

Data Science and Machine Learning Essentials

Lab 2A – Acquiring Data in Azure Machine Learning

By Stephen Elston and Graeme Malcolm

Overview

In this lab, you will learn how to upload data files to Microsoft Azure Machine Learning (Azure ML), and how to use the **Join** module to combine data from multiple sources.

Note: This lab builds on knowledge and skills developed in Lab 1: *Getting Started with Azure ML*. If you have little experience with Azure ML, and you did not complete Lab 1, you are advised to do so before attempting this lab.

What You'll Need

To complete this lab, you will need the following:

- An Azure ML account
- A web browser and Internet connection
- The lab files for this lab

Note: To set up the required environment for the lab, follow the instructions in the **Setup** document for this course. Then download and extract the lab files for this lab.

Uploading a Data File to Azure ML

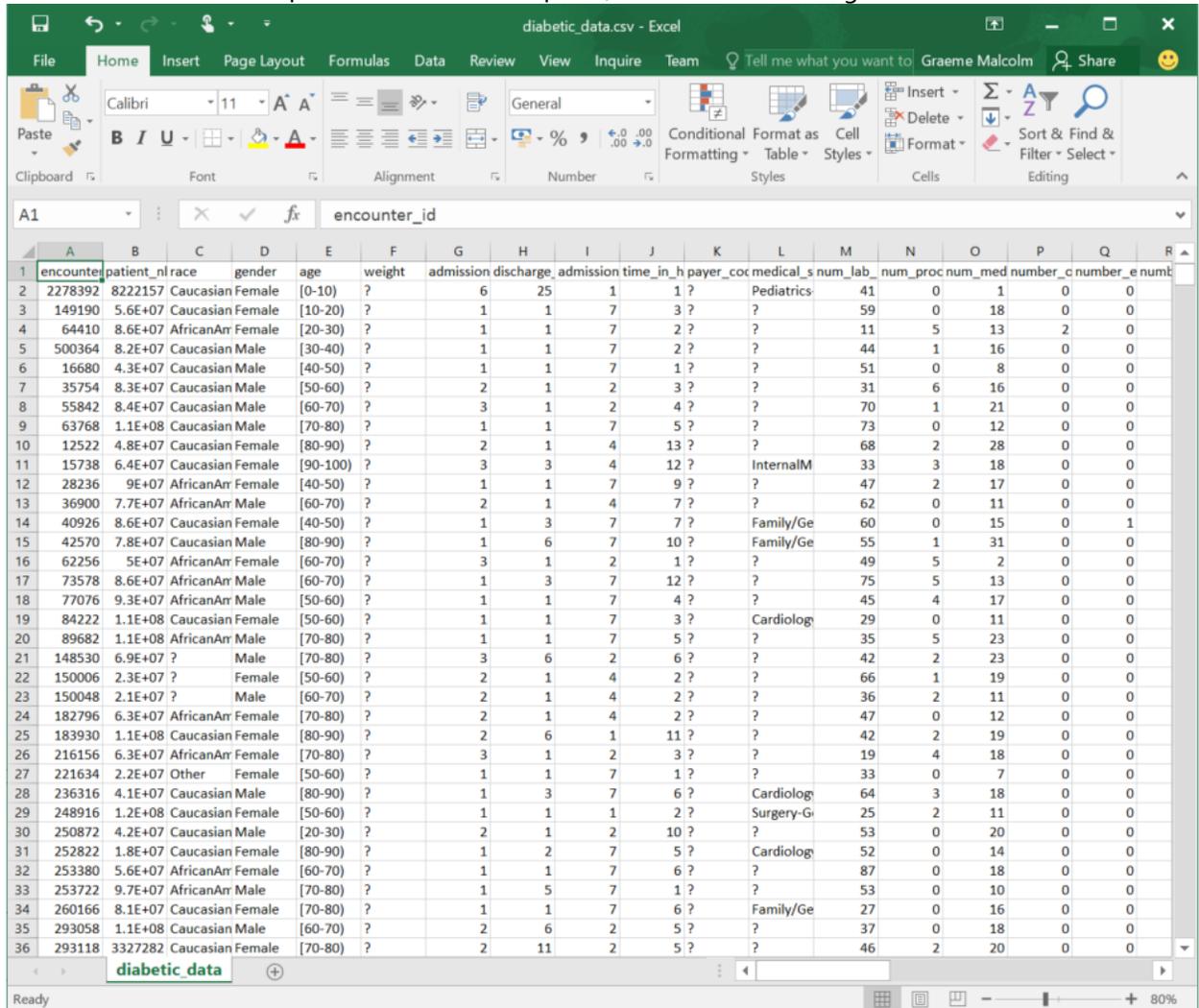
When you need to create an experiment based on your own data, or data you have obtained from a third-party, you must begin by uploading the data to Azure ML. In this exercise, you will upload a file containing data about diabetes patient appointments.

