



"Deploy



Hugging Face Transformers on AWS

in just a few lines of code"



Nicola Procopio 30/05/2023

About me



Nicola Procopio Senior Data Scientist

Contacts



nicola.procopio@finconsgroup.com



https://it.linkedin.com/in/nicolaprocopio



https://github.com/nickprock



https://huggingface.co/nickprock





Background

















Community











What's Hugging Face?

"We are on a mission to democratize **good** machine learning, one commit at a time"

- Hugging Face is the collaboration platform for the machine learning community.
- The Hugging Face Hub works as a central place where anyone can share, explore, discover, and experiment with open-source ML.
- With the fast-growing community, some of the most used open-source ML libraries and tools, and a talented science team exploring the edge of tech, Hugging Face is at the heart of the Al revolution.

Models: > 210K

Datasets: > 36K

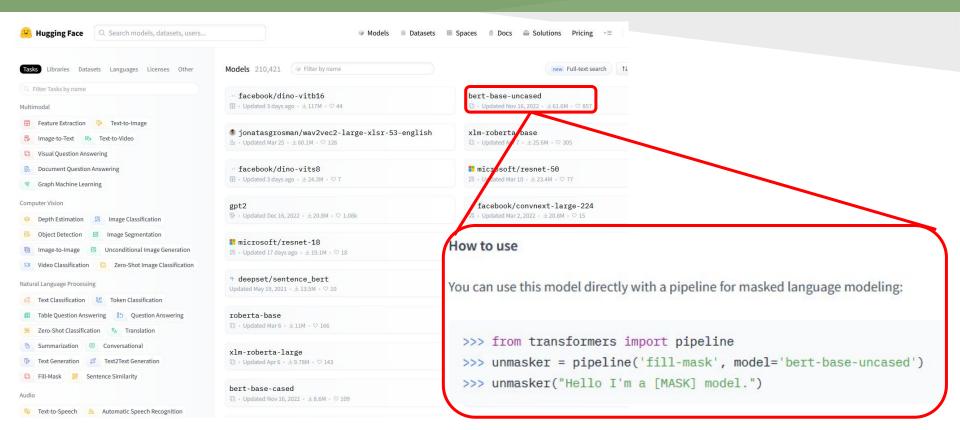
★ Transformers: 102K stars on GitHub

★ Datasets: > 16.3K

★ Diffusers: > 14K

x .

Hub & Libraries



Why Hugging Face on AWS?

Hugging Face Deep Learning Containers are available on AWS Sagemaker beacuse there are many advantages:

- One command is all you need
- Accelerate machine learning from science to production (AWS Toolkit)
- Built-in performance
- Costs

On Sagemaker you can:

- train
- deploy
- accelerate
- ...

your models

```
from sagemaker.huggingface import HuggingFaceModel

# create Hugging Face Model Class and deploy it as SageMaker Endpoint
huggingface_model = HuggingFaceModel(...).deploy()
```

Focus on Deploy

Deploy an Hugging Face model is easy with Inference Toolkit. First step install and setup the environment.

```
pip install sagemaker --upgrade
```

Sagemaker Environment

```
import sagemaker
sess = sagemaker.Session()
role = sagemaker.get_execution_role()
```

Local Environment

```
import sagemaker
import boto3

iam_client = boto3.client('iam')

role = iam_client.get_role(RoleName='role-name-of-your-iam-role-with-right-permissions')['Role']['Arn']
sess = sagemaker.Session()
```

Deploy from Model Hub

- Choose your model from Hugging Face Hub and create a model config.
 - {HF MODEL ID: checkpoint, HF TASK: text-classification}

Create your model class

```
env=hub,
role=role,
transformers_version="4.26",
pytorch_version="1.13",
py_version='py39',
```

Deploy the endpoint

```
# deploy model to SageMaker Inference
predictor = huggingface_model.deploy(
   initial_instance_count=1,
   instance_type="ml.m5.xlarge"
)
```

create Hugging Face Model Class
huggingface model = HuggingFaceModel(

REMEMBER!
Each HF_TASK has an inputs
format!!!

Deploy from S3

- Create a .tar.gz file contains your model
- Create your role and model class

```
import sagemaker
import boto3
   role = sagemaker.get_execution_role()
except ValueError:
   iam = boto3.client('iam')
   role = iam.get role(RoleName='sagemaker execution role')['Role']['Arn']
print(f"sagemaker role arn: {role}")
from sagemaker.huggingface import HuggingFaceModel
# create Hugging Face Model Class
 huggingface model = HuggingFaceModel(
    model data="s3://hf-sagemaker-inference/model.tar.gz", # path to your trained sagemaker model
    role=role, # iam role with permissions to create an Endpoint
    transformers version="4.26", # transformers version used
    pytorch version="1.13", # pytorch version used
    py version="py39", # python version of the DLC
```

Serverless Inference

- Deploy and scale ML models in a Serverless fashion
- Serverless endpoints automatically launch compute resources and scale them in and out depending on traffic similar to AWS Lambda.
- Serverless Inference is ideal for workloads which have idle periods between traffic spurts and can tolerate cold starts.
- Pay-per-use model
- Some Limitations:
 - Memory size
 - Concurrent Invocations
 - Cold Start



Serverless Inference

How the code changes?

