

MINISTERUL EDUCAȚIEI



UNIVERSITATEA TEHNICĂ
DIN CLUJ-NAPOCA
FACULTATEA DE AUTOMATICĂ ȘI CALCULATOARE

DETECTAREA AUTOMATĂ A LIMBAJULUI INSTIGATOR ȘI OFENSATOR PE TWITTER

PROIECT DE SEMESTRU

Student: **Morar Emilian Cristian**

Disciplina: **Sisteme bazate pe Cunoaștere**

2025

Cuprins

1	INTRODUCERE	2
1.1	CONTEXT GENERAL.....	2
1.2	OBIECTIV.....	2
1.3	SPECIFICAȚII.....	2
2	CUNOAȘTEREA ȘI ANALIZA SETULUI DE DATE	4
2.1	DETALII DESPRE MEDIUL DE IMPLEMENTARE.....	4
2.2	ANALIZA SETULUI DE DATE.....	4
3	PREPROCESAREA SETULUI DE DATE	8
3.1	METODELE UTILIZATE PENTRU PREPROCESAREA SETULUI DE DATE.....	8
3.2	REZULTATELE OBTINUTE SI INFORMAȚII NECESARE PENTRU REALIZAREA PROIECTULUI.....	9
4	MODELAREA SISTEMULUI	10
2.2	ANTRENAREMODEL1.PY – ANTRENAREA MODELULUI.....	10
2.2	PREDICTIE_MANUAL.PY - CLASIFICAREA UNUI TWEET INTRODUS MANUAL.....	11
2.2	PRELUARE_DATETw.PY - COLECTAREA DE TWEET-URI DE PE TWITTER.....	11
2.2	PREDICTIE.PY - CLASIFICAREA UNUI SET DE TWEET-URI.....	11
5	CONCLUZII	13
5.1	REZULTATE OBTINUTE.....	13
5.2	DIRECȚII DE DEZVOLTARE.....	13
6	BIBLIOGRAFIE	14

1 Introducere

1.1 Context general

Rețelele sociale au devenit un mediu predominant pentru exprimarea opiniilor, interacțiuni între oameni și distribuire de informații. Aceste platforme au dus și la o creștere a discursului instigator și a limbajului ofensator, putând avea efecte negative asupra sănătății mentale al utilizatorilor sau comunităților. Proiectul de Detectare Automată a Limbajului Instigator și Ofensator pe Twitter , abordează problema instigării și a injuriilor prin dezvoltarea unui model de învățare automată al acestora capabil să identifice și să cenzureze conținutul dăunător.

1.2 Obiective

Dezvoltarea unui Model de învățare automată: Scopul principal este de a dezvolta un model de învățare automată capabil să detecteze limbajul instigator și limbajul ofensator. Acest model va fi antrenat pe un set de date bine structurat , care conține tweet-uri anotate.

Implementarea unei soluții de moderare a conținutului: Crearea unei aplicații care să permită moderarea automată a conținutului. Aceasta va ajuta la identificarea și cenzurarea conținutului în timp real, asigurând un mediu sigur pentru utilizatori.

Cenzură limbajului ofensator: implementarea unei funcționalități care să cenzureze automat limbajul ofensator din postările utilizatorilor.

1.3 Specificații

Proiectul realizat este alcătuit din mai multe componente, fiecare având cerințe specifice. Funcționalitățile aplicației sunt:

1. Detectarea limbajului instigator și ofensator: Aplicația utilizează un model de învățare automată pe un set de date bine structurat pentru a clasifica tweet-urile în trei categorii: limbaj instigator, limbaj ofensator și neutru.

2. Cenzură limbajului ofensator: Utilizând biblioteca better profanity, aplicația va cenzura automat cuvintele ofensatoare din tweet-uri. (Neimplementat , poate fi adăugată această funcționalitate cu ajutorul bibliotecii better-profanity).

3. Calcularea procentajului de tweet-uri: Acesta va calcula și va afișa procentajul de tweet-uri din fiecare categorie (limbaj instigator, limbaj ofensator și neutru) pentru a oferi o imagine de ansamblu asupra conținutului analizat.

4. Partea în care am testat tweet-uri introduse manual.

5. Salvarea rezultatelor: Rezultatele analizei, inclusiv tweet-urile originale, cenzurate și clasificările, vor fi salvate într-un fișier CSV pentru revizuire ulterioară.

