



# Exploring Hive & Spark using Apache Zeppelin Notebooks in Microsoft Azure HDInsight



## Introduction

This class introduces students to Apache Hive and Spark using Apache Zeppelin Notebooks on Azure HDInsight. It helps student to understand the value proposition of Apache Hive and Spark over other Big Data technologies like Hadoop. They should understand the similarities between Hadoop & Spark, their differences and respective nuances. They should be able to decide when to use what and why for a given business use case in a typical enterprise environment.

## Prerequisites

- a) An Azure subscription. See [Get Azure free trial](#).
- b) Download and install Azure Storage Explorer.

# Cluster Credentials:

Generic information:

*Resource Group: MTLHDiLab*

*Storage Account : mtllab*

*Location : East US*

Spark Cluster:

*Cluster Name for Spark Cluster: mtlhdilabsspark<1-12>*

*Cluster URL (Ambari) for Hive Cluster: https://mtlhdilabsspark<1-12>.azurehdinsight.net/*

*Username: admin*

*Password: HDItut@123*

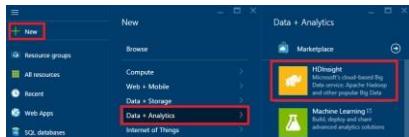
## Class Section 1: Provision an HDInsight Spark cluster

### Access Azure Preview Portal

1. Sign in to the Azure preview portal.

### Create HDInsight Spark cluster

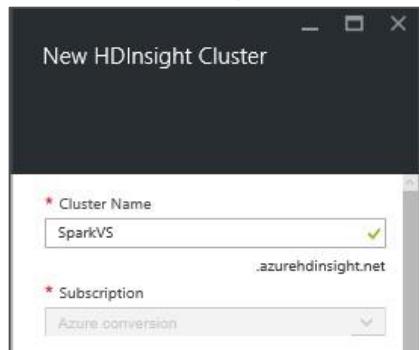
1. Click **NEW**, click **Data + Analytics**, and then click **HDInsight**.



### Provide Cluster Details

1. In the New HDInsight Cluster blade, select the “Custom (size, settings, apps)” option to have more detailed control over the cluster creation.

In the **Basic** tab, enter **Cluster Name**.



A green check mark appears beside the cluster name if it is available.

2. For **Subscription**, select **Azure Conversion**. If you have more than one subscription, click the Subscription entry to select the Azure subscription to use for the cluster.

## Configure Cluster Type

1. Click on **Cluster type**. This will open the Cluster Type configuration blade.

The image shows two overlapping windows from the Azure portal. On the left is the 'Basics' blade, which contains fields for Cluster name (mtllab00), Subscription (Microsoft Azure Internal Consumption (a7C)), Cluster type (Spark, highlighted with a red border), Cluster login username (admin), Cluster login password (redacted), Secure Shell (SSH) username (sshuser), Resource group (Create new, MTLHDiLab selected), and Location (Canada East). On the right is the 'Cluster configuration' blade, which shows the selected Cluster type (Spark), Operating system (Linux), Version (Spark 2.1.0 (HDI 3.6)), Cluster tier (STANDARD), and a detailed list of available and not-available features for Spark.

2. Select **Spark** as **Cluster Type**.
3. **Operating System** will be **Linux** by default.
4. Select **Version** as **Spark 2.1.0 (HDI 3.6)**
5. For **Cluster Tier**, select **STANDARD**
6. Click **SELECT** to complete the configuration settings

## Provide credentials to access cluster

1. Back in the Basic blade, specify the **Credentials**.
  - a. Enter **Cluster Login Password**.
  - b. Enter **SSH Username**.
  - c. Select **PASSWORD** as **SSH Authentication Type**.
  - d. Enter **SSH Password** and confirm it.

\*SSH Username and Password is required to remote session of Spark Cluster

## Provide Resource Group

1. In the **Resource Group** section, you have options:

- a) Click link **Create New** to provide new Resource Group.

Or

- b) Use Existing by selecting appropriate Resource Group from list.

For the purpose of this lab, we will create a new resource group.

2. Enter the name “**MTLHDILab**” for Resource Group

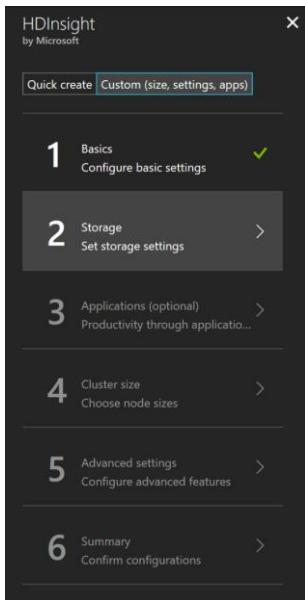
A green check appears beside the cluster name if this name is available.

If an error appears stating that the group already exist, change option to “Use existing and pick the “**MTLHDILab**”

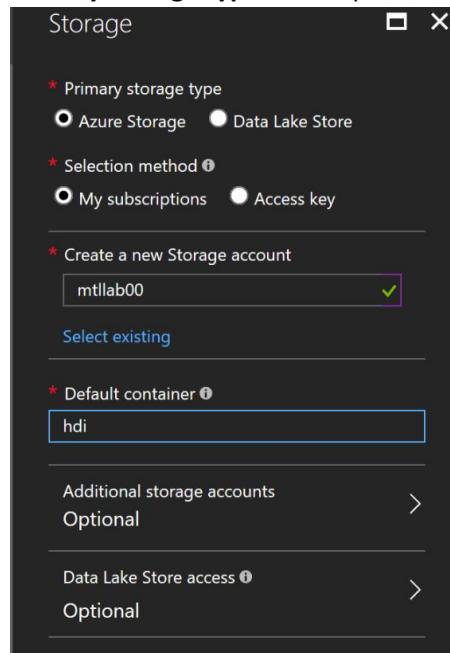
## Provide Data Source

Data source will be used as primary location for most data access, such as job input and log output.

1. Click **Storage tab** present on New HDInsight Cluster blade



2. In the Storage blade, you will need to specify the **primary storage type** and the **Selection Method**. The **Primary storage type** has 2 options:



Option 1 – **Azure Storage** if you want to use blob storage from an Azure storage resource

Option 2 – **Data Lake Store** if you want to use the Azure Data Lake Store resource.

Select **Azure Storage** for purpose of this Lab exercise

As for **Selection Method**, you have 2 options:

Option 1 - Set this to **Access Key** if you want to use existing storage account and you have **Storage Name** and **Access Key** of same, else

Option 2 - select **From all subscriptions** as Selection method.

Select **From all subscriptions** for purpose of this Lab exercise

3. To create new storage account enter a name for new storage account in **Create New Storage Account** input box

**Or**

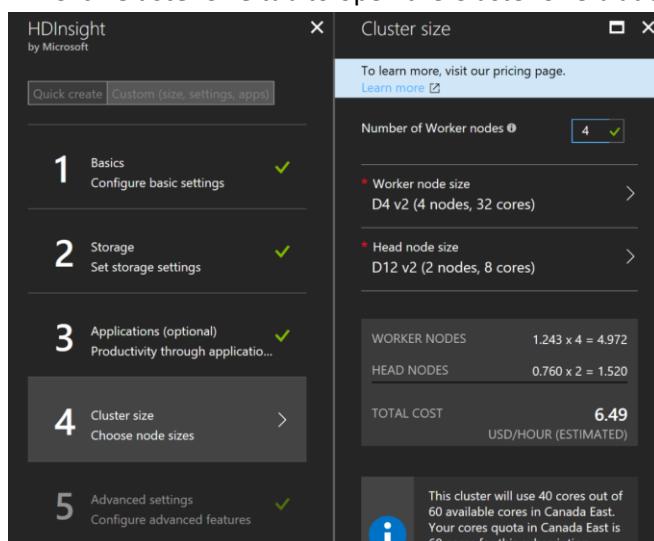
Click on link **Select Existing** to select from existing accounts.

For purpose of this exercise, we will create a new storage account.

4. Enter the name “**hdi**” for default container to be designated for cluster in **Choose Default Container** field.
5. If existing storage account is selected then no need to provide **Location** else select appropriate **Location**. By default, the HDInsight cluster is provisioned in the same data center as the storage account you specify
6. Click **Select** button at the bottom to save the data source configuration.

## Set Pricing

1. Click **Cluster Size** tab to open the Cluster Size blade.



2. The Cluster Size blade provides options to the configure number of nodes in cluster, which will be the base pricing criteria. Enter number of worker node in **Number of Worker nodes** field, set it to **4** for this demo.
3. Leave all other values as default.  
Note that based on the number of worker notes and size, the estimated cost of the cluster is calculated and displayed in USD/HOUR.
4. Click **Select** button to save node pricing configuration.

## Provision cluster

1. After completing all the configuration, in the New HDInsight Cluster blade, make sure to tick on the 'Pin to dashboard' option.
2. Click **Create** button to finalize cluster creation.

This creates the cluster and adds a tile for it to the **Dashboard** of your Azure portal.

The icon will indicate that the cluster is provisioning, and will change to display the HDInsight icon once provisioning has completed.

## 1. DOWNLOAD SAMPLE DATA

Download the sample data for this tutorial:

<https://raw.githubusercontent.com/hortonworks/data-tutorials/master/tutorials/hdp/loadingand-querying-data-with-hadoop/assets/retail-store-logs-sample-data.zip>

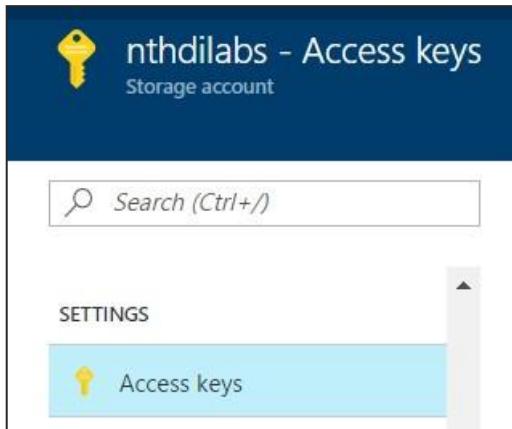
Extract the archive anywhere you'd like – we will upload the contents into our sandbox in the following step. You should see the following files after the archive is extracted:

## 2. LOAD DATASETS FILES TO STORAGE ACCOUNT

In this section, you'll copy the files required for the lab to your storage account. You'll copy the files between two storage account with the help of AzCopy utility. You can download the utility from here <http://aka.ms/downloadazcopy>

To copy the files, follow the below steps.

1. Copy your Azure Storage account access keys. This is required to copy data from the source Azure Storage account to your Azure Storage account. To get your storage account access key, navigate to your storage account on the Azure Management Portal and select **Access keys** under **Settings**.



2. Click on the copy icon to copy **Key1** from the **Access Keys** pane.

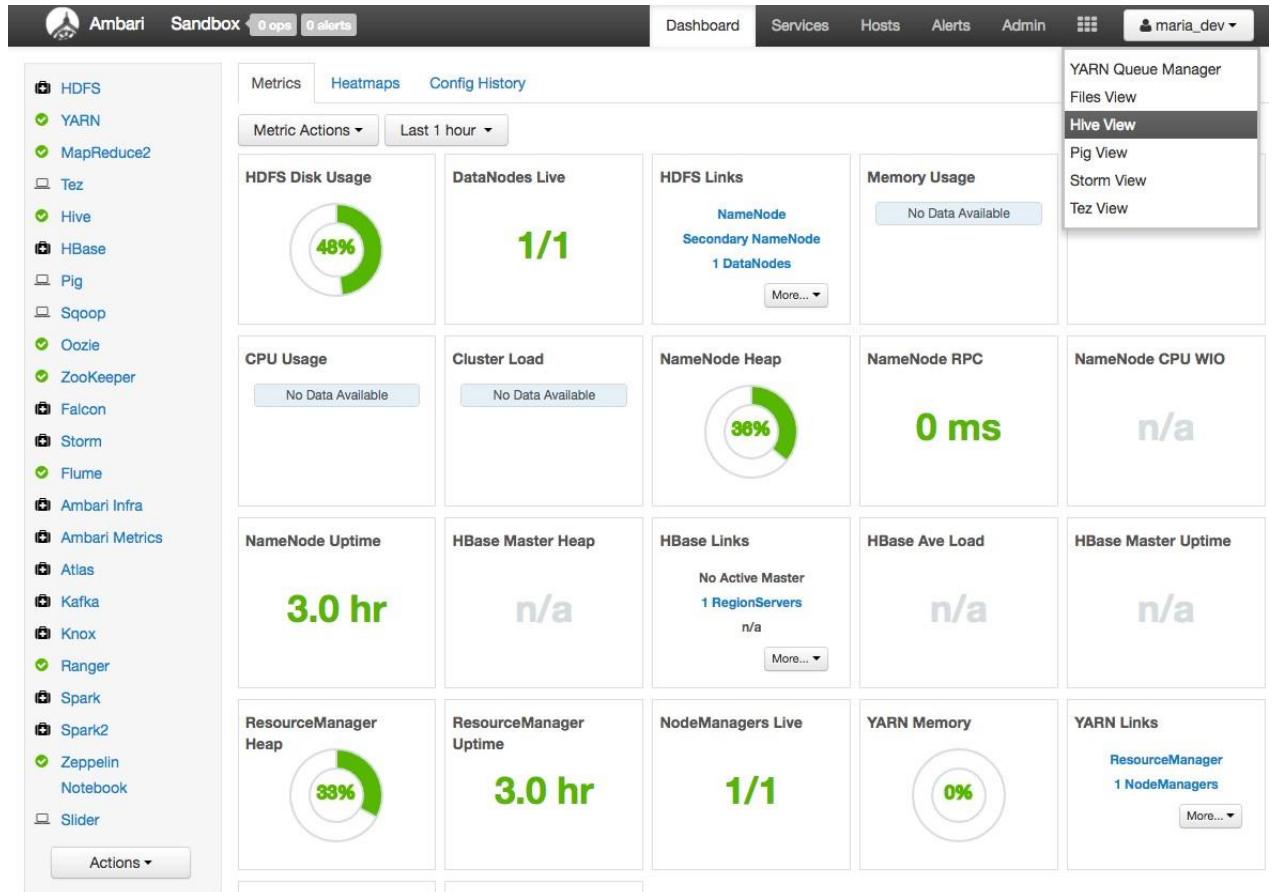
NAME	KEY	
key1	onLLmMI4n9zhQ0OgOzwyabeneOJbEcIVih/nTX+jzh8bpgYpTC	 ...
key2	UDEidwroMarY3lyisw6zKr7s2UQ/yMD1vaKi1D/4cp7EHgKwgin	 ...

