# Sakila Database Analysis Report

## [Git and R]

**Author:** *Atique Ahmad*
**Student ID:** *MSDS25030*
**Date:** *11-Nov-2025*
**Subject:** CS 591: Tools and Techniques for Data Science

---

## Assignment Overview

This project analyzes the **Sakila Database** using **R** and the **data.table** package. The purpose of this analysis is to extract insights from the database, including details about films, customers, rentals, and payments. Additionally, visualization is used to summarize the distribution of film ratings.

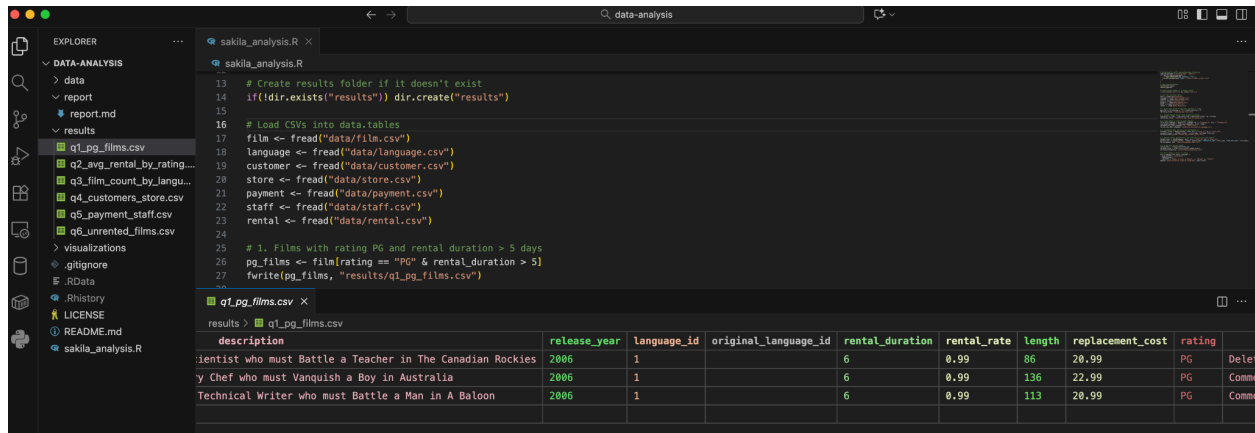All queries were executed using R scripts, and the results were stored as `.csv` files for reproducibility.

---

## 1. Films with Rating PG and Rental Duration > 5 Days

**Objective:** Identify films rated **PG** that have a rental duration greater than 5 days.

**R Code:**

```
pg_films <- film[rating == "PG" & rental_duration > 5]
fwrite(pg_films, "results/q1_pg_films.csv")
```
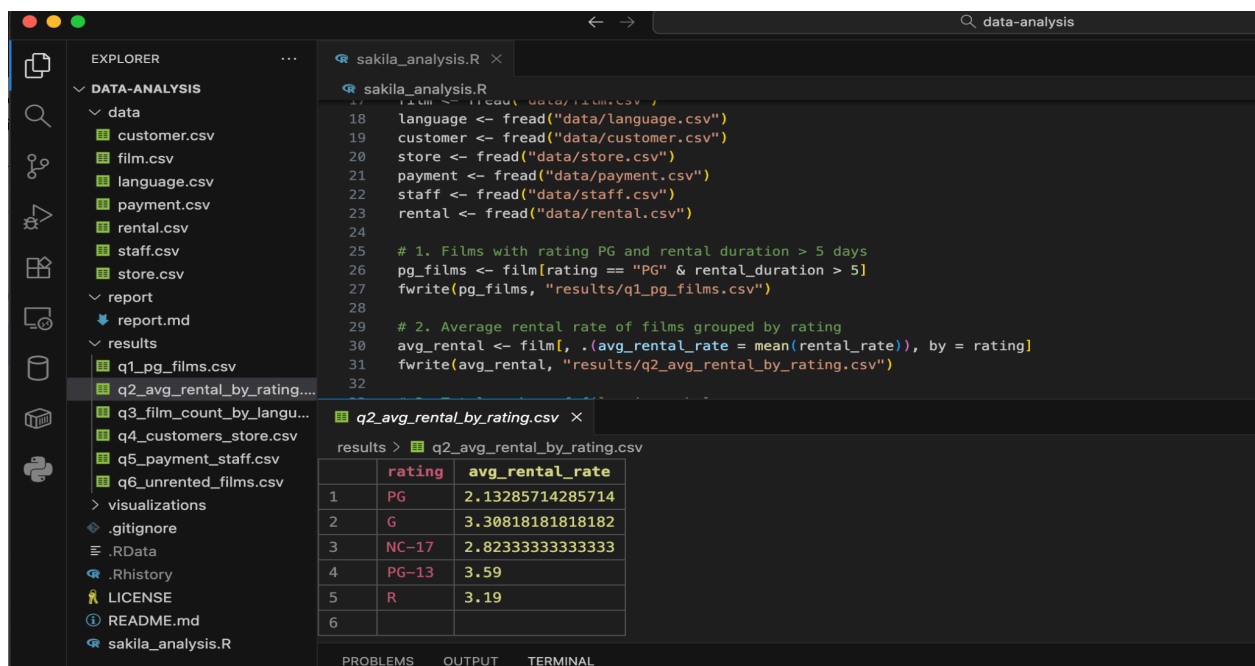
**Output File:** `results/q1_pg_films.csv`

