

SPRINT 4: CREACIÓN DE BASES DE DATOS

NIVEL 1

Descarga los archivos CSV, estúdalos y diseña una base de datos con un esquema de estrella que contenga, al menos, 4 tablas.

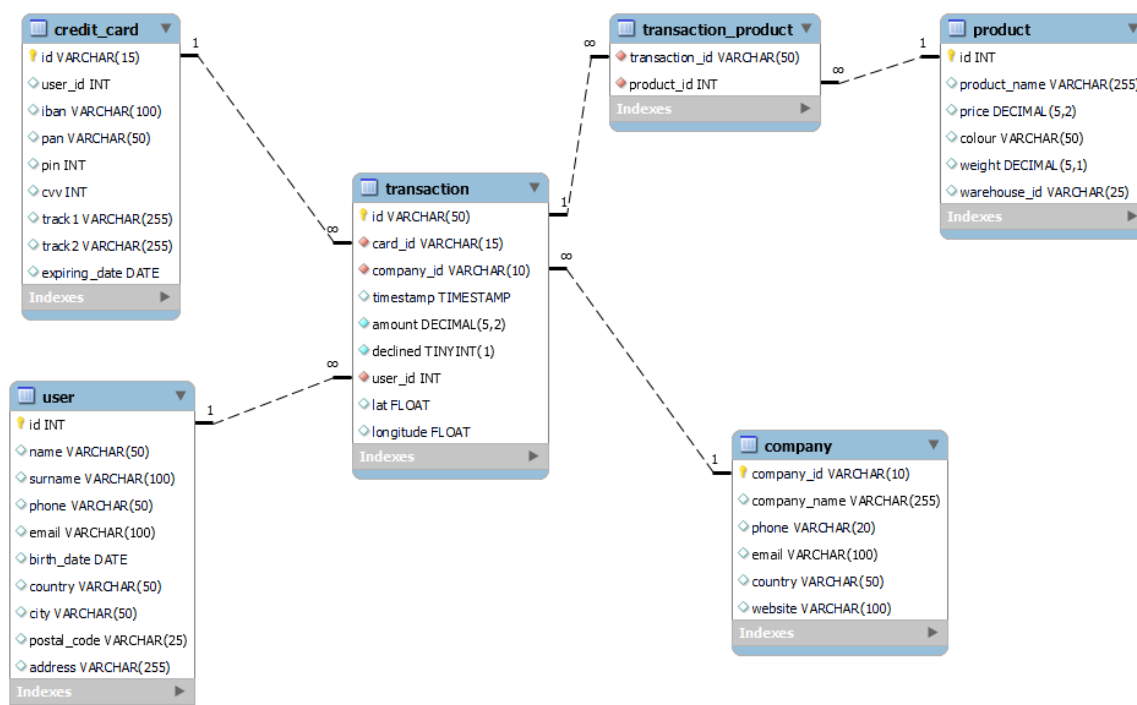
Los archivos CSV contenían registros sobre las transacciones, las empresas que participan en dichas transacciones, las tarjetas de crédito utilizadas, los usuarios que las realizaron (de tres países: Reino Unido, Estados Unidos y Canadá) y los productos vendidos. Se ha creado una base de datos llamada *transactionsdb* y las tablas correspondientes para cargar los datos de los archivos CSV (véase script estructura_db.sql para más detalles).

Los registros de transacciones contenían varios productos vendidos por transacción en una única columna, cosa que complicaba filtrar los datos de un único producto de una lista string. Para solucionarlo, he creado una tabla nueva (*transaction_product*) para los ids de transacciones y los ids de productos comprados. A continuación, he transformado la columna de *product_ids* de *transactions.csv* para separar cada producto de una transacción en una fila (transformación realizada con Power Query, ya que no he conseguido encontrar la manera de hacerlo con MySQL Workbench).

*

Los tipos de datos de las variables se escogieron en función del formato que tenían los datos en los archivos CSV. Después de cargar los datos de los archivos (mediante Table Data Import Wizard, ya que no he conseguido solucionar los errores para poder importarlos con código), he realizado algunas modificaciones de datatype (fechas de VARCHAR a DATE, precio de producto de VARCHAR a DECIMAL; para más detalles, véase script modificaciones_db.sql). **

En este modelo tenemos la tabla de hechos *transaction*, que conecta con las tablas dimensión *credit_card* ***, *user* y *company* con una relación N:1 (transaction N:1 dimensión). Para relacionar *transaction* con dimensión *product* tenemos la tabla intermedia *transaction_product*, que conecta N:1 con *transaction* y N:1 con *product*.



