

# 尚硅谷大数据项目之电商数仓（业务数据采集平台）

（尚硅谷研究院）

版本：V5.0

## 第 1 章 电商业务简介

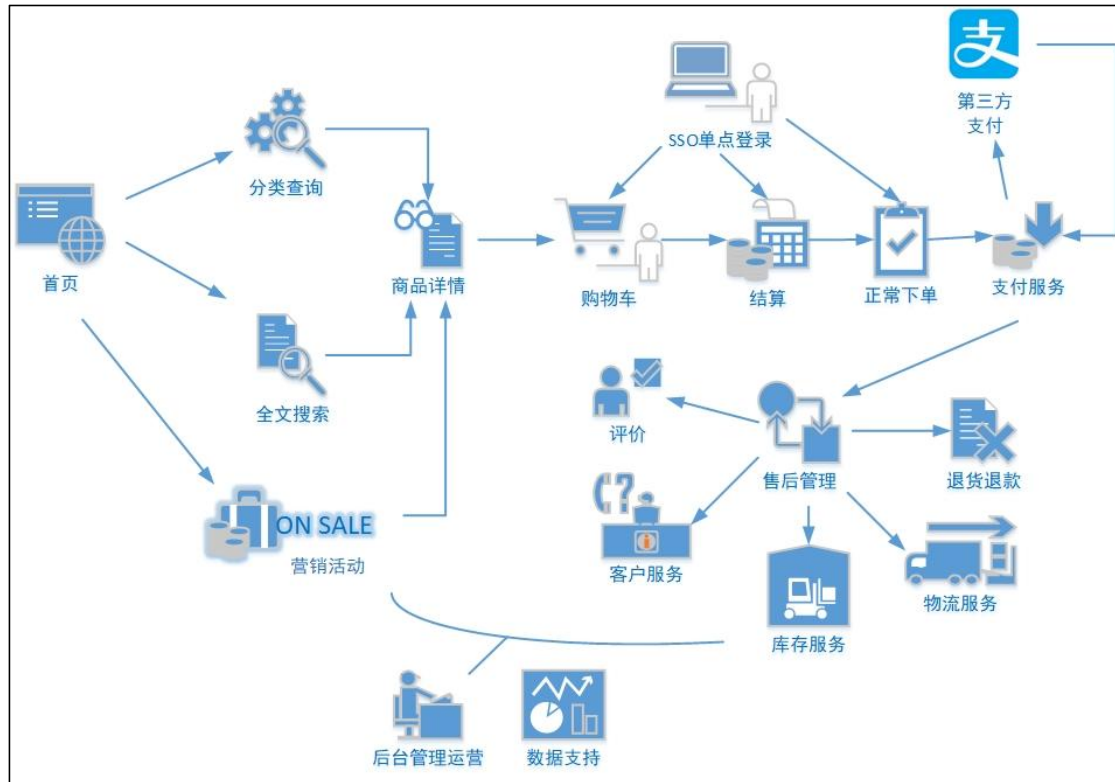
### 1.1 电商业务流程

电商的业务流程可以以一个普通用户的浏览足迹为例进行说明，用户点开电商首页开始浏览，可能会通过分类查询也可能通过全文搜索寻找自己中意的商品，这些商品无疑都是存储在后台的管理系统中的。

当用户寻找到自己中意的商品，可能会想要购买，将商品添加到购物车后发现需要登录，登录后对商品进行结算，这时候购物车的管理和商品订单信息的生成都会对业务数据库产生影响，会生成相应的订单数据和支付数据。

订单正式生成之后，还会对订单进行跟踪处理，直到订单全部完成。

电商的主要业务流程包括用户前台浏览商品时的商品详情的管理，用户商品加入购物车进行支付时用户个人中心&支付服务的管理，用户支付完成后订单后台服务的管理，这些流程涉及到了十几个甚至几十个业务数据表，甚至更多。



## 1.2 电商常识

### 1.2.1 SKU 和 SPU

**SKU=Stock Keeping Unit（库存量基本单位）**。现在已经被引申为产品统一编号的简称，每种产品均对应有唯一的 SKU 号。

**SPU（Standard Product Unit）**：是商品信息聚合的最小单位，是一组可复用、易检索的标准化信息集合。

例如：**iPhoneX 手机就是 SPU**。一台银色、128G 内存的、支持联通网络的 iPhoneX，就是 SKU。



更多 Java - 大数据 - 前端 - python 人工智能资料下载，可百度访问：尚硅谷官网

SPU 表示一类商品。同一 SPU 的商品可以共用商品图片、海报、销售属性等。

## 1.2.2 平台属性和销售属性

### 1.平台属性



### 2.销售属性



## 第 2 章 电商业务数据

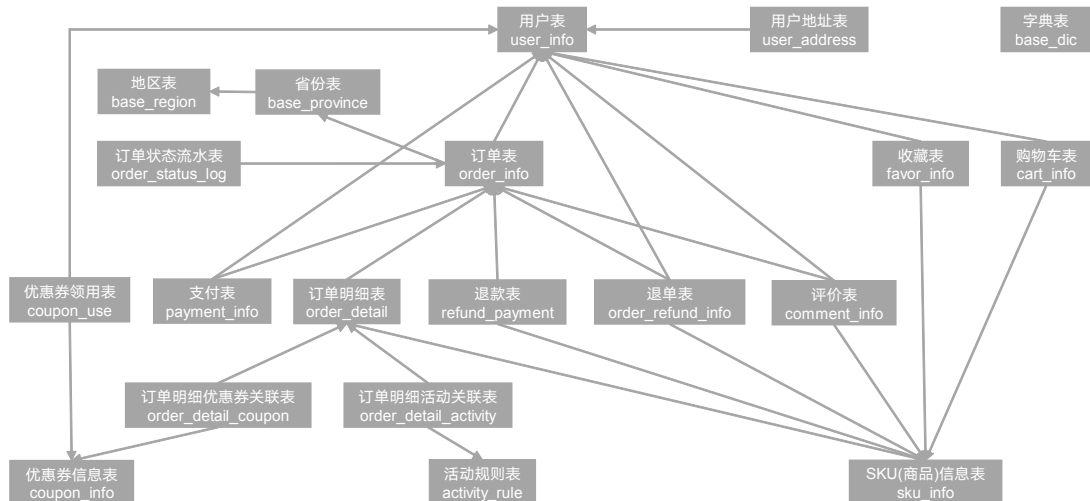
### 2.1 电商系统表结构

以下为本电商数仓系统涉及到的业务数据表结构关系。这 34 个表以订单表、用户表、SKU 商品表、活动表和优惠券表为中心，延伸出了优惠券领用表、支付流水表、活动订单表、订单详情表、订单状态表、商品评论表、编码字典表退单表、SPU 商品表等，用户表提供用户的详细信息，支付流水表提供该订单的支付详情，订单详情表提供订单的商品数量等情况，商品表给订单详情表提供商品的详细信息。本次讲解以此 34 个表为例，实际项目中，业务数据库中表格远远不止这些。

更多 [Java - 大数据 - 前端 - python 人工智能资料下载](#)，可百度访问：[尚硅谷官网](#)



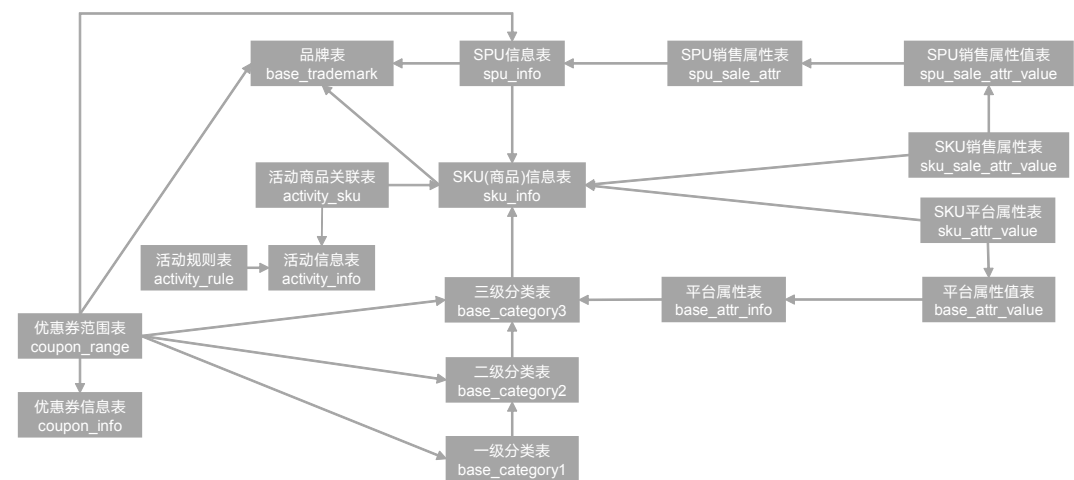
## 电商业务表



让天下没有难学的技术



## 后台管理系统



让天下没有难学的技术

### 2.1.1 活动信息表（activity\_info）

字段名	字段说明
id	活动 id
activity_name	活动名称
activity_type	活动类型（1：满减，2：折扣）
activity_desc	活动描述
start_time	开始时间
end_time	结束时间

create_time	创建时间
-------------	------

### 2.1.2 活动规则表（activity\_rule）

id	编号
activity_id	活动 ID
activity_type	活动类型
condition_amount	满减金额
condition_num	满减件数
benefit_amount	优惠金额
benefit_discount	优惠折扣
benefit_level	优惠级别

### 2.1.3 活动商品关联表（activity\_sku）

字段名	字段说明
id	编号
activity_id	活动 id
sku_id	sku_id
create_time	创建时间

### 2.1.4 平台属性表（base\_attr\_info）

字段名	字段说明
id	编号
attr_name	属性名称
category_id	分类 id
category_level	分类层级

### 2.1.5 平台属性值表（base\_attr\_value）

字段名	字段说明
id	编号
value_name	属性值名称
attr_id	属性 id

### 2.1.6 一级分类表（base\_category1）

字段名	字段说明
id	编号
name	分类名称

### 2.1.7 二级分类表（base\_category2）

字段名	字段说明
id	编号
name	二级分类名称
category1_id	一级分类编号

### 2.1.8 三级分类表（base\_category3）

字段名	字段说明
id	编号
name	三级分类名称
category2_id	二级分类编号

### 2.1.9 字典表（base\_dic）

字段名	字段说明
dic_code	编号
dic_name	编码名称
parent_code	父编号
create_time	创建日期
operate_time	修改日期

### 2.1.10 省份表（base\_province）

字段名	字段说明
id	id
name	省名称
region_id	大区 id
area_code	行政区位码
iso_code	国际编码
iso_3166_2	ISO3166 编码

### 2.1.11 地区表（base\_region）

字段名	字段说明
id	大区 id
region_name	大区名称

### 2.1.12 品牌表（base\_trademark）

字段名	字段说明
id	编号
tm_name	属性值
logo_url	品牌 logo 的图片路径

### 2.1.13 购物车表（cart\_info）

字段名	字段说明
id	编号
user_id	用户 id
sku_id	skuid
cart_price	放入购物车时价格
sku_num	数量
img_url	图片文件
sku_name	sku 名称 (冗余)
is_checked	
create_time	创建时间
operate_time	修改时间
is_ordered	是否已经下单
order_time	下单时间
source_type	来源类型
source_id	来源编号

