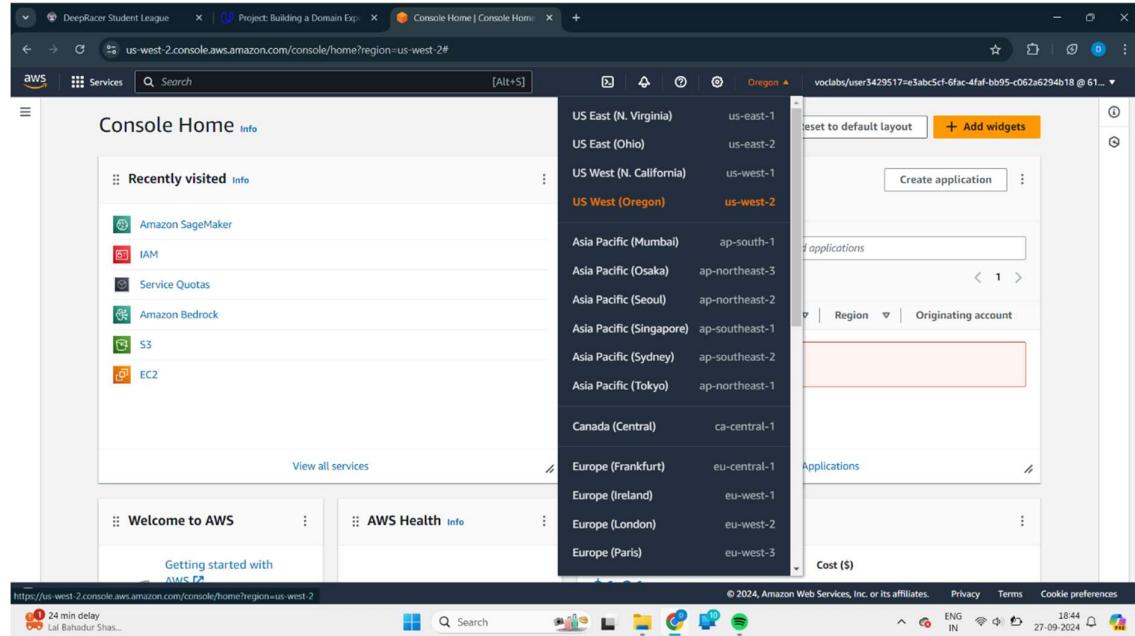
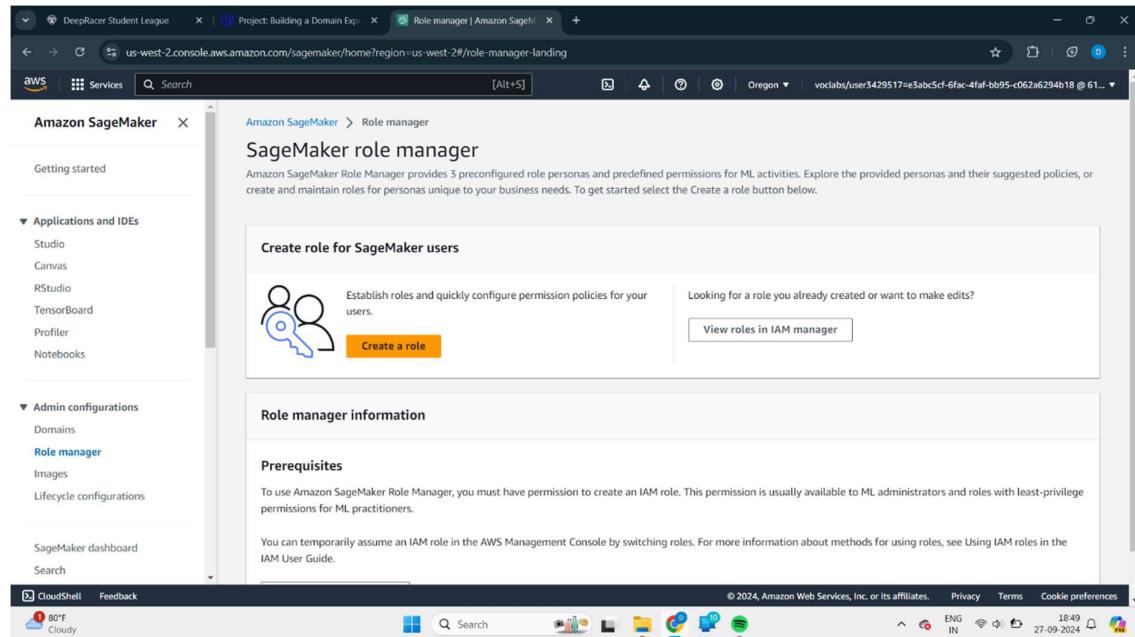


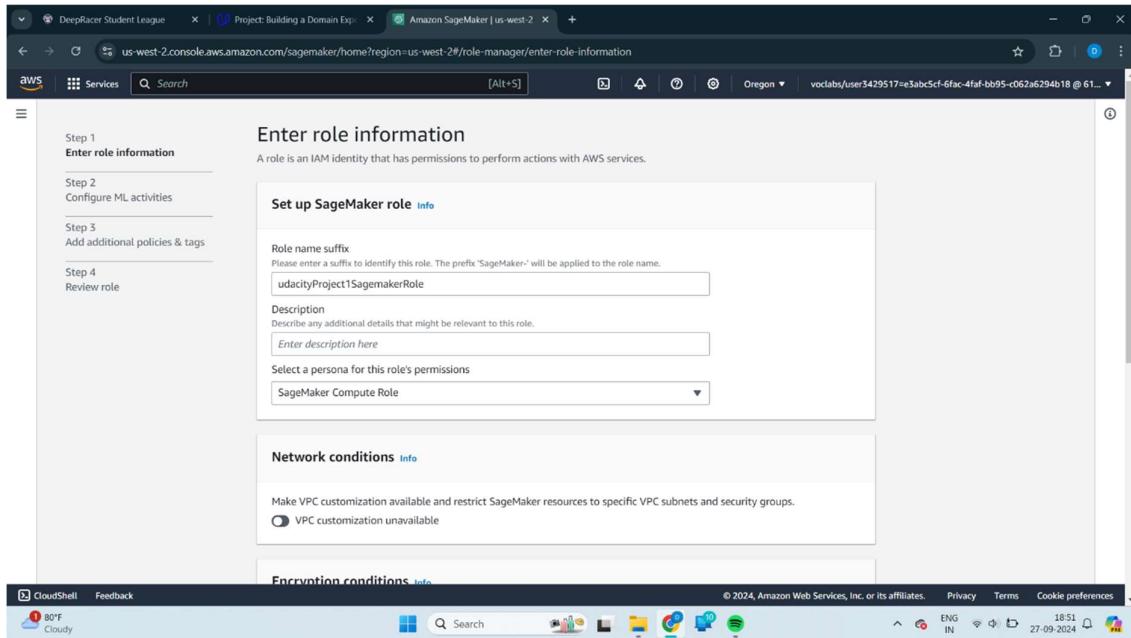
Screenshot of choosing us-west-2 in the region dropdown in the AWS console



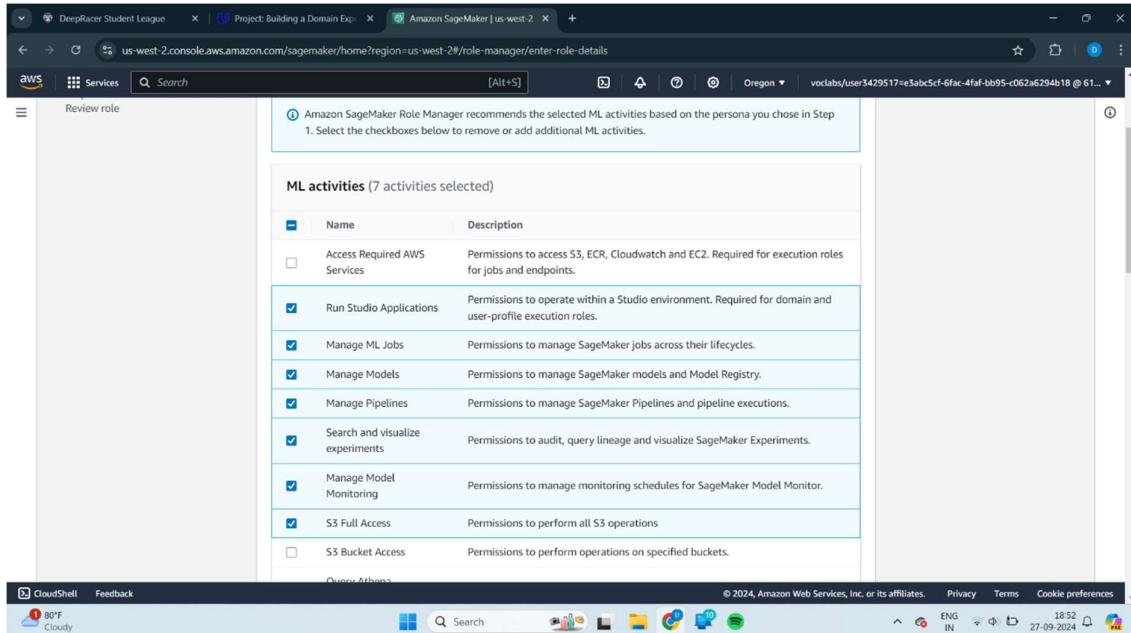
Screenshot of choosing the "create a role" button in the Role manager in AWS Sagemaker



Screenshot of the AWS Sagemaker role configuration with "SageMaker Compute Role" chosen for the role persona.



Screenshot of IAM role configuration with all ML activities checked after "Access required AWS services" up to "S3 full access"



Screenshot of IAM role configuration adding AmazonSageMakerFullAccess to the role

The screenshot shows the 'Add additional policies & tags' step of the IAM role configuration. On the left, a sidebar lists steps: Step 1 (Enter role information), Step 2 (Configure ML activities), Step 3 (Add additional policies & tags), and Step 4 (Review role). The main area is titled 'Add additional policies & tags' with the sub-instruction 'Attach existing IAM policies to this role. These policies will be included when IAM policies have been generated from this workflow.' Below this is a section titled 'Add additional IAM policies to this role' with a 'Info' link. A search bar shows 'AmazonSageMakerFullAccess'. A table lists one policy: 'Name' (AmazonSageMakerFullAccess) and 'ARN' (arn:aws:iam::aws:policy/AmazonSageMakerFullAccess). Below this is a 'Tags' section with an 'Info' link, showing 'No tags associated with the resource.' and a 'Add new tag' button.

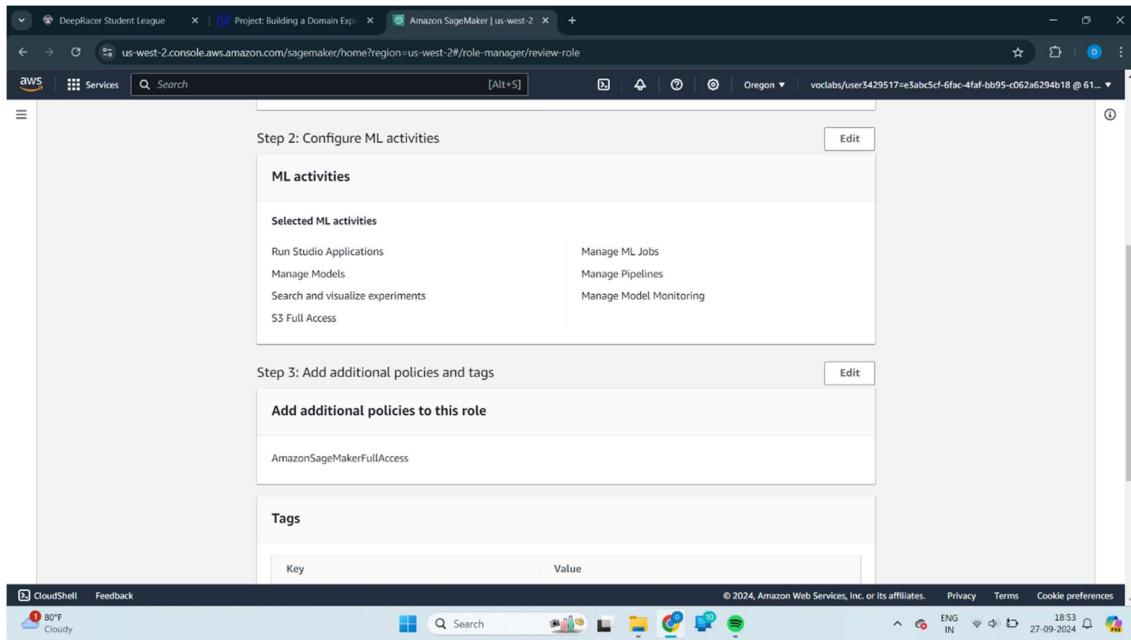
Screenshot of the final configuration of the Sagemaker IAM Role

The screenshot shows the final configuration of the Sagemaker IAM Role. It consists of two main sections: 'Step 1: Enter role information' and 'Step 2: Configure ML activities'.
Step 1: Enter role information
This section includes:

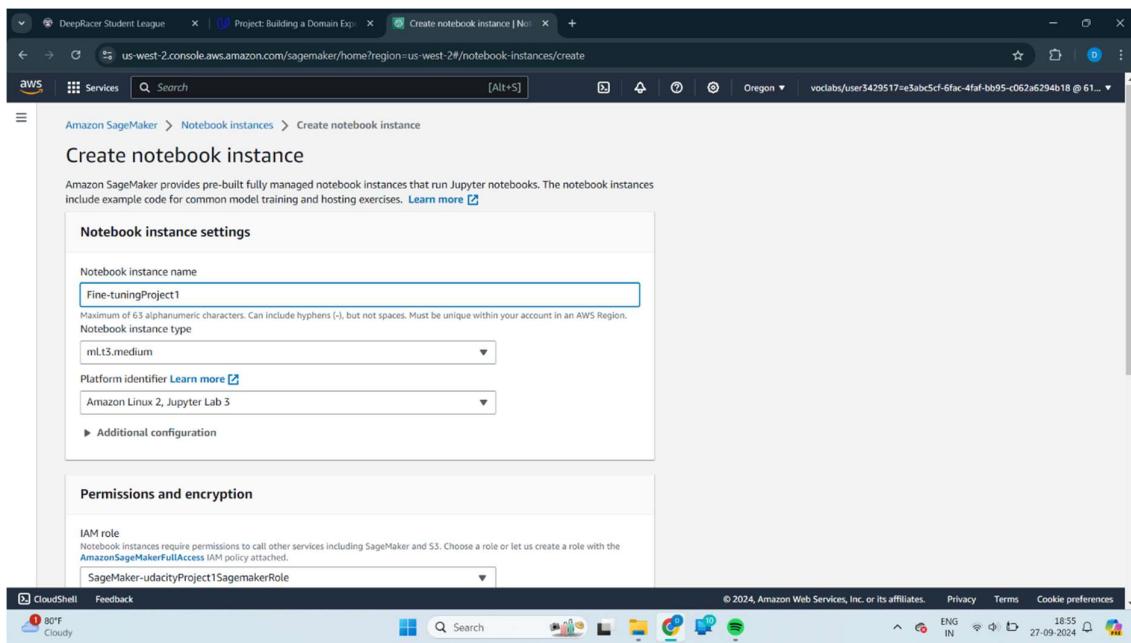
- Role details**: Shows 'Role name' (SageMaker-udacityProject1SagemakerRole) and 'Persona template' (SageMaker Compute Role).
- Network setup**: Shows 'VPC subnet(s)' and 'Security group(s)' both set to '-'.
- Encryption setup**: Shows 'Data encryption key(s)' and 'Volume encryption key(s)' both set to '-'.

