



# IJTIMOYIY-GUMANITAR SOHADA ILMIY-INNOVATSION TADQIQOTLAR

ILMIY METODIK JURNALI

ISSN 3060-5059



**VOL.3 № 4**

**2026**

## **FAYLDAGI OBYEKT LARNI SEMANTIK TAHLIL QILISH ALGORITMLARINI ERGONOMIK METODIKASI**

**Xolmatov Javlon Yusupovich**

O'zbekiston Milliy universiteti Jizzax filiali, katta o'qituvchi

**Mamaraimov Abror Kamalidin o'g'li**

O'zbekiston Milliy universiteti Jizzax filiali, katta o'qituvchi

### **Annotatsiya**

Ushbu maqolada PDF formatidagi elektron adabiyotlar va fayllar ichidan kerakli ma'lumotlarni tezkor topish va ma'lumotlarini semantik jihatdan tahlil qilish masalasi ko'rib chiqiladi. Sun'iy intellekt asosida ishlab chiqilgan semantik qidiruv algoritmi taklif etiladi. Ushbu tadqiqotda matnlarni semantik tahlil qilish uchun zamonaviy axborot qidiruv va tabiiy tilni qayta ishlash usullari qo'llanildi. Olingan natijalar shuni ko'rsatdiki, chuqur o'rganish modellari klassik usullarga nisbatan yuqori aniqlikni ta'minlaydi.

**Kalit so'zlar:** sun'iy intellekt, semantik qidiruv, TF-IDF, Word2Vec, talaba javoblarini tekshirish, PDF qidiruv, matn o'xshashligi, kosinus o'xshashligi.

## **ЭРГОНОМИЧЕСКАЯ МЕТОДОЛОГИЯ АЛГОРИТМОВ ДЛЯ СЕМАНТИЧЕСКОГО АНАЛИЗА ОБЪЕКТОВ В ФАЙЛЕ**

**Холматов Жавлон Юсупович**

Джизакский филиал Национального университета Узбекистана, старший преподаватель

**Мамараймов Аброр Камалидин угли**

Джизакский филиал Национального университета Узбекистана, старший преподаватель

### **Аннотация**

В данной статье рассматривается вопрос быстрого поиска необходимой информации в электронной литературе и файлах в формате PDF, а также семантического анализа их данных. Предложен алгоритм семантического поиска, разработанный на основе искусственного интеллекта. В данном исследовании для семантического анализа текстов использовались современные методы информационного поиска и обработки естественного языка. Полученные результаты показали, что модели глубокого обучения обеспечивают более высокую точность по сравнению с классическими методами.

**Ключевые слова:** искусственный интеллект, семантический поиск, TF-IDF, Word2Vec, проверка ответов студентов, поиск в PDF-файлах, сходство текста, косинусное сходство.

## **ERGONOMIC METHODOLOGY OF ALGORITHMS FOR SEMANTIC ANALYSIS OF OBJECTS IN A FILE**

**Kholmatov Javlon Yusupovich**

Jizzakh Branch of the National University of Uzbekistan, Senior Lecturer

**Mamaraimov Abror Kamalidin ugli**

Jizzakh Branch of the National University of Uzbekistan, Senior Lecturer

### **Abstract**

This article examines the issue of quickly searching for relevant information in electronic

literature and PDF files, as well as the semantic analysis of their data. A semantic search algorithm developed using artificial intelligence is proposed. This study utilized modern methods of information retrieval and natural language processing for semantic text analysis. The results obtained demonstrated that deep learning models provide higher accuracy compared to classical methods.

**Keywords:** artificial intelligence, semantic search, TF-IDF, Word2Vec, student response validation, PDF search, text similarity, cosine similarity.

Zamonaviy ta'lim jarayonida raqamli o'quv qo'llanmalar keng tarqalmoqda. Talabalarga PDF formatida adabiyotlar berilib, ularning shu nazariy ma'lumotlarni o'rganilganlik darajasini bilish uchun shular orqali tuzilgan nazariy savollar va topshiriqlar beriladi. Ammo millionlab talaba javoblarini qo'lda tekshirish imkonsiz[6].

