

Databricks Masterclass Notes

Introduction

Purpose: Learn Databricks from scratch, including foundational concepts, Spark clusters, Delta Lake, and advanced optimization techniques.

Why Databricks?

Independent of cloud environments (works with AWS, Azure, GCP).

Integral for big data processing and analytics.

Constantly evolving with cutting-edge updates (e.g., Unity Catalog, Delta Lake advancements).

Key Components Covered

Spark Cluster and Architecture:

Fundamental to big data processing.

Includes nodes, cluster manager, driver, and worker nodes.

Databricks File System (DBFS):

Provides a distributed file system interface.

Allows seamless data handling across multiple sources.

Magic Commands:

Simplify data exploration and command execution.

Databricks Utilities (DBUtils):

Aid in file system management, secrets management, and job control.

Delta Lake:

Optimized storage layer for Apache Spark.

Features like time travel, data versioning, and Z-Ordering.

Autoloader:

Enables incremental data loading from external sources.

Detailed Explanations

1. Spark Cluster

Definition: A collection of machines (nodes) acting as a single system for parallel data processing.

Components:

Driver Program: Manages the execution of tasks and stores metadata.

Worker Nodes: Execute tasks on data partitions.

Cluster Manager: Allocates resources and schedules tasks.

Cluster vs. Standalone Machine:

Aspect Standalone Machine Cluster

Processing Power Limited by single machine specs Scales horizontally with nodes.

Fault Tolerance Low High

Scalability Limited Virtually unlimited

2. Spark Architecture

Process:

Cluster Manager assigns resources.

Driver Program breaks data tasks into smaller subtasks.

Tasks are distributed to Worker Nodes.

Parallel Processing delivers faster results.

Diagram:

css

Copy code

[User] --> [Cluster Manager]

/

[Driver Program] --> [Worker Nodes]

/

(Assign tasks to partitions)

3. Databricks File System (DBFS)

A layer over cloud storage.

Syntax Examples:

List files: %fs ls /mnt/data

Copy files: %fs cp /source /destination

4. Delta Lake

Features:

Time Travel: Access historical data states.

Versioning: Maintains data updates for auditability.

Tombstoning: Marks deleted files but retains metadata for recovery.

Optimization Commands:

Z-Ordering: OPTIMIZE delta.table_name ZORDER BY (column_name)

Vacuum: VACUUM delta.table_name RETAIN 168 HOURS

5. Databricks Utilities

File Operations:

dbutils.fs.ls("path")

dbutils.fs.rm("path", True)

Secrets Management:

Store and retrieve sensitive credentials securely.

6. Autoloader

Use Case: Continuously ingest new data from cloud sources.

Syntax:

python

Copy code

```
df = spark.readStream.format("cloudFiles") \
.option("cloudFiles.format", "csv") \
.load("s3://path")
```

Important Interview Points

1. Real-Life Scenarios

Scenario: Handling incremental data loads.

Solution: Use Delta Lake Autoloader with streaming.

Scenario: Optimizing large table queries.

Solution: Apply Z-Ordering on frequently queried columns.

2. Common Interview Questions

Explain Spark Architecture.

What is Delta Lake? How does it differ from traditional storage?

What is the difference between Managed and External Delta Tables?

Describe the purpose of the Driver Program.

How does Databricks handle data partitioning?

3. Scenario-Based Questions

Scenario: A dataset stored in AWS S3 needs incremental processing. How would you achieve this in Databricks?

Approach: Use Databricks Autoloader with Delta Lake.

Scenario: Your pipeline requires consistent metadata updates. How do you ensure data integrity during schema evolution?

Approach: Utilize Delta Lake's schema enforcement.

Final Notes

Databricks is versatile, integrating seamlessly with AWS, Azure, and GCP.

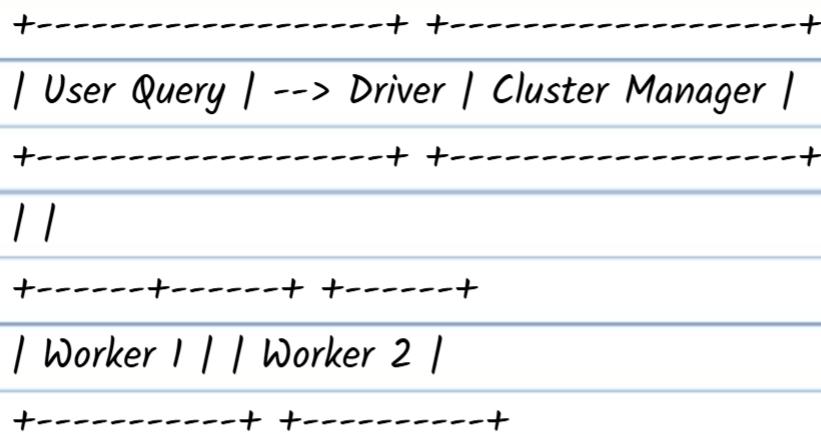
Staying updated with the latest trends (e.g., Unity Catalog, Delta Lake enhancements) is critical for career growth.

Hands-on experience is a must; leverage free Databricks cloud accounts for practice.

Diagram of Spark Architecture

plaintext

Copy code



Part02

Notes and Explanations

1. Spark Architecture in Databricks

Key Components:

Cluster Manager: Responsible for creating the driver program and worker nodes. Its task ends after setting up the cluster.

Driver Program: Orchestrates the tasks, provides instructions, and communicates with worker nodes. It doesn't perform the actual data processing.

Worker Nodes: Perform the heavy lifting by executing tasks assigned by the driver program.

Process Flow:

Cluster Manager initializes the driver program.

Driver program splits the job into smaller tasks.

Cluster Manager creates the required number of worker nodes.

Worker nodes execute the tasks and return results to the driver program.

Communication happens between the driver program and worker nodes, not with the cluster manager after initialization.

Component Responsibility

Cluster Manager Sets up the driver and workers, then exits.

