

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Topics

- What is Machine Learning
- Features of Machine Learning
- Applications of Machine Learning
- Steps of Machine Learning
- Types of Machine Learning
- Mlib in Spark

edureka!

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



My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Objectives

After completing this module, you should be able to:

- Describe What is Machine Learning and Where It is Used
- Describe Different Types of Machine Learning Techniques
- Understand Face Detection : USE CASE
- Describe MLlib
- Explain Features of MLlib and MLlib Tools
- Get an overview of various ML algorithms supported by MLlib



edureka!

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



My Classroom > PySpark Certification Training Course



X PySpark Certification T...



Course Content



Class 7 Recording



Presentation



Deep Dive into Spark ...



Understanding Apache...



Apache Spark Streami...



Apache Spark Streami...



In-class Project



Spark GraphX



Certification Project - ...



Personal Library

Machine Learning using Spark MLlib > Presentation



Presentation

7 / 78



Let us look at the Scenario from Edgeways: A Software Company

edureka!

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



31°C

Haze



Search



My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

Scenario of Edgeways

A professional working in **Edgeways, a software company** checks inbox for new mails or updates and becomes frustrated after seeing unnecessary mails



Hey Dear,
Congratulations!!!!
I am pleased to inform you that
you have won an amount of 100K
in a lottery ticket. You can send
your "account details" by replying
to same mail.
Best Regards
Friend

Hey Dear,
Surprises are waiting at your door!!!!
Reply to the same mail and see the
magic.
Best Regards
Anonymous

many more

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

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

I have been receiving a lot of unnecessary and disturbing mails since a long time
(shows the mails)

To IT Support & Helpdesk

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

Edgeways: Challenge Faced

I have been receiving a lot of unnecessary and disturbing mails since a long time
(shows the mails)

To IT Support & Helpdesk

Sir, these are the "Spam Mails" which can even corrupt your system. We will try to find the solution for it at the earliest

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



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

- ▶ Class 7 Recording
 - 💻 Presentation
 - Deep Dive into Spark ...
 - Understanding Apache...
 - Apache Spark Streami...
 - Apache Spark Streami...
 - In-class Project
 - Spark GraphX
 - Certification Project - ...

Machine Learning using Spark MLlib > Presentation



Now, let us see how IT Support will solve the problem of the Spam Mails

edureka

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

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Machine Learning using Spark MLlib > Presentation



Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Probable Solution

Spam mails can be detected by manually setting a filter on some words

For Example,

If there are words like **lottery**, then mark the mail as a spam mail

Be careful with this message. Similar messages were used to steal people's personal information. Unless you trust the sender, don't click links or reply with personal information. [Learn more](#)

Greetings To You And Your Family,

Our names are John and Lisa Robinson of Munford, Tennessee winner of power-ball lottery game 2016 of the sum of \$528.8 million share of the prize in one lump-sum payment of \$327.8 million. We just commenced our charity donation scheme and we willing to give out cash grant of US\$500,000.00 each to (7) Lucky international recipients worldwide. You received this message because your E-mail ID have been listed as one of our (7) lucky winner in our charity donation.

Our aim is to raise the living standards of people across the world. This may be a surprise to you or a joke or hoax to you due to the scams in the internet this days, please have no doubt as this is very real. To confirm the legitimacy, visit my web page: <http://money.cnn.com/2016/05/01/news/largest-lottery-jackpots/>

After your confirmation, Kindly respond to this message with your Full Names and your contact address for more details on how to receive your cash grant.

Kindly accept my warmest Congratulations.

Warm Regards,
John & Lisa Robinson

SPAM

Such as 'lottery' keyword

edureka!

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



1

31°C

Haze



Search































































































































































































































































































My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

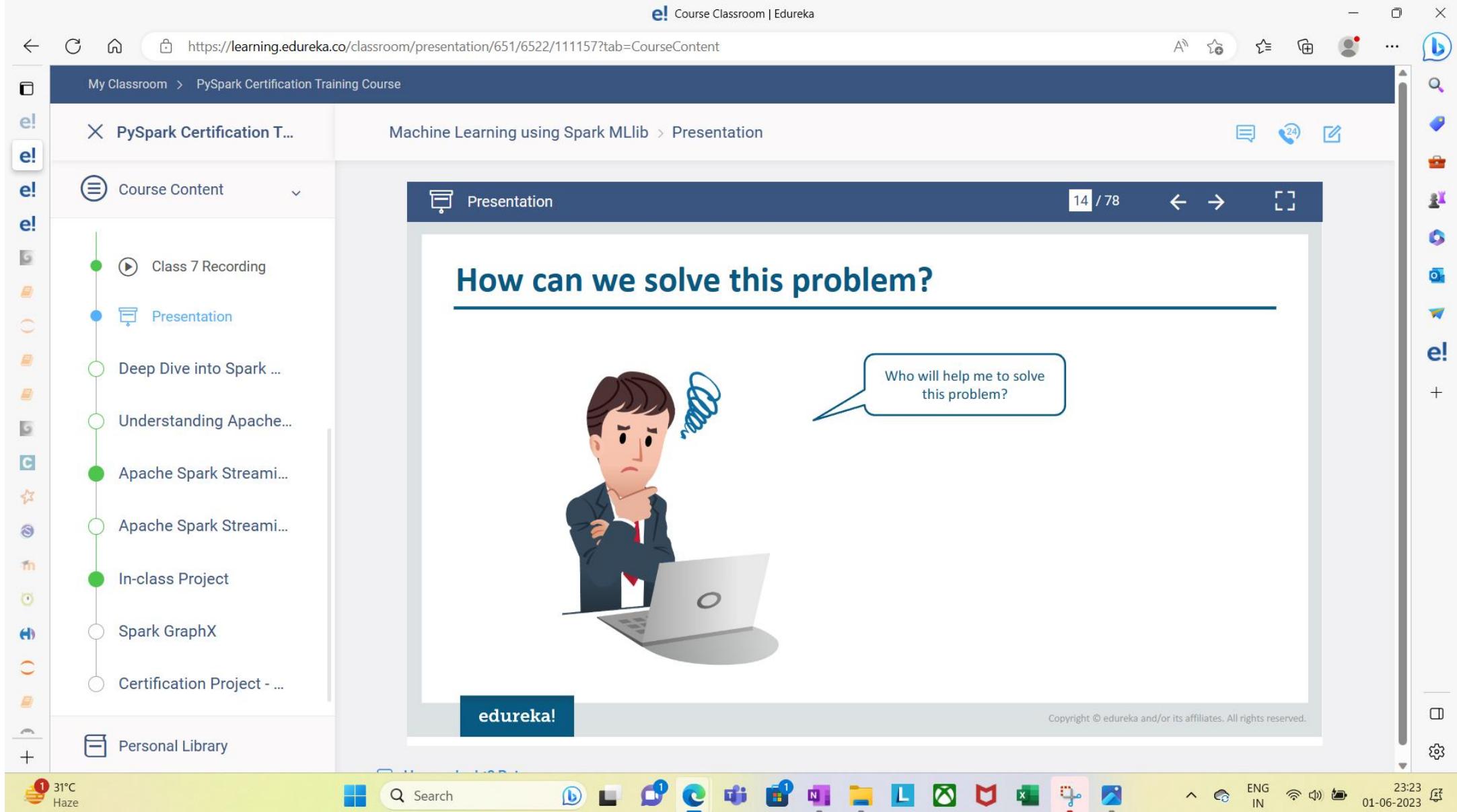
Course Content

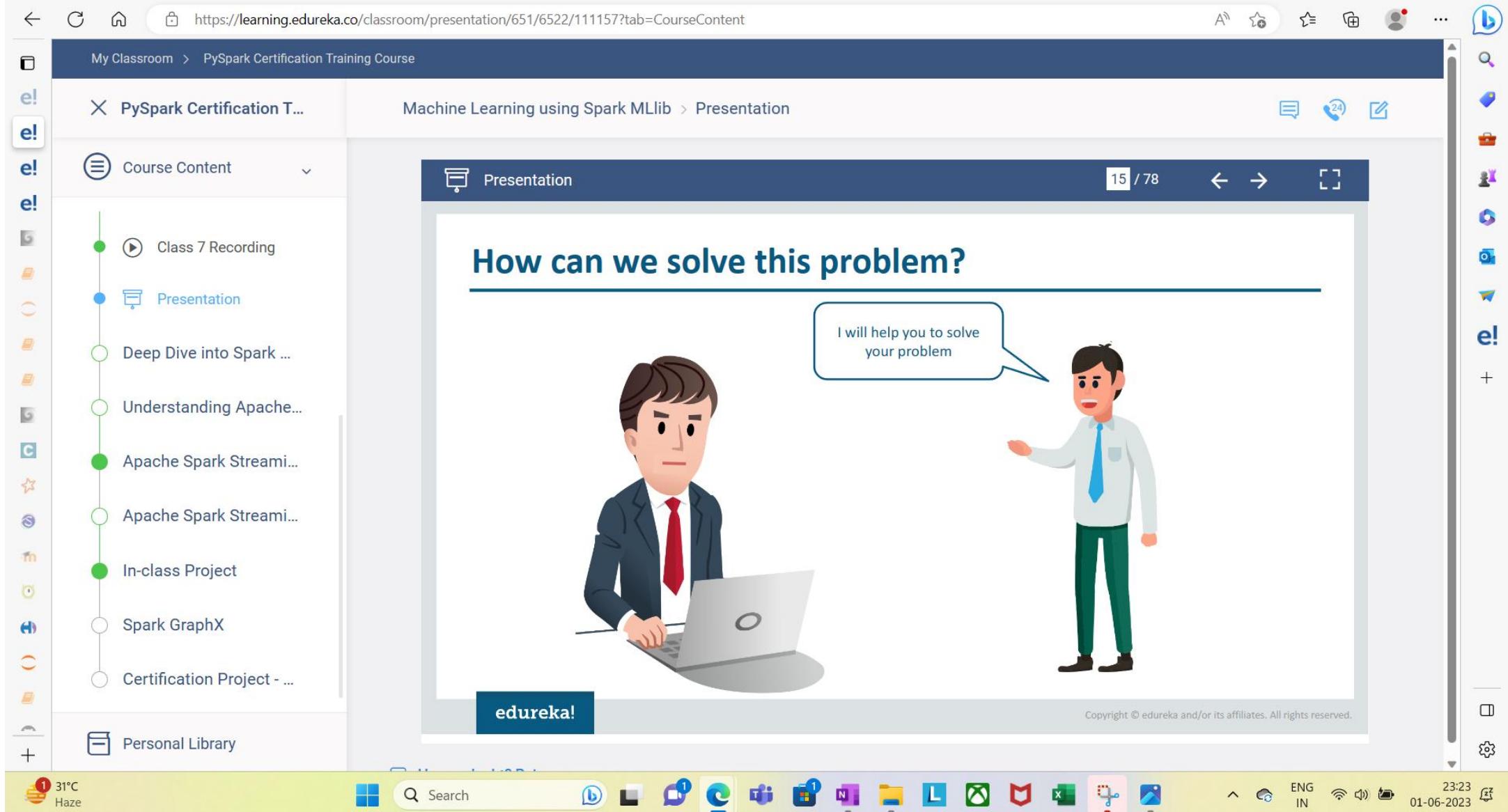
- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

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

13 / 78 ← →





My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Solution

edureka!

