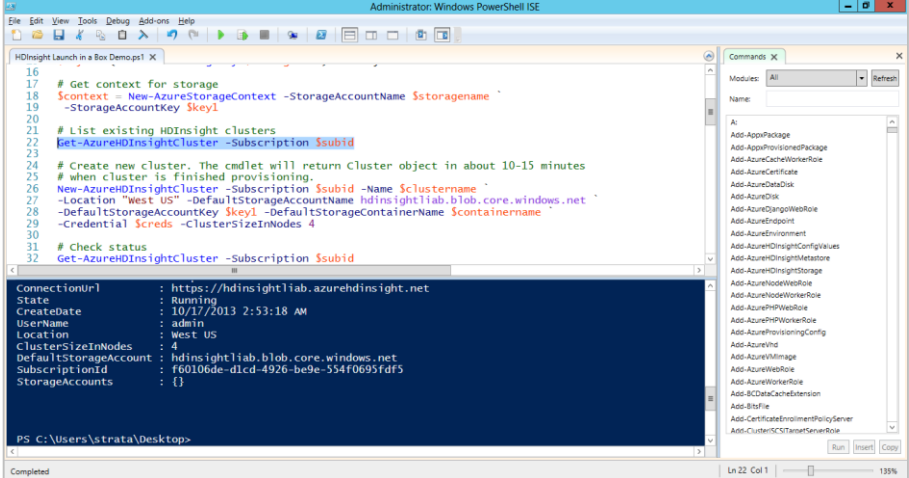# Working with Hive in HDInsight

# Contents

Contents

# Introduction

In this hands-on lab, you will setup a Microsoft Azure storage account along with an HDInsight Hadoop cluster. You will then upload a data file that was processed using a map reduce job that was created using the TR19 Exploring MapReduce in HDInsight using .NET lab. Then, you will use Hive to query the data on your HDInsight cluster help you learn the essentials of Hive.

The hands-on-lab requires a Microsoft Azure account.

During TechReady, please raise your hand and ask a Technical Learning Guide for a pre-configured Azure trial account.

# First Time Setup Instructions

Configure the PowerShell environment as follows:

1. Create a directory called C:\data.

2. Go to the D:\TR19HiveLab folder. Then right click on the file "TR19 Hive Lab.ps1" and choose the Edit command to launch the PowerShell script in the Windows PowerShell ISE.

3. Log into your Microsoft Azure Management portal with your Live ID at https://manage.windowsazure.com. Choose the option to "Keep me signed in".

4. Select the command `Get-AzurePublishSettingsFile` in the editor and press F8 to execute the selection. This will launch a new tab in IE and then prompt you to open or save the .publishsettings file. Select the Save option to copy the file into the c:\data folder.

5. Edit the `Import-AzurePublishSettingsFile` command in the editor include the full file path of the .publishsettings file in your C;\data folder. For example: `Import-AzurePublishSettingsFile "C:\data\Greg_AzurePlatform_Subscription-10-22-2013-credentials.publishsettings"` and then press F8 to execute the command. This completes the one time setup for PowerShell so that it can talk with Microsoft Azure.

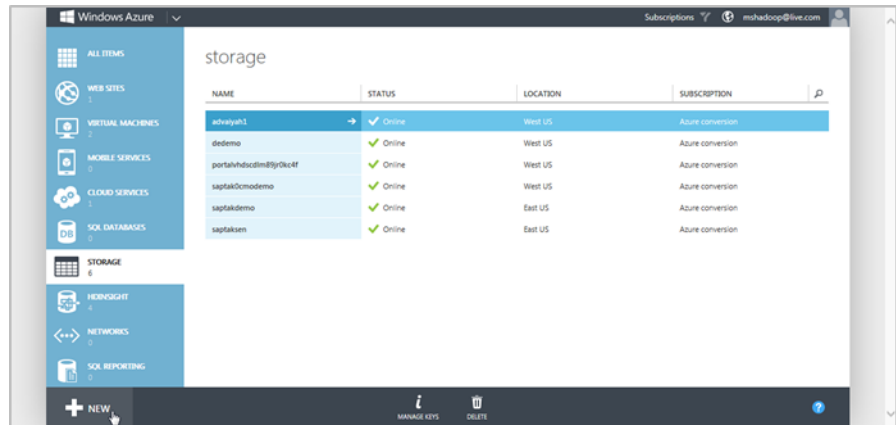6. Verify the subscription by executing the following command: `Get-AzureSubscription –Default`

NOTE: If you completed the TR19 Exploring MapReduce in HDInsight using .NET lab and kept your HDInsight cluster, storage account and student container, you can skip to section **Loading MapReduce data into Azure Storage**.

# Creating an Azure storage account

In this task, you will create a storage account and a container within the account you will use with the HDInsight cluster.

NOTE: Skip this step if you are using a shared account

1. Log into the Microsoft Azure management portal at https://manage.windowsazure.com/.

2. Select the STORAGE page and click NEW in the bottom left corner.



3. Click **QUICK CREATE**.



4. For the URL, you need to provide a unique storage account name that is from 3-15 characters. In this example, the storage account name will be the same as the HDInsight cluster name. You will need to replace this name with one that is unique for your training session. All of the exercises use the name: <storage_account_name>. You can search and replace this document <storage_account_name> with the one you choose.
   Complete the required fields:

   a. URL: <storage_account_name>

b. Location: **West US**.

