

flink sql 指南

说明

flink sql 相关知识整理。

表声明语句

flink 可以通过ddl语句声明一个表。表的声明在flink中分为2个部分：connector和format。

connector复杂读写数据（对接外部存储系统），format 负责解析数据。

通过如下方式可以声明一个表：

```
1 create table tablename(  
2     field1 field_type  
3 ) with (  
4     'key' = 'value'  
5 )
```

数据类型

```
1  STRING  
2  BOOLEAN  
3  BYTES    BINARY and VARBINARY are not supported yet.  
4  DECIMAL      Supports fixed precision and scale.  
5  TINYINT  
6  SMALLINT  
7  INTEGER  
8  BIGINT  
9  FLOAT  
10 DOUBLE  
11 DATE  
12 TIME      Supports only a precision of 0.  
13 TIMESTAMP      Supports only a precision of 3.  
14 TIMESTAMP WITH LOCAL TIME ZONE Supports only a precision of 3.  
15 INTERVAL      Supports only interval of MONTH and SECOND(3).  
16 ARRAY  
17 MULTiset  
18 MAP  
19
```

string 和 varchar 等价。

复制类型：

```
1  -- 定义数组类型
2  arr array<int>
3
4  -- 定义map类型
5  `map` map<string, string>
6
7  -- 定义嵌套类型
8  obj row<
9      id string,
10     name string
11     address row<
12         city string,
13         number int
14     >
15 >
```

连接器

jdbc 连接器

建表sql语句如下：

```
1  create table dim (
2      dim varchar ,
3      channel_eight_role_code varchar ,
4      channel_source_code varchar,
5      CHANNEL_INFO_ID varchar
6  ) with(
7      -- 声明连接器类型。flink会通过spi找到连接器，并且进行参数匹配
8      'connector.type' = 'jdbc',
9
10     -- jdbc的url
11     'connector.url' = 'jdbc:mysql://10.25.76.173:3310/ogg_syncer?useUnicode=true&characterEncoding=UTF-8&useSSL=false',
12
13     -- 表名称
14     'connector.table' = 'epcis_epcisbase_channel_info',
15
16     -- 驱动类型
```

```

18     'connector.driver' = 'com.mysql.jdbc.Driver',
19
20     -- 用过名和密码
21     'connector.username' = 'root',
22     'connector.password' = 'root',
23
24     -- jdbc作为维表的时候，缓存时间。cache默认未开启。
25     'connector.lookup.cache.ttl' = '60s',
26
27     -- jdbc作为维表的时候，缓存的最大行数。cache默认未开启。
28     'connector.lookup.cache.max-rows' = '100000',
29
30     -- jdbc作为维表的时候，如果查询失败，最大查询次数
31     'connector.lookup.max-retries' = '3',
32
33     -- jdbc写入缓存的最大行数。默认值5000
34     'connector.write.flush.max-rows' = '5000',
35
36     -- jdbc 写入缓存flush时间间隔。默认为0，立即写入
37     'connector.write.flush.interval' = '2s',
38
39     -- 写入失败，最大重试次数
40     'connector.write.max-retries' = '3'
    );

```

说明：

1. lookup必须同时设置ttl和max-rows两个参数。
2. jdbc sink 返回的是UpsertStreamSink。

jdbc连接器只支持append/upsert模式，在有些情况下可能无法使用。paic-jdbc基于官方的jdbc开发，

注意添加了如下功能：

1. 添加参数connector.read.sub-query，设置查询jdbc的子查询，如果维度需要自动更新并且需要进行去重
等操作，去重的sql语句放在connector.read.sub-query中。
2. 支持retract模式，设置update-mode为retract，将返回retract sink。

hbase 连接器

实例：

```

1 | CREATE TABLE MyUserTable (
2 |     hbase_rowkey_name rowkey_type,

```

```

3   hbase_column_family_name1 ROW<...>,
4   hbase_column_family_name2 ROW<...>
5 ) WITH (
6   'connector.type' = 'hbase',
7   'connector.version' = '1.4.3',
8
9   -- hbase 表名称
10  'connector.table-name' = 'hbase_table_name', -- required: hbase table
11  name
12
13  -- zk地址
14  'connector.zookeeper.quorum' = 'localhost:2181',
15  -- zk 根节点
16  'connector.zookeeper.znode.parent' = '/base',
17
18  -- buffer 缓存大小。默认 2mb。
19  'connector.write.buffer-flush.max-size' = '10mb',
20
21  -- 缓冲的最大记录数，无默认值
22  'connector.write.buffer-flush.max-rows' = '1000',
23
24  -- flush 时间间隔，默认为0，表示理解刷新到hbase，无缓冲。
25
26  'connector.write.buffer-flush.interval' = '2s'
27 )

