

بخش ۳ - سوال ۴ -

stop word ها کلماتی هستند که شامل اطلاعات مهمی در باره متن نمی شوند و تعداد آن ها در متن زیاد است .

به عبارتی بهتر : کلماتی هستند که یا ناخیز هستند یا تعداد آنقدر زیاد است که سیستم نمی تواند آن ها را پردازش کند

مانند : a an and are as at from
has - - - will with

به طور کلی stop word ها تعداد زیاد کلمات حاوی اطلاعات غیر مفید در یک مدل زبانی هستند که تأثیر زیادی در به دست آوردن احتمال هر کلمه ایجاد نمی کنند زیرا تعداد آن ها زیاد است و در فرمول احتمالی مندرج را با کلماتی که حاوی اطلاعات خاصی نیستند زیاد می کنند .

درست است یا نه ؟ stop words کلمات بدون

اطلاعات ضروری و پرکار را از مجموعه زبانی حذف کنیم تا هم به سرعت و وقت بهتری برسیم.
 (برگرفته از Search در گوگل!)
 (در خزانه اشاره شده)

بخش ۳ - سؤال ۵ -

فرمول Perplexity :

$$\sqrt[N]{\frac{1}{P(\omega_1 \dots \omega_N)}}$$

3-grams برای $PP(w) = \sqrt[N]{\prod_{i=2}^N \frac{1}{P(\omega_i | \omega_{i-1}, \omega_{i-2})}}$

4-grams برای $PP(w) = \sqrt[N]{\prod_{i=3}^N \frac{1}{P(\omega_i | \omega_{i-1}, \omega_{i-2}, \omega_{i-3})}}$

۲.۱ فرمول ۱۱ - $N \leftarrow$ تعداد کلمات - ۱ -

$$P(\omega_i | \omega_{i-1}, \omega_{i-2}) \geq P(\omega_i | \omega_{i-1}, \omega_{i-2}, \omega_{i-3}) \leftarrow \text{میدانیم}$$

← با 4-grams مخرب کربالا کوچکتری شود

← مقدار زیر رادیکال اضافه می شود

← مقدار سرگشتگی بزرگتری شود.

(معلمین سیستم)

بخش ۳، سوال ۶-

$\langle s \rangle B C A B A \langle /s \rangle$

$$P(B|\langle s \rangle) \times P(C|B) \times P(A|C) \times P(B|A) \times$$

$$P(A|B) \times P(\langle /s \rangle | A)$$

$$= \frac{P(\langle s \rangle, B)}{P(\langle s \rangle)} \times \frac{P(B, C)}{P(B)} \times \frac{P(C, A)}{P(C)} \times \frac{P(A, B)}{P(A)}$$

$$\times \frac{P(B, A)}{P(B)} \times \frac{P(A, \langle 15 \rangle)}{P(A)}$$

$$= \frac{1}{5} \times \frac{4}{8} \times \frac{4}{14} \times \frac{5}{11} \times \frac{0}{8} \times \frac{0}{11} = 0$$

به نظرم نباید عبارت بالا ∞ باشد، ولی هر چه جزوه
 رونده می‌کنم هیچ راه کاری برای این شرایط پیدا
 نمی‌کنم. \leftarrow در این شرایط $PP(w) = \infty$

بجسٽي ۽ سوال ۱

T	v	9	0	9	0	ε	0	۲	۱
C	۶	0	۶	0	۶	۳	۲	۱	۲
E	0	ε	۳	۶	۳	۲	۱	۲	۳
R	۳	۳	۲	۳	۲	۱	۲	۳	۳
R	۳	۲	۱	۲	۱	۲	۳	ε	0
O	۲	۱	0	۱	۲	۳	ε	0	۶
C	۱	0	۱	۲	۳	ε	0	۶	v
#	0	۱	۲	۳	ε	0	۶	v	^
	#	C	O	O	R	R	E	C	T

۱ = فاصلو