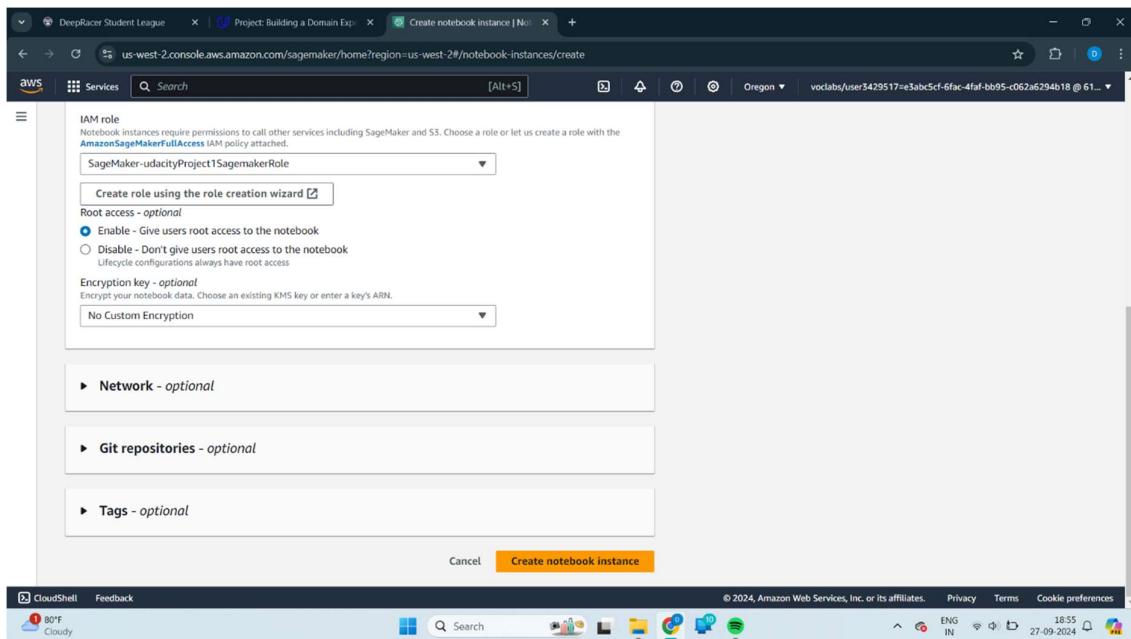
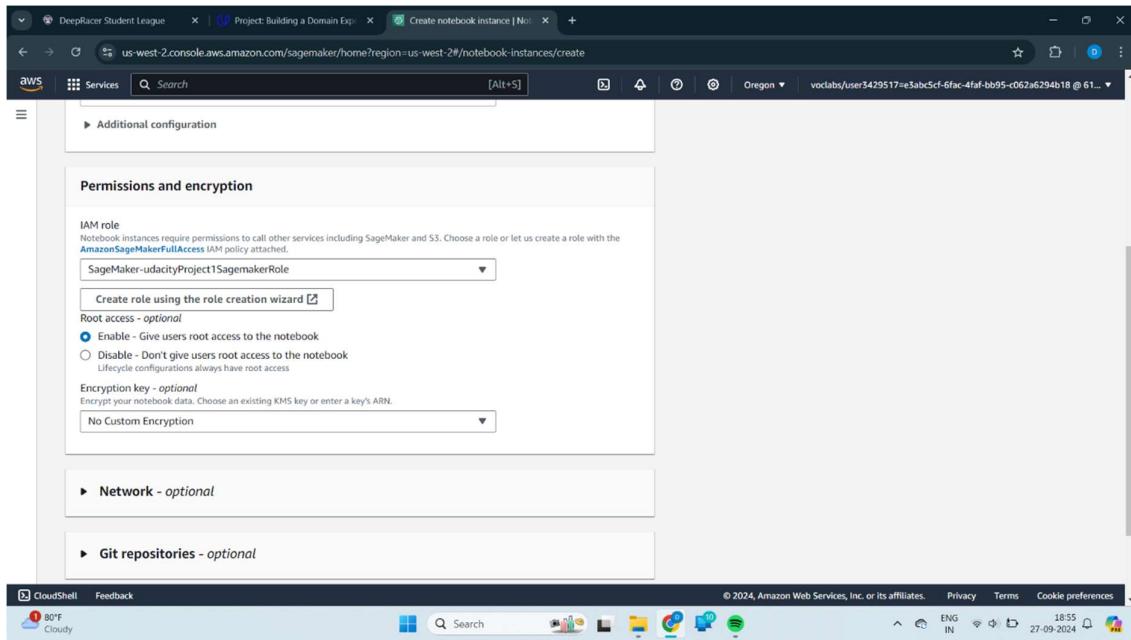
Step 2: Configure ML activities
This section includes:

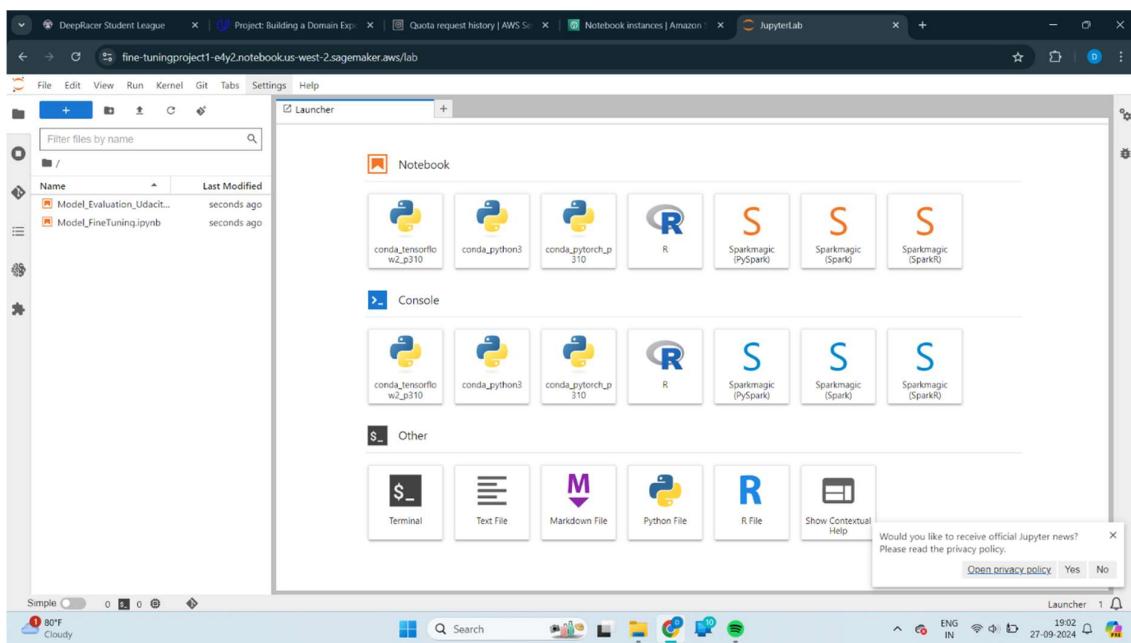
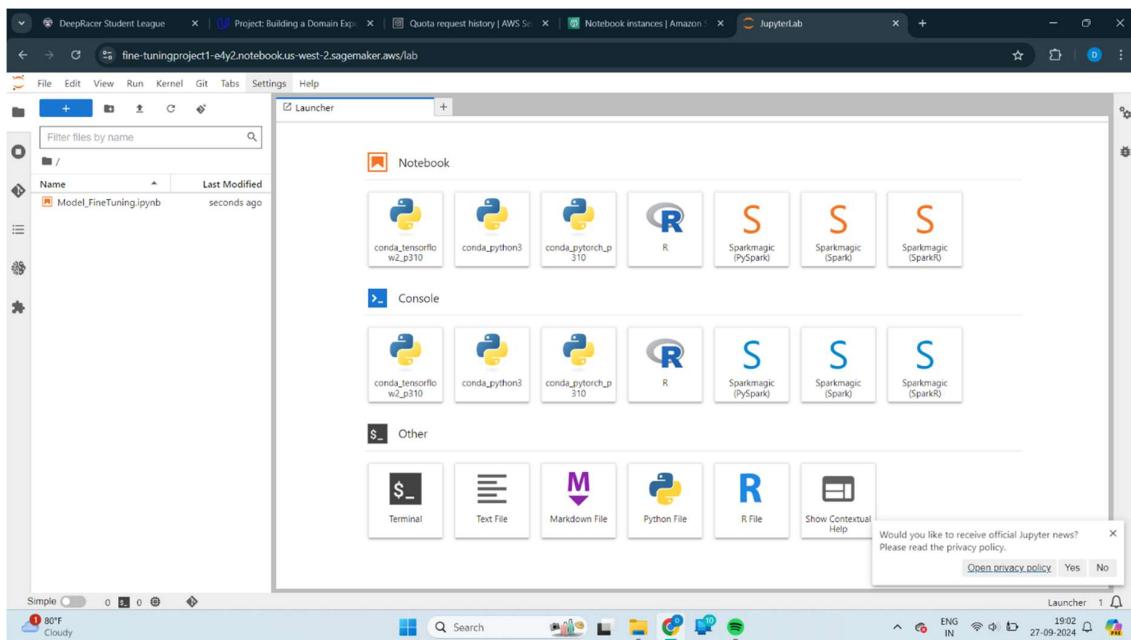
- ML activities**: Shows 'Selected ML activities' (Run Studio Applications, Manage ML Jobs).

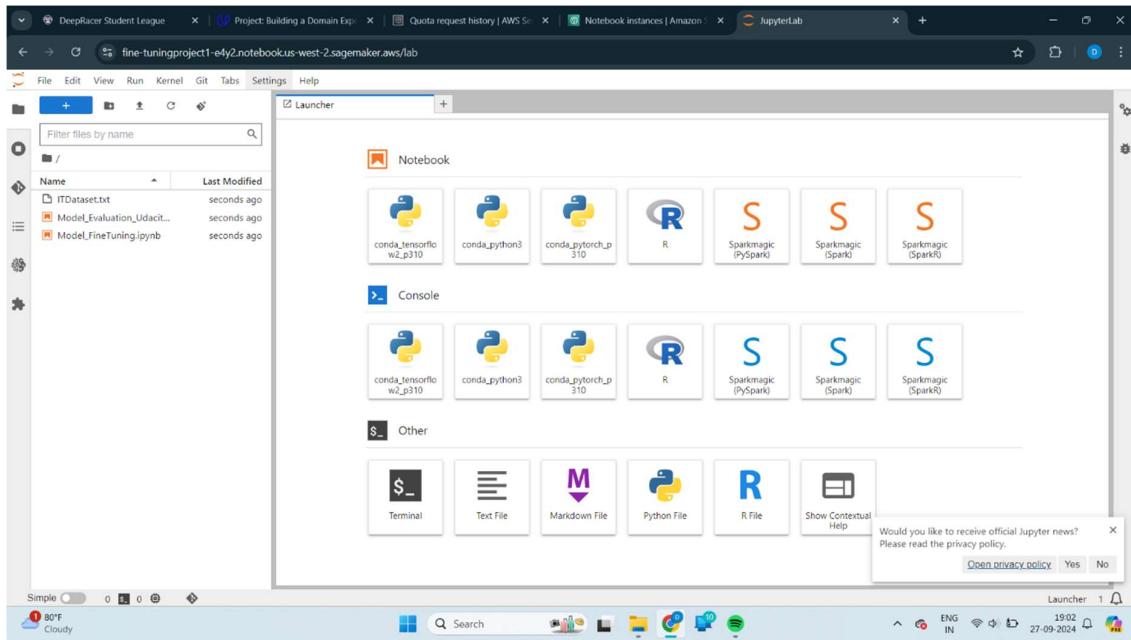


Screenshot from the AWS Sagemaker Notebook instance creation dialog with a custom name and choosing the Sagemaker role created above.

