

Azure Serverless Analytics

Problem Statement:

Businesses and organizations must process and analyze vast amounts of data in a way that is both economical and scalable. Traditional data processing techniques frequently call for significant upfront investment and infrastructure management, which restricts flexibility and hinders agility. This delays the time to insights and results in severe resource waste.

Data Analytics?

Every organization wants to know about the profits they make during the financial year also they want to know about the customer needs and the latest trends so that the company can change itself to make better profits and cater to the customer changing needs.

To get the above valuable information raw data is needed which is collected from **Feedback Forms, Surveys, Social Media Monitoring, etc.** This raw data is then processed and valuable insight/information is taken as an Output which shows/depicts the value that the organization needs for its use.

Processing the raw data to get the valuable information is known as **Data Analytics** in other words the process of extracting the valuable information from the Raw data is Known as **Data Analytics**.

Tools Required For Data Analytics :

1. Data Exploration.
2. Artificial Intelligence.
3. Machine Learning.
4. Data Warehousing.
5. Real-Time BI.
6. Data Integration.

Solution :

Azure Serverless Analytics provides a scalable and cost-effective method for processing and analyzing huge amounts of data. A variety of services and functionalities are available through Azure Serverless Analytics to assist businesses in overcoming the difficulties posed by conventional data processing techniques. The following are some essential elements of the remedy:

- 1) Serverless Computing
- 2) Scalability and Elasticity
- 3) Cost Optimization
- 4) Integration with Azure Services
- 5) Real-time Data Processing
- 6) Ad-hoc Data Exploration

Organizations may transcend the constraints of conventional data processing techniques and realize the full potential of their data by utilizing Azure Serverless Analytics. With agility, cost-effectiveness, and scalability, it enables businesses to process and analyze data at scale, gain insightful knowledge, and make data-driven decisions.

Azure Synapse Analytics:

All the above-mentioned services which are needed for analytics are provided by Azure together in a single platform benefitting users to switch to different platforms and saving the cost.

Suppose a user wants to have a Data Warehouse for which he has to purchase physical storage which will cost him around 500\$ in Azure user can pay according to its needs a user can use the memory and pay 15\$ for a month.

This kind of service is provided by Microsoft which is named **Azure Synapse Analytics**.

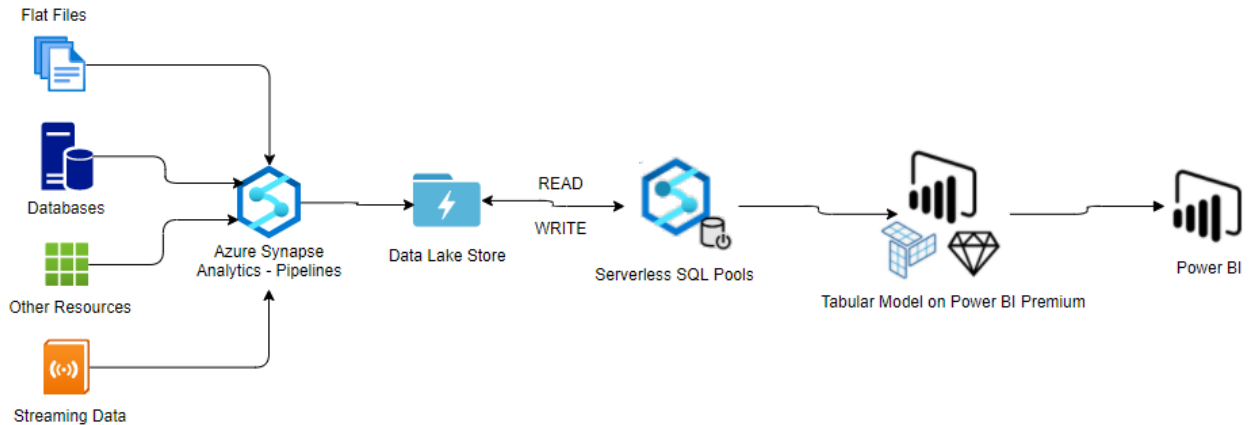


Figure 1: Azure Synapse Analytics

Architecture

Azure Synapse Analytics consists of **Data Lake Store**, **Serverless SQL Pools**, and **Power BI**.

