

LSTM ASOSIDA MATNLARDAN BILIMLAR BAZASINI YARATISH

Muhamediyeva Dilnoz Tulkunovna

Raqamli texnologiyalar va sun'iy intellekt kafedrasi professori, Toshkent irrigatsiya va qishloq xo'jaligini mexanizatsiyalash muhandislari instituti Milliy tadqiqot universiteti

Ungalov Sanjar Sayfullo o'g'li

Raqamli texnologiyalar va sun'iy intellekt kafedrasi izlanuvchisi, Toshkent irrigatsiya va qishloq xo'jaligini mexanizatsiyalash muhandislari instituti Milliy tadqiqot universiteti

Turg'unova Nafisaxon Maxammadjon qizi

Raqamli texnologiyalar va sun'iy intellekt kafedrasi assistenti, Toshkent irrigatsiya va qishloq xo'jaligini mexanizatsiyalash muhandislari instituti Milliy tadqiqot universiteti

Anotatsiya: Ushbu maqolada matnlardan muhim obyektlarni ajratib olish va bilimlar bazasini yaratish uchun chuqur o'rghanishga asoslangan model taklif etiladi. Nomlangan obyektlarni aniqlash (Named Entity Recognition, NER) vazifasi uchun uzoq qisqa muddatli xotira (Long Short-Term Memory, LSTM) modeli qo'llanilgan. Ma'lumotlar oldindan qayta ishlanib, tokenizatsiya va one-hot kodlash usullari orqali raqamli shaklga o'tkaziladi. Model turli obyekt turlarini (shaxs nomlari, sanalar, joy nomlari) ajratib olish uchun o'qitiladi va baholanadi. Eksperimental natijalar modelning samaradorligini ko'rsatadi va turli parametrlarning ta'siri tahlil qilinadi.

Kalit so'zlar: matnlarni qayta ishlash, nomlangan obyektlarni aniqlash, LSTM, mashinaviy o'rghanish, bilimlar bazasi, tokenizatsiya, one-hot kodlash.

СОЗДАНИЕ БАЗЫ ЗНАНИЙ ИЗ ТЕКСТОВ НА ОСНОВЕ LSTM

Аннотация: В данной статье предлагается модель, основанная на глубоком обучении, для извлечения ключевых объектов из текстов и создания базы знаний. Для задачи распознавания именованных сущностей (Named Entity Recognition, NER) используется модель долгой краткосрочной памяти (Long Short-Term Memory, LSTM). Данные предварительно обрабатываются, преобразуются в цифровую форму с помощью токенизации и one-hot кодирования. Модель обучается и оценивается для выделения различных типов объектов (имена людей, даты, географические названия). Экспериментальные результаты демонстрируют эффективность модели, а также анализируется влияние различных параметров.

Ключевые слова: обработка текста, распознавание именованных сущностей, LSTM, машинное обучение, база знаний, токенизация, one-hot кодирование.

CREATING A KNOWLEDGE BASE FROM TEXTS BASED ON LSTM

Annotation: This article proposes a deep learning-based model for extracting key entities from texts and creating a knowledge base. The Long Short-Term Memory (LSTM) model is used for the Named Entity Recognition (NER) task. The data is preprocessed and converted into a digital format using tokenization and one-hot encoding. The model is trained and evaluated to extract various types of entities (e.g., person names, dates, and location names). Experimental results demonstrate the model's effectiveness, and the impact of different parameters is analyzed.

Keywords: Text processing, named entity recognition, LSTM, machine learning, knowledge base, tokenization, one-hot encoding.

Kirish. Hozirgi kunda matnlarni tahlil qilish va ulardan bilimlar bazasini yaratish ko‘plab sohalarda muhim ahamiyat kasb etmoqda. Jumladan, tabiiy tilni qayta ishlash (NLP) usullari yordamida matnlardan muhim ma’lumotlarni ajratib olish imkoniyati mavjud [1]. Nomlangan obyektlarni aniqlash (NER) – bu matndagi shaxs ismlari, sanalar, joy nomlari kabi muhim elementlarni aniqlashga qaratilgan muhim yo‘nalishdir. Ushbu tadqiqotda LSTM modeli asosida matnlardan muhim obyektlarni ajratib olish va ularni bilimlar bazasiga joylashtirish usuli taklif etiladi [2,3].

