

UNIVERSITAT DE BARCELONA

FUNDAMENTAL PRINCIPLES OF DATA SCIENCE MASTER'S
THESIS

Analysis of Content Diversity in News Recommendation Systems

Author:

Rubén BLANCO BORRÁS

Supervisor:

Oriol Pujol, Jordi Vitrià

*A thesis submitted in partial fulfillment of the requirements
for the degree of MSc in Fundamental Principles of Data Science*

in the

Facultat de Matemàtiques i Informàtica

January 17, 2026

UNIVERSITAT DE BARCELONA

Abstract

Facultat de Matemàtiques i Informàtica

Master in Fundamental Principles of Data Science

Analysis of Content Diversity in News Recommendation Systems

by Rubén BLANCO BORRÁS

This Master's Thesis studies informational diversity in news recommendation systems from a dual perspective: producer diversity and consumer-perceived diversity. The main objective is to analyze which diversity metrics are most suitable for evaluating journalistic content and how they can be applied to different datasets.

In a first phase, a study is conducted on various diversity metrics used in recommendation systems, taking as reference those proposed in the literature associated with Microsoft and using the *MIND Small* dataset. This analysis allows the evaluation of classical diversity and coverage metrics applied to news recommendation, as well as an understanding of their advantages and limitations in controlled environments.

In a second phase, these metrics are adapted and applied to a news dataset from the media outlet *3Cat*, aiming to evaluate the diversity of political topics present in the published content. In this context, a distinction is made between producer diversity, related to the ideological and thematic variety of the generated content, and consumer diversity, modeled through an activation metric that approximates the diversity effectively perceived by the reader.

The results allow for a comparison of both diversity perspectives and the analysis of potential mismatches between the informational supply and the diversity experienced by the user. The code developed for data processing and metric calculation can be found in the project repository (Blanco Borrás, 2026).

Acknowledgements

I would like to express my gratitude to my supervisors, Oriol Pujol and Jordi Vitrià, for their guidance, availability, and valuable contributions throughout the development of this project. Their expertise and judgment have been fundamental for the completion of this work.

I would also like to thank my mother for her constant support and encouragement during the Master's program, without which it would not have been possible to complete this academic stage.

Contents

Abstract	iii
Acknowledgements	v
1 Introduction	1
2 Background	3
2.1 Diversity in media	3
2.2 Explanation of the Metrics	3
2.2.1 Calibration (topics)	4
2.2.2 Calibration (complexity)	4
2.2.3 Fragmentation	4
2.2.4 Activation	4
2.2.5 Representation	5
2.2.6 Alternative voices	5
2.3 Classification of Metrics According to TechPolicy	5
3 Diversity measures proposal	7
3.1 Shortcomings of the <i>Alternative Voices</i> Implementation in <i>RADio</i>	7
3.2 Incorporation of Local Large Language Models for Entity Extraction	7
3.3 Improvement of Internal Diversity	7
3.4 Proposed Diversity Metrics	8
4 Experimental Settings	9
4.1 RADio Replication	9
4.1.1 Explanation of the Models	9
4.2 3Cat Dataset - Analysis of Political Diversity and Sentiment	10
4.2.1 Dataset Selection	11
4.2.2 NLP Processing and Entity Extraction	11
4.2.3 Obtaining Reference Lists from Wikidata	11
4.2.4 Ideological Classification and Statistical Aggregation	11
5 Results and discussion	13
5.1 Results on the Radio replication	13
5.1.1 Results and Comparative Analysis of Train Models	13
Direct Comparison between RADio and this Work	13
Selection of Optimal Training	13
Comparative Analysis	14
General Trend Analysis	15
5.1.2 Comparison of Metrics with RADio	15
Metric Ranges	16
Metric Comparison by Model	16
Comparative Analysis	16

5.2 Results on 3Cat	17
6 Conclusions	19
A Technical Requirements	21
A.1 <i>1_train_export</i>	21
A.1.1 Main Compatibility Requirements	21
A.1.2 Environment Installation Procedure	21
A.2 <i>2_metrics_RADio</i>	22
A.3 <i>3_3Cat</i>	22
B Model Metrics Comparison	23
C Images of the Results of 3Cat	25
Bibliography	27

Chapter 1

Introduction

Diversity is a central concept in news recommendation systems, yet it remains one of the most challenging aspects to define and measure precisely. In the media and data science literature, *diversity* has been used to describe a variety of phenomena, from the plurality of topics and viewpoints to the distribution of attention among content producers (Vrijenhoek et al., 2024). This conceptual ambiguity has historically hindered both systematic measurement and integration into algorithmic systems.

Recommendation systems play a crucial role in shaping user access to information on digital platforms. When these systems prioritize popular or similar content, users may be exposed to a limited range of information, potentially leading to filter bubbles or echo chambers. Increased exposure to varied and personalised recommendations has been linked in the literature to improved user engagement and retention, as systems that seek to balance diversity with relevance can influence the long-term value and continued use of a platform (Holtz et al., 2020).

Diversity can be analyzed from multiple perspectives. From the consumer's perspective, it refers to the variety of content to which a user is exposed, including topical diversity, exposure to different viewpoints, novelty, or serendipity. From the producer's perspective, diversity concerns the distribution of visibility across content creators or sources, preventing attention from concentrating on a small set of dominant actors.

