

Java Pair RDD

sortByKey (asc = true)

* sort the RDD by key so that each partition contains
sorted range of elements

ex = Val products = Assume K,V , K is price
V is product record

Val sortedByPrice = products.sortByKey(); ASC
products.sortByKey(false) DESC

RDD Operations

Transformations - return RDD

- lazy evaluation

Actions - returns result to deliver

- immediately cause the previous op (transformations) to be executed.

RDD persistence:

- By default, each transformed RDD will be evaluated every time an action gets subsequently executed
- we can persist an RDD in memory using persist() or cache() method in which case Spark will keep the RDD 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 (mem-only)
- The persists() method allows to specify storage level

MEMORY ONLY - Stored as serialized Java objects in JVM (Default)

MEMORY-AND-DISK - In addition to memory-only, if an RDD does not fit in available cluster memory, the RDD partitions that don't fit are stored on Disk.

MEMORY-ONLY-SER - stored as serialized Ser objects but as one array of bytes (one array per partition)

MEMORY-AND-DISK-SER - spills partitions that don't fit onto disk

DISK-ONLY - stored on disk (not in mem)

MEMORY-ONLY-2 - keep 2 copies of each partition across nodes

MEMORY-AND-DISK-2 - spill onto disk for partitions that don't fit

SPARK PROGRAM EXECUTION (SCALA)

Command Line (Interactive)

> Spark-shell

Params

-- master	MASTER_URL	"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 PYTHON PATH 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 use /etc/conf/spark-defaults

YARN only parameters

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

Spark - Shell in interactive mode

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

* we can type in scala statements for spark on in general on the scala prompt

BATCH Mode

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, SparkConf}
```

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

```
import org.apache.spark.rdd._
```

```
object YourProgram {
```

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

```
    yourStatement ... val conf = new SparkConf()  
      .setAppName(" ")
```

```
}
```

```
    val sc = new SparkContext(conf)
```

```
  }
```

```
}
```

// Package the program into Jar
go to folder that contains simple.sbt

> sbt package → This will create a target folder

Run application using Spark-submit

> Spark-submit
-- class "Your main class"
-- master "local[4]" ← local mode with 4 threads
target/scala / xxx.jar
{ location of your jar }

Batch Mode using Spark-shell

> Spark-shell --class "Your main class"
--jars "path to jar"

1 1

/ /

* SPARK API (PySpark) *

Pyspark is a python API for spark.

Important classes & methods

* SparkContext

* Constructor : SparkContext (master = None , appName = None ,
pyFiles = None , environment = None ---)

* accumulator (value , accum-param = None)

↓ helper object on how to add values

* addfile (path) - add files that gets sent with this job to
every node where this job tasks run .

* broadcast (value) - broadcast read-only variable to the cluster

* default parallelism - default num of partitions for reduce tasks

* new API HadoopFile (path , inputformat Class , keyClass , valueClass ,
keyConverter = None , valueConverter = None ,
conf = None)]
Can provide lambda func here

keyClass = full qualified eg: org.apache.hadoop.io.Text

valueClass = full qualified eg: org.apache.hadoop.io.LongWritable

conf = Hadoop configuration (None by default)

inputformat Class = org.apache.hadoop.mapreduce.lib.input.
TextInputFormat

parallelize(collection, numSlices = None)

↳ distributes a local python collection to form an RDD.

SequenceFile(path, keyClass = None, valueClass = None,

keyConverter = None, valueConverter = None, minSplIt = None)

→ Reads a sequence file from HDFS or LFS and returns RDD
of K, V

stop(): shuts down the SparkContext

textFile(name, numPartitions = None, use_unicode = True)

→ reads a text file from HDFS or LFS and returns RDD
where each record is a line from file.

→ if unicode = False, strings will be kept as str (UTF8)
which is faster.

union(rdd's) - build the union of a list of RDDs

version - Spark Version

wholeTextFiles(path, minPartitions = None, use_unicode = True)

- reads a directory & returns RDD as K, V where
K is filename, V = content of ^{single} file.

* Lambda functions in python only support single statements

* RDD

* aggregate (zeroValue, combinerOp, ReducerOp)

- allows providing zero value
- allows different logic for combiner & reducer
- allows different data types between input records of RDD & output result
- used for overall aggregations (max, min, sum, avg, count)

aggregateByKey (zeroVal, combineOp, ReducerOp), numPartitions = None

- for key aggregations. (rest are similar to aggregate() above)

e.g.: Revenue And Order Count Per Day =

ordersByDate = aggregateByKey

((0.0, 0), lambda acc, val: (acc[0] + value, acc[1] + 1))

✓ ✓
initial value is tuple CombinerOp ReducerOp.
generates tuple of (sum, count) adds up sums & counts. form new tuple.

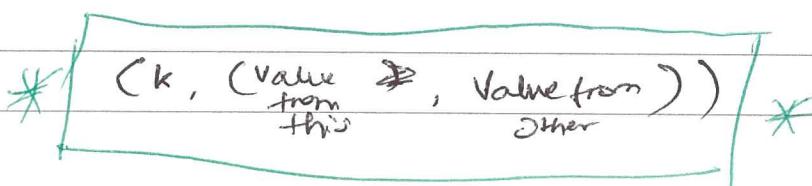
cache() - persists RDD with default storage level (MEMORY_ONLY_SER)

cartesian(other) - returns cartesian product of this RDD with other

coalesce(numPartitions) - returns RDD ~~with~~ ^{that is} reduced to numPartitions

cogroup (other, numPartitions = None)

- for each key in this & other, return RDD that contains tuple with list of values for their key in self as well as other.



* Collect() - returns a list that contains all elements in this RDD

collectAsMap() - returns a K, V RDD into a dictionary.

CombineByKey (zeroValFunc, combineOp, ReducerOp, numPartitions = None)

- Allows specifying a function that gets passed in the first record (from each partition) to generate zero value.
- rest of this are same as aggregateByKey()

Eg:

revAndOrderCountPerDay =

