

# Netflix Data Analysis:

## Comprehensive EDA Report

--Kshitij Thorat

### Executive Summary

This document presents a detailed analysis of Netflix's content catalog using PySpark on Databricks Serverless. The analysis examines 4,103 titles spanning from 1925 to 2021, providing insights into content growth patterns, distribution characteristics, and global content production trends.

## 1. Dataset Overview

### 1.1 Data Structure

The Netflix dataset contains **4,103 records** across **16 columns**, encompassing both movies and TV shows with comprehensive metadata:

#### Core Attributes:

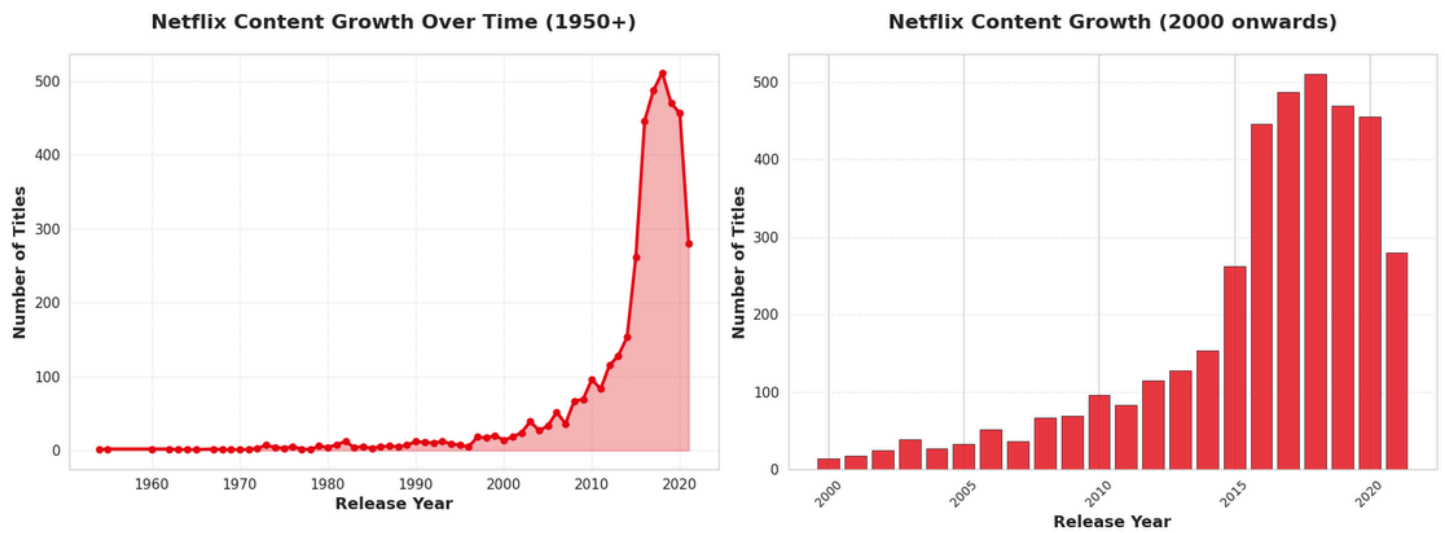
- Identifiers:** show\_id, title, type
- Content Details:** director, cast, description, duration (parsed into duration\_int and duration\_type)
- Classification:** rating, listed\_in (genres)
- Temporal Data:** release\_year, date\_added (parsed into added\_year and added\_month)
- Geographic:** country

#### Data Quality Observations:

- Well-structured schema with appropriate data types (integers for years, strings for categorical data)
- Presence of "Unknown" values in director and country fields indicates incomplete metadata
- Clean separation of duration into numeric and type components facilitates analysis

## 2. Content Growth Over Time Analysis

### 2.1 Release Year Trends



## Key Findings:

1. **Content span:** 96 years (1925-2021)
2. **Peak production year:** 2018 with 511 titles
3. **Average content per year:** 59.5 titles (high standard deviation of 128.2 indicates exponential recent growth)