PDF (Portable Document Format) fayli — bu matn, grafik, tasvir, shriftlar va giper murojaatlarni qurilmadan qat'iy nazar bir xil ko'rinishda saqlaydigan formatdir. U ob'ektlarga asoslangan tuzilmani tashkil etadi. Mavjud avtomatik tekshirish tizimlari faqat kalit so'zlar bo'yicha ishlaydi, semantik ma'noni tushunmaydi. Masalan, talaba “kuchli imperiya” yozsa, kitobda “qudratli davlat” deb yozilgan bo'lsa – tizim xato deb hisoblaydi.

Maqola maqsadi – PDF fayl matni ichidagi javob bilan talaba javobini semantik jihatdan solishtiradigan algoritmlarni tahlil qilish va ergonomik metodikani ishlab chiqish.

#### Adabiyotlar tahlili

Axborot qidiruv tizimlarida klassik usullar:

- Linear Search – barcha matnni ketma-ket skan qilish  $O(n)$
- Binary Search – saralangan ma'lumotlarda  $O(\log n)$
- Inverted Index – kalit so'zlar indeksi  $O(1)$  [1],[5],[7].

#### 1-jadval.

Xususiyat	Linear Search (Chiziqli qidiruv)	Binary Search (Ikkilik qidiruv)	Inverted Index (Teskari indeks)
<b>Ishlash prinsipi</b>	Ma'lumotlar boshidan oxirigacha ketma-ket tekshiriladi	Qidiruv maydoni har qadamda yarmiga bo'linadi va o'rtadagi element tekshiriladi	Har bir kalit so'z uchun alohida indeks yaratiladi: so'z → hujjatlar ro'yxati
<b>Vaqt murakkabligi</b>	$O(n)$	$O(\log n)$	$O(1)$
<b>Talab qilinadigan shart</b>	Maxsus shart yo'q, tartibsiz ma'lumotlarda ham ishlaydi	Ma'lumotlar saralangan bo'lishi kerak	Indeks oldindan yaratilgan bo'lishi kerak
<b>Afzalligi</b>	Algoritm juda oddiy, implementatsiya qilish oson	Juda tez ishlaydi	Juda tez qidiruv, millionlab hujjatlarda ham samarali
<b>Kamchiligi</b>	Katta hajmdagi ma'lumotlarda sekin	Avval ma'lumotni saralash kerak	Indeks yaratish va yangilash uchun qo'shimcha xotira va vaqt kerak
<b>Misol</b>	1 mln so'z ichidan “apple” ni topish uchun 1 mln tekshiruv bo'lishi mumkin	1 mln saralangan so'zda taxminan 20 ta solishtirish yetarli	“apple” so'zi qaysi hujjatlarda borligi indeksdan darhol olinadi

Sun'iy intellekt yondashuvlari:

- TF-IDF – so'z muhimligini hisoblash

- Word2Vec – soʻzlar vektorlari
- BERT – chuqur semantic tahlil [3],[4],[8],[9].

