

Tema 3 – Sumarizarea documentelor
Inteligența Artificială
Nicolescu Radu-Catalin, 343C4

1. Detalii de implementare

Clasa Document va fi folosită pentru a descrie un articol, având 3 variabile: content – conținutul articolului, summary – rezumatul articolului, dtype – tipul articolului (sport, politics, entertainment, tech, business).

Clasa Editor este folosită pentru a prelucra conținutul articolelor, aceasta ocupându-se cu împărțirea unui articol în propoziții, împărțirea unui articol în cuvinte, eliminarea cuvintelor de legătură (stop words) și lematizarea cuvintelor. Aceasta are drept membru o listă de cuvinte de legătură, care este populată prin citirea fișierului stop_words dat în arhiva temei.

Clasa Loader se ocupă cu încărcarea documentelor în memorie și salvarea acestora într-un document în format .csv, care va fi ulterior folosit în clasificare și sumarizare. Aceasta are ca variabile 5 liste de documente, câte una pentru fiecare tip de articole. Clasa poate citi toate articolele și să le salveze local, dar și să le salveze într-un document csv.

Modulul runner.py se ocupă cu sumarizarea și clasificarea documentelor, folosindu-se de documentul .csv generat de Loader. Documentul .csv conține 3 coloane și un număr de linii egal cu numărul de stiri existente. Cele 3 coloane sunt content, summary și type.

Primele două (content, summary) vor fi folosite în algoritmul Naïve Bayes de sumarizare, în timp ce tuplul (content, type) va fi folosit în algoritmul Naïve Bayes de clasificare.

Pentru clasificare, se calculează probabilitățile ca articolul să aparțină fiecărei dintre cele 5 clase, la final alegeându-se clasa cu probabilitatea maximă de apartenență.

Pentru sumarizare, pentru fiecare articol, se calculează probabilitățile ca fiecare propoziție să aparțină sumării, cât și probabilitatea medie a acestora. Propozițiile care au probabilitățile peste această medie vor fi incluse în sumarizarea aferentă articolului.

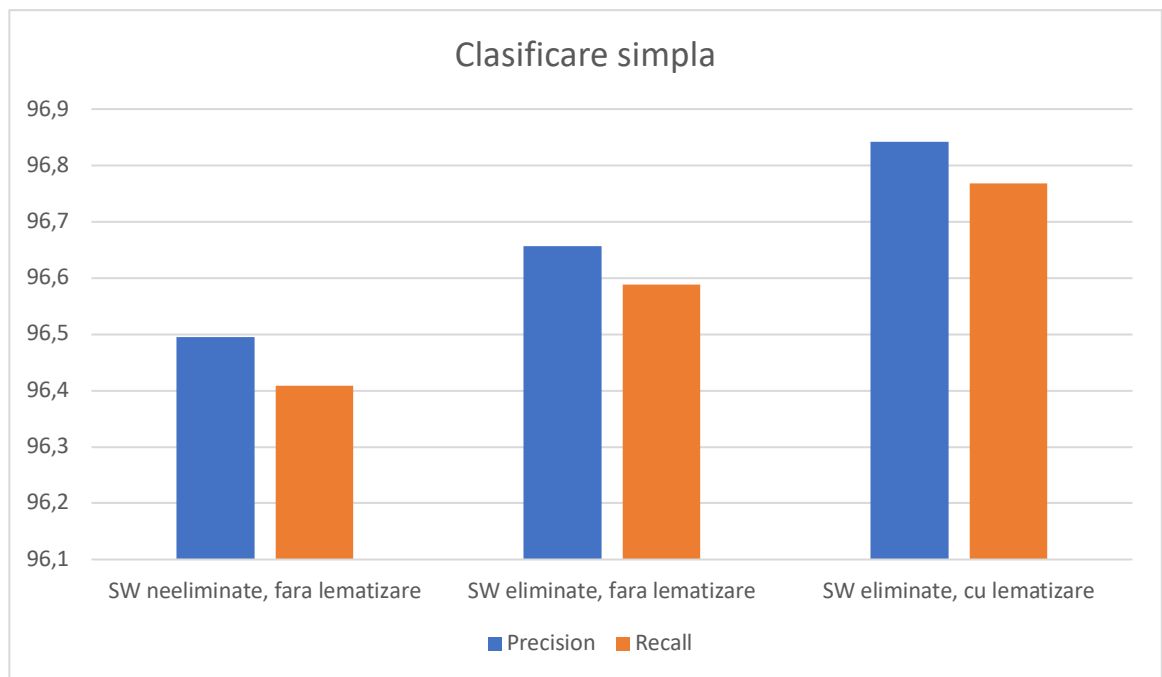
La finalul algoritmului de clasificare, se calculează matricea de confuzie, precizia și recall-ul algoritmului pentru 3 instanțe ale problemei: fără eliminare de cuvinte de legătură și fără lematizare, cu eliminare de cuvinte de legătură și fără lematizare, cu eliminare de cuvinte de legătură și cu lematizare.

