

Data Exploration and Analysis with Apache Spark

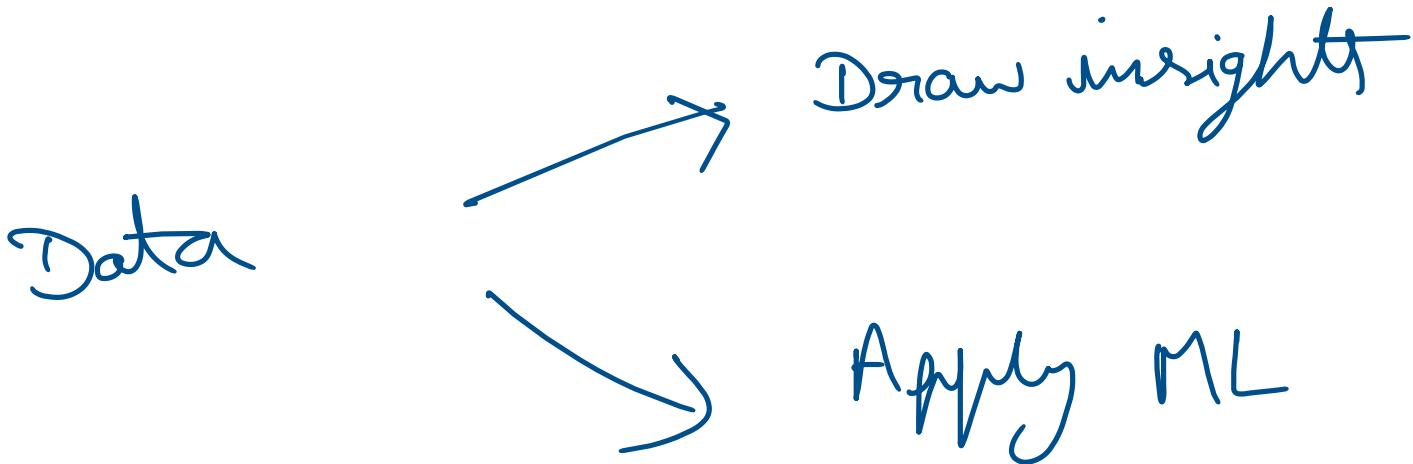
There is no shortage of data today, Public or Private

Public
Comments
Reviews
Tweets
Blogs
Gov Reports



Private
Purchases
Reviews
Clicks
Messages
Surveys
Questionnaires

Untold opportunities:



Doing Something Meaningful with data

- unstructured
- Many
- Semantically Complex

Data processing

Structured
clear
Easy to consume

Data Processing Task

Parsing fields from text

Accounting for missing values

Identifying and investigating anomalies

Summarizing using tables and charts

Small Data
Somewhat
Messy

Spreadsheets
low data collection frequency
10s - 1000s of rows per day
Sometimes involves manual data collection
Many many files

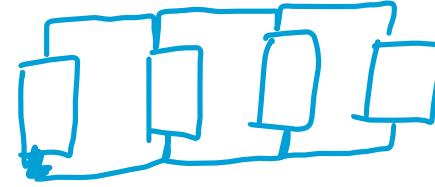
→ 

Medium Data
High data
integrity

Databases
high frequency of collection → SQL
~100K rows per day
programmatically collected
ACID properties

Big Data
very
messy

Distributed
Computing



very high frequency of data collection
millions / Billions of rows per day

files stored across a cluster of machines
Many many files (webpages, log files)



Tools for Data Munging

Small data → XLSX

Medium data → SQL, Python, Java, R

Big data → Hadoop Map Reduce, Spark

developed by open source
power house Apache

Spark
An engine for data processing and analysis

General Purpose

, Interactive

Distributed Computing

General Purpose :- Exploring
cleaning and preparing
Apply Machine Learning
Building Data Applications

Interactive :- Type Command and get output

len([1, 2, 5])

