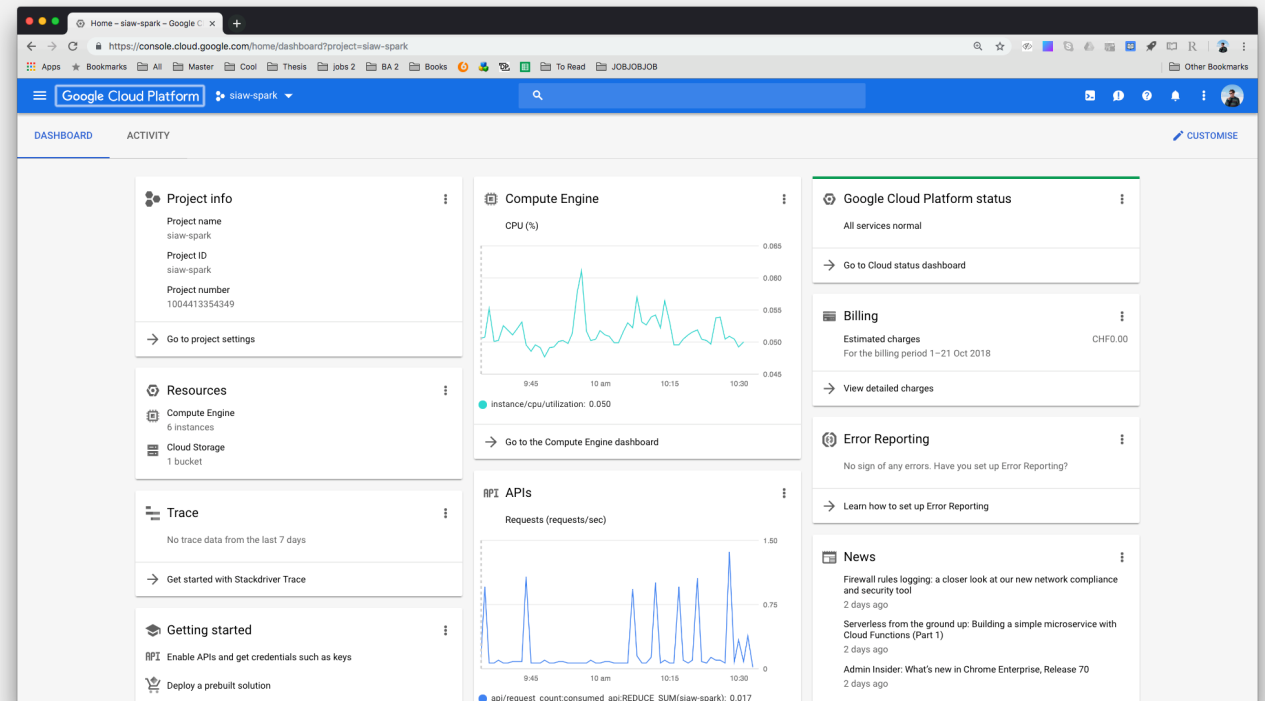# RUNNING A PYSPARK CLUSTER ON A SERVER

This tutorial will show you how to set-up a server cluster with Apache Spark, how to load data into it, and interact with your cluster using a Jupyter Notebook PySpark interface
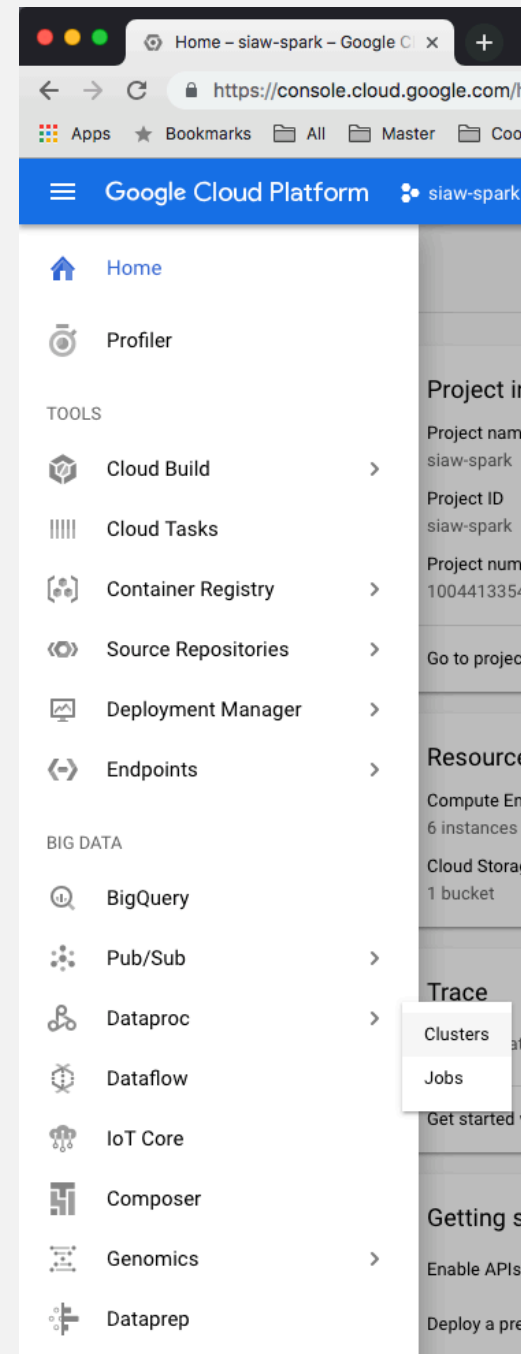
# STEP1: OPEN GOOGLE CLOUD PLATFORM

