# Netflix Data Analysis

# Agenda

## Topics Covered

Data Cleaning with python

Database Connection with MySql

Data manuplation by power bi Dax
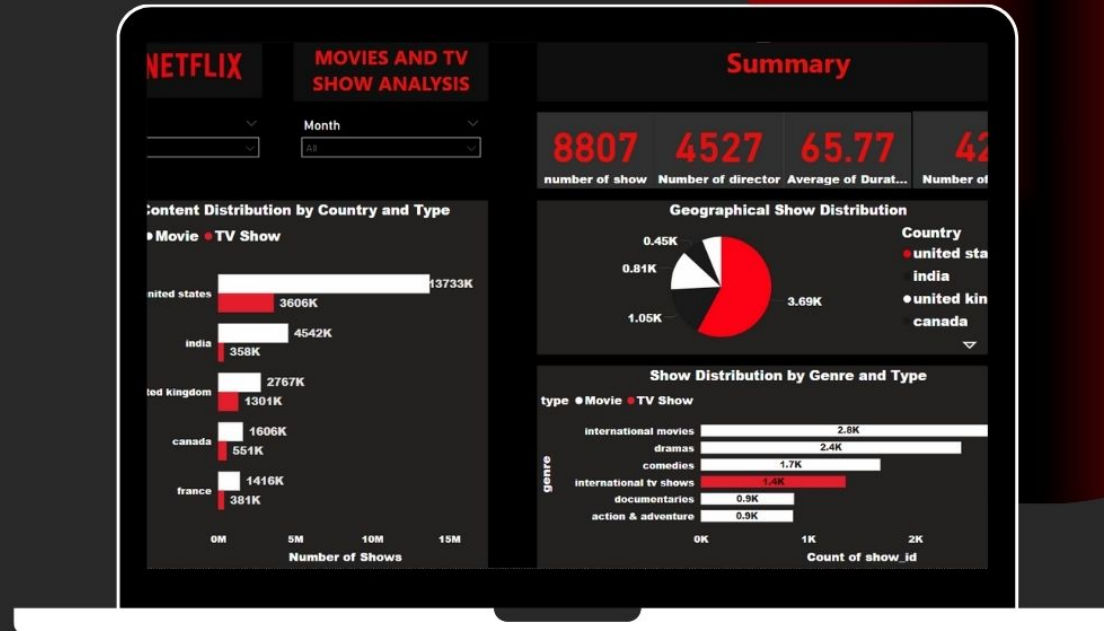
Data Visualization by power bi

## Project Objective:

- **Primary Goal:** To gain insights into viewer preferences and trends on Netflix by analyzing show types, genres, and viewer ratings across different regions.
- **Secondary Goals:** To understand how these trends can influence content creation and marketing strategies for Netflix.

- ## Tools and Methods Used

- Mention the use of Python for data cleaning and manipulation, SQL for data retrieval and management, and Power BI for visualization.

# Data cleaning by python

## What is Data Cleaning

Briefly explain the importance of data cleaning in ensuring accurate, reliable analysis. Mention that clean data helps in making sound decisions based on the analysis.

## Challenges in the Netflix Dataset

Inconsistent Formats: Describe how data entries like dates and show IDs had inconsistent formats that required standardization.

Missing Information: Highlight issues like missing director or cast details and how they impact analysis.

Duplicate Records: Point out any duplicate entries found in the dataset and the need to remove them to avoid skewed results.

# Data cleaning by python

## Import Library

```
import pandas as pd
import numpy as np
```

Pandas: Pandas is a powerful and popular open-source Python library specifically designed for data manipulation and analysis.

