

SISTEMAS INFORMÁTICOS I

PRÁCTICA 2

AUTORES

CARLOS GARCÍA SANTA
EDUARDO JUNOY ORTEGA

SISTEMAS INFORMÁTICOS
GRUPO1321: PAREJA 05

INGENIERÍA INFORMÁTICA
ESCUELA POLITÉCNICA SUPERIOR
UNIVERSIDAD AUTÓNOMA DE MADRID



19/10/2023

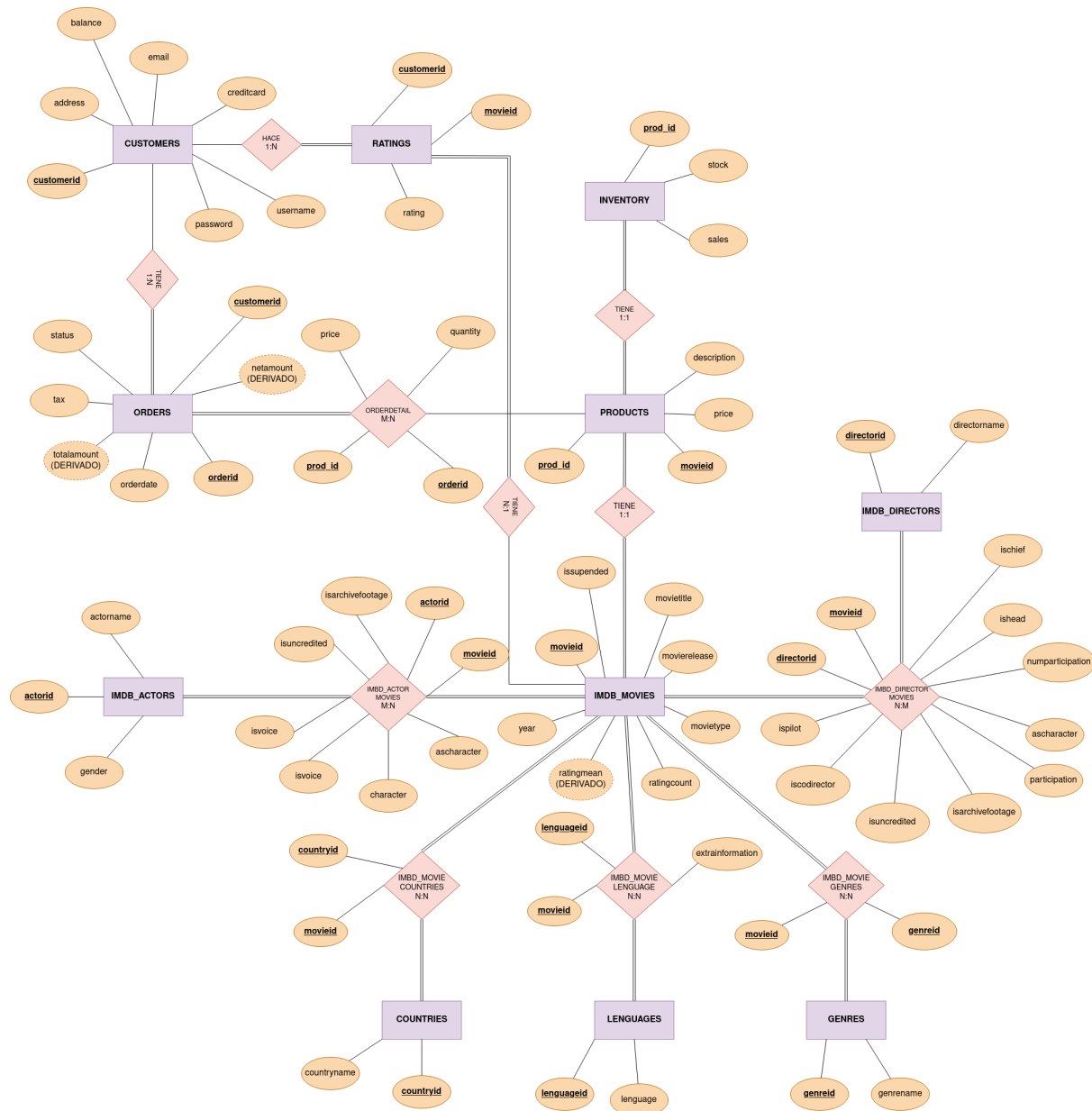
Índice

Índice.....	2
Introducción.....	2
Parte 1: Diseño de la BD.....	3
1. Adición de Claves Foráneas:.....	3
2. Normalización de Datos y Creación de Nuevas Tablas:.....	4
3. Tareas.....	4
4. Funciones.....	5
Parte 2: Optimización.....	8
k) Estudio del impacto de un índice.....	8
l) Estudio del impacto de cambiar la forma de realizar una consulta:.....	11
m) Estudio del impacto de la generación de estadísticas:.....	12
Cuestiones.....	17

Introducción

Esta memoria documenta los conocimientos desarrollados y aplicados en la primera práctica. Esta se ha centrado en explorar las prestaciones y posibilidades avanzadas que provee un sistema gestor de base de datos a través de: el manejo y gestión de bases de datos, programación de consultas, el uso de funciones y la optimización de consultas con el uso de índices y comandos.

Parte 1: Diseño de la BD



1. Adición de Claves Foráneas:

Tablas imdb_actormovies, inventory, orderdetail, y orders:

Se añadieron restricciones de clave foránea (FOREIGN KEY) a estas tablas para establecer relaciones explícitas con las tablas referenciadas (imdb_actors, imdb_movies, products, orders, customers). Esto asegura la integridad referencial, es decir, que no se pueden insertar datos en estas tablas que no tengan correspondencia en las tablas a las que hacen referencia.

2. Normalización de Datos y Creación de Nuevas Tablas:

Creación de las tablas **countries**, **languages**, y **genres**:

Estas tablas se crean para almacenar información única sobre países, idiomas y géneros, respectivamente. Cada una tiene un id como clave primaria y un campo para el nombre, asegurando que los nombres sean únicos y no nulos.