Driver Program Orchestrates the flow, provides task instructions, but does not process data.

Worker Nodes Executes tasks based on the driver's instructions and processes the data.

2. Databricks Overview

Databricks:

A unified analytics platform that manages Spark clusters for efficient big data processing.

Automates cluster management, reducing the need for manual intervention.

Acts as a management layer over Spark, streamlining data engineering tasks.

Benefits:

Ease of Cluster Management: Handles cluster creation, configuration, and scaling.

Integration with Cloud Providers: Works seamlessly with AWS, Azure, and GCP.

Interactive Workspace: Provides notebooks for collaborative development and analysis.

Azure Account Setup for Databricks

Types of Accounts:

Community Edition: Free version for learning, lacks cloud integration.

Cloud Account: Required for production and real-world use cases.

Steps to Create Azure Account:

Visit portal.azure.com and sign up using a valid email.

Enter personal details and billing information (Azure provides \$200 credits for the first 30 days).

Creating Databricks and Data Lake:

Resource Group: Logical container for resources like Databricks workspace and data lakes.

Data Lake: Core storage for structured and unstructured data.

Databricks Workspace: Platform to interact with Spark clusters and process data.

Step Action

Sign up for Azure Create an account on portal.azure.com.

Create Resource Group Group related resources under a common folder-like structure.

Setup Data Lake Configure a scalable storage for big data.

Setup Databricks Create a workspace for data engineering and machine learning tasks.

Key Concepts and Definitions

Cluster Manager:

A central entity responsible for setting up the driver program and workers, then exits once the cluster is operational.

Use Case: Orchestrates the initial setup for distributed processing systems.

Driver Program:

Directs the task execution in Spark.

Interview Tip: Highlight its role as the orchestrator, not the executor of tasks.

Worker Nodes:

Execute tasks assigned by the driver program.

Real-life Example: ETL operations where each worker processes a partition of data.

Important Points

Databricks Independent of Cloud: Learning Databricks on one cloud provider applies across others (Azure, AWS, GCP).

Storage Hierarchy:

Data is stored in data lakes, often in formats like Parquet or Delta.

Delta Lake is built atop Data Lake for ACID compliance and performance.

Azure Specific Terms:

Resource Group: Logical folder for organizing resources.

Service Principal: Provides secure, automated access to resources.

Secret Scope: Manages sensitive data (e.g., credentials) securely.

Real-Life Interview Questions

Scenario-Based:

"How would you optimize a Spark job in Databricks for processing a 1 TB dataset?"

Discuss partitioning, caching, and executor configurations.

Conceptual:

"Explain the difference between a driver program and worker nodes in Spark."

Mention orchestration (driver) vs. execution (workers).

Practical:

"How do you connect Databricks with an Azure Data Lake?"

Mention service principals, secret scopes, and the connection flow.

Debugging:

"A Spark job is taking longer than expected. What steps would you take to troubleshoot?"

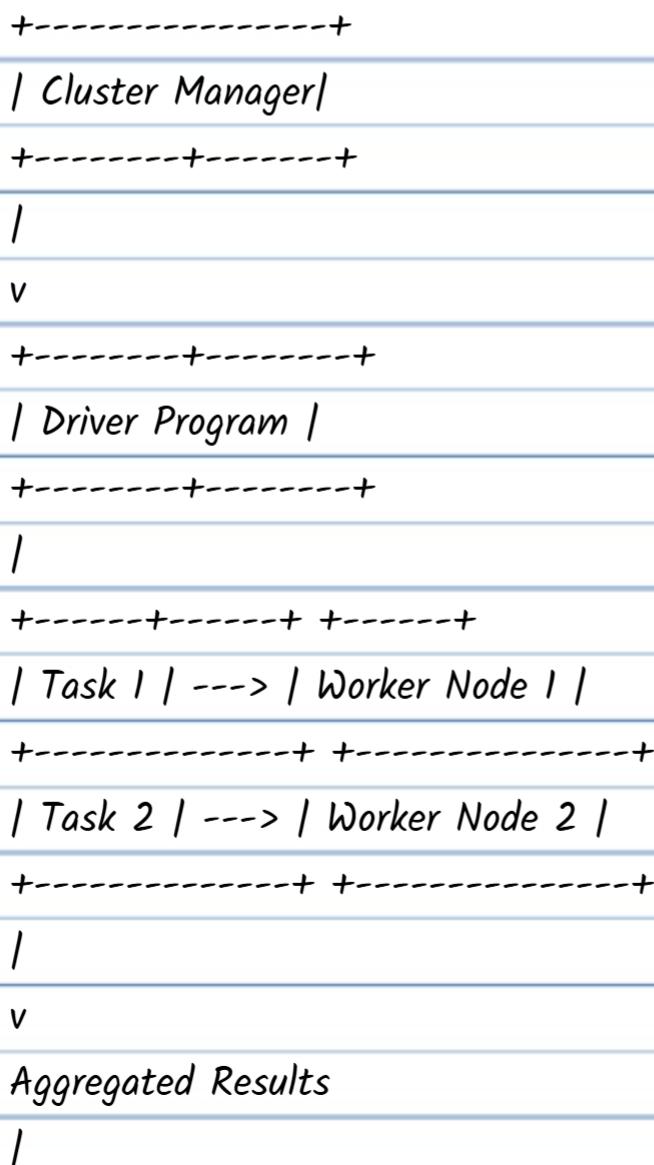
Analyze execution DAG, check skewed partitions, and monitor executor metrics.

Enhanced Visuals

Diagram: Spark Architecture in Databricks

plaintext

Copy code



v

+-----+

| Processed Data |

+-----+

Additional Interview Preparation

Scenario: "How would you configure a Databricks workspace for multi-region support in Azure?"

Discuss resource group setup, VNET peering, and workspace replication.

Optimization: "What are the best practices for partitioning data in Spark?"

Highlight dynamic partitioning, skew handling, and optimal file sizes.

