# An Analysis of Association Measures in Collocation Extraction from a Pedagogical Perspective

KOHEI FUKUDA

## 1. Introduction

Many researchers point out the pervasiveness of collocations in a text and the importance of learning collocations in second language acquisition. Hill (2000: 53) argues that

> Collocation is important because this area of predictability is, as we have seen, enormous. Two, three, four, and even five-word collocations make up a huge percentage of all naturally-occurring text, spoken or written. Estimates vary, but it is possible that up to 70% of everything we say, hear, read, or write is to be found in some form of fixed expression.

A vast number of language texts are composed of collocation; therefore, collocation learning is essential for proficient use of language.

Sinclair (1991) proposed two models of how words occur in a language text: *the open-choice principle* and *the idiom-principle*. The open-choice principle sees language texts as the result of a large number of choices where the only restraint is grammaticalness. This model regards each slot in which an individual word is put as an "open slot." For instance, in a grammatical structure of a transitive verb followed by its object, as in *cause accidents* and *make a cake,* virtually any word can occur in the first slot and the second slot as long as the phrase is grammatically well-formed. On the other hand, according to the idiom principle, there are many more constraints and limitations in the choice of words in a text. The choice of one word determines, to some

extent, the choice of another word. For example, the transitive verb *cause* is usually followed as its object by something bad or unpleasant. Therefore, "*cause accidents*" sounds natural, while "*cause victory*" sounds unnatural. Nation (2013: 483) argues that collocation is an important learning goal because a large amount of language is based on the idiom principle.

Despite its pervasiveness in language texts and its importance in the acquisition of a second language, collocation seems to be often ignored by language teachers and learners in Japan. According to a survey conducted by Kawamura & Ishii (2013), no more than 1.6% of Japanese university students comprehended the concept of collocation, which suggests that few learners in Japan have paid attention to collocation in studying English when they were in junior high school or high school. This is partly because the current Course of Study in Japan does not clearly set collocational competency as a learning objective.

Unlike individual words, it is difficult to choose collocation as an aim of study because teachers themselves simply do not know exactly what collocations are and how they should teach them. Furthermore, since individual words produce a massive number of word combinations in principle, there can be too many collocations for learners to acquire. In order to solve the problem inherent in collocation learning, it is necessary to identify a set of collocations which should be acquired by Japanese learners of English. Koya (2012) argues that it is essential to make a "basic collocations list," which can contribute to clarifying collocation to acquire for learners, and popularizing collocation learning in English education in Japan. Furthermore, L2 collocations can be learned both by incidental and intentional learning, but intentional learning results in bigger and faster gains (Szudarski, 2017: 212). Given that most Japanese people learn English in EFL environment, where they are not exposed to enough input to incidentally learn collocational competence, intentional learning of collocations are of even greater importance, and for that purpose, a collocations list is necessary.

What is essential for creating a collocations list is a set of criteria for identifying collocations for pedagogical purposes. Criteria for identifying collocations are generally divided into two types; the frequency-based view and the phraseological view (Henriksen, 2013). The frequency-based view is an attempt to identify collocations on the basis of statistical measures which assess collocability, generally known as association measures (AMs), using large corpora. The phraseological view employs linguistical classification criteria, such as the degree of semantic opacity, collocational structure, and substitutability of word elements. Granger & Paquot (2008) claimed that researchers should utilize both of these two views in a well-balanced manner in identifying collocation. Therefore, when making a collocations list for pedagogical purposes, researchers should extract collocations using AMs first, and then screen these collocations using the phraseological view from an educational perspective. Selecting collocations based on the phraseological view, however, is a subjective process and requires enormous effort. Thus, making good and efficient uses of AMs in extracting collocations should enhance the reliability and efficiency of the process of choosing an appropriate set of collocations.

Few attempts have been made to explore how to employ AMs in extracting collocations from large corpora from pedagogical perspectives. The main objective of this paper is to explore how AMs of collocations can be used to extract collocations for the purpose of creating a "collocations list" for Japanese learners of English.

## 2.　Review of Related Literature

This section introduces some previous studies on collocational competence of L2 learners and gives an explanation as to why L2 language users need to learn collocation, and what factors have an influence on L2 learners' collocational competence. Furthermore, this section explains widely used AMs of collocations which this study will deal with.

## 2.1 L2 language users' need for collocation competence

Many researchers have investigated the acquisition of collocations by L2 learners, and pointed out the importance of collocation competence for language production and reception. Collocational proficiency enables L2 users to make use of fixed phrases, and therefore strike interlocuters or readers as native-like (Henriksen, 2013). O'Keefee et al. (2007) argues that the use of fixed expressions alleviates the burden on users' cognitive ability when processing language and allows language users to direct cognitive energy into more creative aspect of language use, such as discourse organization and successful interaction. In sum, collocational competence assists L2 learners in communicating in a more natural and creative manner.