### 2.1.14 评价表（comment\_info）

字段名	字段说明
id	编号
user_id	用户 id
nick_name	用户昵称
head_img	
sku_id	skuid
spu_id	商品 id
order_id	订单编号
appraise	评价 1 好评 2 中评 3 差评
comment_txt	评价内容
create_time	创建时间
operate_time	修改时间

### 2.1.15 优惠券信息表（coupon\_info）

字段名	字段说明
id	购物券编号
coupon_name	购物券名称
coupon_type	购物券类型 1 现金券 2 折扣券 3 满减券 4 满件打折券
condition_amount	满额数（3）
condition_num	满件数（4）
activity_id	活动编号
benefit_amount	减金额（1 3）
benefit_discount	折扣（2 4）
create_time	创建时间
range_type	范围类型 1、商品(spuid) 2、品类(三级分类 id) 3、品牌
limit_num	最多领用次数
taken_count	已领用次数
start_time	可以领取的开始日期
end_time	可以领取的结束日期
operate_time	修改时间
expire_time	过期时间
range_desc	范围描述

### 2.1.16 优惠券优惠范围表（coupon\_range）

字段名	字段说明
id	购物券编号
coupon_id	优惠券 id
range_type	范围类型 1、商品(spuid) 2、品类(三级分类 id) 3、品牌
range_id	

### 2.1.17 优惠券领用表（coupon\_use）

字段名	字段说明
id	编号
coupon_id	购物券 ID
user_id	用户 ID
order_id	订单 ID
coupon_status	购物券状态（1：未使用 2：已使用）
get_time	获取时间
using_time	使用时间
used_time	支付时间
expire_time	过期时间



### 2.1.18 收藏表（favor\_info）

字段名	字段说明
id	编号
user_id	用户名称
sku_id	skuid
spu_id	商品 id
is_cancel	是否已取消 0 正常 1 已取消
create_time	创建时间
cancel_time	修改时间

### 2.1.19 订单明细表（order\_detail）

字段名	字段说明
id	编号
order_id	订单编号
sku_id	sku_id
sku_name	sku 名称（冗余）
img_url	图片名称（冗余）
order_price	购买价格(下单时 sku 价格)
sku_num	购买个数
create_time	创建时间
source_type	来源类型
source_id	来源编号
split_total_amount	分摊总金额
split_activity_amount	分摊活动减免金额
split_coupon_amount	分摊优惠券减免金额

### 2.1.20 订单明细活动关联表（order\_detail\_activity）

字段名	字段说明
id	编号
order_id	订单 id
order_detail_id	订单明细 id
activity_id	活动 ID
activity_rule_id	活动规则
sku_id	skuID
create_time	获取时间

### 2.1.21 订单明细优惠券关联表（order\_detail\_coupon）

字段名	字段说明
id	编号
order_id	订单 id
order_detail_id	订单明细 id
coupon_id	购物券 ID
coupon_use_id	购物券领用 id
sku_id	skulD
create_time	获取时间

### 2.1.22 订单表(order\_info)

字段名	字段说明
id	编号
consignee	收货人
consignee_tel	收件人电话
total_amount	总金额
order_status	订单状态
user_id	用户 id
payment_way	付款方式
delivery_address	送货地址
order_comment	订单备注
out_trade_no	订单交易编号（第三方支付用）
trade_body	订单描述(第三方支付用)
create_time	创建时间
operate_time	操作时间
expire_time	失效时间
process_status	进度状态
tracking_no	物流单编号
parent_order_id	父订单编号
img_url	图片路径
province_id	地区
activity_reduce_amount	促销金额
coupon_reduce_amount	优惠券
original_total_amount	原价金额
freight_fee	运费
freight_fee_reduce	运费减免
refundable_time	可退款日期（签收后 30 天）

### 2.1.23 退单表（order\_refund\_info）

字段名	字段说明
id	编号
user_id	用户 id
order_id	订单 id
sku_id	skuid
refund_type	退款类型
refund_num	退货件数
refund_amount	退款金额
refund_reason_type	原因类型
refund_reason_txt	原因内容
refund_status	退款状态（0：待审批 1：已退款）
create_time	创建时间

### 2.1.24 订单状态流水表（order\_status\_log）

字段名	字段说明
id	
order_id	
order_status	
operate_time	

### 2.1.25 支付表（payment\_info）

字段名	字段说明
id	编号
out_trade_no	对外业务编号
order_id	订单编号
user_id	
payment_type	支付类型（微信 支付宝）
trade_no	交易编号
total_amount	支付金额
subject	交易内容
payment_status	支付状态
create_time	创建时间
callback_time	回调时间
callback_content	回调信息

### 2.1.26 退款表（refund\_payment）

字段名	字段说明
id	编号
out_trade_no	对外业务编号
order_id	订单编号
sku_id	
payment_type	支付类型（微信 支付宝）
trade_no	交易编号
total_amount	退款金额
subject	交易内容
refund_status	退款状态
create_time	创建时间
callback_time	回调时间
callback_content	回调信息

### 2.1.27 SKU 平台属性值表（sku\_attr\_value）

字段名	字段说明
id	编号
attr_id	属性 id（冗余）
value_id	属性值 id
sku_id	skuid
attr_name	属性名称
value_name	属性值名称

### 2.1.28 SKU 信息表（sku\_info）

字段名	字段说明
id	库存 id(itemID)
spu_id	商品 id
price	价格
sku_name	sku 名称
sku_desc	商品规格描述
weight	重量
tm_id	品牌(冗余)
category3_id	三级分类 id（冗余）
sku_default_img	默认显示图片(冗余)
is_sale	是否销售（1：是 0：否）
create_time	创建时间

### 2.1.29 SKU 销售属性表（sku\_sale\_attr\_value）

字段名	字段说明
id	id
sku_id	库存单元 id
spu_id	spu_id(冗余)
sale_attr_value_id	销售属性值 id
sale_attr_id	
sale_attr_name	
sale_attr_value_name	

### 2.1.30 SPU 信息表（spu\_info）

字段名	字段说明
id	商品 id
spu_name	商品名称
description	商品描述(后台简述)
category3_id	三级分类 id
tm_id	品牌 id

### 2.1.31 SPU 销售属性表（spu\_sale\_attr）

字段名	字段说明
id	编号(业务中无关联)
spu_id	商品 id
base_sale_attr_id	销售属性 id
sale_attr_name	销售属性名称(冗余)

### 2.1.32 SPU 销售属性值表（spu\_sale\_attr\_value）

字段名	字段说明
id	销售属性值编号
spu_id	商品 id
base_sale_attr_id	销售属性 id
sale_attr_value_name	销售属性值名称
sale_attr_name	销售属性名称(冗余)

### 2.1.33 用户地址表（user\_address）

字段名	字段说明
id	编号

user_id	用户 id
province_id	省份 id
user_address	用户地址
consignee	收件人
phone_num	联系方式
is_default	是否是默认

### 2.1.34 用户信息表（user\_info）

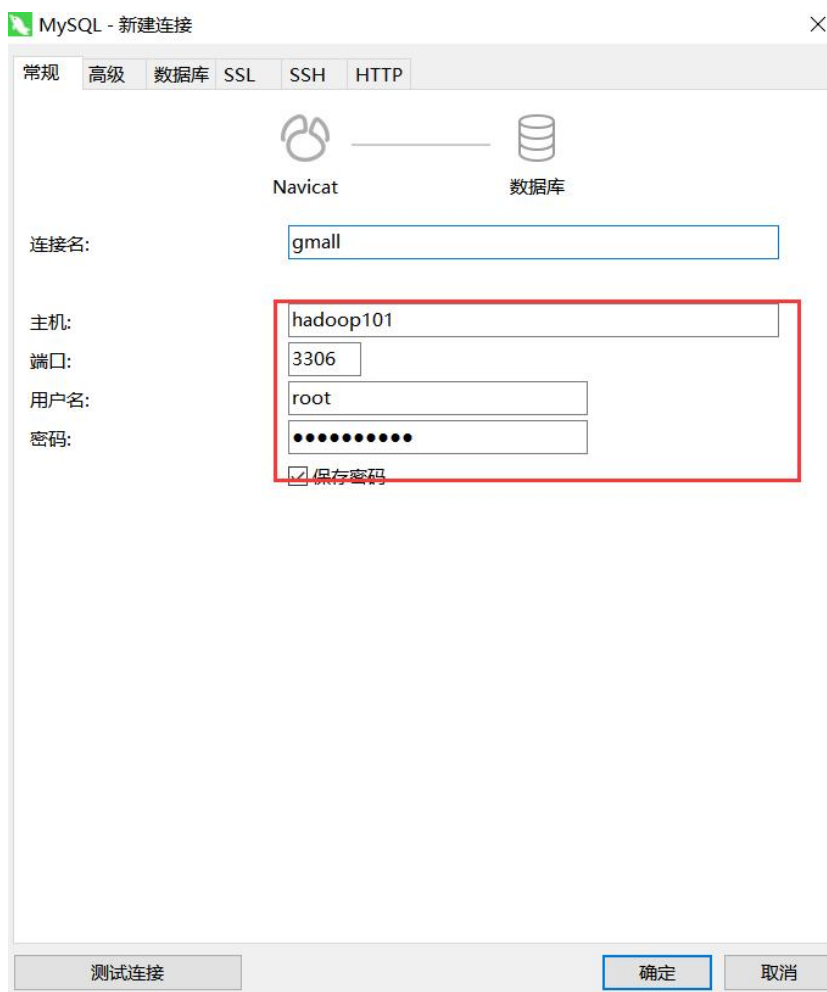
字段名	字段说明
id	编号
login_name	用户名称
nick_name	用户昵称
passwd	用户密码
name	用户姓名
phone_num	手机号
email	邮箱
head_img	头像
user_level	用户级别
birthday	用户生日
gender	性别 M 男,F 女
create_time	创建时间
operate_time	修改时间
status	状态

## 2.2 模拟生成业务数据

### 2.2.1 业务数据生成

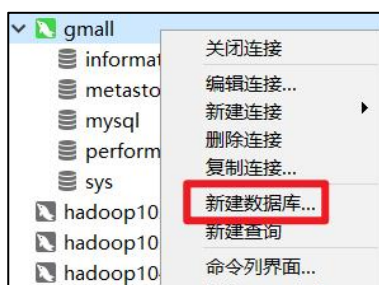
#### 1) 连接 MySQL

通过 MySQL 可视化客户端连接数据库。

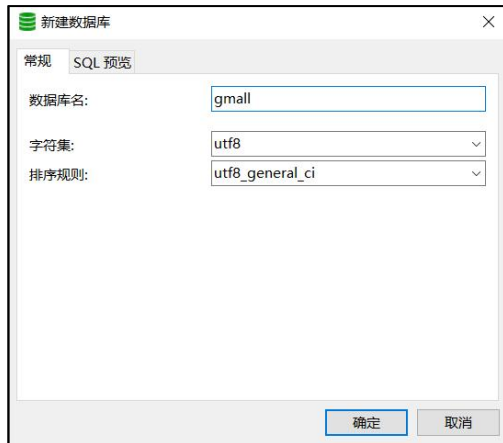


## 2) 建表语句

(1) 通过 SQLyog 创建数据库



(2) 设置数据库名称为 **gmall**，编码为 **utf-8**，**排序规则**为 utf8\_general\_ci



(3) 导入数据库结构脚本 (**gmall.sql**)



注意：完成后，要记得右键，刷新一下对象浏览器，就可以看见数据库中的表了。

### 3) 生成业务数据

(1) 在 hadoop103 的 /opt/module/ 目录下创建 db\_log 文件夹

```
[root@hadoop103 module]# mkdir db_log/
```

