

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 1 / 121 ← →

How to make the best use of Live Sessions

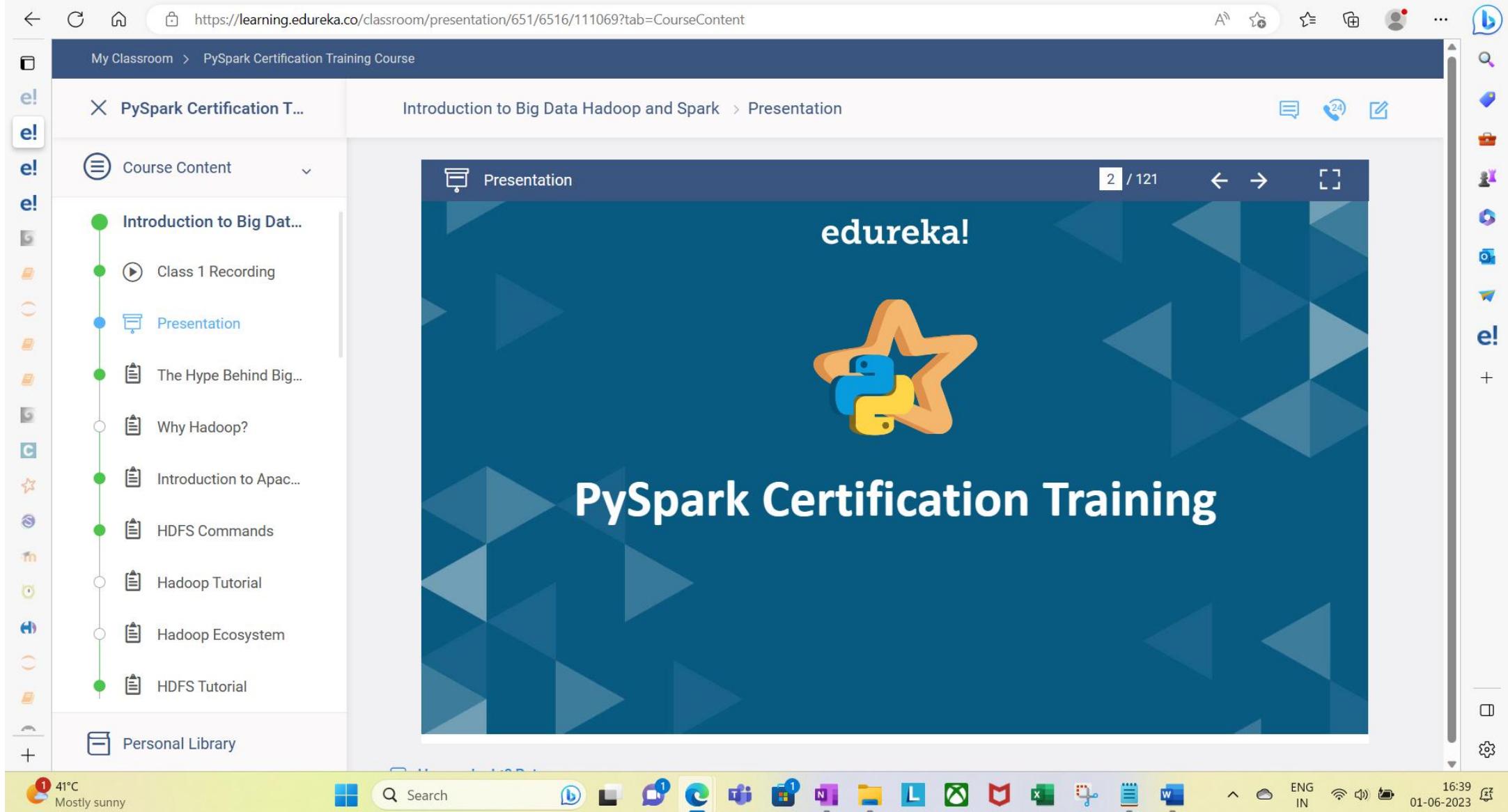
- Please login on time
- Please do a check on your network connection and audio before the class to have a smooth session
- All participants will be on mute, by default. You will be unmuted when requested or as needed
- Please use the “Questions” panel on your webinar tool to interact with the instructor at any point during the class
- Ask and answer questions to make your learning interactive
- Please have the support phone number (US : 1855 818 0063 (toll free), India : +91 90191 17772) and raise tickets from LMS in case of any issues with the tool
- Most often logging off or rejoining will help solve the tool related issues

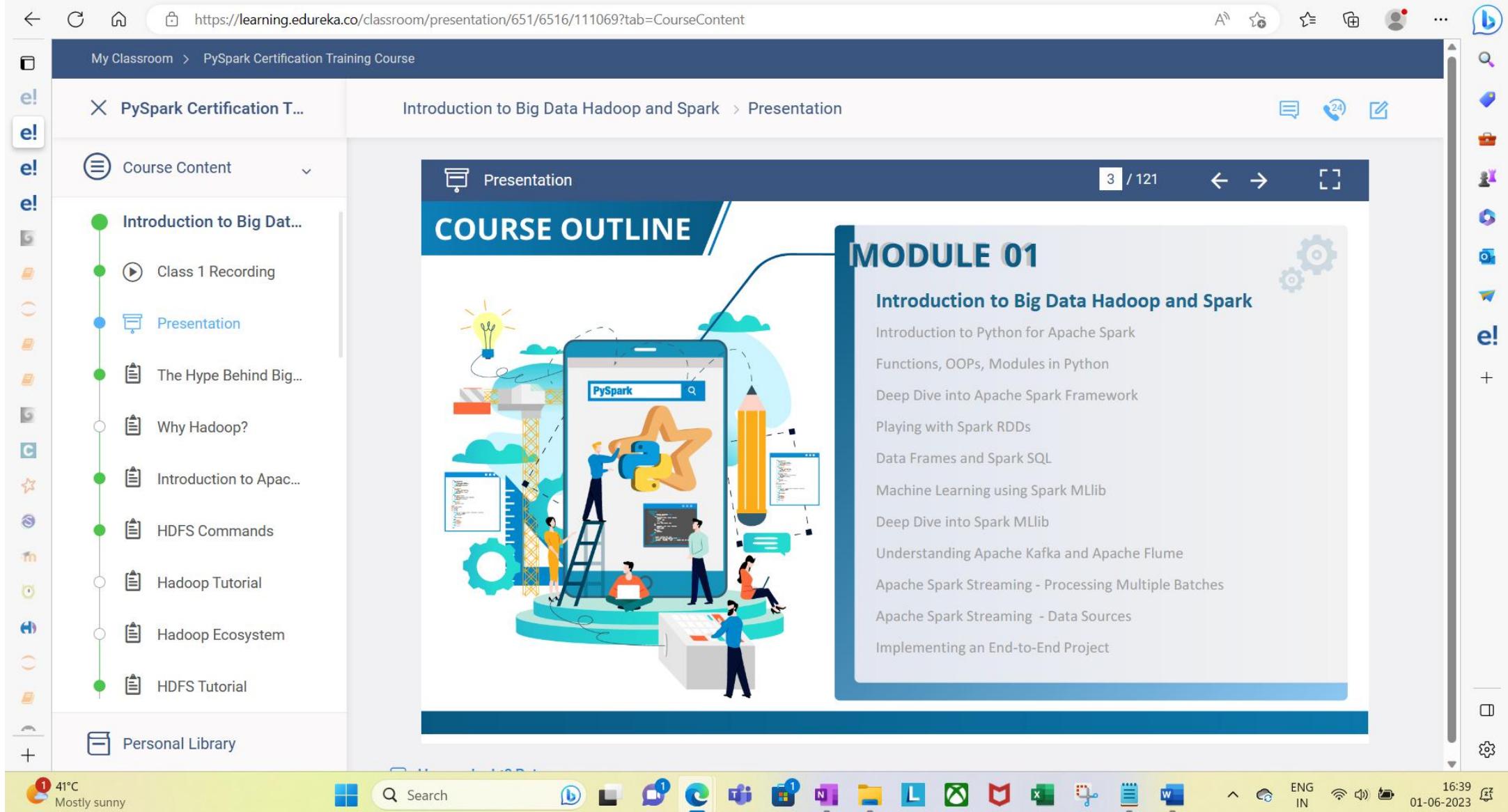
edureka!

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

41°C Mostly sunny Search

16:39 01-06-2023





X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

4 / 121



Objectives

After completing this module, you should be able to:

- Evaluate present Structure of Data and problems
- Describe Big Data and it's domain
- Explain Hadoop Core Components and HDFS
- Understand Replication and Rack Awareness
- Explain YARN and it's components
- Differentiate between Big Data Analytics with Batch and Real-Time Processing
- Understand the need for Spark and explain what is Spark
- Analyse how Spark differs from its competitors
- Recognize Spark's place in Hadoop ecosystem



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

edureka!



My Classroom > PySpark Certification Training Course



X PySpark Certification T...



Course Content



Introduction to Big Dat...



Class 1 Recording



Presentation



The Hype Behind Big...



Why Hadoop?



Introduction to Apac...



HDFS Commands



Hadoop Tutorial



Hadoop Ecosystem



HDFS Tutorial



Personal Library

Introduction to Big Data Hadoop and Spark > Presentation



_PRESENTATION



Presentation

5 / 121



A Tour of Edureka's Cloud Lab

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 6 / 121 ← →

LMS Overview

Once you get the LMS access, you will have this Screen. To login into Cloud Lab, click on "My Lab" option

Getting Started
Pre-recorded Classes
Course Content
Class Recordings
Personal Library
My Lab

Guide to Access Cloud Lab

Cloud Lab Guide - Hue
Edureka Cloud Lab Guide
Cloud Lab Guide

Your Account Balance
₹ 0 edureka cash
REDEEM
Earn upto 3000 credits on every referral!
enter email id's (separated by commas) →
Connect to Learning Manager
Select the time to get a call

edureka!

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

1 41°C Mostly sunny Search ENG IN 16:39 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apache...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial
- Personal Library

7 / 121 ← →

My Lab Contents

After clicking on "My Lab" option, this screen will pop up. Here, you will get the credentials to login into Cloud Lab

Lab

Getting Started

Your Cloud lab is ready for use!

Click on the buttons below to access respective services:

Cloudera Manager Webconsole Hue FTP Jupyter
Namenode Web UI Mapreduce Job ... Spark2.1 History ...

Username: [REDACTED] Copy
Password: [REDACTED] Copy

How to Use Edureka's LMS

You earn an industry trusted course certificate once you complete the course

Preview

Your Account Balance ₹ 0 edureka cash REDEEM

Earn upto 3000 credits on every referral!
Enter email Id's (separated by commas) →

IN G+ TWITTER FACEBOOK

Here, you will find all the services required to work on PySpark

edureka!

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

41°C Mostly sunny

Search

Cloud Lab

L

ENG IN

16:39 01-06-2023

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apache...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

8 / 121

WebConsole

WebConsole is terminal. You have to copy-paste the credentials from the previous window to this window, to login into WebConsole

```
Not secure | bdlabs.edureka.co:50002
ip-20-0-41-93 login: edureka_2
Password:
Last login: Tue Jan 9 09:13:56 on pts/15
[edureka_253770@ip-20-0-41-93 ~]$
```

edureka!

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

41°C Mostly sunny

Search

16:39 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apache...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

HUE

To see all the files present in HDFS, click on "HDFS Browser"

HUE Query Editors Workflows

edureka!

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

41°C Mostly sunny

Search

16:39 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apache...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

FTP Server

Using FTP you can upload or download any file/ folder from your local system to cloud lab

Name	Size	Date	Time
Code		09/01/18	11:32
Codes		09/01/18	09:19
FileContentCount		09/01/18	11:24

edureka!

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

Host: localhost User: edureka_25370 Upload Limit: 1GB

41°C Mostly sunny

Search

16:39 IN 01-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apache...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

11 / 121 ← → []

Jupyter Notebook

Jupyter Notebook is the IDE used to run all the PySpark Programs

Jupyter

Logout Control Panel

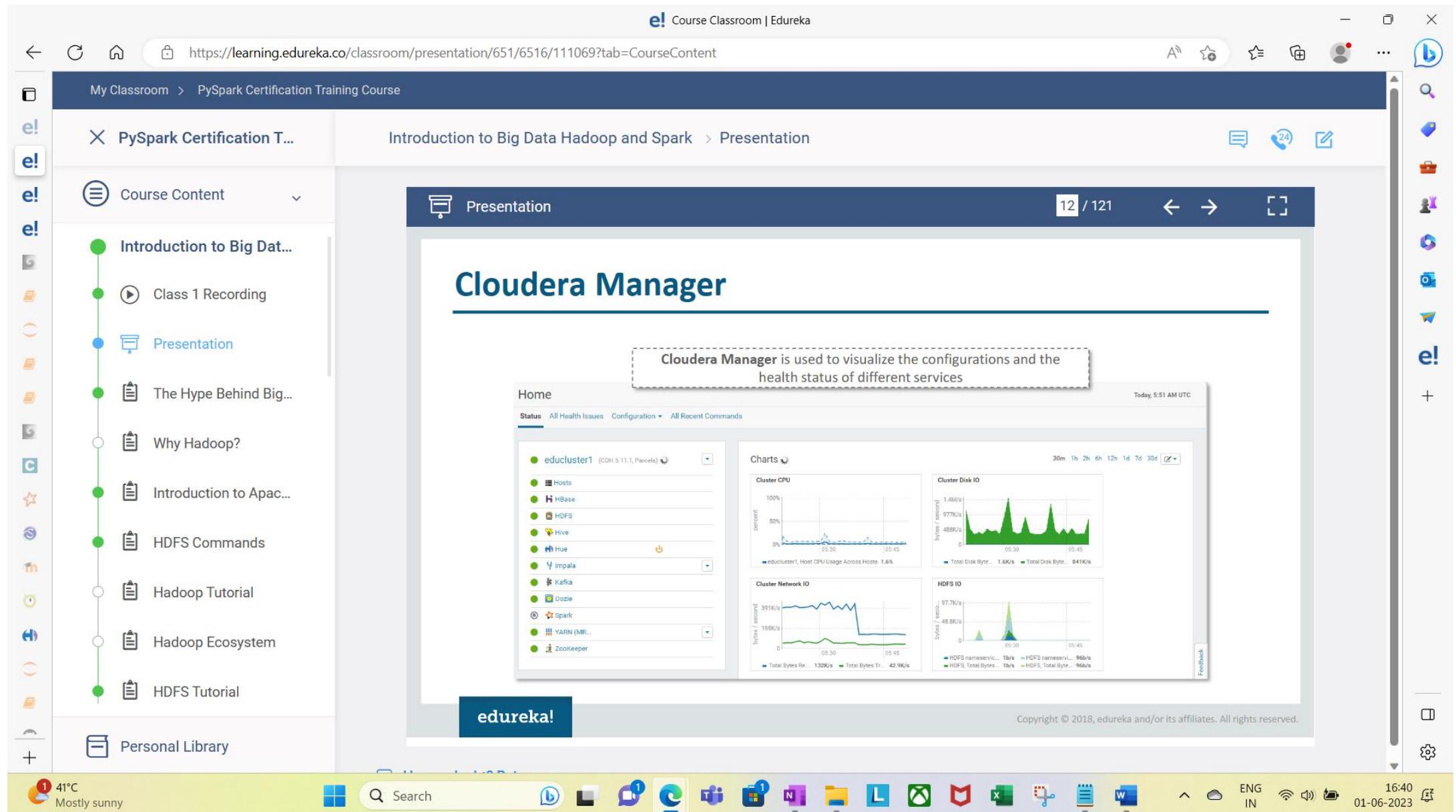
Files Running Clusters

Select items to perform actions on them.

	Name	Last Modified
0	Code	2 days ago
	Codes	2 days ago
	FileContentCount	2 days ago

edureka!

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



https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apache...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

13 / 121

Let us begin with the Structure of Data

edureka!

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

41°C Mostly sunny

Search

16:40 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 14 / 121 ← →

Structure of Data Today

The diagram shows the classification of data into four categories based on structure:

- Structured:** Data containing a defined data type, format, structure. Example: Data from RDBMS.
- Semi-Structured:** Textual data files with a noticeable pattern, enabling parsing. Example: XML data files that are self-describing and defined by an XML schema.
- Quasi-Structured:** Textual data with erratic data formats, can be formatted with effort, tools, and time. Example: Web clickstream data that may contain some inconsistencies in data values and formats.
- Unstructured:** Data that has no inherent structure and is usually stored as different types of files. Example: Images and videos.

edureka!

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

41°C Mostly sunny

Search

16:40 01-06-2023

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 15 / 121 ← →

Unstructured Data is exploding

Data-Distribution

Year	Structured Data (EB)	Un-structured Data (EB)
2007	~1000	~1000
2008	~1500	~1500
2009	~2000	~2000
2010	~2500	~2500
2011	~3000	~3000
2012	~3500	~3500
2013	~4000	~4000
2014	~4500	~4500
2015	~5000	~5000
2016	~5500	~5500
2017	~6000	~6000
2018	~6500	~6500
2019	~7000	~7000
2020	~7500	~7500

This accelerated growth of **UNSTRUCTURED** data is what turned into **BIG DATA**

By 2020, IDC predicted, Unstructured Data to have reached **40 Zettabytes (ZB)**
By 2020, there will be **5,200 GB** of data for every person on Earth

edureka!

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

41°C Mostly sunny

Search

16:40 01-06-2023

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apache...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

16 / 121

Now, let us define what is Big Data

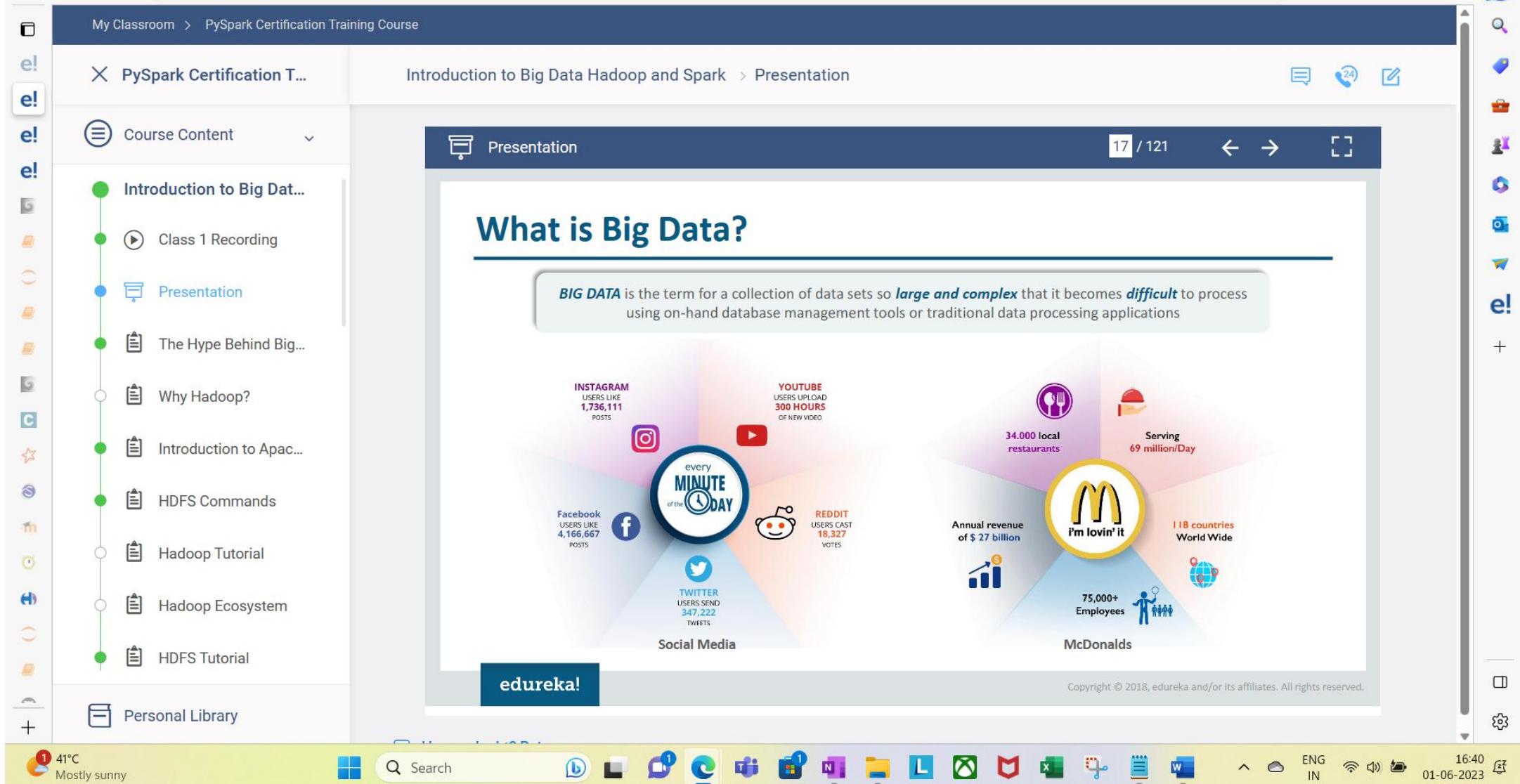
edureka!

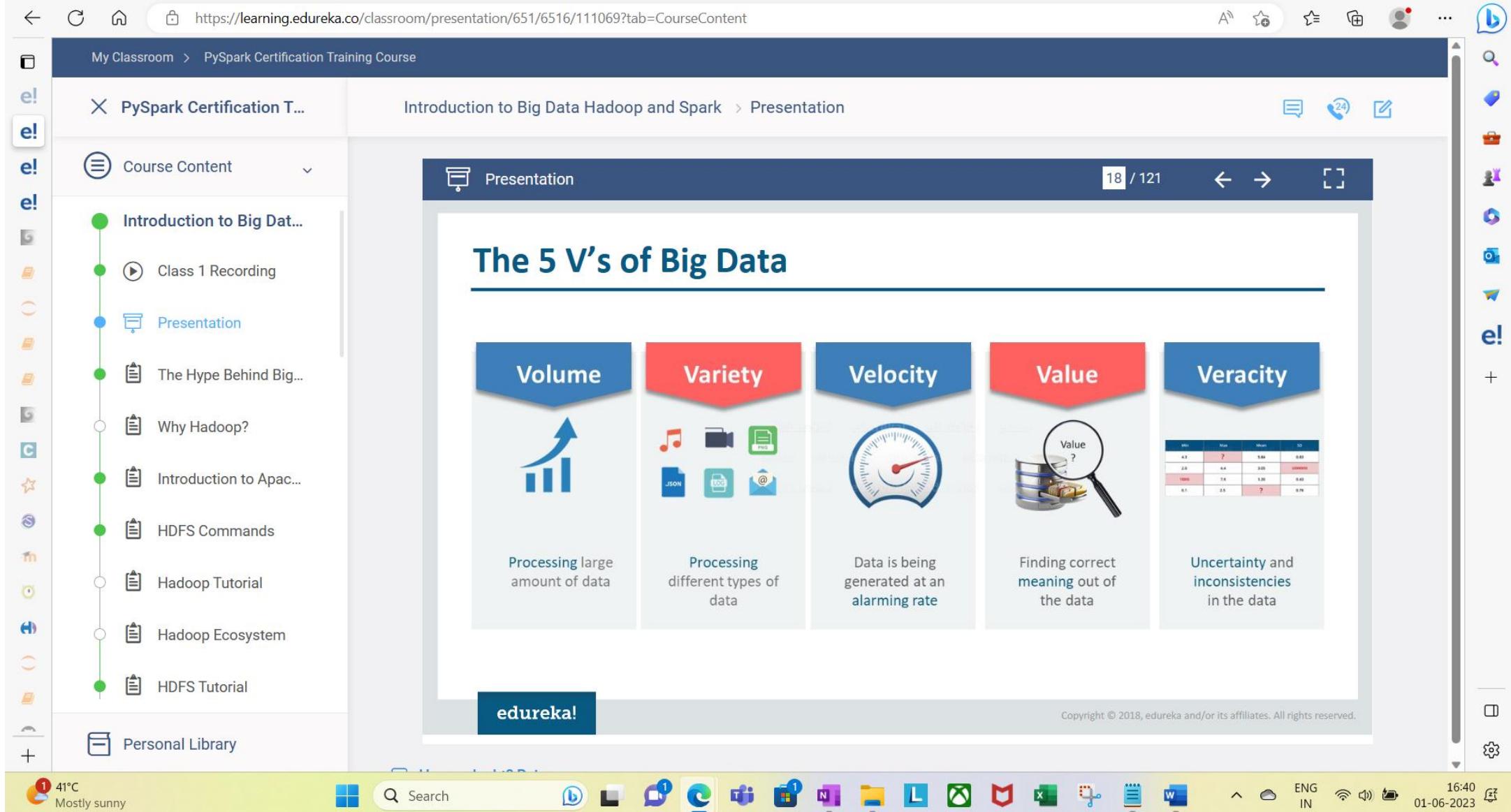
Copyright © 2018, edureka and/or its affiliates. All rights reserved.

