

# ТЕМА 5

## КЛАСИФІКАЦІЯ НА ОСНОВІ ТЕОРІЇ БАЙЕСА

Надія І. Недашківська [n.nedashkivska@gmail.com](mailto:n.nedashkivska@gmail.com)

# Постановка задачі для байесівської класифікації

$X = \{x_1, x_2, \dots, x_m\}$  - незалежні змінні (атрибути)

$y \in Y$  - залежна змінна       $X \times Y$  - ймовірнісний простір з  
щільністю розподілу  $p(x, y)$

**Дано:** Множина об'єктів  $T = \{t_1, t_2, \dots, t_n\}$  навчальна вибірка  
 $t_i \rightarrow (X, y)$  (тестові приклади)

**Знайти:**

функцію класифікації  $f : X \rightarrow Y$  з мінімальною імовірністю помилки

$$p(x, y) = p(x)p(y | x) = p(y)p(x | y)$$

**Максимум апостеріорної імовірності (maximum a posteriori, MAP)**

$$f(x) = \operatorname{argmax}_{y \in Y} p(y | x) = \operatorname{argmax}_{y \in Y} \frac{p(x | y)p(y)}{p(x)} = \operatorname{argmax}_{y \in Y} p(x | y)p(y)$$

# Оптимальний байесівський класифікатор

Області, отримані в результаті класифікації:

$$X_s = \{x \in X \mid f(x) = s\}, \quad s \in Y$$

**Помилка:** об'єкт  $x$  класу  $y$  потрапляє в  $X_s$ ,  $s \neq y$ .

**Штраф за помилку** відомий:  $u_{ys} \geq 0$ ,  $\forall (y, s) \in Y \times Y$

**Середній ризик:** математичне сподівання штрафу для  $f$ :

$$R(f) = \sum_{y \in Y} \sum_{s \in Y} u_{ys} p(X_s, y)$$

Якщо відомі  $p(y)$ ,  $p(x \mid y)$ , то мінімальний середній ризик має байесівський класифікатор\*

$$f(x) = \arg \min_{s \in Y} \sum_{y \in Y} u_{ys} p(x \mid y) p(y)$$

\* <http://www.machinelearning.ru/>

# Оптимальний байесівський класифікатор

Якщо відомі  $p(y)$ ,  $p(x|y)$ ,  $u_{ys} = u_y$ ,  $u_{yy} = 0$ , то мінімальний середній ризик має класифікатор\*

$$f(x) = \arg \max_{y \in Y} u_y p(x|y) p(y)$$

**Підзадача 1:** Дано: навчальна вибірка  $T = \{t_i\}$ ,  $t_i \rightarrow (X, Y)$ .

**Знайти:** оцінки для  $p(y)$ ,  $p(x|y)$ ,  $y \in Y$  за вибіркою.

**Підзадача 2:** Дано: апіорні імовірності  $p(y)$ , функції правдоподібності  $p(x|y)$ ,  $y \in Y$ .

**Знайти:** класифікатор  $f: X \rightarrow Y$ , який мінімізує середній ризик.

\* <http://www.machinelearning.ru/>

# Оцінювання $p(y)$ , $p(x|y)$ , $y \in Y$ за вибіркою

Оцінювання апріорних імовірностей:

$$\hat{p}(y) = \frac{|T_y|}{|T|} \quad T_y = \{x_i \in X : y_i = y\} \quad y \in Y$$

Оцінювання функцій правдоподібності  $p(x|y)$  :

- Параметричне оцінювання:

$$\hat{p}(x|y) = \varphi(x, y, \theta)$$

- Оцінювання (розділ) суміші розподілів:

$$\hat{p}(x|y) = \sum_{k=1}^q \omega_k \varphi(x, y, \theta_k) \quad q \ll n$$

- Непараметричне оцінювання

# Наївний байесівський класифікатор *Naive Bayes*

**Припущення:** нехай  $x_1, x_2, \dots, x_m$  незалежні у сукупності.

Відомі щільності  $p(x_1 | y), p(x_2 | y), \dots, p(x_m | y)$ ,  $y \in Y$ .

