**21st July 2024**

**Report By:  Vinita Jain, Shreya Biswas, Evan Cady, Lakshmialekhya Tatampudi**

**IS6480 – Data Warehousing**

**Group 9 – Summer 2024**

# DVD RENTAL DATA WAREHOUSE

# Table of Contents

# EXECUTIVE SUMMARY

For the purposes of this project, our team has decided to pursue the topic of DVD rentals. We worked with a dataset that contains detailed information from several years of records pertaining to customer rental activities, staff, information about each location, and related information.

When working to provide specific insights for our fictitious business, we focused our efforts on the most important functions that a business of this nature would benefit from. In designing our data warehouse, it was important that we do it in the most efficient way possible. We started with the main requirements and which business questions we wanted to answer. Then, we designed dimensional models that would allow us to more accurately implement the data warehouse to suit our needs listed in requirements and matrix. From there, we set up a Power BI workspace which was then shared with all team members. Within this, we created our dataflows and lakehouse.

When we completed construction of our data warehouse, we worked as a team to create visualizations that answer our business questions and requirements. When we finished our visualizations, we had a report that could be effectively turned into an app and then shared with an audience of our choice. Through completing this step, we were able to provide answers to our business questions in a way that someone who has little knowledge of this fictitious company would rapidly gain insight into its functions and performance.

In conclusion, we were able to provide clear answers to our questions relating to rental activities, details of rentals, business performance, and related items and present this in a way that the average person would be able to make sense of.

# BACKGROUND INFORMATION

### a.  Services Provided

Our data warehouse is to be conceived with the knowledge that our primary dataset is representative of a chain of stores that rent DVDs. With that being said, this is a company that provides video rental services through physical media for a fee to consumers, who are then obligated to return such media to a store within a given timeframe or be billed for late fees.

### b.  Organizational Objectives

Our company wants to build a solution that will allow it to more easily make sense of data that is based on its day-to-day operations. It wants to better understand its overarching trends to better meet the needs of its consumers by understanding what their needs are. To help with this, the company also wants to examine performance of employees and see which employees are completing the most transactions to presumably offer programs that motivate employees to perform better and to allocate talent more efficiently. The company will also want to take a look at performance of each location to see where more resources are needed or where there is an opportunity for expansion. Late returns are a key factor impacting profitability, and understanding the costs associated with customers returning items late is important. Among others, the company also wants to make sense of customer spending to possibly offer a rewards program that might boost revenue.

### c.  Why a Data Warehouse is Useful?

Creating a data warehouse will help the company organize the vast amount of data that it has collected from all areas of its business. This will make it easier and more efficient to analyze the data that has been collected. This is especially true based on the historical data that is in the database and the large amount of data that we have. We also can have a more central data system that pulls from several tables and includes only the items that we find necessary to analyze when seeking to provide answers to our business questions.

# REQUIREMENTS

## SUMMARY OF MAIN REQUIREMENTS

| Requirement Name | Short Description | Status |
|---|---|---|
| **Film Rental** | | |
| Films under each Category | Which film category has highest number of Films | Completed |
| Most rented films | Which films have been most rented and least rented | Completed |
| Track stock levels | Check the availability of DVD titles and ensure that popular titles are always in stock. | Incomplete |
| Language analysis | Which language films are rented the most at every store | Incomplete |
| **Customer Analysis** | | |
| Customers Preference | What film category are the customers most interested and least interested in? | Completed |
| Customer demographics | What is the average age of customer interests according to film genre? | Incomplete |
| Loyalty programs | Analyze customer spending over time. This information can be used to give some special offers to customers frequently renting | Incomplete |
| **Management** | | |
| Store rental | What is the average number of rentals at each store per year? | Incomplete |
| Staff Analysis | Which staff member processed the highest total payment amount in each store in each year? | Incomplete |
| **Financial Analysis** | | |
| Income Analysis | What is the trend of the income that we have made over the time period? | Completed |
| Gross Revenue | Which film category has made highest Income and lowest Income? | Completed |

| Sales Peak analysis | What are the sales peak hours and days at particular store | Incomplete |
|---|---|---|
| Returns | Analyze late returns. This insight can be used to calculate revenue impact. | Incomplete |