Numpy:a Python library that provides support for multidimensional arrays and mathematical functions for scientific computing. It's short for "numerical Python"

# Data cleaning by python

## Steps Taken for Data Cleaning

Standardization: Discuss the methods used to standardize data formats, such as converting all dates to a single format. Changing data type by "astype" and "Dtype"

## Handling Missing Values

```python
1  data_netflix_dummy['country']=data_netflix_dummy['country'].fillna('country not here')
2  data_netflix_dummy.loc[data_netflix.country == 'country not here']
3  data_netflix_dummy.loc[data_netflix_dummy.cast == 'cast not available']
4  data_netflix_dummy['cast']=data_netflix['cast'].fillna('cast not available')
5  data_netflix_dummy['director']=data_netflix_dummy['director'].fillna('director name not available')
6  data_netflix_dummy.loc[data_netflix_dummy.director == 'director name not available'].head(1)
7  # data_netflix_dummy['date_added']=data_netflix_dummy['date_added'].astype('datetime64[s]')
8  data_netflix_dummy['date_added'].isna().sum()
```

# Data cleaning by python

## Normlization

 the process of organizing data in a database to make it more flexible and protect it. It involves: Creating tables, Establishing relationships between tables, Eliminating redundancy, and Removing inconsistent dependency

## Problem in My Data Set

**Genre Column**

Action & Adventure                                                    128
Action & Adventure, Anime Features                              2
Action & Adventure, Anime Features, Children & Family Movies    12
Action & Adventure, Anime Features, Classic Movies            6
Action & Adventure, Anime Features, Horror Movies

As you see above a  Action & Adventure  are  repeated in many rows but when we do group by  where as Action & Adventure  is single  its count not with other so solved this.

# Data cleaning by python

## Solution Of Multivalued Column

Making other table is solution of multivalued is called 1 NF

# Data cleaning by python

## Same as in country column

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 7 | 8 | Movie | Sankofa | Haile Gerima | Kofi Ghanaba, Oyafunmike Ogunlano, Alexandra D... | United States, Ghana, Burkina Faso, United Kin... | |

| Country | |
|---|---|
| afghanistan | 1 |
| albania | 1 |
| algeria | 3 |
| angola | 1 |
| argentina | 91 |
| armenia | 1 |
| australia | 160 |

# Database Connection with MySql

```python
import sqlalchemy
```

```python
engine = sqlalchemy.create_engine('mysql+pymysql://root:Qwer!234@localhost:3306/netflix')
# data_netflix_dummy.to_sql(name= 'root',con=,schema='netflix')
```

# Connection with MySql

```python
1  data_netflix_project.to_sql('netflix_data', engine, if_exists='replace', index=False)
```

8807

+ Code      + Markdown

```python
1  data_netflix_genres.to_sql('netflix_genres', engine, if_exists='replace', index=False)
```

19323

```python
1  netflix_country.to_sql('country',engine,if_exists='replace',index=False)
```
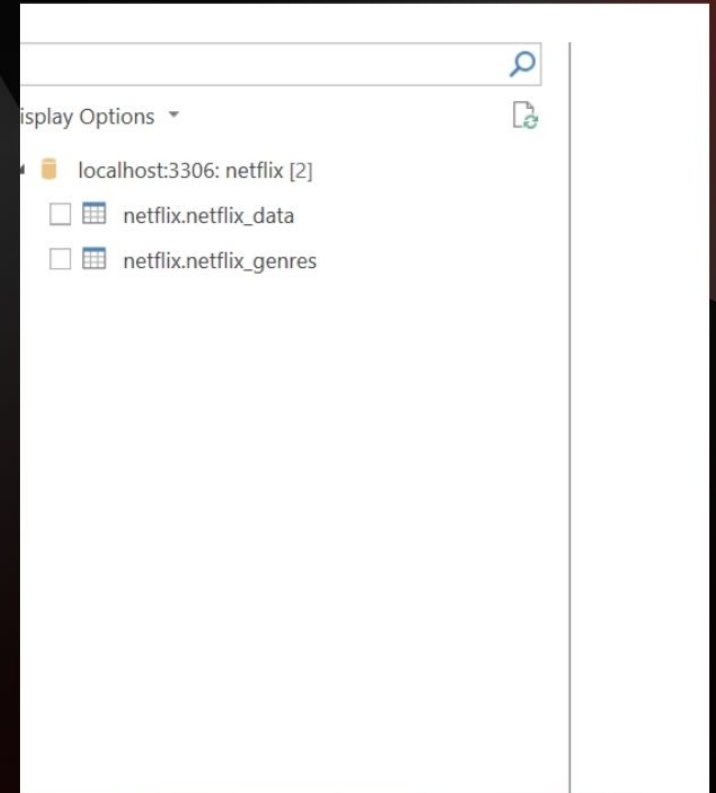
# Data manuplation by power bi Dax

