

Report 1 [IoT]

Student ID: 240581505

Student Name: Shobika Rajeskanna

Introduction:

A pipeline to process Internet of Things (IoT) sensor data in an edge-cloud environment has been created. The operation to be performed are pulling and running EMQX, data injection, data preprocessing, pulling and running RABBITMQ, machine learning model and visualization. Firstly, raw data has been fetched from Urban Observatory and from many metrics, need to store and analysis PM2.5 data with timestamp and values only. Then the filtered data has been sent to EMQX service of Azure lab (edge). Further, collect all PM2.5 data published from EMQX service, next filter out outliers which value is greater than 50. Then, the results (averaged data) have been taken into RabbitMQ service on Azure lab(cloud). After performing data preprocessing, collected daily average PM2.5 data. Furthermore, converted timestamp to date time format and using matplotlib displayed line chart of daily average of PM2.5 data. Finally, using feed averaged PM2.5 data to machine learning model to predict the trend of PM2.5 for next 15 days and visualized the predicted result machine learning classifier model.

Task-1

Design a data injector component by leveraging Newcastle Urban Observatory IoT data streams.

Step 1: Pull Docker image emqx/emqx from Dc:

```
docker pull emqx/emqx
```

Step 2: Run Docker image emqx/emqx from Docker Hub:

```
docker run -d --name emqx_broker -p 1883:1883 emqx/emqx
```

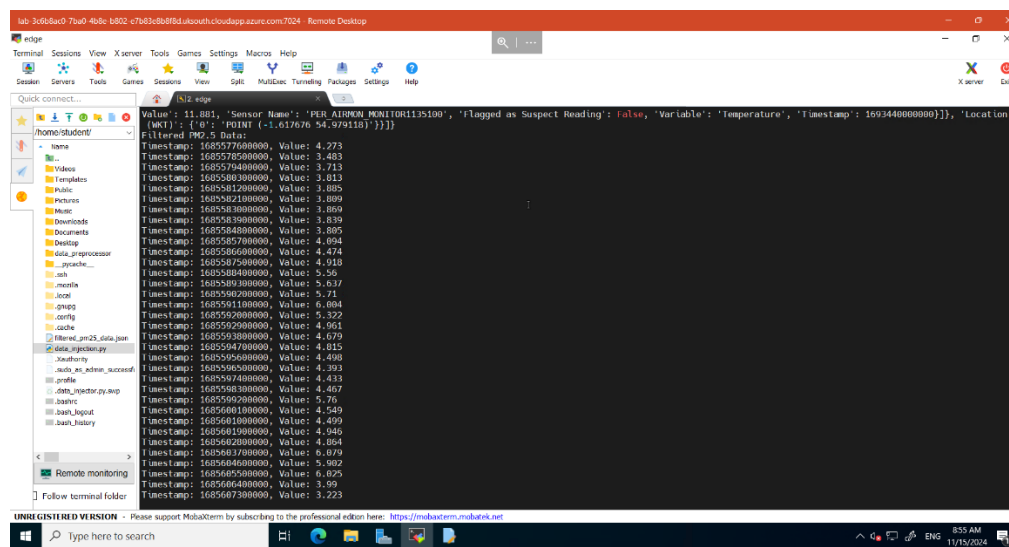
EMQX helps in filtering relevant data (like pm2.5 data) by acting as a local broker for distributing data to data preprocessing operator.

Step 3: Development of data injector code using python in Azure Lab (Edge):

Functions in the code: collect data from Urban Observatory platform by sending HTTP request to the given API. Then fetch the raw data and printed them in the console.

[illegible]

In Urban Observatory API, there are many air quality sensor data available. We need to print PM2.5 data and in their meta-data, we need to print Timestamp and Values alone.