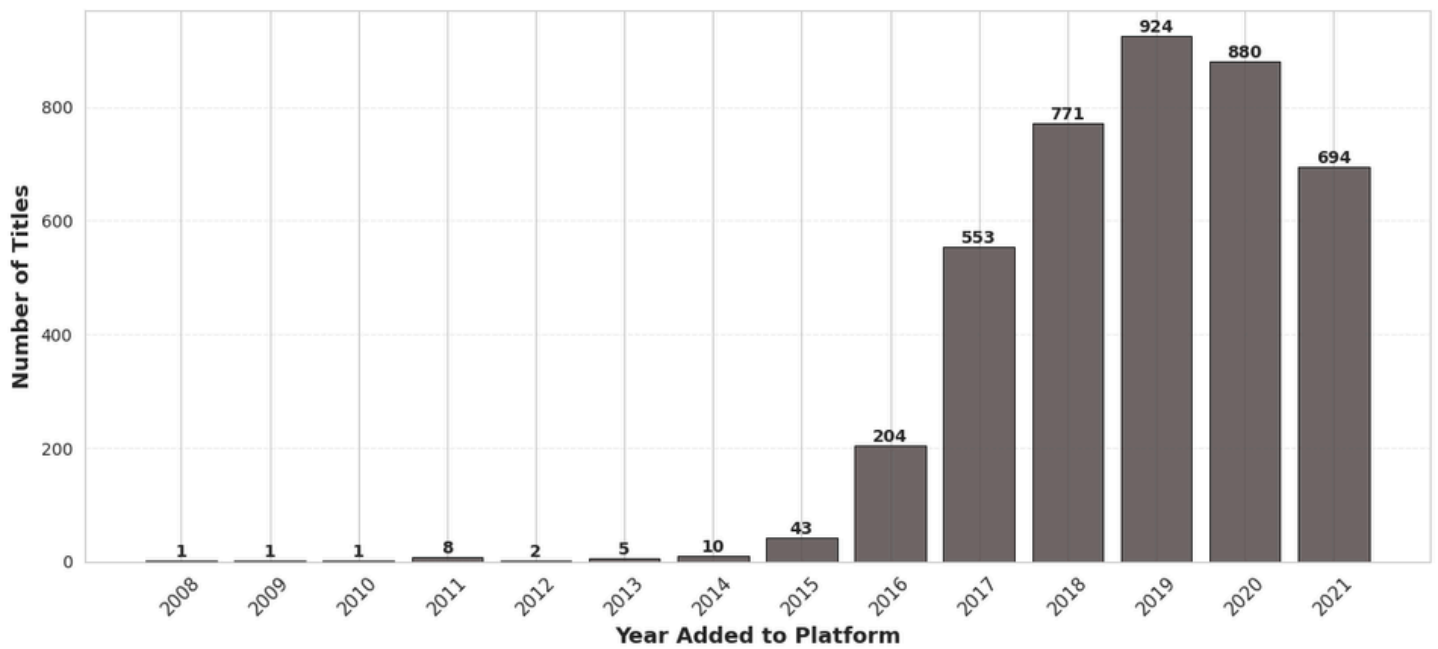
**Growth Pattern Insights:** The data reveals a dramatic acceleration in content production:

1. **Pre-2000 Era:** Minimal representation (< 50 titles annually)
2. **2000-2014:** Gradual increase reflecting early digital content
3. **2015-2018:** Explosive growth period (262 → 511 titles, +95% increase)
4. **2019-2021:** Sustained high production (470, 456, 280 titles)

**Critical Insight:** The 2015-2018 period represents Netflix's aggressive content acquisition and original production phase, coinciding with their global expansion strategy. The drop in 2021 (280 titles) likely reflects partial year data at the time of dataset creation.

