

# Learn SQL from Scratch Capstone Project

Danielle Furlani April 2019



### TABLE OF CONTENTS

**OPTION 2**: Churn Rates with Codeflix

Analyze churn for two groups of Codeflix users and recommend which segment represents a better long-term user base.

Codeflix, a streaming video startup, is interested in measuring their user churn rate

The marketing department is particularly interested in how the churn compares between two segments of users.

- 1. Get familiar with the company.
  - What segments of users exist?
  - How many months has the company been operating? Which months do you have enough information to calculate a churn rate?
- 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) Get familiar with the company, Codeflix

What segments of users exist?

Take a look at the first 100 rows of data in the subscriptions table

• SELECT \* allows me to view data for all the columns in the subscriptions table

• LIMIT is used to look at the first 100 rows of data

1 SELECT \*

2 FROM subscriptions

3 LIMIT **100**;

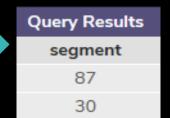
**Query Results** subscription\_start subscription\_end segment 2016-12-01 2017-02-01 87 2016-12-01 2017-01-24 87 2016-12-01 2017-03-07 2017-02-12 2016-12-01 2016-12-01 2017-03-09 2016-12-01 2017-01-19 87 2016-12-01 2017-02-03 2016-12-01 2017-03-02 9 2016-12-01 2017-02-17 10 2016-12-01 2017-01-01 2016-12-01 2017-01-17 12 2016-12-01 2017-02-07 2016-12-01 30

How many different segments do you see?

- From my first query results I saw two different segments, 87 and 30
- SELECT DISTINCT segments helped me determine that 87 and 30 were the only segments in the table

1 SELECT DISTINCT segment

FROM subscriptions;



### 1) Get familiar with the company, Codeflix (cont.)

How many months has the company been operating?

**2**)

Determine the range of months of data provided

- MIN() allows me to see the earliest subscription start date
- MAX() allows me to see the latest subscription start date
- AS allows me to rename each column

- 1 SELECT MIN(subscription\_start) AS min,
- 2 MAX(subscription\_start) AS max
- 3 FROM subscriptions;

From this query we can determine that the company has been operating for 4 months

The range of months for our date is between December 1, 201 and March 30, 2017 Which months will you be able to calculate churn for?

I can calculate the churn rate for the first three months of 2017

We can't calculate for Dec 2016, because there are no subscription\_end values

Query Results			
min	max		
2016-12-01	2017-03-30		
Database Schema			
subscr	riptions 2000 rows		
id	INTEGER		
subscription_start	TEXT		
subscription_end	TEXT		
segment	INTEGER		

1

You'll be calculating the churn rate for both segments over the first 3 months of 2017. To get started, create a temporary table of months

- WITH allows me to create a temporary table, months, to be used in the churn calculation
- UNION allows me to stack one dataset on top of another

Calculating churn rate

- SELECT the first day of each month to be used as a cutoff for subscribers
- SELECT the last day of each month to be used as the cutoff for cancellations

```
1 WITH months AS
2 (SELECT '2017-01-01' AS first_day,
3    '2017-01-31' AS last_day
4 UNION
5 SELECT '2017-02-01' AS first_day,
6    '2017-02-28' AS last_day
7 UNION
8 SELECT '2017-03-01' AS first_day,
9    '2017-03-31' AS last_day)
10 SELECT *
11 FROM months;
```

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

Create a temporary table, cross\_join, from subscriptions and your months

AS allows me to name this new joined table, cross\_join. This query was placed directly after the query that created the months table

```
1 WITH months AS
2 (SELECT '2017-01-01' AS first_day,
3    '2017-01-31' AS last_day
4 UNION
5 SELECT '2017-02-01' AS first_day,
6    '2017-02-28' AS last_day
7 UNION
8 SELECT '2017-03-01' AS first_day,
9    '2017-03-31' AS last_day),
10 cross_join AS
11 (SELECT *
12 FROM subscriptions
13 CROSS JOIN months)
14 SELECT *
15 FROM cross_join
16 LIMIT 50;
```

I joined all columns of the subscriptions table with all columns of the temporary months table using CROSS JOIN

CROSS JOIN simply combines columns from multiple tables together, it does not join data

This step is necessary to calculate churn because we can now view which segment of users correspond with each subscriptions first and last days. This will help to determine which segment of users represent a better long-term user base.