## 2.2 L2 learners' collocational competence

Many studies point out the influence of L1 on L2 collocational competence. Nesselhauf (2003) investigated the use of verb + noun collocation by German learners of English, and suggests that a learners' L1 has an influence on the use of collocations. She drew from her study a conclusion that an explicit instruction of collocations is necessary to enhance learners' proficiency. Granger (1998) explored French L2 learners' use of intensive adverb + adjective collocations (e.g., *completely different*). She found that advanced learners overused certain collocations which were equivalent to their mother tongue. She argued that although learners' unnatural-sounding production of language is generally associated with their lack of prefabricated expressions, it can also be due to their overdependence on certain expressions. Kurosaki (2010) studied the use of verb + noun collocations by Japanese university students, and showed that L1 has an effect on the collocational proficiency of Japanese learners.

Koya (2005) explores the process of the acquisition of verb + noun collocations by Japanese learners of English. She suggests that (1) learners' general vocabulary knowledge correlates with collocational knowledge; (2) knowledge of receptive collocational knowledge is deeper than productive collocational knowledge; (3) productive collo-

cational knowledge is influenced by vocabulary knowledge, semantic opacity, delexicalized verbs, core meaning of nouns and verbs, collocational structure, and L1 equivalence; (4) receptive collocational knowledge is affected by L1 equivalence, delexicalized verbs and core meaning of verbs. She claims that learners at every level should pay attention to collocation and that educators should teach collocation differently to learners at different proficiency levels.

## 2.3   Association measures

Many statistical methods of measuring collocational strength have been developed. Ishikawa (2008) introduces, as widely used AMs, raw frequency, Dice coefficient, t-score, mutual information (MI), Log-likelihood (LL), z-score, and MI3. Raw frequency refers to the number of times when a certain collocation occurs in a corpus.

T-score of 2 or higher is usually considered a statistically significant combination of words, or collocation. MI is used as a measure which shows to what extent a word has information about another word. MI score of 3 or higher can be interpreted as evidence that the combination of the two words is collocation (Hunston, 2002).

McEnery et al. (2006) state that the most commonly used statistical test is the chi-square test and another commonly used statistical test is Log-likelihood (LL). The chi-square test ($\chi^2$) makes a comparison between the observed values and the expected values. LL also compares the observed values and the expected values.

LL is generally preferred, compared to chi-square because (1) it does not presuppose the minimum expected frequencies, (2) it does not overestimate rare cases, and (3) it is not influenced by corpus size (Leech et. al, 2001).

Whereas MI score puts too much emphasis on rare words, MI3 pays more attention to frequent words. Thus, collocations extracted using MI3 are more useful for language learners at the beginning and intermediate level, while those extracted using MI are interesting for a lexicographic purpose (McEnery et al, 2006).

Ishikawa (2008) classifies five measures (raw frequency, t-score, LL,

Dice coefficient, MI) into three categories. Raw frequency, t-score, LL are grouped into frequency-based measures, which put emphasis on high frequent collocations. MI is categorized into non-frequency-based measures. As MI puts weight on low-frequency words which mostly occur with a certain word, MI tends to extract low-frequency collocations. Dice coefficient lies between the two groups of AMs. According to Evert (2008), LL is the best measure in terms of mathematical statistics. T-score is not based on mathematical reasoning, but still it can be useful as a heuristic measure for collocation identification. He argued that it is important to explore what kind of collocations are extracted by different measures.

## 2.4   Research Questions

Some researchers utilize AMs so as to create collocations lists for pedagogical purposes. Ackermann & Chen (2013), in an attempt to make a collocations list for EAP (English for academic purpose), set the criteria for extracting collocations from corpora as follows: (a) raw frequency $\geq 1$ per million; (b) raw frequency $\geq 0.2$  per million in each sub-corpora; (c) MI score $\geq 3$; (d) t-score $\geq 4$. Koya (2015) explores how to select basic collocations for Japanese learners of English to acquire. In the study, she took the noun *time* as an example, and created a collocations list of "verb + *time*". In extracting "verb + *time*" collocations from corpora, she employed t-score, z-score, MI, and Log-likelihood.

As is seen in the study stated above, AMs help create collocations lists for educational purposes. However, which association measures should be best applied to selecting collocations for pedagogical purposes are yet to be explored. The consideration of how effectively each association measure extracts pedagogically useful colocations and the comparison between these measures are an essential process of investigating the usefulness and suitability of these measures. Therefore, the research questions of this paper were formulated as follows:

RQ1:   Which association measure can extract collocations of ped-

agogical use more effectively?

RQ2:    How should each association measure be combined to obtain collocations depending on different proficiency levels of learners?

