

Processing Big Data with Hadoop in Azure HDInsight

Lab 4B – Using Oozie

Overview

In this lab, you will use Oozie to define a workflow of data processing operations that summarize data in Internet Information Services (IIS) web server log files and load the results to a table in Azure SQL Database for further analysis.

What You'll Need

To complete the labs, you will need the following:

- A web browser
- A Microsoft account
- A Microsoft Azure subscription
- A Microsoft Windows computer with the following software installed:
 - Microsoft Azure PowerShell
 - Microsoft Power BI Desktop
- The lab files for this course

Note: To set up the required environment for the lab, follow the instructions in the **Setup** document for this course. Specifically, you must have signed up for an Azure subscription, installed and configured Azure PowerShell, imported the publisher settings for your Azure subscription into PowerShell, and installed Microsoft Power BI Desktop.

When working with cloud services, transient network errors can occasionally cause scripts to fail. If a script fails, and you believe that you have entered all of the required variables correctly; wait a few minutes and run the script again.

Provisioning HDInsight and Azure SQL Database

In this lab, you will build an Oozie workflow that uses Azure HDInsight to process IIS web server log data, and then transfers the processed data to Azure SQL Database. Before running the workflow, you need to provision the required Azure services.

Provision an Azure Storage Account and HDInsight Cluster

Note: If you already have an HDInsight cluster and associated storage account, you can skip this task.

1. In a web browser, navigate to <http://azure.microsoft.com>. Then click **Portal**, and if prompted, sign in using the Microsoft account that is associated with your Azure subscription.
2. In the Microsoft Azure portal, view the **HDInsight** page and verify that there are no existing HDInsight clusters in your subscription.
3. Click **NEW** (at the bottom of the page) and then click **CUSTOM CREATE**. Then use the New HDInsight Cluster wizard to create a new cluster with the following settings. Click the arrows to navigate through all of the wizard pages:
 - **Cluster Name:** *Enter a unique name (and make a note of it!)*
 - **Cluster Type:** Hadoop
 - **Operating System:** Windows Server 2012 R2 Datacenter
 - **HDInsight Version:** 3.2 (HDP 2.2, Hadoop 2.6)
 - **Data Nodes:** 2
 - **Region:** *Select any available region*
 - **Head Node Size:** A3 (4 cores, 7 GB memory)
 - **Data Node Size:** A3 (4 cores, 7 GB memory)
 - **HTTP User Name:** *Enter a user name of your choice (and make a note of it!)*
 - **HTTP Password:** *Enter and confirm a strong password (and make a note of it!)*
 - **Enable the remote desktop for cluster:** Selected
 - **RDP User Name:** *Enter another user name of your choice (and make a note of it!)*
 - **RDP Password:** *Enter and confirm a strong password (and make a note of it!)*
 - **Expires on:** Select tomorrow's date
 - **Enter the Hive/Oozie Metastore:** Unselected
 - **Storage Account:** Create New Storage
 - **Account Name:** *Enter a unique name for your storage account (and make a note of it!)*
 - **Default Container:** *Enter a unique name for your container (and make a note of it!)*
 - **Additional Storage Accounts:** 0
 - **Additional scripts to customize the cluster:** None
4. Wait for the cluster to be provisioned and the status to change to **Running** (this can take a while.)

Provision Azure SQL Database

Note: If you already have an Azure SQL Database named **AnalysisDB**, you can skip this task.

1. In a web browser, navigate to <http://azure.microsoft.com> and click **Portal**, and if prompted, sign in using the Microsoft account that is associated with your Azure subscription.
2. View the **HDInsight** page and note the **Location** of your HDInsight cluster. Then, in the Azure portal, view the **SQL Databases** page.
3. Click **New**, and click **Custom Create**. Then in the **Name** box, type **AnalysisDB** and in the **Server** list, select **New SQL database server**. Then click **Next**.
4. On the **Create Server** page, enter a login name and password of your choice, make a note of the login name and password you have specified, and select the region in which your HDInsight cluster is defined. Then ensure that the **Allow Windows Azure Services to Access the Server** and **Enable Latest SQL Database Updates (V12)** are both selected, and click **Complete**.
5. Wait a minute or so for the database status to become online (you may need to refresh the web page in the browser), and then note the name in the **Server** column – this is the Azure SQL Database server that has been created to host your database.
6. Click the server name, and then click **Configure**.
7. Under **allowed ip addresses**, next to your current IP address, click **Add to the Allowed IP Addresses**. Then at the bottom of the page, click **Save**.

