

## **Databricks:**

Databricks is a cloud-based data analytics platform that provides an interactive workspace for data engineers, data scientists, and analysts to collaborate on big data and machine learning projects. It is built on **Apache Spark** and integrates well with **Azure, AWS, and Google Cloud**.

### **Key Features of Databricks**

1. **Unified Data Platform** – Combines data engineering, data science, and machine learning.
2. **Apache Spark Optimization** – Provides a managed Spark environment with auto-scaling and performance improvements.
3. **Multi-Cloud Support** – Available on **Azure (Azure Databricks), AWS (Databricks on AWS), and GCP**.
4. **Notebooks** – Supports Python, Scala, SQL, and R for interactive data processing.
5. **Delta Lake** – An optimized storage layer that enhances **data reliability and performance**.
6. **Security & Compliance** – Provides enterprise-grade security and compliance support.
7. **Collaborative Workspace** – Allows multiple users to work together in shared notebooks.

Databricks uses its own filesystem for data processing operation, called DBFS (Databricks File system). It is a Linux based file system and we will use Unix commands for various data processing operations.

### **Databricks File System (DBFS)**

DBFS is a **distributed file system** built into Databricks. It allows users to interact with files and directories in **Databricks clusters** as if they were a regular file system.

### **Key Features of DBFS**

- **Mountable Storage** – Can mount cloud storage (Azure ADLS, AWS S3, Google Cloud Storage).
- **Hierarchical File Structure** – Files are organized in directories.
- **Supports Different File Formats** – CSV, JSON, Parquet, Delta, etc.
- **Access Control** – Permissions for files and directories.

- In Databricks, we will use dbutils (databricks utilities) to interact DBFS and for data processing operations.

### Imp points#

```

07:13 PM (1s) 1
...
ctrl+enter --> cursor will be in current cell and all statements will be executed in current cell
shift+enter --> current cell all statements will be executed and cursor will be moved to next cell
ctrl+shift+enter --> only highlighted text will be executed
...

Out[1]: '\nctrl+enter --> cursor will be in current cell and all statements will be executed in current cell\nshift+enter --> current cell all statements will be executed and cursor will be moved to next cell\nctrl+shift+enter --> only highlighted text will be executed\n'
```

dbutils.fs.help() → This command will display the different types of supported commands in DBFS.

```

2 minutes ago (<1s) 1 Python
#DBFS -- DataBricks File System --> It is a linux based file system and we will use linux commands for activities.
dbutils.fs.help()

dbutils.fs provides utilities for working with FileSystems. Most methods in this package can take either a DBFS path (e.g., "/foo" or "dbfs:/foo"), or another FileSystem URI. For more info about a method, use dbutils.fs.help("methodName"). In notebooks, you can also use the %fs shorthand to access DBFS. The %fs shorthand maps straightforwardly onto dbutils calls. For example, "%fs head --maxBytes=10000 /file/path" translates into "dbutils.fs.head("/file/path", maxBytes = 10000)".

mount

mount(source: String, mountPoint: String, encryptionType: String = "", owner: String = null, extraConfigs: Map = Map.empty[String, String]): boolean -> Mounts the given source directory into DBFS at the given mount point
mounts: Seq -> Displays information about what is mounted within DBFS
refreshMounts: boolean -> Forces all machines in this cluster to refresh their mount cache, ensuring they receive the most recent information
unmount(mountPoint: String): boolean -> Deletes a DBFS mount point
updateMount(source: String, mountPoint: String, encryptionType: String = "", owner: String = null, extraConfigs: Map = Map.empty[String, String]): boolean -> Similar to mount(), but updates an existing mount point (if present) instead of creating a new one

fsutils

cp(from: String, to: String, recurse: boolean = false): boolean -> Copies a file or directory, possibly across FileSystems
head(file: String, maxBytes: int = 65536): String -> Returns up to the first 'maxBytes' bytes of the given file as a String encoded in UTF-8
ls(dir: String): Seq -> Lists the contents of a directory
mkdirs(dir: String): boolean -> Creates the given directory if it does not exist, also creating any necessary parent directories
mv(from: String, to: String, recurse: boolean = false): boolean -> Moves a file or directory, possibly across FileSystems
put(file: String, contents: String, overwrite: boolean = false): boolean -> Writes the given String out to a file, encoded in UTF-8
rm(dir: String, recurse: boolean = false): boolean -> Removes a file or directory
```

