

HUJJATLARNI SO‘ROVGA TEGISHLILIK DARAJASIGA KO‘RA
BAHOLASH

Do‘shanov Bekzod Davronbek o‘g‘li

Toshkent Axborot Texnologiyalar Universiteti magistrant,

duschanovbekzod17@gmail.com

Ochilboyev Umidjon Ilxom o‘g‘li

Toshkent Axborot Texnologiyalar Universiteti magistrant, canc41946@gmail.com

Shamuratov Ulug‘bek Alisher uli

Toshkent Axborot Texnologiyalar Universiteti magistrant,

ulugbekshamuratov1219@gmail.com

Annotatsiya: *Ushbu maqola hujjatlarni so‘rovga tegishlilik darajasiga ko‘ra baholash muammosini o‘rganadi. Maqolada tegishlilik baholashning asosiy mezonlari, mavjud usullar va ularning samaradorligi tahlil qilinadi, axborot qidirish tizimlarida hujjat-so‘rov mosligini aniqlashning zamonaviy yondashuvlari ko‘rib chiqiladi.*

Kalit so‘zlar: *Hujjat, so‘rov, indeks, TF-IDF, kosinus o‘xshashlik, vektor, atama.*

Аннотация: *Данная статья исследует проблему оценки документов по степени релевантности запросу. В статье анализируются основные критерии оценки релевантности, существующие методы и их эффективность, а также рассматриваются современные подходы к определению соответствия документа запросу в системах информационного поиска.*

Ключевые слова: *документ, запрос, индекс, TF-IDF, косинусное сходство, вектор, термин,*

Abstract: *This article examines the problem of evaluating documents according to their degree of relevance to a query. The article analyzes the main criteria for relevance assessment, existing methods and their effectiveness, and reviews modern approaches to determining document-query matching in information retrieval systems.*

Keywords: *document, query, index, TF-IDF, cosine similarity, vector, term.*

KIRISH

Hujjatlarni so‘rovga tegishlilik darajasiga ko‘ra tartiblash bu axborot qidiruv tizimlarida foydalanuvchi kiritgan so‘rovga eng mos keladigan hujjatlarni aniqlash va ularni moslik darajasiga qarab tartiblash jarayonidir. Bu jarayon bir necha bosqichlardan iborat bo‘ladi. Avvalo foydalanuvchi so‘rovi tahlil qilinadi, ya‘ni so‘rovdagi kalit so‘zlar va ularning ma‘nosi aniqlanadi. So‘ngra bazadagi hujjatlar orasidan shu kalit so‘zlar uchraydigan yoki mazmunan o‘xshash hujjatlar topiladi. Topilgan hujjatlar har biri uchun so‘rovga qanchalik mos ekanini aniqlash maqsadida tegishlilik darajasi hisoblanadi. Tegishlilikni baholash uchun TF-IDF, yoki semantik tahlil kabi usullar qo‘llaniladi. Har bir hujjatga tegishlilik balli berilgach, ular shu ballar bo‘yicha tartiblanadi, ya‘ni eng yuqori

ball olgan hujjatlar ro‘yxatda yuqorida, past ball olganlari esa pastda joylashadi. Natijada foydalanuvchi uchun eng muhim va kerakli hujjatlar birinchi o‘rinda ko‘rsatiladi. Shu tarzda hujjatlarni so‘rovga tegishlilik darajasiga ko‘ra tartiblash foydalanuvchiga axborotni tez va aniq topish imkonini beradi.

Parametrik indeks bu hujjatlar tarkibida alohida parametrlar yoki atributlar, ya‘ni muallif, sana, mavzu, til, hujjat turi kabi belgilar bo‘yicha tuziladigan indeks hisoblanadi. Parametrik indekslar yordamida foydalanuvchi o‘z so‘rovini aniq chegaralashi va kerakli ma‘lumotni tez topishi mumkin. Masalan, foydalanuvchi “Muallifi A. Bozorov bo‘lgan 2023-yilgi sun‘iy intellekt haqidagi hujjatlar” degan so‘rovni kiritganda tizim parametrik indekslardan foydalanib aynan shu shartlarga mos hujjatlarni topadi. Bu usul ma‘lumotlarni aniq filtrlash va izlash jarayonini tezlashtiradi.

