# READ ME INSTRUCTIONS

The file is divided into 2 parts.
First part explains how we **can create cluster using azure and run the preprocessing code**.
The second part deals **with installations of the important libraries** which will be required to run the questions.

## Part I - CLUSTER CREATION AND PREPROCESSING

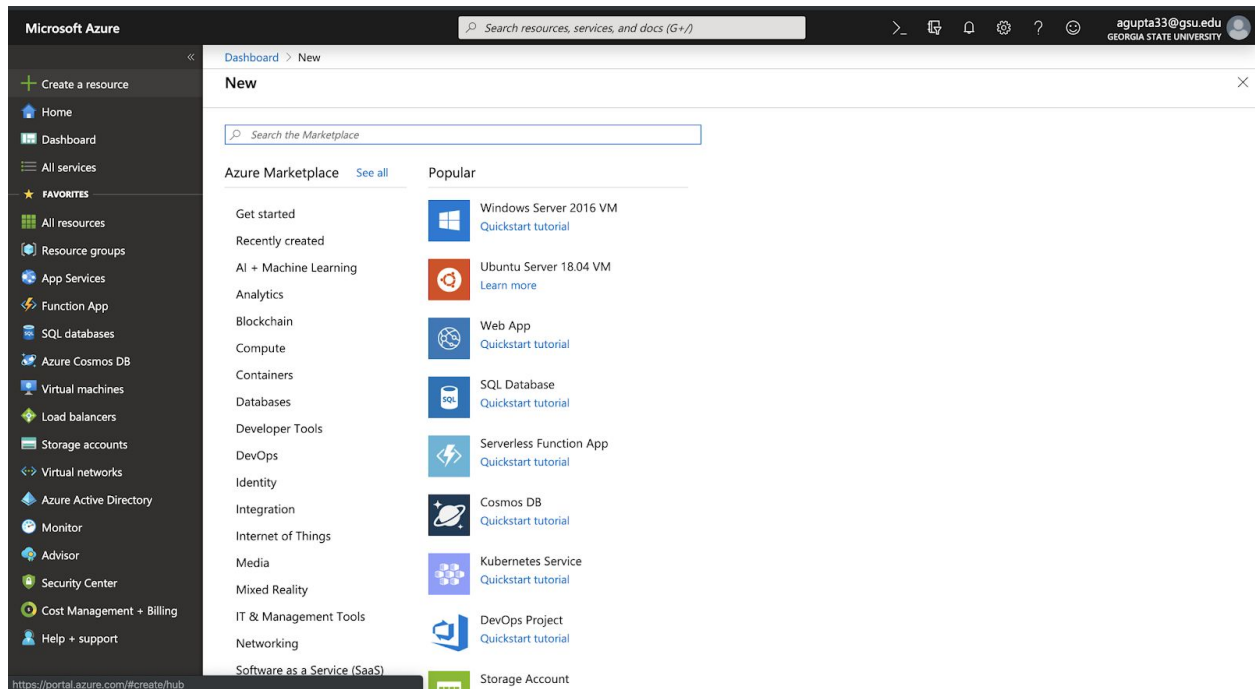Step 1: Download Microsoft Azure Storage Explorer
https://azure.microsoft.com/en-us/features/storage-explorer/
Step 2: Create and Login with your same credentials as you have in Azure Account for
HD Insights.
Step 3: login to Microsoft azure account
https://azure.microsoft.com/en-us/
Step 4: click on create resource button on top left and search for HD insights.

## Step 5: Click on HDInsight and start the creation of the cluster

**Marketplace**

My Saved List

Recently created

Service Providers

**Categories**

Get Started

AI + Machine Learning

Analytics

Blockchain

Compute

Containers

Databases

Developer Tools

DevOps

Identity

Integration

Internet of Things

HDinsights

Pricing : **All**    Operating System : **All**    Publisher : **All**

Showing All Results

**Azure HDInsight**

Microsoft

Cloud-based Big Data service. Apache Hadoop, Spark and other popular big data solutions.

