

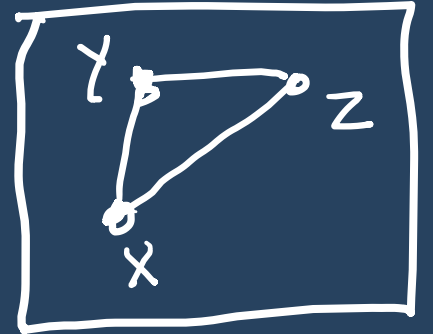
METODA MIN-HASH

① Prz. metryczna : (X, d) , $d: X \times X \rightarrow [0, \infty)$

a. $(d(x, y) = 0) \equiv (x = y)$

b. $d(x, y) = d(y, x)$

c. $d(x, z) \leq d(x, y) + d(y, z)$



Ⓟ \mathbb{R}^n ; $d_e(x, y) = \left(\sum_{l=1}^n (x_l - y_l)^2 \right)^{\frac{1}{2}}$

od. euklidesowa

Ⓟ $C_n = \{0, 1\}^n$; $d_H(x, y) = \sum_{l=1}^n |x_l - y_l|$ $d_H \leftarrow$ od. Hamminga


zadanie: d_H jest od.

$$d_H(x, y) = \sum_{l=1}^n \underbrace{|x_l - y_l|}_{\in \{0, 1\}} = |\{i: x_i = 1 \wedge y_i = 0\}| + |\{i: x_i = 0 \wedge y_i = 1\}|$$

$x, y \in \{0, 1\}^n$

$$X = \{i: x_i = 1\}$$

$$Y = \{i: y_i = 1\}$$

$$\begin{aligned} |X \setminus Y| + |Y \setminus X| &= |X \Delta Y| \\ &= |X \Delta Y| \end{aligned}$$


wniosek:

$$(\mathcal{P}(\Omega), \hat{d})$$

Ω \uparrow
ub. skończ

$$\hat{d}(x, y) = |X \Delta Y|$$

\mathbb{R} prz. metryczna



• modyfikacje metryk : (X, d) - fix.

▶ $\alpha > 0$: $d^*(x, y) = \alpha \cdot d(x, y)$

$(X, \alpha \cdot d)$ - p.metr.

▶ Tw. 2.11.1. Niech (X, d) - p.metr. Niech $f: [0, \infty) \rightarrow [0, \infty)$ będzie taka, że

• $f(0) = 0$

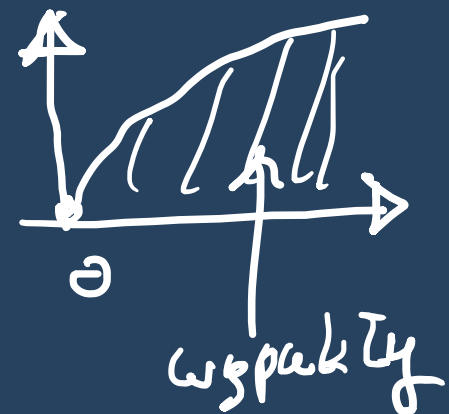
• f jest niemalejąca ($x \leq y \rightarrow f(x) \leq f(y)$)

• f jest WKLĘSKA.

niech

$$\rho(x, y) = f(d(x, y)).$$

wtedy (X, ρ) jest p. metryczna.



Zadanie
z listy.

