

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count
- Personal Library

Presentation

Recap

Spark Components

Spark Architecture and SparkContext

Spark Shell : Python

```
Python 2.7.12 |Anaconda 4.2.0 (64-bit)| (default, Apr 20 2017, 00:16:05)
[GCC 4.2.1 Compatible Apple LLVM 8.0.0 (clang-800.0.42.1)]
Python 2.7.5 |Anaconda 2.3.2 (64-bit)| (default, Aug 10 2017, 17:05:23)
[GCC 4.2.1 Compatible Apple LLVM 7.3.0 (clang-703.0.29)]
Python 3.6.3 |Anaconda 2.3.2 (64-bit)| (default, Oct 24 2017, 17:55:13)
[GCC 4.2.1 Compatible Apple LLVM 8.0.0 (clang-800.0.42.1)]
Python 3.6.5 |Anaconda 2.3.2 (64-bit)| (default, Oct 24 2017, 17:55:13)
[GCC 4.2.1 Compatible Apple LLVM 8.0.0 (clang-800.0.42.1)]"/>
```

What is Sqoop?

- Apache Sqoop is a tool for efficiently transferring data between Apache Hadoop and relational databases.
- Sqoop exports data from a relational database into HDFS.
- Sqoop can also bring the export data from HDFS to a relational database.

How Sqoop Import & Export Works?

Sqoop - Codegen

```
Codegen generates Data Access Objects (DAO) classes that can easily "get" and "set" persistent data.
  1. Codegen generates the DAO classes.
  2. A generated DAO class can write the table definition.
  3. Generated code can be used.
  4. Generated code can be used to generate code.
```

```
Codegen generates Data Access Objects (DAO) classes that can easily "get" and "set" persistent data.
  1. Codegen generates the DAO classes.
  2. A generated DAO class can write the table definition.
  3. Generated code can be used.
  4. Generated code can be used to generate code.
```

edureka!

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

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

Topics

- Challenges in the Existing Computing Methods
- RDD
 - RDD Overview
 - RDD Workflow
 - RDD transformations and actions
 - Shuffling
 - Data Loading and Saving through RDD
 - Key-Value pair in RDD
- RDD Functions
- RDD Persistance
- RDD Lineage



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

edureka!

X PySpark Certification T...

Playing with Spark RDDs > Presentation



Course Content

- Presentation
- Case Study I - Performance Tuning
- Case Study II - Loading Data
- In-Class Demo - Input/Output
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

Personal Library

Presentation

6 / 100



Objectives

After completing this module, you should be able to:

- Analyze the Challenges in Existing Computing Methods
- Describe Probable Solution and How RDD Solves the Problem
- Describe What is RDD, Its Functions, Transformations and Actions
- Describe Data Loading and Saving Through RDDs
- Describe Key-Value Pair RDDs, Other Pair RDDs and Two-pair RDDs
- Describe RDD Lineage and RDD Persistence
- Explain a Spark Program Using RDD Concepts
- Describe RDD Partitioning and How It Helps Achieve Parallelization



edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 7 / 100

Challenges in Existing Computing Methods

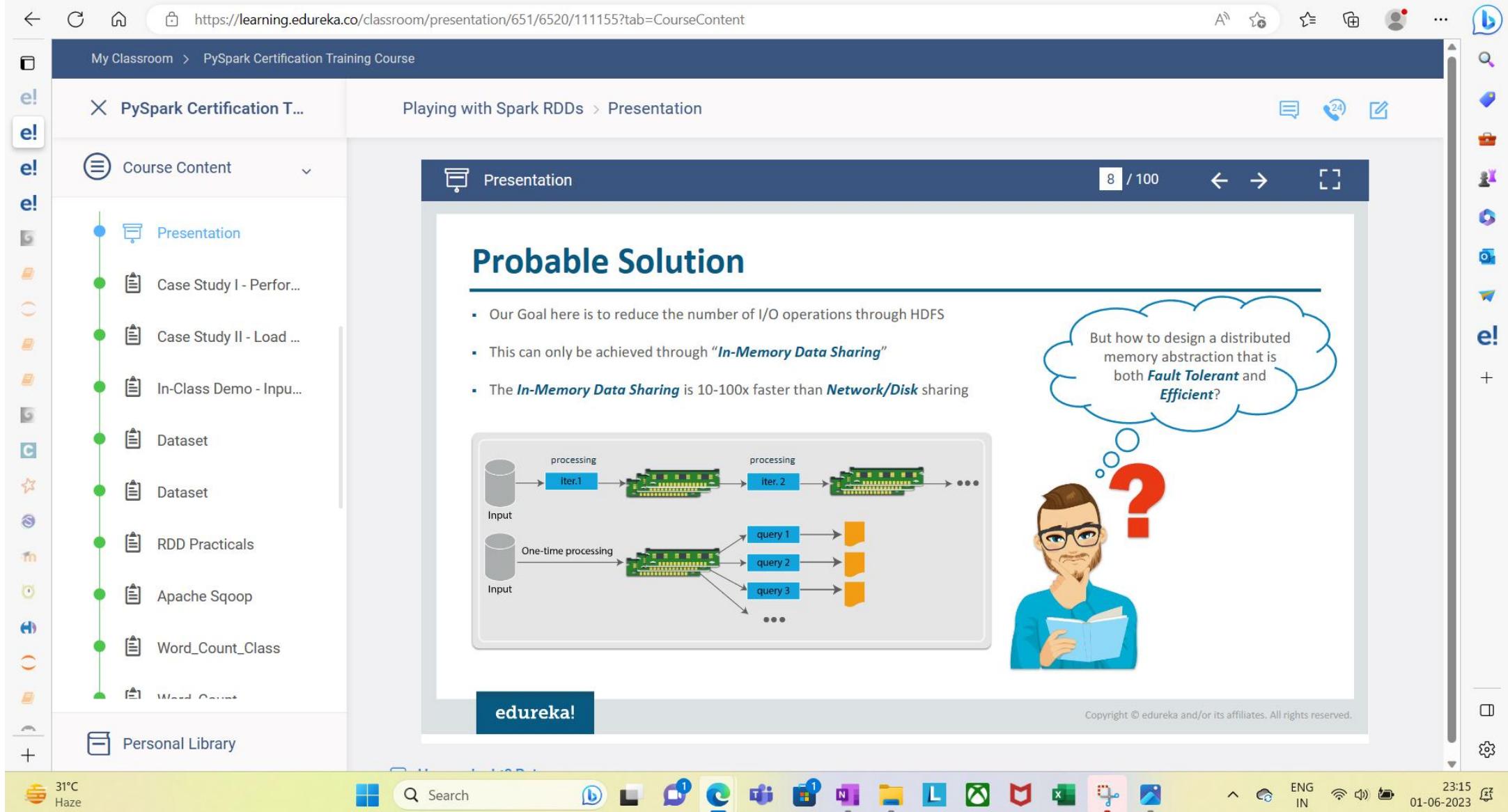
When it comes to iterative distributed computing, i.e. processing data over multiple jobs in computations, we need to reuse or share data among multiple jobs.

There are problems with the current computing method (**MapReduce**).

- We need to store data in some intermediate stable distributed storage such as HDFS
- Multiple I/O operations makes the overall computations of **jobs slower**
- Replications and serializations which in turn makes the **process slower**

edureka!

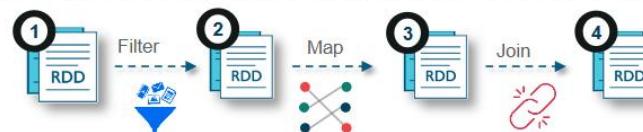
Copyright © edureka and/or its affiliates. All rights reserved.



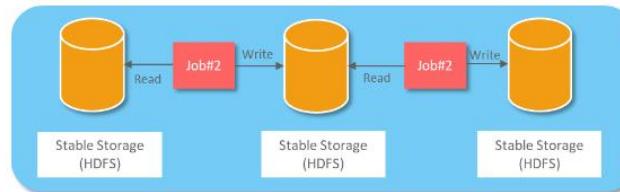
RDDs : The Perfect Solution

- RDDs try to solve all the problems by enabling ***fault tolerant distributed In-memory computations***
- Since RDDs are created over a set of transformations, it logs those transformations, rather than actual data

For example:

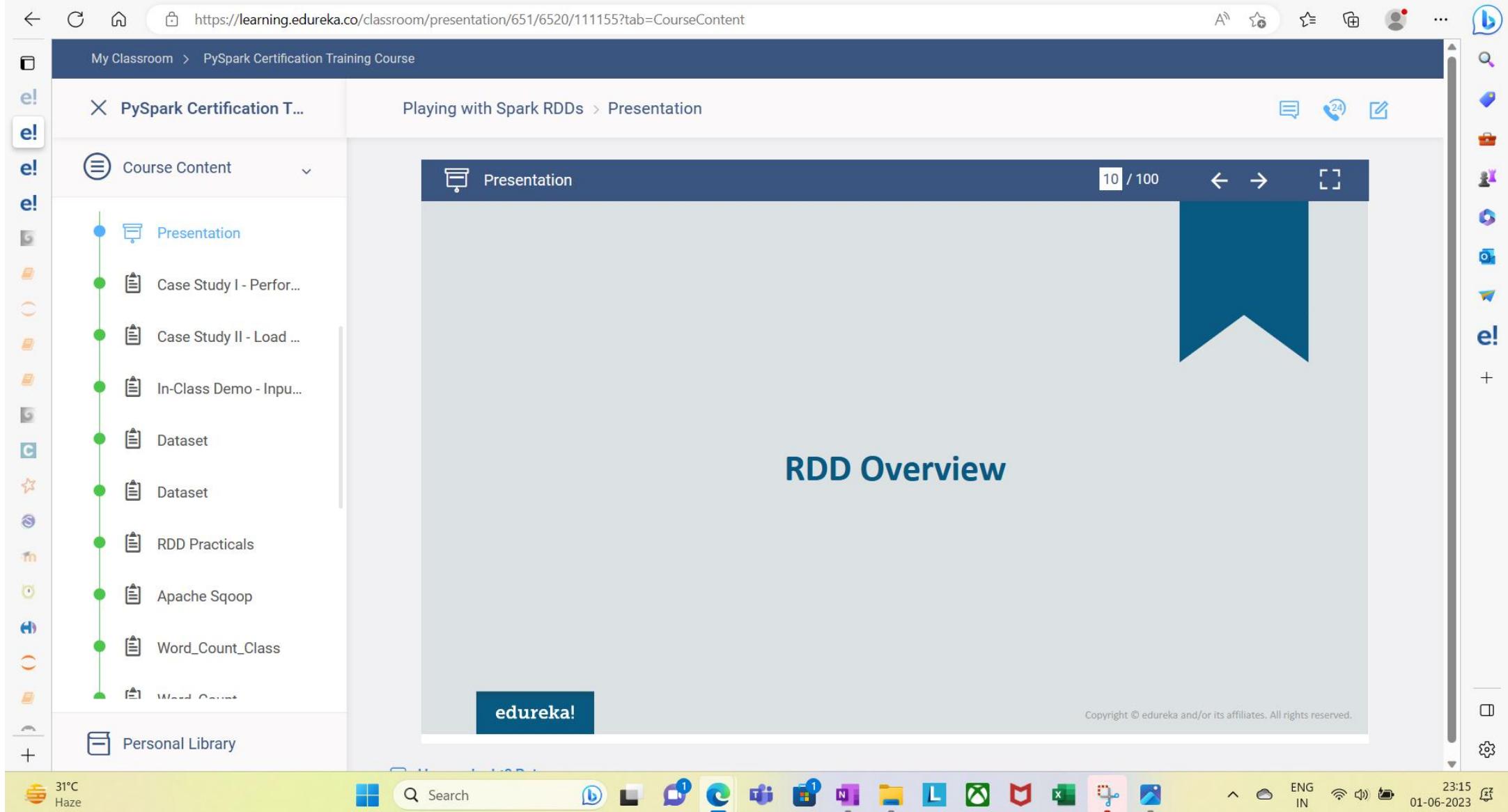


- *In case if we lose some partition of RDD, we can replay the transformation on that partition in LINEAGE to achieve the same computation, rather than doing data replication across multiple nodes*
- This characteristic is biggest benefit of RDD, because it *saves a lot of efforts in data management and replication* and thus achieves *faster computations*



edureka!

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



← ⌛ 🏠 🔍 https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Playing with Spark RDDs > Presentation

e! e!

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Inpu...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

What is RDD?

RDD = Resilient Distributed Datasets

RDD is a distributed memory abstraction which lets programmers perform *in-memory computations* on large clusters in a *fault-tolerant manner*

They are *read-only collection* of objects *partitioned across* a set of machines that *can be rebuilt* if a partition is lost

RDDs can be *created from multiple data sources* e.g. Python collection, local file system, Hadoop, Amazon S3, HBase table etc.

There are several operations performed on RDDs:

- Functions
- Transformations
- Actions

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 12 / 100 ← →

Features of RDD's

- In-Memory Computation: RDDs have a provision of in-memory computation.
- Lazy Evaluations: All transformations are lazy, it doesn't compute their results right away.
- Fault Tolerance: RDDs track data lineage information to rebuild lost data automatically.
- Immutability: Data can be created or retrieved anytime.
- Partitioning: Partitioning is the fundamental unit of parallelism in Spark RDD.
- Persistence: Users can reuse RDDs and choose a storage strategy for them.
- Coarse Grained Operations: Applies to all elements in datasets through maps or filter or group by operation.
- In-Memory Computation: RDDs have a provision of in-memory computation.

edureka!

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

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

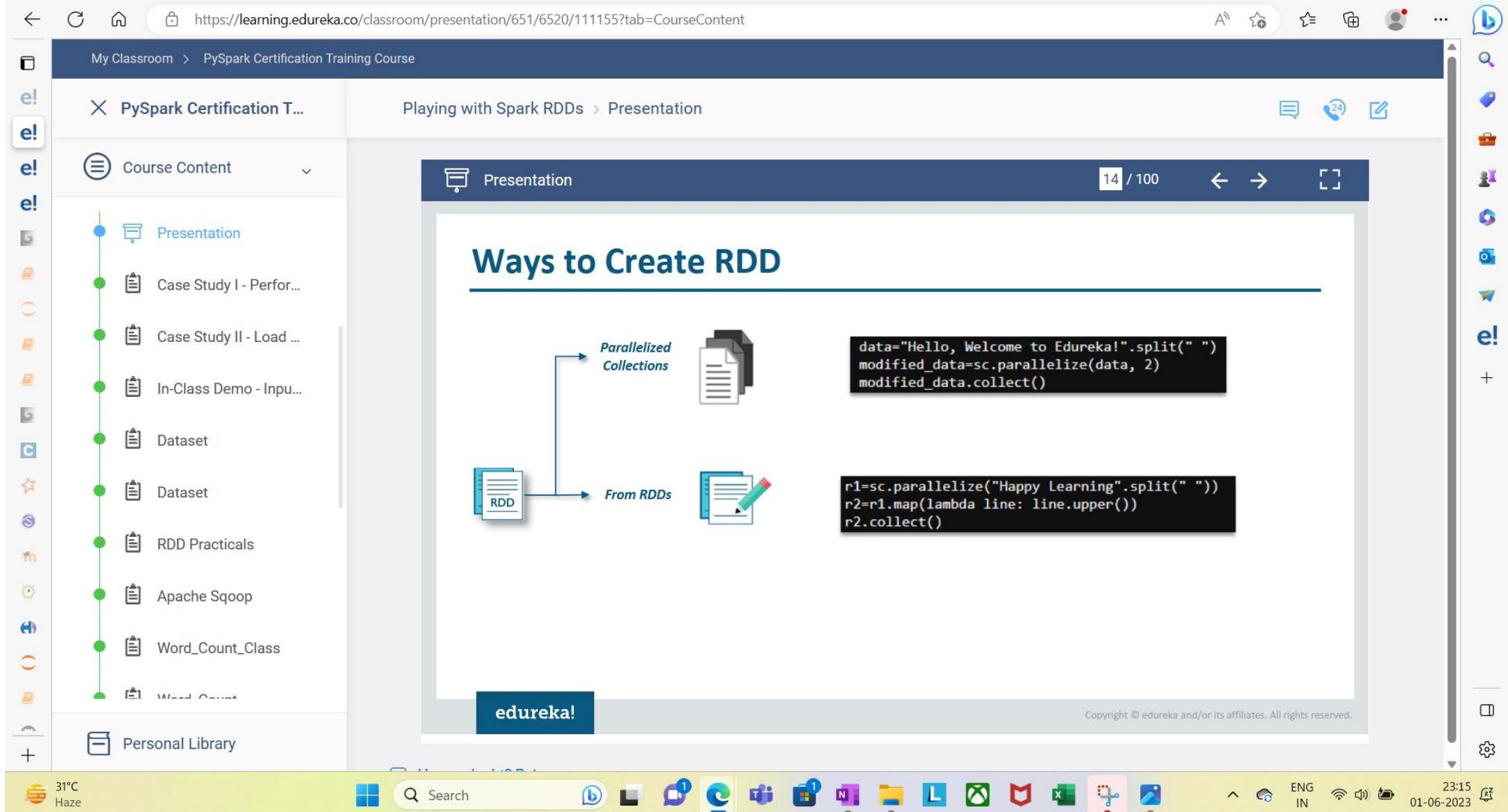
edureka!

