

ALIBABA CLOUD

阿里云

智能数据构建与管理 Dataphin
最佳实践

文档版本：20210303

 阿里云

法律声明

阿里云提醒您在使用或阅读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。

1. 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。
2. 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
 危险	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。	 危险 重置操作将丢失用户配置数据。
 警告	该类警示信息可能会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，恢复业务时间约十分钟。
 注意	用于警示信息、补充说明等，是用户必须了解的内容。	 注意 权重设置为0，该服务器不会再接受新请求。
 说明	用于补充说明、最佳实践、窍门等，不是用户必须了解的内容。	 说明 您也可以通过按Ctrl+A选中全部文件。
>	多级菜单递进。	单击设置>网络>设置网络类型。
粗体	表示按键、菜单、页面名称等UI元素。	在结果确认页面，单击确定。
Courier字体	命令或代码。	执行 <code>cd /d C:/window</code> 命令，进入Windows系统文件夹。
斜体	表示参数、变量。	<code>bae log list --instanceid</code> <i>Instance_ID</i>
[] 或者 [a b]	表示可选项，至多选择一个。	<code>ipconfig [-all -t]</code>
{ } 或者 {a b}	表示必选项，至多选择一个。	<code>switch {active stand}</code>

目录

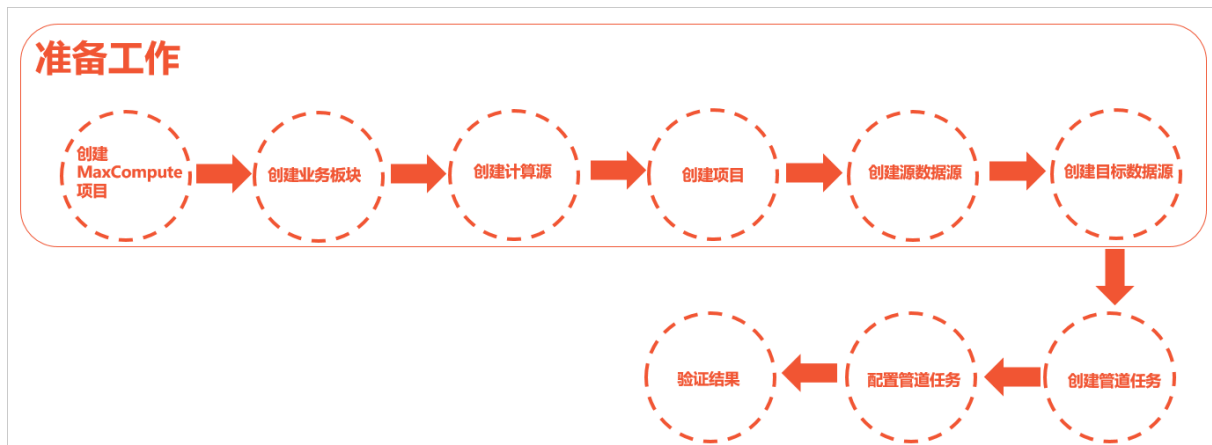
1.数据集成	05
1.1. 一键生成目标表	05
2.数据研发	12
2.1. 规范定义最佳实践	12
2.2. DataX同步数据	14
2.3. 使用Python读文件	20
2.4. Java UDF最佳实践	23

1.数据集成

1.1. 一键生成目标表

在数据集成过程中，当目标数据源类型为MaxCompute时，您可以通过Dataphin提供的一键生成目标表的功能，快速创建目标表。

操作流程



主流程	描述
创建MaxCompute项目	创建Dataphin计算源的MaxCompute项目 (best_practice_dev和best_practice_prod)。
创建业务板块	创建业务板块和数据域。
创建计算源	创建Dev和Prod项目的计算源。
创建项目	创建数据开发的工作空间。
创建源数据源	连接您的业务数据源至Dataphin平台。
创建目标数据源	创建数据集成的目标数据源 (MaxCompute类型的数据源)。
创建离线单条管道	创建用于数据集成的管道任务。
一键生成目标表	一键创建目标数据源的目标表，以配置管道任务。
验证结果	验证目标表是否集成到数据。

准备工作

- 准备数据源，详情请参见[准备数据源](#)。
- 创建MaxCompute项目，详情请参见[创建工作空间](#)。

工作空间名称	选择计算引擎服务
best_practice_dev	MaxCompute

工作空间名称	选择计算引擎服务
best_practice_prod	MaxCompute

- 创建业务板块LD_best_practice和LD_best_practice_dev，创建数据域test，详情请参见[新建业务板块](#)。
- 创建Dev项目计算源（best_practice_dev），配置如下参数，详情请参见[新建MaxCompute计算源](#)。