**HDInsight Spark Monitoring**

Microsoft

HDInsight Spark Log Analytics, Monitoring & Alerting

**HDInsight Storm Monitoring**

Microsoft

HDInsight Storm Log Analytics, Monitoring & Alerting

**Customer Insights**

Microsoft

Deploy highly-available, infinitely-scalable applications and APIs.

**Application Insights**

Microsoft

Application performance, availability and usage information at your fingertips.

**Tidal Migrations -Premium Insights for Database**

Free trial

tidalmigrations.com

Analyze your Databases. Uncover roadblocks to cloud migration Unlock the power of Premium

**Dataiku DSS on HDInsight**

Dataiku

Dataiku DSS is an integrated and collaborative data science platform.

**CASK**

**CDAP for HDInsight**

Cask

CDAP is the first unified integration platform for big data.

Step 6: Fill **project details** as below make sure you are choosing the correct version of spark and the name of your cluster should be unique. Please keep the details as same as possible.

**Project details**

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *              Azure for Students

   Resource group *        BDPDev

Create new

**Cluster details**

Name your cluster, pick a region, and choose a cluster type and version.  Learn more

Cluster name *              mycluster

Region *                    East US

Cluster type *              **Spark**
                            Change

   Version *               Spark 2.4 (HDI 4.0)


Step7: Fill in **cluster credentials .** Give the password and do remember the password.

**Cluster credentials**

Enter new credentials that will be used to administer or access the cluster.

Cluster login username * ⓘ           admin

Cluster login password *

Confirm cluster login password *

Secure Shell (SSH) username * ⓘ      sshuser

Use cluster login password for SSH   ☑

Step8: Once you are on the Storage tab do the following

8.i ) click on create new button to create a new storage.



8.ii ) write the **container name.** and move ahead also copy the container name to some where safe as you will need it later.
Click next, There is no need to change anything in Security + networking

Step 9: Click next then on Configuration + pricing select the configuration of the head node and worker node. And then click next.

**Add application**

| Node type | Node size | Number of ... | Estimated cost/hour |
|---|---|---|---|
| Head node | D12 v2 (4 Cores, 28 GB RAM), 0.37 USD/... ⌄ | 2 | 0.75 USD |
| Worker node | D13 v2 (8 Cores, 56 GB RAM), 0.75 USD/... ⌄ | 4 ✓ | 2.99 USD |

☐ Enable autoscale
Learn more

**Total estimated cost/hour**     **3.74 USD**

---

**Create HDInsight cluster**                                             ✕

⇄ Go to classic create experience

✓ Validation succeeded.

Basics    Storage    Security + networking    Configuration + pricing    **Review + create**

Spark 2.4 (HDI 4.0)

**3.74 USD Total estimated cost/hour**
This estimate does not include subscription discounts or costs related to storage, networking, or data transfer.

**Basics**

| | |
|---|---|
| Subscription | Azure for Students |
| Resource group | BDPDev |
| Region | East US |
| Cluster name | (new) mycluster |
| Cluster type | Spark 2.4 (HDI 4.0) |
| Cluster login username | admin |
| Secure Shell (SSH) username | sshuser |
| Use cluster login password for SSH | Enabled |

**Storage**

| | |
|---|---|
| Primary storage type | Azure Storage |
| Primary storage account | bdpproj1 |
| Container | new |
| Additional Azure storage | None |

Create      « Previous      Next      Download a template for automation

Once the validation is complete click on create button. It will take around 20-25 mins

To get your storage account and access key details do the following :



Click on the storage account name on this page

| | | | | | | |
|---|---|---|---|---|---|---|
| ☐ myclusterhdistorageanitg | Storage account | Storage | BDPDev | East US | Azure for Students | ... |

Then click on Access Keys on the left side as shown below.

- Overview
- Activity log
- Access control (IAM)
- Tags
- Diagnose and solve problems
- Data transfer
- Storage Explorer (preview)

**Settings**

- Access keys
- Geo-replication
- CORS
- Configuration

You need to save  **Account Name, Key** and the container name you have given.

