```
USE imdb;
/* Now that you have imported the data sets, let's explore some of the tables.
To begin with, it is beneficial to know the shape of the tables and whether any column has null
values.
Further in this segment, you will take a look at 'movies' and 'genre' tables.*/
-- Segment 1:
-- Q1. Find the total number of rows in each table of the schema?
-- Type your code below:
select count(*) from director_mapping
select count(*) from genre
select count(*) from movie
select count(*) from names
select count(*) from ratings
select count(*) from role_mapping;
-- Q2. Which columns in the movie table have null values?
-- Type your code below:
select title,
   year,
   date_published,
   duration,
   country,
   worlwide_gross_income,
```

languages,
production\_company

from movie

where year is null

or title is null

or date\_published is null

or country is null

or worlwide\_gross\_income is null

or languages is null

or production\_company is null;

-- Now as you can see four columns of the movie table has null values. Let's look at the at the movies released each year.

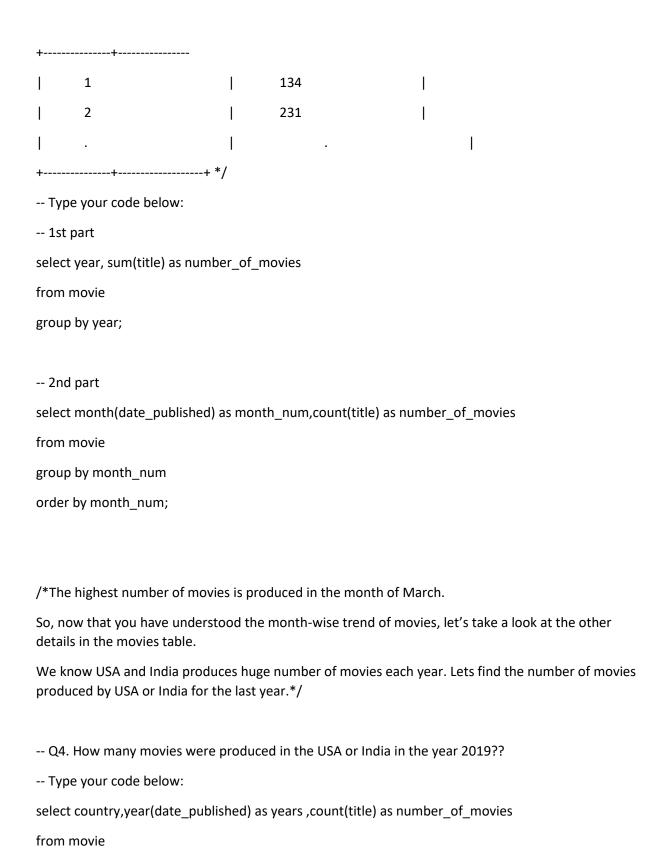
-- Q3. Find the total number of movies released each year? How does the trend look month wise?

(Output expected)

/\* Output format for the first part:

Output format for the second part of the question:

+-----+
| month\_num | number\_of\_movies|



-- as USA was not in the result in 2019 to make sure about that i wrote this code

where country = 'India' or country = 'USA'

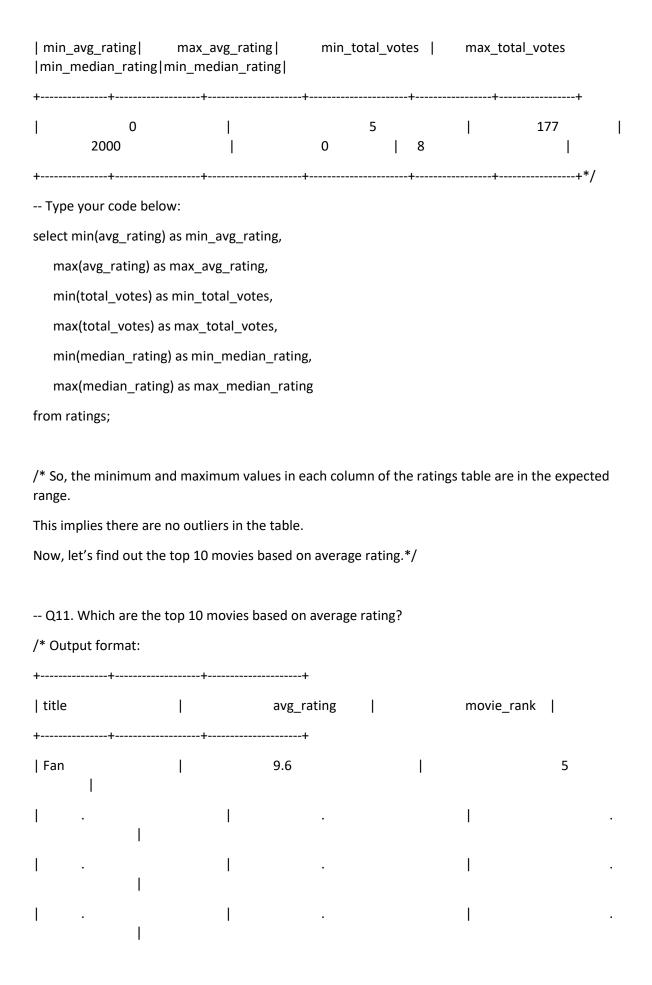
group by years;