```

说明在定义hbase的schema中，唯一的非row类型的字段，会被当做rowkey处理。
一个完整的实例如下：

```

1   -- cf 为列族名称， row里面的是列名称
2   create table hbase_sink(
3     rowkey varchar,
4     cf row<a1 string, a2 string>
5   ) with(
6     'connector.type' = 'hbase',
7     'connector.version' = '1.4.3',
8     'connector.zookeeper.quorum' = '10.25.76.175:2181,10.25.76.173:2181',
9     'connector.zookeeper.znode.parent' = '/hbase',
10    'connector.table-name' = 'xuen',
11    'connector.write.buffer-flush.interval' = '1s'
12  );
13
14  -- row() 用于创建一个row，按位置对应。
15  insert into hbase_sink select rowkey, row(a1, a2) from kfk_source;

```

paic-hbase 基于官方的hbase开发，提供如下功能：

1. 支持retract模式，设置update-mode为retract，将返回retract sink。
2. 添加参数connector.write.null，表示是否写入null值。如果为false，值为null的列将不写入hbase，默认为true

kafka连接器

```
1
2
3 CREATE TABLE MyUserTable (
4     ...
5 ) WITH (
6     'connector.type' = 'kafka',
7     'connector.version' = 'universal',
8
9     -- topic名称
10    'connector.topic' = 'topic_name',
11
12    -- 固定值。必须有
13    'update-mode' = 'append',
14
15    -- 设置kafka集群地址
16    'connector.properties.0.key' = 'bootstrap.servers',
17    'connector.properties.0.value' = 'localhost:9092',
18
19    -- 设置group id
20    'connector.properties.1.key' = 'group.id',
21    'connector.properties.1.value' = 'testGroup',
22
23    -- 设置启动模式。如果指定了checkpoint，将从checkpoint读取offset
24    -- earliest-offset 最早的offset
25    -- latest-offset 最近的offset
26    -- group-offsets group 的offset
27    -- specific-offsets 指定的offset
28    'connector.startup-mode' = 'earliest-offset',
29
30    -- 指定的offset。
31    'connector.specific-offsets.0.partition' = '0',
32    'connector.specific-offsets.0.offset' = '42',
33    'connector.specific-offsets.1.partition' = '1',
34    'connector.specific-offsets.1.offset' = '300',
35
36    -- sink分区器。默认是
37    'connector.sink-partitioner' = '...',
38
39    -- 指定分区器的类。
```

```
40 | 'connector.sink-partitioner-class' = 'org.mycompany.MyPartitioner'
41 | )
```

paic-kafka基于官方的kafka开发，添加的功能如下：

1. 返回的sink为retract sink。

es 连接器

```
1
2
3 CREATE TABLE MyUserTable (
4     ...
5 ) WITH (
6     'connector.type' = 'elasticsearch',
7     'connector.version' = '6',
8
9     -- 定义host/端口/协议类型
10    'connector.hosts.0.hostname' = 'host_name',
11    'connector.hosts.0.port' = '9092',
12    'connector.hosts.0.protocol' = 'http',
13
14    -- 索引名称
15    'connector.index' = 'MyUsers',
16
17    -- es doc-type
18    'connector.document-type' = 'user',
19
20    -- update mode。append 将只有insert操作。
21    'update-mode' = 'append',
22
23    -- 生成文档id的连接符
24    'connector.key-delimiter' = '$',
25
26    -- key null值占位符，默认null
27    'connector.key-null-literal' = 'n/a',
28
29    -- 错误处理handler
30    'connector.failure-handler' = '...',
31
32    -- optional: configure how to buffer elements before sending them in bu
33    lk to the cluster for efficiency
34    'connector.flush-on-checkpoint' = 'true',    -- optional: disables flushi
35    ng on checkpoint (see notes below!)
36
37
38    -- ("true" by default)
39    -- 每个each bulk request的最大操作数量
```

```

39     'connector.bulk-flush.max-actions' = '42',
40
41     -- 缓冲区大小。only MB granularity is supported
42     'connector.bulk-flush.max-size' = '42 mb',
43
44     -- flush 频率
45     'connector.bulk-flush.interval' = '60000',
46
47     -- bulk 重试方式
48     -- optional: backoff strategy ("disabled" by default)
49     -- valid strategies are "disabled", "constant",
50     -- or "exponential"
51     'connector.bulk-flush.back-off.type' = '...',
52
53     -- 最大重试次数
54     'connector.bulk-flush.back-off.max-retries' = '3',
55
56     -- 重试间隔时间
57     'connector.bulk-flush.back-off.delay' = '30000',
58
59     -- optional: connection properties to be used during REST communicatio
60 n to Elasticsearch
61     -- optional: maximum timeout (in milliseconds)
62     -- between retries
63     'connector.connection-max-retry-timeout' = '3',
64
65     -- optional: prefix string to be added to every
66     -- REST communication
67     'connector.connection-path-prefix' = '/v1'
68 )