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

Ways to Create RDD

Parallelized Collections

```
data="Hello, Welcome to Edureka!".split(" ")
modified_data=sc.parallelize(data, 2)
modified_data.collect()
```



My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 15 / 100

Ways to Create RDD

The diagram illustrates three methods for creating RDDs:

- Parallelized Collections:** Represented by a document icon. Example code:

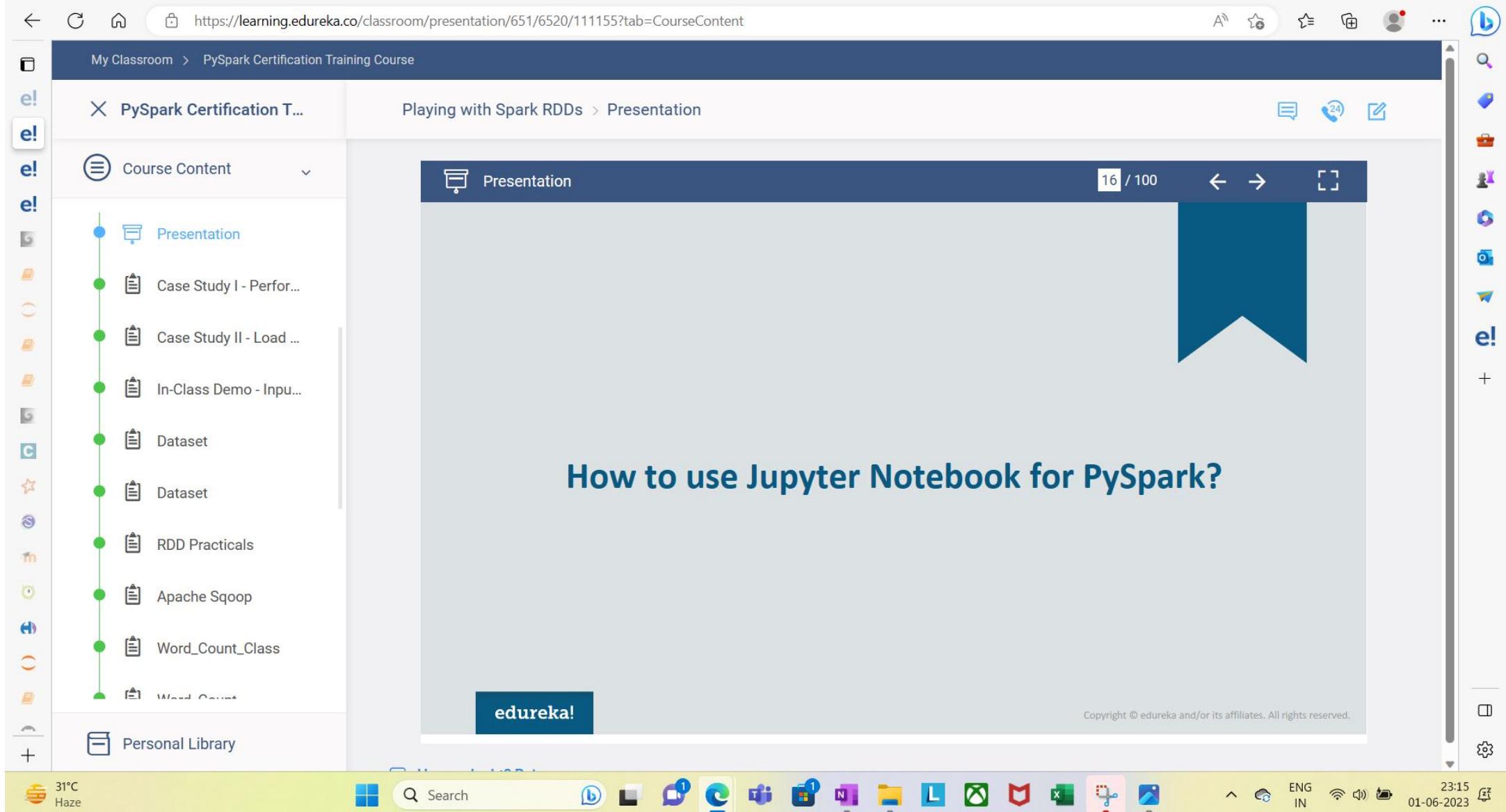
```
data="Hello, Welcome to Edureka!".split(" ") modified_data=sc.parallelize(data, 2) modified_data.collect()
```
- From RDDs:** Represented by a document with a pencil icon. Example code:

```
r1=sc.parallelize("Happy Learning".split(" ")) r2=r1.map(lambda line: line.upper()) r2.collect()
```
- External Data:** Represented by a cylinder icon. Example code:

```
text_file=sc.textFile("/user/edureka_294428/susmita_demo/demo.txt") text_file.map(lambda line: line.upper()).collect()
```

edureka!

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



Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

17 / 100



How to use Jupyter Notebook for PySpark?

```
In [1]: import pyspark  
In [2]: from pyspark import SparkContext  
sc = SparkContext()  
In [3]: data="Hello, Welcome to Edureka".split(" ")  
In [5]: modified_data=sc.parallelize(data,2)  
In [6]: modified_data.collect()  
Out[6]: ['Hello,', 'Welcome', 'to', 'Edureka']
```

You can also perform all the
RDD operations on Webconsole
as well as on Jupyter Notebook

edureka!

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

https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

edureka!

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

31°C Haze

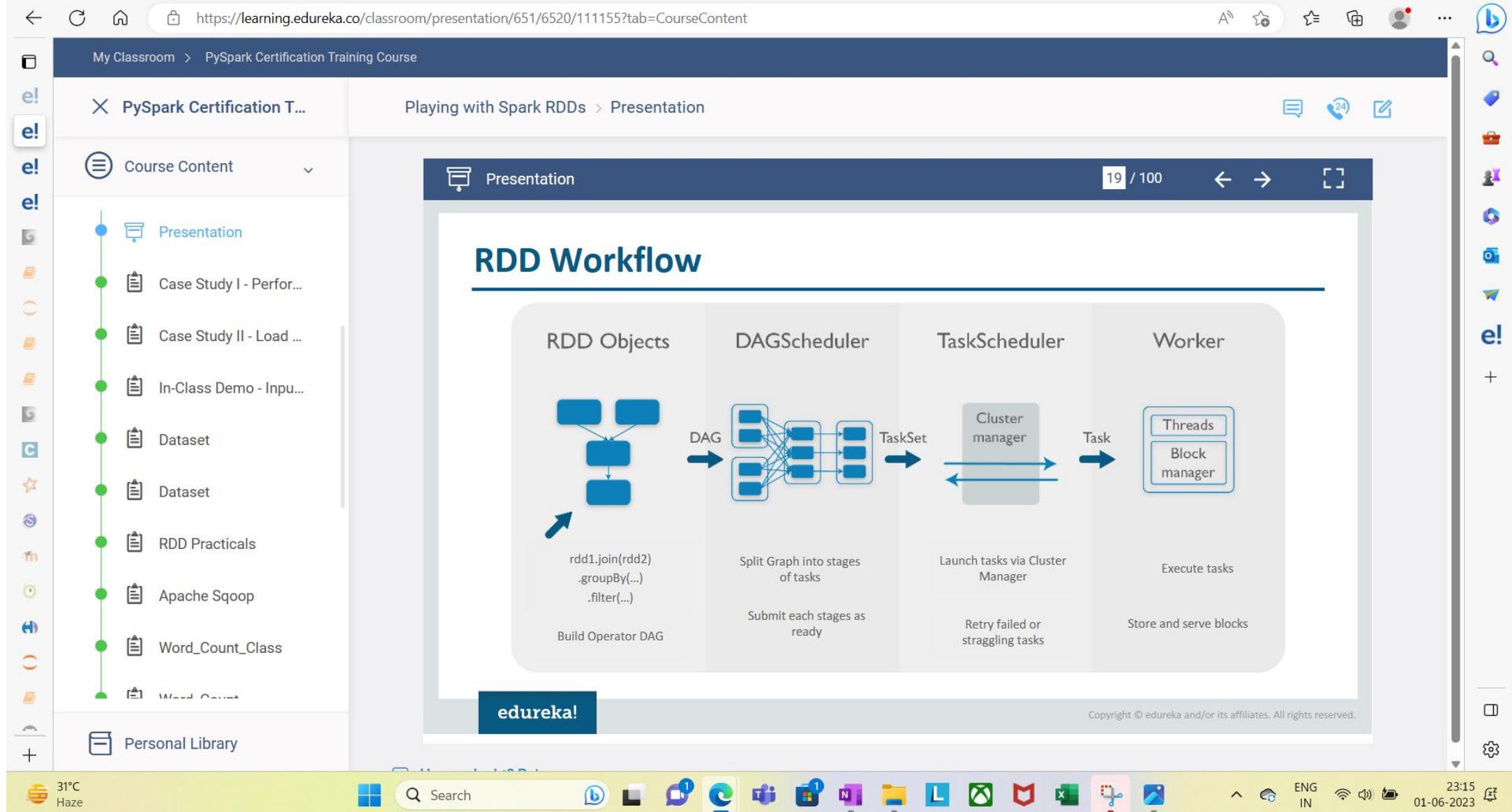
Search

L

ENG IN

01-06-2023

23:15



https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

Presentation

20 / 100

RDD Transformations and Actions

edureka!

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

31°C Haze

Search

L

ENG IN

01-06-2023

23:15

X PySpark Certification T...

Playing with Spark RDDs > Presentation



Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

Personal Library

Transformations

- Transformations *create a new dataset from an existing one*
- All *transformations* in Spark *are lazy*: they do not compute their results right away
- Instead they remember the transformations applied to some base dataset
- This helps in:
 - Optimizing the required calculations*
 - Recovery of lost data partitions*



edureka!

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

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Playing with Spark RDDs > Presentation

e! e! e!

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

Personal Library

Presentation

23 / 100 ← →

Some Popular Transformations

map filter distinct sortBy

Returns a new distributed dataset formed by passing each element of the source through a function

Example

```
>>> data=sc.parallelize("Hello Harry? How was your day".split(" "))  
>>> modified_data=data.map(lambda word: (word, word[0], word.startswith("H")))  
>>> modified_data.collect()  
[('Hello', 'H', True), ('Harry?', 'H', True), ('How', 'H', True), ('was', 'w', False), ('your', 'y', False), ('day', 'd', False)]
```

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

Personal Library

Presentation

Some Popular Transformations

map filter distinct sortBy

Returns a new dataset formed by selecting those elements of the source on which function returns true

Example

```
>>> data=sc.parallelize("Hello Harry? How was your day".split(" "))
>>> modified_data=data.map(lambda word: (word, word[0], word.startswith("H")))
>>> modified_data.filter(lambda answer: answer[2]).take(2)
[('Hello', 'H', True), ('Harry?', 'H', True)]
```

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

Presentation

Some Popular Transformations

map filter distinct sortBy

Returns a new dataset which formed by removing duplicates

Example

```
>>> data=sc.parallelize("Welcome to edureka")
>>> data.distinct().count()
13
```

edureka!

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

Course Classroom | Edureka

https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Inpu...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

edureka!

Some Popular Transformations

map filter distinct sortBy

Returns a new dataset which formed by sorting the elements of the RDD

Example

```
>>> rdd1=sc.parallelize("welcome to edureka".split(" "))  
>>> rdd1.sortBy(lambda line: len(line)*-1).collect()  
['welcome', 'edureka', 'to']
```

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

31°C Haze

Search

edureka!

ENG IN

01-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Playing with Spark RDDs > Presentation



e! Course Content

Presentation

Case Study I - Perform...

Case Study II - Load ...

In-Class Demo - Input...

Dataset

Dataset

RDD Practicals

Apache Sqoop

Word_Count_Class

Word_Count

Personal Library

e! Presentation

27 / 100

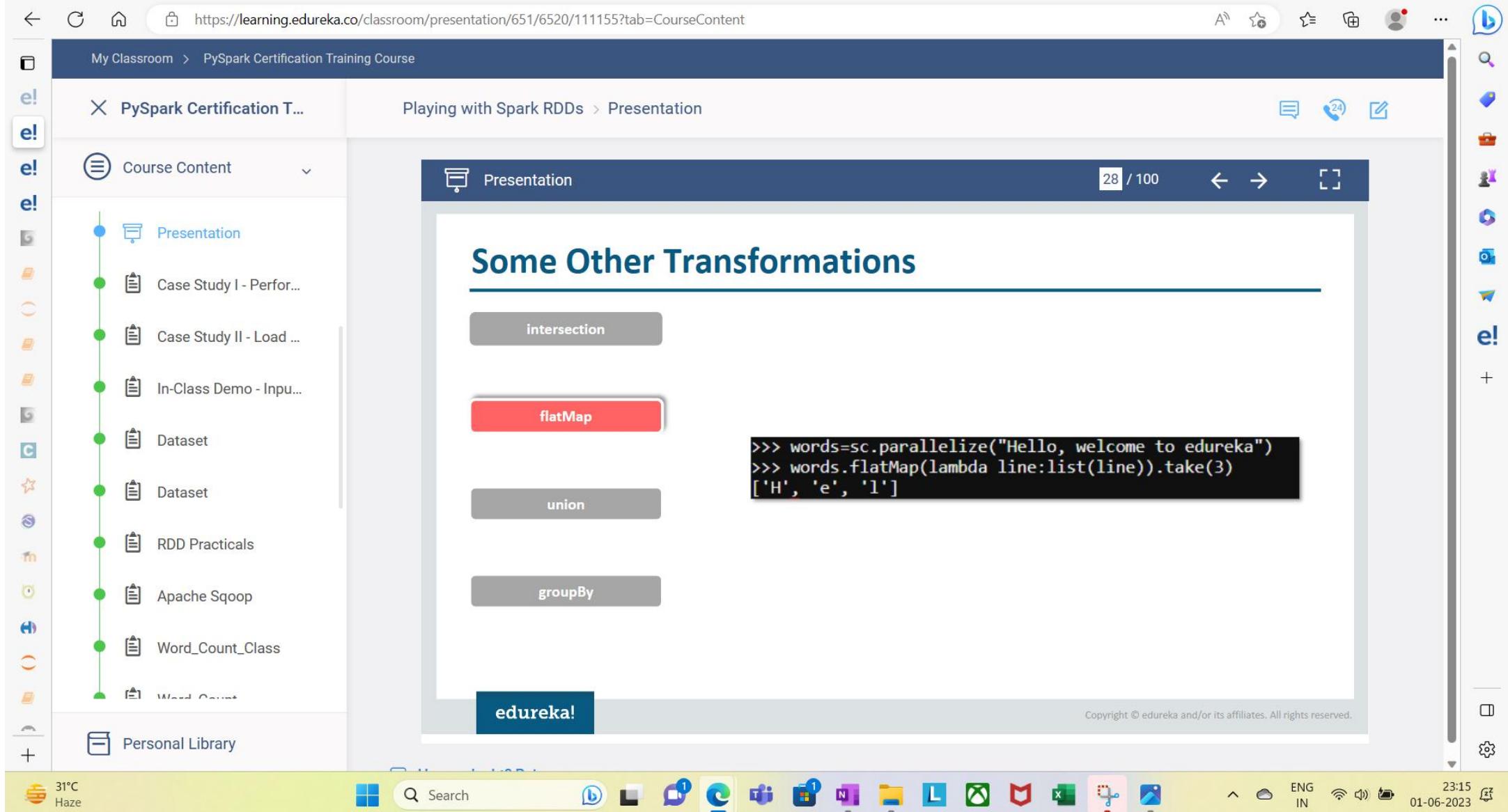


Some Other Transformations

intersection**flatMap****union****groupBy****edureka!**

```
>>> rdd1 = sc.parallelize(((1,"jan",2016),(3,"nov",2014), (16,"feb",2014)))
>>> rdd2 = sc.parallelize(((5,"dec",2014),(1,"jan",2016)))
>>> command = rdd1.intersection(rdd2)
>>> command.collect()
[(1, 'jan', 2016)]
```

