

## Material Suplimentar Curs sap 11: Algoritmul PageRank. Lanțuri Markov absorbante.

### 0.1 Construcția lanțului Markov pe graful WEB. Algoritmul PageRank

Succesul extraordinar și dominația motorului Google se datorează în principal algoritmului PageRank, care exploatează structura linkurilor din WWW pentru a determina un indice de popularitate al fiecărei pagini, independent de interogarea formulată de utilizator.

Documentele de pe WEB (paginile WEB) sunt identificate de aplicațiile software ale motorului, numite roboți sau *crawlere*. Documentele sunt apoi indexate. Modulul de indexare extrage cuvintele cheie, constituind așa numitul sac de cuvinte. Un alt modul, numit *query module* (modulul de interogare), convertește cererea formulată de utilizator, în limbaj natural, într-un vector cerere, cu care consultă indexul de conținut și extrage paginile relevante cererii. Modulul de ierarhizare ordonează descrescător aceste pagini în funcție de coeficienții de popularitate. PageRank-ul este un vector ale cărui coordonate sunt coeficienții de popularitate ai paginilor WEB identificate de crawler. Acest vector este distribuția de echilibru a unui lanț Markov definit pe graful WEB.

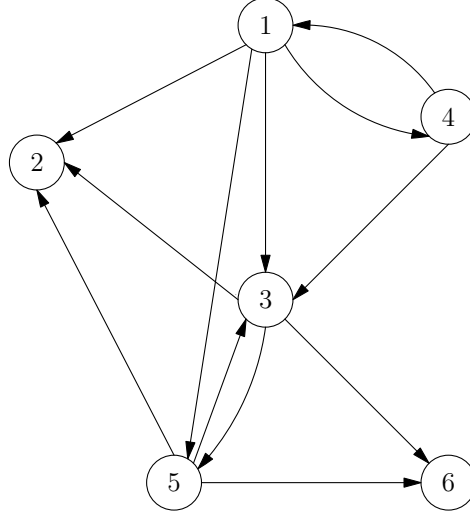
Să definim mai întâi lanțul Markov ce stă la baza algoritmului **PageRank**. Considerăm  $W = \{1, 2, \dots, m\}$  mulțimea tuturor paginilor WEB,  $H = (h_{ij})$  matricea de conectivitate a lui  $W$  sau matricea hyperlink:

$$h_{ij} = \begin{cases} 1, & \text{dacă există link în pagina } i \text{ către pagina } j, \\ 0, & \text{dacă nu există link în pagina } i \text{ către pagina } j. \end{cases}$$

Se spune că  $H$  este o matrice rară, căci are foarte multe zerouri (în medie, 3-10 elemente sunt nenule pe o linie). Suma elementelor de pe linia  $i$  a matricei  $H$  indică numărul de out-linkuri, adică numărul de linkuri din pagina  $i$  către alte pagini sau ea însăși. Notăm această sumă cu  $r_i = \sum_{j=1}^m h_{ij}$ .  $r_i$  se numește ordinul ieșirilor din pagina  $i$ . Suma elementelor de pe coloana  $i$  a matricei hyperlink indică numărul de in-linkuri ale paginii  $i$ , adică numărul de linkuri către pagina  $i$ .

Larry Page și Serghei Brin au definit un mers aleator pe graful WEB considerând că un surfer ajuns în pagina  $i$  alege cu aceeași probabilitate oricare din paginile către care aceasta are linkuri, prin urmare probabilitatea de a trece din pagina  $i$  în pagina  $j$  este:

$$p_{ij} = \begin{cases} \frac{1}{r_i}, & \text{dacă există link în pagina } i \text{ către pagina } j, \\ 0, & \text{dacă nu există link în pagina } i \text{ către pagina } j. \end{cases}$$



**Fig.1:** Graf orientat ilustrând linkurile între 6 pagini WEB.

De exemplu, dacă

$$H = \begin{pmatrix} 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

atunci ordinul de ieșire din pagina 5 este  $r_5 = 3$ , deci probabilitatea de a trece din pagina 5 în oricare din paginile  $\{1, 2, \dots, 6\}$  este  $p_{5j} = h_{5j}/3$ . Altfel spus, cu aceeași probabilitate de  $1/3$  un surfer poate trece din pagina 5 în pagina 2, 3 sau 6.

Vom exemplifica construcția propusă de L. Page și S. Brin prin modelul simplu de rețea izolată de pagini WEB (rețea intranet) din Fig. 1.

Notăm cu  $Q = (p_{ij})_{i,j=\overline{1,6}}$ , matricea probabilităților de tranziție definite mai sus. Se observă din structura grafului de conectivitate că paginile 2 și 6 sunt pagini ce nu conțin linkuri către alte pagini. Acestea se numesc *dangling pages*. De exemplu, fișierele **pdf**, **ps** sau fișierele imagine sunt pagini dangling. Prin urmare, liniile 2 și 6 din matricea de tranziție au toate elementele nule și, astfel,  $Q$  nu este o matrice stohastică, deci nu poate fi interpretată ca matricea de tranziție a unui lanț Markov cu spațiul stărilor  $\{1, 2, 3, 4, 5, 6\}$ .

