Churn Rates with Codeflix

Learn SQL from scratch Patrick Fitzgerald 09.18.2018

Table of Contents

- 1. Getting Familiar with the company
 - How many months has the company been operating and how many months do we have enough information to calculate churn for?
 - What segment of users exist?
- 2. What is the overall churn trend since the company started?
- 3. Compare the churn rates between user segments
 - Which segment of users should the company focus on expanding?

1. Getting familiar with the company

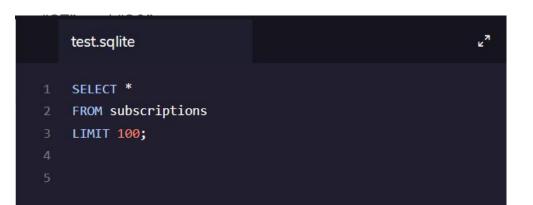
1.1 First we take a look at the data, how many different segments do we see?

For this I just did a simple "Select (*) all columns"

From the "subscriptions" table

And I limited my results to 100 rows

From this we can see there are two user segments



Query Results				
id	subscription_start	subscription_end	segment	
1	2016-12-01	2017-02-01	87	
2	2016-12-01	2017-01-24	87	
3	2016-12-01	2017-03-07	87	
4	2016-12-01	2017-02-12	87	
5	2016-12-01	2017-03-09	87	
6	2016-12-01	2017-01-19	87	
7	2016-12-01	2017-02-03	87	
8	2016-12-01	2017-03-02	87	
9	2016-12-01	2017-02-17	87	
10	2016-12-01	2017-01-01	87	
11	2016-12-01	2017-01-17	87	
12	2016-12-01	2017-02-07	87	
13	2016-12-01	Ø	30	
14	2016-12-01	2017-03-07	30	
15	2016-12-01	2017-02-22	30	
16	2016-12-01	Ø	30	
17	2016-12-01	Ø	30	
18	2016-12-02	2017-01-29	87	
19	2016-12-02	2017-01-13	87	
20	2016-12-02	2017-01-15	87	
21	2016-12-02	2017-01-15	87	
22	2016-12-02	2017-01-24		
23	2016-12-02	2017-01-14		
24	2010 12 02	2017 01 10	07	

1.2 Next we will determine the range of months of data provided. How many months will we be able to calculate churn for?

Here we just did a simple "SELECT MIN" and "SELECT MAX" of the subscription_start column to see that the data starts on 2016-12-01 and the last subscription_start ends on 2017-03-30. We can see that the company have been in business for 4 months

Since all subscriptions must last for at least a month this means we will have no subscription_end in December, which means we will be able to calculate churn for 3 months, Jan, Feb and March.



Query	Results
MIN(subscription_start)	MAX(subscription_start)
2016-12-01	2017-03-30
Database	e Schema
subscr	iptions 2000 rows
id	INTEGER
subscription_start	TEXT
subscription_end	TEXT
segment	INTEGER

2. What is the overall churn trend since the company started?

2.1 First Step, we create a temporary table of months

So our first step here is just to create a table of the beginning and ending days of each month so we have something to use as a marker in our case statements later.

Query R	Query Results			
first_day	last_day			
2017-01-01	2017-01-31			
2017-02-01	2017-02-28			
2017-03-01	2017-03-31			

```
WITH months AS(
SELECT
  '2017-01-01' AS first day,
  '2017-01-31' AS last day
UNION
SELECT
  '2017-02-01' AS first day,
  '2017-02-28' AS last day
UNION
SELECT
  '2017-03-01' AS first day,
  '2017-03-31' AS last day
SELECT *
FROM months;
```

2.2 Create a cross join table from table and months

In this step we create another temporary table which merges the subscriptions table and the newly created months table. We use a cross join which basically merges every row with one table with every row of the other table. As you can see in the query results below we now have a first_day and last_day for each month for every id.

	Query Results				
id	subscription_start	subscription_end	segment	first_day	last_day
1	2016-12-01	2017-02-01	87	2017-01-01	2017-01-31
1	2016-12-01	2017-02-01	87	2017-02-01	2017-02-28
1	2016-12-01	2017-02-01	87	2017-03-01	2017-03-31
2	2016-12-01	2017-01-24	87	2017-01-01	2017-01-31
2	2016-12-01	2017-01-24	87	2017-02-01	2017-02-28
2	2016-12-01	2017-01-24	87	2017-03-01	2017-03-31
3	2016-12-01	2017-03-07	87	2017-01-01	2017-01-31
3	2016-12-01	2017-03-07	87	2017-02-01	2017-02-28

```
test.sglite
WITH months AS
(SELECT
  '2017-01-01' AS first day,
  '2017-01-31' AS last day
UNION
SELECT.
  '2017-02-01' AS first day,
  '2017-02-28' AS last day
UNION
SELECT
  '2017-03-01' AS first day,
  '2017-03-31' AS last day
cross join AS
(SELECT *
FROM subscriptions
CROSS JOIN months)
SELECT *
FROM cross join
LIMIT 100;
```

2.3 Creating a status table

Here we start our first step in creating our status table. We use a CASE statement to create two new columns. is_active_87 gives us a 1 if the user was active at the start of that month and they are part of segment 87 and a 0 otherwise. Similarly, is_active_30 gives us a 1 if the user was active at the start of that month and they are part of segment 30 and a 0 otherwise.