# Connection Sql with Power Bi

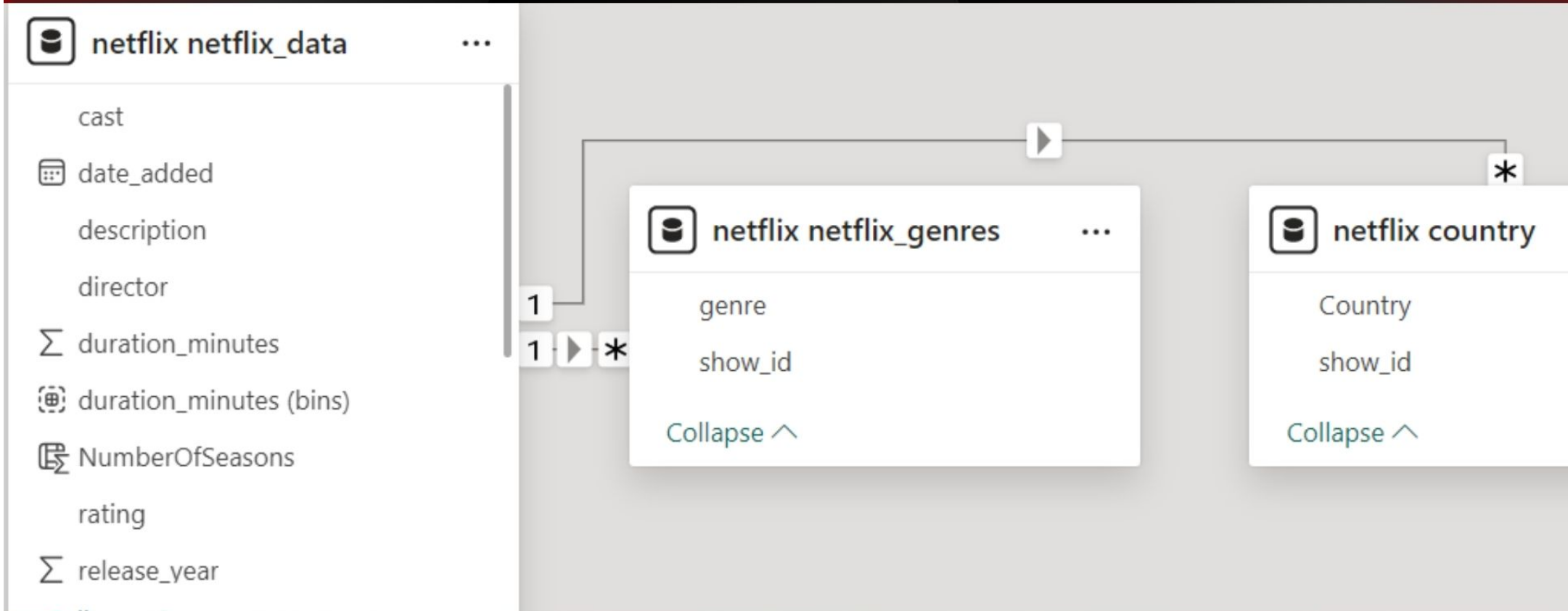Benifit of these changing in data set or real time data set we work easily.

