

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 4 / 93

Recap

Challenges in Existing Computing Methods

When we talk about distributed computing, in processing data over multiple interconnected servers, there are various challenges involved.

- There are different ways to implement distributed computing like MapReduce.
- The need to store data across multiple machines makes it difficult to access and process data.
- MapReduce is a slow and inefficient way to process data.
- It is also difficult to handle large amounts of data.

What is RDD?

An RDD is a distributed dataset that can be partitioned into multiple partitions. It is a collection of elements that are distributed across multiple machines.

- It is a distributed collection of elements that can be partitioned into multiple partitions.
- The elements in an RDD are immutable.
- It can be used for parallel processing of data.

Some Popular Transformations

map, flatMap, filter, distinct, reduce, etc.

Function Passing

It is a function that takes one or more arguments and returns a value. It can be applied to data structures.

Passing Functions to Spark

It is a function that takes one or more arguments and returns a value. It can be applied to data structures.

Shared Variables – Broadcast Variable

A broadcast variable is a variable that is replicated across all nodes in a cluster. It is used to share data between multiple workers.

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:19 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation

Topics

- Why Spark SQL ?
- Spark SQL Success Story
- What is Spark SQL?
- Spark SQL Features
- Comparison between Hive, Impala & Spark SQL
- Spark SQL Architecture
- Spark SQL Libraries
- Data Source API
- DataFrame API
- SQL interpreter and Optimizer
- SQL Service
- What is SQL context
- Creating Dataframe
- User Defined Function
- Inferring the Schema using Reflection
- Hive integration with Spark SQL
- Analysis Using Spark SQL

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:19 01-06-2023

X PySpark Certification T...

DataFrames and Spark SQL > Presentation



Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Personal Library

Presentation

6 / 93



Objectives

After completing this module, you should be able to:

- Describe What is Spark SQL
- Understand Spark SQL Architecture
- Describe SQL Context in Spark SQL
- Work with Data Frames
- Interoperate with RDDs
- Understand JSON and Parquet File Formats
- Integrate Spark with Hive



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

7 / 93

Why Spark SQL?

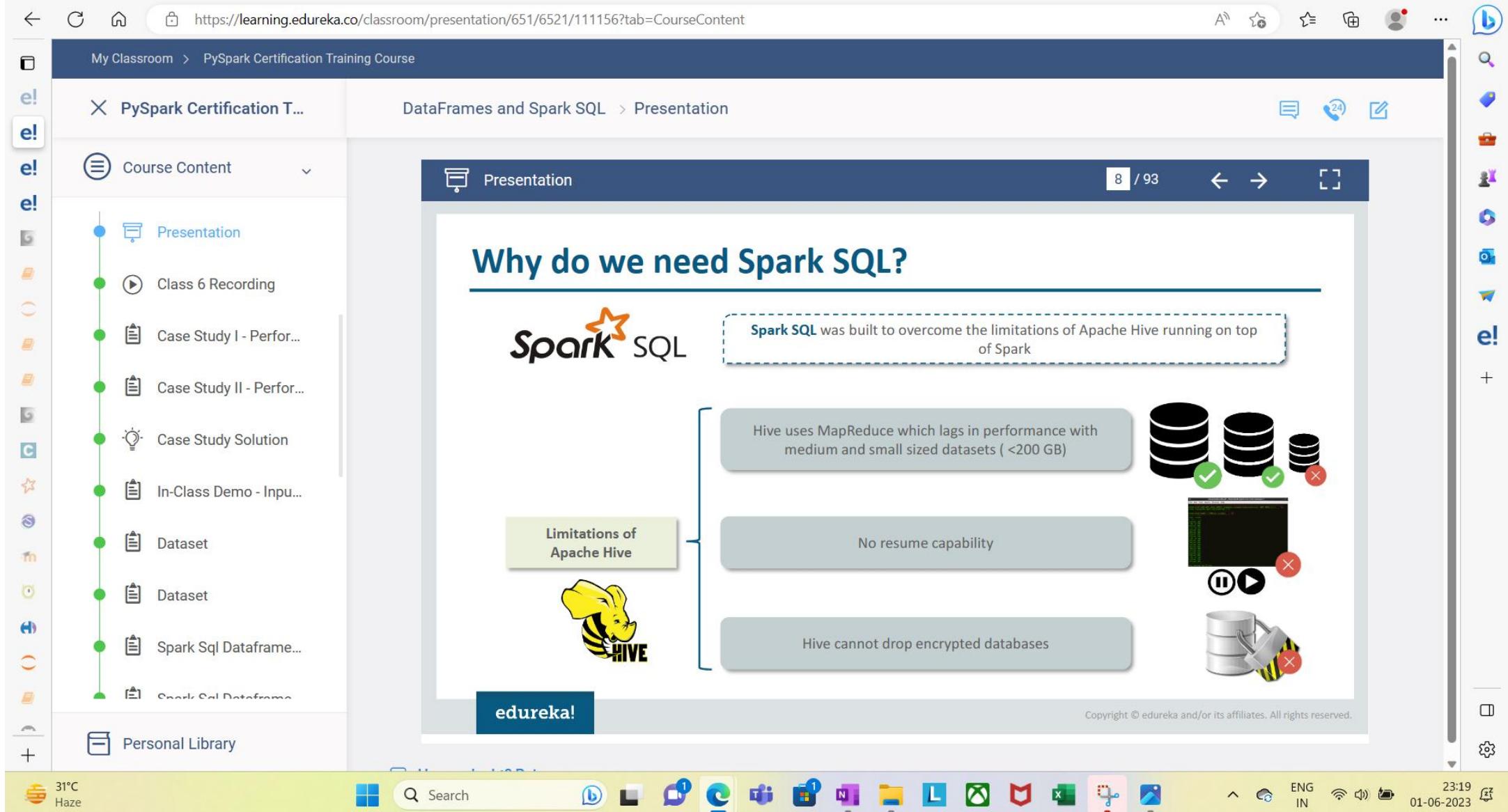
Copyright © edureka and/or its affiliates. All rights reserved.

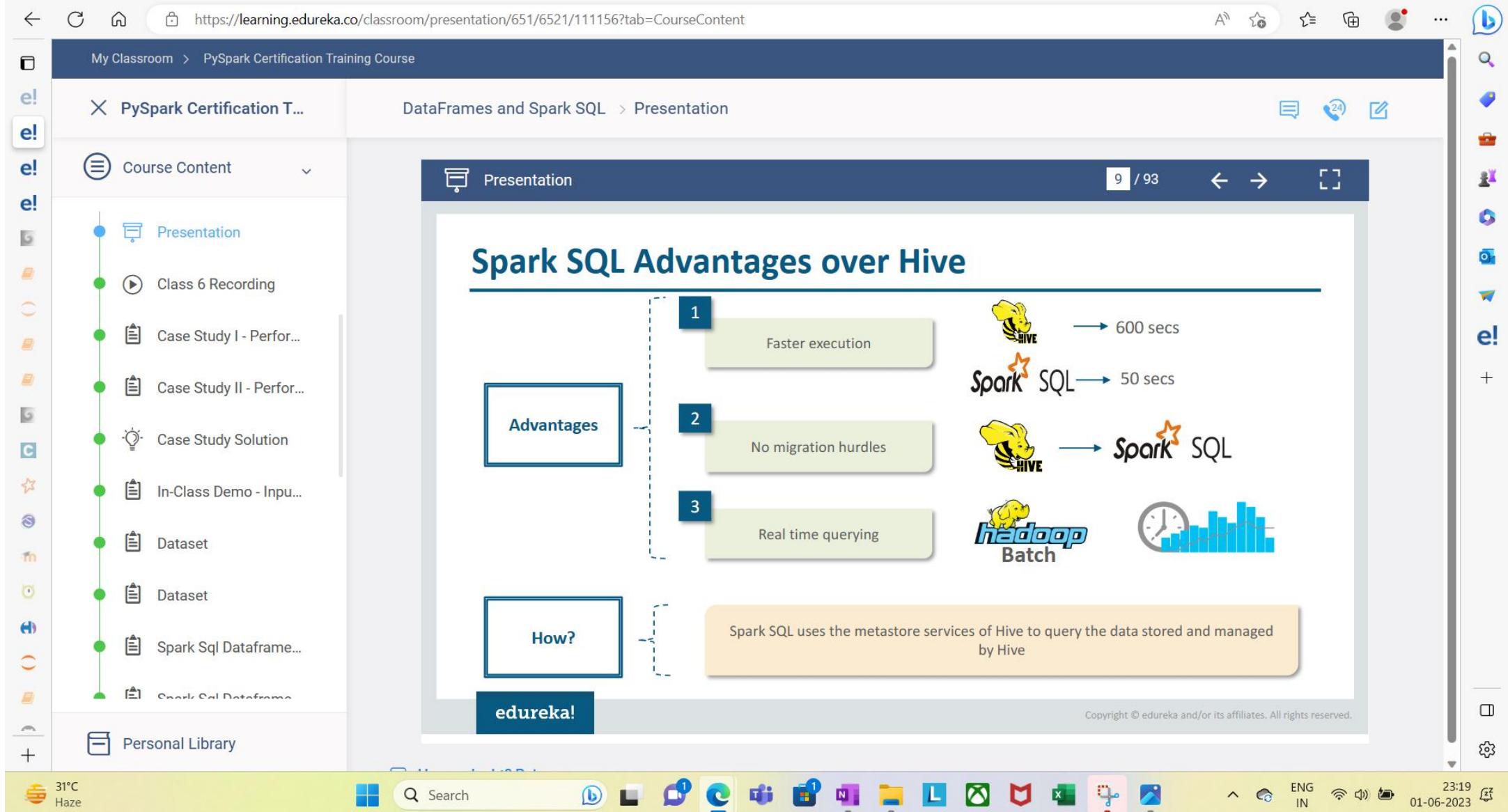
31°C Haze

Search

edureka!

23:19 01-06-2023





https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:19 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 11 / 93

Spark SQL Success Story

 Twitter Sentiment Analysis With Spark SQL

Trending Topics can be used to create campaigns and attract larger audience

Sentiment helps in crisis management, service adjusting and target marketing

 NYSE: Real Time Analysis of Stock Market Data

 Banking: Credit Card Fraud Detection

 Genomic Sequencing

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

e! Course Classroom | Edureka

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 12 / 93 ← → []

edureka!

What is Spark SQL?

Copyright © edureka and/or its affiliates. All rights reserved.

Personal Library

31°C Haze

Search

edureka!

ENG IN

01-06-2023

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Personal Library

What is Spark SQL?

Spark SQL is a Spark module for structured data processing

Spark SQL Capabilities

- It provides a DataFrame abstraction in Python, Java, and Scala to simplify working with structured datasets. DataFrames are similar to tables in a relational database
- It can read and write data in a variety of structured formats (e.g., JSON, Hive Tables, and Parquet)
- It lets you query the data using SQL, both inside a Spark program and from external tools

