

World Wide Web

Caracteristicile rețelei și ale motoarelor de căutare

decembrie 2018

Conținutul acestui curs

- Motoare de căutare web: scurt istoric, tehnologii de bază
- Estimarea nr. de documente indexate de motoarele de căutare web
- Eliminarea documentelor dupicat de către algoritmii de indexare

Caracteristici ale rețelei World Wide Web

1. Foarte mare (rețea globală, de nivel planetar)
2. Se extinde într-un mod necoordonat
3. Diversitate mare de utilizatori/participanți

Tehnologii pe care se bazează WWW:

HTML (Hypertext Markup Language)

HTTP (Hypertext Transfer Protocol)

URL (Universal Resource Locator)

Exemplu de URL:

`http://www.stanford.edu/home/atoz/contact/html`

- ▶ `http`: - protocol de transmitere a datelor
- ▶ `www.stanford.edu` - domeniu pt. o ierarhie de pagini web
- ▶ `/home/atoz/contact/html` - cale la o pagina web în ierarhia domeniului

Arhitectură client-server:

- comunicarea server-client prin protocolul http de transmitere asincronă a informațiilor de diverse feluri (text, imagini, fișiere video/audio etc.) formatare în limbajul html
- clientul (de obicei, un browser), transmite o cerere http la un server web:
 - se specifică un URL, de exemplu
`http://www.stanford.edu/home/atoz/contact.html`
 - serverul transmite un răspuns codificat html, care conține
 - hiperlinkuri
 - conținut formatat cu taguri html.
- browserul are reguli de afișare a conținutului răspunsului, în funcție de formatarea html

Capabilitățile primelor browsere web

- Crearea ușoară de documente formatare html, și adresarea lor cu url
 - proces rapid de învățare
- Vizualizarea documentelor formatare html

Consecințe:

- publicarea masivă de informații diverse de către practic oricine
⇒ WWW a devenit cea mai mare sursă de informații din lume
- **PROBLEMĂ NOU ÎVITĂ:** WWW ca sursă enormă de informații necesită motoare performante de căutare a informațiilor de interes.

Extragerea informațiilor de pe WWW

Capabilitățile primei generații de motoare de căutare

Descoperirea și accesarea informațiilor publicate, care sunt de interes pentru consumatori. Mai întâi au apărut 2 tipuri de motoare de căutare:

- ① Căutare de cuvinte cheie în texte, bazată pe utilizarea de indeces inversați (calculați în prealabil) și mecanisme de calcul al gradului de importanță (engl. *ranking*). Exemple: Altavista, Excite și Infoseek
- ② Căutare de documente organizate într-o ierarhie arborescentă de categorii. Exemplu: Yahoo!
 - **Avantaje:** metodă intuitivă de căutare, ușor de utilizat
 - **Dezavantaje:** clasificarea documentelor în categorii este un proces uman costisitor
 - de regulă, muncă făcută de editori experți
 - greu de efectuat pentru colecții mari de date (WWW)

Motoare de căutare

Caracteristicile primei generații de motoare de căutare

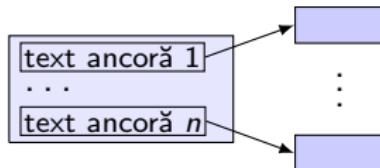
- Bazat pe tehnici clasice de găsire a răspunsurilor la întrebări de căutare, precum
 - ▶ liste de indești inversați
 - ▶ algoritmi de calcul al gradului de importanță al unui răspuns
- Volumul mare de documente publicate pe WWW ($10^6 \approx 10^7$) a determinat:
 - ▶ indexarea unui număr tot mai mare de documente, stocate și procesate pe un număr mare de mașini \Rightarrow timpi de răspuns la întrebări < 1 secundă
 - ▶ calitatea și relevanța rezultatelor căutării lăsau de dorit:
 - \Rightarrow inventarea unor algoritmi mai buni de calcul al gradului de relevanță (engl. *ranking algorithms*)
 - \Rightarrow tehnici de detecție și evitare a spam-ului
 - \Rightarrow tehnici de detecție al **gradului de autoritate** al unui document (de ex., pe baza site-ului web de unde provine)