(2) 把 gmall2020-mock-db-2022-03-29.jar 和 application.properties 上传到 hadoop103 的 /opt/module/db\_log 路径上。

(3) 根据需求修改 application.properties 相关配置

```
logging.level.root=info

spring.datasource.driver-class-name=com.mysql.jdbc.Driver
spring.datasource.url=jdbc:mysql://hadoop101:3306/gmall?useUnicode=true&characterEncoding=utf-8&useSSL=false&serverTimezone=GMT%2B8
spring.datasource.username=root
spring.datasource.password=000000

logging.pattern.console=%m%n

mybatis-plus.global-config.db-config.field-strategy=not_null
```

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：尚硅谷官网



```
mybatis.mapperLocations=classpath:mapper/*.xml

#业务日期
mock.date=2020-06-14
#是否重置，首日须设置为 1
mock.clear=1
#是否重置用户，首日须设置为 1
mock.clear.user=1

#生成新用户数量
mock.user.count=200
#男性比例
mock.user.male-rate=20
#用户数据变化概率
mock.user.update-rate:20

#收藏取消比例
mock.favor.cancel-rate=10
#收藏数量
mock.favor.count=100

#每个用户添加购物车的概率
mock.cart.user-rate=10
#每次每个用户最多添加多少种商品进购物车
mock.cart.max-sku-count=8
#每个商品最多买几个
mock.cart.max-sku-num=3

#购物车来源 用户查询，商品推广，智能推荐，促销活动
mock.cart.source-type-rate=60:20:10:10

#用户下单比例
mock.order.user-rate=30
#用户从购物中购买商品比例
mock.order.sku-rate=50
#是否参加活动
mock.order.join-activity=1
#是否使用购物券
mock.order.use-coupon=1
#购物券领取人数
mock.coupon.user-count=100

#支付比例
mock.payment.rate=70
#支付方式 支付宝：微信：银联
mock.payment.payment-type=30:60:10

#评价比例 好：中：差：自动
mock.comment.appraise-rate=30:10:10:50

#退款原因比例：质量问题 商品描述与实际描述不一致 缺货 号码不合适 拍错 不想买了
其他
mock.refund.reason-rate=30:10:20:5:15:5:5

logging.level.com.atguigu.gmall2020.mock.db.mapper=debug
```

更多 [Java](#) -[大数据](#) -[前端](#) -[python](#) 人工智能资料下载，可百度访问：[尚硅谷官网](#)

(4) 并在该目录下执行，如下命令，生成 2020-06-14 日期数据：

```
[root@hadoop103 db_log]# java -jar
gmall2020-mock-db-2022-03-29.jar
```

(5) 查看 gmall 数据库，观察是否有 2020-06-14 的数据出现

## 2.2.2 业务数据梳理工具

可借助 EZDML 这款数据库设计工具，来辅助我们梳理复杂的业务表关系。

### 1) 下载地址

[http://www.ezdml.com/download\\_cn.html](http://www.ezdml.com/download_cn.html)

### 2) 使用说明

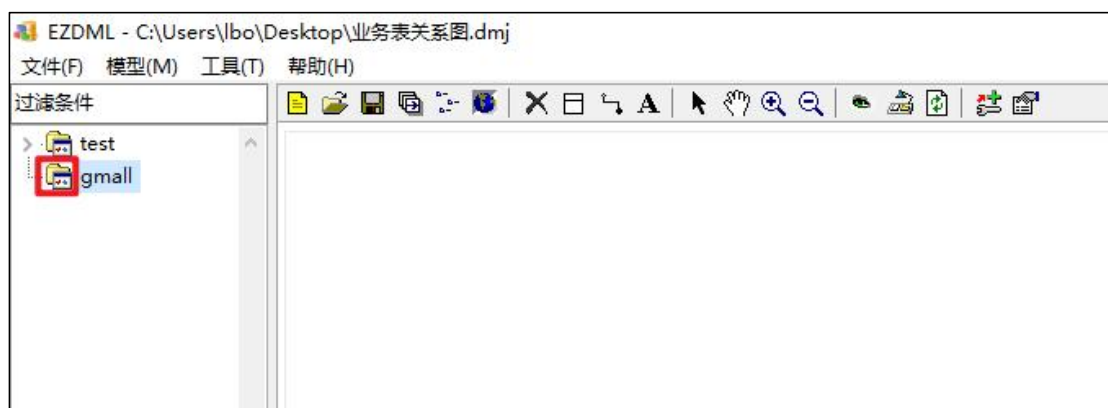
#### (1) 新建模型



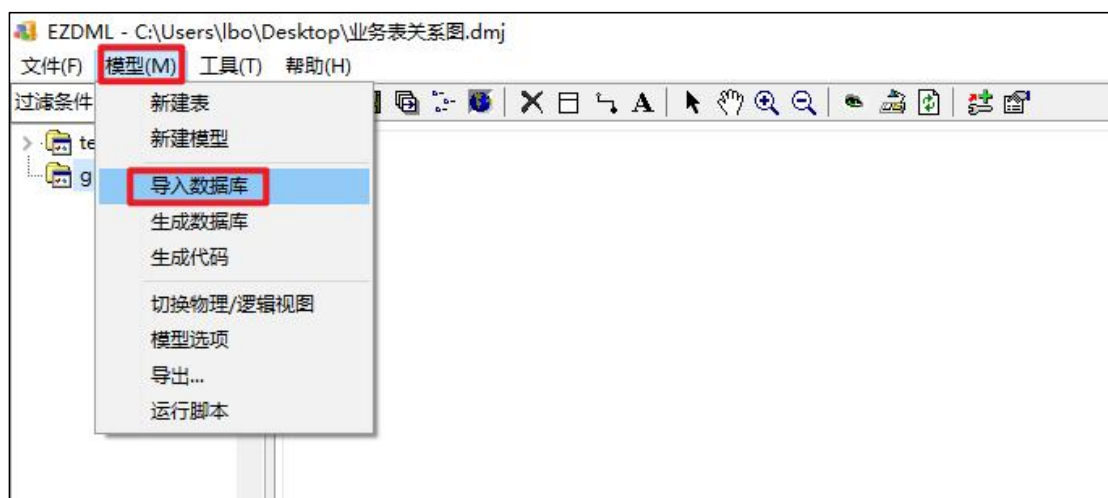
#### (2) 命名模型



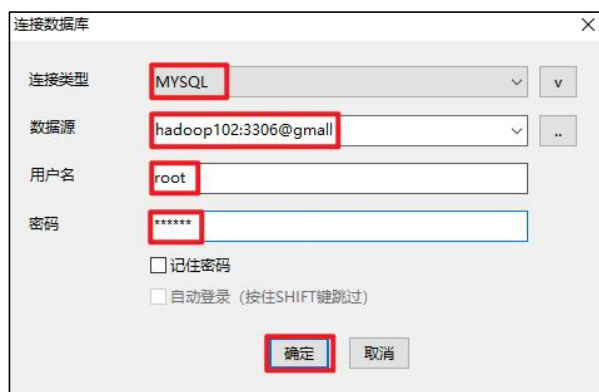
#### (3) 点击图标，选中模型



#### （4）导入数据库



#### （5）配置数据库连接



#### （6）选择导入的表（标注红点的表不需要导入）

连接数据库

[MYSQL]hadoop102:3306@gmail

选择用户

gmail

过滤条件

选择对象

☒ activity\_info
 ☒ activity\_rule
 ☒ activity\_sku
 ☒ base\_attr\_info
 ☒ base\_attr\_value
 ☒ base\_category1
 ☒ base\_category2
 ☒ base\_category3
 ☒ base\_category\_view
 ☒ base\_dic
 ☒ base\_frontend\_param
 ☒ base\_province
 ☒ base\_region
 ☒ base\_sale\_attr
 ☒ base\_trademark
 ☒ cart\_info
 ☒ cms\_banner

☒ comment\_info
 ☒ coupon\_info
 ☒ coupon\_range
 ☒ coupon\_use
 ☒ favor\_info
 ☒ financial\_sku\_cost
 ☒ order\_detail
 ☒ order\_detail\_activity
 ☒ order\_detail\_coupon
 ☒ order\_info
 ☒ order\_refund\_info
 ☒ order\_status\_log
 ☒ payment\_info
 ☒ refund\_payment
 ☒ seckill\_goods
 ☒ sku\_attr\_value
 ☒ sku\_image

☒ sku\_info
 ☒ sku\_sale\_attr\_value
 ☒ spu\_image
 ☒ spu\_info
 ☒ spu\_poster
 ☒ spu\_sale\_attr
 ☒ spu\_sale\_attr\_value
 ☒ user\_address
 ☒ user\_info
 ☒ ware\_info
 ☒ ware\_order\_task
 ☒ ware\_order\_task\_detail
 ☒ ware\_sku