Objectives

The objectives of this thesis are to:

- Analyze and evaluate different diversity metrics applied to news recommendation systems.
- Explicitly distinguish between producer diversity and consumer-perceived diversity.
- Adapt and apply these metrics to both the *MIND Small* dataset (MSNews, 2020) and a temporal subset of the *3Cat* news dataset, capturing political topic diversity and reader activation.
- Combine multiple diversity signals into a single objective metric that reflects both producer and consumer perspectives.

Contributions

This work contributes by:

- Replicating the RADio methodology (Vrijenhoek et al., 2024) and comparing its diversity metrics with our proposed measures.

- Introducing an alternative diversity metric that addresses observed limitations in RADio.
- Integrating local Large Language Models (LLMs) in the analysis pipeline to extract entities, sentiment, and other components necessary for the calculation of these metrics.

Layout of the Thesis

The remainder of the thesis is organized as follows:

- **Section 2: Background** – Reviews diversity in media and news recommendation, summarizing approaches from RADio and related literature.
- **Section 3: Diversity Measures Proposal** – Presents a critical analysis of RADio metrics and introduces the diversity measures proposed in this work, including the integration of LLMs for content analysis.
- **Section 4: Experimental Settings** – Describes the experimental design, including the replication of RADio (data, models, hyperparameters, and metrics) and the application to the 3Cat dataset.
- **Section 5: Results and Discussion** – Analyzes results on both RADio replication and 3Cat, discussing differences, limitations, and insights.
- **Section 6: Conclusions** – Summarizes findings and suggests directions for future work.

Chapter 2

Background

2.1 Diversity in media

Diversity in recommendation systems is a complex and multifaceted concept that can be approached from different perspectives, including consumer-perceived diversity and producer diversity. To evaluate and compare the effectiveness of recommendation models, it is necessary to use metrics that quantify these different dimensions of diversity. This section presents the metrics selected for this work, describing their theoretical foundation, practical application, and how they allow measuring both the thematic and viewpoint variety to which a user is exposed, as well as the distribution of attention among different content producers. The analysis of these metrics provides the basis for understanding how recommendations can promote a more diverse and balanced informational environment.

2.2 Explanation of the Metrics

This work uses several diversity metrics inspired by *RADio*, which attempts to capture five normative concepts in news recommendation: *Calibration* (*topics* and *complexity*), *Fragmentation*, *Activation*, *Representation*, and *Alternative voices*.

To measure diversity, article distributions are compared using two classical divergence metrics: KL-divergence and Jensen-Shannon Divergence (JSD). KL-divergence measures how much one distribution differs from another but can be unstable if there are zeros in the reference distribution. Therefore, JSD is used, which is bounded between 0 and 1 and does not explode when some probabilities are zero.

- $JSD = 0$ indicates that the compared distributions are identical.
- $JSD \rightarrow 1$ indicates that the distributions are completely different.

For all metrics, the interpretation of high or low values depends on the context:

- For metrics comparing with the impression of articles available to the user for clicking (*Activation*, *Representation*, *Alternative voices*), low values are better, as they indicate that the recommendation does not amplify biases and reflects the available content in a balanced way.

- For *Fragmentation*, a high value is positive, indicating that different users receive different recommendations, increasing per-user diversity.

- For *Calibration* (*topics* and *complexity*), low values are desired, indicating that the recommendation faithfully reflects the user's interests and reading complexity without relevant variations or biases.

2.2.1 Calibration (topics)

This metric evaluates the topical calibration of recommendations by calculating the frequency of each article's *category* in the recommended list and comparing this distribution with the user's *reading history*. Its goal is to measure whether the topic distribution in recommendations adequately reflects the user's interests.

The associated code simply counts the proportion of articles in each *category* and optionally weights these values according to the ranking position (higher-ranked articles have more weight). This aligns with the theoretical intent of the metric: ensuring that the topics most important to the user are proportionally represented in the most visible positions.

2.2.2 Calibration (complexity)

This metric measures the diversity in reading complexity of recommended articles, evaluated using the Flesch-Kincaid Reading Ease index:

$$\text{Reading Ease} = 206.835 - 1.015 \cdot \frac{\text{total words}}{\text{total sentences}} - 84.6 \cdot \frac{\text{total syllables}}{\text{total words}}$$

Where:

- Long sentences and words with many syllables decrease the score (more difficult text).
- Short sentences and simple words increase the score (easier text).

The metric calculates the reading complexity of recommended articles and optionally gives more importance to articles appearing in top ranking positions. It compares this distribution with the user's *reading history* to ensure that recommendations are similar to what the user typically reads in terms of complexity.

2.2.3 Fragmentation

The theoretical intent of *Fragmentation* is to measure whether recommended news for a user covers multiple distinct events or if several articles address the same event.

In the implementation, the recommendation stories for a user are compared with 5 *pools* of articles shown to other users using the same recommendation algorithm.

- Low value → multiple users receive similar recommendations (low inter-user diversity).
- High value → different users receive distinct recommendations (positive diversity).

2.2.4 Activation

Activation quantifies the variety of affective intensity of recommended articles using the absolute value of a sentiment analysis score. Recommended articles are compared with the pool of articles shown to the user.

- Low value → the recommendation maintains a similar sentiment distribution relative to the impression shown to the user.
- High value → the distributions differ, indicating some bias.

2.2.5 Representation

Representation analyzes the presence of political actors in articles and evaluates fairness in representation. In this work, it focuses on the U.S. with the *MIND Dataset*, considering only two parties: Republican and Democrat, calculating the percentage of mentions of each. Recommended articles are compared with the pool of articles shown to the user. Low JSD values indicate that the recommended content preserves the distribution of the user impression and does not add potential biases.