参数	描述
计算类型	默认为MaxCompute，不支持修改。
计算源名称	输入best_practice_dev。
计算源描述	输入Dev项目的计算源。
Endpoint	默认为 <code>http://service.cn.maxcompute.aliyun.com/api</code> ，不支持修改。
Project Name	输入best_practice_dev。
Access ID	访问密钥中的AccessKey ID，您可以通过 用户信息管理 页面获取。
Access Key	访问密钥中的AccessKey Secret，您可以通过 用户信息管理 页面获取。

- 创建Prod项目计算源（best_practice_prod），配置如下参数，详情请参见[新建MaxCompute计算源](#)。

参数	描述
计算类型	默认为MaxCompute，不支持修改。
计算源名称	输入best_practice_prod。
计算源描述	输入Prod项目的计算源。
Endpoint	默认为 <code>http://service.cn.maxcompute.aliyun.com/api</code> ，不支持修改。
Project Name	输入best_practice_prod。
Access ID	访问密钥中的AccessKey ID，您可以通过 用户信息管理 页面获取。
Access Key	访问密钥中的AccessKey Secret，您可以通过 用户信息管理 页面获取。

- 创建项目best_practice和best_practice_dev，详情请参见[新建项目](#)。



参数	描述
公用名称	输入最佳实践。
公用英文名	输入best_practice。
业务板块	选择LD_best_practice。
空间类型	选择应用层。

参数	描述
项目1: Dev	选择离线计算源为best_practice_dev。
项目2: Prod	选择离线计算源为best_practice_prod。
描述	输入简单的描述。
沙箱白名单	添加沙箱白名单： <ol style="list-style-type: none"> i. 单击新建。 ii. 在访问地址输入框中，输入RDS ID.mysql.rds.aliyuncs.com。在访问地址输入框中，输入3306。 iii. 单击图标。

- 创建源数据源dataphin，详情请参见[新建MySQL数据源](#)。
- 创建目标数据源，配置如下参数，详情请参见[新建MaxCompute数据源](#)。

参数	描述
数据源类型	选择MaxCompute。
数据源名称	输入为target。
数据源描述	输入目标数据源。
数据源配置	选择生产+开发数据源。
JDBC URL	输入http://service.cn.maxcompute.aliyun.com/api。
Project Name	开发环境输入为best_practice_dev，生产环境输入为best_practice_prod。
Access ID	访问密钥中的AccessKey ID，您可以通过 用户信息管理 页面获取。
Access Key	访问密钥中的AccessKey Secret，您可以通过 用户信息管理 页面获取。

步骤一：创建管道任务

1. 登录Dataphin控制台。
2. 在Dataphin控制台页面，选择工作区地域后，单击进入Dataphin>>。
3. 在Dataphin首页，单击顶部菜单栏的研发。
4. 在数据开发页面，单击项目名称后的图标，单击Dev页签，选择best_practice_dev为数据开发的项目空间。
5. 在数据开发页面，鼠标悬停至顶部菜单栏中的开发上，单击集成。
6. 在数据集成页面，鼠标悬停在图标，单击离线单条管道。
7. 在创建管道开发脚本对话框，配置参数。

参数	描述
管道名称	输入test。
调度类型	调度类型选择手动节点。
描述	填写对离线单条管道的简单描述。
选择目录	默认目录为离线管道。

8. 单击确定。

步骤二：配置管道任务（一键生成目标表）

1. 在离线单条管道开发页面，单击右上方的组件库。
2. 单击输入前的 > 图标后，单击MySQL组件拖动到左侧的管道画布中。
3. 鼠标悬停至组件框内并右键单击，选择属性配置后，配置参数。

MySQL输入配置
组件说明 v x

* 步骤名称

* 数据源

* 来源表量 单表 多表

* 表

切分键

输入过滤

输出字段 字段管理

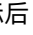
Q 请输入关键词

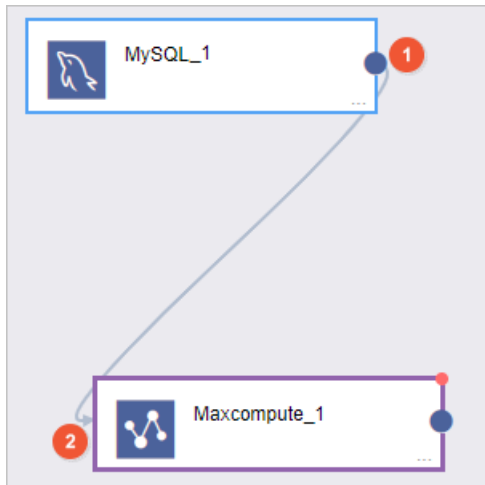
序号	字段	类型	操作
1	order_id	String	🗑
2	report_date	Date	🗑
3	customer_name	String	🗑
4	order_level1	String	🗑
5	order_number	Double	🗑
6	order_amt	Double	🗑

取消
确认

参数	描述
步骤名称	保持默认。
数据源	选择数据源（dataphin）。
来源表量	选择单表。
表	选择来源表company_sales_record_copy。
切分键	无需添加切分键。
输入过滤	无需添加过滤条件。
输出字段	查看输出字段。