- Se extraen datos de las columnas **country**, **language**, y **genre** de las tablas **imdb_moviecountries**, **imdb_movi_languages**, y **imdb_moviegenres**, respectivamente, y se insertan en las nuevas tablas creadas.
- Se añaden nuevas columnas (**countryid**, **languageid**, **genreid**) a las tablas **imdb_moviecountries**, **imdb_movi_languages**, y **imdb_moviegenres**.
- Se establecen relaciones de clave foránea con las nuevas tablas (**countries**, **languages**, **genres**).
- Se actualizan las tablas originales para que el id correspondiente reemplace el nombre del país, idioma o género.
- Las columnas originales **country**, **language**, y **genre** se eliminan de las tablas **imdb_moviecountries**, **imdb_movi_languages**, y **imdb_moviegenres**.

Estos cambios mejoran la estructura de la base de datos. La adición de claves foráneas refuerza la integridad referencial, mientras que la normalización reduce la redundancia de datos y mejora la eficiencia de la base de datos.

3. Tareas

Modificación de la Tabla **customers**

Se añade un nuevo campo **balance** de tipo **NUMERIC** a la tabla **customers**, este campo está destinado a almacenar el saldo de cada cliente.

Creación de la Tabla **ratings**

Se crea una nueva tabla **ratings** para almacenar las valoraciones que los clientes dan a las películas, la tabla tiene los campos **movieid**, **customerid**, y **rating**. Se establece una restricción **CHECK** en **rating** para asegurar que los valores estén entre 1 y 10, y se añaden claves foráneas (**fk_movie**, **fk_customerid**) para vincular **movieid** y **customerid** con las tablas **imdb_movies** y **customers**, respectivamente, además se establece una restricción **UNIQUE** en la combinación (**movieid**, **customerid**) para evitar que un mismo cliente valore dos veces la misma película.

Modificaciones en la Tabla **imdb_movies**

Se añaden dos nuevos campos a la tabla **imdb_movies**: **ratingmean** (valoración media) y **ratingcount** (número de valoraciones), estos campos están

destinados a almacenar la valoración media y el total de valoraciones recibidas para cada película.

