

https://learning.edureka.co/classroom/presentation/651/6526/111474?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming – Data Sources > Presentation

Course Content

- Apache Spark Streami...
- Presentation
- Case Study I - Create...
- Case Study II - Pipeli...
- Dataset
- Dataset
- In-Class Demo
- Class 10 Recording
- In-class Project
- Spark GraphX
- Personal Library

Presentation

6 / 29

Objectives

After completing this module, you should be able to:

- Learn different streaming data sources such as Kafka, Twitter and Flume

edureka!

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

38°C Sunny

Search

Windows Start Menu

Icons for various Microsoft applications: Edge, Mail, Photos, OneDrive, OneNote, Word, Excel, Powerpoint, Teams, Planner, Lists, XBOX, etc.

12:49 02-06-2023



Course Content

Apache Spark Streami...

Presentation

Case Study I - Create...

Case Study II - Pipeli...

Dataset

Dataset

In-Class Demo

Class 10 Recording

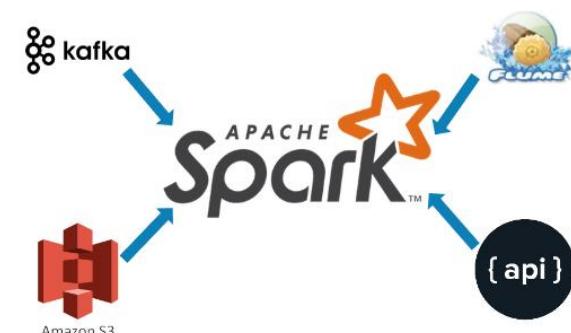
In-class Project

Spark GraphX

Personal Library

Introduction to Data Sources

The Spark engine has the ability to access data from a variety of sources, including HDFS, Amazon S3, Kafka, Flume, Third Party API's. This data can be received as streaming input or can be read from a stored location. We can further utilize this data to perform analysis and other relevant operations.



edureka!

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

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Apache Spark Streaming – Data Sources > Presentation



Course Content

Apache Spark Streami...

Presentation

Case Study I - Create...

Case Study II - Pipeli...

Dataset

Dataset

In-Class Demo

Class 10 Recording

In-class Project

Spark GraphX

Personal Library

Types of Data Sources

You can also categorize the various kind of data sources based on the kind of data they push into Spark engine. Data Sources can be majorly divided in two categories:

Streaming Data Sources that continuously push data inside the Spark Engine.



Streaming Data Sources

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

edureka!

X PySpark Certification T...

Apache Spark Streaming – Data Sources > Presentation



e! Course Content

Apache Spark Streami...

Presentation

Case Study I - Create...

Case Study II - Pipeli...

Dataset

Dataset

In-Class Demo

Class 10 Recording

In-class Project

Spark GraphX

Personal Library

Presentation

9 / 29



Types of Data Sources

You can also categorize the various kind of data sources based on the kind of data they push into Spark engine. Data Sources can be majorly divided in two categories:



Static Data Sources that store data and Spark Engine fetches from this data when needed.



Static Data Sources

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

edureka!

https://learning.edureka.co/classroom/presentation/651/6526/111474?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming – Data Sources > Presentation

Course Content

- Apache Spark Streami...
- Presentation
- Case Study I - Create...
- Case Study II - Pipeli...
- Dataset
- Dataset
- In-Class Demo
- Class 10 Recording
- In-class Project
- Spark GraphX
- Personal Library

