

2017年71月19日 16:34:55







# 大信橱柜好不好

# 



请扫描二维码联系

**2** 400-660-01

■ QQ客服

关于 招聘 广告服务 <sup>©</sup> 1999-2018 CSDN版权所有 京ICPiF09002463号

经营性网站备案信息 网络110报警服务 中国互联网举报中心 北京互联网违法和不良信息举报中心

#### 博主最新文章

#### 区块链常用架构

我们对比了GitHub上8800个开测习项目,并选出了其中的Top30

只用200行Go代码写一个自己的

数据仓库数据分层

写测试用例注意事项

#### 文章分类

#### 文章存档

2018年2月

2018年1月

2017年9月

2017年6月

2017年5月

2017年4月

展开~

## 博主热门文章

Zookeeper的几个应用场景 © 10966

只用200行Go代码写一个自己的 1169

PHP android ios相互兼容的AES

# ∩ 418

# 

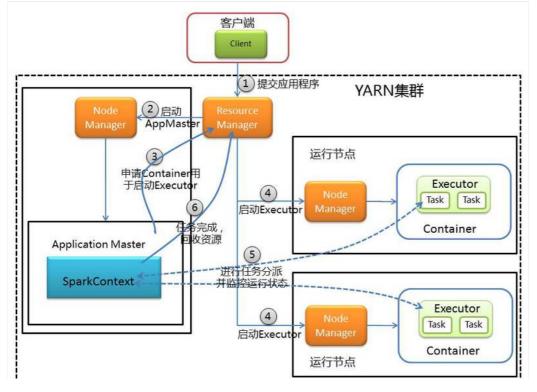
《图解Spark 为技术与案例实战》一书以Spark2.0版本为基础进行编写,系统介绍了Spark核心及其生态圈组件技术。其内容包括Spark生态圈、实战环境搭建和编程模型等,重点介绍了作业调度、容错执行、监控管理、存储管理 云行架构,同时还介绍了Spark生态圈相关组件,包括了Spark SQL的即席查询、Spark Streaming的实时流处理、MLlib的机器学习、GraphX的图处理和Alluxio的分布式内存文件系统等。下面介绍京东预测系统如何进行资源调度,并描述如何使用Spark存储相关知识进行系统优化。

大数据|Spark技术在京东智能供应链预测的应用案例深度剖析

# 4.1 结合系统中的应用

在图解Spark书的第六章描述了Spark运行架构,介绍了Spark集群资源调度一般分为粗粒度调度和细粒度调度两种模式。粗粒度包括了独立运行模式和Mesos粗粒度运行模式,在这种情况下以整个机器作为分配单元执行作业,该模式优点是由于资源长期持有减少了资源调度的时间开销,缺点是该模式中无法感知资源使用的变化,易造成系统资源的闲置,从而造成了资源浪费。而细粒度包括了Yarn运行模式和Mesos细粒度运行模式,该模式的优点是系统资源能够得到充分利用,缺点是该模式中每个任务都需要从管理器获取资源,调度延迟较大、开销较大。

由于京东Spark集群属于基础平台,在公司内部共享这些资源,所以集群采用的是Yarn运行模式,在这种模式下可以根据不同系统所需要的资源进行灵活的管理。在YARN-Cluster模式中,当用户向YARN集群中提交一个应用程序后,YARN集群将分两个阶段运行该应用程序:第一个阶段是把Spark的SparkContext作为Application Master在YARN集群中先启动;第二个阶段是由Application Master创建应用程序,然后为它向Resource Manager申请资源,并启动Executor来运行任务集,同时监控它的整个运行过程,直到运行完成。下图为Yarn-Cluster运行模式执行过程:



加入CSDN,享受更精准的内容推荐,与500万程序员共同成长!