Wnoszele: (X, d) - prz. metryczna

$$f(x) = \min\{x, 1\}$$

$$\tilde{d}(x, y) = f(d(x, y)) =$$
$$\uparrow = \min\{d(x, y), 1\}$$

metryka

PP

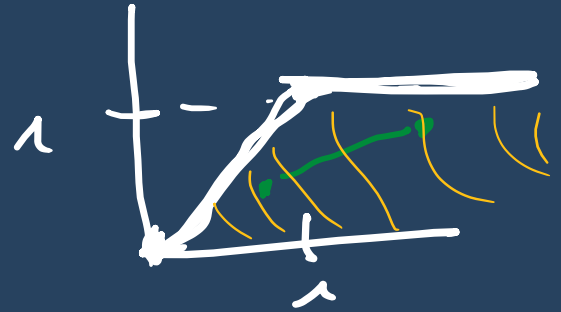
$$g(x) = \sqrt{x}$$

d -metryka

$\tilde{d}(x, y) = \sqrt{d(x, y)}$ - to też
jest metryką

$$g'(x) = (x^{\frac{1}{2}})' = \frac{1}{2} x^{-\frac{1}{2}} > 0$$
$$g''(x) = -\frac{1}{4} x^{-\frac{3}{2}} < 0$$

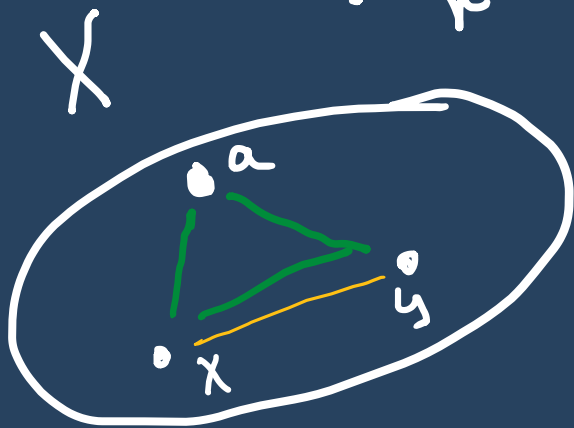
$x > 0$



► Tw (Steinhaus) $\text{w} \bar{X}$. ie (X, d) jest p. metr.
 Niech $a \in X$. Na $X \setminus \{a\}$ określamy
 funkcję

$$\rho(x, y) = \frac{2 \cdot d(x, y)}{d(x, a) + d(y, a) + d(x, y)}$$

wtedy ρ też jest p. metryką (na $X \setminus \{a\}$).



Uwaga:

$$\underline{d(x, a) + d(a, y) + d(x, y)} \geq d(x, y) + d(x, y) = 2 \cdot d(x, y)$$

$$\text{ zatem } \rho(x, y) \leq \frac{2 \cdot d(x, y)}{2 \cdot d(x, y)} = 1.$$

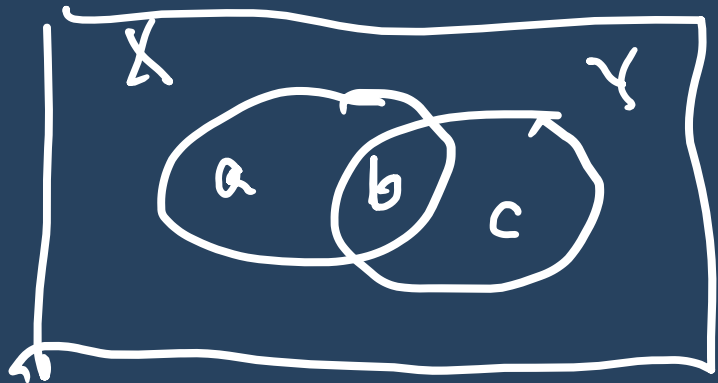
\square -d ← zadanie z listy.

Zaśt. tw. Steinhause

Ustawamy Ω . Dla $X, Y \subseteq \Omega$: $d(X, Y) = |X \Delta Y|$.