4. 单击确认，完成输入组件的属性配置。

- 单击输出前的  图标后，单击MaxCompute组件拖动到左侧的管道画布中。
- 单击下图中输入组件（MySQL）①处后拖动并指向输出组件（MaxCompute）的②处，形成有向连线。



- 鼠标悬停至输出组件框内并右键单击，选择属性配置，配置参数。

Maxcompute输出配置
组件说明

* 步骤名称 Maxcompute_1

* 表 datax_test

* 分区 ds=\${bizdate}

* 数据源 目标数据源

* 加载策略 追加数据

一键生成目标表

* 字段映射

输入字段	输出字段	映射关系																																													
<div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;">Q 请输入关键词</div> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th>字段</th> <th>类型</th> <th>操作</th> </tr> </thead> <tbody> <tr><td>customer_n...</td><td>String</td><td>☰</td></tr> <tr><td>order_level1</td><td>String</td><td>☰</td></tr> <tr><td>order_number</td><td>Double</td><td>☰</td></tr> <tr><td>back_point</td><td>Double</td><td>☰</td></tr> <tr><td>shipping_type</td><td>String</td><td>☰</td></tr> </tbody> </table>	字段	类型	操作	customer_n...	String	☰	order_level1	String	☰	order_number	Double	☰	back_point	Double	☰	shipping_type	String	☰	<div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;">Q 请输入关键词</div> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th>字段</th> <th>类型</th> <th>操作</th> </tr> </thead> <tbody> <tr><td>order_name</td><td>String</td><td>🗑</td></tr> </tbody> </table>	字段	类型	操作	order_name	String	🗑	<div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;">Q 请输入关键词</div> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th>输入</th> <th>输出</th> <th>操作</th> </tr> </thead> <tbody> <tr><td>report_date</td><td>report_date</td><td>🗑</td></tr> <tr><td>order_amt</td><td>order_amt</td><td>🗑</td></tr> <tr><td>area</td><td>area</td><td>🗑</td></tr> <tr><td>province</td><td>province</td><td>🗑</td></tr> <tr><td>city</td><td>city</td><td>🗑</td></tr> <tr><td>product_type</td><td>product_type</td><td>🗑</td></tr> </tbody> </table>	输入	输出	操作	report_date	report_date	🗑	order_amt	order_amt	🗑	area	area	🗑	province	province	🗑	city	city	🗑	product_type	product_type	🗑
字段	类型	操作																																													
customer_n...	String	☰																																													
order_level1	String	☰																																													
order_number	Double	☰																																													
back_point	Double	☰																																													
shipping_type	String	☰																																													
字段	类型	操作																																													
order_name	String	🗑																																													
输入	输出	操作																																													
report_date	report_date	🗑																																													
order_amt	order_amt	🗑																																													
area	area	🗑																																													
province	province	🗑																																													
city	city	🗑																																													
product_type	product_type	🗑																																													

参数	说明
步骤名称	保持默认。
数据源	选择目标数据源。

参数	说明
表	<p>创建目标表：</p> <ol style="list-style-type: none"> i. 单击一键生成目标表。 ii. 在代码输入框中，输入建表语句。 <pre style="background-color: #f0f0f0; padding: 10px;">CREATE TABLE IF NOT EXISTS datax_test (order_id bigint comment '订单号', `area` string comment '区域', province string comment '省份', city string comment '城市', product_type string comment '类型', order_name string comment '客户名称', report_date datetime comment '日期', order_amt double comment '销售额') PARTITIONED BY (`ds` STRING);</pre> <ol style="list-style-type: none"> iii. 单击新建。 <p>? 说明 无需选中是否在生产建表。</p>
加载策略	选择 追加数据 。
分区	输入ds=\${bizdate}。
输入字段	根据上游的输入，为您展示输入字段。
输出字段	为您展示输出字段。
快速映射	<p>映射关系选择为同名映射的操作步骤：</p> <ol style="list-style-type: none"> i. 单击快速映射后的图标。 ii. 选择同名映射。 iii. 在提醒对话框中，单击确定。

8. 单击**确认**，完成输出组件的属性配置。

步骤三：验证结果

1. 单击管道开发页面左上方的**预览**。
2. 在对话框中，**bizdate**填写为20200819，单击**确定**。

运行日志										
MySQL_1										
order_id	report_date	customer_name	order_level1	order_number	order_amt	back_point	shipping_type	profit_amt	shipping_cos	sag
13729	2013-01-01...		其它	9	872.48	0.08	空运	-342.91	35	
28774	2013-01-01...		高级	33	180.36	0.1	火车	-111.8	4.69	
37537	2013-01-02...		低级	43	4083.19	0.07	大卡	-1049.85	45	
37537	2013-01-02...		低级	32	4902.38	0.05	火车	1438.49	7.07	
37537	2013-01-02...		低级	4	1239.06	0	大卡	-193.08	48.8	

