

**МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ
ДОНЕЦЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ ім. ВАСИЛЯ СТУСА
УКРАЇНСЬКИЙ МОВНО-ІНФОРМАЦІЙНИЙ ФОНД НАН УКРАЇНИ
ЦЕНТР ТЕОРЕТИЧНОЇ, ПРИКЛАДНОЇ ТА КОМП'ЮТЕРНОЇ ЛІНГВІСТИКИ
ПЕДАГОГІЧНА КОРПОРАЦІЯ
«УНІВЕРСИТЕТ – ШКОЛА – УНІВЕРСИТЕТ»
КАФЕДРА ЗАГАЛЬНОГО ТА ПРИКЛАДНОГО МОВОЗНАВСТВА І
СЛОВ'ЯНСЬКОЇ ФІЛОЛОГІЇ**

ЛІНГВОКОМП'ЮТЕРНІ ДОСЛІДЖЕННЯ

Збірник наукових праць

Випуск 11

Вінниця – 2018

УДК 004:81'32+33

ББК Ш12=411.4_я43+Ш11_я43

Рекомендовано до друку Вченою радою Донецького національного університету імені Василя Стуса (протокол №8 від 23 лютого 2018р.)

Лінгвоком'ютерні дослідження : зб. наук. праць / Донецький національний університет ім. Василя Стуса / Укл.: А. Загітко (наук. і відп. ред.). – Вінниця : ДонНУ, 2018. – Вип. 11. – 161с.

ISBN 978-966-639-457-9

ISSN 2307-0544

Свідоцтво про держреєстрацію КВ 22548-12448ПР

Редакційна рада: **Загітко Анатолій**, д.філол.н., професор, член-кореспондент НАН України (науковий редактор) (Україна); **Вихованець Іван**, д.філол.н., професор, член-кореспондент НАН України (Україна); **Городенська Катерина**, д.філол.н., професор (Україна); **Клименко Ніна**, доктор філологічних наук, професор (Україна); **Соколова Світлана**, д.філол.н., старший науковий співробітник (Україна); **Всеволодова Майя**, д.філол.н., професор (Росія); **Конюшкевич Марія**, д.філол.н., професор (Білорусь); **Рагавцов Василь**, д.філол.н., професор (Білорусь); **Сарновський Міхал**, д.філол.н., проф. (Польща); **Пісарек Лариса**, д.філол.н., професор (Польща); **Андерш Йозеф**, д.філол.н., професор (Чехія); **Ляхур Чеслав**, д.філол.н., професор (Польща); **Бранднер Алеш**, д.філол.н., професор (Чехія); **Жажа Станіслав**, д.філол.н., професор (Чехія); **Попович Людмила** д.філол.н., професор (Сербія).

Редакційна колегія: **Демська Орія**, д.філол.н., професор (Україна); **Космеда Тетяна**, д.філол.н. професор (Польща); **Кочан Грина**, д.філол.н., професор (Україна); **Селіванова Олена**, д.філол.н., професор (Україна); **Скаб Марія**, д.філол.н., професор (Україна); **Шитик Людмила**, д.філол.н., доцент (Україна); **Кравченко Елла**, д.філол.н., доцент (Україна); **Краснобаєва-Чорна Жанна**, д.філол.н., доцент (Україна); **Данилюк Ілля**, к.філол.н., доцент (Україна); **Ситар Ганна**, к.філол.н., доцент (Україна).

Відповідальний редактор:

Загітко Анатолій

д.філол.н., професор (Донецький національний університет ім. Василя Стуса)

Рецензенти:

Левченко Олена

д.філол.н., професор (Національний університет «Львівська політехніка»)

Кондратенко Наталія

д.філол.н., професор (Одеський національний університет імені І.І. Мечникова)

Розглянуто фундаментальні проблеми теоретичної та прикладної лінгвістики, простежено актуальні напрями автоматичного опрацювання мовної інформації, традиційного та машинного перекладу, схарактеризовано особливості створення та використання лінгвістичних баз даних, з'ясовано особливості й тенденції сучасної лексикографії та дистанційної освіти.

Для мовознавців, лінгвістів-програмістів, викладачів вищої школи, інформатиків, аспірантів, студентів, учителів.

