

Big Data Laboratory

Assignment 7 Deployment of ML model using Fast Api and monitoring

by

Anik Bhowmick
AE20B102

Instructor: Sudarshan
Teaching Assistant: Jashaswimalya Acharjee



Contents

1	Introduction	1
1.1	The problem statement	1
2	Result and conclusion	2
A	Code section	5

1

Introduction

As part of MLOps, this assignment aims to deploy the ML model on a web server using Fast Api and monitor it with Prometheus and Grafana. FastAPI is a modern, fast (high-performance), web framework for building APIs with Python 3.8+ based on standard Python type hints. As a part of this assignment, we had to perform versioning of our project using Git Hub. To have GitHub access to all the files **click here**.

1.1. The problem statement

We have to set up a web interface to deploy an ML model trained on digit classification tasks. FastAPI has to be used to create the API, and Swagger UI for the webserver. On top of that, we have to use Prometheus and Grafana for server health monitoring. The monitoring metrics are like: API run time, API T/L time, API memory utilization, API CPU utilization (rate), API network I/O bytes (and rate).

Next, we have to dockerize and spawn a number of FastAPI servers using different ports.

2

Result and conclusion

Below are some pictures of Grafana UI.



Figure 2.1: Grafana UI for memory utilization, API runtime, etc



Figure 2.2: Grafana UI for CPU utilization, Client IP count etc

Prometheus UI.