### 0.1. CONSTRUCȚIA LANȚULUI MARKOV PE GRAFUL WEB.ALGORITMUL PAGERANK3

$$Q = \begin{array}{c|cccccc} & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 1 & 0 & 1/4 & 1/4 & 1/4 & 1/4 & 0 \\ 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 3 & 0 & 1/3 & 0 & 0 & 1/3 & 1/3 \\ 4 & 1/2 & 0 & 1/2 & 0 & 0 & 0 \\ 5 & 0 & 1/3 & 1/3 & 0 & 0 & 1/3 \\ 6 & 0 & 0 & 0 & 0 & 0 & 0 \end{array}$$

Pentru a remedia această problemă, Page și Brin au propus ca vector de probabilitate dintr-o pagină dangling  $i$  distribuția uniformă, considerând

$$p_{ij} = 1/m, \quad j = \overline{1, m}.$$

Adică, în mod artificial, se adaugă linkuri dintr-o pagină dangling către toate paginile WEB sau, echivalent, ajuns într-o pagină dangling, un navigator poate apoi alege cu o probabilitate uniformă orice pagină din WWW. Astfel, matricea stochastică obținută din matricea  $Q$  este:

$$\tilde{Q} = \begin{array}{c|cccccc} & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 1 & 0 & 1/4 & 1/4 & 1/4 & 1/4 & 0 \\ 2 & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} \\ 3 & 0 & 1/3 & 0 & 0 & 1/3 & 1/3 \\ 4 & 1/2 & 0 & 1/2 & 0 & 0 & 0 \\ 5 & 0 & 1/3 & 1/3 & 0 & 0 & 1/3 \\ 6 & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} \end{array}$$

Lanțul Markov definit de matricea stochastică  $\tilde{Q}$  nu este în general ireductibil (adică nu există drum de linkuri între orice două pagini sau, echivalent, graful WEB nu este tare conex) și pot exista traiectorii periodice, adică surferul, navigând conform matricei de tranziție  $\tilde{Q}$ , ar putea fi prins, ca într-o capcană, într-o mișcare aleatoare ciclică. Din acest motiv, dar și pentru că la un moment dat, și în realitate, orice surfer renunță să navigheze urmând linkurile din pagini, L. Page și S. Brin au introdus ipoteza că doar cu probabilitatea  $\alpha \in (0, 1)$  surferul navighează conform matricei  $\tilde{Q}$  și cu probabilitatea  $1 - \alpha$  ignoră linkurile și alege cu probabilitate uniformă oricare din paginile de pe WEB, introducând adresa URL în linia de comandă a browser-ului. Probabilitatea  $\alpha$  se numește *factor de damping*. În lucrarea inițială a fondatorilor Google,  $\alpha$  era menționat ca având valoarea 0.85. Cu această modificare matricea de tranziție este:

$$G = \alpha \tilde{Q} + (1 - \alpha) \underbrace{\begin{bmatrix} 1/m & 1/m & \dots & 1/m \\ 1/m & 1/m & \dots & 1/m \\ & & \vdots & \\ 1/m & 1/m & \dots & 1/m \end{bmatrix}}_E.$$

Matricea  $G$  se numește matricea Google, iar matricea  $E$ , de elemente identice  $1/m$ , se numește *matricea de teleportare*, deoarece surferul se teleportează din navigarea aleatoare urmând linkuri într-o "navigare artificială". Evident că și matricea  $E$  este o matrice stohastică, iar  $G$  fiind o combinație convexă de astfel de două matrice este o matrice stohastică (vezi proprietățile matricelor stohastice). Mai mult,  $G(i, j) > 0, \forall i, j = \overline{1, m}$ , deci matricea Google este ireductibilă și aperiodică.

Se presupune că matricea Google este cea mai "uriasă" matrice cu care se lucrează în vreo aplicație la ora actuală.

Lanțul Markov având:

- spațiul stărilor constituit din mulțimea paginilor WEB, de cardinal  $m$
- matricea de tranziție de tipul  $G$ , cu  $\alpha$  fixat
- distribuția inițială de probabilitate  $\pi_0$  (distribuția uniformă sau oricare alta)

este un lanț ireductibil și aperiodic, deci are o unică distribuție de echilibru  $\pi$ , numită vectorul PageRank.