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



My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 29 / 100 ← →

Some Other Transformations

intersection

flatMap

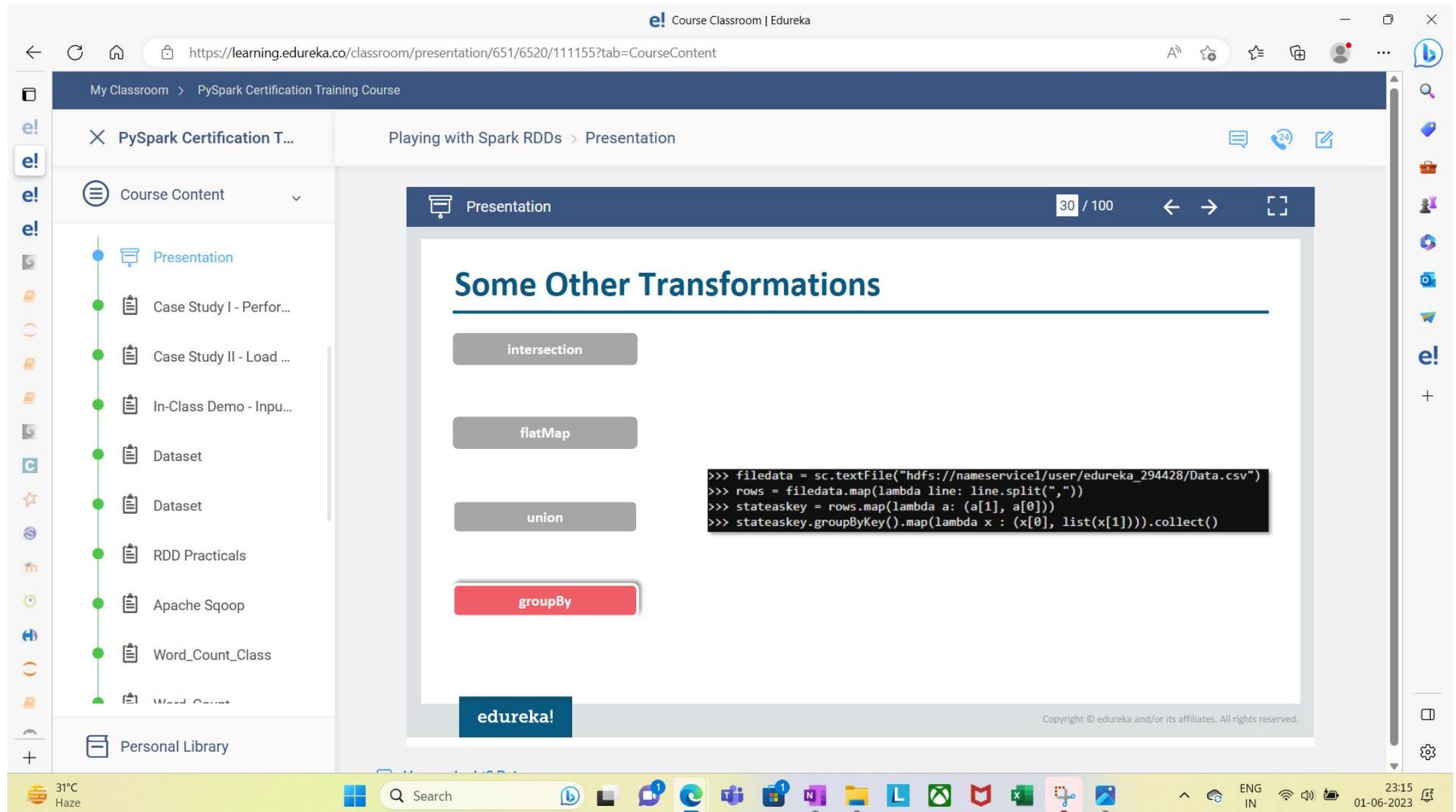
union

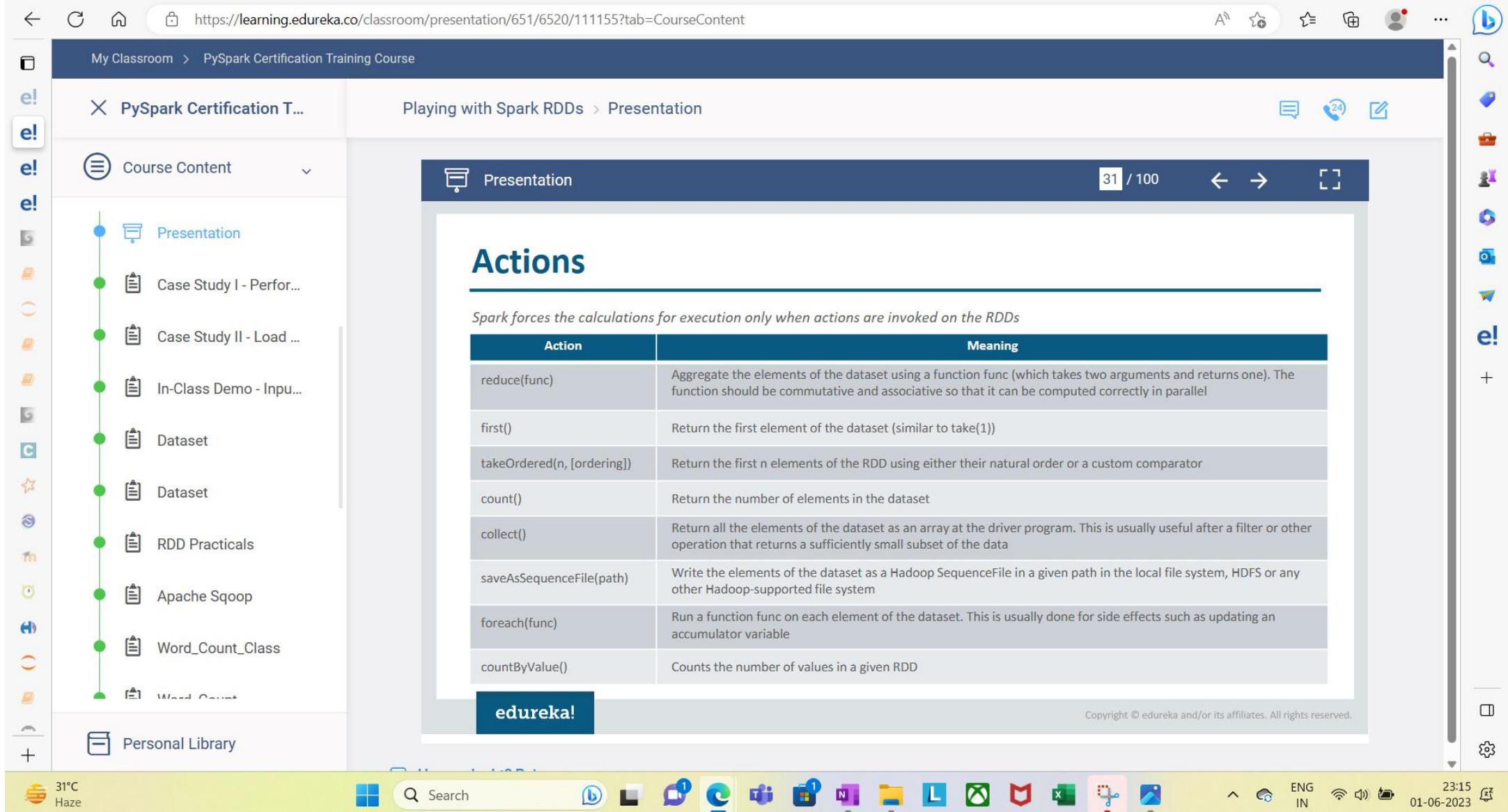
groupBy

edureka!

```
>>> inputRDD = sc.textFile("hdfs://nameservice1/user/edureka_294428/demo.txt")
>>> HelloRDD = inputRDD.filter(lambda x: "Hello" in x)
>>> WelcomeRDD = inputRDD.filter(lambda x: "Welcome" in x)
>>> LinesRDD = HelloRDD.union(WelcomeRDD)
```

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





My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 32 / 100 ← →

Some Popular Actions

The diagram illustrates four popular Spark actions:

- reduce: Represented by a red circle with a red arrow pointing to a red dashed box labeled "reduce".
- count: Represented by a green circle with a blue arrow pointing to a blue dashed box labeled "count".
- first: Represented by a blue circle with a blue arrow pointing to a blue dashed box labeled "first".
- countByValue: Represented by an orange circle with a blue arrow pointing to a blue dashed box labeled "countByValue".

Aggregate the elements of the dataset using a function func (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel

```
>>> sc.parallelize(range(1,50)).reduce(lambda a, b: a+b)
1225
```

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 33 / 100

Some Popular Actions

The diagram illustrates four popular RDD actions:

- reduce
- count
- first
- countByValue

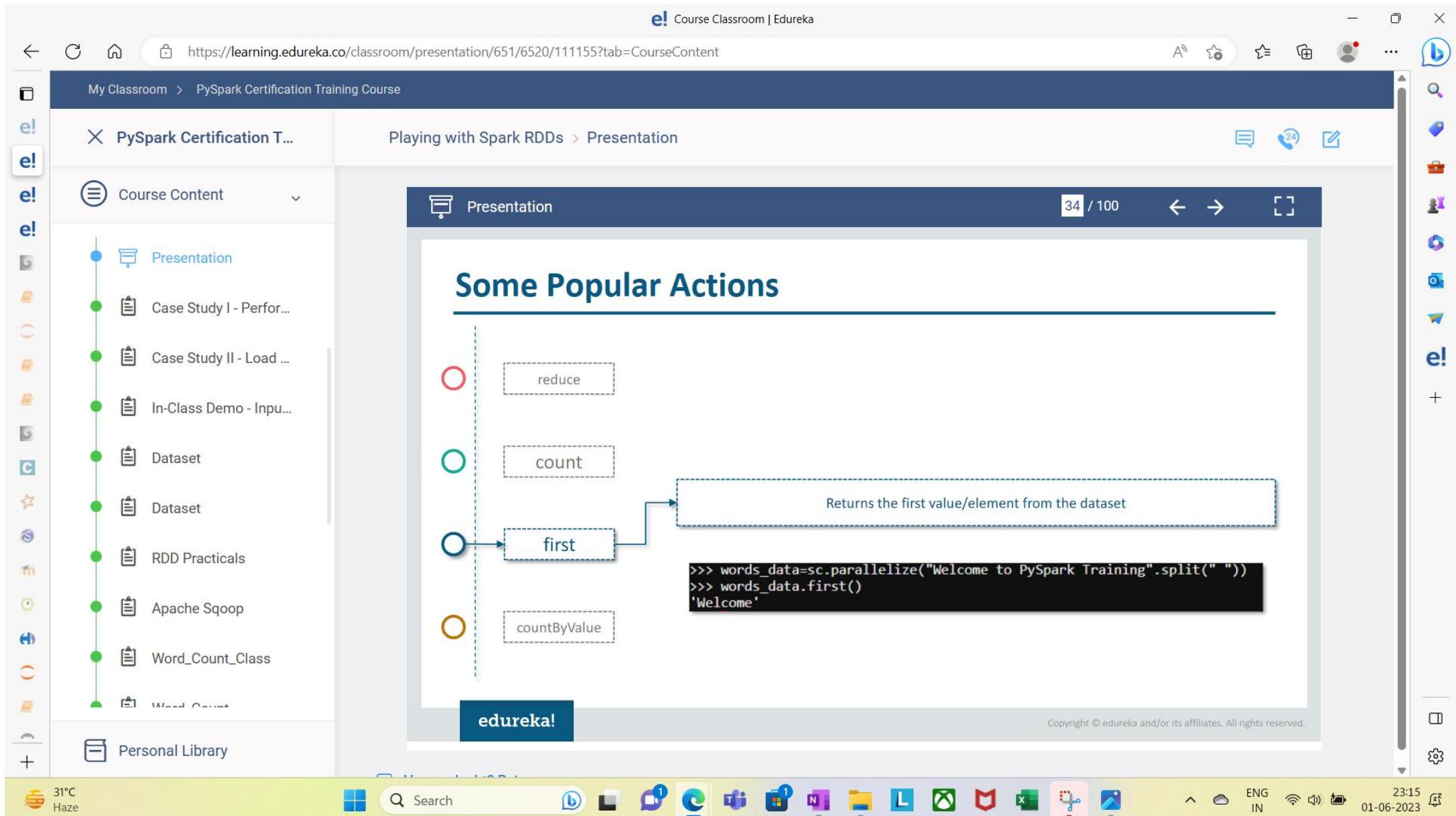
Arrows show the flow from the first action to the count action.

Returns the number of rows in the dataset(i.e. RDD)

```
>>> words_data=sc.parallelize("Welcome to PySpark Training".split(" "))  
>>> words_data.count()  
4
```

edureka!

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



Presentation

Some Popular Actions

reduce

count

first

countByValue

Returns the values present in the RDD

```
>>> words_data=sc.parallelize("Welcome to PySpark Training".split(" "))
>>> words_data.countByValue()
defaultdict(<type 'int'>, {'to': 1, 'Training': 1, 'Welcome': 1, 'PySpark': 1})
```

edureka!

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

https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

edureka!

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

31°C Haze

Search

L

ENG IN

01-06-2023

23:15

https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

Personal Library

Presentation

37 / 100 ← →

Writing on File Using saveAsTextFile and flatMap

```
rdd.flatMap {lambda line: line.split(' ') }.saveAsTextFile(outputPath)
```

The memory usage of this is actually fairly straightforward. In the case of a text file, it streams the file in, line by line applies your flatMap operation and immediately writes (possibly buffering) it out to disk. Therefore you're not really doing anything terribly memory intensive and Spark doesn't really do anything fancy in terms of coordination

edureka!

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

31°C Haze

Search

L

ENG IN

01-06-2023

Writing on File Using saveAsTextFile and flatMap

```
rdd.flatMap {lambda line: line.split(' ') }.saveAsTextFile(outputPath)
```

If you have two files on disk that are read in by this application, it will read in the two files (partitions) in parallel and write them out in parallel. No coordination is required and this is what is known as an "embarrassingly parallel" process.

edureka!

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

My Classroom > PySpark Certification Training Course

Playing with Spark RDDs > Presentation



Course Content

-  Presentation
 -  Case Study I - Perform...
 -  Case Study II - Load ...
 -  In-Class Demo - Inpu...
 -  Dataset
 -  Dataset
 -  RDD Practicals
 -  Apache Sqoop
 -  Word_Count_Class

Reading File and Counting Word Frequency

```
rdd.flatMap { lambda line : line.split(' ') }.map((_, 1)).reduceByKey((a, b) => a + b).collect()
```

Lets see the working of this line of code

Note: This code splits each line by whitespace, turns each word into a row and each row into a tuple of (word, 1) which then gets aggregated to result in (word, # of occurrences for that word)

X PySpark Certification T

Playing with Spark RDDs > Presentation



 Course Content

-  Presentation
 -  Case Study I - Performance
 -  Case Study II - Load Data
 -  In-Class Demo - Input
 -  Dataset
 -  Dataset
 -  RDD Practicals
 -  Apache Sqoop
 -  Word_Count_Class
 -  Word_Count

Reading File and Counting Word Frequency

```
rdd.flatMap { lambda line : line.split(' ') }.map((_, 1)).reduceByKey((a, b) => a + b).collect()
```

In this case, you're actually going to begin by performing the same streaming, parallel map operation that takes place in the previous example, but things get a little more complex after that.

edureka!