W tw. Steinh. na Ω podstawiamy ϕ .

$$p(X, Y) = \frac{2 \cdot |X \Delta Y|}{|X \Delta \phi| + |Y \Delta \phi| + |X \Delta Y|} = \frac{2 \cdot |X \Delta Y|}{|X| + |Y| + |X \Delta Y|}$$



$$\begin{aligned} a &= |X \setminus Y| \\ b &= |X \cap Y| \\ c &= |Y \setminus X| \end{aligned}$$

$$\begin{aligned} &= \frac{2 \cdot |X \Delta Y|}{(a+b) + (b+c) + (a+c)} \\ &= \frac{2 \cdot |X \Delta Y|}{2 \cdot (a+b+c)} = \frac{|X \Delta Y|}{|X \cup Y|} \end{aligned}$$

CZTLI : na $\xi \in \mathcal{P}(\Omega) \setminus \emptyset$ mamy
określoną odległość

$$d_J(X, Y) = \frac{|X \Delta Y|}{|X \cup Y|}$$

ODLEGŁOŚĆ JACCARD'a.

2. Podobieństwo obiektów. // similarity

$$(X, s) \quad s: X \times X \rightarrow [0, 1]$$

- $s(X, X) = 1$

- $s(X, Y) = s(Y, X)$

Idea: jeśli $s(X, Y) \approx 1$ to X jest bardzo podobne do Y .

Dobre podobieństwo:

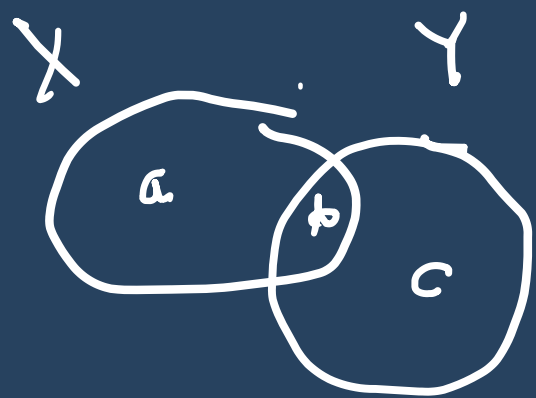
takie s , że $1 - s$ jest metryką
(ograniczone przez 1)

$$d(X, Y) = 1 - s(X, Y) \leftarrow \text{metryka}$$

① $(P^+(\Omega), d_f)$

$s(x, y) = 1 - d_f(x, y)$ ← podobieństwo

$$s(x, y) = 1 - \frac{|x \Delta y|}{|x \cup y|} = \frac{|x \cap y| - |x \Delta y|}{|x \cup y|}$$



$$= \frac{(a+b+c) - (a+c)}{|x \cup y|} = \frac{b}{|x \cup y|}$$

$$s(x, y) = \frac{|x \cap y|}{|x \cup y|}$$

podobieństwo Jaccarda

(P)

Ω - zbiór słów, D_1, D_2 - dokumenty

$D_1^*, D_2^* \leftarrow$ słowa występujące w D_i .

$$J(D_1^*, D_2^*) = \frac{|D_1^* \cap D_2^*|}{|D_1^* \cup D_2^*|}$$

• $D \subseteq (\Sigma)^k$ $D = (i_1 i_2 i_3 \dots)$
 Σ -ascii k -shungles (k -ramiennie)

