

BIG DATA :

- dane, które nie mieszczą się w pamięci
- średnie dane, ale prosty sposób nie mieści się w pamięci

$$|P(\Omega)| = 2^{|\Omega|}$$



¹² "Hello World" dla Big Data

word-count problem

$X =$

| | | | | |
|----|-------|-------|-----|-----|
| ab | label | label | ccc | ... |
|----|-------|-------|-----|-----|

 tekst

- 1) podziel X na słowa
- 2) wyznacz częstotliwość wszystkich słów
- 3) wyznacz z tego N najczęściej występujących słów ($N \approx 100$)

zadanie: (4) wygenerować z (3) "word cloud"

algorytm Big Data

$Z = ["ala", "ka", "kota", "i", "ala", "na", "psie"]$

$W = [("ala", 2), ("ka", 2), \dots]$

Python: GroupBy ← lista posortowana

TF-IDF

$\mathcal{D} = \{D_1, D_2, \dots, D_n\}$ ← kolekcje $\neq \emptyset$ dokumentów

$T =$ zbiór słów wyst. w \mathcal{D} .

Term frequency: $tf(w, d) = \frac{n_{w,d}}{\sum_{d \in T} n_{d,d}}$

$n_{w,d} =$ liczba występień.

$$d = [w_1, w_2, \dots, w_n]$$

$$n_{w,d} = d.\text{count}(w) :$$

$tf(w, d) =$ wyst w w dok d .

$$idf(t, \mathcal{D}) = \log_2 \left(\frac{|\mathcal{D}|}{|\{d \in \mathcal{D} : t \in \mathcal{D}\}|} \right)$$

inverse
doc.
freq

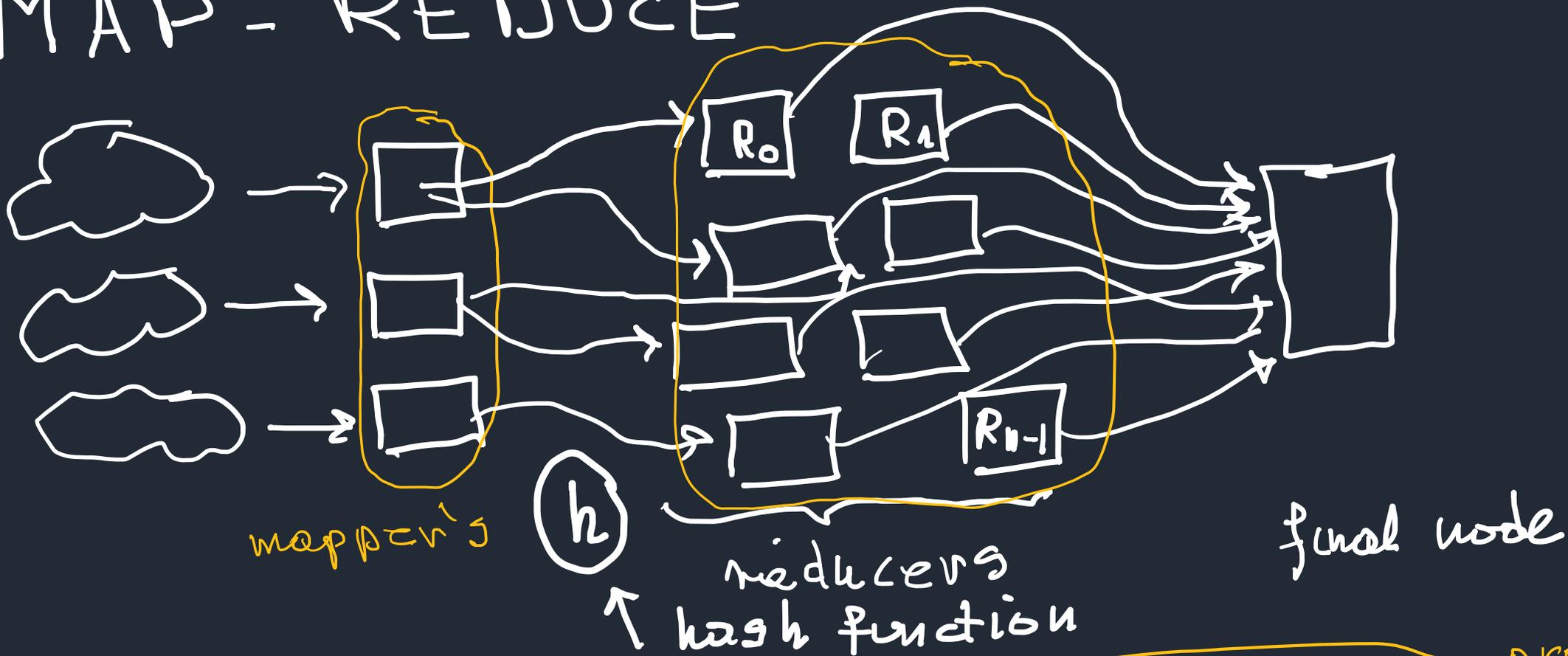
$$1 = \frac{|\mathcal{D}|}{|\mathcal{D}|} \leq \frac{|\mathcal{D}|}{|\{d : t \in d\}|} \leq \frac{|\mathcal{D}|}{1} = |\mathcal{D}|$$

