# Big Data Assignment-5

**TEAMMATES:**
Rahavi Selvarajan - 1007346445
rahavi.selvarajan@mail.utoronto.ca
Hruday Vishal Kanna Anand - 1006874517
vishal.kanna@mail.utoronto.ca

## Part-a

1.

Azure data lake- It is the landing zone for all types of raw data (structured, semi structured, unstructured). It can also be used as our primary storage for all the required data.

Azure Databricks- It is databricks hosted on the azure cloud. It uses spark cluster to transform data in the form of RDDs and data frames.

Azure data factory- is used to transfer data from one service to another or one location to another with the help of linked services.

Azure synapse analytics- is a unified data warehouse platform built on top of azure data lake which takes care of ingestion of data and the analysis of it with the help of SQL or Spark engine.

Azure cosmos db- is a no-SQL data warehouse that is globally distributed and has very low latency.

Ingest data- azure data factory, azure synapse analytics

Data store- azure data lake, Azure cosmos db (noSQL data)

Prepare and transform data- azure databricks, azure synapse analytics

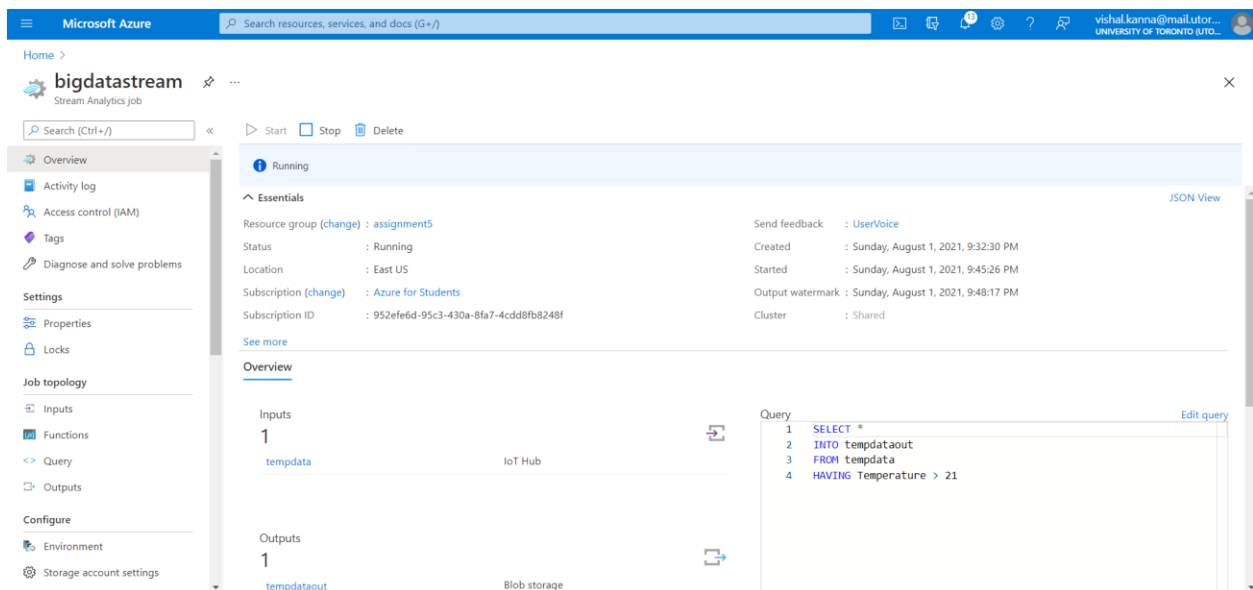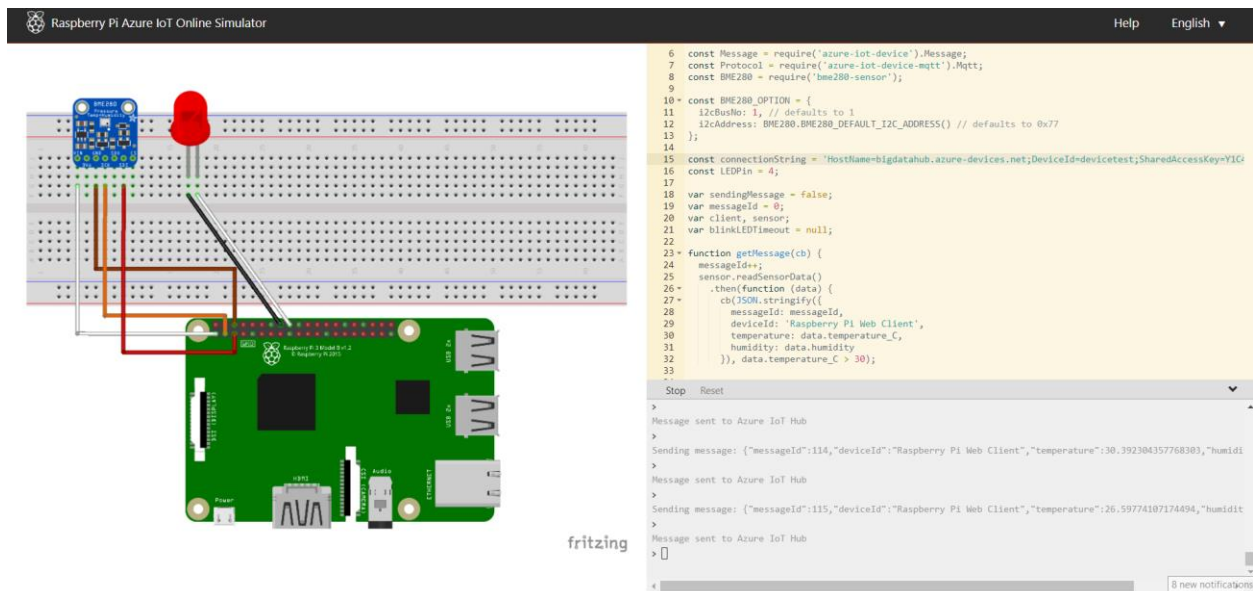Model and serve data- azure cosmos db, azure synapse analytics

