# Spark Movie Query Application Tutorial

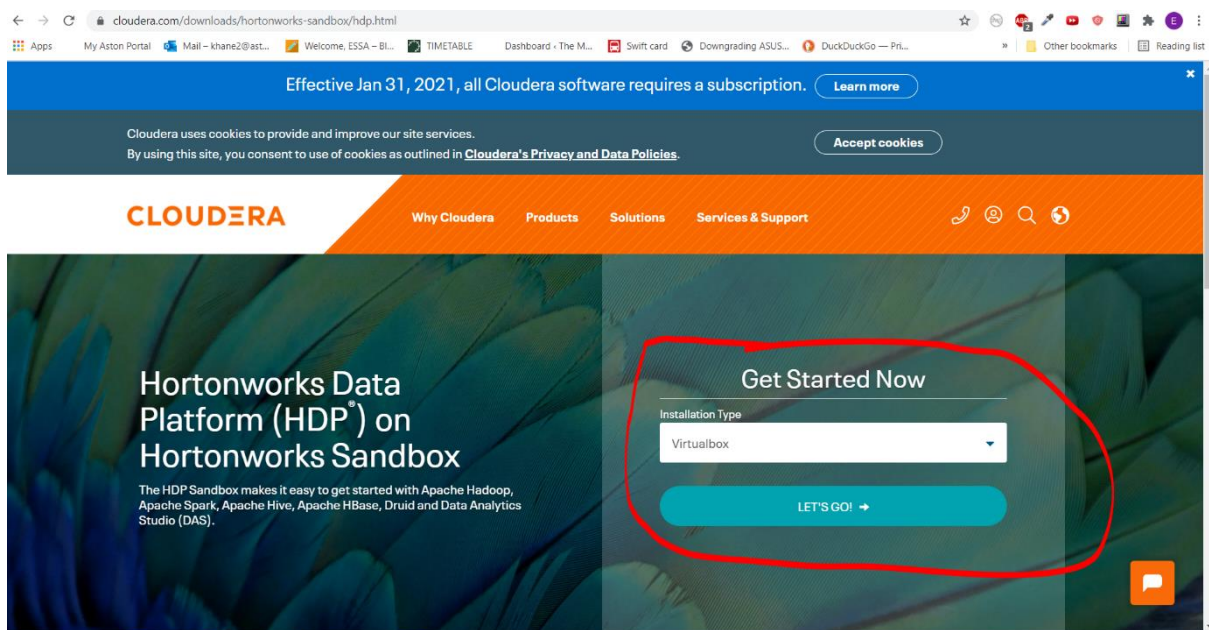| | |
|---|---|
| Student Name | Essa Umar Khan |
| University | Aston University |
| Module Code | CS3800 |
| Module Name | Advanced Data Base Systems |
| Student Number | 170077653 |
| Contact Email | khane2@aston.ac.uk |

# Contents

# 1.    Pre-requisites

1) Computer running Windows 10 with 16 GB ram
2) Internet connection
3) Oracle Virtual Box (application was tested on v6.1)
4) Download the entire project source code files from either:
    a. The GitHub repository
    b. The .ZIP submission contains the files from the GitHub repository.

# 2.    Hortonworks Sandbox 2.6.5
## 2.1    Download

1) Head over to https://www.cloudera.com/downloads/hortonworks-sandbox/hdp.html [1] and select 'VirtualBox' as the instillation type.



---

[1] Alternatively, whilst logged into an Aston outlook email account you can download the sandbox from here: https://liveastonac-my.sharepoint.com/:u:/g/personal/khane2_aston_ac_uk/EXMFaW8Clw5LpZ0mMHQN_6YBavbLoT92 V4ScCMpU-x4kCg?e=z9bRdH

2) A pop-up will appear and to access the download link you have to complete the form on this pop-up, once you have filled in the information accept all the policies and press submit.



3) Click on 2.6.5 link to start download. It is a very large .ova file so it will take some time to download.