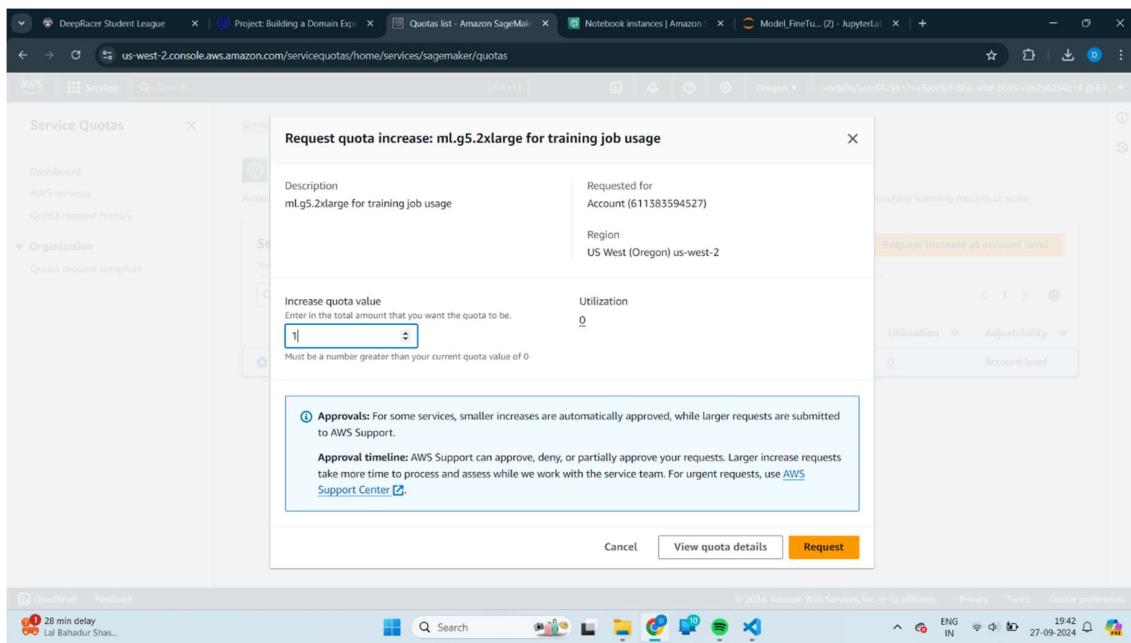








Screenshot of the service quota request form increasing the ml.g5.2xlarge for training job usage request to 1 instance.



Take a screenshot of the Model_Evaluation.ipynb file with the cell output as proof you completed this step of the project

This screenshot shows the JupyterLab interface with the 'Model_Evaluation.ipynb' notebook open. The left sidebar displays a file tree with 'Model_Evaluation_UdacityGenAIAWS.ipynb' selected. The main area contains a section titled 'Step 3: LLM Model Evaluation' with a sub-section 'Set Up'. A code cell at the bottom of this section contains the following Python code:

```
[1]: !pip install ipywidgets==7.0.0 --quiet
!pip install --upgrade sagemaker datasets --quiet
```

A note below the code cell says: 'Restart the notebook kernel now after running the above cell and before you run any cells below!' The status bar at the bottom indicates the kernel is 'Fully initialized'.

This screenshot shows the JupyterLab interface with the 'Model_Evaluation.ipynb' notebook open. The left sidebar displays a file tree with 'Model_Evaluation_UdacityGenAIAWS.ipynb' selected. The main area contains a code cell with the following Python code:

```
[1]: import sagemaker_boto3, json
from sagemaker.session import Session

sagemaker_session = Session()
aws_role = sagemaker_session.get_caller_identity_arn()
aws_region = boto3.Session().region_name
sess = sagemaker.Session()
print(aws_role)
print(aws_region)
```

The code cell has been executed, and the output shows the AWS role and region:

```
sagemaker.config INFO - Not applying SDK defaults from location: /etc/dg/sagemaker/config.yaml
sagemaker.config INFO - Not applying SDK defaults from location: /home/ec2-user/.config/sagemaker/config.yaml
arn:aws:iam::611383594527:role/SageMaker-udacityProjectisSagemakerRole
us-west-2
<sagemaker.session.Session object at 0x7faccc8fc3d0>
```

Below the code cell, a section titled '2. Select Text Generation Model Meta Llama 2 7B' is visible. The status bar at the bottom indicates the kernel is 'Fully initialized'.

DeepRacer Student League | Project: Building a Domain Exp... | Quota request history | AWS S... | Notebook instances | Amazon | Model_Evaluation_UdacityGenAIAWS.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Launcher x Model_Evaluation_UdacityGe x + conda_pytorch_p310

2. Select Text Generation Model Meta Llama 2 7B

Run the next cell to set variables that contain the values of the name of the model we want to load and the version of the model.

```
[1]: (model_id, model_version) = ("meta-textgeneration-llama-2-7b", "2.*")
```

Running the next cell deploys the model. This Python code is used to deploy a machine learning model using Amazon SageMaker's JumpStart library.

1. Import the `JumpStartModel` class from the `sagemaker.jumpstart.model` module.
2. Create an instance of the `JumpStartModel` class using the `model_id` and `model_version` variables created in the previous cell. This object represents the machine learning model you want to deploy.
3. Call the `deploy` method on the `JumpStartModel` instance. This method deploys the model on Amazon SageMaker and returns a `Predictor` object.

The `Predictor` object can be used to make predictions with the deployed model. The `deploy` method will automatically choose an endpoint name, instance type, and other deployment parameters. If you want to specify these parameters, you can pass them as arguments to the `deploy` method.

The next cell will take some time to run. It is deploying a large language model, and that takes time. You'll see dashes (-) while it is being deployed. Please be patient! You'll see an exclamation point at the end of the dashes (-!) when the model is deployed and then you can continue running the next cells.

You might see a warning "For forward compatibility, pin to model_version..." You can ignore this warning, just wait for the model to deploy.

```
[2]: from sagemaker.jumpstart.model import JumpStartModel
model = JumpStartModel(model_id=model_id, model_version=model_version, instance_type="ml.g5.2xlarge")
predictor = model.deploy()
```

sagemaker.config INFO - Not applying SDK defaults from location: /etc/xdg/sagemaker/config.yaml
sagemaker.config INFO - Not applying SDK defaults from location: /home/ec2-user/.config/sagemaker/config.yaml
For forward compatibility, pin to model_version="2.*" in your JumpStartModel or JumpStartEstimator definitions. Note that major version upgrade may have different EULA acceptance terms and input/output signatures.

Simple 0 1 2 3 4 5 6 7 8 9 Fully initialized 19:25 27-09-2024

DeepRacer Student League | Project: Building a Domain Exp... | Quota request history | AWS S... | Notebook instances | Amazon | Model_Evaluation_UdacityGenAIAWS.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Launcher x Model_Evaluation_UdacityGe x + conda_pytorch_p310

```
[2]: from sagemaker.jumpstart.model import JumpStartModel
model = JumpStartModel(model_id=model_id, model_version=model_version, instance_type="ml.g5.2xlarge")
predictor = model.deploy()
```

sagemaker.config INFO - Not applying SDK defaults from location: /etc/xdg/sagemaker/config.yaml
sagemaker.config INFO - Not applying SDK defaults from location: /home/ec2-user/.config/sagemaker/config.yaml
For forward compatibility, pin to model_version="2.*" in your JumpStartModel or JumpStartEstimator definitions. Note that major version upgrade may have different EULA acceptance terms and input/output signatures.
Using vulnerable JumpStart model "meta-textgeneration-llama-2-7b" and version "2.1.8".
Using model "meta-textgeneration-llama-2-7b" with wildcard version identifier "2.*". You can pin to version "2.1.8" for more stable results.
Note that models may have different input/output signatures after a major version upgrade.

Invoke the endpoint, query and parse response

The next step is to invoke the model endpoint, send a query to the endpoint, and receive a response from the model.

Running the next cell defines a function that will be used to parse and print the response from the model.

```
[3]: def print_response(payload, response):
    print(payload["inputs"])
    print("=> ", response[0][1])
    print("-----\n")
```

The model takes a text string as input and predicts next words in the sequence, the input we send it is the prompt.

The prompt we send the model should relate to the domain we'd like to fine-tune the model on. This way we'll identify the model's domain knowledge before it's fine-tuned, and then we can run the same prompts on the fine-tuned model.

Replace "inputs" in the next cell with the input to send the model based on the domain you've chosen.

```
Simple 0 1 2 3 4 5 6 7 8 9 Fully initialized 19:26 27-09-2024
```

The screenshot shows a Jupyter Notebook interface with a sidebar containing file navigation and a search bar. The main area displays a code cell with Python code:

```
[3]: def print_response(payload, response):
    print(payload["inputs"])
    print(f"> {response[0]['generation']}")
    print("=====\n")
```

Below the code, there is explanatory text and a list of prompts for financial domain:

The model takes a text string as input and predicts next words in a sequence, the input we send it is the prompt.

The prompt we send the model should relate to the domain we'd like to fine-tune the model on. This way we'll identify the model's domain knowledge before it's fine-tuned, and then we can run the same prompts on the fine-tuned model.

Replace "inputs" in the next cell with the input to send the model based on the domain you've chosen.

For financial domain:

"inputs": "Replace with sentence below"

- "The investment tests performed indicate"
- "the relative volume for the long out of the money options, indicates"
- "The results for the short in the money options"
- "The results are encouraging for aggressive investors"

For medical domain:

"inputs": "Replace with sentence below"

- "Myeloid neoplasms and acute leukemias derive from"
- "Genomic characterization is essential for"
- "Certain germline disorders may be associated with"
- "In contrast to targeted approaches, genome-wide sequencing"

For IT domain:

"inputs": "Replace with sentence below"

```
Simple 0 1 0 Fully initialized conda_pytorch_p310 | Idle Mode: Command Ln 1, Col 1 Model_Evaluation_UdacityGenAIAWS.ipynb 0 1926 27-09-2024
```

The screenshot shows a Jupyter Notebook interface with a sidebar containing file navigation and a search bar. The main area displays a code cell with Python code:

```
[4]: payload = {
    "inputs": "outline the key aspects of ubiquitous computing from a data management perspective.",
    "parameters": {
        "max_new_tokens": 64,
        "top_p": 0.9,
        "temperature": 0.6,
        "return_full_text": False,
    },
}
try:
    response = predictor.predict(payload, custom_attributes="accept_eula=true")
    print_response(payload, response)
except Exception as e:
    print(e)
```

Below the code, there is explanatory text and a note about the prompt:

outline the key aspects of ubiquitous computing from a data management perspective.

> This paper is motivated by the observation that the data management challenges in ubiquitous computing are quite similar to those in the traditional database management. However, they are much more complex due to the diversity of the data sources, the diversity of the applications, and the diversity of the users.

=====

The prompt is related to the domain you want to fine-tune your model on. You will see the outputs from the model without fine-tuning are limited in providing insightful or relevant content.

Use the output from this notebook to fill out the "model evaluation" section of the project documentation report

Take a screenshot of this file with the cell output for your project documentation report. Download it with cell output by making sure you used Save on the notebook before downloading.

After you've filled out the report, run the cells below to delete the model deployment

```
Simple 0 1 0 Fully initialized conda_pytorch_p310 | Idle Mode: Command Ln 1, Col 1 Model_Evaluation_UdacityGenAIAWS.ipynb 0 1926 27-09-2024
```

The screenshot shows a Jupyter Notebook interface with the title 'Model_Evaluation_UdacityGenAIAWS.ipynb'. The notebook contains the following text:

```
# Delete the SageMaker endpoint and the attached resources
predictor.delete_model()
predictor.delete_endpoint()
```

Below the code, there is a note: 'Verify your model endpoint was deleted by visiting the Sagemaker dashboard and choosing endpoints under 'Inference' in the left navigation menu. If you see your endpoint still there, choose the endpoint, and then under "Actions" select Delete.'

Take a screenshot of the Model_FineTuning.ipynb file with the cell output as proof you completed this step of the project

The screenshot shows a Jupyter Notebook interface with the title 'Model_FineTuning.ipynb'. The notebook contains the following text:

Step 3: Model Fine-tuning

In this notebook, you'll fine-tune the Meta Llama 2 7B large language model, deploy the fine-tuned model, and test its text generation and domain knowledge capabilities.

Fine-tuning refers to the process of taking a pre-trained language model and retraining it for a different but related task using specific data. This approach is also known as transfer learning, which involves transferring the knowledge learned from one task to another. Large language models (LLMs) like Llama 2 7B are trained on massive amounts of unlabeled data and can be fine-tuned on domain-specific datasets, making the model perform better on that specific domain.

Input: A train and an optional validation directory. Each directory contains a CSV/JSON/TXT file. For CSV/JSON files, the train or validation data is used from the column called 'text' or the first column if no column called 'text' is found. The number of files under train and validation should equal to one.

- You'll choose your dataset below based on the domain you've chosen

Output: A trained model that can be deployed for inference.

After you've fine-tuned the model, you'll evaluate it with the same input you used in project step 2: model evaluation.

Set up

Install and import the necessary packages. Restart the kernel after executing the cell below.

```
[11]: pip install --upgrade sagemaker datasets
```

Requirement already satisfied: sagemaker in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (2.232.1)
Requirement already satisfied: datasets in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (3.0.1)
Requirement already satisfied: attrs>24,>=23.1.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker (23.2.0))
Requirement already satisfied: boto<2.0,>=1.34.142 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker (23.2.0))

A screenshot of a Jupyter Notebook interface. The notebook has two tabs open: 'Model_Evaluation_UdacityGeX' and 'Model_FineTuning.ipynb'. The 'Model_FineTuning.ipynb' tab is active, showing a code cell with the command `!pip install --upgrade sagemaker datasets`. The output of this command is a long list of dependency requirements, including packages like sagemaker, datasets, boto3, numpy, packaging, pandas, pathos, platformdirs, protobuf, psutil, pyyaml, requests, sagemaker-core, and schema. The terminal also shows the environment variable `conda_pytorch_p310`.

```
[11]: !pip install --upgrade sagemaker datasets
Requirement already satisfied: sagemaker in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (2.232.1)
Requirement already satisfied: datasets in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (3.0.1)
Requirement already satisfied: attr>=24,>=23.1.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (23.2.0)
Requirement already satisfied: boto3<2.0,>=1.34.142 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (1.35.16)
Requirement already satisfied: cloudpickle<=2.2.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (2.2.1)
Requirement already satisfied: docker in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (7.1.0)
Requirement already satisfied: google-pasta in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (0.2.0)
Requirement already satisfied: importlib-metadata<7.0,>=1.4.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (6.11.0)
Requirement already satisfied: jsonschema in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (4.23.0)
Requirement already satisfied: numpy<2.0,>=1.9.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (1.26.4)
Requirement already satisfied: packaging>=20.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (24.1)
Requirement already satisfied: pandas in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (1.5.3)
Requirement already satisfied: pathos in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (0.3.2)
Requirement already satisfied: platformdirs in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (4.2.2)
Requirement already satisfied: protobuf<5.0,>=3.12 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (4.25.4)
Requirement already satisfied: psutil in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (6.0.0)
Requirement already satisfied: pyyaml<6.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (6.0.1)
Requirement already satisfied: requests in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (2.32.3)
Requirement already satisfied: sagemaker-core<2.0.0,>=1.0.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (1.0.4)
Requirement already satisfied: schema in /home/ec2-user/anaconda3/envs/avtorch_p310/lib/python3.10/site-packages (from sagemaker) (0.7.7)
```

