# Project Opeationalizing Machine Learning on SageMaker

## Geraldo Margjini

## June 15, 2025

**Abstract**

This document presents the write-up of the project, including answers and explanations to the assigned questions.

# Contents

# 1 Step 1: Training and Deployment on Sagemaker

## 1.1 Initial Steps

I selected an ml.t3.medium instance because the dataset I'm working with is small. This instance type is among the most cost-effective options and is sufficient for the task at hand. Moreover, it falls within my AWS Free Tier usage.
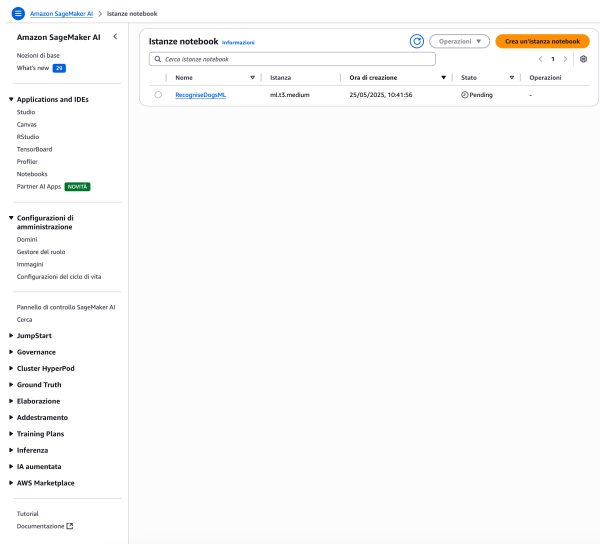


Figure 1: screenshot showing a created notebook instance on Sagemaker.

## 1.2 Download the data in the S3 Bucket

So the next step is to create a new S3 Bucket that i called *reco-dogs* and run the three first cell to download the data of different dog breads ( 133 dog breads, 6680 training, 835 validation and 836 test images).
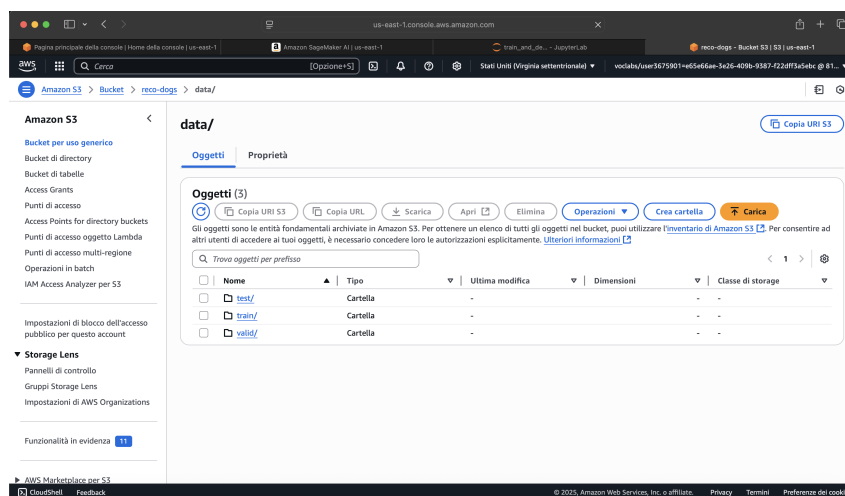


Figure 2: screenshot showing a created and downloaded the Data in my S3 Bucket.

## 1.3 Single Instance Training: Training and Deployment

I created a tuning job with the following parameter and got

```
{'batch_size': 256, 'learning_rate': '0.002896533293199136'}
```

as best hyperparameters.

```
tuner = HyperparameterTuner(
    estimator,
    objective_metric_name,
    hyperparameter_ranges,
    metric_definitions,
    max_jobs=2,
    max_parallel_jobs=1,  # you once have one ml.g4dn.xlarge instance available
    objective_type=objective_type
)
```

So then i trained using this parameter with a ml.m5.xlarge:

```
estimator = PyTorch(
    entry_point='hpo.py',
    base_job_name='dog-pytorch',
    role=role,
    instance_count=1,
    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,
)
```

So following we have the screenshot of the Endpoint created:

## 1.4 Multy Instance Training: Training and Deployment

For the Multy instance i used the value for the

```
instance_count=4
```
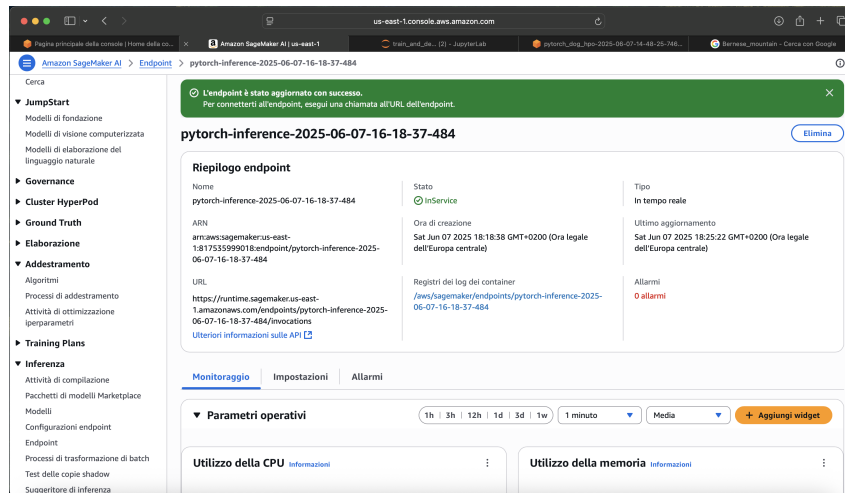
to train four instances at time.

