

Data Mining: Fundamentals

Group 12

Bruno Barbieri, Noemi Dalmasso, Gaia Federica Francesca Ferrara

The aim of this report is to display an analysis carried out on the IMDb dataset; the analysis has been conducted making use of data mining methodologies. After the data understanding and preparation phase, clustering, classification, and pattern mining techniques have been applied.

Contents

1	Data Understanding and Preparation	2
1.1	Data Semantics	2
1.2	Distribution of the variables and statistics	2
1.2.1	Discrete attributes	2
1.2.2	Continuous attributes	4
1.3	Data Quality	4
1.3.1	Syntactic Inconsistencies	5
1.3.2	Missing Values	5
1.4	Variable Transformation	5
1.5	Pairwise correlations and elimination of variables	6
2	Clustering	8
2.1	Centroid-based methods	8
2.2	Density-based methods	8
2.3	Analysis by hierarchical clustering	8
2.4	General considerations	8
3	Classification	9
4	Regression	10
5	Pattern Mining	11

1. Data Understanding and Preparation

1.1 Data Semantics

The dataset *train.csv* contains 16431 titles of different forms of visual entertainment that have been rated on IMDb, an online database of information related to films, television series etc. Each record is described by 23 attributes, both numerical and non-numerical. All the variables of the dataset are introduced and explained in Table 1.1 and Table 1.2.

Attribute	Type	Description
originalTitle	Nominal	Title in its original language
rating	Ordinal	IMDB title rating class The range is from (0,1] to (9,10]
worstRating	Ordinal	Worst title rating
bestRating	Ordinal	Best title rating
titleType	Nominal	The format of the title
canHaveEpisodes	Nominal (Binary)	Whether or not the title can have episodes True: can have episodes; False: cannot have episodes
isRatable	Nominal (Binary)	Whether or not the title can be rated by users True: it can be rated; False: cannot be rated
isAdult	Nominal (Binary)	Whether or not the title is for adults 0: non-adult title; 1: adult title
countryOfOrigin	Nominal	The country(ies) where the title was produced
genres	Nominal	The genre(s) associated with the title

Table 1.1: Description of non-numerical attributes

1.2 Distribution of the variables and statistics

This section will give an overview about the distribution of variables that has been carried on to understand patterns, detect meaningful statistics and assess their relevance to the project.

1.2.1 Discrete attributes

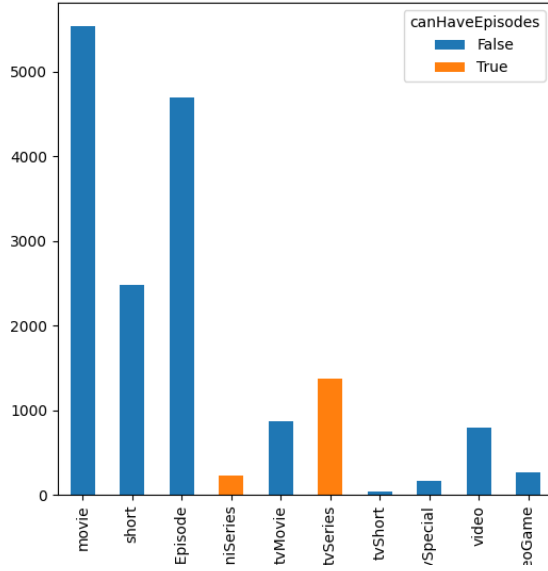
DA CAMBIARE IN BASE AI GRAFICI CHE DECIDIAMO DI TENERE In this paragraph, the most informative discrete attributes of the dataset are examined to provide an

Attribute	Type	Description
runtimeMinutes	Numeric	Runtime of the title expressed in minutes
startYear	Interval	Release/start year of a title
endYear	Interval	TV Series end year
awardWins	Numeric	Number of awards the title won
numVotes	Numeric	Number of votes the title has received
totalImages	Numeric	Number of Images on the IMDb title page
totalVideos	Numeric	Number of Videos on the IMDb title page
totalCredits	Numeric	Number of Credits for the title
criticReviewsTotal	Numeric	Total Number of Critic Reviews
awardNominationsExcludeWins	Numeric	Number of award nominations excluding wins
numRegions	Numeric	The regions number for this version of the title
userReviewsTotal	Numeric	Number of User Reviews
ratingCount	Numeric	The total number of user ratings for the title