## 2.2 Platform Addition Timeline

**Content Added to Netflix by Year**



#### Year-over-Year Growth Analysis:

##### Peak Growth Years:

- 2011: 700% growth (8 titles from 1)
- 2015: 330% growth (43 titles from 10)
- 2016: 374% growth (204 titles from 43)
- 2017: 171% growth (553 titles from 204)

##### Maturation Phase:

1. 2018-2019: Sustained addition of 771 and 924 titles respectively
2. 2020: First decline (-4.8%) to 880 titles
3. 2021: Further decline (-21.1%) to 694 titles

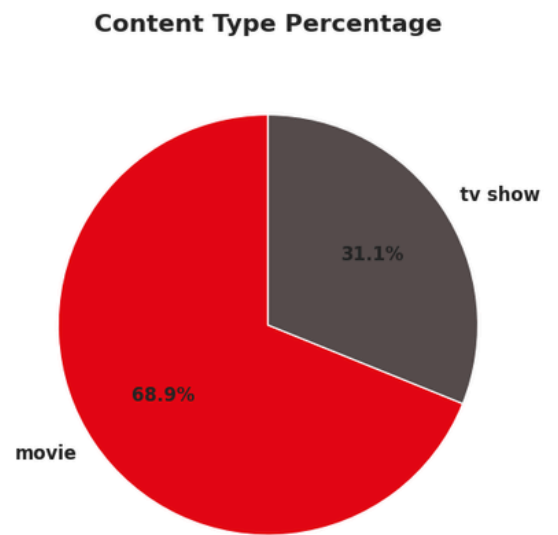
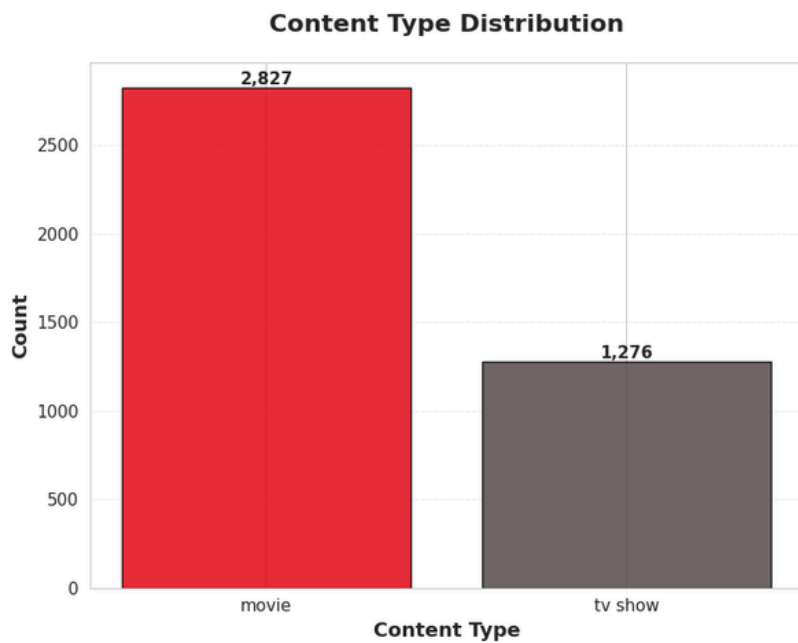
**Strategic Interpretation:** The deceleration in 2020-2021 suggests a shift from aggressive catalog expansion to quality curation and original content focus. This aligns with industry reports of Netflix prioritizing profitable, high-engagement content over volume.