Labelled Data

The illustration shows a user profile icon pointing to a central area labeled 'Labelled Data'. This area contains several colorful envelope icons (red, yellow, green) representing emails. A robot character stands next to the data, with a thought bubble showing gears, symbolizing the machine learning process. Arrows point from the data area to two boxes: 'INBOX' (containing a green envelope) and 'SPAM' (containing a red envelope). Below the illustration, a caption reads: "Machine is programmed to learn from labelled data to create rules based on which we are able to classify the mails".

16 / 78 ← →

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

https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Solution

To get the required solution, Machine Learning is used. Let us understand what is Machine Learning



edureka!

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

31°C Haze

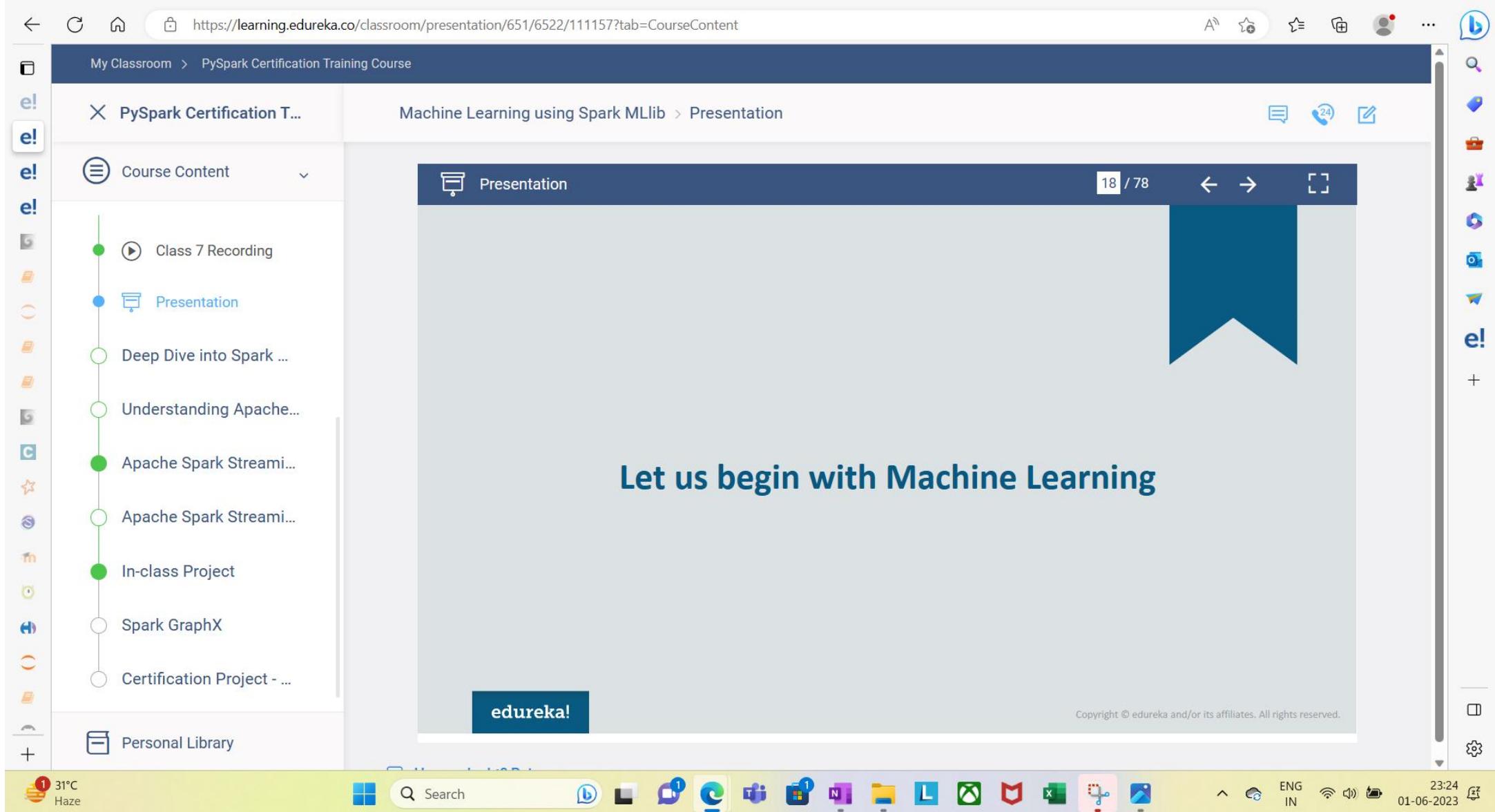
Search

L

ENG IN

23:24

01-06-2023



https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...
- Personal Library

19 / 78 ← →

Presentation

What is Machine Learning?

"Machine learning is a method of *data analysis* that *automates analytical model building*. Using algorithms that *iteratively learn from data*, machine learning allows computers to *find hidden insights without being explicitly programmed where to look*"

Customers	Last Purchase	Products in Cart	Recent Searches
Customer 1	Shoes	Socks, Sneakers	Wedges, Block Heels
Customer 2	Laptop	Laptop Charger, Laptop Bag	Mouse
Customer 3	Mobile Phone	Mobile Charger, Tempered Glass	Back Cover, HandHold
Customer 4	T-Shirt	Shorts, Jeans	Leggings, Kurti
Customer 5	Hand Bag	Sling Bag, Wallets	Clutches
Customer 6	SunGlasses	Cover, Specs	Camera, Handycam
Customer 7	Grocery	Glass, Bowl	Spoons, Cupset
Customer 8	Accessories	Ear-Ring, Rings	Nosepin,
Customer 9	Make-Up Kit	Eyeliner, Lipstic	Kajal, LipGloss
Customer 10	Wrist Watch	Mi-Band, Bracelet	Wall Clocks

Training Data



edureka!

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

31°C Haze

Search

Windows Start button

Icons for various Microsoft applications: Edge, File Explorer, OneDrive, Mail, Teams, OneNote, Powerpoint, Word, Excel, etc.

Network, Battery, Volume, and Date/Time (23:24, 01-06-2023) icons.

X PySpark Certification T...

Machine Learning using Spark MLlib > Presentation



Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Presentation

20 / 78



What is Machine Learning?

"Machine learning is a *method of data analysis* that *automates analytical model building*. Using algorithms that *iteratively learn from data*, machine learning allows computers to *find hidden insights without being explicitly programmed* where to look"

Customers	Last Purchase	Products in Cart	Recent Searches
Customer 1	Shoes	Socks, Sneakers	Wedges, Block Heels
Customer 2	Laptop	Laptop Charger, Laptop Bag	Mouse
Customer 3	Mobile Phone	Mobile Charger, Tempered Glass	Back Cover, HandHold
Customer 4	T-Shirt	Shorts, Jeans	Leggings, Kurti
Customer 5	Hand Bag	Sling Bag, Wallets	Clutches
Customer 6	SunGlasses	Cover, Specs	Camera, Handycam
Customer 7	Grocery	Glass, Bowl	Spoons, Cupset
Customer 8	Accessories	Ear-Ring, Rings	Nosepin,
Customer 9	Make-Up Kit	Eyeliner, Lipstic	Kajal, LipGloss
Customer 10	Wrist Watch	Mi-Band, Bracelet	Wall Clocks

Training Data



Customers	Last Purchase	Products in Cart
Customer 1	Shoes	Socks, Sneakers
Customer 2	Laptop	Laptop Charger, Laptop Bag
Customer 3	Mobile Phone	Mobile Charger, Tempered Glass
Customer 4	T-Shirt	Shorts, Jeans
Customer 5	Hand Bag	Sling Bag, Wallets
Customer 6	SunGlasses	Cover, Specs
Customer 7	Grocery	Glass, Bowl
Customer 8	Accessories	Ear-Ring, Rings

Testing Data

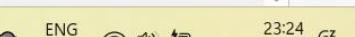
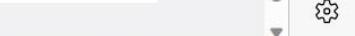
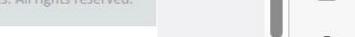
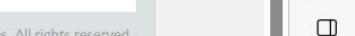
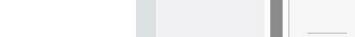
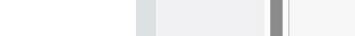
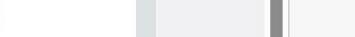
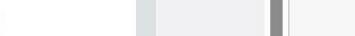
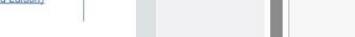
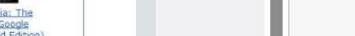
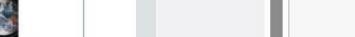
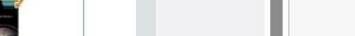
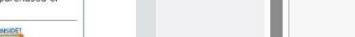
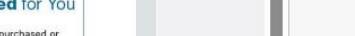
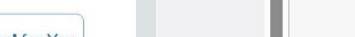
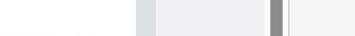
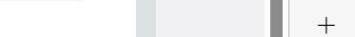
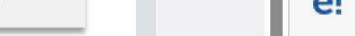
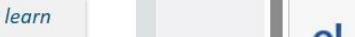
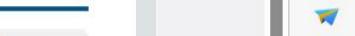
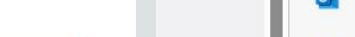
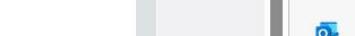
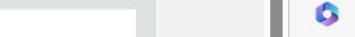
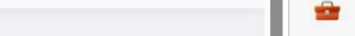
Accuracy

edureka!

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

X PySpark Certification T...

Machine Learning using Spark MLlib > Presentation



- Course Content
- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...
- Personal Library

What is Machine Learning?

"Machine learning is a *method of data analysis* that *automates analytical model building*. Using algorithms that *iteratively learn from data*, machine learning allows computers to *find hidden insights without being explicitly programmed where to look*"

Customers	Last Purchase	Products in Cart	Recent Searches
Customer 1	Shoes	Socks, Sneakers	Wedges, Block Heels
Customer 2	Laptop	Laptop Charger, Laptop Bag	Mouse
Customer 3	Mobile Phone	Mobile Charger, Tempered Glass	Back Cover, HandHold
Customer 4	T-Shirt	Shorts, Jeans	Leggings, Kurti
Customer 5	Hand Bag	String Bag, Wallets	Clutches
Customer 6	SunGlasses	Cover, Specs	Camera, Handycam
Customer 7	Grocery	Glass, Bowl	Spoons, Cupset
Customer 8	Accessories	Ear-Ring, Rings	Nosepin,
Customer 9	Make-Up Kit	EyeLiner, Lipstick	Kajal, LipGloss
Customer 10	Wrist Watch	Mi-Band, Bracelet	Wall Clocks

New Input



amazon.com

Recommended for You

Amazon.com has new recommendations for you based on items you purchased or told us you own.

