



Scala

val

→ creates immutable variable
→ its like final keyword in Java

var

→ creates mutable variable.

Data Types

String, Int, short, Byte
char, Float, Double

var name: String = "Pramod"

var name = "ABC"

var (s1, s2, s3, s4) = (10, 20, 30, 40)

Triple quotes for multiline string.

val address = """#26, jpnagar
Bangalore
India""

Interpolator

s & f are the interpolators

s "fname - \$lname, how are you?"

f "\$fname* \$lname, come here!"

f "\$p1% . \$p2%"

Pattern Guarding

var day = "Monday"

var someday = "Friday"

day match {

case someday ⇒ if (someday == "Saturday") {
 println("Working")}

Case Classes

val x: Int = Random.nextInt(10)

x match

case 0 ⇒ "zero"

case 1 ⇒ "one"

case - ⇒ "other"

Pattern guards

Pattern guards are boolean expressions which are used to make cases more specific.

Eg: notification match

case Email(sender, -) if
importantPeopleInfo.contains(sender)
⇒ "You got an email from
someone"

case SNS(number, -) if
importantPeopleInfo.contains
(number) ⇒ "You got an sms from
special someone"

Eg: var day: Any = 10

day match {

case day: Int ⇒

println("Integer")

case day: Double ⇒

println("Double")

case day: String ⇒

println("String")

case - ⇒ println("default")

}

Passing function as parameter

def sum (sqrFn : Int \Rightarrow Int,

a:Int, b:Int) : Long = {

var sum = 0

for (i < a to b) {

sum += sqrFn(i)

}

sum

3

def sqrFn (i:Int) = i*i

println (sum(sqr, 3, 30))

Type casting in Scala

a. asInstance Of [Int]

char
double

object variables

To convert to string
above one won't work
vars = a.toString()

Before Doing Examples

Scala main topics.

May 30 2023
3.3.0

→ Latest Scala version 3.3.0

→ Type Inference (No need to mention
data type)

→ Executing Scala

- scalac test.scala
- scala test

→ For loop

Object MainObject {

def main(args: Array[String]) {

for (a < -1 until 10) {

println(a)

3 3 3

• for (a < -1 to 10 if a % 2 == 0) {

yield

var result = for (a < -1 to 10) yield a

for (i < -result)
println(i)

O/P
60
10

var list = List(1, 2, 3, 4, 5)

→ for (i < -list)
println(i)

→ list.foreach {

println;

3

→ list.foreach(println)

→ list.foreach(element: Int)

⇒ print(element + " ")

→ In Scala we don't have primitive variables

→ Scala has seamless Java interop

→ Type Inference

→ Concurrency & Distribution

→ Traits

→ Pattern Matching

→ Higher Order Fns.

→ Twitter announced that it had switched large portions of its backend from Ruby to Scala

→ App Inc, Walmart

→ Singleton Object

• Scala has no static variables or methods.

• It uses singleton object which is essentially class with only 1 object in source file.

→ Immutability using val
- This helps to manage concurrency control.

→ Lazy Computation

- We can declare lazy variable by using `lazy` keyword.
- Only vals can be lazy → `lazy val a = 10`

→ Case Classes & Pattern matching.

→ Concurrency control.

→ String Interpolation
`s, f & raw`

→ If, else,

Scala Object & Class

Object	Class
→ It has state	→ Has methods
- data values of an object	can be defined here
→ static like instance of Scala	→

Singleton Object

→ object APP {
 def main(args: Array[String]) {
 SingletonObject.hello()
 }

(3) // Create 1 singleton object with hello method.

Companion Object

Object created to create an instance to a class.

Scala Object & Classes

Eg: of class

```
class Student {  
  var id: Int = 0;  
  var name: String = null;  
}
```

```
object APP {  
  def main(args: Array[String]) {  
    var s = new Student()  
    println(s.id + " " + s.name)  
  }  
}
```

Eg 2:

```
class Student(id: Int,  
             name: String)
```

```
{  
  def show() {  
    println(id + " " + name)  
  }  
}
```

object -

```
var s = new Student(100, "Narain")
```

Anonymous Object

```
class Arithmetic {  
  def add(a: Int, b: Int) {  
    ...  
  }  
}
```

```
new Arithmetic().add(10, 10);
```

Primary Constructor

```
class Student(id: Int, name: String) {  
  def show() {  
    ...  
  }  
}
```

Secondary constructor

```
class Student(id: Int, name: String) {  
  def show() {  
    ...  
  }  
}
```


Exception Handling.

```
class ExceptionExample {
    def divide(a:Int, b:Int) = {
        a/b
        println("Rest of the code is
2           executing")
    }
}
```

```
object MainObject {
    def main(args: Array[String]) {
        var e = new ExceptionExample()
        e.divide(100, 0)
    }
}
```

```
try {
    a/b
} catch {
    case e: ArithmeticException
        ⇒ println(e)
}
```

```
try {
    a/b
    var arr = Array(1, 2)
    arr(10)
} catch {
    case e: ArithmeticException
        ⇒ println(e)
    case ex: Throwable ⇒ println
        ("found unknown exception")
}
```

```
finally {
    try {
        a/b
    } var arr = Array(1, 2)
    arr(10)
}
```

catch {

case e: ArithmeticException ⇒
 println(e)

case ex: Exception ⇒ println(ex)

case th: Throwable ⇒ println

(found a unknown exception " + th)

Data Defn language

```
val df = spark.read.format("json")
    .option("inferSchema", true)
    .schema("ca_dt schema DDL")
    .load("filestore/...")
```

sc.parallelize(seqL)

sc.textFile

Datatypes.

→ spark.read.format(" ") .option
 .schema
 .load

schema = "customized STRING,
 `orderId` INT, `products`
 ARRAY<STRUCT<`orderId` STRING,
 `quantity` BIGINT,
 `soldprice` DOUBLE

→ STRING INT BIGINT
DOUBLE

→ `orderId` STRING, `OrderId` BIGINT
 `products` DOUBLE

→ `orderId` STRING
 `orderId` STRING, `OrderId` BIGINT
 `products` DOUBLE
 `Order` STRING

Scala Singleton & Companion Object

Objects
Classes
Constructor
OODS

Scala singleton

- declared by using object keyword.
- There is no static concept in scala, so scala creates a singleton object to point for your program execution.
- Methods created inside a singleton object are accessible globally.
- A singleton object can extend classes & traits.

Companion object

- when you have an object with same name as companion class.

Primary, Secondary constructor & Constructor Overloading

Default Primary Constructor

class Student { }

```
printin ("Hello from default constructor")
```

3

Primary constructor

```
class Student (id: Int, name: String) { }
```

```
def showDetails () { }
```

```
printin (id, name)
```

3

main { }

```
vars = new Student (101, "Ramo")
```

3

→ As usual primary constructor gets executed during the object creation.

