

УДК 81'373.7:81'32

**СТАТИСТИЧНИЙ АНАЛІЗ ПРИСЛІВ'ІВ І ПРИКАЗОК:
ПОКАЗНИК АСОЦІАЦІЇ *MUTUAL INFORMATION*
(НА МАТЕРІАЛІ УКРАЇНСЬКОГО НАЦІОНАЛЬНОГО ЛІНГВІСТИЧНОГО КОРПУСУ)¹**

Стаття продовжує цикл публікацій, присвячених статистичному аналізу фразеологічних і фразеологізованих одиниць української мови. У ній з'ясовано ступінь невинності поєднання компонентів у складі українських прислів'їв і приказок за допомогою обчислення показника асоціації mutual information (MI).

Отримані результати обчислень для 53 прислів'їв і приказок, виконаних за даними Українського національного лінгвістичного корпусу, доводять, що всі проаналізовані одиниці мають високий ступінь невинності поєднання слів (MI перебуває в діапазоні від 24,5 до 95,27), що є кількісним підтвердженням стійкості їхнього зв'язку.

Зафіксовано статистично вірогідний зв'язок між кількістю компонентів прислів'я або приказки і величиною показника асоціації MI. Наведені результати загалом узгоджуються із статистичними даними, отриманими на попередніх етапах дослідження для інших типів фразеологічних одиниць – лексичних і синтаксичних фразеологізмів.

Ключові слова: показник асоціації, фразеологічна одиниця, mutual information, прислів'я, приказка, статистика, українська мова.

Постановка проблеми, актуальність дослідження. Сучасна лінгвістика позиціонує корпуснозорієнтованість як необхідну умову мовознавчого дослідження. Статистичний аналіз фразеологічних і фразеологізованих одиниць належить до актуальних завдань лінгвістичної статистики, оскільки за допомогою математичних методів і прийомів покликаний підтвердити або спростувати належність певної мовної одиниці до класу стійких. Процедура такого аналізу на матеріалі синтаксичних фразеологізмів української мови запропоновано у працях (Syтар, “Statystychni Kryteriyi Analizu Syntaktychnykh Frazelohizmiv”; Syтар, “Statystychni analiz frazeolohizovanykh rechen...”; Syтар, “Syntaktychni frazeolohizmy v rozrizi konstruktivnoy hramatyky”).

Обчислення показників асоціації як метод визначення невинності поєднання компонентів може бути застосований для різних типів конструкцій. Цю статтю присвячено аналізу прислів'їв і приказок, які в межах широкого підходу до розуміння обсягу фразеології кваліфікують як один із типів стійких одиниць (В. Л. Архангельський, А. М. Баранов, Д. О. Добровольський, Т. О. Туліна та ін.).

У розумінні прислів'їв і приказок спираємось на усталений в українському мовознавстві погляд, згідно з яким прислів'я визначають як «стійкий вислів переважно фольклорного походження, в якому зафіксований практичний досвід народу та його оцінка різних подій і явищ. Прислів'я на відміну від приказок, – це самостійні судження, граматично та інтонаційно оформлені як прості («Дружній череді і вовк не страшний») або складні («Біда тому воліві, котрого корова коле») речення» (Українська мова: Entsyklopediia: 530, автор статті М. Т. Демський). Оскільки в збірниках прислів'їв і приказок їх подають без розмежування, до статистичного аналізу залучаємо їх разом як групу умовно однорідних одиниць.

Матеріал і методи дослідження. Об'єктом статистичного аналізу стали 53 прислів'я і приказки, дібрані з авторитетних джерел (“Українські прыказки, prysliv'ia i take inshe”; “Prysliv'ia ta pryказky”), серед них 9 трикомпонентних, 17 – чотирикомпонентних, 14 – п'ятикомпонентних, 6 – шестикомпонентних і 7 – семикомпонентних одиниць. Вірогідність одержаних кількісних даних забезпечено виконанням обчислень на матеріалі значного за обсягом й індексованого корпусу текстів – Українського національного лінгвістичного корпусу (далі УНЛК) Українського мовно-інформаційного фонду НАН України.