- Zeci de limbi și mii de dialecte
 - ▶ extragerea și indexarea conținutului poate fi făcută doar cu algoritmi lingvistici avansați, pentru fiecare limbaj în parte
- Pagini web cu variații mari de structură, fonturi, și culori
 - ▶ unele pagini web nu au text indexabil, ci doar doar imagini
- Libertatea de a publica orice ⇒ pagini ce conțin adevăruri, minciuni, contradicții, bănuieri, etc.
 - ▶ cum putem măsura gradul de încredere al unei pagini web?
- Întrebarea „cât de mare este WWW?” nu are un răspuns simplu:
 - (+) putem estima nr. de pagini indexate de către un motor de căutare
 - (?) cum recunoaștem și cum tratăm paginile web generate dinamic? (De ex., pagina generată dinamic despre zborul AA129)



Graful web

$G = (V, E)$ cu mulțimea de noduri V =pagini web statice; mulțimea E de arce=hiperlinkuri între pagini web.



NOTIUNI AUXILIARE:

- **link in** al unei pagini A : un arc cu destinația A
- **grad in** al paginii A =nr. de arce cu destinația A
- **link out** al unei pagini A : un arc cu sursa A
- **grad out** al paginii A : nr. de arce cu sursa A

Se estimează că

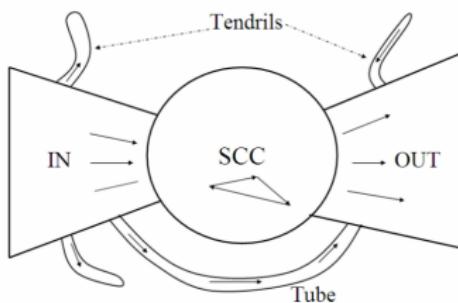
- În medie, **gradul in** al unui nod în G este între 8 și 15.
- Nr. de noduri n cu in-grad(n) = i este proporțional cu $1/i^\alpha$, unde $\alpha \approx 2.1$

Proprietăți ale grafului web

Sunt 3 categorii majore de pagini web: IN, OUT și SCC, astfel încât

- există căi de la orice pagină IN la orice pagină SCC
- există căi de la orice pagină SCC la orice pagină OUT
- există căi de la orice pagină SCC la orice pagină SCC
- Nu există căi $SCC \rightsquigarrow IN$, și nici căi $OUT \rightsquigarrow SCC$
- nr. pagini IN \approx nr. pagini OUT

⇒ graful web G are o structură de papion



Spam și tehnici de evitare a lui

Spam = tehnici de manipulare a conținutului paginilor web a.î. să aibă scor mare de importanță pentru căutări după anumite cuvinte cheie

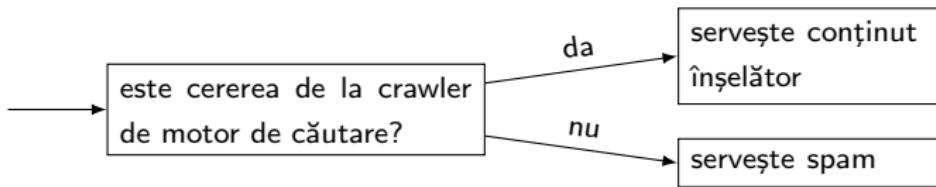
- activitate determinată, de obicei, de motive economice.

Exemplu de spam timpuriu:

- ▶ pagini web cu conținut repetat de anumiți termeni afișați în culoarea de background a paginii; create de vânzători ce fac reclamă la produsele lor, și de către agenții lor
- ▶ antidot: motoare de căutare care detectează nr. mare de repetiții de termeni în pagini web, și nu le atribuie un scor mare de importanță

Cloaking = tehnică mai sofisticată de creat spam:

- serverul web al creatorului de spam păcălește componenta **crawler** a motoarelor de căutare, verificând de la cine vine cererea de servire:



Producătorii de optimizatoare de motoare de căutare (SEO) oferă servicii de consultanță pentru clienții interesați să aibă pagini web cu scor mare de importanță pentru anumite cuvinte cheie