Zona indeksi esa hujjat matnini ma‘lum zonalarga yoki bo‘limlarga ajratish va har bir zonasini alohida indekslash usulidir. Odatda hujjat sarlavha, annotatsiya, matn va kalit so‘zlar kabi zonalarga bo‘linadi. Bu usulning afzalligi shundaki, tizim hujjatdagi so‘zning joylashuviga qarab uning ahamiyatini aniqlay oladi. Masalan, so‘rovdagi kalit so‘z sarlavhada uchrasa, u matn ichidagi so‘zga qaraganda muhimroq hisoblanadi. Shu yo‘l bilan hujjatning foydalanuvchi so‘roviga tegishlilik darajasi yanada aniq baholanadi.

Vaznli zona reytingi esa zona indeksleri asosida hujjatning tegishlilik darajasini vaznli koeffitsientlar orqali baholash usulidir. Har bir zona uchun alohida vazn belgilanadi, masalan sarlavha zonasi uchun 0.5, annotatsiya uchun 0.3, matn uchun esa 0.2 kabi. So‘ngra har bir zonada so‘rovdagi so‘zlar uchrashish chastotasiga qarab ballar hisoblanadi va umumiy reyting aniqlanadi. Shu tariqa, tizim eng muhim zonalarda kerakli so‘zlar uchraydigan hujjatlarni yuqoriroq o‘ringa qo‘yadi. Parametrik indeks hujjatlarni atributlar bo‘yicha tartiblash imkonini beradi, zona indeksi matn ichidagi so‘zlarning joylashuvini hisobga oladi, vaznli zona reytingi esa shu zonalarning ahamiyatini turlicha baholab hujjatning umumiy tegishlilik darajasini aniqlaydi.

Axborot qidiruv tizimlarida foydalanuvchi so‘roviga mos hujjatlarni aniqlashda atamalarning hujjatdagi ahamiyatini hisoblash muhim ahamiyatga ega. Buning uchun ko‘pincha TF-IDF (Term Frequency - Inverse Document Frequency) modeli qo‘llaniladi. Ushbu model ikki asosiy ko‘rsatkichdan iborat: atamalar chastotasi (TF) va teskari hujjat chastotasi (IDF).

Atamalar chastotasi (Term Frequency, TF)

Bu ko‘rsatkich ma‘lum bir atamaning (so‘zning) hujjatdagi uchrashish chastotasini bildiradi. Atama hujjatda qanchalik ko‘p uchrasa, u so‘zning shu hujjatdagi ahamiyati shunchalik yuqori bo‘ladi.

$$TF(t, d) = \frac{f_{t,d}}{\sum_k f_{k,d}}$$

bu yerda:

$f_{t,d}$ - t atamaning d hujjatda uchrashish soni;

$\sum_k f_{k,d}$ - hujjatdagi jami so‘zlar soni.

15-May, 2026-yil

Masalan, agar “sun’iy” so‘zi 100 so‘zli hujjatda 5 marta uchrasa, u holda:

$$TF = \frac{5}{100} = 0.05$$

Atamaning teskari hujjat chastotasi (Inverse Document Frequency, IDF). Bu ko‘rsatkich so‘zning butun hujjatlar to‘plamidagi noyoblik darajasini aniqlaydi. Ko‘p hujjatlarda uchraydigan so‘zlar kam ahamiyatga ega, kam hujjatlarda uchraydigan so‘zlar esa muhimroq hisoblanadi.

$$IDF(t) = \log \frac{N}{n_t}$$

bu yerda:

N - umumiy hujjatlar soni,

n_t - t atama uchragan hujjatlar soni.

Agar bir so‘z barcha hujjatlarda uchrasa, $\frac{N}{n_t} = 1$ bo‘ladi va IDF qiymati kichik bo‘ladi; agar u kam hujjatda uchrasa, IDF qiymati kattalashadi.

TF va IDF kombinatsiyasi (TF-IDF vaznlash). So‘zning hujjatdagi umumiy ahamiyatini aniqlash uchun TF va IDF qiymatlari bir-biriga ko‘paytiriladi:

$$TF-IDF(t, d) = TF(t, d) * IDF(t)$$

Bu formula orqali har bir so‘zning hujjatdagi vazni aniqlanadi. Natijada, hujjatda ko‘p uchraydigan, lekin umumiy korpusda kam uchraydigan so‘zlar yuqori bahoga ega bo‘ladi. Natijada tizim bu hujjatni foydalanuvchi so‘roviga eng mos deb hisoblaydi. TF hujjat ichidagi so‘zning ahamiyatini, IDF esa so‘zning butun korpusdagi noyoblik darajasini o‘lchaydi. Ularning kombinatsiyasi bo‘lgan TF-IDF modeli esa hujjatlarni reytinglashda eng samarali vaznlash usullaridan biri bo‘lib, axborot qidiruv tizimlari, matn tahlili va tabiiy tilni qayta ishlashda keng qo‘llaniladi.