PageRank-ul,  $\pi$ , este limita șirului  $(\pi_n)$ , definit prin  $\pi_n^T = \pi_0^T G^n$ . Limita este aceeași indiferent de distribuția inițială de probabilitate  $\pi_0$ , adică indiferent cu ce probabilitate surferul alege pagina din care începe navigarea.  $\pi(j)$  se numește PageRank-ul paginii  $j$  și reprezintă șansa asimptotică pe care o are pagina  $j$  de a fi vizitată de navigatorul aleator sau proporția din timpul de navigare pe care surferul ar petrece-o vizitând pagina  $j$ . Deci,  $\pi(j)$  este un indice de popularitate al paginii.

Când un utilizator introduce cuvinte cheie în bara de căutare, motorul Google caută paginile ce conțin cuvintele cheie și le afișează în ordinea descrescătoare a PageRank-ului lor.

Remarcăm că PageRank-ul unei pagini este independent de interogarea formulată de utilizator. Ea depinde doar de structura grafului WEB și se poate calcula offline. PageRank-ul se calculează la intervale regulate de timp. Până în 2008 se calcula lunar, dar acum se actualizează la intervale mai scurte de timp.

Vectorul PageRank se calculează numeric, folosind așa numita metodă a puterii, adică se calculează recursiv. Pornind de la  $\pi_0$  și  $G$ , se determină distribuțiile (sau PageRank-ul la  $n$  pași de navigare)  $\pi_n^T = \pi_{n-1}^T G$ . Se consideră că metoda a atins stadiul de convergență (adică s-a ajuns la echilibru) într-o etapă  $n$  în care  $\|\pi_n - \pi_{n-1}\| < \varepsilon$ , unde  $\varepsilon$  este un număr pozitiv foarte mic, prescris.

Pseudocodul algoritmului de calcul al PageRank-ului este:

```

1: function PageRank(G, m);
2:    $\pi = [1/m, 1/m, \dots, 1/m]$ ; //Distributia initiala de probabilitate
3:    $\text{eps} = 10^{-7}$ ;
4:   do
5:      $\pi' = \pi$ ;
6:      $\pi = \pi' * G$ ;
7:   while ( $\|\pi - \pi'\| \geq \text{eps}$ );
8:   return  $\pi$ ;
9: end function
```

S-a demonstrat că viteza de convergență a metodei puterii este aceeași cu rata de convergență a lui  $\alpha^n$ , unde  $\alpha$  este factorul de damping.

**Implicații asupra PageRank-ului.** Din punctul de vedere al vitezei de convergență ar fi preferabil un factor  $\alpha$  cât mai apropiat de zero. În acest caz, ținând seama că matricea Google este  $G = \alpha\tilde{Q} + (1 - \alpha)E$ , ar rezulta că se acordă o pondere redusă,  $\alpha$ , navigării conform linkurilor din graful WEB (cu modificarea pentru pagini dangling) și o pondere mai mare navigării artificiale, conform matricei de teleportare  $E$ . Cu alte cuvinte, în acest caz PageRank-ul asociat nu ar reflecta popularitatea reală a paginilor WEB. De aceea o valoare rezonabilă, așa cum a fost ea aleasă inițial de Larry Page și Sergei Brin,  $\alpha = 0.85$ , conduce la rezultate mai apropiate de realitate și la o viteză de convergență suficient de bună (un reprezentant Google a declarat că metoda puterii converge după 100-200 de iterații). Dacă vreți să aflați PageRank-ul unor pagini WEB intrați aici:

[http://www.prchecker.info/check\\_page\\_rank.php](http://www.prchecker.info/check_page_rank.php)

### 0.1.1 Pagerank-ul personalizat

Pentru o ierarhizare personalizată a paginilor WEB, matricea  $E$  se calculează luând în considerare vectorul personalizat  $w$ , care este un vector probabilist  $w = [a_1, a_2, \dots, a_m]^T$  ale cărui coordonate reprezintă probabilitatea ca surferul, ce iese din navigarea conform linkurilor, să aleagă pagina  $1, 2, \dots, m$  din WEB. Cu alte cuvinte, el nu alege o pagină în mod uniform, ci are anumite preferințe, identificate de motor în decursul timpului. Astfel, matricea de teleportare va fi  $E = ew^T$ , unde  $e = [1, 1, \dots, 1]^T$ , iar matricea Google corespunzătoare este

$$G = \alpha\tilde{Q} + (1 - \alpha)ew^T.$$

Distribuția de echilibru corespunzătoare este PageRank-ul personalizat.

## 0.2 Lanțuri Markov absorbante

Procese într-un sistem de operare, arhitectura unui sistem software, protocoalele rețelelor wireless (WLAN), execuția protocoalelor de rețele (LAN), overflow-ul în traficul printr-un sistem de comunicație se modelează prin lanțuri Markov cu stări absorbante, al căror comportament asimptotic (la limită, când  $n \rightarrow \infty$ ) este diferit de cel al lanțurilor Markov ireductibile și aperiodice.

Pentru a caracteriza un lanț Markov absorbant clasificăm mai întâi stările unui lanț Markov general.

