

Writing SQL Like a Pro: Advanced Techniques Showcased in a Real-Life Scenario

Merging two history tables using advanced SQL techniques with Google's BigQuery.



Furcy Pin · Follow

Published in Learning SQL · 9 min read · Mar 5, 2024

740

8



...



Click to download Learning SQL's free interview guide



Image by author, generated with Stable Diffusion using the prompt "A pixel art landscape representing advanced SQL techniques, high quality". I've been told that having a nice image is better for SEO, so here it is.

Last week, Arnaud Stephan, a fellow Analytics Engineer, asked for help with an advanced use case that I found very interesting. After spending a couple of hours helping him, I thought it would make a perfect subject for a blog article, and I asked for his permission to discuss this advanced real-life problem in public and publish the associated code.

So, let's saddle up and get ready for some advanced hands-on SQL with my favorite engine: Google BigQuery!

(Note: To preserve anonymity, I replaced the IDs and values with dummy ones.)

About the author:

I am a freelance Analytics Engineer / Data Platform Engineer with more than 12 years of experience in SQL and Big Data. I love to play with complex SQL/DataFrame topics, so if you're stuck on challenging problems that you can't crack with SQL, don't hesitate to reach out; I might be able to help.

The Problem

First, let's give some context: We have data from a grocery retailer selling multiple articles in multiple shops. For each article inside each shop, they track the pricing history of these articles. The corresponding data is contained in the table `sell_price_history`, which looks like this:

Query results						
JOB INFORMATION		RESULTS	CHART	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	article_id	shop_id		sell_price	start_datetime	end_datetime
1	yVxrDmtw	dir29L9xF		null	2022-05-10 11:45:19 UTC	2022-05-10 11:53:49 UTC
2	yVxrDmtw	dir29L9xF		null	2022-05-10 11:53:49 UTC	2022-05-10 13:08:38 UTC
3	yVxrDmtw	dir29L9xF		null	2022-05-10 13:08:38 UTC	2022-05-10 13:08:40 UTC
4	yVxrDmtw	dir29L9xF		null	2022-05-10 13:08:40 UTC	2022-05-10 13:08:42 UTC
5	yVxrDmtw	dir29L9xF		null	2022-05-10 13:08:42 UTC	2022-05-10 13:08:45 UTC
6	yVxrDmtw	dir29L9xF		3.99	2022-05-10 13:08:45 UTC	2022-05-10 13:23:18 UTC
7	yVxrDmtw	dir29L9xF		3.99	2022-05-10 13:23:18 UTC	2022-05-11 04:20:13 UTC
8	yVxrDmtw	dir29L9xF		3.99	2022-05-11 04:20:13 UTC	2023-06-14 10:46:07 UTC
9	yVxrDmtw	dir29L9xF		3.99	2023-06-14 10:46:07 UTC	2024-02-14 00:59:59 UTC

I decided to keep the first rows for this example because this is the kind of weird stuff you get when working with "real" data. The person I helped explained to me that some other columns (that they did not include in the example) did change, so it was important for them to keep those lines, even when the price was null. Image by author.

But each shop can also create *promotions* for any article, in which case the `promotion_value` would replace the original article price. A second table named `promotion_history` has the corresponding data:

Query results						
JOB INFORMATION		RESULTS	CHART	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	shop_id	article_id	promotion_id	promotion_value	start_datetime	end_datetime
1	dir29L9xF	yVxrDmtw	-38_P0evh	2.99	2022-05-11 00:00:00 UTC	2022-05-11 05:03:04 UTC
2	dir29L9xF	yVxrDmtw	-38_P0evh	2.49	2022-05-11 05:03:04 UTC	2022-05-18 00:00:00 UTC
3	dir29L9xF	yVxrDmtw	2GLCMItlyw	2.99	2023-06-14 15:05:31 UTC	2023-06-16 11:29:38 UTC
4	dir29L9xF	yVxrDmtw	2GLCMItlyw	2.99	2023-06-16 11:29:38 UTC	2023-06-19 08:58:40 UTC
5	dir29L9xF	yVxrDmtw	2GLCMItlyw	2.99	2023-06-19 08:58:40 UTC	2023-06-22 13:31:46 UTC
6	dir29L9xF	yVxrDmtw	2GLCMItlyw	1.99	2023-06-22 13:31:46 UTC	2023-06-27 23:59:59.999000 UTC

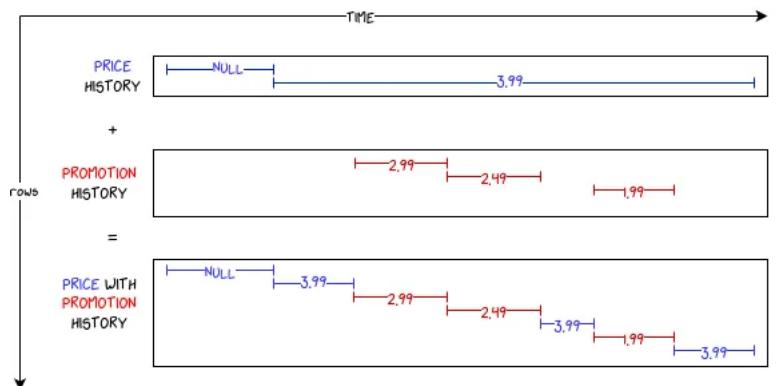
Image by author.

The problem then follows: We want to “merge” these two event tables into one single table that will give us, for each article, the history of prices, *including the effects of promotions.*

In other words, the expected result we want should look like this:

Image by author.

To give a more visual illustration, we want to merge two tables describing ranges of events.



The problem we want to solve. Image by author.

The (Incorrect) Solution

To make things a bit more lively, I will walk you through how I solved the problem. I had a rough idea that to achieve that, I needed to decompose each table into a kind of “event timeline” and then merge the two timelines and fill in the blanks. This might sound a bit fuzzy for now, but so did it when I started solving the problem. Don’t worry; things will get clearer soon.