LSTM – chuqur o‘rganish usullaridan biri bo‘lib, u ketma-ketliklarni o‘rganish va uzoq muddatli bog‘liqliklarni saqlashda samarali ishlaydi. Ushbu model turli sohalarda, jumladan, matnni qayta ishlash va tabiiy til tahlilida muvaffaqiyatlari qo‘llanilmoqda. Tadqiqotimizda LSTM yordamida matnlardan nomlangan obyektlarni ajratib olish jarayoni amalga oshiriladi.

2.Usullar va yondashuvlar. Matnlardan nomlangan obyektlarni ajratib olish uchun bir necha bosqich bajariladi [4]. Dastlab, matnlarni raqamli shaklga o‘tkazish uchun tokenizatsiya jarayoni amalga oshiriladi. Bu jarayonda matn alohida so‘zlarga ajratilib, har bir so‘zga indeks tayinlanadi. Keyinchalik, matnlarning uzunligi bir xil darajaga keltirish uchun padding usuli qo‘llaniladi. Kichikroq uzunlikdagi matnlarga maxsus to‘ldirish belgilarini qo‘sish orqali barcha matnlarning uzunligi tenglashtiriladi. Shundan so‘ng, obyektlarni modelga moslashtirish maqsadida one-hot kodlash amalga oshiriladi. Bu bosqichda obyektlarga mos raqamli ifodalar tayinlanadi va ular modelga kiritish uchun tayyorlanadi.

Nomlangan obyektlarni aniqlash jarayonida har bir so‘zga mos obyekt turi belgilanadi. Bu jarayonda model matndagi shaxs nomlari, sanalar va joy nomlarini ajratib oladi. Ajratilgan obyektlar turli kategoriyalar bo‘yicha tasniflanadi va natijalar model tomonidan qayta ishlanadi. Bu bosqichda shaxs nomlari ("Alice", "Bob", "Charlie"), sanalar ("January 5th, 2023", "March 15th, 2023") va joy nomlari ("New York") kabi obyektlarni aniqlash maqsad qilib qo‘yiladi.

Nomlangan obyektlarni aniqlash uchun chuqur o‘rganish usuli – LSTM modelidan foydalilanadi. Model bir necha asosiy qatlamlardan tashkil topgan bo‘lib, birinchi qatlam so‘zlarni raqamli vektorlarga o‘tkazuvchi embedding qatlami hisoblanadi. Ushbu qavatdan keyin LSTM qatlami kelib, u vaqt bo‘ylab so‘zlarning o‘zaro bog‘liqligini saqlaydi va tahlil qiladi. Modelning haddan tashqari moslashuvini oldini olish uchun dropout qatlami qo‘shiladi. Ushbu qatlam neyronlarning bir qismini vaqtincha o‘chirib, modelni umumlashuvchan qilishga yordam beradi. Modelning yakuniy qatlami esa softmax funksiyasiga asoslangan bo‘lib, har bir so‘zning tegishli kategoriya ehtimolini hisoblab chiqadi.

Modelni o‘qitish jarayonida categorical_crossentropy yo‘qotish funksiyasi qo‘llaniladi. Bu funksiya modelning nomlangan obyektlarni to‘g‘ri tasniflash darajasini oshirishga xizmat qiladi. Model vaznlari Adam optimizer yordamida yangilanadi va model 10 epoch davomida o‘qitiladi. Har bir o‘qitish qadamida mini-batch usuli qo‘llanilib, har bir qadamda 2 ta namunadan iborat ma’lumotlar ishlatiladi.

Model samaradorligini baholash uchun aniqlik (accuracy), F1-score va xatolik matritsasi hisoblab chiqiladi. Aniqlik modelning umumiyligi to‘g‘ri tasnifangan obyektlar ulushini o‘lchashga yordam beradi. F1-score esa aniqlik va to‘g‘rilik muvozanatini hisoblash imkonini beradi. Xatolik matritsasi orqali modelning noto‘g‘ri va to‘g‘ri tasniflash darajasi tahlil qilinadi. Bu natijalar asosida modelning samaradorligi aniqlanadi va uning ishlash sifati baholanadi.