Cambio en la Tabla customers

Se altera el tipo de dato del campo **password** a CHARACTER VARYING(96).

Creación de Procedimiento Almacenado

Se crea un procedimiento almacenado llamado **setCustomersBalance** que inicializa el campo balance de la tabla customers con un valor aleatorio entre 0 y un valor máximo especificado (initialBalance).

4. Funciones

setPrice.sql: La consulta setPrice.sql actualiza la columna price en la tabla orderdetail, utilizando un join para unir las tablas orders y products de donde se obtienen las fechas de los pedidos y los precios actuales, con ello se calcula el nuevo precio incrementándolo un 2% anualmente desde la fecha del pedido, esto se logra con `POW(1.02, EXTRACT(YEAR FROM age(o.orderdate)))`, que aplica el incremento basado en los años transcurridos desde la fecha del pedido. El resultado se multiplica por el precio actual del producto y se redondea a dos decimales, esta actualización afecta solo a las tuplas con fechas de pedido hasta la fecha actual.

Antes de ejecutar setPrice
hay que cambiarla

	123orderid	123prod id	123price	123quantity
1	62,397	1,147	[NULL]	1
2	62,397	5,981	[NULL]	1
3	62,397	6,027	[NULL]	1
4	62,397	3,401	[NULL]	1
5	62,397	2,904	[NULL]	1
6	62,397	611	[NULL]	1
7	62,397	5,390	[NULL]	1
8	62,397	2,479	[NULL]	1

Después de la ejecución de setPrice

	123orderid	123prod id	123price	123quantity
1	62,397	1,147	12.14	1
2	62,397	5,981	17.67	1
3	62,397	6,027	20.98	1
4	62,397	3,401	14.57	1
5	62,397	2,904	16.56	1
6	62,397	611	11.04	1
7	62,397	5,390	26.5	1
8	62,397	2,479	21.2	1

getTopSales.sql: En la función getTopSales, se calculan inicialmente las ventas totales de cada película por año utilizando `SUM(od.quantity)`, lo que suma la cantidad total de unidades vendidas de cada película en lugar de contar simplemente las apariciones del movieid. Luego para determinar las películas más vendidas por año, se aplica la función `RANK()` a estas sumas, asignando un rango a cada película dentro de cada año, donde el rango 1 indica la película más vendida.


```
select * from
getTopActors('Drama')
```

	actorname	nummovies	debutyear	debutfilm	directorname
1	Duval, Robert (I)	26	1.962	To Kill a Mockingbird (1962)	Mulligan, Robert
2	Jackson, Samuel L.	26	1.988	School Daze (1988)	Lee, Spike
3	De Niro, Robert	24	1.971	Born to Win (1971)	Passer, Ivan
4	Walsh, M. Emmet	23	1.969	Midnight Cowboy (1969)	Schlesinger, John
5	Kaitel, Harvey	22	1.976	Taxi Driver (1976)	Scorsese, Martin
6	Hitchcock, Alfred (I)	21	1.927	Lodger, The (1927)	Hitchcock, Alfred (I)
7	Nicholson, Jack	21	1.969	Easy Rider (1969)	Hopper, Dennis
8	Turturro, John	21	1.980	Raging Bull (1980)	Scorsese, Martin
9	Walsh, J.T.	21	1.987	Good Morning, Vietnam (1987)	Levinson, Barry (I)
10	Walsh, J.T.	21	1.987	Tin Men (1987)	Levinson, Barry (I)

```
select
    ia.actorname, ia2.actorid,
    count(ia2.actorid)
from
    imdb_actors ia
    join imdb_actormovies ia2 on
ia.actorid = ia2.actorid
    join imdb_movies im on
ia2.movieid = im.movieid
    join imdb_moviegenres im2 on
im.movieid = im2.movieid
    join genres g on im2.genreid =
g.genreid
where
    g.genrename = 'Drama' and
    ia.actorname = 'Jackson, Samuel L.'
group by
    ia2.actorid, ia.actorname
order by
    count(ia2.actorid)
desc
```

