



# NETFLIX MOVIES AND TV SHOWS DATA ANALYSIS USING SQL

# ABOUT

---

## PROJECT OVERVIEW

This project provides a comprehensive analysis of Netflix's movie and TV show data using SQL. The primary goal is to extract valuable insights regarding content types, ratings, geographical distribution, and various other aspects of the Netflix content library. The analysis is structured to address key business questions and deliver meaningful insights based on the dataset.



---

# OBJECTIVES

1. Analyze the distribution of content types (movies vs TV shows).
  2. Identify the most common ratings for movies and TV shows.
  3. Analyze content based on release years, countries, and durations.
  4. Explore and categorize content using specific criteria and keywords.
-

# DATASET

```
1  -- Netflix Project
2
3  ✓ CREATE TABLE netflix
4  (
5      show_id      VARCHAR(5),
6      type         VARCHAR(10),
7      title        VARCHAR(250),
8      director     VARCHAR(550),
9      casts        VARCHAR(1050),
10     country      VARCHAR(550),
11     date_added   VARCHAR(55),
12     release_year INT,
13     rating       VARCHAR(15),
14     duration     VARCHAR(15),
15     listed_in    VARCHAR(250),
16     description  VARCHAR(550)
17 );
18 SELECT * FROM netflix;
19 SELECT COUNT(*) FROM netflix;
20
```

# BUSINESS PROBLEMS AND SQL QUERIES

## 01. COUNT THE NUMBER OF MOVIES VS TV SHOWS

```
-- 01.Count the Number of Movies vs TV Shows  
SELECT type, COUNT(*) as total_content FROM netflix GROUP BY 1;
```

### **PURPOSE:**

- THIS QUERY COUNTS THE TOTAL NUMBER OF CONTENT ITEMS FOR EACH TYPE (MOVIES VS TV SHOWS) ON NETFLIX.
- TO DETERMINE THE DISTRIBUTION OF CONTENT TYPES (MOVIES VS. TV SHOWS) IN THE NETFLIX DATASET.

### **REASONING:**

THIS HELPS IN UNDERSTANDING WHICH TYPE OF CONTENT DOMINATES NETFLIX'S LIBRARY AND AIDS IN STRATEGIC DECISIONS ON CONTENT INVESTMENTS.

---

## 02. FIND THE MOST COMMON RATING FOR MOVIES AND TV SHOWS

```
24 -- 02.Find the Most Common Rating for Movies and TV Shows
25 ✓ WITH RatingCounts AS (
26     SELECT
27         type,
28         rating,
29         COUNT(*) AS rating_count
30     FROM netflix
31     GROUP BY type, rating
32 ),
33 RankedRatings AS (
34     SELECT
35         type,
36         rating,
37         rating_count,
38         RANK() OVER (PARTITION BY type ORDER BY rating_count DESC) AS rank
39     FROM RatingCounts
40 )
41 SELECT
42     type,
43     rating AS most_frequent_rating
44 FROM RankedRatings
45 WHERE rank = 1;
```

### PURPOSE:

- THIS QUERY IDENTIFIES THE MOST FREQUENTLY OCCURRING RATING FOR EACH TYPE OF CONTENT (MOVIES AND TV SHOWS).
- TO IDENTIFY THE MOST FREQUENTLY OCCURRING RATING FOR MOVIES AND TV SHOWS SEPARATELY.

### REASONING:

UNDERSTANDING COMMON RATINGS CAN HELP GAUGE THE CONTENT'S TARGET AUDIENCE AND INFORM DECISIONS ABOUT CONTENT AGE RESTRICTIONS AND CLASSIFICATIONS.

---

---



### 03. FIND THE MOST COMMON RATING FOR MOVIES AND TV SHOWS

```
-- 03.List All Movies Released in a Specific Year (e.g., 2020)
SELECT *
FROM netflix
WHERE release_year = 2020;
```

#### **PURPOSE:**

- TO LIST ALL MOVIES RELEASED IN 2020.
- RETRIEVE ALL MOVIES THAT WERE RELEASED SPECIFICALLY IN THE YEAR 2020.