2.2.6 Alternative voices

Alternative voices aims to measure the presence of minority versus majority voices. In the theoretical implementation, a minority is any person identified by the NLP pipeline who does not have a Wikipedia or Wikidata page.

Recommended articles are compared with the pool of articles shown to the user. Low JSD values indicate that the recommended content preserves the distribution of the user impression and does not add potential biases.

2.3 Classification of Metrics According to TechPolicy

According to the TechPolicy article (TechPolicy Press, 2023), diversity metrics can be classified into four main categories:

- **Consumer Diversity:** reflects the user experience regarding the variety of content they consume. Example: *Calibration (topics)*, *Calibration (complexity)*, *Activation*.
- **Producer Diversity:** measures whether different producers or sources have equitable presence in recommendations. Example: *Representation*.
- **Content-/item-based Diversity:** evaluates the heterogeneity of recommended items based on attributes such as topic, complexity, or news type. Example: *Calibration (topics)*, *Calibration (complexity)*, *Fragmentation*.
- **Society- or platform-level Diversity:** analyzes the representation of different social, political, or cultural groups within the recommendation set. Example: *Representation*, *Alternative voices*.

The implemented metrics allow evaluating different aspects of diversity in news recommendations: from the user perspective (*Calibration (topics)*, *Calibration (complexity)*, *Activation*), the producer perspective (*Fragmentation*), and society-level (*Representation* and *Alternative voices*). The implementation generally aligns with the theoretical explanation of each metric, except *Alternative voices*, which has practical limitations.

Chapter 3

Diversity measures proposal

While RADio represents a relevant contribution to the measurement of normative diversity in news recommendation systems, a closer inspection reveals certain conceptual and methodological inconsistencies in the way diversity is operationalised. These limitations may affect the interpretability and robustness of the resulting metrics. This section therefore proposes an alternative formulation aimed at providing a more consistent and expressive measure of content diversity.

3.1 Shortcomings of the *Alternative Voices* Implementation in RADio

At the implementation level of *Alternative voices*, it is unrelated to the theoretical part, as it simply measures how many people have a *givenname* in their Wikidata information, and since the vast majority do (28,590 out of 29,977), the obtained metric values are very low, limiting its practical usefulness due to algorithmic limitations, as the dataset itself already identifies people and their corresponding Wikidata *QID*. For the metric to be useful, an NLP algorithm would need to identify people in the body of each news article (this part is absent in the *MIND Dataset* due to legal restrictions) and determine whether they exist in Wikipedia, creating a new *entities* column instead of using the default one where all have QID, but this is beyond the scope of this work, making the metric of limited practical value.

3.2 Incorporation of Local Large Language Models for Entity Extraction

In order to enrich the textual representation of news articles, local Large Language Models (LLMs) can be utilised for Natural Language Processing tasks. Specifically, Named Entity Recognition (NER) can be applied to each news article in order to extract relevant entities directly from the article text. These entities constitute an intermediate semantic layer that enables subsequent processing steps, including sentiment analysis, ideological attribution, and diversity-related computations. The use of local LLMs ensures greater control over the processing pipeline, reproducibility of results, and independence from external APIs.

3.3 Improvement of Internal Diversity

Instead of comparing with an external *pool/history/impression sample*, internal diversity of the recommendation could be measured using:

- Normalized Shannon entropy, where values close to 1 indicate maximum entropy/diversity and 0 minimum.

- ILS Diversity, based on similarity between articles within the recommendation, also normalized between 0 and 1, where values close to 1 indicate maximum diversity and 0 minimum.

This would allow evaluating the “pure” diversity of the recommended list, independent of the available pool diversity. Current metrics (KL or JSD) measure how much the recommendation preserves or distorts the original content distribution, but do not measure the actual internal diversity of the list. If the *pool/history/sample* is not diverse, comparing with it is meaningless. It does not measure pure diversity but rather relative to available content.

3.4 Proposed Diversity Metrics

Diversity in the news set can be assessed by calculating three complementary metrics using two ideological bins:

- **Normalized Shannon Entropy (H)**, which measures coverage diversity across ideological bins: $H = -\sum_i p_i \log p_i / \log 2$, where p_i is the proportion of mentions of bin i . The value is normalized so that 0 represents minimum entropy and 1 maximum.
- **Weighted Mean Sentiment (S)**, which represents the overall sentiment across both bins, to later penalize deviations from balance (0.0):

$$S = \frac{\sum_i (\text{count}_i \cdot \text{mean_sentiment}_i)}{\text{total_count}}$$

- **Normalized Polarization (P)**, which measures the sentiment difference between the two bins, independent of the number of mentions:

$$P = \frac{|S_{\text{left}} - S_{\text{right}}|}{2}$$

Finally, the three metrics are combined into a global diversity metric $HPS - D$:

$$HPS - D = H \cdot (1 - P) \cdot (1 - |S|)$$

where H corresponds to producer diversity, and S and P are variants of the consumer activation metric, capturing respectively the mean bias and polarization. This combination allows identifying hidden manipulations: for example, a positive mean sentiment in both bins could indicate bias, even if apparent polarization is low, and vice versa, if the mean sentiment is zero but both bins have S_i of opposite signs, generating sympathy toward one ideology while creating antipathy toward the other.

Chapter 4

Experimental Settings

This section describes the experimental settings used to evaluate the proposed diversity metrics.