Inheritance in Class

→ single, Multilevel, Hierarchical (A → B → C) is allowed for classes in scala

→ whereas using Traits Multiple & Hybrid (JNL) can be achieved.

→ Use override keyword to override a val type variable as below

override val speed: Int = 150, But var defined variables can be changed just like that

Scala → left out topics.

3.2.2r as of January 30 2023.

Features of Scala

- ① Type Inference
- ② singleton object
- ③ Immutability
- ④ Lazy computation
- ⑤ case classes & pattern matching.
- ⑥ Concurrency control
- ⑦ String Interpolation
- ⑧ Higher order fns
Taking fn as argument & returning function.
- ⑨ Traits
- ⑩ Rich set of collections

Concurrency Control - Future

A Future represents a value which may or may not currently be available, but will be available at some point or exception if not available.

In single threaded world you bind the result of a method call to a variable like this.

→ def aShortRunningTask(): Int = 42

val x = aShortRunningTask()

// Here 42 is immediately bound to x.

→ def aLongRunningTask(): Future[Int] = ???

val x = aLongRunningTask()

// So here aLTask() takes undefined amt of time to return the value of x which will be available in future

→ def aShortRunningTask(): Int =

Thread.sleep(500)

42

val x = aShortRunningTask()

println("Hello")

// println stmt is printed almost immediately once aShortRunningTask is completed

// If aShortRunningTask is created as Future println is printed almost immediately.

→ Akka is a toolkit for building highly concurrent, distributed, & resilient message-driven applications for Java and Scala.

English meaning of concurrent →

"to be served at the same time"
"mentioning regards one point"

PointIn (s "before the 'when' block")
Future → Handle relative slow computation code on another
thread & call me back with result when you are done

Eg: Running multiple features & joining their results

Note: some Future methods
→ filter
→ flatmap
→ map

And ~~these~~ these callback methods
→ onComplete
→ andThen
→ foreach

Transformation methods
→ fallbackTo
→ recover
→ recoverWith

```
import scala.concurrent.Future
import scala.concurrent.ExecutionContext.Implicits.global
import scala.util.{Failure, Success}
val startTime = System.currentTimeMillis()
def delta() = System.currentTimeMillis() - startTime
def sleep(millis: Long) = Thread.sleep(millis)
```

@main def multipleFutures =

pointIn (s "creating the futures : \${delta()}")

```
val f1 = Future { sleep(800); 1 }
```

```
val f2 = Future { sleep(600); 2 }
```

```
val f3 = Future { sleep(400); 3 }
```

val result =

```
for
  r1 <- f1
  r2 <- f2
  r3 <- f3
```

//join the futures in for expression

yield

pointIn (s "In the 'yield' : \${delta()}")
(r1 + r2 + r3)

//process the result

result.onComplete {

case Success(x) =>

pointIn (s "in the success case : \${delta()}")

pointIn (s "result = \$x")

case Failure(e) =>

e.printStackTrace

Println (is "before the 'sleep(3000)': \$ {delta(3)}")
Sleep(3000)

O/P

creating the futures: 1
before the 'sleep(3000)': 2
in the 'yield': 806
in the Success case: 806
result=6

→ As the o/p shows the futures are created rapidly & in 2 seconds the print stmt above sleep(3000) is executed showing the full code is run on the JVM's main thread.

- Then at 806 ms the 3 futures complt & the code in the yield block is run
→ The code goes to the success case in the onComplete method.

→ The 806 o/p is a key to seeing that the 3 computations are run in parallel. If run sequentially it would be 1400ms

Key Points about Futures

- Futures to run tasks off the main thread
- The value inside a Future is always an instance of one of the Try types: Success or Failure

b)

Scala Running Pgm

Object Scala Examples

```
def main(args: Array[String]) {  
    println("Hello Scala")  
}
```

3

3

Save as ScalaExample.scala

cmd ↳ scalac ScalaExample.scala
scala ScalaExample

Data Types

→ Boolean Byte short Char Int Long Float Double String.
1bit 1Byte 2B 2B 4B 8B 4B 8B → Any

Latest version of Scala.

3.2.2v as of January 2023

spark supports Scala 2.12/2.13

Python 3.7+ versions (3.9 is deprecated)

SPARK latest version
3.4.0v was released
April 14 2023

b)

* Why we use = (rough answer)

Python vs Scala

Python → Guido van Rossum
in 1989
Netherlands

→ Python is a dynamically typed language

Eg → $x = 6$

// This will store 6 in the memory & binds name x to it. After it runs, type of x will be int
`print(type(x))`

→ $x = \text{Hello}$

// This will store hello at some location in memory & binds the name x to it. After it runs type of x, it will be string.

→ Easy to learn

→ No need to specify objects since dynamic

→ Huge community

→ Has such libraries ML, NLP

→ Used for small-scale projects

→ doesn't provide scalable feature support

→ Martin Odersky
in 2003
Switzerland

actions in Scala

→ statically type language.

Advantages of statically type

→ Dynamically typed language is fast, But these functions don't give any info about what arguments to pass u will end up passing junk arguments

JS is dynamic so TS is added for static typing.

→ Difficult to learn, 10 times faster than Python

→ Need objects.

→ Okay community.

→ No such tools.

→ Large-scale projects.

→ Provide scalable feature support.

+b)

based on the above tool for select windows



- Why we use = (equal) operator in Scala functions

Pattern Matching in Scala

Functions
in

- ① In java we have switch case

String day = "Monday"

```
switch (day) {
```

```
case "Monday":
```

stmts

break;

```
case :
```

stmts

break

```
default:
```

}

Scala

```
var a =
```

```
a match {
```

```
case 1 => println("one")
```

```
case 2 =>
```

```
case - =>
```

}

→ In Scala there is no break statement can be done as below using breakable method.

import scala.util.control.Breaks._

object Mainobj {

```
def main(args: Array[String]) {
```

```
breakable {
```

=

break

else

}

• F

var

• Pa

0

F

A

- Case class
- Why we use = (equal) operator in Scala functions
 - we can create fn with or without equal operator
 - If used returns values
 - If not used it will return anything & acts as subroutine.

```
def fn(a:Int):Int = {
  var b = a*a;
  b
}
```

Note: From Scala 3

= operator on all the functions is mandatory even if it returns anything or not

- Function Parameter default value in Scala

```
def sum(a:Int=0, b:Int=0):Int = {
```

3
Also here 0,0 is the default values if not passed it considers the value 0.

- Function named parameter

```
vara = sum(a=15, b=10) // while passing params u can pass using the already defined variable names.
```

