

大数据人才培养规划教材

以解决实际问题为**学习目标**

以实战案例贯穿为**学习手段**



Hadoop

大数据开发基础

Hadoop Big Data Development

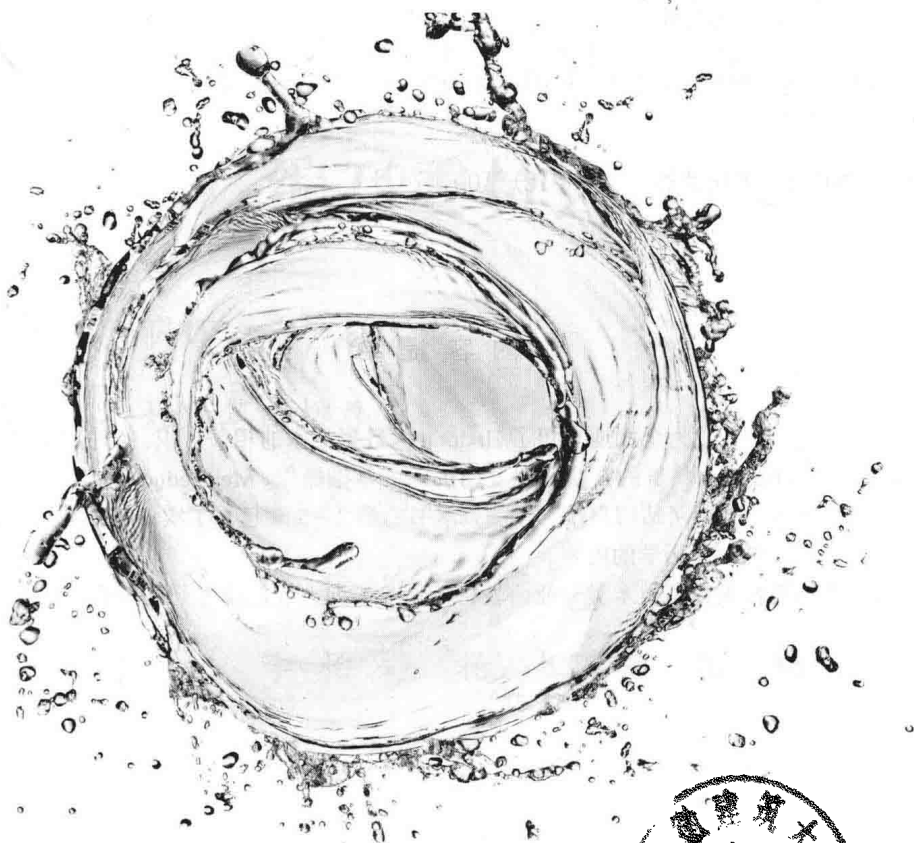
余明辉 张良均 ● 主编
高杨 陈浩 樊哲 ● 副主编



中国工信出版集团



人民邮电出版社
POSTS & TELECOM PRESS



Hadoop

大数据开发基础

Hadoop Big Data Development

余明辉 张良均 ● 主编
高杨 陈浩 樊哲 ● 副主编

人民邮电出版社
北京

图书在版编目 (C I P) 数据

Hadoop大数据开发基础 / 余明辉, 张良均主编. —
北京: 人民邮电出版社, 2018.2
大数据人才培养规划教材
ISBN 978-7-115-37066-2

I. ①H… II. ①余… ②张… III. ①数据处理软件—
教材 IV. ①TP274

中国版本图书馆CIP数据核字(2018)第001474号

内 容 提 要

本书以任务为导向, 较为全面地介绍了 Hadoop 大数据技术的相关知识。全书共 6 章, 具体内容
包括 Hadoop 介绍、Hadoop 集群的搭建及配置、Hadoop 基础操作、MapReduce 编程入门、MapReduce
进阶编程、项目案例: 电影网站用户性别预测。本书的第 2~5 章包含了实训与课后练习, 通过练习
和操作实践, 帮助读者巩固所学的内容。

本书可以作为高校大数据技术类专业的教材, 也可作为大数据技术爱好者的自学用书。

-
- ◆ 主 编 余明辉 张良均
 - 副 主 编 高 杨 陈 浩 樊 哲
 - 责任编辑 左仲海
 - 责任印制 马振武

 - ◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路 11 号
邮编 100164 电子邮件 315@ptpress.com.cn
网址 <http://www.ptpress.com.cn>
三河市潮河印业有限公司印刷

