

Big Data - zapisienia statystyczne

~ 1950-60 : badania zjawisk paranormalnych

10 000 uczestników:

odgrywanie ciągu 0/1 długości 10

• wynik I : ok 10 osób odgadły

• wynik II : poprawnie cały ciąg
nikt nie odgadł poprawnie

WYJAŚNIENIE : są osoby o zdol. paranormal.

ALE : tracę je gdy dowiedzą
się o uchu.

$$X_1, \dots, X_{10} \sim \text{Ber}\left(\frac{1}{2}\right) \quad P_{\theta}[X_L] = \begin{cases} \frac{1}{2} & ; X_L = 1 \\ \frac{1}{2} & ; X_L = 0 \end{cases}$$

$$Y = X_1 + \dots + X_n$$

$$\Pr[Y=10] = \left(\frac{1}{2}\right)^{10} \approx \frac{1}{1000}$$

$$Y_1, \dots, Y_{10000}$$

$$S_L = \mathbb{1}_{\{Y_L=10\}}$$

$$L = S_1 + \dots + S_{10000}$$

$$E[L] = 10000 \cdot \Pr[Y_L=10] \approx \frac{10000}{1000} \approx 10$$

□

P.

$\underbrace{011011011011 \dots}_{n \approx 10^7}$

java jest słabo. $\frac{0}{0}$ zadanie

cał: to jest
cał losowy dr

Q. Ile ~~osób~~ osób dziennie

