

# Proiect 2 Inteligenta Artificiala - ML Applied

Dumitrascu Filip-Teodor 333CA

## Continut

- 1. [Introducere](#)
- 2. [EDA](#)
- 3. [Preprocesare Date](#)
- 4. [Algoritmi](#)
- 5. [Comparatie](#)

## Introducere

În activitatea curentă a unui inginer sau cercetător din domeniul inteligenței artificiale, în special al învățării automate, apar frecvent următoarele trei componente esențiale:

- Vizualizarea și analiza exploratorie a datelor (EDA – Exploratory Data Analysis)
- Identificarea și extragerea caracteristicilor relevante ale datelor, utile pentru obiectivul propus (precum clasificare, regresie sau detecția anomaliilor)
- Testarea și compararea mai multor modele, în vederea alegerii celei mai eficiente soluții pentru problema abordată.

Astfel, în acest proiect sunt analizate două seturi de date (Air pollution si News popularity) care sunt parcurse prin etapele menționate anterior, având ca scop aprofundarea înțelegerii procesului de învățare automată.

## EDA

### Air pollution

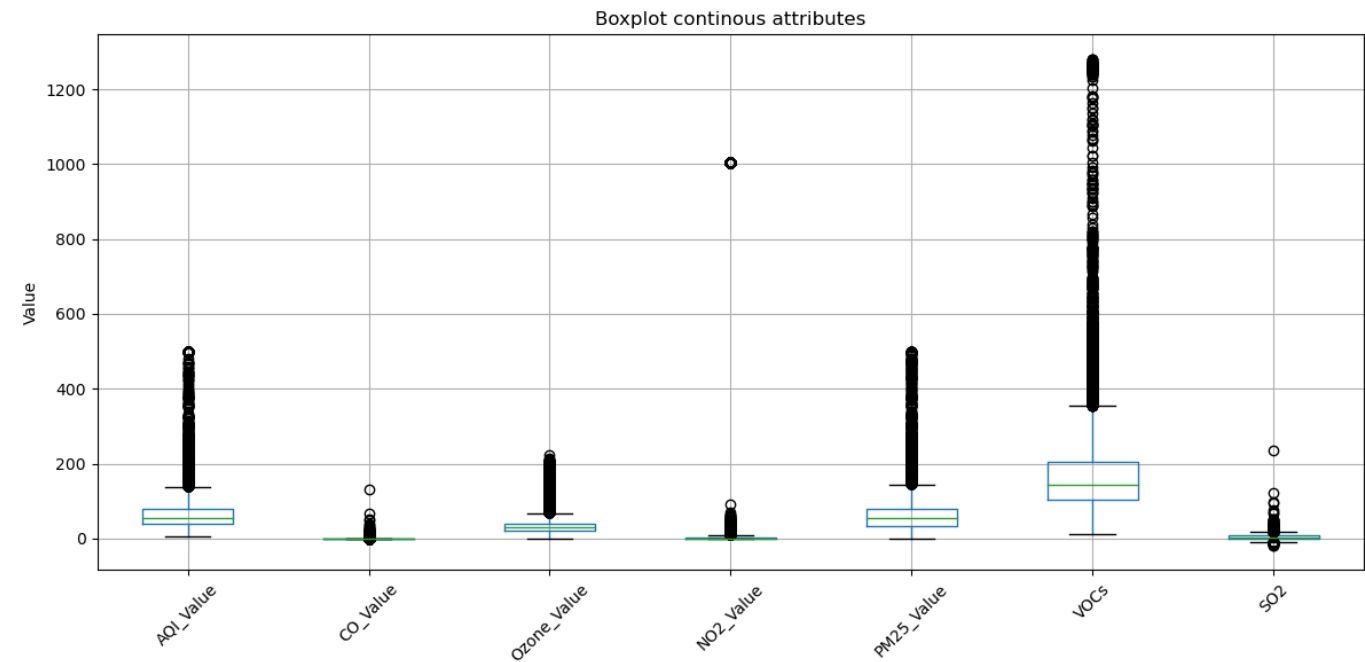
#### Tipul atributelor

Având în vedere tipul valorilor pe care le conțin și intervalele în care acestea variază, attributele din setul de date pot fi clasificate în următoarele categorii: (se ignora **Unnamed: 0**)

- attribute continue: **AQI\_Value, CO\_Value, Ozone\_Value, NO2\_Value, PM25\_Value, VOCs, SO2** (numerice, valori unice > 20, range de valori > 50)
- attribute discrete: **Country, City** (restul care nu sunt continue si ordinale)
- attribute ordinale: **AQI\_Category, CO\_Category, Ozone\_Category, NO2\_Category, PM25\_Category, Emissions** (contin obiecte ordonate: Good-Moderate, Level0-Level5)

#### Atribute numerice continue analiza

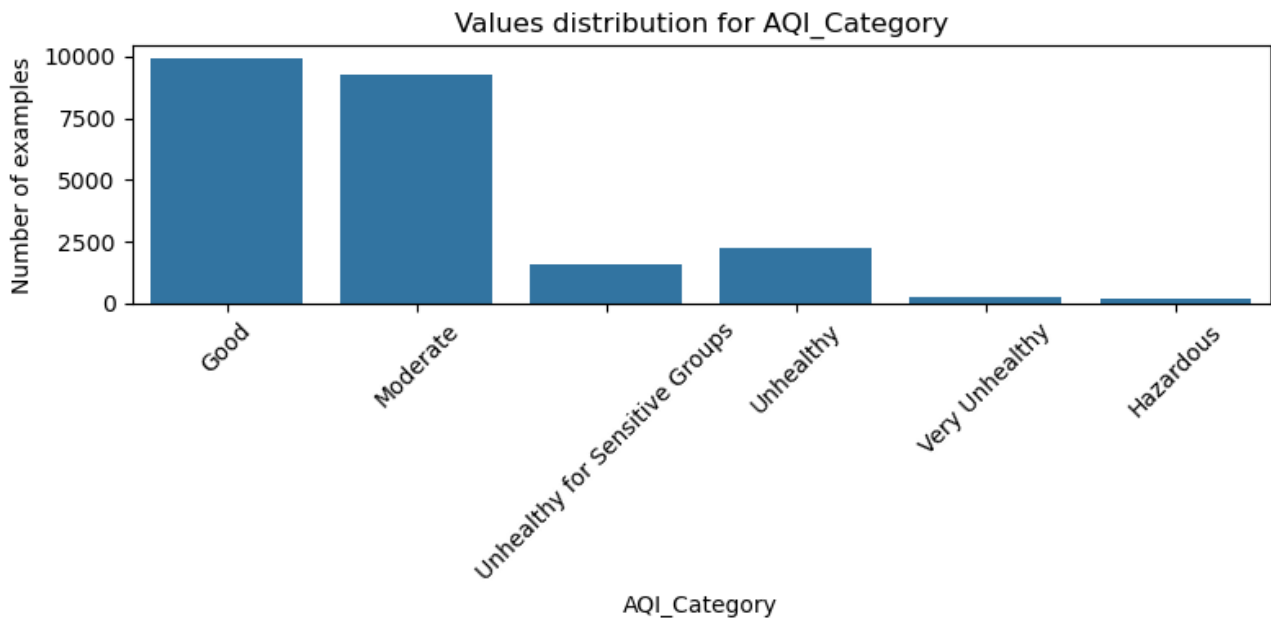
Attribute	Count	Mean	Std Dev	Min	25%	50%	75%	Max
AQI_Value	23463.0	72.01	56.06	6.00	39.00	55.00	79.00	500.00
CO_Value	23463.0	1.37	1.83	0.00	1.00	1.00	1.00	133.00
Ozone_Value	21117.0	35.24	28.15	0.00	21.00	31.00	40.00	222.00
NO2_Value	23463.0	43.08	196.08	0.00	0.00	1.00	4.00	1003.06
PM25_Value	23463.0	68.52	54.80	0.00	35.00	54.00	79.00	500.00
VOCs	23463.0	185.05	140.49	12.42	103.27	142.97	204.23	1280.99
SO2	23463.0	4.45	5.95	-18.53	0.74	4.29	7.92	234.69



Boxplot-ul evidențiază faptul că majoritatea valorilor pentru atributele continue din setul de date Air Pollution (precum **SO2**, **CO\_Value**, **NO2\_Value**) sunt concentrate în intervale relativ mici, în timp ce există numeroși outlieri care depășesc cu mult aceste valori. Atributele **VOCs**, **PM25\_Value** și **AQI\_Value** se remarcă prin plaje mari de valori și un număr ridicat de valori extreme, indicând o distribuție dezechilibrată și posibilă nevoie de tratament al outlierilor.

Atribute ordinale analiza

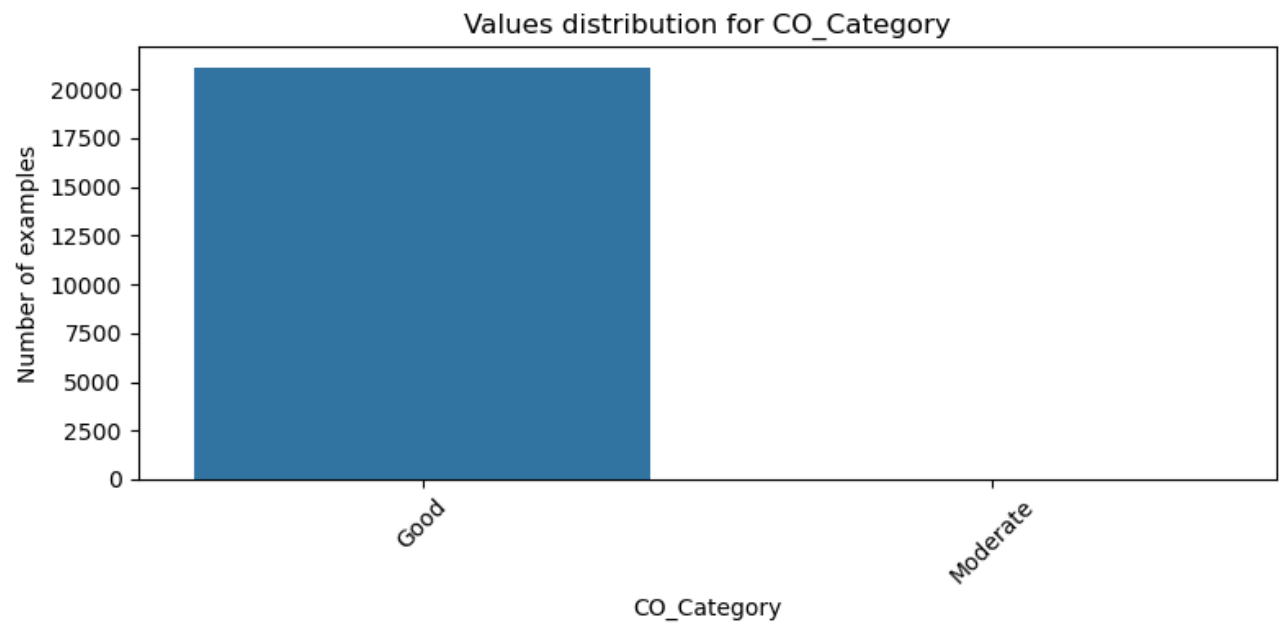
- AQI\_Category — Number of examples non-null: 23463
- AQI\_Category — Number of unique values: 6



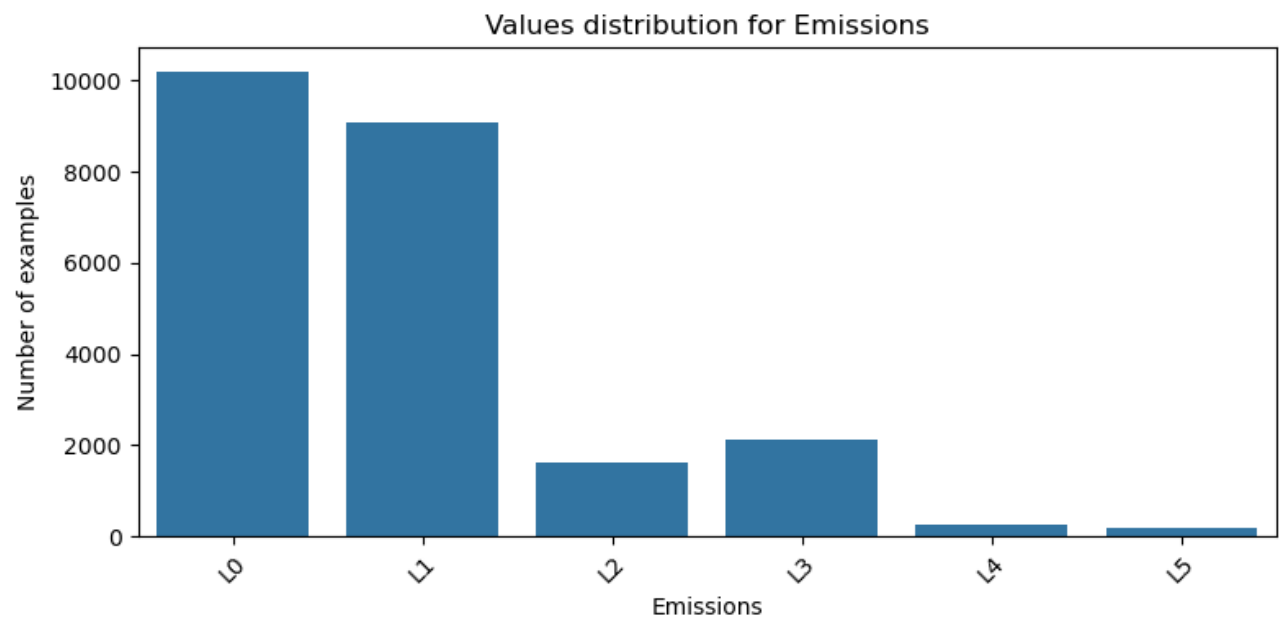
AQI\_Category prezintă o distribuție puternic dezechilibrată, cu majoritatea articolelor etichetate drept Good sau Moderate. Clasele Very Unhealthy și Hazardous sunt semnificativ sub-reprezentate.

- CO\_Category — Number of examples non-null: 21117

- CO\_Category — Number of unique values: 2

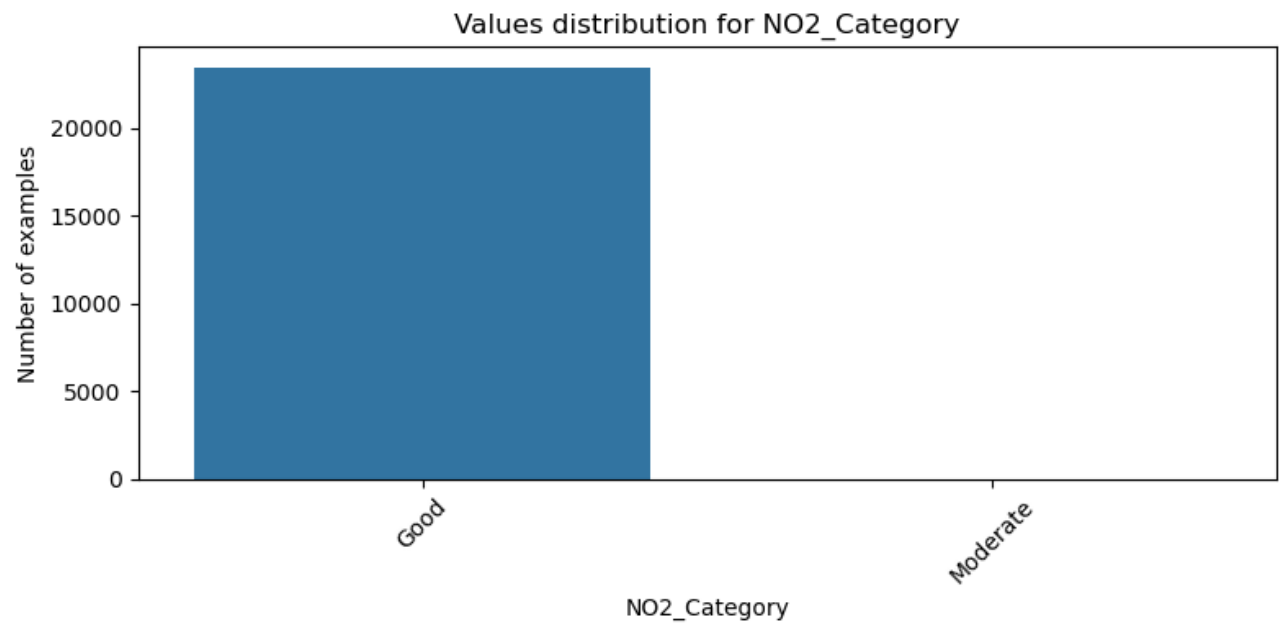


- Ozone\_Category — Number of examples non-null: 23463
- Ozone\_Category — Number of unique values: 5



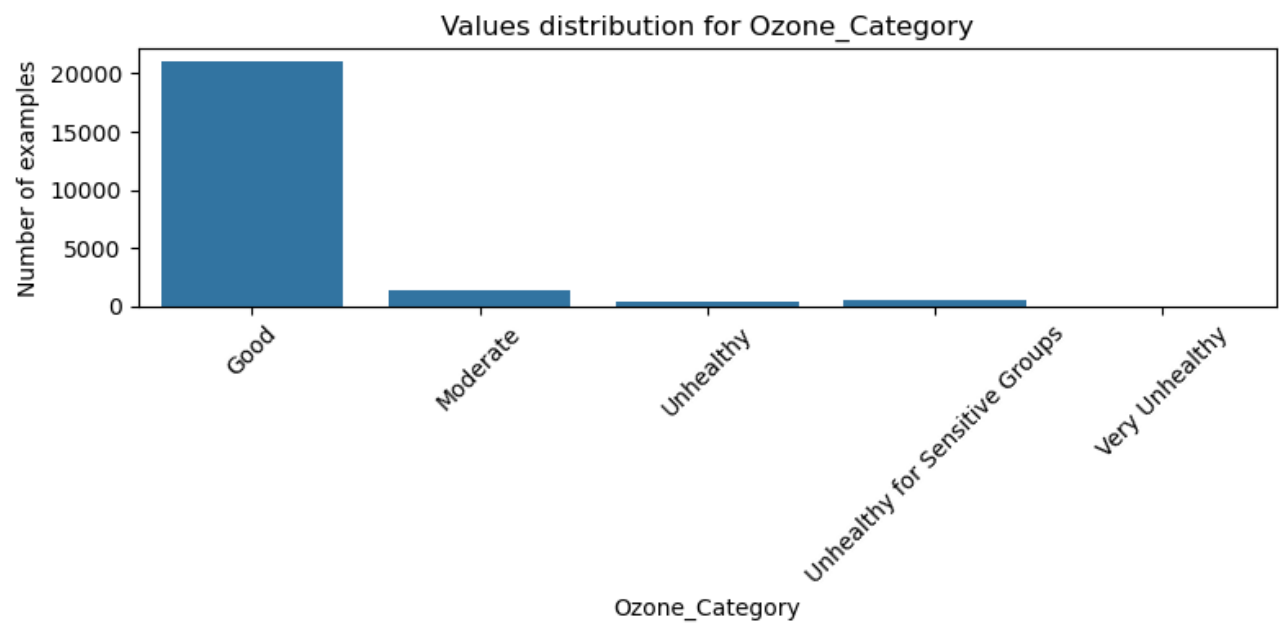
- NO2\_Category — Number of examples non-null: 23463

- NO2\_Category — Number of unique values: 2



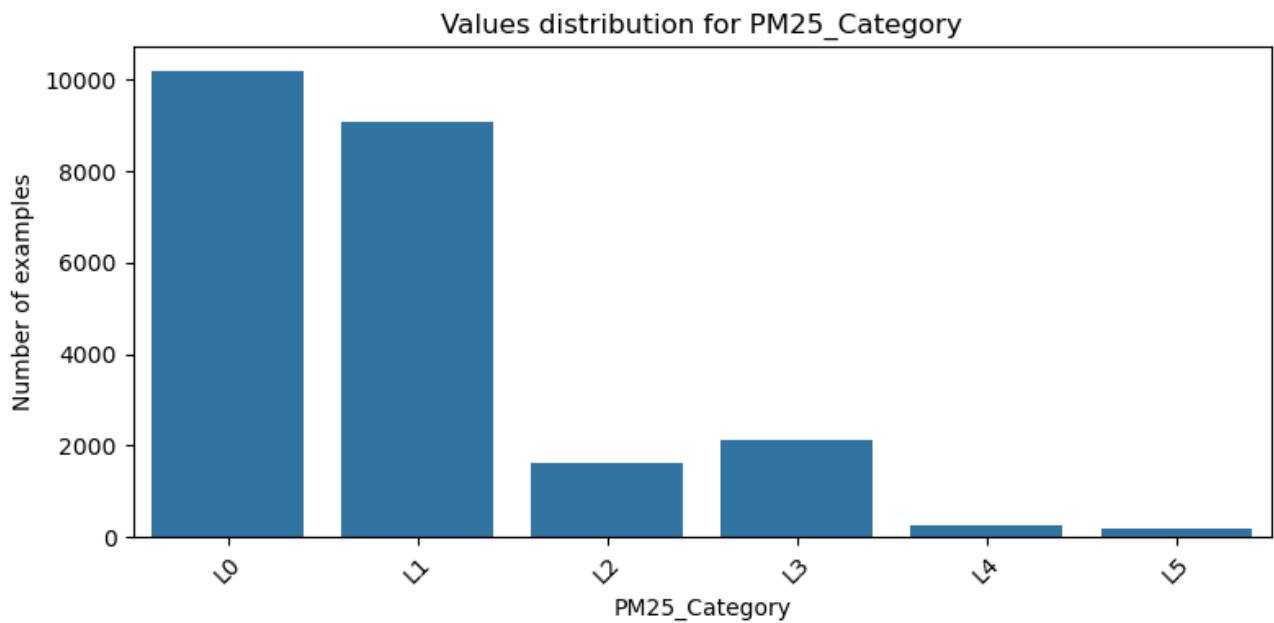
CO\_Category, NO2\_Category și Ozone\_Category sunt dominate de categoria Good, celelalte clase fiind aproape absente. Acest dezechilibru poate afecta performanța algoritmilor de clasificare.

