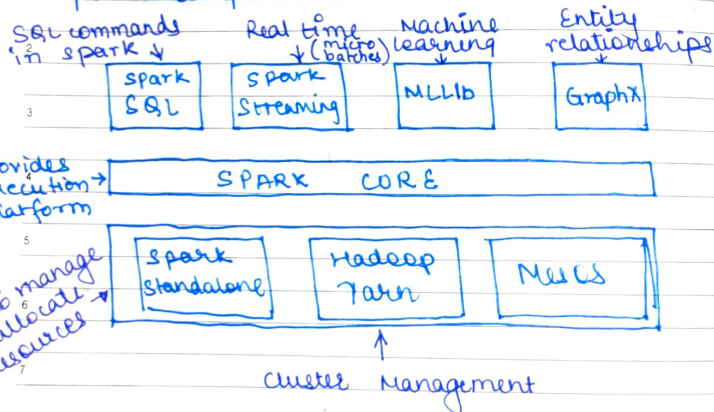


★ SPARK

- It is an independent framework.
- It is a low latency, lightning fast cluster computing platform.
- Uses in-memory storage to process the data.
- Can perform batch as well as stream processing.
- Reads data from HDFS | HIVE | DBS | cloud
- processes & stores in HDFS | DBS

1. Spark Components



1. Spark Core

- Foundation of parallel & distributed processing of huge dataset.
- Believes speed by providing in-memory computation capability.

→ Features :

- In charge of essential I/O functionalities.
- Task dispatching
- Fault recovery
- Embedded with RDDs (Resilient distributed dataset).
- handles partitioning data across all nodes in cluster.
- It holds them in memory pool of the cluster as a single unit.

2. Spark SQL

- works to access structured & semi-structured info.
- Enables powerful, interactive, analytical app. across both streaming & historical data.

→ Features :

- Cost based optimizer
- Mid query fault tolerant
- Full compatibility with existing hive data

3. Spark streaming

- Add on to core spark API which allows scalable, high-throughput, fault-tolerant stream processing of live data streams.
- Access data from sources like kafka, Flume, kinesis or TCP socket

→ It uses Micro-batching for real-time streaming

4. Spark MLlib

→ It is a scalable ML library that performs high-quality algorithm & high speed.

5. GraphX

→ API for graphs & graph parallel execution.
→ It is network graph analytics engine & data store.

• 12 Features of spark

1. Swift processing

→ high data processing speed, made possible by reducing no. of read-write to disk.

2. Dynamic in nature

→ easily develop a parallel application

3. In-Memory computation in spark

→ Processing speed is increased.
→ Data is cached so we need not fetch data from disk every time, thus time is saved.

4. Fault tolerant

- Provided through RDDs.

5. Real-Time Stream processing

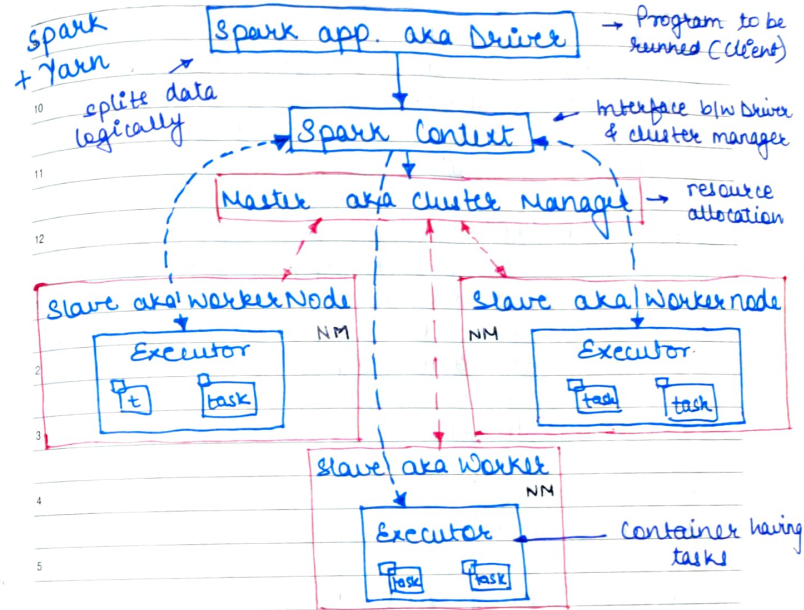
6. Lazy Evaluation

- All the transformation we make in RDD are lazy in nature.

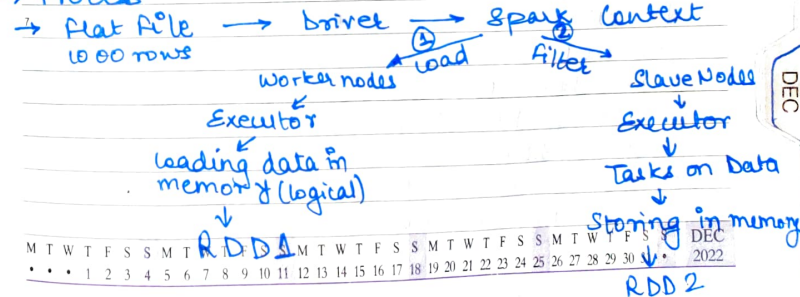
It does not give result until action command are called.