In **Data Lake Store**, all types of files, whether a doc, pdf, .csv, etc., can be stored and queried by **Serverless SQL Pools** to obtain data for further processing and valuable output.

Serverless SQL Pool is the query service that is used to access the data present in the data lake using **T-SQL** query. It consisted of two pools:

1. **Dedicated SQL Pool:** Workplace can have more than one dedicated SQL Pool, which is used for dedicated models.
2. **Serverless SQL Pool:** It is used for serverless models and every workplace has at least one serverless SQL Pool.

Power BI is used for processing the data and getting the valuable output which can be used by the organization to make more profits and cater to the changing/dynamic needs of customers.

Apache Spark :

Spark is used for processing large data sets or big data sets. Here Azure Synapse Analytics has inbuilt Apache Spark for processing large datasets within minutes.

Apache Spark consists of the following parts:

1. Spark for Synapse.
2. Spark Application Spark pool.
3. Job definition for Spark.
4. Notebook.

Synapse Pipeline :

It consists of the following characteristics :

1. Integration Of Data.
2. Data Stream.
3. Pipeline.
4. Activity.
5. Trigger.
6. Combined Dataset.

Synapse Studio :

It consists of secured architecture and it has trustworthy collaboration boundaries for performing cloud-based analytics in Azure.

It can be easily deployed in specific regions, however, it has collaborated with the ADLS Gen2 account for temporary file storage.

