

БРАНИСЛАВА М. ДИЛПАРИЋ¹

УНИВЕРЗИТЕТ У ПРИШТИНИ СА ПРИВРЕМЕНИМ СЕДИШТЕМ
У КОСОВСКОЈ МИТРОВИЦИ, ФИЛОЗОФСКИ ФАКУЛТЕТ
КАТЕДРА ЗА ЕНГЛЕСКИ ЈЕЗИК И КЊИЖЕВНОСТ

НИНА С. ПЕРОВИЋ²

УНИВЕРЗИТЕТ У БЕОГРАДУ
ФИЛОЛОШКИ ФАКУЛТЕТ
МОДУЛ: ЈЕЗИК

МОДЕЛИ ЛЕКСИЧКЕ СЕМАНТИКЕ У АЛГОРИТМИМА ЗА ОБРАДУ ПРИРОДНОГ ЈЕЗИКА

САЖЕТАК. У раду се представљају резултати испитивања које је имало за циљ да утврди да ли се у структури алгоритама за обраду природног језика могу препознати карактеристичне особености појединих приступа лексичке семантике. Присуство карактеристика декомпозиционог, холистичког и релационог модела лексичке семантике испитивано је у структури две врсте алгоритама: алгоритму *word2vec* и алгоритмима рекурентних неуронских мрежа. Добијени компаративном анализом претходно описаних модела и алгоритама, резултати показују да постоји значајно преклапање у испитиваним лингвистичким и програмским техникама за обраду природног језика.

КЉУЧНЕ РЕЧИ: лексичка семантика; рачунарска лингвистика; декомпозициони, холистички и релациони модел; обрада природног језика (ОПЈ); алгоритам; *word2vec*; рекурентна неуронска мрежа (РНМ).

¹ branislava.dilparic@pr.ac.rs

² stracimir@yandex.com

Рад је примљен 17. фебруара 2021, а прихваћен за објављивање на састанку Редакције Зборника одржаном 24. марта 2021.

УВОД

Двадесет први век може се сматрати периодом дигиталне ренесансе имајући у виду колики је успех човек постигао у остварењу комуникације са машином помоћу алгоритама вештачке интелигенције (ВИ).³ Ова комуникација, заправо, почела се развијати од појаве првих рачунара и у почетку је била формална и математичка. Стога је требало пронаћи медијум помоћу ког би се та комуникација могла „продубити“; другим речима, требало је пронаћи начин за обраду природног језика (ОПЈ), тј. за анализирање корпуса природних људских језика, као и њихово аутоматско произвођење и разумевања.

Темељи рачунарске лингвистике и ОПЈ постављени су средином прошлог века када су се, независно једна од друге, одиграле револуције на пољу рачунарства и лингвистике. Најпре је 1950. године британски математичар и отац модерног рачунарства и вештачке интелигенције Алан Тјуринг објавио рад *Computing Machinery and Intelligence*, у ком предлаже и детаљно описује тзв. Тјурингов тест⁴ и по први пут у историји званично поставља питање о томе „да ли машине могу да мисле“ (Turing, 1950, стр. 433). Неколико година касније, тачније 1957. године, амерички лингвиста и један од највећих мислилаца данашњице Ноам Чомски објављује књигу *Syntactic Structures*, која је била револуционарна на више нивоа. У њој је Чомски први пут представио своје „бездојне зелене идеје [које] бесно спавају“⁵, односно раздвојио значење од синтаксе и, још важније, поставио темеље формалних граматика и тиме започео еволуцију програмских језика и ОПЈ. Од тог времена математика и лингвистика постале су нераскидиво повезане у тежњи да изнова и изнова померају границе лингвистичке и технолошке револуције.

Данас се софтвери ОПЈ користе свакодневно а да човек тога можда није ни свестан. Сваки пут када се нешто укуца у претраживач, активира се неки део ОПЈ, нпр. анализирање префикса и су-

³ Вештачка интелигенција је област рачунарства која се бави дизајнирањем паметних машина, тј. рачунарских система способних да обављају задатке који изискују људску интелигенцију попут разумевања (природног људског) језика, учења, закључивања, решавања проблема итд. (Barr & Feigenbaum, 1981, стр. 3).

⁴ Тјурингов тест је експеримент у ком се утврђује способност машине да демонстрира интелигентно понашање еквивалентно људском (више у: Turing, 1950).

⁵ Према енгл. *Colorless green ideas sleep furiously* (Chomsky, 2002, стр. 15).

фикса речи, аутоматско допуњавање речи, исправљање правописних грешака и др. Овакви софтвери представљају тзв. ОПЈ софтвере ниског нивоа пошто се баве простим проблемима и не користе технику машинског учења. На другој страни клатна, тј. у самом врху ОПЈ налазе се софтвери који користе технике машинског учења попут дубинског учења и рекурентних неуронских мрежа (РНМ), донекле моделованих према начину на који учи човек. Из области препознавања говора и машинског превођења изродили су се софтвери као што су *Siri*, *Cortana* и *Amazon Alexa*, виртуелни асистенти, редом, Епла, Микрософта и Амазона; *Amazon Transcribe*, Амазонов софтвер за транскридовање у реалном времену; *Google Translate* и *Google Assistant* итд.

