

Databases

Ivin Lee

October 8, 2020

1 Lecture 1

Database Systems in course

- HyperSQL
- DOCTOR Who
- Neo4j

Interface

This course will deal with Database Management Systems from the perspective of an application developer, and will teach the interfaces of these systems.

Database Management System

Implements CRUD operations

- Create
- Read
- Update
- Delete

Implements ACID transactions for concurrent updates

- Atomicity: either all actions carried out or none carried out
- Consistency: database is consistent before and after transactions
- Isolation: every transaction happens independently of other transactions: allows for concurrency

- Durability: if transaction happens successfully, then its effects will persist

SQL

HyperSql - just a java file, overhead of using is low

NoSQL

Not only SQL - non-relational, distributed, open-source, horizontally scalable.

- Scalability: data stored on multiple machines
- Fault Tolerance: service can survive failure of some machines
- Lower Latency: data located closer to widely distributed users

Distributing Data

- Replication
- Partitioning

CAP Concepts

- Consistency: all reads return data up-to-date
- Availability: all clients can find some replica of the data
- Partition Tolerance: system can still operate despite loss or failure of part of the system

Cannot achieve all of these, some balance.

CAP Principle

In a highly distributed system,

- Assume that network partitions and other connectivity problems will occur
- Implementing ACID transactions is very difficult and slow
- trade-off between availability and consistency

Eventual consistency - if update activity ceases, then the system will reach a consistent state.

BASE

- BA: basically available
- S: soft state
- E: eventual consistency

This is an area of ongoing research.

Polyglot Persistence

Using various Database Management Systems to manage different sets of data

Conclusion: there will be a lot of changes in Database Management Systems.

2 SQL Commands

Aggregate Commands

```

1 select position, count(*) as total
2 from has_position
3 group by position
4 order by total desc;
```

SQL creates new rows where each row 'contains' all the rows which have the same position. Then, each column in the new table can only be aggregate commands like count(*), min(...).

Join

```

1 select title, genre
2 from movies
3 join has_genre on has_genre.movie_id = movies.movie_id
4 join genres on genres.genre_id = has_genre.genre_id
5 where year = 2017
6 limit 20;
```

| TITLE | GENRE |
|----------------------|-----------|
| Wonder Woman | Fantasy |
| Wonder Woman | Action |
| Wonder Woman | Adventure |
| It | Horror |
| The Greatest Showman | Drama |
| The Greatest Showman | Musical |

```

9 The Greatest Showman      Biography
10 The Foreigner            Action
11 The Foreigner            Thriller
12 A Dog's Purpose          Drama
13 A Dog's Purpose          Comedy
14 A Dog's Purpose          Adventure
15 Blade Runner 2049        Drama
16 Blade Runner 2049        Action
17 Blade Runner 2049        Mystery
18 Spider-Man: Homecoming    Action
19 Spider-Man: Homecoming    Sci-Fi
20 Spider-Man: Homecoming    Adventure
21 Coco                     Animation
22 Coco                     Adventure

```

Joining duplicates entries. Since some movies might have more than 1 genre, identifying these movies requires work as SQL evaluates each row separately so `select * from movies join genres where genre = genre_1 and genre = genre_2` doesn't work.

Multiple Joins

The way around this problem is to just join the `genres` table multiple times so each row has two `genres` columns!

```

1 select title, year, rating, votes
2 from movies as m
3 join has_genre as hg1 on hg1.movie_id = m.movie_id
4 join has_genre as hg2 on hg2.movie_id = m.movie_id
5 join genres as g1 on g1.genre_id = hg1.genre_id
6 join genres as g2 on g2.genre_id = hg2.genre_id
7 where m.votes > 100000 and g1.genre = 'Romance' and g2.
      genre = 'Comedy'
8 order by votes desc;

```

If you are discerning enough, you'll realise that there are actually 4 rows for each original row! This results in a bigger search time complexity. Example rows are shown below.