41°C Mostly sunny

Search

16:40 01-06-2023





https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

19 / 121

Let us look at Common Big Data Customers/
Domain Scenarios

edureka!

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

41°C Mostly sunny

Search

16:40 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 20 / 121 ← →

Common Big Data Customer/Domain Scenarios

- Web and e-tailing**
 - Recommendation Engines
 - Ad Targeting
 - Search Quality
 - Abuse and Click Fraud Detection
- Government**
 - Fraud Detection and Cyber Security
 - Welfare Schemes
 - Justice
- Telecommunications**
 - Customer Churn Prevention
 - Network Performance Optimization
 - Calling Data Record (CDR) Analysis
 - Analysing Network to Predict Failure
- Healthcare and Life Sciences**
 - Health Information Exchange
 - Gene Sequencing
 - Serialization
 - Healthcare Service Quality Improvements
 - Drug Safety

ebay
Alibaba Group

at&t
Jio

Fortis
Apollo Pharmacy

edureka!

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

41°C Mostly sunny

Search

16:40 01-06-2023

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Dat...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

21 / 121



Common Big Data Customer/Domain Scenarios

- Banks and Financial services

- Modeling True Risk
- Threat Analysis
- Fraud Detection
- Trade Surveillance
- Credit Scoring and Analysis



- Retail

- Point of Sales Transaction Analysis
- Customer Churn Analysis
- Sentiment Analysis



- Transportation

- Surge Pricing
- Ride pooling
- Revenue Management
- Traffic control
- Route planning
- Logistics



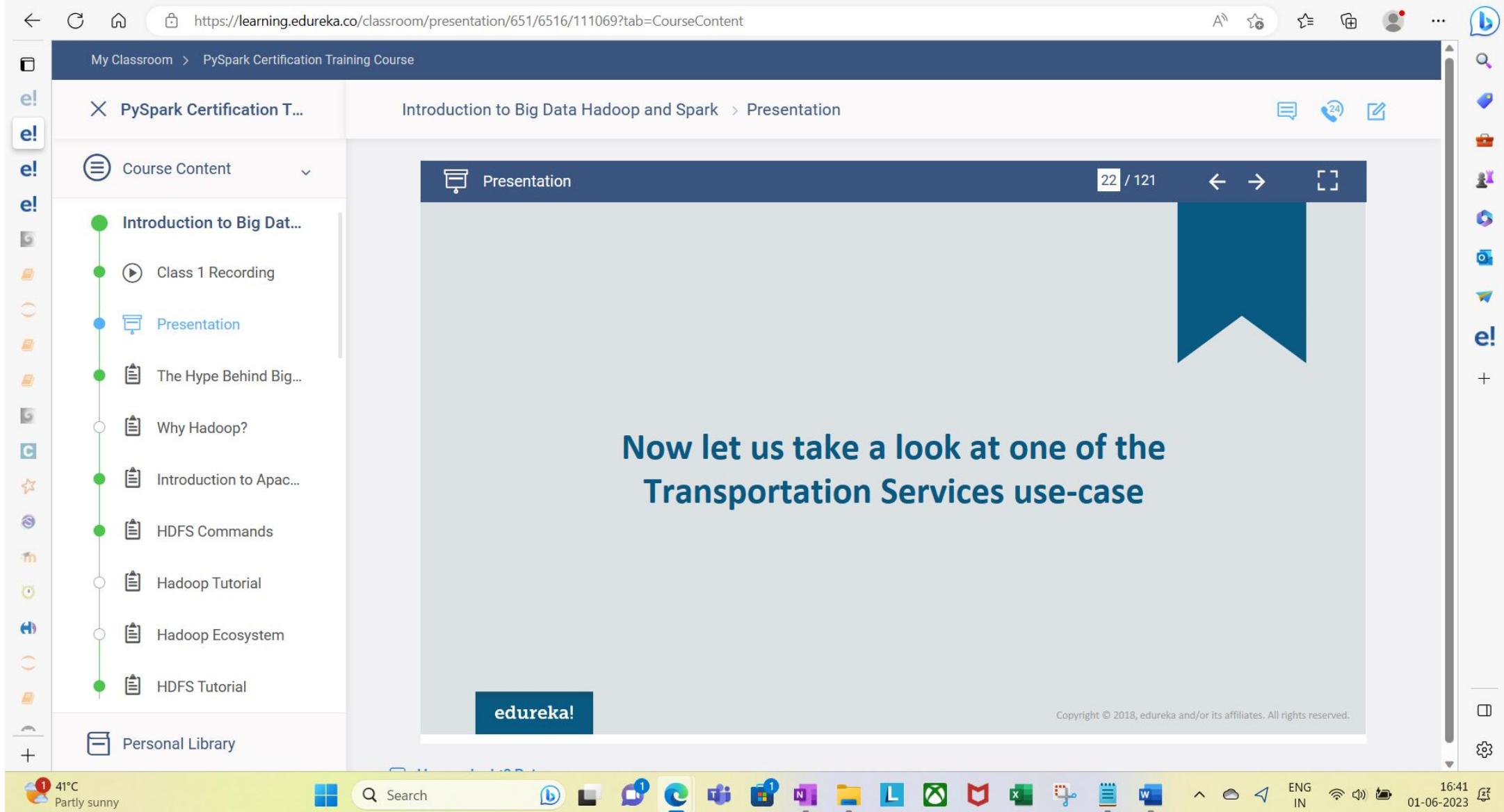
- Hotels and Food Delivery Services

- Customer Demands
- Details of Customers
- Availability and Seasonal Data Changes



edureka!

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



X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation 23 / 121 ← → []

Problems : Old Uber Infrastructure (2014)

PROBLEMS

The diagram illustrates the Old Uber Infrastructure architecture and the four main problems it faced:

- Data Sources:** Kafka Logs, Key - Val DB, RDBMS DB'S.
- Storage:** S3, EMR, Vertical Data Warehouse.
- Processing Applications:** ETL (Business Ops, A/B Experiments), Adhoc Analytics (City Ops, Data Science).

Problems identified:

- Data locality:** Using a single place for storage led to **data movement problems** and also increased the risk of **Single point of Failure**.
- Scale issues:** As they started growing there were many problems, **mainly with the batch upload problem** as they started collecting many sorts of data like "Car Speed, traffic data, user data, fare".
- Storage:** With the increasing amount of data it became difficult to store the data.
- Handling unstructured data:** With inconsistency in data (unstructured), the processing problems also increased.

edureka!

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

← ⌛ 🏠 🔍 https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent ⌛ 🏠 🔍 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

e! e! e! e!

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apache...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

24 / 121 ← →

Probable Solution to the Problems

- We need something which can *store large amounts of data*, and *can be scaled easily* if we get more data
- There shouldn't be any difficulties in *batch-upload of data*
- Data processing and analysis should be fast, *both for static as well as live streaming data*

 Fast Processing

- Handling and processing unstructured data from different sources into a *structured format*

edureka!

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

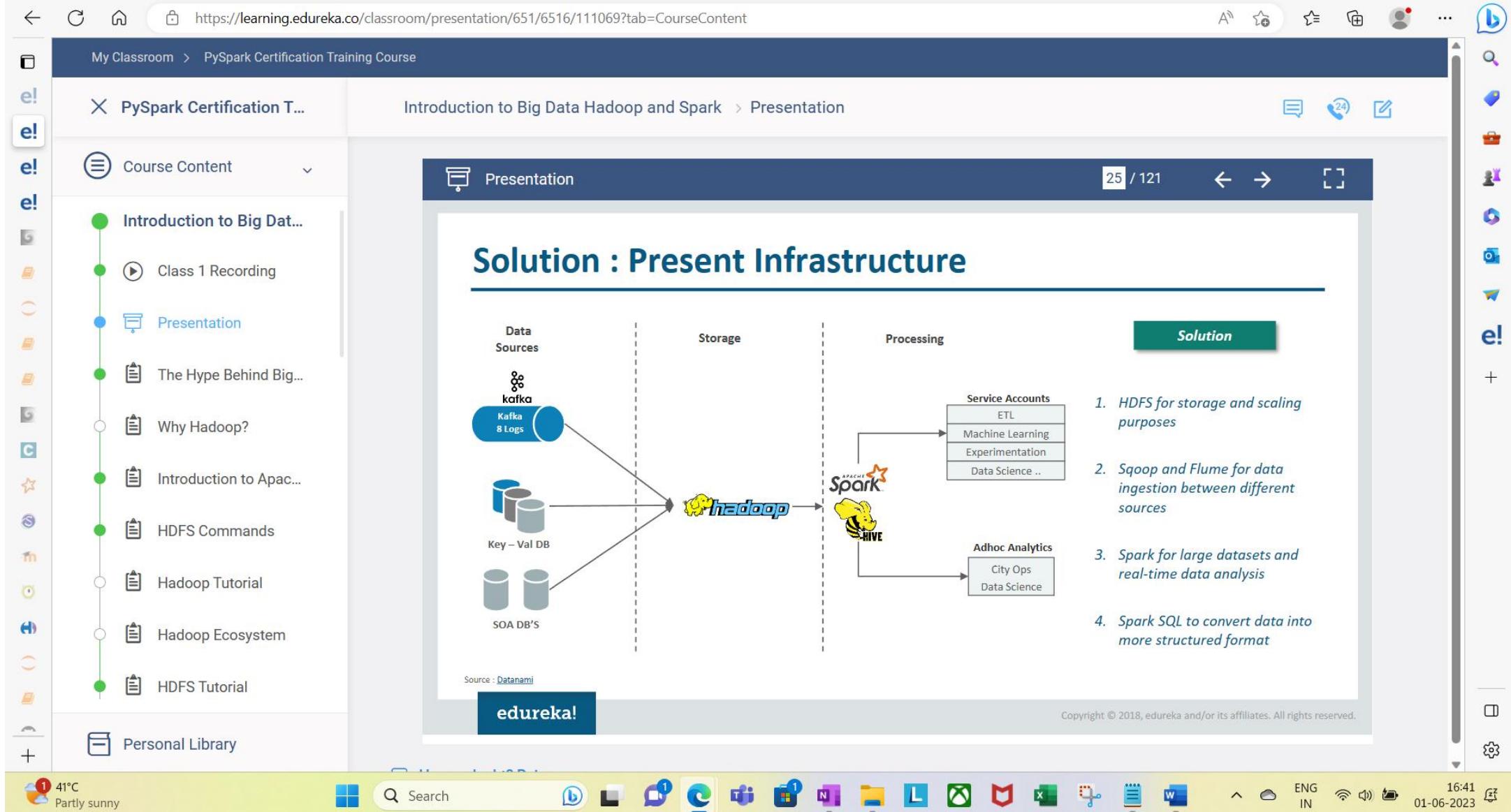
41°C Partly sunny

Search

Windows Start Menu

Cloud

16:41 ENG IN 01-06-2023





My Classroom > PySpark Certification Training Course

X PySpark Certification T...

- Introduction to Big Data
 - Class 1 Recording
 - Presentation
 - The Hype Behind Big Data
 - Why Hadoop?
 - Introduction to Apache Hadoop
 - HDFS Commands
 - Hadoop Tutorial
 - Hadoop Ecosystem
 - HDFS Tutorial

Introduction to Big Data Hadoop and Spark > Presentation



26 / 121



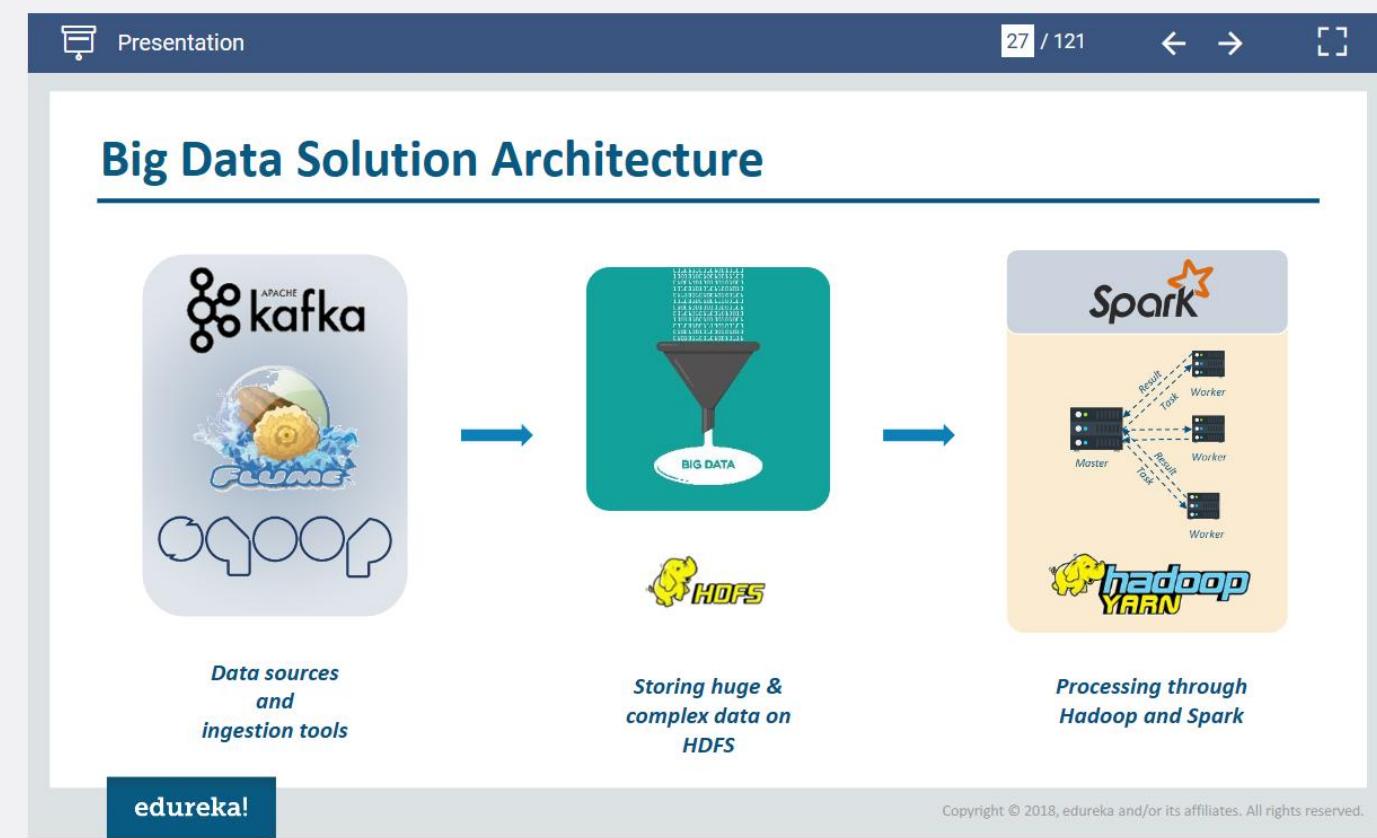
Let's see a generalized Big Data Solution Architecture

edureka

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



- e! Course Content
- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial
- Personal Library



← ⌛ 🏠 🔍 https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent ⌛ 🏠 🔍 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

e! e!

Course Content

- Introduction to Big Dat...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

28 / 121 ← → []

Let's start our journey with Hadoop

edureka!

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

Have a doubt? Raise a query.

41°C Partly sunny

Search

Windows Start

Microsoft Edge

OneDrive

OneNote

PowerPoint

Excel

Word

PowerPoint

OneDrive

OneNote

PowerPoint

Excel

Word

Cloud

ENG IN

16:41 01-06-2023

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



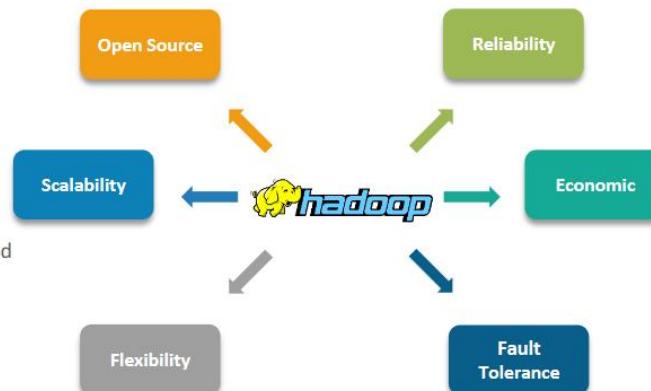
Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

What is Hadoop?

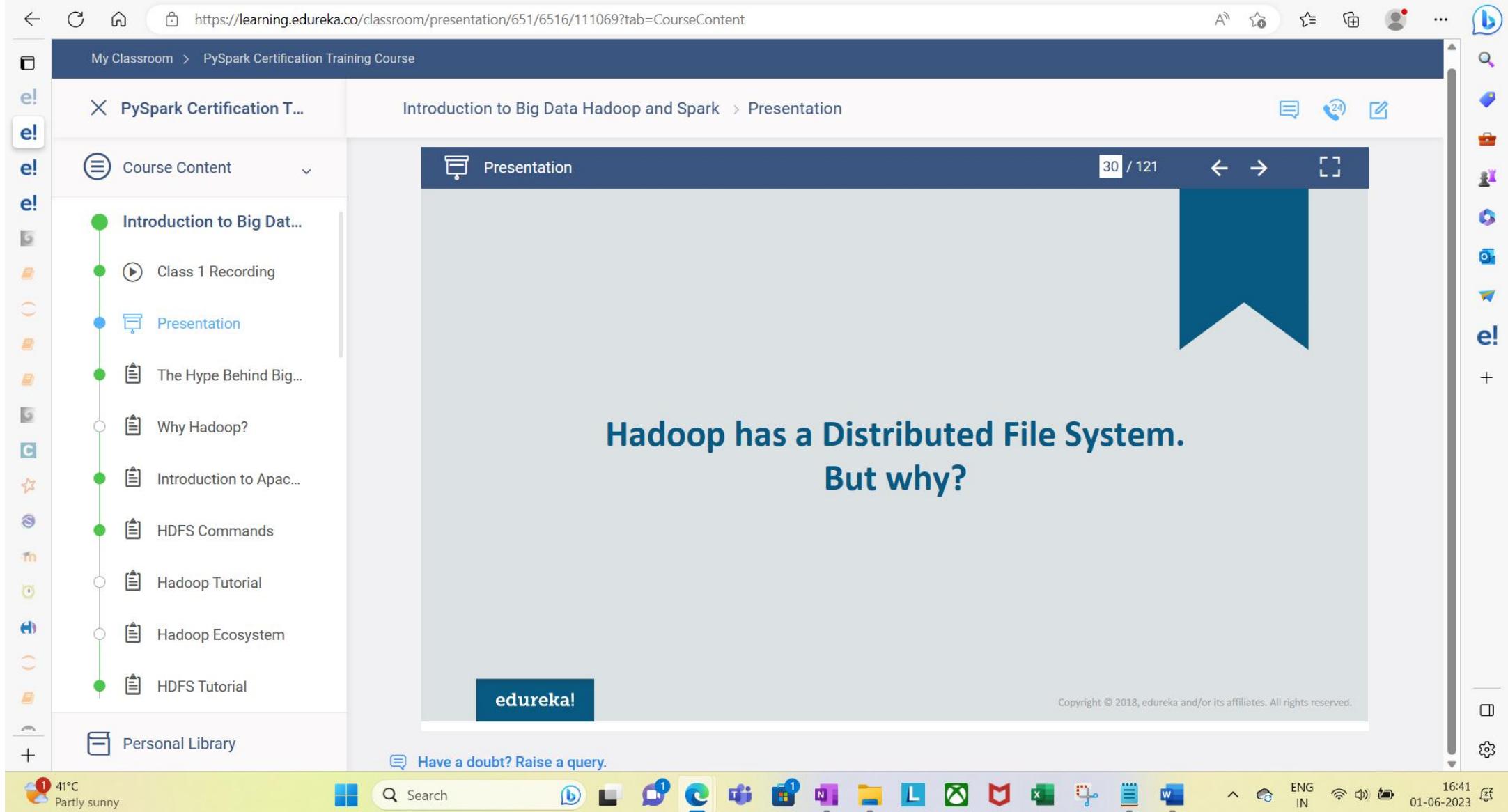
- Apache Hadoop is an **open-source framework** that allows distributed processing of large data sets across clusters of commodity machines using a simple programming model
- It enables **data management** with **scalable storage** and **distributed processing**



edureka!

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

Have a doubt? Raise a query.



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

31 / 121



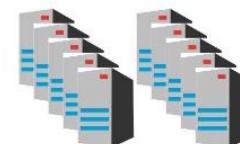
DFS : Distributed File System

Distributed & Parallel Computation

Read 1 TB data



1 Machine

4 I/O channels
Each channel - 100 MB/s

10 Machine

4 I/O channels
Each channel - 100 MB/s

43 Minutes

Time taken = ?

edureka!

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

Have a doubt? Raise a query.

41°C
Partly sunny

Search

ENG
IN

01-06-2023

16:41

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 32 / 121

DFS : Distributed File System

Distributed & Parallel Computation

Read 1 TB data

1 Machine
4 I/O channels
Each channel – 100 MB/s

4.3 Minutes

10 Machine
4 I/O channels
Each channel – 100 MB/s

4.3 Minutes

edureka!

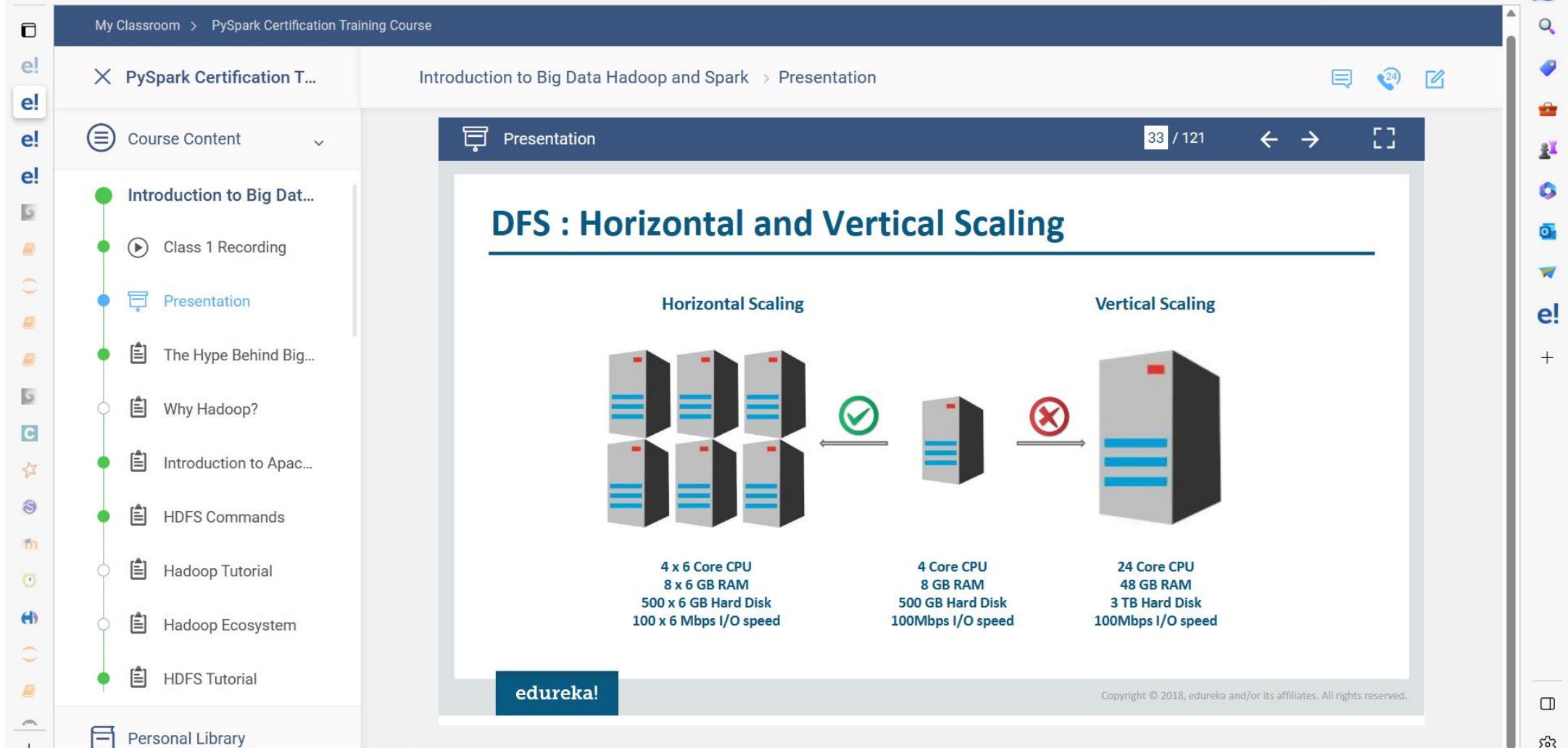
Copyright © 2018, edureka and/or its affiliates. All rights reserved.

