

SPARK

- A distributed computing engine, which distributes our data to process it.

Ex

(1GB DATA)

Architecture
(where it distributes data among machines)

IPC → Increase resources but upto a limit

[Cluster]

Connect all your computer classes machines → tell them to act as 1 machine

→ Connect them together

→ no saturation point / no limit

cluster of machines

(Node ↔ Computer in spark language)

[Worker Node]

Driver program

Cluster Manager

Worker Node

will orchestrate the tasks
(break down the code into smaller tasks & send that info to →)

will create worker nodes based on info & those →

will actually process our data.

(Manage Resources)

① Create a driver node on our cluster

transformations)
processing/
everything

② When we submit our code, it will go to driver node

③ All breakdown info which we provided to (CM) from (D-P) →

will go to these 2 workers

- ④ Asses/break code into stages/jobs/tasks & hand it back to Cluster Manager
- ⑤ Requests (C-M) to create (2 worker Nodes) based on info

⑥ Execute all transformations

⑦ Hand over the data to these 2 workers

Why spark?

- Hadoop uses Disk memory which took lot of processing time
- Spark uses "in memory computation" (RAM) → faster
- Fault Tolerance
- Partitioning
- Lazy Evaluation

[HADOOP]

Disk Based memory

- Like a Notebook, you write results down, then look at them again next time.
- Writing / Reading takes time.

[SPARK]

In memory

- Like a whiteboard (RAM)

you can write, read, erase, update instantly
w/o going back to your notebook

RAM is designed for temp storage & ultra fast access

↳ Can run on your PC / Single machine / No point

↳ Can run on cluster of machines → On premise / cloud

Each machine has its own RAM, CPU

Spark distributes data across nodes

[On premise]

↳ hardware & servers are physically located at your company

Your (IT dept) manages network, security, storage, electricity, cooling

→ scaling is slower, limitations, expensive

[Cloud]

↳ Computing resources are hosted & managed by cloud provider

Expensive

Lazy Evaluation

Ex:

[DATAFRAME]



And we perform some transformations.

Transf 1 (selecting col's)

Transf 2 (adding 1 row)

Transf 3 (deleting 2 row)

Spark will store all these into a Logical Plan
to optimize it.

↳ Once we trigger the action → it will execute thru logical plan.
performing all transactions

• Remembers to build a plan

(combine steps together)

Choose the fastest path (DAG Optimization)

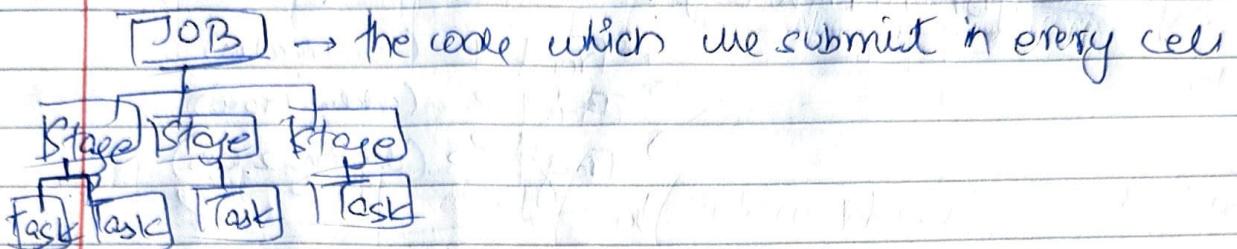
When we record a plan

- ① When we write transformations (filter, group) ↗ 1 direction
No loops
- ② Spark records them as a plan → Directed Acyclic Graph.
DAG is basically a map (spark's map) of all steps
it needs to take.
- ③ Spark's optimizer → looks at the plan & asks
"Is there any better/faster way to get same result."

DAY 2

7/10/25

Hierarchical Structure to execute jobs



SPARK API's

allow us to write code in our native language
(depending on what developer is comfortable with)

A set of rules/definitions/tools, that allows different software systems to interact with each other.