- PM25\_Category — Number of examples non-null: 23463
- PM25\_Category — Number of unique values: 6



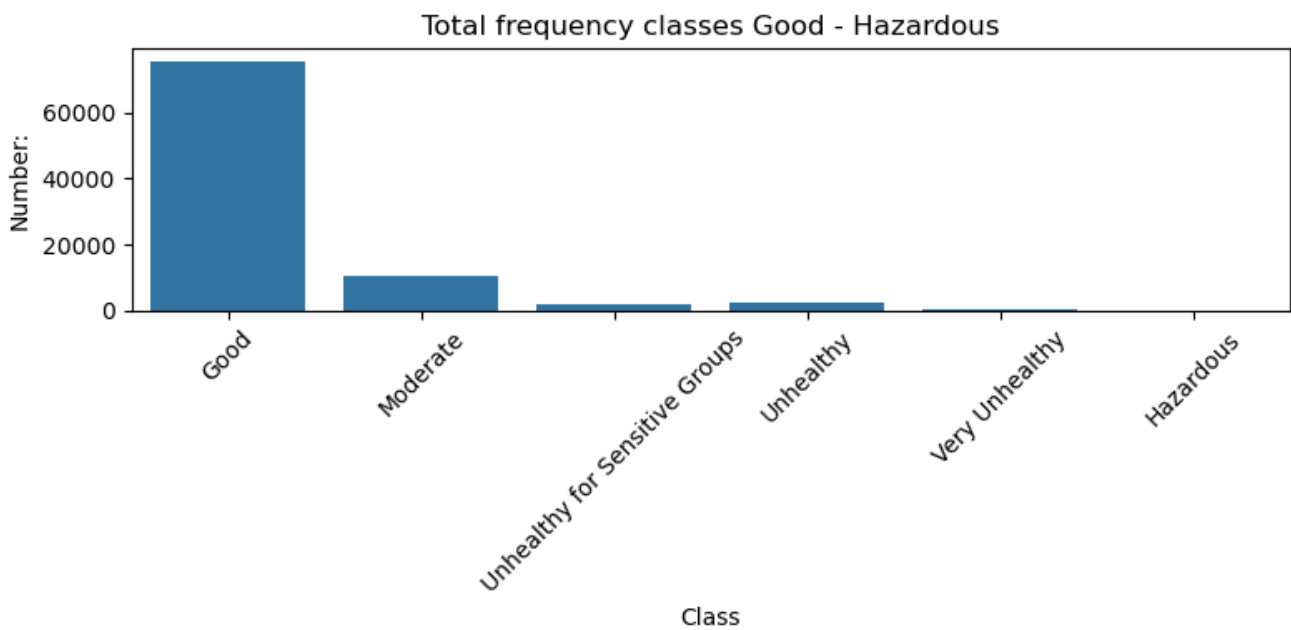
- Emissions — Number of examples non-null: 23463

- Emissions — Number of unique values: 6

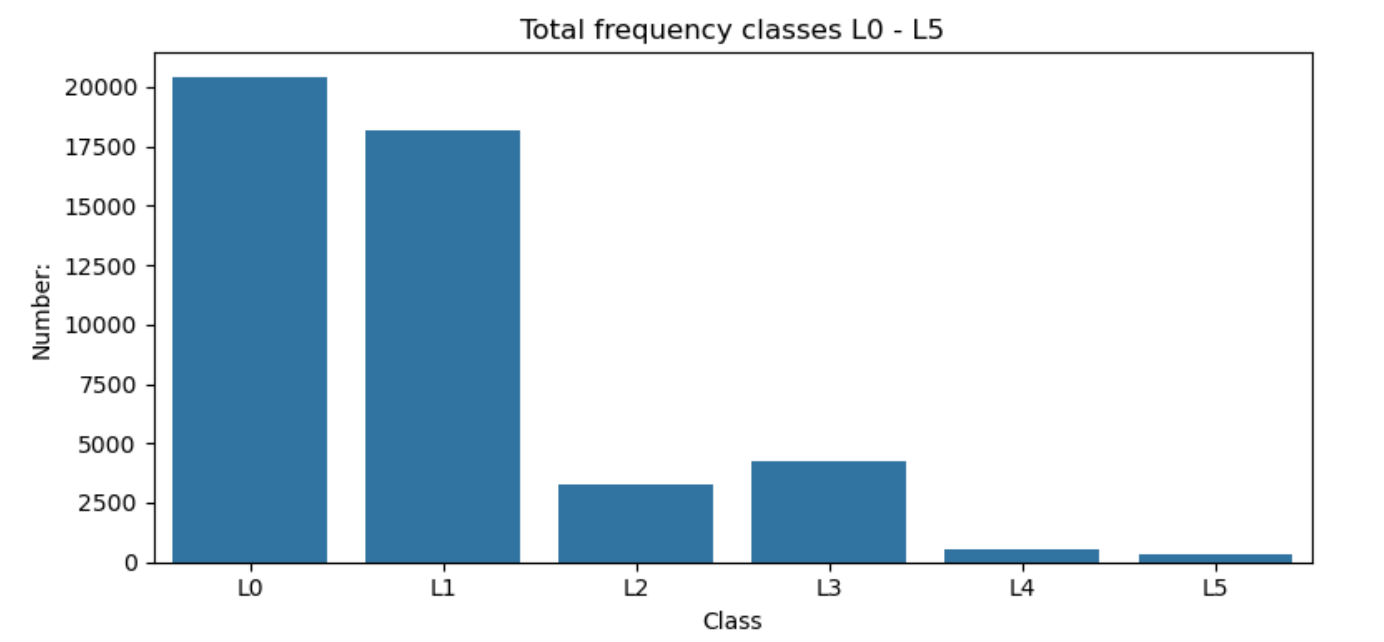


Emissions și PM25\_Category prezintă o distribuție mai variată, dar tot cu o concentrație ridicată în clasele L0 și L1, indicând un dezechilibru moderat spre sever în cadrul claselor ordinale de tip L.

Echilibrul Claselor



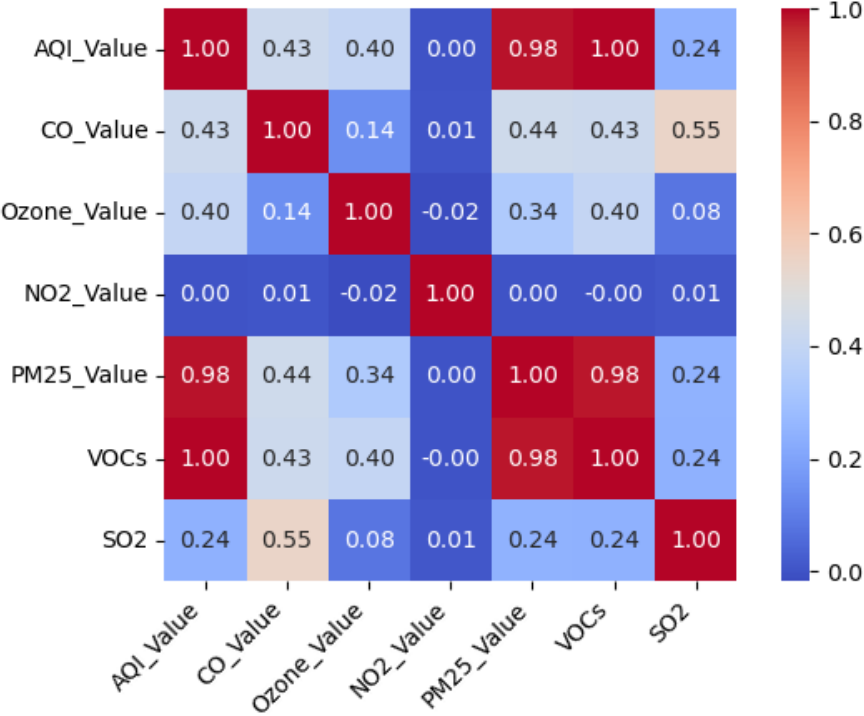
Graficul evidențiază un dezechilibru semnificativ între clasele categoriei AQI. Cea mai mare parte a exemplelor sunt etichetate ca „Good” (peste 70.000), urmată la mare distanță de „Moderate” și apoi de celelalte clase („Unhealthy for Sensitive Groups”, „Unhealthy”, „Very Unhealthy” și „Hazardous”), care sunt subreprezentate. Acest dezechilibru poate afecta performanța algoritmilor de clasificare, care vor înclina spre clasele dominante.



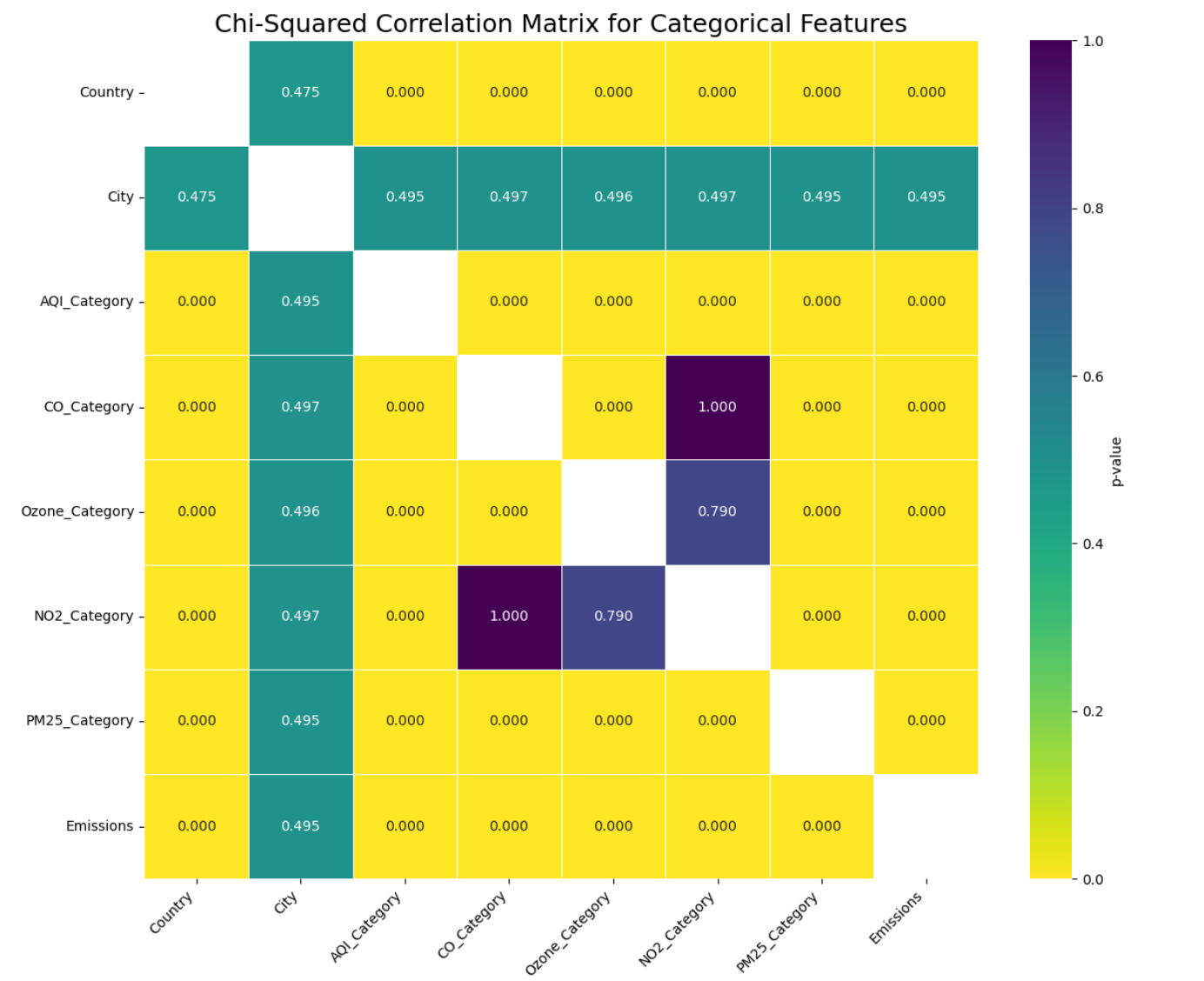
Atributele ordinale legate de nivelul emisiilor (L0–L5) prezintă de asemenea un dezechilibru: clasele L0 și L1 sunt cele mai frecvente, în timp ce L4 și L5 apar rar. Acest tip de distribuție poate influența negativ capacitatea modelului de a învăța să recunoască corect cazurile mai puțin frecvente (niveluri mari de emisii).

Corelatia intre atribute

Pearson Correlation Matrix for Continuous Attributes Air



Această matrice evidențiază corelațiile lineare între atributele numerice continue din setul de date Air. Se observă o corelație aproape perfectă ( $\approx 1.00$ ) între AQI\_Value, PM25\_Value și VOCs, indicând redundanță informațională. În schimb, NO2\_Value are corelații apropiate de 0 cu toate celelalte atribute, sugerând că este independent din punct de vedere liniar față de restul.



Matricea arată valorile p obținute în urma testului Chi-pătrat aplicat perechilor de variabile categorice. Valori apropiate de 0 indică o asociere statistic semnificativă (corelație), iar valorile mari (ex. CO\_Category și NO2\_Category) sugerează lipsa de asociere. De exemplu, AQI\_Category are legături semnificative cu toate celelalte atribute, ceea ce reflectă importanța sa centrală în evaluarea calității aerului.

News pollution

Tipul atributelor

Având în vedere tipul valorilor pe care le conțin și intervalele în care acestea variază, atributele din setul de date pot fi clasificate în următoarele categorii: (se ignora url)

- attribute continue: days\_since\_published, content\_word\_count, unique\_word\_ratio, non\_stop\_word\_ratio, unique\_non\_stop\_ratio, external\_links, internal\_links, image\_count, video\_count, keyword\_worst\_min\_shares, keyword\_worst\_max\_shares, keyword\_worst\_avg\_shares, keyword\_best\_min\_shares, keyword\_best\_max\_shares, keyword\_best\_avg\_shares, keyword\_avg\_min\_shares, keyword\_avg\_max\_shares, keyword\_avg\_avg\_shares, ref\_min\_shares, ref\_max\_shares, ref\_avg\_shares, engagement\_ratio, content\_density (numerice, valori unice > 20, range de valori > 50)
- attribute discrete: title\_word\_count, avg\_word\_length, keyword\_count, topic\_0\_relevance, topic\_1\_relevance, topic\_2\_relevance, topic\_3\_relevance, topic\_4\_relevance, content\_subjectivity, content\_sentiment, positive\_word\_rate, negative\_word\_rate, non\_neutral\_positive\_rate, non\_neutral\_negative\_rate, avg\_positive\_sentiment, min\_positive\_sentiment, max\_positive\_sentiment, avg\_negative\_sentiment, min\_negative\_sentiment, max\_negative\_sentiment,

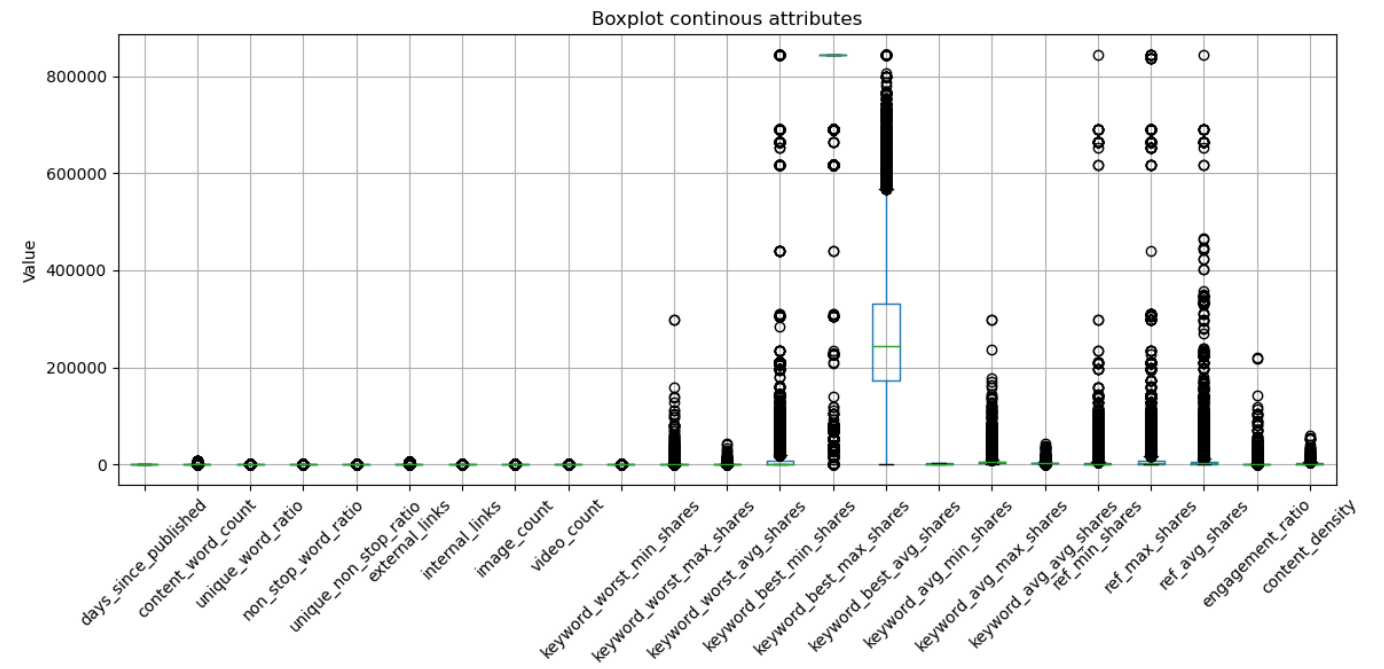
title\_subjectivity, title\_sentiment, title\_subjectivity\_magnitude, title\_sentiment\_magnitude  
(numerie, restul)