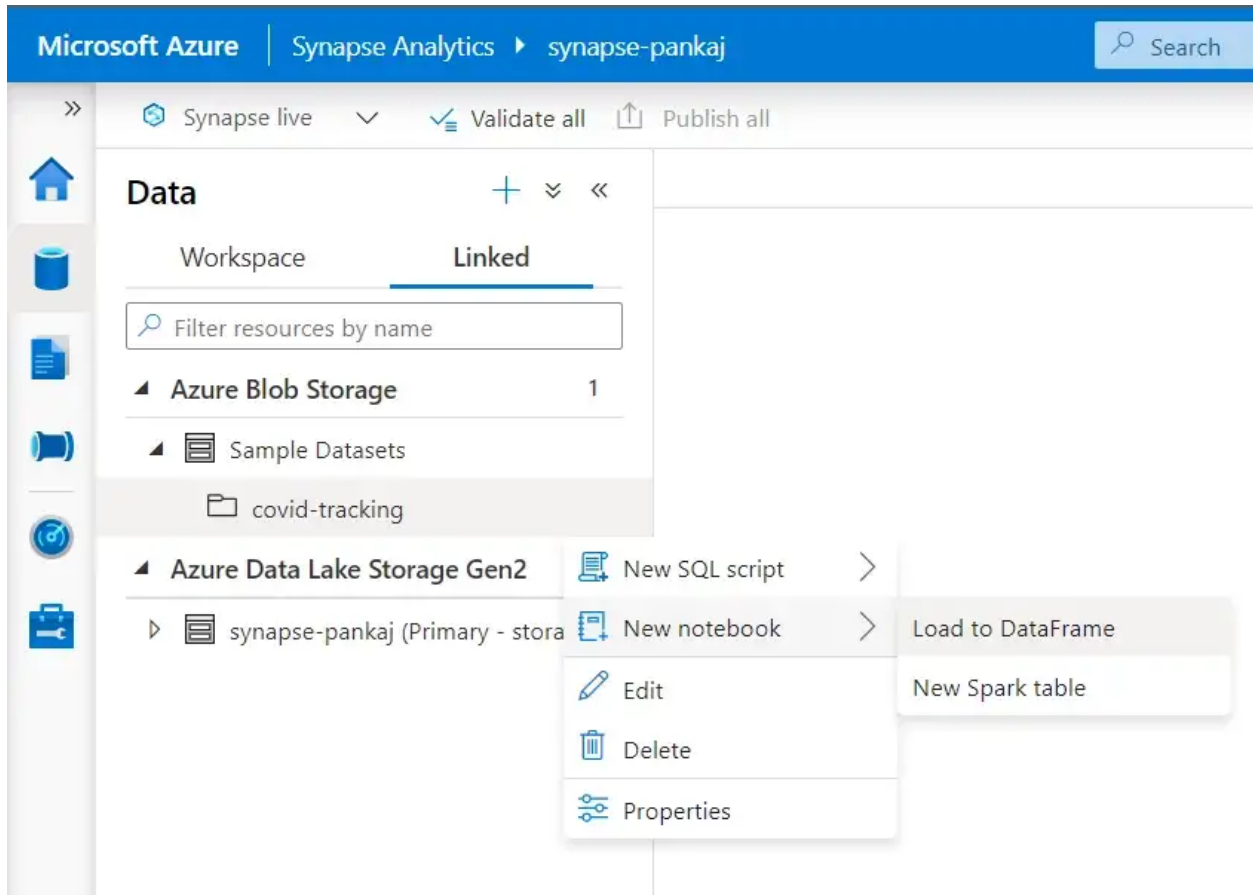
Features Of Azure Synapse Analytics

We assume you've already installed and set up the Azure Synapse Analytics workspace in your Azure subscription.

The screenshot shows the Synapse workspace interface for a user named 'dayosbigdata'. On the left is a navigation sidebar with icons for Home, Data, Develop, Orchestrate, Monitor, and Manage. The main content area features a header with the workspace name and a 'New' button. Below this are four primary action cards: 'Ingest' (copy data tool), 'Explore' (navigate and interact with data), 'Analyze' (use SQL or Spark for insights), and 'Visualize' (build interactive reports with Power BI). At the bottom, there are two sections: 'Resources' (with 'Recent' and 'Pinned' tabs, currently showing 'No recent resources') and 'Useful links' (containing links for Getting started, Synapse Analytics overview, Pricing, Documentation, and Give feedback).

Load and Analyze Data using Spark

1. Using Data Hub View the covid-tracking dataset and the pictures.
2. Open the new Spark notebook with the following dataset.



The notebook includes the default code listed below, along with additional features for using the Spark backbone framework to analyze the dataset:

```

1 %%pyspark
2 blob_account_name = "pandemicdatalake"
3 blob_container_name = "public"
4 blob_relative_path = "curated/covid-19/covid_tracking/latest/covid_tracking.parquet"
5 blob_sas_token = ""
6 # Allow SPARK to read from Blob remotely
7 wasbs_path = 'wasbs://%s@%s.blob.core.windows.net/%s' % (blob_container_name, blob_account_name, blob_relative_path)
8
9 spark.conf.set(
10     'fs.azure.sas.%s.%s.blob.core.windows.net' % (blob_container_name, blob_account_name),
11     blob_sas_token)
12 df = spark.read.parquet(wasbs_path)
13 display(df.limit(10))

```

Command executed in 2mins 24s 38ms by pjanani_jumia on 02-21-2021 16:12:10.948 +05:30

Job execution Succeeded Spark 2 executors 8 cores [View in monitoring](#) [Open Spark UI](#)

Progress indicator is out of sync. Reason: Error: Network Error

View: Table Chart

date	state	positive	negative	pending	hospitalized_cur...	hospitalized_cu...	in_icu_currently	in_icu_cumulative	on_ventilator_cu...	on_ventilato
2021-02-19	AK	55198	null	null	34	1243	null	null	4	null
2021-02-19	AL	485212	1867861	null	951	44767	null	2619	null	1488