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:19 01-06-2023

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

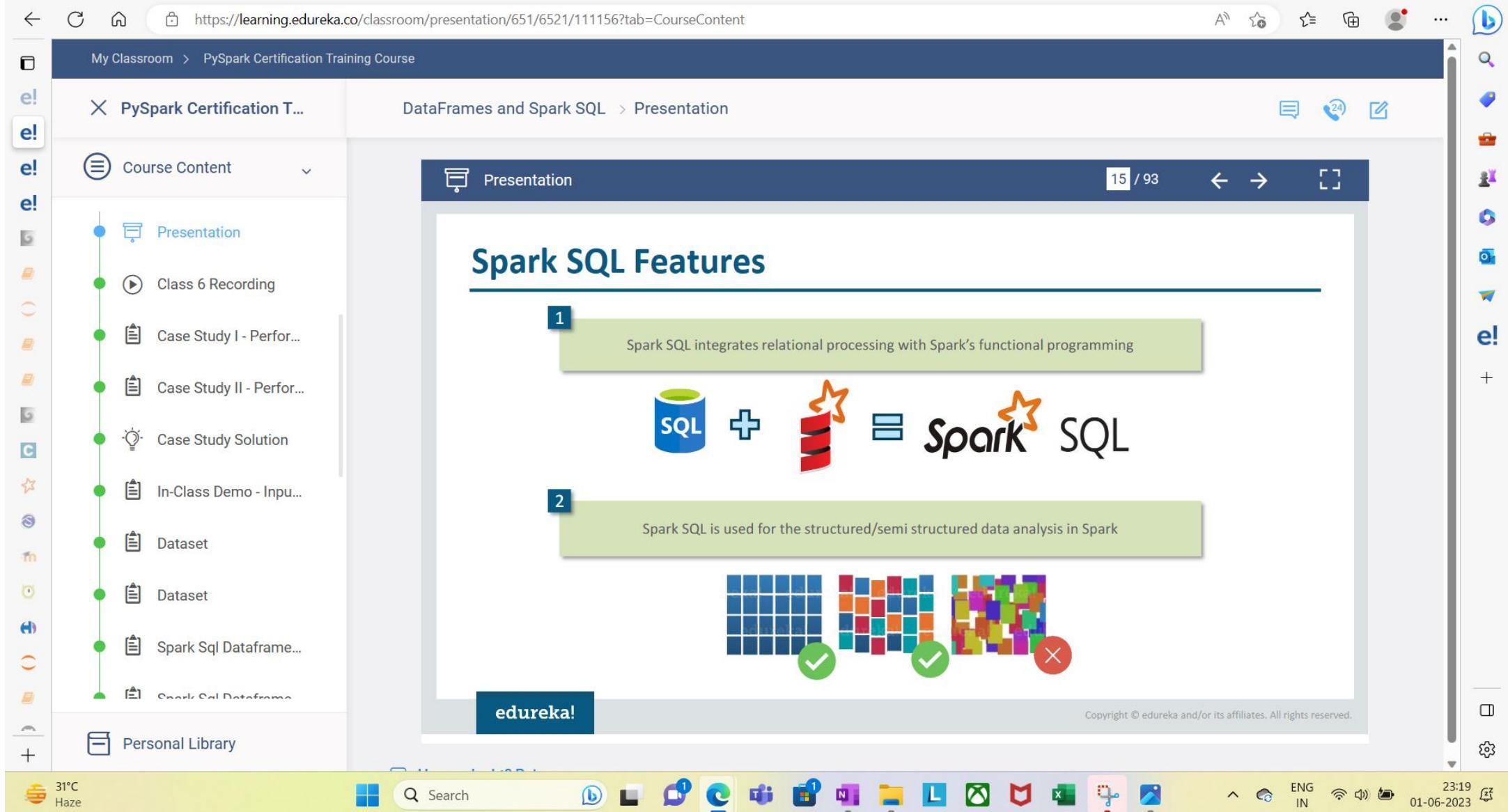
31°C Haze

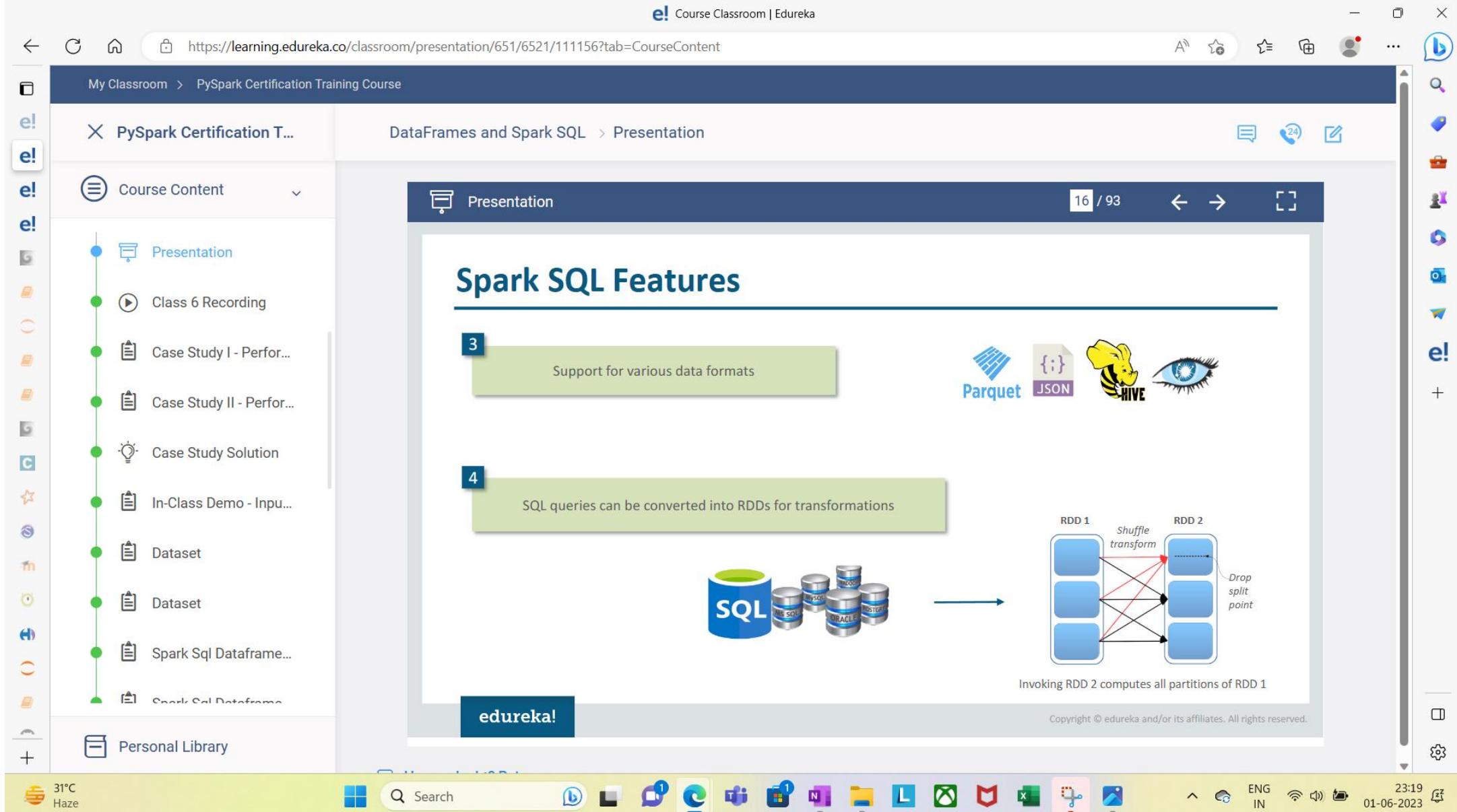
Search

L

ENG IN

23:19 01-06-2023





X PySpark Certification T...

DataFrames and Spark SQL > Presentation



Course Content

Presentation

Class 6 Recording

Case Study I - Perform...

Case Study II - Perform...

Case Study Solution

In-Class Demo - Inpu...

Dataset

Dataset

Spark Sql Dataframe...

Spark Sql Dataframe...

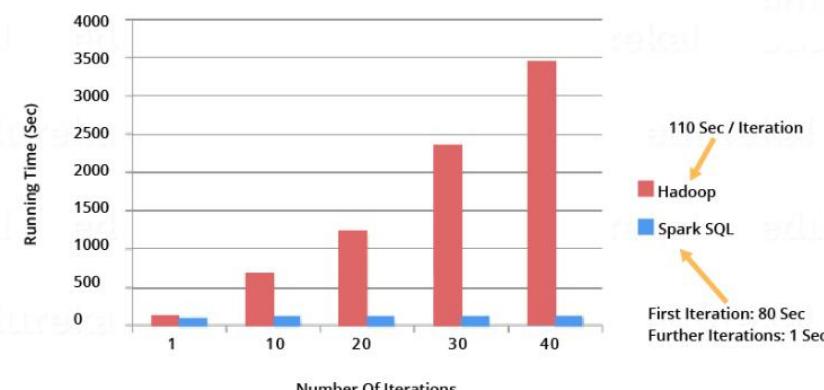
Personal Library

Spark SQL Features

5

Performance And Scalability

Performance: Spark SQL Vs Hadoop



Copyright © edureka and/or its affiliates. All rights reserved.

X PySpark Certification T...

DataFrames and Spark SQL > Presentation



Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Personal Library

Presentation

18 / 93



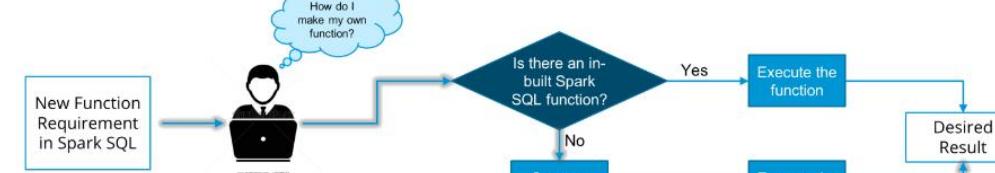
6

Standard JDBC/ODBC Connectivity



7

User Defined Functions lets users define new Column-based functions to extend the Spark vocabulary



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:19 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 20 / 93 ← →

Choosing the right SQL Engine



	Apache Hive	Apache Impala	Apache Spark SQL
Users	ETL Developers	Business Analysts	Data Engineers & Data Scientists
Strengths	<ul style="list-style-type: none">Built for very long running ETL, data processing or batch processingSupports custom file formatsHandles massive ETL sorts with joins	<ul style="list-style-type: none">Scales to high-concurrencySupports high-performance interactive SQLCompatible with BI tools & skillsHadoop integration and usability	<ul style="list-style-type: none">Easily embed SQL into Java, Scala or Python applicationsSimple language for common operationsSeamlessly mix SQL & Spark code within a single applicationsAutomatic performance optimizations

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:19 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 22 / 93

Spark SQL Architecture

Architecture Of Spark SQL

```
graph TD; DSL[DataFrame DSL] --- API[DataFrame API]; API --- DSAPI[Data Source API]; DSAPI --- CSV[CSV]; DSAPI --- JSON[JSON]; DSAPI --- JDBC[JDBC]; DSAPI --- Parquet[Parquet]; DSAPI --- External[External Source]
```

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

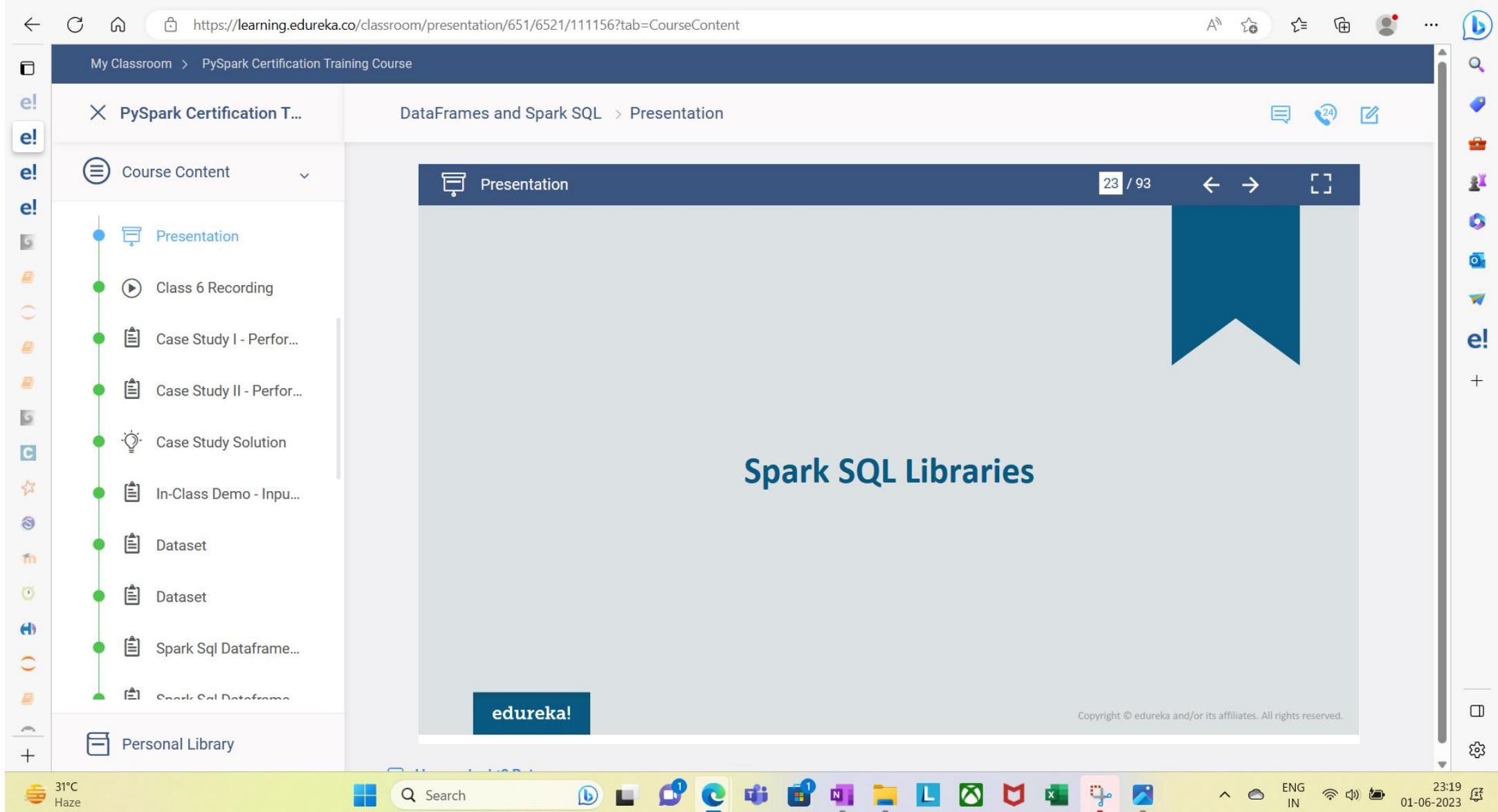
31°C Haze

Search

L

ENG IN

01-06-2023



My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 24 / 93

Spark SQL Libraries

Spark SQL has the following libraries:

- 1 Data Source API
- 2 DataFrame API
- 3 Interpreter & Optimizer
- 4 SQL Service



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

X PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

Presentation

25 / 93

← →

Data Source API

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 26 / 93

Data Source API

Data Source API is used to *read* and *store structured* and *semi-structured data* into *Spark SQL*

Features:

- Can handle structured/ semi-structured data
- Can load files in multiple formats
- 3rd party integration can be done

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

01-06-2023

https://learning.edureka.co/course/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