- attribute ordinale: popularity\_category (contin obiecte ordonate: Slightly Popular, Moderatately Popular)
- attribute categorice: channel\_lifestyle, channel\_entertainment, channel\_business, channel\_social\_media, channel\_tech, channel\_world, day\_monday, day\_tuesday, day\_wednesday, day\_thursday, day\_friday, day\_saturday, day\_sunday, is\_weekend, publication\_period (restul de obiecte)

Attribute numerice continue analiza

Attribute	Count	Mean	Std Dev	Min	25%	50%	75%	Max
days_since_published	39644.0	354.53	214.16	8.00	164.00	339.00	542.00	731.00
content_word_count	39644.0	546.51	471.11	0.00	246.00	409.00	716.00	8474.00
unique_word_ratio	39644.0	0.55	3.52	0.00	0.47	0.54	0.61	701.00
non_stop_word_ratio	39644.0	1.00	5.23	0.00	1.00	1.00	1.00	1042.00
unique_non_stop_ratio	39644.0	0.69	3.26	0.00	0.63	0.69	0.75	650.00
external_links	39644.0	192.25	905.42	0.00	4.00	8.00	15.00	6078.62
internal_links	39644.0	3.29	3.86	0.00	1.00	3.00	4.00	116.00
image_count	39644.0	4.54	8.31	0.00	1.00	1.00	4.00	128.00
video_count	39644.0	1.25	4.11	0.00	0.00	0.00	1.00	91.00
keyword_worst_min_shares	39644.0	26.11	69.63	-1.00	-1.00	-1.00	4.00	377.00
keyword_worst_max_shares	39644.0	1153.95	3857.99	0.00	445.00	660.00	1000.00	298400.00
keyword_worst_avg_shares	39644.0	312.37	620.78	-1.00	141.75	235.50	357.00	42827.86
keyword_best_min_shares	39644.0	13612.35	57986.03	0.00	0.00	1400.00	7900.00	843300.00
keyword_best_max_shares	39644.0	752324.07	214502.13	0.00	843300.00	843300.00	843300.00	843300.00
keyword_best_avg_shares	39644.0	259281.94	135102.25	0.00	172846.88	244572.22	330980.00	843300.00
keyword_avg_min_shares	39644.0	1117.15	1137.46	-1.00	0.00	1023.64	2056.78	3613.04
keyword_avg_max_shares	39644.0	5657.21	6098.87	0.00	3562.10	4355.69	6019.95	298400.00
keyword_avg_avg_shares	39644.0	3135.86	1318.15	0.00	2382.45	2870.07	3600.23	43567.66
ref_min_shares	39644.0	3998.76	19738.67	0.00	639.00	1200.00	2600.00	843300.00
ref_max_shares	39644.0	10329.21	41027.58	0.00	1100.00	2800.00	8000.00	843300.00
ref_avg_shares	39644.0	6401.70	24211.33	0.00	981.19	2200.00	5200.00	843300.00
engagement_ratio	39644.0	1054.07	3496.61	0.04	220.00	465.50	900.00	221200.00
content_density	35680.0	1986.56	2209.10	32.76	746.07	1314.55	2506.71	58857.97

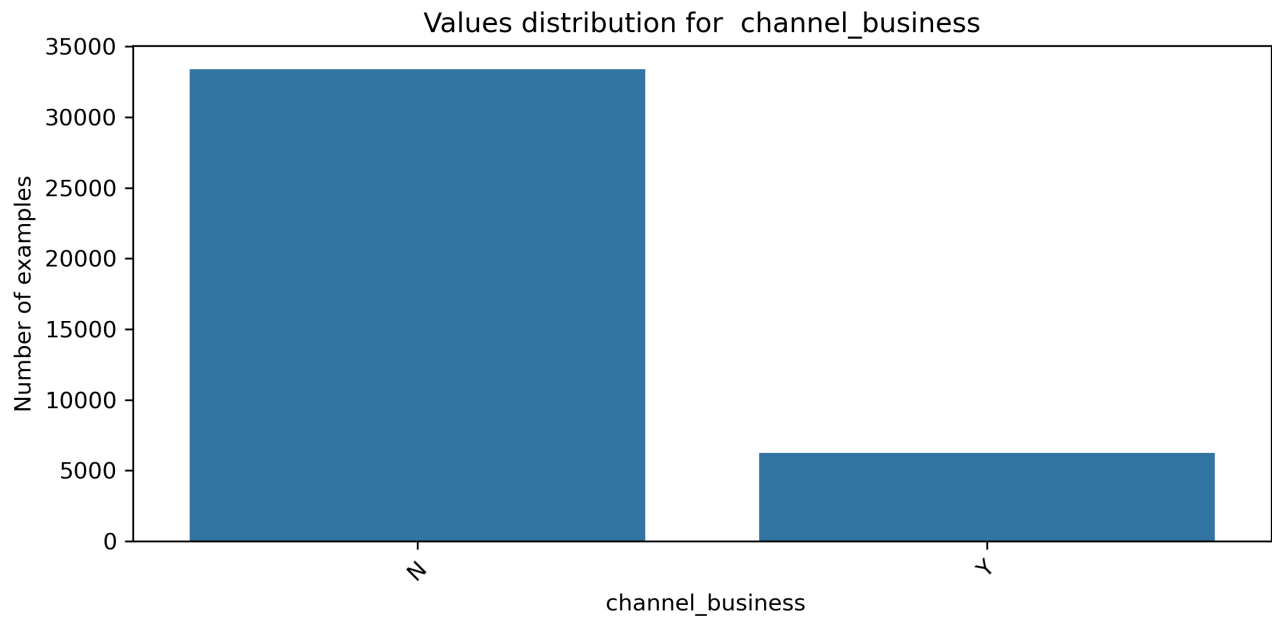




Boxplot-ul pentru setul News Popularity arată o variabilitate mare între atribute, unele dintre ele având valori extrem de ridicate (peste 800.000), precum **keyword\_best\_max\_shares** sau **ref\_max\_shares**. Multe dintre atribute sunt puternic dezechilibrate, cu mediane joase și prezența unor outlieri extremi. Acest lucru sugerează că este necesară o etapă de tratare a valorilor extreme și o eventuală standardizare a datelor pentru a evita influențarea algoritmilor de ML.

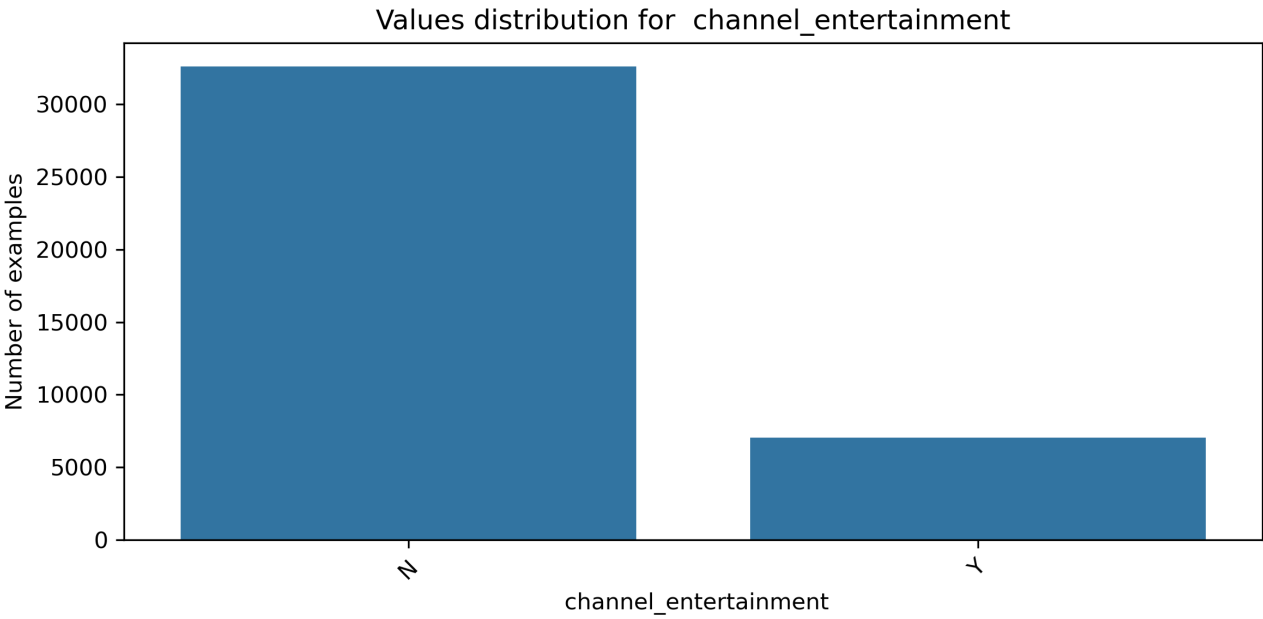
Atribute ordinale/catgorice analiza

- channel\_business — Number of non-null examples: 39644
- channel\_business — Number of unique values: 2

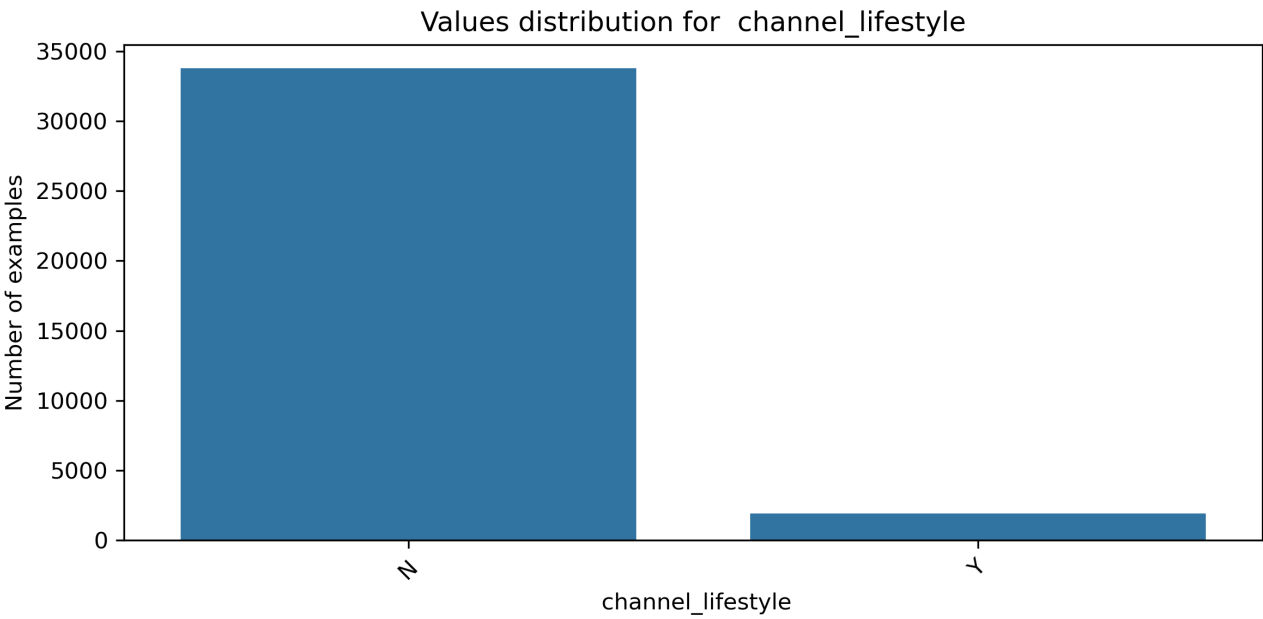


- channel\_entertainment — Number of non-null examples: 39644

- channel\_entertainment — Number of unique values: 2

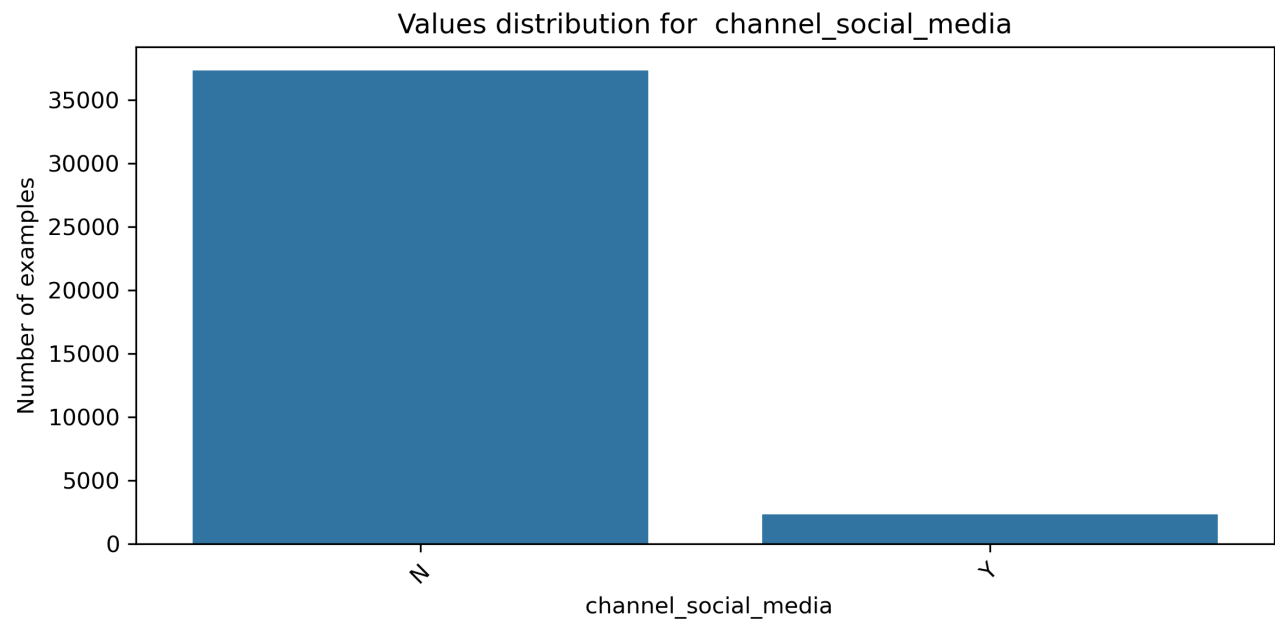


- channel\_lifestyle — Number of non-null examples: 35680
- channel\_lifestyle — Number of unique values: 2

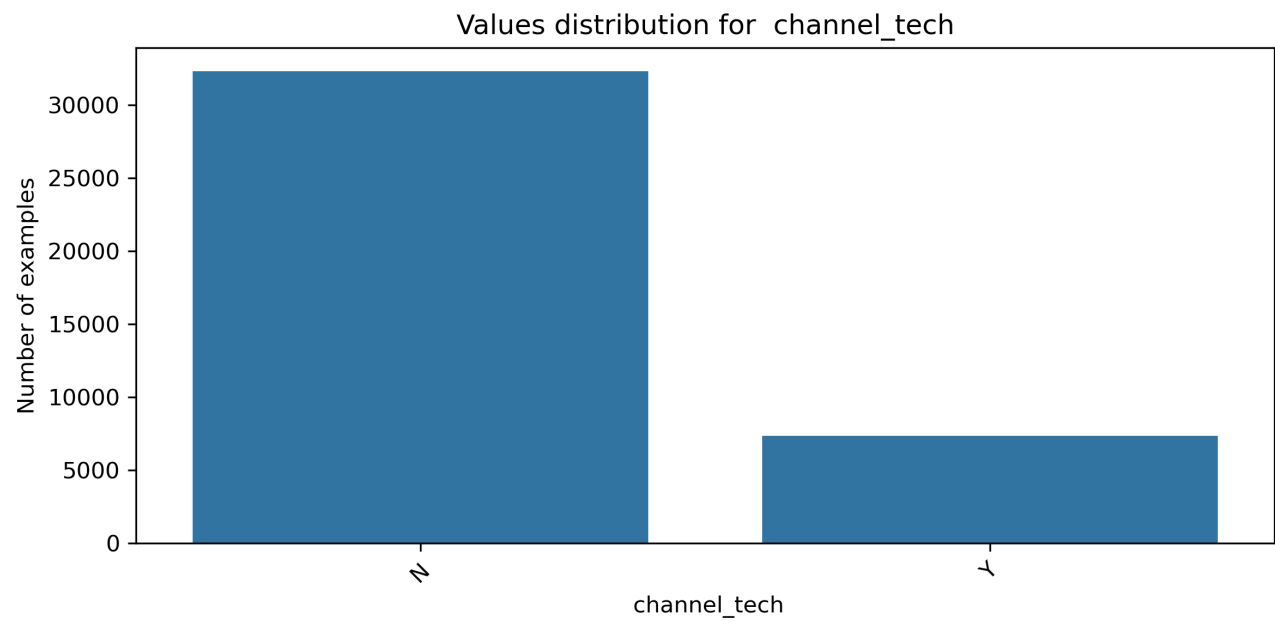


- channel\_social\_media — Number of non-null examples: 39644

- channel\_social\_media — Number of unique values: 2

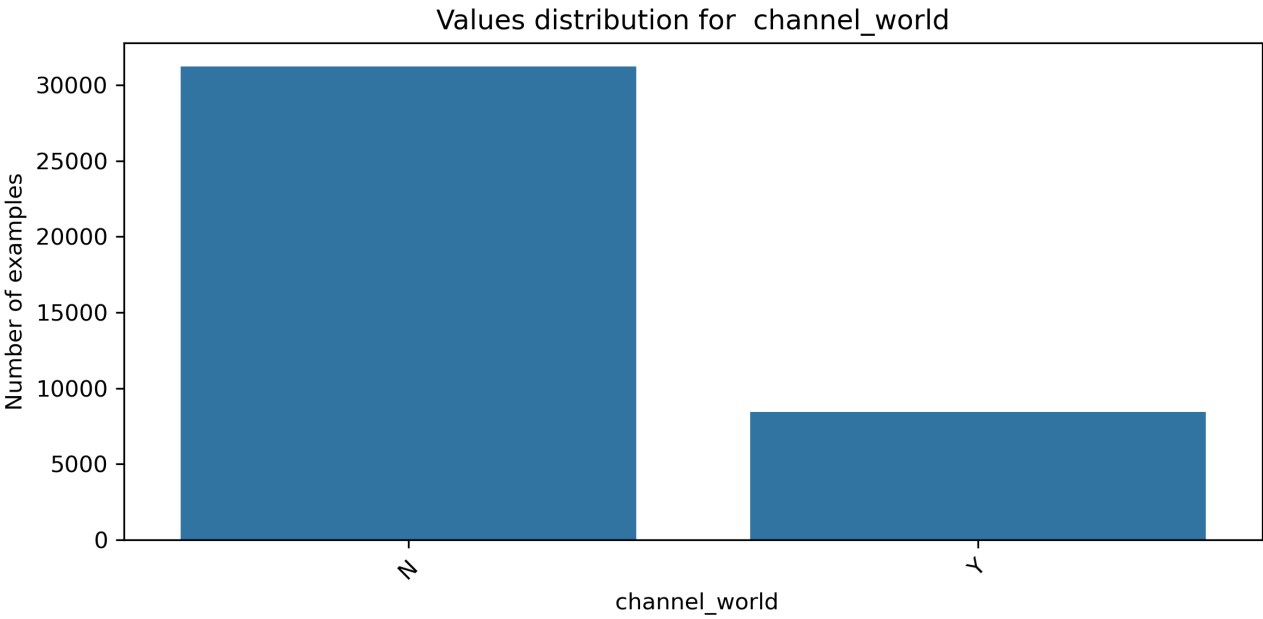


- channel\_tech — Number of non-null examples: 39644
- channel\_tech — Number of unique values: 2



