阿里云 DataWorks

最佳实践

文档版本: 20190808

为了无法计算的价值 | []阿里云

<u>法律声明</u>

阿里云提醒您在阅读或使用本文档之前仔细阅读、充分理解本法律声明各条款的内容。如果您阅读 或使用本文档,您的阅读或使用行为将被视为对本声明全部内容的认可。

- 您应当通过阿里云网站或阿里云提供的其他授权通道下载、获取本文档,且仅能用于自身的合法 合规的业务活动。本文档的内容视为阿里云的保密信息,您应当严格遵守保密义务;未经阿里云 事先书面同意,您不得向任何第三方披露本手册内容或提供给任何第三方使用。
- 未经阿里云事先书面许可,任何单位、公司或个人不得擅自摘抄、翻译、复制本文档内容的部分 或全部,不得以任何方式或途径进行传播和宣传。
- 3. 由于产品版本升级、调整或其他原因,本文档内容有可能变更。阿里云保留在没有任何通知或者 提示下对本文档的内容进行修改的权利,并在阿里云授权通道中不时发布更新后的用户文档。您 应当实时关注用户文档的版本变更并通过阿里云授权渠道下载、获取最新版的用户文档。
- 4. 本文档仅作为用户使用阿里云产品及服务的参考性指引,阿里云以产品及服务的"现状"、"有缺陷"和"当前功能"的状态提供本文档。阿里云在现有技术的基础上尽最大努力提供相应的介绍及操作指引,但阿里云在此明确声明对本文档内容的准确性、完整性、适用性、可靠性等不作任何明示或暗示的保证。任何单位、公司或个人因为下载、使用或信赖本文档而发生任何差错或经济损失的,阿里云不承担任何法律责任。在任何情况下,阿里云均不对任何间接性、后果性、惩戒性、偶然性、特殊性或刑罚性的损害,包括用户使用或信赖本文档而遭受的利润损失,承担责任(即使阿里云已被告知该等损失的可能性)。
- 5. 阿里云网站上所有内容,包括但不限于著作、产品、图片、档案、资讯、资料、网站架构、网站画面的安排、网页设计,均由阿里云和/或其关联公司依法拥有其知识产权,包括但不限于商标权、专利权、著作权、商业秘密等。非经阿里云和/或其关联公司书面同意,任何人不得擅自使用、修改、复制、公开传播、改变、散布、发行或公开发表阿里云网站、产品程序或内容。此外,未经阿里云事先书面同意,任何人不得为了任何营销、广告、促销或其他目的使用、公布或复制阿里云的名称(包括但不限于单独为或以组合形式包含"阿里云"、Aliyun"、"万网"等阿里云和/或其关联公司品牌,上述品牌的附属标志及图案或任何类似公司名称、商号、商标、产品或服务名称、域名、图案标示、标志、标识或通过特定描述使第三方能够识别阿里云和/或其关联公司)。
- 6. 如若发现本文档存在任何错误,请与阿里云取得直接联系。

通用约定

格式	说明	样例				
•	该类警示信息将导致系统重大变更甚至 故障,或者导致人身伤害等结果。	禁止: 重置操作将丢失用户配置数据。				
A	该类警示信息可能导致系统重大变更甚 至故障,或者导致人身伤害等结果。	▲ 警告: 重启操作将导致业务中断,恢复业务所需 时间约10分钟。				
Ê	用于补充说明、最佳实践、窍门等,不 是用户必须了解的内容。	道 说明: 您也可以通过按Ctrl + A选中全部文件。				
>	多级菜单递进。	设置 > 网络 > 设置网络类型				
粗体	表示按键、菜单、页面名称等UI元素。	单击 确定。				
courier 字体	命令。	执行 cd /d C:/windows 命令,进 入Windows系统文件夹。				
##	表示参数、变量。	bae log listinstanceid Instance_ID				
[]或者[a b]	表示可选项,至多选择一个。	ipconfig [-all -t]				
	表示必选项,至多选择一个。	<pre>swich {stand slave}</pre>				

目录

法律声明	I
通用约定	I
1 数据迁移	1
1.1 Hadoop数据迁移MaxCompute最佳实践	1
1.2 Kafka数据迁移MaxCompute最佳实践	18
1.3 JSON数据从OSS迁移至MaxCompute	29
1.4 JSON数据从MongoDB迁移至MaxCompute	37
1.5 Elasticsearch数据迁移至MaxCompute	47
2数据开发	54
2.1 设置调度依赖最佳实践	54
2.2 Eclipse Java UDF开发最佳实践	60
2.3 使用MaxCompute分析IP来源最佳实践	74
2.4 在PyODPS任务中调用第三方包	81
2.5 分支节点实现特定时间执行任务	84
2.6 DataWorks数据服务对接DataV最佳实践	92
2.7 天依赖分钟任务最佳实践	105
2.8 邮件外发最佳实践	111
2.9 PyODPS节点实现结巴中文分词	112
2.10 基于AnalyticDB构建企业数仓	123
3数据安全	139
3.1 实现指定用户访问指定UDF最佳实践	139
3.2 子账号仅从特定IP登录DataWorks	147

1数据迁移

1.1 Hadoop数据迁移MaxCompute最佳实践

本文将为您介绍如何通过DataWorks数据同步功能,将HDFS上的数据迁移至MaxCompute,或 从MaxCompute将数据迁移至HDFS。无论您使用Hadoop还是Spark,均可以 与MaxCompute之间的数据进行双向同步。

环境准备

1. Hadoop集群搭建

进行数据迁移前,您需要保证自己的Hadoop集群环境正常。本文使用阿里云EMR服务自动化 搭建Hadoop集群,详细过程请参见创建集群。

本文使用的EMR Hadoop版本信息如下:

EMR版本: EMR-3.11.0

集群类型: HADOOP

软件信息: HDFS2.7.2 / YARN2.7.2 / Hive2.3.3 / Ganglia3.7.2 / Spark2.2.1 / HUE4.1.0 / Zeppelin0.7.3 / Tez0.9.1 / Sqoop1.4.6 / Pig0.14.0 / ApacheDS2.0.0 / Knox0.13.0

Hadoop集群使用经典网络,区域为华东1(杭州),主实例组ECS计算资源配置公网及内网IP,高可用选择为否(非HA模式)。

集群信息						
ID: 地域: cn-hangzhou 开始时间: 2018-09-03 17:28:2	25	软件配置: 10优化: 是 高可用: 否 安全候式: 标 准	付赉类型: 按量付累 当前状态: 空闲 运行时间: 2天23小		引导操作/软件配置: EMR-3.110 ECS应用角色: AliyunEmrEcsDefaultRole	
软件信息				网络信息		
EMR版本: EMR-3.11.0 集群类型: HADOOP 软件信息: HDFS2.7.2 / YARN2. ApacheDS2.0.0 / Knox0.13.0	2.7.2 / Hive2.3.3 / Ganglia3.7	.2 / Spark2.2.1 / HUE4.1.0 / Zeppelin0.7.3 / Tez0.9.1 / Sqoo	p1.4.6 / Pig0.14.0 /	区城ID: cn-hangzl 网络类型: classic 安全组ID:	hou-f	
主机信息	C	主实例组 🛃				
主实例组(MASTER)	按量付费	ECS ID 状	5	公网	内网	创建时间
主机数量:1 公网 CPU:4核 内存 数据母配管:SSD元母80G8*114	羽带宽: 8M 霁: 8GB	****************	正常		10.80.63.61	2018-09-03 17:28:34
XAMBOR. 5502200000 1X						
核心实例组(CORE)	按量付费					

2. MaxCompute

<mark>开通MaxCompute</mark>并创建好项目。本文以在华东1(杭州)区域创建项目bigdata_DOC为例、同时启动DataWorks相关服务。

数据准备

- 1. Hadoop集群创建测试数据
 - a. 进入EMR Hadoop集群控制台界面,使用交互式工作台,新建交互式任务doc。本例 中HIVE建表语句如下。

```
CREATE TABLE IF NOT

EXISTS hive_doc_good_sale(

    create_time timestamp,

    category STRING,

    brand STRING,

    buyer_id STRING,

    trans_num BIGINT,

    trans_amount DOUBLE,

    click_cnt BIGINT

    )

    PARTITIONED BY (pt string) ROW FORMAT

DELIMITED FIELDS TERMINATED BY ',' lines terminated by '\n'
```

b. 选择运行,出现Query executed successfully提示,则说明成功在EMR Hadoop集

群上创建了表hive_doc_good_sale。

管理控制台 🛛 🔛 华东1	(杭州)▼					搜索			消息	費用	工单	备案	企业	支持与	服务	= î	简体中文	
E-MapReduce管理控制台	交互式功能关联的集群最少3台机器,最	低配置4core8GB,	EMR-2.3以及以	止版本														
概范	交互式工作台																	
集群	交互式任务列表	doc (HIVE)	980412a4-	-5a76-4b96-	-af41-6020	c9bf08430											•	全屏
交互式工作台	bigdata hive	▶ 文件	▶ 视图▼	▶运行全部	部								类型:	HIVE	关联集群	: 8	bigdata1	2 🕶
表管理	DOC1																	
作业	doc												P	保存段落	- 121	藏结果	×m	除
 执行针划 較限开发(呼呼) ● 报警 帮助 		 CREATE cre cat bra buy tra cli) PAR 运行结果: Query exe 状态: FINISI 	TABLE IF NG ate_time tin geory STRIN egory STRIN er_id STRING, er_id STRING, er_id STRING, er_id STRING provide STRING er, id S	T EXISTS hi restamp, , , , , , , , , , , , , ,	ive_doc_goo	od_sale(AT DELIMITEN s : -1 :45:26 PM) FIELDS TER	RMINATED	ВҮ','	lines t	erminato	ed by '	'n'					4

c. 插入测试数据。您可以选择从OSS或其他数据源导入测试数据,也可以手动插入少量的测试

数据。本文中手动插入数据如下。

```
insert into
hive_doc_good_sale PARTITION(pt =1 ) values('2018-08-21','外套','品
牌A','lilei',3,500.6,7),('2018-08-22','生鲜','品牌B','lilei',1,303
,8),('2018-08-22','外套','品牌C','hanmeimei',2,510,2),(2018-08-22
,'卫浴','品牌A','hanmeimei',1,442.5,1),('2018-08-22','生鲜','品牌D
','hanmeimei',2,234,3),('2018-08-23','外套','品牌B','jimmy',9,2000,
7),('2018-08-23','生鲜','品牌A','jimmy',5,45.1,5),('2018-08-23','外
```

套','品牌E','jimmy',5,100.2,4),('2018-08-24','生鲜','品牌G','peiqi', 10,5560,7),('2018-08-24','卫浴','品牌F','peiqi',1,445.6,2),('2018-08-24','外套','品牌A','ray',3,777,3),('2018-08-24','卫浴','品牌G',' ray',3,122,3),('2018-08-24','外套','品牌C','ray',1,62,7);

d. 完成插入数据后,您可以执行select * from hive_doc_good_sale where pt =1
 ;语句,检查Hadoop集群表中是否已存在数据可以用于迁移。

						保存段落	= 隐藏结果 × 删除
> select * from	n hive_doc_good_sa	le where pt =1;					
▶运行							
运行结果:							
hive_doc_good_s ale.create_time	hive_doc_good_s ale.category	hive_doc_good_s ale.brand	hive_doc_good_s ale.buyer_id	hive_doc_good_s ale.trans_num	hive_doc_good_s ale.trans_amount	hive_doc_good_s ale.click_cnt	hive_doc_good_s ale.pt
2018-08-21 00:00:0 0.0	外套	品牌A	lilei	3	500.6	7	1
2018-08-22 00:00:0 0.0	生鮮	品牌B	lilei	1	303.0	8	1
2018-08-22 00:00:0 0.0	外套	品牌C	hanmeimei	2	510.0	2	1
null	卫浴	品牌A	hanmeimei	1	442.5	1	1
2018-08-22 00:00:0 0.0	生鲜	品牌D	hanmeimei	2	234.0	3	1
2018-08-23 00:00:0 0.0	外套	品牌B	jimmy	9	2000.0	7	1

2. 利用DataWorks新建目标表

- a. 登录DataWorks控制台,单击相应工作空间操作栏下的进入数据开发。
- b. 进入DataStudio(数据开发)页面,选择新建>表。



- c. 在新建表对话框中, 填写表名, 并单击提交。
- d. 进入新建表页面,选择DDL模式。
- e. 在DDL模式对话框中输入建表语句,单击生成表结构,并确认操作。本示例的建表语句如下 所示:

```
CREATE TABLE IF NOT EXISTS hive_doc_good_sale(
    create_time string,
    category STRING,
    brand STRING,
    buyer_id STRING,
    trans_num BIGINT,
    trans_amount DOUBLE,
    click_cnt BIGINT
    )
```

```
PARTITIONED BY (pt string) ;
```

在建表过程中,需要考虑HIVE数据类型与MaxCompute数据类型的映射,当前数据映射关系请参见与Hive数据类型映射表。

由于本文使用DataWorks进行数据迁移,而DataWorks数据同步功能暂不支 持TIMESTAMP类型数据。因此在DataWorks建表语句中,将create_time设置 为STRING类型。上述步骤同样可通过odpscmd命令行工具完成,命令行工具安装和配置请 参见安装并配置客户端。



```
1 说明:
```

考虑到部分HIVE与MaxCompute数据类型的兼容问题,建议在odpscmd客户端上执行以 下命令。

set odps.sql.type.system.odps2=true;



f. 完成建表后,单击左侧导航栏中的表管理,即可查看当前创建的MaxCompute表。



数据同步

1. 新建自定义资源组

由于MaxCompute项目所处的网络环境与Hadoop集群中的数据节点(data node)网络通常不可达,您可以通过自定义资源组的方式,将DataWorks的同步任务运行在Hadoop集群的 Master节点上(Hadoop集群内Master节点和数据节点通常可达)。

a. 查看Hadoop集群data node

进入EMR控制台,选择首页>集群管理>集群>主机列表。

E-MapReduce	集群管理数据开发	系统维护 操作日志	5. 帮助					
bigdata12 👻	前页 〉 集群管理 〉 集群 (C-)	> 主机列表					
88 集群基础信息	主机列表 当前集群: C-	/ bigdata	12					同步主机信息
③ 集群与服务管理	ECS InstanceID	主机名	内岡中				直询	
◎ 主机列表	ECS ID	主机名	IP信息	角色 🏹	所屬机器组	付盡类型	规档	到期时间
 2 集計脚本 > 访问链接与端口 △ 用户管理 	i-bp1hif	emr-header-1	内岡:10.80.63.61 外岡:12	MASTER	MASTER	按量付费	CPU-4 核 内存-8G ECS 規慎ecs.n4 xlarge 数据曲配置_SSD云曲 80 X 1块 系统曲配置_SSD云曲 20 X 1块	
◎ 弹性伸缩	i-bp1emov	emr-worker-2	内丽:10.81.78.209	CORE	CORE	按量付费	CPU4 核 内存 8G ECS 規指 ecs.n4 xlarge 数据盘配置 SSD云曲 80 X 4块 系统盘配置 SSD云曲 80 X 1块	
	i-bp1de7	emr-worker-1	内岡:10.31.122.189	CORE	CORE	按量付费	CPU-4 核 内存-8G ECS 規悟 ecs.n4.xlarge 数据金配置。SSD云盘 80 X 4块 系统金配置。SSD云盘 80 X 1块	

您也可以通过单击上图中Master节点的ECS ID,进入ECS实例详情页。然后单击远程连接进入ECS,执行hadoop dfsadmin -report命令查看data node。

DFS Used:: 0.05% Under replicated blocks: 0 Blocks with corrupt replicas: 0 Missing blocks: 0 Missing blocks (with replication factor 1): 0 Live datanodes (2): Name: 10.31.122.189:50010 (emr-worker-1.cluster-74503) Hostname: emr-worker-1.cluster-74503 Decommission Status : Normal Configured Capacity: 333373341696 (310.48 GB) DFS Used: 155725824 (148.51 MB) Non DFS Used: 325541888 (310.46 MB) DFS Remaining: 332892073984 (310.03 GB) DFS Used:: 0.05% DFS Remaining%: 99.86% Configured Cache Capacity: 0 (0 B) Cache Used: 0 (0 B) Cache Remaining: 0 (0 B) Cache Used%: 100.00% Cache Remaining%: 0.00% Xceivers: 1 Last contact: Thu Sep 06 19:41:01 CST 2018 Name: 10.81.78.209:50010 (emr-worker-2.cluster-74503) Hostname: emr-worker-2.cluster-74503 Decommission Status : Normal Configured Capacity: 333373341696 (310.48 GB) DFS Used: 155725824 (148.51 MB) Non DFS Used: 325451776 (310.38 MB) DFS Remaining: 332892164096 (310.03 GB) DFS Used:: 0.05% DFS Remaining%: 99.86% Configured Cache Capacity: 0 (0 B) Cache Used: 0 (0 B) Cache Remaining: 0 (0 B) Cache Used:: 100.00% Cache Remaining%: 0.00% kceivers: 1 Last contact: Thu Sep 06 19:41:02 CST 2018

本示例的data node只具有内网地址,很难与DataWorks默认资源组互通,所以需要设置自定义资源组,将master node设置为执行DataWorks数据同步任务的节点。

b. 新建自定义资源组

A. 进入数据集成 > 资源组页面,单击右上角的新增自定义资源组。关于自定义资源组的详细 信息请参见新增任务资源。



目前仅专业版及以上版本方可使用此入口。

B. 添加服务器时,需要输入ECS UUID和机器IP等信息(对于经典网络类型,需要输入服务 器名称。对于专有网络类型,需要输入服务器UUID)。目前仅DataWorks V2.0华东2区 支持经典网络类型的调度资源添加,对于其他区域,无论您使用的是经典网络还是专有网 络类型,在添加调度资源组时都请选择专有网络类型。

机器IP需要填写master node公网IP(内网IP有可能不可达)。ECS的UUID需 要进入master node管理终端,通过命令dmidecode | grep UUID获取(如果您 的hadoop集群并非搭建在EMR环境上,也可以通过该命令获取)。



- C. 添加服务器后,需要保证master node与DataWorks网络可达。如果您使用的是ECS服务器,需要设置服务器安全组。
 - ·如果您使用的内网IP互通,请参见添加安全组。
 - ·如果您使用的是公网IP,可以直接设置安全组公网出入方向规则。本文中设置公网入 方向放通所有端口(实际应用场景中,为了您的数据安全,强烈建议设置详细的放通规 则)。

<	bigdata12						教我设置	3 返回 添加安全组规则	快速创建规则 添加ClassicLink安全组规则
安全组规则 安全组内实例列表	内网入方向内	网出方向 公网)	方向公网出方	6)					▲ 导入规则 ▲ 导出全部规则
安全组内弹性网卡	□ 授权策略	协议类型	第日范围	授权类型	授权对象	描述	优先级	创建时间	操作
	□ 允许	全部	-1/-1	地址段访问	0.0.0.0/0		1	2018年9月4日 14:35	修改 売編 撤除

D. 完成上述步骤后,按照提示安装自定义资源组agent。当前状态显示为可用时,则新增自 定义资源组成功。

管理资源组 - hdfs				×
			刷新	新加服务器
服务器名称/ECS UUID	服务器IP	当前状态	已使用DMU	操作
F631D86C-		可用	0	修改 删除

如果状态为不可用,您可以登录master node,执行tail -f/home/admin/

alisatasknode/logs/heartbeat.log命令查看DataWorks与master node之间心

跳报文是否超时。

[root@emr-header-1 logs]# hdf:	s dfs -ls /user/hive/warehouse/hive_doc_good_sale/
Found 1 items	
drwxr-xx - hive hadoop	0 2018-09-03 17:46 /user/hive/warehouse/hive_doc_good_sale/pt=1
[root@emr-header-1 logs]# tai]	l -f /home/admin/alisatasknode/logs/heartbeat.log
2018-09-06 21:47:34,440 INFO	[pool-6-thread-1] [HeartbeatReporter.java:104] [] - heartbeat start, current status:2
2018-09-06 21:47:34,465 INFO	[pool-6-thread-1] [HeartbeatReporter.java:133] [] - heartbeat end∎ cost time:0.025s
2018-09-06 21:47:39,465 INFO	[pool-6-thread-1] [HeartbeatReporter.java:104] [] - heartbeat start, current status:2
2018-09-06 21:47:39,491 INFO	[pool-6-thread-1] [HeartbeatReporter.java:133] [] - heartbeat end∎ cost time:0.026s
2018-09-06 21:47:44,491 INFO	[pool-6-thread-1] [HeartbeatReporter.java:104] [] - heartbeat start, current status:2
2018-09-06 21:47:44,515 INFO	[pool-6-thread-1] [HeartbeatReporter.java:133] [] - heartbeat end∎ cost time:0.024s
2018-09-06 21:47:49,516 INFO	[pool-6-thread-1] [HeartbeatReporter.java:104] [] - heartbeat start, current status:2
2018-09-06 21:47:49,538 INFO	[pool-6-thread-1] [HeartbeatReporter.java:133] [] - heartbeat end∎ cost time:0.022s
2018-09-06 21:47:54,539 INFO	[pool-6-thread-1] [HeartbeatReporter.java:104] [] - heartbeat start, current status:2
2018-09-06 21:47:54,555 INFO	[pool-6-thread-1] [HeartbeatReporter.java:133] [] - heartbeat end∎ cost time:0.016s

2. 新建数据源

DataWorks新建工作空间后,默认设置自己为数据源odps_first。因此只需要添加Hadoop集 群数据源。更多详情请参见配置HDFS数据源。

- a. 进入数据集成页面,选择同步资源管理 > 数据源,单击新增数据源。
- b. 在新增数据源弹出框中,选择数据源类型为HDFS。

⑤ ○ 数据集成													
三 ▼ 任务列表	数据源类型: 全部		新增数据源					×		3 刷新 多库多	表搬迁 批量	19F182017 2	1993.CC
👑 高线同步任务			关系型数据库										
💂 同步资源管理	数据源名称	数据源类型 链接		ð	(I)	ORACLE'		- 8	连遍状态	连通时间	适用环境	操作	选择
1 august 1	ordine firret	Endj 项目	MySQL.	SQL Server	PostgreSQL PostgreSQL	Oracle	DM	1			开发		
		Endy 项目		5	o:Åv•	6		- 8			生产		
★ 批量上云		95-10	DRDS	POLARDR	Hubrid DR for MrSOI	AnalyticDR for		- 8					
		sourn 实例 User	2: 100	FOLMOD	Hybridob for My342	PostgreSQL		- 8	成功	2019/07/04 16:04:32	开发	整库迁移批量配 编辑 删除	
	rds_workshop_log	MySQL 数据 实例 User	大数读存储 车2 Aan	×	\diamond	47	\odot				生产	编辑 删除	
	ors workshop log	Acc Buc End	sst et : MaxCompute (ODPS) oin	DataHub	AnalyticDB (ADS)	Lightning	Data Lake Analytics(DLA)	1	成功	2019/07/04 16:04:49	开发	1948 MIN	
	cos_conksnop_log	Acco Buci Endj	est 半结构化存储 et のin	jçe	3 2						生产	145 BIN	
			OSS	HDFS	FTP								

c. 填写HDFS数据源的各配置项。

新增HDFS数据源		×
* 数据源名称:	自定义名称	
数据源描述:		
*适用环境:	✔ 开发 生产	
* DefaultFS :	格式:hdfs://ServerIP:Port	?
测试连通性:	测试连通性	
	上一步	記載

配置	说明
数据源名称	数据源名称必须以字母、数字、下划线组合,且不能以数 字和下划线开头。
数据源描述	对数据源进行简单描述,不得超过80个字符。

配置	说明
适用环境	可以选择开发或生产环境。 道 说明: 仅标准模式工作空间会显示此配置。
DefaultFS	对于EMR Hadoop集群而言,如果Hadoop集群为HA集 群,则此处地址为hdfs://emr-header-1的IP:8020。 如果Hadoop集群为非HA集群,则此处地址为hdfs:// emr-header-1的IP:9000。 本实验中的emr-header-1与DataWorks通过公网连 接,因此此处填写公网IP并放通安全组。

d. 完成配置后,单击测试连通性。

e. 测试连通性通过后,单击完成。

〕 说明:

如果EMR Hadoop集群设置网络类型为专有网络,则不支持连通性测试。

3. 配置数据同步任务

a. 进入数据开发页面,选择新建 > 数据集成 > 数据同步。



- b. 在新建节点对话框中,输入节点名称,单击提交。
- c. 成功创建数据同步节点后,单击工具栏中的转换脚本按钮。



d. 单击提示对话框中的确认,即可进入脚本模式进行开发。



e. 单击工具栏中的导入模板按钮。

Di write	e_result	•				
Ľ	\odot	Þ	1	ե	🔂 💽 🗱	
1 2 3	{	"type "step	": "j(s": [ob",	导入模板	

f. 在导入模板对话框中,选择来源类型、数据源、目标类型及数据源,单击确认。

导入模板			×
	* 来源类型	HDFS V	?
	* 数据源		
		新唱数据源	
	* 日标类型	ODPS V	?
	* 数据源	odps_first (odps)	
		确认	取消

g. 新建同步任务完成后,通过导入模板已生成了基本的读取端配置。此时您可以继续手动配置数据同步任务的读取端数据源,以及需要同步的表信息等。本示例的代码如下所示,更多详情请参见配置HDFS Reader。

```
{
    "configuration": {
        "reader": {
            "plugin": "hdfs",
            "parameter": {
                "path": "/user/hive/warehouse/hive_doc_good_sale/",
                "path": "/user/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/warehouse/hive/wa
```

```
"datasource": "HDFS1",
                    "column": [
                             {
                                      "index": 0,
                                       "type": "string"
                            },
                             {
                                       "index": 1,
                                       "type": "string"
                             },
                             {
                                      "index": 2,
"type": "string"
                            },
                             {
                                      "index": 3,
"type": "string"
                            },
                             {
                                      "index": 4,
"type": "long"
                            },
{
                                       "index": 5,
                                       "type": "double"
                            },
                             {
                                       "index": 6,
                                       "type": "long"
                            }
                   ],
"defaultFS": "hdfs://121.199.11.138:9000",
"fieldDelimiter": ",",
"encoding": "UTF-8",
"filleTuro": "text"
          }
},
"writer": {
    "ugin":
          "plugin": "odps",
           "parameter": {
                   "partition": "pt=1",
                    "truncate": false,
                    "datasource": "odps_first",
                    "column": [
                             "create_time",
                             "category",
                             "brand".
                             "buyer_id"
                             "trans_num",
                             "trans_amount",
                            "click_cnt"
                   ],
"table": "hive_doc_good_sale"
         }
},
"setting": {
    setting": {
        imit
        imit

          "errorLimit": {
                   "record": "1000"
         },
"speed": {
    "sbrottle"

                    "throttle": false,
                    "concurrent": 1,
                    "mbps": "1",
          }
```