$^{\prime *}$ USA and India produced more than a thousand movies(you know the exact number!) in the year 2019.
Exploring table Genre would be fun!!
Let's find out the different genres in the dataset.*/
Q5. Find the unique list of the genres present in the data set?
Type your code below:
select genre
from genre
group by genre
/* So, RSVP Movies plans to make a movie of one of these genres.
Now, wouldn't you want to know which genre had the highest number of movies produced in the last year?
Combining both the movie and genres table can give more interesting insights. */
Q6.Which genre had the highest number of movies produced overall?
Type your code below:
join is starting from
select count(id), genre
from movie m inner join genre g
on g.movie_id=m.id
group by genre;
/* So, based on the insight that you just drew, RSVP Movies should focus on the 'Drama' genre.
But wait, it is too early to decide. A movie can belong to two or more genres.
So, let's find out the count of movies that belong to only one genre.*/
Q7. How many movies belong to only one genre?
Type your code below:

```
select count(*) ,genre
from genre
group by genre;
/* There are more than three thousand movies which has only one genre associated with them.
So, this figure appears significant.
Now, let's find out the possible duration of RSVP Movies' next project.*/
-- Q8. What is the average duration of movies in each genre?
-- (Note: The same movie can belong to multiple genres.)
/* Output format:
                    avg_duration |
     thriller |
                           105
+----+*/
-- Type your code below:
select genre,avg(duration) as avg_duration
from movie m inner join genre g
on g.movie_id=m.id
group by genre;
/* Now you know, movies of genre 'Drama' (produced highest in number in 2019) has the average
duration of 106.77 mins.
```

Lets find where the movies of genre 'thriller' on the basis of number of movies.\*/

Q9.What is th movies produce		'thriller' g	genre of movies	among al	I the genres in terms of number	of
(Hint: Use the	Rank function	on)				
/* Output forma						
++   genre	1		movie_count	I	genre_rank	
+ drama	I	2312		I	2	I
Type your coo		+	+*/			
select genre, co			_		ward.	
			novie_id) desc) a	is genre_	rank	
from movie m in on g.movie_id=i		reg				
group by genre;						
group by genie,						
/*Thriller movie	es is in top 3 a	among all g	genres in terms	of numbe	er of movies	
In the previous	segment, yo	u analysed	the movies and	genres t	ables.	
In this segment	, you will and	alyse the ra	atings table as w	ell.		
To start with let	s get the mir	and max	values of differe	nt colum	ns in the table*/	
Segment 2:						
Q10. Find the movie_id colum		nd maximu	um values in ead	ch columi	n of the ratings table except the	
/* Output forma	at:					
+	·	+	+		++	



++	<del>-</del>	*/	
Type your code below:			
It's ok if RANK() or DENSE_F	RANK() is used to	o	
select title ,avg_rating as avg_	_rating,		
rank() over (order by avg_rat	ing desc) as mov	ie_rank	
from movie m inner join ratin	gs a		
on a.movie_id=m.id			
group by avg_rating;			
/* Do you find you favourite r please check your code again		top 10 movies with an a	average rating of 9.6? If not,
So, now that you know the to these movies?	p 10 movies, do	you think character acto	rs and filler actors can be from
Summarising the ratings table insight.*/	e based on the m	ovie counts by median r	ating can give an excellent
FYI fan is not my favourite r	novie		
Q12. Summarise the ratings	s table based on	the movie counts by me	dian ratings.
/* Output format:			
++			
median_rating	movie_count	1	
+			
1	I	105	1
] .	1		