```

update-mode=append, es将使用es自动生成的文档ID, 也就是只有insert操作。
 update-mode=upsert, 将使用group-by的字段值作为文档ID进行put操作。

format

目前官方提供的连接器中, 只有kafka是需要format的。
 这里介绍json format。

```

1
2
3 CREATE TABLE MyUserTable (
4     ...
5 ) WITH (
6     'format.type' = 'json',

```

```

7
8   -- optional: flag whether to fail if a field is missing or not, false by
9   default
10  'format.fail-on-missing-field' = 'true',
11
12  -- required: define the schema either by using a type string which pars
13  es numbers to corresponding types
14  'format.fields.0.name' = 'lon',
15  'format.fields.0.type' = 'FLOAT',
16  'format.fields.1.name' = 'rideTime',
17  'format.fields.1.type' = 'TIMESTAMP',
18
19  -- or by using a JSON schema which parses to DECIMAL and TIMESTAMP
20  'format.json-schema' =
21    '{
22      "type": "object",
23      "properties": {
24        "lon": {
25          "type": "number"
26        },
27        "rideTime": {
28          "type": "string",
29          "format": "date-time"
30        }
31      }
32    }',
33
34  -- use the table's schema
35  'format.derive-schema' = 'true'
36 )

```

通常使用'format.derive-schema' = 'true'，不在with中单独定义schema。

flink-json的时间 attr只支持utc类型的timestamp，这可不太好用。平台开发了新的json格式，使用如下：

```

1  'format.type' = 'text',
2  'format.udf' = 'com.paic.bentley.flink.sql.format.udf.JSONForNestedUdf'
3  ,
4  'format.derive-schema' = 'true',

```

增加功能如下：

1. json解析失败不会导致任务结束，会返回一个null的row
2. 嵌套的json 定义为string，会将这个嵌套对象进行序列化，返回string。方便使用paic_explode_map进行展开，以规避udf

无法处理嵌套类型的问题。

3. 支持多种类型的time attr

数字类型的毫秒时间戳

utc时间

cst时间

yyyy-[m]m-[d]d hh:mm:ss[.f...] 格式的时间

4. flink 无法定义array类型，请定义为map<int,string>来规避此问题。

flink sink 更新模式

append 模式

没有聚合操作或者有状态的操作，可以使用append模式。历史消息不会更新，只有追加的操作

upsert 模式

需要有group by操作或者append only 为false才可使用。group by 的字段值就是flink更新状态的unique key。

upsert 模式的消息是一个tuple: (Boolean, Row)。

append 消息: (true, Row)

delete 消息: (false, Row), 表示删除消息

upsert 消息: (true, Row), 表现已经存在的唯一key的状态发生了变更。

retract

retract 是通用的类型，任务情况下都可以使用。retract也会group by 的字段值就是flink更新状态的unique key。

retract 模式的消息是一个tuple: (Boolean, Row)

append 消息: (true, Row)

delete 消息: (false, Row), 表示删除消息

update 消息: (false, Row): 表示删除这个key, row的值是之前的状态; (false, Row)表示插入这个key, row的值是现在状态

3中模式的使用

一个查询语句需要insert到sink中的时候，flink 会进行更新模式的教程。主要是判断查询sql 是否有如下状态：

1. appendOnly的是。如果sql查询不包含有状态操作，没有group by，appendOnly=true。

2. 是否有unique key, 通常就是group by的字段。

如果appendOnly为true: 可以使用append, upsert模式

如果有unique key, 可以使用upsert模式。

retract模式, 无显著条件。

需要注意的使用 upsert模式中, group 的字段必须出现在select中, 否则会报错。如例子:

```
1 | -- word 必须出现在select中。  
2 | select word, count(*) from t group by word
```

flink sql 时间属性 (time attr)

在flink 中使用group window 必须定义时间字段。

目前时间字段值能定义在table定义中, 查询语句是无法定义时间属性的

定义process time

schema.位置, 表示引用字段定义的schema。位置从上到下, 从0开始。

schema.0, 表示第一个字段。

```
1 | create table t(  
2 |     ts timestamp  
3 | ) with(  
4 |     -- 表示此字段为进程处理时间。flink会自动填充值。  
5 |     'schema.0.proctime' = 'true'  
6 | )
```