## BUS MATRIX

| Business processes | Dimensions | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | DimFilm | DimStore | DimLanguage | DimCustomer | DimStaff | DimInventory | DimDate | DimFilmCategory | DimCustomerAge |
| **Film Rental** | | | | | | | | | |
| Films under each Category | x | | | | | | | | |
| Most rented films | x | | | x | | | x | | |
| Track stock levels | x | x | | | | x | x | | |
| Language analysis | x | x | x | | | | x | | |
| **Customer Analysis** | | | | | | | | | |
| Customers Preference | x | | | x | | | x | x | |
| Customer demographics | x | x | | | | | x | x | x |
| Loyalty programs | x | x | | x | | | x | | |
| **Management** | | | | | | | | | |
| Store rental | x | x | | | | | x | | |
| Staff Analysis | x | x | | | x | | x | | |
| **Financial Analysis** | | | | | | | | | |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Income Analysis | x | | | x | | | x | | |
| Gross Revenue | x | | | x | | | x | x | |
| Sales Peak analysis | x | x | | x | | | x | | |
| Returns | x | x | | | | | x | | |

# DIMENSIONAL MODEL



*Figure 1: This model is structured to support querying and reporting on various aspects of the DVD rental business, such as customer activities, rental patterns, and store performance.*

## DIMENSIONAL MODEL DESCRIPTION

**Fact Table**

FactRental:

    a.   rental_date_key: Foreign key to DimDate.
    b.   customer_key: Foreign key to DimCustomer.
    c.   film_key: Foreign key to DimFilm.
    d.   store_key: Foreign key to DimStore.
    e.   rental_date_return_key: Another foreign key to DimDate for the return date.
    f.   inventory_id: Identifier for the inventory.
    g.   payment_date: Date of payment.
    h.   amount: Amount paid.
    i.   rental_id: Identifier for the rental.
    j.   staff_id: Identifier for the staff.
    k.   payment_id: Identifier for the payment.

**Dimension Tables**

DimDate:

    a.   date_key: Primary key for the date dimension.
    b.   date: Full date value.
    c.   dt: Date and time.
    d.   time: Time portion of the date.
    e.   month: Month portion of the date.
    f.   year: Year portion of the date.
    g.   quarter: Quarter portion of the year.
    h.   day_of_week: Day of the week.

DimCustomer:

    a.   customer_key: Primary key for the customer dimension.
    b.   customer_id: Identifier for the customer.
    c.   store_id: Identifier for the store.
    d.   firstname: First name of the customer.
    e.   lastname: Last name of the customer.
    f.   email: Email address of the customer.
    g.   address_id: Identifier for the address.
    h.   activebool: Boolean indicating if the customer is active.
    i.   active: Another indicator for the active status of the customer.

DimStore:

    a.   store_key: Primary key for the store dimension.
    b.   store_id: Identifier for the store.
    c.   manager_staff_id: Identifier for the manager staff.
    d.   address_id: Identifier for the address.

DimFilm:

    a.   film_key: Primary key for the film dimension.
    b.   film_id: Identifier for the film.
    c.   title: Title of the film.
    d.   description: Description of the film.
    e.   release_year: Year the film was released.
    f.   rental_duration: Duration for which the film can be rented.
    g.   rental_rate: Rate at which the film can be rented.
    h.   length: Length of the film.
    i.   film_category: Category of the film.
    j.   film_language: Language of the film.

## Relationships

    a.   fact_most_rented is connected to DimDate via rental_date_key and rental_date_return_key.
    b.   fact_most_rented is connected to DimCustomer via customer_key.
    c.   fact_most_rented is connected to DimFilm via film_key.
    d.   fact_most_rented is connected to DimStore via store_key and staff_key.