2.

Azure stream analytics has three major steps-

1. Ingest/ input
2. processed
3. Store/ output

In the first step stream analytics ingests data with the help azure event hubs, azure iot hub, and azure blob storage in a continuous manner. This data is then processed; this can be simple SQL queries or

more complex ML services. This processed data might have several types of outputs; it might be in the form of alerts/ actions which will be sent to azure functions,  service bus, event hub. It may be used for real time dashboarding; then the data is sent to power BI. The prcessed data may also be stored in a warehouse or normal storage such as data lake or cosmos db.

3.

## Azure services

Create a resource · Resource groups · Cost Management ... · Data factories · Reservations · Quickstart Center · Virtual machines · App Services · Storage accounts · More services

## Recent resources

| Name | Type | Last Viewed |
| --- | --- | --- |
| bigdatastream | Stream Analytics job | 4 minutes ago |
| bigdatahub | IoT Hub | 7 minutes ago |
| assignment5 | Resource group | 18 minutes ago |
| vishalstore | Storage account | 19 minutes ago |
| Azure for Students | Subscription | 2 weeks ago |

---

...

Move ∨    Delete    Refresh

### IoT Hub Usage

- Messages used today: 3
- Daily messages quota: 8000 ⓘ
- IoT Devices: 1

### Number of messages used

100
80
60
40
20
0
9 PM    9:15 PM    9:30 PM    UTC+05:30

Total number of messages used (Max)
bigdatahub
95

---

Home > vishalstore > output >

## output
Container

Search (Ctrl+/)

- Overview
- Diagnose and solve problems
- Access Control (IAM)

### Settings

- Shared access tokens
- Access policy
- Properties
- Metadata

Upload    Change access level    ···

Authentication method: Access key (Switch to Azure AD User Account)
Location: output