Nivelul de performanță atins de model are o acuratețe de minimum 88.94% în clasificările tweet-urilor.

Limitele modelului sunt incapacitatea lui de detectare a unui limbaj instigator sau ofensator subtil sau contextual. De asemenea, calitatea datelor de intrare poate influența performanța modelului. Altă limitare a modelului este antrenarea modelului cu date noi, putând fi implementată ulterior.

2 Cunoașterea și analiza setului de date

2.1 Detalii despre mediul de implementare

Setul de date train.csv este esențial pentru antrenarea modelului, acesta conține tweet-uri anotate, fiecare având etichete care indică natura conținutului. Setul de date conține următoarele coloane:

- tweet: Textul tweet-ului. Aceasta este informația principală care va fi analizată și clasificată.
- hate speech count: Numărul de notări care clasifică tweetul ca limbaj instigator.
- offensive language count: Numărul de notări care clasifică tweetul ca limbaj ofensator.
- neither count: Numărul de notări care clasifică tweetul ca fiind nici limbaj instigator sau ofensator.
- count: Numărul total de notări pentru tweet. Acesta este suma tuturor notărilor din celelalte coloane.

Bibliotecile și instrumentele utilizate pentru crearea aplicației:

- Pandas: utilizată pentru manipularea și analiza datelor, permite citirea seturilor de date CSV și efectuarea de operații de prelucrare a datelor.
- NumPy: folosită pentru operații numerice și gestionarea array-urilor.
- Scikit Learn: folosită pentru implementarea algoritmilor de învățare automată, inclusiv regresia logistică și evaluarea modelului.
- NLTK: utilizată pentru procesarea limbajului natural, inclusiv pentru eliminarea cuvintelor frecvențe și reducerea dimensionalității prin aducerea cuvintelor la forma lor de bază.
- Matplotlib și Seaborn: folosite pentru vizualizarea datelor, graficelor de distribuție și matricea de corelație.
- Joblib: utilizată pentru salvarea și încărcarea modelului antrenat.

Proiectul a fost dezvoltat în Python 3.11.9, folosind ca IDE VSCode. De asemenea, resursele necesare NLTK, cum ar fi stopwords și wordnet au fost descărcate pentru a sprijini prelucrarea textului.

2.2 Analiza setului de date

Statisticile descriptive ale setului de date, acesta având 27,783 de tweeturi sunt prezentate mai jos în figura.1, incluzând numărul total de notări, media, deviația standard, valorile minime și maxime și procentele. Și distribuția numerică pentru fiecare clasă este afișată în Figura.2.

Statistici descriptive:						
	count	hate_speech_count	offensive_language_count	neither_count	class	clasa
count	24783.000000	24783.000000	24783.000000	24783.000000	24783.000000	24783.000000
mean	3.243473	0.280515	2.413711	0.549247	1.110277	1.110277
std	0.883060	0.631851	1.399459	1.113299	0.462089	0.462089
min	3.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	3.000000	0.000000	2.000000	0.000000	1.000000	1.000000
50%	3.000000	0.000000	3.000000	0.000000	1.000000	1.000000
75%	3.000000	0.000000	3.000000	0.000000	1.000000	1.000000
max	9.000000	7.000000	9.000000	9.000000	2.000000	2.000000

Figura 2.1:Statisticile descriptive

```

Distribuția claselor:
1: 19190
2: 4163
0: 1430

```

Figura 2.2: Distribuția tweet-urilor pentru fiecare clasa.

În figura 2.3 este prezentă distribuția datelor în setul de antrenare și setul de testare:

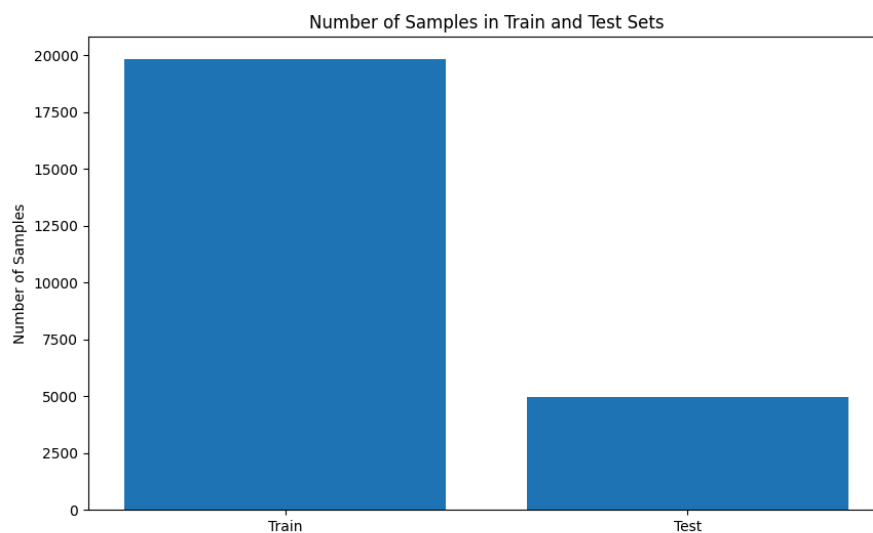


Figura 2.3:Distribuția datelor de train și test

Matricea de confuzie este utilă pentru evaluarea performanței modelului de clasificare. Aceasta reprezintă vizual predicțiile modelului comparativ cu valorile reale. (Figura 2.4)

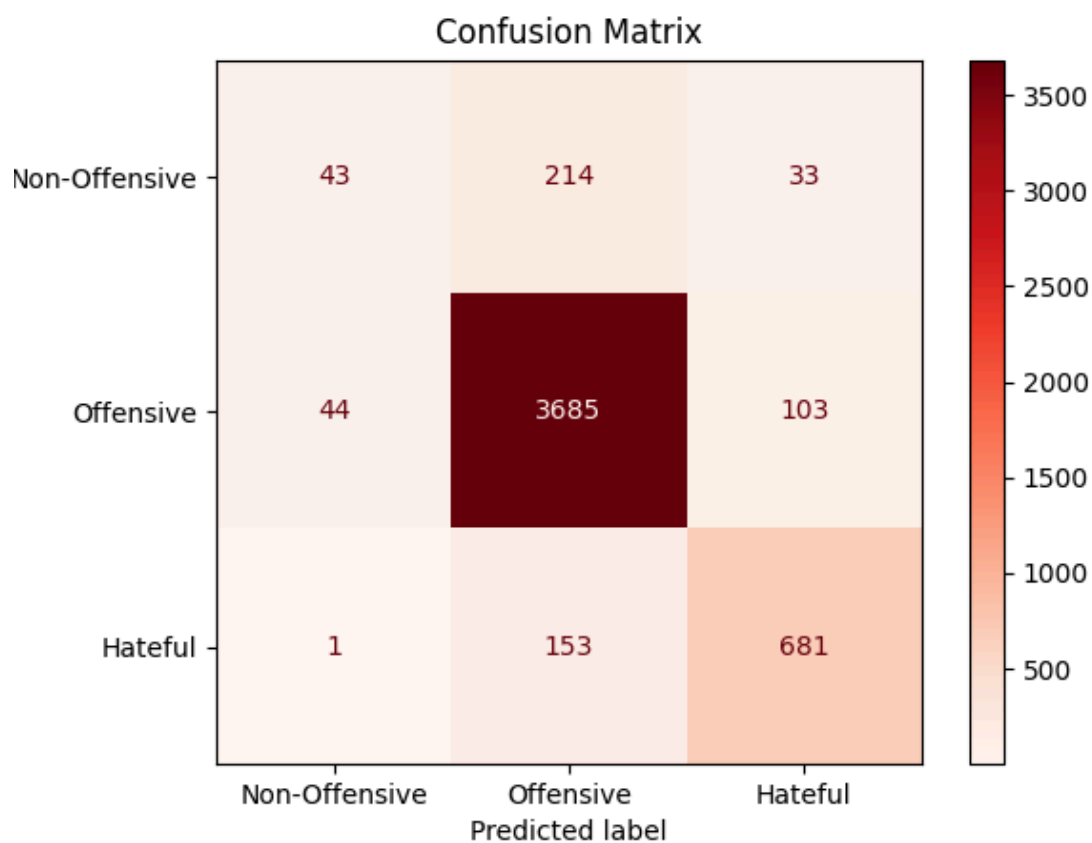


Figura 2.4: Matricea de confuzie

Analiza corelațiilor măsoară corelția dintre două variabile, indicând dacă și cât de mult se influențează reciproc. Valorile apropiate de 1 sugerează o corelație pozitivă puternică și -1 reprezintă o relație inversă. Analiza de corelație a setului de date train.csv este reprezentată în Figura 2.5.

