# Azure Mini Project: Spark on HDInsight

**Estimated Time: 2-3 Hours**

## Overview

Today, data is being generated and stored in staggering amounts at unprecedented velocities. To make things more complicated, this data is generated in an ever-expanding variety of formats.

In order to gain actionable insights into Big Data sources, new tools need to be leveraged that allow the data to be cleaned, analyzed, and visualized quickly and efficiently. Azure HDInsight provides one solution to this problem. HDInsight makes it simple to create high-performance computing clusters equipped with Apache Spark and Spark ecosystem tools. This tool allows you to focus your time on the data itself rather than deploying, installing, configuring and maintaining the resources required to analyze the data.

Apache Spark is an open-source parallel-processing platform that excels at running large-scale data analytics jobs. Spark's combined use of in-memory and disk data storage delivers performance improvements that allow it to process some tasks up to 100 times faster than Hadoop. Azure makes deploying Spark clusters relatively simple and allows you to work on your data analysis sooner.

In this mini-project, you'll get hands-on experience with Apache Spark for Azure HDInsight. After provisioning a Spark cluster, you will use the Microsoft Azure Storage Explorer to add a Jupyter Notebook to the cluster. You will then use the notebook to explore and analyze Walmart stock data. Your goal is to learn how to create and utilize your own Spark clusters, provision them with Azure, and get a working introduction to Spark data analytics.

