

## Structured API's session- 5:-

rdd vs dataFrame vs datasets

### rdd

when we deal with rdd. we deal with low level code

map

filter

flatMap

reduceByKey

This low level code is not developer's friendly.

rdd lacks some of the basic optimizations.

### dataFrames :-

Spark 1.3 version

higher level constructs which makes the developer's life easy.

### challenges with dataFrames:-

1. DataFrames do not offer strongly typed code.  
type error won't be caught at compile time rather we surprise at runtime
2. Developers felt that there flexibility (like to write own anonymous function)  
has become limited.

There was a way that the dataFrames can be converted to rdd.

df.rdd  $\Rightarrow$  whenever we want more flexibility and type safety.

1. This conversion from DataFrames to rdd is not seamless.
2. If we work with raw rdd by converting dataframe to rdd. we will miss out on some of the major optimizations.

## Dataset API:-

Dataset came into picture to address these challenges of DataFrame.

Spark 1.6 version.

### 1. compile time safety

with datasets type errors are caught at compiler level.

### 2. we get more flexibility in terms of using lower level code.

Conversion from DataFrames to datasets is seamless.

we won't lose of any of the optimizations.

Before Spark 2 Both dataframes and datasets are 2 different things.

In Spark 2 they merged these 2 into a single unified Spark dataset API (Structured API)

→ DataFrame is nothing but a Dataset[Row]

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

In case of DataFrames the datatypes are bound at runtime

however Dataset[Employee]

the type will be bound at compile time.

Dataset[Row] → dataframe (type error are caught at runtime)

Dataset[Employee] → dataset (compile time type safety)

How to convert dataframe to a dataset?

Q  
if we replace generic row with specific object then it becomes a dataset.