

SPARK Architecture

- * Spark is a fast, distributed and general purpose cluster computing system
- * Spark offers - lazy computations - optimize job before execution
 - In memory data caching
 - Efficient pipelining
- * Spark Applications run as independent sets of processes on a cluster, coordinated by the SparkContext object of our main program (also called driver program.)

TWO Main Abstractions in Spark

① Resilient Distributed Datasets (RDD)

- * collection of data items split into partitions & stored in the memory of worker nodes of the cluster.
- * its the interface for doing data transformations.
- * RDDs can store certain data from HDFS, HBase, Cassandra or in cache (memory, mem + disk) or another RDD.
- * RDDs automatically recover from node failures.
- * Partitions are recomputed on failure or cache eviction
- * Meta Data stored in RDDs consists of
 - o Partitions - sets of data splits associated with this RDD.
 - o Dependencies - list of parent RDDs involved in the computation
 - o Compute - function to compute partition of the RDD given parent partitions from the dependencies
 - o Preferred Locations - where is the best place to put computations on this partition (data locality)

RDD Operations

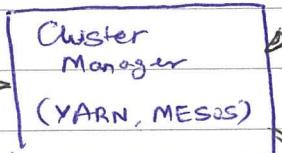
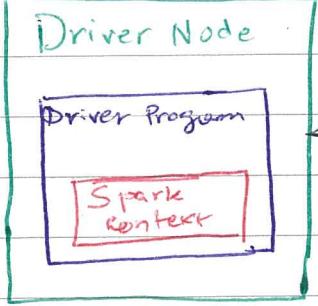
- ① Transformations - do not evaluate immediately (lazy evaluation)
- return new (transformed) RDD
 - only cause change in metadata.

- ② Actions - trigger immediate execution of related previous transformation
- return result to the driver program.

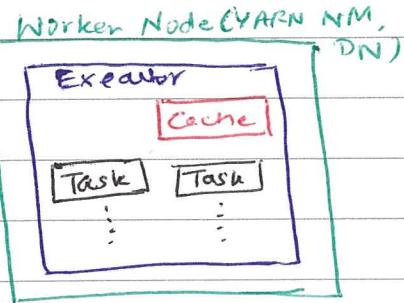
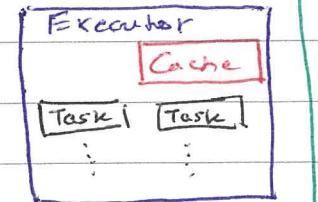
(2) Direct Acyclic Graph (DAG)

- * Sequence of computations performed on data.
- * Node \rightarrow RDD partition
- * Edge \rightarrow transformation
- * Acyclic \rightarrow graph cannot return to the older partition
- * Direct \rightarrow transformation is a transition from A to B for a data partition

SPARK Architecture



Worker Node (YARN NM, DN)



* Driver

- Entry point to the Spark shell (Scala, Python, R)
- Is the class that has the main function.
- Place where Spark Context is created
- Translates RDD into the execution graph.
- Schedules tasks & controls their execution. Due to this reason the driver should run close to worker nodes.
- Stores metadata about ^{all} RDD & partitions..

* SparkContext

- Main entry for Spark functionality.
- A SparkContext represents the connection to Spark cluster. & is used to create RDD, broadcast Variables & accumulators.
- Only one SparkContext may be active per JVM.

* Executor

- Each Spark application will get its own executor process which stays up for the duration of the application and runs tasks in multiple threads.

This results in isolation on both scheduling side (driver schedules its own tasks) and executor side (task from different application run in different JVMs)

- * Stores data in cache in JVM heap or Hard disk
- * Read data from external sources
- * Writes data to external sources
- * Performs all the data processing.

Application :

- * Single instances of Spark Context that stores some data processing logic and can schedule series of jobs, sequentially or in parallel
- * Job is a complete set of transformations on RDD that finishes with an Action or data savings, triggered by the driver application
- * Stage - set of transformations that can be pipelined and executed by single independent worker. Usually it is the transformations between "read", "shuffle", "action" "save".
- * Task - execution of the stage on a single data partition.
- * Spark Application is agnostic of the cluster manager. As long as it can acquire executor processes and these processes can communicate with each other it can run on a cluster which support other non-spark applications

RDD Operations

Transformation

- return RDD
- lazy evaluation

Actions

- return result to driver
- immediately cause the previous op (transformation) to be executed.