3. 单击管道开发页面左上方的执行。

4. 在对话框中，bizdate填写为20200819，单击确定。查看数据是否同步至MaxCompute组件。

运行日志						
步骤度量						
预览结果						
步骤名称	线程号	输入 (条记录)	输出	读取	写入	时间
MySQL_1	0	17136	0	0	17136	2020-08-20 18:36:53
Maxcompute_1	0	0	17136	17136	0	2020-08-20 18:36:53

2. 数据研发

2.1. 规范定义最佳实践

基于Dataphin建模理论和业务需求，明确并规范定义统计指标，以便设计出易于业务使用的数据仓库。

背景信息

规范定义是指以维度建模作为理论基础，构建总线矩阵，划分并定义数据域、业务过程、维度、原子指标、统计周期和派生指标。

在您开始使用Dataphin进行数仓模型设计前，需要完成业务调研、需求分析、构建总线矩阵（从业务数据中抽象出业务过程和维度）、明确并定义统计指标。本教程中假设已完成需求调研、业务分析和构建总线矩阵，带您体验明确并规范定义统计指标，帮助您快速理解如何基于Dataphin设计数仓模型。

基本概念

名词	描述
业务板块	<p>业务板块定义了数据仓库的多种命名空间，是一种系统级的概念对象。当数据的业务含义存在较大差异时，您可以创建不同的业务板块，让各成员独立管理不同的业务，后续数据仓库的建设将按照业务板块进行划分。</p> <p>在Dataphin中，项目可以归属至业务板块以实现规范建模功能，同一个业务板块中可能包含多个不同的项目，所以业务板块与项目的关系为1:N。</p>
数据域	<p>数据域即主题域，是对某个主题分析后确定的主题边界。例如，商品域、交易域、会员域等。</p>
业务过程	<p>业务过程即企业的业务活动事件，通常为不可拆分的事件。创建业务过程，即从顶层视角，规范业务中的事务内容的类型及唯一性。例如电商订单是一个业务过程，业务过程由下单、支付、发货和确认收货等不可拆分的事件组成，每个事件就是一个业务过程。</p>
统计周期	<p>统计的时间范围，也可以称为时间周期。例如最近1天、最近30天等（类似于SQL中Where后的时间条件）。</p>
统计粒度	<p>统计分析的对象或视角，定义数据需要汇总的程度，可以理解为聚合运算时的分组条件（类似于SQL中Group By的对象）。粒度是维度或维度的某些属性的组合。例如，地域（维度）和客户性别（维度属性）组合成统计粒度。</p> <p>在定义粒度时，您需要充分考虑到业务和维度的关系。通常用于派生指标构建，是汇总表是唯一性识别方式。</p>
业务限定	<p>统计的业务范围，用于筛选出符合业务规则的记录（类似于SQL中Where后的条件，不包括时间区间）。</p>
度量	<p>事实就是度量，通常是对某个业务事件的衡量，通常为数字，如某笔订单的金额。</p> <div style="border: 1px solid #ccc; background-color: #e6f2ff; padding: 5px; margin-top: 10px;"> <p> 注意 请注意区分度量和原子指标。任何数据仓库都有维度和度量重要概念，但指标是业务分析中的概念。</p> </div>

名词	描述
维度	维度是分析业务的角度，是对应业务流程中的业务对象。例如客户、商品、部门等都可以作为分析业务的角度。
派生指标	派生指标是分析业务的指标。由原子指标、统计周期、统计粒度和业务限定组成。例如原子指标为支付金额，最近1天海外买家支付金额则为派生指标（最近1天为时间周期，海外为业务限定、买家为维度）。
原子指标	原子指标定义了业务分析的度量和统计方法（类似于SQL中Select后的聚合表达式，例如Sum）。

案例

A电商公司，销售某品牌多种零食。

商品种类	单价
干果类	10元/500g
膨化类	8元/袋
饮品类	15元/瓶

买家和卖家可以通过电商平台进行交易。



订单ID	买家ID	买家姓名	商品种类	支付金额	支付方式
29296	1001	张三	干果类	34	花呗
29297			干果类	67	花呗
29298	1003	李四	膨化类	56	支付宝

本案例中，明确及规范定义指标如下。

定义指标	业务数据
业务板块	电商业务
数据域	交易域
维度	商品种类
业务过程	下单购买

定义指标	业务数据
业务限定	商品种类为干果类
时间周期	最近1天
原子指标	销售总额
派生指标	最近1天干果类商品销售总额

2.2. DataX同步数据

DataX是异构数据源离线同步的工具，支持多种异构数据源之间高效的数据同步。Dataphin系统内嵌了DataX组件，支持通过构建Shell任务调用DataX，实现数据同步。本教程以MySQL数据库为例，为您介绍基于Dataphin如何调用DataX同步数据。

前提条件

开通RDS MySQL，详情请参见[RDS实例购买指南](#)。创建RDS MySQL实例过程中，需要您记录数据库名称、用户名和密码。

背景信息

DataX支持同步数据的数据源包括 MySQL、Oracle、SQL Server、PostgreSQL、HDFS、Hive、HBase等。有关DataX更多信息，请参见[DataX](#)。