Note: The data used in this exercise was obtained from the University of California machine learning repository.

Lichman, M. (2013). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.

Examine the Data

1. Open the **diabetic_data.csv** file in the folder where you extracted the lab files, using either a spreadsheet application such as Microsoft Excel, or a text editor such as Microsoft Windows Notepad.
2. View the contents of the file, noting that it contains data on over 101,000 admissions and readmissions of diabetic patients at 130 US hospitals, as shown in this image:

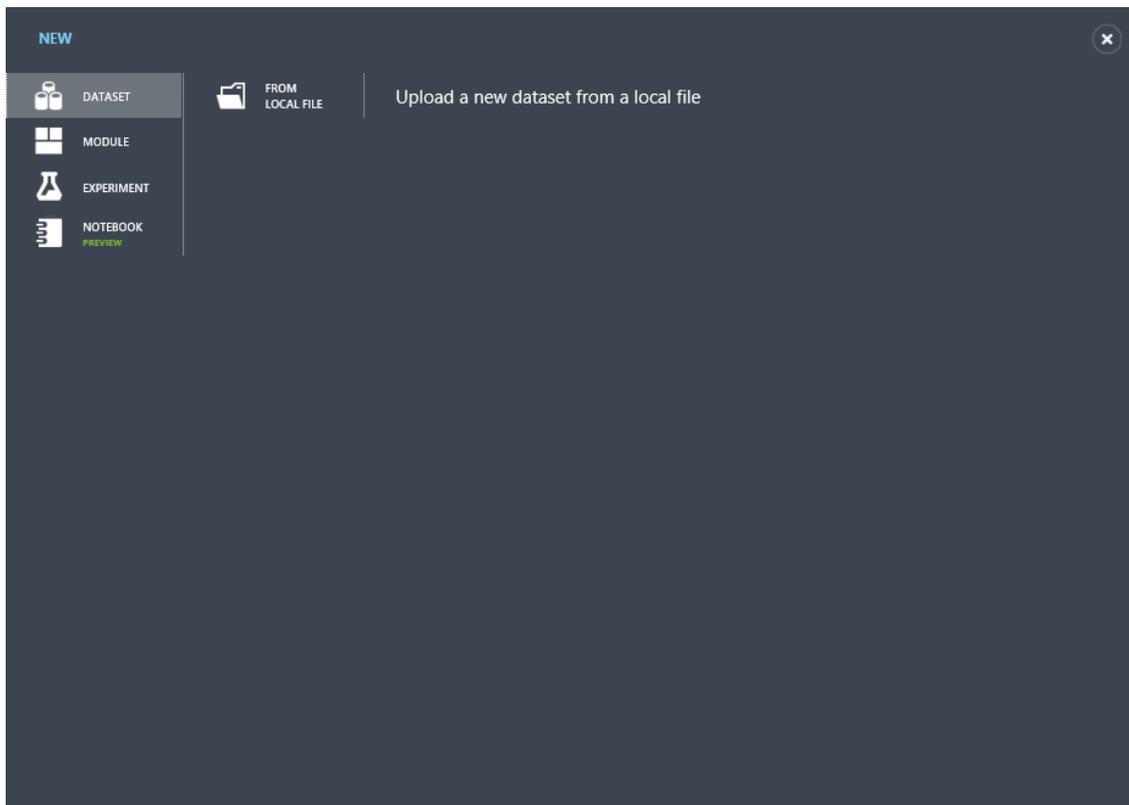


1	encounter_id	patient_id	race	gender	age	weight	admission_discharge	admission_time_in_h	payer_code	medical_specialty	num_lab	num_proc	num_med	number_c	number_e		
2	2278392	8222157	Caucasian	Female	[0-10]	?	6	25	1	1	?	Pediatrics	41	0	1	0	0
3	149190	5.6E+07	Caucasian	Female	[10-20]	?	1	1	7	3	?	?	59	0	18	0	0
4	64410	8.6E+07	AfricanArr	Female	[20-30]	?	1	1	7	2	?	?	11	5	13	2	0
5	500364	8.2E+07	Caucasian	Male	[30-40]	?	1	1	7	2	?	?	44	1	16	0	0
6	16680	4.3E+07	Caucasian	Male	[40-50]	?	1	1	7	1	?	?	51	0	8	0	0
7	35754	8.3E+07	Caucasian	Male	[50-60]	?	2	1	2	3	?	?	31	6	16	0	0
8	55842	8.4E+07	Caucasian	Male	[60-70]	?	3	1	2	4	?	?	70	1	21	0	0
9	63768	1.1E+08	Caucasian	Male	[70-80]	?	1	1	7	5	?	?	73	0	12	0	0
10	12522	4.8E+07	Caucasian	Female	[80-90]	?	2	1	4	13	?	?	68	2	28	0	0
11	15738	6.4E+07	Caucasian	Female	[90-100]	?	3	3	4	12	?	InternalM	33	3	18	0	0
12	28236	9E+07	AfricanArr	Female	[40-50]	?	1	1	7	9	?	?	47	2	17	0	0
13	36900	7.7E+07	AfricanArr	Male	[60-70]	?	2	1	4	7	?	?	62	0	11	0	0