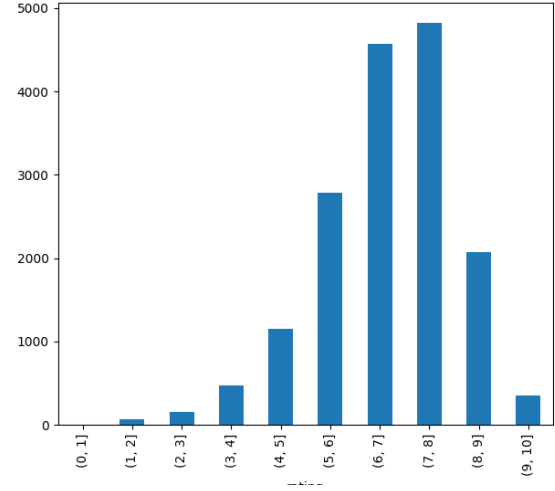
Table 1.2: Description of numerical attributes

overview of their statistics and frequencies.

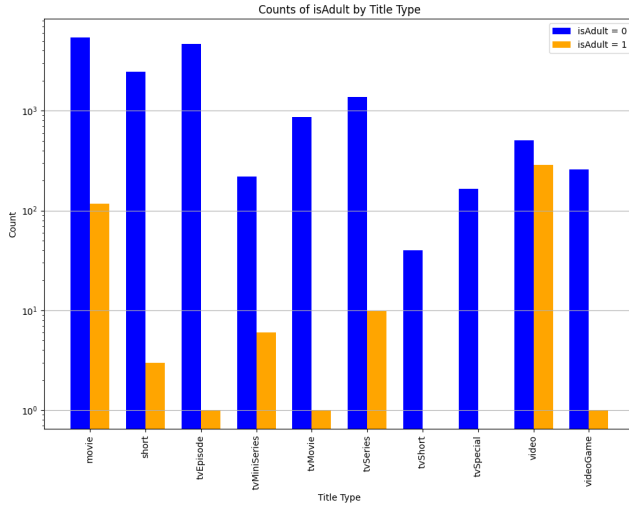
From Fig.1.1(a), it is observed that the classes of the **titleType** attribute are unbalanced, with *movie* being the most frequent class (5535 records) and *tvShort* the least frequent (40 records). By analyzing the **canHaveEpisodes** attribute within these **titleType** values, it is found that only *tvSeries* and *tvMiniSeries* can have episodes, as expected. As shown in Fig.1.1(b), the frequency of rating classes is slightly skewed toward higher values, with the most frequent rating class being (7, 8], which is the rating of 4822 titles. Another important aspect is that all 16341 titles are ratable and the vast majority of them (16005) are non-adults contents, as shown in Fig.1.1(c) Finally, as indicated in Fig.1.1(d), an analysis of the **genres** variable across different **titleType** values reveals that *Drama* and *Comedy* are the most common genres, as they appear in the top 3 genres of nearly every *titleType* category.



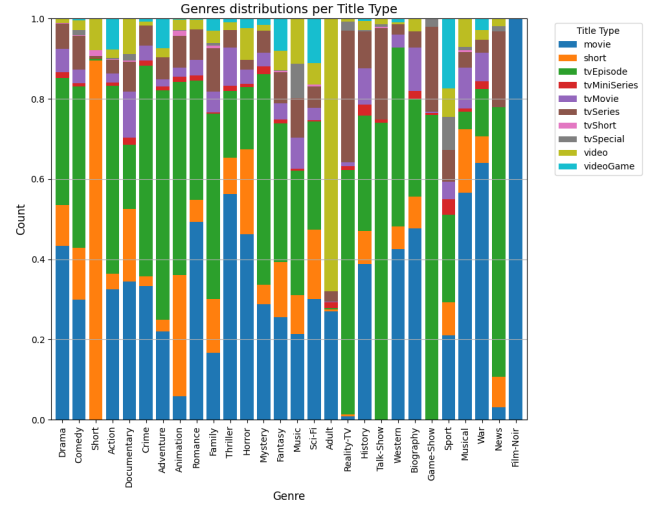
(a) Counting of the title types frequencies



(b) Counting of ratings frequencies



(c) Counting of the adult and non-adult per type



(d) types per genre

Figure 1.1: Bar chart of the discrete attributes.

