# Operationalizing an AWS Machine Learning Project

1) Write a short justification for why you chose the SageMaker notebook instance type you did (consider cost computing power, and speed of launching for each instance type)

   **Answer:** I used the **ml.t3.medium** instance which is a general purpose instance. It provides average computing power and speed at low cost. Considering the nature of this project, there was no need for extra high performance computing, hence the choice of this instance to save cost.

2) After launching your SageMaker instance, take a screenshot of your sagemaker dashboard's Notebooks > Instances section to show what you have done.
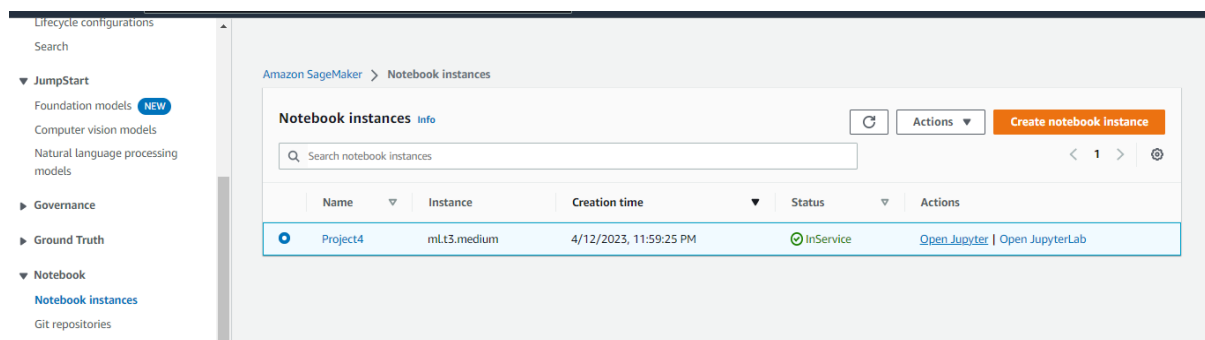


Fig 1: SageMaker notebook instance

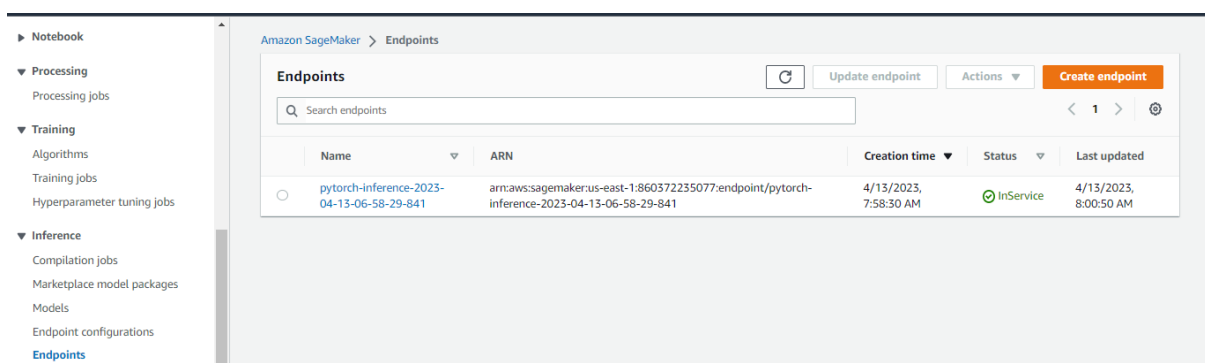3) Take a screenshot of your deployed endpoint (single instance training) and note its name



Fig 2: Deployed endpoint after single instance training

4) Write a justification for the type of EC2 instance you created.

**Answer**: Initially, I used the t2.medium instance to save cost. However, when I tried to run the solution.py file, I encountered a numpy import error and I tried to address it by activating pytorch in EC2. The EC2 t2.medium instance was unable to activate pytorch and I switched to a g3s.xlarge instance as suggested by the error message. Relevant screenshots are attached below.
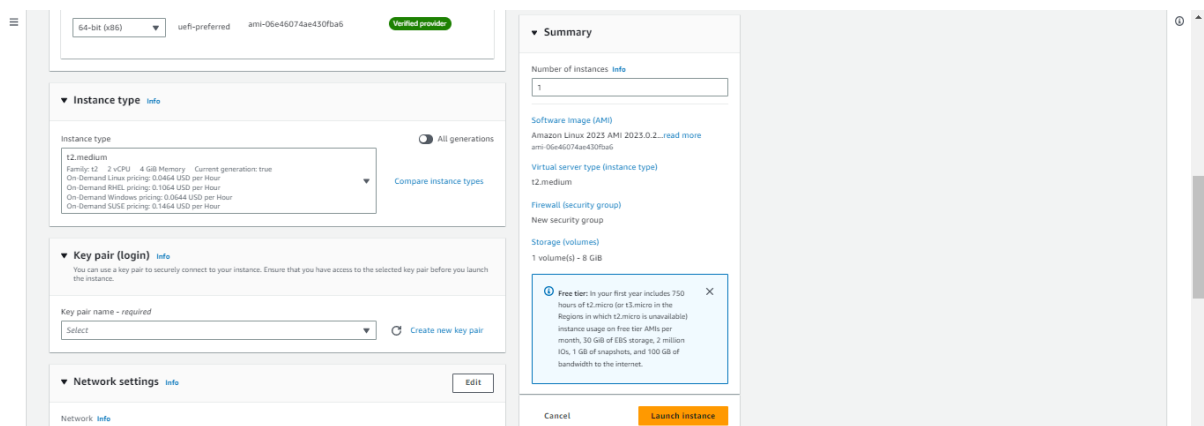


