

Limitations of MapReduce methodology and Hadoop's MapReduce computing engine

Because MapReduce is mainly meant for batch processing jobs, interactive jobs and models are impossible to implement.

Furthermore, jobs that are interdependent cannot be parallelized, which is not achievable with MapReduce. It struggles to tackle problems that aren't easily recombinable.

The MapReduce framework in Hadoop reads and writes data to and from disc. The data is read from and written to the disc at each stage of the processing process. Because disc seeks require time, the entire operation takes a long time. Hadoop is relatively slow for processing tiny amounts of data. Large data sets are suitable for it. Hadoop's performance for real-time processing is limited due to the batch processing engine at its core. Hadoop's Map-Reduce framework is incapable of handling real-time data. Hadoop works in batches to process data. The user must first load the file into HDFS. The user then does a map-reduce operation using the file as input.

References

- A Statistical Density-Based Analysis of Graph Clustering Algorithm Performance - <https://arxiv.org/pdf/1906.02366.pdf>
- MapReduce Algorithms for k-means Clustering https://stanford.edu/~rezab/classes/cme323/S16/projects_reports/bodoia.pdf
- K-Means clustering on MapReduce - https://web2.qatar.cmu.edu/~mhhamou/15440-f19/recitations/Project4_Handout.pdf
- Elbow Method for optimal value of k in KMeans - <https://www.geeksforgeeks.org/elbow-method-for-optimal-value-of-k-in-kmeans/>
- Hadoop – Pros and Cons - <https://www.geeksforgeeks.org/hadoop-pros-and-cons/>