   Note, there are many to choose from, but you must select a location that supports HDInsight. These include West US, East US and North Europe at this time.
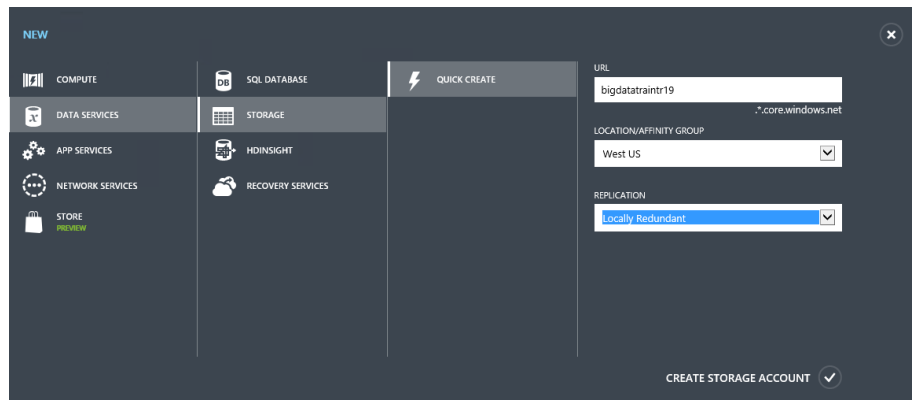
c. GEO-REPLICATION: Deselect it.

d. Click CREATE STORAGE ACCOUNT to complete the operation.



This operation can take between 1 to 5 minutes to complete.

# Creating an HDInsight Cluster

In this section, you will create an HDInsight cluster that aligns with the storage account you created in the previous exercise.

**Key Points**

The Azure Management Portal provides a user interface for creating a basic HDInsight cluster.

**Demonstration**

In this first demo, you will create a new HDInsight cluster. This step takes about 10 minutes to complete.

1. Navigate to the management portal – https://manage.windowsazure.com.

2. Select the HDINSIGHT page and click on the NEW button in the lower left corner.



3. You have two options: QUICK CREATE or CUSTOM CREATE. As of June 10, 2014, the default version created with QUICK CREATE is 3.0.
   NOTE: Do not select the 3.1 preview edition as there are some incompatibilities with the PowerShell cmdlets.

   Select CUSTOM CREATE with the following parameters:

   a. CLUSTER NAME: <cluster_name>

   b. DATA NODES: 1

   c. HDINSIGHT VERSION: 3.0 (HDP 2.0, Hadoop 2.2)

d. REGION: West US (needs to match your storage account)

e. Click next page arrow.



4. For the Configure Cluster User page, enter in the following values:

a. USER NAME: ADMIN

b. PASSWORD: Pass@word12

c. CONFIRM PASSWORD: Pass@word12

d. Click the next page arrow.

NEW HDINSIGHT CLUSTER

## Configure Cluster User

USER NAME

ADMIN

PASSWORD                           CONFIRM PASSWORD

••••••••••••                      ••••••••••••

Enter the Hive/Oozie Metastore ☐ ⊙

1                                                    ← →    3

5. For the Storage Account page, enter in the following:

    a. STORAGE ACCOUNT: Use Existing Storage

    NOTE: If you don't see your storage account available for selection, your default location may be not be the same as West US. In which case, use the CUSTOM CREATE option.

    b. ACCOUNT NAME: select your account name from the region.

    c. DEFAULT CONTAINER: Create default container

    d. ADDITIONAL STORAGE ACCOUNTS: 0 (this is used to link multiple storage accounts to your cluster)

    e. Click the completion check box.

This operation takes about 10 minutes to complete

6.  Review the progress by clicking on the moving bar chart to
    see details. When the operation is compete, the status for the
    cluster shows as Running.

### Summary

Using the Azure management portal is a convenient way to create a
cluster. The CUSTOM CREATE operation offers more flexibility.

# Loading MapReduce data into Azure Storage

In this exercise, you upload the MapReduce part-00000 file that was initially created during the TR19 Exploring MapReduce in HDInsight using .NET lab. A copy of the file is located in the D:\TR19HiveLab directory.

## Key Points

- The Azure Storage cmdlets are the recommended way for working with data files with your storage account because you can create and use pseudo directories within a container.