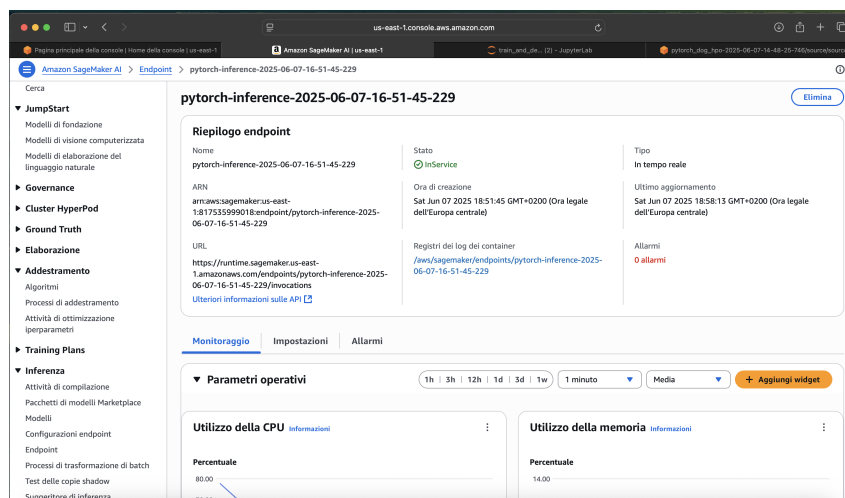Figure 3: screenshot showing a created Endpoint for the single instance.



Figure 4: screenshot showing a created Endpoint for the multy instance.

# 2 Step 2: EC2 Training

For this section i run an EC2 with training image:

```
amazon/Deep Learning AMI Neuron (Amazon Linux 2023) 20250520
```

using an m5.2xlarge because it is a good tradeoff between the task that i have to accomplish and the costs.

So I started it with the WAS UI and connected it with an SSH:

```
ssh -i "ger.pem" ec2-user@ec2-54-158-209-129.compute-1.amazonaws.com
```

Once we are connected to our own EC2 we have to do the following step to: Download your own data, Create the target folder for the trained model, so we can copy the code for the train in a file called *solution.py* and run it.
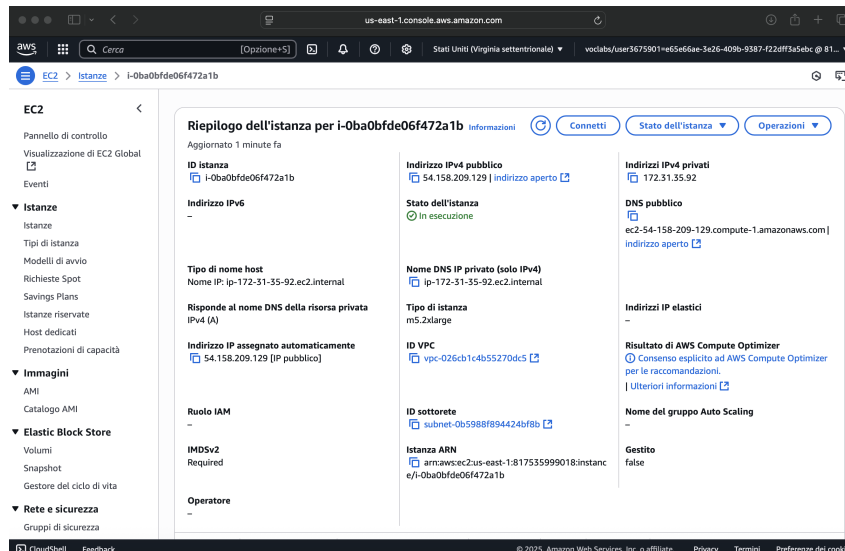
Figure 5: screenshot showing a created EC2.

```
wget https://s3-us-west-1.amazonaws.com/udacity-aind/dog project/dogImages.zip
unzip dogImages.zip
mkdir TrainedModel
nano solution.py
python solution.py
```



Figure 6: screenshot showing a finished training model.

I found some problem with the env that i had to create and install different missing pips.

## 2.1 Differances between the Sagemaker Notebook and EC2

In the Sagemaker we create a Nootebook instance where we write the code to create and instance to perform the training but in the EC2 we directly run the script code "localy" in the EC2.

- In the Sagemaker the hyperparameters are passed to hpo.py as an argument, but in the EC2 we run directly the hpo in the script.