Firstly, need to setup MQTT broker information like Host's IP, Port number and channel topic (in my case the topic name is sensor/pm2.5) where we send our data. Then, connecting to MQTT broker:

Finally sending data to EMQX:

Code snippet for performing the operation:

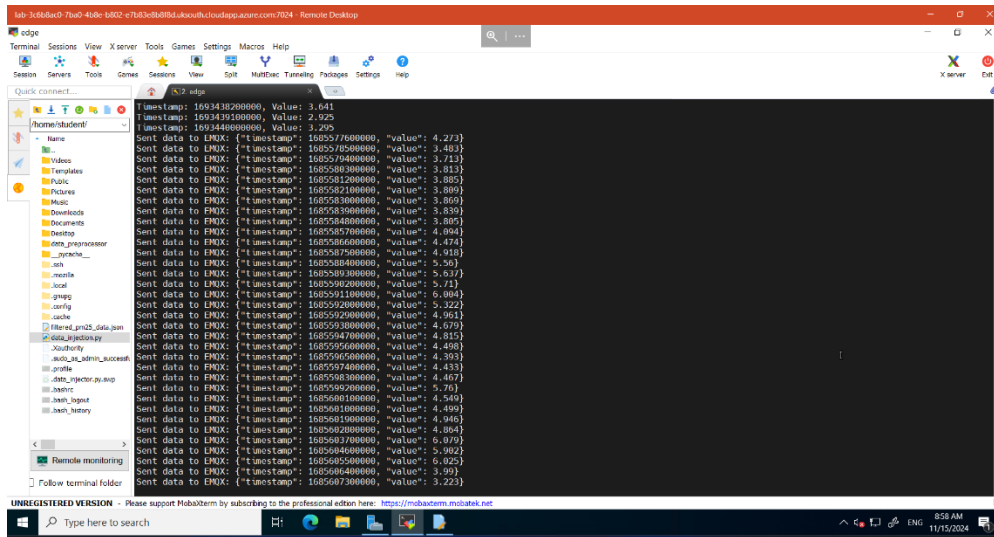
```
for data_point in pm25_data:
```

```

message = json.dumps(data_point)
client.publish(MQTT_TOPIC, message)
print(f"Sent data to EMQX: {message}")
time.sleep(0)

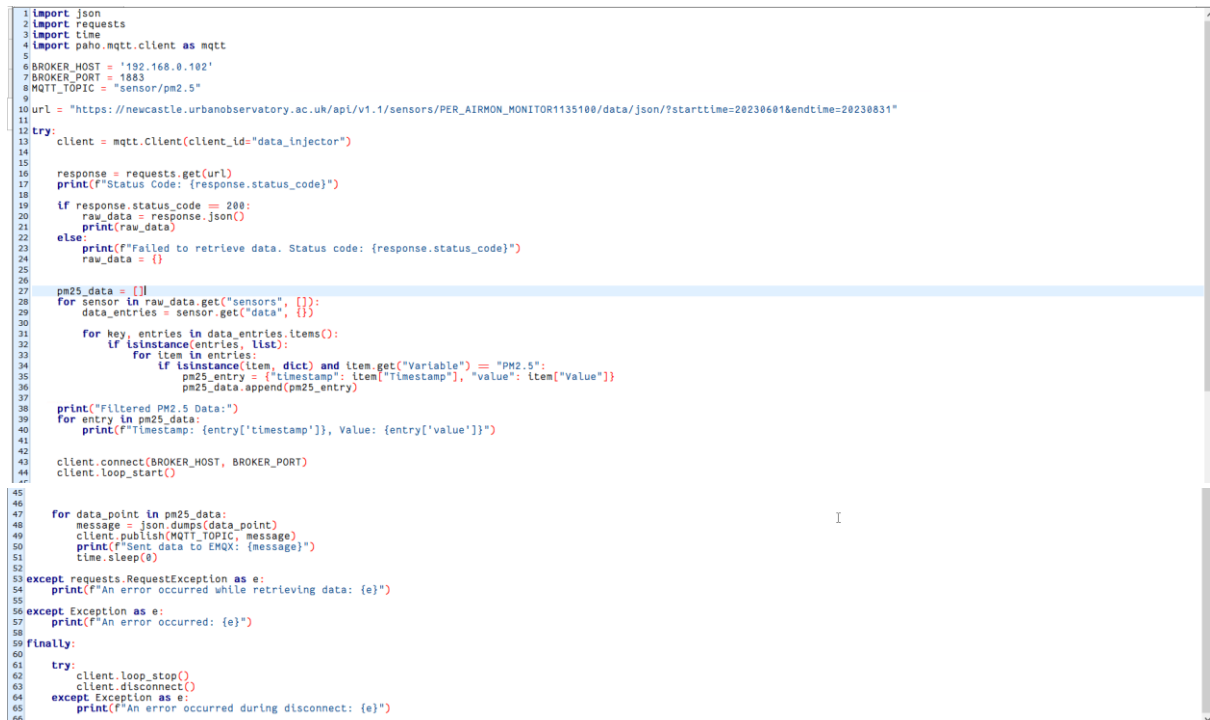
```