Considerăm un lanț Markov omogen ce are spațiul stărilor  $S = \{1, 2, \dots, m\}$ , iar matricea de tranziție este  $Q$ .

Definim pe mulțimea  $S$  relația: starea  $i$  intercomunică cu starea  $j$ , și notăm  $i \rightleftharpoons j$ , dacă probabilitatea de tranziție de la  $i$  la  $j$  într-un număr de pași este pozitivă, la fel ca și probabilitatea de tranziție de la  $j$  la  $i$ , adică există  $n, k \in \mathbb{N}^*$  astfel încât  $Q^n(i, j) > 0$  și  $Q^k(j, i) > 0$ .

Relația  $\Rightarrow$  definită mai sus este tranzitivă și simetrică, dar nu este reflexivă, căci pot exista stări care nu intercomunică cu ele însele. De exemplu, pentru lanțul Markov ce are spațiul stărilor  $S = \{1, 2, 3\}$  și matricea de tranziție

$$Q = \begin{pmatrix} 0 & 1/2 & 1/2 \\ 0 & 1/3 & 2/3 \\ 0 & 1/4 & 3/4 \end{pmatrix},$$

starea 1 nu este reflexivă, adică nu intercomunică cu ea însăși.

În continuare, vom presupune că toate stările unui lanț Markov sunt reflexive, adică intercomunică cu ele însele, caz în care relația  $\Rightarrow$  este o relație de echivalență pe spațiul stărilor.

- Clasa de comunicare a unei stări  $i$ , notată cu  $C(i)$ , este mulțimea stărilor cu care  $i$  intercomunică, adică  $C(i) = \{j \in S : i \Rightarrow j\}$ . Dacă există  $i \in S$  astfel încât  $C(i) = S$ , atunci lanțul Markov este ireductibil (oricare două stări intercomunică).

- O clasă  $C$  se numește *clasă închisă* dacă pentru orice stare  $i \in C$  și  $j$  în complementul său, notat cu  $\overline{C}$ , probabilitatea de tranziție de la  $i$  la  $j$  este 0, adică  $Q^n(i, j) = 0, \forall n \in \mathbb{N}^*$ .

- Dacă o clasă închisă conține o singură stare  $i$ , starea  $i$  se numește *stare absorbantă*. În acest caz  $p_{ij} = 0$ , pentru orice  $j \neq i$ , deci  $p_{ii} = 1$ . Cu alte cuvinte, dacă lanțul Markov ajunge în starea  $i$ , atunci rămâne în acea stare veșnic (cu probabilitatea 1). Pe linia stării absorbante  $i$  a matricei de tranziție  $Q$ ,  $Q(i, i) = 1$ , iar restul elementelor sunt nule.

- O clasă de comunicare  $C$  în care există o stare  $i$  ce comunică într-un singur pas cu complementul clasei se numește *clasă tranzitorie*, adică există  $i \in C$  și  $j \in \overline{C}$  astfel încât  $Q(i, j) > 0$ . Stările unei clase tranzitorii se numesc stări tranzitorii, deoarece cu o probabilitate pozitivă un lanț ce pornește dintr-o astfel de stare eventual părăsește clasa.

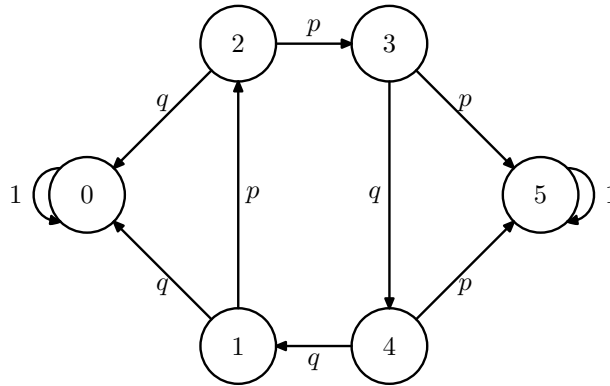
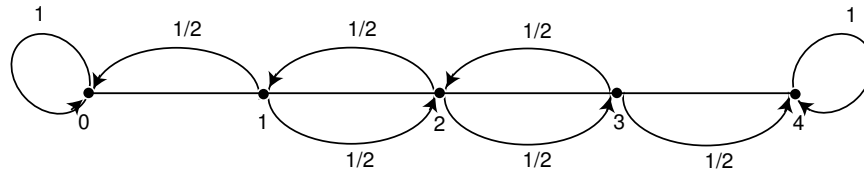
Vom studia în continuare un lanț Markov cu un număr finit de stări, ce poate conține una sau mai multe stări absorbante și o clasă tranzitorie.

Un lanț Markov cu număr finit de stări, în care fiecare stare este fie stare absorbantă, fie stare tranzitorie, se numește *lanț Markov absorbant*.

În Fig. 2 este ilustrat graful unui lanț Markov cu două stări absorbante, 0 și 5, iar clasa  $C = \{1, 2, 3, 4\}$  constă din stări tranzitorii.