## Demonstration

1. Go to your open "TR19 Hive Lab.ps1" using the PowerShell ISE and scroll down to line 18.

2. Change the line `$storagename = "<storage_account_name>"` to match your storage account name. Then change the line `$clustername = "<cluster_name>"` to match your cluster name.

```
TR19 Hive Lab.ps1* X
13   # Setup variables used for Azure storage
14   $subid = (Get-AzureSubscription –Current).SubscriptionId
15   $subscriptionname = (Get-AzureSubscription –Current).SubscriptionName
16   Select-AzureSubscription $subscriptionname
17
18   $storagename = "bigdatatraintr19"
19   $clustername = "bigdatatraintr19"
20
21   # Get context for storage
22   $key1 = (Get-AzureStorageKey $storagename).Primary
23   $context = New-AzureStorageContext -StorageAccountName $storagename `
24     -StorageAccountKey $key1
```

3. Select the following lines of code in the editor and press F8 to setup the storage variables and set content for the storage account.

```
TR19 Hive Lab.ps1* X
13    # Setup variables used for Azure storage
14    $subid = (Get-AzureSubscription -Current).SubscriptionId
15    $subscriptionname = (Get-AzureSubscription -Current).SubscriptionName
16    Select-AzureSubscription $subscriptionname
17
18    $storagename = "bigdatatraintr19"
19    $clustername = "bigdatatraintr19"
20
21    # Get context for storage
22    $key1 = (Get-AzureStorageKey $storagename).Primary
23    $context = New-AzureStorageContext -StorageAccountName $storagename `
24     -StorageAccountKey $key1
25
26    # Create your student storage account container
```

4. Locate the line - `$containername = "<student_id>"` and replace the student number with your assigned number or just use `"student01"`.

```
TR19 Hive Lab.ps1* X
25
26    # Create your student storage account container
27    $containername = "student01" # Example: student01
28    # SKIP this command if you kept your environment from the
29    # TR19 Exploring MapReduce in HDInsight using .NET lab
30    New-AzureStorageContainer -Name $containername -Context $context `
31     -Permission Off
32
```

5. Create the container for the census files by selecting the following lines of code and pressing F8 to execute them.
   NOTE: Skip this step if you retained your HDInsight cluster environment from the TR19 Exploring MapReduce in HDInsight using .NET lab.

```
TR19 Hive Lab.ps1* X
25
26    # Create your student storage account container
27    $containername = "student01" # Example: student01
28    # SKIP this command if you kept your environment from the
29    # TR19 Exploring MapReduce in HDInsight using .NET lab
30    New-AzureStorageContainer -Name $containername -Context $context `
31     -Permission Off
32
```

6. Select the following lines of code to copy the MapReduce part-00000 file to your student container.
   NOTE: Skip this step if you retained your HDInsight cluster environment from the TR19 Exploring MapReduce in HDInsight using .NET lab.

```
TR19 Hive Lab.ps1* ✕

31    -Permission Off
32
33    # Copy the MapReduce file from local workstation to Blob container
34    # SKIP this command if you kept your environment from the
35    # TR19 Exploring MapReduce in HDInsight using .NET lab
36    Set-AzureStorageBlobContent -File "D:\TR19HiveLab\part-00000" `
37       -Container $containername `
38       -Blob "data/mr-results/part-00000" -context $context
39
40    $rootpart = "wasb://$containername@$storagename.blob.core.windows.net"
41
42    #Create table and run test query
43    $querystring = `
44    "CREATE EXTERNAL TABLE working_te_census_info(" + `
45       "state STRING,"+`
```

```
PS D:\TR19HiveLab> $containername = "student01" # Example: student01

PS D:\TR19HiveLab> Set-AzureStorageBlobContent -File "D:\TR19HiveLab\part-00000" `
 -Container $containername `
 -Blob "data/mr-results/part-00000" -context $context


   Container Uri: https://bigdatatraintr19.blob.core.windows.net/student01

Name                    BlobType    Length          ContentType         LastModified
----                    --------    ------          -----------         ------------
data/mr-results/p...    BlockBlob   1583499         application/octe... 6/10/2014 7:37:2..


PS D:\TR19HiveLab>
```

## Summary

You now have the environment prepared to run Hive queries.

# Running Hive commands as a job

In this exercise, you will learn the essentials how to execute HiveQL queries against HDInsight using the PowerShell interfaces. You will learn the similarities of HiveQL with SQL Server Transact-SQL.

## Initial setup for this exercise

Go to the Windows PowerShell ISE with TR19 Hive Lab.ps1 already opened.

There are several ways to run Hive commands within PowerShell. In this first example, you will learn how to run a set of HiveQL commands as a job.

In this exercise, you will create a Hive table and perform a select statement against the table using the following sequence:

- Create a $querystring with the commands. Each command uses a semi-colon end of line delimiter.

- Create the Hive Job Definition.

- Start the HDInsight job

- Wait for the results

- View the job output.

## Running Hive statements as a job

1. Select the following command to tell Hive where the data resides on Azure blob storage. By default, Hive will expect to locate files in the default container for the HDInsight cluster.

```
TR19 Hive Lab.ps1* ×
39
40   $rootpart = "wasb://$containername@$storagename.blob.core.windows.net"
41
42   #Create table and run test query
43   $querystring = `
44   "CREATE EXTERNAL TABLE working_te_census_info(" + `
45       "state STRING "+ `
```

2. Select the following command and press F8 to create the query variable that will be used to create the Hive table based on the map-reduce results **part-00000** file.

```
TR19 Hive Lab.ps1* ×
41
42   #Create table and run test query
43   $querystring = `
44   "CREATE EXTERNAL TABLE working_te_census_info(" + `
45       "state STRING,"+ `
46       "county STRING,"+ `
47       "agegrp STRING,"+ `
48       "total_population STRING) "+ `
49       "ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' "+ `
50       "STORED AS TEXTFILE LOCATION '$rootpart/data/hive/working_te_census_info';"+ `
51   "LOAD data inpath '$rootpart/data/mr-results/part-00000' " + `
52       "OVERWRITE into table working_te_census_info;" + `
53   "SELECT * from working_te_census_info ORDER BY total_population DESC limit 10;"|
54
55   ### Create a Hive job definition
```

In reviewing the syntax, you will see that the table was defined with the EXTERNAL option. This tells Hive to keep the data if the table is dropped. You'll also see that the data is stored as a TEXTFILE with comma separated fields. You can submit multiple HiveQL statements as long as they are separated by a semi-colon. Notice that the STRING data type does not include a length.