Only import a module, the object initialization and a few parameters settings

```
from sagemaker.huggingface.model import HuggingFaceModel
from sagemaker.serverless import ServerlessInferenceConfig
# Hub Model configuration. <a href="https://huggingface.co/models">https://huggingface.co/models</a>
hub = {
    'HF MODEL ID': 'distilbert-base-uncased-finetuned-sst-2-english',
    'HF TASK': 'text-classification'
# create Hugging Face Model Class
huggingface model = HuggingFaceModel(
                                 # configuration for loading model from Hub
   env=hub.
   role=role.
                 # iam role with permissions to create an Endpoint
   transformers version="4.26", # transformers version used
   pytorch_version="1.13", # pytorch version used
   py_version='py39',
                       # python version used
# Specify MemorySizeInMB and MaxConcurrency in the serverless config object
serverless config = ServerlessInferenceConfig(
    memory size in mb=4096, max concurrency=10,
# deploy the endpoint endpoint
predictor = huggingface model.deploy(
    serverless inference config=serverless config
```

- You need to use libraries not directly supported by the toolkit (es. SetFit, Sentence Transformers, SpanMarker, ...)
- Your output format isn't the default
- Other reasons why you cannot use the inference endpoint as it is



In our model folder we must add:

- 1. requirements.txt
- 2. handler.py
- 3. code/inference.py

Create a deploy.py, the script which is creating your endpoint.

requirements.txt is the file that contains the dependencies, in my case it had only one line: **setfit.**

For S3 endpoint put it into code folder.

handler.py defines:

- how the endpoint must load the model
- what the input looks like
- what should return as output

```
from typing import Dict, List, Any
from setfit import SetFitModel
class EndpointHandler:
    def __init__(self, path=""):
        # load model
        self.model = SetFitModel.from_pretrained(path)
        # ag_news id to label mapping
        self.id2label = {0: "World", 1: "Sports", 2: "Business", 3: "Sci/Tech"}
    def __call__(self, data: Dict[str, Any]) -> List[Dict[str, Any]]:
         data args:
              inputs (:obj: 'str')
        Return:
              A :obj: 'list' | 'dict': will be serialized and returned
        # get inputs
        inputs = data.pop("inputs", data)
        if isinstance(inputs, str):
            inputs = [inputs]
        # run normal prediction
        scores = self.model.predict_proba(inputs)[0]
        return [{"label": self.id2label[i], "score": score.item()} for i, score in enumerate(scores)]
```

The code folder contains inference.py:

```
def model_fn(model_dir):
model fn
                           model = SetFitModel.from pretrained(model dir)
                           return model
                    def input_fn(input_data, content_type):
input fn
                          decoded input data = decoder encoder.decode(input data, content type)
                          return decoded input data
                     def predict fn(data, model):
                       setFit model = model
                       inputs = data.pop("inputs", data)
predict fn
                       if isinstance(inputs. str):
                          inputs = [inputs]
                       # run normal prediction
                       scores = model.predict_proba(inputs)[0]
                       return [{"label": model id2label[i], "score": round(score.item(),4)} for i, score in enumerate(scores)]
```

output_fn → decode the result

Create the endpoint, file **deploy.py**:

```
hf_model = HuggingFaceModel(
       model data="s3://call-your-model/model1.tar.gz",
       role=role,
       transformers version="4.6.1",
       pytorch_version="1.7.1",
                                      def predict(input):
       source dir="code",
       py version="py36",
                                          predictor = HuggingFacePredictor(
       source dir="model1/code"
                                              endpoint_name="huggingface-endpoint")
       entry point="inference.py",
                                          resp = predictor.predict({"inputs": input})
                                          print(resp)
                                          out = [data['label'] for data in resp if data['score'] >= 0.5]
                                          return out
 predictor = huggingface model.deploy(
     serverless_inference_config = ServerlessInferenceConfig(memory_size_in_mb=6144)
```

Notes

- Code repository:
 - Hugging Face Sagemaker workshop series
 - Hugging Face notebooks on Sagemaker
- Follow Philipp Schmid (Hugging Face ML Eng & AWS ML Hero):
 - o <u>github</u>
 - website
- Blog:
 - The Partnership: Amazon SageMaker and Hugging Face
 - <u>Deploy Hugging Face models easily with Amazon SageMaker</u>



Thanks!





Remember to join Hugging Face