


Apache Spark

| Start at 9:05

- * Agenda →
- What is Spark
 - Distributed Computing
 - PySpark

- ⇒ • Typical Data Size →
- { → Few mb → 1-100 mb → 10,000 Rows
- { → 3-5 GB → 1M-30M Rows

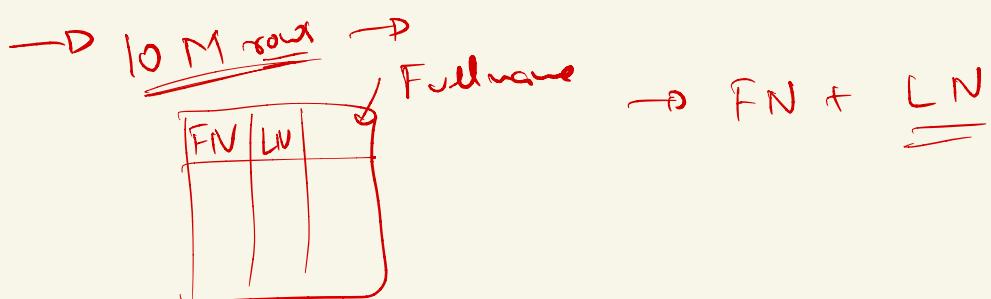
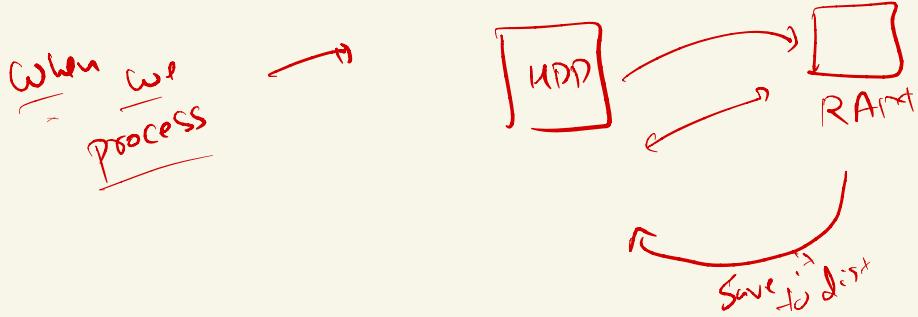
- { → Amazon Order Table → 100M Rows

- Constraint → Time → Disk I/O Performance
- Memory Constraint

→

- ⇒ RAM → 16 GB -
- HD → 1TB ←

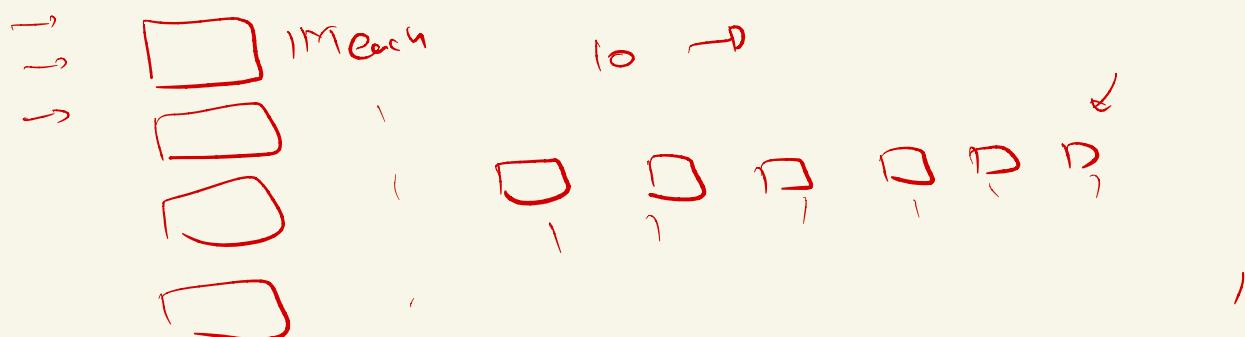
10GB File → Store it on Hard Disk →



→ 1 TB RAM >

(60 TB HDD)