## 3.   Method

In order for language policymakers and practitioners to incorporate collocation learning into a classroom, a collocations list is necessary, as is a word list for vocabulary learning. Extracting collocations from corpora is an essential process of creating a collocations list, and AMs should be employed so as to obtain collocations in an efficient way and on an objective scale. Therefore, the main purpose of this paper is to explore how AMs, such as Dice coefficient, Log-likelihood, t-score, z-score, MI, and MI3, should be used to extract collocations from corpora with pedagogical applications in mind.

### 3.1   Materials and corpora used in the study

In this study, it was hypothesized that collocations which appear in a published study book for collocations are presupposed to be those of high pedagogical value. Although there are several study books for collocations in Japan and the rest of the world, one of the most widely used is *English Collocations in Use -intermediate (second edition)* (McCarthy & O'Dell, 2017). The present study used collocations found in this book in order to explore the validity of association measures.

The study book is organized into 60 two-page units. Collocations are presented in typical contexts, and each unit focuses on a certain topic, such as weather, music, sport, business, money, time, talking about success and failure, and so forth so that you can learn collocations in a meaningful context. The right-hand page provides a series of exercises so you can check that you have understood the collocations you've studied on the left-hand page.

McCarthy & O'Dell (2017) pay attention mainly to two things when selecting collocations which would be most useful for learners to

study. The first thing is that the authors of the book put emphasis on the collocations which many users of English are likely to use in their speech or writing. "So, in the unit on Eating and drinking we include, for example, *have a quick snack* and *processed food* but not *cocoa butter*, which is a very strong collocation, but one which has very limited use for most people" (McCarthy & O'Dell, 2017: 4). Second, the authors carefully selected semantically opaque collocations, which learners of English might have difficulty in decoding, based on the analysis of the Cambridge Learner Corpus.

In this study, information about AMs was obtained by using the British National Corpus (BNC). The present study made use of BNC-web (Sebastian & Evert) because this interface enables users to extract collocations automatically by seven measures of association, including raw frequency, Dice coefficient, Log-likelihood, t-score, z-score, MI, MI3. Therefore, BNCweb is suited for the current study.

## 3.2   Corpus processing and data analysis

This study focuses on verb + noun collocations. Verb + noun collocations were listed up from McCarthy & O'Dell (2017), and 637 collocations were identified. Verbs were considered to be a node word, and nouns were viewed as its collocate. If verbs appeared in more than nine collocations, they were selected as target verbs. Those collocations in which the target verbs were used were chosen as an object of investigation. As a result, the following eleven verbs were selected: *get, do, have, give, take, make, keep, win, raise, change,* and *cause.* In total, 250 collocations were identified for these verbs from McCarthy & O'Dell (2017), which was called the "target collocations" in this study (see Table 1).

After selecting verbs as node words to investigate in this study, collocates for those node words were extracted from the BNC by employing the six association measures (AMs), Dice coefficient, Log-Likelihood, t-score, z-score, MI3, and MI. The collocation search span was set to +4 (within the four words in the right context), and collocations were extracted and listed as lemmas. As for each node word, the top

Table 1.   The list of "verb + noun" target collocations investigated in this study

| VERB + NOUN | Number | Example |
|---|---|---|
| do + NOUN | 35 | do activities, do aerobics, do an assignment |
| make + NOUN | 45 | make a breakthrough, make an allegation |
| have + NOUN | 40 | have a think, have a break, have a conversation |
| give + NOUN | 25 | give credit, give the impression, give a laugh |
| cause + NOUN | 11 | cause damage, cause concern, cause pain |
| change + NOUN | 11 | change doctors, change jobs, change the subject |
| win + NOUN | 11 | win respect, win case, win praise |
| get + NOUN | 10 | get a job, get a place, get the impression |
| keep + NOUN | 10 | keep the pace, keep a record, keep secrets |
| raise + NOUN | 9 | raise a question, raise money, raise taxes |

100 collocations were extracted from the BNC by using association measures.

To explore the validity of respective AMs for pedagogical purposes, the author investigated the extent to which the top 100 noun collocates from the BNC matched the target collocations selected from McCarthy & O'Dell (2017). Moreover, the author examined the rank of the target collocations by each association measure in question. Different measures returned different results and values, and therefore the direct comparison across the AMs was not possible as they were. However, the rank order of the target collocations by each AM made it possible to compare the different AMs from one another.

## 4.   Results

The results first show how many of the target collocations were covered by each AM. Second, a comparison was made across the AMs, and it is explored how different measures evaluated collocations, and how they were classified in terms of similarities.

## 4.1 The coverage of the target collocations by the measures and their rank order

Table 2 shows the coverage of the target collocations by the top 100 collocations extracted by using each AM in question. The results show that on average, the percentage of the target collocations extracted using each association measure were as follows: Dice coefficient: 55%, Log-Likelihood: 56%, t-score: 55%, z-score: 53%, MI3: 59%, and MI: 26%. Approximately 90% of "*cause* + NOUN" collocations were covered, which was extremely high compared to the other; however, only 30% of "*do* + NOUN" collocations were covered on average. Overall, the coverage of MI was much lower than the other five measures.