	actorname	actorid	count
1	Jackson, Samuel L.	305,494	26

```
select
    id2.directorname
from
    imdb_movies im
    join imdb_directormovies id on
im.movieid = id.movieid
    join imdb_directors id2 on
id.directorid = id2.directorid
where
    im.movietitle = 'To Kill a
Mockingbird (1962)'
```

	directorname
1	Mulligan, Robert

Parte 2: Optimización

k) Estudio del impacto de un índice

- a) Crear una consulta, estadosDistintos.sql, que muestre el número de estados (columna state) distintos con clientes que tienen pedidos en un año dado usando el formato YYYY (por ejemplo 2017) y que además pertenecen a un país (country) determinado, por ejemplo, Peru.

```
SELECT COUNT(DISTINCT c.state)
FROM public.customers c
JOIN public.orders o ON c.customerid = o.customerid
WHERE c.country = 'Peru' AND EXTRACT(YEAR FROM o.orderdate) = 2017;
```

- b) Mediante la sentencia EXPLAIN estudiar el plan de ejecución de la consulta.

```
EXPLAIN
SELECT COUNT(DISTINCT c.state)
FROM public.customers c
JOIN public.orders o ON c.customerid = o.customerid
WHERE c.country = 'Peru' AND EXTRACT(YEAR FROM o.orderdate) = 2017;
```

```

QUERY PLAN
-----
Aggregate  (cost=4821.98..4821.99 rows=1 width=8)
  -> Gather  (cost=1529.04..4821.97 rows=5 width=118)
        Workers Planned: 1
        -> Hash Join  (cost=529.04..3821.47 rows=3 width=118)
              Hash Cond: (o.customerid = c.customerid)
              -> Parallel Seq Scan on orders o  (cost=0.00..3291.03 rows=535 width=4)
                    Filter: (EXTRACT(year FROM orderdate) = '2017'::numeric)
              -> Hash  (cost=528.16..528.16 rows=70 width=122)
                    -> Seq Scan on customers c  (cost=0.00..528.16 rows=70 width=122)
                          Filter: ((country)::text = 'Peru'::text)

(10 filas)

count
-----
  185
(1 fila)
```

Utilizando la instrucción EXPLAIN, se ha calculado un costo total aproximado de 4821.99 para la ejecución de la consulta. Este costo incluye tanto la localización como el procesamiento de una única fila. Se identificó que la estrategia de unión de las tablas es mediante un 'Hash Join'. Inicialmente, se realiza un escaneo secuencial en paralelo (Parallel Seq Scan) en la tabla "orders", lo que representa un

costo considerable de 3291.03. Este proceso conlleva un escaneo completo de la tabla "orders", aplicando un filtro basado en las condiciones de la cláusula WHERE. A continuación, se efectúa un escaneo secuencial (Seq Scan) en la tabla "customers" que tiene un costo proyectado de 528.16, donde se filtran los registros que coinciden con el país 'Peru'. La combinación de estos conjuntos de datos se realiza mediante un "Hash Join" que tiene un costo estimado de 3821.47, juntando la información según la coincidencia en los identificadores de cliente (customerid).

c,d) Identificar un índice que mejore el rendimiento de la consulta y crearlo (si ya existía, borrarlo). Estudiar el nuevo plan de ejecución y compararlo con el anterior.

```
CREATE INDEX INDEX_ORDERDATE ON public.orders (EXTRACT (YEAR FROM
ORDERDATE)) ;
```

```

QUERY PLAN
-----
Aggregate  (cost=2030.23..2030.24 rows=1 width=8)
-> Hash Join  (cost=548.50..2030.22 rows=5 width=118)
    Hash Cond: (o.customerid = c.customerid)
    -> Bitmap Heap Scan on orders o  (cost=19.46..1498.79 rows=909 width=4)
        Recheck Cond: (EXTRACT(year FROM orderdate) = '2017'::numeric)
        -> Bitmap Index Scan on index_orderdate  (cost=0.00..19.24 rows=909 width=0)
            Index Cond: (EXTRACT(year FROM orderdate) = '2017'::numeric)
    -> Hash  (cost=528.16..528.16 rows=70 width=122)
        -> Seq Scan on customers c  (cost=0.00..528.16 rows=70 width=122)
            Filter: ((country)::text = 'Peru'::text)

(10 filas)

count
-----
185
(1 fila)

