



## Open IIT Data Analytics Hackathon: Problem Statement - 3

# Problem Statement: Netflix Content Analytics & Strategic Insights Dashboard

## Background

Netflix has revolutionized the entertainment industry as the world's leading streaming service, offering thousands of movies and TV shows to over 230 million subscribers globally. Understanding content trends, viewer preferences, and strategic content acquisition patterns is crucial for maintaining competitive advantage and driving subscriber engagement.

This challenge focuses on **data analytics and strategic visualization** rather than machine learning model development. Participants will explore, analyze, and visualize Netflix's content catalog to uncover patterns, trends, and insights that can inform content strategy, regional expansion, and user engagement initiatives. The goal is to transform raw catalog data into compelling visual narratives and actionable business intelligence.

## Problem Description

You will work with a comprehensive dataset containing Netflix's movies and TV shows catalog with detailed metadata including titles, directors, cast, countries, release years, ratings, genres, and descriptions. Your challenge is to conduct in-depth exploratory data analysis (EDA), create sophisticated visualizations, and derive strategic insights that can guide Netflix's content decisions.

**This is NOT a machine learning competition** – the focus is on business analytics, trend analysis, and visual storytelling that drives strategic decision-making.

### Dataset Components:

The dataset contains **8,800+ titles** with the following key attributes:

1. **Content Identifiers**
  - Show ID (unique identifier)
  - Title
  - Type (Movie or TV Show)
  - Description
2. **Production Information**



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Director(s)
  - Cast members
  - Country of production (often multiple)
  - Release year
  - Date added to Netflix
- 3. Content Classification**
- Rating (TV-MA, PG-13, TV-14, etc.)
  - Duration (minutes for movies, seasons for TV shows)
  - Listed in (genres/categories)
- 4. Temporal Data**
- Release year of original content
  - Date added to Netflix platform
  - Historical trends (2008-2021)

## Objectives

### Primary Analysis Tasks:

- 1. Content Landscape Analysis**
- Overall composition: Movies vs. TV Shows distribution
  - Genre distribution and popularity patterns
  - Rating categories and their prevalence
  - Content volume trends over time
  - Geographic distribution of content production
  - Language and cultural diversity analysis
- 2. Temporal Trend Analysis**
- Netflix catalog growth over the years
  - Content addition patterns (monthly, yearly, seasonal)
  - Shift in content strategy over time (older vs. newer content)
  - Release year vs. addition year analysis
  - Peak addition periods and strategic timing
  - Content aging analysis (lag between release and Netflix addition)
- 3. Geographic Content Analysis**
- Top content-producing countries
  - Regional content preferences and patterns
  - International vs. domestic content ratio
  - Country-specific genre preferences
  - Global expansion strategy through content
  - Multi-country productions analysis
- 4. Genre & Category Intelligence**
- Most common genres and categories
  - Genre evolution over time
  - Cross-genre analysis (content with multiple categories)



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Niche vs. mainstream genre balance
- Genre-rating correlations
- Emerging genre trends

### 5. Content Creator Analysis

- Most prolific directors on Netflix
- Cast member frequency analysis
- Collaboration networks (directors, actors)
- Celebrity content patterns
- International vs. domestic talent distribution
- Director-genre specializations

### 6. Rating & Audience Targeting

- Distribution across rating categories
- Rating evolution over time
- Correlation between rating and genre
- Family-friendly vs. mature content balance
- Regional rating patterns
- Strategic positioning for different demographics

### 7. Duration & Format Analysis

- Movie length distributions and trends
- TV show season patterns
- Binge-worthy content analysis (season counts)
- Duration preferences by genre
- Evolution of content length over time

### Advanced Analysis (Bonus):

- **Text Analytics:** Analyze descriptions for themes, sentiment, keywords
- **Network Analysis:** Visualize collaboration networks between cast/directors
- **Competitive Positioning:** Compare content strategy shifts over years
- **Gap Analysis:** Identify underrepresented genres or regions
- **Content Lifecycle:** Analyze how long content stays on platform
- **Strategic Pivot Points:** Identify when Netflix strategy changed significantly
- **Recommendation Insights:** What combinations of factors define popular content?