Un exemplu de lanț Markov absorbant, ce a inspirat folosirea acestui tip de lanțuri Markov în știința și ingineria calculatoarelor, este mersul bețivului (Fig.3). Între bar (starea 0) și casă (starea 4) există 3 colțuri, 1, 2, 3. Când bețivul ajunge la un colț  $i$ ,  $i \in \{1, 2, 3\}$ , cu probabilitatea  $1/2$  o ia spre stânga, adică trece în  $i - 1$  și cu probabilitatea  $1/2$  o ia spre dreapta, deci trece în starea  $i + 1$ . Barul și casa sunt stări absorbante. O dată ajuns în una din aceste stări, rămâne sigur acolo.

- Lanț Markov absorbant este și lanțul definit de graful rețelei WWW, fără modificările efectuate de Larry Page și Serghei Brin. Și anume, considerăm matricea hyperlink asociată rețelei intranet considerate în prezentarea algoritmului PageRank, în care pagina 2 și pagina 6 sunt pagini dangling. În loc să luăm 0 pe liniile 2 și 6, considerăm  $Q(2, 2) = 1$  și  $Q(6, 6) = 1$ , adică odată ajuns în pagina 2, respectiv 6, un surfer rămâne blocat sigur (cu probabilitatea 1) în aceste pagini:

**Fig.2:** Graf al unui lanț Markov absorbant.**Fig.3:** Graful lanțului Markov ce modelează mersul bețivului.

	1	2	3	4	5	6
1	0	1/4	1/4	1/4	1/4	0
2	0	1	0	0	0	0
3	0	1/3	0	0	1/3	1/3
4	1/2	0	1/2	0	0	0
5	0	1/3	1/3	0	0	1/3
6	0	0	0	0	0	1

- Rețelele P2P (*peer-to-peer networks*)

<http://en.wikipedia.org/wiki/Peer-to-peer>

sunt modelate prin lanțuri Markov absorbante. O astfel de rețea este constituită dintr-un număr de routere, indexate prin  $1, 2, \dots, n - 1$ . Un pachet de informație este forwardat de către un router către altul conform unui anumit protocol, în scopul de a fi direcționat către destinație. Destinația, indexată prin  $n$ , este o stare absorbantă. Acestea i se mai adaugă o stare absorbantă, numită în limbaj de specialitate *drop state*. Ajunse în această stare pachetele fie sunt eronate iremediabil, fie sunt detectate ca având conținut sau scop distructiv.

Problemele uzuale care se pun relativ la un lanț Markov absorbant sunt:

- 1) Care este probabilitatea ca lanțul Markov ce pornește dintr-o stare tranzitorie  $i$  să fie absorbit de starea absorbantă  $j$ ?
- 2) În medie, în câți pași este absorbit lanțul Markov?
- 3) Care este timpul mediu pe care lanțul Markov îl petrece într-o stare tranzitorie?

Răspunsurile la aceste întrebări depind în general de starea din care pornește lanțul Markov, precum și de probabilitățile de tranziție.

Pentru a răspunde acestor întrebări, formulate relativ la lanțuri Markov absorbante, se descompune spațiul stărilor  $S$ , cu  $|S| = m$ , în două submulțimi disjuncte  $S = S_a \cup S_t$ , unde  $S_a$  este submulțimea stărilor absorbante, de cardinal  $n_a$ , iar  $S_t$  este submulțimea stărilor tranzitorii, de cardinal  $n_t$ ,  $n_a + n_t = m$ . Printr-o reordonare (renumerotare) a mulțimii stărilor, considerăm că stările  $1, 2, \dots, n_a$  sunt absorbante, iar stările  $n_a + 1, \dots, m$  sunt tranzitorii. Astfel, matricea de tranziție a lanțului Markov absorbant conține 4 blocuri de forma

$$Q = \begin{pmatrix} I & O \\ R & T \end{pmatrix}, \quad (1)$$

unde:

- $I$  este matricea unitate de tip  $n_a \times n_a$ ;
- $O$  este o matrice nulă de tip  $n_a \times n_t$ ;

Justificarea pentru blocurile  $I$  și  $O$  vine din faptul că stările  $1, 2, \dots, n_a$  sunt absorbante și deci  $p_{ii} = 1$ ,  $p_{ij} = 0$ ,  $\forall i \in \{1, 2, \dots, n_a\}$ ,  $j \in S$ ,  $j \neq i$ .

- Matricea  $R$  este matricea de tranziție de la stările tranzitorii la cele absorbante. Ea este de tip  $n_t \times n_a$ . Cu alte cuvinte,  $R(i, j) = P(X_{n+1} = j \in S_a | X_n = i \in S_t)$ ,  $n \in \mathbb{N}$ ;

- Matricea  $T$  dă probabilitățile de trecere între stările tranzitorii. Ea este de tip  $n_t \times n_t$  și  $T(i, j) = P(X_{n+1} = j \in S_t | X_n = i \in S_t)$ ,  $n \in \mathbb{N}$ .

