



Unlocking Viewer Insights in Linear Television: A Machine Learning Approach

Javier Carreno¹(✉), Khuong An Nguyen¹, Zhiyuan Luo¹, and Andrew Fish²

¹ Royal Holloway University of London, Surrey TW20 0EX, UK

Javier.Carreno.2023@live.rhul.ac.uk,

{Khuong.Nguyen,Zhiyuan.Luo}@rhul.ac.uk

² University of Liverpool, Liverpool L69 3BX, UK

Andrew.Fish@liverpool.ac.uk

Abstract. Amidst the digital transformation, traditional linear TV faces major challenges, including fragmented viewership, fixed schedule, and inaccurate targeting. Therefore, this paper proposes a novel Machine Learning framework to understand the audience's demographics from their viewing behaviour. By employing state-of-the-art classification models on an extensive TV first-party dataset, we achieved an average 88.6% accuracy in correctly identifying each household demographics. Our result offers promising outcomes for refining strategies within linear TV to improve viewer engagement, content programming, and market insights.

Keywords: Audience Insights · Machine Learning · Household Classification

1 Introduction

Linear television (TV) represents the traditional approach to broadcasting, where TV networks adhere to predetermined schedules, airing specific content at scheduled times for viewers [2]. Unlike modern on-demand or streaming services, linear TV restricts viewer control over content access, compelling them to tune in during scheduled broadcasts. While this model has historically been effective in reaching broad audiences, it presents challenges in gaining detailed insights into viewership demographics [5, 7].

Understanding household demographics in linear TV goes beyond advertising effectiveness. With a nuanced understanding of their audience's demographics, broadcasters can curate content that resonates deeply with specific segments, fostering viewer satisfaction, loyalty, and engagement [3]. These are crucial for maintaining audience retention and competitiveness.

Integrating first-party data sourced from Freeview TV addresses limitations in traditional TV demographics. Acquired directly through viewers' connected TVs using the standard *Hybrid Broadcast Broadband TV* [10], this data offers valuable insights into viewers' habits and preferences.

Thus, this paper aims to utilise Machine Learning (ML) techniques to analyse first-party data and understand viewer profiles. In doing so, we aim to address the research question: **How accurately can ML be employed with just first-party TV data to classify household demographics?**

The contributions of our paper are:

- A **Machine Learning-ready dataset** with around 20,000 categorised devices from TV first-party data. It also includes a detailed household taxonomy and is available publicly at <https://github.com/carrenyo/TV-Viewer-Demographics-Machine-Learning> for further research.
- A **detailed pipeline** describing the entire process from initial *data collection* to the *final classification* of devices, for further improvements in the field.
- The **baseline results achieved** with state-of-the-art machine learning models to quantitatively demonstrate the feasibility of our approach.

The rest of the paper is organised as follows. Section 2 elaborates the *analytical pipeline* to collect the dataset and extract insights from it. Section 3 describes the dataset, and Sect. 4 details various ML algorithms and the baseline results on such dataset. Section 5 explains the related work. Finally, Sect. 6 summarises key insights and future research.

2 Pipeline

This section introduces a analytical pipeline involving 5 specific steps to extract valuable insights from TV first-party data (see Fig. 1).

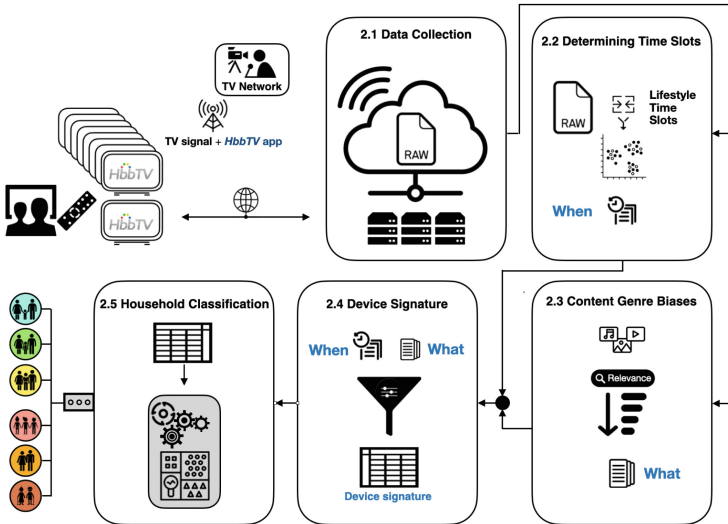


Fig. 1. Our proposed pipeline from data collection to device classification.

2.1 Data Collection

The first step involves collecting first-party data from Hybrid Broadcast Broadband TV (HbbTV), an open standard that integrates traditional broadcast TV (terrestrial, cable, or satellite) with internet connectivity, enhancing the TV experience to an interactive level [10]. With HbbTV, viewers using compatible devices can access additional content and interact with features using their remote controls [17]. The data collection from connected TVs involves an HbbTV application in the transport stream by TV networks. This embedded HbbTV app is unique as it does not need viewer installation. Viewer interactions while accessing TV network channels are systematically recorded and sent via the internet return channel to the cloud (see Fig. 2).