Analyze the Data in Serverless SQL Pool

Again import the sample dataset from the Data Hub but make a New SQL Script as shown below:

Synapse live Validate all Publish all

Data + ≡ <<

Workspace Linked

Filter resources by name

- Azure Blob Storage 1
- Sample Datasets
- covid-tracking
- Azure Data Lake Storage Gen2
 - New SQL script > Select TOP 100 rows
 - New notebook > Create external table
- synapse-pankaj (Primary - stora)
 - fsyssynapsepankaj (Primary)
 - Edit
 - Delete
 - Properties

You will then have access to a free-form script to use native T-SQL queries to investigate the data set.

SQL script 1 • fsyssynapsepankaj • SQL script 2 • Notebook 1

Run Undo Publish Query plan Connect to Built-in Use database master

```

1 -- This is auto-generated code
2 SELECT
3     TOP 100 *
4 FROM
5     OPENROWSET(
6         BULK 'https://pandemicdatalake.blob.core.windows.net/public/curated/covid-19/covid_tracking/latest/covid_tracking.parquet',
7         FORMAT = 'parquet'
8     ) AS [result];

```

Results Messages

View Table Chart Export results

Date	State	Positive	Negative	Pending	Hospitalized_c...	Hospitalized_c...	In_icu_currently	In_icu_cumulat...
2021-02-19T00:...	AK	55198	NULL	NULL	34	1243	NULL	NULL
2021-02-19T00:...	AL	485212	1867861	NULL	951	44767	NULL	2619
2021-02-19T00:...	AR	314713	2348207	NULL	630	14500	237	NULL

