

# Understanding of the Data

## TITLE.BASICS.TSV:

This report provides a formal assessment of the title.basics.tsv dataset structure and identifies necessary preparatory steps for robust Exploratory Data Analysis (EDA). The analysis is crucial for ensuring data integrity and optimizing performance before statistical and temporal investigations.

Statistic	Value	Implication for EDA
Total Records	12,097,750	Due to the high row count, initial processing (such as grouping and aggregation) should prioritize efficient data types to prevent memory overflow.
Memory Footprint	\$\approx\$ 4.07\$ GB	Suggests that optimizing object (string) columns and avoiding memory-intensive operations (like large cross-joins) without prior filtering is necessary.

## Columns Data And its Understanding:

Column Name	Current Type	Data Quality Issue	Required Transformation
runtimeMinutes	object	<b>Critical Non-Numeric Data.</b> The column contains non-convertible characters (typically \N), preventing direct numerical analysis.	<b>Coercion:</b> Must be converted to a numeric type (float64 or Int64) with errors coerced to NaN.
startYear	float64	Presence of \$\text{NaN}\$ values	<b>Casting:</b> Convert to Int64 (nullable)

Column Name	Current Type	Data Quality Issue	Required Transformation
		forces a non-integer type.	integer) after ensuring all non-\$text{NaN}\$ values are whole numbers.
genres	object	<b>Multi-Label String Format.</b> Values are delimited by commas (e.g., "Action,Drama"), rendering direct frequency counting impossible.	<b>Parsing:</b> Must be split and "exploded" to create a standard, single-label row format for accurate genre distribution analysis.
titleType	category	Highly optimized categorical data.	<b>None.</b> Ideal for immediate segmentation (e.g., comparing runtime of 'movie' vs. 'short').
isAdult	int8	Binary flag (0 or 1).	<b>Casting:</b> Convert to a boolean type (bool) for cleaner filtering and logical operations.

## Missing Data:

Handling missing values is paramount, as different columns exhibit different—and expected—levels of data incompleteness.

Column	Approximate Null Count	Missingness Percentage	Interpretation and Strategy
runtimeMinutes	\$\approx 7.8\$ Million	\$\approx 64\%\$	<b>High Information Loss.</b> Analysis involving duration must be conducted on a substantial <b>subset</b> of

Column	Approximate Null Count	Missingness Percentage	Interpretation and Strategy
			the data. Imputation (filling with mean/median) is generally inadvisable due to the vast difference between runtimes of "shorts" and "features."
endYear	\$\approx 11.9\$ Million	\$\approx 98\%\$	<b>Expected Incompleteness.</b> This primarily applies to ongoing TV series and movies. <b>Strategy:</b> This column should only be utilized when calculating the <i>lifespan</i> of completed TV shows; it should otherwise be ignored for general analysis.
genres	\$\approx 530,000\$	\$\approx 4.4\%\$	<b>Low Incompleteness.</b> Missing values here can be safely handled by either dropping the records or assigning them a categorical placeholder (e.g., "Unspecified") without compromising overall genre statistics.

### Title.Rating.TSV:

Statistic	Value	Implication for EDA
Total Records	\$1,605,930\$	Only \$\approx 13\%\$ of the \$12\$ million titles in the BASICS file have a rating. This confirms that the majority of titles (shorts, episodes, etc.) are unrated.

Statistic	Value	Implication for EDA
Memory Footprint	\$\approx 120\$ MB	Excellent memory efficiency, making this table ideal for direct loading and joining operations.
Missing Data	Zero	The dataset is exceptionally clean, eliminating the need for missing value imputation before analysis.

## Column Analysis:

Column Name	Data Type	Analytical Role & Integration Strategy
tconst	category	<b>Primary Join Key.</b> This column acts as the foreign key required to link rating metrics back to descriptive features (e.g., genre, runtime) in the TITLE.BASICS.TSV file.
averageRating	float32	<b>Quality Metric.</b> This is the simple arithmetic mean of all votes. <b>EDA Tip:</b> This metric is easily skewed by low vote counts and must be weighted (see Section 3).
numVotes	int32	<b>Reliability Metric.</b> Indicates the statistical reliability of the averageRating.