```

Después de implementar el nuevo índice, lo hemos ejecutado en conjunto con la instrucción EXPLAIN para evaluar su impacto en la eficiencia y otros aspectos relevantes del proceso. Los resultados de las pruebas han confirmado una notable disminución en los costos de ejecución, pasando de 4821.99 a 2030.23. Esto representa una mejora considerable en términos de rendimiento y optimización del uso de recursos. Esta mejora se atribuye al hecho de que el índice fue creado específicamente para optimizar las columnas referenciadas en la cláusula WHERE de nuestra consulta.

e) Probad distintos índices y discutid los resultados

```
CREATE INDEX INDEX_COUNTRY ON public.customers (COUNTRY) ;
```

```

QUERY PLAN
-----
Aggregate  (cost=1681.90..1681.91 rows=1 width=8)
-> Hash Join  (cost=200.17..1681.89 rows=5 width=118)
    Hash Cond: (o.customerid = c.customerid)
    -> Bitmap Heap Scan on orders o  (cost=19.46..1498.79 rows=909 width=4)
        Recheck Cond: (EXTRACT(year FROM orderdate) = '2017'::numeric)
        -> Bitmap Index Scan on index_orderdate  (cost=0.00..19.24 rows=909 width=0)
            Index Cond: (EXTRACT(year FROM orderdate) = '2017'::numeric)
    -> Hash  (cost=179.83..179.83 rows=70 width=122)
        -> Bitmap Heap Scan on customers c  (cost=4.83..179.83 rows=70 width=122)
            Recheck Cond: ((country)::text = 'Peru'::text)
            -> Bitmap Index Scan on index_country  (cost=0.00..4.81 rows=70 width=0)
                Index Cond: ((country)::text = 'Peru'::text)

(12 filas)

count
-----
    185
(1 fila)

```

Hemos introducido un nuevo índice denominado "INDEX_COUNTRY" y lo hemos puesto a prueba en comparación con el índice previamente creado. Iniciamos el proceso eliminando el índice antiguo mediante un comando DROP. Posteriormente, al ejecutar nuevamente la consulta con el nuevo índice, el resultado mostró un costo de 1681.90 en la operación de agregación.

Al comparar este resultado con los obtenidos anteriormente, se deduce que el índice "INDEX_COUNTRY" logra una reducción más significativa en los costos de ejecución y otros factores relacionados. Por lo tanto, concluimos que el uso del índice "INDEX_COUNTRY" resulta en una optimización más efectiva del rendimiento y los recursos en comparación con el índice "INDEX_ORDERDATE".

f) Modificar el script estadosDistintos.sql, añadiendo las sentencias de creación de índices y de planificación.

```

CREATE INDEX INDEX_ORDERDATE ON public.orders (EXTRACT(YEAR FROM
ORDERDATE));
CREATE INDEX INDEX_COUNTRY ON public.customers (COUNTRY);

--DROP INDEX INDEX_ORDERDATE;
--DROP INDEX INDEX_COUNTRY;

EXPLAIN
SELECT COUNT(DISTINCT c.state)
FROM public.customers c
JOIN public.orders o ON c.customerid = o.customerid
WHERE c.country = 'Peru' AND EXTRACT(YEAR FROM o.orderdate) = 2017;

```

```
SELECT COUNT(DISTINCT c.state)
FROM public.customers c
JOIN public.orders o ON c.customerid = o.customerid
WHERE c.country = 'Peru' AND EXTRACT(YEAR FROM o.orderdate) = 2017;
```

I) Estudio del impacto de cambiar la forma de realizar una consulta:

a) Estudiar los planes de ejecución de las consultas alternativas mostradas en el Anexo 1 y compararlos.

