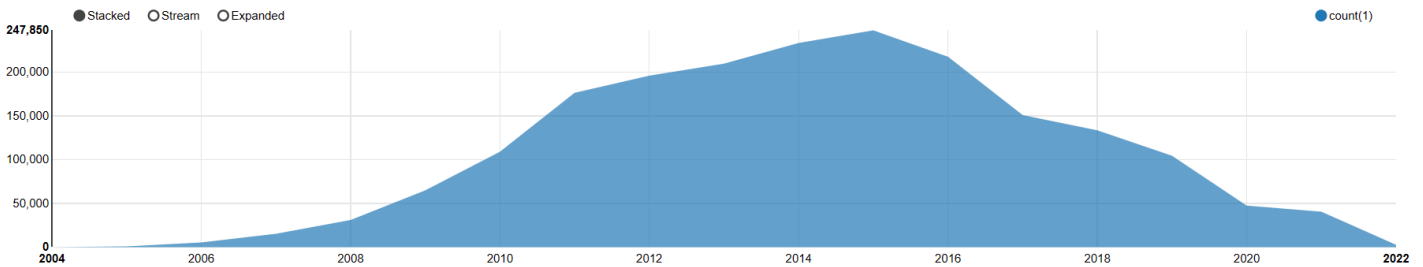
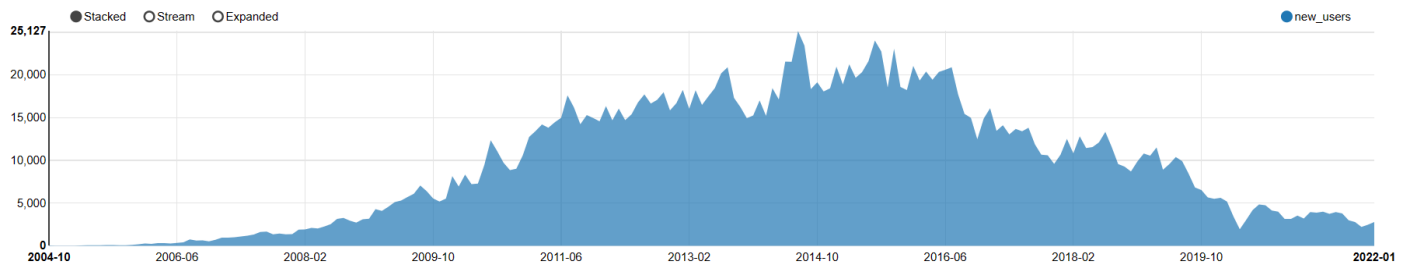


# 用户统计及分析模块Spark

```
1 %pyspark
2
3 # 编写者:凌晨
4 # 功能:分析每年加入的用户数量
5 # 时间:2024.1.3
6 from pyspark.sql import HiveContext
7
8 hc = HiveContext(sc)
9 result = hc.sql("SELECT year(to_date(user_yelping_since)) as year, COUNT(*)
    FROM users GROUP BY year(to_date(user_yelping_since)) ORDER BY year DESC")
10 z.show(result)
```



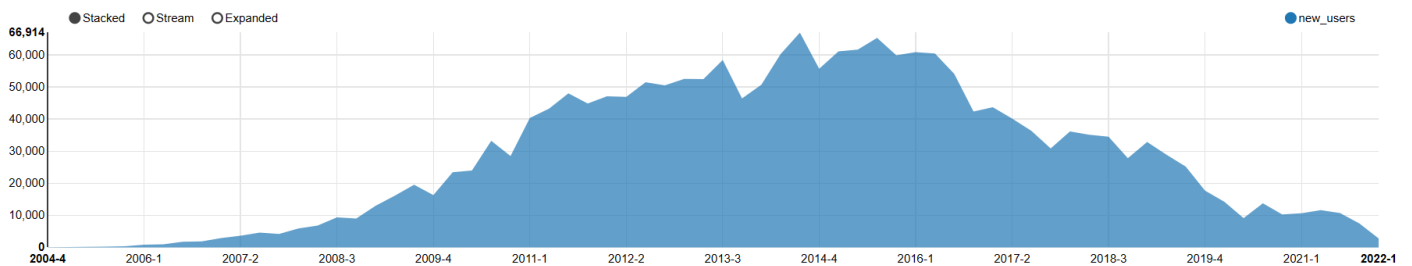
```
1 %pyspark
2
3 # 编写者:凌晨
4 # 功能:分析每月加入的用户数量
5 # 时间:2024.1.3
6 from pyspark.sql import HiveContext
7
8 hc = HiveContext(sc)
9 result = hc.sql("SELECT DATE_FORMAT(user_yelping_since, 'yyyy-MM') AS month,
    COUNT(*) AS new_users FROM users GROUP BY DATE_FORMAT(user_yelping_since,
    'yyyy-MM') ORDER BY month")
10 z.show(result)
```



```

1 pyspark
2
3 # 编写者:凌晨
4 # 功能:分析每季度加入的用户数量, 每三个月视为一个季度。如2004-4代表2004的第四个季度
5 # 时间:2024.1.3
6 from pyspark.sql import HiveContext
7
8 hc = HiveContext(sc)
9 result = hc.sql('''
10 SELECT
11     CONCAT(YEAR(user_yelping_since), '-', CEIL(MONTH(user_yelping_since) / 3))
12     AS quarter,
13     COUNT(*) AS new_users
14 FROM users
15 GROUP BY CONCAT(YEAR(user_yelping_since), '-', CEIL(MONTH(user_yelping_since)
16 / 3))
17 ORDER BY quarter
18 ''')
19 z.show(result)