- Passing function as parameters
or Higher Order Function in Scala.

```
def multiplyBy(a:Int):Int = {
  a*2
}
```

```
def main(args: Array[String]) {
  output(25, multiplyBy)
}
```

```
3
def output(a:Int, f: Int  $\Rightarrow$  AnyVal): Unit = {
  println(f(a))
}
```

- Function composition

```
example
var result = multiplyBy(addBy(10)) // function composition
```

- Anonymous lambda function → can be created by using \Rightarrow or -

```
var result1 = (a:Int, b:Int)  $\Rightarrow$  a+b // anonymous function by using  $\Rightarrow$ 
var result2 = (-:Int) + (-:Int) // anonymous function by using - wildcard
println(result1(10, 10))
println(result2(10, 10))
```

Based on the above code select which one is correct

→ ~~class~~ - classes in Scala - int

- Multiplication expression in Scala
def add(a:Int, b:Int) = {
 a + b
}

- Scala Function currying.

Object MainObject{
 ...
}

```
def add(a:Int)(b:Int) = {  
    a + b  
}
```

```
var result = add(10)(10)  
print(result);
```

println ("O/P is " + result)
var add74

```
var add10 = add(10);  
var result;
```

```
var result2 = addit(3)  
println(result2)
```

Println ("O/p is " + result)

3

• Nested function in Scala •

def add(a:Int, b:Int, c:Int) = {

def add(x:Int, y:Int) = {

$2C+g$

$\text{add}_2(a, \text{add}_2(b, c))$

```
def main (args: Array[String]) = {
```

```
var result = add(10, 10, 10)
```

۲

→ Case classes in Scala - Imp.

• regular classes which are immutable by default & decomposable through pattern matching.

- It uses equal method to compare instances structurally.
- It does not use new keyword to instantiate the object.

→ case class className(params)

Eg:

```
case class Airport(a:Int, b:Int)
object MainObject {
    def main(args: Array[String]) {
        var airp = Airport(10, 10)
        print(airp.a + " " + airp.b)
    }
}
```

Eg2:

~~trait SuperTrait~~ Imp.

~~case class~~

trait Flights

case class Boeing(a:Int, b:Int) extends Flights

case class Mig(a:Int) extends Flights

case object F16 extends Flights

object MainObject {

```
    def main(args: Array[String]) {
        callCase(Boeing(10, 10))
        callCase(Mig(10))
        callCase(F16)
    }
}
```

def callCase(f: SuperTrait) = f match {

case Boeing(a, b) ⇒ println("a = " + a + " b = " + b)

case Mig(a) ⇒ println(a)

case F16 ⇒ println("No Args")

Based on notes above made by ~~select with nulls~~



→ Scala file handling

- Handling file operations
- Predefined methods to deal with the file
- Create, open, write & read the file
- Scala has a complete package for file handling
Scala.io

```
import scala.io.Source
object MainObj {
    def main(args: Array[String]) {
        val filename = "Scalafile.txt"
        val filesource = Source.fromFile(filename)
        for (line <- filesource.getLines) {
            println(line)
        }
        filesource.close()
    }
}
```

→ Multidimensional Array

var arr = Array.ofDim [Array TYPE] (R, C) ^{"row, column"}

var arr = Array (Array (1, 2, 3), Array (4, 5, 6))

→ String interpolation

So, we can use \$ inside string sentence to show the value of variable
%s, %d etc

raw → prints as it is ignore in it

e.g. var s = "ABC\6 DEF \t GHI" O/P with subsitution
var s = raw(" ", ,) O/P as it is.

→ Exceptions in Scala (throw, throws)

• Scala makes "checked" vs "unchecked" very simple. It doesn't have checked exceptions. All exceptions are unchecked in Scala.
Scala pgm without exception handling.

class Example {

def divide (a:Int, b:Int) = {
 a/b

3 Point in ("Rest of the code is executing...")

Object MainObject {

def main (args: Array [String]) {

var e = new Example()

e.divide (100, 0)

> won't get exec.

O/P

- java.lang.ArithmeticException

Exception Propagation

40-50

Collections

50-60

Companion Object

- There is no static concept in Scala
- companion class & companion object must be defined in same source file.

singleton object

• Can call the methods inside this without creating any object.

based on the above hadoop select with nulls

→ Anonymous object in Scala

```
class Arithmetic {  
    def add(a:Int, b:Int)  
        var sum = a+b  
    } Print(sum)
```

```
object Main {
```

```
    def main(args: Array[String]) {
```

```
        new Arithmetic.add(10, 10)
```

→ Method Overriding in Scala

Extend class & override the method.

"extends"

```
override def run():
```

→ final

Used to prevent Inheritance

final class can't be inherited itself

→ Abstract class

- has both abstract & non abstract methods
- using extends keyword other classes can give implementation.
- if all methods are undefined define class as abstract

→ Scala Trait

- all abstract or some abstract & non abstract methods

• Trait Mixins

↳ means u can extend any number of traits with a class or abstract class

↳ u can extend only traits or
combo of traits & class or
traits & abstract class

→ Access modifiers

only 3	No modifier/public	class	Companion	Subclass	Package	world
	Protected	yes	yes	yes	yes	yes
	Private	yes	yes	yes	no	no

Companions is a singleton object named same as the class

→ Array (1,2,3,4,5)

```
var arr = var arr: Array[Int] = new Array[Int](10)
```

```
var arr = new Array[Int](10)
```

```
var arr: Array[Int] = new Array(10)
```

```
var arr = Array(5, 4, 3, 2, 1)
```

based on above Hadoop select with nulls

Final variable

class Vehicle {

 final val speed:Int = 60
 3

class Bike extends Vehicle {

 override val speed:Int = 100

 def show() {

 println(speed)

 3

//Error value speed cannot override final member.

Final Method

class Vehicle {

 final def show() {
 3 println("Vehicle is running")

class Bike extends Vehicle {

 override val speed:Int = 100

 override def show() {

 3 println("Bike is running")

 3 //Error

Final class

→ Final class can't be inherited.

Abstract class

→ object of abstract class can't be created.

→ If a class extend abstract class it must provide implementation of its all abstract methods.

Eg:

abstract class Bike(a:Int)

3

var b:Int = 20

var c:Int = 25

def run()

def performance() {

abstract method

3 println("Performance awesome")

3

Scala Trait

→ same Interface, collection of abstract - non abstract methods.

