Desafíos en el Manejo de Fechas en SQL:

Un Caso Práctico con Datos de Empleados

Introducción

Como parte de mi desarrollo en SQL, he estado trabajando con una base de datos que contiene información valiosa sobre empleados. Este conjunto de datos incluye columnas criticas como salario, fechas de cumpleaños, fecha de inicio y fecha de finalización. Durante este proceso, enfrente varios desafíos en la gestión de las fechas y los datos en general, y me gustaría compartir como los aborde.

1. Estandarización de nombres de columnas

Una de las primeras tareas que lleve a cabo fue la limpieza de las columnas para asegurar la coherencia y facilitar futuras consultas. Esto es crucial en la gestión de bases de datos, ya que nombres confusos pueden llevar a errores en las consultas. Realice los siguientes cambios en la tabla:

```
    ALTER TABLE limpieza CHANGE COLUMN `Id?empleado` id_emp varchar(20);
    ALTER TABLE limpieza CHANGE COLUMN `Name` name varchar(20);
    ALTER TABLE limpieza CHANGE COLUMN `Apellido` last_name VARCHAR(20);
    ALTER TABLE limpieza CHANGE COLUMN `género` gender VARCHAR(20);
```

Al renombrar las columnas, mejoré la legibilidad de la base de datos, facilitando la interacción con los datos y reduciendo el riesgo de cometer errores en futuras consultas.

2. Manejo de Fechas: Birth_date y Start_date

Una vez que los nombres de las columnas estaban estandarizados, me enfoque en trabajar con las columnas de fecha de nacimiento (birth_date) y fecha de inicio (start_date). Ambas columnas contenían valores en diferentes formatos, y en algunos casos, contenían valores nulos.

El primer paso fue asegurarme de que los datos estuvieran en un formato de fecha adecuado. Utilice STR_TO_DATE para transformar las fechas que estaban en formato texto a un formato estándar:

En esta primera parte, utilizamos una consulta SELECT para transformar los valores de birth_date en un formato de fecha adecuado.

- **SUBSTRING_INDEX(birth_date, '/', 1)**: Este fragmento extrae el primer componente de la fecha (el día o el mes) al dividir la cadena birth_date en base al carácter /.
- CASE: Aquí se utiliza una declaración CASE para determinar el formato de fecha. Si el primer componente (día o mes) es mayor que 12, asumimos que el formato es DD/MM/YYYY y lo convertimos usando STR_TO_DATE con el formato correspondiente. Si no, se asume el formato MM/DD/YYYY.
- STR_TO_DATE: Esta función convierte la cadena en un objeto de fecha, utilizando el formato especificado.

El resultado de esta consulta muestra cómo se verían los datos formateados, pero no se guarda en la base de datos en este paso.

```
ALTER TABLE limpieza ADD COLUMN Birth_date_formatted DATE;

UPDATE limpieza

UPDATE limpieza

SET Birth_date_formatted = CASE

WHEN SUBSTRING_INDEX(birth_date, '/', 1) > 12 THEN STR_TO_DATE(birth_date, '%d/%m/%Y')

ELSE STR_TO_DATE(birth_date, '%m/%d/%Y')

END;
```

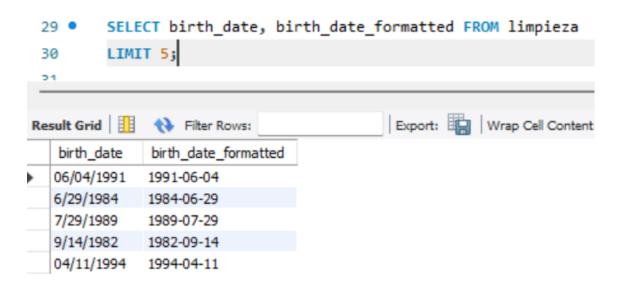
La segunda parte de la consulta es esencial para la estructuración de nuestra base de datos. Aquí, utilizo la instrucción ALTER TABLE para agregar una nueva columna llamada Birth_date_formatted a la tabla limpieza.

• **ADD COLUMN**: Esto permite almacenar las fechas formateadas en un nuevo campo, lo cual es crucial para mantener la integridad de los datos y facilitar futuras consultas o análisis.

Finalmente, utilizo la instrucción UPDATE para llenar la columna recién creada con los valores convertidos.

- **SET**: Aquí se aplica nuevamente la lógica de CASE para determinar el formato correcto y actualizar la columna Birth_date_formatted con los valores correspondientes.
- La utilización de la misma lógica de formateo asegura que todos los registros en la nueva columna sean consistentes y correctos, mejorando la calidad de los datos en la base de datos.