La primera consulta emplea 'NOT IN' para identificar aquellos customerid en la tabla customers que no coinciden con ningún customerid en los pedidos marcados como 'Paid'. Esta operación se lleva a cabo mediante una subconsulta en la tabla orders.

En la segunda consulta, se combina el uso de 'UNION ALL' para fusionar los customerid de customers y aquellos de orders con un estado 'Paid'. Posteriormente, aplica 'GROUP BY' y 'HAVING' para filtrar y mostrar solo aquellos customerid que aparecen una sola vez.

La tercera consulta, por su parte, recurre a 'EXCEPT' para obtener los customerid de customers excluyendo aquellos que se encuentran en los pedidos con estado 'Paid', identificados a través de una subconsulta en orders.

i) ¿Qué consulta devuelve algún resultado nada más comenzar su ejecución?

La consulta que inmediatamente empieza a entregar resultados es aquella que emplea la cláusula EXCEPT, ya que esta operación selecciona directamente aquellos registros de la tabla customers que no tienen correspondencia en la tabla orders. Este método es eficiente porque identifica y excluye los elementos de customers que no se encuentran en orders sin la necesidad de procesar completamente la tabla orders. En otras palabras, la consulta con EXCEPT es eficaz en generar resultados de forma más rápida al enfocarse específicamente en los registros no coincidentes entre las dos tablas mencionadas.

ii) ¿Qué consulta se puede beneficiar de la ejecución en paralelo?

En lo que respecta a la ejecución en paralelo, la consulta que más probablemente se beneficie de esta técnica es aquella que incluye operaciones de agrupación GROUP BY y combina varias subconsultas mediante UNION ALL. Estas operaciones son particularmente aptas para ser ejecutadas de manera paralela, lo que puede acelerar significativamente el proceso de tratamiento de los datos. Esta ventaja depende, no obstante, de que tanto el sistema de gestión de bases de datos

como su configuración actual permitan y sean compatibles con la ejecución paralela. En resumen, la capacidad de realizar tareas en paralelo puede optimizar notablemente el rendimiento de consultas complejas que involucren agrupaciones y uniones de múltiples conjuntos de datos.

m) Estudio del impacto de la generación de estadísticas:

b,c) Partir de la base de datos suministrada para este apartado (recién creada y cargada de datos). Estudiar con la sentencia EXPLAIN el coste de ejecución de las dos consultas indicadas en el Anexo 2.

El plan de ejecución revela que esta consulta realiza una operación de agregación para calcular el recuento de filas en la tabla "orders" donde el campo "STATUS" es NULL. Se utiliza una exploración secuencial de la tabla (Seq Scan), y se aplica un filtro para seleccionar las filas deseadas. El costo estimado de esta consulta es de 3507.17, pero el tiempo de ejecución real es sorprendentemente bajo, con un promedio de 12.125 ms. Esto se debe a que el filtro selecciona una minoría de las filas en la tabla, lo que reduce significativamente la cantidad de datos a procesar.

```

QUERY PLAN
-----
Aggregate  (cost=3507.17..3507.18 rows=1 width=8) (actual time=12.099..12.100 rows=1 loops=1)
-> Seq Scan on orders  (cost=0.00..3504.90 rows=909 width=0) (actual time=12.096..12.096 rows=0 loops=1)
    Filter: (status IS NULL)
    Rows Removed by Filter: 181790
Planning Time: 0.119 ms
Execution Time: 12.125 ms
(6 filas)

```

En esta segunda consulta, nuevamente se realiza una operación de agregación para calcular el recuento de filas en la tabla "orders", pero esta vez se filtra en base a la condición de que el campo "STATUS" sea igual a 'Shipped'. Al igual que en la primera consulta, se utiliza una exploración secuencial de la tabla (Seq Scan). El costo estimado de esta consulta es de 3961.65, y el tiempo de ejecución real es de 30.886 ms. La diferencia en el tiempo de ejecución en comparación con la primera consulta se debe al filtro más restrictivo que selecciona un subconjunto menor de filas, lo que resulta en un mayor costo estimado.

```

QUERY PLAN
-----
Aggregate  (cost=3961.65..3961.66 rows=1 width=8) (actual time=30.858..30.859 rows=1 loops=1)
-> Seq Scan on orders  (cost=0.00..3959.38 rows=909 width=0) (actual time=0.011..24.706 rows=127323 loops=1)
    Filter: ((status)::text = 'Shipped'::text)
    Rows Removed by Filter: 54467
Planning Time: 0.109 ms
Execution Time: 30.886 ms
(6 filas)