Тоді  $m$ -вимірна щільність – це добуток 1-вимірних щільностей:

$$p(x | y) = p(x_1 | y) \cdot p(x_2 | y) \cdot \dots \cdot p(x_m | y) \quad y \in Y$$

Тоді класифікатор з мінімальним середнім ризиком:

$$f(x) = \arg \max_{y \in Y} u_y p(x | y) p(y) = \arg \max_{y \in Y} u_y p(y) \prod_{i=1}^m p(x_i | y)$$

$$f(x) = \arg \max_{y \in Y} \left( \ln(u_y \hat{p}(y)) + \sum_{i=1}^m \ln \hat{p}(x_i | y) \right) \quad \begin{array}{l} \text{нехай} \\ u_y = 1 \end{array}$$

# Алгоритм *Naive Bayes*: приклад

Чи виграє “Динамо” при наступних умовах (вектор  $x$ ):

- $x_1$  = в гостях,
- $x_2$  = суперник нижче в турнірній таблиці,
- $x_3$  = температура в нормі,
- $x_4$  = дощу немає

Такого прикладу немає в навчальній вибірці.

$$f(x) = \arg \max_{y \in Y} p(y | x) = \arg \max_{y \in Y} p(y) \prod_{i=1}^m p(x_i | y)$$

Треба розрахувати дві апостеріорні імовірності класів:

$p(\text{перемога}=\text{так} \mid x)$

$p(\text{перемога}=\text{ні} \mid x)$

# Приклад: результати попередніх ігор

Де грає	Суперник	Температура	Дощ	Перемога
Вдома	Вище	Висока	Так	Ні
Вдома	Нижче	Норма	Ні	Так
В гостях	Нижче	Норма	Так	Так
В гостях	Нижче	Висока	Так	Ні
Вдома	Вище	Висока	Ні	Так
Вдома	Нижче	Висока	Так	Ні
В гостях	Нижче	Висока	Ні	Ні
В гостях	Вище	Норма	Ні	Так



# Алгоритм *Naive Bayes*: приклад

Шукані апостеріорні імовірності класів:

$$\begin{aligned} p(\text{перемога=так} \mid x) &= p(\text{перемога=так}) * \\ &\quad p(x_1=\text{в гостях} \mid \text{перемога=так}) * \\ &\quad p(x_2=\text{нижче} \mid \text{перемога=так}) * \\ &\quad p(x_3=\text{норма} \mid \text{перемога=так}) * \\ &\quad p(x_4=\text{ні} \mid \text{перемога=так}) / p(x) \end{aligned}$$

$$\begin{aligned} p(\text{перемога=ні} \mid x) &= p(\text{перемога=ні}) * \\ &\quad p(x_1=\text{в гостях} \mid \text{перемога=ні}) * \\ &\quad p(x_2=\text{нижче} \mid \text{перемога=ні}) * \\ &\quad p(x_3=\text{норма} \mid \text{перемога=ні}) * \\ &\quad p(x_4=\text{ні} \mid \text{перемога=ні}) / p(x) \end{aligned}$$

# Алгоритм *Naive Bayes*: приклад

**Функції правдоподібності**  $p(x_1 | y), p(x_2 | y), \dots, p(x_m | y), y \in Y$ :

$$p(x_1=\text{в гостях} \mid \text{перемога}=\text{так})=2/4$$

$$p(x_1=\text{в гостях} \mid \text{перемога}=\text{ні})=2/4$$

$$p(x_2=\text{нижче} \mid \text{перемога}=\text{так})=2/4$$

$$p(x_2=\text{нижче} \mid \text{перемога}=\text{ні})=3/4$$

$$p(x_3=\text{норма} \mid \text{перемога}=\text{так})=3/4$$

$$p(x_3=\text{норма} \mid \text{перемога}=\text{ні})=0/4$$

$$p(x_4=\text{ні} \mid \text{перемога}=\text{так})=3/4$$

$$p(x_4=\text{ні} \mid \text{перемога}=\text{ні})=1/4$$

**Апріорні імовірності класів**  $p(y), y \in Y$ :

$$p(\text{перемога}=\text{так})=4/8$$

$$p(\text{перемога}=\text{ні})=4/8$$

# Алгоритм *Naive Bayes*: приклад

Шукані апостеріорні імовірності класів:

$$p(\text{перемога=так} \mid x) = 2/4 * 2/4 * 3/4 * 3/4 * 4/8 / p(x)$$

$$p(\text{перемога=ні} \mid x) = 2/4 * 3/4 * 0 * 1/4 * 4/8 / p(x)$$

$$p(\text{перемога=так} \mid x) > p(\text{перемога=ні} \mid x)$$

**Тому при заданих умовах матч буде виграно.**