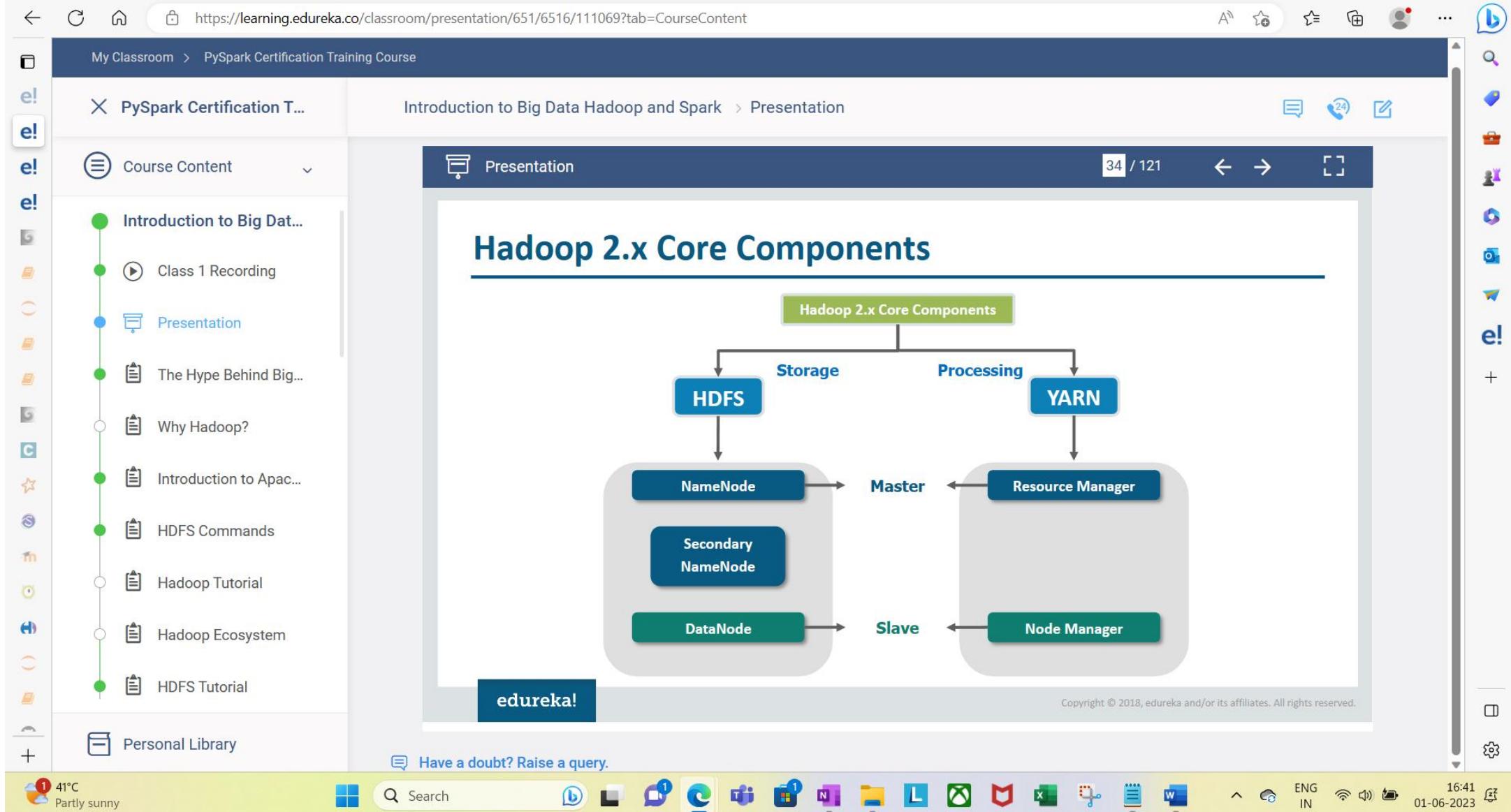
Have a doubt? Raise a query.

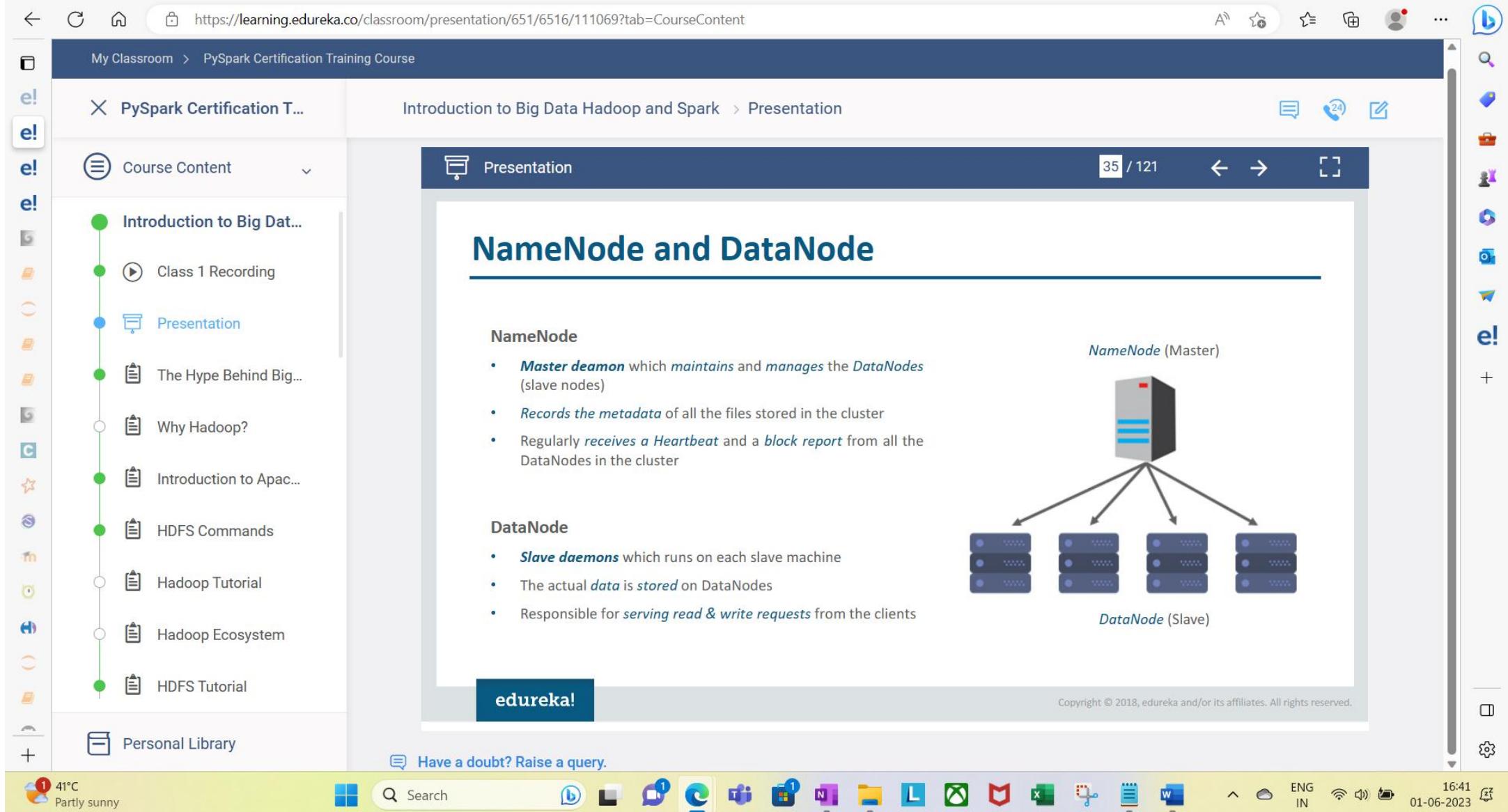
41°C Partly sunny

Search

16:41 IN 01-06-2023







https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 36 / 121

Secondary NameNode & Checkpointing

Secondary NameNode

- Job of **Secondary NameNode** is to contact **NameNode** in a periodic manner after certain time interval(by default 1 hour) and *pulls copy of metadata* information out of NameNode
- Checkpointing** is a process of *combining edit logs with FslImage*
- Secondary NameNode takes over the responsibility of checkpointing, therefore making NameNode more available

edureka!

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

Have a doubt? Raise a query.

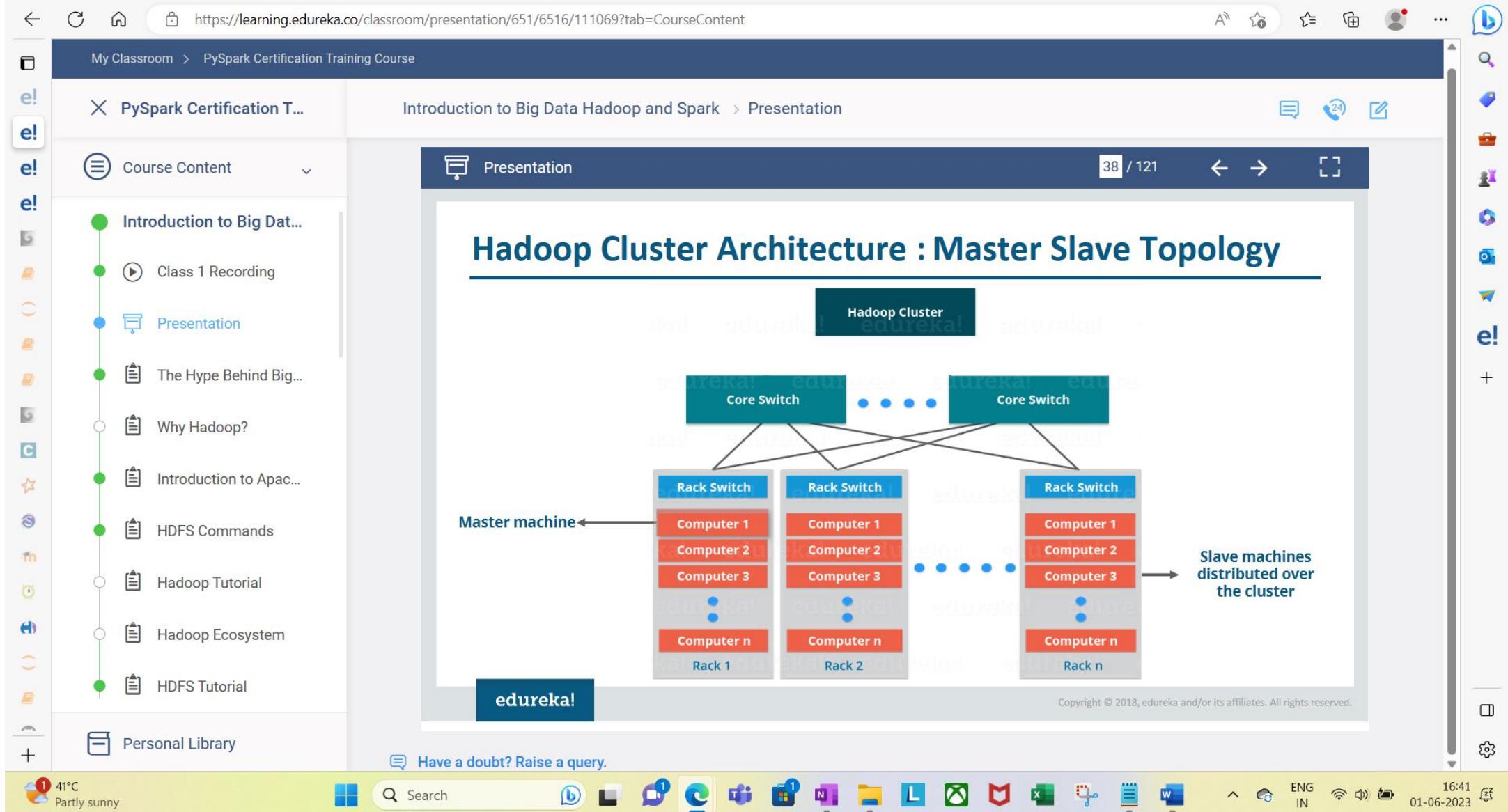
41°C Partly sunny

Search

L

ENG IN

16:41 01-06-2023



https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 39 / 121 ← → []

Typical Hadoop Cluster Configuration

NameNode

- RAM: 64 GB
- Hard disk: 1 TB
- Processor: Xeon with 8 Cores
- Ethernet: 3 x 10 GB/s
- OS: 64-bit CentOS / Linux
- Power: Redundant Power Supply

Secondary NameNode

- RAM: 32 GB
- Hard disk: 1 TB
- Processor: Xeon with 4 Cores
- Ethernet: 3 x 10 GB/s
- OS: 64-bit CentOS / Linux
- Power: Redundant Power Supply

DataNode

- RAM: 16GB
- Hard disk: 6 x 2TB
- Processor: Xeon with 2 cores
- Ethernet: 3 x 10 GB/s
- OS: 64-bit CentOS / Linux

edureka!

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

Have a doubt? Raise a query.

41°C Partly sunny

Search

Cloud

ENG IN

16:41 01-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

40 / 121



Let us talk about, how data is
stored on HDFS?

edureka!

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

Have a doubt? Raise a query.

1 41°C
Partly sunny



Search



ENG
IN



01-06-2023

16:41

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 41 / 121

HDFS Blocks

- Each file is stored on HDFS as blocks
- The default size of each block is 128 MB in Apache Hadoop 2.x
- Let us say, I have a file "example.txt" of size 248 MB. Below is the representation of how it will be stored on HDFS

How many blocks will be created if a file of size 514 MB is copied to HDFS ?

edureka!

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

Have a doubt? Raise a query.

41°C Partly sunny

Search

edureka!

16:41 01-06-2023

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 42 / 121

HDFS Blocks

- Each file is stored on HDFS as blocks
- The default size of each block is 128 MB in Apache Hadoop 2.x
- Let us say, I have a file "example.txt" of size 248 MB. Below is the representation of how it will be stored on HDFS

How many blocks will be created if a file of size 514 MB is copied to HDFS?

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

edureka!

Have a doubt? Raise a query.

41°C Partly sunny

Search

edureka!

16:41 01-06-2023

← ⌛ 🏠 🔍 https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent ⌛ 🏠 🔍 ... 

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Dat...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

43 / 121 ← → []

Is it safe to have just 1 copy of each block?
What do you think?

edureka!

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

Have a doubt? Raise a query.

41°C Partly sunny

Search                    16:41 01-06-2023

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 44 / 121 ← → []

Block Replication

By Default the replication factor is 3

The diagram illustrates the concept of block replication in a distributed system. It shows four nodes (Node 1, Node 2, Node 3, Node 4) each containing three copies of five different blocks. The blocks are represented by colored rectangles: Block 1 (orange), Block 2 (yellow), Block 3 (green), Block 4 (red), and Block 5 (blue). Each node has a small server icon above it. A callout box at the top right states "By Default the replication factor is 3".

edureka!

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

Have a doubt? Raise a query.

41°C Partly sunny

Search

Edureka!

16:42 01-06-2023

← ⌛ 🏠 🔍 https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent ⌛ ⌚ 🗑️ 📁 🌐 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

e! e! e! e!

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

49 / 121 ← → []

Rack Awareness

Rack	DataNodes
Rack - 1	DN1, DN5, DN9
Rack - 2	DN6, DN10
Rack - 3	DN2, DN3, DN4, DN7, DN8, DN11, DN12

edureka!

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

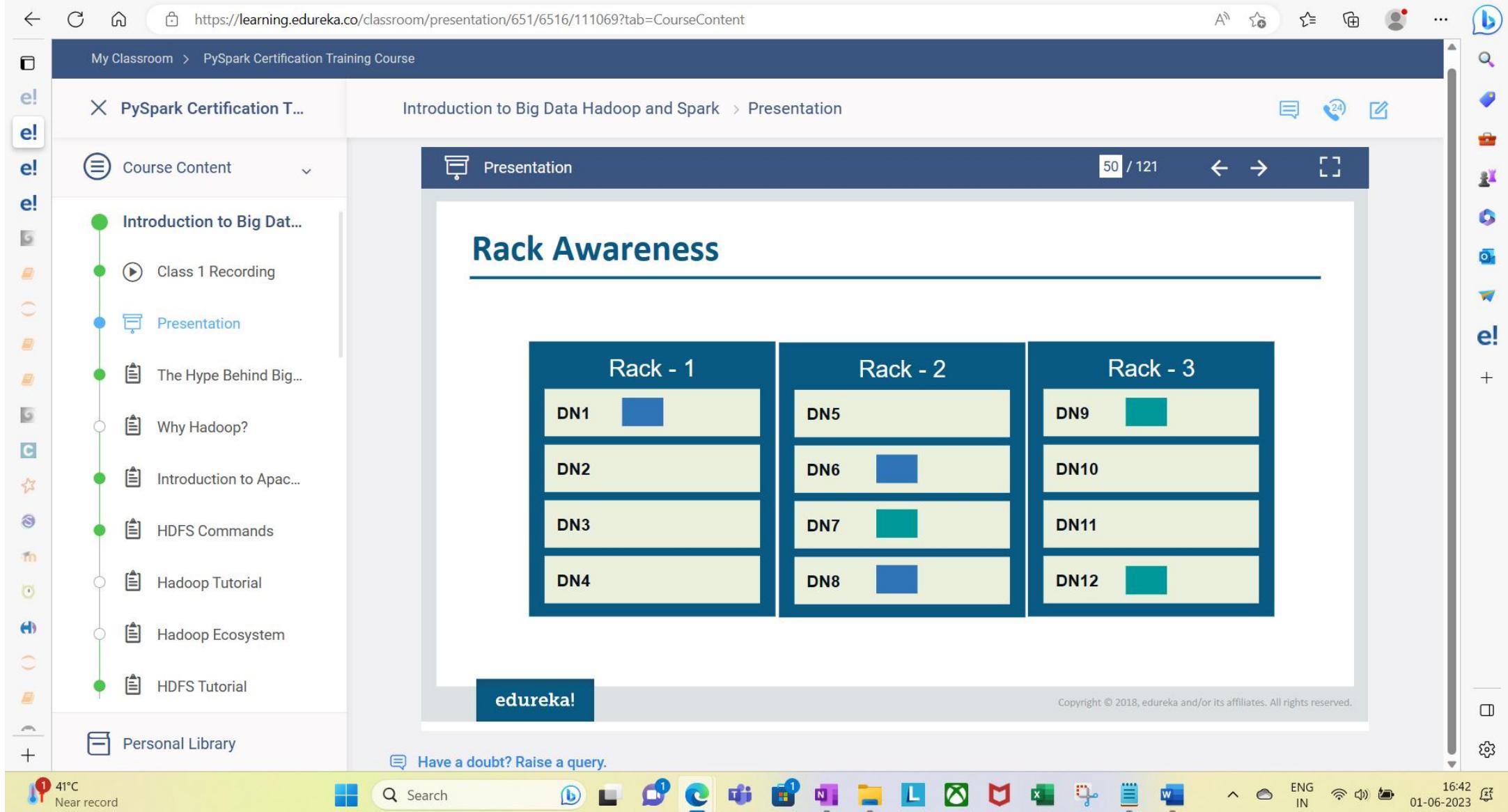
Have a doubt? Raise a query.