14	40926	8.6E+07	Caucasian	Female	[40-50]	?	1	3	7	7	?	Family/Ge	60	0	15	0	1
15	42570	7.8E+07	Caucasian	Male	[80-90]	?	1	6	7	10	?	Family/Ge	55	1	31	0	0
16	62256	5E+07	AfricanArr	Female	[60-70]	?	3	1	2	1	?	?	49	5	2	0	0
17	73578	8.6E+07	AfricanArr	Male	[60-70]	?	1	3	7	12	?	?	75	5	13	0	0
18	77076	9.3E+07	AfricanArr	Male	[50-60]	?	1	1	7	4	?	?	45	4	17	0	0
19	84222	1.1E+08	Caucasian	Female	[50-60]	?	1	1	7	3	?	Cardiology	29	0	11	0	0
20	89682	1.1E+08	AfricanArr	Male	[70-80]	?	1	1	7	5	?	?	35	5	23	0	0
21	148530	6.9E+07	?	Male	[70-80]	?	3	6	2	6	?	?	42	2	23	0	0
22	150006	2.3E+07	?	Female	[50-60]	?	2	1	4	2	?	?	66	1	19	0	0
23	150048	2.1E+07	?	Male	[60-70]	?	2	1	4	2	?	?	36	2	11	0	0
24	182796	6.3E+07	AfricanArr	Female	[70-80]	?	2	1	4	2	?	?	47	0	12	0	0
25	183930	1.1E+08	Caucasian	Female	[80-90]	?	2	6	1	11	?	?	42	2	19	0	0
26	216156	6.3E+07	AfricanArr	Female	[70-80]	?	3	1	2	3	?	?	19	4	18	0	0
27	221634	2.2E+07	Other	Female	[50-60]	?	1	1	7	1	?	?	33	0	7	0	0
28	236316	4.1E+07	Caucasian	Male	[80-90]	?	1	3	7	6	?	Cardiology	64	3	18	0	0
29	248916	1.2E+08	Caucasian	Female	[50-60]	?	1	1	1	2	?	Surgery-Gi	25	2	11	0	0
30	250872	4.2E+07	Caucasian	Male	[20-30]	?	2	1	2	10	?	?	53	0	20	0	0
31	252822	1.8E+07	Caucasian	Female	[80-90]	?	1	2	7	5	?	Cardiology	52	0	14	0	0
32	253380	5.6E+07	AfricanArr	Female	[60-70]	?	1	1	7	6	?	?	87	0	18	0	0
33	253722	9.7E+07	AfricanArr	Male	[70-80]	?	1	5	7	1	?	?	53	0	10	0	0
34	260166	8.1E+07	Caucasian	Female	[70-80]	?	1	1	7	6	?	Family/Ge	27	0	16	0	0
35	293058	1.1E+08	Caucasian	Male	[60-70]	?	2	6	2	5	?	?	37	0	18	0	0
36	293118	3327282	Caucasian	Female	[70-80]	?	2	11	2	5	?	?	46	2	20	0	0