• trait Printable {

 def print()

3

class App extends Printable {

 def print() {

 3 println("Hello")

3

trait Printable {

 def print()

3

abstract class A4 extends Printable {, Java }

 def printA4() {

 3 println("Hello this is A4 sheet")

3

trait Printable {

 def print()

 def show() {

 3 println("ABC")

3

• extends with

trait ABC

trait DEF

class GHI extends ABC

with DEF

github link

③ Download winutils from based on above Hadoop package ① select winutils

Scala Access Modifiers
No Access package package class subclass companion
makeHer Yes Yes Yes Yes
Protected No No Yes Yes
Private No No Yes Yes
No No No No No

Scala Array

```
var arrName : Array[Type]
= new Array[Type](size)
```

```
var arrName = new Array[Object]
(size) or
```

Eg:

```
var arr = Array(1, 2, 3, 4, 5)
```

```
var arr = new Array[Int](6)
```

Multi-dimensional array

```
var arr1 = Array(Array(1, 2, 3, 4, 5),
                 Array(6, 7, 8, 9, 10))
```

```
var arr3 = Array.ofDim[Int](2, 5)
```

Apache Spark

Categories of Data

- ① Structured Data (tabular)
- ② Semi-structured Data (JSON)
- ③ Unstructured Data

Text, audio, video graphs
Logs

RDBMS handles

Yr → 2002 → Big data was

introduced to handle

Unstructured Data

→ Big Data Technologies

The software tools used
to manage all types of datasets
to transform them into business
insights.

we can categorize them into
we can categorize them into
MongoDB, Cassandra

Data Storage → Hadoop, Cassandra

Data Mining

Data Analytics

Data Visualization

→ Map Reduce

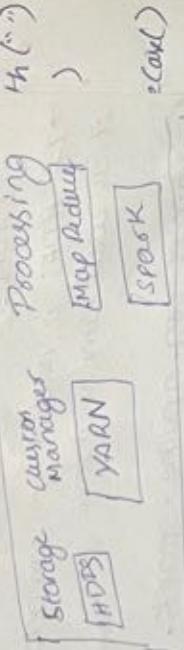
How Map Reduce Flow works

eg: word count

Input | splitting | mapping | shuffling | reduce | final output

(reduces)

In 2007 → YARN (reduces separate
negotiates) - Linu



Main Disadvantage of MapReduce

hdfs

job

every time the o/p is stored
onto HDFS. & retrieved for
the next job

→ But spark jobs are executed
in DAG format

Map Reduce

• Supports only Java

• Slow

Map Reduce

• Supports so many
other languages

Python, SQL, R, Scala (etc)

Scala

• Fast

Java Scala Spark

Java SDK 8

① JAVA-HOME /jdk1.8/bin

Java Path /jdk1.8/bin
click Java-version on cmd prompt

② Download Spark & Search
click Java-version on cmd prompt

Download Spark release, Package type
select Spark release, Package type
Pre-built Apache Hadoop 3.7

The older versions need to go to
Archived versions.

① Download utilities from any github
② Download utilities from apache hadoop
③ based on windows

Lite Line

Flame Retardant Cable 1100
PVC Insulated (Copper Conductor)

⑨ $C://$ → bigdata

↳ spark (standard)
↳ not for
↳ counts.

⑥ set 3 env variables (o)

→ SPARK_HOME

c/bigdata/spark

→ HADOOP_HOME ← counts

c/bigdata/spark

↳ try putting
windows inside

spark if this

wont work

→ Inside Path variable
c/bigdata/spark/bin

> spark - shell

scala >

• spark - Scala

• RDDs, RDDs, RDDs

-
- 18 High Level Syntax
- in RDD
- ① rdd = sc.parallelize(Listⁿ, partitions)
- = sc.textFile(path, partitions)
- ② rdd = rdd.map($i \Rightarrow i, i$)
 $(i \Rightarrow (i, i))$
 $(i \Rightarrow (i, i))$
- ③ rdd = rdd.flatMap($i \Rightarrow (i, spark("n"))$)
- ④ rdd = rdd.take(12)
- ⑤ rdd = rdd.takeOrdered(4)
- ⑥ rdd = rdd.head()
- ⑦ rdd = rdd.countByValue()
= rdd.countByValue()
= rdd.count()
- ⑧ rdd = rdd.union(rdd2)
- ⑨ rdd = rdd.intersection(rdd2)
- ⑩ rdd = rdd.subtract(rdd2)
- ⑪ rdd = rdd.saveAsTextFile("path")
- ⑫ rdd.cache()
- ⑬ rdd.broadcast("rdd")
- ⑭ sc.broadcast(DataBroadcast)
- ⑮ val dt = spark.read.json(path)
- ⑯ dt.printSchema()
- ⑰ dt.show()
- ⑱ dt.createOrReplaceTempView("student_table")
- ⑲ spark.sql("select *
from student_table")
- ⑳ rdd = rdd.saveAsTextFile("path")

val rdd = sc.add.map($i \Rightarrow$ $(i, (\text{scala.math.pow}(i, 0.5)))$) \rightarrow Broadcast Cache

// Broadcast

val names = sc.parallelize($\text{list}(\text{"spark"}, \text{"Hadoop"}, \text{"Mahout"})$)

• map($i \Rightarrow (i, i.\text{length})$)

• sortby($i \Rightarrow i.\text{length}$).collect

// • sortby key().collect
// sortby key(true).collect
// sortby(i $\Rightarrow i.\text{length}$).collect.read them from the disk.

// sortby($i \Rightarrow i.\text{length}$, true).collect

Cachel

// sortby($i \Rightarrow i.\text{length}$, true).map

val words = sc.broadcast($i.\text{cache}$)

// filter($i \Rightarrow \text{startsWith}(i, "s")$)

• collect

filterBy Range

\Rightarrow sc.parallelize($\text{list}("1", "2", "3")$)

• map($i \Rightarrow (i.\text{length}, i)$)

• sortby range(3, 10)

• filter($i \Rightarrow \text{contains}(i, "spark")$)

• filter($i \Rightarrow \text{equals}(i, "spark")$)

• filter($i \Rightarrow \text{length}(i) > 3$)

• collect

• map($i \Rightarrow (i, i.\text{length})$)

• filter($i \Rightarrow i.\text{contains}("spark")$)

• filter($i \Rightarrow i.\text{length} > 3$)

• collect

③

val words = Input.filter

• map . filter

val a = sc.broadcast(words)

• Broadcast is used to supply

a copy of small sized dataset

to each & every executor.

so whenever the executors need

this data set they dont need to

read them from the disk.

• sortBy uses least execution

for computation to enhance the

performance.

• A lineage is maintained, that

is the composition of steps to build

the RDD

• when we perform any action

on an RDD all the steps in

the lineage will be proceeded.

• Spark by default will

reprocess these steps everytime

we call an action. which is

inefficient.

• In order to avoid this we

call cache() on an RDD which

will store the processed data.