LOOK INSIDE! Google Apps Desciphered: Compute in the Cloud to Streamline Your Desktop

LOOK INSIDE! Google Apps Administrator Guide: A Private Label Web Workspace

LOOK INSIDE! Google Apps: The Ultimate Google Resource (3rd Edition)

Predicted Output

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

https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

22 / 78 ← →

Features of Machine Learning

01 It uses the data to *detect patterns* in a dataset and *adjust program actions accordingly*

02 It *focuses on the development of computer programs* that can teach themselves to *grow and change* when *exposed to new data*

edureka!

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

31°C Haze

Search

L

ENG IN

23:24 01-06-2023

← ⌛ 🏠 🔒 https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent ⌛ ⌚ 🗑️ 🌐 🌐 🌐

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Machine Learning using Spark MLlib > Presentation

e! e!

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

23 / 78 ← → []

I wonder, is Machine Learning useful in Industry?

Let's have a look at trends

edureka!

The slide features a blue background with two cartoon characters. On the left, a man with a beard and glasses, wearing a blue shirt, has his hand to his chin in a thinking pose. A thought bubble above him contains the text 'I wonder, is Machine Learning useful in Industry?'. On the right, another man with a beard and glasses, wearing a dark suit, stands with his arms crossed, and a speech bubble next to him contains the text 'Let's have a look at trends'. The word 'edureka!' is written in large white letters at the bottom center of the slide. The top right corner of the slide shows '23 / 78' and navigation arrows. The top left corner has a presentation icon. The top right corner also has a message icon with '24', a phone icon, and a pen icon.

https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

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

Applications of Machine Learning

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

Machine Learning Applications



Siri

Apple claims that the software adapts to user's individual preferences overtime and personalizes results

Marketing and Sales

This ability to capture data, analyse it and use it to personalize a shopping experience (or implement a marketing campaign) is the future of retail

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Presentation 27 / 78 ← →

Machine Learning Applications



Healthcare

Machine Learning involves the use of wearable devices and sensors to assess a patient's health in real time



Financial Services

Financial Industry use machine Learning technology to identify insights in data and prevent fraud

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

Machine Learning Applications

edureka!

Biometrics

Biometrics use machine learning to identify individual credentials and prevent banking fraud

Transportation

Analyzing data to identify patterns and trends is key to the transportation industry, which relies on making routes more efficient and predicting potential problems to increase profitability

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

https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

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

31°C Haze

Search

L

ENG IN

23:24

01-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Machine Learning using Spark MLlib > Presentation



Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

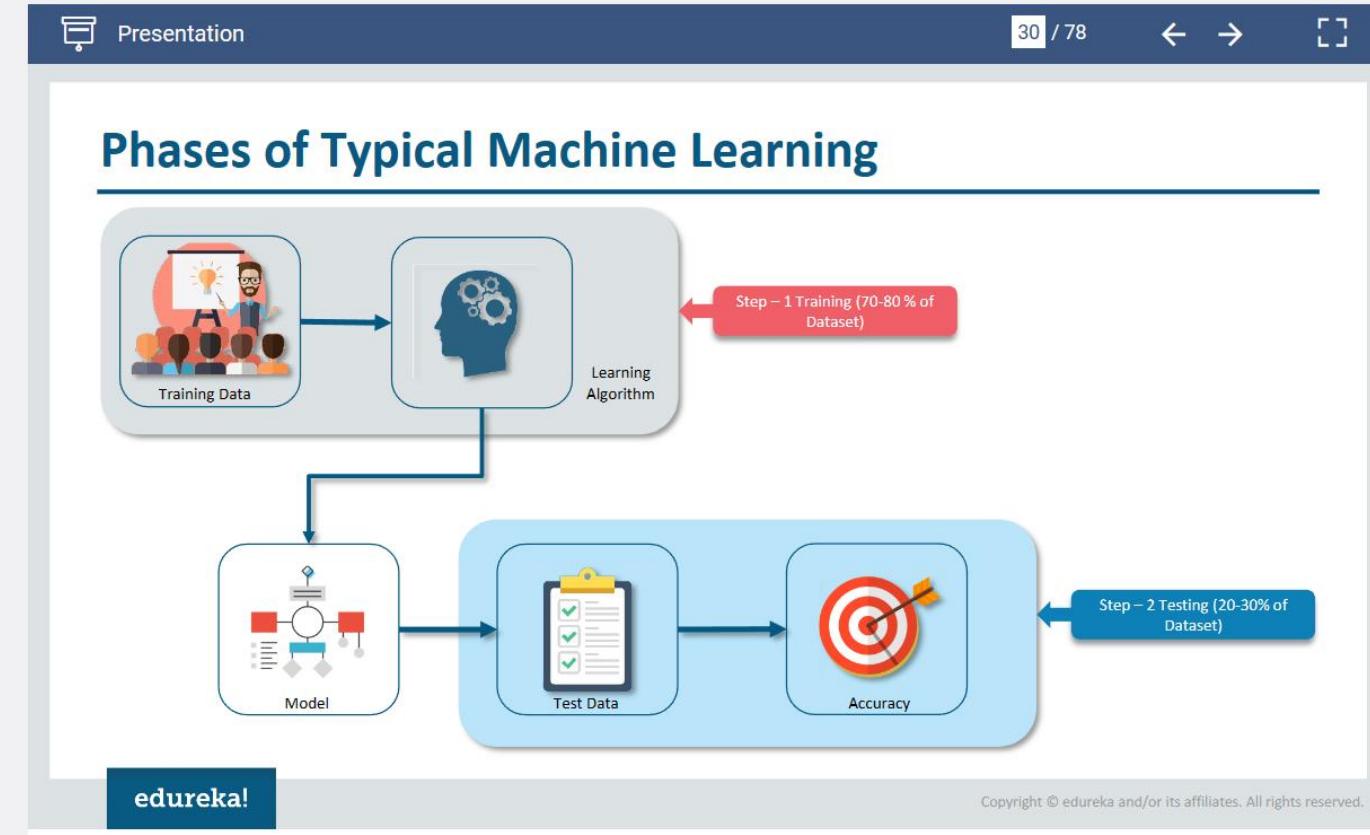
Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

e! Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Machine Learning using Spark MLlib > Presentation