定义 event time

```
1 | create table t1(  
2 |     ts timestamp  
3 | ) with (  
4 |  
5 |     -- 声明字段来源  
6 |     'schema.2.from' = 'ts',  
7 |  
8 |     -- 声明时间来源于字段  
9 |     'schema.2.rowtime.timestamps.type' = 'from-field',  
10 |  
11 |     -- 字段名称  
12 |     'schema.2.rowtime.timestamps.from' = 'ts',  
13 |
```

```

14  -- 定义watermark 类型, periodic-bounded表示周期行生成边界
15  'schema.2.rowtime.watermarks.type' = 'periodic-bounded',
16
17  -- watermark 最大延迟时间。
18  'schema.2.rowtime.watermarks.delay' = '60000'
19 );

```

group window 函数

group window操作和返回的字段类型，都必须是timestamp类型。

窗口函数，下面这些函数必须出现在group by中，表示按窗口聚合：

函数	
TUMBLE(time_attr, interval)	
HOP(time_attr, interval, interval)	
SESSION(time_attr, interval)	

获取窗口的开始时间（包含）：

```

1  TUMBLE_START(time_attr, interval)
2  HOP_START(time_attr, interval, interval)
3  SESSION_START(time_attr, interval)

```

获取窗口结束时间（不包含）：

```

1  TUMBLE_END(time_attr, interval)
2  HOP_END(time_attr, interval, interval)
3  SESSION_END(time_attr, interval)

```

在group by 中使用了窗口函数，select 必须出现一个START/END函数，这个group by的字段

必须出现在select中是同一个道理，如下实例：

```

1  insert into yp_audit_stats_um
2  select
3      TUMBLE_START(execTime, INTERVAL '1' day) as dt,
4      count(distinct userUM) as um_num,
5      'day' as dim
6  from um_log
7  group by TUMBLE(execTime, INTERVAL '1' day);

```

级联窗口

Rowtime列在经过窗口操作后，其Event Time属性将丢失。您可以使用辅助函数TUMBLE_ROWTIME、HOP_ROWTIME或SESSION_ROWTIME获取窗口中的Rowtime列的最大值max(rowtime)作为时、间窗口的Rowtime，其类型是具有Rowtime属性的TIMESTAMP，取值为 window_end - 1

```
1  SELECT
2      -- 使用TUMBLE_ROWTIME作为二级Window的聚合时间
3      TUMBLE_ROWTIME(ts, INTERVAL '1' MINUTE) as rowtime,
4      username,
5      COUNT(click_url) as cnt
6  FROM user_clicks
7  GROUP BY TUMBLE(ts, INTERVAL '1' MINUTE), username;
8
9  -- 时间窗口二次聚合。
10 INSERT INTO tumble_output
11 SELECT
12     TUMBLE_START(rowtime, INTERVAL '1' HOUR),
13     TUMBLE_END(rowtime, INTERVAL '1' HOUR),
14     username,
15     SUM(cnt)
16 FROM one_minute_window_output
17 GROUP BY TUMBLE(rowtime, INTERVAL '1' HOUR), username
```

TUMBLE_PROCTIME(time_attr, interval)

HOP_PROCTIME(time_attr, interval, interval)

SESSION_PROCTIME(time_attr, interval) 这3个函数是针对的proctime的。和上面的功能一样。

watermark

schema.#.rowtime.watermarks.type 定义水印类型，有如下3中：

1. periodic-ascending：事件的最大时间戳 -1 ,基本相当于无延迟
2. periodic-bounded ， 需要设置最大延迟时间delay，水印的大小为最大时间戳 - delay
3. from-source 保留源中的水印。

默认情况下， watermark到达窗口结束后，完成聚合操作，只会执行1次，相关于如果定义了1天的时间窗口， 1天之后才能看下结果。这个时候如果需要实时看到结果，需要定义触发器

watermark到达窗口结束前的发射策略是否开启: `table.exec.emit.early-fire.enabled`, 默认false

`table.exec.emit.early-fire.delay`, 窗口结束前的发射间隔, 单位毫秒。=0, 无间隔, >0 间隔时间, <0 非法值。无默认值

watermark到达窗口结束后的发射策略是否开启 `table.exec.emit.late-fire.enabled`, 默认false

`table.exec.emit.late-fire.delay`, 设置间隔时间

设置实例:

```
1 | -- set 是平台的功能, 非flink本身的。
2 | set table.exec.emit.early-fire.enabled = true;
3 | set table.exec.emit.early-fire.delay = 1.s;
```

