

Pyspark

→ What is spark?

• It is general purpose in-memory computation engine

1. General purpose:

Data cleaning - Pig

Query - HIVE

ML operation - Mahout

...

When spark introduced;
each & every operation; we
can perform in spark itself

2. Computation: (Spark)

Hadoop vs Spark X wrong comparison

MapReduce vs Spark ✓ right comparison

→ Hadoop provide 3 components

- i. HDFS - storage
- ii. Resource manager
- iii. Map Reduce - computation

we can even use
Spark with HDFS

3. In Memory: (RAM)

↓ problem with map reduce

• It stores the intermediate results like aggregation, joins, filtering, functions; everything gets stored in DISK

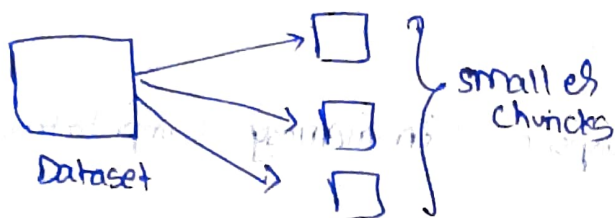
• So; for everytime, if the results are needed; need to go to DISK & fetch from DISK

• It's time taking; So Spark has introduced.

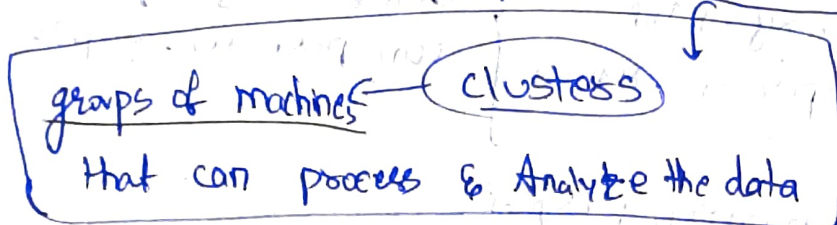
• It stores intermediate results in memory (RAM)

• So spark is faster than mapreduce

Spark Architecture

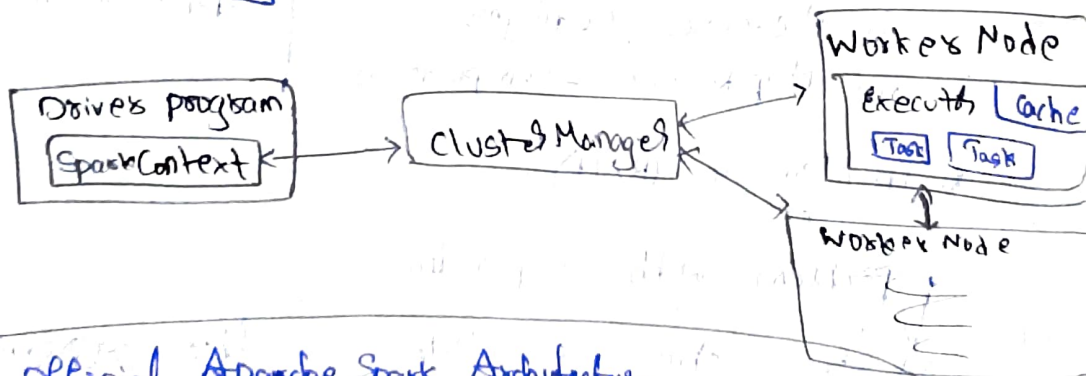


- It follows master-slave Architecture.
- It works on the concept of clusters.

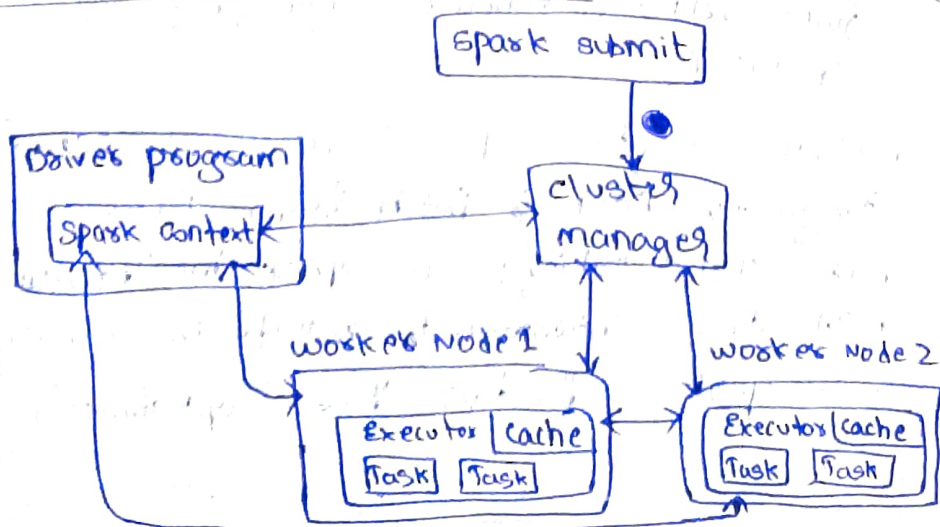


- Spark distributes the data & computations across multiple nodes in the cluster, allowing for parallel processing & faster data processing.