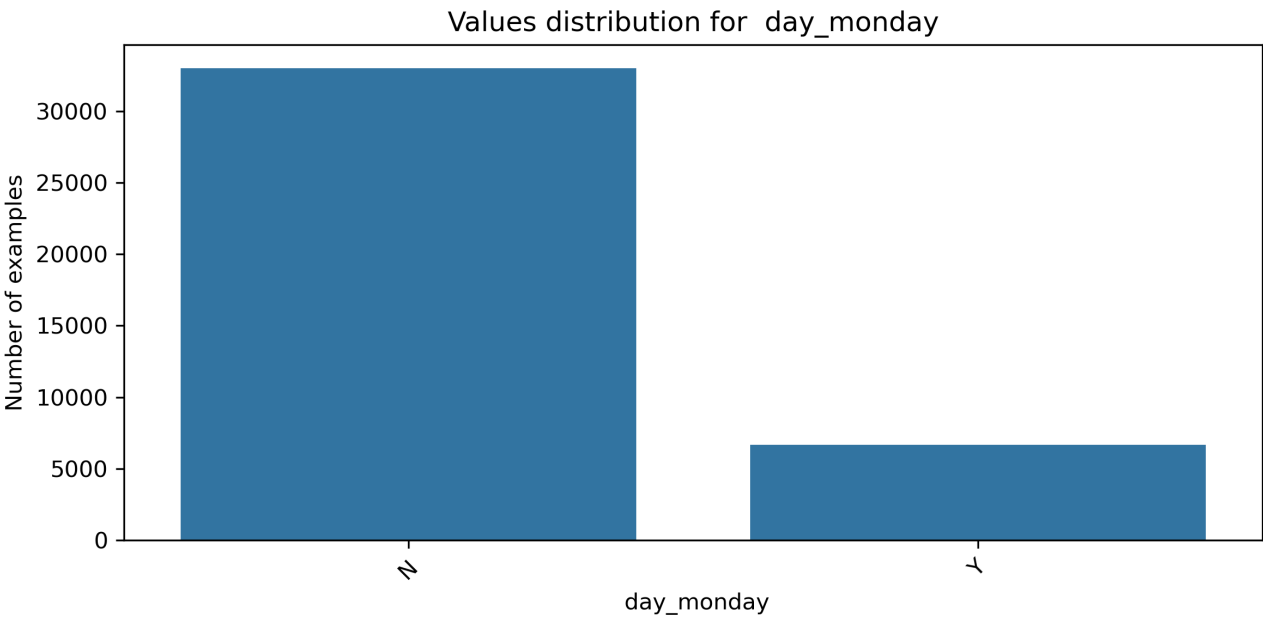
- channel\_world — Number of non-null examples: 39644

- channel\_world — Number of unique values: 2



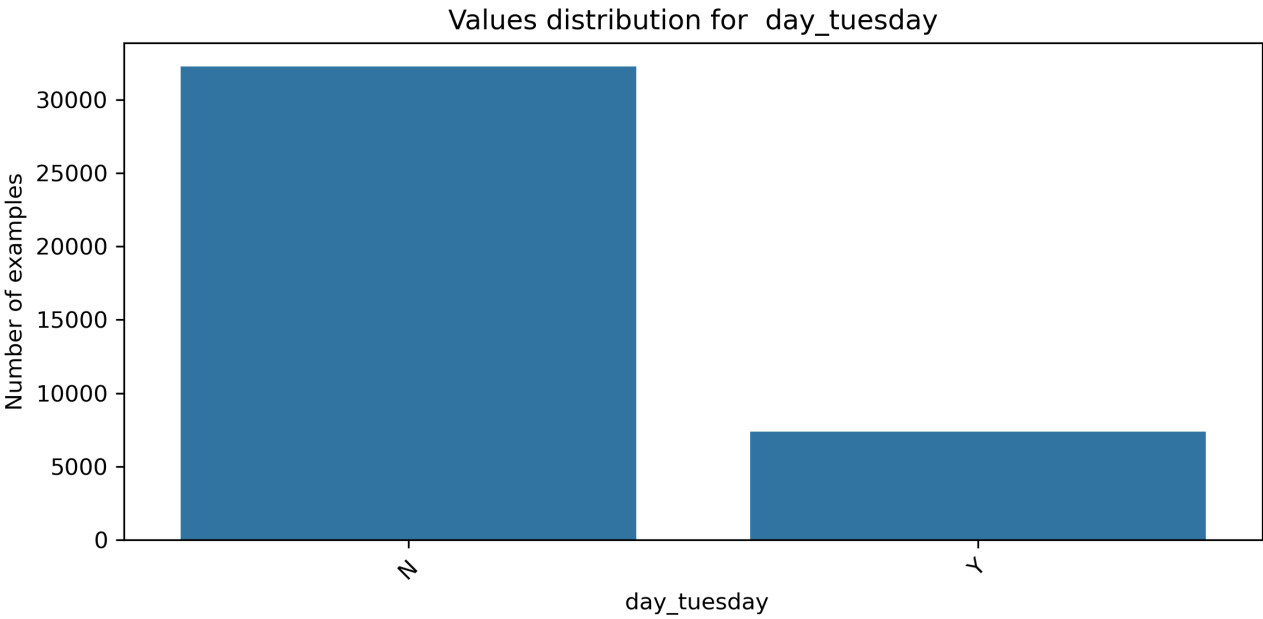
Atributele binare precum channel\_business, channel\_entertainment, channel\_lifestyle, channel\_social\_media, channel\_tech, channel\_world arată dacă un articol aparține unui anumit canal tematic. Distribuția este dezechilibrată în toate cazurile: predomină valorile 'N' (nu aparține canalului), iar doar o fracțiune relativ mică sunt articole marcate cu 'Y' (da). Acest lucru poate influența performanța clasificatorilor, deoarece modelul învață mai greu din clasele rare.

- day\_monday — Number of non-null examples: 39644
- day\_monday — Number of unique values: 2

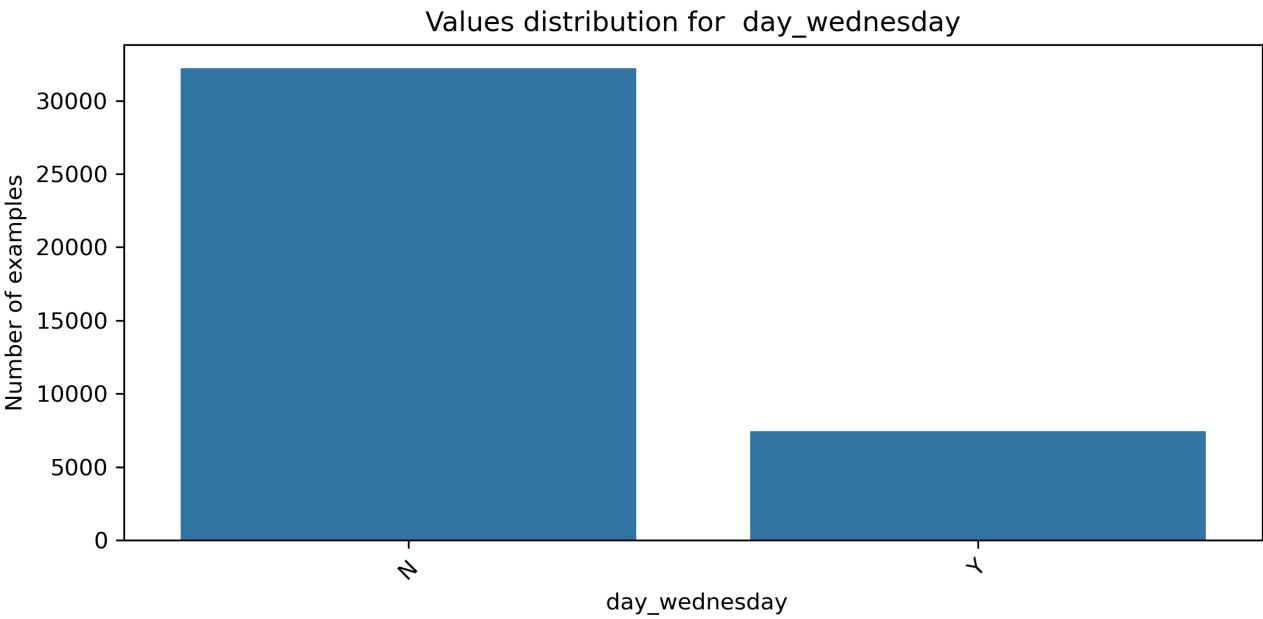


- day\_tuesday — Number of non-null examples: 39644

- day\_tuesday — Number of unique values: 2

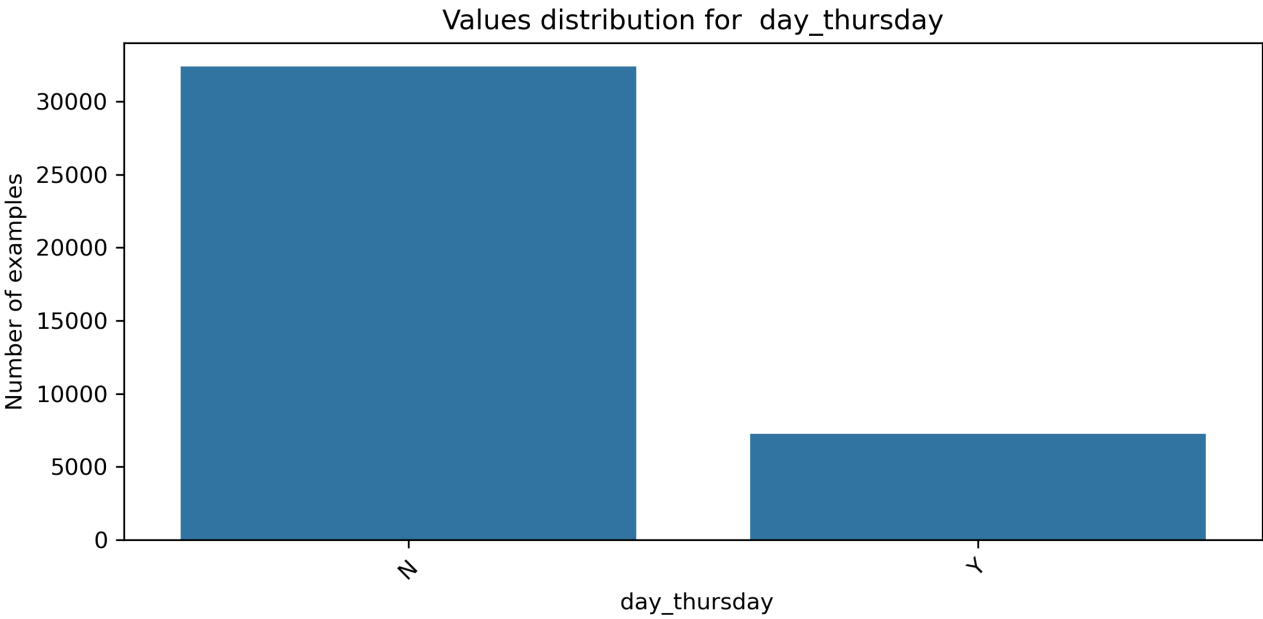


- day\_wednesday — Number of non-null examples: 39644
- day\_wednesday — Number of unique values: 2

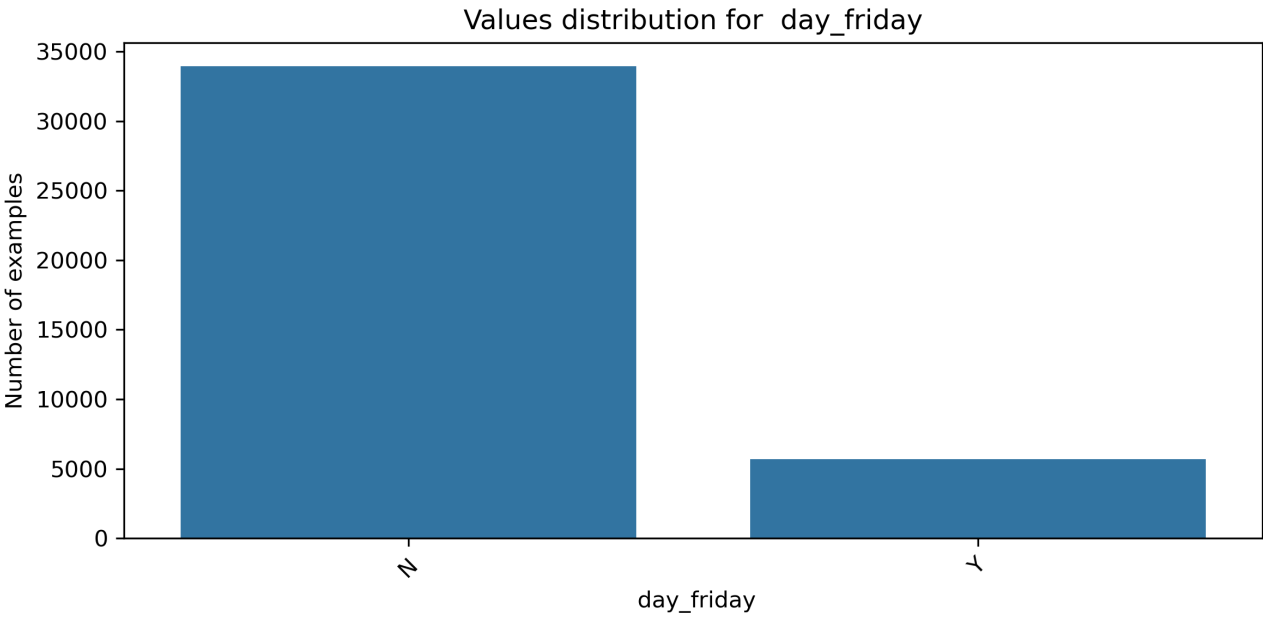


- day\_thursday — Number of non-null examples: 39644

- day\_thursday — Number of unique values: 2

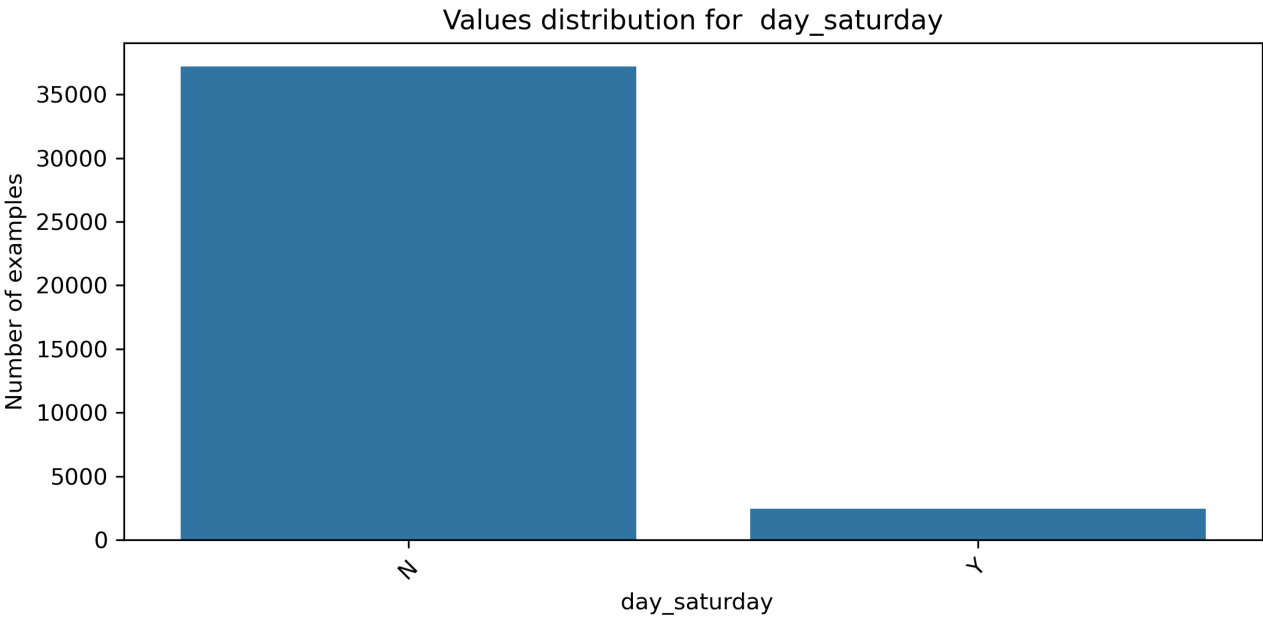


- day\_friday — Number of non-null examples: 39644
- day\_friday — Number of unique values: 2