ЦИЉ, ПРЕДМЕТ И МЕТОДОЛОГИЈА ИСТРАЖИВАЊА

У овом раду анализирају се и пореде карактеристике модела из домена рачунарске лингвистике (тј. алгоритама за ОПЈ), с једне, и модела лексичке семантике, с друге стране. Циљ истраживања био је да се утврди да ли се у архитектури алгоритама *word2vec* и алгоритама РНМ могу препознати карактеристичне црте појединих приступа лексичке семантике у тумачењу значења речи.⁶ Основна претпоставка рада је да *word2vec* носи карактеристике декомпозиционог и релационог модела лексичке семантике, а да се у алгоритмима РНМ могу препознати особине холистичког приступа проучавању лексичког значења. У раду се најпре дескриптивном методом представљају наведени лингвистички модели, дефинишу се њихови принципи и објашњавају начини манифестације њихових својстава у језику; потом се, истоветно, описују и два споменута алгоритама и анализира се начин на који они обрађују текстуалне корпуре. Напоследку, компаративном анализом се утврђују евентуална преклапања у разматраним лингвистичким и математичким техникама за ОПЈ.

КАРАКТЕРИСТИКЕ ИСПИТИВАНИХ МОДЕЛА ЛЕКСИЧКЕ СЕМАНТИКЕ

Значење у природном људском језику представља један од најизазовнијих проблема у ВИ данас. Колико је комплексно говорити

⁶ Томаш Миколов, предводник инжењерског тима који је дизајнирао *word2vec*, у приватној преписци са ауторкама рада наглашава да концепт програма *word2vec* није настао под утицајем било које од лингвистичких теорија о значењу, већ да алгоритам има посве математичку позадину.

о значењу језичких јединица казује и чињеница да се овим концептом данас баве све гране традиционалне лингвистике, потом рачунарска лингвистика, когнитивна лингвистика, филозофија, психологија итд. Иако се о значењу може говорити и на свим нивоима језичке структуре, у раду се овај широки опсег сужава на поље лексичке семантике која се бави дефинисањем значења речи. На први поглед, дефинисање лексичког значења делује као једноставан задатак; међутим, поставља се питање да ли је икада било могуће прецизно дефинисати неку категорију и све њене потенцијалне припаднике. Лудвиг Витгенштајн је 1953. године потресао класични закон о идентитету и увео појам тзв. нејасних граница:

„Примера ради, размислимо о низу активности које називамо *играма* као што су *грушћивене игре*, *игре са картицама*, *игре са лоптом*, *Олимпијске игре* итд. Шта им је заједничко? Не можемо само рећи да им нешто мора бити заједничко јер се иначе не би називале *играма*, већ треба *погледајући* и *уочијући* има ли ичег заједничког код свих њих. Јер ако их заиста погледамо, нећемо уочити нешто што је *свима* заједничко, већ само сличности и везе међу њима, и то читав један низ“⁷ (Wittgenstein, 1986, стр. 31).

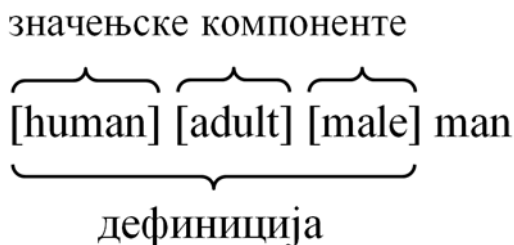
На овом примеру Витгенштајн у ствари показује колико раскошан може бити спектар припадника једне једине категорије и колико је истовремено веома тешко утврдити која је то карактеристика која их ставља у исти оквир. Ове нејасне границе данас су постале још нејасније с обзиром на то да човек живи у времену у коме филозофија флуидности прожима све сфере његове делатности и дивствовања, од идентитета, преко религије до уметности. Стога је лексичка семантика, покушавајући да одговори на овакве изазове, изродила бројне теорије о тумачењу значења речи.⁸ За потребе овог рада у фокус се стављају три приступа: декомпозициони, холистички и релациони.

⁷ “Consider for example the proceedings that we call *games*. I mean board-games, card-games, ball-games, Olympic games, and so on. What is common to them all? Don't say: 'There *must* be something common, or they would not be called *games*'. – But *look and see* whether there is anything common to all. – For if you look at them you will not see something that is common to *all*, but similarities, relationships, and a whole series of them at that.”

⁸ Разуђеност теорија и метода у проучавању лексичког значења од средине XIX века до данашњих дана, заједно са истакнутим публикацијама и истраживачима који су умногоме одредили правце еволуције лексичке семантике, илуструје изузетна књига Дирка Херартса *Theories of Lexical Semantics* (2010).

ДЕКОМПОЗИЦИОНИ МОДЕЛ ЛЕКСИЧКОГ ЗНАЧЕЊА

Основне претпоставке декомпозиционог модела јесу те да се значење речи може изучавати у изолацији, ван било каквог контекста и да се значење речи може раставити на семантичке или значењске компоненте, при чему скуп тих градивних блокова даје дефиницију појма (Katz & Fodor, 1963, стр. 191–192). Принципи ове теорије засновани су на класичној, аристотеловској теорији категоризације и претпоставки да је могуће дефинисати појмове у погледу неопходних и довољних услова: „Како би се анализирано било које референтно значење ... потребно је дефинисати сва *неојходна* и *довољна* обележја која значење једне форме разликују од значења сваке друге форме која може заузимати место на истој семантичкој територији“⁹ (Nida, 1975, стр. 32).