*** NOTA 1:** Cada vez que se necesite actualizar la tabla de hechos con nuevas transacciones, este modelo no sería eficiente, ya que se necesitaría separar en Power Query los productos por fila; lo ideal sería encontrar la manera de hacer esta transformación en el propio MySQL y hacer que se actualice la tabla intermedia transaction_product con nuevos registros.

**** NOTA 2:** En este modelo he actualizado (sobrescrito) los datos existentes en las columnas por los mismos datos, pero en otro formato para poder trabajar con fechas en DATE o precio en DECIMAL; modificar directamente el datatype de las columnas existentes podría resultar problemático, ya que cada actualización en estas tablas podría generar errores y complicaciones adicionales.

Por este motivo, considero que sería mejor crear una nueva columna para el datatype deseado y mantener intacto el formato de los datos de origen; así, podríamos definir un Trigger que se active cada vez que haya nuevos registros en las columnas originales. Este Trigger actualizaría automáticamente los valores en las nuevas columnas con el formato de datos deseado, utilizando los datos de las columnas originales como referencia.

***** NOTA 3:** La tabla credit_card contiene el identificador de los usuarios propietarios de las tarjetas de crédito. No he establecido una relación entre las tablas credit_card y user ya que se formaría una relación cíclica en mi base de datos.

Ejercicio 1. Realiza una subconsulta que muestre a todos los usuarios con más de 30 transacciones utilizando al menos 2 tablas.

Primero de todo debemos conocer qué usuarios realizaron más de 30 transacciones, mediante una subconsulta que devuelva los id del usuario y la cantidad de transacciones (filtradas por > 30).

Una vez sabemos los identificadores de estos usuarios, los relacionamos con la tabla user para conocer sus datos personales.

```
SELECT u.*, user_trans.cant_trans
FROM user u
JOIN (SELECT user_id, COUNT(id) AS cant_trans
      FROM transaction
      GROUP BY user_id
      HAVING cant_trans > 30) user_trans
ON u.id=user_trans.user_id;
```

Obtenemos un total de cuatro usuarios que cumplen el criterio establecido: Lynn Riddle (id 92), Ocean Nelson (id 267), Hedwig Gilbert (id 272) y Kenyon Hartman (id 275). Sus datos personales y la cantidad de transacciones realizadas por cada uno de ellos se muestran en la tabla a continuación.

id	name	surname	phone	email	birth_date	country	city	postal_code	address	cant_trans
92	Lynn	Riddle	1-387-885-4057	vitae.aliquet@outlook.edu	1984-09-21	United States	Bozeman	61871	P.O. Box 712, 7907 Est St.	39
267	Ocean	Nelson	079-481-2745	aenean@yahoo.com	1991-12-26	Canada	Charlottetown	85X 3P4	Ap #732-8357 Pede, Rd.	52
272	Hedwig	Gilbert	064-204-8788	sem.eget@icloud.edu	1991-04-16	Canada	Tuktoyaktuk	Q4C 3G7	P.O. Box 496, 5145 Sapient Road	76
275	Kenyon	Hartman	082-871-7248	convallis.ante.lectus@yahoo.com	1982-08-03	Canada	Richmond	R8H 2K2	8564 Facilisi. St.	48

Ejercicio 2. Muestra el promedio de las transacciones por IBAN de las tarjetas de crédito en la compañía Donec Ltd. utilizando al menos 2 tablas.

En la tabla de transacciones tenemos id_company pero no su nombre, por lo que tenemos que recurrir a la tabla company mediante una subconsulta para conocer el identificador de la compañía Donec Ltd.

Una vez sabemos el id, podemos calcular el promedio de las transacciones de esta compañía en función de las tarjetas de crédito utilizadas (GROUP BY card_id).