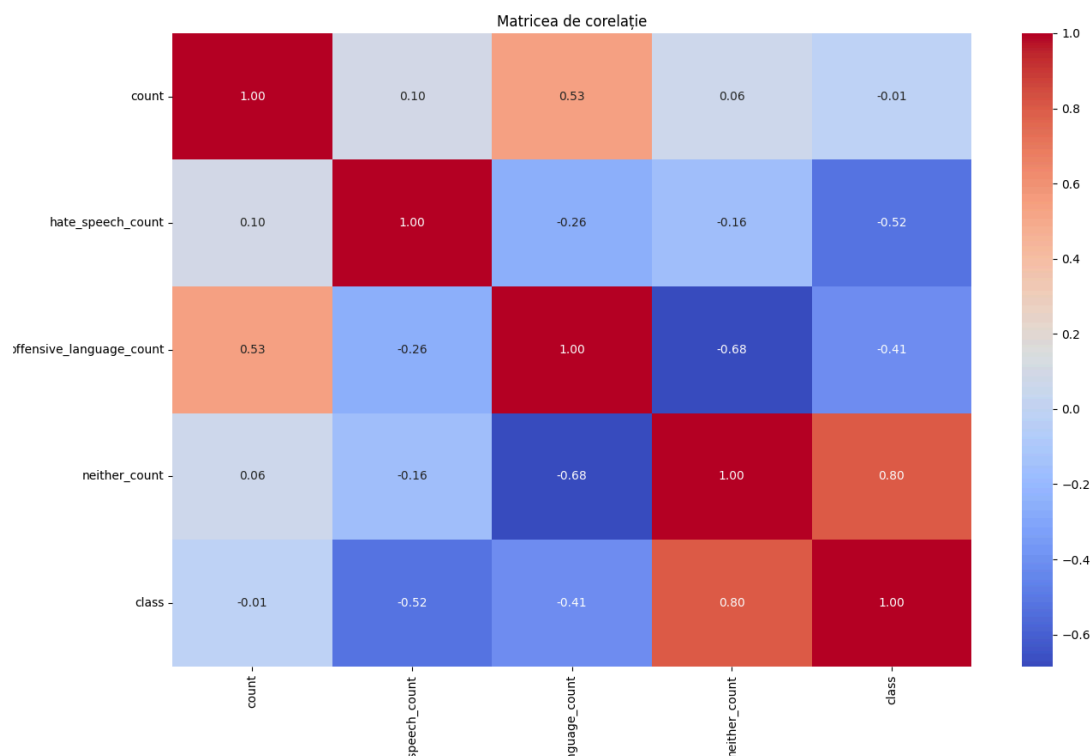


Figura 2.5: Matricea de corelație

Rezultate observate în toate graficele , datele din train.csv se axează mai mult pe limbajul ofensator. Așadar, ne așteptam ca acuratețea următoarelor experimente să nu funcționeze foarte bine pentru mesaje neutre și mesaje instigatoare la ura, dar o să funcționeze foarte bine pentru mesaje ofensatoare.

Modelul are o acuratețe de: 88.94%. Folosind formula:

$$\text{Acuratete} = \text{Nr. Predictii Corecte} / \text{Nr. total de predictii} \times 100;$$

In cod :

```
accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy * 100:.2f}%")
```

3 Preprocesarea setului de date

3.1 Metodele utilizate pentru preprocesarea setului de date

Eliminarea zgomotului din text este un pas esențial în preprocesarea datelor pentru modelele de procesare a limbajului natural, contribuind la îmbunătățirea acuratetii clasificării. Eliminarea manerelor de utilizator (@user) ajuta la reducerea informațiilor relevante, prevenind influența numelor specifice asupra analizei sentimentelor sau clasificării textului. Eliminarea link-urilor și etichetelor HTML reduce conținutul inutil, asigurând că modelul analizează doar textul esențial și nu elemente externe ne semnificative. Inlaturarea semnelor de punctuație și a caracterelor speciale ajuta la uniformizarea textului și previne confuziile generate de simboluri fără valoare semantica în contextul clasificării. Stop Words Urile sunt eliminate pentru a reduce dimensiunea textului procesat și a îmbunătăți performanța modelului, pastrand doar cuvintele relevante pentru interpretarea semantica.

