

Week 10 Big data integration and processing						
09	September 2016					
Wk	M	T	W	T	F	S
37	1	2	3	4		
38	5	6	7	8	9	10
39	12	13	14	15	16	17
40	19	20	21	22	23	24
	26	27	28	29	30	
10	October 2016					
Wk	M	T	W	T	F	S
40/45	31					
41	1	2	3	4	5	6
42	10	11	12	13	14	15
43	17	18	19	20	21	22
44	24	25	26	27	28	29

*Components of Spark

Spark core is the underlying general execution engine for August Wednesday

spark platform that all other functionality

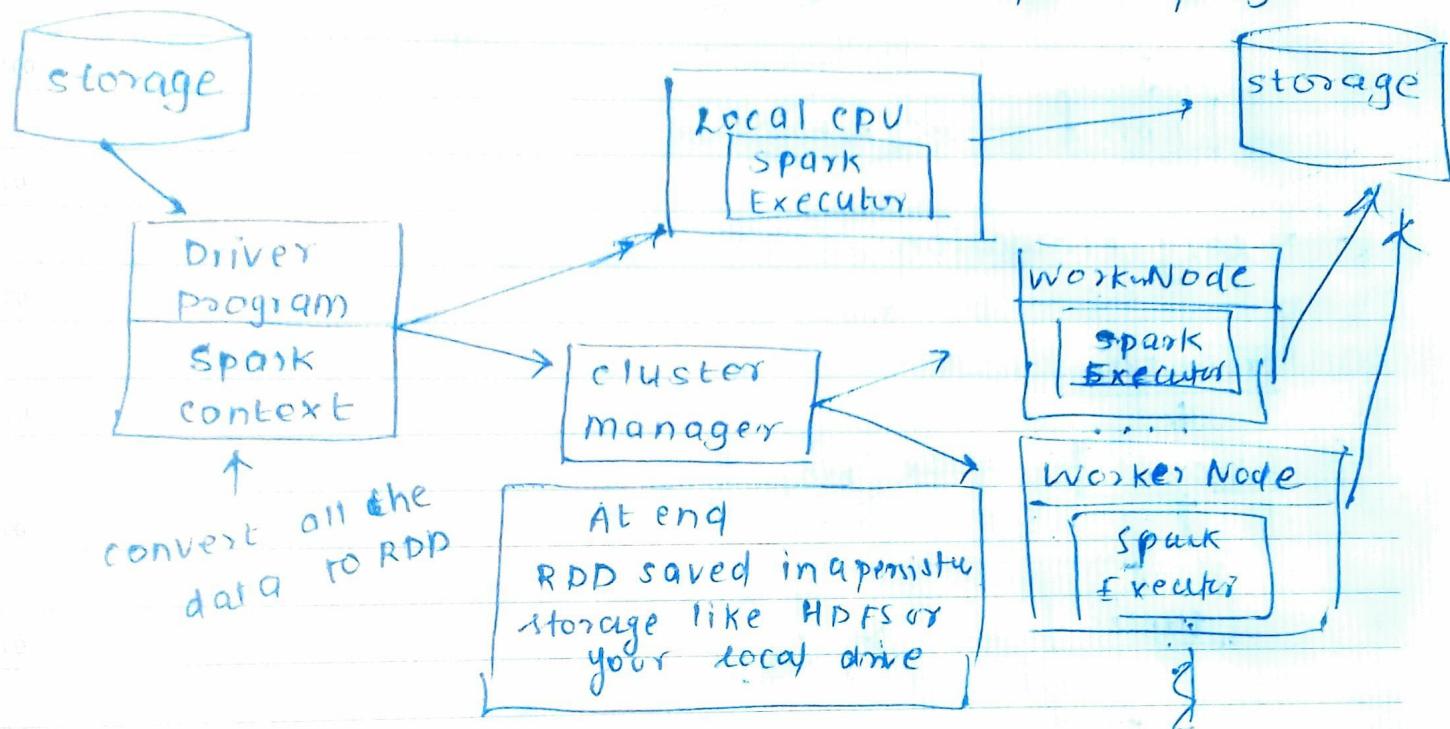
23Q-136 WK 34

17

Video 1:-

Spark Core:- Programming in Spark is built upon.