准备工作

在您开始同步数据前，需要登录至RDS MySQL实例，创建同步数据的源表和目标表、获取实例的外网地址。同时，为了实现RDS MySQL实例和Dataphin之间的网络互通，需要将RDS MySQL实例的外网地址加入项目空间的沙箱白名单中，Dataphin的IP加入至 RDS MySQL实例的白名单中。

- 使用DMS工具连接实例，详情请参见[方法一：使用DMS工具连接实例（推荐）](#)，获取RDS MySQL的外网地址。



- 使用命令行方式连接实例，详情请参见[方法三：使用命令行方式连接实例](#)。创建同步数据的源数据表和目

标数据表。

创建源数据表的代码示例如下。

```
create table 'datax_test1'( 'area' varchar(255),'province' varchar(255) );
insert into datax_test1 values('华北','山东省'),('华南','河南省');
```

创建目标数据表的代码示例如下。

```
create table 'datax_test2'( 'area' varchar(255),'province' varchar(255) );
```

- 添加访问地址外网地址和端口3306至项目空间的沙箱白名单，;详情请参见[新建项目](#)。
RDS ID为RDS MySQL的实例ID。
- 添加Dataphin的IP至RDS MySQL的白名单，详情请参见[设置IP白名单](#)。
 - 如果数据源使用了VPC网络，则Dataphin的IP白名单见下表。

地域	IP白名单
华东2（上海）	100.104.228.128/26、100.104.115.192/26
华南1（深圳）	100.104.48.128/26
华北2（北京）	100.104.238.64/26
所有地域	100.104.0.0/16

- 如果数据源公网数据库或数据源使用了公网IP，则Dataphin的IP白名单见下表。

地域	IP白名单
华东2（上海）	47.102.151.182
华南1（深圳）	119.23.173.65
华北2（北京）	123.56.104.202

- 请您下载DataX任务代码的配置模板（[JSON文件](#)）。

步骤一：编辑JSON文件

编辑已下载的JSON文件，配置模板的代码。

```
{
  "job": {
    "content": [
      {
        "reader": {
          "name": "mysqlreader",
          "parameter": {
            "column": [
```



```
        "{columnname1}",
        "{columnname2}"
    ],
    "connection": [
        {
            "jdbcUrl":
["jdbc:mysql://{Public Endpoint}:3306/{DatabaseName}"],
            "table": ["{table_name1}"]
        }
    ],
    "password": "{password}",
    "username": "{username}"
}
},
"writer": {
    "name": "mysqlwriter",
    "parameter": {
        "column": [
            "{columnname1}",
            "{columnname2}"
        ],
        "connection": [
            {
                "jdbcUrl":
"jdbc:mysql://{Public Endpoint}:3306/{DatabaseName}",
                "table": ["{table_name2}"]
            }
        ],
        "password": "{password}",
        "username": "{username}"
    }
}
}
},
"setting": {
    "speed": {
        "channel": "1"
    }
}
}
}
```


其中,

- {username}: 登录数据库的用户名
- {password}: 登录数据库的密码
- {{Public Endpoint}}: 外网地址
- {DatabaseName}: 数据库名称
- {table_name1}: 源数据表的表名
- {table_name2}: 目标数据表的表名
- {columnname1}、{columnname2}: 字段名

以上参数需要替换为准备工作中已创建数据库的对应信息, 同时根据业务情况可以增删同步的字段, 保存修改后的JSON文件至本地。

步骤二: 上传JSON文件

1. 登录Dataphin控制台。
2. 在Dataphin控制台页面, 选择工作区地域后, 单击进入Dataphin>>。
3. 进入资源管理页面。
 - i. 在Dataphin首页, 单击研发。
 - ii. 在数据开发页面, 单击数据处理。
 - iii. 在左侧导航栏, 单击  资源管理图标。
4. 在资源管理页面, 单击资源管理后的  图标。
5. 在新建资源对话框中, 配置参数。

