

CSC8112 Report1 Boyan Li (IoT)

Student ID: 23010958 (c30109586)

Student Name: Boyan Li

During the demonstration, some problems were found in the implementation of Task2's "data preprocessing operator" part, which has been corrected now:

1. data preprocessing, which should be implemented to filter different types of data into different data sets.
2. Send data to RabbitMQ in a package instead of sending each piece of data individually.

To fix the above problems, I improved the python code in the following ways:

1. Clearer naming
2. Try to implement all logic in one main function per file to avoid unnecessary bugs caused by parameter passing
3. Correct the code logic of "data preprocessing operator" in Task2
4. Send the data to the RabbitMQ message queue in the correct way
5. Confirm the right format of 'Timestamp'.

Part A. Experimental Results Show

Content of Part A including:

1. Detailed response to each task and related sub-tasks.
2. Screenshots of running services in the Docker Environment.
3. Screenshots of Functional Code (Files)

“task1_iot_edge.py”: the data injector component;

“task2_edge_cloud.py”: the data preprocessing operator;

“task3_cloud_processor.py”: time-series data prediction and visualization;

(with task2_edge_cloud.py, there are Dockerfile, docker-compose.yml, and requirements.txt)

Task 1: Design a data injector component by leveraging Newcastle Urban Observatory IoT data streams (by using the back up url)

1. Pull and run the Docker image "emqx/emqx" from Docker Hub in the virtual machine running on Azure lab (Edge). Perform this task first using the command line interface (CLI).
 - pull the Docker image "emqx/emqx" by using: “docker pull emqx/emqx”
 - run the Docker image "emqx/emqx" by using: "docker run -d --name emqx -p 1883:1883 -p 18083:18083 emqx/emqx"

(eg. "1883:1883", the left one is for the server, and the right one is for the docker)

```
student@edge:~$ docker ps
CONTAINER ID   IMAGE       COMMAND                  CREATED             STATUS              PORTS
                                                               NAMES
086f3553c9be   emqx/emqx   "/usr/bin/docker-ent..."   2 days ago        Up 33 hours   4370/tcp, 0.0.0.0
:1883->1883/tcp, ::1883->1883/tcp, 5369/tcp, 8083-8084/tcp, 8883/tcp, 0.0.0.0:8081->8081/tcp,
:::8081->8081/tcp, 0.0.0.0:18083->18083/tcp, ::18083->18083/tcp, 11883/tcp   emqx
student@edge:~$
```

2. Develop a data injector component with the following functions (Code) in Azure Lab (Edge):

“task1_iot_edge.py”: the data injector component;

- (a) Collect data from Urban Observatory platform by sending HTTP request. Following that, please print out the raw data streams that you collected on the console.

```
# Collect raw data from Urban Observatory (IoT Layer Back up url)
url = "https://gist.github.com/ringosham/fbd66654dc53c40bd4581d2828acc94e/raw/d56a0fcfd27ff7ea31e2aec3765eb2c5d64a
response = requests.get(url)
raw_data_dict = json.loads(response.text)
print(raw_data_dict)
```

```
Variable': 'PM2.5', 'Sensor Name': 'PER_AIRMON_MONITOR1135100', 'Flagged as Suspect Reading': False,
    'Value': 3.021}, {'Timestamp': 1689182100000, 'Units': 'ugm -3', 'Variable': 'PM2.5', 'Sensor
Name': 'PER_AIRMON_MONITOR1135100', 'Flagged as Suspect Reading': False, 'Value': 3.208},
{'Timestamp': 1689183000000, 'Units': 'ugm -3', 'Variable': 'PM2.5', 'Sensor Name':
'PER_AIRMON_MONITOR1135100', 'Flagged as Suspect Reading': False, 'Value': 2.395}, {'Timestamp':
1689183900000, 'Units': 'ugm -3', 'Variable': 'PM2.5', 'Sensor Name': 'PER_AIRMON_MONITOR1135100',
'Flagged as Suspect Reading': False, 'Value': 2.074}, {'Timestamp': 1689184800000, 'Units': 'ugm
-3', 'Variable': 'PM2.5', 'Sensor Name': 'PER_AIRMON_MONITOR1135100', 'Flagged as Suspect
Reading': False, 'Value': 2.955}, {'Timestamp': 1689185700000, 'Units': 'ugm -3', 'Variable':
'PM2.5', 'Sensor Name': 'PER_AIRMON_MONITOR1135100', 'Flagged as Suspect Reading': False,
    'Value': 2.58}, {'Timestamp': 1689186600000, 'Units': 'ugm -3', 'Variable': 'PM2.5', 'Sensor Name':
'PER_AIRMON_MONITOR1135100', 'Flagged as Suspect Reading': False, 'Value': 2.169}, {'Timestamp':
1689187500000, 'Units': 'ugm -3', 'Variable': 'PM2.5', 'Sensor Name': 'PER_AIRMON_MONITOR1135100',
'Flagged as Suspect Reading': False, 'Value': 1.831}, {'Timestamp': 1689188400000, 'Units': 'ugm
-3', 'Variable': 'PM2.5', 'Sensor Name': 'PER_AIRMON_MONITOR1135100', 'Flagged as Suspect
Reading': False, 'Value': 1.839}, {'Timestamp': 1689189300000, 'Units': 'ugm -3', 'Variable':
'PM2.5', 'Sensor Name': 'PER_AIRMON_MONITOR1135100', 'Flagged as Suspect Reading': False,
    'Value': 2.134}, {'Timestamp': 1689190200000, 'Units': 'ugm -3', 'Variable': 'PM2.5', 'Sensor Name':
'PER_AIRMON_MONITOR1135100', 'Flagged as Suspect Reading': False, 'Value': 2.677}, {'Timestamp':
1689191100000, 'Units': 'ugm -3', 'Variable': 'PM2.5', 'Sensor Name': 'PER_AIRMON_MONITOR1135100',
'Flagged as Suspect Reading': False, 'Value': 2.121}, {'Timestamp': 1689192000000, 'Units': 'ugm
```

