

Start Lab 02:00:00

MapReduce in Dataflow (Python)

2 hours Free ★★★★★ Rate Lab

- Overview
- Introduction
- Setup
- Identify Map and Reduce operations
- Execute the pipeline
- Use command line parameters
- What you learned
- End your lab

Duration is 1 min

In this lab, you learn how to use pipeline options and carry out Map and Reduce operations in Dataflow.

What you need

You must have completed Lab 0 and have the following:

- Logged into GCP Console with your Qwiklabs generated account

What you learn

In this lab, you learn how to:

- Use pipeline options in Dataflow
- Carry out mapping transformations

Introduction

Duration is 1 min

The goal of this lab is to learn how to write MapReduce operations using Dataflow.

Setup

For each lab, you get a new Google Cloud project and set of resources for a fixed time block.

2. Note the lab's access time (for example, 02:00:00) and make sure you can finish in that time block.

There is no pause feature. You can restart if needed, but you have to start at the beginning.

3. When ready, click **START LAB**.

4. Note your lab credentials. You will use them to sign in to the Google Cloud Console.

Open Google Console

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked. [Learn more](#).

Username

Password

TG959yrKDX

GCP Project ID

qwiklabs-gcp-0855e773352d3560

New to labs? View our introductory video!

5. Click **Open Google Console**.

6. Click **Use another account** and copy/paste credentials for **this lab** into the prompts.

If you use other credentials, you'll get errors or **incur charges**.

7. Accept the terms and skip the recovery resource page.

Do not click **End Lab** unless you are finished with the lab or want to restart it. This places your work and removes the account.

Activate Cloud Shell

Cloud Shell is a virtual machine that is loaded with development tools. It offers a persistent 5GB home directory and runs on the Google Cloud. Cloud Shell provides command-line access to your Google Cloud resources.

In the Cloud Console, in the top right toolbar, click the **Activate Cloud Shell** button.



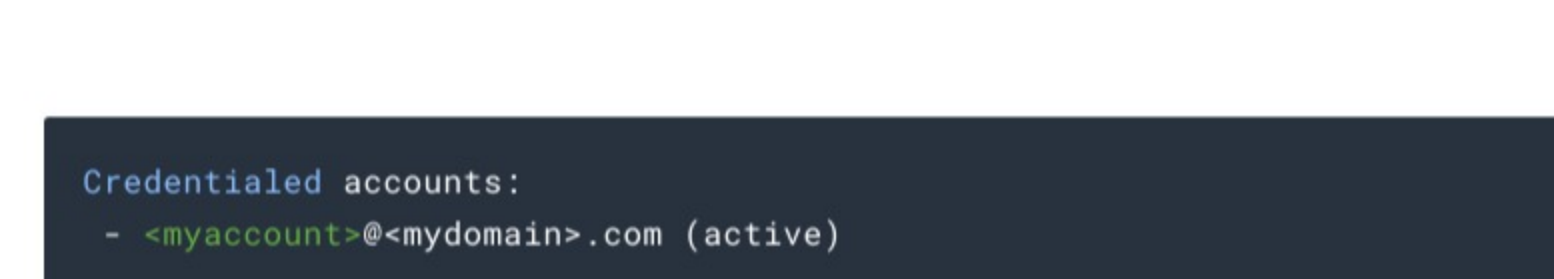
Click **Continue**.

Cloud Shell

Google Cloud Shell provides you with command-line access to your cloud resources directly from your browser. You can easily manage your compute and resource without leaving the browser.

Continue

It takes a few moments to provision and connect to the environment. When you are connected, you are already authenticated, and the project is set to your **PROJECT_ID**. For example:



gcloud is the command-line tool for Google Cloud. It comes pre-installed on Cloud Shell and supports tab-completion.

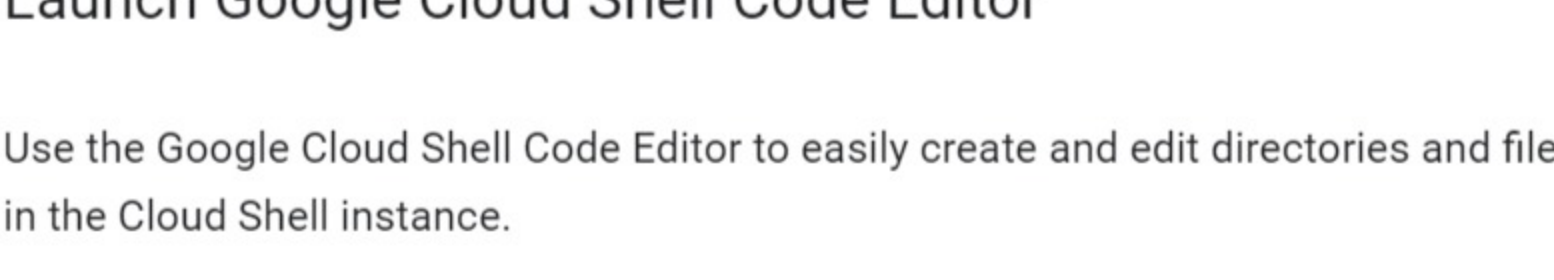
You can list the active account name with this command:



(Example output)



(Output)



(Example output)

For full documentation of gcloud see the [gcloud command-line tool overview](#).

Launch Google Cloud Shell Code Editor

Use the Google Cloud Shell Code Editor to easily create and edit directories and files in the Cloud Shell instance.

