# 📖 05a_retail_forecast_vertex_ai_notebook_execution.md

# 🔗Retail Forecast using Serverless Spark through Vertex AI managed notebooks

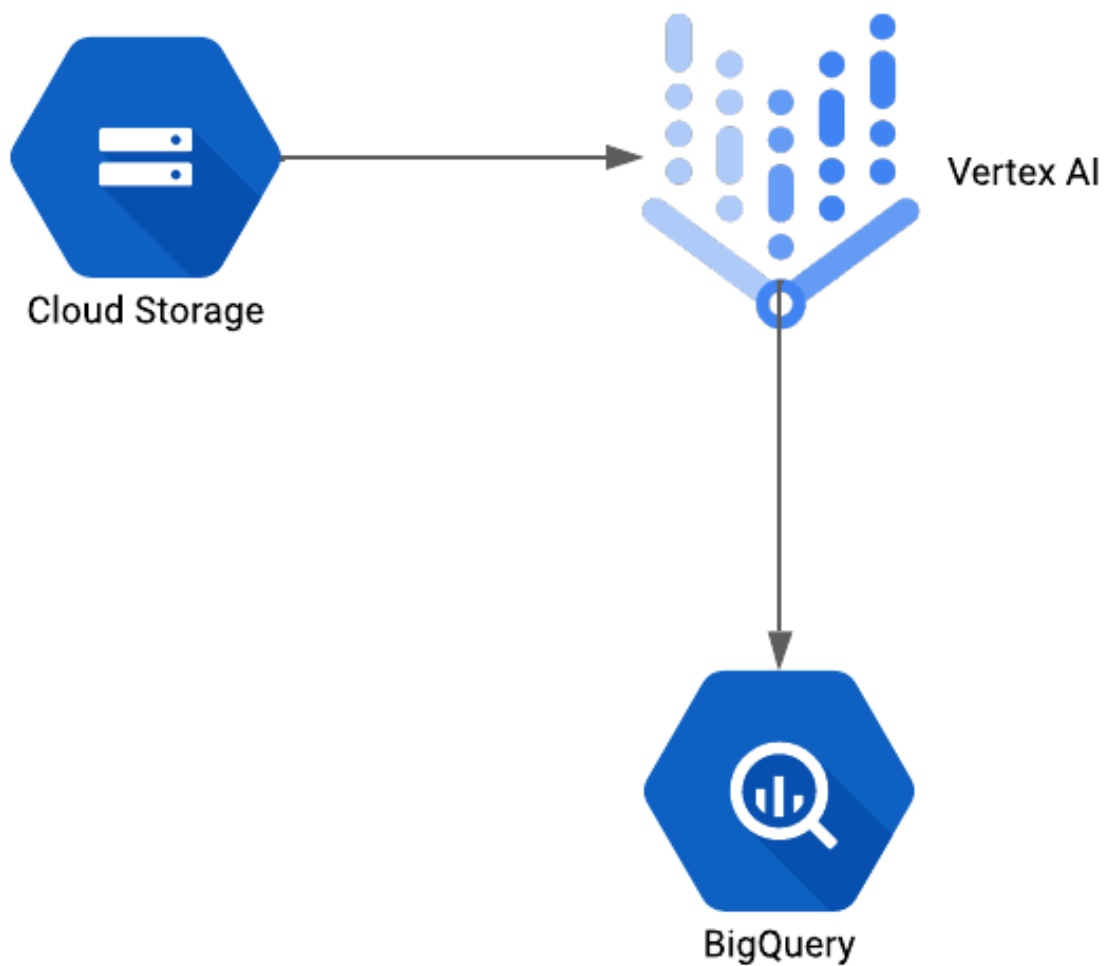Following are the lab modules:

# 🔗1. Understanding Data

The datasets used for this project are

1. [Aisles data](#).
2. [Departments data](#) .
3. [Orders data](#).
4. [Products data](#).
5. [Order_products__prior](#).
6. [Order_products__train](#).