La finalul algoritmului de sumarizare, se calculează precizia și recall-ul pentru 3 instanțe ale problemei: fără eliminare de cuvinte de legătură și fără lematizare, cu eliminare de cuvinte de legătură și fără lematizare, cu eliminare de cuvinte de legătură și cu lematizare, în 2 cazuri: folosind unigrame și folosind bigrame.

Pentru realizarea bonusului, am impartit setul de date in 5 parti egale, si am folosit pe rand cate o parte pentru testare, iar restul de 4 parti pentru antrenare. La final, am calculat media si deviatia standard pentru precizie si recall, pentru toate cazurile mentionate mai sus, atat pentru clasificare, cat si pentru sumarizare.

2. Rezultate experimentale

Clasificare simpla



Matricile de confuzie:

a) SW neeliminate, fara lematizare

	business	entertainment	politics	sport	tech
business	110	0	3	0	4
entertainment	1	96	2	0	3
politics	4	0	97	0	0
sport	0	0	0	134	0
tech	0	0	3	0	100

b) SW eliminate, fara lematizare

	business	entertainment	politics	sport	tech
business	110	0	3	0	4
entertainment	1	97	1	0	3
politics	4	0	97	0	0
sport	0	0	0	134	0
tech	0	0	3	0	100

c) SW eliminate, cu lematizare

	business	entertainment	politics	sport	tech
business	110	0	3	0	4
entertainment	1	96	2	0	3
politics	4	0	97	0	0
sport	0	0	0	134	0
tech	0	0	1	0	102

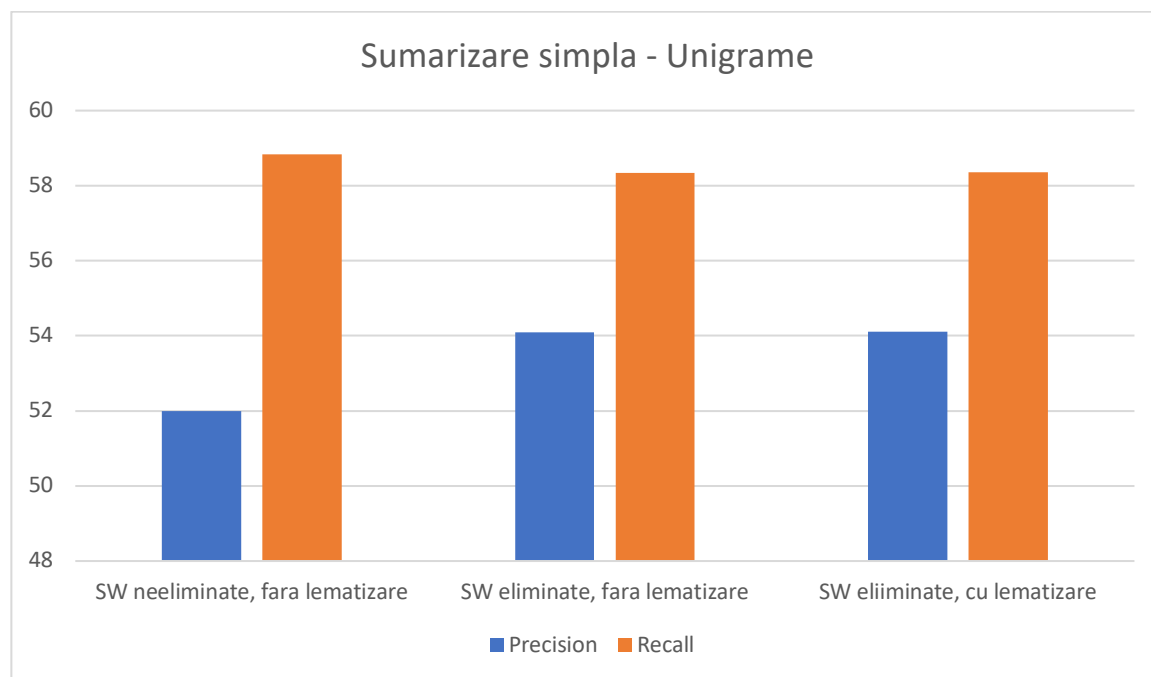
Comparand matricile de confuzie pentru cele trei cazuri, este usor de observat imbunatatirile de rezultate de la o implementare la alta: numarul de elemente de pe diagonala principala (true positives) creste, ceea ce duce la cresterea preciziei si recall-ului de la o implementare la alta.