4.1 RADio Replication

For the replication of RADio, each news article includes the following fields:

- **category**: the main topical category of the article.
- **subcategory**: a more specific subcategory within the main category.
- **story**: the title of the news article.
- **text**: the main body of the news article.
- **sentiment**: the overall polarity of the article.
- **entities**: key actors extracted via Named Entity Recognition (NER).

In the replication of RADio, the diversity of the news set is assessed using the metrics originally proposed by the authors in their article. These include **calibration**(*topic* and *complexity*), **representation**, **activation**, **alternative voices**, and **fragmentation**. A few modifications were made compared to the base code in the `WikidataHandler` class and in `representation.py` to prevent some errors. The rest is identical to the base repository, except for the prior file generation part, as these files are no longer available online and had to be generated from scratch.

4.1.1 Explanation of the Models

Below is a description of the models used in this news recommendation work:

- **Random**: generates a random order of recommendations. This model acts as a baseline reference value that other models must exceed.
- **Most Popular**: ranks news according to total number of views, prioritizing the most popular items.
- **LSTUR (Neural News Recommendation with Long- and Short-term User Representations)**: Captures both long-term user preferences and short-term interests. It:
 - Uses user ID embeddings to learn long-term representations.

- Processes recently read news through a GRU to learn short-term representations.
- **NRMS (Neural News Recommendation with Multi-Head Self-Attention):** A content-based news recommendation approach. Its main functionality includes:
 - Using multi-head self-attention to learn news representations, modeling interactions between words.
 - Learning user representations by capturing relationships among the news they have read.
 - Applying additive attention to select the most important words and news for informative representations.
- **NAML (Neural News Recommendation with Attentive Multi-View Learning):** A multi-view neural approach integrating diverse news and user information:
 - Uses title, body, category, and subcategory to obtain news representation.
 - Uses user behavior history to learn their representation.
 - Applies additive attention to learn informative representations, selecting the most relevant words and news.
 - Due to legal restrictions on the *MIND dataset*, summaries are used instead of full news bodies.
- **NPA (Neural News Recommendation with Personalized Attention):** A content-based method using personalized attention to improve news and user representations:
 - Learns news representations through a CNN.
 - Learns user representations from the news they clicked.
 - Applies personalized word-level attention to highlight important words for each user.
 - Applies personalized news-level attention to emphasize historically relevant news for each user.

4.2 3Cat Dataset - Analysis of Political Diversity and Sentiment

This section describes the process of analyzing political diversity with its associated sentiment carried out on the *3Cat* dataset. The objective is to quantitatively evaluate the ideological diversity of news published by this outlet and its affective activation in the reader, using a natural language processing approach with the *HSP – D* diversity metric proposed earlier, adapted to this case. For processing the *3Cat* dataset, each news article is represented using the following fields:

- **titol:** the headline or title of the news article.
- **cos:** the main body of the article.
- **entities:** key actors (political people or parties) extracted via Named Entity Recognition (NER).

4.2.1 Dataset Selection

For this study, a full month from the 3Cat dataset was selected in order to analyze the political diversity of the news. Due to computational limitations, a sample of 500 news articles was used, which allowed optimization of calculations and reduced processing time. However, with more powerful hardware resources, it would be possible to process the entire month without sampling.

4.2.2 NLP Processing and Entity Extraction

A natural language processing pipeline was built to identify, within each news article, entities of type *political person* and *political party*. Each entity was assigned a normalized sentiment value in the range $[-1, +1]$, corresponding to the context in which it appears in the news. For this, a local language model, Ollama Qwen3:8b, was used, taking approximately 30 seconds per news article on a mid-range 2025 computer.

Subsequently, the entities were processed to normalize their format, removing capital letters, accents, and other special characters to facilitate subsequent processing and comparison with external lists.

4.2.3 Obtaining Reference Lists from Wikidata

To correctly classify the detected entities, the following were downloaded from Wikidata:

- The list of political parties in Spain, along with their ideology (left, center, right, far-left/far-right) and official abbreviations.
- The list of political persons in Spain, including ministers, deputies, and other relevant positions.

These lists allowed verification of whether the extracted text entities corresponded to Spanish parties or political persons, and in the latter case, to associate them with their political party whenever possible.

4.2.4 Ideological Classification and Statistical Aggregation

Each entity classified as a party or political person was assigned to one of five ideological bins: far-left, left, center, right, and far-right. For each bin, the following were calculated:

- Number of mentions.
- Mean sentiment (S_i) in the range $[-1, +1]$.
- Standard deviation of sentiment.

To simplify the analysis, the left and far-left bins were merged into a single bin, the two right bins were also merged, and the center bin was removed to avoid bias, obtaining two final bins: *left* and *right*.

Chapter 5

Results and discussion

5.1 Results on the Radio replication

5.1.1 Results and Comparative Analysis of Train Models

This section describes the training process of the recommendation models used in this work, with the main objective of generating models that allow replicating and understanding the diversity metrics proposed by the authors of the *RADio paper* (Vrijenhoek et al., 2024) and analyzing their behavior in a controlled news recommendation environment. Different recommendation algorithms were applied to the ‘small MIND dataset’ (MSNews, 2020) using two training configurations: one with two *epochs* and another with four. For each model, the highest performance achieved across both configurations was selected as the final reference. For comparison purposes, a table is included summarizing the values reported in *RADio* along with those obtained in this work, highlighting the best result for each architecture.