Main commands used are:

mkdirs → To create a new directory

ls → It is used to list files in a specified directory

cp → To copy a file from one directory to another directory

mv → used to move the file from one directory to other directory

head → used to display the content of the file

put → used to put the content into a file

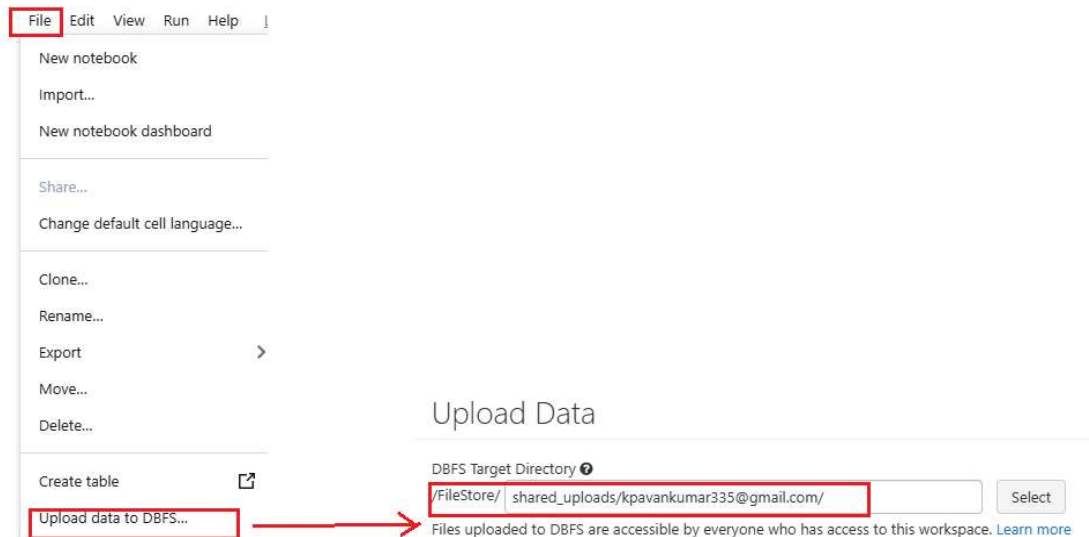
rm → used to delete the file permanently

## 1. Create a new directory:

#Syntax → `dbutils.fs.mkdirs('file system path')`

The path of the file system, can be taken by below process.

Click on File → Upload data to DBFS → Then take the path



```
07:23 PM (<1s) 3
#Creating a new inbound directory
#Syntax --> dbutils.fs.mkdirs('Filepath')
#mkdirs --> Is a linux command used to create new directory.
dbutils.fs.mkdirs('/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound')
dbutils.fs.mkdirs('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound')
dbutils.fs.mkdirs('/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles')
Out[3]: True
```

3 directories are created. Inbound, outbound, srcfiles.

## 2. I uploaded the files using above process in a pic, lets list the files in a inbound folder.

#Syntax → `dbutils.fs.ls('path to the directory')`

```
07:26 PM (<1s) 4
#list files in a specific directory
#Syntax: dbutils.fs.ls('filesystem path')
#ls --> It is a command used to list all files in a dirctory.
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound')
Out[4]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Customers.txt', name='Customers.txt', size=61, modificationTime=1740145943000),
FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Orders.txt', name='Orders.txt', size=69, modificationTime=1740145943000),
FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Sales.txt', name='Sales.txt', size=47, modificationTime=1740145943000),
FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Welcome.txt', name='Welcome.txt', size=0, modificationTime=1740145944000)]
```