#### **REASONING:**

USEFUL FOR TREND ANALYSIS, HELPING UNDERSTAND WHAT KIND OF CONTENT WAS RELEASED IN A SPECIFIC YEAR.

---

---

## 04.FIND THE TOP 5 COUNTRIES WITH THE MOST CONTENT ON NETFLIX

```
52 -- 04.Find the Top 5 Countries with the Most Content on Netflix
53 v SELECT *
54 FROM
55 (
56     SELECT
57         UNNEST(STRING_TO_ARRAY(country, ',')) AS country,
58         COUNT(*) AS total_content
59     FROM netflix
60     GROUP BY 1
61 ) AS t1
62 WHERE country IS NOT NULL
63 ORDER BY total_content DESC
64 LIMIT 5;
```

### PURPOSE:

TO IDENTIFY THE TOP 5 COUNTRIES THAT HAVE CONTRIBUTED THE MOST CONTENT TO NETFLIX.

### REASONING:

UNDERSTANDING GEOGRAPHICAL CONTENT CONTRIBUTIONS HELPS IN ANALYZING REGIONAL CONTENT STRATEGIES AND THE DIVERSITY OF THE NETFLIX LIBRARY.

---



## 05. IDENTIFY THE LONGEST MOVIE

```
66
67  -- 05.Identify the Longest Movie
68  ✓ SELECT
69      *
70  FROM netflix
71  WHERE type = 'Movie'
72  ORDER BY SPLIT_PART(duration, ' ', 1)::INT DESC;
73
```

### PURPOSE:

FIND THE LONGEST MOVIE BASED ON ITS DURATION.

### REASONING:

IDENTIFYING THE LONGEST MOVIE PROVIDES INSIGHTS INTO NETFLIX'S INVESTMENT IN LENGTHY CONTENT, WHICH MIGHT CATER TO SPECIFIC AUDIENCE PREFERENCES.

---

## 06. FIND CONTENT ADDED IN THE LAST 5 YEARS

```
74
75 -- 06.Find Content Added in the Last 5 Years
76 ✓ SELECT *
77 FROM netflix
78 WHERE TO_DATE(date_added, 'Month DD, YYYY') >= CURRENT_DATE - INTERVAL '5 years';
79
80
```

### **PURPOSE:**

TO RETRIEVE CONTENT THAT HAS BEEN ADDED TO NETFLIX IN THE LAST FIVE YEARS..

### **REASONING:**

ANALYZING RECENTLY ADDED CONTENT HELPS IN UNDERSTANDING NETFLIX'S RECENT CONTENT ACQUISITION TRENDS AND STRATEGIES.

---

## 07. FIND ALL MOVIES/TV SHOWS BY DIRECTOR 'RAJIV CHILAKA'

```
80
81 -- 07.Find All Movies/TV Shows by Director 'Rajiv Chilaka'
82 v SELECT *
83 FROM (
84     SELECT
85         *,
86         UNNEST(STRING_TO_ARRAY(director, ',')) AS director_name
87     FROM netflix
88 ) AS t
89 WHERE director_name = 'Rajiv Chilaka';
90
```

### **PURPOSE:**

TO LIST ALL MOVIES AND TV SHOWS DIRECTED BY RAJIV CHILAKA.

### **REASONING:**

THIS QUERY SPLITS THE DIRECTOR FIELD INTO INDIVIDUAL NAMES IN CASE THERE ARE MULTIPLE DIRECTORS, AND THEN IT FILTERS TO FIND CONTENT DIRECTED BY 'RAJIV CHILAKA'. USING UNNEST() ENSURES THAT CONTENT WITH MULTIPLE DIRECTORS IS ALSO CORRECTLY IDENTIFIED.

---

## 08.LIST ALL TV SHOWS WITH MORE THAN 5 SEASONS

```
91
92  -- 08.List All TV Shows with More Than 5 Seasons
93  ✓ SELECT *
94  FROM netflix
95  WHERE type = 'TV Show'
96  AND SPLIT_PART(duration, ' ', 1)::INT > 5;
97
```

**PURPOSE:**

TO RETRIEVE ALL TV SHOWS WITH MORE THAN FIVE SEASONS.