**Learning Goals**
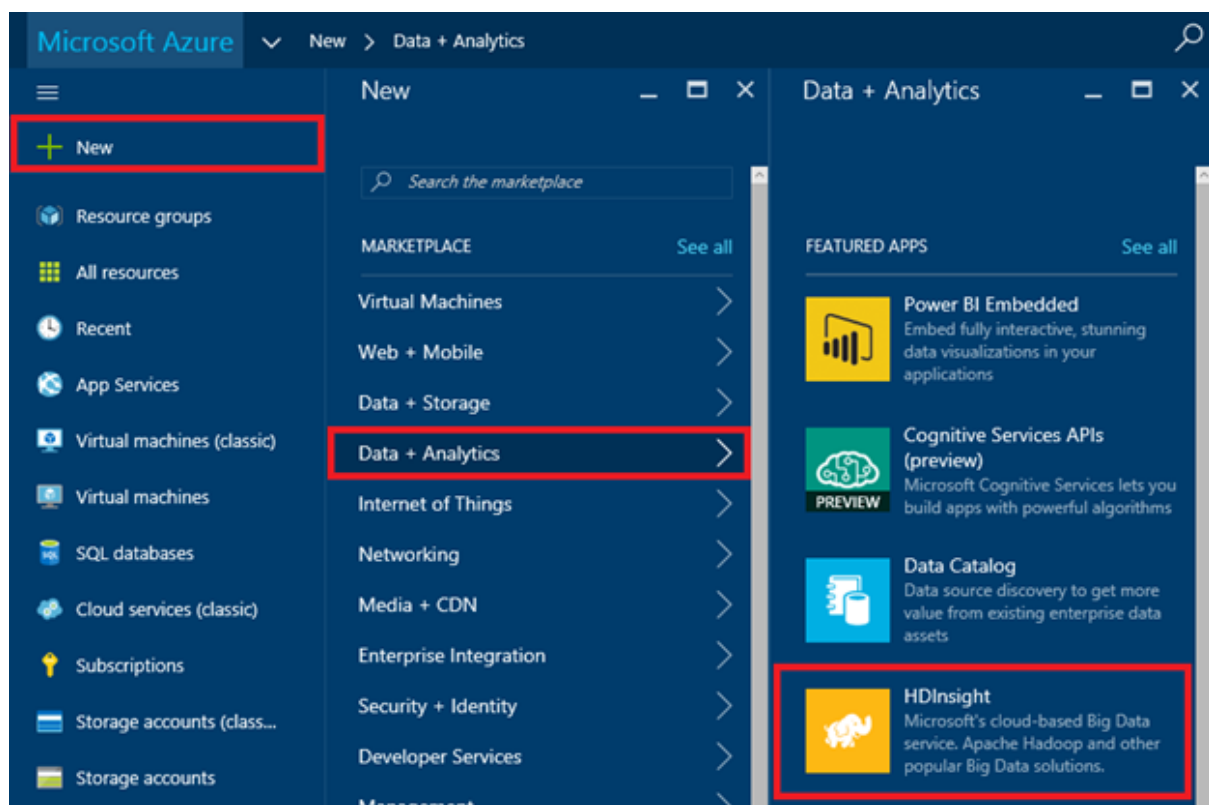
In this mini-project, you will learn how to:

- Deploy an HDInsight Spark cluster
- Work with content stored in Azure Blob Storage and accessed by the Spark cluster as an HDFS volume
- Use a Jupyter Notebook to interactively explore a large dataset
- Delete a Spark cluster to avoid incurring unnecessary charges

Check the Deliverables at the end of this document before you begin working.

# Exercise 1: Create a Spark Cluster on HDInsight

First create an HDInsight cluster running Apache Spark.

1. Go to the Azure Portal and sign in using your Microsoft account.

2. Click **+ New** in the upper-left corner. Then click **Data + Analytics**, followed by **HDInsight**.



*Creating an HDInsight cluster*

3. In the **Cluster Name** box, enter a unique DNS name for the cluster and make sure a green check mark appears next to it indicating that the name is valid and unique.

*Make your cluster name as unique as possible by including birth dates, initials, and anything else you care to add to ensure no one else chooses the same name. The name you entered may be unique right now, but it might NOT be unique a few minutes into the deployment.*



*Specifying the cluster name*

4.  Click **Select Cluster Type** to open a "Cluster Type configuration" blade. In that blade, select **Spark** as the **Cluster Type**. Under **Version**, select the latest Spark version, which is **1.6.1 (HDI 3.4)** at the time of this writing. Then click **Standard** to select a **Cluster Tier**, and finish up by clicking the **Select** button at the bottom of the blade.



*Specifying the cluster type*

5.  Click **Credentials** to open a "Cluster Credentials" blade. Leave **Cluster Login Username** set to "admin" and make sure **SSH Authentication Type** is set to **Password**. Enter "sshuser" (without quotation marks) for the **SSH Username** and "A4rsparkdemo!"

(again without quotation marks) for the **Cluster Login Password** and **SSH Password**. Then click the **Select** button at the bottom of the blade.



*Specifying cluster credentials*

6.  Click **Data Source** to open a "Data Source" blade. Leave **Selection Method** set to **From all subscriptions** and enter a unique storage-account name in the box below **Create a new storage account**. Once more, make the name unique by including birth dates or other values unlikely to be used by someone else. For **Choose Default Container**, enter "sparklab" (without quotation marks). Select the **Location** nearest you, and then click the **Select** button at the bottom of the blade.

*Specifying the data source*

7.  Click **Node Pricing Tiers** to open a "Pricing" blade. Make sure **Number of Worker nodes** is set to **4** and accept the default values everywhere else. Then click the **Select** button at the bottom of the blade.

*Specifying the node pricing tier*

8. Select **Create new** under **Resource Group** and enter the resource-group name "SparkLabResourceGroup" (without quotation marks). Then click the **Create** button at the bottom of the blade to start deploying the cluster.
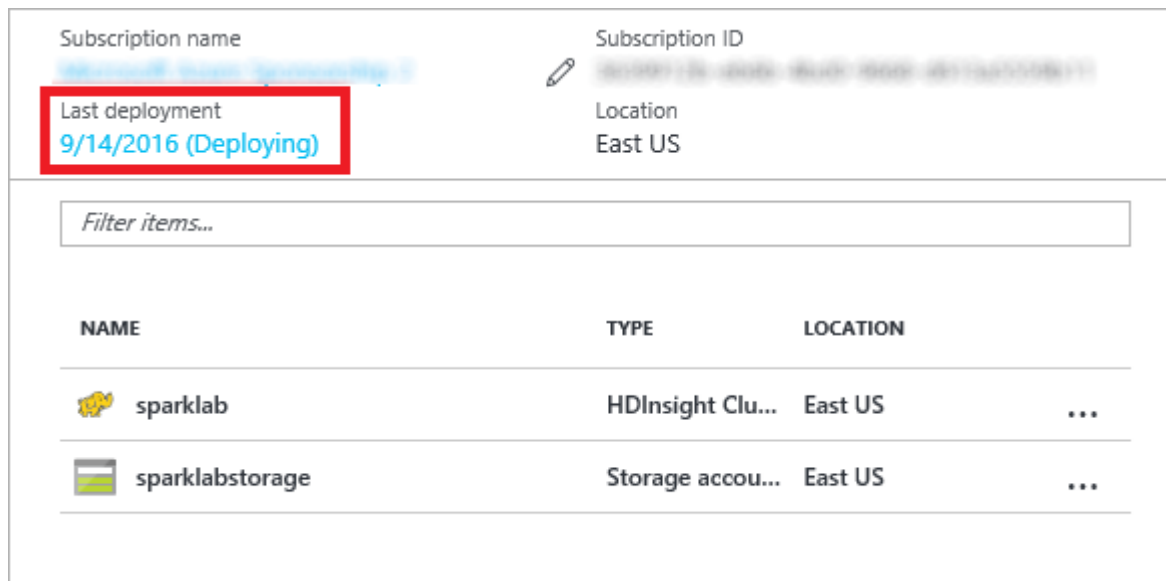