Presentation

33 / 78

Collecting Data

This stage involves the collection of all relevant data from various sources

1 Collecting Data

2 Data Wrangling

3 Analyze Data

4 Train Algorithm

5 Test Algorithm

6 Deployment

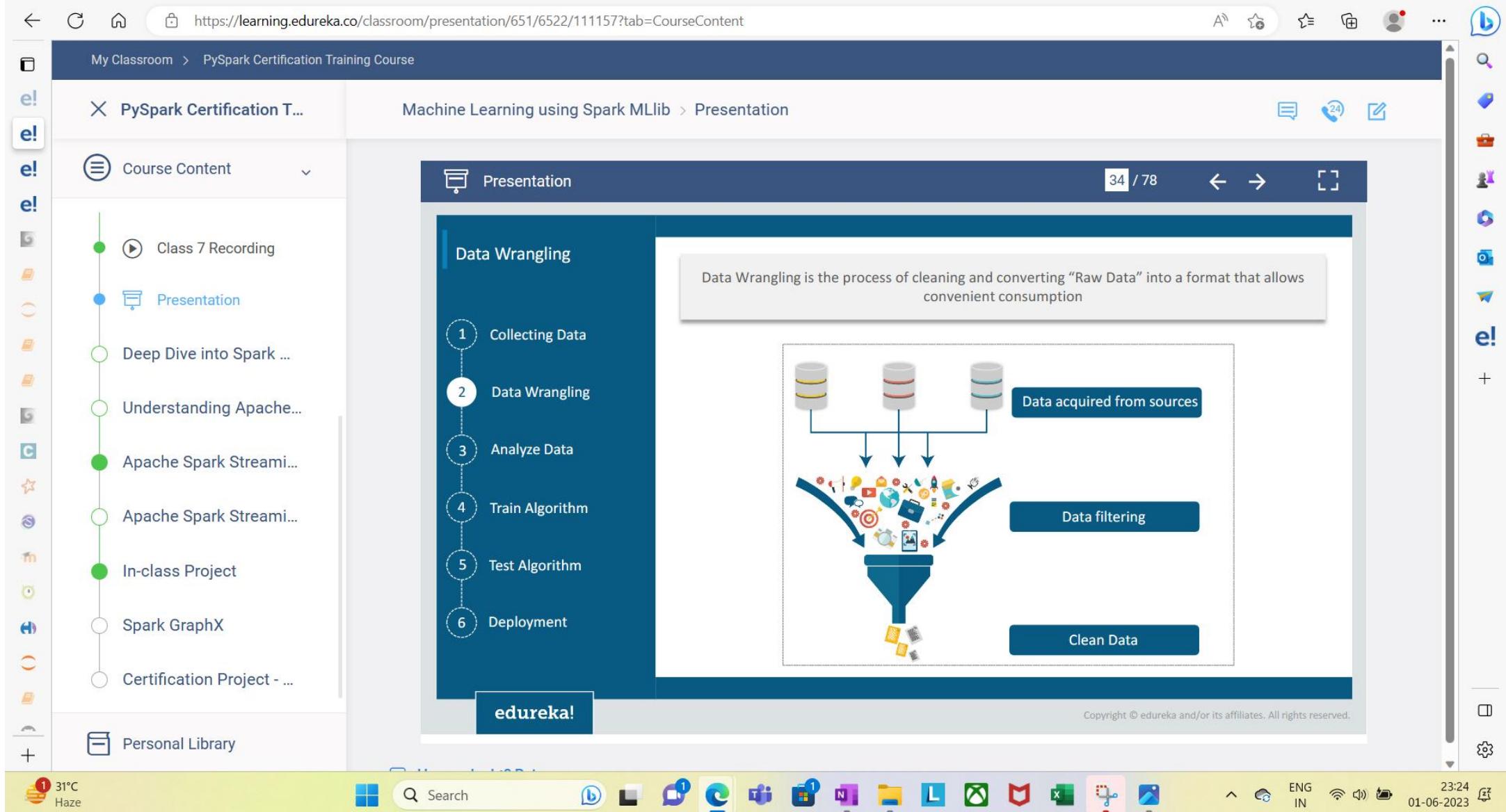
Data Sources

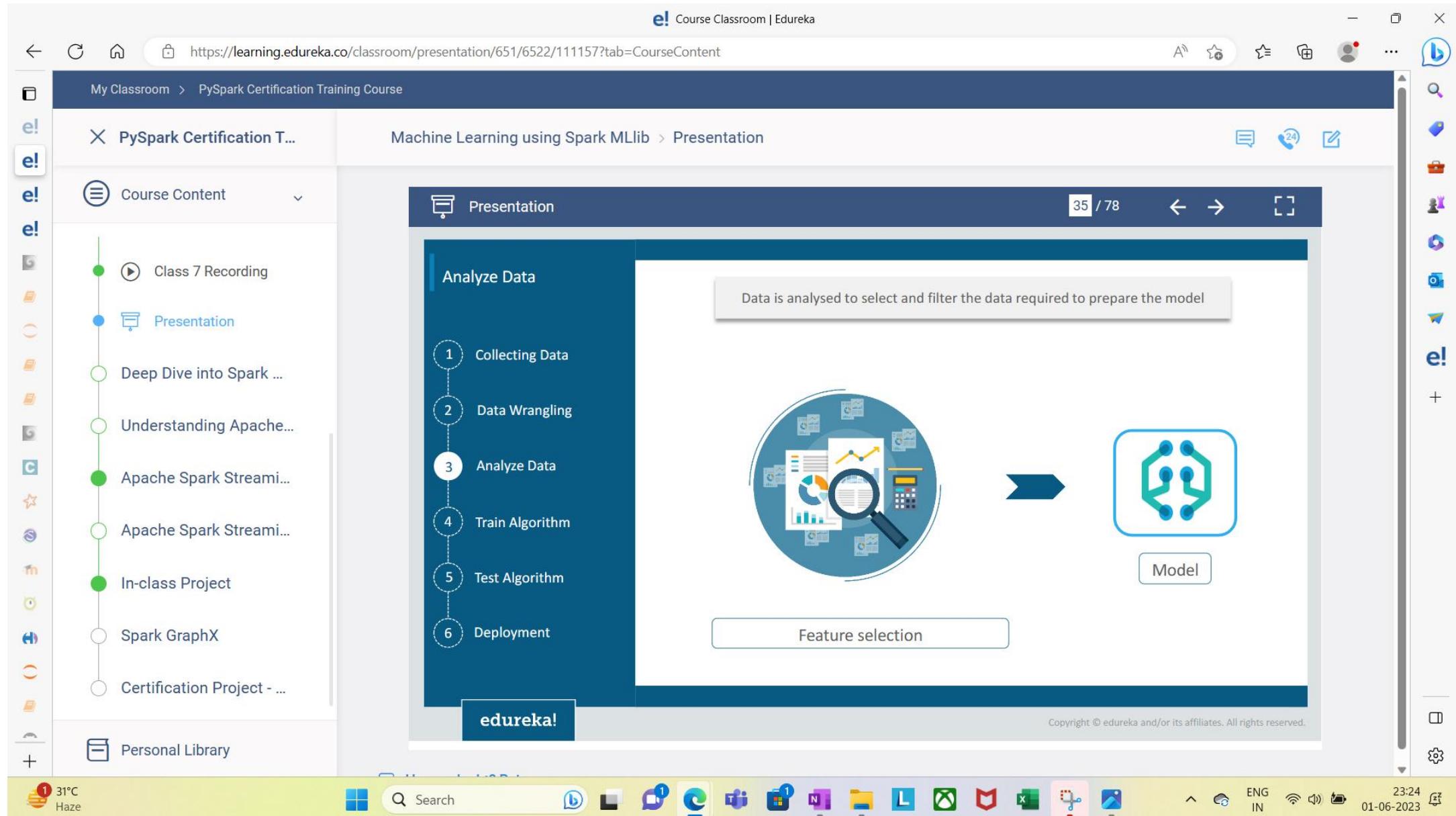
Data is collected from various sources in a server

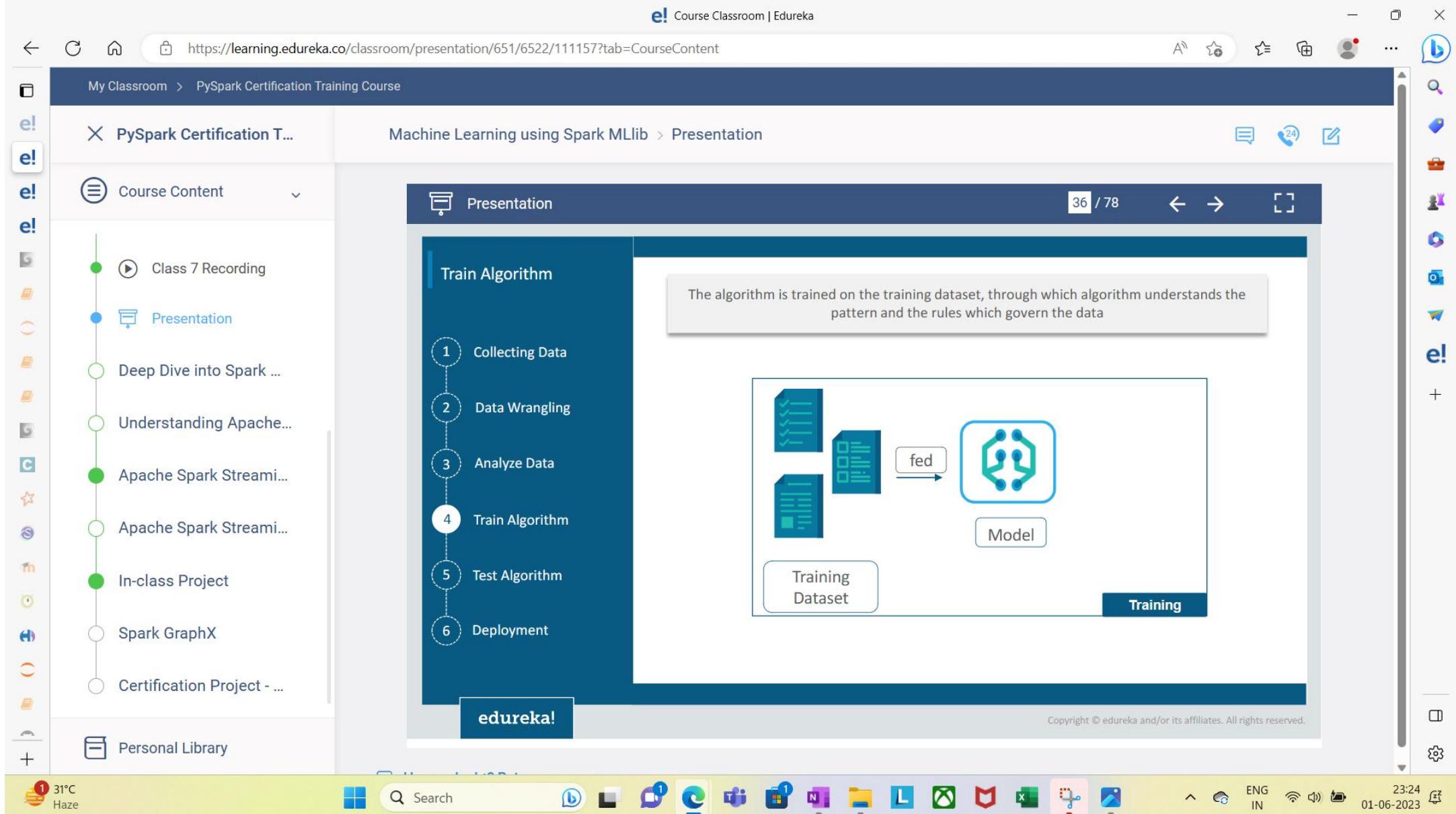
Server

edureka!

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







My Classroom > PySpark Certification Training Course

X PySpark Certification T... Machine Learning using Spark MLlib > Presentation

e! e!

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

Test Algorithm

- Collecting Data
- Data Wrangling
- Analyze Data
- Train Algorithm
- Test Algorithm
- Deployment

The testing dataset determines the accuracy of our model

edureka!

37 / 78 ← → []

Testing Dataset

fed

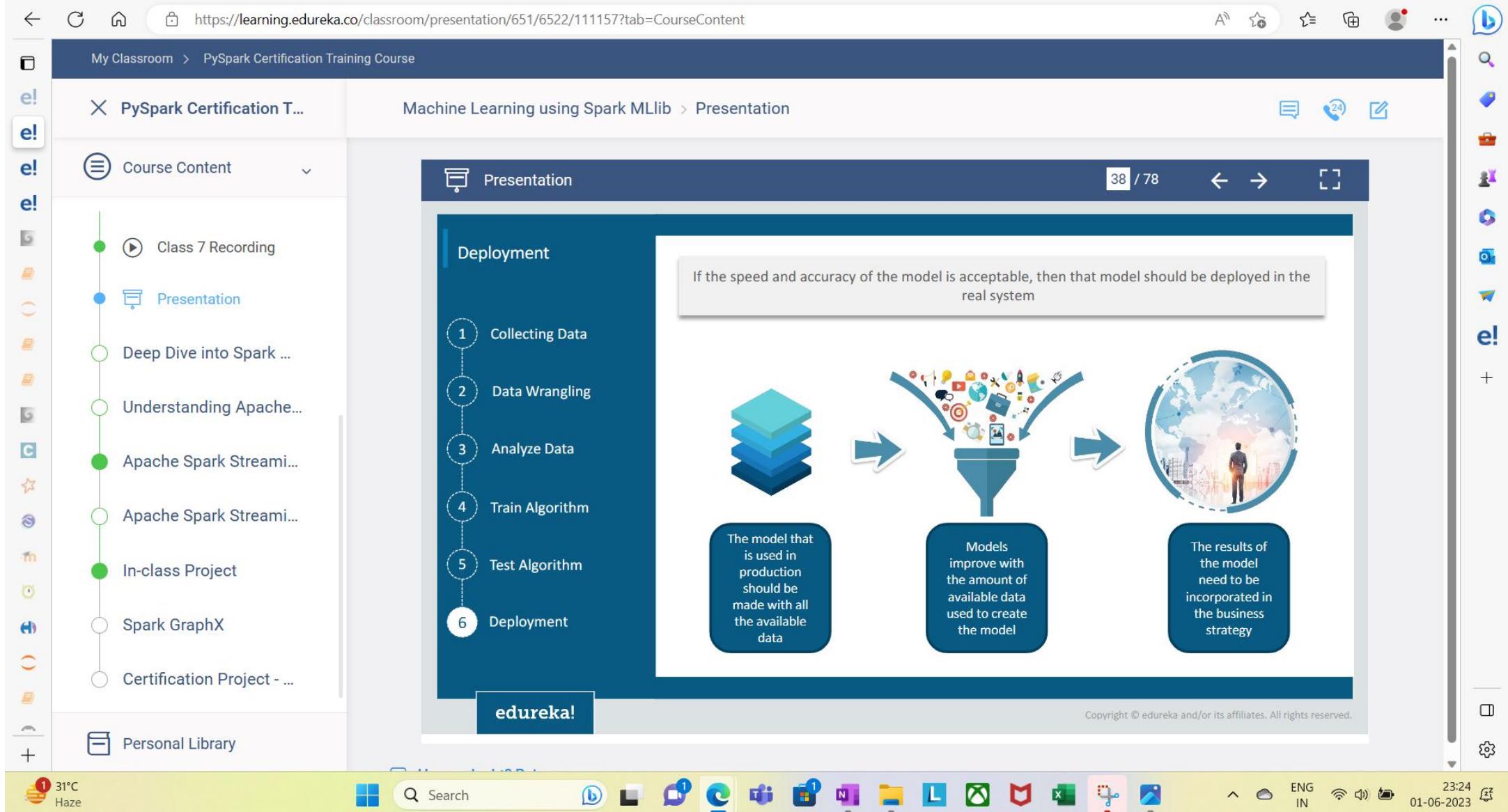
Model

Predicted Output

Accuracy??