Tool Comparison: "Why choose Databricks over a standalone Spark setup?"

Mention automation, collaboration, and performance enhancements.

=====

Part03

Setting Up Azure Resources for Databricks

Step 1: Create a Resource Group

Navigate to Azure Portal.

Go to "Resource Groups" in the menu or search for it.

Click "Create" and fill in:

Resource Group Name: Example - RG_DatabricksMasterClass.

Region: Choose any Azure region close to your location.

Validation Check: Ensure no errors appear, then click "Review + Create", followed by "Create".

Important Notes:

A resource group is a container that holds related Azure resources for better organization.

Interview Tip: What is a Resource Group?

It's a logical container to manage resources, ensuring operational consistency.

Step 2: Create a Storage Account

In the Resource Group, click "Create", then search for "Storage Account".

Click on "Create" for Microsoft Storage Accounts.

Fill in:

Account Name: Unique, e.g., datalake12345.

Region: Same as the resource group.

Performance Tier: Choose Standard.

Redundancy: Pick LRS (Locally Redundant Storage) for cost-effective replication.

Before clicking "Review + Create", check the Hierarchical Namespace box to create a Data Lake Gen2.

What is Hierarchical Namespace?

Enables folder-like structure for better data organization within a Data Lake.

Supports nested directories (e.g., folders inside folders).

Storage Account Types Comparison Table:

Type	Description	Use Case
Blob Storage	Flat structure; stores unstructured data like images and	

logs. File storage, web apps.

Data Lake Hierarchical; supports nested directories and file formats. Big data analytics, machine learning.

Step 3: Upload Files to Storage Account

Go to Containers in the storage account and create two containers:

Source: Holds input data (e.g., sales.csv).

Destination: Stores transformed/output data.

In the Source Container, upload the sales.csv file.

Use the Upload Button in the Azure interface.

Setting Up Databricks Workspace

Step 1: Create the Workspace

Go to the Resource Group.

Click "Create", then search for "Azure Databricks".

Select "Create", then configure:

Workspace Name: Example - Databricks_Master.

Pricing Tier: Choose Premium for full features (Trial available).

Azure creates a Managed Resource Group for cluster management.

Avoid modifying resources in this group as Databricks manages it automatically.

Exploring the Databricks Interface

Key Sections in Databricks UI

Workspace:

Stores notebooks, dashboards, and scripts.

Best Practice: Create folders to organize resources (e.g., DatabricksMasterClass).

Data Catalog:

A unified repository for metadata (schemas, databases).

Use Case: Monitor data sources and schema consistency.

Workflows:

Automates pipeline orchestration, including dependencies.

Scenario: Automating ETL tasks with notebooks.

Clusters:

Provides scalable computing for processing tasks.

Supports auto-scaling to optimize resource usage.

Diagram: Databricks UI Overview

css

Copy code

```
+-----+  
| Home | Workspace | Data | Workflows |  
+-----+
```

```
+-----+  
| [Recent Files] [Notebooks] [Catalog] |  
+-----+
```

```
| [Cluster Info] [Jobs] [Tables] |  
+-----+
```

Key Features of Databricks

1. Delta Lake

A storage layer enabling ACID transactions on Apache Spark.

Supports real-time streaming and batch workloads.

Interview Question: How does Delta Lake improve data reliability?

Guarantees data consistency, schema enforcement, and time travel for historical analysis.

2. Notebooks

Interactive documents for writing and executing code (supports Python, SQL, Scala).

Use Case: Perform ETL and visualize data within a single interface.

Notebook Example (PySpark):

python

Copy code

```
from pyspark.sql import SparkSession
```

```
# Create Spark Session
```

```
spark = SparkSession.builder.appName("DeltaLakeExample").getOrCreate()
```

```
# Read data from Azure Data Lake
```

```
df = spark.read.csv("path/to/sales.csv", header=True, inferSchema=True)  
df.show()
```

Interview Points and Scenario-Based Questions

Conceptual Questions

What is the difference between Blob Storage and Data Lake Gen2?

Explain the redundancy options (LRS, GRS, ZRS, GZRS).

Scenario-Based Questions

You need to reduce costs for a non-critical project on Azure Databricks.

Which storage and redundancy option would you choose and why?

How would you design a pipeline to process daily sales data stored in Azure Data Lake using Databricks?

Best Practices

Always choose LRS for testing environments to save costs.

Use Databricks Workflows for orchestrating multi-step data pipelines.

Regularly monitor cluster performance and enable auto-scaling.

part04

1. Overview of the Interface

Tabs and Packages:

Compute Tab: Manage clusters for computations.

Workflows: Orchestrate notebook executions in a sequence.

Data Engineering, Machine Learning Tabs:

Tailored tools for specific roles.

Example: Connectors for data engineering or model building for machine learning.

Marketplace: Access third-party connectors or services.

Partner Connect: Seamlessly integrate third-party BI tools like Power BI, Tableau.

Cluster Management:

Simplifies provisioning.

Managed entirely by Databricks.

Diagram: Databricks Interface Flow

mathematica

Copy code

Main Menu

- └── Compute Tab
- └── Workflow Tab
- └── Machine Learning Tools
- └── Data Engineering Connectors
- └── Marketplace
- └── Partner Connect

2. Cluster Configuration

Steps to Create a Cluster:

Navigate to the Compute Tab.

Click on Create Compute and configure:

Policy: Unrestricted or as defined.

Node Type: Single-node for small tasks; multi-node for production.

Runtime Version: Long-term supported versions like 14.3 LTS.

Termination Policy: Configure to save costs (e.g., 30 minutes of inactivity).

Key Definitions:

Cluster: Group of virtual machines for running workloads.

Terminate After: Automatically shuts down inactive clusters to reduce costs.

Interview Question: How do you optimize cluster costs in Databricks?

3. Workspace and Notebooks

Creating a Notebook:

Go to Workspace > Click Create > Select Notebook.