RDD Persistence

- By default each transformed RDD will be evaluated everytime an action gets subsequently executed.
- We can persist an RDD in memory using `persist()` or `cache()` method in which case Spark will keep the RDDs in memory around the cluster for much faster access the next time performs operations on this RDD.
- Caching is a key tool for iterative algorithms
- The `cache()` method allows only default storage level (memory-only)
- The `persist()` method allows to specify storage level

~~MEMORY ONLY - stored as serialized Java objects in JVM~~

~~MEMORY AND DISK - Serialized Java objects in JVM. If an RDD does not fit in available memory store partitions that don't fit onto the disk.~~

~~MEMORY ONLY SER - stored as serialized Java objects as an array of byte (one byte array per partition)~~

~~MEMORY AND DISK SER - spill partitions that don't fit onto disk.~~

~~DISK ONLY - only stored on Disk (not in mem)~~

~~MEMORY ONLY 2~~

~~MEMORY AND DISK 2 - keep 2 copies of each partitions across nodes.~~

SPARK PROGRAM EXECUTION (SCALA)

Command Line

> Spark-shell

parameters

- master MASTER_URL Values "client" or "local"
"yarn", spark://host:port
- class CLASS-NAME Application main class
- name "Name" Application name
- jars JARS CSV list of local jars to include in driver & executor class path.
- py-files PYFILES CSV list of zip, egg or py files to place on PYTHONPATH for python apps
- files FILES CSV list of files to be placed in working directories of each executor
- conf PROP=VALUE Arbitrary Spark-config property
- properties-file FILE Path to file from which to load extra properties.
By default uses /etc/conf/spark-default.conf

YARN-only parameters

- archives ARCHIVES CSV list of archives to be extracted into working directories of each executor

Spark-shell in Interactive mode

- we don't specify --class & jar file.

We can type in Scala Statements for Spark or in general on the scala prompt.

Spark-shell in Batch mode

There are two choices in Batch mode.

- ① Spark-Shell
- ② SBT (Scala build tool)

→ Common Steps

* Create Directory Structure

```
.
./
./simple.sbt
./src
./src/main/scala/Yourprogram.scala
```

Simple.sbt

name := "Project Name"

version := "1.0"

scalaVersion := "2.10.4"

libraryDependencies += "org.apache.spark" % "spark-core" % "1.2.1"

YourProgram.scala

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark.SparkConf
import org.apache.spark.RDD._