Setup and Integrate a Pipeline

1. Select the Pipeline from Synapse Analytics Studio's Integrate hub.

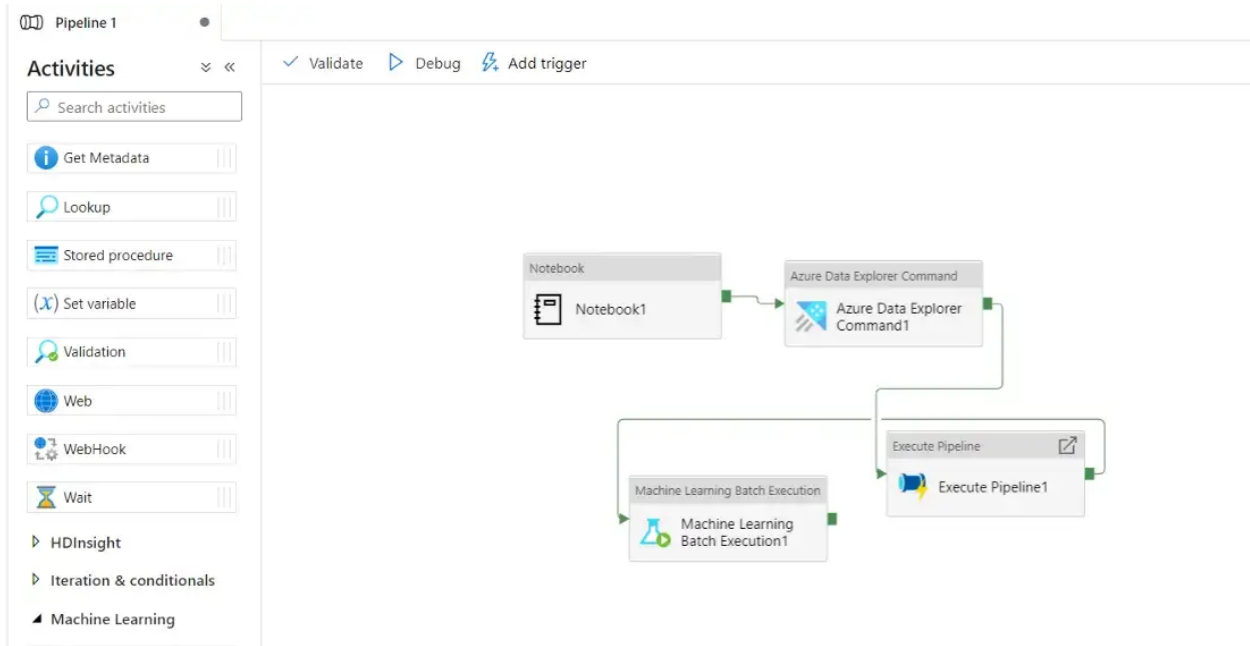
Synapse live Validate all Publish all

Integrate

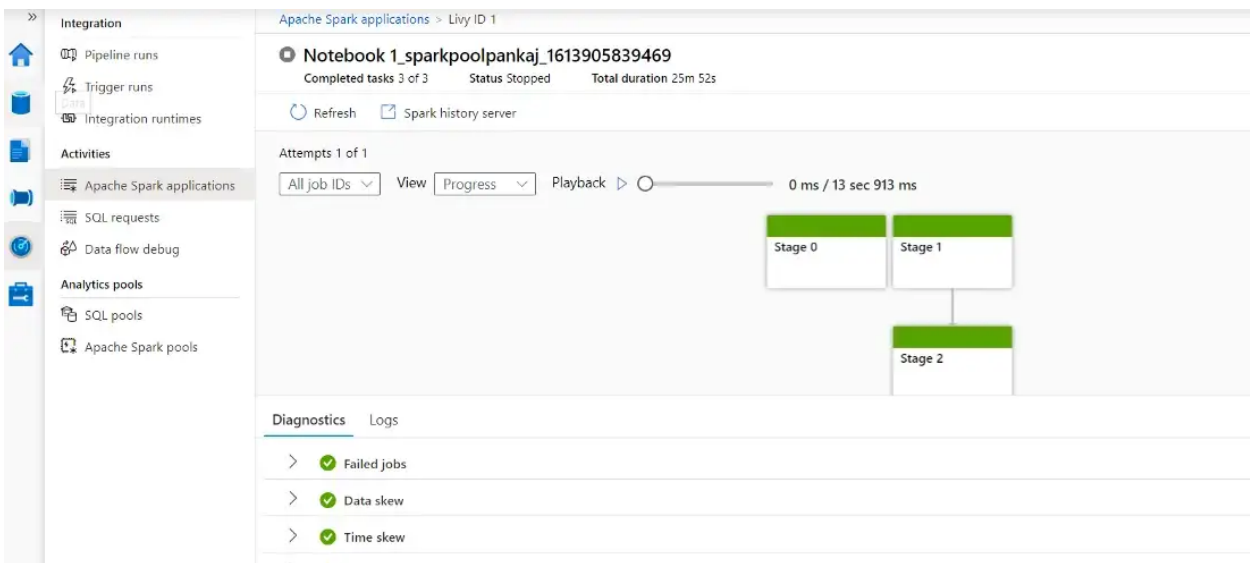
Filter resources by name

- +
- ≡
- «
- Pipeline
- Copy Data tool
- Browse gallery

2. After the pipeline has been created, you can add process activities based on the current issue.



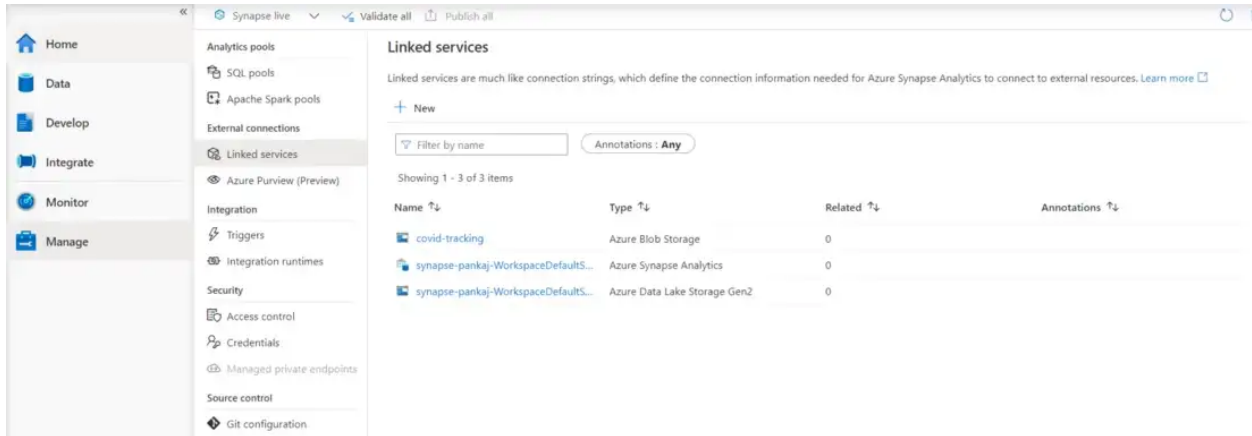
3. To respond to an event or manual execution of the Pipeline workflow, you can add trigger conditions. To track the status of any pipeline execution, click the Monitor hub and then choose the Pipeline runs option.



Integrate Linked Services

To integrate and enable multiple Linked Services, such as Power BI, use Synapse Analytics Studio.

1. Utilize Manage Hub to view the currently linked services.



2. Create a link to the PowerBI workspace as well to make use of data reporting and visualization.

Azure Synapse Service for industries

1. Financial Services :

Utilize industry-leading features to guarantee data security. As it continues to maintain a competitive edge by applying a cutting-edge approach to handling big data, data warehousing, building tailored customer experiences, and putting in place robust compliance and governance mechanisms to protect consumer data.

2. Manufacturing Service :

Gaining scalable real-time insights by using Azure Synapse Analytics. Real-time access to both fresh and historical data is made possible by Industry 4.0's integration of operational and analytical technology.

3. Retail Service :

By combining data from several channels and gaining real-time insights, you can use an end-to-end analytics service to better understand your customers and create a dependable supply chain.

4. Healthcare Service :

The healthcare industry is under pressure from a scarcity of carers, legislative limitations, and changing patient expectations. Deliver tailored care, protect patient information, and give care teams more freedom.

Challenges of Azure Serverless Analytics :