NOV MTWTFSSMTWTFSSMTWTFSSMTWTFSSMTWTFSS
2022

• Architecture



→ Process



MTWTFSSMTWTFSSMTWTFSSMTWTFSSMTWTFSS
• • • 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 DEC 2022

RDD 2

• Limitations

- No inbuilt optimization engine.
- Handling structured data
 - need to specify the schema.
- Performance limitation
 - As it involves the overhead of garbage collection & Java serialization which are expensive when data grows.
- Storage limitation
 - spill over, if size of RDD is larger than RAM then it uses disk for remaining file RDD size.

• Ways to Create RDDs in Spark.

1. Parallelized collection (parallelizing)
 - by taking an existing collection in the program
 - passing it to SparkContext's parallelize() method.
 - Used in initial stage, as it creates RDD quickly
2. External Dataset (Referencing a dataset).
 - We can use textFile method.
 - It takes URL of file & reads it as a collection of line.
 - Loading CSV, JSON, textFile.
3. Creating RDD from existing RDD.
 - Transformation is the way to create an RDD from already existing RDD.

NOV 2022 MTWTFSSMTWTFSSMTWTFSSMTWTFSSMTWTFSS
 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31

• RDD Persistence & Caching

- Persistence is an optimization technique in which saves the result of RDD evaluation.
- Using this we save the intermediate result so that we can use it if required.
- It reduces the computation overhead.
- When we use the cache() method we can store all the RDD in-memory.
- We can persist the RDD in memory & use it efficiently across parallel operations.
- The diff. b/w cache() & persist() is that using cache() the default storage level is MEMORY_ONLY while using persist() we can use various storage levels.
- When we persist RDD each node stores any partition of it that it computes in memory & makes it retrievable for future use. This process speeds up the computation.
- When the RDD is computed for first time, it is kept in the memory on the node. The cache memory of the spark is fault tolerant so whenever any partition of RDD is lost, it can be recovered by transformation operation that originally created it.

MTWTFSSMTWTFSSMTWTFSSMTWTFSSMTWTFSSMTWTFSSMTWTFSSMTWTFSSMTWTFSS
 • • • 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 • DEC 2022

• Storage levels of persist()

1. Memory-only (spill over in disk)
2. Memory-and-disk (serialized Java object)
3. Memory-only-set
4. Memory-and-disk-set
5. Disk-only.

• Unpersist RDD

- Spark automatically drops out the old data partition in LRU (least recently used) fashion.
- We can remove manually using RDD.unpersist().

• RDD Operations

1. Transformation

- Produces new RDD from existing RDD.
- Applying transformation built an RDD lineage with the entire parent RDDs of the final RDD dependency graph.
- Are lazy in nature i.e. they get executed when we call an action.
- Two types of transformation:

1. Narrow transformation

- All elements that are req. to compute the records in single partition live in the single partition of parent RDD.
- map(), filter(), flatMap(), mapPartitions(), sample(), union().

2. Wide transformation

- All elements that req. to compute the records in the single partition may live in many partitions of parent RDD.
- eg. groupByKey(), reduceByKey(), Join(), intersection(), distinct(), cartesian(), repartition(), coalesce().

2. Action

- Works with actual data when action is performed.
- Does not form new RDD.
- The values of action are stored to ~~spark~~ drivers or to the external storage system.
- Action is one of the way of sending data from Executor to the driver.

4. count() - no. of elements is returned

collect() - returns the entire RDD's content to driver program.

take(n) - returns n no. of elements from RDD.

top() - extract top elements from RDD.

countByValue() - returns, occurrence of each element

[reduce() - operation like addition, takes two elements as input.

fold() - Take zero value as input.

diff. - reduce throws an exception for empty collection, but fold is defined for empty collection.

aggregate(), foreach()

★ RDD Lineage

- When we create new RDD from an existing RDD, that new RDD also carries a pointer to the parent RDD.
- All the dependencies b/w RDDs those are logged in a graph, rather than actual data, called lineage graph.
- It is a graph of all the parent RDDs of an RDD.

★ Dataframe

- Data organised into named columns.
- Similar to table in RDBMS.
- Can say that it is a relational table with good optimization technique.
- Processes large amount of structured data.
- Contains schema (illustration of the structure of data).
- Immutability, in-memory, resilient, distributed computing capability.
- Sources data from structured data file, tables in hive, external dbs or existing RDDs.
- Available in scala, Java, Python & R.

→ Data Frame over RDD

- ① Provides memory management
 - data is stored in off heap memory in binary format. Serialization is avoided

② Optimized Execution plan

- query optimizer, where execution plan is created for the execution of a query.

• Limitation of RDD

- Does not have any built-in optimization engine.
- No provision to handle structured data.

• Features of Dataframe

- distributed collection of data organized in named column, equivalent to table in RDBMS.
- Deals with structured & semi-structured data formats. eg. Avro, CSV, elastic search, Cassandra.
- Deals with storage systems - HDFS, HIVE, MySQL etc.
- Catalyst supports optimization, 4 phases:
 1. Analyze logical plan to solve references
 2. Logical plan optimization
 3. Physical planning
 4. Code generation to compile part of query to java bytecode.

• Limitation

- Does not provide provision for compile time type safety. So, we need to make a structure in order to manipulate.

★ Dataset

- 8 - Strongly typed & is map to relational schema.
- Represents structured queries with encoders.
- 9 - Provides both type safety & object-oriented programming interface.

• Features

- Optimized Query - (Catalyst Query Optimizer).
- 12 - Analysis at compile time. - Check syntax & analysis at compile time.
- 1 - Persistent storage - serializable & queryable.
- Faster computation (than RDD).
- 2 - Less memory consumption - structure of data in dataset is known.