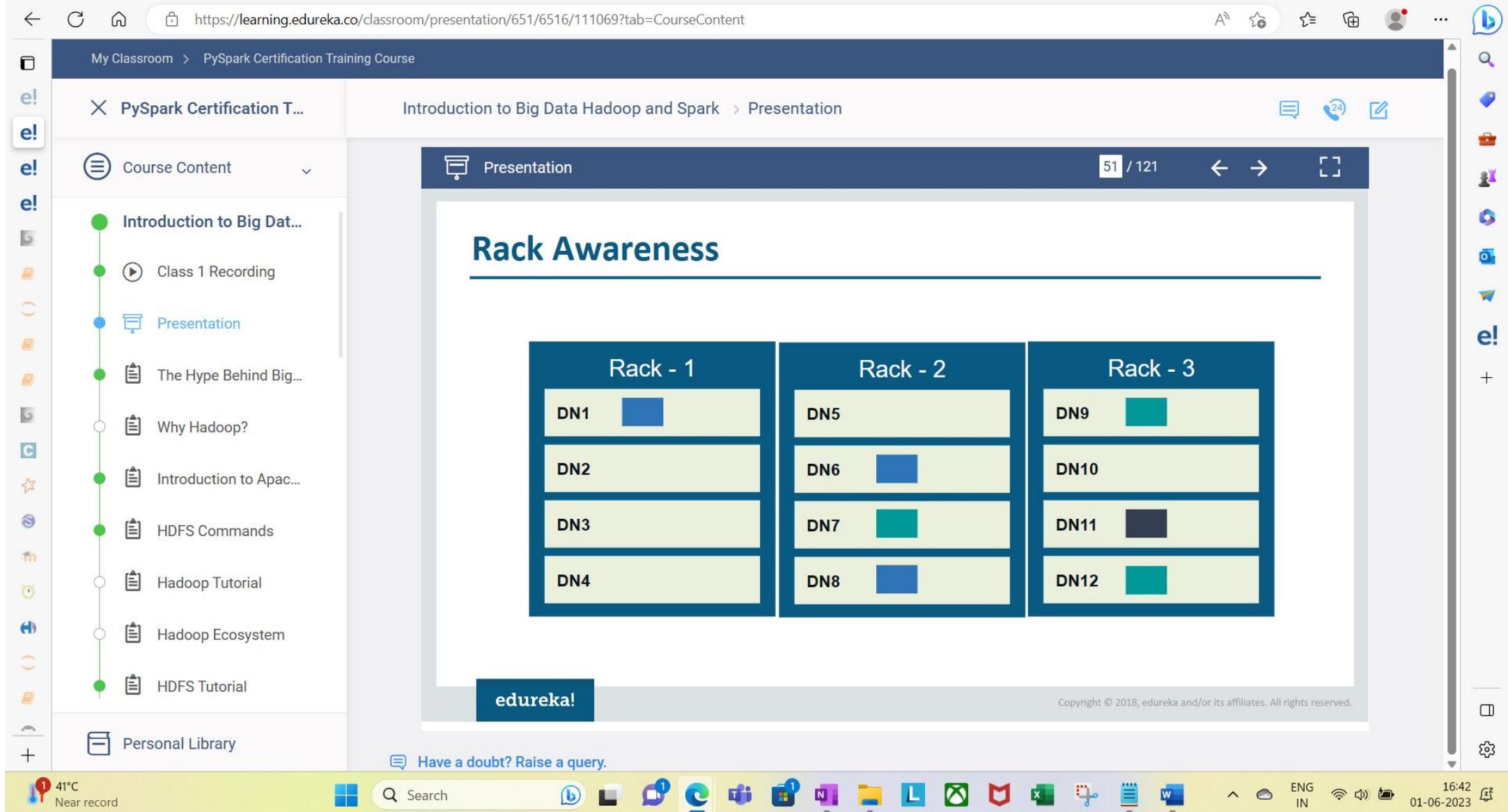
41°C Partly sunny

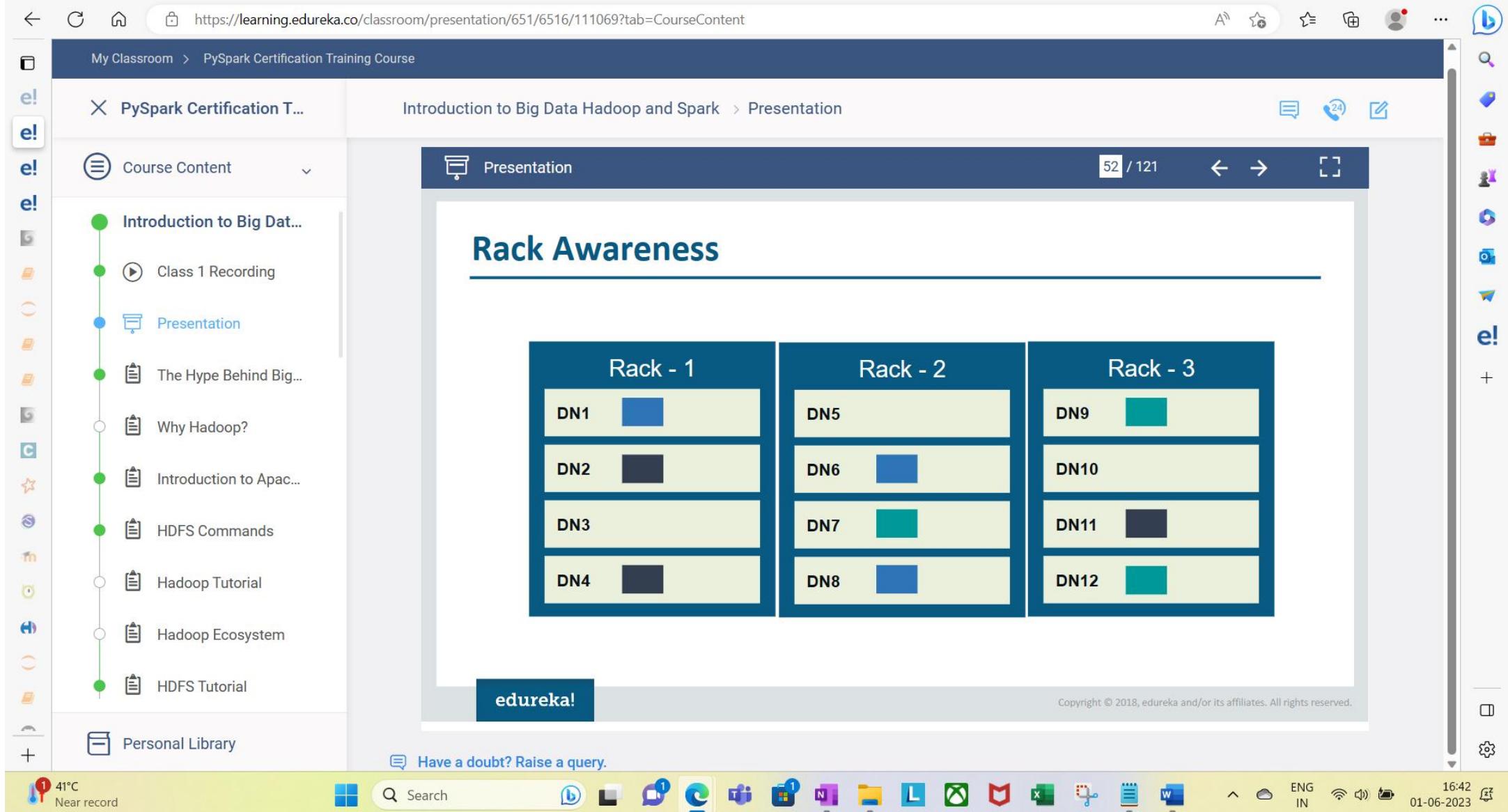
Search

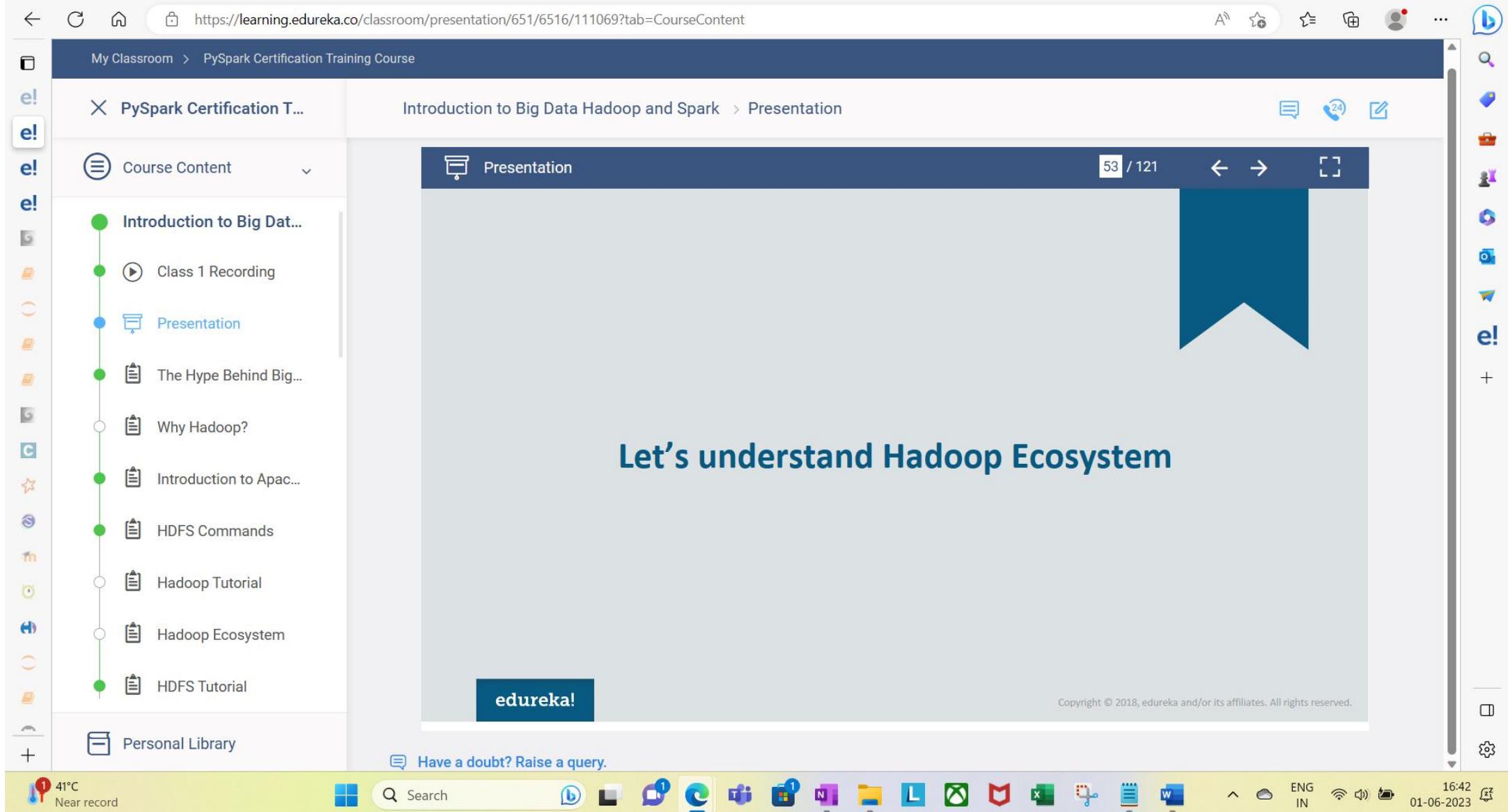
Windows Start

Icons: Edge, File Explorer, Task View, Mail, OneDrive, Microsoft Store, Notepad, File Explorer, Lync, Excel, Word, Powerpoint, Cloud, Network, Battery, Volume, Signal, ENG IN, 01-06-2023, 16:42









My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

edureka!

Hadoop Ecosystem

The diagram illustrates the Hadoop Ecosystem, showing the integration of various Apache projects:

- MapReduce**: Processing using different languages.
- Hive & Drill**: Analytical SQL-on-Hadoop.
- Mahout & Spark MLlib**: Machine Learning.
- PIG**: Scripting.
- HBase**: NoSQL Database.
- ZOOKEEPER & AMBARI**: Management & Coordination.
- SPARK (In-Memory Data Flow Engine)**
- KAFKA**: Streaming.
- Apache Spark MLlib**
- SOLR & LUCENE**: Searching & Indexing.
- OOZIE**: Scheduling.
- YARN**: Resource Management.
- HDFS**: Storage.

External interfaces:

- Flume**: Unstructured or Semi-structured Data.
- Sqoop**: Structured Data.

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

Have a doubt? Raise a query.

41°C Near record

Search

16:42 01-06-2023

Components of Hadoop Ecosystem



Apache Hadoop provides a **scalable solution** to **store and process huge data sets** in parallel and distributed fashion



Apache Spark is an **in-memory data processing engine** that allows us to efficiently execute streaming, machine learning or SQL workloads



Apache Flume is a **distributed, reliable and available service** for **efficiently collecting, aggregating, and moving large amounts of log data**



Apache Sqoop is a tool designed for efficiently **transferring bulk data** between Apache Hadoop and **structured datastores** such as RDBMS

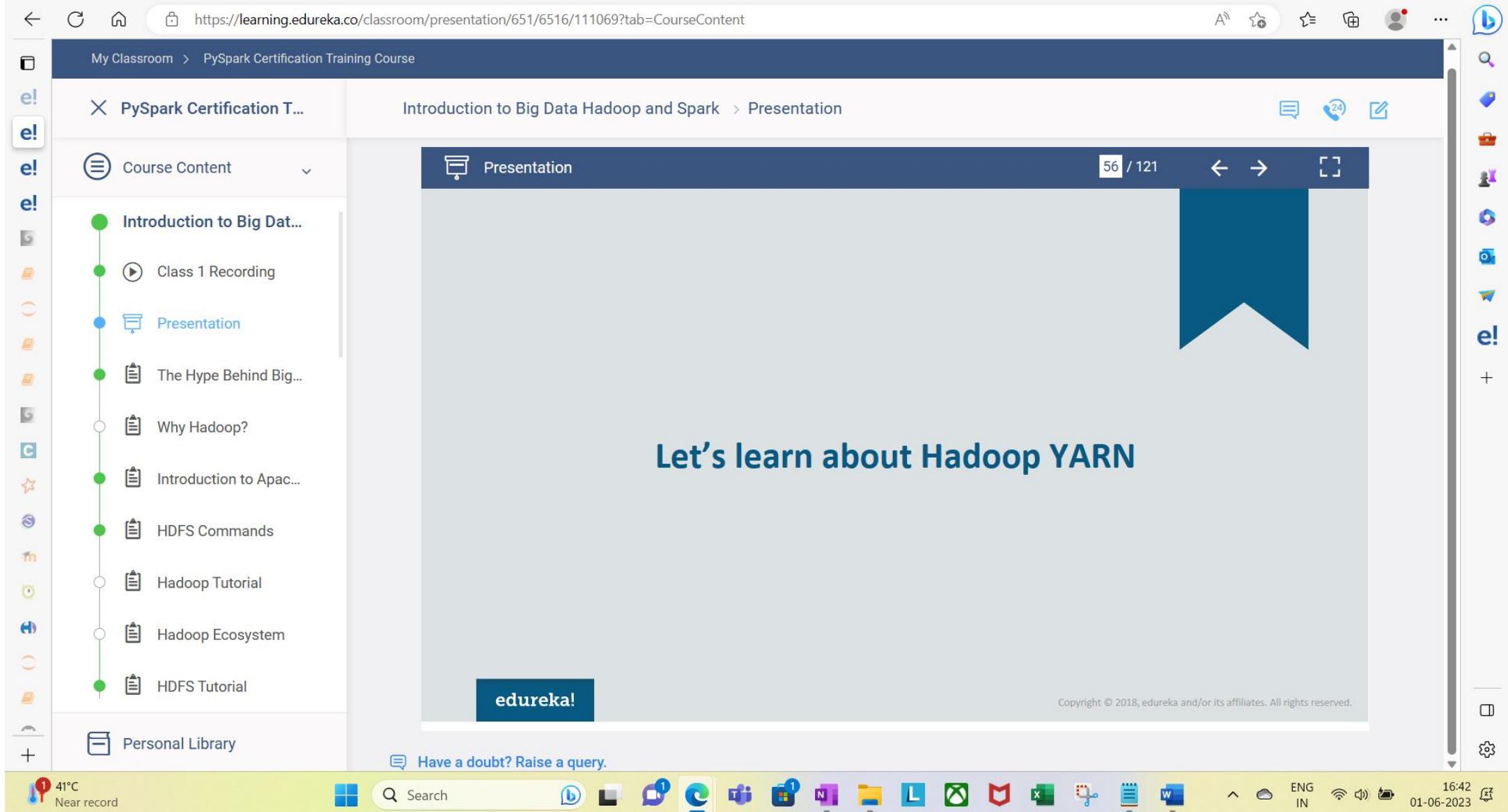


Apache Kafka is a **fast, scalable, durable and fault-tolerant publish-subscribe messaging system**

edureka!

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

... Have a doubt? Raise a query.



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

- Introduction to Big Data...
- (▶) Class 1 Recording
- (▶) Presentation
- (📄) The Hype Behind Big...
- (📄) Why Hadoop?
- (📄) Introduction to Apac...
- (📄) HDFS Commands
- (📄) Hadoop Tutorial
- (📄) Hadoop Ecosystem
- (📄) HDFS Tutorial

Personal Library

Presentation

57 / 121



YARN Use Case - YAHOO



Yahoo has close to 35,000 nodes running Apache Hadoop

Performs 80,000 MapReduce jobs per day

Over 150 PB of storage

YARN at Yahoo helped them increase the load on the *most heavily used Hadoop cluster*. Now performs *125,000 jobs a day* when compared to *80,000 jobs a day* which is close to *50% increase*

edureka!

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

Have a doubt? Raise a query.

My Classroom > PySpark Certification Training Course

Introduction to Big Data Hadoop and Spark > Presentation



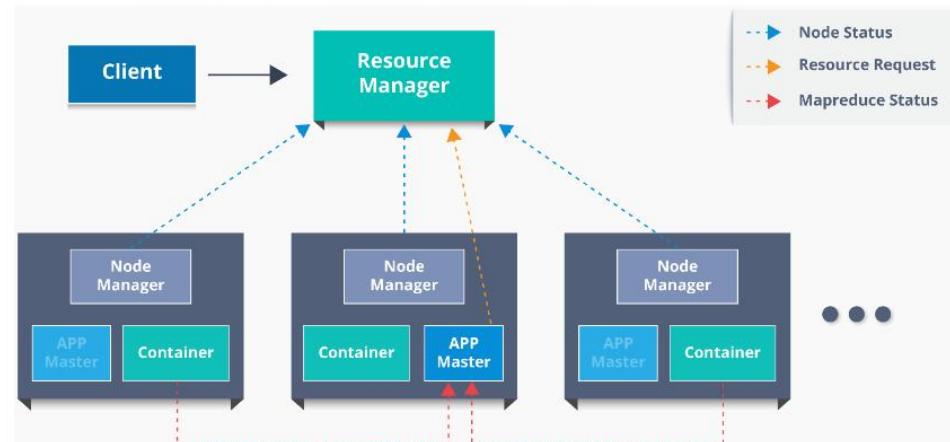
58 / 121

← -

5

What is YARN?

- Hadoop 2.0 came up with a new framework **YARN** (*Yet Another Resource Negotiator*), which provides ability to run *Non-MapReduce application*
 - YARN framework is *responsible for doing Cluster Resource Management*



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

edureka!

 Have a doubt? Raise a query

← ⌛ 🏠 🔍 https://learning.edureka.co/course/presentation/651/6516/111069?tab=CourseContent ⌛ 🏠 🔍 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

e! e!

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

59 / 121 ← → []

YARN Components

edureka!

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

Have a doubt? Raise a query.

41°C Near record

Search

edureka!

16:42 01-06-2023



Main Components of YARN

ResourceManager:

- Master daemon that manages all other daemons & accepts job submission
- Allocates first container for the AppMaster

Resource Manager

Node Manager

App Master

Container

NodeManager:

- Responsible for containers, monitoring their resource usage i.e. (CPU, memory, disk, network) & reports the same to RM

AppMaster:

- One per application
- Coordinates and manages MR jobs
- Negotiates resources from RM

Container:

- Allocates certain amount of resources (memory, CPU etc.) on a slave node (NM)

edureka!

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

Have a doubt? Raise a query.

← ⌛ 🏠 🔍 https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent ⌛ 🏠 🔍 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

61 / 121 ← → []

Application Submission in YARN

edureka!

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

Have a doubt? Raise a query.

41°C Near record

Search

Windows Start

Microsoft Edge

OneDrive

OneNote

PowerPoint

Excel

Word

PowerPoint

OneDrive

OneNote

PowerPoint

Excel

Word

File Explorer

Libraries

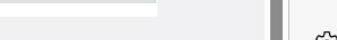
Task View

Cloud

ENG IN

16:42

01-06-2023



X PySpark Certification T...

Course Content

Introduction to Big Dat...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

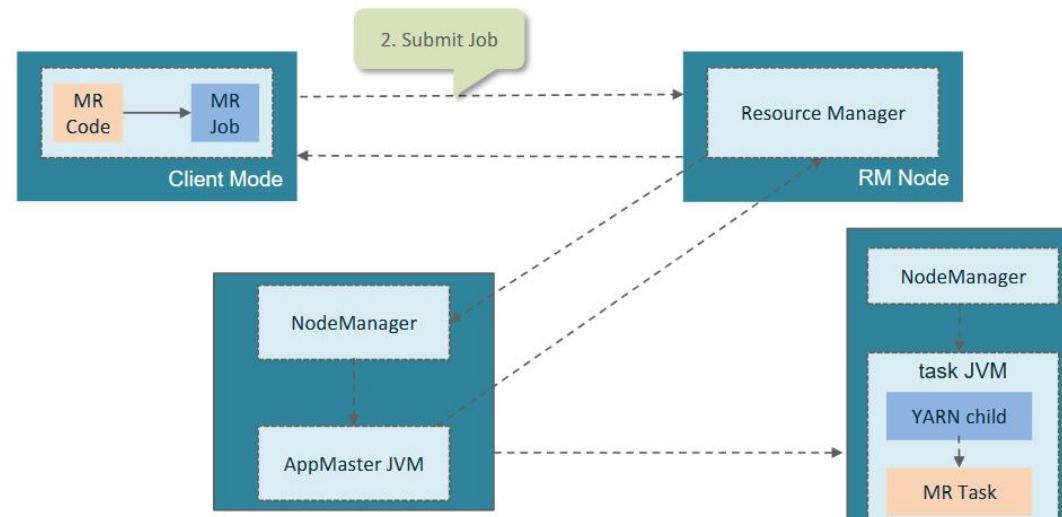
Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

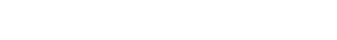
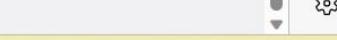
Application Submission in YARN



edureka!

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

Have a doubt? Raise a query.



X PySpark Certification T...

Course Content

Introduction to Big Dat...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

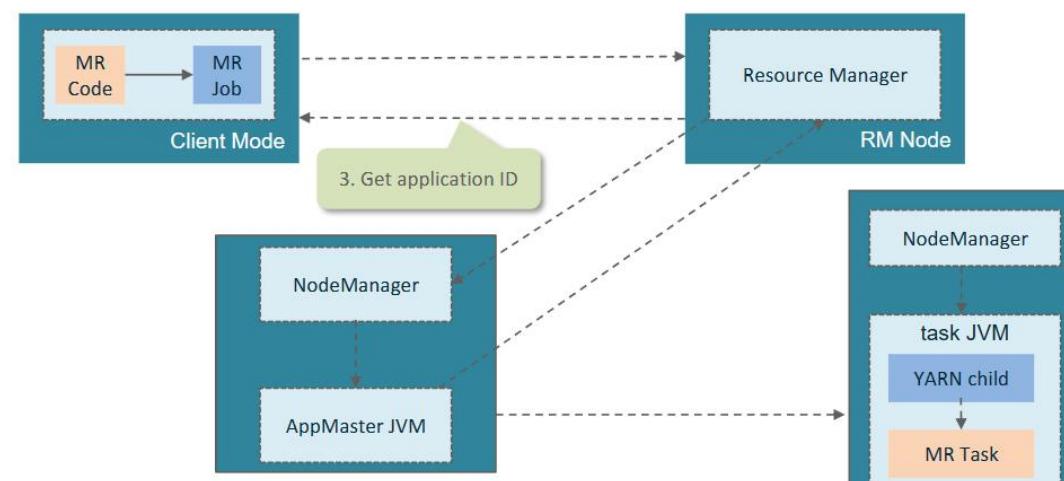
Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

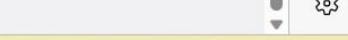
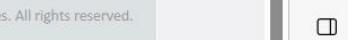
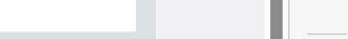
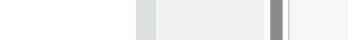
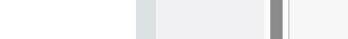
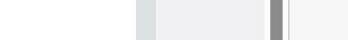
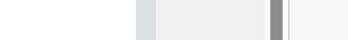
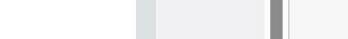
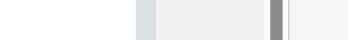
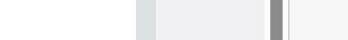
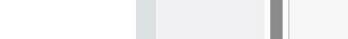
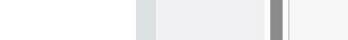
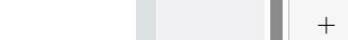
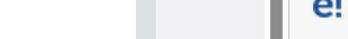
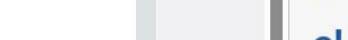
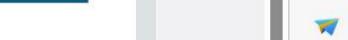
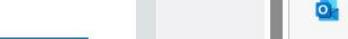
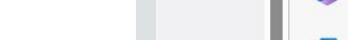
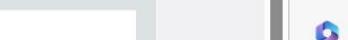
Application Submission in YARN



edureka!

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

Have a doubt? Raise a query.



X PySpark Certification T...

Course Content

Introduction to Big Dat...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

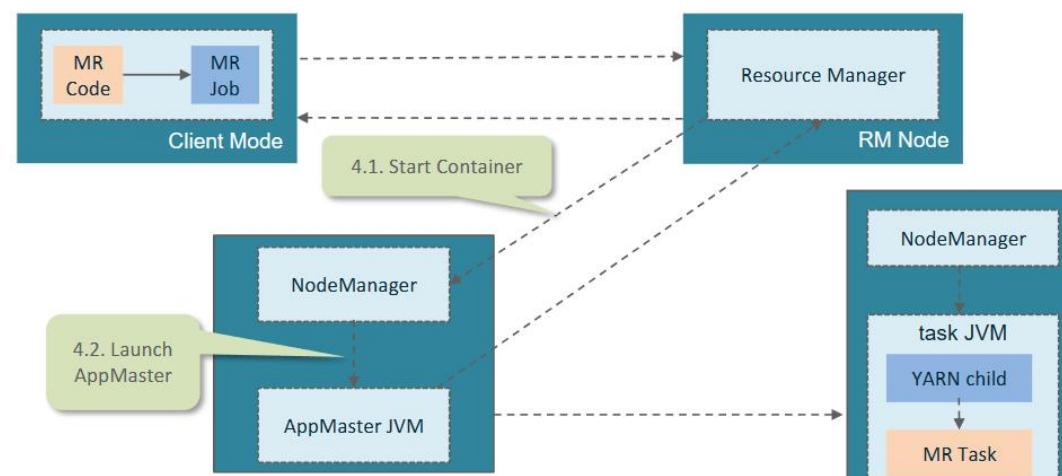
Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

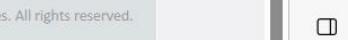
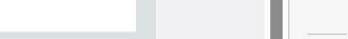
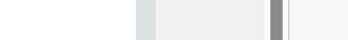
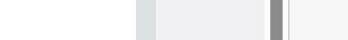
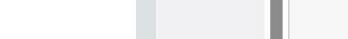
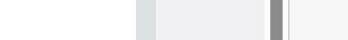
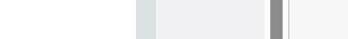
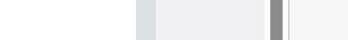
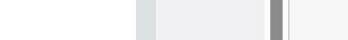
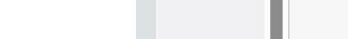
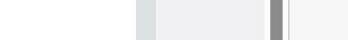
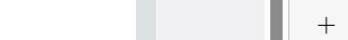
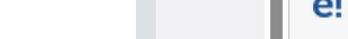
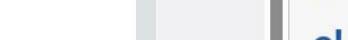
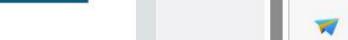
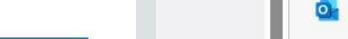
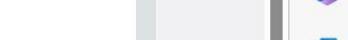
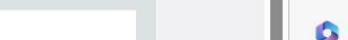
Application Submission in YARN



edureka!