时间单位

```
1 | private[this] val timeUnitLabels = List(
2 |     DAYS          -> "d day",
3 |     HOURS         -> "h hour",
4 |     MINUTES       -> "min minute",
5 |     SECONDS       -> "s sec second",
6 |     MILLISECONDS -> "ms milli millisecond",
7 |     MICROSECONDS -> "µs micro microsecond",
8 |     NANOSECONDS  -> "ns nano nanosecond"
9 | )
```

注意事项

time attr 使用 group window的时候不能使用函数。

`TUMBLE_START(fun(timestamp), INTERVAL '1' hour);`

timestamp 失去了时间属性, 不能使用TUMBLE

`fun(TUMBLE_START(timestamp , INTERVAL '1' hour))`

fun后, 失去了时间属性, 和group by字段不能匹配, 会判定为主键不完整

flink 维表自动更新

flink 目前使用了look up的方式来自动更新维表, 目前只是blink planer支持。

维表的字段更新目前只有jdbc, 使用维表的自动更新, 需要指定

`connector.lookup.cache.ttl`, 和`connector.lookup.cache.max-rows`2个参数。

hbase也支持维表自动更新, 但是没有使用缓存, 每次都会查询hbase。

如果使用维表自动更新

流水表需要定义一个proctime字段：

```
1 | create table t(  
2 |     ts timestamp  
3 | ) with(  
4 |     -- 表示此字段为进程处理时间。flink会自动填充值。  
5 |     'schema.0.proctime' = 'true'  
6 | )
```

维表正常定义，定义好lookup相关参数，不需要定义时间参数，在join维表使用如下语法：

```
1 | select a.id from a  
2 | left join diw FOR SYSTEM_TIME AS OF a.ts b on agr.id = b.id
```

join的表名称后面跟：FOR SYSTEM_TIME AS OF a.ts 加表别名称。

维表自动更新的原理

jdbc的JDBCLookupFunction 就是继承了TableFunction。

JDBCTableSource会继承LookupableTableSource，source的getLookupFunction会返回JDBCLookupFunction

blink planer 碰见FOR SYSTEM_TIME AS OF a.ts，会调用getLookupFunction，这是一个表函数，会返回多行。

JDBCLookupFunction 在创建的时候，会创建一个guava cache：private transient Cache<Row, List> cache;

过期时间为ttl设置的值，最大大小为max-rows设置的值。这个cache的key其实是join的全部字段的值，value的值是对应join字段的值在jdbc中的全部记录。

在eval方法中，传递进来join字段的值，判断cache中是否存在这个记录，如果存在返回。

如果不存在，在jdbc中查找指定join字段值的记录（不会查询全部），保存到缓存中，返回。

说明

如果jdbc维表使用了distinct等有状态操作，是无法使用 FOR SYSTEM_TIME AS 语法的，解析会报错。

在paic-jdbc的封装中，可以用connector.read.sub-query参数，传递一个字查询，在这个

子查询中使用

distinct语法，在实现中这个配置做为子查询：

```
1 | return "SELECT " + selectExpressions + " FROM (" + subQuery + ") " +  
2 |       quoteIdentifier(tableName) + (conditionFields.length > 0 ? " WHERE  
   | E " + fieldExpressions : "");
```

使用方式如下：

```
1 | create table dw_dim_department_source(  
2 | ) with(  
3 |   'connector.read.sub-query' =  
4 |     'select  
5 |       distinct department_code,  
6 |       sec_department_code,  
7 |       third_department_code,  
8 |       fourth_department_code  
9 |     from  
10 |       dw_dim_department'  
11 |  
12 | );
```

flink 1.9.1 问题和bug

bug

1. 无法定义 array类型的数据。sql解析失败
2. cast(a as string), 会报错，只能使用varchar
3. decimal类型在sink和source中，会报类型不匹配的错误。

需要注意的的地方

1. flink udf 无法支持嵌套数据类型（row类型）
2. flink sql 区分大小写
3. flink 不会进行自动类型转换。‘1’ * 0.1，会报错，请使用cast强制类型转换。
4. row()生产一个嵌套类型，只支持写字段名称。写函数，或者加库名称都是不行的
5. map类型的访问，目前只支持：`map['filed']` 这种方式。
6. 数组的下标是从1开始，：`riskGroupInfoList[1]`
7. insert into 不支持部分字段。

8. 如果自己定义了factroy，flink lib 目录下也有factory，java 包需要放到 flink/lib 目录下，否则无法加载自定义的factory。
9. flink kafka 0.11以上版本和其他版本存在冲突，只能引入一个。

自定义sink/source/factory