Search blobs by prefix (case-...
Show deleted blobs
Add filter

Name
0_79bcdd7d72004fc2a14970e...    ···

### 0_79bcdd7d72004fc2a14970edc1177bf7_1.json
Blob

Save    Discard    Download    Refresh    Delete

Overview    Versions    Snapshots    Edit    Generate SAS

```
1   {"messageId":42,"deviceId":"Raspberry Pi Web Client","temperature":22.486323058803681,"humidity":77.7594454...
2   {"messageId":43,"deviceId":"Raspberry Pi Web Client","temperature":21.929799908042948,"humidity":73.4810844...
3   {"messageId":44,"deviceId":"Raspberry Pi Web Client","temperature":29.569653248795415,"humidity":78.6570864...
4   {"messageId":45,"deviceId":"Raspberry Pi Web Client","temperature":24.6194021861521,"humidity":64.2231368020...
5   {"messageId":46,"deviceId":"Raspberry Pi Web Client","temperature":25.138108867027242,"humidity":71.9368571...
6   {"messageId":47,"deviceId":"Raspberry Pi Web Client","temperature":29.672159609209903,"humidity":76.4116831...
7   {"messageId":48,"deviceId":"Raspberry Pi Web Client","temperature":25.248522099249495,"humidity":79.3085097...
8   {"messageId":49,"deviceId":"Raspberry Pi Web Client","temperature":25.055973001492944,"humidity":67.3227522...
9   {"messageId":50,"deviceId":"Raspberry Pi Web Client","temperature":31.256836008244356,"humidity":75.6093696...
10  {"messageId":51,"deviceId":"Raspberry Pi Web Client","temperature":31.045529017317861,"humidity":67.6985354...
11  {"messageId":52,"deviceId":"Raspberry Pi Web Client","temperature":26.419180362512087,"humidity":68.5382658...
12  {"messageId":53,"deviceId":"Raspberry Pi Web Client","temperature":31.923979997936684,"humidity":75.5815351...
13  {"messageId":54,"deviceId":"Raspberry Pi Web Client","temperature":26.9597001886166,"humidity":78.9160877077...
14  {"messageId":55,"deviceId":"Raspberry Pi Web Client","temperature":21.4435288086205,"humidity":69.2575675524...
15  {"messageId":56,"deviceId":"Raspberry Pi Web Client","temperature":27.0270166512053,"humidity":79.1931912246...
16  {"messageId":57,"deviceId":"Raspberry Pi Web Client","temperature":30.764071133556232,"humidity":70.0000165...
```

Json    Preview

17  {"messageId":58,"deviceId":"Raspberry Pi Web Client","temperature":31.887429099256433,"humidity":77.4390802120
18  {"messageId":59,"deviceId":"Raspberry Pi Web Client","temperature":23.733589243233098,"humidity":65.16477481⁹
19  {"messageId":60,"deviceId":"Raspberry Pi Web Client","temperature":24.527051891653819,"humidity":68.2522138268
20  {"messageId":61,"deviceId":"Raspberry Pi Web Client","temperature":30.5562086339128,"humidity":66.8708695521⁹8
21  {"messageId":62,"deviceId":"Raspberry Pi Web Client","temperature":24.663730123455945,"humidity":76.1673456420
22  {"messageId":63,"deviceId":"Raspberry Pi Web Client","temperature":28.623881163545793,"humidity":72.194303744⁹
23  {"messageId":64,"deviceId":"Raspberry Pi Web Client","temperature":21.507452893501736,"humidity":66.078687039⁷
24  {"messageId":65,"deviceId":"Raspberry Pi Web Client","temperature":23.954071724847008,"humidity":66.620100446⁷
25  {"messageId":68,"deviceId":"Raspberry Pi Web Client","temperature":29.825566797998267,"humidity":60.340044649⁷
26  {"messageId":69,"deviceId":"Raspberry Pi Web Client","temperature":30.412646579948508,"humidity":72.716832589⁸
27  {"messageId":70,"deviceId":"Raspberry Pi Web Client","temperature":27.758016975940595,"humidity":70.3170700991
28  {"messageId":71,"deviceId":"Raspberry Pi Web Client","temperature":27.40736970400112,"humidity":78.752502937⁰¹
29  {"messageId":72,"deviceId":"Raspberry Pi Web Client","temperature":30.756845959498492,"humidity":69.226669825⁴
30  {"messageId":73,"deviceId":"Raspberry Pi Web Client","temperature":21.891296329873242,"humidity":78.221896690⁹
31  {"messageId":74,"deviceId":"Raspberry Pi Web Client","temperature":23.841488923314554,"humidity":66.966216159³
22  {"messageId":76,"deviceId":"Raspberry Pi Web Client","temperature":31.779236107648849,"humidity":62.0891888180