### Expected Deliverables

#### 1. Comprehensive Strategic Analysis Report (25-35 pages)

##### Executive Summary (3-4 pages)

- Key strategic findings and recommendations
- Top 10 insights for Netflix leadership
- Critical trends and opportunities



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Data quality assessment

### Content Landscape Overview (5-7 pages)

- Catalog composition and structure
- High-level statistics and distributions
- Key trends visualization
- Comparative analysis (movies vs. TV shows)

### Temporal Analysis Section (5-6 pages)

- Historical growth patterns
- Content addition strategy evolution
- Seasonal and cyclical patterns
- Time-lag analysis (release to addition)
- Strategic inflection points

### Geographic Intelligence (5-6 pages)

- Global content footprint
- Regional production hubs
- International expansion strategy
- Cultural diversity metrics
- Market-specific insights

### Genre & Content Strategy (4-5 pages)

- Genre distribution and evolution
- Category mixing patterns
- Niche vs. mainstream balance
- Rating strategy by genre
- Emerging opportunities

### Creator & Talent Analysis (3-4 pages)

- Key contributors to catalog
- Collaboration patterns
- Talent diversity
- Celebrity content impact

### Strategic Recommendations (4-5 pages)

- Content acquisition priorities



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Geographic expansion opportunities
- Genre diversification strategies
- Audience targeting insights
- Competitive positioning recommendations
- Investment allocation suggestions

### 2. Visualization Portfolio

Create **30-40 distinct visualizations** organized into themed collections:

#### Overview & Summary Visualizations

- Netflix catalog dashboard (single-page overview)
- Content composition pie/donut charts
- Key metrics cards (total titles, countries, genres)
- Timeline of Netflix growth
- Interactive summary infographic

#### Temporal Visualizations

- Line charts showing catalog growth over time
- Stacked area charts for movie vs. TV show trends
- Calendar heatmaps for content additions
- Release year distributions
- Time-lag analysis (release to addition)
- Seasonal pattern analysis
- Animated timeline showing evolution

#### Geographic Visualizations

- World map showing content by country
- Choropleth maps for production volume
- Top countries bar charts
- Regional comparison dashboards
- Multi-country production analysis
- Geographic diversity index

#### Genre & Category Visualizations

- Genre distribution treemaps
- Word clouds for categories
- Sankey diagrams for genre relationships
- Genre evolution over time (animated)



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Cross-genre analysis matrix
- Bubble charts for genre popularity

### Content Creator Visualizations

- Top directors/actors bar charts
- Network graphs for collaborations
- Director-genre heatmaps
- Cast frequency distributions
- Celebrity content timelines

### Rating & Audience Visualizations

- Rating distribution charts
- Age group targeting analysis
- Rating by genre heatmaps
- Family-friendly content metrics
- Demographic targeting dashboard

### Duration & Format Visualizations

- Movie duration histograms
- TV show season distributions
- Box plots for duration by genre
- Duration trends over time
- Format comparison charts

### Advanced Visualizations

- Text analysis word clouds from descriptions
- Sentiment distribution across genres
- Network graphs (cast/director collaborations)
- Parallel coordinates plots
- Interactive drill-down dashboards

## 3. Interactive Strategic Dashboard

Develop a fully functional, executive-ready dashboard with multiple tabs:

### Tab 1: Executive Overview

- KPI cards (total titles, growth rate, diversity metrics)
- High-level summary charts
- Key findings highlights



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Strategic recommendations preview

### Tab 2: Content Explorer

- Search and filter functionality
- Detailed content browser
- Sort by multiple criteria
- Quick facts display
- Export functionality

### Tab 3: Trend Intelligence

- Interactive time-series charts
- Trend comparison tools
- Seasonal pattern analysis
- Growth projections
- Historical milestone markers

### Tab 4: Geographic Insights

- Interactive world map
- Country comparison tool
- Regional deep dives
- Production hub analysis
- Market opportunity scanner

### Tab 5: Genre & Category Intelligence

- Genre explorer with filters
- Category co-occurrence matrix
- Trend analysis by genre
- Opportunity identification
- Competitive gap analysis