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 41 / 100 ← →

Reading File and Counting Word Frequency

```
rdd.flatMap { lambda line : line.split(' ') }.map((_, 1)).reduceByKey((a, b) => a + b).collect()
```

In order to count all of the given words which may exist across an entire dataset in Spark, each partition must aggregate all the counts of words within that partition, but then it must also sum across partitions.

Step #1

edureka!

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

31°C Haze

Search

L

ENG IN

01-06-2023

Reading File and Counting Word Frequency

```
rdd.flatMap { lambda line : line.split(' ') }.map((_, 1)).reduceByKey((a, b) => a + b).collect()
```

The process of moving the data from partition to partition in order to aggregate, join, match up, or spread out in some other way, is known as **shuffling**.

Step #2

Step #1

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 43 / 100 ← →

Reading File and Counting Word Frequency

```
rdd.flatMap { lambda line : line.split(' ') }.map((_, 1)).reduceByKey((a, b) => a + b).collect()
```

The following operations are examples of shuffle inducing operations for RDDs:

- groupBy/subtractByKey/foldByKey/aggregateByKey/reduceByKey
- cogroup
- any of the join transformations
- distinct

Step #1 Step #2 Step #3 : Shuffling

edureka!

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

https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

edureka!

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

31°C Haze

Search

L

ENG IN

01-06-2023

23:15

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 45 / 100 ← →

Data Loading in RDD

- RDD is *re-computed* every time when it is *materialized*
- So it is a good idea to *improve performance* by *caching* a RDD if it is accessed frequently
- One of the easiest way to *load data in RDD* is to load from a *Python collection*:

```
>>> rdd=sc.parallelize([1, 2, 3, 4, 5])
>>> rdd.collect()
[1, 2, 3, 4, 5]
```
- SparkContext provides parallelize function, which converts the Python collection into the RDD of the same type:

```
dataRDD: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <console>:12
```

The diagram shows a teal cylinder labeled "Data on disk" connected by an arrow labeled "HDFS read" to an orange rectangular block labeled "Distributed Memory". Three arrows labeled 1, 2, and 3 point from the "Distributed Memory" block to three separate blue rectangular boxes, each labeled "RDD".

edureka!

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

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Inpu...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

Data Loading in RDD (Contd.)

- Another simple way of *loading data* is *loading text from a file*
- In *local mode* (single node operation), you just need to *specify the file location*
- In *distributed mode*, files should be *accessible from all nodes* of the cluster
- Spark's *addFile* functionality is used to *copy the file* to designated location to all machines in the cluster

```
>>> from pyspark.files import SparkFiles  
>>> sc.addFile("/mnt/home/edureka_294428/demo.txt")  
>>> sc.textFile(SparkFiles.get("/mnt/home/edureka_294428/demo.txt"))
```

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 47 / 100 ← →

Saving Data Through RDD

- While the methods for loading an RDD are largely found in the `SparkContext` class
- Methods for saving an RDD are defined on the RDD classes
- In Python, implicit conversion exists so that an RDD that can be saved as a sequence file is converted to the appropriate type

```
>>> words=sc.parallelize("Hello, welcome to edureka")
>>> words.saveAsTextFile("/user/edureka_294428/output.txt")
```

The diagram illustrates the saving process of an RDD. It starts with a "Data on disk" icon, which is connected via an arrow labeled "HDFS Read" to a "Distributed Memory" folder icon. This folder contains three numbered RDD icons (1, 2, 3). Each RDD is connected via a dashed arrow to a "Query" button, which then points to a "Result" box. Finally, each result box has an arrow pointing to a second "Data on disk" icon.

edureka!

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

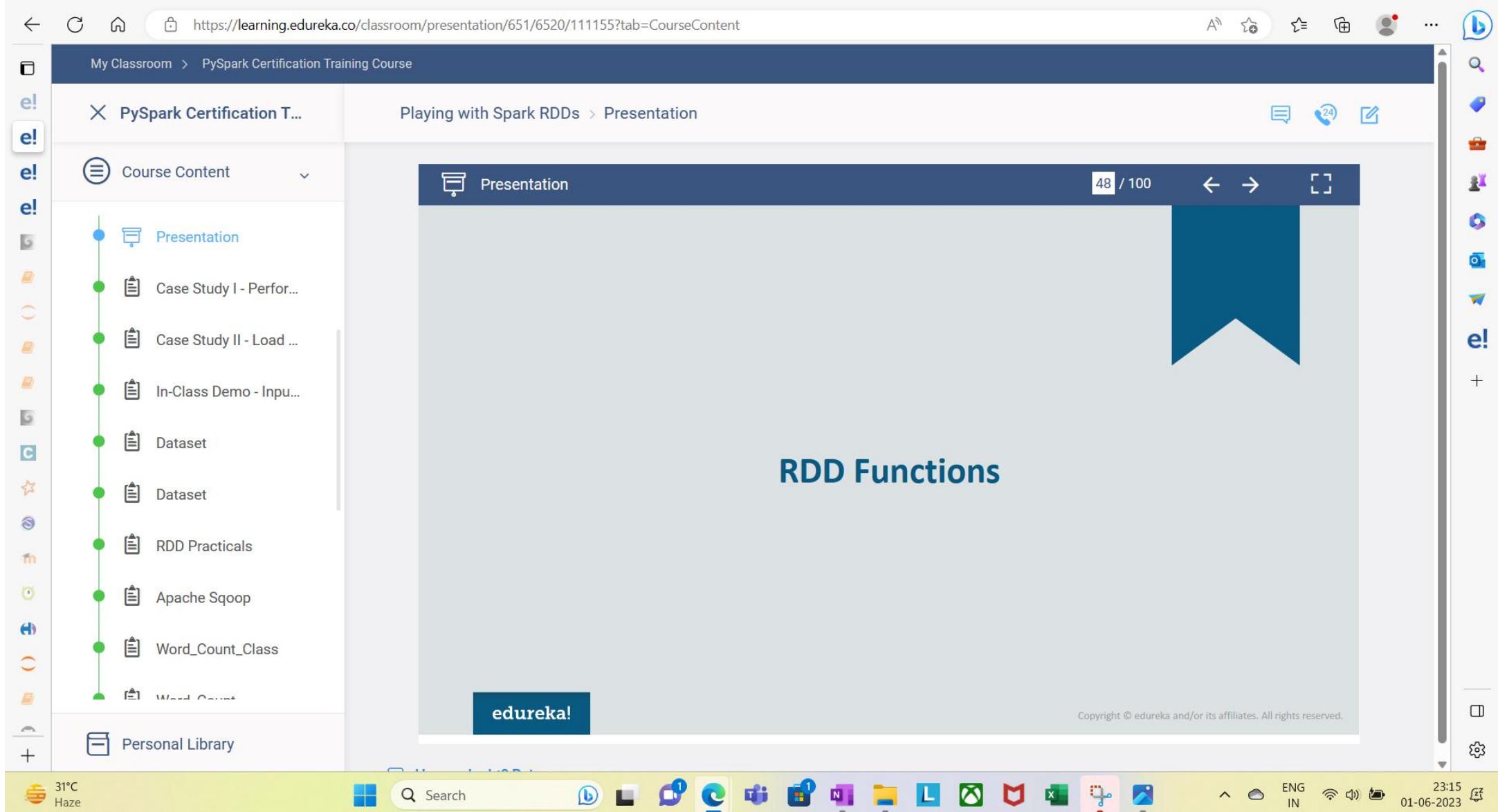
31°C Haze

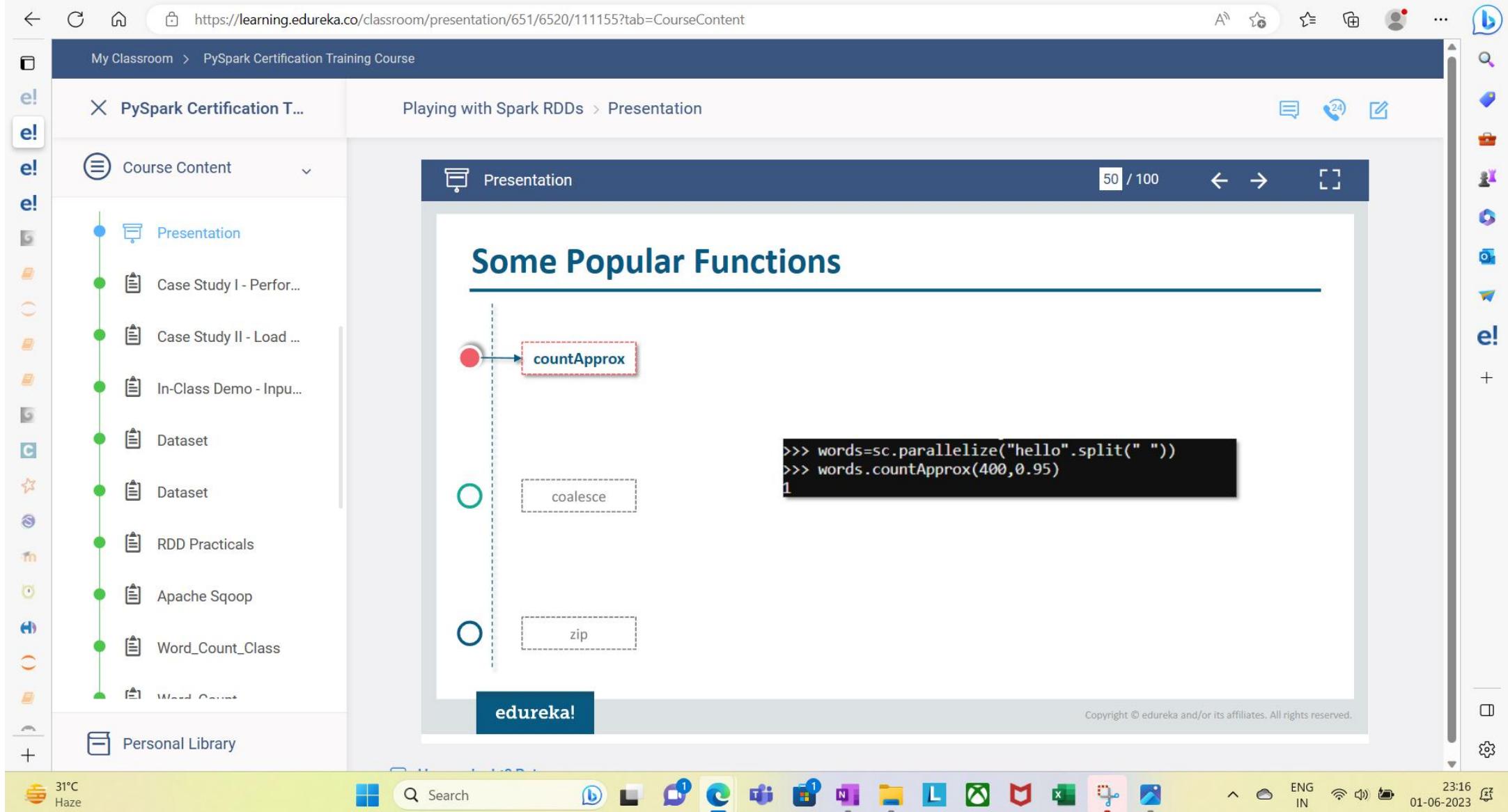
Search

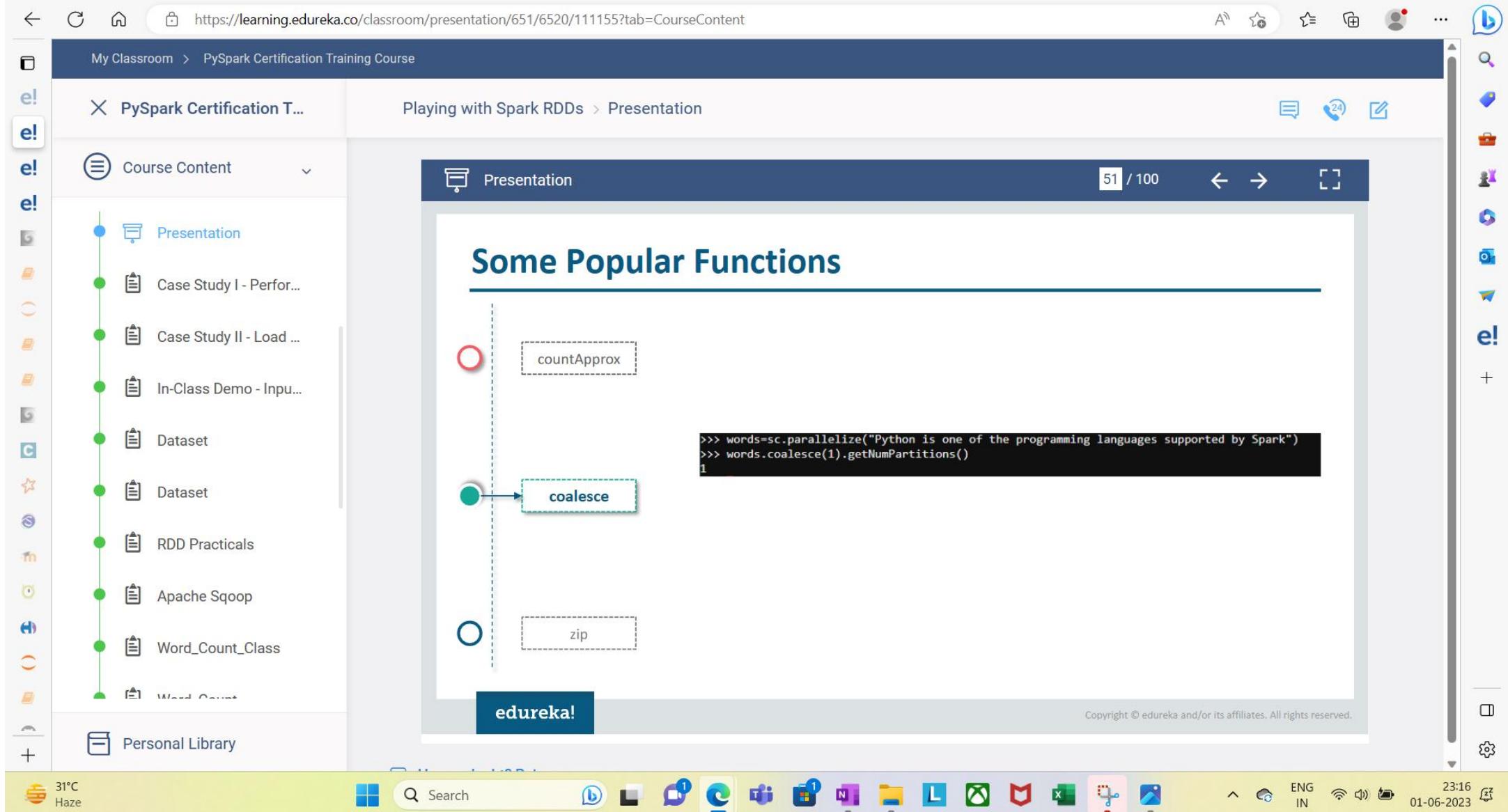
Cloud

ENG IN

23:15 01-06-2023









X PySpark Certification T...

Playing with Spark RDDs > Presentation



Some Popular Functions

countApprox

```
>>> words=sc.parallelize("Python is one of the programming languages supported by Spark".split(" "))
>>> numRange = sc.parallelize(range(10), 2)
>>> words.zip(numRange).collect()
```

coalesce

zip

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

https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

edureka!