Step 1: Co-locating the Data

The very first thing I did before even trying to think about how to solve the problem was join the two tables together and group the data by *shop_id* and *article_id*, like this:

v1.price_and_promotion_history		Run	Save Query	Share	Schedule	More	This query will	
Job Information	Results	Chart	JSON	Execution Details	Execution Graph	Logs	Save Results	Explode
Row	shop_id = * d12345ef	article_id = * priceHistory	pri.article_id = * pri.article_id	pri.article_id = * pri.article_id	pri.history_start_datetime = * 2022-05-10 11:45:19 UTC	pri.history_end_datetime = * 2022-05-10 11:53:49 UTC	promotion_id = * 38_Peoch	promotion_value = * 2.99
				null	2022-05-10 11:45:49 UTC	2022-05-10 11:45:49 UTC	38_Peoch	2.99
				null	2022-05-10 11:46:38 UTC	2022-05-10 11:46:38 UTC	20LCLM1yfw	2.99
				null	2022-05-10 11:46:40 UTC	2022-05-10 11:46:40 UTC	20LCLM1yfw	2.99
				null	2022-05-10 11:48:42 UTC	2022-05-10 11:48:45 UTC	20LCLM1yfw	2.99
				3.99	2022-05-10 11:48:45 UTC	2022-05-10 11:23:31 UTC	20LCLM1yfw	1.99
				3.99	2022-05-10 12:32:18 UTC	2022-05-10 12:40:31 UTC		
				3.99	2022-05-11 04:02:13 UTC	2022-06-14 04:06:07 UTC		
				3.99	2023-04-14 10:46:07 UTC	2024-02-14 10:49:59 UTC		

Image by author

For those of you not familiar with nested structures in BigQuery, this might look a bit weird at first, so here is a detailed explanation of how the data structures are organized:

Row	shop_id	article_id	on_sale_since	price_history_start_datetime	price_history_end_datetime	promotion_id	promotion_id	promotion_value	promotion_start_datetime	promotion_end_datetime
1	d1c92ef	y10e0mte	2020-05-10 13:19:15UTC	2020-05-10 13:19:15UTC	2020-05-10 13:19:15UTC	28_Powh	28_Powh	2.95	2020-05-11 00:00:00 UTC	2022-05-11 00:00:00 UTC
			null	2020-05-10 13:05:49UTC	2020-05-10 13:08:30UTC			2.49	2020-05-11 00:03:44 UTC	2022-05-11 00:00:00 UTC
			null	2020-05-10 13:08:38UTC	2020-05-10 13:08:40UTC	25GLMltppw		2.99	2020-06-14 05:03:17 UTC	2023-06-16 11:20:38 UTC
			null	2020-05-10 13:08:40UTC	2020-05-10 13:08:42UTC	25GLMltppw		2.99	2020-06-16 11:29:30 UTC	2023-06-19 08:50:40 UTC
			null	2020-05-10 13:08:42UTC	2020-05-10 13:08:45UTC	25GLMltppw		2.99	2020-06-19 08:56:02 UTC	2023-06-22 13:31:46 UTC
			3.99	2020-05-10 13:08:45UTC	2020-05-10 13:21:18UTC	25GLMltppw		1.99	2020-06-22 13:31:46 UTC	2023-06-27 23:59:59 9999000 U
			3.99	2020-05-10 13:22:18UTC	2020-05-10 20:13:07UTC					
			3.99	2020-05-11 04:20:13UTC	2020-06-14 10:46:07UTC					
			3.99	2020-06-14 10:46:07UTC	2020-06-22 14:00:59 5999000 U					

Image by author

There are multiple reasons for me to prefer this way of arranging things:

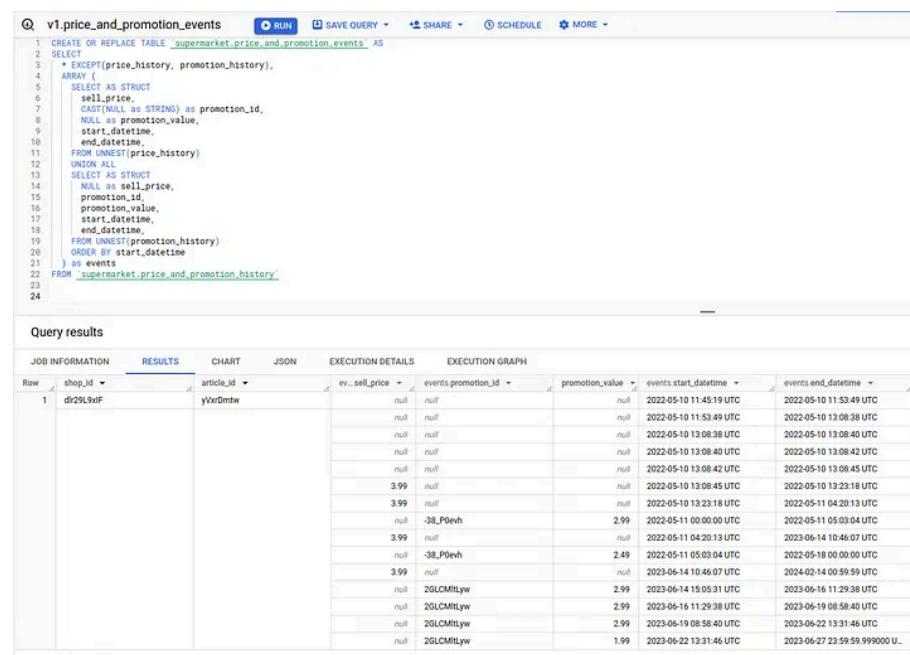
- Simplicity: Co-locating the data for each article allows me to reason on a single-row basis.
 - Performance: Thanks to this pre-grouping, all the computations I will perform will be made at the row level. For those familiar with MapReduce, there will only be Map steps, no Reduce, and most of all, no Shuffle, which is generally the most resource-intensive part.
This is not “crucial” on BigQuery, as query price is only based on the amount of data read and not processed, but it would be on other engines, and I like to optimize my code.
 - Storage efficiency: Since *shop_id* and *article_id* are not repeated at each line, this representation takes less storage space for the same quantity of information. This doesn’t make much sense on a toy dataset like this, but to get an idea, the total size of the two input tables was 720b, but the size of the output table (which contains the very same information) was only 412b. Multiply this by billions of lines, and you might notice a price difference.

The main drawbacks are that most people aren't used to it, so they might require some effort to learn how to query it properly, and many metadata tools (data catalogs, data diff, etc.) don't support nested structures very well yet.

But enough digressing. That subject alone would deserve a whole article, so let's go back to the task at hand.

