

Vitae

Natalia Darchuk, Doctor of Philology, professor, professor of Department of Ukrainian Language and Applied Linguistic in Taras Shevchenko National University of Kyiv. Her areas of research interests include applied linguistics and computational linguistics.

Correspondence: nataliadarchuk@gmail.com

Ганна Ситар

DOI 10.31558/1815-3070.2018.36.24

УДК 81'373.7:81'32

СТАТИСТИЧНИЙ АНАЛІЗ ЦІЛІСНИХ СЛОВОСПОЛУЧЕНЬ: НА МАТЕРІАЛІ УКРАЇНСЬКОГО НАЦІОНАЛЬНОГО ЛІНГВІСТИЧНОГО КОРПУСУ¹

Стаття продовжує цикл публікацій, присвячених статистичному аналізу стійких одиниць української мови. За даними Українського національного лінгвістичного корпусу визначено ступінь невідповідності поєднання словоформ дво-, три- і чотирикомпонентних цілісних словосполучень української мови шляхом обчислення показника асоціації mutual information (MI).

Для всіх обстежених цілісних одиниць властива невідповідність поєднання словоформ (результати MI перебувають у діапазоні від 8,64 до 44,63). У межах одного корпусу текстів величина MI залежить від таких чинників, як абсолютна частота конструкції, абсолютна частота її компонентів, кількість компонентів і тип цілісного словосполучення.

Ключові слова: показник асоціації, фразеологічна одиниця, mutual information, статистика, цілісне словосполучення, українська мова.

Постановка проблеми, актуальність дослідження. Сучасна лінгвістика визначає статистичні дослідження як виключно корпуснобазовані, тобто вчені вважають, що вірогідні статистичні результати можна одержати тільки на підставі аналізу репрезентативного корпусу текстів. Статистичний аналіз стійких одиниць різних типів, виконаний на корпусному матеріалі, є важливим завданням лінгвістичної статистики, здатним дати чіткі кількісні критерії зарахування мовних одиниць до класу стійких, що можуть бути використані для їх автоматичної ідентифікації в корпусі текстів. Процедура такого аналізу на матеріалі синтаксичних фразеологізмів української мови викладено у працях (Syta, "Statystychni Kryteriyi Analizu Syntaksychnykh Frazеolohizmiv"; Syta, "Statystychni analiz frazeolohizovanykh rechen..."; Syta, "Syntaksychni frazeolohizmy v rozrizi konstruktsiinoi hramatyky").

Обчислення показників (індексів) асоціації (англ. association measures, measures of association) як метод визначення випадковості / невідповідності поєднання певних одиниць може бути застосований для різних типів конструкцій. Пропонована стаття присвячена статистичному аналізу цілісних словосполучень.

У трактуванні словосполучення маємо опертям традиційний підхід, згідно з яким його визначають як непередикативну синтаксичну одиницю, «компонентами якої є слово та форма слова або кілька форм слів, з'єднаних між собою підрядним синтаксичним зв'язком» (Zahnitko, "Slovnyk suchasnoyi linhvistyky: ponyattya i terminy", 1040).

У розгалуженій класифікації словосполучень цілісні² одиниці посідають особливе місце й охоплюють кілька структурних і семантичних різновидів. Вони становлять один із трьох типів словосполучень, які виділяють за ступенем злиття компонентів. За цією ознакою з-поміж словосполучень Анатолій Загнітко розмежує:

- 1) вільні словосполучення;
- 2) синтаксично зв'язані словосполучення;
- 3) фразеологічно зв'язані словосполучення (Zahnitko, "Teoretychna hramatyka ukraiyins'koyi movy: Syntaksys", 63).

На думку мовознавця, визначальними ознаками синтаксично зв'язаних (або нечленованих, неподільних, цілісних) словосполучень є такі: виконання ролі одного члена речення, наявність структури і граматичного

¹ Дослідження виконано в межах наукового проекту «Об'єктивна і суб'єктивна мовносоціумна граматики: комунікативно-когнітивний та прагматико-лінгвокомп'ютерний виміри» (0118U003137) Донецького національного університету імені Василя Стуса.

² В україністиці, крім терміна «цілісні словосполучення», у межах різних підходів до кваліфікації цих одиниць та створення різних класифікацій дослідники використовують також терміни «неподільні словосполучення», «синтаксично неподільні словосполучення», «семантично неподільні словосполучення», «нерозкладні словосполучення», «нечленовані словосполучення», «синтаксично нечленовані словосполучення» і под. (Zahnitko, Balko, Maksymiuk, Lychuk та ін.).