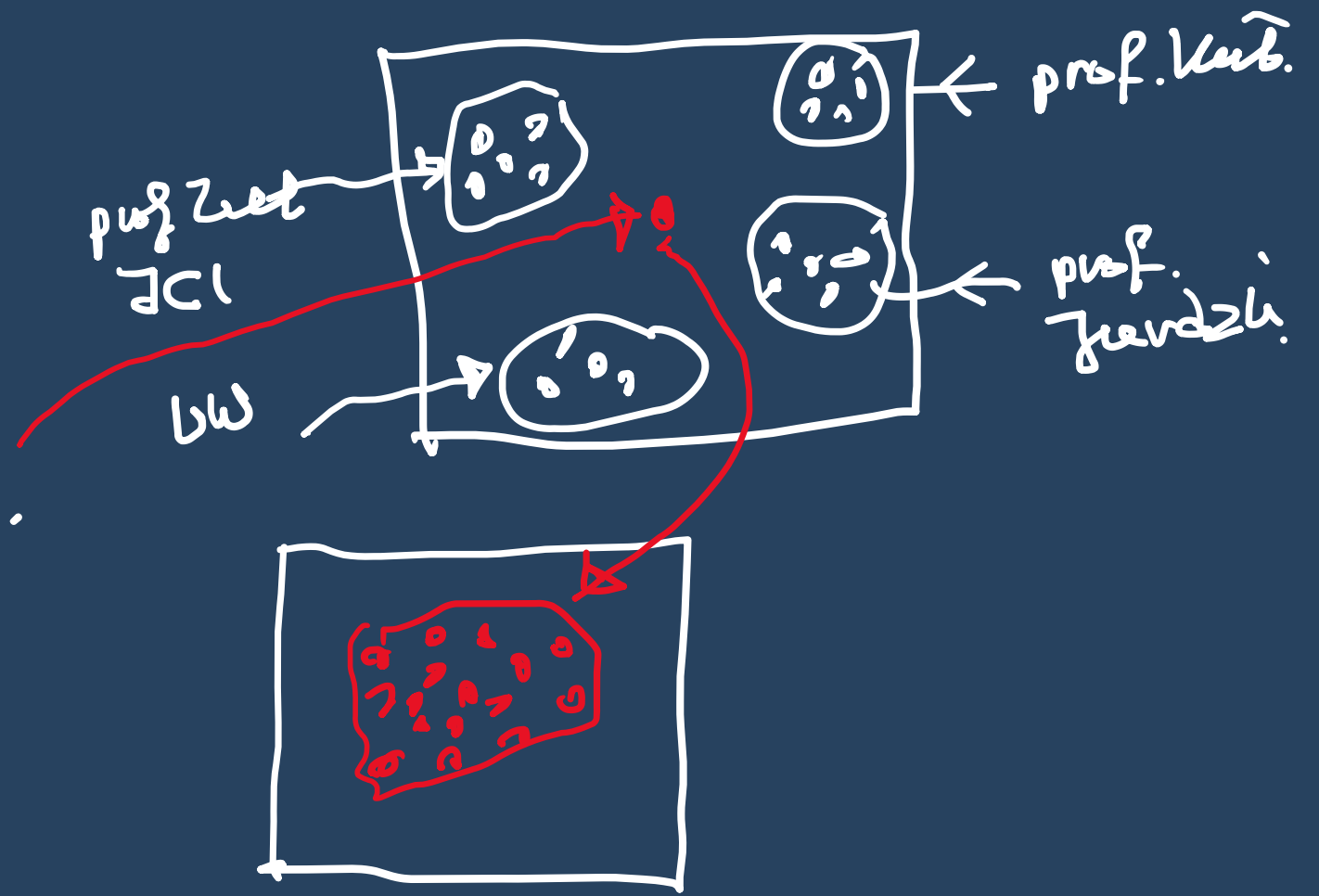
$D =$  $\} \text{-sh}(D) = \{i_1 i_2 i_3, i_2 i_3 i_4, i_3 i_4 i_5, \dots\}$

PODOB: $J(\text{3-sh}(X), \text{3-sh}(Y)) \leftarrow$ DUŻA SIŁA KLASYFIK.