object YourProgram {
    
```

def main(args: Array[String]) {

Your Statement
For Val conf = new SparkConf()
.setAppName(" ")
.setMaster("local")

}

}

Val sc = new SparkContext(conf)

// Package the program into Jar

go to the folder that contains Simple.sbt (* .sbt)

> ~~sbt~~ package sbt package. - This will create a target folder

Run applications using Spark-submit

> spark-submit --class "Your main class"

-- master "local[4]"

target / Scala / * - jar

{ your location of jar }

Running application using Spark-shell

> spark-shell --class "Your main class"

-- jars "path to jar"

5th S&B (

Developer → Architect (ACS → NS)

Canada (Lianne)

US (Lianne)

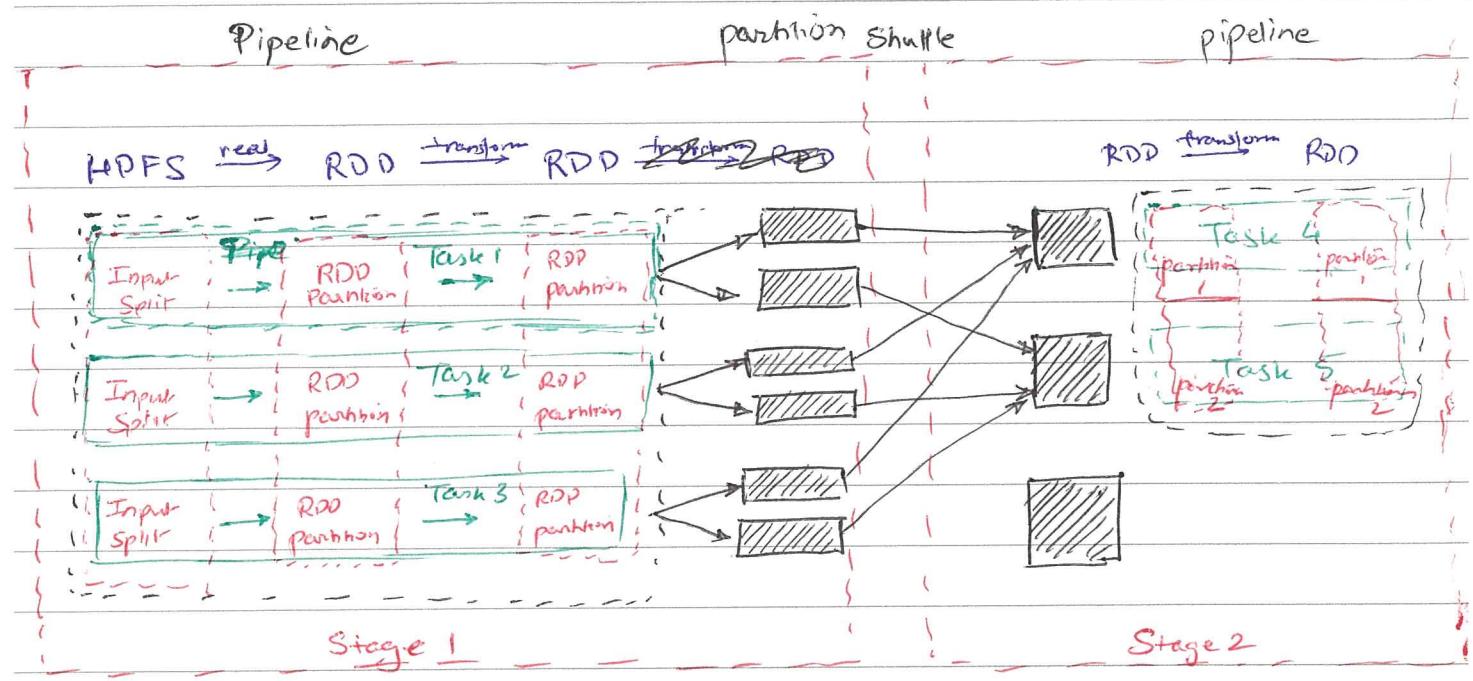
Mobile solution (@newton/)

Liesl

Bekar

NTIDart

New Contact (Empowered with new designation & salary review.)



RDD persistence

- * By default, each transformed RDD will be evaluated everytime an action gets subsequently executed.
- * We can persist an RDD in memory using `cache()` or `persist()` method.
- * Caching is key tool for iterative algorithms
- * The `cache()` method only supports in-memory persistence.

1 MEMORY-ONLY - deserialized Java objects in JVM. if the RDD does not fit in mem, some partitions will not be cached & will be recomputed on the fly. This default.

2 MEMORY-AND-DISK - if RDD does not fit in memory, store partitions that don't fit on the disk & read back when needed

3 MEMORY-ONLY-SER - store as "serialized Java Objects" (one byte array per partition)

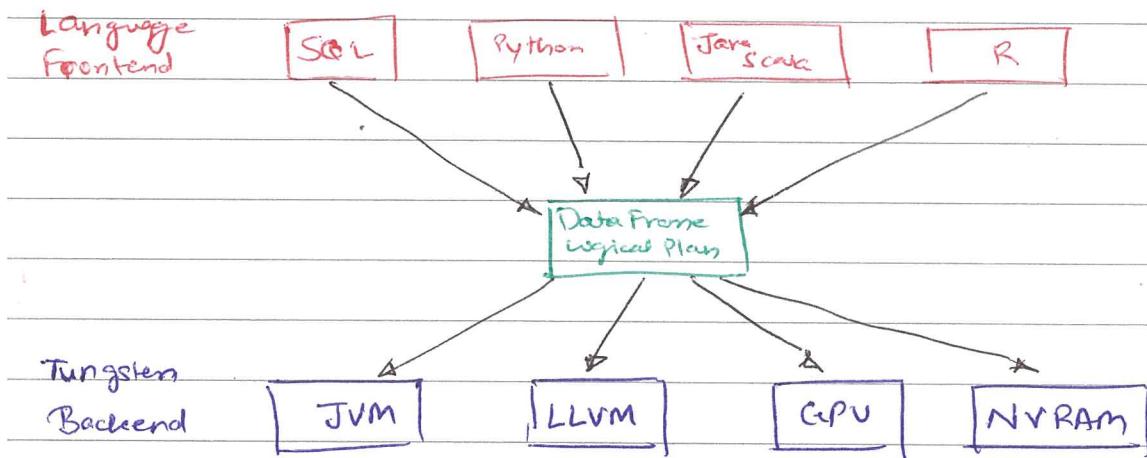
4 MEMORY-AND-DISK-SER - same as 3 above but uses disk

5 DISK-ONLY

6 MEMORY-ONH-2 } Same as top/last but replicate each
7 MEMORY-AND-DISK-2 partition on two cluster nodes

- * Spark considers memory as cache with LRU eviction rules
- * if disk is involved, data is evicted to disk.

SPARK Dataframes :



* Data frame is an RDD with schema - field names, field data types and statistics

* Unified transformation interface in all languages

* Can be accessed as RDD, in this case transformed to RDDs of Row objects.

* Data Frames allows better compression ratio than standard RDD.

* Each column in each partition stores min-max values for partition pruning.

Shared Variables

- * can be used for parallel operations
- * By default when Spark function runs in parallel across nodes it ships a copy of each variable used in the function to each task.

- * Sometimes however we need to share variables across tasks. hence we use `shared variables`

- * TWO Types

(1) Broadcast Variables

- * Used to cache a value in memory of all nodes
- * created by calling `broadcast` function on `SparkContext`
`val bvar = sc.broadcast(v);`
- * After creation it should be used instead of original variable 'v'

(2) Accumulator :

- * variables that are only added such as counters or sums.
- * Created with initial value (v)
`val acc = sc.accumulator(v);`
- * Tasks running on cluster can ~~call~~ add() or
+ = operator on the acc.

SPARK API (Scala)

Important Namespaces

- * org.apache.spark [SparkContext, SparkConf]
- * org.apache.spark.rdd - (Root)
- * org.apache.io - (Compression)
- * org.apache.Spark.sql - (SQLContext)
- * org.apache.spark.sql.hive - (HiveContext)

Important classes & their methods

SparkConf - configuration of Spark Application
- Created using new SparkConf()

getAll(): gets all config settings as K,V

get(String): String - gets the value of specified config property

getDouble(s), getLong(c) - typed methods for config properties

setAppName(string) - set name of our spark application

setJars(Array[String]) - set jar files to distribute to cluster

setMaster(string) - master URL to connect to