значення, пов'язаність компонентів одним з різновидів підрядного зв'язку, наприклад: *три роки, один із нас, дівчина з гарними очима* (Zahnitko, "Teoretychna hramatyka ukrajins'koyi movy: Syntaksys", 63-64).

Оксана Максим'юк звертає увагу на інформативну недостатність, синсемантичність стрижневих слів, що входять до складу цілісних словосполучень, а також на комплетивні (доповнювальні) відношення як типові для словосполучень цього різновиду, наприклад: *зеряя лебедів, тарілка супу* (Maksymiuk, 8 і далі). За нашими спостереженнями, крім комплетивних відношень, цілісним словосполученням властиві також атрибутивні відношення, наприклад: *чоловік високого зросту, дівчина з карими очима, суддя міжнародної категорії*.

Класифікація цілісних словосполучень на сьогодні залишається дискусійною проблемою синтаксису. У цьому дослідженні враховуємо 10 типів цілісних словосполучень, виділених Мариною Балко:

- 1) словосполучення з кількісним значенням: *два хлопці, кілька дівчат*;
- 2) словосполучення зі значенням вибірковості: *дехто із студентів, четверо з них*;
- 3) словосполучення зі значенням сумісності: *Микола з другом, діти з учителем*;
- 4) словосполучення характеризувальної семантики: *людина інтелектуальної праці, чоловік високого зросту*;
- 5) словосполучення фонові семантики: *тим часом, зимового вечора*;
- 6) словосполучення якісної і станової семантики: *надзвичайна подія, добрий господар*;
- 7) словосполучення з обмежувальною семантикою: *у період з 2014 року по 2018 рік; дорога з міста до села*;
- 8) сполуки фазових, модальних дієслів, слів категорії стану з інфінітивом (складені присудки): *мав працювати, перестати ходити*;
- 9) словосполучення із семантикою невизначеності: *щось біле, щось радикальне*;
- 10) метафоричні й перифрастичні конструкції: *сурма великого горя, геній українського народу* (Balco, 164-168).

Матеріал і методи дослідження. Об'єктом статистичного аналізу стали 56 цілісних словосполучень різних структурних (за кількістю компонентів, їх частининомовним вираженням) і семантичних типів. Серед них 20 двокомпонентних одиниць, 20 – трикомпонентних і 16 – чотирикомпонентних. Вірогідність одержаних результатів забезпечено виконанням обчислень на матеріалі значного за обсягом й індексованого корпусу текстів – Українського національного лінгвістичного корпусу (далі УНЛК) Українського мовно-інформаційного фонду НАН України. Загальна кількість слововживань у корпусі в період здійснення підрахунків становила 189200000 одиниць.

У статистиці відомим є індекс асоціації *mutual information* (англ. взаємна, спільна інформація, далі МІ). За допомогою МІ визначають невідповідність (залежність) послідовності певних явищ або подій, у нашому випадку – словоформ у корпусі текстів (Fano). Вперше в лінгвістичних дослідженнях його застосували Кеннет Ворд Чарч (Kenneth Ward Church) та Патрік Генкс (Patrick Hanks) для виявлення залежності поєднання двох слів в англійських корпусах текстів (Church, Hanks).

Через багатоконпонентність значної частини цілісних словосполучень використовуємо формулу (1), наведену у працях (Petrović, Snajder, Basic, Kolar: 323; Yagunova, Pivovarova: 586) і призначену для конструкцій з будь-якою кількістю компонентів:

$$(1) \quad MI = \log_2 \frac{f(c_1, c_2, \dots, c_i) \times N^{(i-1)}}{f(c_1) \times f(c_2) \times \dots \times f(c_i)}$$

де МІ – коефіцієнт *mutual information*;

i – кількість компонентів конструкції;

*c*₁ – перша лексична одиниця;

*c*₂ – друга лексична одиниця;

*c*_{*i*} – *i*-а лексична одиниця;

f(*c*₁, *c*₂, ..., *c*_{*i*}) – абсолютна частота вживання конструкції *c*₁, *c*₂, ..., *c*_{*i*} в корпусі (з урахуванням порядку одиниць усередині конструкції);

f(*c*₁) – абсолютна частота *c*₁ в корпусі;

f(*c*₂) – абсолютна частота *c*₂ в корпусі;

f(*c*_{*i*}) – абсолютна частота *c*_{*i*} в корпусі;

N – загальна кількість слововживань у корпусі;

\log_2 – логарифм числа за основою 2.