3. Select the following four commands in the script and press F8 to submit the Hive job and return the results.

```
TR19 Hive Lab.ps1* ×
54
55   ### Create a Hive job definition
56   $HiveJobDefinition = New-AzureHDInsightHiveJobDefinition `
57       -Query $querystring
58
59   ### Submit the job to the cluster
60   $HiveJob = Start-AzureHDInsightJob `
61       -Cluster $clustername -JobDefinition $HiveJobDefinition
62
63   ### Wait for the Hive job to complete
64   $HiveJob | Wait-AzureHDInsightJob `
65       -WaitTimeoutInSeconds 3600
66
67   ### Print the standard error and the standard output of the Hive job.
68   Get-AzureHDInsightJobOutput -Cluster $clustername `
69       -JobId $HiveJob.JobId -StandardOutput
70
```

```
Waiting for jobDetails : job_1402420269563_0003.
  Running : .
```

```
### Wait for the Hive job to complete
$HiveJob | Wait-AzureHDInsightJob `
     -WaitTimeoutInSeconds 3600

### Print the standard error and the standard output of the Hive job.
Get-AzureHDInsightJobOutput -Cluster $clustername `
     -JobId $HiveJob.JobId -StandardOutput
```

You could have piped the commands together similar to the way you executed the streaming map-reduce job in the prior exercise.

4. Scroll through the console results to see the output.

```
TR19 Hive Lab.ps1* ×
   64   $HiveJob | Wait-AzureHDInsightJob
   65          -WaitTimeoutInSeconds 3600
   66
   67   ### Print the standard error and the standard output of the Hive job.
   68   Get-AzureHDInsightJobOutput -Cluster $clustername `
   69          -JobId $HiveJob.JobId -StandardOutput
   70
   71   # Note that the order by population was by a string versus number
```

```
PercentComplete : map = 100%,  reduce = 100%
Query           : CREATE EXTERNAL TABLE working_te_census_info(state STRING,county STRING,agegrp
                  STRING,total_population STRING) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' STORED AS
                  TEXTFILE LOCATION 'wasb://student01@bigdatatraintr19.blob.core.windows.net/data/hive/working_
                  te_census_info';LOAD data inpath
                  'wasb://student01@bigdatatraintr19.blob.core.windows.net/data/mr-results/part-00000'
                  OVERWRITE into table working_te_census_info;SELECT * from working_te_census_info ORDER BY
                  total_population DESC limit 10;
State           : Completed
StatusDirectory : 421ab3b7-7a42-4f1b-bfdf-3c137350bd4e
SubmissionTime  : 6/10/2014 7:53:45 PM
JobId           : job_1402420269563_0003

Maine       Kennebec County 45-49Yrs      9999
Michigan    Berrien County  5-9Yrs   9999
California   Napa County 50-54Yrs      9995
Florida Bay County  10-14Yrs      9993
California   El Dorado County     35-39Yrs      9993
Missouri    Cooper County    35-39Yrs     999
Michigan    Iosco County     30-34Yrs     999
Mississippi Holmes County    35-39Yrs     999
Minnesota   Waseca County    20-24Yrs     999
Tennessee   Crockett County 15-19Yrs      999
```

Running Hive statements as a job like this allows you to fire off a long running query and then view the results later by skipping the Wait-AzureHDInsightJob command.

# Running Hive commands using Invoke-Hive

The Invoke-Hive cmdlet is a shortcut for defining a hive job, submitting the job and then waiting for the results.

## Using Invoke-Hive

1.  To use the Invoke-Hive, you need to first set the execution context to your HDInsight cluster if you have more than one cluster under your subscription. Select the following command and press **F8** to set the context.



2.  To view the definition of a Hive table, you can use the DESCRIBE command. Select the following command and press **F8** to see the table definition for the *working_te_census_table*.



Notice the total_population is a string. This is why the sort order for the first query did not make sense initially. Scroll through the other table information to see how the table is stored.