- -- Type your code below:
- -- Order by is good to have

select median\_rating , count(median\_rating) as movie\_count

from ratings

```
group by median_rating
order by median_rating;
/* Movies with a median rating of 7 is highest in number.
Now, let's find out the production house with which RSVP Movies can partner for its next project.*/
-- Q13. Which production house has produced the most number of hit movies (average rating > 8)??
/* Output format:
+----+
|production_company|movie_count | prod_company_rank|
+-----+
The Archers
                       1
                                                            1
                                                                          1
+-----+*/
-- Type your code below:
select production_company, count(distinct id) as movie_count,
rank() over (order by count(distinct id) desc) as prod_company_rank
from movie m inner join ratings r
on r.movie_id=m.id
where avg_rating > 8 and production_company is not null
group by production_company;
-- It's ok if RANK() or DENSE_RANK() is used too
-- Answer can be Dream Warrior Pictures or National Theatre Live or both
-- No there are no dream warrior pictures or national theare or both in full list
-- Q14. How many movies released in each genre during March 2017 in the USA had more than
1,000 votes?
/* Output format:
        movie_count
genre
```

+	+										
th	riller		105			I					
			1					1			
			1					1			
+	+	+	*/								
Type you	ır code bel	ow:									
select g.ge	nre,count(i	id) as m	ovie_cou	ınt							
from movi	e m inner jo	oin genr	e g on g.	movie <sub>.</sub>	_id=m.id						
inner join	ratings r on	g.movi	e_id=r.m	ovie_i	d						
where cou	ntry = 'USA s > '1000'	d' and m	onth(dat	e_pub	lished) = '	3' and y	ear(date	e_publishe	ed) = '2017	' and	
group by g	enre										
order by n	ovie_coun	t;									
Lets try	o analyse v	with a u	nique pro	oblem	statemen	t.					
	d movies o	f each g	enre tha	t start	with the v	vord 'Tl	ne' and v	which have	e an avera	ge ratin	g >
8?											
/* Output											
	+										
title		-			_rating			genre	I		
	+	+			-+						
Theeran		ı		8.3			I		Thriller	I	
	1		I		•			I			•
.			1					1			
	1										
			I					I			
	+				. * /						
			·		-+ ' /						
	ır code bel										
	tle,r.avg_ra				اد! دماله:						
Irom movi	e m inner jo	oırı genr	H N UN Ø	MOVID	1/1-m 1/1						
inno-!-!-	ratings r on	_			_						

where title like 'The%' and avg\_rating > 8 order by avg\_rating desc;

- -- You should also try your hand at median rating and check whether the 'median rating' column gives any significant insights.
- -- Q16. Of the movies released between 1 April 2018 and 1 April 2019, how many were given a median rating of 8?
- -- Type your code below:

select count(title) as movie\_count

from movie m inner join genre g on g.movie\_id=m.id

inner join ratings r on g.movie\_id=r.movie\_id

where median\_rating = '8' and date\_published between '2018-04-01' and '2019-04-01'

- -- Once again, try to solve the problem given below.
- -- Q17. Do German movies get more votes than Italian movies?
- -- Hint: Here you have to find the total number of votes for both German and Italian movies.
- -- Type your code below:

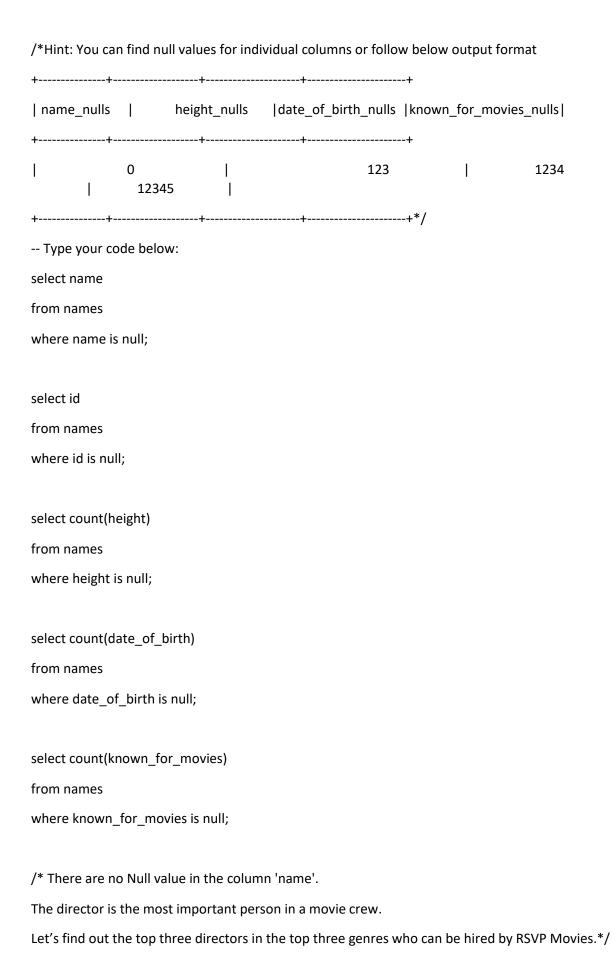
select languages,count(total\_votes) as total\_votes from movie m inner join genre g on g.movie\_id=m.id inner join ratings r on g.movie\_id=r.movie\_id where languages='german' or 'italian' group by total\_votes;