4) Once the download is completed, double click the downloaded .ova file:

4

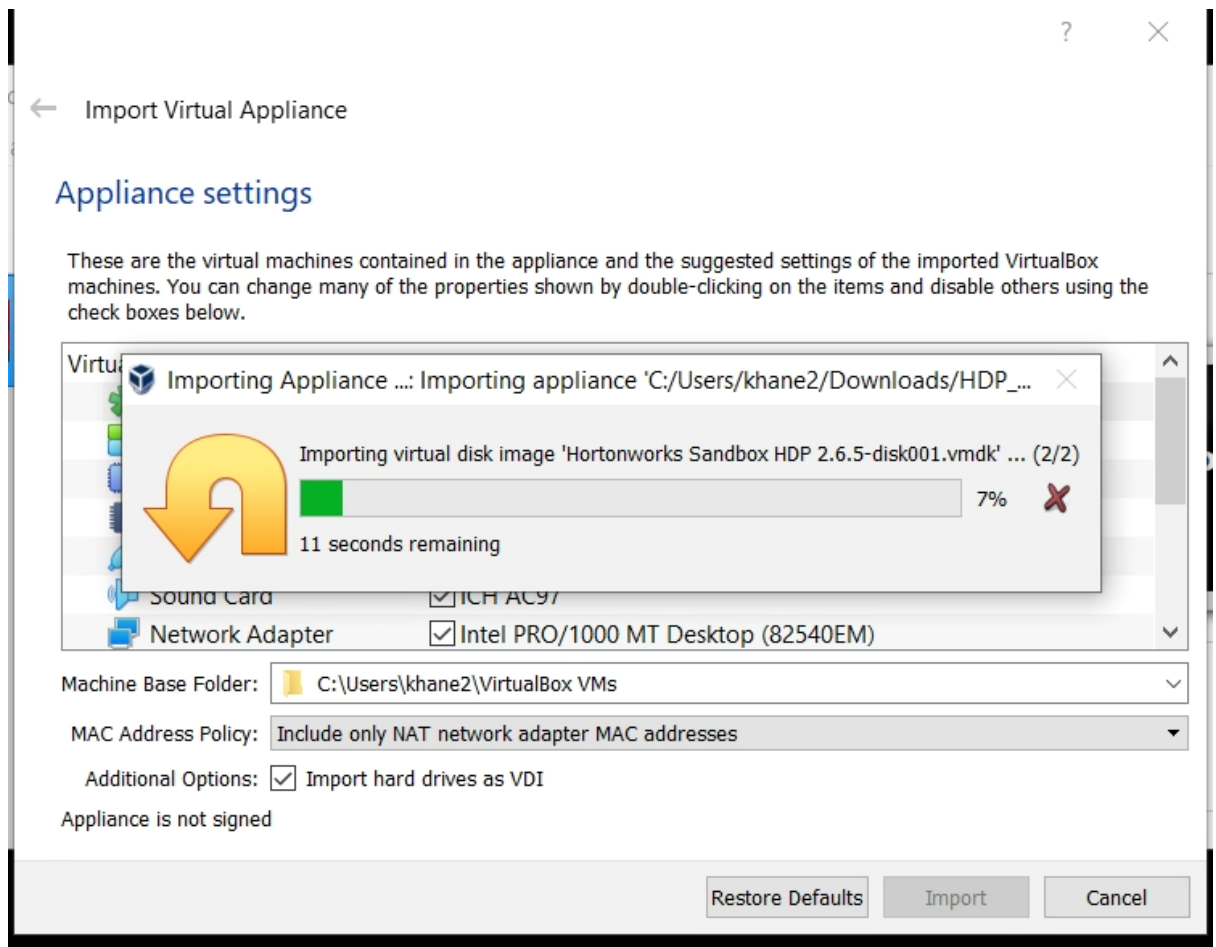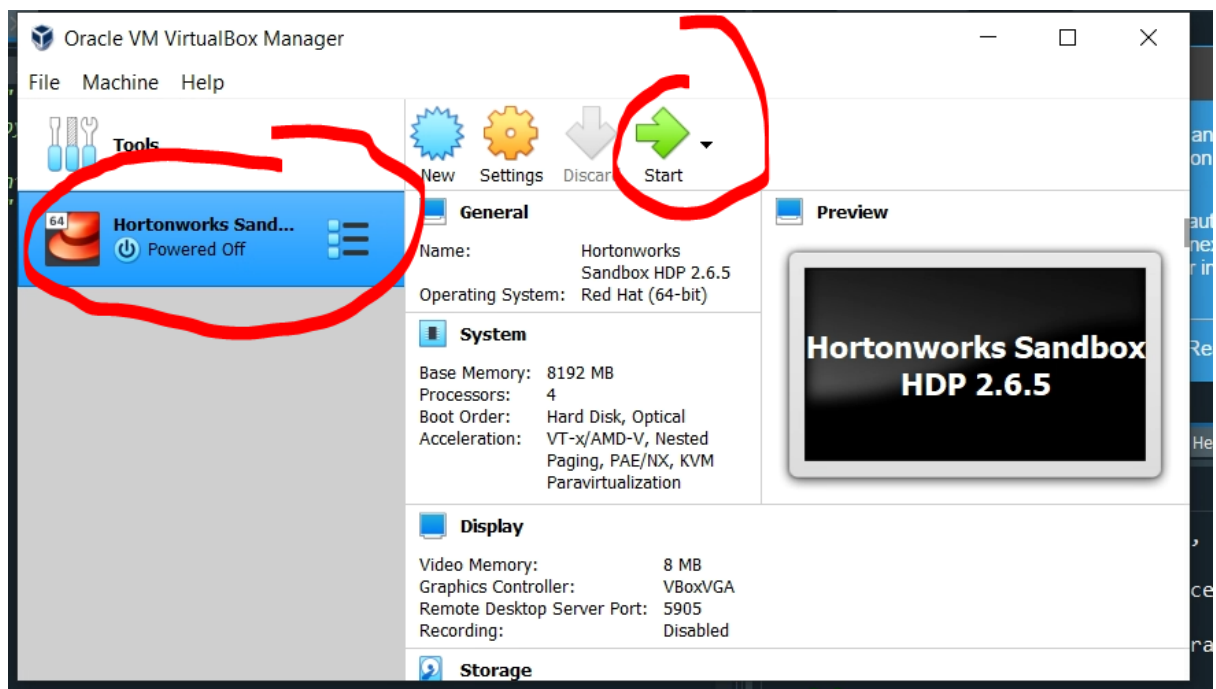5) *'Oracle Virtual Box'* software will open with a pop-up like this, just click the import button with the default settings.