3. HiveQL syntax includes the CAST() function for converting data types. Select the following commands and press F8 to show the correct results.



```
TR19 Hive Lab.ps1* X
79
80    # Use CAST for desired result
81    $querystring = `
82    "SELECT * from working_te_census_info " +
83    "ORDER BY CAST(total_population as DOUBLE) DESC limit 20;"
84    Invoke-Hive $querystring
85
86    # CASE statement example
87    $querystring = `
```

```
Invoke-Hive $querystring
Submitting Hive query..
Started Hive query with jobDetails Id : job_1402420269563_0006
Hive query completed Successfully


California  Los Angeles County  Above59Yrs   1517935
Illinois      Cook County Above59Yrs   877289
California  Los Angeles County   25-29Yrs    759602
California  Los Angeles County   15-19Yrs    753630
California  Los Angeles County   20-24Yrs    752788
California  Los Angeles County   30-34Yrs    716129
California  Los Angeles County   35-39Yrs    715635
California  Los Angeles County   40-44Yrs    714691
California  Los Angeles County   45-49Yrs    706742
California  Los Angeles County   10-14Yrs    678845
California  Los Angeles County   50-54Yrs    662205
Arizona Maricopa County Above59Yrs   652489
California  Los Angeles County   0-4Yrs   645793
California  Los Angeles County   5-9Yrs   633690
California  Los Angeles County   55-59Yrs    560920
Texas     Harris County   Above59Yrs   507254
California  San Diego County    Above59Yrs   500736
```

4. HiveQL syntax includes the CASE statement. Select the following commands and press F8 to see how the CASE statement can be used to provide additional grouping of age categories.

```
TR19 Hive Lab.ps1* ✕
    85
    86    # CASE statement example
    87    $querystring = `
    88    "SELECT te.state, te.county, te.agegrp, "+`
    89    "CASE te.agegrp "+`
    90    "when '0-4Yrs' then 'Infant' "+`
    91    "when '5-9Yrs' then 'Kid' "+`
    92    "when '10-14Yrs' then 'Teenager' "+`
    93    "when '15-19Yrs' then 'Teenager' "+`
    94    "when '20-24Yrs' then 'Adult' "+`
    95    "when '25-29Yrs' then 'Adult' "+`
    96    "when '30-34Yrs' then 'Adult' "+`
    97    "when '35-39Yrs' then 'Adult' "+`
    98    "when '40-44Yrs' then 'Middle-Aged Adult' "+`
    99    "when '45-49Yrs' then 'Middle-Aged Adult' "+`
   100    "when '50-54Yrs' then 'Middle-Aged Adult' "+`
   101    "when '55-59Yrs' then 'Seasoned Citizen' "+`
   102    "else 'Senior Citizen' "+`
   103    "END agegrptitle ,"+`
   104    "te.total_population "+`
   105    "FROM working_te_census_info te limit 10;"
   106    Invoke-Hive $querystring
```

```
Alabama Autauga County  20-24Yrs   Adult          3000
Alabama Autauga County  25-29Yrs   Adult          3157
Alabama Autauga County  30-34Yrs   Adult          3330
Alabama Autauga County  35-39Yrs   Adult          4157
Alabama Autauga County  40-44Yrs   Middle-Aged Adult   4086
Alabama Autauga County  45-49Yrs   Middle-Aged Adult   4332
Alabama Autauga County  5-9Yrs   Kid 3991


PS D:\TR19HiveLab>
```

5. To abstract queries, HiveQL includes the CREATE VIEW statement. To create a view based on the prior query that also includes casting the total_population column as a number, select the following commands and press F8.

```
TR19 Hive Lab.ps1* X
108    # Create a view
109    $querystring =
110    "CREATE VIEW working_agegroup_view AS "+
111    "SELECT te.state, te.county, te.agegrp, "+
112    "CASE te.agegrp "+
113    "when '0-4Yrs' then 'Infant' "+
114    "when '5-9Yrs' then 'Kid' "+
115    "when '10-14Yrs' then 'Teenager' "+
116    "when '15-19Yrs' then 'Teenager' "+
117    "when '20-24Yrs' then 'Adult' "+
118    "when '25-29Yrs' then 'Adult' "+
119    "when '30-34Yrs' then 'Adult' "+
120    "when '35-39Yrs' then 'Adult' "+
121    "when '40-44Yrs' then 'Middle-Aged Adult' "+
122    "when '45-49Yrs' then 'Middle-Aged Adult' "+
123    "when '50-54Yrs' then 'Middle-Aged Adult' "+
124    "when '55-59Yrs' then 'Seasoned Citizen' "+
125    "else 'Senior Citizen'  "+
126    "END agegrptitle ,"+
127    "cast( te.total_population as double) as total_population "+
128    "FROM working_te_census_info te"
129    Invoke-Hive $querystring
```

```
FROM working_te_census_info te
Invoke-Hive $querystring
Submitting Hive query..
Started Hive query with jobDetails Id : job_1402420269563_0010
Hive query completed Successfully


PS D:\TR19HiveLab>
```