登录

淮册

# 4.2 结合系统的优化

# 4.2.1 参数调金

- 减少num itors , 调大executor-memory , 这样的目的是希望Executor有足够的内存可以使用。
- 查看日志发现没有足够的空间存储广播变量,分析是由于Cache到内存里的数据太多耗尽了内存,于是我们将Cache的级别适当调成MEMORY\_ONLY\_SER和DISK\_ONLY。
- 针对某些任务关闭了推测机制,因为有些任务会出现暂时无法解决的数据倾斜问题,并非节点出现问题。
- 调整内存分配,对于一个Shuffle很多的任务,我们就把Cache的内存分配比例调低,同时调高Shuffle的内存比例。 4.2.2 修改设计

#### 参数的调整虽然容易做,但往往效果不好,这时候需要考虑从设计的角度去优化:

- 原先在训练数据之前会先读取历史的几个月甚至几年的数据,对这些数据进行合并、转换等一系列复杂的处理,最终生成特征数据。由于数据量庞大,任务有时会报错。经过调整后当天只处理当天数据,并将结果保存到当日分区下,训练时按天数需要读取多个分区的数据做union操作即可。
- 将"模型训练"从每天执行调整到每周执行,将"模型参数选取"从每周执行调整到每月执行。因为这两个任务都十分消耗资源,并且属于不需要频繁运行,这么做虽然准确度会略微降低,但都在可接受范围内。
- 通过拆分任务也可以很好的解决资源不够用的问题。可以横向拆分,比如原先是将100个品类数据放在一个任务中进行训练,调整后改成每10个品类提交一次Spark作业进行训练。这样虽然整体执行时间变长,但是避免了程序异常退出,保证任务可以执行成功。除了横向还可以纵向拆分,即将一个包含10个Stage的Spark任务拆分成两个任务,每个任务包含5个Stage,中间数据保存到HDFS中。

#### 4.2.3 修改程序逻辑

为了进一步提高程序的运行效率,通过修改程序的逻辑来提高性能,主要是在如下方面进行了改进:避免过多的Shuffle、减少Shuffle时需要传输的数据和处理数据倾斜问题等。

# 1. 避免过多的Shuffle

Spark提供了丰富的转换操作,可以使我们完成各类复杂的数据处理工作,但是也正因为如此我们在写Spark程序的时候可能会遇到一个陷阱,那就是为了使代码变的简洁过分依赖RDD的转换操作,使本来仅需一次Shuffle的过程变为了执行多次。我们就曾经犯过这样一个错误,本来可以通过一次groupByKey完成的操作却使用了两回。业务逻辑是这样的:我们有三张表分别是销量(s)、价格(p)、库存(v),每张表有3个字段:商品id(sku\_id)、品类id(category)和历史时序数据(data),现在需要按sku\_id将s、p、v数据合并,然后再按category再合并一次,最终的数据格式是:[category,[[sku\_id, s,p, v], [sku\_id, s,p, v], [...],[...]]]。一开始我们先按照sku\_id + category作为key进行一次groupByKey,将数据格式转换成[sku\_id, category,[s,p, v],然后按category作为key再groupByKey一次。后来我们修改为按照category作为key只进行一次groupByKey,因为一个sku\_id只会属于一个category,所以后续的map转换里面只需要写一些代码将相同sku\_id的s、p、v数据group到一起就可以了。两次groupByKey的情况:



大信橱柜好不好



与可用州 各类均衡的

#### 联系我们



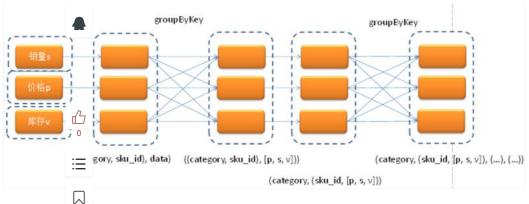
请扫描二维码联系

■ webmaster@

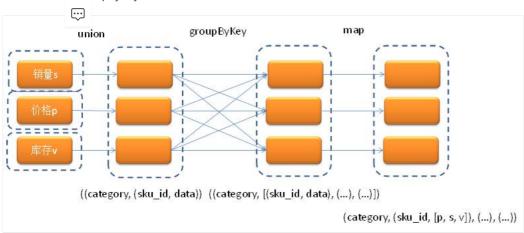
QQ客服 ● ?

