**Kazakh-British Technical University**

**Cloud Computing**
**Assignment 4**

**Student: Khanfiyeva Elnara**
**Professor: Serek Azamat**

**2024**

**Table of Contents**

**Executive Summary**

In this report it is provided an in depth analysis and description of the implementing a Big Data and machine learning pipeline, also the security and compliance measures with the help of Google Cloud. This assignment included the writing of the report, also the technical part was to use data ingestion,preprocessing, ML model training, deployment and monitoring processes along with the configuration of secure infrastructure.

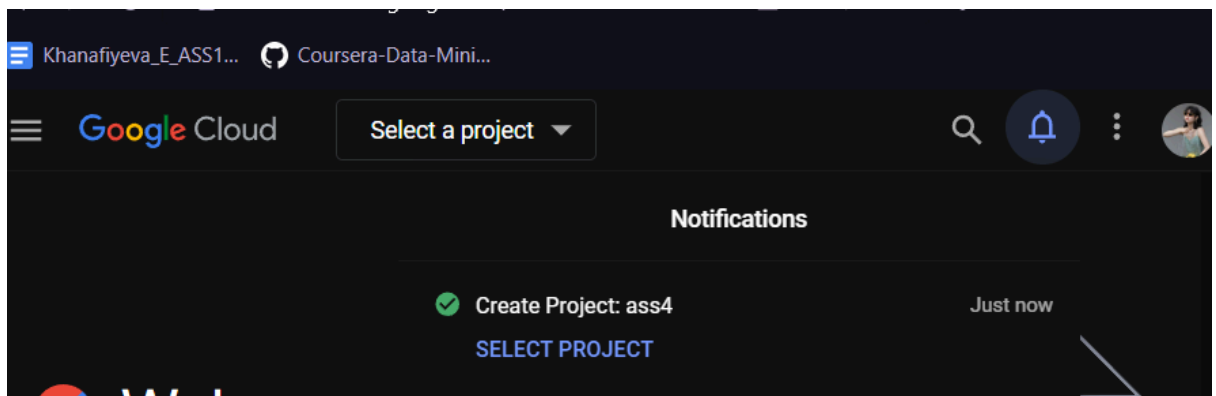I will discuss the process of doing the project in details.

**Introduction**

In modern computing, solving problems come down to (often complex) problems that big data and machine learning can help to solve. In a cloud environment, as important as computing data and building models is security and compliance. I built a pipeline for big data processing and machine learning and its related security and compliance practices on the Google Cloud. Through this report, I describe my method, considering the most characteristic findings, and the most convenient discovering.
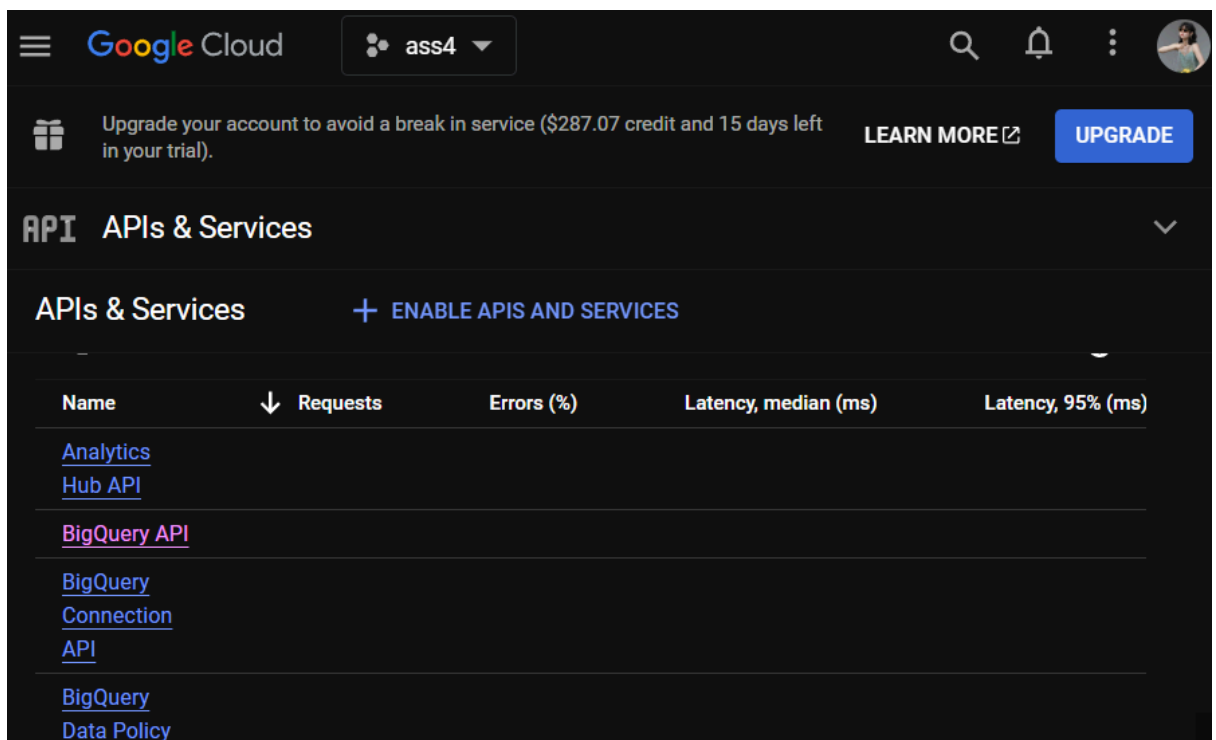
**Big Data and Machine Learning Pipeline**
**Exercise 1: Big Data and Machine Learning on Google Cloud**