6. Creating a view does not improve performance, but it does make repeated queries easier. To run a view, select the following commands and press F8.



```
TR19 Hive Lab.ps1* X
130
131    # Query against the view
132    $querystring =
133    "SELECT * FROM working_agegroup_view " +
134    "ORDER BY total_population DESC LIMIT 20;"
135    Invoke-Hive $querystring
136
137    # Instead of a CASE, we can do a join
138    # Upload text file with age group information
139    Set-AzureStorageBlobContent -File "D:\TR19HiveLab\AgeGroupDetails.txt"
140      -Blob "data/AgeGroupDetails.txt" -context $context
```

```
California  Los Angeles County  25-29Yrs   Adult    739002.0
California  Los Angeles County  15-19Yrs   Teenager     753630.0
California  Los Angeles County  20-24Yrs   Adult   752788.0
California  Los Angeles County  30-34Yrs   Adult   716129.0
California  Los Angeles County  35-39Yrs   Adult   715635.0
California  Los Angeles County  40-44Yrs   Middle-Aged Adult   714691.0
California  Los Angeles County  45-49Yrs   Middle-Aged Adult   706742.0
California  Los Angeles County  10-14Yrs   Teenager     678845.0
California  Los Angeles County  50-54Yrs   Middle-Aged Adult   662205.0
Arizona Maricopa County Above59Yrs  Senior Citizen  652489.0
California  Los Angeles County  0-4Yrs   Infant  645793.0
California  Los Angeles County  5-9Yrs   Kid 633690.0
California  Los Angeles County  55-59Yrs   Seasoned Citizen     560920.0
Texas    Harris County    Above59Yrs  Senior Citizen  507254.0
California  San Diego County    Above59Yrs  Senior Citizen  500736.0
California  Orange County    Above59Yrs  Senior Citizen  496404.0
Florida Miami-Dade County    Above59Yrs  Senior Citizen  476233.0
Illinois    Cook County 25-29Yrs    Adult    435510.0


PS D:\TR19HiveLab>
```

## Performing a join query

1. From a database design perspective, rather than using a CASE statement, it might be better to create a lookup table and then use the HiveQL join syntax. To upload the data for the age group lookup, select the following command and press **F8**.

```
TR19 Hive Lab.ps1* ×
136
137    # Instead of a CASE, we can do a join
138    # Upload text file with age group information
139    Set-AzureStorageBlobContent -File "D:\TR19HiveLab\AgeGroupDetails.txt" -Container $containername `
140     -Blob "data/AgeGroupDetails.txt" -context $context
141

PS D:\TR19HiveLab> # Upload text file with age group information
Set-AzureStorageBlobContent -File "D:\TR19HiveLab\AgeGroupDetails.txt" -Container $containername `
 -Blob "data/AgeGroupDetails.txt" -context $context


   Container Uri: https://bigdatatraintr19.blob.core.windows.net/student01

Name                    BlobType    Length              ContentType           LastModified           SnapshotTime
----                    --------    ------              -----------           ------------           ------------
data/AgeGroupDeta...    BlockBlob   278                 application/octe...   6/10/2014 8:27:3...
```

2. To create the lookup table and see the results in Hive, select the following commands and press F8.

```
TR19 Hive Lab.ps1* ×
142    # Create the join table
143    $querystring = `
144    "CREATE EXTERNAL TABLE working_te_agegroup_info (agegrp string, agegrp_title string)"+`
145    "ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' "+`
146    "STORED AS TEXTFILE LOCATION '$rootpart/data/hive/working_te_agegroup_info';"+`
147    "LOAD data inpath '$rootpart/data/AgeGroupDetails.txt' OVERWRITE into table working_te_agegroup_info;" +`
148    "Select * from  working_te_agegroup_info ;"
149    Invoke-Hive $querystring
150
151    $querystring = `

Invoke-Hive $querystring
Submitting Hive query..
Started Hive query with jobDetails Id : job_1402420269563_0013
Hive query completed Successfully


0-4Yrs    Infant
5-9Yrs    Kid
10-14Yrs      Teenager
15-19Yrs      Teenager
20-24Yrs      Adult
25-29Yrs      Adult
30-34Yrs      Adult
35-39Yrs      Adult
40-44Yrs      Middle-Aged Adult
45-49Yrs      Middle-Aged Adult
50-54Yrs      Middle-Aged Adult
55-59Yrs      Senior Citizen
Above59Yrs    Senior Citizen
```

3. To run a select statement with a join, select the following commands and press F8. Notice the different PowerShell string syntax you can use to avoid string concatenation errors.

```
TR19 Hive Lab.ps1* ×
150
151   $querystring = @"
152     SELECT te.state, te.county, te.agegrp, ag.agegrp_title, total_population
153     FROM working_te_census_info te JOIN working_te_agegroup_info ag
154     ON (te.agegrp = ag.agegrp) LIMIT 20;
155   "@
156   Invoke-Hive $querystring
157
158   $querystring = @"
159     SELECT te.state, te.county, te.agegrp, ag.agegrp_title, total_population
160     FROM working_te_census_info te JOIN working_te_agegroup_info ag
161     ON (te.agegrp = ag.agegrp) WHERE te.county = 'Los Angeles County' LIMIT 20;
162   "@
163   Invoke-Hive $querystring
164
165   # GROUP BY example
166   $querystring =
```

```
Alabama Autauga County   10-14Yrs      Teenager    4290
Alabama Autauga County   15-19Yrs      Teenager    4290
Alabama Autauga County   20-24Yrs      Adult   3080
Alabama Autauga County   25-29Yrs      Adult   3157
Alabama Autauga County   30-34Yrs      Adult   3330
Alabama Autauga County   35-39Yrs      Adult   4157
Alabama Autauga County   40-44Yrs      Middle-Aged Adult   4086
Alabama Autauga County   45-49Yrs      Middle-Aged Adult   4332
Alabama Autauga County   5-9Yrs    Kid    3991
Alabama Autauga County   50-54Yrs      Middle-Aged Adult   3873
Alabama Autauga County   55-59Yrs      Senior Citizen 3083
Alabama Autauga County   Above59Yrs    Senior Citizen 9323
Alabama Baldwin County   0-4Yrs     Infant 11158
Alabama Baldwin County   10-14Yrs      Teenager    11926
Alabama Baldwin County   15-19Yrs      Teenager    11600
Alabama Baldwin County   20-24Yrs      Adult   9449
```