3. Close the data file without saving any changes.

Upload the Data File to Create a New Dataset in Azure ML

1. Open a browser and browse to <https://studio.azureml.net>. Then sign in using the Microsoft account associated with your Azure ML account.
2. Create a new blank experiment with the title **Diabetes Data**.
3. With the **Diabetes Data** experiment open, at the bottom left, click **NEW**. Then in the **NEW** dialog box, click the **DATASET** tab as shown in the following image.



- Click **FROM LOCAL FILE**. Then in the **Upload a new dataset** dialog box, browse to select the **diabetic_data.csv** file from the folder where you extracted the lab files on your local computer and enter the following details as shown in the image below, and then click the OK icon.
 - This is a new version of an existing dataset:** Unselected
 - Enter a name for the new dataset:** Diabetic Data
 - Select a type for the new dataset:** Generic CSV file with a header (.csv)
 - Provide an optional description:** Diabetes patient appointments.

- Wait for the upload of the dataset to be completed, and then on the experiment items pane, expand **Saved Datasets** and **My Datasets** to verify that the **Diabetic Data** dataset is listed.

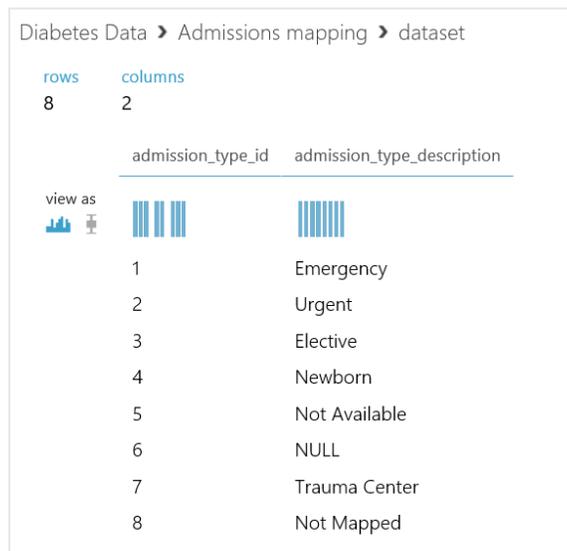
Visualize the Dataset in Azure ML

1. Drag the **Diabetic Data** dataset to the canvas for the **Diabetes Data** experiment.
2. Right-click the output port for the **Diabetic Data** dataset on the canvas and click **Visualize** to view the data in the dataset.
3. Verify that the dataset contains the data you viewed in the source file, and then close the dataset.

Using Reference Data

Occasionally, you may need to combine data from multiple data sources. In this case, the diabetic appointment data includes a code that indicates the admission type for the appointment, but it does not include a human readable description of the admission type. In this exercise, you will use the Azure ML **Join** module to combine rows in the **Diabetes Data** experiment with rows in a second dataset, which provides the admission type description values for the admission type codes.