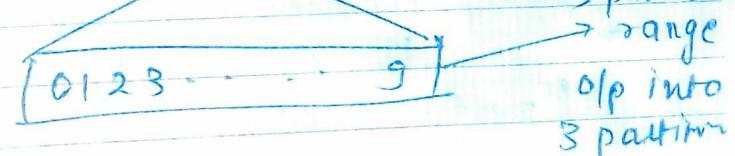
- two methods to create RDDs in Spark
- Explain with immutable means
- Interpret a spark program as a pipeline of transformations and actions
- Lists the steps to create a spark program



lines = sc.textFile("hdfs://user/cloudera/words.txt")
 lines = sc.parallelize(["big", "data"])

using parallelize

numbers = sc.parallelize(range(10), 3)



numbers.collect()

[0, 1, 2], [3, 4, 5], [6, 7, 8, 9]

18

and referencing
datasets in external
storage systems.

Thursday

August

231-135 • WK 34

07

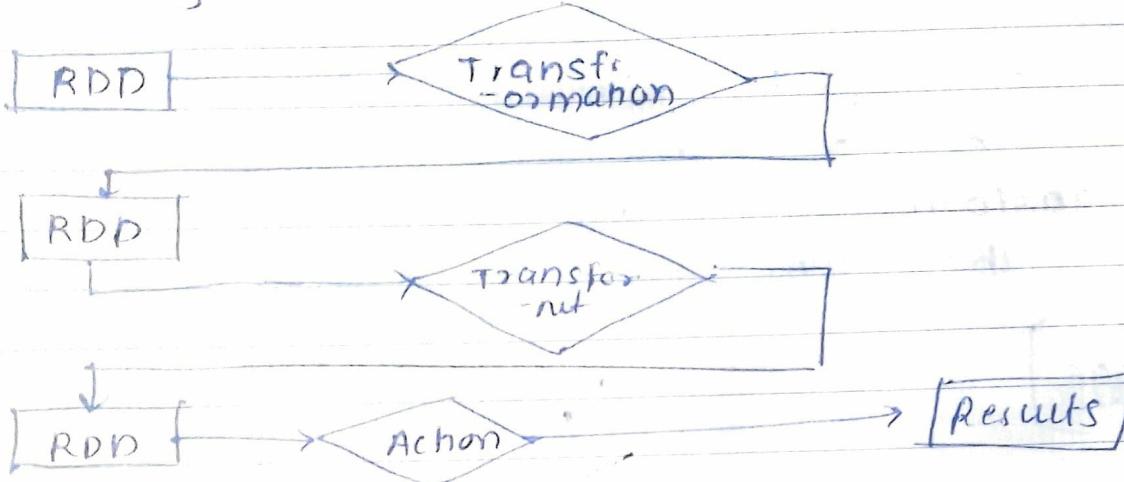
	M	T	W	T	F	S	S
28	4	5	6	7	8	9	10
29	11	12	13	14	15	16	17
30	18	19	20	21	22	23	24
31	25	26	27	28	29	30	31

July 2016

08

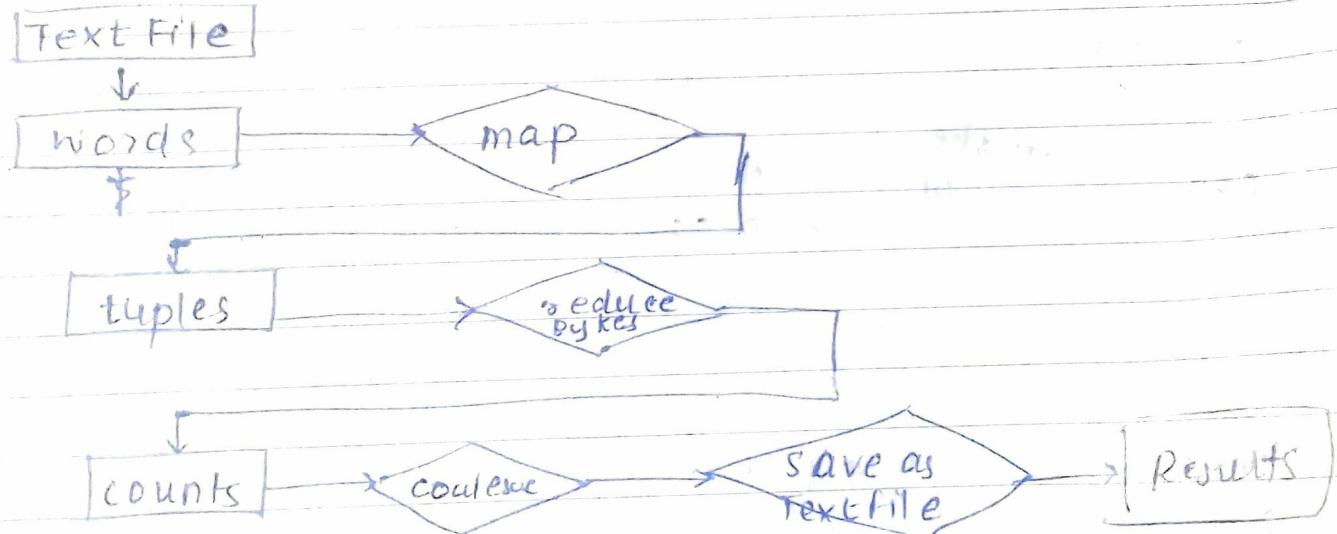
wk	M	T	W	T	F	S	S
32	1	2	3	4	5	6	7
33	8	9	10	11	12	13	14
34	15	16	17	18	19	20	21
35	22	23	24	25	26	27	28
36	29	30	31				