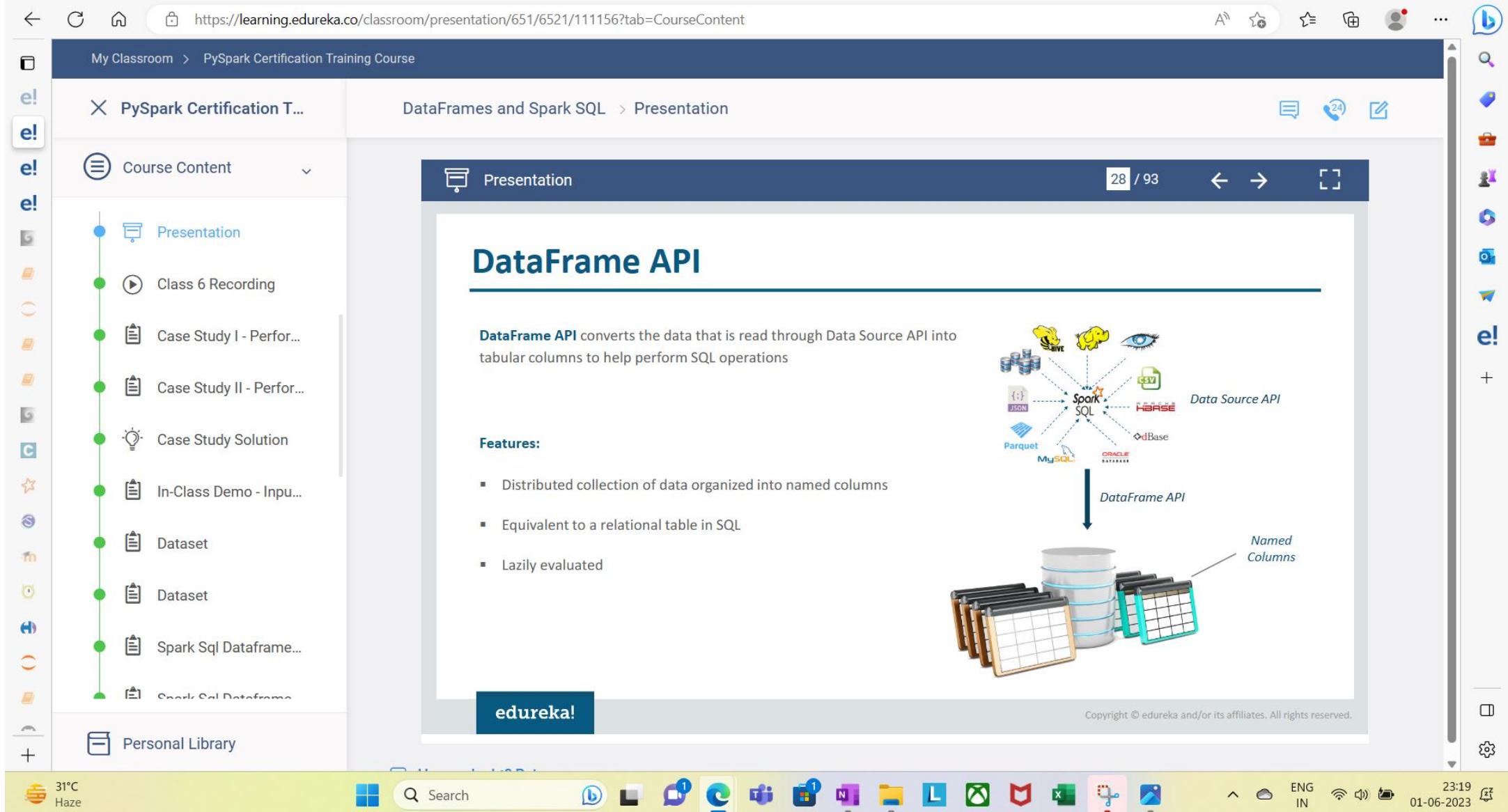
31°C Haze

Search

L

ENG IN

23:19 01-06-2023



https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:19 01-06-2023

X PySpark Certification T...

DataFrames and Spark SQL > Presentation



Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Personal Library

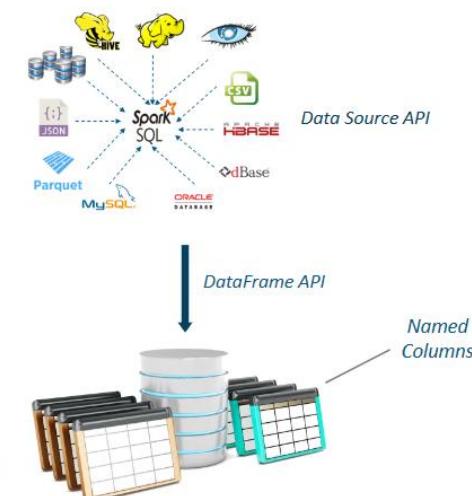
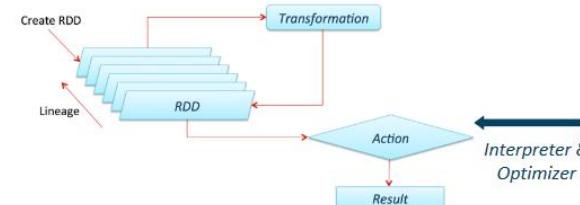
SQL Interpreter and Optimizer

SQL Interpreter & Optimizer handles the *functional programming* part of *Spark SQL*. It transforms the *DataFrames RDDs* to get the *required results* in the *required formats*

Features:

- Functional programming
- Transforming trees
- Faster than RDDs
- Processes all size data

e.g. Catalyst: A modular library for distinct optimization



Copyright © edureka and/or its affiliates. All rights reserved.

edureka!

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31 / 93

SQL Service

edureka!

31°C Haze

Search

Cloud

File Explorer

OneDrive

PowerPoint

Word

Excel

Teams

Outlook

OneNote

Power BI

Xbox

Bookmarks

Snipping Tool

Calculator

File

ENG IN

23:19

01-06-2023

X PySpark Certification T...

DataFrames and Spark SQL > Presentation



Course Content

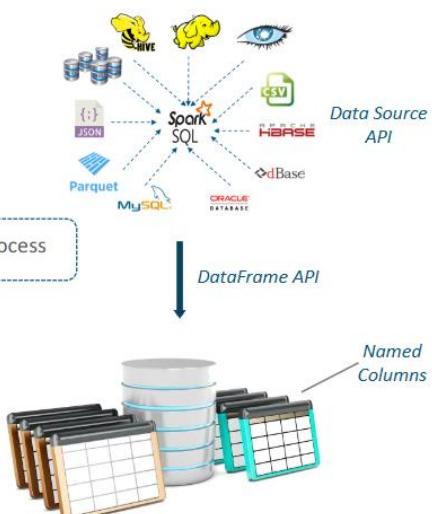
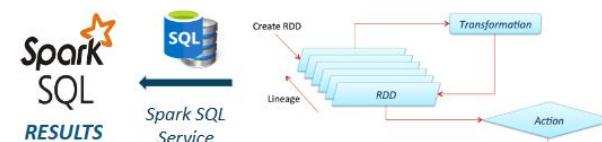
- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Personal Library

SQL Service

- SQL Service is the entry point for working along **structured data** in Spark
- SQL is used to fetch the result from the interpreted & optimized data

We have thus used all the four libraries in sequence. This completes a Spark SQL process



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

SQL Context

- **SQLContext** is a class and is used for initializing the functionalities of Spark SQL
- **SparkContext** class object (**sc**) is required for initializing **SQLContext** class object



The following command is used for initializing the **SparkContext** through **Pyspark-shell**:

```
$ pyspark2
```

- By default, the **SparkContext** object is initialized with the name **sc** when the **Pyspark-shell** starts

Use the following command to create **SQLContext**:

```
sqlContext=SQLContext(sc)
```



My Classroom > PySpark Certification Training Course

DataFrames and Spark SQL > Presentation



 Course Content

-  Presentation
 -  Class 6 Recording
 -  Case Study I - Perform.
 -  Case Study II - Perform.
 -  Case Study Solution
 -  In-Class Demo - Inpu..
 -  Dataset
 -  Dataset
 -  Spark Sql Dataframe..
 -  Case Study 1 - f

SQL Context

edureka

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

01-06-2023

23:19

X PySpark Certification T...

DataFrames and Spark SQL > Presentation



Course Content

-  Presentation
 -  Class 6 Recording
 -  Case Study I - Perform...
 -  Case Study II - Perform...
 -  Case Study Solution
 -  In-Class Demo - Input...
 -  Dataset
 -  Dataset
 -  Spark Sql Dataframe...
 -  SQL - Q1 P1 & f...

Spark Session - New Entry Point for Spark

- In earlier versions of spark, **Spark Context** was entry point for Spark
 - All functionality previously available through Spark Context, SQLContext or HiveContext in early versions of Spark are now available via **Spark Session**



 So in rest of this module, we will implement different functionalities of Spark SQL using Spark Session

edureka



My Classroom > PySpark Certification Training Course

DataFrames and Spark SQL > Presentation



 Course Content

-  Presentation
 -  Class 6 Recording
 -  Case Study I - Perform...
 -  Case Study II - Perform...
 -  Case Study Solution
 -  In-Class Demo - Inpu...
 -  Dataset
 -  Dataset
 -  Spark Sql Dataframe...
 -  Cloud Sql Dataframe...

Starting up Spark Session

You can either use “`spark`” or “`sqlContext`” to start working on Spark SQL

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Presentation

39 / 93

Demo 1 – Creating a Dataframe

Refer to the file Module-6 Demo 1 provided in the LMS for all the steps in detail

edureka!

Copyright © 2018, edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

Cloud

ENG IN

23:19

01-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

Presentation

40 / 93

← →

Running SQL Queries Programmatically

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

- e!
- e!
- e!
- e!
- [Course Content](#)
- [Presentation](#)
- [Class 6 Recording](#)
- [Case Study I - Perform...](#)
- [Case Study II - Perform...](#)
- [Case Study Solution](#)
- [In-Class Demo - Inpu...](#)
- [Dataset](#)
- [Dataset](#)
- [Spark Sql Dataframe...](#)
- [Spark Sql Dataframe...](#)
- [Personal Library](#)



Presentation

41 / 93

← →

[]

Running SQL Queries Programmatically

The `sql` function on a `SparkSession` enables applications to run SQL queries programmatically and returns the result as a `DataFrame`

```
>>> employee_df.createOrReplaceTempView("EMP")
>>> sqlDF = spark.sql("SELECT * FROM EMP")
>>> sqlDF.show()
```

// Register the DataFrame as a SQL temporary view

ID	Name	Salary
1	John	20000
2	Rohit	15000
3	Parth	14600
4	Rishabh	20500
5	Daisy	34000
6	Annie	23000
7	Sushmita	50000
8	Kaivalya	20000
9	Varun	70000
10	Shambhavi	21500
11	Johnson	25500
12	Riya	17000
13	Krish	17000
14	Akanksha	20000
15	Rutuja	21000

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

Cloud

File

OneDrive

OneNote

PowerPoint

Excel

Word

Power BI

Teams

Outlook

Skype

OneDrive

OneNote

PowerPoint

Excel

Word

Power BI

Teams

Outlook

Skype

Cloud

File

31°C Haze

ENG IN

23:19

01-06-2023

https://learning.edureka.co/course/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

42 / 93

JSON Dataset

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:19

01-06-2023



My Classroom > PySpark Certification Training Course



X PySpark Certification T...



Course Content



Presentation



Class 6 Recording



Case Study I - Perform...



Case Study II - Perform...



Case Study Solution



In-Class Demo - Inpu...



Dataset



Dataset



Spark Sql Dataframe...



Spark Sql Dataframe...



Personal Library

DataFrames and Spark SQL > Presentation



JSON Data - Loading File

43 / 93



Spark SQL can automatically infer the schema of a JSON dataset and load it as a Dataset[Row]. This conversion can be done using `SparkSession.read.json()` on a JSON file

```
>>> employee_df=spark.read.json("hdfs://nameservice1/user/edureka_294428/Employees.json")
```

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C
Haze

Search



23:19

01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 44 / 93

JSON Data - Parquet File

Parquet is a *columnar format* and Spark SQL provides support for both *reading* and *writing Parquet files* that automatically *preserves the schema* of the original data

```
>>> peopleDF=spark.read.json("hdfs://nameservice1/user/edureka_294428/people.json")
>>> peopleDF.write.parquet("hdfs://nameservice1/user/edureka_294428/people8.parquet")
>>> parquetFile=spark.read.parquet("people8.parquet")
>>> parquetFile.createOrReplaceTempView('parquetFile')
>>> teenagers=spark.sql("SELECT name FROM parquetFile WHERE age>=13 AND age<=19")
>>> teenagers.show()
+---+
| name|
+---+
|justin|
+---+
```