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

31°C Haze

Search

L

ENG IN

23:16 01-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 54 / 100

Key-Value Pair RDD

- Pair RDDs can be created by running a map() function that returns key-value pairs. The procedure to build the key-value RDDs differs by language
- In Python, for the functions on keyed data to be available, we also need to return tuples

Function	Parameter options	Explanation	Return Type
foldByKey	(zeroValue) (func(V,V)=>V)	Merges the values using the provided function. Unlike a traditional fold function over a list, the zeroValue can be added an arbitrary number of times.	RDD[K,V]
reduceByKey	(func(V,V)=>V, numTasks)	Parallel version of reduce that merges the values for each key using the provided function and returns an RDD.	RDD[K,V]
groupByKey	(numPartitions)	Groups elements together by key.	RDD[K,Seq[V]]

edureka!

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

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

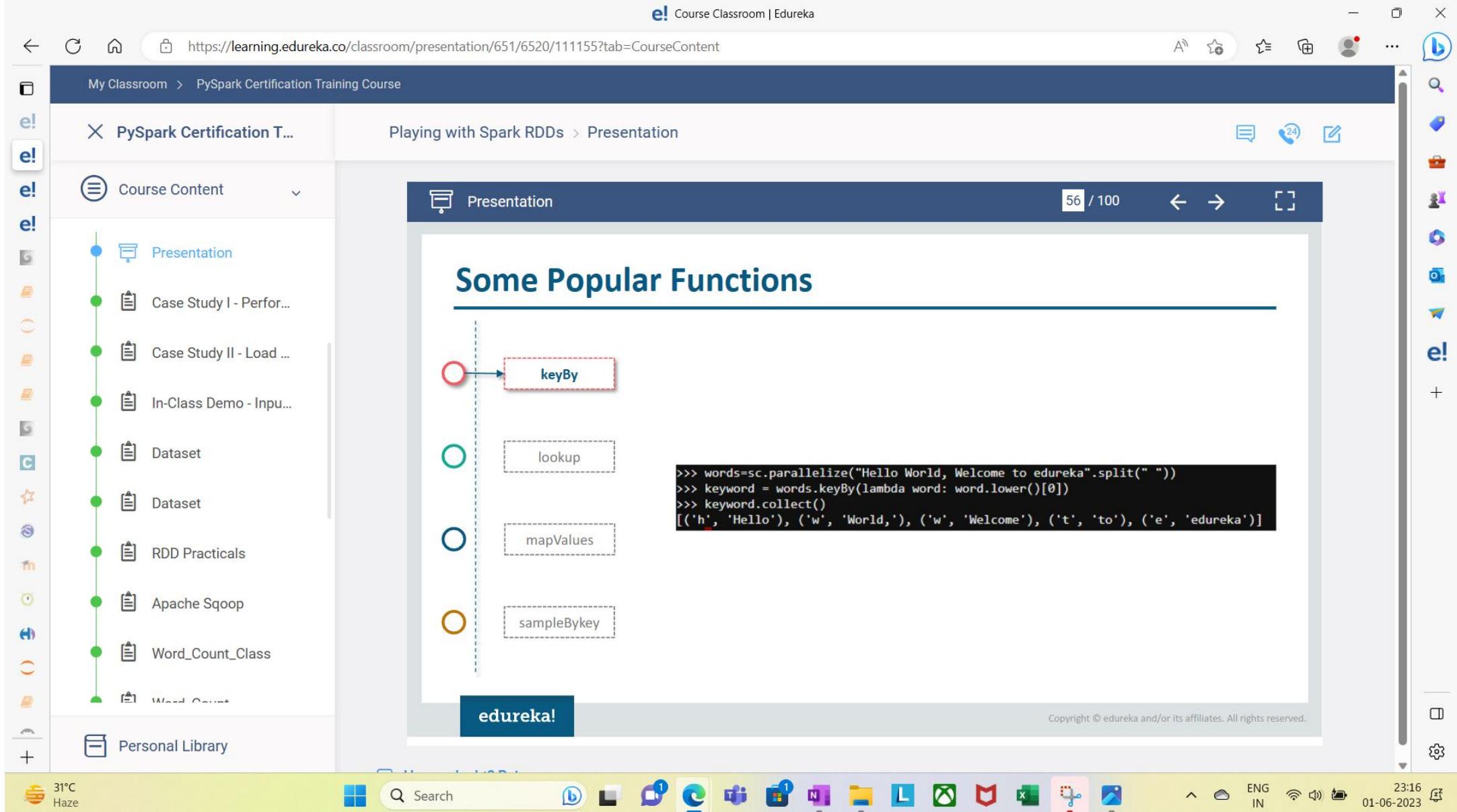
Presentation 55 / 100 ← →

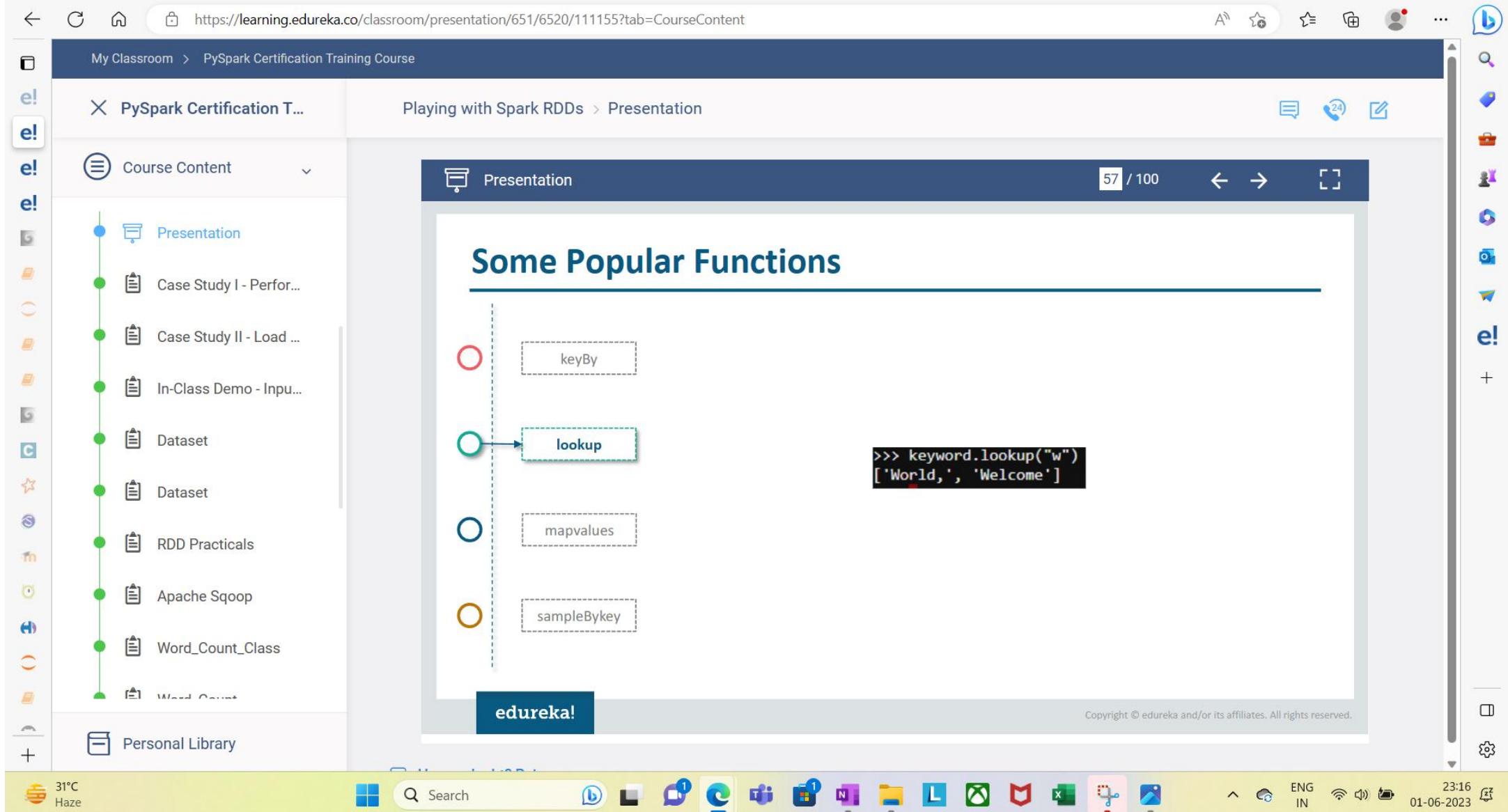
Other Pair RDD Functions

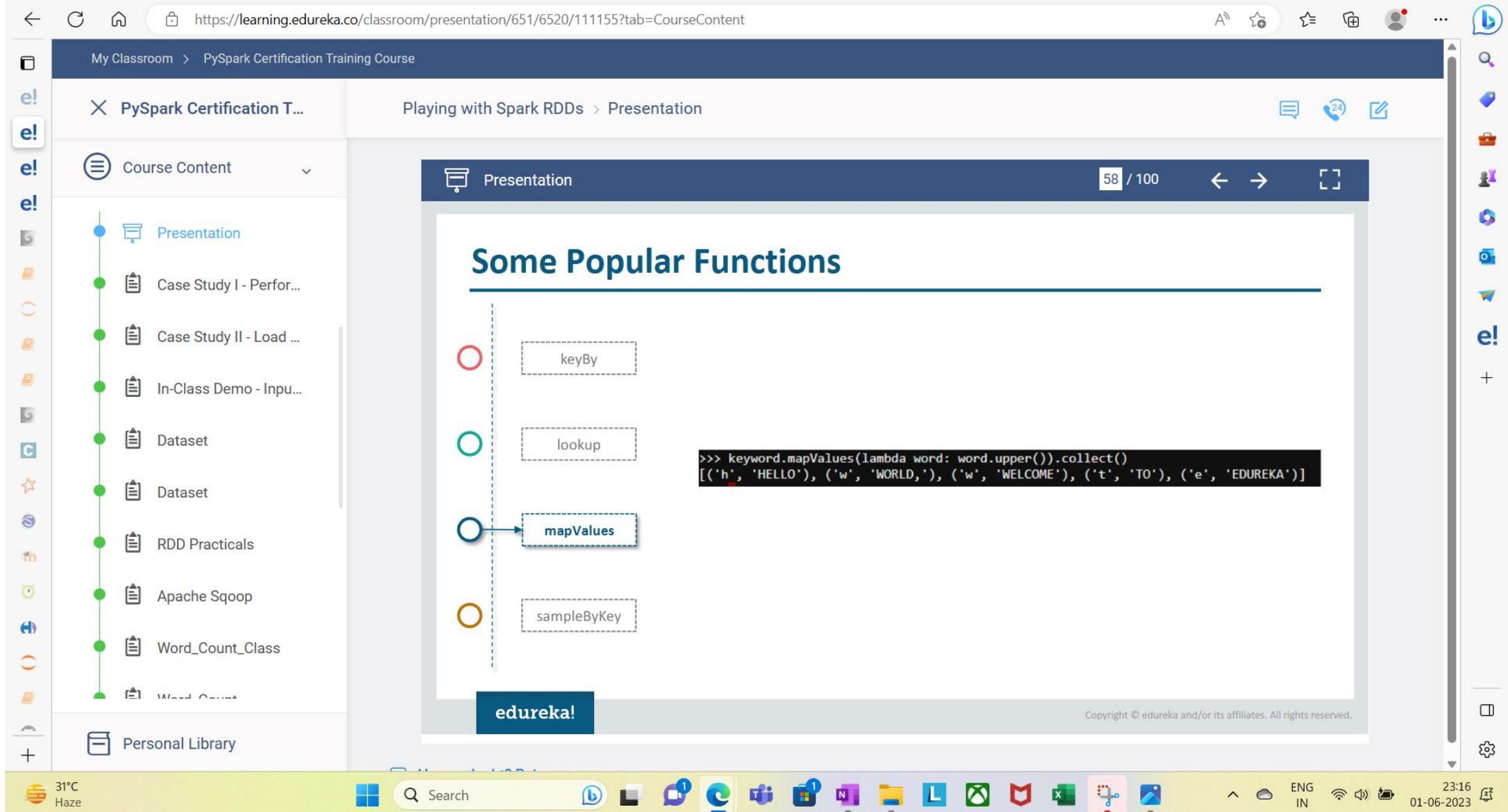
Function	Parameter options	Explanation	Return Type
lookup	(key: K)	Looks up a specific element in the RDD. Uses the RDD's partitioner to figure out which partition(s) to look at.	Seq[V]
mapValues	(f: v=>u)	A specialized version of map for PairRDD when you only want to change the value of the key-value pair. If you need to make your changes based on both key and value, you must use one of the normal RDD Map functions.	RDD[K,U]
collectAsMap	()	Takes an RDD and returns a concrete map. Your RDD must be able to fit into the memory.	Map[K, V]
countByKey	()	Counts the number of elements for each key in RDD.	Map[K, Long]
partitionBy	(partitioner:P, mapSideCombine: Boolean)	Returns a new RDD with the same data but partitioned by the new Partitioner, and mapSideCombine controls Spark group values with the same key together before repartitioning. Defaults to false.	RDD[K,V]
flatMapValues	(f:V => TraversableOnce[U])	Similar to mapValues. A specialized version of flatMap for PairRDDs when you only want to change the value of the key-value pair. Takes the provided Map function and applies it to the value. The resulting sequence is then "flattened".	RDD[K, U]

edureka!