Direct Comparison between RADio and this Work

Table 5.1 shows the NDCG@10 values reported by RADio and those obtained in this project after training the models for two and four *epochs*. Using a single comparative table allows for quick visualization of differences, the evolution with additional training, and the proximity of the resulting values to those published by RADio. Each row highlights in color the best result among the two training configurations applied in this work.

Algorithm	RADio	2 epochs	4 epochs
LSTUR	0.4134	0.4074	0.3943
NAML	0.4091	0.4019	0.4080
NPA	0.4068	0.3561	0.3306
NRMS	0.4163	0.3967	0.4006
Most Popular	0.2750	0.2882	0.2882
Random	0.2949	0.3263	0.3263

TABLE 5.1: Comparison of NDCG@10 between RADio and the models trained in this work. The best value obtained between 2 and 4 *epochs* is highlighted.

Selection of Optimal Training

Since the models were trained with two different configurations (2 and 4 *epochs*), the criterion adopted for final evaluation is to select, for each model, the highest value

among both. This allows a fair assessment of which architecture converges faster and which gains real benefit from extended training.

The observed patterns are as follows:

- **LSTUR:** slightly worsens from 0.4074 to 0.3943, being similar to the reference target set by RADio.
- **NRMS:** shows a clear trend of benefiting from additional training, achieving its best performance with four epochs: 0.4006.
- **NAML:** achieves its best value with four epochs (0.4080), almost identical to that published by RADio.
- **NPA:** is the model most affected by epoch variation: it drops from 0.3561 to 0.3306, indicating no benefit from additional training.

Consequently, each model is analyzed based on its best result, as shown in the table above, and for subsequent predictions, the model with the most optimal number of epochs will be used.

Comparative Analysis

Although the models were trained only for two and four *epochs*, the results obtained are remarkably close to the values published by RADio, especially for the more robust models. This suggests that part of the representational capacity of these architectures is acquired quickly, even with limited training, which is consistent with previous studies on attention-based architectures.

NAML: Practically Identical Results

The NAML model achieves a value of **0.4080** with 4 *epochs* in this work, practically identical to 0.4091 reported by RADio and to the 2 *epochs* value (0.4019). This may be due to:

- **Architectural robustness:** the hierarchical attention structure and use of multiple embeddings allow capturing useful patterns from early training phases.
- **Fast convergence:** word- and news-level attention models typically learn relevant relationships in a few iterations.
- **Broad semantic information representation:** combining title, body, category, and subcategory provides greater contextual diversity.

This behavior indicates that NAML is particularly efficient and stable, even with limited training.

NRMS and LSTUR: Moderate Differences

The values obtained for NRMS and LSTUR show that each model reaches its best performance under different configurations:

- NRMS: improves to **0.4006** with 4 *epochs* (vs. 0.3967 with 2 *epochs* and 0.4163 in RADio).

- LSTUR: achieves its best value with 2 *epochs*, **0.4074** (vs. 0.3943 with 4 *epochs* and 0.4134 in RADio).

The differences can be explained by:

1. NRMS may require more epochs to stabilize parameters.
2. LSTUR depends on the combination of long- and short-term user representations, making performance decrease if training is prolonged too much.
3. NRMS, based purely on self-attention, improves with additional training, though it does not reach RADio’s reference value.

NPA: The Most Epoch-Sensitive Model

NPA shows the largest relative difference between training configurations: **0.3561** with 2 *epochs* vs. 0.3306 with 4 *epochs* (RADio: 0.4068). This reflects its internal properties:

- Uses personalized word- and news-level attention.
- Requires a significant volume of user interactions to stabilize personalized embeddings.
- Is more susceptible to premature overfitting when trained without sufficient regularization or data.

Therefore, increasing to four epochs not only does not improve performance but may even degrade it.

General Trend Analysis

The joint analysis of RADio and the results of this work allows observing different patterns:

- **Models with greater multi-scale or hierarchical attention capacity stabilize earlier:** LSTUR and NRMS show relatively consistent values close to the best performance among configurations (LSTUR: 0.4074; NRMS: 0.4006).
- **Models dependent on personalized representations require more iterations:** NAML and NPA show larger variations between 2 and 4 *epochs* (NAML: 0.4019 → 0.4080; NPA: 0.3561 → 0.3306), reflecting sensitivity to the number of epochs.
- **Performance between 2 and 4 *epochs* does not always increase; in some models, it even decreases:** clearly observable in NPA and partially in LSTUR, where values decline if training exceeds the optimal duration.

5.1.2 Comparison of Metrics with RADio

This section compares the results obtained in this work on the *MIND small dataset* with those reported by the RADio repository on GitHub (Vrijenhoek, 2020), not from their paper, as values differ since the *paper* used the *large dataset* with many more iterations, and the repository used the *demo dataset*, even smaller than the *small dataset*. Values calculated in this work correspond to a subset of 10,000 users from the full

test set, which contains over 73,000 users. Parts of another repository (Halwesit, 2020) were also used in the implementation. To facilitate comparison and obtain more intuitive values, only *Jensen-Shannon divergence* is used, as it is normalized and allows direct distribution comparison.

Metric Ranges

Metric	RADio (min-max)	This work (min-max)
Calibration Topic	0.000 – 0.994	0.000 – 0.994
Calibration Complexity	0.000 – 0.994	0.000 – 0.994
Fragmentation	0.355 – 0.994	0.440 – 0.994
Activation	0.000 – 0.982	0.000 – 0.982
Representation	0.000 – 0.948	0.000 – 0.914
Alternative Voices	0.000 – 0.608	0.000 – 0.138

TABLE 5.2: Metric value ranges for RADio and this work.