*Specifying a resource group and creating the cluster*

9.  Deploying an HDInsight cluster can take 20 minutes or more. You can monitor the status of the deployment by opening the resource group's blade. Click **Resource group** in the ribbon on the left side of the portal, and then click the resource group name ("SparkLabResourceGroup") to open the blade. "Deploying" will change to "Succeeded" when the deployment has completed successfully.

*Click the browser's **Refresh** button every few minutes to update the deployment status. Clicking the **Refresh** button in the resource-group blade refreshes the list of resources in the resource group, but does not reliably update the deployment status.*

*Monitoring the deployment*

**Summary**

In this exercise, you learned how to provision an HDInsight Spark cluster on Azure and some of the available configurations you can choose from. Wait for the deployment to finish, and then proceed to the next exercise.
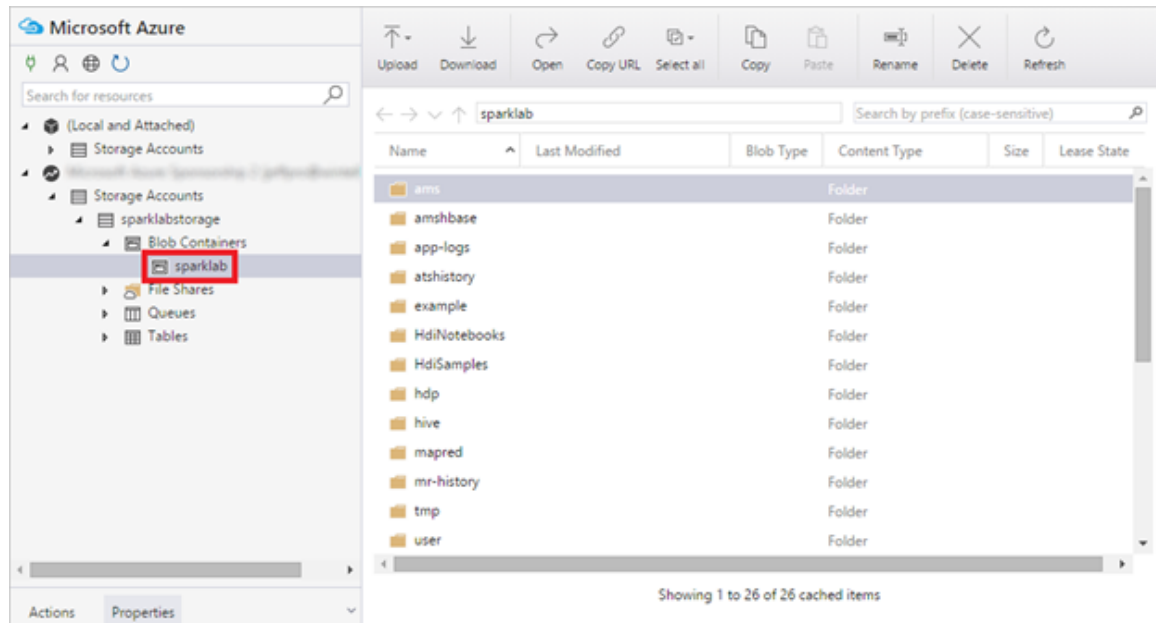
# Exercise 2: Upload Jupyter Notebook to the cluster

You will be using Jupyter Notebooks to do the data-exploration and analysis portions of this lab. A notebook and data file have been prepared for you **(download them here: .ipynb file, data file)** and need to be uploaded to your cluster. In this exercise, you'll use the cross-platform Microsoft Azure Storage Explorer to upload the notebook. If Storage Explorer isn't installed on your computer, take the time to install it now.

