# Getting started with Databricks and Spark

Ladle Patel

# Step 1:

Go to the below URL and click on **GET STARTED** under Community Edition.

**NOTE: Please don't create FREE TRIAL account**

https://databricks.com/try-databricks

### DATABRICKS PLATFORM – FREE TRIAL

For businesses looking for a zero-management cloud platform built around Apache Spark

- Unlimited clusters that can scale to any size
- Job scheduler to execute jobs for production pipelines
- Fully interactive notebook with collaboration, dashboards, REST APIs
- Advanced security, role-based access controls, and audit logs
- Single Sign On support
- Integration with BI tools such as Tableau, Qlik, and Looker
- 14-day full feature trial (excludes cloud charges)

### COMMUNITY EDITION

For students and educational institutions just getting started with Apache Spark

- Single cluster limited to 6GB and no worker nodes
- Basic notebook without collaboration
- Limited to 3 max users
- Public environment to share your work

**GET STARTED**

**GET STARTED**

# Step 2:

Fill up the details and click on the Sign Up button.

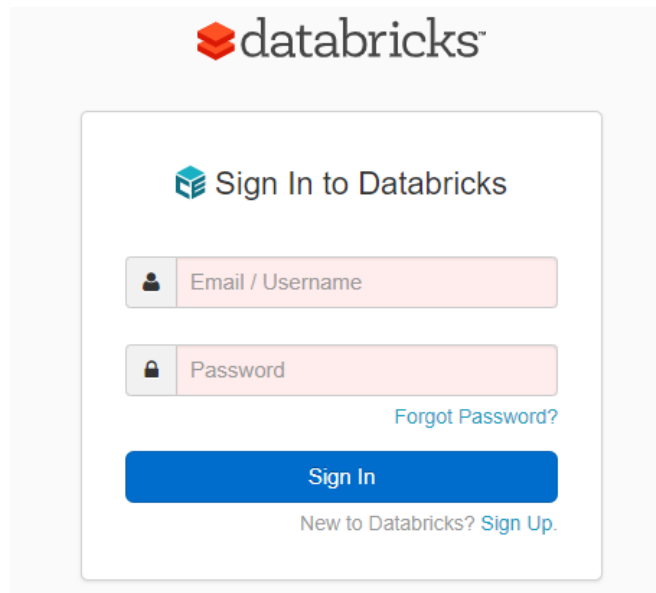NOTE: If you don't have (or) work in a company, then type NA in the Company Name field.

# Step 3:

You will receive a link through an email to verify your account,click on the link to complete the account set up.Use below link to login to your account.
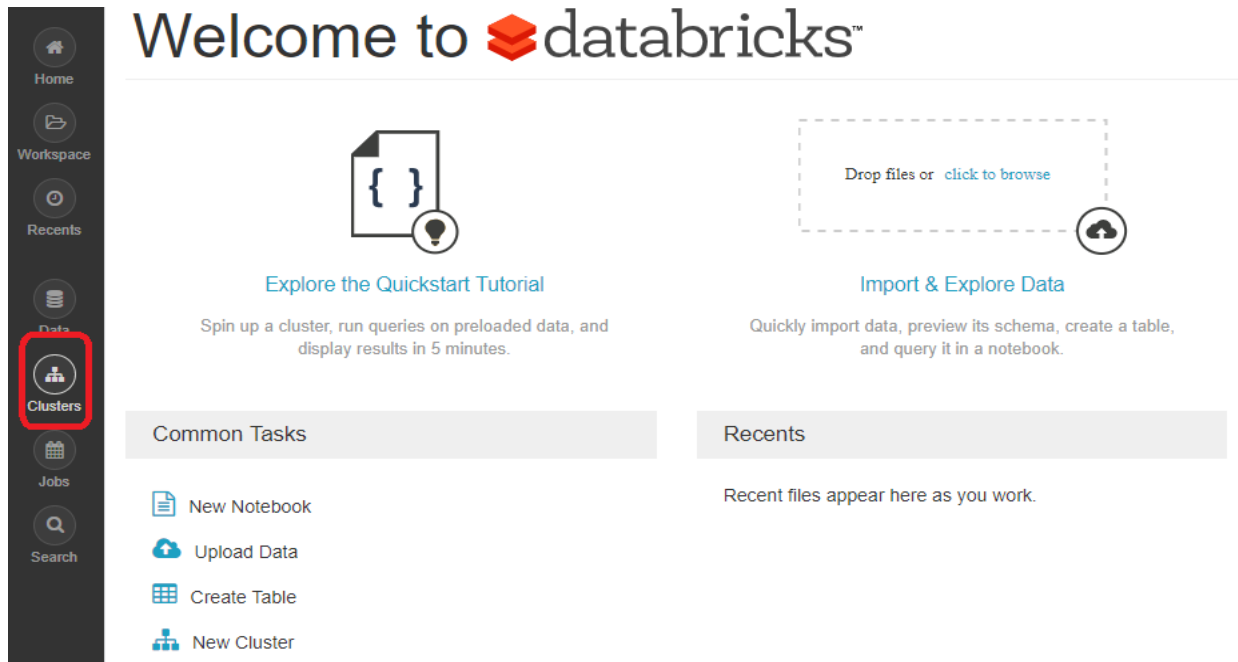https://community.cloud.databricks.com/login.html

# Step 4:

Creating New Cluster.

# Step 4: Cont..

# Step 5:

Enter cluster name

# Step 6:

Upload data from your laptop to the cluster.

# Step 6: Cont..

Upload data from your laptop to the cluster.

# Step 7:

Copy the path.

# Step 8:

Create new notebook

# Step 9:

Provide the name of the notebook and chose the language.

# Step 10:

Read the data in spark

# Further Reading

https://docs.databricks.com/spark/latest/training/index.html

https://spark.apache.org/

https://github.com/databricks/Spark-The-Definitive-Guide

https://databricks.com/sparkaisummit/north-america/sessions?eventName=Summit%202018

https://databricks.com/sparkaisummit/north-america/sessions

https://www.youtube.com/user/TheApacheSpark/playlists