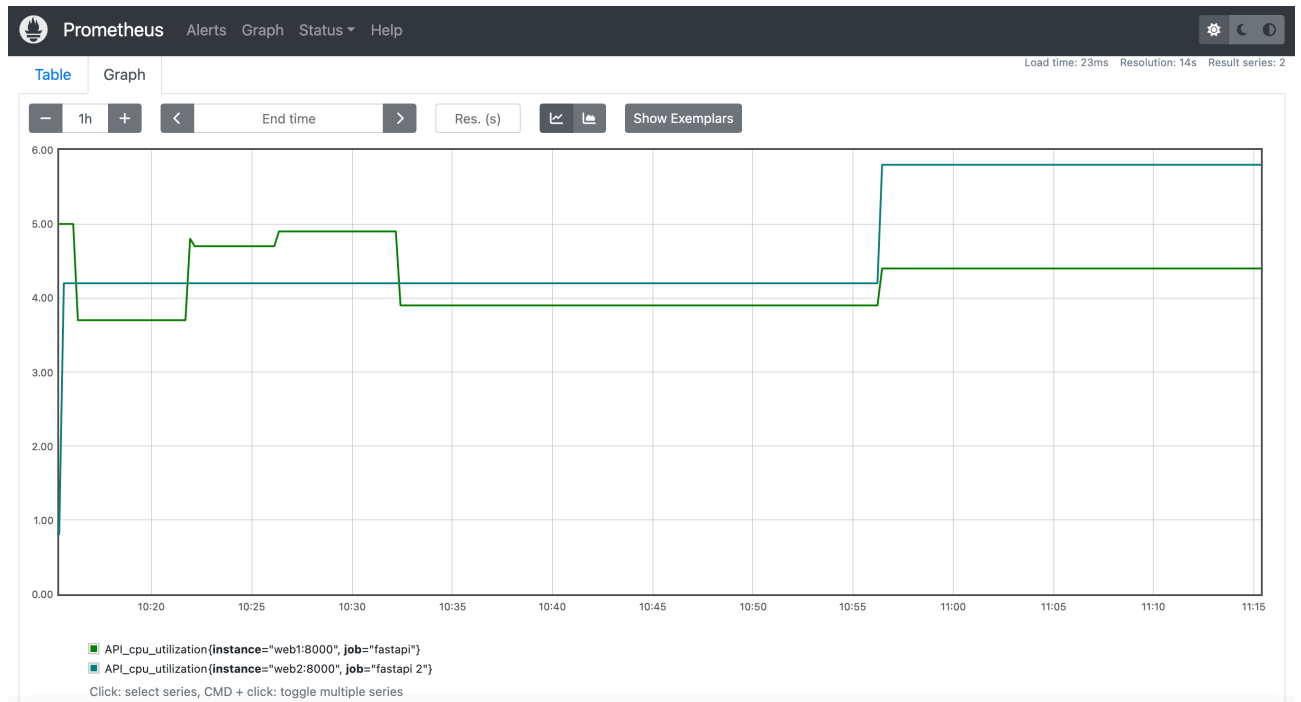


Figure 2.3: Prometheus UI for API CPU utilization

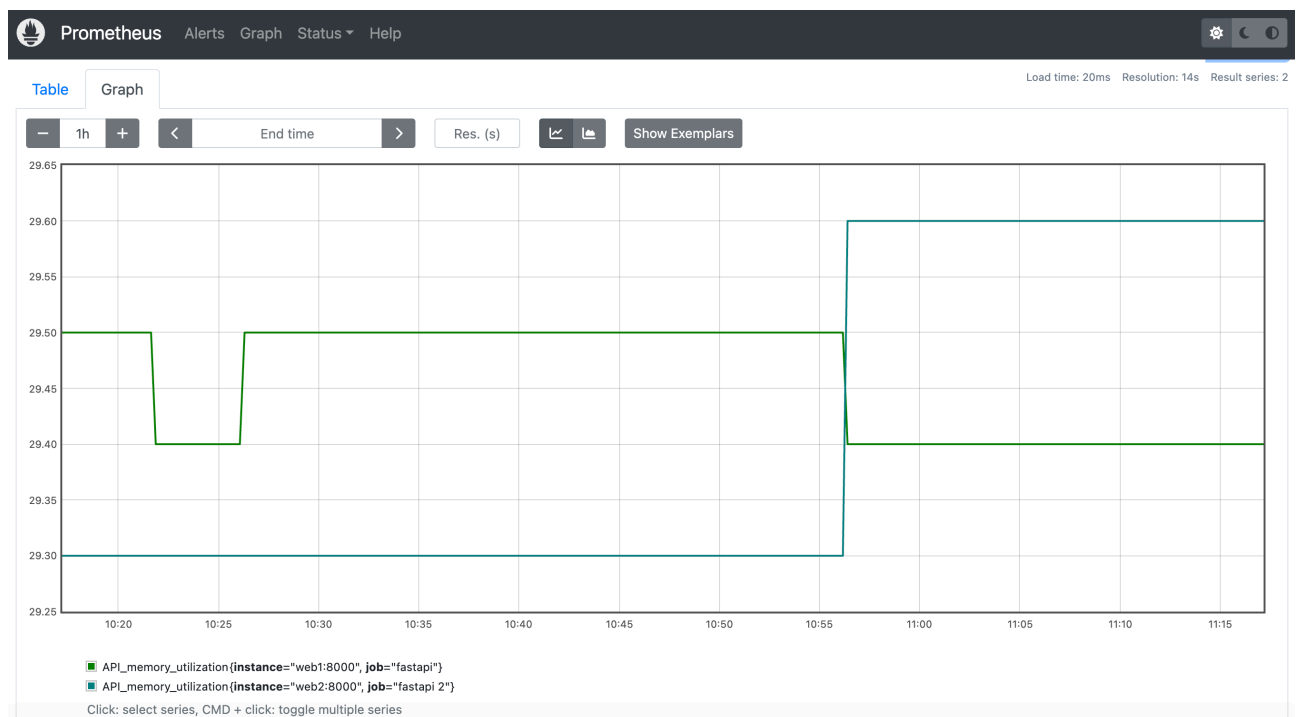
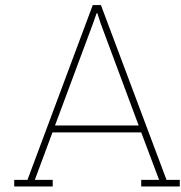


Figure 2.4: Prometheus UI for API memory utilization

For custom handwritten digits:



Figure 2.5: Prometheus UI for API runtime



Code section

Fast Api code

```
1 import os
2 from typing import List
3 from fastapi import FastAPI, UploadFile, File, HTTPException, Request
4 import numpy as np
5 from tensorflow.keras.models import load_model
6 import tensorflow as tf
7 from PIL import Image, ImageOps
8 from io import BytesIO
9 import sys
10 import time
11 import base64
12 import psutil
13 from prometheus_fastapi_instrumentator import Instrumentator
14 from prometheus_client import Counter, Gauge
15
16
17 # Create the FastAPI app
18 app = FastAPI()
19
20 # Load the model from the specified path
21 def get_model(path: str):
22     return load_model(path)
23
24 # Load the MNIST model
25
26 model_path = "MNIST_model.keras"
27 #load the pretrained model
28 model = get_model(model_path)
29 # set the model in inference mode
30 model.trainable=False
31 #counts the number of times a client ip address visit
32 request_counter = Counter("Client_IP_count", "Total_number_of_requests", ["client_ip"])
33 #Calculates the run time of API
34 inference_time_gauge = Gauge("API_runtime", "Inference_time_in_seconds")
35 #processing time of the text, per character calculation
36 processing_time_per_char_gauge = Gauge("processing_time_per_char_microseconds", "Processing_time_per_character_in_microseconds")
37
38 # get API network I/O bytes
39 network_receive_bytes = Gauge("network_receive_bytes", "Total_network_receive_bytes")
40 network_transmit_bytes = Gauge("network_transmit_bytes", "Total_network_transmit_bytes")
41
42 # get the API memory utilization
43 memory_utilization = Gauge("API_memory_utilization", "API_memory_utilization")
44 # get the CPU utilization
45 cpu_utilization = Gauge("API_cpu_utilization", "API_CPU_utilization")
46
47
48 # Function to preprocess image and make prediction
49
50
51 def predict_digit(model, data_point):
52
53     # get the prediction containg the score
```

```

54     pred = model.predict(data_point)
55     # get the class label
56     prediction = tf.argmax(pred,axis=-1)
57     c_score = np.max(pred)# store the confidence score
58     return str(prediction[0].numpy()),str(c_score)
59
60 # API endpoint to accept image upload and return prediction
61 def format_image(image):
62     """
63     get a pillow image
64     """
65     # resize the image in 28X28 format
66     return image.resize((28,28))
67
68 Instrumentator().instrument(app).expose(app)
69
70 @app.post('/predict')
71 async def predict(request: Request, file: UploadFile = File(...)):
72     # load the image in the byte format
73     content = await file.read()
74     accepted_formats = ['.jpeg', '.jpg', '.png']
75     file_format = os.path.splitext(file.filename)[1].lower()
76     # check for the image file is valid or not
77     if file_format not in accepted_formats:
78         # raise the error message if the file is wrong in format
79         raise HTTPException(status_code=400, detail="Bad file format. Accepted formats are .
            jpeg, .jpg, .png")
80     file_name = os.path.splitext(file.filename)[0]
81     image = Image.open(BytesIO(content))
82     #convert the image to gray scale first
83     image = image.convert('L')
84     if image.size!=(28,28):
85         # if the image is not 28 by 28 resize it
86         image = format_image(image)
87     flat = np.array(image,dtype='float32').reshape(-1)/255.0# flatten the image and normalize
            in 0 to 1 scale
88     flat = flat[None,:]
89
90     client_ip = request.client.host
91     request_counter.labels(client_ip=client_ip).inc()
92
93     start_time = time.time()
94     output, c_score = predict_digit(model, flat)
95     end_time = time.time()
96     inference_time_gauge.set(end_time - start_time)
97
98
99     # Get memory utilization and update Prometheus metric
100    memory_utilization.set(psutil.virtual_memory().percent)
101
102    # Get CPU utilization and update Prometheus metric
103    cpu_utilization.set(psutil.cpu_percent())
104
105    input_length = len(file_name)
106    processing_time_per_char = (end_time - start_time) / input_length * 1e6 # microseconds
            per character
107    processing_time_per_char_gauge.set(processing_time_per_char)
108
109    net_io = psutil.net_io_counters()
110    network_receive_bytes.set(net_io.bytes_recv)
111    network_transmit_bytes.set(net_io.bytes_sent)
112    return {
113        "actual":file_name,
114        "predicted": output,
115        "confidence":c_score}

```