1.  Start the Microsoft Azure Storage Explorer. If you're prompted for credentials, sign in with the username and password for your Microsoft account.
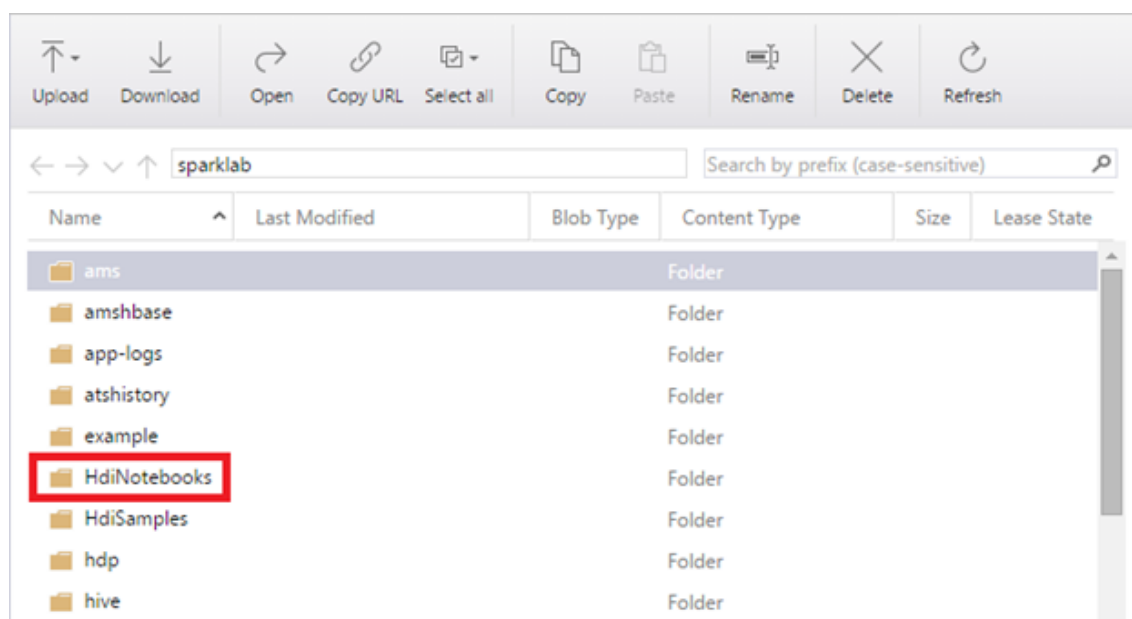
2. In the Storage Explorer window, find the storage account you created when you deployed the Spark cluster (in Step 6 of the previous exercise). Expand the list of items underneath that storage account and click the small arrow next to **Blob Containers** to show a list of containers. Then click the container named "sparklab." It contains several folders created during the provisioning process. sparklab is the root folder for your Spark cluster.
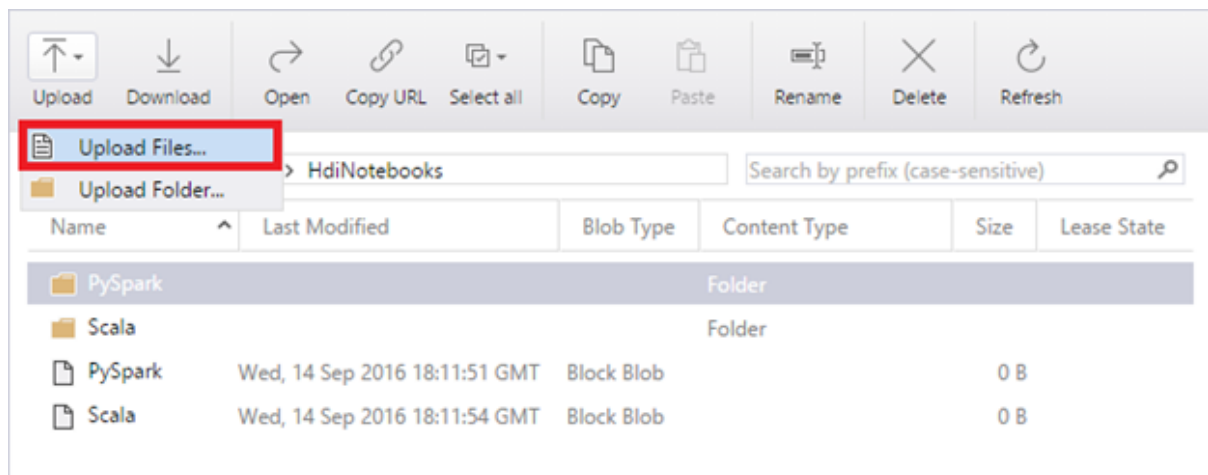
*Contents of the "sparklab" container*

