

Low Level API | High Level API

RDD,

1.1X

Dataframe , Dataset

2.X - 2.1, 2.4 , 2.X

1.5 / 1.6

Not used

\* DF → RDD (Tight Control) → DF \*

Bank

DF & Dataset is  
on top of RDD

Legacy Code

Control  
(RDD)

RNO	Name	Age
1	A	10
2	B	12
3	C	14
:	:	:
;	;	;

part 1

part 2

part 3

Acnt	Deb/Cr
123	D 12
123	D 15
123	C 12
125	D 16

part 1

part 2

RDD  $\Rightarrow$  Resilient distributed dataset  
 $\rightarrow$  Recover quickly  
 $\rightarrow$  Imutable Val

Distributed Dataset

Var  $\Rightarrow$  var  $a = 5$  {• Mutable  
   $a = 6$  ✓ {• can change content}

Val  $\Rightarrow$  val  $a = 5$  {• Immutable  
   $a = 6$  ✗ {• can not change content}

$\rightarrow$  Partitioned

Spark

$\rightarrow$  1.x, Not much used  
  2.x, 3.x

$\rightarrow$  Programmers has to take care of Optimizat"

Optimization

RDD A      RDD B

join ②

filter ①

Transformation :

- \* Transform / Modify
- \* filter (Where)

Lazy : \* Transformation will not be executed until you execute action

Action :

- \* Causes execution of Transformation
- \* Done to a file / Display on Terminal.