# ARCHITECTURE



*Figure 2: Data flows from the Azure Blob Storage (Source Systems) into the Lakehouse (ETL System/Staging Area), where various transformations and stages are managed. From there, the data moves to the Warehouse (Presentation Area) and is used for generating report*

# DATA PIPELINES



*Figure 3: Df_Date : rental table is loaded from azure blob. It is referenced as Silver_rental_date and Silver_rental_return_date. Both the referenced tables are merged into Silver_Silver_date and duplicates are removed. Then DimDate is created from Silver_Silver_Date*

*Figure 4: Df_Store: store table is loaded from azure blob storage and referenced to create DimStore.*

*Figure 5: Df_Customer: customer table is loaded from azure blob storage and referenced to create DimCustomer.*

*Figure 6: Df_Film: film, category, film_category and language tables are loaded from azure blob storage. Film_category and category tables are left joined on category_id to get category name in film_category table(Silver_film_category_bridge).Then film and Silver_film_category_bridge tables are left joined on film_id to get film category name in film table(Silver_Film). DimFilm is then created from Silver_Film table.*
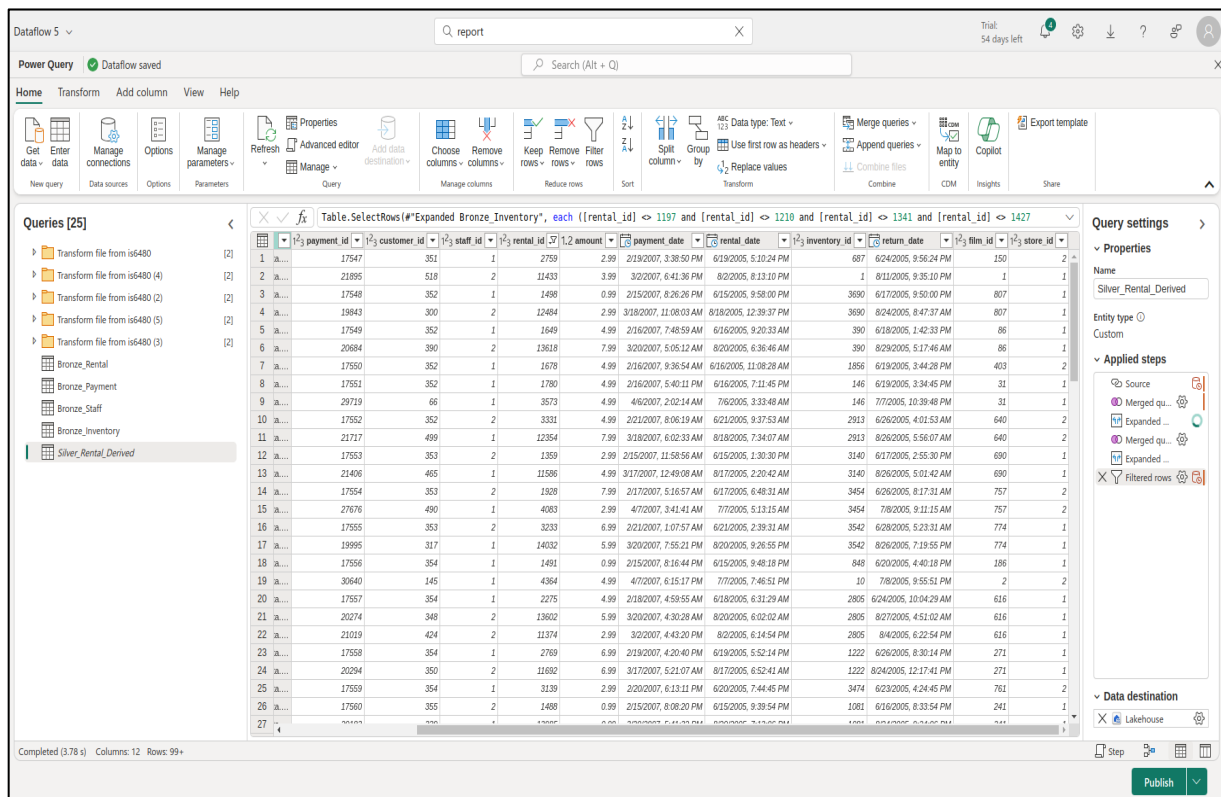
*Figure 7: Df_rental_derived: Silver_rental_derived is created from rental, payment, staff and inventory tables.*