Мета пропонованого дослідження – установити ступінь невідповідності поєднання компонентів цілісних словосполучень шляхом обчислення показника асоціації МІ. Для досягнення поставленої мети розв'язано такі завдання:

- 1) укладено робочий варіант реєстру цілісних словосполучень, що охоплює одиниці з різною кількістю компонентів, різним морфологічним наповненням та нетотожними семантико-синтаксичними відношеннями;
- 2) з УНЛК отримано частотні дані для цілісних словосполучень;
- 3) виконано обчислення за формулою МІ для багатоконпонентних одиниць;

- 4) проаналізовано отримані результати та виявлено взаємозв'язки величини МІ й типу цілісного словосполучення.

Для коректного встановлення абсолютної частоти конструкції та абсолютної частоти окремих словоформ, що входять до її складу, в пошуковій формі УНЛК було задано визначений порядок словоформ та передбачено пошук словоформи, а не слова з урахуванням його парадигми.

Наведемо приклади виконаних обчислень. Для визначення ступеня не випадковості поєднання складників цілісного словосполучення *другого дня* з УНЛК було отримано такі кількісні дані: абсолютна частота конструкції становить 723, абсолютна частота словоформи *другого* – 3656; *дня* – 3623. Підставляючи ці дані до формули (1), отримуємо:

$$MI(\text{другого дня}) = \log_2 \frac{723 \times 189200000}{3656 \times 3623} = 13,334169239 \approx 13,33$$

Коефіцієнт МІ обраховували з точністю до двох знаків після коми. Отримані результати МІ для двокомпонентних цілісних словосполучень подано в таблиці 1.

Таблиця 1

Показник асоціації МІ для двокомпонентних цілісних словосполучень за даними УНЛК

№ з/п	Цілісне словосполучення	Абсолютна частота вживання цілісного словосполучення	Абсолютна частота вживання словоформ-компонентів цілісного словосполучення	Показник асоціації МІ
1	<i>Багато грошей</i>	355	<i>багато 4499; грошей 2537</i>	12,52
2	<i>два хлопці</i>	80	<i>два 4778; хлопці 1887</i>	10,71
3	<i>дещо дивне</i>	14	<i>дещо 3103; дивне 1327</i>	9,33
4	<i>добрий господар</i>	67	<i>добрий 2391; господар 1773</i>	11,55
5	<i>другого дня</i>	723	<i>другого 3656; дня 3623</i>	13,33
6	<i>згряя птахів</i>	17	<i>згряя 709; птахів 1262</i>	11,81
7	<i>зимового вечора</i>	54	<i>зимового 571; вечора 2135</i>	13,03
8	<i>кілька книжок</i>	110	<i>кілька 3673; книжок 1496</i>	11,89
9	<i>літр молока</i>	25	<i>літр 249; молока 1320</i>	13,81
10	<i>мав працювати</i>	21	<i>мав 3731; працювати 2663</i>	8,64
11	<i>надзвичайна подія</i>	71	<i>надзвичайна 580; подія 1471</i>	13,94
12	<i>повинен відповісти</i>	27	<i>повинен 2980; відповісти 1812</i>	9,89
13	<i>почав робити</i>	163	<i>почав 3114; робити 3498</i>	11,47
14	<i>почав ходити</i>	231	<i>почав 3114; ходити 2106</i>	12,70
15	<i>тим часом</i>	2508	<i>тим 5127; часом 4209</i>	14,42
16	<i>троє дівчат</i>	37	<i>троє 1959; дівчат 1680</i>	11,05
17	<i>хочу ніти</i>	80	<i>хочу 2868; ніти 2341</i>	11,14
18	<i>шматок хліба</i>	371	<i>шматок 1488; хліба 1962</i>	14,55
19	<i>щось біле</i>	98	<i>щось 3759; біле 1409</i>	11,77
20	<i>щось радикальне</i>	11	<i>щось 3759; радикальне 200</i>	11,43