4.  To run a query with a combination where and join clause, select
    the following commands and press **F8**.



```
TR19 Hive Lab.ps1* ×
156   Invoke-Hive $querystring
157
158   $querystring = @"
159     SELECT te.state, te.county, te.agegrp, ag.agegrp_title, total_population
160     FROM working_te_census_info te JOIN working_te_agegroup_info ag
161     ON (te.agegrp = ag.agegrp) WHERE te.county = 'Los Angeles County' LIMIT 20;
162   "@
163   Invoke-Hive $querystring
164
165   # GROUP BY example
166   $querystring = @"
167   "SELECT te.state, te.county, ag.agegrp_title, " + `
```

```
Started Hive query with jobDetails id : job_140242020903_0010
Hive query completed Successfully

California  Los Angeles County  0-4Yrs     Infant 645793
California  Los Angeles County  10-14Yrs      Teenager    678845
California  Los Angeles County  15-19Yrs      Teenager    753630
California  Los Angeles County  20-24Yrs      Adult  752788
California  Los Angeles County  25-29Yrs      Adult  759602
California  Los Angeles County  30-34Yrs      Adult  716129
California  Los Angeles County  35-39Yrs      Adult  715635
California  Los Angeles County  40-44Yrs      Middle-Aged Adult   714691
California  Los Angeles County  45-49Yrs      Middle-Aged Adult   706742
California  Los Angeles County  5-9Yrs    Kid    633690
California  Los Angeles County  50-54Yrs      Middle-Aged Adult   662205
California  Los Angeles County  55-59Yrs      Senior Citizen 560920
California  Los Angeles County  Above59Yrs    Senior Citizen 1517935
```

5.  The GROUP BY clause is similar to Transact-SQL as well. There is a
    slight change in how ORDER BY works. To run a group by query
    with a join, select the following commands and press **F8**.

```
TR19 Hive Lab.ps1* X
165    # GROUP BY example
166   ⊟$querystring = @"
167    SELECT te.state, te.county, ag.agegrp_title,
168         sum(te.total_population) as sum_pop
169    FROM working_te_census_info te
170    JOIN working_te_agegroup_info ag ON (te.agegrp = ag.agegrp)
171    GROUP BY te.state, te.county, ag.agegrp_title
172    ORDER BY sum_pop DESC LIMIT 20;
173    "@
174    # Note column number in order by not supported,
175    # but you can use an alias
176    Invoke-Hive $querystring
```

```
Illinois    Cook County    Adult   1300481.0
California  Los Angeles County    Teenager    1432475.0
Texas     Harris County    Adult   1265881.0
Illinois    Cook County    Senior Citizen 1189655.0
Arizona Maricopa County   Adult   1073967.0
Illinois    Cook County    Middle-Aged Adult   1064443.0
California  San Diego County      Adult   952684.0
Arizona Maricopa County   Senior Citizen 860950.0
California  Orange County     Adult   840734.0
Texas     Harris County    Middle-Aged Adult   830036.0
New York     Kings County     Adult   799145.0
Arizona Maricopa County   Middle-Aged Adult   762594.0
Texas     Dallas County    Adult   734758.0
Texas     Harris County    Senior Citizen 732653.0
Florida Miami-Dade County     Adult   702023.0
Illinois    Cook County  Teenager    699766.0
New York     Queens County     Adult   687360.0
```

## Running a script of HiveQL statements

HDInsight allows you to run Hive commands that are in a text file which is stored in a storage container. In this exercise, you will learn how to run HiveQL statements using the Invoke-Hive command. For larger jobs, you will want to use the New-AzureHDInsightHiveJobDefinition cmdlet with the –File parameter to schedule the job and check back later with the results.