3. Press Window + R to open the run window. Type cmd and press enter to open a new command console window.
4. Change the directory to C:\Program Files (x86)\Microsoft SDKs\Azure\AzCopy.
5. Copy and paste the following command on the console window to transfer **all spark lab assets needed** from the source storage account to your storage account.

```
AzCopy
/Source:https://mtlworkshop.blob.core.windows.net/hdi/HiveSpark/
/Dest:https://mtllab<1-12>.blob.core.windows.net/hdi/Data/
/SourceKey:s2eYqDZOhB00v0glTcneXYC7t5j1d58rP28BddfdP4Mv4/hI+pArEUAFj
VgWORDyUC1Kxxro0o144vmm8QfKRw== /DestKey:<KEY1> /S
```

## 4. CREATE HIVE TABLES

Let's create some Hive tables for our data. Open the **Hive View** from the Views menu on the navigation bar.



Once there, create three tables by copy/pasting each query (execute the queries one after the other) below into the query editor and clicking “**Execute**”.

```

1 CREATE TABLE users (swid STRING, birth_dt STRING, gender_cd CHAR(1))
2 ROW FORMAT DELIMITED
3 FIELDS TERMINATED BY '\t'
4 STORED AS TEXTFILE
5 TBLPROPERTIES ("skip.header.line.count"="1");
6
7 CREATE TABLE products (url STRING, category STRING)
8 ROW FORMAT DELIMITED
9 FIELDS TERMINATED BY '\t'
10 STORED AS TEXTFILE
11 TBLPROPERTIES ("skip.header.line.count"="1");
12
13 CREATE TABLE omniturelogs (col_1 STRING,col_2 STRING,col_3 STRING,col_4 STRING,col_5 STRING,col_6 STRING,col_7 STRING,col_8 STRING,col_9 STRING,col_10 STRING,col_11 STRING,col_12 STRING,col_13 STRING,col_14 STRING,col_15 STRING,col_16 STRING,col_17 STRING,col_18 STRING,col_19 STRING,col_20 STRING,col_21 STRING,col_22 STRING,col_23 STRING,col_24 STRING,col_25 STRING,col_26 STRING,col_27 STRING,col_28 STRING,col_29 STRING,col_30 STRING,col_31 STRING,col_32 STRING,col_33 STRING,col_34 STRING,col_35 STRING,col_36 STRING,col_37 STRING,col_38 STRING,col_39 STRING,col_40 STRING,col_41 STRING,col_42 STRING,col_43 STRING,col_44 STRING,col_45 STRING,col_46 STRING,col_47 STRING,col_48 STRING,col_49 STRING,col_50 STRING,col_51 STRING,col_52 STRING,col_53 STRING)
14 ROW FORMAT DELIMITED
15 FIELDS TERMINATED by '\t'
16 STORED AS TEXTFILE
17 TBLPROPERTIES ("skip.header.line.count"="1");

```

Create the tables users, products and omniturelogs:

```

CREATE TABLE users (swid STRING, birth_dt STRING, gender_cd CHAR(1))
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '\t'
STORED AS TEXTFILE
TBLPROPERTIES ("skip.header.line.count"="1");

```

```

CREATE TABLE products (url STRING, category STRING)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '\t'
STORED AS TEXTFILE
TBLPROPERTIES ("skip.header.line.count"="1");

```

```

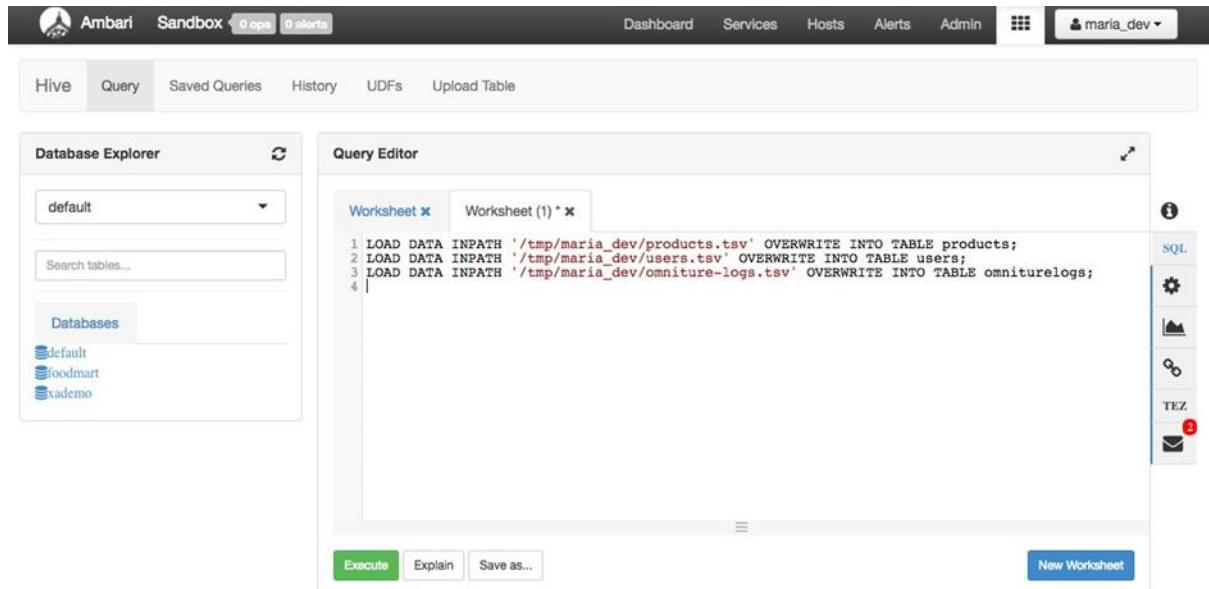
CREATE TABLE omniturelogs (col_1 STRING,col_2 STRING,col_3 STRING,col_4 STRING,col_5 STRING,col_6 STRING,col_7 STRING,col_8 STRING,col_9 STRING,col_10 STRING,col_11 STRING,col_12 STRING,col_13 STRING,col_14 STRING,col_15 STRING,col_16 STRING,col_17 STRING,col_18 STRING,col_19 STRING,col_20 STRING,col_21 STRING,col_22 STRING,col_23 STRING,col_24 STRING,col_25 STRING,col_26 STRING,col_27 STRING,col_28 STRING,col_29 STRING,col_30 STRING,col_31 STRING,col_32 STRING,col_33 STRING,col_34 STRING,col_35 STRING,col_36 STRING,col_37 STRING,col_38 STRING,col_39 STRING,col_40 STRING,col_41 STRING,col_42 STRING,col_43 STRING,col_44 STRING,col_45 STRING,col_46 STRING,col_47 STRING,col_48 STRING,col_49 STRING,col_50 STRING,col_51 STRING,col_52 STRING,col_53 STRING)
ROW FORMAT DELIMITED
FIELDS TERMINATED by '\t'
STORED AS TEXTFILE
TBLPROPERTIES ("skip.header.line.count"="1");

```

## 5. LOAD DATA INTO TABLES

Let's run a simple query to take the data we stored in HDFS and populate our new Hive tables. To do so, execute the following query:

```
LOAD DATA INPATH '/tmp/maria_dev/products.tsv' OVERWRITE INTO TABLE products;
LOAD DATA INPATH '/tmp/maria_dev/users.tsv' OVERWRITE INTO TABLE users;
LOAD DATA INPATH '/tmp/maria_dev/omniture-logs.tsv' OVERWRITE INTO TABLE omniturelogs;
```



The screenshot shows the Ambari Hive Query Editor interface. At the top, there is a navigation bar with links for Ambari, Sandbox, Dashboard, Services, Hosts, Alerts, Admin, and a dropdown for 'maria\_dev'. Below the navigation bar, there are tabs for Hive, Query, Saved Queries, History, UDFs, and Upload Table. The main area is divided into two panes: 'Database Explorer' on the left and 'Query Editor' on the right. In the Database Explorer pane, the 'default' database is selected, and there is a search bar labeled 'Search tables...'. Under 'Databases', three databases are listed: default, foodmart, and vademo. In the Query Editor pane, there is a 'Worksheet' tab and a 'Worksheet (1)' tab. The code from the previous block is pasted into the Worksheet (1) tab. On the right side of the Query Editor pane, there is a sidebar with icons for SQL, Tez, and a message center with a red notification badge showing '2'. At the bottom of the Query Editor pane, there are buttons for 'Execute', 'Explain', 'Save as...', and 'New Worksheet'.

You can verify that the data was loaded properly by browsing information about the different Hive tables in the “**Database Explorer**” tab. Click on the default database. The newly created tables will appear in addition to pre-existing tables on the sandbox. Click on the right-hand side icon to get a sample of 100 rows of the table “**Users**” for example.

The screenshot shows the Ambari Sandbox interface. At the top, there are tabs for Dashboard, Services, Hosts, Alerts, Admin, and a user dropdown for 'maria\_dev'. Below the tabs, there are sub-tabs: Hive, Query (which is selected), Saved Queries, History, UDFs, and Upload Table. The main area has two panes: 'Database Explorer' on the left and 'Query Editor' on the right. In the Database Explorer pane, the 'default' database is selected, and a search bar shows 'Search tables...'. Below it, a 'Databases' section lists several databases: default, omniturelogs, products, sample\_07, sample\_08, users, foodmart, and xademo. In the Query Editor pane, a query titled 'users sample' is displayed: 'SELECT \* FROM users LIMIT 100;'. Below the query, there are buttons for Execute, Explain, and Save as..., and a 'New Worksheet' button. To the right of the Query Editor is a sidebar with icons for SQL, Settings, Help, and TEZ, with a red notification badge showing '10'. Below the Query Editor is a 'Query Process Results' section with a status of 'SUCCEEDED'. It has tabs for Logs and Results, and a 'Save results...' button. The results table shows three rows of data:

users.swid	users.birth_dt	users.gender_cd
0001BDD9-EABF-4D0D-81BD-D9EABFC0D07D	8-Apr-84	F
00071AA7-86D2-4EB9-871A-A786D27EB9BA	7-Feb-88	F
00071B7D-31AF-4D85-871B-7D31AFFD852E	22-Oct-64	F

## 6. SAVE AND EXECUTE A QUERY

Suppose we want to write a query, but not necessarily execute it immediately or perhaps we want to save it for future multiple uses. In the “**QUERY**” tab, copy/paste the following and click “**Save As**”. Save the query as omniture-view.

```
CREATE VIEW omniture AS
SELECT col_2 ts, col_8 ip, col_13 url, col_14 swid, col_50 city, col_51
country, col_53 state FROM omniturelogs
```

To view your saved queries, navigate to the “SAVED QUERIES” tab. For now, let’s open our saved query by clicking either on the first or second column:

The “Query Editor” should automatically open up, with your saved query preloaded. Click “Execute” to run this query, which will create a Hive view named “omniture”, a refined subset of data with only a handful of fields.

If you refresh the Database Explorer using the icon, you will see the “omniture” view added in the list of objects. You can click on the icon next to it to see a data sample.

## 7. JOIN DATA FROM MULTIPLE TABLES

Let's play with our data further, taking specific fields from different tables and creating a custom table from them.

```
CREATE TABLE webloganalytics AS
SELECT to_date(o.ts) logdate, o.url, o.ip, o.city, upper(o.state) state,
o.country, p.category, CAST(datediff(from_unixtime(unix_timestamp()), 
from_unixtime(unix_timestamp(u.birth_dt, 'dd-MMM-yy'))) / 365 AS INT)
age, u.gender_cd FROM omniture o
INNER JOIN products p
ON o.url = p.url
LEFT OUTER JOIN users u
ON o.swid = concat('{', u.swid, '}');
```

You can click on the “**New Worksheet**” button and paste the above query there, then click on “**Execute**”:

The screenshot shows the Ambari Hive View interface. On the left, the Database Explorer lists databases including 'default', 'omniture', 'omniturelogs', 'products', 'sample\_07', 'sample\_08', 'users', 'webloganalytics', 'foodmart', and 'xademo'. The 'webloganalytics' database is highlighted. On the right, the Query Editor displays a SQL query:

```

1 CREATE TABLE webloganalytics AS
2 SELECT to_date(o.ts) logdate, o.url, o.ip, o.city, upper(o.state) state,
3 o.country, p.category, CAST(datediff(from_unixtime(unix_timestamp()), from_unixtime(unix_time
4 FROM omniture o
5 INNER JOIN products p
6 ON o.url = p.url
7 LEFT OUTER JOIN users u
8 ON o.swid = concat('{', u.swid, '}');

```

Below the editor are buttons for 'Execute', 'Explain', and 'Save as...', and a 'New Worksheet' button.

After refreshing the “Database Explorer”, the new “webloganalytics” object will appear there. You can click on the icon next to it to get a data sample and browse the “Results” tab at the bottom.

The screenshot shows the Ambari Hive View interface. The 'webloganalytics' database is selected in the Database Explorer. The Query Editor shows a sample query:

```

1 SELECT * FROM webloganalytics LIMIT 100;

```

Below the editor, the 'Query Process Results' table shows the following data:

Logs	Results	Save results...	
Filter column...		previous next	
webloganalytics.logdate	webloganalytics.url	webloganalytics.ip	webloganalytics.city
2012-03-15	http://www.acme.com/SH55126545/VD55177927	69.76.12.213	coeur d alene
2012-03-15	http://www.acme.com/SH55126545/VD55166807	67.240.15.94	queensbury
2012-03-15	http://www.acme.com/SH55126545/VD55149415	67.240.15.94	queensbury
2012-03-15	http://www.acme.com/SH55126545/VD55179433	98.234.107.75	sunnyvale
2012-03-15	http://www.acme.com/SH55126545/VD55179432	74.85.146.98	san diego

## 8. SUMMARY

Excellent! We've learned how to upload data into HDFS, create tables and load data into them, and run queries for data refinement and enrichment by using Ambari's convenient and robust Hive View.

