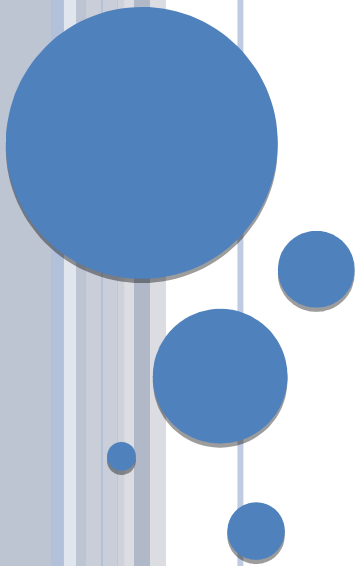


НАЇВНИЙ БАЄСІВСЬКИЙ КЛАСИФІКАТОР

Онищук Роман



КЛАСИФІКАЦІЯ

- Мета класифікації:
 - Відношення значення до певного класу
 - Класифікація даних

- Області застосування
 - банківський бізнес (кредити)
 - маркетинг
 - медицина



ЗАГАЛЬНИЙ ВИГЛЯД КЛАСИФІКАТОРА

- Побудова моделі: опис набору визначених класів
 - Кожен з атрибутів виборки належить до заздалегіть визначеного класу
 - Певний набір кортежів використовується для побудови моделі: training set
- Використання моделі: класифікація даних
 - Оцінювання точності заданої моделі
 - Значення отримане з тренувального набору порівнюється з відповідним значенням тестового набору
 - Оцінювана точність – відношення (зазвичай у відсотках) значень тестового набору, котрі були “коректно” класифіковані
 - Бажано, щоб тренувальний набір не залежав від тестового



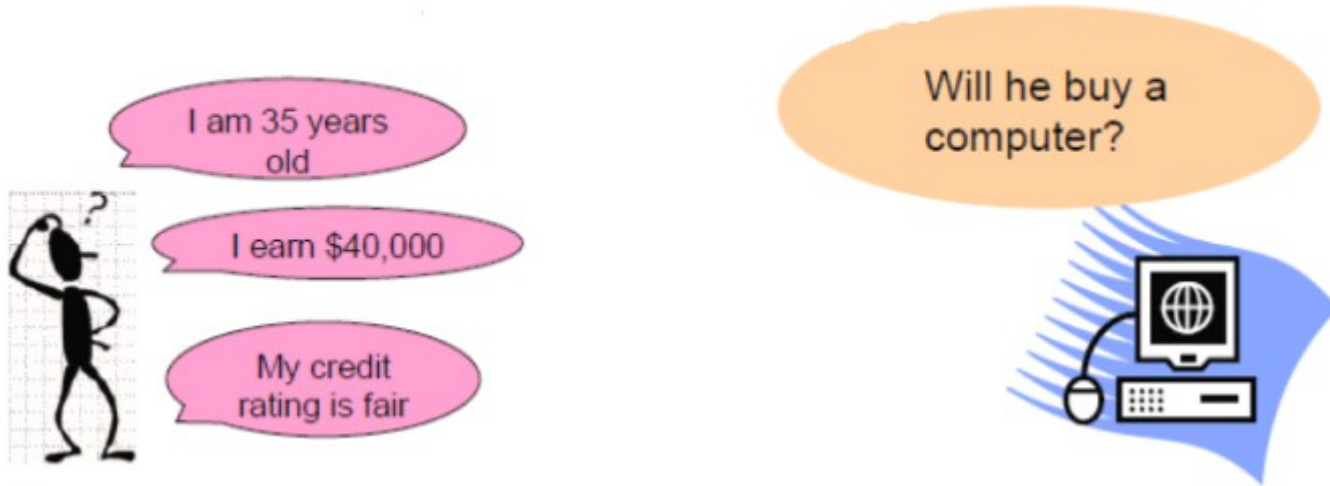
ВВЕДЕННЯ В БАЄСІВСЬКУ КЛАСИФІКАЦІЮ

○ Отож, що це?

- Статистичний метод класифікації.
- Один з базових методів використовуваних в машинному навчанні та інтелектуальному аналізі.
- В його основі лежить теорема Байєса.
- Вирішує проблеми, пов'язані як з дискретними так і з номінальними атрибутами.



Найпоширеніший приклад



X: 35 – річний чоловік, з зарплатою 40 тис. і хорошою кредитною історією

H: Гіпотеза, що цей чоловік купить комп'ютер.



ТЕОРЕМА БАЄСА

- Теорема Баєса:
 - $P(H|X) = P(X|H) P(H) / P(X)$
- $P(H|X)$: умовна ймовірність, того, що клієнт X буде купувати комп'ютер за умови, що ми знаємо вік і дохід клієнта. (Постеріорна ймовірність H)
- $P(H)$: Ймовірність того, що клієнт буде купувати комп'ютер, незалежно від віку та доходів (Апріорна ймовірність H)
- $P(X|H)$: Ймовірність того, що особа 35 років, заробляє \$40,000 точно купила комп'ютер
- $P(X)$: Ймовірність того, що людина з нашої виборки клієнтів 35 років і заробляє \$ 40,000.



BAYESIAN CLASSIFIER

- Нехай $X = \{x_1, x_2, \dots, x_n\}$ – вибірка, компоненти якої являють собою значення, n -мірні вектори.
- Припустимо існує m класів: C_1, C_2, \dots, C_m
- Баєсівський класифікатор передбачає, що X належить до класу C_i , якщо $P(C_i|X) > P(C_j|X)$, при $1 \leq i \leq m, i \neq j$.
- Отримуємо, максимальну апостеріорну ймовірність:

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)}$$

- Враховуючи наївне припущення взаємонезалежності, отримуємо, що $P(X) = \text{const}$.



