

* Data migration to the cloud

⇒ solⁿ around the digital workspace information management archiving & migration space

* Spark Core & SQL

⇒ Hadoop, MR, Ecosystem Component
↳ Spark
↳ philosophy, features, components.

Core programming → Python → pyspark
↳ (Major focus)

* Use Cases :-

* Hadoop ⇒ Open Source framework.
↳ processing big data.
↳ Additional layer on top of any existing system.

* Core Components

⇒ HDFS → Storage
⇒ MapReduce → process data.
↳ processing framework

⇒ Dist Appⁿ in form of cluster.

*] Components of Hadoop

Ecosystem →

Hive	→ Run SQL like queries on data in HDFS → <u>Query data</u>
Pig	→ Scripts for Data Cleansing (Data pre-processed)
Sqoop	→ Database Migrator
Mahout	→ Machine Learning Lib

*] Spark

⇒ (2006) → Hadoop ⇒ HDFS + MapReduce

(2006 - 2018) ⇒ Hadoop → Best Batch processing in the market.

*] Limitations of Hadoop MapReduce

HDFS → Data Storage.

- ① High Latency → Batch processing Slowing down
(Wait) (Not developed)
- ② Architecture → Different operation → Different tools.
- ③ Processing → Disk based ⇒ Involves DISK I/O ops

* MR ⇒ Disk based Processing
→ High Latency.

* Spark ⇒ General purpose, Lightning fast
Cluster computing framework.

↳ Open source, Wide range Data
processing Engine.

* Spark = Batch processing & Streaming processing

* Types of processing

↳ Batch
↳ Real time (Streaming)
↳ Interactive
↳ Iterative.

* General purpose ⇒

* Features ⇒ Single library for SQL
→ Graph Analytics, Streaming &
ML Algo.

* Spark Competent to Hadoop
↳ Memory based processing
→ Open source computing framework

Spark ⇒ Core API

↳ Programming in Spark.

Spark ⇒ Speed
100x faster
than traditional
frameworks
Hadoop MR

⇒ Powerful Caching mechanism

[*] Deployment

⇒ Mesos, Hadoop, via YARN

Spark's own Cluster Manager.

[*] Real Time ⇒ Low Latency by using in-memory processing (Computations)

[*] Ployglot ⇒ Scala, Java, Python, R prog.
(Spark)

Python + Spark ⇒ PySpark

[*] Why Spark (2010) = Fame → 2014

[*] Components of Spark ⇒

① Spark Core

② Spark SQL

③ Spark Streaming

④ Spark ML Lib

⑤ Spark GraphX

⑥ Spark R

+ R programming → Data visualization.

* Spark Core & Spark SQL

- * Module 1 ⇒
- ① Spark Architecture
 - ② Spark Core API → ^{RDD} Prog.
 - ③ Operate on RDDs → Transf / Actions
 - ④ RDP Persistence
 - ⑤ Advanced Spark Concepts → Shared variables

* Module 2 + ① PySpark SQL - 20th Aug
②

* Spark Architecture

Hadoop → Dist. Env. → Multiple nodes
↳ Clusters

→ 3 Components

① Data Storage

② API

③ Resource Management → Spark can be deployed as stand alone service.
↳ Dist. Computed framework.

Cluster Mgr → service → Master node.

* Spark Context ⇒ Entry pt of Spark functionality

* App ⇒ Driver program → (SparkContext object)
Spark → Standalone, YARN, Mesos. ↳ Cluster Manager
we school
S.P. MANDAL'S
Wellngkar Education