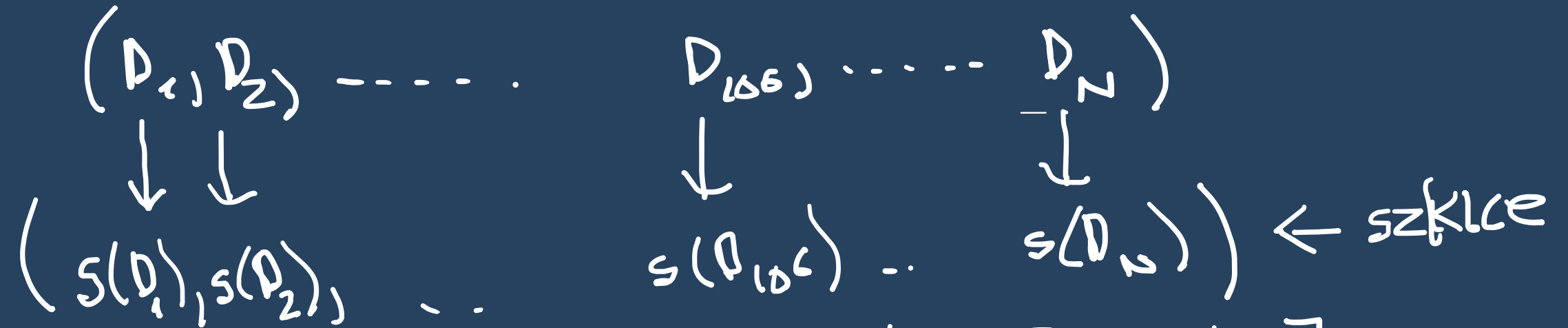
PRZYKŁAD - $\mathcal{D} = \{D_1, \dots, D_{500}\}$ - pluce w pdf z informatyki

$\delta(D_i^*, D_j^*) \quad i \neq j$
 \Rightarrow klasteryzacji

$\mathcal{D} \cup$ kolekcja elementów nauki doświ.
100



JAK WYZNACZAĆ EFEKTYWNIĘ $f(x, y)$?

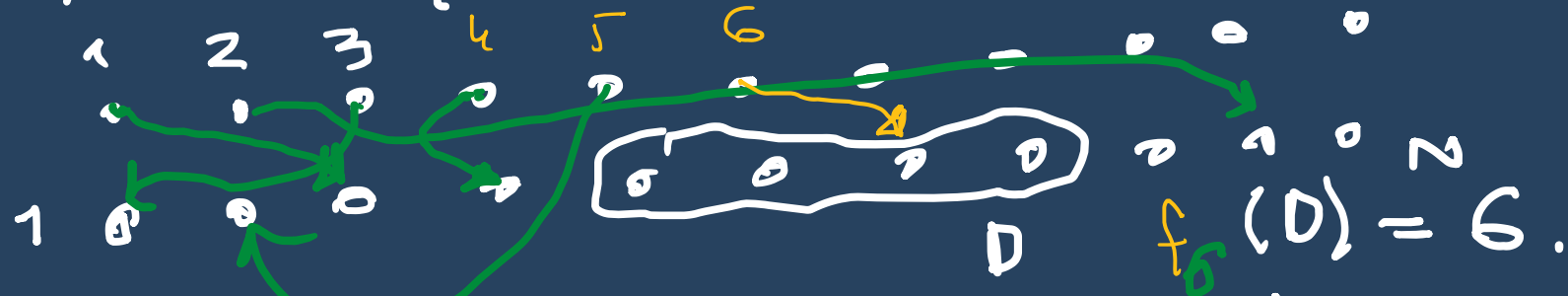


• mają to być wektory $[n_1, \dots, n_k]$
k-nierobyt dane (np. 100, ..., 1000)

Myślimy, że $\mathcal{D} \subseteq \{\omega_1, \omega_2, \dots, \omega_N\}$
 $\subseteq \{1, 2, \dots, N\} = \Omega$.

• Dla $\sigma \in \text{Sym}(\Omega)$ definiujemy

$$f_\sigma(\mathcal{D}) = \min \{k : \sigma(k) \in \mathcal{D}\}$$



• Na $\text{Sym}(\Omega)$ rozważamy rozkład jednostajny;

