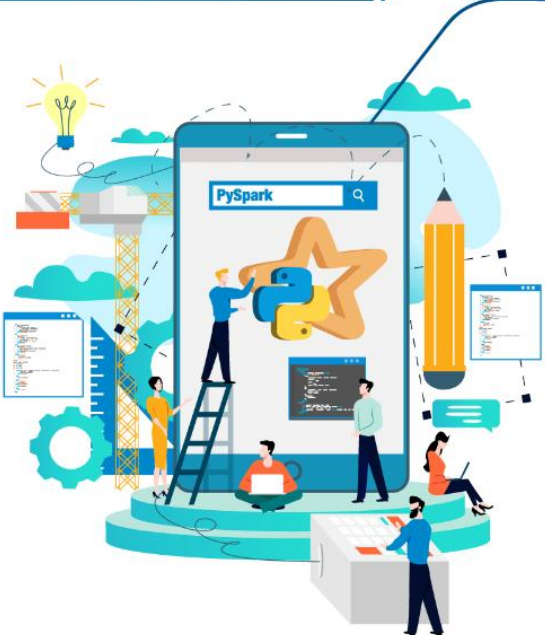


- Course Content
- Presentation
- Dataset
- Dataset
- Code
- Code
- In-Class Project-I Sol...
- In-Class Project-II So...
- Spark GraphX
- Certification Project - ...

Presentation 3 / 20

COURSE OUTLINE



MODULE 12

- Introduction to Big Data Hadoop and Spark
- Introduction to Python for Apache Spark
- Functions, OOPs, Modules in Python
- Deep Dive into Apache Spark Framework
- Playing with Spark RDDs
- Data Frames and Spark SQL
- Machine Learning using Spark MLlib
- Deep Dive into Spark MLlib
- Understanding Apache Kafka and Apache Flume
- Apache Spark Streaming – Processing Multiple Batches
- Apache Spark Streaming - Data Sources
- Implementing an End-to-End Project**

Objectives

After completing this module, you should be able to:

- Create Spark based end-to-end Solutions for Domain Specific data using core Spark functionalities



edureka!

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

Consideration

- You are going to appear for an Interview at World's largest Film and Media Production firm



edureka!

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

Background Information

- You are going to appear for an Interview at India's largest Film and Media Production firm
- Since you have appeared as a candidate for Data Analyst profile you have to provide the firm with insights on the data that they provide

insights on the data that they provide

edureka!

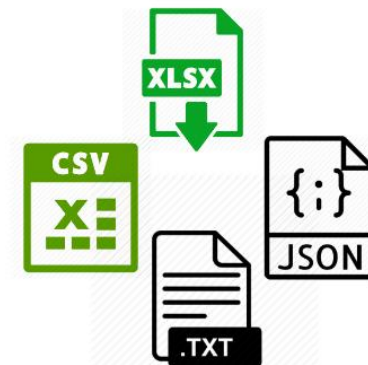
Copyright © 2018, edureka and/or its affiliates. All rights reserved.

Dataset Overview

- You are going to appear for an Interview at India's largest Film and Media Production firm
- Since you have appeared as a candidate for Data Analyst profile you have to provide the firm with insights on the data that they provide
- In the next slide we will discuss about the dataset and the insights that you have to provide

and the insights that you have to provide

- In the next slide we will discuss about the dataset

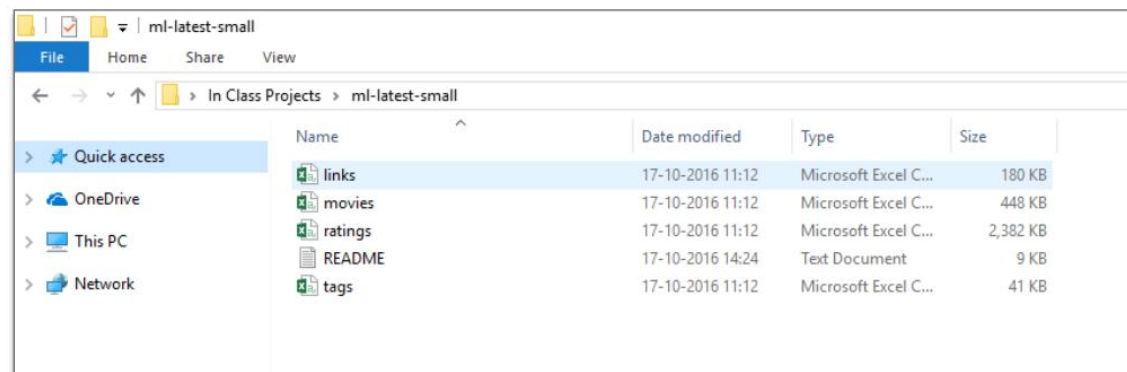


edureka!