Rename your notebook for better management.

Notebook Features:

File Menu:

Create/import notebooks.

Export options (DBC Archive, Python Notebook, HTML).

Editor Themes: Choose between dark or light themes.

Connect Button: Links the notebook to a cluster.

Magic Commands: Enable multiple languages (e.g., Python, SQL, R, Markdown) in a single notebook.

%python for Python cells.

%sql for SQL execution.

%md for Markdown.

Real-life Scenario: Use Markdown to organize notebooks with headings and comments, improving readability and debugging efficiency.

Interview Question: Explain the advantages of using magic commands in a collaborative environment.

4. Databricks File System (DBFS)

What is DBFS?

A distributed file system that abstracts data lake paths for simplicity.

Instead of URLs, use file semantics like file:/mnt/container_name/file.csv.

Key Use Case:

Simplifies code and reduces redundancy when accessing raw data stored in a data lake.

Diagram: DBFS Abstraction

`javascript`

Copy code

Data Lake

```
| └── Raw Container
|   └── File.csv
└── DBFS Layer
  └── file:/mnt/raw/file.csv
```

5. Creating DataFrames

Steps to Create a Custom DataFrame:

Define data as a list of tuples.

Specify schema (column names and types).

Use `spark.createDataFrame()` to create the DataFrame.

Code Example:

`python`

Copy code

```
data = [(1, "Alice", 90), (2, "Bob", 85), (3, "Charlie", 92)]
```

```
schema = "id INT, name STRING, marks INT"
```

```
df = spark.createDataFrame(data, schema=schema)
```

```
display(df)
```

Output:

```
id name marks
```

```
1 Alice 90
```

```
2 Bob 85
```

```
3 Charlie 92
```

6. Visualization in Databricks

Visualization with `display()`:

Quickly visualize data.

Supports graphs like bar, scatter, and pie charts.

Steps:

Use `display(df)` to load data.

Click on the visualization icon.

Interview Question: Describe a scenario where Databricks visualizations helped in decision-making.

7. Additional Interview Points

Cluster Policies:

Enforce security and compliance rules.

Databricks Runtime:

Runtime includes Spark versions, libraries, and optimizations.

Data Lake Integration:

Seamless connection with Azure Data Lake or S3.

Scenario-Based Interview Questions

Compute Optimization:

"You are tasked with analyzing log data for a week-long event using Databricks. How would you configure the cluster to optimize cost and performance?"

Notebook Collaboration:

"How do you ensure efficient collaboration in a Databricks notebook shared among team members?"

Debugging in Databricks:

"Your notebook cell execution is stuck. How do you debug the issue?"

part05

Notes on Databricks Workflow & Key Concepts

1. Vault Access Policy for Storing Secrets in Azure

Key Vault Creation:

To store sensitive information like API keys or credentials, Azure Key Vault is used.

You need to create a Vault Access Policy instead of RBAC to assign access permissions.

Steps:

Go to "Review + Create" and click "Create" after defining the Key Vault. Assign yourself as a Key Vault Administrator (necessary to create secrets).

Use Access Control (IAM) to assign the Key Vault Administrator role to your user.

Once created, navigate to Secrets and click on Generate/Import to store your secret (e.g., app secret).

After saving the secret, link Databricks with Key Vault by creating a secret scope.

2. Creating a Secret Scope in Databricks

Scope Creation:

Go to the Home tab in Databricks and search for your workspace URL, then append #secrets/createScope to it.

Fill out the required fields:

Name: e.g., uncope

Managed Principal: Select All Workspace Users

DNS Name and Resource ID: Grab these from the Azure Key Vault.

Click Create to finalize the scope.

3. Using Secrets in Databricks

dbutils.secrets:

Databricks provides a utility dbutils.secrets to retrieve secrets from the Azure Key Vault.

Commands:

dbutils.secrets.list(scope="uncope"): Lists all secrets under the specified scope.

dbutils.secrets.get(scope="uncope", key="app-Secret"): Retrieves the secret value securely (value will be redacted to prevent exposure).

4. Data Reading in Databricks

Spark Reader API:

To read data from cloud storage, we use spark.read API.

Syntax:

python

Copy code

```
df = spark.read.format("csv").option("header",  
"true").option("inferSchema", "true").load("<path_to_data>")
```

Key Parameters:

header: Treats the first row as the header.

inferSchema: Automatically detects the schema (column types).

Load Data: You can specify the data format (CSV, JSON, Delta, etc.) and path to the file in your cloud storage (e.g., Azure Data Lake).

Example:

python

Copy code

```
df = spark.read.option("header", "true").option("inferSchema",
```

```
"true").csv("abfss://<container>@<storage>.dfs.core.windows.net/sales.csv")  
df.show() # To display the dataframe
```

5. Data Transformations in Databricks (Using PySpark)

Using PySpark Transformations:

`withColumn`: Adds or modifies a column in a DataFrame.

`split`: Splits a string into multiple columns based on a delimiter.

Transformation Example: Create a list by splitting a string on spaces:

python

Copy code

```
from pyspark.sql.functions import split, col
```

```
df_transformed = df.withColumn("item_list", split(col("item_type"), " "))
```

```
df_transformed.show()
```

Type Casting Example:

python

Copy code

```
df_transformed = df.withColumn("item_visibility",  
col("item_visibility").cast("string"))
```

```
df_transformed.show()
```

Creating Constants in Data (using `lit` function):

python

Copy code

```
from pyspark.sql.functions import lit
```

```
df_transformed = df.withColumn("flag", lit("yes"))
```

```
df_transformed.show()
```

6. Pyspark Libraries

Importing Required Libraries:

python

Copy code

```
from pyspark.sql import SparkSession  
from pyspark.sql.functions import col, split, lit  
from pyspark.sql.types import StringType
```

7. Delta Lake Integration

Delta Lake is a storage layer that brings ACID transactions to Apache Spark and big data workloads. It helps in managing and versioning data over time.