-- Answer is Yes

/\* Now that you have analysed the movies, genres and ratings tables, let us now analyse another table, the names table.

Let's begin by searching for null values in the tables.\*/

- -- Segment 3:
- -- Q18. Which columns in the names table have null values??



Q19. Who are the to > 8?	op three director	rs in the top thr	ee genres	whose mo	ovies have an a	verage rating
(Hint: The top three	genres would h	ave the most n	umber of	movies wit	h an average ra	ating > 8.)
/* Output format:						
+	+					
director_name	movie	e_count	1			
+						
James Mangold	1	4		1		
1 .	1					
1 .	1					
+	+ */					
Type your code belo	)W:					
select name as directo	or_name ,count(	known_for_mc	vies) as m	ovie_coun	it	
from names n inner jo	in genre g on g.r	movie_id=n.kno	own_for_r	novies		
inner join ratings r on	g.movie_id=r.m	ovie_id				
where avg_rating > '8'	ı					
group by director_nar	ne					
order by movie_count	: desc;					
/* James Mangold car 'Logan' and 'The Wolv		director for RS	VP's next	project. Do	you remeber	his movies,
Now, let's find out the	top two actors.	*/				
Q20. Who are the to	op two actors wh	nose movies ha	ve a medi	an rating >	= 8?	
/* Output format:						
+	+					
actor_name	movie_count	I				
+						
Christain Bale	10		I			
1 .	1			I		

```
+----+*/
-- Type your code below:
select name as actor_name,count(median_rating) as movie_count
from names n inner join role_mapping d on n.id = d.name_id
inner join ratings r on r.movie_id = d.movie_id
where median_rating >= '8'
group by actor_name
order by movie_count desc;
/* Have you find your favourite actor 'Mohanlal' in the list. If no, please check your code again.
RSVP Movies plans to partner with other global production houses.
Let's find out the top three production houses in the world.*/
-- Q21. Which are the top three production houses based on the number of votes received by their
movies?
/* Output format:
+-----+
|production_company|vote_count
                                                           prod_comp_rank|
+-----+
| The Archers
                                830
                                                    1
                                                                  1
+-----+*/
-- Type your code below:
select production_company as production_company , total_votes as vote_count ,
rank() over (order by total_votes desc) as prod_comp_rank
from movie m inner join ratings r on r.movie_id = m.id
group by production_company
order by vote_count desc;
```

/\*Yes Marvel Studios rules the movie world.

So, these are the top three production houses based on the number of votes received by the movies they have produced.

Since RSVP Movies is based out of Mumbai, India also wants to woo its local audience.

RSVP Movies also wants to hire a few Indian actors for its upcoming project to give a regional feel. Let's find who these actors could be.\*/

- -- Q22. Rank actors with movies released in India based on their average ratings. Which actor is at the top of the list?
- -- Note: The actor should have acted in at least five Indian movies.
- -- (Hint: You should use the weighted average based on votes. If the ratings clash, then the total number of votes should act as the tie breaker.)