Як видно з таблиці 1, коефіцієнт МІ для двокомпонентних цілісних словосполучень становить від 8,64 (*мав працювати*) до 14,55 (*шматок хліба*). При цьому показово, що різні типи цілісних словосполучень мають нетотожні результати, зокрема, низький індекс МІ мають складені присудки, високий – словосполучення фонові (часові) семантики та конструкції з партитивним значенням (останній тип цілісних словосполучень виділяє Оксана Максим'юк (Maksymiuik, 8). Важливим моментом також вважаємо вищі абсолютні частоти як цілісних словосполучень загалом, так і їх складників, порівняно з відповідними частотами лексичних фразеологізмів, прислів'їв та приказок (пор. результати, наведені у працях (Syta, "Syntaksychni frazeolohizmy v rozrizi konstruktsiinoi hramatyky"; Syta, "Statystychnyy analiz prysliv"yiv i prykazok...").

Контрольна величина, починаючи від якої кваліфікуємо поєднання слів як не випадкове, залежить від низки чинників: абсолютної частоти конструкції, абсолютної частоти її окремих складників і від розміру корпусу. Для Українського національного лінгвістичного корпусу, розмір якого під час виконання підрахунків складав 189 200 000 слововживань, ця контрольна величина дорівнює 7,56 (процедуру виведення контрольної величини викладено у праці (Sytar, “Syntaksychni frazeolohizmy v rozrizi konstruktsiinoi hramatyky”, 310-311)):

$$\log_2 189 = 7,56377 \approx 7,56$$

Відповідно результати MI, одержані для всіх проаналізованих двокомпонентних цілісних словосполучень, засвідчують не випадковість поєднання компонентів у їх складі.

Покажемо приклад обчислень для трикомпонентних словосполучень. УНЛК дає такі кількісні дані для конструкції *ніхто з нас*: абсолютна частота конструкції становить 354, абсолютна частота словоформи *ніхто* – 3438, *з* – 6301, *нас* – 4291. Підставляючи ці дані до формули (1), отримуємо:

$$MI(\text{ніхто з нас}) = \log_2 \frac{354 \times (189200000)^2}{3438 \times 6301 \times 4291} = 27,0224621 \approx 27,02$$

Статистичні дані, одержані для три- і чотирикомпонентних цілісних словосполучень, наведено в таблицях 2 і 3 відповідно.

Таблиця 2

Показник асоціації MI для трикомпонентних цілісних словосполучень за даними УНЛК

№ з/п	Цілісне словосполучення	Абсолютна частота вживання цілісного словосполучення	Абсолютна частота вживання словоформ-компонентів цілісного словосполучення	Показник асоціації MI
1	<i>Батько з матір'ю</i>	30	<i>батько</i> 2703; <i>з</i> 6301; <i>матір'ю</i> 577	26,70
2	<i>борошно вищого ґатунку</i>	4	<i>борошно</i> 603; <i>вищого</i> 1528; <i>ґатунку</i> 332	28,80
3	<i>геній українського народу</i>	4	<i>геній</i> 684; <i>українського</i> 2653; <i>народу</i> 2776	24,76
4	<i>дехто зі студентів</i>	5	<i>дехто</i> 1866; <i>зі</i> 4872; <i>студентів</i> 1656	23,50
5	<i>діти шкільного віку</i>	9	<i>діти</i> 2979; <i>шкільного</i> 763; <i>віку</i> 2659	25,67
6	<i>до пізнього вечора</i>	212	<i>до</i> 6278; <i>пізнього</i> 759; <i>вечора</i> 2135	29,47
7	<i>кращий з учнів³</i>	0	-	-
8	<i>людина інтелектуальної праці</i>	3	<i>людина</i> 3660; <i>інтелектуальної</i> 763; <i>праці</i> 3467	23,40
9	<i>ми з другом</i>	12	<i>ми</i> 5060; <i>з</i> 6301; <i>другом</i> 1175	21,57
10	<i>найкращий з усіх</i>	30	<i>найкращий</i> 1612; <i>з</i> 6301; <i>усіх</i> 4321	24,54
11	<i>на другий день</i>	838	<i>на</i> 6328; <i>другий</i> 4156; <i>день</i> 4066	28,06

³ У випадку можливої, але не зафіксованої в УНЛК конструкції (абсолютна частота 0), частоти окремих компонентів не наводимо через те, що обчислення MI не має смислу, оскільки логарифму 0 не існує. Обчислення показників асоціації для таких конструкцій не здійснювали, тому у відповідній графі таблиці стоїть знак «-».