Axborot-qidiruv tizimlarida asosiy maqsad foydalanuvchi so‘roviga eng mos keladigan hujjatlarni topishdir. Buning uchun hujjatlar va foydalanuvchi so‘rovi orasidagi o‘xshashlik darajasi aniqlanadi. O‘xshashlikni baholash hujjatlar orasidagi mazmuniy yaqinlik yoki matematik masofani o‘lchash orqali amalga oshiriladi. Hujjatlar odatda vektorli model asosida taqqoslanadi. Bu modelda har bir hujjat va so‘rov atamalardan tashkil topgan vektor ko‘rinishida tasvirlanadi. Har bir o‘lcham (komponent) hujjatda ma’lum bir atamaning vaznini ifodalaydi. Vektor komponentlari odatda TF-IDF usuli orqali hisoblangan vaznlardir.

$$\text{cosine_similarity}(A, B) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

A - hujjat vektori, B - so‘rov vektori, (bu vektorlar matnlar, so‘zlar yoki jummalarning savoldagi vector qiymatlari bo‘lishi mumkin)

$A \cdot B$ - ikki vector skalyar ko‘paytmasi

$\|A\| * \|B\|$ - ikki vector yevklid uzunligi

$$A * B = \sum t(TF_{t,A} \times TF_{t,B})$$

$$\|A\| = \sqrt{\sum t(TF_{t,A})^2} = \sqrt{n_A \times (TF_A)^2}$$

Kosinus o‘xshashligi natijasi har doim 0 dan 1 gacha bo‘lgan qiymatni oladi. Agar qiymat 1 ga yaqin bo‘lsa, hujjat so‘rovga juda o‘xshash; 0 ga yaqin bo‘lsa, o‘xshashlik past degani.

Amaliy qism

Hujjat 1:

Machine learning is a method of data analysis that automates analytical model building.

Hujjat 2:

Machine learning is a branch of artificial intelligence and based on the idea that systems can learn from data.

Hujjat 3:

Machine learning is a branch of Artificial intelligence focused on building computer systems that learn from data.

Hujjat 1 – so‘zlar soni 13

Hujjat 2 – haqida, so‘zlar soni 19

Hujjat 3 – so‘zlar soni 17

$$TF(t,d) = f(t,d) / N = 1/13 = 0.0769$$

$$TF(t,d) = f(t,d) / N = 1/19 = 0.0526$$

$$TF(t,d) = f(t,d) / N = 1/17 = 0.0588$$

Kosinus o‘xshashlik formulasi:

$$\cos(\theta) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

A – hujjat vektori, B – so‘rov vektori

A*B – ikki vector skalyar ko‘paytmasi

$\|A\| * \|B\|$ - ikki vektor yevklid uzunligi

$$A * B = \sum t(TF_{t,A} \times TF_{t,B})$$

$$\|A\| = \sqrt{\sum t(TF_{t,A})^2} = \sqrt{n_A \times (TF_A)^2}$$

Har bir juftlikda umumiy so‘zlar sonini topamiz.

(1 - 2) umumiy so‘zlar machine, learning, is, a, of, data, that, jami 7 ta so‘z

$$A \cdot B = 7 * (0.0769 * 0.0526) = 7 * 0.004037 = 0.02834$$

$$\|A\| = \sqrt{13 * (0.0769)^2} = \sqrt{0.07735} = 0.27735$$

$$\|B\| = \sqrt{19 * (0.0526)^2} = \sqrt{0.05265} = 0.22942$$

(1 - 3) umumiy so‘zlar machine, learning, is, a, of, data, that, jami 7 ta so‘z

$$A \cdot B = 7 * (0.0769 * 0.0588) = 7 * 0.005171 = 0.03620$$

$$\|A\| = \sqrt{13 * (0.0769)^2} = \sqrt{0.07735} = 0.27735$$

$$\|B\| = \sqrt{17 * (0.0588)^2} = \sqrt{0.05883} = 0.24254$$

(2 - 3) umumiy so‘zlar machine, learning, is, a, of, data, that, branch, artificial, intelligence, on, systems, learn, from, jami 14 ta so‘z

$$A \cdot B = 14 * (0.0526 * 0.0588) = 14 * 0.003097 = 0.04334$$

$$\|A\| = \sqrt{19 * (0.0526)^2} = \sqrt{0.05265} = 0.22942$$

$$\| B \| = \sqrt{17 * (0.0588)^2} = \sqrt{0.05883} = 0.24254$$

1-jadval.