Presentation

10 / 29

Demo – Apache Kafka and Spark Streaming Integration

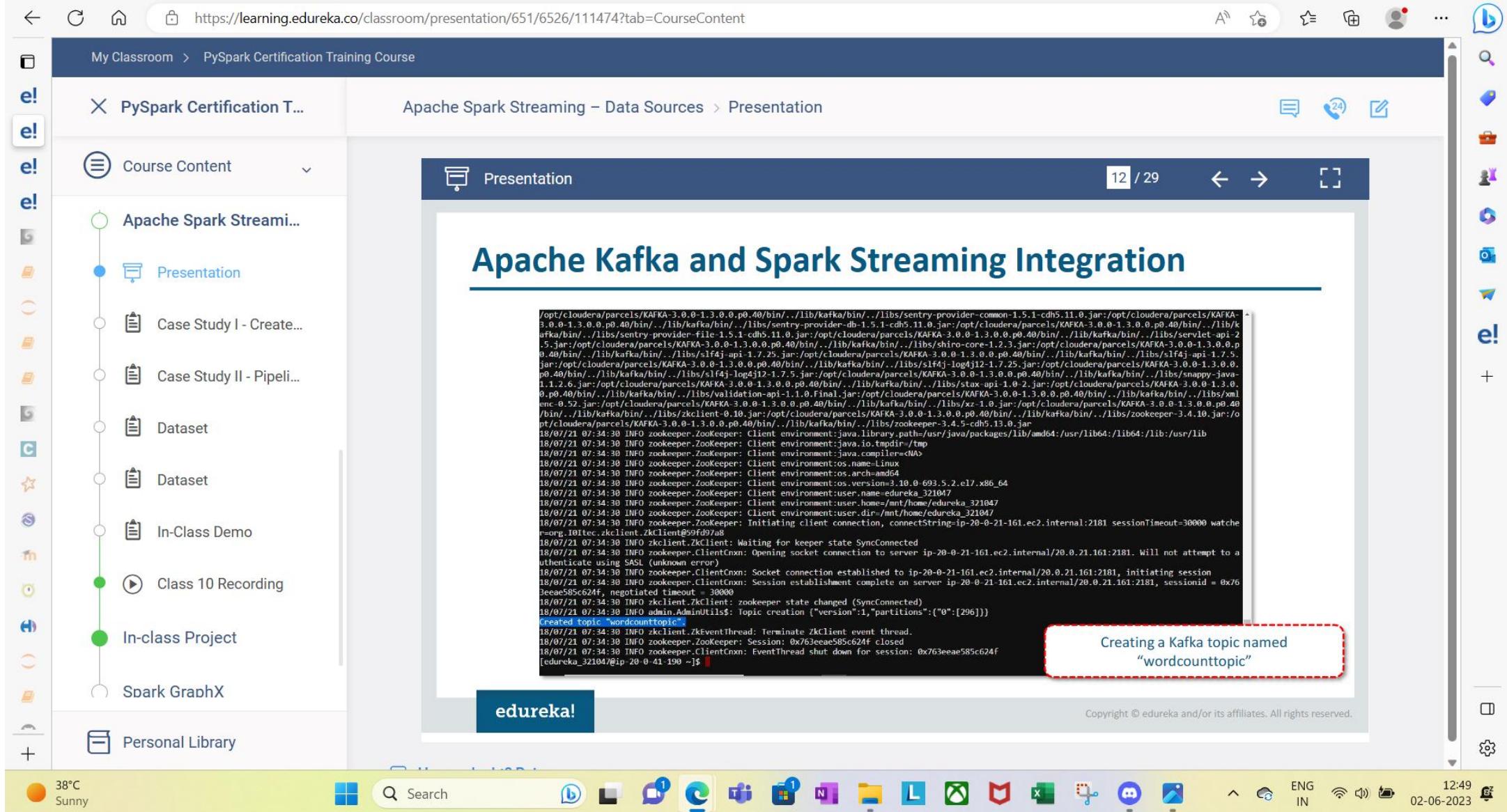
edureka!

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

38°C Sunny

Search

12:49 02-06-2023



X PySpark Certification T...

Apache Spark Streaming – Data Sources > Presentation



Course Content

Apache Spark Streami...

Presentation

Case Study I - Create...

Case Study II - Pipeli...

Dataset

Dataset

In-Class Demo

Class 10 Recording

In-class Project

Spark GraphX

Personal Library

Presentation

13 / 29



Apache Kafka and Spark Streaming Integration

```
1  from __future__ import print_function
2  import sys
3  from pyspark import SparkContext
4  from pyspark.streaming import StreamingContext
5  from pyspark.streaming.kafka import KafkaUtils
6
7  if __name__ == "__main__":
8      if len(sys.argv) != 3:
9          sc = SparkContext(appName="PythonStreamingKafkaWordCount")
10         ssc = StreamingContext(sc, 10)
11
12         zkQuorum, topic = sys.argv[1:3]
13         kvs = KafkaUtils.createStream(ssc, zkQuorum, "spark-streaming-consumer", {topic: 1})
14         lines = kvs.map(lambda x: x[1])
15         counts = lines.flatMap(lambda line: line.split(" ")).map(lambda word: (word, 1)).reduceByKey(lambda a, b: a+b)
16         counts.pprint()
17
18         ssc.start()
19         ssc.awaitTermination()
```

Create a Kafka wordcount Python program

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

edureka!

