# Predict Patient diabetes using Microsoft Azure ML Work Bench

**Prerequisites**

- ➢ Valid Azure Subscription
- ➢ [Azure Machine Learning Workbench](#)

**Steps Involved**

- ➢ Download and Understanding DataSet
- ➢ Build Model

**Objective**

The objective of the dataset is to diagnostically predict whether a patient has diabetes, based on certain diagnostic measurements included in the dataset. Several constraints were placed on the selection of these instances from a larger database. All patients here are females at least 21 years old of Pima Indian heritage.

**The Tool**

Azure Machine Learning services (preview) aka Azure ML workbench are an integrated, end-to-end data science and advanced analytics solution. It helps **professional data scientists** prepare data, develop experiments, and deploy models at cloud scale.

**Download and Understanding Dataset**

The datasets consist of several medical predictor variables and one target variable, Outcome. Predictor variables includes the number of pregnancies the patient has had, their BMI, insulin level, age, and so on
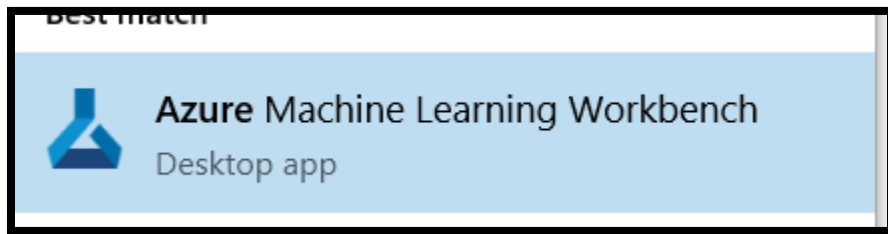
// Please download dataset from https://aka.ms/aidevdaysds2

| Column Name | Description | Data Type |
|---|---|---|
| Pregnancies | Number of times pregnant | Numeric |
| Glucose | Plasma glucose concentration a 2 hour in an oral glucose tolerance test | Numeric |
| BloodPressure | Diastolic blood pressure (mm Hg) | Numeric |
| SkinThickness | Triceps skin fold thickness (mm) | Numeric |
| Insulin | 2-Hour serum insulin (mu U/ml) | Numeric |
| BMI | Body mass index (weight in kg/(height in m)^2) | Numeric |
| DiabetesPedigreeFunction | Diabetes pedigree function | Numeric |
| Age | Age (years) | Numeric |
| Outcome | Class variable (0 or 1) | Numeric |

***Note:-*** Datasource Reference: - https://www.kaggle.com/uciml/pima-indians-diabetes-database/data

**Evaluating Data Fields**

➢ Open "**Azure Machine Learning Workbench**"



➢ Click **+** and **New Project** to add new project



➢ Provide **Project Name**, **Project directory**, **Project Description**, **Selected Workspace** and select **Blank Project**. Click **Create**

➢ Once created it will look like below



➢ Click **Data** (next icon after home) and Click **+** and **Add Data Source**

➢ Select **Text Files** and **Next**

## Add Data Source

1. **Data Store**

2. **File Selection**

3. File Details

4. Data Types

5. Sampling

6. Path Column

### Where does the data come fro

Specify the data store where the data comes from.

Text Files (*.csv, *.json, *.txt, ...)

Parquet

Previous     **Next**

➢ Select the file which downloaded in earlier step, Click **Next**

**Browse to find the file**
Browse to the path of the file you would like to use.

1. Data Store
2. **File Selection**
3. File Details
4. Data Types
5. Sampling
6. Path Column

**Path**

| Local | ▼ | C:\sudhirawdata\SudhirData\EVENTS\AI Dev Days\India |

Previous    Next    Finish

➢ On step number **3,** leave it default and click **Next**

- ➢ Check the datatype on step number **4**. All should be numeric except first column (**Path**).
  Click **Next**

> Click **Next** on step number **5**

➢ On step number **6** (**Path Column**), Select "**Do Not Include Path Column**" and click **Finish**

Path column handling
You can choose to include a column containing source file paths.