**2-jadval.**

<b>Xususiyat</b>	<b>TF-IDF</b>	<b>Word2Vec</b>	<b>BERT</b>
<b>Asosiy gʻoya</b>	Soʻzning hujjat ichida qanchalik muhimligini statistik hisoblash	Soʻzlarni zich vektorlar (embedding) ga aylantirish	Kontekstga bogʻliq chuqur semantik vektorlar yaratish
<b>Ishlash prinsipi</b>	Soʻz chastotasi (TF) va hujjatlar orasidagi noyoblik (IDF) asosida baholanadi	Oʻxshash kontekstda kelgan soʻzlar oʻxshash vektor oladi	Transformer modeli yordamida soʻzlar kontekst bilan birga tahlil qilinadi
<b>Asosiy formula / model</b>	$TF-IDF = TF \times IDF$ $IDF = \log(N / (df + 1))$	Neyron tarmoq asosidagi embedding modeli (CBOW, Skip-gram)	Transformer arxitekturasi asosidagi chuqur neyron tarmoq
<b>Semantikani tushunishi</b>	Past (faqat statistik)	Oʻrtacha (semantik oʻxshashlikni ushlaydi)	Juda yuqori (kontekstni tushunadi)
<b>Kontekstni hisobga olish</b>	Yoʻq	Yoʻq (har soʻz uchun bitta vektor)	Ha (har jumlada boshqa vektor boʻlishi mumkin)
<b>Hisoblash murakkabligi</b>	Juda tez	Oʻrtacha (modelni oʻqitish kerak)	Ogʻir va resurs talab qiladi
<b>Afzalliklari</b>	Oddiy, tez, oʻqitish talab qilmaydi	Semantik oʻxshashlikni yaxshi aniqlaydi	Eng kuchli semantik tushunish
<b>Kamchiliklari</b>	Sinonimlar va kontekstni tushunmaydi	Koʻp maʼnoli soʻzlar uchun bitta vektor beradi	Juda katta model va sekin ishlaydi
<b>Tipik qoʻllanishlar</b>	Qidiruv tizimlari, kalit soʻz ajratish	Semantik oʻxshashlik, tavsiya tizimlari	NLP vazifalari: sentiment, NER, savol-javob
<b>Misol</b>	“apple” soʻzi hujjatda koʻp uchrasa muhim deb baholanadi	king – man + woman ≈ queen	“bank” soʻzi kontekstga qarab turli maʼno oladi
<b>Qachon ishlatiladi</b>	Tez, oddiy kalit soʻzli qidiruv	Soʻzlar oʻrtasidagi semantik oʻxshashlik kerak, lekin resurs cheklangan	Eng aniq semantik tushunish, murakkab matn tahlili

**Metodika.** Foydalanuvchilar tomonidan ushbu tizimlarni oʻrganishda ayrim jihatlarni hisobga olish talab etiladi. Biz bu yerda ushbu modellarni oʻrganishda “Ekspert baholash metodi” dan foydalanishni taklif qilamiz. [2]

Ekspert baholash metodi quyidagicha amalga oshiriladi: dastlab muammoga taʼsir qiluvchi asosiy omillar yoki muammoning yechimi boʻlgan bir nechta variantlari tanlab olinadi. Mazkur muammoning yechimlari yoki unga taʼsir qilayotgan asosiy faktorlarni ekspertlarning fikrini jamlash orqali aniqlashda ekspertlarning fikrini olish uchun erkin suhbat yoki savol-javob koʻrinishidagi intervyu hamda anketalashtirishdan foydalaniladi, ushbu jarayonda har bir ekspert taqqoslanayotgan faktorlarga yoki muqobil variantlarga miqdoriy baho beradi yaʼni ularni

tabaqalaydi. Soʻngra ekspert guruhlarini qatnashchilarning baholari jamlanadi [2].

Ushbu metod yordamida 1- va 2-jadvallar ergonomikasini koʻrib quyidagi xulosaga kelamiz: talaba javoblarini tekshirish uchun semantik qidiruv eng samarali hisoblanadi.

Taklif etilayotgan gʻoya:

1. PDF dan matn olinadi
2. Talaba javobi olinadi
3. Har ikkisi tozalanadi
4. Vektorga oʻtkaziladi (TF-IDF yoki Word2Vec)
5. Oʻxshashlik hisoblanadi
6. Natija chiqariladi

Quyida ushbu algoritmnining ishlashini aniq bir misol sifatida koʻramiz:

Savol: Mulohaza deb nimaga aytiladi?

1. PDF matnda ushbu javob turibdi: A – Mulohaza deb, faqat chin yoki yolgʻon qiymat qabul qiluvchi darak gaplarga aytiladi.