新建资源
✕

* 类型 others ▾

* 名称 datax.json

* 描述 DataX test

10/128

上传文件 请点击选择文件上传

📎 datax

文件大小 1.59KB

* 计算类型 无归属引擎 ▾



选择目录 资源管理 ▾

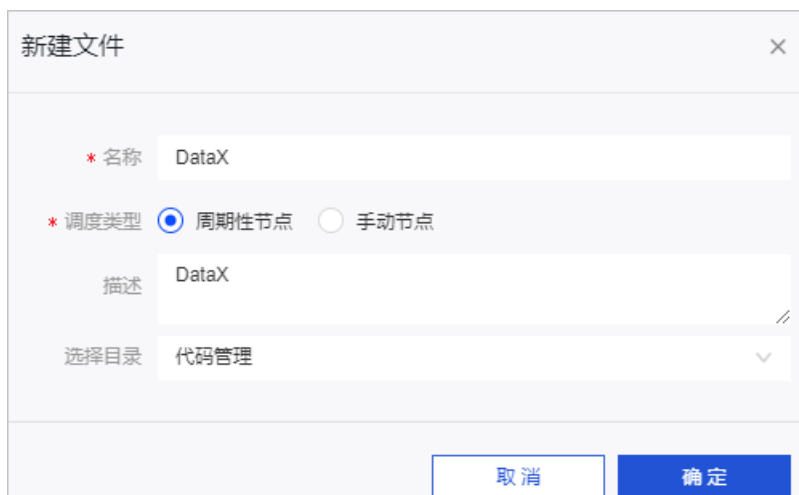
取消
提交

参数	描述
类型	选择others。
名称	填写资源名称，例如datax.json。 <div style="background-color: #e0f2f7; padding: 5px; margin-top: 5px;">🔔 注意 名称必须以.json结尾。</div>
描述	填写资源的描述，例如DataX test。
上传文件	选择本地的datax.json文件。
计算类型	选择无归属引擎。 <div style="background-color: #e0f2f7; padding: 5px; margin-top: 5px;">🔔 注意 DataX的资源存储至Dataphin系统，因此仅支持选择无归属引擎。</div>
选择目录	默认为资源管理。

6. 单击提交，完成资源的提交。
7. 在提交备注对话框，填写备注信息。
8. 单击确定并提交。

步骤三：创建DataX任务

1. 在数据处理页签，单击左侧导航栏计算任务图标。
2. 在计算任务页面，单击计算任务后的图标，选择通用脚本 > SHELL。
3. 在新建文件对话框，配置参数。



新建文件对话框包含以下配置项：

- * 名称: DataX
- * 调度类型: 周期性节点 手动节点
- 描述: DataX
- 选择目录: 代码管理

底部有“取消”和“确定”按钮。

参数	描述
名称	填写计算任务的名称为DataX。
调度类型	选择任务的调度类型为周期性节点。
描述	填写对任务的简单描述。
选择目录	选择DataX任务所在的目录。

4. 单击确定。

步骤四：编写并运行DataX任务代码

1. 在代码编写页面，编写运行DataX任务的代码。代码如下。运行代码节点的默认资源大小为256 MB。

```
@resource_reference{"datax.json"}
python $DATAX_HOME/bin/datax.py datax.json
```

如果DataX任务的资源大于256 MB，建议使用如下代码，自定义所需资源大小。

```
@required_resource{required_memory=2GB;required_cpus=1.0}
@resource_reference{"datax.json"}
python $DATAX_HOME/bin/datax.py --jvm '-Xms2g -Xmx2g' datax.json
```

其中，

- required_resource{}，自定义所需的资源大小。
- 系统已内置DataX_Home为DataX的安装目录（DATAX_HOME/bin/datax.py），DataX入口在DataX安装的bin文件下。


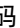
- --jvm '-Xms2g -Xmx2g'定义Dat aX实际运行时的虚拟机内存，建议内存大小与required_resource{}中的required_memory保持一致。
2. 单击页面右上角的执行，即可运行Dat aX任务代码。运行结果显示的读写失败总数为0时，表示Dat aX同步数据成功。

```
2020-11-24 13:33:25.564 [job-0] INFO JobContainer -
任务启动时刻      : 2020-11-24 13:33:14
任务结束时刻      : 2020-11-24 13:33:25
任务总计耗时      :                11s
任务平均流量      :                1B/s
记录写入速度      :                0rec/s
读出记录总数      :                2
读写失败总数      :                0
```

(可选)

(可选) 调度运维

如果您需要定期同步数据，则需要配置Dat aX同步任务的调度参数并发布至生产环境，参与生产环境的调度。

1. 在代码编写页面，单击页面上方的调度配置，配置调度参数，详情请参见[调度配置](#)。
2. 保存、提交和发布Dat aX同步任务。
 - i. 单击页面右上方的图标，保存代码。
 - ii. 单击页面右上方的图标，提交代码。
 - iii. 在提交备注对话框，填写备注信息。
 - iv. 单击确定并提交。
 - v. (可选) 发布Dat aX同步任务至生产环境。
 - 如果您的开发模式是Dev-Prod模式，则需要发布Dat aX同步任务至生产环境，详情请参见[管理发布任务](#)。
 - 如果您的开发模式是Basic模式，则提交成功的Dat aX同步任务即可进入生产环境。

2.3. 使用Python读文件



Dataphin仅支持开发基于Python的脚本，不支持开发依赖第三方组件的脚本。开发基于第三方组件的脚本，需要通过pip install下载第三方组件。本文为您介绍基于Dataphin如何通过构建Shell任务调用Python读取第三方文件。

前提条件

- 添加访问地址mirrors.aliyun.com和端口*至项目空间的沙箱白名单，详情请参见[设置IP白名单](#)。
- 已准备Python支持读取的文件，例如TXT、CSV、XLS、XLSX或PDF等格式文件。