→ applied line by line | row by row

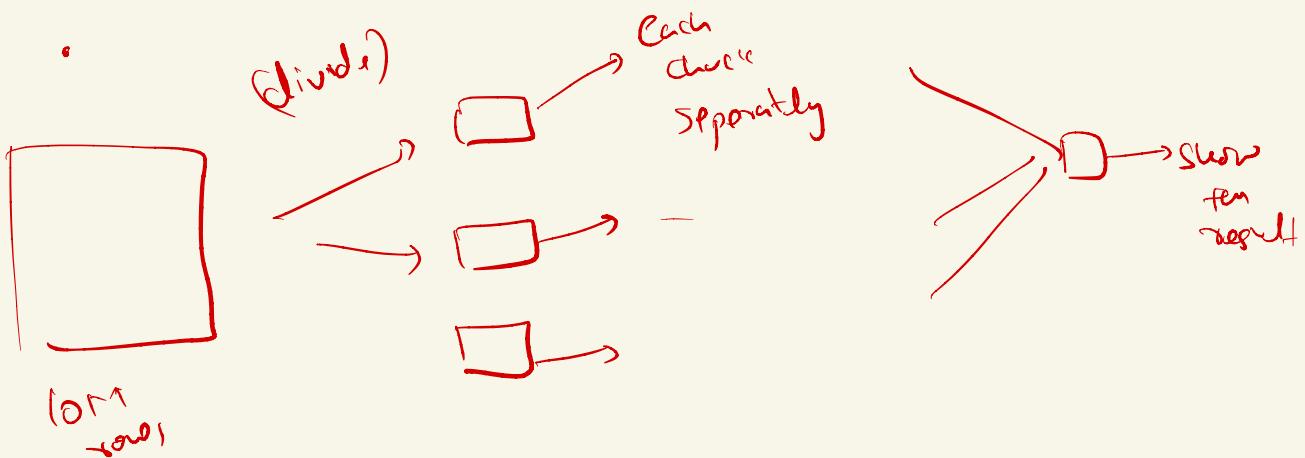


→ Lack of parallelism →

→ python (pandas) → only using

→ Big Data → unstructured data.