Output:



In below screenshot I have added the python program that performs all the functions I mentioned earlier.

data_injector.py



Task 2

Data preprocessing operator design

Step1:

In cloud pull and run the RabbitMQ, receive data from the edge (via EMQX) and queue it for further processing and visualizing.

Command to pull RabbitMQ:

```
docker pull rabbitmq:management
```

Command to run RabbitMQ:

```
docker run -d --name rabbitmq -p 5672:5672 -p 15672:15672 rabbitmq:3-management
```

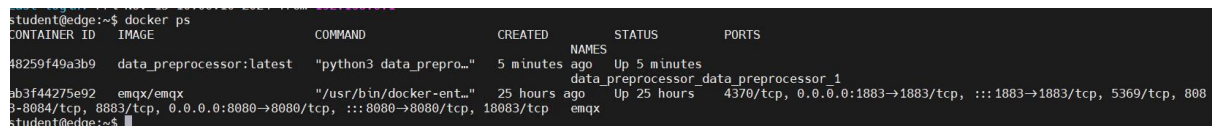
Step2:

Now subscribe all the PM2.5 data that is been sent from task 1 from EMQX service. Containerize this operator using Docker container, so the logs can be seen in the docker logs.

For containerization we need docker-compose, which will help us to start a container in correct order.

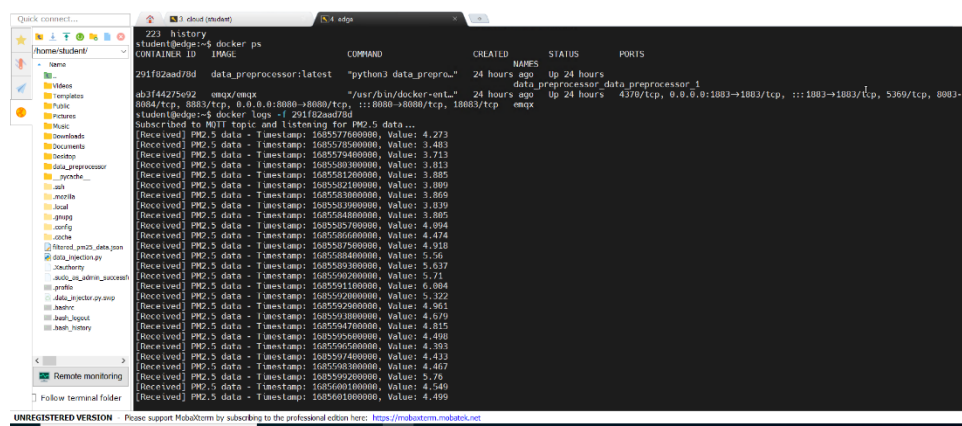
Command for docker-compose: `docker-compose up`

Containers in edge:



CONTAINER ID	IMAGE	COMMAND	CREATED	STATUS	PORTS
48259f49a3b9	data_preprocessor:latest	"python3 data_prepro..."	5 minutes ago	Up 5 minutes	
ab3f44275e92	emqx/emqx	"/usr/bin/docker-ent..."	25 hours ago	Up 25 hours	4370/tcp, 0.0.0.0:1883→1883/tcp, :::1883→1883/tcp, 5369/tcp, 8083-8084/tcp, 8883/tcp, 0.0.0.0:8080→8080/tcp, :::8080→8080/tcp, 18083/tcp

Docker log in edge:



The above output is fetching data from EMQX service and as I containerized my data preprocessor it is been logged and logs can be seen automatically.