Електронна версія збірника доступна на Донецькому лінгвістичному порталі за адресою: mova.dn.ua, на порталі ДонНУ ім. Василя Стуса за адресою: donnu.edu.ua, а також на сайті спеціальності «Прикладна лінгвістика» ДонНУ ім. Василя Стуса за адресою: pldonnu.pp.ua.

Адреса видавця і редакційної ради: 21007, м. Вінниця, вул. академіка Янгеля, 4, ауд. 320.

© Автори статей, 2018

©Донецький національний університет, 2018

ЗМІСТ

ВСТУП	5
РОЗДІЛ I. ТЕОРЕТИКО-ПРИКЛАДНІ ПРОБЛЕМИ ФУНКЦІЙНОЇ СЕМАНТИКИ ЛЕКСИЧНИХ І ФРАЗЕОЛОГІЧНИХ ОДИНИЦЬ	9
Каріна Адамусік. СЛОВНИК ТЕРМІНІВ ЛІНГВІСТИЧНОЇ ЕКСПЕРТИЗИ: СТРУКТУРА ТА ПРИНЦИПИ УКЛАДАННЯ	9
Максим Дерев'янченко. ЕЛЕКТРОННИЙ СЛОВНИК ШАХОВОЇ ЛЕКСИКИ В УКРАЇНСЬКІЙ МОВІ: ДИРЕКТОРІЯ ВЕБ-САЙТУ ТА ПРОГРАМНИЙ КОД	15
Дарія Дубенко. ЧАСТОТА ВЖИВАННЯ СЛЕНГУ ОСОБАМИ ВІКОМ ВІД 14 ДО 21 РОКУ	51
Анна Каліта. АНГЛІЙСЬКІ ФРАЗЕОЛОГІЧНІ ОДИНИЦІ ТА СПОСОБИ ЇХ ПЕРЕКЛАДУ: на прикладі роману Дж. К. Ролінг «Гаррі Поттер і філософський камінь»	57
Жанна Краснобаєва-Чорна. АКСІОФРАЗЕМНА ПРАГМАТИКА КОЛЬОРУ	65
Наталя Матвеева. СОМАТИЧНИЙ КОД ЗДОРОВ'Я ТА КРАСИ В УКРАЇНСЬКІЙ І АНГЛІЙСЬКІЙ ФРАЗЕОЛОГІЇ	71
Лілія Федорюк. ЛЕКСИКО-СЕМАНТИЧНІ ЗАСОБИ РЕАЛІЗАЦІЇ ЦІННІСНОГО СКЛАДНИКА КОНЦЕПТУ СМЕРТЬ	80
Інна Шматко. ОСОБЛИВОСТІ ТВОРЕННЯ ТЕРМІНІВ- СЛОВОСПОЛУЧЕНЬ УКРАЇНСЬКОЇ БДЖІЛЬНИЦЬКОЇ ТЕРМІНОЛОГІЇ	89

РОЗДІЛ II. ТЕОРЕТИКО-ПРИКЛАДНІ ПРОБЛЕМИ МОРФОЛОГІЇ ТА СИНТАКСИСУ	95
Олена Поліщук. ОДНОСКЛАДНЕ РЕЧЕННЯ В ІДЮСТИЛІ ЄВГЕНА ГУЦАЛА (на матеріалі оповідання «Запах кропу»)	95
Ольга Хруленко. СТРУКТУРНІ ТИПИ ПОХІДНИХ СУЧАСНИХ КОМП'ЮТЕРНИХ ТЕРМІНІВ (на матеріалі сучасної української мови)	103
РОЗДІЛ III. ТЕОРЕТИКО-ПРИКЛАДНІ ПРОБЛЕМИ ЛІНГВОПЕРСОНОЛОГІЇ	109
Ілля Данилюк. МОДЕЛЮВАННЯ ФОРМАЛЬНО- ЗМІСТОВОГО РІВНЯ ДІЯЛЬНОСТІ МОВНОЇ ОСОБИСТОСТІ	109
Анатолій Загнітко. ПРИНЦИП КІНЧНОСТІ В ПОРІВНЯЛЬНО-ІСТОРИЧНІЙ ЛІНГВОПАРАДИГМІ	121
РОЗДІЛ IV. ТЕОРЕТИКО-ПРИКЛАДНІ ПРОБЛЕМИ КОРПУСНОЇ ЛІНГВІСТИКИ, ПЕРЕКЛАДОЗНАВСТВА, ДИСКУРСОЛОГІЇ	135
Юлія Калимон. АБЗАЦ ЯК КАТЕГОРІЯ СТРУКТУРНОГО ПОДІЛУ ТЕКСТУ	135
Наталя Лотоцька. МЕТАТЕКСТОВЕ РОЗЗНАЧЕННЯ ТЕКСТІВ (на матеріалі творів Романа Іваничука)	141
Данило Островський. БАЗА ДАНИХ «СПОСОБИ ВІДТВОРЕННЯ АНГЛІЙСЬКИХ ВЛАСНИХ НАЗВ УКРАЇНСЬКОЮ ТА РОСІЙСЬКОЮ МОВАМИ» (на матеріалі роману Дж. К. Ролінг «Гаррі Поттер і філософський камінь»)	149
Єлизавета Тімченко. ПІДХОДИ ДО МАШИННОГО ПЕРЕКЛАДУ: ІСТОРИЧНИЙ АСПЕКТ	155