Así quedaría antes y después de la consulta:



'LIMIT' se utiliza para limitar la consulta, en este caso para mostrarme los primeros 5 registros.

En cuanto a la columna 'Start_date' hice exactamente lo mismo, ya que presentaba las mismas inconsistencias:

```
43 ● UPDATE limpieza

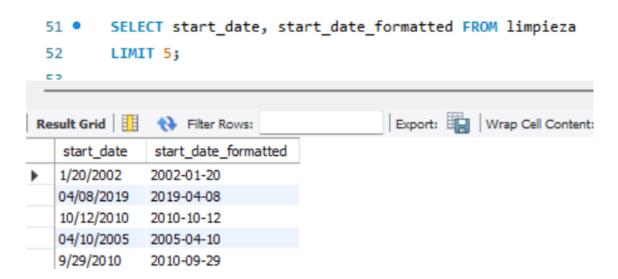
44 ⊖ SET star_date_formatted = CASE

45 WHEN SUBSTRING_INDEX(star_date, '/', 1) > 12 THEN STR_TO_DATE(star_date, '%d/%m/%Y')

ELSE STR_TO_DATE(star_date, '%m/%d/%Y')

END;
```

Quedando antes y después de la siguiente manera:



Esto es fundamental, ya que trabajar con formatos inconsistentes puede llevar a errores en el análisis de datos. Al asegurarnos de que todas las fechas estén en el mismo formato, se puede llevar a cabo un análisis más eficaz y preciso.

Otro desafio importante fue la conversión de fechas en un formato de texto a un formato de fecha que SQL pudiera reconocer. Inicialmente, los registros de 'finish_date' estaban en el formato 2029-10-29 06:09:38 UTC. Para solucionarlo, utilice la función STR_TO_DATE, que me permitió convertir estos valores de texto en un formato de fecha estándar:

```
UPDATE limpieza

SET finish_date = STR_TO_DATE(finish_date, '%Y-%m-%d %H:%i:%s UTC')

WHERE finish_date IS NOT NULL AND finish_date != '';
```

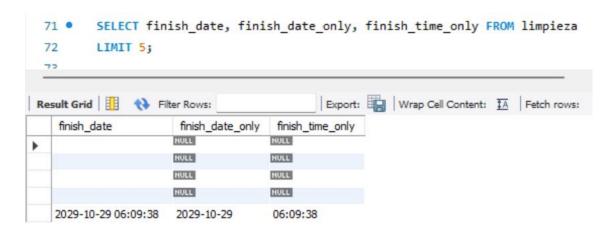
Sin embargo, encontré errores debido a valores vacíos ya que esos registros hacen referencias a empleados que todavía siguen vigentes en la empresa. Ajustando la consulta para manejar estos casos, pude asegurarme de que los registros vacíos se convirtieran en NULL.

3. Separación de Fecha y Hora

Con las fechas correctamente formateadas, decidí dividir la columna 'finish_date' en dos: una para la fecha y otra para la hora. Esto mejoro la legibilidad y facilito análisis posteriores. La consulta fue:

Esta separación permite realizar análisis mas específicos, como calcular la duración de empleo de manera mas precisa.

Quedando finalmente así:



4. Manejo de Salarios en formato incorrect

Adicionalmente, la columna de salario presentaba datos en un formato que incluía símbolo de dólar y comas, lo que complica su análisis. Utilice una consulta para limpiar estos datos:

```
94 • UPDATE limpieza

95 SET salary = CAST(REPLACE(REPLACE(salary, '$', ''), ',', '') AS DECIMAL(10, 2))

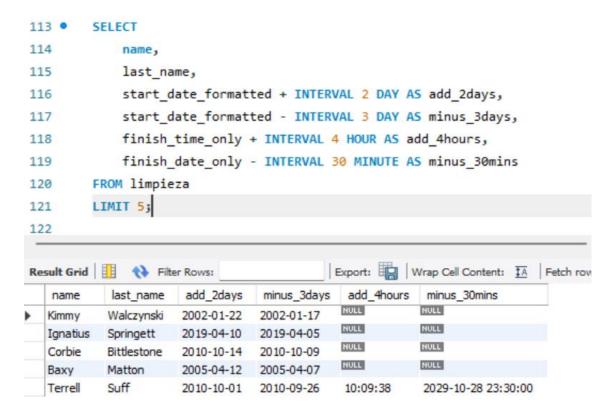
96 WHERE salary IS NOT NULL;
```

Esta limpieza fue crucial para asegurar que los salarios se pudieran utilizar en cálculos y reportes sin errores. Quite el símbolo de dólar y quite la coma. Además de eso, deje la columna 'salary' en formato DECIMAL ya que la misma estaba en formato texto.