其中,path参数为数据在Hadoop集群中存放的位置.您可以在登录master node后,执行 hdfs dfs -ls /user/hive/warehouse/hive_doc_good_sale命令确认。对于分区 表,您可以不指定分区,DataWorks数据同步会自动递归到分区路径。

```
[root@emr-header-1 logs]# hdfs dfs -ls /user/hive/warehouse/hive_doc_good_sale/
Found 1 items
drwxr-x--x - hive hadoop 0 2018-09-03 17:46 /user/hive/warehouse/hive_doc_good_sale/pt=1
```

h. 完成配置后,单击运行。如果提示任务运行成功,则说明同步任务已完成。如果运行失败,可以通过日志进行排查。

验证结果

- 1. 单击左侧导航栏中的临时查询。
- 2. 选择新建 > ODPS SQL。

		临时查询	윤 🗟 🖸	C O	🔤 create_table_ddl 🗙 🚑 workshop 🗙 📾 运行日志 🗙 💽 ftp_数据	同步
	数据开发	文件名称/创建人	2	文件夹	· · · · · · · · ·	
*	组件管理	✔ 临时查询		新建	DOPS SQL TABLE IF NOT EXISTS ods_user_info_d (
R	临时直询 1				SHELL STRING COMMENT '用户ID', 55 55 55 55 55 55 55 55 55 55 55 55 55	
Ē.	运行历史				20 age_range STRING COMMENT '年龄段', 21 zodiac STRING COMMENT '星座'	
	手动业务流程 New				22) 23 PARTITIONED BY (
=	公共表				24 dt STRING 25);	
R	表管理					

3. 编写并执行SQL语句,查看导入hive_doc_good_sale的数据。SQL语句如下所示:

```
--查看是否成功写入MaxCompute
select * from hive_doc_good_sale where pt=1;
```

您也可以在odpscmd命令行工具中输入select * FROM hive_doc_good_sale where

pt =1;, 查询表结果。

MaxCompute数据迁移到Hadoop

如果您想实现MaxCompute数据迁移至Hadoop,步骤与上述步骤类似,不同的是同步脚本内的reader和writer对象需要对调,具体实现脚本如下。

```
{
    "configuration": {
        "reader": {
            "plugin": "odps",
            "parameter": {
            "partition": "pt=1",
            "pt=1",
            "pt=1",
            "partition": "pt=1",
            "pt=1",
```

```
"isCompress": false,
    "datasource": "odps_first",
    "column": [
       "create_time",
       "category",
       "brand"
    "buyer_id"
    "trans_num",
    "trans_amount",
    "click_cnt"
  "table": "hive_doc_good_sale"
  }
},
"writer": {
"plugin": "hdfs",
  "parameter": {
  "path": "/user/hive/warehouse/hive_doc_good_sale",
  "fileName": "pt=1",
"datasource": "HDFS_data_source",
  "column": [
    {
       "name": "create_time",
       "type": "string"
    },
     {
       "name": "category",
       "type": "string"
    },
{
       "name": "brand"
       "type": "string"
    },
     {
       "name": "buyer_id",
       "type": "string"
    },
     {
       "name": "trans_num",
       "type": "BIGINT"
    },
     {
       "name": "trans_amount",
       "type": "DOUBLE"
    },
    {
       "name": "click_cnt",
       "type": "BIGINT"
    }
  ],
  "defaultFS": "hdfs://47.99.162.100:9000",
"writeMode": "append",
"fieldDelimiter": ",",
  "encoding": "UTF-8",
  "fileType": "text"
  }
"errorLimit": {
    "record": "1000"
"speed": {
  "throttle": false,
  "concurrent": 1,
```

```
"mbps": "1",
}
},
"type": "job",
"version": "1.0"
}
```

您需要参见配置HDFS Writer,在运行上述同步任务前,对Hadoop集群进行设置。在运行同步任务后,手动复制同步过去的文件。

1.2 Kafka数据迁移MaxCompute最佳实践

本文将为您介绍如何使用DataWorks数据同步功能,将Kafka集群上的数据迁移至阿里云大数据计 算服务MaxCompute。

前提条件

・搭建Kafka集群

进行数据迁移前,您需要保证自己的Kafka集群环境正常。本文使用阿里云EMR服务自动化搭 建Kafka集群,详细过程请参见Kafka快速入门。

本文使用的EMR Kafka版本信息如下:

- EMR版本: EMR-3.12.1
- 集群类型: Kafka
- 软件信息: Ganglia 3.7.2 ZooKeeper 3.4.12 Kafka 2.11-1.0.1 Kafka-Manager 1.3.3.16

Kafka集群使用专有网络,区域为华东1(杭州),主实例组ECS计算资源配置公网及内 网IP,具体配置如下图所示。

首页 > 集群管理 > 集群	(C-EE) > 详	情										
集群基础信息				· • • • • • • • • • • • • • • • • • • •	源变配 🖌	● 网络管理 🖌	■ 費用管理 ~	日 实例状态	管理 🖌			
集群信息												
集群名称: kafka2mc		IO优化:是	Ŧ	开始时间: 2019年5月27日 11:48:32	2	统一元数据	否					
集群ID: C-EE		高可用: 否	ſ	封费类型:按量付费		引导操作/\$	次件配置: EMR-3.12.1					
地域: cn-hangzhou		安全模式: 标准	i	运行时间:1天5小时13分6秒		ECS应用角1	色: AliyunEmrEcsDefau	ltRole				
当前状态: 空闲												
软件信息				网络信息								
EMR版本: EMR-3.12.1				区城ID: cn-hangzhou-h								
集群类型: Kafka				网络类型: vpc								
软件信息: Ganglia 3.7.2 2	ZooKeeper 3.4.12 Kafka 2.11-1.0.1	Kafka-Manager 1.3.3.16		安全组ID: sg								
				专有网络/交换机: vpc-								
主机信息	C	主实例组 🕿										
主实例组(MASTER)	按量付费	ECS ID	组件部署状态	5. 公网	内	网	创建时间					
主机数量:1	CPU: 4核 数据盘配置: 高效云盘	i- 🕜 🖓] ●正常	47.	19	2.168.1.155	2019年5月27日	∃ 11:48:38				
内存:16GB	60GB * 4块								#1条			
核心实例组(CORE)	按量付费								1			

· 创建MaxCompute项目

开通MaxCompute服务并创建好项目,本文中在华东1(杭州)区域创建项 目bigdata_DOC,同时启动DataWorks相关服务,如下图所示。详情请参见开 通MaxCompute。

	概览 项目列表	调度资源列表
G DataWorks	数据集成・数据开发・MaxCompute	
快速入口		
数据开发	数据集成	运维中心
项目		全部项目
bigdata_DOC 华东1	MaxCompute_DOC 华东2	PAltest 华东2
创建时间:2018-09-02 10:26:59 计算引擎:MaxCompute 服务模块数据开发数据集成数据管理运维中心	创建时间:2018-07-19 09:12:37 计算引擎:MaxCompute 服务模块数据开发数据集成数据管理运维中心	创建时间:2018-05-23 13:32-29 计算引擎: MaxCompute PAI计算引擎 服务模块数据开发数据集成数据管理 运维中心
项目配置 进入数据开发 进入数据集成	项目配置 进入数据开发 进入数据集成	项目配置 进入数据开发 进入数据集成
常用功能 译 创建项目 × 一键CDN		

背景信息

Kafka是一款分布式发布与订阅的消息中间件,具有高性能、高吞量的特点被广泛使用,每秒能处理上百万的消息。Kafka适用于流式数据处理,主要应用于用户行为跟踪、日志收集等场景。

一个典型的Kafka集群包含若干个生产者(Producer)、Broker、消费者(Consumer)以及一 个Zookeeper集群。Kafka集群通过Zookeeper管理自身集群的配置并进行服务协同。

Topic是Kafka集群上最常用的消息的集合,是一个消息存储逻辑概念。物理磁盘不存储Topic ,而是将Topic中具体的消息按分区(Partition)存储在集群中各个节点的磁盘上。每个Topic可 以有多个生产者向它发送消息,也可以有多个消费者向它拉取(消费)消息。

每个消息被添加到分区时,会分配一个offset(偏移量,从0开始编号),是消息在一个分区中的唯 一编号。

操作步骤

1. 准备测试表与数据

a) Kafka集群创建测试数据

为保证您可以顺利登录EMR集群Header主机及MaxCompute和DataWorks可以顺利 和EMR集群Header主机通信,请您首先配置EMR集群Header主机安全组,放行TCP 22及TCP 9092端口。

A. 登录EMR集群Header主机地址

进入EMR Hadoop控制台集群管理 > 主机列表页面,确认EMR集群Header主机地 址,并通过SSH连接远程登录。

首页 > 集群管理 >	集群(、 ・ 主机列表				
主机列表					
ECS实例ID	主机名	内网IP	外网IP		查询
主机名	ECS ID	IP信息	角色 🍸	所属机器组	付费类型
emr-worker-2	i- f Ca	内网:192.168.1.157	CORE	CORE	按量付费
emr-worker-1	i-	内网:192.168.1.156	CORE	CORE	按量付费
emr-header-1	i	内网:192.168.1.155 外网:47.	MASTER	MASTER	按量付费

B. 创建测试Topic

执行如下命令创建测试所使用的Topic testkafka。

```
[root@emr-header-1 ~]# kafka-topics.sh --zookeeper emr-header
-1:2181/kafka-1.0.1 --partitions 10 --replication-factor 3 --
topic testkafka --create
Created topic "testkafka".
```

执行如下命令查看已创建的Topic。

```
[root@emr-header-1 ~]# kafka-topics.sh --list --zookeeper emr-
header-1:2181/kafka-1.0.1
__consumer_offsets
_emr-client-metrics
_schemas
connect-configs
connect-offsets
connect-status
```

testkafka

C. 写入测试数据

```
您可以执行如下命令,模拟生产者向Topic testkafka中写入数据。由于Kafka用于处理
流式数据,您可以持续不断的向其中写入数据。为保证测试结果,建议您写入10条以上的
数据。
```

```
[root@emr-header-1 ~]# kafka-console-producer.sh --broker-list
emr-header-1:9092 --topic testkafka
>123
>abc
>
```

为验证写入数据生效,您可以同时再打开一个SSH窗口,执行如下命令,模拟消费者验证 数据是否已成功写入Kafka。当数据写入成功时,您可以看到已写入的数据。

```
[root@emr-header-1 ~]# kafka-console-consumer.sh --bootstrap-
server emr-header-1:9092 --topic testkafka --from-beginning
123
abc
```

b) 创建MaxCompute表

为保证MaxCompute可以顺利接收Kafka数据,请您首先在MaxCompute上创建表。本例 中为测试便利,使用非分区表。

A. 登录DataWorks创建表,详情请参见表管理。

X DataStudio	data DOC data_DOC	~							c
数据开发 2 园 🗗	С Ф Ф	itestkafka :	×	10.00					
Q 文件名称/创建人	∑	DDL模式	从生产环境加载	提交到生产环t					
> 解决方案									
▼ 业务流程				表名	testkafka				
🗸 轟 123			写》	入该表的业务流程					
> ≓ 数据集成		基本屋性							
> 🕢 数据开发									
▼ 🔳 表			中文名	testkafka					
			一级主题	请选择		二级主题	请选择	新建主题	C
particular production of the second se			描述						
🧮 testkafka 🛛	ops.bigdata_DC <								
▶ 💋 资源									
> ∱≥ 函数		物理模型设计							
> 🔚 算法			分区类型	● 分区表 🌘	非分区表	生命周期			
> S scione database da	atav test		层级	请选择		物理分类	请选择	新建层级	C
> A works	atuv_test		表举型						
 ▼ 旧版工作流									
▶ 🖿 任务开发		表结构设计							
		添加字段	上移 下移						

您可以单击DDL模式进行建表,建表语句如下。

```
CREATE TABLE `testkafka` (
`key` string,
`value` string,
`partition1` string,
`timestamp1` string,
```

```
`offset` string,
`t123` string,
`event_id` string,
`tag` string
);
```

其中的每一列,对应于DataWorks数据集成Kafka Reader的默认列,您可以自主命名。 详情请参见配置Kafka Reader:

A.__key__表示消息的key。

B. __value__表示消息的完整内容。

C.__partition__表示当前消息所在分区。

D.__headers__表示当前消息headers信息。

E.__offset__表示当前消息的偏移量。

F. __timestamp__表示当前消息的时间戳。

2. 数据同步

a) 新建自定义资源组

由于当前DataWorks的默认资源组无法完美支持Kafka插件,您需要使用自定义资源组完成 数据同步。自定义资源组详情请参见新增任务资源。

在本文中,为节省资源,直接使用EMR集群Header主机作为自定义资源组。完成后,请等 待服务器状态变为可用。

3	Co 数据集成	biqdata DOC bigdata_DOC	~				
		资源组管理 輸入調度溶液系					
Ŧ			管理资源组 - kafka2mc				×
4		资源组名称				周閉	f 添加服务器
Ŧ		默认资源组	服务器名称/ECS UUID	服务器IP	当前状态	已使用DMU	操作
ዯ	数据源	hadoop_to_odps	4676F6E1-644C-4D	47.	可用	0	修改删除
¢							
~		test_h					
1	批量上云	kafka2mc					

b) 新建并运行同步任务

A. 在您的业务流程中右键单击数据集成,选择新建数据集成节点 > 数据同步。



B. 新建数据同步节点后,您需要选择数据来源的数据源为Kafka,数据去向的数据源 为ODPS,并且使用默认数据源odps_first。选择数据去向表为您新建的testkafka。完 成上述配置后,请单击下图框中的按钮,转换为脚本模式。

Di kafka2mc 🔹					
01 选择数据源	数据来源		数据去向		
	任这里配直数据的未就属相与入属;可以定款。	从时数据游,也可以走您 回 建时日午	政績隊已有文持的奴括未認。		
* 数据源 Kafka	→ 送择掘源库 → ?	* 数据源	odps ~	odps_first	~ ?
			testkafka		~
此数据源不支持向导模式,需要依 点击转换为脚本	e用脚本模式配置同步任务,	分区信息	无分区信息		
		清理规则	写入前清理已有数据 (Insert	Overwrite)	
		空字符串作为null	● 是 ● 否		

C. 脚本配置如下,代码释义请参见配置Kafka Reader。

```
{
       "type": "job",
      "steps": [
             {
                    "stepType": "kafka",
"parameter": {
    "server": "47.xxx.xxx.xxx:9092",
                           "kafkaConfig": {
    "group.id": "console-consumer-83505"
                           },
"valueType": "ByteArray",
                           "column": [
"__key__",
"__value__",
"__partition__",
"__timestamp__",
"__offset__",
                                 "'123'",
"event_id",
"tag.desc"
                          ],
"topic": "testkafka",
"keyType": "ByteArray",
"waitTime": "10",
"beginOffset": "0",
"...doffset": "3"
                    },
"name": "Reader",
"read
                    "category": "reader"
             },
{
                    "stepType": "odps",
                    "parameter": {
                           "partition": "",
                           "truncate": true,
                           "compress": false,
"datasource": "odps_first",
                           "column": [
                                 "key",
                                  "value",
                                 "partition1",
                                  "timestamp1",
                                  "offset",
                                  "t123",
"event_id",
                                  "tag"
                           ],
```

```
"emptyAsNull": false,
                "table": "testkafka"
            },
            "name": "Writer".
            "category": "writer"
        }
    ],
    "version": "2.0",
    "order": {
        "hops": [
            {
                "from": "Reader",
                "to": "Writer"
            }
        ]
    "errorLimit": {
            "record": ""
        },
        "speed": {
            "throttle": false,
            "concurrent": 1,
        }
    }
}
```

您可以通过在Header主机上使用kafka-consumer-groups.sh --bootstrap

-server emr-header-1:9092 --list命令查看group.id参数,及消费者

的Group名称。

```
[root@emr-header-1 ~]# kafka-consumer-groups.sh --bootstrap-
server emr-header-1:9092 --list
Note: This will not show information about old Zookeeper-based
consumers.
_emr-client-metrics-handler-group
console-consumer-69493
console-consumer-83505
console-consumer-21030
console-consumer-45322
console-consumer-45322
console-consumer-14773
以console-consumer-83505为例,您可以根据该参数在Header主机上使用
kafka-consumer-groups.sh --bootstrap-server emr-header-1:9092
--describe --group console-consumer-83505命令确认beginOffset及
endOffset参数。
```

```
[root@emr-header-1 ~]# kafka-consumer-groups.sh --bootstrap-
server emr-header-1:9092 --describe --group console-consumer-
83505
Note: This will not show information about old Zookeeper-based
consumers.
Consumer group 'console-consumer-83505' has no active members.
TOPIC PARTITION CURRENT-OFFSET LOG-
END-OFFSET LAG CONSUMER-ID
HOST CLIENT-ID
```

testkafka	Θ	_	6	Θ	0
- test	Ŭ		6	- 3	3
-	Θ	-		_	
testkafka	Θ	-	Θ	Θ	Θ
- testkafka	Θ	_	1	1	1
- testkafka	Ũ		5	- 0	0
-	Θ	-		-	

完成脚本配置后,请首先切换任务资源组为您刚创建的资源组,然后单击运行。

Di kafka	2mc	• (Di test		E	test	kafka	1	test12	34	Sq test	Di test123	📄 数	据集	咸			<	> ≡
	\odot	Þ	ſ	۵.			F	23											运维↗
1 2	2	"type	": "jo	b",									档 >	×			配置任务资源组	帮助文档	调度
3		"step	s": [聖
4															任务资源组	默认资源组			–
5			"s1	серТур	be ": "	kafka													
6			"pa	aramet	er": -	{							1						版本
7				"ser	rver":	"loc	alhos	st:9093	",							kafka2mc			**

D. 完成运行后,您可以在运行日志中查看运行结果,如下为成功运行的日志。

运行日志	
读写失败总数	
2019-05-29 11:10:	24 INFO
2019-05-29 11:10:	24 INFO Exit code of the Shell command 0
2019-05-29 11:10:	24 INFO Invocation of Shell command completed
2019-05-29 11:10:	24 INFO Shell run successfully!
2019-05-29 11:10:	24 INFO Current task status: FINISH
2019-05-29 11:10:	24 INFO Cost time is: 114.983s
/home/admin/alisa	ntasknode/taskinfo//20190529/diide/11/08/28/mjuicuxu5slfbv3xu7m8csqy/T3_0242504015.log-END-EOF
Exit with SUCCESS	
2019-05-29 11:10:	28 [INFO] Sandbox context cleanup temp file success.
2019-05-29 11:10:	28 [INFO] Data synchronization ended with return code: [0].
2019-05-29 11:10:	28 INFO
2019-05-29 11:10:	28 INFO Exit code of the Shell command 0

3. 结果验证

您可以通过新建一个数据开发任务运行SQL语句,查看当前表中是否已存在从Kafka同步过来的数据。本例中使用select * from testkafka;语句,完成后单击运行即可。



执行结果如下,本例中为保证结果,在testkafka Topic中输入了多条数据,您可以查验是否和 您输入的数据一致。

e	5 select * from	ı testkafka;						
逆	衍日志	吉果22 ×						
	A	В	С	D	E	F	G	н
1	key 🗸	value 🗸 🗸	partition1 🗸 🗸	timestamp1 🗸 🗸	offset 🗸 🗸	t123 🗸	event_id 🗸 🗸	tag 🗸 🗸
2	\N	123	3	1559100458698	0	123	\N	\N
3	N/N	234	9	1559100458028	0	123	\N	\N
4	\N	567	0	1559100466891	0	123	\N	\N
5	\N	123	7	1559050808437	0	123	\N	N/
6	N/	567	1	1559100457401	1	123	\N	N/

1.3 JSON数据从OSS迁移至MaxCompute

本文将为您介绍如何通过DataWorks的数据集成功能,将JSON数据从OSS迁移

至MaxCompute,并使用MaxCompute內置字符串函数GET_JSON_OBJECT提取JSON信息的 最佳实践。

准备工作

开始将JSON数据从OSS迁移至MaxCompute的操作前,您需要首先将JSON文件重命名为后缀是 txt的文件,并上传至OSS。

本文中使用的JSON文件为applog.txt,将其上传至OSS,本文中OSS Bucket位于华东2区。

对象存储	docgood2	读写权限 公共读写 🛆 类型 标准存储	区域 华东 2 创建时间 2018-11-01 16:42 删除 Bucket
概览	概览 文件管理 基础设置 域名管理 图片处理 事件通知 函数计算	基础数据 热点统计 API统计 文件访问统计	
存储空间 + O J I II Bucket 名称 Q	上作文件 新建目录 神片管理 授权 批品操作 メ 服新		⑦ 通过 SDK 管理文件 输入文件名前缀匹配 Q
 caffe-test002 	文件名(Object Name)	文件大小 存储类型	更新时间 操作
 docgood 	demo/		删除
docgood2 emr-demo	applog.txt	0.85KB 标准存储	2018-11-13 13:38 預览 更多 ∨
 intelligent-speech-i 	userlog1.txt	0.033KB 标准存储	2018-11-12 21:21 預览 更多 🗸
 iinabucket001 			
{ "sto }, "exp	<pre>pre": { "book": [{ "category": "reference "author": "Nigel Rees "title": "Sayings of " "price": 8.95 }, { "category": "fiction" "author": "Evelyn Waug "title": "Sword of Hom "price": 12.99 }, { "category": "fiction "author": "J. R. R. " "title": "The Lord o "isbn": "0-395-19395 "price": 22.99 }], "bicycle": { "color": "red", "price": 19.95 } pensive": 10 </pre>	e", ", the Century", gh", nour", ", Tolkien", f the Rings", -8",	

}

通过DataWorks将JSON数据从OSS迁移至MaxCompute

- 1. 新增OSS数据源
 - a. 以项目管理员身份进入DataWorks控制台,单击对应工作空间操作栏中的进入数据集成。
 - b. 选择同步资源管理 > 数据源, 单击新增数据源。

⑤ Co 数据集成		■ ~									ಲ್ಕೆ 👳	-
= ▼ 任务列表	数据源类型: 全部	~	新增数据源				>	< [2月新 多库多	表搬迁 批量	2 III	en e
意线同步任务	_		关系型数据库					*环境配置信息				
→ 同步资源管理	数据源名称	数据源类型 链接信用	<i>A</i>	*	(F)	OPACLE		连遷状态	连通时间	适用环境	操作	选择
↑ 数据源		Endpoin 项目名#	MySQL: MySQL	SQL Server	PostgreSQL PostgreSQL	Oracle	DM			开发		
**** ***	odps_tirst	ODPS Endpoin 顶目名#	0	~	- Å-•					生产		
🚀 批量上云			00	₩	°‡Xo	ŝ						
	rds workshop log	数据库结 实例名: Useman MySOI	DRDS	POLARDB	HybridDB for MySQL	AnalyticDB for PostgreSQL		成功	2019/07/04 16:04:32	开发	整车迁移批量配置 编辑 删除	
		数据库结 实例名: Useman	\mathbf{v}	×	\bigcirc	47	\bigcirc			生 ^{pe}	编辑 删除	
		Access Bucket Endpoin	MaxCompute (ODPS)	DetaHub	AnalyticDB (ADS)	Lightning	Data Lake Analytics(DLA)	成功	2019/07/04 16:04:49	开发	编辑 删除	
	oss_workshop_log	USS Access Bucket Endpoin	半結ね化存储 のSS	HDFS	FTP					生产	编辑 删除	

- c. 在新增数据源弹出框中,选择数据源类型为OSS。
- d. 填写OSS数据源的各配置项。

新增OSS数据源		×
* 数据源名称:	OSS	
数据源描述:	OSS数据源	
* 适用环境:	✔ 开发 生产	
* Endpoint :	http://	?
* Bucket :		?
* AccessKey ID :		?
* AccessKey Secret :		
测试连通性:	测试连通性	
	上一步	完成

配置	说明
数据源名称	数据源名称必须以字母、数字、下划线组合,且不能以数字和 下划线开头。
数据源描述	对数据源进行简单描述,不得超过80个字符。
适用环境	可以选择开发或生产环境。
	说明:(Q标准模式工作空间会显示此配置。)
Endpoint	OSS Endpoint信息,本示例为http://oss-cn-shanghai. aliyuncs.com或http://oss-cn-shanghai-internal. aliyuncs.com。OSS各地域的外网、内网地址请参见OSS开 通Region和Endpoint对照表。
	 说明: 由于本文中OSS和DataWorks项目处于同一个区域中,所以 本文选用后者,通过内网连接。
Bucket	相应的OSS Bucket信息,指存储空间,是用于存储对象的容 器。
	您可以创建一个或多个存储空间,每个存储空间可添加一个或
	多个文件。
	您可以在数据同步任务中查找此处填写的存储空间中相应的文
	件,没有添加的存储空间,则不能查找其中的文件。
AccessKey ID/ AceessKey Secret	访问秘匙(AccessKeyID和AccessKeySecret),相当于登录 密码。

e. 单击测试连通性。

f. 测试连通性通过后, 单击完成。

2. 新建数据同步任务

在DataWorks上新建数据同步节点,详情请参见数据同步节点。

新建节点		×
节点类型:	数据同步	
节点名称:	数据同步	
目标文件夹:	请选择	
		提交取消

新建的同时,在DataWorks新建一个建表任务,用于存放JSON数据,本示例新建表名为mqdata。

▼ 业务流程	21	"fileFormat": "binary",	
✓ ▲ test			
∨ ☴ 数据集成			
• Di MQ2MaxCompute 我锁定 1			
• Di test22 我锁定 11-09 15:15		"categor 新建表	×
• Di test223 表演定 11-09 15:20			
>		"stepTyp	
		"paramet 数据库类型: • MaxCompute	
▶ ■ 衣		"par	
🗰 mqdata odps.MaxCompute_		"isC 表名: mqdata	
2010 次泊		"tru	
> 🖉 瓦砾		"dat	
> 🗾 函数			
、 == 笛注			
※ N 3年14			

表参数可以通过图形化界面完成。本例中mqdata表仅有一列,类型为string,列名为MQ data。
基本属性								
中文名	G: MQ 数据存放							
一级主题	回: 请选择		二级主题:	请选择		新建主题	С	
描述								
物理模型设计								
分区类型	型: 🔵 分区表 📀 非分	区表	生命周期:					
层级	3: 请选择		物理分类:	请选择		新建层级	С	
表类型	2: 💿 内部表 🔵 外部	 а						
表结构设计								
添加字段 上移 「	下移							
字段英文名	字段中文名	字段类型		长度/设置	描述		主键 ⑦	操作
MQdata N	MQ数据	string		string			否	

3. 配置同步任务参数

完成上述新建后,您可以在图形化界面配置数据同步任务参数,如下图所示。选择目标数据源名称为odps_first,选择目标表为刚建立的mqdata。数据来源类型为OSS,Object前缀可填写文件路径及名称。

DI MQ2MaxCompute ×	🗰 kafka1 x 🌐 tt1 x 🧰 jd	x Di test223 x Di test22	x 🔄 abc x 全部解决方案 x 嚞 test	× 📰 🕏 <
	<u>۲</u> ه 🗊 🗄 ال			
01 选择数据源	数据来源		数据去向	
	在这里配置数据的来源端和写入端;	可以是默认的数据源,也可以是您创建的自得	有数据源查看支持的数据来源类型	
* 数据源:	OSS v OSS_userlog v	(?) * 数据源:	0DPS v odps_first v	?
* Object前缀:	applog.txt	*表:	mqdata ~	
	添加Object +		一键生成目标表	
* 文本类型:	text ~	分区信息:	无分区信息	
< * 列 分隔符 :		清理规则:	写入前清理已有数据 (Insert Overwrite) ~	
编码格式:	UTF-8	压缩:	• 不压缩 ○ 压缩	
null值:		空字符串作为null:	○是 • 否	
* 压缩格式:	None			
* 是否包含表头:	No ~			
	数据预览			

🗾 说明:

列分隔符使用TXT文件中不存在的字符即可,本文使用(^)。对于OSS中的TXT格式数据 源,Dataworks支持多字符分隔符,您可以使用(%&%#^\$\$^%)这种很难出现的字符串作 为列分隔符。

映射方式选择默认的同行映射即可。

(Di MQ2MaxComput	×	kafka1	🌐 π1		🗮 jd	Di test223	Di test22		sq abc	全部解决方案	🛔 test		;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;	
		[↑]	ل	1											ìŻ
								J 10H117/JIGH.	_ ~ (-					
	* 压缩材	乱: N	one												
	* 是否包含著	€头: N	0												
				数据预											
	02 字段映射			源头					目标表						
													同夕岫射		
			置/值	类型	Ć	?				目标表字段	类型		同行映射		
<		第	0列	string		•				mqdata	STRING		取消映射		
٦.															