3. Natijalar. Modelni o‘qitish va sinov jarayonlaridan so‘ng, uning samaradorligi bir nechta mezonlar asosida baholandi. Aniqlik (Accuracy), F1-score va xatolik matritsasi tahlil qilindi. Modelni o‘qitish davomida aniqlik oshib bordi va yo‘qotish funksiyasi qiymati kamaydi. 10 epoch davomida o‘qitish natijalariga ko‘ra, trening aniqligi 92.4% ni, test aniqligi esa 89.7% ni tashkil etdi. Shu bilan birga, trening yo‘qotish funksiyasi qiymati 0.19, test yo‘qotish funksiyasi esa 0.23

ga teng bo'ldi. Bu natijalar shuni ko'rsatadiki, model trening to'plami bo'yicha yaxshi o'rgangan, test to'plamida esa biroz pastroq natija qayd etgan. Biroq, farq katta emas va bu modelning ortiqcha moslashib qolmaganligini (overfitting yo'qligini) bildiradi.

Model matndagi shaxs nomlari, sanalar va joy nomlarini ajratib olishda quyidagi natijalarga erishdi: shaxs nomlarini aniqlash aniqligi 94.1%, sanalarni aniqlash aniqligi 90.3%, joy nomlarini aniqlash aniqligi esa 87.5% ni tashkil etdi. Shaxs nomlarini ajratib olishdagi aniqlik eng yuqori bo'ldi. Bu shuni anglatadiki, model shaxs nomlarini yaxshi farqlaydi. Sanalar va joy nomlarining aniqlik darajasi biroz pastroq bo'lsa ham, qoniqarli natija qayd etildi.

Xatolik matritsasiga ko'ra, model shaxs nomlarini deyarli mukammal aniqlagan, biroq sanalarni va joy nomlarini noto'g'ri tasniflash holatlari mavjud. Ayniqsa, ba'zi sanalar joy nomlari sifatida noto'g'ri belgilangan. Masalan, "March 15th, 2023" iborasi ayrim holatlarda joy nomi sifatida tasniflangan. Bu xatoliklar modelning o'quv ma'lumotlari hajmi va ba'zi matnlardagi murakkab ifodalarga bog'liq bo'lishi mumkin. Ushbu muammoni bartaraf etish uchun o'quv to'plamini kengaytirish yoki sanalarni yanada aniq aniqlash uchun maxsus qoidalarni (rule-based features) qo'shish tavsiya etiladi.

Modelning umumiylashtirish sifatini baholash uchun F1-score quyidagicha bo'ldi: shaxs nomlari uchun 93.8%, sanalar uchun 88.7%, joy nomlari uchun esa 85.9%. Bu natijalarga asoslanib, model shaxs nomlarini mukammal tasniflashga yaqin natijaga ega, sanalar va joy nomlari bo'yicha esa yana takomillashtirish talab qilinishi mumkin.

Quyidagi jadval modelning shaxs nomlari, sanalar va joy nomlarini ajratib olishdagi samaradorligini ifodalaydi:

1-jadval

Model natijalari

Kategoriylar	Aniqlik (%)	F1-score (%)	Xatolik (%)
Shaxs nomlari	94.1	93.8	5.9
Sanalar	90.3	88.7	9.7
Joy nomlari	87.5	85.9	12.5

Jadvaldan ko'rinish turibdiki, shaxs nomlarini aniqlashning aniqlik va F1-score ko'rsatkichlari eng yuqori bo'ldi. Sanalar va joy nomlarini aniqlash natijalari nisbatan pastroq bo'lib, bu ba'zi xatoliklar bilan bog'liq.

4. Xulosa. Ushbu tadqiqotda LSTM asosida matnlardan shaxs nomlari, sanalar va joy nomlarini ajratib olish modeli yaratildi va uning samaradorligi tahlil qilindi. Tahlillar shuni ko'rsatadiki, model o'quv jarayonida matnning kontekstual ma'nosini yaxshi o'zlashtirgan bo'lsada, ba'zi sanalar va joy nomlarini ajratishda xatoliklarga yo'l qo'ygan. Bu xatoliklar asosan ma'lumotlar to'plamidagi o'ziga xosliklar yoki modelning ba'zi murakkab kontekstlarni tushunishdagi chekllovleri bilan bog'liq. Ushbu ish natijalari matnlarni avtomatik qayta ishlash sohasida amaliy ahamiyatga ega bo'lib, turli sohalarda, jumladan, hujjatlar tahlili, huquqiy hujjatlar va ilmiy maqolalar bo'yicha ma'lumotlarni saralash uchun foydali bo'lishi mumkin.