☐ 导入数据类型名称
 ☐ 自动大小写
 ☐ 注释转为逻辑名

☒ 覆盖模型中的同名表

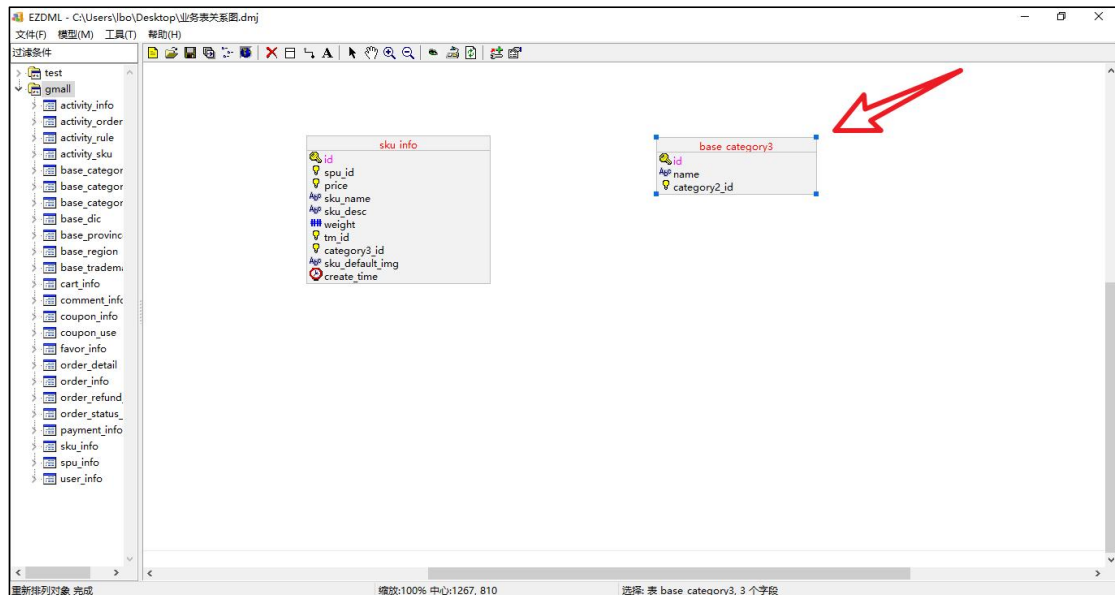
执行进度

0/0

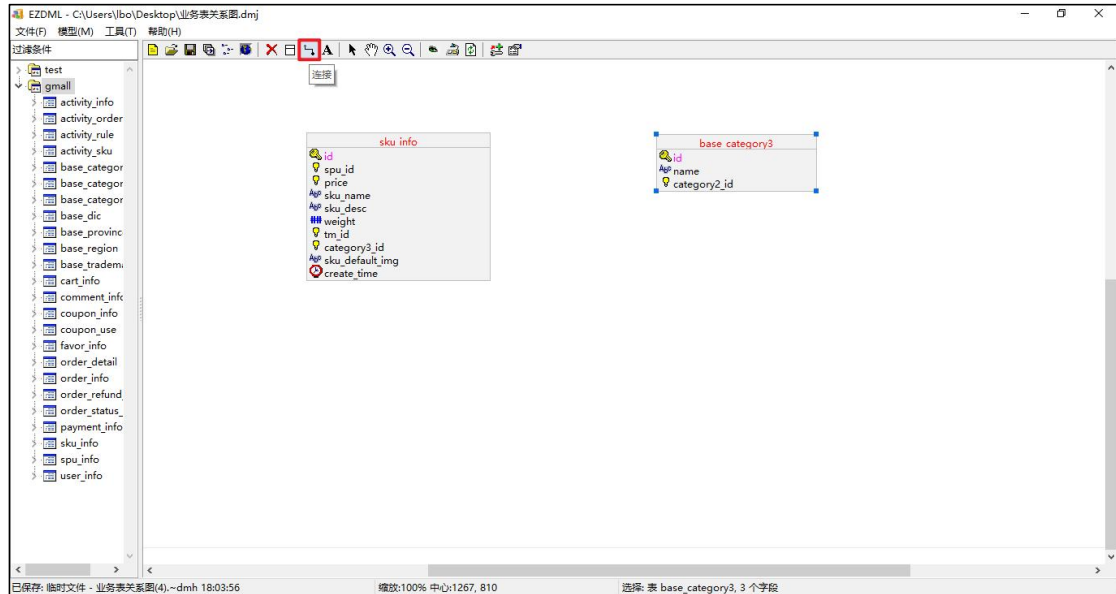
开始 关闭

## （7）建立表关系

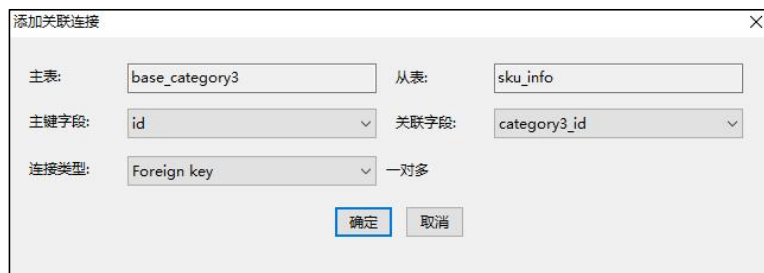
第一步：点击选中主表（主键所在的表）



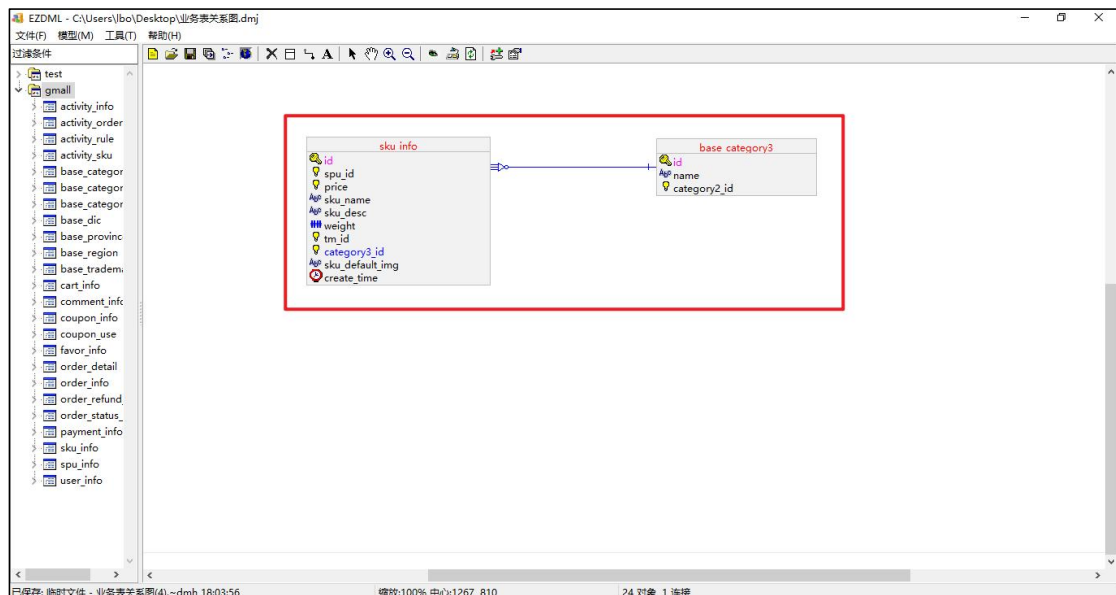
第二步：点击连接按钮



第三步：点击从表，配置连接条件

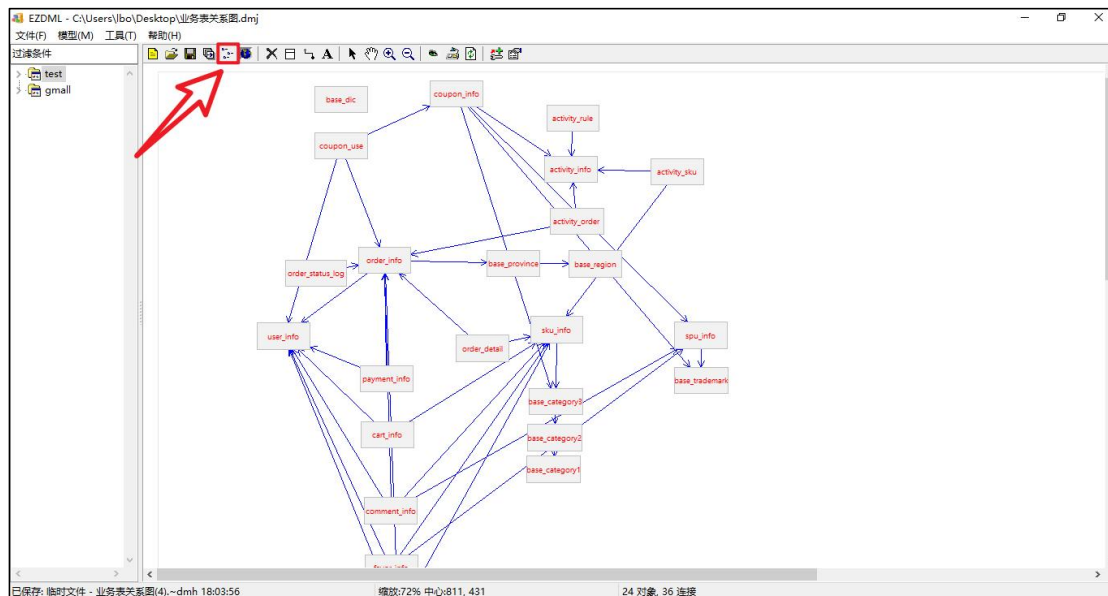


第四步：效果展示



### 3) 使用技巧

#### (1) 缩略图



## （2）热键

按住 shift 键，用鼠标点击表，进行多选，可实现批量移动

按住 ctrl 键，用鼠标圈选表，也可进行多选，实现批量移动

## 第 3 章 业务数据采集模块

### 3.1 业务数据同步概述

#### 3.1.1 数据同步策略概述

业务数据是数据仓库的重要数据来源，我们需要每日定时从业务数据库中抽取数据，传输到数据仓库中，之后再对数据进行分析统计。

为保证统计结果的正确性，需要保证数据仓库中的数据与业务数据库是同步的，离线数仓的计算周期通常为天，所以数据同步周期也通常为天，即每天同步一次即可。

数据的同步策略有**全量同步**和**增量同步**。

**全量同步**，就是每天都将业务数据库中的全部数据同步一份到数据仓库，这是保证两侧数据同步的最简单的方式。

## 全量同步

日期：2020-06-16

业务数据库

id	name
1	张三
2	李小四
3	王五
4	赵六
5	田七

数据仓库

2020-06-14		2020-06-15		2020-06-16	
id	name	id	name	id	name
1	张三	1	张三	1	张三
2	李小四	2	李小四	2	李小四
		3	王五	3	王五
				4	赵六
				5	田七

让天下没有难学的技术

**增量同步**，就是每天只将业务数据中的新增及变化数据同步到数据仓库。采用每日增量同步的表，通常需要在首日先进行一次全量同步。

## 增量同步

日期：2020-06-16

业务数据库

id	name
1	张三
2	李小四
3	王五
4	赵六
5	田七

数据仓库

2020-06-14		2020-06-15		2020-06-16	
id	name	id	name	id	name
1	张三	2	李小四	4	赵六
2	李小四	3	王五	5	田七

让天下没有难学的技术

### 3.1.2 数据同步策略选择

两种策略都能保证数据仓库和业务数据库的数据同步，那应该如何选择呢？下面对两种策略进行简要对比。

同步策略	优点	缺点
全量同步	逻辑简单	在某些情况下效率较低。例如某张表数据量较大，但是每天数据的变化比例很低，若对其采用每日全量同步，则会重复同步

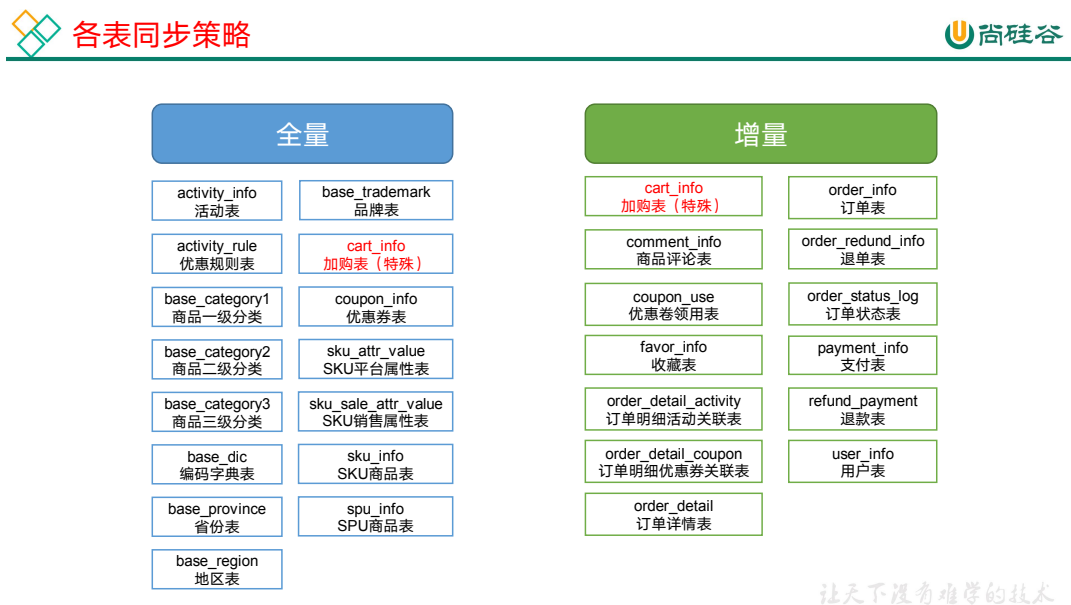


		和存储大量相同的数据。
增量同步	效率高，无需同步和存储重复数据	逻辑复杂，需要将每日的新增及变化数据同原来的数据进行整合，才能使用

根据上述对比，可以得出以下结论：

若业务表数据量比较大，且每天数据变化的比例比较低，这时应采用增量同步，否则可采用全量同步。

下图为各表同步策略：



注：由于后续数仓建模需要，cart\_info 需进行全量同步和增量同步，此处暂不解释，后续章节会作出解释。

## 3.1.3 数据同步工具概述

数据同步工具种类繁多，大致可分为两类，一类是以 DataX、Sqoop 为代表的基于 Select 查询的离线、批量同步工具，另一类是以 Maxwell、Canal 为代表的基于数据库数据变更日志（例如 MySQL 的 binlog，其会实时记录所有的 insert、update 以及 delete 操作）的实时流式同步工具。