单击左上方的切换脚本按钮,切换为脚本模式。修改fileFormat参数为"fileFormat":"

binary"。脚本模式代码示例如下。

```
{
     "type": "job",
     "steps": [
          {
                "stepType": "oss",
                "parameter": {
                     "fieldDelimiterOrigin": "^",
                     "nullFormat": "",
"compress": "",
                     "datasource": "OSS_userlog",
                     "column": [
                           {
                                "name": 0,
"type": "string",
                                "index": 0
                           }
                     ],
                     "skipHeader": "false",
"encoding": "UTF-8",
"fieldDelimiter": "^",
                     "fileFormat": "binary",
                     "object": [
                           "applog.txt"
                     ]
                },
"name": "Reader",
". "reader"
                "category": "reader"
          },
{
                "stepType": "odps",
"parameter": {
                     "partition": ""
                     "isCompress": false,
```

```
"truncate": true,
                     "datasource": "odps_first",
                     "column": [
"mqdata"
                     ],
                     "emptyAsNull": false,
                     "table": "mqdata"
                },
               "name": "Writer",
"category": "writer"
          }
     ],
"version": "2.0",
     {
                     "from": "Reader",
                     "to": "Writer"
                }
          ]
    },
"setting": {
    "errorLimit": {
        "record": ""
          },
"speed": {
    "speed": {
               "concurrent": 2,
                "throttle": false,
          }
     }
}
```

▋ 说明:

该步骤可以保证OSS中的JSON文件同步到MaxCompute之后存在同一行数据中,即为一个字

段,其他参数保持不变。

完成上述配置后,单击运行即可。运行成功日志示例如下所示。

运行日志
2018-11-13 16:58:08 : com.alibaba.cdp.sdk.exception.CDPException: RequestId[075ba938-7d6c-471a-9286-8d864b135e6b] Enror: Run intance encounter problems, reason:
Exit with SUCCESS.
2018-11-13 16:58:08 [INFO] Sandbox context cleanup temp file success.
2018-11-13 16:58:08 [INFO] Data synchronization ended with return code: [0].
2018-11-13 16:58:08 INFO ====================================
2018-11-13 16:58:08 INFO Exit code of the Shell command 0
2018-11-13 16:58:08 INFO Invocation of Shell command completed
2018-11-13 16:58:08 INFO Shell run successfully!
2018-11-13 16:58:08 INFO Current task status: FINISH
2018-11-13 16:58:08 INFO Cost time is: 43.248s
/home/admin/alisatasknode/taskinfo//20181113/datastudio/16/57/23/uv7deija7u8j4wyhzm82sgsr/T3_0616594516.log-END-EOF

JSON数据从OSS迁移至MaxCompute结果验证

1. 在您的业务流程中新建一个ODPS SQL节点。



2. 查看当前mqdata表中数据, 输入SELECT * from mqdata;语句。

ક્વ JSC)Ndata	×	Sq mqd	lata	×	Sq test		× (DI MQ2MaxComput	×	🛱 kafka1	×	π1	×	🗮 jd	:	×	DI test223	×	Di test22	× [
	Þ	ᡗ	ե		\odot	:	\$														
1 2 3																					
4 5					11-13 ****																
6			from n	ıqdata	a;																
																				不	
																				K 3	
运行	旧志		结果	41]	×																
1	ndata										А										
2 {	"store	:{	"book": [{	"cate	gory": "re	ferenc	e", "author":	'Nige	el Rees",	'title":	"Sayings of the	e Century	/", "F	orice": 8.	95	}, {		"category": "f	iction",



3. 确认导入表中的数据结果无误后,使用SELECT GET_JSON_OBJECT(mqdata.MQdata,'\$.

expensive') FROM mqdata;获取JSON文件中的expensive值。

Sq JSC)Ndata	\bullet	Sq mqda	ata	× [Sq test	×	<	Di MQ	2MaxCo	ompute	×	Ħ	kafkaʻ	1	×	Ħ	tt1
	Þ	[↑]	ե		lacksquare	:	\$											
1 2 3 4 5	od ** au cr	ps s **** thor eate ****	ql ******* :dtplus time:2 ******	***** _docs 018-11 *****	***** 1-13 ****	****** 18:56: *****	***** 45 *****											
6	SELE	CT G	ET_JSON	_OBJE(CT(mo	data.M	Qdata,	'\$.expen	nsive') FRO	Mr	nqda	ta;				
运行	日志		结果	[1]	×	结果	2] ×	۲										
1 _(2 1)	c0 0	A																

更多信息

在进行迁移后结果验证时,您可以使用MaxCompute内建字符串函数GET_JSON_OBJECT获取您 想要的JSON数据。

1.4 JSON数据从MongoDB迁移至MaxCompute

本文将为您介绍如何通过DataWorks的数据集成功能,将从MongoDB提取的JSON字段迁移 至MaxCompute。

准备工作

1. 账号准备

在数据库内新建用户,用于DataWorks添加数据源。本示例执行如下命令。

db.createUser({user:"bookuser",pwd:"123456",roles:["root"]})

新建用户名为bookuser, 密码为123456, 权限为root。

2. 数据准备

首先您需要将数据上传至您的MongoDB数据库。本示例使用阿里云的云数据 库MongoDB版,网络类型为VPC(需申请公网地址,否则无法与DataWorks默认资源组互 通),测试数据如下。

```
{
                                  "store": {
                                  "book": [
                                  {
                                  "category": "reference",
"author": "Nigel Rees",
"title": "Sayings of the Century",
                                  "price": 8.95
                                  },
                                  {
                                  "category": "fiction",
"author": "Evelyn Waugh",
"title": "Sword of Honour",
                                  "price": 12.99
                                  },
                                  Ł
                                  "price": 22.99
                                  }
                                  ],
                                  "bicycle": {
                                  "color": "red",
"price": 19.95
                                  }
                                  },
                                  "expensive": 10
```

}

登录MongoDB的DMS控制台,本示例使用的数据库为admin,集合为userlog。您可以在查询窗口执行如下命令,查看已上传的数据。

```
db.userlog.find().limit(10)
```



通过DataWorks将JSON数据从MongoDB迁移至MaxCompute

- 1. 新增MongoDB数据源
 - a. 以项目管理员身份进入DataWorks控制台,单击对应工作空间操作栏中的进入数据集成。
 - b. 选择同步资源管理 > 数据源,单击新增数据源。

=								_					
▼ 任务列表	数据源类型: 全部		新增数据源					×		3 刷新 多声多	表搬迁批	副新塔数据 2 新	增数据源
🖕 离线同步任务			大数据存储										
↓ 同步资源管理	数据源名称	数据源类型 链	爱信 意	el.	\sim	12	9		连通状态	连通时间	适用环境	操作	选择
▲ 数据源		En 顶	dpoin 물침위 MaxCompute (ODF	PS) DataHub	AnalyticDB (ADS)	Lightning	Data Lake				开发		
新新祖	odps_first	ODPS En	dpoin 日本名 半结构化存储				Analytics(ULA)				开车		
✓ 批母上云		iy.		d a							Ð		
		数 实 Us MySOI	照库2 列名: eman OSS	HDFS	FTP				成功	2019/07/04 16:04:32	开发	整库迁移批量配置 编辑 删除	
		anijout 數 实 Us	居库? 列名: eman NoSQL								生产	编辑 删除	
		Ac Bu En	cket :	3 🗳	8	Table Store (OTS)			成功	2019/07/04 16:04:49	开发	编辑 删除	
	oss_workshop_log	Ac Bu En	ccess) cket : dpoin		The state	1886 51676 (515)					生产	编辑 删除	
			LogHub										
							Ę	UHI I					

- c. 在新增数据源弹出框中,选择数据源类型为MongoDB。
- d. 填写MongoDB数据源的各配置项。

新增MongoDB数据源		×
* 数据源类型:	连接串模式(数据集成网络可直接连通)	
* 数据源名称:	自定义名称	
数据源描述:		
* 适用环境:	✔ 开发 生产	
* 访问地址:	host:port	
	添加访问地址	
* 数据库名 :	请输入MongoDB集合名称	
* 用户名:		
* 密码 :		
测试连通性:	测试连通性	
0	如果您使用的是云数据库MongoDB版 出于安全策略的考虑,数据集成仅支持使用MongoDB数据库对应账号进行连接 请避免使用root作为访问账号	
	上一步	完成

配置	说明
数据源类型	由于本文中MongoDB处于VPC环境下,因 此数据源类型需选择连接串模式(数据集成 网络可直接连通)。
数据源名称	数据源名称必须以字母、数字、下划线组 合,且不能以数字和下划线开头。
数据源描述	对数据源进行简单描述,不得超过80个字 符。
适用环境	可以选择开发或生产环境。
	〕 说明: 仅标准模式工作空间会显示此配置。

配置	说明					
访问地址	访问地址及端口号可以通过单 击MongoDB控制台中的实例名称	你获取。				
	Col: Col: <td< th=""><th>20. III.46</th></td<>	20. III.46				
数据库名	该数据源对应的数据库名称。					
用户名/密码	数据库对应的用户名和密码。					

e. 单击测试连通性。

f. 测试连通性通过后, 单击完成。

2. 新建数据同步任务

在DataWorks上新建数据同步节点,详情请参见数据同步节点。

新建节点		×
节点类型:	数据同步	
节点名称:	数据同步	
目标文件夹:	请选择	
		取消

新建的同时,在DataWorks新建一个建表任务,用于存放JSON数据,本示例新建表名为mqdata。

▼ 业务流移	4	21	"fileFormat": "binary",	
· · • •				
~ .	≓ 数据集成			
	• Di MQ2MaxCompute 我锁定 1			
			"categor 新建表	×
> <	(1) 数据开发		"stepTyp	
× 1	■ 表		"paramet 或順序突至。 MaxCompute	-
	🗰 mqdata odps.MaxCompute_D		"isC 表名: mqdata	
> 🤞	资源		"dat	TRAK
> -	f* 函数	运行日志		
> 🚦	算法			

表参数可以通过图形化界面完成。本例中mqdata表仅有一列,类型为string,列名为MQ data。

基本属性									
中	文名: MQ 数据存放								
—级:	主题: 请选择		二级主题:	请选择		新建主题	C		
ł	描述:								
物理模型设计 									
分区	类型: 🔵 分区表 💿 非分	区表	生命周期:						
J	层级: 请选择		物理分类:	请选择		新建层级	С		
表	类型: 💿 内部表 🔵 外部	诔							
表结构设计									
添加字段 上移									
字段英文名	字段中文名	字段类型		长度/设置	描述		主鍵 ⑦	操作	
MQdata	MQ数据	string		string			否		

3. 配置同步任务参数

完成上述新建后,您可以在图形化界面配置数据同步任务参数,如下图所示。选择数据来源类型为MongoDB,来源表为mongodb_userlog。目标数据源名称为odps_first,目标表为刚新建的mqdata。

谢	Þ	♪	٤		-	⟨♪								
电译数据	源					数据来源					数据去向			
					在	这里配置数据的来源端和	写入端;	可以是默	认的数据源 , 也可以是	您创建的自有	有数据源查看支持的数据来源类	型		
													_	
	*数据》	泉: M	longoDB			mongodb_userlog		(?)		*数据源:	ODPS ~	odps_first	?)	
			-1	+						*表:	mqdata			
(此 数据 点击转	源不文 换为脚		む, 斋	要使用腳	₩个模式配直向步任务,				分区信息:	无分区信息			
										清理规则:	写入前清理已有数据 (Insert C	verwrite)		
										压缩:	• 不压缩 ○ 压缩			
									空字符	事作为null:	○是 • 否			

由于MongoDB数据源不支持向导模式开发,您直接点击转换为脚本,即可跳转至脚本模式进行 配置。

```
"type": "document.String" //非一层子属性以最终获取的类型
为准。假如您选取的JSON字段为一级字段,如本例中的expensive,则直接填写string即
可。
                 }
               ],
             "collectionName //集合名称": "userlog"
             },
        "name": "Reader",
        "category": "reader"
        },
        {
             "stepType": "odps",
            "parameter": {
"partition": ""
             "isCompress": false,
             "truncate": true,
             "datasource": "odps_first",
             "column": [
"mqdata" //MaxCompute表列名。
             ],
             "emptyAsNull": false,
             "table": "mqdata"
             },
             "name": "Writer",
             "category": "writer"
             }
             ],
             "version": "2.0",
             "order": {
"hops": [
             "from": "Reader",
             "to": "Writer"
             }
             ]
             },
             "setting": {
             "errorLimit": {
             "record": ""
            },
"speed": {
    current
             "concurrent": 2,
            "throttle": false,
             }
             }
        }
```

完成上述配置后,单击运行即可。运行成功日志示例如下所示。



JSON数据从MongoDB迁移至MaxCompute结果验证

1. 在您的业务流程中新建一个ODPS SQL节点。



2. 输入SELECT * from mqdata;语句, 查看当前mqdata表中数据。



UNDERSECTION DEFINITION OF STREET, UNDERSECTION DEFINITION OF STREET, UNDERSECTION OF STREET, UNDERSE

1.5 Elasticsearch数据迁移至MaxCompute

本文将为您介绍如何通过DataWorks数据同步功能,将阿里云Elasticsearch集群上的数据迁移 至MaxCompute。

前提条件

· 搭建阿里云Elasticsearch集群

进行数据迁移前,您需要保证自己的阿里云Elasticsearch集群环境正常。搭建阿里 云Elasticsearch集群的详细过程请参见Elasticsearch快速入门。

本示例中阿里云Elasticsearch的具体配置如下:

- 地域: 华东2(上海)
- 可用区:上海可用区B
- 版本: 5.5.3 with Commercial Feature

	预付费	后付费							
	地域	华东1 (杭州)	华北2 (北京)	华东2 (上海)	华南1 (深圳)	印度 (孟买)	新加坡	当前配置	
調査		香港 印度尼西亚(推加达)	美国 (硅谷) 华北1 (青岛)	马来西亚 (吉隆坡) 华北3 (张家口)	德国 (法兰克福)	日本 (东京)	澳大利亚 (悉尼)	地域: 可用区: 版本:	华东2(上海) 上海可用区B 5.5.3 with Commercial Feature
	可用区	上海可用区B	*					网络类型: 可用区数量: 专有网络:	专有网络 单可用区
2018日	资源组	全部	• 1	机资源组	•			虚拟交换机: 规格族:	云盘型
	版本	5.5.3 with Commercial Feature	6.3 with Commercia Feature	6.7 with Commercial Feature				实例规格: 数量: 专有主节点: 协调节点:	1板2G 3 否 否
	网络类型	专有网络						冷数据节点: 存储类型: 单节点存储空间:	否 SSD云盘 20

· 创建MaxCompute项目

开通MaxCompute服务并创建好项目,详情请参见开通MaxCompute。本示例中在华东1(杭州)区域创建项目bigdata_DOC,同时启动DataWorks相关服务。

背景信息

Elasticsearch是一个基于Lucene的搜索服务器,它提供了一个分布式多用户功能的全文搜索引擎。Elasticsearch是遵从Apache开源条款的一款开源产品,是当前主流的企业级搜索引擎。

阿里云Elasticsearch提供Elasticsearch 5.5.3 with Commercial Feature、6.3.2 with Commercial Feature、6.7.0 with Commercial Feature及商业插件X-pack服务,致力于数据 分析、数据搜索等场景服务。在开源Elasticsearch基础上提供企业级权限管控、安全监控告警、 自动报表生成等功能。

操作步骤

1. 创建测试表

 kibana
 Conside Search Profiler Grok Debugger

 Oncode Search Profiler Grok Debugger

 interent

 i

将阿里云上的数据导入至阿里云Elasticsearch(离线),具体操作请参见云上数据导入。

2. 创建接收表

为保证MaxCompute可以顺利接收Elasticsearch数据,您需要在MaxCompute上创建表。 本示例使用非分区表。

a) 登录DataWorks控制台进行创建表的操作,详情请参见表管理。

💸 DataSt	tudio	biqdata DOC bigdata_DOC					
数据开发	윤 🗟	മേ⊕ക	elastic2mc_bankdata >				
Q 文件名称	称/创建人	V.	DDL模式 从生产	不境加载 二提			
> 🧭	ý 控制						
> 🚣 ci	lone_databa	ise_datav_test		表名	elastic2mc_bankdata		
🗸 🛃 el	lasticsearch	_reader		该表的业务流程	elasticsearch_reader		
~ 😑	数据集成						
	1 1 1 1						
-	in the			elastic2mc_b	ankdata		
> 🕡	数据开发		一级主题	请选择		二级主题 请选择	新建主题 C
~ 🔳	表		描述				
	=						
	i elastic	:2mc_bankdata o					
> 🧭	2 资源		物理模型设计				
> 🏦	函数				◎ 非分区表	生命周期	
> 📒	算法		,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,				
> 🧕	2 控制		层积			物理分类 加达样	
> 🏯 w	rorks		表类型				
→ 旧版工作	流						
			表结构设计				

b) 为表添加字段。请确保添加的字段与需要被同步的Elasticsearch表字段相对应。

elastic2mc_bankdata	×					
表结构设计						
添加字段 上移	下移					
字段英文名	字段中文名	字段类型	长度/设置	描述	主鍵 ⑦	操作
age	age	bigint				E ŧ
job	job	string				Ê
marital	marital	string				Ê
education	education	string				ÊÊ
default	default	string				Ê
housing	housing	string				ÊÊ
Ioan	loan	string				Ê
contact	contact	string				ÊÊ
month	month	string				ÊÊ
day of week	day of week	stripn			西	

3. 同步数据

a) 在业务流程中,右键单击数据集成,选择新建数据集成节点>数据同步。



b) 单击下图框中的按钮, 转换为脚本模式。

	6		
01 izifatia	数据来源	数据去向	
在这里配置数	19的来源端和写入墙;可以是默认的数1820,也可以是绘创	建的自有数据源音看支持的数据来源类型	
* 数据源 数据源类型		夏数据源类型 > 选择据源库 > (9
02 字段映射	源头表	目标表	
	▲ 请先选择数据源与表后,才会显示学和	<u>Qiçat</u>	

c) 脚本配置如下所示。代码释义请参见配置Elasticsearch Reader。

```
{
    "type": "job",
    "steps": [
         {
             "stepType": "elasticsearch",
             "parameter": {
                  "retryCount": 3,
                  "column": [
                      "age",
"job",
                      "marital",
                      "education",
                      "default",
                      "housing",
                      "loan",
                      "contact",
                      "month",
"day_of_week",
                      "duration",
                      "campaign",
                      "pdays",
                      "previous",
                      "poutcome",
                      "emp_var_rate",
                      "cons_price_idx",
                      "cons_conf_idx",
                      "euribor3m"
                      "nr_employed",
                      "y"
```

```
],
                   "scroll": "1m",
"index": "es_index",
                   "pageSize": 1,
                   "sort": {
                        "age": "asc"
},
                   "type": "elasticsearch",
                   "connTimeOut": 1000,
                   "retrySleepTime": 1000,
                   "endpoint": "http://es-cn-xxxx.xxxx.xxxx.com:
9200",
                   "password": "xxxx",
                   "search": {
                        "match_all": {}
                   },
"readTimeOut": 5000,
"..."
                   "username": "xxxx"
              },
"name": "Reader",
"'' "reader"
               "category": "reader"
         },
{
              "stepType": "odps",
               "parameter": {
                   "partition": "",
                   "truncate": true,
"compress": false,
                   "datasource": "odps_first",
                   "column": [
                        "age"
                        "job",
                        "marital",
                        "education",
                        "default",
"housing",
                        "loan",
                        "contact",
                        "month",
"day_of_week",
                        "duration",
"campaign",
                        "pdays",
                        "previous",
                        "poutcome",
                        "emp_var_rate",
                        "cons_price_idx",
                        "cons_conf_idx",
                        "euribor3m"
                        "nr_employed",
                        "v"
                   ],
"emptyAsNull": false,
"elastic2mc_"
                   "table": "elastic2mc_bankdata"
              },
"name": "Writer",
"* "writ
              "category": "writer"
         }
    ],
"version": "2.0",
     "order": {
         "hops":
                  E
              {
                   "from": "Reader",
```

```
"to": "Writer"
}
]
},
"setting": {
    "errorLimit": {
        "record": "0"
        },
        "speed": {
            "throttle": false,
            "concurrent": 1,
            "dmu": 1
        }
}
```

您可以在创建的阿里云Elasticsearch集群的基本信息中,查看公网地址和公网端口,并将这 些信息填入上面的代码中。

= (-)阿里云	账号全部资源 ▼ 华乐2 (上海) ▼ Q 搜索	费用 工单 备塞 企业 支持与服务 🖸 🛕 📜
<	es-cn-	Kibana控制台 集群监控 重启实例
基本信息	基本信息	
ES集群配置 插件配置	实例D: es-cn-	创建时间:
集群监控	名称: es-cn-编辑 Elasticsearch 版本: 5.5.3 with Commercial Feature	状态 ● 正常 付费类型: 后付费
日志查询	区域: 华东2	可用区: cn-shanghai-b
安全配置	专有网络 一种 一种 一种 一种 一种 一种	VSwitch信息:
数据备份	内网地址: es-cn-	内网端口: 9200
▼ 智能运维	公网地址: es-cncom	公网端口: 9200

d) 完成脚本配置后,单击运行。

i elast	tic2mo	:_bankda	ita	Di ela	sticsea	rch 🔵			
•	\odot	Þ					8		
1 2 3		"type "step	": "j s": [ob",				配置任务资源组	

e) 在运行日志中查看运行结果。



4. 验证结果

a) 新建一个数据开发任务运行SQL语句, 查看当前表中是否已存在从Elasticsearch同步过来的数据。



b) 执行结果如下所示。

2]	5 0	۲) (J	Ê	\$	⊙	1								运维↗
6				elast:	ic2mc_b	anko	data;								调度 配置
															血缘关系
															版本
ìž	行E	志	ŝ	吉果[1]	×									8 D	————————————————————————————————————
		А						с	[H B
1	age		~	job		v 1	narital	~	education		default	housing	🗸 loan	contact	
2				student			ingle		basic.9y			yes		cellular	
3				student			ingle		unknown			no	yes	cellular	
4				student		s	ingle		basic.9y			unknown	unknow	cellular	
5				student			ingle		unknown			yes		cellular	
6				student		s	ingle		basic.9y			yes		cellular	
7	18			student		s	ingle		high.schoo	1		yes	yes	cellular	

