

Klasifikácia YouTube komentárov pomocou neurónovej siete

MAROŠ KOPEC

Masarykova Univerzita
487595@muni.cz

20. júna 2019

I. ÚVOD

Klasifikácia spamu je problémom, ktorému sa venuje množstvo odborných článkov a prác. Za spam sa považuje nevyžiadaná správa rozosielená veľkému počtu adresátov alebo rozosielená na mnoho miest, zväčša za účelom reklamy. Pre množstvo internetových stránok je práve spam veľkým problémom, pretože ich používatelia sú zahltený obsahom, ktorý im znepríjemňuje skúsenosť s ich obsahom. Konkrétne YouTube je terčom veľkého množstva spamu od používateľov, ktorí sa týmto spôsobom snažia zviditeľniť vlastné video či kanál.

II. SÚVISIACE PRÁCE

Táto práca je inšpirovaná prácou TubeSpam[1] z Federálnej Univerzity Sao Carlos. V tejto práci akademici porovnávali niekoľko metód klasifikácie spamu, menovite rozhodovacie stromy, K-najbližších susedov, logistickú regresiu, Bernoulliho naivný Bayes, Gaussov naivný Bayes, Multinomiálny naivný Bayes, náhodné lesy, Support vector machines s lineárnym kernelom, Support vector machines s polynomiálnym kernelom a Support vector machines s Gaussovým kernelom. Práca dopĺňa vyššie zmienenu štúdiu o modely postavené na rekurentných neurónových sieťach.

V pôvodnej práci autori klasifikovali spam oddelene pre každé video. Neurónová sieť

Tabuľka 1: Kompozícia datasetu

Dataset	YouTube ID	# Spam	# Ham	Total
Psy	9bZkp7q19f0	175	175	350
KatyPerry	CevxZvSJLk8	175	175	350
LMFAO	KQ6zr6kCPj8	236	202	438
Eminem	uelHwf8o7_U	245	203	448
Shakira	pRpeEdMmmQ0	174	196	370

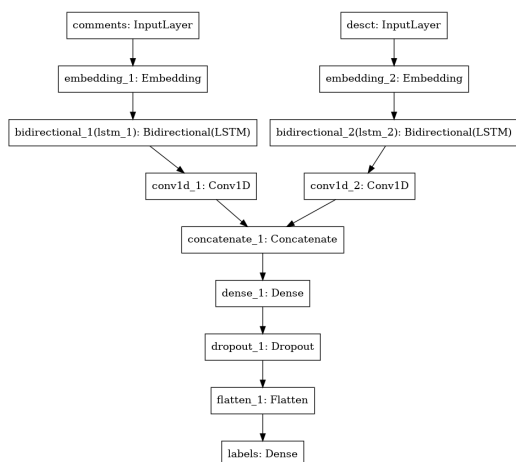
potrebuje veľké množstvo záznamov, aby bola schopná naučiť sa rozpoznávať spam. Dataset však obsahuje len obmedzený počet záznamov. Preto sa dataset spojil a všetky experimenty sa vykonávali nad zmiešanými záznamami.

III. DATASET

Pre trénovanie a testovanie modelov sú použité rovnaké dáta ako boli použité v práci TubeSpam¹. Dataset obsahuje 1956 záznamov z 5 najpozeranejších YouTube videí vo formáte csv. Kompozícia datasetu je znázornená v tabuľke 1. Pre potreby experimentu bol dataset rozšírený o popis videa.

Pre predstavu obsahu záznamov boli vygenerované slovné mapy, ktoré sú zobrazené na obrázkoch 1 a 2. Nad záznamami bol vykonaný experiment, pri ktorom bolo náhodne vybraných 100 záznamov z videa od autora Eminem, ktoré boli znovu označované. Z týchto sto záznamov bolo 8.9% označených odlišne.

¹<https://archive.ics.uci.edu/ml/datasets/YouTube+Spam+Collection>



Obr. 4: Model č.2

chybe nepodarilo zachovať. K dispozícií sú len dáta z optimalizácie hyperparametrov. Tieto výsledky boli zohľadnené pri experimentoch s druhým modelom 4.

i. Metódy evaluácie

Pre zhodnotenie modelu boli použité štatistické metriky *presnosť*, *chytený spam*, *blokováný ham*, *F-measure*, *Matthews korrelačný koeficient*, ktoré boli použité aj pri evaluácii vo vyššie spomínanej práci TubeSpam[1]. Presnosť vyjadruje koľko percent spamu klasifikátor úspešne označil za spam. Chytený spam a blokováný ham udávajú percentuálny podiel označeného spamu k neoznačenému a ham komentáre označené za spam. Najlepšie výsledky zobrazené v tabuľke 2 sa podarilo dosiahnuť s parametrami vymenovanými nižšie.

- Dropout: 0.3
- Activačná funkcia: elu
- Počet neurónov Dense vrsty: 15
- Kernel inicializácia: normal
- Posledná aktivačná funkcia: sigmoid
- Optimizér: Nadam
- Počet epoch: 12

Vysvetlivky pre tabuľku 2 a 3: acc - presnosť [%], sc - chytený spam [%], hb - blokováný ham [%], F1 - F-measure, MCC - Matthews korrelačný koeficient.

Tabuľka 2: Najlepšie výsledky modelu č.2

acc	sc	hb	F1	MCC
69.9%	84.7%	49.4%	0.82	0.008

Tabuľka 3: Najlepšie výsledky práce TubeSpam

acc	sc	hb	F1	MCC
97.73%	95.77%	0.00%	0.978	0.955

Zatiaľ čo presnosť je pomerne dostačujúca, počet chybné označeného hamu za spam je alarmujúco vysoký, až takmer 50%. Keďže žiaden z experimentov nedosiahol oveľa lepšie výsledky môžeme z toho vyvodit', že žiaden z experimentov vhodne nerieši problém klasifikácie spamu.

Naproti tomu práca TubeSpam[1] dosiahla oveľa lepšie výsledky, zobrazené v tabuľke 3. Táto skutočnosť môže byť spôsobená nedostatkom vstupných dát pre učenie rekurentnej neurónovej siete. Model sa nebol schopný naučiť rozoznávať spam na takmer 2000 záznamoch. Predpoklám, že ak by sa vstupné dáta zvýšili 100-násobne (t.j. aspoň na 200000), boli by výsledky oveľa lepšie. Ďalším faktorom, ktorý mohol negatívne ovplyvniť výsledky môže byť vstup krátkeho charakteru. Komentáre sú často len krátke pár slovné heslá. Pre túto úlohu by mohlo byť vhodnejšie spracúvať vstup nie po slovách ale po znakoch. To má však opäť rovnaké limitácie vo forme malého datasetu.

VI. PRÍLOHA

i. Návod na spustenie

1. Je potrebné stiahnuť Dataset ³, Embedding ⁴.
2. Nainštalujte si nástroj:
pipenv
3. Spustite príkaz:

³<https://archive.ics.uci.edu/ml/machine-learning-databases/00380/>

⁴<https://nlp.stanford.edu/projects/glove/>

```
# pipenv install
```

4. Pre spustenie tréovania modelu použite príkaz:

```
# python classifier.py
```

Kde NUMBER je číslo experimentu a meno zložky výstupov.

LITERATÚRA

- [1] Alberto, T.C., Lochter J.V., Almeida, T.A. *TubeSpam: Comment Spam Filtering on YouTube*. Proceedings of the 14th IEEE International Conference on Machine Learning and Applications (ICMLA'15), 1-6, Miami, FL, USA, December, 2015. (preprint)