- Aisles: This table includes all aisles. It has a single primary key (aisle_id)
- Departments: This table includes all departments. It has a single primary key (department_id)
- Products: This table includes all products. It has a single primary key (product_id)
- Orders: This table includes all orders, namely prior, train, and test. It has single primary key (order_id).
- Order_products_train: This table includes training orders. It has a composite primary key (order_id and product_id) and indicates whether a product in an order is a reorder or not (through the reordered variable).
- Order_products_prior : This table includes prior orders. It has a composite primary key (order_id and product_id) and indicates whether a product in an order is a reorder or not (through the reordered variable).

# 🔗2. Solution Architecture

# 🔗3. Execution

## 🔗3.1. Run the Batch by creating session.

### 🔗Creating Notebook in Vertex AI

Select Workbench from the left scroll bar of the Vertex AI main page. Select the Managed Notebooks tab. In the Managed Notebooks tab , click the New Notebook icon.

**⌬Next, fill in the following values in the Notebook creation window as shown in the images below:**

- **Notebook Name** - A unique identifier for your Notebook
- **Region** - The region name provided by the Admin team
- **Permission Type** - Single User Only (Single user only mode restricts access to the specified user)

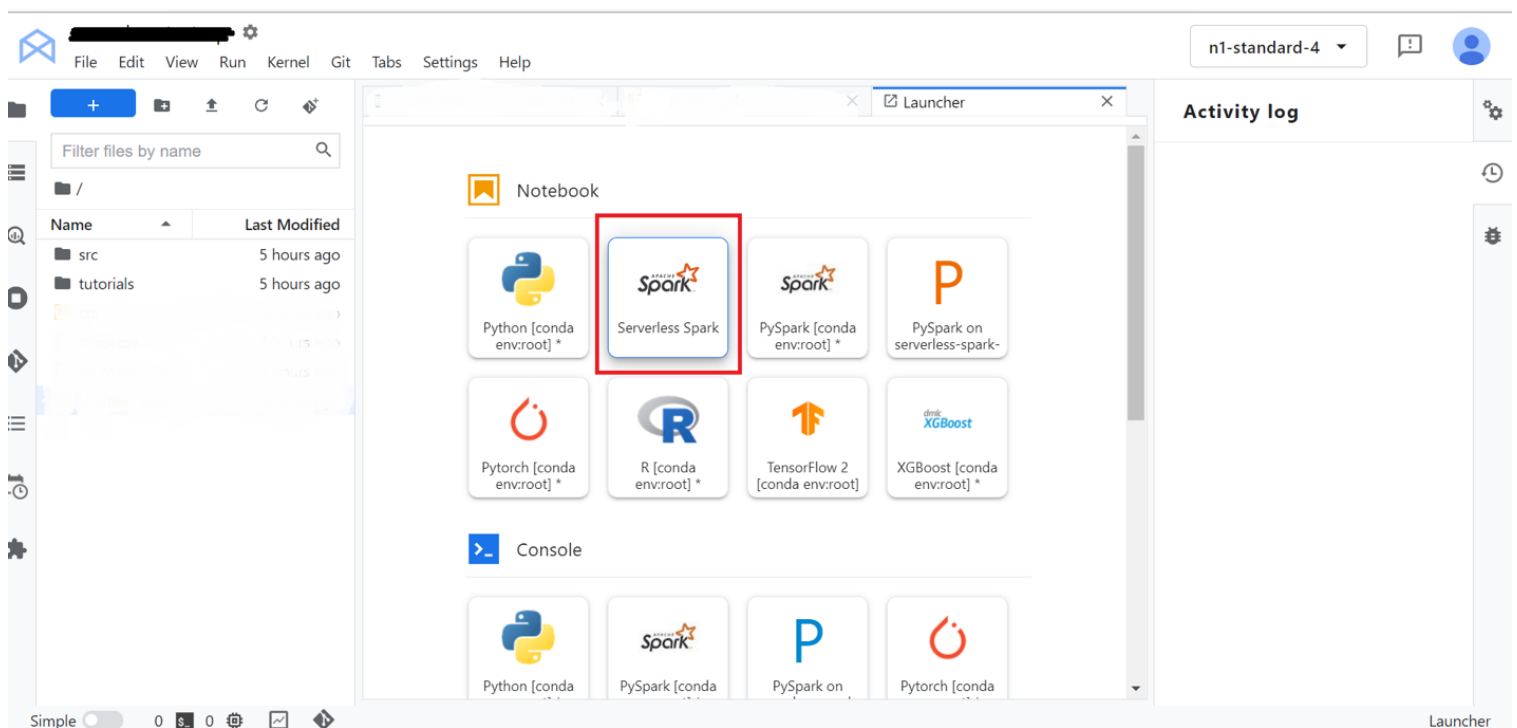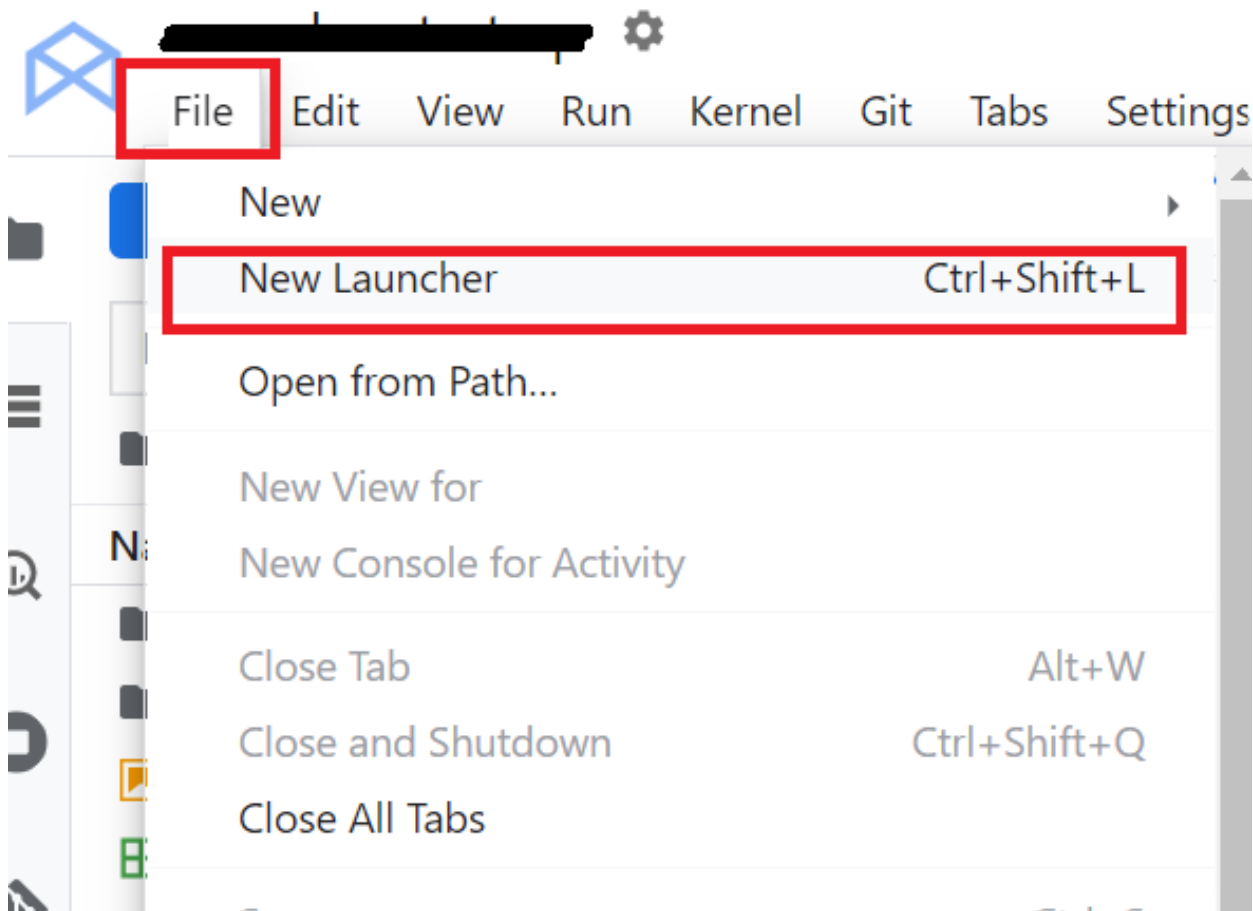- Provide a name and region to the notebook and select 'Single User Only' and click 'Create'. We will let the 'Advanced Settings' remain as the default values.