Processing RDDs



→ lazy evaluation

Word Count



Create RDDs

Apply transformations

perform actions

NOTES

February 2016						
S	M	T	W	F	S	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29						

03

wk	M	T	W	T	F	S	S
10	1	2	3	4	5	6	
11	7	8	9	10	11	12	13
12	14	15	16	17	18	19	20
13	21	22	23	24	25	26	27
14	28	29	30	31			

March 2016

Weeks

Video 2:-

Big data

2016

Integration of Thursday

January

Processing

21

021-345 • WK 04

Spark Core :- Transformations

→ RDD gets generated from external datasets and get partitioned.

RDD = immutable ⇒ they can't be changed in place even partially.

RDD → transformation → New RDD.

They need a transformation operation applied to them and get

* This is essential for keeping track of all process that has been applied to our dataset providing the ability to keep a linear chain of RDDs.

→ In pipeline:- several transformation steps, many other RDDs as intermediate products get executed until we got to our final result.

SPARK → transformation are lazy.

It will not execute immediately

→ Narrow Transmission vs Wide Transmission

Map, flatmap, filter and coalesce

⇒ Narrow Transmission.

NOTES

Home network ⇒ It applies the function to each partition or element of an RDD

⇒ one to one transformation

⇒ element wise transformation

* In some cases, even if you started with even partitions within the worker nodes, the RDD size significantly vary across the workers after a filter operation, then this happens is a good idea to join some of partitions to increase performance and even out processing across clusters. This transformation is called coalesce.

→ It helps with balancing the data partition numbers and sizes.

100

500

600

700

OTES

September 2016						
S	M	T	W	T	F	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30					

October 2016						
Wk	M	T	W	T	F	S
40/45	31			1	2	
				3	4	5
				6	7	8
				9		
	41			10	11	12
				13	14	15
				16	17	18
	42			19	20	21
				22	23	24
	43			25	26	27
				28	29	30
	44					

2016
Wednesday
August

17

230-136 • WK 34

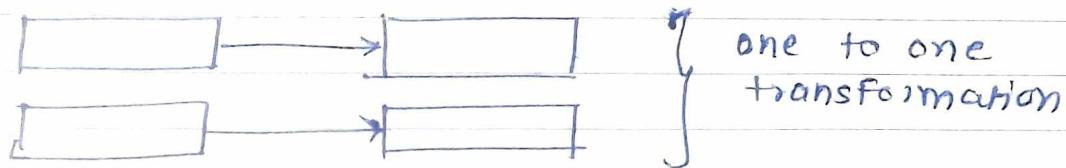
Transformations :- RDD = immutable

→ Difference b/w a narrow transformation and wide transformation

→ Define map, flatmap, filter and coalesce as narrow transformations.

→ List two wide transformations.

① map : Apply function to each element of RDD



def lower(line):

return line.lower()

lower case letters → lower-text-RDD = text-RDD.map(lower)

② flatmap : map then flatten o/p



NOTES

18

2016

Thursday
August

231-135 • WK 34

07

July 2016

wk	M	T	W	T	F	S	S
27					1	2	3
28	4	5	6	7	8	9	10
29	11	12	13	14	15	16	17
30	18	19	20	21	22	23	24
31	25	26	27	28	29	30	31

08

August

wk	M	T	W	T	F	S
32	1	2	3	4	5	6
33	8	9	10	11	12	13
34	15	16	17	18	19	20
35	22	23	24	25	26	27
36	29	30	31			

1. map(func) :- Returns a new distributed dataset formed by passing each element of source through a function func.

2. filter(func) :- Returns a new dataset formed by selecting those elements of the source on which func returns true.

3. flatMap(func) :- Similar to map, but each ip item can be mapped to 0 or more op items (so func should return a seq. rather than a single item)

4. mapPartitions(func) :- similar to map, but runs separate on each partition (block) of the RDD, so func must be of type Iterator<T> => Iterator<U> when running on an RDD of type

5. mapPartitionsWithIndex(func) :- similar to mapPartitions, but also provides func with an integer value representing the index of partition, so func must be of type (Int, Iterator<U>) => Iterator<U> when running

NOTES

T/W/T/F/S/SU/

Date.....

