

ClickHouse->Apache Doris(Local File)

1、环境配置

0、前置工作

- 搭建Doris
- 准备好HDFS/对象存储
- 在Doris建好对应的表（字段名称大小写需和ClickHouse相同）

1、下载ClickHouse导出工具

```
git clone https://github.com/LOVEGISER/clickhouse\_export.git
```

2、检查环境是否有python3

2、配置工具

```
1 #进入目录
2 cd clickhouse_export/clickhouse_python_sink/
3 #打开配置文件
4 vim config.py
```

```
1 # -*- encoding=utf8 -*-
2
3 """
4 -----
5 @author: "wanglei@flywheels.com"
6 @file: config.py
7 @time: 2022-07-27
8 @desc: clickhouse_python_sink Server Config
9 -----
10 """
11 from log_utils import logger
12 #1. which table should want to been export
13 export_table_list = [
14     {
15         "db": "ssb", #指定数据库
```

```

16     "table": "lineorder", #指定表
17     "format": "Parquet", #Parquet/CSVWithNames #指定导出格式
18     "filenameExtension": "Parquet", #Parquet/csv #指定导出文件后缀名
19     "mode": "partition", #导出模式, 可选按照partition分区导出或者all全量导出
20     "partition_expr": "toYear(LO_ORDERDATE)", #指定的分区字段
21     "upper_condition": "toYear(LO_ORDERDATE)<=999912", #结束分区
22     "lower_condition": "toYear(LO_ORDERDATE)>=000001", #开始分区
23     "partition_split_filed": "LO_TAX", #对数据量过大的partition按照split_filed再次分区
24     "partition_split_filed_model": "continuous", #continuous:连续型 (数据按照字段连续)
25     "partition_split_filed_type": "long" #datetime/long/date
26 }
27 # 可以同时指定多个表
28 # ,{
29 #     "db": "default",
30 #     "table": "trips_np",
31 #     "format": "Parquet",
32 #     "mode": "all",
33 #     "partition_expr": "",
34 #     "upper_condition": "",
35 #     "lower_condition": "",
36 #     "partition_split_filed": "",
37 #     "partition_split_filed_model": "",
38 #     "partition_split_filed_type": ""
39 # }
40 ]
41 #例子
42 '''
43 example:
44 export_table_list = [
45     {
46         "db": "default",
47         "table": "trips",
48         "mode": "all/partition",
49         "partition_expr": "toYYYYMM(pickup_date)",
50         "upper_condition": "toYYYYMM(pickup_date)<=201508",
51         "lower_condition": "toYYYYMM(pickup_date)>=201506",
52     }
53 ]
54 '''
55 #example ./clickhouse-client --host=<host> --port=<port> --user=<user> --password=<password>
56 clickhouse_connect_command = "clickhouse-client --host=172.16.70.243 --port=9000"
57 #2.thread number use
58 process_number = 10
59 sub_partition_max_size=200000
60 #指定clickhouse的userfiles路径, 导出的文件在此
61 user_files_path = "/mnt/data/clickhouse/user_files"

```

3、运行导出

```
python3 scheduler.py
```

4、数据上传HDFS/对象存储

由clickhouse-client导出的Parquet文件，并不能直接导入到Doris中，需要上传到HDFS/对象存储再由Doris进行Broker Load/S3 Load导入

上传HDFS可用HDFS命令

1、上传对象存储

对象存储上传文件参考如下

腾讯COS: <https://cloud.tencent.com/document/product/436/10976>

阿里OSS: https://help.aliyun.com/document_detail/50451.html

此处以OSS为例:

```
1 #先找到clickhouse的user_file目录或者在python脚本中自定义的数据文件目录
2 #用ossutil64上传整个目录到存储桶
3 ./ossutil64 cp /mnt/data/clickhouse/user_files/ssb/lineorder oss://ck-doris/li
  neorder
```

2、上传HDFS

5、Doris S3 Load/Broker Load 数据

1、S3 Load

S3 Load针对ck导出的数据存到对象存储上

参考: <https://doris.apache.org/zh-CN/docs/dev/data-operate/import/import-way/s3-load-manual>

```
1 LOAD LABEL ssb.lineorder2
2 (
3     DATA INFILE("s3://ck-doris/lineorder/*")
4     INTO TABLE lineorder
5     FORMAT AS "parquet"
6 )
7 WITH S3
8 (
```

```
9      "AWS_ENDPOINT" = "*",
10     "AWS_ACCESS_KEY" = "*",
11     "AWS_SECRET_KEY"="*",
12     "AWS_REGION" = "*"
13   )
14  PROPERTIES
15  (
16     "timeout" = "3600"
17  );
```

2、Broker Load

Broker Load 针对ck导出的数据存到HDFS里

6、数据验证

可以跑一下count(*) 或者常用分析SQL对比CK验证下结果