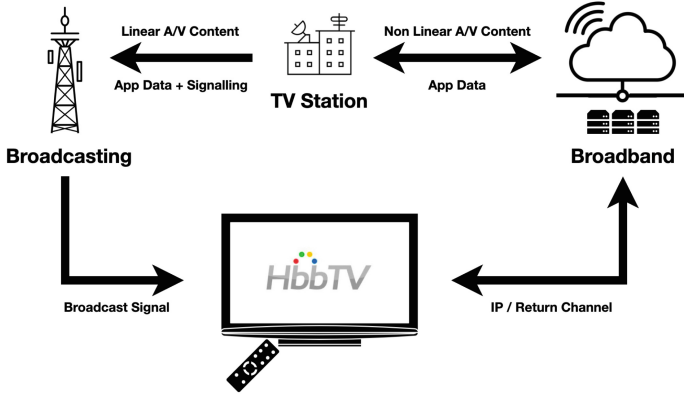


Fig. 2. HbbTV deployment diagram.

Our dataset includes essential details like *device ID*, *interaction timestamps*, *channel ID*, *IP Address*, and contextual information such as *device type* and *operating system version*. Most importantly, the data collection process includes a *consent management system* to comply with *GDPR regulations*. Viewers are explicitly asked for consent before their data is used for tasks beyond essential technical functions. Table 1 shows an excerpt of the raw data.

Table 1. A snapshot of first-party raw data.

did	start	end	dur	active	userAgent	resolution	IP	chID
0846...	1693853858596	1693853958427	100	true	Mozilla/5.0 (Web0...	1920, 1080	91.11...	4032
11ca...	1693853413068	1693853690133	111	true	Mozilla/5.0 (Web0...	1280, 720	154.4...	4032
1406...	1693852329199	1693853956001	1501	false	Mozilla/5.0 (Linu...	1280, 720	95.2...	4032
1a4a...	1693852958450	1693853058459	100	true	HbbTV/1.2.1 (+DRM...	1920, 1080	83.49...	4032
220d...	1693853109400	1693853686558	550	true	Mozilla/5.0 (Web0...	1920, 1080	92.17...	4032
257e...	1693853428111	1693853528019	100	true	HbbTV/1.4.1 (+DRM...	1920, 1080	170.2...	4032

2.2 Determining Viewer Time Slots

Using the raw first-party data above, this step aims to uncover the diverse TV viewership patterns, which are essential for effective audience understanding. A session starts when a device tunes in and ends when the viewer switches channels. Sessions lasting less than 300 s are labeled as *zapping* and are excluded from analysis, as they do not provide meaningful insights into viewer preferences, especially in linear TV where programmes typically have longer durations. Similarly, sessions lasting longer than 10,800 s are considered *extreme* and are also excluded (atypical behaviour).

The study then categorises devices based on their viewing consumption, considering *session frequency* and *total viewing time* on a weekly basis. Devices meeting specific criteria receive scores based on their weekly behaviour, with those having 20 or more sessions totaling at least 54,000 s getting a score of **3**, those with 5 or more lasting at least 7,200 a score of **2**, and those with 1 or more lasting at least 1,500 a score of **1**. These aggregated scores determine each device’s viewership level over the sample period, as follows.

- **fan** if the cumulative score exceeds 2.5 times the study’s duration in weeks.
- **regular** if the score surpasses 1.5 times.
- **occasional** if the score is more than 0.25 times.
- **no viewer** otherwise.

To capture the diverse viewer routines and lifestyles accurately, it is essential to establish cultural time slots linked to *meal times*, *work/school hours*, and *leisure periods*. Traditional 9am-5pm time slots are often too broad, making it difficult to discern specific viewing patterns like *lunchtime* or *mornings*. Hence, we divide each day into **seven time slots** for a more nuanced breakdown, effectively differentiating between *weekdays* and *weekend* (see Table 2).

Table 2. Our proposed finer-grained time slots for better distinction between weekdays and weekend for one specific Mediterranean European region.

Time Slot	weekdays	weekend
breakfast	07:00 - 08:30	07:00 - 10:00
morning	08:30 - 13:00	10:00 - 13:30
lunchtime	13:00 - 15:00	13:30 - 15:00
afternoon	15:00 - 17:00	15:00 - 17:00
evening	17:00 - 20:30	17:00 - 20:30
dinner	20:30 - 22:00	20:30 - 22:00
afterdinner-night	22:00 - 07:00	22:00 - 07:00

Once the time slots are established, each legitimate session is associated with one of them, and aggregated weekly per device. Our analysis involved two cluster

analyses: *weekdays* and *weekend* [19]. We used the *Silhouette score* method to determine the optimal number of clusters, focusing on cohesion and separation within clusters [16]. The goal was to maintain high scores consistently across weeks. Only devices with **regular** viewership patterns were included to avoid skewed results, excluding *occasional* and *fan* viewers. Applying the bisecting k-means algorithm over four weeks consistently identified **6 clusters** as optimal for weekdays (see Fig. 3).