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

Have a doubt? Raise a query.



X PySpark Certification T...

Course Content

Introduction to Big Dat...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

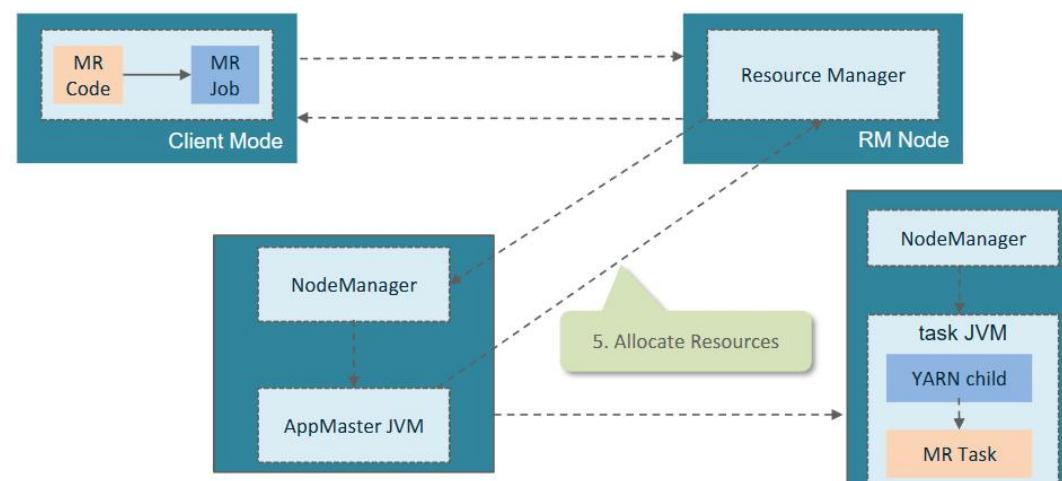
Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

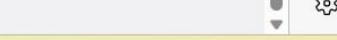
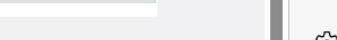
Application Submission in YARN



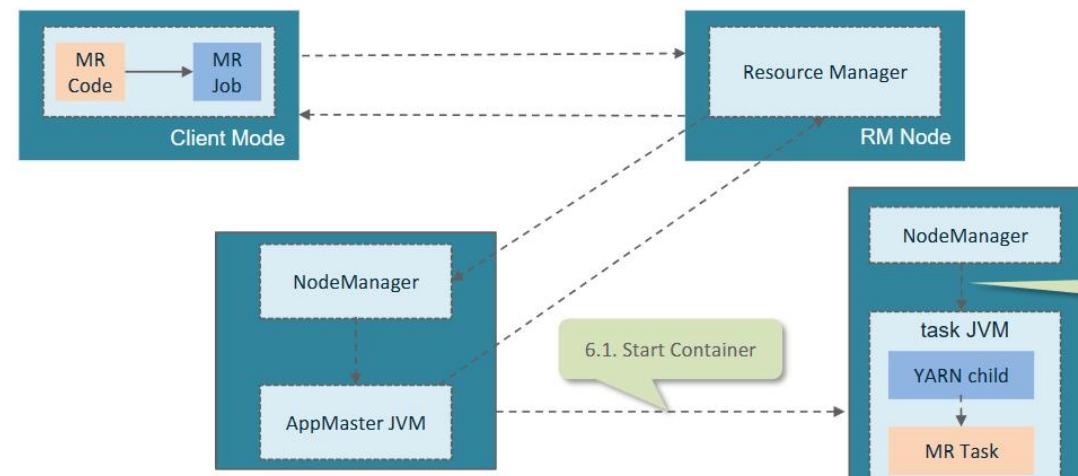
edureka!

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

Have a doubt? Raise a query.



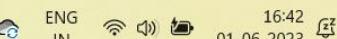
Application Submission in YARN



edureka!

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

Have a doubt? Raise a query.



X PySpark Certification T...

Course Content

Introduction to Big Dat...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

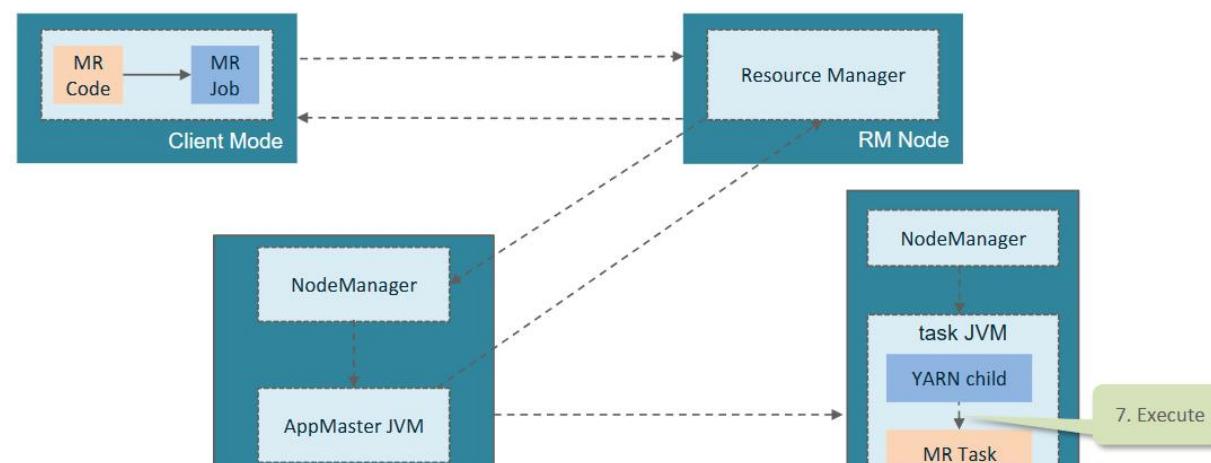
Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Application Submission in YARN



edureka!

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

Have a doubt? Raise a query.

← ⌛ 🏠 🔍 https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent ⌛ ⌚ 🗑️ 🗑️ ... 📡

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

e! e! e! e!

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

edureka! 69 / 121 ← → []

YARN Application Workflow

➤ Execution Sequence :

- Client submits an application

```
graph LR; Client -- "1" --> RM; RM --> NM; NM --> AM;
```

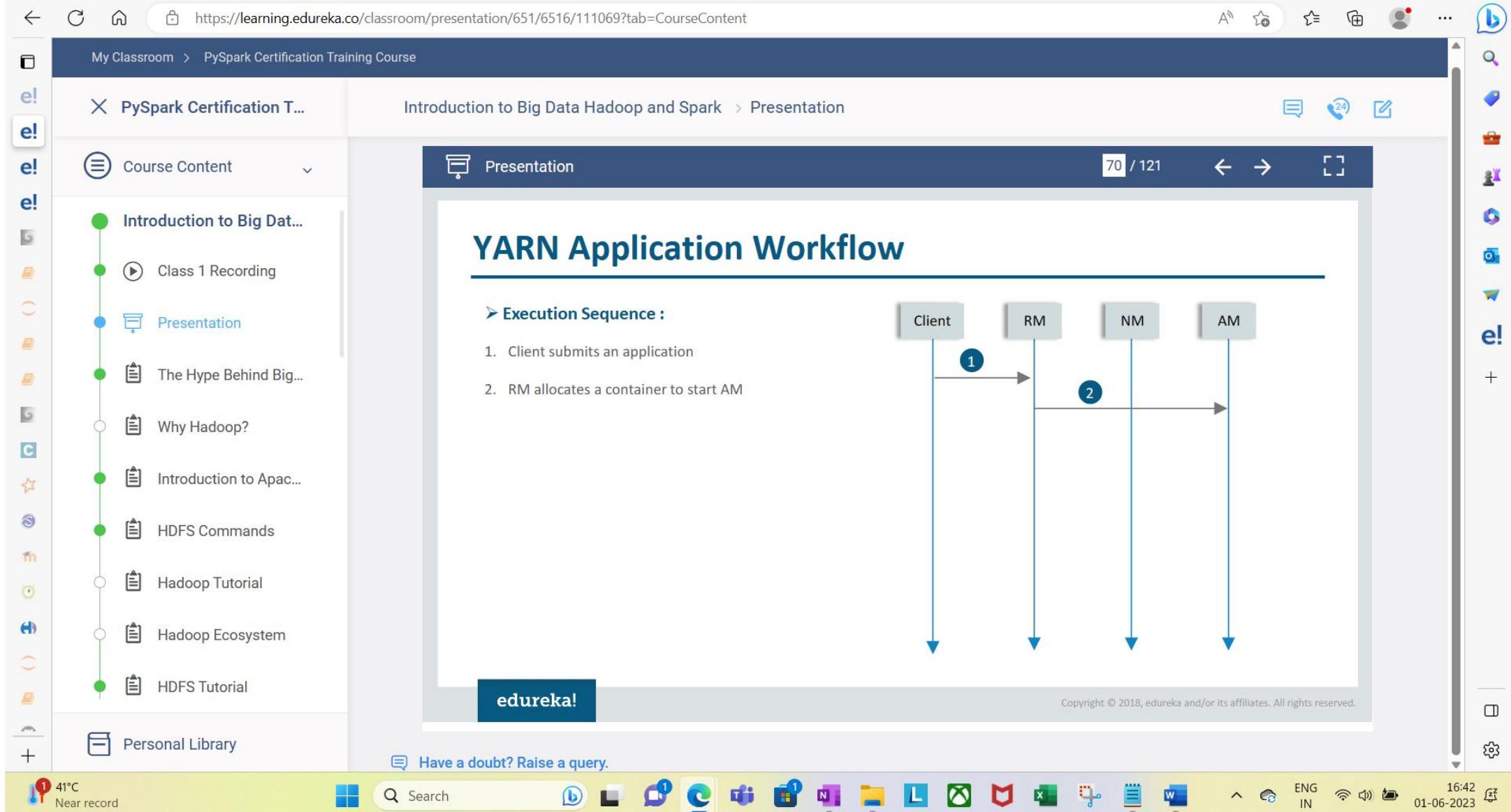
Copyright © 2018, edureka and/or its affiliates. All rights reserved.

Have a doubt? Raise a query.

41°C Near record

Search

edureka! 16:42 ENG IN 01-06-2023





YARN Application Workflow

➤ Execution Sequence :

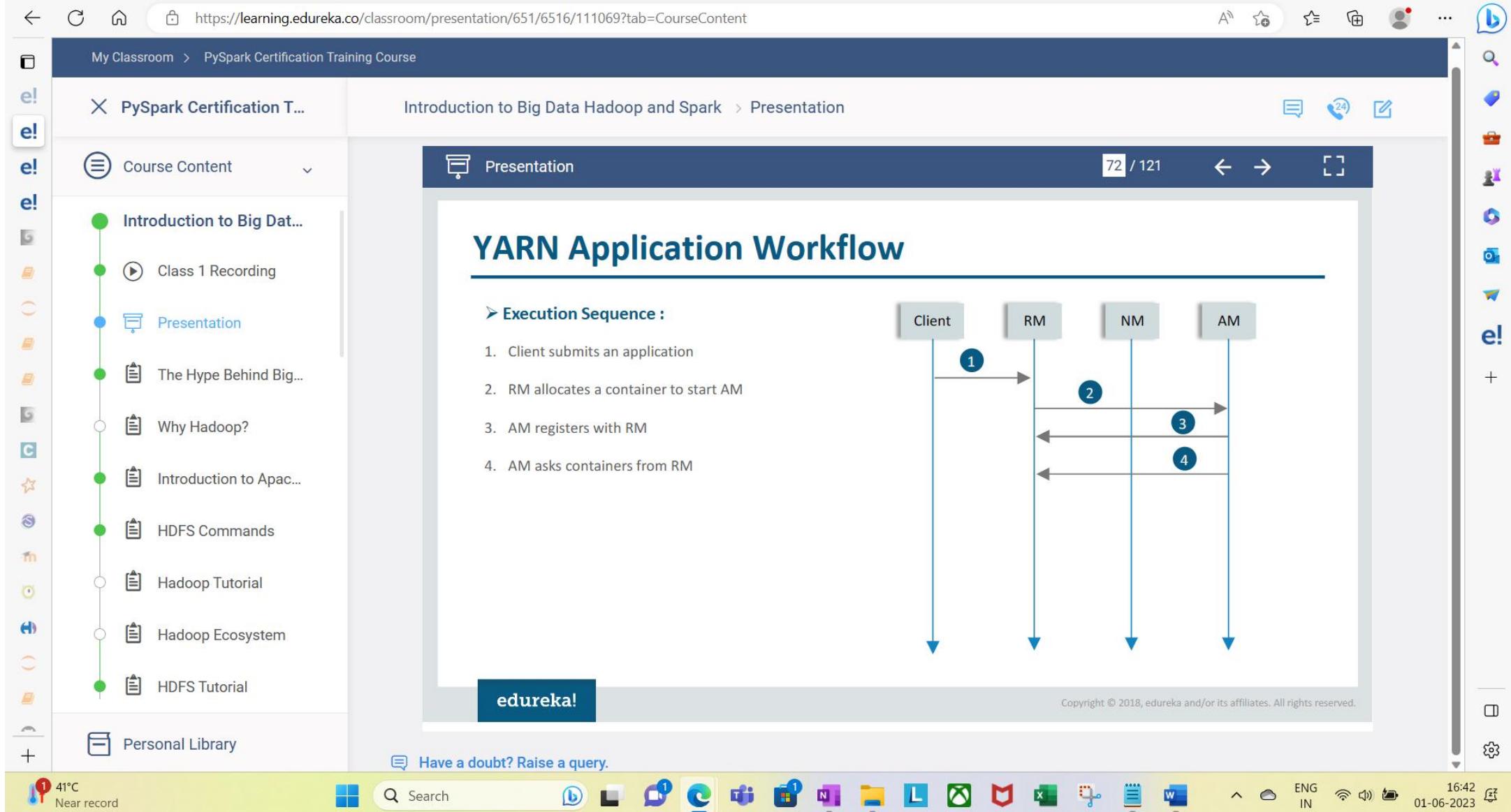
1. Client submits an application
2. RM allocates a container to start AM
3. AM registers with RM

```
sequenceDiagram
    participant Client
    participant RM
    participant NM
    participant AM
    Client->>RM: Submit application
    RM->>AM: Allocate container
    AM->>RM: Register with RM
```

edureka!

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

Have a doubt? Raise a query.



← ⌛ 🏠 🔍 https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent ⌛ ⌚ 🗑️ 🗑️ ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

e! e! e! e!

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apache...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

edureka!

Have a doubt? Raise a query.

73 / 121 ← → []

YARN Application Workflow

Execution Sequence :

- Client submits an application
- RM allocates a container to start AM
- AM registers with RM
- AM asks containers from RM
- AM notifies NM to launch containers

```
sequenceDiagram
    participant Client
    participant RM
    participant NM
    participant AM
    Client->>RM: Submit application
    RM->>AM: Allocate container
    AM->>RM: Register with RM
    AM->>NM: Ask containers
    NM->>AM: Launch containers
    AM->>NM: Notify NM
```

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

41°C Near record

Search

edureka!

Have a doubt? Raise a query.

ENG IN

16:42 01-06-2023

X PySpark Certification T...

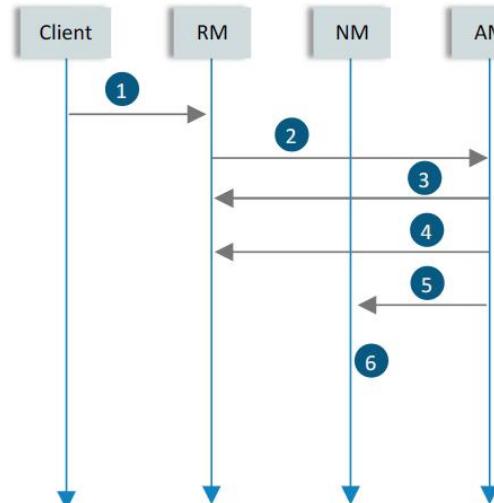
Introduction to Big Data Hadoop and Spark > Presentation



YARN Application Workflow

➤ Execution Sequence :

1. Client submits an application
2. RM allocates a container to start AM
3. AM registers with RM
4. AM asks containers from RM
5. AM notifies NM to launch containers
6. Application code is executed in container



edureka!

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

Have a doubt? Raise a query.

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



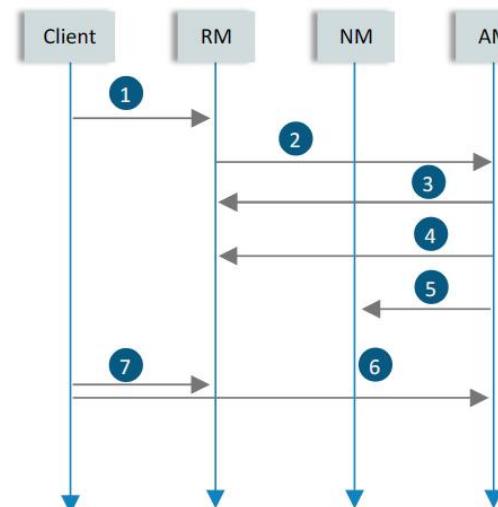
Course Content

- Introduction to Big Data
 - Class 1 Recording
 - Presentation
 - The Hype Behind Big Data
 - Why Hadoop?
 - Introduction to Apache Hadoop
 - HDFS Commands
 - Hadoop Tutorial
 - Hadoop Ecosystem
 - HDFS Tutorial

YARN Application Workflow

➤ Execution Sequence

1. Client submits an application
 2. RM allocates a container to start AM
 3. AM registers with RM
 4. AM asks containers from RM
 5. AM notifies NM to launch containers
 6. Application code is executed in container
 7. Client contacts RM/AM to monitor application's status



edureka!

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

 Have a doubt? Raise a query

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 76 / 121 ← → []

YARN Application Workflow

Execution Sequence :

1. Client submits an application
2. RM allocates a container to start AM
3. AM registers with RM
4. AM asks containers from RM
5. AM notifies NM to launch containers
6. Application code is executed in container
7. Client contacts RM/AM to monitor application's status
8. AM unregisters with RM

```
sequenceDiagram
    participant Client
    participant RM
    participant NM
    participant AM
    Client->>RM: 1. Client submits an application
    RM->>AM: 2. RM allocates a container to start AM
    AM->>RM: 3. AM registers with RM
    AM->>RM: 4. AM asks containers from RM
    RM->>NM: 5. AM notifies NM to launch containers
    NM->>AM: 6. Application code is executed in container
    Client->>RM: 7. Client contacts RM/AM to monitor application's status
    AM->>RM: 8. AM unregisters with RM
```

edureka!

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

Have a doubt? Raise a query.

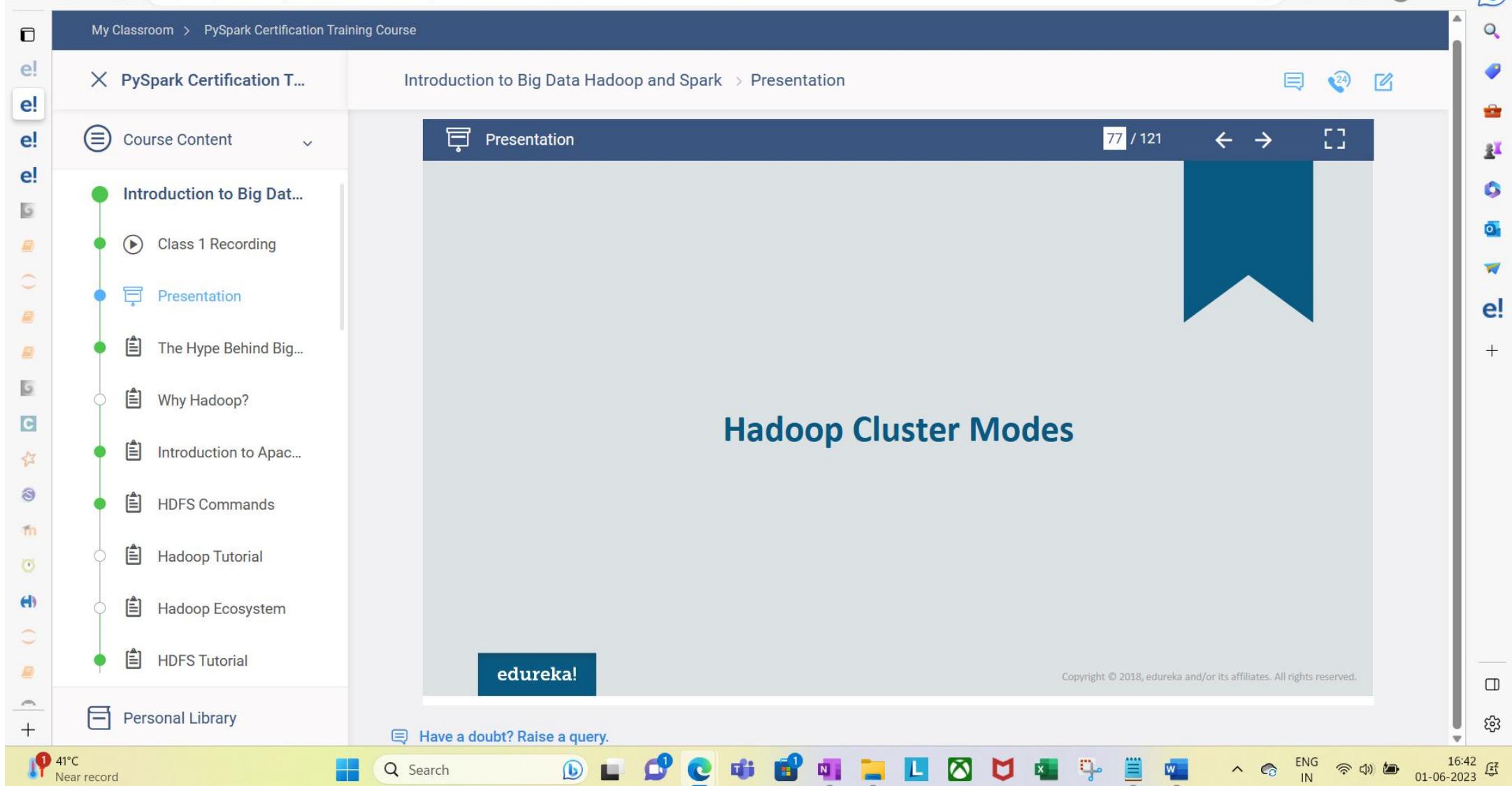
41°C Near record

Search

Cloud

ENG IN

16:42 01-06-2023



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

78 / 121



Hadoop Cluster Modes

Standalone (or Local) Mode

- No daemons, everything runs in a single JVM
- Suitable for running MapReduce programs during development
- Has no DFS or Distributed File System

Pseudo Distributed Mode

- All Hadoop daemons run on the local machine

Multi-Node Cluster Mode

- Hadoop daemons run on a cluster of machines

edureka!

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

Have a doubt? Raise a query.

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

- Introduction to Big Data...
- (▶) Class 1 Recording
- (▶) Presentation
- The Hype Behind Big...
- (📄) Why Hadoop?
- (📄) Introduction to Apac...
- (📄) HDFS Commands
- (📄) Hadoop Tutorial
- (📄) Hadoop Ecosystem
- (📄) HDFS Tutorial

Personal Library

Presentation

79 / 121



Real-Time Hadoop Cluster Deployment

edureka!

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

Have a doubt? Raise a query.

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

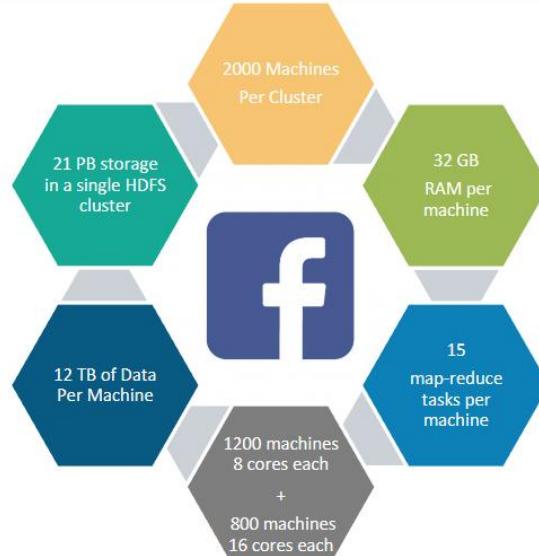
My Classroom > PySpark Certification Training Course

PySpark Certification T... Course Content

Introduction to Big Data Hadoop and Spark > Presentation

Presentation 80 / 121

Hadoop Cluster : Facebook Use Case



That's a total of *more than 21 PB* of *configured storage capacity*, this is larger than the previously known Yahoo!'s cluster of 14 PB

edureka!

