

Proiect P&S, Tema 2

Proiectarea cache-ului pentru conținut internet

Distribuția Zipf este o distribuție discretă ce modelează popularitatea conținutului internet. În general distribuțiile internetului (Zipf și Pareto sunt distribuții cu coadă lungă *long tail distributions* sau *heavy tail distributions*).

Ce înseamnă această coadă? Dacă X este o variabilă aleatoare cu coadă lungă și x o valoare pozitivă foarte mare, $P(X > x)$ deși este neglijabilă, este nenulă. Cu alte cuvinte variabilele cu coadă lungă iau valori foarte mari, dar cu probabilitate mică (vezi imaginile din Notebook-ul inclus).

Distribuția de probabilitate este în procent foarte mare repartizată valorilor mici ale variabilelor, iar restul valorilor foarte mari.

În acest proiect studiați distribuția Zipf în următorul context:

Presupunem că lucrezi la startup-ul *Viață veselă*, ce pune la dispoziție conținut generat de utilizatori http://en.wikipedia.org/wiki/User-generated_content.

Printre altele utilizatorii pot să posteze filme și în același timp să le urmărească online. Ți se dă task-ul de proiectare a cache-ului pt cererea celor $n = 200$ de filme de care dispune până în prezent *Viață veselă*. Filmele sunt codificate $1, 2, \dots, 200$. Într-un fișier `DateFilme.txt` sunt înregistrate codurile filmelor urmărite online într-o perioadă de timp.

Pentru a decide popularitatea filmelor și a salva în cache filmele ce sunt cerute/vizionate de 75% dintre utilizatori, trebuie să identifiți din date distribuția Zipf ce aproximează distribuția de probabilitate a popularității filmelor.

Distribuția de probabilitate a popularității filmelor, "se învață" din datele din fișier astfel:

- Citești datele din fișier într-un vector `fv` (filme vizionate). Fie L lungimea vectorului.
- Definești tipul de date `Film` astfel:

```
typedef struct Film{
    int cod; // codul filmului
    int rank; // popularitatea filmului
    int nr; // nr de vizionari ale filmului avand codul cod
    double pre; // probabilitatea de vizionare a filmului
}Film;
```

și un pointer la `Film`:

```
Film *film;
```

Aloci memorie pt $n=200$ de filme:

```
film=malloc(n*sizeof(Film));
```

La început se inițializează doar membrul `cod`, și anume: `film[k].cod=k+1`, $k = 0, 1, \dots, n-1$.

`film[k].rank` reprezintă rank-ul filmului având codul $k+1$, în ierarhizarea filmelor în ordinea descrescătoare a popularității lor, `film[k].nr`, numărul de vizionări și respectiv `film[k].pre` este probabilitatea experimentală de vizionare a filmului cu codul $k+1$.

• Determinați din analiza coordonatelor vectorului `fv` nr de vizionări ale fiecărui film. Cum? Într-un caz simplu, dacă vectorul `fv` are coordonatele:

$$fv = [6, 3, 2, 7, 6, 1, 3, 3, 4, 3, 2, 6, 5, 6]$$

înseamnă că în perioada monitorizată au fost vizionate 7 filme, codificate 1,2,3,4,5,6,7:

- filmul 1 a fost vizionat o dată ;
- filmul 2 a fost vizionat 2 ori;
- filmul 3 a fost vizionat de 4 ori;

etc

Numărul de vizionări ale fiecărui film este calculat de o funcție
`void, frecventa(int *fv, int L, int n, Film *film);`

care determină de câte ori a fost vizionat filmul k și setează corespunzător `film[k].nr`, $k = 0, 1, \dots, n-1$.

- se sortează apoi filmele în ordinea descrescătoare a numărului de vizionări.

Inițial avem `film[0].cod=1`, `film[1].cod=2` ... `film[n-1].cod=n`

Dacă de exemplu, în urma sortării rezultă că filmul cu codul 32 este cel mai popular, apoi filmul 12, etc și ultimul este filmul 81, atunci se setează `film[0].cod=32`, `film[1].cod=17`, ... `film[n-1].cod=81` (adică se resetează membrul `cod` pt fiecare film);

În cursul sortării se setează și rank-ul fiecărui film:

`film[0].rank=1`, `film[1].rank=2`, ... `film[n-1].rank=n`

Deci ATENȚIE cum implementați sortarea!!!

- Calculați apoi `Nrtv`, numărul total de vizionări ale celor n de filme, adică: $Nrtv = \sum_{k=0}^{n-1} film[k].nr$. Ea trebuie să fie egală cu nr L al înregistrărilor din fișierul `DateFilme.txt`

- Probabilitatea experimentală de vizionare a fiecărui film din cele $n = 200$ se calculează ca fiind $film[k].pre = \frac{film[k].nr}{Nrtv}$, $k = 0, 1, \dots, n-1$.

După aceste calcule am definit practic o variabilă aleatoare ce ia ca valori popularitatea filmelor, și probabilitățile sunt probabilitățile deduse din frecvențele de vizionare:

$$X = \begin{pmatrix} 1 & 2 & \dots & n \\ film[0].pre & film[1].pre & \dots & film[n-1].pre \end{pmatrix}$$

Dorim să studiem dacă această variabilă aleatoare este aproximată suficient de bine de o v.a Zipf:

$$Z = \begin{pmatrix} 1 & 2 & \dots & n \\ \frac{1}{c} & \frac{2}{c} & \dots & \frac{n}{c} \\ \frac{1}{1^\alpha} & \frac{1}{2^\alpha} & \dots & \frac{1}{n^\alpha} \end{pmatrix}$$

Teoretic distribuția variabilei Z , adică $p_k = P(Z = k)$, se exprimă astfel $p_k = \frac{c}{k^\alpha}$, $k = 0, 1, \dots, n-1$.