2 数据开发

2.1 设置调度依赖最佳实践

DataWorks V2.0中调度依赖在配置时,需要根据本节点输出名称作为关联项来给任务间设置依赖 关系。本文将为您介绍如何配置任务调度依赖的输入输出。

配置任务的本节点输入

您可以通过以下两种方式,配置本节点输入。

- · 使用代码自动解析功能,解析出任务的依赖。
- · 手动输入任务依赖(手动输入父节点的本节点的输出名称)。

	[↑]	ե		谢	Þ		С	\checkmark		22		\$							运维	
		1 ****** time:2					X 调题	夏依赖	0 -	<u> </u>									调度配置	
INSEF SELEC	атоу атоу ат * р	erwrii Frwrii	re TAB	LE pm	table	**** 2_a 2 b	日司)解析: 納上游	。 是 节点		解 (父节)	所 输入输出 点输出名称或	輸出表			使用项目	目根节点		血缘关系	
					-		父节点输出名称					父节点输出表 名		节点 名	父节点 D	责任 人	来源	操作	版本	
								1	-			-		-		-	自动解析	Û	结构	5

▋ 说明:

手动输入上游节点时,输入的是父节点的本节点的输出名称。如果父节点的任务名称和父节点的本 节点输出名称不一致,请务必正确输入本节点输出名称。

配置上游节点时,自动解析出来的上游节点,是一个无效的上游依赖。您可以通过查看解析出来的 上游依赖,在父节点ID这一列是否显示有值,来识别设置的依赖是否有效的方法。

🖳 : 🕢 🗄 🕼									运维
<pre>select \${today},\${month} from dual ;</pre>	×								调度配
	依赖的上游节点	请输入父	节点输出名	称或输出表			使用项目根节点		直
@resource_reference{"test123.py"}	父节点输 出名称	父节点输 出表名	节点	洺	父节点ID	责任人	来源	操作	血缘关系
add py test123.py;	1_root		1_ro	ot		4	手 动 添 加		版本
	本节点的输出	请输入节点	輸出名称						结构
	輸出名称		輸出表 名	下游节点 名称	后 下游节 点ID	责任 人	来源	操作	
	© -		- Ø	-	-	-	手动添加	Û	

任务依赖的配置,实质上是给两个节点设置节点间的依赖关系。只有真实存在的节点,才能够设置 有效的依赖关系,任务依赖才能设置成功。

无效的上游依赖

无效的上游依赖通常有以下两种情况:

・父节点不存在。

	<u>ት</u>] [ኔ]	⊘	Þ		C	\checkmark		23	:	\$							运	维
odps **** auth crea ****	sql ******** or te time ******	:****** :2018-: :*****	******* 10-23 (******	****** 19:50:1	***** 54 ****	× 调加 _{自动}	夏 依赖 解析:(⑦ — • 是(解析	输入输出								调度配置血
SELECT	* proje	ect_b_na	ame.pm	_table_	_a _b	依赖	的上游,	京	请输入	父节点	輸出名称或輸	出表名			使用项目	目根节点			
						ک ا	(节点输	出名称			父节点输出表 名	ŧ =	节点 名	父节点 D	责任 人	来源	操作		版本
						pi _l	roject_b b	_name.	pm_tabl	le						自动解 析			结构
						±#			F10.5										
						本节	点的输出	H i	輸入节	「点輸出:	名称			+					

・ 父节点输出不存在。

提交			×	
只能提交自己的节点,一共提交2个节点		10-1-11-10-		
节点 "客服导出用户销售数据数据同步"	 依赖的父节点输	<u>最大大政</u> 出		×
	all_fmys_order_in ofile,yishou_data 不能提交本节点	ozzouz_out nfos,yishou .all_fmys_o , 请先提交:	_data_user_ rder不存在 父节点开始	pr
		0		

无效的上游依赖通常是由于解析出来的父节点输出名称不存在。本示例中,可能是由于表 project_b_name.pm_table_b没有产出任务,或该表产出任务的本节点输出的配置不正确,导 致无法解析出来。您可以通过以下两种方案解决该问题:

- ・确认这个表是否有产出任务。
- ·确认这个表产出任务的本节点输出名称,将该本节点的输出名称手动输入到依赖的上游节点中。

	C 🛛 E 🗱	: (\$)						运维
odps sql ***********************************	く 调度依赖 ⑦ _{自动解析:} ● 是 ○ 否	解析输入输出		戦到以后, ③				调 度 配 置
INSERT OVERWRITE TABLE pm_table_a SELECT * FROM project b name.pm table b	依赖的上游节点 请输入父	节点输出名称或输出表	铭 🗸		使用项目	根节点		血 缘 关 系
;	父节点输出名称	父节点输出表 名	节点 名	父节点) D	责任 人	来源	操作	版
	project_b_name.pm_table _b					自动解 析		结
找到此表:	的产出任务,查看产	出任务的输出名	各称					14
\backslash	本节点的输出 请输入节点	輸出名称						
输出名称这里找	输出名称	輸出表名	下游节点 名称	下游节 点ID	责 任 人	来源	操作	
	test_pm_01.pm_table _a ⊘	test_pm_01.pm_t able_a				自动 解析	Ē	



说明:

手动输入上游节点时,输入的是父节点的本节点的输出名称。如果父节点的任务名称和父节点的本 节点输出名称不一致,请务必正确输入本节点输出名称。

例如,上游节点A的本节点输出名称是A1,下游节点B需要依赖A,此时应该在依赖上游节点的输 入框中输入A1,并单击右侧的加号进行添加。

配置上游依赖

如果您的表是从源库抽取出来,不存在上游,您可以通过单击使用工作空间根节点获得上游依赖。

× 调度配置							调度配
调度依赖 ⑦							Ĩ
自动解析: 🧿 是 🤇	否 解析输入输出						版本
依赖的上游节点: 请输入公	计点输出名称或输出表名	× +	使用工作空间根节点]			
父节点输出名称	父节点输出表名	节点名	父节点ID	责任人	来源	操作	
		没有数据					

配置任务的本节点输出

本节点名称、本节点输出名称和本节点输出表名共用同一个名称,可以高效配置本节点输出。

- · 能够快速的知道这个任务操作的是哪个表。
- · 能够快速知道这个任务失败后造成的影响范围有多大。
- ·使用自动解析配置任务依赖时,只要本节点输出符合三名合一的规则,自动解析的精准性能得到 大大的提升。

自动解析

自动解析:通过代码自动解析调度依赖关系。

自动解析的实现原理:代码中只能拿到表名,自动解析功能可以根据表名解析出对应的产出任务。

类型节点代码如下所示。

```
INSERT OVERWRITE TABLE pm_table_a SELECT * FROM project_b_name.
pm_table_b ;
```

解析出来的依赖关系如下。

(∱] [J]	÷ ()				23		\$						运	堆
					··· × 调	▼本∜ 度依赖	т <u>,⊾</u> 襘 ⑦ —	*#								调度配置
INSERT SELECT FROM	OVERWRIT * project	E TABLE	pm_table pm_table	2_a 2_b	自道	动解析:(颜的上游节	• 是() 香 清輸入:	解析输入输出 父节点输出名称 a	成输出表名		使	用项目根节点			血缘关系
						父节点输	出名称		父节点输出 表名	节点 名	父节点 ID	责任 人	来源	操作		版本
						project_b ble_b	_name.p	om_ta					自动解 析			结构

DataWorks会自动解析本节点,即依赖project_b_name产出pm_table_b的节点,同时本节点 最终产出 pm_table_a,因此父节点输出名称为project_b_name.pm_table_b,本节点的输 出名称为project_name.pm_table_a(本工作空间名称为test_pm_01)。

- ·如果您不想使用从代码解析到的依赖,则选择否。
- ·如果代码中有很多表是临时表:如t_开头的表为临时表。则该表不会被解析为调度依赖。可以通 过项目配置,定义以什么形式开头的表为临时表。
- ・如果代码中的一个表,既是产出表又是被引用表(被依赖表),则解析时只解析为产出表。
- ・如果代码中的一个表,被多次引用或者被多次产出,则解析时只解析一个调度依赖关系。

॑ 说明:

默认情况下,表名为t_开头的会被当成临时表,自动解析不解析临时表。如果t_开头不是临时 表,请联系自己的项目管理员到项目配置页面进行修改。

6	X 配置中心	11212	• •		
ŧ¥ŧ	三 配置中心				
	项目配置		分	区日期格式:	
□ ◆	模板管理		分[∑字段命名∶	
۲	层级管理		1	临时表前缀·	
3	项目备份恢复				
			上传表(导)	∖表)前缀∶	
				旬内容脱敏:	未开遭 数据保护 全 "模块 ,请点击前往开启

删除某表的输入输出

您在进行数据开发时,经常会用到静态表(数据通过本地文件上传到表中),这部分静态数据其实 是没有产出任务的。在配置依赖时,需要您删除静态表的输入:如果静态表不满足t_的格式,不会 被处理为临时表,此时您需要删除静态表的输入。

您在代码中右键单击表名,选择删除输入。

Sq test	_sql_01	×	Sq sql_t	ask :	× Sq	p_sql	_task	×	Sq	absda	dnas	×		测试新	业务流移	≣ ×
	Ē	<u>[</u>	[ح]		€	Þ		C	1 2	\checkmark		ć	8	:	\$	
	odp	os s	ql													
	***	kxxx Ebon	******	*****	*****	****	cakoakoakoa ko	k skrak st	colicolicolic:	*****	o ne ne ne ne ne n	AC AC AC	ak ak ak a	CACAR MEAN OF	cikojkojkojkojko	**
	au	unor aate	: time:2	018-10	-23 1	0.50.	54									
	>	**	*****	*****	*****	*****	****	kakaka)	****	****	****	***	****	*****	*****	**
	INSE	кт о	VERWRIT	E TABLI	E pm_	table	_a									
	SELE	ст														
	FROM		project	_b_nam	e.pm_	table	h									
	;						添加	俞入								
							添加	諭出								
							删除	前出								
							删除	俞入								
							转到	定义		Ct	rl+F12					
							<u> </u>	定义 		A	lt+F12					
							更改)	新有[匹配功	页 C	trl+F2					
							剪切									
							复制									

如果您是从DataWorksV1.0升级至DataWorks V2.0的用户,则您迁移过来的DataWorks任务的本节点输出默认置为项目名.节点名。

Ш	数据开发 2 🗟 🛱 Ċ 🔂 🖸	Sq absdadnas ×	🛃 测试新业务流程 🗙									≡
(7)	文件名称/创建人		l 🗴 🗇 🕑									运维
*	▼ 📄 任务开发		~									
R	> D	2 add py 3	へ									
	> `		自动解析: 〇 是 💿 否	解析输入输出								
B	>											ш
	> D		依赖的上游节点 请输	1入父节点输出名称	或輸出表名							缘关
_	Sc 07-06 17:		公共占给中安委	公共占检出主体	z 节占夕		公共占回) ±4	ı	本道	提作	*
	• 😭 11-22 17:4		又口派調出自称				× 17///10				258(1)-	
B	• 🕵		test_pm_01_root		test_pm	_01_root		80434		手动添加		
17	● 函 absdadnas 我锁定 10-23 18:5											
244	Sh		本节点的输出	节点输出名称								构
Ť												
	• Ja		输出名称		輸出表名	下游节点名称		下游节点ID	责任人	来源	操作	
	56											
	• 📦		test_pm_01.absdadna	is Ø	- @					手动添加		

注意事项

当任务依赖配置完成后,提交的窗口会有一个选项:当输入输出和代码血缘分析不匹配时,是否确 认继续执行提交操作。 该选项的前提是您已经确认依赖关系正确。如果不能确认,则可以按照上述方法确认依赖关系。

确认依赖关系无误后,直接勾选我确认继续执行提交操作,并单击确认。

2.2 Eclipse Java UDF开发最佳实践

本文将为您介绍如何使用Eclipse开发工具,配合ODPS插件进行Java UDF开发的全流程操作。

准备工作

开始使用Eclipse进行Java UDF开发前,您需要进行如下准备工作:

1. 使用Eclipse安装ODPS插件。

2. 创建ODPS Project。

a. 在Eclipse中选择File > New > ODPS Project, 输入项目名称, 单击Config ODPS console installation path, 配置odpscmd客户端安装路径。

New ODPS Project Wizard	
Create ODPS project	
Project name: ODPS JAVA UDF	
Vse default location	
Location: C:\Users\furui.fr\eclipse-workspace\OI	DPS JAVA UDF Browse
Config ODPS console installation path	
 Use default ODPS console installation path Specify ODPS console installation path Version: 0.29.4 	Config ODPS console installation path Browse
? Sack	Next > Finish Cancel

b. 输入客户端整体安装包的路径后,单击Apply,ODPS插件会为您自动解析出客户端Version。

Preferences		X
ODPS Settings	Config ODPS console installation path	
	Config ODPS console installation path	
	C:\odpscmd_public Brow	wse
	Version: 0.29.4	
	Run Mode	
	● Local ○ Remote	
	limit record count of downloaded	
	100 (0~10000)	
	Retain local job temp directory	
	Restore Defaults Ap	ply
?	Apply and Close Cance	əl

c. 单击Finish,即可完成项目的创建。

开发步骤

- 1. 在ODPS Project中创建Java UDF。
 - a. 在左侧Package Exploer中右键单击新建的ODPS Java UDF项目,选择New > UDF。



b. 输入UDF的Package名称(本例中为com.aliyun.example.udf)和Name(本例中 为Upper2Lower),单击Finish,即可完成UDF的创建。

New UDF		
New UDF		
Create a new l	JDF implementation.	
Source folder:	ODPS JAVA UDF/src	Browse
Package:	com.aliyun.example.udf	Browse
Name:	Upper2Lower	
Superclass:	com.aliyun.odps.udf.UDF	Browse
Interfaces:		Add
		Remove
?	Fin	ish Cancel

完成UDF创建后,您即可看到生成的默认Java代码,请注意不要改变evaluate()方法的名称。



2. 实现UDF类文件中的evaluate方法。

将您想要实现的功能代码写到evaluate方法中,且不要改变evaluate()方法的名称。此处以实现大写字母转化为小写字母为例。

```
WordCount.java
                  Upper2Lower.java
                                      TestUpper2Lower
                                                           ☑ Upper2Lower.java ⋈
 1 package com.aliyun.example.udf;
  2
 3 import com.aliyun.odps.udf.UDF;
 4
 5 public class Upper2Lower extends UDF {
  6⊜
        public String evaluate(String s) {
 7
            if (s == null) { return null; }
 8
            return s.toLowerCase();
 9
        }
 10 }
```

```
package com.aliyun.example.udf;
import com.aliyun.odps.udf.UDF;
public class Upper2Lower extends UDF {
    public String evaluate(String s) {
        if (s == null) { return null; }
            return s.toLowerCase();
        }
}
```

代码编写完成后,请及时保存。

测试Java UDF

为了测试Java UDF代码,您可以先在MaxCompute上存放一些大写字母作为输入数据。您可以 在odpscmd客户端使用SQL语句create table upperABC(upper string);,新建一个名 为upperABC的测试表格。



使用SQL语句insert into upperABC values('ALIYUN');,在表格中插入测试用的大写字 母ALIYUN。

完成测试数据准备后,单击Run > Run Configurations,配置测试参数。



配置测试参数:Project一栏中填写创建的Java ODPS Project名称,Select ODPS project中填 写您的MaxCompute项目名称(请注意与odpscmd客户端当前连接的MaxCompute项目保持一 致),Table一栏填写刚才创建的测试表格名称。完成配置后单击Run进行测试。

Run Configurations		PRAKENET	
Create, manage, and run config	urations		
Image: Second	Name: Upper2Lower G UDF UDTF UDAF A JRE Classpath E Envir Project: ODPS JAVA UDF UDF UDTF UDAF class: com.aliyun.example.udf.Upper2Lower Select ODPS project MaxCompute_DOC example_project Input Table Table: upperABC Partitions: Columns:	ronment Common	fault all partitions)
Filter matched 10 of 10 items		Revert	Apply
2		Run	Close

您可以在Console中查看测试结果。

						~
					•	
💦 Problems @ Javadoc 😣 Declaration 🗐 Con	nsole ¤				• ×	‰ [
<terminated> Upper2Lower [ODPS UDF UDTF U</terminated>	JDAF] C:\P	Program Files\Ja	va\jre1.8.0_1	92\bin\javaw	.exe (2	2018£
[INFO]Finished to write table scheme [INFO]Start to download table: 'MaxCon [INFO]Tunnel DownloadSession ID is : 2 [INFO]Start to write table: MaxCompute [INFO]Finished write table: MaxCompute	: MaxCom mpute_DO 20181214 e_DOC.up e_DOC.up	mpute_DOC.upp DC.upperABC', 417544782dcdb operABC>C:\ operABC>C:\	erABC>C:\ download m 0b0f817516 Jsers\furui Jsers\furui	Users\furu ode:AUTO .fr\eclips .fr\eclips	ii.fr\ se-wor	ecli kspa kspa
aliyun						



测试结果只是Eclipse获取表格中的数据后在本地转换的结果,并不代表MaxCompute中的数据 已经转换为小写的aliyun了。

使用Java UDF

确定测试结果正确后,即可开始使用Java UDF,操作步骤如下:

1. 导出Jar包

eclipse-workspace	- C	DPS JAVA UDF/src/com/al	iyun/example/udf/Upper2Lower.ja	ava - Eclipse	
File Edit Source I	Refa	actor Navigate Search	Project Run Window Help		
🔁 🖛 🖩 🕼 👎 🥖 🤅	≥ ℝ	₽ 🔲 🔳 🗣 🕶 💽 🕶	♀ ▼ 🔮 🮯 ▼ 🤔 🗁 🔗 ▼ 🖗 י	▼ 初 ▼ ⁽) ▼ ○ ▼	
😫 Package Explorer	ß	□ 😫 😂 🗢 🗆	🛿 WordCount.java 🔹 🛽 Upper2l	Lower.java 🖹 TestUpper2Lower 🖸 Upper2Lower.java 🛛 🖓 🗖	₽
MaxCompute_DOC			1 package com.aliyun.exam	nple.udf;	
A 💯 ODPS JAVA UDI)F		2	NHE LIDE.	
 ▲ src ▲ com.aliyu ▷ Upper2 ▷ JRE System I ▷ A Referenced ▷ A Referenced ▷ examples ▷ temp ▷ temp ▷ temp 		New Co.Into	,	idi.obr,	4
		Go Into		• extends UDF {	
		Open in New Window		te(String s) {	
		Open Type Hierarchy		Case():	
		Show In	Alt+Shift+W		
		Copy	Ctrl+C		
		Paste	Ctrl+V		
		Delete	Delete		
		Remove from Context	Ctrl+Alt+Shift+Down		
		Build Path	+		
		Source	Alt+Shift+S •		
	2	Refactor	Alt+Shift+T ►		
		Import			
	4	Export			
	S	Refresh	F5		
		Close Project			
		Close Unrelated Projects			
		Assign Working Sets			
		Coverage As	•		
		Run As	•		
		Debug As	•		
		Validate			
		Team	n y		
		Compare With	•	► E	
		Configure	•	ation 🖳 Console 🛛 🔲 🗶 💥 🗎	e 🔓
		Properties	Alt+Enter	UDF UDTF UDAF] C:\Program Files\Java\jre1.8.0_192\bin\javaw.exe (2018年	12)
			[INFO]Finished to write ta	<pre>ble scheme : MaxCompute_DOC.upperABC>C:\Users\furui.fr\eclip</pre>	se-
			[INFO]Start to download ta	ble: 'MaxCompute_DOC.upperABC', download mode:AUTO	
			[INFO]Start to write table	: MaxCompute DOC.upperABC>C:\Users\furui.fr\eclipse-workspac	e\(

在左侧新建的ODPS Project上右键单击,选择Export。

在弹框中选择JAR file, 单击Next。
Export	_ D X
Select Export resources into a JAR file on the local file system.	2
Select an export wizard:	
type filter text	
🕨 🗁 Install	
 ▲ Java ↓ JAR file ↓ Javadoc 	
📮 Runnable JAR file	=
Run/Debug	
▷ 🗁 Tasks	
▷ ➢ XML	-
? < Back Next > Finish	Cancel

在对话框中JAR file处填写Jar包名称,单击Finish,即可导出至当前workspace目录下。

JAR Export	
JAR File Specification (1) The export destination will be relative	e to your workspace.
Select the resources to export:	 ✓ I.classpath ✓ I.project
 Export generated class files and reso Export all output folders for checked Export Java source files and resource Export refactorings for checked proj Select the export destination: 	urces l projects es ects. <u>Select refactorings</u>
JAR file: upper.jar Options: Compress the contents of the JAR file Add directory entries Overwrite existing files without warn	✓ Browse e ing
? < Back Nex	xt > Finish Cancel

2. 使用DataWorks引用Jar包

登录DataWorks控制台,进入同一个项目(本例中为项目MaxCompute_DOC)的数据开发页 面。选择业务流程 > 资源 > 新建资源 > JAR,新建一个JAR类型资源。



在弹窗中上传您刚导出的JAR资源。

新建资源			×
* 资源名称:	upper.jar		
目标文件夹:			
资源类型:	JAR	~	
	✓ 上传为ODPS资源本次上传,资源会同步上传至ODPS [®]	中	
上传文件:	upper.jar (47.66K)	×	
		确定	取消

刚才只是将JAR资源上传至DataWorks,接下来您需要单击进入JAR资源,单击提交并解锁(提交)按钮,将资源上传至MaxCompute。

d b f	
上传资源	
已保存文件:	upper.jar
资源唯一标识:	OSS-KEY-i2397ptr0u3id1k39lp3cmrl
	✓ 上传为ODPS资源本次上传,资源会同步上传至ODPS中
重新上传:	点击上传

完成上传后,您可以在odpscmd客户端使用list resources命令查看您上传的JAR资源。

3. 创建资源函数

现在Jar资源已经存在于您的MaxCompute项目中了,接下来您需要单击业务流程 > 函数 > 新 建函数,新建一个与Jar资源对应的函数,本例中函数名称为upperlower_Java。完成后,依次 单击保存和提交并解锁(提交)

名称/创建人	
Sq JSONdata 我锁走 11-24 12:1	
Py Pytest 我锁定 12-14 14:58	注册函数
● Sp test 我锁定 12-01 11:01	函数名: upperlower_java
● Mr testMR 我锁定 10-25 11:56	a ₩dz.
● VI vi 我锁定 11-16 10:51	* 중수. com anyon example our upperzetower
> 🔠 表	* 资源列表: upperjar
✔ 🧭 资源	
Fi abc.py 我锁定 10-18 14:21	摘述:
Py ipint.py 可编辑 11-27 19:37	
Ja mapreduce-examples.jar 設帐	
Ja testJAR.jar 我锁定 10-2510	命令楷式:
Ja upper.jar 我微定 12-14 11:32	
✔ 🔂 函数	参数说明:
Fx ipint 可编辑 11-27 19:33	
▶ upperlower_java 影號定 1:	

完成提交后,您可以在odpscmd客户端使用list functions命令,查看已注册的函数。至此,您使用Eclipse开发工具注册的Java UDF函数upperlower_Java已经可用了。

使用Java UDF结果验证

打开您的odpscmd命令行界面,运行select upperlower_Java('ABCD') from dual;命

令,可以观察到该Java UDF已经可以转换字母的大小写了,函数运行正常。

+	-+				
l_c0	1				
labcd	-+ -+				
1 records (a	t most 10000	supported)	fetched by	y instance	tunnel.