Query Results				
id	month	is_active_87	is_active_30	
1	2017-01-01	1	0	
1	2017-02-01	1	0	
1	2017-03-01	0	0	
2	2017-01-01	1	0	
2	2017-02-01	0	0	
2	2017-03-01	0	0	
3	2017-01-01	1	0	
3	2017-02-01	1	0	
3	2017-03-01	1	0	
4	2017-01-01	1	0	
4	2017-02-01	1	0	
4	2017-03-01	0	0	

```
status AS
(SELECT id, first day AS month,
CASE
  WHEN (subscription start < first_day)
   AND (segment = 87)
   AND (
     subscription end >= first day
     OR subscription end IS NULL
   ) THEN 1
   ELSE 0
 END as is active 87,
 CASE
  WHEN (subscription start < first day)
   AND (segment = 30)
   AND (
     subscription end >= first day
     OR subscription_end IS NULL
   ) THEN 1
  ELSE 0
 END as is active 30
FROM cross join)
 SELECT *
 FROM status
LIMIT 100;
```

2.4 Adding the is_canceled columns

In this segment we use additional CASE statements to create additional columns to create values for users of each segment that have canceled within that month.

Query Results					
id	month	is_active_87	is_active_30	is_canceled_87	is_canceled_30
1	2017-01-01	1	0	0	0
1	2017-02-01	1	0	1	0
1	2017-03-01	0	0	0	0
2	2017-01-01	1	0	1	0
2	2017-02-01	0	0	0	0
2	2017-03-01	0	0	0	0
3	2017-01-01	1	0	0	0
3	2017-02-01	1	0	0	0
3	2017-03-01	1	0	1	0
4	2017-01-01	1	0	0	0
4	2017-02-01	1	0	1	0
4	2017-03-01	0	0	0	0

```
CASE
  WHEN subscription end
   BETWEEN first day AND last day
   AND segment = 87
   THEN 1
  ELSE 0
END as is canceled 87,
 CASE
  WHEN subscription end
   BETWEEN first day AND last day
   AND segment = 30
   THEN 1
  ELSE 0
END as is canceled 30
FROM cross join)
```

2.5 Create a status_aggregate table

Finally we create another temporary table that gives us the sum of all the 1's in each of the columns. We create the status_aggregate table by using the SUM command on each of the columns from our status table, then group those sums by month.

```
FROM cross_join),
status_aggregate AS
(SELECT month,
SUM(is_active_87) as sum_active_87,
SUM(is_canceled_87) as sum_canceled_87,
SUM(is_active_30) as sum_active_30,
SUM(is_canceled_30) as sum_canceled_30
FROM status
GROUP BY month)
SELECT *
FROM status_aggregate;
```

		Query Results		
month	sum_active_87	sum_canceled_87	sum_active_30	sum_canceled_30
2017-01-01	279	70	291	22
2017-02-01	467	148	518	38
2017-03-01	541	258	718	84

3. Compare the churn rate between user segments

3.1 Churn comparisons between user segments

Here we see the final results:

	User Segment 87	User Segment 30
January	25%	7.5%
February	31.6%	7.3%
March	47.6%	11.6%

Query Results			
month	churn_rate_87	churn_rate_30	
2017-01-01	0.25089605734767	0.0756013745704467	
2017-02-01	0.316916488222698	0.0733590733590734	
2017-03-01	0.476894639556377	0.116991643454039	

```
61 SELECT
62 month,
63 1.0 * sum_canceled_87/sum_active_87 AS churn_rate_87,
64 1.0 * sum_canceled_30/sum_active_30 AS churn_rate_30
65 FROM status_aggregate;
```

Conclusion: Here we can see that the churn rate for user segment 87 is much higher, showing that up to 50 percent of users have canceled during March. The company should clearly focus on expanding the growth of user segment 30.

Bonus

How would we modify this code to support a large

number of segments?

Here we add the segment column to our status table. Then we add it to the status_aggreate table. Then just add it to the final table and group by segment as well as month.

Query Results				
month	segment	churn_rate		
2017-01-01	30	0.0756013745704467		
2017-01-01	87	0.25089605734767		
2017-02-01	30	0.0733590733590734		
2017-02-01	87	0.316916488222698		
2017-03-01	30	0.116991643454039		
2017-03-01	87	0.476894639556377		

```
(SELECT id, segment, first day as month,
CASE
 WHEN (subscription start < first day)
   AND (
      subscription end >= first day
     OR subscription end IS NULL
    ) THEN 1
 ELSE 0
END as is active,
CASE
 WHEN subscription end BETWEEN first day AND last day THEN 1
 ELSE 0
END as is canceled
FROM cross join),
status aggregate AS
(SELECT month, segment,
 SUM(is active) as active,
 SUM(is canceled) as canceled
FROM status
GROUP BY month, segment)
SELECT.
 month, segment,
 1.0 * canceled/active AS churn rate
FROM status aggregate;
```