Metric Comparison by Model

The average metrics for each model, comparing results obtained in this work with those published by *RADio*, are presented in **Appendix B** (values rounded to 3 decimals). The comparison is performed on a subset of 10,000 users from the test set (over 73,000 *behaviors*).

Comparative Analysis

Comparing each model with *RADio* reveals the following trends:

- **Calibration Topic and Complexity:** In general, values in this work are slightly lower than *RADio* but maintain the same relative trend across models. This indicates that topic distribution and article complexity are comparable.
- **Fragmentation:** In all models, obtained values are higher than those reported by *RADio*, suggesting that recommendations in this work tend to concentrate on a few stories per user, showing lower event diversity within each list.
- **Activation:** Mean values are slightly lower than *RADio*, indicating that the emotional intensity of recommended articles is slightly more moderate.
- **Representation:** This work shows significantly lower values than *RADio*, reflecting reduced presence of political actors in recommendations.
- **Alternative Voices:** Values are very low compared to *RADio*, consistent with implementation limitations: most identified people have Wikidata pages, reducing the proportion of alternative voices. Differences between *RADio* and this work arise from different preprocessing of the *MIND Dataset* in generating the *articles.entities* column in this work compared to *RADio*.

Overall, although there are differences in *Fragmentation*, *Representation*, and *Alternative Voices*, the general metric trends are maintained and allow coherent comparison of this work’s results with *RADio*’s on a representative subset of the dataset.

5.2 Results on 3Cat

For the selected sample of 500 news articles from August 2024, the following results were obtained per ideological bin:

Ideological bin	Number of mentions	Mean sentiment	Standard deviation
Far-left	7	0.13	0.53
Left	162	0.07	0.53
Center	15	-0.04	0.59
Right	60	-0.09	0.46
Far-right	8	-0.41	0.46

TABLE 5.3: Sentiment statistics by ideological bin

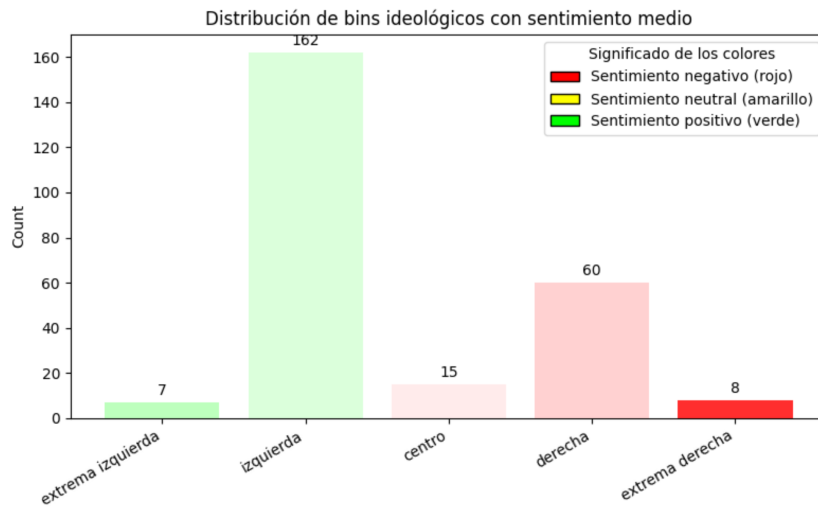


FIGURE 5.1: Distribution of mentions and mean sentiment per ideological bin

After merging the respective bins into *left* and *right* and removing the *center* bin, we obtain:

Ideological bin	Number of mentions	Mean sentiment
Left	169	0.07
Right	68	-0.13

TABLE 5.4: Sentiment statistics by ideological bin

Applying the formulas for entropy, mean sentiment, polarization, and diversity gives:

$$H = 0.8647, \quad S = 0.0151, \quad P = 0.1001, \quad D = 0.7665$$

Appendix C contains detailed figures supporting the analysis of the 3Cat dataset results.

These results indicate that, although the coverage between left and right is relatively balanced (H high), there is slight sentiment polarization (S and P), reflecting the usefulness of combining multiple metrics to accurately assess informational diversity.

The analysis of the 3Cat dataset reveals a clear polarization and imbalance between the two ideological bins. The left bin comprises 169 articles with a mean sentiment of 0.07, whereas the right bin contains only 68 articles with a mean sentiment of -0.13. Although the Shannon entropy (H) is relatively high, indicating moderate coverage balance, the disparity in article counts and sentiment highlights a notable skew in informational exposure. This underscores the importance of combining multiple metrics—entropy (H), polarization (P), and sentiment (S)—to fully capture the diversity of a news set.

It is important to note that normalized Shannon entropy is not linear with respect to the proportions of the bins. As H_n approaches 1, the difference between the two proportions decreases, tending toward perfect balance (50/50), as shown in Table 5.5. This illustrates how even slight imbalances can result in a high entropy value, while true balance requires H close to 1.

H_n (Normalized Entropy)	Proportion Bin 1	Proportion Bin 2
0.90	31.58%	68.42%
0.95	36.91%	63.09%
0.98	41.25%	58.75%
1.00	50.00%	50.00%

TABLE 5.5: Proportions of the two bins corresponding to different values of normalized entropy H_n .