1.2.2 Continuous attributes

...content...content...content...content...content...content...

1.3 Data Quality

In this phase, a proper evaluation of the observed data was conducted in preparation for the analysis. Once having checked that there are no duplicates and no incomplete rows in the dataset, attention was given at identifying missing values and outliers within the columns.

1.3.1 Syntactic Inconsistencies

In the exploration of the dataset it has been noticed that `awardWins` was the only feature having missing values identified with NaN. However, there were missing values also in other columns (`endYear`, `runtimeMinutes` and `genres`), but they were indicated with the string "\N". To avoid this inconsistency causing problems during data preparation, these values have been replaced with NaN. By doing so, any cell in the `endYear`, `runtimeMinutes` and `genres` column that previously contained the string "\N" is now considered a proper missing value, detectable and manageable using Pandas' functions.

1.3.2 Missing Values

Once having solved the above-mentioned inconsistency, the resulting total amount of the missing values are the following, also represented in percentages for a better understanding:

- `endYear`: it is the feature with the highest number of NaN values (15617; about 95%). To handle them it has been decided to **COSA FARE?**;
- `runtimeMinutes`: it has 4852 missing values (29.5%) that have been handled by grouping the records by `titleType` and substituting the NaN value with the median of each group;
- `awardWins`: this feature has 2618 NaN values (about 16%). Since the mode associated with this variable is 0, it has been decided to substitute the missing values with 0;
- `genres`: it is 382 missing values (2.3%). Having dealt this variable with a multi-label one-hot encoding process (as will be described in the *Variable Transformation* section), a vector of all zeros is assigned to record with missing genres values.

1.4 Variable Transformation

As the first step in the variable transformation process, the `CountryOfOrigin` and `genres` variables were converted from strings into lists of strings to facilitate further analysis. This transformation was necessary because some records contain multiple genres or countries as values for these variables. Multi-label one-hot encoding was applied to the `genres` column; each unique genre was represented as a binary feature, allowing records that belong to multiple genres simultaneously to maintain this information. Furthermore, it has been decided to extract the ceiling value for each entry in the `rating` column, in order to use it as an integer for further analysis.

For the numeric attributes, it was observed that some variables required stronger transformations due to their highly positively skewed distributions. Specifically, a log-transformation was applied to variables with a skewness greater than 10, **while a square root transformation was deemed**

reasonable for the variable with a skewness of 4 (although this step has not yet been implemented). This approach was chosen because the standard normalization techniques, such as `MinMaxScaler` and `StandardScaler` provided by `sklearn.preprocessing`, were insufficient to achieve an adequate normalization of the dataset. These scalers are more effective when applied after stronger preprocessing steps aimed at addressing extreme skewness.

1.5 Pairwise correlations and elimination of variables

The plot in figure 1.2 is a Pearson's correlation matrix that takes into account the continuous numerical variables of the dataset. For what can be observed, `numVotes` and `ratingCount` have a perfect positive correlation, and so it would be redundant to keep them both. For this reason it has been decided to drop **DECIDERE QUALE DROPPARE**. On the other hand, regarding categorial attributes, after having analized them, **METTERE VAR CHE TOGLIAMO: es. best/worst rating, isatable, endyear(???)** were found to have limited contribution **explain why** based on their distributions and were excluded to the dataset. **DA FINIRE**

