

Cloud-based Data Analytics

MIE1628

Assignment 5

MIE1628

Assignment on Azure Cloud Platform

Due Date: Aug 15, 2025

- *No submission is accepted via email. We have no exception.*
- *Students are responsible for submitting the correct files on time.*
- *Assignments submitted up to 48 hours late will incur a 20% penalty.*
- *Your grade will be zero if you submit your answer after 48 hours.*

*** Instructions ***

1. Note:

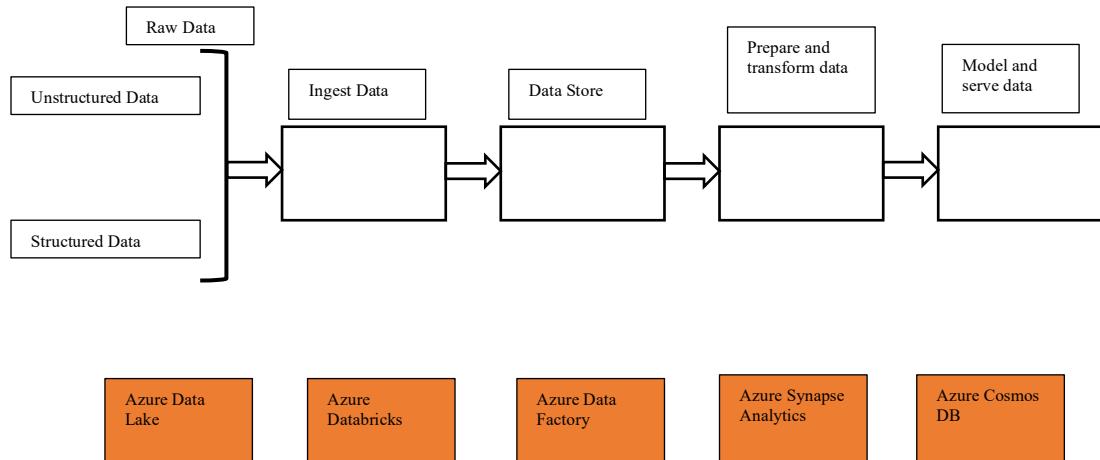
Part B of this assignment can be done in a group of two students or individually. Both students need to submit the assignment for both parts and provide both names, email, and student IDs at the top of the assignment.

Submit your complete project, including your Python Notebook with markdown explanations, and a comprehensive PDF document. The PDF should contain clear screenshots of input/output commands with results, images of your deployed Azure portal resources, detailed step-by-step explanations for each process, and final output screenshots. Additionally, include a section in the PDF answering the provided questions (to be specified separately).

Contact your TA for any questions related to this assignment or post clarification questions to the Piazza platform.

PART A:

1. [Marks: 5] Explain below the 5 components shown in orange boxes. Explain which Azure components you will use where in this big data architecture and why.



Cloud-based Data Analytics

MIE1628

Assignment 5

2. [Marks: 5] Explain how Stream Analytics works in Azure. Mention at least two common use cases or applications for this service.
3. [Marks: 10] Deploy all the resources in Azure Portal. Implement a Stream Analytics job by using the Azure portal. See this for reference - <https://learn.microsoft.com/en-us/azure/stream-analytics/stream-analytics-quick-create-portal>

For query use below:

```
SELECT *
INTO BlobOutput
FROM IoTHubInput
HAVING Temperature > 25
```

See the below screenshot and show the top 30 results for your output.

