

INFORMATION RETRIEVAL SYSTEMS
EVALUATION METHODOLOGIES

*METODOLOGII DE EVALUARE A SISTEMELOR
DE INVESTIGARE A INFORMAȚIEI*

Viorica LUPU

*„Dacă nu poți măsura, nu poți îmbunătăți”
 (“If you can not measure it, you can not improve it.”)*

Lord Kelvin

Abstract: *Evaluation is an important step in the development of information retrieval systems, as it allows for the successful measurement of a system designed to help users meet their information goals. Several scholars have conducted a series of tests to calculate the efficiency of information retrieval systems. Among the developed methodologies, the Cranfield model remains the dominant approach in evaluating information systems. This paper presents the steps to be taken to test an information system, the indicators used to evaluate the effectiveness of information systems (recall and precision). Studying the evaluation and its role in the functioning of information retrieval tools can essentially contribute to improving the performance of these tools and the quality of their results.*

Key words: *information retrieval system, evaluation, efficiency, relevance, precision, recall*

Dezvoltarea rapidă a noilor tehnologii de informare și comunicare, precum și creșterea impresionantă a volumului de informații, date și documente disponibile face ca identificarea informației de care avem nevoie să devină o operațiune extrem de importantă, cu implicații majore în activitatea cotidiană, dar mai ales în cea academică ori de cercetare. Astfel, este imperios să găsim în termene optime conținuturi de o relevanță maximă, care să răspundă cel mai bine nevoilor noastre cognitive. Anume în acest scop au fost elaborate instrumentele de investigare a informațiilor (baze de date, sisteme de informare, cataloage etc.), pe care bibliotecile le implementează activ în activitățile lor informaționale.

În cadrul proiectului Tempus „Servicii informaționale moderne pentru îmbunătățirea calității studiilor” (MISISQ), șapte biblioteci universitare din Republica Moldova și-au unit eforturile și au creat catalogul colectiv partajat al bibliotecilor universitare, un sistem menit să asigure accesul operativ și de calitate al utilizatorilor la informație.

Una dintre cerințele înaintate față de acest produs informațional este descoperirea și livrarea informației de o plenitudine și precizie maxime. Descoperirea informațiilor presupune obținerea, cu ajutorul sistemului de informare, a resurselor relevante dintr-o colecție de resurse informaționale. Alcătuit din cataloagele a șapte biblioteci, funcționarea catalogului colectiv trezește mai multe întrebări și suspiciuni, în special în ceea ce privește performanța acestuia și relevanța rezultatelor pe care le oferă. Iată de ce este nevoie de a aprecia măsura în care sistemul răspunde necesităților de informare ale utili-

zatorilor. Fără o evaluare adecvată în acest sens nu putem stabili cât de bine funcționează un sistem, nu putem compara în mod obiectiv nivelul acestuia de investigare a informației cu cel al altor sisteme.

Scopul articolului de față nu este de a evalua eficiența unui sistem concret de informare, ci de a identifica și analiza unele metodologii de evaluare a instrumentelor de investigare a informației, metodologii care, în perspectivă, ar putea fi aplicate la măsurarea eficienței sistemelor de acest fel din Republica Moldova.

Evaluarea sistemelor de investigare a informației este una dintre cele mai mari provocări pentru specialiștii în științele informării. Determinarea performanței unui sistem este condiționată de aprecierea gradului de relevanță a documentelor furnizate de sistem în raport cu nevoile de informare ale utilizatorului. Evaluarea permite cuantificarea și măsurarea succesului unui sistem de investigare a informațiilor, fiind importantă pentru proiectarea, dezvoltarea și menținerea unor sisteme eficiente de informare.

Există două modalități de evaluare: evaluarea nemijlocită a sistemului și evaluarea bazată pe utilizator. Prima modalitate presupune măsurarea capacității sistemului de a clasifica documentele în funcție de relevanță, în timp ce a doua se axează pe percepția utilizatorului și măsoară gradul de satisfacție al acestuia în legătură cu funcționarea sistemului. Abordarea centrată pe utilizator examinează sarcina de căutare a informațiilor în contextul comportamentului uman, urmărind a înțelege complexitatea interacțiunii utilizatorului cu un sistem informatic. Se pornește de la premisa că înțelegerea comportamentului utilizatorului facilitează proiectarea mai eficientă a sistemului și stabilește criteriile de utilizare în evaluarea relațiilor dintre utilizator și sistem.