12	<i>ніхто з нас</i>	354	<i>ніхто 3438; з 6301; нас 4291</i>	27,02
13	<i>один із нас</i>	78	<i>один 4977; із 5792; нас 4291</i>	24,43
14	<i>речі великого розміру</i>	1	<i>речі 3121; великого 3198; розміру 1265</i>	21,44
15	<i>система шкільних закладів</i>	0	-	-
16	<i>суддя міжнародної категорії</i>	3	<i>суддя 1071; міжнародної 1385; категорії 1755</i>	25,30
17	<i>сурма великого горя</i>	1	<i>сурма 233; великого 3198; горя 1290</i>	25,15
18	<i>четверо з них</i>	69	<i>четверо 1282; з 6301; них 5527</i>	25,72
19	<i>чоловік високого зросту</i>	7	<i>чоловік 3019; високого 2169; зросту 888</i>	25,36
20	<i>я з подругою</i>	2	<i>я 4815; з 6301; подругою 389</i>	22,53

Таблиця 3

Показник асоціації МІ для чотирикомпонентних цілісних словосполучень за даними УНЛК

№ з/п	Цілісне словосполучення	Абсолютна частота вживання цілісного словосполучення	Абсолютна частота вживання словоформ-компонентів цілісного словосполучення	Показник асоціації МІ
1	<i>Багато з моїх знайомих</i>	1	<i>багато 4499 з 6301; моїх 2395; знайомих 1452</i>	36,00
2	<i>дехто з моїх друзів</i>	5	<i>дехто 1866; з 6301; моїх 2395; друзів 1885</i>	39,21
3	<i>дехто з моїх колег</i>	4	<i>дехто 1866; з 6301; моїх 2395; колег 1082</i>	39,69
4	<i>дівчина з синіми очима</i>	3	<i>дівчина 2035; з 6301; синіми 585; очима 2904</i>	39,76
5	<i>дівчина з карими очима</i>	3	<i>дівчина 2035; з 6301; карими 232; очима 2904</i>	41,10
6	<i>жінки з довгими ногами</i>	1	<i>жінки 2672; з 6301; довгими 1286; ногами 2294</i>	36,99
7	<i>з вересня по грудень</i>	3	<i>з 6301; вересня 1888; по 6024; грудень 660</i>	38,64
8	<i>з Києва до Вінниці</i>	1	<i>з 6301; Києва 1979; до 6278; Вінниці 434</i>	37,54
9	<i>з міста до села</i>	6	<i>з 6301; міста 3063; до 6278; села 2396</i>	37,03
10	<i>з січня по жовтень</i>	6	<i>з 6301; січня 1703; по 6024; жовтень 806</i>	39,50
11	<i>із села до міста</i>	7	<i>із 5792; села 2396; до 6278; міста 3063</i>	37,37
12	<i>людина з великим серцем</i>	4	<i>людина 3660; з 6301; великим 2738; серцем 1702</i>	37,88
13	<i>людина з добрим серцем</i>	2	<i>людина 3660; з 6301; добрим 1513; серцем 1702</i>	37,73
14	<i>майстер спорту міжнародного класу</i>	13	<i>майстер 1303; спорту 826; міжнародного 1446; класу 2075</i>	44,63
15	<i>хтось з наших людей</i>	2	<i>хтось 3186; з 6301; наших 3266; людей 4272</i>	35,49
16	<i>чоловік з високим чолом</i>	1	<i>чоловік 3019; з 6301; високим 1942; чолом 709</i>	37,91

Дані, наведені в таблицях 2 і 3, засвідчують, що коефіцієнт МІ для всіх цілісних словосполучень є високим: для трикомпонентних одиниць він перебуває в межах від 21,44 (*речі великого розміру*) до 29,47 (*до пізнього вечора*), тобто у 2,8 – 3,9 разу більший за 7,56; для чотирикомпонентних – від 35,49 (*хтось з наших людей*) до 44,63 (*майстер спорту міжнародного класу*), тобто в 4,7 – 5,9 разу більший за контрольну величину.

З-поміж трикомпонентних одиниць низькі результати МІ мають словосполучення зі значенням сумісності, у складі яких є займенник, середні – словосполучення зі значенням вибірковості, високі – словосполучення темпоральної семантики, зокрема внаслідок високої абсолютної частоти конструкцій у корпусі текстів. Цікаво, що метафоричні словосполучення як одиничні, часто оказіональні, утворення мають низьку абсолютну частоту, проте досить високий показник МІ, очевидно, внаслідок порівняно невисоких частот їх складників.