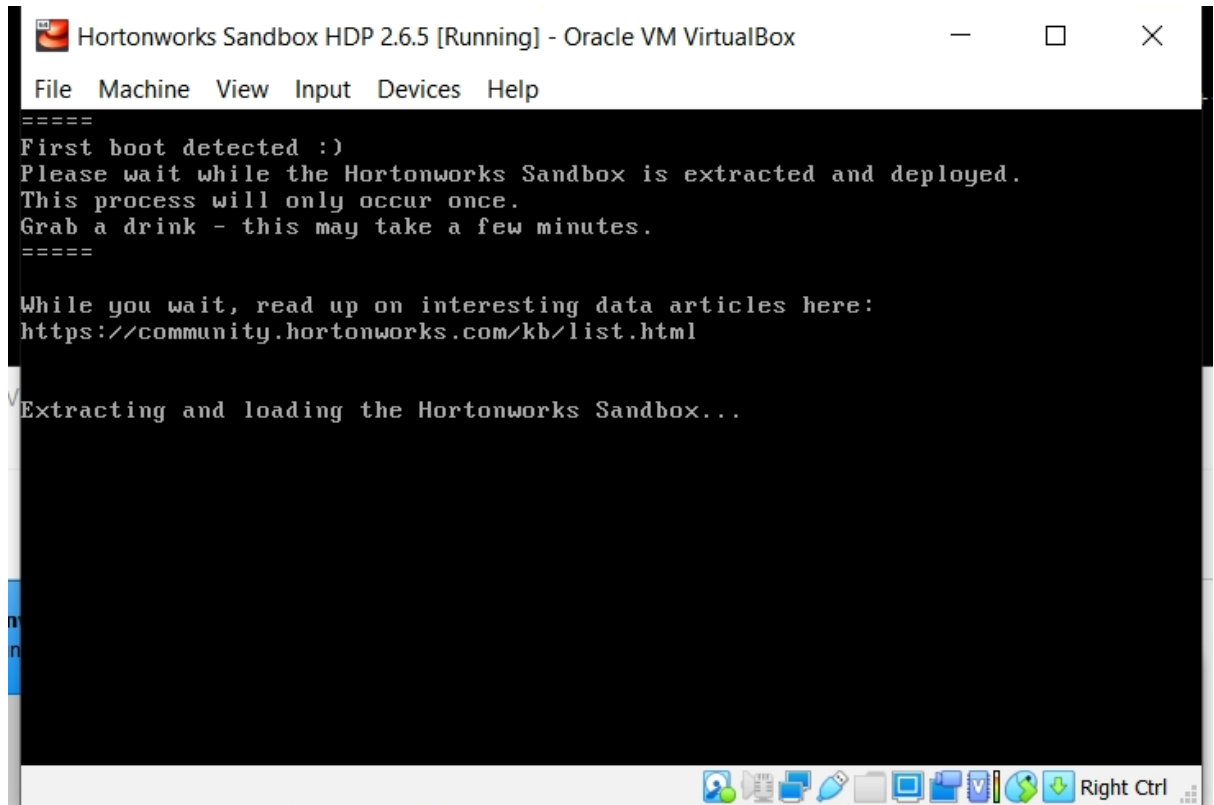Once you have pressed import then the image will be imported into virtual box. This will take some time to be copied over:
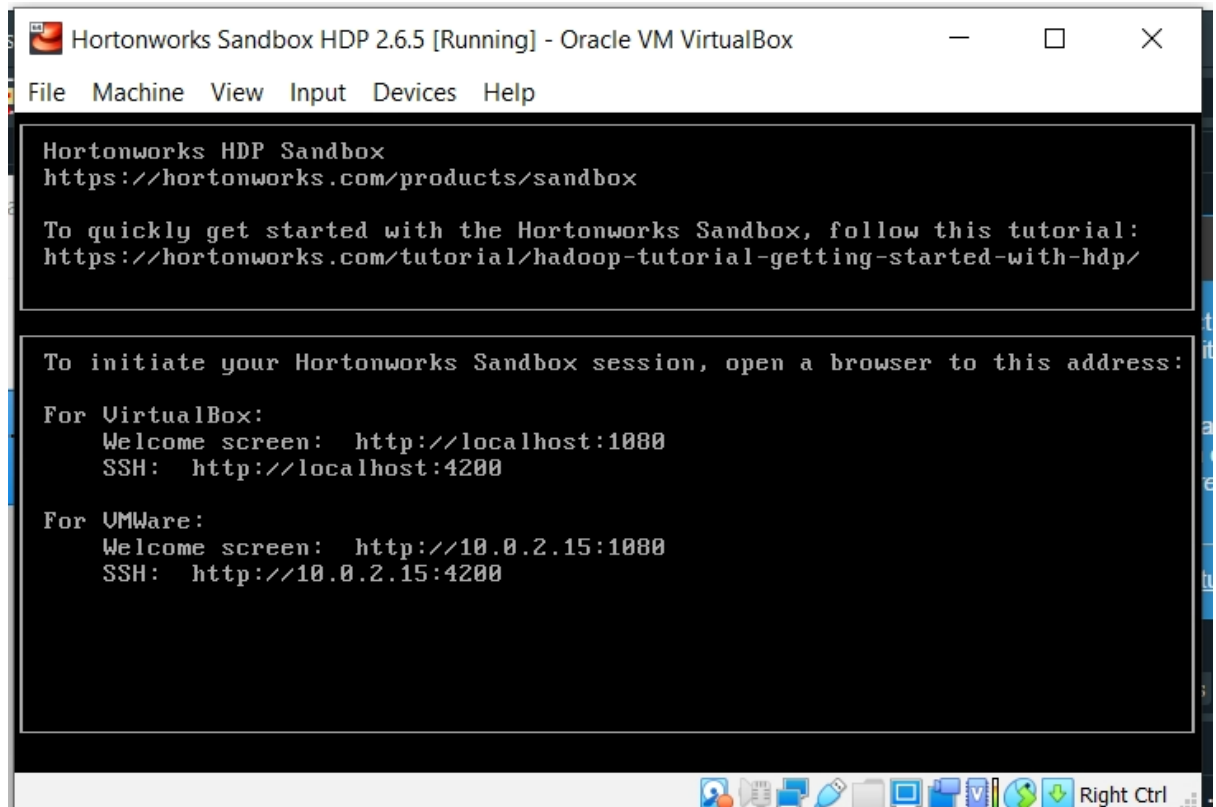
6) Once the import has completed, select the image and press run virtual box's software main window:

7) There is a 'first time launch set-up' boot up process that will now occur, let this complete and do not close the window. You don't have to do anything but wait. This usually takes around 15 minutes.

8) Once the set-up in step 7 has completed you will see a screen like the one below.
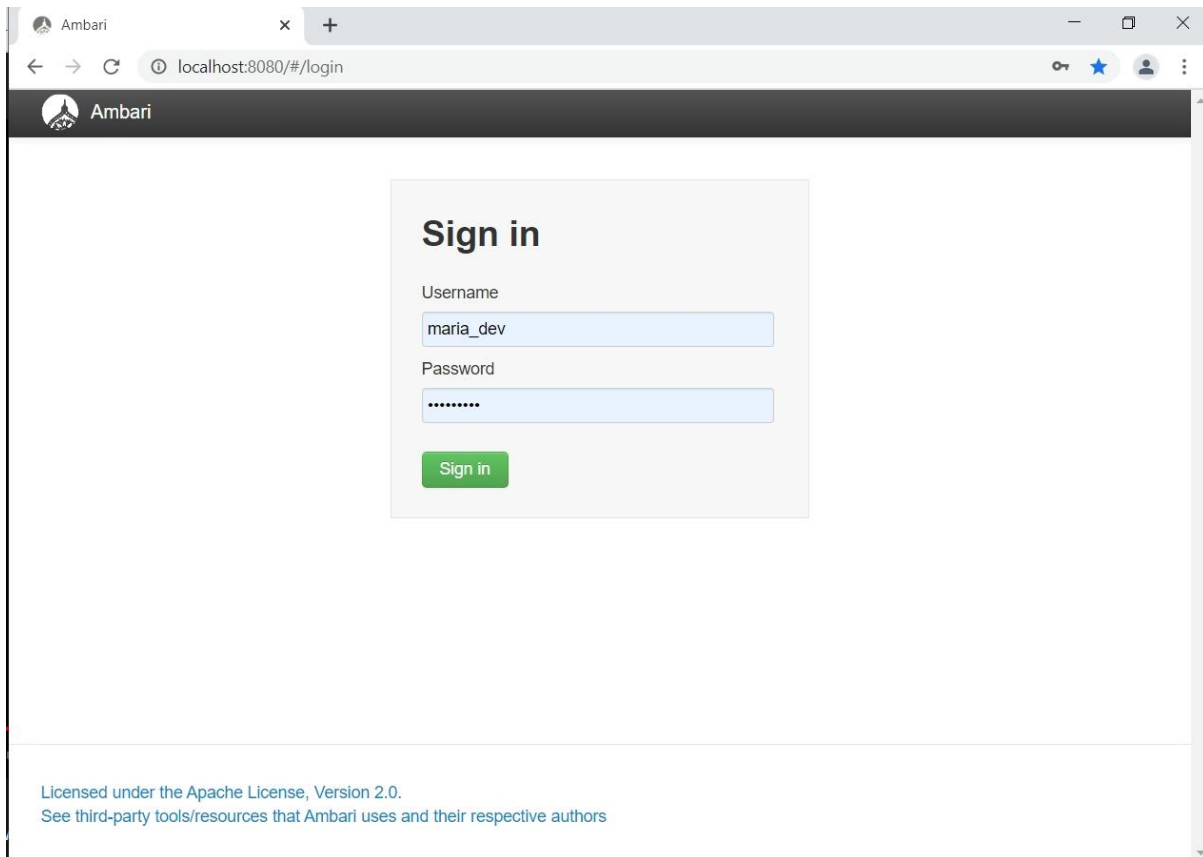