Cu toate avantajele pe care le prezintă, evaluarea bazată pe utilizator este extrem de costisitoare și dificil de realizat corect. Pentru o apreciere corespunzătoare este nevoie de un eșantion suficient de mare și reprezentativ de utilizatori efectivi ai sistemului de informare (a cărui rutină zilnică va fi întreruptă de procesul de evaluare), fiecare dintre sistemele supuse evaluării trebuie să fie la fel de complete și dezvoltate, cu o interfață optimă de căutare a informației, prietenoasă cu utilizatorul, iar fiecare subiect/utilizator trebuie să fie instruit la fel de bine în toate sistemele (Voorhees, E. M. 2002). Astfel de considerente îi determină pe cercetători să utilizeze opțiunea mai puțin costisitoare, adică evaluarea centrată pe sistem.

Testele de laborator sunt mult mai puțin costisitoare decât evaluările bazate pe utilizator, oferind în special informații de diagnosticare a comportamentului sistemului.

Primele experimente riguroase vizând evaluarea sistemelor de informare au fost realizate în cadrul Bibliotecii Colegiului Aeronautic Cranfield din Marea Britanie, sub conducerea bibliotecarului și omului de știință în domeniul informaticii Cyril Cleverdon. Acesta a inițiat o serie de experimente, numite Cranfield, prin care au fost puse bazele cercetării privind evaluarea sistemelor de investigare a informațiilor. Autorii și-au propus îmbunătățirea eficienței sistemelor studiate prin dezvoltarea unor limbaje și metode de indexare mai eficiente. Experimentele au fost efectuate în două etape principale. Prima etapă, numită în mod obișnuit Cranfield 1, a cuprins perioada anilor 1957–1961, iar cea de-a doua, Cranfield 2, s-a desfășurat între 1963 și 1966 (Cleverdon, Cyril 1967). Cleverdon și colegii săi s-au confruntat cu provocări serioase în realizarea experimentelor. Menționăm că experimentele nu au fost computerizate, ci efectuate manual, destul de laborios, cu indicii scrise manual, ca și cataloagele de scoruri și căutări efectuate prin încrucișarea informațiilor despre acești indici de scor.

Experimentele Cranfield au stabilit metodologia experimentală standard în domeniu: evaluarea sistemului în baza unei colecții fixe de teste – un set de documente, un set

de cereri de informare și un set de hotărâri de relevanță (evaluări ale documentelor).

La baza concepției lui Cleverdon stă noțiunea de relevanță, clasificată în două tipuri: relevanța semantică, ce presupune corespunderea conținutului informației livrate cu conținutul cererii de informare a utilizatorului, și relevanța formală, axată pe corespunderea modelului de investigare a documentului cu modelul de investigare a cererii de informare.

Experimentele de la Cranfield au culminat cu instituirea indicatorilor de precizie și plenitudine, care au devenit criterii clasice de evaluare a sistemelor de regăsire a informației.

Precizia exprimă capacitatea sistemului de a livra documentele relevante pentru o cerere, adică informația ce coincide cu conținutul cererii de informare. De exemplu, dacă precizia este de 50%, înseamnă că printre documentele găsite jumătate sunt relevante și jumătate irelevante.

Pentru calcularea preciziei informației livrate se propune următoarea formulă:

$$Pr = \frac{a}{a+b} \times 100\%$$

unde a este cantitatea de documente relevante livrate, iar b – cantitatea de documente livrate de către sistem.

Plenitudinea este raportul dintre numărul de documente relevante livrate de sistem și numărul total de documente relevante nelivrate de sistem:

$$Pl = \frac{a}{a+c} \times 100\%$$

De exemplu, dacă plenitudinea este de 50%, înseamnă că jumătate din documentele relevante nu sunt găsite de sistem.

Între indicatorii de precizie și plenitudine există o corelație strânsă. Odată cu creșterea excesivă a preciziei scade plenitudinea și invers, dacă ridicăm plenitudinea peste 90% precizia scade brusc, de aceea este foarte importantă existența unui echilibru între acești doi indicatori. Pentru a stabili acest echilibru este necesar mai întâi de toate de a determina nivelul de profunzime necesară indexării. Stabilirea acestui nivel trebuie să se realizeze în funcție de categoriile de utilizatori pentru care este destinat sistemul și de tipul bibliotecii în care este exploatat sistemul. Creșterea plenitudinii este posibilă numai în baza scăderii preciziei și invers.