→ One of the node acts as master, remaining act as worker nodes



Official Apache Spark Architecture



- 1) When a user submits a Spark Application;
- 2) The cluster manager ~~will~~ takes the spark-context & assigns a node & shares it.
- 3) That node will act as Driver ~~or~~ Master Node.
- 4) Cluster manager (or) resource manager will basically be installed on top of all the nodes.
- 5) ~~Driver~~ Driver Node communicates with cluster manager (YARN, Mesos) whose work is to manage the resources.

- Driver program calls the main application & SparkContext is created which is the entry point of spark functionality.

Note • It would communicate with cluster manager whose job is to allocate executors inside worker nodes.

- Driver has components like DAG scheduler, Task scheduler etc.
- Cluster manager would send status & details of executors to Driver.
- Then Driver would schedule job.
- It would take RDD (Resilient Distributed Dataset) and translate into execution graph called DAG (Directed Acyclic Graph).
- These jobs would further split into different stages.

- In each stage, it would take Split the work for each executors called task & send to executors for processing.
- executor would perform that task execution given by driver & send the result back to driver.

RDD:

• Resilient Distributed Dataset

{ Dataset is nothing but the }
data that we provided.

- The distribution means the input data is stored across all worker nodes.
- Resilient means fault tolerance.

• RDD is distributed collection of dataset in Memory

So, file1 is in HDFS

rdd1 = load file from hdfs

rdd2 = rdd1.filter()

rdd3.collect()

(Lazy loading)
This will not be happen
untill

action got it.

→ There are 2 type of operations happens in RDD

1. transformation

2. action

DAG:

• (Directed Acyclic Graph)

Spark practical setup

1. use databricks (✓)
2. use jupyter notebook for spark



• Anaconda
• Java bdk
• Apache spark
• winutils

→ create dataframe:

Dataframe:

→ Dataframe is two-dimensional size-mutable, potentially heterogeneous tabular data structure with labeled axes (rows & columns)

df =

```
df = spark.read.csv('—')
```

```
df.show()
```

```
display(df)
```

⇒ Dataframe with headers

```
df1 = spark.read.csv('—',  
                      ; header = True)
```

```
df1 = spark.read.csv('—', headers = True)
```

→ How to print/get schema of a dataframe?

df.printSchema()

Dataframe with headers and proper schema:

```
df = spark.read.csv('—',  
                    headers = True,  
                    inferSchema = True)
```

if we did not specify it,
By Default values in the schema
(columns)
will be considered as 'string'
↓
df.printSchema()

(or)

→

```
df = spark.read.format("csv") \  
    .option('headers', True) \  
    .option('inferSchema', True) \  
    .load('—')
```

(or)

(or)

```
df = spark.read.format('csv') \
```

```
    .options(headers=True, inferSchema=True) \
```

```
    .load(—)
```

↳ How to read json file?

```
df = spark.read.json(—)
```

Then this will be okay.

if we have json files single-line

```
{ 'a' : 'A', { 'b' : 'B', { ... }, { ... } }
```

But,

how to read multi-line json file?

```
[ { ... }, { ... }, { ... }, { ... }, ... ]
```

```
df = spark.read.json(—, multiLine = True)
```

```
df = spark.read.format('json') \
```

```
    .option('multiline', True) \
```

```
    .load(—)
```

```
df = spark.read.option('multiline', True) \
```

```
    .json(—)
```


select:

select() function is used to select single, multiple ~~complex~~ column by index, all columns from the list & the nested columns from a DataFrame

→ `df.select("col1", "col2").display()`

(or)

`df.select(df.col1, df.col2, df.col3)`

select() using Column by index:

Emp	name	salary
0	1	2

 ⇒ `df.select(df.columns[index])`

`df.select(df.columns[2])`

`df.select(df.columns[1:4])`

withColumn:

PySpark withColumn() is a transformation function of DataFrame which is used to change the value, convert the datatype of an existing column, create a new columns & more.

