**Name: Abhishek Kanoujia**

**DATA ENGINEERING BATCH 1**

**DAY 13 ASSIGNMENT**

**Class hand written notes:-**

Pyspark

Date _____

05/02/2024 ( Monday )

1) Introduction to pyspark and setting up the environment
2) Hands -on Exercise:
3)
4)

Pyspark: its is python Apache spark library written in python to run python ap using Apache spark capability.

# Pyspark is python Api.

Apache spark:

⇒ It is an open source unified analytics engine used for large data processing.

⇒ spark can run on the single node as well as multiple node.
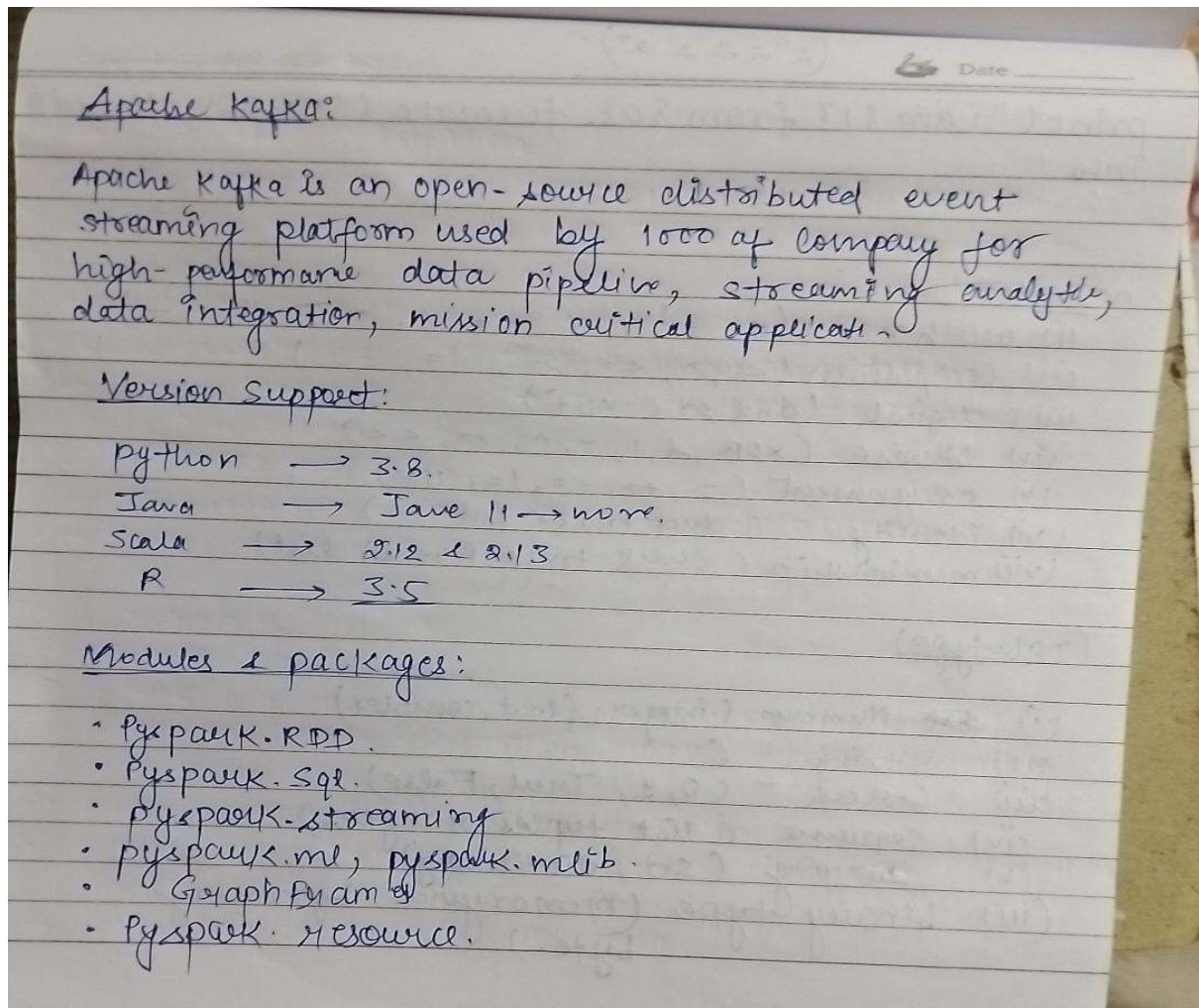
## for development:

→ Anaconda distribution
→ Spyder IDE
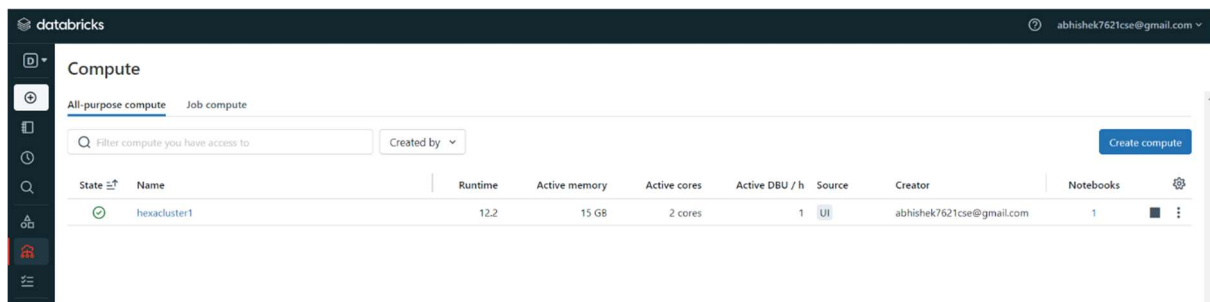→ Jupyter.

} <u>All can used</u>

## Pyspark Features

→ fault-tolerance
→ Immutable
→ Lazy evalution
→ Cache & persistance.
→ Inbuild-optimization when using dataframes.
→ Support ANSI SOL.

→ In-memory computation
→ Distributed processing.

## Advantages of Pyspark:

⇒ Pyspark is a general-purpose, in-memory, distributed processing engine that allow process data efficiently.

⇒ Running 100x-faster than normal.
⇒ get benefit using pyspark for data ~~engi~~ intrigation pipeline
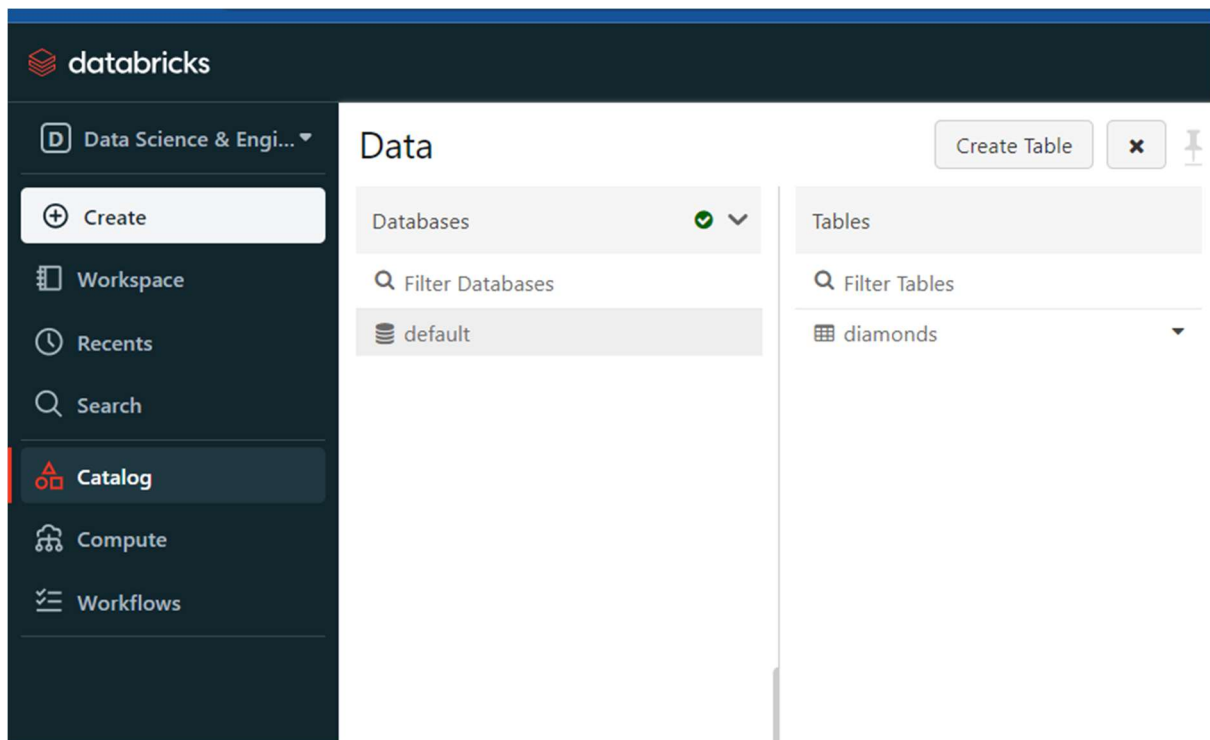⇒ used to process real-life data.