Therefore, to ensure fair and balanced media representation, it is recommended that outlets achieve a minimum HSP-D score of 0.98. This can be obtained if both mean sentiment (S) and polarization (P) are close to zero, and the Shannon entropy is around 0.98, corresponding to approximate proportions of 41/59 between the ideological bins. This criterion provides a quantitative standard of fairness: even with unequal numbers of articles, a high HSP-D can be achieved if coverage is sufficiently broad and sentiment biases are minimized.

Chapter 6

Conclusions

In this work, we first replicated the RADio study, reproducing its methodology for measuring informational diversity in news recommendation systems. Building on this foundation, we proposed improvements to the original diversity metrics, aiming to capture both producer coverage and consumer-perceived diversity more accurately. We then applied the enhanced metric to the 3Cat dataset, allowing a detailed evaluation of ideological and topical diversity across multiple news outlets. The results highlight the usefulness of composite metrics such as HSP-D, which provide a robust, interpretable, and comparable measure of diversity across different media sources and topics.

Looking ahead, future work could extend this approach in several directions. Incorporating user behavior data would enable the design of recommendation systems that actively maximize diversity, rather than only measuring it. Additionally, applying the metric over consecutive months would allow temporal analyses, offering insights into how diversity evolves over time and how producer coverage and consumer-perceived diversity change across different periods. These extensions could help better understand the dynamics of news ecosystems and guide interventions to promote more balanced information exposure.

Appendix A

Technical Requirements

A.1 *1_train_export*

To ensure the correct reproduction of the experiments and avoid incompatibilities between the different libraries used, it was necessary to build a specific execution environment to guarantee the proper functioning of the *Recommenders* package (Recommenders Team, 2020) version 1.2.1, thus avoiding incompatibilities and errors derived from unsupported dependencies. In particular, the recommendation models used (*LSTUR*, *NRMS*, *NAML*, and *NPA*) depend on Microsoft's *recommenders* library, which is strongly tied to specific versions of *TensorFlow*, *CUDA*, and *cuDNN*. Any deviation from these versions causes compilation errors, model loading failures, or incompatibilities with *NumPy*.

A.1.1 Main Compatibility Requirements

- **Python < 3.10.** TensorFlow 2.10, the last version compatible with GPU via *pip* installation, only works correctly with Python 3.9 or lower.
- **TensorFlow 2.10.0.** This is the last version that includes official GPU support with precompiled binaries. Higher versions require custom builds.
- **CUDA 11.2 and cuDNN 8.2.0.** These are the exact versions that TensorFlow 2.10 expects. Using different versions (e.g., CUDA 12 or cuDNN 9) causes failures when initializing the GPU.
- **NumPy < 2.0.** Recent NumPy versions break compatibility with TensorFlow 2.x, so version 1.26.4 is fixed.
- **Binary installation of blis, thinc, and spacy.** These libraries may fail to compile from source on Windows systems, so they are installed exclusively in binary format.

A.1.2 Environment Installation Procedure

The complete process used to create the functional *reco_env* environment is detailed below:

```
--- 1. Creation and Activation of the Environment
>>> conda create -n reco_env python=3.9 -y
>>> conda activate reco_env

--- 2. Installation of C/C++ Packages
>>> pip install blis==0.7.9 --only-binary :all:
```

```
>>> pip install thinc==8.1.10 --only-binary :all:
>>> pip install spacy==3.5.4 --only-binary :all:

--- 3. Installation of 'recommenders' with GPU support
>>> pip install recommenders[gpu]

--- 4. Compatible NumPy Installation (version < 2.0)
>>> pip install numpy==1.26.4

--- 5. Register Environment in Jupyter Notebook
>>> conda install ipykernel -y
>>> python -m ipykernel install --user --name=reco_env
--display-name="Python3.9 (reco_env)"

--- 6. Installation of TensorFlow 2.10 and exact CUDA/cuDNN dependencies
>>> pip install tensorflow==2.10.0
>>> conda install -c conda-forge -c nvidia cudatoolkit=11.2
cudnn=8.2.0 -y
```

This set of versions ensures full compatibility between TensorFlow, CUDA/cuDNN, NumPy, and the recommenders library. In this way, the models can be trained and evaluated without errors and with GPU acceleration.

It will also be necessary to download the files 'glove.6B.300d.txt' (Pennington, Socher, and Manning, 2014) and 'MIND_small' (MSNews, 2020).

A.2 2_metrics_RADio

To correctly run the metrics repository, it is necessary to install the following packages using the commands:

```
>>> pip install bs4 community elasticsearch gensim lxml nano
python-louvain stop_words textblob textstat minimock textblob_nlp
pathlib
>>> pip install https://github.com/explosion/spacy-models/releases
/download/en_core_web_sm-3.0.0
/en_core_web_sm-3.0.0.tar.gz#egg=en_core_web_sm
```

A.3 3_3Cat

To correctly run the 3_3Cat section, it is necessary to have an environment with support for local language models via Ollama, as well as the corresponding orchestration libraries. The installation of requirements is performed using the following commands:

```
>>> pip install langchain langchain_community sparqlwrapper
```

Additionally, it is necessary to install and configure Ollama on the system and download, for mid-range computers, the model Qwen3:8b. Once Ollama is installed, the model can be obtained by selecting it in the dropdown menu and typing any message in the console, which will trigger automatic download.

After these steps, the environment is ready for executing the 3_3Cat components based on language models and workflows defined via LangChain.

Appendix B

Model Metrics Comparison

This appendix includes the full tables of diversity metrics for all models analyzed in this work. These tables complement the main text by providing detailed values for comparison with RADio.