# Data manuplation by power bi Dax

# Relationship between tables

# NETFLIX

## MOVIES AND TV SHOW ANALYSIS

## Summary

**Year**
All

**Month**
All

| 8807 | 4527 | 65.77 | 42 |
|---|---|---|---|
| number of show | Number of director | Average of Durat... | Number of genre |

### Content Distribution by Country and Type

type ● Movie ● TV Show



**Country**

- united states: 13733K (Movie), 3606K (TV Show)
- india: 4542K (Movie), 358K (TV Show)
- united kingdom: 2767K (Movie), 1301K (TV Show)
- canada: 1606K (Movie), 551K (TV Show)
- france: 1416K (Movie), 381K (TV Show)

Number of Shows

### Geographical Show Distribution

**Country**
- ● united states
- india
- ● united kingdom
- canada

0.45K
0.81K
1.05K
3.69K

### Show Distribution by Genre and Type

type ● Movie ● TV Show

| genre | Count of show_id |
|---|---|
| international movies | 2.8K |
| dramas | 2.4K |
| comedies | 1.7K |
| international tv shows | 1.4K |
| documentaries | 0.9K |
| action & adventure | 0.9K |

# NETFLIX

# Analysis Data By Duration

## Distribution of Content Length by Type



Count of type

- 3K — 2.7K
- 2K
- 1K — 1.3K, 0.9K
- 0.0K 0.1K 0.3K 0.4K 0.1K 0.0K 0.0K
- 0K

duration_minutes (bins): 0, 50, 100, 150, 200

## Total Runtime by Rating Category



Sum of duration_minu...

- 200K
- 150K
- 100K
- 50K
- 0K

rating: PG-13, R, TV-14, TV-MA, TV-PG

## Number of Titles per Rating



| rating | Count of show_id |
|---|---|
| TV-MA | 3207 |
| TV-14 | 2160 |
| TV-PG | 863 |
| R | 799 |
| PG-13 | 490 |
| TV-Y7 | 334 |
| TV-Y | 307 |
| PG | 287 |

Count of show_id: 0, 500, 1,000, 1,500, 2,000, 2,500, 3,000

## Average Content Length by Type



1.8

99.5

type
- ● Movie
- ● TV Show

# NETFLIX

Year
All

Month
All

## Top Directors by Number of Shows

director

| Director | Number of Shows |
|----------|-----------------|
| Rajiv Chilaka | 19 |
| Raúl Campos, Jan Suter | 18 |
| Marcus Raboy | 16 |
| Suhas Kadav | 16 |
| Jay Karas | 14 |

Number of Shows
0  5  10  15  20

## Top Cast Members by Number of Shows

type ●Movie ●TV Show

cast

| Cast | Movie | TV Show |
|------|-------|---------|
| David Attenborough | 5 | 14 |
| Samuel West | 10 | |
| Craig Sechler | 6 | |
| Jay O. Sanders | 4 | |
| Sonal Kaushal, Rup... | 3 | 1 |

Number of Shows
0  5  10  15  20

## Show Distribution by Genre

2.68K
6.13K
6.13K
6.13K
6.13K

genre
● comedies
● documentaries
● dramas
● international ...
● international t...

## Show Count by Type and Genre

type

| Type | Count |
|------|-------|
| Movie | 13.2K |
| TV Show | 6.1K |

0K  5K  10K  15K

Genre

Country   ✕     type   ✕     title   ✕     rating   ✕     💡 season

united states     TV Show     Sin senos sí hay paraíso     TV-MA

**united states**
3690

**Movie**
2752

**Sin senos sí hay para...**
2

**TV-MA**
2

**3 Seasons**
2

**india**
1046

**TV Show**
938

**#blackAF**
1

**country not here**
831

**(Un)Well**
1

**united kingdom**
806

**100 Humans**
1

**Count of show_id**
10850

**canada**
445

**13 Reasons Why**
1

**france**
393

**13 Reasons Why: Bey...**
1

**japan**
318

**1983**
1

**spain**
232

**30 Rock**
1

# Insights from Netflix Data Analysis

## Dominance of United States Content:

The United States offers significantly more content compared to other countries, showcasing its leading role in the content availability on Netflix.

## Preference for Movies Over TV Shows in the US:

Analysis data by duration using different charts.

## Popularity of International Movies:

International movies constitute a substantial portion of Netflix's catalog, highlighting its diverse global audience appeal.

## Notable Figures:

- David Attenborough is prominently featured across various content, emphasizing a strong viewer interest in nature and documentary genres.
- Rajiv Chilaka holds the title for directing the most shows

# Thank you!

Kush Shukla