

Interview

- ① Tell me about Yourself
- ② Difference between conceptual, logical and physical data model.
- ③ what is OLTP System
- ④ Normalization Example
- ⑤ What is OLAP System
- ⑥ What is Data Vault
- ⑦ What is SCD
- ⑧ What is DQS
- ⑨ What is Data lake
- ⑩ Different type of databases
- ⑪ Python
- ⑫ Shell
- ⑬ Hive
- ⑭ PySpark
- ⑮ Data Architecture

PYSPARK

- ① Introduction to Spark
- ② Why Spark
- ③ memory management in Spark
- ④ Context in Spark
- ⑤ Initiating Spark Job and Submit Job
- ⑥ Parameters to be submitted to Spark Job
- ⑦ Spark Architecture
- ⑧ Spark Features
- ⑨ Spark Support Core, SQL, ML, Stream
- ⑩ RDD Dataframe, SQL Dataset
- ⑪ Types of RDD Action / Transformation
 - Select
 - Count
 - Filter
 - Distinct
 - groupby
 - Ordering
 - windowing
 - Top record
 - writing
 - Joins
 - Sampling
 - Set Operation
 - repartitioning
 - Coalescing
- ⑫ Ways to create RDD
 - parallelize
 - from other RDD
 - Python list
- ⑬ Partition And repartition.
- ⑭ Extraction to HDFS / Save

(15)

Performance

RDD Persistence

Caching

Spark SQL Architecture

Broadcast Accumulator variable

Catalyst Optimizer

Tungsten Optimizer Execution engine

prior 2.0 Volcano iterator model

After 2.0 whole stage code generation

Tungsten Engine

Debugging

to Debug String

explain.

①

- Introduction to Spark | Why spark
- Distributed Parallel Processing
 - Horizontally Scalable
 - Unified Technology can be used for SQL, ML, Graph, Stream
 - Support multiple languages
- In memory processing
- Map reduce only supported map reduce so it was miss
 - Other support filter / count
 - So it was Stateless, Spark support the same.
 - Spark separate storage and processing and hence can run on different cluster manager

	Map Reduce	Spark
1	Stateless	support filter
2	Support Java	support all very
3	lot of tools	in memory processing
4	Different tool make maintenance difficult as they have their own cluster management	Unified cluster management
5	Support Batch Processing	Support Batch plus Streaming

(3)

memory management Spark
Spark support 5 cluster manager
Kubernetes

MESOS

ECS

Standalone

Yarn

Spark Core is main component of Spark
and on top of it it has RDD API
to support

Spark SQL

Spark ML

Spark GraphX

Spark Streaming

(4)

Spark Session / Spark Context

Spark 1

Spark one has
Independent Context

Spark Context

SQL Context

hive Context

Spark 2

Spark has Spark
Session Objectwhich has all
three Context

SparkContext

SQLContext

hive Context

When we start PySpark

it gives sc as

SparkContext

when we start
PySpark2 it gives

spark as spark

session object.

to get spark context
you do sparkSession.
sparkContext.

⑤ | ⑥

Spark-submit job is used to submit Spark job to Spark cluster for execution

Below are important command line parameter of spark-submit job

- master
- deploy-mode (client | cluster)
- conf (configuration file)
- Properties-file (to submit our own cont file)

- driver-memory
- executor-memory
- num-executors

- jars (connector jars odbc)
- packages (Pyro package)
- py-files (passing Python file to copy on cluster)

⑧ 19

Resilient Distributed Dataset (RDD)

RDD has below features

- Resilient
- Distributed
- Immutable
- Lazy evaluation
- Lineage
- Fault tolerant
- Compile time safety

When to Use RDD?

Performance is not an issue

Need control of processing

For semistructured data

When schema of data is not known

RDD can't be optimized by Spark because
RDD data + operation is opaque to
Spark and hence SQL and DataFrame
are designed on top of RDD.