Urmărind logica testelor Cranfield, modalitatea desfășurării lor și indicatorii folosiți, credem că această metodologie ar putea fi aplicată cu succes la evaluarea sistemelor de informare existente la noi.

Pentru a testa un sistem de informare și a compara strategiile de căutare este nevoie de a parcurge mai multe etape, prima dintre care este *pregătirea unei colecții de documente pentru a permite testarea semnificativă a rezultatelor*. Sistemele de informare indexează documente care sunt identificate ca răspuns la cererile utilizatorilor. O colecție de teste trebuie să conțină un set static de documente care să reflecte tipurile de documente ce pot fi găsite în setările sau domeniile

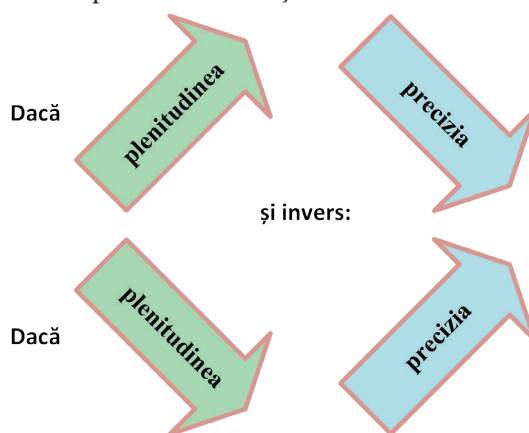


Fig. 1. Relația dintre precizie și plenitudine

operaționale. Acest lucru ar putea implica colecții de biblioteci digitale, cataloage sau seturi de pagini web. Noțiunea de colectare a documentelor statice este importantă deoarece asigură posibilitatea de a reproduce rezultatele la reutilizarea colecției de teste. Întrebările specifice care ar putea fi luate în considerare la colectarea documentelor includ:

- Câte documente ar trebui colectate? În cazul colecțiilor de biblioteci digitale mai mici, toate elementele ar putea fi incluse în colecția de teste. Totuși, în situații care implică colecții mai mari sau colecții în continuă schimbare, de exemplu pagini web, trebuie de colectat un eșantion de documente.
- Ce elemente ar trebui să fie eșantionate pentru a crea colecția de documente? Elementele din colecția de testare ar trebui să le reprezinte cât mai fidel posibil pe cele găsite într-un cadru operațional. Dacă elementele trebuie să fie selectate pentru o colecție de documente, atunci poate fi oportună selectarea unui eșantion aleatoriu sau a unui anumit subset de documente.

Modelarea unei aplicații reale de utilizator, cu nevoi realiste de informare, este o altă etapă de evaluare, în cadrul căreia se stabilește un set de cereri de informare tipice ale utilizatorilor. În baza acestora se creează modelul de investigare a informației (traducerea cererilor de informare într-un limbaj de indexare și întocmirea unui șir de cuvinte-cheie în baza cărora se va efectua căutarea informației). Fiecare cerere trebuie să reprezinte o necesitate reală de informare și trebuie să fie exprimată corect și fără ambiguitate. În practică, pentru obținerea unor rezultate de evaluare fiabile ar trebui incluse cel puțin 50 de cuvinte-cheie sau vedete de subiect.

Următoarea etapă, *investigarea informației și evaluarea relevanței*, presupune accesarea sistemului și adresarea cererilor de informare. În rezultat, pentru fiecare cerere sistemul va livra o cantitate de documente ce trebuie evaluată din punctul de vedere al relevanței.

Logica cercetării este următoarea: în orice masiv informațional al unui sistem, față de o cerere de informare există o parte de documente relevante. O parte din aceste documente este livrată utilizatorului, iar altă parte, din cauza unor imperfecțiuni tehnice, rămâne în sistem.