 - ◆ 开本: 787×1092 1/16
印张: 12.5 2018年2月第1版
字数: 283千字 2018年2月河北第1次印刷
-

定价: 39.80 元

读者服务热线: (010)81055256 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147 号

大数据专业系列图书

编写委员会

编委会主任：余明辉 聂哲

编委会成员（按姓氏笔画排序）：

王玉宝	王宏刚	王海	石坤泉	冯健文
刘名军	刘晓玲	刘晓勇	许昊	麦国炫
李红	李怡婷	杨坦	杨征	杨惠
肖永火	肖刚	肖芳	吴勇	邱伟绵
何小苑	何贤斌	何燕	汪作文	张玉虹
张红	张良均	张健	张凌	张敏
张澧生	陈胜	陈浩	林志章	林昆
林碧娴	欧阳国军	易琳琳	周龙	周东平
郑素铃	官金兰	赵文启	胡大威	胡坚
胡洋	钟阳晶	施兴	姜鹏辉	敖新宇
莫芳	莫济成	徐圣兵	高杨	郭信佑
黄华	黄红梅	梁同乐	焦正升	雷俊丽
詹增荣	樊哲			



序

PREFACE

随着大数据时代的到来，移动互联网和智能手机迅速普及，多种形态的移动互联网应用蓬勃发展，电子商务、云计算、互联网金融、物联网等不断渗透并重塑传统产业，大数据当之无愧地成为了新的产业革命核心。

未来 5~10 年，我国大数据产业将会是一个飞速发展时期，社会对大数据相关专业人才有着巨大的需求。目前，国内各大高校都在争相设立或准备设立大数据相关专业，以适应地方产业发展对战略性新兴产业的人才需求。

人才培养离不开教材，大数据专业是 2016 年才获批的新专业，目前还没有成套的系列教材，已有教材也存在企业案例缺失等急需解决的问题。由广州泰迪智能科技有限公司和人民邮电出版社策划，校企联合编写的这套图书，犹如大旱中的甘露，可以有效解决高校大数据相关专业教材紧缺的困境。

实践教学是在一定的理论指导下，通过引导学习者的实践活动，从而传承实践知识、形成技能、发展实践能力、提高综合素质的教学活动。目前，高校教学体系的设置有诸多限制因素，过多地偏向理论教学，课程设置与企业实际应用切合度不高，学生无法把理论转化为实践应用技能。课程内容设置方面看似繁多又各自为“政”，课程冗余、缺漏，体系不健全。本套图书的第一大特点就是注重学生的实践能力培养，根据高校实践教学中的痛点，首次提出“鱼骨教学法”的概念。以企业真实需求为导向，学生学习技能紧紧围绕企业实际应用需求，将学生需掌握的理论知识，通过企业案例的形式进行衔接，达到知行合一、以用促学的目的。

大数据专业应该以大数据技术应用为核心，紧紧围绕大数据应用闭环的流程进行教学，才能够使学生从宏观上理解大数据技术在行业中的具体应用场景及应用方法。高校现有的大数据课程集中在如何进行数据处理、建模分析、调整参数，使得模型的结果更加准确。但是，完整的大数据应用却是一个容易被忽视的部分。本套图书的第二大特点就是围绕大数据应用的整个流程，从数据采集、数据迁移、数据存储、数据

分析与挖掘，最终到数据可视化，覆盖完整的大数据应用流程，涵盖企业大数据应用中的各个环节，符合企业大数据应用真实场景。

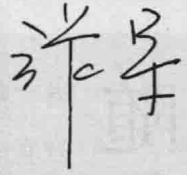
我很高兴看到这套书的出版，也希望这套书能给更多的高校师生带来教学上的便利，帮助读者尽快掌握本领，成为有用之才！

教育部长江学者特聘教授

国家杰出青年基金获得者

IEEE Fellow

华南理工大学计算机与工程学院院长



2017年12月



前言

FOREWORD

2007年9月，Apache基金会整合 Doug Cutting（Hadoop 创始人）以及其他 IT 公司（如 Facebook 等）的贡献成果，开发并正式推出了第一个 Hadoop 系统版本。Hadoop 是一个可以搭建在廉价 x86 服务器上的分布式集群系统架构，它具有可用性高、容错性高和可扩展性好等优点。由于它提供了一个开放式的平台，用户可以在完全不了解底层实现细节的情形下，开发适合自身应用的分布式程序。经过十多年的发展，目前 Hadoop 已经成长为一个全栈式的大数据技术生态圈，并在事实上成为应用最广泛、最具有代表性的大数据技术。

如何从零基础开始学习 Hadoop 大数据技术，并能够理论结合实践，运用相关知识解决一些实际的业务需求，正是本书致力解决的问题。

本书特色

本书是定位于 Hadoop 大数据技术从入门到应用的简明系统教程，主要包括 Hadoop 基本原理与架构、集群安装配置、MapReduce 编程、完整项目案例等精选内容。本书涉及的知识点简要精到，实践操作性强，能有效指导读者对 Hadoop 大数据技术的学习理解及开发应用。

本书采用了以任务为导向的教学模式，按照解决实际任务的工作流程路线，逐步展开介绍相关的理论知识点，推导生成可行的解决方案，最后落在任务实现环节。全书大部分章节紧扣任务需求展开，不堆积知识点，着重于解决问题时思路的启发与方案的实施。通过对从任务需求到实现这一完整工作流程的体验，帮助读者真正理解与消化 Hadoop 大数据技术。