- se bazează pe descifrarea tehnicilor secrete ale motoarelor de căutare, care se folosesc pentru a calcula scorul de importanță al paginilor web
 - ⇒ conflict perpetuu între producătorii de SEO și producătorii de motoare de căutare
 - ⇒ numeroase motoare de căutare (începând cu Google) folosesc o tehnică, numită **analiza link-urilor**, pentru a combate spamming ca se bazează pe manipularea textului din pagini web

Reclama ca model economic

Branding=prezentarea de reclame (pe pagini web de pe site-uri populare) care produc impresii pozitive despre compania pentru care se face reclama.

Modele economice pentru reclamele făcute pe pagini web:

- ❶ **cost per mil (CPM)**: costul suportat de companie pt. a afișa 100 de reclame
- ❷ **cost per click (CPC)**: costul pt. numărul de click-uri pe reclama respectivă pt. a ajunge la pagina web a companiei pt. care se face reclama, de unde se pot efectua cumpărături

Modelul economic de **căutare sponsorizată** al companiei Goto (preluat apoi de majoritatea motoarelor de căutare):

- acceptă licitații de la companii pentru a le afișa paginile web ca **reclame de căutare** pentru anumiți termeni de căutare q
- afișarea reclamelor de căutare în ordinea descrescătoare a sumelor licitate
- Compania Goto era plătită după modelul economic CPC

Tipuri de căutare

Căutarea algoritmică și căutarea sponsorizată

Motoarele actuale de căutare afișează două tipuri de răspunsuri la o întrebare q :

- rezultate ale **căutării algoritmice** (sau **căutare pură**) ca răspuns primar la solicitarea utilizatorului
- rezultate ale **căutării sponsorizate**: de obicei acestea sunt afișate separat, în dreapta rezultatelor algoritmice
 - căutarea sponsorizată se bazează pe versiuni mai sofisticate ale modelului Goto de reclamă sponsorizată

Tipuri de căutare

Exemplu ilustrat

YAHOO! SEARCH

Web | Images | Video | Local | Shopping | [more »](#)

A320 Advanced Search

Search Results

1 - 10 of about 5,050,000 for **A320** - 0.22 sec. ([About this page](#))

Also try: [airbus a320](#), [a320 family](#), [airbus industrie a320](#), [a320 type rating](#) [More...](#)

1. Airbus A320 family - Wikipedia, the free encyclopedia
... more than 3,000 aircraft of the **A320** family built, it is the second best ...
Airbus intends to relocate Toulouse **A320** final assembly activity to Hamburg
as ...
en.wikipedia.org/wiki/Airbus_A320 - 112k - [Cached](#)

2. Airbus A320 - Airliners.net
Offers a history, specifications, photos, and performance data of the Airbus
A320.
www.airliners.net/info/stats.main?id=23 - 26k

3. Airbus: A320 Family
From the official Airbus site, featuring information on the Airbus A318, A319,
A320, and A321.
www.airbus.com/en/aircraftfamilies/a320 - 18k

SPONSOR RESULTS

WADS A320 - Refurbished
A320 - on sale for \$293.25.
20A 240V 3P - free UPS ground.
www.relectric.com

Bluetooth Stereo USB - Jabra A320s
Connect Your PC to Your Bluetooth Stereo Headset with the Jabra **a320s**.
www.helldirect.com/jabra-a320s

Căutarea orientată pe cerințele utilizatorilor

Majoritatea utilizatorilor WWW sunt *neprofesioniști*:

- nu sunt interesați să folosească sintaxa sofisticată a limbajelor artificiale de interogare web (operatori booleeni, wildcard-uri, etc.); tind să folosească 2 sau 3 cuvinte de căutare
- ⇒ pentru a fi atractive și profitabile, motoarele de căutare s-au adaptat la cerințele utilizatorilor

Strategia de succes a lui Google:

- 1 A îmbunătățit **precizia** primelor rezultate afișate
- 2 A scurțat **timpul de localizare** a informațiilor de interes, prin afișarea rezultatelor în format text, cu foarte puține elemente grafice