Include File Paths in Data?

Do Not Include Path Column

Previous    Finish

1. Data Store
2. File Selection
3. File Details
4. Data Types
5. Sampling
6. Path Column

➤ Once setup, screen will look like

➢ Click on **Matrix.** Data Profiling will help you identify problems in your data such as invalid dates, missing values etc.

| Column | Profile | Max value | Min value | Count | Quantile at 50% | Median value | Kurtos |
|---|---|---|---|---|---|---|---|
| PregnantTimes | | 17 | 0 | 768 | 3 | 3 | 0.159 19777 |
| PlasmaGlucose | | 199 | 0 | 768 | 117 | 117 | 0.640 79820 |
| DbloodPressure | | 122 | 0 | 768 | 72 | 72 | 5.180 56560 |

*__Note__: - Since data is clean we are not going to do any transformation

➢ Click **Data**

⊞ Data    ⅆ Prepare

Column metrics for current sample      Choose M

| viation | Variance | Quantile at 25% | Is num |
|---|---|---|---|
| 0627 | 11.354 056320 | 1 | |

➢ Right click on **Data Source** and click **Generate Data Access Code File**



➢ After generating data access code file, the screen will look like



**Building Machine Learning Model using** [scikitlearn](scikitlearn)

➢ Copy below code and paste it in editor

```
# import references
import pickle
import pandas as pd
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
import numpy as np

# Use the Azure Machine Learning data source package
from azureml.dataprep import datasource
```

➢ The screen will look like below after adding above code

```
1  # import references
2  import pickle
3  import pandas as pd
4  from sklearn.metrics import accuracy_score
5  from sklearn.model_selection import train_test_split
6  from sklearn.tree import DecisionTreeClassifier
7  import numpy as np
8
9  # Use the Azure Machine Learning data source package
10 from azureml.dataprep import datasource
11
12 # Use the Azure Machine Learning data collector to log various metrics
13 from azureml.logging import get_azureml_logger
14 logger = get_azureml_logger()
15
16 # This call will load the referenced data source and return a DataFrame.
17 # If run in a PySpark environment, this call returns a
18 # Spark DataFrame. If not, it will return a Pandas DataFrame.
19 df = datasource.load_datasource('IndianFDiabitiesDataSet.dsource')
20
21 # Remove this line and add code that uses the DataFrame
22 df.head(10)
23
24
25
```

➤ Comment code **# df.head(10)** and copy below code and paste it to the editor. In this piece of code we are defining columns which algorithm used to define patterns

```
# Capture features from the given dataset
features=['PlasmaGlucose', 'bmi', 'age', 'tricepsskinfold', 'PregnantTimes',
'diabetespedigreefunction',
    'DbloodPressure', '2hourseruminsulin']
X = df[features].copy()
print("Features")
print(X)
```

➤ After changes

```
 9  # Use the Azure Machine Learning data source package
10  from azureml.dataprep import datasource
11
12  # Use the Azure Machine Learning data collector to log various metrics
13  from azureml.logging import get_azureml_logger
14  logger = get_azureml_logger()
15
16  # This call will load the referenced data source and return a DataFrame.
17  # If run in a PySpark environment, this call returns a
18  # Spark DataFrame. If not, it will return a Pandas DataFrame.
19  df = datasource.load_datasource('IndianFDiabitiesDataSet.dsource')
20
21  # Remove this line and add code that uses the DataFrame
22  # df.head(10)
23
24  # Capture features from the given dataset
25  features=['PlasmaGlucose', 'bmi', 'age', 'tricepsskinfold', 'PregnantTimes', 'diabetespedigreefunction',
26          'DbloodPressure', '2hourseruminsulin']
27  X = df[features].copy()
28  print("Features")
29  print(X)
```

➢ Copy below code and paste it in editor. In this piece of code, we are defining the target column