# Spark Lab – Introduction to Zeppelin

## 1. INTRODUCTION

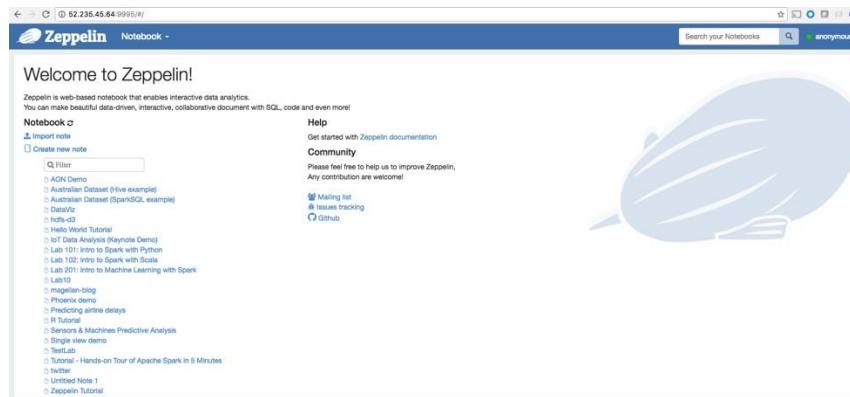
The HDP Sandbox includes the core Hadoop components, as well as all the tools needed for data ingestion and processing. You are able to access and analyze data in the sandbox using any number of Business Intelligence (BI) applications.

In this tutorial, you will access and explore Apache Zeppelin core capabilities. You will also interact with Hadoop cluster using various language interpreters embedded in zeppelin notebook. Finally, you reuse existing notebook to run an end-to-end prebuilt data pipeline.

*Caution! Please don't copy and paste in the notebooks, it will generate some errors in your syntax if you do so.*

## 2. ACCESS AND BROWSE ZEPPELIN

Open a web browser and enter the following URL: <https://mtlhdilabsspark<1-12>.azurehdinsight.net/zeppelin/>



Create a new note named Introduction to Zeppelin with default interpreter as Spark



### 3. EXPLORE INTERPRETER BINDINGS

At the top right click on the gear icon to see interpreter binding. Your administrator has enabled an interpreter called “spark” which is configured for the HDInsight cluster you are using. Click on Save button to return to your notebook page.

#### 4. WORKING WITH ZEPPELIN INTERPRETERS

Find the values for Spark version and the Spark home directory. When you type the commands, run them either by pressing the Shift + Enter keys, or by clicking on the Play icon to the right of the word Ready.

```
%spark  
sc.version  
sc.getConf.get("spark.home")
```



```
sc.version  
sc.getConf.get("spark.home")
```

The screenshot shows a Zeppelin notebook cell with the following code:  
`sc.version  
sc.getConf.get("spark.home")`

A red oval highlights the "READY" button, which is located to the right of the code input field. The status bar at the bottom of the window also shows "READY".

While processing, Zeppelin will display a status of RUNNING. It will also display a Pause icon should it become necessary.



```
sc.version  
sc.getConf.get("spark.home")
```

The screenshot shows a Zeppelin notebook cell with the following code:  
`sc.version  
sc.getConf.get("spark.home")`

The status bar at the bottom of the window shows "RUNNING 0%" and a pause icon.

The output may vary slightly from the screenshot below, but should look something like this when processing is completed:



```
sc.version  
sc.getConf.get("spark.home")
```

The screenshot shows a Zeppelin notebook cell with the following code:  
`sc.version  
sc.getConf.get("spark.home")`

The status bar at the bottom of the window shows "FINISHED 0%" and a refresh icon.

```
res0: String = 1.6.0  
res1: String = /usr/hdp/2.4.0.0-169/spark
```

Zeppelin can be instructed to use multiple languages in an interactive fashion within the same notebook. Simply specify the desired language prior to the command.

Run the following commands to demonstrate this flexibility using **Shell**, **Python**, **Scala**, **Markdown**, and **Spark SQL**. Execute each command by clicking on the Play icon or pressing **Shift + Enter** when you are finished typing.

### Shell:

```
%sh  
echo "Introduction to Zeppelin"
```

```
| %sh echo "Introduction to Zeppelin"  
Introduction to Zeppelin
```

### Python:

```
%pyspark  
print "Introduction to Zeppelin"
```

```
| %pyspark  
| print "Introduction to Zeppelin"  
Introduction to Zeppelin
```

### Scala (default, so no need to specify prior to running command):

```
%spark  
val s = "Introduction to Zeppelin"
```

```
| val s = "Introduction to Zeppelin"  
|  
| s: String = Introduction to Zeppelin
```

### Markdown:

```
%md Introduction to Zeppelin
```

```
| %md Introduction to Zeppelin  
|  
| Introduction to Zeppelin
```

## Spark SQL:

```
%sql  
show tables
```

The screenshot shows a Zeppelin notebook interface. At the top, there is a code cell with the command "%sql show tables". Below it is a table with two columns: "tableName" and "isTemporary". The table is currently empty, indicating no tables have been listed.

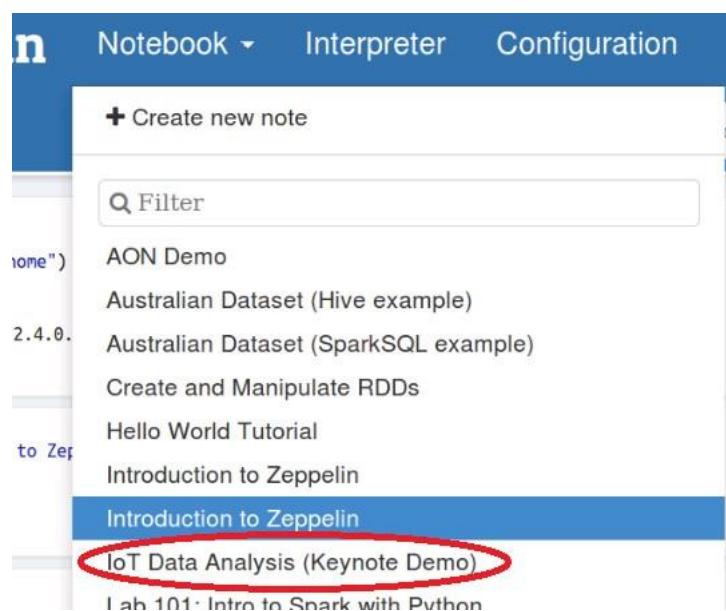
tableName	isTemporary

## 5. USE A PREBUILT NOTEBOOK TO EXPLORE ZEPPELIN CAPABILITIES

Zeppelin has four major functions:

- data ingestion
- data discovery
- data analytics
- data visualization

One of the easiest ways to explore these functions is with a preconfigured notebook, many of which are available by default. Click on Notebook at the top of the browser window and find and select the notebook labeled IoT Data Analysis (Keynote Demo) in the resulting drop-down menu.



For the purposes of this lab, all necessary code has already been entered for you in the saved notebook. All you have to do is scroll to the appropriate section and click the Play icon or press Shift + Enter.

```
Using Zeppelin for Data Science Tasks: Data Ingestion, Data Formatting, Exploratory Analysis and Model Building.

Data Science involves a typical sequence of tasks: acquiring data, cleaning it, analyzing it for relationships, and then building a model. Zeppelin allows you to do all these from one unified interface.

Play 2 seconds
```

```
First, let's load the data into HDFS and make sure we can access it.
```

```
%sh
whoami
curl -sSL -O "https://www.dropbox.com/s/ggj1robwxpl9vrt/iotdemo-notebook-data.zip"
unzip iotdemo-notebook-data.zip
hadoop fs -mkdir -p /user/zeppelin/iotdemo
hadoop fs -copyFromLocal -f trainingData /user/zeppelin/iotdemo/
hadoop fs -copyFromLocal -f enrichedEvents /user/zeppelin/iotdemo/
hadoop fs -ls /user/zeppelin/iotdemo/
zeppelin
Archive: iotdemo-notebook-data.zip
```

The first major block of code ingests data from an online source into HDFS and then displays those files using the shell scripting interpreter. Find and run that code.

**Note** that the label to the left of the Play icon says FINISHED, but this will not prohibit you from running the code again on this machine.

**Also note:** this notebook uses a deprecated command, hadoop fs, rather than the more updated hdfs dfs command we used in the previous lab. This should not affect the functionality of the demo.

```
%sh
whoami
curl -sSL -O "https://www.dropbox.com/s/ggj1robwxpl9vrt/iotdemo-notebook-data.zip"
unzip iotdemo-notebook-data.zip
hadoop fs -mkdir -p /user/zeppelin/iotdemo
hadoop fs -copyFromLocal -f trainingData /user/zeppelin/iotdemo/
hadoop fs -copyFromLocal -f enrichedEvents /user/zeppelin/iotdemo/
hadoop fs -ls /user/zeppelin/iotdemo/
```

When the code has finished, the output at the bottom should look like this:

```
| hadoop fs -ls /user/zeppelin/iotdemo/
zeppelin
Archive: iotdemo-notebook-data.zip
Found 2 items
-rw-r--r-- 3 zeppelin zeppelin      63570 2016-05-27 16:50 /user/zeppelin/iotdemo/enrichedEvents
-rw-r--r-- 3 zeppelin zeppelin     33084 2016-05-27 16:50 /user/zeppelin/iotdemo/trainingData
```

The next section of the notebook once again uses the shell scripting interpreter to view some of the raw data in one of the downloaded files. Scroll down and run this code, then view its output.

```
%sh
hadoop fs -cat /user/zeppelin/iotdemo/enrichedEvents | tail -n 10

Overspeed,"Y","hours",45,2773,-90.07,35.68,0,1,1
Lane Departure,"Y","hours",45,2773,-90.04,35.19,1,1,0
Normal,"Y","hours",45,2773,-90.68,35.12,1,0,0
Normal,"Y","hours",45,2773,-91.14,34.96,0,0,0
Normal,"Y","hours",45,2773,-91.93,34.81,0,0,0
Normal,"Y","hours",45,2773,-92.31,34.78,0,1,0
Normal,"Y","hours",45,2773,-92.09,34.8,0,0,0
Normal,"Y","hours",45,2773,-91.93,34.81,0,0,0
Normal,"Y","hours",45,2773,-90.68,35.12,0,0,0
```

The next section of the notebook performs actions necessary to import and use this data with Spark SQL. You may note that the status to the left of the Play icon is shown as ERROR. This is due to the fact that the file being manipulated did not exist at the time the notebook was opened on this system. Run this code and view the output.

```
val sqlContext = new org.apache.spark.sql.SQLContext(sc)
val eventsFile = sc.textFile("hdfs:///user/zeppelin/iotdemo/enrichedEvents")

case class Event(eventType: String,
                 isCertified: String,
                 paymentScheme: String,
                 hoursDriven: Int,
                 milesDriven: Int,
                 lat: Float,
                 long: Float,
                 isFoggy: Int,
                 isRainy: Int
```

The output should look like this:

```
eventsRDD.toDF().registerTempTable("enrichedEvents")

sqlContext: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@312d2a12
eventsFile: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[3] at textFile at <console>:29
defined class Event
eventsRDD: org.apache.spark.rdd.RDD[Event] = MapPartitionsRDD[5] at map at <console>:35
res4: Long = 1359
```

The next block of code utilizes Spark SQL to view this data. Run this code and examine the output.

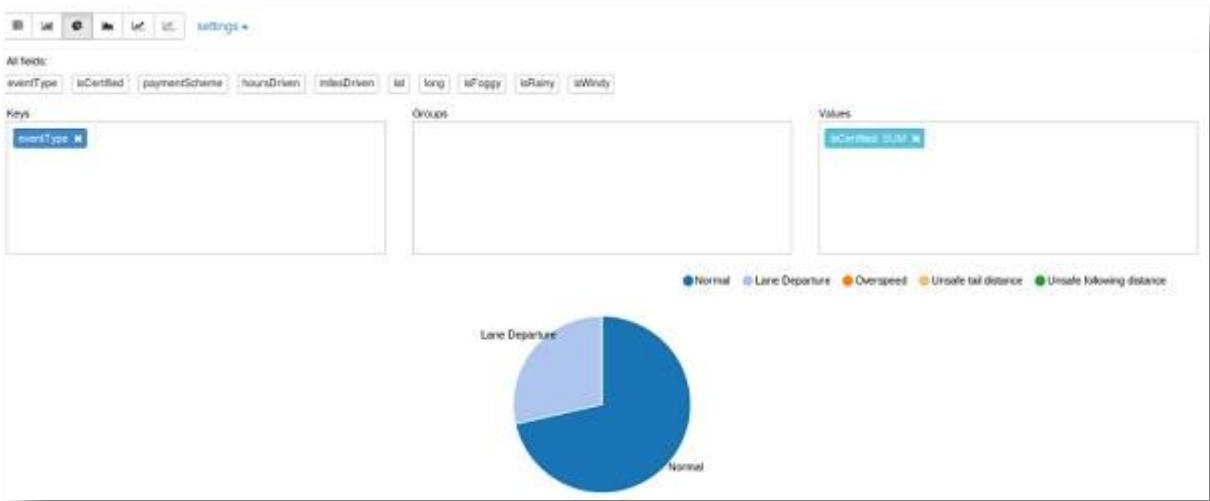
A screenshot of a Jupyter Notebook cell. The cell starts with the command `%sql`. Below it is the SQL query `select * from enrichedEvents order by hoursDriven desc limit 10`. The result is a table with the following data:

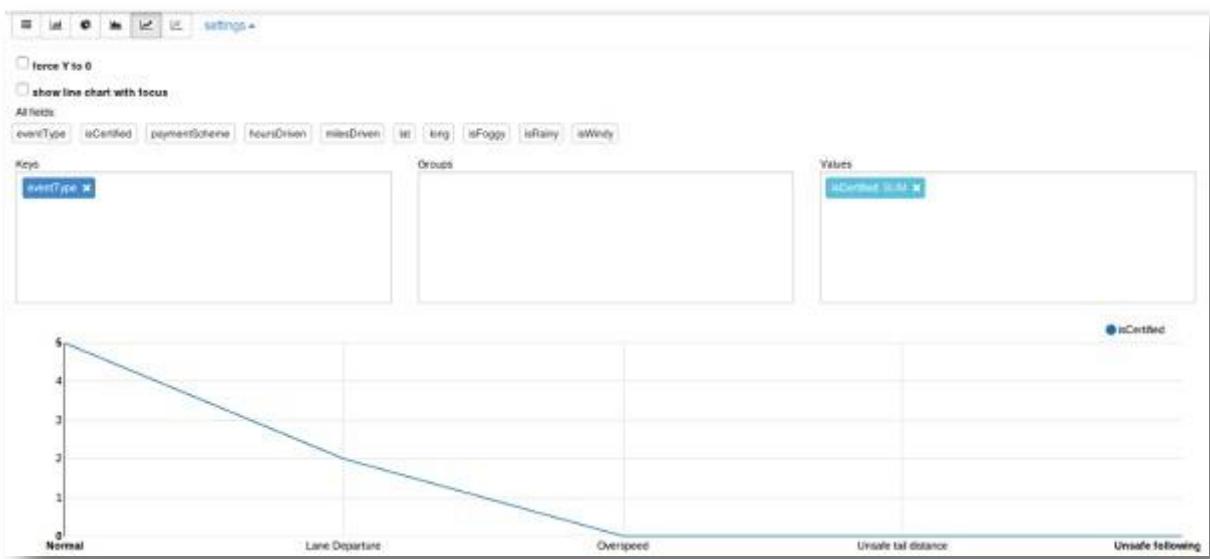
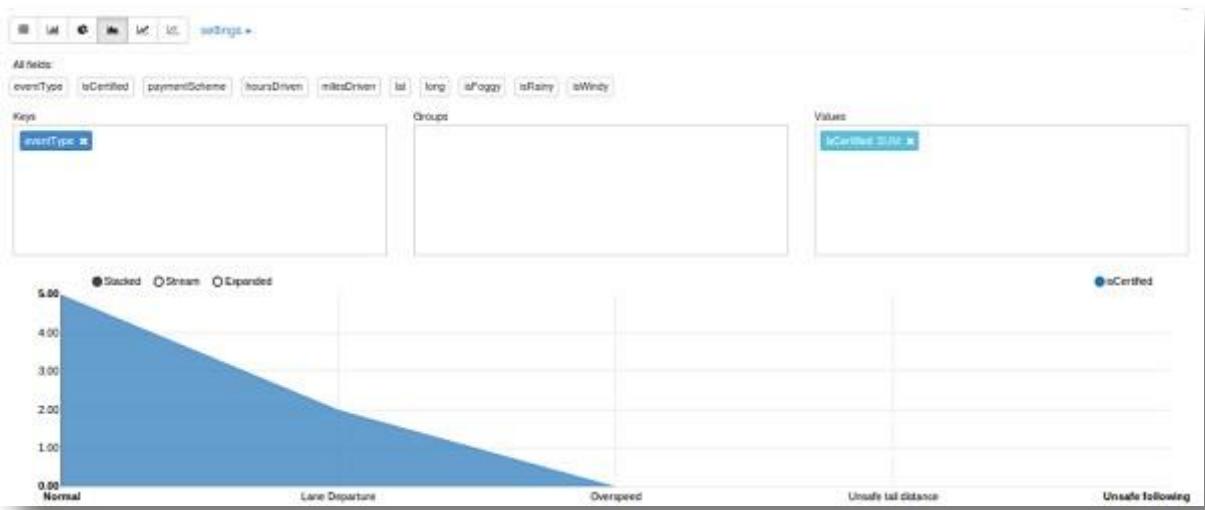
	eventType	isCertified	paymentScheme	hoursDriven	milesDriven	lat	long	isFoggy	isRainy
1	Normal	N	miles	90	4,300	-90.29	40.96	0	0
2	Lane Departure	N	miles	90	4,300	-88.42	41.11	1	1
3									

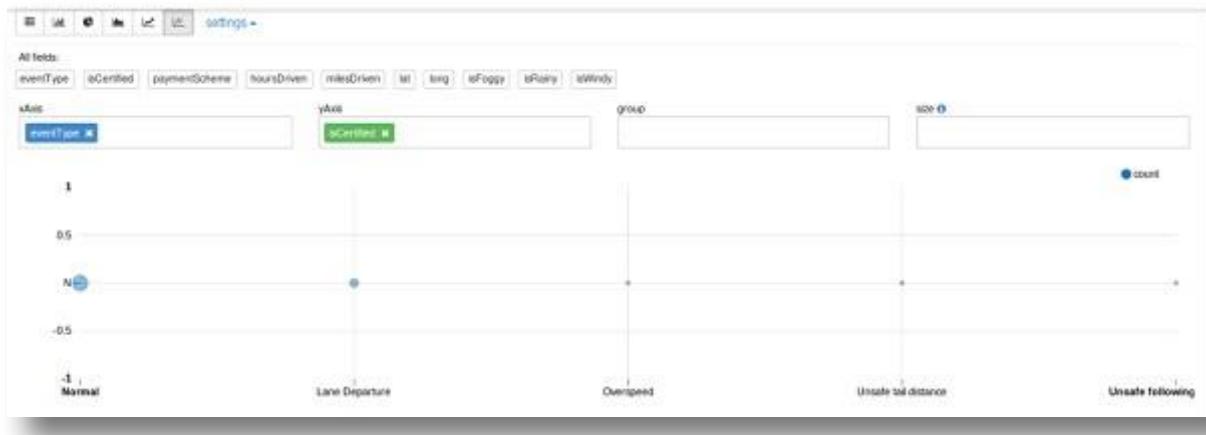
Note that at the top of the results there are six buttons that allow you to display the results using six different visualizations. Click on each one to view the differences between them.

A screenshot of a Jupyter Notebook cell showing the same SQL query and table as the previous image. At the top of the results, there are six small buttons with icons: a grid, a line graph, a bar chart, a pie chart, a scatter plot, and a histogram. The table data is identical to the previous screenshot.

eventType	isCertified	paymentScheme	hoursDriven	milesDriven	lat	long	isFoggy	isRainy	isWindy
Normal	N	miles	90	4,300	-90.29	40.96	0	0	1
Lane Departure	N	miles	90	4,300	-88.42	41.11	1	1	1
Normal	N	miles	90	4,300	-89.91	40.86	0	0	0
Overspeed	N	miles	90	4,300	-93.04	41.71	1	0	0
Unsafe tail distance	N	miles	90	4,300	-87.67	41.87	1	1	1
Normal	N	miles	90	4,300	-89.52	40.7	0	0	0
Normal	N	miles	90	4,300	-91.05	41.72	0	0	1
Normal	N	miles	90	4,300	-91.47	41.74	0	0	0
Lane Departure	N	miles	90	4,300	-91.59	41.7	1	0	0







**TIP:** In this lab you ran each section of code, known as a paragraph, individually. The entire notebook could have been played at once, however, by clicking the Play icon labeled Run all paragraphs directly to the right of the notebook title at the top of the browser.



## 6. SUMMARY

Congratulation! You have successfully created a Zeppelin notebook and demonstrated Zeppelin's ability to interpret multiple languages, and used a pre-built Zeppelin notebook to briefly explore Zeppelin's ability to ingest, view, analyze, and visualize data.

# Spark Lab - Data Visualization, Reporting and using Zeppelin (Scala)

## 1. INTRODUCTION

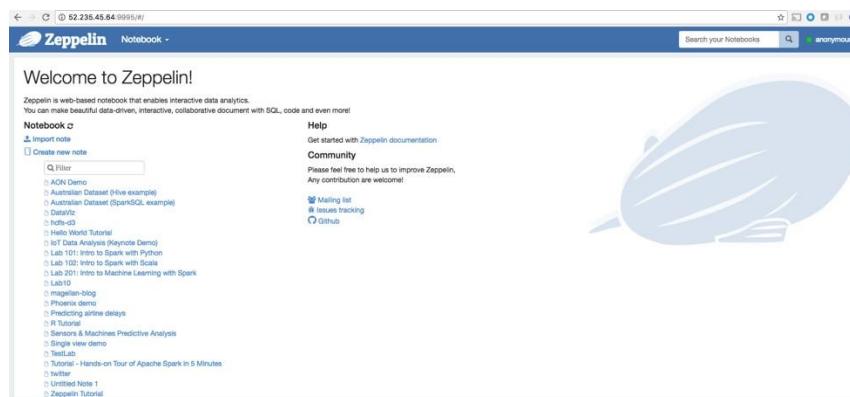
In this lab, we are going to learn to how to use Zeppelin to perform data visualizations, collaborate, and integrate visualization into reports.

We will describe how to ingest data into HDFS, to create table and manipulates tables with apache spark SQL and perform analytic queries on those tables, visualize the query results in real-time. Finally, we will show how collaborate on data project using zeppelin to create dynamic data visualization assets and collaborate with other users.

**Caution! Please don't copy and paste in the notebooks, it wil generate errors in the code syntax if you do so.**

## 2. CREATE DATA VISUALIZATIONS NOTEBOOK

Open a web browser and enter the following URL: <https://mtlhdilabsspark<1-12>.azurehdinsight.net/zeppelin/>



Create a new note named Data Visualization

The screenshot shows the Zeppelin web-based notebook interface. At the top, there's a navigation bar with links for 'Notebook', 'Search your Notebooks', and 'anonymous'. Below the header, a 'Welcome to Zeppelin!' message is displayed. On the left, a sidebar lists various notebooks and datasets, including 'AON Demo', 'Australian Dataset (Hive example)', 'Australian Dataset (SparkSQL example)', 'DataViz', 'hdfs-d3', and 'Hello World Tutorial'. A central panel shows a 'Create new note' dialog box with a 'Note Name' field containing 'Data Visualization'. Below the dialog, there's a note about contributing to Zeppelin and links for 'Mailing list' and 'Issues tracking'.

### 3. CREATE A DATAFRAME FROM HDFS FILE AND SAVE IT AS AN HIVE TABLE

Use the bankdata3.orc file to create a DataFrame named bankdata , a temporary table named banktemp , and a Hive table named bankdataperm .

```
%spark
val sqlContext = new org.apache.spark.sql.hive.HiveContext(sc)
val bankdata = sqlContext.read.format("orc").load("/Data/bankdata3.orc")
bankdata.registerTempTable("banktemp")
sqlContext.sql("create table bankdataperm as select * from banktemp")
```

```
val sqlContext = new org.apache.spark.sql.hive.HiveContext(sc)
val bankdata = sqlContext.read.format("orc").load("/user/zeppelin/bankfiles/bankdata3.orc")
bankdata.registerTempTable("banktemp")
sqlContext.sql("create table bankdataperm as select * from banktemp")
```

Use SQL to show the tables available and confirm that bankdataperm is available.

```
%sql
show tables
```

```
%sql  
show tables
```



tableName	isTemporary
health_table	true
bankdataperm	false
table1hive	false

#### 4. ANALYSE YOUR DATA

Use SQL to select and display all rows and columns from bankdataperm .

```
%sql  
select * from bankdataperm
```