Benefits:

ACID transactions

Schema enforcement

Time travel (versioning)

Unified batch and streaming data processing

8. Important Databricks Interview Points

Databricks Secrets Management:

Use of dbutils for securely handling sensitive data in Databricks notebooks.

Integration with Azure Key Vault for secret storage and management.

Spark Reader API:

Importance of reading data from cloud storage (Azure Data Lake, ADLS Gen2, etc.).

Handling different file formats (CSV, JSON, Delta, etc.).

Transformations in PySpark:

Importance of transformations like withColumn, split, lit, and cast for data wrangling.

Familiarity with PySpark functions and DataFrame APIs.

Scenario-Based Interview Questions for Databricks

Scenario 1:

Problem: You need to read a CSV file stored in Azure Data Lake, infer its schema, and load it into a Databricks DataFrame. How would you do that?

Answer: Use the spark.read method to read the data, specifying the format as CSV and enabling schema inference.

python

Copy code

```
df = spark.read.option("header", "true").option("inferSchema", "true").csv(" <path_to_data>")
```

df.show()

Scenario 2:

Problem: How would you securely store an application secret in Databricks without hardcoding it?

Answer: Store the secret in Azure Key Vault, create a secret scope in Databricks, and retrieve the secret using dbutils.secrets.get().

python

Copy code

```
secret_value = dbutils.secrets.get(scope="uncope", key="app-Secret")
```

Scenario 3:

Problem: You have a DataFrame with columns containing JSON data. How would you parse these JSON columns into structured data?

Answer: Use from_json function to parse JSON columns into structured columns.

python

Copy code

```
from pyspark.sql.functions import from_json
```

```
df_parsed = df.withColumn("json_parsed", from_json(df["json_column"], schema))
df_parsed.show()
```

Scenario 4:

Problem: How do you handle large datasets and ensure that the job is optimized for performance in Databricks?

Answer: Use partitioning, caching, and Delta Lake optimizations to handle large datasets. Ensure proper schema management and limit shuffling by using repartition() and coalesce().

=====

part06

Delta Lake Overview

Delta Lake is a storage layer that adds ACID transactions and schema enforcement to big data workloads.

It is not a file format by itself. Instead, it works on top of existing formats (such as Parquet), and adds additional features like transaction logs.

Key Concepts Explained

Delta Lake vs. Parquet Format

Parquet Format: A columnar-based file format optimized for fast reads and efficient storage, commonly used for big data.

Delta Lake: Involves Delta transaction logs and metadata to manage the data effectively, enhancing the performance and providing additional features over Parquet.

Why Delta Lake?

It is built on top of Parquet, meaning it maintains the same structure but adds powerful features like ACID transactions, data versioning, and schema enforcement.

Transaction Log in Delta Lake

Transaction Log: A key component of Delta Lake, it stores the metadata (such as schema, columns, data types, and other data) for every file written.

Metadata is stored in JSON files within the Delta log folder.

Unlike traditional file formats, where metadata is stored within each file, Delta Lake stores metadata in a separate log, saving time by reading metadata from a central location rather than multiple individual files.

Key Commands and Concepts

Creating Notebooks:

Use %run to import and execute an entire notebook, saving time and simplifying workflows in Databricks.

Example:

python

Copy code

%run /Workspace/DeltaLakeNotebook

Saving Data in Delta Format:

Use the DataFrame Writer API to write data to a Delta table.

python

Copy code

```
df.write.format("delta").mode("overwrite").save("/path/to/delta-table")
```

The write modes are:

Append: Adds data to the existing table.

Overwrite: Replaces the existing table data with the new data.

ErrorIfExists: Throws an error if the table already exists.

Ignore: Does nothing if the table already exists.

Path: The path in Databricks where the Delta table is saved.

python

Copy code

```
df.write.format("delta").mode("append").option("path",  
"/mnt/data/delta").save()
```

Delta Table Types: Managed vs. External

Managed Delta Table:

Metadata and data are managed by Databricks.

Deleting the table will also remove the underlying data.

Example:

sql

Copy code

```
CREATE TABLE sales_table USING DELTA LOCATION '/mnt/data/sales'
```

External Delta Table:

The data is stored in a location external to Databricks (e.g., an S3 bucket, Azure storage).

Deleting the table does not remove the underlying data.

Example:

sql

Copy code

```
CREATE EXTERNAL TABLE sales_table USING DELTA LOCATION  
'/mnt/external_data/sales'
```

Real-Life Use Case for Delta Tables

Managed Tables are ideal when you want Databricks to handle both the

metadata and data storage.

External Tables are useful when you want to control the data storage separately (such as in a specific cloud storage account) and ensure that deleting the table does not delete the data.

Example Scenario

Transforming Data from Parquet to Delta Format

Source Folder: Contains data in Parquet format.

Destination Folder: Data is written in Delta format.

Command to move data:

python

Copy code

```
df.write.format("delta").mode("overwrite").save("/path/to/destination")
```

Delta Lake Architecture Overview

Files in Delta Lake:

Files are still in Parquet format.

Delta Lake adds a transaction log that stores metadata separately in JSON and CRC files.

Metadata Handling:

Delta Lake reads metadata from the transaction log, not the individual Parquet files, saving computational time.

Additional Interview Points & Questions

Key Concepts for Interviews:

ACID Transactions in Delta Lake:

Delta Lake provides Atomicity, Consistency, Isolation, and Durability for big data workloads.

Use case: In scenarios where multiple operations need to be executed atomically (e.g., during data migrations or updates).

Schema Enforcement:

Ensures that data adheres to a specific schema, preventing corrupt or invalid data from being inserted into Delta tables.

Versioning:

Delta Lake maintains data versioning, allowing users to access older versions of the data and perform time travel queries.

Scenario-Based Interview Questions for Databricks

Scenario: You have a large dataset in Parquet format, and you want to start using Delta Lake for your transformations. How would you convert the data to Delta format in Databricks?