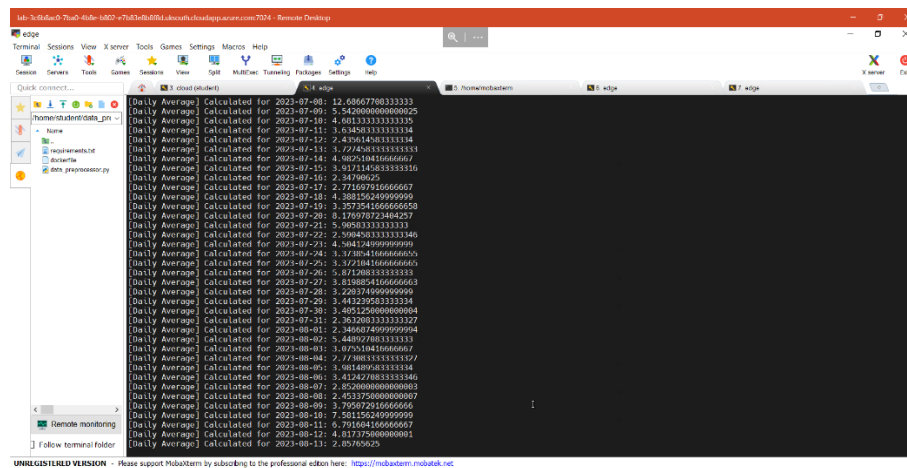
Step 3: Outliers

Perform a function in the code to find the values greater than 50.

Step 4:

Average value of PM2.5 data in daily basis has been created and the result has been printed in docker logs console.

Output:



Step 5:

Finally, the averaged PM2.5 data has been published into RabbitMQ service on Azure lab (Cloud). Connection has been made, for the connection we need RabbitMQ broker information like Host's IP, and rabbitmq_queue (in my code processed_pm25_data is queue).

Connection code snippet:

```
def publish_to_rabbitmq(data):  
  
    connection = pika.BlockingConnection(pika.ConnectionParameters(host=RABBITMQ_HOST))  
  
    channel = connection.channel()  
  
    channel.queue_declare(queue=RABBITMQ_QUEUE)  
  
    for item in data:  
  
        message = json.dumps(item)  
  
        channel.basic_publish(exchange="", routing_key=RABBITMQ_QUEUE,  
body=message)  
  
        print(f"[Published] to RabbitMQ - Timestamp: {item['timestamp']}, PM2.5 Average:  
{item['value']}")  
  
    connection.close()
```

Output:

So, in this task we have created Dockerfile to migrate data preprocessor and docker-compose to run the container.

docker-compose:

```
1 version: "3"
2 services:
3   data_preprocessor:
4     build: ./subscriber
5     image: data_preprocessor:latest
6     environment:
7       rabbitmq_broker: "192.168.0.100"
8       exq_broker: "192.168.0.102"
9 networks:
10  default:
11    driver: bridge
```

Dockerfile:

```
1 FROM python:3.8.12
2
3 USER root
4
5 ADD . /usr/local/source
6
7 WORKDIR /usr/local/source
8
9 RUN pip3 install -r requirements.txt
10
11 CMD ["python3", "data_preprocessor.py"]
12
```

Code:

Now, the screenshot of entire code that performs task 2 has been added below.