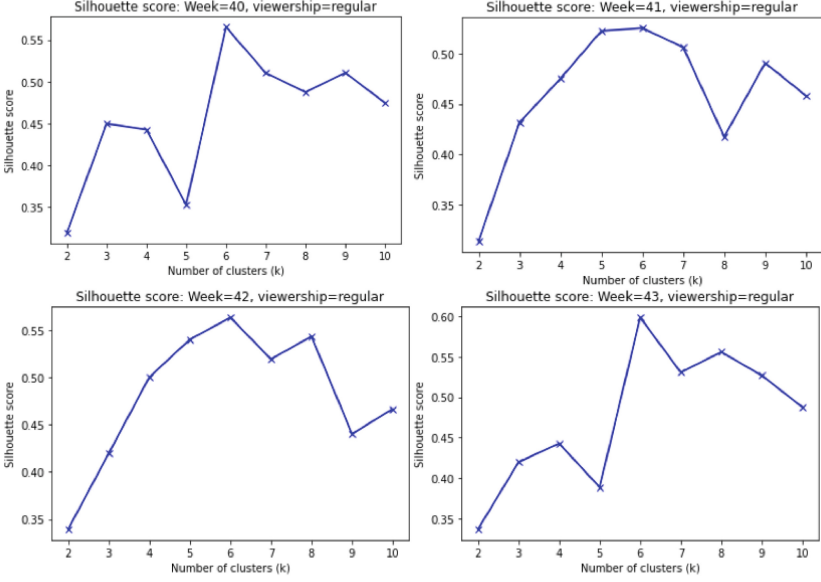


Fig. 3. Silhouette scores for weekdays cluster analysis, which indicated an optimal number of 6 clusters.

Then, we calculated **time slot ratios** for each cluster to highlight unique viewing behaviours.. The ratios were calculated by dividing the number of viewing sessions in each time slot by the total sessions for that cluster, showing what percentage of viewing occurs during different times of the day. Clusters with consistent viewing habits across the four-week period were merged into **prototypes**, while inconsistent patterns were excluded from further analysis. Ultimately, 5 *prototypes* emerged: **morning**, **lunchtime**, **afternoon**, **evening**, and **dinner**. Similar steps were taken for the second cluster analysis. Optimal results for *week-end* were also found with $k = 6$, yielding **five prototypes**: **morning**, **lunch**, **evening**, **dinner**, and **afterdinner-night**.

However, aligning weekly viewing behaviours with the identified clusters posed a challenge. The traditional k-means algorithm assigns all data points to centroids, which deviate from the primary objective of this analysis. Our goal

was to label viewing patterns only when closely resembling the defined prototypes. To address this challenge, a customised threshold was created, incorporating *standard deviation* and *Euclidean distance* calculations. Data points within 2 times the standard deviation (Z-score) were labeled as **close**, while others were marked as **far**. This approach ensured that only data points closely matching the prototypes were assigned to clusters, guaranteeing accurate categorisation.

Once common consumption patterns are identified and matched with device usage on a weekly basis, it becomes possible to determine **when** the viewers watch TV. The next step is to establish **what** they watch.

2.3 Content Genre Biases: Unveiling Viewer Preferences

This step aimed to uncover biases towards content genres, crucial for categorising devices into household types, as genres heavily influence audience preferences [2]. Weekly analysis quantified viewing volume per genre per device, based on manual content classification in Table 3. Merging genres formed cohesive groups, enhancing identification of viewer engagement and reducing sparse viewing tracking.

Table 3. Content taxonomy used in this study.

Category	Description
action & adventure	High-energy content with intense sequences and physical feats.
biography & historical events	Narratives about real-life individuals and historical occurrences.
children & family	Content suitable for family viewing, including kid-friendly themes.
comedy	Light-hearted, humorous content designed to entertain and amuse.
cooking & wellness	Programmes focused on culinary arts and general well-being.
crime & mystery & horror	Content centered on criminal investigations, mysteries, and horror.
current affairs & social issues	Programmes addressing contemporary events and societal concerns.
drama	Engaging emotional narratives delving into human experiences.
folk culture & heritage	Depictions of traditional customs, folklore, and cultural heritage.
game show & quizzes	Entertainment involving game formats and intellectual challenges.
magazines & talk shows	Shows featuring discussions, interviews, and topical segments.
media & popular culture	Exploration of media trends and elements shaping popular culture.
music & arts	Programmes dedicated to musical and artistic expressions.
news & politics	Informational content covering current news and political affairs.
romantic	Content centred around love stories and romantic relationships.
science & nature & animals	Exploration of scientific topics, nature, and wildlife.
sports	Coverage of various sports and related events.
thriller & suspense	Suspenseful and gripping narratives designed to captivate viewers.
travel & lifestyle	Programmes featuring travel destinations and diverse lifestyles.

For a reliable assessment, substantial data covering multiple weeks is crucial to minimise the impact of seasonal fluctuations in traditional TV content. These fluctuations arise from factors such as holiday-themed programming, viewer

shifts during major events, network schedule changes, and variations in viewer behaviour across seasons [2].