Fig 3: Initial EC2 instance



Fig 4: EC2 error message

5) Open the TrainedModels directory in your EC2 instance and take a screenshot of the model that has been saved in it.



Fig 5: Model in the TrainedModels directory

6) Write on the differences between the code in "train_and_deploy-solution.ipynb" and "ec2train1.py"

**Answer:** In ec2train1.py, the parameters, estimators, and tuners are not explicitly defined. ec2train1.py is optimized by building functions to handle model training and testing. These functions can be called on data contained in EC2 when needed. In the train_and_deploy-solution notebook, however, the data is imported and model training and testing is done using explicitly defined parameters, estimators, and tuners. In addition, libraries such as boto3 and sagemaker cannot be used in ec2train1.py.

7) Take a screenshot of the lambda function after running the test event, and provide the list of numbers given in the output.



Fig 6: Lambda function test output

```
[-0.3194071054458618,       0.06982807070016861,       -0.10188700258731842,
0.05443248152732849,        0.23411345481872559,        0.2769347131252289,
0.016361981630325317,       0.09543728828430176,       -0.23329131305217743,      -
0.1326563060283661,         0.08965423703193665,        0.0579850859940052,        -
0.15090619027614594,        0.10275612771511078,        0.1422945261001587,
0.02920246124267578,       -0.02602965384721756,       -0.39369359612464905,       -
0.20661674439907074,        0.05201426148414612,        0.020475119352340698,
0.1092258021235466,         0.14406603574752808,        0.16800859570503235,       -
0.2569163739681244,        -0.4271831512451172,         0.18519136309623718,       -
0.3146352469921112,         0.034361548721790314,       -0.2248910516500473,
0.07397255301475525,        0.01564362645149231,       -0.2630861699581146,        -
0.2257373183965683,        -0.19176961481571198,       -0.13752171397209167,       -
0.03265161812305450,       -0.15763114392757416,       -0.059963926672935486,      -
0.03800392150878906,       -0.2107524275779724,        -0.07816909998655319,       -
0.1020075805187225,        -0.35978806018829346,        0.0466545969247818,        -
0.13159333169460297,        0.02251681685447693,        0.10969046503305435,       -
0.12295417487621307,       -0.048323310911655426,      -0.026371128857135773,      -
0.07860969007015228,       -0.2702891230583191,         0.07395516335964203,       -
0.31811490654945374,        0.15498119592666626,        0.020910531282424927,      -
0.13895606994628906,       -0.1581280380487442,         0.056791454553604126,      -
0.1681289672851562 5,      -0.001905910670757293 7,     -0.3697312176227569 6,      -
0.4153045117855072,        -0.12665320932865143,       -0.47466224431991577,       -
0.034660764038562775,      -0.3707352876663208,        -0.09482190012931824,       -
0.1700029969215393,         0.2272071689367294 3,       -0.180681973695755,         -
0.30132168531417847,       -0.1123449057340622,        -0.23300659656524658,
0.09635794162750244,       -0.29833337664604187,       -0.25851720571517944,
0.1241689175367355 3,      -0.09933248162269592,       -0.031318724155426025,      -
0.015487030148506165,      -0.10894891619682312,        0.1430465281009674,        -
```
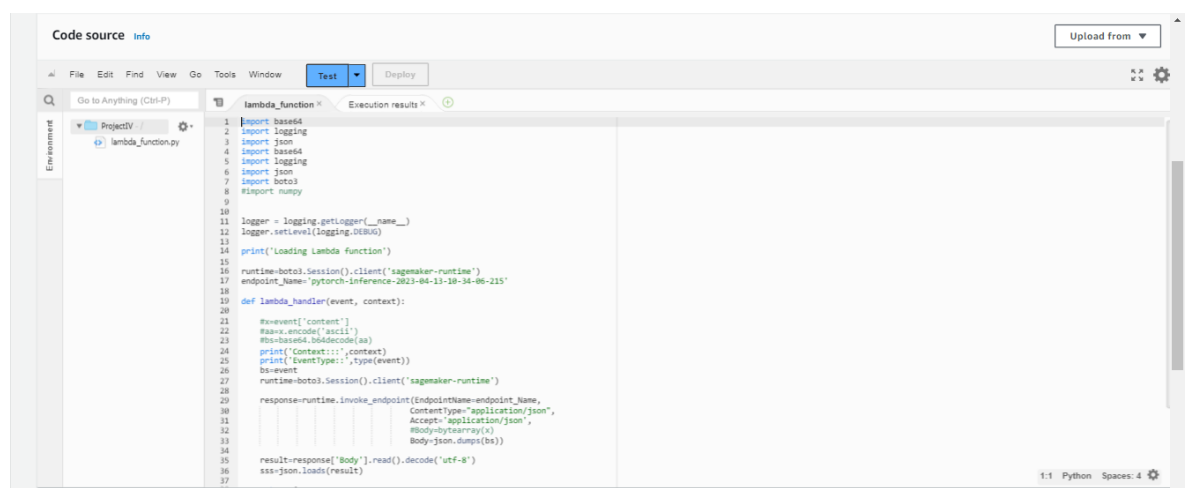
0.24369531869888306,      -0.22964990139007568,      -0.17065121233463287,      -
0.1537196934223175,          -0.10664661228656769,          -0.0409412682056427,
0.030008748173713684,      0.006851717829704285,      -0.4325335621833801,      -
0.12947550415992737,      -0.056942954659461975,      -0.1904909610748291,      -
0.17198391258716583,      0.133949875831604,      -0.1790093183517456,      -
0.1203477531671524,      0.06557272374629974,      -0.41081827878952026,      -
0.039602264761924744,          -0.3460289239883423,          -0.27233603596687317,
0.06590510904788971,          -0.17482388019561768,          -0.1722773313522339,
0.02484789490699768,          -0.3098023235797882,          -0.17887285351753235,
0.005011945962905884,          -0.15574386715888977,          -0.27694201469421387,
0.11342388391494751,      -0.2868531346321106,      0.013332098722457886,      -
0.19056998193264008,      -0.05950509474277749634,      -0.0750817060470581,      -
0.37999147176742554,          -0.29933661222457886,          -0.0784531682729721,
0.05784786492586136,      -0.36373087763786316,      -0.18035843968391418,      -
0.1644153892993927,      -0.4166147708892822,      -0.03248770534992218,      -
0.3103395402431488,      -0.4522474408149719,      -0.25467294454574585,      -
0.284565806388855]