## 

```
rank() over (order by sum(avg_rating*total_votes)/sum(total_votes) desc) as actor_rank
from movie m inner join
 ratings r on r.movie_id=m.id
inner join role_mapping t on t.movie_id = m.id
inner join names n on t.name_id=n.id
where country = 'India'
group by actor_name
having movie_count >4;
-- Top actor is Vijay Sethupathi
-- Q23. Find out the top five actresses in Hindi movies released in India based on their average
ratings?
-- Note: The actresses should have acted in at least three Indian movies.
-- (Hint: You should use the weighted average based on votes. If the ratings clash, then the total
number of votes should act as the tie breaker.)
/* Output format:
+-----+
actress name total votes
                                             movie count
      actress_avg_rating | actress_rank |
      Tabu
                                      3455 | 11 |
                                                                        8.42
   -----+*/
-- Type your code below:
select name as actress_name,
```

```
sum(total_votes) as total_votes ,
count(m.id) as movie_count,
sum(avg_rating*total_votes)/sum(total_votes) AS actress_avg_rating,
rank() over (order by sum(avg_rating*total_votes)/sum(total_votes) desc) as actress_rank
from movie m inner join
 ratings r on r.movie_id=m.id
inner join role_mapping t on t.movie_id = m.id
inner join names n on t.name_id=n.id
where country = 'India' and category = 'actress'
group by actress_name
having movie_count > 4;
/* Taapsee Pannu tops with average rating 7.74.
Now let us divide all the thriller movies in the following categories and find out their numbers.*/
/* Q24. Select thriller movies as per avg rating and classify them in the following category:
                       Rating > 8: Superhit movies
                       Rating between 7 and 8: Hit movies
                        Rating between 5 and 7: One-time-watch movies
                       Rating < 5: Flop movies
-- Type your code below:
select
 genre,avg_rating,
 case
   when avg_rating >= 8 then 'Superhit Movies'
    when avg_rating between 7 and 8 then 'Hit Movies'
      when avg_rating between 5 and 7 then 'One time watch Movies'
      else 'Flop Movies'
        end as film_type
```