Determining viewer engagement with content genres often involves calculating the time ratio for each genre against total viewing time (*Relative Frequency Viewed*). However, unlike online platforms with *Video on Demand*, linear TV networks do not always offer all genres uniformly, curating their lineup for specific audiences throughout the day. Moreover, live events such as sports or award shows may also vary in availability, potentially biasing engagement metrics.

To address this challenge, we devised a genre ‘**Relevance**’ metric. It compares the time spent watching each genre against the total viewing time and normalises it by the proportion of time each genre was broadcasted compared to the total broadcasted time within the week. This method helps assess the significance of content genres relative to their availability, refining our understanding of viewer engagement.

2.4 Device Signature: Capturing Essential Features

Constructing a *device signature* involves condensing all the activity and traits of each device into a single representative row. To address the variability in TV consumption, it is important to consider that a device may exhibit different behaviours from week to week.

In relation to viewer engagement, while the ‘*Relevance*’ metric effectively gauges engagement with specific weekly genres, deriving ‘**Total Relevance**’ by summing these scores across weeks has significant limitations, especially with infrequently broadcasted genres. For instance, if a genre like *adventure* is scarcely broadcasted, available only for 6 weeks, and a viewer engages with the genre only once during those weeks, displaying a high ‘*Relevance*’ score in one week. When computing the ‘*Total Relevance*’ by summing the weekly values, the resulting ‘*Total Relevance*’ disproportionately impacts the analysis. To address this challenge, a new metric called the ‘**Relative Audience Engagement Index**’ (**RAE Index**) was introduced. The ‘*RAE Index*’ standardises the ‘*Total Relevance*’, mitigating this distortion.

$$\text{RAE Index} = \frac{\text{Weeks Genre Viewed}}{\text{Total Weeks Viewed}} \bigg/ \frac{\text{Weeks Genre Broadcasted}}{\text{Total Weeks Broadcasted}} \quad (1)$$

By using the ‘*RAE Index*’, the analysis ensures a more accurate representation of the viewer’s genuine content preferences across all genres, regardless of variations in their broadcast availability. The ‘**Weighted Relative Audience Engagement**’ (**WRAE**) was introduced to refine the relevance assessment, calculated by multiplying the ‘*Total Relevance*’ by the ‘*Relative Audience Engagement Index*’: $\text{WRAE} = \text{Total Relevance} \times \text{RAE Index}$

Additionally, a ‘**Normalised Weighted Relative Audience Engagement**’ (**nWRAE**) is computed for each genre by dividing it by the ‘**Total Weighted Relative Audience Engagement**’. This new metric is a valuable

tool for assessing the level of audience engagement and contributes to refining the understanding of viewer preferences.

$$nWRAE = \frac{\text{Weighted RAE of Genre}}{\text{Total Weighted RAE of all Genres}} \quad (2)$$

2.5 Household Categorisation: Timing and Content

The final step involves categorising devices based on their usage timing and preferred content genres, aligning with an established sociodemographic approach. This categorisation includes ‘Couple with young kids (0–8 years)’, ‘Couple with teenagers (9–17 years)’, ‘Couple with adult children (18+ years)’, ‘Only young adults (18–35 years)’, ‘Only middle-aged adults (36+ years)’, ‘Seniors (elderly/retired adults)’ [8, 15]. Devices are classified based on their usage time and preferred content genres.

3 Data Analysis and Feature Selection

Our data came from the first-party data of an European regional Freeview linear TV channel, serving a population of over 5 million with a high internet penetration rate of over 80%. The one-year sample consists of approximately 700 classified TV programmes.

This study focused on the top 19,386 devices, with 994 labeled as **fan** and 18,392 as **regular**. Table 4 summarises key dataset details: sample size (devices), data source, served population, collection duration, recorded sessions, unique devices, TV programmes, genres, and classification features.

Table 4. Summary of the dataset.

Number of Samples (Devices)	19,386
Population Served	Over 5 millions
Duration of Data Collection	52 weeks
Total Sessions Recorded	62,911,754
Unique Devices Connected via HbbTV	352,987
Number of TV Programmes	Approximately 700
Number of Genres	19
Total Features	138
Subset of Features Used for Classification	31

The normalised weighted audience measurement ($nWRAE$) reveals varied audience preferences across genres. Genres such as *news & politics*, *cooking*, *drama*, and *science* showed higher mean engagement levels, indicating sustained audience interest. Conversely, genres like *action & adventure*, *children & family*, exhibited lower mean engagement levels, suggesting comparatively subdued

viewer interest. Interesting patterns emerge in genres such as *crime & horror*, *game show*, and *romantic* genres. Despite lower mean engagement, these genres showed occasional high peaks in their maximum engagement values. These peaks suggest sporadic yet intensified audience interest in specific content within these genres (see Fig. 4).

3.1 Classification Results

Devices were categorised into household groups based on viewer behaviours and preferences using predetermined thresholds and business rules, complemented by insights from market research surveys and audience measurement panels [3, 7]. This heuristic approach facilitated the analysis of TV viewing patterns, providing practical insights for understanding viewer engagement. Table 5 illustrates the dataset’s distribution across various household labels.