2. Talaba javobi: Masalan, B – 1 yoki 0 qiymat qabul qiladigan darak gaplar mulohaza deyiladi (yoki, yana boshqacha javob boʻlishi mumkin)

3. Har ikkisi tozalanadi. Tozalash – bu matndan “keraksiz elementlarni olib tashlab”, uni algoritm uchun tayyorlash.

Tozalash quyidagicha amalga oshiriladi:

1. Kichik harfga oʻtkazish:
  - A1 – mulohaza deb, faqat chin yoki yolgʻon qiymat qabul qiluvchi darak gaplarga aytiladi;
  - B1 – 1 yoki 0 qiymat qabul qiladigan darak gaplar mulohaza deyiladi;
2. Belgilarni olib tashlash:
  - A2 – mulohaza deb faqat chin yoki yolgʻon qiymat qabul qiluvchi darak gaplarga aytiladi;
  - B2 – 1 yoki 0 qiymat qabul qiladigan darak gaplar mulohaza deyiladi;
3. Tokenizatsiya - soʻzlarga ajratish:
  - A3 – “mulohaza”, “deb”, “aytiladi”, “faqat”, “chin”, “yoki”, “yolgʻon”, “qiymat”, “qabul”, “qiluvchi”, “darak”, “gaplarga”;
  - B3 – “1”, “yoki”, “0”, “qiymat”, “qabul”, “qiladigan”, “darak”, “gaplar”, “mulohaza”, “deyiladi”;
4. Stop-soʻzlarni olib tashlash – kelishik qoʻshimchalari, “-lar”, “deb”, “edi” va hokazo kabi soʻzlar, qoʻshimchalarni olib tashlash:
  - A4 – “mulohaza”, “aytiladi”, “faqat”, “chin”, “yolgʻon”, “qiymat”, “qabul”, “qiluvchi”, “darak”, “gap”;
  - B4 – “1”, “0”, “qiymat”, “qabul”, “qiladigan”, “darak”, “gap”, “mulohaza”, “deyiladi”;
5. Stemming/normalizatsiya-soʻzlarni ildizga keltirish:
  - A5 – “mulohaza”, “aytil”, “faqat”, “chin”, “yolgʻon”, “qiymat”, “qabul”, “qil”, “darak”, “gap”;
  - B5 – “1”, “0”, “qiymat”, “qabul”, “qil”, “darak”, “gap”, “mulohaza”, “deyil”;
6. Boʻsh joylarni olib tashlash:
  - A6 – “mulohaza”, “aytil”, “faqat”, “chin”, “yolgʻon”, “qiymat”, “qabul”, “qil”, “darak”, “gap”;
  - B6 – “1”, “0”, “qiymat”, “qabul”, “qil”, “darak”, “gap”, “mulohaza”, “deyil”.

**4. Vektorga oʻtkaziladi (TF-IDF yoki Word2Vec)** – A va B matnlar lugʻatda umumlashtiriladi. A va B matnlar lugʻatdagi har bir soʻz qiymati 0/1 ga teng boʻlgan lugʻat-vektorga ajratiladi. Vektor har bir qiymati matndagi mavjud soʻzlarga qarab, mavjud boʻlsa – 1, mavjud boʻlmasa – 0, qiymat qabul qiladi.

Lugʻat quyidagi koʻrinishda boʻladi:

$D = \{\text{"mulohaza"}, \text{"aytil"}, \text{"faqat"}, \text{"chin"}, \text{"yolgon"}, \text{"qiymat"}, \text{"qabul"}, \text{"qil"}, \text{"darak"}, \text{"gap"}, \text{"1"}, \text{"0"}, \text{"deyil"}\}$

Har bir matn vektorga aylantiriladi:

$V_A = [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0]$

$V_B = [1, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1]$

**5. Oʻxshashlik hisoblanadi** – Toʻgʻridan-toʻgʻri kalit soʻzlar (indeks) oʻxshashligini tekshirilganda:

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

$$A \cdot B = 6$$

$$|A| = 3,3$$

$$|B| = 3$$

$$\cos \theta = \frac{A \cdot B}{|A| \cdot |B|} = \frac{6}{3,3 \cdot 3} \approx 0,6$$