Running an Oozie Workflow

Now that you have the required infrastructure, you are ready to initiate an Oozie workflow to process the log data.

Examine the Oozie Workflow

1. In the C:\HDILabs\Lab04B\oozieworkflow folder, open **workflow.xml** in Visual Studio or Notepad. This XML document defines an Oozie workflow that consists of multiple actions.
2. Note the following element, which directs the workflow to start at the **PrepareHive** action:

```
<start to="PrepareHive"/>
```

3. Review the following element, which defines the **PrepareHive** action:

```
<action name='PrepareHive'>
  <hive xmlns="uri:oozie:hive-action:0.2">
    <job-tracker>${jobTracker}</job-tracker>
    <name-node>${nameNode}</name-node>
    <configuration>
      <property>
        <name>mapred.job.queue.name</name>
        <value>default</value>
      </property>
    </configuration>
    <script>${DropTableScript}</script>
    <param>CLEANSED_TABLE_NAME=${cleansedTable}</param>
    <param>SUMMARY_TABLE_NAME=${summaryTable}</param>
  </hive>
  <ok to="CleanseData"/>
  <error to="fail"/>
</action>
```

The **PrepareHive** action is a Hive action that runs the HiveQL script specified in a workflow configuration setting named **DropTableScript**. Two parameters named **CLEANSED_TABLE_NAME** and **SUMMARY_TABLE_NAME** are passed to the script based on the values in the **cleansedTable** and **summaryTable** workflow configuration settings. If the action succeeds, the workflow moves on to the **CleanseData** action.

4. Keep workflow.xml open, and in the C:\HDILabs\Lab04B\oozieworkflow folder, open **DropHiveTables.txt** in Notepad. This is the script that will be specified in the **DropTableScript** workflow configuration setting. Note that it contains HiveQL statements to drop the two tables specified by the parameters if they exist. Then close DropHiveTables.txt without saving any changes.
5. In workflow.xml, review the following element, which defines the **CleanseData** action:

```
<action name="CleanseData">
  <pig>
    <job-tracker>${jobTracker}</job-tracker>
    <name-node>${nameNode}</name-node>
    <script>${CleanseDataScript}</script>
    <param>StagingFolder=${stagingFolder}</param>
    <param>CleansedFolder=${cleansedFolder}</param>
  </pig>
  <ok to="CreateTables"/>
  <error to="fail"/>
</action>
```

The **CleanseData** action is a Pig action that runs the Pig Latin script specified in a workflow configuration setting named **CleanseDataScript**. Two parameters named **StagingFolder** and **CleansedFolder** are passed to the script based on the values in the **stagingFolder** and **cleansedFolder** workflow configuration settings. If the action succeeds, the workflow moves on to the **CreateTables** action.

6. Keep workflow.xml open, and in the C:\HDILabs\Lab04B\oozieworkflow folder, open **CleanseData.txt** in Notepad. This is the script that will be specified in the **CleanseDataScript** workflow configuration setting. Note that it contains Pig Latin statements to load the data in the path specified by the **stagingFolder** parameter as a single character array for each row, filter it to remove rows that begin with a "#" character, and store the results in the path defined by the **cleansedFolder** parameter. Then close CleanseData.txt without saving any changes.
7. In workflow.xml, review the following element, which defines the **CreateTables** action:

```
<action name='CreateTables'>
  <hive xmlns="uri:oozie:hive-action:0.2">
    <job-tracker>${jobTracker}</job-tracker>
    <name-node>${nameNode}</name-node>
    <configuration>
      <property>
        <name>mapred.job.queue.name</name>
        <value>default</value>
      </property>
    </configuration>
    <script>${CreateTableScript}</script>
    <param>CLEANSED_TABLE_NAME=${cleansedTable}</param>
    <param>CLEANSED_TABLE_LOCATION=${cleansedFolder}</param>
    <param>SUMMARY_TABLE_NAME=${summaryTable}</param>
    <param>SUMMARY_TABLE_LOCATION=${summaryFolder}</param>
  </hive>
  <ok to="SummarizeData"/>
  <error to="fail"/>
</action>
```