关于 招聘 广告服务 📸 ©1999-2018 CSDN版权所有 京ICPiF09002463号

经营性网站备案信息 网络110报警服务 中国互联网举报中心 北京互联网违法和不良信息举报中心



修改后变为一次groupByKey的情况:



多表join时,如果key值相同,则可以使用union+groupByKey+flatMapValues形式进行。比如:需要将销量、库存、价格、促销计划和商品信息通过商品编码连接到一起,一开始使用的是join转换操作,将几个RDD 彼此join在一起。后来发现这样做运行速度非常慢,于是换成union+groypByKey+flatMapValue形式,这样做只需进行一次Shuffle,这样修改后运行速度比以前快多了。实例代码如下:

```
rdd1 = rdd1.mapValues(lambda x: (1, x))
rdd2 = rdd2.mapValues(lambda x: (2, x))
rdd3 = rdd3.mapValues(lambda x: (3, x))

def dispatch(seq):
    vbuf, wbuf, xbuf = [], [], []
    for (n, v) in seq:
        if n == 1: vbuf.append(v)
        elif n == 2: wbuf.append(v)
        elif n == 3: xbuf.append(v)
        return [(v, w, x) for v in vbuf for w in wbuf for x in xbuf]

new_rdd = sc.union([rdd1, rdd3, rdd3]).groupByKey().flatMapValues(lambda x: dispatch(x))
```

如果两个RDD需要在groupByKey后进行join操作,可以使用cogroup转换操作代替。比如 ,将历史销量数据按品类进行合并,然后再与模型文件进行join操作,流程如下:

```
[ sku_id, category, sales ]
-> [category, [ [sku_id, sales], [...] ] ]
-> [category, [ [ [sku_id, sales], [...] ], [model1, model2]] ]
```

使用cogroup后,经过一次Shuffle就可完成了两步操作,性能大幅提升。

2. 减少Shuffle时传输的数据量

在Shuffle握作前尽量将不需要的数据讨渡掉

加入CSDN,享受更精准的内容推荐,与500万程序员共同成长!



大信橱柜好不好



#### 联系我们



请扫描二维码联》

■ webmaster@

A QQ客服 ●

关于 招聘 广告服务 🐉 ©1999-2018 CSDN版权所有 京ICP证09002463号

经营性网站备案信息 网络110报警服务 中国互联网举报中心 北京互联网违法和不良信息举报中心 comebineyeByKey属于聚合类操作,由于它支持map端的聚合所以比groupByKey性能好,又由于它的map端与reduce端可以设置成不一样的逻辑,所以它支持的场景比reduceByKey多,它的定义如下:

 $\begin{tabular}{ll} $\sf def & combine By Key (self, create Combiner, & merge Value, & merge Combiners, \\ num Partitions) \end{tabular}$ 

reduceByKey,oupByKey内部实际是调用了comebineyeByKey,

```
def createCombiner(x):
    return [x]
    def mergeValue(xs, x):
        xs.append(x)
    return xs
    def mergeCombiners(a, b):
        a.extend(b)
    return a

return self combineByKey(createCombiner, mergeValue, mergeCombiners, numPartitions).mapValues(lambda x: ResultIterable(x))
```

我们之前有很多复杂的无法用reduceByKey来实现的聚合逻辑都通过groupByKey来完成的,后来全部替换为comebineyeByKey后性能提升了不少。

#### 3. 处理数据倾斜

有些时候经过一系列转换操作之后数据变得十分倾斜,在这样情况下后续的RDD计算效率会非常的糟糕,严重时程序报错。遇到这种情况通常会使用repartition这个转换操作对RDD进行重新分区,重新分区后数据会均匀分布在不同的分区中,避免了数据倾斜。如果是减少分区使用coalesce也可以达到效果,但比起repartition不足的是分配不是那么均匀。

# 5. 小结

虽然京东的预测系统已经稳定运行了很长一段时间,但是我们也看到系统本身还存在着很多待改进的地方,接下来我们会在预测准确度的提高、系统性能的优化、多业务支持的便捷性上进行改进。未来,随着大数据、人工智能技术在京东供应链管理中的使用越来越多,预测系统也将发挥出更大作用,对于京东预测系统的研发工作也将是充满着挑战与乐趣。

○ 目前您尚未登录,请 登录 或 注册 后进行评论

#### 大数据|Spark技术在京东智能供应链预测的应用案例深度剖析(一)

大数据|Spark技术在京东智能供应链预测的应用案例深度剖析(一) 2017-03-27 11:58 浏览次数:148 1. 背景 前段时间京东公开了面向第二个十二年的战略规划,表...