Juftlik	A*B	A	B	$\cos(\theta)$
1-2	0.02834	0.27735	0.22941	0.4454
1-3	0.03619	0.27735	0.24253	0.5381
2-3	0.04334	0.22941	0.24253	0.7790

Eng yuqori o‘xshashlik: 2-hujjat va 3-hujjat - 0.7790. Ya’ni ular bir-biriga eng o‘xshash. 1-hujjat 2- va 3-hujjatlarga nisbatan kamroq o‘xshash (0.4454 va 0.5381). Sababi: 2- va 3-hujjatlarda birgalikda uchraydigan (va 1-hujjatda yo‘q) atamalar ko‘proq bo‘lgan (masalan: branch, artificial, intelligence, systems, learn, from), shuning uchun ularning TF-vektorlari bir-biriga yaqinroq.

Xulosa

Hujjatlarni to‘g‘ri baholash metodologiyasi qidiruv tizimlarining aniqligini sezilarli darajada oshiradi. Ushbu maqolada hujjatlarni so‘rovga tegishlilik darajasiga ko‘ra baholashning asosiy usullari va yondashuvlari ko‘rib chiqildi. Hujjatlarning so‘rovlarga tegishlilikini hisoblash davomida TF-IDF va kosinus o‘xshashlik kabi vektorli modellar axborot qidirish tizimlarida hujjat-so‘rov mosligini aniqlashda samarali vosita ekanligi isbotlandi. Indekslov jarayoni va atama chastotasiga asoslangan baholash usullari qidiruv natijalarining aniqligini sezilarli darajada oshirishi olingan natijalar yordamida isbotlandi.

Olingan natijalar shuni ko‘rsatadiki, tegishlilik baholashning to‘g‘ri metodologiyasini tanlash tizim samaradorligiga bevosita ta’sir qiladi. Kelajakda ushbu sohada neyron tarmoq va sun’iy intellekt asosidagi yondashuvlarni qo‘llash orqali yanada yuqori natijalar qo‘lga kiritish mumkin. O‘tkazilgan tadqiqot natijalari TF-IDF modeli hamda jumlar o‘rtasidagi o‘xshashlikni aniqlash usullarining resurs jihatidan cheklangan tillarda, jumladan, o‘zbek tilida ham yuqori samaradorlik ko‘rsatishini tasdiqlaydi. Mazkur metodologiya davlat boshqaruvi, ta’lim tizimi va tibbiyot kabi turli ijtimoiy sohalarida keng miqyosda qo‘llanilish imkoniyatiga ega. Bundan tashqari, tizim imkoniyatlarini kengaytirish maqsadida bilim grafiklarini integratsiya qilish va transformer arxitekturasiga asoslangan modellardan foydalanish kelgusidagi eng istiqbolli tadqiqot yo‘nalishlari sifatida belgilab olingan.

FOYDALANILGAN ADABIYOTLAR:

1. Vektor fazo modeli hamda jumlar o‘xshashligi o‘lchovlariga asoslangan savol – javob tizimi ishlab chiqish, Musayev Muhammadjon Mahmudovich, Ochilov Mannon Musinovich, Xolmatov Orzimurod, Abjalolovich, Narzullayev Oybek Otabek o‘g‘li, Raqamli Transformatsiya va Sun’iy Intellekt ilmiy jurnali, VOLUME 3, ISSUE 1, FEBRUARY 2025, ISSN: 3030-3346

2. Fan, H., and Qin, Y. (2018). “Research on Text Classification Based on Improved TF-IDF Algorithm” International Conference on Network, Communication, Computer Engineering (NCCE 2018), vol. 147

15-May, 2026-yil

3. Understanding Inverse Document Frequency: On theoretical arguments for IDF, Stephen Robertson, Microsoft Research, 7 JJ Thomson Avenue Cambridge CB3 0FBUK, (and City University, London, UK).
4. Manning, C.D., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge: Cambridge University Press.
5. Joachims, T. (1997). *A Probabilistic Analysis of the Rocchio Algorithm with TF-IDF for Text Categorization*. *ICML*, 143–151.
6. Mannon Ochilov, Dilshod Jurayev, Oybek Narzullayev, Matndagi tinish belgilarini tiklash muammolari va yechimlari, *Computer linguistics: problems, solutions*. 2024.
7. Zhao, K., & Mao, X.-L. Integrating ROUGE into Sentence Similarity Measures, *Journal of Chinese Information Processing*, 2017.
8. Beel, J., Gipp, B., Langer, S., Breiting, C. "Research-paper recommender systems: a literature survey", *International Journal on Digital Libraries*, 2016.