### Tab 6: Creator & Talent Hub

- Director/actor search
- Collaboration network visualization
- Portfolio analysis
- Talent diversity metrics
- Rising stars identification



## Open IIT Data Analytics Hackathon: Problem Statement - 3

### Tab 7: Strategic Recommendations

- Actionable insights summary
- Priority recommendations
- Investment opportunities
- Risk areas identification
- Next steps roadmap

### Technical Features:

- Responsive design (desktop, tablet)
- Real-time filtering across all visualizations
- Export capabilities (PDF, PNG, CSV)
- Bookmark favorite views
- Shareable dashboard links
- Dark/light mode toggle
- Tooltips with detailed information
- Drill-down capabilities

## 4. Executive Presentation (15-20 slides)

### Presentation Structure:

1. **Title & Context** (1 slide): Project overview
2. **Key Findings** (2-3 slides): Top strategic insights
3. **Content Landscape** (2 slides): Catalog overview
4. **Temporal Trends** (2-3 slides): Growth and addition patterns
5. **Geographic Strategy** (2 slides): Global footprint analysis
6. **Genre Intelligence** (2 slides): Category insights
7. **Content Creators** (1-2 slides): Talent analysis
8. **Strategic Recommendations** (3-4 slides): Actionable next steps
9. **Data Methodology** (1 slide): Analysis approach
10. **Appendix** (2-3 slides): Additional details

### Presentation Requirements:

- Executive-friendly language (minimal jargon)
- Visual-first approach (charts on every slide)
- Clear action items
- Data-driven recommendations
- Professional design aesthetic



## Open IIT Data Analytics Hackathon: Problem Statement - 3

### 5. Analysis Code & Documentation

Jupyter Notebook / R Markdown containing:

- Well-documented analysis workflow
- Data cleaning and preprocessing steps
- Exploratory data analysis code
- Statistical analysis functions
- Visualization generation code
- Custom functions for reusability
- Inline explanations and insights
- Reproducible workflow

**Code Quality Requirements:**

- Clean, PEP 8 compliant (Python) or tidyverse style ®
- Comprehensive comments
- Function documentation
- Clear variable naming
- Modular code structure
- Error handling
- Performance optimization

## Evaluation Criteria

### Business Impact & Strategic Thinking (30%)

- **Insight Quality:** Actionable, relevant to Netflix strategy
- **Strategic Depth:** Understanding of business implications
- **Recommendation Clarity:** Clear, specific, implementable
- **Market Understanding:** Awareness of streaming industry context
- **Innovation:** Novel perspectives and creative analysis

### Data Analysis Excellence (25%)

- **Thoroughness:** Comprehensive exploration of dataset
- **Statistical Rigor:** Appropriate analytical methods
- **Data Quality:** Proper handling of missing data, outliers
- **Insight Generation:** Ability to extract meaningful patterns
- **Validation:** Cross-checking findings multiple ways

### Visualization & Design (25%)

- **Clarity:** Easy to understand and interpret



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Professional Quality:** Publication-ready aesthetics
- Appropriateness:** Right chart types for data
- Consistency:** Unified visual language
- Interactivity:** Engaging, functional dashboard elements
- Storytelling:** Visual narrative flow

### Communication & Presentation (15%)

- Executive Readiness:** Suitable for C-suite presentation
- Narrative Flow:** Logical progression of insights
- Writing Quality:** Professional, clear, concise
- Audience Awareness:** Business-focused communication
- Actionability:** Clear implications and next steps

### Technical Execution (5%)

- Dashboard Functionality:** Bug-free, smooth operation
- Code Quality:** Clean, documented, reproducible
- Tool Mastery:** Effective use of analytics tools
- Documentation:** Clear setup and usage instructions

## Dataset Information

### Primary Dataset

#### Netflix Movies and TV Shows

- Source:** Kaggle
- Link:** <https://www.kaggle.com/datasets/shivamb/netflix-shows/data>
- Description:** Comprehensive Netflix catalog with 8,800+ titles
- Columns:** 12 attributes including title, type, director, cast, country, date\_added, release\_year, rating, duration, listed\_in, description
- Time Period:** Content from 2008-2021
- Update Frequency:** Regularly updated