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

The slide is titled 'Test Algorithm' and shows a vertical list of six steps: Collecting Data, Data Wrangling, Analyze Data, Train Algorithm, Test Algorithm, and Deployment. To the right, a diagram illustrates the 'Testing' phase. It shows three stacked rectangular boxes labeled 'Testing Dataset'. An arrow labeled 'fed' points from the datasets to a circular icon containing two stylized human figures, which is labeled 'Model'. An arrow points from the 'Model' to a box labeled 'Predicted Output'. A final arrow points from the 'Predicted Output' to a box labeled 'Accuracy??'. A callout box above the diagram states, 'The testing dataset determines the accuracy of our model'.



My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Presentation 39 / 78

Deployment

- 1 Collecting Data
- 2 Data Wrangling
- 3 Analyze Data
- 4 Train Algorithm
- 5 Test Algorithm
- 6 Deployment

After the model is deployed based upon its performance the model is updated and improved, if there is a dip in performance the model is retrained

The model that is used in production should be made with all the available data

Models improve with the amount of available data used to create the model

The results of the model need to be incorporated in the business strategy

edureka!

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



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

- ▶ Class 7 Recording
 - 💻 Presentation
 - Deep Dive into Spark ...
 - Understanding Apache...
 - Apache Spark Streami...
 - Apache Spark Streami...
 - In-class Project
 - Spark GraphX
 - Certification Project - ...

Machine Learning using Spark MLLib > Presentation



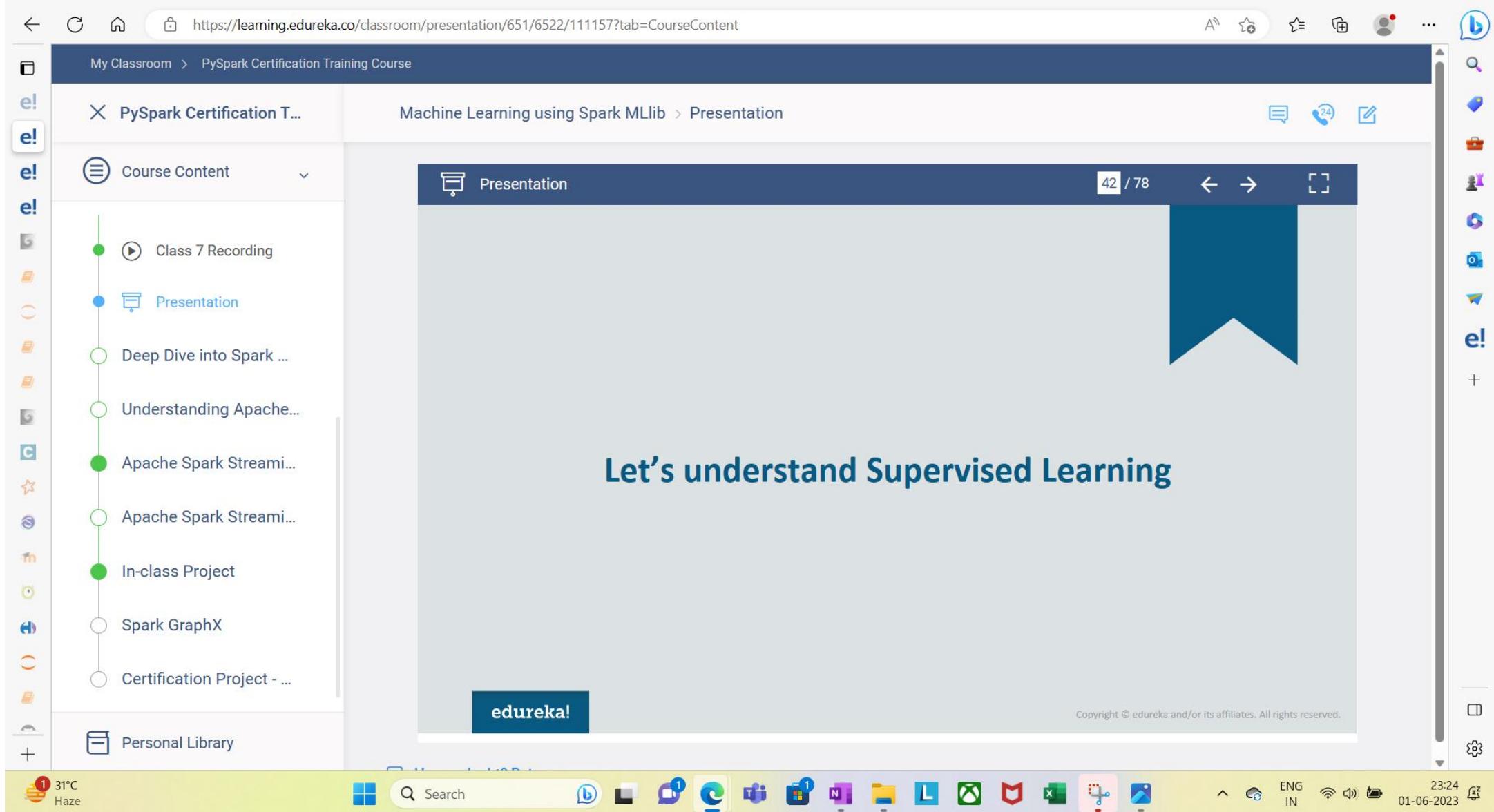
Let's have a Look at Different Types of Machine Learning

edureka



My Classroom > PySpark Certification Training Course







X PySpark Certification T...

Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

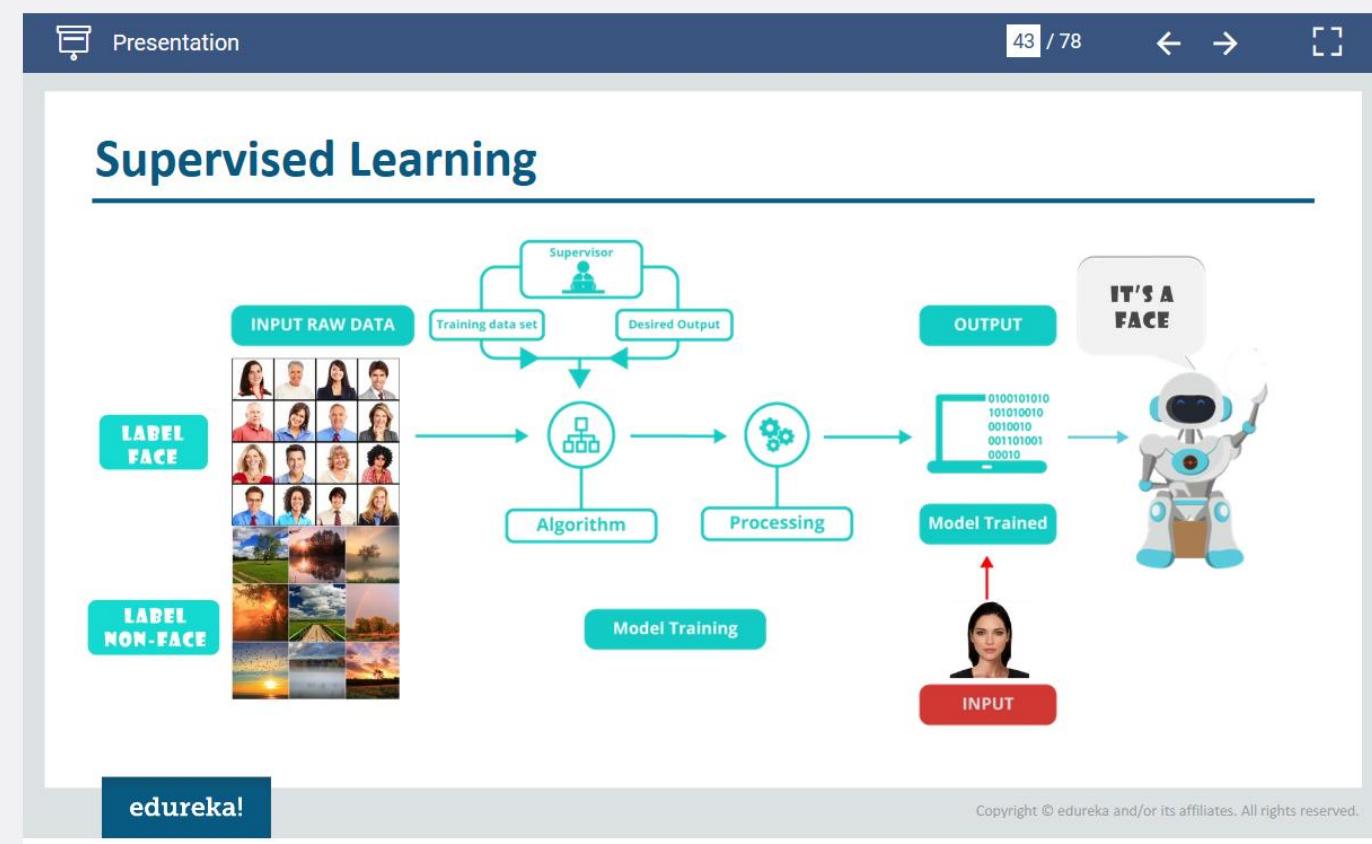
Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library





My Classroom > PySpark Certification Training Course

X PySpark Certification T...

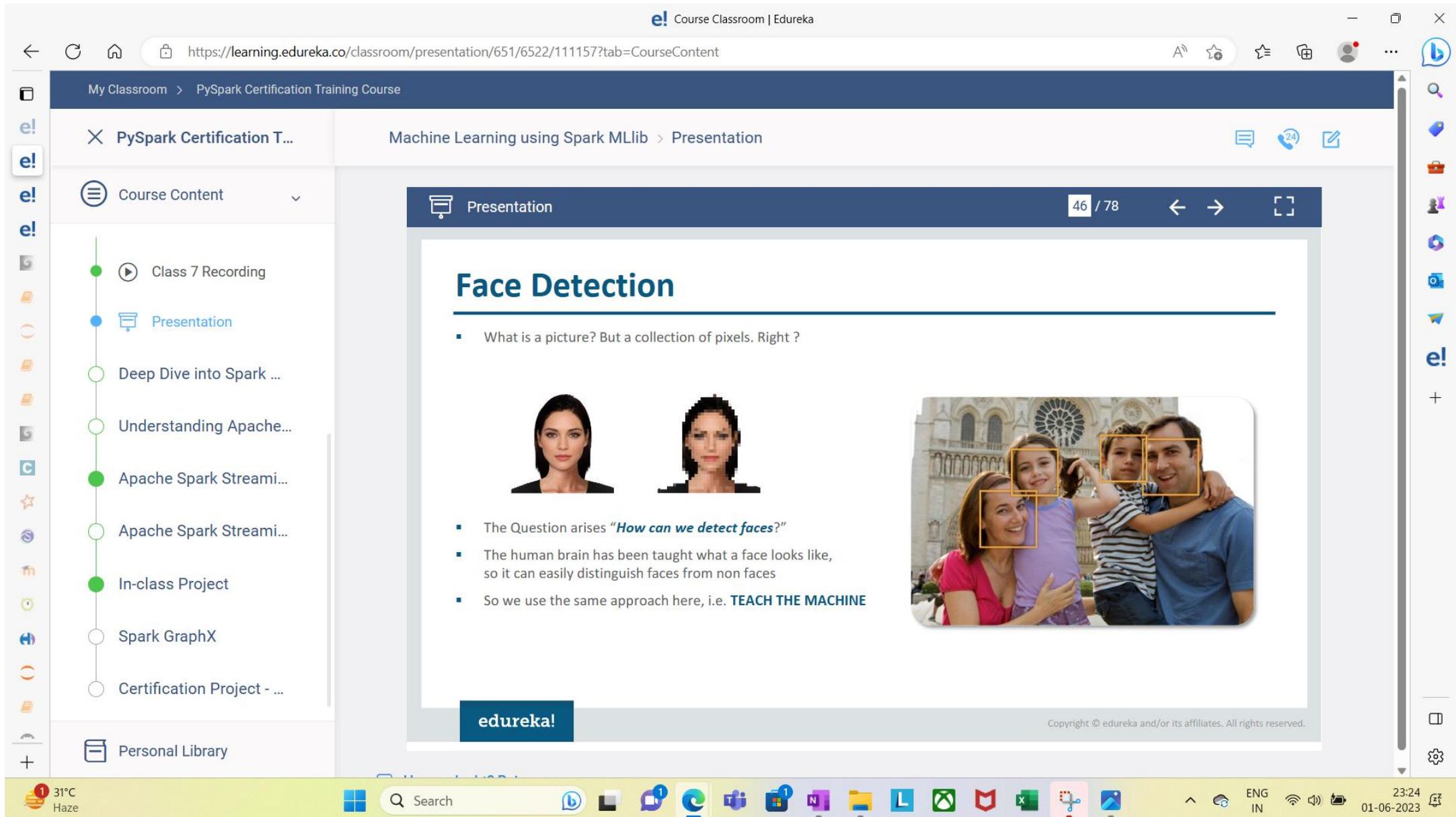
- ▶ Class 7 Recording
 - 💻 Presentation
 - Deep Dive into Spark ...
 - Understanding Apache...
 - Apache Spark Streami...
 - Apache Spark Streami...
 - In-class Project
 - Spark GraphX
 - Certification Project - ...