本书适用对象

- 开设有大数据相关课程的学生。

目前国内不少高校将大数据技术引入教学中，在计算机、数学、自动化、电子信息、金融等专业开设了与大数据技术相关的课程，但缺乏适合课堂教学的相关教材。而本书提供了大数据相关技术的介绍、原理、实践、企业应用等，能有效指导高校学生学习大数据相关技术原理，为以后工作和学习打下良好基础。

- 大数据开发技术人员。

本书由浅及深、系统地介绍了 Hadoop 大数据开发技术，并且每一模块有对应的动手实践，对于初级开发人员有较强的指导作用。

- 关注大数据技术的各行业技术人员。

本书不仅对 Hadoop 大数据的相关技术进行了理论性的介绍及讲解，还提供了多个行业实践任务与大数据技术相结合的综合案例。各行业技术人员可以通过学习书中案例的解决思路与实现方法，尝试以新技术解决行业中的相关问题。

代码下载及问题反馈

为方便读者的实践与练习，书中提供全部实例的数据文件及源代码，读者可登录人民邮电出版社教育社区（www.ryjiaoyu.com）或“泰迪杯”全国数据挖掘挑战赛网站（www.tipdm.org/tj/1233.jhtml）下载。为方便广大教师授课需要，本书也提供了教学课件 PPT。有需要的教师可通过泰迪大数据挖掘微信公众号（TipDataMining）或者热线电话（40068-40020）进行在线咨询获取。



我们已经尽最大努力避免在文本和代码中出现错误，但是由于水平有限，编写时间仓促，书中难免出现一些疏漏和不足的地方。如果您有更多的宝贵意见，欢迎发送邮件至邮箱 13560356095@qq.com，期待能够得到您真挚的反馈。同时，本书更新内容将及时在“泰迪杯”全国数据挖掘挑战赛网站上发布，读者可以登录网站或关注泰迪大数据挖掘微信公众号（TipDataMining）查阅相关信息。

编者
2017年9月

目 录 CONTENTS

第 1 章 Hadoop 介绍	1
1.1 Hadoop 概述	1
1.1.1 Hadoop 简介	1
1.1.2 Hadoop 的发展历史	2
1.1.3 Hadoop 的特点	3
1.2 Hadoop 核心	4
1.2.1 分布式文件系统——HDFS	4
1.2.2 分布式计算框架——MapReduce	7
1.2.3 集群资源管理器——YARN	9
1.3 Hadoop 生态系统	12
1.4 Hadoop 应用场景	14
小结	15
第 2 章 Hadoop 集群的搭建及配置	16
任务 2.1 安装及配置虚拟机	17
2.1.1 创建 Linux 虚拟机	17
2.1.2 设置固定 IP	25
2.1.3 远程连接虚拟机	27
2.1.4 虚拟机在线安装软件	29
2.1.5 任务实现	32
任务 2.2 安装 Java	32
2.2.1 在 Windows 下安装 Java	33
2.2.2 在 Linux 下安装 Java	35
2.2.3 任务实现	36
任务 2.3 搭建 Hadoop 完全 分布式集群	36
2.3.1 修改配置文件	36
2.3.2 克隆虚拟机	41
2.3.3 配置 SSH 免密码登录	43
2.3.4 配置时间同步服务	44
2.3.5 启动关闭集群	46
2.3.6 监控集群	47
小结	50
实训	50
实训 1 为 Hadoop 集群增加一个节点	50
实训 2 编写 Shell 脚本同步集群时间	51
课后练习	51
第 3 章 Hadoop 基础操作	53
任务 3.1 查看 Hadoop 集群的 基本信息	54
3.1.1 查询集群的存储系统信息	55
3.1.2 查询集群的计算资源信息	58
任务 3.2 上传文件到 HDFS 目录	59
3.2.1 了解 HDFS 文件系统	59
3.2.2 掌握 HDFS 的基本操作	62
3.2.3 任务实现	65
任务 3.3 运行首个 MapReduce 任务	67
3.3.1 了解 Hadoop 官方的示例程序包	67
3.3.2 提交 MapReduce 任务给集群运行	68
任务 3.4 管理多个 MapReduce 任务	71
3.4.1 查询 MapReduce 任务	72
3.4.2 中断 MapReduce 任务	74
小结	76
实训	77
实训 1 统计文件中所有单词的平均长度	77
实训 2 查询与中断 MapReduce 任务	77
课后练习	78
第 4 章 MapReduce 编程入门	80
任务 4.1 使用 Eclipse 创建 MapReduce 工程	81
4.1.1 下载与安装 Eclipse	81