• `withColumnRenamed()` to rename a DataFrame column, we often need to rename one column (or) multiple (or all) columns on pySpark DataFrame.

withColumn()

1. add new column based on existing column
2. add new column
3. change datatype
4. update the value of an existing column

→ `df.withColumn("total_cost", df.price * df.quantity)`

`df.withColumn("new_column", lit("val"))`

↓
PySpark function used to add a
{ constant (literal) value }

`df.withColumn("salary", col("Salary").cast("Integer"))`

`df.withColumnRenamed("old_col", "new_col_name")`

`df.withColumn("salary", col("Salary") - 5000)`

filter():

- `pySpark filter()` function is used to create a new Dataframe by filtering the elements from an existing Dataframe based on the given condition.

```
df.filter(df.col1 == 'val1')
```

```
df.filter(df.col1 != 'val1')
```

```
df.filter((df.col1 == 'val1') & (df.col2 == 'val2'))
```

```
df.filter((df.col1 == 'val1') | (df.col2 == 'val2'))
```

```
df.filter(col("col1").startswith("a"))
```

```
df.filter(col("col1").endswith("d"))
```

```
df.filter(col("col1").like("%oul%"))
```

Distinct() and dropDuplicates()

distinct() : transformation used to drop/remove the duplicate rows (all columns) from dataframe.

dropDuplicates() : used to drop rows based on selected (one or multiple) columns.

⇒ df = df.distinct()

df = df.dropDuplicates(['col1', 'col2'])

if there are duplicates in the attributes then those will be removed with their respective record.

dropDuplicates() : removes rows that are identical across ALL columns

dropDuplicates(['col1', 'col2']) : removes duplicates based only on selected columns.

dropDuplicates(subset=['col1']) : same as above

→ how to deduplicate using latest ~~timestamp~~ timestamp
{ window + row-number }

RDD:

• RDD (Resilient Distributed Dataset) is a spark's low-level Distributed data structure that stores data across multiple machines & processes it in parallel.

Why RDD exists?

• Before DataFrames existed, Spark needed:

- Distributed data
- Fault tolerance
- Parallel Computation

• RDD was the first abstraction of spark

Note:

DataFrames internally use RDDs but add schema + optimizations

Feature	RDD	DataFrame
API Level	low-level	high-level
Schema	X	✓
Optimization	X	✓
speed	slow	faster
SQL support	X	✓
usage today	Rare	90% jobs

```
rdd = spark.sparkContext.parallelize([1,2,3,4,5])
```

```
rdd = spark.sparkContext.parallelize([ — ])
```

Python collection

From a file

```
rdd = spark.sparkContext.textFile(" — ")
```

RDD operations

```
rdd2 = rdd.map(lambda x: x*2)
```

```
rdd3 = rdd.filter(lambda x: x>3)
```

Actions

```
rdd.collect()
```

```
rdd.count()
```

```
rdd.take(3)
```

when should we use RDD? (rare)

Use RDD only if:

→ we need low-level control

→ we have unstructured data

→ Complex logic not supported in DataFrame

```
rdd.mapPartitions( — )
```

DataFrame → RDD

```
rdd = df.rdd
```

RDD → DataFrame

initially need to add schema

①

```
rdd = spark.sparkContext.parallelize([  
    (1, 'A'), (2, 'B')])
```

```
df = rdd.toDF(['id', 'name'])
```

②

```
df = spark.createDataFrame(rdd, ['id', 'name'])
```

to-datecs (vs) date-format()

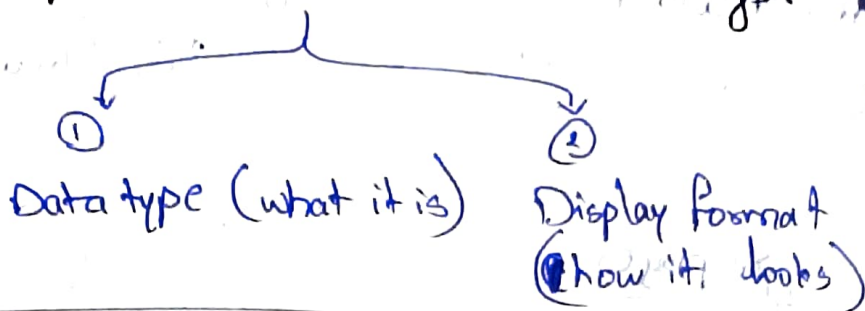
① to-datecs



to-datecs converts string → Date

• It doesn't control how the date is displayed

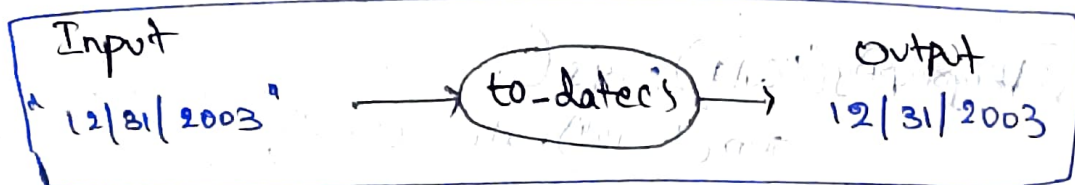
→ spark has two different things:



Note:

to-datecs

• Only changes the data type,
not the display format



↓
Converted correctly

But format will not be preserved

② date-formatcs

↳ This converts Date to String

date-formatcs ("___", "MM/dd/yyyy")

"dd/MM/yyyy"