加入CSDN,享受更精准的内容推荐,与500万程序员共同成长!



大信橱柜好不好



#### 联系我们



请扫描二维码联系

■ webmaster@

■ QQ客服 ● ?

关于 招聘 广告服务 ©1999-2018 CSDN版权所有 京ICP证09002463号

经营性网站备案信息 网络110报警服务 中国互联网举报中心 北京互联网违法和不良信息举报中/

#### 实例讲解spark在京东智能供应链预测系统的应用

woddle 2018年02月25日 18:08 🕮 166

问题导读: 1. 京东的供应链是什么样的呢? 2. 预测技术在京东的供应链起着什么样的作用呢? 3. 京东整个预测系统的架构是什 么样的呢? 4. 预测系统不同层面的技术选型分别为什么? 5...

# 这样做,运动员技术提升100%!



# Spark技术产言东智能供应链预测的应用

 **javastart** 2017年03月10日 19:05 □ 630

Spark技术在京为 Livi Lipi链预测的应用原创 2017-03-06 杨冬越 郭景瞻 大数据杂谈 大家晚上好,做一个简单的介绍:我叫郭景 瞻,来自京东,著有《图解Spark:核心..

# 大数据分析-京东慧眼

🐶 zhuhengv 2015年10月23日 09:14 🕮 1615

1.通过市场分析,知道哪些地区卖哪些产品比如冷暖空调在南方比较畅销,单冷型空调在北方比较畅销2.通过用户分析,能确定 将哪些商品卖给哪些人 比如不会把烟酒推荐给孕妇 3.通过商品属性分...

## 区块链在供应链领域的应用

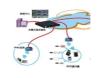
Chenhaifeng2016 2018年02月04日 13:18 🕮 1066

近年来,区块链作为一种新兴的应用模式被不同行业广泛应用。在包括金融、物联网、社会公益、供应链等领域中,出现了很多应 用落地的探索和尝试。其中,供应链领域由于具有市场规模大,及多信任主体、多方协作等特点,...

#### 物联网学习网

中国物联网工程师培训网

百度广告



# 大数据在工业4.0智能工厂中的应用

**(ingmax54212008** 2017年05月04日 10:14 🚇 594

随着近些年国家工业信息化进程脚步的不断加快,以及国际社会在工业现代化、工业4.0等方面的不断演进,使得大数据技术在工 业行业以及制造业方面也进行了比较深度的技术融合和应用融合,我们就来聊聊在上述领域的大...

## 公安大数据应用之情报分析与关联挖掘

近年来,随着信息化时代大数据的广泛运用,各类犯罪频发,犯罪分子的作案手段多变,反侦查意识进一步提升,社会不稳定因素 加剧。将大数据应用于公安领域势在必行,为适应新形势的要求,全国各地公安机关正在全面实施...

## 实战:供应链如何应用大数据

dongzhumao86 2015年03月27日 13:14 🔘 1858

摘要: 随着供应链变得越来越复杂,必须采用更好的工具来迅速高效地发挥数据的最大价值。供应链作为企业的核心网链,将彻底 变革企业市场边界、业务组合、商业模式和运作模式等。 第三产业供应链协同应用市场进...

## 浅谈大数据与智能电厂

₩ u014449866 2015年06月12日 22:56 🕮 3569

浅谈大数据与智能电厂——2015\6\12 前段时间,IBM刚刚宣布了一项新技术HyRef,用于能源电力...

### 直播:京东大数据的应用!

CSDN TG229dvt5I93mxaQ5A6U 2017年09月30日 00:00 🕮 880

前言:由CSDN主办的SDCC 2017之大数据技术实战线上峰会将在CSDN学院举行。作为SD系列技术峰会的一部分,本次线上峰 会秉承干货实料(案例)的内容原则,将邀请圈内顶尖的布道师、技术专家和技术引...