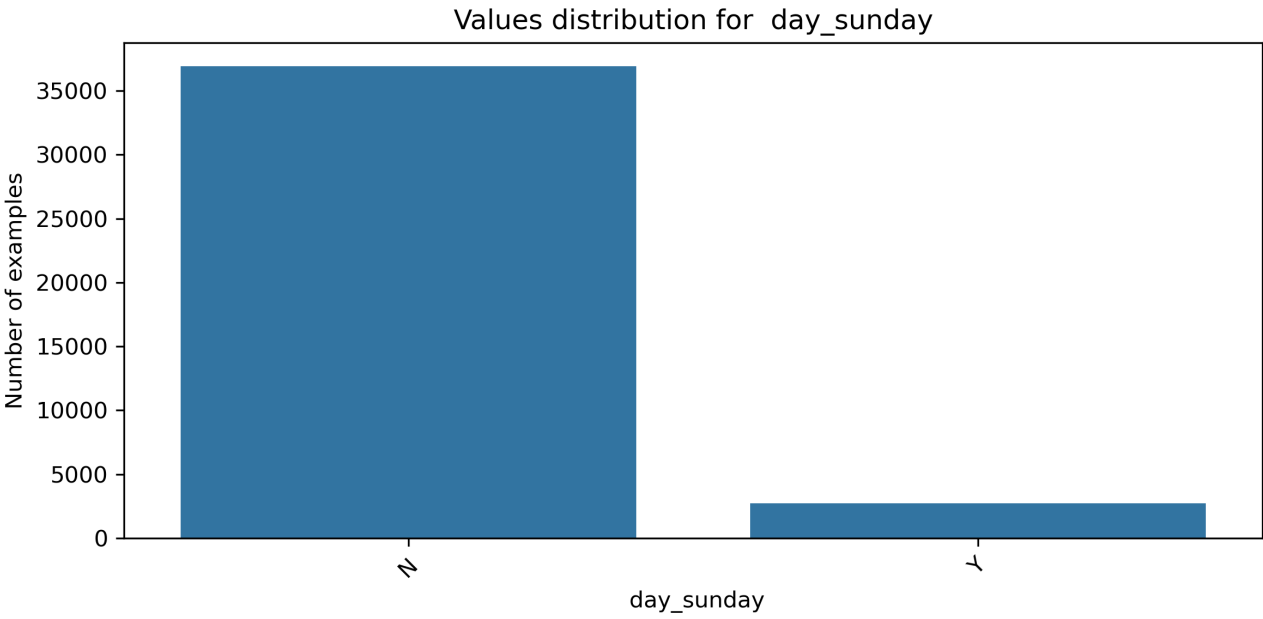


- day\_saturday — Number of non-null examples: 39644

- day\_saturday — Number of unique values: 2

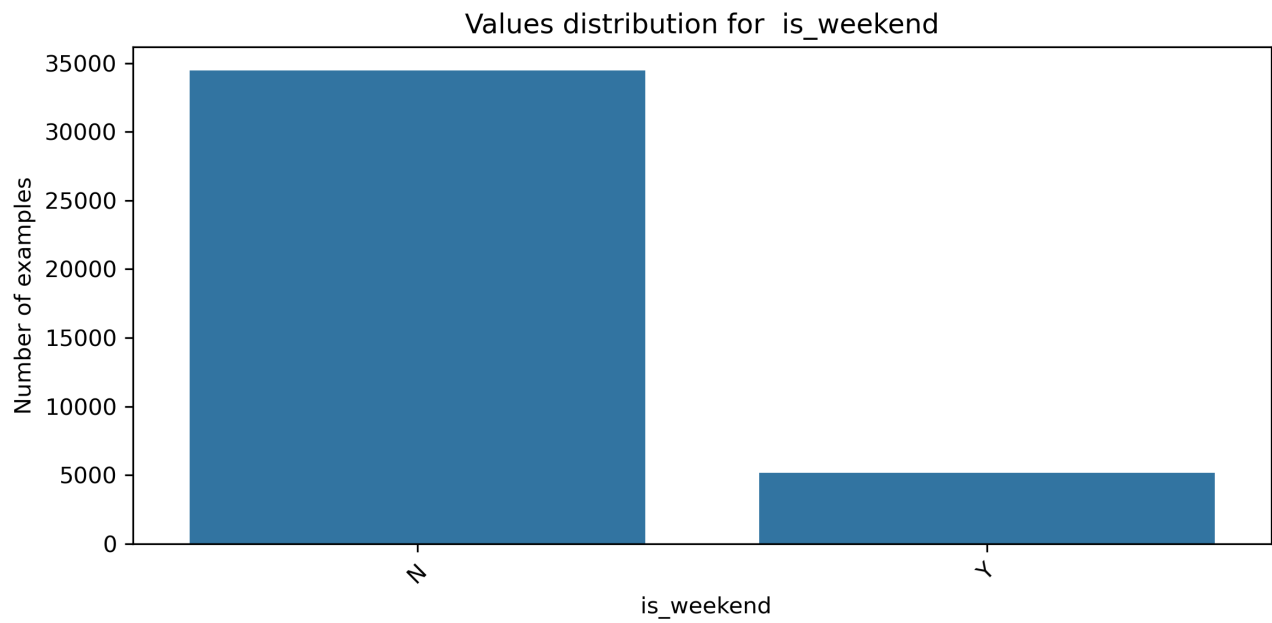


- day\_sunday — Number of non-null examples: 39644
- day\_sunday — Number of unique values: 2



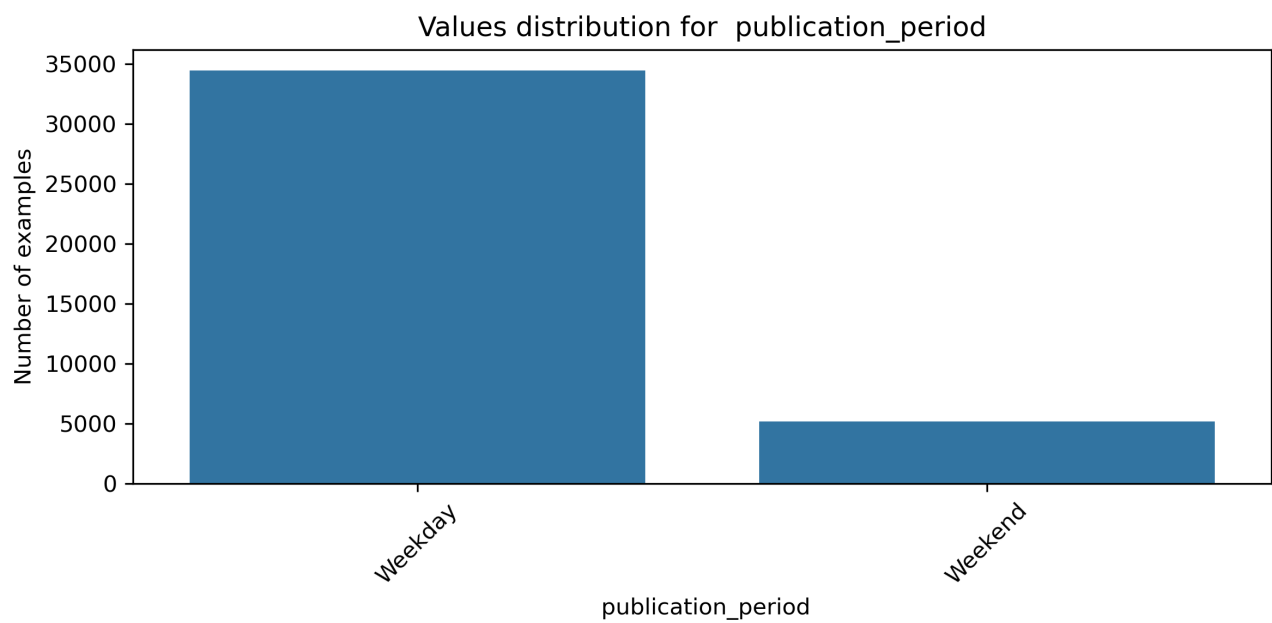
- is\_weekend — Number of non-null examples: 39644

- is\_weekend — Number of unique values: 2



Atributele binare de tipul day\_monday, day\_friday, day\_saturday etc. indică dacă articolul a fost publicat într-o anumită zi. Distribuția este tot dezechilibrată — majoritatea valorilor sunt 'N', în timp ce 'Y' apare doar când articolul a fost publicat în acea zi. Se observă o ușoară variație între zilele săptămânii și weekend.

- publication\_period — Number of non-null examples: 39644
- publication\_period — Number of unique values: 2

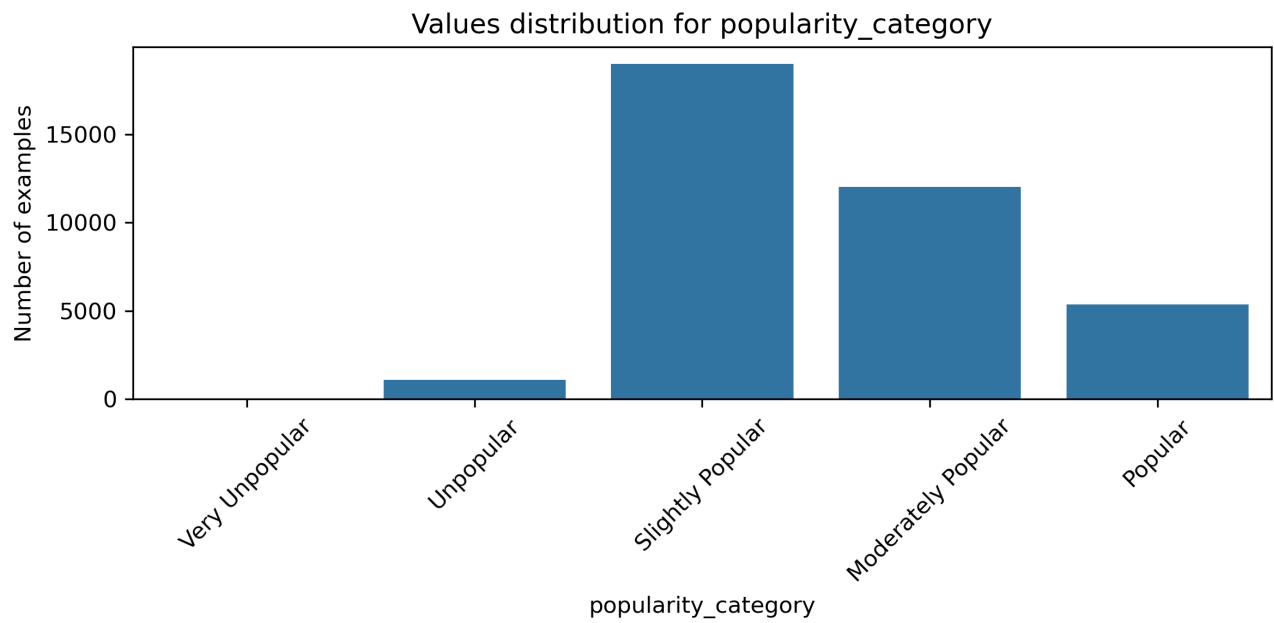


Majoritatea articolelor sunt publicate în timpul săptămânii (Weekday), cu o proporție mai mică în Weekend.

- popularity\_category — Number of non-null examples: 39644

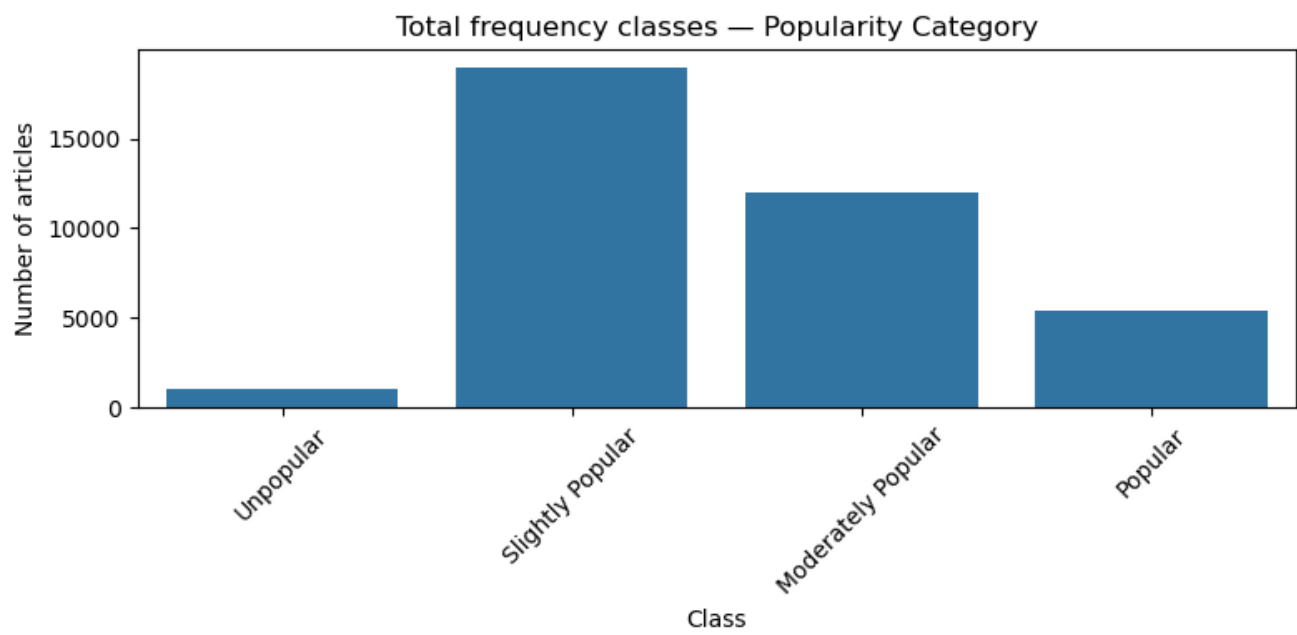


- popularity\_category — Number of unique values: 5

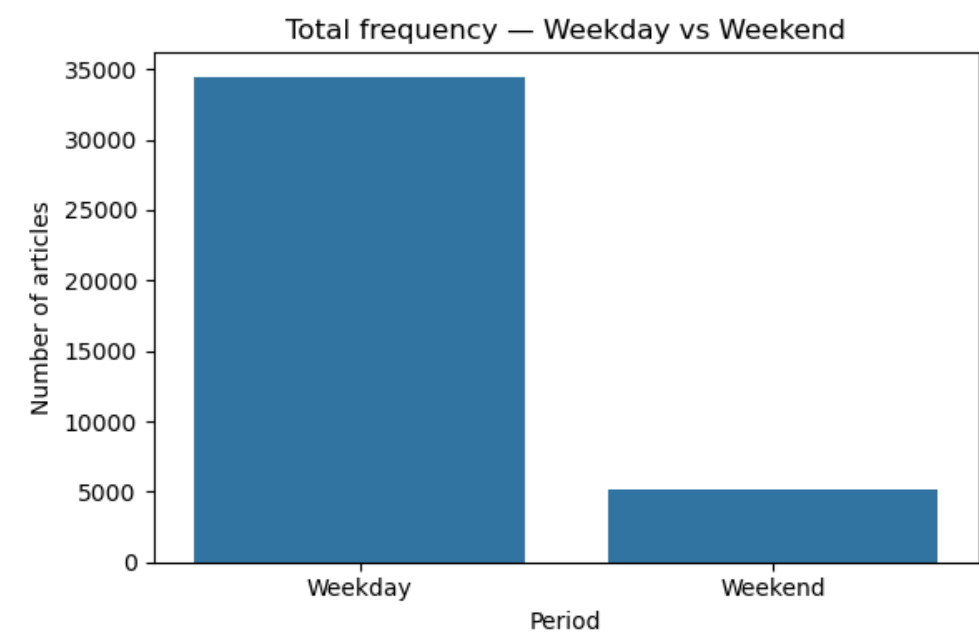


Cea mai frecventă clasă este Slightly Popular, urmată de Moderately Popular. Clasele extreme, precum Very Unpopular, sunt rare, ceea ce sugerează o etichetare dezechilibrată și potențială dificultate în clasificare.

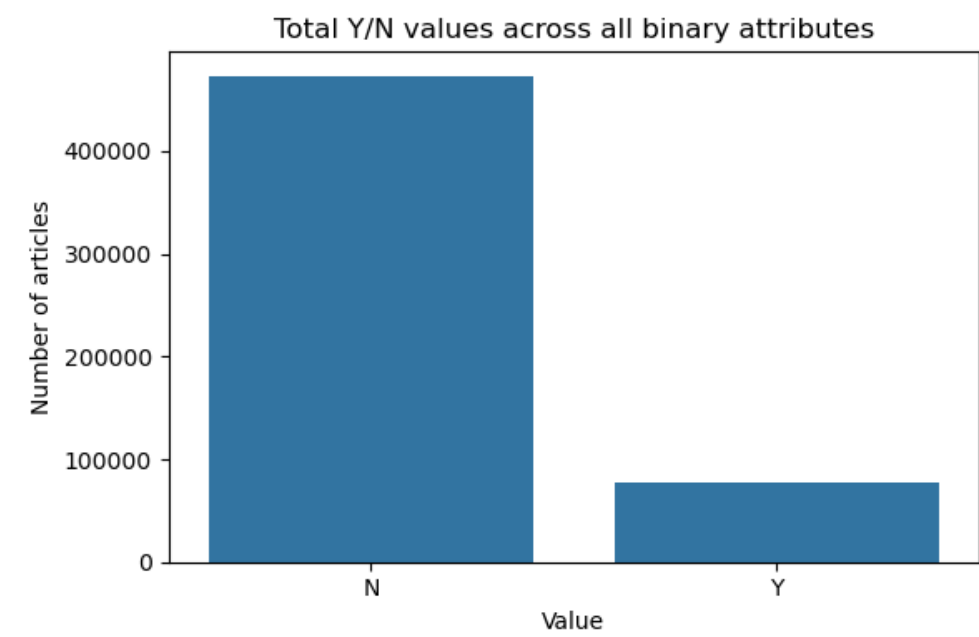
Echilibrul Claselor



Graficul arată un dezechilibru semnificativ între clase: cele mai multe articole sunt etichetate ca Slightly Popular și Moderately Popular, în timp ce clasele extreme (Very Unpopular, Popular) sunt subreprezentate. Acest dezechilibru poate afecta algoritmii de clasificare, care tind să favorizeze clasele dominante.

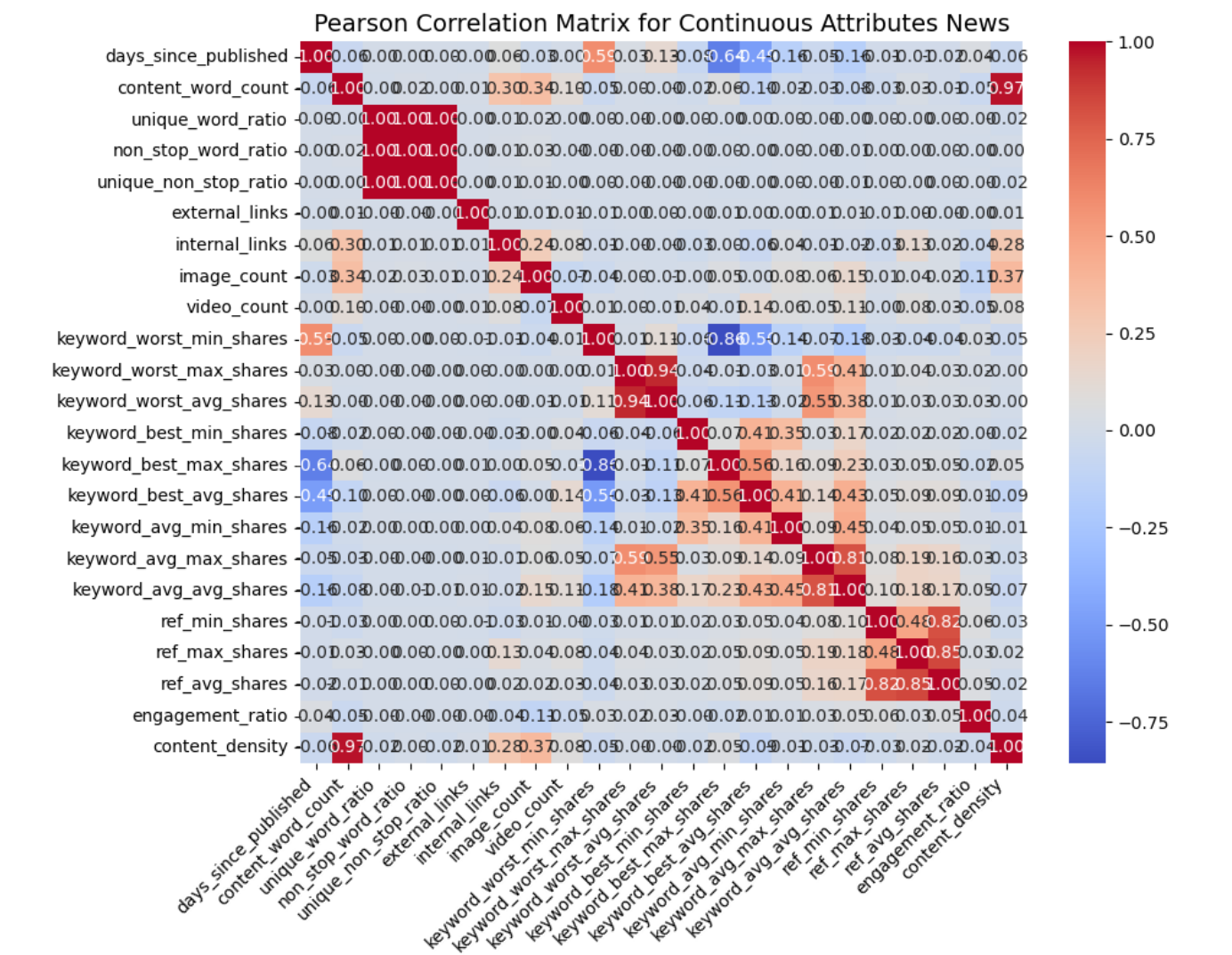


Majoritatea articolelor sunt publicate în timpul săptămânii (Weekday), ceea ce reflectă o practică obișnuită în media online. Articolele din Weekend reprezintă un procent mult mai mic, sugerând un potențial efect de sezon sau comportament diferit al publicului.

