

L1: Hadoop 基础与操作



李丰

feng.li@cufe.edu.cn

中央财经大学统计与数学学院

本节知识要点

- 1 基本信息
- 2 分布式存储与计算的意义
- 3 Hadoop 概述
- 4 安装 Hadoop
- 5 上机实践

基本信息

- 讲师: 李丰 <feng.li@cufe.edu.cn>
- 参考书目
 - 《大数据分布式计算与案例》李丰 著 (2016) 中国人民大学出版社
- 讲义下载
<http://feng.li/distcomp/>
- 分布式计算案例
<https://github.com/feng-li/Distributed-Statistical-Computing/>
- 其他参考书目
 - Holmes, Alex. Hadoop in practice. Manning Publications Co., 2012.
 - White, Tom. Hadoop: The definitive guide, Third Edition. "O'Reilly Media, Inc.", 2012.
 - 陆嘉恒. Hadoop 实战. 机械工业出版社, 2012.

- 大数据带来的两个基础挑战（价值）：
 - 如何灵活地操作海量数据？
 - 如何高效地从海量数据中获取价值？
- 分布式系统（包括分布式存储系统和分布式计算系统）为这两个问题的解决提供了桥梁
- 目前广泛使用的有 Hadoop、Spark 等

我们身边的分布式计算 I

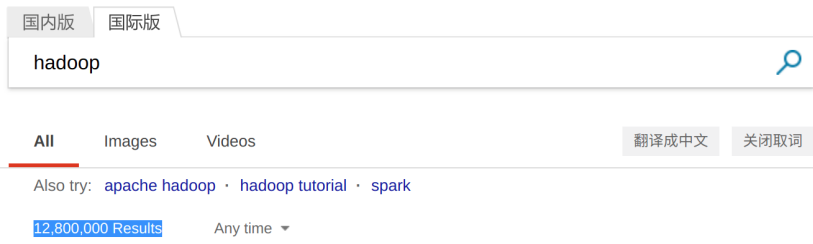


Figure: 从海量互联网数据中响应用户搜索

我们身边的分布式计算 II



海量选股公式

2

搜索文档

为您推荐: [选股公式经典](#) [选股公式集合](#) [精准选股公式](#) [简单选股公式](#) [条件选股公式](#)

 [海量选股公式1000个](#)

质量 4.5 分

0001){[选股公式](#)} ZLC:= EMA((WINNER(CLOSE) * 70),3); SHC:= EMA(((WINNER((CLOSE * 1.1)) - WINNER
((CLOSE * 0.9))) * 80),3); 主力控盘:...

2013-07-24 | 共108页 | 251次下载 | 1下载券 | 贡献者: .ysfdgd

马上下载

Figure: 海量选股

我们身边的分布式计算 III



240x180 | 170x170

把图片拖拽到此区域

本地上传

粘贴网址

相似图片



Figure: 相似图片识别 → 人脸识别 → 天网犯罪分子监控

Hadoop 简史

- Hadoop 项目由 Doug Cutting 等于 2003 年基于谷歌分布式文件系统的论文的开源实现。
- 谁在用 Hadoop? ——几乎所有的数据科学相关前沿企业
- 为什么是 Hadoop?——廉价、高效、易用、可扩展

Hadoop 概述 I

- Hadoop 提供了分布式存储和计算能力。
- Hadoop 的 Master-Worker 架构包含了 Hadoop 分布式存储系统 (**HDFS**) 和分布式计算框架 (**MapReduce**)

Hadoop 概述 II

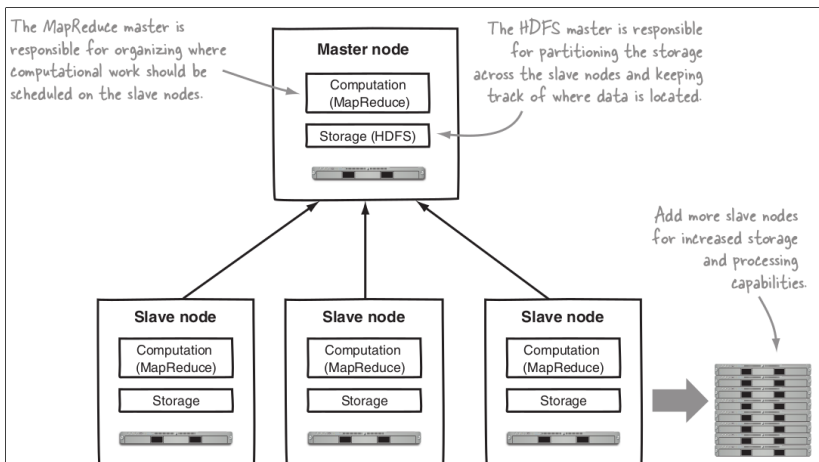


Figure: Hadoop 架构

Hadoop 核心组件：HDFS I

- **HDFS** 是 Hadoop 的存储组件
- HDFS 的文件系统操作包括：读取文件、创建目录、移动文件、删除数据、列出目录等。绝大多数命令与 Linux 文件系统操作类似。可以输入：

```
$ hadoop fs -help
```