A screenshot of a Jupyter Notebook interface, identical to the one above, showing the execution of the same `!pip install --upgrade sagemaker datasets` command. The output shows the same list of dependency requirements, indicating that the environment is fully initialized.

```
Requirement already satisfied: tbllib4>=1.7.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (3.0.0)
Requirement already satisfied: tqdm in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (4.66.4)
Requirement already satisfied: urllib3<3.0.0,>=1.26.8 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker) (2.2.2)
Requirement already satisfied: filelock in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from datasets) (3.15.4)
Requirement already satisfied: pyarrow<15.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from datasets) (17.0.0)
Requirement already satisfied: dill<0.3.9,>=0.3.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from datasets) (0.3.8)
Requirement already satisfied: xxhash in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from datasets) (3.5.0)
Requirement already satisfied: multiprocess in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from datasets) (0.7.0_16)
Requirement already satisfied: fsspec<=2024.6.1,>=2023.1.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from fsspec[http]<=2024.6.1,>=2023.1.0>datasets) (2024.6.1)
Requirement already satisfied: aiohttp in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from datasets) (3.10.6)
Requirement already satisfied: huggingface-hub<0.22.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from datahub[http])
Requirement already satisfied: botocore<1.36.0,>=1.35.16 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from boto3<2.0,>=1.34.142>sagemaker) (1.35.16)
Requirement already satisfied: jmespath<2.0.0,>=0.7.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from boto3<2.0,>=1.34.142>sagemaker) (1.0.1)
Requirement already satisfied: s3transfer<0.11.0,>=0.10.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from boto3<2.0,>=1.34.142>sagemaker) (0.10.2)
Requirement already satisfied: aiohttpwebsockets<=2.3.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from aiohttp>datasets) (2.4.2)
Requirement already satisfied: aiosignal<1.1.2 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from aiohttp>datasets) (1.3.1)
Requirement already satisfied: frozenlist<=1.1.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from aiohttp>datasets) (1.4.1)
Requirement already satisfied: multidict<7.0,>=4.5 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from aiohttp>datasets) (6.1.0)
Requirement already satisfied: yarl<2.0,>=1.12.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from aiohttp>datasets) (1.13.0)
Requirement already satisfied: asyncio-timeout<5.0,>=4.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from aiohttp>datasets) (4.0.3)
```

A screenshot of a Jupyter Notebook interface. The main area shows a code editor with a large list of dependency requirements. The requirements list includes packages like typing-extensions, zipp, charset-normalizer, idna, certifi, pydantic, rich, mock, jsonschema, referencing, rpdsp, six, google-pasta, python-dateutil, pytz, pftt, pox, annotated-types, pydantic-core, and pydantic. The file browser on the left shows files like ITDataset.txt, Model_Evaluation_Udacit..., and Model_FineTuning.ipynb. The bottom status bar indicates the environment is conda_pytorch_p310 and the notebook is fully initialized.

```
Requirement already satisfied: typing-extensions>=3.7.4.3 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from huggingface-hub>=22.0->datasets (4.12.2))
Requirement already satisfied: zipp>=0.5 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from importlib-metadata<7.0,>=1.4.0->sagemaker (3.19.2))
Requirement already satisfied: charset-normalizer<4,>=2 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from requests->sagemaker (3.3.2))
Requirement already satisfied: idna>4,>=2.5 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from requests->sagemaker (3.7))
Requirement already satisfied: certifi>=2017.4.17 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from requests->sagemaker (2024.7.4))
Requirement already satisfied: pydantic<3.0.0,>=1.7.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker<0.0.0,>=0.0.0->sagemaker) (3.19.1)
Requirement already satisfied: rich<0.0.0,>=13.0.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker<0.0.0,>=0.0.0->sagemaker) (4.0.3)
Requirement already satisfied: mock<5.0.0,>=4.0.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from sagemaker<0.0.0,>=0.0.0->sagemaker) (4.0.3)
Requirement already satisfied: jsonschema<specifications>=2023.03.6 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from jsonschema->sagemaker) (2023.12.1)
Requirement already satisfied: referencing<0.28.4 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from jsonschema->sagemaker) (0.35.1)
Requirement already satisfied: rpdsp-py>=0.7.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from jsonschema->sagemaker) (0.19.1)
Requirement already satisfied: six in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from google-pasta->sagemaker) (1.16.0)
Requirement already satisfied: python-dateutil>=2.8.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from pandas->sagemaker) (2.9.0)
Requirement already satisfied: pytz>=2020.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from pandas->sagemaker) (2024.1)
Requirement already satisfied: pftt>=1.7.6.8 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from pathos->sagemaker) (1.7.6.8)
Requirement already satisfied: pox>=0.3.4 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from pathos->sagemaker) (0.3.4)
Requirement already satisfied: annotated-types>=0.4.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from pydantic<3.0.0,>=1.7.0->sagemaker<core>2.0.0,>=1.0.0->sagemaker) (0.7.0)
Requirement already satisfied: pydantic-core>=2.0.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages (from pydantic<3.0.0,>=1.7.0->sagemaker<core>2.0.0,>=1.0.0->sagemaker) (2.20.1)
```

A screenshot of a Jupyter Notebook interface. The code editor contains a partially filled code cell. The first few lines of the code are:

```
model_id, model_version = "meta-textgeneration-llama-2-7b", "2.0"
```