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming – Data Sources > Presentation

Course Content

Apache Spark Streami... Presentation

Case Study I - Create... Case Study II - Pipeli...

Dataset Dataset

In-Class Demo

Class 10 Recording In-class Project

Spark GraphX Personal Library

Presentation

14 / 29 ← →

Apache Kafka and Spark Streaming Integration

```
[edureka_321047@ip-20-0-41-199 ~]$ spark2-submit --master yarn kafka wordcount.py ip-20-0-21-161.ec2.internal:2181 wordcounttopic
18/07/21 09:46:20 INFO spark.SparkContext: Running Spark version 2.1.0.cloudera2
18/07/21 09:46:21 INFO spark.SecurityManager: Changing view acls to: edureka_321047
18/07/21 09:46:21 INFO spark.SecurityManager: Changing modify acls to: edureka_321047
18/07/21 09:46:21 INFO spark.SecurityManager: Changing view acls groups to:
18/07/21 09:46:21 INFO spark.SecurityManager: Changing modify acls groups to:
18/07/21 09:46:21 INFO spark.SecurityManager: authentication disabled; ui acls disabled; users with view permissions: Set(edureka_321047); groups with view permissions: Set(); users with modify permissions: Set(edureka_321047); groups with modify permissions: Set()
18/07/21 09:46:22 INFO util.Utils: Successfully started service 'sparkDriver' on port 38294.
18/07/21 09:46:22 INFO spark.SparkEnv: Registering MapOutputTracker
18/07/21 09:46:22 INFO spark.SparkEnv: Registering BlockManagerMaster
18/07/21 09:46:22 INFO storage.BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
18/07/21 09:46:22 INFO storage.BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
18/07/21 09:46:22 INFO storage.DiskBlockManager: Created local directory at /tmp/blockmgr-e3ba80c7-15df-4942-be4e-b2fff95c3b3c
18/07/21 09:46:22 INFO memory.MemoryStore: MemoryStore started with capacity 93.3 MB
18/07/21 09:46:22 INFO spark.SparkEnv: Registering OutputCommitCoordinator
18/07/21 09:46:23 INFO yarn.Client: Requesting a new application from cluster with 3 NodeManagers
18/07/21 09:46:23 INFO yarn.Client: Verifying our application has not requested more than the maximum memory capability of the cluster (4096 MB per container)
18/07/21 09:46:23 INFO yarn.Client: Will allocate AM container, with 896 MB memory including 384 MB overhead
18/07/21 09:46:23 INFO yarn.Client: Setting up container launch context for our AM
18/07/21 09:46:23 INFO yarn.Client: Preparing resources for our AM container
18/07/21 09:46:25 INFO yarn.Client: Uploading resource file:/tmp/spark-2116b936-698d-4a33-8db0-e9a1cc921637/_spark_conf_3308945078087702336.zip -> hdfs://nameservice1/user/edureka_321047/.sparkStaging/application_1528714825862_16231/_spark_conf_.zip
18/07/21 09:46:25 INFO spark.SecurityManager: Changing view acls to: edureka_321047
18/07/21 09:46:25 INFO spark.SecurityManager: Changing modify acls to: edureka_321047
18/07/21 09:46:25 INFO spark.SecurityManager: Changing view acls groups to:
18/07/21 09:46:25 INFO spark.SecurityManager: Changing modify acls groups to:
```

Running the Kafka wordcount program on the shell

edureka!

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

38°C Sunny

Search

12:49 IN 02-06-2023

X PySpark Certification T...

Apache Spark Streaming – Data Sources > Presentation



Course Content

Apache Spark Streami...

Presentation

Case Study I - Create...

Case Study II - Pipeli...

