

MapReduce was not sufficient
It takes lots of time for
processing data as it
works on batch processing
So, the Spark came to over-
come all the problems

Spark uses in memory
computations i.e. It takes all
the data in RAM &
process there where as
in Hadoop first of all
data is stored in hard-
disk and then it is loaded
to RAM for processing, which
takes lots of time. If
the size of RAM is less
than data then most of
the data is stored in RAM &
rest in the hard disk.

In Hadoop, if we want to
do different tasks we have
to install different tool
where as Spark is can do

different processing like

1st) Batch processing

2nd) machine learning

3rd) Streaming data

or Graph data processing

Spark can take data from all
most all the platform &
can process data of any
type.

Components of Spark

- * Spark Core :- Contains the
basic functionality of Spark
like task scheduling, fault
recovery etc.
- * Spark SQL :- It is Spark package
for working with structured
data.
- * Spark Streaming :- Enable processing
of live streams.
- * MLlib :- Spark library which
provide multiple machine-
learning algo.
- * GraphX :- Library for manipulating
graphs.

Basics of Spark

RDD :- It is a Resilient Distributed
Data set. RDD is a Dataset which
is distributed among clusters &
have the power of recovery of
the fault.

Page:

Date: / /

Features of RDD :-

- 1.) Fault tolerant, If the RDD is lost then we can gain data from previous RDD.
- 2.) RDD is immutable i.e. we can not change the Data of RDD.
- 3.) Lazy Evaluation :- Spark uses lazy Evaluation it does not execute anything until it is asked to do so, it makes only DAG.