Metric	RADio	This work
Calibration Topic	0.665	0.656
Calibration Complexity	0.523	0.497
Fragmentation	0.746	0.933
Activation	0.374	0.298
Representation	0.290	0.048
Alternative Voices	0.096	0.001

TABLE B.1: Metric comparison for the LSTUR model.

Metric	RADio	This work
Calibration Topic	0.642	0.642
Calibration Complexity	0.497	0.497
Fragmentation	0.746	0.933
Activation	0.374	0.299
Representation	0.290	0.034
Alternative Voices	0.096	0.001

TABLE B.2: Metric comparison for the NAML model.

Metric	RADio	This work
Calibration Topic	0.699	0.699
Calibration Complexity	0.494	0.494
Fragmentation	0.746	0.929
Activation	0.374	0.295
Representation	0.290	0.048
Alternative Voices	0.096	0.002

TABLE B.3: Metric comparison for the NPA model.

Metric	RADio	This work
Calibration Topic	0.647	0.647
Calibration Complexity	0.493	0.493
Fragmentation	0.746	0.933
Activation	0.374	0.291
Representation	0.290	0.037
Alternative Voices	0.096	0.001

TABLE B.4: Metric comparison for the NRMS model.

Metric	RADio	This work
Calibration Topic	0.774	0.701
Calibration Complexity	0.540	0.498
Fragmentation	0.752	0.945
Activation	0.386	0.300
Representation	0.293	0.058
Alternative Voices	0.083	0.001

TABLE B.5: Metric comparison for the Most Popular model.

Metric	RADio	This work
Calibration Topic	0.736	0.701
Calibration Complexity	0.552	0.497
Fragmentation	0.810	0.946
Activation	0.421	0.300
Representation	0.339	0.060
Alternative Voices	0.102	0.001

TABLE B.6: Metric comparison for the Random model.

Appendix C

Images of the Results of 3Cat

This appendix presents visualizations illustrating the distribution, sentiment, and diversity metrics of the 3Cat dataset.

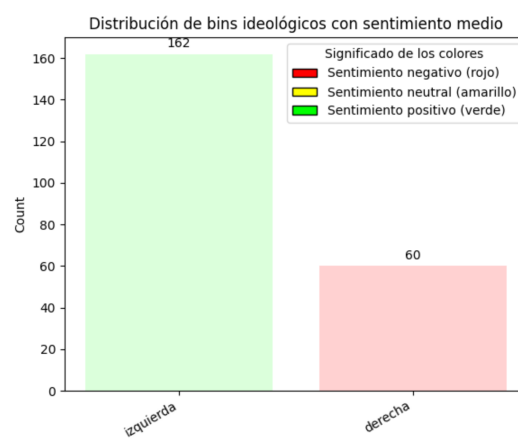


FIGURE C.1: Distribution of mentions and mean sentiment after merging ideological bins

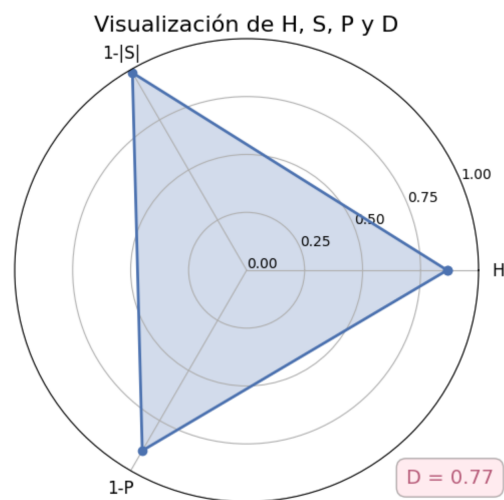


FIGURE C.2: Combined diversity metric D and components of sentiment and polarization

Bibliography

- Blanco Borrás, Rubén (2026). *Analysis of Content Diversity in News Recommendation Systems*. Software repository. URL: <https://github.com/Rbb93/TFM>.
- Halwesit (2020). *News-Recommender-MIND – Base YAML files*. Experimental configuration. URL: <https://github.com/halwesit/News-Recommender-MIND/blob/main/results/utils/lstur.yaml>.
- Holtz, David et al. (2020). *The Engagement-Diversity Connection: Evidence from a Field Experiment on Spotify*. URL: <https://arxiv.org/abs/2003.08203>.
- MSNews (2020). *MIND Small & Large Datasets – Download*. Dataset. URL: <https://msnews.github.io/>.
- Pennington, Jeffrey, Richard Socher, and Christopher D. Manning (2014). *GloVe 6B 300d*. Pre-trained word embeddings file. URL: <https://www.kaggle.com/datasets/thanakomsn/glove6b300dtxt>.
- Recommenders Team (2020). *Recommenders*. Software repository. URL: <https://github.com/recommenders-team/recommenders>.
- TechPolicy Press (2023). *What is Media Diversity and Do Recommender Systems Have It?* Popular science article. URL: <https://www.techpolicy.press/what-is-media-diversity-and-do-recommender-systems-have-it/>.
- Vrijenhoek, Sanne (2020). *Metrics calculation (RADio)*. Metrics calculation notebook. URL: https://github.com/svrijenhoek/RADio/blob/main/metrics_calculation.ipynb.
- Vrijenhoek, Sanne et al. (2024). *RADio* – An Introduction to Measuring Normative Diversity in News Recommendations*. Popular science article. URL: <https://dl.acm.org/doi/10.1145/3636465>.