→ • Parallel Computation



→ Spark → Very easy to use • mature
◦ Fast

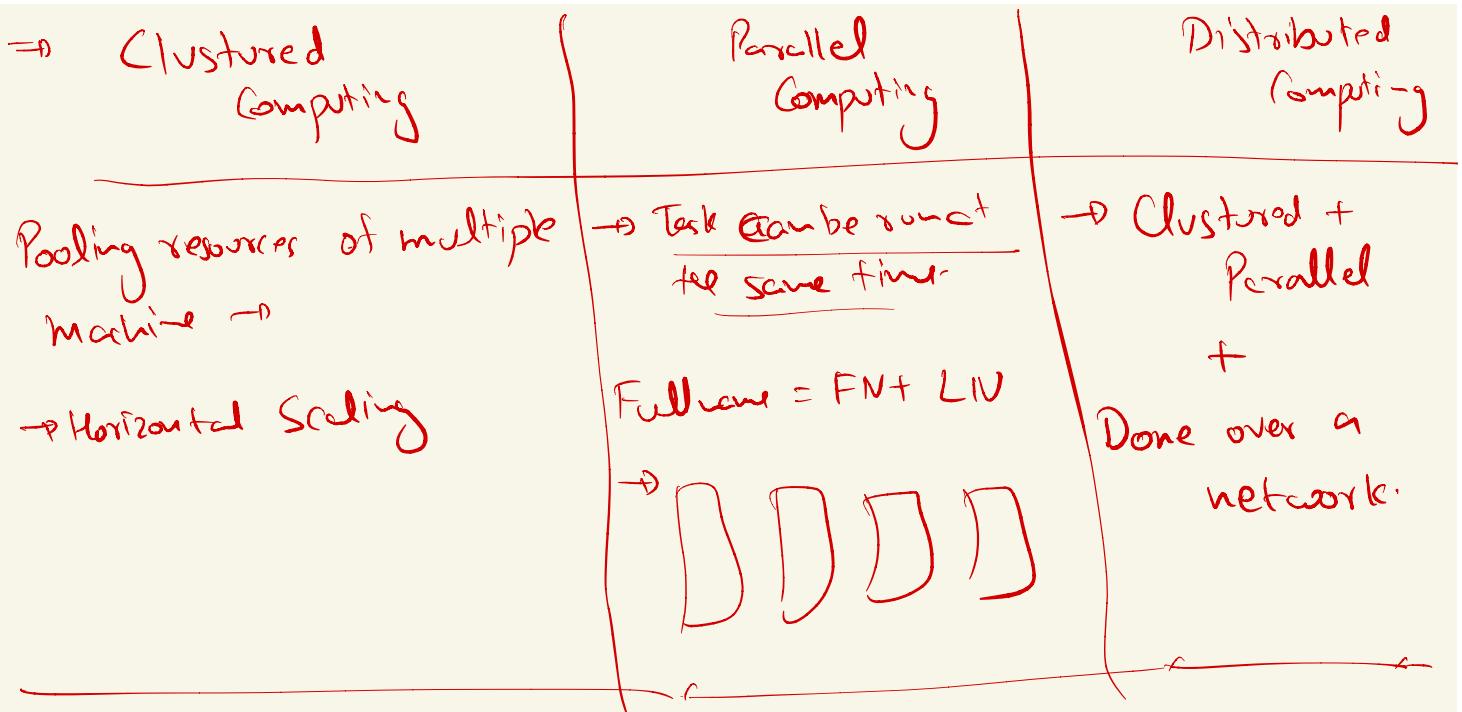
→ open source, distributed Computing framework



→ 100x → Traditional system (Pandas)

◦ Spark → Connect to almost all data source

↳ almost
all
major language
→ Java
• Python



\rightarrow RDD \rightarrow Resilient Distributed Dataset

\hookrightarrow How Spark stores the data

\rightarrow Pyspark \rightarrow built on top of Spark framework

- Handle much more data than pandas
- ML Task + Python Task

\rightarrow Stepper Learning Curve

• Setting up Spark is very painful

• Pandas is sufficient for most small cases

Desire \rightarrow even 5x Pandas

Hadoop \rightarrow Storage framework

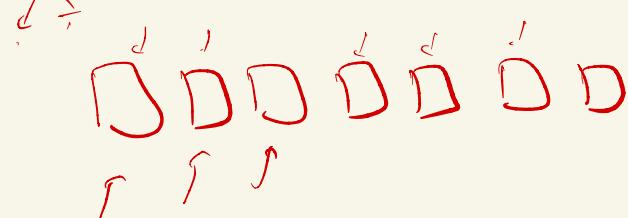
→
Other
framework

→ Mapreduce
→ Disk based
processing
 $99 \times \text{I/O}$
Slow

Spark
→ In memory processing
→ Fast

→ Spark → Store the data in distribution

→ • Multiple storage computer
• Duplication → $(3 \times)$ →



• S?
• your own service
• HDFS

→ 10 TB → 10 machines →
→ 1 TB data to each machine → normal
→ 3 TB data to each machine →
→ duplication
replication