| 1 | TITLE | GENRE | GENRE |
|---|----------------------|-----------|-----------|
| 2 | ----- | ----- | ----- |
| 3 | In the Mood for Love | Drama | Drama |
| 4 | In the Mood for Love | Drama | Romance |
| 5 | In the Mood for Love | Romance | Drama |
| 6 | In the Mood for Love | Romance | Romance |
| 7 | Chicken Run | Animation | Animation |

| | | | |
|----|-------------|-----------|-----------|
| 8 | Chicken Run | Animation | Comedy |
| 9 | Chicken Run | Animation | Adventure |
| 10 | Chicken Run | Comedy | Animation |
| 11 | Chicken Run | Comedy | Comedy |
| 12 | Chicken Run | Comedy | Adventure |

Nesting

Possible to nest queries, i.e. search a sample space generated from another select statement.

```

1  select m1.title, m1.year, g.genre, m1.rating, m1.votes
2  from movies as m1
3  join has_genre as hg on hg.movie_id = m1.movie_id
4  join genres as g on g.genre_id = hg.genre_id
5  where m1.votes > 100000 and (not (g.genre = 'Romance' or
      g.genre = 'Comedy'))
6      and m1.movie_id in
7          (select m2.movie_id
8           from movies as m2
9           join has_genre as hg1 on hg1.movie_id = m2.
              movie_id
10          join has_genre as hg2 on hg2.movie_id = m2.
              movie_id
11          join genres as g1 on g1.genre_id = hg1.genre_id
12          join genres as g2 on g2.genre_id = hg2.genre_id
13          where g1.genre = 'Romance' and g2.genre = '
              Comedy')
14  order by m1.votes desc
15  limit 10;

```

Views

Instead of nesting select statements, it is possible to create a view, which can be used in the same way as a table.

```

1  drop view if exists romcom_ids;
2
3  create view romcom_ids as
4      select m.movie_id as movie_id
5      from movies as m
6      join has_genre as hg1 on hg1.movie_id = m.movie_id
7      join has_genre as hg2 on hg2.movie_id = m.movie_id
8      join genres as g1 on g1.genre_id = hg1.genre_id
9      join genres as g2 on g2.genre_id = hg2.genre_id
10     where g1.genre = 'Romance' and g2.genre = 'Comedy';
11

```

```

12 select m.title, m.year, g.genre, m.rating, m.votes
13 from romcom_ids as r
14 join has_genre as hg on hg.movie_id = r.movie_id
15 join genres as g on g.genre_id = hg.genre_id
16 join movies as m on m.movie_id = r.movie_id
17 where m.votes > 100000 and (not (g.genre = 'Romance' or
    g.genre = 'Comedy'))
18 order by m.votes desc
19 limit 10;

```

Three-valued Logic

SQL allows columns to be null. This means that statements can evaluate to true, false and null.

```

1 select count(*)
2 from people
3 where deathYear = null
4
5 returns 0

```

* = null returns null, and select ... where only returns records when the where clause returns true.

To solve this, SQL introduced is null.

```

1 select count(*)
2 from people
3 where deathYear is null
4
5 returns 5919

```

Left Join

The way join works is to strike off the rows which cannot match the other table. However, in some situations we would like to preserve these rows, and we can do so using left join or right join.

```

1 select name, position
2 from people as p
3 join has_position as c on p.person_id = c.person_id
4 join movies as m on c.movie_id = m.movie_id
5 where title = 'Silver Linings Playbook';

```

| NAME | POSITION |
|----------------|----------|
| ----- | ----- |
| Robert De Niro | actor |
| Bruce Cohen | producer |

```

5 Bradley Cooper      actor
6 Donna Gigliotti    producer
7 Jonathan Gordon    producer
8 David O. Russell   director
9 Jacki Weaver       actor
10 Jennifer Lawrence  actor
11 Matthew Quick     writer

```

Joining this table with the roles table gives:

```

1 select name, position, role
2 from people as p
3 join has_position as c on p.person_id = c.person_id
4 join movies as m on c.movie_id = m.movie_id
5 join plays_role as r on r.movie_id = m.movie_id and r.
   person_id = c.person_id
6 where title = 'Silver Linings Playbook';