Actions :- Following table gives a list of Actions which returns value

1. reduce(func)

→ 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 in parallel.

2. collect()

→ Returns all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.

3. count() :- Returns the no. of elements in the dataset

4. first() :- Returns the first element of the dataset (similar to take(1))

5. take(n) :- Returns an array with the first n elements of the dataset

6. takeSample (withReplacement, num, [seed])

→ Returns an array with a random sample of num elements of the dataset, with or without replacement, optionally pre-specifying a random number generator seed.

Priority

7. takeOrdered (n, [ordering])

Returns the first n elements of the RDD using either their natural order or a custom comparator.

8. saveAsTextFile (path)

Writes the elements of the dataset as a text file (or sequence of text files) in a given directory in the local filesystem HDFS or any other Hadoop-supported file system. Spark calls toString on each element to convert it to a line of text in the file.

9. saveAsSequenceFile (path) (Java & Scala)

→ Writes the elements of the dataset as a Hadoop SequenceFile in a given path in the local filesystem, HDFS or any other Hadoop-supported file system. This is available on RDDs of key-value pairs that implements Hadoop's writeable interface.

10. saveAsObjectFile (path) (Java & Scala)

Writes the elements of the dataset in a simple format using Java serialization, which can then be loaded using SparkContext.objectFile()

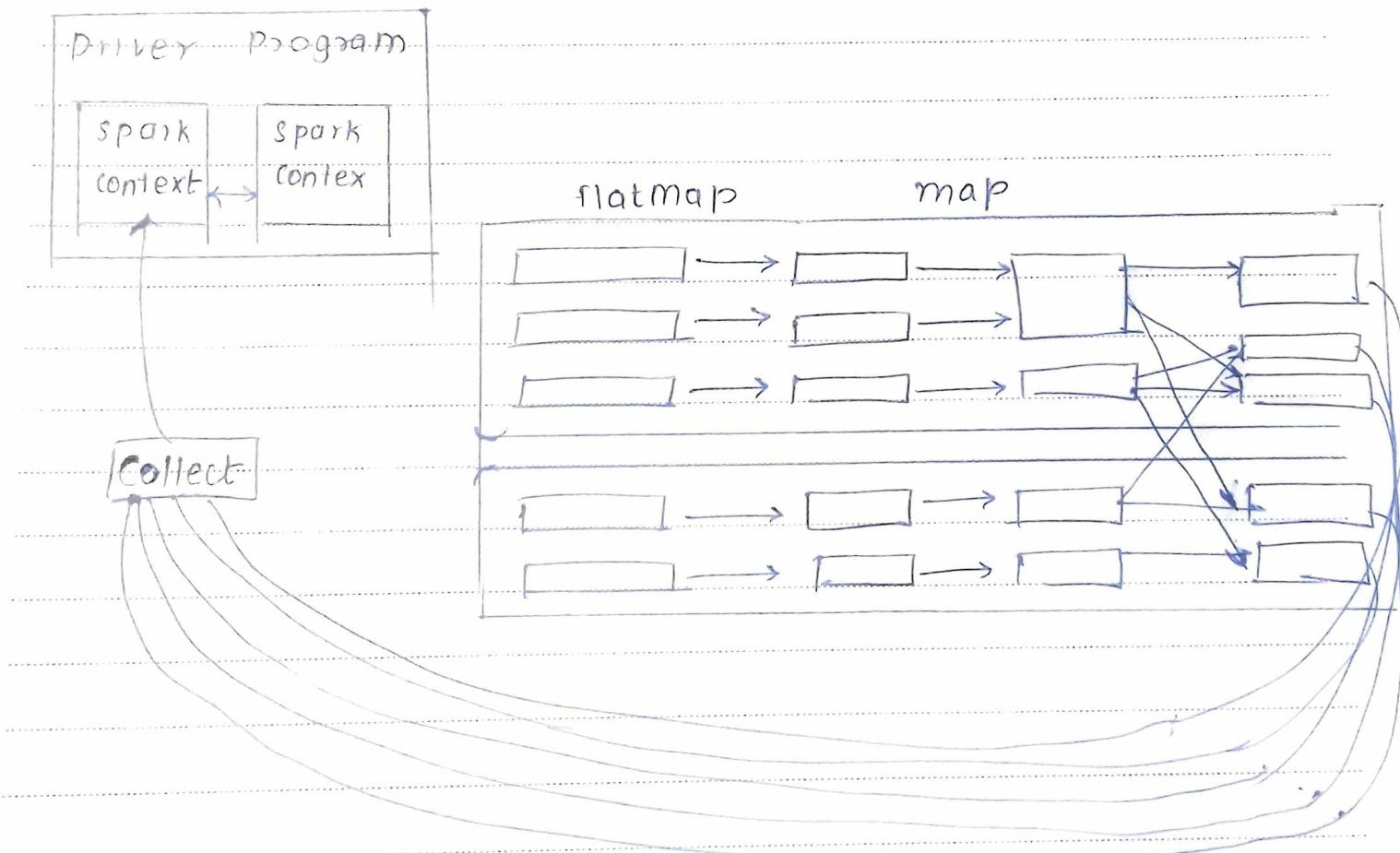
11. countByKey () - Only available on RDD of type (K, V)

