

O ROZKŁADZIE ODWOŁAŃ DO STRON WWW

1. ZNACZĄCOŚĆ STRUKTURY ODWOŁAŃ W SIECI WWW POZWALA ZOPTYMALIZOWAĆ JEJ BUDOWĘ. BADANIA EMPIRYCZNE POKAZAŁY, ŻE ROZKŁAD ODWOŁAŃ DO STRON WWW ODBYWA SIĘ ZGODNIE Z PRAWEM ZIPFA:

NIECH Ω BĘDZIE ZBIOREM ZASOBÓW WWW. $|\Omega| = N$

WPROWADZIMY PORZĄDEK NA TYCH ZASOBACH: $\Omega = \{z_1, z_2, \dots, z_N\}$

PRZY CZYM INDEKS i PRZY z_i OZNACZA MIEJSCE NA LIŚCIE RANKINGOWEJ NADPOPULARNIEJSZYCH ZASOBÓW (z_1 - NADBARDZIEJ POPULARNY; z_N - NAJMIEJ POPULARNY). OZNACZMY PRZEZ P_i PRAWDOPODOBIEŃSTWO, ŻE DANE ODWOŁANIE JEST ODWOŁANIEM DO ZASOBU z_i . ZACHODZI ZALEŻNOŚĆ:

$$P_i = \frac{c}{i^\alpha}$$

GDZIE c JEST STAŁĄ NORMUJĄCĄ I WYNOŚI $\frac{1}{H_N^\alpha}$, GDZIE $H_N^\alpha = \sum_{i=1}^N \frac{1}{i^\alpha}$. (DLA UPROSZCZENIA $H_N^1 = H_N$). BADANIA WYKAZAŁY, ŻE $0 < \alpha \leq 1$. WARTOŚĆ α ZALEŻY OD POPULACJI LUDZI, KTÓRZY ODWOŁUJĄ SIĘ DO STRON.

2. JAKIE JEST ZNACZENIE WSPÓŁCZYNNIKA α ?

ZASTANÓWIMY SIĘ NAD $\phi(k)$ - PRAWDOPODOBIEŃSTWEM, ŻE NASTĄPI ODWOŁANIE DO JEDNEJ Z k NADPOPULARNIEJSZYCH STRON

$$\phi(k) = \sum_{i=1}^k \frac{c}{i^\alpha} = c \cdot H_k^\alpha = \frac{H_k^\alpha}{H_N^\alpha}$$

ŻEBY MIEĆ LEPSZE WYOBRAZENIE O TEJ WARTOŚCI, WYKONAMY SZACOWANIE

$$\int_1^N \frac{1}{x^\alpha} dx \leq H_n^\alpha \leq 1 + \int_1^N \frac{1}{x^\alpha} dx$$

$$\int_1^m \frac{1}{x^\alpha} dx = \begin{cases} \ln x \Big|_1^m = \ln m \\ \frac{x^{1-\alpha}}{1-\alpha} \Big|_1^m = \frac{m^{1-\alpha}}{1-\alpha} - \frac{1}{1-\alpha} = \frac{m^{1-\alpha} - 1}{1-\alpha} \end{cases}$$

$$\frac{k^{\alpha-1}}{N^{\alpha-1}-1} = \frac{k^{\alpha-1}-1}{\alpha-1} \leq \frac{H_k^\alpha}{H_N^\alpha} \leq \frac{1 + \frac{k^{\alpha-1}-1}{\alpha-1}}{\frac{N^{\alpha-1}-1}{\alpha-1}} = \frac{k^{\alpha-1}-\alpha}{N^{\alpha-1}-1} \leq (1+\epsilon) \left(\frac{k}{N}\right)^{\alpha-1} + \delta \quad (2)$$

$$(1+\epsilon) \left(\frac{k}{N}\right)^{\alpha-1} - \delta$$

CZYLI $\frac{H_k^\alpha}{H_N^\alpha} = \Omega\left(\left(\frac{k}{N}\right)^{\alpha-1}\right)$.