步骤一：上传文件

1. 登录[Dataphin控制台](#)。
2. 在Dataphin控制台页面，选择工作区地域后，单击进入Dataphin>>。
3. 进入资源管理页面。
 - i. 在Dataphin首页，单击研发。

- ii. 在数据开发页面，单击数据处理。
 - iii. 在左侧导航栏，单击  资源管理图标。
4. 在资源管理页面，单击资源管理后的  图标。
5. 在新建资源对话框中，配置参数。

新建资源
✕

* 类型 others ▾

* 名称 test.xlsx

* 描述 test

0/128

上传文件 📁 请点击选择文件上传

 test.xlsx

文件大小 8.13KB

* 计算类型 无归属引擎 ▾

选择目录 资源管理 ▾



取消
提交

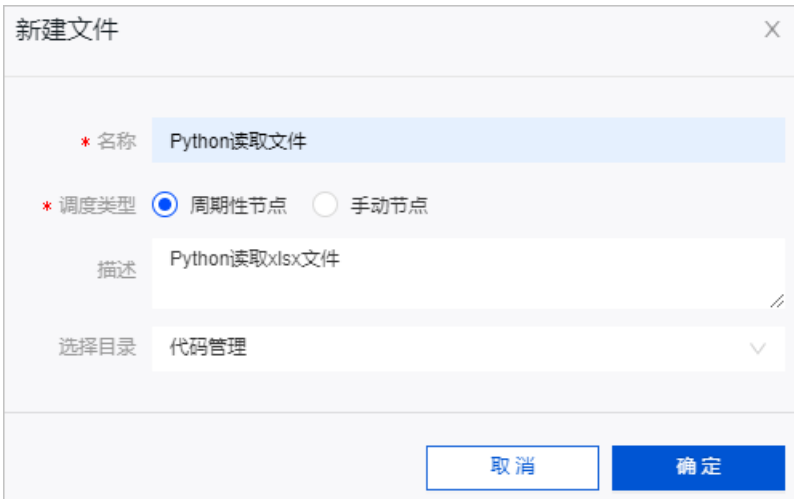
参数	描述
类型	选择others。
名称	上传文件的名称需要以文件类型结尾。例如test.xlsx。
描述	填写资源的描述。
上传文件	选择本地的文件，例如test.xlsx。
计算类型	选择无归属引擎。 <div style="background-color: #e0f2f7; padding: 5px; margin-top: 5px;"> 🔔 注意 文件资源存储至Dataphin系统，因此仅支持选择无归属引擎。 </div>
选择目录	默认为资源管理。

6. 单击提交，完成资源的提交。
7. 在提交备注对话框，填写备注信息。

8. 单击确定并提交。

步骤二：创建Shell任务

1. 在数据处理页签，单击左侧导航栏计算任务图标。
2. 在计算任务页面，单击计算任务后的图标，选择通用脚本 > SHELL。
3. 编写DataX任务代码。
 - i. 在新建文件对话框，配置参数。



新建文件对话框，包含以下配置项：

- * 名称: Python读取文件
- * 调度类型: 周期性节点 手动节点
- 描述: Python读取xlsx文件
- 选择目录: 代码管理

底部有取消和确定按钮。

参数	描述
名称	填写计算任务的名称，例如Python读取文件。
调度类型	选择任务的调度类型为周期性节点。
描述	填写对任务的简单描述。
选择目录	系统自动选择为代码管理。

ii. 单击确定。

步骤三：编写并运行Shell任务代码

1. 在代码编写页面，编写代码。

```
#在Dataphin的Linux服务器上新建目录。
mkdir -p /tmp/chars/ && \
#指定目录/tmp/chars/为python源。
pip install -i https://mirrors.aliyun.com/pypi/simple/\
--target=/tmp/chars/\
openpyxl
#指定的python源写入至openfile.py。
cat >openfile.py <<EOF
@resource_reference{"test.xlsx"}
#-*- coding:utf -*-
import os
import sys
sys.path.append('/tmp/chars/')
import openpyxl
print '===== python execute ok ====='
print("start=====")
args = sys.argv
# 打开excel文件，获取sheet名
wb = openpyxl.load_workbook(args[1])
# wb.get_sheet_names 这个方法已过时 会有一个警告
print(wb.worksheets[0])
EOF
#python中调用文件。
python openfile.py test.xlsx
```

其中，test.xlsx参数需要替换为您已上传的文件。

- 单击页面右上角的执行，即可运行任务代码。运行结果的状态为SUCCESS，表示读取文件成功。

```
2020-11-25 14:52:29 =====
2020-11-25 14:52:29 Current task status: SUCCESS
2020-11-25 14:52:29 Elapsed time: 4.997 s
2020-11-25 14:52:29 ----- voidemort task ends -----
```