The cell below the code cell is a text input field where users can enter dataset text for domain expert models. Below this, there are sections for creating finance, medical, and IT domain expert models, each with a bulleted list of training datasets. The code editor also shows imports for JumpStartEstimator and boto3, and a call to estimator.set hyperparameters. The bottom status bar indicates the environment is conda_pytorch_p310 and the notebook is fully initialized.

```
from sagemaker.jumpstart.estimator import JumpStartEstimator
import boto3

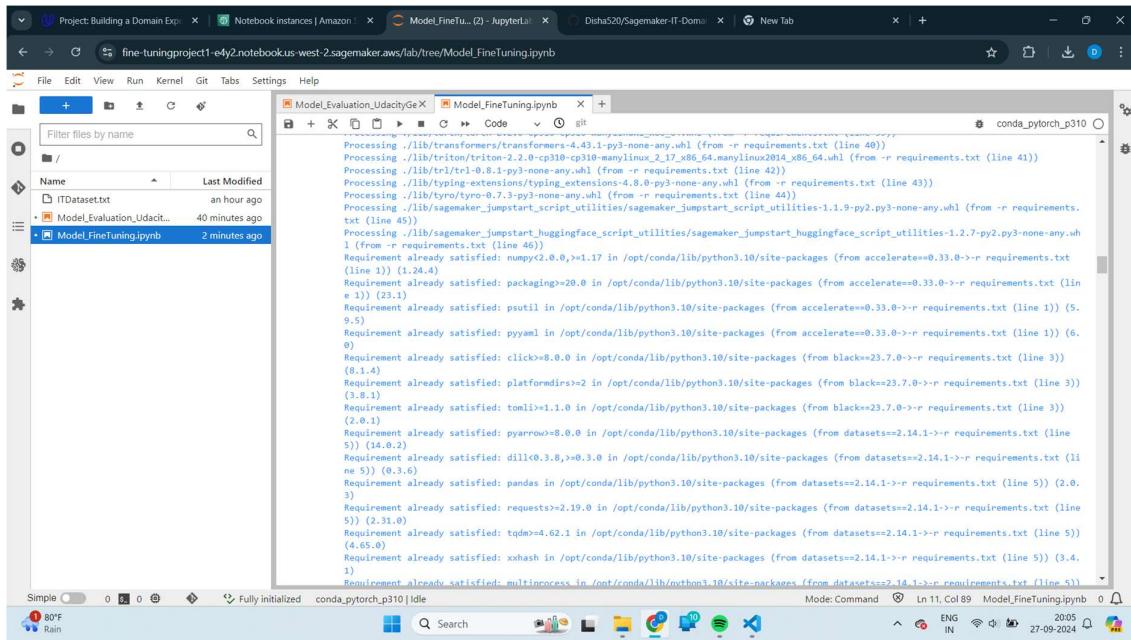
estimator = JumpStartEstimator(model_id=model_id, environments={"accept_eula": "true"}, instance_type="ml.g5.2xlarge")

estimator.set_hyperparameters(instruction_tuned=False, epochs=5)

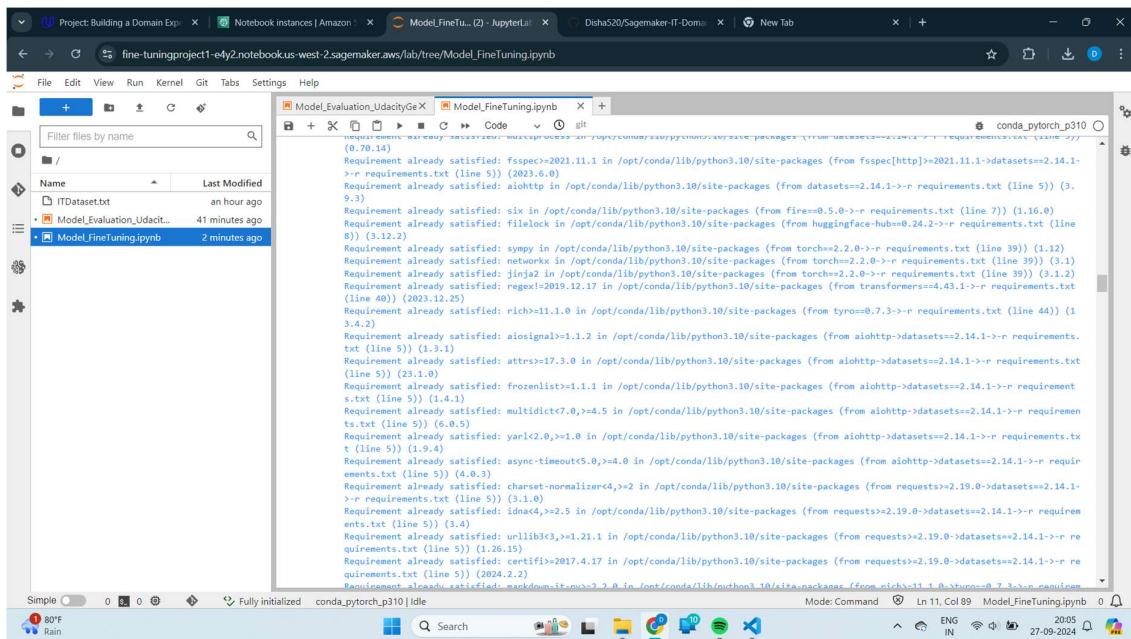
#Fill in the code below with the dataset you want to use from above.
#example: estimator.fit("training", f"s3://genaiwithawsproject2024/training/datasets/finance")
estimator.fit("training": "f3://genaiwithawsproject2024/training/datasets/it1")

INFO:sagemaker:Found credentials from IAM Role: BaseNotebookInstanceRole
INFO:sagemaker:Creating training-job with name: meta-textgeneration-llama-2-7b-2024-09-27-14-00-010
2024-09-27 14:13:01 Starting - Starting the training job
2024-09-27 14:13:01 Pending - Training job waiting for capacity...
2024-09-27 14:13:27 Pending - Preparing the instances for training...
2024-09-27 14:14:01 Downloading - Downloading input data.....
2024-09-27 14:19:25 Training - Training Image download completed. Training in progress..bash: cannot set terminal process group (-1): Inappropriate ioctl for device
last-child job contains in this shell
bash: job control not supported
2024-09-27 14:19:34,818 sagemaker-training-toolkit INFO Imported framework sagemaker.pytorch.contains.training
2024-09-27 14:19:34,839 sagemaker-training-toolkit INFO No Neurons detected (Normal) if no neurons installed
2024-09-27 14:19:34,849 sagemaker_pytorch_container.training INFO Block until all host DNS lookups succeed.
2024-09-27 14:19:34,853 sagemaker_pytorch_container.training INFO Invoking user training script.
2024-09-27 14:19:44,265 sagemaker-training-toolkit INFO Installing dependencies from requirements.txt:
/opt/conda/bin/python3.10 -m pip install -r requirements.txt
Processing /lib/accelerate/accelerate-0.33.0-py3-none-any.whl (from -r requirements.txt (line 1))
Processing /lib/awswrangler/awswrangler-2.17.0-py3-none-any.whl (from -r requirements.txt (line 2))
Processing /lib/black/black-23.7.0-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 3))
Processing /lib/brentt/brentt-1.0.9-cp310-cp310-manylinux_2_5_x86_64.manylinux1_2_5_x86_64.manylinux_2_12_x86_64.manylinux2010_x86_64.whl (from -r requirements.txt (line 4))
Processing /lib/datasets/datasets-2.14.1-py3-none-any.whl (from -r requirements.txt (line 5))
Processing /lib/doctering-parser/doctstring_parser-0.16-py3-none-any.whl (from -r requirements.txt (line 6))
Processing /lib/fire/fire-0.5.0.tar.gz
Preparing metadata (setup.py): started
Preparing metadata (setup.py): finished with status 'done'
```

```
Processing /lib/huggingface-hub/huggingface_hub-0.24.2-py3-none-any.whl (from -r requirements.txt (line 8))
Processing /lib/inflated4/inflated4-0.1-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 9))
Processing /lib/lowlevel/lowlevel-0.1.1-py3-none-any.whl (from -r requirements.txt (line 10))
Processing /lib/multivolumefile/multivolumefile-0.2.3-py3-none-any.whl (from -r requirements.txt (line 11))
Processing /lib/mypy-extensions/mypy_extensions-1.0.0-py3-none-any.whl (from -r requirements.txt (line 12))
Processing /lib/nvidia/cuda/cublas-cu12/nvidia_cublas_cu12-12.1.3-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 13))
Processing /lib/nvidia/cuda/cupti-cu12/nvidia_cuda_cupti_cu12-12.1.05-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 14))
Processing /lib/nvidia/cuda/cvrte-cu12/nvidia_cuda_nvrtc_cu12-12.1.105-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 15))
Processing /lib/nvidia/cuda/cuda-runtime-cu12/nvidia_cuda_nvrtc_cu12-12.1.105-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 16))
Processing /lib/nvidia/cudnn-cu12/nvidia_cudnn_cu12-9.2.26-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 17))
Processing /lib/nvidia/cufft-cu12/nvidia_cufft_cu12-11.0.2.54-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 18))
Processing /lib/nvidia/curand-cu12/nvidia_curand_cu12-10.3.2.106-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 19))
Processing /lib/nvidia/cusolver-cu12/nvidia_cusolver_cu12-11.4.5.107-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 20))
Processing /lib/nvidia/cusparse-cu12/nvidia_cusparse_cu12-12.1.06-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 21))
Processing /lib/nvidia/cucl-ncc1-cu12/nvidia_cucl_ncc1_cu12-19.3-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 22))
Processing /lib/nvidia/cnviflink-cu12/nvidia_nviflink_cu12-12.3.101-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 23))
Processing /lib/nvidia/nvtx-cu12/nvidia_nvtx_cu12-12.1.105-py3-none-manylinux_x86_64.whl (from -r requirements.txt (line 24))
Processing /lib/nvidia/petc-petc-0.4.0-py3-none-any.whl (from -r requirements.txt (line 25))
Processing /lib/nvidia/pycurl/pycurl-0.20.5-py3-none-any.whl (from -r requirements.txt (line 26))
Processing /lib/pycvc/pycbc-1.0.1-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 27))
Processing /lib/pycvc/pycbc-1.0.1-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 28))
Processing /lib/pycryptodome/pycryptodome-3.18.0-cp35-ab13-manylinux_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 29))
Processing /lib/pyppmd/pyppmd-1.0.0-cp310-cp310-manylinux_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 30))
Processing /lib/pystd/pystd-0.15.9-cp310-cp310-manylinux_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 31))
Processing /lib/safetensors/safetensors-0.4.2-cp310-cp310-manylinux_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 32))
Processing /lib/scipy/scipy-1.11.1-cp310-cp310-manylinux_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 33))
Processing /lib/shibak/shibak-1.7.1-py3-none-any.whl (from -r requirements.txt (line 34))
Processing /lib/tencolor/tencolor-2.3.0-py3-none-any.whl (from -r requirements.txt (line 35))
Processing /lib/texttable/texttable-1.6.7-py2,py3-none-any.whl (from -r requirements.txt (line 36))
Processing /lib/tokenize-rt/ tokenize-rt-5.1.0-py2,py3-none-any.whl (from -r requirements.txt (line 37))
Processing /lib/tokenizers/tokenizers-0.19.1-cp310-cp310-manylinux_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 38))
Processing /lib/torch/torch-2.2.0-cp310-cp310-manylinux_x86_64.whl (from -r requirements.txt (line 39))
```



```
Processing ./lib/transfomers/transfomers-4.43.1-py3-none-any.whl (from -r requirements.txt (line 40))
Processing ./lib/triton/triton-2.2.0-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (from -r requirements.txt (line 41))
Processing ./lib/typing_extensions/typing_extensions-4.8.0-py3-none-any.whl (from -r requirements.txt (line 42))
Processing ./lib/tairo/tairo-0.7.3-py3-none-any.whl (from -r requirements.txt (line 44))
Processing ./lib/sagemaker_jumpstart_s3_script_utilities/sagemaker_jumpstart_s3_script_utilities-1.1.9-py2.py3-none-any.whl (from -r requirements.txt (line 45))
Processing ./lib/sagemaker_jumpstart_huggingface_script_utilities/sagemaker_jumpstart_huggingface_script_utilities-1.2.7-py2.py3-none-any.whl (from -r requirements.txt (line 46))
Requirement already satisfied: numpy<2.0.0,>=1.17 in /opt/conda/lib/python3.10/site-packages (from accelerate==0.33.0->-r requirements.txt (line 1)) (1.24.4)
Requirement already satisfied: packaging>=20.0 in /opt/conda/lib/python3.10/site-packages (from accelerate==0.33.0->-r requirements.txt (line 1)) (23.1)
Requirement already satisfied: psutil in /opt/conda/lib/python3.10/site-packages (from accelerate==0.33.0->-r requirements.txt (line 1)) (5.9.5)
Requirement already satisfied: pyyaml in /opt/conda/lib/python3.10/site-packages (from accelerate==0.33.0->-r requirements.txt (line 1)) (6.0)
Requirement already satisfied: click>=8.0.0 in /opt/conda/lib/python3.10/site-packages (from black==23.7.0->-r requirements.txt (line 3)) (8.1.4)
Requirement already satisfied: platformdirs>=2 in /opt/conda/lib/python3.10/site-packages (from black==23.7.0->-r requirements.txt (line 3)) (3.8.1)
Requirement already satisfied: tomlii>=1.1.0 in /opt/conda/lib/python3.10/site-packages (from black==23.7.0->-r requirements.txt (line 3)) (2.0.1)
Requirement already satisfied: pyarrow>=8.0.0 in /opt/conda/lib/python3.10/site-packages (from datasets==2.14.1->-r requirements.txt (line 5)) (14.0.2)
Requirement already satisfied: dill<0.3.8,>=0.3.0 in /opt/conda/lib/python3.10/site-packages (from datasets==2.14.1->-r requirements.txt (line 5)) (0.3.6)
Requirement already satisfied: pandas in /opt/conda/lib/python3.10/site-packages (from datasets==2.14.1->-r requirements.txt (line 5)) (2.0.3)
Requirement already satisfied: requests>=2.19.0 in /opt/conda/lib/python3.10/site-packages (from datasets==2.14.1->-r requirements.txt (line 5)) (2.31.0)
Requirement already satisfied: tqdm>=46.2.1 in /opt/conda/lib/python3.10/site-packages (from datasets==2.14.1->-r requirements.txt (line 5)) (46.5.0)
Requirement already satisfied: xxhash in /opt/conda/lib/python3.10/site-packages (from datasets==2.14.1->-r requirements.txt (line 5)) (3.4.1)
Requirement already satisfied: multineuronc_in /opt/conda/lib/python3.10/site-packages (from datasets==2.14.1->-r requirements.txt (line 5))
```



```
Requirement already satisfied: fsspec>=2021.11.1 in /opt/conda/lib/python3.10/site-packages (from fsspec[http]>=2021.11.1->datasets==2.14.1->r requirements.txt (line 5)) (0.70.14)
Requirement already satisfied: aiohttp in /opt/conda/lib/python3.10/site-packages (from fsspec[http]>=2021.11.1->datasets==2.14.1->r requirements.txt (line 5)) (2023.6.0)
Requirement already satisfied: aiotoys in /opt/conda/lib/python3.10/site-packages (from datasets==2.14.1->-r requirements.txt (line 5)) (0.3.1)
Requirement already satisfied: six in /opt/conda/lib/python3.10/site-packages (from fire<0.5.0,>=r requirements.txt (line 7)) (1.16.0)
Requirement already satisfied: filelock in /opt/conda/lib/python3.10/site-packages (from huggingface-hub==0.24.2->-r requirements.txt (line 8)) (3.12.2)
Requirement already satisfied: sympy in /opt/conda/lib/python3.10/site-packages (from torch<2.2.0,>=r requirements.txt (line 39)) (1.12)
Requirement already satisfied: networks in /opt/conda/lib/python3.10/site-packages (from torch<2.2.0,>=r requirements.txt (line 39)) (3.1)
Requirement already satisfied: jinja2 in /opt/conda/lib/python3.10/site-packages (from torch<2.2.0,>=r requirements.txt (line 39)) (3.1.2)
Requirement already satisfied: regex<2019.12.17 in /opt/conda/lib/python3.10/site-packages (from transformers==4.43.1->-r requirements.txt (line 40)) (2023.12.25)
Requirement already satisfied: rich>=11.1.0 in /opt/conda/lib/python3.10/site-packages (from tyro==0.7.3->-r requirements.txt (line 44)) (13.4.2)
Requirement already satisfied: aiosignal>=1.1.2 in /opt/conda/lib/python3.10/site-packages (from aiohttp->datasets==2.14.1->-r requirements.txt (line 5)) (1.3.1)
Requirement already satisfied: attrs>=17.3.0 in /opt/conda/lib/python3.10/site-packages (from aiohttp->datasets==2.14.1->-r requirements.txt (line 5)) (23.1.0)
Requirement already satisfied: frozenlist>=1.1 in /opt/conda/lib/python3.10/site-packages (from aiohttp->datasets==2.14.1->-r requirements.txt (line 5)) (1.4.1)
Requirement already satisfied: multidict>7.0,>=4.5 in /opt/conda/lib/python3.10/site-packages (from aiohttp->datasets==2.14.1->-r requirements.txt (line 5)) (6.0.5)
Requirement already satisfied: yarl<2.0,>=1.0 in /opt/conda/lib/python3.10/site-packages (from aiohttp->datasets==2.14.1->-r requirements.txt (line 5)) (1.9.4)
Requirement already satisfied: async-timeout<5.0,>=4.0 in /opt/conda/lib/python3.10/site-packages (from aiohttp->datasets==2.14.1->-r requirements.txt (line 5)) (4.0.3)
Requirement already satisfied: charset-normalizer<4,>=2 in /opt/conda/lib/python3.10/site-packages (from requests>=2.19.0->datasets==2.14.1->-r requirements.txt (line 5)) (3.1.0)
Requirement already satisfied: idna>=2.5 in /opt/conda/lib/python3.10/site-packages (from requests>=2.19.0->datasets==2.14.1->-r requirements.txt (line 5)) (3.2.0)
Requirement already satisfied: urllib3<3,>=1.21.1 in /opt/conda/lib/python3.10/site-packages (from requests>=2.19.0->datasets==2.14.1->-r requirements.txt (line 5)) (1.26.15)
Requirement already satisfied: certifi>=2017.4.17 in /opt/conda/lib/python3.10/site-packages (from requests>=2.19.0->datasets==2.14.1->-r requirements.txt (line 5)) (2024.2.2)
Requirement already satisfied: markdown_it<0.9.2.0 in /opt/conda/lib/python3.10/site-packages (from rich>=11.1.0->tyro==0.7.3->-r requirements.txt (line 5)) (0.9.0)
```

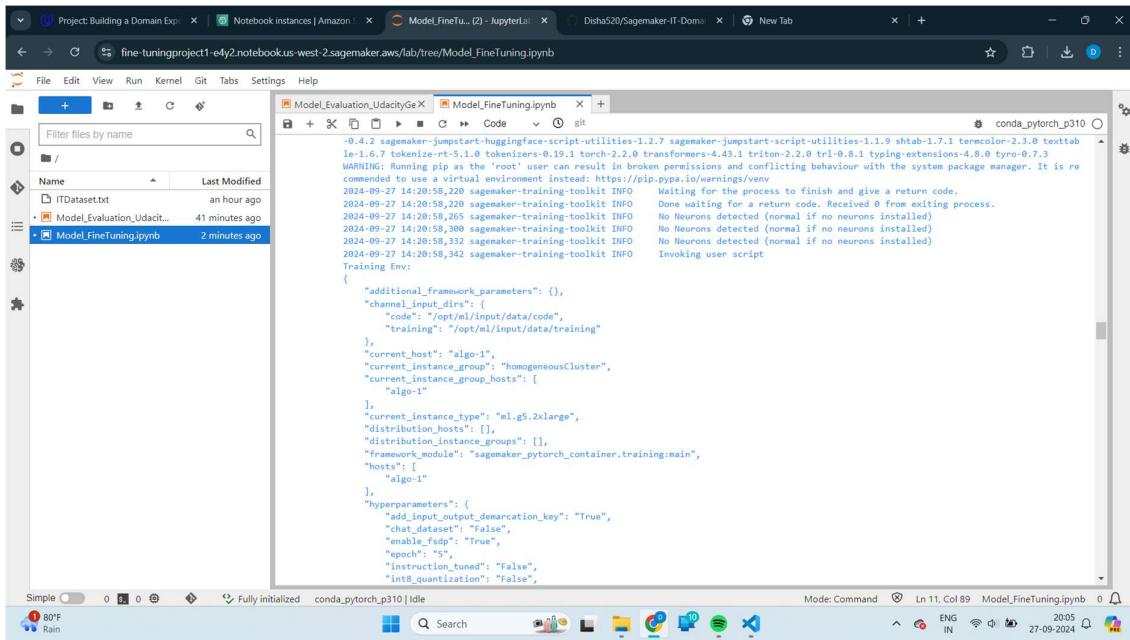
A screenshot of a Jupyter Notebook interface on a Windows desktop. The main window shows a file tree with 'Model_Evaluation_Udacity.ipynb' selected. Below it is a terminal window titled 'conda_pytorch_p310' with the following command history:

```
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /opt/conda/lib/python3.10/site-packages (from rich>=11.1.0->tyro==0.7.3->r requirements.txt [line 44]) (2.15.1)
Requirement already satisfied: MarkupSafe>=2.0 in /opt/conda/lib/python3.10/site-packages (from jinja2>torch>2.2.0->r requirements.txt [line 39]) (2.1.3)
Requirement already satisfied: python-dateutil>=2.8.2 in /opt/conda/lib/python3.10/site-packages (from pandas>datasets==2.14.1->r requirements.txt [line 55]) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /opt/conda/lib/python3.10/site-packages (from pandas>datasets==2.14.1->r requirements.txt [line 55]) (2023.3)
Requirement already satisfied: mpmath>=0.19 in /opt/conda/lib/python3.10/site-packages (from sympy>torch>2.2.0->r requirements.txt [line 39]) (1.3.0)
Requirement already satisfied: mdurl>=0.1 in /opt/conda/lib/python3.10/site-packages (from markdown-it-py>=2.2.0->r rich>=11.1.0->tyro==0.7.3->r requirements.txt [line 44]) (0.1.0)
scipy is already installed with the same version as the provided wheel. Use --force-reinstall to force an installation of the wheel.
Building wheels for collected packages: fire
Building wheel for fire (setup.py): started
Building wheel for fire (setup.py): finished with status 'done'
Created wheel for fire: filename=fire-0.5.0-py2.py3-none-any.whl size=116932 sha256=8bbda492975dea7ad36130349e5e28a6a1b2a16169c2e3797b0461dae594479
Stored in directory: /root/.cache/pip/wheels/db/3d/41/7e69da5f61e37d109a4457082ffcc6ed85ab63bbafded38
Successfully built fire
Installing collected packages: Brotli, bitsandbytes, typing-extensions, triton, tokenizers, termcolor, shtab, sagemaker-jumpstart-t-script-utilities, sagemaker-jumpstart-huggingface-script-utilities, safetensors, pyzstd, pycryptodomex, pycbc, pathspec, nvidia-nvtx-cu12, nvidia-nvjlalink-cu12, nvidia-curand-cu12, nvidia-cufft-cu12, nvidia-cuda-runtime-cu12, nvidia-cuda-nvrtc-cu12, nvidia-cuda-cupti-cu12, nvidia-cublas-cu12, mypy-extensions, multivolumefile, lorainb, inflate64, docstring-parser, py7zr, nvidia-cusparse-cu12, nvidia-cudnn-cu12, huggingface-hub, fire, black, tyro, tokenizers, nvidia-cusolver-cu12, transformers, torch, datasets, accelerate, tr1, pefc
Attempting uninstall: typing-extensions
Found existing installation: typing_extensions 4.7.1
Uninstalling typing_extensions-4.7.1:
Successfully uninstalled typing_extensions-4.7.1
Attempting uninstall: triton
Found existing installation: triton 2.0.0.dev20221202:
Uninstalling triton 2.0.0.dev20221202:

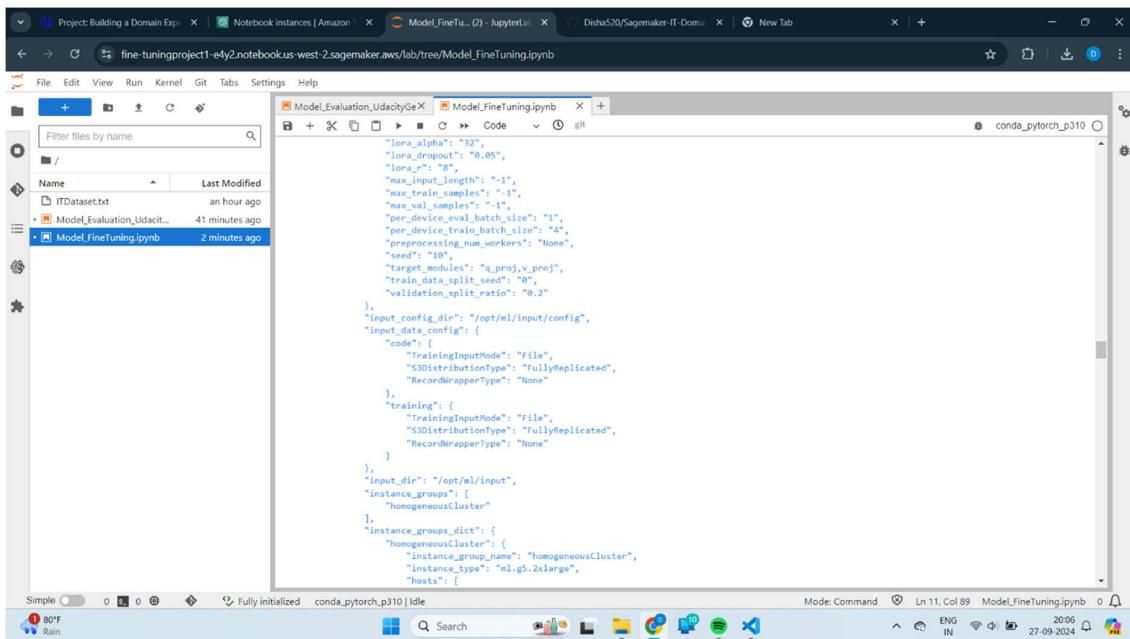
```