Step 2: Merging the Two Histories

Next, we perform a *union* between the two history tables like this:



```

1 CREATE OR REPLACE TABLE `supermarket.price_and_promotion_events` AS
2 SELECT
3   * EXCEPT(price_history, promotion_history),
4   ARITHMETIC(price_history,
5     SELECT AS STRUCT
6       sell_price,
7       CAST(NULL as STRING) as promotion_id,
8       NULL as promotion_value,
9       start_datetime,
10      end_datetime,
11     FROM UNNEST(price_history)
12   UNION ALL
13   SELECT AS STRUCT
14     NULL as sell_price,
15     promotion_id,
16     promotion_value,
17     start_datetime,
18     end_datetime,
19     FROM UNNEST(promotion_history)
20   ORDER BY start_datetime
21 ) as events
22 FROM `supermarket.price_and_promotion_history`
23
24

```

Query results

Row	shop_id	article_id	ev_sell_price	events.promotion_id	promotion_value	events.start_datetime	events.end_datetime
1	dr29L9xif	yVwDmtw	null	null	null	2022-05-10 11:45:19 UTC	2022-05-10 11:53:49 UTC
			null	null	null	2022-05-10 11:53:49 UTC	2022-05-10 13:08:38 UTC
			null	null	null	2022-05-10 13:08:38 UTC	2022-05-10 13:08:40 UTC
			null	null	null	2022-05-10 13:08:40 UTC	2022-05-10 13:08:42 UTC
			3.99	null	null	2022-05-10 13:08:42 UTC	2022-05-10 13:08:45 UTC
			3.99	null	null	2022-05-10 13:08:45 UTC	2022-05-10 13:23:18 UTC
			3.99	null	null	2022-05-10 13:23:18 UTC	2022-05-10 13:20:13 UTC
			3.99	-38_P0evh	2.99	2022-05-11 00:00:00 UTC	2022-05-11 05:00:04 UTC
			3.99	null	null	2022-05-11 04:20:13 UTC	2023-06-14 10:46:07 UTC
			3.99	-38_P0evh	2.49	2022-05-11 05:03:04 UTC	2022-05-18 00:00:00 UTC
			3.99	null	null	2023-06-14 10:46:07 UTC	2024-02-14 00:59:59 UTC
			null	2GLCMtlyw	2.99	2023-06-14 15:05:31 UTC	2023-06-16 11:29:38 UTC
			null	2GLCMtlyw	2.99	2023-06-16 11:29:38 UTC	2023-06-19 00:58:40 UTC
			null	2GLCMtlyw	2.99	2023-06-19 08:58:40 UTC	2023-06-22 13:31:46 UTC
			null	2GLCMtlyw	1.99	2023-06-22 13:31:46 UTC	2023-06-27 23:59:59.999000 UTC

To understand this query, imagine that `price_history` is like a mini-table contained inside the top-level row. Then `"SELECT ... FROM UNNEST(price_history)"` is a way to query this mini-table while staying inside the top-level row. `"SELECT AS STRUCT"` is just a way to wrap the result of the select inside a STRUCT. Image by author.

Step 3: Filling in the Blanks

Finally, we use the *last_value* window function to keep the last non-null value in the event orders, to fill in the blanks, or should I say the *NULLs*.

Open in app ↗



```

1  /* NAME: MY_RELATIVE_TABLES_SUPERMARKET_PRICE_AND_PROMOTION_EVENTS */
2  SELECT * FROM supermarket_price_and_promotion_events
3  EXCEPT(events),
4  ARRAY(
5    SELECT AS STRUCT
6      LAST_VALUE(price IGNORE NULLS) OVER (ORDER BY start_datetime) AS sell_price,
7      LAST_VALUE(promotion_id IGNORE NULLS) OVER (ORDER BY start_datetime) AS promotion_id,
8      LAST_VALUE(promotion_value IGNORE NULLS) OVER (ORDER BY start_datetime) AS promotion_value,
9      start_datetime,
10     LEAD(start_datetime) OVER (ORDER BY start_datetime) AS end_datetime
11   FOR EACH ROW
12 ) AS price_with_promotion
13  FROM supermarket_price_and_promotion_events
14

```

Query results

JOB INFORMATION		RESULTS	CHART	JSON	EXECUTION DETAILS	EXECUTION GRAPH	
Row	shop_id	article_id	pr..sell_price	pr..price_with_p..promotion_id	promotion_value	price_with_p..start_datetime	price_with_p..end_datetime
1	dr29L9xif	yVxrDmtw	null	null	null	2022-05-10 11:45:19 UTC	2022-05-10 11:53:49 UTC
			null	null	null	2022-05-10 11:53:49 UTC	2022-05-10 13:08:38 UTC
			null	null	null	2022-05-10 13:08:38 UTC	2022-05-10 12:08:40 UTC
			null	null	null	2022-05-10 13:08:40 UTC	2022-05-10 13:08:42 UTC
			null	null	null	2022-05-10 13:08:42 UTC	2022-05-10 13:08:45 UTC
			3.99	null	null	2022-05-10 13:08:45 UTC	2022-05-10 13:23:18 UTC
			3.99	null	null	2022-05-10 13:23:18 UTC	2022-05-11 00:00:00 UTC
			3.99	-38_P0vh	2.99	2022-05-11 00:00:00 UTC	2022-05-11 04:20:13 UTC
			3.99	-38_P0vh	2.99	2022-05-11 04:20:13 UTC	2022-05-11 05:03:04 UTC
			3.99	-38_P0vh	2.49	2022-05-11 05:03:04 UTC	2023-06-14 10:46:07 UTC
			3.99	-38_P0vh	2.49	2023-06-14 10:46:07 UTC	2023-06-14 15:05:31 UTC
			3.99	2GLCMltiyw	2.99	2023-06-14 15:05:31 UTC	2023-06-16 11:29:38 UTC
			3.99	2GLCMltiyw	2.99	2023-06-16 11:29:38 UTC	2023-06-19 08:58:40 UTC
			3.99	2GLCMltiyw	2.99	2023-06-19 08:58:40 UTC	2023-06-22 13:31:46 UTC
			3.99	2GLCMltiyw	1.99	2023-06-22 13:31:46 UTC	null

The results look good but aren't. Look closely at the highlighted line. Image by author.

While this result does look like the result we are looking for, sharp readers might already have caught that it is, in fact, incorrect.

Why This Does Not Work:

This would have worked if the promotions went on without any interruption, but if we look more closely at our input data, we see that this is not the case: We do have an interruption between the first and the second promotions. This interruption is not visible in the previous screenshot, as we can see on the line highlighted in blue.

```

1  SELECT * FROM supermarket_promotion_history ORDER BY start_datetime
2

```

Query results

JOB INFORMATION		RESULTS	CHART	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	shop_id	article_id	promotion_id	promotion_value	start_datetime	end_datetime
1	dr29L9xif	yVxrDmtw	-38_P0vh	2.99	2022-05-11 00:00:00 UTC	2022-05-11 05:03:04 UTC
2	dr29L9xif	yVxrDmtw	-38_P0vh	2.49	2022-05-11 05:03:04 UTC	2022-05-18 00:00:00 UTC
3	dr29L9xif	yVxrDmtw	2GLCMltiyw	2.99	2023-06-14 15:05:31 UTC	2023-06-16 11:29:38 UTC
4	dr29L9xif	yVxrDmtw	2GLCMltiyw	2.99	2023-06-16 11:29:38 UTC	2023-06-19 08:58:40 UTC
5	dr29L9xif	yVxrDmtw	2GLCMltiyw	2.99	2023-06-19 08:58:40 UTC	2023-06-22 13:31:46 UTC
6	dr29L9xif	yVxrDmtw	2GLCMltiyw	1.99	2023-06-22 13:31:46 UTC	2023-06-27 23:59:59.999000 UTC

The input table "promotion_history" where we highlighted the interruption between promotions. The first promotion ended on 2022-05-18, while the second promotion only started on 2023-06-14, more than one year later. Image by author.