>>> 3

REPL : Read - Evaluate - Print-Loop

Interactive environment
for both Python and Scala

Process data
of Machines
across a cluster
Integrate with Hadoop

Distributed Computing :-

To work with Spark, we have Spark API

Spark APIs

Scala

Python

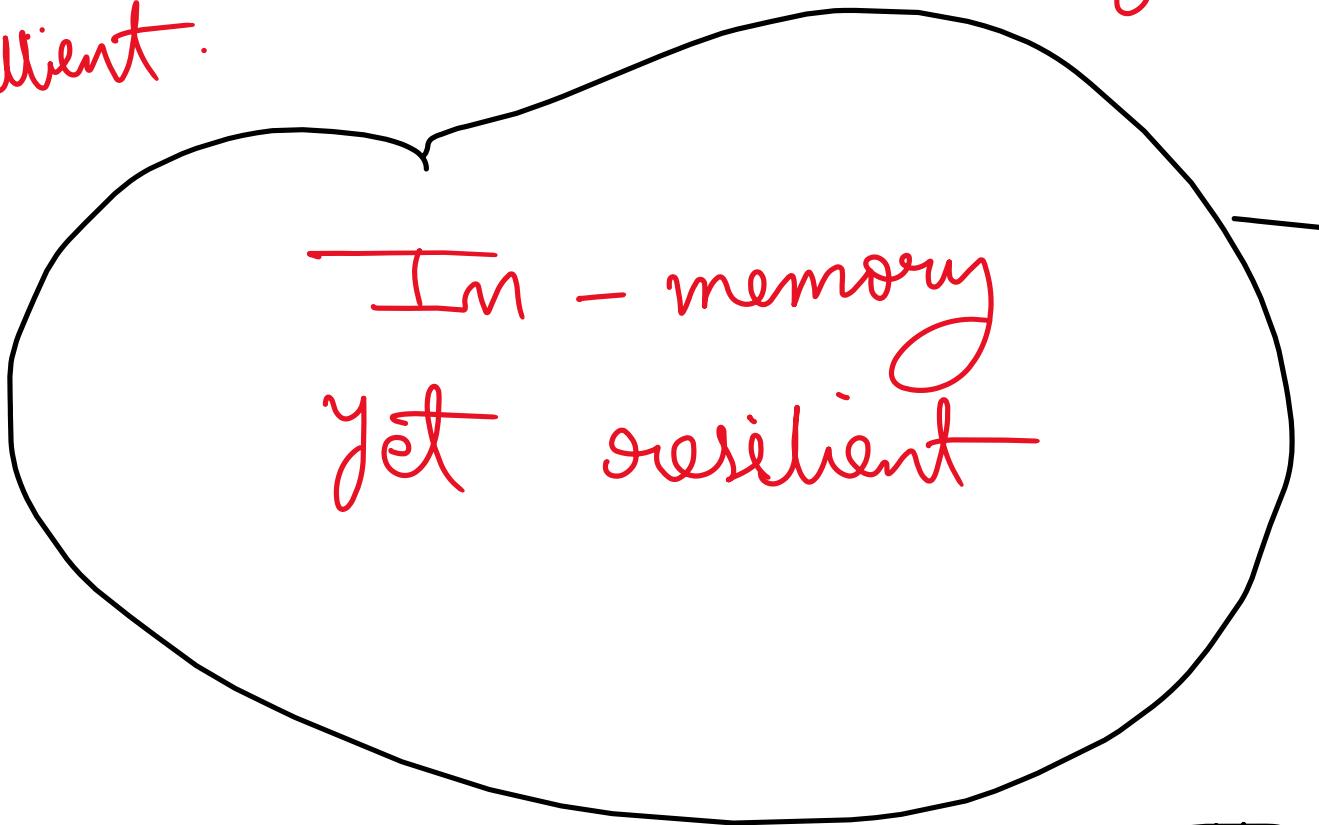
Java

We will be using Spark Python API

Almost all data is processed using
Resident Distribution Datasets (RDDs)

- RDDs are the main programming abstraction in Spark. They are how spark load data into memory and they processes them.
- RDDs are in-memory collections of objects. Each object in the collection represents one record in your dataset. and all records in the dataset are actually kept in memory across a cluster of machines.

what makes RDDs cool is that they are in-memory, yet they are resilient.



resilient
(able to recoil or
spring back into shape
after bending , stretching
or being compressed)

They are tolerant to falls
or crashes

With RDDs, you can interact and play with billions of rows of data

..... without caring about any of the complexities involved in resource management or fault tolerance :-