- If you already have an account simply log-in

- If you do not have a Google Cloud Platform (GCP) account go to https://cloud.google.com/ and create a free account

- If this is your first account, you will receive 300 USD of free platform credit so you should not be incurring any costs for the sake of this tutorial

- As soon as you have your account set up create a new project. The name used for this example will be called "siaw-spark"

- You should now be able to see your dashboard

# 2. CLUSTER CREATION

- Go to the menu in the top left corner and look for the item "Dataproc" under the subcategory "Big Data"

- Click on "Clusters"

# 2. CLUSTER CREATION

- In the new menu you will see no clusters running so far.

- If prompted, first enable the Cloud Dataproc API by following the instructions provided

- Next, click on "create cluster"

# 2. CLUSTER CREATION

- Assign a name to your cluster. For this installation example we will be using "pyspark-cluster"

- Change the region to a location close to you, I will choose Europe West

- You can change many other preferences here but for now we will ignore most of them

- Open "advanced options" at the bottom and click on "add initialisation action"

- Insert the link as explained on the right. The link directs to a script which allows us to use the Jupyter interface later on

- Click: "create"*

* Depending on which size you choose your cluster to be, you might have to enable billing first. Enabling Billing, however, does not mean that you lose your credit. You will only be charged should you surpass your available credit

*Insert this link here:*

*gs://dataproc-initialization-actions/jupyter/jupyter.sh*

# 3. SETTING UP THE JUPYTER INTERFACE

- Our cluster will a few minutes to be launched. Its status will change to "Running" when ready

- Click on the menu in the top left again and click on "Compute Engine"

- Click on the three vertical dots to the right of your cluster with the suffix "-m" (this is the master node)

- Click on "View Network Details"

# 3. SETTING UP THE JUPYTER INTERFACE

- From the left, select "Firewall Rules"

- Select "Create Firewall Rule"

- Use the following settings:
  - Name: jupyter
  - Target tags: http-server
  - Source IP ranges: your v4 IP
  - tcp: 8123 (the script sets up Jupyter on this port)

- Once created, you can use this rule for all future spark clusters

- Click "Create"

# 3. SETTING UP THE JUPYTER INTERFACE

- Go back to the overview of your compute engines and click on your master node

- Click on edit

- Tick the box stating: "Allow HTTP traffic"

- Click "save"

