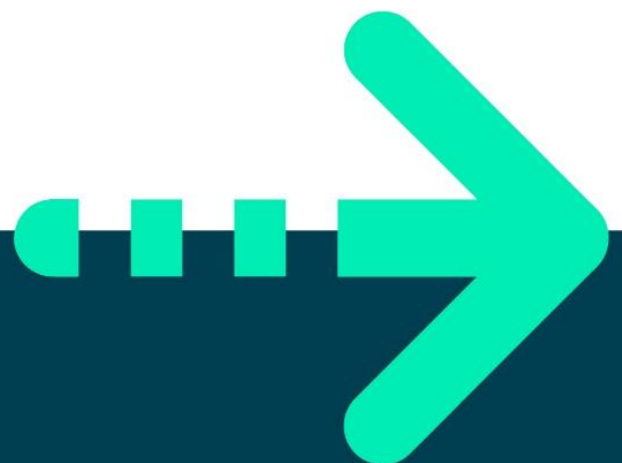# Implementing a Data Science Solution with Azure

**By Sadique Mohammed – DFEDATA4**

# Part 1: Problem Domain Understanding

**Question 1**

Given the nature of the business you have been asked to describe the benefits of implementing the data science solution in the cloud compared to using on premise/hybrid data solutions.

The evaluation criteria should consider the following:

- security
- compliance
- scalability
- efficiency
- reliability
- fidelity
- flexibility
- portability

Please see table below for benefits of using a cloud data solution compared to on premise/hybrid data solution.

| Cloud | On Premise/Hybrid |
|---|---|
| <ul><li>On-demand scalability</li><li>Cost efficiency</li><li>Bundled capabilities such as IAM and analytics</li><li>Security</li><li>System uptime and availability</li></ul> | <ul><li>Complete control over the tech stack</li><li>Local speed and performance</li><li>Governance and regulatory compliance</li></ul> |

# Part 2: Data Understanding

**Question 1**

Complete all the tasks below in Microsoft Azure Machine Learning Studio using the dataset identified for your project and discuss what the implications of the results for at least two of the tasks on the business questions that are candidates for your data science solution.