Усі чотирикомпонентні цілісні словосполучення мають низьку абсолютну частоту вживання в корпусі текстів. Діапазон індексу МІ для цих одиниць є невеликим, близькі результати МІ зафіксовано для словосполучень зі значенням вибірковості та словосполучень з обмежувальною семантикою, вищі – для словосполучень характеризувальної семантики за умови невисокої частоти хоча б одного з її компонентів.

Статистично вірогідного зв'язку між результатами МІ та типом переданого семантико-синтаксичного відношення не зафіксовано.

Висновки. Отримані результати обчислень для цілісних словосполучень, виконаних за даними Українського національного лінгвістичного корпусу, доводять, що для всіх проаналізованих одиниць властива не випадковість поєднання словоформ: коефіцієнт МІ перебуває в діапазоні від 8,64 до 44,63. У межах одного корпусу текстів величина МІ залежить від таких чинників, як абсолютна частота конструкції, абсолютна частота її компонентів, кількість компонентів і тип цілісного словосполучення.

Наведені статистичні дані корелюють із результатами, отриманими для інших типів фразеологічних одиниць – лексичних фразеологізмів, синтаксичних фразеологізмів, прислів'їв і приказок. Результати МІ для цілісних словосполучень є прогнозовано нижчими порівняно з іншими стійкими одиницями через меншу кількість компонентів у їх складі.

Перспективи. Подальший етап дослідження передбачає виконання статистичного аналізу інших типів стійких одиниць, зокрема складених прийменникових еквівалентів та ін.

References

Balko, Maryna. “*Semantyko-syntaksychni i strukturni aspekty tsilisnykh slovospoluchen' suchasnoyi ukrayins'koyi movy (Semantic-syntactical and Structural Aspects of Indivisible Word-combinations of Modern Ukrainian Language)*”: Diss. Zaporizhzhya National U, 2004. Abstract. Print.

Balko, Maryna. *Aktual'ni problemy teorii slovospoluchennya suchasnoyi ukrayins'koyi movy (Actual Problems of the Word-combination Theory of Modern Ukrainian Language)*: [monohrafiya]. Dnipropetrovs'k: Svidler, 2014. Print.

Church, Kenneth Ward, and Hanks, Patrick. “Word Association Norms, Mutual Information, and Lexicography.” *Computational Linguistics* 16(1) (1990): 22–29. Print.

Everitt, B. S. *The Cambridge Dictionary of Statistics*. 2nd edition. Cambridge: Cambridge University Press, 2002. Print.

Fano, Robert M. *Transmission of Information: A Statistical Theory of Communications*. The Technology Press, M.I.T., and John Wiley & Sons, Inc., New York, 1961. Print.

Lychuk, Mariya. “Syntaksychno nechlenovani slovospoluchennya: ustalenist' termina, istoriya doslidzhennya (Syntactically Nondivided Word-combinations: Term Sustainability, History of Research)” *Linguistic Bulletin*, 21 (2016): 142–148. Print.

Maksymiuk, Oksana. “*Koreferentnist' nerozkladnykh komponentiv u strukturi rechennya (Co-reference of Stable Components in the Structure of the Sentence)*.” Diss. Chernivtsi National U, 2005. Abstract. Print.

Petrovic, S., Snajder, J., Basic, B.D., Kolar, M. “Comparison of collocation extraction for document indexing.” *Journal of Computing and Information Technology*, 14 (4) (2006): 321–327. Print.

Sytar, Hanna. “Statystychni Kryteriyi Analizu Syntaksychnykh Frazeholohizmiv (Statistical Criteria of Analysis of Syntactic Idioms)” *Visnyk Donets'koho Natsional'noho Universytetu. Seriya B. Humanitarni Nauky (The Bulletin of Donetsk National University. Series B. Humanities)* 1-2 (2015): 245–256. Print.

Sytar, Hanna. “Statystychnyi analiz pryslyv'yiv i prykazok: pokaznyk asotsiatsiyi mutual information (na materialy Ukrayins'koho natsional'noho linhvistychnoho korpusu) (Statistical Analysis of Proverbs and Sayings: Association Measure of Mutual Information (on material of Ukrainian National Linguistic Corpus))” *Lingvisticnyi smydyi / Linguistic Studies* 35 (2018): 170–177. Print.