OBSERWACJA:

DLA α BLISKIEGO ZERU PRAWDOPODOBIEŃSTWO, ŻE ODWOLEAMY SIĘ DO KTÓREGOŚ Z k PIERWSZYCH STRON JEST MNIEJSZE NIŻ, GDY α JEST BLISKIE 1.

3. JAKI JEST WSPÓŁCZYNNIK TRAFIEŃ?

A GDY CACHE JEST NA TYLE DUŻE, ŻE NIE TRZEBA BYĆO NIC USUWAĆ.

ZAKŁÓŻMY, ŻE WYKONAŁIMY JUŻ R ODWOŁAŃ I ZASTANAWIAMY SIĘ, JAKA JEST SZANSA, ŻE $R+1$. ODWOŁANIE JEST TRAFIENIEM.

$$\begin{aligned} H(R) &= \Pr\left(\bigcup_{i=1}^N \text{"WYLOSOWANO } i\text{-TY ZASÓB ZA } R+1 \text{ RAZEM I WCZEŚNIEJ TEŻ GO WYLOSOWANO}\right) \\ &= \sum_{i=1}^N \Pr(\text{"WYLOSOWANO } i\text{-TY ZASÓB ZA } R+1 \text{ RAZEM I WCZEŚNIEJ TEŻ GO WYLOSOWANO}) \\ &= \sum_{i=1}^N p_i \cdot \Pr(\text{"WCZEŚNIEJ WYLOSOWANO } i\text{-TY ZASÓB}) \\ &= \sum_{i=1}^N p_i \cdot (1 - \Pr(\text{"WCZEŚNIEJ NIE WYLOSOWANO } i\text{-TEGO ZASÓBU})) \\ &= \sum_{i=1}^N p_i \cdot \left(1 - \prod_{k=1}^R \Pr(\text{"W } k\text{-TEJ RUNDZIE NIE WYLOSOWANO } i\text{-TEGO ZASÓBU})\right) \\ &= \sum_{i=1}^N p_i \cdot \left(1 - \prod_{k=1}^R (1 - p_i)\right) = \sum_{i=1}^N p_i \cdot (1 - (1 - p_i)^R) \end{aligned}$$

PODSTAWIAMY TERAZ TO, CO WYNIKA Z PRAWA ZIPFA

$$S = \sum_{i=1}^N \frac{1}{i \ln N} \left(1 - \left(1 - \frac{1}{i \ln N}\right)^R\right)$$

Oszacujemy teraz:

(3)

$$1 - \left(1 - \frac{1}{i \ln N}\right)^R$$

przez 1 dla $i \leq \frac{R}{\ln N}$ oraz przez $1 - \left(1 - \frac{1}{i \ln N}\right)^R$ dla $i > \frac{R}{\ln N}$

Dostaniemy

$$S \leq \sum_{i=1}^{\frac{R}{\ln N}} \frac{1}{i \ln N} + \underbrace{\sum_{i=1}^N \frac{1 - \left(1 - \frac{1}{i \ln N}\right)^R}{i \ln N} - \sum_{i=1}^{\frac{R}{\ln N}} \frac{1 - \left(1 - \frac{1}{i \ln N}\right)^R}{i \ln N}}$$

Wartości $\left(1 - \frac{1}{i \ln N}\right)^R$ dla odpowiednio dużych R bliska jest $\frac{1}{i \ln N}$

Możemy zatem przyjąć, że $1 - \left(1 - \frac{1}{i \ln N}\right)^R \approx d$. Mamy zatem

$$S \approx \frac{\ln \frac{R}{\ln N}}{\ln N} + d - d \left(\frac{\ln \frac{R}{\ln N}}{\ln N}\right) \approx D \ln R + C$$

Czyli obserwując zmienność R powinniśmy dostać wykres logarymiczny

(B) GÓY ROZMIAR CACHE JEST OGRANICZONY.