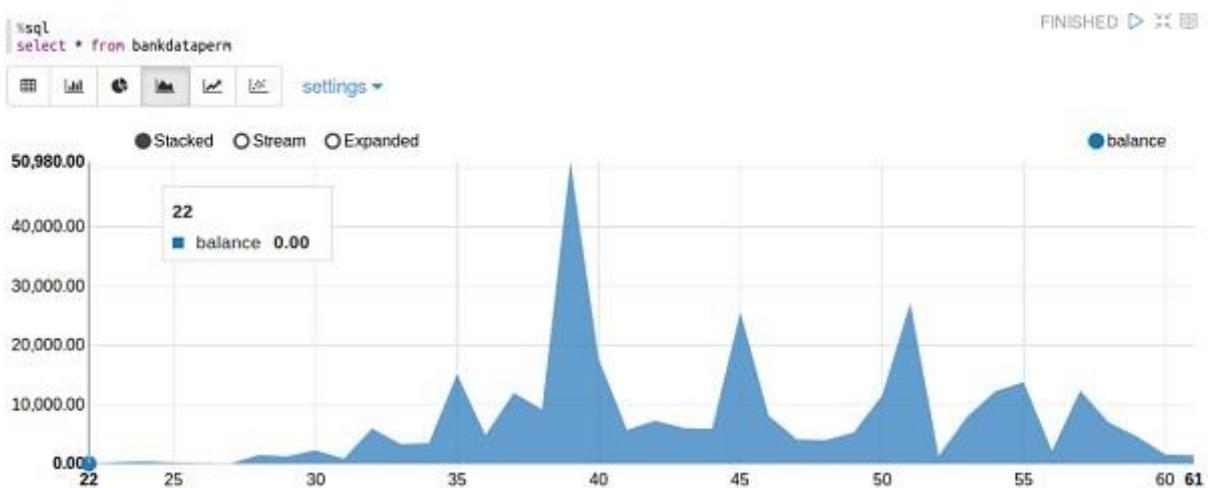
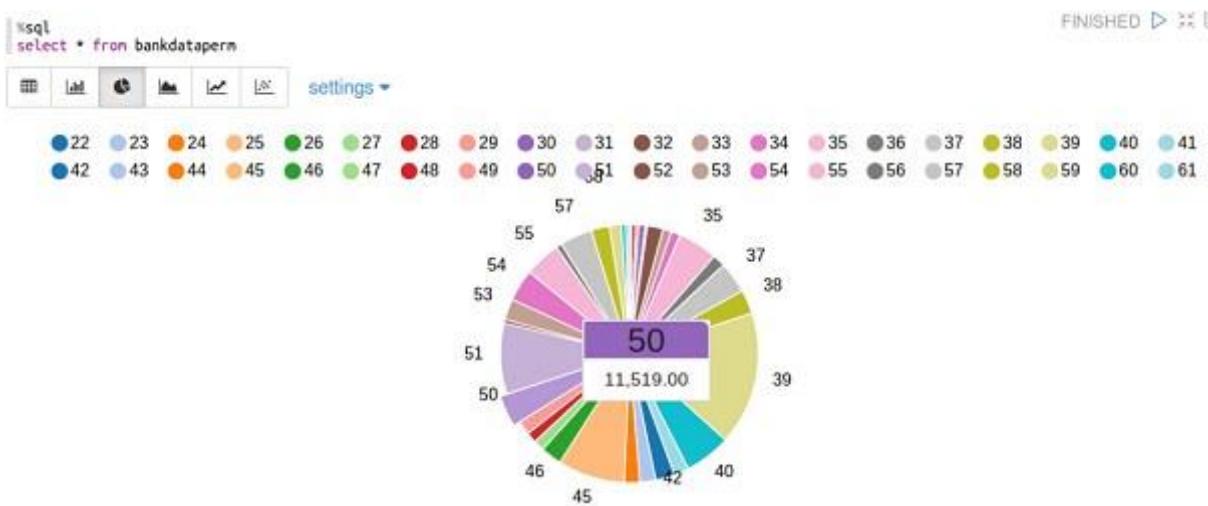
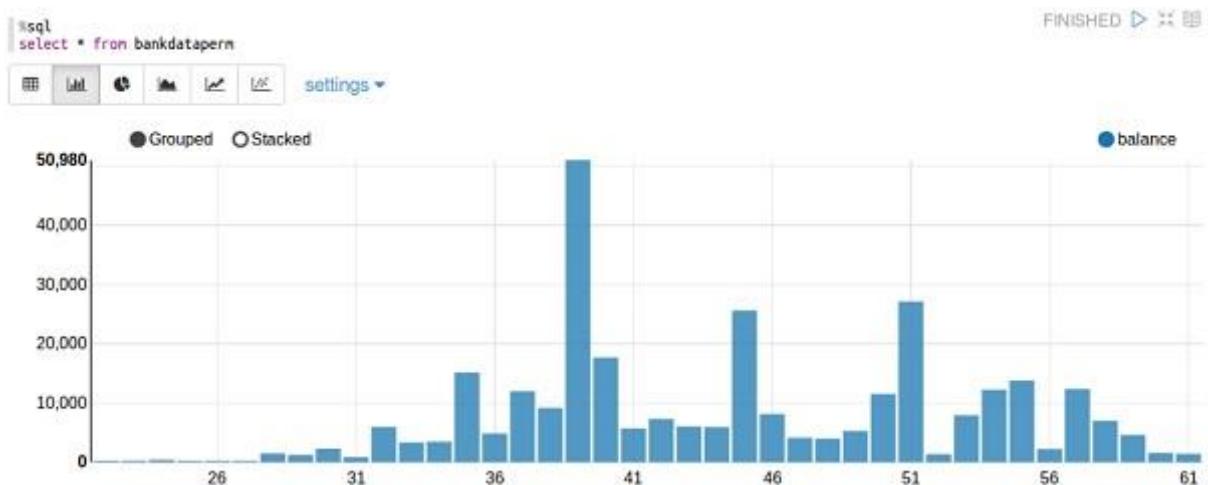
```
%sql  
select * from bankdataperm
```

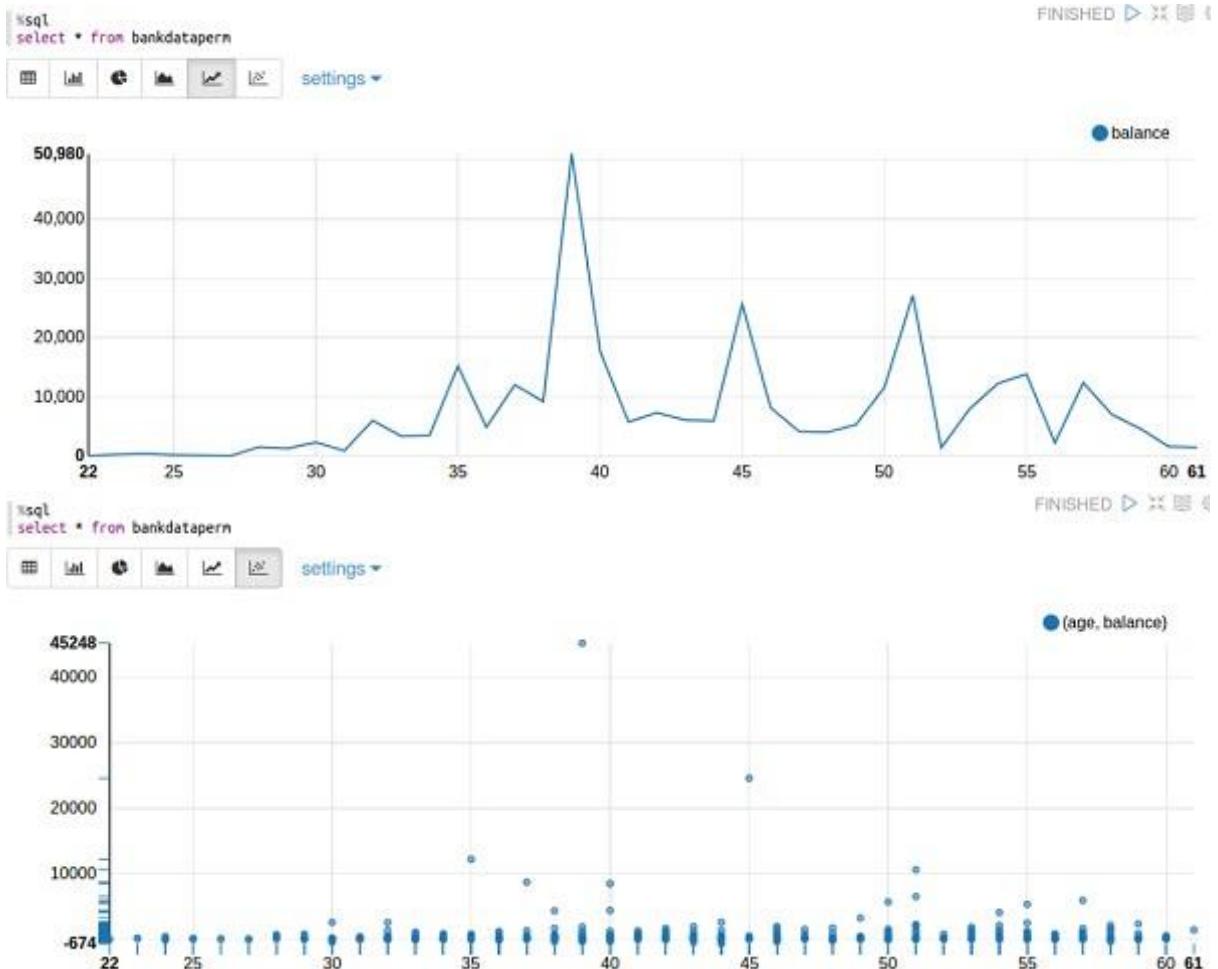


age	balance	marital
58	2,143	married
44	29	single
33	2	married
47	1,506	married
33	1	single

Quickly browse through the five data visualizations available by default in Zeppelin.

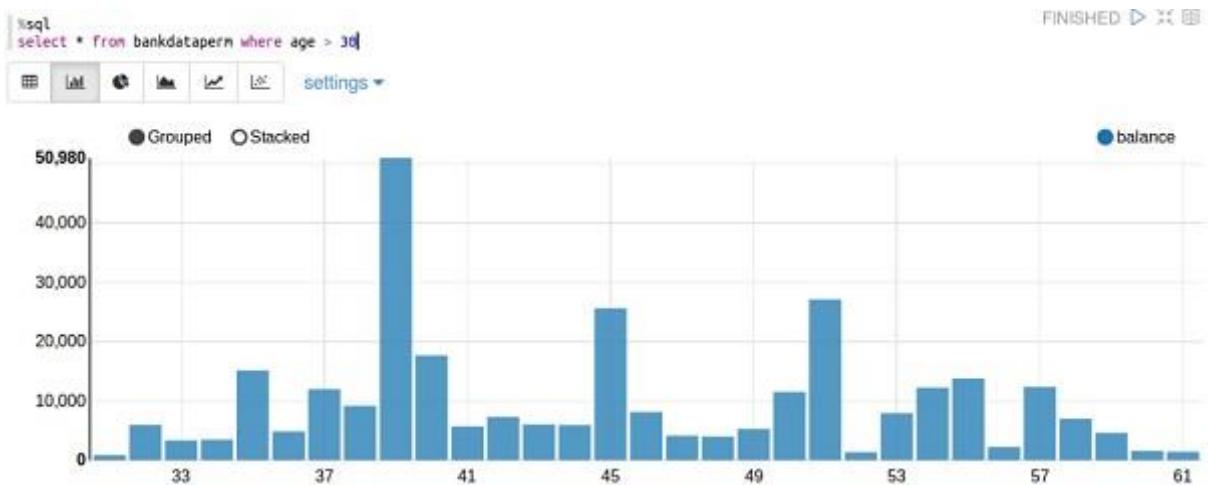






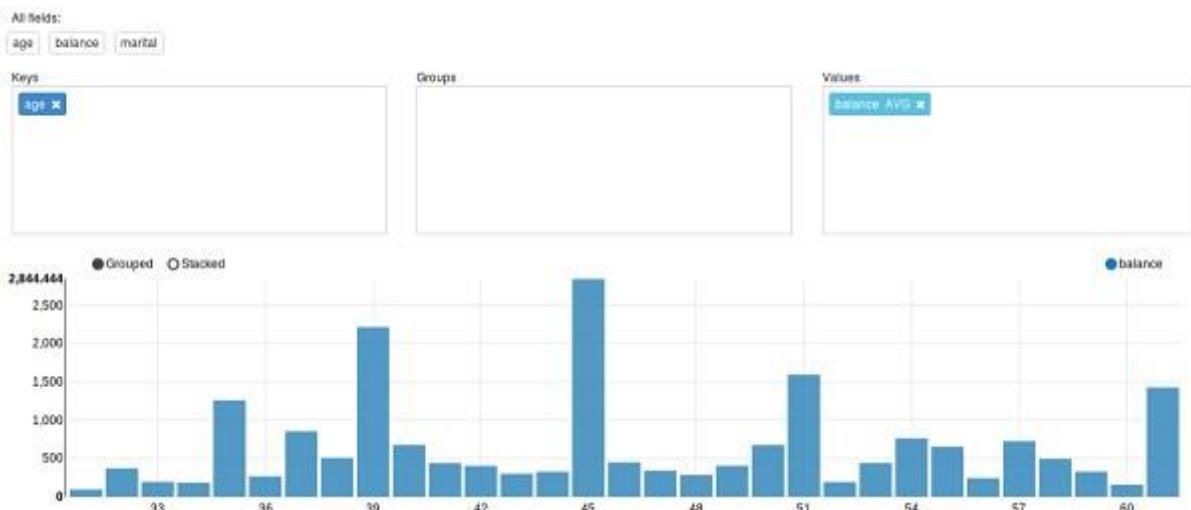
Now, go back to the bar chart view. Then, edit your SQL query so that it only shows data for individuals over the age of 30. Run the query and note the change in the chart.

```
%sql
select * from bankdataperm where age > 30
```



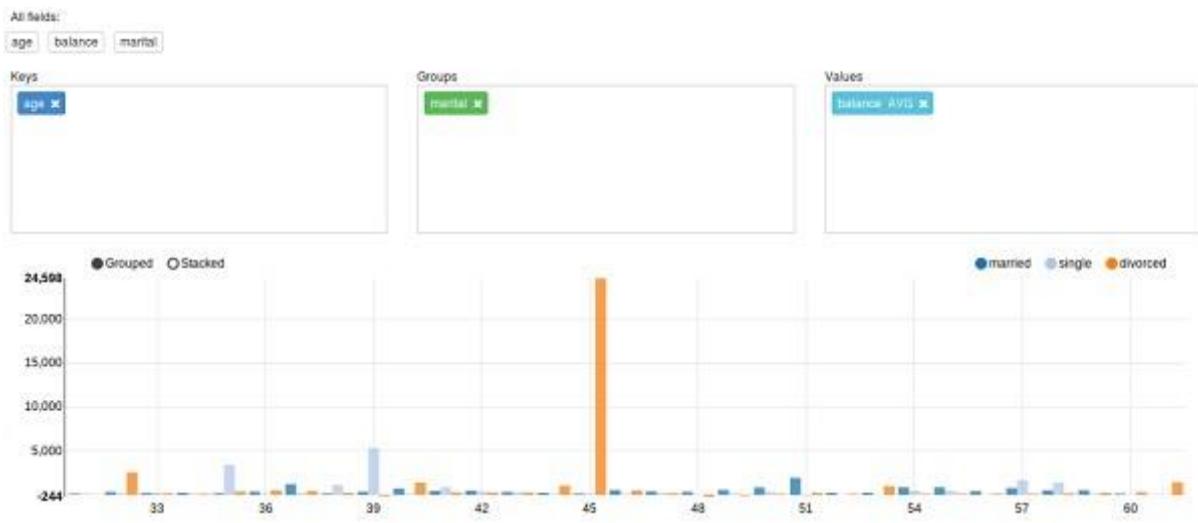
Click on the settings link and notice that Zeppelin has selected the age column as the key column and is showing the sum of the balances for all individuals in each age bracket. Display the average balance instead of the sum of balances.





Click and drag the available **marital** field into the **Groups** category to modify the visualization so that data is shown not only by age, but also grouped by **marital status**. When you are finished, click the settings link again to close the pivot chart options.