---

# 2. Average Rental Rate of Films Grouped by Rating

**Objective:** Calculate the average rental rate for each film rating category.

**R Code:**

```r
avg_rental <- film[, .(avg_rental_rate = mean(rental_rate)), by = rating]
fwrite(avg_rental, "results/q2_avg_rental_by_rating.csv")
```
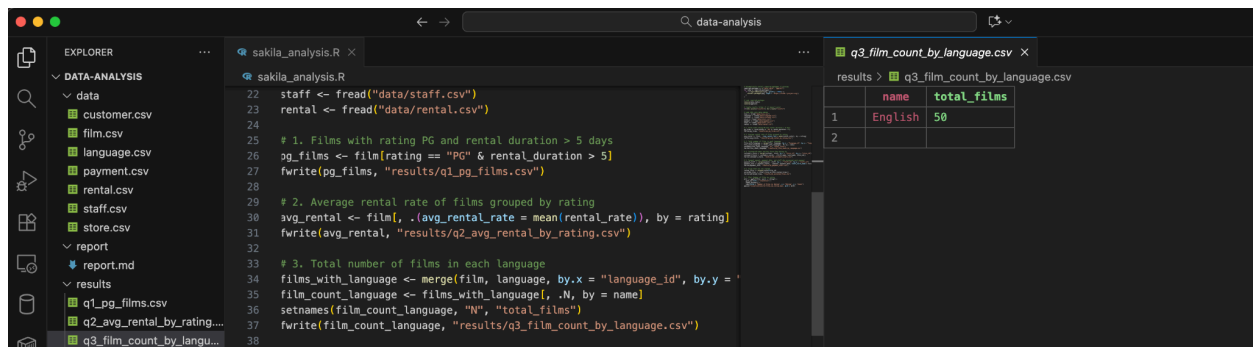
**Output File:** `results/q2_avg_rental_by_rating.csv`

# 3. Total Number of Films in Each Language

**Objective:** Count how many films exist for each language.

**R Code:**

```
films_with_language <- merge(film, language, by = "language_id")
film_count_language <- films_with_language[, .N, by = name]
setnames(film_count_language, "N", "total_films")
fwrite(film_count_language, "results/q3_film_count_by_language.csv")
```

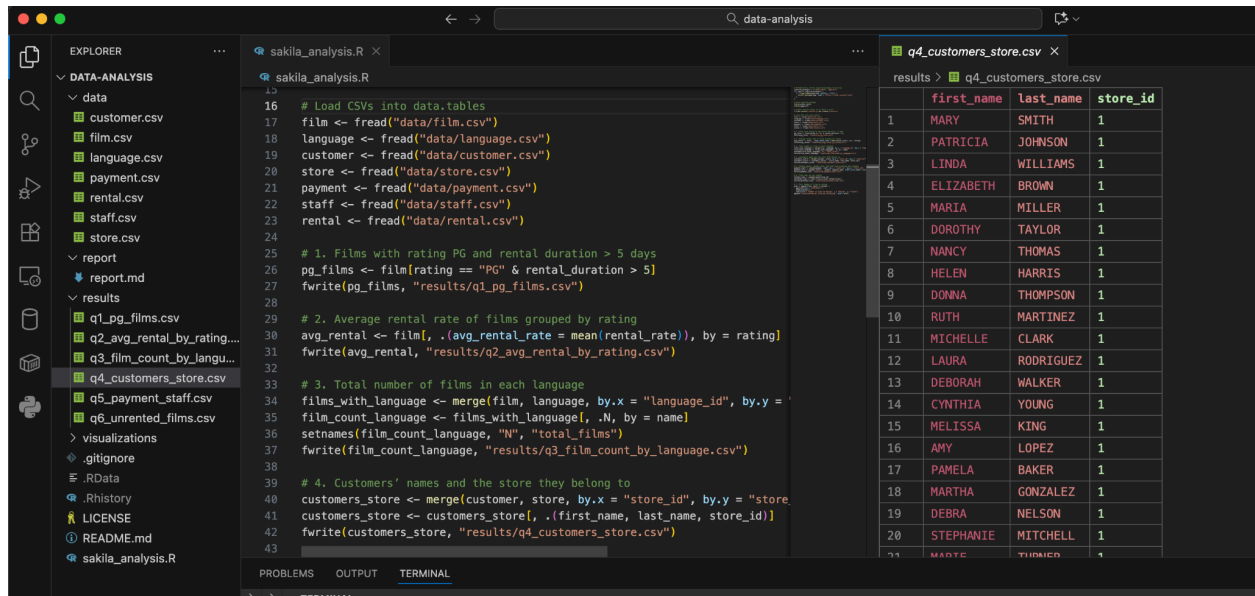**Output File:** `results/q3_film_count_by_language.csv`



