

Tema 3 IA 2018

Sumarizare de texte

Vlad Bogolin

Scopul temei

Scopul temei este acela de familiarizare cu tehnici de baza din prelucrarea limbajului natural. Prelucrarea limbajului natural este una dintre principalele direcții de cercetare în momentul actual. Sumarizarea textelor consta în crearea unui text scurt și concis care sa reprezinte rezumatul unui document mai lung. Asadar, sumarizarea automata este o problema din ce in ce mai importanta datorita numarului mare de texte (de exemplu, stiri) ce se gasesc online. Metodele de sumarizarea automata ajuta la descoperirea informatiei relevante.

Enunt

Dandu-se un articol de stiri (titlu + text), sa se produca un rezumat al acestuia care sa fie cat mai relevant. Algoritmul va avea urmatoorii pasi:

1. Preprocesarea articolului - eliminarea informatiei nerelevante (cuvinte de stop, numere, etc). Pentru acest pas se gasesc online liste ce contin cuvintele de stop.
2. Definire vocabular - pentru a limita numarul de cuvinte folositi lematizare/stemming si puteti ignora cuvintele ce apar de putine ori. De asemenea, puteti lua in calcul doar substantivele (vezi cerinta 3).
3. Calculare Term Frequency (TF) pentru fiecare cuvant din vocabular
4. Calculare TF-IDF (Term Frequency - Inverse Document Frequency) pentru fiecare cuvant din vocabular
5. Calculare scor pentru fiecare propozitie din articolul de interes
6. Creare rezumat folosind propozitiile cu cel mai mare scor

Cerinta 1 [2p] - Curatare text

Eliminare cuvinte de stop, semne de punctuatie, eventuale spatii suplimentare. Aceasta cerinta se va implementa intr-o functie separata care primeste ca parametru textul unui articol si intoarce textul preprocesat.

Cerinta 2 [3p] - Calculare TF-IDF

Calculare TF-IDF pe baza textelor din setul de date de antrenare pentru fiecare cuvânt din vocabular. Aveți grijă să preprocesati textul înainte (lematizare/stemming - pentru lematizare/stemming puteți folosi orice resursă, etc).

Cerinta 3 [3p] - Calculare scor per propozitie

Folosind valorile TF-IDF calculate până acum, se va calcula un scor pentru fiecare propoziție. Pentru împartirea în propoziții puteți folosi nltk sau orice altceva. Scorul pentru o propoziție va fi dat de suma TF-IDF-ului pentru fiecare cuvânt. Totuși, pentru a avea un scor cât mai bun se vor parcurge următorii pași:

1. Se va calcula scorul doar pentru substantive. Pentru a identifica părțile de vorbire puteți folosi orice resurse (nltk, etc.). Adică scorul propoziției va fi afectat doar de scorul substantivelor continute. Normalizați scorul propoziției.
2. Calculare similaritate cu titlul. Mai exact, se va adăuga o valoare adițională dacă propoziția are cuvinte ce apar în titlu. Această valoare este egală cu numărul de cuvinte din propoziție care se regăsesc în titlu împartit la numărul total de cuvinte din titlu. La final, această valoare se va pondera cu o constantă care îi definește importanța (de exemplu 0.1) și se va aduna la scorul calculat la pasul 1.
3. La final se va pondera fiecare propoziție cu o pondere (între 0 și 1) corespunzătoare cu poziția ei în text. De exemplu, pentru un text cu 10 propoziții, ponderea pentru cea de-a nouă va fi 0.9. Acest lucru se întâmplă deoarece se tinde ca propozițiile cele mai importante să fie spre finalul articolului (de exemplu concluzia). Definiți propria metrică de ponderare a propoziției în funcție de locația în document.
4. Se vor returna primele 3 propoziții cu cel mai mare scor.

Cerinta 4 [2p] - Evaluare rezultate

Pentru evaluare se va folosi ground truth-ul pentru fiecare propoziție din dataset. Având în vedere că sumarizarea nu folosește informația de ground truth, evaluare și învățarea TF-IDF-ului se pot face pe același dataset.

Pentru evaluare se vor folosi două metrici: $BLEU@n$ și $ROUGE@n$. În general metrica BLEU măsura precizia, iar ROUGE recall-ul. Ele sunt definite după cum urmează:

- $BLEU@n = \frac{nr\ n\text{-grame comune între text și ground truth}}{nr\ n\text{-grame în text}}$
- $ROUGE@n = \frac{nr\ n\text{-grame comune între text și ground truth}}{nr\ n\text{-grame în ground truth}}$

Realizati grafice sau tabele pentru a scoate in evidenta cum influenteaza scorul fiecare pas din algoritm (cel putin urmatoarele: cum se comporta algoritmul fara fiecare din cei 3 pasi, cum influenteaza ponderea similaritatii cu titlul performanta, cum influenteaza politica de ponderare a propozitiei in raport cu locatia ei in text performanta). In comparatie se vor folosi metricile pentru unigrame, bigrame si 4-grame, adica (BLEU@1, BLEU@2, BLEU@4 si ROUGE@1, ROUGE@2 si ROUGE@4). Pentru a calcula scorul pe tot datasetul se va face media scorurilor pentru fiecare text in parte.

Bonus [max 2p]

Aduceti imbunatatiri proprii algoritmului. Punctajul pentru bonus se va acorda in functie de complexitatea imbunatatirilor aduse cat si de performanta acestora.

Descriere set de date

Setul de date pentru evaluare este in format csv si contine urmatoarele coloane:

- headlines - titlul
- text - sumarizarea textului (ground truth, se va folosi doar pentru evaluare)
- ctext - textul complet

Datasetul se gaseste aici:

<https://drive.google.com/open?id=1jNqEdgHYcrX9aa41guTvJtUxyB5-fnuuy>

Alte precizari

In afara de cazurile specificate in enunt, nu se pot folosi resurse externe pentru a realiza tema.