```

Como conclusión de estas dos consultas podemos determinar que la primera consulta es más eficiente en términos de tiempo de ejecución debido a una

condición de filtro menos restrictiva (NULL en lugar de una cadena específica), lo que resulta en un menor costo estimado y un tiempo de ejecución más rápido.

d,e) Crear un índice en la tabla orders por la columna status. Estudiar de nuevo la planificación de las mismas consultas.

```
CREATE INDEX INDEX_STATUS ON public.orders (STATUS)
```

Como hemos observado en el apartado anterior, la diferencia en el tiempo de ejecución en comparación con la primera consulta se debe al filtro más restrictivo en la segunda consulta, que selecciona un subconjunto menor de filas, lo que resulta en un costo estimado más alto. De cualquier forma, se observa cómo se han reducido significativamente los tiempos de ejecución después de crear el índice.

```

QUERY PLAN
-----
Aggregate  (cost=3507.17..3507.18 rows=1 width=8) (actual time=12.198..12.198 rows=1 loops=1)
-> Seq Scan on orders  (cost=0.00..3504.90 rows=909 width=0) (actual time=12.193..12.194 rows=0 loops=1)
    Filter: (status IS NULL)
    Rows Removed by Filter: 181790
Planning Time: 0.075 ms
Execution Time: 12.227 ms
(6 filas)

QUERY PLAN
-----
Aggregate  (cost=3961.65..3961.66 rows=1 width=8) (actual time=20.105..20.106 rows=1 loops=1)
-> Seq Scan on orders  (cost=0.00..3959.38 rows=909 width=0) (actual time=0.004..16.404 rows=127323 loops=1)
    Filter: ((status)::text = 'Shipped'::text)
    Rows Removed by Filter: 54467
Planning Time: 0.039 ms
Execution Time: 20.121 ms
(6 filas)

```

El código utilizado, en conjunto, es el siguiente:

```

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS IS NULL;

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Shipped';

CREATE INDEX INDEX_STATUS ON public.orders (STATUS);
--DROP INDEX INDEX_STATUS;

EXPLAIN ANALYZE
SELECT COUNT(*)

```

```
FROM public.orders
WHERE STATUS IS NULL;

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Shipped';
```

f,g,h) Ejecutar la sentencia ANALYZE para generar las estadísticas sobre la tabla orders. Estudiar de nuevo el coste de las consultas y comparar con el coste anterior a la generación de estadísticas. Comparar el plan de ejecución y discutir el resultado. Comparar con la planificación de las otras dos consultas proporcionadas y comentar los resultados.

```
ANALYZE VERBOSE public.orders;

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS IS NULL;

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Shipped';

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Paid';

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Processed';
```

```

----- QUERY PLAN -----
Aggregate  (cost=4.32..4.33 rows=1 width=8) (actual time=0.006..0.007 rows=1 loops=1)
->  Index Only Scan using index_status on orders  (cost=0.29..4.31 rows=1 width=0) (actual time=0.004..0.004 rows=0 loops=1)
      Index Cond: (status IS NULL)
      Heap Fetches: 0
Planning Time: 0.059 ms
Execution Time: 0.019 ms
(6 filas)

----- QUERY PLAN -----
Aggregate  (cost=3011.06..3011.07 rows=1 width=8) (actual time=8.838..8.838 rows=1 loops=1)
->  Index Only Scan using index_status on orders  (cost=0.29..2690.21 rows=128338 width=0) (actual time=0.009..5.390 rows=127323 loops=1)
      Index Cond: (status = 'Shipped'::text)
      Heap Fetches: 0
Planning Time: 0.029 ms
Execution Time: 8.847 ms
(6 filas)

----- QUERY PLAN -----
Aggregate  (cost=420.72..420.73 rows=1 width=8) (actual time=1.504..1.504 rows=1 loops=1)
->  Index Only Scan using index_status on orders  (cost=0.29..376.16 rows=17821 width=0) (actual time=0.017..0.944 rows=18163 loops=1)
      Index Cond: (status = 'Paid'::text)
      Heap Fetches: 0
Planning Time: 0.026 ms
Execution Time: 1.512 ms
(6 filas)

----- QUERY PLAN -----
Aggregate  (cost=836.91..836.92 rows=1 width=8) (actual time=2.915..2.916 rows=1 loops=1)
->  Index Only Scan using index_status on orders  (cost=0.29..747.84 rows=35631 width=0) (actual time=0.010..1.829 rows=36304 loops=1)
      Index Cond: (status = 'Processed'::text)
      Heap Fetches: 0
Planning Time: 0.025 ms
Execution Time: 2.927 ms
(6 filas)