- Using the EC2 we limited just to this instance but if we used the Sagemaker Notenook we can used the coded seen befored to perform distribuited trainig

- The final model, in the EC2 we have the model saved in then target folder, if we would like to save it in the could we should run the code wo upload it in the S3 bucket but using the Sagemaker Nootebook the "artifact" with is the model and the hpo are directly saved in a S3 target foolder.

# 3   Step 3: Lambda Fuction

So the Lambda function created takes the endpoint created in the previous steps and makes inferences passing the S3 url of it. So in this i passed the name of my endpoint to my lambda function:   *pytorch-inference-2025-06-15-15-57-32-476*.

# 4   Step 4: Lambda Security and Testing

Now i sended to the Endpoint this image path to do an inference:

`https://s3.amazonaws.com/cdn-origin-etr.akc.org/wp-content/uploads/2017/11/20113314/Ca`

Which is the same one used to do inference in the previous steps, for the security consideration i created a restricted policy to have just the specific Endpoint Access so this LambdaFuction can't perform other action to my other Endpoints.



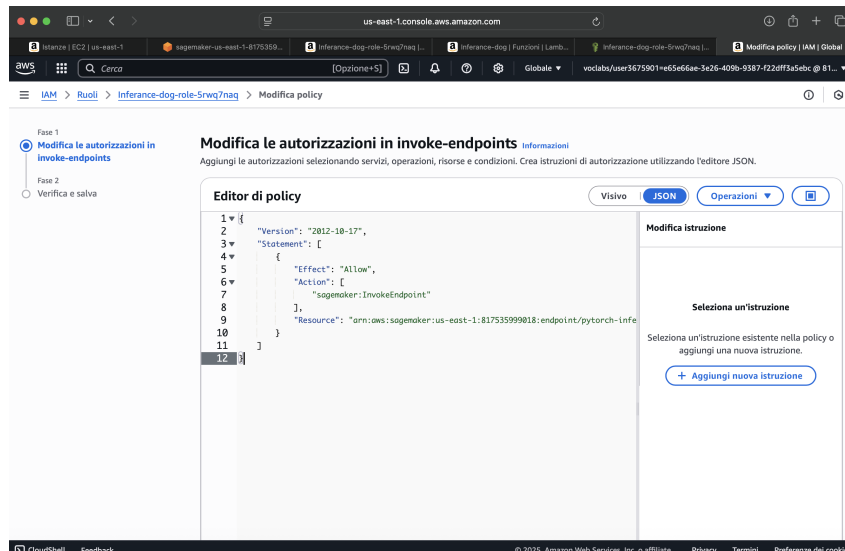Figure 7: screenshot showing a inference made with my LambdaFuction.

Figure 8: screenshot showing the specific policy for my LambdaFunction.

# 5 Step 5: Concurrency and Auto-Scaling

## 5.1 Concurrency

Concurrency refers to the ability of Lambda functions to handle multiple requests simultaneously.

We can configure either Reserved or Provisioned concurrency for our function. Provisioned concurrency provides faster response times but comes with higher cost.

Because our expected request volume is low, we don't need very high concurrency. I have set Provisioned concurrency to 3, which is sufficient for our needs, and Reserved concurrency to 100.

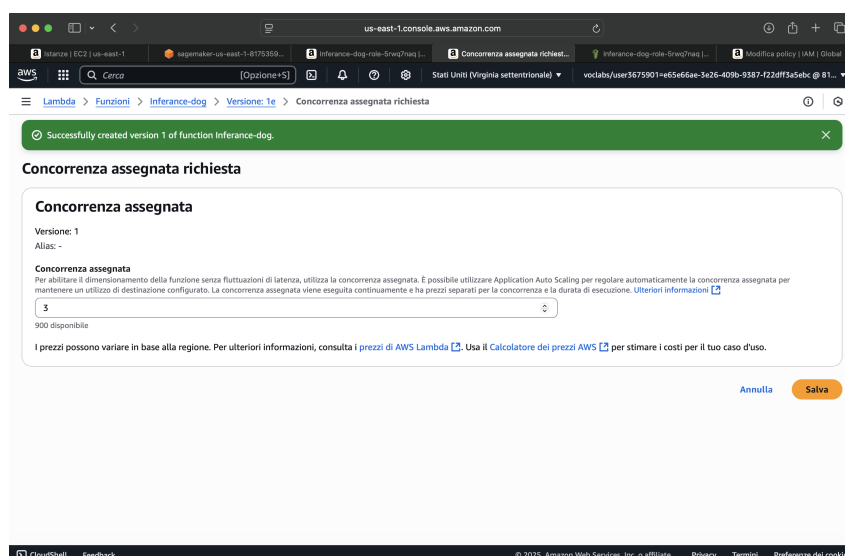Below is a screenshot of the Lambda concurrency settings:



Figure 9: screenshot showing the concurrency settins.

## 5.2   Auto-scaling

Auto-scaling refers to the ability of endpoints to service multiple lambda function requests at once. I chose to autoscale endpoints to 4 instances maximum, with scale in cool down time of 30 seconds and scale out cool down time of 300 seconds. These settings are sufficient for our project needs and workload.