w Polsce równo widzi czarnego kota

a potem (w ciągu dnia) traci robotę? $\frac{2}{0}$

MAP - REDUCE

① inverted index of word

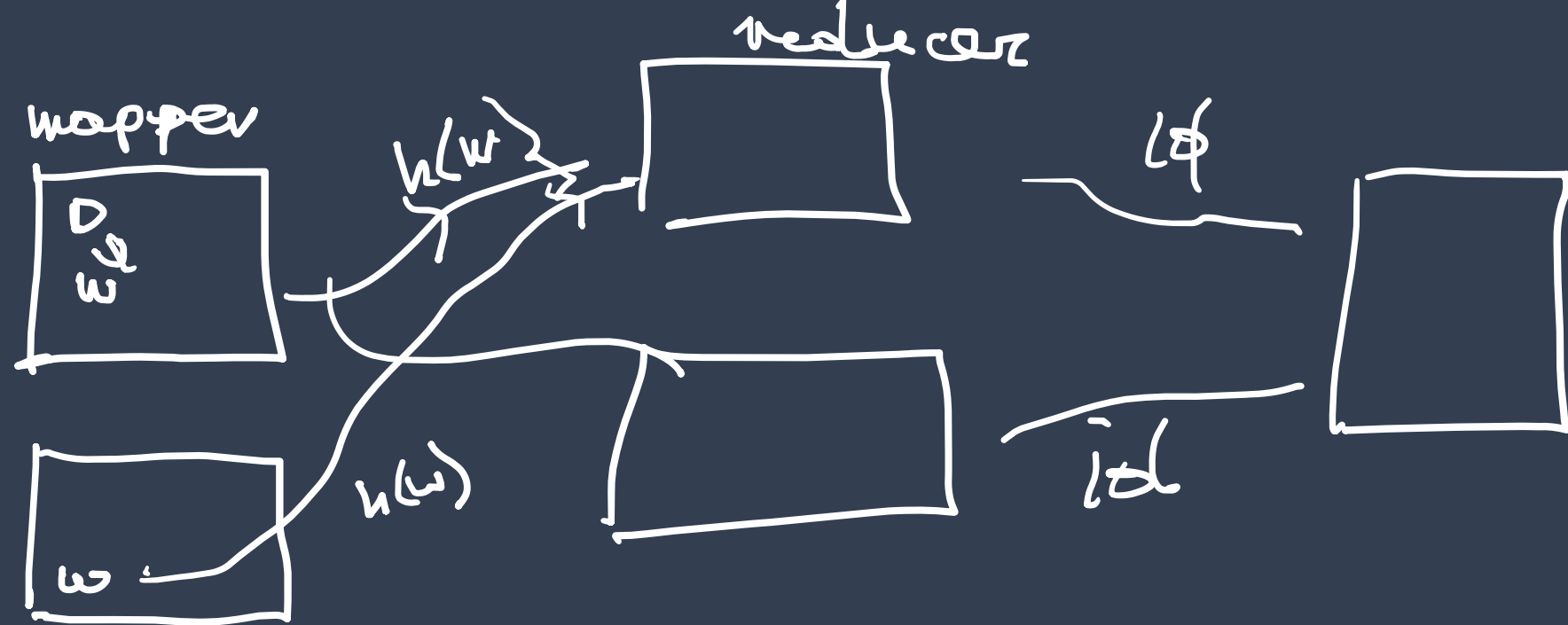


CEL: dla danego słowa w
wyznacz listę dokumentów
w których w występuje.

```
map(doc, Doc) {  
  for all  $w \in Doc$  do {  
    emit( $w, doc$ );  
  }  
}
```

```
reduce( $w, L$ ) {  
  emit( $w, L$ );  
}
```

$[w_1, L_1], [w_2, [id_1^2, id_2^2, \dots, id_k^2]]$



② Agregaty : $[x_1, \dots, x_n]$ \rightsquigarrow suma
 minimum
 maksimum
 ? $x \rightsquigarrow$ klucz
 $\varphi(x) = 10$ ostatnich bitów $(x)_2$

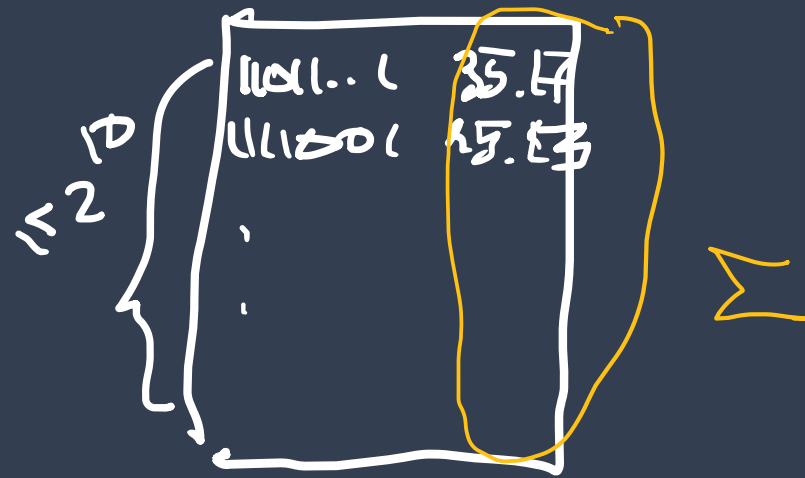
sum

map(x) {
 emit($\varphi(x), x$);
}

reduce(k, L) {
 emit(k, ΣL);
}

reducer

["110110001", [y₁, y₂, ..., y_L]]



$$\sum_i \sum_j a_{ij} =$$

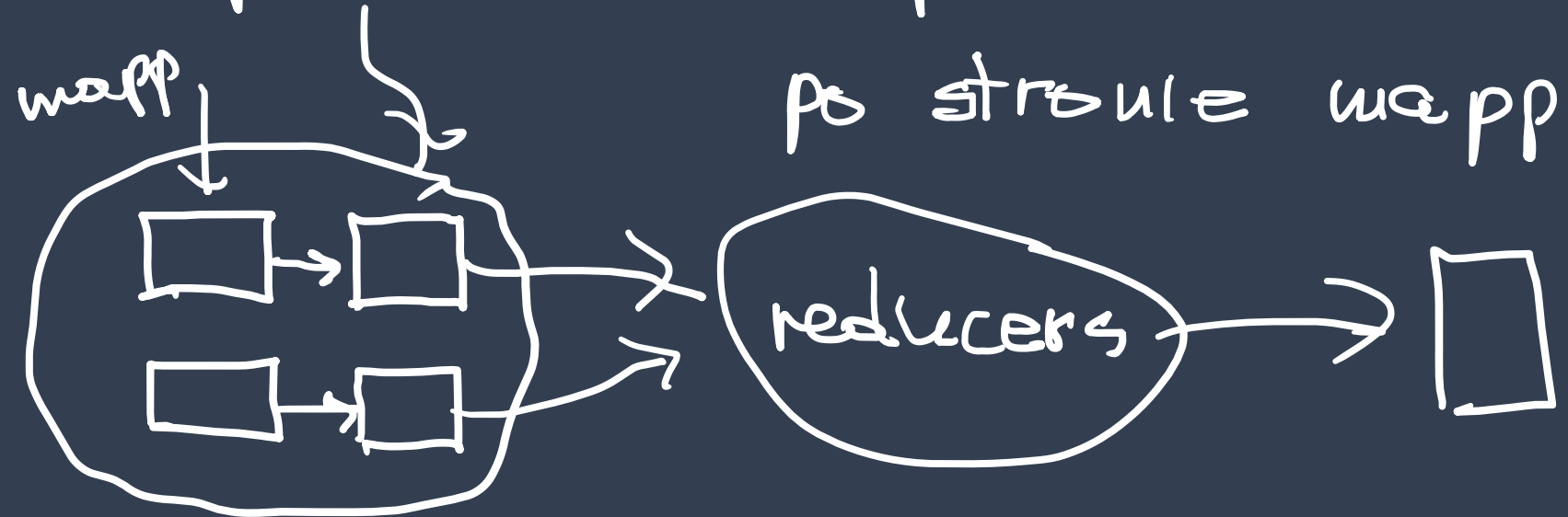
$$= (a_{11} + a_{12}) + (a_{21} + a_{22}) + \dots + (a_{m1} + a_{m2})$$