#### 程序员不会英语怎么行?

北大猛男教你:不背单词和语法,一个公式学好英语





大信橱柜好不好



#### 联系我们



请扫描二维码联

■ webmaster@ **2** 400-660-010

■ QQ客服

招聘 广告服务 ©1999-2018 CSDN版权所有

京ICP证09002463号 经营性网站备案信息 网络110报警服务 中国互联网举报中心

北京互联网违法和不良信息举报中心

案例主要关注三个问题:数据从哪里来?数据如何存储?数据如何计算?来自《Hadoop权威指南》的案例 1. Last.fm 1.1 背景 创建于2002年,提供网络电台和网络音乐服务的...

## 【智能物流】罗戈研究院京东物流《数字化供应链综合研究报告》;AI智慧仓储和物流...

来源丨罗戈研究院平台基于个人消费习惯进行智能推荐商品、定制化门槛不再高,交易可以通过移动端应用完成、交付逐渐快至半 小时……这个世界 / 3 "数" 发生巨大的变化: • 一切都在数字化• 一切都将通过云共…

# -基于Spark的大数据实时处理及应用技术

培训要点 互联网占未数据、传感数据、日志文件、具有丰富地理空间信息的移动数据和涉及网络的各类评论,成为了海量信息的 多种形式。当数技 💬 百上千TB不断增长的时候,我们在内部交易系统的历史信息之外,需要...

Shenmanli 2016年03月11日 14:57 □ 1804

# 大数据在京东的典型应用:京东用户画像技术曝光

http://www.36dsj.com/archives/16090 大数据在京东的典型应用:京东用户画像技术曝光 数控小V 2014-11-05 9:02:34 大 数...

👣 giezikuaichuan 2015年09月18日 19:51 🚇 1668

#### 射频技术系列谈:射频技术及其在供应链管理中的应用之一

在《IT经理世界》与IDC公司联合举办的"2003年度优秀CIO评选"的颁奖大会上,IDC的首席研究官John Gantz 就信息技术的 发展趋势做了精彩的讲演。其中,他特别强调了射频技术将是未来IT...

🦜 zhaoyang17 2004年10月02日 19:08 🕮 2220

# 国内十大猎头公司排名

十大猎头公司

百度广告



#### 京东供应链溯源防伪平台

🍞 tigerking1017 2018年01月17日 15:51 🔘 388

PPT整理自:"别人在忙挖矿,京东架构师却悄悄用区块链搞了件大事!" http://blog.csdn.net/dev\_csdn/article/details/790 62081 仅供参考学习...

#### 京东商城大数据面试题

👜 qq\_26442553 2017年12月05日 12:42 🕮 1520

京东商城 - 大数据(1) Java篇 1、JVM, GC(算法,新生代,老年代),JVM结构 2、hashcode, hashMap, list, hashSe t, equals (结构原理), A exten...

# 京东DNN实验室: 大数据、深度学习与计算平台的实践 🦠 jdbc 2015年09月09日 18:45 🕮 1722

7月26日-27日,2015中国人工智能大会(CCAI 2015)在北京召开,深度学习毫无意外地成为与会嘉宾热议的一个话题。来自京 东DNN实验室的四位专家,核心科学家李成华、张晓鑫,以及京东智能通讯部...

#### 电商大数据——用数据驱动电商和商业案例解析

电商大数据——用数据驱动电商和商业案例解析(国内第1本将大数据与电商完美结合的权威之作!) 雪鹰传奇 著 ISBN 978-7-121-22556-7 2014年3月出版 定价:...

broadview2006 2014年04月17日 15:22 Q 3498

### 【智能制造】爱(AI)在新工业

■ np4rHI455vg29y2 2017年12月06日 00:00 □ 275

本文是工业4.0俱乐部秘书长杜玉河老师在2017国际人工智能大会发表演讲时所用的PPT。人工智能是我们近期比较关注的领域,

加入CSDN,享受更精准的内容推荐,与500万程序员共同成长!



大信橱柜好不好



#### 联系我们



请扫描二维码联

™ webmaster@ **2** 400-660-010

QQ客服 ● ?

招聘 广告服务

©1999-2018 CSDN版权所有 京ICP证09002463号

经营性网站备案信息 网络110报警服务 中国互联网举报中心 北京互联网违法和不良信息举报中心