

سوال 1 عملی

در فایل نوت بوک قرار داده شده، توضیحات هر بخش کامل نوشته شده است.

سوالات تئوری

(1

الف)

We need to compute $p(\text{text} \mid \text{category}) * p(\text{category})$ for each category and then compare them to find most proper category.

$$P(\text{'a very close race' } \mid \text{politics}) * p(\text{politics}) =$$

$$p(\text{'a' } \mid \text{politics}) * p(\text{'very' } \mid \text{politics}) * p(\text{'close' } \mid \text{politics}) * p(\text{'race' } \mid \text{politics}) * p(\text{politics})$$

$$P(\text{'a very close race' } \mid \text{not politics}) = ?$$

$$p(\text{'a' } \mid \text{not politics}) * p(\text{'very' } \mid \text{not politics}) * p(\text{'close' } \mid \text{not politics}) * p(\text{'race' } \mid \text{not politics}) * p(\text{not politics})$$

$$p(\text{politics}) = 3/5$$

$$p(\text{not politics}) = 2/5$$

$$p(\text{'a' } \mid \text{politics}) = 1/12$$

$$p(\text{'very' } \mid \text{politics}) = 1/12$$

$$p(\text{'close' } \mid \text{politics}) = 0$$

$$p(\text{'race' } \mid \text{politics}) = 1/12$$

$$p(\text{'a' } \mid \text{not politics}) = 2/12$$

$$p(\text{'very' } \mid \text{not politics}) = 0$$

$$p(\text{'close' } \mid \text{not politics}) = 1/12$$

$$p(\text{'race' } \mid \text{not politics}) = 1/12$$

$$p(\text{'a very close race' } \mid \text{politics}) * p(\text{politics}) = 1/12 * 1/12 * 0 * 1/12 * 3/5 = 0$$

$$p(\text{'a very close race' } \mid \text{not politics}) * p(\text{not politics}) = 2/12 * 0 * 1/12 * 1/12 * 2/5 = 0$$

we can not decide about category !! we need laplace smoothing

(ب)

We add a value of 0.5 to each word count in order to smooth!

$$P(\text{politics}) = 3/5$$

$$P(\text{not politics}) = 2/5$$

$$|V| = 17$$

$$P('a' | \text{politics}) = (1 + 0.5) / (12 + 0.5 * 17) = 1.5 / 20.5 = 0.0731707317$$

$$P('very' | \text{politics}) = (1 + 0.5) / (12 + 0.5 * 17) = 1.5 / 20.5 = 0.0731707317$$

$$P('close' | \text{politics}) = 0.5 / 20.5 = 0.0243902439$$

$$P('race' | \text{politics}) = (1 + 0.5) / (12 + 0.5 * 17) = 1.5 / 20.5 = 0.0731707317$$

$$P('a' | \text{not politics}) = (2 + 0.5) / (12 + 0.5 * 17) = 2.5 / 20.5 = 0.1219512195$$

$$P('very' | \text{not politics}) = 0.5 / (12 + 0.5 * 17) = 0.5 / 20.5 = 0.0243902439$$

$$P('close' | \text{not politics}) = (1 + 0.5) / (12 + 0.5 * 17) = 1.5 / 20.5 = 0.0731707317$$

$$P('race' | \text{not politics}) = (1 + 0.5) / (12 + 0.5 * 17) = 1.5 / 20.5 = 0.0731707317$$

$$0.0731707317 * 0.0731707317 * 0.0243902439 * 0.1219512195 * (2/5) >$$

$$0.0731707317 * 0.0731707317 * 0.0243902439 * 0.0731707317 * (3/5)$$

As shown, the probability to be in 'non politics' category is more than 'politics' category!

Yes, as shown in this part Laplace smoothing helps us add some small possibility to words not seen in train data of a category.

سوال (2)

(الف)

$$\text{Precision} = 70 / 70 + 30 = 0.7$$

$$\text{Recall} = 70 / 70 + 70 = 0.5$$

(ب)

کاربر اول نمیخواهد لزوماً همه ایمیل های درست را بگیرد. او بیشتر ترجیح می دهد که از بین ایمیل های دریافت شده اکثریت درست باشند نه اینکه همه لزوماً دریافت شده باشند. به عبارتی دقت (precision) برایش مهم تر است. در مقابل، کاربر دوم ترجیح می دهد تا جایی که میشود همه درست ها را ببیند حتی شده به قیمت دیدن تعداد زیادی اسپم میان پیشبینی شده ها. پس recall بیشتر برایش مهم است.