(4)

ZAKŁADAMY, ŻE W CACHE MIEŚCI SIĘ S STRON METALEŻNIE OD ICH ROZMIARU. ZAKŁADAMY TEŻ, ŻE W CACHE ZNAJDUJE SIĘ S NAJPOPULARNIEJSZYCH STRON. WSPÓŁCZYNNIK TRAFIENIA

$$H(S) = \Pr(\text{WYŁOŻYLIŚMY ZASÓB SPOŚRÓD } S \text{ NAJBARDZIEJ POPULARNYCH}) = \sum_{i=1}^S p_i$$

MAMY

$$\frac{S^{1-\alpha} - 1}{1-\alpha} \leq H_S^\alpha \leq 1 + \frac{S^{1-\alpha} - 1}{1-\alpha}$$

LUB

$$\ln S \leq H_S^2 \leq 1 + \ln S$$

CZYLI

$$\frac{\frac{S^{1-\alpha} - 1}{1-\alpha}}{\frac{N^{1-\alpha} - 1}{1-\alpha}} \leq H(S) \leq \frac{\frac{S^{1-\alpha} - 1}{1-\alpha}}{\frac{N^{1-\alpha} - 1}{1-\alpha}}$$

$$\frac{S^{1-\alpha} - 1}{N^{1-\alpha} - 1} \leq H(S) \leq \frac{S^{1-\alpha} - 1}{N^{1-\alpha} - 1}$$

LUB

$$\frac{\ln S}{1 + \ln N} \leq H(S) \leq \frac{1 + \ln S}{\ln N}$$

CZYLI $H(S) = \Omega(c \ln S)$, gdzie $c = H_N$

4. ILE CZASU UPYTYWA MIĘDZY KOLEJNYMI ODWOŁANIAMI DO TEJ SAMEJ STRONY?

OBLICZYMY PRAWDOPODOBIEŃSTWO ŻE DOPIERO ZA k -TYM ODWOŁANIEM DOSTANIEMY ZASÓB

$$\begin{aligned} d(k) &= \sum_{i=1}^N \Pr(\text{UZYSKAŁYŚMY } i\text{-TY ZASÓB I DOSTANIEMY GO DOPIERO ZA } k\text{-TYM ODWOŁANIEM}) \\ &= \sum_{i=1}^N p_i \Pr(\text{Przez } k-1 \text{ nie ma } z_i \text{ oraz że } k\text{-tym razem jest } z_i \mid \text{ten był } z_i) = \\ &= \sum_{i=1}^N p_i (1-p_i)^{k-1} \cdot p_i = \sum_{i=1}^N p_i^2 (1-p_i)^{k-1} \end{aligned}$$

SPRÓBUJEMY DO WYZKANEGO WZORU WSTAWIĆ PRAWO ZIPFA I OSZACOWAĆ ZA POMOCĄ CAŁKI:

$$\int_1^N \left(\frac{1}{x \ln N}\right)^2 \left(1 - \frac{1}{x \ln N}\right)^{k-1} dx = (-1)^k \sum_{j=0}^{k-1} \frac{(-1)^j \binom{k-1}{j}}{(j+1) x^{j+1} \ln^{j+2} N} \Big|_{x=1}^{x=N} = \tag{1}$$

$$\sum_{j=0}^{k-1} \frac{(-1)^j \binom{k-1}{j}}{(j+1) x^{j+1} \ln^{j+2} N} z^{j+1} \Big|_{z=1} = \int_0^z \left(1 - \frac{z}{x \ln N}\right)^{k-1} dz \Big|_{z=1} =$$

$$= \left(\frac{(z - x \ln(N)) \left(1 - \frac{z}{x \ln N}\right)^k}{k+1} + \frac{x \ln(N)}{k+1} \right) \Big|_{z=1} =$$

$$= \frac{(1 - x \ln N) \left(1 - \frac{1}{x \ln N}\right)^k + x \ln N}{k+1}$$

PO ZASTOSOWANIU DO (1) MAMY:

$$(-1)^k \frac{1}{(\ln N) \ln N} \left((1 - N \ln N) \left(1 - \frac{1}{N \ln N}\right)^k + N \ln N - (1 - \ln N) \left(1 - \frac{1}{\ln N}\right)^k - \ln N \right)$$

WNIOSKI:

MOŻEMY ZAŁOŻYĆ, ŻE N JEST STĄCĄ WTEDY CAŁOŚĆ ZACHOWUJE SIĘ $\approx \frac{1}{k}$.