### Dataset Schema:

Column	Description	Data Type
show_id	Unique identifier	String



## Open IIT Data Analytics Hackathon: Problem Statement - 3

type	Movie or TV Show	Categorical
title	Name of content	String
director	Director(s) name	String (comma-separated)
cast	Cast members	String (comma-separated)
country	Production country	String (comma-separated)
date_added	Date added to Netflix	Date
release_year	Original release year	Integer
rating	Content rating	Categorical
duration	Length (minutes or seasons)	String
listed_in	Genres/categories	String (comma-separated)
description	Plot summary	Text

### Complementary Datasets (For Enhanced Analysis)

#### 1. Netflix Movies and TV Shows (Updated 2025)

- Link:  
<https://www.kaggle.com/datasets/bhargavchirumamilla/netflix-movies-and-tv-shows-till-2025>
- Description: Extended dataset with more recent titles



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Use: Trend validation and future projection
- 2. **Netflix Prize Dataset** (Historical)
  - Link: <https://www.kaggle.com/netflix-inc/netflix-prize-data>
  - Description: Historical viewer ratings
  - Use: Understanding viewer preferences
- 3. **IMDb Dataset**
  - Link: <https://www.kaggle.com/datasets/ashirwadsangwan/imdb-dataset>
  - Description: Movie/TV show ratings and metadata
  - Use: Cross-reference quality metrics
- 4. **The Movies Dataset**
  - Link: <https://www.kaggle.com/datasets/rounakbanik/the-movies-dataset>
  - Description: Movie metadata, ratings, revenues
  - Use: Contextual analysis

## Analytical Approach & Methodology

### 1. Data Preparation

```
# Suggested workflow
- Load dataset and examine structure
- Handle missing values appropriately
- Parse comma-separated fields (cast, director, country, genres)
- Convert date_added to datetime
- Extract additional features (month, quarter, year)
- Clean text fields
- Create derived metrics
```

### 2. Exploratory Data Analysis Techniques

#### Univariate Analysis:

- Distribution of content types
- Rating frequency
- Genre popularity
- Country distribution
- Year-wise content volume

#### Bivariate Analysis:

- Type vs. Rating
- Genre vs. Rating
- Country vs. Content Type
- Release Year vs. Addition Year
- Duration patterns by genre



## Open IIT Data Analytics Hackathon: Problem Statement - 3

### Multivariate Analysis:

- Genre-Country-Rating relationships
- Temporal-Geographic patterns
- Content type-Genre-Duration interactions
- Complex trend analysis

### Text Analysis:

- Description keyword extraction
- Genre co-occurrence patterns
- Sentiment analysis of descriptions
- Theme identification

## 3. Statistical Analysis Methods

### Descriptive Statistics:

- Central tendency (mean, median, mode)
- Dispersion (variance, standard deviation)
- Distribution shape (skewness, kurtosis)
- Percentiles and quartiles

### Comparative Analysis:

- Chi-square tests for categorical variables
- T-tests for comparing groups
- ANOVA for multiple group comparisons
- Correlation analysis

### Time Series Analysis:

- Trend detection
- Seasonal decomposition
- Growth rate calculations
- Moving averages
- Year-over-year comparisons

### Text Analytics:

- Word frequency analysis



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- TF-IDF for important terms
- N-gram analysis
- Topic modeling (optional)

### 4. Visualization Best Practices

#### Chart Selection Guidelines:

- **Bar charts:** Categorical comparisons
- **Line charts:** Trends over time
- **Pie/Donut:** Composition (use sparingly)
- **Heatmaps:** Two-dimensional patterns
- **Treemaps:** Hierarchical data
- **Maps:** Geographic distribution
- **Sankey:** Flow and relationships
- **Network graphs:** Connections and collaborations
- **Word clouds:** Text frequency (supplementary only)
- **Box plots:** Distribution comparisons
- **Scatter plots:** Relationships between variables

#### Design Principles:

- Use Netflix brand colors (red #E50914, black, white)
- Consistent color scheme throughout
- Clear, large labels
- Minimal chart junk
- Appropriate aspect ratios
- Accessible color palettes
- Professional typography

