# 基于加权PageRank的Github项目活跃度分析

## 1.数据预处理

### 1.1基于开发者的 GitHub 行为数据计算其对项目的贡献度

$A_d = \sum w_i c_i$ ( 其中的 $A_d$ 为开发者对项目的贡献度，而 $c_i$ 为行为事件由该开发者触发的发生次数，$w_i$ 为该行为事件的加权比例)

| 行为事件 | 赋分 |
|---|---|
| Issue 评论 | 1 |
| Issue | 2 |
| PR | 3 |
| PR Review | 4 |

```sql
SELECT   SUM(t.score) AS score, actor_id, repo_id
FROM
  (SELECT CASE WHEN type = 'IssueCommentEvent' THEN 1
   WHEN type = 'IssuesEvent' THEN 2 when  type = 'PullRequestEvent' THEN 3
   WHEN type = 'PullRequestReviewCommentEvent' THEN 4
   ELSE 0 END AS score ,actor_id,repo_id
   FROM ods_github_log
   WHERE pt="20211201" ) t
GROUP BY actor_id,repo_id
```
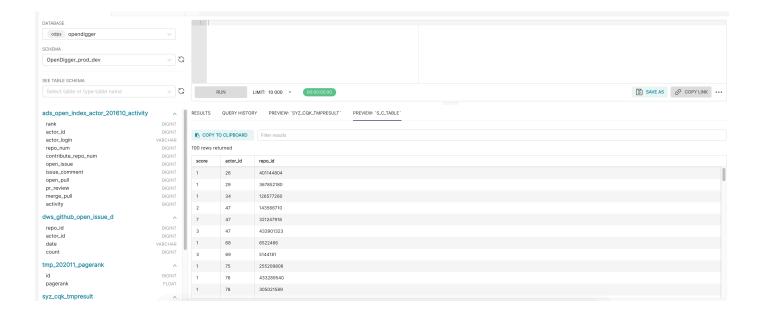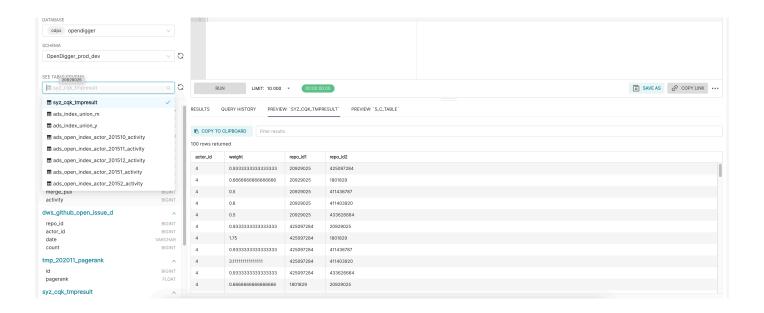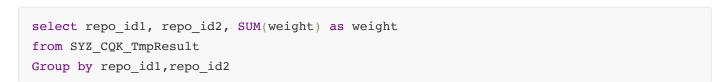
将统计好的数据存入中间表s_c_table中：

## 2.构建开源协作网络

**筛选数据：**

```sql
SELECT *
FROM s_c_table
WHERE actor_id not in
(SELECT actor_id FROM s_c_table GROUP BY actor_id HAVING COUNT(*) = 1)
```
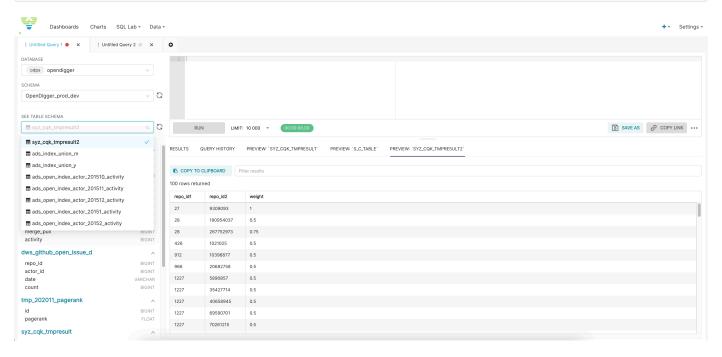
**计算项目与项目之间的协作关联度：**

```sql
select s1.actor_id, s1.score*s2.score/(s1.score+s2.score) as weight, s1.repo_id  AS repo_id1,s2.repo_id AS repo_id2
FROM syz_cqk_tmptable s1 JOIN syz_cqk_tmptable s2
ON s1.actor_id = s2.actor_id AND s1.repo_id != s2.repo_id
```

**汇总：**

```sql
select repo_id1, repo_id2, SUM(weight) as weight
from SYZ_CQK_TmpResult
Group by repo_id1,repo_id2
```



边的条数：648亿 顶点个数: 687万

```
ACCESS_ID = 'LTAI5tHDArybZRnXaPS3pdkJ'
SECRET_ACCESS_KEY = 'lXXLM2U1dB3ExgKquGUwdGb88WHqkN'
ODPS_PROJECT = 'OpenDigger_prod_dev'
ODPS_ENDPOINT = 'http://service.cn-shanghai.maxcompute.aliyun.com/api'

o = ODPS(ACCESS_ID, SECRET_ACCESS_KEY,
         project=ODPS_PROJECT, endpoint=ODPS_ENDPOINT)
options.tunnel.limit_instance_tunnel = False
# options.read_timeout = 3600000

result = o.execute_sql('SELECT * FROM SYZ_CQK_TmpResult2',hints={'odps.sql.allow.fullscan': 'true'})

with result.open_reader() as reader:
    print(reader.count/2)
    # for record in reader:
    #     # o.write_table(table,record)
    #     print(record)

# #读取SQL执行结果。
# with result.open_reader() as reader:
# github_log = DataFrame(o.get_table('ods_github_log'))
# print(github_log.dtypes)
```

```
/usr/local/bin/python3.9 "/Users/sunyinzheng/2021课程 /社会计算/test.py"
64842429194.0

Process finished with exit code 0
```



```
options.tunnel.limit_instance_tunnel = False
# options.read_timeout = 3600000

result = o.execute_sql('SELECT distinct repo_id1 FROM SYZ_CQK_TmpResult2',hints={'odps.sql.allow.fullscan': 'true'})

with result.open_reader() as reader:
    print(reader.count)
    # for record in reader:
    #     # o.write_table(table,record)
    #     print(record)

# #读取SQL执行结果。
# with result.open_reader() as reader:
# github_log = DataFrame(o.get_table('ods_github_log'))
# print(github_log.dtypes)
```

```
/usr/local/bin/python3.9 "/Users/sunyinzheng/2021课程 /社会计算/test.py"
6871158
```

# 3.后续计划

通过pagerank算法来计算项目在开源协作网络中的影响力。