• print

Reposition & coalesce

```
var n = sc.parallelize(
```

$$\text{val } n_1 = n \cdot \text{represents}(o) \\ \text{val } n_2 = n \cdot \text{coalesce}_3$$

```

    all words = input · filter( $i \Rightarrow i \neq !$   

 $\text{newline}$ ). map( $x \Rightarrow (\text{lc-split}$   

 $(\text{`||`}\text{word}(\text{`}))\text{`})) \cdot \text{filter}$   

 $(-\text{eq}`value`ignoreCase(`separ`))$ )

```

• count()

UNION
INTERSECTION

SUBTRACT ZIP WITH INDEX

my odd partition

11 gives the solution.

\rightarrow `sc.parallelize(case)`

- 22 -

$\rightarrow \text{val } R = \text{odd} : \text{int} \cap \text{even}$
 $\rightarrow \quad \quad \quad .\text{collect}$
 $\rightarrow \text{odd} : \text{intersection}(\text{odd} \&$

= *xddi*. *subtract*(*rade*);

卷之三

• groupByKey()

cert. makes
calo " makes
named exec
Scalo
class Eco
class d

S →

- val `input` = sc.textFile("...")
- val `coords` = `input.flatMap(x ⇒`

$x \Rightarrow x.split(" ")$

\mapsto map

flatMap

If x is Array
if used on top, from Array to
of split in now

$\text{def flatMap}(i \Rightarrow i.split(" "))$

op

[Project, Greenberg]

[Alice's Adventures in Wonderland]

② flatMap

$\text{def flatMap}(i \Rightarrow i.split(" "))$

of project
GuruBabu
Alice's
Adventures
In
Wonderland

- `SortByKey()`
- `CountByValue()`

reduceByKey

op

$\text{map}(x \Rightarrow x.split(" ") \rightarrow \text{map}(x \Rightarrow$

$\text{toDouble})$

$)$

\mapsto reduceByKey($x -> x ->$

$x \text{ if } x > y \Rightarrow x \text{ else } y)$

③ map map

$\text{def flatMap}(i \Rightarrow i.split(" "))$

op

$\text{var j} = i \cdot \text{split}(...)$

$y \cdot \text{reduce}((x,y) \Rightarrow \text{if}$
 $(x.-y < y.-x) x \text{ else } y)$

④ map map

$\text{def flatMap}(i \Rightarrow i.split(" "))$

op

$\text{input.flatMap}(x \Rightarrow x.split$

$("\\n")) \cdot \text{map}(x \Rightarrow ((x \cdot \text{split}$

$("\\n", ",") \cdot \text{distinct}().\text{count}$

$)$

⑤ flatMap

$\text{def flatMap}(x \Rightarrow \text{map}(x \Rightarrow$

$\text{map}(x \Rightarrow ((x \cdot \text{split}$

$("\\n", ",") \cdot \text{map}(x \Rightarrow \text{if}$

$(x \cdot \text{split}(" ") \cdot \text{toDouble}))$

$\cdot \text{sortByKey}()$

$\cdot \text{reduceByKey}(-+ -)$

Life!

Flame Retardant
PVC Insulation

Spark Session

object JSON {

def main(args: Array[String]) {

Logger.getLogger("log")

• setLevel(Level.ERROR)

val spark = SparkSession

- builder
- appName ("sparkSQL")
- master ("local[4]")
- config ("spark.sql.aurhouse.Urls", "file:///C:/temp")
- getOrCreate()

val data = spark.read.json(path)

data.printSchema

data.show()