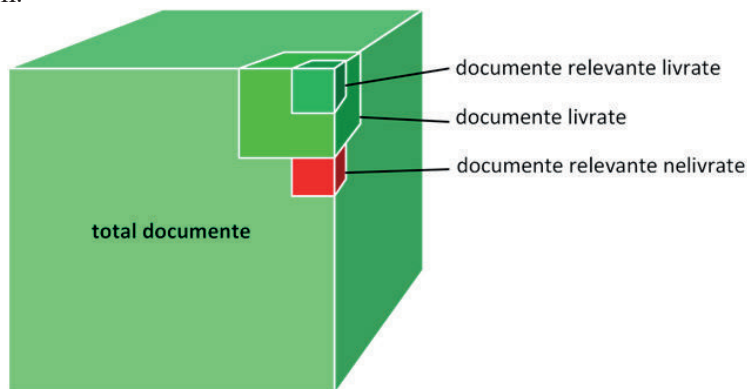


Fig. 2. Categoriile de documente în sistemul de informare
(figură adaptată după Hinkelmann, Knut)

Determinarea relevanței este extrem de importantă pentru evaluarea de ansamblu a sistemelor de regăsire a informației și este efectuată în permanență de experți în domeniu. Judecățile emise în cadrul acestui proces sunt afectate de multe caracteristici ale documentelor și utilizatorilor, precum și de factori situaționali (Harter, S.P. 1996).

Într-un studiu recent al procesului de evaluare a relevanței, cercetătorii Cuadra și Katter au identificat patru tipuri principale de factori care pot afecta rezultatul unei judecăți de relevanță: tipul de document care este evaluat/judecat, inclusiv subiectul, nivelul de dificultate, stilul etc.; condițiile în care trebuie să fie pronunțate judecățile, adică timpul disponibil, ordinea de prezentare și mărimea setului de documente, tipul de specificație a sarcinii etc.; declarația care specifică cerința de informare ce determină relevanța; evaluatorii folosiți pentru a face judecățile, adică experiența, cunoștințele, atitudinea lor etc. Aceste variabile sunt rezumate în tabelul din figura 3.

Trebuie să admitem totuși că aprecierea relevanței presupune o mare doză de subiectivism, întrucât chiar și aceeași persoană poate aprecia diferit în momente diferite, fiind influențată de numărul de documente ce trebuie evaluate, de ordinea în care ele sunt prezentate, de cunoștințele în cauză ale celui ce apreciază și, eventual, chiar de dispoziția lui.

Relevanța din perspectiva umană este subiectivă (depinde de opinia unui anumit utilizator), situațională (se referă la nevoile curente ale utilizatorului), cognitivă (depinde de percepția umană) și dinamică (modificări în timp) (Zuva, K., Zuva, T. 2012).

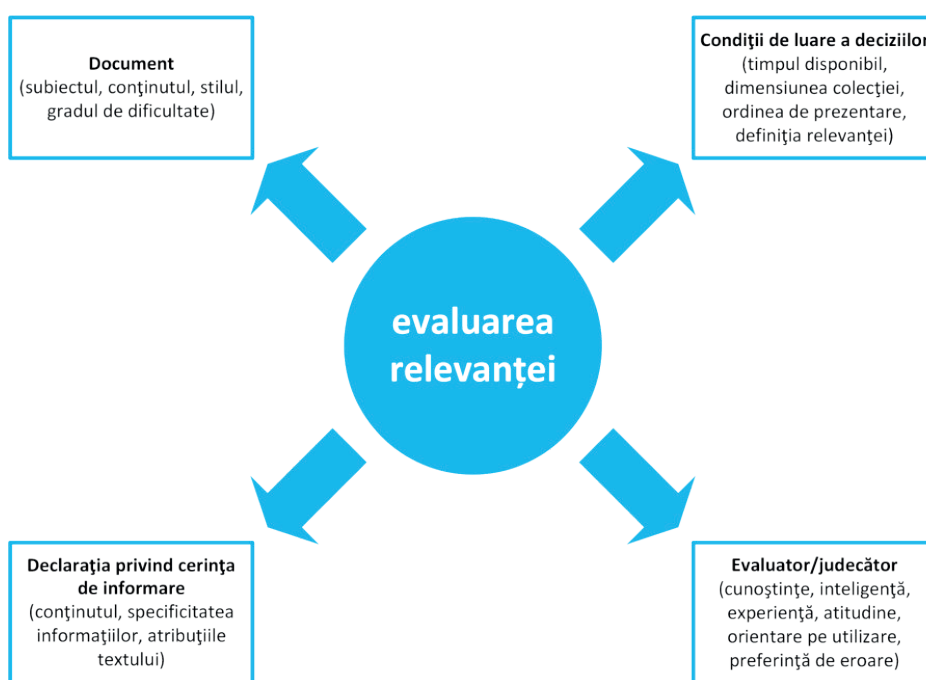


Fig. 3. Variabile privind evaluarea relevanței (Lesk, M.E., Salton, G. 1968)

Evaluatorii de relevanță vor emite judecăți de relevanță pentru fiecare document, acestea fiind ulterior utilizate pentru a calcula eficacitatea sistemului. Generarea unor asemenea evaluări este adesea extrem de consumatoare de timp și de intensitate a forței de muncă. Așa cum am arătat mai sus, judecățile de relevanță tind să difere în funcție de cine face evaluarea, iar valorile de plenitudine și de precizie obținute prin utilizarea acestor evaluări se pot dovedi apoi instabile.