Storage account name

    myclusterhdistorageanitg

**key1** ↺

Key

    o2arJVsYvYtDjiUUhEq9euuKwnrGAEao+818hw54rU8g7/qT0vfHcTv+WOxiV7d6c9q2stOqeJusD1RmNTL/oA==

Connection string

**AZURE STORAGE EXPLORER**

Step 10: After the cluster is created go to Azure Storage Explorer and refresh you will see your storage account as below :

Step 11: Then in your container first create a folder and then upload the **october months data**. from your local machine to the storage and upload **bash.sh** from the zipped folder

**CLUSTER DASHBOARD**

Step 12: Once the cluster is setup go-to resource and click on Ambari home

## Step 13: Then go to Spark 2



## STEP 14: Click on config and change the below-mentioned parameters

Step 15: You need to update 4 parameters so your kernel doesn't timeout while the code is running

Step 15. i) Update livy.server.session.timeout from 36000000 to 180000000: make sure you put right number of zeros. (we did this mistake every single time we created the clusters)
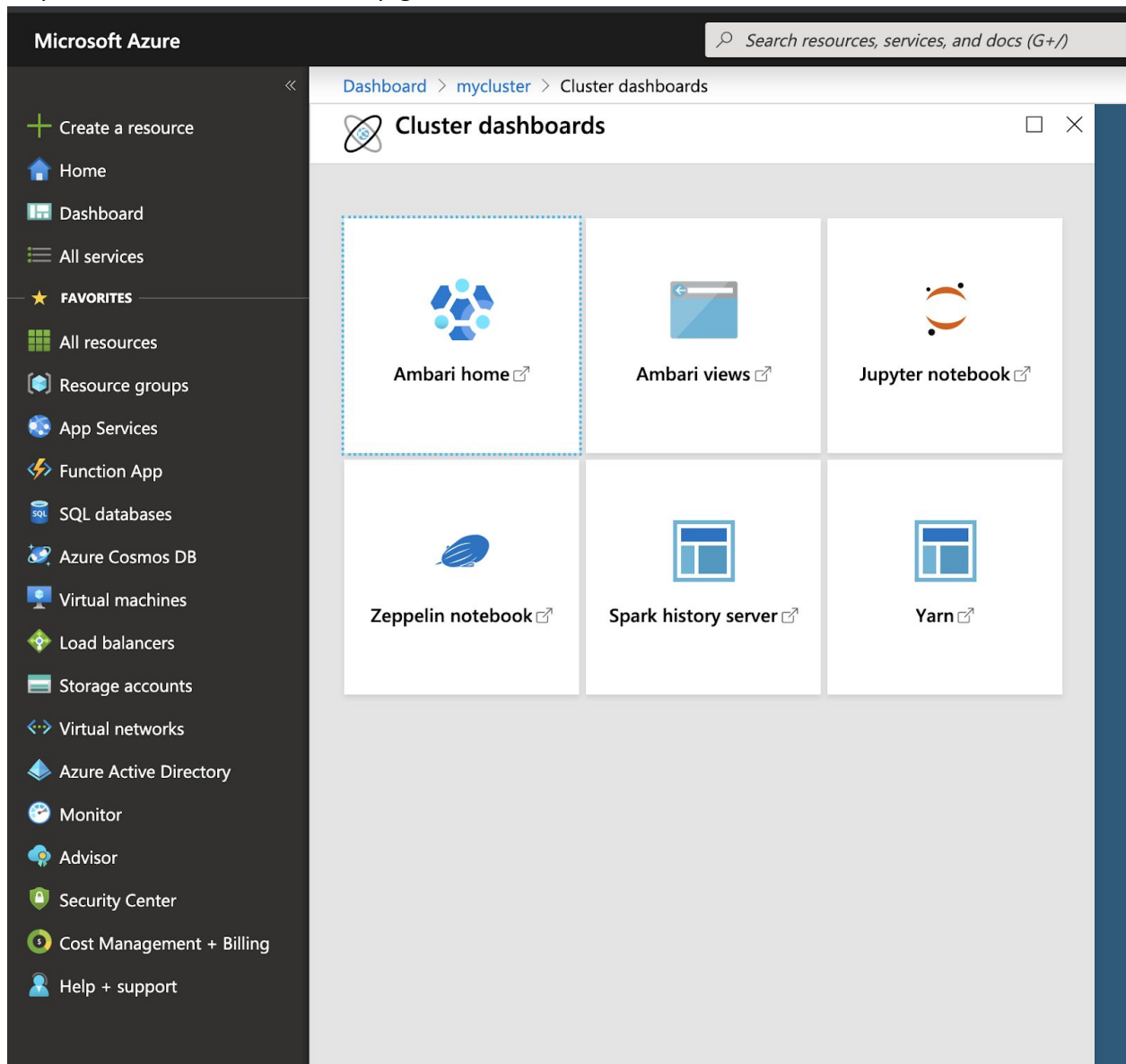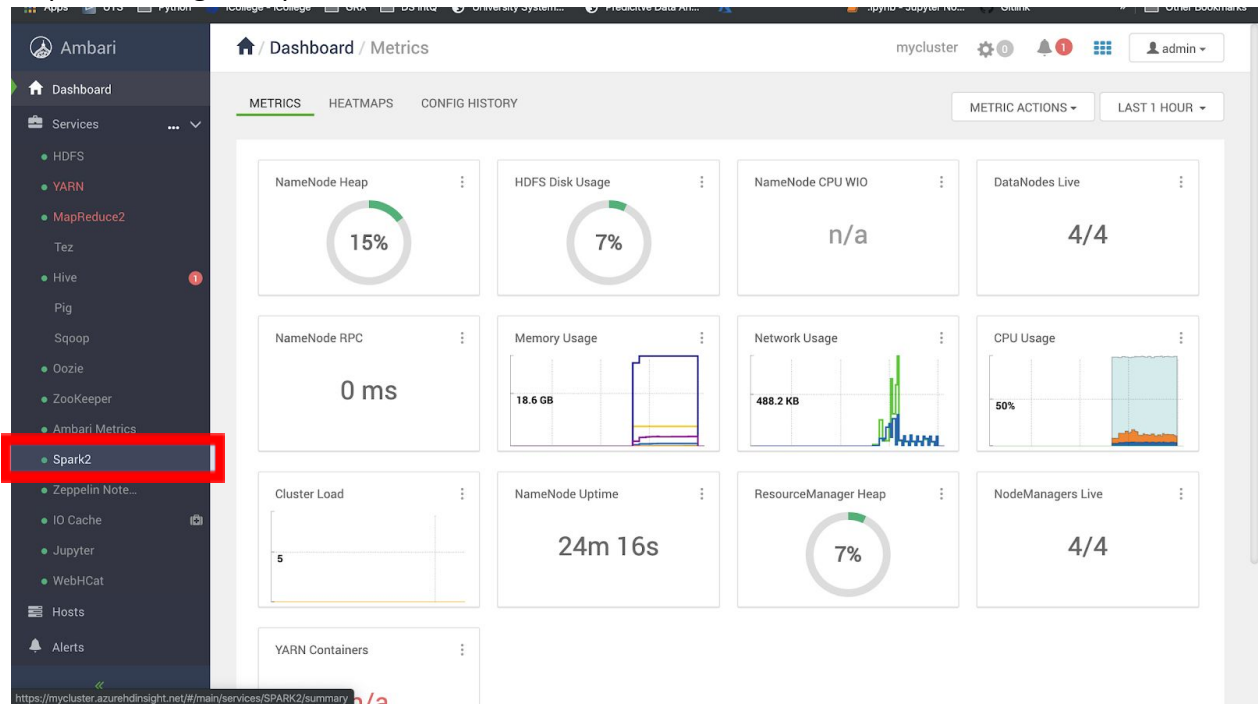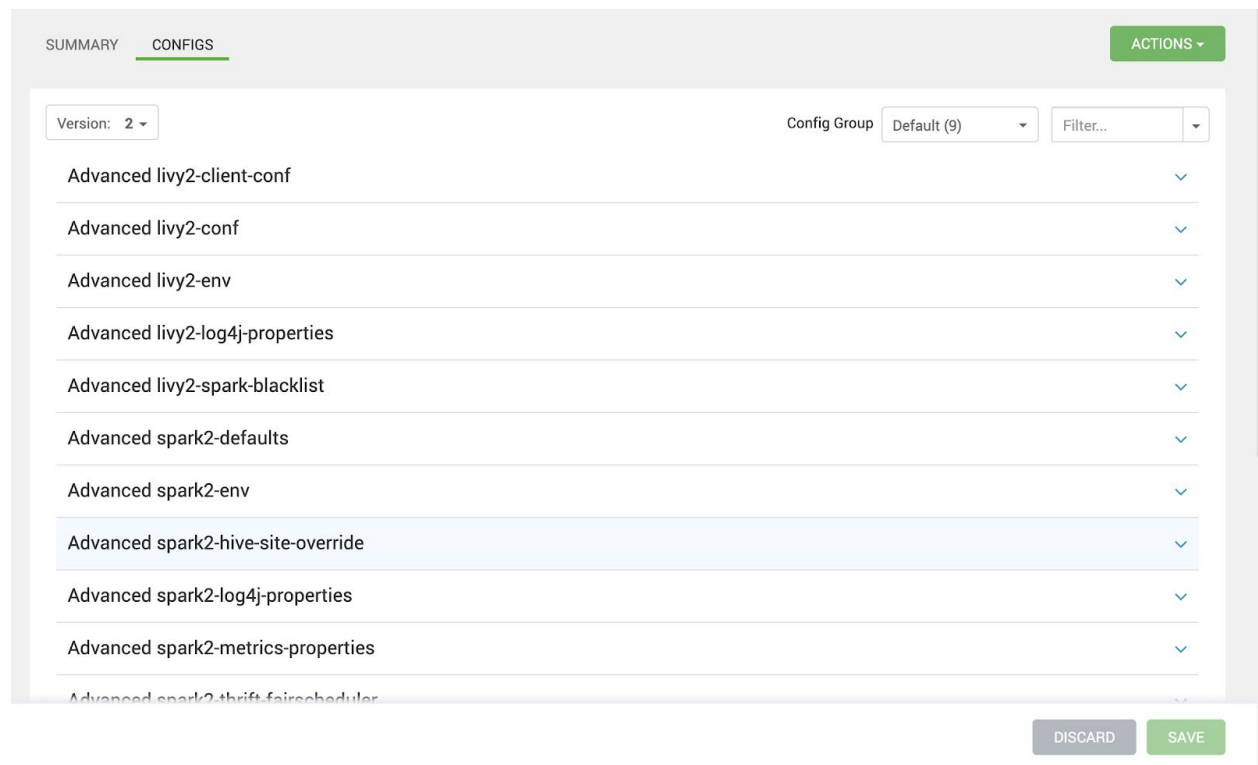
## Advanced livy2-conf