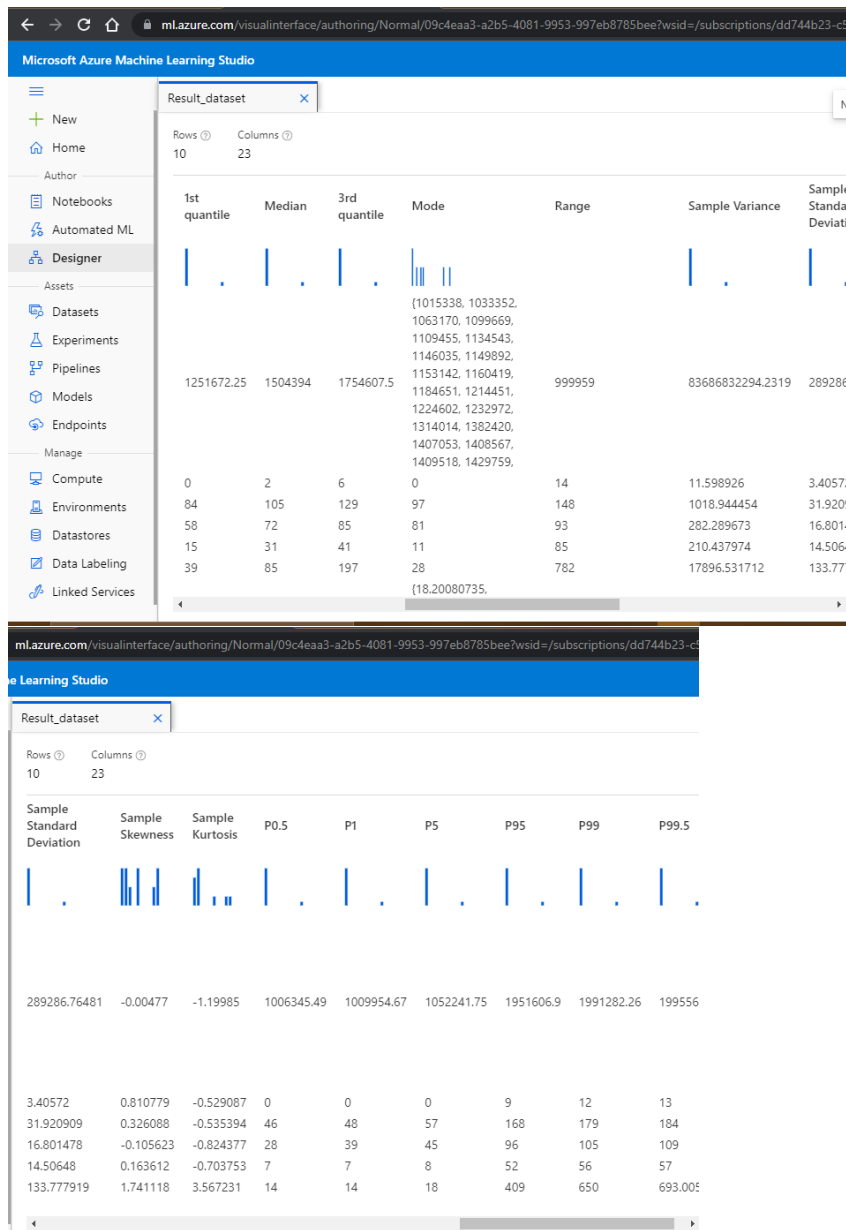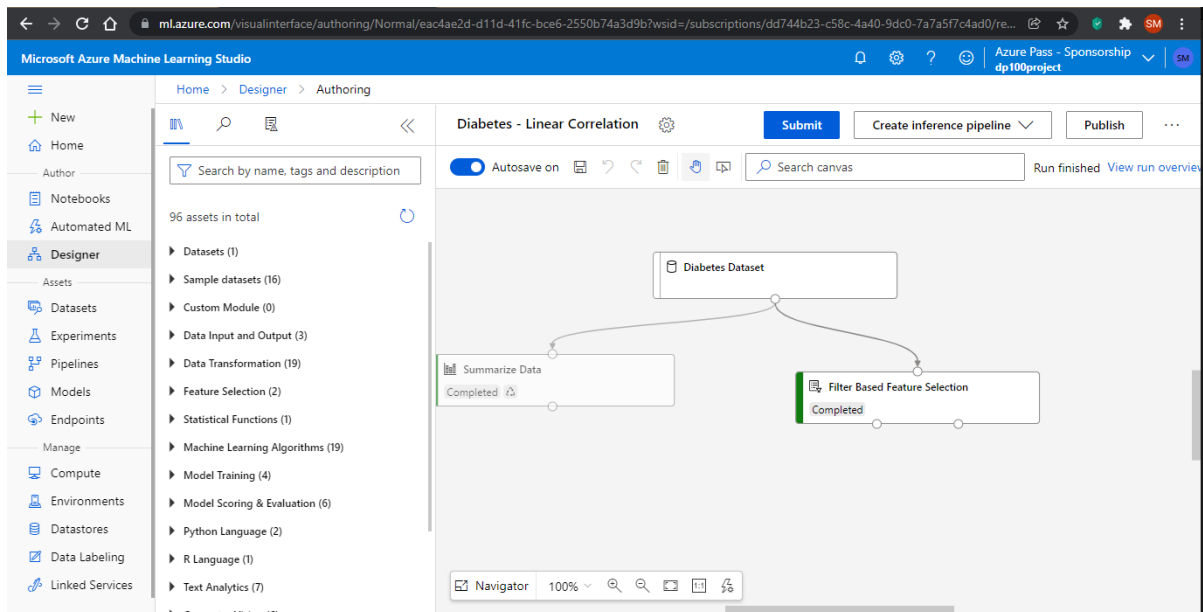**Task 1:** Generate a descriptive statistics report for the columns in your dataset

