

Documentație Proiect ML

Pentru acest proiect am plecat de la codul făcut în timpul laboratoarelor, adică în prima fază am separat fiecare cuvânt din fiecare text și am aflat numărul de apariții în textul respectiv. Textul pe care îl avem de analizat se afla într-un document de tip .csv, iar pentru a citi acest format folosim librăria “pandas” în modul următor: am creat o variabilă “corpus” în care se stochează datele din coloana text, iar apoi cu ajutorul funcției “preprocess” am creat un dicționar cu cele mai folosite cuvinte care apar în .csv.

Funcția “preprocess” este ultima parte din faza de data_cleaning a textului. Pentru data_cleaning am folosit librăriile “string”, “re”, “nltk”, prin care mi am transformat toate cuvintele în cuvinte cu litere mici, apoi am șters din texte datele de care nu aveam nevoie, de exemplu punctuația, userii, URL-urile, hashtag-urile, numerele, cuvintele cele mai des folosite în limba italiană (inițial pentru acest aspect am folosit comanda “stopword = stopwords.words('italian')”), dar am observat că îmi mai apar niște cuvinte care nu erau în adăugate în librăria de la “nltk”, am încercat să le adaug manual, dar tot îmi apăreau așa că am decis să caut pe net aceste cuvinte și să le adaug într-un vector), apoi am folosit un tokenizator, pentru a îmi delimita cuvintele între ele și în ultima fază am făcut ca acele cuvinte scrise greșit să fie identificate cu cuvintele scrise corect.

După partea de data_cleaning am folosit funcția “Counter” din librăria “counter” cu rolul de a-mi afișa frecvența de apariție a cuvintelor, cu aceste informații creez două matrici, una care reprezintă cele mai întâlnite cuvinte, iar a doua care reprezintă un bag of words pentru fiecare tweet. Rezultatele obținute care se afișează în matrici le folosesc atât pe textele de antrenare cât și pe textele de testare, astfel voi folosi un clasificator ca să pot detecta tweeturile misogine.

Pentru partea de clasificatori am folosi clasificatorul Complement Naive Bayes cu care am obținut un scor pe testele finale de 0.73860 și clasificatorul Gaussian Naive Bayes cu care am obținut un scor de 0.72424. Acești doi clasificatori i-am introdus în codul meu cu ajutorul librăriei “sklearn.naive_bayes”.

Complement Naive Bayes:

- în loc să calculăm probabilitatea unui element aparținând unei anumite clase, calculăm probabilitatea ca elementul să aparțină tuturor claselor.

O prezentare step-by-step a algoritmului:

- I. Pentru fiecare clasă calculez probabilitatea ca instanța dată să nu îi aparțină.
- II. După ce am făcut toate calculele pentru toate clasele, verific toate valorile calculate și selectez cea mai mică valoare.
- III. Selectez cea mai mică valoare, deoarece are cea mai mică probabilitate de a nu fi o clasă particulară, ceea ce implică că are cea mai mare probabilitate de a aparține clasei respective.

Nu selectez cea mai mare valoare deoarece calculez complementul probabilității, iar cea mai mare valoare este cel mai probabil să fie clasa căreia îi aparține elementul.

Gaussian Naive Bayes:

- Valorile continue asociate cu fiecare caracteristică sunt presupuse să fie distribuite în conformitate cu o distribuție Gaussiană. O distribuție Gaussiană se mai numește distribuție Normală (distribuția de probabilitate continuă). Când plotează are forma unei curbe în forma de clopot care este simetrică cu media valorilor caracteristici.