## 5. USE DYNAMIC VARIABLES IN YOUR ANALYSIS

It appears that we have a single outlier that is skewing the data significantly. We can easily see that the vast majority of average balances are well below \$5,000. Add a dynamic form to the SQL query that allows you to filter out data where the maximum balance for any individual exceeds a certain threshold, but set the default to 1,000,000 so that it doesn't immediately modify the chart. Rerun the query with this new code, then use this dynamic form to adjust the maximum balance to \$10,000 and \$5,000 and note the effects on the visualization.

```
%sql
select * from bankdataperm where age > 30 and
balance <= ${Maximum Balance=1000000}
```

```
%sql
select * from bankdataperm where age > 30 and balance <= ${Maximum Balance=1000000}
```

```
sql  
select * from bankdataper where age > 30 and balance == $[Maximum Balance-2000000]
```

FINISHED ▶ ✎

Maximum Balance



Maximum Balance



```
sql  
select * from bankdataper where age > 30 and balance == $[Maximum Balance-1000000]
```

FINISHED ▶ ✎

Maximum Balance





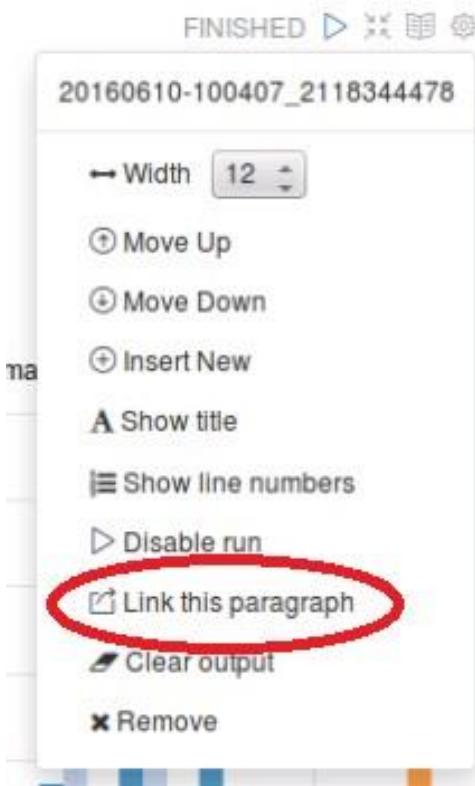
#### QUESTIONS:

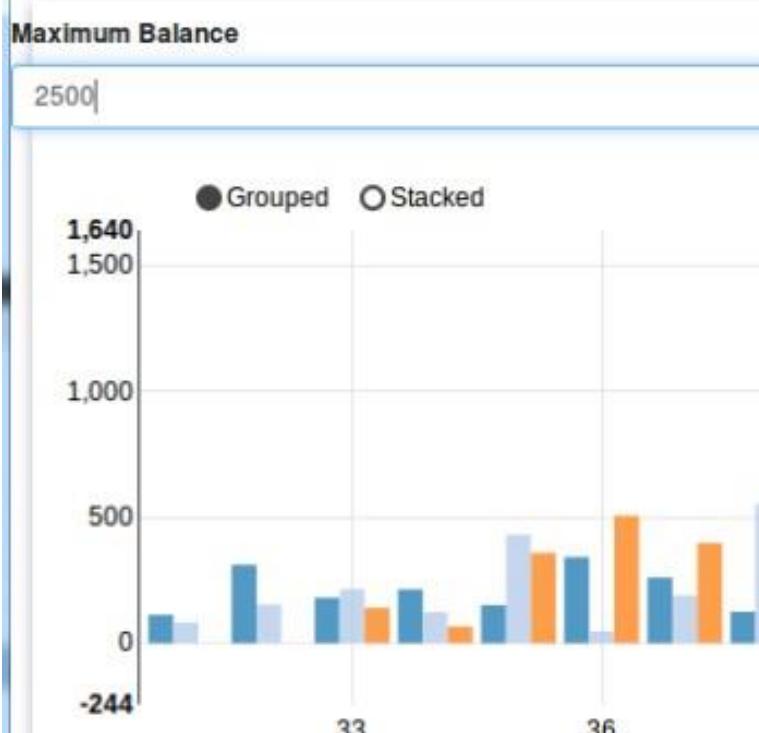
Why do you think changing the maximum balance from \$10,000 to \$5,000 had so little effect on the chart?

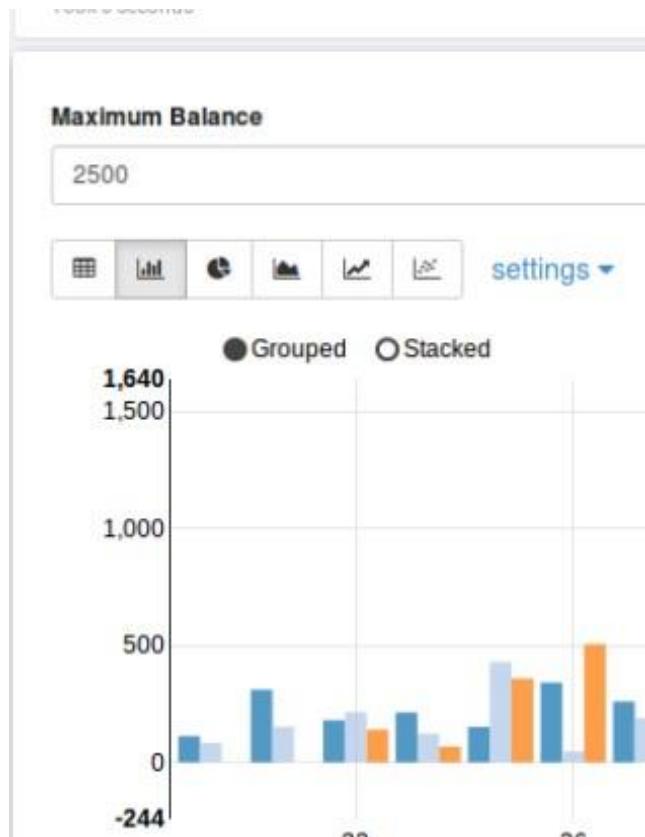
What group (married, single, or divorced) had the most change based on changing the maximum balance?

#### 6. SHARE YOUR FINDINGS WITH OTHERS

Create a URL that allows you to share this chart with others without giving them access to the code or the Zeppelin note. Use the linked page to change the maximum balance to \$2,500, then return to your note and observe the effects the change had at the source.

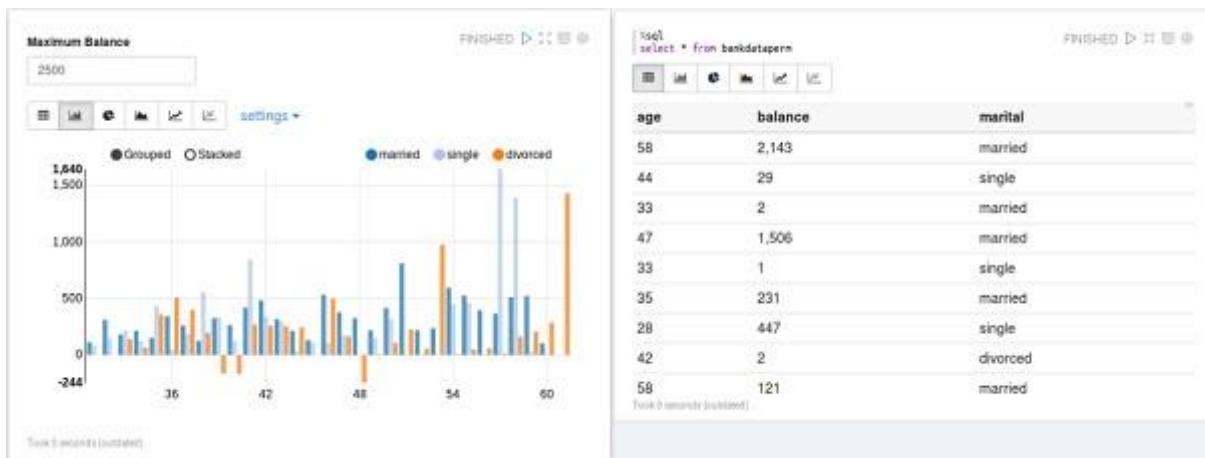
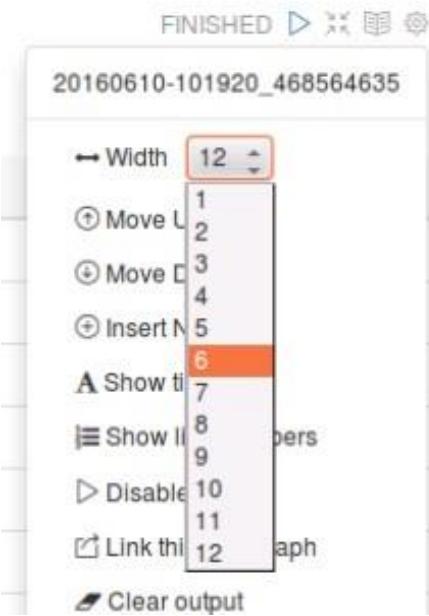






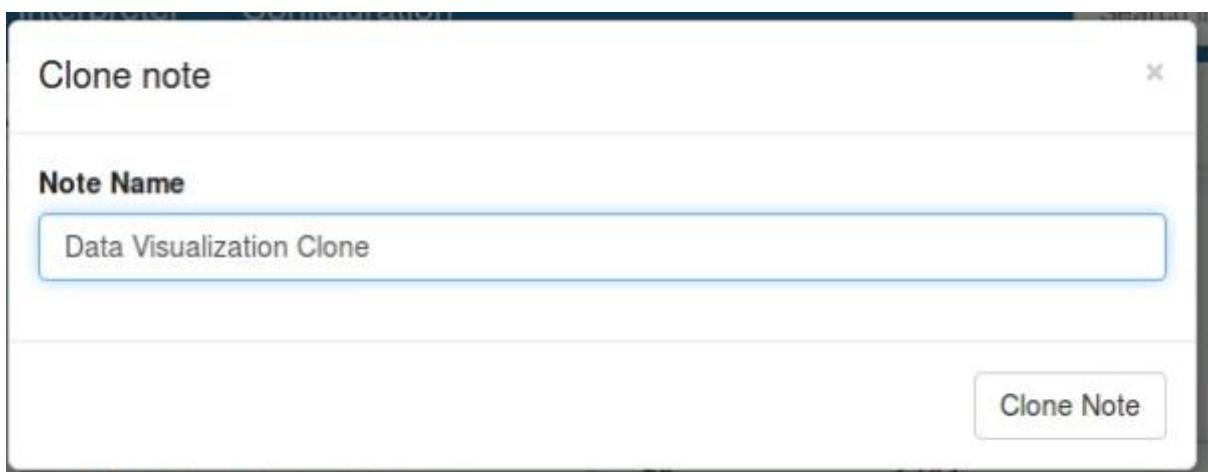
In the paragraph below this one, run the SQL command to read all data from bankdataperm. Then adjust the width of the two paragraphs so that they both appear on the same line.

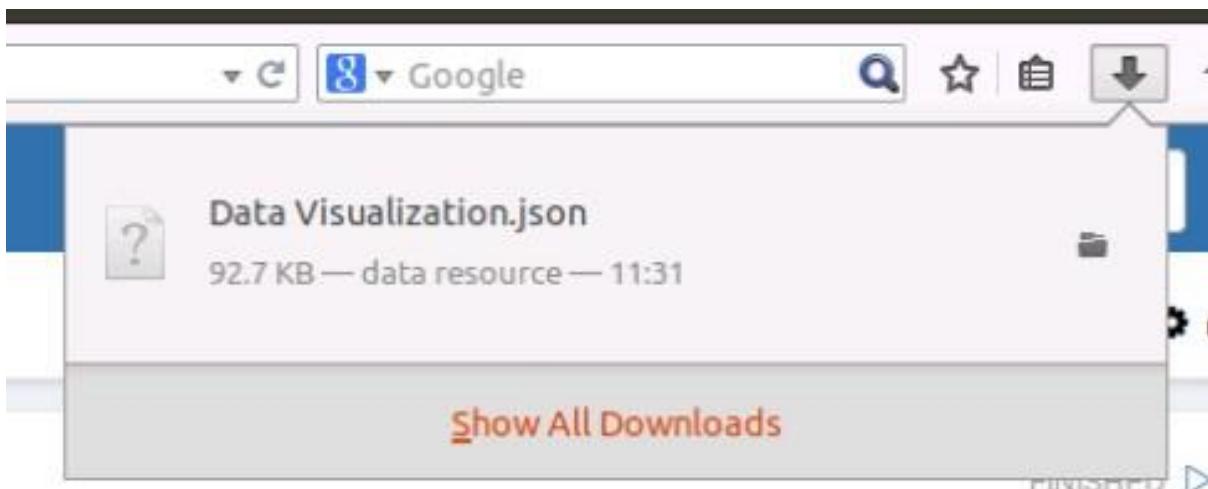




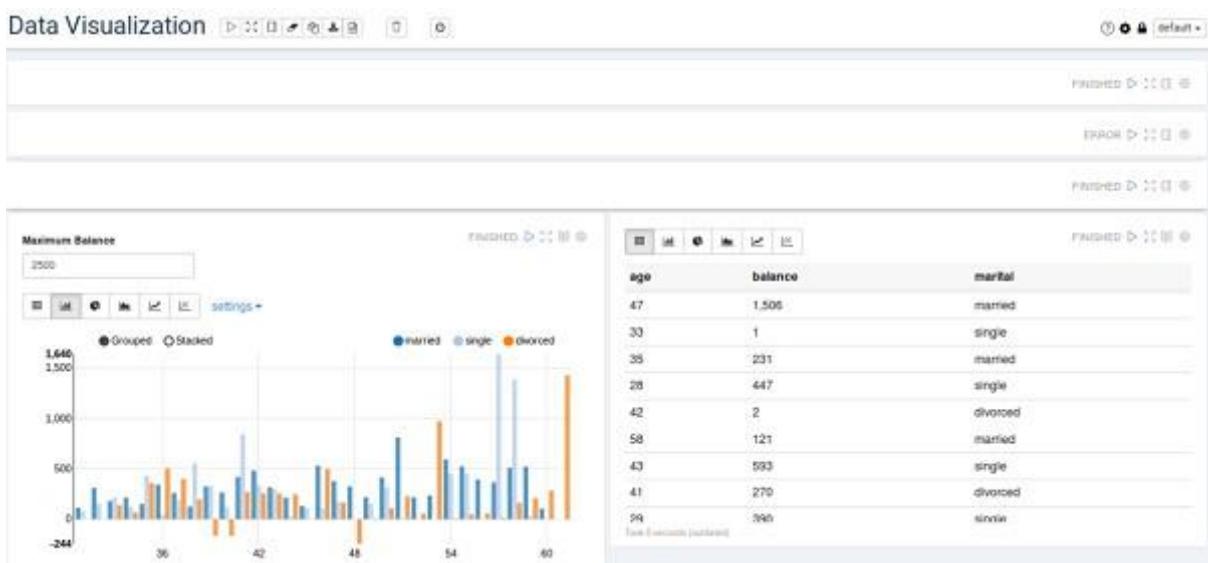
We are now ready to prepare this note for sharing. Create a clone copy of this note named Data Visualization Clone. Also export a copy of the note.