*Above, you can see how I have obtained a statistical report for columns in the dataset in the Microsoft Azure Machine Learning Studio Designer.*

**Task 2:** Calculate a single statistical measure for each column which is useful for determining central tendency, dispersion, and shape of your dataset or whatever you deem necessary.
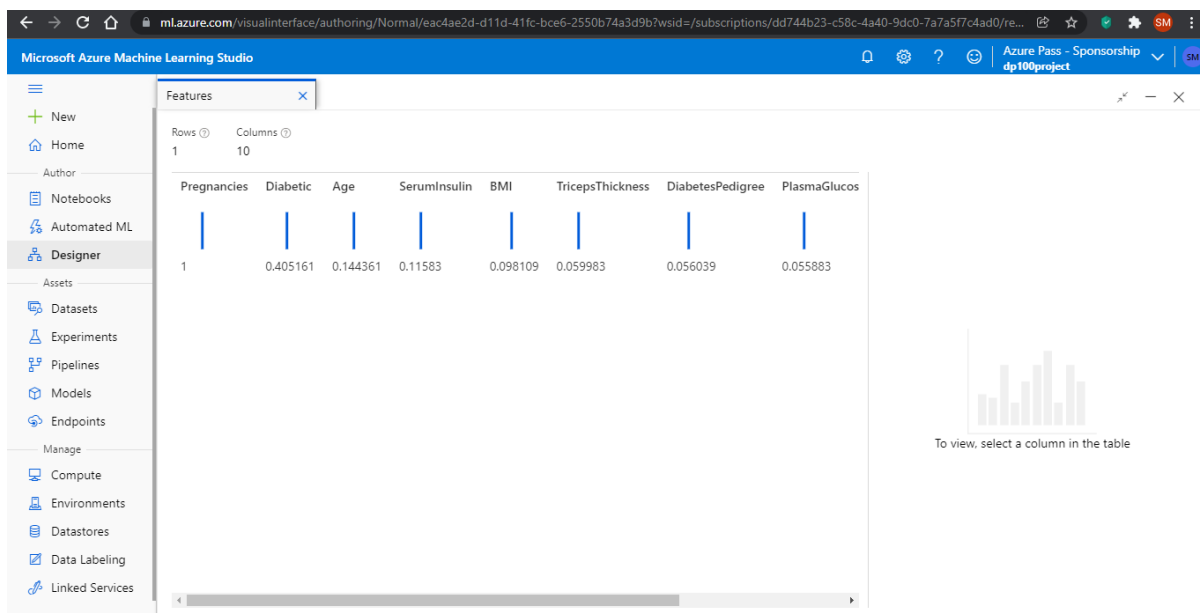




*Please also see task 1 for single statistical measures for each column of dataset.*

**Task 3:** Calculate the linear correlation between column values in your dataset

Above model is used to obtain the linear correlation between column values in dataset.



As you can see from above screenshot that Diabetic, Age, Seruminsulin and BMI are all positively correlated.

**Task 4:** Identify any data quality issues that may impact analysis

There are currently no data quality issues with this data sets, it is a pure dataset. There are no missing values or outliers to list. Also, there are no Nulls in this dataset.

**Task 5:** Conduct exploratory visual analysis of the dataset and comment on at least 2 of your key findings.

From the visual analysis above, all columns in dataset have a positive skewness. The maximum BMI value is 56.03 and minimum BMI is 18.2. Also, the mean age of the sample in this data is 30.13 yrs, which is 30 yrs.

**Note: Tasks may be completed by identifying and using the relevant modules in Designer.**

# Part 3: Data Preparation and Transformation

**Question 1:**

Discuss and apply the appropriate data processing techniques to address the data quality issues you identified in your exploratory analysis of the data.

No data quality issues were identified in the exploratory analysis of the data. All values in dataset are numerical and there is not any categorical values in dataset.