Json  ⌄     ✎ Preview

31  {"messageId":74,"deviceId":"Raspberry Pi Web Client","temperature":23.841488923314554,"humidity":66.966216159³
32  {"messageId":76,"deviceId":"Raspberry Pi Web Client","temperature":31.779236107648849,"humidity":62.089188818⁰
33  {"messageId":77,"deviceId":"Raspberry Pi Web Client","temperature":27.825611760378621,"humidity":61.332545484⁹
34  {"messageId":78,"deviceId":"Raspberry Pi Web Client","temperature":23.595794033027467,"humidity":71.097904818⁹
35  {"messageId":79,"deviceId":"Raspberry Pi Web Client","temperature":22.686222167849614,"humidity":66.904505974³
36  {"messageId":80,"deviceId":"Raspberry Pi Web Client","temperature":23.563615038479927,"humidity":75.907096602⁷
37  {"messageId":81,"deviceId":"Raspberry Pi Web Client","temperature":27.915955882678542,"humidity":67.916067745⁷
38  {"messageId":82,"deviceId":"Raspberry Pi Web Client","temperature":30.727778753966149,"humidity":72.801938404⁰
39  {"messageId":83,"deviceId":"Raspberry Pi Web Client","temperature":29.565481945985383,"humidity":61.212127547³
40  {"messageId":84,"deviceId":"Raspberry Pi Web Client","temperature":24.996982815569975,"humidity":61.702192290⁰
41  {"messageId":85,"deviceId":"Raspberry Pi Web Client","temperature":30.125597644991956,"humidity":64.534437428³
42  {"messageId":86,"deviceId":"Raspberry Pi Web Client","temperature":28.196548645180645,"humidity":69.633966779⁰
43  {"messageId":87,"deviceId":"Raspberry Pi Web Client","temperature":25.22306879709144,"humidity":71.1153657415⁹
44  {"messageId":88,"deviceId":"Raspberry Pi Web Client","temperature":26.559600721636993,"humidity":71.602662458⁰
45  {"messageId":89,"deviceId":"Raspberry Pi Web Client","temperature":23.497132710919931,"humidity":61.041120087²

44  {"messageId":88,"deviceId":"Raspberry Pi Web Client","temperature":26.559600721636993,"humidity":71.602662458⁰
45  {"messageId":89,"deviceId":"Raspberry Pi Web Client","temperature":23.497132710919931,"humidity":61.041120087²
46  {"messageId":90,"deviceId":"Raspberry Pi Web Client","temperature":25.718827026151615,"humidity":67.112418543⁸
47  {"messageId":91,"deviceId":"Raspberry Pi Web Client","temperature":25.740181131691514,"humidity":61.466858942⁰
48  {"messageId":92,"deviceId":"Raspberry Pi Web Client","temperature":24.70892037869254,"humidity":77.0497111131⁴
49  {"messageId":93,"deviceId":"Raspberry Pi Web Client","temperature":26.518076218294148,"humidity":79.788552838⁰
50  {"messageId":94,"deviceId":"Raspberry Pi Web Client","temperature":30.031066108175537,"humidity":73.964708416⁰
51  {"messageId":95,"deviceId":"Raspberry Pi Web Client","temperature":24.775206107394613,"humidity":77.304864939⁸