+ jest przewidziane i Σ czuje.



← największy fragment

Composers \equiv grup. reducera działający po stronie mappera



SUMY :

$$\left\{ \begin{array}{l} \text{map}(x) \quad \{ \text{emit}(\varphi(x), x) \} \\ \text{compose}(K, L) \quad \{ \text{emit}(k, \Sigma L); \} \\ \text{reduce}(L, L) \quad \{ \text{emit}(k, \Sigma L); \} \end{array} \right.$$

Przykłady łącznych operacji: Σ , min, max

uwaga: * - łączne $\ln(x_1 \cdot \dots \cdot x_n) = \sum_{L=1}^n \ln(x_L)$

monoidalne:

$$\text{map}(x) \quad \{ \text{emit}(\varphi(x), \text{lee}(x)) \}$$

Średnia : nie jest łączna

map(x) { emit ($\varphi(x)$, x, 1); }

reducer { (k, L) } // $L = [(x_1, 1), (x_2, 1), \dots, (x_n, 1)]$

emit ($k, \sum \pi_i(L), \text{length}(L)$)

}

$\left[\begin{array}{ccc} k_1 & S_1 & L_1 \\ k_2 & S_2 & L_2 \\ \vdots & \vdots & \vdots \\ k_m & S_m & L_m \end{array} \right]$

\rightarrow

$$\frac{S_1 + \dots + S_m}{L_1 + \dots + L_m}$$

Średnia

ZADANIE : zobowiązań to z
compositional $(x_1, L_1) \otimes (x_2, L_2)$

~~Task~~ (P) $(x_1, \dots, x_N) \rightsquigarrow \frac{\sum x_L}{N}, \underbrace{\frac{1}{N} \sum_{L=1}^N (x_L - \mu)^2}_{\text{var}}$

$$\frac{1}{N} \sum_{L=1}^N (x_L - \mu)^2 = \dots = \frac{1}{N} \sum x_L^2 - \left(\frac{\sum x_L}{N} \right)^2$$

zadanie

odpow: $\text{var}(X) = E[X^2] - (E[X])^2$

$$\text{map}(x) \rightsquigarrow \text{emit}(\psi(x), (1, x, x^2)); \}$$

$$\begin{array}{ccc} \downarrow & \downarrow & \downarrow \\ N & \sum x_L & \sum x_L^2 \end{array}$$