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

Dataset

The given dataset has various files that you can download from [here](#). You will only need the [ml-latest-small.zip](#) and [ml-100k.zip](#) files. Once you have downloaded the files your **ml-latest-small** download will contain some files like this:



Name	Date modified	Type	Size
links	17-10-2016 11:12	Microsoft Excel C...	180 KB
movies	17-10-2016 11:12	Microsoft Excel C...	448 KB
ratings	17-10-2016 11:12	Microsoft Excel C...	2,382 KB
README	17-10-2016 14:24	Text Document	9 KB
tags	17-10-2016 11:12	Microsoft Excel C...	41 KB

Course Classroom | Eureka

←

↻

🏠

🔒

https://learning.edureka.co/classroom/presentation/651/6527/111478?tab=CourseContent

🔊

🌟

🌟

🔖

👤

⋮

💬

My Classroom > PySpark Certification Training Course

PySpark Certification T...

In-class Project > Presentation

💬

📞24

✍️

Course Content

📄 Presentation

📄 Dataset

📄 Dataset

📄 Code

📄 Code

📄 In-Class Project-I Sol...

📄 In-Class Project-II So...

📄 Spark GraphX

📄 Certification Project - ...

Personal Library

Presentation

10 / 20

⏪ ⏩

🖼️

Dataset

This is how your **ml-100k** folder will look after it is downloaded. You may notice that previous folder contained CSV files while this has files with .data, item etc extensions. These files are pipe separated files and contain movie records of older releases.

ml-100k

File Home Share View

← → ↶ ↷ ↵

In Class Projects > ml-100k

Name	Date modified	Type	Size
allbut.pl	19-07-2000 16:09	PL File	1 KB
mku.sh	19-07-2000 16:09	SH File	1 KB
README	29-01-2016 14:26	File	7 KB
u.data	19-07-2000 16:09	DATA File	1,933 KB
u.genre	19-07-2000 16:09	GENRE File	1 KB
u.info	19-07-2000 16:09	INFO File	1 KB
u.item	19-07-2000 16:09	ITEM File	231 KB
u.occupation	19-07-2000 16:09	OCCUPATION File	1 KB
u.user	19-07-2000 16:09	VisualStudio.user....	23 KB
u1.base	08-03-2001 12:33	BASE File	1,550 KB
u1.test	08-03-2001 12:32	TEST File	384 KB

edureka!

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

38°C Sunny

🪟 🔍 Search

💬

🖼️

📞1

🌐

👤

📅1

📧N

📁

📄L

🎮

📖

📊

🗺️

🔊

🔌

📶

📶

🔊

🔌

📶

🔌

🔌

ENG IN 12:52 02-06-2023

Problem Statement

Now, below mentioned is the problem statement along its specification. We will perform all necessary operations and find out the required result:

Top 10 Most Popular movies on the following dataset:

1. ml-100k using **Spark RDD's**
2. ml-small-latest using **Spark Data Frames**

Consideration

- A leading financial bank is trying to broaden the financial inclusion for the unbanked population by providing a positive and safe borrowing experience



edureka!

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

Background Information

- A leading financial bank is trying to broaden the financial inclusion for the unbanked population by providing a positive and safe borrowing experience
- In order to make sure this underserved population has a positive loan experience, it makes use of a variety of alternative data--including telco and transactional information--to predict their clients' repayment abilities



Problem Statement

Now, below mentioned is the problem statement along its specification. We will perform all necessary operations and find out the required result:

The bank has asked you to develop a solution to ensure that clients capable of repayment are not rejected and that loans are given with a principal, maturity, and repayment calendar that will empower their clients to be successful

Problem Statement

Note : You can download the dataset from the additional files folder on LMS.

Summary

End-to-End In-Class Project 1

End-to-End In-Class Project 2

Loan Repayment Prediction and Analysis

edureka!

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