2.4. Java UDF最佳实践

为了满足复杂的数据开发场景，Dataphin智能研发版支持自定义Java UDF函数。本教程以Java自带函数（toLowerCase）为例，为您介绍如何基于Dataphin自定义Java UDF函数。

前提条件

下载JAR包。

背景信息

本教程基于下载的JAR包自定义的Java UDF函数，实现大写字母转换为小写字母。您也可以编写Java UDF代码，以实现更多的功能，请参见[IntelliJ IDEA Java UDF开发最佳实践](#)。



本教程中的JAR包的代码如下。

```
package org.alidata.odps.udf.examples;
import com.aliyun.odps.udf.UDF;
public final class javaudf extends UDF {
    public String evaluate(String s) {
        if (s == null) {
            return null;
        }
        return s.toLowerCase();
    }
}
```

其中，

- JAR包路径为org.alidata.odps.udf.examples。
- class文件名为javaudf。

步骤一：上传JAR包

1. 登录Dataphin控制台。
2. 在Dataphin控制台页面，选择工作区地域后，单击进入Dataphin>>。
3. 进入资源管理页面。
 - i. 在Dataphin首页，单击研发。
 - ii. 在数据开发页面，单击数据处理。
 - iii. 在左侧导航栏，单击资源管理图标。
4. 在资源管理页面，单击资源管理后的图标。
5. 在新建资源对话框中，配置参数。

新建资源
✕

* 类型 ▼
jar

* 名称 ✎
javaudf.jar

* 描述 8/128 ✎
JavaUDF

上传文件 请点击选择文件上传

📎 javaudf.jar

文件大小 254.00B

* 计算类型 ▼
MaxCompute



选择目录 ▼
资源管理

取消
提交

参数	描述
类型	选择jar。
名称	上传文件的名称需要以文件类型结尾。例如javaudf.jar。
描述	填写资源的描述。
上传文件	选择本地JAR文件，例如javaudf.jar。
计算类型	选择MaxCompute。
选择目录	选择用于存放JAR包的目录。系统默认为资源管理，保持默认即可。

6. 单击提交，完成资源的提交。
7. 在提交备注对话框，填写备注信息。
8. 单击确定并提交。
9. （可选）发布资源至生产环境。
 - 如果您的开发模式是Dev-Prod模式，则需要发布资源至生产环境，详情请参见[管理发布任务](#)。
 - 如果您的开发模式是Basic模式，则提交成功的资源，即可进入生产环境。

步骤二：创建MAXC函数

1. 在数据处理页签，单击左侧导航栏的  函数管理图标。
2. 单击函数管理后的  图标，选择MAXC函数。

3. 在新建函数对话框，配置参数。




参数	描述
名称	填写函数的名称，例如java
选择资源	选择已上传的资源javaudf.jar。
类名	类名的格式为JAR包路径.class文件名。填写org.alidata.odps.udf.examples.javaudf。
类型	函数的类型。选择字符串。
命令格式	定义引用函数的格式。填写to_char(string i)。
使用文档	填写函数的使用描述，例如javaudf。
选择目录	默认为MAXC函数-用户自定义函数，保持默认。

4. 单击提交，完成资源的提交。
5. 在提交备注对话框，填写备注信息。
6. 单击确定并提交。
7. （可选）发布函数至生产环境。
 - 如果您的开发模式是Dev-Prod模式，则需要发布函数至生产环境，详情请参见[管理发布任务](#)。
 - 如果您的开发模式是Basic模式，则提交成功的函数，即可进入生产环境。

步骤三：新建SQL任务

1. 在数据处理页签，单击左侧导航栏计算任务图标。

2. 在计算任务页面，单击计算任务后的图标，选择MAXC任务 > MAX_COMPUTE_SQL。
3. 在新建文件对话框，配置参数。

参数	描述
名称	填写计算任务的名称，例如javaudf。
调度类型	选择任务的调度类型为周期性节点。
描述	填写对任务的简单描述。
选择目录	系统自动选择为计算任务。

4. 单击确定。

步骤四：使用Java UDF函数



1. 在SQL任务的代码编写页面，编写代码，例如 `select java('ABCGDfagHH');`。
2. 单击页面右上方的执行，查看运行结果。



(可选)

调度运维

如果需要定期的运行SQL任务，则需要配置SQL任务的调度参数并发布至生产环境，参与生产环境的调度。

1. 在代码编写页面，单击页面上方的**调度配置**，配置调度参数，详情请参见[调度配置](#)。
2. 保存、提交和发布SQL任务。
 - i. 单击页面右上方的图标，保存代码。
 - ii. 单击页面右上方的图标，提交代码。
 - iii. 在**提交备注**对话框，填写备注信息。
 - iv. 单击**确定并提交**。
 - v. （可选）发布SQL任务至生产环境。
 - 如果开发模式是Dev-Prod模式，则需要发布SQL任务至生产环境，详情请参见[管理发布任务](#)。
 - 如果开发模式是Basic模式，则提交成功的SQL任务，即可进入生产环境。