DE RETINUT: calitatea răspunsurilor unui motor de căutare la o întrebare q se estimează după două criterii:

precizie: procentul de răspunsuri relevante din totalul răspunsurilor primite

senzitivitate: procentul de răspunsuri relevante din totalul documentelor relevante

(vezi slide-ul următor)

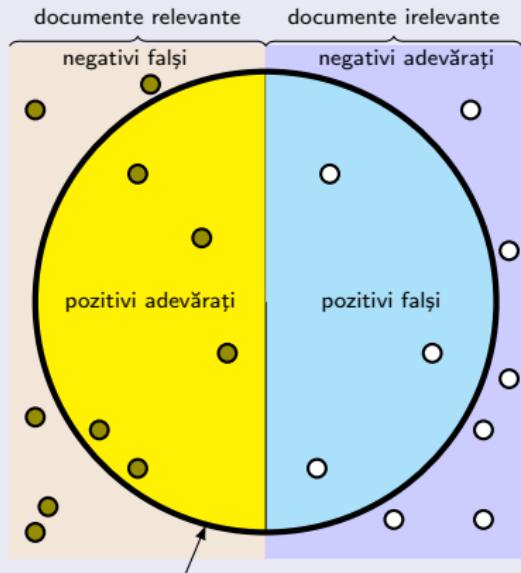
Motore de căutare

Estimarea calității căutării: precizie și senzitivitate

Exemplu ilustrat

nr. total documente $N = 19$; nr. documente relevante $N_r = 10$;

nr. documente selectate $N_s = 8$



● : document relevant

○ : document irrelevant

$$\text{precizie} = \frac{\text{pozitivi adevărați}}{N_s} = \frac{5}{8} = 0.625$$

$$\text{senzitivitate} = \frac{\text{pozitivi adevărați}}{N_r} = \frac{5}{10} = 0.5$$

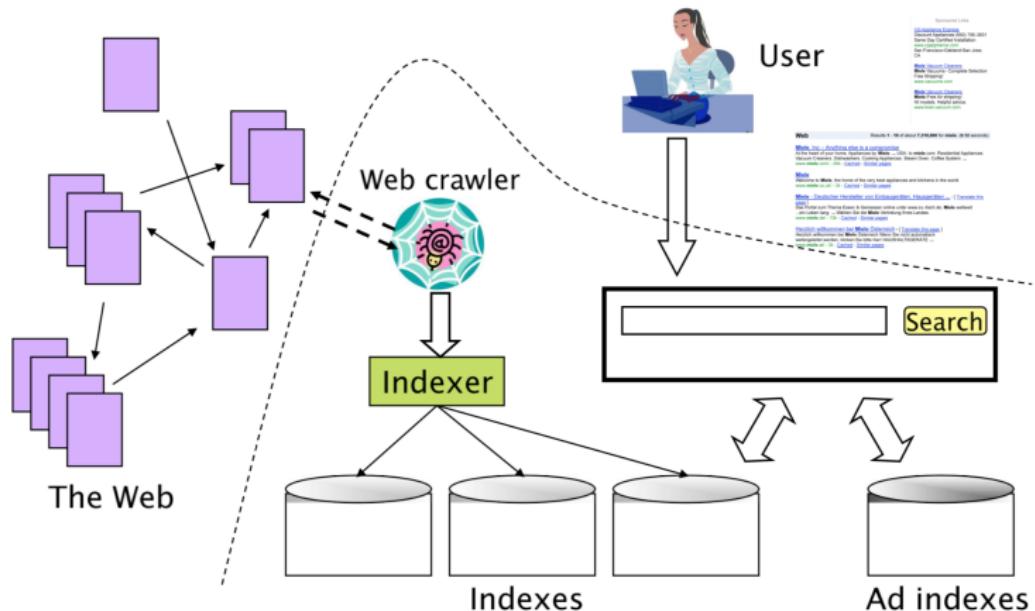
Precizia măsoară exactitatea sau calitatea răspunsurilor.

Senzitivitatea măsoară completitudinea sau sau cantitatea răspunsurilor.