- Once the notebook is running, click the 'OPEN JUPYTERLAB' option next to the Notebook name as shown below



- Follow the on screen instructions to launch the JupyterLab session

## 🔗 Create Serverless Spark Session

- Click on the File and the New launcher and select Serverless Spark

# ⌘Follow the on screen instructions to create Session

## ⌘3.2. Provide the details for the Session

Next, fill in the following values in the session creation window as shown in the images below:

- **Session Name** - A unique identifier for your session
- **Region** - The region name provided by the Admin team
- **Language** - Pyspark
- **Autoshutdown** - 24 hours
- **Service Account** - <UMSA_NAME>@<PROJECT_ID>.iam.gserviceaccount.com
- **Network Configuration** - Select the network and subnetwork provided by the Admin team
- **History Server Cluster** - projects/<PROJECT_ID>/regions/<REGION_NAME>/clusters/<HISTORY_SERVER_NAME>
- **Properties** - spark.jars=gs://spark-lib/bigquery/spark-bigquery-with-dependencies_2.12-0.22.2.jar

- Click the **SUBMIT** button to create the session.

# Create Serverless Spark Session  PREVIEW

## Basic info

Session name *

Up to 128 lowercase letters, numbers, or underscores.

Language

PySpark ▼

Region

us-west1

Autoshutdown

24h

The session will automatically shutdown after 24 hours.

## Execution configuration

Service Account

Enter your service account

If not provided, the default GCE service account will be used. Learn More ⬈

## Network configuration

Private IP Google Access must be enabled on the network.

🔘 Networks in this project

⚪ Networks shared from host project: "undefined"

| Network | Subnetwork |
|---|---|
| serverless-spark-west ▾ | ss-w-subnet ▾ |

∧ ADVANCED OPTIONS

**History server cluster**

Choose a history server cluster to store logs in.

projects/█████████████████████/regions/us-central1/clusters/sp

**Properties**

Input parameters (optional)

spark.jars=gs://spark-lib/bigquery/spark-bigquery-with-
dependencies_2.12-0.22.2.jar

Each parameter needs to be separated by commas (Example:a=x,b=y)

**Labels**

Input parameters (optional)

