

# 任务1

## 连接数据库

```
%load_ext sql
```

```
1 import pymysql
2 pymysql.install_as_MySQLdb()
3 %sql mysql://stu2000013058:stu2000013058@162.105.146.37:43306
```

```
%sql use stu2000013058;
```

- 创建testIndex

```
1 %%sql
2 DROP TABLE IF EXISTS testIndex;
3 SET @i = 0;
4 CREATE TABLE testIndex SELECT (@i :=@i + 1) AS id, userId AS A, movieId AS B,
  tag AS C FROM dataset.tags;
```

## 针对A列的分组和自连接操作，观察A列上建立索引前后的性能差异

- 建立索引前：分组操作

```
1 %%time
2 %%sql
3
4 SELECT min(B) FROM testIndex GROUP BY A;
```

- 建立索引前：自连接操作

```
1 %%time
2 %%sql
3
4 SELECT * FROM testIndex AS tmp1 JOIN testIndex AS tmp2 ON tmp1.A = tmp2.A;
```

- 建立索引

```
1 %%sql
2
3 CREATE INDEX my_A ON testIndex(A);
```

- 建立索引后：分组操作

```
1 %%time
2 %%sql
3
4 SELECT min(B) FROM testIndex GROUP BY A;
```

- 建立索引后：自连接操作

```
1 %%time
2 %%sql
3
4 SELECT * FROM testIndex AS tmp1 JOIN testIndex AS tmp2 ON tmp1.A = tmp2.A;
```

- 我们发现建立索引之后效率反而变低了，我们推测这是因为
  - 确实使用了索引进行查询
  - 索引的开销超过了顺序读的开销

## 针对select B where A类型的查询，观察基于(A, B)的组合索引相对于A上的单列索引的性能提升

- 使用A上的单列索引进行查询

```
1 %%time
2 %%sql
3
4 SELECT B FROM testIndex WHERE A = 424;
```

- 删除A上的索引，建立(A, B)上的组合索引

```
1 %%sql
2
3 DROP INDEX my_A ON testIndex;
4 CREATE INDEX my_A ON testIndex(A, B);
```

- 使用基于(A, B)的组合索引进行查询

```
1 %%time
2 %%sql
3
4 SELECT B FROM testIndex WHERE A = 424;
```

- 在这个查询中，使用组合索引优于单列索引

## 对字符串的子字符串建立索引，观察函数索引的作用

- 建立索引前的三次查询(查询前3个字母、前两个字母和查询全部)

```
1 %%time
2 %%sql
3
4 SELECT * FROM testIndex WHERE SUBSTRING(C, 1, 3) = "pow";
```

```
1 %%time
2 %%sql
3
4 SELECT * FROM testIndex WHERE SUBSTRING(C, 1, 2) = "po";
```

```
1 %%time
2 %%sql
3
4 SELECT * FROM testIndex WHERE C = "powerful ending";
```

- 建立对前三个字母的索引

```
1 %%sql
2
3 CREATE INDEX my_C ON testIndex(C(3));
```

- 重新进行查询

```
1 %%time
2 %%sql
3
4 SELECT * FROM testIndex WHERE SUBSTRING(C, 1, 3) = "pow";
```

```
1 %%time
2 %%sql
3
4 SELECT * FROM testIndex WHERE SUBSTRING(C, 1, 2) = "po";
```

```
1 %%time
2 %%sql
3
4 SELECT * FROM testIndex WHERE C = "powerful ending";
```

- 删除索引

```
1 %%sql
2
3 DROP INDEX my_C ON testIndex;
```

- 可以看到，前3个字母和总字符串的查询的时间均发生了改变，意味着会从索引处进行查询，而前2个字母的查询时间没有发生太大变化，我们推测此时没有使用索引

## 任务2

- 连接数据库

```
1 %load_ext sql
2
3 import pymysql
4 pymysql.install_as_MySQLdb()
5 %sql mysql://stu2100013111:stu2100013111@162.105.146.37:43306
6
7 %sql use stu2100013111;
```

基于集合：

首先，使用C1，将会话的开始时间和结束时间以及对应的类型从 dbo.Sessions 表中选择出来。（其中+1表示开始，-1表示结束）

其次，在C2中，基于C1的结果，计算每个时间点上的并发会话数量，通过对类型进行累积求和并按照应用程序和时间戳进行排序。

最后，在C2的结果上进行聚合操作，按照应用程序分组，并选择每个应用程序的最大并发数量。

```
1 %%time
2 %%sql
3
4 WITH TimePoints AS
5 (
6     SELECT app, starttime AS ts FROM dbo.Sessions
7 ),
8 Counts AS
9 (
10    SELECT app, ts,
11           (SELECT COUNT(*)
12            FROM dbo.Sessions AS S
13            WHERE P.app = S.app
14                  AND P.ts >= S.starttime
15                  AND P.ts < S.endtime) AS concurrent
16    FROM TimePoints AS P
17 )
18 SELECT app, MAX(concurrent) AS mx
19 FROM Counts
20 GROUP BY app;
```