Hadoop 大数据开发基础

4.1.2 配置 MapReduce 环境	82	5.3.1 自定义键值类型	124
4.1.3 新建 MapReduce 工程	84	5.3.2 初步探索 Combiner	128
任务 4.2 通过源码初识 MapReduce 编程	86	5.3.3 浅析 Partitioner	130
4.2.1 通俗理解 MapReduce 原理	86	5.3.4 自定义计数器	132
4.2.2 了解 MR 实现词频统计的执行流程	88	5.3.5 任务实现	134
4.2.3 读懂官方提供的 WordCount 源码	89	任务 5.4 Eclipse 提交日志文件统计程序	137
任务 4.3 编程实现按日期统计访问次数	94	5.4.1 传递参数	137
4.3.1 分析思路与处理逻辑	94	5.4.2 Hadoop 辅助类 ToolRunner	139
4.3.2 编写核心模块代码	95	5.4.3 Eclipse 自动打包并提交任务	140
4.3.3 任务实现	97	小结	144
任务 4.4 编程实现按访问次数排序	99	实训	144
4.4.1 分析思路与处理逻辑	99	实训 1 统计全球每年的最高气温和最低气温	144
4.4.2 编写核心模块代码	100	实训 2 筛选气温在 15~25℃ 之间的数据	145
4.4.3 任务实现	102	课后练习	146
小结	104	第 6 章 项目案例：电影网站用户性别预测	151
实训	104	任务 6.1 认识 KNN 算法	152
实训 1 获取成绩表的最高分记录	104	6.1.1 KNN 算法简介	152
实训 2 对两个文件中的数据进行合并与去重	105	6.1.2 KNN 算法原理及流程	152
课后练习	107	任务 6.2 数据预处理	154
第 5 章 MapReduce 进阶编程	110	6.2.1 获取数据	154
任务 5.1 筛选日志文件并生成序列化文件	111	6.2.2 数据变换	155
5.1.1 MapReduce 输入格式	111	6.2.3 数据清洗	160
5.1.2 MapReduce 输出格式	113	6.2.4 划分数据集	163
5.1.3 任务实现	113	任务 6.3 实现用户性别分类	167
任务 5.2 Hadoop Java API 读取序列化日志文件	115	6.3.1 实现思路	167
5.2.1 FileSystem API 管理文件夹	115	6.3.2 代码实现	169
5.2.2 FileSystem API 操作文件	119	任务 6.4 评价分类结果的准确性	179
5.2.3 FileSystem API 读写数据	121	6.4.1 评价思路	179
5.2.4 任务实现	123	6.4.2 实现分类评价	180
任务 5.3 优化日志文件统计程序	124	6.4.3 寻找最优 K 值	184
		小结	188
		参考文献	189



第 1 章 Hadoop 介绍



学习目标

- (1) 认识 Hadoop。
- (2) 了解 Hadoop 的核心组件。
- (3) 了解 Hadoop 的生态系统。
- (4) 了解 Hadoop 的应用场景。



任务背景

随着时代的发展，“大数据”已经成为一个耳熟能详的词汇。与此同时，针对大数据处理的新技术也在不断的开发和运用中，逐渐成为数据处理挖掘行业广泛使用的主流技术之一。本章就来简要介绍一款非常有代表性的大数据处理框架——Hadoop。

在大数据时代，Hadoop 作为处理大数据的分布式存储和计算框架，得到了国内外大、中、小型企业的广泛应用，学习 Hadoop 技术是从从事大数据行业工作必不可少的一步。本章将从以下几个方面了解 Hadoop 的框架理论。首先了解 Hadoop 的发展历史与特点，然后进一步讲解 Hadoop 的两大核心——HDFS 和 MapReduce，以及用于资源与任务调度的 YARN 框架。接下来对 Hadoop 生态系统中的组件进行简单的了解，包括组件的特点和应用。最后简要介绍了 Hadoop 的使用情况和应用场景。

1.1 Hadoop 概述

1.1.1 Hadoop 简介

随着移动设备的广泛使用和互联网的快速发展，数据的增量和存量快速增加，硬件发展跟不上数据发展，单机设备很多时候已经无法处理庞大的甚至 TB、PB 级别的数据。如果一头牛拉不动货物，显然找几头牛一起拉会比培育一头更强壮的牛容易。同理，对于单机无法解决的问题，综合利用多个普通机器要比打造一台超级计算机更加可行，这就是 Hadoop 的设计思想。

Hadoop 是一个由 Apache 基金会所开发的可靠的、可扩展的用于分布式计算的分布式系统基础架构和开发开源软件。Apache Hadoop 软件库是一个框架，允许使用简单的编程模型在计算机集群中对大规模数据集进行分布式处理。它的目的是从单一的服务器扩展到成千上万的机器，将集群部署在多台机器中，每台机器提供本地计算和存储，并且将存储的数据备份在多个节点，由此提升集群的可用性，而不是通过硬件提升。当一台机器宕机时，其他节点依然可以提供备份数据和计算服务。

Hadoop 大数据开发基础

Hadoop 框架最核心的设计是 HDFS (Hadoop Distributed File System) 和 MapReduce。HDFS 是可扩展、高容错、高性能的分布式文件系统，负责数据的分布式存储和备份，文件写入后只能读取，不能修改。MapReduce 是分布式计算框架，包含 Map (映射) 和 Reduce (规约) 两个过程。