В арсеналі сучасної статистики існує низка статистичних критеріїв (коефіцієнтів), об'єднаних терміном «показники асоціації» (англ. association measures, measures of association). За Кембриджським словником статистики Брайана Еверітта (Brian S. Everitt), «Показники асоціації – числові індекси, що обчислюють силу статистичної залежності двох або більше квалітативних змінних» (Everitt: 241).

У цьому авторитетному лексикографічному виданні термін «асоціація» витлумачено як «загальний термін, що використовується для опису відношення між двома змінними. Значною мірою є синонімічним до кореляції» (Everitt: 20).

У словнику статистики й методів дослідження Американської асоціації психологів (за ред. Sheldon Zedeck) асоціацію визначено як «ступінь статистичної залежності або відношення між двома або більше явищами», а кореляцію – як «ступінь відношення (зазвичай лінійного) між двома змінними, який може бути обрахований як коефіцієнт кореляції, сила асоціації» (APA: 65).

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

Статистична залежність послідовності словоформ у корпусі визначається за допомогою показника асоціації *mutual information* (далі МІ) (буквально – взаємна, спільна інформація). Поняття МІ ввів у теорію інформації Роберт Маріо Фано (Fano). У лінгвістичних дослідженнях його вперше застосували Кеннет Ворд Чарч (Kenneth Ward Church) і Патрік Генкс (Patrick Hanks) (Church, Hanks). Сутність спільної інформації вчені визначили так: «спільна інформація порівнює ймовірність спостереження x та y разом (поєднана ймовірність) з ймовірностями спостереження x та y незалежно (випадкова)» (Church, Hanks: 23). Відповідно у дослідників мова йшла про невинновідповідність поєднання двох слів у тексті і про потребу залучення цього методу для лексикографії, укладання конкордансів, вивчення сполучуваності слів та ін.

Оскільки обраний об'єкт дослідження – прислів'я та приказки – є багатокомпонентними (три- і більше) одиницями, постає потреба врахувати у формулі МІ більшу кількість компонентів. Тому обчислення здійснено за формулою (1), виведеною у працях (Petrovic, Snajder, Basic, Kolar: 323; Yagunova, Pivovarova: 586).

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

де МІ – коефіцієнт *mutual information*;

i – це кількість компонентів конструкції;

c_1 – перша лексична одиниця;

c_2 – друга лексична одиниця;

c_i – i -а лексична одиниця;

$f(c_1, c_2, \dots, c_i)$ – абсолютна частота вживання конструкції c_1, c_2, \dots, 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 словоформ у корпусі;

\log_2 – логарифм числа за основою 2.

Мета цього дослідження – визначити ступінь невинновідповідності поєднання компонентів у складі українських прислів'їв і приказок за допомогою обчислення показника асоціації МІ. Для досягнення поставленої мети розв'язано такі **завдання**:

1) укладено реєстр прислів'їв і приказок, що охоплює одиниці з різною кількістю компонентів і різну тематику;

2) з УНЛК отримано частотні дані для прислів'їв і приказок;

3) виконано обчислення за формулою МІ для багатокомпонентних одиниць;

4) проаналізовано отримані результати.

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. Оскільки цей корпус текстів є динамічним, зазначимо, що частотні дані подаємо станом на лютий 2018 року. Загальна кількість слововживань у корпусі в період здійснення підрахунків становила 189 200 000 одиниць.

Покажемо приклад здійснених підрахунків. Для обчислення ступеня невинновідповідності поєднання словоформ у межах прислів'я *Терпи, козаче, отаманом будеш* з УНЛК було отримано такі кількісні дані: абсолютна частота прислів'я становить 11, абсолютна частота словоформи *терпи* – 228; *козаче* – 307; *отаманом* – 239; *будеш* – 1968. Підставляючи ці дані до формули (1), отримуємо:

$$MI(\text{Терпи, козаче, отаманом будеш}) = \log_2 \frac{11 \times (189200000)^3}{228 \times 307 \times 239 \times 1968} = 51,007076 \approx 51,01.$$

Коефіцієнт МІ обраховували з точністю до двох знаків після коми. Отримані результати МІ для трикомпонентних прислів'їв і приказок подано в таблиці 1.

