

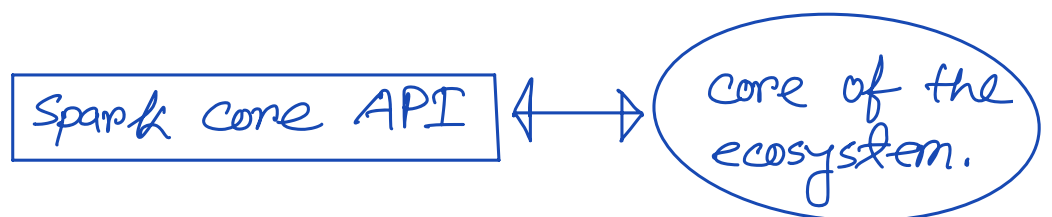
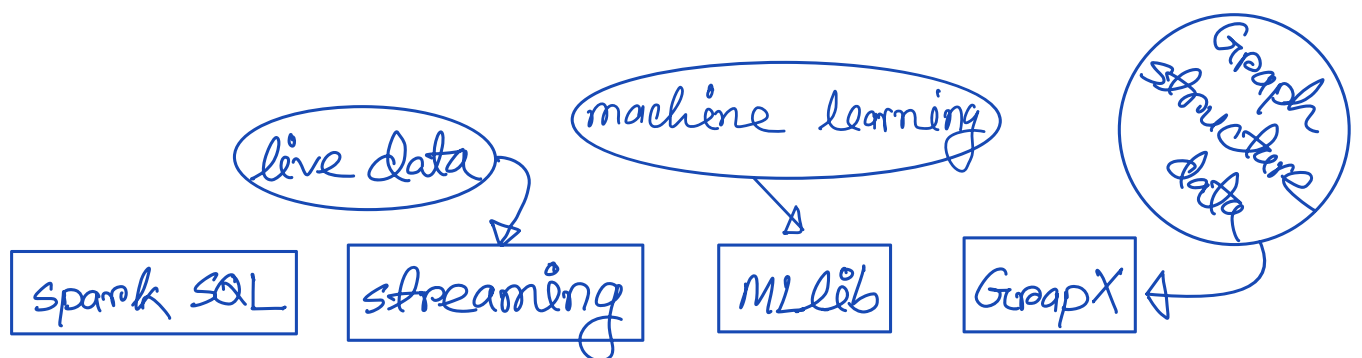
Course: Apache PySpark by Example

Practice environment: Google Colaboratory.  
Spark 2.3

## Why Ryspark :

- ① Speed.
- ② Easy to development.
- ③ Couple of language APIs.
- ④ Own eco-system.

## Apache Spark Eco-system :



standalone

## Responsibilities of Spark Core API:

- (i) Task scheduling.
- (ii) Memory management.
- (iii) Fault recovery.
- (iv) Interacting with storage system.

## Responsibilities of Spark SQL and dataframe:

- (i) Allow dataframe programming abstraction to execute query and make visualization.
- (ii) Spark SQL act as a distributed SQL query engine to intermix SQL's queries with the programmer data manipulation to solve complex analytics with SQL.

## Responsibilities of Spark streaming:

- (i) Process real time data.
- (ii) Analyze streaming and historical data.
- (iii) Use similar code for batch data and real-time data.