Al tener ya el promedio por tarjeta, solo falta relacionarlo con la tabla credit_card para tener el número de IBAN asociado a la tarjeta.

```
WITH avg_trans_card AS
(
  SELECT card_id, ROUND(AVG(amount),2) AS avg_amount
  FROM transaction
  WHERE company_id = (SELECT company_id
                      FROM company
                      WHERE company_name='Donec Ltd')
  GROUP BY card_id
)
SELECT iban, atc.avg_amount AS promedio_transaccion
FROM credit_card cc
JOIN avg_trans_card atc
ON cc.id=atc.card_id;
```

En este caso, la compañía Donec Ltd. tenía tan solo dos transacciones y las dos se realizaron con la misma tarjeta de crédito. Por tanto, como resultado obtenemos un número iban (mostrado en la tabla) con un promedio de 203.72 euros.

iban	promedio_transaccion
PT87806228135092429456346	203.72

NIVEL 2

Crea una nueva tabla que refleje el estado de las tarjetas de crédito basado en si las últimas tres transacciones fueron declinadas.

Para crear una tabla que muestre el estado de las tarjetas, debemos conseguir las últimas tres transacciones realizadas con las tarjetas. Para ello, seleccionamos información relevante, como es el id de la tarjeta, timestamp (fecha y hora de la transacción) y declined (0: aceptada, 1: rechazada).

Dado que el número de transacciones por tarjeta varían, debemos aplicarle un “índice” interno a cada fila de transacción de una determinada tarjeta con la función ROW_NUMBER ().

Al usar PARTITION BY card_id, estamos indicando a la función que separe las asignaciones de números secuenciales para cada tarjeta individualmente. Es decir, cuando terminan las transacciones de una tarjeta y comienzan las de otra, el índice se reinicia a 1 para la nueva tarjeta. Esto asegura que cada tarjeta tenga su propio conjunto de índices únicos para sus transacciones.

Además, al especificar ORDER BY timestamp DESC, estamos ordenando las transacciones dentro de cada grupo de tarjeta en orden descendente según su fecha de realización. Esto significa que la asignación de índices se realiza de las transacciones más recientes a las más antiguas.

Así, obtenemos 587 transacciones, cada una con un índice asignado en función de la tarjeta. Este índice comienza desde 1 para la transacción más reciente de cada tarjeta y aumenta secuencialmente hacia las transacciones más antiguas.

Para considerar si una tarjeta está activa o no, hay que considerar las 3 transacciones más recientes de cada tarjeta ($\text{row_transaction} \leq 3$); por tanto, en la tabla `card_status` únicamente incluiremos aquellas tarjetas que tengan al menos 3 transacciones registradas ($\text{COUNT card_id} = 3$). La razón detrás de este requisito es asegurarnos de tener suficiente información para tomar una decisión precisa sobre el estado de la tarjeta.

Esto quiere decir que aquellas tarjetas que tengan un total de 1 o 2 transacciones no cumplen con el criterio por el siguiente motivo:

- Dos transacciones: aunque ambas transacciones sean rechazadas, no podemos concluir que la tarjeta está inactiva. Faltaría saber si la tercera transacción que se realizará con esa tarjeta sería aceptada o no.
- Una transacción: aunque la única transacción esté rechazada, desconocemos si es debido a que la tarjeta está inactiva u otro motivo. Faltarían dos transacciones más para poder evaluar con mayor precisión el estado de la tarjeta.

En la tabla creada tendremos dos columnas, la primera con el identificador de la tarjeta y la segunda con el estado de la tarjeta: 'operativa' o 'inoperativa' (asignamos el estado con CASE).

- Si las tres transacciones son rechazadas, la suma de declined sería 3, por tanto 'inoperativa'
- Si hay una o dos transacciones rechazadas de las tres a considerar, la suma de declined sería ≤ 2 , por tanto 'operativa'

Todas aquellas tarjetas que actualmente tengan una o dos transacciones se podrán añadir más adelante a la tabla `card_status` si obtenemos registros de dos o una transacción más, respectivamente.

```
CREATE TABLE IF NOT EXISTS card_status AS
(WITH trans_card AS
(
SELECT card_id,
timestamp,
declined,
ROW_NUMBER() OVER (PARTITION BY card_id ORDER BY timestamp DESC) AS row_transaction
FROM transaction
)
SELECT card_id AS id_tarjeta,
CASE
WHEN SUM(declined) <= 2 THEN 'operativa'
ELSE 'inoperativa'
END AS estado_tarjeta
FROM trans_card
WHERE row_transaction <= 3
GROUP BY card_id
HAVING COUNT(card_id) = 3
);
```

Ejercicio 1. ¿Cuántas tarjetas están activas?