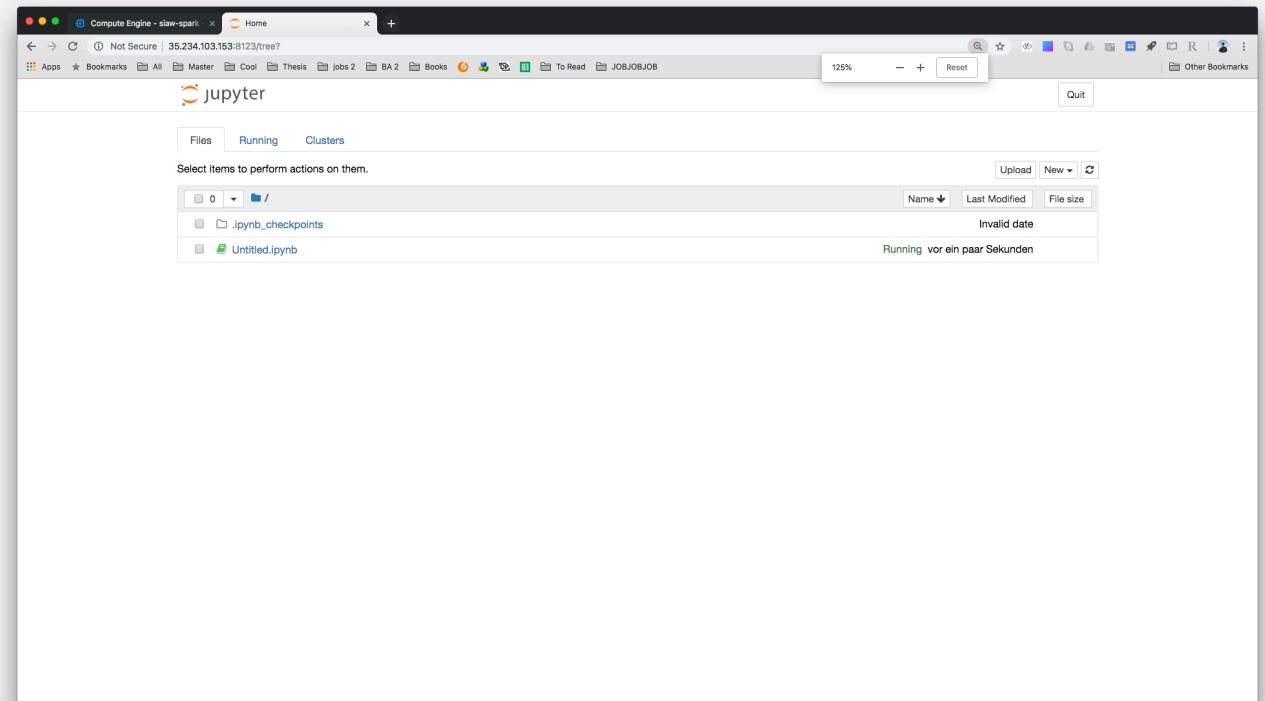# 3. SETTING UP THE JUPYTER INTERFACE

- Go back to the overview of your master node and copy your master nodes IP address into your browser

- add ":8123"

- This is port where the jupyter notebook is available
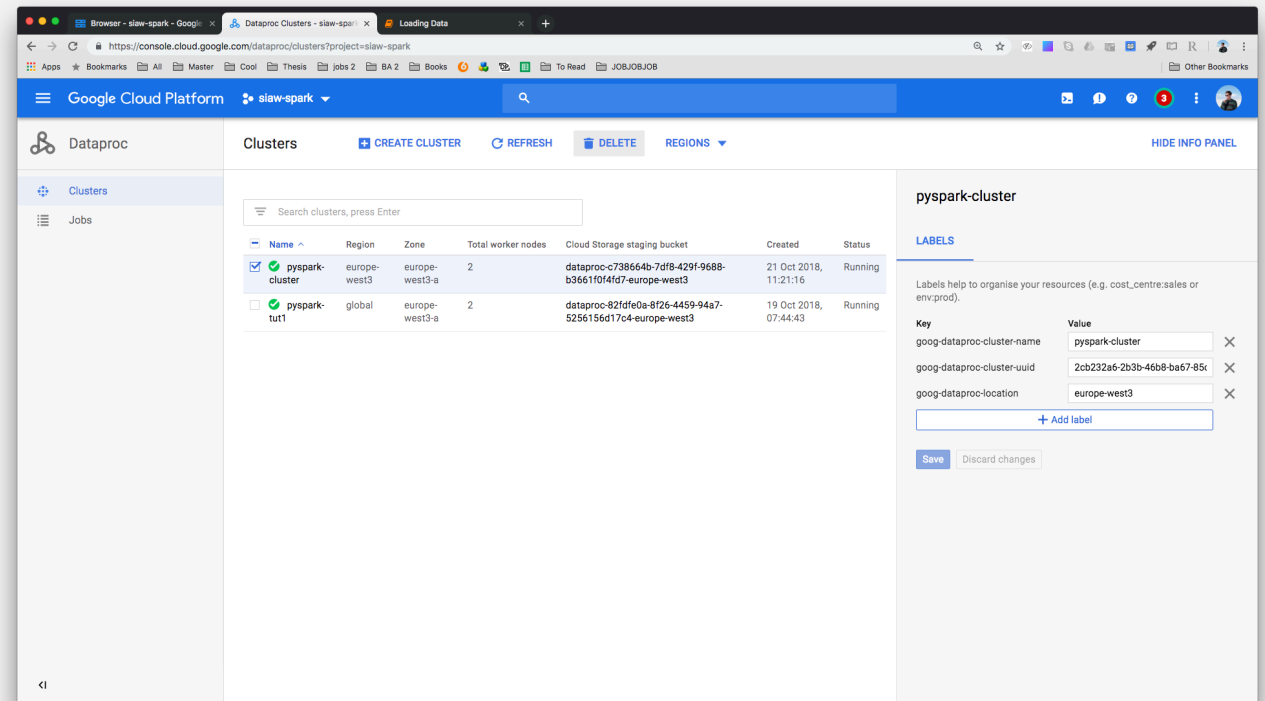
- Hit enter

