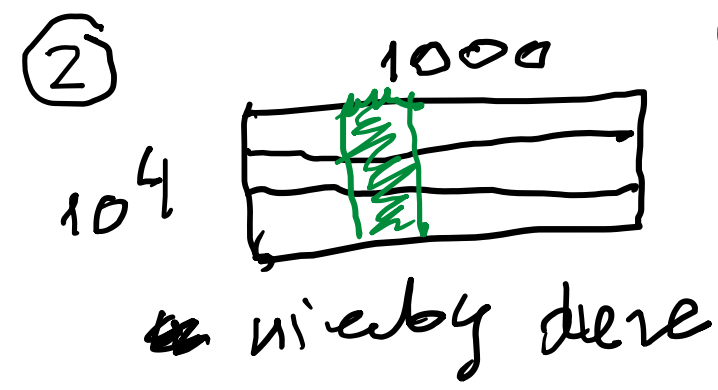
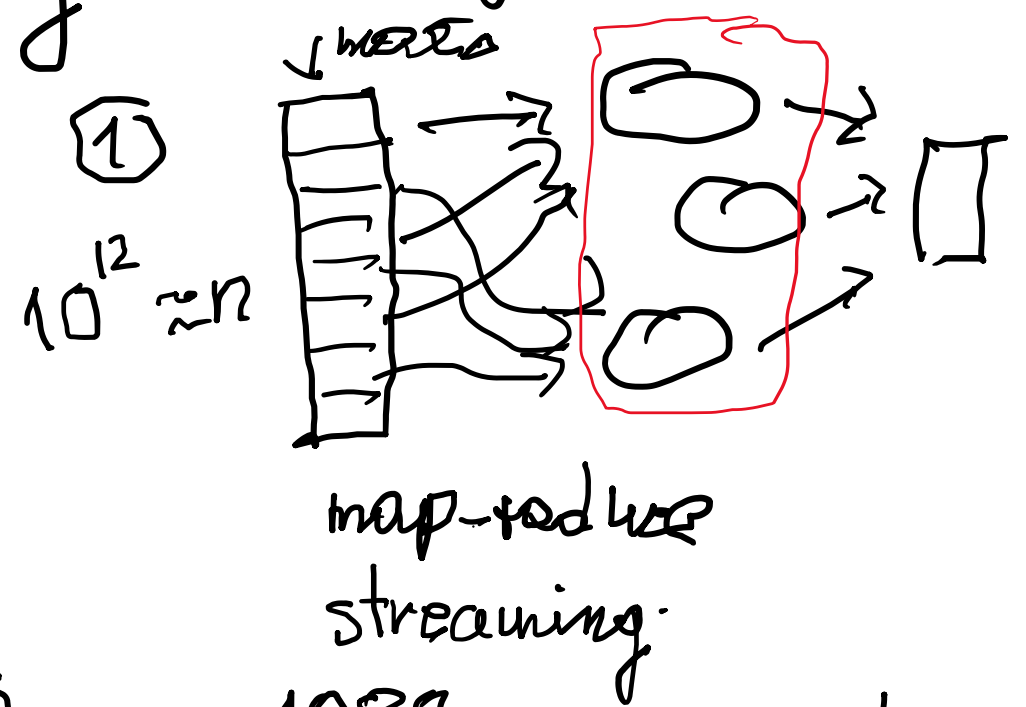
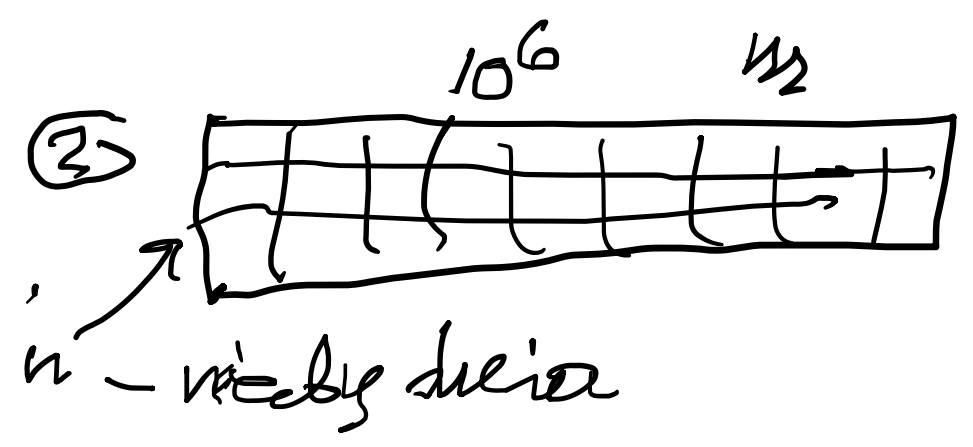


