

第1章Hadoop介绍

(源自:<https://biglab.site>)

(版本:Ver2.0-20230910)

视频学习

[在线学习视频](#)

Hadoop概述

简介

大数据特点4V

- 大数据的4V特征：规模性 (Volume)、高速性 (Velocity)、多样性 (Variety)、价值性 (Value)

自然界中，哪种数据类型最多：非结构化数据

发展历史

- 2002.Doug Cutting创建Nutch
- 2003.Google发布了GFS和MapReduce论文
 - Google-Bigtable中文版_1.0
 - Google-File-System中文版_1.0
 - Google-MapReduce中文版_1.0
- 2004.Doug Cutting研究出Nutch的GFS和MapReduce
- 2006.Doug Cutting加入Yahoo,Hadoop成立
- 2007.纽约时报在Ec2上转换4TB图片数据
- 2008.Cloudera成立;Facebook团队开发Hive;910个节点对1TB数据进行排序209s
- 2009.Yahoo对1TB数据排序62s
- 2011.Yahoo的Hortonworks成立
- 2012.Hortonworks的Yarn出v1
- 2013.Hortonworks完全开源
- 2016.Hadoop生态圈广泛应用

特点

HDFS特点

- A. 高可靠性
- B. 高效性
- C. 高容错性
- D. 成本低
 - Hadoop开源
- E.高扩展性
- F.基本框架是Java编写的
- G.可构建在廉价机器上

Hadoop核心

分布式文件系统HDFS

HDFS架构及简介

- HDFS文件系统主要包括一个NameNode, 一个 Secondary NameNode和多个DataNode
- 组成
 - 元数据
 - 三类信息
 - 文件和目录属性
 - 文件名、目录名、大小、时间等
 - 文件内容存储的相关信息
 - 文件块、副本数、副本所在的DataNode
 - 所有DataNode的信息
 - NameNode
 - 名称节点 (NameNode)负责管理分布式文件系统的命名空间, 保存了两个核心的数据结构, 即fsImage和edits(EditLog)
 - 文件以块形式进行存储
 - 以128MB的数据块切割文件
 - 备份副本, 默认3个
 - 修改块大小
 - 修改hdfs-site.conf配置文件
 - hadoop3.x配置dfs.blocksize=134217728
 - 参考:<https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-hdfs/hdfs-default.xml>
 - SN:Secondary NameNode
 - DataNode
 - 数据块

HDFS分布式原理

- 利用多个节点共同协作完成一项或多项具体业务功能的系统就是分布式系统
- 分布式文件系统是分布式系统的一个子集，其解决的问题就是数据存储
- HDFS作为分布式文件系统体现在三个方面
 - HDFS分布在多个集群节点上的文件系统
 - 文件存储时被分布在多个节点上

HDFS特点

优点

- 高容错性
- 适合大数据的处理
- 流式数据访问

缺点

- 不适合低延时数据访问
- 无法高效存储大量小文件
- 不支持多用户写入及任务修改文件

调度策略

- FIFO:先进先出

分布式计算框架MapReduce

- 简介
- 工作原理

集群资源管理器YARN

简介

- <https://hadoop.apache.org/docs/stable/hadoop-yarn/hadoop-yarn-site/YARN.html>
- YARN:Yet Another Resource Negotiator:另外一种资源管理协调工具

基本架构

- ResourceManager(RM)
- NodeManager(NM)
- ApplicationMaster(AM)
- ClientApplication

任务流程

Hadoop3.x与2.x

- Hadoop2.x
 - Job Tracker是Map-reduce框架的中心，他需要与集群中的机器定时通信heartbeat,需要管理哪些程序应该跑在哪些机器上，需要管理所有job失败、重启等操作

- TaskTracker是Map-Reduce集群中每台机器都有的一个部分，他做的事情主要是监视自己在机器的资源情况
- Hadoop3.x
 - YARN 替代了单独存在的JobTracker 与 TaskTracker
 - 取代(Job Tracker,TaskTracker)的是ResourceManager ApplicationMaster NodeManager三个部分

Hadoop生态系统

Nutch

- <https://nutch.apache.org/>
- Nutch 是一个开源Java实现的搜索引擎。它提供了我们运行自己的搜索引擎所需的全部工具。包括全文搜索和Web爬虫