Table 3 shows the rank orders of the target collocations according to the top 100 collocations list extracted from the BNC by employing

Table 2.　The coverage of the target collocations

| Node | Total | D | LL | T | Z | MI3 | MI | D(%) | LL (%) | T (%) | Z (%) | MI3 (%) | MI (%) |
|------|-------|----|----|----|----|-----|----|------|--------|-------|-------|---------|--------|
| make | 45 | 29 | 31 | 29 | 28 | 32 | 9 | 64% | 69% | 64% | 62% | 71% | 20% |
| have | 43 | 19 | 21 | 20 | 21 | 21 | 10 | 44% | 49% | 47% | 49% | 49% | 23% |
| take | 40 | 21 | 23 | 19 | 22 | 23 | 7 | 53% | 58% | 48% | 55% | 58% | 18% |
| Do | 35 | 12 | 12 | 12 | 12 | 14 | 6 | 34% | 34% | 34% | 34% | 40% | 17% |
| give | 25 | 13 | 14 | 14 | 13 | 14 | 6 | 52% | 56% | 56% | 52% | 56% | 24% |
| cause | 11 | 10 | 10 | 10 | 10 | 10 | 8 | 91% | 91% | 91% | 91% | 91% | 73% |
| change | 11 | 5 | 5 | 7 | 4 | 5 | 3 | 45% | 45% | 64% | 36% | 45% | 27% |
| Win | 11 | 11 | 9 | 9 | 9 | 10 | 7 | 100% | 82% | 82% | 82% | 91% | 64% |
| Get | 10 | 3 | 2 | 3 | 2 | 3 | 0 | 30% | 20% | 30% | 20% | 30% | 0% |
| keep | 10 | 7 | 7 | 7 | 7 | 7 | 6 | 70% | 70% | 70% | 70% | 70% | 60% |
| raise | 9 | 8 | 7 | 8 | 4 | 8 | 4 | 89% | 78% | 89% | 44% | 89% | 44% |
| total/average | 250 | 138 | 141 | 138 | 132 | 147 | 66 | 55% | 56% | 55% | 53% | 59% | 26% |

Notes: D = Dice coefficient, LL = Log-likelihood, T = t-score, Z = z-score,
　　　 MI = mutual information

Table 3.   The rank orders of the target MAKE + NOUN collocations

| Collocate | D | LL | T | Z | MI3 | MI |
|---|---|---|---|---|---|---|
| decision | 2 | 3 | 2 | 3 | 3 | 75 |
| mistake | 5 | 2 | 6 | 1 | 2 | 15 |
| way | 6 | 20 | 4 | 46 | 11 | - |
| point | 7 | 14 | 7 | 27 | 13 | - |
| contribution | 8 | 6 | 8 | 5 | 6 | 41 |
| effort | 9 | 7 | 9 | 11 | 7 | 97 |
| money | 11 | 23 | 10 | 40 | 16 | - |
| progress | 12 | 8 | 16 | 9 | 9 | 56 |
| change | 14 | 29 | 15 | 57 | 23 | - |
| profit | 15 | 12 | 19 | 16 | 14 | - |
| choice | 20 | 18 | 24 | 24 | 21 | - |
| arrangement | 24 | 16 | 25 | 20 | 18 | - |
| note | 26 | 25 | 28 | 30 | 25 | - |
| comment | 28 | 19 | 29 | 21 | 22 | - |
| impact | 33 | 26 | 38 | 29 | 27 | - |
| friend | 36 | 89 | 35 | - | 66 | - |
| start | 37 | 32 | 40 | 39 | 34 | - |
| sound | 40 | 54 | 42 | 68 | 48 | - |
| time | 41 | - | 20 | - | 76 | - |
| reference | 43 | 50 | 45 | 67 | 50 | - |
| speech | 45 | 48 | 48 | 60 | 52 | - |
| film | 49 | 70 | 54 | 88 | 67 | - |
| demand | 51 | 74 | 55 | 99 | 69 | - |
| appointment | 56 | 49 | 63 | 56 | 55 | - |
| assumption | 60 | 56 | 66 | 61 | 59 | - |
| comparison | 64 | 58 | 73 | 62 | 63 | - |
| case | 68 | - | 59 | - | - | - |

| list | 74 | - | 74 | - | 90 | - |
| loss | 82 | - | 80 | - | 94 | - |
| adjustment | - | 62 | - | 53 | 70 | 88 |
| allegation | - | - | - | - | - | - |
| breakthrough | - | - | - | - | - | 89 |
| commitment | - | 99 | - | - | 99 | - |
| detour | - | 84 | - | 51 | - | 20 |
| excuse | - | 64 | - | 55 | 72 | 93 |
| headline | - | 95 | - | 92 | - | - |
| improvement | - | - | - | - | - | - |
| modification | - | - | - | - | - | - |
| observation | - | 88 | - | - | 100 | - |
| photocopy | - | - | - | - | - | - |
| preparation | - | - | - | - | - | - |
| recording | - | - | - | - | - | - |
| redundant | - | - | - | - | - | - |
| reservation | - | - | - | - | - | - |
| withdrawal | - | - | - | - | - | - |

the AMs. Numbers indicates the rank orders in each measure, and the unmarked cells (shown by "-") shows that the target collocations were not found in the top 100 lists. It is obvious that MI covers a very small numbers of the target collocations, and thus it seems that MI is not suitable for identifying pedagogically useful collocations.