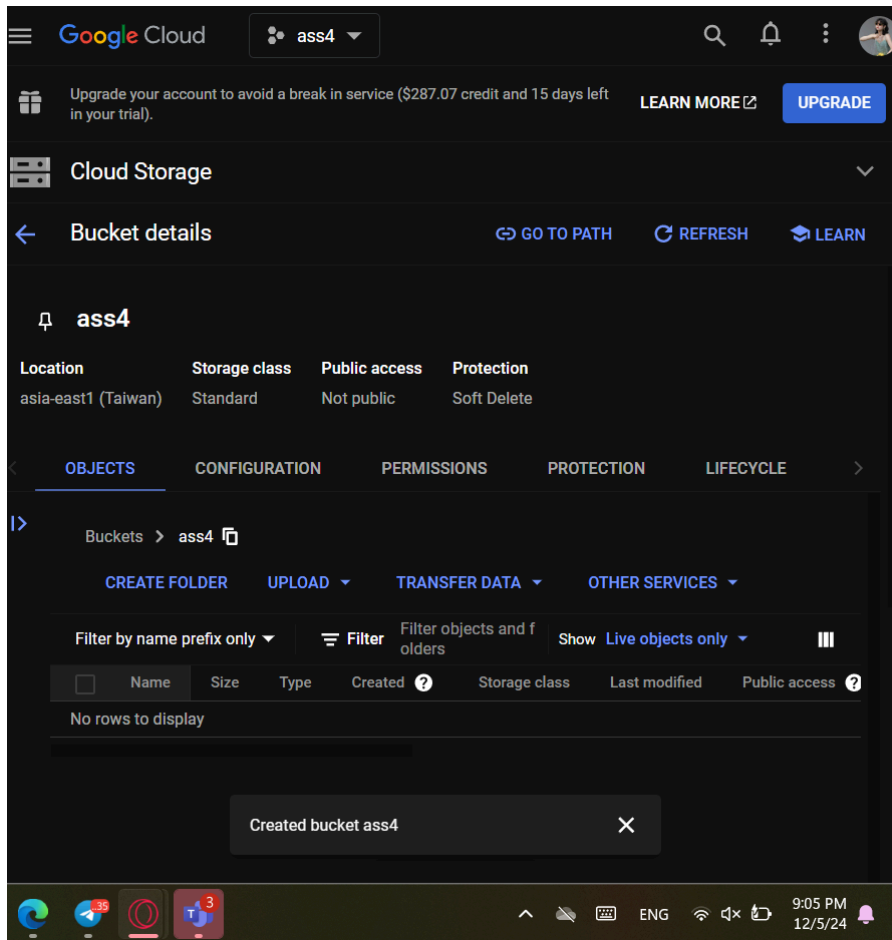First, I create the project in GCC:



Then I enable all of the required API's that are gonna be used in the assignment here:
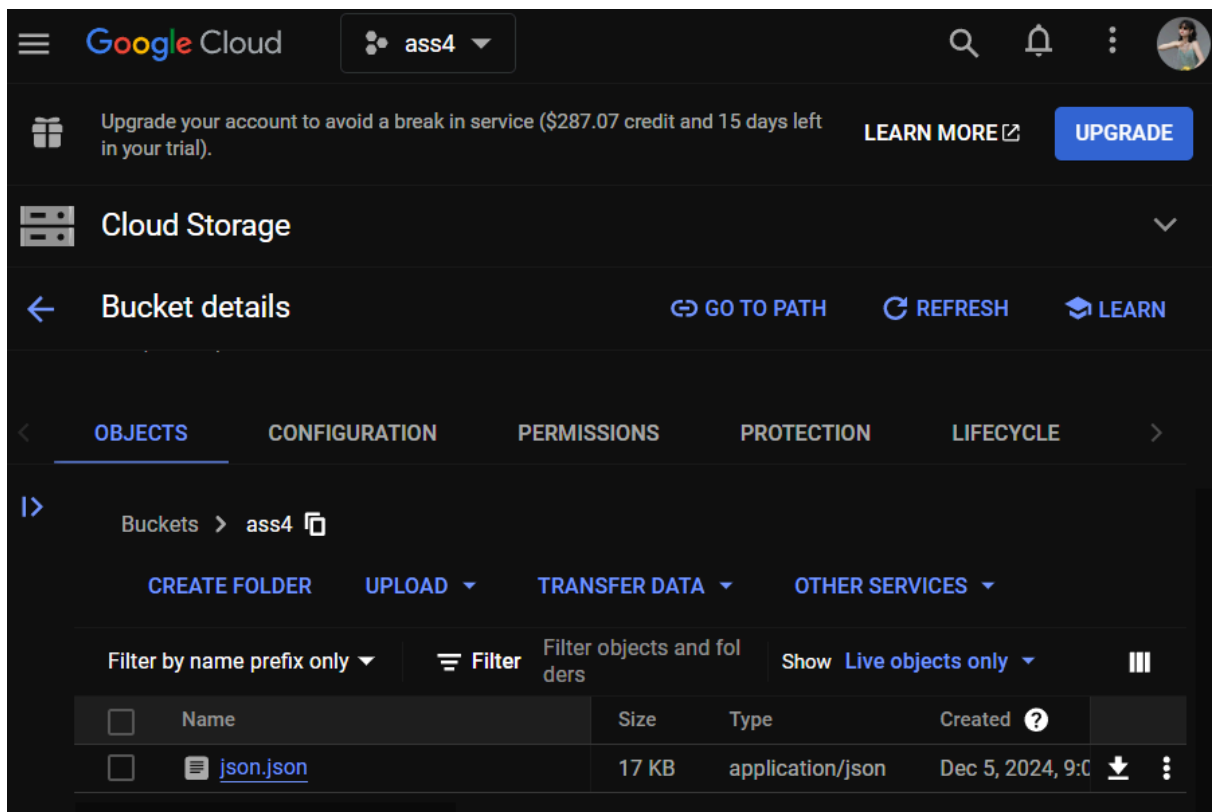