**Question 2:**

Conduct a comparative analysis of the mathematical techniques below for data normalization and apply the relevant one to your dataset.

- z-score



Above shows process of obtaining the Z Score.

- min-max

Above shows process of obtaining MinMax.

- log-normal



Above shows process of obtaining log-normal

- tanh

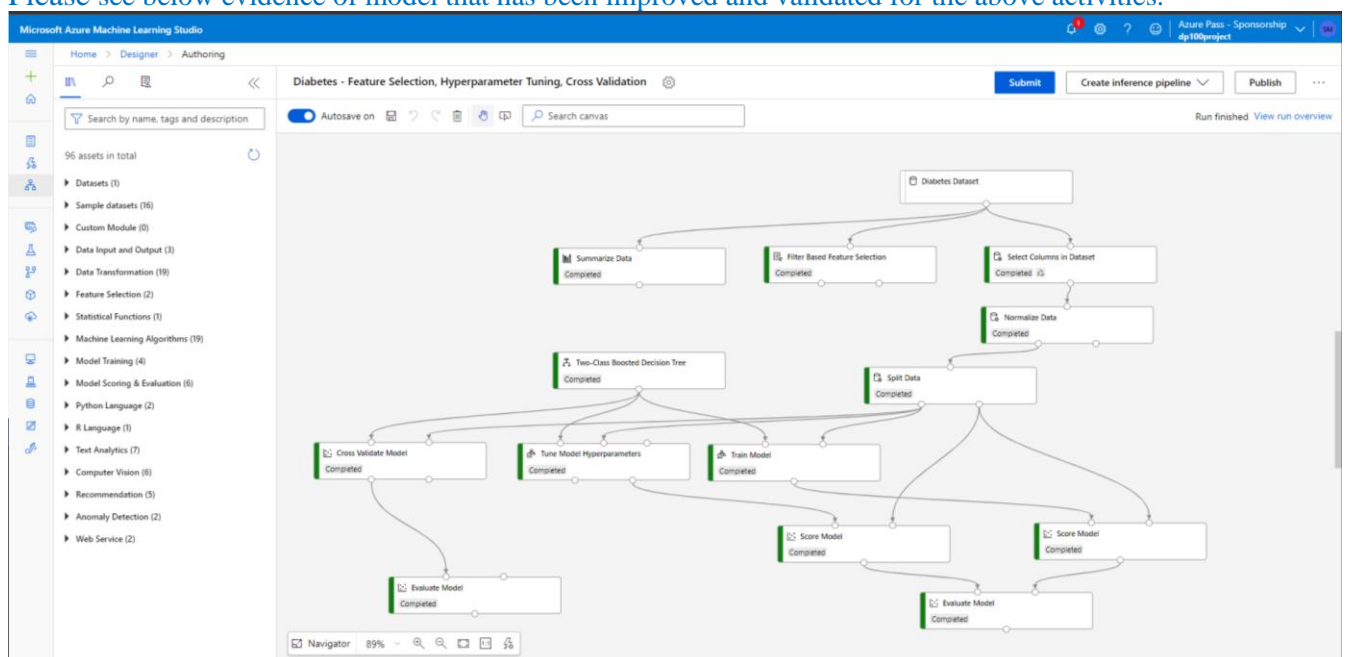Above shows process of obtaining tanh.

# Part 4: Model Validation and Improvement

**Question 1:**

Discuss and evidence how your model can be improved and validated through the following activities.