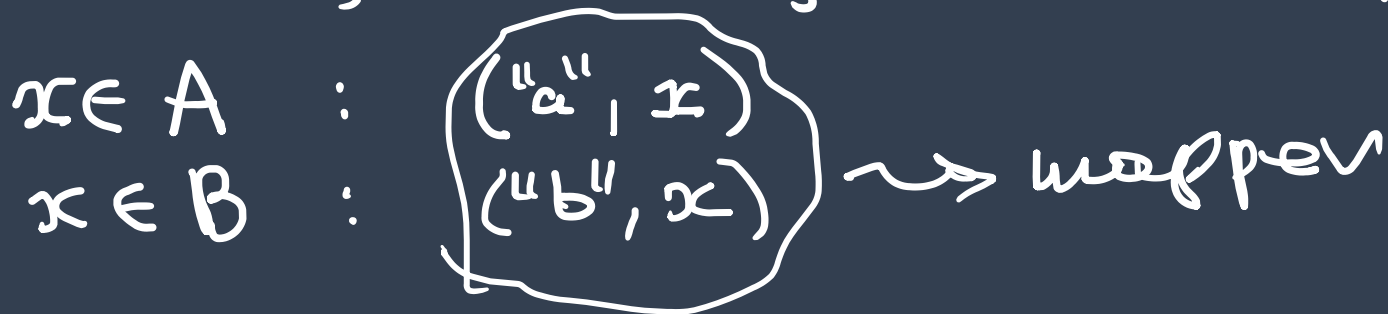
SŁKA MAP-REDUCE



SŁKA \equiv rozbić zadania na
niezaw. wyk. zadania

OPERACJE ZBIOROWE

domy: $A, B \leftarrow$ zbiory cel: $A \cup B, A \cap B, A \setminus B$



map(k, x) {emit(x, k);}

suma:

reduce(x, L) {emit(x);}

iloczyn

reduce(x, L) { if |L|=2 then emit(x); }

roznica (A \ B)

reduce(x, L) { if (L = ["a"]) then emit(x); }

Operacje baz-danych:

Relacja: $R(a_1, \dots, a_n)$, $\theta(a_1, \dots, a_n) \leftarrow$ Boolean

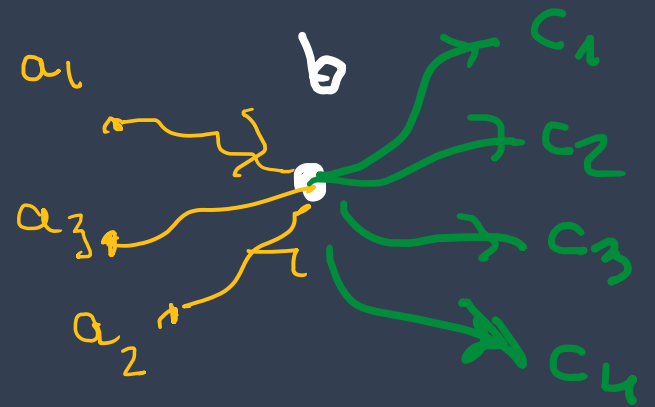
SELECT a_1, \dots, a_n FROM R WHERE $\theta(a_1, \dots, a_n)$

map(a_1, \dots, a_n) { if $\theta(\underbrace{a_1, \dots, a_n}_{\vec{a}})$ then emit $(\psi(\vec{a}), \vec{a})$
 $\psi(a_1, \dots, a_n) = a_1 || a_2 || \dots || a_n$

JOIN : $R(a,b), S(b,c)$ relational BD

{ SELECT $R.a, S.c$ FROM R, S
WHERE $R.b = S.b$

done : (rel, x, y)
 ↑
 R/S



map (r, x, y) {
 if $(r = "R")$ { emit $(y, ("R", x))$
 else { emit $(x, ("S", y))$
 }
}

reduce (b, L) {

1) $L_1 = \text{posortuj } L \text{ po 1}$
wspórz.

2) rozdzielanie

$$L_1 = L_{11} \parallel L_{12}$$

↑
pieniz
R

↑
pieni S

3) forall $x \in \pi_2(L_{11})$ do
forall $y \in \pi_2(L_{12})$ do
emit (x, y);

}}}

$$L = [(r_1, x_1), (r_2, x_2), \dots, (r_m, x_m)]$$

"R" "x" "S"

$R \times S$

$$L_1 = [\underbrace{(R_1, x_1), (R_1, x_2), \dots, (R_1, x_k)}_{(S, y_1), \dots, (S, y_e)}]$$

GRAPH : $G = (V, E)$ ← graph properties
 $E \subseteq [V]^2$

$N(x) = \{y : \{x, y\} \in E\}$

$N_2^+(x) = \{z : (\exists y) (\{x, y\} \in E \wedge \{y, z\} \in E)\}$



E

x	y1
x	y2
x	y3
...	...

$E^* = E \cup \text{flip}(E)$

Map JOIN $E^* \rightarrow E^*$

(Z) napisz to w zrodle