Returns a hashmap of (K, Int) pairs with the count of each key!

12. foreach(func) :- Runs a function func on each element of the dataset: This is usually done for side effects such as updating an accumulator or interacting with external storage systems.

Courses :- Spark Core - Actions

- Explain the steps of spark pipeline ending with a collect actions
- List four common actions operations in spark.



some common Actions

- ① collect
- ② take(n)
- ③ reduce(func)
- ④ ~~map ?~~ → save to local file or HDFS

- Learning is a treasure that will follow its owner everywhere

Priority

2

December 2016
T W T F S S
1 2 3 4
5 6 7 8 9 10 11
12 13 14 15 16 17 18
19 20 21 22 23 24 25
26 27 28 29 30 31

01 January 2017
wk M T W T F S S
01/01 30 31 1
02 2 3 4 5 6 7 8
03 9 10 11 12 13 14 15
04 16 17 18 19 20 21 22
05 23 24 25 26 27 28 29

provides a common query language · 2016 Saturday November top of Spark Core that introduces a new data abstraction called SchemaRDD.
which provides support for structured and semi-structured data.

12 WK 49

video 4

Spark SQL :- Spark SQL is the component of SQL

that enables querying structured and unstructured data through a common language.

→ It can connect to many data sources and provide APIs to convert the query results to RDDs in python, Scala, Java programs.

→ Spark SQL gives a mechanisms for SQL users to deploy SQL queries on Spark.

→ Business intelligence tools to connect to spark using protocols like JDBC and ODBC.

→ Spark SQL also provides API to convert query data into DataFrames to hold (distributed) data named

↓ look like relational databases columns

S1 : To run any SQL spark is to create a From pyspark.sql import SQLContext
SQLContext. sqlContext = SQLContext(sc)

S2 : Once you have an SQL context, you want to leverage it to create a DataFrame so you can deploy complex operations on the data set.

existing RDD

DataFrames

Hive Tables

any other data source

Relational Operations :-

Sunday 13

→ perform Relational processing such as
NOTES declarative Queries

→ Embed SQL queries inside Spark program.

leverages spark core's fast scheduling capability to perform streaming analytics

14

31-07 - WK47

November		December	
10	11	12	13
2016 data in mini-batch	2017 data in mini-batch	2017 data in mini-batch	2017 data in mini-batch
Monday and performs RDD			
November (Resilient distributed datasets)			

RDDs can be converted to DataFrame of data.

- ① Convert each line into a row.
- ② Data is in DataFrame, \Rightarrow perform all sort of transformation

Show, printSchema, select, filter, groupBy, data streaming data

Now spark reads streaming data

List different sources of streaming data

Spark's streaming windows

spark streaming provides scalable processing for real-time data and runs on top of spark core.
continuous \Rightarrow grouped \Rightarrow parallel
data streams \Rightarrow discrete RDD

spark streaming provides APIs for Scala; Java and Python, like other spark products.

spark streaming can read data from many different types of sources including Kafka and Flume.

Kafka \rightarrow High throughput published subscribed messaging system.

Flume \rightarrow collects and aggregates log data.

spark streaming can also read from batch I/P data source such as HDFS, S3, and many other non data sources.

NOTES spark streaming reads streaming data, and converts it into micro batches with size called streams which is short for discretized stream.

