# upGrad

## Spark Setup on EMR Cluster

In this document, you will learn how to launch an EMR cluster with Spark installed on it. You will also see how to run a Jypyter notebook on the EMR cluster.

## Steps to launch an EMR Cluster with spark installed on it

1. First, go to the EMR services section and click on the "**Create Cluster"** button.



2. You will be redirected to the below window, here click on the "**Go to advanced options"** button.

3. Under the **Release** section and select the checkbox **Spark.**



4. After selecting Spark, you need to scroll down and press "**Next".** Here you don't need to add any step in **"Add step"** section.

5. **Important -** Now change the **Master** and **Core** to **m4.large** and press **Next**. Ensure you are doing this else you would incur huge cost with the default settings.



6. Name your cluster as you can see we have named it "**Demo_Cluster**". After naming the cluster, press **Next**.

7. In this step select an existing key pair and click on **Create cluster** button.You can proceed further with the existing EC2 key pair or any other keypair for which the corresponding PPK file is available. In case you don't have a PPK file for any of the existing key pairs then you need to create a new key pair, you can refer to previous modules where creating a new pair is demonstrated to you Finally hit the **"Create cluster"** button.



8. This will take some time to create the cluster and then you will find cluster status as **Waiting**.

# Running Jupyter notebook on EMR cluster

9. Now you need to go to the **"Notebook"** section as shown in the image below.



10. Hit on the **Create notebook** button.



11. Name the notebook as per your desire. As you can see in our case the name is "**Demo_notebook**"

12. Then click on **choose** in the cluster section.



13. Now, select the existing cluster as shown in the image below.

14. Once you choose the cluster, then you will have to hit the "**Create notebook**" button at the bottom of the window. It will take some time to create a notebook.



15. Now click on the "**Open in Jupyter**" button. This will land you to the jupyter notebook tab.

16. You can upload your own notebook from the "**upload**" button on the top right part of the window.



.

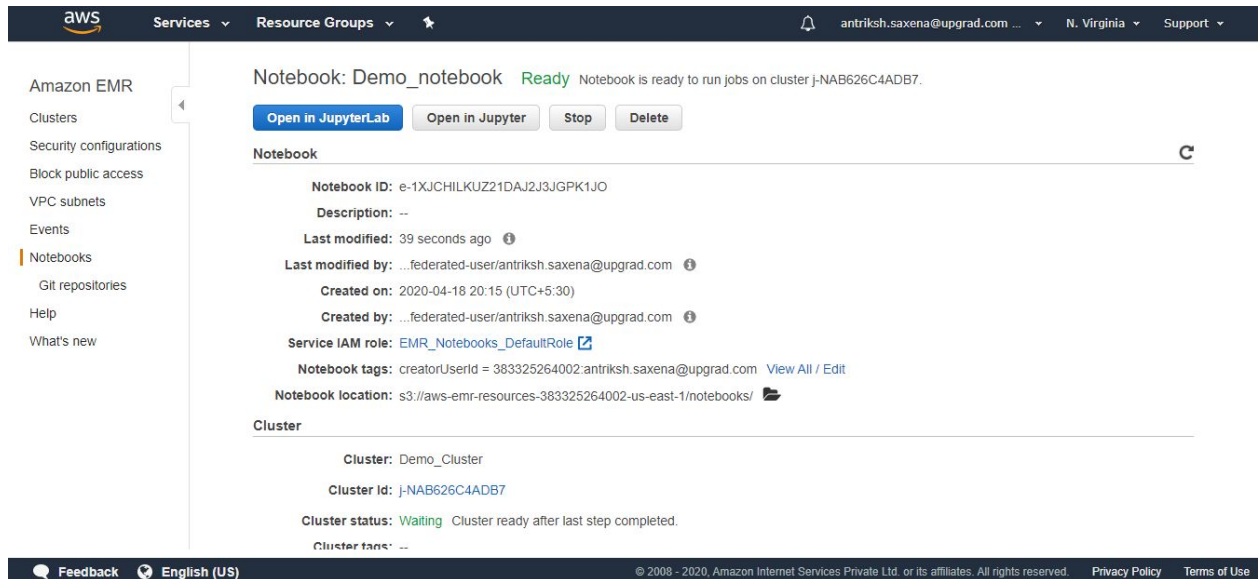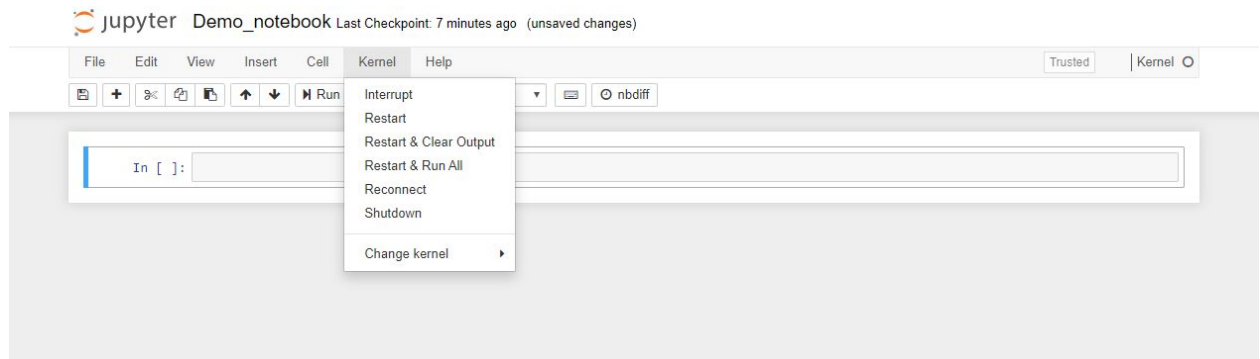17. There is also an option to create a new notebook with the desired kernel .



18. There is an option in the notebook to change the kernel also. From which you can select the desired kernel.

In this way you can run the Spark application on the EMR cluster.