РОЗДІЛ ІІІ. ТЕОРЕТИКО-ПРИКЛАДНІ ПРОБЛЕМИ ЛІНГВОПЕРСОНОЛОГІЇ

Ілля Данилюк

(к. філол. н., доц., докторант)

Донецький національний університет імені Василя Стуса **МОДЕЛЮВАННЯ ФОРМАЛЬНО-ЗМІСТОВОГО РІВНЯ ДІЯЛЬНОСТІ МОВНОЇ ОСОБИСТОСТІ**

У статті описано методи комп'ютерної лінгвістики, покладені в основу проекту «Комунікативно-прагматична і дискурсивно-граматична лінгвоперсонологія: структурування мовної особистості та її комп'ютерне моделювання», який реалізує кафедра загального та прикладного мовознавства та слов'янської філології ДонНУ. Подано опис ключових методів, як-от машинне навчання, синтез мовлення, моделювання почерку.

Ключові слова: лінгвоперсонологія, мовна особистість, моделювання, методи, машинне навчання, синтез мовлення.

Лінгвоперсонологія – наука, що має предметом вивчення власне-людську мовну особистість – спирається на вже практично столітні теоретичні засади, закладені у працях В. Гумбольдта, Й. Л. Вайсгербера, І. Бодуена де Куртене, В. Вундта, О. Потебні, В. Виноградова, згодом поглиблені та чіткіше окреслені в роботах В. Карасика, Ю. Караулова, В. Нерознака та ін. (теоретичні засади лінгвоперсонології є також предметом розгляду першої публікації циклу (Danylyuk, I. "Teoretychni zasady i metody lnhvopersonolohiyi")). З іншого боку, розвиток інформаційних технологій, здобутки квантитативної та корпусної лінгвістики, можливості автоматичного опрацювання природного мовлення озброїли лінгвоперсонологію новим інструментарієм і, відповідно, накреслили нові завдання та перспективи.

Метою статті є опис та узагальнення методології фундаментального дослідження з лінгвоперсонології. Завдання, які має бути розв'язано, включають визначення підходів до моделювання 1) формально-змістового рівня діяльності мовної особистості; 2) з'ясування обсягу поняття мовної особистості; 2) моделювання формально-звукового рівня діяльності мовної особистості – за допомогою систем синтезу мовлення; 3) моделювання формально-графічного рівня діяльності мовної особистості – почерку.

Актуальність нашого дослідження, крім усього, цілком мотивована розвитком надзвичайного теоретично-наукового і практичного інтересу до можливостей опрацювання величезного обсягу мовних даних, які генерує людина у повсякденному професійному та особистому житті в електронних формах комунікації (e-mail, sms, голосовий зв'язок, аудіо- та відеоблоги, соціальні мережі тощо).