3. Double-click the folder named "HdiNotebooks." This is the root folder for the cluster's Jupyter notebooks.
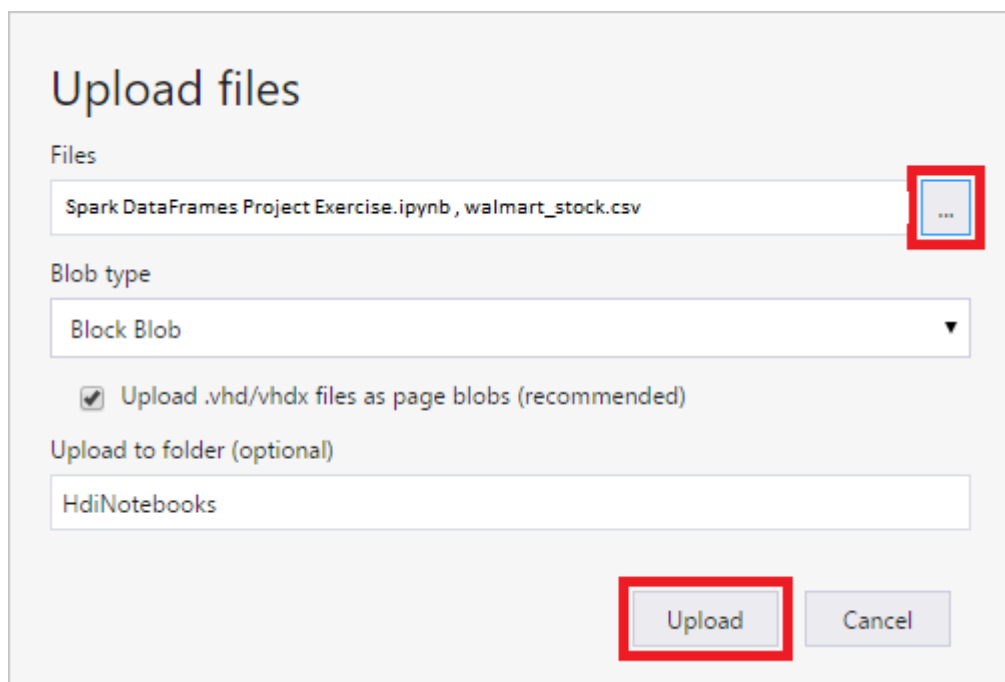
*Opening the "HdiNotebooks" folder*