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







Some Popular Functions

keyBy

lookup

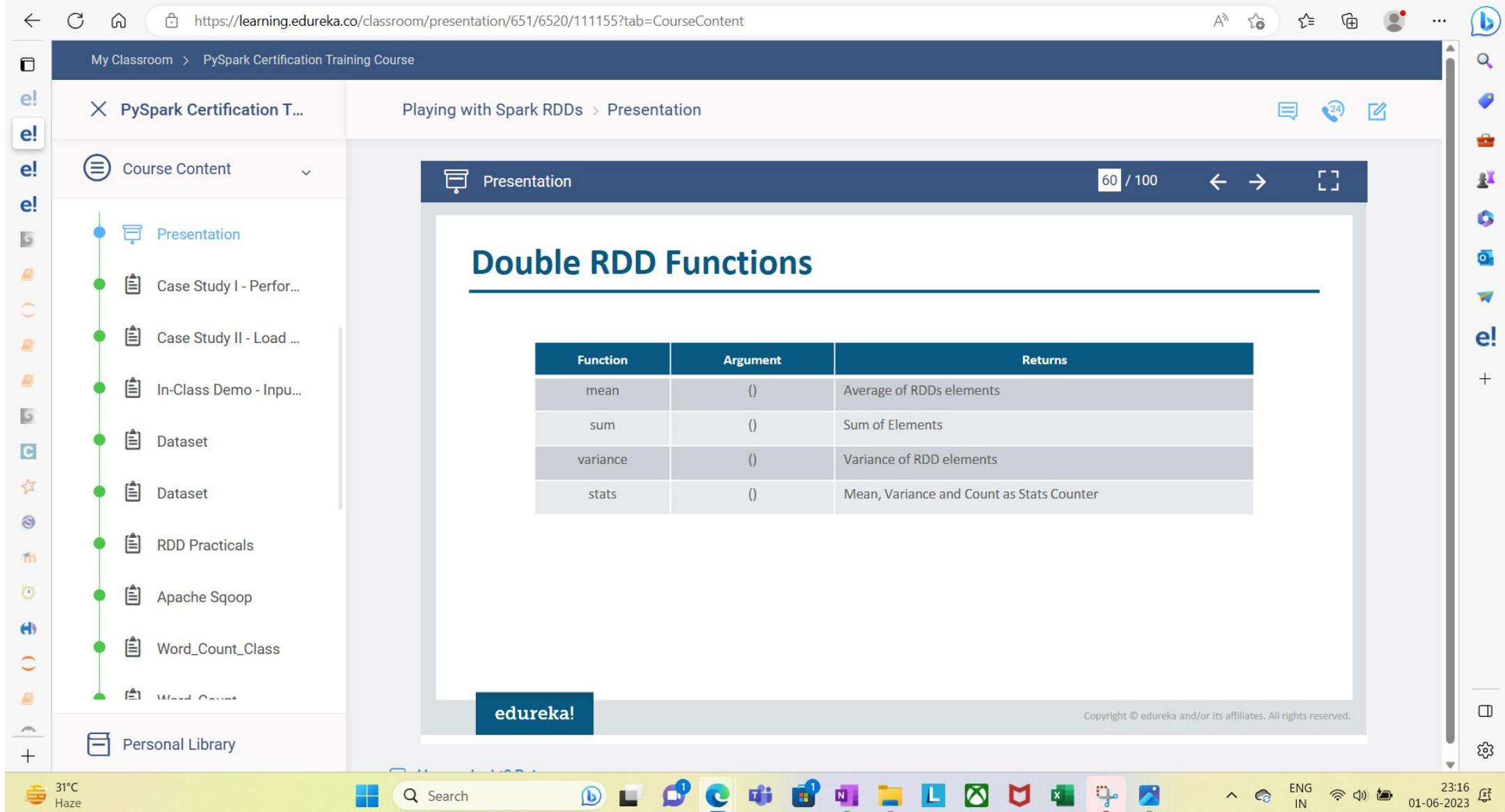
mapValues

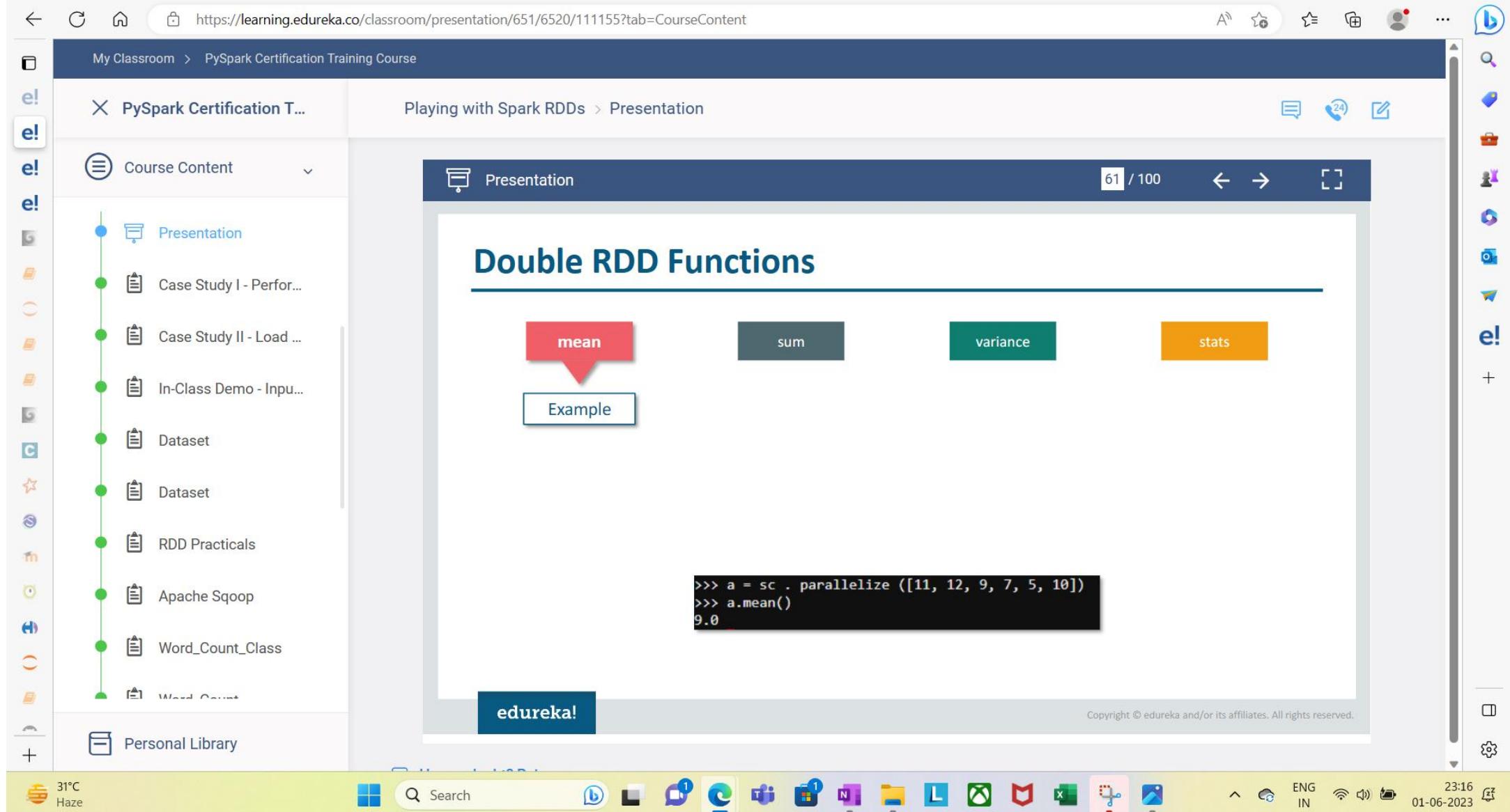
sampleByKey

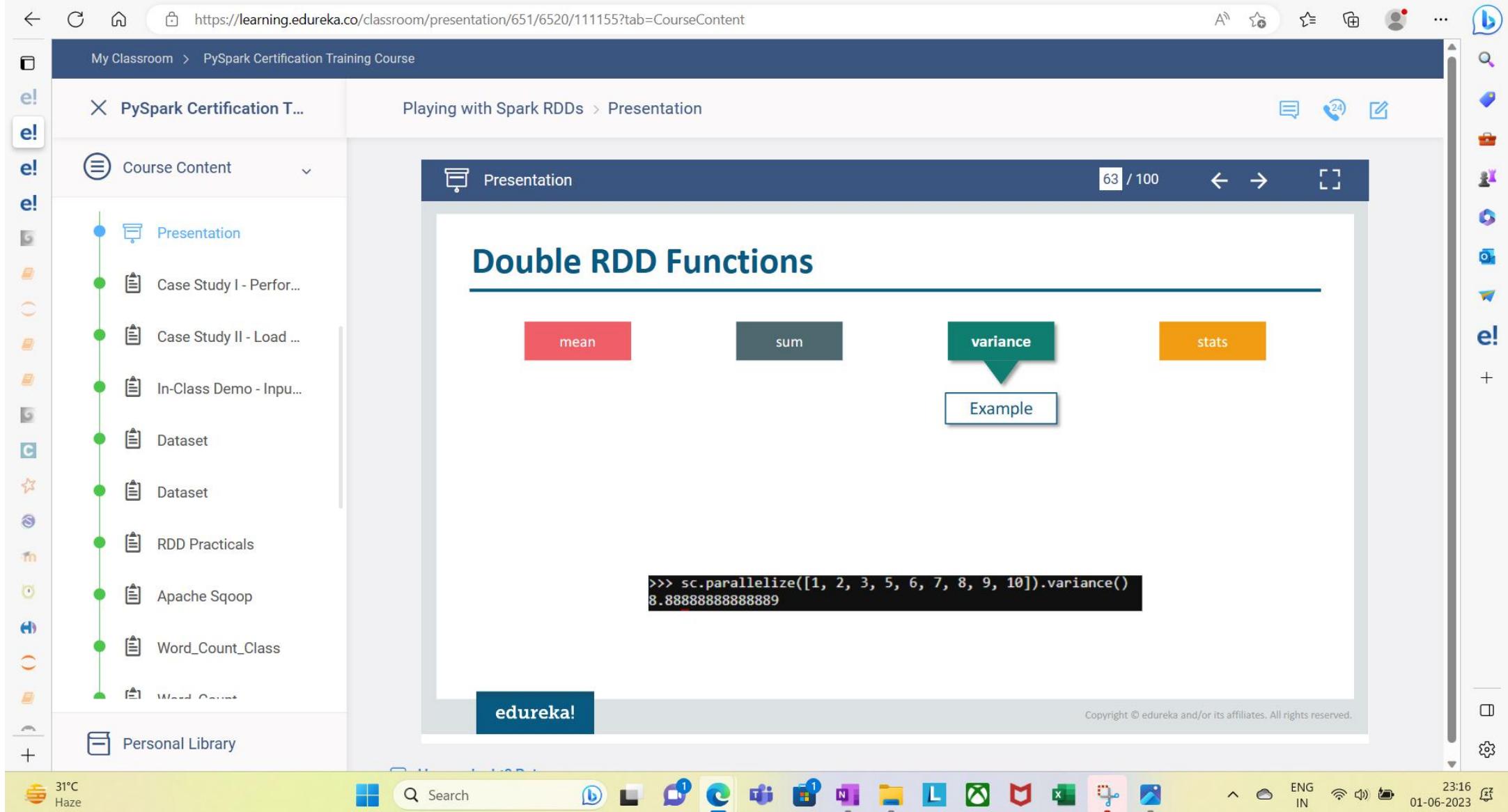
```
>>> import random
>>> distinctChars = words.flatMap(lambda word: list(word.lower().distinct().collect()))
>>> sampleMap = dict(map(lambda c: (c, random.random()), distinctChars))
>>> words.map(lambda word: (word.lower()[0], word)).sampleByKey(True, sampleMap, 6).collect()
[('e', 'edureka')]
```

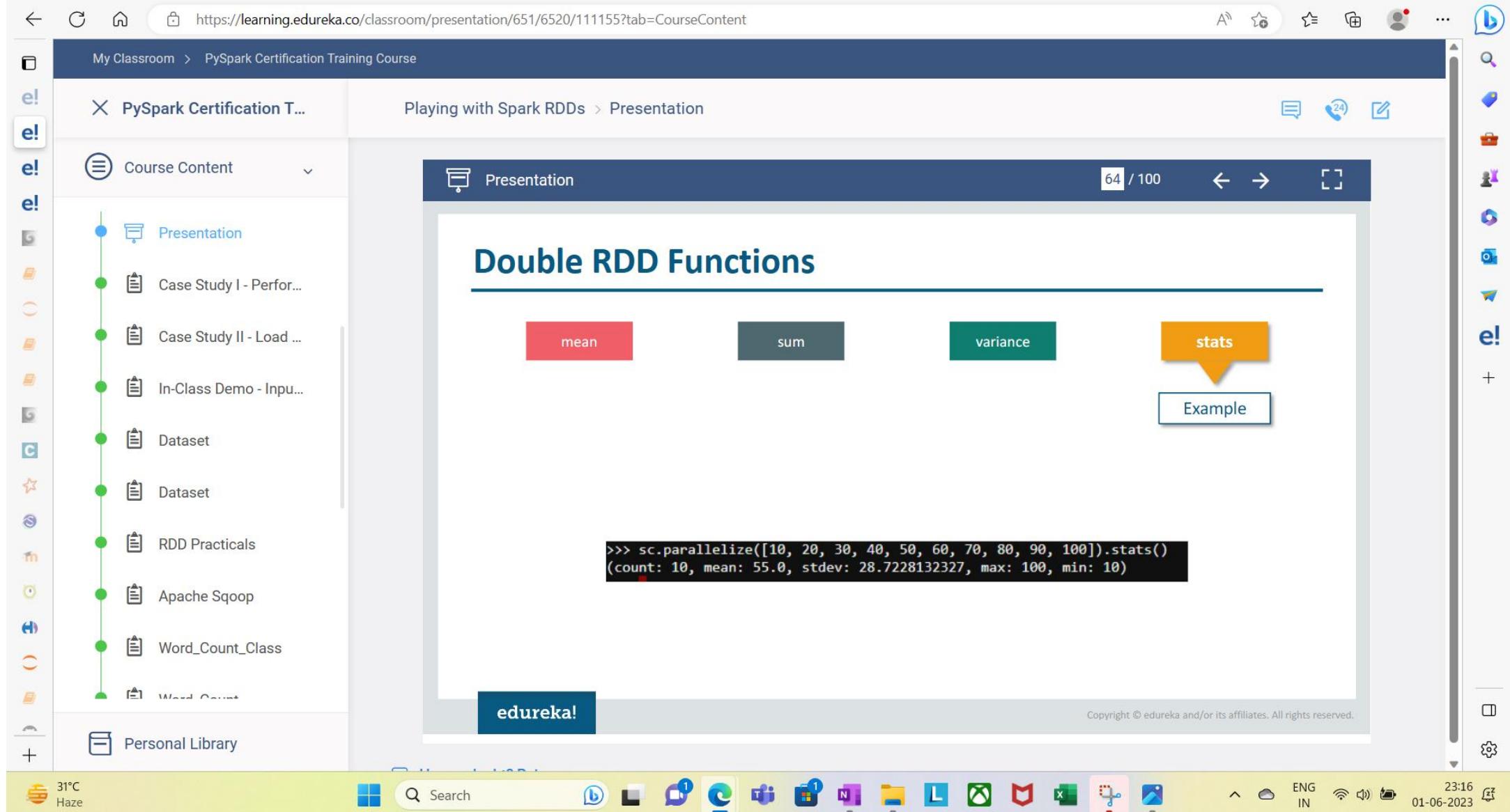
edureka!

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









X PySpark Certification T...

Playing with Spark RDDs > Presentation



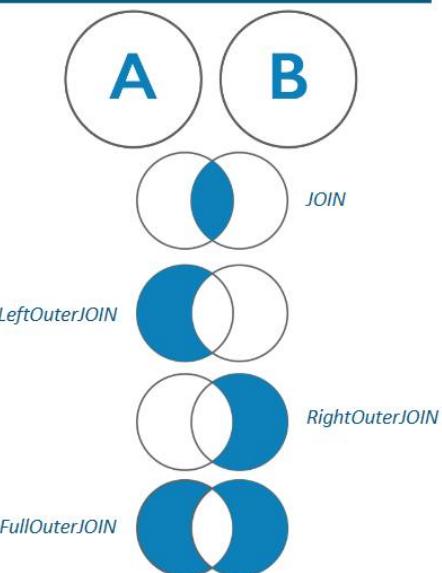
e! Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

