(High level of coding)

(High level of coding)

(Low level of coding)

$\rightarrow$ | Dataset |

| Dataframes |

[Dataset [Row]]

Dataset [Employee]

| RDD |

Dataframes is called as Dataset with generic type

Dataset is of specific type

Resilient Distributed Dataset

$\Rightarrow$ Advantage of RDD ( Resilient Distributed Dataset)

① Type errors are caught at run time (Compile type safety)

② Developers have flexibility in writing code, for example anonymous function.

$\Rightarrow$ Disadvantage of RDD, which gave birth to Dataframes.

① Low level of coding is used like using of map, flat map, filter etc. function.

② Optimisation was not there

③ Data was not stored in Tabular format

Developers were facing issue while writing a code for mugging up the low level code

# Dataframes

Data was stored in tabular format in the form of table rows & columns like a RDBMS table.

⇒ Advantage of Dataframes over RDD

① Code was written in high level construct i.e using where, groupby, select etc function

② in case of dataframes, since code is written in high level construct first it gets converted to low level coding at driver side, compiler performs some optimisation on it and then it sends to executor for execution. But if we have low level code, it directly sends to executor.

③ Dataframes provide Optimisation

⇒ Disadvantage of dataframes

① Compile time safety is missing, because of this runtime surprises comes

② no flexibility provided to developer in terms of like writing anonymous function

To cover the disadvantage of the dataframes, Spark developer brought to convert dataframes (Df) to Rdd (Rda) using df.rdd but conversion was not seamless and loosing the optimisation which was give by catalyst optimiser.

So Developer Optimised the Dataframes instead of rdd and converted to Dataset

⇒ Datasets

Datasets is basically take care of compile time safety as well as the flexibility of the code, where we can easily write Low level coding like lambda or anonymous function

so if we talk about Dataframes & Dataset there is slight diff. b/w both of them

Dataframes ⇒ Dataset [Row]

where Row is nothing but a generic type in which type will be bound at runtime.

Dataset [Employee] ⇒ Dataset

Dataset where type is bound at compile time