

# SQL for Data Science Bootcamp

## Final Project

Congratulations! You have completed this bootcamp and qualified to take this final project examination.

In this final project, you are asked to work mainly on two different datasets.

### **Method:**

- You are going to provide your SQL queries in a written format in a docs like file for each question.
- You can use BigQuery to try your queries before summiting your answers.
- When you write or try your queries, you need project name, dataset name and table name. You can call you project name as what ever you would like to, but please create a dataset under your project, and give **final\_project** name to this dataset. Then create your tables under it.

### **Question 1:**

You have three tables:

Customers table:

<https://drive.google.com/file/d/17oRPW3jA2mCsTOXKKt6Xx7wSgD78twq8/view?usp=sharing>

Orders table:

<https://drive.google.com/file/d/1TWBB9rqGRnWOJsDcwkUE2HkJNEafawQm/view?usp=sharing>

Payments table

[https://drive.google.com/file/d/1vakf\\_7F9cb7ZKOeCofyI370Paiswg3JX/view?usp=sharing](https://drive.google.com/file/d/1vakf_7F9cb7ZKOeCofyI370Paiswg3JX/view?usp=sharing)

- A. Which payment\_type has the fastest average delivery day where delivery day can be found by using order\_delivered\_customer\_date and order\_purchase\_timestamp columns? You should be careful about choosing these dates that shouldn't be null.  
(Hint: The result should show only one payment\_type with the average\_delivery\_day info)
- B. What is the maximum payment\_value of each state where order status is neither unavailable nor canceled? Additionally, the max\_payment\_revenue column should have 1 decimal point.  
(Hint: The result table will only have state and max\_payment\_revenue columns)
- C. In the payments table, how many distinct order\_ids that exist more than once?  
(Hint: Try window functions!)
- D. For each row, find the payment\_value of previous row and leading row in the partition of order\_id and order of payment\_sequential ascendingly.  
(Hint: Again window functions!  
The result table should have all columns in the payments table and two new columns: previous\_payment\_value and leading\_payment\_value.  
Also, we don't want null values for these two new fields)

- E. Get all the columns in the customers table and add three new columns:
- \* order\_id
  - \* order\_purchase\_date (Hint: convert the timestamp to the date value)
  - \* payment\_value
  - \* payment\_value\_int (Hint: payment\_value but in the integer format. You should use a function to convert it type)
- (In this join, we want all the rows from customers table, and if there is corresponding data in either orders or payments, we can have it. But we definitely want to have all info from customers table. Think about your join!)
- F. Get all columns from the payments table excluding the payment\_type column, and create payment\_method column instead of that. This new column should have this format:
- If payment\_type is “not\_defined”, then payment\_method will be 0  
 If payment\_type is “credit\_card”, then payment\_method will be 1  
 If payment\_type is “voucher”, then payment\_method will be 2  
 If payment\_type is “debit\_card”, then payment\_method will be 3  
 If payment\_type is “boleto”, then payment\_method will be 4

### **Question 2:**

You have one table:

Netflix\_movies: [https://drive.google.com/file/d/1VRdGqa3KmgDj\\_BEX4\\_OBEtDa-JRmE1LH/view?usp=sharing](https://drive.google.com/file/d/1VRdGqa3KmgDj_BEX4_OBEtDa-JRmE1LH/view?usp=sharing)

- A. How many times each word is used in the table, that you get when you split the movie\_name by each space?  
 (Hint: The result table only have 2 columns: words\_used\_in\_movie\_names, number\_of\_occurence.  
 Ex. You can imagine that asks for how many times you see Rome word in this table)
- B. If a movie\_name has “**Vol. 1**”, then change it to “Part 1”; if has “**Vol. 2**”, then change it to “Part. 2”. Else movie\_name. After these changes, name show this column as “movie\_name” instead of showing the original movie\_name column.
- C. In this table, we have some movie\_names that occur multiple times. Create a new column in the result table, call this column as “**previous\_year**” and show the previous year of each movie if they exist multiple times in that table. In the result table, there shouldn’t be any null previous\_year values.