```
ordersByDate.combineByKey(lambda t: (t, 1),  
                           lambda acc, val: (acc[0] + val, acc[1] + 1),  
                           lambda acc, val: (acc[0] + val[0], acc[1] + val[1]))
```

the return value of
this function forms the
Accumulator

Context - the Spark context that this RDD was created on.

Count() - return number of elements in the RDD.

* countByKey() : returns number of unique values per key as a dictionary of (key, count(v))

Eg: numItemsPerOrder = ordersKV.countByKey()

distinct (numPartitions = None)

- returns RDD containing distinct elements in this RDD
- if the elements are tuples then each different tuple (ie any field different) is considered unique.

eg: `CancelledOrders = orders.filter(lambda x: x[3] == "CANCELED")`

* filter (Op) - return new RDD containing elements that match the Op. (ie op returns true)
- always filter as early as possible before any transformation

`first()` - returns first element in this RDD

* flatMap (Op, preserves Partitioning = False)

- returns an RDD by first applying the function & flattening the results.

eg:

`wordCount = book.flatMap`

(`lambda line: line.split(" ")`)

• `map (lambda w: (w, 1))`

• `reduceByKey (lambda x, y: x+y)`

`flatMapValues (Op)`

- parses the value of each key to the function which should return a collection which is then flattened with the key.
- usually used for sorting & ranking after `groupByKey()` or `groupBy()`

eg: `productsCatGroup = products.groupBy (lambda x: x[1])`

`top3PriceByCat = productsCatGroup.`

flatMapValues

(`lambda x: sorted(x, key = lambda i: i[2])`)

↑

These are

all products (list)

for a

category

sort using this function

↓

[0:3]

↑

4th
file
is
in

reverse = True

↑
slice
for take first 3
elements

fold(zeroVal, op)

- aggregate elements of each partition & results of all partitions given the associative function
- Similar to reduce() but can specify a zero value
- used for overall aggregations.

~~foreach (op)~~

foldByKey (zeroVal, func)

- similar to fold but operates on list of values per key.

~~foreach (func)~~ - applies a function to all elements of RDD

- its an Action which does not return anything
- used mainly for side effects (printing on screen)

eg: ~~products.~~ foreach (lambda x: print(x))

getNumPartitions() - gets the ~~number of partitions for this RDD
size of the file to which this RDD~~

getStorageLevel()

glom() - returns RDD created by coalescing all elements within each partition into a list.

* groupBy (func, numPartitions = None)

- returns a RDD with K, V where K is the value returned by the function & value is the element of the RDD.

eg: productsCatGrp = products.groupBy (lambda x : x [1])

group by 2nd element
of tuple resulting is
(key, (tuple))

* groupByKey (numPartitions = None)

- groups a K, V RDD to return an RDD of (K, list(V))

eg:

productsCatGrp = products.map (lambda x : (x [1], x)).groupByKey ()

groupWith (other, *others)

- cogroup with multiple RDDs.

id() - unique id for this RDD.

keyBy (func)

- generates K, V RDD ~~from~~ for each element in source RDD by applying the function

eg:

ordersKV = ordersRDD.keyBy(lambda x: x[0])

↑
generates (k, v) from v
where k is orderid (x[0])

~~map~~ keys() - return RDD with the keys of a this k,v RDD

lookup(key) - returns list of values in the RDD for the given key.

~~* map (func, preservesPartitioning = False)~~

- transforms element U of this RDD into V .
- returns an RDD(V)
- This is PUSH model (elements are pushed one by one into V)

// Transforms Orders Tuple to Orders KV while removing $x[4]$

ordersKV = orders.map(lambda x: (x[0], (x[1], x[2], x[3])))

mapPartitions (func, preservePartitioning = False)

- the function func is applied to each partition.
- The elements of each partition are accessible to the func as iterator.
- This is PULL model.
- Useful if we want to execute piece of code only once per partition. (map() will execute the func per element).

// Transforms Orders Tuple to Orders KV

Sum of orderSubtotals per order

mapPartitionsWithIndex (func, preservesPartitioning = False)

- similar to mapPartitions() but the index of original partition ^{is} also available to func.

~~* mapValues(func)~~

- passes the V of a K,V RDD to the func
- The func can return single value (in which case the resulting RDD will be flat structure (K,U))
but it can also ~~also~~ return collection (in which case the resulting RDD will be nested structure (K, list(U)))

Eg: SomeRDD.mapValues(lambda v: v*x) V is single value
Input - Flat Output - Flat U is single value (v*x)

↓
transform

SomeRDD.mapValues(lambda v: [v,v]) V is single value
Input - Flat Output - Nested U is list of values
↓
transform/explode

SomeRDD.mapValues(lambda v: sum(v)) V is iterator (collection)
Input - Nested Output - Flat U is single value
↓
aggregation

SomeRDD.mapValues(lambda v: for x in v: x+x) V is iterator
Input - Nested Output - Nested U is list.
↓
transform.

~~* max(key=None)~~

- Find max item in this RDD according to the function

Eg: highestProductPrice = products.max(lambda x: x[4])

~~min(key=None)~~

sum(.)

~~persist(storageLevel)~~

- set this RDD storage level.

~~performs~~

~~* reduce(func)~~

- overall aggregation
- no zero value can be specified
- the func serves as both combine & reduce op.

~~* reduceByKey(func, numPartitions=None)~~

- similar to reduce() but does per key aggregation
- values for key are grouped together & the values are passed one by one into the func.

eg:

wordcount = words.map(lambda x: (x, 1))

• reduceByKey(lambda acc, val: acc + val)

repartition(numPartitions)

- returns new RDD that has exactly numPartitions

setName(name)

- Assign a name to this RDD

* sortBy (func), asc = True , numPartitions = None)

- sorts this RDD according to the given func.
- used for global sorting. ~~& ranking~~.

eg: // sort by product price

priceAsc = products.sortBy(lambda x: x[4])

priceDesc = products.sortBy(lambda x: x[4], False)

* sortByKey (asc = True , numPartitions = None)

- similar to sortBy but requires K,V. & sorts by key.
- used for global sorting.

priceAsc = products.map(lambda x: (x[4], x))
 • sortByKey()

priceDesc = products.map(lambda x: (x[4], x))
 • sortByKey(False)

* takeOrdered (num , key = None)

- used for global ranking.
- sorts the elements in Asc and picks only top "num"
- can specify a sort key function.

// get top 3 least priced products

bottom5price = products.takeOrdered(3, key = lambda x: x[4])

OR using K,V

bottom5price = products.map(lambda x: (x[4], x))
 • takeOrdered(3)

* top (num , key = None)

- works exactly opposite of takeOrdered . default DESC
- picks "num"
- used for global ranking.

top5price = products.top(3, lambda x: x[4])

top5price = products.map(lambda x: (x[4], x)).top(3)

~~X~~ take (num)

- takes first num elements in RDD.

unpersist ()

- mark RDD as non-persistent

values () - returns an RDD of V from K, V RDD

union (Other) - return union of this RDD with other

intersection (Other) - return intersection of this RDD with other

~~X~~ join (Other, numPartitions = None)

- returns RDD containing pair of matching elements (ie elements with matching keys.)
- this & other must be a K, V RDD.
- returned RDD as (k, (v1, v2)) tuple

ordersKV = orders.map(lambda x : (x[0], x))

orderItemsKV = orderItems.map(lambda x : (x[1], x))

ordersJoin = ordersKV.join(orderItemsKV)

leftOuterJoin (other, numPartition = None)

fullOuterJoin (other, numPartition = None)

rightOuterJoin (other, numPartition = None)

Save As HDFS New API Hadoop file (path, key Converter = None
Value Converter = None)

K Save As New API Hadoop File (path, outputformat Class , key Class = None
Value Class = None , key Converter = None
Value Converter = None)

outputformat Class = org.apache.hadoop.mapreduce.lib.output.

SequenceFileOutputFormat}

keyClass = org.apache.hadoop.io.Text

value Class = org.apache.hadoop.io.IntWritable.

Save As Sequence File (path , compression Codec Class = None)

Save As Text File (path)

SparkConf

contains(key)

get(key, defaultValue = None)

getAll() - returns list of KV pairs config settings

* setAppName(name) - set app Name

* setMaster(value) - set master URL

PACKAGE PySparkSQL

Class StructField (name, datatype, nullable = True)

~~* class SQLContext (SparkContext, sqlcontext = None)~~

- main entry point for SparkSQL functionality

- SchemaRDD - RDD with schema (Row objects)

applySchema(rdd, schema)

→ applies the given schema to an RDD.

schema = StructType([StructField("field1", IntegerType(), False),
StructField("field2", StringType(), False)])

rdd = sc.parallelize([(1, "One"), (2, "Two")])

SchemaRDD = sqlContext.applySchema(rdd, schema)

sqlContext.registerRDDAsTable(schemarDD, "table1")

This is method on sqlContext to create table

sqlContext.sql("Select * from table1").collect()

cacheTable (tableName) - cache table in mem.

inferSchema (rdd)

- infer & apply schema to an RDD of Row

OrdersRDD = orders.map(lambda x :
Row(order_id = int(x[0]),
order_date = x[1],
order_cust_id = int(x[2]),
order_status = x[3]))

OrdersSchemaRDD = ~~ordersRDD~~ SqlContext.inferSchema(OrdersRDD)

orderschema RDD.registerTempTable ("orders")

↙ sqlContext.sql("select * from orders")

This method is on SchemaRDD to create a table

* jsonfile (path, schema = None)

eg: deptJSON = sqlcontext.jsonfile(path)

deptJSON.registerTempTable("deps")

depRDD = sqlContext.sql("select * from deps")

parquetfile (path)

- loads a parquet file & return result as SchemaRDD

register RDD As Table (rdd, tableName)

- registers an RDD as temp table in catalog

sql (sqlquery)

- returns an SchemaRDD representing results of the given query.

* SQLContext is used to define a structure on top of data in HDFS
It is not Hive aware

* HiveContext knows the structure from Hive metastore.

* Class SchemaRDD

cache(), count(), distinct(), intersection(), limit(), persist(), take(), saveAsParquetFile().

registerTempTable(name)

- registers the RDD as temporary table with given name

limit(nnn)

- limits the result to specified number

toJSON()

- Convert schemaRDD to MappedRDD of JSON documents

* Row - used to create SchemaRDD.

Constructor

Row (name = "Alice", age = 21)

Class HiveContext

- Variant of SQLContext that is aware of Hive metastore
- ie query Hive tables & execute Hive commands like show table, describe formatted etc.
- Also can set hive properties.

Eg: `hiveContext.sql("set hive.execution.engine = mr spark.sql")`

`hiveContext.sql("use retail_db")`

`hiveContext.sql("select * from departments")`

Spark Execution Python

① Interactive Mode - directly type python statements at prompt

> pyspark .

--master MASTER_URL "yarn", "local", spark://host:port
--name AppName
-py-files PY_FILES CSV list of egg, zip or .py files
to place on PYTHONPATH
--files FILES CSV list of files to be placed in
working directory of each executor
--conf prop = Value
--properties-file FILE path of file to load extra properties
default: /etc/spark/conf/spark-defaults.conf
--version

② Batch mode - create python script & submit for execution

SimpleApp.py