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

01-06-2023

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

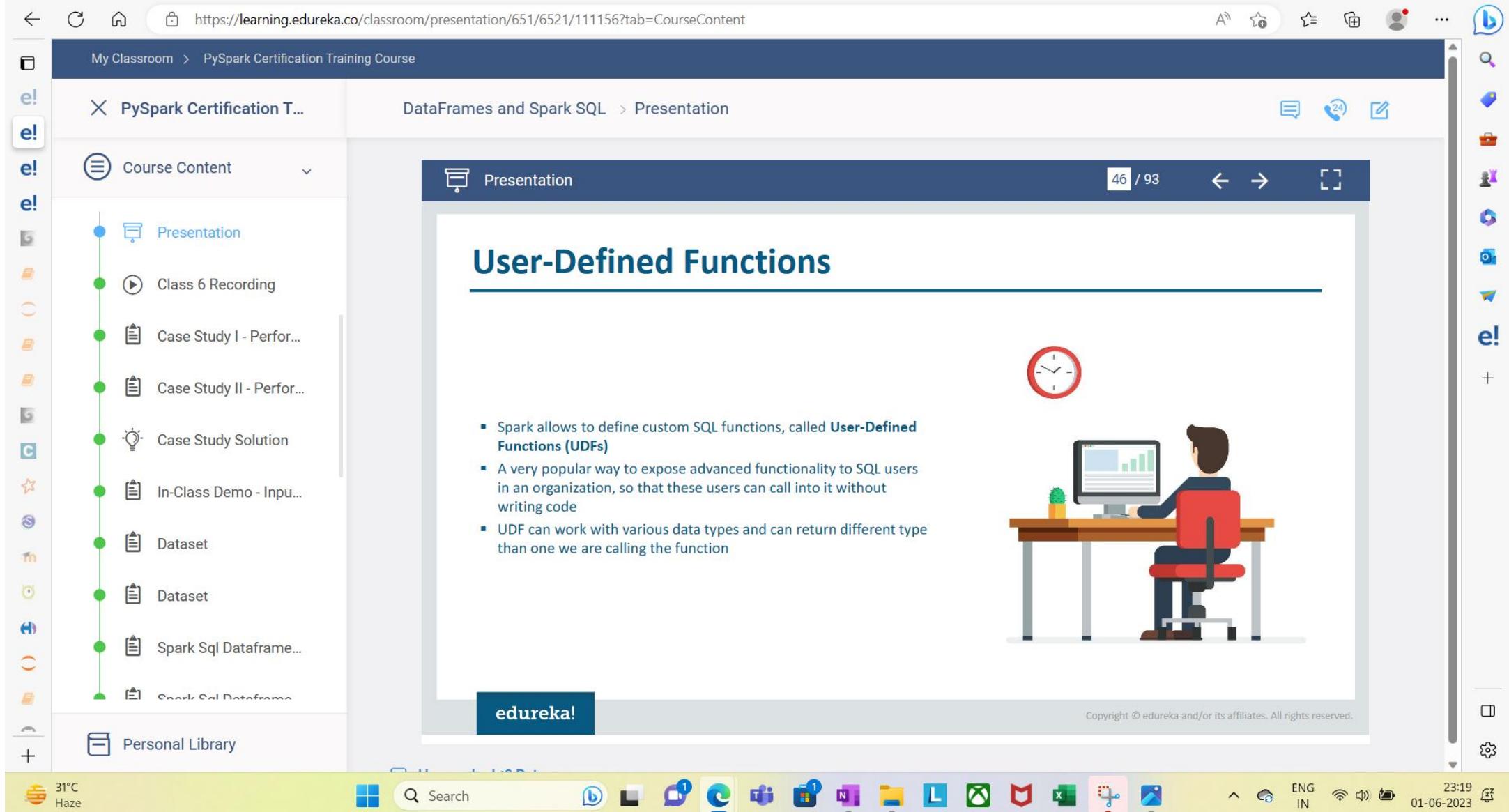
31°C Haze

Search

L

ENG IN

23:19 01-06-2023



https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Presentation

47 / 93

Demo 2 – User-defined Functions

Refer to the file Module-6 Demo 2 provided in the LMS for all the steps in detail

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

Cloud

OneDrive

Microsoft Edge

Teams

OneNote

File Explorer

Libraries

Xbox

Bookmarks

Excel

PowerPoint

PowerPoint

23:19

ENG IN

01-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

DataFrames and Spark SQL > Presentation



Course Content

Presentation

Class 6 Recording

Case Study I - Perform...

Case Study II - Perform...

Case Study Solution

In-Class Demo - Inpu...

Dataset

Dataset

Spark Sql Dataframe...

Spark Sql Dataframe...

Personal Library

Presentation

48 / 93



UDFs are Blackbox



UDFs are great when built-in SQL functions aren't sufficient, but should be used sparingly because they're not performant

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 50 / 93

Interoperating with RDDs

Spark SQL supports two different methods for converting existing RDDs into Datasets

1 Inferring the Schema Using Reflection

This method uses reflection to infer the schema of an RDD that contains specific types of objects. This reflection based approach works well when we already know the schema while writing a spark application

2 Programmatically Specifying the Schema

This second method for creating Datasets is through a programmatic interface that allows us to construct a schema and then apply it to an existing RDD

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

X PySpark Certification T...

DataFrames and Spark SQL > Presentation



e! Course Content

Presentation

Class 6 Recording

Case Study I - Perform...

Case Study II - Perform...

Case Study Solution

In-Class Demo - Inpu...

Dataset

Dataset

Spark Sql Dataframe...

Spark Sql Dataframe...

e! Personal Library

Presentation

52 / 93



Inferring Schema

- Data exists in multiple formats including TXT, CSV, DOC, PDF etc.
- At times you might be aware of the schema (order of data) of the data.
- When the schema is already known we use the concept of Reflection to directly infer the schema



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation

53 / 93

Demo: Inferring Schema

John, 20
Rohit, 15
Parth, 14
Rishabh, 25
Daisy, 34
Annie, 23
Sushmita, 50
Kaivalya, 26
Varun, 16
Shambhavi, 21
Johnson, 22
Riya, 17
Krish, 19
Akanksha, 28
Rutuja, 33

Consider we have a text file of the following format. Please note that here we are considering that schema of the file is already known and is stored in a TXT format (Non-Columnar)

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

Cloud

File Explorer

OneDrive

PowerPoint

Word

Excel

Teams

Outlook

OneNote

Power BI

Xbox

Bookmarks

Calculator

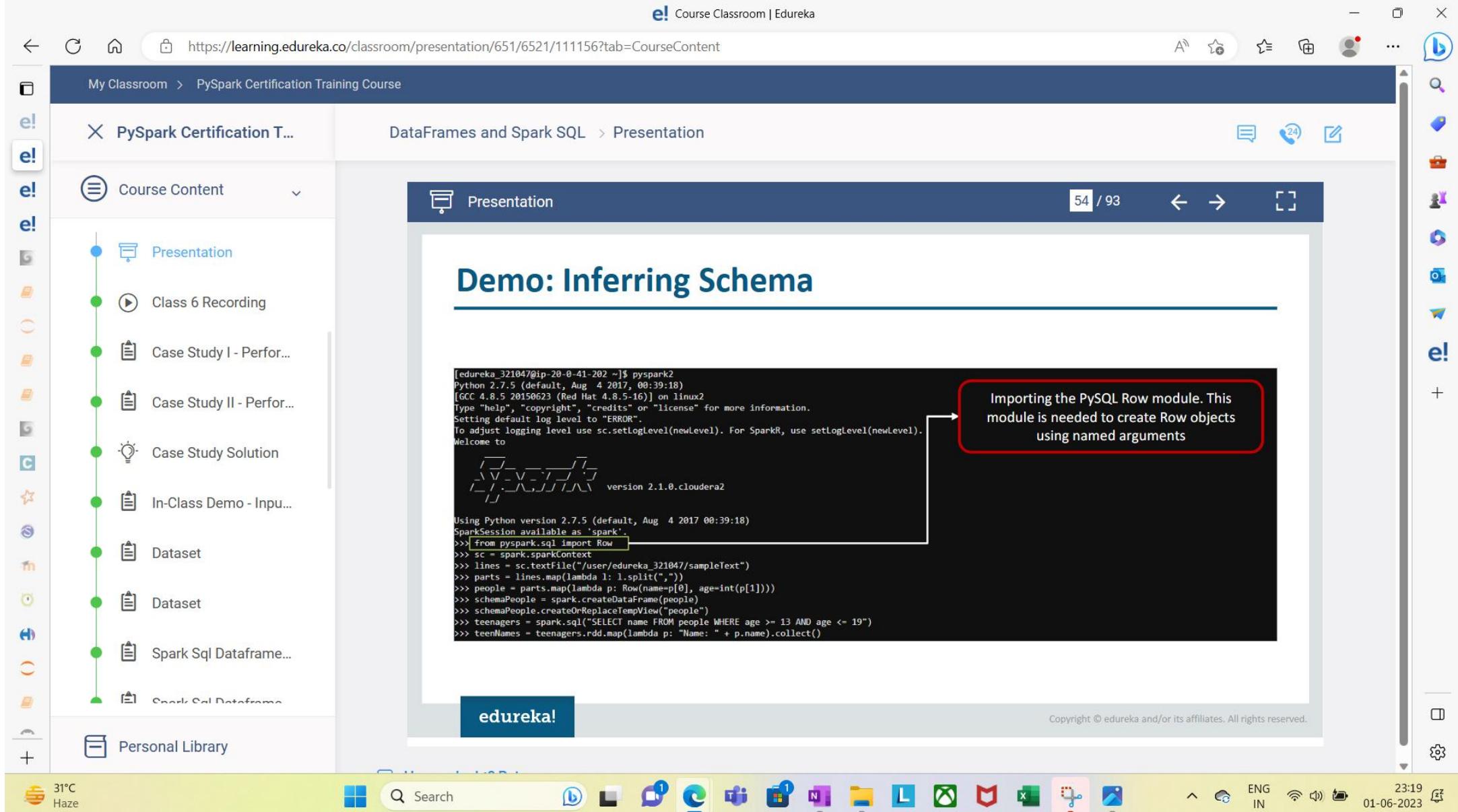
Snipping Tool

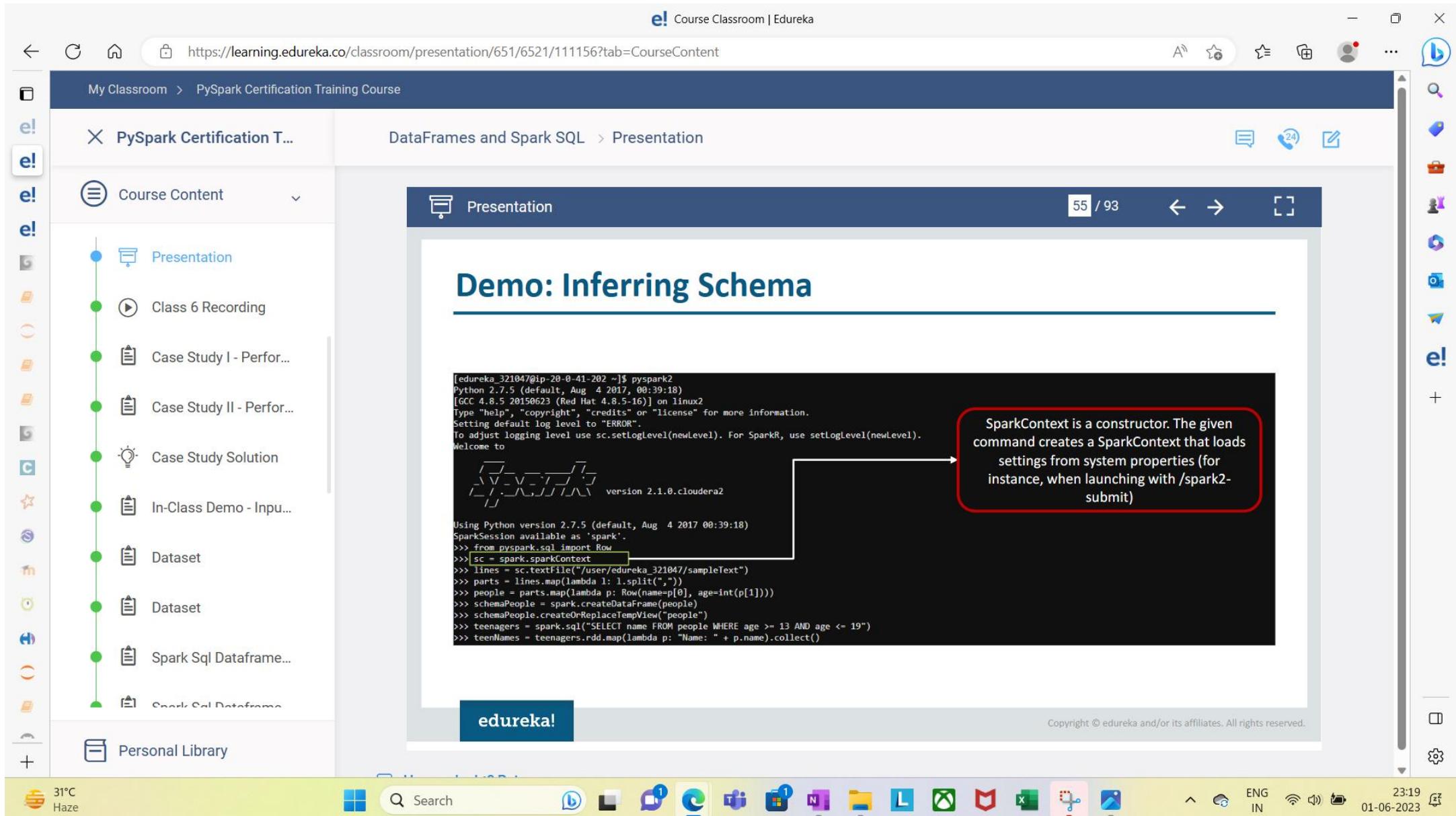
File

ENG IN

23:19

01-06-2023







My Classroom > PySpark Certification Training Course

DataFrames and Spark SQL > Presentation



 Course Content

-  Presentation
 -  Class 6 Recording
 -  Case Study I - Perform...
 -  Case Study II - Perform...
 -  Case Study Solution
 -  In-Class Demo - Inpu...
 -  Dataset
 -  Dataset
 -  Spark Sql Dataframe...
 -  Spark Sql Dataframe...

