

۴-۴-۱. یک سند مانند d ممکنه چندین کلاس داشته باشه. مثلا هم خبری باشه و هم اقتصادی. فرض اینکه سند d از طریق c نوشته شده

که برای آن سند، احتمال اینکه c از طریق d نوشته شده باشه $p(c|d) = \frac{p(d|c)p(c)}{p(d)}$ (قانون بیز در مورد متن)

چون هدف یافتن کلاس با بیشترین احتمال است و از طرفی، سندها ثابت در محاسبات همجوشی کلاسها به نظر نمیانند

$$c_{MAP} = \text{arg max}_c p(d|c) \times p(c)$$

سند d در اینجا از معادله c که سبب شده است. $c_{MAP} = \text{arg max}_c p(u_1, \dots, u_n | c) \times p(c)$

پس اگر فرض کنیم تعداد کلمات n است: u_1, \dots, u_n و در ترتیب این کلمات $n!$ حالت مختلف داریم و چون c تا n کلاس داریم تعداد بالا ترصا $O(n! \cdot |c|)$ خواهد بود و این به $order$ نامیده است. با در نظر گرفتن شرط مستقل بودن کلمات

از یکدیگر، (که در عمل همیشه اینطور نیست و این به وفن غیر واقعیه است) c $order$ را به طور قابل علاقه n کاهش میدهیم.

(فرض کنیم c هم داریم اونم اینکه کلمات به هم ربط bag of $word$ هستند یعنی فرض کنیم ترتیب کلمات مهم نیست.)

$$c_{MAP} = \text{arg max}_c p(u_1, \dots, u_n | c) \times p(c)$$

$$p(u_1, \dots, u_n | c) = p(u_1 | c) \times p(u_2 | c) \times \dots \times p(u_n | c)$$

ما فرض استقلال کلمات

$$c_{NIB} = \text{arg max}_c p(c) \prod_{x \in X} p(x | c)$$

این فرض استقلال کلمات در قانون Bayes را در نظر بگیریم میشه معین $Naive Bayes$

پس علت تفاوت این از قول اینست که فرض دوم ما فرض مستقل بودن کلمات است تا تعداد پارامترها و حجم مورد نیاز به یادگیری ما

۵- از کل مجموعه V vocabulary c training set را استخراج می‌کنیم، بازاری که در آن از کلمات c استفاده می‌کنیم.

$$\hat{P}(c_j) = \frac{\text{doc count}(c=c_j)}{N_{\text{doc}}}$$

مربط به این کلمات c شمارم و تقسیم بر تعداد کل document های کنیم.

احتمال هر کلمه در کلاس c_j برابر است با $\hat{P}(c_j)$ ، یعنی تعداد آن کلمه در کلاس c_j تقسیم بر تعداد همه کلمات کلاس c_j :

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

ما توجه به جدول بالا، اگر کلمه از کلمات test set در کلاس c_j train مربوط به این موضوع دیده شده، احتمال

مربط به این موضوع c_j است. دل باید به این باشد که کلمات دیده شده، لانز احتمال وقوع (هد) - یعنی

ممکن است کلمه w_i در train نباشد ولی دلیل نمی‌شود که همان موضوعی تعلق نداشته باشد. برای جلوگیری

از 0 شدن احتمالات، از روش لاپلاس در ردیف Bayes استفاده می‌کنیم:

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j) + 1}{\left[\sum_{w \in V} \text{count}(w, c_j) \right] + |V|} = \frac{\text{count}(w_i, c_j) + 1}{\sum_{w \in V} (\text{count}(w, c_j) + 1)}$$

$P(c) = \frac{3}{8}$ $P(j) = \frac{1}{8}$

اقبال کلاس با
 table c

اقبال کلاس با
 table j

$\hat{p}(w|c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V|}$

chinese chinese chinese Tokyo japan

$P(\text{chinese}|c) = \frac{3+1}{8+4} = \frac{4}{12} = \frac{1}{3}$

$P(\text{chinese}|j) = \frac{1+1}{3+6} = \frac{2}{9}$

$P(\text{Tokyo}|c) = \frac{0+1}{8+4} = \frac{1}{12}$

$P(\text{Tokyo}|j) = \frac{1+1}{3+6} = \frac{2}{9}$

$P(\text{japan}|c) = \frac{0+1}{8+4} = \frac{1}{12}$

$P(\text{japan}|j) = \frac{1+1}{3+6} = \frac{2}{9}$

$P(c|d) = \frac{3}{4} \times \frac{1}{12} \times \frac{1}{12} \times \left(\frac{3}{7}\right)^3$

سہارا و درہ این طبعہ

≈ 0.0003 ✓

اقبال نوز
 کلاس c

$P(\text{japan}|c)$

$P(\text{Tokyo}|c)$

$P(\text{chinese}|c)$

کلاس c
 انتخاب ہو گا

$P(j|d) = \frac{1}{8} \times \frac{2}{9} \times \frac{2}{9} \times \left(\frac{2}{9}\right)^3 \approx 0.0001$

-6

$P(\text{باشگاه و مجلس و پرسولیس | سیاسی بودن}) = ?$

ضرب صورت ہا = 9

$P(\text{باشگاه و مجلس و پرسولیس و سیاسی}) = \frac{3}{7} \times \frac{1}{8} \times \frac{3}{8} \times \frac{1}{8}$

$P(\text{باشگاه و مجلس و پرسولیس و ورزشی}) = \frac{4}{7} \times \frac{3}{8} \times \frac{0}{8} \times \frac{3}{8}$

ضرب صورت ہا = 0

$P(\text{باشگاه و مجلس و پرسولیس | سیاسی بودن}) = \frac{9}{9+0} = \frac{9}{9} = 1$

$P(\text{باشگاه و مجلس و پرسولیس | ورزشی}) = \frac{0}{9+0} = \frac{0}{9} = 0$

-7