Ex: "local" to run locally in one thread

"local[4]" to run locally in 4 cores

"yarn" to run on YARN cluster

SparkContext :

- * Main entry point for Spark functionality.
- * Represents the connection to Spark cluster
- * used to create RDD, Acco., & broadcast var.
- * only one spark context per JVM.
- * created using new SparkContext (appName, sparkConf)

* addFile(path) - add file to be downloaded with this Spark Job
on every node

* addJar(path) - adds a dependency JAR for all tasks to
be executed on this context in the future

* appName - get App Name

* broadcast(obj) - create broadcast variables

* accumulator(obj) - create accumulator.

* getConf() - returns configurations (sparkConf) associated with
this context,

* new API Hadoopfile (path, class [InputFormat], class [k], class [v])

ex: import org.apache.hadoop.mapreduce.lib.input -

val depseq = sc.newAPIHadoopfile[IntWritable, Text, SequenceFileInputFormat]
(outputPath,
classOf[SequenceFileInputFormat[IntWritable, Text]],
classOf[IntWritable],
classOf[Text]);

* sequencefile (path, class [k], class [v]) - reads sequence file

eg: val depseq = sc.sequencefile[Int, String]

("/user/joudheri/datasets/dep-seq");

* `textFile(path)` - reads text file from HDFS/LFS

* `wholeTextFiles(path)` - reads directory of text files from HDFS/LFS
returns `RDD[String, String]`
 \downarrow file path \downarrow file content

Important hadoop namespaces

`org.apache.hadoop.mapreduce.lib.input.XXXInputFormat`

`org.apache.hadoop.mapreduce.lib.output.XXXOutputFormat`

`org.apache.hadoop.io.` { all Writable classes & interfaces }

`org.apache.hadoop.io.compression` { all compression writers }

RDD Operations

cache() - persists an RDD in with default storage level (MEM_ONLY)

coalesce(n) - returns new RDD with specified number of partitions. Usually used to reduce number of partitions after filter operation.

collect(f), collect() - returns new RDD by applying specified function
→ returns ^{array} RDD containing all RDD elements

context - returns SparkContext that the RDD is created on

count() - number of elements in this RDD.

CountByValue() (Ordering: fn)

- returns count of each unique value in this RDD as a map.
- For single fields it's a unique value
- For tuples its every unique combination of fields.

distinct(n) = n - num of partitions

- return RDD containing distinct elements

filter(f: (T) ⇒ Boolean) : returns RDD containing elements that satisfy the specified predicate

first() - return first element

flat

flatMap [U] (f : (T) \Rightarrow U)

- Returns an RDD by first applying the function to all elements of the RDD & then flattening the results.
- This requires the predicate to return a collection instead of single element.

ex : Wordcount = buck.flatMap (line \Rightarrow line.split(" "));

fold (zeroValue : T) (op : (T, T) \Rightarrow T)

- Aggregates the elements ~~of each~~ partition & results of each partition
- Similar to reduce but allows specifying zero value which gets passed as first argument to predicate during first call ~~in~~ ~~to~~ each partition

foreach (f : (T) \Rightarrow Unit)

- Applies predicate to all elements of RDD.
- This is an action & not transformation
- It does not return anything
- used for side-effects like printing elements of RDD

groupBy ($f : (T) \Rightarrow K$) returns $RDD[(K, Iterable(T))]$

- groups a tuple by the specified ~~field~~ predicate
(the predicate returns a field which becomes the grouping key)
- this function creates a nested structure
- After this function, we can use mapValues for doing aggregations (max, min, sum)
- After this function, we can use flatMapValues for doing sorting / ranking.