```
if __name__ == "__main__":  
    def myfunc(s):  
        words = s.split(" ")  
        return len(words)
```

```
sc = SparkContext(...)
```

```
sc.textFile("file.txt").map(myfunc)
```

> spark-submit --master local[4] SimpleApp.py .

Spark Configuration

Spark.app.name

Spark.master

Spark.driver.extraClassPath

spark.driver.extraLibraryPath

Spark.executor.extraClasspath

spark.executor.extraLibraryPath

Spark.shuffle.compress - true - Whether to compress map output

Spark.ui.port 4040 Port for Application Dashboard

Spark.rdd.compress false Whether to compress serialized partitions

Spark.io.compression.codec snappy

Spark.default.parallelism

Spark.task.maxFailures 4

Spark.speculation False Speculative execution

SPARK SQL

Spark.sql.shuffle.partitions 200 Num partitions when shuffling data.

* Types are always implicitly derived

* Everything is Obj.

SCALA FOR SPARK

* Pass by reference default

Important class & functions.

i: Int = 1

val i = 1 OR val ~~i~~ ~~Int~~ = 1

Val - immutable

Var str = "" OR var str: String = "ABC"

Var - mutable

// Formatted

* printf ("%d %.2f", i, j, k)

// Unformatted

* println (i, j, k)

// Arrays are fixed size but mutable (ie values within the array can be changed)

var myArr = Array(1, 2, 3, 4, 5)

// Access

myArr(2)

// List - size can change (add/remove elements) but cannot change existing value

var myList1 = List(1, 2, 3, 4, 5)

myList1(2) = 5 // Invalid

myList1 = List(7, 8, 9) // Valid - changing object ref.

// Mutable List

var myList2 = collection.mutable.MutableList(1, 2, 3, 4)

myList2(3) = 100 // Valid

myList2 = (7, 8, 9) // Invalid since val type

* for loops

```
for (i ← myList) {  
    println(i)  
}
```

// excludes last number

* for ($i \leftarrow 0$ until str.length()) {
 print(str(i) + " ")
}

* for ($i \leftarrow 0$ to 5) { // 0 to 5 inclusive
 print(i)
}

// Yield -

* var evenlist = for { $i \leftarrow 1$ to 20 if $i \% 2 == 0$ } yield i;
// gets all even values between 1 & 20.

// Printing Simple String or Expressions

* println(s "Hello \$name") // where val name = "Derek"

// Printing with Formatted String or Expressions

println(f" I am \${age}")

println(f" I am \${age+1} and weight \${weight%.2f}")

printf("I am %.d and weight %.2f", age, weight)

Strings: same as Java - `String`

Var `str = "Something important"`

* `str.length`

* `str.concat(" for me")`

* `str.equals("first other string")`

* `str.contains(" something")`

* `str.startsWith()`

* `str.endsWith()`

* ~~String~~ `String.format("%d, %f", age, name)`

`str.indexOf()`

* `str.replace(old, new)`

* `str.split(",")`

* `str.substring(beginIndex, endIndex)`

* `str.trim()`

Arrays:

Var `friends = Array("Bob", "Tom")`

// Access `friends[0]`

Var `friends2 = ArrayList<String>()`; ArrayList can vary in size

`friends2.insert(0, "Mark")` // insert

`friends2 += "Phil"` // append

`friends2 += Array("Suey", "Ben")` // Append an array.

`friends2.insert(1, "Mike", "Sally", "Sam");` // Insert multiple elements at 1

`friends2.remove(2, 3)` // remove elements from 2 to 3

`friends2.foreach(println)` // print

// Transform.