Як видно з таблиці 1, коефіцієнт МІ для трикомпонентних прислів'їв і приказок перебуває у межах від 24,5 (*Хліб усьому голова*) до 35,34 (*Горбатого могила виправить*).

Показник асоціації МІ для трикомпонентних прислів'їв і приказок за даними УНЛК

№ з/п	Прислів'я або приказка	Абсолютна частота вживання прислів'я або приказки	Абсолютна частота вживання словоформ-компонентів прислів'я або приказки	Показник асоціації МІ
1	<i>Береженого Бог береже</i>	38	<i>береженого 89;</i> <i>Бог 2 556;</i> <i>береже 569</i>	33,29
2	<i>Гол, як сокол</i>	4	<i>гол 457;</i> <i>як 6 274;</i> <i>сокол 53;</i>	29,81
3	<i>Голий, як кістка</i>	2	<i>голий 920;</i> <i>як 6 274;</i> <i>кістка 362</i>	25,03
4	<i>Горбатого могила виправить</i>	25	<i>горбатого 138;</i> <i>могила 963;</i> <i>виправить 155</i>	35,34
5	<i>Змерз, як собака</i>	2	<i>змерз 262;</i> <i>як 6 274;</i> <i>собака 1 465</i>	24,83
6	<i>На двох стільцях</i>	65	<i>на 6 331;</i> <i>двох 4 584;</i> <i>стільцях 283</i>	28,08
7	<i>Сила соломі ломить</i>	7	<i>сила 3 032;</i> <i>соломі 517;</i> <i>ломить 200</i>	29,57
8	<i>Собака на сні</i>	25	<i>собака 1 465;</i> <i>на 6 331;</i> <i>сні 342</i>	28,07
9	<i>Хліб усьому голова</i>	9	<i>хліб 2 067;</i> <i>усьому 2 111;</i> <i>голова 3 122</i>	24,5

Контрольна величина, починаючи від якої вважаємо зв'язок слів не випадковим, залежить від показників абсолютної частоти конструкції, від абсолютної частоти її окремих складників і від розміру корпусу. Для Українського національного лінгвістичного корпусу, розмір якого в лютому 2018 року становив 189 200 000 слововживань, ця контрольна величина становить 7,56 (детально процедуру виведення контрольної величини викладено у праці (Syta 2017: 310-311)):

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

Відповідно отримані результати можна кваліфікувати як такі, що відбивають високий ступінь не випадковості (зв'язаності) компонентів конструкції, оскільки вони більше ніж утричі перевищують контрольну величину.

Статистичні дані щодо чотири-, п'яти-, шести- й семикомпонентних прислів'їв і приказок наведено в таблицях 2, 3, 4 і 5 відповідно.

Показник асоціації МІ для чотирикомпонентних прислів'їв і приказок за даними УНЛК