## Title.Crew.Tsv:

The CREW file maintains the same dimensionality as the BASICS file, indicating that every title record has a corresponding crew record, even if the crew fields themselves are null.

Statistic	Value	Implication for EDA
Total Records	\$12,097,750\$	This file must be joined with TITLE.BASICS.TSV using the tconst key. Since it has the same number of rows, it suggests a \$text{1:1}\$ relationship with the basics file.

Statistic	Value	Implication for EDA
Memory Footprint	\$\approx 2.03\$ GB	Moderate memory usage for its size, but efficiency is hampered by all three columns being loaded as object (string) types.

## Column Analysis:

Column Name	Current Type	Data Quality Issue	Required Transformation
tconst	object	<b>Suboptimal Type.</b> The unique identifier (join key) is loaded as a generic string.	<b>Casting:</b> Convert to category for memory efficiency and optimal join performance with other IMDB tables.
directors	object	<b>Multi-Value Field.</b> Contains comma-separated IMDB Person IDs (e.g., nm0005690,nm0721526) which are critical for network analysis.	<b>Parsing &amp; Exploding:</b> Must be split into lists and then "exploded" to create one row per title-director relationship, enabling director-centric analysis (e.g., director influence).
writers	object	<b>Multi-Value Field.</b> Contains comma-separated IMDB Person IDs.	<b>Parsing &amp; Exploding:</b> Similar to directors, this field must be processed for writer-centric network analysis.

## Missing Data:

Column	Approximate Null Count	Missingness Percentage	Interpretation and Strategy
directors	\$\approx 5.3\$ Million	\$\approx 44\%\$	<b>Significant Missingness.</b> A large portion of titles lack director information. <b>Strategy:</b> Filter for non-null values before running director network analysis or

Column	Approximate Null Count	Missingness Percentage	Interpretation and Strategy
			calculating director-based statistics.
writers	\$\approx 6.1\$ Million	\$\approx 50\%\$	<b>Significant Missingness.</b> Half the dataset lacks writer data. This is expected, as many shorts or non-narrative titles may not credit a writer.

## NAME.BASICS.TSV:

This is the largest personnel-related file, containing nearly 15 million records, representing the unique individuals tracked in the IMDB database.

Statistic	Value	Implication for EDA
Total Records	\$14,905,821\$	This table acts as the master lookup dictionary for all individuals. It must be joined with the exploded TITLE.CREW.TSV file to transform \$\text{nconst}\$ IDs into readable names.
Memory Footprint	\$\approx 3.11\$ GB	Moderate memory usage for its size. The primary goal is to convert the object type \$\text{nconst}\$ to category for improved join performance.

## Column Analysis:

Column Name	Current Type	Data Quality Issue	Required Transformation
nconst	object	<b>Key Column Suboptimal Type.</b> The unique personnel ID is loaded as a generic string.	<b>Casting:</b> Convert to category for memory efficiency and optimal join performance.

Column Name	Current Type	Data Quality Issue	Required Transformation
<b>primaryName</b>	object	<b>Identification Field.</b> Contains the individual's full name.	<b>None.</b> Ready for use in visualizations, tables, and top-list rankings.
<b>birthYear</b>	float64	<b>Temporal Data.</b> Loaded as float due to \$\\text{NaN}\$ values (e.g., birth year unknown, or living).	<b>Casting:</b> Convert to Int64 (nullable integer). Used to calculate age at death or career span.
<b>deathYear</b>	float64	<b>Temporal Data.</b> Primarily \$\\text{NaN}\$ for living individuals.	<b>Casting:</b> Convert to Int64 (nullable integer). Used to calculate career longevity for deceased individuals.
<b>primaryProfession</b>	object	<b>Multi-Value Field.</b> Contains comma-separated professions (e.g., "actor,director,producer") .	<b>Parsing &amp; Exploding:</b> Must be split and "exploded" to analyze the distribution of individuals across different professions (e.g., how many actors are also directors).