The (Correct) Solution

Now, we will see how the previous (incorrect) solution can be fixed, with more detailed explanations of what is happening.

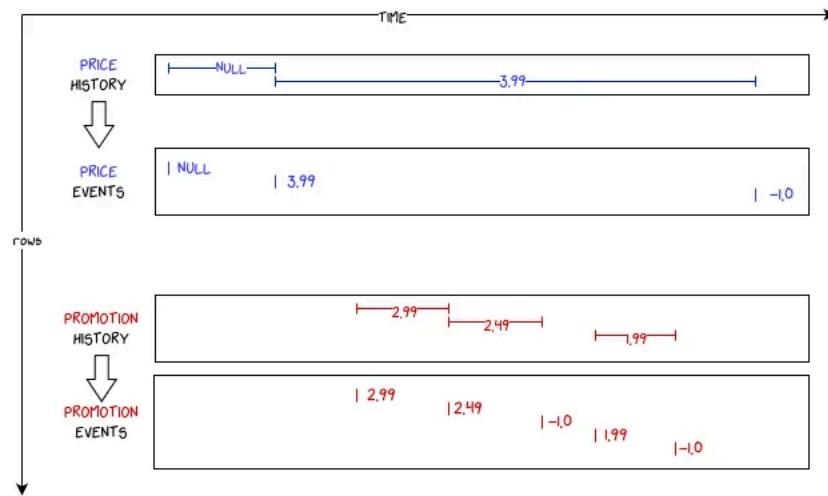
Step 1: Co-locating the Data [Nothing to Change]

Step 1 was just a preparation step. We can keep it as is.

Step 2: Transforming Date Ranges into Timelines and Merging Them

Our intuition of merging the two timelines was good, but we missed the fact that the date ranges could be discontinuous. To fix this, the first step we must achieve is to transform our “history”, which is made of date ranges (from *start_date* to *end_date*), into an event timeline, made only of events (just one *event_date*), without losing any information.

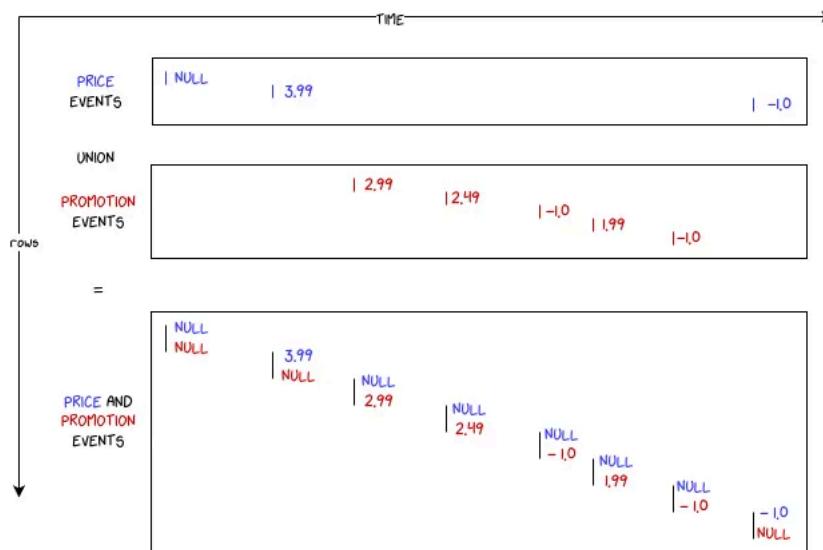
Schematically, this will look like this:



Schematical representation of going from a history (date range) to event representation. Image by author.

Here, the “missing” values have been represented with the *magic number -1*. This is because in step 3. we don’t want to treat them the same way as *NULLS* when we use the window function *last_value(... ignore nulls)*. Arguably, this is a code smell, but we will keep it for the sake of simplicity in this article.

Once this is done, we will simply merge the two timelines with a UNION like this:



Schematical representation of unioning two timelines. Image by author.

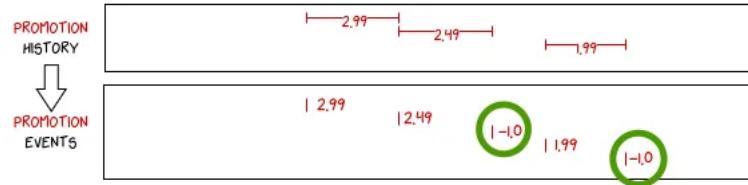
So, what does the code look like? Let's focus only on promotions for a while. The first step is to retrieve the “*start_datetime*” of the next event; this will help us detect when there is a discontinuity:

```
SELECT
...
ARRAY(
  SELECT AS STRUCT
    *,
    LAG(start_datetime) OVER (ORDER BY start_datetime DESC) as next_event_date
  FROM UNNEST(promotion_history)
  ORDER BY start_datetime
) as promotion_history,
...
```

The second step consists of introducing new events whenever such a discontinuity happens, like this:

```
SELECT
...
ARRAY (
  SELECT AS STRUCT
    promotion_value,
    start_datetime,
    FROM UNNEST(promotion_history)
  UNION ALL
  SELECT AS STRUCT
    -1 as promotion_value,
    end_datetime,
    FROM UNNEST(promotion_history)
    WHERE end_datetime < next_event_datetime OR next_event_datetime IS NULL
) as events
...
```

The second part of the union is what reintroduces the previously missing “end-events”:



Circled in green are the newly added end-events. Image by author.

Bringing everything back together, we get a query in two parts:

```
④ v2.price_and_promotion_hist...ime RUN SAVE QUERY SHARE SCHEDULE MORE This query will process
1 CREATE OR REPLACE TABLE `supermarket.price_and_promotion_history_with_next_event_datetime` AS
2 SELECT
3   shop_id,
4   article_id,
5   ARRAY(
6     SELECT AS STRUCT
7       *
8       LAG(start_datetime) OVER (ORDER BY start_datetime DESC) as next_event_datetime
9     FROM UNNEST(price_history)
10    ORDER BY start_datetime
11  ) as price_history,
12  ARRAY(
13    SELECT AS STRUCT
14      *
15      LAG(start_datetime) OVER (ORDER BY start_datetime DESC) as next_event_datetime
16    FROM UNNEST(promotion_history)
17   ORDER BY start_datetime
18  ) as promotion_history,
19 FROM `supermarket.price_and_promotion_history`
20
```