Next, let's start with data ingestion first. I will download datasets from somewhere else in internet or Google ready sets. I have found one here :State of the Cities Baseline Survey 2012-2013 - Kenya (catalog.ihsn.org datacatalog.ihsn.org).
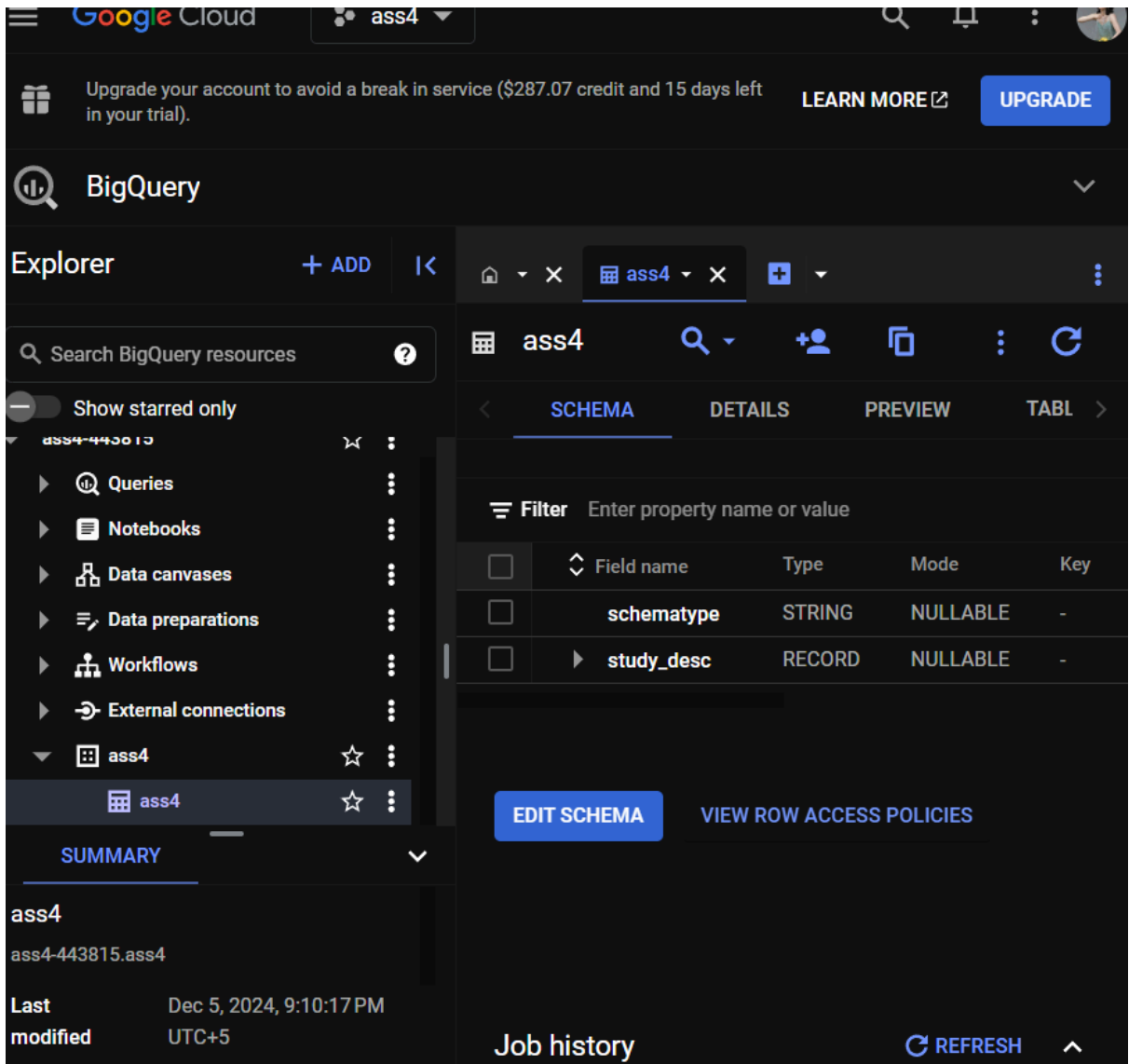
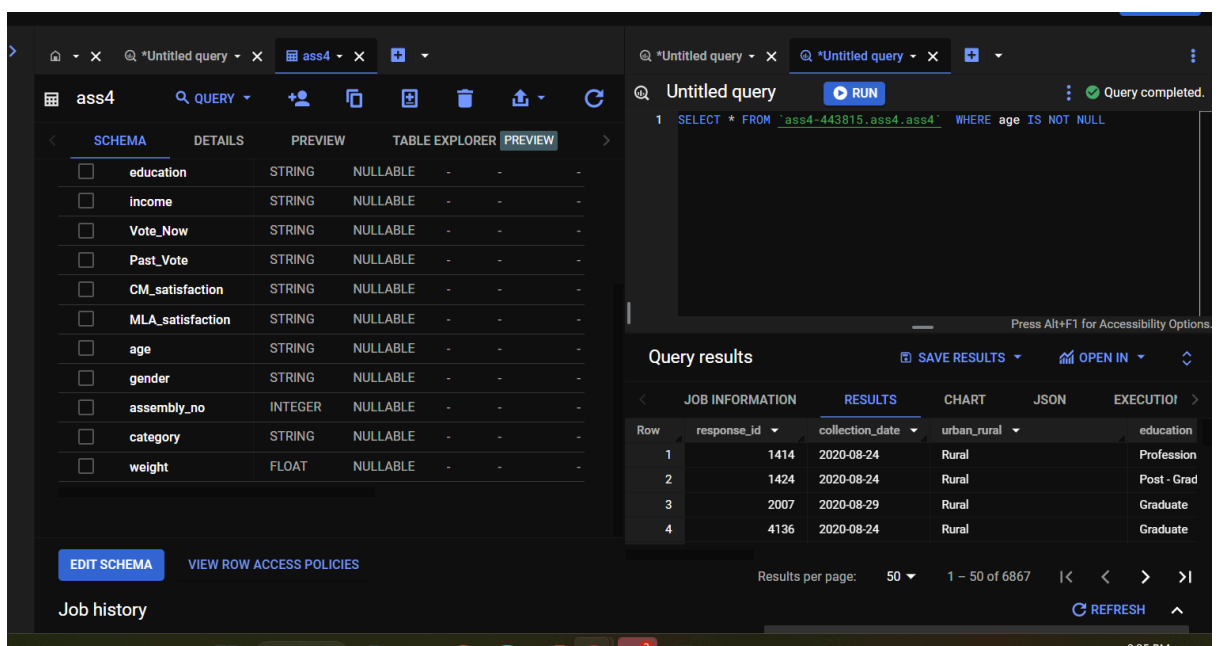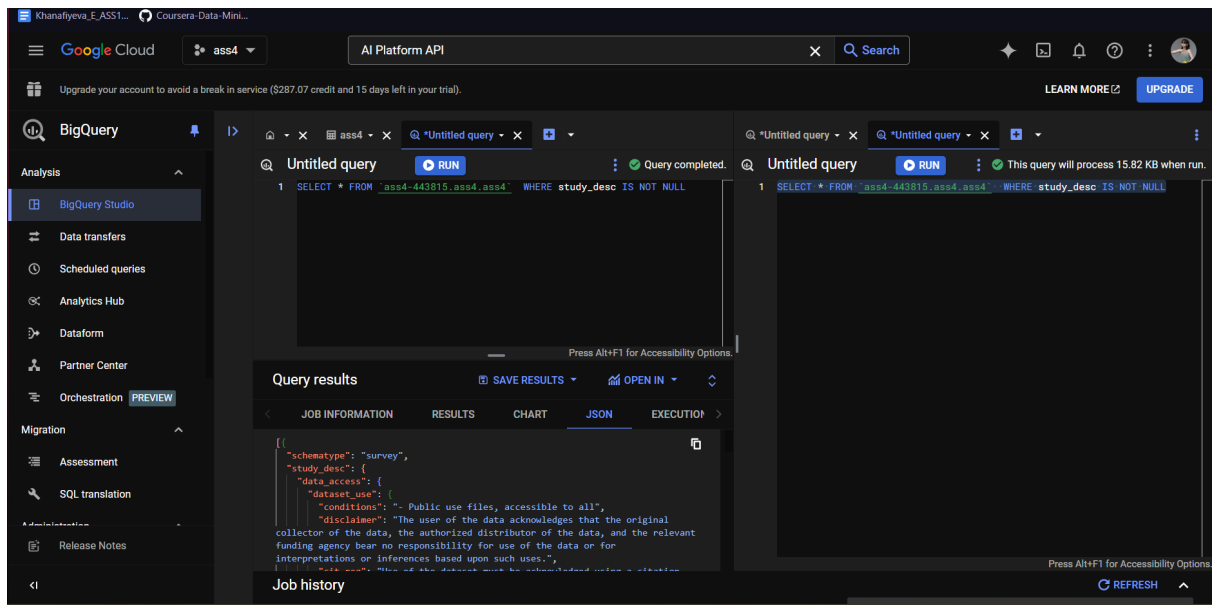Then I upload this dataset to my cloud storage bucket in GC:

Then upload this way:

Then i move to creation of the dataset in BigQuery console, to do so I create a dataset:



I create the dataset, then create the table from the cloud storage bucket file.
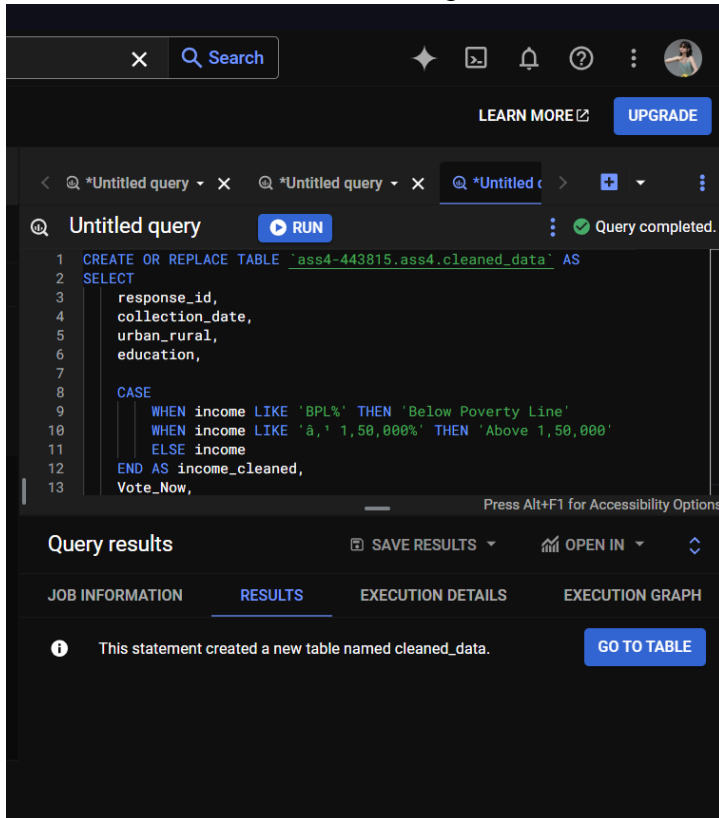
I have changed the dataset to the one that has more numeric data to easily visualize.

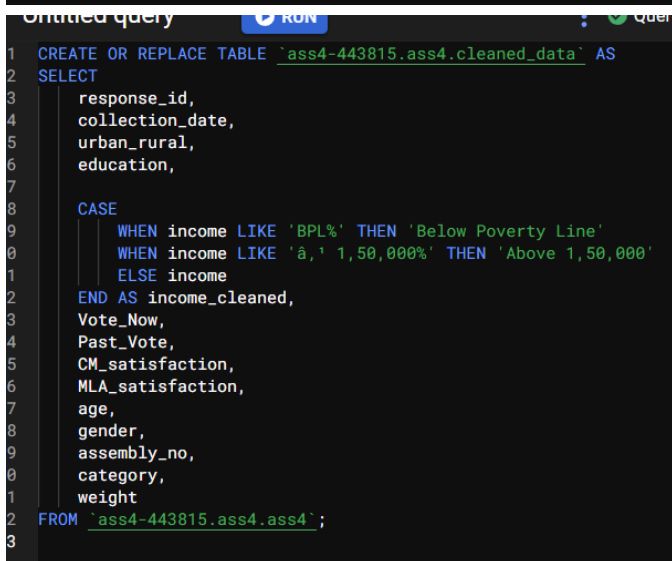https://www.kaggle.com/datasets/shashankshukla9919/sample-dataset/discussion?sort=hotness

I clean the table from corrupted information tot he new table here:





Then for example let's analyze the current voting preference (Vote_Now) by urban vs. rural respondents by this code below:

Here you can see the visualization of this sql response.

Voting Trends by Age Group here:

The results of the visualization of this result:



Then after cleaning the set:

So after we divide the set that will be used for training and the validation sets.

Since now I exported the processed final dataset to the initial bucket and then also wrote a code for training in zip and uploaded in bucket with requirements inside of it.

Then I go to the Vertex AI > Training in the console,



After configuring the data and mode training pipeline, here is the model training.

**Exercise 2: Cloud Security and Compliance**

Tasks:

1. Identity and Access Management (IAM):
   ○ Configure IAM roles and permissions for different users in your project.
   ○ Implement the principle of least privilege for service accounts and users.

In my project, I implemented Identity and Access Management (IAM) to effectively manage user roles and permissions, ensuring that access to resources within the cloud environment is both secure and streamlined. In my project, taking advantage of IAM serves to effectively manage user roles and permissions to provide easy access to resources as well as giving a sense of security with which resources can be accessed in the cloud environment. At first I defined a set of roles based on the needs of specific application, define different users and what responsibilities they have and what degree of access they need.