У лінгвістиці підходи до вивчення мовної особистості включають її: 1) психологічний аналіз; 2) соціологічний аналіз; 3) культурологічний аналіз – моделювання лінгвокультурних типажів – узагальнених відомих представників певних груп суспільства, поведінка яких втілює в собі норми лінгвокультури загалом і впливає на поведінку всіх представників суспільства; 4) лінгвістичний аналіз (опис комунікативної поведінки носіїв елітарної або масової мовної культури, характеристика людей з позицій їхньої комунікативної компетенції, аналіз креативної і стандартної мовного свідомості); 5) прагмалінгвістичний аналіз мовної особистості, в основі якого лежить виділення типів комунікативної тональності, характерної для того чи іншого дискурсу (Karasyk).

Саме лінгвістичний підхід, на нашу думку, набуває нової актуальності через можливість збирати велику кількість

мовленнєвих даних у вигляді корпусів текстів та звукових корпусів мовлення, а також через появу інструментарію для їх автоматичного опрацювання.

У межах проекту «Комунікативно-прагматична і дискурсивно-граматична лінгвоперсонологія: структурування мовної особистості та її комп'ютерне моделювання» однією з опорних точок ж теза про те, що мовносоціумна особистість може бути описана, змодельована і клонована у віртуальному вимірі на основі аналізу породжених реальною особою (донором) усних або писемних текстів. Відповідно, лінгвістичний підхід до вивчення такої особистості включає, на нашу думку, а) *моделювання формально-змістового рівня* діяльності мовної особистості – за допомогою використання розробок корпусної лінгвістики; б) *моделювання формально-звукового рівня* діяльності мовної особистості – за допомогою систем синтезу мовлення; в) *моделювання формально-графічного рівня* діяльності мовної особистості – почерку.

Корпусні технології давно вже стали одним із основних методів лінгвістичних досліджень. Так, ще в 1960-і роки створювався Браунівський корпус (США), який містив 1 млн слів. У 1970-і роки минулого століття стартував LOB корпус (Великобританія, Норвегія), у 1980-ті роки почали створюватися такі корпусу, як: Машинний Фонд російської мови, Упсальський корпус російської мови (Швеція) – обидва по 1 млн слів, The Bank of English, Birmingham, – 20 млн слів. У 1990-і роки було створено British National Corpus, який включав на той час 100 млн слів, а також інші національні корпуси для угорської, італійської, хорватської, чеської, японської мови обсягом по 100 млн слів. На початку XXI ст. створювалися такі корпуси, як American National Corpus і Gigaword corpora (англійська, арабська, китайська) на 1 млрд слів, Національний корпус російської мови, над яким

працюють лінгвісти Москви і Санкт-Петербурга, містить 300 млн слововживань. В Україні проблемою корпусної лінгвістики активно займаються вчені Інституту української мови НАН України, Українського мовно-інформаційного фонду, Інституту філології Київського національного університету ім. Т.Шевченка, Національного університету «Львівська політехніка» та ін.

Корпусний менеджер Manatee/Bonito й розроблені для нього корпуси текстів було докладно описано в (Danylyuk, I. "Korpus tekstiv dlya vyvchennya hramatychnoyi sluzhbovosti.", Danylyuk, I. "Korpus tekstiv dlya vyvchennya hramatychnoyi sluzhbovosti: klasyfikatsiya hramatychnykh klasiv i pidklasiv."). Корпус текстів Юрія Шевельова (Шереха) викладено на corpora.donnu.edu.ua.

Під терміном лінгвістичний, або мовний, корпус текстів сьогодні розуміють великий, представлений в електронному вигляді, уніфікований, структурований, розмічений, філологічно компетентний масив мовних даних, призначений для вирішення конкретних лінгвістичних завдань. Такими завданнями у проєкті «Комунікативно-прагматична і дискурсивно-граматична лінгвоперсонологія: структурування мовної особистості та її комп'ютерне моделювання» є мовна модель, побудована на основі текстів авторства Юрія Шевельова, а також екстракцію знань з цих текстів та дії, спрямовані на розуміння текстів для наповнення бази знань, формування відповідей на запитання і ведення діалогу з моделлю його мовної особистості.