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

Have a doubt? Raise a query.

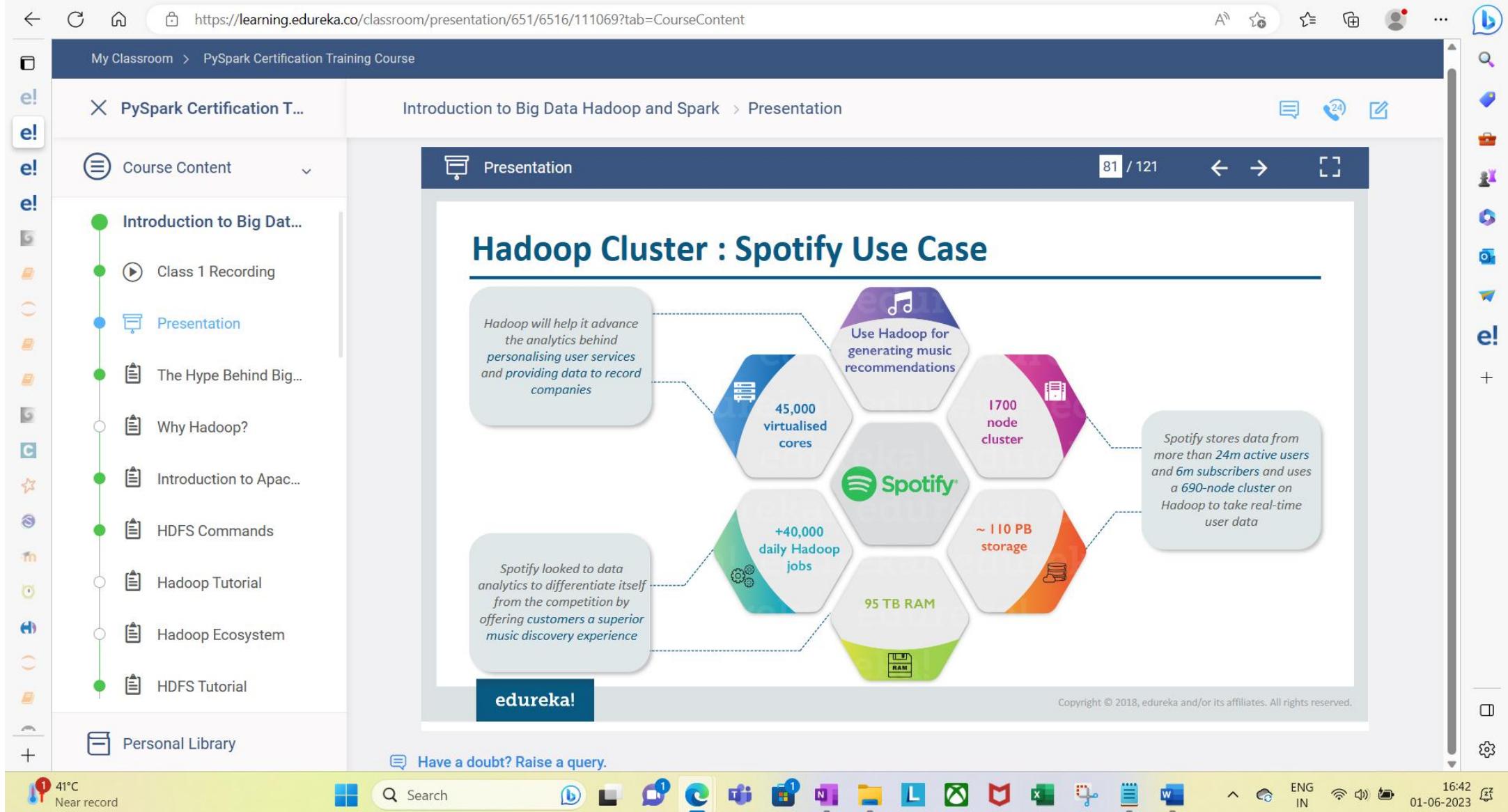
41°C Near record

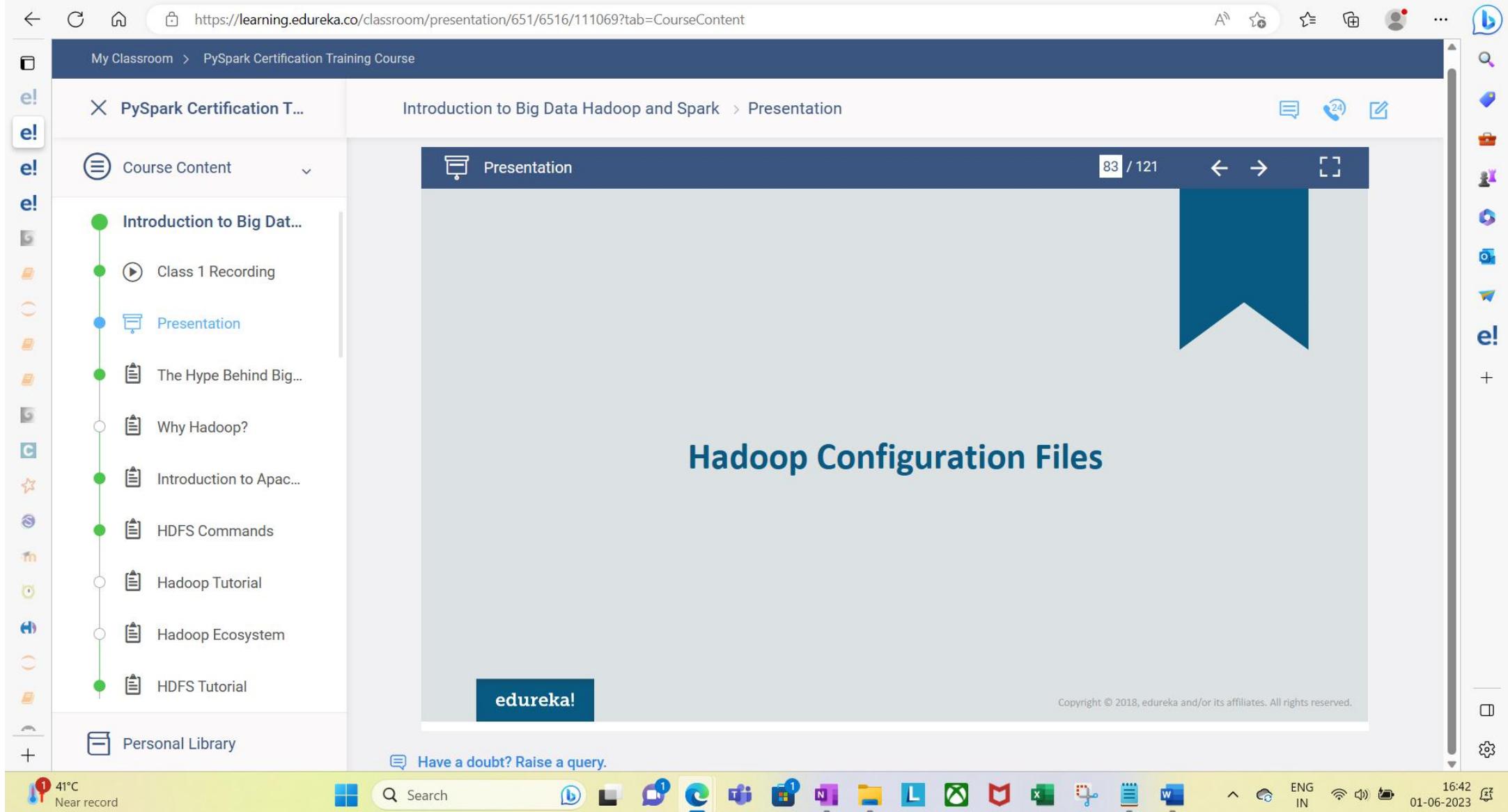
Search

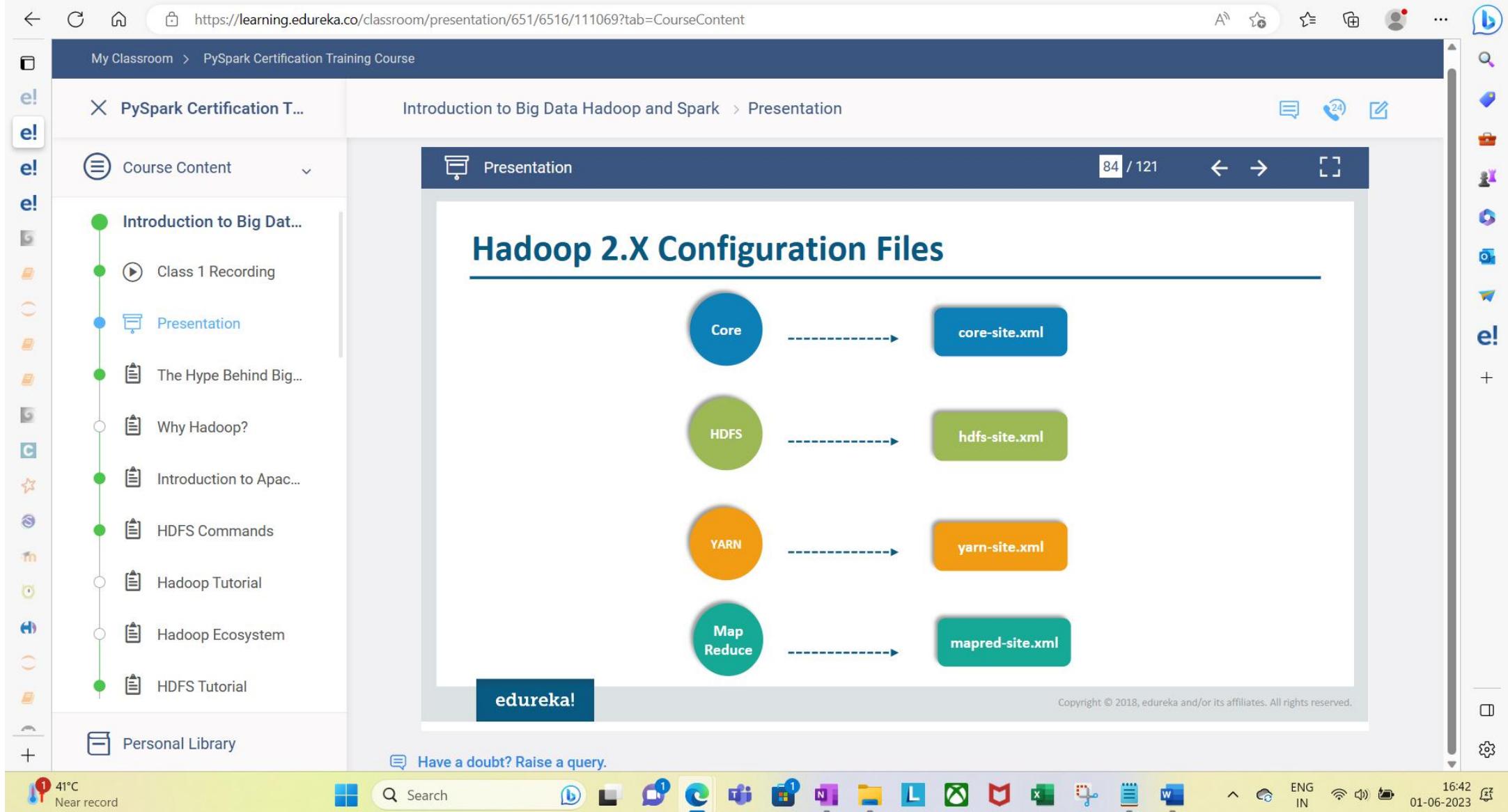
edureka!

ENG IN

01-06-2023 16:42







My Classroom > PySpark Certification Training Course

Introduction to Big Data Hadoop and Spark > Presentation



 Course Content

- Introduction to Big Data
 - Class 1 Recording
 - Presentation
 - The Hype Behind Big Data
 - Why Hadoop?
 - Introduction to Apache Hadoop
 - HDFS Commands
 - Hadoop Tutorial
 - Hadoop Ecosystem
 - HDFS Tutorial

Hadoop 2.X Configuration Files : Description

Configuration Filenames	Description
<code>hadoop-env.sh</code>	Environment variables that are used in the scripts to run Hadoop.
<code>core-site.xml</code>	Configuration settings for Hadoop Core such as I/O settings that are common to HDFS and MapReduce.
<code>hdfs-site.xml</code>	Configuration settings for HDFS daemons, the namenode, the secondary namenode and the data nodes.
<code>mapred-site.xml</code>	Configuration settings for MapReduce Applications.
<code>yarn-site.xml</code>	Configuration settings for ResourceManager and NodeManager.
<code>masters</code>	A list of machines (one per line) that each run a secondary namenode.
<code>slaves</code>	A list of machines (one per line) that each run a Datanode and a NodeManager.

edureka!

 Have a doubt? Raise a query

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apache...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Hadoop Web UI Parts

Service	Servers	Default Used ports	Protocol	Description
NameNode WebUI	Master Nodes (NameNode and any back-up NameNodes)	50070	http	Web UI to look at current status of HDFS, explore file system
DataNode	All Slave Nodes	50075	http	Data Node WebUI to access the status, logs etc.
ResourceManager Web UI	Cluster Level resource manager	8088	http	Web UI for Resource-Manager and for application submissions
NodeManager	Monitors resources on Data Node	8042	TCP	Node information, List of Applications and List of containers
MapReduce JobHistory Server	Get status on finished applications	19888	TCP	Providing logs of important events in MapReduce job execution and associated profiling metrics

edureka!

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

Have a doubt? Raise a query.



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

- Introduction to Big Data
 - Class 1 Recording
 - Presentation
 - The Hype Behind Big Data
 - Why Hadoop?
 - Introduction to Apache Hadoop
 - HDFS Commands
 - Hadoop Tutorial
 - Hadoop Ecosystem
 - HDFS Tutorial

 Presentation

88 / 121

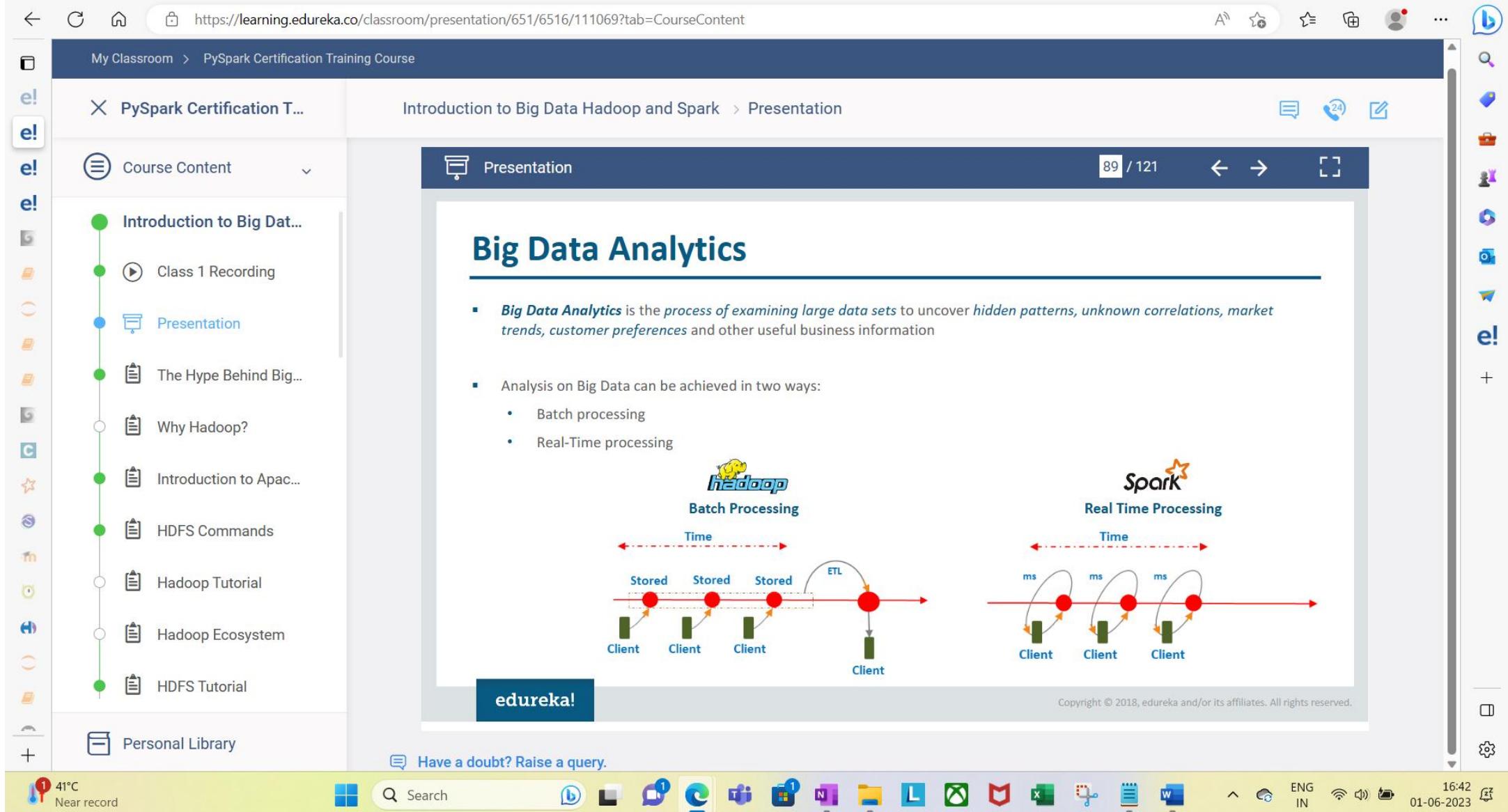


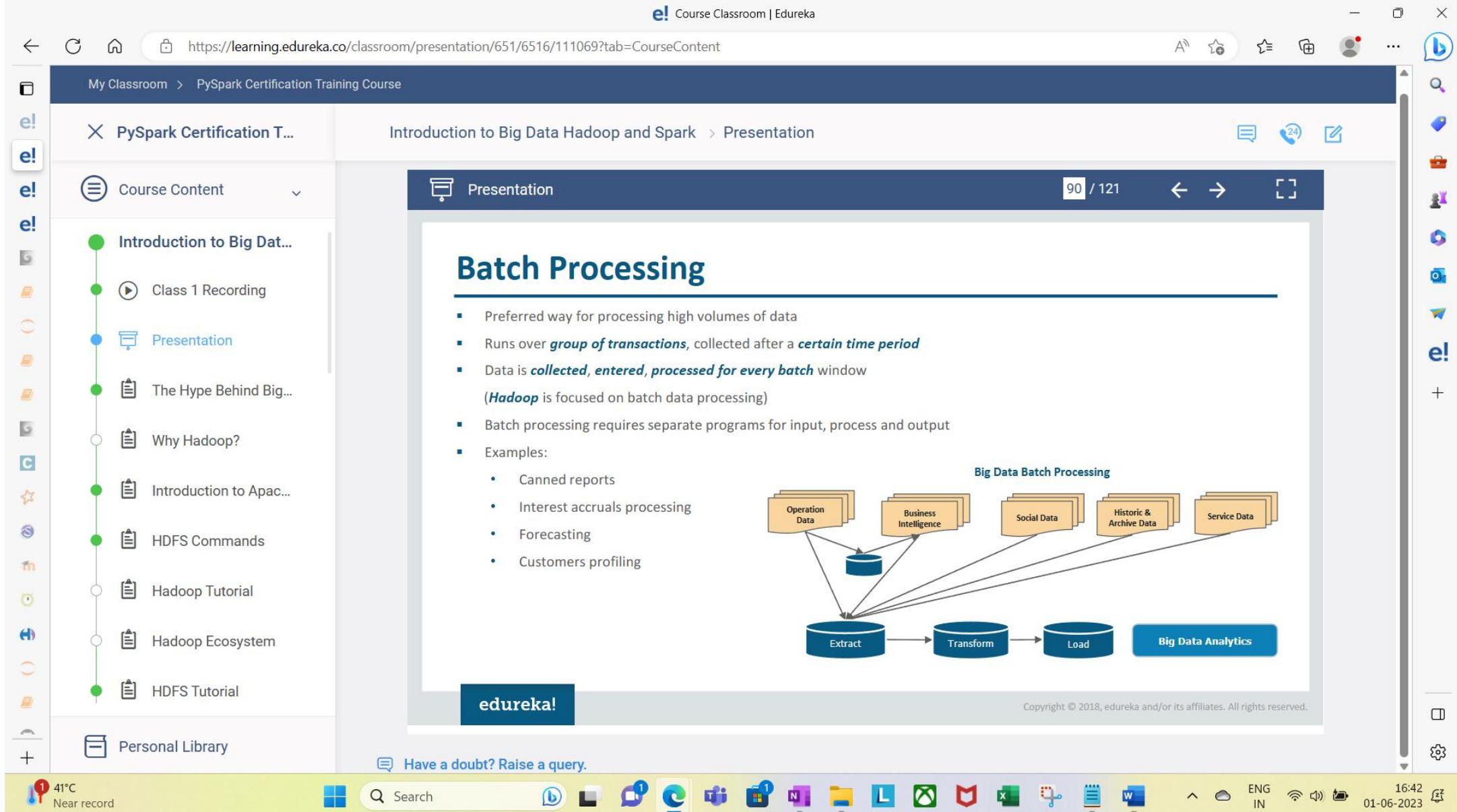
**Now that we have understood Hadoop,
let's see how Big Data Analytics can be achieved
through Batch & Real-Time Processing**

edureka!

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

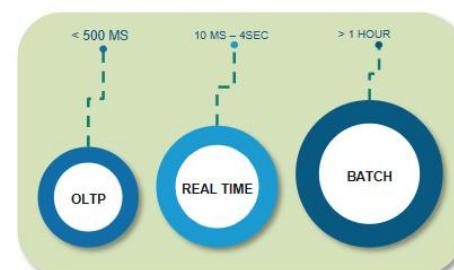
 Have a doubt? Raise a query





Real-Time Processing

- It involves a **continual input**, process and output of data
- Data processing time is much smaller (**in fractions of seconds**)
- One such example is a **complex event processing** (CEP combines data from multiple sources to infer events or patterns that suggest more complicated circumstances) platform
- Another example is **operational intelligence** (OI is a form of real-time dynamic, business analytics that delivers visibility and insight into business operations) platforms
- There are continuously running programs which keep consuming data from streams
- They keep running forever, till they are manually stopped



edureka!

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

Have a doubt? Raise a query.

