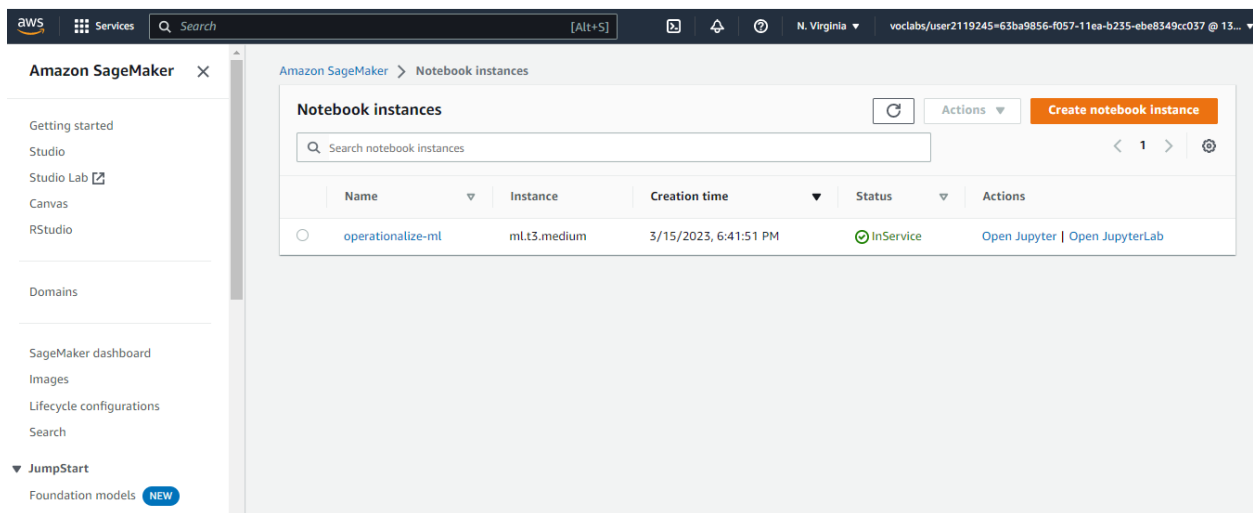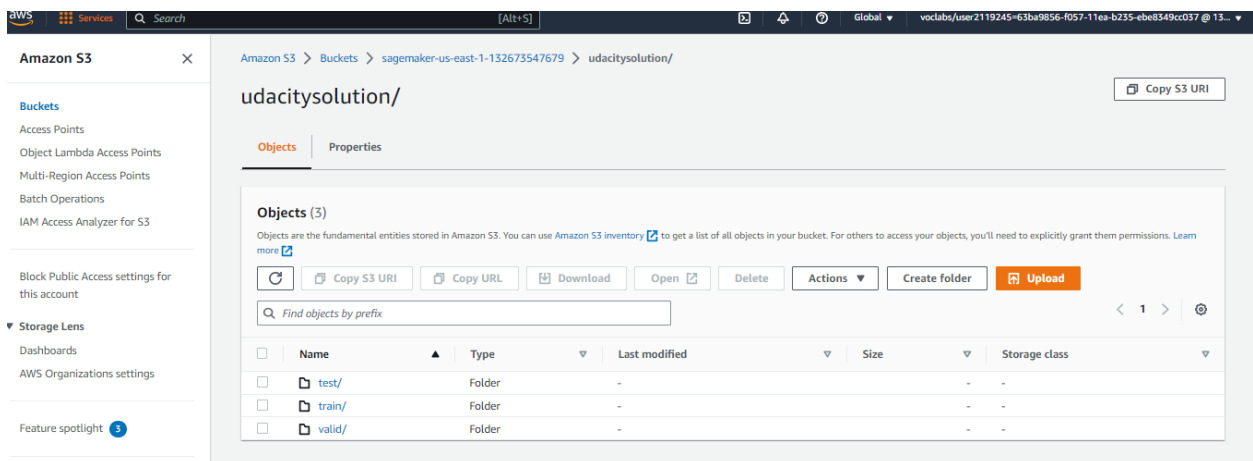# Operationalize AWS ML Project