Answer: Write the data using the `df.write.format("delta")` method and specify the path where the data should be stored.

Scenario: You have been asked to ensure that all new records inserted into a Delta table are validated against a predefined schema. How would you accomplish this in Databricks?

Answer: Implement schema enforcement by specifying the schema when writing data to the Delta table. Use Delta's merge functionality to ensure data is inserted only if it matches the schema.

Scenario: What happens if you delete a managed Delta table in Databricks? What happens if you delete an external Delta table?

Answer: Deleting a managed Delta table deletes both the metadata and the actual data. Deleting an external Delta table removes only the metadata; the data remains in its external storage location.

Diagrams

Example Diagram of Delta Lake Architecture:

Delta Lake Architecture:

Data is stored in Parquet files.

Metadata is managed by the transaction log (in JSON format).

Delta Log is the key to fast metadata reading, improving performance.

Part07

Databricks Managed and External Tables: A Detailed Breakdown

I. Introduction to Managed Tables in Databricks

Managed tables in Databricks are tables where the data is stored and managed by Databricks itself. When a managed table is dropped, both the table's metadata and the data are deleted. These tables are typically created with Delta format.

Creating a Managed Table (SQL Syntax):

sql

Copy code

```
CREATE DATABASE salesDB;
```

-- Creating a Managed Table

```
CREATE TABLE salesDB.managed_table (
    id INT,
    name STRING,
```

```
marks INT  
) USING DELTA;
```

Explanation:

`CREATE DATABASE salesDB;` Creates a database named salesDB.
`CREATE TABLE salesDB.managed_table;` Creates a table in the salesDB database, named managed_table.
`USING DELTA;` Specifies that the table uses the Delta format for storage.

In Databricks, the default table format is Delta. You can choose other formats (like CSV or Parquet) but Delta is typically preferred for its ACID transaction support and efficient performance.

Inserting Data into the Table:

sql

Copy code

```
INSERT INTO salesDB.managed_table VALUES  
(1, 'AA', 30),  
(2, 'BB', 33),  
(3, 'CC', 35),  
(4, 'DD', 40);
```

2. Querying the Managed Table:

You can query the managed table just like you would in traditional SQL:

sql

Copy code

```
SELECT * FROM salesDB.managed_table;
```

When you query the table, Databricks will return the records stored in the managed table from the underlying Delta storage.

3. External Tables in Databricks

External tables in Databricks are those where the data is stored in an external location (like a data lake), and Databricks only manages the table's metadata. Unlike managed tables, when you drop an external table, only the table metadata is deleted, not the underlying data.

Creating an External Table (SQL Syntax):

sql

Copy code

-- Create an External Table

```
CREATE TABLE salesDB.external_table (
    id INT,
    name STRING,
    marks INT
) USING DELTA
LOCATION 'abfss://<container-name>@<storage-
account>.dfs.core.windows.net/salesDB/external_table';
```

Explanation:

LOCATION: Specifies the path where the table's data resides. In this case, it's an external storage account.

4. Inserting Data into External Tables:

Data can be inserted into external tables similarly to managed tables:

sql

Copy code

```
INSERT INTO salesDB.external_table VALUES
(5, 'EE', 45),
(6, 'FF', 50);
```

5. Managing Data for External Tables

When data is inserted or updated in an external table, it's stored in the specified location (e.g., data lake storage).

Advantage of External Tables: You have more control over where your data resides. It does not get deleted when the table is dropped.

6. Dropping Tables

Dropping a managed table deletes both the table's metadata and data.

Dropping an external table deletes only the metadata, and the data remains intact in its original storage.

sql

Copy code

```
-- Drop Managed Table (Deletes both data and metadata)
```

```
DROP TABLE salesDB.managed_table;
```

```
-- Drop External Table (Deletes only metadata, not data)
```

```
DROP TABLE salesDB.external_table;
```

7. Tombstoning Concept in Delta Tables

Tombstoning is a Delta feature that marks deleted rows rather than physically removing them immediately. This allows you to handle deletes more efficiently and maintain ACID guarantees.

Delete Operation:

sql

Copy code

```
DELETE FROM salesDB.external_table WHERE id = 6;
```

When a delete is executed in Delta, the data is not immediately removed but is instead flagged as deleted in the Delta log file. This ensures that the data is recoverable and allows Delta to handle the deletion efficiently in the background.

8. Delta Table Functionalities

Delta tables come with a wide range of functionalities for managing and manipulating data, including:

Insert: Insert new data into a table.

Update: Modify existing data in the table.

Delete: Remove data based on specific conditions.

Merge: Combine new data into existing data.

Example of an Update Operation in Delta:

sql

Copy code

```
UPDATE salesDB.external_table
```

```
SET marks = 55
```

```
WHERE id = 5;
```

Example of a Merge Operation:

sql

Copy code

```
MERGE INTO salesDB.external_table AS target
```

```
USING salesDB.staging_table AS source
```

```
ON target.id = source.id
```

```
WHEN MATCHED THEN
```

```
UPDATE SET target.marks = source.marks
```

```
WHEN NOT MATCHED THEN
```

```
INSERT (id, name, marks) VALUES (source.id, source.name,  
source.marks);
```

Additional Interview Points & Real-Life Scenarios

Managed vs External Tables:

When would you choose an external table over a managed table?

Answer: Choose external tables when you need to control the location of your data, typically in a data lake, and ensure that the data persists even if the table is dropped.

Delta Table Advantages:

What are the advantages of using Delta format for storing tables in Databricks?

Answer: Delta provides ACID transactions, time travel (versioning), schema enforcement, and efficient data handling.

Handling Data Consistency:

How does Delta ensure data consistency in large-scale ETL processes?

Answer: Delta handles consistency by using a write-ahead log (WAL), ensuring that data is never lost and can be rolled back in case of errors.

Real-Life Scenario Question:

You have a large dataset that needs frequent updates. The updates must be applied in a way that older versions of data can be preserved. Which table type would you use and why?