## 4.2   A comparison between the AMs

In Table 3, it seems that the five measures except MI produced apparently similar results. Thus, it is necessary to further explore the differences between the five measures. To this end, the target collocations were compared in terms of the coverage across different AMs. Out of 250, 99 target collocations were covered in all the AMs except

MI. Ten collocations were covered by four measures, twenty-five collocations by three measures, fourteen collocations by two, eight collocations by only one measure.

MI3 covered all the target collocations covered by four or three measures (see Table 4 and Table 5). In table 4, no more than three of the target collocations covered by the four measures were extracted by z-score. As is shown in Table 5 and Table 6, Dice coefficient and t-score produced similar results, and Log-Likelihood and z-sore assessed collocations in a similar way.

This result suggests that Dice coefficient and t-score can be grouped together in terms of collocation selection behavior, and Log-likelihood and z-score can be classified into another group of association measures. To explore the difference between Dice and t-score versus LL and s-score, it can be useful to examine the collocations which were extracted by Dice and t-score, but not by LL and z-score and which are covered by LL and z-score, but not by Dice and t-score in terms of word level of collocates.

The pedagogical importance should be assessed by the frequency of

Table 4.  The target collocations covered by four out of the five measures except MI, and their rank orders

| node | collocate | D | LL | T | Z | MI3 |
|------|-----------|-----|-----|-----|-----|-----|
| make | friend | 36 | 89 | 35 | - | 66 |
| change | place | 49 | 74 | 29 | - | 42 |
| raise | subject | 63 | 85 | 48 | - | 66 |
| raise | capital | 63 | 85 | 48 | - | 66 |
| raise | child | 66 | 92 | 29 | - | 52 |
| give | performance | 68 | 94 | 73 | - | 84 |
| win | praise | 71 | 57 | - | 59 | 64 |
| give | talk | 76 | 85 | 87 | - | 90 |
| take | pleasure | 99 | 70 | - | 89 | 81 |
| give | sigh | - | 47 | 29 | 52 | 45 |

Table 5.   The target collocations covered by three out of the
five measures except MI, and their rank orders

| node | collocate | D | LL | T | Z | MI3 |
|------|-----------|----|----|----|----|-----|
| have | child | 13 | - | 11 | - | 23 |
| do | course | 16 | - | 16 | - | 22 |
| have | word | 30 | - | 30 | - | 45 |
| do | research | 37 | - | 36 | - | 46 |
| get | place | 40 | - | 37 | - | 64 |
| take | course | 41 | - | 34 | - | 70 |
| make | time | 41 | - | 20 | - | 76 |
| have | view | 52 | - | 54 | - | 62 |
| do | duty | 63 | - | 67 | - | 50 |
| make | list | 74 | - | 74 | - | 90 |
| do | hair | 77 | - | 81 | - | 66 |
| make | loss | 82 | - | 80 | - | 94 |
| raise | family | 86 | - | 46 | - | 75 |
| have | game | 92 | - | 94 | - | 95 |
| take | clothes | 97 | 85 | - | - | 91 |
| win | case | 98 | - | 37 | - | 88 |
| do | washing | - | 14 | - | 17 | 33 |
| do | cooking | - | 20 | - | 26 | 53 |
| give | go-ahead | - | 32 | - | 7 | 26 |
| keep | temper | - | 54 | - | 47 | 78 |
| make | adjustment | - | 62 | - | 53 | 70 |
| make | excuse | - | 64 | - | 55 | 72 |
| have | ability | - | 77 | - | 86 | 84 |
| take | prisoner | - | 86 | - | 96 | 98 |
| have | option | - | 87 | - | 96 | 96 |