- ① **Căutări informaționale:** au ca scop găsirea de informații generale despre un subiect larg (de exemplu, istoria imperiului roman); informațiile se extrag din mai multe pagini web.
- ② **Căutări navigaționale:** se caută pagina web a unei anumite entități, de ex. pagina web a liniei aeriene KLM.
- ③ **Căutări tranzacționale:** au ca scop identificarea paginilor web ce ofera anumite servicii, de ex. cumpărarea unui produs, descărcarea unui fișier, efectuarea unei rezervări.

Motoarele de căutare trebuie să poată detecta tipul de căutare dorit de utilizator, și să își adapteze căutarea.

Structura generală a motoarelor moderne de căutare



REMARĂ: Rolul și modul de funcționare al componentei numite **crawler** sunt descrise în alt curs.

După mărime =numărul de pagini web statice indexate

- (-) Mărimea indexului variază rapid și nu este informație publică
- (+) Există algoritmi de estimare a indexului.

Remarcă

Întrebarea: „*Care din motoarele de căutare A sau B are o mărime mai mare a documentelor indexate?*” nu are un răspuns precis, din mai multe motive:

- ▶ De regulă, se indexează doar o mică parte a conținutului paginilor web (câteva mii de cuvinte de la început)
- ▶ Indecșii sunt organizați pe nivele și în partiții distribuite pe mai multe mașini ⇒ estimarea mărimii este un proces dificil

⇒ s-au inventat mai multe metode de măsurare aproximativă a mărimii indecșilor motoarelor de căutare.

Compararea mărimii indecșilor motoarelor de căutare

Metoda capture-recapture

Să se compare mărimele indecșilor motoarelor de căutare E_1 și E_2 :

$$\frac{|E_1|}{|E_2|} = ? \quad \text{unde } |E_i| = \text{mărimea indexului motorului de căutare } E_i$$

Presupunem că

- ① Fiecare motor E_i indexează o fracțiune de documente D_i publicate pe WWW
- ② Fracțiunile de documente D_i sunt alese uniform și aleator.

Metoda experimentală **capture-recapture** este definită astfel:

- ▶ alege aleator o pagină din indexul lui E_1 și verifică dacă este în indexul lui $E_2 \Rightarrow x\%$ din paginile în E_1 sunt în E_2
- ▶ alege aleator o pagină din indexul lui E_2 și verifică dacă este în indexul lui $E_1 \Rightarrow y\%$ din paginile în E_2 sunt în E_1

$$\text{Estimăm că } x |E_1| \approx y |E_2| \Rightarrow \frac{|E_1|}{|E_2|} \approx \frac{y}{x}$$

Metoda capture-recapture

Scenarii de calcul al valorilor x, y a.â. $|E_1|/|E_2| \approx y/x$

- **Cazul simplu:** x și y sunt calculate de către cineva cu acces la indecșii motoarelor de căutare $\Rightarrow x$ și y se pot calcula ușor, alegând aleator documente din indecșii motoarelor E_1 și E_2 .
- **Cazul dificil:** putem alege aleator pagini web doar *din exteriorul indecșilor* lui E_1 și E_2
 - alegerea aleatoare a unei pagini web din colecția totală D de documente publicate pe WWW este o problemă dificilă
 - tehnici aproximative de alegere aleatoare a unui document d din D :
 - ① Căutări aleatoare
 - ② Adrese IP aleatoare
 - ③ drumuri aleatoare (engl. *random walks*)
 - ④ interogări aleatoare (engl. *random queries*)

Metoda capture-recapture

Scenarii de calcul al valorilor x, y a.î. $|E_1|/|E_2| \approx y/x$

Căutarea aleatoare

- se pornește de la un log de căutare al unui grup (de ex., membrii unui grup de cercetare) care a fost de acord să î se înregistreze toate cererile de căutare
- Se alege aleator o cerere q din log și o pagină răspuns p din log pentru q
- Se trimit cererea q la motorul E_1 și se verifică dacă p apare în lista de răspunsuri $\Rightarrow x\%$ din paginile răspuns din log sunt indexate de către E_1
- Se trimit cererea q la motorul E_2 și se verifică dacă p apare în lista de răspunsuri $\Rightarrow y\%$ din paginile răspuns din log sunt indexate de către E_2

DEFICIENȚĂ: calculul lui x și y nu este aleator și uniform:

- ▶ depinde de interesele de căutare ale grupului ales.