A screenshot of a Jupyter Notebook interface on a Windows desktop, similar to the one above. The main window shows a file tree with 'Model_Evaluation_Udacity.ipynb' selected. Below it is a terminal window titled 'conda_pytorch_p310' with the following command history:

```
Uninstalling triton 2.0.0.dev20221202:
Successfully uninstalled triton-2.0.0.dev20221202
Attempting uninstall: huggingface-hub
Found existing installation: huggingface-hub 0.20.3
Uninstalling huggingface-hub-0.20.3:
Successfully uninstalled huggingface-hub-0.20.3
Attempting uninstall: tokenizers
Found existing installation: tokenizers 0.13.3
Uninstalling tokenizers-0.13.3:
Successfully uninstalled tokenizers-0.13.3
Attempting uninstall: transformers
Found existing installation: transformers 4.28.1
Uninstalling transformers-4.28.1:
Successfully uninstalled transformers-4.28.1
Attempting uninstall: torch
Found existing installation: torch 2.0.0
Uninstalling torch-2.0.0:
Successfully uninstalled torch-2.0.0
Attempting uninstall: datasets
Found existing installation: datasets 2.16.1
Uninstalling datasets-2.16.1:
Successfully uninstalled datasets-2.16.1
Attempting uninstall: accelerate
Found existing installation: accelerate 0.19.0
Uninstalling accelerate-0.19.0:
Successfully uninstalled accelerate-0.19.0
Pip's global dependency conflict detection does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency conflicts:
fastai 2.7.12 requires torch<2.1,>=1.7, but you have torch 2.2.0 which is incompatible.
Successfully installed Brotli-1.0.9 accelerate-0.33.0 bitsandbytes-0.39.1 black-23.7.0 datasets-2.14.1 docstring-parser-0.16 fire-0.5.0 huggingface-hub-0.24.2 inflate64-0.3.1 lorainb-0.1 multivolumefile-0.2.3 mypy-extensions-1.0.0 nvidia-cublas-cu12-12.1.3.1 nvidia-cuda-cupti-cu12-12.1.105 nvidia-cuda-nvrtc-cu12-12.1.105 nvidia-cuda-runtime-cu12-12.1.105 nvidia-cudnn-cu12-8.9.2_26 nvidia-cufft-cu12-11.0.2.54 nvidia-curand-cu12-10.3_2.106 nvidia-cusolver-cu12-11.4.5.107 nvidia-cusparse-cu12-12.1.0.106 nvidia-nccl-cu12-2.19.3 nvidia-nvjlalink-cu12-12.3.1 @nvidia-nvtx-cu12-12.1.105 pathspec-0.11.1 peft-0.4.0 py7zr-0.20.5 pycbc-1.0.1 pycryptodomex-3.18.0 pypmd-1.0.0 pyzstd-0.15.9 safetensors-0.4.2 sagemaker-jumpstart-huggingface-script-utilities-1.2.7 sagemaker-jumpstart-script-utilities-1.1.9 shtab-1.7.1 termcolor-2.3.0 texttab-1e-1.6.7 tokenize-rt-5.1.0 tokenizers-0.19.1 torch-2.2.0 transformers-4.43.1 triton-2.2.0 tr1-0.8.1 typing-extensions-4.8.0 tyro-0.7.3
```



```
-0.4.2 sagemaker-jumpstart-huggingface-script-utilities-1.2.7 sagemaker-jumpstart-script-utilities-1.1.9 shtab-1.7.1 termcolor-2.3.0 texttab-1e-4.6.7 tokenize-5.1.0 tokenizers-0.19.1 torch-2.2.0 transformers-4.43.1 triton-2.2.0 trl-0.8.1 typing-extensions-4.8.0 tyro-0.7.3
WARNING: Running pip as the "root" user can result in broken permissions and conflicting behaviour with the system package manager. It is recommended to run pip as a regular user.
2024-09-27 14:20:58,229 sagemaker-training-toolkit INFO Waiting for the process to finish and give a return code.
2024-09-27 14:20:58,230 sagemaker-training-toolkit INFO Done waiting for a return code. Received 0 from exiting process.
2024-09-27 14:20:58,265 sagemaker-training-toolkit INFO No Neurons detected (normal if no neurons installed)
2024-09-27 14:20:58,300 sagemaker-training-toolkit INFO No Neurons detected (normal if no neurons installed)
2024-09-27 14:20:58,332 sagemaker-training-toolkit INFO No Neurons detected (normal if no neurons installed)
2024-09-27 14:20:58,342 sagemaker-training-toolkit INFO Invoking user script
Training Env
{
    "additional_framework_parameters": {},
    "channel_input_dirs": {
        "code": "/opt/ml/input/data/code",
        "training": "/opt/ml/input/data/training"
    },
    "current_host": "algo-1",
    "current_instance_group": "homogeneousCluster",
    "current_instance_group_hosts": [
        "algo-1"
    ],
    "current_instance_type": "ml.g5.2xlarge",
    "distribution_hosts": [],
    "distribution_instance_groups": [],
    "framework_module": "sagemaker_pytorch_container.training:main",
    "hosts": [
        "algo-1"
    ],
    "hyperparameters": {
        "add_end_of_output_decimation_key": "True",
        "chat_dataset": "False",
        "enable_fdp": "True",
        "epoch": "5",
        "instruction_tuned": "False",
        "int8_quantization": "False"
    }
}
```



```
"lora_alpha": "32",
"lora_dropout": "0.05",
"lora_r": "8",
"max_input_length": "-1",
"max_train_samples": "-1",
"max_val_samples": "-1",
"per_device_eval_batch_size": "1",
"per_device_train_batch_size": "4",
"preprocessing_num_workers": "None",
"seed": "42",
"target_modules": "proj, x_proj",
"train_data_split_seed": "0",
"validation_split_ratio": "0.2"
},
"input_config_dir": "/opt/ml/input/config",
"input_data_config": {
    "code": {
        "TrainingInputMode": "File",
        "S3DistributionType": "FullyReplicated",
        "RecordWrapperType": "None"
    },
    "training": {
        "TrainingInputMode": "File",
        "S3DistributionType": "FullyReplicated",
        "RecordWrapperType": "None"
    }
},
"input_dir": "/opt/ml/input",
"instance_groups": [
    "homogeneousCluster"
],
"instance_groups_dict": {
    "homogeneousCluster": {
        "instance_group_name": "homogeneousCluster",
        "instance_type": "ml.g5.2xlarge",
        "hosts": [
            "algo-1"
        ]
    }
}
```

```
Model_Evaluation_UdacityGeV Model_FineTuning.ipynb
{
    "hosts": [
        "algo-1"
    ],
    "is_hetero": false,
    "is_master": true,
    "is_modelparallel_enabled": null,
    "is_smndpmpun_installed": true,
    "job_name": "meta-textgeneration-llama-2-7b-2024-09-27-14-13-00-010",
    "log_level": 20,
    "master_hostname": "algo-1",
    "model_dir": "/opt/ml/input/code/sourcedir.tar.gz",
    "module_name": "transfer_learning",
    "network_interface_name": "eth0",
    "num_cpus": 8,
    "num_gpus": 1,
    "num_neurons": 0,
    "output_data_dir": "/opt/ml/output/data",
    "output_dir": "/opt/ml/output",
    "output_intermediate_dir": "/opt/ml/output/intermediate",
    "resource_config": {
        "current_host": "algo-1",
        "current_instance_type": "ml.g5.2xlarge",
        "current_group_name": "homogeneousCluster",
        "hosts": [
            "algo-1"
        ],
        "instance_groups": [
            {
                "instance_group_name": "homogeneousCluster",
                "instance_type": "ml.g5.2xlarge",
                "hosts": [
                    "algo-1"
                ]
            }
        ]
    },
    "user_entry_point": "transfer_learning.py"
}
```

```
Model_Evaluation_UdacityGeV Model_FineTuning.ipynb
{
    "hosts": [
        "algo-1"
    ],
    "network_interface_name": "eth0",
    "user_entry_point": "transfer_learning.py"
}
Environment variables:
SM_HOSTS=['algo-1']
SM_NETWORK_INTERFACE_NAME=eth0
SM_HPS={"add_input_output_deactivation_key": "True", "chat_dataset": "False", "enable_fadv": "True", "epoch": "5", "instruction_tuned": "False", "int8_quantization": "False", "learning_rate": "0.0001", "lora_alpha": "32", "lora_dropout": "0.05", "lora_r": "8", "max_input_length": "1", "max_train_samples": "-1", "max_val_samples": "-1", "per_device_eval_batch_size": "1", "per_device_train_batch_size": "4", "preprocessing_num_workers": "None", "seed": "10", "target_modules": "'q_proj,v_proj'", "train_data_split_seed": "0", "validation_split_ratio": "0.2"}
SM_USER_ENTRY_POINT=transfer_learning.py
SM_FRAMEWORK_PARAMS={}
SM_INSTANCE_CONFIG=[{"current_group_name": "homogeneousCluster", "current_host": "algo-1", "current_instance_type": "ml.g5.2xlarge", "hosts": ["algo-1"], "instance_group": [{"hosts": ["algo-1"], "instance_group_name": "homogeneousCluster", "instance_type": "ml.g5.2xlarge"}]}, {"network_interface_name": "eth0"}]
SM_INPUT_DATA_CONFIG={"code": "RecordWriterType": "None", "S3DistributionType": "FullyReplicated", "TrainingInputMode": "File"}, {"training": {"RecordWriterType": "None", "S3DistributionType": "FullyReplicated", "TrainingInputMode": "File"}}
SM_OUTPUT_DATA_DIR=/opt/ml/output/data
SM_CHANNELS=["code", "training"]
SM_CURRENT_HOST=algo-1
SM_CURRENT_INSTANCE_TYPE=ml.g5.2xlarge
SM_CURRENT_INSTANCE_GROUP=homogeneousCluster
SM_CURRENT_INSTANCE_GROUP_HOSTS=['algo-1']
SM_INSTANCE_GROUPS=['homogeneousCluster']
SM_INSTANCE_GROUPS_DICT={("homogeneousCluster": {"hosts": ["algo-1"], "instance_group_name": "homogeneousCluster", "instance_type": "ml.g5.2xlarge"})}
SM DISTRIBUTION_INSTANCE_GROUPS=[]
SM_IS_HETERO=False
SM_MODULE_NAME=transfer_learning
SM_LOG_LEVEL=20
SM_MODEL_MODULES=sagemaker_pytorch_container.training:main
SM_INPUT_DIR=/opt/ml/input
```