- 1) **Data Complexity and Variety:** Complex data types include structured, semi-structured, and unstructured data that can be difficult to process and analyze. Careful planning and data preparation may be necessary to handle data quality issues and ensure interoperability with various data formats.
- 2) **Performance Optimization:** It can be difficult to optimize performance in serverless systems. To process data quickly and minimize latency, it is essential to properly design and configure resources, manage data partitions, and optimize query speed.
- 3) **Scalability and Resource Allocation:** Although serverless computing scales resources automatically based on workload needs, resource allocation optimization and assuring efficient scaling can be challenging. To reconcile performance requirements with cost efficiency, much thought must be given.
- 4) **Security and Compliance:** It is essential to safeguard sensitive data and make sure that data privacy laws are followed. To protect data and satisfy compliance standards, it becomes crucial to implement the proper security controls, access restrictions, encryption, and monitoring tools.
- 5) **Integration with Existing Infrastructure:** Integrating Azure Serverless Analytics with existing infrastructure, applications, and data systems can pose challenges. Compatibility issues, data integration complexities, and ensuring seamless data flow across different components require careful planning and integration strategies.
- 6) **Monitoring and Troubleshooting:** The distributed and event-driven architecture of serverless analytics systems makes monitoring and troubleshooting them difficult. Effective error handling, debugging, and system performance visibility can necessitate the use of powerful monitoring tools and techniques.
- 7) **Cost Management:** Although serverless computing reduces costs by scaling resources according to demand, maintaining and controlling expenses can be challenging. Cost-effective implementation depends on ensuring optimal resource allocation, optimizing function execution, and tracking usage.