Концептуально завдання автоматичного опрацювання текстів було розглянуто в роботах Н. Хомського, присвячених граматиці природної мови, в яких було описано ключову парадигму комп'ютерної лінгвістики – контекстно-незалежну граматику (CFG). Перші спроби автоматичної обробки текстів зводилися до розбору із застосуванням такої граматики, побудову дерева розбору і переведення його в певне логічне представлення

знань за допомогою правил і лексикону. Після цього логічне представлення можна було додати в базу знань і виконувати з ним різні операції: шукати інформацію потрібного змісту чи типу, відповідати на запитання, перевіряти твердження тощо. Однак практичне застосування цього підходу було обмежено труднощами, пов'язаними з необхідністю враховувати загальноприйняті знання про світ, усталеної моделі для чого так і не було створено.

З 1990-х років у розпорядженні вчених з'явилися методи машинного навчання і статистичної лінгвістики. У машинному навчанні ефективними були алгоритми класифікації для різних завдань, пов'язаних з обробкою текстів: визначення спаму, сортування документів за тематиками, виділення іменованих сутностей (власних назв). У комп'ютерній лінгвістиці визначення частин мови стало високоточним завдяки таким статистичним методам, як приховані ланцюга Маркова і моделі максимальної ентропії. З'явилися парсери на основі імовірнісних контекстно-незалежних граматик, а в корпорації ІВМ було реалізовано масштабний проект зі статистичного машинного перекладу. Нарешті, було закладено основи *глибинного навчання* – найефективнішого на сьогодні з усіх автоматичних методів.

Глибинне навчання – навчання багаторівневих («глибоких») нейронних мереж на великих обсягах даних, що дозволяють уникнути роботи зі ручного розмічування корпусів текстів для машинного навчання, оскільки система «вчиться» виділяти їх автоматично. До речі, перший загальний робочий алгоритм керованого навчання багаторівневої мережі перцептронів було опубліковано нашими співвітчизниками О. Івахненком і В. Лапою у (Yvakhnenko).

2010 року було запропоновано модель лексикалізованої ймовірнісної граматики, яка дозволила підвищити точність

граматичного розбору до 93%. Точність розбору – це відсоток правильно побудованих граматичних зв'язків, однак імовірність того, що довге речення буде розібрано правильно, зазвичай дуже низька. Одночасно, завдяки новим алгоритмам і підходам, включно з глибинним навчанням, збільшилася швидкість граматичного розбору. Крім того, практично всі провідні алгоритми і моделі сьогодні у відкритому доступі, зокрема й один з найефективніших алгоритм Томаса Міколова (Mikolov).

Сьогодні для моделювання формально-змістового рівня діяльності мовної особистості маємо інструменти, які можна умовно поділити на три класи: *а) методи роботи зі словами, б) методи роботи з реченнями і в) методи для обробки довільних текстів.*

1. Методи роботи зі словами

Традиційно слова обробляли як елементи множини зі словника, обсяг і повнота якого цілковито визначали ефективність такої системи. Однак побудова вичерпного словника – з усіма словоформами чи з включенням професійної і розмовної лексики, жаргону, діалектизмів тощо – надзвичайно важке завдання. На відміну від традиційного підходу, алгоритм word2vec (Mikolov) спирається на ймовірнісну модель мови – кожне слово представлено вектором з дійсних чисел у маленькому (якщо порівняти з розміром повного словника) просторі, наприклад розмірністю в 300 вимірювань. Спочатку векторам присвоюють випадкові значення. Далі в процесі навчання на укладеному корпусі для слова обчислюють вектор, максимально схожий на вектори інших слів, які трапляються у схожих контекстах. За контекст беруть невелике вікно попередніх і наступних слів, наприклад, у п'ять одиниць. У результаті виявляється, що векторно близькі слова виявляються дійсно семантично близькими. Крім того, виявляється, що багато важливих для обробки природного

мовлення відношень закодовано у вектори. Відомий приклад: якщо від вектора слова «Париж» відняти вектор слова «Франція» і додати вектор «Італія», то вийде вектор, дуже близький до вектора «Рим» – відношення «столиця» виявилось закодованим у вектори слів.

Алгоритм word2vec укладається в парадигму глибокого навчання: він сам знаходить ознаки в режимі «навчання без учителя».