### 3. Copy a file between directories.

#Syntax → `dbutils.fs.cp('source path', 'Target path')`

```
Just now (1s) 5

#copy files between inbound and outbound folders
#Syntax: dbutils.fs.cp(source path, target path)
#cp --> It is a command used to copy file between source and target.

dbutils.fs.cp(
    '/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Customers.txt',
    '/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound'
)

Out[8]: True
```

Here we can the customer.txt file copied to outbound folder.

```
Just now (<1s) 6 Python

dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/')

Out[9]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Customers.txt', name='Customers.txt', size=61, modificationTime=1740147386000)]
```

Case:

If we want to copy all files from source to target, we can use recursive argument.

By default, the recursive option will be false. We just mention it as True.

```
Just now (1s) 7

#Copy all files from inbound to outbound using recursive
#By default the recursive option will set to false, we need to make true.
#Syntax --> dbutils.fs.cp('src path', 'target path', True)

dbutils.fs.cp(
    '/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound',
    '/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound', True
)

Out[10]: True
```

Inbound directory:

```
Just now (<1s) 8

dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound')

Out[11]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Customers.txt', name='Customers.txt', size=61, modificationTime=1740145943000),
  FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Orders.txt', name='Orders.txt', size=69, modificationTime=1740145943000),
  FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Sales.txt', name='Sales.txt', size=47, modificationTime=1740145943000),
  FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Welcome.txt', name='Welcome.txt', size=0, modificationTime=1740145944000)]
```

Outbound directory:

```
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound')
```

```
Out[12]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Customers.txt', name='Customers.txt', size=61, modificationTime=1740147756000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Orders.txt', name='Orders.txt', size=69, modificationTime=1740147756000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Sales.txt', name='Sales.txt', size=47, modificationTime=1740147756000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Welcome.txt', name='Welcome.txt', size=0, modificationTime=1740147756000)]
```

#### 4. Move a file between directories.

Source – Inbound

Target – srcfiles

#Syntax – dbutils.fs.mv('source path', 'target path')

```
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles')
```

```
Out[13]: []
```

```
#copy files between inbound and outbound folders
#Syntax: dbutils.fs.mv(source path, target path)
#mv --> It is a command used to move file between source and target.

dbutils.fs.mv(
    '/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Customers.txt',
    '/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles'
)
```

```
Out[14]: True
```

```
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound')
```

```
Out[15]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Orders.txt', name='Orders.txt', size=69, modificationTime=1740145943000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Sales.txt', name='Sales.txt', size=47, modificationTime=1740145943000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Welcome.txt', name='Welcome.txt', size=0, modificationTime=1740145944000)]
```

```
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles')
```

```
Out[16]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles/Customers.txt', name='Customers.txt', size=61, modificationTime=1740148080000)]
```

Case:

If we want to move all files from source to target, we can use recursive argument.

By default, the recursive option will be false. We just mention it as True.

But, be aware that after moving all the files the source directory will be deleted.

List of files in inbound folder before:

```
▶ Just now (<1s) 14
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound')

Out[17]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Orders.txt', name='Orders.txt', size=69, modificationTime=1740145943000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Sales.txt', name='Sales.txt', size=47, modificationTime=1740145943000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound/Welcome.txt', name='Welcome.txt', size=0, modificationTime=1740145944000)]
```

Execution: From inbound to outbound

```
▶ Just now (1s) 15 Python
#Move all files from inbound to outbound using recursive
#By default the recursive option will set to false, we need to make true.
#Syntax --> dbutils.fs.mv('src path','target path',True)

dbutils.fs.mv(
    '/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound',
    '/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound', True
)

Out[18]: True
```

List of files in inbound folder post execution:

```
▶ Last execution failed 16
1 dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/inbound')

> java.io.FileNotFoundException: /FileStore/shared_uploads/kpavankumar335@gmail.com/inbound
```

We can observe that, directory got deleted. There is no inbound folder.