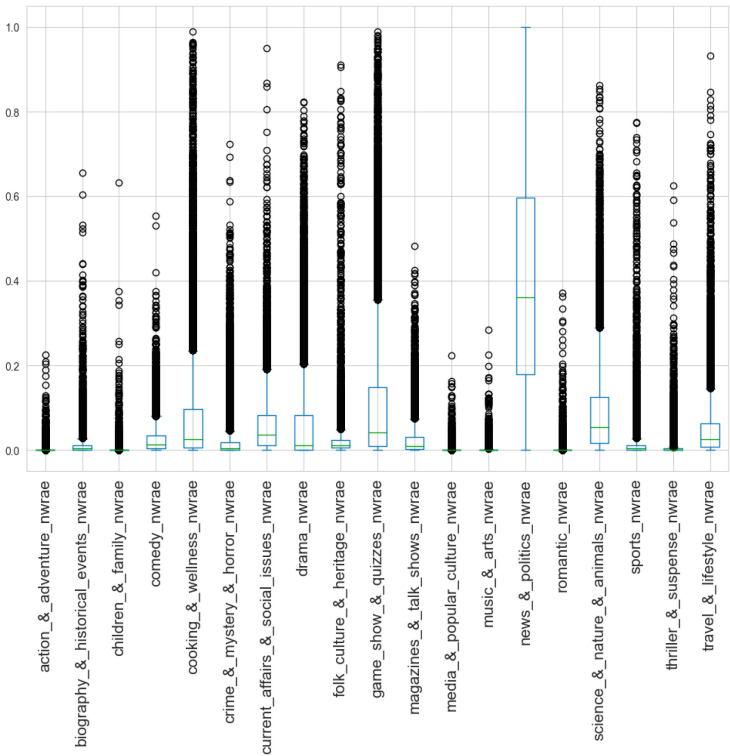


Fig. 4. Boxplots illustrating *nwrae* variation across genres. Only *news_&_politics* shows a distinct central tendency compared to other genres, with lower values and numerous outliers above.

Table 5. Household classification distribution.

Household Classification	Num of Devices	Percentage
Couple with young kids (0–8 years)	143	0.74%
Couple with teenagers (9–17 years)	1,514	7.81%
Couple with adult children (18+ years)	3,567	18.39%
Only young adults (18–35 years)	3,550	18.30%
Only middle-aged adults (36+ years)	6,985	36.02%
Seniors (elderly/retired adults)	3,627	18.74%

3.2 Feature Significance

Feature significance is crucial in Machine Learning, revealing attributes driving predictions. Random Forest model’s feature importance scores indicate each feature’s contribution to predictive performance (see Fig. 5).

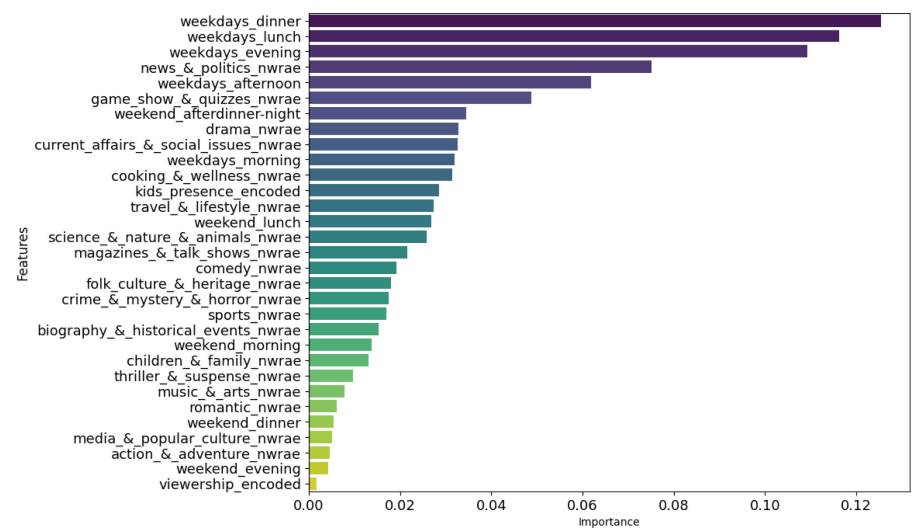


Fig. 5. Feature Importance for 31 ML features, highlighting their contribution to predictive accuracy and model performance.

Primarily, *weekdays_dinner* is the most influential factor, particularly during weekday dinner hours. Additionally, *news_ & politics_nwrae* and *weekdays_afternoon* also contribute significantly to the model’s predictive capacity.

Features related to **various genres** like *drama*, *game shows*, *current affairs*, and *cooking*, *travel* also demonstrate notable importance in predicting household classifications. The presence of certain content themes like *crime*, *sports*, *biographies*, *comedy* and *music* are influential but to a slightly lesser extent.

Per-class Feature Importance. Next, we analyse the importance of individual attributes for each classification category, highlighting the factors that significantly contribute to household differentiation (see Fig. 6).