Додаткові відношення (ознаки), встановлені у процесі роботи word2vec, можуть виявитися корисними для завдань обробки текстів, але не зрозуміло, які відношення дійсно містяться в векторах після навчання і наскільки надійно вони закодовані. Є методи, які дозволяють доповнювати векторні представлення слів онтологіями, котрі гарантують, що відношення будуть надійно закодовані у вектори. Наприклад, у (Xu) дослідники запропонували метод навчання векторів, в якому будь-які відношення і таксономії надійно кодуються у векторні представлення.

Втім, стандартний алгоритм word2vec не дозволяє розв'язати проблеми, пов'язані з омонімією. Модифікація алгоритму з елементами автоматичного визначення омонімів і створення окремих векторів для окремих смислів, а також процедури визначення правильного значення омонімів щодо заданого контексту запропоновано у (Bartunov).

Методи обробки природного мовлення здебільшого використовують тільки представлення, ігноруючи синтаксис і семантику, які можна вивести з синтаксичної структури речень. Така модель подання текстів називається «торбою слів» (bag of words) – простий набір слів без урахування їхнього порядку. Наприклад, використовуючи векторне представлення слів можна об'єднати в кластери вектори слів корпусу, на яких тренується

модель, і використовувати такі кластери для завдань простої класифікації, як-от – до якого стилю належить текст, визначення авторства тексту тощо. Але якщо завдання полягає в добуванні якісніших семантичних представлень або по суті знань, то потрібні інструменти обробки текстів, які працюють з синтаксичною структурою речення і не ігнорують порядок слів у ньому.

2. Методи роботи з реченнями

Перше завдання обробки мовлення на рівні речення – встановлення його синтаксичної структури. Інструменти роботи з синтаксисом помітно прогресували – лексикалізовані імовірнісні граматики значно підвищили якість синтаксичного розбору, а зручні для багатьох випадків граматики залежностей досягли якості, достатньої для розв’язання значного класу завдань обробки текстів. Крім того, за останні кілька років в сотні разів збільшилася швидкість алгоритмів синтаксичного аналізу.

Втім, точність автоматичного синтаксичного аналізу, особливо складних речень, і далі є порівняно невисокою. По-перше, багато чого залежить від якості розпізнавання частин мови, яке має бути дуже високим (97-98%), а в довгих реченнях часто трапляються неправильно визначені граматичні класи, що призводить до помилок розбору. По-друге, сам граматичний розбір дає точність щонайбільше 90-93% (відсоток правильно визначених відношень), а це означає, що в довгому реченні практично завжди будуть помилки розбору. Наприклад, за точності розбору 90% ймовірність розбору речення у 10 слів без жодної помилки складе лише 35%.

Методи глибинного навчання дозволяють інакше підійти до роботи з реченнями – моделювати речення як послідовність векторів, отриманих методом word2vec, і використати його в алгоритмах машинного навчання. Стандартні алгоритми машинного навчання працюють з фіксованим набором атрибутів, і

їх не можна пристосувати до такої моделі. Для таких випадків запропоновано скористатися рекурентними нейронними мережами, які на вході приймають одне слово у векторному представленні і мають кілька внутрішніх рівнів, а на виході будують класифікатор. На відміну від звичайних нейромереж, внутрішні рівні *рекурентної мережі* (а іноді і верхній рівень) підключені назад у мережу, тобто стан мережі, у який вона перейшла на попередньому слові, буде передано в мережу як додатковий вхід на наступному слові. Таким чином, у нейромережі з'являється аналог «пам'яті», що дозволяє їй послідовно обробляти слова в реченні та будувати припущення окремо щодо кожного слова або всього речення цілком. Інакше кажучи, мережі послідовно передають одне слово речення за іншим, а мережа використовує свої попередні стани для визначення поточного кроку. Однак на практиці прості рекурентні мережі працюють не вельми ефективно через те, що пам'ять про попередні слова в реченні швидко втрачається під час тренування та експлуатації мережі. Тому зазвичай використовують спеціальні елементи пам'яті – *LSTM (Long Term Short Memory)*, що є множиною нейронів і керівних елементів, які визначають, коли треба записувати, читати й очищати пам'ять. Ці елементи дозволяють пам'яті не змінюватися під час тривалого послідовного обчислення і правильно атрибутувати помилку під час навчання.

