

Spark working

- It has a small code base and the system is divided into several layers and the layers are independent to each other.

Apache Spark

- It is an open-source which is used for large-scale data processing.
 - Spark is fast, flexible, and easy to use.
 - Spark can run on single-node machines or multi-node machines (clusters).
 - Apache Spark can also process real-time streaming.
 - It provides libraries like Java, Scala, Python etc.
- PySpark features
- ① in-memory computation
 - ② distributed processing using parallelize
 - ③ used with many cluster managers
 - ④ fault-tolerant
 - ⑤ immutable
 - ⑥ lazy evaluation
 - ⑦ Cache / co-partition
 - ⑧ in-built optimization
 - ⑨ supports ANSI SQL

Advantages of PySpark

- It is distributed processing engine that allows us to process data efficiently.
- It is 100X faster.
- It has better benefits for data ingestion pipelines.
- We can process data from Hadoop/HDFS, AWS S3 & Amazon Redshift, and real-time data using streaming and Kafka.

Module or Package

- PySpark RDD (PySpark.RDD)
- PySpark DataFrame & SQL (PySpark.SQL)
- PySpark Streaming (PySpark.Streaming)
- PySpark MLlib (PySpark.mllib.PySpark.MLlib)
- PySpark GraphFrames (GraphFrames)