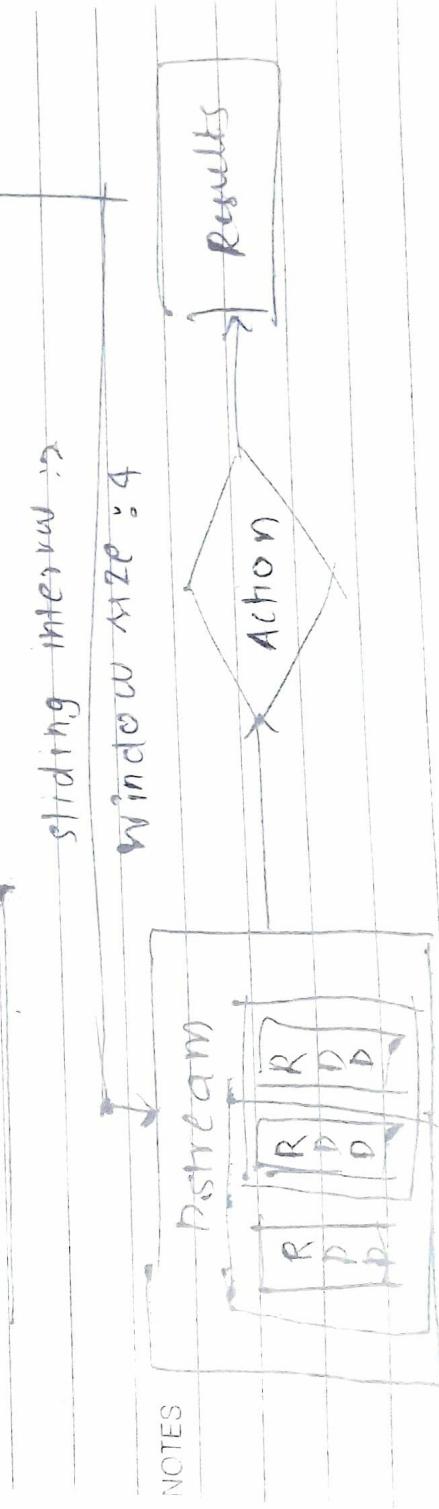
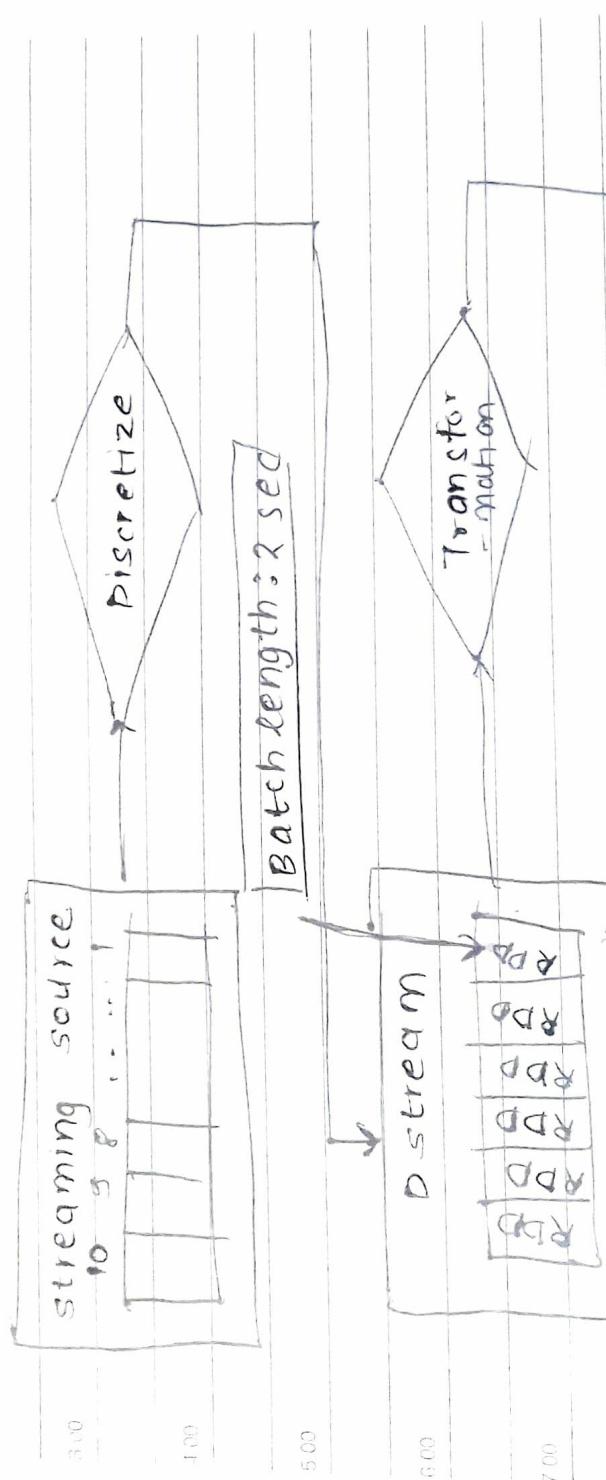
December 2016							January 2017						
S	M	T	W	F	S	S	S	M	T	W	F	S	S
1	2	3	4	5	6	7	8	9	10	11	12	13	14
15	16	17	18	19	20	21	22	23	24	25	26	27	28
29	30	31	1	2	3	4	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21	22	23	24	25
26	27	28	29	30	31	1	2	3	4	5	6	7	8
22	23	24	25	26	27	28	29	30	31	1	2	3	4