data_preprocessor.py


```

1 import json
2 import paho.mqtt.client as mqtt
3 import pika
4 from datetime import datetime, timedelta
5 import time
6
7 BROKER_ADDRESS = "192.168.0.102"
8 MQTT_PORT = 1883
9 MQTT_TOPIC = "sensor/pm2.5"
10
11 RABBITMQ_QUEUE = 'processed_pm25_data'
12 RABBITMQ_HOST = '192.168.0.100'
13
14 pm25_data = []
15 outliers = []
16
17
18 def on_message(client, userdata, msg):
19     global pm25_data
20     data = json.loads(msg.payload.decode())
21     print(f"[Received] PM2.5 data - Timestamp: {data['timestamp']], Value: {data['value']]")
22
23     if data['value'] <= 50:
24         pm25_data.append(data)
25     else:
26         print(f"[Outlier] Filtered out PM2.5 data - Timestamp: {data['timestamp']], Value: {data['value']]")
27
28
29 client = mqtt.Client()
30 client.on_message = on_message
31 client.connect(BROKER_ADDRESS, MQTT_PORT)
32 client.subscribe(MQTT_TOPIC)
33 print("Subscribed to MQTT topic and listening for PM2.5 data...")
34 client.loop_start()
35
36
37 def calculate_daily_averages():
38     global pm25_data
39     if not pm25_data:
40         return []
41
42     pm25_data.sort(key=lambda x: x['timestamp'])
43
44     daily_averages = []
45     current_day = datetime.fromtimestamp(pm25_data[0]['timestamp'] / 1000).date()
46     daily_sum = 0
47     daily_count = 0
48
49     for entry in pm25_data:
50         entry_date = datetime.fromtimestamp(entry['timestamp'] / 1000).date()
51
52         if entry_date == current_day:
53             daily_sum += entry['value']
54             daily_count += 1
55         else:
56             if daily_count > 0:
57                 average = daily_sum / daily_count
58                 daily_averages.append({
59                     'timestamp': int(datetime.combine(current_day, datetime.min.time()).timestamp() * 1000),
60                     'value': average
61                 })
62                 print(f"[Daily Average] Calculated for {current_day}: {average}")
63
64             current_day = entry_date
65             daily_sum = entry['value']
66             daily_count = 1
67
68     if daily_count > 0:
69         average = daily_sum / daily_count
70         daily_averages.append({
71             'timestamp': int(datetime.combine(current_day, datetime.min.time()).timestamp() * 1000),
72             'value': average
73         })
74         print(f"[Daily Average] Calculated for {current_day}: {average}")
75
76     return daily_averages
77
78
79 def publish_to_rabbitmq(data):
80     connection = pika.BlockingConnection(pika.ConnectionParameters(host=RABBITMQ_HOST))
81     channel = connection.channel()
82     channel.queue_declare(queue=RABBITMQ_QUEUE)
83
84     for item in data:
85         message = json.dumps(item)
86         channel.basic_publish(exchange='', routing_key=RABBITMQ_QUEUE, body=message)
87         print(f"[Published] to RabbitMQ - Timestamp: {item['timestamp']], PM2.5 Average: {item['value']]")
88
89     connection.close()
90
91 try:
92     while True:
93         time.sleep(300)
94
95         daily_averages = calculate_daily_averages()
96
97         if daily_averages:
98             publish_to_rabbitmq(daily_averages)
99
100             with open("filtered_pm25_data.json", "w") as f:
101                 json.dump(pm25_data, f, indent=4)
102                 print(f"[Saved] Filtered PM2.5 data to 'filtered_pm25_data.json'.")
103
104             if outliers:
105                 print("Outliers:", outliers, flush=True)
106
107             pm25_data = []
108
109 except KeyboardInterrupt:
110     print("Stopping data processing and disconnecting from MQTT broker.")
111     client.loop_stop()
112     client.disconnect()

```

Task 3

Time series data prediction and visualization

Step 1:

Firstly, I added pre-defined machine learning engine code in cloud according to my prediction code.