Fig 7: List of numbers in Lambda function test output

8) Describe how the Lambda function is written and how it works

**Answer:** The lamda_handler function is the major component of this Lambda function. It takes in a picture (passed in as a url using the JSON file of the Lambda function) which it passes through the model contained in a SageMaker endpoint (the endpoint is defined within the Lambda function as endpoint_Name). The model contained in this specified endpoint is a dog breed classifier. Therefore, the output of this Lambda function is a list containing the probability of the dog in the input picture belonging to the various breeds in the dataset.



Fig 8: Lambda function

9) Take a screenshot of the policies attached to your Lambda function role
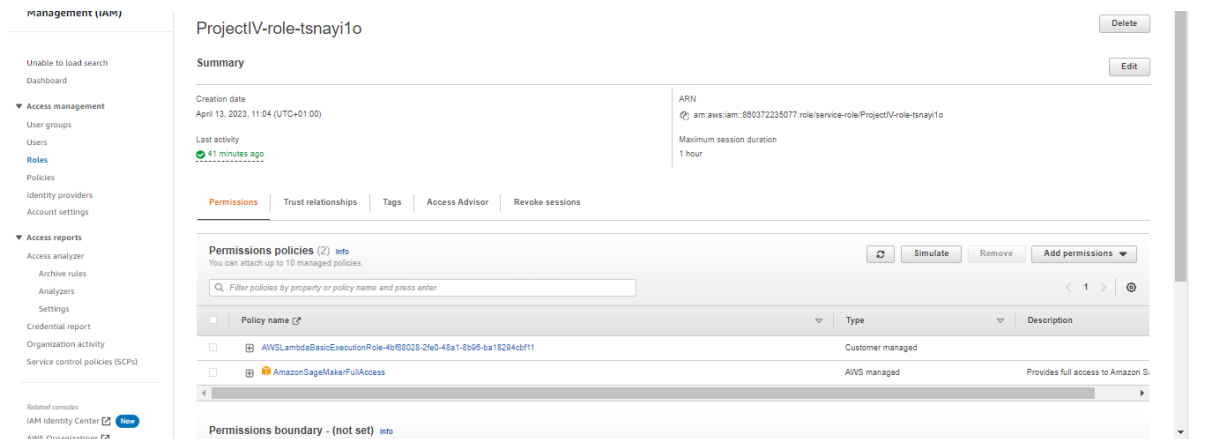
Fig 9: Policies attached to the Lambda function role before security considerations

10) Write about the security vulnerabilities in your AWS workspace

**Answer:** Attaching the AmazonSageMakerFullAccess policy to my Lambda function exposes my account to security risks through the function. This is because the Lambda function has access to all my endpoints in SageMaker.