# 4. DONE

- You can now use a regular Jupyter Notebook interface for working with PySpark

- Note: even though not covered specifically your your cluster comes with a storage bucket upon instantiation. You can upload data into this bucket and load it in your notebook

# 4. DELETION



- Don't remember to use your cluster when finished to avoid incurring high costs despite not using it.

- Open the menu and click on the Dataproc item

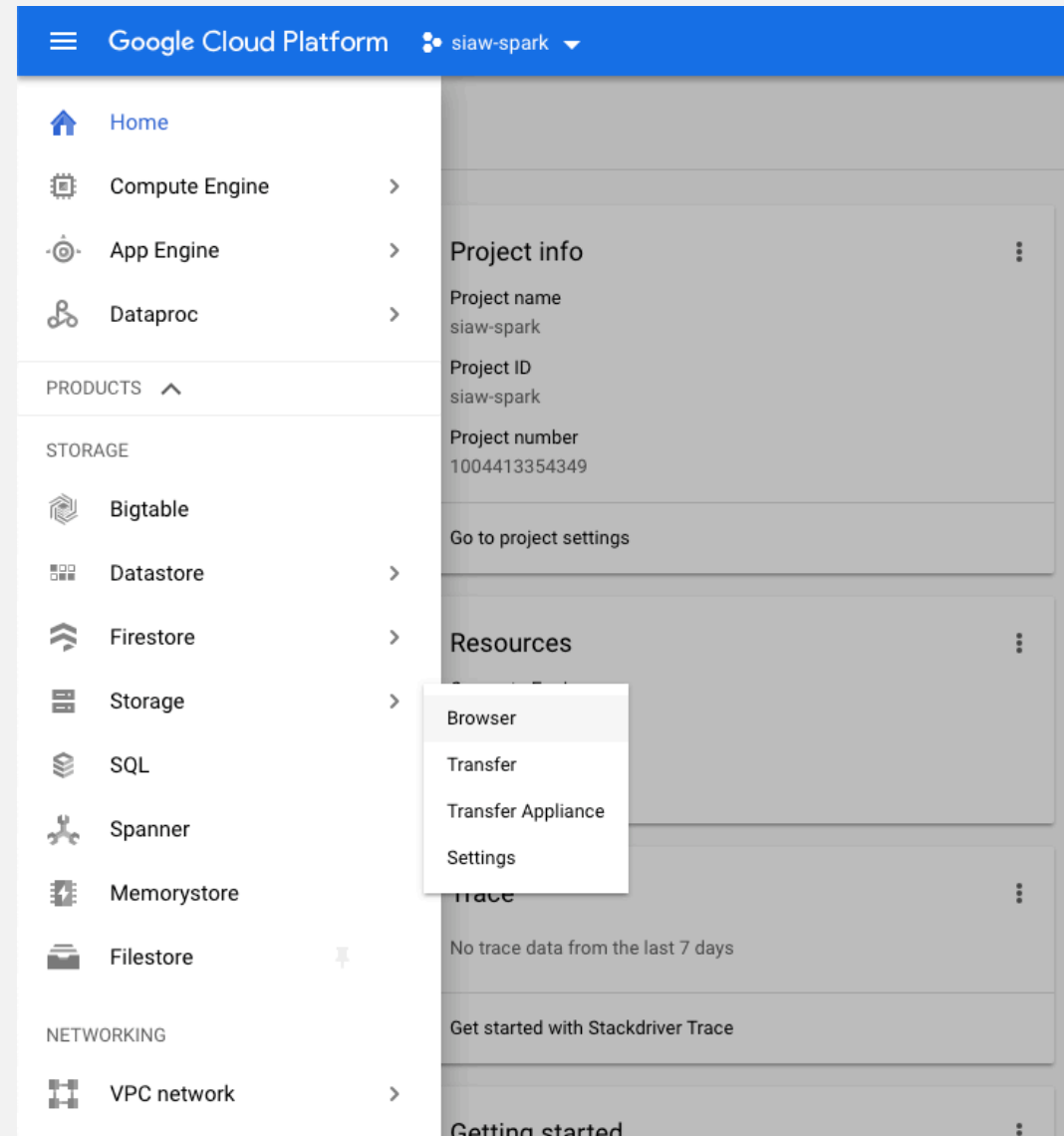- Select your cluster on the left and press "Delete" in the top menu

