



handles structured or semi-structured

Day/Date analysis / ML / streaming / graphs. Features.

Supports data

Function Programming

APACHE SPARK (fast & general engine for large-scale data)

programming → RDD.

• Written in Scala

• Developed by AMPLab UC Berkeley, now by Databricks.co

Spark Core API

R Python Scala SQL Java

• Spark Architecture

• Master-Slaves architecture →

master → 'Driver'

slaves → 'Workers'

• Transformation and actions are executed on worker node.

• Every spark application requires a spark context (main entry point to spark API)

• Spark shell provides preconfigured spark context called SC.

processed across cluster

• RDD (Resilient Distributed Dataset)

↓
data in memory
lost, it can be
recreated

• RDD fundamental units of data in spark (data contains)

• In-memory computation → computed results stored in distributed memory (RAM). Very fast.

• Lazy Evaluation → not performed immediately action then

• Fault tolerance → If failure occurs in any partition of RDD, partition can be re-computed from original fault tolerance input data to create it.

Maxim.....

- A full Spark API can be used with Spark SQL data by accessing underlying RDD.

Day/Date Calling persist or cache does not trigger execution / computation.

- Immutability → value can not be changed
- Partitioning →
 - collection of various data (RDD)
 - can not fit into a single node
- Persistence →
 - save result of RDD evaluation
 - stores intermediate result.

Transformation

Action

- define new RDD based on current one

Return values to driver/master node.

- map() → returns collection
- filter()

- count()
- take(n)
- collect()
- takeOrdered()

- flatMap() → returns flattened results

- saveAsTextFile()
- reduce()
- first()
- saveAsSequenceFile()
- foreach()
- saveAsObjectFile()
- countByKey()

input item can be mapped to 0 or more output (returns seq rather than single)

- RDD organized as directed Acyclic graph (DAG)

DAG track dependencies (lineage)

- nodes are RDDs

- arrows are transformation

- **Pipelining** (spark will perform sequence of transformation by row so no data is lost)

action re-evaluates the lineage transformation starting with base

- **RDD Lineage** (spark maintains each RDD lineage previous RDD on which it depends)

- **Spark performance Tuning** (computations over RDD, transformation embedded in chain)

- dataset loaded from databases computation are performed results are returned.

Q: Call actions on RDD?

- recomputation of transformations each time, increase resource usage. To optimize, see that action is performed again and again.

Maxim.

JVM → Java virtual machine.

Day / Date

cache method is implemented itself as a call to `rdc.persist` (Storage level. • Memory = only)

RDD: cache()

- cache RDD into memory
- By calling cache method, first action say RDD keep values it has calculated in memory. RDD then uses cached values for calculating 2nd action.

Unpersisting

- As more and more RDDs are cached, memory decrease.
 - Spark starts expelling partitions from cache.
 - Zeida JVM garbage collection time → Unavoidable.
- Call `un-persist` method on RDD when caching is no longer needed.

Why not to use Caching (Recomputation is faster as increase memory → more money.)

- if once to read dataset then no point of caching
- depends on
- How many times data is accessed
 - Amount of work involved.

PySpark (interface for spark in python)

↓
pyspark shell.

Spark Dataframe

- Supports parallelization
- Multiple nodes
- Lazy execution
- Immutable
- Distributed & Spark dataframe is faster for large amount of data.

Pandas Dataframe.

- no parallelization
- single node
- eager execution
- Mutable
- not distributed & slow for large data.

Maxim.....

- Simple Algo for frequent elements in streams & bags

Day/Date • Spark SQL is not a replacement for a database
ETL/structured to other application.

• Narrow Transformation Wide Transformation

• each input partition → one output partition

each input → many outputs

• Faster

Slower

• Not require any data shuffling over cluster network

might require data shuffling over cluster network.

• Map, Flatten map, Map partition, Filter, sample, Union

• Intersection, Join, Cartesian, repartition

• returns dataset of (K, V) pairs where value for each key are aggregated.
• map side combine
• same to combine in map reduce().

reduce by key groupby key (returns dataset of $(K, \text{iterable})$ pairs).
• do not map side combine
• cause diverse effect to output.

• Spark DataFrames. (main abstraction in Spark SQL)

- Comparable to RDDs in Core Spark.
- Distributed collection of data organized into named columns.
- are created →
 - existing structured data source
 - existing RDD
 - programmatically defining a schema.

Example:

```
from pyspark.sql import SQLContext entry point
# initialize SparkSession
spark = SparkSession.builder
    .appName("Tashi")
    .getOrCreate()
or
sqlc = SQLContext(sc)
```

Maxim..

sqlCtx.sql("select ...")

Day / Date

deals with dataframe metadata

• DataFrame Basic Operations.

- - schema - schema object describing data.
- - printSchema - displays schema as visual tree
- - cache/persist
- - dtypes - array of (col name, type) pairs.
- - explain - prints debug information about dataframe.
- - columns (names of columns)

• DataFrame Queries. (returns new dataframe similar to RDD transformation)

- - distinct - return new dataframe with different elements.
- - join - 1 dataframe join with another
- - limit -
- - select -
- - filter -

DF.select("age")

DF.where("age > 21")

- - Some queries take one or more cols

1.

agedDF = DF.select(DF.age)

DF.select(DF.name, DF.age + 1)

DF.sort(DF.age.desc())

- - Frequent Pattern mining \Rightarrow df.freqItems.

- - Data can be stored to data source.

- Build in support for JDBC & Parquet file
- Create JDBC Table
- Insert into
- Save as Parquet File
- Save As table

- - DataFrames are built on RDDs.

- Base RDD contains row object.
- use rdd to get underlying rdd.

- - Row RDD have all standard actions and transformations.

Maxim.....