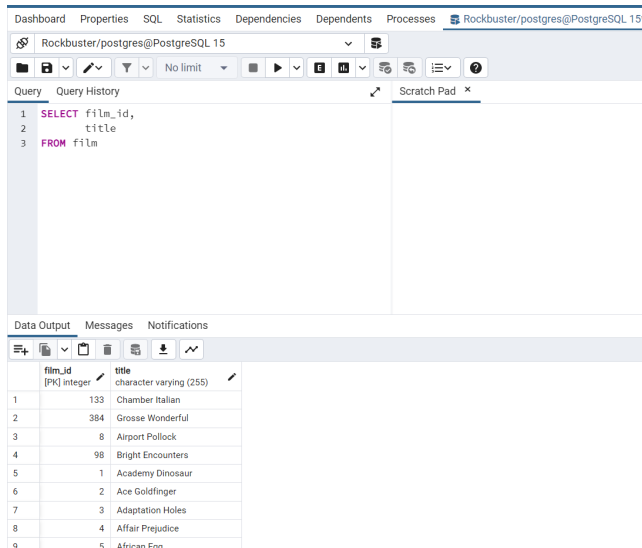


1.) **Refining Your Query: You need to get some data from the “film” table and decide to use the query SELECT \* FROM film.**

- You realize that only the “film\_id” and “title” columns are needed. Write a new query that selects only those 2 columns.



The screenshot shows a PostgreSQL query editor interface. The query editor has a tab labeled "Query" and a "Query History" button. The query text is as follows:

```
1 SELECT film_id,  
2 title  
3 FROM film
```

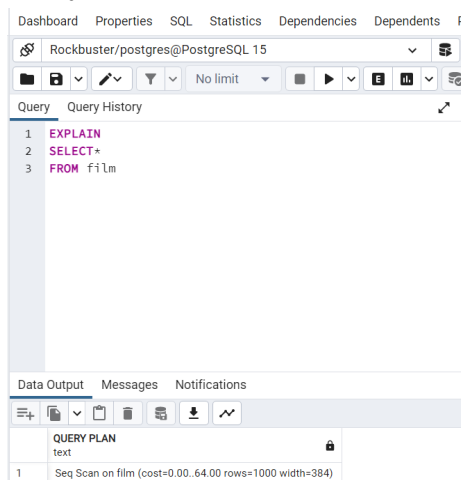
Below the query editor, there is a "Data Output" section with a table of results. The table has two columns: "film\_id" (integer) and "title" (character varying (255)). The results are as follows:

film_id	title
133	Chamber Italian
384	Grosse Wonderful
8	Airport Pollock
98	Bright Encounters
1	Academy Dinosaur
2	Ace Goldfinger
3	Adaptation Holes
4	Affair Prejudice
5	African Fox

- Compare the cost of the original query and the revised query, and write a few sentences explaining the comparison. Can you suggest any ways to optimize this query?

The costs of the respective queries are the same. The respective widths differ, though, with the revised query returning a narrower (i.e. less wide) result. Depending on what type of information is being sought, these queries could be optimized by limiting the number of rows returned with LIMIT or WHERE, to restrict returns based on certain specifications.

### Query 1



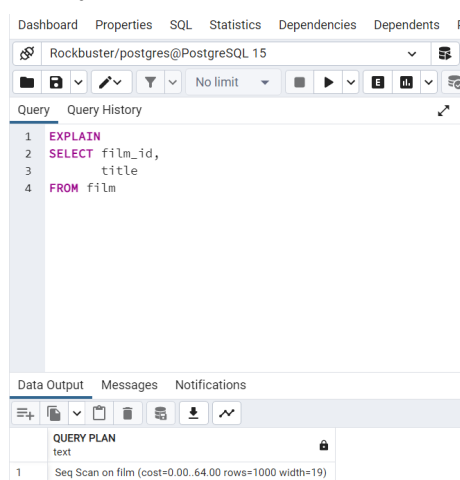
The screenshot shows a PostgreSQL query editor interface. The query editor has a tab labeled "Query" and a "Query History" button. The query text is as follows:

```
1 EXPLAIN  
2 SELECT*  
3 FROM film
```

Below the query editor, there is a "Data Output" section with a table of results. The table has two columns: "QUERY PLAN" (text) and "text". The results are as follows:

QUERY PLAN	text
Seq Scan on film	(cost=0.00..64.00 rows=1000 width=384)

### Query 2



The screenshot shows a PostgreSQL query editor interface. The query editor has a tab labeled "Query" and a "Query History" button. The query text is as follows:

```
1 EXPLAIN  
2 SELECT film_id,  
3 title  
4 FROM film
```

Below the query editor, there is a "Data Output" section with a table of results. The table has two columns: "QUERY PLAN" (text) and "text". The results are as follows:

QUERY PLAN	text
Seq Scan on film	(cost=0.00..64.00 rows=1000 width=19)

## 2.) Ordering the Data:

- In the pgAdmin Query Tool, run a query that selects every film from the “film” table, with the movies sorted by title from A to Z, then by most recent release year, and then by highest to lowest rental rate.

