

## 1. Data Cleaning Process

Section	Description
Process	<b>Data Cleaning</b> is the process of detecting and correcting corrupt or inaccurate records from a record set, table, or database. It involves identifying incomplete, incorrect, inaccurate, or irrelevant parts of the data and replacing, modifying, or deleting the dirty data.
Steps	<ol style="list-style-type: none"><li><b>Data Inspection and Loading:</b> Read the raw dataset into a DataFrame.</li><li><b>Data Type Conversion:</b> Convert columns to the appropriate data type. For instance, the 'date_added' column was converted from an object type to the proper datetime64[ns] type to enable time-series analysis.</li><li><b>Missing Value Handling :</b> Identify missing values using <code>df.isna().sum()</code>. All missing values in the dataset were replaced with the value 0 using <code>df.fillna(0)</code> to ensure a complete dataset for subsequent steps.</li></ol>

Section	Description
	4. <b>Handling of Duplicates:</b> Check for and remove duplicate rows to ensure data uniqueness.
<b>Metrics / Outcome</b>	<p>1. <b>Missing Value Count:</b> The primary metric is the count of missing values per column, aiming for a count of <b>0</b> in all essential columns.</p> <p>2. <b>Data Type Integrity:</b> Verification that columns like 'date_added' are successfully converted to the correct data type</p> <p>3. <b>Data Completeness:</b> Final dataset contains the full expected number of rows and columns</p>

## 2. Normalization Process

Section	Description
<b>Process</b>	<b>Normalization</b> in data preprocessing is the process of adjusting values measured on different scales to a notionally common scale, often to prevent features with larger ranges from dominating the model training

Section	Description
	process. (Note: The uploaded notebook also includes feature encoding steps).
Steps	<ol style="list-style-type: none"> <li><b>1. Data Loading:</b> Read the cleaned dataset for further transformation.</li> <li><b>2. Data Imputation:</b> Fill in missing values for text-based columns, such as replacing missing values in the 'director' column with a placeholder like "Unknown Director".</li> <li><b>3. Categorical Encoding:</b> Create a new feature by mapping a categorical variable to a numerical code. A feature called 'director_code' was created, where a <b>Movie</b> type is mapped to <b>1</b> and a <b>TV Show</b> type is mapped to <b>0</b>.</li> <li><b>4. Numerical Scaling/Standardization:</b> Apply techniques like Min-Max Scaling or Z-Score Standardization to numerical columns if they were not already handled in the Feature Engineering stage.</li> </ol>

Section	Description
<b>Metrics / Outcome</b>	<p><b>1. Missing Value Resolution:</b> Confirmation that all placeholder values (e.g., "Unknown Director") have successfully replaced the original missing values.</p> <p><b>2. Encoding Consistency:</b> Verification that the encoding rule is correctly applied across all rows (e.g., every 'Movie' row has director_code = 1).</p> <p><b>3. Distribution Uniformity (Inferred):</b> For true normalization (scaling), the resulting numerical distribution should have a more uniform range (e.g., 0 to 1 for Min-Max Scaling) or mean 0 and standard deviation 1 (for Standardization).</p>

### 3.Exploratory Data Analysis (EDA) Process

Section	Description
<b>Process</b>	<b>Exploratory Data Analysis (EDA)</b> is an approach to analyzing data sets to summarize their main characteristics, often with visual methods. It is used to discover

Section	Description
	patterns, spot anomalies, test hypotheses, and check assumptions with the help of summary statistics and graphical representations.
Steps	<p>1. <b>Load Cleaned Data:</b> Read and display the cleaned dataset (e.g., netflix_cleaned.csv).</p> <p>2. <b>Structural Check:</b> Examine the dataset's basic structure, including the number of rows/columns and data types (df.info()) to confirm the data is ready for analysis.</p> <p>3. <b>Summary Statistics:</b> Generate and print "Quick Insights" which include mode and frequency counts for key categorical columns.</p> <p>4. <b>Visualization :</b> Create plots (e.g., bar charts, histograms, heatmaps) to visualize distributions and relationships between features.</p>
Metrics / Outcome	<p>1. <b>Data Count:</b> Total record count</p> <p>2. <b>Feature Distribution (Counts):</b> Quantified distribution of data types (e.g., <b>Movies: 6130</b> and <b>TV Shows: 2676</b>).</p>

Section	Description
	<p>3. <b>Mode/Frequency:</b> Identification of the most frequent categories, such as the <b>Most common rating (TV-MA)</b>, <b>Most frequent country (United States)</b>, and <b>Most frequent genre (dramas, international movies)</b>.</p> <p>4. <b>Data Quality Score (Inferred):</b> Insights into data consistency and potential outliers are implicitly measured.</p>

#### 4.Feature Engineering

Section	Description
<b>Process</b>	<b>Feature Engineering</b> is the process of using domain knowledge to select, transform, or create variables (features) that best represent the underlying problem to a predictive model, resulting in improved model accuracy.
<b>Steps</b>	<p>1. <b>Categorical Encoding (Label Encoding):</b> Use techniques like LabelEncoder (from sklearn.preprocessing) to convert</p>

Section	Description
	<p>categorical columns, such as 'type', into numerical representations creating a new feature 'type_encoded'.</p> <p><b>2. Feature Creation (Binning/Grouping):</b> Create a new, derived feature by grouping values from an existing feature. A function was used to categorize the 'duration' column into four distinct bins: 'Short', 'Medium', 'Long', and 'Series', creating the feature 'Content_Length_Category'.</p> <p><b>3. Numerical Scaling (MinMaxScaler):</b> Scale a numerical column to a fixed range. The MinMaxScaler was applied to 'release_year' to scale its values to the range of 0.0 to 1.0, resulting in the new feature 'release_year_scaled'.</p>
<b>Metrics / Outcome</b>	<p><b>1. Range Validation (Scaling):</b> Scaled numerical features (e.g., 'release_year_scaled') must fall within the defined range (e.g., 0.0 to 1.0).</p>

Section	Description
	<p><b>2. Successful Categorization:</b> Verification that all values in the original column ('duration') are correctly mapped to one of the defined categories in the new feature ('Content_Length_Category').</p> <p><b>3. Information Gain :</b> The success of Feature Engineering is ultimately measured by the improvement in the predictive model's performance metrics (e.g., accuracy, precision) when trained on the new features.</p>