

Spark

No.

Date:

1. Spark是什么.

快速且通用的集群计算框架.

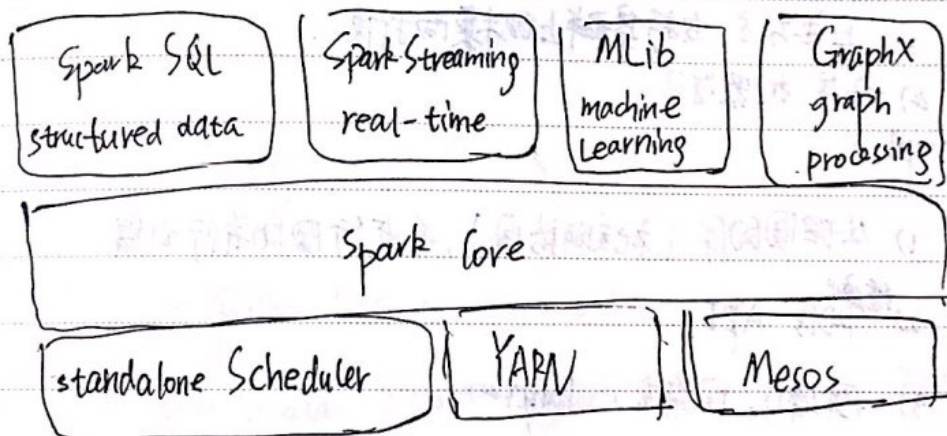
快速: { 1. 扩充 M/R 计算模型.
2. 基于内存.

通用 { 1. 容纳其他分布式系统拥有的功能.
2. 批处理, 迭代计算, 交互查询, 和流处理.

优点: 降低成本.

2. Spark生态

Spark 组件:



① Spark Core

1) 包含 Spark 基本功能, 包含任务调度, 内存管理, 容错机制.

2) 内部定义 RDDs (弹性分布式数据集).

3) APIs \rightarrow RDDs.

4) 为其他组件, 提供底层服务.



② Spark SQL

- 1) 类似 Hive, MySQL
- 2) 场景: 报表统计

③ Spark Stream

- 1) 实时数据流处理组件, 类似 Storm.
- 2) API → 实时数据流.
- 3) 场景: 从 Kafka 接收数据做实时统计.

④ MLlib

- 1) 包含通用机器学习的包.
- 2) 分类, 聚类, 回归等, 模型评估, 和特征导入.
- 3) 上述方法, 支持集群上的横向扩展.
- 4) 场景: 机器学习.

⑤ GraphX

- 1) 处理图的库 (社交网络图), 并进行图的并行计算.
- 2) ~~提供~~ PDD API
- 3) 图操作, 图算法, PageRank.
- 4) 场景: 图计算.

⑥ Cluster Managers

- 1) 自带集群管理: 单独调度器;
- 2) 常见集群管理: YARN, Mesos.



④ 集成优点

- 1) 整体优化.
- 2) 便于部署
- 3) 共享组件.

3. Spark 与 Hadoop

比较: Hadoop: { 1. 离线处理.
2. 时效性不高.

Spark { 1. 实时性
2. 机器学习.

比较: ①. 各自生态.

②. 不兼容 HDFS, 借助

③.

二. Spark 安装

①. Scala 编写 \rightarrow Java7+

②. Python API \rightarrow Python 2.6+ 或 Python 3.4+

③. Spark 1.6.2 \leftrightarrow Scala 2.10, Spark 2.0.0 - Scala 2.11

下载:

①. spark.apache.org/downloads.html

②. 不需要 Hadoop, 如果有, 可下载 Hadoop-Spark 版本.

③ 解压.



目录.

- ① bin \rightarrow Spark Shell
- ② core. streaming, python \rightarrow 源码.
- ③ examples \leftrightarrow 单机 Spark job.