Big Data



Koszyk,
zakupów
 1000
 2^{1000} - podz



$A \in \mathbb{R}^n$, m -dim
dimension reduction

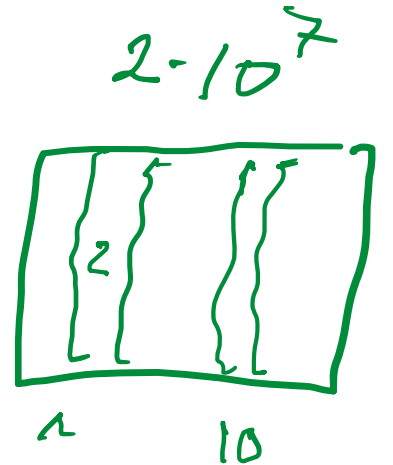


Koincydencje

① Konek pracy : $20 \cdot 10^5$ - stanowisk

długość pracy : 10 lat

ile ludzi zsumę pracy w czasie $\frac{20 \cdot 10^5}{10} = 2 \cdot 10^6$



ile jest węgłowodanów : 10^6 rok

$$P = \frac{25 \cdot 10^3}{20 \cdot 10^6} \approx \frac{1}{10^4} \quad \frac{10^6}{4 \cdot 10^2} = \frac{1}{4} \cdot 10^4 = 0.25 \cdot 10^4 \text{ ~~rocznik~~ ^{skocznic} }$$

$q =$ pr. zobaczenie ~~człowieka~~ ^{bestii} $\approx \frac{1}{10^2}$

wit, uwal :

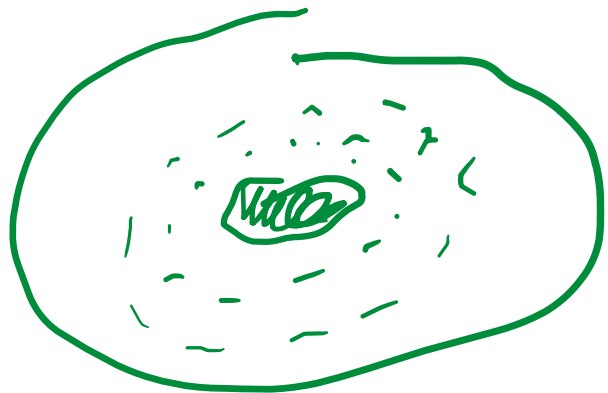
$$P[\text{kot } \wedge \text{wynarciarz}] = \frac{1}{10^4} \cdot \frac{1}{10^2} = \frac{1}{10^6}$$

ile osób duńwie?

$$20 \cdot 10^6 \cdot \frac{1}{10^6} \approx 20 \text{ osób}$$

ile wargczue :

$$20 \cdot 30 \approx 600$$



1950-70 : paraps. Rhöne (E)



próba odgromień, sekce

• osoby : $\approx 10^4$

• efekt : 10 osób bezładnych

• podwójny eksp : 1 osoba

$$\textcircled{1} P[\text{sukces}] = \left(\frac{1}{2}\right)^{10} \approx \frac{1}{1000}$$

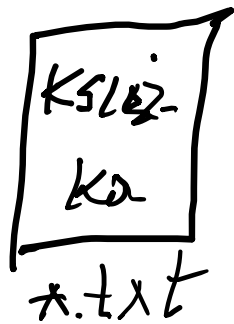
$$E[L_1] = 10^4 \cdot 10^{-3} = 10$$

$$E[L_2] = 10 \cdot 10^3 = \frac{1}{100}$$

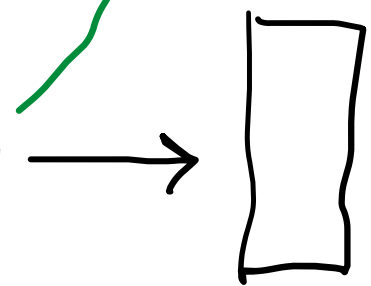
PIERWSZY PROGRAM

Hello world dla BD:

Analiza tekstu:



ile słów



100 najczęstszych słów

chlup, stolo



ktoś, albo coś, ...

