# Experiment 11

**Aim:**

To perform data analytics using **Apache Spark** on the Amazon Food Dataset to identify **pairs of products that are frequently reviewed together** by the same users.

**Theory:**

Apache Spark is a distributed computing framework used for large-scale data processing. It provides RDD and DataFrame APIs to perform parallel operations on large datasets across multiple cluster nodes.

The Amazon Food Review dataset contains user IDs and product IDs reviewed by users.
To find frequently co-reviewed products, we must:

- Transform raw (user, product) data into:
  **user_id → list of product_ids**

- Generate all product pairs reviewed by each user.

- Count how many users reviewed each product pair together.

- Filter pairs with frequency greater than 1.

- Sort by frequency to identify strongly associated product pairs.

Spark transformations used:

- map(), groupByKey(), flatMap()

- reduceByKey(), filter(), sortBy()

- DataFrame functions like collect_list() and explode()

This experiment demonstrates Spark's ability to handle large datasets and compute meaningful relationships at scale.
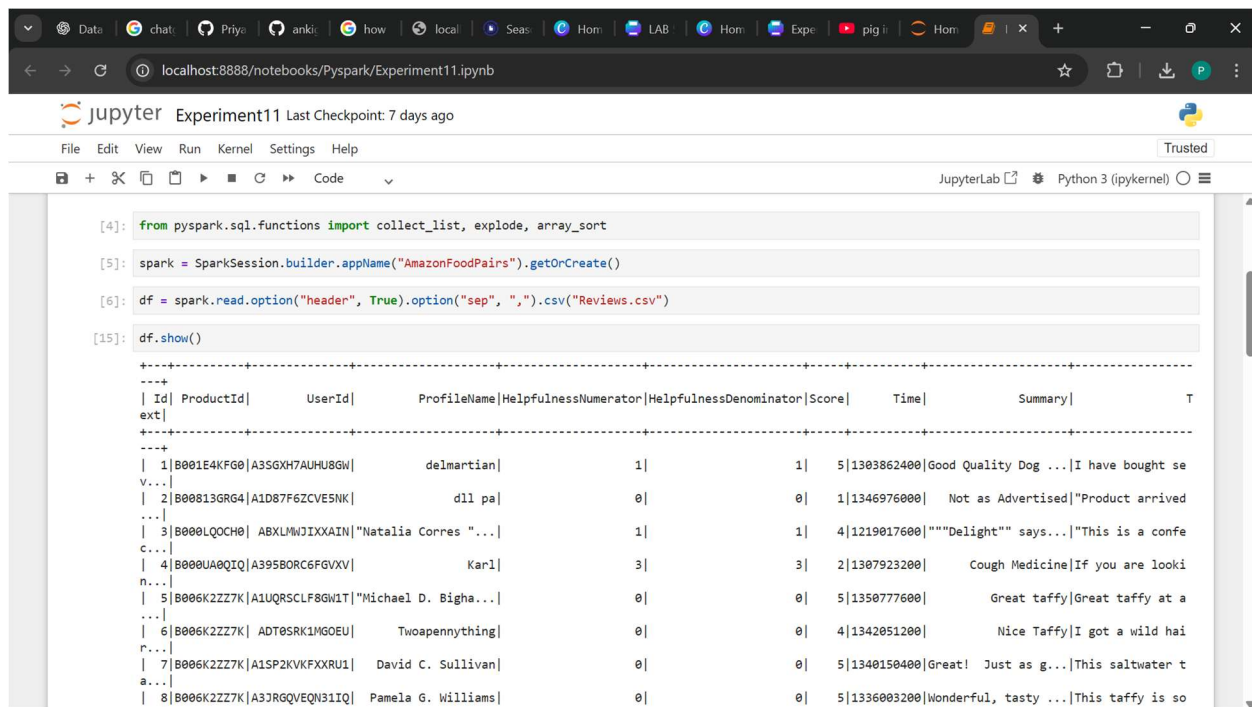
**Procedure:**

**Step 1: Download Dataset**

Download the **Amazon Fine Food Reviews** dataset from:
https://snap.stanford.edu/data/web-FineFoods.html

**Step 2: Load Data into Spark**

Read CSV/TSV file containing UserId and ProductId.

## Step 3: Extract User–Product Pairs

Convert rows into key-value pairs:
(UserId, ProductId)

## Step 4: Group Products by User

Use groupByKey() or DataFrame groupBy().agg(collect_list()) to obtain:
UserId → [P1, P2, P3, ...]

## Step 5: Generate Product Pairs

For each user's product list, generate all 2-item combinations:

(P1, P2), (P1, P3), (P2, P3), ...

## Step 6: Count Pair Frequencies

Use reduceByKey() or DataFrame groupBy().count() to compute:

(ProductA, ProductB) → frequency

## Step 7: Filter & Sort

Keep only those pairs appearing more than once and sort them in descending order of frequency.

## Step 8: Save Output

Write results to HDFS/local storage as CSV/text.

**Output:**

A list of product pairs reviewed together more than once, displayed or saved in sorted order. Example:

(B001E4KFG0, B00813GRG4) → 12

(B00474GHKQ, B00813GRG4) → 9

(B0078LX05C, B00474GHKQ) → 8

```
[11]: result = pairs_df.groupBy("pair").count().filter("count > 1").orderBy("count", ascending=False)

[12]: result.show()
      +--------------------+-----+
      |                pair|count|
      +--------------------+-----+
      |[B002QWP89S, B002...|  682|
      |[B0026RQTGE, B002...|  682|
      |[B002QWHJOU, B002...|  682|
      |[B002QWHJOU, B002...|  682|
      |[B0026RQTGE, B002...|  682|
      |[B0026RQTGE, B002...|  682|
      |[B000UBD88A, B001...|  608|
      |[B001RVFEP2, B001...|  604|
      |[B000VK8AVK, B007...|  604|
      |[B0026KPDG8, B007...|  604|
      |[B000VK8AVK, B001...|  604|
      |[B0026KPDG8, B006...|  604|
      |[B0013NUGDE, B002...|  604|
      |[B000VK8AVK, B006...|  604|
      |[B0026KNQSA, B002...|  604|
      |[B0013NUGDE, B007...|  604|
      |[B0013NUGDE, B006...|  604|
      |[B0013NUGDE, B001...|  604|
      |[B0026KNQSA, B006...|  604|
      |[B0026KPDG8, B007...|  604|
      +--------------------+-----+
      only showing top 20 rows
```

Fig of Result

**Conclusion:**

In this experiment, we successfully applied Apache Spark to analyze a large, real-world dataset. We discovered frequently co-reviewed product pairs using transformations and aggregations. This demonstrates Spark's ability to efficiently handle parallel computation for big data analytics.