- Feature Selection
- Hyperparameter Tuning
- Cross Validation

Please see below evidence of model that has been improved and validated for the above activities.
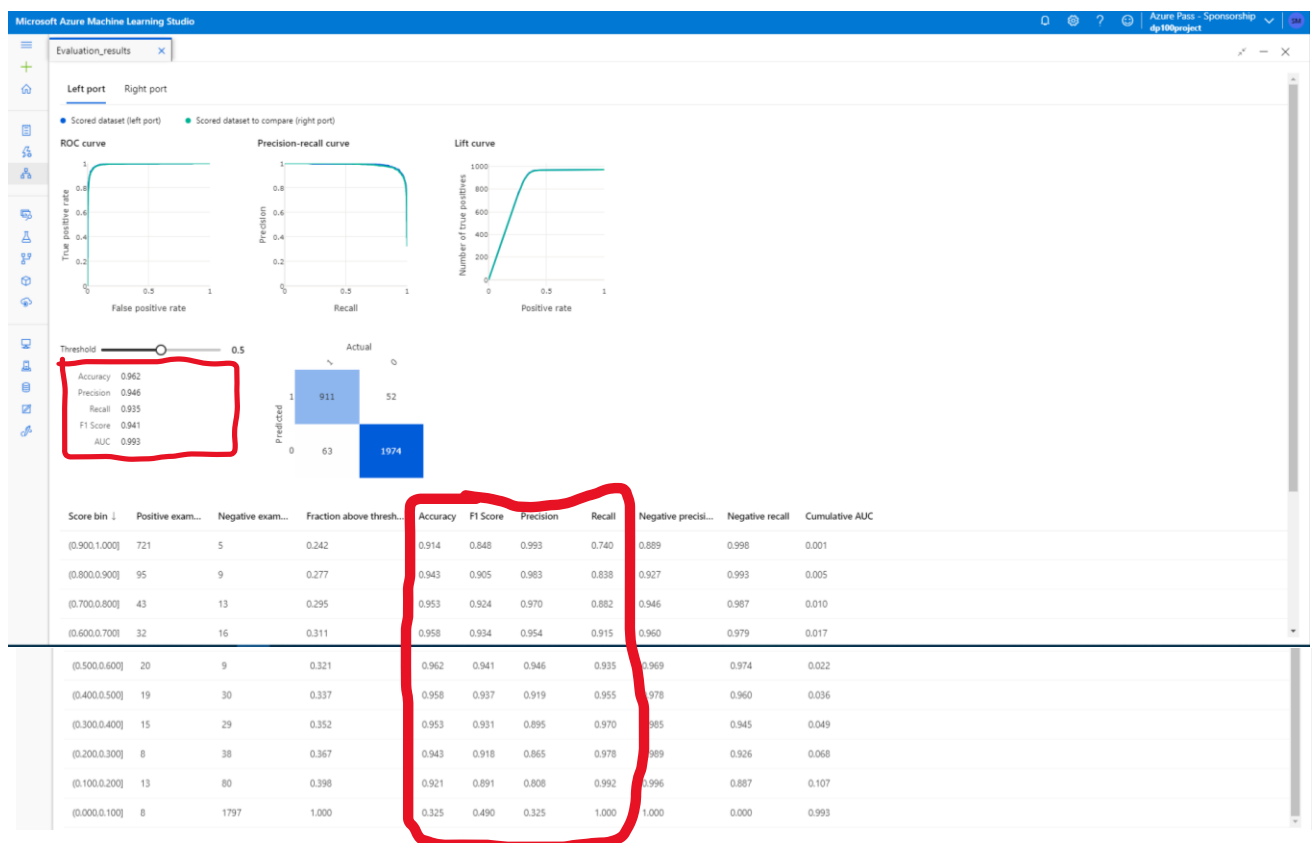
# Part 5: Evaluation and Metrics