(b) Although the raw air quality data you collected from the Urban Observatory API contains many metrics including N O₂, NO, CO₂, PM2.5, and PM10, among others, for the purpose of this coursework you only need to store and analyze PM2.5 data. While many meta-data are available for PM2.5 data, such as sensor name, timestamp, value, and location, you only need to store the metrics related to the Timestamp and Value meta-data fields.

```
# Extract PM2.5 data
pm25_data = {}
for data in raw_data_dict["sensors"]:
    for items in data["data"]["PM2.5"]:
        timestamp = items['Timestamp']
        value = items['Value']
        timestamp = datetime.datetime.utcfromtimestamp(timestamp / 1000).strftime('%Y-%m-%d %H:%M:%S')
        pm25_data[timestamp] = value
print(pm25_data)
```

```

Run task1_iot_edge ×
[{"2023-06-01 00:00:00": 4.273, "2023-06-01 00:15:00": 3.483, "2023-06-01 00:30:00": 3.713,
 "2023-06-01 00:45:00": 3.813, "2023-06-01 01:00:00": 3.885, "2023-06-01 01:15:00": 3.809,
 "2023-06-01 01:30:00": 3.869, "2023-06-01 01:45:00": 3.839, "2023-06-01 02:00:00": 3.805,
 "2023-06-01 02:15:00": 4.094, "2023-06-01 02:30:00": 4.474, "2023-06-01 02:45:00": 4.918,
 "2023-06-01 03:00:00": 5.56, "2023-06-01 03:15:00": 5.637, "2023-06-01 03:30:00": 5.71,
 "2023-06-01 03:45:00": 6.004, "2023-06-01 04:00:00": 5.322, "2023-06-01 04:15:00": 4.961,
 "2023-06-01 04:30:00": 4.679, "2023-06-01 04:45:00": 4.815, "2023-06-01 05:00:00": 4.498,
 "2023-06-01 05:15:00": 4.393, "2023-06-01 05:30:00": 4.433, "2023-06-01 05:45:00": 4.467,
 "2023-06-01 06:00:00": 5.76, "2023-06-01 06:15:00": 4.549, "2023-06-01 06:30:00": 4.499,
 "2023-06-01 06:45:00": 4.946, "2023-06-01 07:00:00": 4.864, "2023-06-01 07:15:00": 6.079,
 "2023-06-01 07:30:00": 5.902, "2023-06-01 07:45:00": 6.025, "2023-06-01 08:00:00": 3.99,
 "2023-06-01 08:15:00": 3.223, "2023-06-01 08:30:00": 2.434, "2023-06-01 08:45:00": 2.07,
 "2023-06-01 09:00:00": 1.968, "2023-06-01 09:15:00": 2.497, "2023-06-01 09:30:00": 3.192,
 "2023-06-01 09:45:00": 3.166, "2023-06-01 10:00:00": 2.994, "2023-06-01 10:15:00": 3.023,
 "2023-06-01 10:30:00": 3.652, "2023-06-01 10:45:00": 2.243, "2023-06-01 11:00:00": 2.652,
 "2023-06-01 11:15:00": 3.137, "2023-06-01 11:30:00": 3.075, "2023-06-01 11:45:00": 3.376,
 "2023-06-01 12:00:00": 3.494, "2023-06-01 12:15:00": 3.639, "2023-06-01 12:30:00": 3.826,
 "2023-06-01 12:45:00": 3.888, "2023-06-01 13:00:00": 4.861, "2023-06-01 13:15:00": 5.634,
 "2023-06-01 13:30:00": 5.097, "2023-06-01 13:45:00": 5.414, "2023-06-01 14:00:00": 6.029,
 "2023-06-01 14:15:00": 5.582, "2023-06-01 14:30:00": 5.885, "2023-06-01 14:45:00": 4.228,
 "2023-06-01 15:00:00": 4.729, "2023-06-01 15:15:00": 4.843, "2023-06-01 15:30:00": 5.195,
}

```

(c) Send all PM2.5 data to be used by Task 2.2 (a) to EMQX service of Azure lab (Edge).

```

# MQTT (Publisher)
mqtt_ip = "192.168.0.102" # edge vm ip address
mqtt_port = 1883
topic = "CSC8112"

# Create a mqtt client object
client = mqtt_client.Client()

# Connect to MQTT service
client.on_connect = on_connect
client.connect(mqtt_ip, mqtt_port)

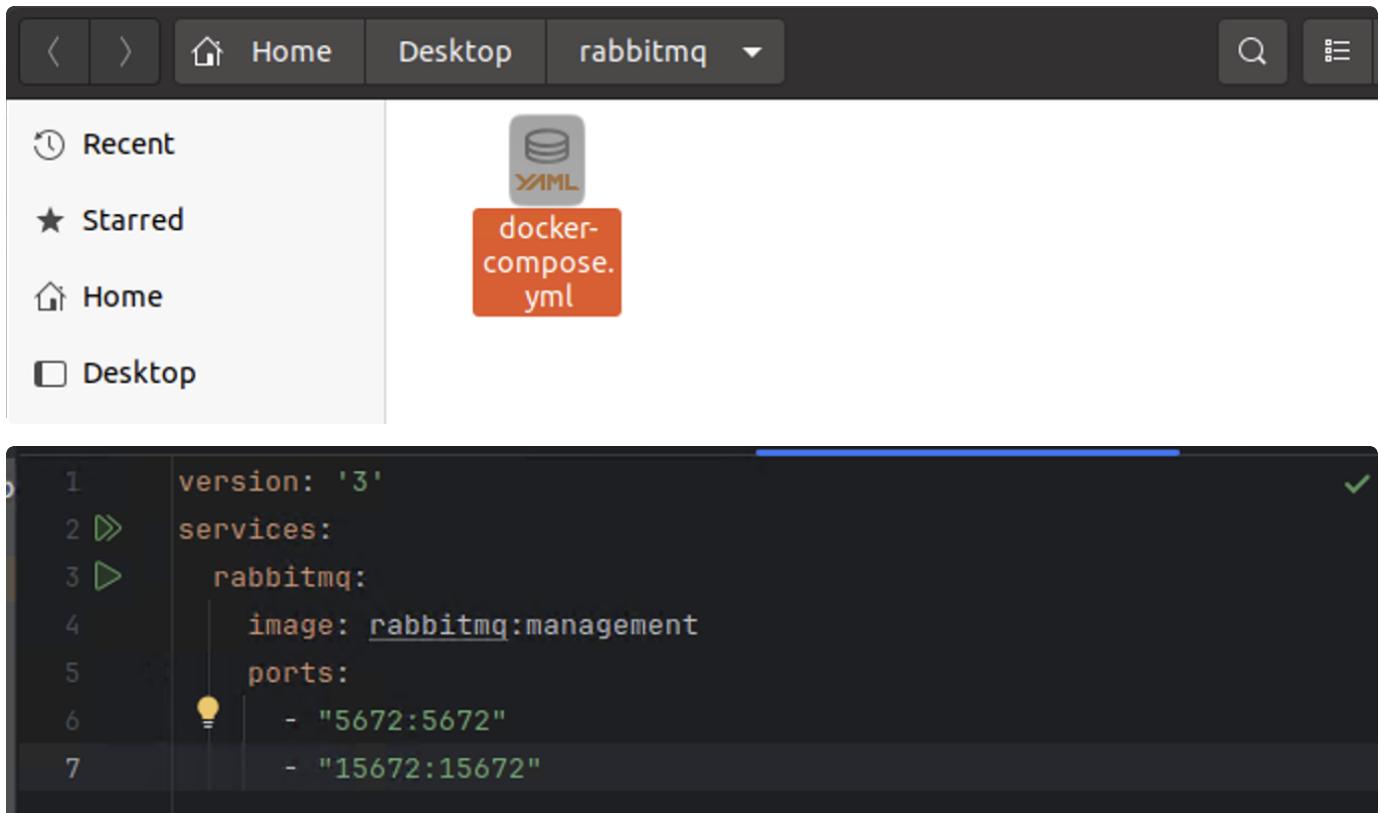
# Publish message to MQTT
# Note: MQTT payload must be a string, bytearray, int, float or None
msg = json.dumps(pm25_data)
client.publish(topic, msg)