*Figure 8: Df_fact_rental: FactRental table is created from the Silver_Rental_Derived and all the dimension tables created.*
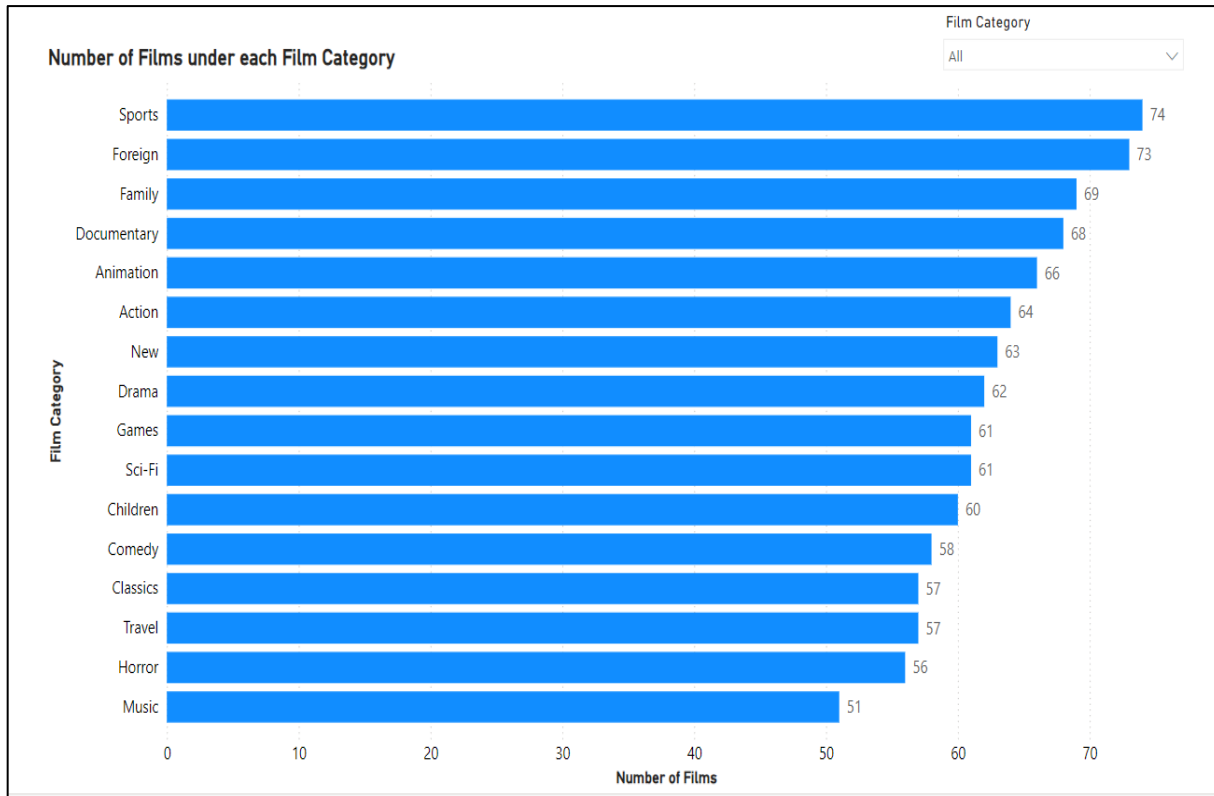
# DASHBOARD AND REPORTS



