

Statystyka i analiza danych 2025/2026

Lista ćwiczeń

Lista 10

Ćwiczenie 59 — Ślad macierzy kwadratowej $A \in \mathbb{R}^{n \times n}$ definiujemy jako: $\text{tr}(A) = \sum_{i=1}^n A_{ii}$.

(a) Udowodnij, że dla macierzy $A \in \mathbb{R}^{m \times n}$ oraz $B \in \mathbb{R}^{n \times m}$ zachodzi

$$\text{tr}(AB) = \text{tr}(BA).$$

(b) Korzystając z (a), pokaż własność cykliczności śladu: dla macierzy A, B, C odpowiednich wymiarów zachodzi

$$\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA).$$

Czy $\text{tr}(ABC) = \text{tr}(BAC)$? Jeśli nie, podaj kontrprzykład.

Ćwiczenie 60 — W oparciu o ćwiczenie 53 dla modelu regresji liniowej $\mathbf{y} = X\boldsymbol{\theta} + \mathbf{e}$ wartości $\hat{\mathbf{y}}$ dopasowane przez model możemy przedstawić jako

$$\hat{\mathbf{y}} = X(X^T X)^{-1} X^T \mathbf{y} = H\mathbf{y},$$

gdzie $H = X(X^T X)^{-1} X^T$ (ang. hat matrix), $X \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^n$. *Leverage statistic* i -tej obserwacji definiuje się jako $h_{ii} = H_{ii}$ (element na przekątnej H).

(a) W oparciu o to, że H jest symetryczna ($H^T = H$) i idempotentna ($H^2 = H$), udowodnij, że $0 \leq h_{ii} \leq 1$.

(b) Korzystając z własności śladu udowodnionych w poprzednim zadaniu, oblicz $\text{tr}(H)$, a następnie pokaż, że średnia wartość *leverage statistic* wynosi

$$\bar{h} = \frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{p}{n}.$$

(c) Podaj intuicyjne uzasadnienie, dlaczego obserwacja z $h_{ii} \gg \bar{h}$ może być uważana za *wpływową* w kontekście dopasowania modelu.

Ćwiczenie 61 — Wyjaśnij jaka jest intuicja związana z odległością Mahalanobisa (zobacz np. [tu](#) i [tu](#)) i jaki jest związek tej odległości z *leverage statistic*.

Ćwiczenie 62 — Wyjaśnij, czym jest [naiwny klasyfikator Bayes'a](#) i kiedy warto go zastosować. Wyjaśnij i udowodnij następujący wzór:

$$p(C_k | x_1, \dots, x_n) = \frac{1}{p(\mathbf{x})} p(C_k) \prod_{i=1}^n p(x_i | C_k),$$

gdzie

$$p(\mathbf{x}) = \sum_k p(\mathbf{x} | C_k) p(C_k) \quad \text{a} \quad \mathbf{x} = (x_1, \dots, x_n).$$

Ćwiczenie 63 — Rozważ problem klasyfikacji wiadomości e-mail na dwie klasy: Spam (C_1) lub Nie-spam (C_2). Dysponujesz dwoma cechami binarnymi:

- x_1 opisuje czy e-mail zawiera słowo *promocja*: $x_1 = 1$ jeśli tak, $x_1 = 0$ jeśli nie,
- x_2 opisuje czy e-mail zawiera słowo *oferta*: $x_2 = 1$ jeśli tak, $x_2 = 0$ jeśli nie.

Klasa Spam zawiera 100 e-maili i każdy z nich zawiera oba słowa: $x_1 = 1$ i $x_2 = 1$. Klasa Nie-spam zawiera 900 e-maili:

- 450 e-maili zawiera tylko słowo *promocja*: $x_1 = 1, x_2 = 0$,

- 450 e-maili zawiera tylko słowo *oferta*: $x_1 = 0, x_2 = 1$,

Korzystając z naiwnego klasyfikatora Bayesa (zakładającego niezależność cech), oblicz prawdopodobieństwa $p(C_1|x)$ i $p(C_2|x)$ dla nowego e-maila, który zawiera oba słowa: $x_1 = 1, x_2 = 1$. Czy wynik jest zaskakujący? Wyjaśnij.

Ćwiczenie 64 — Ocena modelu w problemie klasyfikacji.

- Przedstaw i wyjaśnij metryki często stosowane do oceny modelu w problemie klasyfikacji: *accuracy*, *precision*, *recall* oraz *F1*. Podaj prosty przykład.
- Czym się różni ich wersja *micro* od wersji *macro* (zobacz np. [tutaj](#))? Podaj prosty przykład.
- Jak te metryki mają się do bardziej potocznego rozumienia pojęć *accuracy* i *precision*?

cdn.

J.L.