# ASSIGNMENT 5

**TEAM MATES:**
Rahavi Selvarajan - 1007346445
rahavi.selvarajan@mail.utoronto.ca
Hruday Vishal Kanna Anand - 1006874517
vishal.kanna@mail.utoronto.ca

## PART - B

1. **Explain what problem you are going to solve using this dataset. Provide a brief overview of your problem statement.**

   The dataset chosen is **Gait Classification Dataset**. The dataset was created by calculating the walking parameters of 16 different volunteers aged between 20 and 34 years old. A total of 321 attributes are present and out of which only 26 attributes have proper labels.

   **Problem Statement:** One of the gait attributes is gait variability. Gait variability is the fluctuation of gait measures between steps. Higher variability indicates the risk of falling or frailty or developing a neuro-degenerative disease and is more common in aged people. Lower variability indicates that a person is physically fit. Our target here is gait variability. We are trying to train a model to predict whether a person is prone to the risk of having a neuro degenerative disease or is physically fit with the help of the other gait attributes.

   **Cleaning:** The dataset has a large number of attributes with no proper labels and these attributes are removed.

   **Preprocessing:** We have to map the variability feature to 1's and 0's as it will be our target feature for our models.

   **Modelling:** For the modelling, we chose two algorithms - Logistic Regression and Naive Bayes.

2. **Explain your dataset. Explore your dataset and provide at least 5 meaningful charts/graphs with explanation.**

   1. **Cycle Time Vs Speed**

   Gait cycle time is the time interval between two successive occurrences of one of the repetitive events of the locomotion. Gait Speed is the time one takes to walk a specified distance. From the graph below, we could see that the cycle time remains almost equal for all the volunteers irrespective of the time taken by them to cover a specified distance.

## 2. Speed Vs Cadence

Gait Cadence is the number of steps or cycles completed by a person in a specified period of time i.e., steps or cycles per minute. From the graph, it is evident that as the speed increases the number of steps taken to complete a specified distance decreases. The graph shows a linearly decreasing behaviour.



## 3. Step Angle Grouped by variability

The step angle captures the angular change (ie, rotation) in the sagittal plane of lower limbs during walking. From the graph below, it is evident that people without risks of neurodegenerative diseases have larger step angles (indicated by red) and those with the risks show smaller step angles (indicated by green).

### 4. Speed Vs Loading - Grouped by Variability

Gait loading is the phase where the body absorbs the impact of the foot and the trunk of the body moves forward to align with the line of gravity. The graph shows that the loading is almost zero in people with the risks of developing the neurodegenerative diseases.

**5. Thrust Vs Loading - Grouped by Variability**

Like the previous graph, irrespective of the thrust, loading remains zero.



**3. Do data cleaning/pre-processing as required and explain what you have done for your dataset and why?**

**Loading dataset into notebook and exploring**

**Removing unwanted columns-** while we uploaded the file as dataset to azure ml service we removed most of the unlabeled columns but some still persisted. These columns were removed with the code below.

```
1  df.drop(["Column172", "Column323"], axis = 1, inplace = True)
2  print (df)
```
[16]  ✓ <1 sec

```
    instances  Speed  Variability  Symmetry  Heel Press Time  Cycle Time \
0           0   1.32         4.15      4.00            1.054       1.054
1           0   1.29         0.00      0.90            1.119       1.137
2           0   1.25         5.06     -3.80            1.109       1.109
3           1   1.21         0.00     -6.30            1.185       1.162
4           1   1.20         4.43     -7.90            1.188       1.172
5           1   1.20         0.00     -8.40            1.235       1.213
6           2   1.37         0.00      1.00            1.239       1.215
7           2   1.46         0.00     -4.50            1.160       0.168
8           2   1.36         0.00      2.80            1.185       1.197
9           3   1.29         5.83     -4.80            1.167       1.170
10          3   1.30         7.84     -5.80            1.161       1.158
11          3   1.58         0.00     -7.40            1.114       1.093
12          4   1.45         0.00      0.90            1.029       1.025
13          4   1.47         0.00      3.50            1.039       1.064
14          4   1.49         0.00      1.80            1.050       1.033
15          5   1.21         6.09      1.30            1.140       1.134
16          5   1.15         0.00     -0.70            1.264       1.211
```