命令获取所有命令的详细帮助文件。

- HDFS 又包括两个逻辑组件 NameNode 和 DataNode.
- 文件在分布式系统中的存储是有冗余的，以避免软件或者硬件错误造成数据的毁坏。
- HDFS 不是简单的数据分块存储，允许存储与读写海量数据。
- HDFS 不擅长对许多小文件随机访问。

Hadoop 核心组件: HDFS II

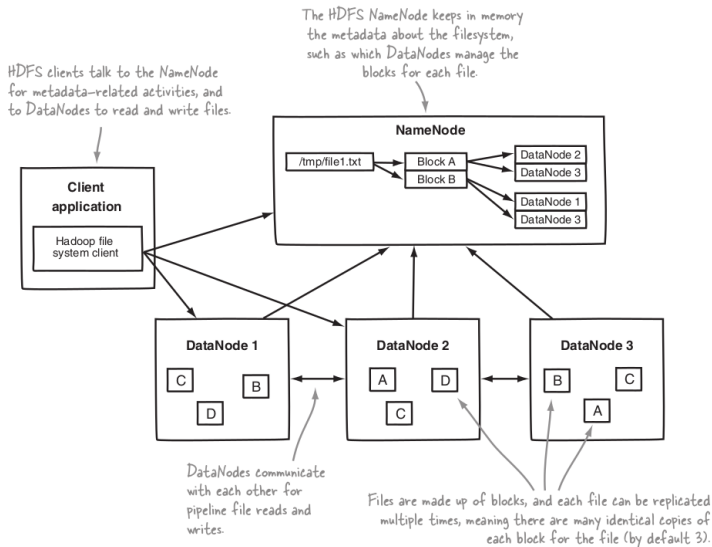


Figure: 文件在 HDFS 存储的结构

Hadoop 核心组件：MapReduce I

- **MapReduce** 是一个分布式计算模式。
- 通过 MapReduce, 大量的数据处理可以被分布式计算。
- Hadoop 提供了标准的 MapReduce 接口。

Hadoop 核心组件: MapReduce II

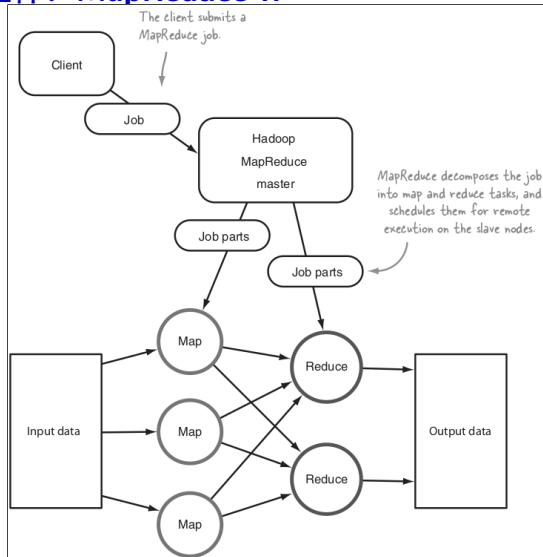


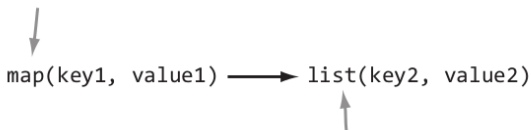
Figure: MapReduce 原理

Hadoop 解放了数据科学家（程序员）的双手 I

- 传统的并行计算需要考虑通信、负载、存储、任务切割等多个专业的计算机领域，即便资深的程序员也不一定能胜任。
- 有了 Hadoop，任何一个数据分析问题只需要定义为一个 Mapper 函数和一个 Reducer 函数。
- Mapper 函数的输出变为 Reducer 的输入，并可多个嵌套。
- Hadoop 的 shuffle 和 sort 机制完美地结合 MapReduce。

Hadoop 解放了数据科学家（程序员）的双手 II

The map function takes as input a key/value pair, which represents a logical record from the input data source. In the case of a file, this could be a line, or if the input source is a table in a database, it could be a row.