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

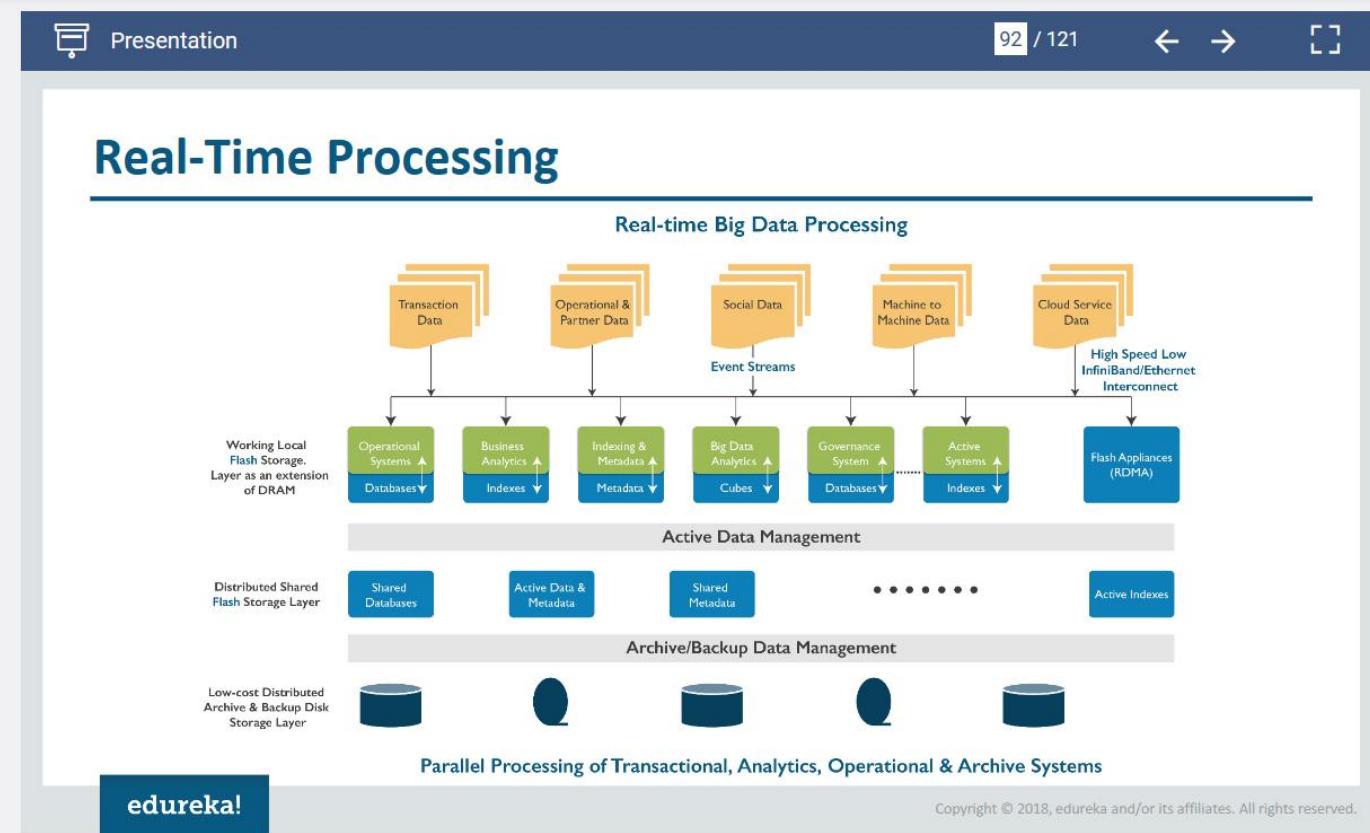
HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library



Have a doubt? Raise a query.

← ⌛ 🏠 🔍 https://learning.edureka.co/course/presentation/651/6516/111069?tab=CourseContent ⌛ 🏠 🔍 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

e! e!

Course Content

- Introduction to Big Dat...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

93 / 121 ← → []

There are other alternatives, then why go for Spark?

edureka!

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

Have a doubt? Raise a query.

41°C Near record

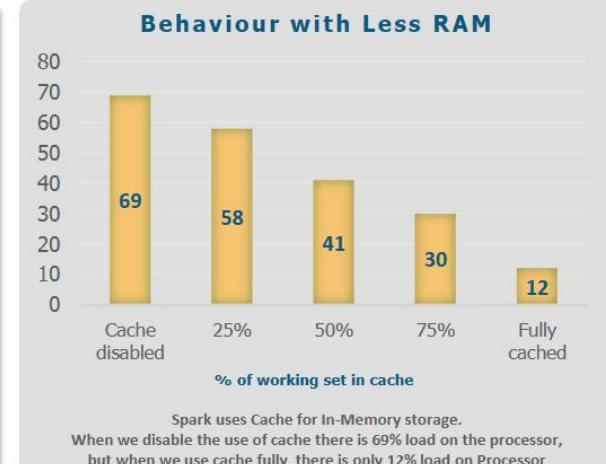
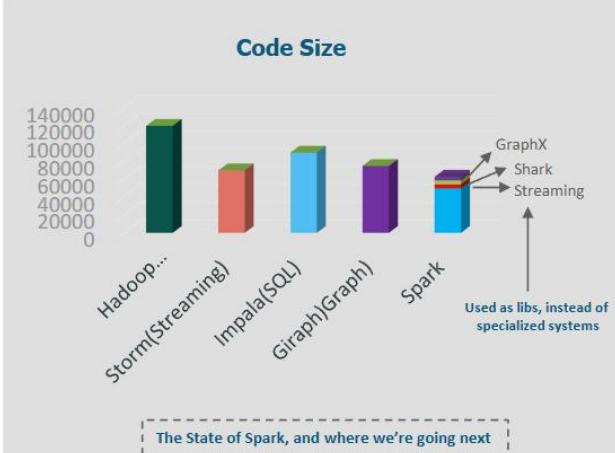
Search

edureka!

16:42 01-06-2023



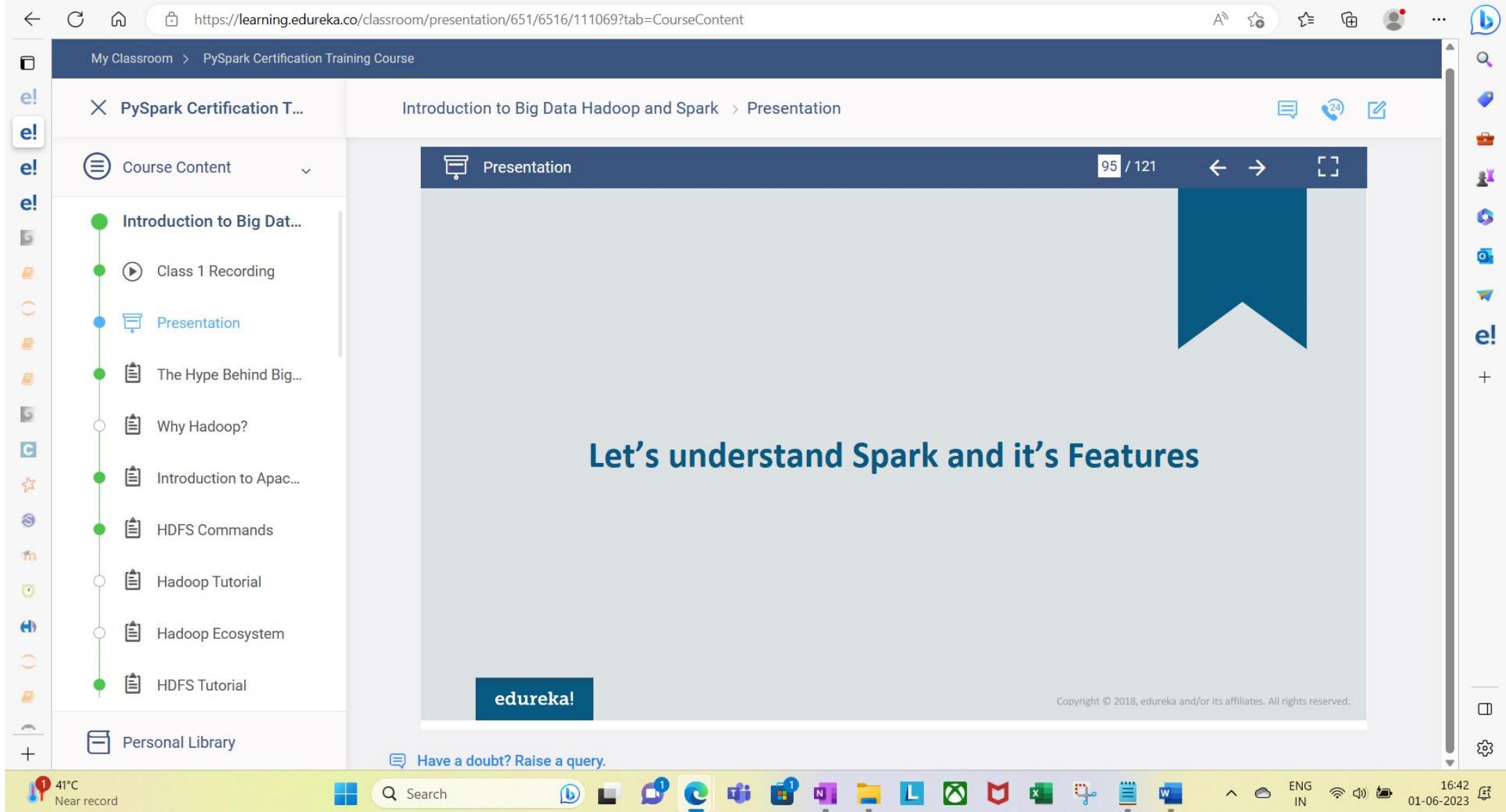
How Spark Differs?



edureka

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

 Have a doubt? Raise a query



Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

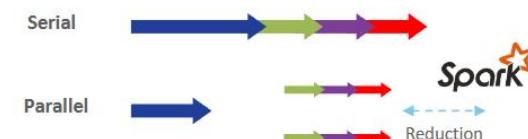
Personal Library

What is Spark?

- Apache Spark is an *open-source cluster-computing framework* for both *Batch* as well as *real time processing*
- Spark provides an interface for *programming entire clusters* with *implicit data parallelism* and *fault-tolerance*



Real Time Processing In Spark

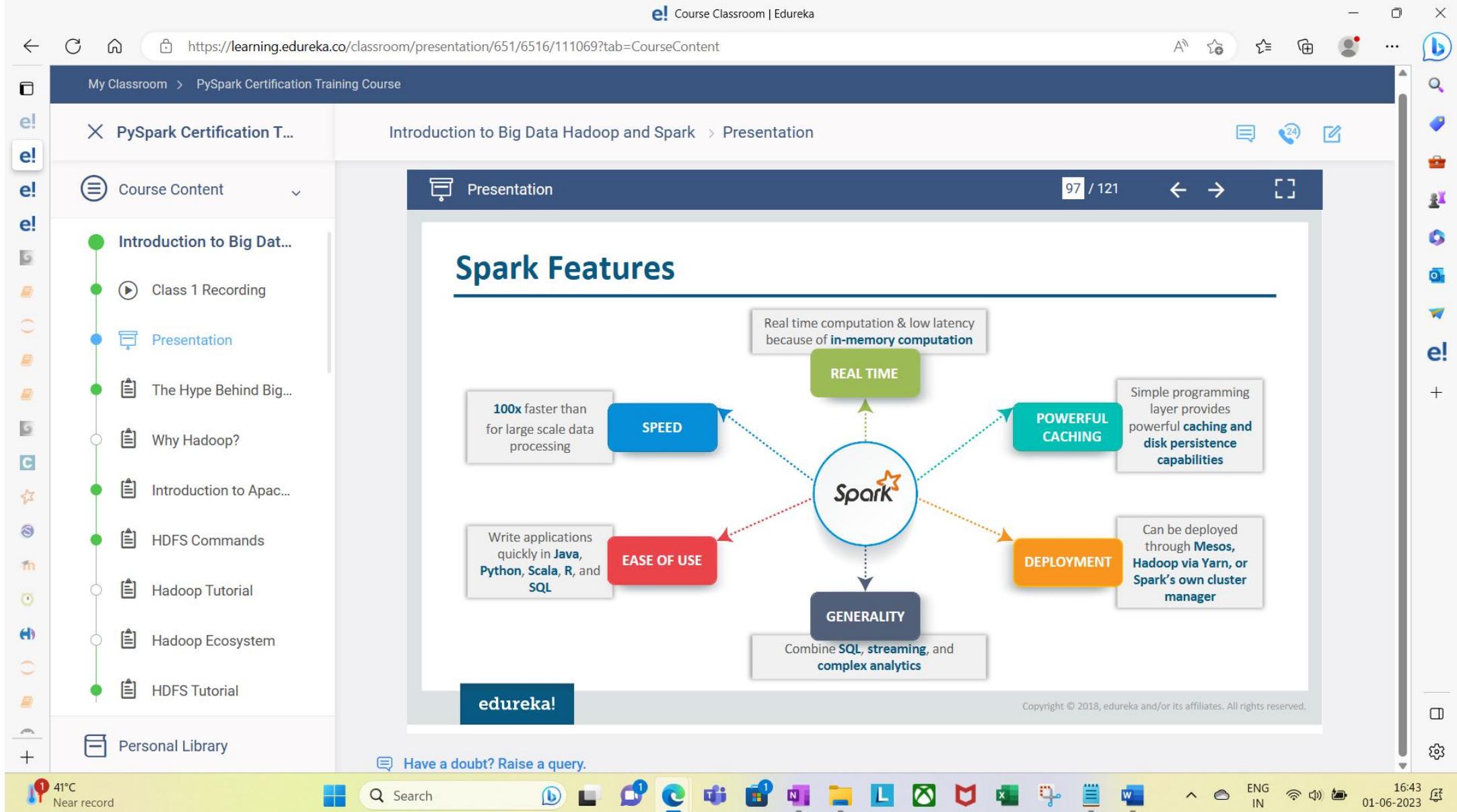


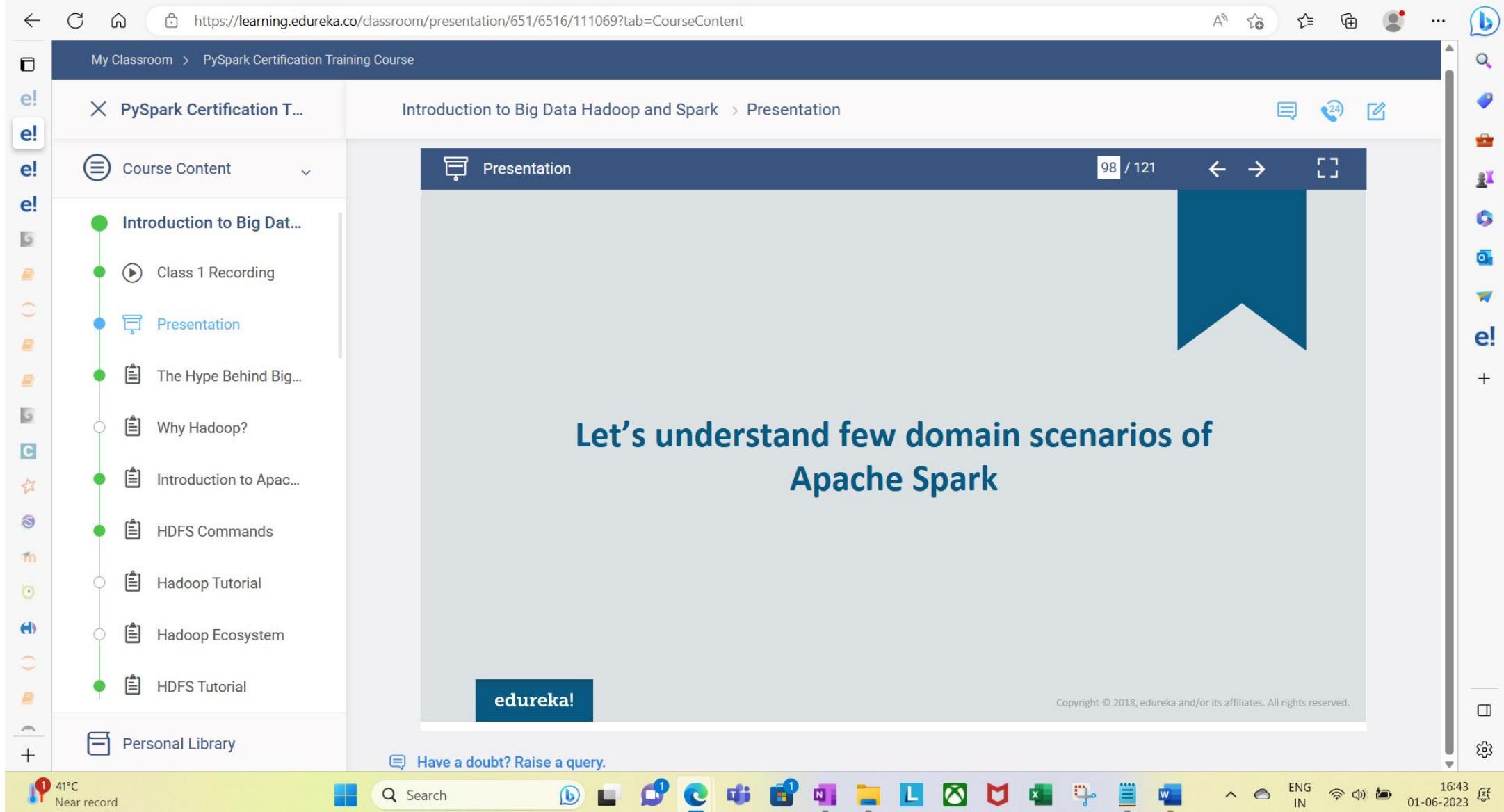
Data Parallelism in Spark

edureka!

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

Have a doubt? Raise a query.





My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

99 / 121



Apache Spark Domain Scenarios

- From small startups to *Fortune 500s* all are adopting Apache Spark to *build, scale and innovate* their big data applications



edureka!

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

Have a doubt? Raise a query.

← ⌛ 🏠 🔍 https://learning.edureka.co/course/presentation/651/6516/111069?tab=CourseContent ⌛ 🏠 🔍 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

e! e!

Course Content

- Introduction to Big Dat...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial

Personal Library

Presentation

100 / 121 ← → []

Let's understand Apache Spark with the help of a use-case from the Media & Entertainment Industry

edureka!

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

Have a doubt? Raise a query.

41°C Near record

Search

edureka!

16:43 ENG IN 01-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

101 / 121



Problems Faced at Yahoo!

Yahoo! properties are highly personalized to maximize relevance

edureka!

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

Have a doubt? Raise a query.



X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

- Introduction to Big Data
 - Class 1 Recording
 - Presentation
 - The Hype Behind Big Data
 - Why Hadoop?
 - Introduction to Apache Hadoop
 - HDFS Commands
 - Hadoop Tutorial
 - Hadoop Ecosystem
 - HDFS Tutorial

Problems Faced at Yahoo!



Yahoo! properties are highly personalized to maximize relevance

- ① The *algorithms* used to provide personalization(targeted advertisement and personalized content) are **highly sophisticated**
 - ② *Relevance Model* must be *updated* frequently as stories, etc with change in time



edureka

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

 Have a doubt? Raise a query



X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

- Introduction to Big Data
 - Class 1 Recording
 - Presentation
 - The Hype Behind Big Data
 - Why Hadoop?
 - Introduction to Apache Hadoop
 - HDFS Commands
 - Hadoop Tutorial
 - Hadoop Ecosystem
 - HDFS Tutorial

Problems Faced at Yahoo!

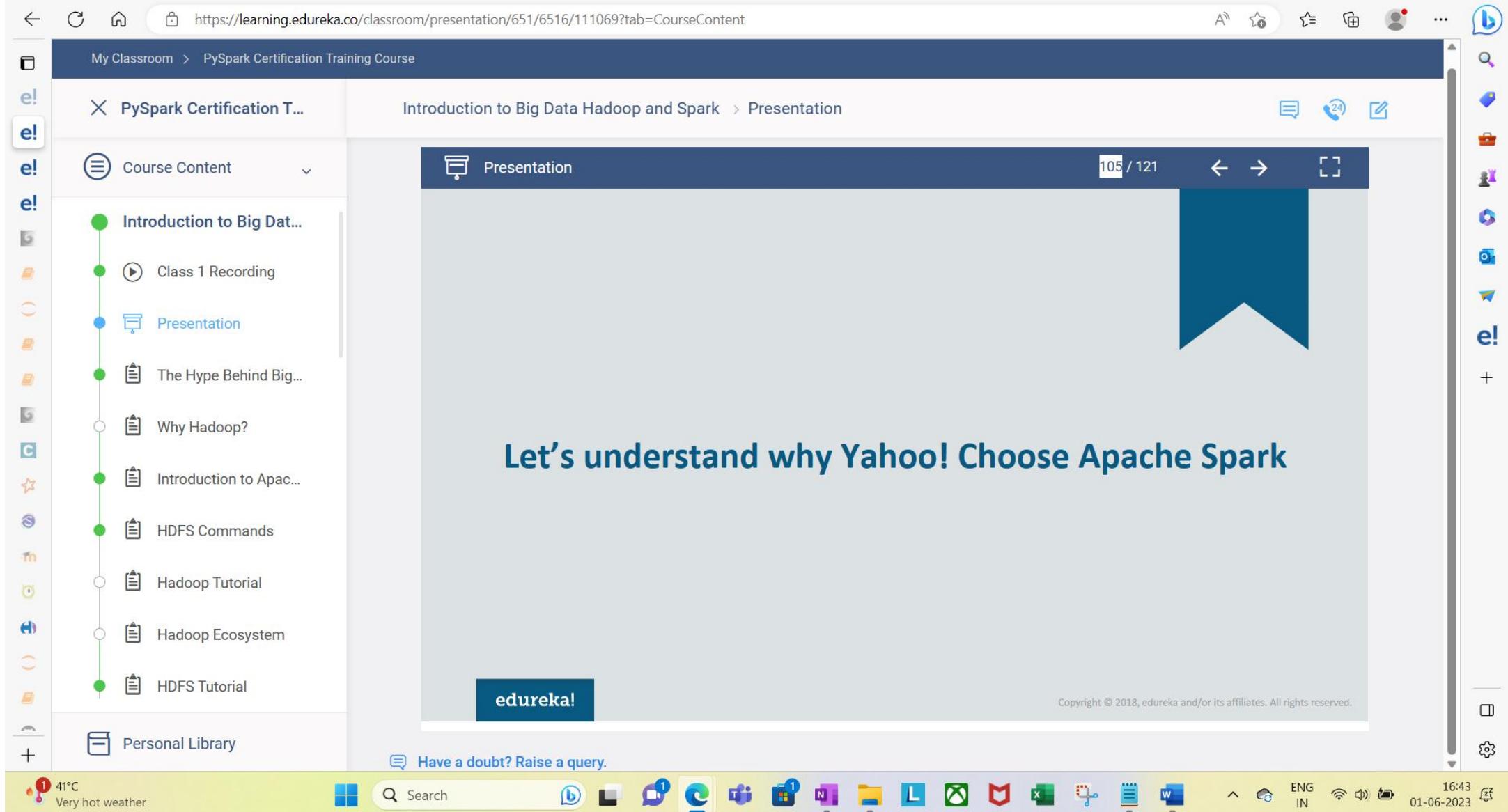


Yahoo! properties are highly personalized to maximize relevance

- ① The *algorithms* used to provide personalization(targeted advertisement and personalized content) are ***highly sophisticated***
 - ② *Relevance Model* must be *updated* frequently as stories, etc with change in time
 - ③ Yahoo! has over *150 petabytes* of data stored on a *35,000-node Hadoop cluster* which should be accessed efficiently to:
 - *Avoid latency* caused by data movement and
 - *Gain insights* from data in cost-effective manner

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

 Have a doubt? Raise a query



My Classroom > PySpark Certification Training Course

Introduction to Big Data Hadoop and Spark > Presentation



106 / 121

← →

1

Why Apache Spark?

Yahoo! looked to Spark to improve *performance* of its *iterative model training*

edureka!

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

 Have a doubt? Raise a query

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data Hadoop and Spark

Class 1 Recording

Presentation

The Hype Behind Big Data

Why Hadoop?

Introduction to Apache Spark

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

107 / 121



Why Apache Spark?

Yahoo! looked to Spark to improve **performance** of its **iterative model training**



The machine learning algorithm for news personalization required **15000 lines of C++ code**

edureka!

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

Have a doubt? Raise a query.

41°C
Very hot weather

Search



ENG IN 16:43 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Course Content

Introduction to Big Data Hadoop and Spark > Presentation

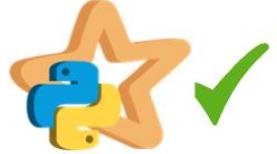
Presentation 108 / 121

Why Apache Spark?

Yahoo! looked to Spark to improve **performance** of its *iterative model training*

The machine learning algorithm for news personalization required **15000 lines of C++ code**

PySpark has reduced the length of code to more than half the size

edureka!

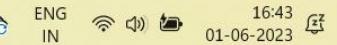
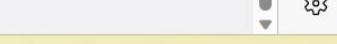
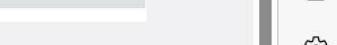
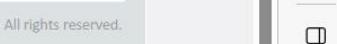
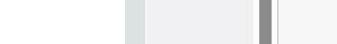
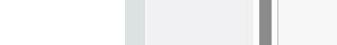
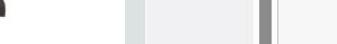
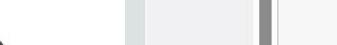
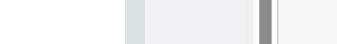
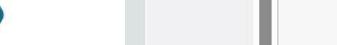
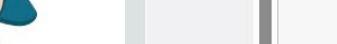
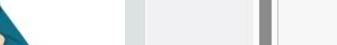
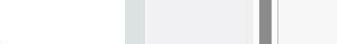
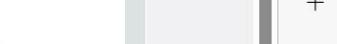
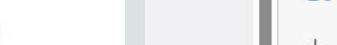
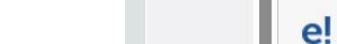
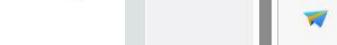
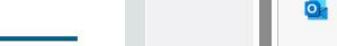
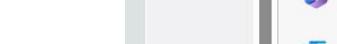
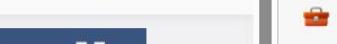
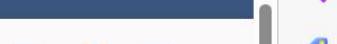
Copyright © 2018, edureka and/or its affiliates. All rights reserved.