№ з/п	Прислів'я або приказка	Абсолютна частота вживання прислів'я або приказки	Абсолютна частота вживання словоформ-компонентів прислів'я або приказки	Показник асоціації МІ
1	<i>Більшому й більше треба²</i>	0	—	—
2	<i>Велике дерево поволі росте</i>	1	<i>велике 3 063; дерево 2 016; поволі 1 629; росте 1 605</i>	38,61
3	<i>Вік живи — вік учись</i>	19	<i>вік 2 382; живи 662; учись 100</i>	48,28
4	<i>Гарна дівка, як маківка</i>	5	<i>гарна 1 697; дівка 613; як 6 274; маківка 125</i>	45,24
5	<i>Гусь свині не товариш</i>	4	<i>гусь 34; свині 750; не 6 308; товариш 1 421</i>	46,75
6	<i>Два українці – три гетьмани</i>	2	<i>два 4 780; українці 1 333; три 4 607; гетьмани 267</i>	40,65
7	<i>Двічі літа не буває</i>	2	<i>двічі 2 140; літа 2 050; не 6 308; буває 2 820</i>	37,34
8	<i>Знай, коза, своє стійло</i>	1	<i>знай 1 118; коза 523; своє 4 234; стійло 66</i>	45,24
9	<i>На двох стільцях сидить</i>	2	<i>на 6 331; двох 4 584; стільцях 283; сидить 2 241</i>	39,42
10	<i>Народ скаже, як зав'яже</i>	4	<i>народ 2 547; скаже 2 038; як 6 274; зав'яже 32</i>	44,56
11	<i>Рідна мова — не полова</i>	1	<i>рідна 1 371; мова 3 101; не 6 308; полова 210</i>	40,13
12	<i>Ситий голодному не вірить</i>	4	<i>ситий 499; голодному 220; не 6 308; вірить 1 337</i>	44,73
13	<i>Ситий голодному не товариш</i>	7	<i>ситий 499; голодному 220; не 6 308; товариш 1 421</i>	45,45
14	<i>Терпи, козаче, отаманом будеш</i>	11	<i>терпи 228; козаче 307; отаманом 239; будеш 1 968</i>	51,01
15	<i>Усі під Богом ходимо</i>	8	<i>усі 4 175; під 5 636; Богом 1 619; ходимо 423</i>	41,61
16	<i>Хвали мене, моя губонько</i>	3	<i>хвали 296; мене 3912; моя 3 214; губонько 6</i>	49,69
17	<i>Язык до Києва доведе</i>	20	<i>язык 1 642; до 6 281; Києва 1 979; доведе 592</i>	43,35

² У випадку можливої, але не зафіксованої в УНЛК конструкції (абсолютна частота 0), частоти окремих компонентів не наводимо через те, що обчислення МІ не має смислу, оскільки логарифму 0 не існує. Обчислення показників асоціації для таких конструкцій не здійснювали, тому у відповідній графі таблиці стоїть знак «—».

Показник асоціації МІ для п'ятикомпонентних прислів'їв і приказок за даними УНЛК

№ з/п	Прислів'я або приказка	Абсолютна частота вживання прислів'я або приказки	Абсолютна частота вживання словоформ-компонентів прислів'я або приказки	Показник асоціації МІ
1	<i>Баба з воза – коням легше</i>	6	<i>баба 1 299; з 6 304; воза 841; коням 330; легше 2 310</i>	60,35
2	<i>Всяка пташка своє гніздо знає</i>	2	<i>всяка 762; пташка 844; своє 4 234; гніздо 953; знає 3 268</i>	58,07
3	<i>Гуртом добре й батька бити</i>	1	<i>гуртом 1 046; добре 4 220; й 6 092; батька 2 547; бити 1 552</i>	53,42
4	<i>Життя прожити — не поле перейти</i>	14	<i>життя 4 809; прожити 965; не 6 308; поле 2 858; перейти 1 968</i>	56,60
5	<i>Козак хороший, та немає грошей</i>	0	–	–
6	<i>Кому весілля, а курці смерть</i>	4	<i>кому 2 995; весілля 1 403; а 6 299; курці 217; смерть 2 986</i>	58,05
7	<i>Розуміється, як вовк на зорях</i>	3	<i>розуміється 854; як 6 274; вовк 1 221; на 6 331; зорях 218</i>	58,56
8	<i>Степ та воля – козацька доля</i>	4	<i>степ 953; та 6 216; воля 1 932; козацька 526; доля 2 727</i>	58,12
9	<i>Тепер життя панам та котам</i>	0	–	–
10	<i>Що село, то й сотник</i>	4	<i>що 6 293; село 2 162; то 5 532; й 6 092; сотник 329</i>	54,92
11	<i>Як дбаєш, так і маєш</i>	6	<i>як 6 274; дбаєш 116; так 6 074; і 6 189; маєш 1 879</i>	57,05
12	<i>Яке зіллячко, таке й сім'ячко</i>	1 ²	<i>яке 4 644; зіллячко 55; таке 4 384; й 6 092; сім'ячко 3</i>	65,76
13	<i>Яке їхало, таке й здибало</i>	7	<i>яке 4 644; їхало 189; таке 4 384; й 6 092; здибало 26</i>	63,67
14	<i>Яке коріння, таке й насіння</i>	11	<i>яке 4644; коріння 1436; таке 4384; й 6 092; насіння 1173</i>	55,91