Logaritmând avem:

$$\ln(p_k) = \ln(c) - \alpha \ln k$$

Inmulțind cu -1 rezultă că $-\ln(p_k) = \alpha \ln(k) - \ln(c)$ și notând $x = \ln k$, $y = -\ln(p_k)$, $b = -\ln(c)$, relația devine:

$$y = \alpha x - b$$

adică pentru orice $k = 0, 1, \dots, n-1$, punctele de coordonate $(\ln(k), -\ln(p_k))$ verifică ecuația unei drepte de pantă α .

Distribuția variabilei X , dedusă din date este aproximativ o distribuție Zipf dacă punctele asociate $(\ln(k), -\ln(\text{film}[k].\text{pre}))$ sunt ușor împrăștiat în jurul unei drepte. Cum $n=200$ puncte obținute din date experimentale este puțin probabil să fie coliniare, determinăm dreapta celor mai mici pătrate asociate acestor puncte și analizăm (vizual și computațional) cât de bine se potrivește modelul liniar, logaritmilor datelor.

Dacă punctele sunt suficient de apropiate de dreapta celor mai mici pătrate, atunci panta drepte va fi parametrul α al distribuției Zipf aproximative a datelor.

Pentru a genera și analiza dreapta celor mai mici pătrate procedați astfel:

- Scrieți într-un fișier **NumeLeader2.txt**, pe câte o linie coordonatele celor $n = 200$ puncte:

$$\ln(\text{film}[k].\text{rank}) - \ln(\text{film}[k].\text{pre})$$

Atenție sunt 2 nr pe o linie, nu cumva să interpretați semnul minus ca diferență!!!!!!!!!!!!!!

- Citiți acest fișier în IPython Notebook-ul **NumeLeaderPr2.ipynb**, desenați punctele, calculați dreapta celor mai mici pătrate, deci coeficienții ce definesc dreapta și un parametru r , numit coeficientul de corelație (îl studiem la curs în săptămâna 11-16 mai). Dacă coeficientul $r > 0.9$ considerăm că modelul liniar este potrivit pentru date și deci panta drepte, α , dă parametrul Zipf al distribuției aproximative a datelor din fișierul **DateFilme.txt**.

Introduceți apoi în codul C/C++ valoarea pentru parametrul α ce caracterizează distribuția Zipf aproximativă a datelor

- Definiți o funcție **double *Zipf(double alpha, int n)** ce returnează vectorul probabilităților $pr[k]$, $k=0, 1, 2, \dots, n-1$ asociate distribuției Zipf de parametru α (luați doar 2 zecimale pt α).

- Salvați într-un fișier **ProbsExp.txt** probabilitățile experimentale $\text{film}[k].\text{pre}$ $k = 0, \dots, n-1$ și într-un alt fișier **ProbsTheor.txt** aceste probabilități ale distribuției teoretice ce aproximează distribuția datelor; (la sfârșit vor fi vizualizate într-un notebook cele două distribuții)

- Calculați coordonatele vectorului F , $F[0] = 0$, $F[1] = pr[0]$, $F[2] = pr[0] + pr[1]$, \dots , $F[n] = pr[0] + pr[1] + \dots + pr[n-1] = 1$

- Deduceți prin calcul până la ce rang de popularitate k trebuie salvate copii ale filmelor în cache, pt a putea răspunde rapid cererilor de vizionare a filmelor din top 25%. Adică determinați cel mai mic k pentru care $P(X \leq k) \geq 0.75$.

- Pentru a testa că alegerea filmelor pentru cache este OK, simulați variabila aleatoare Zipf de parametru α , generând 2000 de valori sau mai mult (adică cereri de vizionare de filme) și numărați ce procent din valorile generate sunt coduri ale filmelor salvate în cache. Afișați acest procent cu un mesaj, pt a înțelege ce reprezintă. Ce părere aveți despre procentul găsit? Ce ilustrează el?

Pentru simularea acestei variabile aleatoare discrete, nu se aplică algoritmul din curs pentru că variabila ia valori $n=200$ valori și ar fi inefficient. Și anume se folosește vectorul F definit mai sus, în care s-au precalculat sume ale probabilităților consecutive (ATENȚIE, spre deosebire de algoritmul din curs aici $F[0] = 0$; adică îl includem și pe 0 în vectorul F). Algoritmul de simulare folosește căutarea binară pentru a determina intervalul $(F[k-1], F[k]]$ în care cade numărul $u = \text{urand}()$:

```

1: function SimulZipf( $n, F$ )
2:    $u = \text{urand}()$ ;
3:    $L=0$ ;  $R=n$ ; //  $L=\text{Left}$ ,  $R=\text{Right}$ 
4:   while ( $L < R$ ) {
5:      $m = (L+R)/2$ ;
6:     if ( $u < F[m]$ )  $R=m$ ;
7:     else  $L=m+1$ ;
8:   }
9:   return  $L$ ; //  $u \in (F[L-1], F[L])$ ; s-a produs un ev de prob  $pr[L-1]$ 
10: end function

```

Atentie! Notebook-ul inclus se ruleaza de 2 ori. Prima data doar pana ce afiseaza dreapta celor mai mici patrute si parametrul α pentru distribuția Zipf aproximantă.

A doua oară rulați din nou de la început până la sfârșit și veți obține vizualizarea comparativă a celor două distribuții: distribuția datelor și distribuția Zipf aproximantă.

Atât prima oară, cât și a doua se generează câte o imagine pe care apoi o includeți în arhiva cu această temă de proiect.