Project: Building a Domain Expertise Model | Notebook instances | Amazon Sagemaker | Model_FineTuning.ipynb | Disha520/Sagemaker-IT-Domains | New Tab

fine-tuningproject1-e4y2.notebook.us-west-2.sagemaker.aws/lab/tree/Model_FineTuning.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Model_Evaluation_UdacityGeX Model_FineTuning.ipynb

SM_IS_HETERO=False
SM_MODULE_NAME=transfer_learning
SM_LOG_LEVEL=2
SM_FRAMEWORK_MODULE=sagemaker_pytorch_container.training:main
SM_INPUT_DIR=/opt/ml/input
SM_INPUT_CONFIG_DIR=/opt/ml/input/config
SM_OUTPUT_DIR=/opt/ml/output
SM_NUM_Cpus=8
SM_NUM_Gpus=1
SM_NUM_Neurons=0
SM_MODEL_DIR=/opt/ml/model
SM_MODULE_DIR=/opt/ml/input/data/code/sourcedir.tar.gz
SM_TRAINING_ENV={"additional_framework_parameters":{}, "channel_input_dirs": {"code": "/opt/ml/input/data/code", "training": "/opt/ml/input/data/training"}, "current_host": "algo-1", "current_instance_group": "homogeneouscluster", "current_instance_type": "ml.g5.2xlarge", "data": "/opt/ml/input/data/training", "data_type": "text", "data_dir": "/opt/ml/input/data/training", "decomposers": [{"host": "algo-1", "model": "algo-1", "module": "transfer_learning", "type": "algo-1"}, {"host": "algo-1", "model": "algo-1", "module": "lora", "type": "algo-1"}], "hosts": [{"algo-1": "algo-1"}], "hyperparameters": {"add_input_output_demarcation_key": "False", "chat_dataset": "False", "enable_fsdp": "True", "epoch": "5", "instruction_tuned": "False", "int8_quantization": "False", "learning_rate": "0.0001", "lora_alpha": "32", "lora_dropout": "0.05", "lora_r": "8", "max_input_length": "1", "max_train_samples": "1", "max_val_samples": "1", "per_device_eval_batch_size": "1", "per_device_train_batch_size": "4", "preprocessing_num_workers": "None", "seed": "10", "target_modules": "q_proj,v_proj", "train_data_split_seed": "0", "validation_split_ratio": "0.2"}, "preprocesssing_num_workers": "None", "seed": "10", "target_modules": "q_proj,v_proj", "train_data_split_seed": "0", "validation_split_ratio": "0.2"}, "input_config_dir": "/opt/ml/input/config", "input_data_config": {"code": "RecordDumperType": "None", "S3DistributionType": "FullyReplicated", "TrainingInputMode": "File", "training": {"RecordDumperType": "None", "S3DistributionType": "FullyReplicated", "TrainingInputMode": "File"}, "input_dir": "/opt/ml/input", "instance_groups": [{"homogeneouscluster": {"hosts": [{"algo-1": "algo-1"}]}}, {"instance_group_name": "homogeneouscluster", "instance_type": "ml.g5.2xlarge"}], "instance_group_type": "MasterSlave", "is_mpiparallel_enabled": "True", "is_standalone": "True", "is_standalone_installer": "True", "job_name": "meta-training", "log_level": "INFO", "log_subdir": "algo-1", "model_dir": "/opt/ml/model", "module_dir": "/opt/ml/input/data/code/sourcedir.tar.gz", "module_name": "transfer_learning", "network_interface_name": "eth0", "num_cpus": "8", "num_gpus": "1", "num_neurons": "0", "output_data_dir": "/opt/ml/output/data", "output_dir": "/opt/ml/output", "output_intermediate_dir": "/opt/ml/output/intermediate", "resource_config": {"current_group_name": "homogeneouscluster", "current_host": "algo-1", "current_instance_type": "ml.g5.2xlarge", "instance_group": "algo-1", "instance_group_name": "homogeneouscluster", "instance_type": "ml.g5.2xlarge"}, "network_interface_name": "eth0", "user_entry_point": "transfer_learning.py"}, "SM_USER_ARGS="--add_input_output_demarcation_key=True --chat_dataset=False --enable_fsdp=True --epoch=5 --instruction_tuned=False --int8_quantization=False --learning_rate=0.0001 --lora_alpha=32 --lora_dropout=0.05 --lora_r=8 --max_input_length=1 --max_train_samples=1 --max_val_samples=1 --per_device_eval_batch_size=1 --per_device_train_batch_size=4 --preprocessing_num_workers=None --seed=10 --target_modules=q_proj,v_proj --train_data_split_seed=0 --validation_split_ratio=0.2", "SM_OUTPUT_INTERMEDIATE_DIR=/opt/ml/output/intermediate
SM_CHANNEL_CODE=/opt/ml/input/data/code

conda_pytorch_p310 | Idle

80°F Rain

Fully initialized

Mode: Command

Ln 11, Col 89 Model_FineTuning.ipynb

27-09-2024

Project: Building a Domain Expertise Model | Notebook instances | Amazon Sagemaker | Model_FineTuning.ipynb | Disha520/Sagemaker-IT-Domains | New Tab

fine-tuningproject1-e4y2.notebook.us-west-2.sagemaker.aws/lab/tree/Model_FineTuning.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Model_Evaluation_UdacityGeX Model_FineTuning.ipynb

SM_CHANNEL_CODE=/opt/ml/input/data/code
SM_CHANNELS_ENCODING=/opt/ml/input/data/training
SM_INPUT_DIR=/opt/ml/input
SM_INPUT_CONFIG_DIR=/opt/ml/input/config
SM_HP_CHAT_DATASET=True
SM_HP_ENABLE_FSDP=True
SM_HP_EPOCHS=5
SM_HP_INSTRUCTION_TUNED=False
SM_HP_INT8_QUANTIZATION=False
SM_HP_LEARNING_RATE=0.0001
SM_HP_LORA_ALPHA=32
SM_HP_MP_INPUT_OUTPUT=0.05
SM_HP_LORA_R=8
SM_HP_MAX_INPUT_LENGTH=1
SM_HP_MAX_TRAIN_SAMPLES=1
SM_HP_MAX_VAL_SAMPLES=1
SM_HP_PER_DEVICE_EVAL_BATCH_SIZE=1
SM_HP_PER_DEVICE_TRAIN_BATCH_SIZE=4
SM_HP_PREFPROCESSING_NUM_WORKERS=None
SM_HP_SEED=42
SM_HP_TARGET_MODULES=q_proj,v_proj
SM_HP_TRAIN_DATA_SPLIT_SEED=0
SM_HP_VALIDATION_SPLIT_RATIO=0.2
PYTHONPATH=/opt/ml/code:/opt/conda/bin:/opt/conda/lib/python3.10:/opt/conda/lib/python3.10/lib-dynload:/opt/conda/lib/python3.10/site-packages
Invoking script with the following command:
/opt/conda/bin/python3.10 transfer_learning.py --add_input_output_demarcation_key True --chat_dataset False --enable_fsdp True --epoch 5 --instruction_tuned False --int8_quantization False --learning_rate 0.0001 --lora_alpha 32 --lora_dropout 0.05 --lora_r 8 --max_input_length 1 --max_train_samples 1 --max_val_samples 1 --per_device_eval_batch_size 1 --per_device_train_batch_size 4 --preprocessing_num_workers None --seed 10 --target_modules q_proj,v_proj --train_data_split_seed 0 --validation_split_ratio 0.2
2024-09-27 14:20:58,384 sagemaker-training-toolkit INFO Exceptions not imported for SageMaker TF as Tensorflow is not installed.

***** BUG REPORT *****
Welcome to bitsandbytes. For bug reports, please run
python -m bitsandbytes
and submit this information together with your error trace to: <https://github.com/TimDettmers/bitsandbytes/issues>

bin /opt/conda/lib/python3.10/site-packages/bitsandbytes/libbitsandbytes.cuda118.so

conda_pytorch_p310 | Idle

80°F Rain

Fully initialized

Mode: Command

Ln 11, Col 89 Model_FineTuning.ipynb

27-09-2024

Project: Building a Domain Expert Model

fine-tuningproject1-e4y2.notebook.us-west-2.sagemaker.aws/lab/tree/Model_FineTuning.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Model_Evaluation_UdacityGe... Model_FineTuning.ipynb + git

Filter files by name /

Name Last Modified

- ITDataset.txt an hour ago
- Model_Evaluation_Udacit... 41 minutes ago
- Model_FineTuning.ipynb 2 minutes ago

bin /opt/conda/lib/python3.10/site-packages/bitsandbytes/libbitsandbytes_cudal18.so
/opt/conda/lib/python3.10/site-packages/bitsandbytes/cuda_setup/main.py:149: UserWarning: WARNING: The following directories listed in your path were found to be non-existent: {PosixPath('/usr/local/nvidia/lib64'), PosixPath('/usr/local/nvidia/lib')}

warn(msg)

CUDA SETUP: CUDA runtime path found: /opt/conda/lib/libcudart.so.11.0

CUDA SETUP: Highest compute capability among GPUs detected: 8.6

CUDA SETUP: Detected CUDA version 118

CUDA SETUP: Loading binary /opt/conda/lib/python3.10/site-packages/bitsandbytes/libbitsandbytes_cudal18.so...

INFO:root:Using pre-trained artifacts in SAGEMAKER_ADDITIONAL_S3_DATA_PATH:/opt/ml/outputs/1

INFO:root:Identify file serving_proportions in the un-tar directory /opt/ml/outputs/1. Copying it over to /opt/ml/model for model deployment after training is finished.

INFO:root:Invoking the training command ["torchrun", "--nproc_per_node", "1", "llama_finetuning.py", "--model_name", "/opt/ml/outputs/1", "--dataset", "ITDataset", "--batch_size_training", "4", "--train_file", "/opt/ml/input/data/training", "--lr", "0.0001", "--do_train", "--output_dir", "saved_peft_model", "--num_epochs", "5", "--use_peft", "--peft_method", "lora", "--max_train_samples", "1", "--max_val_samples", "1", "--seed", "0", "--train_data_eval_batch_size", "1", "--max_input_length", "1", "--preprocessing_num_workers", "None", "--validation_split_ratio", "0.2", "--train_data_split_seed", "0", "--num_workers_dataloader", "0", "--weight_decay", "0.1", "--lora_r", "8", "--lora_alpha", "32", "--lora_dropout", "0.05", "--target_modules", "q_proj,v_proj", "--chat_template", "None", "--enable_fsdp", "--input_output_deprecation_key"]

=====BLUR REPORT=====

Welcome to bitsandbytes. For bug reports, please run

python -m bitsandbytes

and submit this information together with your error trace to: <https://github.com/TimDettmers/bitsandbytes/issues>

bin /opt/conda/lib/python3.10/site-packages/bitsandbytes/libbitsandbytes_cudal18.so
/opt/conda/lib/python3.10/site-packages/bitsandbytes/cuda_setup/main.py:149: UserWarning: WARNING: The following directories listed in your path were found to be non-existent: {PosixPath('/usr/local/nvidia/lib64'), PosixPath('/usr/local/nvidia/lib')}

warn(msg)

CUDA SETUP: CUDA runtime path found: /opt/conda/lib/libcudart.so.11.0

CUDA SETUP: Highest compute capability among GPUs detected: 8.6

CUDA SETUP: Detected CUDA version 118

CUDA SETUP: Loading binary /opt/conda/lib/python3.10/site-packages/bitsandbytes/libbitsandbytes_cudal18.so...

INFO:root:Local rank is 0. Rank is 0

INFO:root:--> Using torch dist debug set to detail

INFO:root:loading the tokenizer.

Simple 0 0 0 0 Fully initialized conda_pytorch_p310 idle Mode: Command Ln 11, Col 89 Model_FineTuning.ipynb 0