| | |
|---|---|
| livy.environment | production |
| livy.impersonation.enabled | true |
| livy.repl.enableHiveContext | true |
| livy.server.access-control.enabled | true |
| livy.server.csrf_protection.enabled | true |
| livy.server.port | 8998 |
| livy.server.recovery.mode | recovery |
| livy.server.recovery.state-store | zookeeper |
| livy.server.recovery.state-store.url | zk1-bdppro.2h5ffrof4nveneybdi1ajpa2qh.bx.internal.cloudapp.net:2181,zk2-bdppro.2h5ffrof |
| livy.server.session.timeout | 18000000 |
| livy.spark.master | yarn-cluster |

Step 15.ii) Change livy.server.yarn.app-loolup-timeout from 2m to 10m as shown below

## Custom livy2-conf

| | |
|---|---|
| livy.server.session.state-retain.sec | 3600000 |
| livy.server.yarn.app-lookup-timeout | 10m |

Add Property ...

Step 15.iii) then in Custome spark2-daefaults you need to add a property

## Add Property ✕

| Type | spark2-defaults.xml | 🏷 🏷 |
|---|---|---|

**Key**    spark.sql.broadcastTimeout

**Value**
```
6000
```

**Property Type**
```
PASSWORD
USER
GROUP
TEXT
```

CANCEL    **ADD**

Step 15 iv) Add property spark.driver.memory as 32g

## Add Property ✕

| Type | spark2-defaults.xml | 🏷 🏷 |
|---|---|---|

**Key**    spark.driver.memory

**Value**
```
32g
```
Ⓖ

**Property Type**
```
PASSWORD
USER
GROUP
TEXT
```

CANCEL    **ADD**

Step 16 i) Go to jupyter configuration and change below 2 parameters

**Add Property**                                                                    ✕

| Type | jupyter-site.xml |
| Key | MappingKernelM0nager.cull_idletimeoutInt |
| Value | 0 |

Property Type
PASSWORD
USER
GROUP
TEXT

CANCEL   ADD

---

**Add Property**                                                                    ✕

| Type | jupyter-site.xml |
| Key | NotebookApp.shutdown_no_ActivitytimeoutInt |
| Value | 0 |

Property Type
PASSWORD
USER
GROUP
TEXT

CANCEL   ADD

MappingKernelManager.cull_idle_timeoutInt - 0
NotebookApp.shutdown_no_activity_timeoutInt - 0

After this save and restart Spark and Jupyter both from the actions on the top right corner of the window.



Step 17) having done that add goto script actions in order to run your imports

Step ) Get bash script URI from azure storage account explorer by right clicking on your file and clicking properties there you can find the URI

Step 18) Then got to Submit new script option. Select script type as Custom. Add the bash.sh raw file URL in `Bash script URI` or any public URL of bash file containing the commands to install the libraries.

In our case
    /usr/bin/anaconda/envs/py35/bin/pip install azure
    /usr/bin/anaconda/envs/py35/bin/pip install pandas==0.19.2

Then select the node types required. (head and worker in our case).
Optional: Select 'Persist this script action when new nodes are added to the cluster' option if you want to install these libraries when new worker nodes are added.

Upon clicking create. Our packages will be installed on all the nodes and we are ready to run our preprocessing code.

Step 19) Then go to Jupyter notebook under `Cluster management interfaces` in the overview section in our cluster home page. Then go ahead and run `SparkAllWords.ipynb`. This will generate the tweets and retweets files inside the Azure storage account.

