlyTrigge	r						
Descr	iption							
								1.
Type	e i i i i i i i i i i i i i i i i i i i							
Sche	duleTrigge	2r						
Start (date * 🕕							
6/1/	22 00:00:0	0						
Time	zone * 🕕							
Coo	rdinated U	niversal Time	(UTC)					~
Recur	rence * 🛈)						
Every	1				Month(s)			~
\sim Ac	lvanced re	currence opt	ions					
_								
0) Month da	ays 🔿 Wee	ek days					
Se	lect day(s) of the mor	<u>ith to execut</u>	<u>e</u>	-		-	_
	1	2	3	4	5	6	7	
	8	9	10	11	12	13	14	
	15	16	17	18	19	20	21	_
	22	23	24	25	26	27	28	
	29	30	31	Last				_
Ex	ecute at th	nese times	1		1			()
Но	ours	12	×					
			~					
M	nutes	0	^					

Figura 45. Trigger mensual

Al realizar la configuración del trigger, también se deben establecer los parámetros a emplear.

info@stratebi.com



Name	Туре	Value
subscriptionId	string	Value
poolName	string	Value
resourceGroup	string	Value
start_time	string	@trigger.startTime

Figura 46. Parámetros del trigger

En este caso, solo se asigna uno de los parámetros y el resto se dejan por defecto. El parámetro start_time sirve para calcular el mes anterior del que se quiere extraer la información. Para obtener los datos del mes y año que se quieren extraer, se utilizan variables. Para asignar estas variables tenemos que modificar el pipeline una última vez añadiendo dos actividades de tipo "Set variable".

Set variable				
(X)	Set month		ľ	

Figura 47. Set variable month

Mediante las expresiones de Synapse, podemos calcular el mes anterior:

```
@formatDateTime(subtractFromTime(formatDateTime(pipeline().parameters.sta
rt_time, 'yyyy/MM/ddThh:mm'), 1, 'MONTH'), 'MM')
```

De manera análoga, calculamos el año correspondiente al mes anterior:

Set variable						
$\left(oldsymbol{\mathcal{X}} ight)$ Set year						

Figura 48. Set variable year

```
@formatDateTime(subtractFromTime(formatDateTime(pipeline().parameters.sta
rt_time, 'yyyy/MM/ddThh:mm'), 1, 'MONTH'), 'yyyy')
```

Finalmente, el pipeline queda de la siguiente forma:

www.stratebi.com

91.788.34.10

strate

info@stratebi.com





7. **CONCLUSIONES**

Hemos visto a lo largo de este documento cómo utilizar Azure Synapse Analytics, con el que hemos hecho referencia a capacidades de almacenamiento, exploración y procesamiento de datos. En lo que respecta al procesamiento, hemos visto que con los pipelines podemos obtener datos diferentes orígenes y cómo tratarlos para que se adapten al mismo formato con notebooks y dataflows. Además, hemos visto cómo gestionar los ADLS Gen 2 en diferentes capas para mantener el dato en los diferentes estados desde su extracción pasando por diferentes transformaciones.

Por otra parte, hemos demostrado cómo integrar Power BI con Synapse Analytics mediante los pools SQL dedicados. Además, para un correcto uso de los servicios de Azure, hemos podido aprovechar las APIs de Azure y de Power BI para automatizar muchos de los procesos como el refresco de datos o la gestión del estado del pool dedicado.

En definitiva, se ha llevado a cabo un paseo por Azure Synapse Analytics, como pieza clave en las arquitecturas modernas de datos cuando hablamos de ecosistema Azure de Microsoft.