   #Target column
   label_features=['isdiabetic']

   Y = df[label_features].copy()

➢ Below is the screenshot after pasting the code

```
24  # Capture features from the given dataset
25  features=['PlasmaGlucose', 'bmi', 'age', 'tricepsskinfold', 'PregnantTimes', 'diabetespedigreefunction',
26          'DbloodPressure', '2hourseruminsulin']
27  X = df[features].copy()
28  print("Features")
29  print(X)
30
31  #Target column
32  label_features=['isdiabetic']
33
34  Y = df[label_features].copy()
35
36
```

➢ Copy below code and paste it in editor. In this example we are splitting the dataset to avoid any overfitting and underfitting scenario. Next we are defining the algorithm (DecisionTreeClassifier) to solve this problem and used fit method to build a decision tree classifier from the training set.

# Split Data into Training and Test
# test_size =.30 defines that training data is 70% and test is 30%

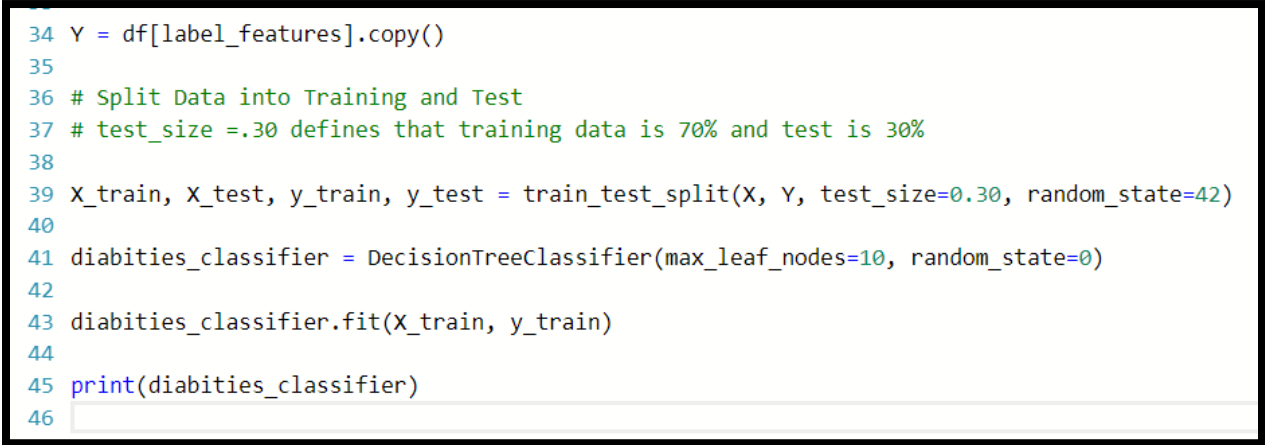X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.30, random_state=42)

diabities_classifier = DecisionTreeClassifier(max_leaf_nodes=10, random_state=0)

diabities_classifier.fit(X_train, y_train)

print(diabities_classifier)

➢ Below is the screenshot after implementing above code

```
34  Y = df[label_features].copy()
35
36  # Split Data into Training and Test
37  # test_size =.30 defines that training data is 70% and test is 30%
38
39  X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.30, random_state=42)
40
41  diabities_classifier = DecisionTreeClassifier(max_leaf_nodes=10, random_state=0)
42
43  diabities_classifier.fit(X_train, y_train)
44
45  print(diabities_classifier)
46
```

➢ Copy below code. Once Model is trained. Next step is to test it against the test dataset

# evaluate the test set

predictions = diabities_classifier.predict(X_test)

print("Prediction")

print(predictions)