e! Personal Library

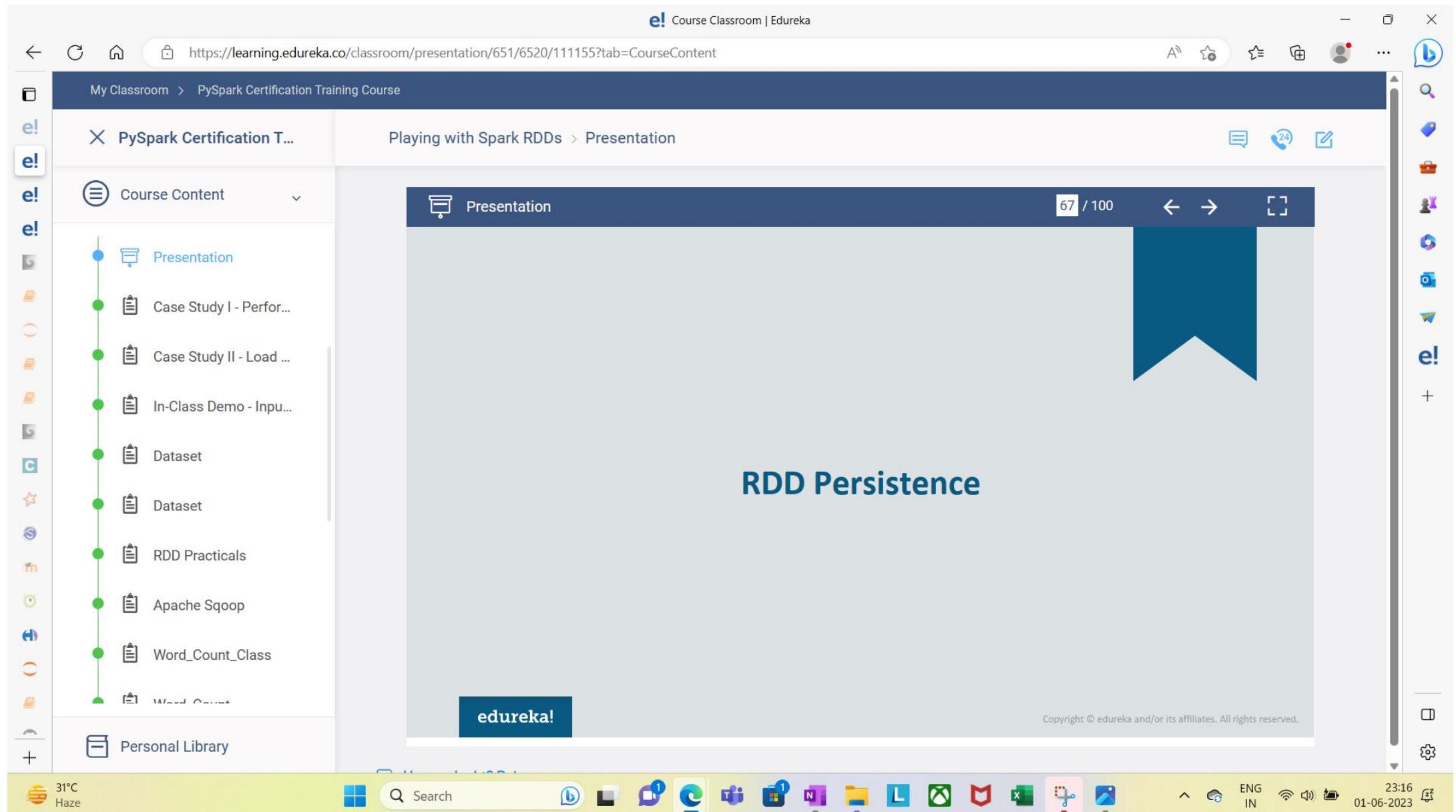
RDD Joins

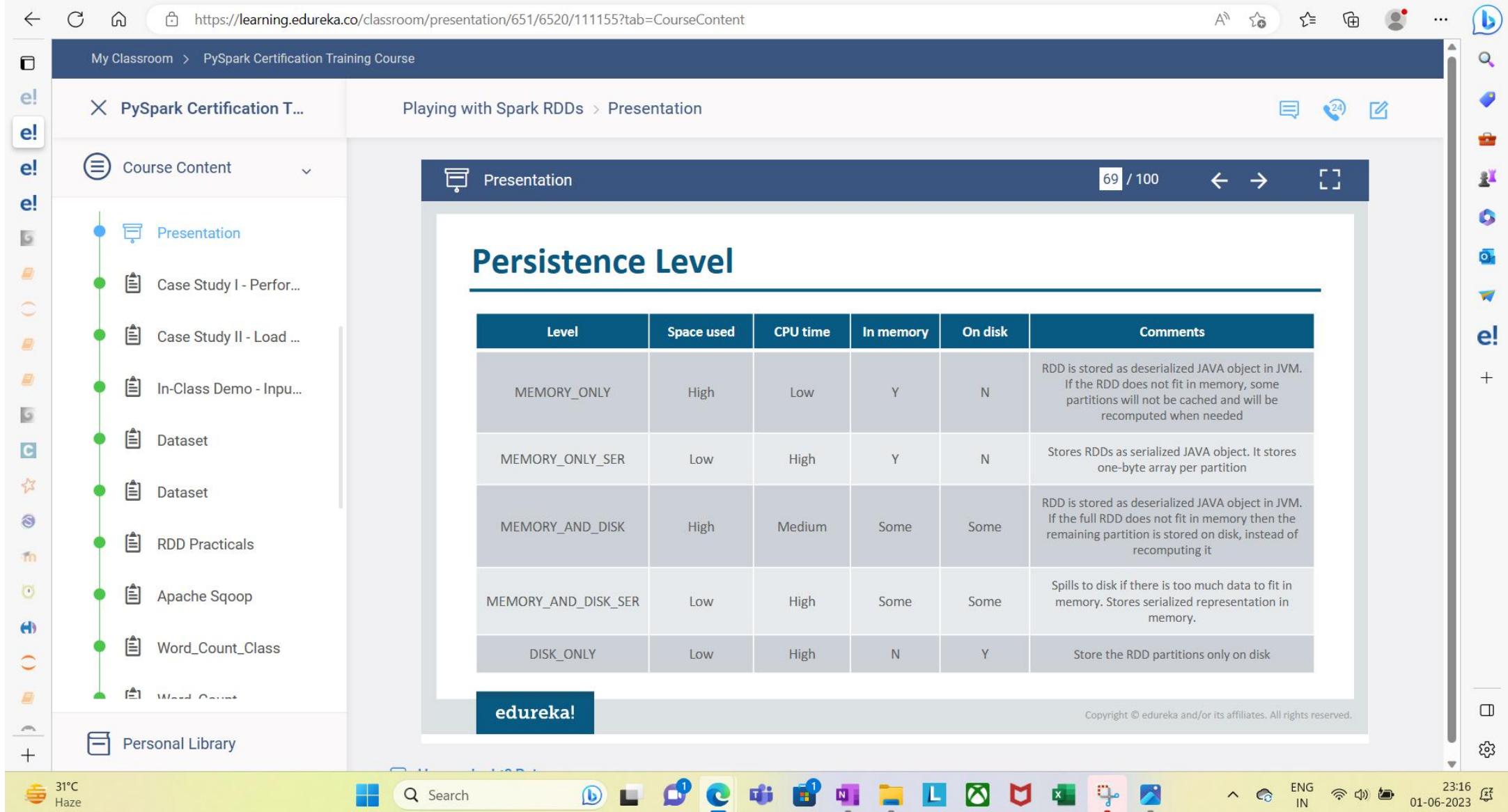
- Spark allows *easy manipulation* of *multiple data sets* with out-of-the-box join operators
- Spark Pair RDDs provide *join*, *leftOuterJoin*, *rightOuterJoin*, *fullOuterJoin* methods to perform respective joins
- Programmers need to create the RDDs carefully, as the *keys* are *automatically chosen by the join operation*
- Pair RDDs also provide the *cartesian* and *cogroup* methods to perform the respective operations



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

edureka!







My Classroom > PySpark Certification Training Course



X PySpark Certification T...



Course Content



 Presentation



 Case Study I - Perform...



 Case Study II - Load ...



 In-Class Demo - Input...



 Dataset



 Dataset



 RDD Practicals



 Apache Sqoop



 Word_Count_Class



 Word_Count



Personal Library

Playing with Spark RDDs > Presentation



 Presentation

71 / 100



Demo 2 – Word Count Program

Refer to the file Module-5 Demo 2 provided in the LMS for all the steps in detail

edureka!

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



31°C

Haze



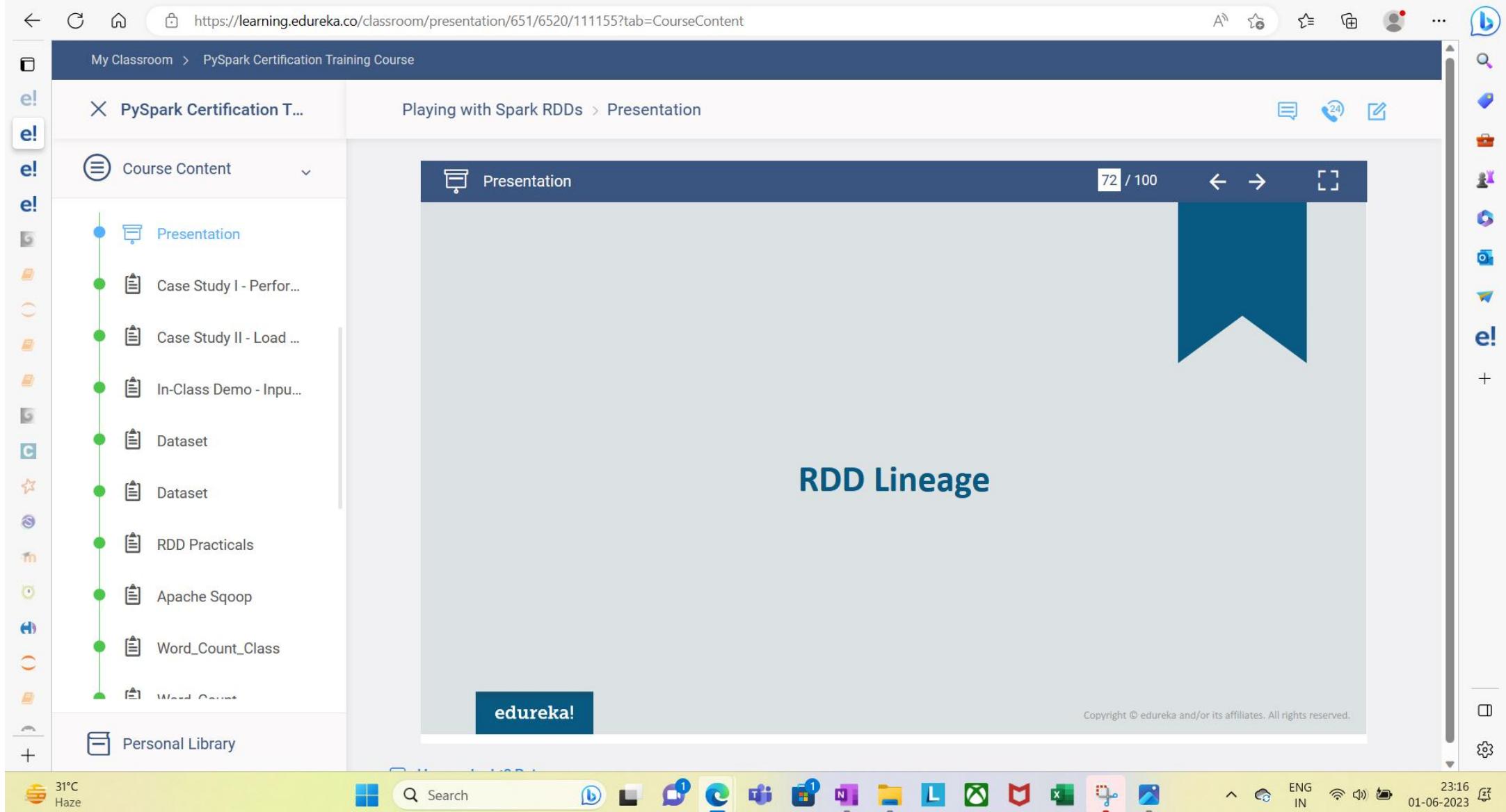
Search



ENG
IN



23:16
01-06-2023



https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

Presentation

73 / 100

RDD Lineage

- RDD Lineage** (RDD dependency graph) is a *graph* of all the *parent RDDs* of a *RDD*. It is built as a result of applying transformations to the RDD and creates a logical execution plan

The given RDD graph could be the result of the following series of transformations:

- A RDD lineage graph is hence a graph of what transformations need to be executed after an action has been called

edureka!

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

31°C Haze

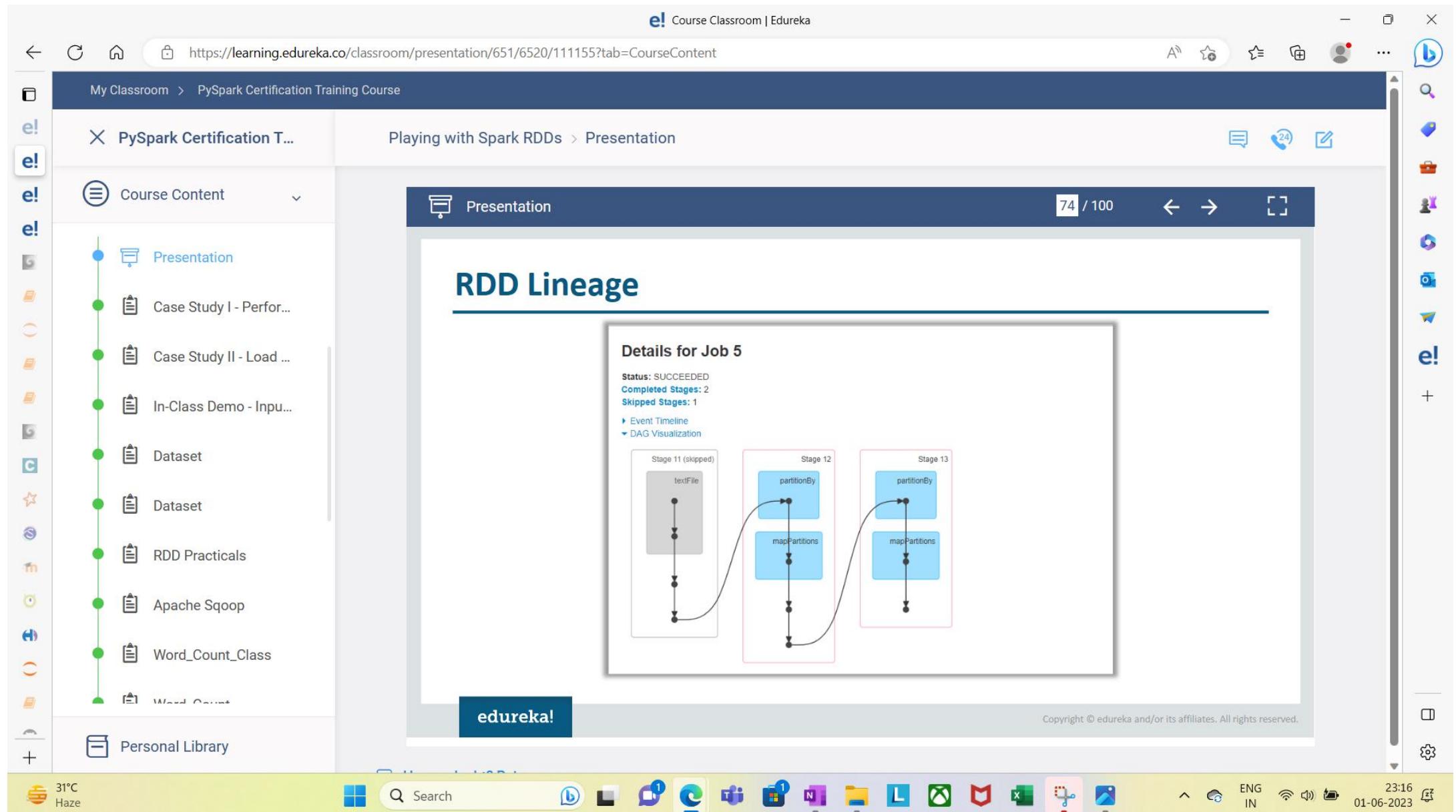
Search

L

ENG IN

23:16

01-06-2023



https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

edureka!

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

31°C Haze

Search

L

ENG IN

23:16 01-06-2023

https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 76 / 100 ← →

Partitioning RDD and Achieving Parallelism

- A **partition** is a *logical chunk* of a *large distributed data set*
- Spark *manages data using partitions* that helps *parallelize distributed data processing* with *minimal network traffic* for sending data between executors
- By default, Spark tries to *read data into an RDD* from the *nodes* that are *close to it*

Partitions of RDD

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

X PySpark Certification T...

Playing with Spark RDDs > Presentation



Course Content

-  Presentation
 -  Case Study I - Performance
 -  Case Study II - Load Testing
 -  In-Class Demo - Input
 -  Dataset
 -  Dataset
 -  RDD Practicals
 -  Apache Sqoop
 -  Word Count Class
 -  Word Count

Partitioning RDD and Achieving Parallelism

- Since Spark usually accesses distributed partitioned data, to optimize transformation operations it creates partitions to hold the data chunks
 - RDDs get partitioned automatically without programmer intervention
 - However, there are times when we'd like to adjust the size and number of partitions or the partitioning scheme according to the needs of our application
 - We use the given method on an RDD to know the number of partitions in this RDD: `rdd.getNumPartitions()`

```
>>> rdd3=sc.parallelize(range(1,25,3))
>>> rdd3.collect()
[1, 4, 7, 10, 13, 16, 19, 22]
>>> rdd3.getNumPartitions()
2
```



X PySpark Certification T...