data.createOrReplaceTempView("ppd")
spark.sql("select * from people").show

m1

jfor

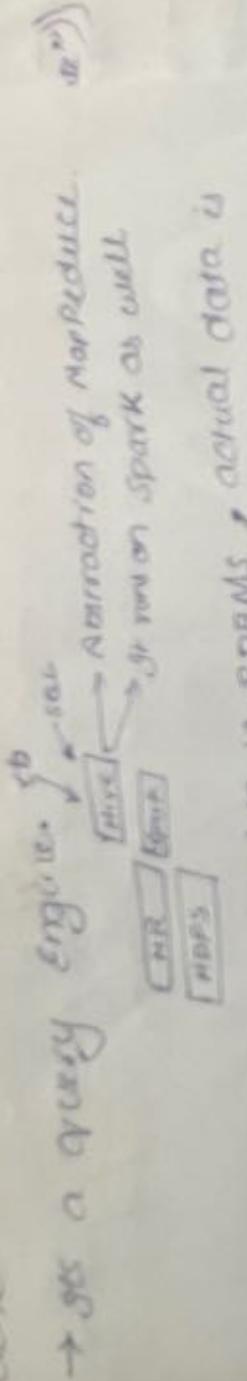
Life+

Flame Retardant

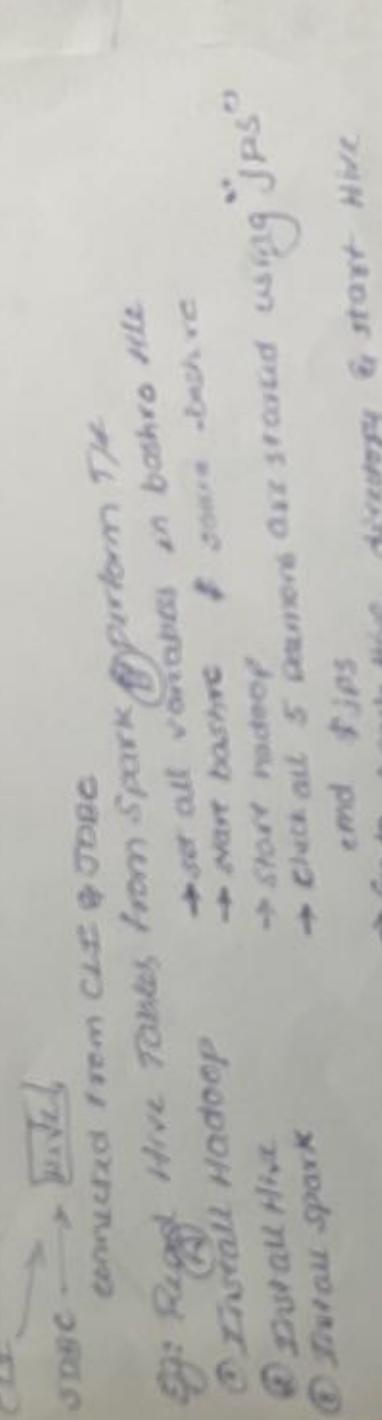
PVC Insulation

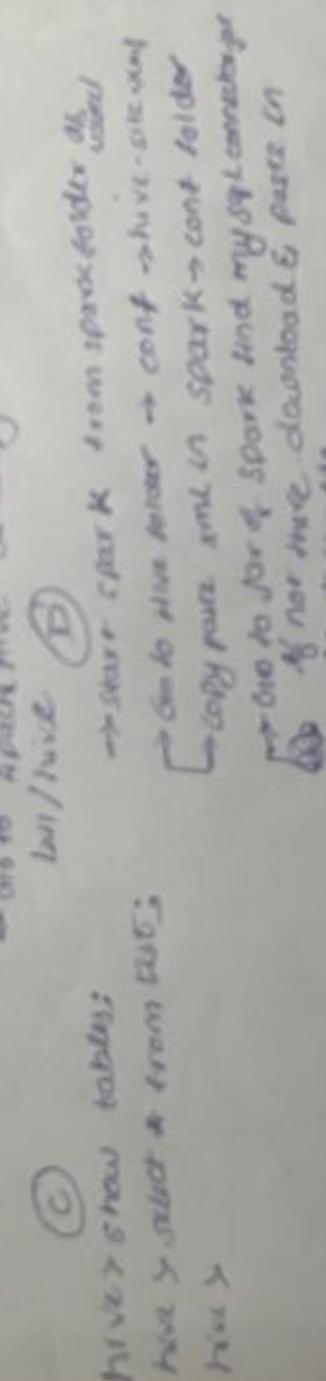
Apache Hive

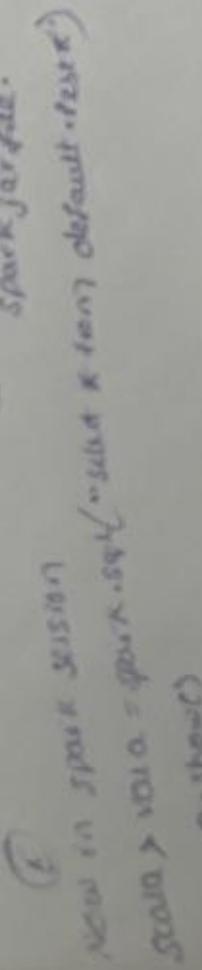
Big Data Tools

- gets a query engine. 
 - ↳ SQL → Abstraction of MapReduce.

- give always sees metadata in RDBMS, actual data is stored in HDFS
- by default has an embedded RDBMS called Derby.

- CLIE → 
 - ↳ derived from CLIE & STAGE
 - ↳ Read file from HDFS from Spark (Perform The ETL)
 - ① External Hadoop → set all variables in bashrc
 - ② External HDFS → start bashrc & source bashrc
 - ③ External HDFS → start hadoop
 - ④ External Spark → Check all 5 command line standard using JPS
 - end \$1\$
 - Go to Apache HIVE directory & start hive bin/hive

- (C) 
 - Hive > create table;
 - Hive > select * from table;
 - Hive >
 - Start SPARK from spark folder then
 - Go to HIVE folder → conf → hive-site.xml
 - ↳ copy path and in SPARK → conf folder
 - Go to jar of SPARK find my sql connector
 - ↳ If not there download & paste in SPARK jar file.

- (D) 
 - New in SPARK SQL
 - scala> val a = spark.sql("select * from default.table")
 - a.show()

(Copper Conductor)



HAVELLS

CABLES THAT DON'T CATCH FIRE*

② HBase

- **HBase** → Hbase is a database management system
- Hbase runs on top of HDFS
- It can store huge amt of data in tabular format for extremely fast read & writes.
- Fast mostly where frequent writing is required.
- Hbase is a type of NoSQL DB

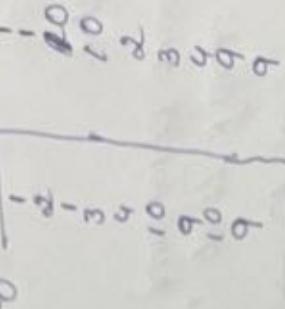
Flame Retardant Cable
1100 Volt PVC Insulated
(Copper Conductor)

HAVELL

Window

val windowspec = Window.partitionBy("ELID").orderBy(desc
df=df.withColumn("SNo", row_number().over
WindowSpec))

row_number → takes the serial numbers of a table or df.
partitionBy → If multiple ELID are there, they will be kept as
it is but in the next consecutive row
Eg: ELID / After partitionBy
1 2 3 4 5 6 7 8 9 10 10 10 10 10 10 10 10 10 10 10



Cache & Persist

- Optimization techniques for interactive & interactive SPARK applications to improve performance.
- cache & persist store the intermediate computation of an RDD Dataframe & Dataset
- cost efficient, Time efficient & Execution time reduces.

Cache

- spark DataFrame & Dataset caching default saves it to storage level "MEMORY-AND-DISK"
- RDD.cache() → "MEMORY-ONLY"

Persist

- It can take storage level argument
- default "MEMORY-AND-DISK"
- def persist(storageLevel, MEMORY_ONLY)
- df.persist(storageLevel, MEMORY_ONLY, 2, MEMORY_ONLY, MEMORY_ONLY, MEMORY_ONLY, DISK_ONLY, MEMORY_AND_DISK_SER, MEMORY_AND_DISK, 2)
- import org.apache.spark.storage.StorageLevel

Unpersist
remove the storage

storageLevel	space used	time used	memory used	disk used	serializable
None	0	0	0	0	0
1	1	1	1	1	1
2	2	2	2	2	2
3	3	3	3	3	3
4	4	4	4	4	4

Recompute some partitions
Y/N

→ All new persist based on parameters provided based on mentioned persist offered level mentioned gives me storage level mentioned

Create a quick dataframe.

val schema =

Partitions & coalesce

- val df = spark.range(0, 100)
- print(df.^{partitions.size})
df.partitions.size
→ df.^{partitions.size}.addPartition. length
+ df.addPartition.^{length} → 8
- df.addPartition.^{length} → 9
- df.coalesce(7)
or
5
- df.coalesce(10, col("id"))
→ repartition based on col("id")
 - repartition it to 10 only.
 - limit shuffle & can only be 10 partitions
- coalesce doesn't involve shuffle or repartition
- need to decrease the number of partitions
- and to display of partitioned dataset^{"dataset"}
- to check shuffle & display of partitioned dataset^{"dataset"}.csv
- df.write.mode(SaveMode.Overwrite).csv("partition.csv")
need to show.

Comparison Operators.