Слика 1: СЕМАНТИЧКА ДЕКОМПОЗИЦИЈА (FIGURE 1: SEMANTIC DECOMPOSITION)

Уколико би се у дефиницији енглеског појма *man* приказаној на Слици 1. изоставила било која од наведених значењских компонента, дефиниција би била непотпуна и отуд нетачна. На пример, ако би се изоставило обележје [HUMAN], којом се мушкарац најпре дефинише као људско биће и тиме се разликује од осталих живих бића [ANIMAL] и [PLANT], онда би се могло говорити о некој животињи која је доживела полну зрелост и која је мушког пола. Дакле, неопходно је излистати све три компоненте. Довољни услови односе се на економичност приликом излиставања компонента за дефинисање значења. Наведеним обележјима би се могло додати још и [KIND], [OUTGOING] или [CAUCASIAN], али се за дефинисање категорије *man* ове компоненте сматрају сувишним.

⁹ “In order to analyze any referential meaning ... one must identify those *necessary* and *sufficient* features that distinguish the meaning of any one form from every other form which might compete for a place within the same semantic territory.”

Како би се упростио и организовао преглед семантичких компонента приликом дефинисања појмова, уведена је тзв. компонентна матрица. Ова матрица је бинарна, тј. садржи две врсте информација, а у датом случају информације се односе на присуство и одсуство неког семантичког обележја. У Табели 1. може се видети на који начин лексичка декомпозиција разликује концепте који се налазе у истом семантичком домену. Примера ради, у појму *husband* биће присутне све семантичке компоненте из матрице, док ће се појам *spouse*, као родно неутралан, одликовати потпуним одсуством семантичког обележја за род.

	HUMAN	ADULT	MALE	MARRIED
<i>man</i>	+	+	+	
<i>woman</i>	+	+	-	
<i>bachelor</i>	+	+	+	-
<i>spinster</i>	+	+	-	-
<i>husband</i>	+	+	+	+
<i>wife</i>	+	+	-	+
<i>spouse</i>	+	+		+

Табела 1: КОМПОНЕНЦИЈАЛНА МАТРИЦА^а (TABLE 1: COMPONENTIAL MATRIX)

^а Матрица је конструисана по узору на Nida (1975).

За разлику од Чомског који је понудио формални запис структуре реченице, декомпозициони приступ, дакле, формализује структуру лексичког значења: дефиниција једног појма такође се записује својеврсним метајезиком (симболима), при чему се и више десетина семантичких обележја могу представити у виду бинарног низа, тј. присуства и одсуства семантичке компоненте.

ХОЛИСТИЧКИ МОДЕЛ ЛЕКСИЧКОГ ЗНАЧЕЊА

Супротно претходном приступу, холизам заступа став да тумачење значења речи искључиво зависи од њеног окружења. Још тридесетих година прошлог века енглески лингвиста Џон Ферт каже: „потпуно значење неке речи увек је контекстуално те се

ниједно испитивање значења не може схватити озбиљно, осим када се врши у оквиру целокупног контекста“¹⁰ (Firth, 1935, стр. 37). Управо на оваквом темељу своју теорију оквира изградио је Чарлс Филмор, објашњавајући њен кључни појам на следећи начин: „Под појмом *оквир* подразумевам било који систем концепата који су повезани тако да разумевање сваког од њих појединачно претпоставља разумевање читаве структуре којој припадају“¹¹ (Fillmore, 1982, стр. 111). На пример, оквир *ѝричешће* одређују основни елементи *свешћеник*, *ѝричешник*, *хлеб* и *вино*. Сви остали појмови који се налазе у значењској димензији концепта *ѝричешће* допуњују тај оквир као што су *свешћа ѡајна*, *исћовесћ*, *молићва*, *крв*, *ћело*, *ѝућир*, *сћојачи* итд. Међутим, комбинација елемената *молићва*, *ѝричешник*, *хлеб*, *вино* не могу одређивати оквир *ѝричешће* јер се тај чин без свештеника не може обавити.

Природа речи која је осетљива на контекст може се илустровати примером декомпозиционе дефиниције енглеског појма *bachelor*. У зависности од тога у ком се окружењу налази, скуп семантичких обележја [HUMAN] + [ADULT] + [MALE] + [UNMARRIED] може, осим нежење, означавати и патријарха, монаха, Николу Теслу, Деда Мраза, хомосексуалца, римокатоличког свештеника итд. Слично томе, предлошке синтагме *from shore to shore* и *from coast to coast* не могу се користити независно од контекста као апсолутни синоними. Иако се чини да имају исто значење, прва синтагма се употребљава у случају када се од обале до обале иде воденим путем, а друга уколико се од обале до обале иде земљаним путем (Fillmore, 1982, стр. 121). Ови примери јасно показују да значење речи зависи од њеног окружења и „позадинског“, тј. енциклопедијског знања о свету. Стога, може се закључити да основна јединица за тумачење значења не може бити реч, не нужно ни реченица, већ читава ситуација, тј. контекст у коме се говори.

РЕЛАЦИОНИ МОДЕЛ ЛЕКСИЧКОГ ЗНАЧЕЊА

Како Лајонс (1977, стр. 268) наводи, основна претпоставка релационог модела је да се значење једне речи може тумачити уз помоћ различитих значењских веза које она образује са другим ре-

¹⁰ “[...] the complete meaning of a word is always contextual, and no study of meaning apart from a complete context can be taken seriously.”

¹¹ “By the term *frame* I have in mind any system of concepts related in such a way that to understand any one of them you have to understand the whole structure in which it fits.”