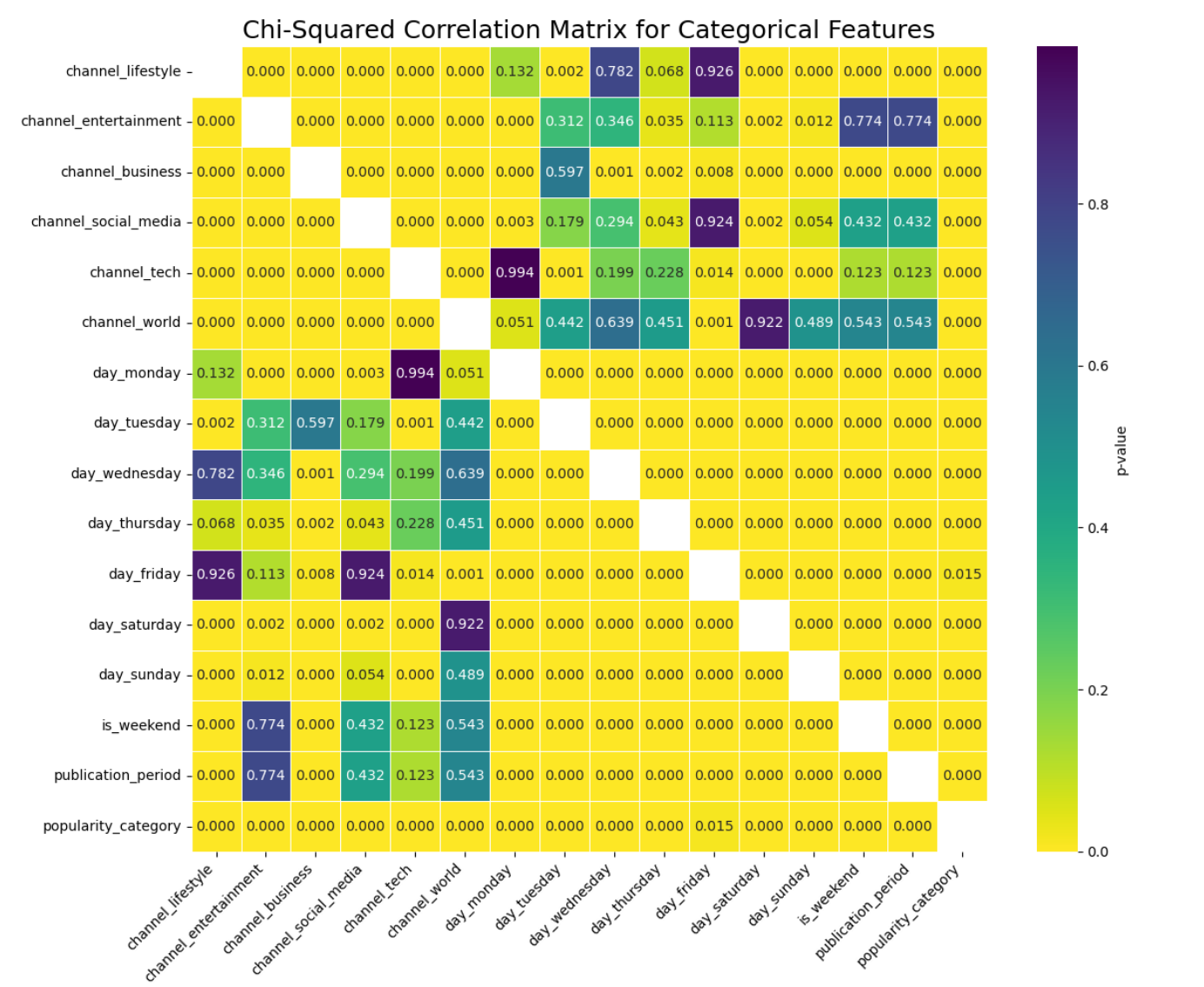


Distribuția totală a valorilor Y și N pentru toate coloanele binare (precum `channel_`, `day_`) este profund dezechilibrată în favoarea valorii N. Asta indică faptul că articolele aparțin rar unui canal sau unei zile anume. Acest dezechilibru trebuie tratat cu atenție în modelele de clasificare (ex. prin ponderare sau sampling).

**Corelatia intre atribute**



Această matrice arată corelațiile Pearson între attributele numerice continue. Se observă corelații foarte puternice între: **content\_word\_count** și **content\_density** (0.97) – articolele mai lungi tind să aibă o densitate de conținut mai mare. Variabilele privind partajările (shares) ale cuvintelor cheie (keyword\_\*) sunt corelate pozitiv între ele (valori între 0.4 și 0.8). **ref\_avg\_shares** și **ref\_max\_shares** sau **ref\_min\_shares** sunt, de asemenea, corelate moderat spre puternic. În general, există grupuri clare de variabile corelate, utile pentru reducerea dimensionalității sau selecția atributelor.



Această matrice exprimă semnificația statistică a relației dintre atributele categorice (inclusiv cele ordinale), folosind testul Chi-pătrat (p-value). Valorile apropiate de 0 (culoare galbenă intensă) indică o relație semnificativă statistic între variabile. Se observă corelații puternice între: Zilele săptămânii (**day\_monday**, **day\_tuesday**, etc.) și **is\_weekend** sau **publication\_period**. Canalele (**channel\_\***) și **publication\_period**, sugerând că unele tipuri de conținut sunt mai frecvente în anumite perioade. **popularity\_category** este semnificativ asociat cu majoritatea celorlalte variabile categorice, sugerând că popularitatea articolelor este influențată de mulți factori discreți.

## Preprocesare Date

Air pollution data set

### Atribute cu date lipsa si imputarea lor

Attributes with missing values: CO\_Category 1893 Ozone\_Value 1870 Country 349

After imputation, missing values in train: CO\_Category 0 Ozone\_Value 0 Country 0

### Valorile extreme si inlocuirea acestora cu cele din imputare

[Train] AQI\_Value: 2354 outliers replaced with NaN [Train] CO\_Value: 6858 outliers replaced with NaN [Train] Ozone\_Value: 1219 outliers replaced with NaN [Train] NO2\_Value: 2038 outliers replaced with NaN [Train] PM25\_Value: 2160 outliers replaced with NaN [Train] VOCs: 2331 outliers replaced with NaN [Train] SO2: 240 outliers replaced with NaN

Replacing outliers with NaN in airDataTest using train IQR... [Test] AQI\_Value: 581 outliers replaced with NaN [Test] CO\_Value: 1738 outliers replaced with NaN [Test] Ozone\_Value: 284 outliers replaced with NaN [Test] NO2\_Value: 516 outliers replaced with NaN

[Test] PM25\_Value: 547 outliers replaced with NaN [Test] VOCs: 576 outliers replaced with NaN [Test] SO2: 56 outliers replaced with NaN

Outliers replaced with NaN and imputed with column mean.

### Variabile redundante

Removed redundant features: ['VOCs', 'PM25\_Value', 'CO\_Category', 'Ozone\_Category', 'NO2\_Category', 'PM25\_Category', 'Emissions', 'City', 'Country']

### Standardizare

Pentru ca Regresia Logistica sa functioneze se standardizeaza valorile atributelor numerice

News pollution data set

### Atribute cu date lipsa si imputarea lor

Attributes with missing values: channel\_lifestyle 3175 content\_density 3145

After imputation, missing values in train: channel\_lifestyle 0 content\_density 0

### Valorile extreme si inlocuirea acestora cu cele din imputare

[Train] days\_since\_published: 0 outliers replaced with NaN [Train] title\_word\_count: 123 outliers replaced with NaN [Train] content\_word\_count: 1539 outliers replaced with NaN [Train] unique\_word\_ratio: 1290 outliers replaced with NaN [Train] non\_stop\_word\_ratio: 2234 outliers replaced with NaN [Train] unique\_non\_stop\_ratio: 1393 outliers replaced with NaN [Train] external\_links: 2665 outliers replaced with NaN [Train] internal\_links: 1674 outliers replaced with NaN [Train] image\_count: 6159 outliers replaced with NaN [Train] video\_count: 2330 outliers replaced with NaN [Train] avg\_word\_length: 1363 outliers replaced with NaN [Train] keyword\_count: 47 outliers replaced with NaN [Train] keyword\_worst\_min\_shares: 3743 outliers replaced with NaN [Train] keyword\_worst\_max\_shares: 2939 outliers replaced with NaN [Train] keyword\_worst\_avg\_shares: 1694 outliers replaced with NaN [Train] keyword\_best\_min\_shares: 4025 outliers replaced with NaN [Train] keyword\_best\_max\_shares: 7542 outliers replaced with NaN [Train] keyword\_best\_avg\_shares: 791 outliers replaced with NaN [Train] keyword\_avg\_min\_shares: 0 outliers replaced with NaN [Train] keyword\_avg\_max\_shares: 1893 outliers replaced with NaN [Train] keyword\_avg\_avg\_shares: 1284 outliers replaced with NaN [Train] ref\_min\_shares: 3979 outliers replaced with NaN [Train] ref\_max\_shares: 3450 outliers replaced with NaN [Train] ref\_avg\_shares: 3362 outliers replaced with NaN [Train] topic\_0\_relevance: 4238 outliers replaced with NaN [Train] topic\_1\_relevance: 4624 outliers replaced with NaN [Train] topic\_2\_relevance: 2864 outliers replaced with NaN [Train] topic\_3\_relevance: 211 outliers replaced with NaN [Train] topic\_4\_relevance: 0 outliers replaced with NaN [Train] content\_subjectivity: 1551 outliers replaced with NaN [Train] content\_sentiment: 682 outliers replaced with NaN [Train] positive\_word\_rate: 415 outliers replaced with NaN [Train] negative\_word\_rate: 1052 outliers replaced with NaN [Train] non\_neutral\_positive\_rate: 1289 outliers replaced with NaN [Train] non\_neutral\_negative\_rate: 390 outliers replaced with NaN [Train] avg\_positive\_sentiment: 1701 outliers replaced with NaN [Train] min\_positive\_sentiment: 2527 outliers replaced with NaN [Train] max\_positive\_sentiment: 0 outliers replaced with NaN [Train] avg\_negative\_sentiment: 704 outliers replaced with NaN [Train] min\_negative\_sentiment: 0 outliers replaced with NaN [Train] max\_negative\_sentiment: 1987 outliers replaced with NaN [Train] title\_subjectivity: 0 outliers replaced with NaN [Train] title\_sentiment: 6319 outliers replaced with NaN [Train] title\_subjectivity\_magnitude: 0 outliers replaced with NaN [Train] title\_sentiment\_magnitude: 1343 outliers replaced with NaN [Train] engagement\_ratio: 3225 outliers replaced with NaN [Train] content\_density: 2355 outliers replaced with NaN

Replacing outliers with NaN in newsDataTrain using train IQR... [Test] days\_since\_published: 0 outliers replaced with NaN [Test] title\_word\_count: 0 outliers replaced with NaN [Test] content\_word\_count: 0 outliers replaced with NaN [Test] unique\_word\_ratio: 0 outliers replaced with NaN [Test] non\_stop\_word\_ratio: 0 outliers replaced with NaN [Test] unique\_non\_stop\_ratio: 0 outliers replaced with NaN [Test] external\_links: 0 outliers replaced with NaN [Test] internal\_links: 0 outliers replaced with NaN [Test] image\_count: 0 outliers replaced with NaN [Test] video\_count: 0 outliers replaced with NaN [Test] avg\_word\_length: 0 outliers replaced with NaN [Test] keyword\_count: 0 outliers replaced with NaN [Test] keyword\_worst\_min\_shares: 0 outliers replaced with NaN [Test] keyword\_worst\_max\_shares: 0 outliers replaced with NaN [Test] keyword\_worst\_avg\_shares: 0 outliers replaced with NaN [Test] keyword\_best\_min\_shares: 0 outliers replaced with NaN [Test] keyword\_best\_max\_shares: 0 outliers replaced with NaN [Test] keyword\_best\_avg\_shares: 0 outliers replaced with NaN [Test] keyword\_avg\_min\_shares: 0 outliers replaced with NaN [Test] keyword\_avg\_max\_shares: 0 outliers replaced with NaN [Test] keyword\_avg\_avg\_shares: 0 outliers replaced with NaN [Test]

ref\_min\_shares: 0 outliers replaced with NaN [Test] ref\_max\_shares: 0 outliers replaced with NaN [Test] ref\_avg\_shares: 0 outliers replaced with NaN [Test] topic\_0\_relevance: 0 outliers replaced with NaN [Test] topic\_1\_relevance: 0 outliers replaced with NaN [Test] topic\_2\_relevance: 0 outliers replaced with NaN [Test] topic\_3\_relevance: 0 outliers replaced with NaN [Test] topic\_4\_relevance: 0 outliers replaced with NaN [Test] content\_subjectivity: 0 outliers replaced with NaN [Test] content\_sentiment: 0 outliers replaced with NaN [Test] positive\_word\_rate: 0 outliers replaced with NaN [Test] negative\_word\_rate: 0 outliers replaced with NaN [Test] non\_neutral\_positive\_rate: 0 outliers replaced with NaN [Test] non\_neutral\_negative\_rate: 0 outliers replaced with NaN [Test] avg\_positive\_sentiment: 0 outliers replaced with NaN [Test] min\_positive\_sentiment: 0 outliers replaced with NaN [Test] max\_positive\_sentiment: 0 outliers replaced with NaN [Test] avg\_negative\_sentiment: 0 outliers replaced with NaN [Test] min\_negative\_sentiment: 0 outliers replaced with NaN [Test] max\_negative\_sentiment: 0 outliers replaced with NaN [Test] title\_subjectivity: 0 outliers replaced with NaN [Test] title\_sentiment: 0 outliers replaced with NaN [Test] title\_subjectivity\_magnitude: 0 outliers replaced with NaN [Test] title\_sentiment\_magnitude: 0 outliers replaced with NaN [Test] engagement\_ratio: 0 outliers replaced with NaN [Test] content\_density: 0 outliers replaced with NaN

Outliers replaced with NaN and imputed with column mean.

Variabile redundante

Removed redundant features: ['content\_density', 'non\_stop\_word\_ratio', 'unique\_non\_stop\_ratio', 'keyword\_best\_avg\_shares', 'keyword\_avg\_avg\_shares', 'ref\_avg\_shares', 'publication\_period', 'is\_weekend', 'channel\_business', 'channel\_entertainment', 'channel\_lifestyle', 'channel\_social\_media', 'channel\_tech', 'channel\_world', 'day\_monday', 'day\_tuesday', 'day\_wednesday', 'day\_thursday', 'day\_friday', 'day\_saturday', 'day\_sunday']

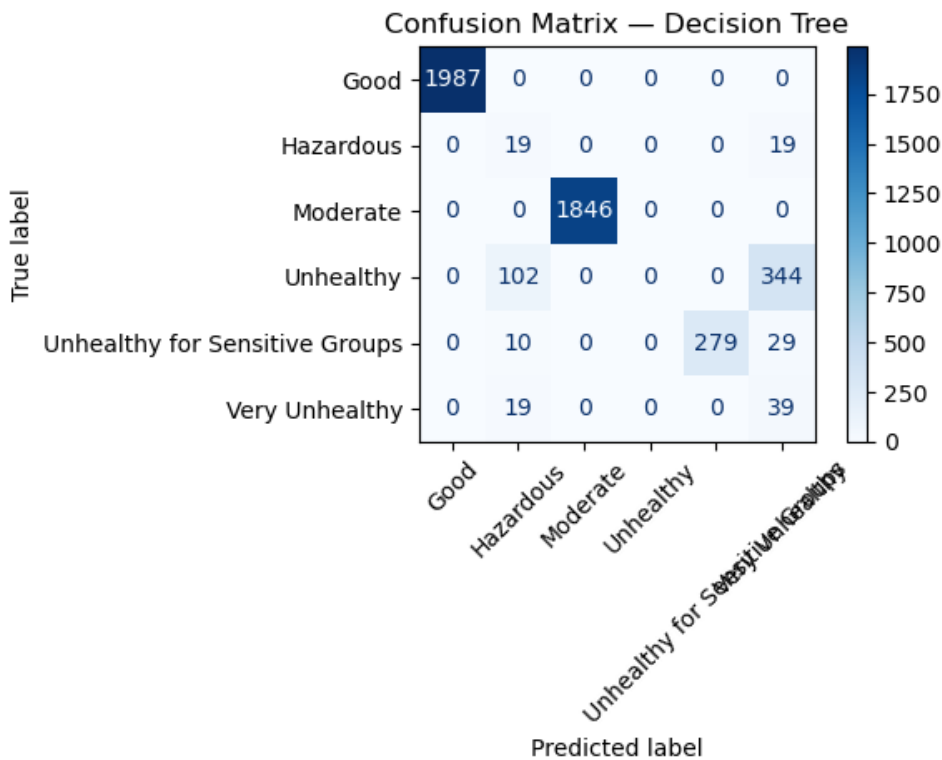
Standardizare

Pentru ca Regresia Logistica sa functioneze se standardizeaza valorile atributelor numerice

Algoritmi

Air pollution data set

Arbori de decizie



Classification Report — Decision Tree

Class	Precision	Recall	F1-score	Support
-------	-----------	--------	----------	---------

Class	Precision	Recall	F1-score	Support
Good	1.00	1.00	1.00	1987
Hazardous	0.13	0.50	0.20	38
Moderate	1.00	1.00	1.00	1846
Unhealthy	0.00	0.00	0.00	446
Unhealthy for Sensitive Groups	1.00	0.88	0.93	318
Very Unhealthy	0.09	0.67	0.16	58
Accuracy			0.89	4693
Macro Avg	0.54	0.67	0.55	4693
Weighted Avg	0.89	0.89	0.88	4693

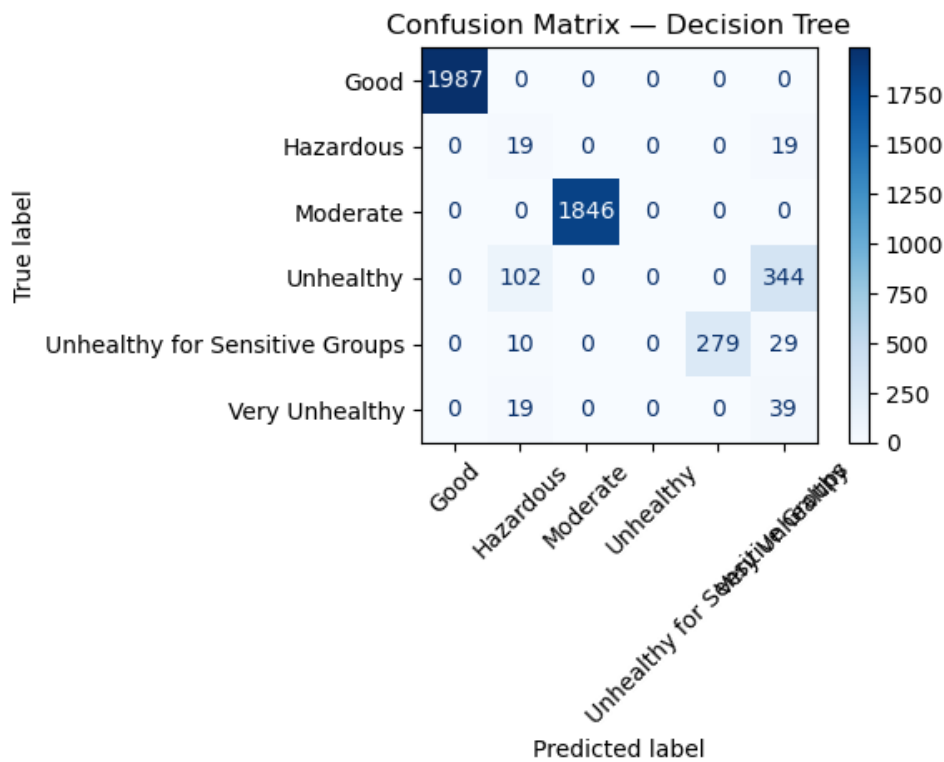
Accuracy (Decision Tree): 0.8886

Un model de arbore de decizie folosind biblioteca scikit-learn, cu scopul de a prezice categoria calității aerului (AQI\_Category). Modelul a fost antrenat pe setul de date procesat (X\_train, y\_train), cu următorii hiperparametri:

- Criteriu: entropy — pentru a maximiza informația câștigată la fiecare split
- Adâncime maximă (max\_depth): 4 — pentru a evita suprapotrivirea
- Număr minim de exemple într-o frunză (min\_samples\_leaf): 5 — pentru a evita frunze rare
- Ponderare a claselor (class\_weight): balanced — pentru a contracara distribuția dezechilibrată a claselor

Deși modelul obține o acuratețe generală ridicată (88.86%), scorurile de precizie și recall pentru clasele subreprezentate (ex. „Hazardous”, „Unhealthy”) sunt foarte mici. Aceasta indică faptul că modelul tinde să favorizeze clasele frecvente, fiind ineficient în identificarea corectă a celor rare. Într-un context real, acest lucru poate fi problematic dacă aceste clase reprezintă condiții critice.