更多信息

更多Java UDF开发示例请参见Java UDF。

如果您要使用IntelliJ IDEA开发工具完成完整的Java UDF开发过程,请参见IntelliJ IDEA Java UDF开发最佳实践。

2.3 使用MaxCompute分析IP来源最佳实践

本文将为您介绍如何在MaxCompute上分析IP来源,包括下载、上传IP地址库数据及编写UDF函

数、编写SQL四个步骤。

背景介绍

淘宝IP地址库的查询接口为IP地址字串,使用示例如下。

)
{"code" 京","coi	:0,"data": {"ip":"114.114.114.114","country":"中国","area":"","region":"江苏","city":"南 unty":"XX","isp":"XX","country_id":"CN","area_id":"","region_id":"320000","city_id":"320100","county_id":"xx","isp_id":"xx"}}
由于在	MaxCompute中禁止使用HTTP请求,目前可以通过以下三种方式,实现在MaxComput
中查询	JIP。
・用S	SQL将数据查询到本地,再发起HTTP请求查询。
ſ	〕 说明:
效	率低下,且淘宝IP库查询频率需小于10QPS,否则拒绝请求。
・下葬	烖IP地址库到本地,进行查询 。
ſ	月 说明·
同	兰 秋羽: 样效率低,且不利于数仓等分析使用。
・将I	P地址库定期维护上传至MaxCompute,进行连接查询。本文重点为您介绍该方式。
ſ	〕 说明:
比	较高效,但是IP地址库需自己定期维护。
]P地址/	库
1. 首纥	先您需要获取地址库数据。地址库您可以自行获取,本文仅提供一个UTF8格式的不完整的地
址四	革demo 。
2. 下载	烖UTF-8地址库数据到本地后,检查数据格式,示例如下。

```
0,16777215,"0.0.0.0","0.255.255.255","","","内网IP","内网IP","内网IP"
16777216,16777471,"1.0.0.0","1.0.0.255","澳大利亚","","","",""
16777472,16778239,"1.0.1.0","1.0.3.255","中国","福建省","福州市","","电信"
```

前四个数据是IP地址的起始地址与结束地址:前两个是十进制整数形式,后两个是点分形式。这 里我们使用整数形式,以便计算IP是否属于这个网段。

说明:

如果您需要使用真实IP地址,请自行下载IP地址库,具体的下载地址和使用方式请参见云栖社 区。

上传IP地址库数据

1. 创建表DDL,您可以使用MaxCompute客户端进行操作,也可以使用DataWorks进行图形化 建表。

```
DROP TABLE IF EXISTS ipresource ;
CREATE TABLE IF NOT EXISTS ipresource
(
    start_ip BIGINT
   ,end_ip BIGINT
   ,start_ip_arg string
   ,end_ip_arg string
   ,country STRING
   ,area STRING
   ,city STRING
   ,county STRING
   ,isp STRING
);
```

2. 使用Tunnel上传下载命令上传您的文件,本例中ipdata.txt.utf8文件存放在D盘。

odps@ workshop_demo>tunnel upload D:/ipdata.txt.utf8 ipresource;

可以通过SQL语句select count(*) from ipresource;查看表中上传的数据条数(由于 地址库有人更新维护,所以条目数会不断增长)。

使用SQL语句select * from ipresource limit 10;查看ipresource表前10条的样本数据,示例如下。

Job Queueing					
start_ip	end_ip	start_ip_arg	end_ip_arg	country area city county isp	
+ 3395369026 3395369027 3395369029 3395369030 3395369031 3395369034 3395369035	3395369026 3395369028 3395369029 3395369030 3395369033 3395369034 2205569035	<pre></pre>	*202.97.56 *202.97.56 *202.97.56 *202.97.56 *202.97.56 *202.97.56 *202.97.56	++++ 6.66″ "中国" "湖南省" "长沙市" "" "电 6.63″ "中国" "黑龙江省" "" "" "他信 6.69″ "中国" "安徽省" "合肥市" "" "电 7.70″ "中国" "湖南省" "长沙市" "" "电 7.73″ "中国" "澜南省" "长沙市" "" "电信 7.74″ "中国" "湖南省" "长沙市" "" "电信	信″ ″ 信″ ″
3395369035 3395369037 3395369038 3395369039	3395369036 3395369037 3395369038 3395369040	202. 97. 56. 75 202. 97. 56. 77" 202. 97. 56. 78" 202. 97. 56. 79"	202.97.56 202.97.56 202.97.56 202.97.56	.76 中国 黒龙江省 1 电语 .77″ ″中国″ ″江苏省″ ″南京市″ ″″ .78″ ″中国″ ″湖南省″ ″长沙市″ ″″ ″电 .80″ ″中国″ ″黑龙江省″ ″″ ″″ ″″	 信″ 信″

编写UDF函数

通过编写Python UDF,将点号分割的IP地址转化为整数类型的IP地址,本示例使用ataWorks的PyODPS节点完成。

1. 右键单击相应业务流程下的资源,选择新建资源 > Python。



- 2. 在新建资源对话框中,填写资源名称,并勾选上传为ODPS资源,单击确定。
- 3. 在新建的Python资源内,编写Python资源代码,示例如下。

```
from odps.udf import annotate
@annotate("string->bigint")
class ipint(object):
    def evaluate(self, ip):
        try:
        return reduce(lambda x, y: (x << 8) + y, map(int, ip.
split('.')))
        except:</pre>
```

return 0

单击提交并解锁。



- 4. 右键单击相应业务流程下的函数,选择新建函数。
- 5. 在新建函数对话框中,填写函数名称,单击提交。
- 6. 编辑函数配置,单击提交并解锁。

Fx ipint		•	Py ipint.py		Di ODPS2		Di json2max_	Di json2max	Di MQ2MaxCompute ×	Sq JSONdat	a x	Mr testMR	Sq a
	[↑]	ß		7)									
注册函	鐓												
					函数名:								
					* 类名:	ipint.ip	pint						
					*资源列表:	ipint.p	у						
					描述:								
					命令格式:								

本示例中, 函数的类名为ipint.ipint, 资源列表填写上文提交的资源的名称。

7. 验证ipint函数是否生效并满足预期,您可以在DataWorks上新建一个ODPS SQL类型节点,执行SQL语句进行查询,示例如下。



您也可以在本地创建ipint.py文件,使用MaxCompute客户端上传资源。

odps@ MaxCompute_DOC>add py D:/ipint.py; OK: Resource 'ipint.py' have been created.

完成上传后,使用客户端注册函数。

odps@ MaxCompute_DOC>create function ipint as ipint.ipint using ipint. py; Success: Function 'ipint' have been created.

完成注册后,即可使用该函数。您可以在客户端运行select ipint('1.2.24.2');进行测试。

蕢 说明:

如果同一主账号下其他项目需要使用这个UDF,您可以进行跨项目授权。

1. 创建名为ipint的package。

odps@ MaxCompute_DOC>create package ipint; OK

2. 将已经创建好的UDF函数加入package。

odps@ MaxCompute_DOC>add function ipint to package ipint; OK

3. 允许另外一个项目bigdata_DOC安装这个package。

odps@ MaxCompute_DOC> allow project bigdata_DOC to install package ipint; OK

4. 切换到另一个需要使用UDF的项目bigdata_DOC,安装package。

```
odps@ MaxCompute_DOC>use bigdata_DOC;
odps@ bigdata_DOC>install package MaxCompute_DOC.ipint;
OK
```

5. 现在您就可以使用这个UDF函数了,如果项目空间bigdata_DOC的用户Bob需要访问这些资

源,那么管理员可以通过ACL给Bob自主授权。

odps@ bigdata_DOC>grant Read on package MaxCompute_DOC.ipint to user aliyun\$bob@aliyun.com; --通过ACL授权Bob使用package

在SQL中使用



本文以一个随机的IP 1.2.24.2地址为例,您在使用时可以用具体表的字段来读入。

测试使用的SQL代码如下,单击运行即可查看结果。

select * from ipresource
WHERE ipint('1.2.24.2') >= start_ip

AND ipint('1.2.24.2') <= end_ip

		۵ 🗊	() :	\$												
		from ipreso	urce													
	WHERE ip	oint('1.2.24.	2') >= st	art_ip												
	AND 1p1r	it('1,2,24,2') <= end_	тр												
															不	
															67	
															К Я	
			_													
运行	行日志	结果[1]	× 结	果[2] ×												
	А															
1	start_ip	✓ end_ip	~	start_ip_arg	~	end_ip_arg	~		~		~	city	 county 	~	isp	~
2	16910592	16941055		'1.2.9.0"		"1.2.127.255"		"中国"		"广东省"		"广州市"			"电信"	

通过为保证数据准确性,您可以定期从淘宝IP库获取数据来维护ipresource表。

2.4 在PyODPS任务中调用第三方包

本文将为您介绍如何使用DataWorks PyODPS类型任务节点调用单文件第三方包。

1. 右键单击相应业务流程下的资源,选择新建资源 > Python。



2. 在新建资源对话框中,填写资源名称,并勾选上传为ODPS资源。

新建资源				×
	*资源名称:	test2.py		
	目标文件夹:			
	资源类型:	Python	~	
		上传为ODPS资源本次上传,资源会同步上传至ODPS中	þ	
			确定	取消

- 3. 单击确定。
- 4. 在新建的Python资源内,粘贴需要引用的第三方包的代码,示例如下。

```
# import os
# print os.getcwd()
# print os.path.abspath('.')
# print os.path.abspath('..')
# print os.path.abspath(os.curdir)
def printname():
    print 'test2'
print 123
```

粘贴代码完成后,单击提交。

1	# import os
2	<pre># print os.getcwd()</pre>
3	<pre># print os.path.abspath('.')</pre>
4	<pre># print os.path.abspath('')</pre>
5	<pre># print os.path.abspath(os.curdir)</pre>
6	
7	
8	def printname():
9	print 'test2'
10	
11	print 123

5. 在您的业务流程内新建一个PyODPS类型节点。

🗸 🛃 works		5 sys. 6 impo
▶ 😑 数据集	成	7 test
▼ ひ 数据	···· 新建数据开发节点 >	8 ODPS SQL
● Sq Ir	新建文件夹	ODPS MR
Vi s	看板	虚拟节点
● Sq t€	引用组件	PyODPS

6. 在节点内输入引用第三方包的代码并测试,示例如下。

```
##@resource_reference{"test2.py"}
import sys
import os
sys.path.append(os.path.dirname(os.path.abspath('test2.py'))) #将资源
引入工作空间
import test2 #引用资源
test2.printname() #调用方法
```

请注意下图中框内的代码,用于引用业务流程中您之前新建的test2.py资源,请不要遗漏。



7. 完成上述操作后,单击运行测试您的代码。您可以在下方的日志中查看运行结果。



2.5 分支节点实现特定时间执行任务

本文将为您介绍分支节点如何实现在特定时间执行任务。

分支节点产生背景

在日常DataWorks的使用过程中,如果您有一个节点,需要每个月的最后一天执行,可以进行如下 设置。

在分支节点出现前,由于Cron表达式无法实现该场景,所以暂时无法支持。

现在,DataWorks已经正式支持分支节点。利用分支节点,您可以套用switch-case编程模型实现 上述需求。

分支节点与其他控制节点

在数据开发页面,您可以看到当前版本的DataWorks支持的各种控制节点,包括赋值、分支、归并 节点等。

各类型控制节点的作用如下:

・赋值节点:可以把自己的结果传给下游。

赋值节点复用了节点上下文依赖的特性,在已有常量/变量两种节点上下文的基础上,赋值节点 自带一种自定义的上下文输出。DataWorks会捕获或打印赋值节点的select结果,并将该结果 以outputs形式作为上下文输出参数的值,供下游节点引用。

- ・分支节点:可以决定哪些下游正常执行。
 - 分支节点复用了DataWorks上依赖关系设置的输入输出的特性。

对于普通节点,节点的输出仅仅是一个全局唯一的字符串。当下游需要设置依赖时,搜索这个全 局唯一的字符串作为节点的输入就能挂到下游节点列表中。

但是,对于分支节点您可以给每个输出关联一个条件:

当下游设置依赖时可以选择性的把某一个条件关联的输出作为分支节点的输出。这样,节点在成 为分支节点下游的同时,也关联到了分支节点的条件上:

- 满足该条件,该输出对应的下游才会被正常执行。
- 其他未满足条件的输出对应的下游节点, 会被置为空跑。

・ 归并节点: 无论上游是否正常执行, 本身都会正常调度。

对于未被分支节点选中的分支,DataWorks会把这个分支链路上所有的节点实例均置为空跑实例,也就是说一旦某个实例的上游有一个实例是空跑的话,它本身也会变为空跑。

DataWorks当前可以通过归并节点来阻止这个空跑的属性无限制的传递下去:对于归并节点实例,无论它的上游有多少个空跑的实例,它都会直接成功并且不会再把下游置为空跑。

您可以从下图看到在有分支节点的情况下,依赖树的逻辑关系。



· ASN: 一个赋值节点, 用于对比较复杂的情况做计算, 为分支节点条件选择做准备。

- · X/Y:分支节点,他们处于赋值节点ASN下游,根据赋值节点的输出做分支的选择。如图中绿色
 线条所示,X节点选择了左边的分支,Y节点选择了左边两个分支:
 - A/C节点由于处于了X/Y节点被选择的输出下游,因此正常执行。
 - B节点虽然处于Y节点被选择的分支下游,但由于X节点未选择这个输出,因此B节点被置为 空跑。
 - E节点由于未被Y节点选中,因此即使有一个普通的Z节点上游,也同样被置为了空跑。
 - G节点由于上游E节点空跑,因此即使C/F都正常执行,G节点同样空跑。
 - 空跑属性什么情况下才能不再向下传递?

JOIN节点是一个归并节点,它的特殊功能就是停止空跑属性的传递,可以看到由于D节点处于JOIN节点下游,因此B节点的空跑属性被阻断了,D节点可以开始正常跑了。

您可以通过利用分支节点配合其他控制节点,满足某个节点只有每个月最后一天运行的需求场景。

使用分支节点

定义任务依赖

首先您需要定义一组任务依赖。



- 根节点赋值节点通过定时时间SKYNET_CYCTIME来计算当前是不是本月的最后一天,如果是则 输出1,不是则输出0。该输出会被DataWorks捕获,传递给下游。
- 2. 分支节点通过赋值节点的输出来定义分支。
- 3. 两个shell节点挂在分支节点下面,分别执行不同的分支逻辑。

定义赋值节点

赋值节点新建时会自带一个outputs,赋值节点的代码支持SQL/SHELL/Python三种。

- · 对于SQL类型, DataWorks捕获最后一条SELECT语句作为outputs的值。
- ・ 对于SHELL/Python类型, DataWorks捕获最后一行标准输出作为outputs的值。

本文采用Python类型作为赋值节点的代码,调度属性和代码设置如下。

・代码设置

Sh 除了,	最后一天之外运行 × Sh 只有最后一天执行 × 🚴 分支_根据最后一天决定… × 🔗 赋值、判断今天是不是最… × 🛃 分支节点DEMO ×
	请选择赋值语言: Python
	import os import time
	<pre>dueTimeStr = os.environ['SKYNET_CYCTIME'] # 20190104000200 dueTime = time.strptime(dueTimeStr[:8], "%Y%m%d") dueMonth = time.strftime("%m", dueTime)</pre>
	<pre>nextDueTimeStamp = int(time.mktime(dueTime))+3600*24 nextDueTime = time.localtime(nextDueTimeStamp) nextDueMonth = time.strftime("%m", nextDueTime)</pre>
	<pre>print ('Current month: %s, next month: %s') % (dueMonth, nextDueMonth)</pre>
	<pre>if nextDueMonth == dueMonth: print 0 else:</pre>
17	print 1

・调度属性配置

Sh 除了	了最后一天之外	运行 × Sh 只	有最后一天执行	× 🔥 分支_桁	親据最后一天决定.	× 🛕 赋值_判断今;	天是不是最	× 晶分	友节点DEMO ×			≡
		C C									发布	
	请	×										调度
1	节 import	autotest_ro	ot			autotest_root	100933548			手动添加		電置
2	import	本节占的输出	请输λ节占输	电夕森								版本
4	dueTim	~ р // р	时间八口/2/191									
5 6	dueTim dueMon	输出名称			输出表名	下游节点名称		下游节点	ID 责任人	来源	操作	
7 8	nextDu	autotest.92	08506_out			分支_根据最后一天决	定下游执行			系统默认添加		
9 10 11	nextDu nextDu	autotest.赋f	值_判断今天是不是	量最后一天 ⑦	- C					手动添加		
12	print											
13	if nex	节点上下文	. @									
15 16	prin else:	本节点输入参数	数 添加									
1/	prin	编号	参数名	取值来》	原	描述		父节点ID	来源	操作		
						没有数据						
		本节点输出参数	数 添加									
		编号	参数名	类型	取值	描述 outputs	: 赋值节点		来源	操作		
		1	outputs	变量	\${outputs}	赋值节点输出值,取值	直由运行时决定		系统默认添	加编辑		

定义分支

分支节点可以用Python语法的表达式定义条件,每个条件会绑定一个输出。当满足这个条件 时,该输出的下游节点会被执行起来,而其他的节点会被置为空跑。

・调度配置

Sh 除	了最后一天之	外运	行 × Sr	只有最后一天执行	ī × 🗼 分3	友_根据最后一天决策	定 × 晶 分支节点DEMO ×					Ξ
			ê C	к К							发布	
	分支逻辑员	>	く 本节点的输	出 请输入节点转	俞出名称							调度配置
			输出名和	尔		输出表名	下游节点名称	下游节点ID	责任人	来源	操作	血
	支		outotoot	+ 0208507 out		(1)	只有最后一天执行		80	マムの事やよし法もn		^嫁 关 系
			autotesi	l.9206507_0ut		- @	除了最后一天之外运行			がまいまたてきること		版本
	2		autotest	t.last_day_cond.is_la	ast	- C				系统默认添加		
			autotest	t.last_day_cond.not_	last	- C				系统默认添加		
			节点上下 ^{本节点输入}	「文 ⑦ 参数 添加 ^{参数 条数 2}	依加	臉赋值节点的	输出:最后一天isLast	为1,其他情况	isLast为0	运 作:		
			1	isLast	autotest.920	8506_out:outputs	赋值节点输出值,取值由运行	时决定 1078764	442 系统默	3本F 3、汤添加 编辑 删除).	
			本节点输出	参数 添加								
			编号	参数名	类型	取值	描述		来源	操作		
				outputs	变量	\${outputs}	分支节点输出值,取值由运	行时决定	系统默认	添加 编辑		

・ 分支配置

分支逻辑定义	X ()			
添加分支	根据isLast定义 ^ス	不同的分支		
分支	条件	关联到节点输出	分支描述	操作
	\${isLast}==0	autotest.last_day_cond.not_last	如果当天不是最后一天,则使用这个分支	编辑 删除
2	\${isLast}==1	autotest.last_day_cond.is_last	如果当前是最后一天,则使用这个输出	编辑 删除

· 调度配置生成条件绑定的输出

Sh 除了	了最后一天	之外运	行 × Sh 只有	有最后一天执行 🗙 🍌 分支	一根据最后一天决策	È × 嚞 分支节点DEMO ×					≡
			ê C							发布	
	分支逻辑	版 文	く 本节点的输出	请输入节点输出名称							调度配置
	分		输出名称		输出表名	下游节点名称	下游节点ID	责任人	来源	操作	血缘
	文		autotort 0200	2507 out	CI.	只有最后一天执行			妥体财计添hn		天系
	1		autotest.9200	5007_00t	- 0	除了最后一天之外运行			<u>ЖЭРим м/икли</u>		版本
	2		autotest.last_	.day_cond.is_last	- C				系统默认添加		
			autotest.last_	.day_cond.not_last	- Ø	S里,默认添加了网个与 -	余件问名的制	н -	系统默认添加		

将执行任务节点挂在不同分支下

最后,给真正执行任务的节点设置依赖时需要注意:您可以看到分支节点已经有三个输出了,按照 过去设置依赖的逻辑,把这三个输出中的任意一个当做输入即可。由于现在分支节点的输出关联了 条件,所以要慎重选择。

· 每月最后一天执行的节点依赖

Sh 除了	。 最后一天之外读	运行 × Sh 只有最后一天执行 × 🙏	分支_根据最后一天决定	× 嚞 分支节点DEMO ×					≡
	E) I	[b] 🗊 🕑 :						发布	
1 2 3	#!/bin #***** ##auth	× 调度依赖 ⑦							调度配置
4 5 6	##crea #***** echo "	自动解析: • 是) 否 解析输入输 依赖的上游节点 请输入父节点输出名	₩ 出 称或输出表名	使用工作空间根节点					血缘关系
		父节点输出名称	父节点输出表名	节点名	父节点ID	责任人	来源	操作	肟
		autotest.last_day_cond.is_last 挂在每月最后一天的分支下面		分支_根据最后一天决定下游执行			手动添加		本

· 每月其他时间执行的节点依赖

Sh 除了最后一天之	外运行 × 🗼 分支_根据最后一天决定	× 🛃 分支节点DEMO	×					≡
] [5] 🗊 🕑 :						发布	
1 #!/bin 2 #***** 3 ##auth 4 ##crea 5 #*****	X 调度依赖 ⑦ 自动解析: ● 是 ○ 否 解析输入输	н						调度配置血
6 echo"	依赖的上游节点 请输入父节点输出名称	你或输出表名 🛛 🖌	+ 使用工作空间根节点					缘 关 系
	父节点输出名称	父节点输出表名	节点名	父节点ID	责任人	来源	操作	版
	autotest.last_day_cond.not_last 挂在非最后一天的分支下面		分支_根据最后一天决定下游执行		80	手动添加		本
	本节点的输出 请输入节点输出名称							

结果验证

完成上述所有配置后,您可以提交并发布任务。完成发布后,可以执行补数据来测试效果:业务日 期选择2018-12-30和2018-12-31,也就是定时时间分别为2018-12-31和2019-01-01,这样第一 批补数据会触发最后一天的逻辑,第二批触发非最后一天的逻辑。两者的区别如下所示。

业务日期2018-12-30(定时时间2018-12-31)

· 分支节点分支选择结果

>			S 7		話最 点	
属性	上下文	运行日志	操作日志		代码	
 ⊘ 01-04 23:26:19 23:26:25 持续时间: 6s Gateway: 	9 ~ 01-04	2019-01-04 23 2019-01-04 23 python output 2019-01-04 23 2019-01-04 23 python output 2019-01-04 23	:26:19.867 :26:21.239 : False :26:21.419 :26:22.424 : True :26:22.435	INFO INFO INFO INFO INFO	- The foll - Started - not meet	owing profiles are active: dev ControllerWrapper in 2.072 seconds (JVM runni the condition! condition:1 == 0
		2019-01-04 23 2019-01-04 23 2019-01-04 23 2019-01-04 23 2019-01-04 23	26:23.436 26:23.437 26:24.070 26:24.071	INFO INFO INFO INFO	- meet the - ===>Outp - cost Tin - job fini	e condition. condition:1==1 out Result: autotest.last_day_cond.is_last ne: 2 .shed!

・节点(最后一天执行)正常执行

» •		\bigcirc	只有最后一天执行 SHELL		除了最后一天… SHELL	
属性	上下文	运行日志	操作日志	代码		
 ✓ 01-04 23:26:29 23:26:29 持续时间: 0s Gateway: 	9 ~ 01-04	2019-01-04 23: 2019-01-04 23: 2019-01-04 23: 2019-01-04 23: 2019-01-04 23: Is last day: 2	26:29 INFO AL 26:29 INFO AL 26:29 INFO AL 26:29 INFO 26:29 INFO 0181231000500	[SA_TASK_EXEC_1 [SA_TASK_PRIOR] - Invoking She]	ARGET=autotest_new_group ITY=1: 1 command line now	
-		2019-01-04 23: 2019-01-04 23: 2019-01-04 23: 2019-01-04 23: 2019-01-04 23: 2019-01-04 23: 2019-01-04 23:	26:29 INFO === 26:29 INFO Ex: 26:29 INFO -== 26:29 INFO Sho 26:29 INFO Cui 26:29 INFO Coi	it code of the - Invocation of ell run success rrent task stat st time is: 0.0	Shell command 0 F Shell command completed sfully! cus: FINISH 007s	

・ 节点(除了最后一天之外运行)被置为空跑

»		\odot	只有最后一天执行 SHELL	\odot	除了最后一天… SHELL	
属性	上下文	运行日志	操作日志	代码		
 ✓ 01-04 23:26:25 → 23:26:25 持续时间: 0s Gateway: 	5 ~ 01-04	It's set condi	ition-skip by task	(9868911318	3-分支_根据最后一天决定下游执行	Ĺ,

业务日期2018-12-31(定时时间2019-01-01)

· 分支节点分支选择结果

>				▶支_根据 分支节系	建最	
属性	上下文	运行日志	操作日志		代码	
 ✓ 01-04 23:26:4² 23:26:52 持续时间:11 Gateway: 	1 ~ 01-04 s	2019-01-04 23: python output: 2019-01-04 23: 2019-01-04 23: python output:	26:48.414 : True :26:48.614 :26:49.622 : False	INFO INFO INFO	- Started - meet the	ControllerWrapper in 3.485 seconds (JVM runni e condition. condition:0==0
_		2019-01-04 23 2019-01-04 23 2019-01-04 23 2019-01-04 23 2019-01-04 23 2019-01-04 23 2019-01-04 23	26:49.634 26:50.635 26:50.636 26:50.981 26:50.981 26:51 INFO	INFO INFO INFO INFO INFO	 - not meet - ===>0utp - cost Tin - job fini	the condition! condition:0==1 but Result: autotest.last_day_cond.not_last ne: 2 .shed!