```
val friends3 = for (f <- friends2) yield f.toUpperCase();
```

// Functions available on Arrays

```
val nums = Array(1, 2, 3, 4, 5)
```

* nums.max

* nums.min

* nums.sum

// Sorting

* nums.sortWith(_ > _) // Desc

nums.sortWith(_ < _) // Asc

// to string in CSV

* nums.mkString(", ")

// Mixed Arrays (Any type)

```
val arr = Array("a", 1, true)
```

```
array.find {x => x == 4}.getOrElse(10);
```

// return 4 if found else 10.

List Immutable (ie all operations return a new List object)

val list = List(1, 2, 3, 4)

val list2 = List(7, 8, 9, 10)

* ~~list~~ list1 ++ list2 // concatenate

// declare empty list var list1 : List[Int] = List()

var list1 = List[Int]() ; // define list of integers.

* list1 ::= 5 // append 5

} returns new list object which
is assigned to list1

* list1 = 6 :: list1 // append 6

list1. drop(2) // drop 2 elements from front

list. init // retrieve everything except last

* list. head // first element

* list. last // last element

* list. slice(2, 5) // retrieves from 2 (inclusive) and 5 (exclusive)

list. splitAt(4) // split at index 4.

* list. take(5) // take the first 5 elements

* list. contains(6) // true if 6 is present

* list. startsWith(Array(1, 2))

* list. endsWith(Array(6, 7, 8))

* list. isEmpty // whether empty

list. indexOf(5) // index of element 5

list. indexOf(5, 2) // index of element 5 from index 2

* list. distinct // returns list of unique values

* list. reverse // reverse the list

* list. min

* list. max

* list. sum

- * `list.map { x => x.length }` // Transform
- * `list.filter { x => x % 2 == 0 }` // filter returns even number list
- * `list.count { x => x % 2 == 0 }` // count ^{num} of even numbers in list
- * `list.exists (x => x % 3 == 0)` // whether multiples of 3 are present
- * `list.flatMap (x => list.fill(x)(x))` // flattens the nested collection
- * `list.foreach (x => print(x + ":"))`

`list.sortWith(f)` → use f for sorting
eg. `list.sortWith(_ > _)` // Desc.

Apache SPARK

- Spark-shell (Scala command line)
- pyspark (python cli)

Spark Context

Sql context - Using sql to create, load & store tasks

Hive context - Use hive metastore to get Hive data.

- Directly run Hive statements

- Instead of running MR Hive queries will be translated to Spark code.

- Copy or create Softlink to /etc/hive/conf/hive-site.xml in /etc/spark/conf directory

- Spark-Sql
 - Allows hive queries directly to be executed from CLI
 - Similar to Hive CLI
 - Does not require Scala or Python knowledge

To connect to RDBMS using JDBC syntax of spark

➢ spark-shell --driver-class-path /usr/share/ -- /mysql-connector.jar
import org.apache.spark.SparkSQLContext