Table 6.    The target collocations covered by two measures out
of the five measures except MI, and their rank order

| node | collocate | D | LL | T | Z | MI3 |
|------|-----------|---|----|---|---|-----|
| make | case | 68 | - | 59 | - | - |
| keep | word | 90 | - | 62 | - | - |
| give | word | 93 | - | 79 | - | - |
| do | ironing | - | 24 | - | 24 | - |
| have | chat | - | 35 | - | 25 | - |
| have | tendency | - | 52 | - | 58 | - |
| have | sympathy | - | 54 | - | 56 | - |
| do | exam | - | 56 | - | 70 | - |
| make | detour | - | 84 | - | 51 | - |
| make | observation | - | 88 | - | - | 100 |
| take | trip | - | 89 | - | - | 97 |
| make | headline | - | 95 | - | 92 | - |
| take | photo | - | 96 | - | 99 | - |
| make | commitment | - | 99 | - | - | 99 |

collocates because the more frequent words are, the more important
they are for learners. Therefore, it is meaningful to examine the fre-
quency of collocates listed in Table 7 and Table 8. The frequency of
collocates is based on the BNC. It is also valuable to investigate the
level of collocates on the basis of CEFR-J Wordlist, because the
wordlist is created in order to display the levels of words from an edu-
cational perspective (Tono, 2013). Referring to the wordlist enables
researchers and educators to know objectively how useful individual
words are for Japanese learners of English.

As is indicated by Table 7 and Table 8, the average frequency of
collocates that make up the target collocations which are extracted by
Dice and t-score is much higher than that by LL and z-score. On top
of that, most of the collocates by Dice and t-score fall into A1, while
more than half of the collocates by LL and z-score are B1 or on a

Table 7. The target collocations covered by Dice and t-score, but not by LL and z-score, their rank orders, and the frequency and CEFR-J level of their collocates.

| node | collocate | D | T | freq | CEFR-J |
|------|-----------|-----|-----|--------|--------|
| have | child | 13 | 11 | **69271** | **A1** |
| do | course | 16 | 16 | **56036** | **A1** |
| have | word | 30 | 30 | **42301** | **A1** |
| do | research | 37 | 36 | **25531** | **A2** |
| get | place | 40 | 37 | **52469** | **A1** |
| take | course | 41 | 34 | **56036** | **A1** |
| make | time | 41 | 20 | **180243** | **A1** |
| have | view | 52 | 54 | **30686** | **A2** |
| do | duty | 63 | 67 | **11648** | **B1** |
| make | case | 68 | 59 | **63148** | **A1** |
| make | list | 74 | 74 | **13661** | **A1** |
| do | hair | 77 | 81 | **14100** | **A1** |
| make | loss | 82 | 80 | **15261** | **B1** |
| raise | family | 86 | 46 | **41889** | **A1** |
| keep | word | 90 | 62 | **42301** | **A1** |
| have | game | 92 | 94 | **20601** | **A1** |
| give | word | 93 | 79 | **42301** | **A1** |
| win | case | 98 | 37 | **63148** | **A1** |

NOTES: freq = frequencies of the collocates in the BNC

higher level.

## 4.3　Summary

It is clear that the coverage of MI was much more restricted than that of the other five measures, and therefore MI did not seem to be suitable for selecting collocations for pedagogical purposes. In terms of coverage, the five collocational measures except MI yielded a simi-

Table 8.   The target collocations covered by LL and z-score, but not by Dice and t-score, their rank orders, and the frequency and CEFR-J level of their collocates.

| node | collocate | LL | Z | freq | CEFR-J |
|------|-----------|----|----|------|--------|
| do | washing | 14 | 17 | **1504** | - |
| do | cooking | 20 | 26 | **1540** | A2 |
| do | ironing | 24 | 24 | **178** | B1 |
| give | go-ahead | 32 | 7 | **271** | - |
| have | chat | 35 | 25 | **944** | B1 |
| have | tendency | 52 | 58 | **3582** | B1 |
| have | sympathy | 54 | 56 | **2304** | B1 |
| keep | temper | 54 | 47 | **1264** | B1 |
| do | exam | 56 | 70 | **1584** | A2 |
| make | adjustment | 62 | 53 | **2109** | B2 |
| make | excuse | 64 | 55 | **2190** | A1 |
| have | ability | 77 | 86 | **10378** | A2 |
| make | detour | 84 | 51 | **238** | - |
| take | prisoner | 86 | 96 | **4507** | B1 |
| have | option | 87 | 96 | **9138** | B1 |
| make | headline | 95 | 92 | **1378** | B1 |
| take | photo | 96 | 99 | **2011** | A1 |

lar result, but closer scrutiny revealed that they were classified into three groups. Dice coefficient and t-score seem to make one group, which tend to represent relatively frequent collocations, and Log-likelihood and z-score can form another group, which ranks highly the collocations whose collocates are intermediate level words. MI3 lies between these two groups. MI3 succeeded in extracting many of the target collocations which were covered either by Dice coefficient and t-score, or by Log likelihood and z-score.

## 5.  Discussion

This section summarizes the major findings of this study, and discusses the findings in light of previous studies. Moreover, it will address the question of how the AMs should be used for extracting collocations for educational purposes.