To address this issue, I replaced the AmazonSageMakerFullAccess policy with a policy that allows my Lambda function access to only the pytorch inference endpoint that it passes data to. Appropriate screenshots are attached below.
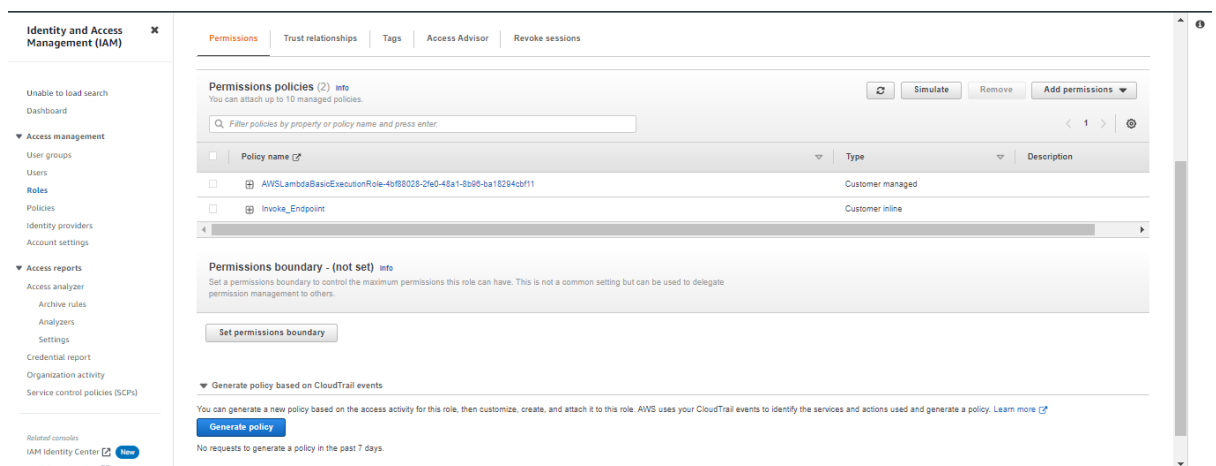


Fig 10: Policies attached to the Lambda function role after security considerations
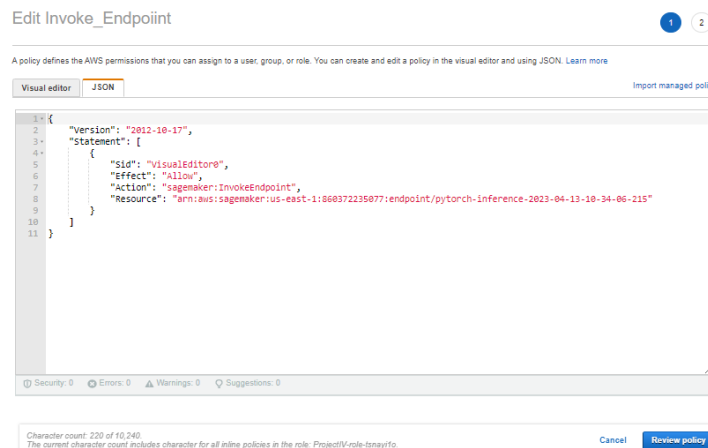
Fig 11: JSON for new policy attached to Lambda function

11) Write about the choices you made during concurrency and autoscaling setup

**Answer:**

Autoscaling: I set the target value to 40. This means that instances will be added when there are greater than or equal to 40 simultaneous requests to this endpoint. I set the value high to ensure that new instances are not activated unless absolutely necessary, ultimately leading to less cost.

I set the maximum instance count to 2 because I do not expect high traffic on my endpoint.

I set 30 for the scale in cool down value to ensure that extra resources are not employed for momentary spikes in request traffic. 30 minutes is enough time to ascertain that there is indeed high traffic, before deploying more resources.

I also set 30 for the scale out cool down value to ensure that resources are not retrieved for temporary reductions in request traffic. 30 minutes is sufficient time to ensure that there is indeed traffic reduction before retrieving deployed resources.

Concurrency: Since I do not expect high traffic on this project and provisioned concurrency needs to be lower than reserved concurrency, I picked 5 reserved instances and 3 provisioned instances to save cost.

Fig 12: Autoscaling configuration