4 Machine Learning Experiments and Results

In this section, we apply state-of-the-art Machine Learning models to our dataset to infer household demographics, including Random Forest, K-Nearest Neighbour, and Gradient Boosting. We aim to address the following questions:

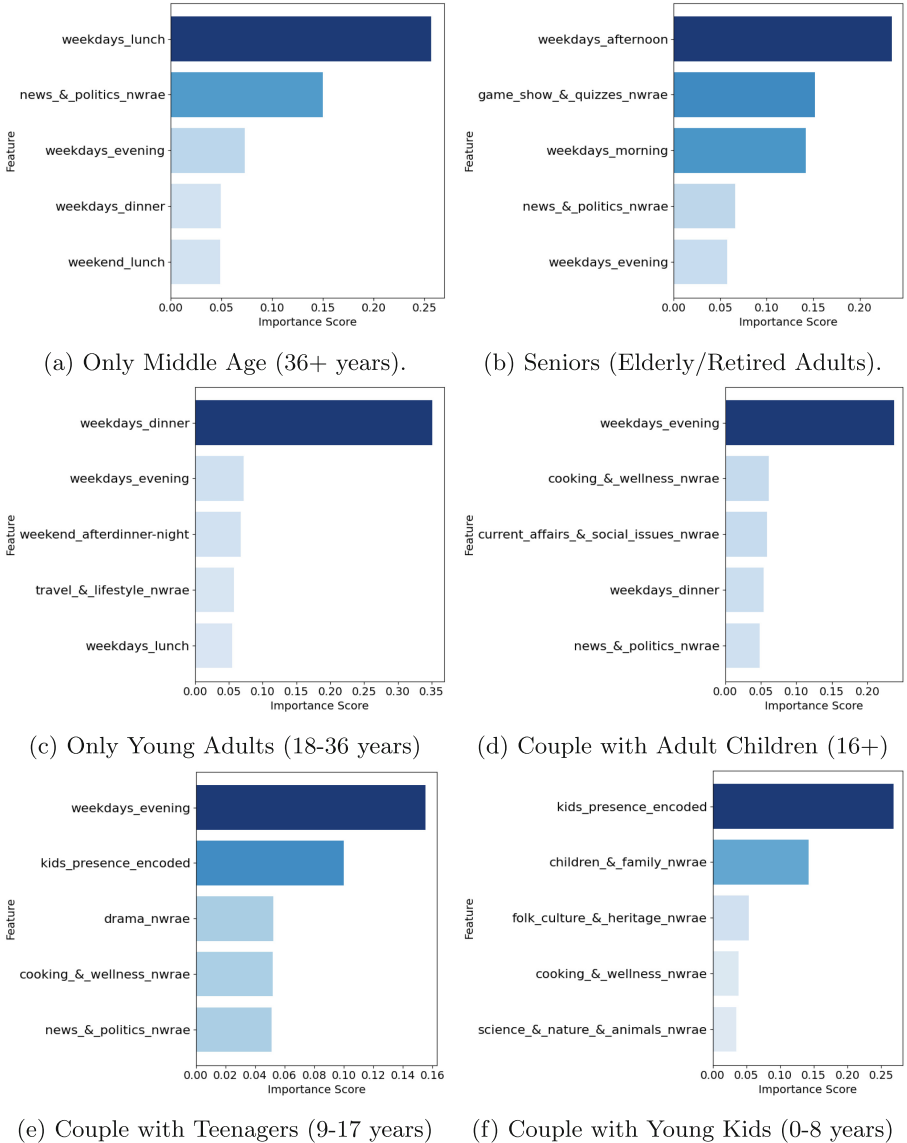


Fig. 6. Top features for each household classification.

- What are the accuracy in classifying household demographics?
- Which demographics are most challenging to classify?

4.1 Error Analysis: Confusion Matrices

In addition to the performance metrics previously discussed, the confusion matrices offers a more granular insight into the classification models' performance.

Random Forest. The model exhibited recurring misclassifications (see Fig. 7a), notably between '*Couple with adult children*' and '*Couple with teenagers*', and '*Only middle age*'. Families with adult children or teenagers might share TV viewing habits, leading to classification based on similar content preferences (**similar content consumption patterns**). Similarly, middle age and young adults may have overlapping lifestyle and viewing preferences, posing challenges in distinguishing between them based solely on viewing behaviours (**age-based similarities**). '*Couple with young kids (0–8 years)*' faces frequent confusion, likely due to diverse content preferences within this demographic (**complex family dynamics**).

K-Nearest Neighbour. The KNN model exhibited similar confusion patterns to other models, especially in misclassifying '*Couple with adult children*' and related categories (see Fig. 7b). However, it demonstrated balanced performance across most classes. Unlike the other models, it encountered difficulties in accurately classifying '*Couple with teenagers*' and '*Couple with young kids*', indicating distinctive challenges in capturing these specific categories accurately.

Gradient Boosting. The Gradient Boosting model showed similar confusion patterns to other models, particularly in distinguishing '*Couple with adult children*' accurately (see Fig. 7c). It performed moderately well with '*Couple with teenagers*' but aligned with other models in misclassification trends, indicating comparable behaviours in classifying different demographic groups.