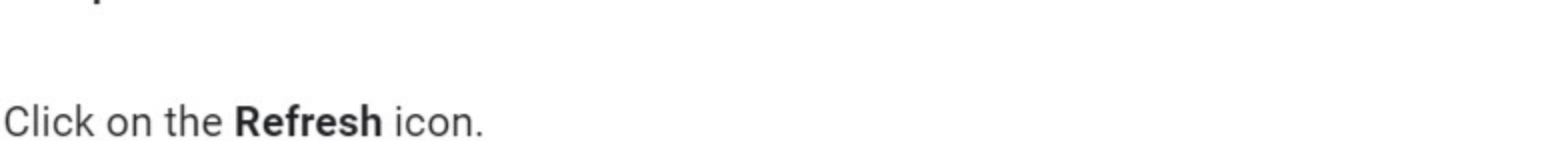
Once you activate the Google Shell, click the **Open editor** button to open the Cloud Shell Code Editor.



Duration is 5 min

Step 1

In CloudShell clone the source repo which has starter scripts for this lab:



Then navigate to the code for this lab.



Step 2

Click on the **Refresh** icon.

browser editor or with the command line using nano:



Step 3

What custom arguments are defined? _____

What is the default output prefix? _____

How is the variable output_prefix in main() set? _____

How are the pipeline arguments such as --runner set? _____

Step 4

What are the key steps in the pipeline?

Which of these steps happen in parallel? _____

Which of these steps are aggregations? _____

Execute the pipeline

Duration is 2 min

Step 1

Install the necessary dependencies for Python dataflow:



If not, open a new CloudShell tab and it should pick up the updated pip.

Step 2

Run the pipeline locally:



Note: If you see an error that says "No handlers could be found for logger: 'oauth2client.contrib.multistore_file'", you may ignore it. The error is simply saying that logging from the oauth2 library will go to stderr.

Step 3

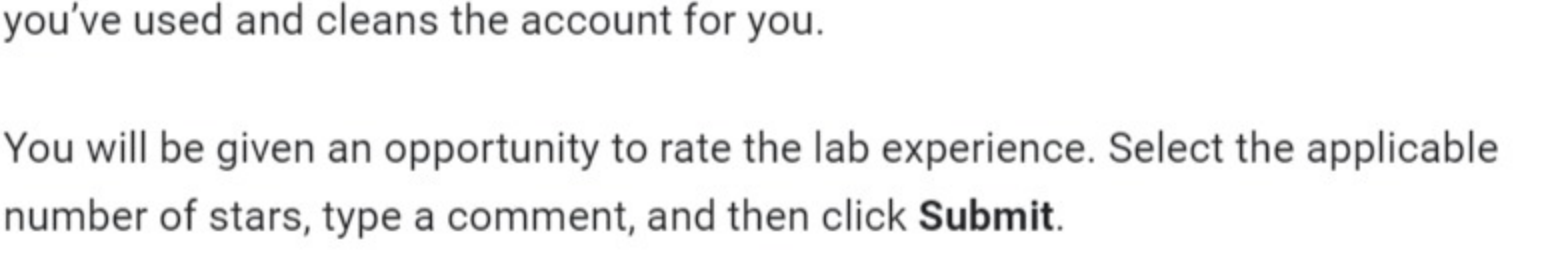


Use command line parameters

Duration is 2 min

Step 1

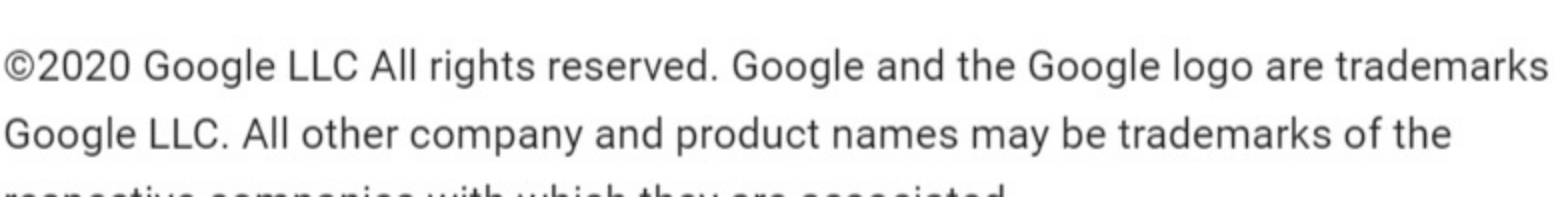
Change the output prefix from the default value:



What will be the name of the new file that is written out?

Step 2

Note that we now have a new file in the /tmp directory:



What you learned

Duration is 1 min

In this lab, you:

- Used pipeline options in Dataflow
- Identified Map and Reduce operations in the Dataflow pipeline

End your lab

When you have completed your lab, click **End Lab**. Qwiklabs removes the resources you've used and cleans the account for you.

You will be given an opportunity to rate the lab experience. Select the applicable number of stars, type a comment, and then click **Submit**.

The number of stars indicates the following:

- 1 star = Very dissatisfied
- 2 stars = Dissatisfied
- 3 stars = Neutral
- 4 stars = Satisfied
- 5 stars = Very satisfied

You can close the dialog box if you don't want to provide feedback.

For feedback, suggestions, or corrections, please use the **Support** lab.

©2020 Google LLC All rights reserved. Google and the Google logo are trademarks of Google LLC. All other company and product names may be trademarks of the respective owners.