## 3. Distribution Analysis

### 3.1 Content Type Distribution

#### Movies vs. TV Shows:

1. Movies: 2,827 titles (68.9%)
2. TV Shows: 1,276 titles (31.1%)



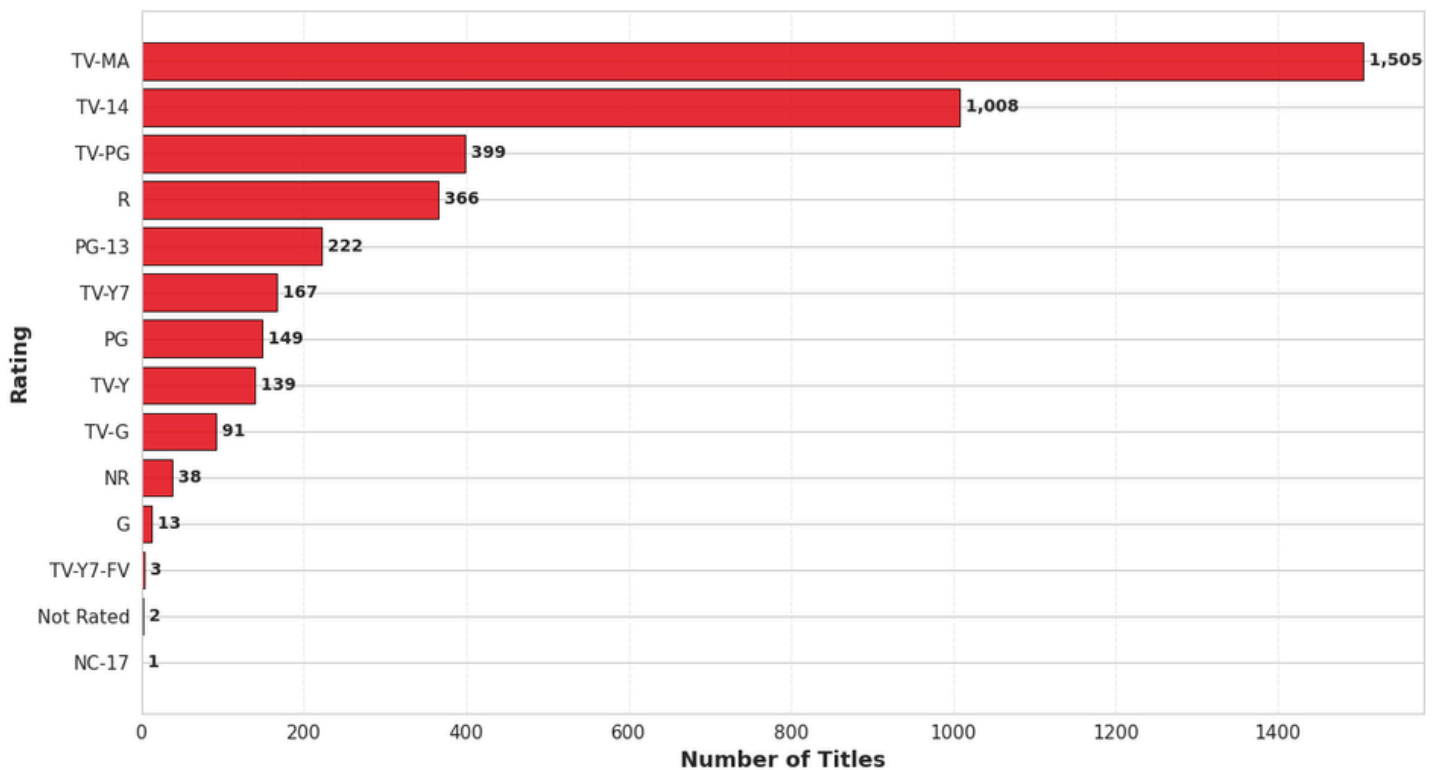
**Insight:** The 2:1 movie-to-TV show ratio reflects Netflix's foundation as a movie streaming service. However, with TV shows comprising nearly one-third of the catalog, this demonstrates significant investment in serialized content, which drives higher subscriber retention through binge-watching behavior.

## 3.2 Rating Distribution

### Top 5 Ratings:

1. **TV-MA (Mature Audiences): 1,505 titles (36.7%)**
2. **TV-14 (Ages 14+): 1,008 titles (24.6%)**
3. **TV-PG (Parental Guidance): 399 titles (9.7%)**
4. **R (Restricted): 366 titles (8.9%)**
5. **PG-13: 222 titles (5.4%)**