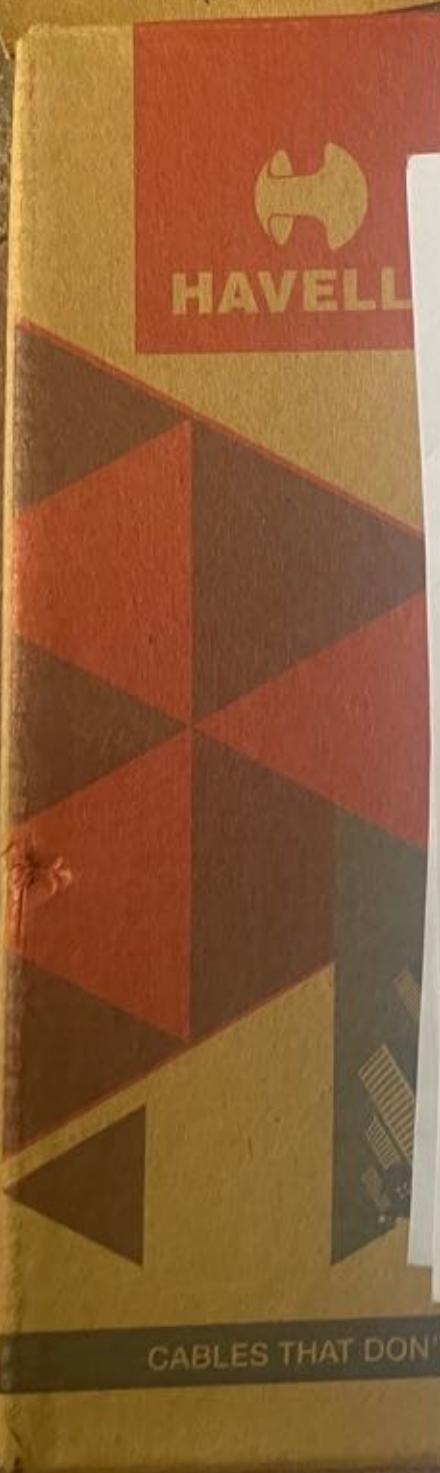
- • `equalTo(F8987)`
- `contains("POSTAGE")`
- `where(col("InvoiceNo") == 536365)`
- `where("InvoiceNo" = 536365)`
- `where("Invoice < 536365")` → other than 536365 it gives its values.
- `where("Invoice > 536365")` → shows above & below values.
- `col("UnitPrice") > 50 & "UnitPrice" <= 100` → including that value above & below 5 values.
- `col(" ") . isin(["DOT", "WHEAT", "PUMPKIN])`
- `and`
- `withColumn("ABC", expr("NOT UnitPrice & UnitPrice <= 250"))` not
- `expr("eq1 & eq2 & eq3 & eq4 & eq5 & eq6")`
- `eqNullSafe("Hello")`
- `eqn(250)`

Spark Streaming

```
→ val streamingDataFrame = spark.readStream
    .schema(schema.ex)
    .option("maxFilesPerTrigger", 1)
    .option("maxFilePerTrigger", 1)
    .format("csv")
    .option("header", "true")
    .option("sep", ",")
    .load("./path")
```

```
// to check if df is still streaming
→ streamingDataFrame.isStreaming // p true or false
```

```
→ streamingDataFrame.format("memory") // memory
    .option("checkpointLocation", "/tmp/checkpoint") // store in-memory
    .writeStream
    .outputMode("complete")
    .queryName("customer-purchases")
    .option("numPartitions", 1) // the name of the in-memory table
    .outputMode("complete") // complete = all the counts should be in the table
```



community cloud Databricks

Notepad Shortcuts.

- ① `hdfs ls /userstore/cabbar/ promod Demo/`
→ lists all the data set files under this path with size

HDFS
file scale
ing.

stem!
using
part

l
num

15

,

12

Flame Retardant
1100 Volt PVC Insulated
(Copper Conductor)



CABLES THAT DON'T

Spark with Cloud.

Sp① Reading from amazon S3 bucket
var df=spark.read.json ("s3://apache-spark-data
/consumers.json")

• R1
you
can
do

Flame Retardant Cable
1100 Volt PVC Insulated
(Copper Conductor)

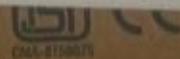


Spark RDD

Resilient Distributed Datasets.

- Fundamental data structure of spark
- RDD is an immutable distributed collection of elements of your data, partitioned across nodes in your cluster, that can be operated in parallel with a low-level API that offers transformations & actions.

CABLES THAT DON'T CATCH FIRE™



CABLES THAT DON'T CATCH FIRE™

Flat
1100 Volt PVC Ins.
(Copper Conductor)



HAVELLS



CABLES THAT DON'T CATCH FIRE®

Spark ML

cmd 14

PP 46-51
Spark diffuser
guide

→ impost org.apache.spark.ml. feature. stringIndexer



UL

CABLES THAT DON'T CATCH FIRE®



Apache Spark Interview questions.

① How is APSP different from Map reduce.

Map Reduce

Apache Spark

- ② spark stores data in memory @ Map Reduce it has to iterate.
- ③ 100 times faster than map Reduce
- ④ provides caching & different types of storage
- ⑤ provides persist level of storage using batches
- ⑥ spark process data in batches as well as in real time.

longer processing

data is highly disk dependent

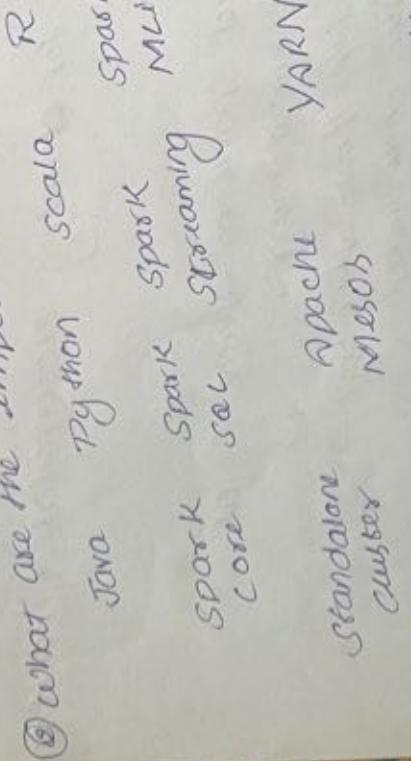
process data in batches only.

process data in batches!

language support

components of spark ecosystem!

support



③ Explain how SPARK runs applications with the help of its architecture.