**Updating target column(Variability)-** to 1 and 0 so that we can proceed to make a valid classification model. All the classification models the target feature needs to be distinct classes for the model to function properly.

```
1  df['Variability'] = (df['Variability'] > 0).astype(int)
2  df
```
[17]  ✓ <1 sec

| | instances | Speed | Variability | Symmetry | Heel Press Time | Cycle Time | Cadence | Posture | Oscillation | Loading | ... | Peak Angle Speed, | Max |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 1.32 | 1 | 4.00 | 1.054 | 1.054 | 1.050 | 1.060 | 0.043 | 0.044 | ... | 1.280 | 0.99 |
| 1 | 0 | 1.29 | 0 | 0.90 | 1.119 | 1.137 | 1.060 | 1.065 | 0.406 | 0.567 | ... | 5.660 | 0.98 |
| 2 | 0 | 1.25 | 1 | -3.80 | 1.109 | 1.109 | 1.105 | 1.115 | 0.048 | 0.056 | ... | 1.265 | 1.00 |
| 3 | 1 | 1.21 | 0 | -6.30 | 1.185 | 1.162 | 1.155 | 1.160 | 0.215 | 0.103 | ... | 1.640 | 1.02 |
| 4 | 1 | 1.20 | 1 | -7.90 | 1.188 | 1.172 | 1.175 | 1.177 | 0.153 | 0.052 | ... | 1.290 | 1.01 |
| 5 | 1 | 1.20 | 0 | -8.40 | 1.235 | 1.213 | 1.165 | 1.165 | 0.567 | 0.406 | ... | 4.345 | 1.01 |
| 6 | 2 | 1.37 | 0 | 1.00 | 1.239 | 1.215 | 1.140 | 1.145 | 0.648 | 0.460 | ... | 4.655 | 1.01 |
| 7 | 2 | 1.46 | 0 | -4.50 | 1.160 | 0.168 | 1.100 | 1.100 | 0.348 | 0.346 | ... | 3.395 | 0.99 |

**Replicating rows-** to create more data for automated ml, as the process requires a minimum of 50 rows to run.



**Saving cleaned and processed file-**



4. **Implement 2 machine learning models, explain which algorithms you have selected and why. Compare them and show success metrics (Accuracy/RMSE/Confusion Matrix) as per your problem. Explain results.**

**Model 1:**
**Algorithm:** Logistic Regression
**Regularization Parameter:** 2.0
**Reason for choosing Logistic Regression:** Logistic Regression is the most common algorithm used for binary classification involving numerical values and is known to prevent overfitting.

**Model 2:**
**Algorithm:** Naive Bayes
**Reason for choosing Naive Bayes:** It is easy to implement and performs well in test dataset.



From the above two models, Logistic Regression performs better compared to the Naive Bayes. The Accuracy, F1 score and AUC scores of the first model are better compared to the model 2 i.e, Naive Bayes. F1 score shows the accuracy of the test set and logistic regression shows better performance on the test set.

**5. Use Automated ML for your data set. Explain best model results.**

**Automated ml run initiated-**



**Automated run completed**

**Models-**



We can see that 42 models have been made and trained and almost all of them have an accuracy of 1.

**Best model-**
The best model with Light GBM with max abs scaler. Light GBM is a gradient boosting framework that uses tree based learning algorithms. Max abs scaler is used to preprocess our features dividing our features by the maximum value to scale down the features to a common scale.