**Distribution of Content Ratings (Top 15)**



### Critical Insights:

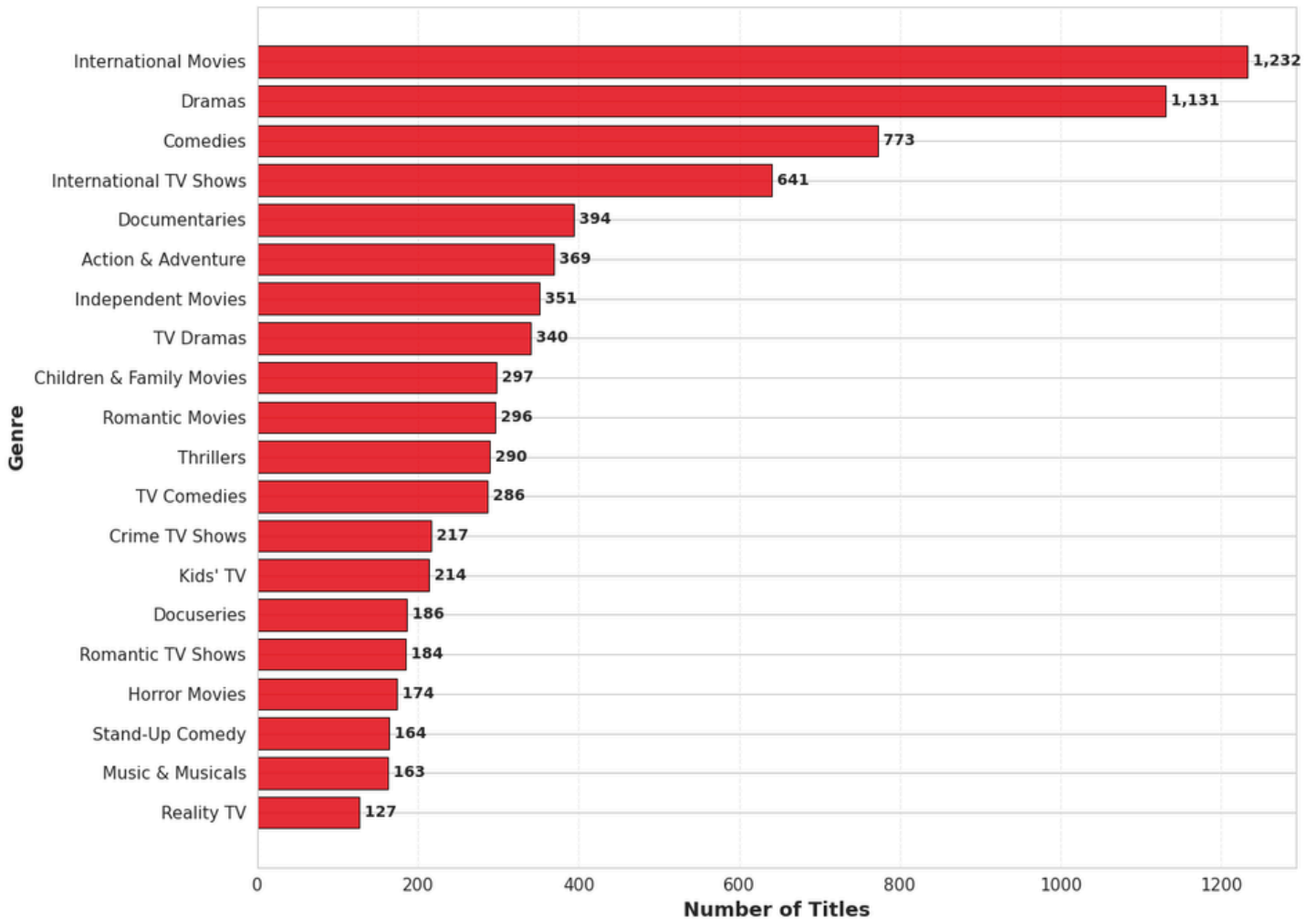
1. 61.3% of content is rated TV-MA or TV-14, indicating a strong skew toward adult audiences
2. Only 12.3% of content is family-friendly (TV-Y, TV-G, G, PG combined)
3. This distribution reflects Netflix's positioning as a premium service targeting 18-49 demographic
4. The dominance of mature content distinguishes Netflix from traditional broadcast networks

## 3.3 Genre Distribution

### Top 5 Genres:

1. International Movies: 1,232 titles (30% of catalog)
2. Dramas: 1,131 titles (27.6%)
3. Comedies: 773 titles (18.8%)
4. International TV Shows: 641 titles (15.6%)
5. Documentaries: 394 titles (9.6%)

**Top 20 Genres on Netflix**



### Strategic Insights:

1. International content dominance (1,873 titles combined) demonstrates Netflix's global content strategy
2. Drama-Comedy combination provides broad appeal across demographics
3. 42 unique genres showcase content diversity, reducing subscriber churn through varied offerings
4. Notable presence of niche categories (Stand-Up Comedy: 164, Horror: 174) indicates micro-targeting strategy.