・ 节点(最后一天执行)被置为空跑

		○ 除了最后 SHE	一天执行			
属性	上下文	运行日志	操作日志	代码		J
 ※ 01-04 23:26:52 ~ 01-04 23:26:52 持续时间: 0s Gateway: 		It's set condi	ition-skip by t	ask(9868911326	5-分支_根据最后一天决定下游执行)	

・ 节点(除了最后一天之外运行)正常执行

>>		○ 除了最后 SHEL	一天	○ 只有最后 SH	一天执行 ELL
属性	上下文	运行日志	操作日志	代码	
 ※ 01-04 23:26:58 注:23:26:58 持续时间: 0s Gateway: ■ 	8 ~ 01-04	2019-01-04 23:26:57 INFO ALISA_TASK_PRIORITY=1: 2019-01-04 23:26:57 INFO Invoking Shell command line now 2019-01-04 23:26:57 INFO Invoking Shell command line now Is not last day: 20190101000900 2019.01-04 23:26:57 INFO			
		2019-01-04 23: 2019-01-04 23: 2019-01-04 23: 2019-01-04 23: 2019-01-04 23: 2019-01-04 23: /home/admin/al	26:57 INFO Exi 26:57 INFO 26:57 INFO She 26:57 INFO Cu 26:57 INFO Cos isatasknode/to	it code of the - Invocation of ell run success rrent task stat st time is: 0.0 askinfo//20190	Shell command 0 5 Shell command completed sfully! cus: FINISH 004s 104/phoenixprod/23/26/52/xktmzpirgs5pe441ggryv

总结

基于分支节点,您已经实现了每个月最后一天执行的这一目标,当然这只是分支节点最简单的使用 方法。将赋值节点与分支节点配合使用,您可以组合出各种各样的条件满足业务上的需求。

最后回顾分支节点使用要点:

- · DataWorks捕获赋值节点的最后一条SELECT语句或者最后一行标准输出流,作为赋值节点的输出,供下游引用。
- · 分支节点的每一个输出都被关联了条件,下游挂分支节点作为上游,需要了解每个输出关联的条件的意义再选择。
- · 未被选中的分支会被置为空跑,并且空跑属性会一直向下传递,直到遇到归并节点。
- ·归并节点除了阻断空跑属性外,还有更多更强大的功能等待您的挖掘。

2.6 DataWorks数据服务对接DataV最佳实践

DataV通过与DataWorks数据服务的对接,可使用DataWorks数据服务开发数据API,快速 在DataV中调用API并展现MaxCompute的数据分析结果。

MaxCompute是阿里巴巴集团自主研究的快速、完全托管的TB/PB/EB级数据仓库解决方案。当 今社会数据收集的方式不断丰富,行业数据大量积累,导致数据规模已增长到传统软件行业无法承 载的海量级别。MaxCompute服务于批量结构化数据的存储和计算,已经连续多年稳定支撑阿里 全部的离线分析业务。

过去,如果您想要通过DataV展示海量数据的分析结果,需要自建一套离线数据计算自动导 入MySQL的任务流程,这个过程非常繁琐,且成本很高。而现在通过DataWorks为您提供的数据 集成 > 数据开发 > 数据服务的全链路数据研发平台,并结合MaxCompute可快速搭建企业数仓。

DataWorks数据服务提供了快速将数据表生成API的能力,通过可视化的向导模式操作,无需代码 便可快速生成API,然后通过DataV调用API并在大屏中展示数据分析结果,高效实现数仓的开发 和数据的展示。

数据服务对接DataV产生背景

本文为您介绍如何实现DataWorks数据服务与DataV联合进行API开发,并通过大屏可视化展现数 据分析结果。

前提条件

要想实现DataWorks数据服务与DataV的对接,您需提前准备好数据源,并开通DataV服务。

新建数据源

数据服务支持丰富的数据源类型,如下所示。

- ·关系型数据库: RDS/DRDS/MySQL/PostgreSQL/Oracle/SQL Server
- ・分析型数据库: AnalyticDB (ADS)
- · NoSQL数据库: TableStore (OTS) /MongoDB
- · 大数据存储: Lightning (MaxCompute)

- 1. 登录DataWorks控制台,单击对应项目后的进入数据服务。
- 2. 在服务开发页面,单击新建按钮,选择新建数据源。

Solution Solution	bipdata.DOC	~
≡	服务开发	₽₽₽€
₩ 服务开发	API名称	生成API >
■数据表	> API列表	注册API 新建数据源
		新建分组

3. 进入数据集成 > 同步资源管理 > 数据源页面,单击右上角的新增数据源。

数据集成											🔍 💷 ex
= ▼ 工作空间概览	数据源 数据源	类型: 全部	新增数据源					×			C 刷新 新増数振振
🖕 任务列表	数据源名称	数据源类型	关系型数据库						连通状态	连通时间	操作
👺 资源消耗监控	odar, Ser.	ODPS	MySQL	SQL Server	PostgreSQL	ORACLE [.]	ø				
- 同步资源管理			MySQL	SQL Server	PostgreSQL	Oracle	DM				
	cdpc.ex	ODPS	8	3	¢¢¢	\otimes		- 1	成功	2018/12/28 13:52:19	编辑 删除
* 2011	HER_Sets_source	HDFS	DRDS	POLARDB	HybridDB for MySQL	HybridDB for PostgreSQL					编辑 删除
	HOFSI	HDFS	大数据存储	¥	\Diamond	42					编辑 激除
	Lightning	Lightning	MaxCompute (0DPS) 半结构化存储	Datahub	AnalyticDB (ADS)	Lightning					編編 删除
			055	HDFS	FIP						
			NoSQL		3		R	8 4			

本文将以Lightning数据源为例,通过Lighning数据源可以直接实时查询MaxCompute中的数据。



说明:

Lightning目前是内测阶段,需要单独申请才能开通,您也可以加入数据服务用户群(钉钉群 号21993540)咨询Lightning服务的开通事项。

4. 单击Lightning,	填写新增Lightning数据源对话框中的配置。
-----------------	--------------------------

新增Lightning数据源		×
* 数据源名称:	Lightning	
数据源描述:	Lightning_test	
* Lightning Endpoint :	lightning on hangchournexcompane eiliyun com	
* Port :	443	
* MaxCompute项目 :	highlinin_DOC	
名称		
* AccessKey ID :	LT2ED&BPQ966-M12	
* AccessKey Secret :		
* JDBC扩展参数:	ssImode=require&prepareThreshold=0	
测试连通性:	测试连通性	
0	确保数据库可以被网络访问	
	确保数据库没有被防火墙禁止	
	确保数据库域名能够被解析	
	上一步	完成

配置	说明
数据源名称	数据源名称必须以字母、数字、下划线组合,且不能以数字 和下划线开头。
数据源描述	对数据源进行简单描述,不得超过80个字符。
Lightning Endpoint	Lightning的连接信息,详情请参见访问域名Endpoint。
Port	默认值为443。
MaxCompute项目名称	填写MaxCompute的项目名称。
AccessKey ID/AccessKey Secret	访问秘匙(AccessKeyID和AccessKeySecret),相当于 登录密码。

配置	说明
JDBC扩展参数	JDBC扩展参数中的sslmode=require&prepareThr eshold=0是默认且不可删除的,否则会无法连接。

5. 单击测试连通性,测试通过后,单击完成。

新建API

数据源创建完成后,进入数据服务页面,本文以向导模式生成API为例,为您介绍如何新建API。

1. 选择服务开发 > 新建 > 生成API > 向导模式。

 	bipdata.DOC	~	
	服务开发	₽₽₽₽₽	
	API名称	生成API >	向导模式
	_	注册API	脚本模式
■ 数据表	> API列表	新建数据源	
		新建分组	

2. 填写生成API对话框中的配置。

生成API			×
* AI	PI名称:)/50
* AI	PI分组:	Demo V	
* AP	PI Path :	/ D 5/	200
		支持英文,数字,下划线,连字符(-),且只能 / 开头,不超过200个字符,如/user	
•	*协议:	✓ НТТР	
*请	求方式:	GET ×	
*返	回类型:	JSON V	
	♥描述:	查询成交金额增长趋势API	
		14/2	000
			取消

配置	说明
API名称	支持汉字、英文、数字、下划线,且只能以英文或汉字开头,4~50个 字符。
API分组	API分组是指针对某一个功能或场景的API集合,也是API网关 对API的最小管理单元。您可以单击新建API分组进行新建。
API Path	API的存放路径,必须以/开头,如/user。
协议	目前仅支持HTTP协议。
请求方式	默认选择GET请求方式。
返回类型	默认返回JSON类型。
描述	对API进行简单描述。

3. 单击确认,进入API配置页面。

Detail	数据服务	R1 04539423 a.tohet					服务开发 服务管理	٩	staffin	中攻	Ż
Ш	服务开发	ନ⊑C⊕	📔 查询成交金额增长趋势	•							
3	API名称		∎ C								
⊞	> API列表		选择表								
	> - 7.000										旺
	> ••••••		• 数据源类型:	Lightning(ODPS)							请9
	> 🖬 alaa		• 数据源名称:	leightning							参数
	Y DEMO		• 数据表名称:	demo_trade_amount							
	😨 gan, Jagán, 20184 206, 1 💽 detaemica. api.demo 🚊 regapirtza		NH 472 45 WH								返回
			近洋梦致								参数
			搜索字段名称								
			设为语文参数	▶ 设为返回参数	李段名	字段类型	字段描述				
		155	879HAP90		TAN	776.4	7 PAILAE				
	🔝 para depe	6.300 M 201.7		✓	date	STRING					
	🔝 markaka				amount	STRING					
		int '			undurk						
		eet t									
		ind_regy									
	🔝 handara (leet.									

配置	说明
选择表	依次选择数据源类型、数据源名称和数据表名称进行表设置。 · 数据源类型:此处选择上一步创建的Lightning数据源。 · 数据源名称:选择对应的数据源名称。 · 数据表名称:选择要查询的MaxCompute表,此处以成交金 额表demo_trade_amount为例,该表中存储了一个月的成 交金额数据。
选择参数	选择表后,会自动展示表的字段列表。然后勾选要作为API请求参数的字段和作为返回参数的字段。
	 说明: 本示例中,为了查询成交金额趋势,需要返回所有数据,即日期 和成交金额都作为返回参数,不设请求参数。

单击右侧的返回参数,设置参数信息。

- 说明:

如果不设置请求参数,则需要开启返回结果分页开关,进行分页查询,以避免单次查询返回数 据量过大影响性能。

🔝 查询成交金额增长趋势	֥							
e c								
选择表		× 返回参数						属性
* 数据源类型:	Lightning(0DPS)	参数名称	绑定字段	參数类型	示例值	描述		请求
 数据源名称: 数据表名称: 	lightning	date	date	STRING V		日期		参数
选择参数		amount	amount	STRING ~		成交金額		返回参数
搜索字段名称		高级配置						
● 设为请求参数	✔ 设为返回参数	✔ 返回结果分页	当返回结果记录数大于500时请	告择分页,不分页则最多	返回500条记录。当无请求参	•数时,必须开启返回结果分页。		
		● 使用过滤器					0	
	✓							

4. 单击工具栏右边的测试,填写API请求参数(由于打开了分页查询开关,系统会自动添加两个分页参数),单击开始测试。

API 测试					×
API Path: /demo/trade 请求参数	e/amount			请求详情 // [INFO] [16:39:19.010] api context init,take time 2 ms	
参数名称	參数类型	是否必填	ųu 值	(INFO) [16:39:19.01] start to test spi(2815): 查询成交金额增长趋势 [INFO] [16:39:19.012] werlfy spi test(22640]. [OK] [INFO] [16:39:13.013] parse test case parameters. [OK]	
pageNum	int	是	1	[INFO] [16:39:19.014] test case parameters: {{rmarmfers?; pageNum*, paramValue*: '1 [INFO] [16:39:19.015] build backend sql aprequest. [OK] [INFO] [16:39:19.017] ready to execute api request. [OK] [INFO] [16:39:19.018] ani gel realer1 : statExect tast as "date" amount AS "amount"	
pageSize	int	是	31	[INFO] [16:39:19.019] query database starting [INFO] [16:39:20.644] query database finished. [OK]	
开始测试 🗹 🕯	自动保存正常返回示例			送回内容 ど回内容 ************************************	9 1
				④ Hot 0.5	

您可以在测试页面看到API延迟,会发现通过Lightning查询MaxCompute表仅花费了1秒

多,比直接通过MaxCompute SQL查询更高效。

发布API

新建API完成后,单击工具栏右侧的发布,即可将API发布。

发布完成后,您可以单击顶部导航栏中的服务管理查看API详情。

如果您要调用API,可进入服务管理API调用页面,数据服务为您提供简单身份认

证(AppCode)和加密签名身份认证(AppKey&AppSecret)两种认证方式,您可以自由选

择。下文将为您介绍如何在DataV中进行数据服务API的调用。

添加数据服务为数据源

- 1. 登录DataV控制台。
- 2. 进入我的数据页面,单击添加数据。

98

3. 填写添加数据对话框中的配置。

WELCOME TO DATAV Empowering Intelligent City	
	添加数据 🚯 🛛 🕹 🗙
◎ 我的可视化 ◎ 我的数据	*类型
2	DataWorks 数据服务 🗸 👻
* 法前期提	自定义数据源名称
	我的数据服务
	•项目
□ataWorks DW数据分析	
	АррКеу
	012+1;y039;11+01038+(s++
	*AppSecret
	·····································

配置	说明										
类型	添加的数据源类型。										
自定义数据源名 称	数据源的显示名称,可以自由命名。										
项目	DataWorks项目(工作空间)。										
AppKey/ AppSecret	拥有DataWorks数据服务中某一项目访问权限的账号 的AppKeyID和AppSecret。										
	说明: 您可以登录DataWorks数据服务控制台,进入服务管理 > API调用页面 进行查看。										

在大屏中调用数据服务API

1. 进入DataV控制台中的我的可视化页面,单击新建可视化。

2. 选择一个模板,单击创建,本文以智能工厂模板为例。



📋 说明:

模板中的组件自带了静态数据,下文将以把模板中间的基本折线图改为调用上文创建好的查询 成交金额增长趋势的API为例,为您介绍如何在组件中使用数据服务API。

- 3. 选中基本折线图组件,切换到数据面板,在数据源类型中选择DataWorks数据服务。
- 4. 选择刚刚创建的数据源和API,并设置查询参数,本示例将pageSize设置为31,以查询一个月的数据。



5. 单击查看数据响应结果,即可查看API的查询结果。

6. 填写字段映射关系,在x中填写date,将日期作为横轴,在y中填写amount,将成交金额作为 纵轴。

<	ılı ♀ ⊵ T	& b ↔ ⊕ h ♥	5	₽ ⋪ 0 0
图层				
£				基本折线图
TITLE		某某工厂车间实时状况 Crathed 2012/11/4		
201712:24	10%			
24	529	2 3 2 4 2 5 196020 75% 75% 98% 98%		
24			× 🞾 🔡	
24	人员供放 Workers Performance	使用过滤器 数据响应结果 点击写制		DataWorks 数据服务
TITLE		1 K 2 "data": { 3 "totalNum": 31		选择已有数据源:
TITLE		4 "pageSize": 31, 5 "rows": [我的政績服務 选择 API: 刷新
0	授發信息 Alern Information	7 "date": "2018/12/1", 8 "amount": "1000" 9 }.		查询成交金额增长趋势 ▼
TITLE		10 { 11 "date": "2018/12/2", 12 "amount": "1200"		
TITLE		13 b		
TITLE				
TITLE				
B				
24				查看数据响应结果
=				

📕 说明:

由上图可见,当前x和y无法匹配到字段。这是因为DataV对数据格式有一定要求,不能识 别结构较深的字段,因此需要添加一个数据过滤器,过滤掉不必要的字段,在本例中直接返 回rows数组即可。

7. 勾选使用过滤器,单击新建图标。这里支持编写JS代码对数据结果进行二次过滤和处理,过滤器的data参数为API返回结果JSON对象。



本示例只需返回API结果中的rows数组,因此您只需输入return data.data.rows;,便可 在下方预览过滤后的结果,并单击完成。

数据响应结果	×
☑ 使用过滤器	教程
+	× -
名称: MyFilter	
<pre>function filter(data) { 1 return data.data.rows; </pre>	
}	
· 预览	取消完成
数据响应结果 点击复制	
1 { 2 { 3 "date": "2018/12/1", 4 "amount": "1000" 5 }, 6 {	
	e la

添加过滤器后,字段便会匹配成功。



] 说明:

但此时的折线图并没有正确展示,由于API返回的日期格式与组件默认的格式不一样,因此还 需要设置一下折线横轴的日期格式。

切换至配置面板,在x轴>轴标签中选择数据种类为时间型,数据格式选择本API所返回的格式2016/01/01,即可看见折线图的正常展示。

۲.																					0
图层			o t															ŧ			Q
1				-														基本折线图			
TITLE	通用标题_P4								某某	「工厂车前	间实时状	況									
201712:24	时间器_e5K7		100		10.01 10.01 10.01			≝⊟r≓≣ ani			A7884		CTERES		Hor Por Polation			✔ ×铀			۲
- 24	单图片_ESyX		200 30				529, 3	2324	₩ 2 5 _{, 196}		* 75		98%)				> 文本			•
- 24	单图片_PjEb		0 40		831											9		数据种类	时间型		\$
24	单图片_iaYb		500		derlars Parto	•	-							15.88.12.10				数据格式 ⊙	2016/01	/01	\$
TITLE	通用标题_Qv		600	14 - 1770 -						RAL BASS MUS				. 1085 (082	-			显示格式 ②	01/01(月	/日)	\$
TITLE	通用标题_ud		700		=	in in	1											留白			
0	基础款饼图_]		800										.1			20		留白距离	0 最小值		1 最大值
TITLE	通用标题_Q>		900	406A.848 (889:403)	CC308887C13 87C13 34⇒5 84⇒5413 [H 1433 [H108	133、14、131、14、1 133、1433日〜厂集 33日〜厂鉄会総会の へ口知道40日本大学家			.1.1		ы.		di i					最大值	auto		
TITLE	通用标题_ix2	уК	1000						RURI RURI RURI R	1.47 1.44 1.474 I.A	837 (0833 (0832 (07	NAR FRANK FRANK AN AN	817 00818 00811		5 53 550	10 200 250		最小值	auto		
TITLE	通用标题_ij9		1100															位移			_
TITLE	通用标题_Qt	FUo	1200															数量			
H	多行文本_KD	g98	1300														-	角度	请选择		\$
24	单图片_K6Lg		1400															> 轴线			ø
-															8		±₽				ø

至此,便完成了通过数据服务将MaxCompute表生成API,然后在DataV数据大屏中进行展示的 所有操作,效果如下图所示。

	设备运行状态 Equipment Running Status	某某工厂车间实时状况 ◎ 2019-01-22 17:45:31	▲ 车间温度:23*
ACHING 1 97 100 <	10.0% 0% 16.7% 15.7% 15.3% 15.3%	当日产量 ㎡ 2 3 2 4 2 5 計却产量 196020 75% 98% 98% 成文金額増长超多 Trade Statistics	Service of the servic
#9	Kall Workers Performance 1 1 2 1 2 1 2 1 2 1 2 1 2 1 3 1 3 1	10000 4000 2000 1000 1000 1000 1000 1000 1000 1000 1000	质量控制 Quality Control
[H338->厂物冷却水入水温度(*C)] 34->14.13 [H338- 300 >厂物冷却水入水温度(*C)] 34->14.13 [H338->厂物冷 300 かた入水温度(*C)] 34->14.13 [H338->厂物冷 500	用?	机器产量统计 production	4UR3 300 9UR4 90 9UR5 200 4UR5 200
規度(*C) 34->14.33 (+338->「労務抑化人水規模(1)] 34- [*C] 34->14.33 (+338->「労務抑化人水規模(*C)] 34- 100 (***********************************	[H338->厂务/冷却水入水温度[['C]] 34->14.13 [H338- >厂务/冷却水入水温度[['C]] 34->14.13 [H338->厂务/冷却水入水温度[['C]] 34->14.13 [H338->厂务/冷却水入水 温度['C]] 34->14.13 [H338->厂务/冷却水入水温度 ['C]] 34->14.13 [H338->厂务/冷却水入水温度[['C]] 34-		

注意事项

DataWorks数据服务与DataV进行无缝对接后,则不需要使用DataV中的API数据源去填写一个 URL调用API,直接新建一个DataWorks数据服务作为数据源,便可直接选用数据服务中的API ,无需每个API都设置AppKev和AppSecret认证信息,且支持通过表单填写API参数,操作快捷

通过数据服务,您可以将MaxCompute中加工好的数据结果,直接在DataV中进行呈现,实现数据开发-数据服务-数据分析展现的全链路开发。在开发过程中,请注意以下事项。

- · DataWorks数据服务向导模式生成API只支持单表简单条件查询,脚本模式支持用户编写查询 SQL语句,支持多表关联查询、函数以及复杂条件。您可以根据自己的需求灵活选择。
- Lightning采用的是PostgreSQL语法,在编写SQL时,需要注意使用PostgreSQL函数,而不 是MaxCompute的UDF。目前Lightning仅支持MaxCompute的UDF max_pt,可用于获取 当前最新分区。连接字符串时使用||。
- · Lightning目前只支持秒级查询,并且查询的MaxCompute不宜过大(控制在GB级),尽量 将分区作为请求参数,尽量避免扫描过多分区,否则会比较慢。
- ·如果您要求毫秒级API查询,则建议采用关系型数据库、NoSQL数据库或AnalyticDB作为数据 源。
DataV组件要求的数据格式是个数组,数据服务生成的API返回结果是带有错误码的完整JSON,因此要使用过滤器对API结果进行处理。您可以选择在DataV中添加过滤器,也可以选择直接在数据服务配置API时添加过滤器。

通常对于未分页查询的API,直接返回data数组即可,对于分页查询的API直接返回data.rows数组。

・若您要在DataV的折线图或柱状图中添加多个系列,通常DataV要求每个系列的数据是一个对象,并通过字段s来区分系列,此时要注意使用过滤器进行格式转换。

2.7 天依赖分钟任务最佳实践

本文将为您介绍每5分钟抽取一次数据,待每天00:00的同步任务抽取完成后,对当天共288次同步 任务抽取的所有数据进行计算。

实现思路

- 1. 创建一个同步任务为上游, 一个SQL任务为下游。
- 2. 将同步任务的调度时间设置为每5分钟调度一次(开始时间00:00,结束时间23:59,时间间隔5 分钟)。
- 3. 配置依赖上一周期-本节点,以形成自依赖。
- 4. 将SQL任务设置为每天00:00调度一次。

实现原理

在DataWorks调度系统中,下游对上游的依赖遵循原则为:下游任务生成的实例会找到当天离自己 最近结束的一个上游实例作为上游依赖,如果上游依赖实例运行成功,才会触发本节点实例运行。 如果上游节点每天生成多个实例,则下游无法识别是哪一个实例离它最近结束,导致必须等上游当 天生成的所有实例运行完成后才会运行。因此,上游必须配置自依赖,SQL任务在00:00的实例才 会准确依赖00:00生成的同步任务实例结束后再运行。

前提条件

开始本实验前,需要首先准备好下述内容。

- ・请确保已拥有阿里云云账号并进行实名认证。详情请参见准备阿里云账号。
- · 创建数据库,并准备好源端数据和目标端表。

步骤一:创建业务流程

1. 登录DataWorks控制台,单击相应工作空间后的进入数据开发。

2. 选择新建 > 业务流程。



3. 填写业务名称和描述,单击新建。

4. 新建节点并配置依赖关系。

∱	⊙		1	»					
~ 节	点组		С						
TT IT		Đ							
~ 数	据集成								
回数	据同步								
~ 数	据开发								
Sc OE	PS Scrip	ot							
Sq OL	PS SQL						Di	rds mysal 2 odps	
<u>ි</u> so	L组件节	点							
Sp OE)PS Spar	ĸ							
Ру Ру	ODPS								
vi 🛓	拟节点								
Mr OE	OPS MR						Sq	workshop_odps	2
Sh Sh	ell								
🔞 Da	ta Lake	Analytic	s						
_									
	说明	:							
节点1	用来料	将每5 分	分钟调	周度一次的My	ySQL数据	同步至Max	xCor	mpute。	
节点2	用来》	C总M	laxCo	ompute接收	的数据。				

步骤二: 配置分钟调度任务

1. 配置数据来源和数据去向。

上游为数据同步任务,将MySQL的数据同步至MaxCompute。根据过滤条件过滤每5分钟更新的数据,目标端的分区根据定时时间的前5分钟创建,保证所有数据都写到同一天的分区中。

Sq workshop_odps ×	DI rds_mysql_2_odps ×	嚞 workshop 🛛 🗙						
		Ø						
01 选择数据源	数	据来源			数据去向			
	在这里配置	数据的来源端和写入端;	可以是默认的数据源	,也可以是您创建的自	有数据源查看支持的数据	居来 源类型		
* 数据源:	MySQL ~	workshop_test000	· ?	* 数据源:	ODPS ~	odps_first		?
*表:	`workshop_test` ×			* 表:	workshop_odps_mi			
		添加数据	ह +				一键生成目标表	
数振过滤:	insert_time>=\${startTime} insert_time<\${endTime}	and	0	* 分区信息:	ds = \${startTime}	?		
				清理规则:	写入前清理已有数据 (In	isert Overwrite)		
切分键:	id		0	空字符串作为null:	〇 是 💿 否			
	数据	顾览						

过滤条件为insert_time>=\${startTime} and insert_time<\${endTime},在右侧调 度配置中,配置参数过滤每间隔5分钟的数据。

2. 字段映射。

源端和目标端都创建了id、name、insert_time三列字段,insert_time为时间列,数据可以 根据时间来过滤。



3. 单击右侧的调度配置,进入调度配置页面进行设置。

为过滤条件中的参数赋值: startTime=[yyyymmddhh24miss-5/24/60] endTime=[yyyymmddhh24miss],每个参数间用空格分隔。

×			Q
基础属性②			配置
节点名:	rds_mysql_2_odps	节点ID: 7(版本
节点类型:	数据同步	表任人: heimei haogtafiyan innen zam	
描述:			
参数:	bizdate=\$bizdate startTime=\$[yyyymmddhh24miss-5/24/60] endTime=\$[y	yyymmddhh24miss]	

开始时间是从00:00点开始,间隔5分钟调度一次,为保证一天的实例都运行完,需要设置自依赖。

时间属性 ?	l
生成实例方示	式: 💿 T+1次日生成 🔵 发布后即时生成
时间属	生: 🧿 正常调度 💿 空跑调度
出错重	ಷೆ: 🗌 🕐
生效日期	期: 1970-01-01 9999-01-01
	注:调度将在有效日期内生效并自动调度,反之,在有效期外的任务将不会自动调度,也不能手动调度。
暫停调	ĝ: 📃
调度周期	明: 分钟 ~
定时调	度: 🔽
开始时间	目: 00:00 ①
时间间	隔: 05 ① 分钟
结束时间	刊: 23:59 ①
cron表达	式: 00 */5 00-23 * * ?
依赖上一周期	
依赖耳	项: 本节点 · · · · · · · · · · · · · · · · · · ·

📕 说明:

- ・上述配置可以保证您一天产生的实例能够依次运行完成,并保存至MaxCompute表同一天的分区中。
- ·如果您的分钟任务中有一个调度任务出错,则设置自依赖后面的实例都不会运行,需要您手动进行处理。

· 为避免出现上述问题,您可以设置数据过滤为insert_time<\${endTime},每次都进行全量同步,只要有成功的便可将endTime数据进行同步,您不需设置自依赖,但会增加数据库的负担。

步骤三: 配置天调度任务

下游是一个SQL节点天任务,将workshop_odps_mi一天分区中的数据都过滤出来插入 到workshop_odps_dd表中。

odps sql
author:
create time:2019-07-31 15:59:11
<pre>insert overwrite table workshop_odps_dd partition (ds=\${yestoday})</pre>
<pre>select id,name,insert_time from workshop_odps_mi where \${startTime}<=ds and ds<\${endTime};</pre>

insert overwrite table workshop_odps_dd partition (ds=\${yestoday})
select id, name,insert_time from workshop_odps_mi where \${startTime}<=
ds and ds<\${endTime};</pre>

```
调度配置中的赋值为startTime=$[yyyymmddhh24miss-1] endTime=$[yyyymmddhh
```

