# Solutions

## 1. Name the top 3 aisles with the most products ordered

**Screenshot - Hive:**



**Screenshot - PySpark:**

## 2. What is the mean and the variance of number of products in an aisle?

**Screenshot - Hive:**



**Screenshot - PySpark:**

## 3. Which aisle products are most bought with products from the "speciality cheeses" aisle?

**Screenshot - Hive:**



**Screenshot - PySpark:**

4. Sales department is making a recommendation that frozen products should be placed next to bakery products. Write 1-2 SQL queries to get the data that would support/oppose this recommendation.

**Screenshot - Hive:**



**Screenshot - PySpark:**

## 5. Name the top 5 products which are the most reordered

**Ravi**

**Screenshot - Hive:**



**Screenshot - PySpark:**

**Aman Belwal**

**Screenshot – Hive:**

**Screenshot – PySpark:**

## 6. Which 3 products should a customer retargeting (to bring customers back) campaign offer discounts on?

**Screenshot - Hive:**



**Screenshot - PySpark:**

7. Do orders from different departments vary over time of the day? Are there morning and evening departments (popular in the morning and popular in the evening)?

**Screenshot - Hive:**



**Screenshot - PySpark:**



—-------- Only needed for 8 ppl groups

## 8. What is the average number of products in an order? What is the max?
### Screenshot - Hive:

Dataproc

Job details   CLONE   DELETE   STOP   REFRESH

Jobs on Clusters

Clusters

Jobs

Workflows

Autoscaling policies

Serverless

Batches

Interactive

Interactive Templates

Metastore Services

Metastore

Federation

Utilities

Component exchange

Workbench

Release Notes

Job ID          job-bf6bd614
Job UUID        3a7aff55-7064-480c-99e0-1ef96a73db82
Type            Dataproc Job
Status          Succeeded

MONITORING      CONFIGURATION

The charts below represent the metrics from the cluster this job ran on, scoped to the time that this job was running. It is possible for more than one job to run on a cluster at a time, so these metrics may not reflect this job's resource usage accurately. Metrics for a job may lag behind the job run by several minutes.

SAVE AS DASHBOARD   RESET ZOOM   1 hour   6 hours   12 hours   1 day   2 days   4 days   7 days   14 days   30 days   Custom 12:03 AM - 12:08 AM

YARN memory                                 YARN pending memory

Output          LINE WRAP: OFF

```
[22;0m[2K--------------------------------------------------------------------------
INFO  : Completed executing command(queryId=hive_20241008040816_3e0befa1-7e92-4fea-a0ff-6282778f1130); Time taken: 39.236 seconds
INFO  : OK
INFO  : Concurrency mode is disabled, not creating a lock manager
+--------------------------+--------------------------+
| avg_products_per_order   | max_products_per_order   |
+--------------------------+--------------------------+
| 10.088883421247614       | 145                      |
+--------------------------+--------------------------+
1 row selected (39.64 seconds)
0: jdbc:hive2://cluster01-m:10000> Closing: 0: jdbc:hive2://cluster01-m:10000
```

Output is complete

EQUIVALENT COMMAND LINE

Job job-bf6bd614 successfully submitted

### Screenshot - PySpark:

Dataproc

Jobs on Clusters

Clusters

Jobs

Workflows

Autoscaling policies

Serverless

Batches

Interactive

Interactive Templates

Metastore Services

Metastore

Federation

Utilities

Component exchange

Workbench

Release Notes

Job details    CLONE    DELETE    STOP    REFRESH

**Job ID**        job-e475a8f8
**Job UUID**      7b0fff6f-93a4-4735-96c0-5c0f6f598f8f
**Type**          Dataproc Job
**Status**        Succeeded

MONITORING    CONFIGURATION

**Output**    LINE WRAP: OFF

```
24/10/11 20:03:03 INFO Configuration: resource-types.xml not found
24/10/11 20:03:03 INFO ResourceUtils: Unable to find 'resource-types.xml'.
24/10/11 20:03:04 INFO YarnClientImpl: Submitted application application_1728670252448_0006
24/10/11 20:03:05 INFO DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at cluster-8d2c-m.us-central1-f.c.geometric-team-436520-j7.internal./10.
24/10/11 20:03:07 INFO RequestTracker: Detected high latency for [url=https://storage.googleapis.com/storage/v1/b/dataproc-temp-us-central1-653692749797-kususifq
24/10/11 20:03:07 INFO GhfsGlobalStorageStatistics: periodic connector metrics: {gcs_api_client_non_found_response_count=1, gcs_api_client_side_error_count=1, gc
24/10/11 20:03:07 INFO GoogleCloudStorageImpl: Ignoring exception of type GoogleJsonResponseException; verified object already exists with desired state.
24/10/11 20:03:08 INFO GoogleHadoopOutputStream: hflush(): No-op due to rate limit (RateLimiter[stableRate=0.2qps]): readers will *not* yet see flushed data for
Maximum number of products in any order:
+--------+------------+
|order_id|Max_products|
+--------+------------+
| 1564244|         145|
+--------+------------+

24/10/11 20:04:25 INFO GoogleHadoopOutputStream: hflush(): No-op due to rate limit (RateLimiter[stableRate=0.2qps]): readers will *not* yet see flushed data for
Average number of products per order: 10.088883421247614
24/10/11 20:04:27 INFO DataprocSparkPlugin: Shutting down driver plugin. metrics=[action_http_patch_request=0, files_created=1, gcs_api_server_timeout_count=0, o
```

Output is complete

Job job-e475a8f8 successfully submitted    ✕

EQUIVALENT COMMAND LINE