基于游标：

通过游标来循环遍历游标中的数据行，跟踪并发的数量，并进行更新，最终将结果取出。

```
1  %%time
2  %%sql
3
4  DECLARE
5      @app AS varchar(10),
6      @prevapp AS varchar (10),
7      @ts AS datetime,
8      @type AS int,
9      @concurrent AS int,
10     @mx AS int;
11
12  DECLARE @AppsMx TABLE
13  (
14      app varchar (10) NOT NULL PRIMARY KEY,
15      mx int NOT NULL
16  );
17
18  DECLARE sessions_cur CURSOR FAST_FORWARD FOR
19      SELECT app, starttime AS ts, +1 AS type
20      FROM dbo.Sessions
21
22      UNION ALL
23
24      SELECT app, endtime, -1
25      FROM dbo.Sessions
26
27      ORDER BY app, ts, type;
28
29  OPEN sessions_cur;
30
31  FETCH NEXT FROM sessions_cur
32      INTO @app, @ts, @type;
33
34  SET @prevapp = @app;
35  SET @concurrent = 0;
36  SET @mx = 0;
37
38  WHILE @@FETCH_STATUS = 0
39  BEGIN
40      IF @app <> @prevapp
41      BEGIN
42          INSERT INTO @AppsMx VALUES(@prevapp, @mx);
43          SET @concurrent = 0;
44          SET @mx = 0;
45          SET @prevapp = @app;
46      END
47
48      SET @concurrent = @concurrent + @type;
49      IF @concurrent > @mx SET @mx = @concurrent;
50
```

```

51     FETCH NEXT FROM sessions_cur
52     INTO @app, @ts, @type;
53 END
54
55 IF @prevapp IS NOT NULL
56     INSERT INTO @AppsMx VALUES(@prevapp, @mx);
57
58 CLOSE sessions_cur;
59
60 DEALLOCATE sessions_cur;
61
62 SELECT * FROM @AppsMx;

```

```

1  %%time
2  %%sql
3
4  WITH C1 AS
5  (
6      SELECT app, starttime AS ts, +1 AS type
7      FROM dbo.Sessions
8
9      UNION ALL
10
11     SELECT app, endtime, -1
12     FROM dbo.Sessions
13 ),
14 C2 AS
15 (
16     SELECT *,
17         SUM(type) OVER(PARTITION BY app ORDER BY ts, type
18             ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW) AS cnt
19     FROM C1
20 )
21 SELECT app, MAX(cnt) AS mx
22 FROM C2
23 GROUP BY app;
24