24miss] yestoday=\$[yyyymmdd-1]。

将workshop_odps_mi一天分区中的数据都过滤出来插入到workshop_odps_dd表中。如下所示:

```
insert overwrite table workshop_odps_dd partition (ds=20190320)
select id, name,insert_time from workshop_odps_mi where 20190320000000
<=ds and ds<20190321000000;</pre>
```

因为天运行定时时间在分钟任务运行结束后,所以插入workshop_odps_dd分区ds时间要和分钟 任务时间在同一天,您将ds时间减1即可。

单击右侧的调度配置,进入调度配置页面。

×			Q
基础属性 ②			配置
节点名:	workshop_odps	节点D:	血缘
节点类型:	ODPS SQL	责任人: v	关系
描述:			版本
参数:	bizdate=\$bizdate_startTime=\$[yyyymmddhh24miss-1] endTime=\$[yyyy	mmddhh24miss] yestoday=\$[yyyymmdd-1] ⑦	<i>/+</i>
			垣构

步骤四: 查看运行结果

同步任务配置完成后,单击运行,即可查看运行结果。



2.8 邮件外发最佳实践

本文将为您介绍如何通过PyODPS节点结合独享资源组的方式,实现邮件外发的需求。

背景信息

DataWorks PyODPS节点和Python脚本并不相同,PyODPS节点主要用于和MaxCompute交互 进行数据分析处理。您可以通过PyODPS节点结合独享资源组的方式,实现从MaxCompute拉取 数据进行邮件外发的场景需求。

通过PyODPS节点+独享资源组的方式实现邮件外发

- 1. 创建独享资源组,独享资源组和DataWorks工作空间的地域保持一致。详情请参见独享资源模
 - 式。

				概览	工作空间列	表资	源列表	计算引擎列	ŧ			
华北2 华东1 华纬	マ2 华南1 香港 💈	é西1 亚太东南 1	美东1 亚太东北1	欧洲中部 1	亚太东南 2	亚太东南 3	中东东部 1	1 亚太南部 1	亚太东南 5	英国		
独享资源 公共资	源 自定义资源组											
	授索											
资源名称	备注	类型	状态	到期时间		资源	ž i	资源使用率	操作			
1.000	foregame.	调度资源	✔ 运行中	2019-07-25 00	:00:00	1	(D	查看信息	扩容 缩容 续费	专有网络绑定	修改归屋工作空间



资源组创建成功后,单击相应资源组后的修改归属工作空间,将其指派给相应的工作空间。

2. 在数据开发面板创建PyODPS节点,填写SMTP发送代码并保存,示例如下:

```
import smtplib
from email.mime.text import MIMEText
from odps import ODPS
mail_host = '<yourHost>' //邮箱服务地址
mail_username = '<yourUserName>' //登录用户名
mail_password = '<yourPassWord>' //登录用户密码
```

```
mail_sender = '<senderAddress>' //发件人邮箱地址
mail_receivers = ['<receiverAddress>'] //收件人邮箱地址
mail_content=""
                          //邮件内容
o=ODPS('access_key', 'access_secretkey', 'default_project_name',
endpoint='maxcompute_service_endpoint')
with o.execute_sql('query_sql').open_reader() as reader:
            for record in reader:
                     mail_content+=str(record['column_name'])+' '+
record['column_name']+'\n'
message = MIMEText(mail_content,'plain','utf-8')
message['Subject'] = 'mail test
message['From'] = mail_sender
message['To'] = mail_receivers[0]
try:
            smtpObj = smtplib.SMTP_SSL(mail_host+':465')
            smtpObj.login(mail_username,mail_password)
            smtpObj.sendmail(
                 mail_sender,mail_receivers,message.as_string())
            smtpObj.quit()
print('mail send success')
except smtplib.SMTPException as e:
            print('mail send error',e)
```

3. 提交PyODPS节点,进入运维中心页面修改任务调度资源组为独享资源组。

6	🤣 运维中心	~				₽ DataSti	udio 🖏 💎
()	医维大屏	搜索 节点名称/节点D Q 解决	方案: 请选择 > 业务流程:	j>选择 > 节点类型	请选择 > 责任人	请选择责任人 >	
-	任务列表	基线 请选择 ~ □ ≸	戏的节点 今日修改的节点 暂停 (冻	结)节点 重置 清空			
6	周期任务						○ 刷新 收起搜索
1	手动任务	名称	节点ID 修改日期↓	任务类型 责任	人 调度类型	资源组 🎧	操作
		send_email_test	700002529558 2019-06-24 20:25:27	PY_ODPS	日调度	ALC: NO. OF THE OWNER OF THE OWNE	DAG图 測试 补数据 ▼ 更多 ▼
-	出方道理						

4. 测试运行PyODPS节点。

搜索: 节点名称/节点D Q 节点类型 请选择	▼ 责任人:	请选择责任人	∨ 运行日期	2019-06-24	曲业务	日期: 请选择日期
运行状态: 调选择 > 基线 调选择	✓ □ 我	的节点 📃 我名	天测试的节点	暂停(冻结)节点	重置	青空
基本信息	属性	上下文	运行日志	操作日志	代码	
 Send_email_test #700002529558 06-24 20-25:42 ~ 20:25:53 (dur 11s) 	⊘ 2019-06-24 20 2019-06-24 20 持续时间:11 Gateway:cn-	:25:42 ~ :25:53 s	begin to connect begin to login s begin to cond or mail send succes 2019-06-24 20:25	t to smtp server smtp server		
	shanghai.3427	749826073218.1	2019-06-24 20:25 2019-06-24 20:25 2019-06-24 20:25 2019-06-24 20:25	5:52 INFO Exit c 5:52 INFO In 5:52 INFO Shell 5:52 INFO Curren	ode of the She vocation of Sh run successful t task status:	<pre>11 command 0 ell command completed ly! FINISH</pre>

总结

通过PyODPS节点+独享资源组的方式实现邮件外发时,独享资源组用户无法登录到对应的机

器,导致无法安装更多Python第三方模块,实现更多的功能。

2.9 PyODPS节点实现结巴中文分词

本文为您介绍如何使用DataWorks的PyODPS类型节点,借助开源结巴中文分词包实现对中文字 段的分词并写入新的表。

前提条件

·请首先确保您已经完成DataWorks工作空间的创建,本例使用简单模式的工作空间,详情请参见创建工作空间。完成工作空间创建后,单击进入数据开发。

= (-)阿里云		Q. 搜索	要	用 工单 备案	企业 支持与服务		简体中文
		概览 工作空间列表	资源列表 计算引擎	列表			
华东1 华东2 华南1 华北2 雪港 授	美西1 亚太东南1 美东1 欧洲中部1 案	亚太东南 2 亚太东南 3 亚太东	に北1 中东东部1 亚太南部	1 亚太东南 5 英	12	创建工	作空间 剧新列号
工作空间名称/显示名	模式	创建时间	管理员	状态	开通服务	操作	
Manufacture Control of	10000-01000-0000	10 x 1 10 10 10 10 10	dista de c	正常	0. 🔨	工作空间配置 进入数据开发 进入数据集成 进入数据服务	修改服务 更多 ▼
The second s	1000-010	10103-0010-010-0	101.01	正常	∞ 🔨	工作空间配置 进入数据开发 进入数据集成 进入数据服务	修改服务 更多 ▼
	000010000	100000 (Color)	1000 c. de c	正常	∞ 🔨	工作空间配置 进入数据开发 进入数据集成 进入数据服务	修改服务 更多 ▼
0.0000	1200-2200-200	100000-000	1000,000	正常	∞ 🔨	工作空间配置 进入数据开发 进入数据集成 进入数据服务	修改服务 更多 ▼
workshop_DOC workshop_DOC	简单横式(单环境)	2019-07-02 16:31:51	dtplus_docs	正常	∞ 🔨	工作空间配置 进入数据开发 进入数据集成 进入数据服务	修改服务 更多 ▼

·请在GitHub下载开源结巴分词中文包。

fxsjy / jieba			• Watch 1,253	★ Star 1	19,220 ¥ F	ork 5,0
Code 🕕 Issues 44	8 🕅 Pull requests 34 🔳 I	Projects 0 💷 Wiki 🕕 Secu	urity 🔟 Insights			
	GitHub is home to review cod	Join GitHub today o over 36 million developers worki de, manage projects, and build so Sign up	ng together to host and ftware together.			Dismiss
巴中文分词 ⑦ 500 commits	پ 2 branches	♥ 24 releases	2 35 contributor	'S	MIT الله	
巴中文分词 ⑦ 500 commits anch: master • New 1	۶ 2 branches ull request	⊗ 24 releases	🎎 35 contributor	s Find Fil	কুঁ MIT le Clone or	download
巴中文分词 ② 500 commits ranch: master • New p New p	ৃ ዖ 2 branches iull request suggest_freq中add_word指向的bug (#	© 24 releases	1 35 contributor Clone with H	s Find Fil ITTPS ③	গ্রুঁত MIT le Clone or	download
巴中文分词 ② 500 commits ranch: master New Iinhx13 and fxsjy 修复 extra_dict	ダ 2 branches full request suggest_freq申add_word指向的bug (# update to v0.33	♥ 24 releases 723)	2 35 contributor Clone with H Use Git or check	s Find Fil ITTPS ③ cout with SVN	호 MIT e Clone or using the web	download
巴中文分词 ⑦ 500 commits ranch: master New Iinhx13 and fxsjy 修复 extra_dict Jieba	》2 branches null request suggest_freq中add_word指向的bug (# update to v0.33 修复suggest_freq中add_word	© 24 releases 723) 指向的bug (#723)	Clone with H Use Git or check	S Find Fil ITTPS ③ cout with SVN ub.com/fxsjy/	∯ MIT ie Clone or using the web	download URL.
巴中文分词 ⑦ 500 commits ranch: master ▼ New p 2] linhx13 and fxsjy 修复 extra_dict jieba test	P 2 branches Suggest_freq中add_word指向的bug (# update to v0.33 修复suggest_freq中add_word fix the error about imoprting 0	◇ 24 releases 723) 指向的bug (#723) ChineseAnalyzer	L 35 contributor	s Find Fil ITTPS ③ cout with SVN ub.com/fxsjy/	화 MIT Ie Clone or using the web 'jieba.git	download URL.

背景信息

PyODPS集成了Maxcompute的Python SDK。您可以在DataWorks的PyODPS节点上直接编 辑Python代码并使用Maxcompute的Python SDK。关于PyODPS节点的详情请参见PyODPS节 点。

操作步骤

- 1. 创建业务流程。
 - a) 右键单击业务流程, 选择新建业务流程。



b) 输入您的业务流程名称后, 单击新建。

新建业务流程		×
业务名称:	jiebafenci	j _i
描述:	请输入业务描述	

- 2. 上传jieba-master.zip包。
 - a) 右键单击资源,选择Archive。



b) 上传您已下载到本地的jieba-master.zip, 勾选上传为ODPS资源, 单击确定。

新建资源			×
资源名称:	jieba-master.zip		
目标文件夹:	业务流程/jiebafenci/资源		~
资源类型:	Archive		
	✓ 上传为ODPS资源本次上传,资源会同步上传至ODPS中		
上传文件:	jieba-master.zip (11.83M)		×
		确定	取消

c) 提交资源。

e الح	
提交上在次酒	
——————————————————————————————————————	
已保存文件:	jieba-master.zip
资源唯一标识:	OSS-KEY-
	☑ 上传为ODPS资源本次上传,资源会同步上传至ODPS中
重新上传:	点击上传

- 3. 创建测试数据表。
 - a) 右键单击表,选择新建表。输入表名jieba_test。

~	🚣 jiebafenci
	▶ <mark> ຸ</mark> 数据集成
	>
	> # ±
	新建表

b) 单击DDL模式, 输入建表DDL语句如下。

本教程准备了两列测试数据,您在后续开发过程中可以选择一列进行分词。

CREATE TABLE	`jieba_test`	(
`chinese`	string,	
`content`	string	

);

c) 单击提交到生产环境。

🧮 jieba_test 🗙 🛄 表	Ar jieba-ma	aster.zip	🛃 jiebafenci		
DDL模式 从生产环境加载	提交到生产环境				
	表名	jieba_test			
	入该表的业务流程	jiebafenci			
基本属性					
	jieba_test				
一级主题	请选择			二级主题	请选择

4. 创建测试结果存放表。

本例仅对测试数据的chinese列做分词处理,因此结果表仅有一列,创建方法同上。DDL语句 如下所示。

```
CREATE TABLE `jieba_result` (
`chinese` string
```

-);
- 5. 上传测试数据。

本例已为您准备好分词测试数据,请点击此处下载。

a) 单击导入。



b) 输入测试数据表名jieba_test, 单击下一步。

数据导入向导			×
选择要导入数据的表: 名称	jieba jieba_result jieba_test		
		没有数据	
			下一步 取消

c) 单击浏览, 上传您下载到本地的jieba_test.csv文件, 单击下一步。

数据导入向导			×
选择数据导入方式:		● 来自数据分析的电子表格	
选择文件: 前		问题。 只支持.txt、.csv和.log文件类型	
选择分隔符:	通号 🗸 🔿		
原始字符集: 0	GBK		
导入起始行: 1	1		
首行为标题: 🗸	2		
数据预览			
Chinese		Content	
数据库备份是为数据库提供证	连续数据保护低成本的备份服务	数据库备份是为数据库提供连续数据保护低成本的备份服务	
数据库备份拥有一套完整的	数据备份和数据恢复解决方案	数据库备份拥有一套完整的数据备份和数据恢复解决方案	
可以通过简单的配置实现数据恢复	關库全量备份增量备份以及数据	可以通过简单的配置实现数据库全量备份增量备份以及数据恢复	
		上一步	取消

d) 勾选按名称匹配,单击导入数据。

数据导入向导	×
选择目标表字段与源字段的匹配方式: 🦳 按位置匹配 💽 按谷	名称匹配
目标字段	源字段
chinese	Chinese
content	Content
	上一步

- 6. 创建PyODPS节点。
 - a) 在业务流程中右键单击数据开发,选择新建数据开发节点 > PyODPS。

▼ 业务流程			
> 🚣 Worksho	р		
🗸 🛃 jiebafeno	i		基本属性
> 😑 数据	耒成		
➤ 🗤 数	新建数据开发节点	>	ODPS SQL
✔ 🛄 表	新建文件夹		ODPS Script
	看板		ODPS Spark
		. [PyODPS
	snc	虚拟节点	

- b) 输入节点名称word_split。
- c) 输入您的PyODPS代码。

输入代码如下,释义请见代码注释。

```
def test(input_var):
    import jieba
    import sys
    reload(sys)
    sys.setdefaultencoding('utf-8')
    result=jieba.cut(input_var, cut_all=False)
    return "/ ".join(result)
hints = {
    'odps.isolation.session.enable': True
}
libraries =['jieba-master.zip'] #引用您的jieba-master.zip压缩包。
iris = o.get_table('jieba_test').to_df() #引用您的jieba_test表中的
数据。
```

```
example = iris.chinese.map(test).execute(hints=hints, libraries=
libraries)
print(example) #查看分词结果, 分词结构为MAP类型数据。
abci=list(example) #将分词结果转为list类型数据。
i = 0
for i in range(i,len(abci)):
    pq=str(abci[i])
    o.write_table('jieba_result',[pq]) #通过循环, 将数据逐条写入您的结
果表jieba_result中。
    i+=1
else:
    print("done")
```

d) 运行代码进行测试。

	Sq query	y • • • word_split ×
		F L 🔂 🖸 🗉 C 🗹 🗏 🗱 :
		def test(input_var):
		import jieba
		import sys
		reload(sys)
		sys.setdefaultencoding('utf-8')
		result=jieba.cut(input_var, cut_all=False)
		<pre>return "/ ".join(result)</pre>
		hints = {
		'odps.isolation.session.enable': True
		}
		libraries =['jieba-master.zip']
		iris = o.get_table('jieba_test').to_df()
		<pre>example = iris.chinese.map(test).execute(hints=hints, libraries=libraries)</pre>
		print(example)
D	16	abci=list(example)

e) 您可以在运行日志查看结巴分词的程序运行结果。

运行日志

数据库/备份/是/为/数据库/提供/连续/数据保护/低成本/的/备份/服务数据库/备份/拥有/一套/完整/的/数据备份/和/数据恢复/解决方案
可以/通过/简单/的/配置/实现/数据库/全量/备份/增量/备份/以及...
为了/节省成本/,/可以/选择/多种/OSS/存储/类型/进行/存储
在/进行/数据恢复/时/,/可以/使用/存储/的/增量/备份/实现/...
为了/降低/在/故障/发生/后/数据/丢失/,/数据库/备份/DBS/...
出于/安全/合规/要求/,/部分/数据/需要/长期/保存
作为/完整/数据库/灾备/方案/,/除了/要/有/本地/数据库/备份/...
数据库/备份/DBS/提供数据/全量/备份/增量/备份/和/数据恢复

7. 在数据开发创建一个ODPS SQL类型节点, 输入select * from jieba_result, 单击运





8. 查看运行结果,验证数据是否已写入结果表中。

6	sele	ect *	fro	om jie	ba_res	ult;						
运	行日志			结果[1]	×							
						А						
1	chinese											~
2	chinese	为了。	/ 降低	/在/故	障/ 发生/	「后/ 数据	/丢失/	, / 数据	库/备份)/ DBS/	Name	e: 5
3	chinese	作为	/ 完整	ダ 数据库	訂 灾备/ フ	5案/ ,/	除了/ 雾	》有/本	地/ 数据	裤/备 份	}/ Na	ime
4	chinese	数据	库/	份/拥有	訂/一套/ Э	完整/的/	数据备代	分/ 和/	如据恢复	/ 解决方	案 Nar	ne
5	chinese	数据	库/	份/ DBS	3/ 提供数	据/ 全量/	备份/ 堦	鳁/备	分/ 和/ 夎	如据恢复	Name:	: 8,
6	chinese	为了	/ 节省	诚本/ ,	/可以/រ	选择/ 多种	P/ OSS/	存储/ 勢	型/进行	テ/存储	Name:	3, (
7	chinese	可以	/ 通过	/ 简单/	的/ 配置/	实现/ 数	据库/ 全	量/备	}/ 増量/	备份/ じ	J及 N	lan
8	chinese	数据	库/ 备	份/是/	为/ 数据	车/ 提供/ 注	连续/ 数	据保护	/低成本	/的/备	分/ 服务	5 N
9	chinese	在/说	进行/∮	数据恢复	夏/时/,/	可以/使	用/ 存储	矿的/增	量/ 备 份	/ 实现/	Nam	ne:
10	chinese	出于	/ 安全	/ 合规/	要求/ ./	'部分/数	据/ 雲要	/长期/	保存 Na	me: 6. d	tvpe: o	bie

2.10 基于AnalyticDB构建企业数仓

本文将为您介绍如何基于AnalyticDB构建企业数仓,并进行运维和元数据管理等操作。

创建工作空间

- 1. 使用主账号登录DataWorks控制台。
- 2. 使用主账号登录DataWorks控制台。

3. 单击控制台概览 > 常用功能下的创建工作空间。

		概览	工作空间列表
🌀 Dat	aWorks 数据	3 集成・数据开发・ 3	数据服务
快速入口			
新产品推荐: ∭ Stream Studio	MEW 快速入门		
数据开发	数据集成	运维中心	数据服务
工作空间			
10000	华东1	4,000,000	4
创建时间:2019-01-30 10:18:52 计算引擎:MaxCompute PAI计算引 服务模块:Internal System Error With	擎 Code module_code_AppStudio 数	创建时间:2018-09-02 10:26:59 计算引擎:MaxCompute 服务模块:Internal System Error With	Code module_code_AppStudio
工作空间配置	进入数据开发	工作空间配置	进入数据开发
进入数据服务	进入数据集成	进入数据服务	进入数据集成
常用功能 分 创建工作空间 X 一键CC	IN		

您也可以进入工作空间列表页面,单击创建工作空间。

		概览 工作空间列表	资源列表 计算引擎列表			
<u>华东1</u> 华东2 华南1 华北2 香港 美西1	亚太东南1 美东1 欧洲中部1	亚太东南 2 亚太东南 3 亚太东北 1	中东东部 1 亚太南部 1 亚太东南	5 英国		创建工作空间 刷新列表
搜索						
工作空间名称/显示名	模式	创建时间	管理员	状态	开通服务	授作
Marca and Annual State	标准模式(开发跟生产隔离)	2019-01-30 10:18:52		正常	00 🔨 度	工作空间配置 进入数据开发 修改服务 进入数据集成 进入数据服务 更多 ▼
Table 11	简单模式(单环境)	2018-09-02 10:26:59	10.00	正常	∞ 🔨	工作空间配置 进入数据开发 修改服务 进入数据集成 进入数据服务 更多 ▼

4. 填写创建工作空间对话框中的配置项,选择地域、计费方式和服务。

如果选择的地域没有购买相关的服务,会直接显示该地域下暂无可用服务,默认选中数据集成、 数据开发、运维中心和数据质量。

创建工作空间
选择计算引擎服务
✓ MaxCompute 按量付费 包年包月 去购买 开发者版 去购买 开通后,您可在DataWorks里进行MaxCompute SQL, MaxCompute MR任务的开发。
□ 2 机器学习PAI · 按量付费 开通后,您可使用机器学习算法、深度学习框架及在线预测服务。使用机器学习PAI,需要使用MaxCompute
□ 纪 实时计算 ○ 共享模式 ○ 独享模式 开通后,您可在DataWorks里面使用Stream Studio进行流式计算任务开发。
选择DataWorks服务
数据集成、数据开发、运维中心、数据质量 您可以进行数据同步集成、工作流编排,周期任务调度和运维,对产出数据质量进行检查等
取当

选项	配置	说明
选择计算引擎服 务	MaxCompute	MaxCompute是一种快速、完全托管的TB/PB级数据 仓库解决方案,能够更快速为您解决海量数据计算问 题,有效降低企业成本,并保障数据安全。
		 说明: 完成创建Dataworks工作空间后,需要关 联MaxCompute项目,否则现执行命令会报project not found的错误。
	机器学习PAI	机器学习是指机器通过统计学算法,对大量的历史数据 进行学习从而生成经验模型,利用经验模型指导业务。

选项	配置	说明
	实时计算	开通后,您可以在DataWorks使用Stream Studio,进 行流式计算任务开发。
选 择DataWorks服 务	数据集成	数据集成是稳定高效、弹性伸缩的数据同步平台。致力 于提供复杂网络环境下、丰富的异构数据源之间数据高 速稳定的数据移动及同步能力。详情请参见数据集成模 块的文档。
	数据开发	该页面是您根据业务需求,设计数据计算流程,并实现 为多个相互依赖的任务,供调度系统自动执行的主要操 作页面。详情请参见数据开发模块的文档。
	运维中心	该页面可对任务和实例进行展示和操作,您可以在此查 看所有任务的实例。详情请参见 <mark>运维中心</mark> 模块的文档。
	数据质量	DataWorks数据质量依托DataWorks平台,为您提供 全链路的数据质量方案,包括数据探查、数据对比、数 据质量监控、SQLScan和智能报警等功能。详情请参 见数据质量模块的文档。

5. 单击下一步,配置新建工作空间的基本信息和高级设置。

创建工作空间	>
基本信息	
工作空间名称:	需要字母开头,只能包含字母下划线和数字
显示名:	如果不填,默认为工作空间名称
* 模式:	标准模式 (开发跟生产隔离) 🛛 🗸
描述:	
高级设置 * 启动调度周期:	^Ħ ⊘
* 能下载select结果:	# ◎
面向 MaxCompute	
* MaxCompute项目名称:	0
* MaxCompute访问身份:	工作空间所有者 🖸 📀
* Quota组切换:	按量付费默认资源组 🗸
	上一步 创建工作空间

分类	配置	说明
基本信息	工作空间名称	工作空间名称的长度需要在3到27个字符,以字母开 头,且只能包含字母下划线和数字。
	显示名	显示名不能超过27个字符,只能字母、中文开头,仅包 含中文、字母、下划线和数字。

分类	配置	说明
	模式	工作空间模式是DataWorks新版推出的新功能,分为 简单模式和标准模式,双项目开发模式的区别请参见简 单模式和标准模式的区别。
		 · 简单模式:指一个Dataworks工作空间对应一个 MaxCompute项目,无法设置开发和生产环境,只 能进行简单的数据开发,无法对数据开发流程以及表 权限进行强控制。 · 标准模式:指一个Dataworks工作空间对应两个 MaxCompute项目,可以设置开发和生产双环 境,提升代码开发规范,并能够对表权限进行严格控 制,禁止随意操作生产环境的表,保证生产表的数据 安全。
	描述	对创建的工作空间进行简单描述。
高级设置	启用调度周期	控制当前工作空间是否启用调度系统,如果关闭则无法 周期性调度任务。
	能下载select结果	控制数据开发中查询的数据结果是否能够下载,如果关闭无法下载select的数据查询结果。
	MaxCompute项目名称	默认与DataWorks工作空间名称一致。
	MaxCompute访问身份	推荐使用工作空间所有者。
	Quota组切换	Quota用来实现计算资源和磁盘配额。