чима, при чему се стварају лексичка поља. Да би се образовало лексичко поље, између речи мора постојати неко семантичко преклапање. Чини се да космос и све што је у њему функционише према овом принципу; све је повезано на један или други начин; све се може објаснити у односу на нешто. Нови правци у уметности често настају као реакција на претходне: максимализам се појавио као реакција на минимализам, а минимализам као реакција на амерички експресионизам и модернизам. Такође, антифашизам је настао као реакција на фашизам, а сродан је пацифизму и антимилитаризму. У овим примерима могу се идентификовати две значењске везе, синонимија и антонимија, а лексичка семантика се, поред њих, бави и другим семантичким релацијама попут хипонимије (*purple – violet, lavender, magenta*), тропонимије (*nibble – gorge*), меронимије (*finger – hand*), полисемије (*bright* 1. “shining”; 2. “intelligent”), логичке импликације (*buy – pay*) итд.

Описујући свет око себе, човек неретко користи читав систем речи, а не само једну. Уколико се описује нека биљка, велика је вероватноћа да ће у том опису бити употребљен читав низ именица и придева: *зелена, једнојодишња, вишејодишња, дрвенаста, зимзелена, семе, корен, изданак, њуљак, стабло, трана, крошња, лист, цвеш, њолен, латице, њлог, фотосинтеза, засађивање, влада, саксија, биљкован, рукавице, маказе, њудриво, ливада, земља, заливање, вода, живо биће, сунце, врш* итд. Овај низ појмова налази се у истој значењској димензији и образује лексичко поље именице *биљка* (в. Слика 2).



СЛИКА 2: ЛЕКСИЧКО ПОЉЕ (FIGURE 2: LEXICAL FIELD)

Лексичко поље није скуп насумично генерисаних концепата, већ су појмови нераскидиво повезани „објективно утврдивим смисаоним везама“ (Prčić, 2016, стр. 138). Један једини појам може бити окидач за читав низ асоцијација, синтагматских или парадигматских. Људи свакодневно тумаче свет око себе на овај начин, повезујући појмове, те би релациони приступ требало разматрати са нарочитом пажњом приликом дефинисања значења концепата.

КАРАКТЕРИСТИКЕ ИСПИТИВАНИХ АЛГОРИТАМА ЗА ОПЈ

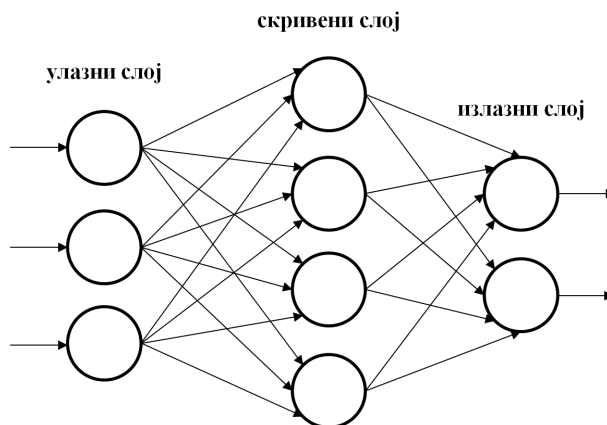
Један од кључних термина у области ВИ јесу неуронске мреже: „Неуронске мреже представљају биолошки инспирисану парадигму (која подражава функционисање мозга сисара) која рачунару омогућава да усвоји људске способности на основу учења посматрањем“¹² (Goyal, Pandey, Jain, 2018, стр. 39). Једну неуронску мрежу чини скуп алгоритама који слојевито обрађују прослеђене им информације. Може се рећи да неуронске мреже одговарају моделу учења са разумевањем, а да је све што им је претходило било учење напамет.

Примера ради, машинско превођење је у најранијим фазама функционисало према статистичком моделу. Направљени су двојезични или вишејезични корпуси и прослеђени су алгоритмима како би они напамет „научили“ на који се начин преводе колокације, фразни глаголи, идиоми, свакодневни изрази и сл. Овиме је решен проблем дословног превођења те се, рецимо, *a piece of cake* почело преводити као *йросџо кд љасуљ*, уместо дословног *љарче џорџе*, или *take time* као *иззвојџиџи време*, уместо *наљравџиџи време*, итд. Ипак, овим приступом није решен проблем контекста. Рачунар није могао препознати када лексички спој *a piece of cake* треба превести дословно, а када идиоматски. Стога су слојевитост значења и проблем контекста принудили програмере да напусте овај модел и пронађу нов начин за ОПЈ. Тада су се окренули неуронским мрежама чији је циљ да на сложенији начин, са разумевањем, приступе тумачењу значења.

Неуронска мрежа је вишеслојна структура у којој су вештачки „неурони уређени по слојевима, при чему су улазни и излазни слој развојени групом скривених слојева“¹³ (Aggarwal, 2018, стр.

¹² “Neural networks are a biologically inspired paradigm (imitating the functioning of the mammalian brain) that enables a computer to learn human faculties from observational data.”

4) (в. Сliku 3). Улазни слој чине информације које се прослеђују за обраду, а излазни слој је резултат који се добија по завршетку анализе. У скривени слој су уписани услови под којима се анализа извршава. Сличност са нервним системом огледа се у томе да су сви неурони повезани и да размењују информације зарад успешног анализирања података.

