# Operationsalizing an AWS ML project

The goal of this project was to train and deploy a dog breed classification model using all possible best practices related to optimization, cost and security.

## Step 1: Training and deployment on Sagemaker

As model training won't be performed directly within the working notebook, a simple ml.t2.medium 5GB instance was selected for this project. The conda_python3 kernel was used to run the main Jupyter notebook.
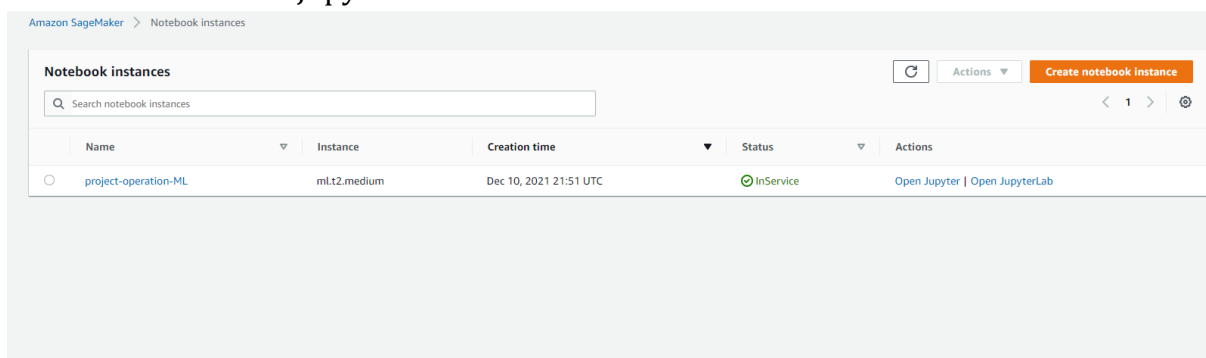


Figure 1: SageMaker notebook instances.

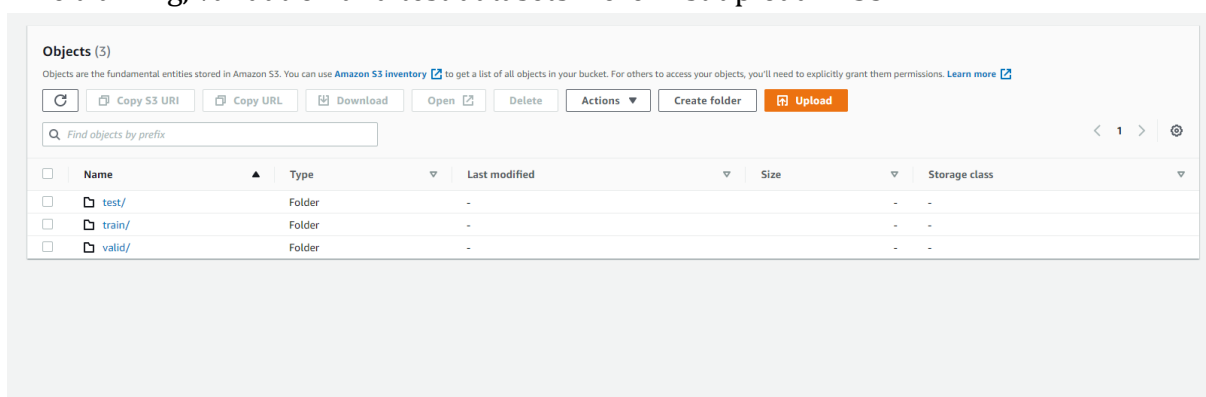The training, validation and test datasets were first upload in S3.



Figure 2: S3 console of the "dog-breed-project" directory.