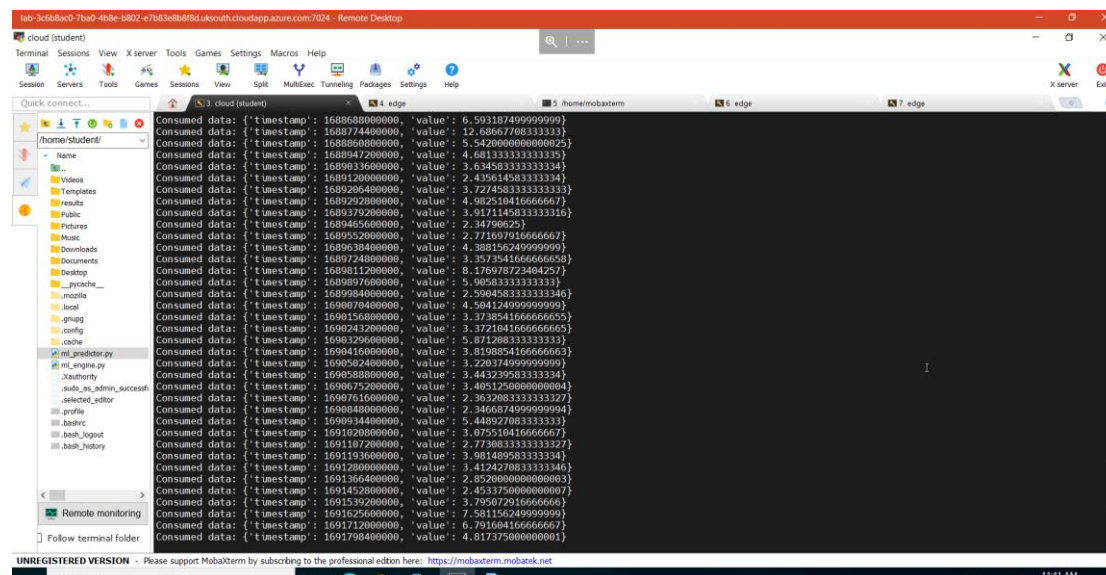
ml_engine.py

```
1 import pandas as pd
2 from prophet import Prophet
3
4 class MLPredictor(object):
5     def __init__(self, data_df):
6
7         self.__train_data = self.__convert_col_name(data_df)
8         self.__trainer = Prophet(Changepoint_prior_scale=12)
9
10    def train(self):
11        self.__trainer.fit(self.__train_data)
12
13    def __convert_col_name(self, data_df):
14
15        data_df.rename(columns={"timestamp": "ds", "value": "y"}, inplace=True)
16
17        data_df['ds'] = pd.to_datetime(data_df['ds'], unit='ms')
18
19        print(f"After renaming and converting columns:\n{data_df.head()}")
20        return data_df
21
22    def __make_future(self, periods=15):
23        future = self.__trainer.make_future_dataframe(periods=periods)
24        return future
25
26    def predict(self):
27        future = self.__make_future()
28        forecast = self.__trainer.predict(future)
29        return forecast
30
31    def plot_result(self, forecast):
32        fig = self.__trainer.plot(forecast, figsize=(15, 6))
33        return fig
34
35
```

Step 2:

Consume all averaged daily PM2.5 data that has been published by RabbitMQ service in task 2.

Output:



```
Consumed data: {'timestamp': 1688680000000, 'value': 6.593187499999999}
Consumed data: {'timestamp': 1688774400000, 'value': 12.686677083333333}
Consumed data: {'timestamp': 1688868800000, 'value': 5.5420000000000025}
Consumed data: {'timestamp': 1688963200000, 'value': 4.601333333333333}
Consumed data: {'timestamp': 1689057600000, 'value': 3.6345833333333334}
Consumed data: {'timestamp': 1689152000000, 'value': 2.4356145833333334}
Consumed data: {'timestamp': 1689246400000, 'value': 3.7274583333333333}
Consumed data: {'timestamp': 1689340800000, 'value': 4.982510416666667}
Consumed data: {'timestamp': 1689435200000, 'value': 3.9171145833333316}
Consumed data: {'timestamp': 1689529600000, 'value': 2.34706625}
Consumed data: {'timestamp': 1689624000000, 'value': 2.771697916666667}
Consumed data: {'timestamp': 1689718400000, 'value': 4.388156249999999}
Consumed data: {'timestamp': 1689812800000, 'value': 3.3573541666666665}
Consumed data: {'timestamp': 1689907200000, 'value': 8.17697072404257}
Consumed data: {'timestamp': 1689996000000, 'value': 5.980833333333333}
Consumed data: {'timestamp': 1690084800000, 'value': 2.5904583333333334}
Consumed data: {'timestamp': 1690178400000, 'value': 4.594124999999999}
Consumed data: {'timestamp': 1690272000000, 'value': 3.373084166666665}
Consumed data: {'timestamp': 1690365600000, 'value': 3.8710416666666665}
Consumed data: {'timestamp': 1690459200000, 'value': 3.871208333333333}
Consumed data: {'timestamp': 1690552800000, 'value': 3.8198854166666663}
Consumed data: {'timestamp': 1690646400000, 'value': 3.220374999999999}
Consumed data: {'timestamp': 1690740000000, 'value': 3.4432395833333334}
Consumed data: {'timestamp': 1690833600000, 'value': 3.4051250000000004}
Consumed data: {'timestamp': 1690927200000, 'value': 2.3632083333333327}
Consumed data: {'timestamp': 1691020800000, 'value': 2.3466874999999994}
Consumed data: {'timestamp': 1691114400000, 'value': 5.448927083333333}
Consumed data: {'timestamp': 1691208000000, 'value': 3.075510416666667}
Consumed data: {'timestamp': 1691301600000, 'value': 2.7730833333333327}
Consumed data: {'timestamp': 1691395200000, 'value': 3.9814895833333334}
Consumed data: {'timestamp': 1691488800000, 'value': 3.4124270833333334}
Consumed data: {'timestamp': 1691582400000, 'value': 2.0520000000000003}
Consumed data: {'timestamp': 1691676000000, 'value': 2.4533750000000007}
Consumed data: {'timestamp': 1691769600000, 'value': 3.7950729166666666}
Consumed data: {'timestamp': 1691863200000, 'value': 7.581156249999999}
Consumed data: {'timestamp': 1691956800000, 'value': 6.791604166666667}
Consumed data: {'timestamp': 1692050400000, 'value': 4.8173750000000001}
```