4.2 Summary of Results

From the above results, **Random Forest** model emerged as the top choice due to its balanced performance across demographic categories (see Fig. 7). Despite **moderate confusion** between '*Couple with Adult Children*' and '*Couple with Teenagers*', it consistently **performed well** across most demographic groups. Moreover, it demonstrated **good stability** and **reliability** during cross-validation (see Table 6).

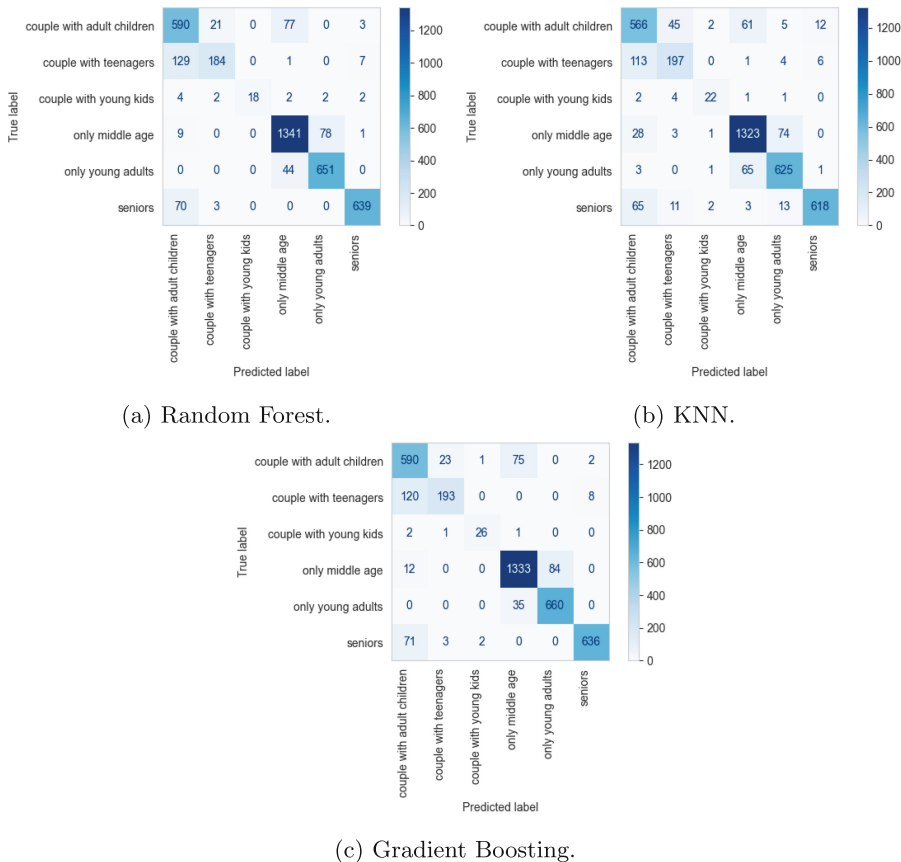


Fig. 7. Confusion matrices of 3 Machine Learning models, illustrating their accuracy and misclassifications, distinguishing false positives from negatives.

Table 6. Machine Learning Model Comparison with Cross-Validation.

Model	Fold Scores (Accuracy)					Mean	Precision	Recall	F1-Score
	Fold-1	Fold-2	Fold-3	Fold-4	Fold-5				
RF	0.8889	0.8723	0.8788	0.8801	0.8873	0.8815	0.89	0.88	0.88
KNN	0.8896	0.8741	0.8816	0.8865	0.8811	0.8826	0.87	0.86	0.86
GB	0.8961	0.8806	0.8814	0.8842	0.8883	0.8861	0.89	0.89	0.89

5 Related Works

The traditional method of studying viewer behaviour in linear TV uses structured panels like the UK’s 5,300-household panel, employing *Peplemeters* to

track viewership and programming. However, the limited panel size and concerns about participant reporting accuracy raise reliability questions.

Research often uses pay TV platforms for direct data collection, bypassing panel-based methods [9, 11]. However, the high costs of pay TV systems introduce biases toward wealthier households that can afford these services.

Conventional audience measurement methods with small panel samples often lack the precision needed for precise ad-specific viewership measurement. Advertisers, prioritising commercial audience over programme content, struggle to accurately assess ad effectiveness [1].

Connected TVs and online devices are reshaping TV advertising by leveraging behavioural data for advanced analysis, aiding digital evolution and the adoption of AI and sophisticated algorithms [13].

Understanding **selective exposure** is key for precise audience segmentation, going beyond generic demographics and predicting behaviour accurately [12, 14].

This study's **significant contribution** lies in its **dataset**, especially in ML and household classification tasks. It assesses the effectiveness and precision of classifying data into household taxonomy using diverse metrics, evaluation methods, and model performance assessments. Relevant literature offers insights into methodologies used to assess efficiency in comparable contexts [6, 18].