The screenshot shows the pgAdmin Query Tool interface. The query editor contains the following SQL query:

```
1 SELECT *
2 FROM film
3 ORDER BY title ASC
4         ,release_year DESC
5         ,rental_rate DESC
```

The Data Output tab shows the results of the query, sorted by title, release year, and rental rate. The table has the following columns: film\_id, title, description, release\_year, language\_id, rental\_duration, rental\_rate, length, replacement\_cost, rating, and last\_update.

film_id	title	description	release_year	language_id	rental_duration	rental_rate	length	replacement_cost	rating	last_update
1	Academy Dinosaur	A Epic Drama of a Feminist And ...	2006	1	6	0.99	86	20.99	PG	2013-05-26 14:50:58.951
2	Ace Goldfinger	A Astounding Epistle of a Datab...	2006	1	3	4.99	48	12.99	G	2013-05-26 14:50:58.951
3	Adaptation Holes	A Astounding Reflection of a Lu...	2006	1	7	2.99	50	18.99	NC-17	2013-05-26 14:50:58.951
4	Affair Prejudice	A Fanciful Documentary of a Fris...	2006	1	5	2.99	117	26.99	G	2013-05-26 14:50:58.951
5	African Egg	A Fast-Paced Documentary of a ...	2006	1	6	2.99	130	22.99	G	2013-05-26 14:50:58.951
6	Agent Truman	A Intrepid Panorama of a Robot ...	2006	1	3	2.99	169	17.99	PG	2013-05-26 14:50:58.951
7	Airplane Sierra	A Touching Saga of a Hunter An...	2006	1	6	4.99	62	28.99	PG-13	2013-05-26 14:50:58.951
8	Airport Pollock	A Epic Tale of a Moose And a Gir...	2006	1	6	4.99	54	15.99	R	2013-05-26 14:50:58.951

- Extract the data output of your query into a CSV file for the film collection department to analyze in Excel. To do this, click the button “Save results to file”

## 3.) Grouping Data: The strategy department has asked you the questions below. Write a SQL query to retrieve the correct answers, then extract your results as a CSV file.

- What is the average rental rate for each rating category?

The screenshot shows the pgAdmin Query Tool interface. The query editor contains the following SQL query:

```
1 SELECT rating,
2        Count(film_id),
3        AVG (rental_rate)
4 FROM film
5 GROUP BY rating
```

The Data Output tab shows the results of the query, grouped by rating. The table has the following columns: rating, count, and avg.

rating	count	avg
PG	194	3.0518556701030928
R	195	2.9387179487179487
NC-17	210	2.970952380952381
PG-13	223	3.034843049327354
G	178	2.888876404494382

PG = 3.0518556701030928

R= 2.9387179487179487

NC-17= 2.970952380952381

PG-13= 3.034843049327354

G= 2.888876404494382

- What are the minimum and maximum rental durations for each rating category?

The screenshot shows a PostgreSQL query editor interface. The query is as follows:

```

1 SELECT rating,
2     Count(film_id),
3     AVG (rental_rate),
4     MAX (rental_duration),
5     MIN (rental_duration)
6 FROM film
7 GROUP BY rating
8

```

Below the query editor, the 'Data Output' tab is active, displaying the results of the query in a table format:

	rating mpaa_rating	count bigint	avg numeric	max smallint	min smallint
1	PG	194	3.0518556701030928	7	3
2	R	195	2.9387179487179487	7	3
3	NC-17	210	2.970952380952381	7	3
4	PG-13	223	3.034843049327354	7	3
5	G	178	2.888876404494382	7	3

Rating	Maximum	Minimum
PG	7	3
R	7	3
NC-17	7	3
PG-13	7	3
G	7	3

4.) **Database Migration: Your team has decided to use an external tool to collect data on user behavior in the new Rockbuster Android app. Data collected from this new source will need to be loaded into the data warehouse before you can analyze it.**

- Can you outline the procedure for migrating the data and who will be responsible for it?

This process is commonly referred to as Extract, Transform, and Load in which the data is downloaded from the application, subsequently transformed into a different format, and then saved in the data warehouse. Typically, this is the domain of the data engineer but analysts should be aware of the general process.

- **What problems do you foresee if you start analyzing the data before it's been loaded into the data warehouse?**

Inconsistent formatting could lead to unintentional analysis (i.e. analyzing the wrong fields or analyzing data that is not yet in the intended format.) There could also be consistency issues and, if uploading data from multiple sources, faulty association of records.