IAM Control setup.

IAM roles and policies were carefully defined with access controls, using the principle of least privileged in order to let only authorized ones access the sensitive information.

**14**

2. Data Encryption:
   - Ensure that data is encrypted at rest and in transit.
   - Utilize Google Cloud KMS for managing encryption keys.
   -

By default, Google Cloud uses encryption techniques like AES-256 to encrypt all data while it is at rest. This eliminates the need for further setting and guarantees that data is safe while at rest. However, we may also generate and maintain encryption keys using Google Cloud Key Management Service. Google Cloud employs TLS  for all connections, including network data transmissions and API calls, to guarantee the security of data while it is in transit. This guarantees that information is encrypted during transmission, shielding it from unwanted access.

3. Network Security:
   - Set up Virtual Private Cloud (VPC) and configure firewall rules to restrict inbound and outbound traffic.
   - Implement private Google access and ensure that sensitive data is not exposed to the public internet.

The creation and configuration of a Virtual Private Cloud (VPC) are critical steps in establishing a secure and efficient networking environment for an application.
Firstly. i determined the IR addresses for IPv4 as 0.0.0.0/0. The default, I also created the rule to allow http traffic on the 5000 port (tcp). Also the main part was enabling load balancers on the VPC settings.  By establishing a well-structured VPC environment and implementing stringent security protocols, organizations can effectively manage their cloud resources while safeguarding against potential threats.

*VPC Network settings.*

4. Audit Logging:
   ○ Enable Cloud Audit Logs to track access and changes to your resources.
   ○ Review logs for unusual activities and set up alerts for suspicious events.



5. Compliance Standards:

- Identify applicable compliance standards (e.g., GDPR, HIPAA) relevant to your project.
- Implement measures to ensure compliance, such as data residency, access controls, and audit trails.
- 

To fulfill compliance need, I briefly identified relevant standards such as GDPR and HIPAA, and one such global compliance standards, is ISO/IEC 27001.And we tend to use ISO/IEC 27001 which is also a known as ISMS (Information Security Management System) which is a widespread international standard on management. Best practices on how to manage offers protection of confidential information from information security risks. This is the standard that Google Cloud follows, using VPCs and firewall rules, strong controls on access, encryption and security on all data and services. I ensured that data storage and processed in line with limitations in terms of data residency. Additionally, access controls and audit trails were identified to support accountability and transparency.

6. Incident Response Planning:
   - Develop an incident response plan outlining the steps to take in case of a security breach.
   - Simulate a security incident and execute the response plan.

To cover potential security breaches I wrote out an extensive incident response playbook. The steps outlined in the playbook were for detection, containment, eradication and recovery. Then I simulated an incident in which unauthorized access was made to a resource. With the playbook steps I was able to detect and fix the issue and there was minimal impact. Secondly, this exercise also helped form the response plan by identifying what can be worked upon.

**Conclusion**

I have successfully designed and implemented a complete Big Data and Machine Learning pipeline in this project using Google Cloud, and covered important security and compliance steps. The project consisted of data ingestion process, data preprocessing, model training and deployment in the cloud environment where the biggest emphasis was made on data security and compliance with industry standards. It is a project that showcases the proper use of google cloud services like BigQuery for storing data, Vertex AI for model training, and a suite of security measures that ensure the data, but also the infrastructure, are kept safe. I also protected sensitive data by using encryption techniques and Identity and Access Management (IAM) ensuring that only authorized users could view sensitive data and that no sensitive was stored at rest or in transit. Further, by configuring Virtual Private Cloud (VPC) and setting very strict firewall rules, they protected the network environment against possible attacks. In addition, I also prioritized compliance with global standards like ISO/IEC 27001, GDPR, and HIPAA, to make sure that your data is dealt out with and processed according to legal and regulatory standards. Preparation for potential breaches made the incident response plan a further enhancement of the security posture in case of incidents. In conclusion, however, this project offered practical experience with using Google Cloud for Big Data and machine learning as well as a reminder of the importance of combining robust security and compliance. This exercise will help build secure, scalable, compliant cloud based solutions in real world scenarios and the skills gained from this are essential.

**References**

Google Cloud *Documentation*, Google Compute Engine. https://cloud.google.com/compute

Google Cloud IAM *Documentation*. https://cloud.google.com/iam