In order to complete model training in a timely manner, an "ml.g4dn.xlarge" instance was selected. This instance type is the cheapest one within the "Accelerated computing" category.
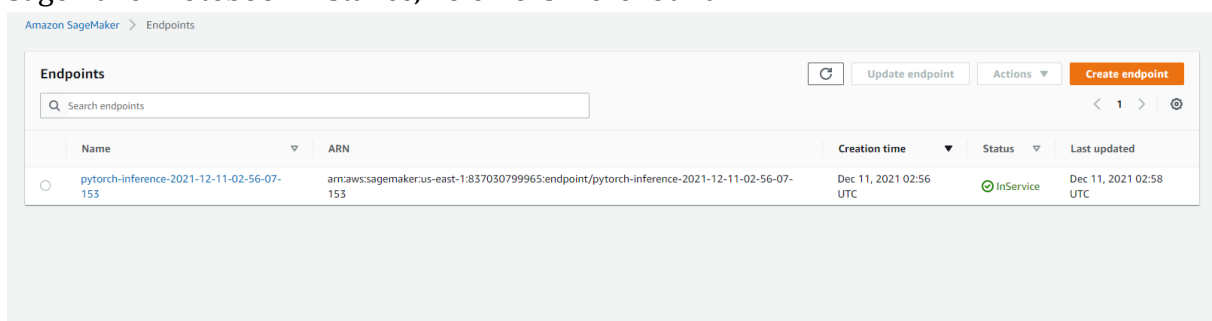
| Accelerated Computing | vCPU | Memory | Price per Hour |
|---|---|---|---|
| ml.p3.2xlarge | 8 | 61 GiB | $3.825 |
| ml.p3.8xlarge | 32 | 244 GiB | $14.688 |
| ml.p3.16xlarge | 64 | 488 GiB | $28.152 |
| ml.p2.xlarge | 4 | 61 GiB | $1.125 |
| ml.p2.8xlarge | 32 | 488 GiB | $8.64 |
| ml.p2.16xlarge | 64 | 732 GiB | $16.56 |
| ml.g4dn.xlarge | 4 | 16 GiB | $0.736 |
| ml.g4dn.2xlarge | 8 | 32 GiB | $0.94 |
| ml.g4dn.4xlarge | 16 | 64 GiB | $1.505 |
| ml.g4dn.8xlarge | 32 | 128 GiB | $2.72 |
| ml.g4dn.12xlarge | 48 | 192 GiB | $4.89 |
| ml.g4dn.16xlarge | 64 | 256 GiB | $5.44 |

Figure 3: Accelerated Computing EC2 instances pricing.

To ensure that the training algorithm can utilize the provided GPU, the hpo.py script has been modified. The model and the inputs/targets are now automatically passed to the available device. With those changes, the hyperparameter optimisation took only 9 minutes.

The model was retrained using the best-hyperparameters, this time, using a spot ml.m5.large instance to drastically reduce cost.

The model was then deployed on an ml.m5.large instance and invoked from the SageMaker notebook instance, no errors were found.



Figure 4: Endpoints console.

The same model was then retrained using multi-instance training, more specifically, 5 ml.m5.large instances were used.

Training time with 1 instance: 22 minutes.
Training time with 5 instance: 24 minutes.

On a personal level, I am quite unsure of the benefits of just multiplying the number of training instances. All it does is to create 5 instances, perform the same training independently and save 5 different models without assembling/aggregating the model weights. As demonstrated in the following images:
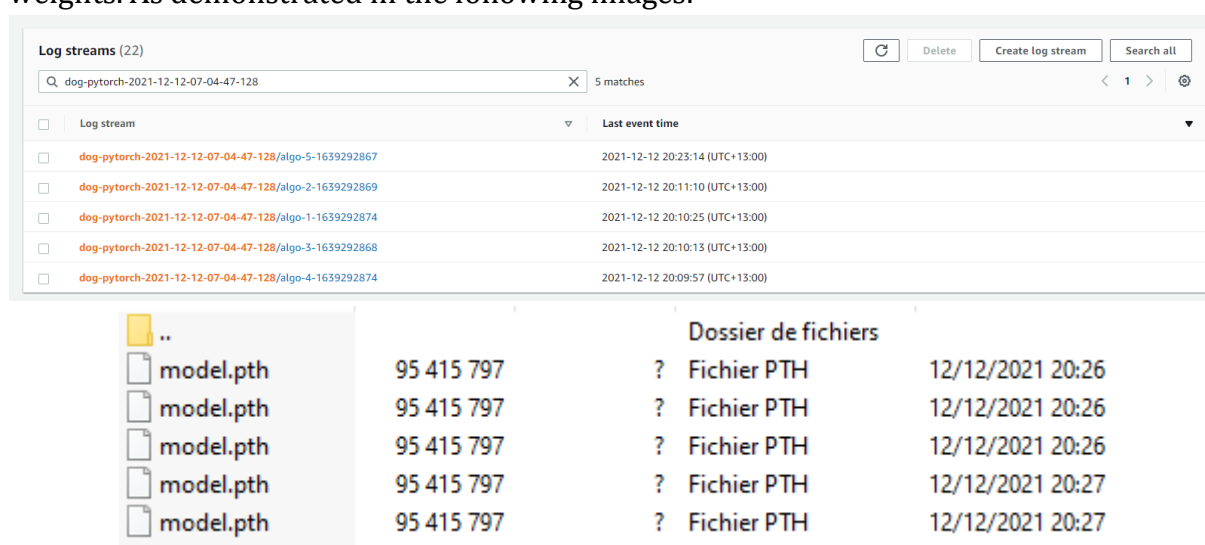


Figure 5: Multi-instance training output.