⑦

Spark Architecture

Scenario

200MB file in hdfs is to be loaded
map the data
filter the data
group the data
store the data

- Client machine JVM Submit the job via spark-submit command depending on deploy-mode and cluster manager assuming cluster and yarn the application will be processed
- Yarn Resource manager consist of Application manager and Resource scheduler.
- Application manager launch the container (Data node) seeking space equal to Driver memory and launch Driver Program inside Application master which is launched inside container.
- The Driver program consist of SparkContext, JVM and SparkApplication

- If spark code has Action transformation
The code is passed to DAG scheduler.

The DAG scheduler first create logical execution plan, in this case

Stage 1 → load file → RDD1

Stage 2 → map → RDD2

Stage 3 → filter → RDD3

Stage 4 → group → RDD4

Stage 5 → store data → RDD5

and also does optimization of logical plan and combine narrow and wide transformation and create physical execution plan.

Stage 1 → RDD1

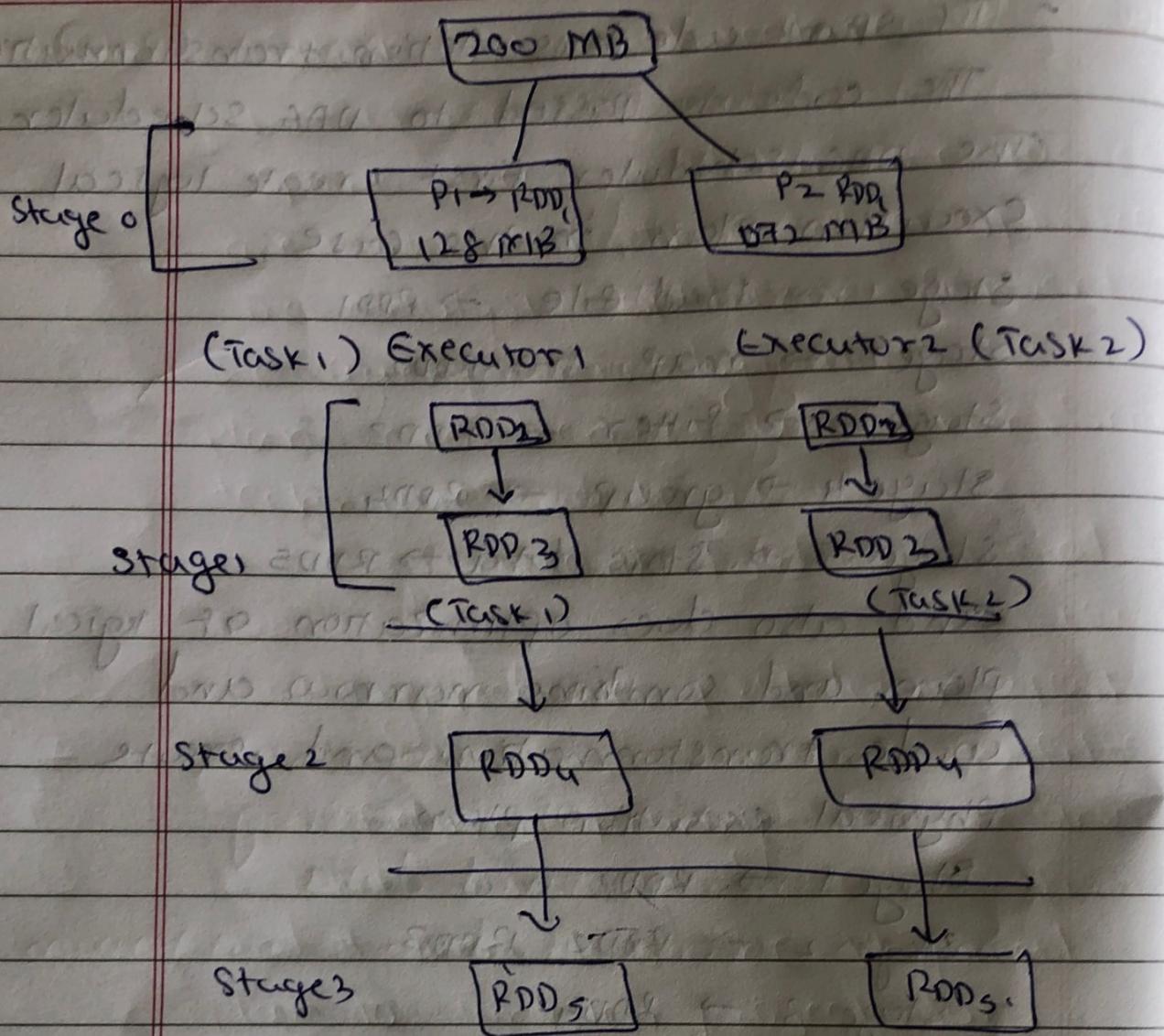
Stage 2 → RDD2, RDD3

Stage 3 → RDD4

Stage 4 → RDD5

- Driver program need two partition to be created as size is of 200MB

so it ask for 2 executor to application master which negotiate with application manager and two executor program are created on worker node which status are passed by Node manager and also the code is passed to the executor along with data



- Once the code is executed and result is output the individual node manager or cluster or the executor passes output to driver program i.e Application master and it passed it to client. The Application master deregister itself from Resource manager.

- (11) RDD :- How to Create RDD
GRAC \Rightarrow Aggregate key
shuffle Combiner concept
- (12) Partition Coalesce
use coalesce to decrease partition
use repartition to increase partition
repartition try to create equal size
partition and hence lot of shuffle happen.
- (13) Hadoop only had text, csv and
Sequence file
- (14) Spark need to define storage level
and then persist once storage level
is set and then it can't be changed.
unpersist method is used to unpersist.