Eliminarea cuvintelor de legatura are ca efect eliminarea din calcule a cuvintelor care nu au informatie relevanta. (nu contribuie la clasificarea catre o anume clasa)

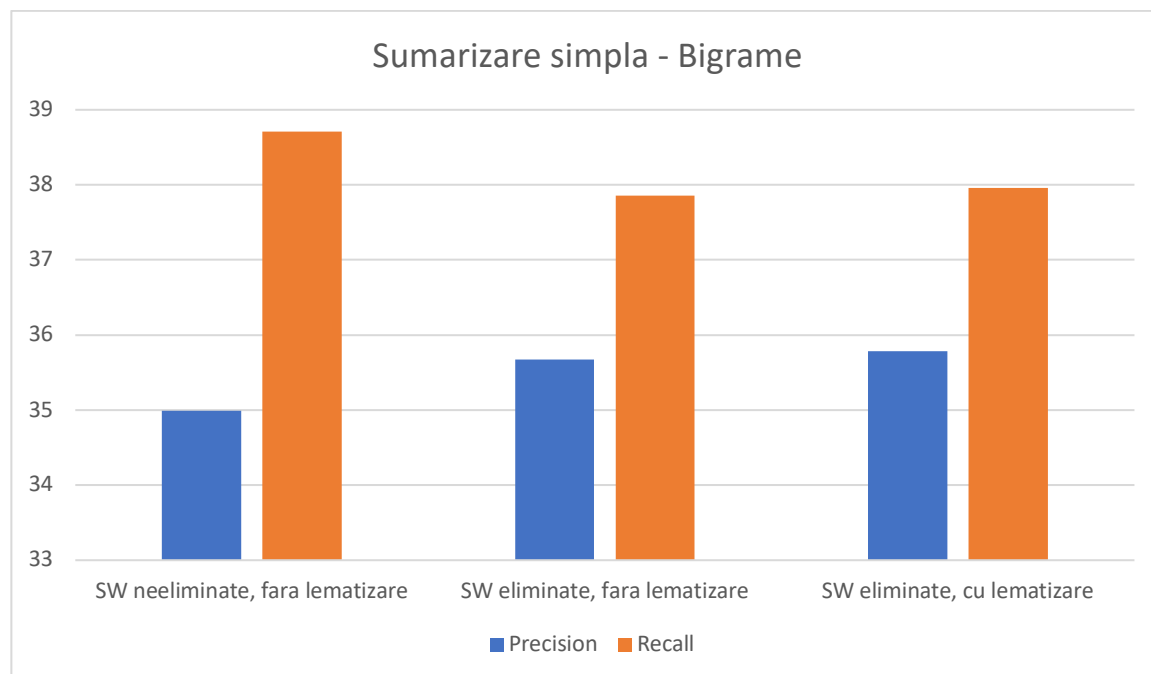
Lematizarea are ca efect cresterea numarului de cuvinte care au informatie relevanta, deoarece sunt eliminate sufixele, pluralurile. Astfel, un cuvânt care apare sub mai multe forme va aparea, dupa lematizare, într-o singura forma.

Sumarizare simpla

Unigrame



Bigrame



Performantele sumarizarii folosind unigrame sunt net superioare sumarizarii folosind bigrame, datorita faptului ca unigramele (formate dintr-un cuvant) au mult mai multe aparitii in texte decat bigramele (formate din doua cuvinte vecine), ceea ce duce la cresterea relevantei celor dintai in sensul includerii unei propozitii intr-un sumar.

Clasificare prin 5-fold cross-validation

	PRECISION		RECALL	
	Mean %	Standard Deviation %	Mean %	Standard Deviation %
SW neeliminate, fara lematizare	97,304	0,004	97,213	0,004
SW eliminate fara lematizare	97,605	0,004	97,528	0,004
SW eliminate cu lematizare	97,696	0,004	97,617	0,004

Sumarizare prin 5-fold cross-validation**UNIGRAME**

	PRECISION		RECALL	
	Mean %	Standard Deviation %	Mean %	Standard Deviation %
SW neeliminate, fara lematizare	51,469	0,005	58,920	0,004
SW eliminate fara lematizare	53,121	0,005	58,412	0,002
SW eliminate cu lematizare	53,329	0,006	58,067	0,002

BIGRAME

	PRECISION		RECALL	
	Mean %	Standard Deviation %	Mean %	Standard Deviation %
SW neeliminate, fara lematizare	34,601	0,007	38,917	0,007
SW eliminate fara lematizare	35,262	0,006	38,077	0,004
SW eliminate cu lematizare	35,274	0,006	38,099	0,003

Rezultatele preciziei si recall-ului obtinute in urma 5-fold cross-validation sunt foarte apropiate de rezultatele obtinute in urma impartirii clasice a setului de date in 75% date de antrenare si 25% date de testare, fiind diferite de sub 1% intre acestea.