# 4. Customers' Names and the Store They Belong To

**Objective:** Display each customer's name along with their store ID.

**R Code:**

```
customers_store <- merge(customer, store, by = "store_id")
customers_store <- customers_store[, .(first_name, last_name, store_id)]
fwrite(customers_store, "results/q4_customers_store.csv")
```

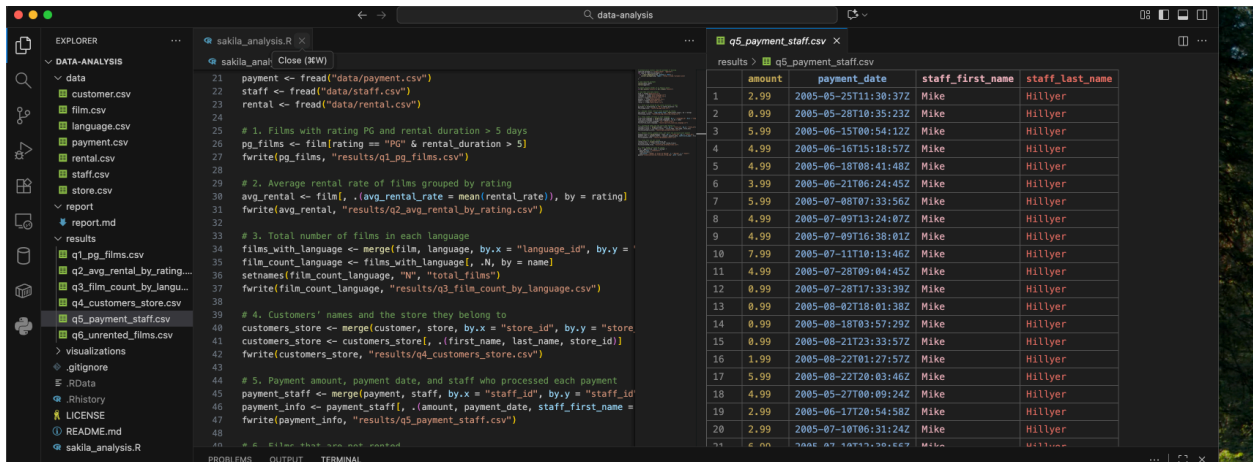**Output File:** `results/q4_customers_store.csv`

---

# 5. Payment Amount, Date, and Staff Who Processed It

**Objective:** Combine payment data with staff information to show transaction details.

**R Code:**

```
payment_staff <- merge(payment, staff, by = "staff_id")
payment_info <- payment_staff[, .(amount, payment_date, staff_first_name = first_name,
staff_last_name = last_name)]
fwrite(payment_info, "results/q5_payment_staff.csv")
```

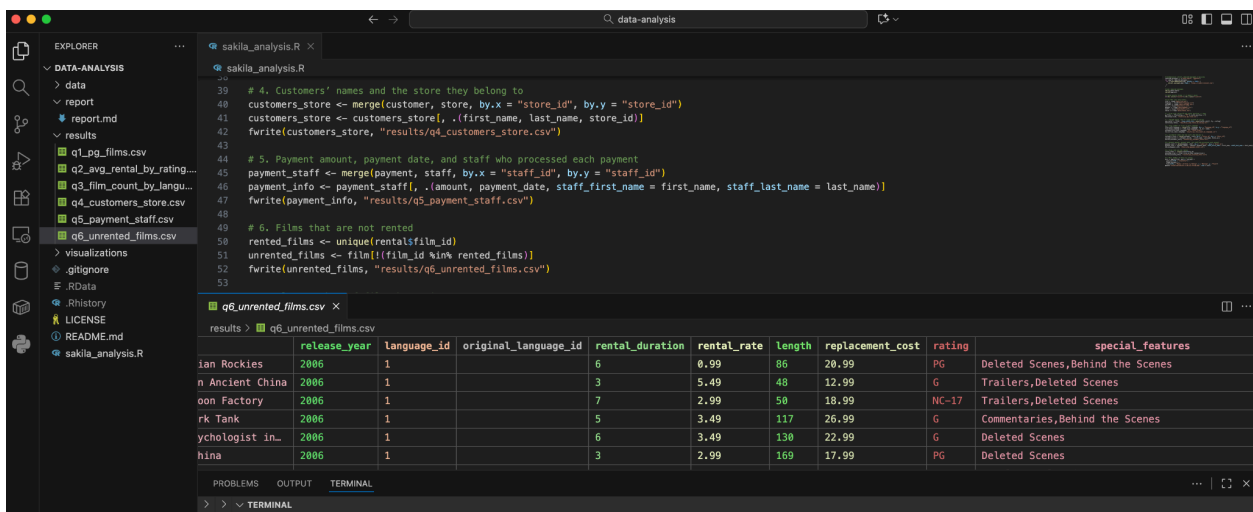**Output File:** results/q5_payment_staff.csv