4. Click the **Upload** button and select **Upload Files** from the menu.
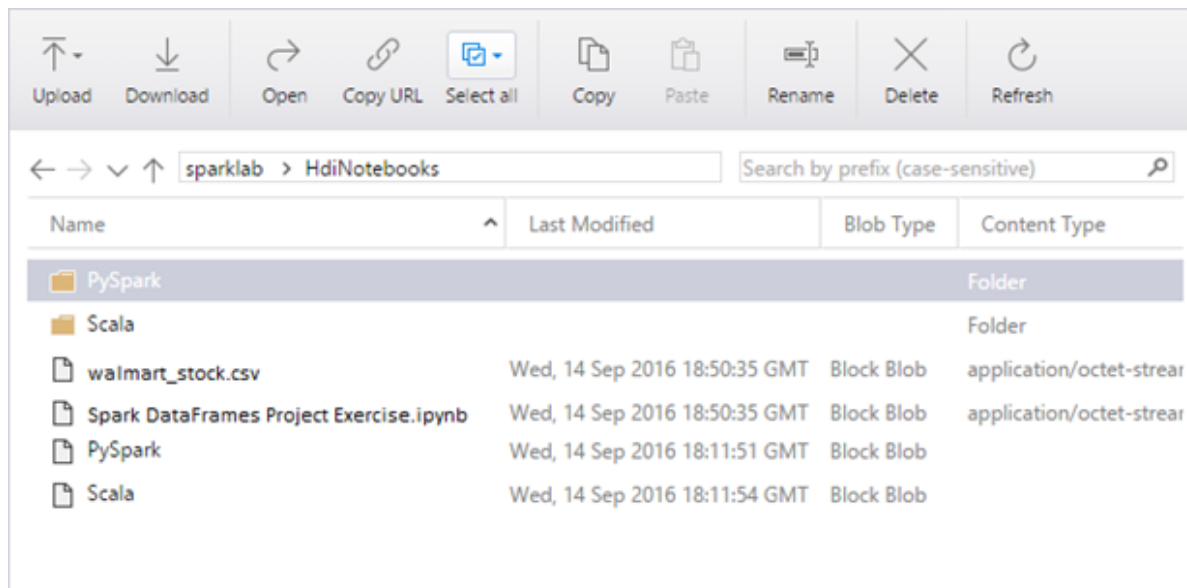


*Uploading notebook to blob storage*

5. Click the **...** button to the right of the field labeled "Files." In the ensuing dialog, navigate to this lab's "resources" subdirectory and select the .ipynb file and the walmart stock data file in that subdirectory. Then close the dialog and click the **Upload** button.



*Uploading .ipynb file and Walmart Stock Data File*

6. Confirm that all three files were uploaded to the "HdiNotebooks" folder.



*Uploaded notebook and walmart stock data*

**The notebook and data is uploaded and ready to go. Let's put them to work!**

# Exercise 3: Work with Jupyter Notebooks

Jupyter Notebooks are extraordinarily useful for data professionals. Jupyter supports several programming languages through the use of installable interpreters called *kernels*. Spark clusters on HDInsight include the Spark and PySpark kernels for Scala and Python, respectively.

In this exercise, you will learn how to access Jupyter notebooks in your Spark cluster and acquire basic skills for using them.

1. In the Azure Portal, return to the blade for the resource group ("SparkLabResourceGroup") that contains the cluster. In the list of resources that belong to the resource group, click the HDInsight cluster.

*Opening the cluster*

2.   In the blade for the HDInsight Spark cluster, click **Cluster Dashboards.**



*Opening the cluster dashboards*

3.   In the ensuing blade, click **Jupyter Notebook**.



*Opening a Jupyter notebook*

4.   When prompted for a username and password, log in with your cluster credentials ("admin" and "A4rsparkdemo!") from Exercise 1, Step 5.

*If you mistype the password and are greeted with a 403 Forbidden error, start a new incognito or private browsing session, go to the Azure Portal, and open the Jupyter notebook again.*

5.   A new browser window (or tab) will open showing the Jupyter notebooks in your cluster. Here, you can manage your notebooks, upload new ones, and more. You can also see which notebooks are currently "running", meaning they are currently consuming resources in your Spark cluster. Confirm that you see the notebook you uploaded in Exercise 2. Then click **Spark Dataframes Project Exercises.ipynb** to open that notebook.