## 4. Geographic Analysis

### 4.1 Content Production by Country

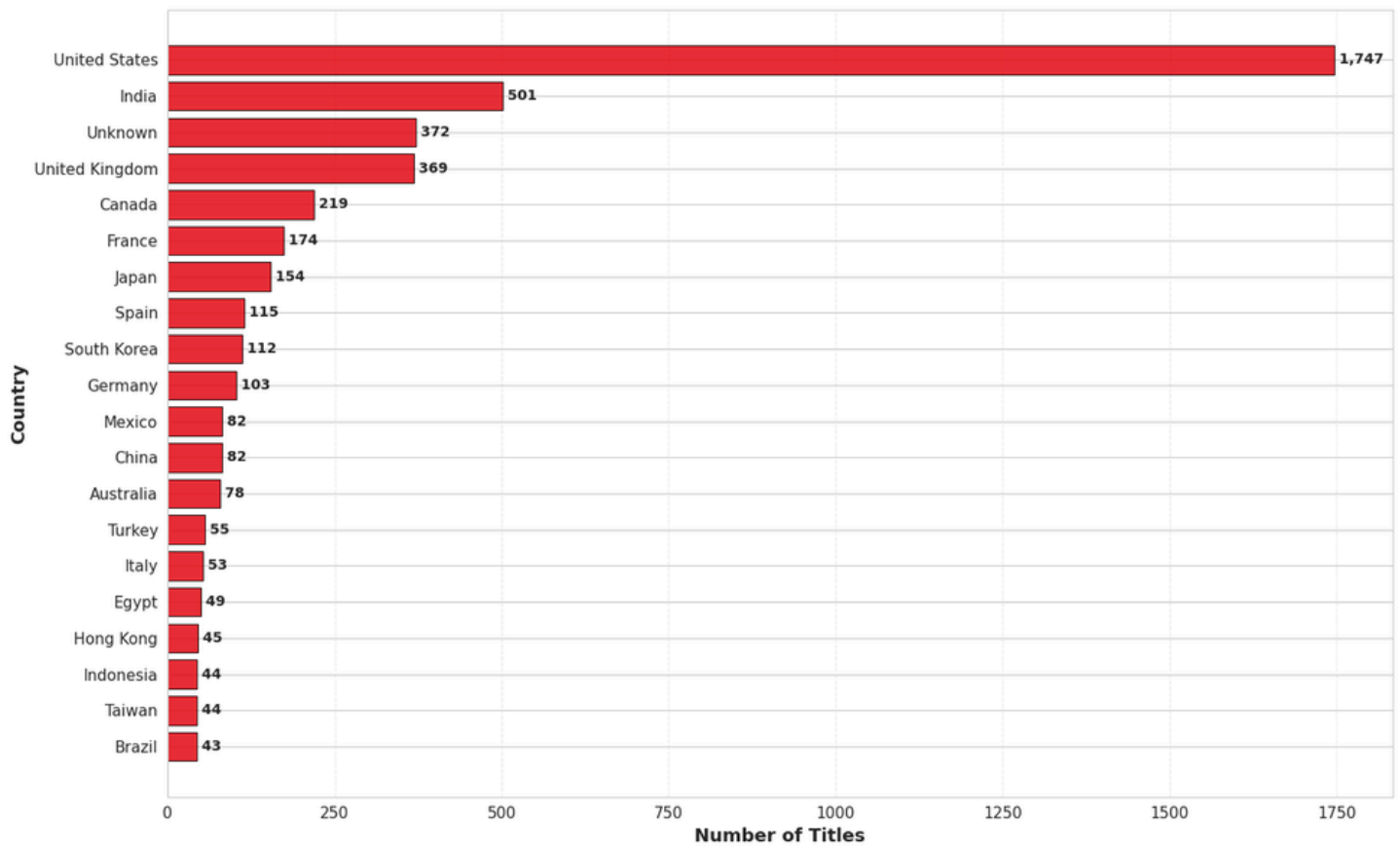
#### Top 5 Content Producers:

1. United States: 1,747 titles (42.6%)
2. India: 501 titles (12.2%)
3. Unknown: 372 titles (9.1%)
4. United Kingdom: 369 titles (9.0%)
5. Canada: 219 titles (5.3%)

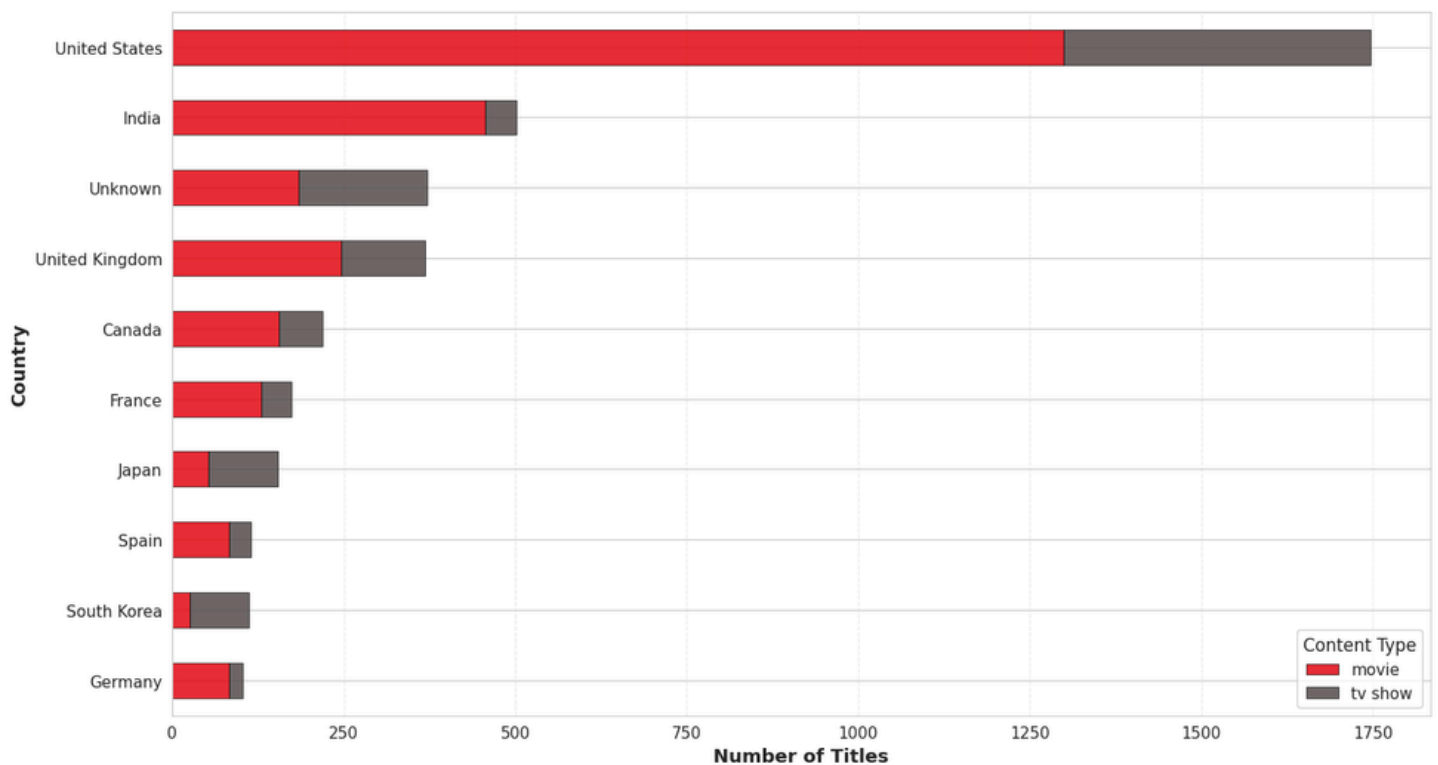
#### Global Footprint:

1. 91 countries represented demonstrates truly global content sourcing
2. Top 10 countries account for ~75% of content
3. 372 "Unknown" entries suggest data quality improvement opportunities