```

Task 2: Data preprocessing operator design

1. Define a Docker compose file which contains the following necessary configurations and

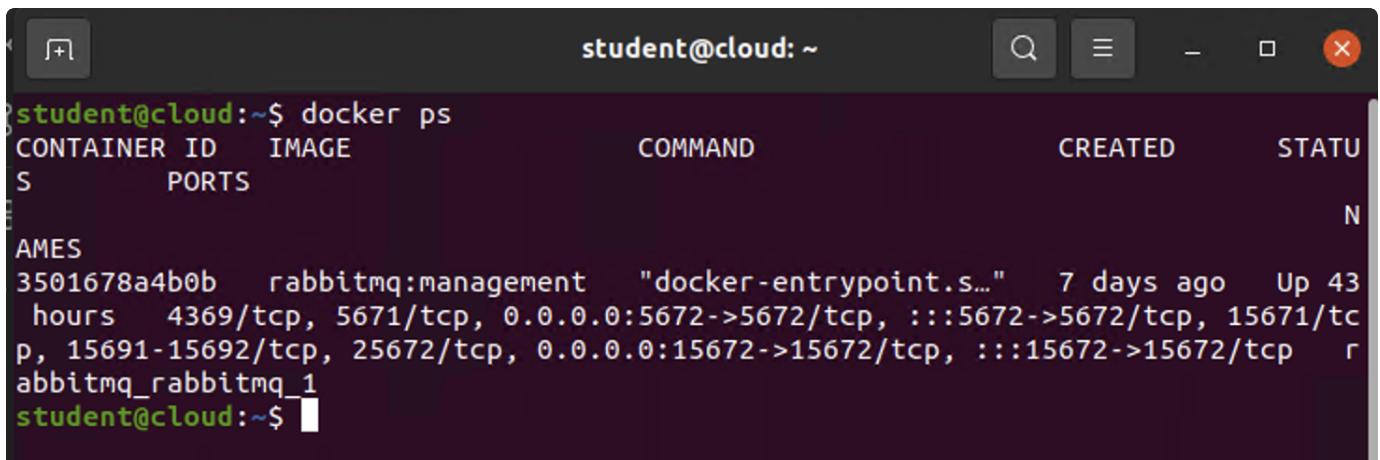
instructions for deploying and instantiating the following set of Docker images (as shown in Figure 1) on Azure lab (Cloud):



The screenshot shows a code editor interface with a dark theme. On the left is a sidebar with icons for Recent, Starred, Home, and Desktop. The main area displays a YAML file named `docker-compose.yml`. The file content is as follows:

```
version: '3'
services:
  rabbitmq:
    image: rabbitmq:management
    ports:
      - "5672:5672"
      - "15672:15672"
```

(a) Download and run RabbitMQ image (rabbitmq:management);



The screenshot shows a terminal window titled `student@cloud: ~`. The command `docker ps` was run, and the output is as follows:

```
student@cloud:~$ docker ps
CONTAINER ID        IMAGE               COMMAND                  CREATED             STATUS              PORTS
S          Ports
5301678a4b0b        rabbitmq:management   "docker-entrypoint.s..."   7 days ago         Up 43 hours
        4369/tcp, 5671/tcp, 0.0.0.0:5672->5672/tcp, :::5672->5672/tcp, 15671/tcp, 15691-15692/tcp, 25672/tcp, 0.0.0.0:15672->15672/tcp, :::15672->15672/tcp
rabbitmq_rabbitmq_1
student@cloud:~$
```

2. Design a data preprocessing operator with the following functions (code) in Azure Lab (Edge):

“task2_edge_cloud.py”: the data preprocessing operator;

- Collect all PM2.5 data published by Task 1.2 (c) from EMQX service, and please print out the PM2.5 data to the console (this operator will run as a Docker

container,
so the logs can be seen in the docker logs console automatically).

```
# MQTT (Subscriber)
# Create a mqtt client object
client = mqtt_client.Client()
# Callback function for MQTT connection
client.on_connect = on_connect
client.connect(mqtt_ip, mqtt_port)
# Subscribe MQTT topic
client.subscribe(topic)
client.on_message = on_message
# Start a thread to monitor message from publisher
client.loop_forever()

# Callback function will be triggered
def on_message(client, userdata, msg):
    print(f"Get message from publisher {json.loads(msg.payload)}")

# Collect raw PM 2.5 data from MQTT
pm25_data = json.loads(msg.payload)
print(pm25_data)
```

```

# Collect raw PM 2.5 data from MQTT
pm25_data = json.loads(msg.payload)
print(pm25_data)

# Collect and Filter PM2.5 data outliers (>50)
if __name__ == '__main__':
    on_message()