Query results											
JOB INFORMATION		RESULTS		CHART		JSON		EXECUTION DETAILS		EXECUTION GRAPH	
Row	shop_id	article_id	sell_price	price_history.start_datetime	price_history.end_datetime	price..next_event_datetime	promotion_id	promotion_value	promotion_start_datetime	promotion_h..end_datetime	prom..next_event_datetime
1	dfr29l9xIf	yVxxDmtw	null	2022-05-10 11:45:19 UTC	2022-05-10 11:53:49 UTC	2022-05-10 11:53:49 UTC	-38_P0evh	2.99	2022-05-11 00:00:00 UTC	2022-05-11 05:03:04 UTC	2022-05-11 05:03:04 UTC
			null	2022-05-10 11:53:49 UTC	2022-05-10 13:08:38 UTC	2022-05-10 13:08:38 UTC	-38_P0evh	2.49	2022-05-11 05:03:04 UTC	2022-05-18 00:00:00 UTC	2023-06-14 15:05:31 UTC
			null	2022-05-10 13:08:38 UTC	2022-05-10 13:08:40 UTC	2022-05-10 13:08:40 UTC	2GLCMILyW	2.99	2023-06-14 15:05:31 UTC	2023-06-16 11:29:38 UTC	2023-06-16 11:29:38 UTC
			null	2022-05-10 13:08:40 UTC	2022-05-10 13:08:42 UTC	2022-05-10 13:08:42 UTC	2GLCMILyW	2.99	2023-06-16 11:29:38 UTC	2023-06-19 08:58:40 UTC	2023-06-19 08:58:40 UTC
			null	2022-05-10 13:08:42 UTC	2022-05-10 13:08:45 UTC	2022-05-10 13:08:45 UTC	2GLCMILyW	2.99	2023-06-19 08:58:40 UTC	2023-06-22 13:31:46 UTC	2023-06-22 13:31:46 UTC
			3.99	2022-05-10 13:08:45 UTC	2022-05-10 13:23:18 UTC	2022-05-10 13:23:18 UTC	2GLCMILyW	1.99	2023-06-19 08:58:40 UTC	2023-06-22 13:31:46 UTC	2023-06-22 13:31:46 UTC
			3.99	2022-05-10 13:23:18 UTC	2022-05-11 04:20:13 UTC	2022-05-11 04:20:13 UTC					
			3.99	2022-05-11 04:20:13 UTC	2023-06-14 10:46:07 UTC	2023-06-14 10:46:07 UTC					
			3.99	2023-06-14 10:46:07 UTC	2024-02-14 00:59:59 UTC	null					

Part 1: Retrieving the “start_datetime” of the next event. Image by author.

v2.price_and_promotion_events

```

1 CREATE OR REPLACE TABLE `supermarket.price_and_promotion_events` AS
2 SELECT
3   * EXCEPT(price_history, promotion_history),
4   ARRAY (
5     SELECT AS STRUCT
6       sell_price,
7       CAST(NULL as STRING) as promotion_id,
8       NULL as promotion_value,
9       start_datetime,
10      end_datetime,
11    FROM UNNEST(price_history)
12  UNION ALL
13  SELECT AS STRUCT
14    -1 as sell_price,
15    CAST(NULL as STRING) as promotion_id,
16    NULL as promotion_value,
17    end_datetime,
18    FROM UNNEST(price_history)
19  WHERE end_datetime < next_event_datetime OR next_event_datetime IS NULL
20  UNION ALL
21  SELECT AS STRUCT
22    NULL as sell_price,
23    promotion_id,
24    promotion.value,
25    start_datetime,
26    FROM UNNEST(promotion_history)
27  UNION ALL
28  SELECT AS STRUCT
29    NULL as sell_price,
30    -1 as promotion_id,
31    -1 as promotion_value,
32    end_datetime,
33    FROM UNNEST(promotion.history)
34  WHERE end_datetime < next_event_datetime OR next_event_datetime IS NULL
35  ORDER BY start_datetime
36 ) as events,
37 FROM `supermarket.price_and_promotion_history_with_next_event_datetime`
38

```

Query results

JOB INFORMATION		RESULTS	CHART	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	shop_id	article_id	ev_sell_price	events.promotion_id	promotion_value	events.start_datetime
1	dr29L9xdF	yVxrDmtw	null	null	null	2022-05-10 11:45:19 UTC
			null	null	null	2022-05-10 11:51:49 UTC
			null	null	null	2022-05-10 13:08:38 UTC
			null	null	null	2022-05-10 13:08:40 UTC
			3.99	null	null	2022-05-10 13:08:42 UTC
			3.99	null	null	2022-05-10 13:08:45 UTC
			3.99	null	null	2022-05-10 13:21:18 UTC
			null	-38_P0evh	2.99	2022-05-11 00:00:00 UTC
			3.99	null	null	2022-05-11 04:20:13 UTC
			null	-38_P0evh	2.49	2022-05-11 05:03:04 UTC
			null		-1.0	2022-05-18 00:00:00 UTC
			3.99	null	null	2023-06-14 10:46:07 UTC
			null	2GLCMlItLyw	2.99	2023-06-14 15:05:31 UTC
			null	2GLCMlItLyw	2.99	2023-06-16 11:29:38 UTC
			null	2GLCMlItLyw	2.99	2023-06-19 08:58:40 UTC
			null	2GLCMlItLyw	1.99	2023-06-22 13:11:46 UTC
			null		-1.0	2023-06-27 23:59:59.999900 U...
			-1.0	null	null	2024-02-14 00:59:59 UTC

Part 2: Adding the end-event and merging the two timelines. We do both in one go which is why we are unioning 4 sub-tables together. Note that we also used a magic value "" for the promotion_id. Image by author.

Step 3: Filling in the Blanks + Replacing Magic Values

Now that this is done, the previous query from step 3 gives us a result that looks almost like the one we want:

