

1. Теоретичні відомості

1.1. Мультиномінальна модель та максимальна апостеріорна ймовірність

В мультиномінальній моделі, на відміну від моделі Бернуллі, враховується кількість входжень кожного слова в документ. $w_k, k = 1, \dots, D$ — всі унікальні слова в корпусі. x^k — кількість повторень слова w_k в документі θ_y^k — ймовірність зустрічі слова w_k на позиції

$$\hat{y}(x) = \operatorname{argmax}_y P(y|x) = \operatorname{argmax}_y P(y)P(x|y) = (1)$$

$$\begin{aligned} &= \operatorname{argmax}_y P(y) \frac{(\sum_{k=1}^D x^k)!}{\prod_{k=1}^D x^k!} \prod_{k=1}^D (\theta_y^k)^{x^k} = (2) \\ &= \operatorname{argmax}_y \ln P(y) + \sum_{k=1}^D (x^k \ln \theta_y^k) \end{aligned}$$

де $P(y) = \frac{N_y}{N}$, $\theta_y^k = \frac{n_{yk} + \alpha}{n_y + \alpha D}$ — емпіричні оцінки ймовірностей. N_y — кількість документів класу y n_y — кількість слів в документах класу y n_{yk} — кількість повторень слова w_k в документах класу y (1)— Байєсівське правило апостеріорної ймовірності (2)— припущення "наївно байєса"

1.2. Припущення про позиційну незалежність

Розглянемо множину об'єктів $D = d_1, \dots, d_m$, кожен з яких володіє набором ознак з множини всіх ознак $F = f_1, \dots, f_q$. Модель наївного байєсівського класифікатора приймає два припущення:

- 1. порядок признаков об'єкта не має значення
- 2. ймовірності признаков не залежать один від одного при заданому класі:
 $P(f_i \cap f_j | c) = P(f_j | c)$.

1.3. Застосування для інформаційного пошуку

Модель байєсівського наївного класифікатора часто використовується для інформаційного пошуку та класифікації текстів, оскільки вона часто дає хороші результати, хоч породжені нею оцінки часто далекі від ідеальних. Хоч зазвичай в моделі, яка називається наївною, ознаки розглядаються як категоріальні, або випадкові, величини з розподіли Бернуллі, варіанти, в яких використовується мультиноміальний розподіл, як правило краще моделюють входження слів в документ. Врахувати якісні ознаки можна двома способами: моделюючи їх як нормально розподілені в межах кожного класу чи за допомогою параметричної оцінки щільності.