80°F Rain ENG IN 27-09-2024 2006

```
algo-1:57:79 [0] ncc1_net_ofi_init:1444 NCC1_WARN NET/OFI Only EFA provider is supported
algo-1:57:79 [0] ncc1_net_ofi_init:1483 NCC1_WARN NET/OFI aws-ofi-ncc1 initialization failed
step 0 is completed and loss is 3.955679178237915
Training Epoch0: 100%[#033|3dm] 0/1 [00:00<0:00, 5.07it/s]
Training Epoch0: 100%[#033|3dm] 0/1 [00:00<0:00, 5.07it/s]
Max CUDA memory allocated was 15 GB
Max CUDA memory reserved was 15 GB
Peak active CUDA memory was 15 GB
Cuda Malloc retires : 0
CPU Total Peak Memory consumed during the train (max): 1 GB
evaluating Epoch: 0%[#033|3dm] 0/0 [00:00<?, ?it/s]
We detected that you're passing `past_key_values` as a tuple and this is deprecated and will be removed in v4.43. Please use an appropriate 'Cache' class (https://huggingface.co/docs/transformers/v4.43/en/internal/generation\_utils#transformers.Cache)
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.09it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.09it/s]
eval_pl-tensor(45.7752, device='cuda:0') eval_epoch_loss=tensor(3.8718, device='cuda:0')
we are about to save the PEFT modules
PEFT modules are saved in saved_peft_model directory
best eval loss on epoch 0 is 3.8717923164367676
Epoch 1: train_perplexity=52.2312, train_epoch_loss=3.9557, epoch time 5.360014476999936s
Training Epoch1: 0%[#033|3dm] 0/0 [00:00<?, ?it/s]
step 0 is completed and loss is 3.9874656884277344
Training Epoch1: 100%[#033|3dm] 0/1 [00:03<0:00, 3.46it/s]
Training Epoch1: 100%[#033|3dm] 0/1 [00:03<0:00, 3.46it/s]
Max CUDA memory allocated was 15 GB
Max CUDA memory reserved was 15 GB
Peak active CUDA memory was 15 GB
Cuda Malloc retires : 64
CPU Total Peak Memory consumed during the train (max): 2 GB
evaluating Epoch: 0%[#033|3dm] 0/0 [00:00<?, ?it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.11it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.10it/s]
eval_pl-tensor(45.7752, device='cuda:0') eval_epoch_loss=tensor(3.8237, device='cuda:0')
we are about to save the PEFT modules
PEFT modules are saved in saved_peft_model directory
best eval loss on epoch 1 is 3.82374238077895s
```

```
algo-1:57:79 [0] ncc1_net_ofi_init:1444 NCC1_WARN NET/OFI Only EFA provider is supported
algo-1:57:79 [0] ncc1_net_ofi_init:1483 NCC1_WARN NET/OFI aws-ofi-ncc1 initialization failed
step 0 is completed and loss is 3.907542894622803
Training Epoch0: 100%[#033|3dm] 0/1 [00:00<0:00, 3.45it/s]
Training Epoch0: 100%[#033|3dm] 0/1 [00:00<0:00, 3.45it/s]
Max CUDA memory allocated was 15 GB
Max CUDA memory reserved was 15 GB
Peak active CUDA memory was 15 GB
Cuda Malloc retires : 128
CPU Total Peak Memory consumed during the train (max): 2 GB
evaluating Epoch: 0%[#033|3dm] 0/0 [00:00<?, ?it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.11it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.10it/s]
eval_pl-tensor(45.7752, device='cuda:0') eval_epoch_loss=tensor(3.7754, device='cuda:0')
we are about to save the PEFT modules
PEFT modules are saved in saved_peft_model directory
best eval loss on epoch 0 is 3.907542894622803
Epoch 2: train_perplexity=49.7727, train_epoch_loss=3.9075, epoch time 3.9417824229999498s
Training Epoch2: 0%[#033|3dm] 0/0 [00:00<?, ?it/s]
step 0 is completed and loss is 3.853542894622803
Training Epoch2: 100%[#033|3dm] 0/1 [00:03<0:00, 3.45s/it]
Training Epoch2: 100%[#033|3dm] 0/1 [00:03<0:00, 3.45s/it]
Max CUDA memory allocated was 15 GB
Max CUDA memory reserved was 15 GB
Peak active CUDA memory was 15 GB
Cuda Malloc retires : 128
CPU Total Peak Memory consumed during the train (max): 2 GB
evaluating Epoch: 0%[#033|3dm] 0/0 [00:00<?, ?it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.11it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.10it/s]
eval_pl-tensor(41.4619, device='cuda:0') eval_epoch_loss=tensor(3.7248, device='cuda:0')
we are about to save the PEFT modules
PEFT modules are saved in saved_peft_model directory
best eval loss on epoch 2 is 3.7248416222076416
Epoch 3: train_perplexity=47.1598, train_epoch_loss=3.8535, epoch time 3.9143312810000452s
Training Epoch3: 0%[#033|3dm] 0/0 [00:00<?, ?it/s]
step 0 is completed and loss is 3.7962462690206909
Training Epoch3: 100%[#033|3dm] 0/1 [00:03<0:00, 3.45s/it]
Training Epoch3: 100%[#033|3dm] 0/1 [00:03<0:00, 3.45s/it]
Max CUDA memory allocated was 15 GB
Max CUDA memory reserved was 15 GB
Peak active CUDA memory was 15 GB
Cuda Malloc retires : 192
CPU Total Peak Memory consumed during the train (max): 2 GB
evaluating Epoch: 0%[#033|3dm] 0/0 [00:00<?, ?it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.11it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.10it/s]
evaluating Epoch: 100%[#033|3dm] 0/1 [00:00<0:00, 3.10it/s]
eval_pl-tensor(41.4619, device='cuda:0') eval_epoch_loss=tensor(3.7248, device='cuda:0')
we are about to save the PEFT modules
PEFT modules are saved in saved_peft_model directory
best eval loss on epoch 3 is 3.724774360567383
Epoch 4: train_perplexity=44.5337, train_epoch_loss=3.7962, epoch time 3.915506209999971s
Training Epoch4: 0%[#033|3dm] 0/0 [00:00<?, ?it/s]
```

```

Epoch 4: train_perplexity=44.5337, train_epoch_loss=3.7962, epoch time 3.915506200999971s
Training Epoch: 0%#03|0m 0#03|0m 0/1 [00:00<?, ?it/s]
step 0 is completed and loss is 3.742390271987915
Training Epoch: 100%#03|0m 1#03|0m 1/1 [00:03<00:00, 3.45s/it]
Training Epoch: 100%#03|0m 1#03|0m 1/1 [00:03<00:00, 3.45s/it]
Max CUDA memory allocated was 15 GB
Max CUDA memory reserved was 15 GB
Peak active CUDA memory was 15 GB
Cuda Malloc retires : 256
CPU Total Peak Memory consumed during the train (max): 2 GB
evaluating: 0%#03|32m 0#03|0m 0/1 [00:00<?, ?it/s]
evaluating: 100%#03|32m 1#03|0m 1/1 [00:00<00:00, 3.18it/s]
evaluating: 100%#03|32m 1#03|0m 1/1 [00:00<00:00, 3.18it/s]
eval_checkpoint(39.70%) eval_epoch_loss=tensor(3.6814, device='cuda:0') eval_epoch_loss=tensor(3.6814, device='cuda:0')
we are about to save the PEFT modules
PEFT modules are saved in saved_peft_model directory
best eval loss on epoch 4 is 3.681448459625244
Epoch 5: train_perplexity=42.1949, train_epoch_loss=3.7423, epoch time 3.915930775999982s
INFO:root:Key: avg_train_prep, Value: 47.178466796875
INFO:root:Key: avg_train_loss, Value: 3.85169472774505615
INFO:root:Key: avg_eval_prep, Value: 43.716985015105625
INFO:root:Key: avg_eval_loss, Value: 3.209526089932
INFO:root:Key: avg_checkpoint_time, Value: 4.28953103159977
INFO:root:Key: avg_checkpoint_time, Value: 0.9383596705999682
INFO:root:Combining pre-trained base model with the PEFT adapter module.
Loading checkpoint shards: 0% [ 0/2 [00:00<?, ?it/s]
Loading checkpoint shards: 50% [ 1/2 [00:29<00:29, 29.62s/it]
Loading checkpoint shards: 100% [ 2/2 [00:35<00:00, 15.66s/it]
Loading checkpoint shards: 100% [ 2/2 [00:35<00:00, 17.76s/it]
INFO:root:Saving the combined model in safetensors format.
INFO:root:Saving the combined model in safetensors format.
INFO:root:Copying tokenizer to the output directory.
INFO:root:Putting inference code with the fine-tuned model directory.
2024-09-27 14:25:02,693 sagemaker-training-toolkit INFO Waiting for the process to finish and give a return code.
2024-09-27 14:25:02,693 sagemaker-training-toolkit INFO Done waiting for a return code. Received 0 from exiting process.
2024-09-27 14:25:02,693 sagemaker-training-toolkit INFO Reporting training SUCCESS

```

```

INFO:root:Putting inference code with the fine-tuned model directory.
2024-09-27 14:25:02,693 sagemaker-training-toolkit INFO Waiting for the process to finish and give a return code.
2024-09-27 14:25:02,693 sagemaker-training-toolkit INFO Done waiting for a return code. Received 0 from exiting process.
2024-09-27 14:25:02,693 sagemaker-training-toolkit INFO Reporting training SUCCESS

2024-09-27 14:25:27 Uploading - Uploading generated training model
2024-09-27 14:26:10 Completed - Training job completed
Training seconds: 728
Billable seconds: 728

Deploy the fine-tuned model

```

Next, we deploy the domain fine-tuned model. We will compare the performance of the fine-tuned and pre-trained model.

```

[15]: ...
# Do not use estimator.deploy() without mentioning the instance_type...
# It's because when you call estimator.deploy() without explicitly setting the instance_type_for_the_endpoint...
# SageMaker selects a default instance type for hosting, which, in this case, is ml.g5.2xlarge.
# However, Udacity doesn't allow instance type more than "ml.*.2xlarge"...
...
finetuned_predictor = estimator.deploy(instance_type="ml.g5.2xlarge", initial_instance_count=1)

INFO:sagemaker:Creating model with name: meta-textgeneration-llama-2-7b-2024-09-27-14-26-27-392
INFO:sagemaker:Creating endpoint-config with name meta-textgeneration-llama-2-7b-2024-09-27-14-26-27-388
INFO:sagemaker:Creating endpoint with name meta-textgeneration-llama-2-7b-2024-09-27-14-26-27-388
-----!

```

Evaluate the pre-trained and fine-tuned model

Next, we use the same input from the model evaluation step to evaluate the performance of the fine-tuned model and compare it with the base pre-trained model.

Project: Building a Domain Expertise Model | Notebook instances | Amazon | Model_FineTuning.ipynb | Model_Evaluation_Udacity.ipynb | Disha520/Sagemaker-IT-Domains | New Tab

fine-tuningproject1-e4y2.notebook.us-west-2.sagemaker.aws/lab/tree/Model_FineTuning.ipynb

Evaluate the pre-trained and fine-tuned model

Next, we use the same input from the model evaluation step to evaluate the performance of the fine-tuned model and compare it with the base pre-trained model.

Create a function to print the response from the model

```
[140]: def print_response(payload, response):
    print(payload["inputs"])
    print(f"> {response}")
    print("-----\n")
```

Now we can run the same prompts on the fine-tuned model to evaluate its domain knowledge.

Replace "inputs" in the next cell with the input to send the model based on the domain you've chosen.

For financial domain:

"inputs": "Replace with sentence below from text"

- "The investment tests performed indicate"
- "the relative volume for the long out of the money options, indicates"
- "The results for the short in the money options"
- "The results are encouraging for aggressive investors"

For medical domain:

"inputs": "Replace with sentence below from text"

- "Myeloid neoplasms and acute leukemias derive from"

Simple Fully initialized conda_pytorch_p310 | Idle Mode: Command Ln 11, Col 89 Model_FineTuning.ipynb 0 27-09-2024

80°F Rain

Project: Building a Domain Expertise Model | Notebook instances | Amazon | Model_FineTuning.ipynb | Model_Evaluation_Udacity.ipynb | Disha520/Sagemaker-IT-Domains | New Tab

fine-tuningproject1-e4y2.notebook.us-west-2.sagemaker.aws/lab/tree/Model_FineTuning.ipynb

Evaluate the pre-trained and fine-tuned model

Next, we use the same input from the model evaluation step to evaluate the performance of the fine-tuned model and compare it with the base pre-trained model.

Create a function to print the response from the model

```
[140]: def print_response(payload, response):
    print(payload["inputs"])
    print(f"> {response}")
    print("-----\n")
```

Now we can run the same prompts on the fine-tuned model to evaluate its domain knowledge.

Replace "inputs" in the next cell with the input to send the model based on the domain you've chosen.

For financial domain:

"inputs": "Replace with sentence below from text"

- "Traditional approaches to data management such as"
- "A second important aspect of ubiquitous computing environments is"
- "because ubiquitous computing is intended to"
- "outline the key aspects of ubiquitous computing from a data management perspective."

```
[177]: payload = {
    "inputs": "outline the key aspects of ubiquitous computing from a data management perspective.",
    "parameters": {
        "max_new_tokens": 64,
        "top_p": 0.9,
        "temperature": 0.6,
        "return_full_text": False,
    },
}
try:
    response = finetuned_predictor.predict(payload, custom_attributes="accept_eula=true")
    print_response(payload, response)
except Exception as e:
    print(e)
```

outline the key aspects of ubiquitous computing from a data management perspective.

> [{"generated_text": "\nThe book is divided into four parts. Part I provides an overview of the field, including an introduction to ubiquitous computing, its history and its relationship to mobile computing, as well as a discussion of the key technologies that enable ubiquitous computing. Part II discusses the data management challenges"}]

Do the outputs from the fine-tuned model provide domain-specific insightful and relevant content? You can continue experimenting with the inputs of the model to test its domain knowledge.

Use the output from this notebook to fill out the "model fine-tuning" section of the project documentation report

After you've filled out the report, run the cells below to delete the model deployment

Simple Fully initialized conda_pytorch_p310 | Idle Mode: Command Ln 11, Col 89 Model_FineTuning.ipynb 0 27-09-2024

80°F Rain

```
Use the output from this notebook to fill out the "model fine-tuning" section of the project documentation report  
After you've filled out the report, run the cells below to delete the model deployment  
IF YOU FAIL TO RUN THE CELLS BELOW YOU WILL RUN OUT OF BUDGET TO COMPLETE THE PROJECT  
[18]: finetuned_predictor.delete_model()  
finetuned_predictor.delete_endpoint()  
INFO:sagemaker:Deleting model with name: meta-textgeneration-lama-2-7b-2024-09-27-14-26-27-392  
INFO:sagemaker:Deleting endpoint configuration with name: meta-textgeneration-lama-2-7b-2024-09-27-14-26-27-388  
INFO:sagemaker:Deleting endpoint with name: meta-textgeneration-lama-2-7b-2024-09-27-14-26-27-388
```

Visit the AWS S3 bucket where your fine-tuned model weights are stored after training and take a screenshot for your submission.

Name	AWS Region	IAM Access Analyzer	Creation date
sagemaker-us-west-2-611383594527	US West (Oregon) us-west-2	View analyzer for us-west-2	September 27, 2024, 19:30:09 (UTC+05:30)

Screenshot of verification of the model has been deleted.

The screenshot shows the AWS SageMaker console interface. The left sidebar is expanded, showing various services like Labeling datasets, Processing, Training, Inference, and Augmented AI. Under the 'Models' section, which is currently selected, there is a sub-menu for Endpoint configurations, Endpoints, Batch transform jobs, Shadow tests, and Inference Recommender. The main content area is titled 'Models' and contains a search bar and a table with columns for Name, ARN, and Creation time. A message at the bottom of the table says 'There are currently no resources.' The top navigation bar shows the URL 'us-west-2.console.aws.amazon.com/sagemaker/home?region=us-west-2#models'. The bottom status bar indicates it's 2021, ENG IN, and the date is 27-09-2024.

Screenshot of verification of the endpoint has been deleted.

The screenshot shows the AWS SageMaker console interface. The left sidebar is expanded, showing various services like Labeling datasets, Processing, Training, Inference, and Augmented AI. Under the 'Endpoints' section, which is currently selected, there is a sub-menu for Endpoint configurations, Endpoints, Batch transform jobs, Shadow tests, and Inference Recommender. The main content area is titled 'Endpoints' and contains a search bar and a table with columns for Name, ARN, Creation time, Status, and Last updated. A message at the bottom of the table says 'There are currently no resources.' The top navigation bar shows the URL 'us-west-2.console.aws.amazon.com/sagemaker/home?region=us-west-2#endpoints'. The bottom status bar indicates it's 2021, ENG IN, and the date is 27-09-2024.