

## Лабораторна робота №2.

Наївний байєсівський аналізатор в задачі класифікації тексту.

1. В якості навчальної вибірки необхідно використати dataset з лабораторної роботи #1 (текстові повідомлення та електронні листи).
2. Реалізувати алгоритм наївного байєсовського класифікатора для задачі класифікації тексту.
3. Для програми розробленої для лабораторної роботи №1 додати:
  - а. поле для задання повідомлення та вибору навчальної вибірки, яке необхідно класифікувати;
  - б. можливість запустити аналізатор та вивести для користувача до якого класу відноситься повідомлення.

Технології та мови програмування: рекомендована python, але може буде використана будь-яка.

### Теоретичний матеріал

Детальніше на 2 лекції

“Наївність” алгоритму полягає в тому, що ми припускаємо незалежність появи слів в повідомленнях. За теоремою Байєса:

$$P(\text{ham} \mid \text{bodyText}) = \frac{P(\text{ham}) * P(\text{bodyText} \mid \text{ham})}{P(\text{bodyText})}$$
$$P(\text{spam} \mid \text{bodyText}) = \frac{P(\text{spam}) * P(\text{bodyText} \mid \text{spam})}{P(\text{bodyText})}$$

Нам необхідно тільки порівняти  $P(\text{ham} \mid \text{bodyText})$  та  $P(\text{spam} \mid \text{bodyText})$ , тому  $P(\text{bodyText})$  - рахувати не потрібно.

$P(\text{ham})$  = кількість повідомлень з категорії ham / загальна кількість повідомлень

$P(\text{spam})$  = кількість повідомлень з категорії spam / загальна кількість повідомлень

Для повідомлення довжиною N:

$$\text{bodyText} = [\text{word1}, \text{word2}, \dots, \text{wordN}]$$
$$P(\text{bodyText} \mid \text{ham}) = P(\text{word1} \mid \text{ham}) * P(\text{word2} \mid \text{ham}) * \dots$$
$$P(\text{bodyText} \mid \text{spam}) = P(\text{word1} \mid \text{spam}) * P(\text{word2} \mid \text{spam}) * \dots$$

, де:

$P(\text{word1} \mid \text{spam}) = \text{кількість word1 в категорії spam} / \text{загальна кількість слів в spam}.$

Згладжування Лапласа - якщо слова немає в навчальній вибірці, ми вважаємо, що слова зустрічається один раз, але нам треба також змінити вірогідності інших слів:

$P(\text{word1} \mid \text{spam}) = (\text{кількість word1 які належать категорії ham} + 1) / (\text{загальна кількість слів, які належать категорії ham} + \text{кількість слів, яких немає в навчальній вибірці})$

Всі обчислення можна привести до логарифмічної форми, щоб не працювати з маленькими числами з плаваючою точкою:  
<http://getpopfile.org/docs/faq:bayesandlogs>