Таблиця 4

Показник асоціації МІ для шестикомпонентних прислів'їв і приказок за даними УНЛК

№ з/п	Прислів'я або приказка	Абсолютна частота вживання прислів'я або приказки	Абсолютна частота вживання словоформ-компонентів прислів'я або приказки	Показник асоціації МІ
1	<i>Боже помози, а сам не лежи</i>	4	<i>Боже 2 257; помози 430; а 6 299; сам 4 306; не 6 308; лежи 276</i>	74,16
2	<i>Дай серцю волю — заведе в неволю</i>	21	<i>дай 2 231; серцю 894; волю 2 349; заведе 309; в 6 333; неволю 465</i>	79,98
3	<i>Де два українці, там три гетьмани</i>	14	<i>де 5 679; два 4 780; українці 1 333; там 4 605; три 4 607; гетьмани 267</i>	73,81
4	<i>Під лежачий камінь вода не тече</i>	20	<i>під 5 636; лежачий 98; камінь 1 939; вода 2 814; не 6 308; тече 1 286</i>	77,39
5	<i>Що можна лялі, не можна мамі</i>	1	<i>що 6 293; можна 5 520; лялі 41; не 6 308; мамі 757</i>	72,45
6	<i>Як неділя, то й сорочка біла</i>	1	<i>як 6 274; неділя 739; то 5 532; й 6 092; сорочка 1 009; біла 1 808</i>	69,53

Таблиця 5

Показник асоціації МІ для семикомпонентних прислів'їв і приказок за даними УНЛК

№ з/п	Прислів'я або приказка	Абсолютна частота вживання прислів'я або приказки	Абсолютна частота вживання словоформ-компонентів прислів'я або приказки	Показник асоціації МІ
1	<i>Береженого Бог береже, а козака шабля стереже</i>	0	–	–
2	<i>Гол, як сокол, а гострий, як бритва</i>	2	<i>гол 457; як 6 274; сокол 53; а 6 299; гострий 1 425; бритва 225</i>	95,27
3	<i>Коли маєш сто кіп, то будеш ніп</i>	1	<i>коли 5 545; маєш 1 879; сто 2 631; кіп 152; то 5 532; будеш 1 968; ніп 539</i>	90,6
4	<i>Коли убогому жениться, то й ніч мала</i>	1	<i>коли 5 545; убогому 97; жениться 173; то 5 532; й 6 092; ніч 2 985; мала 3 642</i>	90,12
5	<i>Пан з паном, а Іван з Іваном</i>	2	<i>пан 2 056; з 6 304; паном 916; а 6 299; Іван 2 018; Іваном 701</i>	86,83
6	<i>Під лежачий камінь і вода не тече</i>	5	<i>під 5 636; лежачий 98; камінь 1 939; і 6 189; вода 2 814; не 6 308; тече 1 286</i>	90,29
7	<i>Що вільно панові, те не можна Іванові</i>	0	–	–