Demo: Inferring Schema

56 / 93

←

5

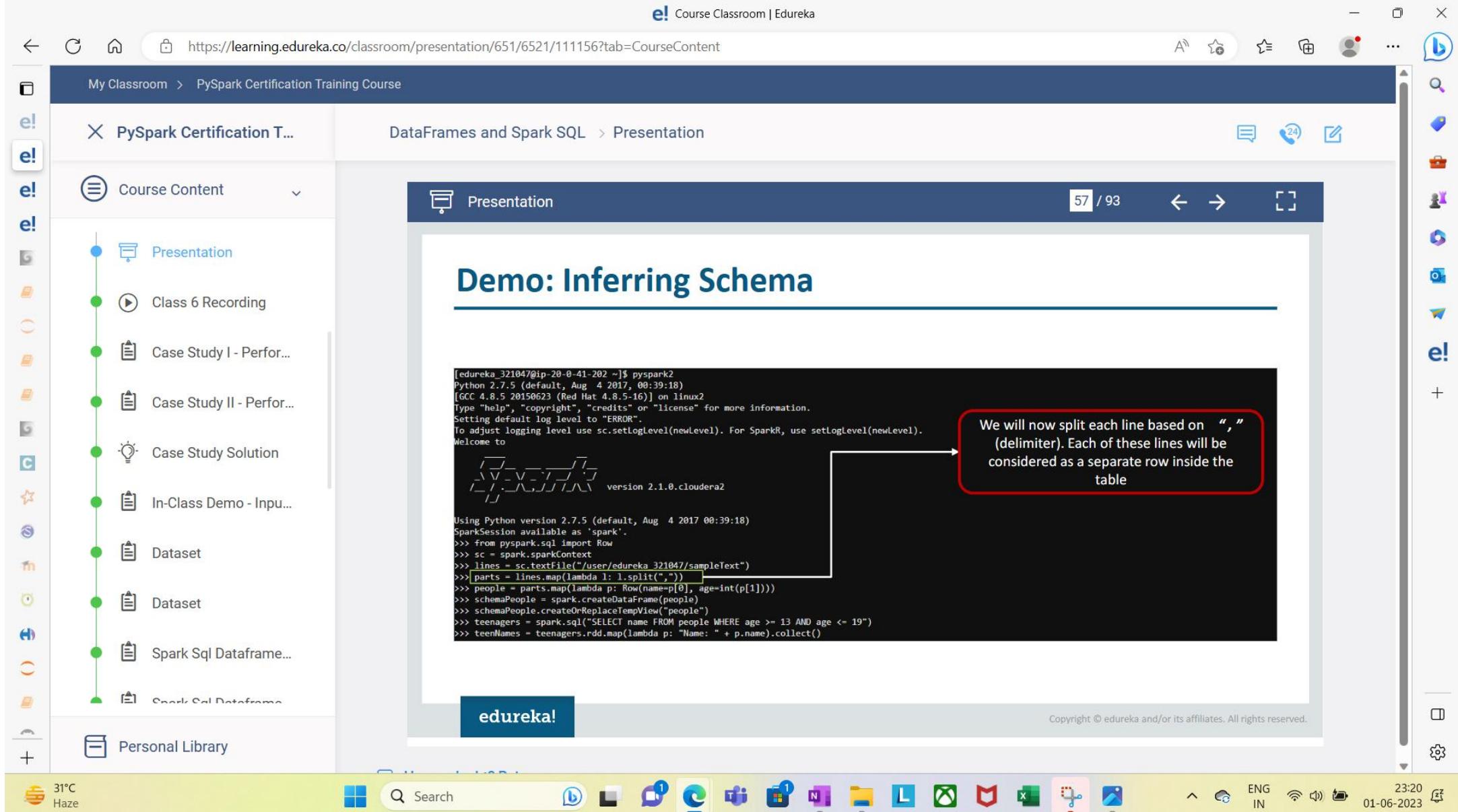
1

```
[edureka_321047@ip-20-0-41-202 ~]$ pyspark2
Python 2.7.5 (default, Aug  4 2017, 00:39:18)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-16)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "ERROR".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel().
Welcome to

   / \ \
  /  \ - V - / \ - / \
 /_ / \_\_/_/ / \_\_ \
 /_ \
version 2.1.0.cloudera2

Using Python version 2.7.5 (default, Aug  4 2017 00:39:18)
SparkSession available as 'spark'.
>>> from pyspark.sql import Row
>>> sc = spark.sparkContext
>>> lines = sc.textFile("/user/edureka_321047/sampleText")
>>> parts = lines.map(lambda l: l.split(","))
>>> people = parts.map(lambda p: Row(name=p[0], age=int(p[1])))
>>> schemaPeople = spark.createDataFrame(people)
>>> schemaPeople.createOrReplaceTempView("people")
>>> teenagers = spark.sql("SELECT name FROM people WHERE age >= 13 AND age <= 19")
>>> teenNames = teenagers.rdd.map(lambda p: "Name: " + p.name).collect()
```

We will now read our sample text file and store it in our lines RDD





My Classroom > PySpark Certification Training Course

DataFrames and Spark SQL > Presentation



 Course Content

-  Presentation
 -  Class 6 Recording
 -  Case Study I - Perform...
 -  Case Study II - Perform...
 -  Case Study Solution
 -  In-Class Demo - Inpu...
 -  Dataset
 -  Dataset
 -  Spark Sql Dataframe...
 -  Cloud Sql Dataframe...

Demo: Inferring Schema

```
[edureka_321047@ip-20-0-41-202 ~]$ pyspark2
Python 2.7.5 (default, Aug 4 2017, 00:39:18)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-16)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "ERROR".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel)
Welcome to

      / \ \
     /   \ - \ \ - / \ \
    /_ \ / . \ \ \ \ \ \ \ \
   / \ \ \ \ \ \ \ \ \ \ \ \ \
version 2.1.0.cloudera2

Using Python version 2.7.5 (default, Aug 4 2017 00:39:18)
SparkSession available as 'spark'.
>>> from pyspark.sql import Row
>>> sc = spark.sparkContext
>>> lines = sc.textFile("/user/edureka_321047/sampleText")
>>> parts = lines.map(lambda l: l.split(","))
>>> people = parts.map(lambda p: Row(name=p[0], age=int(p[1])))
>>> schemaPeople = spark.createDataFrame(people)
>>> schemaPeople.createOrReplaceTempView("people")
>>> teenagers = spark.sql("SELECT name FROM people WHERE age >= 13 AND age < 19")
>>> teenNames = teenagers.rdd.map(lambda n: "Name: " + n.name).collect()
[Name: edureka, Name: edureka, Name: edureka, Name: edureka, Name: edureka]
```

Each of the lines will now be mapped into a table with specific column names. Here we will consider the first column to be Name and second to be Age

My Classroom > PySpark Certification Training Course

DataFrames and Spark SQL > Presentation



 Course Content

-  Presentation
 -  Class 6 Recording
 -  Case Study I - Perform...
 -  Case Study II - Perform...
 -  Case Study Solution
 -  In-Class Demo - Inpu...
 -  Dataset
 -  Dataset
 -  Spark Sql Dataframe...
 -  SQL & SQL Dataframe

Demo: Inferring Schema

```
[edureka_321047@ip-20-0-41-202 ~]$ pyspark2
Python 2.7.5 (default, Aug  4 2017, 00:39:18)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-16)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "ERROR".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
Welcome to

    / \ _ \
   / \ V _ V / \ _ \
  / \ / \ _ \ / \ _ \ / \
 / \ / \ _ \ / \ _ \ / \ _ \
/ \ / \ _ \ / \ _ \ / \ _ \
version 2.1.0.cloudera2

Using Python version 2.7.5 (default, Aug  4 2017 00:39:18)
SparkSession available as 'spark'.
>>> from pyspark.sql import Row
>>> sc = spark.sparkContext
>>> lines = sc.textfile("/user/edureka_321047/sampleText")
>>> parts = lines.map(lambda l: l.split(","))
>>> people = parts.map(lambda p: Row(name=p[0], age=int(p[1])))
>>> schemaPeople = spark.createDataFrame(people)
>>> schemaPeople.createOrReplaceTempView("people")
>>> teenagers = spark.sql("SELECT name, age FROM people WHERE age >= 13 AND age <= 19")
>>> teenNames = teenagers.rdd.map(lambda n: "Name: " + n.name).collect()
Creates a DataFrame named people.
When schema is defined, column names will be inferred.
When schema is not defined, column names are inferred from the first row or named directly.
```

- Creates a DataFrame from an RDD of tuple/list, list or pandas.DataFrame
 - When schema is a list of column names, the type of each column will be inferred from data
 - When schema is None, it will try to infer the schema (column names and types) from data, which should be an RDD of Row, or namedtuple, or dict

edureka



My Classroom > PySpark Certification Training Course

DataFrames and Spark SQL > Presentation



 Course Content

-  Presentation
 -  Class 6 Recording
 -  Case Study I - Perform...
 -  Case Study II - Perform...
 -  Case Study Solution
 -  In-Class Demo - Input...
 -  Dataset
 -  Dataset
 -  Spark SQL Dataframe...
 -  Spark SQL Dataframe...

Demo: Inferring Schema

`createOrReplaceTempView` creates (or replaces if that view name already exists) a lazily evaluated "view" that you can then use like a hive table in Spark SQL. It does not persist to memory unless you cache the dataset that underpins the view.



My Classroom > PySpark Certification Training Course

DataFrames and Spark SQL > Presentation



 Course Content

-  Presentation
 -  Class 6 Recording
 -  Case Study I - Perform...
 -  Case Study II - Perform...
 -  Case Study Solution
 -  In-Class Demo - Input...
 -  Dataset
 -  Dataset
 -  Spark Sql Dataframe...
 -  Spark Sql Dataframe...

Demo: Inferring Schema

61 / 93

←

1

1

```
[edureka_321047@ip-20-0-41-202 ~]$ pyspark2
Python 2.7.5 (default, Aug  4 2017, 00:39:18)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-16)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "ERROR".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel)
Welcome to

   / \ _ \
  / \ \ V _ / \ _ / \
 / _ \ . / \ , / \ / \ \
/ \_ \_ / \_ \_ / \_ \_ / \_ \
version 2.1.0.cloudera2

Using Python version 2.7.5 (default, Aug  4 2017 00:39:18)
SparkSession available as 'spark'.
>>> from pyspark.sql import Row
>>> sc = spark.sparkContext
>>> lines = sc.textFile("/user/edureka_321047/sampleText")
>>> parts = lines.map(lambda l: l.split(","))
>>> people = parts.map(lambda p: Row(name=p[0], age=int(p[1])))
>>> schemaPeople = spark.createDataFrame(people)
>>> schemaPeople.createOrReplaceTempView("people")
>>> teenagers = spark.sql("SELECT name FROM people WHERE age >= 13 AND age <= 19")
>>> teenagersNames = teenagers.rdd.map(lambda p: "Name: " + p.name).collect()
We will run a SQL query on this DataFrame
```

We will run a SQL Query on the created table and store it inside an RDD

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 62 / 93 ← →

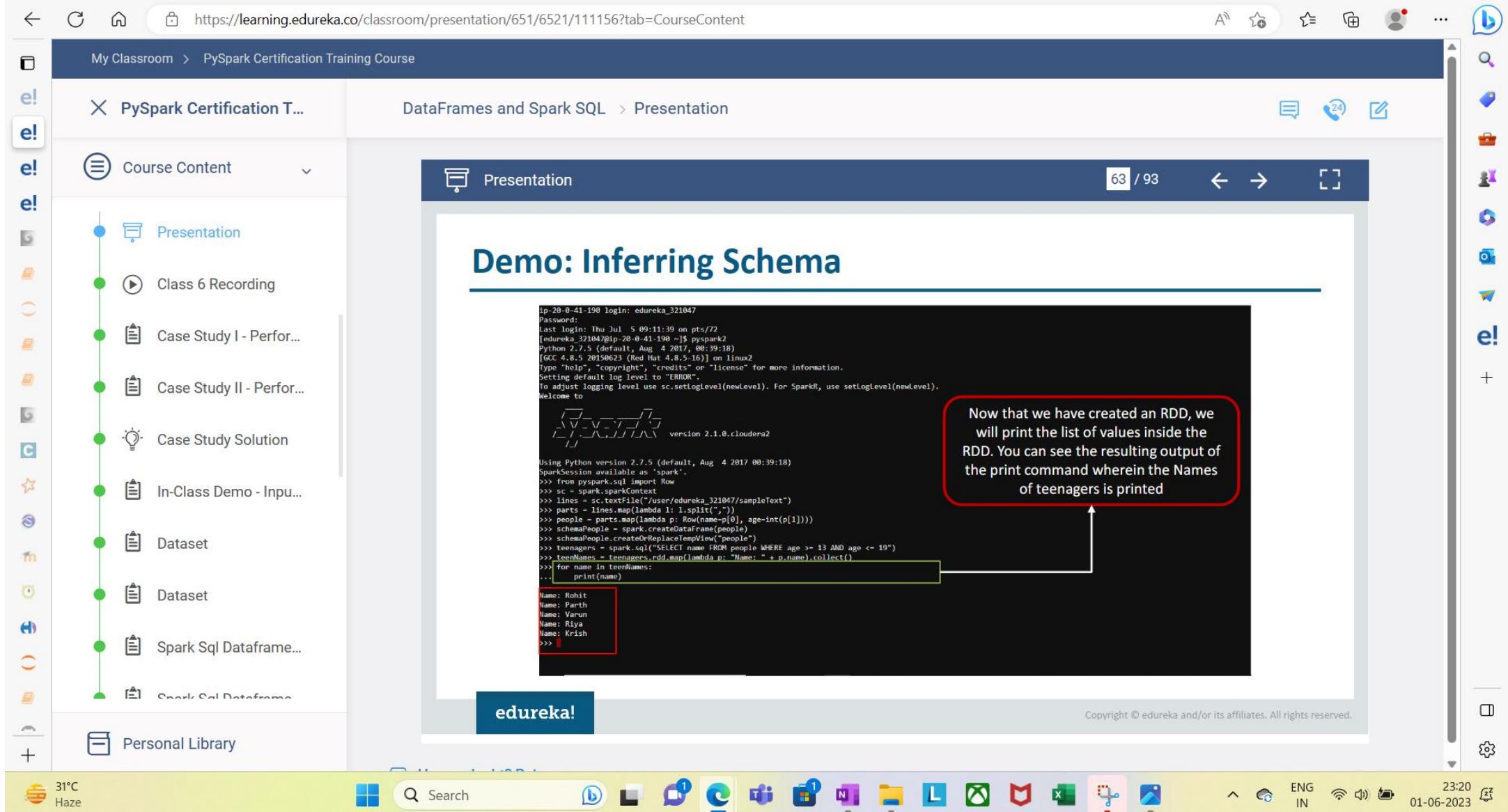
Demo: Inferring Schema

```
[edureka_321047@ip-20-0-41-202 ~]$ pyspark2
Python 2.7.5 (default, Aug 4 2017, 00:39:18)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-16)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "ERROR".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
Welcome to
   __| |
  / \ \_\_V\_/\_/\_/\_/\_
 / \ / .\_\_/\_/\_/\_/\_/
 / \_/
Using Python version 2.7.5 (default, Aug 4 2017 00:39:18)
SparkSession available as 'spark'.
>>> from pyspark.sql import Row
>>> sc = spark.sparkContext
>>> lines = sc.textfile("/user/edureka_321047/sampleText")
>>> parts = lines.map(lambda l: l.split(","))
>>> people = parts.map(lambda p: Row(name=p[0], age=int(p[1])))
>>> schemaPeople = spark.createDataFrame(people)
>>> schemaPeople.createOrReplaceTempView("people")
>>> teenagers = spark.sql("SELECT name FROM people WHERE age >= 13 AND age <= 19")
>>> teenNames = teenagers.rdd.map(lambda p: "Name: " + p.name).collect()
```

You can then map the RDD according to the type of value you want to fetch from it. Here we are fetching the Name of teenagers from the RDD and collect it inside another RDD

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.