1.1.2 Hadoop 的发展历史

Hadoop 是由 Apache Lucence 创始人 Doug Cutting 创建的，Lucence 是一个应用广泛的文本搜索系统库。Hadoop 起源于开源的网络搜索引擎 Apache Nutch，它本身也是 Lucence 项目的一部分。Hadoop 的发展历程如图 1-1 所示。

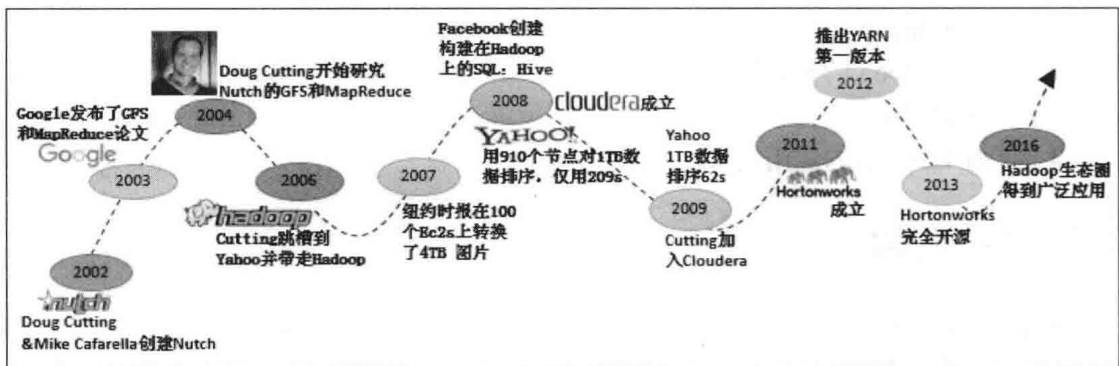


图 1-1 Hadoop 发展历程

Nutch 项目开始于 2002 年，当时互联网第一轮泡沫刚刚结束，Doug Cutting 与好友 Mike Cafarella 认为网络搜索引擎由一个互联网公司垄断十分可怕，他们将掌握信息的入口，因此决定自己开发一个可以代替当时主流搜索产品的开源搜索引擎，该项目命名为 Nutch。Nutch 致力于提供开源搜索引擎所需的全部工具集。但后来，两位开发者发现这一架构的灵活性不够，只能支持几亿数据的抓取、索引和搜索，不足以解决数十亿网页的搜索问题。

2003 年，Google 发表的论文 “The Google File System” 为此提供了帮助，文中描述的是谷歌产品架构，该架构被称为 “谷歌分布式文件系统”，简称 GFS。Nutch 的开发者们发现 GFS 架构能够满足网页抓取和搜索过程中生成的超大文件存储需求，特别关键的是，GFS 能够节省系统管理所使用的大量时间。于是在 2004 年，开发者们借鉴谷歌新技术开始进行开源版本的实现，即 Nutch 分布式文件系统 (NDFS)。不同的是，Google 用的是 C++ 语言，而 Nutch 使用 Java 语言。

2004 年，谷歌又发表了论文 “MapReduce: Simplified Data Processing on Large Clusters”，向全世界介绍他们的 MapReduce 框架。Nutch 的开发者们发现 Google MapReduce 所解决的大规模搜索引擎数据处理问题，又解决了他们当时同样面临并急需解决的问题。Nutch 的开发者们基于 Google 发布的 MapReduce 报告，模仿了 Google MapReduce 框架的设计思路，用 Java 设计并实现了一套新的 MapReduce 并行处理软件系统，在 Nutch 上开发了一个可工作的 MapReduce 应用。

2005 年年初，Nutch 的开发人员在 Nutch 上实现了一个 MapReduce 算法，半年左右的时间，Nutch 的所有主要算法均完成移植，用 MapReduce 和 NDFS 来运行。

2006 年, Doug Cutting 在经过一系列周密考虑和详细总结后, 决定加入优秀的公司进一步完善 Nutch 的性能。IBM 却对他的早期项目 Lucence 更感兴趣, 而雅虎则看中 Nutch 底层 NDFS/MapReduce。在 2006 年 2 月, 开发人员将 NDFS 和 MapReduce 移出了 Nutch, 形成 Lucence 的子项目, 命名为 Hadoop, 这个名字来源于 Doug Cutting 儿子的一只玩具象。随后, Doug Cutting 几经周折加入了 Yahoo 公司, 雅虎为此组织了一个专门的团队和资源, 致力将 Hadoop 发展成为能够处理海量 Web 数据的分布式系统。

当加入 Yahoo 以后, Hadoop 项目逐渐发展并迅速成熟起来。首先是集群规模, 从最开始几十台机器的规模发展到能支持上千个节点的机器, 中间做了很多工程性质的工作, 然后是除搜索以外的业务, Yahoo 逐步将自己广告系统的与数据挖掘相关工作也迁移到了 Hadoop 上, 进一步促进了 Hadoop 系统的成熟与发展。