To perform proper data distributed training (and so increase the training speed) on Deep learning models, one should pick a multi-GPU instance(s) and follow the [SageMaker documentation](#).

## Step 2: EC2 Training

The same model was trained again, this time, using a manually deployed EC2 instance.
To be consistent with the step 1 of this project, a spot m5.large instance with an Deeplearning Amazon Linux 2 AMI was selected.

Unlike creating a training job from a notebook instance, which requires to provide the training data location and where the model's artifacts will be saved (as S3 paths), training a model inside an EC2 instance is similar to training a model on a local computer. Everything happen locally and the respective data/model folders must be manually created.

The following image shows the last bash commands used to train the model and access its artifacts.



```
[root@ip-172-31-20-97 ~]# python solution.py
Downloading: "https://download.pytorch.org/models/resnet50-19c8e357.pth"
100%|
Starting Model Training
saved
[root@ip-172-31-20-97 ~]# ls
dogImages  solution.py  TrainedModels
[root@ip-172-31-20-97 ~]# cd TrainedModels
[root@ip-172-31-20-97 TrainedModels]# ls
model.pth
[root@ip-172-31-20-97 TrainedModels]#
```

Figure 6: EC2 instance CLI.

BONUS: As an extra-challenge, I created an endpoint from this model artifact.

Step 1: Extract model.pth from the EC2 instance.
The model artifact was first moved to S3 using the following command:
aws s3 cp model.pth s3://sagemaker-us-east-1-836049028555593

Step 2: Transform to .tar.gz and upload to S3.
The model artifact was first download from S3.
Then, it was transformed in the appropriate format using the following command:
!tar -cvsf model.tar.gz model.pth
Finally, the compressed model was uploaded back to S3.

Step 3: Create and deploy a Pytorch model.
The same methodology as in the first step was used to deploy the model.

For more details on how the model has been transformed and deployed, look at the "Bonus section: Deploy the model trained on EC2" section of the notebook "train_and_deploy-solution".

## Step 3: Lambda function setup

A lambda function was created in order to easily perform inference from other applications/services.

Lambda function do not automatically come with the sagemaker library, however, one can invoke a deployed endpoint via a boto3 runtime instance. First, the data is extracted from the received event, then, it is passed to the invoke_endpoint function with its type.

The prediction output is an utf-8 encoded json object. It is decoded, then send back as a json object within the "body" section of the function output.

## Step 4: Security and testing

The sagemakerfullaccess policy was added to the function IAM role in order to let it access the deployed endpoint. I did not find a SageMaker policy that would only give access to deployed endpoints. Therefore, I do not think the security can be improved through just the Lambda  IAM role.



Figure 7: Lambda function role summary.

The lambda function was then written and deployed as shown in the following figure.



Figure 8: Lambda function code.

The test event did not raise any errors, and an array of 133 probabilities was obtained.



Figure 9: Lambda function execution result.

# Step 5: Concurrency and auto-scaling

To reduce potential latency due to the lambda function, reserved concurrency was set up. Up to 3 functions can now be instantiated simultaneously in order to meet traffic demand.



**Concurrency**

| | |
|---|---|
| Function concurrency | Reserved concurrency |
| Use reserved concurrency | 3 |

Figure 10: Lambda concurrency parameters.

To reduce potential latency due to the model endpoint, a custom scaling policy was set up. Up to 3 instances can now be instantiated simultaneously in order to meet traffic demand. More specifically, an additional instance will be created if the average CPU usage is greater than 70% for more than 30 seconds.

This custom policy was created from the notebook instance in two steps:

1) Register the endpoint and its associated variant.

2) Set up the scaling with the client.put_scaling_policy() function.

For more details on how the policy has been created and attached to the enpoint, look at the "Configure autoscaling policy" section (the last one) of the notebook "train_and_deploy-solution".

Here are screenshots of the endpoint console that prove the policy has been correctly attached:



Figure 11: Endpoint console.

In a real world situation, one should adapt concurrency/auto-scaling configuration depending on the observed traffic. For example, if after several days I realise that my endpoint instance CPU utilisation tend to stagnate around 80%, I would consider raising the threshold to 90%, otherwise I would be paying another instance for nothing.

One should also look out for where the throughput bottleneck is. If the Lambda function is responsible for the latency, increasing the endpoint auto-scaling capabilities would induce useless additional cost.

Finally, the initial concurrency/auto-scalling configuration should be decided depending on the latency tolerance decided by the stakeholders. If the application is very latency sensitive, I would consider provisioned concurrency for the lambda function and a very short scale-out cooldown for the model endpoint. Only after having properly quantify the traffic patterns, I would consider cost optimistation/strategies.