# Data Skew in Apache Spark

**Suppose we have an ecommerce sales dataset which is partitioned on product**

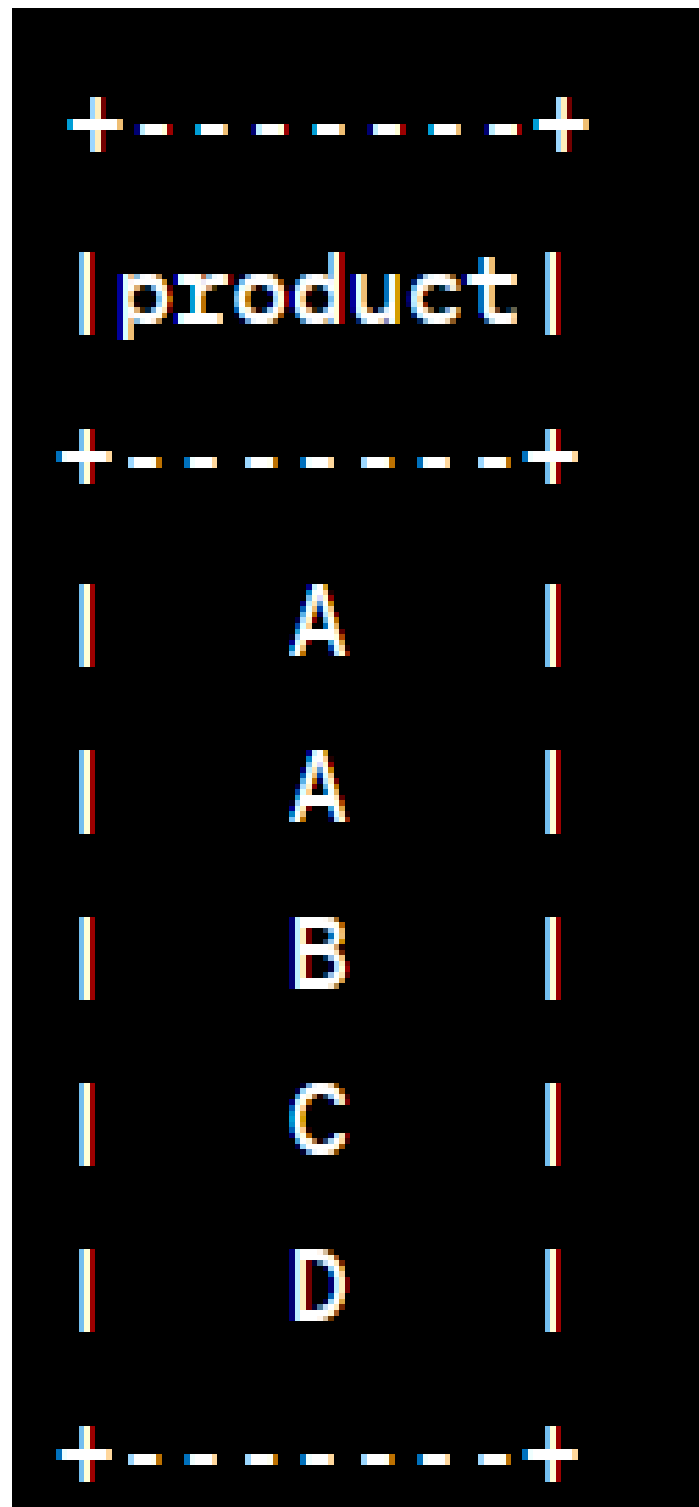| | | |
|---|---|---|
| Product E | ▮ | Partition 5 |
| Product D | ▮ | Partition 4 |
| Product C | ▮ | Partition 3 |
| Product B | ▮ | Partition 2 |
| Product A | ▬▬▬▬▬ | Partition 1 |

Product Sales dataframe after partitioned on Product

**We can see since the data volume for Product A is way higher than others, partition size for it is way larger**

▼

# Data Skew!!

# Mitigation Plan 1: Salting

```
+-------+
|product|
+-------+
|   A   |
|   A   |
|   B   |
|   C   |
|   D   |
+-------+
```

In the input dataset, we can observe data skew, as the key "A" appears more frequently compared to other keys. This uneven distribution of data can lead to data skew and adversely affect the performance of our data processing tasks.

# Mitigation Plan 1: Salting

**To mitigate data skew using the salting technique, we can apply the provided code snippet in PySpark.**

🐍 main.py     🐍 module1.py     🐍 module2.py     🐍 module3.py

```python
 6   data = [["A"],["A"],["B"],["C"],["D"]]
 7   rdd = spark.sparkContext.parallelize(data)
 8   columns=["product"]
 9
10   data = spark.createDataFrame(rdd,columns)
11   data.show(truncate=False)
12
13   # Define the salt length and salt value
14   salt_length = 10
15   salt_value = "salt_"
16
17   # Add the salted column to the DataFrame
18   saltedData = data.\
19       withColumn("salted_product",
20                   concat(lit(salt_value),
21                   rand().cast("string")))
22   saltedData.show()
23   saltedData = saltedData.\
24       withColumn("salted_product",
25                   concat(col("salted_product")
26                       .substr(0, salt_length),
27                   col("product").cast("string")))
28
29   saltedData.show()
```

# Mitigation Plan 1: Salting

**After applying the salting technique, the output DataFrame saltedData would look like this:**

```
+-------+-----------------+
|product|salted_product|
+-------+-----------------+
|      A|     salt_0.567A|
|      A|     salt_0.585A|
|      B|     salt_0.365B|
|      C|     salt_0.796C|
|      D|     salt_0.032D|
+-------+-----------------+
```

**The salting technique adds a random value (salt) to each product key, resulting in new salted keys. These salted keys introduce randomness and ensure that the skewed keys ("A" in this case) are distributed more evenly across partitions.**

## Problem: Additional Computational Overhead for Salting

# Mitigation Plan 2: AQL

**Adaptive Query Language (AQL) is a powerful feature introduced in Apache Spark to combat data skew. It leverages dynamic adaptation techniques during query execution to identify and mitigate data skew effectively.**

## How?

1    **Dynamic Query Optimization**: AQL analyzes the query execution plan and detects potential data skew. It intelligently adapts the plan to redistribute the workload more evenly across the available resources, minimizing the impact of skew on performance.

2    **Range Partitioning**: AQL can leverage range partitioning, dividing skewed data into smaller ranges and distributing them across multiple partitions. By doing so, it ensures that the processing of skewed keys is distributed among different worker nodes, enabling better parallelism.

# Mitigation Plan 2: AQL

3   **Sampling**: AQL uses sampling techniques to estimate the distribution of data across partitions. By analyzing a sample of the data, it can make informed decisions about how to best redistribute the skewed data for optimal load balancing.

4   **Dynamic Repartitioning**: AQL dynamically determines when to repartition data during query execution. It intelligently redistributes data to achieve a more balanced distribution, optimizing resource utilization and reducing the impact of data skew.

# Mitigation Plan 3: Skew Hint

- This can be applied only if you are using databricks
- Employ Databricks specific Skew Hint

*Link* for Databricks Documentation

@shuvajithazra