```

2023-08-28 08:45:00': 3.269, '2023-08-28 09:00:00': 3.232, '2023-08-28 09:15:00': 3.042, '2023-08-28 09:30:00': 2.902, '2023-08-28 09:45:00': 2.93, '2023-08-28 10:00:00': 2.939, '2023-08-28 10:15:00': 2.857, '2023-08-28 10:30:00': 2.745, '2023-08-28 10:45:00': 2.807, '2023-08-28 11:00:00': 2.805, '2023-08-28 11:15:00': 2.598, '2023-08-28 11:30:00': 2.424, '2023-08-28 11:45:00': 2.691, '2023-08-28 12:00:00': 2.488, '2023-08-28 12:15:00': 3.612, '2023-08-28 12:30:00': 2.349, '2023-08-28 12:45:00': 2.031, '2023-08-28 13:00:00': 2.214, '2023-08-28 13:15:00': 2.377, '2023-08-28 13:30:00': 2.551, '2023-08-28 13:45:00': 2.6, '2023-08-28 14:00:00': 2.39, '2023-08-28 14:15:00': 2.423, '2023-08-28 14:30:00': 2.447, '2023-08-28 14:45:00': 2.708, '2023-08-28 15:00:00': 2.427, '2023-08-28 15:15:00': 2.451, '2023-08-28 15:30:00': 2.393, '2023-08-28 15:45:00': 2.629, '2023-08-28 16:00:00': 2.574, '2023-08-28 16:15:00': 2.558, '2023-08-28 16:30:00': 2.497, '2023-08-28 16:45:00': 2.293, '2023-08-28 17:00:00': 2.505, '2023-08-28 17:15:00': 2.081, '2023-08-28 17:30:00': 2.099, '2023-08-28 17:45:00': 2.534, '2023-08-28 18:00:00': 2.242, '2023-08-28 18:15:00': 2.359, '2023-08-28 18:30:00': 2.852, '2023-08-28 18:45:00': 2.516, '2023-08-28 19:00:00': 2.83, '2023-08-28 19:15:00': 2.695, '2023-08-28 19:30:00': 2.373, '2023-08-28 19:45:00': 2.506, '2023-08-28 20:00:00': 2.763, '2023-08-28 20:15:00': 3.616, '2023-08-28 20:30:00': 3.901, '2023-08-28 20:45:00': 4.267, '2023-08-28 21:00:00': 4.104, '2023-08-28 21:15:00': 4.64, '2023-08-28 21:30:00': 4.788, '2023-08-28 21:45:00': 5.074, '2023-08-28 22:00:00': 4.956, '2023-08-28 22:15:00': 4.789, '2023-08-28 22:30:00': 5.063, '2023-08-28 22:45:00': 4.872, '2023-08-28 23:00:00': 4.924, '2023-08-28 23:15:00': 5.143, '2023-08-28 23:30:00': 5.25, '2023-08-28 23:45:00': 5.085, '2023-08-29 00:00:00': 5.299, '2023-08-29 00:15:00': 5.259, '2023-08-29 00:30:00': 5.246, '2023-08-29 00:45:00': 5.137, '2023-08-29 01:00:00': 5.107, '2023-08-29 01:15:00': 5.256, '2023-08-29 01:30:00': 5.092, '2023-08-29 01:45:00': 5.23, '2023-08-29 02:00:00': 5.25, '2023-08-29 02:15:00': 5.263, '2023-08-29 02:30:00'

- b. Filter out outliers (the value greater than 50), and please print out outliers to the console (docker logs console).

```

# Collect and Filter PM2.5 data outliers (>50)
pm25_outliers = []
for key, value in pm25_data.items():
    if value > 50:
        pm25_outliers.append(key)
        print(f"Timestamp: {key}, Value: {value}")
for key in pm25_outliers:
    pm25_data.pop(key) # Filter the outliers!

```

The screenshot shows the IntelliJ IDEA interface with the following details:

- Title Bar:** IntelliJ IDEA Community Edition - 24 Nov 12:12
- Project Structure:** UrbanAirMonitor (~/), showing .idea, venv, docker-compose, Dockerfile, ml_engine.py, and README.md.
- Code Editor:** task1_iot_edge.py (active tab) and task2_edge_cloud.py are open. task2_edge_cloud.py contains code for collecting PM2.5 data from MQTT and filtering outliers. A terminal window shows the execution of task2_edge_cloud.py.
- Terminal Output:** Log entries from task2_edge_cloud.py, including timestamped PM2.5 values and Docker logs for task2_edge_cloud_1.

```
# Collect raw PM 2.5 data from MQTT
pm25_data = json.loads(msg.payload)
print(pm25_data)

# Collect and Filter PM2.5 data outliers (>50)
if __name__ == '__main__':
    on_message()

-30 18:30:00': 4.719, '2023-08-30 18:45:00': 3.298, '2023-08-30 19:00:00': 3.936, '2023-08-30 19:15:00': 4.692, '2023-08-30 19:30:00': 4.193, '2023-08-30 19:45:00': 4.923, '2023-08-30 20:00:00': 4.787, '2023-08-30 20:15:00': 6.005, '2023-08-30 20:30:00': 5.039, '2023-08-30 20:45:00': 5.225, '2023-08-30 21:00:00': 5.994, '2023-08-30 21:15:00': 6.484, '2023-08-30 21:30:00': 5.442, '2023-08-30 21:45:00': 8.705, '2023-08-30 22:00:00': 3.995, '2023-08-30 22:15:00': 3.992, '2023-08-30 22:30:00': 7.768, '2023-08-30 22:45:00': 7.387, '2023-08-30 23:00:00': 6.199, '2023-08-30 23:15:00': 5.769, '2023-08-30 23:30:00': 3.641, '2023-08-30 23:45:00': 2.925, '2023-08-31 00:00:00': 3.295}
task2_edge_cloud_1 | Timestamp: 2023-06-15 03:45:00, Value: 52.17
task2_edge_cloud_1 | Timestamp: 2023-06-15 04:00:00, Value: 69.76
task2_edge_cloud_1 | Timestamp: 2023-06-15 04:15:00, Value: 76.49
task2_edge_cloud_1 | Timestamp: 2023-06-15 04:30:00, Value: 67.01
task2_edge_cloud_1 | Timestamp: 2023-06-15 04:45:00, Value: 59.91
```

c. Since the original PM2.5 data readings are collected every 15 mins, so please implement