- Confirm

# UPLOADING DATA TO GCP

This section will explain how to upload data into a Google Cloud Platform storage bucket and read it from a PySpark script

# 1. UPLOADING DATA

- Open your GCP console, navigate to the hamburger menu and click on "Storage"
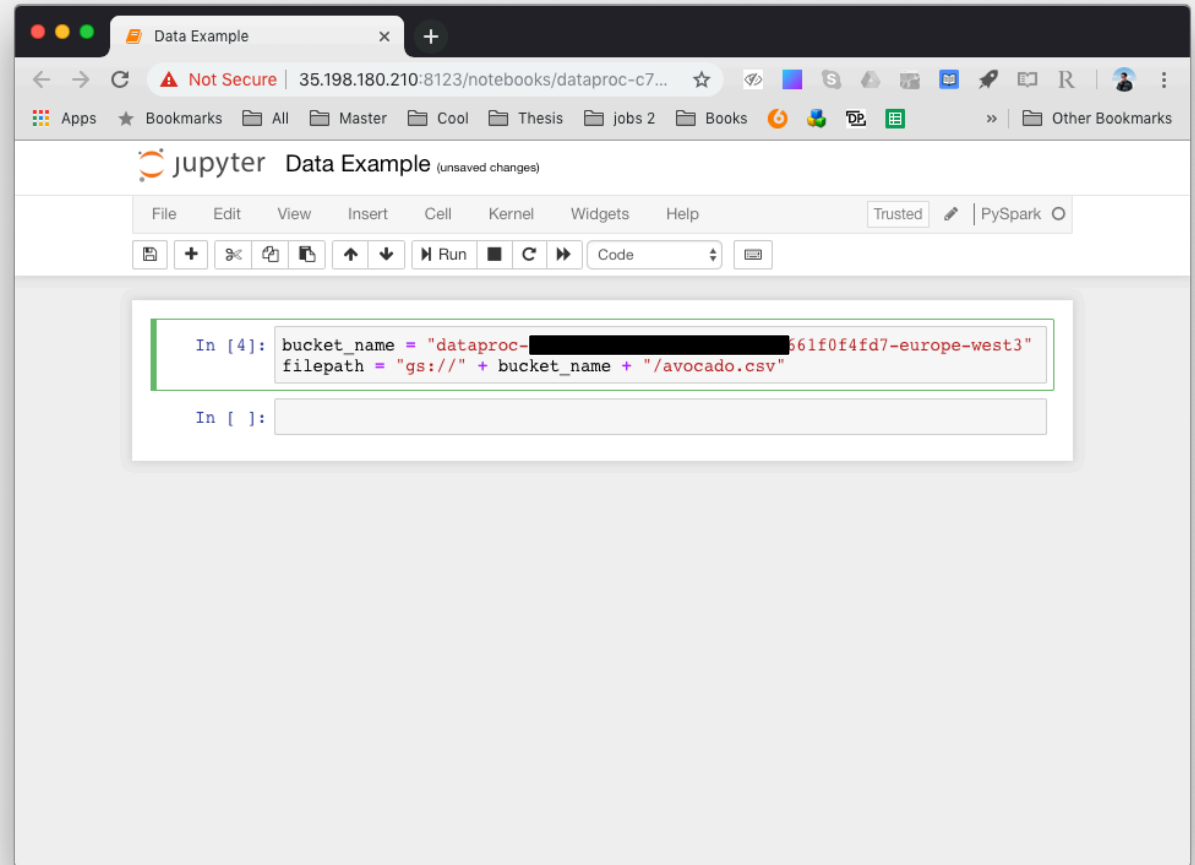
# 1. UPLOADING DATA

- You should be seeing a storage bucket that was automatically created when you set up your Spark Cluster. Click on it.