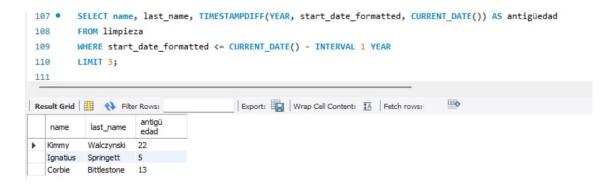
5. Uso de operadores de comparación y funciones SQL avanzadas

Aparte de la limpieza de datos, también implemente operadores de comparación y funciones avanzadas. Utilice el operador INTERVAL.

La función INTERVAL se utiliza para manejar intervalos de fecha y hora sumando y restando intervalos como "3 días", "5 horas", "45 minutos".



Por ejemplo, la utilice para calcular la antigüedad de los empleados.



TIMESTAMPDIFF(YEAR, start_date, CURRENT_DATE()) AS antigüedad: Calcula la diferencia en años entre la fecha de start_date y la fecha actual, y la renombra como antigüedad.

También explore el uso de la función EXTRACT (Extrae un componente especifico como año, mes, día, hora o minuto), que permite obtener partes especificas de una fecha. Esto es particularmente útil para análisis de rendimiento basados en años o meses:

```
123 •
         SELECT
124
              EXTRACT(YEAR FROM start_date_formatted) AS año_inicio,
125
              COUNT(*) AS cantidad empleados
         FROM limpieza
126
127
         GROUP BY año inicio
         HAVING COUNT(*) > 10
128
                                              Export: Wrap Cell Content: $\frac{1}{4}
Result Grid
               Filter Rows:
              cantidad empleados
   año inicio
  2002
             434
  2019
             408
  2010
             423
  2005
             432
  2018
             448
```

En esta consulta EXTRACT(YEAR FROM start_date_formatted) AS año_inicio extra el año de la columna start_date_formatted y lo renombra como año_inicio, COUNT(*) as cantidad_empleados cuenta el numero total de empleados para cada año, GROUP BY año_inicio agrupa los resultados por el año extraído y HAVING COUNT(*) > 10 filtra los grupos para mostrar solo aquellos años con mas de 10 empleados.

Conclusión

En el proceso de trabajar con la base de datos de empleados, he tenido la oportunidad de enfrentar diversos desafíos que me han permitido profundizar mis conocimientos en SQL y demostrar mis habilidades en la manipulación y limpieza de datos. A lo largo de este artículo, hemos explorado no solo cómo formatear fechas, sino también cómo abordar la calidad de los datos desde un enfoque integral.

Desde la estandarización de los nombres de las columnas hasta la conversión de formatos de fecha, cada paso ha sido fundamental para garantizar que la información se presente de manera consistente y efectiva. He aprendido que, más allá de realizar consultas básicas, es crucial entender la lógica detrás de cada transformación y cómo estas pueden impactar el análisis posterior. Por ejemplo, el uso de

CASE para identificar y convertir formatos de fecha no solo optimiza la limpieza de datos, sino que también prepara el terreno para análisis más complejos y significativos en el futuro.

Asimismo, trabajar con columnas que contienen valores nulos y vacíos ha sido un recordatorio constante de que, en el análisis de datos, la calidad es más importante que la cantidad. Aprender a manejar estos casos con elegancia, asegurando que la integridad de los datos se mantenga, es un signo de un analista consciente y profesional.

Además, he integrado conceptos más avanzados, como operadores de comparación y funciones como INTERVAL y EXTRACT, que han enriquecido aún más mi entendimiento de SQL. Estos elementos son cruciales para realizar análisis más sofisticados, permitiendo insights que pueden informar decisiones estratégicas en un entorno empresarial.

Finalmente, la capacidad de comunicar efectivamente lo aprendido a través de un artículo como este es esencial. La documentación clara y convincente de mis procesos no solo demuestra mis habilidades técnicas, sino que también muestra mi compromiso con la mejora continua y el aprendizaje en el ámbito de los datos. Cada desafío superado y cada consulta escrita son un paso más hacia la maestría en SQL, y estoy emocionado por aplicar todo este conocimiento en futuros proyectos y análisis.

A medida que continuo mi camino en el análisis de datos, sé que la práctica constante y la disposición para enfrentar nuevos desafíos serán clave para mi crecimiento profesional. Espero que este artículo no solo sea un reflejo de mis habilidades actuales, sino también un indicativo de mi ambición y determinación para seguir aprendiendo y evolucionando en el fascinante mundo de los datos.