Dimensionality reduction reprezintă un proces important în preprocesarea datelor de text, care are rolul de a reduce numărul de caracteristici într-un set de date, păstrând totuși informația relevantă. În cazul acestui cod, reducerea dimensionalității poate fi realizată folosind tehnici precum PCA (Principal Component Analysis) sau TruncatedSVD, aplicate pe vectorul de caracteristici obținut de la vectorizer. Aceasta ajută la accelerarea procesului de clasificare și la îmbunătățirea performanței modelului prin eliminarea zgomotului din datele de intrare. În cazul unor modele de clasificare text, reducerea dimensionalității poate contribui la îmbunătățirea acurateții, prin eliminarea redundanței dintre caracteristici. Un avantaj major al acestei tehnici este că, prin reducerea numărului de dimensiuni ale datelor de intrare, se poate evita overfitting ul, contribuind astfel la generalizarea mai bună a modelului pe date noi.

Librăriile folosite pentru preprocesare:

- Pandas: O bibliotecă esențială pentru manipularea și analiza datelor. Este utilizată pentru a citi fișier CSV și pentru a gestiona datele într-un format tabular (DataFrame).

-
- NLTK: O bibliotecă pentru procesarea limbajului natural. Este folosită pentru tokenizare, eliminarea cuvintelor de legătură și alte sarcini de prelucrare a textului.
 - re: O bibliotecă standard în Python pentru manipularea textului folosind expresii regulate. Este utilizată pentru a căuta și a înlocui modele în text, cum ar fi eliminarea URL-urilor și a mențiunilor de utilizator.
 - NumPy: O bibliotecă pentru manipularea eficientă a array-urilor și a calculului numeric. Este adesea utilizată împreună cu Pandas.

3.2 Rezultatele obținute și informații necesare pentru realizarea proiectului

Se importă librăriile necesare pentru preprocesarea datelor și antrenarea modelului de clasificare, cum ar fi nltk și sklearn. Datele sunt curățate prin eliminarea handle-urilor de pe rețelele sociale, URL-urilor, semnelor de punctuație și a cuvintelor de legătură folosind funcții dedicate. Setul de date este împărțit în seturi de antrenament și testare, iar textul este vectorizat folosind TfidfVectorizer pentru a-l transforma într-un format numeric. Modelul de regresie logistică este antrenat pe datele vectorizate și evaluat pe setul de testare prin acuratețe și raport de clasificare. La final, modelul și vectorizatorul sunt salvate pentru a fi utilizate ulterior în predicții pe date noi. Rezultate după preprocesarea textului este textul filtrat care asigură performanțe mai bune.

4 Modelarea sistemului

4.1 antrenareModel1.py – Antrenarea Modelului

Se utilizează mai multe librării esentiale:

- pandas – pentru citirea și manipularea datelor
- nltk – pentru preprocesarea textului
- sklearn – pentru modelarea și evaluarea performanțelor modelului
- matplotlib și seaborn – pentru vizualizarea datelor
- joblib – pentru salvarea și încărcarea modelului antrenat

Preprocesarea datelor:

- Citirea setului de date: Datele sunt încărcate într-un fișier CSV (train.csv).
- Curatarea textului: Se elimina elemente nedorite precum:
- Mentiuni de utilizatori (@user)
- URL-uri
- Punctuație excesivă
- HTML/unicode entities
- Stop Words (cuvinte fără semnificație semantică importantă)

Tokenizare și vectorizare:

- Se utilizează TfidfVectorizer pentru a transforma textul într-o reprezentare numerică

Antrenarea modelului:

- Se împarte setul de date în date de antrenare și date de testare
- Se antrenează un model de clasificare folosind Logistic Regression
- Se evaluează performanța modelului utilizând:

Acuratețea (accuracy_score)

Raportul de clasificare (classification_report)

Matricea de confuzie

Matricea de corelație

Salvarea modelului:

Modelul antrenat și vectorizatorul sunt salvate în fișiere `model_clasificare.joblib` și `vectorizator.joblib` pentru a fi utilizate ulterior

4.2 `predictie_manual.py` - Clasificarea unui tweet introdus manual

Acest script permite utilizatorului să introducă manual un tweet pentru a fi clasificat de model.