## Apache Kafka:

Apache Kafka is an open-source distributed event streaming platform used by 1000 of company for high-performane data pipline, streaming analytic, data integration, mission critical applicati.

### Version Support:

Python $\longrightarrow$ 3.8.

Java $\longrightarrow$ Java 11 $\longrightarrow$ more.

Scala $\longrightarrow$ 2.12 & 2.13

R $\longrightarrow$ 3.5

### Modules & packages:

- Pyspark.RDD.
- Pyspark.Sql.
- Pyspark.streaming
- pyspark.ml, pyspark.mlib.
-   Graph Fram e
- Pyspark.resource.

**Create cluster in databricks:-**



**Create table using csv file in databricks:-**

**Step:1**

**Step:2**

**Step:3**

## Specify Table Attributes

Specify the Table Name, Database and Schema to add this to the data UI for other users to access

| Table Name ⊘ | | Table Preview | | |
|---|---|---|---|---|
| currency_csv | | _c0 | _c1 | _c2 |
| **Create in Database** ⊘ | | STRING ⌄ | STRING ⌄ | STRING ⌄ |
| default ⇅ | | Code | Symbol | Name |
| **File Type** ⊘ | | AED | ﺩ.إ | United Arab Emirates d |
| CSV ⇅ | | AFN | ؋ | Afghan afghani |
| **Column Delimiter** ⊘ | | ALL | L | Albanian lek |
| , | | AMD | AMD | Armenian dram |
| ☐ First row is header ⊘ | | ANG | ƒ | Netherlands Antillean gu |
| ☐ Infer schema ⊘ | | | | |
| ☐ Multi-line ⊘ | | | | |
| ⊞ Create Table | | | | |

**Step:4**



**Table created successfully.**

**First program :-**

**1st program** Python ⌄ ☆
File  Edit  View  Run  Help    Last edit was 10 minutes ago    New cell UI: OFF ⌄                    ▶ Run all    ● hexacluster1 ⌄    Share    Publish    ⌃

Cmd 1

```
                                                                                      Python  ▶▾  ⌄  —  ✕
1    # Import SparkSession
2    from pyspark.sql import SparkSession
3
4
5    # Create SparkSession
6    spark = SparkSession.builder \
7            .master("local[1]") \
8            .appName("SparkByExamples.com") \
9            .getOrCreate()
10   dataList = [("Java", 20000), ("Python", 100000), ("Scala", 3000)]
11   rdd=spark.sparkContext.parallelize(dataList)
12   rdd.collect()
```

▸ (1) Spark Jobs

Out[1]: [('Java', 20000), ('Python', 100000), ('Scala', 3000)]

Command took 9.53 seconds -- by abhishek7621cse@gmail.com at 2/5/2024, 4:40:11 PM on hexacluster1

---

▶ Run all    ● hexacluster1 ⌄    Share    Publish    ⌃

```
                                                                                      Python  ▶▾  ⌄  —  ✕
1    # Import SparkSession
2    from pyspark.sql import SparkSession
3
4
5    # Create SparkSession
6    spark = SparkSession.builder \
7            .master("local[1]") \
```

```
In [1]:  import pyspark
```

```
In [2]:  from pyspark.sql import SparkSession
```

```
In [3]:  spark=SparkSession.builder.appName("Jupyter Notebook").getOrCreate()
```

```
In [4]:  spark
```

Out[4]:  **SparkSession - in-memory**
         **SparkContext**

         Spark UI
         **Version**
         v3.5.0
         **Master**
         local[*]
         **AppName**
         Jupyter Notebook

```
In [5]:  df=spark.read.csv("friends.csv")
```

```
In [6]:  df
```

Out[6]:  DataFrame[_c0: string, _c1: string, _c2: string, _c3: string]

```
In [7]:  df.show()
```

```
+----+------+-----+----------+
| _c0|   _c1|  _c2|       _c3|
+----+------+-----+----------+
|NULL|  name|marks|      city|
|   0| harry|   92|    rampur|
|   1| rohan|   34|   kolkata|
|   2|skillf|   24|  bareilly|
|   3| shubh|   17|antarctica|
+----+------+-----+----------+
```

```
In [8]:  from pyspark.sql import SparkSession
```

```
In [9]:  spark=SparkSession.builder.appName(" Notebook").getOrCreate()
```

```
In [10]:  dataList = [("Java", 20000), ("Python", 100000), ("Scala", 3000)]
```

```
In [12]:  rdd=spark.sparkContext.parallelize(dataList)
```

```
In [13]:  result=rdd.collect()
```

```
In [14]:  result
```

Out[14]:  [('Java', 20000), ('Python', 100000), ('Scala', 3000)]

```
In [ ]:
```