```

Luego de ejecutar el comando "ANALYZE" en PostgreSQL, se han observado reducciones significativas en los costos de ejecución tanto en la primera como en la segunda consulta. Esta mejora se debe a que:

El comando "ANALYZE" se utiliza para recopilar estadísticas precisas sobre las tablas y los índices en la base de datos. Estas estadísticas proporcionan al optimizador de consultas información valiosa sobre cómo están distribuidos los datos, el tamaño de las tablas y otros detalles relevantes.

Al disponer de información actualizada sobre la distribución de los datos y la cantidad de filas en las tablas, el optimizador puede tomar decisiones más acertadas sobre cómo acceder y unir los datos requeridos para una consulta específica. Esto puede llevar a la selección de planes de ejecución más eficientes, lo que se traduce en una disminución potencial del tiempo de ejecución y de los costos asociados con la consulta.

i) Crear un script countStatus.sql, conteniendo las consultas, la creación de índices y las sentencias ANALYZE.

El contenido del archivo es el siguiente:

```

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders

```

```
WHERE STATUS IS NULL;

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Shipped';

CREATE INDEX INDEX_STATUS ON public.orders(STATUS);
--DROP INDEX INDEX_STATUS;

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS IS NULL;

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Shipped';

ANALYZE VERBOSE public.orders;

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS IS NULL;

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Shipped';

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Paid';

EXPLAIN ANALYZE
SELECT COUNT(*)
FROM public.orders
WHERE STATUS = 'Processed';
```


Cuestiones

¿Qué hace el generador de estadísticas?

El generador de estadísticas en PostgreSQL, activado por la sentencia ANALYZE, recopila información sobre las tablas de la base de datos y sus contenidos. Estas estadísticas incluyen datos como el número de filas, la distribución de los valores en las columnas, la frecuencia de los valores, etc. Esta información es crucial para el optimizador de consultas de PostgreSQL, ya que le permite tomar decisiones informadas sobre el mejor plan de ejecución para una consulta. Por ejemplo, puede decidir si usar un índice, qué tipo de join realizar, y en qué orden unir las tablas. Cuanto más precisas sean las estadísticas, más eficiente puede ser la planificación de las consultas.

¿Por qué la planificación de las dos consultas es la misma hasta que se generan las estadísticas?

Antes de generar estadísticas con ANALYZE, el planificador de consultas de PostgreSQL tiene que hacer suposiciones basadas en información limitada o desactualizada sobre los datos. Esto puede llevar a que el planificador elija el mismo plan de ejecución para diferentes consultas que podrían beneficiarse de planes diferentes. Una vez que se actualizan las estadísticas con ANALYZE, el planificador tiene información más precisa sobre los datos, lo que le permite tomar decisiones más informadas y potencialmente diferentes para cada consulta. Así, la planificación de las consultas puede diferir notablemente después de generar estadísticas actualizadas, lo que a menudo resulta en una mayor eficiencia y rendimiento.