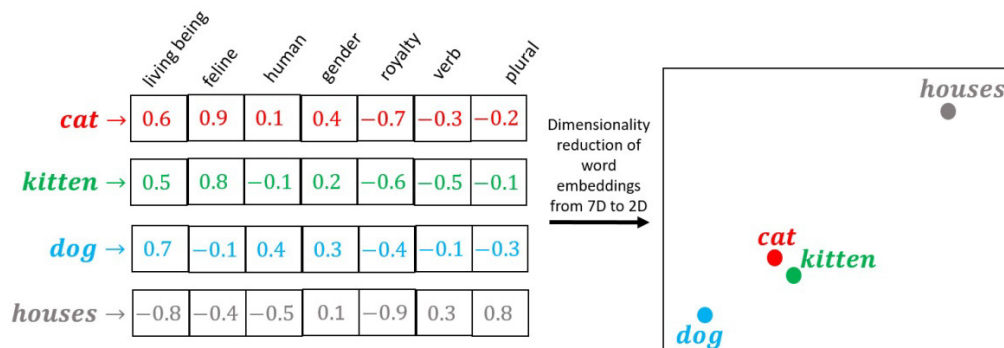


СЛИКА 3: НЕУРОНСКА МРЕЖА (FIGURE 3: NEURAL NETWORK)

АЛГОРИТАМ *WORD2VEC*

Уз помоћ неуронских мрежа алгоритам *word2vec* или *word to vector* (Mikolov & al., 2013a; Mikolov & al., 2013b) претвара речи у математичке векторе. Алгоритму се проследи велики текстуални корпус у коме он покушава да пронађе образац, у овом случају речи које се заједно појављују у текстовима. На пример, мало је вероватно да ће се у текстовима о квантној физици говорити о летовању, кућним љубимцима или викторијанској књижевности, а врло је вероватно да ће се заједно наћи речи као што су *квантна механика, честичице, атоми, теорија, сјање, фреквенција, правилица, Ајнштајн* и сл. На овај начин се групишу речи према блискости, односно контексту у коме се јављају, и формира се векторски простор (в. приказ векторског простора енглеског глагола *fly* на Слици 4).

¹³ “[In multi-layer neural networks], the neurons are arranged in layered fashion, in which the input and output layers are separated by a group of hidden layers.”



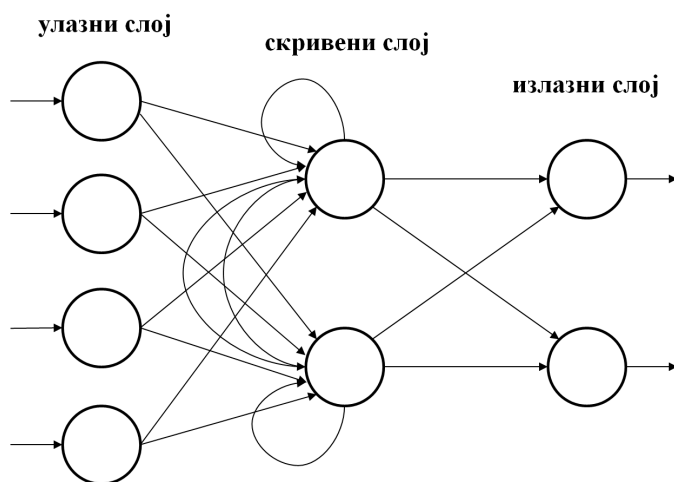
Слика 5: ДИМЕНЗИЈЕ ВЕКТОРСКОГ ПРОСТОРА (FIGURE 5: VECTOR SPACE DIMENSIONS)
 (Преузето са: <https://medium.com/@hari4om/word-embedding-d816f643140>
 [НОВ. 2020].)

Уколико се даље анализирају ове вредности, може се приметити да оне нису насумичне. Према овим бројевима најпре се може закључити да се вредности крећу од 0 до 1 и од 0 до -1. Чини се да је алгоритам исправно израчунао да ће, на пример, *cat* имати негативне вредности за димензије *royalty*, *verb* и *plural*, а високу вредност за *living being* и скоро максималну вредност за *feline*. Слично томе, појам *houses* нема максималну негативну вредност за димензију *living being* као што би се могло очекивати. То значи да алгоритам из текстуалног корпуса успешно препознаје да *houses* као појам није у потпуности лишена живота, већ да се користи за живљење (становање) и да у њој неко живи (људи, биљке, кућни љубимци). На овај начин *word2vec* дефинише различите семантичке везе које постоје између речи, а уз то може вршити и предвиђања у различитим семантичким аналогијама. Mikolov & al. (2013b, стр. 3111) наводе пример недовршене аналогије *Spain : Madrid = France : ?*, коју алгоритам успешно може довршити именицом *Paris*.

Међутим, иако се на овај начин успешно записује слојевито значење речи, и даље није решен проблем вишезначности и контекста. *Word2vec* производи фиксне векторе речи у којима се бројеви не мењају, што би значило да реч увек има једно апсолутно значење. Насупрот томе, речи су по природи полисемичне те су програмери наставили да развијају софтвере у покушају да са владају и ове изазове.

АЛГОРИТМИ РЕКУРЕНТНИХ НЕУРОНСКИХ МРЕЖА

Неуронска мрежа о којој је досада било речи је ациклична. То значи да је смер извршавања анализе једносмеран: улазни слој → скривени слој → излазни слој. Стога *word2vec* даје фиксне векторе који се не могу мењати. Насупрот томе, РНМ је циклична (Mikolov & al., 2010, стр. 1045) и способна да мења вредности вектора у зависности од њиховог окружења (в. Сliku 6). Ово је омогућено тзв. дугом краткорочном меморијом (Boden, 2001, стр. 22) уз помоћ које неурони памте одлуке које су донели раније и то знање користе за тумачење нових информација.