**REASONING:**

THE QUERY EXTRACTS THE NUMERIC PART FROM THE DURATION FIELD (E.G., '6 SEASONS'), CONVERTS IT TO AN INTEGER, AND THEN CHECKS WHETHER IT'S GREATER THAN 5, ENSURING THE RETRIEVAL OF TV SHOWS WITH MORE THAN 5 SEASONS.

---

## 09.COUNT THE NUMBER OF CONTENT ITEMS IN EACH GENRE

```
98
99 -- 09.Count the Number of Content Items in Each Genre
100 ✓ SELECT
101     UNNEST(STRING_TO_ARRAY(listed_in, ',')) AS genre,
102     COUNT(*) AS total_content
103 FROM netflix
104 GROUP BY 1;
105
```

### **PURPOSE:**

TO COUNT HOW MANY CONTENT ITEMS BELONG TO EACH GENRE.

### **REASONING:**

THE LISTED\_IN FIELD MAY CONTAIN MULTIPLE GENRES (E.G., 'ACTION, DRAMA'). THE QUERY SPLITS THESE VALUES INTO INDIVIDUAL GENRES USING UNNEST() AND COUNTS THE TOTAL NUMBER OF CONTENT ITEMS FOR EACH GENRE.

---

## 10.FIND EACH YEAR AND THE AVERAGE NUMBERS OF CONTENT RELEASE IN INDIA ON NETFLIX

```
106
107 -- 10.Find each year and the average numbers of content release in India on netflix
108 v SELECT
109     country,
110     release_year,
111     COUNT(show_id) AS total_release,
112     ROUND(
113         COUNT(show_id)::numeric /
114         (SELECT COUNT(show_id) FROM netflix WHERE country = 'India')::numeric * 100, 2
115     ) AS avg_release
116 FROM netflix
117 WHERE country = 'India'
118 GROUP BY country, release_year
119 ORDER BY avg_release DESC
120 LIMIT 5;
```

### **PURPOSE:**

TO IDENTIFY THE TOP FIVE YEARS WITH THE MOST CONTENT RELEASES IN INDIA.

### **REASONING:**

THE QUERY FILTERS CONTENT PRODUCED IN INDIA, GROUPS IT BY THE RELEASE YEAR, AND COUNTS THE NUMBER OF RELEASES PER YEAR. THE TOP 5 YEARS WITH THE HIGHEST RELEASE COUNT ARE RETURNED.

---

## 11.LIST ALL MOVIES THAT ARE DOCUMENTARIES

```
122
123  -- 11.List All Movies that are Documentaries
124  ✓ SELECT *
125  FROM netflix
126  WHERE listed_in LIKE '%Documentaries';
127
128
```

### **PURPOSE:**

TO LIST ALL MOVIES CLASSIFIED AS DOCUMENTARIES.

### **REASONING:**

THIS QUERY CHECKS WHETHER THE LISTED\_IN FIELD (WHICH CONTAINS GENRE INFORMATION) INCLUDES THE WORD 'DOCUMENTARIES'. IT RETRIEVES ALL CONTENT THAT FALLS UNDER THIS GENRE.

---



## 12.FIND ALL CONTENT WITHOUT A DIRECTOR

```
-- 12.Find All Content Without a Director
✓ SELECT *
  FROM netflix
 WHERE director IS NULL;
```

### **PURPOSE:**

TO LIST ALL CONTENT THAT DOES NOT HAVE A DIRECTOR ASSOCIATED WITH IT.

### **REASONING:**

THIS QUERY SIMPLY FILTERS OUT ALL CONTENT WHERE THE DIRECTOR FIELD IS NULL, HELPING TO IDENTIFY CONTENT WITH MISSING DIRECTOR INFORMATION.

---

### 13.FIND HOW MANY MOVIES ACTOR 'SALMAN KHAN' APPEARED IN THE LAST 10 YEARS

```
135 -- 13.Find How Many Movies Actor 'Salman Khan' Appeared in the Last 10 Years
136 v SELECT *
137 FROM netflix
138 WHERE casts LIKE '%Salman Khan%'
139 AND release_year > EXTRACT(YEAR FROM CURRENT_DATE) - 10;
140
```

#### **PURPOSE:**

TO COUNT THE NUMBER OF MOVIES FEATURING 'SALMAN KHAN' RELEASED IN THE PAST 10 YEARS.