**SparkAllWords.ipynb file changes that you need to do as per the storage Account.**

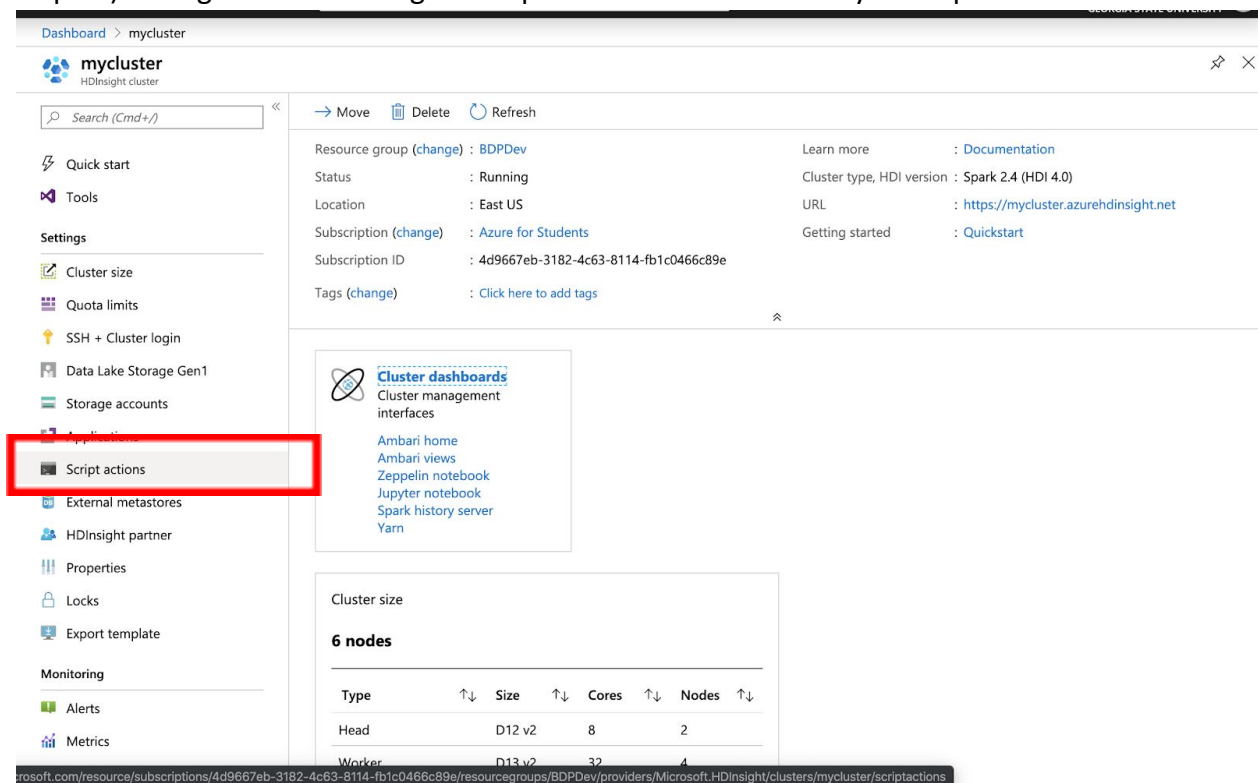You would need to add your paths and storage account details in the file as we have saved them above.
You need to update below details as per you saved above.

```
-----------------------------------------------------------------------------------------------------------------------------
accountName = "bdpproj1"
accountKey=
"qDRRza9BkroEhjC8rOJcEZ70WU9fFbXOylwcjtWT6mjGd1AVZ9O2dgd+s0+m5vzL39Fix+TvxDwSABtvadjQEg==
"
  containerName = "new"
-----------------------------------------------------------------------------------------------------------------------------
```

You need to set the path of the folder that you created in Step 11()

```
-----------------------------------------------------------------------------------------------------------------------------
blobService.create_blob_from_text('new/anit',file_name, output)
-----------------------------------------------------------------------------------------------------------------------------
```

You need to add path to your data set here and add * if you want an entire folder to be read.

```
-----------------------------------------------------------------------------------------------------------------------------
dataset_path = "wasbs:///anit/10/01/*/*"
-----------------------------------------------------------------------------------------------------------------------------
```

Once done you should download the generated CSVs from the azure storage explorer and put it where you keep the final code for all the codes.

# PART II - Question Specific installations

Run the below commands on the local spark, not on Azure spark since all the rest of the code is to be run on local spark.

Initial spark setup:
1. pip install pyspark

(we downloaded the spark-2.4.4-bin-hadoop2.7 version)
2. Install Java 8
3. In MAC OS, edit .bash_profile with following OR In Windows environment variables add:
   ```
   export JAVA_HOME=$(/usr/libexec/java_home)
   export SPARK_HOME=~/spark-2.4.4-bin-hadoop2.7
   export PATH=$SPARK_HOME/bin:$PATH
   export PYSPARK_PYTHON=python3
   ```

4. Start pyspark by running below command in terminal -
   pyspark

Question 1
------------------------------------------------------------------------
pip install pygal
------------------------------------------------------------------------
Question 2
------------------------------------------------------------------------
pip install geopy
pip install plotly
------------------------------------------------------------------------

Question 3
------------------------------------------------------------------------
pip install nltk
nltk.download('stopwords')
nltk.download('wordnet')
------------------------------------------------------------------------
Question 4
------------------------------------------------------------------------
pip install pyldavis
------------------------------------------------------------------------

Question 5 - Introduction
------------------------------------------------------------------------
pip install pygal

----------------------------------------------------------------------

## Question 6 - Time Series Analysis
----------------------------------------------------------------------
pip install statsmodels

pip install datetime

pip install ipython

pip install ipywidgets

pip install strings

pip install seaborn

----------------------------------------------------------------------