Spark Shell

- ① 处理分布的集群上的数据
- ② 把数据加载到内存中, 秒级
- ③ 迭代式计算, 实时查询, 分析在 shells 中完成.
- ④ Python Shells, Scala Shells.

Python Shell:

bin / pyspark.

Scala Shell:

bin / ~~spark~~ spark-shell

```
var lines = sc.textFile (".../testfile/helloSpark")
```

```
lines.count()
```

```
lines.first()
```

修改日志级别

```
conf / log4j , cp log4j.properties.
```

```
vi log4j.properties. log4j.rootCategory=WARN,console
```



开发 Spark 环境搭建

①. idea.

②. 插件安装: Scalar, 重启 idea.

③. 新建 → Scalar → SBT 打包, scala 1.8.0-40, 0.13.8, 2.10.5.
Spark 1.6.2

④. 常见问题:

脚本兼容

ssh 配置

①. ssh-keygen.

②. cat rsa-pub-key → authorized_keys.

③. chmod 600 authorized_keys.

Spark 程序:

①. scalar class → word count → Object.

②. 代码: 1) conf =, 2) sc = 3) input =, 4) lines, count = 4) output

③. project Structure

Artifacts → jar → From modules with dependencies

Module (2 个), → Main class → copy entry class → ok

Build Artifacts → build

④. 启动集群: 1) 启动 Master: ./bin/start-master.sh
2) 启动 worker: ./bin/spark-class
3) 提交作业: ./bin/spark-submit



Spark web

localhost:8080.

Spark master at spark11 - $\leq \leq$ 7.077./bin/spark-class org.apache.spark.deploy.worker.Worker $\leq \leq$ 7.077./bin/spark-submit --master $\leq \leq$ 7.077 --class WordCount --class

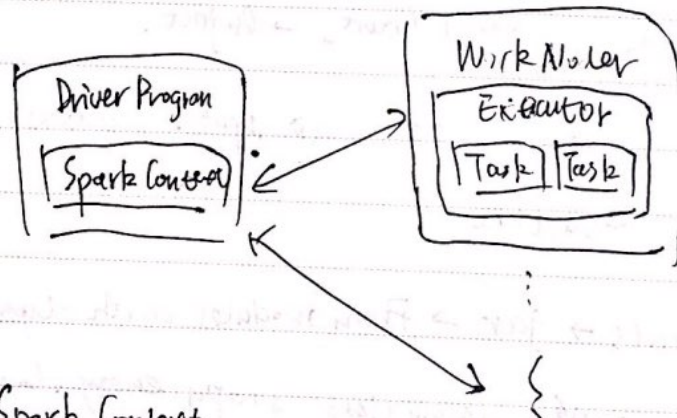
localhost:4040/jobs/ 任务进度条.

四. RDDs.

1. 介绍.

1) Driver Program:

包含 main(), RDDs 定义, 操作.



2) Spark Context

• 访问 Spark

• 一个和 Spark 集群的连接

• 在 shell 中, SC 已经创建

• RDDs 分发数据和计算的基础



RDDs:

- 一个RDD是一个不可改变的分布式集合对象。

- Spark中, 所有计算 \rightarrow RDDs 创建, 转换, 操作完成。

- 一个RDD内部由许多partitions (分片) 组成。

分片:

partitions 在不同机器。

spark 并行处理单元, spark 顺序的并行处理分片。

RDDs 的创建:

把一个已存在的集合传给 SparkContext 的 parallelize 方法。

```
var rdd = sc.parallelize(Array(1, 2, 2, 4), 4)
```

第一个参数: 待并行处理的集合, 第二个参数: 分区个数。

Scalar 变量声明。

• val, var

• val 变量值不可修改, 一旦声明不能重名。

• var 声明后, 可指向两类型的值。

Scalar 匿名函数和类型推断。

• lines.filter(line \Rightarrow line.contains("world"))

