

Извор: <https://www.aclweb.org/anthology/L18-1506.pdf>

ТЕРМИНИ:

The Recall-Oriented Understudy for Gisting Evaluation (Rouge)-score -> скаларна вредност во ранг од $[0,1]$. ја прикажува сличноста помеѓу еден и еден или повеќе токенизирани документи. Го оценува квалитетот на моделите за сумирање.

Сентиментален - нешто што одма ми ги разгорело луѓето. Пример, ако напишеш твит, Donald Trump се самоубил, тоа е сентиментален твит.

ВОВЕД

Во воведот се објаснува за претходни истражувања во врска со текст сумаризацијата. Проблемот настанувал поради мали податочни множества за проценка кој е подобар модел за сумаризација. Пример за твитовите, авторите морале сами да создаваат податочни множества. Иако имало множество од твитови со преку 230 000, ипак имало проблем со референците до нив што влијаеле врз проценките на моделот.

Овој научен труд уствари претставува текст сумаризација, прикажувајќи ново множество на твитови кои обработуваат 6 теми од социјален аспект. За да покажат дека овој dataset е уствари потенцијално употреблив, користат хибриден TF-IDF за извлекување на нови твитови, за да се докаже дека има компетитивен ROUGE скор.

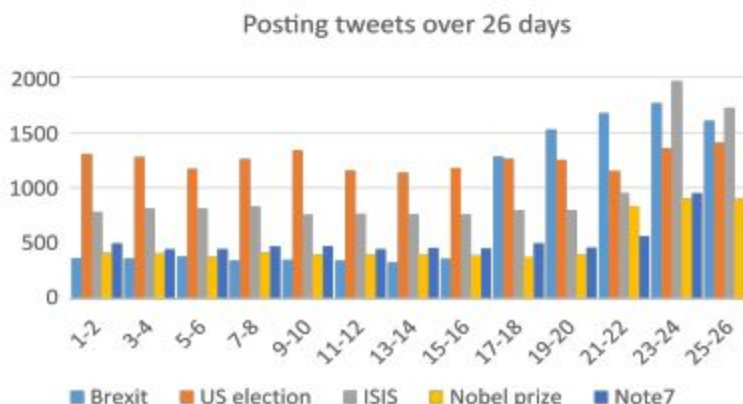
СУМАРИЗАЦИЈА:

Се дели на два чекора: data creation и моделот вклучувајќи и tweet scoring, selection

Data Creation

Значи за да се направи под. М-во потребно е да листа од области кои треба задоволуваат некои услови. Како на пример, да бидат преземени од различни извори, не се постари од 30 дена и мора да бидат

некои популарни (импресивни) да речеме..Потоа за секоја област (topic) се доделува листа со клучни зборови, поединечно.



Event	#tweets	#hashtags	Keywords
Brexit	10,978	9,705	brexit, #brexitshambles, #Brexiters, #BrexitCentral, England Europe exit.
US election	17,714	8,566	Donald Trump, Trump, Hilary Clinton, Clinton #debate, election.
ISIS	13,047	9488	ISIS, IS Syria, IS Mosul, IS Iraq, ISIS Aleppo, ISIS US, ISIS Rusia.
SS Note 7	7,362	7,465	Galaxy Note 7, #note7, #GalaxyNote7, thegalaxynote7, #SamsungGalaxyNote7.
Nobel prize	6,812	2,780	Nobel prize 2016, Nobel peace, Nobel chemistry, Nobel economy, Nobel physics.
SpaceX	4,982	2,417	Facebook SpaceX, SpaceX Explosion, Falcon 9 exploded, Falcon 9 explosion.

Ова е распределба на твитови. Покажува статистика за користење на клучните зборови.