The **CreateTables** action is a Hive action that runs the HiveQL script specified in a workflow configuration setting named **CreateTableScript**. Four parameters named **CLEANSED_TABLE_NAME**, **CLEANSED_TABLE_LOCATION**, **SUMMARY_TABLE_NAME**, and **SUMMARY_TABLE_LOCATION** are passed to the script based on the values in the **cleansedTable**, **cleansedFolder**, **summaryTable**, and **summaryFolder** workflow configuration settings. If the action succeeds, the workflow moves on to the **SummarizeData** action.

8. Keep workflow.xml open, and in the C:\HDILabs\Lab04B\oozieworkflow folder, open **CreateHiveTables.txt** in Notepad. This is the script that will be specified in the **CreateTableScript** workflow configuration setting. Note that it contains HiveQL statements to create Hive tables based on the table name and location parameters. Then close CreateHiveTables.txt without saving any changes.
9. In workflow.xml, review the following element, which defines the **SummarizeData** action:

```
<action name='SummarizeData'>
  <hive xmlns="uri:oozie:hive-action:0.2">
    <job-tracker>${jobTracker}</job-tracker>
    <name-node>${nameNode}</name-node>
    <configuration>
      <property>
        <name>mapred.job.queue.name</name>
        <value>default</value>
      </property>
    </configuration>
  </hive>
  <ok to="finish"/>
  <error to="fail"/>
</action>
```

```

        </property>
    </configuration>
    <script>${SummarizeDataScript}</script>
    <param>CLEANSED_TABLE_NAME=${cleansedTable}</param>
    <param>SUMMARY_TABLE_NAME=${summaryTable}</param>
</hive>
<ok to="TransferData"/>
<error to="fail"/>
</action>

```

The **SummarizeData** action is a Hive action that runs the HiveQL script specified in a workflow configuration setting named **SummarizeDataScript**. Two parameters named **CLEANSED_TABLE_NAME** and **SUMMARY_TABLE_NAME** are passed to the script based on the values in the **cleansedTable** and **summaryTable** workflow configuration settings. If the action succeeds, the workflow moves on to the **TransferData** action.

10. Keep workflow.xml open, and in the C:\HDI\Lab04B\oozieworkflow folder, open **SummarizeData.txt** in Notepad. This is the script that will be specified in the **SummarizeDataScript** workflow configuration setting. Note that it contains HiveQL statements to insert data into the table specified by the **SUMMARY_TABLE_NAME** parameter based on the results of a query from the table specified by the **CLEANSED_TABLE_NAME** parameter. Then close SummarizeData.txt without saving any changes.
11. In workflow.xml, review the following element, which defines the **TransferData** action:

```

<action name="TransferData">
  <sqoop xmlns="uri:oozie:sqoop-action:0.2">
    <job-tracker>${jobTracker}</job-tracker>
    <name-node>${nameNode}</name-node>
    <configuration>
      <property>
        <name>mapred.compress.map.output</name>
        <value>true</value>
      </property>
    </configuration>
    <arg>export</arg>
    <arg>--connect</arg>
    <arg>${sqlConnectionString}</arg>
    <arg>--table</arg>
    <arg>${sqlTable}</arg>
    <arg>--export-dir</arg>
    <arg>${sourceDir}</arg>
    <arg>--input-fields-terminated-by</arg>
    <arg>\t</arg>
  </sqoop>
  <ok to="end"/>
  <error to="fail"/>
</action>

```

The **TransferData** action is a Sqoop action that exports data from the storage location specified by the **sourceDir** workflow configuration setting to a table in a database specified by the **sqlTable** and **sqlConnection** workflow configuration settings. If the action succeeds, the workflow moves on to the **end** terminator.

12. In workflow.xml, review the **kill** and **end** elements. These elements are terminators for the workflow. If all actions succeed, the workflow ends with the end terminator. If an action fails, the workflow is redirected to the fail terminator and the most recent error is reported.
13. Close workflow.xml without saving any changes.