8) **Skill Set and Expertise:** Azure Serverless Analytics adoption and implementation call for technical knowledge in cloud computing, data processing, and analytics. To ensure successful implementation and continued management, organizations may need to spend money on training or hiring qualified personnel.

Business Benefits of Azure Serverless Analytics :

1) **Cost Savings:** A pay-as-you-go business model is used by Azure Serverless Analytics, where you only pay for the resources used to run your analytics workloads. As a result, there is no longer a need for initial infrastructure investment, and the cost of idle resources is decreased. The automated scaling function makes sure that resources are used to their full potential, further optimizing expenses.

2) **Scalability and Flexibility:** Your analytics workloads can accommodate changing data quantities and processing needs thanks to the automatic scaling features offered by Azure Serverless Analytics. You can handle and analyze massive volumes of data effectively without worrying about infrastructure limitations because it can scale up or down fluidly based on workload requirements.

3) **Faster Time-to-Insights:** Azure Serverless Analytics uses serverless computing to do away with the requirement for infrastructure deployment and management so you can concentrate on data analysis and insights. As a result, it is possible to set up and run analytics workloads rapidly, cycle through data processing pipelines, and gain insightful knowledge from your data more quickly.

4) **Real-Time and Near Real-Time Analytics:** Data streaming and processing in real-time are supported by Azure Serverless Analytics. This makes it possible for you to gather, process, and analyze streaming data from a variety of sources, including IoT devices, social media feeds, and logs. You can react fast to shifting data trends, make timely decisions, and gain fresh insights thanks to real-time analytics capabilities.

5) **Agility and Iterative Analysis:** Ad-hoc data exploration and analysis may be done in a flexible context with Azure Serverless Analytics. You can query and alter your data as needed without having to create a schema in advance or keep up infrastructure. Through iterative analysis, which is made possible by this agility, you can improve your analytics workflows, try out new ideas, and gain deeper insights from your data.

6) **Integration with Azure Ecosystem:** Azure Data Lake Storage, Azure Synapse Analytics, Azure Stream Analytics, and other Azure services are all effortlessly

integrated with Azure Serverless Analytics. By utilizing the capabilities of several Azure services, this integration enables the creation of complete and integrated data analytics ecosystems.

7) **Reduced Operational Complexity:** Azure Serverless Analytics does away with the necessity for infrastructure management and upkeep by utilizing serverless computing. This lessens operational complexity, frees up resources, and enables your team to concentrate on data analysis and insight generation rather than infrastructure management.

8) **Improved Collaboration and Productivity:** For data analysts, data scientists, and other stakeholders to collaborate on analytics projects, Azure Serverless Analytics offers a platform. It allows version control, permits sharing and collaboration on data processing pipelines, and offers tools for monitoring and debugging. Because of increased productivity and teamwork, data analysis and decision-making procedures become more effective.

Organizations may save money, and increase scalability, agility, and speed of insight by utilizing Azure Serverless Analytics, allowing them to make data-driven decisions and achieve a competitive edge in today's data-driven business environment.

Conclusion :

In this succinct overview of Azure Synapse Analytics, I've waded across the shallow seas to comprehend the managed service's fundamental functional capabilities.

For data engineers, Azure Synapse is the one-stop shop where they can find a comprehensive end-to-end data pipeline. Additionally, Synapse can manage all company information in a shorter amount of time. As a result, you do not need to invest money in additional technology platforms to combine data from many platforms in one location. If you're persuaded and want to experience Synapse's commercial advantages, we're here to help.

We now know what Synapse Analytics is, why it's crucial to achieving organizational goals, and just a little bit about how its capabilities can be used for different analytical use cases.