*Figure 9: This Visual shows the number of films for each film Category. The sports Category has the highest number of Films.*
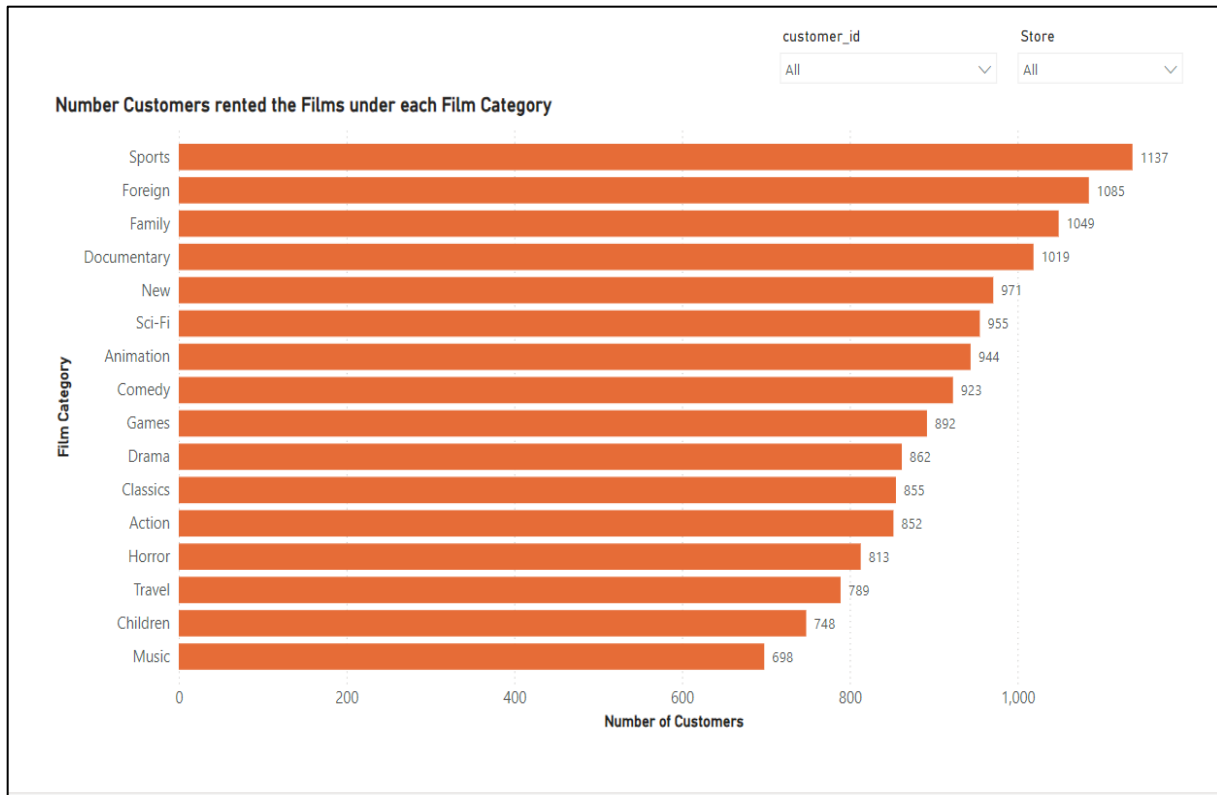
*Figure 10: This Visual shows Number of Customers that have rented the films by Film Category. The highest rented Film category is Sports.*