#### **REASONING:**

THE QUERY CHECKS IF 'SALMAN KHAN' IS LISTED IN THE CASTS FIELD AND FILTERS THE CONTENT BASED ON THE RELEASE YEAR BEING WITHIN THE LAST 10 YEARS. THIS HELPS TRACK THE ACTOR'S PRESENCE ON THE PLATFORM.

---

## 14. FIND THE TOP 10 ACTORS WHO HAVE APPEARED IN THE HIGHEST NUMBER OF MOVIES PRODUCED IN INDIA

```
141 -- 14. Find the Top 10 Actors Who Have Appeared in the Highest Number of Movies Produced in Ind
142 ✓ SELECT
143     UNNEST(STRING_TO_ARRAY(casts, ',')) AS actor,
144     COUNT(*)
145 FROM netflix
146 WHERE country = 'India'
147 GROUP BY actor
148 ORDER BY COUNT(*) DESC
149 LIMIT 10;
```

### **PURPOSE:**

TO IDENTIFY THE TOP 10 ACTORS WITH THE MOST APPEARANCES IN MOVIES PRODUCED IN INDIA.

### **REASONING:**

THE CASTS FIELD MAY CONTAIN MULTIPLE ACTORS, SO THE QUERY SPLITS THESE ENTRIES INTO INDIVIDUAL NAMES USING UNNEST(). IT THEN GROUPS BY ACTOR AND COUNTS THEIR APPEARANCES IN INDIAN MOVIES, RETURNING THE TOP 10 ACTORS WITH THE HIGHEST COUNTS.

---

## 15.CATEGORIZE CONTENT BASED ON THE PRESENCE OF 'KILL' AND 'VIOLENCE' KEYWORDS

```
152 -- 15.Categorize Content Based on the Presence of 'Kill' and 'Violence' Keywords
153 v SELECT
154     category,
155     COUNT(*) AS content_count
156 FROM (
157     SELECT
158         CASE
159             WHEN description ILIKE '%kill%' OR description ILIKE '%violence%' THEN 'Bad'
160             ELSE 'Good'
161         END AS category
162     FROM netflix
163 ) AS categorized_content
164 GROUP BY category;
```

### PURPOSE:

TO CATEGORIZE CONTENT AS 'GOOD' OR 'BAD' BASED ON THE PRESENCE OF SPECIFIC KEYWORDS LIKE 'KILL' OR 'VIOLENCE' IN THE CONTENT DESCRIPTION.

### REASONING:

THIS QUERY CHECKS THE DESCRIPTION FIELD FOR KEYWORDS SUCH AS 'KILL' OR 'VIOLENCE'. IF ANY OF THESE WORDS ARE FOUND, THE CONTENT IS CLASSIFIED AS 'BAD'; OTHERWISE, IT IS CLASSIFIED AS 'GOOD'. THE RESULTS ARE THEN GROUPED AND COUNTED BASED ON THESE CLASSIFICATIONS.

---

# CONCLUSION

THIS SQL ANALYSIS OF NETFLIX'S DATASET REVEALED KEY INSIGHTS:



- Content Type Distribution: Grouping by type showed the balance between movies and TV shows, aiding in content strategy decisions.
  - Rating Insights: Using RANK() revealed the most frequent ratings, helping understand Netflix's target audience.
  - Temporal Analysis: Filtering by release\_year provided insights into content trends over specific years.
  - Geographical Distribution: UNNEST() allowed us to identify the top content-producing countries, showing regional diversity.
  - Longest Content and Seasons: Parsing duration identified the longest movies and TV shows with over 5 seasons.
  - Recent Additions: Date functions (TO\_DATE()) highlighted content added in the last 5 years.
  - Director/Actor Insights: UNNEST() revealed content associated with specific directors and actors like Rajiv Chilaka and Salman Khan.
  - Genre and Keyword Categorization: Genre analysis and keyword filtering categorized content for genre diversity and sensitivity.
-



# GOT ANY QUESTIONS?

These SQL queries provided clear insights into Netflix's content catalog, supporting strategic decision-making on content development and acquisition.