IMAAD IMRAN HAJWANE

202101132 / 21

BIG DATA ANALYTICS – LAB 5

## Hive Partitioning

Partitioning is a method of dividing a large dataset into smaller, manageable parts, based on the values of one or more columns. It is an optimization technique that helps improve query performance in Hive by reducing the amount of data read during query execution. When a query is run on partitioned data, only the relevant partitions are scanned, leading to faster processing.

## Types of Partitioning:

1. **Static Partitioning**:
   - In static partitioning, the partitions are explicitly defined by the user when loading data into a table.
   - For example, consider a table storing sales data where the partition is based on the `year` column. The user must specify the partition value when inserting data, such as `year=2020`.
   - It requires more manual intervention but allows for precise control over how data is organized.
2. **Dynamic Partitioning**:
   - Dynamic partitioning automatically creates partitions based on the data being loaded.
   - Instead of manually specifying the partition key, Hive determines the partitions dynamically based on the column values at runtime.
   - For example, if the sales data includes different years (2018, 2019, 2020) in one dataset, Hive will automatically create partitions for each year.
   - It is useful when data for multiple partitions is being loaded simultaneously.

## Advantages of Partitioning:

- **Improved Query Performance**: When querying partitioned data, Hive scans only the required partitions instead of the entire dataset. This reduces the input/output (I/O) operations and speeds up query execution.
- **Organized Data**: Partitioning helps organize large datasets logically, making it easier to manage and retrieve specific data subsets.

## Hive Bucketing

Bucketing is another technique in Hive to optimize query performance by dividing data into more manageable segments (buckets). Unlike partitioning, bucketing distributes data within a partition into equally-sized buckets using a hashing function. Each bucket can then be treated as a separate file.

## How Bucketing Works:

- In bucketing, Hive divides the data into buckets based on the hash value of the bucketed column (often an integer).
- For example, if a table is bucketed by the `emp_id` column into 4 buckets, Hive applies a hash function to `emp_id` values to distribute the rows across 4 different buckets. Rows with similar hash values end up in the same bucket.
- This makes operations like joins more efficient, as similar data is stored in the same bucket across tables.

## Advantages of Bucketing:

- **Efficient Joins**: Bucketing is particularly useful in join operations, as data with the same bucket value can be co-located. This minimizes the amount of data shuffled across the cluster during joins.
- **Sampling**: Bucketing enables Hive to perform sampling efficiently. Since the data is evenly distributed across buckets, Hive can select a subset of data by sampling specific buckets.

- Initialising, Displaying & Creating table in database



- Loading data using "LOAD" commands and displaying data stored



- Final display of stored data

## Dynamic

- Creation, storing & Execution



## BUCKETING

- Initializing, Creating & Displaying buckets

- Execution of bucketing



- Final Display of buckets values