Pentru a genera valori coerente de plenitudine și de precizie se va efectua mai întâi examinarea conținutului documentelor în vederea aprecierii gradului de coincidență a

informației livrate de sistem cu conținutul cererii de informare. Ulterior se va face raportarea documentelor livrate la cele nelivrate, precum și a documentelor livrate relevante la cele nerelevante.

În vederea asocierii relevanței cu livrarea informației poate fi utilizat tabelul de mai jos.

| | documente relevante | documente nerelevante | |
|----------------------------|---------------------|-----------------------|---------|
| documente livrate | a | b | a+b |
| documente nelivrate | c | d | c+d |
| | a+c | b+d | a+b+c+d |

Utilizând acest tabel putem determina indicatorii de eficiență a sistemului de informare (precizie și plenitudine).

Să presupunem că pentru cererea de informare C1, sistemul a livrat un număr total de 100 de documente, dintre care 65 sunt relevante și 35 nerelevante. În rezultatul căutării, în masivul de documente s-a depistat că referitor la cerința respectivă în colecție mai există un număr de 20 de documente relevante, dar pe care sistemul din mai multe motive nu le-a livrat. Astfel, pentru cererea de informare C1 plenitudinea va constitui 76,4%, iar precizia – 65%.

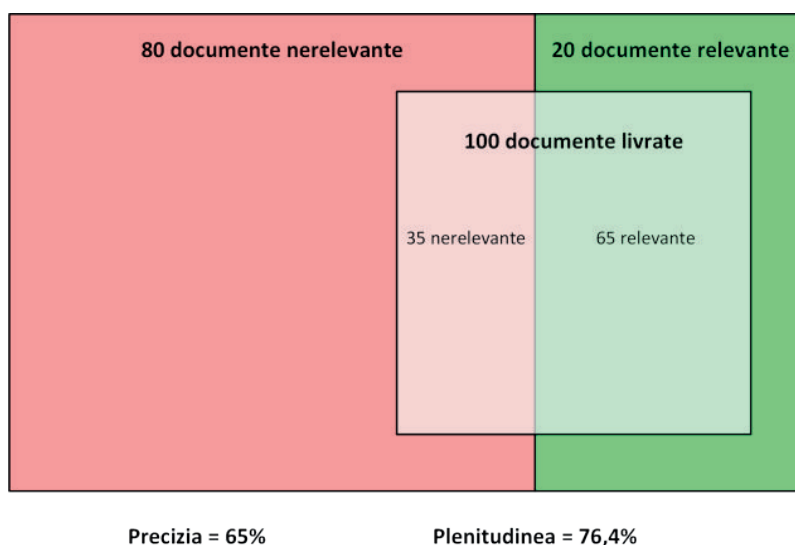


Fig. 4. Matricea precizie-plenitudine (adaptată după James D. McCaffrey 2016)

Testele Cranfield au introdus o nouă paradigmă în evaluarea sistemelor de informare și continuă să fie un model experimental dominant utilizat în eforturile de evaluare ale mai multor grupuri internaționale precum SMART (System for the Mechanical Analysis and Retrieval of Text), Text Retrieval Conference (TREC), Cross-Language Evaluation Forum (CLEF) și NII-NACIS Test Collection for IR Systems (NTCIR).

Proiectul SMART a adaptat această metodologie la mediul computerizat și a promovat un model de cercetare și publicare bazat pe un experiment automatizat. Acesta a fost inițiat de Gerard Salton și colaboratorii săi la începutul anilor 1960, mai întâi la Universitatea Harvard, apoi la Universitatea Cornell, și a continuat într-o formă sau alta până la sfârșitul anilor 1990. Evaluarea SMART a folosit modelul dezvoltat la Cranfield, iar după ce a devenit disponibilă în formă ușor de citit, colecția Cranfield a fost inclusă în suita de testare SMART. Noutatea pe care o aduce modelul SMART este importantă atât pentru metoda de regăsire, cât și pentru metodologia de evaluare și presupune că răspunsul la o căutare ar trebui să fie mai degrabă o listă de documente clasificate în funcție de gradul (sau probabilitatea) de relevanță decât un set neordonat de potriviri exacte.