Oʻxshashlik 0,6 (60%) ga teng boʻlmoqda, va bu natijalarni tekshirishda ayrim xatoliklarga olib keladi. Bu yerda asosiy muammo vektorning qiymatlari sifatida qabul qilinayotgan soʻzlar – sonlar maʼno jihatidan bir xil boʻlsada, lekin qiymat jihatidan har xil hisoblanayapti. Yaʼni: “chin”=“1”, “yolgon”=“0”, “deyil”=“aytil” soʻzlar hisobga olinmayapti.

Endi vektorlarni boshqatdan yaratamiz, yaʼni yuqoridagi maʼno jihatdan bir xil soʻzlarni inobatga olgan holda ushbu  $V_A$ ,  $V_B$  vektorlarni yaratamiz:

Yangi lugʻat quyidagi koʻrinishda boʻladi:

$D = \{\text{"mulohaza"}, \text{"aytil"}, \text{"faqat"}, \text{"chin"}, \text{"yolgon"}, \text{"qiymat"}, \text{"qabul"}, \text{"qil"}, \text{"darak"}, \text{"gap"}\}$

$V_A = [1, 1, 1, 1, 1, 1, 1, 1, 1, 1]$

$V_B = [1, 1, 0, 1, 1, 1, 1, 1, 1, 1]$

Soʻzlarning semantik tahlil asosida sinonimlarining bir xilligi hisobga olingan holda oʻxshashligini tekshirilganda:

$$A \cdot B = 9$$

$$|A| \approx 3,2$$

$$|B| = 3$$

$$\cos \theta = \frac{A \cdot B}{|A| \cdot |B|} = \frac{9}{3,2 \cdot 3} \approx 0,95$$

Eʼtibor qaratsak, bu yerda, oʻxshashlik darajasi 0,89 (89%) ni tashkil qilmoqda. Bu bilan B javobning A javobga juda yaqinligini aniqlaymiz.

**6. Natija chiqariladi:** Natija quyidagi baholash mezonlari asosida qoʻyiladi.

Baholash mezonlari:

Oʻxshashlik	Baho	Tavsif
$\geq 0.80$	5	Toʻliq toʻgʻri
0.65-0.79	4	Asosiy maʼno toʻgʻri
0.50-0.64	3	Qisman toʻgʻri
0.30-0.49	2	Notoʻgʻri
$< 0.30$	1	Mavzudan chetga chiqqan

Bu yerda dasturda ushbu Word2Vec modeli yaratilib, matnlar asosida oʻqitiladi:

model = Word2Vec(sentences, vector\_size=50, window=3, min\_count=1, seed=42)

bu yerda, sentences – oʻquv maʼlumotlari:

sentences = [  
tokens\_A,

```
tokens_B,
["mulohaza", "fikr"],
["chin", "togri"],
["yolgon", "notogri"],
["aytil", "deyil"],
["qiluvchi", "qiluvchi"],
["darak", "gap"]
]
```

bu:

- model nimani o'rganishini belgilaydi;
- qancha katta bo'lsa, model shuncha aniq bo'ladi.

vector\_size=50 – vektor uzunligi

har bir so'z nechta sondan iborat bo'ladi, ya'ni, so'zlarni sonlarga aylantirib beradi.

Masalan: “mulohaza” = [0.12, -0.34, 0.22, ..., 0.10] (50 ta son).