ment a python code to calculate the averaging value of PM2.5 data on daily basis (every 24 hours), please pick the first start date of a day or a 24 hours interval as the new timestamp of averaged PM2.5 data, and please print out the result to the console (docker logs console)

```

# Collect valid daily data
daily_data = {}
for key, value in pm25_data.items(): # Outliers already filtered.
    key = datetime.strptime(key, '%Y-%m-%d %H:%M:%S')
    # Converts the timestamp to a date object, extracting only the date part(year, month, and day).
    date = key.date()
    if date not in daily_data:
        daily_data[date] = {'Timestamp': key, 'Value': [], 'Average': None}
    daily_data[date]['Value'].append(value)

# Calculate the daily average and use the first datetime of each day as the new timestamp
average_data = {}
for date, day_data in daily_data.items():
    day_data['Average_Value'] = sum(day_data['Value']) / len(day_data['Value'])
    day_data['Timestamp'] = datetime(date.year, date.month, date.day)
    average_data[f"{day_data['Timestamp']}"] = day_data['Average_Value']
    print(f"date:{day_data['Timestamp']} 24h average:{day_data['Average_Value']:.3f}")

```

The screenshot shows the IntelliJ IDEA interface with the following details:

- Project:** UrbanAirMonitor
- File:** task2_edge_cloud.py
- Code Editor:** Displays the Python script for collecting PM2.5 data and calculating daily averages.
- Terminal:** Shows the execution of the script and its output, which lists daily average PM2.5 values for Task 2 over several days.

```

task2_edge_cloud_1 | date:2023-06-02 00:00:00 24h average:7.366
task2_edge_cloud_1 | date:2023-06-03 00:00:00 24h average:8.778
task2_edge_cloud_1 | date:2023-06-04 00:00:00 24h average:6.630
task2_edge_cloud_1 | date:2023-06-05 00:00:00 24h average:5.261
task2_edge_cloud_1 | date:2023-06-06 00:00:00 24h average:4.187
task2_edge_cloud_1 | date:2023-06-07 00:00:00 24h average:3.730
task2_edge_cloud_1 | date:2023-06-08 00:00:00 24h average:4.797
task2_edge_cloud_1 | date:2023-06-09 00:00:00 24h average:7.573
task2_edge_cloud_1 | date:2023-06-10 00:00:00 24h average:11.063
task2_edge_cloud_1 | date:2023-06-11 00:00:00 24h average:11.868
task2_edge_cloud_1 | date:2023-06-12 00:00:00 24h average:16.162
task2_edge_cloud_1 | date:2023-06-13 00:00:00 24h average:13.977
task2_edge_cloud_1 | date:2023-06-14 00:00:00 24h average:10.226
task2_edge_cloud_1 | date:2023-06-15 00:00:00 24h average:12.477
task2_edge_cloud_1 | date:2023-06-16 00:00:00 24h average:8.826
task2_edge_cloud_1 | date:2023-06-17 00:00:00 24h average:9.853
task2_edge_cloud_1 | date:2023-06-18 00:00:00 24h average:13.412
task2_edge_cloud_1 | date:2023-06-19 00:00:00 24h average:5.143
task2_edge_cloud_1 | date:2023-06-20 00:00:00 24h average:4.903
task2_edge_cloud_1 | date:2023-06-21 00:00:00 24h average:5.317

```

- d. Transfer all results (averaged PM2.5 data) to be used by Task 3.2 (a) into RabbitMQ

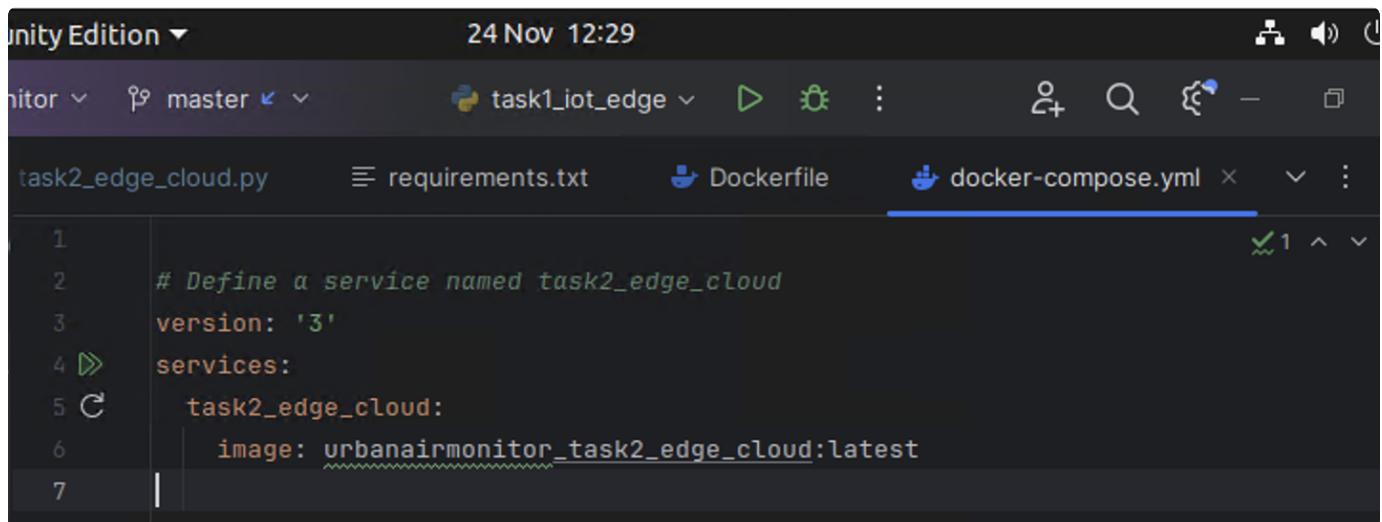
service on Azure lab (Cloud).

```
# RabbitMQ (Producer)
rabbitmq_ip = "192.168.0.100"
rabbitmq_port = 5672
# Queue name
rabbitmq_queue = "CSC8112"
# Connect to RabbitMQ service
connection = pika.BlockingConnection(pika.ConnectionParameters(host=rabbitmq_ip, port=rabbitmq_port))
channel = connection.channel()
# Declare a queue
channel.queue_declare(queue=rabbitmq_queue)
# Produce message
channel.basic_publish(exchange='',
                      routing_key=rabbitmq_queue,
                      body=json.dumps(average_data))
connection.close()
```

3. Define a Dockerfile to migrate your "data preprocessing operator" source code into a Docker image and then define a docker-compose file to run it as a container locally on the Azure

lab (Edge).