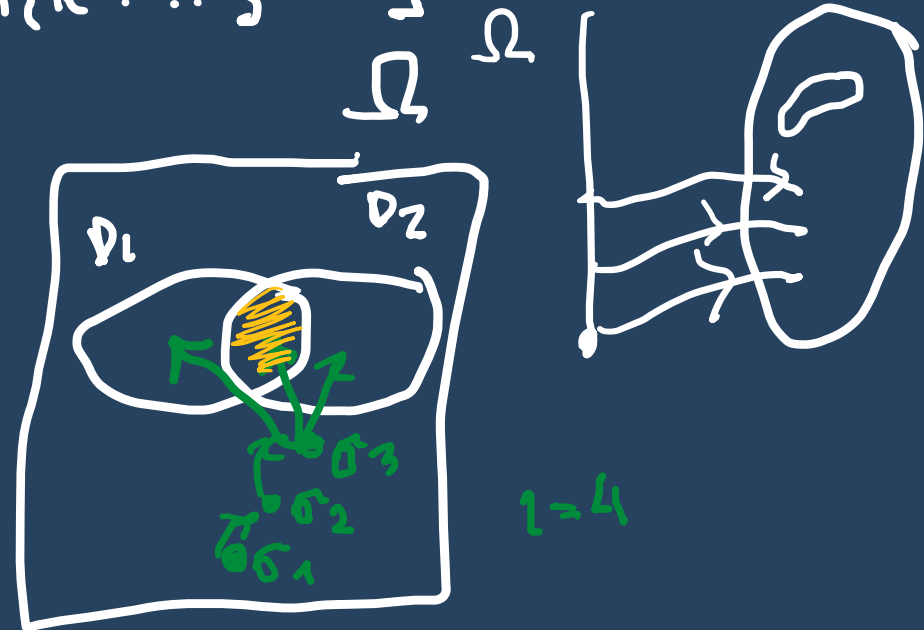
$$A \subseteq \text{Sym}(\Omega): \Pr[A] = \frac{|A|}{|\Omega|!}.$$

Ustawa $D_1, D_2 \subseteq \Omega$, $D_1, D_2 \neq \emptyset$.

$$\begin{aligned} & \Pr[h_\sigma(D_1) = h_\sigma(D_2)] = \\ &= \sum_{i=1}^N \Pr[h_\sigma(D_1) = h_\sigma(D_2) \mid \min\{k: \sigma(k) \in D_1 \cup D_2\} = i] \cdot \Pr[\min\{k: \dots\} = i] \\ &= \sum_{i=1}^N \frac{|D_1 \cap D_2|}{|D_1 \cup D_2|} \cdot \Pr[\min\{k: \dots\} = i] \end{aligned}$$

$N = |\Omega|$

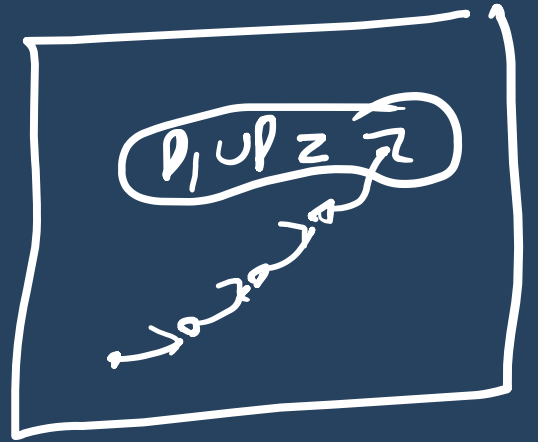
$h_\sigma(D) = \min\{k: \sigma(k) \in D\}$



$2 \rightarrow 4$

$$= \frac{|D_1 \cap D_2|}{|D_1 \cup D_2|} \sum_{i=1}^N \Pr [\min \{k : \sigma(k) \in D_1 \cup D_2\} = i]$$

$$= f(D_1, D_2)$$



$$\Pr [h_G(D_1) = h_G(D_2)] = f(D_1, D_2)$$

Q

: jak wylosować $\delta \in \text{Sym}(L)$?

! czym to ^{wzrosty} zastąpić!

Porządek $F \subset \mathbb{R}$

Haszujące \mathbb{C} .