2007 年, 纽约时报在 100 个亚马逊的虚拟机服务器上使用 Hadoop 转换了 4TB 的图片数据, 更加深了人们对 Hadoop 的印象。

2008 年, 一位 Google 的工程师发现要把当时的 Hadoop 放到任意一个集群中去运行是一件很困难的事情, 所以就拉上了几个好朋友成立了一个专门商业化 Hadoop 的公司 Cloudera。同年, Facebook 团队发现大多数分析人员编写 MapReduce 程序时有难度, 而对 SQL 很熟悉, 所以他们就在 Hadoop 之上开发了一个叫作 Hive 的数据仓库工具, 专门把 SQL 转换为 Hadoop 的 MapReduce 程序。

2008 年 1 月, Hadoop 已经成为 Apache 的顶级项目。

2008 年 4 月, Hadoop 打破世界纪录, 成为最快的 TB 级数据排序系统。在一个 910 节点的集群, Hadoop 在 209s 内完成了对 1TB 数据的排序, 击败了前一年的 297s。到 2009 年 5 月, 有报道称雅虎有一个团队使用 Hadoop 对 1TB 数据进行排序只花了 62s。

2011 年, Yahoo 将 Hadoop 团队独立出来, 成立了一个子公司 Hortonworks, 专门提供 Hadoop 相关的服务。

2012 年, Hortonworks 在 Hadoop 发展上推出了与原框架有很大不同的 YARN 框架的第一版本, 从此对 Hadoop 的研究又迈进一个新的层面。

2013 年, 大型 IT 公司, 如 EMC、Microsoft、Intel、Teradata、Cisco 都明显增加了 Hadoop 方面的投入, Hortonworks 宣传要 100% 开源软件, Hadoop 2.0 转型基本上无可阻挡。

2014 年, Hadoop 2.x 的更新速度非常快, 从 2.3.0 到 2.6.0, 极大地完善了 YARN 框架和整个集群的功能。很多 Hadoop 的研发公司如 Cloudera、Hortonworks 都与其他企业合作共同开发 Hadoop 新功能。

2015 年, 在 Hadoop 创新阶段走向错误道路的供应商逐步退出, 与其他 Hadoop 版本发行企业整合, 数据的实时访问处理是一个关注的重点。

2016 年, Hadoop 及其生态圈 (包括 Spark 等) 在各行各业落地并且得到广泛的应用, YARN 将持续发展以支持更多应用。

1.1.3 Hadoop 的特点

Hadoop 是一个能够让用户轻松架构和使用的分布式计算平台。用户可以轻松地在 Hadoop 上开发和运行处理海量数据的应用程序。其优点主要有以下几个。

(1) 高可靠性

数据存储有多个备份, 集群设置在不同机器上, 可以防止一个节点宕机造成集群损坏。

Hadoop 大数据开发基础

当数据处理请求失败后，Hadoop 会自动重新部署计算任务。Hadoop 框架中有备份机制和检验模式，Hadoop 会对出现问题的部分进行修复，也可以通过设置快照的方式在集群出现问题时回到之前的一个时间点。

(2) 高扩展性

Hadoop 是在可用的计算机集群间分配数据并完成计算任务的。为集群添加新的节点并不复杂，所以集群可以很容易进行节点的扩展，扩大集群。

(3) 高效性

Hadoop 能够在节点之间动态地移动数据，在数据所在节点进行并发处理，并保证各个节点的动态平衡，因此处理速度非常快。

(4) 高容错性

Hadoop 的分布式文件系统 HDFS 在存储文件时会在多个节点或多台机器上存储文件的备份副本，当读取该文档出错或者某一台机器宕机了，系统会调用其他节点上的备份文件，保证程序顺利运行。如果启动的任务失败，Hadoop 会重新运行该任务或启用其他任务来完成这个任务没有完成的部分。

(5) 低成本

Hadoop 是开源的，即不需要支付任何费用即可下载并安装使用，节省了软件购买的成本。

(6) 可构建在廉价机器上

Hadoop 不要求机器的配置达到极高的水准，大部分普通商用服务器就可以满足要求，它通过提供多个副本和容错机制来提高集群的可靠性。

(7) Hadoop 基本框架用 Java 语言编写

Hadoop 带有用 Java 语言编写的框架，因此运行在 Linux 生产平台上是非常理想的，Hadoop 上的应用程序也可以使用其他语言编写，比如 C++。

1.2 Hadoop 核心

1.2.1 分布式文件系统——HDFS

1. HDFS 架构及简介

HDFS 是以分布式进行存储的文件系统，主要负责集群数据的存储与读取。HDFS 是一个主/从 (Master/Slave) 体系结构的分布式文件系统，从某个角度看，它就和传统的文件系统一样。HDFS 支持传统的层次型文件组织结构，用户或者应用程序可以创建目录，然后将文件保存在这些目录里。文件系统名字空间的层次结构和大多数现有的文件系统类似，可以通过文件路径对文件执行创建、读取、更新和删除操作。但是由于分布式存储的性质，它又和传统的文件系统有明显的区别。它的基本架构如图 1-2 所示。