· *docker-compose.yml*



A screenshot of a terminal window titled "Unity Edition" showing the contents of a "task1_iot_edge" directory. The current file is "docker-compose.yml". The terminal shows the following code:

```
task2_edge_cloud.py      requirements.txt      Dockerfile      docker-compose.yml
1
2      # Define a service named task2_edge_cloud
3      version: '3'
4      services:
5          task2_edge_cloud:
6              image: urbanairmonitor_task2_edge_cloud:latest
7
```

· *Dockerfile*

A screenshot of a code editor window titled "Unity Edition". The status bar shows the date and time as "24 Nov 12:31". The editor tabs include "task1_iot_edge", "Dockerfile" (which is currently selected), and "docker-compose.yml". The code in the Dockerfile is as follows:

```
1 # Use an official Python runtime as a parent image
2 FROM python:3.8.10
3 USER root
4 RUN mkdir -p /usr/local/source
5 ADD ./task2_edge_cloud.py /usr/local/source
6 WORKDIR /usr/local/source
7
8 # Install any needed packages specified in requirements.txt
9 RUN pip install --no-cache-dir -r requirements.txt
10
11 # Run task_2.py when the container launches
12 CMD python3 ./task2_edge_cloud.py
```

• *requirements.txt*

A screenshot of a code editor window titled "Unity Edition". The status bar shows the date and time as "24 Nov 12:32". The editor tabs include "task1_iot_edge", "requirements.txt" (which is currently selected), and "Dockerfile". The code in the requirements.txt file is as follows:

```
1 paho-mqtt~=1.6.1
2 pika~=1.3.2
3 pandas~=2.0.3
4 matplotlib~=3.7.4
5 prophet~=1.1.5
```

Task 3: Time-series data prediction and visualization

1. Download a pre-defined Machine Learning (ML) engine code from [https://github.com/ncl-iot-team/CSC8112_MLEngine]. And put it in the same directory with `task3_cloud_operator.py`.

The screenshot shows the IntelliJ IDEA interface with the following details:

- Title Bar:** Cloud on LAB000001 - Virtual Machine Connection
- Menu Bar:** File, Action, Media, Clipboard, View, Help
- Toolbar:** Includes icons for file operations like Open, Save, and Run.
- Activity Bar:** Shows "Activities" and "IntelliJ IDEA Community Edition".
- Project Bar:** CSC8112-IoT (selected), master branch, task3, m_Engine.py (current file).
- Code Editor:** Displays Python code for `m_Engine.py`. The code defines a class with methods for training, making predictions, and visualizing data.

```
def __init__(self, data_df):
    """
    :param data_df: Dataframe type dataset
    """
    self.__train_data = self.__convert_col_name(data_df)
    self.__trainer = Prophet(changepoint_prior_scale=12)

    3 usages ▲ Boyan
    def train(self):
        self.__trainer.fit(self.__train_data)

    1 usage ▲ Boyan
    def __convert_col_name(self, data_df):
        data_df.rename(columns={"Timestamp": "ds", "Value": "y"}, inplace=True)
        print(f"After rename columns \n{data_df.columns}")
        return data_df

    1 usage ▲ Boyan
    def __make_future(self, periods=15):
        future = self.__trainer.make_future_dataframe(periods=periods)
        return future

    3 usages ▲ Boyan
    def predict(self):
        future = self.__make_future()
        forecast = self.__trainer.predict(future)
        return forecast
```

- Status Bar:** CSC8112-IoT > m_Engine.py, 1:1 LF UTF-8 4 spaces.

2. Design a PM2.5 prediction operator with the following functions (code) in Azure Lab (Cloud) or the Azure Lab localhost:

"task3_cloud_processor.py": time-series data prediction and visualization;

- a. Collect all averaged daily PM2.5 data computed by Task 2.2 (d) from RabbitMQ service, and please print out them to the console.
- b. Convert timestamp to date time format (year–month–day hour:minute:second), and please print out the PM2.5 data with the reformatted timestamp to the console.

```

def callback(ch, method, properties, body):
    print(f"Got message from producer msg: {json.loads(body)}")
    average_data = json.loads(body)
    Timestamp = []
    Value = []
    for key, value in average_data.items():
        average_value = int(value)
        Value.append(average_value)
        Timestamp.append(f"{key}")
    data = {
        'Timestamp': Timestamp,
        'Value': Value
    }
    print(data)

