# Lab 3:  RAG application using Amazon OpenSearch
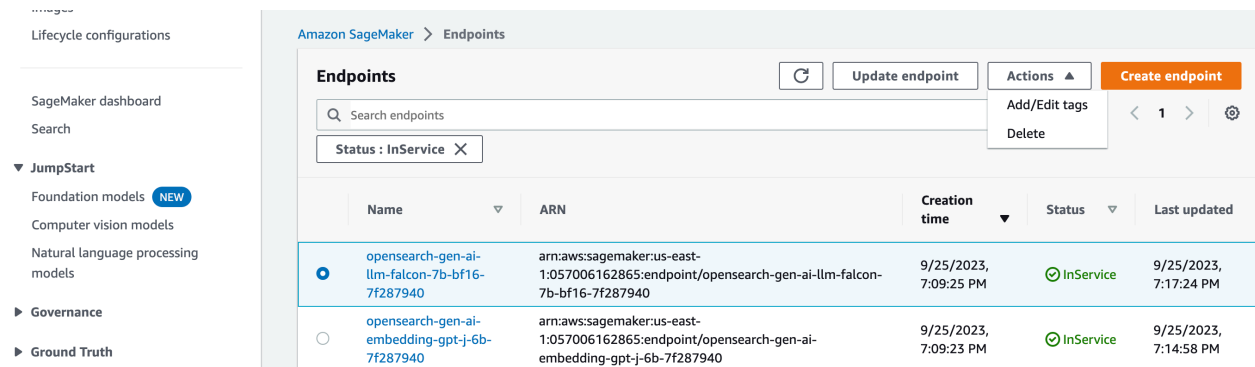
1.  Complete the pre-requisites for the lab
    a.  Delete both the endpoints created in the previous labs
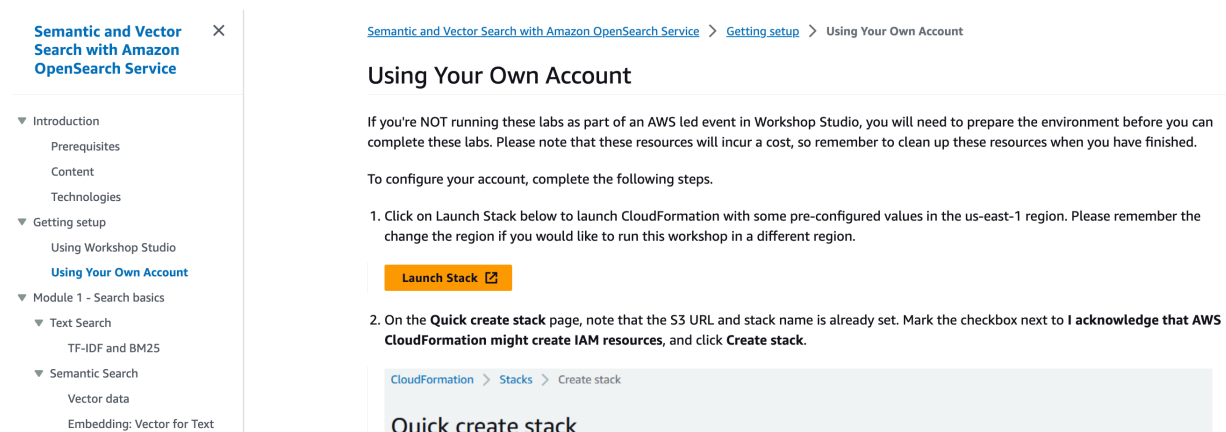
Note:  the name of the end point will be different in your case



b.  Download the dataset required for running the lab.  The dataset is available in https://github.com/thandavm/rag_sm_js/blob/main/data/winemag-data-130k-v2.json

2.  Launch the Cloud formation stack from https://catalog.workshops.aws/semantic-search/en-US/setup/using-own-account



3.  Select "I Acknowledge…." And click on "Create Stack".  The stack creation takes ~15 mins to complete

4. Once the stack is created. Go to SageMaker -> Notebook Instances.



5. Open "semantic-search-nb" and launch "Open Jupyter"



6. Click on "Upload" and load the downloaded "winemag-data-130k-v2.json"



7. Open the Notebook "Module 7 - Retrieval Augmented Generation.ipynb" and add the endpoints created as part of the cloud formation template. Search for the cell below and add the SM falcon end point name here

```
In [ ]: # If you already deployed a model,
        # uncomment the following lines and add your endpoint name below

        from sagemaker.huggingface import HuggingFacePredictor
        sagemaker_session = sagemaker.Session()
        llm_endpoint_name = "opensearch-gen-ai-llm-falcon-7b-bf16-7f287940"
        llm_predictor = HuggingFacePredictor(endpoint_name=llm_endpoint_name, sagemaker_session = sagemaker_session)
```

8. Start executing the Notebook
9. Do not execute Step 12, because we have already deployed the model

### 12. Deploy the Large Language Model for Retrieval Augmented Generation

This module uses the Falcon 7B model to create recommendations based on a given wine review. The next cell deploys a model endpoint into your environment that will be called by subsequent steps. For more information on the Falcon LLM, see HuggingFace's announcement regarding the model.

```
In [ ]: image_uri = get_huggingface_llm_image_uri(
            backend="huggingface", # or lmi
            region=region
        )

        model_name = "falcon-7b-" + time.strftime("%Y-%m-%d-%H-%M-%S", time.gmtime())

        hub = {
            'HF_MODEL_ID':'tiiuae/falcon-7b',
            'HF_TASK':'question-answering',
            'SM_NUM_GPUS':'1',
            'HF_MODEL_QUANTIZE':'bitsandbytes'
        }

        model = HuggingFaceModel(
            name=model_name,
            env=hub,
            role=role,
```

10. Continue and complete the lab!!!