intersection (other: RDD)

- returns an RDD which contains common elements to this and other RDD

keyBy [K] ($f : (T) \Rightarrow K$) returns $RDD[(K, T)]$

- creates a tuple of K, V where K is based on field or fields that^{is} returned by the predicate

ex: depRDD.keyBy($x \Rightarrow x.dep_id$)

will return $RDD[(dep_id, dep_Record)]$

* $\text{map}[U](f: (T) \Rightarrow U)$ returns $\text{RDD}(U)$

- returns a new RDD by applying the predicate to elements of this RDD .
- used for doing transformations on whole RDD (not by key)
- This is a **PUSH MODEL** where elements in source RDD are pushed one by one into the predicate function call.

$\text{mapPartitions}[V](f: \text{Iterator}(T) \Rightarrow \text{Iterator}(U))$

$\text{mapPartitionsWithIndex}[U](f: (\text{Int}, \text{Iterator}(T)) \Rightarrow \text{Iterator}(U))$

- Similar to map but applies the function at partition level instead of each element.
- All the elements of a partition are passed in as an Iterator through which we can traverse
- $\text{mapPartitionsWithIndex}()$ additionally passes the partition number.
- These functions are meant for doing any initialization at partition level (instead of each element, can enhance performance)
- Since these functions pass in an Iterator (instead of element as in $\text{map}()$) and we have to put elements within the predicate this is a **PULL MODEL**
- These functions convert all elements from U to V and return an $\text{Iterator}(V)$. (\Leftrightarrow which means we can also have logic to drop certain elements)

`max()` (Ordering func optional)

`min()` (Ordering Func optional)

`persistent (newLevel)` - sets RDDs storage level to persist value across operations

`unpersist()` - make RDD non-persistent & remove all blocks from memory & disk

* reduce($f: (\tau, \tau) \Rightarrow \tau$)

* used to perform overall/global aggregations (sum, min, max)

* can be applied to single elements, tuples etc.

* cannot be used for "per key" aggregations

ex: Find max price of product

`products.reduce((x, y) \Rightarrow if (x.price > y.price) x else y));`

`repartition (numPartitions)` - return new RDD that contain exactly numPartitions

`sameAsObjectFile (path)`

`saveAsTextFile (path)`

`saveAsNewAPIHadoopFile (path, class [key], class [value])`

`class [outputFormat [key, value]]]);`

ex: `import org.apache.hadoop.mapreduce.lib.output.`

`depRDD.saveAsNewAPIHadoopFile (outputPath,`

`classOf [IntWritable], classOf [Text],`

`classOf [SequenceFileOutputFormat [IntWritable, Text]]);`

* SortBy ($f : (T) \Rightarrow K$, $asc = true$, $numPartitions = default$)

- * used for global sorting of RDD by any field.
- * can be used with single elements, tuples etc
- * predicate returns key to sort on)
- * can be used to sort on multiple fields

ex: Sort products by category-id in ASC (2nd field)
and Price in DESC (5th field)

products.sortBy($x \Rightarrow (x._2, -x._5)$)

subtract (other: RDD)

- returns RDD that contains elements in this, which are not present in other.

RANKING

take(n) - take first n elements from RDD

- this is an action that returns an Array:

takeOrdered (N) (Ordering func)

- first sorts the elements (default ASC & on first field)
or according to the specified ordering function.
and then takes only the first N elements

ex: products.takeOrdered(3); take ^{5 smallest} bottom 3 based on first field.

products.takeOrdered(3)(Ordering[Int].on(x => x._1))

The above two func calls are equivalent.

// Top 3 based on third field in tuple
products . takeOrdered(3) (Ordering[Int]) , reverse . on ($x \Rightarrow x_{-3}$)

// Top 3 based on ordering with multiple columns
first column is Int which is 2nd field (Asc order)