Sytar, Hanna. “Statystychnyi analiz frazeolohizovanykh rechen: pokaznyk asotsiatsii mutual information (Statistical Analysis of Sentences with Phraseological Structures: Association Measure of Mutual Information)” *Ukrainske movoznavstvo (Ukrainian Linguistics)*. 1(46) (2016): 103–125. Print.

Sytar, Hanna. *Syntaksychni frazeolohizmy v rozryzi konstruktsiinoi hramatyky (Syntactic Idioms in the Context of Construction Grammar)*. Vinnytsya: TOV «Nilan-LTD», 2017. Print.

Ukrainska mova: Entsyklopediia (Ukrainian language: Encyclopedia). Redkol.: Rusanivskiy V. M. (spivholova), Taranenko O. O. (spivholova), Ziabliuk M. P. ta in. 2-he vyd., vypr. i dop. Kyiv: Vyd-vo "Ukrainska entsyklopediia" im. M. P. Bazhana, 2004. Print.

Yagunova, Ye.V., Pivovarova, L.M. "Ot kollokatsiy k konstruksiyam (From Collocations to Constructions)". *АКТА LINGUISTICA PETROPOLITANA. Works of the Institute of Linguistic Researches of RAS, Russkiy yazyk: grammatika konstruksiy i leksiko-semanticheskie podkhody (The Russian Language: Construction Grammar and Lexical and Semantic Approaches)*: X, part 2. (2014) 568-617. Print.

Zahnitko, Anatoliy. *Slovnnyk suchasnoyi linhvistyky: ponyattya i terminy (Dictionary of Modern Linguistics: Concepts and Terms)*. Donets'k: DonNU, 2013. Print.

Zahnitko, Anatoliy. *Teoretychna hramatyka ukrayins'koyi movy: Syntaksys (Theoretical Grammar of the Ukrainian Language: Syntax)*. Donets'k: DonNU, 2001. Print.

Надійшла до редакції 23 листопада 2018 року.

STATISTICAL ANALYSIS OF INDIVISIBLE WORD-COMBINATIONS: ON MATERIAL OF UKRAINIAN NATIONAL LINGUISTIC CORPUS

Hanna Sytar

Department of General and Applied Linguistics and Slavonic Philology,
Vasyl' Stus Donetsk National University, Vinnytsia, Ukraine

Abstract

Background: The article is devoted to the statistical analysis of indivisible word-combinations in the Ukrainian. The research was performed on the material of the Ukrainian National Linguistic Corprs of the Ukrainian Language and Information Fund of the National Academy of Sciences of Ukraine. The object of analysis is 56 homogeneous word-combinations with different components, different morphological content and non-identical semantic-syntactic relations. Among them, there are 20 two-component units, 20 three-component, and 16 – four-component ones.

Purpose: The purpose of the research is to establish the degree of indivisible word-combinations components sequence non-randomness in Ukrainian by calculating the association measure *mutual information* (MI).

Results: the frequency data for word-combination from the Ukrainian National Linguistic Corpus is received, the MI for multicomponent units is computed, the obtained results are analyzed, the correlation between the MI value and the type of indivisible word-combinations is revealed.

For the analyzed indivisible word-combinations, the MI association measure is in the range of 8.64 (Ukr. *мав працювати*) to 44.63 (Ukr. *майстер спорту міжнародного класу*), that means that the components combination in all these units is non-randomness. Within a single text corpus, the MI value depends on such factors as the absolute construction frequency, the absolute frequency of its components, the number of components and the type of indivisible word-combinations.

Discussion: MI results for indivisible word-combinations are predicted to be lower than other phraseological units due to the smaller number of components in their composition. It is important to find out that the predefined MI value depends on the type of indivisible word-combinations. The next stage of the study involves performing a statistical analysis for other types of phraseological units, in particular, composed prepositional equivalents.

Keywords: association measure, phraseological units, mutual information, indivisible word-combination, statistics, the Ukrainian language.

Vitae

Hanna Sytar is PhD of Philology, Associate Professor, Associate Professor of Department of General and Applied Linguistics and Slavonic Philology at Donetsk National University named after Vasyl' Stus. Her areas of research interests include syntax, semantics, pragmatics, construction grammar, applied linguistics.

Correspondence: h.v.sytar@donnu.edu.ua