```



```

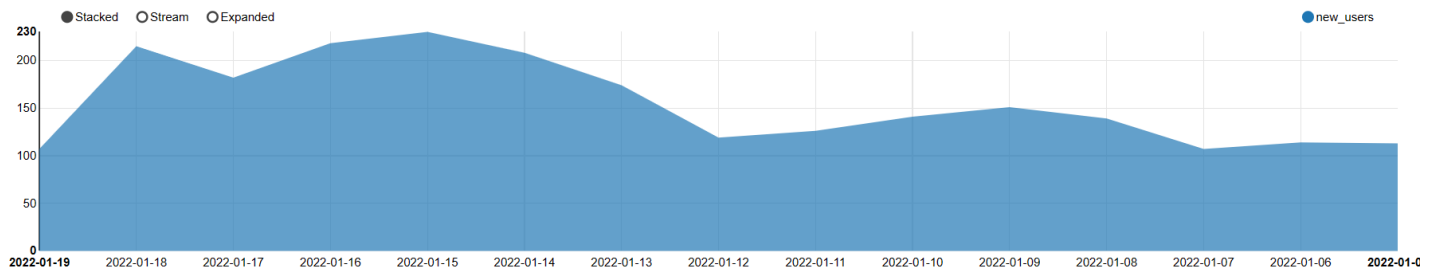
1 %pyspark
2
3 # 编写者:凌晨
4 # 功能:分析最近15天内的加入的用户数量
5 # 时间:2024.1.3

```

```

6
7 from pyspark.sql import HiveContext
8
9 hc = HiveContext(sc)
10 result = hc.sql("SELECT DATE(user_yelping_since) AS date, COUNT(*) AS
    new_users FROM users GROUP BY DATE(user_yelping_since) ORDER BY date DESC
    LIMIT 15")
11 z.show(result)

```

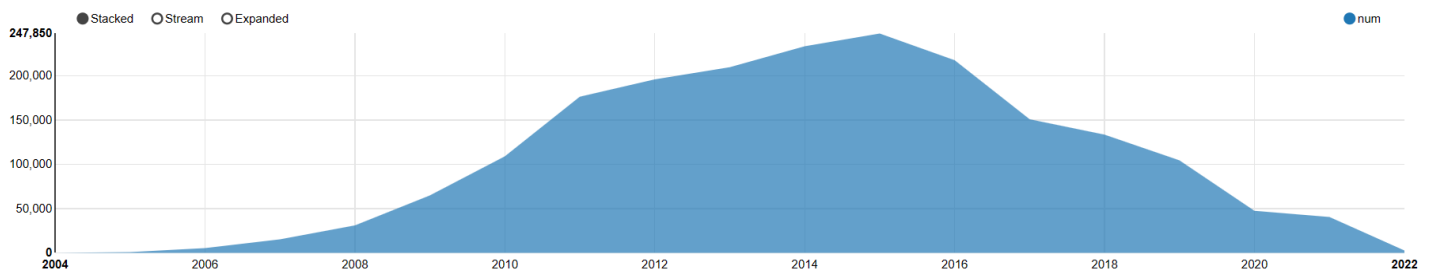


Task # can I get undisturbed anonymous at January 02 2024 3:27:04 PM (outdated)

```

1 %pyspark
2
3 # 编写者:凌晨
4 # 功能:分析每年加入的用户数量,按照用户数排序
5 # 时间:2024.1.3
6 from pyspark.sql import HiveContext
7
8 hc = HiveContext(sc)
9 result = hc.sql("SELECT year(to_date(user_yelping_since)) as year, COUNT(*) AS
    num FROM users GROUP BY year ORDER BY num DESC")
10 z.show(result)

```



```

1 %pyspark
2
3 # 编写者:凌晨

```

```

4 # 功能:找出最常见的前20的用户名
5 # 时间:2024.1.3
6
7 from pyspark.sql import HiveContext
8
9 hc = HiveContext(sc)
10 result = hc.sql("SELECT user_name, COUNT(user_name) as name_count FROM users
    GROUP BY user_name ORDER BY name_count DESC LIMIT 20")
11
12 z.show(result)

```

user_name	name_count
John	18719
Michael	16645
David	15967
Chris	14591
Mike	13390
Jennifer	13275
Sarah	11893
Jessica	11688

```

1 %pyspark
2
3 # 编写者:凌晨
4 # 功能:统计已注册的全部用户数量
5 # 时间:2024.1.3
6
7 from pyspark.sql import HiveContext
8
9 hc = HiveContext(sc)
10 result = hc.sql("SELECT COUNT(*) as num_of_user FROM users")
11
12 z.show(result)

```

num_of_user
1987897

```
1 %pyspark
2
3 # 编写者:钟鹏
4 # 功能:通过id查找name
5 # 时间:2024.1.4
6
7 from pyspark.sql import HiveContext
8 from pyspark.sql.functions import *
9
10 hcx = HiveContext(sc)
11
12 df = hcx.table('users')
13
14 user_id = 'Js4nQlGRjS1Bd5vsAQQfgA'
15 # 需要查找的 user_id
16
17 result = df.filter(col('user_id') == user_id).select('user_name')
18
19 z.show(result)
```

user\_name

Tom