Run the Oozie Workflow

1. In the C:\HDILabs\Lab04B folder, rename **Run Oozie Workflow.txt** to **Run Oozie Workflow.ps1** (you may need to modify the *View* options for the folder to see file extensions). Then right-click **Run Oozie Workflow.ps1** and click **Edit** to open the script in the Windows PowerShell interactive script environment (ISE).
2. Change the values assigned to the **\$clusterName** **\$storageAccountName**, **\$containerName**, **\$hdUser** and **\$hdPassword** variables to match the settings of your HDInsight cluster.
3. Change the values assigned to the following variables based on the configuration of your Azure SQL Database server:
 - **\$sqlServer**: Your Azure SQL server name
 - **\$sqlLogin**: Your SQL login name
 - **\$sqlPassword**: Your SQL Login Password (note that you must prefix reserved PowerShell characters such as **\$** and **#** with a **`** character, so the value *SecurePa`\$`\$wOrd* assigns the password *SecurePa\$\$wOrd*).
4. Save the script and then review the rest of the code, noting that it performs the following actions:
 - a. Prepares the Azure SQL Database by dropping (if necessary) and creating a table named **logdata**, to which the summarized log data generated by the workflow will be transferred.
 - b. Uploads the files in the local **iislogs_gz** subfolder to **/data/iislogs/stagedlogs** in your Azure storage container.
 - c. Uploads the files in the local **oozieworkflow** subfolder to **/data/iislogs/oozieworkflow** in your Azure storage container.
 - d. Defines the Oozie workflow configuration settings as an XML structure, specifying the values that will be passed to the workflow.
 - e. Initiates an Oozie job using the HTTP REST interface.
 - f. Waits for the job to complete, displaying the job status every 30 seconds.
5. On the toolbar, click **Run Script**, and wait for the script to finish (which can take several minutes), observing the status information displayed in the console window.
6. When the script has completed, close Windows PowerShell ISE.

Access the Processed Data in Azure SQL Database

1. Start Power BI Desktop and close the welcome page if it opens.

Note: Power BI Desktop is the released version of the Power BI Designer preview tool used in the demonstrations for this course. The tool has been renamed and updated, and looks cosmetically different from the preview version; but still provides the same functionality as shown in the demonstrations.
2. On the **Home** tab, in the **Get Data** list, click **SQL Server**.
3. In the **Microsoft SQL Database** page, in the **Server** box, type the fully-qualified name of your Azure SQL Database server in the format *your_azure_sql_server_name.database.windows.net*, and in the Database box, type **AnalysisDB**. Then click **OK**.
4. If the **Access a Microsoft SQL Database** dialog box is displayed, view the **Database** page; and in the **Username** box, type your Azure SQL Database login name, in the **Password** box type your Azure SQL Database login password, and click **Connect**.
5. In the **Navigators** window, select the **logdata** table. Then click **Load**.
6. In the **Visualizations** pane, select **Line and Stacked Columns Chart**.
7. In the **Fields** pane, expand the **logdata** table, and select **InboundBytes** and **OutboundBytes**. Then drag **log_date** to the **Shared Axis** area in the **Visualizations** pane, and drag **Requests** to the **Line Values** area in the **Visualizations** pane.

8. View the chart, which shows the summarized web server log data that was generated and transferred to Azure SQL Database by the Oozie workflow.
9. Close Power BI Desktop without saving the report.

Cleaning Up

Now that you have finished this lab, you can delete the HDInsight cluster and storage account.

Delete the HDInsight Cluster and Azure SQL Database Server

If you no longer need the HDInsight cluster and Azure SQL Database server used in this lab, you should delete them to avoid incurring unnecessary costs (or using credits in a free trial subscription).

1. In the Azure portal, click the **HDInsight** page.
2. Select the row containing your HDInsight cluster, and then at the bottom of the page, click **Delete**. When prompted to confirm the deletion, click **Yes**.
3. Wait for your cluster to be deleted, and then click the **Storage** page, and if necessary refresh the browser to view the storage account that was created with your cluster.
4. Select the row containing the storage account, and then at the bottom of the page, click **Delete**. When prompted to confirm the deletion, enter the storage account name and click **OK**.
5. In the Azure portal, click the **SQL Databases** page.
6. Click the **Servers** tab.
7. Select the row containing your SQL Database server, and then at the bottom of the page, click **Delete**. When prompted to confirm the deletion, enter the server name and click **OK**.
8. Close the browser.