## Missing Data:

Column	Approximate Null Count	Missingness Percentage	Interpretation and Strategy
birthYear	\$\approx 660,000\$	\$\approx 4.4\%\$	<b>Low Missingness.</b> Missing records should be excluded from age-based cohort analysis but can be tolerated for name-based lookups.
deathYear	\$\approx 14.65\$ Million	\$\approx 98\%\$	<b>Expected Incompleteness.</b> The vast majority of individuals in the dataset are assumed to be living or their death date is unrecorded.
primaryProfession	\$\approx 2.97\$ Million	\$\approx 20\%\$	<b>Significant Missingness.</b> One-fifth of the records lack a profession tag. <b>Strategy:</b> Exclude these records from professional distribution analyses, or categorize them as 'Unspecified.'

### TITLE.PRINCIPALS.TSV:

Statistic	Value	Implication for EDA
Total Records	\$\approx 500,000\$	This is a significantly smaller, curated list compared to the \$\approx 12\$ million records in TITLE.BASICS.TSV. It likely only lists the top 10 or so principals per title.

Statistic	Value	Implication for EDA
Columns	6	It is a highly detailed linking table, containing two foreign keys (tconst and nconst) and four descriptive fields for the relationship.

## Column Analysis:

Column Name	Current Type	Data Quality Issue	Required Transformation
tconst	object	<b>Title Join Key.</b> The unique title identifier is loaded as a generic string.	<b>Casting:</b> Convert to category for efficient joining with TITLE.BASICS.TSV.
nconst	object	<b>Personnel Join Key.</b> The unique person identifier is loaded as a generic string.	<b>Casting:</b> Convert to category for efficient joining with NAME.BASICS.TSV.
ordering	int64	<b>Rank/Importance.</b> Indicates the order in which the principal is credited.	<b>Casting:</b> Convert to a smaller integer type (e.g., int8) if the maximum ordering value is small, to save memory.
characters	object	<b>String Array.</b> Contains character names, often encapsulated in an array-like string (e.g., ["Self"]).	<b>Parsing &amp; Cleaning:</b> Must clean/parse the array brackets ([ ]) and quotes) and split the string to extract readable character names for analysis.
job	object	<b>Specific Role.</b> Provides further detail beyond category (e.g., if category is 'producer', job might be 'executive producer').	<b>Cleaning &amp; Casting:</b> Handle \$\"text{NaN}\$ values and convert to category to analyze job distribution.
category	object	<b>Broad Role.</b> Specifies the broad type of	<b>Casting:</b> Convert to category for efficient

Column Name	Current Type	Data Quality Issue	Required Transformation
		principal (e.g., 'actor', 'director', 'writer').	analysis of role distribution.

## Relationship Analysis:

Relationship	Join Type	Records Involved	Records Matched	Interpretation
<b>Basics</b> \$\rightarrow\$ <b>Ratings</b>	Primary Key Join	\$12.1M \$\rightarrow\$ 1.6M	\$1.6\$ Million	Confirms that only \$13.2\%\$ of all titles have associated rating data.
<b>Basics</b> \$\rightarrow\$ <b>Crew</b>	Left Join Potential	\$12.1M \$\rightarrow\$ 12.1M	\$12.1\$ Million	Confirms a <b>1:1</b> relationship. Every title in the basics file has a corresponding entry in the crew file (even if director/writer fields are null).
<b>People ID</b> (\$\text{nconst}\$)	Lookup Key	\$14.9\$ Million Total	N/A	This key will be used to enrich the final merged dataset by joining with name_basics to retrieve primary name and profession.

The three-step merge strategy prioritizes quality and measurable performance over simple coverage of all titles.

Step	Operation	Key	Output Dataset	Analytical Justification
1	<b>title_basics INNER JOIN ratings</b>	tconst	\$\approx 1.6\$ Million Rows	<b>Filters for Quality.</b> The INNER JOIN discards ~87% of the data, retaining only titles for which quality and popularity metrics (\$text{averageRating}\$, \$text{numVotes}\$) are available.
2	<b>Result LEFT JOIN crew</b>	tconst	\$\approx 1.6\$ Million Rows	<b>Adds Personnel.</b> Ensures all 1.6 million rated titles retain their records and gain director/writer IDs. Missing director/writer IDs will appear as \$NaN\$ in the crew columns.
3	<b>Result JOIN name_basics</b>	nconst	Varies (Post-Expansion)	<b>Enriches Data.</b> This final step is essential for converting numerical IDs into readable names, making the crew analysis usable (e.g., director nm0005690 becomes 'William K.L. Dickson').