## Upload Data

Select a folder from /FileStore/ ⓘ

shared_uploads	kpavankumar335@gmail.com	outbound
tables		srcfiles



Files in outbound folder:

```
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound')
```

```
Out[21]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Customers.txt', name='Customers.txt', size=61, modificationTime=1740147750000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Orders.txt', name='Orders.txt', size=69, modificationTime=1740148231000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Sales.txt', name='Sales.txt', size=47, modificationTime=1740148232000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Welcome.txt', name='Welcome.txt', size=0, modificationTime=1740148232000)]
```

## 5. Head- used to display the content of the file

#Syntax- dbutils.fs.head('file path')

```
#to display the content of the file.
#syntax - dbutils.fs.head('file path')

dbutils.fs.head('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Customers.txt')
```

```
Out[25]: 'Cust_id,Cust_name\r\n1, Srinivas\r\n2, Pavan\r\n3, Phani\r\n4, Naresh'
```

## 6. Put – To input content to the file.

#Syntax – dbutils.fs.put('file path', 'content to be added')

```
#to create new file and input content to the file
#syntax - dbutils.fs.put('file path', 'content')
dbutils.fs.put('dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles/welcome1.txt',
              'Welcome to Databricks Community'
              )
```

```
Wrote 31 bytes.
Out[41]: True
```

We can see below, content of the welcome1.txt file.

```
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles')
```

```
Out[42]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles/Customers.txt', name='Customers.txt', size=61, modificationTime=1740148080000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles/welcome.txt', name='welcome.txt', size=0, modificationTime=1740149612000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles/welcome1.txt', name='welcome1.txt', size=31, modificationTime=1740149710000)]
```

```
dbutils.fs.head('dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/srcfiles/welcome1.txt')
```

```
Out[43]: 'Welcome to Databricks Community'
```

## 7. To delete the files in a directory.

# syntax - dbutils.fs.rm('filepath')

Before execution:

```
1 minute ago (<1s) 20
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound')

Out[31]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Customers.txt', name='Customers.txt', size=61, modificationTime=1740147755000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Orders.txt', name='Orders.txt', size=69, modificationTime=1740148231000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Sales.txt', name='Sales.txt', size=47, modificationTime=1740148232000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Welcome.txt', name='Welcome.txt', size=0, modificationTime=1740148232000)]
```

Execution:

```
1 minute ago (<1s) 21 Python
# to delete the files
# syntax - dbutils.fs.rm('filepath')
dbutils.fs.rm('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Welcome.txt')

Out[33]: True
```

Post Execution:

```
1 minute ago (<1s) 22
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound')

Out[34]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Customers.txt', name='Customers.txt', size=61, modificationTime=1740147755000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Orders.txt', name='Orders.txt', size=69, modificationTime=1740148231000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Sales.txt', name='Sales.txt', size=47, modificationTime=1740148232000)]
```

Case:

If we want to delete all files from directory, we can use recursive argument.

By default, the recursive option will be false. We just mention it as True.

But, be sure to check while deleting in production environment.

List of files in outbound folder before:

```
3 minutes ago (<1s) 22 Python
dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound')

Out[34]: [FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Customers.txt', name='Customers.txt', size=61, modificationTime=1740147755000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Orders.txt', name='Orders.txt', size=69, modificationTime=1740148231000),
          FileInfo(path='dbfs:/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound/Sales.txt', name='Sales.txt', size=47, modificationTime=1740148232000)]
```

Execution step:



```
Just now (<1s) 23

#Deletion of all files in a directory
#Syntax: dbutils.fs.rm(source path, recursive argument)

dbutils.fs.rm(
  '/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound',
  True
)

Out[35]: True
```

Post deletion, all the files were deleted including the directory as well.

```
Last execution failed 24

1 dbutils.fs.ls('/FileStore/shared_uploads/kpavankumar335@gmail.com/outbound')

> java.io.FileNotFoundException: /FileStore/shared_uploads/kpavankumar335@gmail.com/outbound
```