Have a doubt? Raise a query.

41°C Very hot weather

Search

16:43 01-06-2023



X PySpark Certification T...

Course Content

Introduction to Big Dat...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

109 / 121



Why Apache Spark?

Yahoo! looked to Spark to improve **performance** of its **iterative model training**



PySpark has reduced the length of code to more than half the size



The algorithm was ready for production use in just **30 minutes** of training, on a **hundred million datasets**

Its rich API is available in several programming languages, has resilient in-memory storage options and is compatible with Hadoop through YARN and the Spark-YARN project

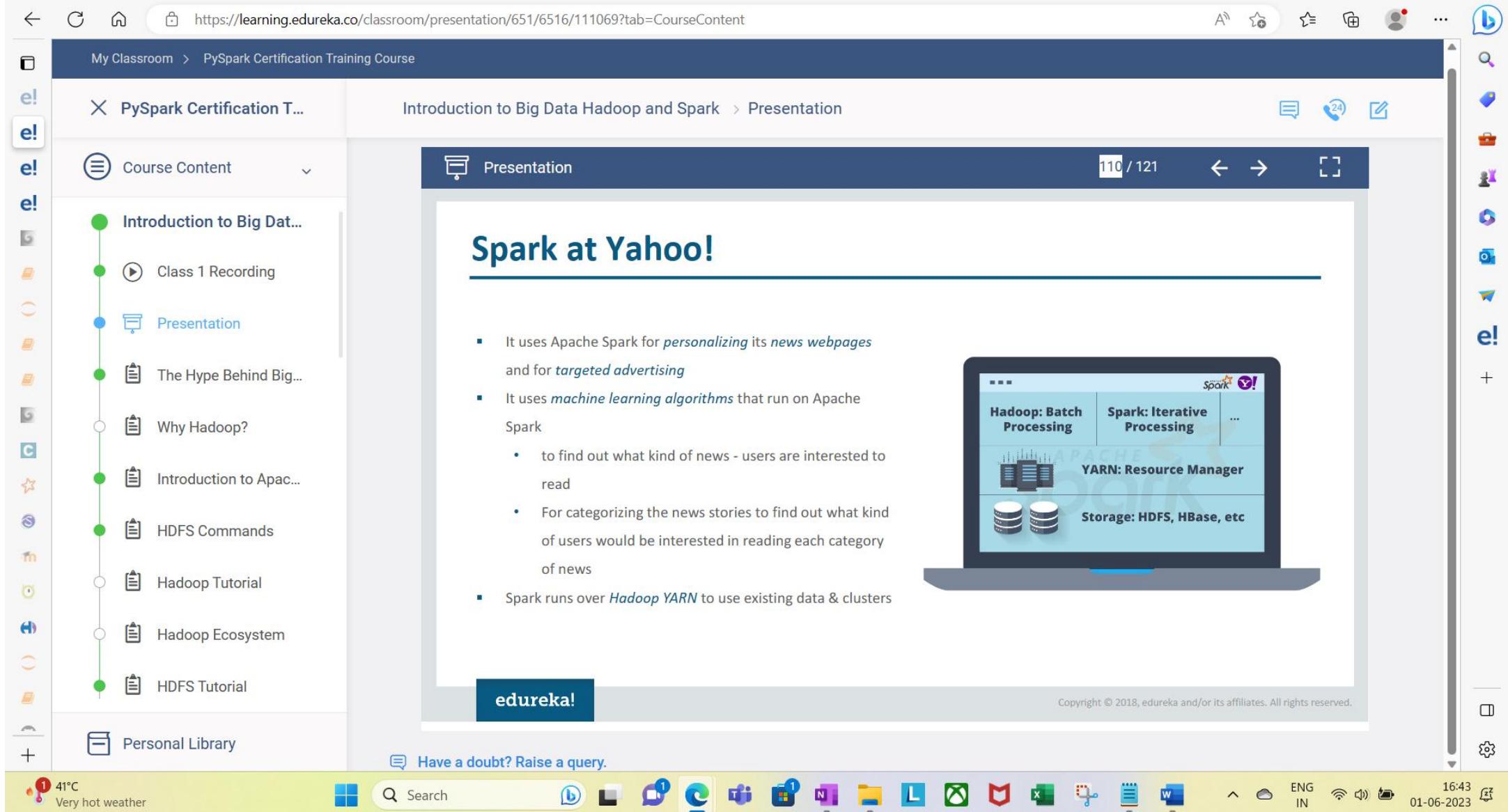
Copyright © 2018, edureka and/or its affiliates. All rights reserved.

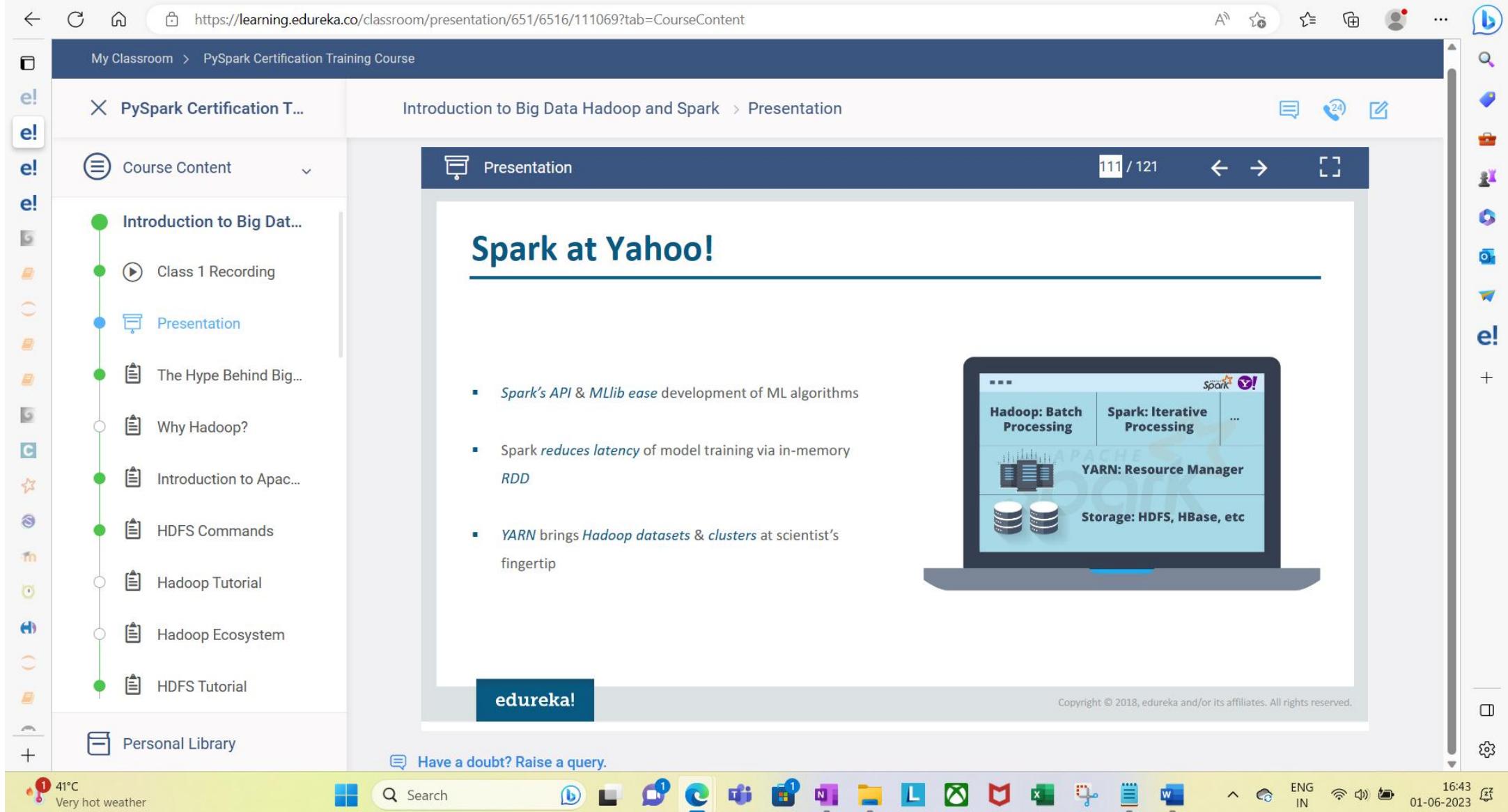
Have a doubt? Raise a query.

 41°C
Very hot weather

Search

ENG
IN16:43
01-06-2023





My Classroom > PySpark Certification Training Course

PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Presentation 112 / 121 ← → []

Spark at Yahoo! (Ad Analytics)

- Yahoo! Ads wanted *interactive BI* on *terabytes of data*
- Chose Shark (Hive on Spark) to provide this through standard *Hive server API + Tableau*
- **Result:** Were able to *query their ad visit data interactively*

LARGE HADOOP CLUSTER

Hadoop (Pig, Hive, MR) | Spark

YARN

Historical DW (HDFS)

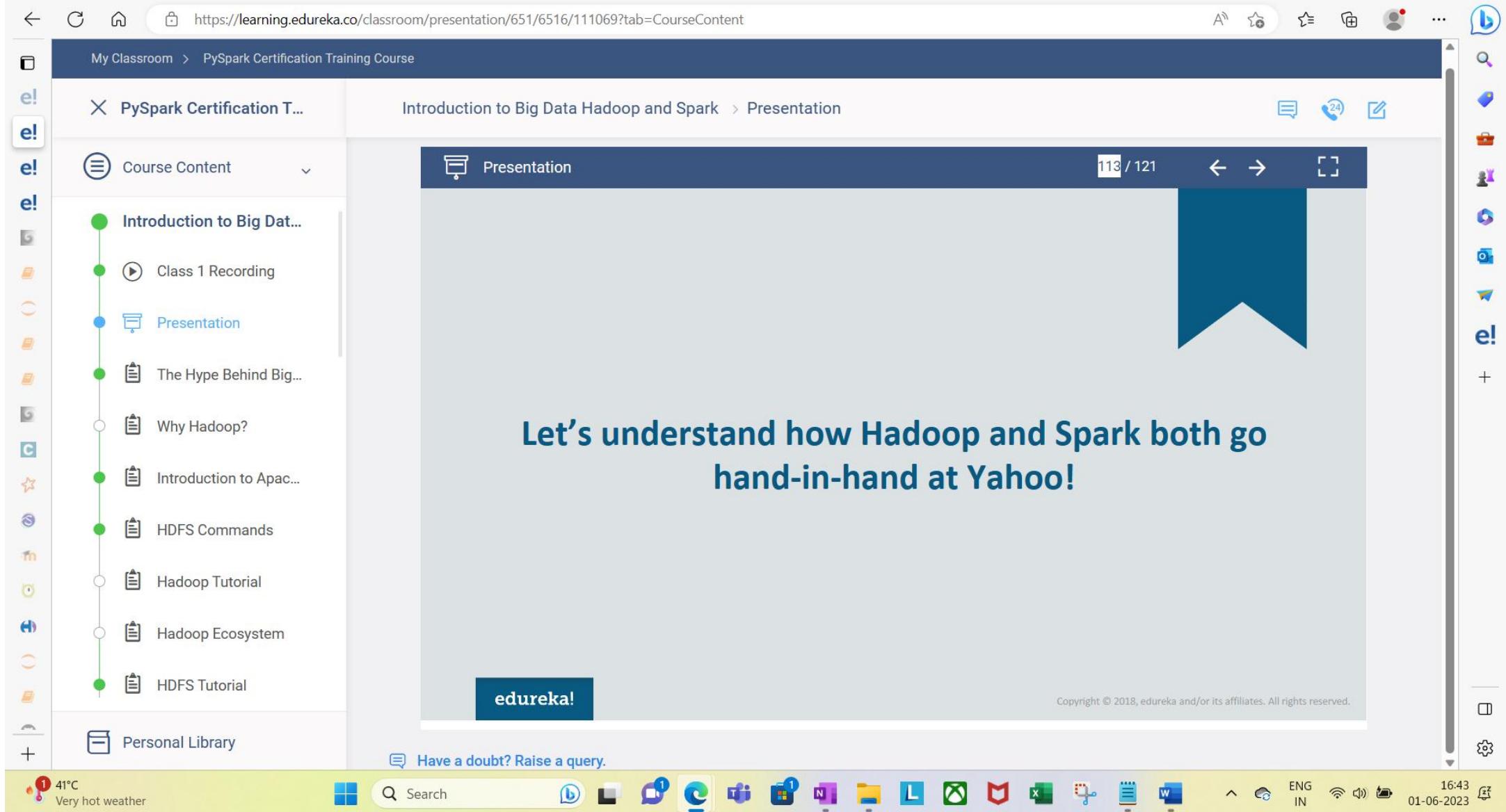
Satellite Shark Cluster

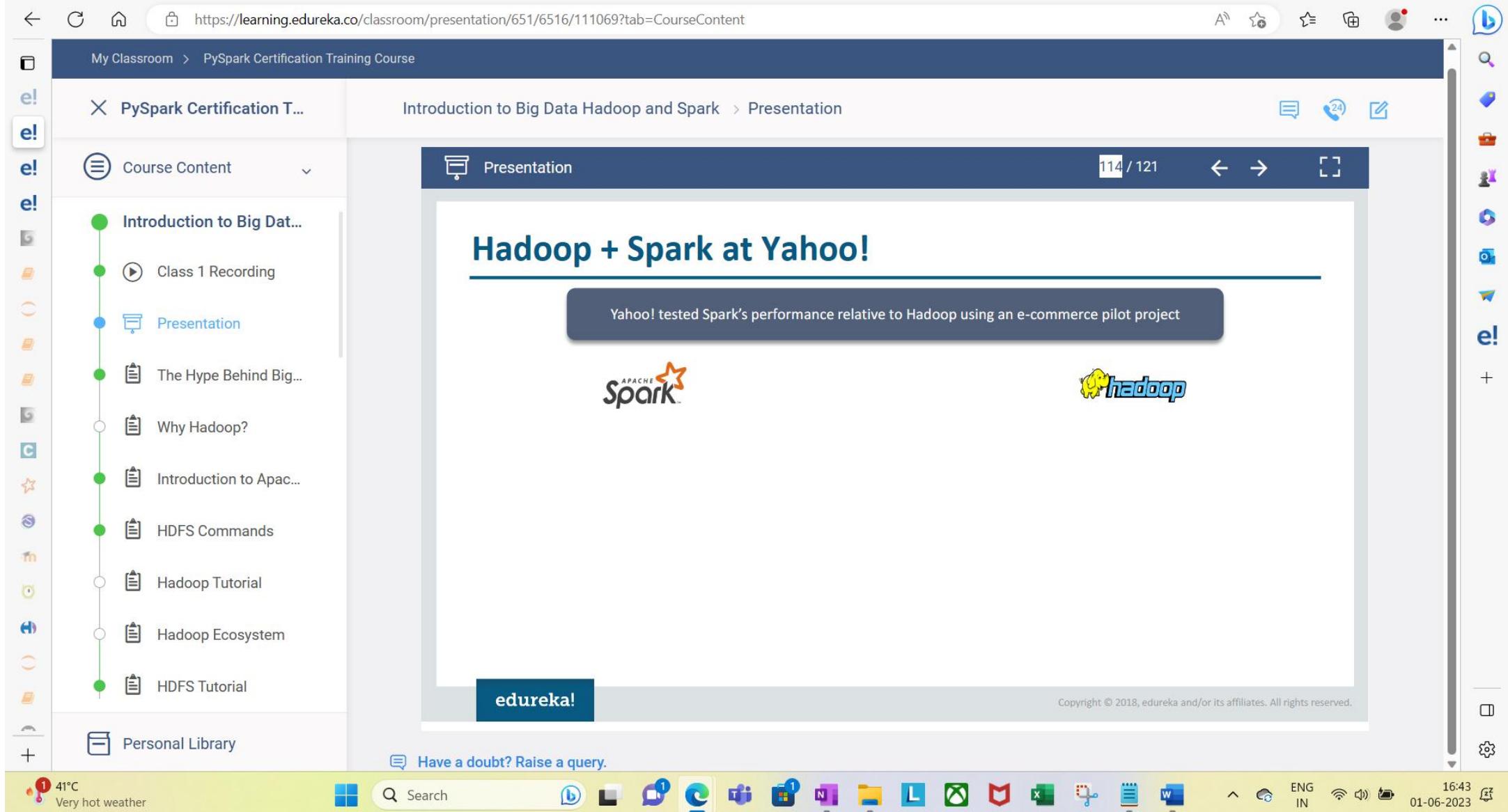
Satellite Shark Cluster

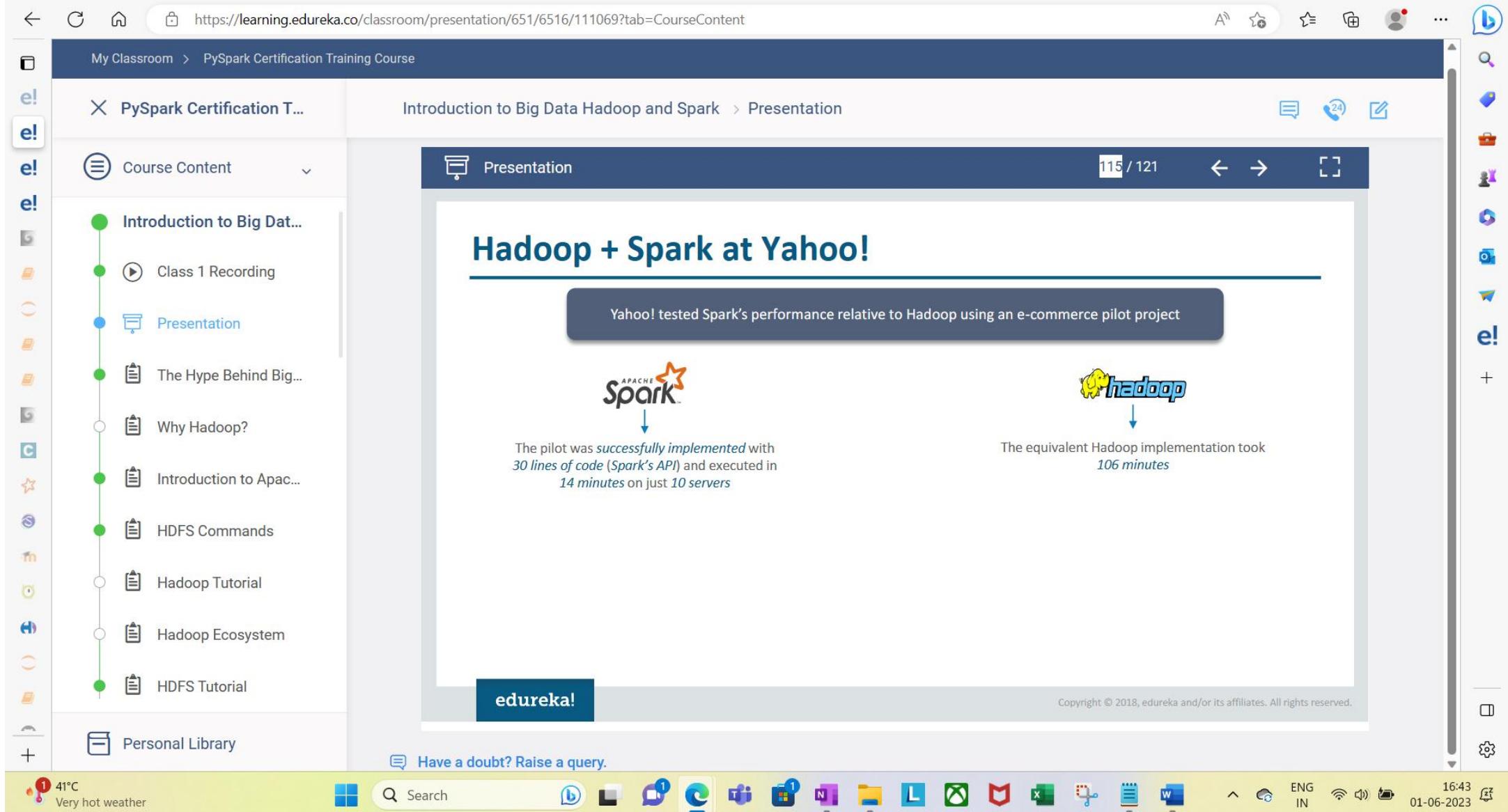
edureka!

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

Have a doubt? Raise a query.







X PySpark Certification T...

Introduction to Big Data Hadoop and Spark > Presentation



Course Content

Introduction to Big Data...

Class 1 Recording

Presentation

The Hype Behind Big...

Why Hadoop?

Introduction to Apac...

HDFS Commands

Hadoop Tutorial

Hadoop Ecosystem

HDFS Tutorial

Personal Library

Presentation

116 / 121



Hadoop + Spark at Yahoo!



Yahoo! tested Spark's performance relative to Hadoop using an e-commerce pilot project

The pilot was successfully implemented with 30 lines of code (Spark's API) and executed in 14 minutes on just 10 servers



The equivalent Hadoop implementation took 106 minutes



- While these improvements are impressive, Yahoo! isn't abandoning its Hadoop cluster for Spark
- There is a clear *need for both* types of workloads
- *Spark* will be the *preferred* technology for *iterative processing*, while *Hadoop* continues to fulfil its niche for *batch data processing* tasks
- What's interesting is that both types of tasks run on the same Hadoop cluster through YARN

edureka!

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

Have a doubt? Raise a query.

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Introduction to Big Data Hadoop and Spark > Presentation

Course Content

- Introduction to Big Data...
- Class 1 Recording
- Presentation
- The Hype Behind Big...
- Why Hadoop?
- Introduction to Apac...
- HDFS Commands
- Hadoop Tutorial
- Hadoop Ecosystem
- HDFS Tutorial
- Personal Library

117 / 121 ← → []

Summary

What is Big Data?

Big data is the term for a collection of data sets so large and complex that it becomes difficult to process using traditional database management tools or conventional data processing applications.

What is Hadoop?

- Apache Hadoop is an open-source framework that allows distributed processing of large data sets across clusters of commodity computers.
- It provides distributed storage and distributed processing.

Hadoop 2.x Core Components

```
graph TD; RM[ResourceManager] --> NN[NameNode]; RM --> NM[NodeManager]; NN --> HMaster[HMaster]; NM --> HMaster; NM --> DataNode[DataNode]; NM --> SecondaryNameNode[Secondary NameNode]; HMaster --> HDFS[HDFS]; HDFS --> Storage[Storage]; HDFS --> Processing[YARN];
```

Block Replication

Hadoop/HDFS Commands

```
Check the version of Hadoop: hadoop -version  
Check HDFS健康: hdfs dfs -testracks /<path>
```

What is Spark?

- Analyse Spark is an open-source cluster computing framework for real-time processing.
- Spark provides an interface for programming entire clusters with implicit data parallelism and fault tolerance.

edureka!

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

Have a doubt? Raise a query.

41°C Very hot weather

Search

16:43 01-06-2023

https://learning.edureka.co/classroom/presentation/651/6516/111069?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Course Content

Introduction to Big Data Hadoop and Spark > Presentation

Presentation 118 / 121

Further Reading

- Big Data and Hadoop Tutorials
 - <https://www.edureka.co/blog/hadoop-tutorial/>
- HDFS Architecture
 - <https://www.edureka.co/blog/apache-hadoop-hdfs-architecture/>
- Crack Hadoop Interviews
 - <https://www.edureka.co/blog/interview-questions/top-50-hadoop-interview-questions-2016/>
- PySpark Tutorial – Learn Apache Spark using Python
 - <https://www.edureka.co/blog/pyspark-tutorial/>
- Apache Spark with Hadoop
 - <https://www.edureka.co/blog/apache-spark-with-hadoop-why-it-matters/>

edureka!

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

Have a doubt? Raise a query.

41°C Very hot weather

Search

16:43 01-06-2023