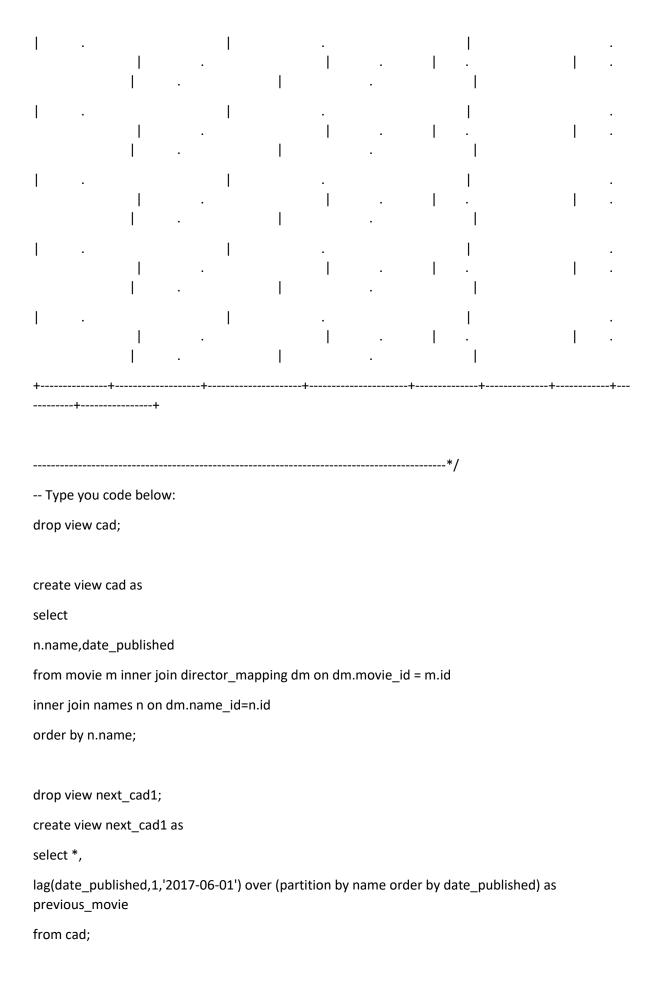
```
from genre g inner join ratings r on g.movie_id = r.movie_id
 where genre = 'Thriller'
/* Until now, you have analysed various tables of the data set.
Now, you will perform some tasks that will give you a broader understanding of the data in this
segment.*/
-- Segment 4:
-- Q25. What is the genre-wise running total and moving average of the average movie duration?
-- (Note: You need to show the output table in the question.)
/* Output format:
+-----+
                  avg_duration | running_total_duration | moving_avg_duration |
genre
+-----+
                                     145
comdy |
                                                           106.2
                                                                     -
128.42
-- Type your code below:
with duration as
select genre as genre,
sum(duration)/count(duration) as avg_duration
from genre g
inner join
movie m on g.movie_id=m.id
group by genre
```

```
)
select *,
 sum(avg_duration) over w1 as running_total_duration,
 avg(avg_duration) over w2 as moving_avg_duration
from duration
window w1 as (order by avg_duration rows unbounded preceding),
w2 as (order by avg_duration rows 13 preceding)
;
-- Round is good to have and not a must have; Same thing applies to sorting
-- Let us find top 5 movies of each year with top 3 genres.
-- Q26. Which are the five highest-grossing movies of each year that belong to the top three genres?
-- (Note: The top 3 genres would have the most number of movies.)
/* Output format:
genre
                              year
                                                            movie_name
|worldwide_gross_income|movie_rank |
                                             2017
       comedy
                                                                indian
$103244842
                              1
-- Type your code below:
-- Top 3 Genres based on most number of movies
with grossing as
(
```

```
select genre, year(date_published) as year,
      title as movie_name,
         worlwide_gross_income as worldwide_gross_income,
      rank() over (partition by year order by worlwide_gross_income desc) as movie_rank
                      from genre g inner join movie m on g.movie_id = m.id
                      where genre in ('Drama', 'comedy', 'thriller')
)
select * from grossing
where movie_rank < 6;
-- Finally, let's find out the names of the top two production houses that have produced the highest
number of hits among multilingual movies.
-- Q27. Which are the top two production houses that have produced the highest number of hits
(median rating >= 8) among multilingual movies?
/* Output format:
+-----+
|production company |movie count
                                                           prod comp rank
| The Archers
                                     830
                                                                          1
-- Type your code below:
with production as
select production_company as production_company , count(id) as movie_count,
row_number() over (order by count(id) desc) as prod_comp_rank
from movie m inner join ratings r on r.movie_id = m.id
where languages like '%,%' and median_rating > 8 and production_company is not null
group by production_company
```

```
)
select * from production
where prod_comp_rank < 3;
-- Multilingual is the important piece in the above question. It was created using POSITION(',' IN
languages)>0 logic
-- If there is a comma, that means the movie is of more than one language
-- Q28. Who are the top 3 actresses based on number of Super Hit movies (average rating >8) in
drama genre?
/* Output format:
+-----+
| actress_name | total_votes
                                    | movie_count
|actress_avg_rating |actress_rank |
      Laura Dern | 1
                                         1016 | 1 |
                                                                             9.60
-- Type your code below:
with actress as (
select name as actress_name,
sum(total_votes) as total_votes ,
count(m.id) as movie_count,
sum(avg_rating*total_votes)/sum(total_votes) AS actress_avg_rating,
row_number() over (order by sum(avg_rating*total_votes)/sum(total_votes) desc) as actress_rank
from movie m inner join
 ratings r on r.movie_id=m.id
inner join role_mapping t on t.movie_id = m.id
inner join names n on t.name_id=n.id
inner join genre g on g.movie_id = t.movie_id
where category = 'actress' and avg_rating > 8 and genre = 'drama'
```

```
group by actress_name
)
select * from actress
where actress_rank < 3;
/* Q29. Get the following details for top 9 directors (based on number of movies)
Director id
Name
Number of movies
Average inter movie duration in days
Average movie ratings
Total votes
Min rating
Max rating
total movie durations
Format:
----+
| director_id | director_name |
                                          number_of_movies | avg_inter_movie_days |
                     | total_votes | min_rating | max_rating | total_duration |
nm1777967
                            A.L. Vijay
          177
                                     5.65
                                                  1754
                                                                3.7
                                                                                      6.9
                            613
```



```
drop view final2;
create view final 2as
select *,
  datediff(date_published,previous_movie) as nextmovdiff
from next_cad1;
select dm.name_id as director_id,
n.name as director_name,
count(distinct(m.id)) as number_of_movies,
avg(nextmovdiff) as avg_inter_movie_days,
avg(avg_rating) as avg_rating,
sum(total_votes) as total_votes,
min(avg_rating) as min_rating,
max(avg_rating) as max_rating,
sum(duration) as total_duration
from movie m inner join genre g
on g.movie_id = m.id
inner join ratings r on r.movie_id = m.id
inner join director_mapping dm on dm.movie_id=m.id
inner join names n on dm.name_id = n.id
inner join final f1 on f1.name=n.name
group by n.name
order by count(distinct(m.id)) desc
limit 9;
```