```
sqlContext.load("jdbc", Map(  
    "url" → url,  
    "dbtable" → "departments"))  
.collect().foreach(println);
```

➢ Spark can run in local mode with being aware of Master

Spark 1.2.1 does work properly with YARN

to start Spark 1.2.1 use

~~spark-shell~~ ➤ spark-shell --master local

Scala App

(1) Write scala file (spark code)

(2) Generate jar file using SBT (Scala build tool)

(3) use spark-submit command to submit Spark job to cluster

• μ is the mean of the sample
• σ is the standard deviation of the population

• n is the sample size

• σ/\sqrt{n} is the standard error of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

• σ/\sqrt{n} is the standard deviation of the sampling distribution of the mean

1. $\frac{1}{2} \times 10 = 5$

2. $10 \times 10 = 100$

3. $10 \times 10 = 100$

4. $10 \times 10 = 100$

5. $10 \times 10 = 100$

6. $10 \times 10 = 100$

7. $10 \times 10 = 100$

8. $10 \times 10 = 100$

9. $10 \times 10 = 100$

10. $10 \times 10 = 100$

11. $10 \times 10 = 100$

12. $10 \times 10 = 100$

13. $10 \times 10 = 100$

14. $10 \times 10 = 100$

15. $10 \times 10 = 100$

16. $10 \times 10 = 100$

17. $10 \times 10 = 100$

18. $10 \times 10 = 100$

19. $10 \times 10 = 100$

20. $10 \times 10 = 100$

21. $10 \times 10 = 100$

22. $10 \times 10 = 100$

23. $10 \times 10 = 100$

24. $10 \times 10 = 100$

25. $10 \times 10 = 100$

26. $10 \times 10 = 100$

27. $10 \times 10 = 100$

28. $10 \times 10 = 100$

29. $10 \times 10 = 100$

30. $10 \times 10 = 100$

31. $10 \times 10 = 100$

32. $10 \times 10 = 100$

33. $10 \times 10 = 100$

34. $10 \times 10 = 100$

35. $10 \times 10 = 100$

36. $10 \times 10 = 100$

37. $10 \times 10 = 100$

38. $10 \times 10 = 100$

39. $10 \times 10 = 100$

40. $10 \times 10 = 100$

41. $10 \times 10 = 100$

42. $10 \times 10 = 100$

43. $10 \times 10 = 100$

44. $10 \times 10 = 100$

45. $10 \times 10 = 100$

46. $10 \times 10 = 100$

47. $10 \times 10 = 100$

48. $10 \times 10 = 100$

49. $10 \times 10 = 100$

50. $10 \times 10 = 100$

- * Sudo yum -y install sbt
- * Env Variable: ~/.bash-profile

Important classes -

org.apache.spark.SparkContext
SQLContext
HiveContext

org.apache.spark.executor

Need to also know

Hadoop filesystem API
hadoop.io - package

Name spaces

- * org.apache.spark [SC, S]
- * org.apache.spark.rdd [RDDs]

org.apache.spark.io

↳ compression

org.apache.spark.sql

↳ SQLContext

org.apache.spark.sql.hive

↳ HiveContext

SPARK

RDD :

Driver Program : runs the main() function and executes various parallel operations on a cluster.

(Resilient Distributed Datasets)

RDD : main abstraction in Spark.

- it's a collection of elements partitioned across nodes of a cluster that can be operated on parallel.
- RDDs automatically recompute from node failures.

Shared Variables : → can be used by parallel operations.

- By default Spark's function runs in parallel across nodes. It ships a copy of each variable used in the function to each task.
- Sometimes we need to share variables across tasks. & hence the shared variables are used.
- Two types

- ① Broadcast variables : used to cache a value in memory of all nodes.
- created by calling broadcast function in spark context
`val bavar = sc.broadcast(Array(1, 2, 3))`
- After creation it should be used instead of the object passed in to the broadcast function.

② Accumulators : variables that are only added to such as counters or sums.

- created with initial value(v) by using function
`val accum = sc.accumulator(v);`
- Tasks running on the cluster can either add() method or += operator

Important Spark API

SparkConf : Configuration for spark application
: created using new `SparkConf()`

get() : get value for a configuration key

getAll() : get all config settings as `K, V`

getDouble(), getLong() - type getter methods for config values

set (String key, String value) : SparkConf
- set config parameter

setAppName (String) - name of our spark application

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

Spark Context

Spark Context :

- Main entry point for Spark functionality
- Represents the connection to a spark cluster
- Used to create RDD, Accumulators & Broadcast variables
- Only one spark context should be active in a single JVM.

* 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` - name of the application

* `broadcast(value:T)` - create broadcast variables

* `getConf`

* `new APIHadoopFile(path, class[Format], class[K], class[V])`

ex : `import org.apache.hadoop.mapreduce.lib.input._`
`val depseq = sc.new APIHadoopFile[IntWritable, Text, SequenceFileFormat[IntWritable, Text]]`
`(outputpath, classof[SequenceFileInputFormat[IntWritable, Text]],`
`classof[IntWritable], classof[Text]);`

* `SequenceFile(path, class[K], class[V])` : reads sequence files.

ex : `val depseq = sc.SequenceFile[Int, String]`

`("/user/donneral/Pointcut/itversity/spark/`
`2013-09-13-logs-with-vertices-and-simple-app Seq");`

* `textfile (path)` - reads textfile from HDFS / LFS

* `wholeTextFiles (path)` - Reads directory of text files from HDFS / LFS
returns `RDD [String, String]`

\downarrow file path \downarrow file content

* A pair is a type of tuple with 2 fields (k, v)

RDD Operations

Transformations - are operations that return another RDD.

- it does not evaluate the statement immediately.

- statements are executed only on Action operation.

cache() - Persists the RDD with default storage level (mem only)

coalesce() - returns a new RDD with ~~for~~ specified number of partitions

- used to reduce number of partitions after a filter operation when the dataset becomes sufficiently small

collect(f:) - returns new RDD by applying specified functions.

collect() - returns RDD containing all elements of this RDD.

context - returns Spark context that this RDD is created on.

count() - number of elements in this RDD.

Count By Value() (Ordering)

- returns count of each unique value in this RDD as a map.

- for single fields its unique value.

- for tuples its unique combination of fields

~~distinct(n)~~ : n = num Partitions

- returns RDD containing distinct elements

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

first(1): returns first element

flatMap : 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 ~~function~~ predicate return a collection instead of single element.

ex: WordCount - ~~the book~~ $\text{flatMap}(\text{line} \Rightarrow \text{line.split(" "))}$;

fold(zeroValue) ($Op: (T, T) \Rightarrow T$)

→ Aggregates the elements and their result of each partition

→ the zero value is passed as first argument to the predicate during the first call.

($(\text{Op})-\text{neutral} \& (\text{Op})-\text{associative}$) + \forall reduction from

foreach ($f: T \Rightarrow \text{Unit}$)

- Applies the predicate to all elements of RDD

- This is an action & not transformation, & does not return anything

→ used for side-effects like printing elements of RDD to screen/file

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

- groups a tuple by the specified predicate (the predicate should return a field)

- this function creates a nested structure RDD

- we can use the mapValues (for aggregations) flatmapValues (for sorting/ranking) to process further

reduceByKey

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 KV pair where the key is derived based on applying the predicate

ex: `depRDD.keyBy (x \Rightarrow x.depId)`

will return `RDD[(depId, depRec)]`

map - returns a new RDD by applying functor to all elements of this RDD.

`map[U] (f: T \Rightarrow U)`

- used for doing transformations on whole RDD (not by key)
- This is push model since elements are printed into function one by one

`mapPartitions (f: Iterator(T) \Rightarrow Iterator(U))`

`mapPartitionsWithIndex[U] (f: (Int, Iterator(T)) \Rightarrow Iterator(U))`

- Similar to map but applies the functor at partition level.

- So if we need to perform any initializations only once per partition

Ques - `mapPartitionsWithIndex()` case suppose the partition number, in addition to an Iterable (T) used to iterate over the elements

- These functions returns an Iterable (U) where 'U' is the transformation from V.

- In contrast to map() function, this is a pull model since our function needs to pull data from the Iterable

max() (Ordering) → finds maximum value in RDD
min() (Ordering)

persistent (hawtlevel) - Set RDD's storage level to persist values across operations.

reduce(f : (T, T) ⇒ T)

- used to perform overall/global aggregations (sum, max, min)
- can be applied to single elements, tuples etc.
- cannot be used for per key aggregations.

ex: find max price of product

product RDD. reduce ((x, y) ⇒ if (x.price > y.price) x else y);

repartition(numPartitions): return new RDD that has exactly numPartitions