As soon as this screen shows all the technologies are loading to be (Hadoop, spark etc.) ready for use – even though it does not indicate this on the above screen. Until this process is completed do not proceed to section 2.3. This process usually takes 15 minutes max to complete, so you can wait for 15 minutes and then proceed to the next section of this tutorial. Alternatively, as this process can finish a lot faster than 15 minutes it is recommended that you visually track the progress of this process please see the forthcoming optional section 2.2: *'Optional: How to visually track apache Ambari services start-up progress.'*
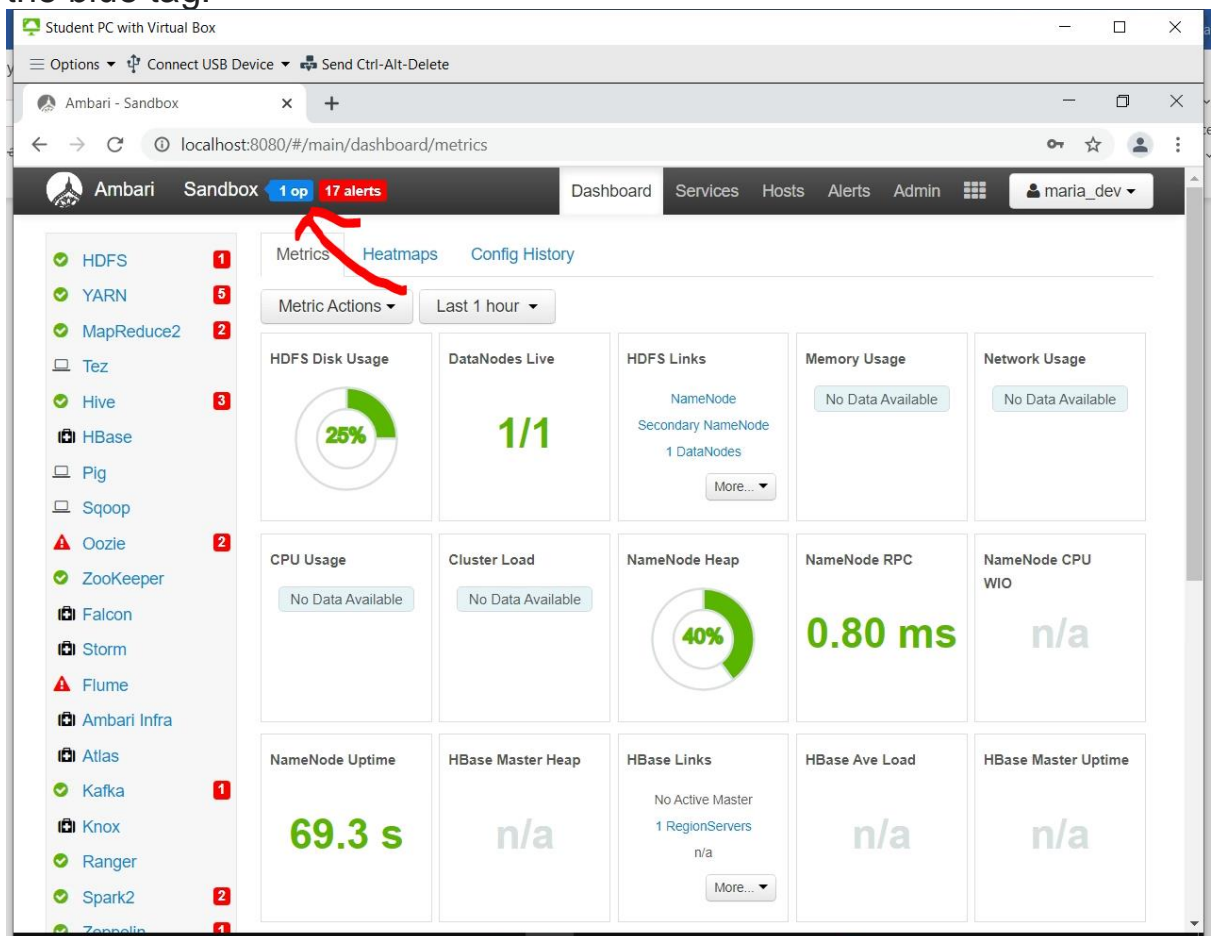
## 2.2   Optional: How to visually track Apache Ambari services start-up progress.