Enter your comma separated parameters (Example:a=x,b=y)
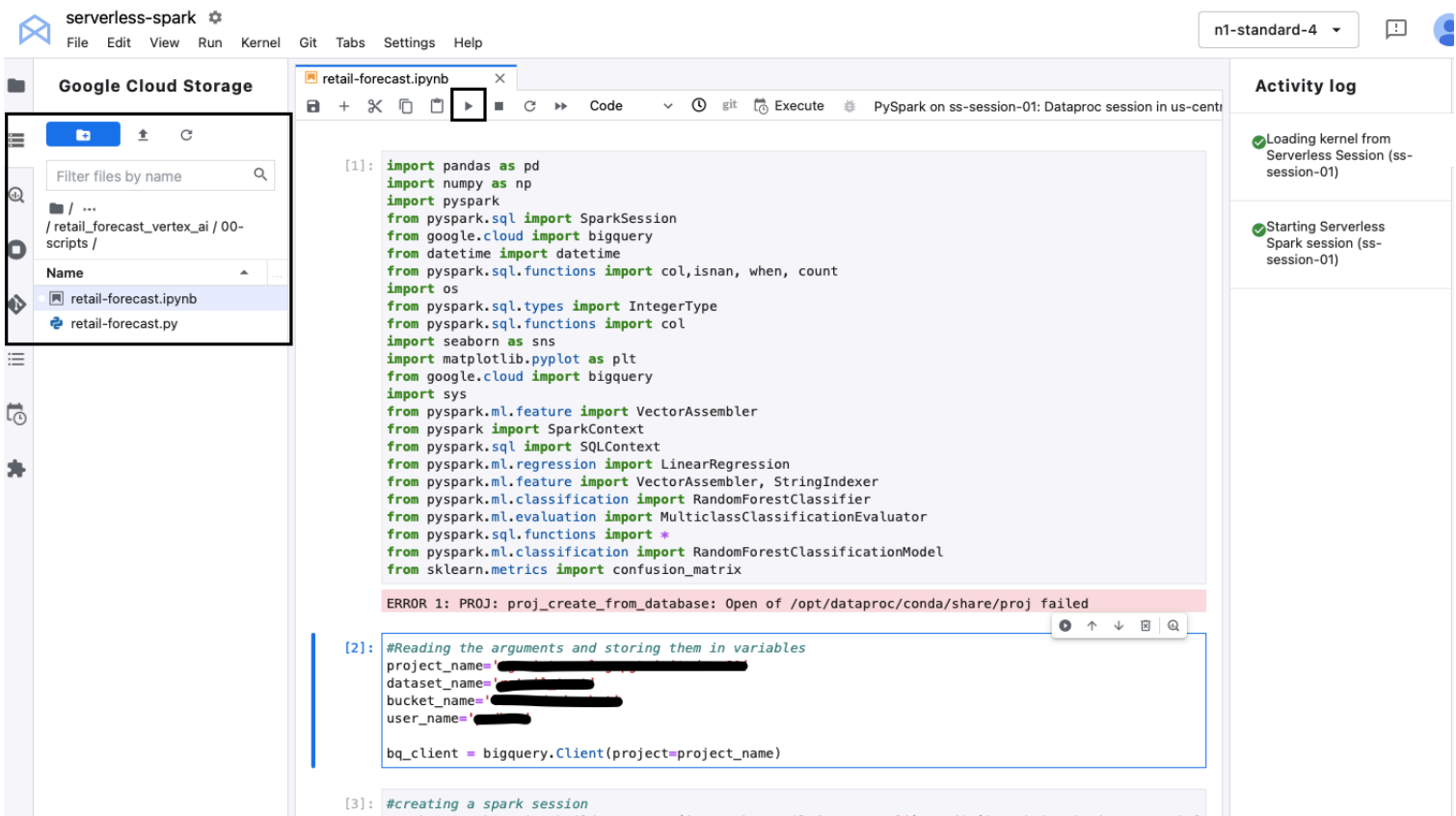
Each parameter needs to be separated by commas (Example:a=x,b=y)

**SUBMIT** CANCEL

- Once the Session is created select 'No Kernel' from the kernel dropdown list and then delete the notebook

- Next, using the browser option from JupyterLab, navigate to the Notebook file located at:
  <bucket_name> > 'retail_forecast_vertex_ai' > 00-scripts > retail-forecast.ipynb
- From the kernel dropdown list, select the kernel for the session created in section 3.2
- Pass the values to the variables project_name, dataset_name, bucket_name as provided by the Admin
  and replace user_name by your username
- Next, hit the **Execute** button as shown below to run the code in the notebook.



## 3.3. Check the output table in BQ

Navigate to BigQuery Console, and check the **retail_forecast** dataset.
Once the data preparation batch is completed, four new tables '<your_name_here>_train_data', '<your_name_here>_test_data', '<your_name_here>_predictions_data' and '<your_name_here>_eval_output' will be created:

To view the data in these tables -

- Select the table from BigQuery Explorer by navigating 'project_id' **>** 'dataset' **>** 'table_name'
- Click on the **Preview** button to see the data in the table

**Note:** If the **Preview** button is not visible, run the below queries to view the data. However, these queries will be charged for the full table scan.

```
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_train_data` LIMIT 1000
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_test_data` LIMIT 1000
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_predictions_data` LIMIT 1000
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_eval_output` LIMIT 1000
```

**Note:** Edit all occurrences of <project_name> and <dataset_name> to match the values of the variables PROJECT_ID, and BQ_DATASET_NAME respectively



# ∞4. Logging

## ∞4.1 Persistent History Server logs

To view the Persistent History server logs, click the 'View History Server' button on the Sessions monitoring page and the logs will be shown as below:

As the session is still in active state , we will be able to find the logs in show incomplete applications.