The map function produces zero or more output key/value pairs for that one input pair. For example, if the map function is a filtering map function, it may only produce output if a certain condition is met. Or it could be performing a demultiplexing operation, where a single input key/value yields multiple key/value output pairs.

Figure: Mapper 函数

Hadoop 解放了数据科学家（程序员）的双手 III

The shuffle and sort phases are responsible for two primary activities: determining the reducer that should receive the map output key/value pair (called partitioning); and ensuring that, for a given reducer, all its input keys are sorted.

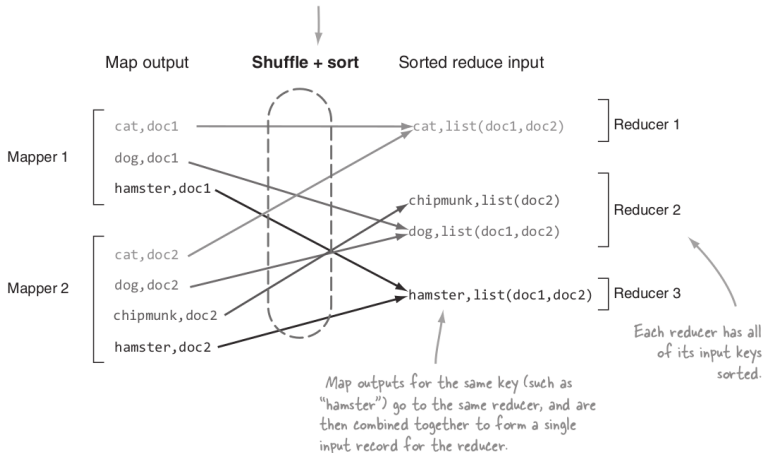


Figure: MapReduce shuffle and sort.

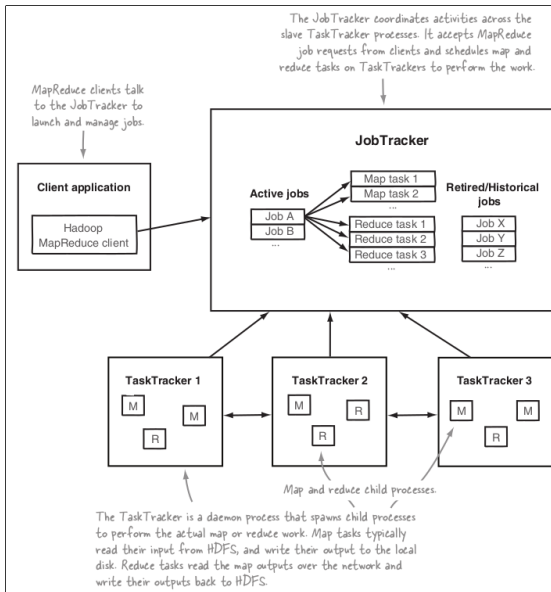


Figure: MapReduce 拓扑架构

运行一个词频统计的示例 I

- Mapper: Linux 自带程序 `/usr/bin/cat`
concatenate files and print on the standard output.
- Reducer: Linux 自带程序 `/bin/wc`
print newline, word, and byte counts for each file, and a total line if more than one file is specified.
- 测试文件: 存储在 HDFS 的文件
- 语法

```
$ $HADOOP_HOME/bin/hadoop jar \  
  $HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-2.5.  
  -input /stocks.txt \  
  -output wcoutfile \  
  -mapper "/bin/cat" \  
  -reducer "/usr/bin/wc" \  
  -
```

Hadoop 运行模式

- 单机模式
- 伪分布模式
- 全分布模式

安装伪分布式 Hadoop

- 依赖条件
 - Linux OS
 - JDK
 - 配置 SSH
 - 安装 Open SSH Server
 - 秘钥
- Hadoop 的配置文件位置 `$HADOOP_HOME/hadoop/ect/hadoop/.`
- Hadoop 在线文档 <http://hadoop.apache.org/docs/current/>.

- 实现一个简单的 MapReduce
- 命令行与 Hadoop 交互：对 HDFS 上传下载一个文件以及基本的分布式文件操作。
- 遇到问题怎么办？——查看在线帮助文档。
- 熟悉 Hadoop 配置文件以及交互界面 (configuration files, logs, http interfaces...)