**Top 20 Countries Contributing Content to Netflix**



**Content Type Distribution by Top 10 Countries**



## 4.2 Content Type by Country Analysis

### Key Patterns:

#### Movie-Dominant Markets:

1. India: 91% movies (456/501) - reflects Bollywood's film-centric industry
2. France: 74% movies (129/174) - traditional cinema culture
3. United States: 74% movies (1,300/1,747)

#### TV Show-Dominant Markets:

1. South Korea: 78% TV shows (87/112) - K-drama phenomenon
2. Japan: 66% TV shows (101/154) - anime and J-drama culture
3. Unknown: 51% TV shows (188/372)

### Strategic Implications: This distribution reveals Netflix's adaptive content strategy:

1. Acquiring film libraries in cinema-strong markets (India, France, US)
2. Investing in serialized content from markets with strong TV production (South Korea, Japan)
3. Balancing global content to serve diverse viewing preferences

## 5. Critical Insights & Business Implications

### 5.1 Content Strategy Evolution

#### Phase 1 (2008-2014): Foundation Building

1. Limited catalog (<25 titles/year)
2. Focus on establishing technical infrastructure

#### Phase 2 (2015-2018): Aggressive Expansion

1. 374% average annual growth
2. Global market penetration
3. Original content initiatives (Orange Is the New Black, Stranger Things era)

#### Phase 3 (2019-2021): Strategic Curation

1. Stabilization around 700-900 annual additions
2. Quality over quantity approach
3. Increased investment in local/regional content

### 5.2 Audience Targeting

#### Primary Demographic: Adults 18-49

1. 61% of content rated TV-MA/TV-14
2. International and drama content dominance
3. Limited family/children content (12%)



**Recommendation:** Consider increasing family-friendly content (currently at 12%) to capture multi-generational households and reduce churn in family subscriber segments.

## 5.3 Competitive Positioning

### Strengths:

1. **Global content diversity** (91 countries, 42 genres)
2. **Strong international presence** (45% non-US content)
3. **Balanced movie-TV show portfolio** (69%-31% split)

### Opportunities:

1. **Data quality improvement** (372 "Unknown" country entries)
2. **Family content gap** (only 12% of catalog)
3. **Regional content expansion** in underrepresented markets

## 5.4 Technical Implementation Excellence

The PySpark implementation demonstrates several best practices:

1. **Efficient data processing** using distributed computing
2. **Comprehensive visualizations** for stakeholder communication
3. **Modular analysis structure** enabling iterative insights
4. **Data quality checks** throughout the pipeline

# 6. Recommendations

## 6.1 Content Acquisition

1. **Increase family-friendly content** from 12% to 20% to capture underserved family demographics
2. **Expand regional content** in emerging markets (Southeast Asia, Africa, Latin America)
3. **Address data quality issues** by completing "Unknown" metadata fields

## 6.2 Strategic Focus Areas

1. **Double down on successful formats:** K-dramas (South Korea), Anime (Japan), Bollywood (India)
2. **Invest in mid-tier content ratings** (PG, PG-13) to broaden appeal
3. **Leverage documentary growth** as cost-effective, award-winning content

## 6.3 Analytics Enhancements

1. **Add temporal analysis** of genre trends over time
2. **Implement viewer engagement metrics** correlation with content attributes
3. **Develop predictive models** for content performance based on metadata

# 7. Conclusion

This exploratory data analysis reveals Netflix's transformation from a US-centric movie service to a global entertainment platform with sophisticated content strategy. The data demonstrates:

- **Mature platform evolution** with strategic pivots from growth to curation
- **Global content leadership** with 91 countries and 42 genres represented
- **Adult-focused positioning** with 61% mature content
- **Data-driven decision making** evidenced by genre and geographic diversification

The analysis provides actionable insights for content acquisition, audience targeting, and strategic planning while highlighting opportunities in family content and emerging markets. The technical implementation showcases enterprise-grade data engineering practices suitable for production analytics pipelines.

**Dataset Coverage:**

4,103 titles | 96 years of content | 91 countries | 42 genres | 14 rating categories