Рекурентні нейромережі з LSTM добре себе зарекомендували для розв'язання різних завдань – моделювання мови, машинний переклад, – але у цього класу мереж є істотний недолік, пов'язаний з тим, що вони використовують тільки порядок слів у реченні і не працюють з граматичними структурами, отриманими традиційними інструментами, як-от автоматичним морфологічним аналізом. По суті, рекурентні мережі для кожного завдання з нуля «вчаться» граматики мови.

Крім того, рекурентна мережа не буде представлення для проміжних фраз, тому для завдань, у яких потрібні якісні представлення різних фраз у складі речення, використовують *рекурсивні нейронні мережі*.

На відміну від рекурентних, рекурсивні мережі працюють не з ланцюжком слів у реченні, а на основі граматики залежностей – для кожного речення будується бінарне дерево для його розбору. Роботу рекурсивної мережі можна описати ось як. Спочатку вона обробляє листочки дерева розбору (листочки дерева – вказівники на два слова речення і на тип граматичної залежності між ними), заміщаючи листочки отриманим вектором тієї ж розмірності, що і вектори слів. І продовжує працювати далі, але тепер листочки вже об'єднують фрази, а не слова – будуються векторні представлення фраз речення. Отже, маючи дерево розбору, можна побудувати рекурсивну мережу з такою ж топологією, як і дерево, замінивши кожен вузол дерева на нейронну мережу. Природно, всі розмножені у такий спосіб мережі мають спільні параметри, тобто під час навчання та експлуатації робота йде з однією мережею.

Під час навчання рекурсивна мережа може навчитися робити якісні представлення не тільки для повних речень, але і для всіх фраз речення. Водночас нейромережа може послабити ефект помилок граматичного розбору. Таким чином, мережа дозволяє визначити міру семантичної близькості як для слів, так і для всіх фраз у реченні. Якщо в рекурсивну нейронну мережу додати елементи пам'яті LSTM, то можна отримати дуже якісні векторні представлення (Tai).

Інший підхід для отримання векторів речення полягає в тому, що для кожного речення, параграфа або цілого документа тренується окремий вектор, який також бере участь в прогнозі контексту кожного слова речення або параграфа, і в процесі навчання вибираються вектори, які найбільшою мірою

поліпшують передбачення. За якістю отриманих векторів цей метод (його зазвичай називають doc2vec) змагається з рекурсивними нейромережами, водночас для навчання не потрібна розмічена навчальна вибірка. Щоправда, у цього методу є суттєві недоліки: йому потрібні великі речення або цілі параграфи – він не працює на рівні коротких фраз; і він вимагає значних обчислювальних потужностей.

Ще один підхід до моделювання слів і речення – нейромережі, що працюють із символічними представленнями слів або змішаними представленнями. Ці мережі успішно використовують для доповнення векторних представлень слів. Наприклад, у навчальній вибірці не було слова «побігати», але слово «бігати» траплялося часто, і було отримано його якісне представлення. Тоді нейромережа зі змішаним представленням зможе отримати вектор слова «побігати», застосувавши префікс «по». Тобто нейромережі із символічними представленнями навчаються використовувати морфологію слів, а не тільки працювати з певними словами як неподільними сутностями.

Отже, за допомогою рекурентних і рекурсивних нейромереж можна ефективно розв'язувати прості завдання, пов'язані з автоматичною обробкою текстів: класифікації, визначення тональності, виділення іменованих сутностей, простих фактів тощо.

3. Методи для обробки довільних текстів

Обробка текстів, що складаються з декількох речень, які потрібно розглядати не як незалежні сутності, а як взаємопов'язаний ряд висловлювань, для всіх наявних технологій є суттєвою проблемою. У цьому разі виникає семантичний контекст, який збагачується і модифікується наступними реченнями, і моделювати його дуже складно. У комп'ютерній лінгвістиці вже багато років не має остаточного прийнятного рішення завдання

корелюваності або анафоричних відношень. Тільки значне обмеження прикладної сфери дозволяє будувати прийнятні семантичні моделі для текстів, що складаються з декількох речень, і діалогів.