```

基于窗口函数；

使用窗口函数来减少游标的额外开销，进行关于app的排序进行筛选，比游标进一步提高效率

## 结论

- 基于集合的方案

运行时间随着数据量增加，近似指数爆炸式增长，效率最差

- 基于游标的方案

运行时间成一个线性函数增长，效率一般

- 基于窗口函数的方案

运行时间平稳，增长极为缓慢，效率高

集合	游标	窗口
----	----	----

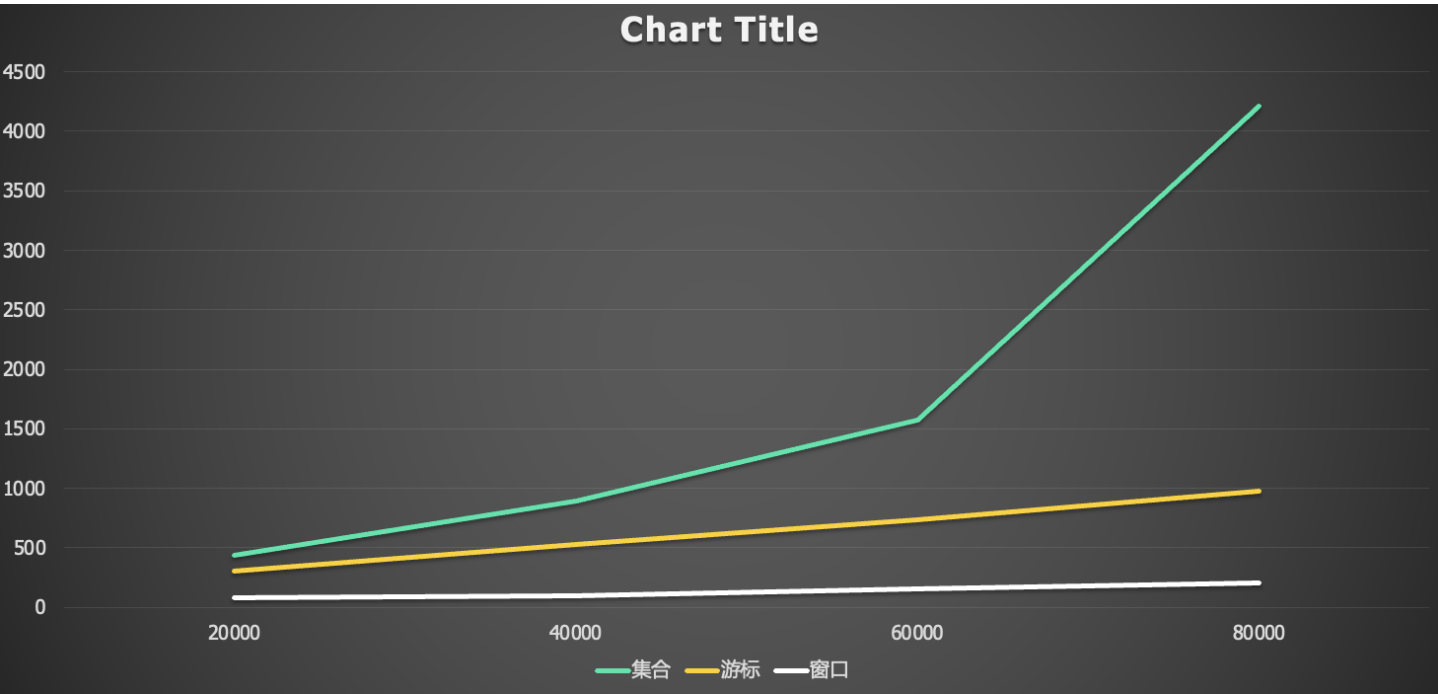
20000 435 301 77

40000 890 528 96

60000 1573 737 154

80000 4209 978 205

100000 7199 763 229



## 任务3

- 连接数据库

```
1 db = pymysql.connect(host='162.105.146.37', user='stu2000013058',
2   passwd='stu2000013058', port=43306, db = 'stu2000013058')
3 cursor = db.cursor()
4 %sql mysql://stu2000013058:stu2000013058@162.105.146.37:43306
5 %sql use stu2000013058;
```

- 创建skyline数据表，同时，建立一个分区的my\_skyline数据表

--

```

1 %%sql
2
3 DROP TABLE IF EXISTS skyline;
4 CREATE TABLE skyline(
5     id INT PRIMARY KEY,
6     x INT,
7     y INT
8 );
9
10 DROP TABLE IF EXISTS my_skyline;
11 CREATE TABLE my_skyline(
12     id INT,
13     x INT,
14     y INT
15 )
16 PARTITION BY RANGE(x)
17 (
18     PARTITION m0 VALUES LESS THAN(10),
19     PARTITION m1 VALUES LESS THAN(20),
20     PARTITION m2 VALUES LESS THAN(30),
21     PARTITION m3 VALUES LESS THAN(40),
22     PARTITION m4 VALUES LESS THAN(50),
23     PARTITION m5 VALUES LESS THAN(60),
24     PARTITION m6 VALUES LESS THAN(70),
25     PARTITION m7 VALUES LESS THAN(80),
26     PARTITION m8 VALUES LESS THAN(90),
27     PARTITION m9 VALUES LESS THAN(101)
28 );

```

- 生成随机数据

```

1 import random
2
3 cur_sql = "INSERT INTO skyline(id, x, y) VALUES (%d, %d, %d)"
4 my_cur_sql = "INSERT INTO my_skyline(id, x, y) VALUES (%d, %d, %d)"
5 n_samples = 1000
6 max_val = 100
7
8 for i in range(n_samples):
9     x = random.randint(0, max_val)
10    y = random.randint(0, max_val)
11    cursor.execute(cur_sql % (i, x, y))
12    cursor.execute(my_cur_sql % (i, x, y))
13
14 db.commit()

```

- 运行原始skyline查询，看看运行效率如何，同时简单观察一下原始集规模大小和相关性对结果集大小的影响

```

1 %%time
2 %%sql

```



```

3
4 SELECT *
5 FROM skyline h
6 WHERE NOT EXISTS (
7     SELECT *
8     FROM skyline h1
9     WHERE h1.x <= h.x
10    AND h1.y <= h.y
11    AND (h1.x < h.x OR h1.y < h.y)
12 );

```

- 1000条数据运行时间为19.4 ms, 10000条数据39.1 ms
- 我们对表中的x, y均建立索引

```

1 %%sql
2
3 CREATE INDEX my_x ON skyline(x);
4 CREATE INDEX my_y ON skyline(y);

```

- 然后再运行上述查询

```

1 %%time
2 %%sql
3
4 SELECT *
5 FROM skyline h
6 WHERE NOT EXISTS (
7     SELECT *
8     FROM skyline h1
9     WHERE h1.x <= h.x
10    AND h1.y <= h.y
11    AND (h1.x < h.x OR h1.y < h.y)
12 );

```

- 最后, 我们对分区表进行查询, 查看其性能

```

1 %%time
2 %%sql
3
4 SELECT *
5 FROM my_skyline h
6 WHERE NOT EXISTS (
7     SELECT *
8     FROM my_skyline h1
9     WHERE h1.x <= h.x
10    AND h1.y <= h.y
11    AND (h1.x < h.x OR h1.y < h.y)
12 );

```

- 在本数据集上，分区表执行的时间要慢于原表