Matricea de tranziție a unui lanț Markov absorbant scrisă în forma (1) se numește *matrice de tranziție în forma standard*.

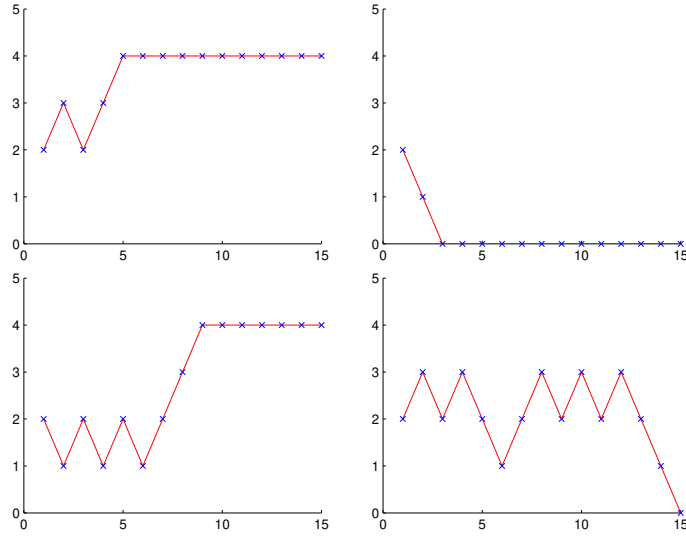
Matricea de tranziție a modelului pentru mersul la întâmplare al bețivului este

$$Q = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0.5 & 0 & 0.5 & 0 & 0 \\ 0 & 0.5 & 0 & 0.5 & 0 \\ 0 & 0 & 0.5 & 0 & 0.5 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}.$$

În figura Fig. 4 sunt ilustrate câteva traiectorii simulate ale bețivului. Folosim acest exemplu pentru a ilustra aspectele teoretice ce le discutăm.

Permutând vectorul  $(0, 1, 2, 3, 4)$  în  $(0, 4, 1, 2, 3)$ , obținem matricea de tranziție în forma standard, notată cu  $Q' = P_\pi Q P_\pi^T$ , unde  $\pi \in S_5$  este permutarea descrisă mai sus, adică  $\pi = \begin{pmatrix} 0 & 1 & 2 & 3 & 4 \\ 0 & 4 & 1 & 2 & 3 \end{pmatrix}$ , iar  $P_\pi$  este matricea permutare corespunzătoare,





**Fig.4:** Traectorii ale bețivului în 15 pași, ce pornesc din colțul 2. El este "absorbit" fie de bar (cod 0), fie de casă (cod 4). Stările sunt marcate pe axa ordonatelor, iar timpul pe axa absciselor.

$$Q' = \begin{array}{c|ccccc} & 0 & 4 & 1 & 2 & 3 \\ \hline 0 & 1 & 0 & 0 & 0 & 0 \\ 4 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0.5 & 0 & 0 & 0.5 & 0 \\ 2 & 0 & 0 & 0.5 & 0 & 0.5 \\ 3 & 0 & 0.5 & 0 & 0.5 & 0 \end{array}.$$

Să calculăm puterea  $n$  a matricei de tranziție (1) a unui lanț Markov absorbant. Produsul a două matrice de blocuri se efectuează ca produsul matricelor în general:

$$Q^2 = \begin{pmatrix} I & O \\ R & T \end{pmatrix} \begin{pmatrix} I & O \\ R & T \end{pmatrix} = \begin{pmatrix} I & O \\ R + TR & T^2 \end{pmatrix},$$

$$Q^3 = Q^2 Q = \begin{pmatrix} I & O \\ R + TR & T^2 \end{pmatrix} \begin{pmatrix} I & O \\ R & T \end{pmatrix} = \begin{pmatrix} I & O \\ R + TR + T^2 R & T^3 \end{pmatrix}.$$

Prin inducție se arată că

$$Q^n = \begin{pmatrix} I & O \\ R + TR + T^2 R + \dots + T^{n-1} R & T^n \end{pmatrix} = \begin{pmatrix} I & O \\ (I + T + T^2 + \dots + T^{n-1})R & T^n \end{pmatrix}.$$

Scopul nostru este să determinăm  $\lim_{n \rightarrow \infty} Q^n$  și să interpretăm semnificația blocurilor din matricea limită.