*Figure 11: This visual shows the highest and least rented films with their Film Categories and Titles.*
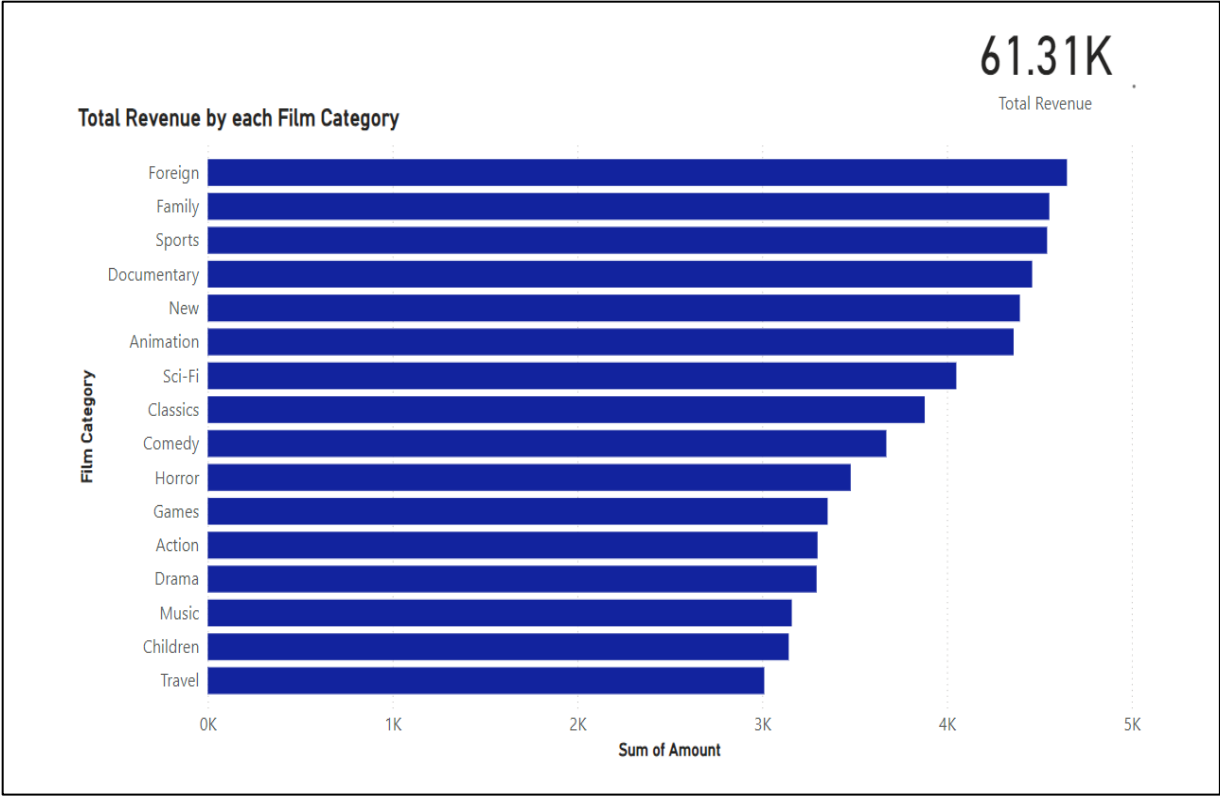
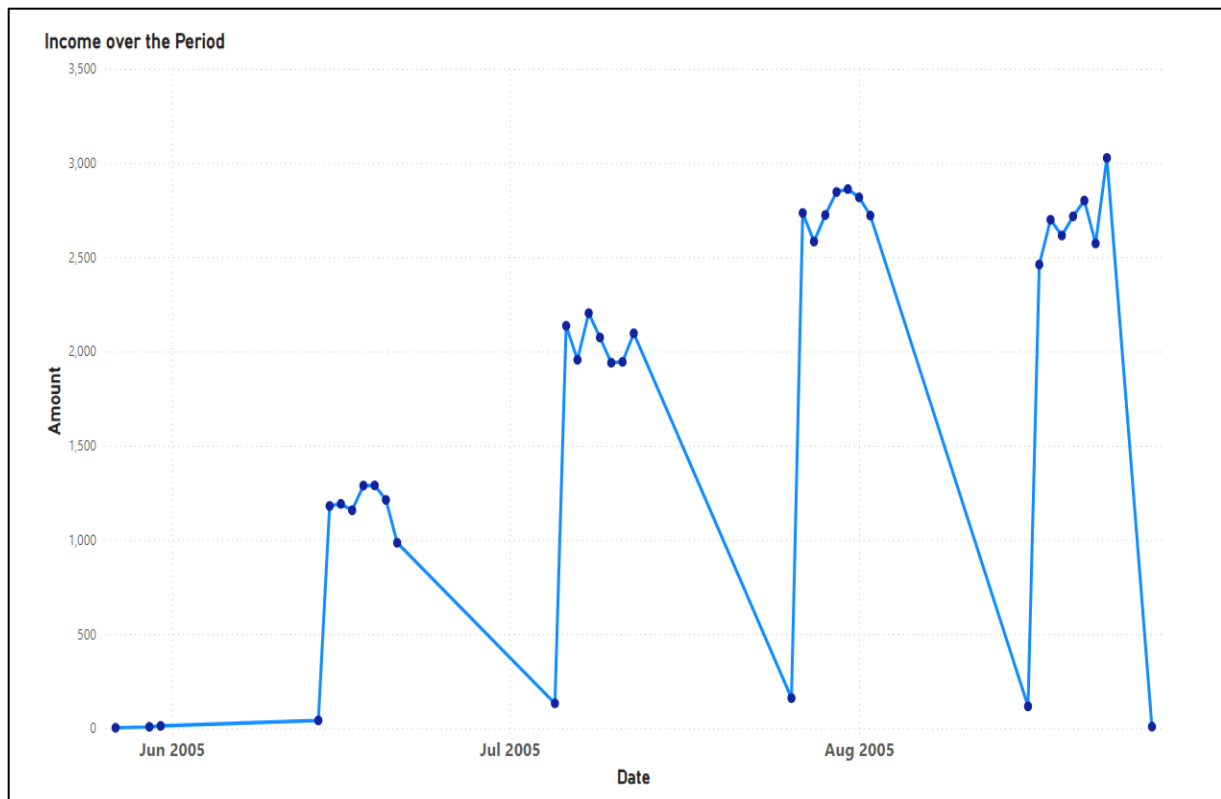*Figure 12: This Visual shows the Total Revenue each Film Category has made.*