СЛИКА 6: РЕКУРЕНТНА НЕУРОНСКА МРЕЖА (FIGURE 6: RECURRENT NEURAL NETWORK)

Примера ради, када је у питању машинско превођење, РНМ ће у сваком тренутку „гледати“ неколико речи испред и неколико речи иза оне која је у том тренутку у фокусу, те у реалном времену прилагођавати вредности вектора у зависности од контекста. Овакав приступ умногоме је побољшао „вештачко“ разумевање значења.

АНАЛИЗА И ДИСКУСИЈА

Описане кључне карактеристике модела лексичке семантике и алгоритама упоређиване су према критеријуму присуства, односно одсуства четири елемената: семантичких обележја, матрице, значењског поља и контекста. На основу ових фактора говориће

се о уоченим сличностима и разликама између поменутих техника за ОПЈ, а преглед резултата компаративне анализе приказани су у Табели 2.

	СЕМАНТИЧКА ОБЕЛЕЖЈА	МАТРИЦА	ЗНАЧЕЊСКО ПОЉЕ	КОНТЕКСТ
ЛЕКСИЧКА ДЕКОМПОЗИЦИЈА	+	+	-	-
РЕЛАЦИОНИ ПРИСТУП	-	-	+	-
ХОЛИСТИЧКИ ПРИСТУП	-	-	+	+
WORD2VEC	+	+	+	-
РЕКУРЕНТНА НЕУРОНСКА МРЕЖА	+	+	+	+

ТАБЕЛА 2: КОМПАРАТИВНА ТАБЕЛА (TABLE 2: COMPARISON TABLE)

Значајно преклапање у присуству/одсуству компаративних елемената може се осматрети у лексичкој декомпозицији и алгоритму програма *word2vec*. Један од кључних концепата у декомпозиционом приступу је компоненцијална матрица уз помоћ које се значење речи може записати у облику бинарне информације. Уколико се плусеви и минуси замене бројевима, тј. јединицама и нулама, при чему ће јединица означавати присуство, а нула одсуство неког семантичког обележја¹⁴, добија се вектор (в. Табелу 3). Вектор је ништа друго до матрица са бројчаним вредностима. У овом случају, примери су два вектора, *man* и *woman*, који су тродимензионални (HUMAN, ADULT, MALE).

	HUMAN	ADULT	MALE
<i>man</i>	1	1	1
<i>woman</i>	1	1	0

ТАБЕЛА 3: НАСТАЈАЊЕ ВЕКТОРА (TABLE 3: VECTOR FORMATION)

Семантичка обележја су у потпуности аналогна векторским димензијама у алгоритму *word2vec*, а компоненцијална матрица истоветна је векторској матрици тог алгоритма. Оба модела за

¹⁴ Сличан принцип користи се у бинарном систему из информатике и рачунарства, при чему се јединицама и нулама конвенционално записује присуство и одсуство напона (Manojlović, 2007, стр. 54).

циљ имају раслојавање значења речи на његове компоненте, с том разликом што *word2vec* то чини на знатно већој скали од декомпозиционог приступа – од педесет до триста значењских компоненти по једној речи.

Додатна сличност ових модела огледа се у изостављању контекста као технике којом се обрађује језик. У случају лексичке декомпозиције број семантичких обележја је ограничен и не може се променити, као што се не може променити ни вредност вектора у алгоритму *word2vec*. То значи да се уз помоћ њих не може изучавати вишезначност речи.

Разлика између ова два модела огледа се у одсуству значењског, односно лексичког поља у декомпозиционом приступу. Када алгоритам *word2vec* произведе векторе речи према критеријуму заједничког појављивања у текстовима, формира се вишедимензионални векторски простор, који се потом редукује у дводимензионални како би се могао изучавати. Иако одсутан у декомпозиционом приступу, овај модел је у потпуности истоветан лексичком пољу из релационог приступа лексичке семантике који на математички начин потврђује његово присуство и чак израчунава колико је семантичко преклапање између речи. Векторске вредности нису генерисане случајно, већ се у њима огледају постојеће различите семантичке везе између речи. Слично семантичким обележјима, векторске димензије се изражавају позитивним и негативним вредностима, што се може тумачити као присуство или одсуство неке значењске компоненте.

Холистички приступ и алгоритам *word2vec* одликују се значајним разликама у погледу компаративних елемената. Док *word2vec* почива на присуству семантичких обележја и компонентијалне матрице, холистички приступ као теоријски модел који не подржава изучавање значења у изолацији одликује се одсуством ових елемената и присуством контекста. Овде, ипак, долази до преклапања карактеристика када је у питању значењско поље које се у холистичком моделу и програму *word2vec* у језику манифестује, редом, кроз теорију оквира и векторски простор.