		Query Res	ults			^
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	ı
3	2016-12-01	2017-03-07	87	2017-01-01	2017-01-31	ı
3	2016-12-01	2017-03-07	87	2017-02-01	2017-02-28	ı
3	2016-12-01	2017-03-07	87	2017-03-01	2017-03-31	ı
4	2016-12-01	2017-02-12	87	2017-01-01	2017-01-31	ı
4	2016-12-01	2017-02-12	87	2017-02-01	2017-02-28	ı
4	2016-12-01	2017-02-12	87	2017-03-01	2017-03-31	ı
5	2016-12-01	2017-03-09	87	2017-01-01	2017-01-31	
5	2016-12-01	2017-03-09	87	2017-02-01	2017-02-28	
5	2016-12-01	2017-03-09	87	2017-03-01	2017-03-31	
6	2016-12-01	2017-01-19	87	2017-01-01	2017-01-31	
6	2016-12-01	2017-01-19	87	2017-02-01	2017-02-28	
6	2016-12-01	2017-01-19	87	2017-03-01	2017-03-31	
7	2016-12-01	2017-02-03	87	2017-01-01	2017-01-31	
7	2016-12-01	2017-02-03	87	2017-02-01	2017-02-28	
7	2016-12-01	2017-02-03	87	2017-03-01	2017-03-31	
8	2016-12-01	2017-03-02	87	2017-01-01	2017-01-31	
8	2016-12-01	2017-03-02	87	2017-02-01	2017-02-28	
8	2016-12-01	2017-03-02	87	2017-03-01	2017-03-31	
9	2016-12-01	2017-02-17	87	2017-01-01	2017-01-31	
9	2016-12-01	2017-02-17	87	2017-02-01	2017-02-28	

Create a temporary table, *status*, from cross\_join. This table should contain:

5)

- id selected from cross\_join
- month as an alias of first\_day
- is\_active\_87 and is\_active 30 created using a CASE WHEN to find any users from each segment who existed prior to the beginning of the month. This is 1 if true and 0 otherwise.

The results show us the months a user had an active subscription and which segment the user is from

1 represents a true statement in a CASE WHEN subscription\_start is prior to first\_day AND subscription\_end is after first\_day OR subscription\_end IS NULL

O represents if the CASE WHEN statement is false

This data will be used as the denominator in our churn calculation for each segment

```
FROM subscriptions
 CROSS JOIN months),
status AS
(SELECT id, first_day AS month,
CASE
 WHEN (segment = 87)
 AND (subscription start < first day)
AND (subscription end > first day
      OR subscription end IS NULL) THEN 1
  ELSE 0
END AS is active 87,
CASE
 WHEN (segment = 30)
 AND (subscription start < first day)
 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 50;
```

Query Results			^	
id	month	is_active_87	is_active_30	
1	2017-01-01	1	0	
1	2017-02-01	0	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	
5	2017-01-01	1	0	
5	2017-02-01	1	0	
5	2017-03-01	1	0	
6	2017-01-01	1	0	
6	2017-02-01	0	0	
6	2017-03-01	0	0	
7	2017-01-01	1	0	
7	2017-02-01	1	0	
7	2017-03-01	0	0	

Add an is\_canceled\_87 and an is\_canceled\_30 column to the status temporary table. This should be 1 if the subscription is canceled during the month and 0 otherwise.

A CASE WHEN statement can be used to determine whether a subscription was cancelled and in what month for each segment

```
OR subscription end IS NULL)
      THEN 1
      ELSE 0
    END AS is_active_30,
    CASE
     WHEN segment = 87
     AND (subscription_end BETWEEN first_day
          AND last day) THEN 1
      ELSE @ END AS is canceled 87,
    CASE
      WHEN segment = 30
     AND (subscription end BETWEEN first day
          AND last day) THEN 1
      ELSE 0
     END as is canceled 30
    FROM cross join)
    SELECT *
    FROM status
46 limit 50;
```