Pașii executați:

- Încărcarea modelului antrenat: Se încarcă modelul și vectorizatorul utilizând `joblib.load()`
- Preprocesarea textului introdus: Se aplică aceleași transformări ca în `antrenareModel1.py` (eliminarea mențiunilor, URL-urilor, stopwords etc.)
- Vectorizarea textului: Textul este transformat în vector numeric folosind `TfidfVectorizer`
- Predicția: Se utilizează `model.predict()` pentru a clasifica tweetul
- Afișarea rezultatului: Tweetul este clasificat în una dintre cele trei categorii: Non-Offensive, Offensive, Hateful.

4.3 `preluare_dateTw.py` - Colectarea de tweet-uri de pe Twitter

În acest script se folosește API-ul de la platforma X pentru a extrage Tweet-uri.

Etapele procesului:

- Citirea credențialelor: Credențialele sunt încărcate din `credentials.csv` pentru autentificare
- Configurarea API-ului Twitter: Se autentifică accesul utilizând cheile API
- Extragerea tweet-urilor: Se colectează tweet-uri în limba engleză folosind `tweepy.Cursor().Tweet`-urile sunt salvate într-un fișier CSV (`date_noi_tweeter.csv`)

4.4 predictie.py - Clasificarea unui set de tweet-uri

Datele colectate de pe platforma X vor ajunge în csv-ul date noi tweeter.csv de unde cu ajutorul script-ului predictie.py vom categoriza fiecare tweet.

Pașii executați:

- Încărcarea datelor din date noi tweeter.csv
- Preprocesarea textului:Aplicarea aceleiasi metode de curatare ca in antrenareModel1.py
- Vectorizarea textului:Transformarea textului in reprezentare numerica
- Clasificarea tweet-urilor:Aplicarea modelului antrenat pentru predictie
- Afișarea rezultatului:Tweet-urile sunt etichetate și aflate alături de predicțiile respective

5 Concluzii

5.1 Rezultate obținute

Prin testări succesive ale aplicației în diverse condiții și utilizând scriptul predicție manual.pe, care permite introducerea manuala a propozitiilor, am observat ca aceasta are o performanta foarte buna în detectarea cuvintelor obscene în limba engleza. Cu toate acestea, cea mai mare provocare intampinata de model consta în diferențierea dintre limbajul instigator și cel ofensator. În unele cazuri, propozitii simple sunt clasificate eronat ca fiind ofensatoare sau instigatoare. O soluție eficienta pentru îmbunătățirea acuratetii modelului este antrenarea acestuia pe un set de date mai extins și mai diversificat, care sa includă un număr mai mare de exemple de limbaj neutru și instigator, contribuind astfel la o clasificare mai precisă și mai nuanțată.

5.2 Direcții de dezvoltare

Pentru a imbunatati performanta sistemului de clasificare a tweet-urilor, se pot explora mai multe directii. Se poate extinde setul de date prin colectarea de tweet-uri suplimentare din diverse surse și limbi. Un model mai complex, precum RNN sau BERT, ar permite o analiza mai detaliată a textului. Combinarea analizei textuale cu imagini și metadata (detecție multimodala) ar duce la o clasificare mai precisă. Un sistem de feedback ar ajuta la imbunatatirea continua a modelului, iar optimizarea performanței ar reduce timpul de procesare. Introducerea unor categorii suplimentare, ca sarcasm sau dezinformare, ar oferi o analiza mai nuanțată. Aceste îmbunătățiri ar crește acuratețea și fiabilitatea sistemului.

6 Bibliografie

Api Twitter: <https://docs.x.com/x-api/tools-and-libraries/overview>

Antrenarea modelului:

<https://www.kaggle.com/code/saramoha88/hate-speech-and-offensive-language-detection>

Libraria NLTK: <https://www.nltk.org/>

Librarie joblib: <https://joblib.readthedocs.io/en/stable/>