## SageMaker Notebook instance:

I have chosen an ml.t3.medium instance because it has sufficient computing power (vCPU =2, Memory = 4 GiB) and it's cheap (Price per Hour = $0.05). So, it does the job while maintaining a low cost.
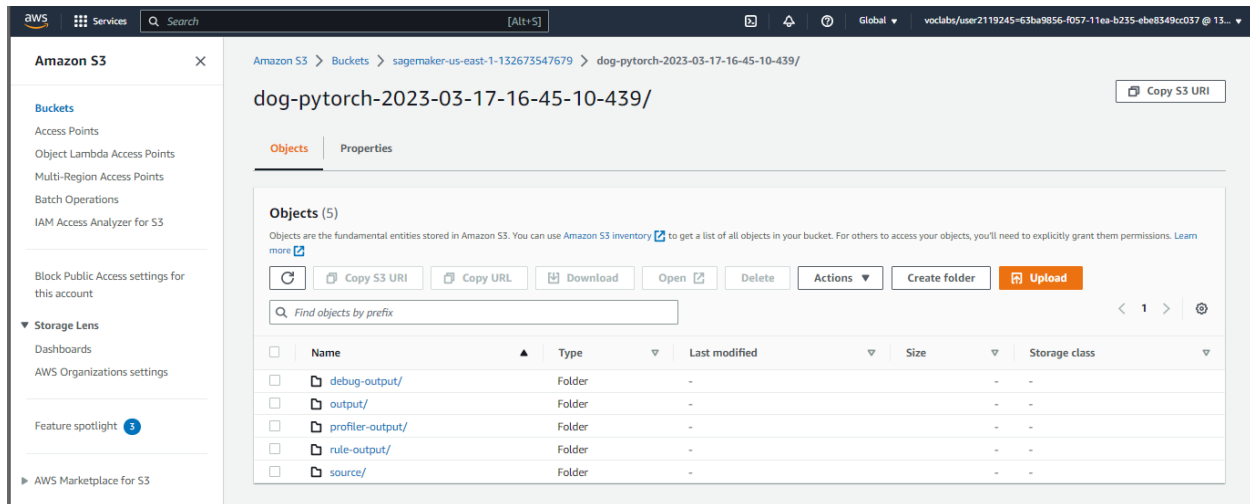


## S3 Bucket:

## The Training Dataset

# The Output of Training



# Training and Deployment:

# Hyperparameter Training Jobs

# For the training I started training using one instance and another time using multi-instance.

## - One instance

- **Multi instances**



**And here is the deployed end points (one using training with one instance and the other with multi-instance)**



**EC2:**

**I have used a t2.micro instance because it's eligible for the free tier and provide the needed** computation power.



**Difference between EC2 training code and SageMaker's:**

- **In EC2 code there is no calling for any Estimator or Tuner functions. The code in the EC2 script is responsible for saving the model to the local path. While in the SageMaker scripts this was handled internally by SageMaker where the model data was stored to a S3 location.**
- **In the EC2 code, all the variables already mentioned in the code itself.**
- **In the EC2 the training happens on the same while in the SageMaker the training job runs on a separate container than the one on which the SageMaker notebook is running.**

**Lambda function:**

**Lambda functions are used for invoking our deployed endpoints. And the end point used is the one of the mult-instance training (pytorch-inference-2023-03-17-18-02-16-600). Also we attached an Amazon SageMaker full Accesses policy to be able to interact with SageMaker successfully with no errors.**

**In the following images there are the role and result of the testing**

## Security:

**It's not a good idea to give full access permission because it may result in a security preach, but always choose the right policies and permission and remove them when the task is done. Below is the policies attached to my SageMaker execution role.**

## Concurrency and Auto Scaling:

I have configured a provisioned concurrency after publishing a version for lambda. Also I have configured Auto Scaling to cope with the traffic requests.

**JumpStart**
Foundation models  NEW
Computer vision models
Natural language processing models

▶ Governance

▶ Ground Truth

**Notebook**
Notebook instances
Git repositories

▶ Processing

**Training**
Algorithms
Training jobs
Hyperparameter tuning jobs

**Inference**
Compilation jobs
Marketplace model packages
Models
Endpoint configurations
Endpoints
Batch transform jobs
Shadow tests

▶ Edge Manager

▶ Augmented AI

▶ AWS Marketplace

Tutorials

No widget on this dashboard.

### Endpoint runtime settings

[ Update weights ] [ Update instance count ] [ Configure auto scaling ]

| | | Variant name ▲ | Current weight ▽ | Desired weight | Elastic Inference | Instance type ▽ | Current instance count ▽ | Desired instance count ▽ | Instance min - max | Automatic scaling |
|---|---|---|---|---|---|---|---|---|---|---|
| ○ | P | AllTraffic | 1 | 1 | - | ml.m5.large | 1 | 1 | 1 - 5 | Yes |

### Endpoint configuration settings

[ Change ] [ Clone ]

**Endpoint configuration**

| Name | ARN | Encryption key | Creation time |
|---|---|---|---|
| pytorch-inference-2023-03-17-18-02-16-600 | arn:aws:sagemaker:us-east-1:132673547679:endpoint-config/pytorch-inference-2023-03-17-18-02-16-600 | - | 3/17/2023, 8:02:16 PM |

**Data capture**

| Enable data capture | Data capture options | S3 location to store data collected | Capture content type |
|---|---|---|---|