2016 | 15
 Tuesday
 November | 320-046 • WK 47

Transformations \Rightarrow Dstreams (map, reduce, filter)

- Spark streaming is spark's library to work with streaming data in near real time.
- Dstreams can be used just like any other RDD and can go through the same transformations as batch datasets
- Dstreams can create a sliding window to perform calculations on a window of time.

Creating and Processing Dstreams



NOTES

16 leverages spark core's fast scheduling framework ab
mlib is a distributed ~~for~~ ml framework
because of the 2016 distributed memory-based Wednesday spark architecture. November It is according to benchmarks, done by the ~~the~~ mlib developers against CS benchmarks, done by mlib (APIs) implementation.

Main Take-Aways
Spark uses Dstreams to make discrete RDDs

- Spark uses Dstreams to make discrete RDDs
- same transformations and calculations applied from streaming data
- same transformations and calculations applied to batch RDDs can be applied.
- Dstreams can create a sliding window to perform calculations on a window of time.

video

Spark mllib - What mllib is
List main categories of techniques available in mllib
Code segments containing mllib algo.

video

Spark mllib - scalable ml library
Provides distributed implementations of common ml algo and utilities
Has APIs for Scala, Java, Python, and R.
ml Algo & techniques:
• ML Learning

- classification
- regression
- clustering

- Statistics
- summary statistics
- sampling
- Utilities
- dimensionality reduction
- transformation
- Evaluation metrics

Spark mllib is nine times as fast as the Hadoop disk-based version of Apache Mahout

notes

Mlib

Ex. → summary statistics

- compute column summary statistics

```
from pyspark.mllib.stat import Statistics
# Data as RDD of Vectors.
datamatrix = sc.parallelize([[1,2,3],[4,5,6],[7,8,9],[10,11,12]])
# compute column summary statistics

1.00
summary = statistics.colStats(datamatrix)
print(summary.mean())
print(summary.variance())
print(summary.numNonzeros())
3.00
```

Mlib Ex. - classifier

* Build decision tree model for classification
→ from pyspark.mllib.tree import DecisionTree,

DecisionTreeModel

```
>> from pyspark.mllib.util import MlUtils
7.00
# Read and parse data.
```

```
data = sc.textFile("data.txt")
```

Decision tree for classification

NOTES

```
model = DecisionTree.trainClassifier(parsedData, numClasses=2)
print(model.toDebugString())
model.save(sc, "decisiontreeModel")
```

23

328-038 • WK 48

2016

Wednesday
November

October 2016

wk	M	T	W	T	F	S	S
40/45	31				1	2	
41	3	4	5	6	7	8	9
42	10	11	12	13	14	15	16
43	17	18	19	20	21	22	23
44	24	25	26	27	28	29	30

November 2016

wk	M	T	W	T	F	S	S
45		1	2	3	4	5	6
46	7	8	9	10	11	12	13
47	14	15	16	17	18	19	20
48	21	22	23	24	25	26	27
49	28	29	30				

MLlib Example - clustering

- Build k-means model for clustering

```
>> from pyspark.mllib.clustering import KMeans,
   KMeansModel
>> from numpy import array
# Read and parse data
>> data = sc.textFile("data.txt")
>> parsedData = data.map(lambda line:
   array([float(x) for x in line.split('')]))
1.00
```

k-means model for clustering

```
>> clusters = KMeans.train(parsedData, K=3)
>> print(clusters.centers)
```

3.00

→ MLlib is spark's machine learning library

- Distributed implementations

• Main categories :-

5.00

6.00

7.00

NOTES

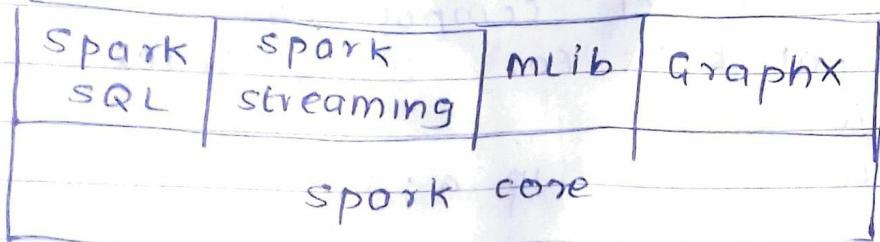
6

	S	S
1	2	
6	7	8
12	13	14
19	20	21
25	26	27

graphX is a distributed graph-processing
on top of Spark. It provides an API for
expressing graph computation that can model
the user-defined graphs by using Pregel abstraction
Week 28 Day 187 Date 09/07/15 Monday
Video :- API. It also provides an optimized API for this abstraction.

