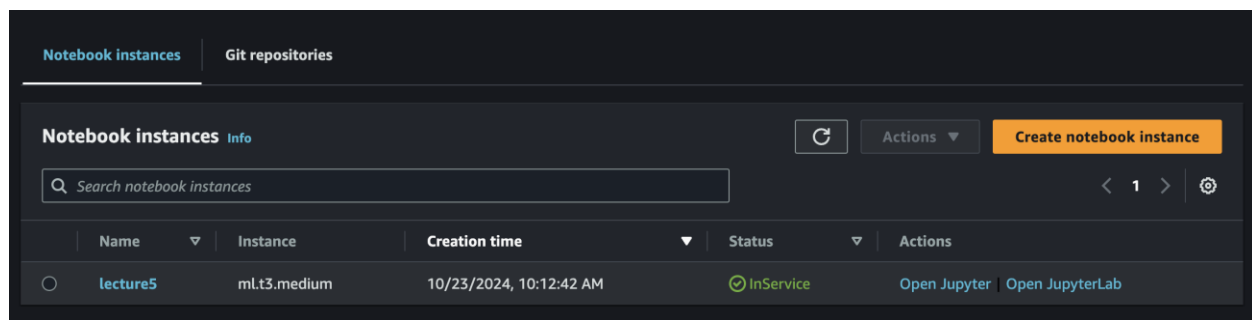# Operationalizing ML on Sagemaker
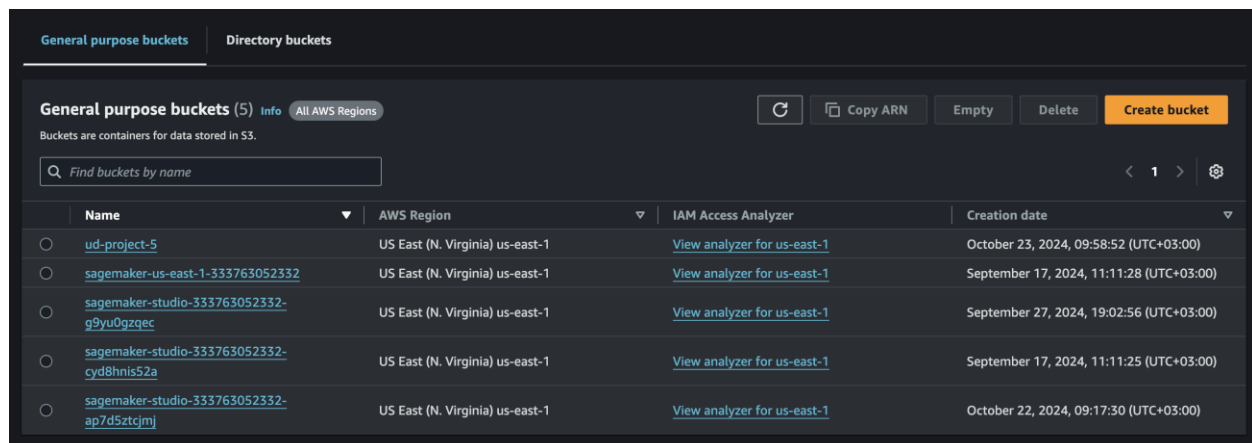
## Training and Deployment on Sagemaker

### Initial Setup

I selected the ml.t3.medium instance as it balances cost and performance for small to medium training tasks. It offers adequate CPU, memory, and network performance to efficiently train and deploy a computer vision model without high expenses. The instance also launches quickly, making it suitable for tasks requiring decent computational speed without heavy GPU use.



### S3 Setup

For storing and managing data, I set up an Amazon S3 bucket. This provides a scalable and secure location for the large datasets required for computer vision tasks. Using S3, I can easily access and manage the data from Sagemaker, facilitating seamless integration between storage and computation.



## Training and Deployment

# Hyperparameter Tuning

Even though 2 of my trainings failed at first, one of them was a success and I used it.

## Training job status counter

| Completed 1 | In Progress 0 | Stopped 0 | Failed 2 ( Retryable: 2, Non-retryable: 0 ) |
|---|---|---|---|

### Training jobs

Sorting by objective metric value will display only jobs that have metric values.

