

Documentație proiect

“Misogynistic Italian Tweets Detector”

Scopul proiectului este acela de a realiza un algoritm de machine learning pentru detectarea tweeturilor misogine în italiană. Etapele realizării acestui proiect sunt descrise în pași următori.

1.Setul de date:

- Setul de date este format din 5000 de tweet-uri în italiană care reprezintă partea de antrenare a modelului și 1000 de tweet-uri pentru partea de testare a modelului antrenat.
- Datele de antrenare sunt transmise în format CSV.
- Aspecte ale datelor ce trebuie menționate sunt următoarele:
 - Tweeturile conțin “**twitter handle**” sub forma “@ + text”, **hyperlinkuri** de forma “https:// + text” și **emoji-uri în forma hexazecimală**.
 - Datele de antrenare sunt așezate într-un format necorespunzător întrucât acestea au numai texte cu label-ul “1” la început, urmate de un număr foarte mare de texte cu label-ul “0” după care urmează o alternanță.

2.Organizarea proiectului:

- Proiectul este împărțit în patru scripturi cu roluri distincte:
 - “**text_preprocessing.py**” cu scopul de a aplica toate modificările de preprocesare.
 - “**data_partition.py**” cu scopul de a procesa datele din corpus train
 - “**prediction_to_csv.py**” cu scopul de a transforma predicțiile modelelor într-un format potrivit pentru încărcare pe Kaggle.
 - “**main.py**” este scriptul principal care se folosește de restul scripturilor pentru a face predicții.

3.Preprocesarea datelor:

- Pentru preprocesarea datelor ne folosim de librăriile `nltk.corpus`, `nltk.stem.snowball`, `nltk.tokenize`, `re`, `preprocessor` și `string`.
- Pentru a putea avea un text cât mai relevant trebuie ca acesta să fie curățat de numele de utilizatori, de linkuri, de cuvintele de legătură și de punctuație.
- Ne creem variabilele **tokenizer (TweetTokenizer)**, **stemmer (SnowballStemmer("italian"))**, **stop_words** pentru a prelucra textul.
- Folosind funcția **"clean"** din biblioteca **preprocessor**, pe fiecare propoziție din corpus o să eliminăm **emojiurile** în forma hexazecimală. În aceeași manieră folosim **"re.sub(r"http\S+", "", sentence)"** pentru a elimina link-urile de forma **"http"** urmate de un string (cu ajutorul wild cardului `\S+`).
- Cu ajutorul `TweetTokenizer`-ului o să transformăm tweet-urile în tokenuri și scăpăm și de numele de utilizatori specificând **"strip_handles=True"**.
- Parcurgând lista de tokenuri putem elimina cuvintele comune, și semnele de punctuație folosind o listă nouă de tokenuri la care adăugăm cuvintele care nu se află în lista de semne de punctuație cât nici în lista de stop words.
- Înainte de a le adăuga la lista de tokenuri procesate o să le trecem și prin `stemmer` pentru a le reduce la forma cea mai simplă.
- La finalul acestei proceduri propozițiile sunt în cea mai simplă formă astfel obținând date cât mai relevante.

4.Procesarea datelor:

- Procesarea datelor se face cu ajutorul scriptului **"data_partition.py"**.
- Am definit două funcții **"to_TFIDF_bow()"**. Această funcție este folosită pentru a transforma cuvintele în bag of words folosind algoritmul **TFIDF** format din **Term Frequency + Inverse Data Frequency**.
- **Term Frequency** reprezintă de câte ori a apărut un cuvânt în cadrul unei propoziții împărțit la numărul total de cuvinte dintr-o propoziție. Fiecare document(propoziție) are un **TF** propriu.
- **Inverse Data Frequency** reprezintă logaritm din numărul total de documente împărțit la numărul de documente care conțin un anumit cuvânt. Logaritmul se folosește pentru a ameliora importanța cuvintelor cu o frecvență foarte mare.
- **TFIDF** reprezintă produsul dintre **TF** și **IDF**.
- Am ales ca parametrul **"max_features"** să fie egal cu 1000 prin observații din aplicarea mai multor valori

5. Modelele folosite:

a. Multinomial Naive Bayes:

- i. Am importat din librăria `sklearn.naive_bayes` modelul `MultinomialNB` pentru a antrena datele. Am observat că atunci când există un număr mai mic de `"max_features"`(1000) în vectorizatorul `TFIDF`, modelul `MultinomialNB()` funcționează mai bine.

max_features	acuratetea medie
2500	0,8122
2000	0,8110
1500	0,8290
1000	0,8372

- ii. Pentru acest model hiperparametrii au fost cei default.
- iii. În toate cazurile datele antrenate și cele de test au fost prelucrate cu metodele prezentate în secțiunea **"Preprocesarea datelor"** ⇔ scoatem stop words, punctuație, numele de utilizatori, linkurile web, emojiurile și le trecem printr-un stemmer.
- iv. Am antrenat și testat datele în manieră **10 fold cross validation**, datele au fost împărțite în 10 intervale egale. Folosind această metodă împărțim datele în 10 intervale și acestea sunt pe rând atât date de antrenare cât și date de validare. La final toate rezultatele sunt adunate pentru a obține matricea de confuzie

```
Matricea de confuzie este:  
[[2347  472]  
 [ 316 1865]]
```

- v. Folosind clasificatorul de mai sus cu parametrii de la 10 fold cross validation obținem următoarele scoruri și acuratețe medie.

```
Iteratia numarul 1  
Acuratetea modelului este: 0.836  
  
Iteratia numarul 2  
Acuratetea modelului este: 0.85  
  
Iteratia numarul 3  
Acuratetea modelului este: 0.836  
  
Iteratia numarul 4  
Acuratetea modelului este: 0.848  
  
Iteratia numarul 5  
Acuratetea modelului este: 0.85
```

```
Iteratia numarul 6  
Acuratetea modelului este: 0.858  
  
Iteratia numarul 7  
Acuratetea modelului este: 0.83  
  
Iteratia numarul 8  
Acuratetea modelului este: 0.852  
  
Iteratia numarul 9  
Acuratetea modelului este: 0.808  
  
Iteratia numarul 10  
Acuratetea modelului este: 0.856  
  
Media acuratetii pe model: 0.8424
```

- vi. Deși pe datele de test local acuratețea medie este **0.8424** pe Kaggle pe leaderboard-ul public acuratețea este **0.76842** și pe cel privat este **0.77935**.

- vii. Antrenarea datelor a durat **2.5 secunde**. Din librăria `time` am folosit funcția `time` pentru a afla timpul de antrenare.

b. SVM/SVC:

- i. Importăm din librăria `sklearn.svc` modelul de antrenare de date **SVM**. La fel ca la modelul de mai devreme un număr mic de “**max_features**”(1500) rezultă într-o acuratețe mai mare de prezicere a modelului pe datele de validare.

max_features	acuratetea medie
2500	0,8596
2000	0,8528
1500	0,8606
1000	0,8538

- ii. Datele pe acest model au fost prelucrate precum cele de dinainte.
iii. Hiperparametrii modelului sunt cei default.
iv. Matricea de confuzie pentru datele de validare ale acestui model este:

```
Matricea de confuzie este:  
[[2264 298]  
 [ 399 2039]]
```

- v. Datele de validare au da o **acuratețe medie** de **~0.8606** însă pe Kaggle pe **leaderboard-ul privat** acuratețea este **0.74446** și pe cel **public** este de **0.74042**
vi. Antrenarea datelor a durat **0.3 secunde**.