SaveAsObjectFile(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 ($f : (T) \Rightarrow K$, ascending = true/false, numPartitions = ?)

- used for global sorting of RDD.
- can be used with single fields, pairs, tuples
- the predicate specifies the key to sort on.
- can be used to sort on multiple fields.
using (k_1, k_2)

ex: `productRDD.sortBy(x => (x._2, -x._5))`

Sorts the productRDD by 2nd field ASC, then by fifth field in DESC

subtract (other: RDD)

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

Ranking Global

`take(nun)` - take first n elements.

- this is an action that returns an Array.

takeOrdered(N)(Ordering)

- first sorts the elements either by implicit ordering or by the specified ordering. & then takes first n elements.
- by default Ascending order.
- can specify ordering function for DESC order.

`productRDD.take(3); ASC - Default order (first field in tuple.)`

`productRDD.take(3)(Ordering[Int].on(x => x._1))`

(equivalent to previous bottom 3 based on first field)

products RDD. takeOrdered (3) (Ordering[Int]). reverse on ($x \Rightarrow x - 3$)
// Top 3 based on third field

products RDD. takeOrdered (3).
(Ordering[Int, Float]). on ($x \Rightarrow (x_0 - 2, -x_0 - 5)$)

Top 3 based on ordering with multiple columns.

first column is Int which is 2nd field in tuple

Second column is float which is 5th field in tuple

top (N) (Ordering)

- works in exactly opposite way of takeOrdered by sorting elements in Desc order (based on first field)
- we can specify Ordering function to override this behaviour (desc, or other) field on multiple fields) by providing ordering function.

products RDD. top (3) (Ordering[Float]). on ($x \Rightarrow x_0 - 2$)

// top 3 records based on sorting by field 2

Union (other: RDD)

- Return the RDD which is Union of this with other

unpersist (blocking = true)

- make an RDD non-persistent & remove all blocks for it from memory & disk

PAIR RDD FUNCTIONS

- only available for tuples with two fields (k, v) pair

↙ reducer function

reduceByKey (func: $(v, v) \Rightarrow v$)

- passed in an accumulator & value (one after the other) into the predicate & the predicate should return a single tuple $(value)$ which is associative & cumulative.
- Similar to reduce() but the values will be first grouped by the key & then the values for each key will be passed into the predicate.
- Is used for doing aggregations (sum, max, min) per key basis.
- is a transformation ~~&~~ Uses a combiner which has same logic as Reducer functions

Ex: totalOrderAmtPerOrder

= orders . reduceByKey

$((acc, value) \Rightarrow (acc + value))$;

[Assuming Orders is $RDD[Int, Float]$]

\downarrow \downarrow
orderid order subtotal

CombineByKey [C] (combineFunc: $(acc, value) \Rightarrow Value$, func: $(acc, value) \Rightarrow Value$)

↑ combiner function

$\Rightarrow Value$

↓
reduce function

- Similar to reduceByKey but used when the combine logic & reduce logic is different
- The input to the combiner & output of combine & reduce should be same datatype

aggregateByKey ($\text{zerovalue}: V$) ($\text{combineop}: (V, V) \Rightarrow V$), $\text{reduceop}: (V, V) \Rightarrow V$

- more flexible than reduceByKey & combineByKey.
- Used when we need a different combine & reduce logic & further we need that the input values (elements of Values for each key) is of different datatype (V) and the output of combine & reduce function needs to be of different datatype.
- It also allows us to specify a zero value which gets passed in to combineOp as first parameter per key.
- * This zero value also forms the type that we are expecting after the combine & reduce operations.

ex: Calculating Avg revenue per day.

$\text{val revenueAvgPerDay} = \text{orderRevMap}.aggregateByKey((0, 0, 0))$

$((0, 0, 0))$, ← we are expecting tuple of (float, float, Int)

$((\text{acc}, \text{val})) \Rightarrow \{(\text{acc} - 1 + \text{value}, \text{acc} - 2 + 1)\}$

↑
Combine op. ↑ we are accumulating the total
in first field of tuple & count
in 2nd field.

$((\text{acc}, \text{val})) \Rightarrow \{(\text{acc} - 1 + \text{value} - i, \text{acc} + 2 + \text{val} - 2)\}$

↑
reduce op ↑ we are adding up all totals
input parameter & all counts
Match combine up
output is a tuple (Int, float)
so the acc & val both are tuples

foldByKey ($\text{zerovalue}: V$) ($\text{func}: (V, V) \Rightarrow V$)

- Similar to ~~foldByKey()~~ aggregateByKey that we can specify a zero value but this does not allow separate combine & reduce logic and also does not allow changing type of elements & output.

groupByKey (numPartitions = default)

- groups ~~as~~ a K, V flat structure & returns a key with Iterable (values) nested structure.
- Normally used for sorting & ranking within a key.
- Should not be used for calculating aggregates since this is inefficient (does not use combiner)

ex: `val products = products.groupByKey();`

Assuming products is K, V pair

where K = category id.

mapValues [U](func: (V) ⇒ U)

- used to perform a function over all values ~~per~~ key
- ~~partition~~ can perform transformations (aggregates) to return single value or sorting, ranking of the values to return a collection. (However this will result in a nested K, V structure.)
- For each key the values will be passed in ~~as Iterable~~ to the predicate (PULL model) ~~as~~ In case
- Normally used after grouping of data where V is Iterable

ex: `val grouped = orders.groupByKey()`

Assuming K, V where K is order id &

V is order subtotal.

`val revenuePerOrder`

= `grouped.mapValues (v ⇒ v.sum)`

→ V could be Iterable or single value

flatMapValues [U] (func: [V] ⇒ U)

- Similar to mapValues but expects U to be a collection instead of single value so that it flattens the result.
- used mostly for sorting & ranking ^{values} within a key.

ex: // Top 3 products per category by price

val prodGroup = products.groupBy[Key] () OR groupBy[---K]

Assume products is KV, K is prod cat & V is tuple
(product record)

val top3 = prodGroup.flatMapValues

(x ⇒ x.toList.sortBy[---4].take(3))

This will result in collection of
top 3 records sorted by
price DESC.

flatMapValue will convert this nested structure

K, Iterable(V) into

flat structure

K, V1

K, V2

K, V3

Count-By Key ()

- will return the number of values for each key in a local map.
- NOTE: This does not require grouping prior to this call i.e. the K, V could be a flat structure.

ex: val orderItems = ... (Assume orderItems is KV with

Order id as key & the
order item Record (tuple) as value

val numOrdersPerOrderID = orderItems.countByKey()

- Nested Tuples (Tuple inside a tuple) is not considered to be nested structure. It still is regarded as flat structure

Keys() - returns an RDD containing keys of each tuple

values() - returns an RDD containing values of each tuple

lookup(key:k) - returns a list of values for a given RDD key

- This is an action

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

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

- joins two RDDs & return a tuple containing the key & the corresponding records from this RDD & other as a tuple

- generates a flat structure.

~~ex:~~ - performs an inner join (filters out records from this & other that don't match)

ex: val joined = orders.join(orderItems)

Assuming orders is K,V where K is orderid (int)

V is order record string

orderItems is K,V where K is orderid (int)

V is orderItem record string

then joined will be

RDD: [(Int, (String, String))]]

↓ ↓ ↓
key order OrderItem
Record Record Record

cogrouped. flatMapValues ($V \Rightarrow \{$

for ($v1 \leftarrow V_{i-1}$, $v2 \leftarrow V_{i+2}$) yield ($v1, v2$);

)
Value from first iterable ↓
Value from 2nd iterable
combine them into a tuple

leftOuterJoin (other: RDD) - same as join but performs left outer join

- which means the 2nd field in the value tuple is

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

- i.e. the 1st or 2nd field in the value tuple of the result

will be empty for missing records in this or other respectively

- generates a flat structure.

cogroup (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 with K as stock symbol & V as stock record

Val Companies

- Assume K, V with K as stock symbol & V as company details

String
String
String
String
String
String
String
String

Val cogrouped = Stocks.cogroup(Companies)

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

Key
StockSymbol
↓
iterable for
iterating through
all stock records
for this key
↓
iterable for
iterating through
all companies
records for
this key

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

cogrouped.flatMapValues (

JAVA Pair RDD

SortByKey (asc : boolean)

Sort the RDD by key. So that each partition contains sorted range of elements