Metoda capture-recapture

Scenarii de calcul al valorilor x, y a.â. $|E_1|/|E_2| \approx y/x$

Adrese IP aleatoare

- Se generează adrese IP aleatoare și se trimit o cerere la serverul web de la adresa aleasă \Rightarrow se colectează toate paginile web de pe serverul respectiv
- Se combină colecțiile de pagini web culese de pe servere într-o colecție D
- se estimează
 - x : ce procent de pagini din D sunt indexate de către E_1 ;
 - y : ce procent de pagini din D sunt indexate de către E_2

DEFICIENȚE

- ① (determinată de virtual hosting): host-uri diferite de pagini web au aceeași adresă IP
- ② Unele site-uri web au un nr. mic de pagini web \Rightarrow documentele din D nu sunt alese aleator și uniform

Metoda capture-recapture

Scenarii de calcul al valorilor x, y a.î. $|E_1|/|E_2| \approx y/x$

Drumuri aleatoare (engl. *random walks*)

Tehnică adekvată pentru grafuri orientate puternic conexe (\exists drum de la orice nod la orice alt nod)

- Un drum aleator care pornește de la o pagină arbitrar aleasă converge la o distribuție stabilă.

DEFICIENȚE:

- Graful web nu este graf orientat puternic conex
- Chiar și dacă ar fi, timpul necesar pt. a aproxima bine starea de distribuție stabilă este, în general mare și dificil de estimat.

Metoda capture-recapture

Scenarii de calcul al valorilor x, y a.â. $|E_1|/|E_2| \approx y/x$

Interogări aleatoare (engl. *random queries*)

Idee de bază: alegerea (aproape) uniform aleatoare a paginilor web indexate de un motor de căutare se face trimițând interogări aleatoare la motorul de căutare; de exemplu:

- ▶ din primele 100 răspunsuri ale lui E_1 la o interogare conjunctivă aleatoare q , alegem aleator o pagină p .
- ▶ verificăm dacă pagina p este indexată de E_2 , astfel:
 - alegem 6-8 termeni cu frecvență joasă de apariții în p , și îi folosim pt. a construi o interogare conjunctivă pentru E_2

Evitarea indexării documentelor duplicate

- Estimare: aprox. 40% din paginile indexate web sunt duplicate ale altor pagini.
- Se dorește indexarea copiilor multiple ale aceluiași conținut
⇒ reducere a spațiului alocat pt. liste de indecesi inversați, și a timpilor de preprocesare.

Amprentarea (engl. *fingerprinting*)

- calculul, pt. fiecare pagină web, a unei amprente, care este un rezumat (digest) succint (de ex., 64 biți) al conț. de caractere al paginii respective.
- dacă două documente au aceeași amprentă, se verifică dacă sunt duplicate.

DEFICIENȚE:

- Nu detectează ca fiind duplicate documentele care diferă prin modificarea câtorva caractere (engl. *near duplication*)
- Această deficiență poate fi evitată cu **shingling**

Tehnici de detecție a documentelor duplicate

Shingling

Pt. un întreg $k > 1$, mulțimea de ***k-shingles*** a unui document d este muș. de secvențe consecutive de 4 termeni care apar în d .

- Dacă d conține **a rose is a rose is a rose**, atunci muș. de 4-shingles a lui d este
 - a rose is a
 - rose is a rose
 - is a rose is

Intuiție: două documente sunt aproape duplicate (engl. *near duplicates*) dacă mulțimile de *shingles* generate din ele sunt similare.

- avem nevoie de metode eficiente de calcul și comparare a muș. de shingles ale paginilor web.

Christopher D. Manning, Prabhakar Raghavan, Hinrich Schütze:
Capitolul 19: *Web search basics* din **AN INTRODUCTION TO
INFORMATION RETRIEVAL**.
Ediție online (c) 2009 Cambridge UP.