Ex → B/w your code (client) & website (server)

Python + Spark → PySpark library

Modern way of making table

→ Delta Lake is default storage format for tables in databricks

→ It takes raw data → stores it in optimized, reliable format

∴ Creating a Table → means creating of Delta Table

Other way

→ Create a Volume (secure storage layer)

(we get a temporary df in memory)

→ upload CSV → & then read

DATA3

8/10/25

Code

- `spark.read.format("csv")`
- `-option("inferSchema", True)`

Automatically detects datatypes (int, float, string)

- `-option("header", True)`

Treats first row as column names

- `-load('Volumes/workpace---path')`

To trigger a action

displays text based preview

- `df.show()`

// it will actually execute this time

- `df.display()`

// display beautiful, interactive table
letting you to sort/filter/charts etc.

For JSON

- `-option('multiLine', False)`

Schema

→ df.printSchema()

Create

`my-dll-schema = "`

`Item_weight STRING,`

`Item-MRP INT,`

`"`

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

~~option~~ `.schema(my-dll-schema)\`

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

`.load()`

Select

from pyspark.sql.Junctions import col

df_sel = df.select('Item_Weight', 'Item_MRP')
OR
display

df_sel = df.select(col('Item_Weight'), col('Item_MRP')).display

Alias

df.select(col('Item_weight').alias('G')).display()

Filter

df.filter(col('Item_Fat') == 'Regular').display

df.filter(col('Item_type') == 'Coke' & (col('Item_MRP') <= 10))

df.filter((col('Size').isNotNull() & col('Outlet').isin('G'))
OR
col('Size').isNull() & col('Outlet').isin('Tier 1', 'Tier 2'))

9 Oct
Day 4

Rename

df. with column renamed ('Weight', 'Weight')

Create / Modify column

df = df. with column ('Flag', lit('new'))

↳

Flag
new
new
new

df. with column ("multiply", col("Weight") * col("MRP"))
• display()

Replace / Edit content

df. with column (del('Fat'), regexp_replace(col('Fat'), "Regular", "Reg"))
• with column ('fat', regexp_replace(col('Fat'), 'Low fat', 'LF'))

WEEK

Day 5

Type Casting

`df = df.withColumn("Weight", col("Weight").cast(StringType))
 .cast(IntegerType))`

Sort / Order By

~~`df.sort(@"Weight")`~~ `df.sort("Weight")`
`df.sort(col("Weight").desc())`

`df.filter(col("Weight").isNotNull())
 .sort(col("Weight").asc())`

`df.sort(col("Weight").asc(), col("MRP").desc())`

limit

`df.limit(10).display()`

Oct 11
Day 6

Drop

df.drop (col('MRP')), display()

df.dropDuplicates(), display

→ df.dropDuplicates (subset = ['MRP']), display ()
OR
→ df.distinct (), display()

df1 = [(1, 'Guru'), (2, 'Hari')]
schema1 = id INT, name STRING
df1 = create DataFrame (df1, schema1)
(df2
schema2
df2)

→ df1.union(df2), display()

df1.unionByName (df2)

df1.unionByName (df2), allowMissingColumns=True

Oct 12

Jay String

df.select(initcap('Item-type')).display()

lower ('')

UPPER ('')

y.with('date-format')

.alias('upper-Item-type')

Date

df.withColumn('Current Date', current_date())

df.withColumn('Week After', date_add('Current Date', 7))

df.withColumn("Week-Before", date_sub('Current Date', 7))

OR

date.add('Current Date', -7)

df.withColumn('Date-Diff', datediff('Current Date', 'Week After'))

modify existing

df.withColumn('week-before', date_format('week-before', 'dd-MM-yyyy'))

(yyyy-MM-dd)
ORIGINAL

drop Na

df.dropNa('all')

df.dropNa('any')

that rows which had

'remove null' from all
cols

→ any col'n

→ df.dropNa(subset = ['order-type'])

Filling Nulls

df.fillna('NA')

df.fillna('NA', subset = ['Outlet-type'])

Split Inverting

① df.withColumn('Outlet', split('Outlet', ','))
↓ ['Outlet-A', 'Type-A'] (delimiter)