Step 3:

Convert timestamp to date time format.

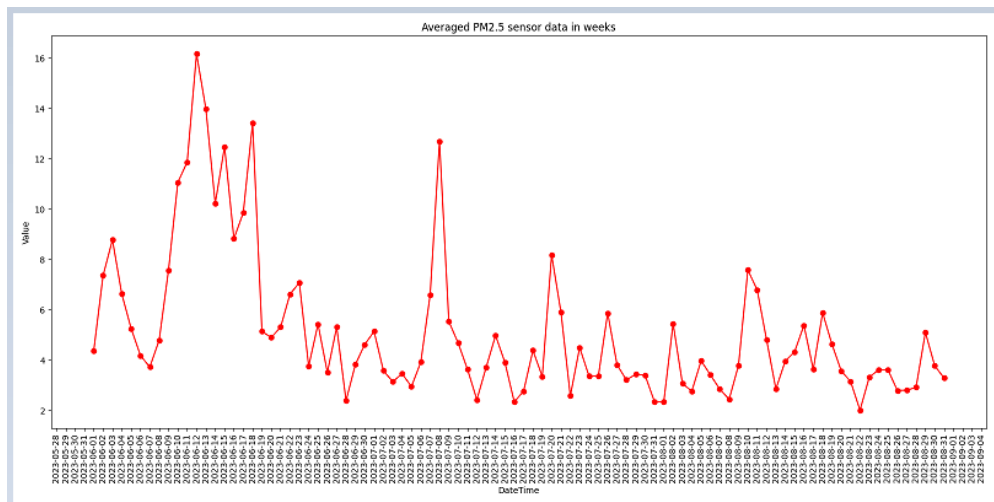
Output:

```
0 2023-06-01 4.385896
1 2023-06-02 7.366385
2 2023-06-03 8.777760
3 2023-06-04 6.630323
4 2023-06-05 5.260854
```

Step 4:

Using matplotlib, visualize the averaged PM2.5 daily data.

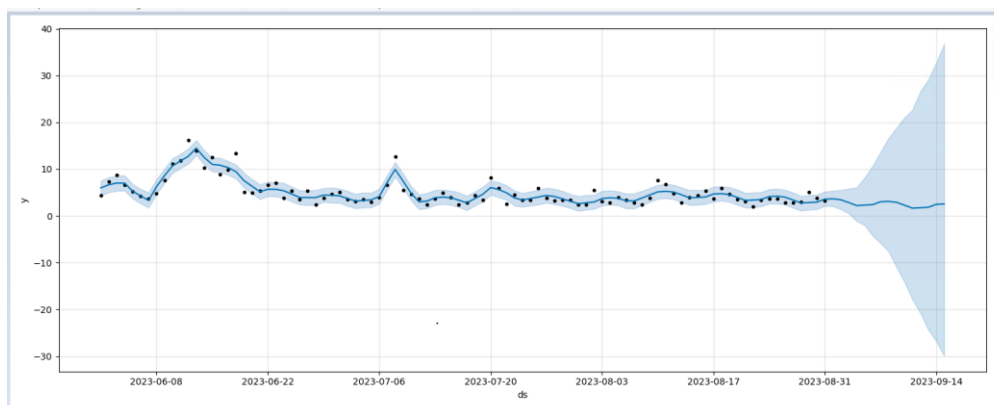
Line chart:



Step 5:

Finally, I feed the averaged data to machine learning model to predict the trend of PM2.5 data for next 15 days. Then, the prediction has been visualized and stored in the directory.