- (15) Broadcast variable are variable, constant
which can be shared across
executor once for node. It need
not be changed once broadcasted.
- (16) Accumulator variable are used for
commutative or additive or associative
operation. This is used for counter.
Accumulator variable value are send to
Driver program.
Broadcast and Accumulator variable
work only on action.

(16)

Spark SQL Architecture

Spark SQL or DataFrame queries are submitted to Catalyst optimizer. The catalyst optimizer converts the query into logical execution plan and then it uses transform techniques such as constant folding (literal identification), predicate pushdown → push filter.

Column pruning → remove unused column from scan and then pass the physical plan to Tungsten project for execution.

Tungsten execution engine executes query in parallel.

Prior 2.0 Tungsten execution engine was not used because Volcano iterator model used to execute in parallel as it was working independently and also no dependency on step 1 and 2, it has disadvantage of extensive memory use hence whole stage code generation Tungsten engine was used which uses

cost based optimizer, and it
optimize all stages of execution and
consider dependency.

- ①
 - ②
 - ③
 - ④
 - ⑤
 - ⑥
 - ⑦
 - ⑧
 - ⑨
- to obtain correct result with following
order
- ①
 - ②
 - ③
 - ④
 - ⑤
 - ⑥
 - ⑦
 - ⑧
 - ⑨
- and last executed query two stages
- ①
 - ②
 - ③
 - ④
 - ⑤
 - ⑥
 - ⑦
 - ⑧
 - ⑨
- return value to dependent query
- ①
 - ②
 - ③
 - ④
 - ⑤
 - ⑥
 - ⑦
 - ⑧
 - ⑨
- and then return final result
- ①
 - ②
 - ③
 - ④
 - ⑤
 - ⑥
 - ⑦
 - ⑧
 - ⑨

DataFrame

Dataset + Schema = DataFrame

①

Introduction

②

How to Create DataFrame

③

Selecting data in DataFrame

Sqrl

table

④

Reading data from sources into df

⑤

UDF

⑥

Catalog

⑦

Datatypes and Row, Column function

⑧

Practice

⑨

Built in Functions

- ① DataFrame = Dataset + schema
- ② Command to create DataFrame
- Version
 - range

method CreateDataFrame

- Python list

- Dict

- RDD

- Row object + Parallelize

- Pandas DataFrame + 2D list

- ③ Spark SQL from TempView and Global

TempView

Global tempview is accessible to all sessions

RDD + DataFrame + tempview + SparkSQL.