, - często

1) kodowanie → UTF-8

2) UPPER-CASE

3) eliminacja "STOP-WORDS"

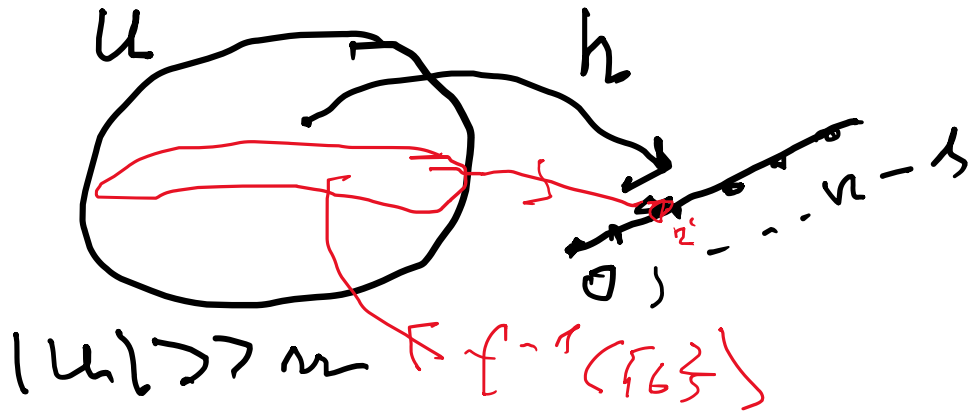
4) usunięcie słów: !, ,, ,

(ala, pola, ala, mes, ala) → groupBy [[ala, 35], [pola, 1], ...]

filterNot

sortBy

Hash - functions



n - dieser Rest

$$a_0 a_1 a_2 \dots a_m \rightarrow c \text{ (e.g. } 11)$$

$$\sum_{i=1}^m a_i \cdot 256^i$$

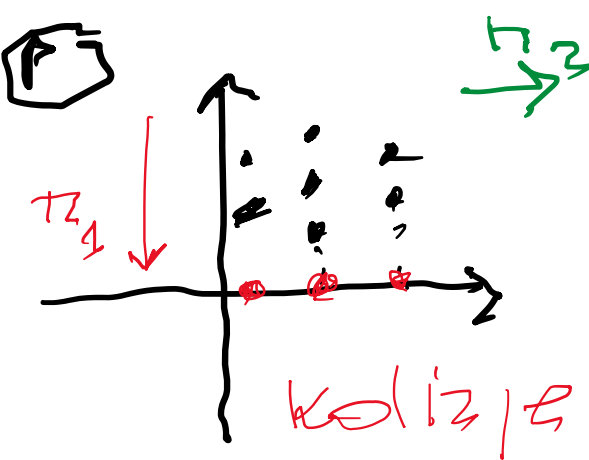
• $h: U \rightarrow [n]$

$\rightarrow (\exists i) |h^{-1}(\{z\})| \geq \frac{|U|}{n}$

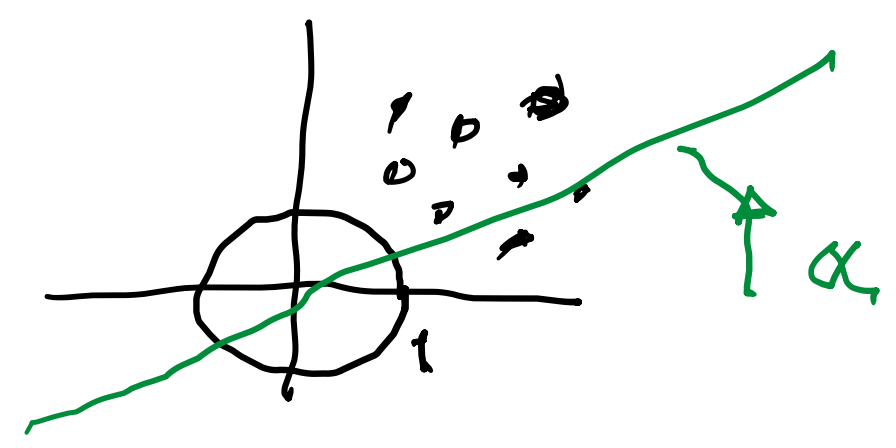
$U = \bigcup_{z=0}^{n-1} h^{-1}(\{z\})$; $\text{golgebly } |h^{-1}(\{z\})| < \frac{|U|}{n}$

