

1 بنام خدا - محمد محبت صادقی - پروژه سسم خوش مصنوعی - شماره دانشجویی: 86199483

3 1- Stemming و Lemmatization ، دو روش برای normalize کردن کلمات می باشد.

5 در روش Stemming ، چند حرف آخر کلمه حذف می شود (حتی اگر حرف باقی مانده دارای معنی خاصی نباشد).

7 وقتی dataset بزرگ باشد و به Performance بالا نیاز داشته باشیم استفاده بهتر است

9 در Lemmatization ، کلمه base را می یابیم که دارای معنی باشد و به قواعد زبان بستگی دارد.

11 نیاز به Lookup table دارد، implement کردن اش سخت تر است. لذا ابراهیم بیشتر

13 با تفسیر بالا ، ما از Library حفظ و Lemmatization برای تبدیل کلمات به ریشه آن ها استفاده می کنیم.

15 مثال از هر دورا می توانید در قسمت "Lemmatization vs stemming" مشاهده کنید.

18 تابع heizm.Stopwordslist ، 389 تا از کلمات رایج را ~~از لیست~~ به ما می دهد و ما آن ها را

19 رایج از داده ها حذف می کنیم.

$$P(C|x) = \frac{P(x|C) P(C)}{P(x)}$$

22 2- Bayes rule :

24 مجموعه text چند از x_1 تا x_n : x_1, \dots, x_n : موضوع چند C : $topic$

1 $P(C|x) :$

2 احتمال اینکه خبر از یک topic خاص باشد به شرطی که متن خبر x باشد.
3 که هدف ما یافتن این احتمال می باشد.

4
5 $P(x|C) :$

6 احتمال اینکه متن خبر x باشد به شرطی که از topic C باشد.
7 ما $P(x_i|C)$ را با می بر تعداد $accoring$ ها هر کلمه می باشد.

8
9 $P(C) :$

10 احتمال اینکه خبر از topic C باشد. چون در صورت
11 مسئله گفته شده که
12 از تمام شان به یک اندازه دار شده. این احتمال را $\frac{1}{N}$ در نظری می گیریم.

13 $P(x) :$

14 احتمال اینکه متن خبر x باشد. که برابر این کا هم جمع تمام x ها در خبر
15 ها دیده شده است به دست می آید.

16
17
$$P(C|x_1, \dots, x_n) = P(C) P(x_1|C) P(x_2|C) \dots P(x_n|C)$$

18
19 بدلیل استفاده از $Naive\ Bayes$ می توانی از نتیجه لیست بالا استفاده کنی.

20
21 ۳- مثال: میت روز و طلا افزایش یافت

22 روز دزدی من سستم

23
24 در اینجا $bigram$ ها هست. گاهی کلمات همراه کلمات دیگر $Context$ معنا می گیرند.

25
26 من هم در مدل خود از هر دو $bigram$ و $unigram$ استفاده کردم به علت

27
28 کم بودن تعداد را به $bigram$ ها، به آن ها $weight$ نیز اضافه کردم.

1. Precision % در این حجم ما توجهی به False negative ها نمی کنیم.

(تا 100)

2. یعنی فرض کنید تعداد زیاد ~~میانگین~~ خبر سیاسی داریم و ما یکی

3. را سیاسی تشخیص می دهیم. $Precision = 100\%$ خواهد شد که باز هم به

4. نتایج خندان صعب است.

5. F_1 از میانگین $harmonic$ استفاده می کنند. اگر False Positive

6. False negative برابر ما هزینه بزرگی داشته باشد، استفاده از Accuracy

7. خوب است. ولی در عمل ~~صحیح~~ زمانی که این هزینه ها متفاوت

8. استفاده از F_1 اهمیت زیاد پیدا می کند.

9. از میانگین F_1 تمام مدل ها استفاده می کنیم و $macro (macro-F_1)$

10. $micro$ % Accuracy $\frac{Correct}{total}$ که حاصل برابر است

11. $weighted$ به هر موضوع، به اندازه تعداد خبر هر در آن مدل،

12. به آن $weight$ اختصاص می دهیم. انواع $weighted-F_1$, $Precision$, ...

دریم

Blue: Additive Smoothing

Red: without additive Smoothing

10. نتایج را در جدول مشاهده کنید

ربر

	Health	Political	Sports	Technology	Art	Accidents	All Classes
Precision	0.195 0.105	0.192 0.101	0.197 0.105	0.196 0.101	0.187 0.101	0.198 0.101	-
Recall	0.197 1	0.192 1	0.198 1.0	0.196 0.106	0.196 1	0.188 0.175	-
F1-score	0.196 0.101	0.192 0.101	0.1978 0.101	0.196 0.103	0.191 0.102	0.193 0.102	-
Accuracy	-	-	-	-	-	-	0.194
Macro Avg	-	-	-	-	-	-	0.194
Micro Avg	-	-	-	-	-	-	0.194
Weighted Avg	-	-	-	-	-	-	0.194

8 Micro & Macro Accuracy ✓ 1

$$\text{Accuracy} = \frac{\text{Total} - \text{Wrong}}{\text{Total}} = \frac{1115 - 80}{1115} = 94\%$$

$$\text{Macro F}_1 \text{ average} = \frac{0.96 + 0.92 + 0.978 + 0.96 + 0.91 + 0.93}{6} = 94\%$$

Micro Averages 94%

$$+ 200 \times 0.92 + 180 \times 0.91$$

$$\text{Weighted Average} = \frac{168 \times 0.96 + 190 \times 0.978 + 200 \times 0.93 + 177 \times 0.96}{1115} = 94\%$$

۱۱- همانطور که مشاهده می کنید، این احتمال بسیار بالا وجود دارد که

۱۲- حداقل در یک خبر که استفاده شود که در dataBase ما وجود ندارد، به همین

۱۳- خاطر مشاهده می کنیم که با افزودن Adaboost Smoothing، نتایج به طرز چشمگیری

۱۴- بهتر می شود.

۱۵- می توانید مشاهده کنید. مثلاً علی Predict استباه:

۱۶- ۱- استفاده از training Set کوچک

۱۷- ۲- Naive Bayes قرار دادن کلمات نسبت بهم

۱۸- ۳- در برخی مثال ها از برخی کلمات استفاده شده که کمتر در اخبار کاربرد دارند

۱۹- مثلاً در در مثال اول، استفاده از لفظ هر رئیس جمهور، وزیر

۲۰- چند بار امور خارج، انعکاس باعث شده که اخبار فرهنگی هنر را سیاسی

۲۱- بیش بینی کند که کاریش همیشه لرد!