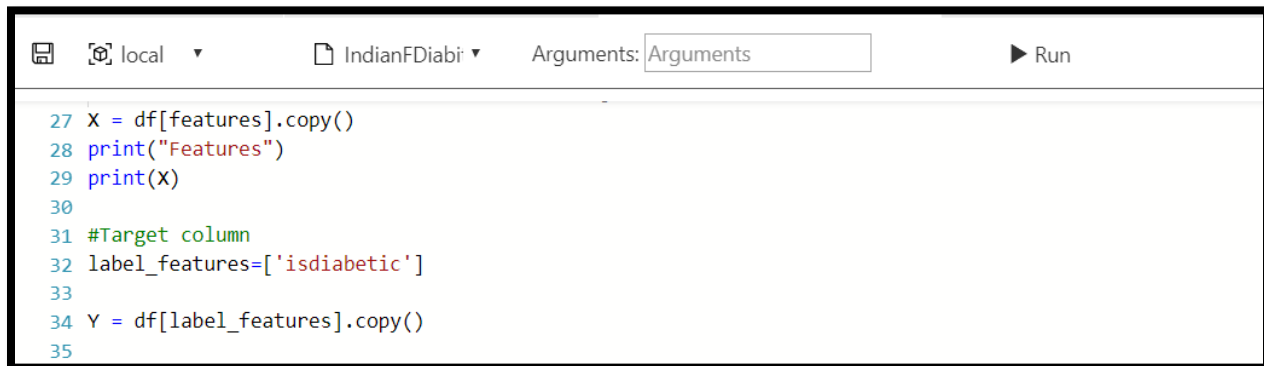accuracy_scored = accuracy_score(y_true = y_test, y_pred = predictions)

print ("Accuracy is {}".format(accuracy_scored))

➢ Below is the screenshot after implementing above code

```
43  diabities_classifier.fit(X_train, y_train)
44
45  print(diabities_classifier)
46
47  # evaluate the test set
48  predictions = diabities_classifier.predict(X_test)
49  print("Prediction")
50  print(predictions)
51  accuracy_scored = accuracy_score(y_true = y_test, y_pred = predictions)
52  print ("Accuracy is {}".format(accuracy_scored))
53
```

➢ Let's save and run the code to test the accuracy of the model. Click "b" to save the file and click **Run**



```
27  X = df[features].copy()
28  print("Features")
29  print(X)
30
31  #Target column
32  label_features=['isdiabetic']
33
34  Y = df[label_features].copy()
35
```

➢ On the right side of the screen it will show the progress of the job. Wait till job gets completed

➤ Click on the completed Job



➤ Scroll down and click on **driver_log**

➢ It will show the accuracy around 76. You can try to add more random data and tweaking algorithm parameter to test for better accuracy. For now we are good with 76.
➢ Click on python file to build model file.



➢ Copy below code and paste it. Below code will serialize the model by inserting pickle file in output folder.

```
# serialize the model on disk in the special 'outputs' folder
print ("Export the model to model.pkl")
f = open('./outputs/model.pkl', 'wb')
pickle.dump(diabities_classifier, f)
f.close()
```

➢ Below is the screenshot after implementing above code.

```
50 print(predictions)
51 accuracy_scored = accuracy_score(y_true = y_test, y_pred = predictions)
52 print ("Accuracy is {}".format(accuracy_scored))
53
54 # serialize the model on disk in the special 'outputs' folder
55 print ("Export the model to model.pkl")
56 f = open('./outputs/model.pkl', 'wb')
57 pickle.dump(diabities_classifier, f)
58 f.close()
59
```

➤ Copy below code and past it in editor. In below code, we are loading pickle file in pass sample data to predict the outcome

# load the model back from the 'outputs' folder into memory

print("Import the model from model.pkl")

f2 = open('./outputs/model.pkl', 'rb')

diabities_classifier2 = pickle.load(f2)