Spark is made up of a few different components : —

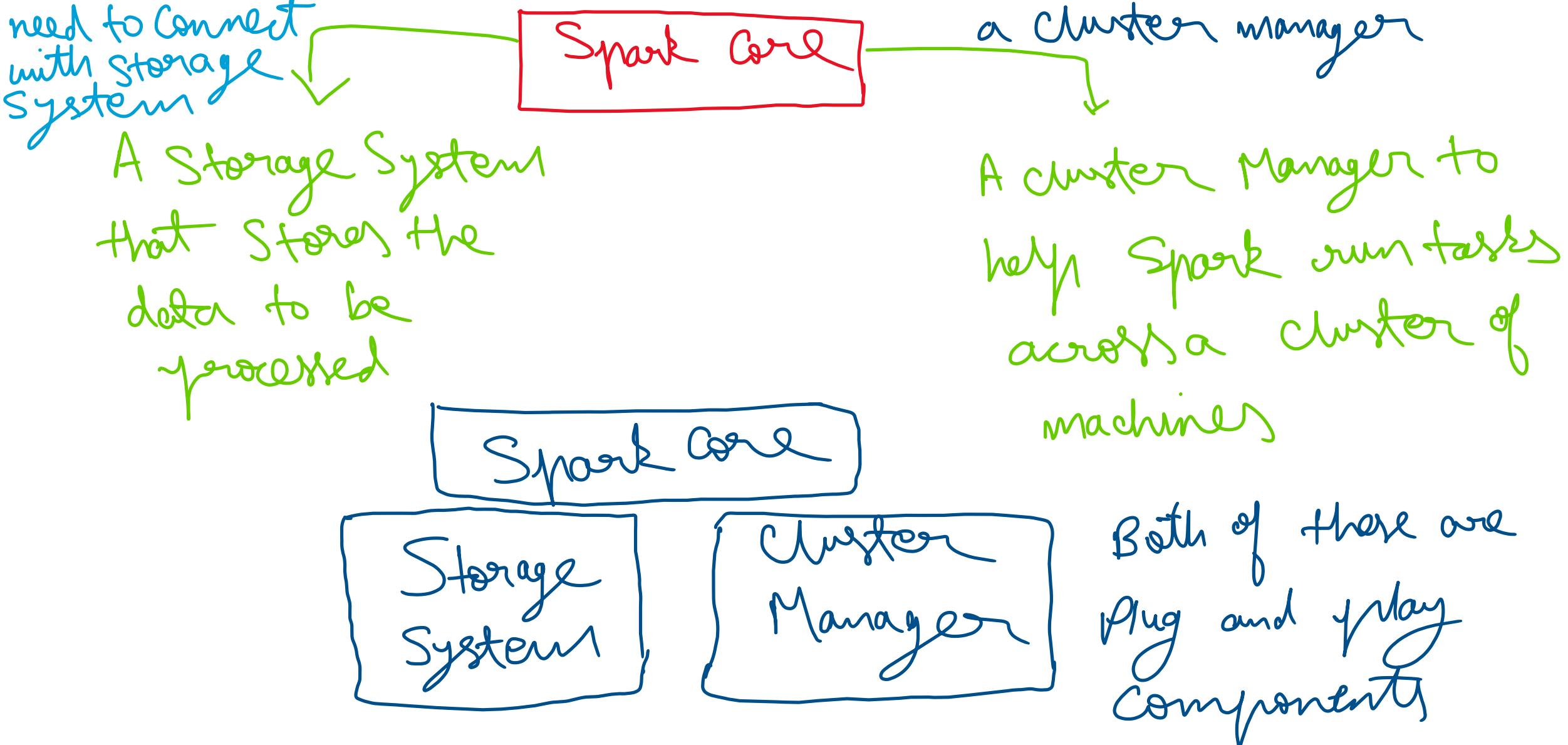
Spark Core



Handles basic functionality of Spark
including creating and processing
RDDs.

Spark Core is just a computing engine

Spark needs two additional components because otherwise
it is just a computing engine



Storage System

: Local file system, used in a standalone mode.

Cluster Manager

- Built-in Cluster Manager (comes in the Standalone mode)

- YARN (if you are integrating with Hadoop)

HDFS → which can be plugged in if you are integrating with Hadoop.

This plug and play capability makes Spark really easy to integrate with any existing systems that might be already using Hadoop.

Installing Spark

Prerequisites :- Java 7 or above

Scala

- Python (Anaconda)

after prerequisites → Go and download spark
binaries
update environment variables

Spark Environment variables

SPARK_HOME

Point to the folder where Spark has been extracted

PATH

· · SPARK_HOME · · /bin

Configure ipython Notebook for Spark

export SPARK_DRIVER_PYTHON = ipython

export SPARK_DRIVER_PYTHONOPTS = "notebook"

Then go to terminal → pyspark --master local[2]

Open Jupyter notebook
→ check sc command