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	0	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
5	2017-01-01	1	0	0	0
5	2017-02-01	1	0	0	0
5	2017-03-01	1	0	1	0
6	2017-01-01	1	0	1	0
6	2017-02-01	0	0	0	0
6	2017-03-01	0	0	0	0
7	2017-01-01	1	0	0	0
7	2017-02-01	1	0	1	0
7	2017-03-01	0	0	0	0

The results now show us the months a user had an active subscription, if they cancelled their subscription and when, as well as which segment the user is from

For the new is\_canceled columns, 1 represents a true statement in a CASE WHEN subscription\_end is BETWEEN first\_day AND last\_day 0 represents if the CASE WHEN statement is false

This data will be used as the numerator in our churn calculation for each segment

Create a status\_aggregate temporary table that is a SUM of the active and canceled subscriptions for each segment, for each month.

The resulting columns should be:

sum\_active\_87

sum\_active\_30

sum\_canceled\_87

sum\_canceled\_30

SUM()

 SUM () is an aggregate function, meaning that it can perform calculations on multiple rows

 With each calculation we have found the sum of the values in each column of our status table

Calculating churn rate

- GROUP BY month allows me to find the SUM() for each month
- We now have our numerators and denominators to calculate the churn for each segment

```
END AS is active 30,
CASE
 WHEN segment = 87
 AND (subscription end BETWEEN first day
      AND last day) THEN 1
  ELSE @ END AS is canceled 87,
CASE
  WHEN segment = 30
 AND (subscription end BETWEEN first day
      AND last day) THEN 1
  ELSE 0
 END as is canceled 30
FROM cross join),
status aggregate AS
(SELECT month.
 SUM(is active 87) AS sum_active_87,
 SUM(is active 30) AS sum active 30,
 SUM(is canceled 87) AS sum canceled 87,
 SUM(is canceled 30) AS sum canceled 30
FROM status
GROUP BY month)
SELECT *
FROM status aggregate
limit 100;
```

		Query Results		
month	sum_active_87	sum_active_30	sum_canceled_87	sum_canceled_30
2017-01-01	278	291	70	22
2017-02-01	462	518	148	38
2017-03-01	531	716	258	84

### 3) Compare the churn rates between user segments

8) Calculate the churn rates for the two segments over the three month period

To calculate churn we must divide the number of cancellations for each month by the number of active subscribers at the beginning of each month and multiply by 1 to force a float result instead of an integer.

We must do this for each segment

### Segment 87

 $1 \times 70/278 = January 2017 churn$ 

 $1 \times 148/462 = February 2017 churn$ 

 $1 \times 258/531 = March 2017 churn$ 

### Segment 30

 $1 \times 22/291 = January 2017 churn$ 

 $1 \times 38/518 = February 2017 churn$ 

 $1 \times 84/716 = March 2017 churn$ 

```
AND last day) THEN 1
  ELSE @ END AS is canceled 87,
CASE
  WHEN segment = 30
 AND (subscription end BETWEEN first day
      AND last day) THEN 1
  ELSE 0
 END as is canceled 30
FROM cross join),
status_aggregate AS
(SELECT month,
 SUM(is active 87) AS sum active 87,
SUM(is active 30) AS sum active 30,
SUM(is canceled 87) AS sum canceled 87,
SUM(is canceled 30) AS sum canceled 30
FROM status
GROUP BY month)
SELECT month,
  1.0 * sum_canceled_87/sum_active_87 AS churn_87,
 1.0 * sum canceled 30/sum active 30 AS churn 30
FROM status aggregate;
```

# 3) Compare the churn rates between user segments (cont.) Which segment of users should the company focus on expanding?

8)

Which segment has a lower churn rate?

We've found our churn rates for each segment over the course of 3 months.

January 2017: 25% for segment 87 & 7% for segment 30 February 2017: 32% for segment 87 & 7% for segment 30

March 2017: 48% for segment 87 & 11% for segment 30

Query Results			
month	churn_87	churn_30	
2017-01-01	0.251798561151079	0.0756013745704467	
2017-02-01	0.32034632034632	0.0733590733590734	
2017-03-01	0.485875706214689	0.11731843575419	

By comparing segment 87 and 30 for each month, we see that the churn is significantly lower for segment 30 for each month

After analyzing the churn for the two groups of Codeflix users, I would highly recommend that segment 30 represents a better long-term user base

1