HBase

- 针对非结构化数据的可伸缩、高可靠、高性能、分布式和面向列的动态模式数据库
- 针对BigTable的开源实现
- 名称来源于Hadoop database
- 访问接口
 - 1. Native Java API, 最常规和高效的访问方式, 适合Hadoop MapReduce Job并行批处理HBase表数据
 - 2. HBase Shell, HBase的命令行工具, 最简单的接口, 适合HBase管理使用
 - 3. Thrift Gateway, 利用Thrift序列化技术, 支持C++, PHP, Python等多种语言, 适合其他异构系统在线访问HBase表数据
 - 4. REST Gateway, 支持REST 风格的Http API访问HBase, 解除了语言限制
 - 5. Pig, 可以使用Pig Latin流式编程语言来操作HBase中的数据, 和Hive类似, 本质最终也是编译成MapReduce Job来处理HBase表数据, 适合做数据统计
 - 6. Hive, 可以使用类似SQL语言来访问HBase

Thrift

- <https://thrift.apache.org/>

thrift是一个软件框架，用来进行可扩展且跨语言服务的开发。thrift允许定义一个简单的定义文件中的数据类型和服务接口，以作为输入文件。thrift结合了功能强大的软件堆栈和代码生成引擎，以构建在 C++、Java、Python、PHP、Ruby、Erlang、Perl、Haskell、C#、Cocoa、JavaScript、Node.js、Smalltalk、and OCaml这些编程语言间无缝结合的、高效的服务。

Lucene

- <https://lucene.apache.org/>

Lucene是apache软件基金会4 jakarta项目组的一个子项目，是一个开放源代码的全文检索引擎工具包，但它不是一个完整的全文检索引擎，而是一个全文检索引擎的架构，提供了完整的查询引擎和索引引擎，部分文本分析引擎

Hive

- 是建立在Hadoop上的数据仓库基础构架。它提供的工具可存储、查询、分析存储中大规模数据。HQL可转换为复杂的MapReduce运行。相比Pig更像数据库

Pig

- 是基于Hadoop的大规模数据分析框架，Pig Latin语言也能转换成MapReduce运行，并且使用关系进行存储，减少文件输出
- Pig与Hive区别
 - https://blog.csdn.net/weixin_33721344/article/details/92049273

Sqoop

- 开源工具，主要用于Hadoop(Hive)与传统数据(MySQL、PostgreSQL等)间的数据传递

Flume

- <https://flume.apache.org/>
- 海量日志采集、聚合和传输系统
- 下载
 - <https://flume.apache.org/download.html>
 - wget <https://dlcdn.apache.org/flume/1.9.0/apache-flume-1.9.0-bin.tar.gz>
- 解压
 - tar -xvf ./apache-flume-1.9.0-bin.tar.gz
 - mv ./apache-flume-1.9.0-bin /usr/local/flume
- 使用帮助
 - <https://flume.apache.org/releases/content/1.9.0/FlumeUserGuide.html>
 - 启动
 - bin/flume-ng agent -n \$agent_name -c conf -f conf/flume-conf.properties.template

Oozie

- <https://oozie.apache.org/>
- 中文:驯象人;功能:是Hadoop调度器，可调度MapReduce,Pig,Hive,Shell,JAR任务等

ZooKeeper

- 分布式环境下的数据管理问题：统一命名、状态同步、集群管理、配置同步等

Mahout

- 中文:管象的人
- 机器学习经典算法：聚类、分类、推荐引擎等数据挖掘方法；还包含数据输入输出工具、与其他存储系统，包括 MongoDB集成等数据挖掘支持

Solr

- Solr (读作“solar”) 是Apache Lucene项目的开源企业搜索 (英语: Enterprise search) 平台

Avro

- <https://avro.apache.org/>
- 文档: <https://avro.apache.org/docs/current/gettingstartedpython.html>
- 一个数据序列化的系统

Ambari

- Apache Ambari是一种基于Web的工具, 支持Apache Hadoop集群的供应、管理和监控。Ambari 已支持大多数Hadoop组件, 包括HDFS、MapReduce、Hive、Pig、Hbase、Zookeeper、Sqoop和Hcatalog等

Apache Ambari 支持HDFS、MapReduce、Hive、Pig、Hbase、Zookeeper、Sqoop和Hcatalog 等的集中管理。也是5个顶级hadoop管理工具之一

Hadoop应用场景

- 在线旅游
- 移动数据
- 电子商务
- 能源开采
- 图像处理
- 诈骗检测
- IT安全
- 医疗保健
- 搜索引擎
- 社交平台

版本历史

- Ver1.1-20220121
 - 初始版本
- Ver1.2-20230828
 - 增加markdown版本
- Ver2.0-20230910
 - 修改hadoop版本为3.1.4重新编写

[下一章节](#)