$|U| = \sum_{z=0}^{n-1} |h^{-1}(\{z\})| \leq \frac{|U|}{n} \cdot n = |U|$

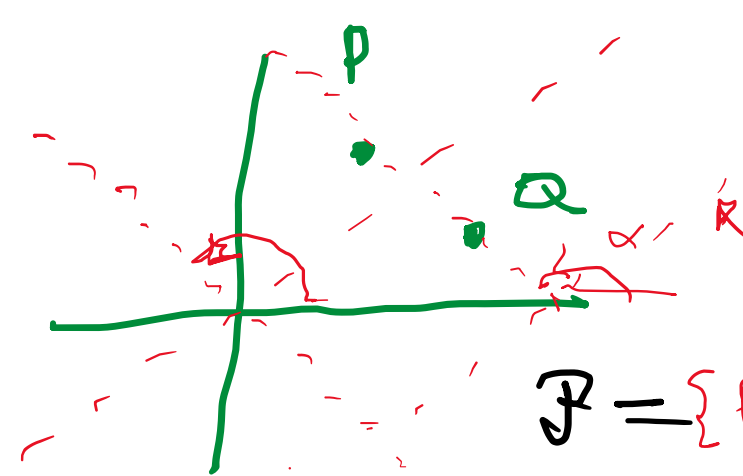
$x \neq y \wedge h(x) = h(y) \leftarrow \text{KOLLISION}$



brak kolizije



$\alpha \in [0, \pi)$



$\pi_\alpha(x, y) = \cos\alpha \cdot x + \sin\alpha \cdot y$

$\beta_{PQ} \quad \beta_{QP} = \text{"zty kat"}$

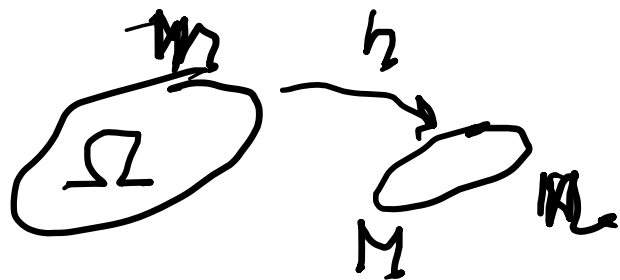
$P = \{P_1, \dots, P_m\} : 0 \leq j \leq m$

$\text{skor} \rightarrow \{ \beta_{P_i P} : 0 \leq i < j \leq m \} : \binom{m}{2}$

WYGENERUJ LOSOWO $\alpha \in [0, \pi)$; $\omega \in \mathbb{Z} \pi_\alpha$

$P[\pi_\alpha \text{ ma kolizje dla}] = 0$

UNIVERSAL HASHING



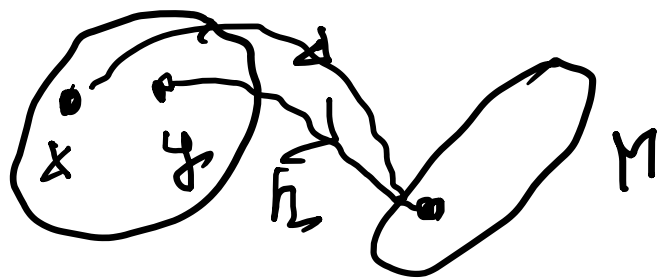
$$\mathcal{H} = \dot{M}^{\Omega} ; \text{pr. prob.}$$

total, good.

$$P(\omega) = \frac{1}{|\mathcal{H}|}$$

$$\omega \in \mathcal{H}$$

$$P[h(x) = h(y)] = \frac{|\{h \in \mathcal{H} : h(x) = h(y)\}|}{|\mathcal{H}|}$$



Fix x, y ; $x \neq y$; $LDSUJEM$
 $h \in \mathcal{H}$

$$\{h : h(x) = h(y)\} = \bigcup_{a \in M} \{h : h(x) = h(y) = a\}$$

$$L(x) = \sum_{a \in M} |\mathcal{H}|^{\Omega-2} = |\mathcal{H}|^{\Omega-1}$$

$$= \frac{|\mathcal{H}|^{\Omega-1}}{|\mathcal{H}|^{\Omega-1}} = \frac{1}{|\mathcal{H}|}$$

DEF. $\mathcal{H} \subseteq \mathcal{M}^\Omega$ jet univ. hash' e je co

(L)

$$(\forall x, y \in \Omega) (x \neq y \rightarrow \Pr_{\mathcal{H}}[h(x) = h(y)] \leq \frac{1}{|\mathcal{M}|})$$