## Question 1:

Critically evaluate the metrics for classification models below and identify which of the statistical measures are suitable for the evaluation of your data science solution. The selected metric should be used in the evaluation of your model.
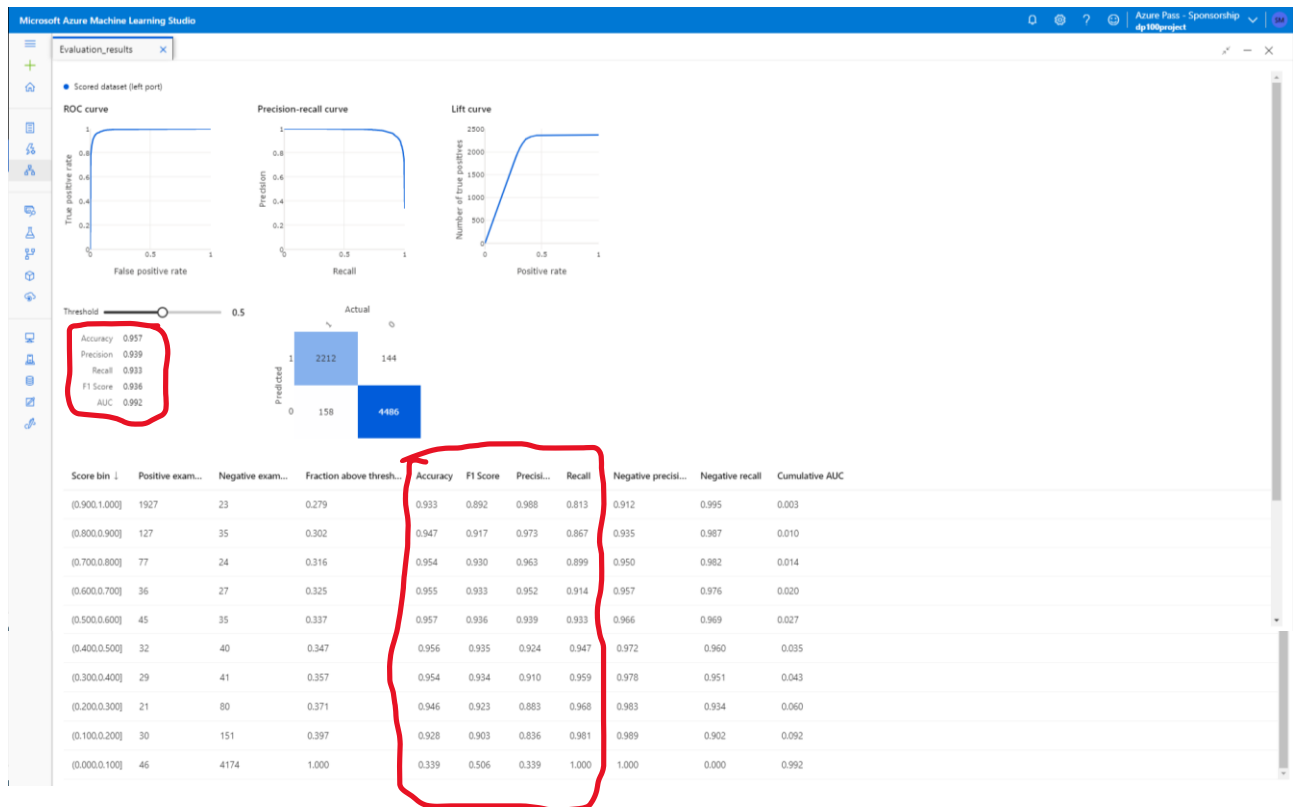
**Classification Metrics:**

- Accuracy
- Precision
- Recall
- F1 Score

Using the designer model from Part 4, I was able to obtain the above classification metrics via the evaluate model.

Below shows a screenshot of classification metrics via the hyperparameter tuning evaluate model.