1. In Azure ML Studio, ensure that the **Diabetes Data** experiment you created in the previous exercise is open.
2. Add a new dataset by uploading the **admissions_mapping.csv** file in the folder where you extracted the lab files. Name the new dataset **Admissions Mapping**.
3. In the **Diabetes Data** Azure ML experiment, search for the **Admissions Mapping** dataset and drag it to the canvas, next to the **Diabetic Data** dataset. Do not connect it to any other modules at this point, but visualize the dataset output port and note that the dataset contains rows that map **admission_type_id** codes to **admission_type_description** text values as shown in the following image:

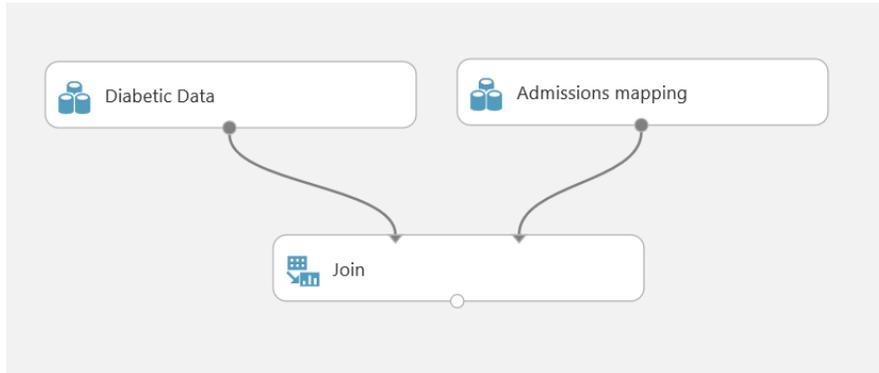


admission_type_id	admission_type_description
1	Emergency
2	Urgent
3	Elective
4	Newborn
5	Not Available
6	NULL
7	Trauma Center
8	Not Mapped

4. Close the dataset.
5. Search for the **Join** module, and drag it to the canvas below both the **Diabetic Data** dataset and the **Admission Mapping** dataset. Then connect the output ports from the **Diabetic Data** dataset and the **Admission Mapping** dataset to the **Dataset1** and **Dataset2** input ports of the **Join** module respectively.
6. Select the **Join** module, and in the **Properties** pane, set the following properties:
 - a. **Join key columns for L:** Launch the column selector and select **admission_type_id**.
 - b. **Join key columns for R:** Launch the column selector and select **admission_type_id**.
 - c. **Join type:** Left Outer Join
 - d. **Keep right key column:** *Unselected*

Note: The key columns in this example have the same name, but this is not a requirement to join columns from two datasets. By using a left outer join, The **Join** module will retain any rows in the left dataset (the diabetic data) that do not have a matching **admission_type_id** column in the right dataset (the admissions mapping data).

- Verify that your experiment resembles the following image, and then save and run it.



- When the experiment has finished running, visualize the **Results dataset** output port of the **Join** module, and verify that it contains the cleansed patient admission data from the original **Diabetic Data** dataset and the **admission_type_description** column from the **Admissions Mapping** dataset as shown in the following image.

mepiride-oglitazone	metformin-rosiglitazone	metformin-pioglitazone	change	diabetesMed	readmitted	admission_type_description
	No	No	No	No	NO	NULL
	No	No	Ch	Yes	>30	Emergency
	No	No	No	Yes	NO	Emergency
	No	No	Ch	Yes	NO	Emergency
	No	No	Ch	Yes	NO	Emergency
	No	No	No	Yes	>30	Urgent
	No	No	Ch	Yes	NO	Elective
	No	No	No	Yes	>30	Emergency
	No	No	Ch	Yes	NO	Urgent
	No	No	Ch	Yes	NO	Elective
	No	No	No	Yes	>30	Emergency
	No	No	Ch	Yes	<30	Urgent
	No	No	Ch	Yes	<30	Emergency
	No	No	No	Yes	NO	Emergency
	No	No	No	Yes	>30	Elective
	No	No	Ch	Yes	NO	Emergency
	No	No	Ch	Yes	<30	Emergency
	No	No	No	Yes	NO	Emergency

- Close the results dataset.
- If necessary, save the experiment. Then close your browser. The datasets you uploaded will remain available in your Azure ML account.

Summary

This lab was designed to teach you how to upload data to Azure ML and to use the Join module to combine data from multiple data sources.

Note: The experiment created in this lab is available in the Cortana Analytics library at <http://gallery.cortanaanalytics.com/Collection/5bfa7c8023724a29a41a4098d3fc3df9>.