Дані, наведені в таблицях 2–5 дають змогу констатувати, що коефіцієнт МІ для чотирикомпонентних прислів'їв і приказок перебуває в межах від 37,34 (*Двічі літа не буває*) до 51,01 (*Терпи, козаче, отаманом будеш*), тобто в 5–6 разів більший за 7,56; для п'ятикомпонентних – від 53,42 (*Гуртом добре й батька бити*) до 65,76 (*Яке зіллячко, таке й сім'ячко*), тобто в 7–8 разів більший за контрольну величину; для шестикомпонентних – від 69,53 (*Як неділя, то й сорочка біла*) до 79,98 (*Дай серцю волю – заведе в неволю*), тобто в 9–10 разів вищий від контрольної величини; для семикомпонентних – від 86,83 (*Пан з паном, а Іван з Іваном*) до 95,27 (*Гол, як сокол, а гострий, як бритва*), тобто більший в 11–12 разів за контрольну величину.

Висновки. Отримані результати обчислень для 53 прислів'їв і приказок, виконаних за даними Українського національного лінгвістичного корпусу, доводять, що всі проаналізовані одиниці мають високий ступінь не випадковості поєднання словоформ: коефіцієнт МІ перебуває в діапазоні від 24,5 до 95,27 (тобто є

втричі – удванадцятьоро більшим, ніж контрольна величина), що є кількісним підтвердженням стійкості зв'язку словоформ у складі відповідних одиниць.

Зафіксовано статистично вірогідний зв'язок між кількістю компонентів прислів'я/приказки й величиною показника асоціації MI. Так, для трикомпонентних одиниць результат MI становить від 24,5 (*Хліб усьому голова*) до 35,34 (*Горбатого могила виправить*); для чотирикомпонентних – від 37,34 (*Двічі літа не буває*) до 51,01 (*Терпи, козаче, отаманом будеш*); для п'ятикомпонентних від 53,42 (*Гуртом добре й батька бити*) до 65,76 (*Яке зіллячко, таке й сім'ячко*); шестикомпонентних – від 69,53 (*Як неділя, то й сорочка біла*) до 79,98 (*Дай серцю волю – заведе в неволю*); семикомпонентних – від 86,83 (*Пан з паном, а Іван з Іваном*) до 95,27 (*Гол, як сокол, а гострий, як бритва*).

Серед нерозв'язаних на сьогодні проблем статистичного аналізу варто відзначити омонімію, зокрема, потребу залучення людини-експерта для розмежування випадків типу *гол* (пор. *перший гол* і *гол, як сокол*) або *на двох стільцях* (пор. стійка сполука *сидіти на двох стільцях* і вільний (нефразеологізований) вияв у реченні *Посеред кімнати на двох стільцях стоїть маленька з сірої бляхи ванночка* (В. Винниченко. Записки Кирпатого Мефістофеля).

Наведені результати загалом узгоджуються із статистичними даними, отриманими на попередніх етапах дослідження для інших типів фразеологічних одиниць – лексичних і синтаксичних фразеологізмів. Водночас вони засвідчують вищий ступінь не випадковості поєднання словоформ саме для прислів'їв і приказок, що є, очевидно, наслідком їхньої багатоконпонентності, стійкості й цілісності їхнього сприйняття носіями української мови.

Перспективним вважаємо статистичний аналіз інших типів стійких одиниць і зіставлення відповідних даних з результатами, отриманими для синтаксичних і лексичних фразеологізмів, прислів'їв і приказок.

References

APA Dictionary of Statistics and Research Methods. Sheldon Zedeck, PhD, editor in chief. Washington, DC: American Psychological Association, 2014. Print.

Church, Kenneth Ward, and Patrick Hanks. "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.

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 Frazеологізмів (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 frazeологізовanykh 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 frazeологізми v rozrizi konstruktivnoi 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 konstruktivnyam (From Collocations to Constructions)". *ACTA LINGUISTICA PETROPOLITANA. Works of the Institute of Linguistic Researches of RAS, Russkiy yazyk: grammatika konstruktivnykh i leksiko-semanticheskie podkhody (The Russian Language: Construction Grammar and Lexical and Semantic Approaches)*: X, part 2. (2014) 568–617. Print.