My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Personal Library

Presentation

Specifying Schema Programmatically

When case classes cannot be defined ahead of time (for example, the structure of records is encoded in a string, or a text dataset will be parsed and fields will be projected differently for different users), a DataFrame can be created programmatically with three steps

1. Create an RDD of Rows from the original RDD
2. Create the schema represented by a StructType matching the structure of Rows in the RDD created in Step 1
3. Apply the schema to the RDD of Rows via createDataFrame method provided by SparkSession



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation

66 / 93

Demo: Inferring Schema

John, 20
Rohit, 15
Parth, 14
Rishabh, 25
Daisy, 34
Annie, 23
Sushmita, 50
Kaivalya, 26
Varun, 16
Shambhavi, 21
Johnson, 22
Riya, 17
Krish, 19
Akanksha, 28
Rutuja, 33

Consider we have a text file of the following format. Please note that here we are considering that schema of the file is already known and is stored in a TXT format (Non-Columnar)

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

Cloud

File Explorer

OneDrive

PowerPoint

Word

Excel

Teams

Outlook

OneNote

Power BI

Xbox

Bookmarks

Calculator

Snipping Tool

File

ENG IN

23:20 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 68 / 93 ← →

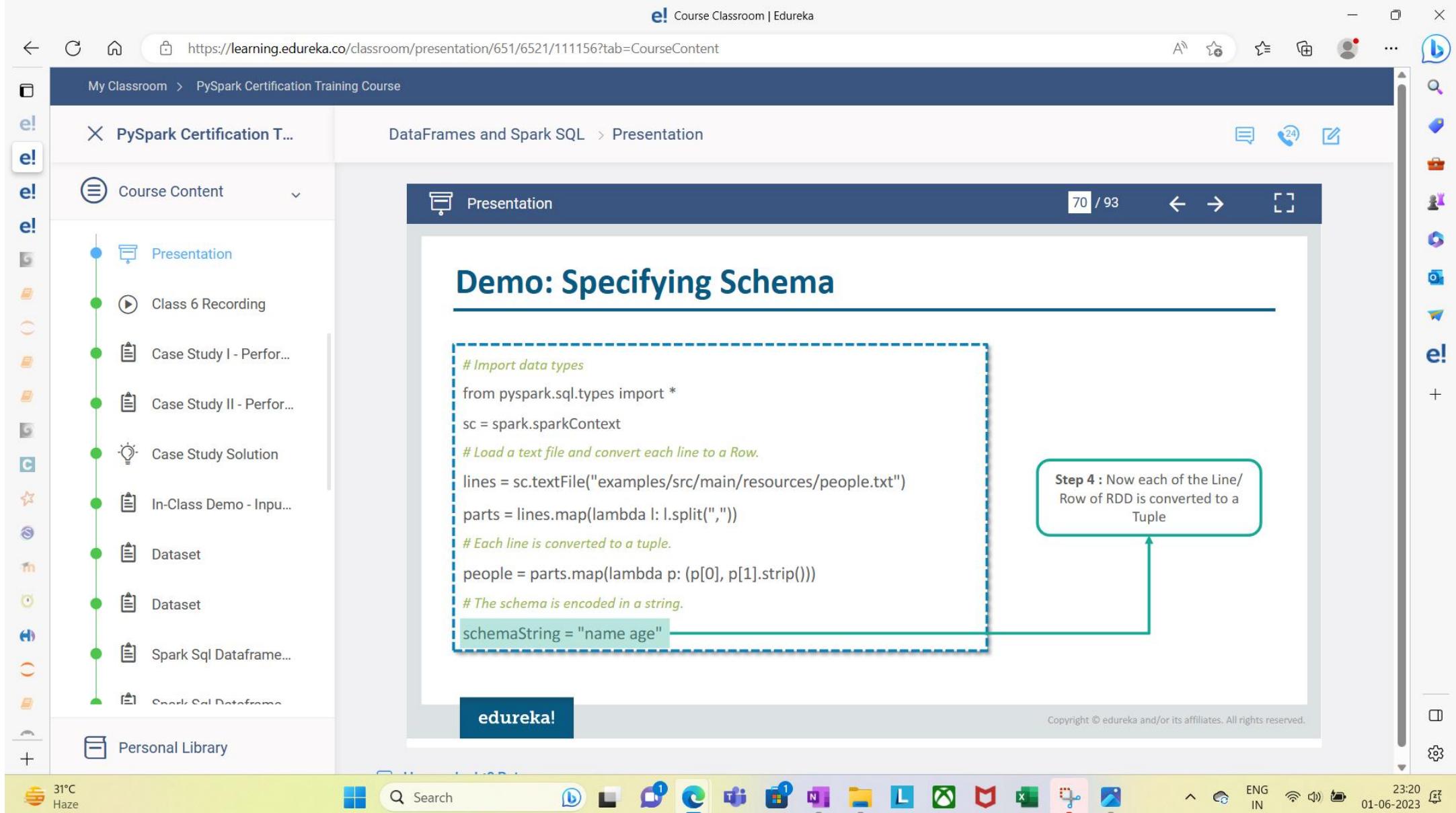
Demo: Specifying Schema

```
# Import data types
from pyspark.sql.types import *
sc = spark.sparkContext
# Load a text file and convert each line to a Row.
lines = sc.textFile("examples/src/main/resources/people.txt")
parts = lines.map(lambda l: l.split(","))
# Each line is converted to a tuple.
people = parts.map(lambda p: (p[0], p[1].strip()))
# The schema is encoded in a string.
schemaString = "name age"
```

Step 2 : Load the necessary file and convert each of these lines into a row of RDD

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.



Course Content

Presentation

Class 6 Recording

Case Study I - Perform...

Case Study II - Perform...

Case Study Solution

In-Class Demo - Inpu...

Dataset

Dataset

Spark Sql Dataframe...

Spark Sql Dataframe...

Personal Library

Demo: Specifying Schema

```
fields = [StructField(field_name, StringType(), True) for field_name in  
schemaString.split()]  
  
schema = StructType(fields)  
  
# Apply the schema to the RDD.  
  
schemaPeople = spark.createDataFrame(people, schema)  
  
# Creates a temporary view using the DataFrame  
  
schemaPeople.createOrReplaceTempView("people")
```

Step 5 : We will now generate a schema based on the Schema String we created in previous step

Copyright © edureka and/or its affiliates. All rights reserved.

edureka!

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Presentation 72 / 93 ← →

Demo: Specifying Schema

```
fields = [StructField(field_name, StringType(), True) for field_name in schemaString.split()]
schema = StructType(fields)
# Apply the schema to the RDD.
schemaPeople = spark.createDataFrame(people, schema)
# Creates a temporary view using the DataFrame
schemaPeople.createOrReplaceTempView("people")
```