# 6. Films That Are Not Rented

**Objective:** Identify films that have never been rented.

**R Code:**

```r
rented_films <- unique(rental$film_id)
unrented_films <- film[!(film_id %in% rented_films)]
fwrite(unrented_films, "results/q6_unrented_films.csv")
```

**Output File:** `results/q6_unrented_films.csv`

# 7. Visualization — Number of Films by Rating

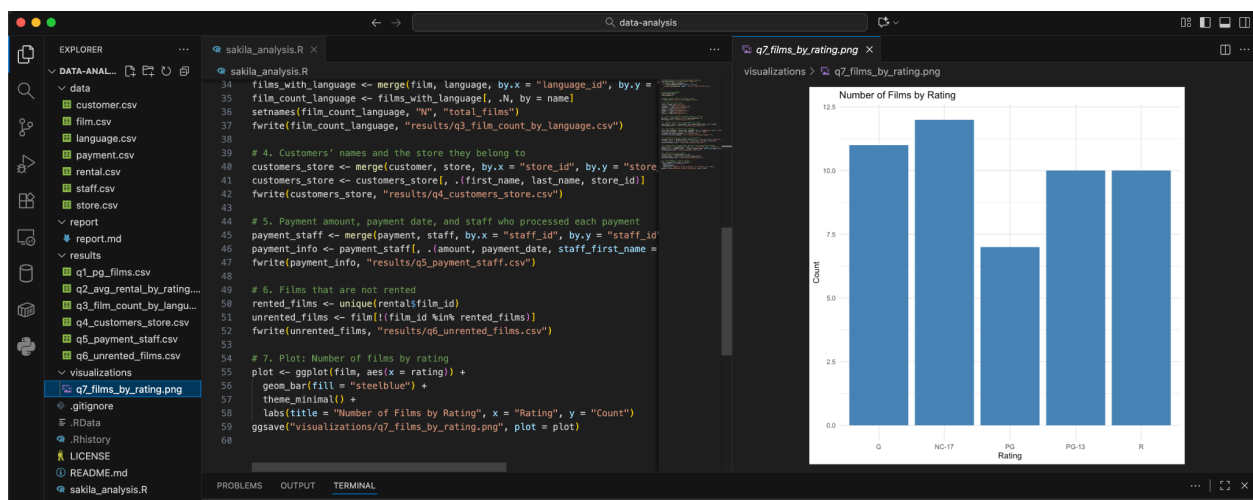**Objective:** Create a bar chart to visualize the number of films per rating.

**R Code:**

```
plot <- ggplot(film, aes(x = rating)) +
  geom_bar(fill = "steelblue") +
  theme_minimal() +
  labs(title = "Number of Films by Rating", x = "Rating", y = "Count")
ggsave("visualizations/q7_films_by_rating.png", plot = plot)
```

**Generated Plot:**
*Number of Films by Rating*
(Stored in `visualizations/q7_films_by_rating.png`)



# 8. Use of Git in the entire assignment.

Git was used throughout the project to manage version control, ensure collaboration, and maintain a clean development workflow. The project repository is publicly available on GitHub.

**Repository Link:** https://github.com/atique-ahmad-01/data-analysis/

**Documentation:**
A detailed project report is available in Markdown format, maintained through Git commits.

**Report Link:** https://github.com/atique-ahmad-01/data-analysis/blob/main/report/report.md

# Conclusion

The Sakila database analysis provided insights into:

- Films with **PG** rating and longer rental durations.
- **Average rental rates** grouped by film rating.
- Film distribution across **languages**.
- Association between **customers and stores**.
- Tracking of **payments and responsible staff**.
- Identification of **unrented films**.
- A visualization summarizing the **film ratings distribution**.

All outputs are saved in the `results/` and `visualizations/` folders for future reference.