Scala → object oriented
+
Spark SQL
↓
supports data
analysis /ML/ streaming/graphs.
data.
function Programming
features.
handles structured or
semi-structured

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

programming → RDD.
operations

- Written in Scala

- Developed by AMPLab UC Berkely, now by Databricks.com

| Spark Core API |
|---|
| R    Python    Scala    SQL    Java. |

- Spark Arkitecture

  - Master-Slave arkitecture → master → 'Driver'
    slaves → 'Workers'

  - Transformation and actions are executed on
    worker node.

- Every spark application requires a spark
  Context ( main entry point to spark A2 )

- 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 containers

  - 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.

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

- — Persistance →
  - 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 | • Count () · take sample () |
| filter () | returns ← take (n) · take ordered () |
| first n rows | collect () |
| input item can be mapped to 0 or more output (returns Seq rather than single) ← flatmap ( ) → returns flattened result. | collect all the rows. Save as Text file () reduce () Save as object file () first () Count by key () |
| | Save as Sequence file () foreach () |

- RDD organized as Directed Acyclic graph (DAG)

  DAG track dependencies (Lineage)

  - — nodes are RDDs
  - — arrows are transformation

2.

- Pipelining (Spark will perform sequences of transformation by row so no data is lost).

  ↓

  action re-executes 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.

JVM → Java virtual machine.

Cache method is implemented itself as a call to rdd.persist (Storagelevel. .Memory-only)

## RDD: cache()

- — Cache RDD into memory
- — By calling cache method, first action say RDD keep values it has calculated in memory, RDD than uses cached values for calculating $2^{nd}$ action.

## Unpersisting

- — As more and more RDD's are cached, memory decrease.
  - Spark starts expelling partition from cache.
  - Zaida JVM garbadge collection time → Unsvoidable.
- — Call un-persist method on RDD when caching is no longer needed.

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

| ⬇ | ⬇ |
|---|---|
| if once to reed 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 | Pandas Dataframe.. |
|---|---|
| • Supports parallelization | no parallelization |
| • Multiple nodes | single node |
| • Lazy execution | Eager execution |
| • Immutable | Mutable |
| • Distributed & Spark dataframe is faster for large amount of data. | not distributed & slow for large data. |

Maxim

- Simple Algo for frequent elements in streams & bags

- Spark SQL is not a replacement for a database ETL/structured to other application.

- **Narrow Transformation**          **Wide Transformation**

| 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 required data shuffling over cluster network. |
| - Map, flatten map, Map partition, filter, sample, Union | - Intersection, Join, Cartesian, repartition |

- returns dataset of $(k, V)$ pairs where values for each key are aggregated.
- map side Combine
- Same to Combiner in map reduce ().

reduce by key
group by key (returns dataset of $(k, iter-bleVrs)$ pairs).
- do not map side combine)
cause diverse effect to output.

freq 2×7

- **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 $\boxed{SQL Context}$ ← entry point

# intilize Spark session
spark = SparkSession . builder \
.appname ("Tashi") \
. getOr Created ( )

or
SqlC= SQL Context (sc)

sqlCtx.sql("Select ....)

deal1 with dataframe metadata

- **Dataframe Basic Operations.**
  - Schema — Schema object describing data.
  - printSchema — displays schema as visual tree
  - Cache/persist
  - Columns (names of columns)
  - dtypes — array of (col name, type) pairs.
  - explain — prints debug information about dataframe.

- **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

    ageDF = DF. Select (DF.age)

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

    DF. Sort (DF.age . desc())

  - frequent Pattern mining => df. freqItems.
  - Data can be stored to data source.
    - Built in support for JDBC & Parquet file
    - Create JDBCTable
    - Insert into
    - Saveas Parquet File
    - Save As table
  - Dataframes are built on RDDs.
    - Base RDD contains row object.
    - Use rdd to get underlying rdd.

  - Row RDDs have all standard action and transformations

*Maxim*