Machine Learning using Spark MLlib > Presentation



Let us understand Supervised Learning using a Use Case: FACE DETECTION

edureka





48 / 78



Providing Input to Machine With Labels

We have a *collection of photos (Faces and Non-Faces)*, out of that we take **70-80%** of images and provide as input with Labels to the machine



Label : "Face"



Label : "Non-Face"



The *machine scans the images and finds all the pixel features* of a picture that are particular to faces and *creates a MODEL*.

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

edureka!



My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Machine Learning using Spark MLlib > Presentation



Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Relevant and Irrelevant Features

- We have more than **160,000** features within a window to be calculated
- We don't need all the features, rather just the useful ones
- So using *Machine Learning* it makes a *collection of strong features* and are finally *combined to form a MODEL*



Relevant Feature

Irrelevant Feature

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



1

31°C

Haze



Search



23:24

01-06-2023

← ⌛ 🏠 🔒 https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent ⌛ 🏠 🔒 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Machine Learning using Spark MLlib > Presentation

e! e!

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

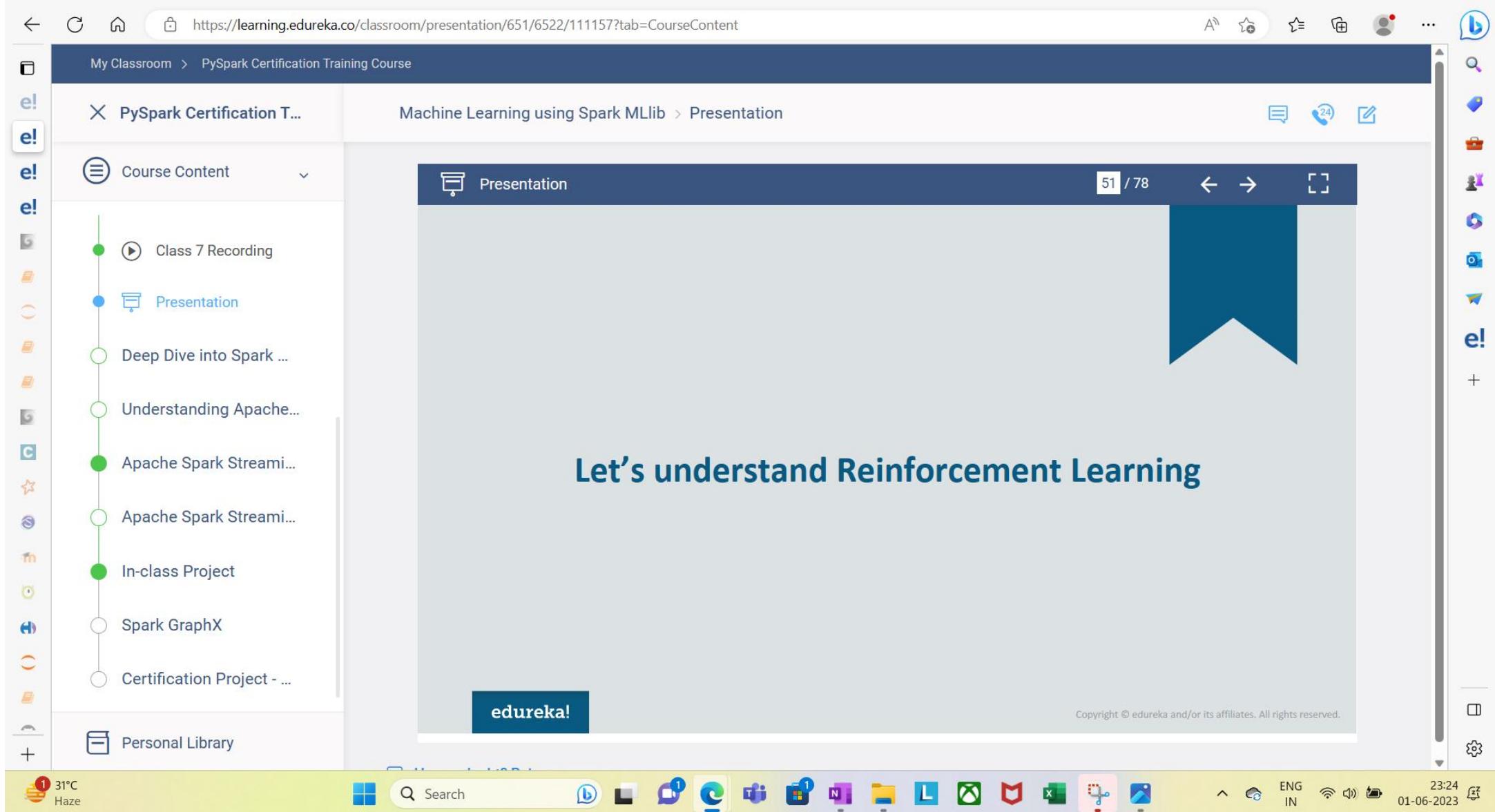
50 / 78 ← → []

Getting Results

- After the model is created, the *remaining 20-30%* of the *images* are *used for testing* purposes
- These *images* are *provided to the model* to get the *desired results*

edureka!

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



My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Presentation 55 / 78 ← →

Self Driving Cars at UBER : USE CASE

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

56 / 78 ← →

Self Driving Cars at UBER : USE CASE

ENVIRONMENT

Roads

Other Vehicles

People and Obstacle

edureka!

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

31°C Haze

Search

Cloud

ENG IN

23:24

01-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Machine Learning using Spark MLlib > Presentation



Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Presentation 57 / 78 ← → []

Self Driving Cars at UBER : USE CASE

ACTIONS and REWARD

edureka!

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

https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

59 / 78

Unsupervised Learning

- Unsupervised Learning is used against data that has no historical labels. The system is not told the "right answer." **The algorithm must figure out what is being shown**
- The goal is to explore the data and find some structure within. Unsupervised learning works well on transactional data
- Popular techniques include *self-organizing maps, nearest-neighbour mapping, K-MEANS CLUSTERING and singular value decomposition*
- These algorithms are also used to segment text topics, recommend items and identify data outliers



The Machine learns the patterns all by itself and divides them into categories. There are no labels provided in advance

edureka!

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

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Machine Learning using Spark MLlib > Presentation



e! Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

e! Personal Library

Presentation

61 / 78



Unsupervised Learning: USE CASE



So, how will I arrange them?

edureka!

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



31°C

Haze



Search



My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Presentation 63 / 78 ← →

Unsupervised Learning: USE CASE

So now you will take another physical character such as colour, and size

- RED COLOR AND BIG SIZE: **apple**
- YELLOW COLOR AND BIG SIZE: **bananas**
- GREEN COLOR AND SMALL SIZE: **grapes**

The task has been done

Here you did not learn anything before ,means no train data and no target variable

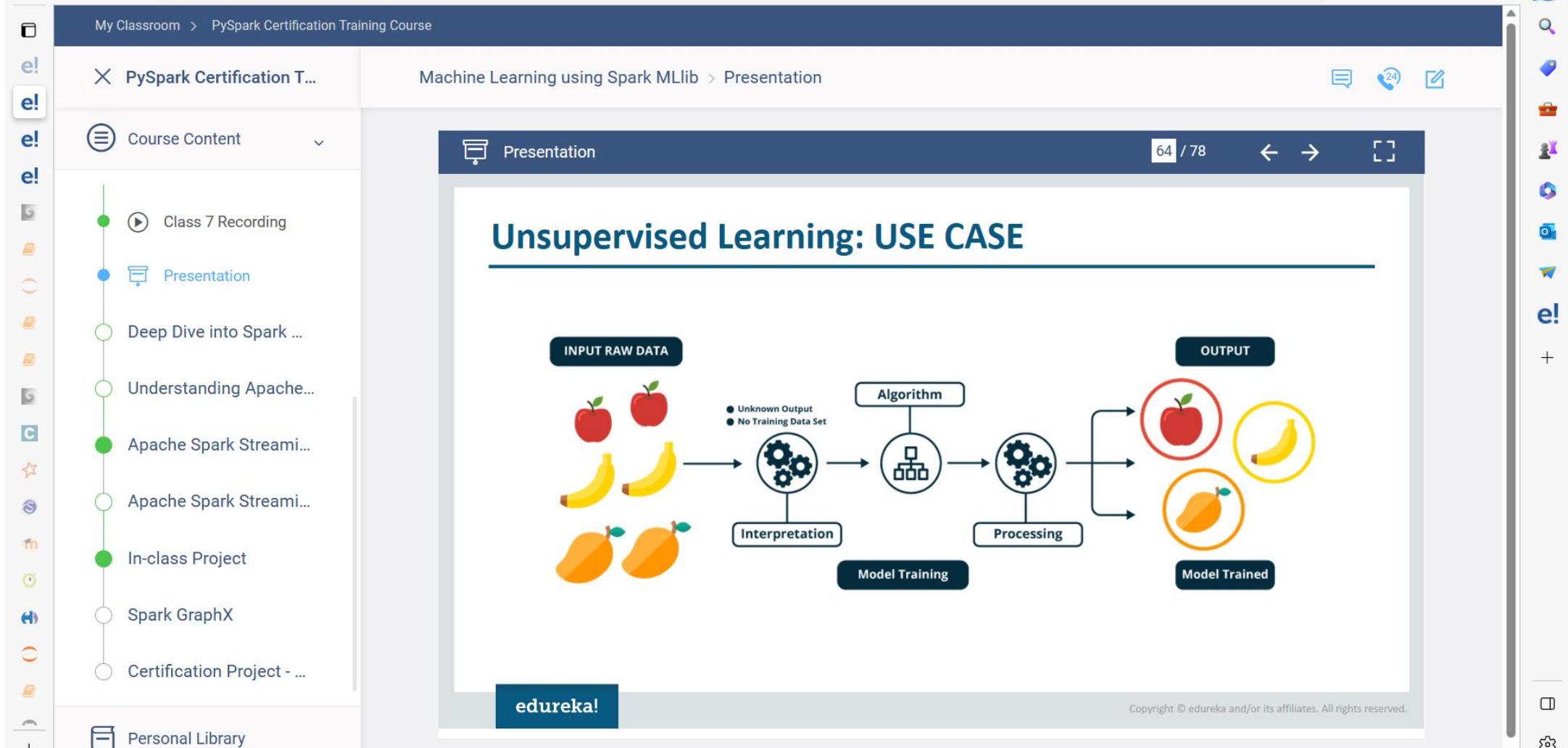
In Machine Learning, this kind of learning is known as **Unsupervised Learning**

The problem has following characteristics:

- Unlabelled learning data is available
- Output is dynamic to the input values, upon input of new values, output might change
- No predefined output classes. It can only be grouped into clusters based on the characteristics by the machine at runtime

edureka!