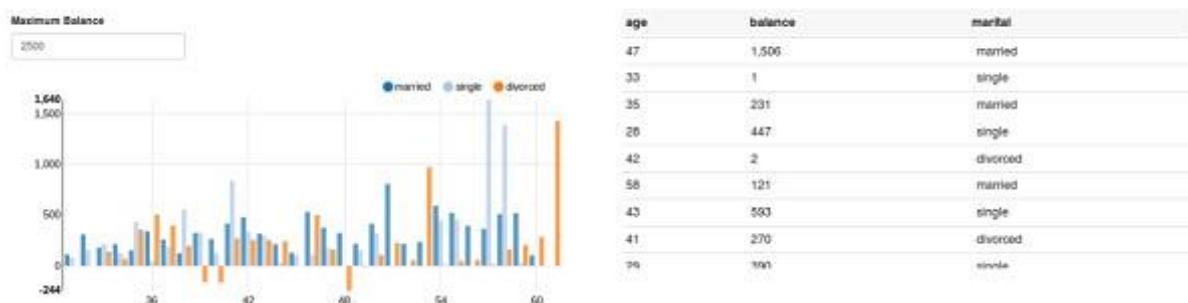


On the Data Visualization note we are going to share, hide the code for all paragraphs. Then hide the output for every paragraph except for the two that are on the same line.



Next, convert this from the default view to report view. Now the URL to this note is ready to share with your stakeholders.

### Data Visualization



Import the copy of this note you made earlier and name the new note Data Visualization Imported. Confirm that the copy contains all original code and formatting.

Zeppelin is web-based notebook that enables interactive data analytics.  
You can make beautiful data-driven, interactive, collaborative document with SQL, code, and more.

Notebook

Import note  

Create new note

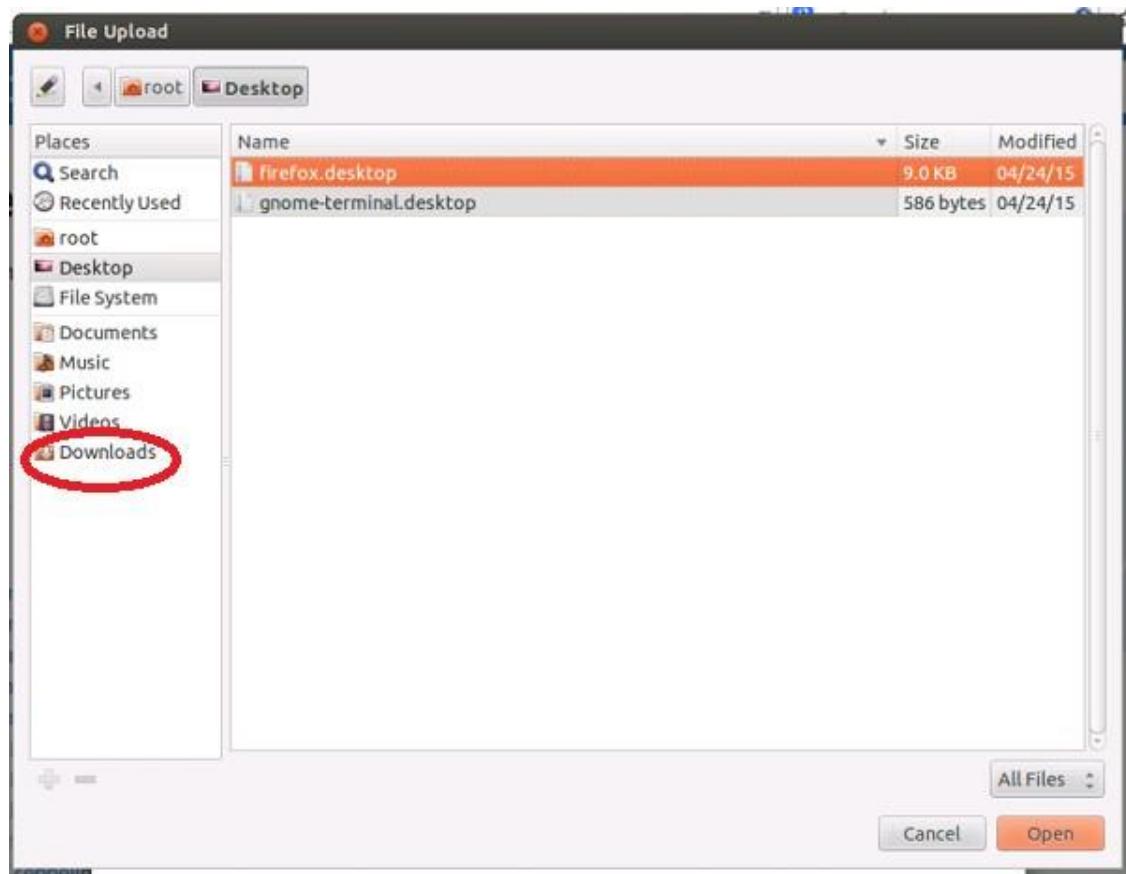
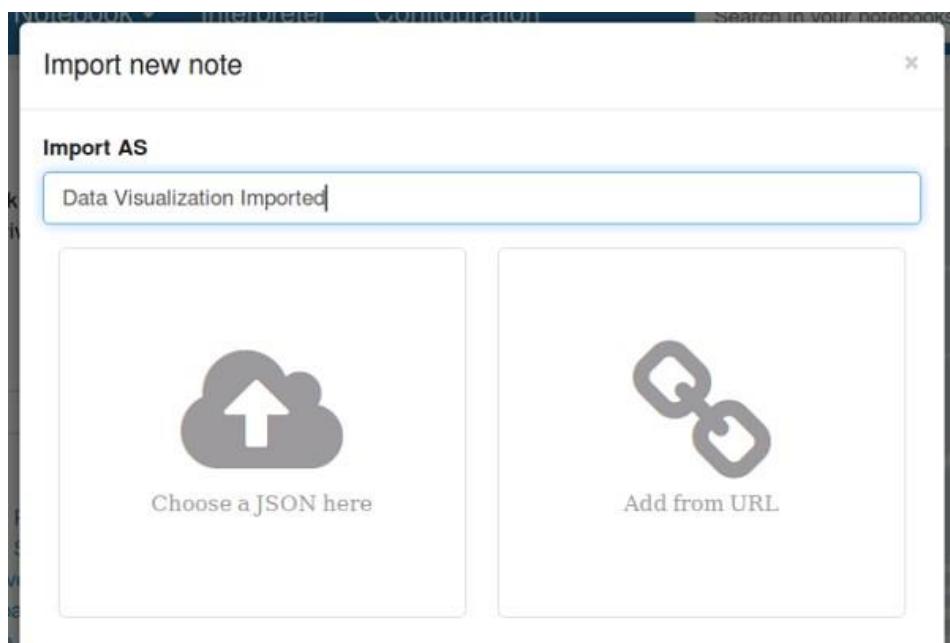
Filter

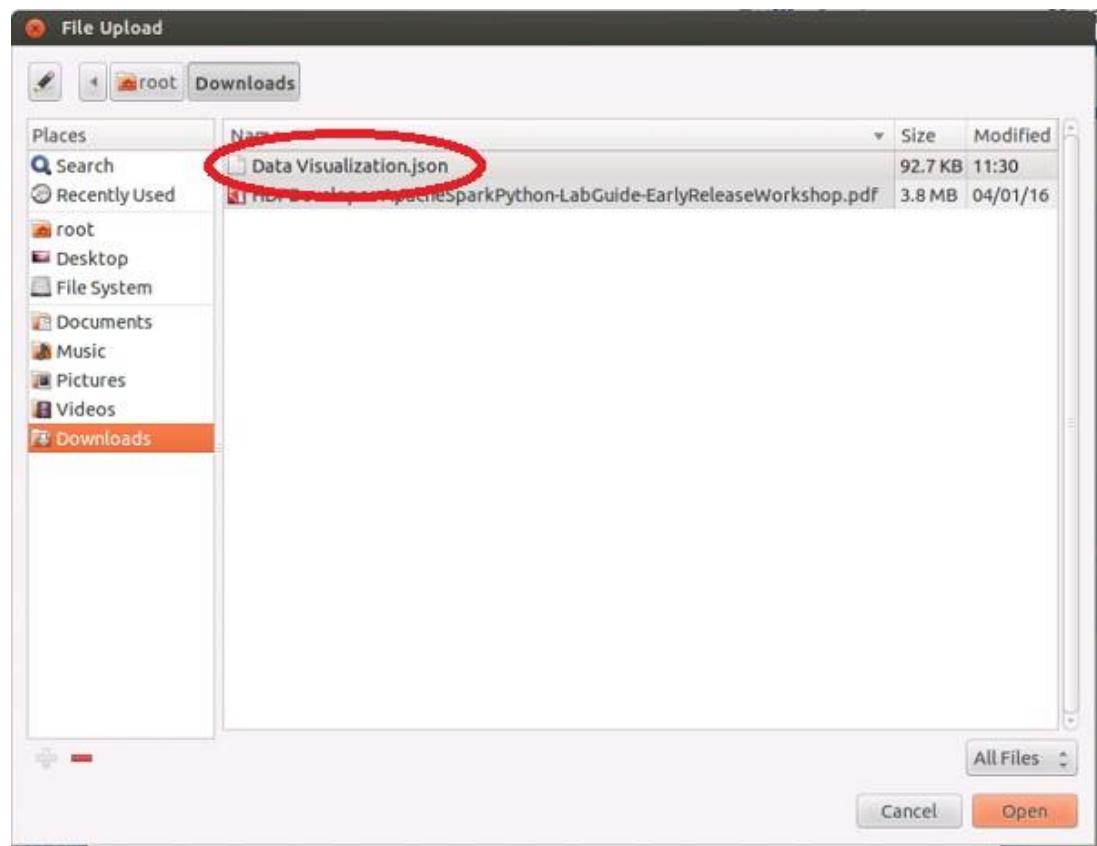
Help

Get started with [Zeppelin documentation](#)

Community

Please feel free to help us to improve this page.





- [Data Visualization](#)
- [Data Visualization Clone](#)
- [Data Visualization Imported](#)
- [Hello World Tutorial](#)

Data Visualization Imported

The screenshot shows the Zeppelin interface with a top navigation bar and several cells. The first cell contains a table with three rows: 'maximun\_value' (true), 'bankdataperm' (false), and 'table1hive' (false). The second cell contains a search bar with 'Maximum Balance' and a value of '2500'. The third cell contains a SQL query: 'select \* from bankdataperm'.

## 8. SUMMARY

Congratulation! You have successfully created and manipulated Zeppelin visualizations, made them available for collaboration, and used Zeppelin to create a shareable report.

# SparkSQL Lab – Using SparkSQL on Hive Tables

## 1. INTRODUCTION

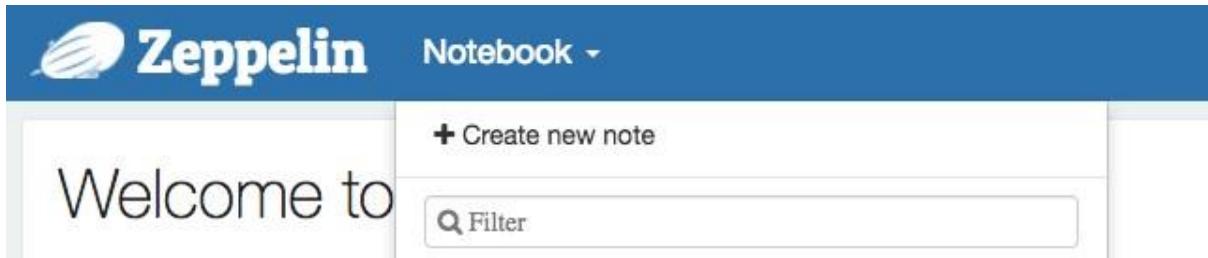
In this lab, you will use Zeppelin to query data using Spark SQL on Hive Table. You will learn to visualize data, and write data after processing with SparkSQL back to a Hive Table or a text file on HDFS.

*Caution! Please don't copy and paste in the notebooks, it will generate errors in the code syntax if you do so.*

## 2. CREATE A NEW ZEPPELIN NOTEBOOK

Open a web browser and enter the following URL: <https://mtlhdilabsspark<1-12>.azurehdinsight.net/zeppelin/>

Expand the Notebook drop down and click on “+Create new note”:



Create a new note called for example: “**Hive Data Exploration with SparkSQL**”.