Paduri Aleatoare



Class	Precision	Recall	F1-score	Support
-------	-----------	--------	----------	---------

Class	Precision	Recall	F1-score	Support
Good	1.00	1.00	1.00	1987
Hazardous	0.08	0.71	0.14	38
Moderate	0.98	0.94	0.96	1846
Unhealthy	0.87	0.26	0.40	446
Unhealthy for Sensitive Groups	1.00	0.88	0.93	318
Very Unhealthy	0.12	0.36	0.19	58
Accuracy			0.89	4693
Macro Avg	0.67	0.69	0.60	4693
Weighted Avg	0.96	0.89	0.91	4693

Accuracy (Random Forest): 0.8881

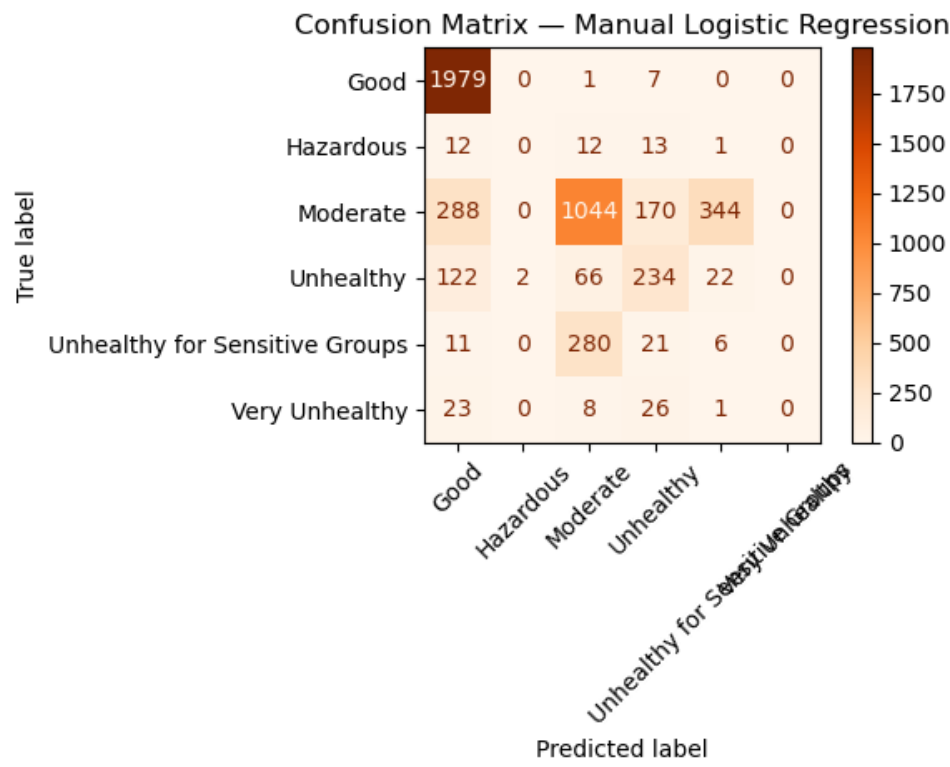
Un model Random Forest a fost antrenat folosind biblioteca scikit-learn pentru a prezice categoria calității aerului (AQI\_Category). Modelul a fost antrenat pe datele procesate (X\_train, y\_train) și a utilizat următorii hiperparametri:

- Criteriu: entropy — maximizează câștigul informațional la fiecare split;
- Adâncime maximă (max\_depth): 4 — limitează complexitatea fiecărui arbore pentru a preveni suprapotrivirea;
- Număr minim de exemple într-o frunză (min\_samples\_leaf): 5 — previne formarea frunzelor cu puține instanțe;
- Ponderare a claselor (class\_weight): balanced — ajustează automat greutatea în funcție de distribuția claselor, fiind util pentru seturi dezechilibrate;
- Număr de estimatori (n\_estimators): 100 — utilizează 100 de arbori de decizie individuali;
- Proporție de eșantion pentru fiecare estimator: implicită (bootstrap=True) — fiecare arbore este antrenat pe un subset bootstrap diferit al datelor;
- Proporție de attribute utilizate de fiecare arbore: implicită (max\_features='sqrt') — fiecare arbore selectează aleator  $\sqrt{n}$  attribute la fiecare split, ceea ce ajută la reducerea corelației dintre arbori și îmbunătățirea generalizării.

Modelul obține o acuratețe globală de 88.9%, dar performanța pe clasele dezechilibrate rămâne limitată. De exemplu, pentru clasa „Hazardous”, modelul are un scor de recall de 71%, dar precizia este doar 8%, indicând un număr ridicat de fals pozitive. De asemenea, clase precum „Unhealthy” sau „Very Unhealthy” au dificultăți în a fi clasificate corect. Această situație sugerează că, deși Random Forest îmbunătățește ușor performanțele față de arborele de decizie simplu, provocările legate de dezechilibrul claselor persistă. În aplicații critice (ex. poluare extremă), aceste limitări pot afecta capacitatea sistemului de a reacționa adecvat la condiții periculoase.

## Regresie Logistica





Class	Precision	Recall	F1-score	Support
Good	0.81	1.00	0.90	1987
Hazardous	0.00	0.00	0.00	38
Moderate	0.74	0.57	0.64	1846
Unhealthy	0.50	0.52	0.51	446
Unhealthy for Sensitive Groups	0.02	0.02	0.02	318
Very Unhealthy	0.00	0.00	0.00	58
Accuracy			0.70	4693
Macro Avg	0.34	0.35	0.34	4693
Weighted Avg	0.68	0.70	0.68	4693

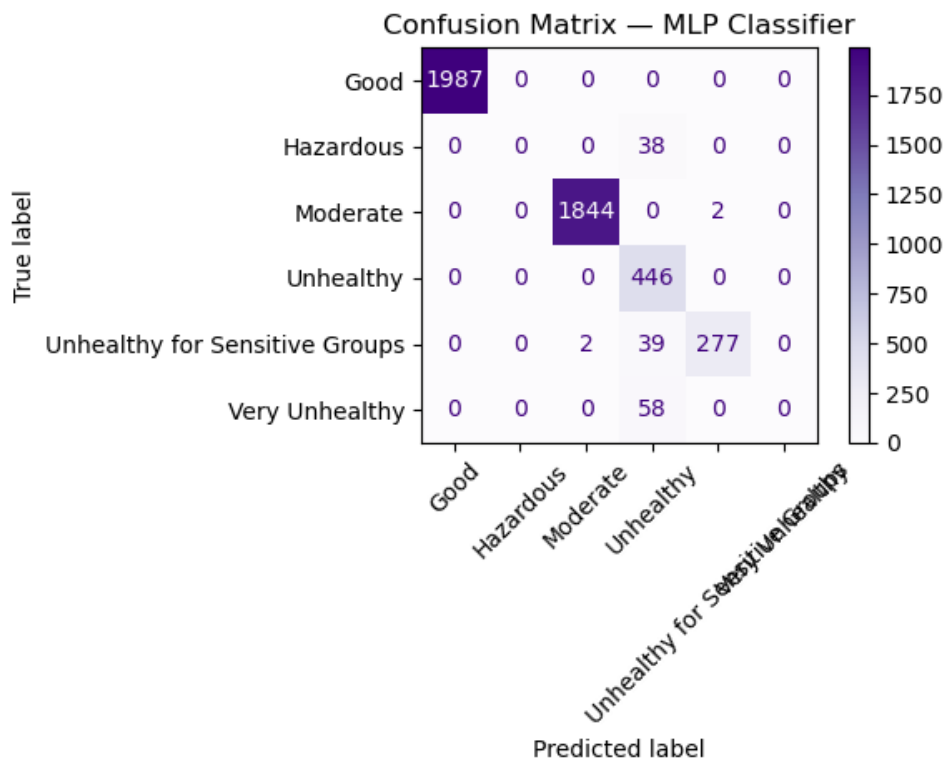
Accuracy (Logistic Regression): 0.6953

Un model de regresie logistică a fost implementat manual și antrenat pentru a clasifica calitatea aerului (AQI\_Category) folosind un optimizator de tip gradient descent. Acest model a fost testat pe un set de date cu mai multe clase, inclusiv „Good”, „Moderate”, „Unhealthy” și „Hazardous”. Pentru preprocesarea datelor:

- Atributele categorice au fost transformate folosind One-Hot Encoding, aplicat exclusiv pe coloanele cu cardinalitate redusă (prin ColumnTransformer).
- Atributele numerice au fost păstrate nemodificate (remainder='passthrough'). Algoritmul de optimizare a utilizat:
- Gradient Descent full-batch
- Learning rate: 0.1
- Număr epoci: 3000
- Inițializare greutăți: distribuție normală (N(0,1))
- Regularizare: L2 (cu  $\lambda = 0.01$ ) pentru a penaliza coeficienții mari și a reduce overfitting-ul Modelul a obținut o acuratețe globală de 69.5%, ceea ce indică un rezultat rezonabil pentru un model liniar aplicat pe un set dezechilibrat. Cu toate acestea, performanța per clasă evidențiază limitele:
- Clasa „Good” este clasificată excelent (recall 100%, f1-score 0.90)

- Clasele „Very Unhealthy” și „Hazardous” nu sunt identificate corect (precizie și recall 0), din cauza dezechilibrului și a numărului redus de exemple Această distribuție sugerează că regresia logistică are dificultăți în a separa corect clasele mai rare sau care se suprapun în spațiul atributelor. Pentru aplicații critice de mediu, acest lucru limitează utilitatea modelului fără metode suplimentare de balansare.

MLP



Class	Precision	Recall	F1-score	Support
Good	1.00	1.00	1.00	1987
Hazardous	0.00	0.00	0.00	38
Moderate	0.99	0.98	0.98	1846
Unhealthy	0.72	0.99	0.83	446
Unhealthy for Sensitive Groups	1.00	0.84	0.91	318
Very Unhealthy	0.00	0.00	0.00	58
Accuracy			0.96	4693
Macro Avg	0.62	0.63	0.62	4693
Weighted Avg	0.95	0.96	0.95	4693

Accuracy (MLP): 0.9593

Un model de rețea neurală de tip MLP (Multi-Layer Perceptron) a fost antrenat pentru a prezice clasa de calitate a aerului (AQI\_Category). Modelul a fost antrenat pe setul de date procesat cu următorii hiperparametri:

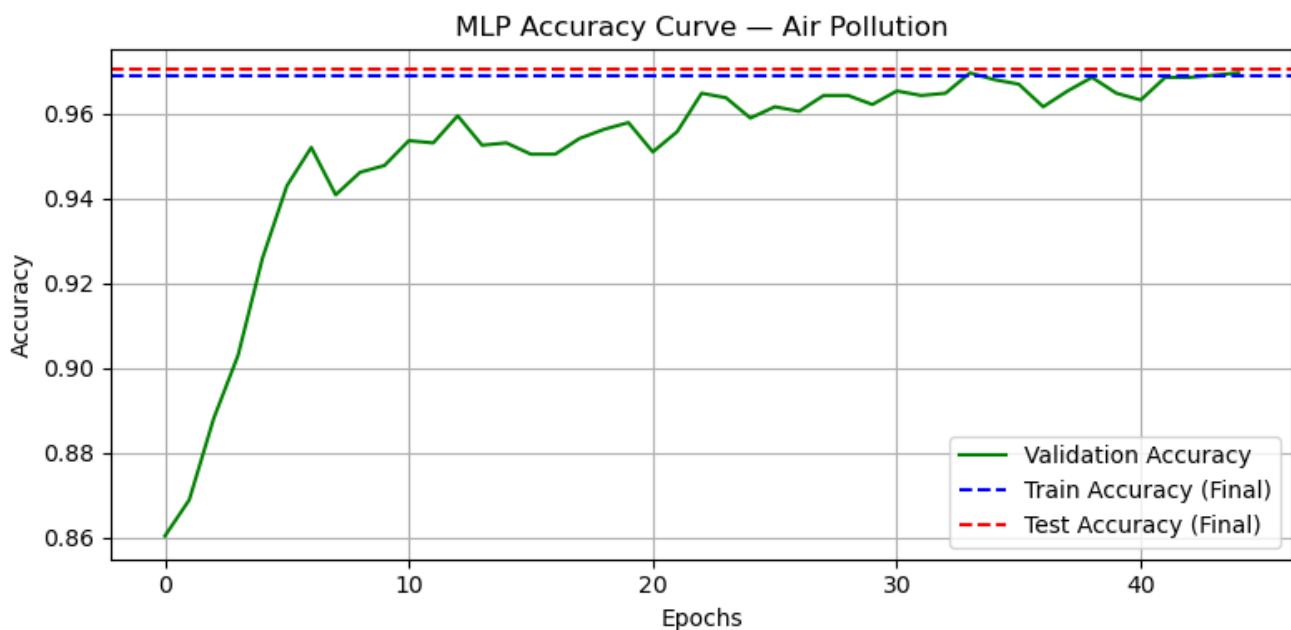
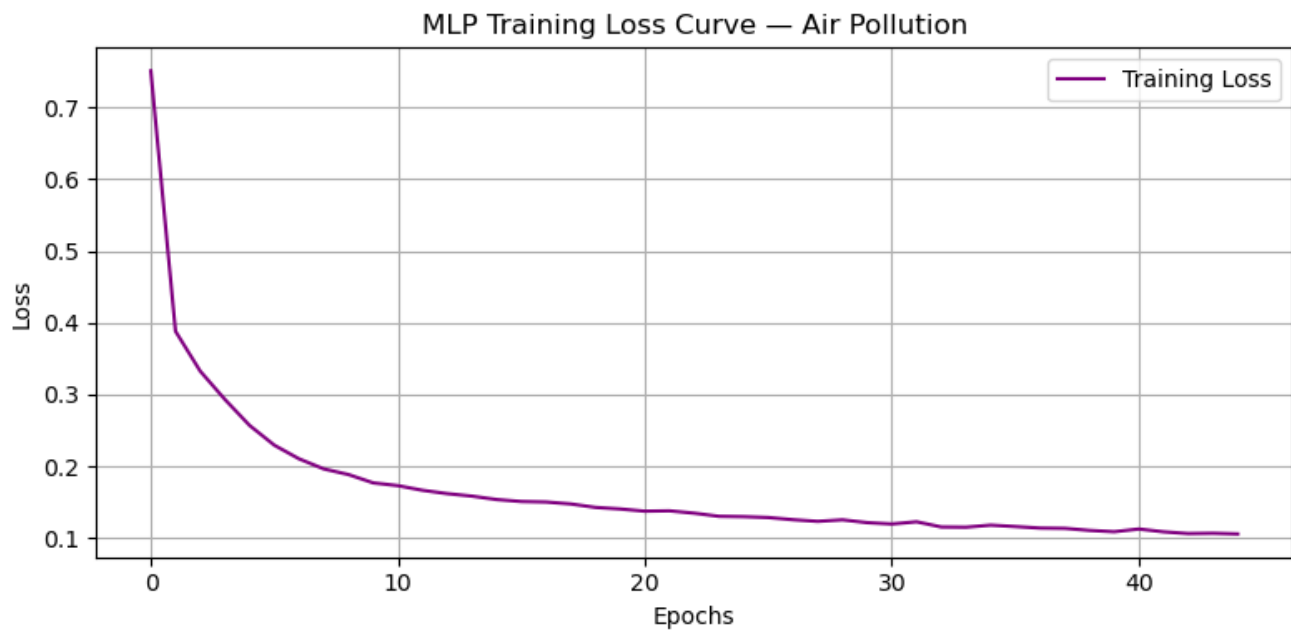
- Arhitectură: un singur strat ascuns cu 100 de neuroni
- Funcția de activare: ReLU
- Optimizator: Adam
- Learning rate: 0.001

- Număr maxim de epoci: 200
- Criteriu de oprire timpurie (early\_stopping=True) pentru prevenirea suprapotrivirii
- Coeficient de regularizare:  $\alpha=0.0001$

Modelul a atins o acuratețe globală ridicată de 95.93%, dar performanța este puternic dezechilibrată între clase:

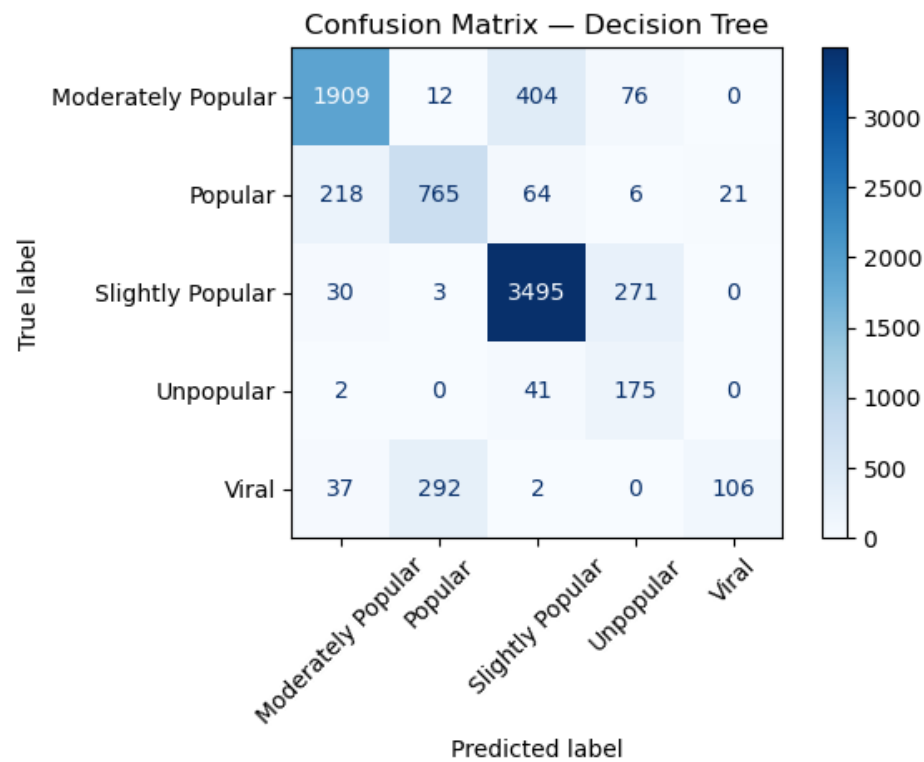
- Clasele majore precum Good, Moderate sau Unhealthy for Sensitive Groups sunt clasificate corect în proporție foarte mare.
- Clasele rare precum Hazardous și Very Unhealthy au precizie și recall 0, ceea ce înseamnă că modelul nu le identifică deloc corect.
- Clasa Unhealthy are un recall foarte bun (0.99), dar o precizie scăzută (0.72), ceea ce indică existența multor fals pozitive.