全量同步通常使用 DataX、Sqoop 等基于查询的离线同步工具。而增量同步既可以使用 DataX、Sqoop 等工具，也可使用 Maxwell、Canal 等工具，下面对增量同步不同方案进行简要对比。

增量同步方案	DataX/Sqoop	Maxwell/Canal
对数据库的要求	原理是基于查询，故若想通过 select 查询获取新增及变化数	要求数据库记录变更操作，例如 MySQL 需开启 binlog。



	据，就要求数据表中存在 create_time、update_time 等字段，然后根据这些字段获取变更数据。	
数据的中间状态	由于是离线批量同步，故若一条数据在一天中变化多次，该方案只能获取最后一个状态，中间状态无法获取。	由于是实时获取所有的数据变更操作，所以可以获取变更数据的所有中间状态。

本项目中，全量同步采用 DataX，增量同步采用 Maxwell。

## 3.1.4 数据同步工具部署

### 1) DataX



尚硅谷大数据技术  
之DataX.docx

### 2) 基于 KerberosMaxwell




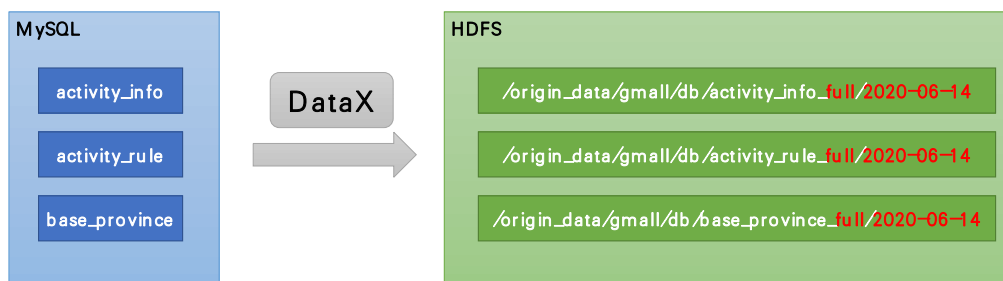
尚硅谷大数据技术  
之Maxwell.docx

## 3.2 全量表数据同步

### 3.2.1 数据通道

全量表数据由 DataX 从 MySQL 业务数据库直接同步到 HDFS，具体数据流向如下图所示。

 全量表数据流向



注：目标路径中表名须包含后缀full，表示该表为全量同步  
目标路径中包含一层日期，用以对不同天的数据进行区分

让天下没有难学的技术

### 3.2.2 DataX 配置文件

我们需要为每张全量表编写一个 DataX 的 json 配置文件，此处以 activity\_info 为例，配置文件内容如下：

```
{
  "job": {
    "content": [
      {
        "reader": {
          "name": "mysqlreader",
          "parameter": {
            "column": [
              "id",
              "activity_name",
              "activity_type",
              "activity_desc",
              "start_time",
              "end_time",
              "create_time"
            ],
            "connection": [
              {
                "jdbcUrl": [
                  "jdbc:mysql://hadoop102:3306/gmall"
                ],
                "table": [
                  "activity_info"
                ]
              }
            ],
            "password": "000000",
            "splitPk": "",
            "username": "root"
          }
        },
        "writer": {
          "name": "hdfswriter",
          "parameter": {
            "column": [
              {
                "name": "id",
                "type": "bigint"
              },
              {
                "name": "activity_name",
                "type": "string"
              },
              {
                "name": "activity_type",
                "type": "string"
              },
              {
                "name": "activity_desc",
                "type": "string"
              }
            ]
          }
        }
      }
    ]
  }
}
```

```
{
    {
        "name": "start_time",
        "type": "string"
    },
    {
        "name": "end_time",
        "type": "string"
    },
    {
        "name": "create_time",
        "type": "string"
    }
],
"compress": "gzip",
"defaultFS": "hdfs://hadoop102:8020",
"fieldDelimiter": "\t",
"fileName": "activity_info",
"fileType": "text",
"path": "${targetdir}",
"writeMode": "append"
}

}

],
"setting": {
    "speed": {
        "channel": 1
    }
}
}
```

注：由于目标路径包含一层日期，用于对不同天的数据加以区分，故 `path` 参数并未写死，需在提交任务时通过参数动态传入，参数名称为 **targetdir**。

### 3.2.3 DataX 配置文件生成脚本

方便起见，此处提供了 DataX 配置文件批量生成脚本，脚本内容及使用方式如下。

#### 1) 在 `~/bin` 目录下创建 `gen_import_config.py` 脚本

```
[root@hadoop103 bin]# vim ~/bin/gen_import_config.py
```

脚本内容如下

```
# coding=utf-8
import json
import getopt
import os
import sys
import MySQLdb

#MySQL 相关配置，需根据实际情况作出修改
mysql_host = "hadoop101"
mysql_port = "3306"
mysql_user = "root"
mysql_passwd = "000000"
```

```
#HDFS NameNode 相关配置, 需根据实际情况作出修改
hdfs_nn_nameservice = "nameservice1"
hdfs_ha_namenodes_nameservice = "namenode34,namenode41"
hdfs_nn_name1 = "namenode34"
hdfs_nn_name2 = "namenode41"
hdfs_nn_name1port = "hadoop101:8020"
hdfs_nn_name2port = "hadoop102:8020"
#生成配置文件的目标路径, 可根据实际情况作出修改
output_path = "/opt/module/datax/job/import"

def get_connection():
    return MySQLdb.connect(host=mysql_host, port=int(mysql_port),
user=mysql_user, passwd=mysql_passwd)

def get_mysql_meta(database, table):
    connection = get_connection()
    cursor = connection.cursor()
    sql = "SELECT COLUMN_NAME,DATA_TYPE from
information_schema.COLUMNS WHERE TABLE_SCHEMA=%s AND TABLE_NAME=%s
ORDER BY ORDINAL_POSITION"
    cursor.execute(sql, [database, table])
    fetchall = cursor.fetchall()
    cursor.close()
    connection.close()
    return fetchall

def get_mysql_columns(database, table):
    return map(lambda x: x[0], get_mysql_meta(database, table))

def get_hive_columns(database, table):
    def type_mapping(mysql_type):
        mappings = {
            "bigint": "bigint",
            "int": "bigint",
            "smallint": "bigint",
            "tinyint": "bigint",
            "decimal": "string",
            "double": "double",
            "float": "float",
            "binary": "string",
            "char": "string",
            "varchar": "string",
            "datetime": "string",
            "time": "string",
            "timestamp": "string",
            "date": "string",
            "text": "string"
        }
    return mappings[mysql_type]

    meta = get_mysql_meta(database, table)
    return map(lambda x: {"name": x[0], "type":
type_mapping(x[1].lower())}, meta)
```

```
def generate_json(source_database, source_table):
    job = {
        "job": {
            "setting": {
                "speed": {
                    "channel": 3
                },
                "errorLimit": {
                    "record": 0,
                    "percentage": 0.02
                }
            },
            "content": [{
                "reader": {
                    "name": "mysqlreader",
                    "parameter": {
                        "username": mysql_user,
                        "password": mysql_passwd,
                        "column": get_mysql_columns(source_database,
source_table),
                        "splitPk": "",
                        "connection": [{
                            "table": [source_table],
                            "jdbcUrl": ["jdbc:mysql://" + mysql_host +
":" + mysql_port + "/" + source_database]
                        }]
                    },
                    "writer": {
                        "name": "hdfswriter",
                        "parameter": {
                            "defaultFS": "hdfs://" + hdfs_nn_nameservice ,
                            "hadoopConfig":{
                                "dfs.nameservices": hdfs_nn_nameservice,
                                "dfs.namenode.rpc-address."+hdfs_nn_nameservice+"."+hdfs_nn_name1:
hdfs_nn_name1port,
                                "dfs.namenode.rpc-address."+hdfs_nn_nameservice+"."+hdfs_nn_name2:
hdfs_nn_name2port,
                                "dfs.ha.namenodes."+hdfs_nn_nameservice:
hdfs_ha_namenodes_nameservice,
                                "dfs.client.failover.proxy.provider."+hdfs_nn_nameservice:
"org.apache.hadoop.hdfs.server.namenode.ha.ConfiguredFailoverProx
yProvider"
                            },
                            "fileType": "text",
                            "path": "${targetdir}",
                            "fileName": source_table,
                            "column": get_hive_columns(source_database,
source_table),
                            "writeMode": "append",
                            "fieldDelimiter": "\\t",
                            "compress": "gzip",
                            "haveKerberos": "true",
                            "kerberosKeytabFilePath": "/var/lib/hive/hive.keytab",
                            "kerberosPrincipal": "hive/hive@HADOOP.COM",
```

```

    }
    }
    }]
}
}
if not os.path.exists(output_path):
    os.makedirs(output_path)
    with open(os.path.join(output_path, ".".join([source_database,
source_table, "json"])), "w") as f:
        json.dump(job, f)

def main(args):
    source_database = ""
    source_table = ""

    options, arguments = getopt.getopt(args, '-d:-t:', ['sourcedb=',
'sourcetbl='])
    for opt_name, opt_value in options:
        if opt_name in ('-d', '--sourcedb'):
            source_database = opt_value
        if opt_name in ('-t', '--sourcetbl'):
            source_table = opt_value

    generate_json(source_database, source_table)

if __name__ == '__main__':
    main(sys.argv[1:])

```

注：

### （1）安装 Python Mysql 驱动

由于需要使用 Python 访问 Mysql 数据库，故需安装驱动，命令如下：

```

[root@hadoop101 bin]# cd /opt/software/mysql-lib
[[root@hadoop101 mysql-lib]# yum install -y MySQL-python

```

### （2）脚本使用说明

```
python gen_import_config.py -d database -t table
```

通过-d 传入数据库名，-t 传入表名，执行上述命令即可生成该表的 DataX 同步配置文件。

## 2) 在~/bin 目录下创建 gen\_import\_config.sh 脚本

```
[root@hadoop103 bin]# vim ~/bin/gen_import_config.sh
```

脚本内容如下

```

#!/bin/bash

python ~/bin/gen_import_config.py -d gmall -t activity_info
python ~/bin/gen_import_config.py -d gmall -t activity_rule
python ~/bin/gen_import_config.py -d gmall -t base_category1
python ~/bin/gen_import_config.py -d gmall -t base_category2
python ~/bin/gen_import_config.py -d gmall -t base_category3
python ~/bin/gen_import_config.py -d gmall -t base_dic
python ~/bin/gen_import_config.py -d gmall -t base_province
python ~/bin/gen_import_config.py -d gmall -t base_region
python ~/bin/gen_import_config.py -d gmall -t base_trademark