String 推断。



2. PDDs & Transformations:

`map` → `val lines = { c.parallelize` ^{Array} (`"hello"`, `"space"`, `"hello world"`, `!"`)
↓
`并行处理`

7 lines. foreach (println)

```
> val line2 = lines.map { word => (word, 1)}
```

→ line 2. `foreach(prmIn).`

```
filter > val line = lines.filter(word => word.contains("hello"))
```

```
7 wine}.foreach (println)
```

```
flatMap > val inputs = sc.textFile("x.txt")
```

7 inputs.foreach (println)

```
> val lines = inputs.flatMap(line => line.split(" "))
```

> likes. foreach (print).

3. 集句: 长集, 文集

→ val radl = sc.parallelize(数据集)

> val rad_distance = rdd1.distIn(1)

7. 打呼

```
val rddUnion = rdd1.union(rdd2)
```

~~var rdd-sub = rdd1.subtract(rdd2)~~

4. RDDs - Action

> reduce()

> val rdd = sc.parallelize(Array(1, 2, 2, 3))

> rdd.collect()

> rdd.reduce((x, y) => x + y)

> rdd.take(3)

> rdd.top(2)

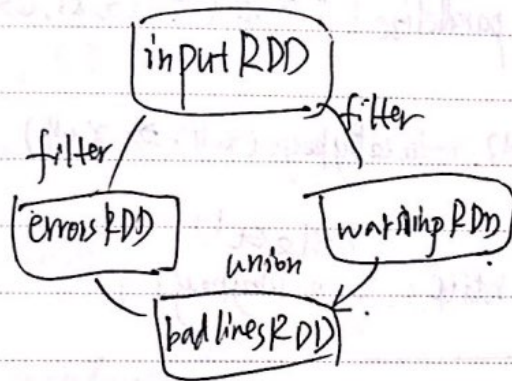
> rdd.foreach(println)

5. RDDs 特性

① 血缘关系图

1) Spark 维护着 RDD 之间的依赖关系和创建关系，叫做血缘关系图。

2) Spark 使用它来计算 RDD 的需求和恢复丢失的数据。



② 延迟计算

1) 第一次使用 action.

2) 减少数据冗余

3) metadata 减少 transform 操作已明白

4) 加载数据，100%.



③. RDDs 持久化.

`rdd.persist()` 重复利用

`rdd.unpersist()` 解除.

Memory, seq, Disk.

占用 CPU, 磁盘, 内存.

6. keyValues 对 RDDs.

`> val rdd = sc.textFile("x.txt")`

`> val rdd-map = rdd.map(line => (line.split(" ")[0], line))`

`> val rdd3 = sc.parallelize(key-value (1, 2), (3, 4), (5, 6))`

`> val rdd4 = rdd3.reduceByKey((x, y) => x + y) // (1, 2), (3, 4), (5, 6)`

`> val rdd5 = rdd4.groupByKey().collect()`

`> var rdd6 = rdd3.keys`

`> var rdd7 = rdd3.values`

`> var rdd8 = rdd3.sortByKey()`



1. combineBy key(): 聚合函数
 返回类型可以与输入类型不一致

key → 对: ~~create~~ mergeValue
 → 未对: createCombiner() \searrow partition 中已存在
 mergeCombiners 函数

求平均值

> val scores = sc.parallelized(Array(Jack, 80, Jack 90, Jack 85, Mike 80
 Mike 85, Mike 80))

> val scoreL = score.combineBy key (score \Rightarrow (1, score), (c1: (Int, Double),
 new Score) \Rightarrow (c1._1 + 1, c1._2 + newScore); (c1: (Int, Double), c2: (Int,
 Double) \Rightarrow (c1._1 + c2._1, c1._2 + c2._2))

(Jack, (3, 255.0))

(Mike, (3, 267.0))

> var average = scores2.map { (case (name, num, score)) \Rightarrow (name,
 score/num) }

⊗ (mik, 85)

(jack, 87.5)