The results of the present study show that MI-score covers much less of the target collocations than the other five association measures. The findings clearly indicate that MI-score is not suitable for selecting collocations from a pedagogical perspective. This result corresponds to McEnery et al. (2006).

Z-score extracted only three of the collocations which were retrieved by the other four AMs except MI, and MI3 extracted all the collocations. In other words, z-score could not extract the collocations which the other measures rated highly. In addition to that, given the fact that z-score extracted only 53% of the target collocations, which was lower than the other four measures except MI, z-score seems to be less suitable for extracting collocations for pedagogical purposes. On the other hand, MI3 succeeded in extracting all the target collocations covered by the other four or three measures, which means that MI3 can reliably evaluate the collocations that other collocational measures rank highly. Therefore, MI3 seems to be the best measure in selecting collocations which are educationally valuable if you try to employ a single measure instead of combining two or more measures.

The investigation into the target collocations which were covered by two or three measures reveals that Dice coefficient and t-score assess collocations in a similar way, and Log-likelihood and z-score produce a similar result. The collocates which comprise the collocations extracted by both Dice coefficient and t-score are much more frequent and fall into more basic levels according to CEFR-J Wordlist than those extracted by Log-likelihood and z-score. These results suggest that Dice coefficient and t-score are suitable for selecting collocations for learners at an elementary level, while Log-likelihood and z-score are appropriate to the needs of learners at an intermediate or advanced level. This result is inconsistent with Ishikawa (2008), who classified

Log-likelihood and t-score into the same group on the based of the correlation with raw frequency of collocations.

Since there is no consensus on how to judge collocations by using association measures, it is necessary to examine how the AMs can be combined to select collocations for educational purposes. Three suggestions can be made from the results of this study.

Firstly, the combination of Dice coefficient and t-score can be utilized so as to extract collocations for learners at an elementary level. Dice coefficient and t-score put emphasis on collocations whose components are relatively frequent in corpora and therefore on a more basic level in terms of the CEFR-J. Secondly, Log-likelihood and z-score enable researchers and educators to select collocations which are useful for learners at an intermediate or advanced level. These two measures place relatively higher value on collocations whose collocates are relatively infrequent, and many of the collocates of the target collocations which can be extracted only by these two statistics fall into B1 level or higher on the basis of the CEFR-J. Thirdly, MI3 can be said to be a well-balanced association measure of collocation. MI3 covered many of the target collocations extracted by the four measures except MI. Therefore, MI3 can be an efficient measure when attempting to select collocations which are worth learning for a wide range of students.

It is important to consider how to apply these findings to create a collocations list. Since there are potentially a vast number of collocations, it is vital to reduce the number of collocations in a collocations list to manageable numbers by selecting pedagogically relevant collocations depending on different proficiency levels of learners.

Based on the characteristics of AMs found in this study, the author would suggest a method of selecting collocations. The first step is to use Dice, t-score, and MI3, extract the top 100 lists respectively (this number can be changed according to the size of the intended collocations list), and identify collocations found in all of the three lists. These collocations are meant for learners at a basic level. The second step is to extract the lists using LL, z-score, and MI3. Collocations

extracted by these three measures are intended for lower-intermediate learners. The final step is to identify collocations found in the lists extracted by LL and z-score, not by MI3, and these collocations are meant for upper-intermediate learners. This is just one of the potential methods of creating a collocations list for learners at a basic or intermediate level by using association measures. How to use different association measures is open to discussion. We also have to consider the phraseological approach to collocation and the effect of learners' L1 on the acquisition of collocation in order to select pedagogically relevant collocations.

## 6.  Conclusion

This study aimed to explore how association measures can be utilized in extracting collocations from corpora from a pedagogical perspective. It is important to conduct this kind of research because finding methods to extract educationally useful collocations by using statistics is necessary for making a collocations list efficiently and objectively.

This study explores the usefulness and characteristics of six association measures (Dice coefficient, Log-likelihood, t-score, z-score, MI3, and MI) by investigating how many of the collocations which appeared in a learning book for collocation, *English Collocation in Use*, were covered by the top 100 lists extracted from the BNC using each of the six AMs.

The result of this study suggests that MI is not appropriate for selecting collocations for pedagogical purposes, as was expected from previous studies. A significant finding is that the other five measures in question can be classified into three groups. The first group includes Dice coefficient and t-score, which could be suitable for selecting collocations for learners at an elementary level. The second group is composed of Log-likelihood and z-score, which can be useful in extracting collocations for intermediate learners. MI3 lies between these two groups, and this seems to be a well-balanced measure which can be appropriate for choosing collocations for a wide range of learn-

ers.

The present study has some methodological limitations. First, this study regards collocations which appeared in *English Collocations in Use* as target collocations. Future research should consider collocations which are not included in the textbook. Second, the collocation patterns which this study dealt with were limited to "verb + noun" collocations only. Other collocational patterns, such as "adjective + noun," "adverb + verb," and "adverb + adjective" should also be explored because different collocational patterns could yield different results.