Query results

Row	shop_id	article_id	pr.sell_price	price_and_pr_promotion_id	promotion_value	price_and_start_datetime	price_and_pr_end_datetime
1	dr29L9xF	yVrDmTw	null			2022-05-10 11:45:19 UTC	2022-05-10 11:53:49 UTC
			null			2022-05-10 11:53:49 UTC	2022-05-10 13:08:38 UTC
			null			2022-05-10 13:08:38 UTC	2022-05-10 13:08:40 UTC
			null			2022-05-10 13:08:40 UTC	2022-05-10 13:08:42 UTC
			null			2022-05-10 13:08:42 UTC	2022-05-10 13:08:45 UTC
			3.99			2022-05-10 13:08:45 UTC	2022-05-10 13:23:18 UTC
			3.99			2022-05-10 13:23:18 UTC	2022-05-11 00:00:00 UTC
			3.99	-38_P0evh	2.99	2022-05-11 00:00:00 UTC	2022-05-11 04:20:13 UTC
			3.99	-38_P0evh	2.99	2022-05-11 04:20:13 UTC	2022-05-11 05:03:04 UTC
			3.99	-38_P0evh	2.49	2022-05-11 05:03:04 UTC	2022-05-18 00:00:00 UTC
			3.99	-1.0	2.99	2022-05-18 00:00:00 UTC	2023-06-14 10:46:07 UTC
			3.99	-1.0	2.99	2023-06-14 10:46:07 UTC	2023-06-14 15:05:31 UTC
			3.99	2GLCMltlyw	2.99	2023-06-14 15:05:31 UTC	2023-06-16 11:29:38 UTC
			3.99	2GLCMltlyw	2.99	2023-06-16 11:29:38 UTC	2023-06-19 08:58:40 UTC
			3.99	2GLCMltlyw	2.99	2023-06-19 08:58:40 UTC	2023-06-22 13:31:46 UTC
			3.99	2GLCMltlyw	1.99	2023-06-22 13:31:46 UTC	2023-06-27 23:59:59.999000 U...
			3.99	2GLCMltlyw	1.99	2023-06-27 23:59:59.999000 U...	2024-02-14 00:59:59 UTC
			-1.0	2GLCMltlyw	-1.0	2023-06-27 23:59:59.999000 U...	2024-02-14 00:59:59 UTC
			-1.0	2GLCMltlyw	-1.0	2024-02-14 00:59:59 UTC	null

The execution graph illustrates the transformation process. It starts with 'PRICE AND PROMOTION EVENTS' (raw data). An arrow labeled 'LAST_VALUE(... IGNORE NULLS)' points to a second stage of 'PRICE AND PROMOTION EVENTS'. A third stage follows, labeled 'LEAD(event_datetime) as event_end_datetime'. Finally, the data is shown in the 'Query results' table.

This is the same query as in our first (incorrect) solution. But the inputs are different, so the output looks much better.

Image by author.

The window function `last_value(... ignore nulls)` allows us to fill in the blanks, and the window function `lead` lets us retrieve the `end_dates` to re-form our date ranges.

All we need to do next is perform some cleaning to remove the magic values `-1` and `“”`:

Query results

Row	shop_id	article_id	pr.sell_price	price_and_pr_promotion_id	promotion_value	price_and_start_datetime	price_and_pr_end_datetime
1	dr29L9xF	yVrDmTw	null			2022-05-10 11:45:19 UTC	2022-05-10 11:53:49 UTC
			null			2022-05-10 11:53:49 UTC	2022-05-10 13:08:38 UTC
			null			2022-05-10 13:08:38 UTC	2022-05-10 13:08:40 UTC
			null			2022-05-10 13:08:40 UTC	2022-05-10 13:08:42 UTC
			null			2022-05-10 13:08:42 UTC	2022-05-10 13:08:45 UTC
			3.99			2022-05-10 13:08:45 UTC	2022-05-10 13:23:18 UTC
			3.99			2022-05-10 13:23:18 UTC	2022-05-11 00:00:00 UTC
			3.99	-38_P0evh	2.99	2022-05-11 00:00:00 UTC	2022-05-11 04:20:13 UTC
			3.99	-38_P0evh	2.99	2022-05-11 04:20:13 UTC	2022-05-11 05:03:04 UTC
			3.99	-38_P0evh	2.49	2022-05-11 05:03:04 UTC	2022-05-18 00:00:00 UTC
			3.99	-1.0	2.99	2022-05-18 00:00:00 UTC	2023-06-14 10:46:07 UTC
			3.99	-1.0	2.99	2023-06-14 10:46:07 UTC	2023-06-14 15:05:31 UTC
			3.99	2GLCMltlyw	2.99	2023-06-14 15:05:31 UTC	2023-06-16 11:29:38 UTC
			3.99	2GLCMltlyw	2.99	2023-06-16 11:29:38 UTC	2023-06-19 08:58:40 UTC
			3.99	2GLCMltlyw	2.99	2023-06-19 08:58:40 UTC	2023-06-22 13:31:46 UTC
			3.99	2GLCMltlyw	1.99	2023-06-22 13:31:46 UTC	2023-06-27 23:59:59.999000 U...
			3.99	2GLCMltlyw	1.99	2023-06-27 23:59:59.999000 U...	2024-02-14 00:59:59 UTC
			-1.0	2GLCMltlyw	-1.0	2023-06-27 23:59:59.999000 U...	2024-02-14 00:59:59 UTC

Image by author.

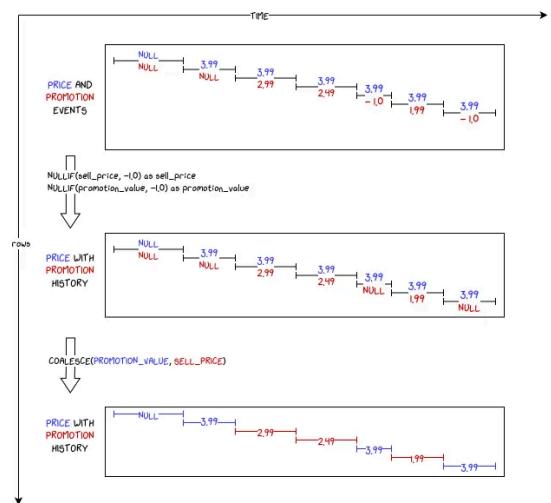
We perform one final query to compute the `price_with_promotion_included` and one final `unnest` since we want the final results to be flat (`shop_id` and

The screenshot shows a SQL query being run in a cloud-based SQL editor. The query is as follows:

```

1 SELECT
2   shop_id,
3   article_id,
4   sell_price,
5   promotion_id,
6   promotion_value,
7   (sell_price * promotion_value) AS price_with_promotion_included,
8   start_datetime,
9   end_datetime
10 FROM supermarket.price_with_promotion_history
11 LEFT JOIN UNNEST([price_and_promotion])
12
  
```

The results table has columns: Row, shop_id, article_id, sell_price, promotion_id, promotion_value, price_with_promotion_included, start_datetime, and end_datetime. The data shows various price and promotion events over time.



Removing the magic values with NULLIF and performing one final COALESCE to compute the price_with_promotion_included. Images by author.

And that's it, we are done!

Conclusion

We managed to compute what we wanted in four or five steps. I think this problem is a good example of something that looks simple at first sight but is rather tricky to do with SQL. This is the kind of advanced transformation where I recommend considering the trade-off between writing the logic in pure SQL vs. writing a User Defined Function (UDF) in another language (Python, Scala, Javascript, or whatever), if possible. The UDF has the advantage of making the code more generic and safe (with unit tests). However, it has the drawback of being generally less efficient and more costly compared to pure SQL.