```

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：[尚硅谷官网](#)

```
python ~/bin/gen_import_config.py -d gmall -t cart_info
python ~/bin/gen_import_config.py -d gmall -t coupon_info
python ~/bin/gen_import_config.py -d gmall -t sku_attr_value
python ~/bin/gen_import_config.py -d gmall -t sku_info
python ~/bin/gen_import_config.py -d gmall -t sku_sale_attr_value
python ~/bin/gen_import_config.py -d gmall -t spu_info
```

### 3) 为 gen\_import\_config.sh 脚本增加执行权限

```
[root@hadoop103 bin]# chmod +x ~/bin/gen_import_config.sh
```

### 4) 执行 gen\_import\_config.sh 脚本，生成配置文件

```
[root@hadoop103 bin]# gen_import_config.sh
```

### 5) 观察生成的配置文件

```
[root@hadoop101 bin]# ll /opt/module/datax/job/import/
总用量 60
-rw-rw-r-- 1 atguigu atguigu 957 10 月 15 22:17 gmall.activity_info.json
-rw-rw-r-- 1 atguigu atguigu 1049 10 月 15 22:17 gmall.activity_rule.json
-rw-rw-r-- 1 atguigu atguigu 651 10 月 15 22:17 gmall.base_category1.json
-rw-rw-r-- 1 atguigu atguigu 711 10 月 15 22:17 gmall.base_category2.json
-rw-rw-r-- 1 atguigu atguigu 711 10 月 15 22:17 gmall.base_category3.json
-rw-rw-r-- 1 atguigu atguigu 835 10 月 15 22:17 gmall.base_dic.json
-rw-rw-r-- 1 atguigu atguigu 865 10 月 15 22:17 gmall.base_province.json
-rw-rw-r-- 1 atguigu atguigu 659 10 月 15 22:17 gmall.base_region.json
-rw-rw-r-- 1 atguigu atguigu 709 10 月 15 22:17 gmall.base_trademark.json
-rw-rw-r-- 1 atguigu atguigu 1301 10 月 15 22:17 gmall.cart_info.json
-rw-rw-r-- 1 atguigu atguigu 1545 10 月 15 22:17 gmall.coupon_info.json
-rw-rw-r-- 1 atguigu atguigu 867 10 月 15 22:17 gmall.sku_attr_value.json
-rw-rw-r-- 1 atguigu atguigu 1121 10 月 15 22:17 gmall.sku_info.json
-rw-rw-r-- 1 atguigu atguigu 985 10 月 15 22:17 gmall.sku_sale_attr_value.json
-rw-rw-r-- 1 atguigu atguigu 811 10 月 15 22:17 gmall.spu_info.json
```

## 3.2.4 测试生成的 DataX 配置文件

以 activity\_info 为例，测试用脚本生成的配置文件是否可用。

### 1) 创建目标路径

由于 DataX 同步任务要求目标路径提前存在，故需手动创建路径，当前 activity\_info 表的目标路径应为 origin\_data/gmall/db/activity\_info\_full/2020-06-14。

```
[root@hadoop103 ~]# kinit hive/hive
[root@hadoop103 bin]# hadoop fs -mkdir /origin_data/gmall/db/activity_info_full/2020-06-14
```

### 2) 执行 DataX 同步命令

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：[尚硅谷官网](#)

```
[root@hadoop103 bin]# python /opt/module/datax/bin/datax.py
-p"-Dtargetdir=/origin_data/gmall/db/activity_info_full/2020-06-14"
/opt/module/datax/job/import/gmall.activity_info.json
```

### 3) 观察同步结果

观察 HFDS 目标路径是否出现数据。

## 3.2.5 全量表数据同步脚本

为方便使用以及后续的任务调度，此处编写一个全量表数据同步脚本。

### 1) 在~/bin 目录创建 mysql\_to\_hdfs\_full.sh

```
[root@hadoop103 bin]# vim ~/bin/mysql_to_hdfs_full.sh
```

脚本内容如下

```
#!/bin/bash
kinit -kt /var/lib/hive/hive.keytab hive/hive

DATAX_HOME=/opt/module/datax

# 如果传入日期则 do_date 等于传入的日期，否则等于前一天日期
if [ -n "$2" ] ;then
    do_date=$2
else
    do_date=`date -d "-1 day" +%F`
fi

#处理目标路径，此处的处理逻辑是，如果目标路径不存在，则创建；若存在，则清空，目的是保证同步任务可重复执行
handle_targetdir() {
    hadoop fs -test -e $1
    if [[ $? -eq 1 ]]; then
        echo "路径$1 不存在，正在创建....."
        hadoop fs -mkdir -p $1
    else
        echo "路径$1 已经存在"
        fs_count=$(hadoop fs -count $1)
        content_size=$(echo $fs_count | awk '{print $3}')
        if [[ $content_size -eq 0 ]]; then
            echo "路径$1 为空"
        else
            echo "路径$1 不为空，正在清空....."
            hadoop fs -rm -r -f $1/*
        fi
    fi
}

#数据同步
import_data() {
    datax_config=$1
    target_dir=$2

    handle_targetdir $target_dir
    python $DATAX_HOME/bin/datax.py -p"-Dtargetdir=$target_dir"
    $datax_config
}
```

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：尚硅谷官网



```
}

case $1 in
"activity_info")
    import_data
/opt/module/datax/job/import/gmall.activity_info.json
/origin_data/gmall/db/activity_info_full/$do_date
;;
"activity_rule")
    import_data
/opt/module/datax/job/import/gmall.activity_rule.json
/origin_data/gmall/db/activity_rule_full/$do_date
;;
"base_category1")
    import_data
/opt/module/datax/job/import/gmall.base_category1.json
/origin_data/gmall/db/base_category1_full/$do_date
;;
"base_category2")
    import_data
/opt/module/datax/job/import/gmall.base_category2.json
/origin_data/gmall/db/base_category2_full/$do_date
;;
"base_category3")
    import_data
/opt/module/datax/job/import/gmall.base_category3.json
/origin_data/gmall/db/base_category3_full/$do_date
;;
"base_dic")
    import_data /opt/module/datax/job/import/gmall.base_dic.json
/origin_data/gmall/db/base_dic_full/$do_date
;;
"base_province")
    import_data
/opt/module/datax/job/import/gmall.base_province.json
/origin_data/gmall/db/base_province_full/$do_date
;;
"base_region")
    import_data /opt/module/datax/job/import/gmall.base_region.json
/origin_data/gmall/db/base_region_full/$do_date
;;
"base_trademark")
    import_data
/opt/module/datax/job/import/gmall.base_trademark.json
/origin_data/gmall/db/base_trademark_full/$do_date
;;
"cart_info")
    import_data /opt/module/datax/job/import/gmall.cart_info.json
/origin_data/gmall/db/cart_info_full/$do_date
;;
"coupon_info")
    import_data /opt/module/datax/job/import/gmall.coupon_info.json
/origin_data/gmall/db/coupon_info_full/$do_date
;;
"sku_attr_value")
    import_data
/opt/module/datax/job/import/gmall.sku_attr_value.json
/origin_data/gmall/db/sku_attr_value_full/$do_date
```

```
;;
"sku_info")
    import_data      /opt/module/datax/job/import/gmall.sku_info.json
/origin_data/gmall/db/sku_info_full/$do_date
;;
"sku_sale_attr_value")
    import_data
/opt/module/datax/job/import/gmall.sku_sale_attr_value.json
/origin_data/gmall/db/sku_sale_attr_value_full/$do_date
;;
"spu_info")
    import_data      /opt/module/datax/job/import/gmall.spu_info.json
/origin_data/gmall/db/spu_info_full/$do_date
;;
"all")
    import_data
/opt/module/datax/job/import/gmall.activity_info.json
/origin_data/gmall/db/activity_info_full/$do_date
    import_data
/opt/module/datax/job/import/gmall.activity_rule.json
/origin_data/gmall/db/activity_rule_full/$do_date
    import_data
/opt/module/datax/job/import/gmall.base_category1.json
/origin_data/gmall/db/base_category1_full/$do_date
    import_data
/opt/module/datax/job/import/gmall.base_category2.json
/origin_data/gmall/db/base_category2_full/$do_date
    import_data
/opt/module/datax/job/import/gmall.base_category3.json
/origin_data/gmall/db/base_category3_full/$do_date
    import_data      /opt/module/datax/job/import/gmall.base_dic.json
/origin_data/gmall/db/base_dic_full/$do_date
    import_data
/opt/module/datax/job/import/gmall.base_province.json
/origin_data/gmall/db/base_province_full/$do_date
    import_data      /opt/module/datax/job/import/gmall.base_region.json
/origin_data/gmall/db/base_region_full/$do_date
    import_data
/opt/module/datax/job/import/gmall.base_trademark.json
/origin_data/gmall/db/base_trademark_full/$do_date
    import_data      /opt/module/datax/job/import/gmall.cart_info.json
/origin_data/gmall/db/cart_info_full/$do_date
    import_data      /opt/module/datax/job/import/gmall.coupon_info.json
/origin_data/gmall/db/coupon_info_full/$do_date
    import_data
/opt/module/datax/job/import/gmall.sku_attr_value.json
/origin_data/gmall/db/sku_attr_value_full/$do_date
    import_data      /opt/module/datax/job/import/gmall.sku_info.json
/origin_data/gmall/db/sku_info_full/$do_date
    import_data
/opt/module/datax/job/import/gmall.sku_sale_attr_value.json
/origin_data/gmall/db/sku_sale_attr_value_full/$do_date
    import_data      /opt/module/datax/job/import/gmall.spu_info.json
/origin_data/gmall/db/spu_info_full/$do_date
;;
esac
```

## 2) 为 mysql\_to\_hdfs\_full.sh 增加执行权限

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：尚硅谷官网

```
[root@hadoop103 bin]# chmod 777 ~/bin/mysql_to_hdfs_full.sh
```

### 3) 测试同步脚本

```
[root@hadoop103 bin]# mysql_to_hdfs_full.sh all 2020-06-14
```

### 4) 检查同步结果

查看 HDFS 目录路径是否出现全量表数据，全量表共 15 张。

## 3.2.6 全量表同步总结

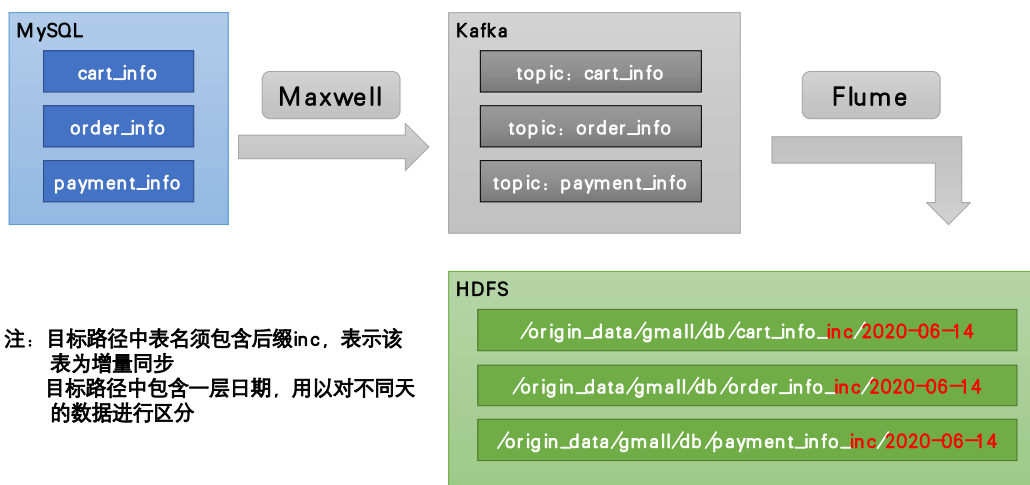
全量表同步逻辑比较简单，只需每日执行全量表数据同步脚本 `mysql_to_hdfs_full.sh` 即可。