#### Interactive Elements:

- Hover tooltips with details
- Click-through drill-downs
- Filter panels
- Date range selectors
- Geographic zoom
- Dynamic legends

## Success Metrics

#### Quantitative Indicators:



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- **Visualization Count:** 30-40 distinct charts/graphs
- **Insight Depth:** 20-25 strategic findings
- **Dashboard Pages:** 5-7 interactive tabs
- **Report Length:** 25-35 pages with visuals
- **Code Documentation:** 100% of major functions commented
- **Analysis Breadth:** Cover all 7 primary analysis areas

### Qualitative Indicators:

- **Executive Readiness:** Suitable for C-suite presentation
- **Strategic Value:** Insights lead to actionable recommendations
- **Visual Impact:** Charts are memorable and compelling
- **Business Alignment:** Analysis addresses real business questions
- **Innovation:** Unique perspectives beyond obvious observations

## Strategic Questions to Address

Your analysis should provide answers to questions like:

### Content Strategy:

- Is Netflix focusing more on movies or TV shows over time?
- Which genres are underrepresented in the catalog?
- How quickly does Netflix add new releases vs. catalog content?
- What's the optimal content mix by rating category?

### Geographic Strategy:

- Which countries are emerging as content production hubs?
- Is Netflix diversifying away from US content?
- What regions show the most growth potential?
- How does international content perform?

### Temporal Strategy:

- When are the peak content addition periods?
- How has the content addition strategy evolved?
- What's the average lag between release and Netflix addition?
- Are there seasonal patterns in content strategy?

### Audience Strategy:

- Is Netflix balancing family vs. mature content?



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Which rating categories dominate?
- How are different demographics being targeted?

### Competitive Strategy:

- How has Netflix's strategy shifted over the years?
- What strategic inflection points can be identified?
- Where are the competitive gaps?
- What's the next frontier for content expansion?

## Technical Requirements

### Required Python Libraries:

```
# Core Data Analysis
pandas, numpy

# Visualization
matplotlib, seaborn, plotly, altair

# Dashboard
streamlit / dash / panel

# Text Processing
nltk, wordcloud, textblob (optional)

# Geographic
folium, geopandas (for maps)

# Network Analysis (optional)
networkx

# Utilities
datetime, collections, itertools
```

### Required R Libraries (if using R):

```
# Data Manipulation
tidyverse, dplyr, tidyr, stringr

# Visualization
ggplot2, plotly, ggmap

# Dashboard
shiny, flexdashboard, shinydashboard
```



## Open IIT Data Analytics Hackathon: Problem Statement - 3

```
# Text Processing  
tidytext, wordcloud
```

```
# Utilities  
lubridate, scales
```

## Submission Guidelines

### Deliverable Structure:

```
submission/  
|   └── report/  
|       ├── strategic_analysis_report.pdf  
|       ├── executive_summary.pdf  
|       └── appendices/  
|           ├── detailed_tables.pdf  
|           └── methodology_notes.pdf  
  
|   └── visualizations/  
|       ├── overview_charts/  
|       ├── temporal_analysis/  
|       ├── geographic_insights/  
|       ├── genre_intelligence/  
|       └── creator_analysis/  
  
|   └── dashboard/  
|       ├── app.py (or app.R)  
|       ├── requirements.txt  
|       ├── data/ (processed data files)  
|       ├── assets/ (images, css)  
|       └── README.md  
  
|   └── code/  
|       ├── netflix_analysis.ipynb  
|       ├── data_preprocessing.py  
|       ├── visualization_functions.py  
|       ├── text_analysis.py  
|       └── statistical_analysis.py  
  
|   └── presentation/  
|       └── netflix_insights_presentation.pptx
```



## Open IIT Data Analytics Hackathon: Problem Statement - 3

└─ README.md

## Evaluation Timeline

- Data exploration and EDA
- Deep analysis and pattern identification
- Visualization development
- Dashboard creation
- Report writing and refinement
- Final polish and presentation prep
- Submission and presentations

## Judging Panel

Solutions will be evaluated by:

- Arijit Hajra, aka Robot Man of India, CEO, Think Again Lab

## Case Study Examples

### High-Quality Analysis Would Include:

#### Example 1: Geographic Expansion Insight

“Analysis reveals that Netflix added 300% more Korean content between 2018-2021, coinciding with global K-drama popularity. Recommendation: Accelerate investment in Southeast Asian content production.”