As an exercise, if you want some practice, I recommend trying to rewrite this exercise yourself with one of the following variants:

- Rewrite this without using magic values.
(Hint: Add an *event_status* column that can take the value “ON” or “OFF”.)
- Rewrite this without using ARRAYS or STRUCT.
(Hint: It is quite doable, but the windows might need an extra “PARTITION BY shop_id, article_id”.)
- Rewrite it using your favorite dbt-like tool.
(I voluntarily did not use dbt here to keep this blog post accessible to most.)

- Try to write this in a fully generic way: The query should work with any number of columns in the input tables and keep working even if we add a new column in the input table without updating the query.
(Hint: I'm not sure this one is possible in pure-SQL because when I tried, I noticed that BigQuery seems unable to automatically cast NULLs into complex types when performing a union. This might be a nice challenge to do with [bigquery-dataframes](#) or [bigquery-frames](#), though).
- With DataFrames, it's probably even possible to make a fully generic transformation that would work automatically with any number of history tables.

Finally, if you want to learn more about how to work with arrays, you can start with [this awesome page](#) from BigQuery's documentation, and if you would like to see more articles on that specific subject from me, please leave a comment to let me know.

That's all for today. Thank you for reading! I hope you enjoyed reading this as much as I did writing it.

I leave you with the full query used in this article (rewritten with CTEs) and the full explanation graph. The code and the data samples [are available here](#). Enjoy ;)

Links

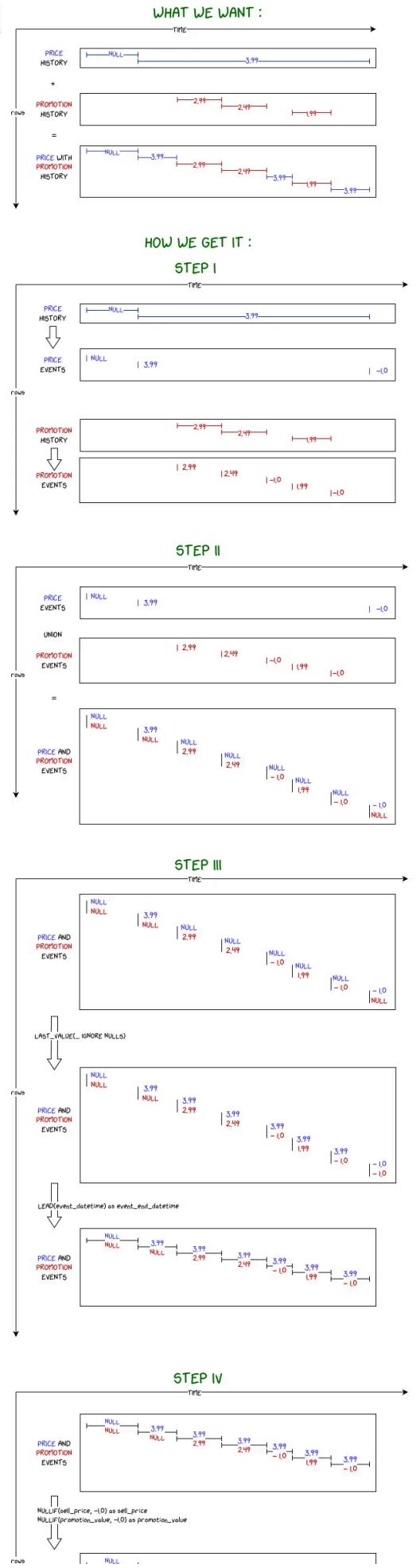
Here are all the links shared in this article:

- [python-bigquery-dataframes: Google's official project of Python-based DataFrames that runs on BigQuery.](#)
- [bigquery-frames: My very own project of Python-based DataFrames that runs on BigQuery.](#)
- [Google BigQuery's documentation: Working with arrays.](#)
- [Git gist with this article's code.](#)


```

    v2.full      RUN  SAVE QUERY  SHARE  SCHEDULE  MORE
1 WITH price_history_per_shop_per_article AS (
2     SELECT
3         shop_id,
4         article_id,
5         ARRAY_AGG(SELECT AS STRUCT T.* EXCEPT(shop_id, article_id)) ORDER BY start_datetime) as price_history
6     FROM `supermarket.sell_price_history` T
7     GROUP BY shop_id, article_id
8 )
9 , promotion_history_per_shop_per_article AS (
10    SELECT
11        shop_id,
12        article_id,
13        ARRAY_AGG(SELECT AS STRUCT T.* EXCEPT(shop_id, article_id)) ORDER BY start_datetime) as promotion_history
14     FROM `supermarket.promotion_history` T
15     GROUP BY shop_id, article_id
16 )
17 , price_and_promotion_history AS (
18    SELECT
19        shop_id,
20        article_id,
21        price_history,
22        promotion_history,
23     FROM price_history_per_shop_per_article AS price
24     LEFT JOIN promotion_history_per_shop_per_article AS promotion
25     USING (shop_id, article_id)
26 )
27 , price_and_promotion_history_with_next_event_datetime AS (
28    SELECT
29        shop_id,
30        article_id,
31        ARRAY(
32            SELECT AS STRUCT
33                *,
34                LAG(start_datetime) OVER (ORDER BY start_datetime DESC) as next_event_datetime
35            FROM UNNEST(price_history)
36            ORDER BY start_datetime
37        ) as price_history,
38        ARRAY(
39            SELECT AS STRUCT
40                *,
41                LAG(start_datetime) OVER (ORDER BY start_datetime DESC) as next_event_datetime
42            FROM UNNEST(promotion_history)
43            ORDER BY start_datetime
44        ) as promotion_history,
45     FROM price_and_promotion_history
46 )
47 , price_and_promotion_events AS (
48    SELECT
49        * EXCEPT(price_history, promotion_history),
50        ARRAY (
51            SELECT AS STRUCT
52                sell_price,
53                CAST(NULL as STRING) as promotion_id,
54                NULL as promotion_value,
55                start_datetime,
56                FROM UNNEST(price_history)
57                UNION ALL
58            SELECT AS STRUCT
59                -1 as sell_price,
60                CAST(NULL as STRING) as promotion_id,
61                NULL as promotion_value,
62                end_datetime,
63                FROM UNNEST(price_history)
64                WHERE end_datetime < next_event_datetime OR next_event_datetime IS NULL
65                UNION ALL
66            SELECT AS STRUCT
67                NULL as sell_price,
68                promotion_id,
69                promotion_value,
70                start_datetime,
71                FROM UNNEST(promotion_history)
72                UNION ALL
73            SELECT AS STRUCT
74                NULL as sell_price,
75                "" as promotion_id,
76                -1 as promotion_value,
77                end_datetime,
78                FROM UNNEST(promotion_history)
79                WHERE end_datetime < next_event_datetime OR next_event_datetime IS NULL
80                ORDER BY start_datetime
81        ) as events,
82        FROM price_and_promotion_history_with_next_event_datetime
83 )
84 , price_with_promotion_history_before_cleaning AS (
85    SELECT
86        * EXCEPT(events),
87        ARRAY(
88            SELECT AS STRUCT
89                LAST_VALUE(sell_price IGNORE NULLS) OVER (ORDER BY start_datetime) as sell_price,
90                LAST_VALUE(promotion_id IGNORE NULLS) OVER (ORDER BY start_datetime) as promotion_id,
91                LAST_VALUE(promotion_value IGNORE NULLS) OVER (ORDER BY start_datetime) as promotion_value,
92                start_datetime,
93                LEAD(start_datetime) OVER (ORDER BY start_datetime) as end_datetime
94                FROM UNNEST(events)
95                ORDER BY start_datetime
96            ) as price_and_promotion
97        FROM price_and_promotion_events
98 )
99 , price_with_promotion_history AS (
100    SELECT
101        * EXCEPT(price_and_promotion),
102        ARRAY(
103            SELECT AS STRUCT
104                NULLIF(sell_price, -1) as sell_price,
105                NULLIF(promotion_id, "") as promotion_id,
106                NULLIF(promotion_value, -1) as promotion_value,
107                start_datetime,
108                end_datetime
109                FROM UNNEST(price_and_promotion)
110                WHERE start_datetime < end_datetime
111        ) as price_and_promotion

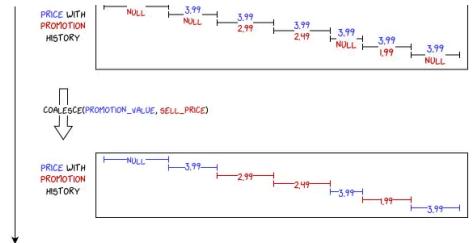
```



```

112 | FROM price_with_promotion_history_before_cleaning
113 )
114 SELECT
115   shop_id,
116   article_id,
117   sell_price,
118   promotion_id,
119   promotion_value,
120   COALESCE(promotion_value, sell_price) as price_with_promotion_included,
121   start_datetime,
122   end_datetime
123 FROM price_with_promotion_history
124 LEFT JOIN UNNEST(price_and_promotion)
125