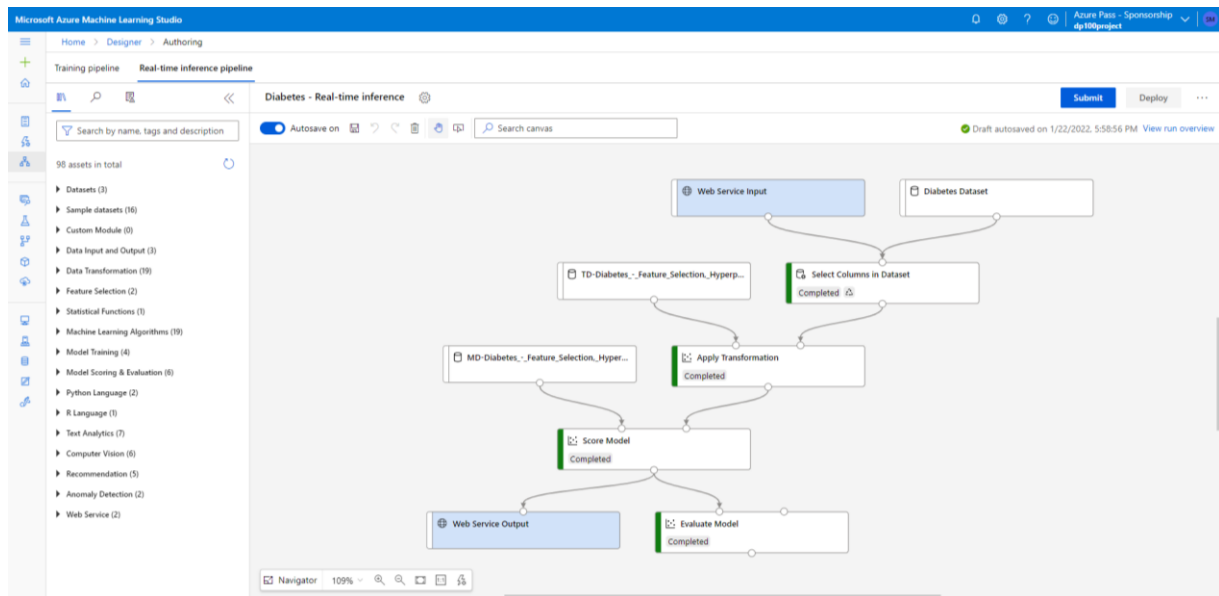Below shows a screenshot of classification metrics via the cross validation evaluate model.

# Part 6: Deployment

**Question 1:**

Deploy your model as a real-time inferencing service.

Below shows a model for a inference pipeline for real-time inferencing to be deployed.



Below shows deployment of Diabetes dataset as a real time inferencing service.

## Question 2:

Write a simple test script in Jupyter Notebook to consume the service.

Please see below simple test script in Jupyter Notebook to consume the service.



### Python Script in Jupyter Notebook

```python
import urllib.request
import json
import os
import ssl

def allowSelfSignedHttps(allowed):
    # bypass the server certificate verification on client side
    if allowed and not os.environ.get('PYTHONHTTPSVERIFY', '') and getattr(ssl, '_create_unverified_context', None):
        ssl._create_default_https_context = ssl._create_unverified_context
```

```python
allowSelfSignedHttps(True) # this line is needed if you use self-
signed certificate in your scoring service.

# Request data goes here
data = {
    "Inputs": {
        "WebServiceInput0":
        [
            {
                'PatientID': "1354778",
                'Pregnancies': "5",
                'PlasmaGlucose': "171",
                'DiastolicBloodPressure': "100",
                'TricepsThickness': "34",
                'SerumInsulin': "23",
                'BMI': "43.50972593",
                'DiabetesPedigree': "1.213191354",
                'Age': "60",
                'Diabetic': "0",
            },
        ],
    },
    "GlobalParameters": {
    }
}

body = str.encode(json.dumps(data))

url = 'http://0cceb6b4-5531-4c07-ba35-
7617a0c5f0cf.uksouth.azurecontainer.io/score'
api_key = '0kAlmBsCYoF37wnhtRWNdSxAhuPdX2ak' # Replace this with the API key for t
he web service
headers = {'Content-
Type':'application/json', 'Authorization':('Bearer '+ api_key)}

req = urllib.request.Request(url, body, headers)

try:
    response = urllib.request.urlopen(req)

    result = response.read()
    print(result)
except urllib.error.HTTPError as error:
    print("The request failed with status code: " + str(error.code))

    # Print the headers - they include the requert ID and the timestamp, which are
 useful for debugging the failure
    print(error.info())
    print(json.loads(error.read().decode("utf8", 'ignore')))
```