List of Sources

Prysliv'ia ta prykazky (Proverbs and Sayings). Ukl. M. Paziak. Red. Myshanych S. V.; red. Berezovskiy I. P.; red. Hordiichuk M. M.; red. Zubkov S. D.; red. Sushko L. D.; red. Kuz V. P. Kyiv: Nauk. dumka, 1991. Print.

Ukrainski prykazky, prysliv'ia i take inshe (Ukrainian Sayings, Proverbs etc). Ukl. M. Nomys. Kyiv: Lybid, 1993. Print.

Надійшла до редакції 20 березня 2018 року.

STATISTICAL ANALYSIS OF PROVERBS AND SAYINGS: ASSOCIATION MEASURE OF MUTUAL INFORMATION (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 author examines the statistical analysis of proverbs and sayings on the material of the Ukrainian National Linguistic Corpus of Ukrainian Lingua-Information Fund, NAS of Ukraine. Corpus-oriented statistical research of Ukrainian proverbs and sayings has not been carried out yet. The object of the analysis is 53 proverbs and sayings, selected from authoritative sources: *Ukrainski prykazky, pryslivia i take inshe* (Ukrainian proverbs, sayings etc) / Ukl. M.Nomys. Kyiv: Lybid, 1993; *Pryslivia ta prykazky* (Proverbs and sayings) / ukl. M. Paziak Kyiv: Nauk. dumka, 1991. Among them, there are 9 three-component, 17 – four-component, 14 – five-component, 6 – six-component and 7 – seven-component units.

Purpose: The purpose of this study is to determine the degree of non-randomness of the components combination in Ukrainian proverbs and sayings by means of the calculating the association measure of *mutual information* (hereinafter MI).

Results: The obtained results of the calculations for 52 proverbs and sayings, done with the help of the Ukrainian National Linguistic Corpus data, prove that all analysed units have high degree of non-randomness combination of word forms (MI is in the range from 24.5 to 95.27), which is a quantitative confirmation of their connection stability. The reference value, from which the connection of words is considered to be non-random, depends not only on the indicators of the absolute frequency of the construction and its individual components but also on the size of the corpus. For the Ukrainian National Linguistic Corpus, which consisted of 189 200 000 words in February 2018, the reference value is 7.56.

There is a statistically probable connection between the number of components in the proverb or saying and the size of association measure *mutual information*. Thus, for the three-component units the result of the MI is from 24.5 (*Хліб усьому голова / Khlib usomu holova*) to 35.34 (*Горбатого могила виправить / Horbatoho mohyla vypravyt*); for four-component units – from 37.34 (*Двічі літа не буває / Dvichi lita ne buvaie*) to 51.01 (*Терпи, козаче, отаманом будеш / Terpu, kozache, otamanom budesh*); for five-component from 53.42 (*Гуртом добре й батька бити / Hurtom dobre y batka byty*) to 65.76 (*Яке зіллячко, таке й сім'ячко / Yake zilliachko, take y simiachko*); for six-component – from 69,53 (*Як неділя, то й сорочка біла / Yak nedilia, to y sorochka bila*) to 79.98 (*Дай серцю волю – заведе в неволю / Dai sertsiu voliu – zavede v nevoliu*); for seven-component – from 86.83 (*Пан з паном, а Іван з Іваном / Pan z panom, a Ivan z Ivanom*) to 95.27 (*Гол, як сокол, а гострий, як бритава / Hol, yak sokol, a hostryi, yak brytva*).

Discussion: The presented results are broadly consistent with the statistical data obtained in the previous stages of the study of other types of phraseological units – lexical and syntactic idioms. At the same time, they prove the higher degree of non-randomness of the combination of word forms for proverbs and sayings, which, obviously, are the consequence of their multicomponent, as well as the stability and integrity of their perception by the speakers of the Ukrainian language. It is considered perspective to involve other types of phraseological units to the statistical analysis.

Keywords: association measure, phraseological units, mutual information, proverb, saying, 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