**Proprietatea 1.** Matricea  $T^n$  tinde la matricea nulă când  $n \rightarrow \infty$ .

**Interpretare:**  $T^n(i, j) = P(X_n = j \in S_t | X_0 = i \in S_t)$  reprezintă probabilitatea ca la momentul  $n$  lanțul să fie în starea tranzitorie  $j$  știind că la momentul inițial era în starea tranzitorie  $i$ . Deoarece  $T^n \rightarrow O$ , rezultă că fiecare element al matricei tinde la zero, deci pentru  $n$  foarte mare probabilitatea ca lanțul să mai fie într-o stare tranzitorie, știind că a pornit dintr-una tranzitorie, este foarte mică, adică, aproape sigur (cu probabilitatea 1) o traiectorie ce pornește dintr-o stare tranzitorie este absorbită de o stare din  $S_a$ .

**Proprietatea 2.** Matricea  $I - T$  este inversabilă și

$$\lim_{n \rightarrow \infty} (I + T + \dots + T^{n-1}) = (I - T)^{-1}. \quad (2)$$

**Demonstrație:** Presupunem că matricea  $I - T$  nu este inversabilă, deci sistemul liniar și omogen  $(I - T)\mathbf{x} = 0$  admite și soluții nebanale  $\mathbf{x} \neq 0$ . Prin urmare, avem

$$\mathbf{x} - T\mathbf{x} = 0 \Leftrightarrow T\mathbf{x} = \mathbf{x}.$$

Succesiv, se obține:

$$T\mathbf{x} = \mathbf{x}, \quad T^2\mathbf{x} = T(T\mathbf{x}) = T\mathbf{x} = \mathbf{x}, \dots, T^n\mathbf{x} = \mathbf{x}.$$

Trecând la limită când  $n \rightarrow \infty$  în relația  $T^n\mathbf{x} = \mathbf{x}$  și ținând seama că  $T^n \rightarrow O$ , obținem  $0 = \mathbf{x}$ , ceea ce este imposibil, deci presupunerea făcută, că matricea  $I - T$  nu este inversabilă, este falsă. Se poate verifica relația:

$$(I - T)(I + T + \dots + T^{n-1}) = I - T^n.$$

Înmulțind la stânga cu  $(I - T)^{-1}$ , obținem

$$I + T + \dots + T^{n-1} = (I - T)^{-1}(I - T^n).$$

Trecând la limită când  $n$  tinde la infinit, avem

$$\lim_{n \rightarrow \infty} (I + T + \dots + T^{n-1}) = (I - T)^{-1}(I - \lim_{n \rightarrow \infty} T^n) = (I - T)^{-1}(I - O) = (I - T)^{-1}.$$

□

Notăm cu  $N = (I - T)^{-1}$ . Matricea  $N$  se numește *matricea fundamentală* a lanțului Markov absorbant.

**Distribuția staționară a stărilor:** Din considerațiile de mai sus rezultă că matricea  $Q^n$  a unui lanț Markov absorbant este convergentă și limita sa este:

$$\Pi = \begin{pmatrix} I & O \\ NR & O \end{pmatrix}.$$

Linia  $i$  din matricea  $\Pi$  se numește distribuția staționară a stării  $i$ . Remarcăm că spre deosebire de lanțurile Markov ireductibile și aperiodice, unde fiecare stare avea aceeași

distribuție staționară  $\pi$  (vezi cursul precedent), în cazul lanțurilor Markov absorbante fiecare stare are altă distribuție staționară (de echilibru).

**Interpretarea elementelor matricei fundamentale  $N$ :** Elementul  $n_{ij} = N(i, j)$ ,  $i, j \in S_t$ , reprezintă numărul mediu de vizite pe care lanțul Markov ce pornește din starea  $i \in S_t$  îl face stării tranzitorii  $j \in S_t$  înainte de a fi absorbit.

**Exemplul 1.** Să se determine numărul mediu de treceri ale bețivului prin colțul 2 înainte de a ajunge acasă sau la bar, știind că a pornit din colțul 1.

**Rezolvare:** Din analiza matricei de tranziție, scrisă în forma standard, deducem că matricea  $T$  este

$$T = \begin{pmatrix} 0 & 0.5 & 0 \\ 0.5 & 0 & 0.5 \\ 0 & 0.5 & 0 \end{pmatrix}.$$

Matricea fundamentală a lanțului este

$$N = (I_3 - T)^{-1} = \begin{pmatrix} 1 & -0.5 & 0 \\ -0.5 & 1 & -0.5 \\ 0 & -0.5 & 1 \end{pmatrix}^{-1} = \begin{pmatrix} 1.5 & 1 & 0.5 \\ 1 & 2 & 1 \\ 0.5 & 1 & 1.5 \end{pmatrix}.$$

Cum stările tranzitorii ale lanțului Markov sunt 1, 2, 3 (în indexarea inițială), rezultă că numărul mediu de treceri prin colțul 2 înainte de a ajunge acasă sau la bar, știind că a pornit din colțul 1, este  $n_{12} = 1$ .