Întregul potențial pentru evaluarea la scară largă a procesului de regăsire a informației nu a fost realizat însă până la producerea de colecții de teste TREC. Proiectul TREC, inițiat în 1992 la Institutul Național de Standarde și Tehnologie (NIST, o agenție guvernamentală a Statelor Unite), a avut drept motivație principală necesitatea de a crea colecții de dimensiuni realiste și de a consolida și extinde cercetarea în tehnologia de regăsire a informației prin utilizarea datelor experimentale comune și de înaltă calitate, precum și a tehnicilor standard de evaluare (Prange, 1996; Harman, 1992b). A fost un exercițiu experimental de mare anvergură, colaborativ și comparativ, cu participarea a peste 250 de grupuri internaționale de cercetare (Voorhees și Harman, 2005a), având rezultate inovatoare nu doar în ceea ce privește mărimea colecției, ci și cu referire la metoda experimentală. Impactul TREC asupra cercetării de căutare a informațiilor a fost la fel de semnificativ. Au fost construite colecții mari de teste atât pentru regăsirea tradițională ad-hoc, cât și pentru noi sarcini, cum ar fi regăsirea în limbajul transversal, regăsirea în limbaje controlate (vocabulare). Proiectul TREC a standardizat metodologia de evaluare utilizată pentru a evalua calitatea rezultatelor de regăsire, atât prin dimensiunea colecției de date, cât și prin eficacitatea metodologiei (Voorhees, Ellen M. 2006).

În lipsa unui consens la aplicarea metodelor alternative de evaluare a sistemelor informative, modelul Cranfield rămâne abordarea de bază în domeniul regăsirii informației, fiind aplicată în majoritatea proiectelor la această sferă.

Referințe bibliografice:

1. cLEVERDON, Cyril (1967). The Cranfield tests on index language devices. In: *Aslib Proceedings*, Vol. 19, Issue 6, pp. 173-194. doi:<https://doi.org/10.1108/eb050097>.
2. HINKELMANN, Knut. Information Retrieval and Knowledge Organisation [online]. Available at: http://www.hinkelmann.ch/knut/lectures/irko/IRKO-InformationRetrieval_eval.pdf.
3. HARTER, Stephen P. (1996). Variations in Relevance Assessments and the Measurement of Retrieval Effectiveness. In: *Journal of the American Society for Information Science*, vol. 47(1), pp. 37-49. Available at: https://asis.org/Publications/JASIS/Best_Jasist/1997Harter.pdf
4. LESK, M. E., SALTON, G. (1968). Relevance assessments and retrieval system evaluation. In: *Information Storage and Retrieval* [online], vol. 4, issue 4, pp. 343-359. Available at: [https://doi.org/10.1016/0020-0271\(68\)90029-6](https://doi.org/10.1016/0020-0271(68)90029-6).
5. SALTON, Gerard (). The Smart environment for retrieval system evaluation-advantages and problem areas [online]. Available at: http://sigir.org/files/museum/Information_Retrieval_Experiment/pdfs/p316-salton.pdf

6. VOORHEES, Ellen M. (2002) The Philosophy of Information Retrieval Evaluation. In: PETERS C., BRASCHLER M., GONZALO J., KLUCK M. (eds) *Evaluation of Cross-Language Information Retrieval Systems*. CLEF 2001. Lecture Notes in Computer Science [online], vol 2406. Springer, Berlin, Heidelberg. Available at: https://link.springer.com/chapter/10.1007/3-540-45691-0_34#citeas
7. VOORHEES, Ellen M. (2006). TREC: Improving information access through evaluation. In: *Bulletin of the American Society for Information Science and Technology* [online], vol. 32, pp. 16-21. doi:10.1002/bult.2003.1720320105. Available at: <https://onlinelibrary.wiley.com/doi/full/10.1002/bult.2003.1720320105>
8. ZUVA, K., ZUVA, T. (2012). EVALUATION OF INFORMATION RETRIEVAL SYSTEMS. In: *International Journal of Computer Science & Information Technology* (IJCSIT), Vol 4, No 3, pp. 35-43. DOI : 10.5121/ijcsit.2012.4304 35