

UNIVERSITATEA DIN BUCURESTI

FACULTATEA DE MATEMATICA SI INFORMATICA

DOCUMENTATIE

PROIECT INTELIGENTA ARTIFICIALA

„Automatic Misogyny Identification”

Coordonatori:

Conf.dr. Bogdan Alexe

Asist.(det.) drd. Sergiu Nisioi

Student:

Damian-Luca Rares

Gr. 352

Pre-procesarea datelor

În vederea obținerii unor predicții cât mai precise, datele inițiale trebuie să treacă prin mai multe etape de pre-procesare pentru a se ajunge la setul de antrenare. (1)

Pășii următori au fost preluați de la utilizatorul Kaggle Rageeni Sah (2):

Inițial am încărcat datele din `train.csv`, `test.csv`, și am plotat numărul de intrări cu eticheta 0 și 1, pentru a observa dacă datele sunt balansate. Au fost preluate coloanele "id" și "text" și încărcate în dataframe-uri pandas. După aceasta am generat trei imagini word cloud, pentru a vedea cele mai comune cuvinte din toate textele, textele cu eticheta 0 și textele cu eticheta 1, pentru a le elimina ulterior.

Urmărește partea de data cleansing pentru înlăturarea părților irelevante. (3). Au fost definite funcții pentru fiecare pas și apoi aplicate coloanelor de text respective utilizând funcții lambda pentru a crea o nouă coloană. Aceste funcții au rolul de a:

- Înlăturarea punctuației- `def remove_punct(text):`
- Tokenizarea textului- `def tokenization(text):`
- Înlăturarea cuvintelor de oprire- `def remove_stopwords(text):`
- Stemming (aducerea la rădăcină)- `def stemming(text):`

Corpus-ul, reprezentând textul trecut prin toți pașii de cleansing și având forma de array de cuvinte, este utilizat pentru a obține vocabularul. Urmărește obținerea dicționarului `wd2idx` (atribuirea unui index primelor `n` cele mai comune cuvinte) și formarea modelului bag of words (textul fiind reprezentat ca un multiset alcătuit din propriile cuvinte, păstrându-se multiplicitatea) (4)

Submisia nr. 1 (MultinomialNB)

Metodele Naive Bayes sunt un set de algoritmi de invatare supervizati bazati pe aplicarea teoremei lui Bayes cu presupunerea de independenta conditionala intre fiecare pereche de caracteristici fiind data valoarea variabilei clasa. (5)

Clasificarea tweet-urilor:

- X_1, X_2, \dots, X_{100} sunt aparitiile celor 100 cele mai des intalnite cuvinte
- $c \in \{0, 1\}$ (clasifica daca textul este misogin sau nu)

Teorema lui bayes: $P(c/X) = P(X/c) \times P(c) / P(X)$

$P(c/X)$ - probabilitatea sa avem textul din clasa c dandu-se reprezentarea bag of words X

$P(X/c)$ - probabilitatea sa avem reprezentarea bag of words X stiind ca textul apartine clasei c

$P(c)$ – probabilitatea a-priori

$P(X)$ – probabilitatea sa observam reprezentarea bag of words X

$P(c = 0/X) = P(X/c = 0) \times P(c = 0) / P(X)$

$P(c = 1/X) = P(X/c = 1) \times P(c = 1) / P(X)$

Intrucat avem numitori egali, putem renunta la acestia

$P(c = 0/X) \text{ d.p. } P(X/c = 0) \times P(c = 0)$

$P(c = 1/X) \text{ d.p. } P(X/c = 1) \times P(c = 1)$

Regula de clasificare:

$$c^* = \operatorname{argmax} P(X | c = i) \times P(c = i)$$

Aceasta alege clasa care maximizeaza numaratorul

Probabilitatea a-priori $P(c = i)$ se calculeaza numarand cate exemple din multimea de antrenare au clasa i.