ex : Val products = Assume K, V K is price &

V is product record.

val sortedByPrice = products . sortByKey () . ASC

products . sortByKey (false) . DESC

Maps

val emptyMap = Map[("One", 1), ("Two", 2), ("Three", 3)]

val emptyMap = Map[Int, String]()

val numMap = Map(1 → "One", 2 → "Two", 3 → "Three")

* numMap.contains(1)

numMap.get(1) //Access

* numMap(2) //Access

foreach, count(), find(), map(), filter(), flatMap(), exists() ~~same as list~~
min, max, mkString, isEmpty same as list

* filterKeys (f(k) ⇒ Boolean)

* keys() list of keys

* values() list of values

* mapValues(f: (V) ⇒ C) ⇒ Map[K, C]

Tuples - Ordered collection of fields

val tupleMerge = (103, "Simpson", 10.25)

* tupleMerge._1, tupleMerge._2, tupleMerge._3 // Access

tupleMerge.productIterator // Iterator

tupleMerge.toString()

// Assign multiple variables using tuples

* var (a, b, c, d) = (1, 2, 3, 4)

SETS (Unique elements only)

val set1 = Set(1, 2, 1, 4, 5, 3, 2) // Immutable

val set2 = collection.mutable.HashSet(1, 1, 2, 1, 3, 5) // Mutable

* set1.exists(1)

* set1 -= 2 // remove element

* set1 --= Array(1, 2) // remove elements.

// Most methods come as list

Mutable collections

var arrayBuffer = ArrayBuffer(1, 2, 3)

var listBuffer = collection.mutable.ListBuffer(1, 2, 3)

var hashSet = collection.mutable.Set(0, 1, 1, 2, 2)

var hashMap = collection.mutable.Map("one" → 1, "two" → 2)

arrayBuffer = new ArrayBuffer[Int]

listBuffer = new collection.mutable.ListBuffer()

hashSet = new collection.mutable.~~Set~~ HashSet()

hashMap = new collection.mutable.HashMap()

For loop

Range (1 to 10)

Range (1 until 10)

R

Range (1, 10)

Range (1, 10, 2)

// gives list of to 10 inclusive

// gives list of to 10 exclusive

// 1 to 10

// steps of 2

1 to 10 // 1 to 10 inclusive

1 until 10 // 1 to 10 exclusive

1 to 10 by 2 // Steps of 2

10 to 1 by -2

Transformations using yield

for (i < 1 to 10) yield i*i // transform

for (i < 1 to 10 if i%2==0) yield // conditional transform

for (i < 1 to 5; j < 6 to 10) yield (i, j); multiple generations

// Paren matching

Var list = List ((1, 2), (3, 4), (5, 6))

for ((n1, n2) < list) println(n1 + n2)

for ((1, n2) < list) yield n2; select 2nd field where first field is 1

Function

def doSomething(a: Int) : Unit = {
 }
 ↗ a is automatically forced to val
 ↗ a is pass by reference
 ↘ no return type \Rightarrow

def getSum (num1: Int, num2: Int) : Int = {
 }
 ↗

getSum(5, 2) // Normal call

get(num2 = 2, num1 = 1) // Nested argument

def getSum2(args: Int*) : Int = {
 // Variable args
 for (item < args)
 {
 sum += item
 return sum;
 }

getSum2(1, 2, 3, 4)

var mList = (1, 2, 3, 4)

getSum2(mList) // This will not work since we need one

getSum2(mList: _*) // This works. ✅

Higher Order Func is a function which takes a func as parameter

def higherfunc(func: (Int) \Rightarrow Double, num: Int) = {
 }

func(num)

{

↑ This parameter is a function

which takes in Int & returns Double

def divisors = (num: Double) \Rightarrow num / 5;

higherfunc(num: Int \Rightarrow num / 5, 15)

Anonymous Functions

def add1(x:Int, y:Int) = x+y; //method (named)

val add1 = (x:Int, y:Int) => x+y; //anonymous methods

def doSomeOp(f: (Int, Int) => Int) = ? // This is higher
f(6, 2) order function
}

//Addition	doSomeOp((x, y) => x+y)	// Anonymous
//Subtraction	doSomeOp((x, y) => x-y)	// Function
//Multiplication	doSomeOp((x, y) => x*y)	
	doSomeOp(- * -)	// concise syntax

~~DateTime (joda)~~
~~import org~~

Date Time (Joda)

import org.joda.time. -

import org.joda.time.format. -

val currentDT = new DateTime()

val specificDT = new DateTime(2015, 5, 24, 18, 1, 34,

DateTimezone.getDefault())

val unixTimeMillis = DateTimeutils.currentTimeMillis()

val dt = new DateTime(unixTimeMillis, DateTimezone.getDefault())

val unixTimeMillis = specificDT.getMillis()

val formatter1 = DateTimeformat.forPattern("dd MMM yyyy hh:mm:ss")

val formatter2 = DateTimeformat.forPattern("yyyy-mm-dd hh:mm:ss")

val dt = formatter1.parseDateTime("26 OCT 2015 03:15:05.245 AM")

val str = formatter2.print(dt)

OR

val str = dt.toString("yyyy-MM-dd hh:mm:ss")

Datetime Java import java.util.*;

val currentDT = new Date()

val calendar
~~// Specific~~ = Calendar.getInstance(TimeZone.getDefault())
calendar.set(2015, 9, 17, 15, 7, 3)

val specificDT = calendar.getTime()

val unixTimeMillis = System.currentTimeMillis()

val dtFromUnix = new Date(unixTimeMillis)

val unixTimeMillis = specificDT.getTime()

//Formatting

import java.text.SimpleDateFormat

val fmt1 = new SimpleDateFormat("dd MMM yyyy hh:mm:ss a")

val dt = fmt1.parse("15 OCT 2015 03:15:25 PM")

val str

val fmt2 = new SimpleDateFormat("yyyy-MM-dd HH:mm:ss")

val str = fmt2.format(dt);

// Pattern Matching / switch case

n match {

// Simple match

case 0 => 1;

// Conditional match

case x if x > 0 => n * fact(n-1)

// anything else

case _ => -1;

Python

- interpreted as well can be compiled
- object oriented + scripting + functional
- interactive , batch

↓
⟩ python

⟩ directly write python statements
at the prompt

⟩ python myprog.py

Built functions

abs(x) - x (int, float)

* bool(x) - string to boolean

chr(i) - i is int.

* filter(func, iterable) - returns iterator for which the func returns True for elements in iterable

* float(x) - string to float

format (same as str())

`hex(x)` - x is int to hex formar string "0x"

input(prompt) -

`str = input("Enter something")`

→ This is some text

`print(str)`

* `int(x)` - convert string to int

* `len(s)` - s is string.

* `list(iterable)` - convert iterable to list

* `set(iterable)` - convert iterable to set (This is used to remove duplicates)

* `map(func, iterable)` - transform iterable to new sequence by applying the function.

* `max(iterable, key, default)`

↓
opt
↓ opt

`open(file, mode='r')`

mode → 'r', 'w', 'a', 'b', 't'

* `print(*objects, sep=' ', end='\n')`
`print(x, y, z...)`

* `range(start, stop)` - returns immutable sequence type

* `reversed(seq)` - returns a reversed iterator

* `slice(start, stop, step)`
 ↓ opt

* sorted (iterable , key , reverse)



↳ True or False

func to specify key to sort on

sorted(values , key = lambda x : x[6] , reverse = True)

* sum (iterable) - returns sum of all elements

* tuple (iterable) - convert to tuple

print()

print(x)

print(x, y, z)

print("%.s %.d %.0.5f" % (str, x, y))

print("The value is " , x)

* String is kind of sequence

Sequences - List, Tuple, Range, String

* Common Operations

Immutable Sequence Common operations (~~tuples~~, tuples, Keys of Dictionaries)

$x \in \text{seq}$ True if x present in seq.

$x \notin \text{seq}$

$\text{seq1} + \text{seq2}$ - Concatenate of seq1 & seq2

* For string concatenation the other types must be converted to string explicitly
eg: "The num is " + str(index)

$\text{seq}[i]$ - access i th element of the sequence

$\text{seq}[i:j]$ - take slice from i to j (exclusive)

$\text{seq}[i:j:k]$ - take slice from i to j in steps of k

$\text{len}(\text{seq})$ - length of sequence.

$\text{min}(\text{seq})$ - smallest element of seq.

$\text{max}(\text{seq})$ - largest element of seq

$\text{seq.count}(x)$ - total number of occurrences of x in seq.

Mutable Sequence Common Operations (List, Dictionary Values, Strings)

$\text{seq}[i] = x$ → assign x to i th element of seq.

$\text{seq}[i:j] = t$ → slice of seq from i to j

seq.append(x) → appends x to end of sequence

seq.clear() → removes all elements from seq.

seq.insert(i, x) → inserts x into i th position

seq.pop([i]) → retrieves & removes i th element

seq.remove(x) → removes first occurrence of x

seq.reverse() → reverses seq. in place

Text Sequence - str

- Strings can be enclosed in single, double or triple quotes
- String in triple quotes can span multiple lines.

`str(obj)` → returns string version of object

`str.startswith(str, start, end)`

`str.endswith(str, start, end)`

`str.find(substr, start, end)`

`str == "Some String"` → returns true

`str.format(*args, **kwargs)`

- perform string formatting operation

"The sum of {0} & {1} is {2}" `format(1, 2, 1+2)`

`str.lower()`

`str.upper()`

`str.split(sep, maxsplit)`

`str.strip()` - removes leading & trailing whitespace

Set Types

`set(iterable)` - converts iterable to set
→ removes duplicates

~~len(set)~~ len(s)

$x \in s$

`set.add(elem)`

`set.remove(elem)`

`set.pop()`

`set.clear()`

Map Type - dict

`dict1 = dict("one"=1, "two"=2, "three"=3)`

`len(dict)`

`dict[key]` → returns value corresponding to key else `KeyError`

`dict[key] = value` → sets value for key.

`key in dict` - whether k is present in dictionary

`iter(dict)` - returns an iterator over the keys

`dict.clear()`

`dict.copy()`

`dict.keys()`

`dict.values()`

Formatting

