

NLP HW4. Answer key

Yutong Wangyan.

1.

$$(a). (i) = \text{Accuracy} = \frac{4+8}{15} = \frac{4}{5}$$

$$(ii) = \text{precision} = \frac{4}{5}$$

$$(iii) = \text{recall} = \frac{4}{6} = \frac{2}{3}$$

$$(iv). 2 \times \frac{P \times R}{P+R}$$

$$= 2 \times \frac{\frac{4}{5} \times \frac{2}{3}}{\frac{4}{5} + \frac{2}{3}} = 0.727$$

(b).

		predicted	
		P	N
actual	P	4	2
	N	1	8

(c). $P(\text{gender_actual} = \text{"women"})$ = probability of an observation of being women actually.

$P(\text{predicted women} | \text{actual women})$ = probability of an observation being predicted as women given the observation is women actually.

(d). Recall should be prioritized because recall cares about how many predicted values are positive out of all actual positives. A high recall indicates most women clothings are categorized correctly.

2.

(A): (i). $\frac{1}{3}$

(iv): $P(y|x) = \frac{P(x|y) P(y)}{P(x)} = \frac{\frac{1}{3} \times \frac{1}{2}}{\frac{1}{3}} = \frac{1}{2}$

(ii): $\frac{2}{6} = \frac{1}{3}$

(iii): $\frac{2}{6} = \frac{1}{3}$

(B). if "love" and "movie" are independent,

then $P(\text{love and movie})$ should equal to $P(\text{love}) \times P(\text{movie})$.

$P(\text{love and movie}) = \frac{1}{6}$

$P(\text{love}) \times P(\text{movie}) = \frac{1}{6} \times \frac{4}{6} = \frac{1}{9} \neq \frac{1}{6}$

.. so "love" and "movie" are not independent.

3.

(a).	trendy	jeans	old	blue	red	woof
tfdf	2	$\frac{7}{4}$	2	2	$\frac{5}{2}$	$\frac{5}{2}$
tf(A)	1	1	0	0	0	0
tf(B)	0	1	2	1	0	0
tf(C)	1	1	1	1	1	1
tfidf(A)	2	$\frac{7}{4}$	0	0	0	0
tfidf(B)	0	$\frac{7}{4}$	4	2	0	0
tfidf(C)	2	$\frac{7}{4}$	2	2	$\frac{5}{2}$	$\frac{5}{2}$

(b)

$$\begin{array}{ccccccc} q & 0 & 1 & 1 & 0 & 0 & 0 \\ f(q) & 0 & \frac{7}{4} & 2 & 0 & 0 & 0 \end{array}$$

$$\cos(q, B) = \frac{\frac{7}{4} \times \frac{1}{4} + 2 \times 4}{\sqrt{\left(\frac{7}{4}\right)^2 + 2^2} \cdot \sqrt{\left(\frac{1}{4}\right)^2 + 2^2}} = 0.867$$

$$\cos(q, c) = \frac{\frac{3}{4} \times \frac{1}{4} + 2 \times 2}{\sqrt{\left(\frac{3}{4}\right)^2 + 2^2 + 2^2} + \sqrt{\left(\frac{1}{4}\right)^2 + 2^2}} = 0.506 \approx 0.867$$

Therefore, we recommended

Therefore, we should recommend product B.

4.
(A)

|A>.

1. I love go to store
 2. He love work at restaurant
 3. Stone is close today
 4. today I am worl
 5. He is go to restaurant

THERMOPHILIC BACTERIA

(B).

$$P(X_1 = I \mid x_0 = \text{start}) = 0.2$$

$$P(X_2 = \text{love} \mid X_1 = I) = 0.5$$

$$P(X_3 = \text{work} \mid X_2 = \text{love}) = 0.5$$

$$P(X_4 = \text{End} \mid X_3 = \text{work}) = 0.5$$

$$0.2 \times 0.5 \times 0.5 \times 0.5 = 0.025$$

(C).

- ①. stemming / lemmatization \rightarrow reduce different forms of the same word / token
- ②. removal of stopwords so the probability of target word following other words is precise.
- ③. Perplexity : we can't compare sentences with different lengths as longer sentence has smaller probability. We can apply perplexity formula $\sqrt[N]{\prod p(\text{sentence})}$ to reduce the effect of sentence length.
minimizing the perplexity \Leftarrow maximizing the probability of sentence to be in natural language)