```

There are two lines of data below in the console. The first line is for "a." and the second line is for "b".

```

task3_cloud_operator x
:
/home/student/IdeaProjects/CSC8112-IoT/venv/bin/python /home/student/IdeaProjects/CSC8112-IoT/task3_cloud_operator.py -b
Got message from producer msg: {'2023-06-01 00:00:00': 4.3858958333333334, '2023-06-02 00:00:00': 7.366333333333333, '2023-06-03 00:00:00': 5.100291666666666, '2023-06-04 00:00:00': 3.7789042553191488, '2023-06-05 00:00:00': 3.295}
task3_cloud_operator x
:
{'2023-07-18 00:00:00': 4.388156249999999, '2023-07-19 00:00:00': 3.3573541666666658, '2023-07-20 00:00:00': 8.125, '2023-07-21 00:00:00': 5.100291666666666, '2023-07-22 00:00:00': 3.7789042553191488, '2023-07-23 00:00:00': 3.295, '2023-07-24 00:00:00': 4.388156249999999, '2023-07-25 00:00:00': 3.3573541666666658, '2023-07-26 00:00:00': 8.125, '2023-07-27 00:00:00': 5.100291666666666, '2023-07-28 00:00:00': 3.7789042553191488, '2023-07-29 00:00:00': 3.295, '2023-07-30 00:00:00': 4.388156249999999, '2023-07-31 00:00:00': 3.3573541666666658, '2023-08-01 00:00:00': 8.125, '2023-08-02 00:00:00': 5.100291666666666, '2023-08-03 00:00:00': 3.7789042553191488, '2023-08-04 00:00:00': 3.295}
task3_cloud_operator x
:

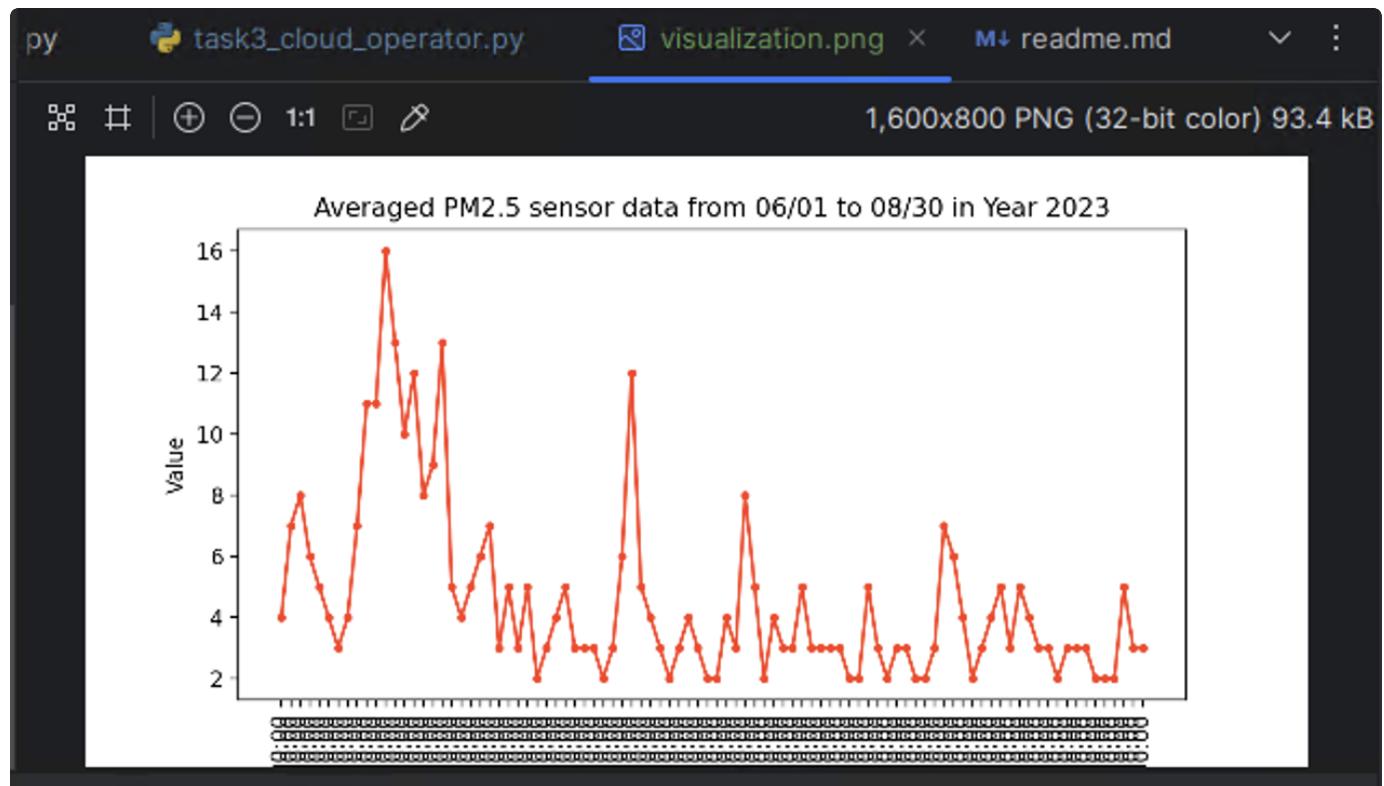
```

- c. Use the line chart component of matplotlib to visualize averaged PM2.5 daily data, directly display the figure or save it as a file.

```

# Visualization
data_df = pd.DataFrame(data)
# Initialize a canvas
plt.figure(figsize=(8, 4), dpi=200)
# Plot data into canvas
plt.plot(data_df["Timestamp"], data_df["Value"], color="#FF3B1D", marker='.', linestyle="--")
plt.title("Averaged PM2.5 sensor data from 06/01 to 08/30 in Year 2023")
plt.xlabel("DateTime")
plt.ylabel("Value")
# Rotate the scale label to display vertically
plt.xticks(Timestamp, rotation='vertical')
# Save as file
plt.savefig("visualization.png")