→ Spark Backtrace

↳ ① Distributed Computing

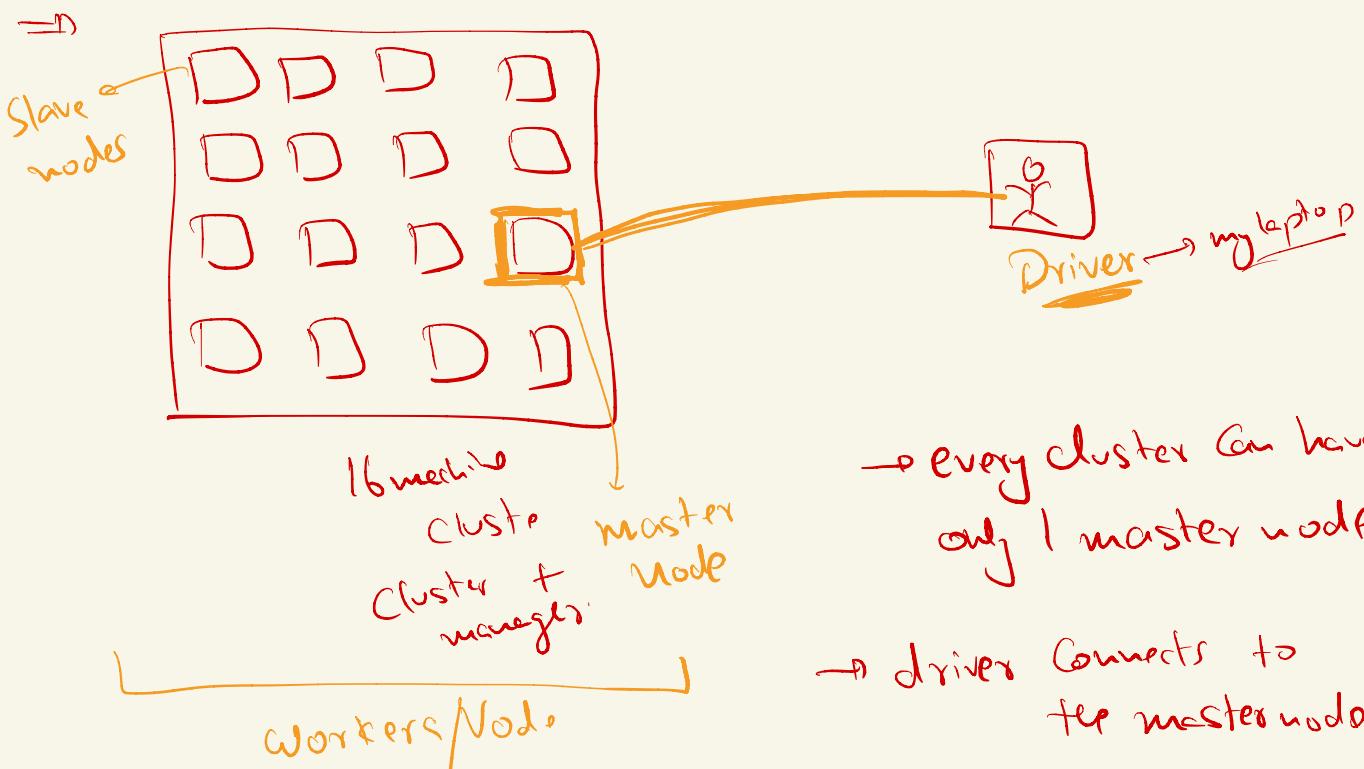
② Make replication of Data

⇒ • 30 years of data → 1

→ We will only require small part of this data on our local machine

→ • all your computing → on network / cloud cluster

Results get send to local machine.