Література

1. Bartunov, Sergey et al. "Breaking Sticks and Ambiguities with Adaptive Skip-gram." arXiv preprint arXiv:1502.07257 (2015). Web. 10 Sep. 2016. **2. Xu Chang et al.** "Rc-net: A general framework for incorporating knowledge into word representations." Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management. ACM, 2014. Web. 10 Sep. 2016. **3. Graves Alex.** "Generating sequences with recurrent neural networks." arXiv preprint arXiv:1308.0850 (2013). Web. 10 Sep. 2016. **4. Jin Zeyu et al.** "Cute: A concatenative method for voice conversion using exemplar-based unit selection." 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2016. Web. 10 Sep. 2016. **5. Tai Kai Sheng, Richard Socher and Christopher D. Manning.** "Improved semantic representations from tree-structured long short-term memory networks." arXiv preprint arXiv:1503.00075 (2015). Web. 10 Sep. 2016. **6. Mikolov, Tomas et al.** "Efficient estimation of word representations in vector space." arXiv preprint arXiv:1301.3781 (2013). Web. 10 Sep. 2016. **7. Danylyuk I.** "Korpus tekstiv dlya vyvchennya hramatychnoyi sluzhbovosti." (*Text corpora to study of a grammatical auxiliary*) Lihvistychni studiyi (*Linguistic Studies*) 26 (2013): 224-230. Print. **8. Danylyuk I.** "Korpus tekstiv dlya vyvchennya hramatychnoyi sluzhbovosti: klasyfikatsiya hramatychnykh klasiv i pidklasiv." (*Text corpora for studying a grammatical auxiliary: classification of grammatical classes and subclasses*) Lihvistychni studiyi (*Linguistic Studies*) 27 (2013): 221-229. Print. **9. Danylyuk I.** "Teoretychni zasady i metody lindhvopersonolohiyi" (*Theoretical Principles And Methods of*

Lingvopersonology) *Linhvistychni studiyi (Linguistic Studies)* 31 (2016): 63-66. Print. **10. Yvakhnenko A. H., and V. H. Lapa.** "Kybernetycheskye predskazyvayushchye ustroystva." (*Cybernetic predicting device*). Kyiv: Naulova dumka (1965). Print. **11. Karasyk, V.** *Yazykovoy kruh: lychnost', kontsepty, dyskurs (Linguistic circle: personality, concepts, discourse)*. Volgograd: Peremena, 2002. Print.

Анатолій Загнітко

(чл.-кор НАН України, д. філол. н., проф.)

Донецький національний університет імені Василя Стуса

ПРИНЦИП ІКОНІЧНОСТІ В ПОРІВНЯЛЬНО-ІСТОРИЧНІЙ ЛІНГВОПАРАДИГМІ

Розглянуто особливості реалізації принципу іконічності в порівняльно-історичній лінгвопарадигмі із застосуванням окремих елементів функційно-когнітивної лінгвопарадигми, проаналізовано п'ятикрокову методику дослідження реалізації принципу іконічності в лексиці індоєвропейської прамови (перший – встановлення особливостей структури слова в індоєвропейській прамові; другий – аналіз структурної організації лексики індоєвропейської прамови; третій – аналіз формально-змістових імплікацій у коренях індоєвропейської прамови; четвертий – аналіз фонологічно-змістових контрастів в аспекті мовних змін; п'ятий – розгляд фонологічно-змістових контрастів у мовній картині світу), які обґрунтовані Т. Козловою в низці студійовань. Звернуто увагу на функційне співвідношення принципу іконічності та принципу економності в лексичному та граматичному ладі мови.

Ключові слова: принцип іконічності, принцип економності, порівняльно-історична лінгвопарадигма, функційно-когнітивна лінгвопарадигма, моносемія, полісемія, омонімія.

ЛІНГВОКОМП'ЮТЕРНІ ДОСЛІДЖЕННЯ

Збірник наукових праць

Випуск 11

Макетування, технічне оформлення: *Ілля Данилюк*

Адреса редколегії:

21007, м. Вінниця, вул. академіка Янгеля, 4, ауд. 320

Підписано до друку: 01.03.2018. Формат 60x84 1/16.

Вид. друк. арк. 8,5. Друк лазерний. Зам. № 16. 300 прим.