Dataset

Dataset

In-Class Demo

Class 10 Recording

In-class Project

Spark GraphX

Personal Library

Presentation

15 / 29



Apache Kafka and Spark Streaming Integration

```
ip 20-0-41-190 login: edureka_321047
Password:
Last login: Sat Jul 21 09:44:10 on pts/33
[edureka_321047@ip-20-0-41-190 ~]$ kafka-console-producer --broker-list ip-20-0-31-221.ec2.internal:9092 --topic wordcounttopic
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/opt/cloudera/parcels/KAFKA-3.0.0-1.3.0.0.p0.40/lib/kafka/libs/slf4j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/opt/cloudera/parcels/KAFKA-3.0.0-1.3.0.0.p0.40/lib/kafka/libs/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
18/07/21 09:48:48 INFO producer.ProducerConfig: ProducerConfig: ProducerConfig values:
acks = 1
batch.size = 16384
bootstrap.servers = [ip-20-0-31-221.ec2.internal:9092]
buffer.memory = 33554432
client.id = console-producer
compression.type = none
connections.max.idle.ms = 540000
enable.idempotence = false
interceptor.classes = null
key.serializer = class org.apache.kafka.common.serialization.ByteArraySerializer
linger.ms = 1000
max.block.ms = 60000
max.in.flight.requests.per.connection = 5
max.request.size = 1048576
metadata.max.age.ms = 300000
metric.reporters = []
metrics.num.samples = 2
metrics.recording.level = INFO
metrics.sample.window.ms = 30000
partitioner.class = class org.apache.kafka.clients.producer.internals.DefaultPartitioner
receive.buffer.bytes = 32768
reconnect.backoff.ms = 1000
reconnect.backoff.ms = 50
request.timeout.ms = 1500
```

Running the Kafka Producer on shell for pushing the data to the topic

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

The screenshot shows a web browser window with the URL <https://learning.edureka.co/classroom/presentation/651/6526/111474?tab=CourseContent>. The page title is "Course Classroom | Edureka". The main content area is titled "Apache Spark Streaming – Data Sources > Presentation". The slide has a title "Apache Kafka and Spark Streaming Integration" and contains a large block of log output from a Kafka wordcount program. A red callout box highlights a portion of the log output: "We can see, the Kafka wordcount program displaying each word along with its frequency". The left sidebar shows a navigation tree for the "PySpark Certification Training Course". The bottom right corner of the slide footer contains the "edureka!" logo.

Course Classroom | Edureka

My Classroom > PySpark Certification Training Course

PySpark Certification T...

Course Content

Apache Spark Streami...

Presentation

Case Study I - Create...

Case Study II - Pipeli...