```



- d. Feed averaged PM2.5 data to machine learning model to predict the trend of PM2.5 for the next 15 days (this predicted time period is a default setting of provided machine learning predictor/classifier model).

```

# Format Timestamp to %Y-%m-%d
formatted_Timestamp = [datetime.strptime(timestamp, '%Y-%m-%d %H:%M:%S').strftime('%Y-%m-%d') for timestamp in Timestamp]
print(formatted_Timestamp)

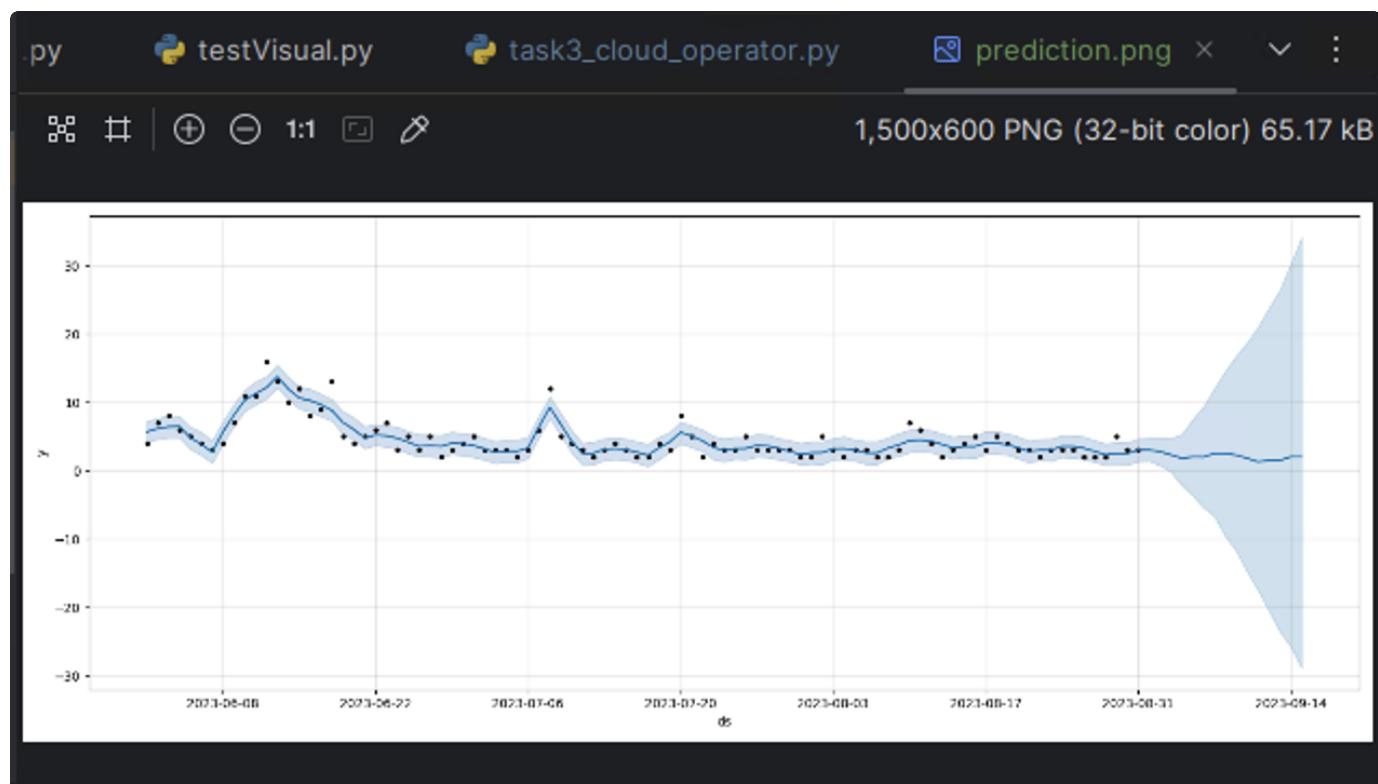
# Prepare data
data = {
    'Timestamp': formatted_Timestamp,
    'Value': Value
}
data_df = pd.DataFrame(data)

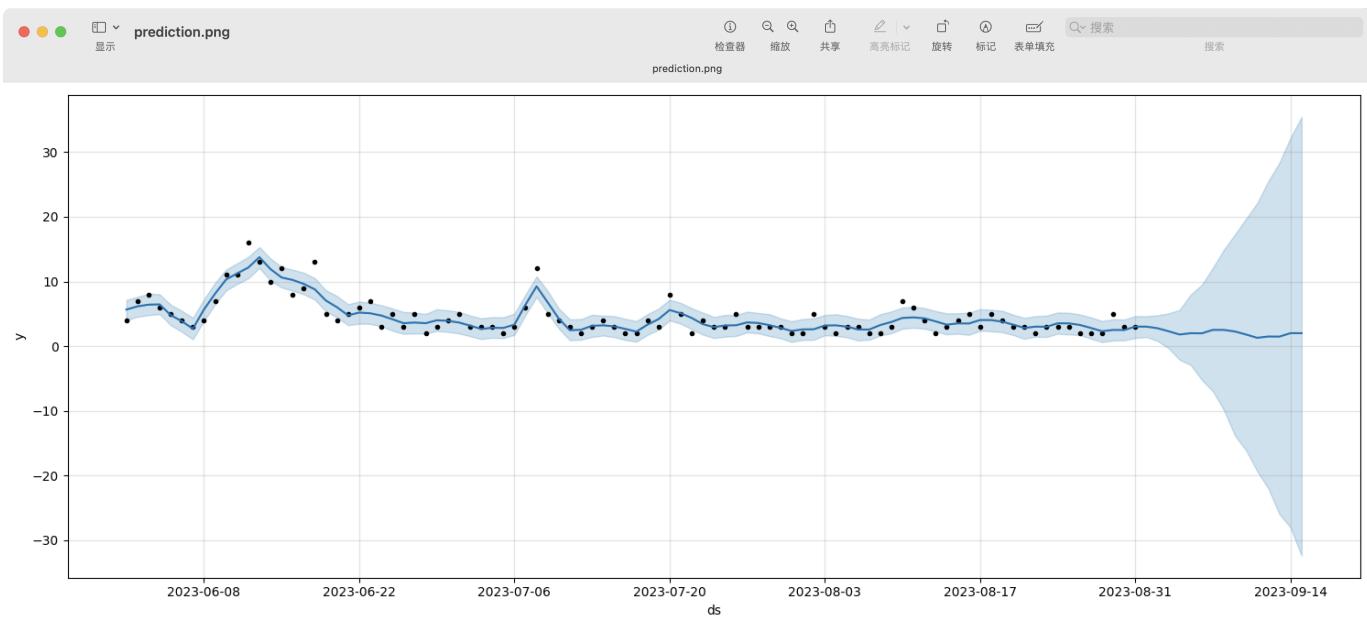
# Create ML engine predictor object
predictor = MLPredictor(data_df)
# Train ML model
predictor.train()
# Do prediction
forecast = predictor.predict()

# Get canvas
fig = predictor.plot_result(forecast)
fig.savefig("prediction.png")
fig.show()

```

- e. Visualize predicted results from Machine Learning predictor/classifier model, directly display the figure or save as it a file (pre-defined in the provided Machine Learning code)





Part B: Analytical discussion of the results and related conclusions.

PM2.5 refers to particulate matter with a diameter of 2.5 microns or less. These particles can penetrate deep into the lungs and even enter the bloodstream, posing a serious threat to human health. PM2.5 can be produced from various sources, including vehicle exhaust, industrial emissions, volcanic eruptions, and dust storms. Due to the significant health concerns associated with PM2.5, forecasting its trends is crucial for issuing health recommendations and taking preventive measures for at-risk populations.

The input dataset has two attributes: a Timestamp and a Value column, as observed from "prediction.png." The axis labeled "ds" corresponds to the "Timestamp," and "y" corresponds to the "Value," where "ds" represents the date/time and "y" is the metric being predicted. In the figure, black dots represent historical data points, the dark blue line indicates the predicted trend, and the light blue shaded area denotes the uncertainty interval.

It is apparent that the prediction model detects change points, indicating significant shifts in the trend of PM2.5 levels at positions corresponding to "2023-06-10" and "2023-07-07." These shifts could be attributed to various factors. The historical input data show that PM2.5 values have remained between 0 and 10, except at the change points. The forecast for the next 15 days suggests that PM2.5 levels are expected to decrease slightly in the short term and then remain relatively stable towards the end of the forecast period.

However, due to the small size of the input dataset and the presence of many unpredictable factors, the uncertainty in the predictions increases rapidly for data points further out in time.