## 3.3 增量表数据同步

### 3.3.1 数据通道



#### 增量表数据通道



让天下没有难学的技术

### 3.3.2 Maxwell 配置

按照规划，有 `cart_info`、`comment_info` 等共计 13 张表需进行增量同步，默认情况下，Maxwell 会同步 binlog 中的所有表的数据变更记录，因此我们需要对 Maxwell 进行配置，另其只同步这特定的 13 张表。

另外，为方便下游使用数据，还需对 Maxwell 进行配置，另其将不同表的数据发往不同的 Kafka Topic。Maxwell 最终配置如下：

#### 1) 修改 Maxwell 配置文件 `config.properties`

```
[root@hadoop103 conf]# cd /opt/module/maxwell/  
[root@hadoop103 maxwell]# vim /opt/module/maxwell/config.properties
```

#### 2) 全部配置参数如下

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：尚硅谷官网

```
log_level=info

producer=kafka
kafka.bootstrap.servers=hadoop103:9092,hadoop104:9092

#kafka topic 动态配置
kafka_topic=topic_db
#kerberos 相关配置
kafka.security.protocol=SASL_PLAINTEXT
kafka.sasl.kerberos.service.name=kafka

# mysql login info
host=hadoop101
user=maxwell
password=maxwell
jdbc_options=useSSL=false&serverTimezone=Asia/Shanghai
```

### 3) 重新启动 Maxwell

```
[root@hadoop103 maxwell]# mxw.sh restart
```

### 4) 通道测试

#### (1) 启动 Zookeeper 以及 Kafka 集群

#### (2) 启动一个 Kafka Console Consumer，消费任一 topic 数据

```
[root@hadoop103 conf]# export
KAFKA_OPTS="-Djava.security.auth.login.config=/var/lib/hive/jaas.
conf"
[root@hadoop103 maxwell]#
/opt/cloudera/parcels/CDH/bin/kafka-console-consumer
--bootstrap-server hadoop103:9092 --topic topic_db
--consumer.config /etc/kafka/conf/consumer.properties
```

#### (3) 生成模拟数据

```
[root@hadoop103 bin]# cd /opt/module/db_log/
[root@hadoop103 db_log]# java -jar gmall2020-mock-db-2022-03-29.jar
```

#### (4) 观察 Kafka 消费者是否能消费到数据

```
{"database":"gmall","table":"cart_info","type":"update","ts":1592
270938,"xid":13090,"xoffset":1573,"data":{"id":100924,"user_id":"
93","sku_id":16,"cart_price":4488.00,"sku_num":1,"img_url":"http:
//47.93.148.192:8080/group1/M00/00/02/rBHu8l-sklaALrngAAHGDqdpFtU
741.jpg","sku_name":"华为 HUAWEI P40 麒麟 990 5G SoC 芯片 5000 万超感知
徕卡三摄 30 倍数字变焦 8GB+128GB 亮黑色全网通 5G 手机
","is_checked":null,"create_time":"2020-06-14
09:28:57","operate_time":null,"is_ordered":1,"order_time":"2021-1
0-17
09:28:58","source_type":"2401","source_id":null},"old":{"is_order
ed":0,"order_time":null}}
```

## 3.3.3 Flume 配置

### 1) Flume 配置概述

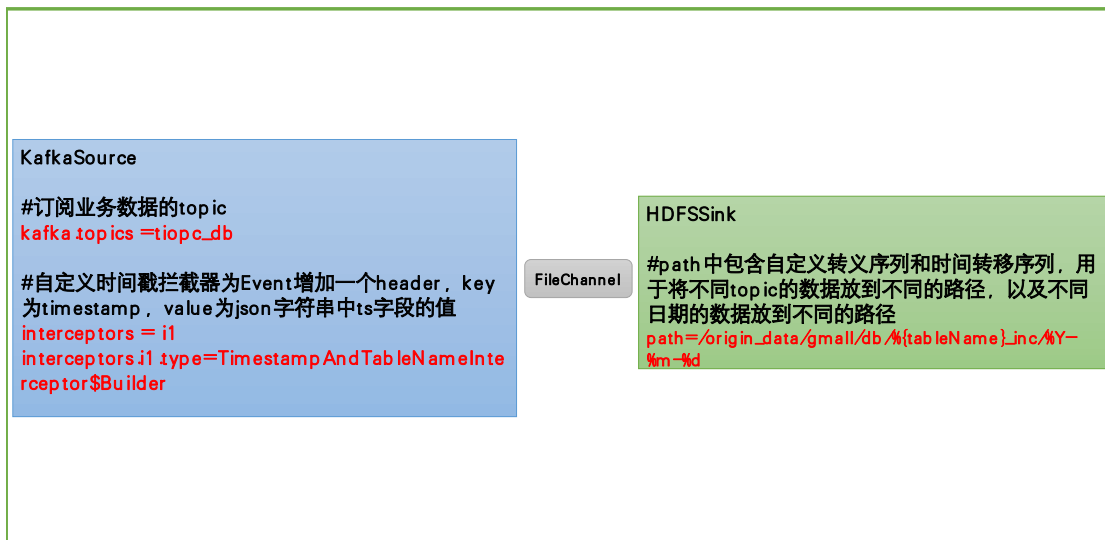
Flume 需要将 Kafka 中 topic\_db 主题的数据传输到 HDFS，故其需选用 KafkaSource 以及 HDFSSink，Channel 选用 FileChannel。

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：[尚硅谷官网](#)

需要注意的是，HDFS Sink 需要将不同 mysql 业务表的数据写到不同的路径，并且路径中应当包含一层日期，用于区分每天的数据。关键配置如下：



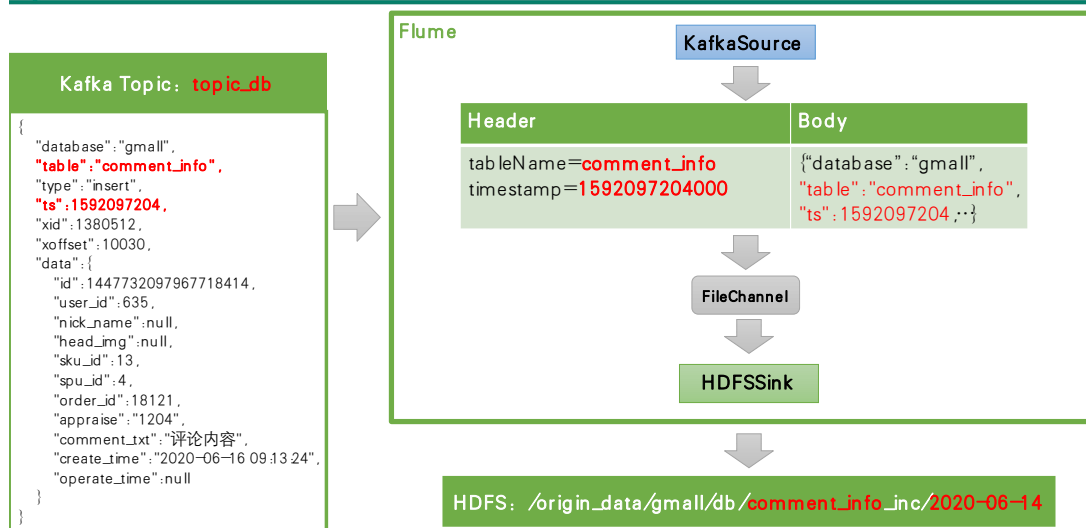
## Flume配置重点



具体数据示例如下：



## Flume数据示例



让天下没有难学的技术

## 2) Flume 配置实操

### 创建一个 Flume 服务

Cloudera Manager

将服务添加到 Cluster 1

选择您要添加的服务类型。

服务类型	说明
<input type="radio"/> ADLS Connector	The ADLS Connector service provides key management for accessing ADLS Gen1 accounts and ADLS Gen2 containers from CDH services.
<input type="radio"/> Accumulo	The Apache Accumulo sorted, distributed key/value store is a robust, scalable, high performance data storage and retrieval system. This service only works with releases meant to run on top of CDH6.
<input checked="" type="radio"/> Flume	Flume 从几乎所有来源收集数据并将这些数据聚合到永久性存储（如 HDFS）中。
<input type="radio"/> HBase	Apache HBase 提供对大型数据集的随机、实时的读/写访问权限（需要 HDFS 和 ZooKeeper）。

Cloudera Manager

将 Flume 服务添加到 Cluster 1

Select Dependencies

自定义角色分配

审核更改

命令详细信息

汇总

Select Dependencies

HDFS	Kafka	ZooKeeper
<input type="radio"/>	Kafka	ZooKeeper
<input type="radio"/> HDFS		ZooKeeper
<input checked="" type="radio"/> HDFS	Kafka	ZooKeeper

Cloudera Manager

将 Flume 服务添加到 Cluster 1

Select Dependencies

**自定义角色分配**

审核更改

命令详细信息

汇总

自定义角色分配

您可以在此处自定义新服务的角色分配，但请注意，如果分配不正确（例如，分配到某个主机上的角色太多），性能将受到影响。

还可以按主机查看角色分配。 [按主机查看](#)

Agent × 1 新建

hadoop105

返回

继续




The screenshot shows the Cloudera Manager interface for the Flume-2 service. The '实例' (Instances) tab is selected, showing a table of instances. The '角色类型' (Role Type) column is highlighted with a red box, showing 'Agent' as the selected role type. The table lists the following details for the Agent instance:

角色类型	状态	主机	授权状态	角色组
Agent	已停止	hadoop105	已授权	Agent Default Group

## (1) 创建 Flume 配置文件

a2

## (2) 配置文件内容如下

```
a2.sources = r1
a2.channels = c1
a2.sinks = k1

a2.sources.r1.type = org.apache.flume.source.kafka.KafkaSource

a2.sources.r1.batchSize = 5000
a2.sources.r1.batchDurationMillis = 2000
a2.sources.r1.kafka.bootstrap.servers = hadoop103:9092,hadoop104:9092
a2.sources.r1.kafka.topics = topic_db
a2.sources.r1.kafka.consumer.group.id = flume
a2.sources.r1.setTopicHeader = true
```

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：尚硅谷官网

```
a2.sources.r1.topicHeader = topic
a2.sources.r1.interceptors = i1
a2.sources.r1.interceptors.i1.type =
com.atguigu.gmall.flume.interceptor.db.TimestampAndTableNameInter
ceptor$Builder

a2.channels.c1.type = file
a2.channels.c1.checkpointDir = /opt/flume/checkpoint/behavior2
a2.channels.c1.dataDirs = /opt/flume/data/behavior2/
a2.channels.c1.maxFileSize = 2146435071
a2.channels.c1.capacity = 1000000
a2.channels.c1.keep-alive = 6

## sink1
a2.sinks.k1.type = hdfs
#a2.sinks.k1.hdfs.proxyUser=hive
a2.sinks.k1.hdfs.kerberosPrincipal=hive/hive@HADOOP.COM
a2.sinks.k1.hdfs.kerberosKeytab=/var/lib/hive/hive.keytab

a2.sinks.k1.hdfs.path =
/origin_data/gmall/db/{tableName}_inc/%Y-%m-%d
a2.sinks.k1.hdfs.filePrefix = db
a2.sinks.k1.hdfs.round = false

a2.sinks.k1.hdfs.rollInterval = 10
a2.sinks.k1.hdfs.rollSize = 134217728
a2.sinks.k1.hdfs.rollCount = 0

a2.sinks.k1.hdfs.fileType = CompressedStream
a2.sinks.k1.hdfs.codeC = gzip

## 拼装
a2.sources.r1.channels = c1
a2.sinks.k1.channel= c1
```

