# Big Data Graded Assignment 4
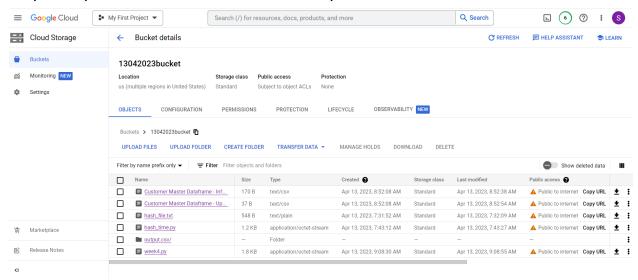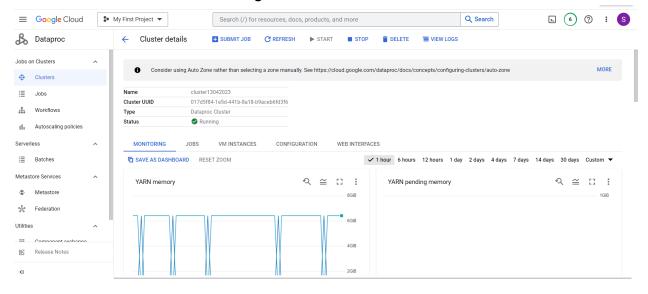
❖ Step 1 : Upload all files under bucket with public access
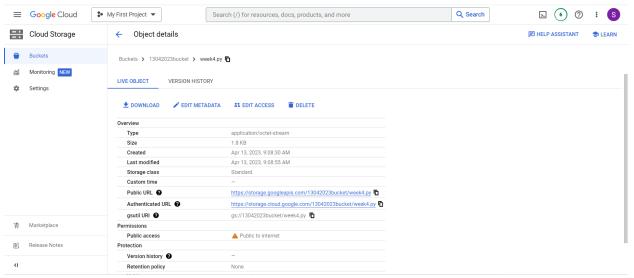


❖ Python file

## ❖ Create a cluster which is running



## ❖ Get gsutil url from file

## ❖ Job succeeded with correct output