HDFS 文件系统主要包括一个 NameNode、一个 Secondary NameNode 和多个 DataNode。

(1) 元数据

元数据不是具体的文件内容，它有三类重要信息：第一类是文件和目录自身的属性信息，例如文件名、目录名、父目录信息、文件大小、创建时间、修改时间等；第二类记录文件内容存储的相关信息，例如文件分块情况、副本个数、每个副本所在的 DataNode 信息

等；第三类用来记录 HDFS 中所有 DataNode 的信息，用于 DataNode 管理。

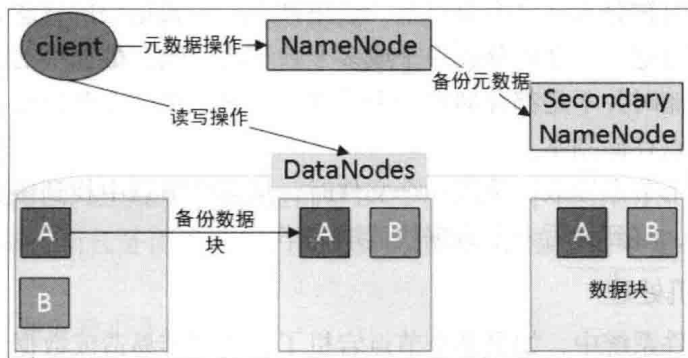


图 1-2 HDFS 架构图

(2) NameNode

NameNode 用于存储元数据以及处理客户端发出的请求。在 NameNode 中存放元信息的文件是 `fsimage` 文件。在系统运行期间，所有对元数据的操作都保存在内存中，并被持久化到另一个文件 `edits` 中。当 NameNode 启动的时候，`fsimage` 会被加载到内存，然后对内存里的数据执行 `edits` 所记录的操作，以确保内存所保留的数据处于最新的状态。

(3) Secondary NameNode

Secondary NameNode 用于备份 NameNode 的数据，周期性将 `edits` 文件合并到 `fsimage` 文件并在本地备份，将新的 `fsimage` 文件存储到 NameNode，取代原来的 `fsimage`，删除 `edits` 文件，创建一个新的 `edits` 继续存储文件修改状态。

(4) DataNode

DataNode 是真正存储数据的地方。在 DataNode 中，文件以数据块的形式进行存储。当文件传到 HDFS 端时以 128MB 的数据块将文件进行切割，将每个数据块存到不同的或相同的 DataNode 并且备份副本，一般默认 3 个，NameNode 会负责记录文件的分块信息，确保在读取该文件时可以找到并整合所有块。

(5) 数据块

文件在上传到 HDFS 时根据系统默认文件块大小把文件分成一个个数据块，Hadoop 2.x 默认 128MB 为一个数据块，比如存储大小为 129MB 的文件，则被分为两个块来存储。数据块会被存储到各个节点，每个数据块都会备份副本。

2. HDFS 分布式原理

什么是分布式系统？分布式系统会划分成多个子系统或模块，各自运行在不同的机器上，子系统或模块之间通过网络通信进行协作，实现最终的整体功能。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。换句话说，它是横跨在多台计算机上的存储系统。存储在分布式文件系统上的数据自动分布在不同的节点上。

HDFS 作为一个分布式文件系统，主要体现在以下三个方面。

(1) HDFS 并不是一个单机文件系统，它是分布在多个集群节点上的文件系统。节点之间通过网络通信进行协作，提供多个节点的文件信息，让每个用户都可以看到文件系统

Hadoop 大数据开发基础

的文件，让多机器上的多用户分享文件和存储空间。

(2) 文件存储时被分布在多个节点上。这里涉及一个数据块的概念，数据存储不是按一个文件存储，而是把一个文件分成一个或多个数据块存储，数据块的概念在前面已经描述过。数据块在存储时并不是都存储在一个节点上，而是被分布存储在各个节点中，并且数据块会在其他节点存储副本。

(3) 数据从多个节点读取。读取一个文件时，从多个节点中找到该文件的数据块，分布读取所有数据块，直到最后一个数据块读取完毕。

3. HDFS 宕机处理

数据存储的文件系统中，如果某个节点宕机了，就很容易造成数据流失，HDFS 针对这个问题提供了有效的保护措施。

(1) 冗余备份

在数据存储的过程中，对每个数据块都进行了副本备份，副本个数可以自行设置。

(2) 副本存放

仅仅对数据进行冗余备份还不够，假设所有的备份都在一个节点上，那么该节点宕机后，数据一样会丢失，因此 HDFS 要有一个更有效的副本存放策略。目前使用的策略是，以 `dfs.replication=3` 为例，在同一机器的两个节点上各备份一个副本，然后在另一台机器的某个节点上再放一个副本。前者防止该机器的某个节点宕机，后者防止某个机器宕机。