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



https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

Recap of the Topics

The diagram illustrates the Apache Hadoop ecosystem components and their integration with YARN and HDFS. Components shown include:

- 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)
- SOLR & LUCENE (Searching & Indexing)
- OOZIE (Scheduling)
- YARN (Resource Management)
- Storage (HDFS)
- Flume (Unstructured or Semi-structured Data)
- Sqoop (Structured Data)

A callout box on the right states: "Let us recall the concepts of Machine Learning that we have learnt, since they will be used in further module where we will be working on Spark MLlib".

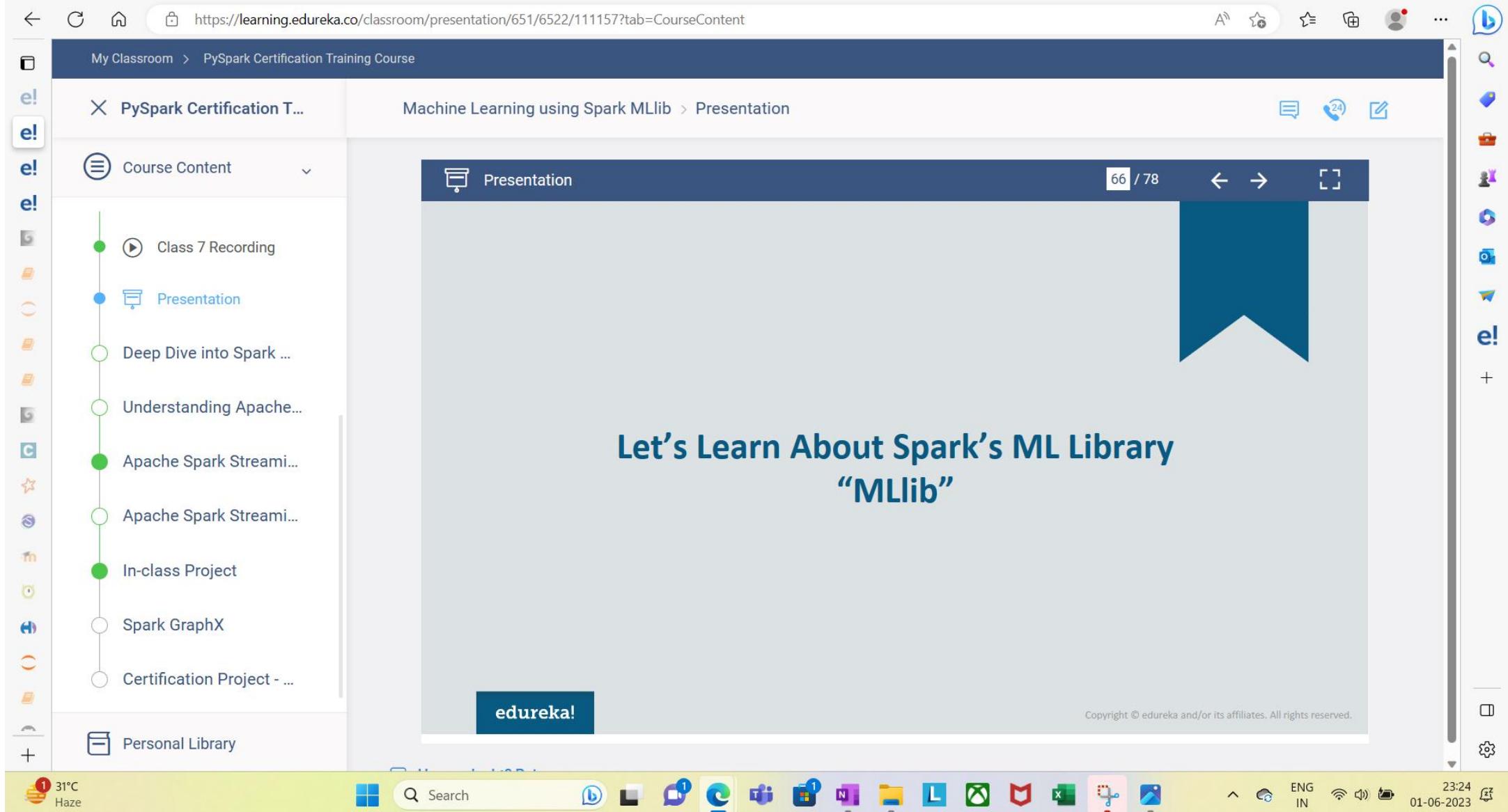
edureka!

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

31°C Haze

Search

Cloud, Mail, OneDrive, Teams, Edge, File Explorer, Task View, Taskbar icons, Volume, Network, Taskbar settings, Date/Time, Language, Battery



X PySpark Certification T...

Machine Learning using Spark MLlib > Presentation



Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

Personal Library

Presentation

67 / 78



MLlib

- MLlib is *Spark's machine learning (ML) library*
- It consists of *common learning algorithms and utilities*, including *classification, regression, clustering, collaborative filtering, dimensionality reduction*, as well as *lower-level optimization primitives* and *higher-level pipeline APIs*
- It divides into two packages:
 - **spark.mllib** contains the original API built on top of **RDDs**
 - **spark.ml** provides higher-level API built on top of **DataFrames** for constructing ML pipelines



edureka!

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



1

31°C

Haze



Search

ENG
IN

01-06-2023

23:24



My Classroom > PySpark Certification Training Course

PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Presentation 68 / 78 ← →

Why MLlib?

Simplicity and Compatibility: Simple APIs familiar to data scientists coming from tools like R and Python

Scalability: Ability to run the same ML code on your laptop and on a big cluster seamlessly without breaking down

Streamlined end-to-end: Building MLlib on top of Spark makes it possible to tackle distinct needs with a single tool instead of many disjointed ones

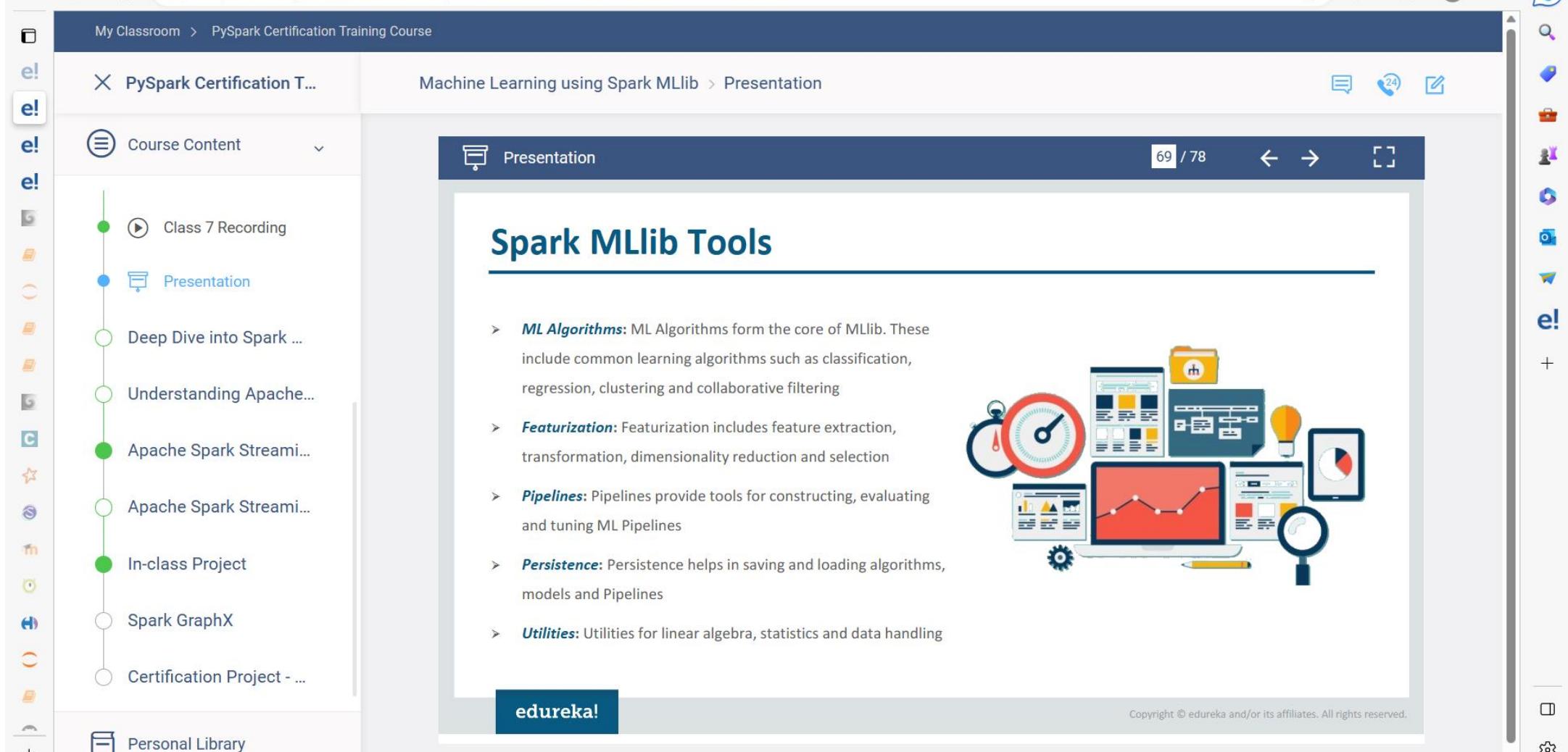
Faster Processing: Run Programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk

APACHE  **Spark**
MLlib

edureka!

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





X PySpark Certification T...

Machine Learning using Spark MLlib > Presentation



e! Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

Apache Spark Streami...