Step 6 : Once the RDD is created based on schema we will create a Data Frame

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

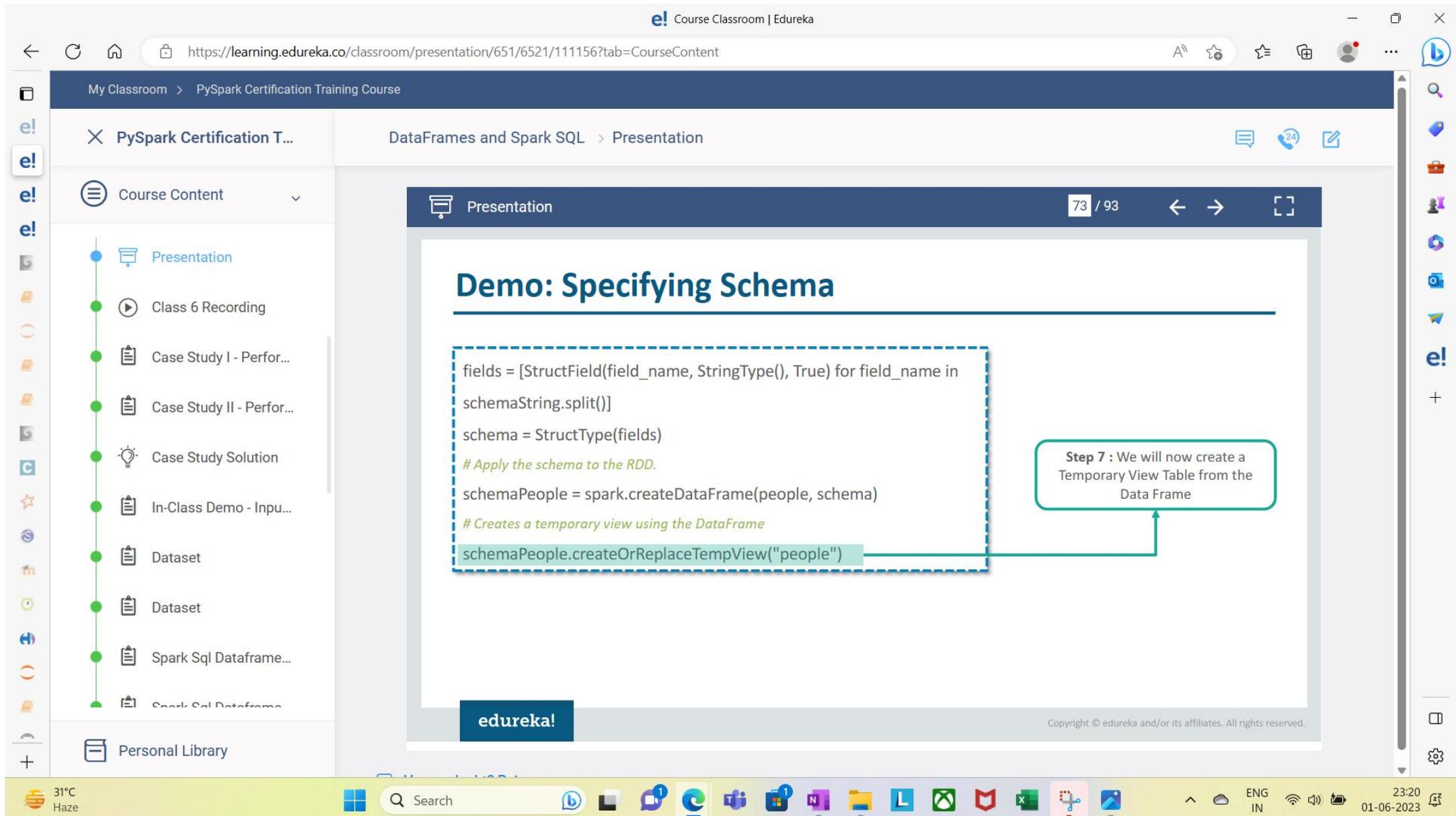
31°C Haze

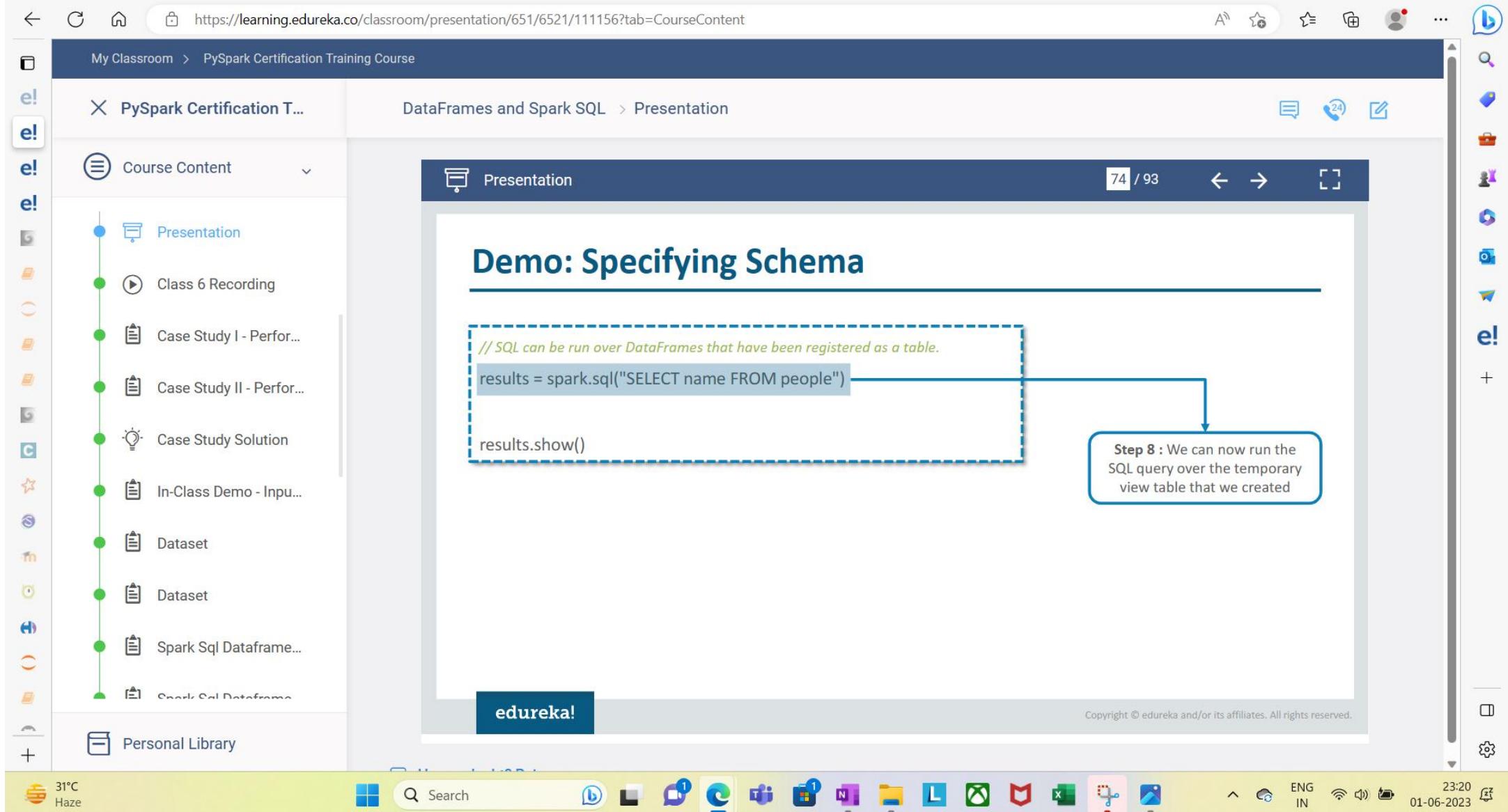
Search

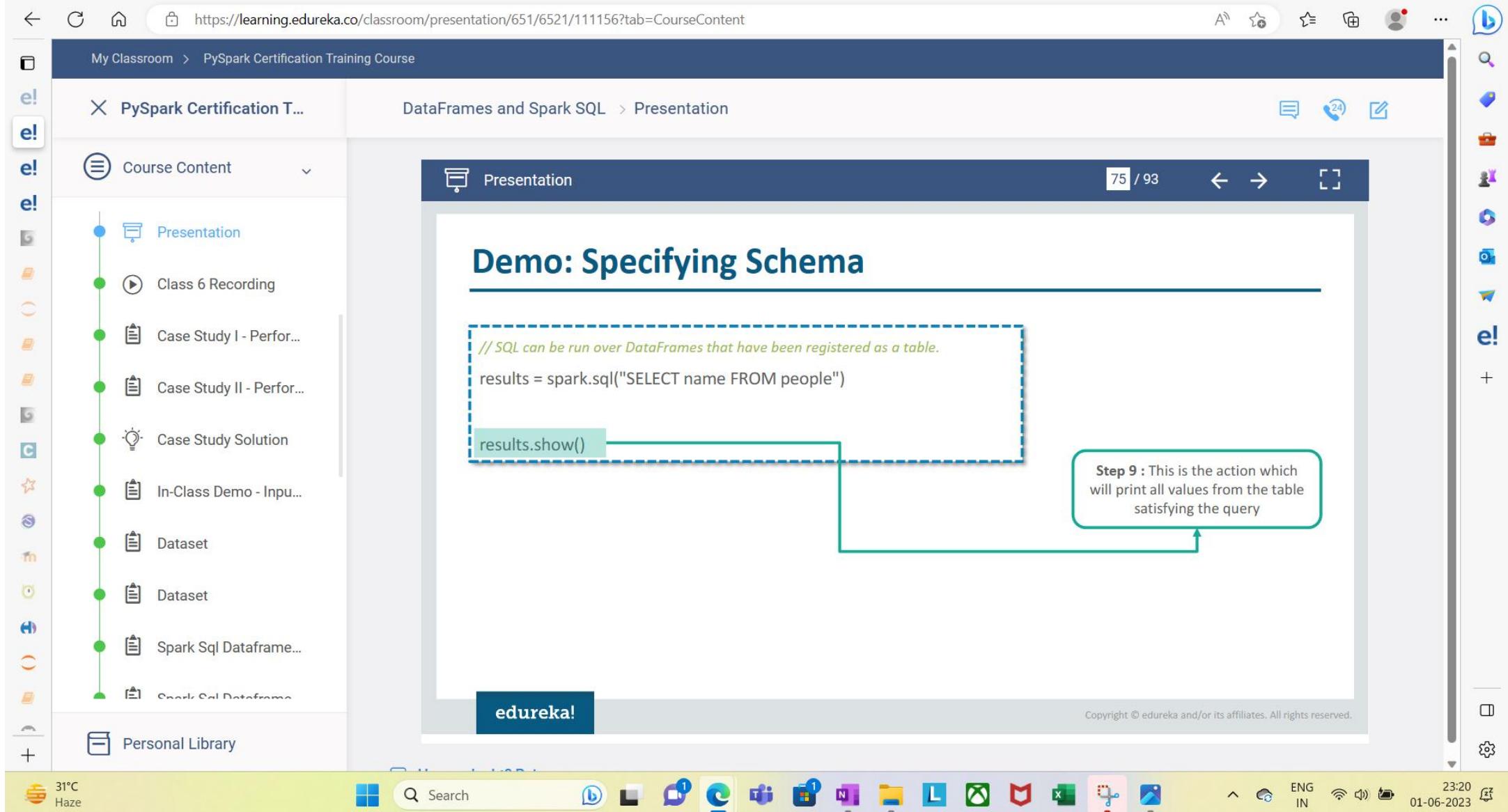
L

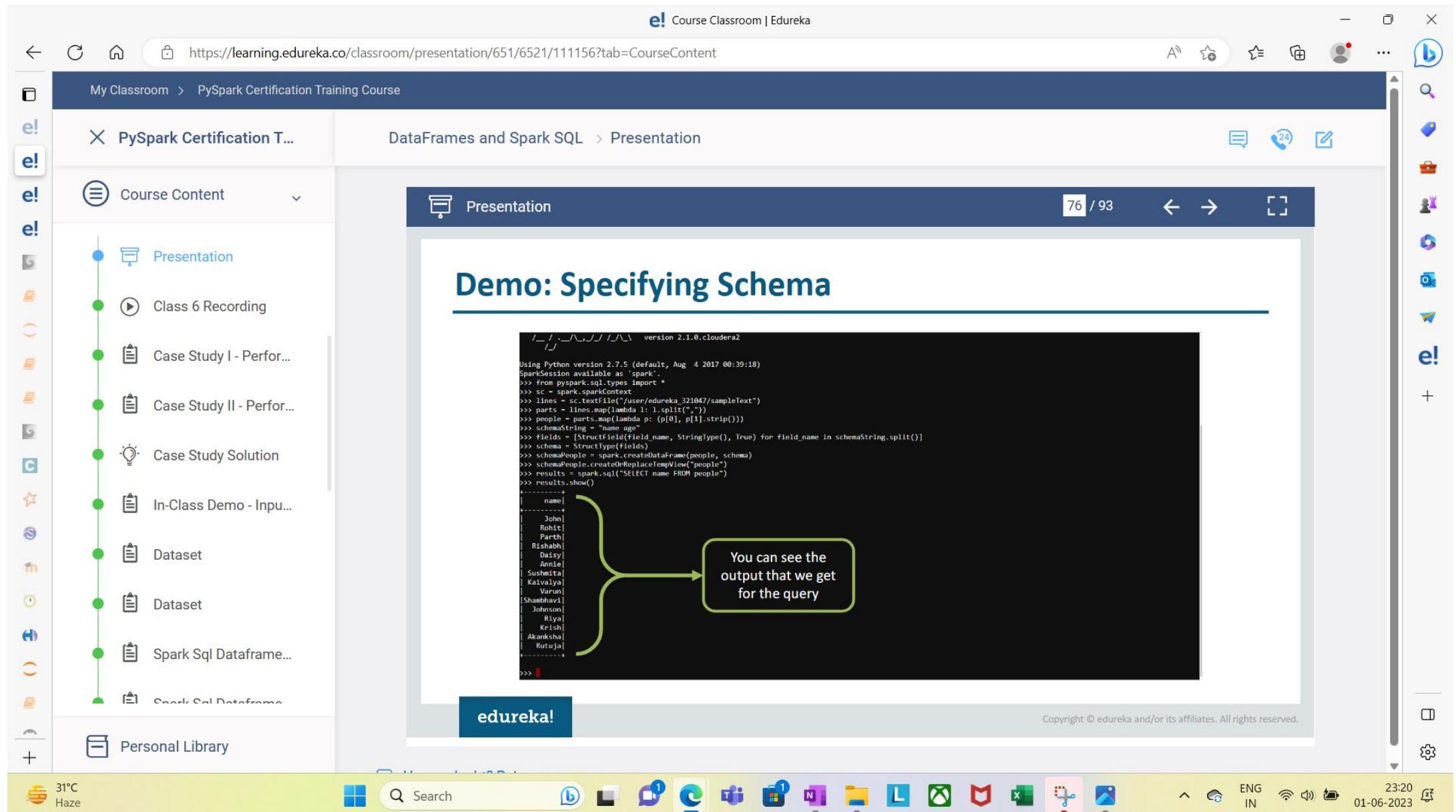
ENG IN

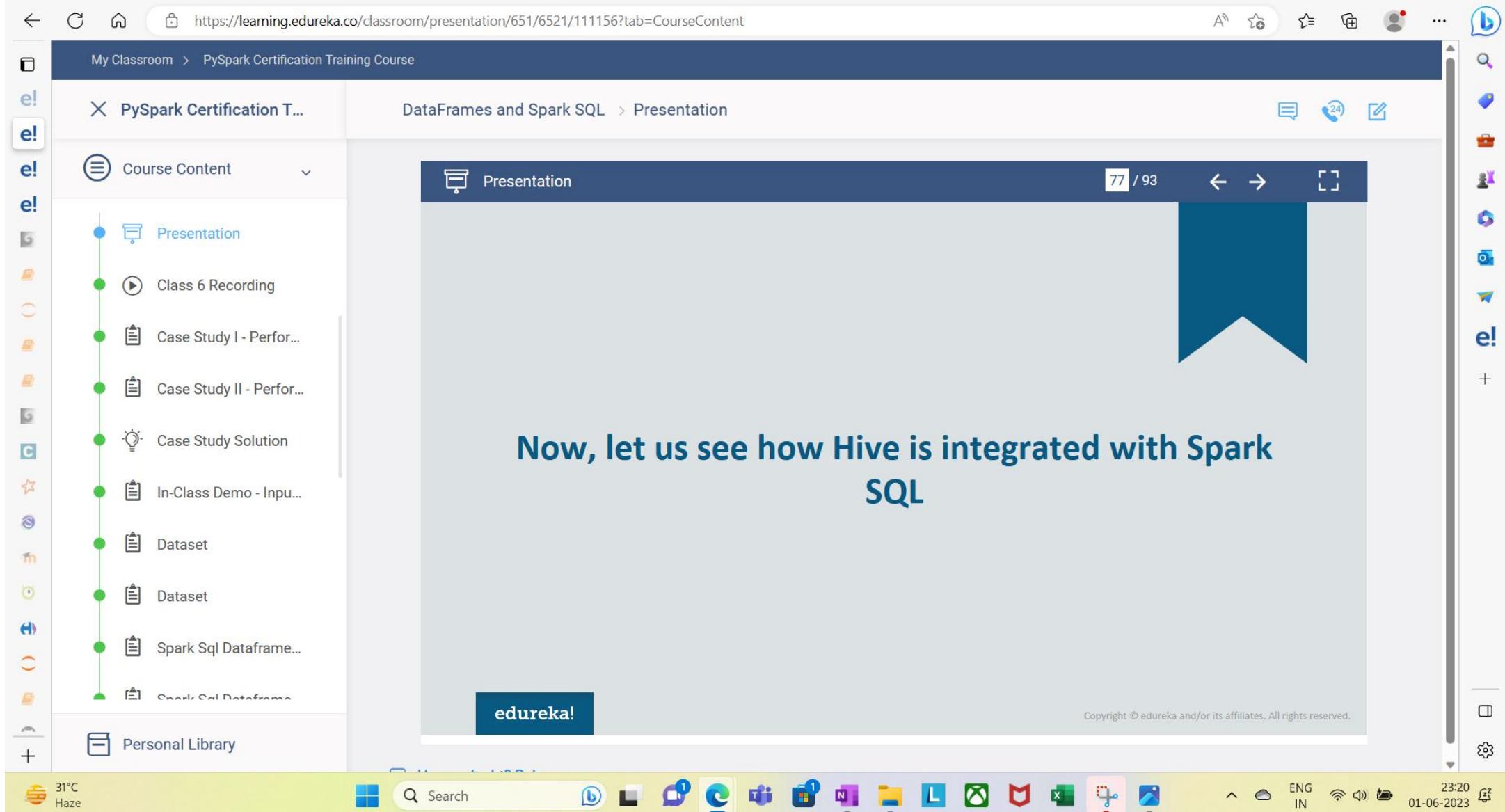
23:20 01-06-2023

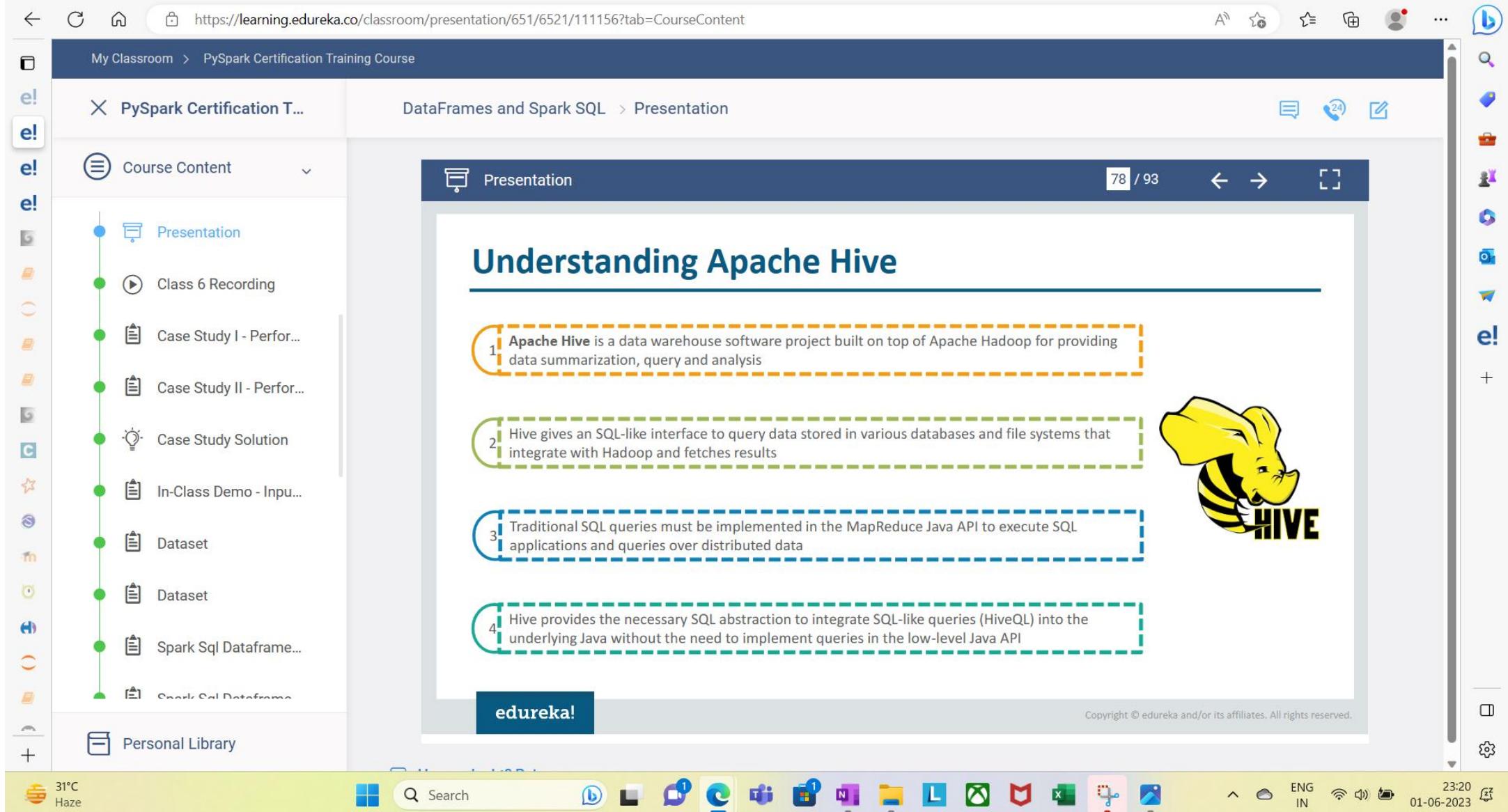












Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Personal Library

Features of Apache Hive

Indexing to provide acceleration, index type including compaction and bitmap index as of 0.10, more index types are planned



Metadata storage in a relational database management system, significantly reducing the time to perform semantic checks during query execution.

Operating on compressed data stored into the Hadoop ecosystem using algorithms including DEFLATE, BWT, snappy, etc

Built-in user-defined functions (UDFs) to manipulate dates, strings, and other data-mining tools. Hive supports extending the UDF set to handle use-cases not supported by built-in functions

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

DataFrames and Spark SQL > Presentation

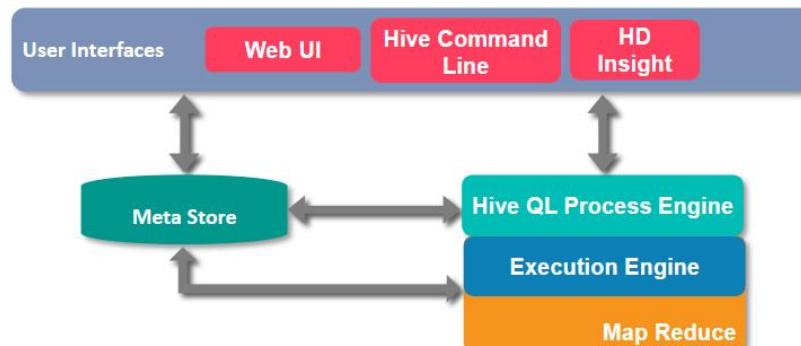


Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Personal Library

Hive Architecture



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Presentation

81 / 93

Demo 3 – Integrating Apache Hive and Spark

Refer to the file Module-6 Demo 3 provided in the LMS for all the steps in detail

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

Cloud

OneDrive

Microsoft Edge

Teams

OneNote

File Explorer

PowerPoint

Xbox

Bookmarks

Excel

Snipping Tool

Calculator

23:20 01-06-2023

https://learning.edureka.co/course/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:20 01-06-2023

X PySpark Certification T...

e! Course Content

Presentation

Class 6 Recording

Case Study I - Perform...

Case Study II - Perform...

Case Study Solution

In-Class Demo - Inpu...

Dataset

Dataset

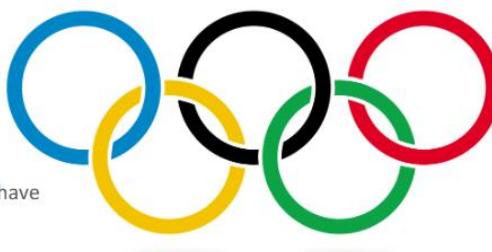
Spark Sql Dataframe...

Spark Sql Dataframe...

e! Personal Library

Use Case: Problem Statement

- Analysis on sports data is very important
- It provides relevant insights into the data and helps us garner better results out of it
- One such important event is the Olympics Games
- The Olympic Games are held every four years, with the Summer and Winter Games alternating by occurring every four years but two years apart
- You will have to perform multiple computations on a given given set of data and provide insights on it