Answer: Use a Delta table because it supports versioning, ACID transactions, and efficient incremental updates.

Additional Code Examples and Syntax

Create a Managed Table:

sql

Copy code

```
CREATE TABLE salesDB.managed_table (
    id INT,
    name STRING,
```

marks INT
) USING DELTA;

Insert Data into Managed Table:

sql

Copy code

```
INSERT INTO salesDB.managed_table VALUES (1, 'Alice', 90), (2, 'Bob', 85);
```

Create an External Table:

sql

Copy code

```
CREATE TABLE salesDB.external_table (  
id INT,  
name STRING,  
marks INT  
) USING DELTA  
LOCATION 'abfss://container-name@storage-  
account.dfs.core.windows.net/path';
```

Real-Life Interview Questions:

Delta Table Use Case:

How does Delta's ACID transaction support benefit the ETL pipeline in a multi-user environment?

Managing External Tables:

What would you do if an external table becomes corrupted? How can you restore the data?

Tombstoning in Delta:

Explain tombstoning in Delta tables. How does it affect data retrieval and performance?

Diagram: Managed vs External Table Comparison

Feature Managed Table External Table

Data Storage Managed by Databricks Stored externally (e.g., in Data Lake)

Dropping Table Deletes both metadata and data Deletes only metadata, data remains

Data Control No control over data Full control over data location

Default Format Delta Delta (with a specified location)

=====

part08

Delta Tables and Data Processing in Databricks

In Databricks, Delta tables offer powerful capabilities for data versioning, time travel, and soft deletion. Let's go over the key concepts discussed:

1. Data Ingestion and Delta Log Management

Delta Log: When reading data in Delta tables, Databricks first reads the Delta log. It doesn't matter if there are 300 partitions or 300,000. The Delta log contains metadata that tells Databricks which partitions to read and what data to include.

Partitioned Data: Delta tables manage partitions internally. The system reads partitioned files based on Delta log instructions, ignoring any files that are not referenced in the log.

For example, if a partition labeled 42b and c92 is added to the Delta table, the Delta log will specify those partitions. Even if there are many partitions, the system focuses on the ones listed in the Delta log.

2. Soft Deletion and Tombstoning

Soft Deletion: When a delete operation is performed, Delta tables do not physically delete the records from the partition. Instead, they create a new partition without the deleted records and update the Delta log to reflect these changes.

Tombstoning: This process is also known as tombstoning, where old partitions are marked as "ghosted" or "blacklisted," meaning they are no longer used. The actual data still exists, but Delta doesn't read from it anymore.

Why not delete directly? This approach ensures that data can be recovered if needed (through time travel). It's a safety feature of Delta tables, allowing the system to maintain historical data for auditing or rollback purposes.

Example:

If you delete a record ($ID = 8$), it doesn't disappear from the partition immediately. Instead, Delta creates a new partition without that record. The old partition is "tombstoned," and the system reads from the updated partition.

3. Time Travel and Data Versioning

Time Travel: Delta Lake allows you to "travel back in time" to a previous version of the table, which is useful if you accidentally delete or update data.

How does it work? Every operation (e.g., insert, delete, update) in Delta Lake is versioned. Each version of the table is stored in the Delta log.

Example Command:

sql

Copy code

```
DESCRIBE HISTORY sales_db.external_table
```

This command shows the history of changes, including table creation, inserts, deletes, and other modifications.

Restoring Data: If you accidentally delete data (like ID = 8), you can restore it by rolling back to an earlier version of the table.

Example Command:

sql

Copy code

```
RESTORE TABLE sales_db.external_table TO VERSION AS OF 2
```

This command restores the table to version 2, where the record with ID = 8 was still present.

4. Vacuum Command and Data Cleanup

Vacuum Command: This command is used to physically delete obsolete data files (e.g., partitions) that are no longer referenced by any version of the table.

Usage: It is useful when you want to clean up data after performing time travel or other data operations that leave old files.

Syntax:

sql

Copy code

```
VACUUM sales_db.external_table
```

Guideline: Do not run this command in production unless you're sure you don't need the old data for time travel or rollback purposes.

Important Consideration:

Default Retention Period: By default, the vacuum command will not delete files less than 7 days old.

Modify Retention Period: You can modify the retention period with the RETAIN option. For example, to delete immediately:

sql

Copy code

```
VACUUM sales_db.external_table RETAIN 0 HOURS
```

5. Delta Table Optimization

Optimize Command: The OPTIMIZE command is used to improve the performance of Delta tables by compacting small files into larger ones.

This operation can greatly improve query performance, especially with larger datasets.

Syntax:

sql

Copy code

```
OPTIMIZE sales_db.external_table
```

Result: This will compact files and reduce the number of small files, which can lead to faster read times and more efficient data processing.

Example Performance Increase: After optimizing the table, query performance can improve by 30-40% for large datasets.

Z-Order Optimization: The ZORDER BY command allows you to optimize the layout of data files for efficient querying by clustering data based on specific columns (e.g., ID).

Syntax:

sql

Copy code

```
OPTIMIZE sales_db.external_table ZORDER BY (ID)
```

Performance Benefit: It groups the data by the specified column (e.g., ID), reducing the scan time for queries that filter on that column.

Real-Life Interview Questions for Databricks (Scenario-Based)

Scenario: You're working with a Delta table in Databricks. A large batch of records was mistakenly deleted. How would you recover the deleted records?

Answer: You can use Delta Lake's time travel feature to restore the table to a previous version before the deletion.

sql

Copy code

```
RESTORE TABLE sales_db.external_table TO VERSION AS OF  
<version_number>
```

Scenario: A Delta table is experiencing slow query performance due to many small files. How would you optimize the table?

Answer: Use the OPTIMIZE command to compact small files into larger ones for better performance.

sql

Copy code

`OPTIMIZE sales_db.external_table`

Scenario: You want to permanently delete partitions that are no longer needed from a Delta table. What command would you use?