S, NA ILE MOŻEMY ZBLIŻYĆ SIĘ DO ROZKŁADU ZIPFA?

NIECH P BĘDIE CIĄGIEM KOŁEDNYCH ODWOŁAŃ DO ZASOBÓW WNN.

NIECH $n_p(i)$ BĘDIE LICZBĄ WYSTĄPIEŃ ODWOŁANIA DO z_i W P .

NIECH $E(i)$ BĘDIE WARTOŚCIĄ OCZEKIWANĄ LICZBY WYSTĄPIEŃ ODWOŁANIA DO z_i .

$$E(i) = \sum_{k=0}^{|P|} k \binom{|P|}{k} p_i^k (1-p_i)^{|P|-k} = \frac{|P|}{i \ln N}$$

UZYJEMY TERAZ METODY MERÓWNOJU CHERNOJA:

NIECH X_1, \dots, X_n BĘDĄ NEZALEZNYMI ZMIENNYMI LOSOWYMI TŁE

$$\Pr(X_i = 1) = p$$

$$\Pr(X_i = 0) = 1 - p$$

NIECH $X = \sum_{i=1}^n X_i$, NIECH $EX = pn$. WTEDY

$$\Pr(|X - pn| > \alpha) \leq 2e^{-\frac{2\alpha^2}{n}}$$

TO PASUJE DO NASZEJ SYTUACJI, BO $X_j = \begin{cases} 1 & \text{gdy } j\text{-TE ODWOZANIE JEST DO Z;} \\ 0 & \text{gdy } j\text{-TE ODWOZANIE NIE JEST DO Z} \end{cases}$

$$\Pr\left(\left|n p(i) - \frac{|P|}{\ln N}\right| > \alpha\right) \leq 2e^{-\frac{2\alpha^2}{|P|}}$$

MOŻEMY CHIEĆ SIĘ DOWIEDIEĆ, JAK DŁUGIEGO P TRZEBA, ABY UZYSKAĆ ODPOWIEDNIO MAŁE PRAWDOPODOBIEŃSTWO, ŻE WYDARZY SIĘ OD WARTOŚCI OCZEKIWANEJ.

$$2e^{-\frac{2\alpha^2}{|P|}} \leq \epsilon$$

$$-\frac{2\alpha^2}{|P|} \leq \ln \frac{\epsilon}{2}$$

$$|P| \geq \frac{2\alpha^2}{\ln 2 - \ln \epsilon}$$

GDYBYŚMY CHIEĆ DWA PRZYKŁADY, ŻE $n_p(i)$ JEST BLIŻEJ $E(i)$ NIŻ $E(i+1)$, TO PRZYJĘLIŚMY

$$\alpha < \frac{E(i) - E(i+1)}{2} = \frac{|P|}{2 \ln N} \left(\frac{1}{i} - \frac{1}{i+1} \right) = \frac{|P|}{2 \ln N} \cdot \frac{1}{i(i+1)}$$

DLA DOWOLNEGO i , CZYLI

$$\alpha < \frac{|P|}{2 \ln N} \cdot \frac{1}{N(N+1)}$$

WTEDY

$$|P| \geq \frac{2 \frac{|P|^2}{4 h^2 N \cdot N^2 (N+1)^2}}{\ln \frac{2}{\varepsilon}}$$

$$\ln \frac{2}{\varepsilon} \geq \frac{2 |P|}{4 h^2 N \cdot N^2 \cdot (N+1)^2}$$

$$|P| \leq 2 h^2 N \cdot N^2 (N+1)^2 \cdot \ln \frac{2}{\varepsilon}$$

Wniosek: przy bardzo długich próbkach możemy się pogubić.