*Opening Notebook*

6.   Jupyter notebooks consist of a series of cells into which you can insert commands, HTML, or Markdown text. The notebook you opened contains the remaining instructions for this exercise. Follow the instructions in the notebook to complete Exercise 3.
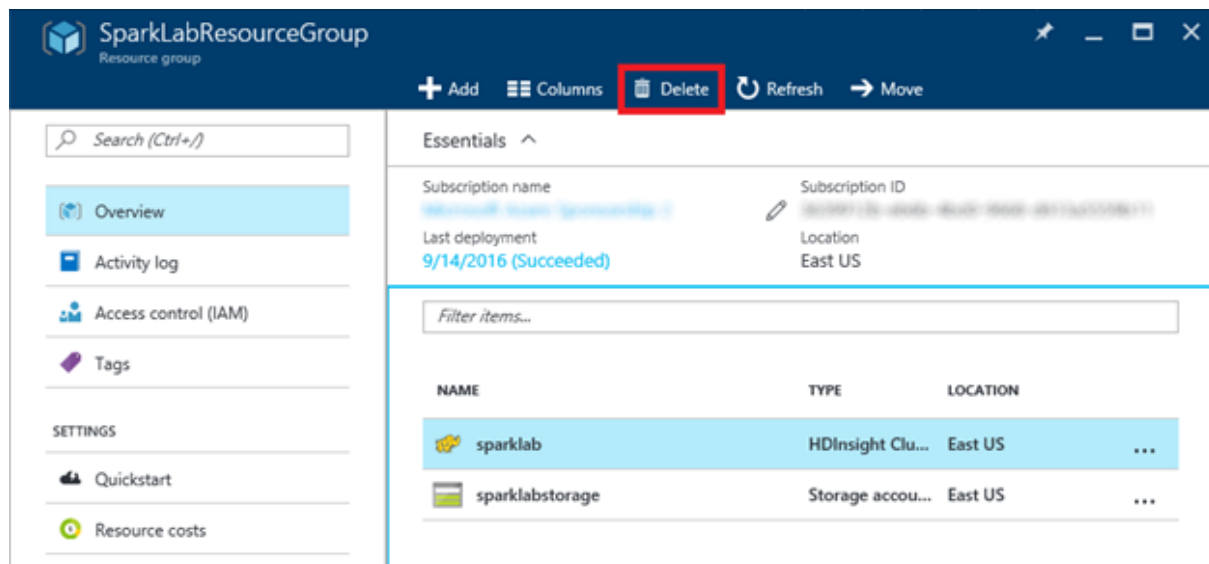
*Working with Jupyter notebook*

Once you have completed the steps in the notebook, save and checkpoint the notebook then close the browser window in which the notebook is displayed and return to the Azure Portal.

7.  Download the **Spark Dataframes Project Exercises.ipynb** to your computer by right clicking the file in Azure Storage Explorer and clicking 'Download'. You will submit the Jupyter Notebook to your mentor .

# Exercise 4: Remove the HDInsight Spark cluster

When you are finished using an HDInsight Spark cluster, you should delete it. **You are charged for it while it exists, regardless of whether it's doing any work.** In this exercise, you will delete the resource group created in Exercise 1 when you created the cluster. Deleting the resource group deletes everything in it and prevents any further charges.

1.  In the Azure Portal, open the blade for the "SparkLabResourceGroup" resource group that holds the cluster. Then click the **Delete** button at the top of the blade.



*Deleting a resource group*

2.  For safety, you are required to type in the resource group's name. **Once deleted, a resource group cannot be recovered.** Type the name of the resource group. Then click the **Delete** button to remove all traces of this lab from your account.

After a few minutes, the cluster and all of its resources will be deleted. Billing stops when you click the **Delete** button, so you're not charged for the time required to delete the cluster. Similarly, billing doesn't start until a cluster is fully and successfully deployed.

## Summary

Here is a summary of what you learned in this lab:

- Apache Spark for Azure HDInsight is Microsoft Azure's implementation of Hadoop, Spark, and several other Big Data tools
- The Azure Portal makes it easy to create, configure, and delete HDInsight Spark clusters
- HDInsight Spark clusters come with Jupyter preinstalled
- Jupyter Notebooks provide a powerful means for querying and analyzing data
- HDInsight Spark clusters should be deleted when they're no longer needed to avoid incurring unwanted charges

With Apache Spark for Azure HDInsight, high-performance computing clusters with all the tools you need to handle big data are just a few button clicks away. It's just one example of why cloud computing is changing the face of research.

# Deliverables

1. Jupyter Notebook with code to achieve stated objective in each cell

2. A PPT/Word-doc with screenshots of each step as performed according to this lab