Pobabilitatea likelihood $P(X/c = i)$ se calculeaza considerand caracteristicile independente din cauza numarului limitat de date antrenare.

$$P(X/c = i) = P(X_1=x_1, X_2=x_2, \dots, X_{100}=x_{100} / c = i) = \\ = P(X_1 = x_1 / c = i) \times P(X_2 = x_2 / c = i) \times \dots \times P(X_{100} = x_{100} / c = i)$$

Regula de clasificare pentru naive Bayes:

$$c^* = \operatorname{argmax}(\prod P(X_j = x_j | c = i)) \times P(c = i)_{(6)}$$

Am utilizat clasificatorul MultinomialNB din biblioteca scikit-learn, o varianta des folosita pentru clasificarea textelor. Acesta implementeaza algoritmul naive Bayes pentru date distribuite multinomial. Repartizarea este parametrizata de vectori $\theta_c = (\theta_{c1}, \dots, \theta_{cn})$ pentru fiecare clasa c , unde n este dimensiunea vocabularului si θ_{ci} este probabilitatea $P(x_i/y)$ a caracteristicii i de a apare intr-o mostra apartinand clasei c . (7)

Pentru antrenare am iterat cu ajutorul unui for valori pentru alfa (utilizat pentru caracteristici absente in datele de invatare). Clasificatorul este initializat la fiecare iteratie cu o noua valoare pentru alfa si este initializat un vector gol pentru a ajuta la calculul scorului f1 mediu. Apoi au loc 1000 de antrenari, preziceri si calcule de scor f1 10-fold cross-validate. Timpul de antrenare este de 490 secunde, adica 8 min 10 s . Scorul f1 obtinut pe Kaggle pe 40% din date este 0.72351.

Date 10 fold cross-validation si matrice de confuzie:

```
10 fold cv si matrice de confuzie
acc: 84.5 f1: 0.8319407426349429
Matricea de confuzie:
1919   357
 418  2306
```

Submisia nr. 2 (KNN)

Algoritmul de invatare supervizata a vecinilor cei mai apropiati pe date cu etichete discrete rezolva probleme de clasificare. Principiul de functionare este gasirea unui numar predefinit de probe de antrenament cat mai apropiate de punctul nou si prezicerea etichetei pe baza acestora. Distanța poate fi orice metrica, iar in general este folosita distanta Euclidiană. Acest clasificator este bazat pe instanta, adica nu construiește un model intern generalizat, ci retine instantele datelor de antrenare. (8)

Textul a fost tokenizat utilizand `nlk.WordPunctTokenizer().tokenize()`. Datele de train au fost impartite in antrenare si validare cu un split random de 30%. Pentru antrenare am iterat cu ajutorul unui for valori pentru `n_neighbours`. Clasificatorul este initializat la fiecare iteratie cu un numar de vecini nou, apoi are loc antrenarea, prezicerea si calcularea acuratetii si scorului f1. Aceste ultime doua valori calculate, impreuna cu numarul vecinilor din for sunt adaugate in vectori pentru a putea fi folosite la realizarea unor grafice. Acestea ajuta la descoperirea valorilor celor mai avantajoase pentru numarul de vecini. Antrenarea dureaza 6.82 secunde. Scorul f1 obtinut pe Kaggle pe 40% din date este 0.70613.

Bibliografie

1. https://en.wikipedia.org/wiki/Data_pre-processing
2. <https://www.kaggle.com/ragnisah/text-data-cleaning-tweets-analysis>)
3. https://en.wikipedia.org/wiki/Data_cleansing
4. https://en.wikipedia.org/wiki/Bag-of-words_model
5. https://scikit-learn.org/stable/modules/naive_bayes.html#naive-bayes
6. Cursul 4 de IA - Conf.dr. Bogdan Alexe
7. https://scikit-learn.org/stable/modules/naive_bayes.html#naive-bayes
8. <https://scikit-learn.org/stable/modules/neighbors.html#neighbors>