→ every cluster can have only 1 master node

→ driver connects to the master node

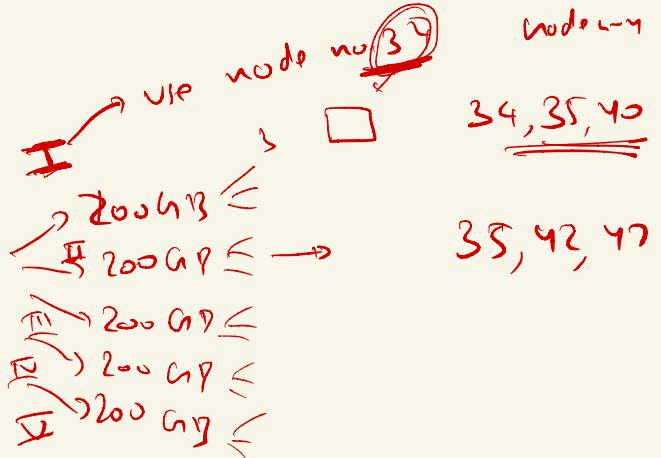
→ master node distributes work to other workers via execution plan

- Data Splitting → divide into multiple chunks.
 - Which column
 - How many parts

- Task Distribution →

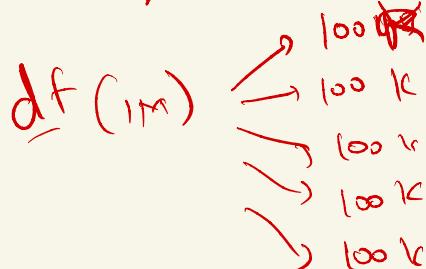
- In-memory Processing →

→ 1 TB data → Sort



- * RDD ⇒ immutable

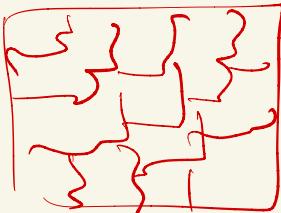
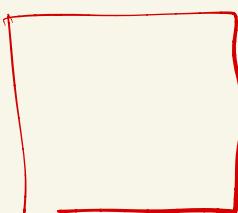
Resilient
Distributed
Dataset



Process

- Aggregation → Merging data from each worker node.