1. There are some cases where you must run your HiveQL statements using a file. The D:\TR19HiveLab\m3case1.hql file contains the query shown for the Like query example. To upload the file to Azure storage so that you can execute it, select the following command and press **F8**.

```
179    # Like query example
180  □$querystring = @"
181    SELECT te.state, te.county, te.agegrp, ag.agegrp_title, total_population
182    FROM working_te_census_info te JOIN working_te_agegroup_info ag
183    ON (te.agegrp = ag.agegrp) WHERE te.county LIKE 'Los%' LIMIT 20;
184    "@
185    Invoke-Hive $querystring
186
187    Set-AzureStorageBlobContent -File "D:\TR19HiveLab\m3case1.hql" `
188     -Container $containername
189     -Blob "applications/m3case1.hql" -context $context -Force
190
```

```
PS D:\TR19HiveLab> Set-AzureStorageBlobContent -File "D:\TR19HiveLab\m3case1.hql" `
 -Container $containername
 -Blob "applications/m3case1.hql" -context $context -Force


   Container Uri: https://bigdatatraintr19.blob.core.windows.net/student01

Name                      BlobType     Length            ContentType          LastModified
----                      --------     ------            -----------          ------------
applications/m3ca...  BlockBlob    205               application/octe...  6/10/2014 8:53:0...
```

2. To run the file with the HiveQL statement, select the following command and press **F8**.



```
186
187    Set-AzureStorageBlobContent -File "D:\TR19HiveLab\m3case1.hql" `
188     -Container $containername
189     -Blob "applications/m3case1.hql" -context $context -Force
190
191    Invoke-Hive -File "$rootpart/applications/m3case1.hql"
192  |
193    # To test running Hive queries using Hive queries, navigate to the
194    # management console for your cluster.
195    # For example: https://<cluster_name>.azurehdinsight.net
196
197    ### The following three commands will clean up the items created for this lab.
```

```
California  Los Angeles County  20-24Yrs    Adult    732788
California  Los Angeles County  25-29Yrs    Adult    759602
California  Los Angeles County  30-34Yrs    Adult    716129
California  Los Angeles County  35-39Yrs    Adult    715635
California  Los Angeles County  40-44Yrs    Middle-Aged Adult   714691
California  Los Angeles County  45-49Yrs    Middle-Aged Adult   706742
California  Los Angeles County  5-9Yrs    Kid     633690
California  Los Angeles County  50-54Yrs    Middle-Aged Adult   662205
California  Los Angeles County  55-59Yrs    Senior Citizen 560920
California  Los Angeles County  Above59Yrs    Senior Citizen 1517935
New Mexico  Los Alamos County  0-4Yrs    Infant 960
New Mexico  Los Alamos County  10-14Yrs    Teenager    1307
New Mexico  Los Alamos County  15-19Yrs    Teenager    1126
New Mexico  Los Alamos County  20-24Yrs    Adult    498
New Mexico  Los Alamos County  25-29Yrs    Adult    729
New Mexico  Los Alamos County  30-34Yrs    Adult    984
New Mexico  Los Alamos County  35-39Yrs    Adult    1114
```

# Using the Hive console with HDInsight

Starting in June 2014, the HDInsight team brought back the ability of running Hive queries through the Manage Cluster web interface. The console is a useful way of testing out Hive queries before automating execution of them using PowerShell. In this exercise, you will run a query using the console.

1. Navigate to https://<your_cluster>.azurehdinsight.net.

2. Enter in your user name Admin and password Pass@word12 in the credentials dialog.

3. Select the text for the Like query in the TR19 Hive Lab.ps1 shown below and copy the text to the clipboard.

```
TR19 Hive Lab.ps1*  ×
177
178
179     # Like query example
180   ⊟$querystring = @"
181     SELECT te.state, te.county, te.agegrp, ag.agegrp_title, total_population
182     FROM working_te_census_info te JOIN working_te_agegroup_info ag
183     ON (te.agegrp = ag.agegrp) WHERE te.county LIKE 'Los%' LIMIT 20;
184     "@
185     Invoke-Hive $querystring
186
187     Set-AzureStorageBlobContent -File "D:\TR19HiveLab\m3case1.hql"
```

| Cut   | Ctrl+X |
| Copy  | Ctrl+C |
| Paste | Ctrl+V |

4.  Go back to the browser and enter in a **Query Name** such as **HiveLike** and then paste the HiveQL statement into the editor region.



5.  Click **Submit** to run the query. HDInsight displays the progress of the query in the **Status** column.

6.  Once the **Status** value shows **Succeeded**, click on the Title link for the query to view the results.

7. HDInsight displays a status page as shown below that lists the Job Query, Job Output and Job Log.



## Summary

In this set of exercises, you learned the different ways of running HiveQL statements and the similarities of HiveQL to Transact-SQL. For a complete HiveQL language reference, go to https://cwiki.apache.org/confluence/display/Hive/LanguageManual.

# Terms of Use

INCLUDING ALL WARRANTIES AND CONDITIONS OF MERCHANTABILITY, WHETHER EXPRESS, IMPLIED OR STATUTORY, FITNESS FOR A PARTICULAR PURPOSE, TITLE AND NON-INFRINGEMENT.  MICROSOFT DOES NOT MAKE ANY ASSURANCES OR REPRESENTATIONS WITH REGARD TO THE ACCURACY OF THE RESULTS, OUTPUT THAT DERIVES FROM USE OF THE VIRTUAL LAB, OR SUITABILITY OF THE INFORMATION CONTAINED IN THE VIRTUAL LAB FOR ANY PURPOSE.

**DISCLAIMER**

This lab contains only a portion of new features and enhancements in Microsoft SQL Server 2014. The exercises are built on the General Availability release as of October 25, 2013 and updated with the April 2014 release that includes HDInsight version 3.0. Some of the features might change in future releases of the product. In this lab, you will learn about some, but not all, new features.