6. 配置完成后,单击创建工作空间。

工作空间创建成功后,即可在工作空间列表页面查看相应内容。

配置AnalyticDB数据源

- 1. 单击相应工作空间操作栏中的进入数据集成。
- 2. 选择同步资源管理 > 数据源,单击新增数据源。

⑤ ○ 数据集成		•						থ্	
=	数据源类型: 全部		> 数据源名称:			C刷新	多库多表搬迁	批量新增数据源	新鴬数据渡
				由田教探護的开发环境都要信息 (4名发去到生产环)	自法行时会使用生产环境配置信息				
▲ ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ● ●	数据源名称	数据源类型	御接信息	数据源描述	创建时间 冻通状态	许通!	时间 适用现	這 爆作	选择
→ 同步资源管理			Endpoint	connection	2010/07/02				
▲ 数据源	a da a farat	00.00	项目名称	from odps calc engine 81548	17:10:26		开发		
☆ 資源組	odps_tirst	ODPS	Endpoint: 10日名称:	connection from odos celo	2019/07/02		牛产		
✓ 批量上云				engine 81547	17:10:13		2		

3. 在新建数据源弹出框中,选择数据源类型为AnalyticDB(ADS)。

4. 填写AnalyticDB数据源的各配置项。

新增AnalyticDB (ADS)	数据源	×
* 数据源名称:	自定义名称	
数据源描述:		
* 适用环境:	✔ 开发 生产	
* 连接Url :	格式:Address:Port	
* 数据库:		
* AccessKey ID :		?
* AccessKey Secret :		
测试连通性:	测试连通性	
	上一步	完成
首 说明:		
 ・ 请配置外网IP。 		

- ·如果使用的是AnalyticDB2.0版本,通过用户AK信息进行身份验证。
- ·如果使用的是AnalyticDB3.0版本,通过数据库的用户名和密码进行身份验证(开通ADB3. 0数据库后,首先在控制台创建用户和密码)。
- 5. 单击测试连通性。
- 6. 测试连通性通过后,单击完成。

提供测试连通性能力,可以判断输入的信息是否正确。

设置白名单

・ 设置DataWorks白名単

您需要在DataWorks中,将AnalyticDB的连接URL和端口设置为白名单,调度引擎方可连接 上AnalyticDB数据库。

1. 单击右上角的工作空间管理按钮,进入工作空间配置页面。

⑤ 0.数据	車成	11222	•								থ	▼
=		新田海米田・今年		い、教伝道な物・				Г	(1 B) + (2 # 4		工作空间带	
- 任务列表		MANDOLE . LEAP		. MAINUR 177-		-		L	0 10/10	100		1000000
🕛 高线同步任	5			● 标准项目模式下,配置任务均	他用数据源的并发外现配宣信息,也	土势友布到生产环境	UZITU SUUREr	"林硯配宣信思				
	1	数据源名称	数据源类型	链接信息		数据源描述	创建时间	连通状态	连通时间	适用环境	操作	选择
↑ 数据源			0000	Endpoint: 项目名称:		connection from odps calc engine 81548	2019/07/02 17:10:26			开发		
令 資源組		odps_tirst	ODPS	Endpoint: 简目文表·		connection from odes calc	2019/07/02			生产		
🛃 批量上云				ALLEY		engine 81547	17:10:13			10		

2. 在沙箱白名单(配置shell任务可以访问的IP地址或域名)下,单击添加。

G DataWorks	••				<i>₹</i> , ⊽
E ② I作空间配置	显示名:		能下载select结果:		
基 成员管理	负责人: dtplua_docs ~		启用调度周期:		
	状态正常		允许子账号变更自己的节点责任人:)	
√ MaxCompute高级配置	播述:				
	沙箱白名单(配置shell任务可以访问的IP地址或域名)				添加
	IP地址	第日		操作	

3. 在添加沙箱白名单对话框中,填写地址和port。

添加沙箱白名单		×
* 地址:	.ads.aliyuncs.com	
* port :	3306	
	取消	确定

4. 单击确定。

・设置AnalyticDB白名单

由于AnalyticDB3.0版本基于用户名密码访问,因此需要设置客户端白名单,才允许连接数据 库。

1. 获取DataWorks白名单

为了能让DataWorks gateway请求AnalyticDB3.0,需要将DataWorks的机器IP设置 为AnalyticDB3.0的白名单(AnalyticDB2.0不需要设置)。DataWorks文档为您提供了 各地域对应的白名单,您根据自身地域进行复制即可,详情请参见添加白名单。

- 2. 设置AnalyticDB白名单
 - a. 登录AnalyticDB3.0控制台,进入集群列表 > 数据安全页面。

= (-)阿里云	华北2(北京) ▼	Q 搜索
AnalyticDB 控制台	集群信息	〇 tongguyan 210 (运行中)
集群列表	账号管理	
告警规则	数据安全	集群属性
	监控信息 •	集群ID
	备份恢复	版本 3.0
		集群类型常规

b. 单击添加白名单分组,将复制的DataWorks白名单粘贴至AnalyticDB中。

集群信息 账号管理	🔷 tonggua	n-2ze (运行中)	这回集群列表		2 登录数据库 2 操	作指引 〇 刷新		
数据安全 监控信息	数据安全							
备份恢复	网络隔离模式:通用白名单模式,以下白名单不区分经典网络及专有网络。							
						ビ修改 面删除		
	100.106.48.0/24	47.93.110.0/24	100.64.0.0/8	11.193.99.0/24	10.152.167.0/2	24		
	47.94.185.0/24	11.193.75.0/24	47.95.63.0/24	10.152.168.0/24	11.193.50.0/24	1		
	11.193.82.0/24	11.197.231.0/24	11.195.172.0/24	182.92.144.0/24	47.94.49.0/24			
						□修改 ☆清空		

新建业务流程

1. 单击左上角的图标,选择全部产品 > DataStudio(数据开发)。

2. 右键单击业务流程,选择新建业务流程。

Data	DataStudio	DataWroks演示项目 💙 🛛 🗸	
	=	数据开发 2 🗟 🕻 С 🕀 🗗	
iii)	数据开发 1	文件名称/创建人	Ē
*	组件管理	➤ 解決方案	
R	临时查询	> 业务流程 2 新建业务流程	3
Ē	运行历史	全部业务流程看板	Z
Ň	手动业务流程 New		
#	公共表		
R	表管理		
£×	函数列表		
Ū	回收站		

- 3. 在新建业务流程对话框中,填写业务流程名称和描述。
- 4. 单击新建。

创建数据同步任务

1. 右键单击新建业务流程下的数据集成,选择新建数据集成节点 > 数据同步。



2. 在新建节点对话框中,填写节点名称,单击提交。

新建节点		×
节点类型:	数据同步	
节点名称:	1.0.40.00	
目标文件夹:	\$155-CB\$18	
	提交	取消

3. 选择数据来源和数据去向,单击下一步。

۳	\odot	Þ	♪	٤.		<u>_</u>	 ひ						
01	选择数据	部原					数据来源			数据去向			
在这里配置数据的未源操和写入端;可以是默认的数据源,也可以是您创建的自有数据源查看支持的数据未源类型													
		* 数据	源 M	ySQL			<pre>v rds_tongguan_adb_test </pre>	?	* 数据源	ADS ~	adb_dw_tongguan	?	
			表	product						product			
									*导入模式	实时导入			
		数据过	滤 <mark>pr</mark>	oid>100				?	* 批量插入条数	5000			
		切分	键 pr	oid				?					
							数据预览						

4. 字段映射。

选择字段的映射关系。左侧的源头表字段和右侧的目标表字段为一一对应关系。单击添加一行可以增加单个字段,鼠标放至需要删除的字段上,即可单击删除图标进行删除。

02 字段映射		源头表		目标表			
							同夕咖餅
	源头表字段	类型	Ć		目标表字段	类型	同行映射
	proid	INT	•		proid	bigint	取消映射
	proname	VARCHAR	•		proname	varchar	
	price	DOUBLE	•	••	price	double	
	prodesc	VARCHAR	•		prodesc	varchar	
	添加一行 +						

5. 通道控制。

单击下一步, 配置作业速率上限和脏数据检查规则。

03 通道控制			
		您可以配置作业的传输速率和错误纪录数来控制整个数据同步过	程:数据同步文档
•任务期望最大并发	发 3	~ 0	
* 同步速	🚈 🔵 不限流 🧿 限流 📘	0 MB/s	
错误记录数超	过 脏数据条数范围,默认允许脑	数据	条,任务自动结束 🥎
任务资源	且 默认资源组		

配置	说明
任务期望最大并发数	数据同步任务内,可以从源并行读取或并行写入数据存储端的最大线 程数。向导模式通过界面化配置并发数,指定任务所使用的并行度。
同步速率	设置同步速率可以保护读取端数据库,以避免抽取速度过大,给源库 造成太大的压力。同步速率建议限流,结合源库的配置,请合理配置 抽取速率。
错误记录数	错误记录数,表示脏数据的最大容忍条数。

配置	说明
任务资源组	任务运行的机器,如果任务数比较多,使用默认资源组出现等待资源 的情况,建议购买独享数据集成资源或添加自定义资源组,详情请参 见DataWorks独享资源和新增任务资源。

6. 单击右侧的调度配置,为节点配置调度属性。

" • •	ি ঠে			运	维↗
01 选择数据源		× 调度配置			调
	左 边田配罢购据	基础属性 ② -			度配量
		节点名:	etLrds_adb_001		Ľ
* 数据源	MySQL 🗸	节点D:			版本
*表	[°] product' ×	节点类型:	数据同步		
		责任人:			
数据过滤	proid>100	描述:			
		参数:	bizdate=Sbizdate 0		
切分键	proid				
()))) ())		时间属性 🕐 –			
		生成实例方式:	○ T+1次日生成 ● 发布后即时生成		
02 字段映射		时间屋性:	① 正常调度 ② 空胞调度		
		出错重试:	✓ ②		
	源头表字段	生效日期:	1970-01-01 . 9999-01-01		
	proid				
	proname	暫停调度:			
	price	调度周期:			
		定时调度:			

7. 配置完成后,单击保存并提交。

新建数据开发任务

1. 右键单击新建业务流程下的数据开发,选择新建数据开发节点 > AnalyticDB for MySQL。



2. 在新建节点对话框中,填写节点名称,单击提交。

新建节点		×
节点类型:	AnalyticDB for MySQL	
节点名称:	节点名称	
目标文件夹:	业务流程/works/数据开发	
		提交取消

3. 选择相应的数据源后,根据AnalyticDB for MySQL支持的语法,编写SQL语句。通常支持DML语句,您也可以执行DDL语句。



4. 单击右侧的调度配置,为节点配置调度属性。

	B								
洗择数	× ***	调度配置							
据源	机处								
1 INS	SERT INTO income(pro	父节点输出名称	父节点输出表名	3 节点名		父节点ID	责任人	来源	操作
2 SEI	LECT	tongguan_dab_dw.500134826_out		etl_rds_a	db_001			手动添加	
3 a.µ 4 ,S	pro_id TR_TO_DATE('\${bizdat	null.product						自动解析	
5 <mark>, SI</mark> 6 fro	um(b.price*a.num) in om orders a	null.orders						自动解析	
7 let 8 whe	ft join product b or ere a.order_id<'100'	本节点的输出: 请输入节点输出名称	ĸ						
10 gro	oup by a.pro_1d;	输出名称		输出表名	下游节点名	3称 下游节点ID	责任人	来源	操作
		tongguan_dab_dw.500135517_out		- C				系统默认添加	
		tongguan_dab_dw.adb_dw_pruduct_d	C	- C				手动添加	
		节点上下文 ②							
		本节点输入参数 添加							

5. 配置完成后,单击保存按钮,将其保存至服务器。然后单击运行按钮,即可立即执行编辑的SQL语句。

数据运维

提交并发布新建的节点任务后,单击左上角的图标,选择全部产品 > 运维中心,即可进行数据运维 操作。详情请参见运维中心模块的文档。

⑤ 袋 运维中心			& DataStudio থ্	· ····
=				
③ 运维大屏	实例执行概览		任务类型: 全	部 ~
→ 任务列表		4		
局期任务				
家 手动任务				
ᇦ 任务运维				
同期实例		3.		A la
F 手动实例				$= 1/N_{\odot}$
副 测试实例	💼 未运行 👘 等待资源 🛑 运行失时	ά.		$i \rightarrow i$
計数据实例	🧰 等待时间 👥 运行中 📰 运行成1	2 00 01 02 03 04 05 06 07 08 09 10 11 12	13 14 15 16 17 18 19	20 21 22 23
▶ 智能监控	24日 100 万史平均			
	任务运行情况	任务节点执行		地行時代
	11 <u>2</u>	DIMINE .	Date Mark	DADADES
	0.9 -		暂无数据	
	·		<i>∂</i> DataStudio ଝି	
 (¹) 运维大屏	授業			
- 任务列表	基後 読売経 ✓ ✓ 契約市点 今日修改的市点 暂停 (冻结) 市点 重置 清空			
[2] 周期任务				○ 刷新 收起搜索
(家) 手动任务	名称节点ID	生产环境,请谨慎	喿作	୯ ଝ ୧ ୧ ଅ
_ 任务运维	2100002308	19		
- 同期实例	2100002308	16		
○ =====	2100002299	18		
		eti_rds_adb_001		
「智能監持		< > adb_dw_prudu AnalyticB for MrSqL	展开父节点	
			查看代码 三层	
			編編10点 查看实例 四层	
			查看血缘 五层	
			測试 六层 补数据 >	
	更多 🔻 < 1/1 >		暂停 (冻结)	
A 法供由 A A A			2 DataStudio	🔊 🔍 analyticdh sunn
	基本信息			g
③ 运维大屏	规则管理规则名利	調整の規則名称		
▶ 任务列表	规则名称: 请输入规则: 对象类型	1: 任务节点 ~	接收人 : 请选择	
▶ 任务运维	触发条件: 🔽 完成 🔽 未完成 規则対象	·· 序号 任务名称 责任人 工作空间		新建自定义规
▼ 智能监控	规则由规则的		接收人	操作
慧线实例	20417 节点弧		任务责任人	详情 关闭
基线管理	20418 节点成	请输入任务节点名称/ID	④ 任务责任人	详情 关闭
8 事件管理	20415 全局費 触发方式		事件责任人	详情 关闭
😑 规则管理	20416 全局基 触发条件	: 靖选择	基线责任人	详情 关闭
	报警行为			
	最大报警次费	t: 3 次		共4条 < 🚹
	最小小探察消息	: 30 分钟		
		- 00:00 X no.an		
	9071370871			
	报警方式	: □ 粒目 □ 即任 拨音力式差妙如子段		
	接收人	: 任务责任人		

元数据管理

您可以单击左上角的图标,选择全部产品 > 数据地图,进行元数据管理操作。详情请参见数据地 图模块的文档。

3数据安全

3.1 实现指定用户访问指定UDF最佳实践

本文将为您介绍如何实现将具体的某个资源(表、UDF等)设置为仅能被指定的用户访问。 此UDF涉及到数据的加密解密算法,属于数据安全管控范畴。

前提条件

您需要提前安装MaxCompute客户端,以实现指定UDF被指定用户访问的操作。详情请参见安装 并配置客户端。

常见方案

· Package方案,通过打包授权进行权限精细化管控。

Package通常是为了解决跨项目空间的共享数据及资源的用户授权问题。当通过Package授予 用户开发者角色后,用户则拥有所有权限,风险不可控。

- 首先,用户熟知的DataWorks开发者角色的权限如下所示。

odps@ sz	_mc>desc role role_project_dev;		
Authorization Type: Policy			
A	projects/sz_mc: *		
A	projects/sz_mc/instances/*: *		
A	projects/sz_mc/jobs/*: *		
A	<pre>projects/sz_mc/offlinemodels/*: *</pre>		
A	projects/sz_mc/packages/*: *		
A	<pre>projects/sz_mc/registration/functions/*: *</pre>		
A	projects/sz_mc/resources/*: *		
A	projects/sz_mc/tables/*: *		
А	projects/sz_mc/volumes/*:		

由上图可见,开发者角色对工作空间中的Package、Functions、Resources和Table默认 有全部权限,明显不符合权限配置的要求。

- 其次,通过DataWorks添加子账号并赋予开发者角色,如下所示。



由此可见,通过打包授权和DataWorks默认的角色都不能满足我们的需求。例如将子账号RAM \$xxxxx.pt@aliyun-test.com:ramtest授予开发者角色,则默认拥有当前工作空间中所 有Object的所有操作权限,详情请参见用户授权。
· 在DataWorks中新建角色来进行高级管控。

您可以进入DataWorks控制台中的工作空间配置 > MaxCompute高级配置 > 自定义用户角 色页面,进行高级管控。但是在MaxCompute高级配置中只能针对某个表/某个项目进行授 权,不能对资源和UDF进行授权。

· Role Policy方案,通过Role Policy自定义Role的权限集合。

```
通过Policy可以精细化地管理到具体用户针对具体资源的具体权限粒度,下文将为您详细描述如何通过Role Policy方案实现指定UDF被指定用户访问。
```

▋ 说明:

为了安全起见,建议初学者使用测试项目来验证Policy。

通过Policy自定义Role的权限集合

- 1. 创建默认拒绝访问UDF的角色。
 - a. 在客户端输入create role denyudfrole;, 创建一个role denyudfrole。
 - b. 创建Policy授权文件,如下所示。

```
{
"Version": "1", "Statement"
[{
  "Effect":"Deny",
  "Action":["odps:Read","odps:List"],
  "Resource":"acs:odps:*:projects/sz_mc/resources/getaddr.jar"
},
{
  "Effect":"Deny",
  "Action":["odps:Read","odps:List"],
  "Resource":"acs:odps:*:projects/sz_mc/registration/functions/
getregion"
}
```

] }

c. 设置和查看Role Policy。

在客户端输入put policy /Users/yangyi/Desktop/role_policy.json on role denyudfrole;命令,设置Role Policy文件的存放路径等配置。

通过get policy on role denyudfrole;命令, 查看Role Policy。

dps@sz_mc>get policy on role denyudfrole;
"Statement": [{
"Action": ["odps:Read",
"odps:List"],
"Effect": "Deny",
"Resource": ["acs:odps:*:projects/sz_mc/resources/getaddr.jar"]},
{
"Action": ["odps:Read",
"odps:List"],
"Effect": "Deny",
"Resource": ["acs:odps:*:projects/sz_mc/registration/functions/getre
ſion"]}],
"Version": "1">

d. 在客户端输入grant denyudfrole to RAM\$xxxx.pt@aliyun-

test.com:ramtest;, 添加子账号至role denyudfrole。

2. 验证拒绝访问UDF的角色是否创建成功。

以子账号RAM\$xxxx.pt@aliyun-test.com:ramtest登录MaxCompute客户端。

a. 登录客户端输入whoami;确认角色。

```
odps@ sz_mc>whoami;
Name: RAM$y_______pt@aliyun-test.com:ramtest
End_Point: http://service.odps.aliyun.com/api
Tunnel_End_Point: http://dt.cn-shanqhai.maxcompute.aliyun.com
Project: sz_mc
```

b. 通过show grants;查看当前登录用户权限。

odps@ sz_mc>show grants;					
[roles]					
role_project_dev, denyudfrole					
Authorization Type: Policy					
[role/denyudfrole]					
projects/sz_mc/registration/functions/getregion: List	Read				
projects/sz_mc/resources/getaddr.jar: List Read					
[role/role_project_dev]					
A projects/sz_mc: *					
A projects/sz_mc/instances/*: *					
A projects/sz_mc/jobs/*: *					
<pre>A projects/sz_mc/offlinemodels/*: *</pre>					
A projects/sz_mc/packages/*: *					
<pre>A projects/sz_mc/registration/functions/*: *</pre>					
A projects/sz_mc/resources/*: *					
A projects/sz_mc/tables/*: *					
A projects/sz_mc/volumes/*: *					

通过查询发现该RAM子账号有两个角色,一个是role_project_dev(即DataWorks默认的 开发者角色),另一个是刚自定义创建的denyudfrole。

c. 验证自建UDF以及依赖的包的权限。



通过上述验证发现,该子账号在拥有DataWorks开发者角色的前提下并没有自建UDF: getregion的读权限。但还需要结合Project Policy来实现该UDF只能被指定的用户访问。

3. 配置Project Policy。

a. 编写Policy。

```
{
  "Version": "1", "Statement":
  [{
  "Effect":"Allow",
  "Principal":"RAM$yangyi.pt@aliyun-test.com:yangyitest",
  "Action":["odps:Read","odps:List","odps:Select"],
  "Resource":"acs:odps:*:projects/sz_mc/resources/getaddr.jar"
  },
  {
    "Effect":"Allow",
    "Principal":"RAM$xxxx.pt@aliyun-test.com:yangyitest",
    "Action":["odps:Read","odps:List","odps:Select"],
    "Resource":"acs:odps:*:projects/sz_mc/registration/functions/
  getregion"
```

}] }

b. 设置和查询Policy。

通过put policy /Users/yangyi/Desktop/project_policy.json;命令设 置Policy文件的存放路径。

通过get policy;命令查看Policy。

odps	@ sz_mc	get policy;
{		
	"Stateme	ent": [{
		"Action": ["odps:Read",
		"odps:List".
		"odps:Select"].
		"Effect": "Allow".
		"Principal": ["RAM\$;;;;;::::::::::::::::::::::::::::::::
		"Resource": ["acs:odps:*:projects/sz_mc/resources/getaddr.jar"]},
		"Action": ["odps:Read",
		"odps:List",
		"odps:Select"].
		"Effect": "Allow".
		"Principal": ["RAM\$vanavi.pt@alivun-test.com:vanavitest"].
		"Resource": ["acs:odps:*:projects/sz mc/registration
	"Version	n": "1"}

c. 通过whoami;和show grants;进行验证。



d. 运行SQL任务, 查看是否仅指定的RAM子账号能够查看指定的UDF和依赖的包。

odps@_sz	z_mc>sel	.ect getregion(172.100.30	⊥');					
ID = 201 Log view http://1 U00DA2MI bil6IjEi Job Queu	L9011409 v: Logview. TU2Myx7I LfQ== weing.	odps.aliyun.cc	2 m/logview/?h	=http:/ alijøbD	//service	e.odţ /hZ(os.aliyun CJdLCJFZm	.com/api& ZlY3QiOiJ	p=sz_n BbGx∨c
M1_job_0		STAGES	STATUS TERMINATED	TOTAL 1	COMPLE	TED 1	RUNNING Ø	PENDING Ø	BACKI
STAGES:	01/01	[100% E	LAPSI	ED TIME:	24.34 s	
Summary: resource inputs: outputs: Job run Job run M1: M1: + _c0 + [美国,	e cost: time: 1 mode: f engine: instand run tim instand input r output	cpu 0.27 Core 8.000 Fuxi job execution eng te count: 1 me: 18.000 te time: min: 16.000, records: records: AdhocSink1: 1 ,]	* Min, memor yine max: 16.000, L (min: 1, m	y 0.53 avg: 1	GB * Mii L6.000 avg: 1)	n			
odps@ sz_mc_d Name Owner Type Comment CreatedTime LastWodifiedT LastUpdator Size MdSsum	esc resource	getaddr.jar; getaddr.ja JAR IDE RESOUR 2018-05-24 2018-05-24 1353716 770497a9f6	n .pt#aliyun-test.com 19:51:16 19:51:16 05e09e198cb166cec7fa00	n me/admin/oxs 8	s-base-biz-pho	enix/ter	πρ/4d3efcczØsl9	eiirin53o5n4/get	addr.jar

总结

关于DataWorks和MaxCompute如何实现指定用户访问指定UDF,总结如下:

- ·如果您不想让其他用户访问工作空间内具体的资源,在DataWorks中添加数据开发者权限 后,再根据Role Policy的操作,在MaxCompute客户端将其配置为拒绝访问权限。
- ·如果您要指定用户访问指定资源,通过Role Policy在DataWorks中配置数据开发者权限后,再 根据Project Policy的操作,在MaxCompute客户端将其配置为允许访问权限。

3.2 子账号仅从特定IP登录DataWorks

在数据开发过程中,部分对权限控制较为严格的用户要求RAM子账号仅能通过公司内某个特定IP进行登录。本文将为您介绍如何实现RAM子账号仅从特定本地IP登录DataWorks。

前提条件

您需要首先参见#unique_80, 创建RAM子账号, 并且给子账号授权。权

限AliyunDataWorksFullAccess为系统默认权限,无法修改,您需要额外创建一个自定义授权策略。

创建自定义策略

- 1. 云账号登录RAM控制台。
- 2. 在左侧导航栏的权限管理菜单下,单击权限策略管理。
- 3. 单击新建权限策略。

4. 在新建自定义权限策略对话框中,填写策略名称为dataworksIPlimit1,勾选配置模式为脚本 配置,配置您的自定义权限。

RAM访问控制	RAM访问控制 / 权限策略管理 / 新建自定义权限策略					
概览	← 新建自定义权限策略					
人员管理へ						
用户组	策略名称					
用户	dataworksIPlimit1					
设置	备注					
权限管理 ヘ						
授权	配置模式					
权限策略管理	○ 可视化配置 ○ 脚本配置					
RAM角色管理	≪					
OAuth应用管理 (公测	导入已有系统策略					
	1 {					
	2 "Version": "1",					
	3 "Statement":					
	4 [{					
	5 "Effect": "Deny",					
	6 "Action": ["dataworks:*"],					
	7 "Resource": ["acs:dataworks:*:*:"],					
	8 "Condition":					
	9 {					
	10 "NotIpAddress":					
	确定 返回					

自定义权限的完整内容如下所示, "acs:SourceIP"后填写的IP, 即为您允许访问DataWorks的IP, 本示例设置为100.1.1.1/32。



5. 单击确认。

授权自定义策略给RAM用户

- 1. 在左侧导航栏的人员管理菜单下,单击用户。
- 2. 在用户登录名称/显示名称列表下,找到目标RAM用户。
- 3. 单击添加权限, 被授权主体会自动填入。
- 4. 在左侧权限策略名称列表下,单击需要授予RAM用户的权限策略。

沿山
「尻明」

在右侧区域框,选择某条策略并单击×,可撤销该策略。

- 5. 单击确定。
- 6. 单击完成。

验证结果

使用不同于100.1.1.1/32的地址登录DataWorks控制台,发现登录失败。

