

Python & Spark → PySpark

1. GCP DataProc
2. Terminal of DataProc
3. Google Colab
4. Databricks

Python map reduce higher orders fn

Reading & Partitioning Data in Spark

Data is read and divided into partition

What is partitioning in Spark?

> Logical chunks of data

Processing

large dataset is divided into smaller chunks
called partition.

1. When already divided
into blocks → hdfs, S3
lake gen 2

2. read data & create partition
based on the config.

Execution of one partition of data = one task
= 1 CPU core required for

Execution \approx Mapper

(executor) \approx Reducer

performs a task

& it can have
more than
1 core as
well

8 blocks \rightarrow 8 partition \rightarrow 8 tasks

($\lg b$)

inc.

dec

i will
get more illism

less parallelism

Spark parallelize

We can create an RDD using
parallelize & then work on it

Subset of data



< 128 mb

↓
2 partition \Rightarrow 2 cores



cores

+ partition \Rightarrow better results

→ Read the data → Parallelize & partition, counting value

→ Transformations → map, reduce, filter, reduceByKey

Reduce By Key

transformation

Count key value

Action

Notebook → functions in spark

reduceByKey 5 keys \rightarrow 5 values

map 100 rows \rightarrow 100

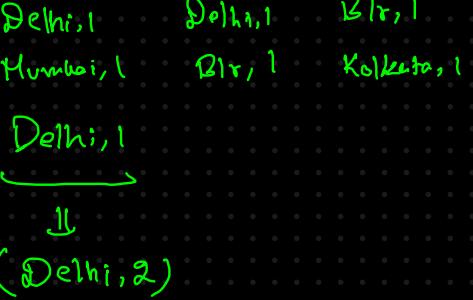
filter 100 rows \rightarrow 0 - 100

This is still the most diff. way of writing code in spark
 \Rightarrow more easy way

• Save as text file (hdfs)

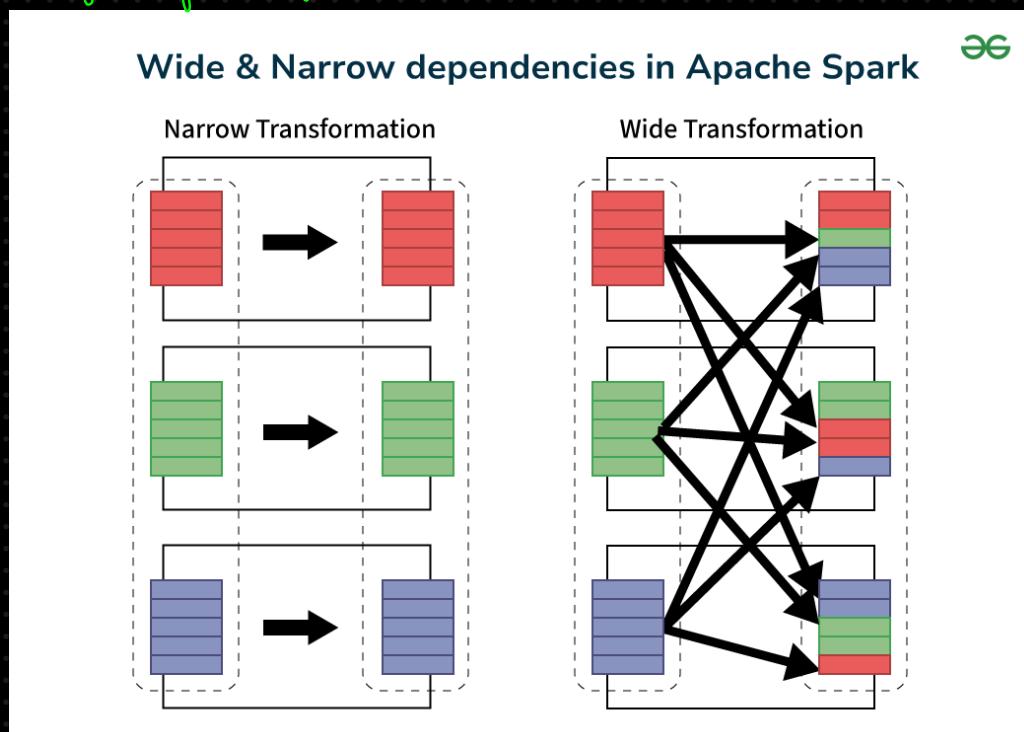
Transformation in Spark : Narrow vs wide

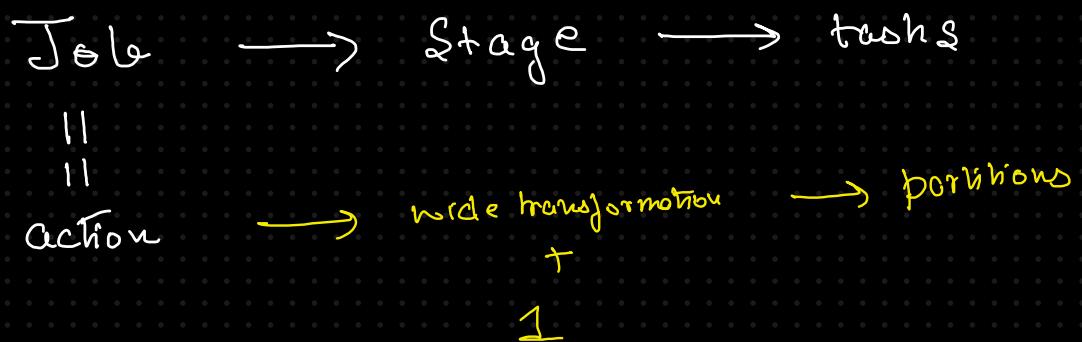
transformations are operation that create a new rdd from existing one

Narrow	Wide	Broad
data is not shuffled across partition	Data is shuffled across our partitions	
faster → no shuffle	Slower - as shuffling ⇒ I/O network	 Delhi,1 Delhi,1 Bihar,1 Mumbai,1 Bihar,1 Bangalore,1 Kolkata,1 ↓ ↓ ↓ (Delhi,2)
map, filter, flatmap dependent on one parent position	map groupByKey reduceByKey SortByKey dependent on multiple position	

either

- try to have very less wide transformation
- narrow down the data so that data getting shuffled is less.





When we do a wide transformation \Rightarrow a stage is created



\rightarrow homework

\rightarrow just run what I did & understand spark UI

\rightarrow be comfortable in pyspark

from next class, lots of code without major error.

Reduce By Key & Group by Key