The screenshot shows the Azure Storage Blob container 'container1'. The file '0_ee2fce9a554748648a45de288e6c623f_1.json' is selected. The JSON content is as follows:

```

1 {"messageId": "775", "deviceId": "Raspberry Pi Web Client", "temperature": 27.56286555465898, "humidity": 77.9
2 {"messageId": "776", "deviceId": "Raspberry Pi Web Client", "temperature": 31.76983389970918, "humidity": 64.1
3 {"messageId": "777", "deviceId": "Raspberry Pi Web Client", "temperature": 30.173582241045128, "humidity": 61
4 {"messageId": "778", "deviceId": "Raspberry Pi Web Client", "temperature": 29.70936239344644, "humidity": 70
5 {"messageId": "779", "deviceId": "Raspberry Pi Web Client", "temperature": 29.537158745343632, "humidity": 73
6 {"messageId": "780", "deviceId": "Raspberry Pi Web Client", "temperature": 30.11726573198575, "humidity": 66
7 {"messageId": "781", "deviceId": "Raspberry Pi Web Client", "temperature": 27.237104885232631, "humidity": 63
8 {"messageId": "782", "deviceId": "Raspberry Pi Web Client", "temperature": 30.541928098495646, "humidity": 66
9 {"messageId": "783", "deviceId": "Raspberry Pi Web Client", "temperature": 30.46121575922222, "humidity": 61
10 {"messageId": "784", "deviceId": "Raspberry Pi Web Client", "temperature": 31.687045217662682, "humidity": 75
11 {"messageId": "785", "deviceId": "Raspberry Pi Web Client", "temperature": 29.006529012500579, "humidity": 70
12 {"messageId": "786", "deviceId": "Raspberry Pi Web Client", "temperature": 28.859692581436892, "humidity": 69
13 {"messageId": "787", "deviceId": "Raspberry Pi Web Client", "temperature": 30.702890613843248, "humidity": 71
14 {"messageId": "788", "deviceId": "Raspberry Pi Web Client", "temperature": 29.466947493481154, "humidity": 63
15 {"messageId": "789", "deviceId": "Raspberry Pi Web Client", "temperature": 31.118801842064, "humidity": 64.0
16 {"messageId": "790", "deviceId": "Raspberry Pi Web Client", "temperature": 27.136869152806462, "humidity": 64
17 {"messageId": "791", "deviceId": "Raspberry Pi Web Client", "temperature": 30.46121575922222, "humidity": 61
18 {"messageId": "792", "deviceId": "Raspberry Pi Web Client", "temperature": 31.687045217662682, "humidity": 75
19 {"messageId": "793", "deviceId": "Raspberry Pi Web Client", "temperature": 29.006529012500579, "humidity": 70
20 {"messageId": "794", "deviceId": "Raspberry Pi Web Client", "temperature": 28.859692581436892, "humidity": 69
21 {"messageId": "795", "deviceId": "Raspberry Pi Web Client", "temperature": 30.702890613843248, "humidity": 62
22 {"messageId": "796", "deviceId": "Raspberry Pi Web Client", "temperature": 27.860822670967593, "humidity": 66
23 {"messageId": "797", "deviceId": "Raspberry Pi Web Client", "temperature": 27.775068973730495, "humidity": 60
24 {"messageId": "798", "deviceId": "Raspberry Pi Web Client", "temperature": 30.5594159223586, "humidity": 69.1
25 {"messageId": "799", "deviceId": "Raspberry Pi Web Client", "temperature": 30.96622244545926, "humidity": 64
26 {"messageId": "800", "deviceId": "Raspberry Pi Web Client", "temperature": 29.4118621360886513782, "humidity": 69
27 {"messageId": "801", "deviceId": "Raspberry Pi Web Client", "temperature": 30.84046401525703, "humidity": 64
28 {"messageId": "802", "deviceId": "Raspberry Pi Web Client", "temperature": 27.634944188044331, "humidity": 62
29 {"messageId": "803", "deviceId": "Raspberry Pi Web Client", "temperature": 31.171862047724645, "humidity": 66
30 {"messageId": "804", "deviceId": "Raspberry Pi Web Client", "temperature": 31.748161497954346, "humidity": 64

```

Part B:

Data Input: Claim a dataset from Piazza - link. If the dataset is too large, you can take a subset of the data as well. No two groups can have the same dataset.

Your selected dataset should meet the following criteria:

1. It must contain a minimum of 1,000 instances (rows or data points).
2. It should include at least six features (columns or attributes).

Using this dataset, you are required to address a substantial and meaningful problem. Your analysis should demonstrate:

1. A clear understanding of the dataset's context and potential applications.
2. The ability to formulate relevant questions or hypotheses based on the data.
3. Appropriate use of data analysis techniques to extract insights.

Cloud-based Data Analytics

MIE1628

Assignment 5

4. The capacity to draw meaningful conclusions that could inform decision-making or further research.

Some problems to consider:

1. Fraud Detection System
2. Customer Churn Rate Prediction
3. Segmentation using Clustering
4. Recommendations with your Dataset
5. Sales Forecasting
6. Stock Price Predictions
7. Human Activity Recognition with Smartphones
8. Wine Quality Predictions
9. Breast Cancer Prediction
10. Sorting of Specific Tweets on Twitter etc.

Implement this part in Azure Machine learning using Azure Notebook

1. [Marks: 15] Clearly define the problem you intend to address using this dataset. Present a comprehensive problem statement that includes:
 - a. A detailed description of the meaningful issue you're tackling
 - b. An outline of all necessary steps, including:
 - i. Data preprocessing
 - ii. Data cleaning
 - iii. Modeling approach

Your problem statement should be thorough, spanning approximately half to one full page. If you determine that data cleaning is unnecessary, please provide a justification for why this dataset doesn't require cleaning. In such a case, allocate more attention to other crucial aspects such as EDA and the modeling process.

Ensure your problem statement is well-structured, coherent, and provides a clear roadmap for your data analysis project.

2. [Marks: 10] Explore your dataset and provide at least 5 meaningful charts/graphs with an explanation.
3. [Marks: 10] Do data cleaning/pre-processing as required and explain what you have done for your dataset and why?
4. [Marks: 15] Implement 2 machine learning models and explain which algorithms you have selected and why. Compare them and show success metrics (Accuracy/RMSE/Confusion Matrix) as per your problem. Explain results.
5. [Marks: 15] Deploy a run-time pipeline for your dataset using Azure Designer Studio.
Or
Do hyperparameter tuning for your algorithms. Explain your results.
Or
Use Automated ML for your data set. Explain the best model results.

Cloud-based Data Analytics
MIE1628
Assignment 5

6. [Marks: 15] Summarize your project's key findings and overall conclusions in a brief paragraph. Ensure your summary is firmly grounded in the data and analysis you've presented throughout your project. Offer meaningful insights that not only encapsulate your work but also lay a foundation for potential future research in this area. Your conclusions should be well-reasoned and directly supported by your results.