② df.withColumn('Name', col('Outlet').getItem(0))
• withColumn('Type', col('Outlet').getItem(1))

or
df.withColumn('Outlet', split('Outlet', ','))
(only give type)

df.withColumn('Flag', array_contains('Outlet', 'Type1'))

groupBy

df.groupby('Item').agg(sum('MRP'))
avg

df.groupby(['Item', 'Size']).agg(sum('MRP')).alias('...')

df.groupby(['Item-type', 'Size']).agg(sum('MRP'), avg('MRP'))

14th Oct

Day 9

[SPARK]

collect list

df.agg().groupBy('user').agg(collect_list('book'))

pivot

df.groupby('Type').pivot('Size').agg(avg('MRP'))

when otherwise

df.withColumn('Veg', when(col('Type') == 'Meat',
 'Non-Veg').otherwise('Veg'))

inner join

df1.join(df2, df1.dept_id == df2.dept_id, 'inner')
 .drop(df2.dept_id)

df1.join(df2, on = "dept_id", how = "inner")

left → All from left | Both cols

left semi → All from left but matching
| And only cols from left

Inner join + left cols

left anti → only what is not in right

16th Oct

Day 11

Window functions

```
from pyspark.sql.window import Window
```

- `df.withColumn('row_col', row_number().over(Window))`

1 `(Window.orderBy('Item')).display()`

Dense Rank	Rank	Row Number
1	1	1
2	1	2
3	2	3
2	2	4
3	3	5
1	4	6
4	4	7
3	5	8
2	6	9
3	7	10
2	8	11
3	9	12
1	10	13
2	10	14
3	11	15
1	12	16
2	12	17
3	13	18
1	14	19
2	14	20
3	15	21

~~df.withColumn('rank', rank().over(Window.orderBy(col('e').desc)))~~

`df.withColumn('l', row_number().over(Window))`

- `partitionBy('Item').orderBy('MRP'))`

18th Oct

Day 13

STOPAGE

Data writing → Azure → Data Lake → AWS → S3

provide a dedicated location

df.write.format('csv')

.save('File Store/tables/csv/date-csv')

→ Append

df.write.format('csv')

df.write.mode('append')

.format('csv').save()

Schema must match (same cols, data types)

df1 data = [(1, "TV", 500)]

schema1 = "id INT, product STRING, price FLOAT"

df1 = spark.createDataFrame(df1, schema1)

→ df1.write.mode("overwrite").parquet("/mnt/date")

{
date2
schema2
df2}

df2.write.mode("append").parquet("/mnt/date")

{ df-all = spark.read.parquet('/mnt/date') }

df-all.show()

1 DataFrame [multiple files]

/mnt/data/

part0000.parquet

(contains 1, TV, 500)

part0001.parquet

(contains 2, Cell, 100)

CSV : human readable → Row Based

Parquet : Binary, not human readable

Faster, for large scale datasets

Fetch only necessary columns ← Columnar storage format

↳ allows for more efficient querying

→ But no transactional control

can't edit → only rewrite

Read Only, Immutable, ↗ CSV ↗ Parquet

Delta : Its a layer built on top of Parquet that adds

↳ ACID transactions

Real D
Base

↳ Optimised file management

↳ Time travel

Because of JSON log

JSON log → A folder containing files representing versions of your tables.

every time you update/delete

↳ delta creates new log file

↳ never overwrites old files

like Git → Version + history

10. DataFrames

Writing modes

↳ append, overwrite, ignore, error

SQL in Pyspark

- We read data in pyspark using dataframe

↳ A python object

- For SQL engine to read the data we need to create temp view table.

df.createTempView("sales")

OR df.createOrReplaceTempView("sales")

- But when notebook session ends → temp view disappears

→ If we want permanent table

↳ df.write.saveAsTable("bigmart.sales_table")

(1) spark.sql("Create Database if not exists bigmart")

- Convert SQL Query to dataframe

df_sql = spark.sql("Select * from view where ...")