NAÏVE BAYESIAN CLASSIFIER...

- Звідси й отримаємо:

$$\begin{aligned} P(X | C_i) &= P(x_1, x_2, \dots, x_n | C_i) \\ &= P(x_1 | C_i) * P(x_2 | C_i) * \dots * P(x_n | C_i) \\ &= \prod_{k=1}^n P(x_k | C_i) \end{aligned}$$



NAÏVE BAYESIAN CLASSIFIER...

- Ймовірності $P(x_1 | C_i), P(x_2 | C_i), \dots, P(x_n | C_i)$ можуть бути легко оцінені по навчальній вибірці, де x_k - це значення атрибуту A_k виборки X . Звідси, ми отримуємо два випадки:
- Якщо значення A_k - дискретна (категорійована) змінна (тобто A_k прийматиме значення з певної фіксованої множини), то $P(x_k | C_i)$ це кількість зразків класу C_i в множині значень навчальної вибірки T , що мають значення x_k для атрибуту A_k , розділене на частоту (C_i, T) , кількості C_i -тих зразів класу C в T .
- Якщо значення A_k безперервна (цифрова) змінна (тобто така, яка може приймати будь-яке значення в безперервному діапазоні), то зазвичай припускається, що змінні мають Гауссівський розподіл з середнім μ і стандартним відхиленням σ , і визначається формулою:

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp - \frac{(x - \mu)^2}{2\sigma^2},$$

$$P(x_k | C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i}).$$



ПРОБЛЕМА "НУЛЬОВОЇ ЧАСТОТИ"

- Що робити, якщо X має значення атрибуту x_k , таке, що для нього не існує відповідного значення атрибуту в класі C_i ? У цьому випадку $P(x_k | C_i) = 0$, що призводить до $P(X | C_i) = 0$, хоча $P(x_k | C_i)$ для всіх інших атрибутів в X може бути великим і матимуть вплив на загальну ймовірність. Існує простий трюк, щоб уникнути цієї проблеми. Ми можемо припустити, що наш навчальний набір настільки великий, що додавши до кожної змінної, значення котрої нам треба підрахувати, одиницю привнесе мізерні зміни в розрахунок ймовірностей, але допоможе оминати отримання нульової ймовірності.



ПРИКЛАД: ВИБІРКА ПОКУПЦІВ КОМП'ЮТЕРІВ

Вік	Дохід	Кредитний статус	Купити комп'ютер
Молодий	Високий	Так	Ні
Молодий	Високий	Ні	Ні
Дорослий	Високий	Так	Так
Пенсіонер	Середній	Так	Так
Пенсіонер	Низький	Так	Так
Пенсіонер	Низький	Ні	Ні
Дорослий	Низький	Ні	Так
Молодий	Середній	Так	Ні
Молодий	Низький	Так	Так
Пенсіонер	Середній	Так	Так
Молодий	Середній	Ні	Так
Дорослий	Середній	Ні	Так
Дорослий	Високий	Так	Так

ПРИКЛАД...

- Нехай $X = (\text{Вік} = \text{Молодий}, \text{Дохід} = \text{Середній}, \text{Кредитний рейтинг} = \text{Так})$
- $P(\text{Купити комп} = \text{Так}) = 9/14 = 0,643$
- $P(\text{Купити комп} = \text{Ні}) = 5/14 = 0,357$

- $P(\text{Вік} = \text{Молодий} \mid \text{Купити комп} = \text{Так}) = 2/9 = 0,222$
- $P(\text{Вік} = \text{Молодий} \mid \text{Купити комп} = \text{Ні}) = 3/5 = 0,6$

- $P(\text{Дохід} = \text{Середній} \mid \text{Купити комп} = \text{Так}) = 4/9 = 0,444$
- $P(\text{Дохід} = \text{Середній} \mid \text{Купити комп} = \text{Ні}) = 2/5 = 0,4$

- $P(\text{Кредитний рейтинг} = \text{Так} \mid \text{Купити комп} = \text{Так}) = 0,667$
- $P(\text{Кредитний рейтинг} = \text{Так} \mid \text{Купити комп} = \text{Ні}) = 0,4$



ПРИКЛАД...

- $P(X \mid \text{Купити комп} = \text{Так}) =$
 - $P(\text{Вік} = \text{Молодий} \mid \text{Купити комп} = \text{Так})^*$
 - $P(\text{Дохід} = \text{Середній} \mid \text{Купити комп} = \text{Так})^*$
 - $P(\text{Кредитний рейтинг} = \text{Так} \mid \text{Купити комп} = \text{Так})$
 - $= 0.222 * 0.444 * 0.667 = 0.066$
- $P(X \mid \text{Купити комп} = \text{Ні}) = 0,096$

*Шукаємо клас, який максимізує $P(X|C_i) * P(C_i)$*

- $P(X \mid \text{Купити комп} = \text{Так}) * P(\text{Купити комп} = \text{Так}) = 0.0424$
- $P(X \mid \text{Купити комп} = \text{Ні}) * P(\text{Купити комп} = \text{Ні}) = 0,034$



Переваги та недолітки

Переваги:

- Легкий в реалізації
- Потребує невеликої тренувальної виборки для навчання
- Дуже високі результати класифікації

Недоліки:

- Через припущення взаємонезалежності можлива втрата точності



Висновки

- ✓ Наївний Баєсовський метод є надзвичайно привабливим через його простоту та надійності.
- ✓ Це один з найстаріших алгоритмів класифікації і навіть у найпростішому вигляді він є на диво ефективним.
- ✓ Він широко використовується в таких областях, як класифікація тексту і фільтрації спаму.
- ✓ Існує велика кількість модифікацій даного алгоритму.