[ C ]  [ View logs ⧉ ]  [ View instance metrics ⧉ ]  [ Stop ]  [ Create model ]

[ Q Search training jobs ]                                                                                    < 1 >  ⚙

| | Name | Status ▽ | Final objective metric value ▽ | Creation time ▽ | Training Duration |
|---|---|---|---|---|---|
| ○ | pytorch-training-241023-0728-003-e25c7c30 | ⊗ Failed | - | 10/23/2024, 10:54:58 AM | 3 minute(s) |
| ○ | pytorch-training-241023-0728-002-13ba49e6 | ⊘ Completed | 150 | 10/23/2024, 10:33:56 AM | 18 minute(s) |
| ○ | pytorch-training-241023-0728-001-99fc3841 | ⊗ Failed | - | 10/23/2024, 10:28:33 AM | 4 minute(s) |

### Describe the tuning results

```
In [12]: exp = HyperparameterTuningJobAnalytics(
             hyperparameter_tuning_job_name='pytorch-training-241023-0728')

         jobs = exp.dataframe()

         jobs.sort_values('FinalObjectiveValue', ascending=0)
```

Out[12]:

| | batch_size | learning_rate | TrainingJobName | TrainingJobStatus | FinalObjectiveValue | TrainingStartTime | TrainingEndTime | TrainingElapsedTimeSeconds |
|---|---|---|---|---|---|---|---|---|
| 1 | "32" | 0.018242 | pytorch-training-241023-0728-002-13ba49e6 | Completed | 150.0 | 2024-10-23 07:34:02+00:00 | 2024-10-23 07:51:41+00:00 | 1059.0 |
| 0 | "512" | 0.050171 | pytorch-training-241023-0728-003-e25c7c30 | Failed | NaN | 2024-10-23 07:55:02+00:00 | 2024-10-23 07:57:33+00:00 | 151.0 |
| 2 | "512" | 0.001703 | pytorch-training-241023-0728-001-99fc3841 | Failed | NaN | 2024-10-23 07:29:31+00:00 | 2024-10-23 07:33:36+00:00 | 245.0 |

# Training

| Log streams | Tags | Anomaly detection | Metric filters | Subscription filters | Contributor Insights | Data protection |
|---|---|---|---|---|---|---|

### Log streams (23)

[ C ]  [ Delete ]  [ Create log stream ]  [ Search all log streams ]

[ Q dog-pytorch-2024-10-23-08-04-36-371                    × ]  4 matches  ☐ Exact match  ☐ Show expired ⓘ Info   < 1 >  ⚙

| ☐ | Log stream | Last event time ▽ | ▼ |
|---|---|---|---|
| ☐ | dog-pytorch-2024-10-23-08-04-36-371/algo-2-1729670728 | 2024-10-23 08:26:33 (UTC) | |
| ☐ | dog-pytorch-2024-10-23-08-04-36-371/algo-3-1729670728 | 2024-10-23 08:26:32 (UTC) | |
| ☐ | dog-pytorch-2024-10-23-08-04-36-371/algo-4-1729670729 | 2024-10-23 08:26:26 (UTC) | |
| ☐ | dog-pytorch-2024-10-23-08-04-36-371/algo-1-1729670728 | 2024-10-23 08:26:23 (UTC) | |

# Endpoint

# Multi Instance Training

I increased the instance count to 4 for my training. I added the code and logs.

**Creating an Estimator - Multi-Instance Training,**

```
In [27]: ###in this cell, create and fit an estimator using multi-instance training
         estimator = PyTorch(
             entry_point='hpo.py',
             base_job_name='dog-pytorch',
             role=role,
             instance_count=4,   # Change this to 2 or more for multi-instance training
             instance_type='ml.m5.xlarge',
             framework_version='1.4.0',
             py_version='py3',
             hyperparameters=hyperparameters,
             ## Debugger and Profiler parameters
             rules=rules,
             debugger_hook_config=hook_config,
             profiler_config=profiler_config,
         )
```

```
In [28]: estimator.fit({"training": "s3://ud-project-5/"}, wait=False)
```

```
INFO:sagemaker.image_uris:image_uri is not presented, retrieving image_uri based on instance_type, framework etc.
INFO:sagemaker:Creating training-job with name: dog-pytorch-2024-10-23-08-04-36-371
```