(3) 宕机处理

① 当一切运行正常时，DataNode 会周期性发送心跳信息给 NameNode（默认是每 3s 一次）。如果 NameNode 在预定的时间内没有收到心跳信息（默认是 10min），它会认为 DataNode 出问题了，把它从集群中移除。对于 HDFS 来说，丢失一个 DataNode 意味着丢失了存储在它的硬盘上的数据块的副本。HDFS 会检测到存储在该硬盘的数据块的副本数量低于要求，且主动对副本数量不符合要求的数据块创建需要的副本，以达到满副本状态。DataNode 可能因为多种原因脱离集群，如硬件故障、主板故障、电源老化和网络故障等。

② 当 HDFS 读取某个数据块时，如果正好该节点宕机了，客户端就会到存储该数据块的其他节点上读取，除非其他节点损坏或者该数据块在存储时损坏，否则依然可以得到该数据块的信息。HDFS 也会检测到该数据块副本个数不符合要求而重新补全副本。

③ 当 HDFS 存储数据时，如果要存放数据的节点宕机，HDFS 会再分配一个节点给数据块，然后备份宕机节点的数据。

4. HDFS 特点

(1) 优点

① 高容错性。

HDFS 上传的数据自动保存多个副本，通过增加副本的数量来增加它的容错性。如果某一个副本丢失，HDFS 机制会复制其他机器上的副本，而我们不必关注它的实现。

② 适合大数据的处理。

HDFS 能够处理 GB、TB 甚至 PB 级别的数据，规模达百万，数量非常大。

③ 流式数据访问。

HDFS 以流式数据访问模式来存储超大文件，“一次写入，多次读取”。文件一旦写入，不能修改，只能增加。这样可以保证数据的一致性。

(2) 缺点

① 不适合低延迟数据访问。

如果要处理一些用户要求时间比较短的低延迟应用请求，则 HDFS 不适合。HDFS 是为了处理大型数据集分析任务而设计的，目的是为达到高的数据吞吐量，这就可能要求以高延迟作为代价。

② 无法高效存储大量小文件。

因为 NameNode 把文件系统的元数据放置在内存中，所以文件系统所能容纳的文件数目是由 NameNode 的内存大小来决定的，即每存入一个文件都会在 NameNode 中写入文件信息。如果写入太多小文件，NameNode 内存会被占满而无法写入文件信息。而与多个小文件大小相同的单一文件只会写入一次文件信息到内存中，所以更适合大文件存储。

③ 不支持多用户写入及任意修改文件。

在 HDFS 的一个文件中只有一个写入者，而且写操作只能在文件末尾完成，即只能执行追加操作。目前 HDFS 还不支持多个用户对同一文件的写操作，以及在文件任意位置进行修改。

1.2.2 分布式计算框架——MapReduce

1. MapReduce 简介

MapReduce 是 Hadoop 的核心计算框架，是用于大规模数据集（大于 1TB）并行运算的编程模型，主要包括 Map（映射）和 Reduce（规约）两部分。当启动一个 MapReduce 任务时，Map 端会读取 HDFS 上的数据，将数据映射成所需要的键值对类型并传到 Reduce 端。Reduce 端接收 Map 端传过来的键值对类型的数据，根据不同键进行分组，对每一组键相同的数据进行处理，得到新的键值对并输出到 HDFS，这就是 MapReduce 的核心思想。

2. MapReduce 工作原理

MapReduce 作业执行流程如图 1-3 所示。

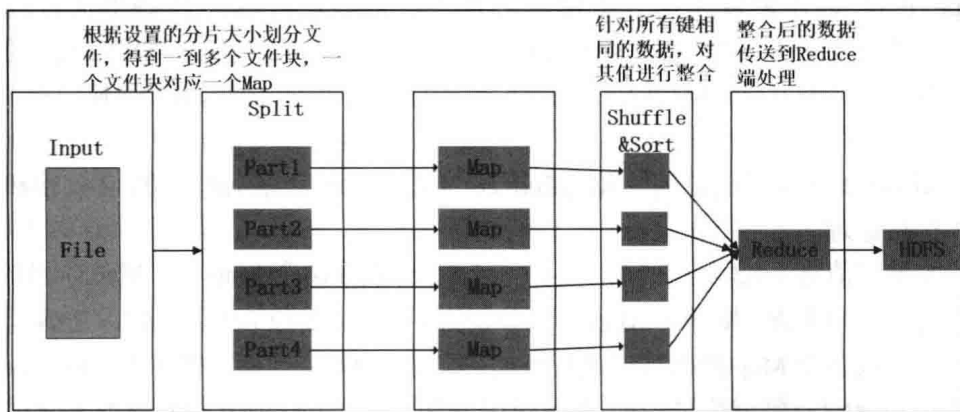


图 1-3 MapReduce 执行流程图