Una vez tenemos creada la tabla con el estado de las tarjetas, contamos cuántas están operativas.

```
SELECT COUNT(id_tarjeta) AS 'cantidad tarjetas activas'
FROM card_status
WHERE estado_tarjeta='operativa';
```

Como resultado, obtenemos 9 tarjeta activas. Cabe destacar que en la tabla card_status tenemos un total de nueve tarjetas, que son las únicas que cumplen con los criterios establecidos. Por tanto, todas las tarjetas presentes en card_status están activas.

cantidad tarjetas
activas

9

NIVEL 3

Ejercicio 1. Necesitamos conocer el número de veces que se ha vendido cada producto.

Como tenemos interés en conocer el número de veces que se vendió cada producto, las transacciones rechazadas no deberían considerarse para el conteo. Es decir, considero que una transacción rechazada no representa una venta exitosa del producto.

Sin embargo, tener información sobre la cantidad de veces que fue rechazada la venta de un producto puede ser útil para otros análisis. Por ejemplo, podría servir para hacer un seguimiento de rechazos, para optimizar los procesos de venta, para estudiar las tendencias de venta, entre otros. Por tanto, en mi informe he querido incluir tanto la cantidad de transacciones aceptadas como la cantidad de rechazadas para cada producto.

Con la tabla transaction_product podemos saber la cantidad total de transacciones que se han registrado para un determinado producto; pero, tal como he diseñado la tabla, no podemos distinguir si esas transacciones fueron aceptadas o rechazadas. Para ello, tenemos que recurrir a la tabla transaction mediante una JOIN.

Para calcular la cantidad de transacciones, utilizo la función SUM con la condición descrita en CASE.

- En la columna ventas.cant_acept aplico la condición de que todas las transacciones aceptadas (declined 0) tengan valor 1 y las rechazadas 0, por lo que sumando las veces que aparece cada producto obtenemos el total de ventas exitosas.
- En la columna ventas.cant_rech aplico la condición de que todas las transacciones rechazadas (declined 1) tengan valor 1 y las aceptadas 0, por lo que la suma daría el total de ventas rechazadas para cada producto.

Una vez tenemos esta información, podemos relacionarla con la tabla product para conocer los detalles de cada producto mediante LEFT JOIN, ya que hay productos sin ninguna transacción registrada. Para evitar tener valores NULL para los productos sin ventas, aplico la función IFNULL ().

```

WITH ventas AS
(
  SELECT product_id,
         SUM(CASE WHEN declined = 0 THEN 1 ELSE 0 END) as cant_acept,
         SUM(CASE WHEN declined = 1 THEN 1 ELSE 0 END) as cant_rech
  FROM transaction_product tp
  JOIN transaction t
  ON tp.transaction_id = t.id
  GROUP BY product_id
  ORDER BY product_id
)
SELECT p.*,
       IFNULL(v.cant_acept,0) AS cantidad_vendida,
       IFNULL(v.cant_rech,0) AS cantidad_rechazada,
       IFNULL(v.cant_acept,0) + IFNULL(v.cant_rech,0) AS cantidad_total
FROM product p
LEFT JOIN ventas v
ON v.product_id=p.id;

```

Con esta consulta, obtenemos una tabla con 100 registros, lo que corresponde a la cantidad total de productos registrados en la tabla product. Al examinar la tabla, vemos que algunos productos tienen cero transacciones registradas. Por ejemplo, el producto con el ID 4, denominado 'Warden South Duel', no tiene transacciones registradas. Esto significa que no se ha realizado ninguna transacción relacionada con este producto hasta el momento.

id	product_name	price	colour	weight	warehouse_id	cantidad_vendida	cantidad_rechazada	cantidad_total
1	Direwolf Stannis	161.11	#7c7c7c	1.0	WH-4	51	10	61
2	Tarly Stark	9.24	#919191	2.0	WH-3	56	9	65
3	duel tourney Lannister	171.13	#d8d8d8	1.5	WH-2	43	8	51
4	warden south duel	71.89	#111111	3.0	WH-1	0	0	0
5	skywalker ewok	171.22	#dbdbdb	3.2	WH-0	42	7	49