**Proprietatea 3.** Numărul mediu de pași ai lanțului Markov absorbant ce pleacă din starea tranzitorie  $i$  înainte de a fi absorbit este:

$$t_i = n_{i,n_a+1} + n_{i,n_a+2} + \cdots + n_{i,n_a+n_t},$$

unde  $\{n_a + 1, n_a + 2, \dots, n_a + n_t\}$  este mulțimea stărilor tranzitorii.

**Explicație:** Dacă stările tranzitorii ale unui lanț Markov absorbant sunt 3, 4, 5, iar 1, 2 sunt stări absorbante, atunci o traiectorie ce pornește din  $i = 3$  poate fi, de exemplu, de forma 3, 5, 5, 4, 3, 5, 4, 1. Numărul de pași parcurși până ce lanțul Markov este absorbit de starea 1 este dat de suma numărului de vizite ale stării 3, ale stării 4 și ale stării 5.

Numărul  $t_i$  se numește *timpul mediu până la absorbția lanțului Markov* ce pleacă din starea  $i$ .

Observăm că vectorul coloană  $\mathbf{t}$ , de coordonate  $t_i$ , se poate obține înmulțind matricea  $N$  cu vectorul coloană  $\mathbf{e}$ , ce are toate elementele 1, adică  $\mathbf{t} = N\mathbf{e}$ .

**Exemplul 2.** Numărul mediu de pași (treceri de la un colț la altul) ai bețivului înainte de a ajunge fie acasă, fie la bar, știind că a pornit din colțul 3, este  $t_3 = n_{31} + n_{32} + n_{33} = 3$ .

**Proprietatea 4.** Probabilitatea ca lanțul Markov ce pornește din starea tranzitorie  $i$  să fie absorbit de starea  $j \in S_a$  este

$$b_{ij} = (NR)(i, j) = \sum_{k \in S_t} N(i, k)R(k, j).$$

**Demonstrație:** Înainte de a fi absorbit, lanțul Markov ce pornește din  $i \in S_t$  parcurge 0 pași, 1 pas, 2 pași,  $\dots$ ,  $n$  pași în mulțimea stărilor tranzitorii  $S_t$ , și apoi într-un pas trece într-o stare absorbantă  $j$ . Notăm cu  $A_n(i, j)$ ,  $n \in \mathbb{N}$ , evenimentul ca lanțul ce pornește din starea tranzitorie  $i$  să facă  $n$  pași în  $S_t$  înainte de a fi absorbit de starea absorbantă  $j$ . Evident că pentru  $n \neq n'$ , evenimentele  $A_n(i, j)$ ,  $A_{n'}(i, j)$  sunt mutual exclusive. Probabilitatea evenimentului  $A_n(i, j)$  este

$$p_{i,j}(n) = P(A_n(i, j)) = \sum_{k \in S_t} T^n(i, k) R(k, j).$$

Probabilitatea ca lanțul ce pornește din starea tranzitorie  $i$  să fie absorbit de starea absorbantă  $j$  este

$$\begin{aligned} b_{ij} &= P\left(\bigcup_{n \in \mathbb{N}} A_n(i, j)\right) = \sum_{n=0}^{\infty} P(A_n(i, j)) \\ &= \sum_{n=0}^{\infty} \sum_{k \in S_t} T^n(i, k) R(k, j) = \sum_{k \in S_t} \left(\sum_{n=0}^{\infty} T^n(i, k)\right) R(k, j) \\ &= \sum_{k \in S_t} \lim_{n \rightarrow \infty} (T^0(i, k) + T^1(i, k) + \dots + T^{n-1}(i, k)) R(k, j) \\ &= \sum_{k \in S_t} N(i, k) R(k, j). \end{aligned}$$

□

Notând cu  $B$  matricea de elemente  $(b_{ij})$ , avem  $B = NR$ .

**Exemplul 3.** Probabilitatea ca bețivul ce pornește din colțul 2 să ajungă la bar (să fie absorbit de bar, codificat cu starea 0) este

$$b_{20} = \sum_{k=1}^3 N(2, k) R(k, 0) = N(2, 1) R(1, 0) + N(2, 2) R(2, 0) + N(2, 3) R(3, 0).$$

Ținând seama că matricea  $R$  este

$$R = \begin{pmatrix} R(1, 0) & R(1, 4) \\ R(2, 0) & R(2, 4) \\ R(3, 0) & R(3, 4) \end{pmatrix} = \begin{pmatrix} 0.5 & 0 \\ 0 & 0 \\ 0 & 0.5 \end{pmatrix},$$

avem că

$$b_{20} = 1 \times 0.5 + 2 \times 0 + 1 \times 0 = 0.5.$$

Atenție la indexarea elementelor matricei  $R$ :  $R(i, j)$ ,  $i \in S_t$ , iar  $j \in S_a$ . În cazul lanțului analizat mulțimea stărilor tranzitorii este  $S_t = \{1, 2, 3\}$ , iar cea a stărilor absorbante  $S_a = \{0, 4\}$ .