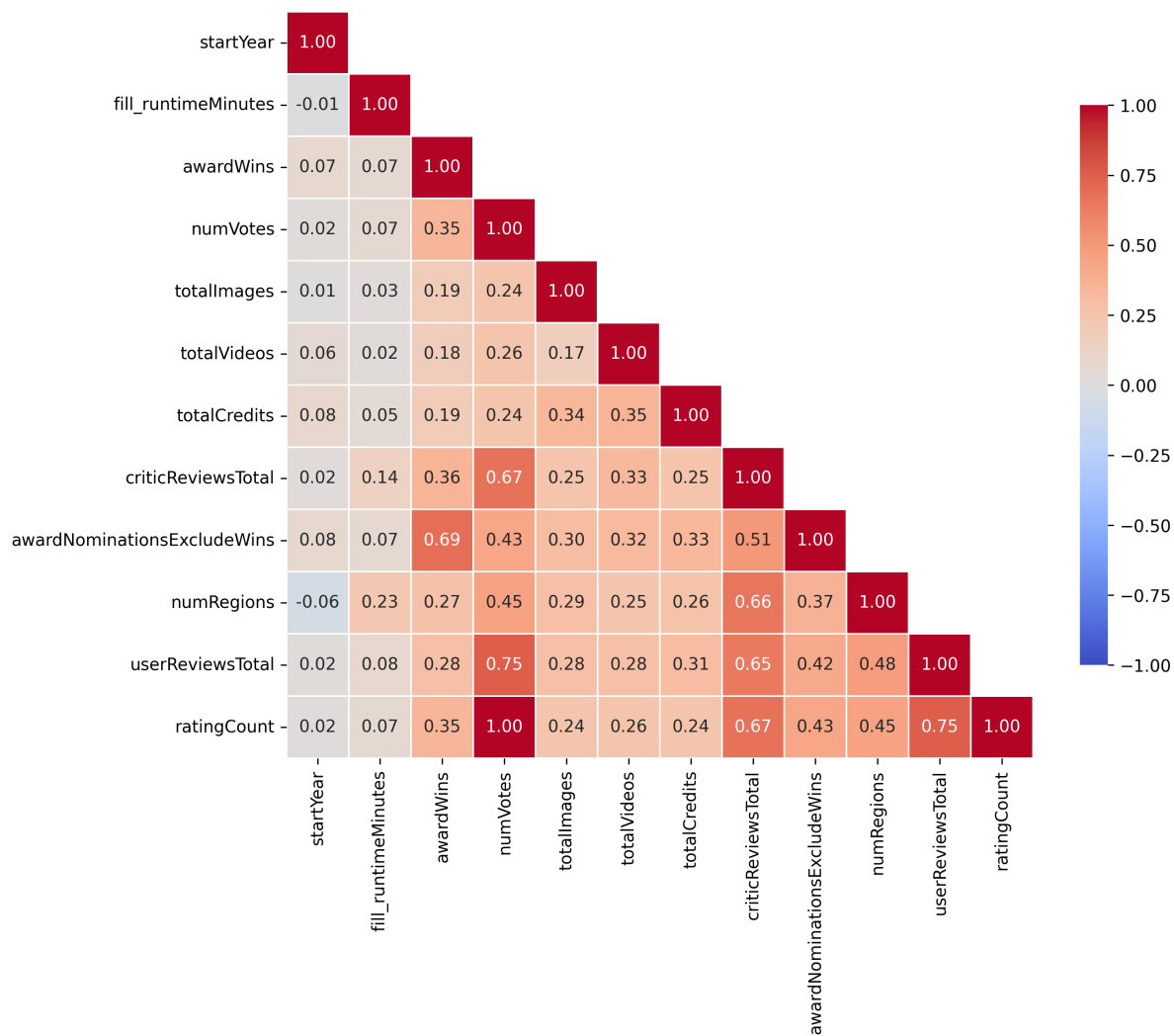


Figure 1.2: Correlation matrix

2. Clustering

2.1 Centroid-based methods

2.2 Density-based methods

2.3 Analysis by hierarchical clustering

2.4 General considerations

3. Classification

4. Regression

5. Pattern Mining