Насупрот свим моделима лексичке семантике и програму *word2vec*, чије вредности варирају од компоненте до компоненте, алгоритми РНМ одликују се присуством свих поредбених елемената. У случају алгоритма *word2vec*, једном када неуронска мрежа обради информацију, добија се непроменљив вектор. Насупрот томе, вредности својих вектора РНМ прилагођава контексту у реалном времену; вредности вектора се мењају заједно са контек-

стом у коме се реч налази. Овим се осигурава да се полисемичност речи успешно идентификује и памти. РНМ посматра реч у односу на све остале речи у корпусу, као и алгоритам *word2vec*, али истовремено посматра и реч и њено непосредно окружење. Оваквим приступом се појачава присуство контекста и на макро нивоу (у читавом корпусу) и на микро нивоу (у непосредном окружењу посматране речи). На пример, посматрајући окружење појма *bachelor* (Katz & Fodor, 1963, стр. 185), РНМ ће препознати да ли ова енглеска именица означава нежењу или особу која има диплому основних академских студија или младог витеза који је у служби другог витеза или младу фоку. Ово је принцип који заговара и семантички холизам – да је значење речи условљено контекстом у коме се јавља и да се стога мењају и семантичка обележја речи у датој ситуацији. РНМ за тумачење значења речи користе реченицу, као и читав корпус, те се због велике ефикасности ове мреже сматрају најнапреднијим моделом за ОПЈ.

Може се унедоглед говорити о томе колико је дефинисање значења речи сложено. Сложено је колико и сам човек, колико и процес учења код човека, колико и начин на који човек доживљава свет, у складу са чим и речима даје нова значења. Из тог разлога и не постоји један апсолутни модел према ком лексичка семантика третира значење. И у оквиру модела напоменутих у овом раду постоје различити приступи – природни семантички метајезик, теорија прототипа, семантика оквира итд. Врло је вероватно да се у структури неких других алгоритама за ОПЈ такође могу препознати карактеристике неког од ових модела за тумачење значења речи.

На основу анализе спроведене у овом раду може се потврдити иницијална претпоставка да се у алгоритму *word2vec* могу идентификовати декомпозициони и релациони модел лексичке семантике, а да се холистички модел може препознати код РНМ.

Занимљиво је уочити да ни један ни други алгоритам не прихватају само један модел лексичке семантике. У основи алгоритма *word2vec* налазе се компоненцијална матрица и лексичко поље. У овом алгоритму декомпозициони и релациони приступи су нераскидиво повезани, што можда и није толико уочљиво када се ова два модела изучавају само у оквирима лексичке семантике. Слично томе, РНМ је надоградња алгоритма *word2vec*, што значи да је поново присутно деловање ова два приступа заједно са семантичким холизмом.

Како су старе теорије о значењу биле оповргаване због великих недостатака у својим моделима, развијале су се нове које би могле надоместити те недостатке. Међутим, овај рад показује да су и поред свих недостатака такве теорије нашле значајно место у рачунарској лингвистици. Исто тако, анализа је показала да се модели морају комбиновати зарад валидних резултата и да, стога, не постоји једна свеобухватна тачна теорија.

ЗАКЉУЧНИ ОСВРТ

Анализа која је спроведена у овом раду требало је да одговори на питање да ли се у алгоритму *word2vec* програма и алгоритмима РНМ могу идентификовати карактеристике три модела лексичке семантике – декомпозициони, холистички и релациони. Приликом упоређивања поменутих лингвистичких и програмских техника, показало се да међу њима постоји значајно преклапање у начинима ОПЈ. Утврђено је да се у алгоритму програма *word2vec* могу препознати својства декомпозиционог и релационог приступа, док се у алгоритмима РНМ могу идентификовати карактеристике холистичког приступа.

Овакви резултати значајни су из више разлога. Пре свега, анализом је уочена изразита испреплетеност лингвистике и рачунарске информатике. Иако структура наведених алгоритама није лингвистички мотивисана, резултати су показали да сличност између математичког поступка за тумачење значења и теоријских модела лексичке семантике није занемарљива. Ова чињеница у великој мери говори о математичкој природи језика и, како је анализа представила, значења. У супротном би се могло говорити само о случајностима када су у питању резултати добијени у овом раду. Потом, показано је да ниједна теорија о значењу није погрешна без обзира на њене мањкавости и, такође, да ниједна теорија није тачнија од друге. Штавише, значење је толико сложено да не може постојати један свеобухватни семантички теоријски модел који би у потпуности могао рашчланити и објаснити све појаве овог феномена. Напослетку, идентификовање карактеристика различитих приступа лексичке семантике у једном алгоритму за ОПЈ показало је да су све теорије о значењу, иако бројне и разнолике у својим теоријским моделима, заиста драгоцене и да се морају удруживати како би се значење могло тумачити на исправан начин.