МОЖЕ ДА БИДАТ ИСКОРИСТЕНИ ЗА ДА СЕ ПРОЦЕНИ КВАЛИТЕТОТ НА ТВИТОВИТЕ, МЕРЕЈЌИ КОЛКУ Е БИТНА ИНФОРМАЦИЈАТА КОЈА Е ГЕНЕРИРАНА ОД КОРИСНИЦИТЕ НА ТВИТОВИТЕ.

Податочна сигментација

Собраните твитови се делат во кластери. Идејата е таква бидејќи дури и бројот на твитови на ден да е мал, директно извлекување на подмножество од овие твитови можеби ќе ги елиминира останатите кои се битни. Значи... со кластерирањето, целта е да твитовите бидат што порепрезентивни.

Е сега за да биде се тоа во реално време од аспект на твит сумаризацијата, се користи алгоритам **Affinity Propagation** , алгоритам за кластерирање.

AP алгоритмот, го идентификува подмножеството од data points како темплејтови и формира кластери така што ги доделува останатите data points во еден од тие темплејтови (exemplars). После кластерирањето и елиминирање на тие денови со малку твитови, се формирале 6 множества

кои кореспондираат на бте настани во тие 26 дена и секој ден вклучува м-во од кластери, кои на некој начин се подобласти.

Standard reference creation

Ова се користи за еволуација на моделот. Се користат 3 методи

1. Luhn кој е хеuristicки метод за извлекување на тие податоци.
2. LexRank е графички-базиран модел, кој гради граф на сличност и ги селектира оние кои се засноваат на нивниот сопствен вектор (ова е од линеарна алгебра и сега не е толку битно да го знаеш ти).
3. DSDR-non се базира на ненегативна линеарна податочна реконструкција.

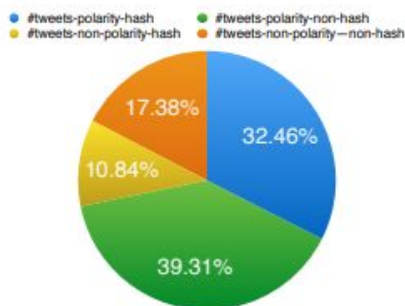
Овие извлечени твитови од овие три методи формираат множества кои се кандидати.

Method	ROUGE-1	ROUGE-2	ROUGE-SU
Luhn	0.556	0.502	0.250
LexRank	0.581	0.516	0.273
DSDR-non	0.419	0.310	0.145

Според ова LexRank е најдобар метод, бидејќи има најголема ROUGE вредност

Обзервација на податоците:

Луѓето последователно користеле емотивни зборови и хаштази во нивните твитови. Емотивните зборови ги прикажувале нивните интереси и хаштазите ги користеле како назначување за кој настан зборува.



Поларитетот е уствари позитивно, негативно, неутрално мислење.

Figure 2: Hashtag and polarity observation on six datasets.

МОДЕЛ

Tweet scoring

$$TF-IDF = tf_{i,j} * \log_2 \frac{N}{df_j} \quad (1)$$

$tf_{i,j}$ - фреквенцијата на терминот T_j во документот D_i , каде N е вкупен број на документи и df_j е бројот на документи кој го содржи терминот T_j

Равенката претсавува мерка која пресметува колку е релевантен еден збор во документ во однос на една колекција од документи.

Понатаму, во трудот се објаснува за недостатоците и предностите од аспект на извлекувањето на твитовите. Поради тие недостатоци, се усвојува нов метод, односно хибриден TF-IDF. Разликата е во тоа што, сите твитови се сметаат како еден единствен документ при пресметување на TF, а секој твит како посебен документ при пресметување на IDF.

$$h_{TFIDF}(t) = \frac{\sum_{i=0}^{\#WordsInTweet} W(w_i)}{nf(t)} \quad (2)$$

$$W(w_i) = tf(w_i) * \log_2(idf(w_i)) \quad (3)$$

$$tf(w_i) = \frac{\#OccurrencesOfWordInAllTweets}{\#WordsInAllTweets} \quad (4)$$

$$idf(w_i) = \frac{\#Tweets}{\#TweetsInWhichWordOccurs} \quad (5)$$

$$nf(t) = \#WordsInTweet \quad (6)$$

w_i е i -тиот термин во твитот t , $W()$ враќа тежина на терминот, $tf()$ враќа TF скор, соодветно $idf()$ и $nf()$ е некаков фактор на нормализација на твитот, бидејќи се долги твитовите.