Answer: Use the VACUUM command to remove outdated files, ensuring that data no longer referenced is deleted.

sql

Copy code

```
VACUUM sales_db.external_table RETAIN 0 HOURS
```

Scenario: How does Delta Lake ensure data consistency when performing a delete operation?

Answer: Delta Lake uses soft deletion (tombstoning), where the deleted records are not physically removed, but a new partition is created that excludes the deleted data. This ensures data consistency and enables time travel.

Diagram: Delta Table with Time Travel and Partition Updates

plaintext

Copy code

```
-----  
| Delta Log | --> | Partition 1 |
```

```
----- | (contains ID=8) |
```

```
-----  
-----  
-----
```

```
-----  
| Delta Log | --> | Partition 2 |
```

```
----- | (contains ID=56, 78) |
```

```
-----  
-----  
-----
```

```
| Tombstoned Partition | |
```

Updated Partition 2

=====

part

Delta Lake Optimization & Z-Ordering

Partitioning and Optimization

Scenario: Let's assume you have 4 GB of data split into 4 partitions of 1 GB each.

The alternative is 2 partitions of 2 GB each.

Better Approach: It's always better to have fewer, larger partitions, as more partitions can cause overhead.

Optimization in Delta Lake:

Optimize Command: When we run the OPTIMIZE command, it consolidates smaller partitions into bigger partitions to enhance performance, reduce overhead, and improve read/write efficiency. The main goal of optimization is to reduce the number of partitions to increase the overall processing speed.

Z-Order by:

Sorting: When we run Z-ORDER BY after optimization, it sorts the data by the specified column (e.g., ID), making it easier for Delta Lake to skip unnecessary data during queries.

Data Skipping: The biggest advantage of Z-ORDER BY is Data Skipping. This means that when filtering based on a sorted column, the engine doesn't need to scan entire partitions, improving query performance. For example, if you query SELECT * WHERE ID = 3, Delta can skip

partitions where the value 3 doesn't exist, improving speed.

Use Case Example:

ID Column Example:

After applying the Z-ORDER BY on the ID column, data is sorted within partitions.

For instance, if the data in the first partition was 1, 2, 3, 4 and in the second partition it was 5, 6, 7, 8, Z-ordering sorts them in ascending or descending order.

Advantages:

Reduces the time to perform queries because Delta skips irrelevant partitions.

Increases performance by consolidating smaller partitions.

Workflows and Data Orchestration in Databricks

Workflow Management:

Workflows in Databricks allow you to orchestrate and automate the execution of notebooks, jobs, and other operations in sequence.

Databricks notebooks can be used for data preparation, processing, and analysis, while workflows help you run these tasks on a scheduled or event-triggered basis.

Incremental Loading with AutoLoader

Problem: In real-world data engineering, you often deal with continuous or incremental data. Instead of reprocessing all the data every time, we need a way to only load new data that hasn't been processed before.

AutoLoader:

What is AutoLoader?

AutoLoader is a Databricks framework designed for streaming data. It allows you to ingest and process files incrementally.

AutoLoader uses streaming data frames, which allow continuous

processing of data.

Key Concepts:

Streaming Data Frame:

A streaming data frame is similar to a regular data frame, but it operates on continuous streams of data (such as incoming files).

Example: Instead of loading data in batch, a streaming data frame continuously reads new data as it arrives.

Cloud Files:

With AutoLoader, you specify the file format as `cloud_files`, and set the format type (e.g., Parquet).

This tells AutoLoader to read files incrementally from a source location and write them to a destination.

Schema Location (Checkpointing):

Checkpointing is essential for AutoLoader because it tracks the state of the data stream. After processing a file, the checkpoint ensures the system knows which files have been processed, so only new data is processed on the next read.

Example Scenario:

Source and Destination:

Source: You have a folder where new data files (e.g., Parquet files) are uploaded periodically.

Destination: Data needs to be transferred to a destination folder (e.g., Delta format).

Code Example:

`python`

Copy code

```
df = (spark.readStream  
.format("cloudFiles")  
.option("cloudFiles.format", "parquet")  
.option("cloudFiles.schemaLocation", "/mnt/destination/checkpoints")  
.load("/mnt/source/"))
```

Writing to Delta

```
df.writeStream  
.format("delta")  
.option("checkpointLocation", "/mnt/destination/checkpoints")  
.start("/mnt/destination/")
```

The above code reads Parquet files incrementally from a source folder, processes them, and writes them in Delta format to a destination.

Triggering:

The trigger defines the frequency of checking for new data (e.g., every 10 seconds).

python

Copy code

```
df.writeStream  
.format("delta")  
.trigger(processingTime='10 seconds')  
.option("checkpointLocation", "/mnt/destination/checkpoints")  
.start("/mnt/destination/")
```

Real-World Interview Questions

Z-Ordering and Optimization:

Explain the benefits of using Z-ORDER BY in Delta Lake.

Can you demonstrate how partitioning affects the performance of queries

in Delta Lake?

What is data skipping, and how does it improve query performance?

Incremental Loading with AutoLoader:

How does AutoLoader work in Databricks, and what are its use cases?

Describe the steps to implement incremental loading with AutoLoader.

What is checkpointing in the context of streaming data in Databricks?

Databricks Workflow:

How would you orchestrate a data pipeline in Databricks using notebooks?

What steps are involved in managing a data pipeline from data ingestion to processing using Databricks workflows?

Scenario-Based Question:

You are processing log files that arrive continuously every minute. How would you implement an incremental loading solution that processes only new log files without reprocessing the old ones?

Additional Notes:

Z-Order by is essential for optimizing queries on large datasets in Delta Lake by reducing unnecessary reads.

Streaming Data Frames are crucial for handling continuous data and automating data ingestion processes.

AutoLoader is a powerful tool for efficient, incremental loading in Databricks, especially for large datasets or real-time data processing.

=====