second column is Float which is 5th field (DESC order)

products . takeOrdered(3)
(Ordering[Int, Float] . on ($x \Rightarrow (x_{-2}, -x_{-5})$))

~~* top(N) (Ordering function)~~

- Works in exact opposite way of takeOrdered by sorting elements in DESC order by default & based on first field.
- This can be overridden by specifying Ordering function (either change to Asc order or choose other fields)

// top 3 requires based on sorting on field 2
products . top(3) (Ordering[Float] , on ($x \Rightarrow x_{-2}$))

Union (other : RDD)

- returns the RDD which is union of this with other (Simply appends other to this)

PAIR RDD Functions

* Only available for tuples with K, V pairs.

* reduceByKey (func : (V, V) => V)

* passes in an accumulator & value (one per predicate call) and the predicate should return one element (tuple or single value) but not a collection. (not suitable for sorting, ranking)

* This is similar to reduce but the values will be first grouped by the key and then the values for each key will be passed one-by-one into the predicate function.

* This method is used to do aggregations (sum, max, min per key basis).

* is a transformation

* Uses a combiner whose logic is same as the predicate (reducer function) case hence this method is efficient way of doing per key aggregation.

* The source RDD could be a flat structure of K,V

* Since values per key are passed to predicate one-by-one this is a PUSH MODEL

ex : total OrderAmt Per Order

= orders . reduceByKey
 $((acc, value) \Rightarrow acc + value);$

Assuming orders is RDD [Int, Float]

\downarrow orderId \downarrow subtotal.

/ /

CombineByKey [C] (Combinerfunc : $(acc, val) \Rightarrow Value$),
func : $(acc, val) \Rightarrow Value$
(reducer) C C C

* Similar to reduceByKey but use when the combiner logic is different than reducer function.

* The input ~~and output~~ to the combiner & output of combine & reduce function should be of same datatype

* AggregateByKey (zeroval: U)

(combineOp : $(U, V) \Rightarrow U$,
reduceOp : $(U, U) \Rightarrow U$)

* more flexible than reduceByKey & combineByKey

* used when we need not only a different combiner & reducer logic but also the input type to combiner & reducer to be different.

* also allows to specify a zero value which gets passed into combineOp (# param) per key.

* This zero value should be of the same type that we expect out of the combiner predicate & reducer predicate

ex: Calculate Avg revenue per day

val RevenueAggPerDay =

zero value we are expecting
tuple of (Float, Int)

OrderPerMap - aggregateByKey ((0.0, 0)).

$\{(acc, val) \Rightarrow \{acc_1 + value, acc_2 + 1\}\}$

↑
combineOp

↑ we are accumulating the total in first field & count in 2nd field

$\{(acc, value) \Rightarrow \{acc_1 + value_1, acc_2 + value_2\}\}$

↑
reducerOp
both input parameters

match combiner output ie Tuple (Float, Int)

↑ we are adding up all intermediate totals from first field &

intermediate counts from 2nd field
so acc & value are both tuples

foldByKey (zeroval: V) (func : (V, V) \Rightarrow V)

- * Similar to reduceByKey but allows to specify a zero value
- * the predicate logic is used for both combiner & reducer.

* groupByKey (numPartitions)

- groups a K,V flat structure & returns a key with Iterable (values) i.e nested structure
- normally used for sorting & ranking of values per key.
- does not use combiner hence should not be used for performing subsequent aggregations.

ex: val productsGroupByCat = products.groupByKey();

Assuming products is K-V pair where K is category-id.

mapValues [U] (func: (V) => U)

- * used to perform a function over all values per key.