*Figure 13: This visual shows the Revenue performance over the period of time.*

# PLANNING AHEAD

1.  Dashboards/reports/analyses that can be supported by the current data warehouse but have yet to be developed

## BELOW MENTIONED ADDITIONAL REPORTS WE CAN CREATE IN FUTURE USING THE CURRENT DATA WAREHOUSE

**Customer Segmentation Analysis**:

a.  Segment customers based on rental behavior, spending patterns, and preferences.
b.  Identify high-value customers and target them with personalized offers.

**Rental Trends and Patterns**:

a.  Analyze daily, weekly, and monthly rental patterns.
b.  Identify peak rental times and seasons.

2.  Dashboards/reports/analyses that would need updates to the dimensional model in order to support

## SUGGESTED DIMENSIONAL MODEL UPDATES:

1.  **DimCustomer**:
    o  Add columns for age, gender, and location.
2.  **DimFilm**:
    o  Add columns for rating, director, and actors.
3.  **DimStore**:
    o  Add detailed address attributes such as city, state, country, and size.
4.  **FactRental**:
    o  Include detailed payment information with attributes like payment_method and payment_status.

# APPENDIX

a. Link to Presentation:
   https://drive.google.com/file/d/1lBdQdjvb7vzmf1G6WCouHVg4QKYj72Lm/view?usp=sharing

b. Power BI Reporting:
   https://app.powerbi.com/Redirect?action=OpenApp&appId=3443816e-f440-43a8-8a7a-108cf100c571&ctid=5217e0e7-539d-4563-b1bf-7c6dcf074f91

c. Data Source and ERD:
   https://www.postgresqltutorial.com/postgresql-getting-started/postgresql-sample-database/