↑
≠ 0

$$0 \leq \text{idf}(t, \mathcal{D}) \leq \log_2 |\mathcal{D}|$$

$$t f i d f (t, d, \mathbb{D}) = t f (t, d) \cdot i d f (t, \mathbb{D})$$

MAP - REDUCE



MAPPER:

on get x :

k - klucz wyzku z x
 y - wartość wyzku z x

output (k, y)

programista

$$h: \Sigma \longrightarrow \{0, \dots, n-1\}$$

$$(k, y) \rightsquigarrow R_h(k)$$

PO ZAKOŃCZENIU PRACY PRZEZ MAPPERY

$$1) (k_1, y_1), (k_2, y_2), (k_3, y_3), \dots$$

$$\left[(k^1, [y_{k_1}^1, \dots, y_{k_{k_1}}^1]), (k^2, [y_{k_1}^2, \dots, y_{k_{k_2}}^2]), \dots \right]$$

$$[(a, \underline{x}), (b, \underline{y}), (a, \underline{z}), (a, \underline{u}), (b, \underline{w}), (c, \underline{v})]$$

$$[(a, [\underline{x}, \underline{z}, \underline{u}]), (b, [\underline{y}, \underline{w}]), (c, [\underline{v}])]$$

2) dla każdej pary (k, L)
wykonywane jest obliczenie
i wyznaczony jest liczbowy wynik).

przebieg

PRZYKŁAD : Mnożenie macierzy
przez wektor

$$\begin{bmatrix} m_{11} & m_{12} & \dots & m_{1n} \\ \vdots & \vdots & & \vdots \\ m_{w1} & m_{w2} & \dots & m_{wn} \end{bmatrix} \cdot \underbrace{\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}}_{\vec{x}}$$

Zał.: \vec{x}
mieści się
w pamięci
mapy wejściowej.

$$M \cdot \vec{x} = \vec{y}$$

INPUT : (i, j, m_{ij})

$$y_i = \sum_j m_{ij} \cdot x_j$$

```

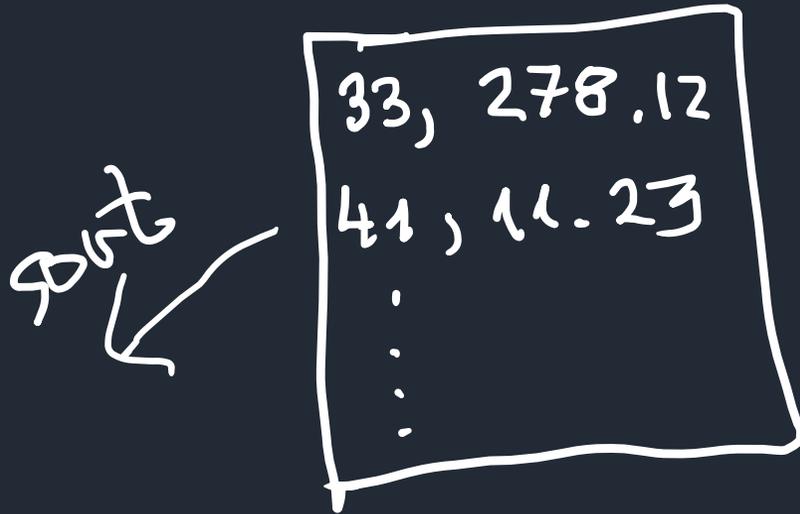
mapper(i, j, m) {
    emit(i, m * x_j);
}

```

```

reduce(i, L) {
    emit(i, sum(L));
}

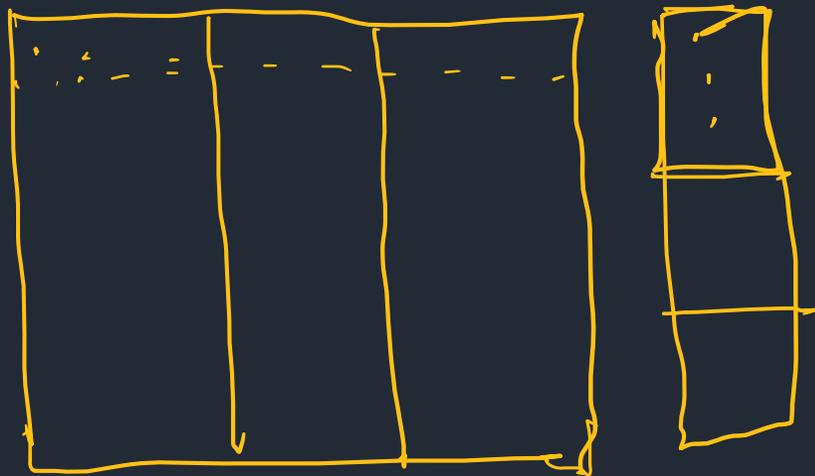
```



$(i, [m_{i_1} \cdot x_1, m_{i_2} \cdot x_2, \dots])$



zadanie: $\text{rot}(\bar{x})$ nie mieści się
w pamięci komputerowej.
ale $\frac{1}{3}|\bar{x}|$ mieści się



jak to zrobić?

Q

word - count.

```
mapper(w) {  
  emit(w, 1);  
}
```

```
reduce(w, L)  
  emit(w, length(L));  
}
```

reducer

(a, 1), (b, 1), (a, 1), (b, 1), (c, 1), ...



[(a, [1, 1, 1]), (b, [1, 1]), (c, [1])]

wyśnik : plik

| | |
|---|---|
| a | 3 |
| c | 1 |
| b | 2 |