6 Conclusions and Future Work

This paper proposed a novel ML approach using extensive first-party linear TV data to unveil viewership behaviour. Analysis of the dataset's demographic distribution revealed a significant prevalence in the '*Only middle age (36+)*' segment, surpassing '*Seniors*', '*Couple with adult children (18+)*', and '*Only young adults (18–35 years)*'. Preferences impact viewer engagement, with timing factors like *weekdays_dinner*, *weekdays_lunch*, *weekdays_evening*, and *weekdays_afternoon* were crucial in household categorisation. This representation contrasts sharply with notably lower numbers in '*Couple with teenagers (9–17 years)*' and particularly '*Couple with young kids (0–8 years)*', each accounting for 1% of recorded devices. These statistics underscore modern viewing trends, especially among younger demographics, highlighting a **sharp decline** in their engagement with linear TV platforms.

Programme preferences emerged as a crucial factor, indicating increased consumption in genres like *news & politics*, aligning with the interests of middle-aged and adult demographics. Despite the overall decline in television viewership, the enduring trust in TV as a dependable information source remains evident, consistently positioning TV as the preferred medium for advertisers [4, 5].

Machine Learning performance varied in accuracy, showing recurring misclassifications, particularly between '*Couple with adult children (18+ years)*' and '*Couple with teenagers (9–17)*', suggesting shared content preferences. The Gradient Boosting model achieved the best result with an average accuracy of 88.7% across 5 data folds, demonstrating the potential of our approach in classifying household demographics with just first-party TV data.

This study explores TV viewership patterns, but future research could expand by incorporating children's programming from other channels for a comprehensive perspective. Adding viewer *surveys* could offer deeper insights beyond consumption patterns.

References

1. Abernethy, A.M.: Television exposure: programs vs. advertising. *Curr. Issues Res. Advertising* **13**(1–2), 61–77 (1991)
2. Barwise, P., Ehrenberg, A.: *Television and Its Audience*, vol. 3. Sage (1988)
3. Bondad-Brown, B.A., Rice, R.E., Pearce, K.E.: Influences on tv viewing and online user-shared video use: demographics, generations, contextual age, media use, motivations, and audience activity. *J. Broadcast. Electron. Media* **56**(4), 471–493 (2012)
4. Bruce, N.I., Becker, M., Reinartz, W.: Communicating brands in television advertising. *J. Mark. Res.* **57**(2), 236–256 (2020)
5. Bulgrin, A.: Why knowledge gaps in measurement threaten the value of television advertising: the best available screen for brand building is at a crossroads. *J. Advert. Res.* **59**(1), 9–13 (2019)
6. Choi, J.A., Lim, K.: Identifying machine learning techniques for classification of target advertising. *ICT Express* **6**(3), 175–180 (2020)
7. Clark, J., Paiement, J.F., Provost, F.: Who's watching TV? *Inf. Syst. Res.* **34**(4), 1622–1640 (2023)
8. Cohen, P.N.: *The family: diversity, inequality, and social change*. (No Title) (2018)
9. Deng, Y., Mela, C.F.: Tv viewing and advertising targeting. *J. Mark. Res.* **55**(1), 99–118 (2018)
10. Fischer, W., Fischer, W.: Broadcast over internet, HbbTV, OTT, streaming. In: *Digital Video and Audio Broadcasting Technology: A Practical Engineering Guide*, pp. 903–913 (2020)
11. Fudurić, M., Malthouse, E.C., Lee, M.H.: Understanding the drivers of cable TV cord shaving with big data. *J. Media Bus. Stud.* **17**(2), 172–189 (2020)
12. Knobloch-Westerwick, S., Meng, J.: Looking the other way: selective exposure to attitude-consistent and counterattitudinal political information. *Commun. Res.* **36**(3), 426–448 (2009)
13. Malthouse, E.C., Maslowska, E., Franks, J.: The role of big data in programmatic TV advertising. In: *Advances in Advertising Research IX: Power to Consumers*, pp. 29–42 (2018)
14. Malthouse, E.C., Maslowska, E., Franks, J.U.: Understanding programmatic TV advertising. *Int. J. Advertising* **37**(5), 769–784 (2018)
15. Scott, J., Treas, J., Richards, M.: *The Blackwell companion to the sociology of families*. Wiley (2008)
16. Shutaywi, M., Kachouie, N.N.: Silhouette analysis for performance evaluation in machine learning with applications to clustering. *Entropy* **23**(6), 759 (2021)
17. Tagliaro, C., Hahn, F., Sepe, R., Aceti, A., Lindorfer, M.: I still know what you watched last sunday: privacy of the HbbTV protocol in the European smart TV landscape. In: *30th Annual Network and Distributed System Security, NDSS 2023* (2023)
18. Vaccari, I., Chiola, G., Aiello, M., Mongelli, M., Cambiaso, E.: MQTTset, a new dataset for machine learning techniques on MQTT. *Sensors* **20**(22), 6578 (2020)
19. Yeo, J.: The weekend effect in television viewership and prime-time scheduling. *Rev. Ind. Organ.* **51**(3), 315–341 (2017)