Tweet selection

Се одбираат $\text{top } m$ рангирани твитови по најголемите скорови кои се селектирани како сумаризација за секој кластер(документ)

Експеримент и резултати

Се споредуваат и други методи

Method	ROUGE-1	ROUGE-2	ROUGE-SU
KL	<u>0.394</u>	<u>0.263</u>	<u>0.146</u>
LSA	0.462	0.368	0.175
Sumbasic	<u>0.444</u>	<u>0.298</u>	0.174
TextRank	0.495	0.418	0.213
Retweet	<u>0.384</u>	<u>0.264</u>	<u>0.129</u>
DSDR-lin	0.460	<u>0.351</u>	0.183
h-TFIDF	<i>0.482</i>	<i>0.384</i>	<i>0.199</i>

ROUGE scores with hashtags We also evaluated all the methods by using hashtags. The intuition is that tweets usually include hashtags, which show important information regarding user’s interests. To do that, we extracted all hashtags of each cluster to form its artificial references.

Table 4: The average ROUGE scores over six datasets.

Method	ROUGE-1	ROUGE-2	ROUGE-SU
KL	0.122	0.033	0.015
LSA	0.111	<i>0.034</i>	<i>0.017</i>
Sumbasic	0.137	<i>0.034</i>	<i>0.018</i>
TextRank	0.100	0.030	0.011
Retweet	0.101	0.024	0.007
DSDR-lin	0.117	0.033	0.011
DSDR-non	<i>0.123</i>	0.036	0.010
Luhn	0.105	0.032	0.011
LexRank	0.118	0.033	0.013
h-TFIDF	0.113	0.031	0.010

Обзервација на поларитетот кај твитовите

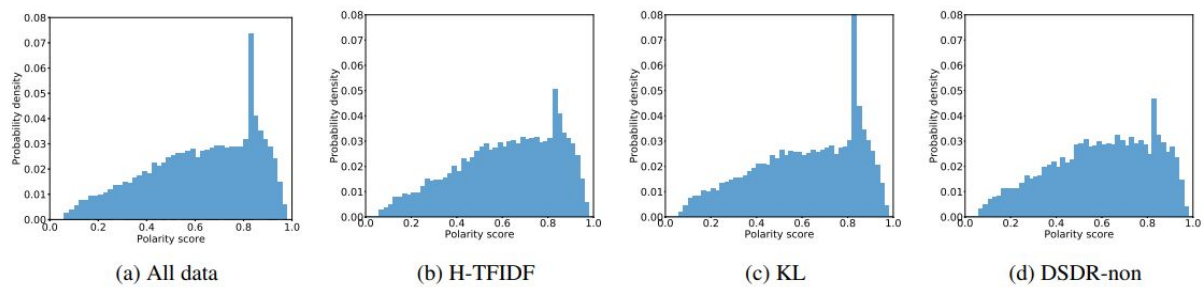


Figure 3: Polarity distribution of original and extracted tweets from three models.

Се тренира модел на класификатор за да види колку и каков е поларитетот на твитовите. Значи се зима м-вото од податоци, се тренира CNN, потоа се предвидува, тој скор што ќе се добие усвари е конвертиран во ранг од $[0, 1]$ и тоа е интентитетот на поларитет, каде твитовите со високи скорови се всушност несентименталните твитови, а оние кои се до приближно 0 се сентиментални.

Различни методи тестирани според интентитетот поларитетот на твитот заедно со нивниот закон на распределба.