Vektor uzunligini tanlashda ushbu jadval asosida tanlansa maqsadga muvofiq hisoblanadi:

Vektor qiymatlari	uzunligi	Natija
10-50		tez, oddiy
100-300		yaxshi aniqlik
300+		kuchli, lekin sekin

Bizni eksperimentimizda vector\_size=50 o'rtacha hisoblanadi.

window=3 – so'zlarning atrofini o'rganish parametri. Model bir so'z atrofidagi nechta so'zni ko'rishini anglatadi. Masalan, “mulohaza” “chin” “yolgon” qiymat” “qabul” berilgan bo'lsa, va agar window=3 bo'lsa: “yolgon” uchun: [“mulohaza”, “chin”, qiymat”, “qabul”] bo'ladi. Quyida window uchun kattaliklar berilgan:

window	Natija
kichik (2-3)	lokal ma'no
katta (5-10)	umumiy ma'no

Bizni eksperiment uchun window=3, bu holat tezkor ishlash uchun va local ma'noni anglash uchun yetarli hisoblanadi.

**min\_count=1 – minimal chastota**

so'z necha marta chiqsa model uni o'rganadi

Qiymat	Natija
1	hamma so'zlar olinadi
2+	kam uchraydiganlar tashlanadi

Ushbu eksperimentda min\_count =1 – barcha so'zlar ishlatiladi

seed=42 – tasodifiylikni boshqarish. model har safar bir xil natija berishi uchun ishlatiladi. Aks holda, dastur har safar ishga tushirilganda vektor qiymatlar o'zgaradi va o'qitilish murakkablashib, vaqt ko'p sarflaydi. Bu yerda – bizning misolimizda 42-shunchaki standart son.

Biz yaratilgan Word2Vec modeli: chin ≈ 1, yolgon ≈ 0, mulohaza ≈ fikr va hokazo kabi o'rganadi. Har bir so'z uchun alohida vektor beriladi. Bu yerda, Word2Vec so'zlarni kontekst orqali o'rganadi. Ya'ni, “chin” va “1” birga chaqirilsa o'zaro yaqinlashadi.

Xulosa. Xulosa qilib aytadigan bo'lsak, taklif etilgan algoritim PDF kitoblardan talaba javoblarini semantik tahlil orqali yuqori aniqlikda tekshiradi. Tizim o'zbek tillarida muvaffaqiyatli ishlaydi.

Ushbu eksperimentdagi Word2Vec(...) modeli quyidagini ta'minlaydi:

- sinonimlarni aniqlash;
- semantik yaqinlik;
- matnni tushunish.

Eng muhimi model sifati bu, sentences sifati. Agar, dataset kichik bo'lsa natija yomon yoki o'rtacha, va agar katta dataset bilan ishlanganda natija yuqori bo'ladi.

Word2Vec modelida vektor o'lchami, kontekst oynasi va minimal chastota kabi parametrlar modelning semantik aniqligiga sezilarli ta'sir ko'rsatadi. Ushbu parametrlarning optimal tanlanishi matnlar orasidagi o'xshashlikni aniqlash sifatini oshiradi.

#### **FOYDALANILGAN ADABIYOTLAR RO'YXATI**

1. Abdurahmonov S. Matnlarni qayta ishlash va tabiiy til texnologiyalari. – Toshkent, 2020.
2. Axmedov J. R. Tizimli tahlil: darslik. – Jizzax: Ilm nuri print MChJ, 2023. – 266 bet.
3. Devlin J. et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding // NAACL-HLT. – 2019.
4. Goldberg Y., Levy O. word2vec Explained: Deriving Mikolov et al.'s Negative-Sampling Word-Embedding Method [Electronic resource]. – 2014. – Mode of access: arXiv.
5. Jones K. S. A statistical interpretation of term specificity and its application in retrieval // Journal of Documentation. – 1972. – Vol. 28, № 1. – P. 11–21.
6. Karaboga D. et al. Automatic Short Answer Grading Using Deep Learning // Computers & Education. – 2020.
7. Karimov A., Tursunov B. Sun'iy intellekt va uning qo'llanilishi. – Toshkent, 2021.
8. Le Q., Mikolov T. Distributed Representations of Sentences and Documents // International Conference on Machine Learning (ICML). – 2014.
9. Mikolov T. et al. Efficient Estimation of Word Representations in Vector Space [Electronic resource]. – 2013. – Mode of access: arXiv:1301.3781.