```
print ("%.3f %d %.5f" %(a, b, c))
```

```
str1 = "{} {} {}".format(a, b, c)
```

Conversion

```
int(x), float(x), str(x)
```

```
list(x), tuple(x), set(x)
```

Flow Control

```
if (cond) :
```

```
    statements
```

```
:
```

```
elif (condn) :
```

```
    statements
```

```
else :
```

```
    statements
```

```
for(x in list1) :
```

```
    print(x)
```

```
for(x in range(0,3)) :
```

```
    for(y in range(0,3)) :
```

```
        print(list1[x][y])
```

Comments # Single line

... Multiline

'''

'''

While (wthdn) :

statements :

functions (Pass by reference by default)

def funcName(arg1, arg2, ...):

 print --

 statements : . . .

variable arguments def funcName(*args):

 for(x in args)

 --
 :
 :

default args : def funcName(arg1, arg2 = 50):

 statements

:

Named args

funcName(arg2 = 30, arg1 = "str")

Anonymous

func1 = lambda x, y: x + y

* lambda can have only one statement

Hence if requires more complex logic define
a function & call within the lambda

OR

use function chaining.

func1(3, 2) gives 6

List Demo

grocery-list = ['Juice', 'Tomatoes']
empty-list = []

Access

grocery-list [0]

modify grocery-list [0] = "Apple"

Slice grocery-list [0:3] =
↳ exclusive

other-list = ['Potatoes']

Append (Combine/concatenation) : grocery-list + other-list

Append : grocery-list.append('Mango')

Extend (Append multiple) : grocery-list.extend(['Onions', 'Grapes'])

Insert : grocery-list.insert(3, 'Oranges')

Remove : grocery-list.remove('Apples');

Reverse : grocery-list.reverse()

Distinct : set(grocery-list) // convert to set

Count : len(grocery-list)

Min : min(grocery-list) // works with numbers

Max : max(numbers)

Sum : sum(numbers)

Print : print(grocery-list)

Membership : if ('Apples' in grocery-list)
if ('Bananas' not in grocery-list)

TupleDemo : enclosed in ()
pi-tuple = (3, 1, 4, 5, 9)

modification : list(pi-tuple) // convert to list & modify
// Tuples are immutable

All other operations same as list

sorted (iterable, func: key, reverse = false)

↳ lambda expr which provides sort key.

reversed (iterable, func: key)

↳ lambda expr which provides reverse key.

Dictionary : Keys immutable (tuples, string, numbers)

Values mutable

: keys are automatically sorted

: represented by {key: value, key: value}

Access : dict[1]

modify : dict[1] = "One"

keys : dict.keys()

values : dict.values()

Insert : dict[8] = "Eight"

Add multiple : dict1. ~~update~~ ^{update} extend ({9: "Nine", 10: "Ten"})

Clear : ~~dict[1]~~ dict.clear()

~~Count~~ : len(dict)

String demo

Format `"{} --- {} b4d {} --- {:.5f} ".format(a,b,c)`

Upper `str1.upper()`

Lower `str1.lower()`

Find `str1.find("Sometext")`

Length `len(str1)`

replace `str1.replace("oldtext", "newtext")`

Trim `str1.strip()`

Split `str1.split(",")`

Eqnals `str1 == "Sometext"`

Startwith `str1.startswith("Sometext")`

endwith `str1.endswith("Sometext")`

Substring `str1[0:4]`
 ↘ exclusive

Concat `str1 + str2`