```



The full query is 124 lines long. Images by author.

[Analytics Engineering](#)
[Bigquery](#)
[Sql](#)
[Data Engineering](#)
[Editors Pick](#)

More from the list: "SQL"

Curated by @reenum

SQL Fundam... in Stackad...

9 Advanced SQL Queries for Data Mastery

◆ 3 min read · Jan 20, 2024

SQL Fundam... in Stackad...

Top 15 Advanced SQL Queries

◆ 5 min read · Feb 24, 2024

SQL Fundam... in DevOps...

Top 10 Advanced SQL Queries

◆ 3 min read · Nov 7, 2023


Les
SQI



11 m

[View list](#)


Written by Furcy Pin

871 Followers · Writer for Learning SQL

[Follow](#)


[Available for freelance work] Data Engineer, Data Plumber, Data Librarian, Data Smithy.

More from Furcy Pin and Learning SQL





Furcy Pin in Towards Data Science



Sarang S. Babu in Learning SQL

2003–2023: A Brief History of Big Data

Summing up 20 years of history of Hadoop and everything related

18 min read · Nov 9, 2022



342



4



+



...



Zach Quinn in Learning SQL

How I Reduced My Query's Run Time From 30 Min. To 30 Sec. In 1...

The query optimization steps a senior data engineer took to reduce the process time of ...

· 6 min read · Mar 13, 2024



413



16



+



...

12 Tips for Optimizing SQL Queries for Faster Performance

Ways to Optimize SQL Queries

8 min read · Mar 6, 2023



120



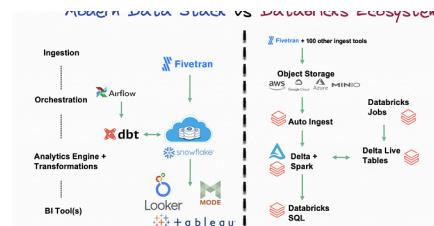
2



+



...



Furcy Pin in Towards Data Science

Modern Data Stack: which place for Spark ?

One year ago, some were already predicting that dbt will one day become bigger than...

7 min read · Jan 25, 2022



692



15



+

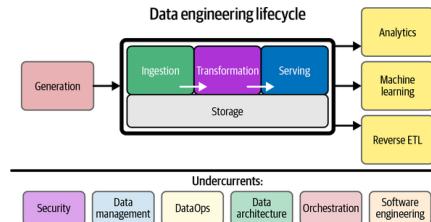


...

[See all from Furcy Pin](#)

[See all from Learning SQL](#)

Recommended from Medium



Mudra Patel

Data Engineering concepts: Part 1, Data Modeling

What is Data Engineering?

5 min read · Feb 11, 2024

332 3

Rosaria Silipo in Low Code for Data Science

Is Data Science dead?

In the last six months I have heard this question thousands of time: "Is data science...

6 min read · Mar 11, 2024

1.7K 36

Lists



The New Chatbots: ChatGPT, Bard, and Beyond

12 stories · 360 saves



ChatGPT

21 stories · 586 saves



Natural Language Processing

1386 stories · 882 saves



Business

41 stories · 96 saves



Kallol Mazumdar in ILLUMINATION

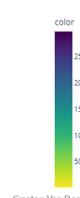
I Went on the Dark Web and Instantly Regretted It

Accessing the forbidden parts of the World Wide Web, only to realize the depravity of...

8 min read · Mar 13, 2024

15.3K 302

Number of Job Openings by Location
Non-remote



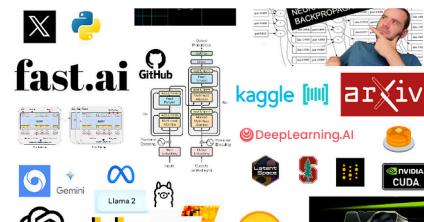
Luna (Van) Doan

Data Job Market 2024: Insights You Need to Boost Your Career

Level up your data career with the insights from 3K+ LinkedIn jobs posts

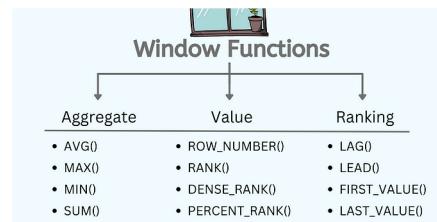
5 min read · Feb 19, 2024

1.1K 21



Benedict Neo in bitgrit Data Science Publication

Roadmap to Learn AI in 2024



Hasan Hüseyin Coşgun

A free curriculum for hackers and programmers to learn AI

11 min read · Mar 10, 2024

10.3K

111



...

Learn SQL Window Functions With Queries

Time to discover the world of Window Functions with examples; lead(), lag(),...

7 min read · Nov 3, 2023

88



...

See more recommendations