# predict on a new sample

# ['PlasmaGlucose', 'bmi', 'age', 'tricepsskinfold', 'PregnantTimes', 'diabetespedigreefunction',

#      'DbloodPressure', '2hourseruminsulin']

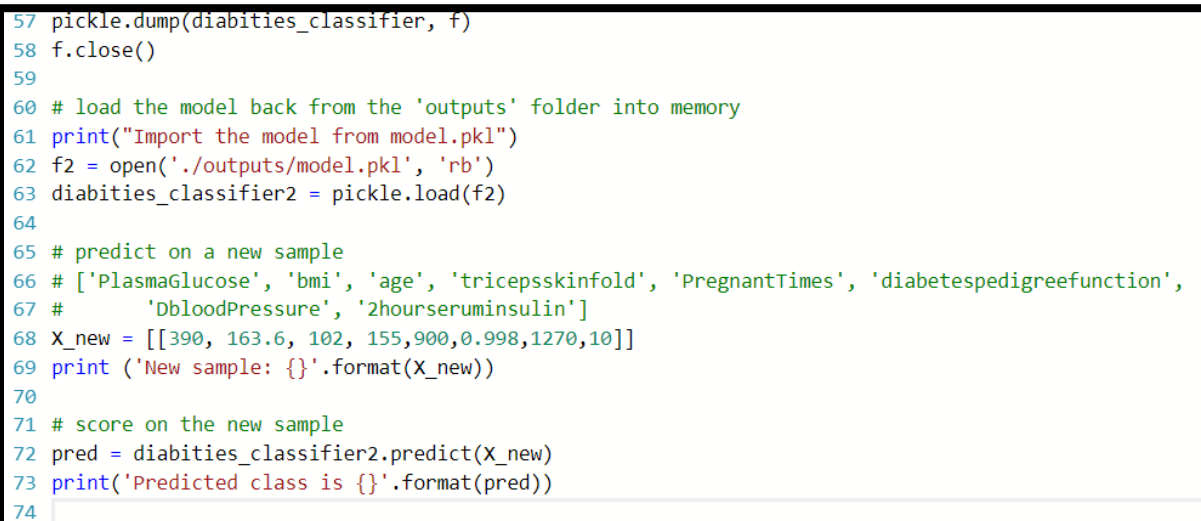X_new = [[390, 163.6, 102, 155,900,0.998,1270,10]]

print ('New sample: {}'.format(X_new))

# score on the new sample
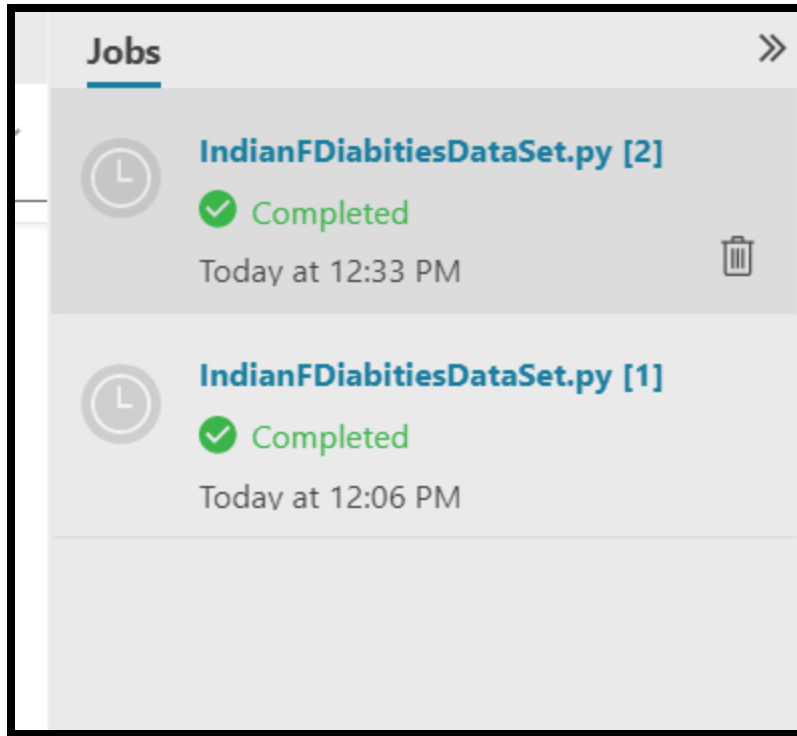
pred = diabities_classifier2.predict(X_new)

print('Predicted class is {}'.format(pred))