**Best model summary**

Algorithm name
MaxAbsScaler, LightGBM

Hyperparameters
View hyperparameters

Accuracy
1.00000    View all other metrics

Sampling
100.00 %

Registered models
No registration yet

Deploy status
No deployment yet

**Hyperparameters of the best model-**

Hyperparameters                                    ×

Data transformation:
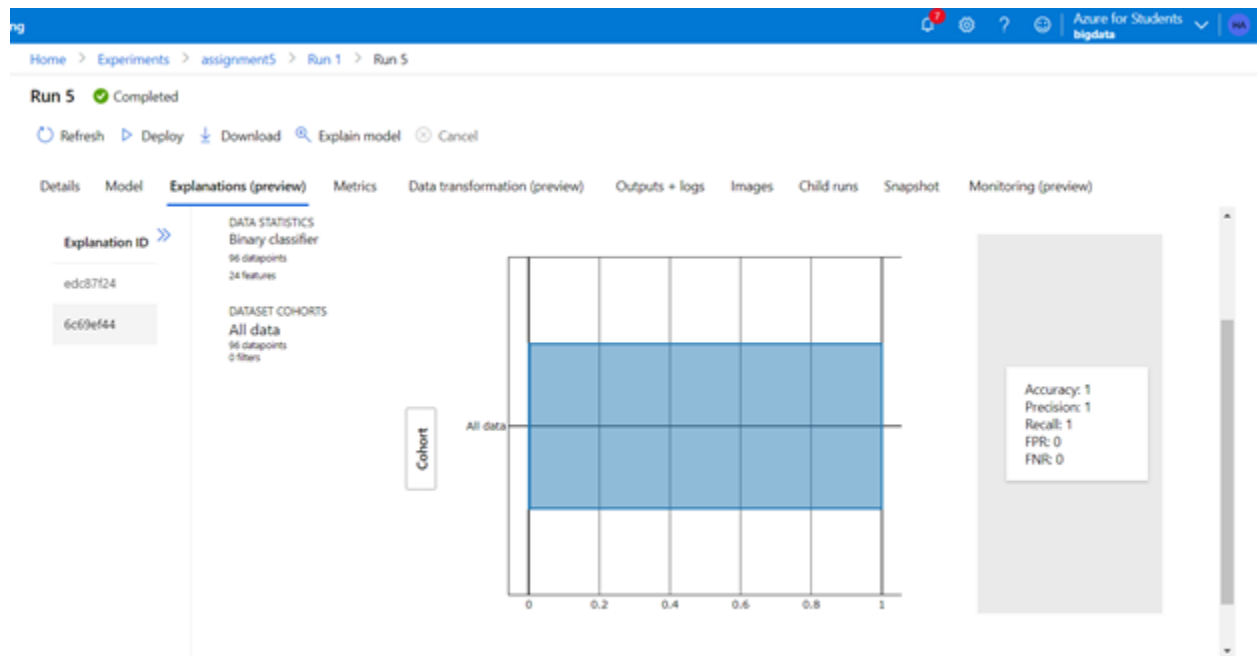
```
1  {
2      "spec_class": "preproc",
3      "class_name": "MaxAbsScaler",
4      "module": "sklearn.preprocessing",
5      "param_args": [],
6      "param_kwargs": {},
7      "prepared_kwargs": {}
8  }
```

Training algorithm:

```
1  {
2      "spec_class": "sklearn",
3      "class_name": "LightGBMClassifier",
4      "module": "automl.client.core.common.model_wrappers",
5      "param_args": [],
6      "param_kwargs": {
7          "min_data_in_leaf": 20
8      },
9      "prepared_kwargs": {}
10 }
```

**Model explanation-**
We can see the model has an accuracy, precision and recall 1.



From this graph we can see our model relies largely on the peak angle speed feature to come up with an accurate output.