of

spark on cluster



CABLES THAT DON'T CATCH FIRE™

④ What are different cluster managers available in Apache spark

Hadoop YARN

Mesos YARN

Apache Spark Mesos

Standalone Node

Flame Retardant PVC Insulated Cables

- ⑤ what is lazy evaluation in spark
transformations in spark are not evaluated until you perform an action which acts in optimizing the overall data processing workflow.
- ⑥ what makes spark good at low latency workloads?
use graph processing & machine learning?
use graph stores data in-memory for faster processing
- Apache spark stores data in-memory for fast iterations
 - & model building.
 - Machine learning algorithms require multiple iterations all
 - Machine learning algorithm → Graph algorithm traverse all to create an optimal model
 - The nodes & graphs.
- ⑦ Automatic cleanups for cleaning accumulated metadata
- spark.cleaner.ttl
- ⑧ How can you connect spark to Apache Mesos.
Mesos is a open source cloud manager.
Mesos is a open source cloud manager.
<https://mesos.apache.org/>
- Step1: Install mesos master on the cluster
Sup1: Start mesos master
./bin/mesos-start.sh
- Step2: Start mesos agents on all the worker nodes
./bin/mesos-agent --cmd
- Step3: Make sure mesos is up & running with all worker nodes
configured
<http://localhost:5050>.
- Step4: Spark binary package accessible to mesos.
Install spark in the same location as that of mesos.
- Step5: Configure the
spark.mesos.executor.memory
pointing to spark location.

⑨ what is parquet file & what are its advantages.

what are different file formats. (parquet, avro ..)

+ Parquet (par) is a columnar format that is supported

+ parquet data processing systems.

by several data columns the columns are

→ you have a columnar file & most of the storage.

→ If you have multiple parquet type of storage
not in use we can go for parquet operations claim parquet file format &

both read & write operations

Spark performs

advantages of parquet

- ① Able to fetch specific columns for access
- ② Consumes less space
- ③ Allows type specific encoding.
- ④ Limited IO operations.

CSV → comma separated values
JSON → javascript object notation
ORC → Optimized Row columnar
Labeled ORC
AVRO → tree include markers
that can be used to split large
datasets into subsets suitable
for Apache MR processing.

⑩ what is shuffling in spark?

when do it occur?

→ unless partition is set manually, spark's sub partition board
on the size of the dataset or on available cores
is the process of redistributing data across executors
→ Shuffling may lead to data movement across partitions
partitions that may lead to data moving
partitions while joining 2 tables or while performing
join operations

Bykey operations

⑪ use of coalesce in spark

partition A: 11, 12
partition B: -
partition C: 6, 7
partition D: 9
partition E: 6, 7
partition F: 9, 10
partition G: 11, 12
partition H: 6, 7, 9
partition I: 6, 7, 9
partition J: 6, 7, 9
partition K: 6, 7, 9
partition L: 6, 7, 9
partition M: 6, 7, 9
partition N: 6, 7, 9
partition O: 6, 7, 9
partition P: 6, 7, 9
partition Q: 6, 7, 9
partition R: 6, 7, 9
partition S: 6, 7, 9
partition T: 6, 7, 9
partition U: 6, 7, 9
partition V: 6, 7, 9
partition W: 6, 7, 9
partition X: 6, 7, 9
partition Y: 6, 7, 9
partition Z: 6, 7, 9
coalesce
remove all
partitions
nos which are
multiple of 10



Q) How can you calculate the executor memory?

Eg: Consider Nodes = 10

Each Node has = 16 cores - 1 for OS

Each Node RAM = 64 GB Ram - 1 for OS.

No of executors = No of cores / concurrent tasks

a)

$$= \frac{16}{5}$$

= 3

b)

$$\text{Total no of executors} = \frac{\text{No of Nodes}}{5} * \frac{\text{no of executors in each node}}{3}$$

$$= 10 * 3$$

$$= 30$$

c)

Q) What are the various functionalities supported by Spark Core

A)

→ spark core is the engine for parallel & distributed processing of large datasets.

- ② scheduling & monitoring jobs
- ⑥ memory management
- ③ fault recovery.
- ④ task dispatching.

Q) What is the significance of RDD

- Building blocks of spark
- immutable, fault-tolerant distributed collection of objects that can be operated in parallel.
- whenever you do transformation on RDD is getting reward.

Q) How do you convert RDD into DataFrame.

→

• toDF
def createDataFrame(RDD, schema: StructType)

Q) Does spark provides an API
→ Yes Apache spark provides an API for making streaming applications resilient to failures.
→ checkpointing is the process of making streaming applications resilient to failures.

16) Transformations

ACTIONS.

- Transformations are operations that are performed on an RDD to create a new RDD containing the results
- eg: map, filter, join, union
- Eg reduce, first, count

18) What is a lineage graph

- Spark does not support data replication in the memory
- If any data is lost it can be rebuilt using RDD lineage
- so it is also called an RDD operator graph
- It is a lineage graph.

→ DAG is a lineage graph.

Q) What do you understand by DStreams in Spark?

- Discretized streams is the basic abstraction provided by Spark Streaming
- It represents a continuous stream of data that is either in row or rdd form.



Q) What is the need of broadcast variables in spark?

- Q1) What is the need of broadcast variables in spark?
- Reduces shuffling.
- Increases speed.
- Since a chunk of data stored variable is stored in all the workers.

- Q2) filter/take) returns a new DStream by selecting only those records

- Q3) Does Apache spark provide check points
- Yes Apache spark provides an API for adding & managing checkpoints sufficient to checkpointing is the process of making streaming applications resilient to failures.

Q5) Data types supported by MLlib Spark

- a) what do you mean by saving window operation
 - controlling the transmission of data packets between multiple computer networks is done by sliding window
 - Spark Streaming library provides windowed computations
 - **ignores** where the transformation on RDDs are applied over a sliding window of data

Q6) Difference b/w
map()

- Transformation applies to every element
 - applies to every element
 - will have more rows
- function to each row in a dataset/pandas & returns new transformed dataset
- will have same no of rows as i/p

Q7) What are the different MLlib tools available in Spark?

ML Algorithms
→ classification, regression, clustering & collaborative filtering.

Feature generation
→ feature extraction, transformation, dimensionality reduction & selection

Pipelines
→ tools for constructing, evaluating & running ML pipelines.

Persistence
→ saving & loading algorithms, models & pipelines

Mathematics
→ linear algebra, statistics & data handling

→ linear algebra, statistics & data handling



CABLES THAT DON'T CA

Q5) Data types supported by MLlib Spark.

- ① Local vector
- ② Labeled point
- ③ Local matrix
- ④ Distributed matrix

(Q6) What is a sparse vector in spark MLlib
Ans) What is a sparse vector which is represented
sparse vector is a type of local vector which is represented
by an index array & a value array.
 $\text{sparse} = \text{SparseVector}(4, [1, 3], [3, 0, 1, 0])$

where
• 4 is the size of the vector
• [1,3] are the ordered indices of the vector
• [3,0] are the value

Understand

→ Spark GraphX
→ Read spark MLlib.