### （3）编写 Flume 拦截器

#### ① 新建一个 Maven 项目，并在 pom.xml 文件中加入如下配置

```
<dependencies>
  <dependency>
    <groupId>org.apache.flume</groupId>
    <artifactId>flume-ng-core</artifactId>
    <version>1.9.0</version>
    <scope>provided</scope>
  </dependency>

  <dependency>
    <groupId>com.alibaba</groupId>
    <artifactId>fastjson</artifactId>
    <version>1.2.62</version>
  </dependency>
</dependencies>

<build>
  <plugins>
```



```
<plugin>
  <artifactId>maven-compiler-plugin</artifactId>
  <version>2.3.2</version>
  <configuration>
    <source>1.8</source>
    <target>1.8</target>
  </configuration>
</plugin>
<plugin>
  <artifactId>maven-assembly-plugin</artifactId>
  <configuration>
    <descriptorRefs>

<descriptorRef>jar-with-dependencies</descriptorRef>
    </descriptorRefs>
  </configuration>
  <executions>
    <execution>
      <id>make-assembly</id>
      <phase>package</phase>
      <goals>
        <goal>single</goal>
      </goals>
    </execution>
  </executions>
</plugin>
</plugins>
</build>
```

② 在 com.atguigu.gmall.flume.interceptor.db 包 下 创 建

TimestampAndTableNameInterceptor 类

```
package com.atguigu.gmall.flume.interceptor.db;

import com.alibaba.fastjson.JSONObject;
import org.apache.flume.Context;
import org.apache.flume.Event;
import org.apache.flume.interceptor.Interceptor;

import java.nio.charset.StandardCharsets;
import java.util.List;
import java.util.Map;

public class TimestampAndTableNameInterceptor implements Interceptor
{
    @Override
    public void initialize() {

    }

    @Override
    public Event intercept(Event event) {

        Map<String, String> headers = event.getHeaders();
        String log = new String(event.getBody(), StandardCharsets.UTF_8);

        JSONObject jsonObject = JSONObject.parseObject(log);
```

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：尚硅谷官网

```
Long ts = jsonObject.getLong("ts");
//Maxwell 输出的数据中的 ts 字段时间戳单位为秒，Flume HDFSSink 要求单
单位为毫秒
String timeMills = String.valueOf(ts * 1000);

String tableName = jsonObject.getString("table");

headers.put("timestamp", timeMills);
headers.put("tableName", tableName);
return event;
}

@Override
public List<Event> intercept(List<Event> events) {

    for (Event event : events) {
        intercept(event);
    }

    return events;
}

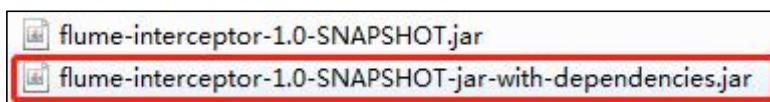
@Override
public void close() {
}

public static class Builder implements Interceptor.Builder {

    @Override
    public Interceptor build() {
        return new TimestampAndTableNameInterceptor ();
    }

    @Override
    public void configure(Context context) {
    }
}
}
```

### ③ 重新打包

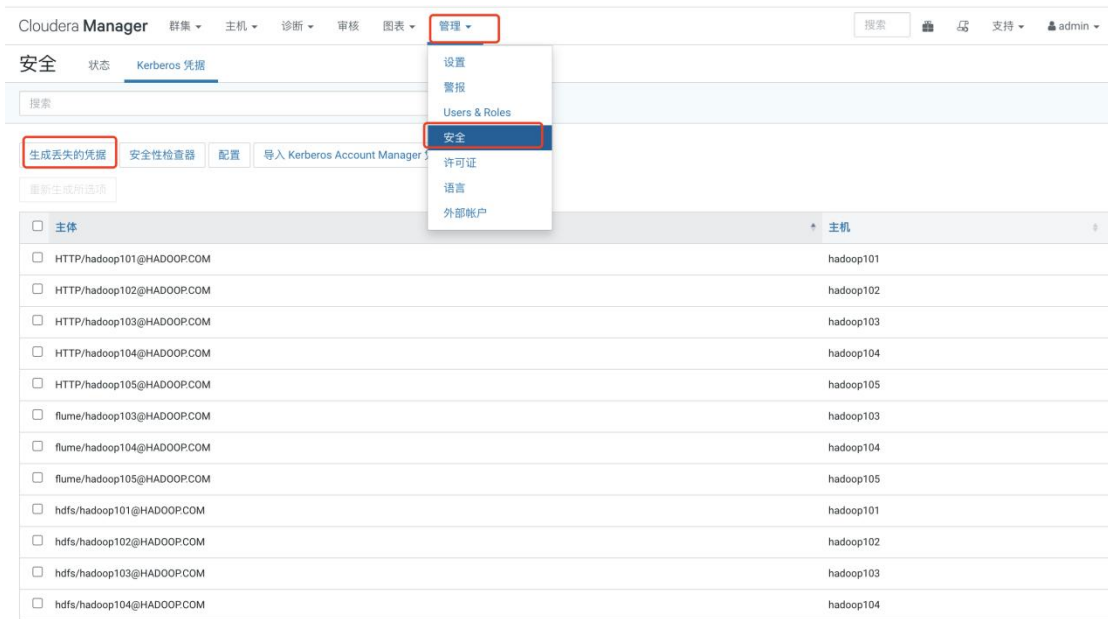


### ④ 将打好的包放入到 hadoop105 的/opt/cloudera/parcels/CDH/lib/flume-ng/lib 文件夹下

```
[root@hadoop105 ~]# cd /opt/cloudera/parcels/CDH/lib/flume-ng/lib
[root@hadoop105 lib]# ls | grep interceptor
flume-interceptor-1.0-SNAPSHOT-jar-with-dependencies.jar
```

## 3) 通道测试

## （1）生成 kerberos 凭证



## （2）启动 Zookeeper、Kafka 集群

## （3）在 hadoop105 上进行认证

```
[root@hadoop105 flume-ng]# kinit hive/hive
```

## （4）启动 hadoop105 上面的 flume

## （5）生成模拟数据

```
[root@hadoop103 flume-ng]# cd /opt/module/db_log/  
[root@hadoop103 db_log]# java -jar gmall2020-mock-db-2022-03-29.jar
```

## （6）观察 HDFS 上的目标路径是否有数据出现

若 HDFS 上的目标路径已有增量表的数据出现了，就证明数据通道已经打通。

## （7）数据目标路径的日期说明

仔细观察，会发现目标路径中的日期，并非模拟数据的业务日期，而是当前日期。这是由于 Maxwell 输出的 JSON 字符串中的 ts 字段的值，是数据的变动日期。而真实场景下，数据的业务日期与变动日期应当是一致的。

此处为了模拟真实环境，对 Maxwell 源码进行了改动，增加了一个参数 `mock_date`，该参数的作用就是指定 Maxwell 输出 JSON 字符串的 ts 时间戳的日期，接下来进行测试。

### ① 修改 Maxwell 配置文件 `config.properties`，增加 `mock_date` 参数，如下

```
#该日期须和/opt/module/db_log/application.properties 中的 mock.date 参  
数保持一致  
mock_date=2020-06-14
```

注：该参数仅供学习使用，**修改该参数后重启 Maxwell 才可生效。**

### ② 重启 Maxwell

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：[尚硅谷官网](#)

```
[root@hadoop103 bin]# mxw.sh restart
```

③ 重新生成模拟数据

```
[root@hadoop103 bin]# cd /opt/module/db_log/  
[root@hadoop103 db_log]# java -jar gmall2020-mock-db-2022-03-29.jar
```

④ 观察 HDFS 目标路径日期是否正常

### 3.3.4 增量表首日全量同步

通常情况下，增量表需要在首日进行一次全量同步，后续每日再进行增量同步，首日全量同步可以使用 Maxwell 的 bootstrap 功能，方便起见，下面编写一个增量表首日全量同步脚本。

1) 在~/bin 目录创建 mysql\_to\_kafka\_inc\_init.sh

```
[root@hadoop103 bin]# vim mysql_to_kafka_inc_init.sh
```

脚本内容如下

```
#!/bin/bash  
  
kinit -kt /var/lib/hive/hive.keytab hive/hive  
  
# 该脚本的作用是初始化所有的增量表，只需执行一次  
  
MAXWELL_HOME=/opt/module/maxwell  
  
import_data() {  
    $MAXWELL_HOME/bin/maxwell-bootstrap --database gmall --table $1  
    --config $MAXWELL_HOME/config.properties  
}  
  
case $1 in  
    "cart_info")  
        import_data cart_info  
        ;;  
    "comment_info")  
        import_data comment_info  
        ;;  
    "coupon_use")  
        import_data coupon_use  
        ;;  
    "favor_info")  
        import_data favor_info  
        ;;  
    "order_detail")  
        import_data order_detail  
        ;;  
    "order_detail_activity")  
        import_data order_detail_activity  
        ;;  
    "order_detail_coupon")  
        import_data order_detail_coupon  
        ;;  
    *)  
        echo "Invalid table name: $1"  
        exit 1  
    esac
```

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：尚硅谷官网

```
"order_info")
    import_data order_info
;;
"order_refund_info")
    import_data order_refund_info
;;
"order_status_log")
    import_data order_status_log
;;
"payment_info")
    import_data payment_info
;;
"refund_payment")
    import_data refund_payment
;;
"user_info")
    import_data user_info
;;
"all")
    import_data cart_info
    import_data comment_info
    import_data coupon_use
    import_data favor_info
    import_data order_detail
    import_data order_detail_activity
    import_data order_detail_coupon
    import_data order_info
    import_data order_refund_info
    import_data order_status_log
    import_data payment_info
    import_data refund_payment
    import_data user_info
;;
esac
```

## 2) 为 mysql\_to\_kafka\_inc\_init.sh 增加执行权限

```
[root@hadoop103 bin]# chmod +x ~/bin/mysql_to_kafka_inc_init.sh
```

## 3) 测试同步脚本

### (1) 清理历史数据

为方便查看结果，现将 HDFS 上之前同步的增量表数据删除

```
[root@hadoop103 ~]# hadoop fs -ls /origin_data/gmall/db | grep _inc
| awk '{print $8}' | xargs hadoop fs -rm -r -f
```

### (2) 执行同步脚本

```
[root@hadoop103 bin]# mysql_to_kafka_inc_init.sh all
```

## 4) 检查同步结果

观察 HDFS 上是否重新出现增量表数据。

## 3.3.5 增量表同步总结

增量表同步，需要在首日进行一次全量同步，后续每日才是增量同步。首日进行全量同步时，需先启动数据通道，包括 Maxwell、Kafka、Flume，然后执行增量表首日同步脚本

更多 [Java](#) - [大数据](#) - [前端](#) - [python](#) 人工智能资料下载，可百度访问：[尚硅谷官网](#)

mysql\_to\_kafka\_inc\_init.sh 进行同步。后续每日只需保证采集通道正常运行即可，Maxwell 便会实时将变动数据发往 Kafka。