- ЛИТЕРАТУРА Aggarwal, Ch. C. (2018). *Neural Networks and Deep Learning*. New York: Springer.
- Barr, A. & Feigenbaum, E. A. (1981). *The Handbook of Artificial Intelligence 1*. Stanford, California: HeurisTech Press – Los Althos, California: William Kaufmann, Inc.
- Boden, M. (2001). A Guide to Recurrent Neural Networks and Backpropagation. In: A. Holst (Ed.), *The DALLAS Project* (17–26). Kista: Swedish Institute of Computer Science.
- Chomsky, A. N. (2002). *Syntactic Structures* (2nd edition). Berlin and New York: Mouton de Gruyter.
- Fillmore, Ch. J. (1982). Frame Semantics. In: The Linguistic Society of Korea (Eds.), *Linguistics in the Morning Calm* (111–137). Seoul: Hanshin Publishing Company.
- Firth, J. R. (1935). The Technique of Semantics. *Transactions of the Philological Society*, 1 (34), 36–72. doi: 10.1111/j.1467-968X.1935.tb01254.x
- Geeraerts, D. (2010). *Theories of Lexical Semantics*. Oxford: Oxford University Press.
- Goyal, P., Pandey, S., Jain, K. (2018). *Deep Learning for Natural Language Processing: Creating Neural Networks with Python*. New York: Apress.
- Katz, J. J. & Fodor, J. A. (1963). The Structure of a Semantic Theory. *Language*, 2 (39), 170–210. doi: 10.2307/411200
- Lyons, J. (1977). *Semantics 1*. Cambridge: Cambridge University Press.
- Manojlović, V. V. (2007). *Osnovi računarske tehnike, Deo 1: Podaci i operacije* (2. izdanje). Beograd: Akademska misao.
- Mikolov, T., Karafiat, M., Burget, L., Cernocky, J., Khundanpur, S. (2010). Recurrent Neural Network Based Language Model. In: T. Kobayashi, K. Hirose, S. Nakamura (Eds.), *11th Annual Conference of the International Speech Communication Association* (1045–1048). Makuhari: Interspeech.
- Mikolov, T., Chen, K., Corrado, G., Dean, J. (2013a). Efficient Estimation of Word Representations in Vector Space. In: *Proceedings of Workshop at International Conference on Learning Representations*. Доступно на: <https://iclr.cc/archive/2013/workshop-proceedings.html>.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G., Dean, J. (2013b). Distributed Representations of Words and Phrases and Their Compositionality. In: C. J. C. Burges, L. Botou, M. Welling, Z. Ghahramani, K. Q. Weinberger (Eds.), *Proceedings of the 26th International Conference on Neural Information Processing Systems, Volume 2* (3111–3119). New York: Curran Associates Inc.
- Nida, E. A. (1975). *Componential Analysis of Meaning*. The Hague: Mouton Publishers.

Patel, K. & Bhattacharyya, P. (2017). Towards Lower Bounds on Number of Dimensions for Word Embeddings. In: G. Kondrak & T. Watanabe (Eds.), *Proceedings of the Eighth International Joint Conference on Natural Language Processing, Volume 2: Short Papers* (31–36). Taipei: Asian Federation of Natural Language Processing.

Prčić, T. (2016). *Semantika i pragmatika reči* (3. elektronsko izdanje). Novi Sad: Filozofski fakultet. Доступно на: <http://digitalna.ff.uns.ac.rs/sites/default/files/db/books/978-86-6065-356-9.pdf>.

Turing, A. M. (1950). Computing Machinery and Intelligence. *Mind*, 59 (236), 433–460. <https://doi.org/10.1093/mind/LIX.236.433>

Wittgenstein, L. (1986). *Philosophical Investigations*. Oxford: Basil Blackwell.

BRANISLAVA M. DILPARIĆ

UNIVERSITY OF PRIŠTINA IN KOSOVSKA MITROVICA
FACULTY OF PHILOSOPHY
DEPARTMENT OF ENGLISH LANGUAGE AND LITERATURE

NINA S. PEROVIĆ

UNIVERSITY OF BELGRADE
FACULTY OF PHILOLOGY

SUMMARY

MODELS OF LEXICAL SEMANTICS IN THE ALGORITHMS
FOR NATURAL LANGUAGE PROCESSING

The aim of this study was to determine whether some of the approaches of lexical semantics for studying word meaning could be identified in *word2vec* and *recurrent neural networks* (RNN), the algorithms for natural language processing (NLP). Linguistic concepts from the field of lexical semantics were *decompositional*, *holistic*, and *relational*. Although it is assumed that algorithms for natural language processing cannot be written only on the basis of mathematical knowledge, but also linguistic, this analysis was carried out so as to determine the exact models used in the above-mentioned algorithms.

First, the aforementioned linguistic models were concisely explained through descriptive research. In describing those, the authors of the paper referred to the studies of the most prominent linguists within the field of lexical semantics such as Fillmore, Firth and Lyons. Then, in a similar fashion, the architecture of NLP algorithms was introduced and described. For this intent, the authors of the paper relied mostly on the studies of Mikolov et al., who designed *word2vec*. Next, a comparative analysis was conducted between

approaches of lexical semantics on one hand and the algorithms in question on the other.

This analysis confirmed the underlying assumption of the paper that the characteristics of decompositional and relational approach to studying word meaning could be recognized in *word2vec* and that the properties of the holistic model could be observed in RNN algorithms. More than that, the analysis showed a considerable overlap in processing natural language, as if the models of lexical semantics were taken and mathematically implemented in the algorithms examined. It should be emphasized that for the purposes of this paper only the basic principles of lexical semantics and NLP algorithms were taken into account. The aim of the paper was not to describe edge cases or to talk in detail about the mechanisms of these structures and their advantages and disadvantages. Essentially, this study sought to examine whether the basic ideas and characteristics of lexical semantics could be found in the architecture of the above-noted algorithms.

KEYWORDS: lexical semantics; computational linguistics; decompositional, holistic, and relational models; natural language processing; algorithm; word2vec; recurrent neural network.



Овај чланак је објављен и дистрибуира се под лиценцом Creative Commons Ауторство-Некомерцијално Међународна 4.0 (CC BY-NC 4.0 | <https://creativecommons.org/licenses/by-nc/4.0/>).

This paper is published and distributed under the terms and conditions of the Creative Commons Attribution-NonCommercial International 4.0 licence (CC BY-NC 4.0 | <https://creativecommons.org/licenses/by-nc/4.0/>).