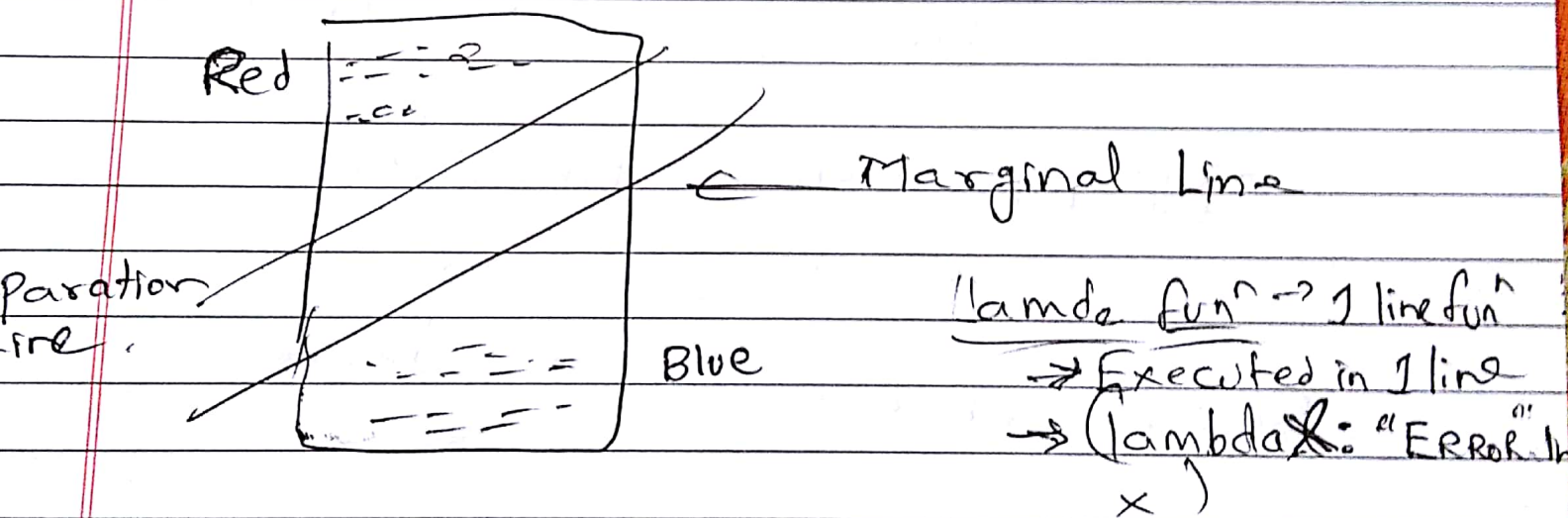
* Spark RDD (^{fault tolerance} Resilient Distributed Dataset)

- Fundamental unit of Spark programming
- Fundamental Data structure
- Immutable collection of objects
- ↳ value of the dataset can be changed.

Dist → Multi Node Existence

↳ Dataset → Represents records of the data.

SVM → predict diabetes, BMI etc



* Bagging Technique

* RDD → Data in Stable storage, Other RDD
 → Parallelizing already.

In memory proc → 100x times faster
 Disk based proc → 10x times faster.

Softmax
Pytorch

map() \Rightarrow Iterates over each line in RDD, split into new RDD

Filter() \Rightarrow Set of values & principles
 Belief, decisions \rightarrow Standup meeting

Virtual machine & Ubuntu

Customer collaboration
Contract negotiation

Desktop

Anaconda download, Left side of the list more items

Spyder (Python 2-7)

Tableau
Ubuntu
Pytorch & Kivy
Anaconda
Virtual machine

Kivy

RDD \Rightarrow In memory computation capability

MapReduce \Rightarrow Operation (Function) \rightarrow RDD
 funⁿ, Appⁿ, Job

* Dataset 1 \rightarrow MapReduce Job - O/P adable

per team
2 ops

\Rightarrow Collaboration & Productivity

\rightarrow Automating

Inv & Interaction

① Infra
② Workflow

process & tools

weschool
Welingkar Education

③ Metrics

Work flow, Comp Docⁿ

④ Testing

Machine Learning

→ News Categorization - Google
Recommendations - Amazon

→ Deep Learning Algo, ^{Complex} ~~Deep~~ Fusion
Clustering

- ML - Computer to work without being explicitly programmed.

* Prediction Accuracy

Mean abs % Error

100 - mean abs % Error = Prediction

GLM → Generalized Linear Model
↳ Family of Algorithms.

Bernoulli → Versicolor

SUM → Petosa Non-setosa

→ separating line ⇒ Greater Segregation

Greater gap → higher prediction activity