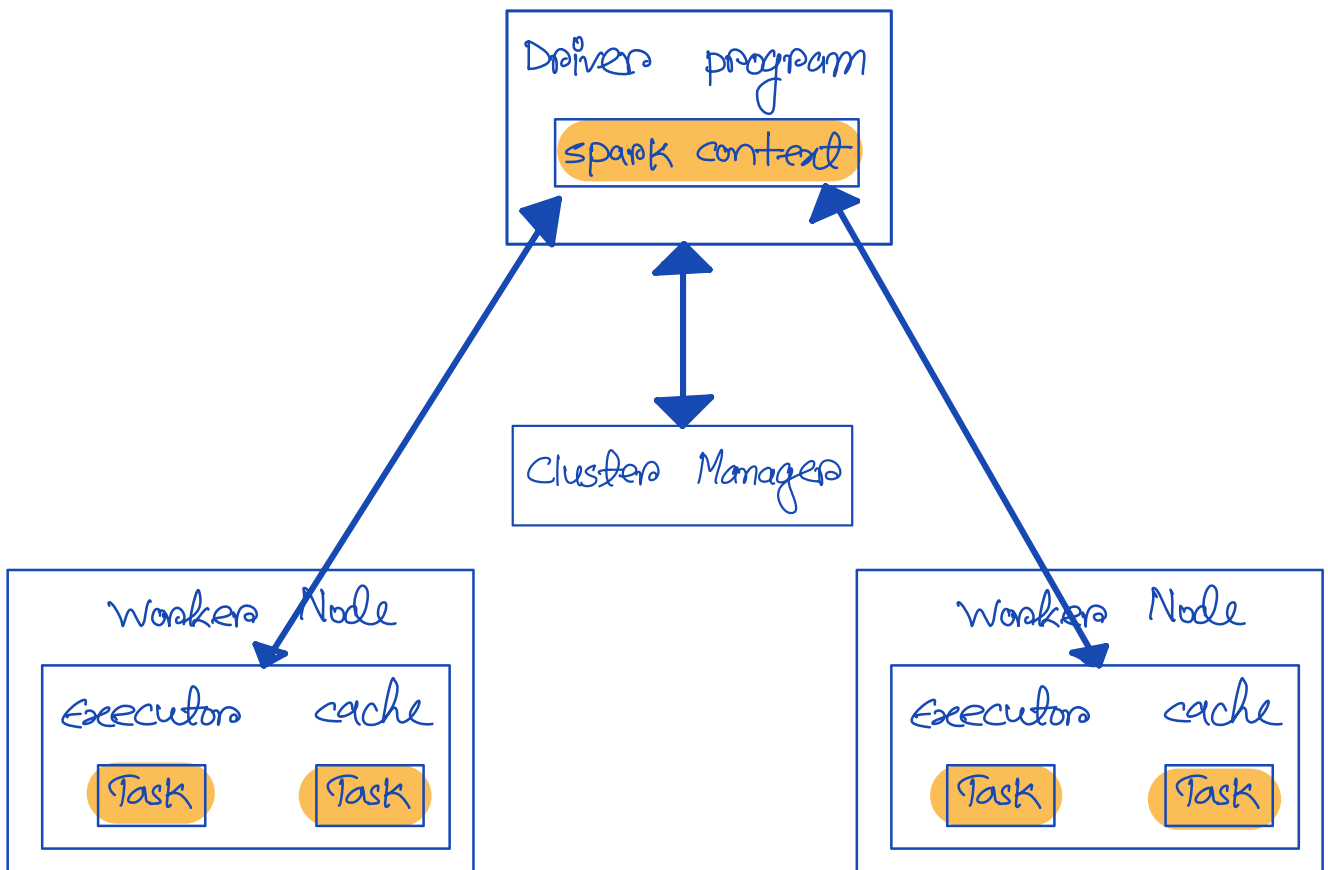
## Responsibilities of MLlib:

- (i) Provide scalable machine learning algorithm.
- (ii) Works in memory

## Responsibilities of GraphX:

- (i) Enable working with graph-structure data.
- (ii) Introduce a new graph abstraction, the directed multigraph.
- (iii) Useful for visualize graph like social network.

## Spark Components:



## ○ Different types of cluster managers:

① standalone.      ② Apache Mesos.

③ Hadoop YARN      ④ Kubernetes.

## ○ Spark workflow:

create spark session.

↳ in background PySpark uses Py4J to launch a JVM and creates a JAVA spark context.

Partitions: As Spark is a distributed system, we want the workers to work in parallel, and that's why Sparks need to break the data into chunks or partitions.

Actions: Three types of actions perform in spark.

① view data (.show()): fetch data to console.

① collect data (.collect()): collect data for drivers.

② write data (.write.format()): to write output sources.

Basic commands:

!wget url... : to download the file from remote source.

!ls : to see the file stored in the directory.

!mv : for renaming files.

This course focus on two major APIs:

① Dataframe API  
high-level API

② Resilient Distributed Dataset (RDDs)  
Low-level API

df.limit(n) : returns a new dataframe whose's

df.head(n) : returns an array

df.printSchema() : returns the schema of the dataset.

`df.column1` or `df['column1']` : to access a column.

```
df.select('column1', 'column2').show(3)
```

`df.withColumnRenamed('Existing Column', 'New column')`

`df.filter(col('colname') > 1)`  $\rightarrow$  conditions

`df.select('column-name').distinct().show()` shows all distinct elements from column-name.

```
df.orderBy(col('colnum'))
```

`df1.union(df2)`: concate two dataframe with same numbers of column and schema.

`df.count()` : return the number of rows.

## Join Operation:

df.join(df2, df.column == df1.column, how = 'inner').

`df.cache()` command for lazy evaluations.