- We will now discuss the outputs that you have to provide from given Dataset and have a look at the kind of data we have



edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

Personal Library

Presentation

Use Case: Data Overview

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	ID	Name	Sex	Age	Height	Weight	Team	NOC	Games	Year	Season	City	Sport	Event	Medal
2	1	A Dijiang	M	24	180	80	China	CHN	1992 Summer	1992	Summer	Barcelona	Basketball	III Men's Br.	NA
3	2	A Lamusi	M	23	170	60	China	CHN	2012 Summer	2012	Summer	London	Judo	's Extra-Lig	NA
4	3	Car Nielsen	M	24	NA	NA	Denmark	DEN	2020 Summer	1920	Summer	Antwerpen	Football	III Men's Fo	NA
5	4	Lindenau	M	34	NA	NA	mark/Swe	DEN	2000 Summer	1900	Summer	Paris	Tug-Of-War	Men's Tu	Gold
6	5	Hee Jacoba	F	21	185	82	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Women's Sp	NA
7	5	Hee Jacoba	F	21	185	82	Netherlands	NED	1988 Winter	1988	Winter	Calgary	Speed Skating	Women's Sp	NA

You can download the dataset from [here](#)

Given is the CSV file containing details of Players, Teams, Sports etc. You have to perform the analysis on this data and provide relevant outputs. Let us now discuss the various operations that you need to perform over this data set

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

DataFrames and Spark SQL > Presentation



Course Content

Presentation

Class 6 Recording

Case Study I - Perform...

Case Study II - Perform...

Case Study Solution

In-Class Demo - Inpu...

Dataset

Dataset

Spark Sql Dataframe...

Spark Sql Dataframe...

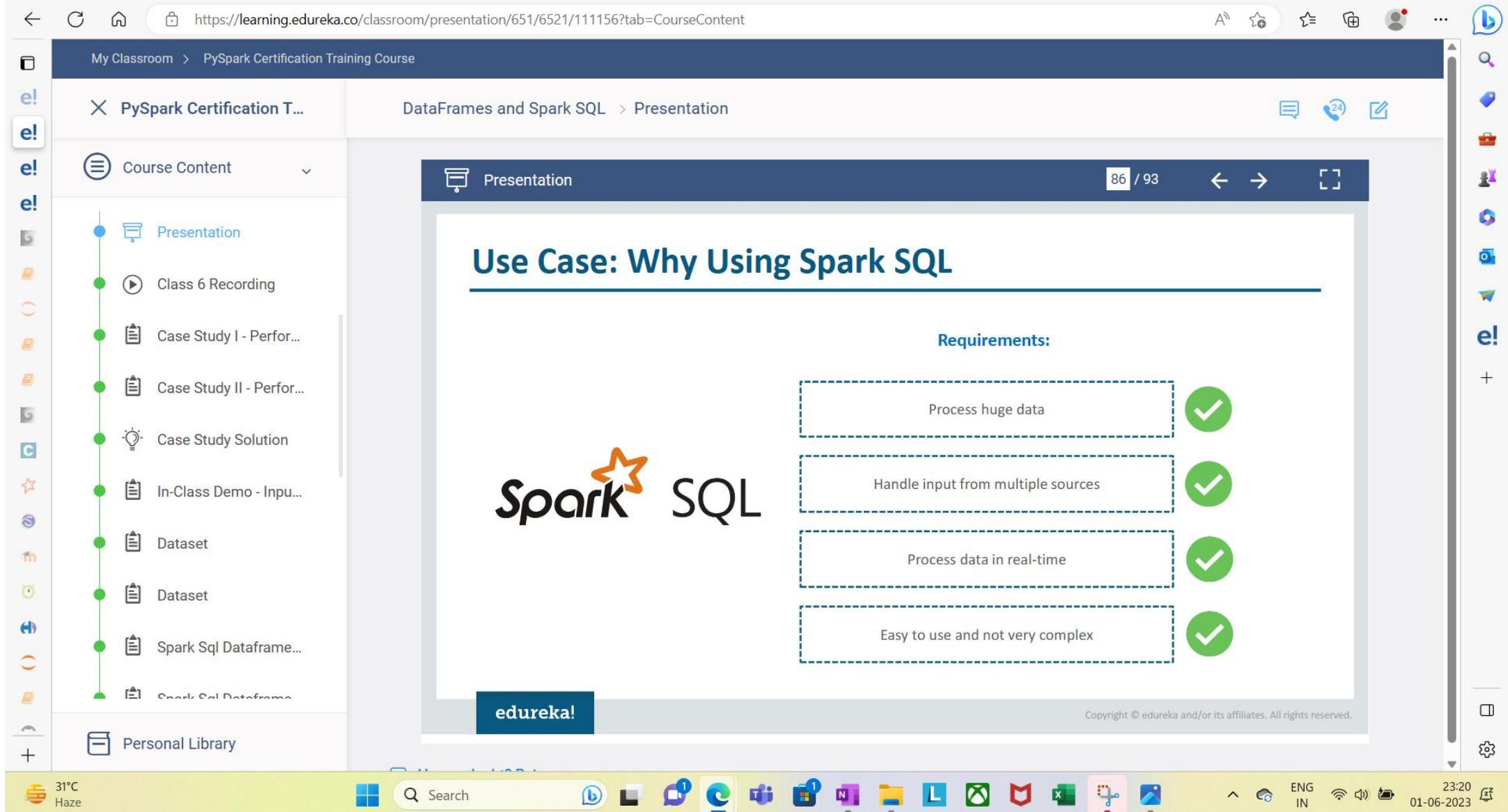
Personal Library

Use Case: Computational Requirements

- 1 Find the List of Male players from United States who Play Basket Ball and there age is greater than 28
- 2 Find the list of All Players who have Played winter Olympics after 1952 in the Athletics Women's High Jump
- 3 Find Name, Age, Team of all players from Denmark who have played Summer Olympics in Rio De Janeiro
- 4 Find Name and Age of all the players who have played Football as a Sport from United States, France and Uruguay
- 5 Find the Name and NOC of players who have played the 2012 Summer and 2006 Winter Olympics

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.



https://learning.edureka.co/classroom/presentation/651/6521/111156?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... DataFrames and Spark SQL > Presentation

Course Content

- Presentation
- Class 6 Recording
- Case Study I - Perform...
- Case Study II - Perform...
- Case Study Solution
- In-Class Demo - Inpu...
- Dataset
- Dataset
- Spark Sql Dataframe...
- Spark Sql Dataframe...

edureka!

Copyright © edureka and/or its affiliates. All rights reserved.

31°C Haze

Search

L

ENG IN

23:20 01-06-2023