⇒



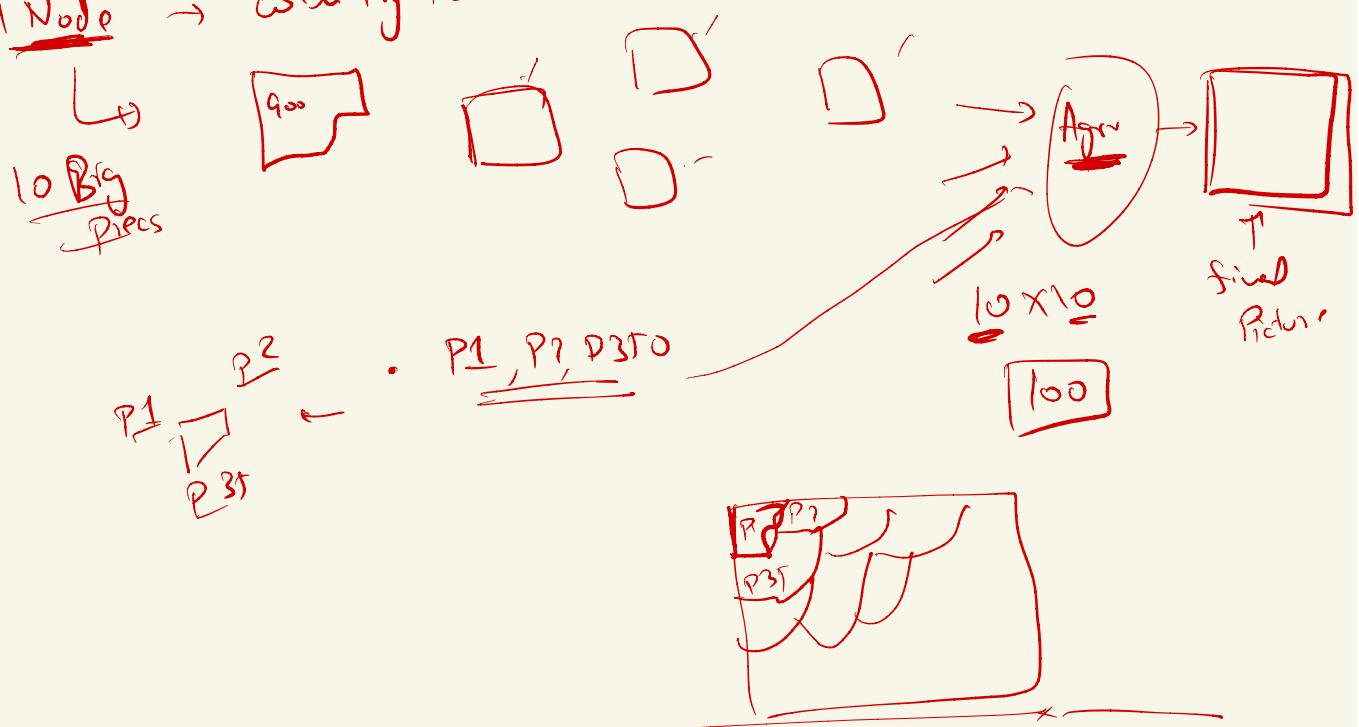
Create the original

→ look pieces to each

$\times 10$ nodes

D D P

1 Node → will try to solve for puzzle info



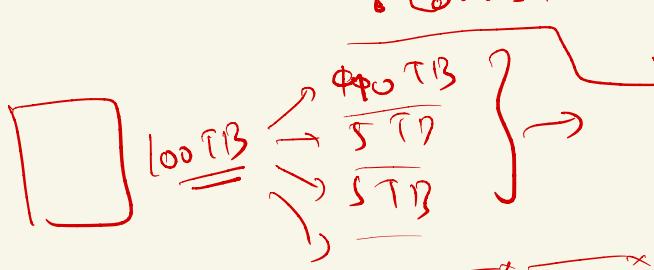
⇒ Partitioning

→ improving Performance

• Consistent data across node

• Help in faults tolerance

• Efficient Storage



→ Flipkart → day wise sales

⇒ Data for each day → Move → finding
28th at this data

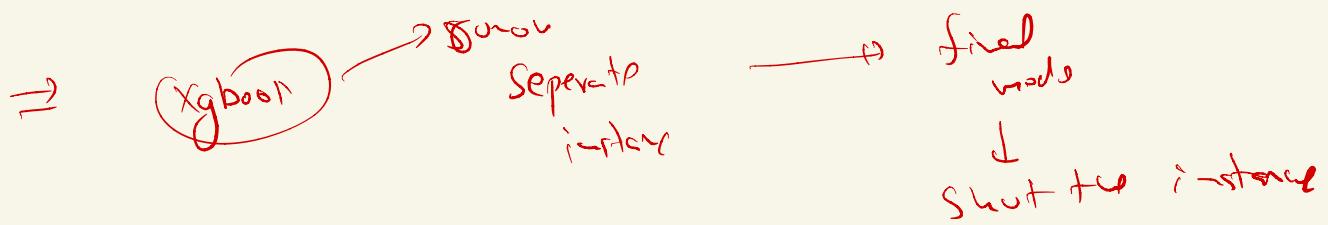
⇒ BB Sales Comp

→ 28th Jun
29th Jun

100 K →
1 Month →
100 K

Doubts

We used sagemaker's xgboost and trained it using a different instance to reduce cost. But for models like KNN, DBSCAN, agglomerative clustering, for which if you need the prediction, you need to retrain the whole model again, how does it work?



→ Instance you running
for Sagemaker → increased size → just normal code

→ • → Sagemaker instance → Pythonic operations
Basic EDA

→ Training instance →

→ deployment → Model → 1 → 100k visits / day
2 → 10k visits / day
3 → 100 visits / day

	day to day Use	Featureview
• Git	✓	-
• Flask	✓	-
• ML System dev'n	-	✓
ML Flow	✓	-
Spark	-	✓
Docker	-	-