# EC2 Training

I selected "Deep Learning OSS Nvidia Driver AMI GPU PyTorch 2.3 (Amazon Linux 2)" on the AMI
list, so I can proceed with the AMI compatible with the Pytorch environment, which we are
using. We chose the G5.xlarge instance type for training our small model, as it performs
excellently due to being a new instance.

```
        #_
  ,  ~\_  ####_         Amazon Linux 2
 ~~   \_#####\
  ~~       \###|        AL2 End of Life is 2025-06-30.
  ~~       \#/ ___
   ~~       V~' '->
     ~~~         /      A newer version of Amazon Linux is available!
      ~~._.   _/
        _/ _/            Amazon Linux 2023, GA and supported until 2028-03-15.
      _/m/'                  https://aws.amazon.com/linux/amazon-linux-2023/

5 package(s) needed for security, out of 5 available
Run "sudo yum update" to apply all updates.
=====================================================================
AMI Name: Deep Learning OSS Nvidia Driver AMI GPU PyTorch 2.3.1 (Amazon Linux 2)
Supported EC2 instances: G4dn, G5, G6, Gr6, G6e, P4, P4de, P5, P5e
* To activate pre-built pytorch environment, run: 'source activate pytorch'
NVIDIA driver version: 550.90.07
CUDA versions available: cuda-12.1
Default CUDA version is 12.1

Release notes: https://docs.aws.amazon.com/dlami/latest/devguide/appendix-ami-release-notes.html
AWS Deep Learning AMI Homepage: https://aws.amazon.com/machine-learning/amis/
Developer Guide and Release Notes: https://docs.aws.amazon.com/dlami/latest/devguide/what-is-dlami.html
Support: https://forums.aws.amazon.com/forum.jspa?forumID=263
For a fully managed experience, check out Amazon SageMaker at https://aws.amazon.com/sagemaker
=====================================================================
[root@ip-172-31-75-205 ~]# ls
dogImages  dogImages.zip  solution.py  TrainedModels
[root@ip-172-31-75-205 ~]#
```

```
[root@ip-172-31-75-205 ~]# ls
dogImages   dogImages.zip   solution.py   TrainedModels
[root@ip-172-31-75-205 ~]# cd TrainedModels/
[root@ip-172-31-75-205 TrainedModels]# ls
model.pth
[root@ip-172-31-75-205 TrainedModels]#
```

We ran a Python script on EC2 rather than a Jupyter Notebook, which requires a distinct environment. Additionally, we didn't leverage Sagemaker's hyperparameter tuning, opting for hard-coded values instead. Apart from these changes, the code closely resembles hpo.py, which was used as our estimator entry point in the notebook.

# Lambda Function Setup

The lambda_function.py file receives an event and context, decodes the event with base64, and forwards the decoded data to a SageMaker endpoint. The function then processes the endpoint's response, converting it into a JSON object that consists of HTTP status codes and headers. Furthermore, it employs the Python logging module to log debug messages during the entire procedure.

# Concurrency and Auto Scaling

I assigned three instances to run my Lambda function concurrently in order to cut costs.

| Provisioned concurrency | | C Edit Remove |
|---|---|---|
| Provisioned concurrency | Status | |
| 3 | ⊘ Ready | |

I set up auto scaling for my endpoint to activate when the number of invocations hits 20 per second, with a scale-in period of 30 seconds and a scale-out period of 30 seconds. This configuration was aimed at optimizing costs while managing high traffic.

| Endpoint runtime settings | | Update weights | Update instance count | Configure auto scaling | | |
|---|---|---|---|---|---|---|
| weight | Elastic Inference | Instance type ▽ | Current instance count ▽ | Desired instance count ▽ | Instance min - max | Automatic scaling |
| - | | ml.m5.large | 1 | 1 | 1 - 3 | Yes |