Playing with Spark RDDs > Presentation



Executing a Parallel Task on Shell

```
>>> rdd1=sc.parallelize(range(1,100,5))
>>> rdd1.collect()
[1, 6, 11, 16, 21, 26, 31, 36, 41, 46, 51, 56, 61, 66, 71, 76, 81, 86, 91, 96]
>>> rdd1.getNumPartitions()
2
>>> 
```

```
sc.parallelize(range(1,100,2)).collect()
```

2 here, represents the number of partitions the RDD is split into which in turn Runs them parallelly

edureka

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

edureka!

Viewing Partitions (Job History Server)

Go to Job History Server and in Search bar enter your app id

History Server

Event log directory: hdfs://nameservice1/user/spark/applicationHistory
Last updated: 6/21/2018, 1:04:28 PM

Show	20	entries	Search:	edureka_29	
App ID	App Name	Started	Spark User	Last Updated	Event Log
application_1528714825862_3904	PySparkShell	2018-06-21 06:46:41	edureka_294428	2018-06-21 06:46:46	Download
application_1528714825862_3086	Spark shell	2018-06-19 13:40:48	edureka_292003	2018-06-19 13:40:47	Download
application_1528714825862_3040	Spark shell	2018-06-19 10:54:50	edureka_292003	2018-06-19 10:54:50	Download
application_1528188402551_0787	Spark shell	2018-06-06 12:55:55	edureka_292003	2018-06-06 14:02:54	Download
application_1517837014215_35537	Spark shell	2018-05-09 10:32:34	edureka_292140	2018-05-09 11:55:36	Download

Showing 1 to 5 of 5 entries (filtered from 398 total entries)

Previous 1 Next

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

← ⌛ 🏠 🔍 https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent ⌛ 🏠 🔍 ...

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Playing with Spark RDDs > Presentation

e! e!

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

Personal Library

Presentation

80 / 100 ← →

Viewing the Partitions (Job History Server)

Spark Jobs (?)

User: edureka_294428
Total Uptime:
Scheduling Mode: FIFO
Completed Jobs: 8

Event Timeline

Completed Jobs (8)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
7	collect at <stdin>:1	2018/06/21 07:06:31	0.1 s	1/1	2/2

Total 2 Partitions are done

edureka!

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

My Classroom > PySpark Certification Training Course

X PySpark Certification T... Playing with Spark RDDs > Presentation

e! e!

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count
- Personal Library

Presentation

81 / 100 ← →

Viewing Parallelism

APACHE Spark 2.1.0 Jobs Stages Storage Environment Executors Spark shell application UI

Details for Stage 0 (Attempt 0)

Total Time Across All Tasks: 17 ms
Locality Level Summary: Process local: 5

DAG Visualization Show Additional Metrics Event Timeline Enable zooming

Scheduler Delay Task Deserialization Time Shuffle Read Time Executor Computing Time Shuffle Write Time Getting Result Time Result Serialization Time

Parallel Execution of 2 completed tasks

1 / ip-20-0-31-221.ec2.internal 900 950 000 050 100 150 200 250 300 260 270 280 290 300 310 320 330 340 350 360 370 380 390 400 410 420 430 440 450 460 470 480 490 500 510 520 530 540 550 560 570 580 590 600 610 620 630 640 650 660 670 680 690 700 710 720 730 740 750 760 770 780 790 800 810 820 830 840 850 860 870 880 890 900 910 920 930 940 950 960 970 980 990 1000 13:41:42 13:41:43

Summary Metrics for 5 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	1 ms	1 ms	1 ms	2 ms	12 ms
GC Time	0 ms	0 ms	0 ms	0 ms	0 ms

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

31°C Haze Search

ENG IN 01-06-2023 23:16

https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

Presentation

82 / 100

Passing Function in Spark

edureka!

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

31°C Haze

Search

L

ENG IN

23:16

01-06-2023

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

Presentation

83 / 100



Function Passing

- Function Passing and Object Sharing is an inherent part of any standalone Software Application
- When you pass functions to another class or function, you implicitly provide access of the function output and members to the calling function or class



edureka!

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

https://learning.edureka.co/classroom/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

edureka!

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

Function Passing

- In Spark, most of the Transformation operations and a considerable number of Action Operations depend on the Function Passing
- Function Passing performs an incomprehensible role in performing Data Computation and defining Program Flow

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Presentation 85 / 100 ← →

Passing Functions to Spark

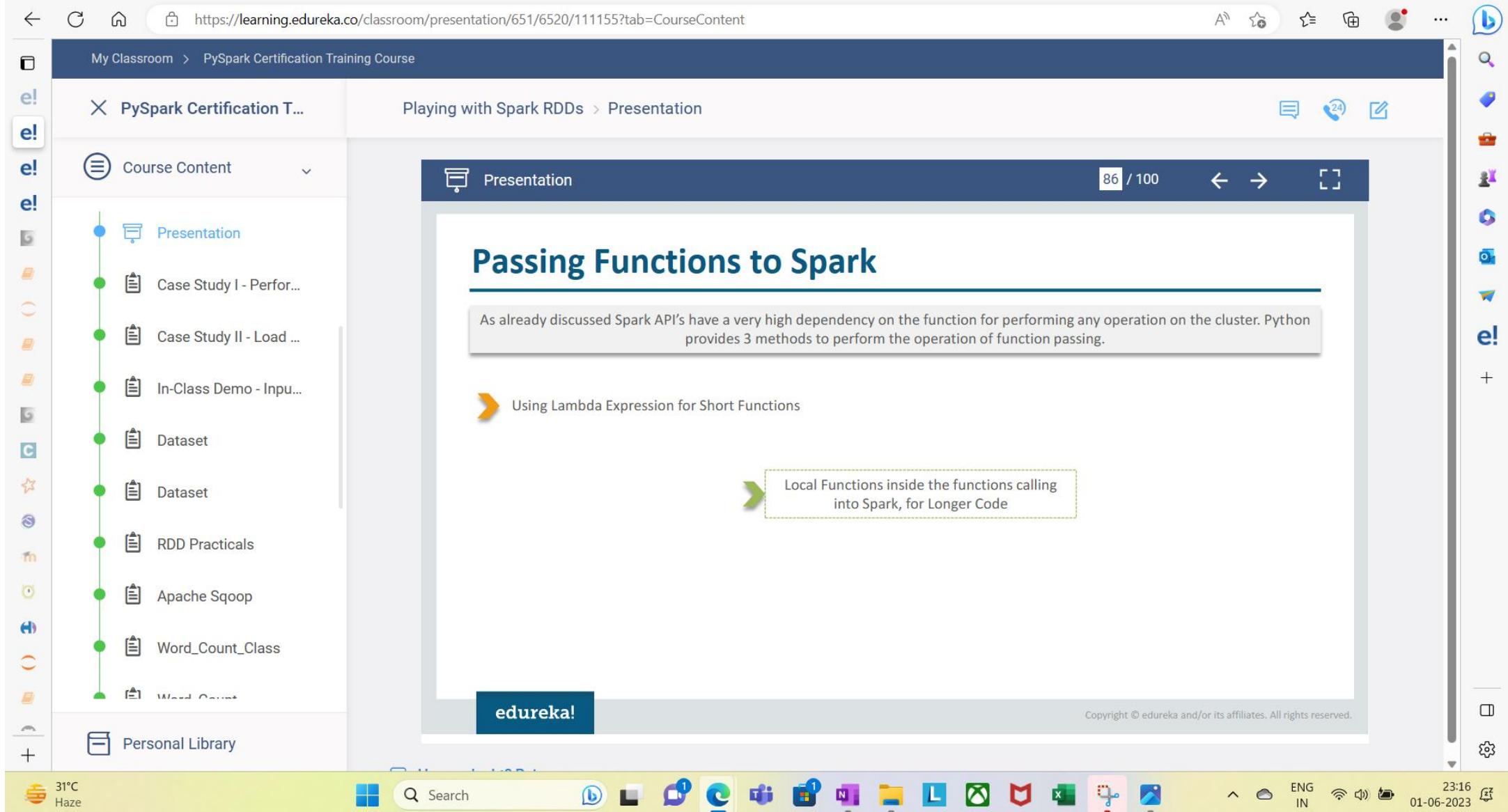
As already discussed Spark API's have a very high dependency on the function for performing any operation on the cluster. Python provides 3 methods to perform the operation of function passing.

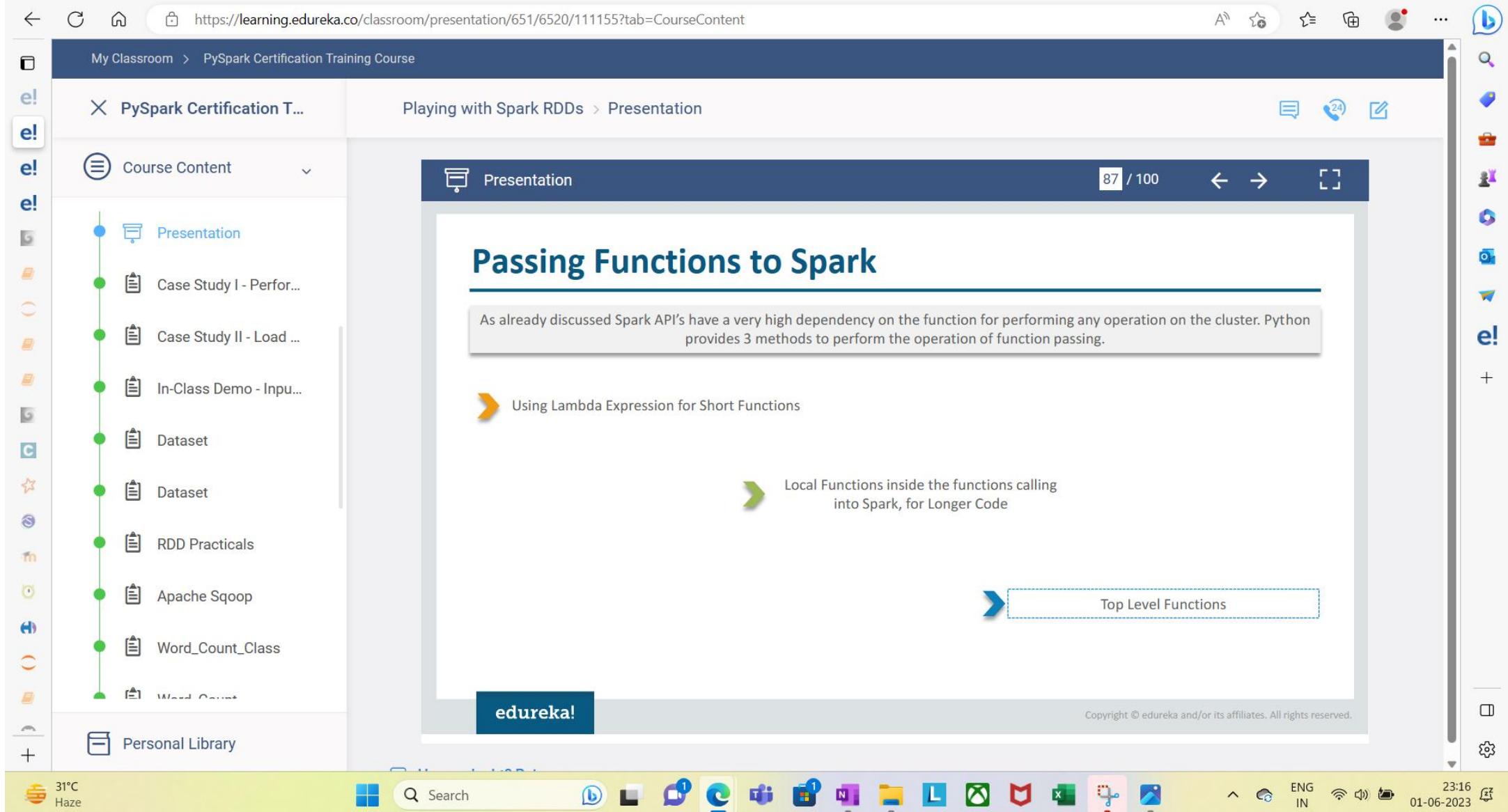
Using Lambda Expression for Short Functions

edureka!

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

Personal Library





My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

Personal Library

Presentation

Using Lambda Expression for Short Function

```
lines = sc.textFile("README.md")
pythonLines = lines.filter(lambda line: "Python" in line)
pythonLines.first()
```

edureka!

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

1. Lambda expressions for simple functions that can be written as an expression
2. Lambdas do not support multi-statement functions or statements that do not return a value

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class

edureka!

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

Using Local Functions

```
def containsError(s):  
    return "error" in s  
word = rdd.filter(containsError)
```

1. We pass the locally created functions inside the pre-defined functions as parameters to create RDDs
2. One of the major characteristics of local functions is that they are not interpreted with the other Top-Level Functions
3. These functions are at non-zero indentation level

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

Personal Library

Presentation

90 / 100



Using Top Level Functions

```
def __init__(self):  
    self.field = "Hello"  
  
class MyClass(object):  
  
    def doStuff(self, rdd):  
        return rdd.map(lambda s: self.field + s)
```

1. We can very clearly see that both the function and the class are at same indentation level and therefore will be executed at the same time.

edureka!

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

Note : Passing Function with Field Reference

```
class SearchFunctions(object):
    def __init__(self, query):
        self.query = query
    def isMatch(self, s):
        return self.query in s
    def getMatchesFunctionReference(self, rdd):
        # Problem: references all of "self" in "self.isMatch"
        return rdd.filter(self.isMatch)
    def getMatchesMemberReference(self, rdd):
        # Problem: references all of "self" in "self.query"
        return rdd.filter(lambda x: self.query in x)
```

1. When you pass a function that is the member of an object, or contains references to fields in an object (e.g., self.field), Spark sends the entire object to worker nodes, which can be much larger than the bit of information you need.
2. Sometimes this can also cause your program to fail, if your class contains objects that Python can't figure out how to pickle.

My Classroom > PySpark Certification Training Course

Playing with Spark RDDs > Presentation



Course Content

-  Presentation
 -  Case Study I - Performance
 -  Case Study II - Load Data
 -  In-Class Demo - Input
 -  Dataset
 -  Dataset
 -  RDD Practicals
 -  Apache Sqoop
 -  Word_Count_Class

Note : Passing Function without Field Reference

```
class WordFunctions(object):
    ...
    def getMatchesNoReference(self, rdd):
        # Safe: extract only the field we need into a local variable
        query = self.query
        return rdd.filter(lambda x: query in x)
```

1. To avoid the previous issue , the simplest way is to copy field into a local variable instead of accessing it externally
 2. We just extract the fields we need from our object into a local variable and pass that in the function

https://learning.edureka.co/course/presentation/651/6520/111155?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Playing with Spark RDDs > Presentation

Course Content

- Presentation
- Case Study I - Perform...
- Case Study II - Load ...
- In-Class Demo - Input...
- Dataset
- Dataset
- RDD Practicals
- Apache Sqoop
- Word_Count_Class
- Word_Count

edureka!

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

31°C Haze

Search

L

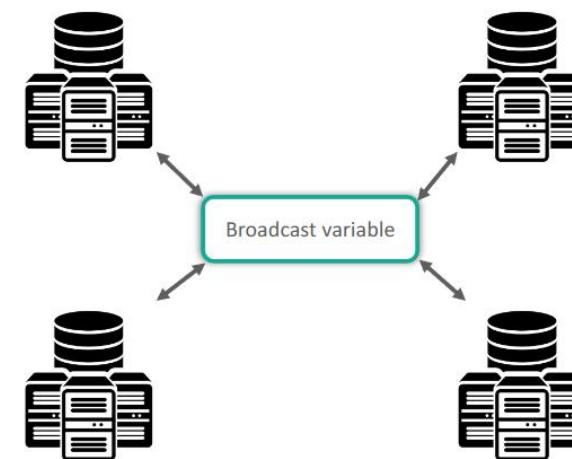
ENG IN

23:16 01-06-2023

Shared Variables – Broadcast Variable

What are Broadcast Variables?

1. Broadcast variables in Apache Spark is a mechanism for sharing variables across executors that are meant to be read-only.
2. Without broadcast variables these variables would be shipped to each executor for every transformation and action, and this can cause network overhead.
3. However, with broadcast variables, they are shipped once to all executors and are cached for future reference.



edureka!

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