- * can perform transformation (aggregates) to return single value or collection. In case of returning collection this function results in a nested structure. K, multiple values $[V_1, V_2, V_3]$

- * For each key the values are passed in ~~as is~~ into the predicate.
- * In case the Value V is iterable, its responsibility of predicate to PULL the values.

- * Normally used after grouping of data

ex: val grouped = orders.groupByKey() OR orders.groupBy(_._1)
Assuming K, V where K is order-id &
V is order subtotal

val revenuePerOrder = grouped.mapValues (v => v.sum);

flatMapValues [U] (func: (V) => U);

V could be iterable or single value

- * Similar to mapValues but the predicate should return a collection since this method flattens the result.

- * used mostly for sorting & ranking values within a key.

ex: // Top 3 products per category by price.

val prodGroup = products.groupByKey() OR groupBy(_._1)

// Assume products is K, V where K is prod.cat

V is product tuple

NOTE take(3) is with

val top3ProductGroup = flatMapValues (x => x.toList.sortBy(_._4).take(3))

↓ This will result in a collection of top 3 rewards sorted by price (4th field) DESC

flatMapValues will convert this nested structure K, Iterable (V) into flat-structure K, V₁, V₂, V₃

* Nested Tuples (Tuple inside Tuple) is not considered to be a nested structure. It's still a flat structure.

keys () - returns RDD containing keys of each tuple

values () - returns RDD containing values of each tuple

lookup (key : K) - returns a list of values for a given RDD key
- this is an action

jp.

JOINS

join [W] (other : RDD [K,W]) returns RDD [(K,(V,W))]

* joins two RDDs and returns a tuple containing the key and the corresponding records from this and other as a tuple

* generates a flat structure

* performs an inner join (filters out records from this & other that don't exist).

ex: val joined = orders.join(orderItems)

Assume orders is K,V where K is order-id (int)

V is order record String

orderItems is K,V where K is order-id (int)

V order Items record String

then join will be

RDD : [(Int, (String, String))]
↓ ↓ ↓
key order orderItem
(order-id) record record

leftOuterJoin (other : RDD)

* same as join but performs left outer join

* which means the 2nd field in the value tuple of the result can be empty for missing records in other

* generates flat structure.

rightOuterJoin (other : RDD)

* same as join but performs right outer join.

* the 1st field in the value tuple of the result could be empty for missing records in this.

* generates flat structure.

fullOuterJoin (other : RDD)

- * performs a full outer join
- * the 1st or 2nd field in the value tuple of the result will be empty for missing records in this OR other respectively.
- * generates flat structure

* ~~cogroup~~ (also called groupwith)

[W] (other : RDD[(K, V)]) : RDD[(K, (Iterable(V), Iterable(W)))]

- performs a full outer join but generates a nested structure because it also groups the result after the join by the key.

ex: val stocks Assume K, V (K is stock symbol & V is stock record) string
val companies Assume K, V (K is stock symbol & V is company record) string

val cogrouped = stocks.cogroup(companies)

Cogrouped will be RDD[(String, (Iterable(String), Iterable(String)))]

Key (stocks symbol) iterator for all stock records for this symbol iterator for all company records for this key.

We can remove nesting by using flatMapValues to combine the two iterables into one Iterable

cogrouped.flatMapValues (

value from 1st iterable
 \Rightarrow for (v1 \leftarrow v._1, v2 \leftarrow v._2)
yield (v1, v2); value from 2nd iterable
Combine them into one tuple.