Dataset

Dataset

In-Class Demo

Class 10 Recording

In-class Project

Spark GraphX

Personal Library

Apache Spark Streaming – Data Sources > Presentation

17 / 29

Presentation

Apache Kafka and Spark Streaming Integration

```
18/07/21 09:51:10 INFO storage.BlockManagerInfo: Added broadcast_62_piece0 in memory on ip-20-0-31-4.ec2.internal:36951 (size: 4.1 KB, free: 366.3 M)
18/07/21 09:51:10 INFO spark.MapOutputTrackerMasterEndpoint: Asked to send map output locations for shuffle 28 to 20.0.31.4:33804
18/07/21 09:51:10 INFO spark.MapOutputTrackerMaster: Size of output statuses for shuffle 28 is 174 bytes
18/07/21 09:51:10 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 112.0 (TID 143) in 91 ms on ip-20-0-31-4.ec2.internal (executor 1) (1/1)
18/07/21 09:51:10 INFO cluster.YarnScheduler: Removed TaskSet 112.0, whose tasks have all completed, from pool
18/07/21 09:51:10 INFO scheduler.DAGScheduler: ResultStage 112 (runJob at PythonRDD.scala:441) finished in 0.091 s
18/07/21 09:51:10 INFO scheduler.DAGScheduler: Job 56 finished: runJob at PythonRDD.scala:441, took 0.293799 s
18/07/21 09:51:10 INFO spark.SparkContext: Starting job: runJob at PythonRDD.scala:441
18/07/21 09:51:10 INFO scheduler.DAGScheduler: Got job 57 (runJob at PythonRDD.scala:441) with 1 output partitions
18/07/21 09:51:10 INFO scheduler.DAGScheduler: Final stage: ResultStage 114 (runJob at PythonRDD.scala:441)
18/07/21 09:51:10 INFO scheduler.DAGScheduler: Parents of final stage: List(ShuffleMapStage 113)
18/07/21 09:51:10 INFO scheduler.DAGScheduler: Missing parents: List()
18/07/21 09:51:10 INFO scheduler.DAGScheduler: Submitting ResultStage 114 (PythonRDD[227] at RDD at PythonRDD.scala:48), which has no missing parent
...
18/07/21 09:51:10 INFO memory.MemoryStore: Block broadcast_63 stored as values in memory (estimated size 7.2 KB, free 93.1 MB)
18/07/21 09:51:10 INFO memory.MemoryStore: Block broadcast_63_piece0 stored as bytes in memory (estimated size 4.1 KB, free 93.1 MB)
18/07/21 09:51:10 INFO storage.BlockManagerInfo: Added broadcast_63_piece0 in memory on 20.0.41.190:42344 (size: 4.1 KB, free: 93.2 MB)
18/07/21 09:51:10 INFO spark.SparkContext: Created broadcast 63 from broadcast at DAGScheduler.scala:991
18/07/21 09:51:10 INFO cluster.YarnScheduler: Submitting 1 missing tasks from ResultStage 114 (PythonRDD[227] at RDD at PythonRDD.scala:48)
18/07/21 09:51:10 INFO scheduler.TaskSetManager: Adding task set 114.0 with 1 tasks
18/07/21 09:51:10 INFO scheduler.TaskSetManager: Starting task 0.0 in stage 114.0 (TID 144, ip-20-0-31-4.ec2.internal, executor 1, partition 1, NODE_LOCAL, 5678 bytes)
18/07/21 09:51:10 INFO storage.BlockManagerInfo: Added broadcast_63_piece0 in memory on ip-20-0-31-4.ec2.internal:36951 (size: 4.1 KB, free: 366.3 M)
...
18/07/21 09:51:10 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 114.0 (TID 144) in 150 ms on ip-20-0-31-4.ec2.internal (executor 1) (1/1)
...
18/07/21 09:51:10 INFO cluster.YarnScheduler: Removed TaskSet 114.0, whose tasks have all completed, from pool
18/07/21 09:51:10 INFO scheduler.DAGScheduler: ResultStage 114 (runJob at PythonRDD.scala:441) finished in 0.151 s
18/07/21 09:51:10 INFO scheduler.DAGScheduler: Job 57 finished: runJob at PythonRDD.scala:441, took 0.155179 s
```

Time: 2018-07-21 09:51:10

```
(u'Spark', 5)
(u'SparkleIlo', 1)
(u'Hello', 5)
```

We can see, the Kafka wordcount program displaying each word along with its frequency

edureka!

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

38°C Sunny

Search

12:49 02-06-2023

https://learning.edureka.co/classroom/presentation/651/6526/111474?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming – Data Sources > Presentation

Course Content

- Apache Spark Streami...
- Presentation
- Case Study I - Create...
- Case Study II - Pipeli...
- Dataset
- Dataset
- In-Class Demo
- Class 10 Recording
- In-class Project
- Spark GraphX
- Personal Library

Presentation

18 / 29

Demo – Flume and Spark Streaming Integration

edureka!

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

38°C Sunny

Search

12:49 02-06-2023

https://learning.edureka.co/classroom/presentation/651/6526/111474?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming – Data Sources > Presentation

Course Content

- Apache Spark Streami...
- Presentation
- Case Study I - Create...
- Case Study II - Pipeli...
- Dataset
- Dataset
- In-Class Demo
- Class 10 Recording
- In-class Project
- Spark GraphX
- Personal Library

19 / 29 ← →

Flume and Spark Streaming Integration

We are working on a solution to monitor any application by analyzing the log files in Apache Spark through streaming module. We will build a flume and spark pipeline to parse the logs.



edureka!

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

38°C Sunny

Search

12:49 02-06-2023

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming – Data Sources > Presentation

Presentation 20 / 29 ← →

Flume and Spark Streaming Integration

```
1 # Name the components on this agent
2 stlog.sources = r1
3 stlog.sinks = k1
4 stlog.channels = c1
5
6 # Describe/configure the source
7 stlog.sources.r1.type = exec
8 stlog.sources.r1.command = tail -F -s 2 /mnt/home/edureka_321047/av/workspace/logs/event2.log
9
10 # Describe the sink
11 stlog.sinks.k1.type = hdfs
12 stlog.sinks.k1.hdfs.path = use_cases/streaming/events/%Y-%m-%d/
13 stlog.sinks.k1.hdfs.filePrefix = events-
14 stlog.sinks.k1.hdfs.fileSuffix = .log
15 stlog.sinks.k1.hdfs.useLocalTimeStamp = true
16 stlog.sinks.k1.hdfs fileType = DataStream
17
18 # Use a channel stlogich buffers events in memory
19 stlog.channels.c1.type = memory
20 stlog.channels.c1.capacity = 1000
21 stlog.channels.c1.transactionCapacity = 100
22
23 # Bind the source and sink to the channel
24 stlog.sources.r1.channels = c1
25 stlog.sinks.k1.channel = c1
26
```