- As of now, only metafiles and a folder for Jupyter notebooks should be living in your notebook.

- Click on "Upload Files" in the top menu and select your file(s) to upload. In our case we will use a csv containing information on Avocado sales in the US
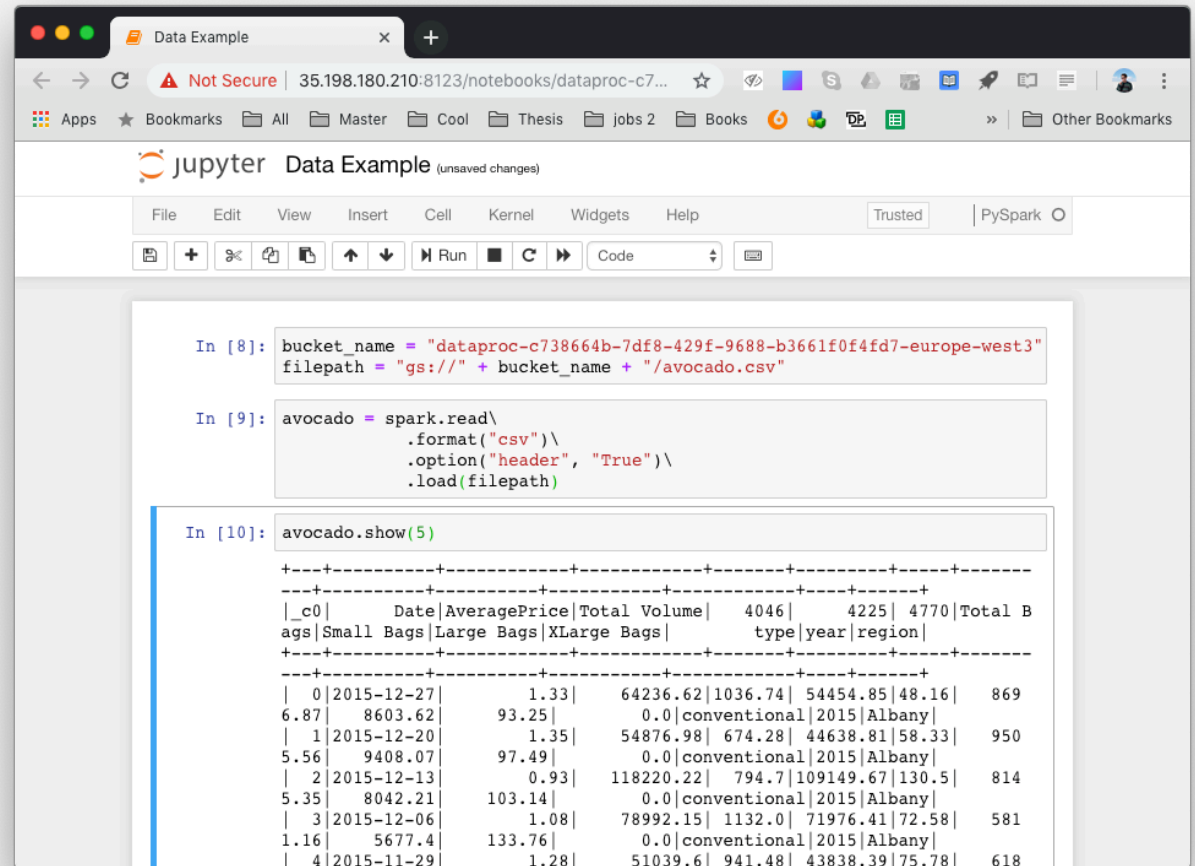
# 2. CALLING THE DATA FROM A NOTEBOOK

- Open a new PySpark Jupyter notebook

- You can access your data under the server file path which consists of the prefix "gs://", your server name and the path of the file in the bucket

- You can copy your server name from the bucket menu in the "Uploading Data" section

# 2. CALLING THE DATA FROM A NOTEBOOK

- You can now access the files at the specified file path with the PySpark commands and use it as intended.