Figure:



Code for task 3:

ml_predictor.py

```

1 import matplotlib
2 matplotlib.use('Agg')
3 import pika
4 import json
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 from ml_engine import MLPredictor
8 import os
9
10 RABBITMQ_QUEUE = 'processed_pm25_data'
11 RABBITMQ_HOST = 'localhost'
12
13 def consume_from_rabbitmq():
14     connection = pika.BlockingConnection(pika.ConnectionParameters(RABBITMQ_HOST))
15     channel = connection.channel()
16     channel.queue_declare(queue=RABBITMQ_QUEUE)
17
18     pm25_data = []
19
20     def callback(ch, method, properties, body):
21         data = json.loads(body)
22         pm25_data.append(data)
23         print(f"Consumed data: {data}")
24         ch.basic_ack(delivery_tag=method.delivery_tag)
25
26     while True:
27         method_frame, _, body = channel.basic_get(queue=RABBITMQ_QUEUE, auto_ack=False)
28         if method_frame:
29             callback(channel, method_frame, None, body)
30         else:
31             break
32
33     df = pd.DataFrame(pm25_data)
34     connection.close()
35     return df
36
37 pm25_df = consume_from_rabbitmq()
38
39 if not pm25_df.empty:
40     pm25_df['timestamp'] = pd.to_datetime(pm25_df['timestamp'], unit='ms')
41
42     if not os.path.exists("results"):
43         os.makedirs("results")
44
45
46     plt.figure(figsize=(16, 8))
47     plt.plot(pm25_df['timestamp'], pm25_df['value'], color='red', marker='o', linestyle='--')
48     plt.title("Averaged PM2.5 sensor data in weeks")
49     plt.xlabel("Dateime")
50     plt.ylabel("Value")
51
52     plt.xticks(rotation=90)
53     plt.gca().xaxis.set_major_locator(plt.MultipleLocator(1))
54     plt.tight_layout()
55
56     daily_plot_path = "results/daily_pm25_averages.png"
57     plt.savefig(daily_plot_path)
58     print(f"Daily PM2.5 plot saved as '{daily_plot_path}'.")
59
60     predictor = MLPredictor(pm25_df)
61     predictor.train()
62     print("Training complete.")
63
64     forecast = predictor.predict()
65     print("Prediction complete.")
66
67     forecast_plot_path = "results/forecast_result.png"
68     fig = predictor.plot_result(forecast)
69     fig.savefig(forecast_plot_path)
70     print(f"Forecast plot saved as '{forecast_plot_path}'.")
71
72 else:
73     print("No data available for training and plotting.")
74
75

```

Cloud Container:

```

student@cloud:~$ docker ps
CONTAINER ID   IMAGE          COMMAND                  CREATED        STATUS        PORTS
f93b42f21617   rabbitmq:management   "docker-entrypoint.s..."   26 hours ago   Up 26 hours   4369/tcp, 0.0.0.0:1567→1567/tcp, :::1567→1567/tcp, 5671/tcp, 1567-1567
2/tcp, 15691-15692/tcp, 0.0.0.0:5672→5672/tcp, :::5672→5672/tcp, 25672/tcp
student@cloud:~$

```

Analytics and Conclusion:

In this IOT system, we work with three tiers, they are IoT tier, edge tier and cloud tier. From IoT tier we are accessing the real time data. Then in edge tier we are fetching the real time data and extracting the set of data be processed and, we are publishing the data to preprocessor and then the preprocessor will subscribe the data and removed outlier and finds the average and push the data to cloud tier. Now the cloud tier displays the averaged data into line chart and then using machine learning it predicts for next 15 days, then the predicted image has been stored.