Adabiyotlar ro'yxati

- Mamatov, N. S., Niyozmatova, N. A., Samijonov, A. N., & Samijonov, B. N. (2022, September). Construction of language models for Uzbek language. In *2022 International Conference on Information Science and Communications Technologies (ICISCT)* (pp. 1-4). IEEE.
- Duppatti, S. K. ., & Babu, A. (2023). Named Entity Recognition for English Language Using Deep Learning Based Bi Directional LSTM-RNN. *International Journal on Recent and Innovation Trends in Computing and Communication*, 11(5), 330–337. <https://doi.org/10.17762/ijritcc.v11i5.6621>

3. A. Mohamed and N. Jaitly, 2013, “Hybrid speech recognition with deep bidirectional LSTM,” IEEE, 2013, pp. 273–278.
4. Ниёзматова, Н. А., Маматов, Н. С., Отахонова, Б. И., Бобоев, Л. Б., & Самижонов, А. Н. Матнларни таснифлашда информатив белгилар мажмуасини аниқлаш усуллари.

МАЪЛУМОТЛАРНИ ИНТЕЛЛЕКТУАЛ ТАҲЛИЛЛАШ УЧУН БАШОРАТЛАШ ЁНДАШУВЛАРИ

Бабомурадов Озод Жураевич

т.ф.д.(DS), проф. Жиззах шаҳридаги ҚФУ филиали

bobomuradov@gmail.com

Хайдаров Озоджон Асламкулович

Малака ошириш йўналиши раҳбари

“Creative Associates International” Ўзбекистондаги ваколатхонаси

haydarov@gmail.com

Аннотация: Ушбу мақола бошқарув тизимларида сифатли қарорлар қабул қилиш учун маълумотлар асосида башорат моделларини ишлаб чиқиши заруратига бағишиланган. Маълумотларни таҳлил қилиш ва башоратлаш жараёнида вақтли қаторлар асосидаги усуллар кўрсатилган. Вақтли қаторлар анализи ва бошқарув стратегиясини шакллантиришдаги аҳамияти, аниқлик ва самарадорликка таъсири ҳамда натижаларнинг ишончлилигини ошириш учун зарур бўлган математик модель ва алгоритмлар таҳлили келтирилган.

Калит сўзлар: NNS, MASE, RF, TREE Баглар, SMAPE нейрон тармоғи, XG boost, PCA, FFORMS – белгилар, SDAE, SAE, MLP, ARIMA, PropNet, BSE, BC, Ada Boost, NYSE, NASDAQ, Extra Trees, GSO, SYM, KNN, LSTM, GRU, NSE, FinBERT, ELM, GB, KNN, DT, SARIMA, тасодифий ўрмон, SVM, чизиқли регрессия.

ПОДХОДЫ К ПРОГНОЗИРОВАНИЮ ДЛЯ ИНТЕЛЛЕКТУАЛЬНОГО АНАЛИЗА ДАННЫХ

Аннотация: Данная статья посвящена необходимости разработки моделей прогнозирования на основе данных для принятия качественных решений в управленических системах. Рассмотрены методы, основанные на временных рядах, в процессе анализа данных и прогнозирования. Приводятся анализ временных рядов, влияние на формирование управленической стратегии, а также математические модели и алгоритмы, необходимые для повышения точности и эффективности, а также улучшения надежности результатов.

Ключевые слова: NNS, MASE, RF (Случайный лес), TREE (Деревья), SMAPE, Нейронная сеть, XGBoost, PCA (Анализ главных компонент), FFORMS – признаки, SDAE, SAE, MLP (Многослойный перцептрон), ARIMA, Prophet, BSE, BC, AdaBoost, NYSE, NASDAQ, Extra Trees, GSO, SYM, KNN (Метод k-ближайших соседей), LSTM, GRU, NSE, FinBERT, ELM (Экстремальная машина обучения), GB (Градиентный бустинг), DT (Дерево решений), SARIMA, Случайный лес, SVM (Метод опорных векторов), Линейная регрессия.

PREDICTIVE APPROACHES FOR INTELLIGENT DATA ANALYSIS

Abstract: This article is dedicated to the necessity of developing predictive models based on data for making quality decisions in management systems. Methods based on time series in the