#### Example 2: Genre Gap Identification

“Documentary content represents only 8% of catalog but 15% of trending titles. Visualization shows this gap widening post-2019. Recommendation: Increase documentary acquisitions by 50%.”

#### Example 3: Temporal Strategy

“Calendar heatmap reveals Netflix adds 40% more content in Q4, likely preparing for holiday viewing. However, engagement data suggests Q1 additions have higher retention. Recommendation: Shift strategy to Q1 loading.”

## Resources and Learning Materials

### Business Context:

- Netflix investor relations materials
- Streaming industry reports



## Open IIT Data Analytics Hackathon: Problem Statement - 3

- Media strategy case studies
- Content licensing economics
- Global entertainment trends

### Data Visualization:

- "Storytelling with Data" by Cole Nussbaumer Knaflic
- "The Visual Display of Quantitative Information" by Edward Tufte
- Netflix Tech Blog visualization articles
- Tableau Public gallery
- Observable HQ examples

### Data Analysis:

- Python for Data Analysis (Wes McKinney)
- R for Data Science (Hadley Wickham)
- Business analytics case studies
- Industry benchmark reports

### Tools Documentation:

- Plotly: <https://plotly.com/python/>
- Streamlit: <https://docs.streamlit.io/>
- Tableau: <https://www.tableau.com/learn>
- Power BI: <https://docs.microsoft.com/power-bi/>

## Important Considerations

### Business Context Awareness:

- **Streaming Wars:** Netflix competes with Disney+, Prime Video, HBO Max, etc.
- **Content Licensing:** Not all content shown is Netflix-owned
- **Regional Variations:** Catalog differs by country due to licensing
- **Cost Considerations:** Original content vs. licensed content economics
- **Subscriber Metrics:** Content drives retention and acquisition

### Analytical Ethics:

- **Data Limitations:** Dataset may not be complete or current
- **Bias Awareness:** Dataset represents what Netflix chose to add, not all content
- **Inference Boundaries:** Don't claim causation without evidence
- **Transparency:** Clearly state assumptions and limitations
- **Responsible Recommendations:** Consider real-world constraints



## Open IIT Data Analytics Hackathon: Problem Statement - 3

### Common Pitfalls to Avoid:

1. **Surface-Level Analysis:** Going beyond “X is the most common genre”
2. **Chart Overload:** Too many similar visualizations
3. **Missing the “So What”:** Insights without business implications
4. **Ignoring Temporal Dimension:** Not analyzing how things change over time
5. **Geographic Oversimplification:** Not considering regional nuances
6. **Text Field Neglect:** Ignoring valuable description/cast/director data
7. **Static Dashboards:** Interactive dashboards that aren’t actually interactive
8. **Jargon Heavy:** Reports too technical for business audience

### Expected Impact

High-quality analysis will:

- **Guide Content Strategy:** Inform what content to license or produce
- **Optimize Investment:** Identify high-ROI content categories
- **Support Global Expansion:** Reveal geographic opportunities
- **Enhance User Engagement:** Understand viewer preferences
- **Competitive Intelligence:** Identify strategic positioning
- **Data-Driven Culture:** Demonstrate value of analytics in entertainment

### Final Notes

This challenge offers an opportunity to apply **business analytics and data visualization** to real entertainment industry data. Your work should be:

- ✓ **Business-Focused:** Every insight should have strategic implications
- ✓ **Visually Compelling:** Dashboard could be shown to executives
- ✓ **Analytically Rigorous:** Findings should be statistically sound
- ✓ **Clearly Communicated:** Non-technical stakeholders should understand
- ✓ **Actionable:** Recommendations should be implementable

Remember: You are **NOT building recommendation engines or ML models**. Focus on understanding Netflix’s content strategy deeply, visualizing it beautifully, and providing strategic insights that could guide a billion-dollar content business.

**Think like a strategy consultant, analyze like a data scientist, present like a storyteller.**

**Good Luck!**