- ④ Table Function

Reading data into DataFrame
read.load is used to load different file formats

CSV

ORC

AVRO

TEXT

Parquet

Hive : since Hive is integrated so SparkSQL JDBC is not there we need to pass

JDBC data read into Spark we can partition using parameter

lower bound upperbound and partition
(upperbound - lowerbound) / partition = stride
part = lowerbound + stride

⑤ UDF in DF.

UDF are user defined function
defined in code

ⓐ Annotation is required to use
pyfunc function in dataframe
however to use pyfunc function
in sparkSQL it need to be register

UDF are not shared across session
newSession is used to initiate
new session in spark

⑥ Catalog gives details about metadata of dataframe

list databases tables columns

is cached cached uncacheTable

clearCache

recreateTempTable refreshable

refreshByParam

dropGlobalTempView

dropTempView

listFunctions

registerFunctions with register Partition

(7)

Data types in Spark

Integer - Int Float Double

String - char varchar string

Boolean

Date - Date Timestamp

Complex - Array, map, struct

All are Types in DataFrame and has alias in Spark SQL

Row - Row class : count, index, toDict()

column - select

orderBy

where

nullCheck - null function

Column function apply on Right side of the column



⑧

Practice till now

Select

where / filter

orderby

SortWith Purification

Union / Unionall

Intersect / Intersectall

UnionByName

ExceptAll (similar to minus)

inner join / join

full Outer join

Left join

Right join

Cross join

left anti join

left semi

self join

multi condition join

multi dataframe join

Aggregation function

summary

avg

min

max

sum

sumDistinct

count

countdistinct

first

last

collect_set

collect_list

Group by API

avg mean count min max sum agg

pivot apply

filter similar to having

no unpivot function so we need to
use stack and expr function in spark

stack, pivot, unpivot

use selectExpr in DataFrame

window function

Ranking

Row_number

(8)

Dense_rank

Dense_rank

Percent_rank

rank

ntile

Cume_dist

Analytical

log
lead

Aggregate

avg, sum, min, max, count, first, last

Other

range between

nowbetween

unbounded following

unbounded preceding

currentrow

first, last, greatest, least, collect-list

Sampling

Sample

sampleBy

⑧

Built in Functions

monotonically_increasing_id()

lit()

- literal static column

expr()

- dynamic column

specific_partition_id()

rand() / randn()

String Functions

split

length

lower, upper, initcap

ltrim, rtrim, trimboth

lpad, rpad

reverse

repeat

hex

concat

concat_ws

substring

substring_index

instr

• locate

translate

• overlay

Regexp Functions

regexp_extract

regexp_replace

rlike()

Date Functions

current_date

current_timestamp

next_day

last_day - of month

day_ofweek

day_ofmonth

day_ofyear

week_of_year

second

minute

hour

month

quarter

year

months_between

date_add()

date_sub()

add_months()

datediff()

date_trunc()

date_format()

unix_timestamp

to_timestamp

from_unixtime

from_utc_timestamp

to_date

Null Functions

(ISnull)

(ISnum)

numvl

Coalesce

+-----+
|-----|

Collection Functions

Array map struct

size

element_at

struct

array

array_max

array_min

array_distinct

array_repeat

flatten

slice

array_position

array_remove

array_sort

Sort_array - we can pass where to be my

array_contains

array_union

array_except

array_intersect

array_join

arrayzip

arrayoverlap

shuffle

union

create_map

map_from_entries

map_from_array

map_keys

map_values

map_concat

Sequence

NA Functions

drop : drops null values

df_na_drop()

fill : fill null values -

replace

Maths Function

abs

exp

factorial

sqrt

cbrt

power - power

log10 - log10

log2 - log2

log - log

log1p - log1p

MOR

explode

explode_outer

posexplode

posexplode_outer

flatten

Formatting

format_number

format_string

JSON Functions

from_json

to_json

Json_tuple

schema_of_json

get_json_object

Distribution

repartition

coalesce

Extraction

CSV txt parquet JSON ORC Hive JDBC