We will first create a config file for flume by the name **stlog.conf** and upload it on FTP. This file has been uploaded in the Additional files folder on LMS.

edureka!

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

38°C Sunny

Search

12:49 02-06-2023

https://learning.edureka.co/classroom/presentation/651/6526/111474?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming – Data Sources > Presentation

Course Content

Apache Spark Streami... Presentation

Case Study I - Create... Case Study II - Pipeli... Dataset Dataset In-Class Demo Class 10 Recording In-class Project Spark GraphX Personal Library

Presentation

Flume and Spark Streaming Integration

21 / 29 ← →

```
from __future__ import print_function
import sys
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
import datetime
import re
from pyspark.sql import Row, SparkSession
from pyspark.sql.functions import regexp_extract

if __name__ == "__main__":
    sc = SparkContext(appName="LogParser-py-Streaming")
    ssc = StreamingContext(sc, 10)
    now = datetime.datetime.now()

    # define the regular expression - <date-time> <log-type> <message>
    regex = re.compile("(^([\\d/]+ [\\d:]+) ([a-zA-Z]+) (.*)")")

    filepath = "/user/edureka_321047/use_cases/streaming/events/" + now.strftime("%Y-%m-%d/")
    print("filepath:", filepath)
    lines = ssc.textFileStream(filepath)

    def getSparkSessionInstance(sparkConf):
        if ("sparkSessionSingletonInstance" not in globals()):
            globals()["sparkSessionSingletonInstance"] = SparkSession \
                .builder \
                .config(conf=sparkConf) \
                .getOrCreate()
        return globals()["sparkSessionSingletonInstance"]

    def process(t, rdd):
        rdd.foreach(lambda line: print(line))
```

We will now create a Streaming Application by the name **logparser.py** and upload it on FTP. This file has been uploaded in the Additional files folder on LMS.

edureka!

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

38°C Sunny

Search

12:49 IN 02-06-2023

e! PySpark Certification T...

e! Course Content

e! Apache Spark Streami...

e! Presentation

e! Case Study I - Create...

e! Case Study II - Pipeli...

e! Dataset

e! Dataset

e! In-Class Demo

e! Class 10 Recording

e! In-class Project

e! Spark GraphX

e! Personal Library

Flume and Spark Streaming Integration

```
[edureka_321047@ip-20-0-41-202 ~]$ spark2-submit logparser_streaming.py --deploy-mode client > /mnt/home/edureka_321047/av/workspace/logs/streaming_output.log 2>&1
```

This command contains the path of the location where your streaming output will be stored. Here the name of the file is **streaming_output.log**.

Run this command on console to start the flume agent :

```
spark2-submit logparser_streaming.py --deploy-mode client > /mnt/home/edureka_321047/av/workspace/logs/streaming_output.log 2>&1
```

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

X PySpark Certification T...

Apache Spark Streaming – Data Sources > Presentation



Course Content

Apache Spark Streami...

Presentation

Case Study I - Create...

Case Study II - Pipeli...