```

| 1 | NAME | POSITION | ROLE |
|---|-------------------|----------|---------|
| 2 | ----- | ----- | ----- |
| 3 | Robert De Niro | actor | Pat Sr. |
| 4 | Bradley Cooper | actor | Pat |
| 5 | Jacki Weaver | actor | Dolores |
| 6 | Jennifer Lawrence | actor | Tiffany |

So some entries in the first table are not shown because only actors have roles. In this case, using a left join will help:

```

1 select name, position, role
2 from people as p
3 join has_position as c on p.person_id = c.person_id
4 join movies as m on c.movie_id = m.movie_id
5 left join plays_role as r on r.movie_id = m.movie_id and
   r.person_id = c.person_id
6 where title = 'Silver Linings Playbook';

```

| 1 | NAME | POSITION | ROLE |
|----|-------------------|----------|---------|
| 2 | ----- | ----- | ----- |
| 3 | Robert De Niro | actor | Pat Sr. |
| 4 | Bruce Cohen | producer | [null] |
| 5 | Bradley Cooper | actor | Pat |
| 6 | Donna Gigliotti | producer | [null] |
| 7 | Jonathan Gordon | producer | [null] |
| 8 | David O. Russell | director | [null] |
| 9 | Jacki Weaver | actor | Dolores |
| 10 | Jennifer Lawrence | actor | Tiffany |

```
11 Matthew Quick          writer      [null]
```

Note that null is displayed as [null] to distinguish between the empty string and actual null.

Select Distinct

SQL is based on multisets of records (i.e. sets containing duplicates, also called bags) so the following query will result in multiple repeated rows:

```
1 select r1.role as role, m1.title as title, m1.year as
   year
2 from plays_role as r1
3 join plays_role as r2 on r2.person_id = r1.person_id
4 join movies as m1 on m1.movie_id = r1.movie_id
5 join movies as m2 on m2.movie_id = r2.movie_id
6 join people as p on p.person_id = r1.person_id
7 where p.name = 'Noomi Rapace'
8        and r1.role = r2.role
9        and m1.movie_id <> m2.movie_id
10 order by m1.title, r1.role, m1.year;
```

| 1 | ROLE | TITLE | YEAR |
|---|------------------|---------------------------------------|-------|
| 2 | ----- | ----- | ----- |
| | ---- | | |
| 3 | Lisbeth Salander | The Girl Who Kicked the Hornet's Nest | 2009 |
| 4 | Lisbeth Salander | The Girl Who Kicked the Hornet's Nest | 2009 |
| 5 | Lisbeth Salander | The Girl Who Played with Fire | 2009 |
| 6 | Lisbeth Salander | The Girl Who Played with Fire | 2009 |
| 7 | Lisbeth Salander | The Girl with the Dragon Tattoo | 2009 |
| 8 | Lisbeth Salander | The Girl with the Dragon Tattoo | 2009 |

The keyword `select distinct` collapses the repeated rows. Note that the table with repeated listings is first generated, then the repeats removed, so there is no improvement in query time (I believe?)

```
1 select distinct r1.role as role, m1.title as title, m1.
   year as year
2 from plays_role as r1
```



```

3 join plays_role as r2 on r2.person_id = r1.person_id
4 join movies as m1 on m1.movie_id = r1.movie_id
5 join movies as m2 on m2.movie_id = r2.movie_id
6 join people as p on p.person_id = r1.person_id
7 where p.name = 'Noomi Rapace'
8        and r1.role = r2.role
9        and m1.movie_id <> m2.movie_id
10 order by m1.title, r1.role, m1.year;

```

| 1 | ROLE | TITLE | YEAR |
|---|------------------|---------------------------------------|-------|
| 2 | ----- | ----- | ----- |
| | ---- | | |
| 3 | Lisbeth Salander | The Girl Who Kicked the Hornet's Nest | 2009 |
| 4 | Lisbeth Salander | The Girl Who Played with Fire | 2009 |
| 5 | Lisbeth Salander | The Girl with the Dragon Tattoo | 2009 |