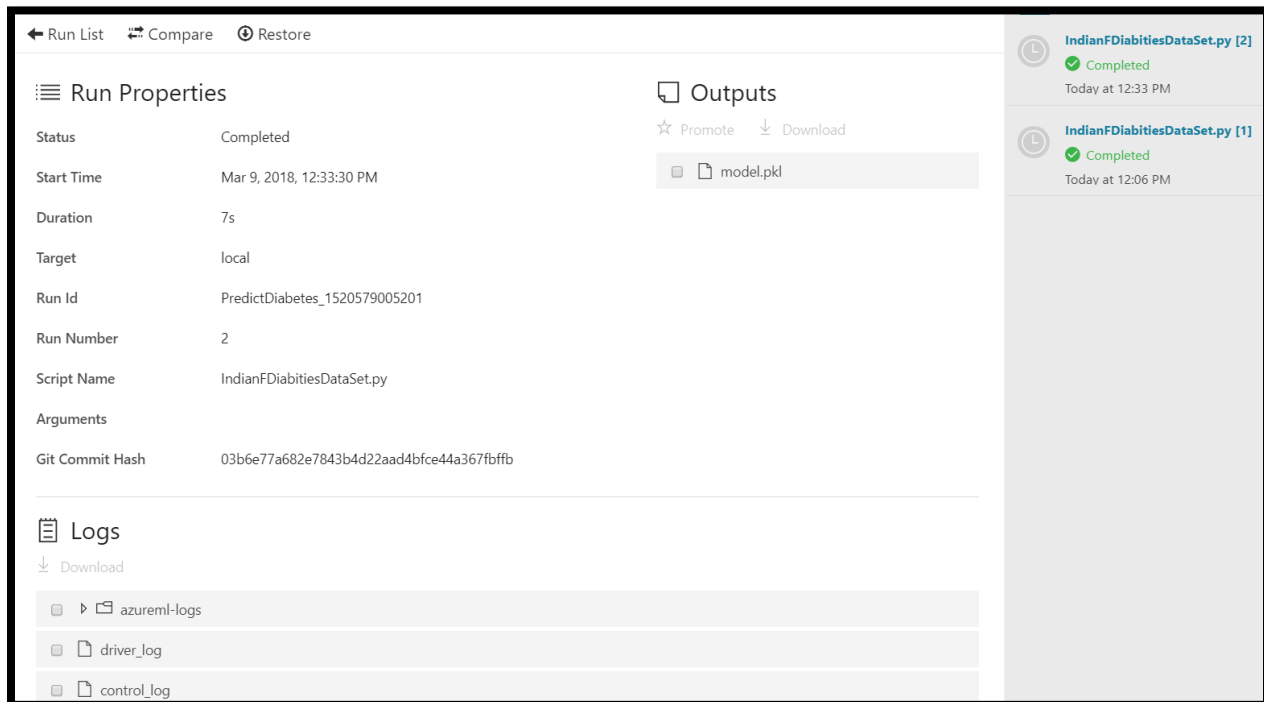➤ Below is the screenshot after pasting above code

```
57 pickle.dump(diabities_classifier, f)
58 f.close()
59
60 # load the model back from the 'outputs' folder into memory
61 print("Import the model from model.pkl")
62 f2 = open('./outputs/model.pkl', 'rb')
63 diabities_classifier2 = pickle.load(f2)
64
65 # predict on a new sample
66 # ['PlasmaGlucose', 'bmi', 'age', 'tricepsskinfold', 'PregnantTimes', 'diabetespedigreefunction',
67 #       'DbloodPressure', '2hourseruminsulin']
68 X_new = [[390, 163.6, 102, 155,900,0.998,1270,10]]
69 print ('New sample: {}'.format(X_new))
70
71 # score on the new sample
72 pred = diabities_classifier2.predict(X_new)
73 print('Predicted class is {}'.format(pred))
74
```
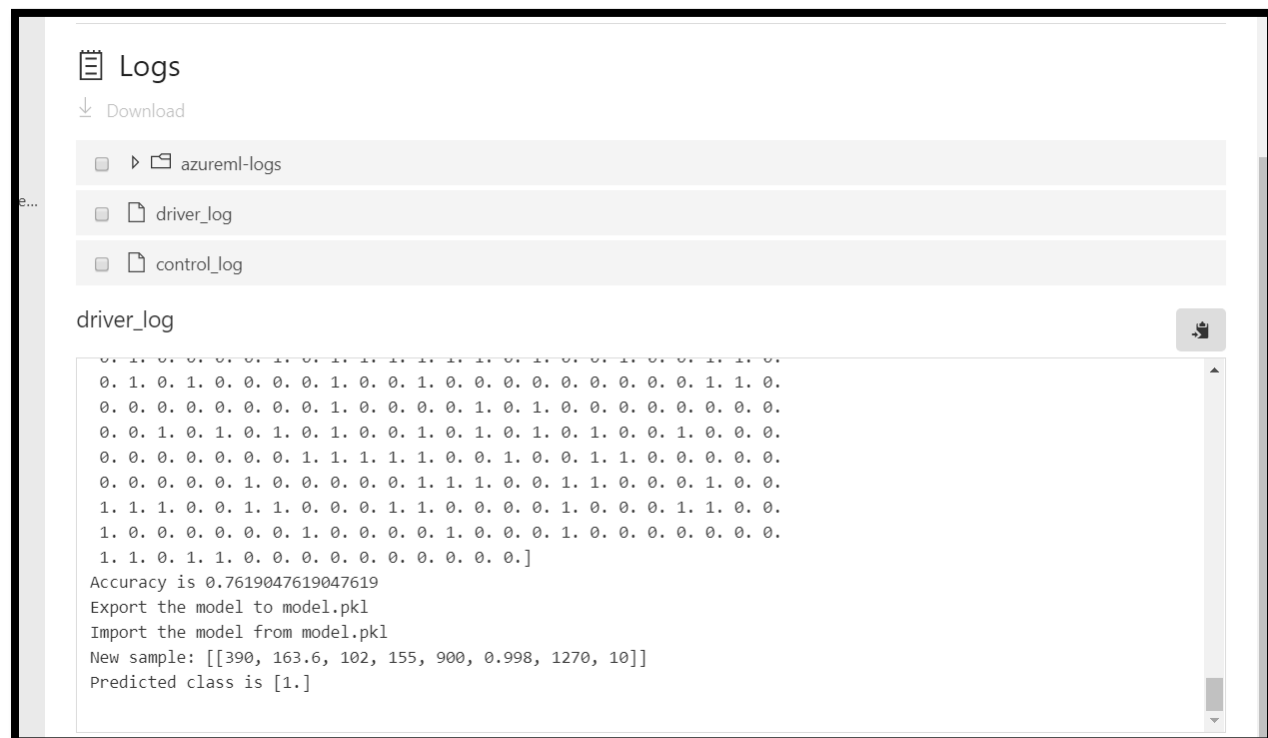
➤ Click **Run**

➢ Click on latest job execution

➢ Click on **driver_log.** In below screenshot you will notice exporting and importing of model file and test it against new sample data



➢ Next Step is to operationalize it which you already did in previous HOL. If time permits you can try it again here. You can also look out advance data preparation to understand various data wrangling techniques.

**Summary**

Well done! In this hands-on lab we solved a problem to predict whether a person going to be a diabetic or not. This will help to take proactive measure and cure patient. You can deploy this model and make it available as a web service so that any application from anywhere can access it. We also learn the capabilities of Azure Machine Learning and how to work with it. You can download data from UCI and Kaggle to solve other problems. Let's move onto the next step.