Dataset

Dataset

In-Class Demo

Class 10 Recording

In-class Project

Spark GraphX

Personal Library

Presentation

24 / 29



Flume and Spark Streaming Integration

```
1 18/09/06 08:08:18 INFO spark.SparkContext: Running Spark version 2.1.0.cloudera2
2 18/09/06 08:08:19 INFO spark.SecurityManager: Changing view acls to: edureka_321047
3 18/09/06 08:08:19 INFO spark.SecurityManager: Changing modify acls to: edureka_321047
4 18/09/06 08:08:19 INFO spark.SecurityManager: Changing view acls groups to:
5 18/09/06 08:08:19 INFO spark.SecurityManager: Changing modify acls groups to:
6 18/09/06 08:08:19 INFO spark.SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view
7 18/09/06 08:08:19 INFO util.Utils: Successfully started service 'sparkDriver' on port 45132.
8 18/09/06 08:08:19 INFO spark.SparkEnv: Registering MapoutputTracker
9 18/09/06 08:08:20 INFO spark.SparkEnv: Registering BlockManagerMaster
10 18/09/06 08:08:20 INFO storage.BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting
11 18/09/06 08:08:20 INFO storage.BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
12 18/09/06 08:08:20 INFO storage.DiskBlockManager: Created local directory at /tmp/blockmgr-3257952b-c8e6-48eb-ab15-1a50dac6:
13 18/09/06 08:08:20 INFO memory.MemoryStore: MemoryStore started with capacity 93.3 MB
14 18/09/06 08:08:20 INFO spark.SparkEnv: Registering OutputCommitCoordinator
15 18/09/06 08:08:21 INFO yarn.Client: Requesting a new application from cluster with 4 NodeManagers
16 18/09/06 08:08:21 INFO yarn.Client: Verifying our application has not requested more than the maximum memory capability of
17 18/09/06 08:08:21 INFO yarn.Client: Will allocate AM container, with 896 MB memory including 384 MB overhead
18 18/09/06 08:08:21 INFO yarn.Client: Setting up container launch context for our AM
19 18/09/06 08:08:21 INFO yarn.Client: Setting up the launch environment for our AM container
20 18/09/06 08:08:21 INFO yarn.Client: Preparing resources for our AM container
21 18/09/06 08:08:23 INFO yarn.Client: Uploading resource file:/tmp/spark-fa2e848b-9f99-41a3-884a-05b0bc166272/_spark_conf_
22 18/09/06 08:08:23 INFO spark.SecurityManager: Changing view acls to: edureka_321047
23 18/09/06 08:08:23 INFO spark.SecurityManager: Changing modify acls to: edureka_321047
24 18/09/06 08:08:23 INFO spark.SecurityManager: Changing view acls groups to:
25 18/09/06 08:08:23 INFO spark.SecurityManager: Changing modify acls groups to:
26 18/09/06 08:08:23 INFO spark.SecurityManager: SecurityManager: authentication di
27 18/09/06 08:08:23 INFO yarn.Client: Submitting application application_152871482
28 18/09/06 08:08:23 INFO impl.YarnClientImpl: Submitted application application_152871482
29 18/09/06 08:08:23 INFO cluster.SchedulerExtensionServices: Starting Yarn extensi
30 18/09/06 08:08:24 INFO yarn.Client: Application report for application_152871482
31 18/09/06 08:08:24 INFO yarn.Client:
```

You can go to the location of streaming_output.log file. Once downloaded you can see the values stored inside the file.

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

https://learning.edureka.co/classroom/presentation/651/6526/111474?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming – Data Sources > Presentation

Course Content

- Apache Spark Streami...
- Presentation
- Case Study I - Create...
- Case Study II - Pipeli...
- Dataset
- Dataset
- In-Class Demo
- Class 10 Recording
- In-class Project
- Spark GraphX
- Personal Library

Presentation

25 / 29

Twitter Sentiment Analysis

edureka!

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

38°C Sunny

Search

12:49 02-06-2023

https://learning.edureka.co/course/presentation/651/6526/111474?tab=CourseContent

My Classroom > PySpark Certification Training Course

PySpark Certification T... Apache Spark Streaming – Data Sources > Presentation

Course Content

- Apache Spark Streami...
- Presentation
- Case Study I - Create...
- Case Study II - Pipeli...
- Dataset
- Dataset
- In-Class Demo
- Class 10 Recording
- In-class Project
- Spark GraphX
- Personal Library

Presentation

Twitter Sentiment Analysis

- You can find the demo doc for this file on LMS

edureka!

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

38°C Sunny

Search

12:49 02-06-2023

My Classroom > PySpark Certification Training Course

X PySpark Certification T...

Apache Spark Streaming – Data Sources > Presentation

...

Course Content

Apache Spark Streami...

Presentation

Case Study I - Create...

Case Study II - Pipeli...

Dataset

Dataset

In-Class Demo

Class 10 Recording

In-class Project

Spark GraphX

Personal Library

Presentation

27 / 29

← →

[]

Questions

edureka!

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