Although this study provides a clue as to how AMs can be used for selecting collocations for educational purposes, it remains to be seen how this result should be used in creating a collocations list for Japanese learners of English. Further research will be needed to investigate how AM-derived collocations list can be applied to actual classroom practice while considering other pedagogical factors such as cognitive or affective domains of Japanese learners of English.

## REFERENCES

Ackermann, Kirsten, and Yu-Hua Chen. "Developing the Academic Collocation List (Acl) – a Corpus-Driven and Expert-Judged Approach." *Journal of English for Academic Purposes* 12.4 2013: 235–47. Print.

Evert, Stefan. "Corpora and Collocations." *Corpus Linguistics. An International Handbook*. Ed. Anke Lüdeling, Merja Kytö. Berlin: Mouton de Gruyter, 2008. 1212–48. Print.

Granger, Sylviane. "Prefabricated Patterns in Advanced EFL Writing: Collocations and Formulae." *Phraseology: Theory, Analysis and Applications*. Ed. Cowie, A.: Oxford University Press, 1998. 145–60. Print.

Granger, Sylviane, and Paquot Magali. "Disentangling the Phraseological Web." *Phraseology: An Interdisciplinary Perspectives*. Ed. Granger, S. & Meunier, F: Benjamins, 2008. 27–49. Print.

Henriksen, Birgit. "Researching L2 Learners' Collocational Competence and Development - a Progress Report." *L2 Vocabulary Acquisition, Knowledge and Use*. Ed. C Bardel, B Laufer & C Lindqvist: Erosla, 2013. 29–56. Print.

Hill, Jimmie. "Revising Priorities: From Grammatical Failure to Collocational Success." *Teaching Collocation: Further Development in the Lexical Approach*. Ed. Lewis, Michael: Language Teaching Publications, 2000. 47–69. Print.

Hoffmann, Sebastian, and Evert Stefan. "BNCweb (CQP-Edition)." Web. http://corpora.lancs.ac.uk/BNCweb/

Hunston, Susan. *Corpora in Applied Linguistics*. Cambridge University Press, 2002. Print.

Ishikawa, Shinichiro. *Eigo corpus to gengo kyouiku [English Corpus and Language Education]*. Taisyukan, 2008. Print.

Kawamura, Akihiro and Ishii Yasutake. "Communication Nouryoku no tameno Goi Shidou. [Developing students' vocabulary knowledge for better communication skills : with special emphasis on politeness and collocation]" *Seijo University, Social innovation studies* 8.2, 2013: 37–68. Print.

Koya, Taeko. "The Acquisition of Basic Collocations by Japanese Learners of English." Waseda University, 2006. Print.

Koya, Taeko. "Eigo Kyouiku to Collocation [English Education and Collocation." *Korekara no Collocation Kenkyuu [Future Research on Collocation]*. Ed. Hori, Masashiro: Hitsuji Shobou, 2012. 23–60. Print.

Koya, Taeko. "Kihon Collocation list Sakusei no tameno Ichi Kousatu [A study for creating basic collocations list]." *Hosei University Koganei Ronsyu* 11, 2015: 19–32. Print.

Kurosaki, Shino. "An Analysis of Japanese L2 Learners' Knowledge of "Verb + Noun" Collocations." *Jissen Women's University FLC journal* 5, 2010: 63–77. Print.

Leech, Geoffrey, Rayson Paul and Wilson Andrew. *Word Frequencies in Written and Spoken English Based on the British National Corpus*. London: Routledge, 2001. Print.

McCarthy, Michael and O'Dell Felicity. *English Collocations in Use Intermediate [Second Edition]*. Cambridge University Press, 2017. Print.

McEnery, Anthony, Xiao Richard and Yukio Tono. *Corpus-Based Language Studies: An Advanced Resource Book*. Routledge, 2006. Print.

Nation, I.S.P. *Learning Vocabulary in Another Language*. Cambridge University Press, 2013. Print.

Nesselhauf, Nadja. "The Use of Collocations by Advanced Learners of English and Some Implications for Teaching." *Applied Linguistics* 24, 2003: 223–42. Print.

O'Keeffe, Anne, McCarthy Michael and Carter Ronald. *From Corpus to Classroom: Language Use and Language Teaching*. Cambridge University Press, 2007. Print.

Sinclair, John. "Corpus, Concordance, Collocation." *Modern Language Journal* 78, 1991. Print.

Szudarski, Pawe. "Learning and Teaching L2 Collocations: Insights from Research." *TESL Canada Journal* 34, 2017: 205–16. Print.

Tono, Yukio (ed). *Can-Do list Sakusei Katsuyou Eigo Toutatsudo Shihyou CEFR-J Guide book [The Creation and Utilization of Can-Do lists CEFR-J Guidebook*. Taishukan, 2013. Print