Spark GraphX :- Describe what GraphX is

- Explain how Vertices and Edges are stored
- Describe how pregel works at a high level

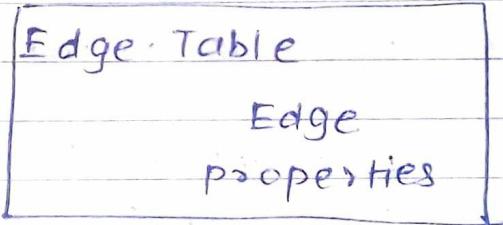
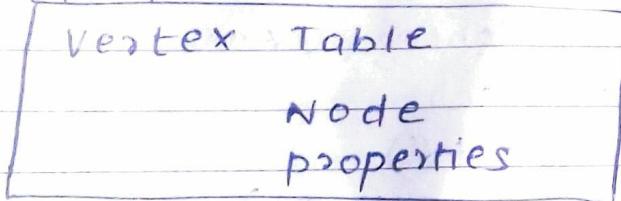


GraphX is Apache Spark's API for graphs and graph-parallel computation.

uses a property graph model

⇒ Both Nodes and Edges can have attributes and values.

properties → Tables

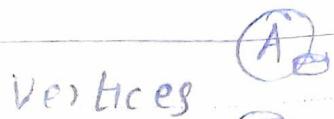


GraphX uses special RDDs

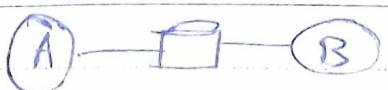
Vertex RDD[A] extends RDD[(VertexID, A)]

Edge RDD[ED, VD] extends RDD[Edge[ED]]

Triplets :- The triplet view logically joins the vertex and edge properties.



Edges



Triplets :-



Pregel API :- Bulk-synchronous parallel messaging

- mechanism constrained to the topology of the graph
- GraphX → Graph Parallel Computations
- Special RDDs for storing vertex and Edge info
- Pregel operator works in a series of super steps

JSON → DataFrame

Read

```
df = sqlContext.read.json ("filename.json")
```

Display

df.show()

RDD

of Row objects → DataFrames

```
From pyspark.sql import SQLContext, Row
```

```
sqlContext = SQLContext(sc)
```

```
lines = sc.textFile ("filename.txt")
```

```
cols = lines.map (lambda l: l.split (" "))
```

```
data = cols.map (lambda l: Row (name=p[0],  
                                zip=int(p[1])))
```

```
df = sqlContext.createDataFrame (data)
```

```
df.registerTempTable ("table")
```

```
Output = sqlContext.sql ("SELECT * FROM table  
WHERE ...")
```

Which part of spark is in charge?

→ Driver program

How does lazy evaluation work in Spark?

→ Transformations are not executed until the action stage.

What are the consequences of lazy evaluation as mentioned in le?

→ Errors sometimes don't show up until the action stage.

Wide transformation:- A transformation that requires data shuffling across node partitions.

Where does the data for each worker node get sent to after a collect function is called?

→ spark context

Dataframes:- A column like data format that can be read by spark SQL.

Can RDD's be converted into Dataframes directly without manipulation?

→ No: lines have to converted into row.

NOTES Functions of Spark SQL:-

① Enables relational queries on spark.

② Deploy business intelligence tools over spark.

③ Connect to variety of databases.

18

2016

Friday

November

323-043 • WK 47

October 2016							November 2016							
10	11	M	T	W	T	F	S	S	M	T	W	T	F	S
wk	40/45	31	1	2	3	4	5	6	7	8	9	10	11	12
41	45	17	18	19	20	21	22	23	24	25	26	27	28	29
42	46	10	11	12	13	14	15	16	17	18	19	20	21	22
43	47	17	18	19	20	21	22	23	24	25	26	27	28	29
44	48	24	25	26	27	28	29	30	28	29	30	1	2	3

9. What is a triplet in GraphX?

→ A type of data to contain the information on connections b/w vertices and edges.

10.00

11. How many tweets have location not null?

→ 6937

12.00

12. How many people have more followers than friends?

→ 5809

13.00

3.00

4.00

5.00

6.00

7.00