Modelul este eficient în identificarea claselor frecvente, dar insuficient pentru aplicații critice unde recunoașterea corectă a stărilor periculoase (ex. Hazardous) este esențială.



News pollution data set

Arbori de decizie



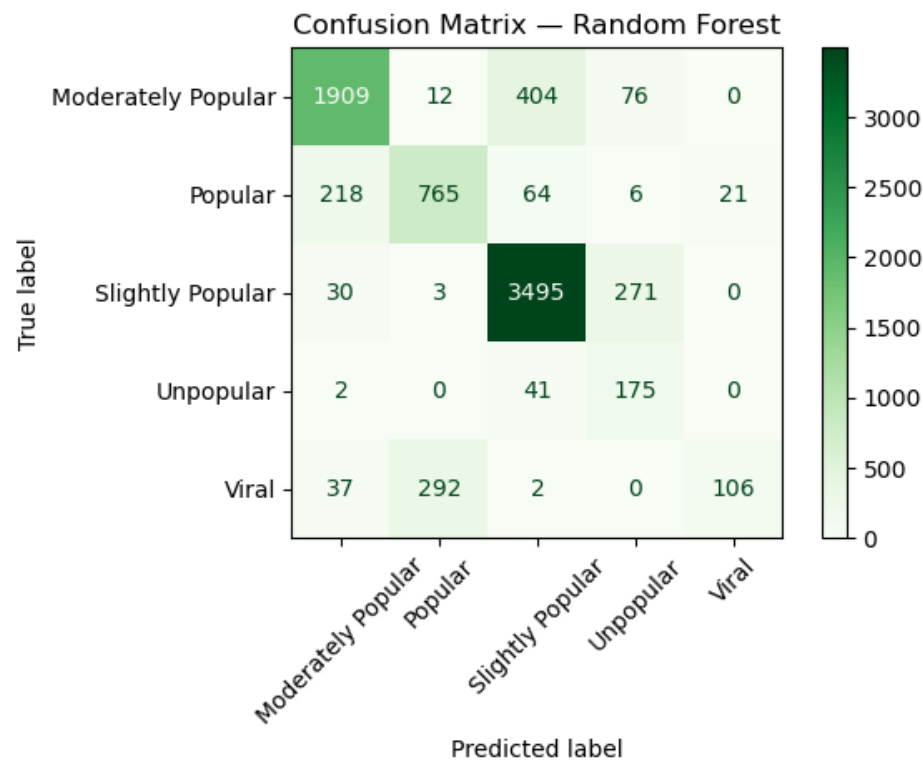
Class	Precision	Recall	F1-score	Support
Moderately Popular	0.87	0.80	0.83	2401
Popular	0.71	0.71	0.71	1074
Slightly Popular	0.87	0.92	0.90	3799
Unpopular	0.33	0.80	0.47	218
Viral	0.84	0.24	0.37	437
Accuracy			0.8135	7929
Macro avg			0.66	7929
Weighted avg			0.81	7929

Modelul de arbore de decizie antrenat pentru clasificarea popularității știrilor a atins o acuratețe generală de 81.35%, un scor foarte bun având în vedere numărul de clase și natura dezechilibrată a setului de date.

- Clasele „Slightly Popular”, „Moderately Popular” și „Popular” sunt bine clasificate, cu f1-scores de peste 0.70 și recall ridicat.
- Clasa „Unpopular”, deși are un număr mic de exemple, este surprinzător de bine captată (recall 0.80), dar cu precizie slabă (0.33), ceea ce indică multe false positive.
- Clasa „Viral” este mai dificil de identificat: are o precizie bună (0.84), dar un recall foarte mic (0.24), ceea ce sugerează că modelul tinde să fie conservator și să nu clasifice decât foarte puține articole drept virale — ceea ce poate fi acceptabil dacă se dorește evitarea alarmelor false. Modelul a fost antrenat cu hiperparametri aleși atent:
- max\_depth=20 pentru flexibilitate crescută,
- min\_samples\_leaf=2 pentru a reduce overfitting,
- class\_weight='balanced' pentru a trata dezechilibrul dintre clase.

Astfel, modelul oferă o performanță solidă și este potrivit ca soluție de bază pentru predicția popularității știrilor. Clasele extreme, cum ar fi „Viral” sau „Unpopular”, pot fi îmbunătățite suplimentar prin oversampling, ensemble methods sau modele mai complexe precum Random Forest sau MLP.

Paduri Aleatoare



Class	Precision	Recall	F1-score	Support
Moderately Popular	0.87	0.80	0.83	2401
Popular	0.71	0.71	0.71	1074
Slightly Popular	0.87	0.92	0.90	3799
Unpopular	0.33	0.80	0.47	218
Viral	0.84	0.24	0.37	437
Accuracy			0.8135	7929
Macro avg	0.73	0.69	0.66	7929
Weighted avg	0.83	0.81	0.81	7929

Accuracy (Random Forest): 0.8135

Modelul Random Forest a obținut o acuratețe generală de 81.35% în clasificarea popularității articolelor de știri, oferind rezultate aproape identice cu arborele de decizie individual, dar cu un plus de robustețe și stabilitate.

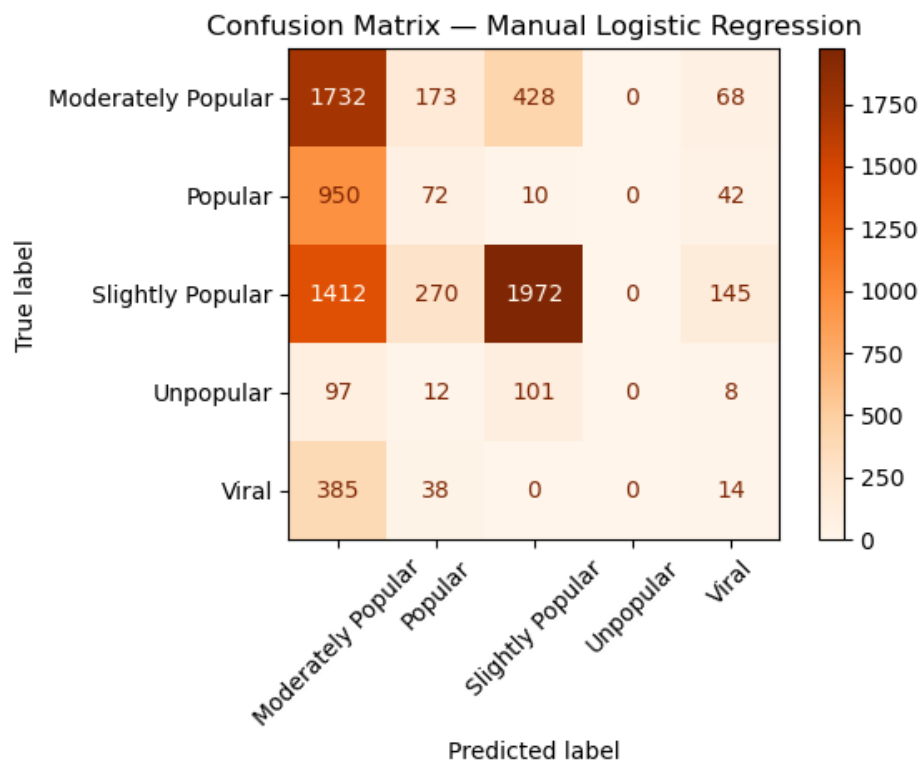
- Performanța este foarte bună pentru clasele „Slightly Popular”, „Moderately Popular” și „Popular”, care reprezintă majoritatea setului. Acestea au f1-score-uri de peste 0.70 și un recall foarte ridicat.
- Clasa „Unpopular” este captată bine în recall (0.80), dar suferă la precizie (0.33), ceea ce sugerează un număr mare de clasificări false pozitive.
- Clasa „Viral” are din nou dificultăți de identificare, cu un recall de doar 0.24, deși precizia este bună (0.84), similar cu Decision Tree.

Modelul a fost antrenat cu următorii hiperparametri:

- n\_estimators=100: o pădure de 100 de arbori,
- max\_depth=20: pentru a permite învățarea relațiilor complexe,
- min\_samples\_leaf=2: pentru a reduce overfitting-ul pe frunze mici,
- criterion='entropy': pentru a maximiza câștigul informațional,- class\_weight='balanced': pentru a contrabalansa distribuția inegală a claselor.

Astfel, Random Forest este un model solid pentru această sarcină, oferind un echilibru foarte bun între precizie și generalizare. Pentru clasele subreprezentate, performanța poate fi îmbunătățită prin tehnici de oversampling, creșterea numărului de estimatori sau ajustarea ponderii claselor.

Regresie Logistica



Class	Precision	Recall	F1-score	Support
Moderately Popular	0.38	0.72	0.50	2401
Popular	0.13	0.07	0.09	1074
Slightly Popular	0.78	0.52	0.63	3799
Unpopular	0.00	0.00	0.00	218
Viral	0.05	0.03	0.04	437
Accuracy			0.48	7929
Macro Avg	0.27	0.27	0.25	7929
Weighted Avg	0.51	0.48	0.46	7929

Accuracy (Manual Logistic Regression): 0.4777

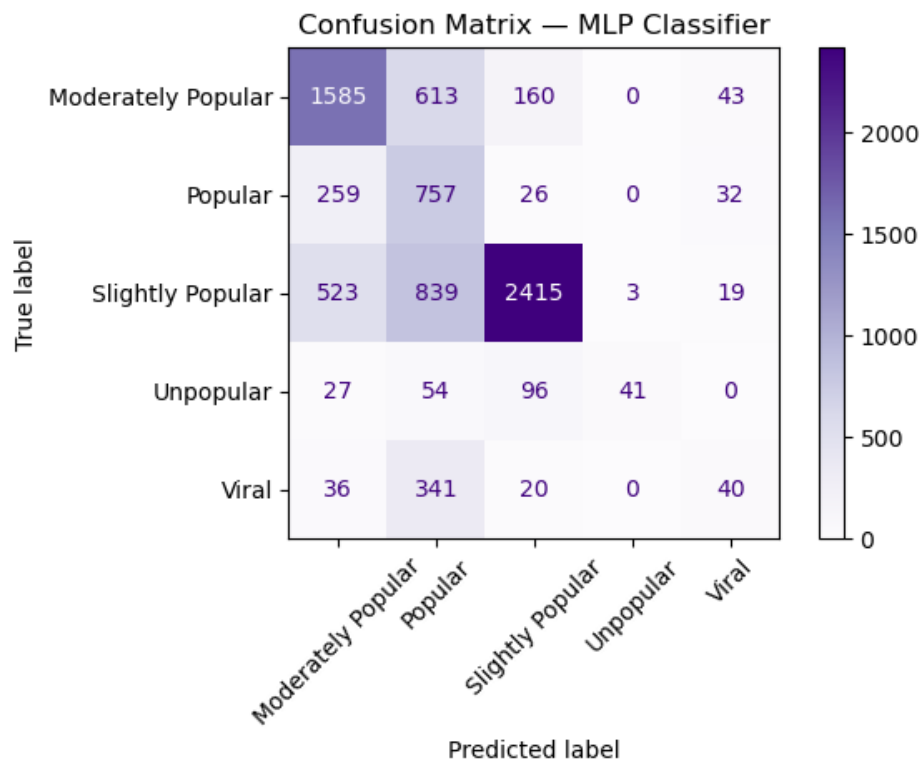
Pentru problema clasificării popularității știrilor, a fost implementat un model de regresie logistică multi-clasă, folosind o strategie One-vs-Rest (OvR). Aceasta presupune antrenarea câte unui model logistic binar pentru fiecare clasă posibilă, modelând fiecare „vs restul”. Encodare: Atributele categorice au fost prelucrate prin One-Hot Encoding (doar pentru cele cu cardinalitate redusă), folosind ColumnTransformer. Atributele numerice au fost păstrate nealterate. Optimizare:

- Algoritm: Gradient Descent standard
- Learning rate: 0.1
- Număr epoci: 3000
- Regularizare L2:  $\lambda = 0.01$  Fiecare model binar din OvR a fost antrenat independent, iar la inferență, scorurile de probabilitate pentru fiecare clasă au fost comparate pentru a alege predicția finală. Modelul a obținut o acuratețe globală de 47.7%, semnificativ mai bună decât o clasificare aleatorie, dar cu performanțe neuniforme:
- Clasele „Moderately Popular” și „Slightly Popular” sunt cele mai bine clasificate (recall 72% și 52%)

- Clasele „Popular” și „Viral” sunt adesea confundate cu clasele dominante
- Clasa „Unpopular” are precizie și recall aproape 0, sugerând confuzie sistematică cu clasele apropiate

Distribuția rezultatelor sugerează că, deși OvR permite aplicarea regresiei logistice la clasificare multi-clasă, separabilitatea liniară limitată și dezechilibrul între clase afectează puternic performanța.

MLP



Class	Precision	Recall	F1-score	Support
Moderately Popular	0.62	0.73	0.67	2401
Popular	0.38	0.64	0.48	1074
Slightly Popular	0.88	0.68	0.77	3799
Unpopular	0.65	0.22	0.33	218
Viral	0.17	0.13	0.15	437

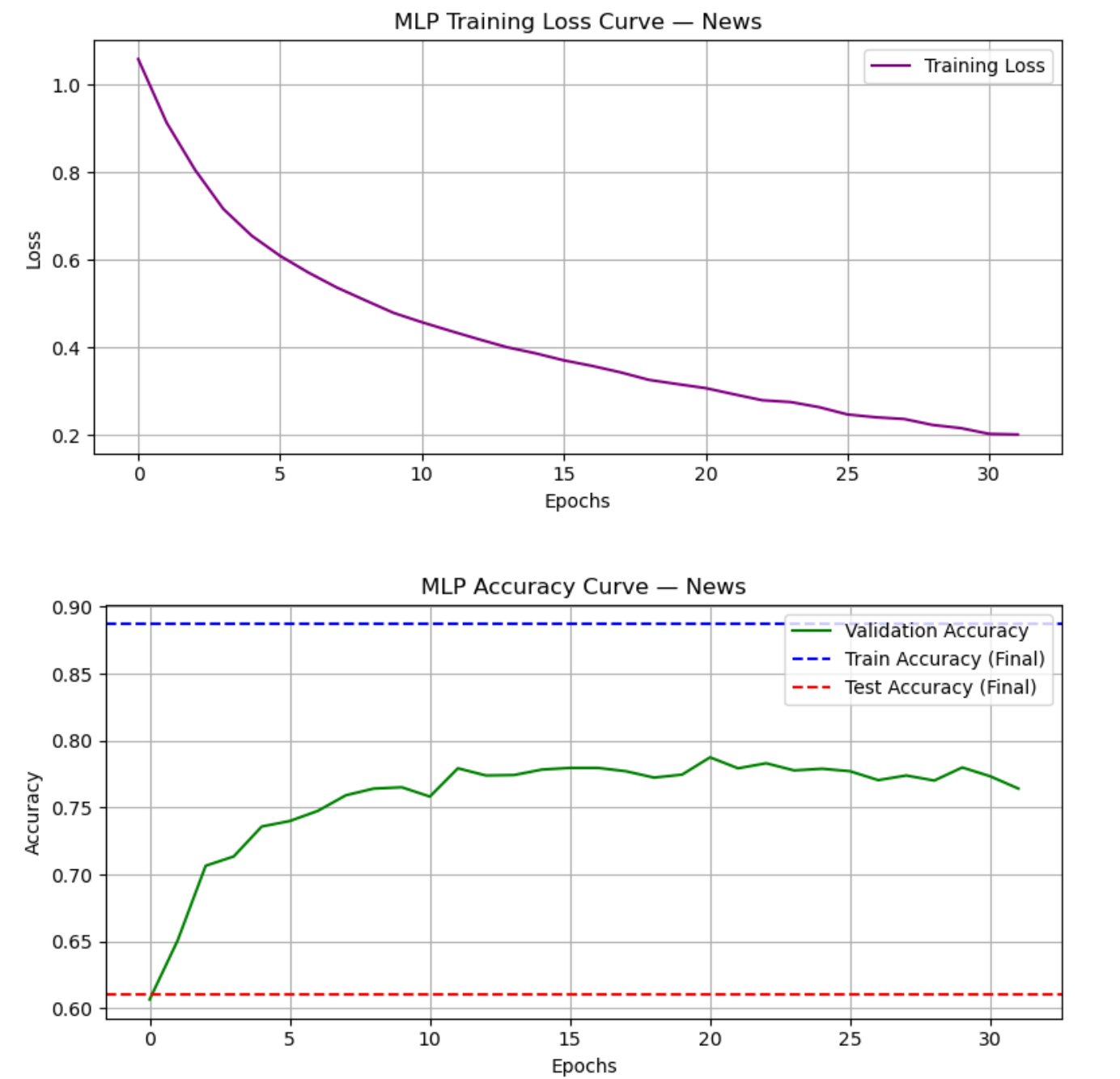
Un model MLP (Multi-Layer Perceptron) a fost antrenat folosind biblioteca scikit-learn pentru a prezice categoria de popularitate a articolelor de știri (popularity\_category). Modelul a fost antrenat pe datele procesate (X\_train, y\_train) și a utilizat următorii hiperparametri: Arhitectura:

- Straturi ascunde (hidden\_layer\_sizes): două straturi ascunde cu 256 și 128 de neuroni;
  - Funcție de activare (activation): relu, aleasă pentru eficiența sa în propagarea gradientului și învățarea non-liniară. Optimizare:
  - Algoritm (solver): adam, optimizator eficient pentru date mari și sparse;
  - Rată de învățare (learning\_rate\_init): 0.001 — valoare implicită adecvată pentru stabilitate;
  - Număr maxim de epoci (max\_iter): 3000 — permite rețelei suficiente cicluri de antrenament pentru convergență.
- Regularizare:
- Early stopping: activat (early\_stopping=True) — antrenarea se oprește când performanța de validare nu mai crește;
  - Regularizare L2 (alpha): 0.0001 — previne suprapotrivirea penalizând ponderile mari.

Modelul a obținut o acuratețe generală de 64.48%, cu următoarele observații importante extrase din matricea de confuzie:

- Slightly Popular și Moderately Popular sunt cel mai bine clasificate (f1-scores de 0.77 și 0.67), fiind și cele mai frecvente clase;
- Popular are o performanță decentă (f1 ≈ 0.48), dar confuziile sunt frecvente cu Moderately Popular;

- Clasele Unpopular și Viral sunt slab clasificate ( $f1 \approx 0.15\text{--}0.33$ ), cu numeroase confuzii către Popular și Slightly Popular.
- Modelul tinde să favorizeze clasele dominante și are dificultăți în a învăța reprezentări clare pentru clasele rare. De asemenea, observăm că Viral este frecvent confundată cu Popular, iar Unpopular este distribuită între mai multe clase.



Comparatie

Method	Accuracy-AIR	Accuracy-NEWS
dt	0.8886	0.8135
rf	0.9111	0.8135
lr	0.6953	0.4780
mlp	0.9704	0.6102