In-class Project

Spark GraphX

Certification Project - ...

e! Personal Library

Presentation

70 / 78



MLlib : Supported Algorithms

- Data types

- Basic statistics

- Summary Statistics
- Correlations
- Stratified Sampling
- Hypothesis Testing
- Random Data Generation

- Classification and regression

- Linear Models (SVMs, logistic regression, linear regression)
- Naive Bayes
- Decision Trees
- Ensembles of Trees (Random Forests and Gradient - Boosted Trees)

- Collaborative filtering

- Alternating Least Squares (ALS)

- Clustering

- k-Means
- Gaussian Mixture
- Power Iteration

- Dimensionality reduction

- Singular Value Decomposition (SVD)
- Principal Component Analysis (PCA)

- Feature extraction and transformation

- Optimization (developer)

- Stochastic Gradient Descent
- Limited-Memory BFGS (L-BFGS)

edureka!

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



31°C

Haze



Search



ENG

IN



01-06-2023



23:24

← ⌛ 🏠 🔒 https://learning.edureka.co/classroom/presentation/651/6522/111157?tab=CourseContent ⌛ 🏠 🔒 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Machine Learning using Spark MLlib > Presentation

e! e!

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

e! e!

Presentation

71 / 78 ← →

Spark ML Concepts

Concept	Explanation
DataFrame	Spark ML uses DataFrame from Spark SQL as an ML dataset, which can hold a variety of data types. E.g., a DataFrame could have different columns storing text, feature vectors, true labels, and predictions
Transformer	A Transformer is an algorithm which can transform one DataFrame into another DataFrame. E.g., an ML model is a Transformer which transforms DataFrame with features into a DataFrame with predictions
Estimator	An Estimator is an algorithm which can be fit on a DataFrame to produce a Transformer. E.g., a learning algorithm is an Estimator which trains on a DataFrame and produces a model
Pipeline	A Pipeline chains multiple Transformers and Estimators together to specify an ML workflow
Parameter	All Transformers and Estimators now share a common API for specifying parameters

edureka!

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



X PySpark Certification T...

e! Course Content

Class 7 Recording

Presentation

Deep Dive into Spark ...

Understanding Apache...

Apache Spark Streami...

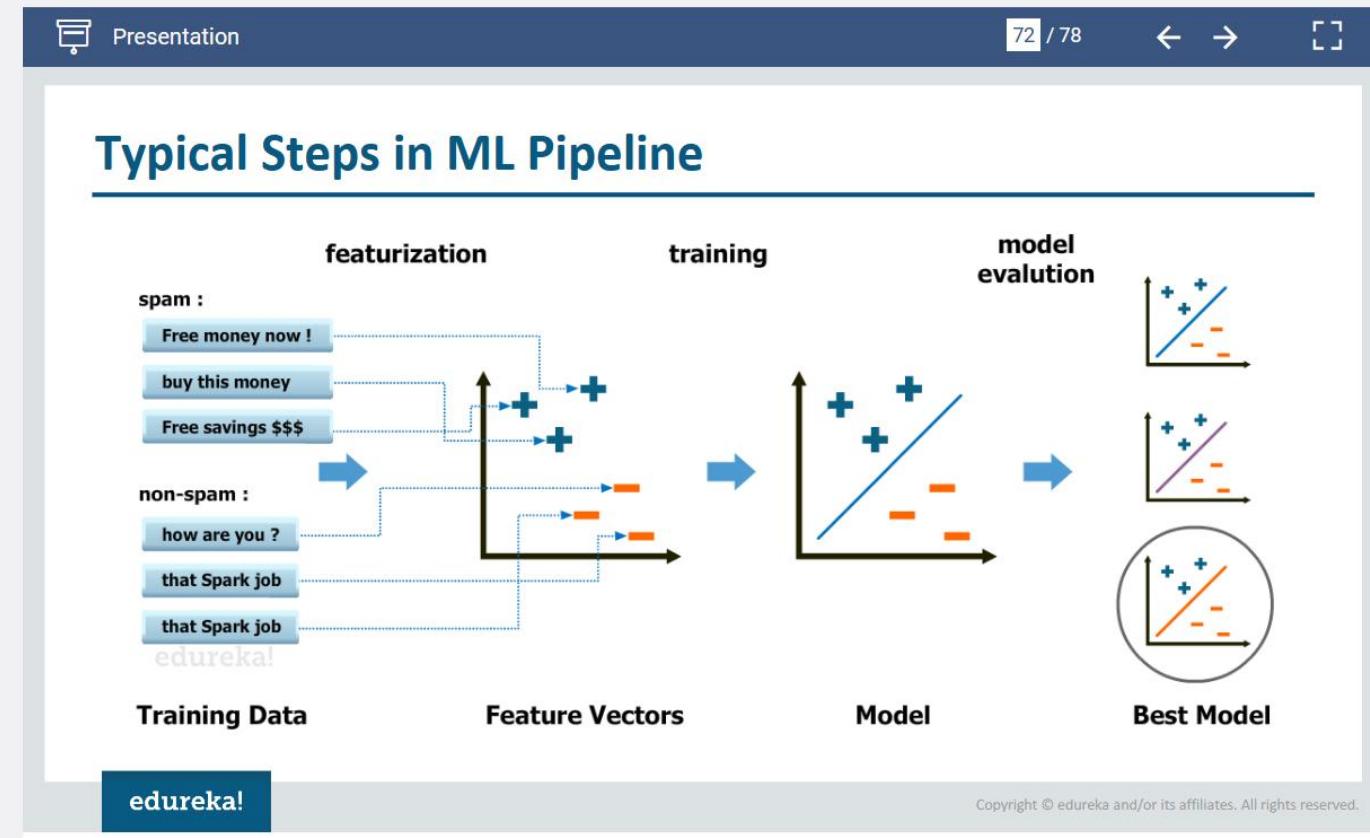
Apache Spark Streami...

In-class Project

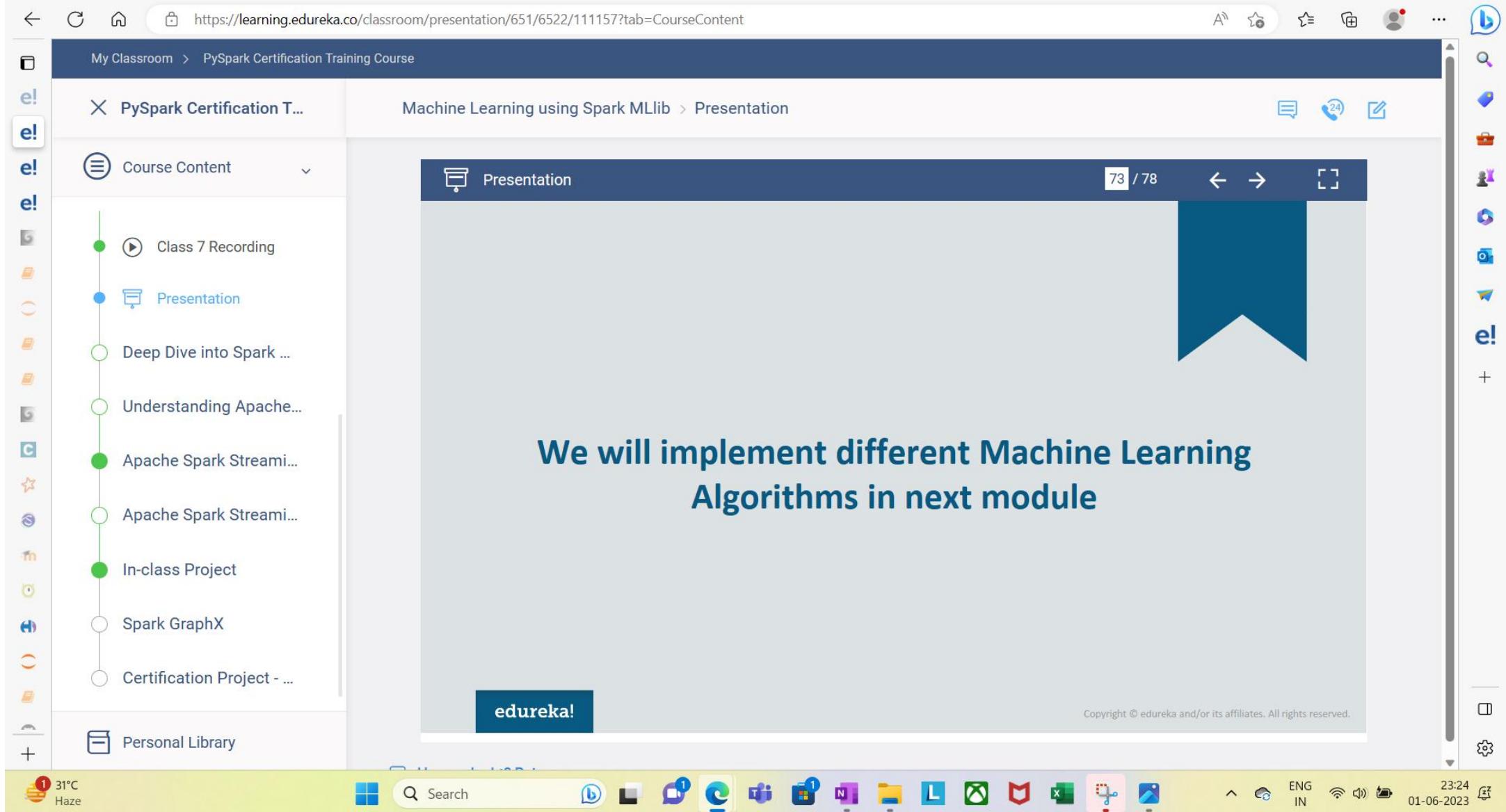
Spark GraphX

Certification Project - ...

Personal Library



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



My Classroom > PySpark Certification Training Course

X PySpark Certification T... Machine Learning using Spark MLlib > Presentation

Course Content

- Class 7 Recording
- Presentation
- Deep Dive into Spark ...
- Understanding Apache...
- Apache Spark Streami...
- Apache Spark Streami...
- In-class Project
- Spark GraphX
- Certification Project - ...

Personal Library

Presentation

Summary

What is Machine Learning?

Machine learning is a method of data analysis that uses computer algorithms to build models that learn from data, much like a human would. This allows machines to make predictions or decisions without being explicitly programmed to do so.

Machine Learning Applications

- SMS:** Predicts if the software adapts to user's individual preferences, our time and saves time results.
- Marketing and Sales:** Analyzes customer data and predicts consumer behavior to personalize shopping experiences or implement marketing campaigns for the future.

Phases of Typical Machine Learning

Steps of Machine Learning

- Collecting Data
- Data Wrangling
- Analyse Data
- Train Algorithm
- Test Algorithm
- Deployment

Types of Machine Learning

- Supervised Learning:** A person provides labeled training data for the algorithm to learn from.
- Unsupervised Learning:** An algorithm finds patterns in unlabeled data.
- Reinforcement Learning:** An algorithm learns by interacting with its environment through trial and error.

MLlib

- MLlib is a machine learning library.
- It contains various learning algorithms, including classification, regression, clustering, and matrix factorization, all running efficiently in parallel on large datasets.
- MLlib includes packages:
 - Apache MLlib:** Contains the original MLlib library.
 - Apache MLLib:** Contains improved MLlib with support for distributed learning.

edureka!

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