## 3. READING A HIVE TABLE IN A SPARK DATAFRAME

In the first paragraph, we are going to invoke the Spark Hive Context sqlc to query the omniture Hive table created in the Hive lab using the spark interpreter, and read the data into a Spark DataFrame object:

```
%spark
val sqlContext = new org.apache.spark.sql.SQLContext(sc)

val df = sqlContext.sql("select country, count(*) as num_visits from
omniture group by country order by num_visits desc")
```

We can make basic Spark Dataframe APIs to print the schema of the dataframe object, and display the dataset (still in the same paragraph):

```
//Basic DataFrame operations
df.printSchema()
df.show()
```

At this point, you can run the paragraph by clicking the play icon (Blue triangle) on the top right side of the paragraph. This should show you the data frame schema as well as the data frame contents:

The screenshot shows a Zeppelin Notebook interface. The title bar says "Zeppelin Notebook". The search bar says "Search your Notebooks". The user is "anonymous". A tooltip "Run this paragraph (Shift+Enter)" points to the play icon in the toolbar. The main area contains the following code and its output:

```
%spark
//Reading SparkSQL query results on a Hive table in a Spark DataFrame object
val df = sqlc.sql("select country, count(*) as num_visits from omniture group by country order by num_visits desc")

//Basic DataFrame operations
df.printSchema()
df.show()

df: org.apache.spark.sql.DataFrame = [country: string, num_visits: bigint]
root
 |-- country: string (nullable = true)
 |-- num_visits: long (nullable = false)
+-----+-----+
|country|num_visits|
+-----+-----+
|  usal |    1136441|
|   pril |      531|
|    ausl |      571|
|    deul |      411|
|    canl |      331|
|    thal |      321|
|    czel |      231|
|    phll |      221|
|    prtl |      191|
|    zafl |      191|
|    芬nl |      171|

```

Below the code, it says "Took 2 sec. Last updated by anonymous at September 21 2017, 9:51:25 AM."

After the processing was done in SparkSQL, we can now write the results back as a Hive table, as well as save the results in a text file on HDFS using the following commands. On the same paragraph, type:

```
//Saving DataFrame as a Hive table
df.write.saveAsTable("default.country_visits")

//Saving DataFrame on HDFS
df.repartition(1).rdd.saveAsTextFile("/tmp/country_visits")
```

Execute the paragraph again by clicking on the Run icon.

**Zeppelin** Notebook Search your Notebooks  anonymous +

## Hive Data Exploration with SparkSQL

```
%spark
//Reading SparkSQL query results on a Hive table in a Spark DataFrame object
val df = sqlc.sql("select country, count(*) as num_visits from omniture group by country order by num_visits desc")

//Basic DataFrame operations
df.printSchema()
df.show()

//Saving DataFrame as a Hive table
df.write.saveAsTable("default.country_visits")

//Saving DataFrame on HDFS
df.repartition(1).rdd.saveAsTextFile("//tmp/country_visits")

df: org.apache.spark.sql.DataFrame = [country: string, num_visits: bigint]
root
 |-- country: string (nullable = true)
 |-- num_visits: long (nullable = false)
+-----+-----+
|country|num_visits|
+-----+-----+
|  usal | 1136441 |
|  pril |    531 |
|  ausl |     57 |
|  deul |     41 |
|  canl |     33 |
|  thal |     32 |
|  czel |     23 |
|  phll |     22 |
|  prtl |     19 |
|  zafl |     19 |
|  fapl |     17 |

```

FINISHED Run this paragraph (Shift+Enter)

Took 5 sec. Last updated by anonymous at September 21 2017, 10:01:07 AM.

If the paragraph ran successfully, we can now check the existence of the Hive table. Navigate to Ambari on a separate tab on your

. Navigate to the Hive view:

The screenshot shows the Ambari interface with the file browser. A file named 'country\_visits' is selected in the '/tmp' directory. A context menu is open, with 'Hive View' being the selected option. Other options in the menu are YARN Queue Manager, Files View, Pig View, Storm View, and Tez View.

Once the Hive view, click on the **Database Explorer**, and expand the **default** database. You should see the **country\_visits** table listed there. Click on the icon on the right hand side to get a sample of the data. In the **Results** tab at the bottom, you should see the contents of the table corresponding to the original data frame.

Database Explorer

Query Editor

Worksheet: country\_visits sample

```
1 SELECT * FROM country_visits LIMIT 100;
```

Databases:

- default
- omniture
- omniturelogs
- products
- sample\_07
- sample\_08
- users
- webloganalytics
- foodmart
- xademo

Execute | Explain | Save as... | New Worksheet

Query Process Results (Status: SUCCEEDED)

country_visits.country	country_visits.num_visits
usa	113644
nri	531

#### 4. EXPLORE AND VISUALIZE HIVE TABLES USING SPARKSQL

Start a new paragraph and invoke the SparkSQL interpreter `%sql`. Type the following query to get a drill down of the number of visits in the US per state, as this seems to be the main geography visiting the web site.

```
%sql
select state, count(*) as num_visits from omniture where country = 'usa'
group by state
order by num_visits desc
```

Click on the Run icon.

Run this paragraph (Shift+Enter)

FINISHED

state	num_visits
ca	14,395
fl	7,745
ny	7,540
tx	7,049
pa	4,610
oh	4,018
il	3,771
va	3,310
mi	3,289

Took 1 sec. Last updated by anonymous at September 21 2017, 10:30:41 AM.

Let's visualize this data in a pie chart. Click on the pie chart icon, and expand the settings drop down. Drag and drop “**state**” field in the Keys section, and **num\_visits** in the value section as follows:

The screenshot shows a data visualization tool interface. At the top, there is a SQL query editor with the following code:

```
%sql  
select state, count(*) as num_visits from omniture where country='usa' group by state order by num_visits desc
```

Below the query editor are several icons and a "settings" button, with two blue arrows pointing to the pie chart icon and the "settings" button respectively.

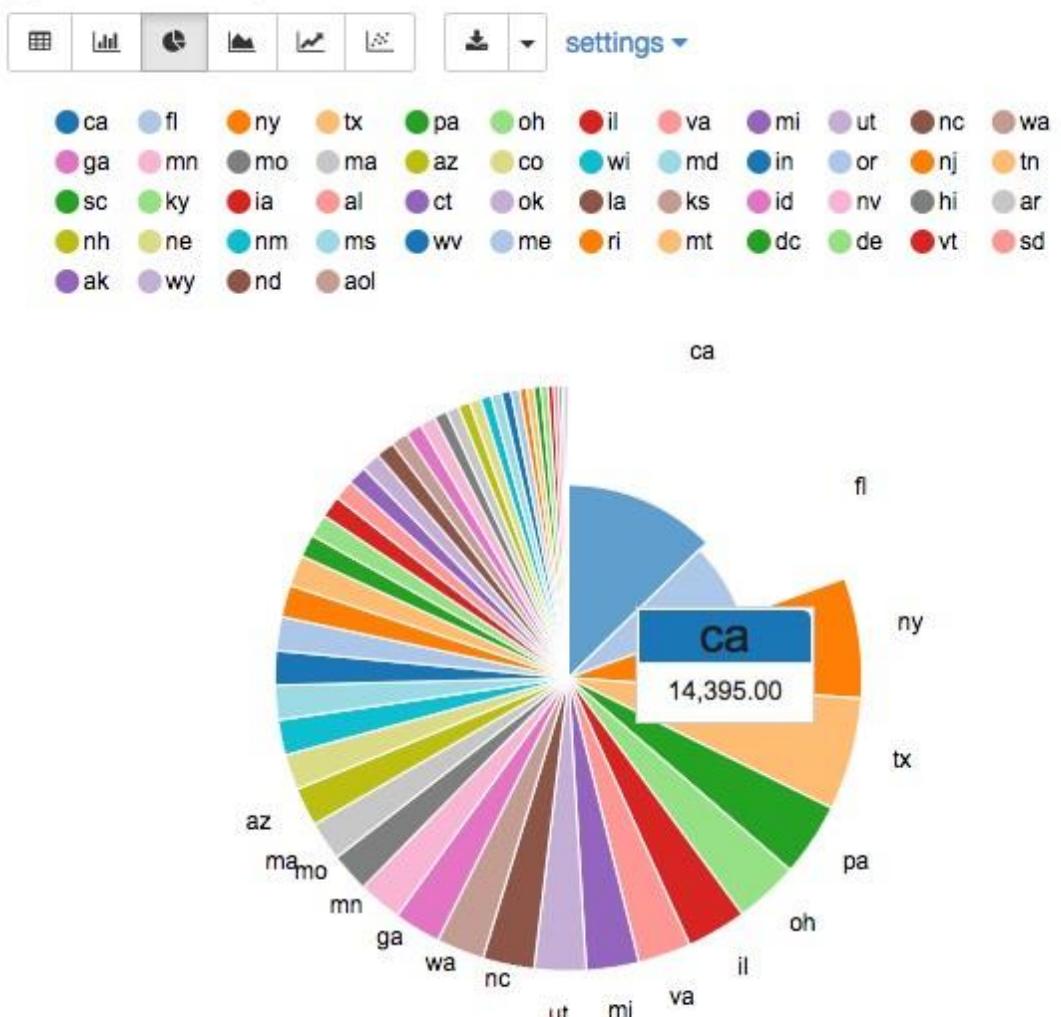
The main interface consists of three sections: "Keys", "Groups", and "Values".

- Keys:** Contains the field "state".
- Groups:** Empty.
- Values:** Contains the expression "num\_visits SUM".

At the bottom right of the interface, there is a "FINISHED" status indicator.

Click on settings again to hide the setting section. Resize the paragraph window appropriately and hover over the pie chart, you should see the different values corresponding to the various states:

```
%sql
select state, count(*) as num_visits from omniture where country='usa' group
by state order by num_visits desc
```

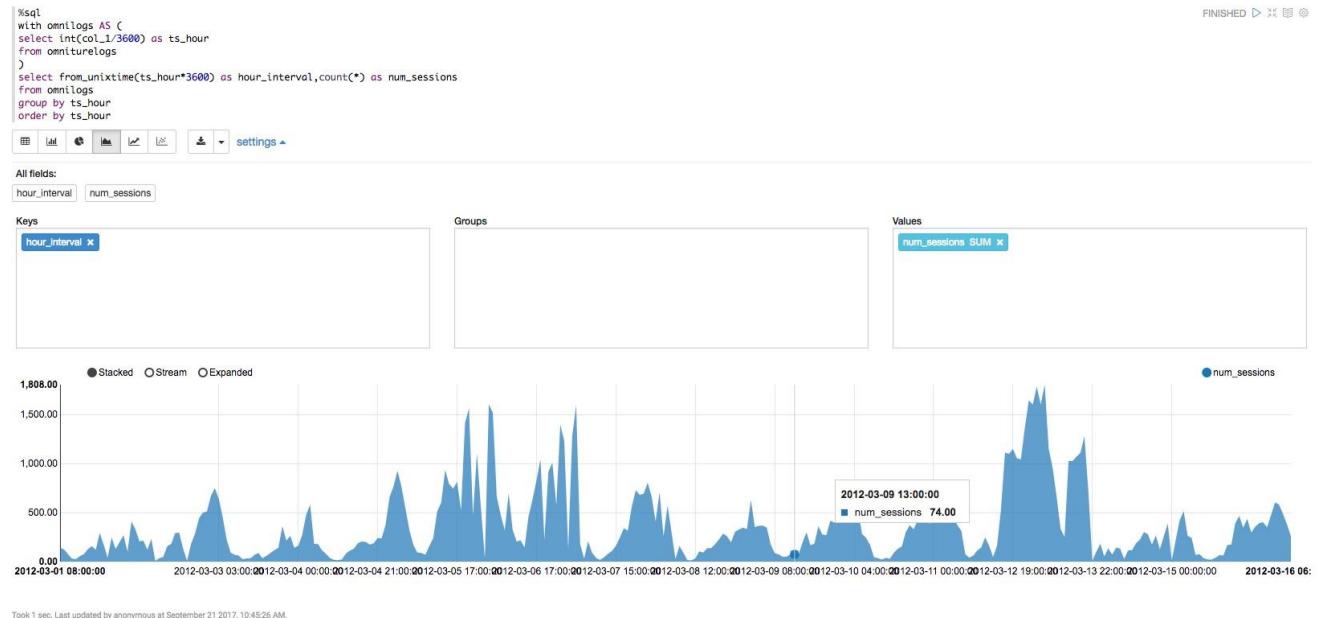


Took 1 sec. Last updated by anonymous at September 21 2017, 10:34:45 AM. (outdated)

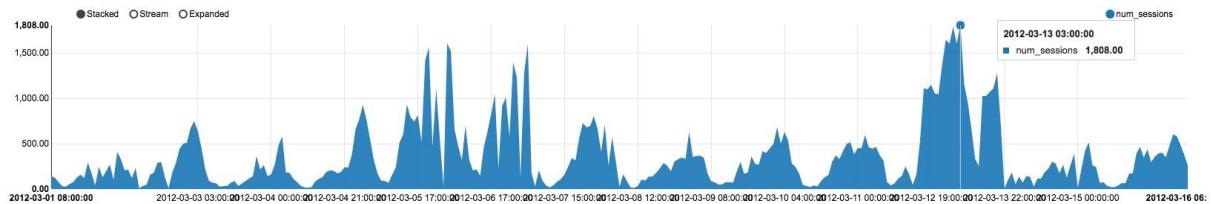
We can also try an interesting visualization which gives us the number of visits per hour from the omniture logs. Start a new paragraph, and type the following query:

```
%sql
with omnilog AS (
    select int(col_1/3600) as ts_hour
    from omniturelogs
)
select from_unixtime(ts_hour*3600) as hour_interval, count(*) as
num_sessions
from omnilog
group by ts_hour
order by ts_hour
```

Click on the run icon. We can visualize the data as a graph as follows. Select the graph icon (in grey below). Use **hour\_interval** as Keys, and **num\_sessions** as Values. You should see the following visualization which shows peak times when the web site was visited.



Hover over the graph to find out the peak from the period captured from the time interval of the web log:



## 5. SUMMARY

Excellent! We've learned how to access data in Hive tables using the SparkSQL API, as well as how to use the SparkSQL interpreter to seamlessly query and visualize Hive tables using Zeppelin.