You have been directed to this optional sub-tutorial, and follows from section 2.1 step 8 of this document.

1) Open an internet browser (chrome is recommended) and go to this link: http://localhost:8080/#/login and login with the credentials:
    a. Username: maria_dev
    b. Password: maria_dev

2) Next, you will see something similar to the screen-shot below, click the blue tag:



3) A popup similar as shown in the screenshot below will open. Once, this progress bar is completed you can proceed from section 2.3.
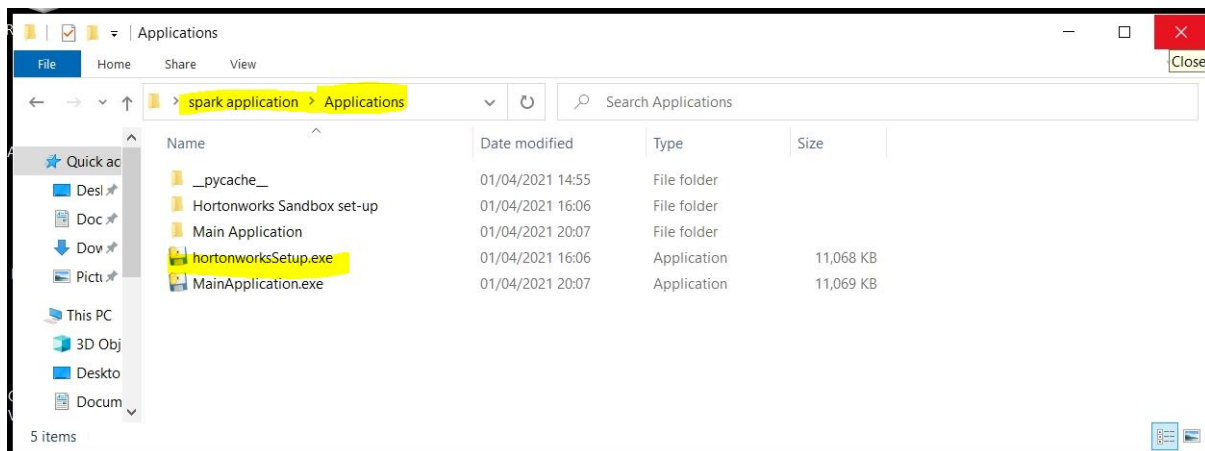
## 2.3   Set-up

The tutorial from this point onwards can also be followed by watching this youtube video that gives a screen recorded demonstration: https://www.youtube.com/watch?v=Dm_yH_McmbI&ab_channel=EssaKhan
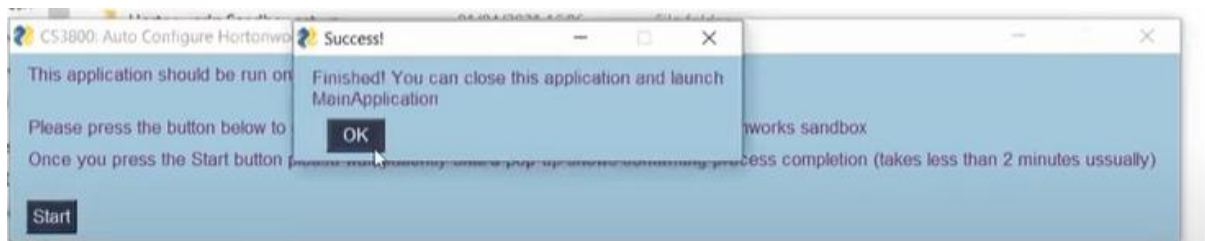
Alternatively you may follow these instructions:

1) Open source-code folder location 'spark-application/applications' and launch the application 'hortonworksSetup.exe'
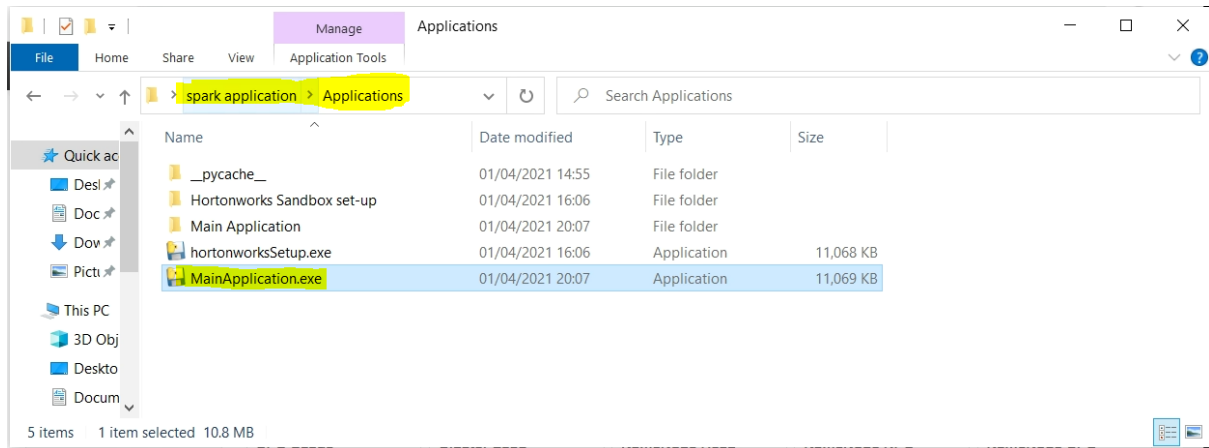


2) On the GUI that opens up, follow the on-screen instructions. You must press the 'start' button once and let the process complete. Note: hortonworksSetup.exe must only only run once per each Hortonworks sandbox virtual machine image else you will get errors.

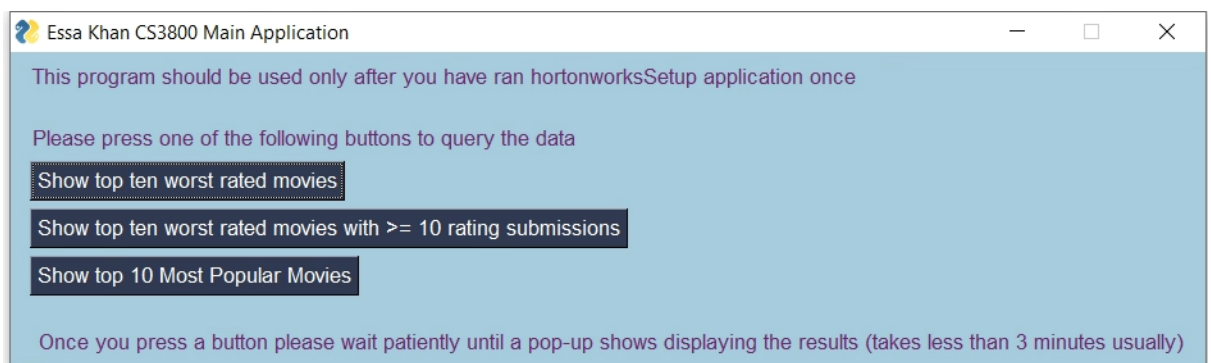   Once process is complete, exit/close this application.

## 3.    Running the Main Application

1) Open source-code folder location 'spark-application/applications' and launch the application 'MainApplication.exe'



2) A GUI like the one below will open, follow the onscreen instructions to see how to operate the application.



You may press each button one by one to query the underlying dataset. The results are shown in a pop-up window. Once one request of results by a button press are output to you then you may proceed to press another button.

## 4.    Optional: Run python .py scripts instead of .exe

This may be helpful alternative if the .exe do not work for some reason. After opening the GUIs this way please follow the respective aforementioned GUIs walkthrough in the previous section.

This has been tested on python 3.7.6

You need to install the relevant packages:

pip install pysimplegui

pip install spur

1) To open the GUIs through their scripts directly instead of the provided .exe files:
   a. run the script 'spark-application\Applications\Hortonworks Sandbox set-up\hortonworksSetup.py'
   b. run the script 'spark-application\Applications\Main Application\MainApplication.py'