Name : Parth Nandedkar
Date : 12 Feb 2024
Topics : PySpark
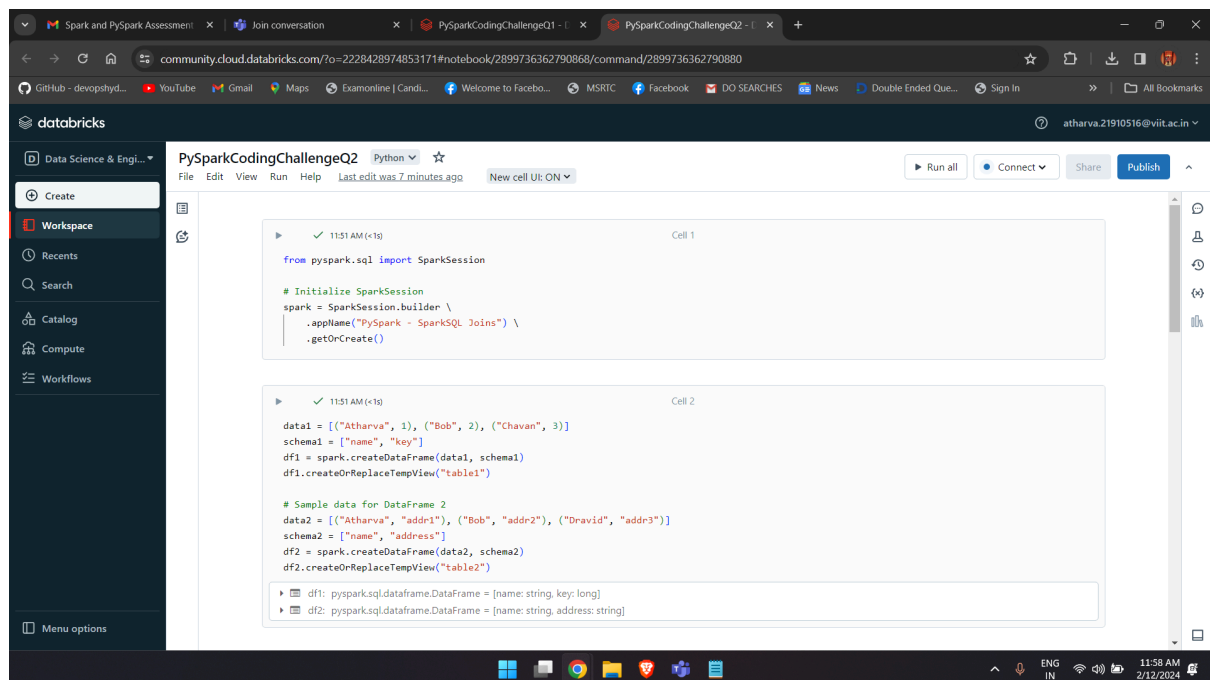Batch : Data Engineering Batch-1

## Q2. Execute Pyspark -sparksql joins & Applying Functions in a Pandas DataFrame.

Spark SQL is a module in Apache Spark that provides a higher-level abstraction for working with structured and semi-structured data. It allows you to run SQL queries and access data using SQL statements, as well as perform various data manipulation and analysis tasks using DataFrame API and SQL functions.

Pandas are used to manipulate/work on big data. Here we are using spark sql and pandas so that data from databases can be handled by pandas.

Here I'm using databricks to run Join queries as well as Pandas functions to get desired output.

**Initialising the Spark Session :**

**Creating Sample Tables in PySpark :**

As we can use data frames as table views and use spark sql queries on those to retrieve data .



Here **data1** and **schema1** are used to create a dataframe and we created a data frame **df1** so as to get data in it .
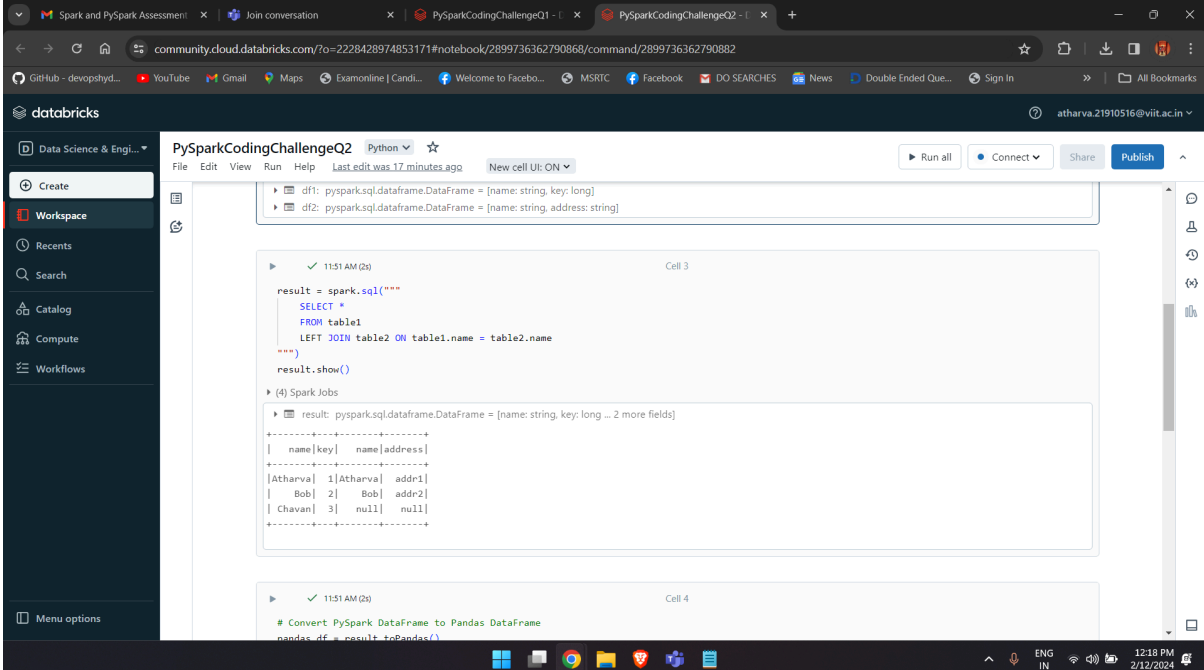And by using the **createOrReplaceTempView()** method we created **table1** which means table1 will be representing the name and code column .

Same procedure is followed to create table two and use the name **table2** to df2 so that table2 can be used to recognize name and address.

**Using Spark SQL joins :**

Spark SQL is an Integrated language to use databases using Spark service.

By using spark.sql() method we can give any query as input to retrieve results using spark.



Here in the above example **table1** and **table2**  are joined on names .
As table1 is the left table and we are using left join all the values from table1
will be considered and table2 common values will be there.

# Functions in Pandas dataframe :

In pandas data set is converted and results are as follows :

Above execution demonstrates how to perform a left join between two Pandas DataFrames (df1 and df2) using the merge() method and then apply a function to a column in the resulting DataFrame (result). Adjust the data and column names as needed based on your specific use case.